

Lecture 1: Introduction

André Martins, Francisco Melo, Mário Figueiredo



Deep Learning Course, Winter 2021-2022

Deep Learning Course

- A new MSc-level course
- Offered jointly by DEEC and DEI
- MSc programs: MEEC, MECD, MEIC-A, MEIC-T
- 475 students enrolled this year!!! (264 DEEC, 211 DEI).



Course Website(s)

<https://fenix.tecnico.ulisboa.pt/disciplinas/AProf-2/2021-2022/1-semester> (DEEC)

<https://fenix.tecnico.ulisboa.pt/disciplinas/AP-Dei/2021-2022/1-semester> (DEI)

There we'll post:

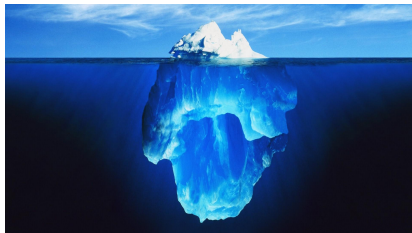
- Syllabus
- Lecture slides
- Literature pointers
- Practical assignments
- Homework assignments
- Announcements
- ...

Instructors

- **Main instructors:** André Martins (DEI Alameda), Francisco Melo (DEI Tagus), Mário Figueiredo (DEEC)
- **Practical classes:** Andreas Wichert, Ben Peters, Gonçalo Faria, João Santinha, José Moreira, Pedro Balage, Rita Ramos, Taisiya Glushkova, Tomás Costa
- **Location & schedule:** see course webpage in Fenix
- **Office hours:** see information in Fenix
- **Communication:**
piazza.com/tecnico.ulisboa.pt/fall2021/c88

Please register in Piazza!!!

What is “Deep Learning”?



- Neural networks?
- Neural networks with many hidden layers?
- Anything beyond shallow (linear) models for machine learning?
- Anything that learns representations?
- A form of learning that is really intense and profound?

Why Did Deep Learning Become Mainstream?

Lots of recent breakthroughs:

- Object recognition
- Speech and language processing (Transformers, BERT, GPT-3)
- Machine translation
- Chatbots and dialog systems
- Self-driving cars
- Solving games (Atari, Go, StarCraft II)
- Protein design (AlphaFold)

No signs of slowing down...



Microsoft's Deep Learning Project Outperforms Humans In Image Recognition



Michael Thomsen, CONTRIBUTOR

I write about tech, video games, science and culture. [FULL BIO](#) ▾

Opinions expressed by Forbes Contributors are their own.



Microsoft's new breakthrough: AI that's as good as humans at listening... on the phone

Microsoft's new speech-recognition record means professional transcribers could be among the first to lose their jobs to artificial intelligence.



By [Liam Tung](#) | October 19, 2016 -- 10:10 GMT (11:10 BST) | Topic: [Innovation](#)

A closer look at Google Duplex

Google's appointment booking AI wowed the crowd and raised concern at I/O

Make a haircut appointment on Tuesday
morning anytime between 10 and 12.
No problem. I'll make you an appointment and
update you soon.



Who is wearing glasses?

man



woman

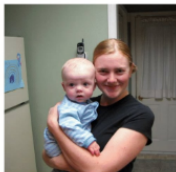


Where is the child sitting?

fridge



arms

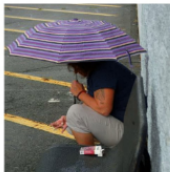


Is the umbrella upside down?

yes



no

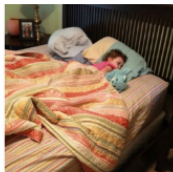


How many children are in the bed?

2



1



The Great A.I. Awakening

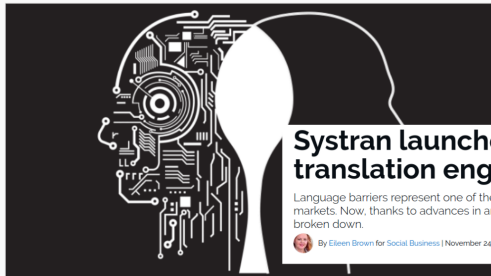
How Google used artificial intelligence to transform Google Translate, one of its more popular services — and how machine learning is poised to reinvent computing itself.

BY GIDEON LEWIS-KRAUS DEC. 14, 2016



Google unleashes deep learning tech on language with Neural Machine Translation

Posted Sep 27, 2016 by [Devin Coldewey](#), Contributor



Systran launches neural machine translation engine in 30 languages

Language barriers represent one of the biggest challenges to develop business strategies among global markets. Now, thanks to advances in artificial intelligence and machine translation, these barriers are being broken down.



By Eileen Brown for Social Business | November 24, 2016 -- 13:49 GMT (13:49 GMT) | Topic: Artificial Intelligence

Siri and Alexa Are Fighting to Be Your Hotel Butler

By **Hui-yong Yu** and **Spencer Soper**

March 22, 2017, 9:00 AM GMT *Updated on* March 22, 2017, 2:13 PM GMT

- Hotels are new frontier for voice-command technologies
- Wynn Las Vegas was first to install Alexa devices in December





AlphaGo Beats Go Human Champ: Godfather Of Deep Learning Tells Us Do Not Be Afraid Of AI

21 March 2016, 10:16 am EDT By [Aaron Mamiit](#) Tech Times



Last week, Google's artificial intelligence program

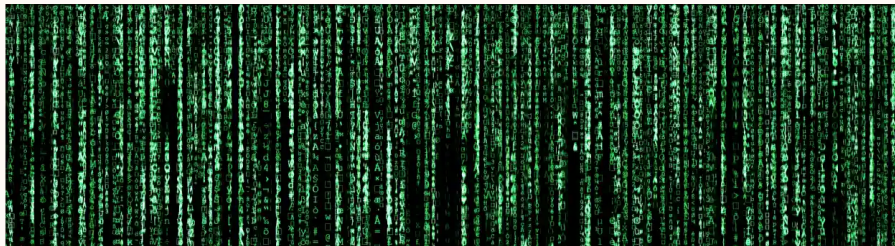
Last week, Google's artificial intelligence program AlphaGo **dominated** its match with South Korean world Go champion Lee Sedol, winning with a 4-1 score.

The achievement stunned artificial intelligence experts, who previously thought that Google's computer program would need at least 10 more years before developing enough to be able to beat a human world champion.

A robot wrote this entire article. Are you scared yet, human?

We asked GPT-3, OpenAI's powerful new language generator, to write an essay for us from scratch. The assignment? To convince us robots come in peace

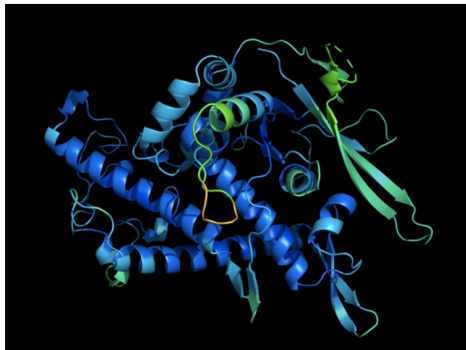
- For more about GPT-3 and how this essay was written and edited, please read our editor's note below



'It will change everything': DeepMind's AI makes gigantic leap in solving protein structures

Google's deep-learning program for determining the 3D shapes of proteins stands to transform biology, say scientists.

Ewen Callaway



A protein's function is determined by its 3D shape. Credit: DeepMind

Why Now?

Why does deep learning work now, but not 30 years ago?

Many of the core ideas were there, after all.

But now we have:

- more data
- more computing power
- (much) better software engineering (e.g. auto-diff)
- some algorithmic innovations (many layers, ReLUs, better learning rates, dropout, CNNs, LSTMs, transformers, etc.)

“But It’s Non-Convex”

For many years (2000–2010), NNs weren’t popular in machine learning because they were hard to learn (e.g. initialization was important)

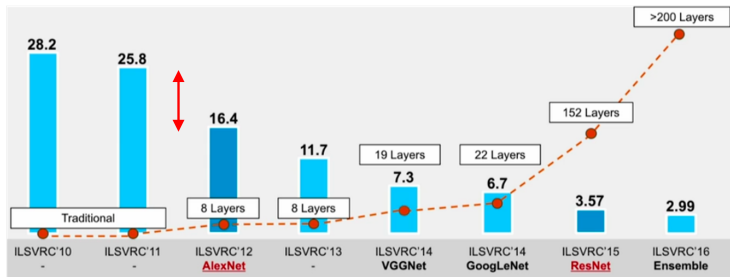
Why does gradient-based optimization work at all in NNs despite the non-convexity?

One possible, partial answer (this is an open research topic)

- there are generally many hidden units
- there are many ways a neural net can approximate the desired input-output relationship
- we only need to find one

One turning point: AlexNet

- Alex Krizhevsky, Ilya Sutskever, Geoffrey Hinton; 2012
- ImageNet: Large Scale Visual Recognition Challenge (14 million images, 20000 categories)

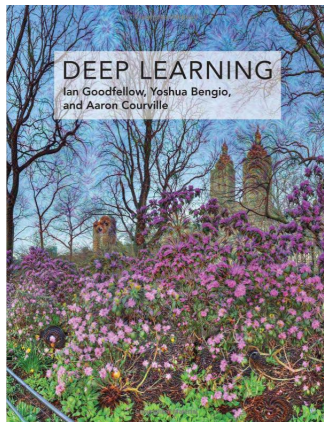


- Large CNN, much deeper than anything else at the time
- Used parallel processing (one of the first uses of GPUs in NNs)
- Convinced many people that deep learning would change the field.

Recommended Books

Main book:

- **Deep Learning.** Ian Goodfellow, Yoshua Bengio, and Aaron Courville. MIT Press, 2016. Chapters available at <http://deeplearningbook.org>



Recommended Books

Secondary books:

- **Artificial Intelligence Engines: A Tutorial Introduction to the Mathematics of Deep Learning.** James Stone. Sebtel Press, 2019.
- **Dive into Deep Learning.** Aston Zhang, Zach Lipton, Mu Li, Alex Smola (<https://d2l.ai/>)
- **Deep Learning with Python.** François Chollet. Manning Publications, 2017.
- **Machine Learning – A Journey to Deep Learning with Exercises and Answers.** Andreas Wichert and Luis Sa-Couto, 2021
- **Machine Learning: a Probabilistic Perspective.** Kevin P. Murphy. MIT Press, 2013.

Tentative Syllabus

- Week 1 Introduction and Course Description
Basic Machine Learning
- Week 2 Linear Classifiers I (linear regression, perceptron)
Linear Classifiers II (logistic regression, regularization)
- Week 3 Neural Networks I
Neural Networks II
- Week 4 Representation Learning and Auto-Encoders
Convolutional Networks
- Week 5 Recurrent Neural Networks and LSTMs
Sequence-to-Sequence Models and Attention Mechanisms
- Week 6 Transformers
Self-Supervised Learning (BERT, GPT3, etc.)
- Week 7 Deep Generative Models (VAEs, GANs)
Interpretability and Fairness

What This Class Is About

- Introduction to **deep learning** (DL)
- **Goal:** after finishing this class, you should be able to:
 - ✓ Understand how DL works (it's not magic)
 - ✓ Understand the math and intuition behind DL models
 - ✓ Apply DL on practical problems (language, vision, ...)
- **Target audience:**
 - ✓ MSc students with basic background in: probability theory, linear algebra, and programming
 - ✓ Preferred (not required): basic background in machine learning.

What This Class Is **Not** About

It's **not** about:

- Just playing with a DL toolkit without learning the fundamental concepts
- Introduction to machine learning (other courses offered by DEEC and DEI)
- Natural language processing (another course offered by DEI)
- Computer vision (another course offered by DEEC)
- Optimization (other courses by DEEC and DEI)
- ...

Prerequisites

- Calculus and basic linear algebra
- Basic probability theory
- Basic knowledge of machine learning (preferred, but not required)
- Programming (Python & PyTorch, preferred but not required)

Course Information

- **Main instructors:** André Martins (DEI Alameda), Francisco Melo (DEI Tagus), Mário Figueiredo (DEEC)
- **Practical classes:** Andreas Wichert, Ben Peters, Gonçalo Faria, João Santinha, José Moreira, Pedro Balage, Rita Ramos, Taisiya Glushkova, Tomás Costa
- **Location & schedule:** see course webpage in Fenix
- **Office hours:** see information in Fenix
- **Communication:**
piazza.com/tecnico.ulisboa.pt/fall2021/c88

Please register in Piazza!!!

Schedule (DEI Alameda/Tagus)

	Seg 12/13	Ter 12/14	Qua 12/15	Qui 12/16	Sex 12/17
07:00					
08:00					
09:00					08:30 - 10:00 L F2
10:00				10:00 - 11:30 L 0 - 21	10:00 - 12:00 T 1 - 2
11:00		11:30 - 13:00 L Q5.1			
12:00		12:00 - 13:00 L Q4.2			12:00 - 13:30 L 1 - 19
13:00					
14:00	14:00 - 15:30 L V1.23	13:30 - 15:00 L C01	14:00 - 16:00 T 1 - 2		13:30 - 15:00 L C12
15:00		15:00 - 17:00 T QA			15:00 - 17:00 T QA
16:00		15:00 - 18:00 L 1 - 15			
17:00		17:00 - 18:00 L 1 - 19		17:00 - 18:30 L Q4.7	17:00 - 18:00 L C12
18:00		17:00 - 18:00 L F2			17:30 - 19:00 L PB
19:00					

Each week:

- 2 theoretical classes (first Mon-Wed, second Thu-Fri)
- 2 practical classes (first Mon-Wed, second Thu-Fri) – **pick your slots and register as a group!** (more later)

Schedule (DEEC Alameda)

	Seg 12/13	Ter 12/14	Qua 12/15	Qui 12/16	Sex 12/17
07:00					
08:00		08:00 - 09:00 L C01	08:00 - 09:00 L E2	08:00 - 09:00 L E1	08:00 - 09:00 L E3
09:00					08:30 - 10:30 T EA1
10:00	09:30 - 11:30 T GA5	09:30 - 11:00 L F8			
11:00		11:00 - 12:00 L C12	11:30 - 13:00 T GA1		10:30 - 12:00 L E4
12:00					12:00 - 13:30 L E4
13:00				12:30 - 14:00 L E1	
14:00	13:30 - 15:00 L E5			14:00 - 16:00 T EA2	14:00 - 15:00 L E1
15:00	13:30 - 15:00 L E8				14:00 - 15:00 L E5
16:00					14:00 - 15:00 L Q4.2

Each week:

- Lecture shifts (shift 1: Mon & Fri; shift 2: Tue & Thu)
- Practical shifts: **pick your slots and register as a group!** (more later!)

Grading

- 2 homework assignments: 50%
 - Minimal grade: 9.5
 - Theoretical questions & implementation
 - Groups of 3 – you need to register your group in Fenix!
 - Some of the practicals will be Q&A about these assignments, so please join the practicals as a group
 - Submission through Fenix
 - No late days allowed
- Final exam: 50%
 - Minimal grade: 9.5

Registering Your Group in Fenix

- Pick a group of 3
- Register in Fenix
- **Deadline: Sunday, December 5**
- Use Piazza to find group mates
- If you can't find a group, let the instructors know by December 6 (in Piazza), and we'll find a solution.

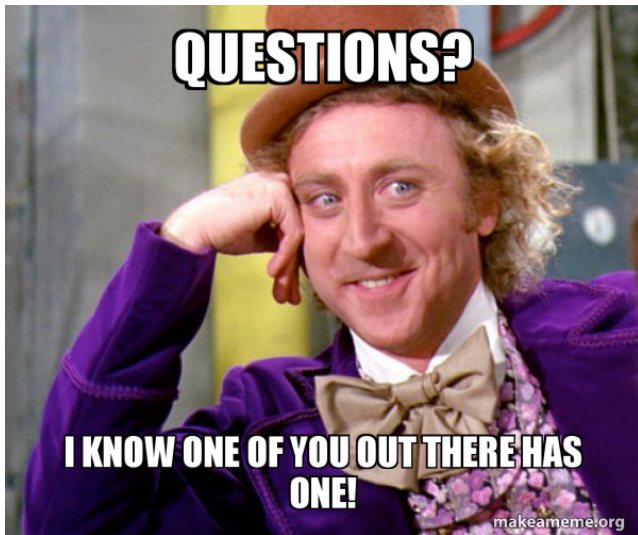
Collaboration Policy

- Assignments should be done within each group
- Students may discuss the questions across groups, as long as they write their own answers and their own code
- If this happens, acknowledge with whom you collaborate!
- Zero tolerance on plagiarism!!
- Always credit your sources!!!

Caveat

- This is the first year we're teaching this class
- Also the first year with trimesters
- ... which means you're the first batch of students taking it :)
- **Constructive feedback will be highly appreciated!**

Questions?



Notation: Matrices and Vectors

- $A \in \mathbb{R}^{m \times n}$ is a **matrix** with m rows and n columns.

$$A = \begin{bmatrix} A_{1,1} & \cdots & A_{1,n} \\ \vdots & \ddots & \vdots \\ A_{m,1} & \cdots & A_{m,n} \end{bmatrix}.$$

- $x \in \mathbb{R}^n$ is a **vector** with n components,

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}.$$

- A **(column) vector** is a matrix with n rows and 1 column.
- A matrix with 1 row and n columns is called a **row vector**.

Matrix Transpose and Matrix Products

- Given matrix $A \in \mathbb{R}^{m \times n}$, its **transpose** A^T is such that $(A^T)_{i,j} = A_{j,i}$.
- Matrix A is **symmetric** if $A^T = A$.
- Given matrices $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times p}$, their **product** is

$$C = AB \in \mathbb{R}^{m \times p} \quad \text{where} \quad C_{i,j} = \sum_{k=1}^n A_{i,k} B_{k,j}$$

- **Inner product** between vectors $x, y \in \mathbb{R}^n$:

$$\langle x, y \rangle = x^T y = y^T x = \sum_{i=1}^n x_i y_i$$

$$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}^T \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = [x_1, \dots, x_n] \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \sum_{i=1}^n x_i y_i$$

Outer and Hadamard Products

- **Outer product** between vectors $x \in \mathbb{R}^n$ and $y \in \mathbb{R}^m$:

$$x y^T \in \mathbb{R}^{n \times m}, \quad \text{where } (x y^T)_{ij} = x_i y_j$$

$$x y^T = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} [y_1, \dots, y_m] = \begin{bmatrix} x_1 y_1 & \cdots & x_1 y_m \\ \vdots & \ddots & \vdots \\ x_n y_1 & \cdots & x_n y_m \end{bmatrix}$$

- **Hadamard/Schur product** between vectors $x, y \in \mathbb{R}^n$: $(x \odot y)_i = x_i y_i$,

$$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \odot \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_1 y_1 \\ \vdots \\ x_n y_n \end{bmatrix}$$

Properties of Matrix Products and Transposes

- Given matrices $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times p}$, their **product** is

$$C = AB \in \mathbb{R}^{m \times p} \quad \text{where} \quad C_{i,j} = \sum_{k=1}^n A_{i,k} B_{k,j}$$

- Matrix product is **associative**: $(AB)C = A(BC)$.
- In general, matrix product is not **commutative**: $AB \neq BA$.
- Transpose of product: $(AB)^T = B^T A^T$.
- Transpose of sum: $(A + B)^T = A^T + B^T$.

Norms

- The **norm** of a vector is (informally) its “magnitude.” Euclidean norm:

$$\|x\|_2 = \sqrt{\langle x, x \rangle} = \sqrt{x^T x} = \sqrt{\sum_{i=1}^n x_i^2}.$$

- More generally, the ℓ_p norm of a vector $x \in \mathbb{R}^n$, where $p \geq 1$,

$$\|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}.$$

- Notable case: the ℓ_1 norm, $\|x\|_1 = \sum_i |x_i|$.
- Notable case: the ℓ_∞ norm, $\|x\|_\infty = \max\{|x_1|, \dots, |x_n|\}$.

Special Matrices

- The **identity matrix** $I \in \mathbb{R}^{n \times n}$ is a square matrix such that

$$I_{ij} = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases} \quad I = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}$$

- Neutral element of matrix product: $AI = IA = A$.
- Diagonal matrix: $A \in \mathbb{R}^{n \times n}$ is diagonal if $(i \neq j) \Rightarrow A_{i,j} = 0$.
- Upper triangular matrix: $(j < i) \Rightarrow A_{i,j} = 0$.
- Lower triangular matrix: $(j > i) \Rightarrow A_{i,j} = 0$.

Eigenvalues, eigenvectors, determinant, trace

- A vector $x \in \mathbb{R}^n$ is an **eigenvector** of matrix $A \in \mathbb{R}^{n \times n}$ if

$$Ax = \lambda x,$$

where $\lambda \in \mathbb{R}$ is the corresponding **eigenvalue**.

- The eigenvalues of a diagonal matrix are the elements in the diagonal.
- Matrix **trace**:

$$\text{trace}(A) = \sum_i A_{i,i} = \sum_i \lambda_i$$

- Matrix **determinant**:

$$|A| = \det(A) = \prod_i \lambda_i$$

- Properties: $|AB| = |A||B|$, $|A^T| = |A|$, $|\alpha A| = \alpha^n |A|$

Matrix Inverse

- Matrix $A \in \mathbb{R}^{n \times n}$ is **invertible** if there is $B \in \mathbb{R}^{n \times n}$ s.t. $AB = BA = I$.
- ...matrix B , such that $AB = BA = I$, denoted $B = A^{-1}$.
- Matrix $A \in \mathbb{R}^{n \times n}$ is invertible $\Leftrightarrow \det(A) \neq 0$.
- Determinant of inverse: $\det(A^{-1}) = \frac{1}{\det(A)}$.
- Solving system $Ax = b$, if A is invertible: $x = A^{-1}b$.
- Properties: $(A^{-1})^{-1} = A$, $(A^{-1})^T = (A^T)^{-1}$, $(AB)^{-1} = B^{-1}A^{-1}$
- There are many algorithms to compute A^{-1} ; general case, computational cost $O(n^3)$.

Quadratic Forms and Positive (Semi-)Definite Matrices

- Given matrix $A \in \mathbb{R}^{n \times n}$ and vector $x \in \mathbb{R}^n$,

$$x^T A x = \sum_{i=1}^n \sum_{j=1}^n A_{i,j} x_i x_j \in \mathbb{R}$$

is called a **quadratic form**.

- A symmetric matrix $A \in \mathbb{R}^{n \times n}$ is **positive semi-definite** (PSD) if, for any $x \in \mathbb{R}^n$, $x^T A x \geq 0$.
- A symmetric matrix $A \in \mathbb{R}^{n \times n}$ is **positive definite** (PD) if, for any $x \in \mathbb{R}^n$, $(x \neq 0) \Rightarrow x^T A x > 0$.
- Matrix $A \in \mathbb{R}^{n \times n}$ is PSD \Leftrightarrow all $\lambda_i(A) \geq 0$.
- Matrix $A \in \mathbb{R}^{n \times n}$ is PD \Leftrightarrow all $\lambda_i(A) > 0$.

Minimizing a function

- We are given a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$.
- Goal: find x^* that minimizes $f : \mathbb{R}^n \rightarrow \mathbb{R}$.
- **Global minimum**: for any $x \in \mathbb{R}^n$, $f(x^*) \leq f(x)$.
- **Local minimum**: for any $\|x - x^*\| \leq \delta \Rightarrow f(x^*) \leq f(x)$.
- Are local minima also global minima?

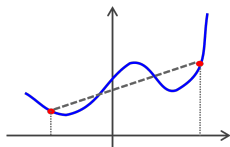
Convex Functions

- Function $f : \mathbb{R}^n \rightarrow \mathbb{R}$
- f is a **convex function**, if, for any $\lambda \in [0, 1]$, and any x, x' ,

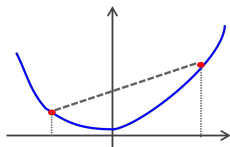
$$f(\lambda x + (1 - \lambda)x') \leq \lambda f(x) + (1 - \lambda)f(x')$$

- f is a **strictly convex function**, if, for any $\lambda \in]0, 1[$, and any x, x' ,

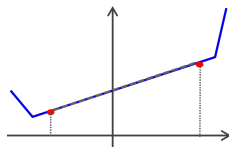
$$f(\lambda x + (1 - \lambda)x') < \lambda f(x) + (1 - \lambda)f(x')$$



non-convex



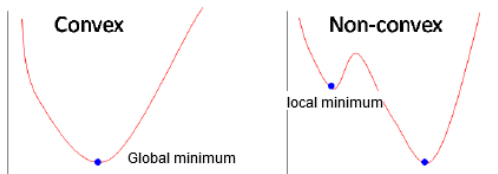
convex
strictly convex



convex, not strictly

Relationship Between Convexity and Minimization

- Goal: find x^* that minimizes $f : \mathbb{R}^n \rightarrow \mathbb{R}$.
- If f is **convex** and x^* is a **local minimizer**, then it is also a **global minimizer**.
- If f is **strictly convex** and x^* is a **local minimizer**, then it is also the **unique global minimizer**.



Gradients and Minimization

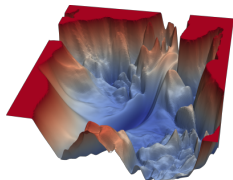
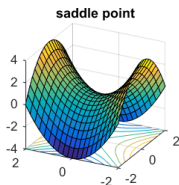
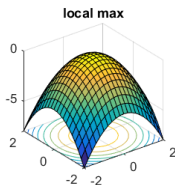
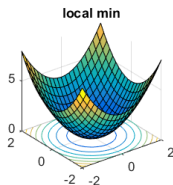
- Given $f : \mathbb{R}^n \rightarrow \mathbb{R}$ (differentiable), the **gradient** of f at x :

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n$$

- Relationship between gradient and minimization

$$x^* \text{ is local minimizer} \Rightarrow \nabla f(x^*) = 0$$

~~\Leftarrow~~



Hessians and Convexity

- Given $f : \mathbb{R}^n \rightarrow \mathbb{R}$ (differentiable), the **Hessian** of f at x :

$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_n^2} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

- Relationship between Hessian and convexity:
 - ✓ Positive semi-definite Hessian $\Leftrightarrow f$ is convex
 - ✓ Positive definite Hessian $\Leftrightarrow f$ is strictly convex.

More on Gradients

- Gradient of quadratic form $f(x) = x^T Ax$: $\nabla f(x) = (A + A^T)x$
- ...if A symmetric: $\nabla f(x) = 2Ax$
- .Particular case: $f(x) = x^T x = \|x\|_2^2$, then $\nabla f(x) = 2x$
- If $f(x) = x^T b = b^T x$, then $\nabla f(x) = b$
- If $g(x) = f(Ax)$, then $\nabla g(x) = A^T \nabla f(Ax)$
- If $g(x) = f(a \odot x)$, then $\nabla g(x) = a \odot \nabla f(a \odot x)$
- In simple cases, we can find minima analytically: $f(x) = \|Ax - y\|_2^2$

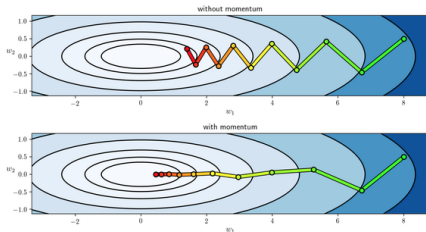
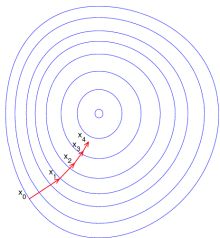
$$\nabla f(x) = 2A^T(Ax - y) = 0 \Rightarrow x^* = (A^T A)^{-1} A^T y$$

Gradient Descent

- Goal: minimize $f : \mathbb{R}^d \rightarrow \mathbb{R}$, for differentiable f
- Take **small steps** in the **negative gradient direction** until a **stopping criterion** is met:

$$x^{(t+1)} \leftarrow x^{(t)} - \eta_{(t)} \nabla f(x^{(t)})$$

- Choosing the **step-size**: crucial for convergence and performance.
- GD may work well, or not so well. There are many ways to improve it.

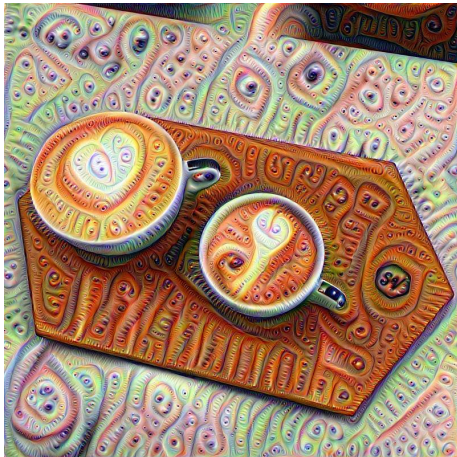


Recommended Reading

- Z. Kolter and C. Do, “Linear Algebra Review and Reference”, Stanford University, 2015 (<https://tinyurl.com/44x2qj4>)

Thank you!

Questions?



Lecture 2: Machine Learning Basics

André Martins, Francisco Melo, Mário Figueiredo



Deep Learning Course, Winter 2021-2022

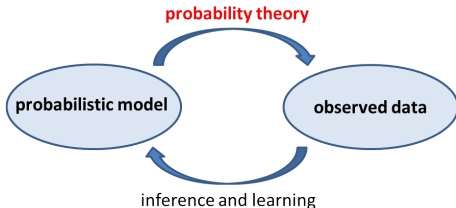
Announcements

- Please don't forget to register your groups in Fenix!!
- **Deadline this Sunday!**
- If needed, use Piazza to find a group.
- Many shifts are still available.

Today's Roadmap

- Probability Refresher
- Introduction to Machine Learning
 - “Deep Learning Superheroes”
 - Supervised, Unsupervised, Reinforcement Learning
 - Classification and Regression
- Naive Bayes Classifier

Probability theory



- “Essentially, all models are wrong, but some are useful”; [G. Box, 1987](#)
- The study of probability has roots in games of chance (dice, cards, ...)
- Natural tool to model [uncertainty, information, knowledge, belief, ...](#)
- ...thus also [learning, inference, ...](#)

What is probability?

- Classical definition: $\mathbb{P}(A) = \frac{N_A}{N}$

...with N mutually exclusive equally likely outcomes,
 N_A of which result in the occurrence of event A .

Laplace, 1814

Example: $\mathbb{P}(\text{randomly drawn card is } \clubsuit) = 13/52$.

Example: $\mathbb{P}(\text{getting 1 in throwing a fair die}) = 1/6$.

- Frequentist definition: $\mathbb{P}(A) = \lim_{N \rightarrow \infty} \frac{N_A}{N}$

...relative frequency of occurrence of A in infinite number of trials.

- Subjective probability: $\mathbb{P}(A)$ is a degree of belief.

de Finetti, 1930s

...gives meaning to $\mathbb{P}(\text{"Tomorrow it will rain"})$.

Key concepts: Sample space and events

- **Sample space** \mathcal{X} = set of possible outcomes of a random experiment.

Examples:

- Tossing two coins: $\mathcal{X} = \{HH, TH, HT, TT\}$
 - Roulette: $\mathcal{X} = \{1, 2, \dots, 36\}$
 - Draw a card from a shuffled deck: $\mathcal{X} = \{A\clubsuit, 2\clubsuit, \dots, Q\heartsuit, K\heartsuit\}$.
- An **event** A is a subset of \mathcal{X} : $A \subseteq \mathcal{X}$.

Examples:

- “exactly one H in 2-coin toss”: $A = \{TH, HT\} \subset \{HH, TH, HT, TT\}$.
- “odd number in the roulette”: $B = \{1, 3, \dots, 35\} \subset \{1, 2, \dots, 36\}$.
- “drawn a \heartsuit card”: $C = \{A\heartsuit, 2\heartsuit, \dots, K\heartsuit\} \subset \{A\clubsuit, \dots, K\heartsuit\}$

Kolmogorov's Axioms for Probability

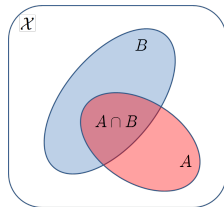
- Probability is a function that maps events A into the interval $[0, 1]$.

Kolmogorov's axioms (1933) for probability \mathbb{P}

- For any A , $\mathbb{P}(A) \geq 0$
- $\mathbb{P}(\mathcal{X}) = 1$
- If $A_1, A_2 \dots \subseteq \mathcal{X}$ are disjoint events, then $\mathbb{P}\left(\bigcup_i A_i\right) = \sum_i \mathbb{P}(A_i)$

- From these axioms, many results can be derived. **Examples:**

- $\mathbb{P}(\emptyset) = 0$
- $C \subset D \Rightarrow \mathbb{P}(C) \leq \mathbb{P}(D)$
- $\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B)$



Conditional Probability and Independence

- If $\mathbb{P}(B) > 0$, $\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$
- Events A, B are independent ($A \perp B$) $\Leftrightarrow \mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B)$.
- Relationship with conditional probabilities:

$$A \perp B \Leftrightarrow \mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} = \frac{\mathbb{P}(A)\mathbb{P}(B)}{\mathbb{P}(B)} = \mathbb{P}(A)$$

- Example: $\mathcal{X} = \text{"52 cards"}$, $A = \{3\heartsuit, 3\clubsuit, 3\diamondsuit, 3\spadesuit\}$, and $B = \{A\heartsuit, 2\heartsuit, \dots, K\heartsuit\}$; then, $\mathbb{P}(A) = 1/13$, $\mathbb{P}(B) = 1/4$

$$\mathbb{P}(A \cap B) = \mathbb{P}(\{3\heartsuit\}) = \frac{1}{52}$$

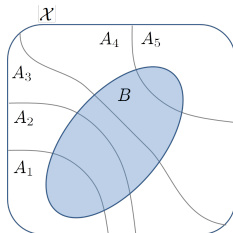
$$\mathbb{P}(A)\mathbb{P}(B) = \frac{1}{13} \frac{1}{4} = \frac{1}{52}$$

$$\mathbb{P}(A|B) = \mathbb{P}(\text{"3"} | \text{"\heartsuit"}) = \frac{1}{13} = \mathbb{P}(A)$$

Law of Total Probability and Bayes Theorem

- Law of total probability: if A_1, \dots, A_n are a partition of \mathcal{X}

$$\begin{aligned}\mathbb{P}(B) &= \sum_i \mathbb{P}(B|A_i)\mathbb{P}(A_i) \\ &= \sum_i \mathbb{P}(B \cap A_i)\end{aligned}$$

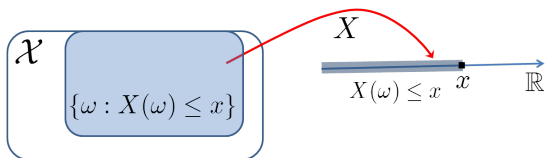


- Bayes' theorem: if $\{A_1, \dots, A_n\}$ is a partition of \mathcal{X}

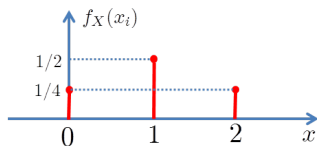
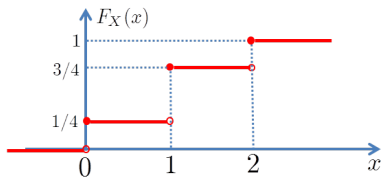
$$\mathbb{P}(A_i|B) = \frac{\mathbb{P}(B \cap A_i)}{\mathbb{P}(B)}$$

Discrete Random Variables

- Distribution function:** $F_X(x) = \mathbb{P}(\{\omega \in \mathcal{X} : X(\omega) \leq x\})$



- Example:** number of heads in tossing 2 coins; $\text{range}(X) = \{0, 1, 2\}$.



- Probability mass function** (discrete RV): $f_X(x) = \mathbb{P}(X = x)$,

$$F_X(x) = \sum_{x_i \leq x} f_X(x_i).$$

Important Discrete Random Variables

- **Uniform:** $X \in \{x_1, \dots, x_K\}$, pmf $f_X(x_i) = 1/K$.
- **Bernoulli RV:** $X \in \{0, 1\}$, pmf $f_X(x) = \begin{cases} p & \Leftarrow x = 1 \\ 1 - p & \Leftarrow x = 0 \end{cases}$

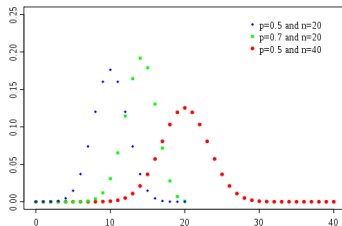
Can be written compactly as $f_X(x) = p^x(1-p)^{1-x}$.

- **Binomial RV:** $X \in \{0, 1, \dots, n\}$ (sum on n Bernoulli RVs)

$$f_X(x) = \text{Binomial}(x; n, p) = \binom{n}{x} p^x (1-p)^{(n-x)}$$

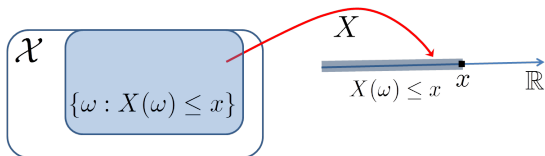
Binomial coefficients
("n choose x"):

$$\binom{n}{x} = \frac{n!}{(n-x)! x!}$$

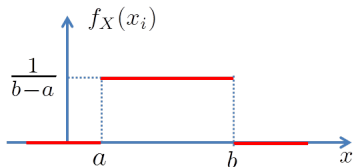
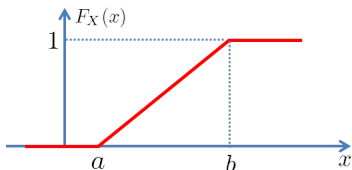


Continuous Random Variables

- Distribution function:** $F_X(x) = \mathbb{P}(\{\omega \in \mathcal{X} : X(\omega) \leq x\})$



- Example:** continuous RV with uniform distribution on $[a, b]$.



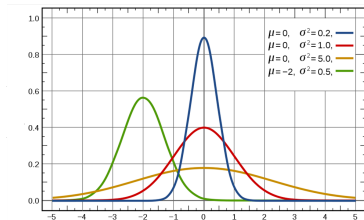
- Probability density function (pdf, continuous RV):** $f_X(x)$

$$F_X(x) = \int_{-\infty}^x f_X(u) du, \quad \mathbb{P}(X \in [c, d]) = \int_c^d f_X(x) dx, \quad \mathbb{P}(X = x) = 0$$

Important Continuous Random Variables

- **Uniform:** $f_X(x) = \text{Uniform}(x; a, b) = \begin{cases} \frac{1}{b-a} & \Leftarrow x \in [a, b] \\ 0 & \Leftarrow x \notin [a, b] \end{cases}$
(previous slide).

- **Gaussian:** $f_X(x) = \mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$



Expectation of Random Variables

- **Expectation:** $\mathbb{E}(X) = \begin{cases} \sum x_i f_X(x_i) & X \in \{x_1, \dots, x_K\} \subset \mathbb{R} \\ \int_{-\infty}^{\infty} x f_X(x) dx & X \text{ continuous} \end{cases}$

- **Example:** Bernoulli, $f_X(x) = p^x (1-p)^{1-x}$, for $x \in \{0, 1\}$.

$$\mathbb{E}(X) = 0(1-p) + 1p = p.$$

- **Example:** Gaussian, $f_X(x) = \mathcal{N}(x; \mu, \sigma^2)$. $\mathbb{E}(X) = \mu$.

- **Linearity of expectation:**

$$\mathbb{E}(X + Y) = \mathbb{E}(X) + \mathbb{E}(Y); \quad \mathbb{E}(\alpha X) = \alpha \mathbb{E}(X), \quad \alpha \in \mathbb{R}$$

Expectation of Functions of Random Variables

- $\mathbb{E}(g(X)) = \begin{cases} \sum g(x_i) f_X(x_i) & X \text{ discrete, } g(x_i) \in \mathbb{R} \\ \int_{-\infty}^{\infty} g(x) f_X(x) dx & X \text{ continuous} \end{cases}$
- **Example:** variance, $\text{var}(X) = \mathbb{E}((X - \mathbb{E}(X))^2) = \mathbb{E}(X^2) - \mathbb{E}(X)^2$
- **Example:** Bernoulli variance, $\mathbb{E}(X^2) = \mathbb{E}(X) = p$, thus $\text{var}(X) = p(1 - p)$.
- **Example:** Gaussian variance, $\mathbb{E}((X - \mu)^2) = \sigma^2$.
- Probability as expectation of indicator, $\mathbf{1}_A(x) = \begin{cases} 1 & \Leftarrow x \in A \\ 0 & \Leftarrow x \notin A \end{cases}$

$$\mathbb{P}(X \in A) = \int_A f_X(x) dx = \int \mathbf{1}_A(x) f_X(x) dx = \mathbb{E}(\mathbf{1}_A(X))$$

Two (or More) Random Variables

- **Joint pmf** of two discrete RVs: $f_{X,Y}(x,y) = \mathbb{P}(X = x \wedge Y = y)$.

Extends trivially to more than two RVs.

- **Joint pdf** of two continuous RVs: $f_{X,Y}(x,y)$, such that

$$\mathbb{P}((X, Y) \in A) = \iint_A f_{X,Y}(x,y) dx dy.$$

Extends trivially to more than two RVs.

- **Marginalization:** $f_Y(y) = \begin{cases} \sum_x f_{X,Y}(x,y), & \text{if } X \text{ is discrete} \\ \int_{-\infty}^{\infty} f_{X,Y}(x,y) dx, & \text{if } X \text{ continuous} \end{cases}$

- **Independence:**

$$X \perp\!\!\!\perp Y \Leftrightarrow f_{X,Y}(x,y) = f_X(x) f_Y(y) \stackrel{\Rightarrow}{\neq} \mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y).$$

Conditionals and Bayes' Theorem

- **Conditional pmf** (discrete RVs):

$$f_{X|Y}(x|y) = \mathbb{P}(X = x|Y = y) = \frac{\mathbb{P}(X = x \wedge Y = y)}{\mathbb{P}(Y = y)} = \frac{f_{X,Y}(x,y)}{f_Y(y)}.$$

- **Conditional pdf** (continuous RVs): $f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)}$

- **Bayes' theorem:** $f_{X|Y}(x|y) = \frac{f_{Y|X}(y|x) f_X(x)}{f_Y(y)}$ (pdf or pmf).

- Also valid in the mixed case (e.g., X continuous, Y discrete).

Joint, Marginal, and Conditional Probabilities: An Example

- A pair of binary variables $X, Y \in \{0, 1\}$, with **joint** pmf:

$f_{X,Y}(x,y)$	$Y=0$	$Y=1$
$X=0$	1/5	2/5
$X=1$	1/10	3/10

- Marginals:** $f_X(0) = \frac{1}{5} + \frac{2}{5} = \frac{3}{5}$, $f_X(1) = \frac{1}{10} + \frac{3}{10} = \frac{4}{10}$,
 $f_Y(0) = \frac{1}{5} + \frac{1}{10} = \frac{3}{10}$, $f_Y(1) = \frac{2}{5} + \frac{3}{10} = \frac{7}{10}$.

- Conditional** probabilities:

$f_{X Y}(x y)$	$Y=0$	$Y=1$
$X=0$	2/3	4/7
$X=1$	1/3	3/7

$f_{Y X}(y x)$	$Y=0$	$Y=1$
$X=0$	1/3	2/3
$X=1$	1/4	3/4

An Important Multivariate RV: Multinomial

- Multinomial:** $X = (X_1, \dots, X_K)$, $X_i \in \{0, \dots, n\}$, such that $\sum_i X_i = n$,

$$f_X(x_1, \dots, x_K) = \begin{cases} \binom{n}{x_1 \ x_2 \ \dots \ x_K} p_1^{x_1} p_2^{x_2} \dots p_K^{x_K} & \Leftarrow \sum_i x_i = n \\ 0 & \Leftarrow \sum_i x_i \neq n \end{cases}$$

$$\binom{n}{x_1 \ x_2 \ \dots \ x_K} = \frac{n!}{x_1! x_2! \dots x_K!}$$

Parameters: $p_1, \dots, p_K \geq 0$, such that $\sum_i p_i = 1$.

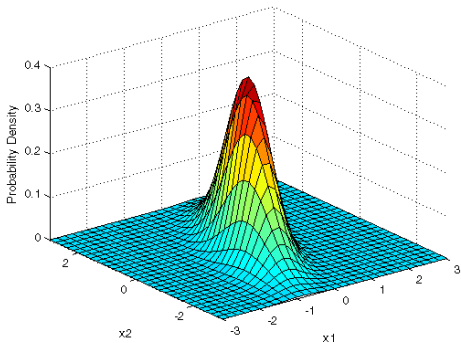
- Generalizes the binomial from binary to K -classes.
- Example:** tossing n independent fair dice, $p_1 = \dots = p_6 = 1/6$.
 x_i = number of outcomes with i dots. Of course, $\sum_i x_i = n$.
- For $n = 1$, sometimes called **categorical** or **multinoulli**.
For $n = 1$, one and only one $x_i = 1$, others are 0, thus $\binom{n}{x_1 \ \dots \ x_K} = 1$.

An Important Multivariate RV: Gaussian

- **Multivariate Gaussian:** $X \in \mathbb{R}^n$,

$$f_X(x) = \mathcal{N}(x; \mu, C) = \frac{1}{\sqrt{\det(2\pi C)}} \exp\left(-\frac{1}{2}(x - \mu)^T C^{-1}(x - \mu)\right)$$

- Parameters: vector $\mu \in \mathbb{R}^n$ and matrix $C \in \mathbb{R}^{n \times n}$.
Expected value: $\mathbb{E}(X) = \mu$. Meaning of C : next slide.



Covariance, Correlation, and all that...

- **Covariance** between two RVs:

$$\text{cov}(X, Y) = \mathbb{E} \left[(X - \mathbb{E}(X)) (Y - \mathbb{E}(Y)) \right] = \mathbb{E}(X Y) - \mathbb{E}(X) \mathbb{E}(Y)$$

- Relationship with variance: $\text{var}(X) = \text{cov}(X, X)$.

- $X \perp\!\!\!\perp Y \Leftrightarrow f_{X,Y}(x,y) = f_X(x) f_Y(y) \stackrel{\Rightarrow}{\neq} \text{cov}(X, Y) = 0$

- **Covariance matrix** of multivariate RV, $X \in \mathbb{R}^n$:

$$\text{cov}(X) = \mathbb{E} \left[(X - \mathbb{E}(X)) (X - \mathbb{E}(X))^T \right] = \mathbb{E}(X X^T) - \mathbb{E}(X) \mathbb{E}(X)^T$$

- Covariance of Gaussian RV, $f_X(x) = \mathcal{N}(x; \mu, C) \Rightarrow \text{cov}(X) = C$

More on Expectations and Covariances

Let $A \in \mathbb{R}^{n \times n}$ be a matrix and $a \in \mathbb{R}^n$ a vector.

- If $\mathbb{E}(X) = \mu$ and $Y = AX$, then $\mathbb{E}(Y) = A\mu$;
- If $\text{cov}(X) = C$ and $Y = AX$, then $\text{cov}(Y) = ACA^T$;
- If $\text{cov}(X) = C$ and $Y = a^T X \in \mathbb{R}$, then $\text{var}(Y) = a^T C a \geq 0$;
- If $\text{cov}(X) = C$ and $Y = C^{-1/2}X$, then $\text{cov}(Y) = I$;
- If $f_X(x) = \mathcal{N}(x; 0, I)$ and $Y = \mu + C^{1/2}X$, then $f_Y(y) = \mathcal{N}(y; \mu, C)$;
- If $f_X(x) = \mathcal{N}(x; \mu, C)$ and $Y = C^{-1/2}(X - \mu)$, then $f_Y(y) = \mathcal{N}(y; 0, I)$.

Entropy and all that...

Entropy of a discrete RV $X \in \{1, \dots, K\}$:

$$H(X) = - \sum_{x=1}^K f_X(x) \log f_X(x)$$

- **Positivity:** $H(X) \geq 0$;
 $H(X) = 0 \Leftrightarrow f_X(i) = 1$, for exactly one $i \in \{1, \dots, K\}$.
- **Upper bound:** $H(X) \leq \log K$;
 $H(X) = \log K \Leftrightarrow f_X(x) = 1/k$, for all $x \in \{1, \dots, K\}$
- Measure of **uncertainty/randomness** of X

Continuous RV X , **differential entropy**:

$$h(X) = - \int f_X(x) \log f_X(x) dx$$

- $h(X)$ can be positive or negative. Example, if $f_X(x) = \text{Uniform}(x; a, b)$, $h(X) = \log(b - a)$.
- If $f_X(x) = \mathcal{N}(x; \mu, \sigma^2)$, then $h(X) = \frac{1}{2} \log(2\pi e\sigma^2)$.

Kullback-Leibler divergence

Kullback-Leibler divergence (KLD) between two pmf:

$$D(f_X \| g_X) = \sum_{x=1}^K f_X(x) \log \frac{f_X(x)}{g_X(x)}$$

Positivity: $D(f_X \| g_X) \geq 0$

$$D(f_X \| g_X) = 0 \Leftrightarrow f_X(x) = g_X(x), \text{ for } x \in \{1, \dots, K\}$$

KLD between two pdf:

$$D(f_X \| g_X) = \int f_X(x) \log \frac{f_X(x)}{g_X(x)} dx$$

Positivity: $D(f_X \| g_X) \geq 0$

$$D(f_X \| g_X) = 0 \Leftrightarrow f_X(x) = g_X(x), \text{ almost everywhere}$$

Recommended Reading

- A. Maleki and T. Do, “Review of Probability Theory”, Stanford University, 2017 (<https://tinyurl.com/pz7p9g5>)
- L. Wasserman, “All of Statistics: A Concise Course in Statistical Inference”, Springer, 2004.

Machine Learning

Tom Mitchell's definition:

- “A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T , as measured by P , improves with experience E .”
- In a nutshell: **learn from data; improve performance with experience**

This comes in many flavours:

- Supervised learning
- Unsupervised learning
- Self-supervised learning
- Reinforcement learning
- Active learning

Formulate the problem; get data; learn the model from the data; evaluate.

Example Tasks

Binary classification: given an e-mail: is it spam or not-spam?

Multi-class classification: given a news article, determine its topic (politics, sports, etc.)

Regression: how much time a person will spend reading this article?

    **AlphaGo Beats Go Human Champ:
Godfather Of Deep Learning Tells Us Do
Not Be Afraid Of AI**

21 March 2016, 10:16 am EDT By Aaron Mamlit Tech Times



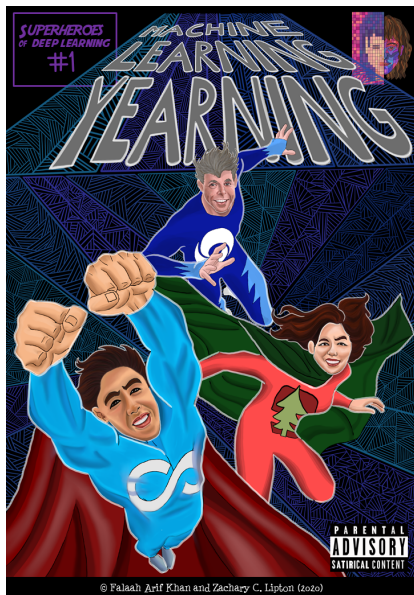
Last week, Google's artificial intelligence program

Last week, Google's artificial intelligence program AlphaGo **dominated** its match with South Korean world Go champion Lee Sedol, winning with a 4-1 score.

The achievement stunned artificial intelligence experts, who previously thought that Google's computer program would need at least 10 more years before developing enough to be able to beat a human world champion.

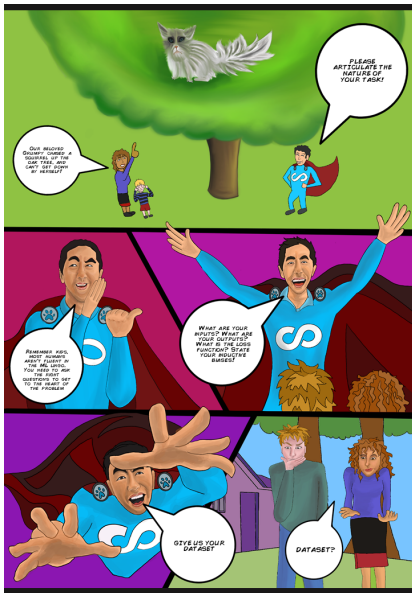


sports
politics
technology
economy
weather
culture



<https://www.approximatelycorrect.com/2020/10/26/>

superheroes-of-deep-learning-vol-1-machine-learning-yearning/



<https://www.approximatelycorrect.com/2020/10/26/>

superheroes-of-deep-learning-vol-1-machine-learning-yearning/



<https://www.approximatelycorrect.com/2020/10/26/>

superheroes-of-deep-learning-vol-1-machine-learning-yearning/



<https://www.approximatelycorrect.com/2020/10/26/>

superheroes-of-deep-learning-vol-1-machine-learning-yearning/

Let's Start Simple

- Example 1 – sequence: $\star \diamond \circ$; label: -1
- Example 2 – sequence: $\star \heartsuit \triangle$; label: -1
- Example 3 – sequence: $\star \triangle \spadesuit$; label: $+1$
- Example 4 – sequence: $\diamond \triangle \circ$; label: $+1$

- New sequence: $\star \diamond \circ$; label ? -1
- New sequence: $\star \diamond \heartsuit$; label ? -1
- New sequence: $\star \triangle \circ$; label ?

Why can we do this?

(Credits: Ryan McDonald)

Let's Start Simple: Machine Learning

- Example 1 – sequence: $\star \diamond \circ$; label: -1
- Example 2 – sequence: $\star \heartsuit \triangle$; label: -1
- Example 3 – sequence: $\star \triangle \spadesuit$; label: $+1$
- Example 4 – sequence: $\diamond \triangle \circ$; label: $+1$
- New sequence: $\star \triangle \circ$; label ?

Label -1

Label $+1$

$$P(-1|\star) = \frac{\text{count}(\star \text{ and } -1)}{\text{count}(\star)} = \frac{2}{3} = 0.67 \text{ vs. } P(+1|\star) = \frac{\text{count}(\star \text{ and } +1)}{\text{count}(\star)} = \frac{1}{3} = 0.33$$
$$P(-1|\triangle) = \frac{\text{count}(\triangle \text{ and } -1)}{\text{count}(\triangle)} = \frac{1}{3} = 0.33 \text{ vs. } P(+1|\triangle) = \frac{\text{count}(\triangle \text{ and } +1)}{\text{count}(\triangle)} = \frac{2}{3} = 0.67$$
$$P(-1|\circ) = \frac{\text{count}(\circ \text{ and } -1)}{\text{count}(\circ)} = \frac{1}{2} = 0.5 \text{ vs. } P(+1|\circ) = \frac{\text{count}(\circ \text{ and } +1)}{\text{count}(\circ)} = \frac{1}{2} = 0.5$$

(Credits: Ryan McDonald)

Machine Learning

- 1 Define a model/distribution of interest
 - 2 Make some assumptions if needed
 - 3 Fit the model to the data
- Model: $P(\text{label}|\text{sequence}) = P(\text{label}|\text{symbol}_1, \dots, \text{symbol}_n)$
 - Prediction for new sequence = $\text{argmax}_{\text{label}} P(\text{label}|\text{sequence})$
 - Assumption (**naive Bayes**—more later):

$$P(\text{symbol}_1, \dots, \text{symbol}_n | \text{label}) = \prod_{i=1}^n P(\text{symbol}_i | \text{label})$$

- Fit the model to the data: count!! (simple probabilistic modeling)

Some Notation: Inputs and Outputs

- Input $x \in \mathcal{X}$
 - e.g., a news article, a sentence, an image, ...
- Output $y \in \mathcal{Y}$
 - e.g., spam/not spam, a topic, the object in the image (cat? dog?); a segmentation of the image (pedestrian; car; grass; background)
- Input/Output pair: $(x, y) \in \mathcal{X} \times \mathcal{Y}$
 - e.g., a **news article** together with a **topic**
 - e.g., a **image** together with an **object**
 - e.g., an **image** partitioned into **segmentation regions**

Many Flavours

- **Supervised learning:** pairs (x, y) are provided at training time (the main focus of this class)
 - Examples: perceptron, SVMs, decision trees, nearest neighbor, neural networks, ...
 - Caveat: the labels y may be hard or expensive to annotate
- **Unsupervised learning:** only x is provided; the model needs to figure out what the labels are without any supervision
 - Examples: clustering, PCA, ...
- **Reinforcement learning:** x is provided, and the model can act on the environment to obtain a reward (but doesn't get to know y)
 - Example: a robot acting on an environment trying to achieve a goal
- **Active learning:** the model requests which data points to label next.

Supervised Learning

- We are given a **labeled dataset** of input/output pairs:

$$\mathcal{D} = \{(x_n, y_n)\}_{n=1}^N \subseteq \mathcal{X} \times \mathcal{Y}$$

- **Goal:** use it to learn a **predictor** $h : \mathcal{X} \rightarrow \mathcal{Y}$ that generalizes well to arbitrary inputs.
- At test time, given $x \in \mathcal{X}$, we predict

$$\hat{y} = h(x).$$

- Hopefully, $\hat{y} \approx y$ most of the time.

Things can go by different names depending on what \mathcal{Y} is...

Regression

Deals with **continuous** output variables:

- **Regression:** $y = \mathbb{R}$
 - e.g., given a news article, how much time a user will spend reading it?
- **Multivariate regression:** $y = \mathbb{R}^K$
 - e.g., predict the X-Y coordinates in an image where the user will click

Classification

Deals with **discrete** output variables:

- **Binary classification:** $\mathcal{Y} = \{\pm 1\}$
 - e.g., spam detection
- **Multi-class classification:** $\mathcal{Y} = \{1, 2, \dots, K\}$
 - e.g., topic classification
- **Structured classification:** \mathcal{Y} exponentially large and structured
 - e.g., machine translation, caption generation, image segmentation
 - **Later in this course!**

Sometimes **reductions** are convenient:

- logistic regression reduces classification to regression
- one-vs-all reduces multi-class to binary
- greedy search reduces structured classification to multi-class

... but other times it's better to tackle the problem in its native form.

More later!

Probabilistic Models

- Let's assume our goal is to model the conditional probability of output labels y given inputs x , i.e. $P(y|x)$
- If we can define this distribution, then classification becomes:

$$\hat{y} = \arg \max_{y \in \mathcal{Y}} P(y|x)$$

Bayes Rule

- One way to model $P(y|x)$ is through **Bayes Rule**:

$$P(y|x) = \frac{P(y)P(x|y)}{P(x)}$$

$$\arg \max_y P(y|x) = \arg \max_y P(y)P(x|y)$$

(since x is fixed!)

- $P(y)P(x|y) = P(x, y)$: a joint probability
- Above is a “generative story”: ‘pick y ; then pick x given y .’”
- Models that consider the joint $P(x, y)$ are called **generative models**, because they come with a generative story.

Naive Bayes

Assume that an input x is partitioned as x_1, \dots, x_L , where $x_k \in \mathcal{X}_k$

Example:

- x is a document of length L
- x_k is the k^{th} token (a word)
- The set $\mathcal{X}_k = \mathcal{V}$ is a fixed vocabulary (all tokens drawn from \mathcal{V})

Naive Bayes Assumption

(conditional independence)

$$P(\underbrace{x_1, \dots, x_L}_x | y) = \prod_{k=1}^L P(x_k | y)$$

Multinomial Naive Bayes

$$P(x, y) = P(y) \underbrace{P(x_1, \dots, x_L | y)}_x = P(y) \prod_{k=1}^L P(x_k | y)$$

- All tokens are conditionally independent, given the topic
- The word order doesn't change $P(x, y)$ (bag-of-words assumption)

Small caveat: we assumed that the document has a fixed length L .

This is not realistic.

How to deal with variable length?

Multinomial Naive Bayes – Arbitrary Length

Solution: introduce a distribution over document length $P(|x|)$

- e.g. a Poisson distribution.

We get:

$$P(x, y) = P(y) \underbrace{P(|x|) \prod_{k=1}^{|x|} P(x_k|y)}_{P(x|y)}$$

$P(|x|)$ is constant (independent of y), so nothing really changes

- the posterior $P(y|x)$ is the same as before.

What Does This Buy Us?

$$P(\underbrace{x_1, \dots, x_L}_x | y) = \prod_{k=1}^L P(x_k | y)$$

What do we gain with the Naive Bayes assumption?

- A huge reduction in the number of parameters!
- If we haven't done any factorization assumption, how many parameters would be required for expressing $P(x_1, \dots, x_L | y)$? $O(|\mathcal{V}|^L)$
- And how many parameters with Naive Bayes? $O(|\mathcal{V}|)$

Less parameters \implies Less computation; less risk of overfitting (more later)

Naive Bayes – Learning

$$P(y)P(\underbrace{x_1, \dots, x_L}_x | y) = P(y) \prod_{k=1}^L P(x_k | y)$$

- Input: dataset $\mathcal{D} = \{(x^{(t)}, y^{(t)})\}_{t=1}^N$ (examples assumed i.i.d.)
- Parameters $\Theta = \{P(y), P(v|y)\}$ for $v \in \mathcal{V}$
- **Objective: Maximum Likelihood Estimation (MLE):** choose parameters that maximize the likelihood of observed data

$$\mathcal{L}(\Theta; \mathcal{D}) = \prod_{t=1}^N P(x^{(t)}, y^{(t)}) = \prod_{t=1}^N \left(P(y^{(t)}) \prod_{k=1}^L P(x_k^{(t)} | y^{(t)}) \right)$$

$$\hat{\Theta} = \arg \max_{\Theta} \prod_{t=1}^N \left(P(y^{(t)}) \prod_{k=1}^L P(x_k^{(t)} | y^{(t)}) \right)$$

Naive Bayes – Learning via MLE

For the multinomial Naive Bayes model, MLE has a **closed form solution!!**

It all boils down to counting and normalizing!!

(The proof is left as an exercise...)

Naive Bayes – Learning via MLE

$$\hat{\Theta} = \arg \max_{\Theta} \prod_{t=1}^N \left(P(y^{(t)}) \prod_{k=1}^L P(x_k^{(t)} | y^{(t)}) \right)$$

$$\hat{P}(y) = \frac{\sum_{t=1}^N [[y^{(t)} = y]]}{N}$$

$$\hat{P}(v|y) = \frac{\sum_{t=1}^N \sum_{k=1}^L [[x_k^{(t)} = v \text{ and } y^{(t)} = y]]}{L \sum_{t=1}^N [[y^{(t)} = y]]}$$

$[[X]]$ is 1 if property X holds, 0 otherwise (Iverson notation)
Fraction of times a feature appears in training cases of a given label

Naive Bayes Example

- Corpus of movie reviews: 7 examples for **training**

Doc	Words	Class
1	Great movie, excellent plot, renown actors	Positive
2	I had not seen a fantastic plot like this in good 5 years. Amazing!!!	Positive
3	Lovely plot, amazing cast, somehow I am in love with the bad guy	Positive
4	Bad movie with great cast, but very poor plot and unimaginative ending	Negative
5	I hate this film, it has nothing original	Negative
6	Great movie, but not...	Negative
7	Very bad movie, I have no words to express how I dislike it	Negative

Naive Bayes Example

- **Features:** adjectives (bag-of-words)

Doc	Words	Class
1	Great movie, excellent plot, renowned actors	Positive
2	I had not seen a fantastic plot like this in good 5 years. amazing !!!	Positive
3	Lovely plot, amazing cast, somehow I am in love with the bad guy	Positive
4	Bad movie with great cast, but very poor plot and unimaginative ending	Negative
5	I hate this film, it has nothing original. Really bad	Negative
6	Great movie, but not...	Negative
7	Very bad movie, I have no words to express how I dislike it	Negative

Naive Bayes Example

Relative frequency:

Priors:

$$P(\text{positive}) = \frac{\sum_{t=1}^N [[y^{(t)} = \text{positive}]]}{N} = 3/7 = 0.43$$

$$P(\text{negative}) = \frac{\sum_{t=1}^N [[y^{(t)} = \text{negative}]]}{N} = 4/7 = 0.57$$

Assume standard pre-processing: tokenization, lowercasing, punctuation removal (except special punctuation like !!!)

Naive Bayes Example

Likelihoods: Count adjective v in class y / adjectives in y

$$\hat{P}(v|y) = \frac{\sum_{t=1}^N \sum_{k=1}^L [[x_k^{(t)} = v \text{ and } y^{(t)} = y]]}{L \sum_{t=1}^N [[y^{(t)} = y]]}$$

$P(\textit{amazing} \textit{positive})$	$= 2/10$	$P(\textit{amazing} \textit{negative})$	$= 0/8$
$P(\textit{bad} \textit{positive})$	$= 1/10$	$P(\textit{bad} \textit{negative})$	$= 3/8$
$P(\textit{excellent} \textit{positive})$	$= 1/10$	$P(\textit{excellent} \textit{negative})$	$= 0/8$
$P(\textit{fantastic} \textit{positive})$	$= 1/10$	$P(\textit{fantastic} \textit{negative})$	$= 0/8$
$P(\textit{good} \textit{positive})$	$= 1/10$	$P(\textit{good} \textit{negative})$	$= 0/8$
$P(\textit{great} \textit{positive})$	$= 1/10$	$P(\textit{great} \textit{negative})$	$= 2/8$
$P(\textit{lovely} \textit{positive})$	$= 1/10$	$P(\textit{lovely} \textit{negative})$	$= 0/8$
$P(\textit{original} \textit{positive})$	$= 0/10$	$P(\textit{original} \textit{negative})$	$= 1/8$
$P(\textit{poor} \textit{positive})$	$= 0/10$	$P(\textit{poor} \textit{negative})$	$= 1/8$
$P(\textit{renowned} \textit{positive})$	$= 1/10$	$P(\textit{renowned} \textit{negative})$	$= 0/8$
$P(\textit{unimaginative} \textit{positive})$	$= 0/10$	$P(\textit{unimaginative} \textit{negative})$	$= 1/8$

Naive Bayes Example

Given a new segment to classify (**test time**):

Doc	Words	Class
8	This was a fantastic story, good , lovely	???

Final decision

$$\hat{y} = \arg \max_y \left(P(y) \prod_{k=1}^L P(x_k|y) \right)$$

$$P(\text{positive}) * P(\text{fantastic}|\text{positive}) * P(\text{good}|\text{positive}) * P(\text{lovely}|\text{positive})$$

$$3/7 * 1/10 * 1/10 * 1/10 = 0.00043$$

$$P(\text{negative}) * P(\text{fantastic}|\text{negative}) * P(\text{good}|\text{negative}) * P(\text{lovely}|\text{negative})$$

$$4/7 * 0/8 * 0/8 * 0/8 = 0$$

So: *sentiment = positive*

Naive Bayes Example

Given a new segment to classify (**test time**):

Doc	Words	Class
9	Great plot, great cast, great everything	???

Final decision

$$P(\text{positive}) * P(\text{great}|\text{positive}) * P(\text{great}|\text{positive}) * P(\text{great}|\text{positive})$$

$$3/7 * 1/10 * 1/10 * 1/10 = 0.00043$$

$$P(\text{negative}) * P(\text{great}|\text{negative}) * P(\text{great}|\text{negative}) * P(\text{great}|\text{negative})$$

$$4/7 * 2/8 * 2/8 * 2/8 = 0.00893$$

So: *sentiment = negative*

Naive Bayes Example

But if the new segment to classify (**test time**) is:

Doc	Words	Class
10	Boring movie, annoying plot, unimaginative ending	???

Final decision

$$P(\text{positive}) * P(\text{boring}|\text{positive}) * P(\text{annoying}|\text{positive}) * P(\text{unimaginative}|\text{positive})$$

$$3/7 * 0/10 * 0/10 * 0/10 = 0$$

$$P(\text{negative}) * P(\text{boring}|\text{negative}) * P(\text{annoying}|\text{negative}) * P(\text{unimaginative}|\text{negative})$$

$$4/7 * 0/8 * 0/8 * 1/8 = 0$$

So: *sentiment* = ???

Laplace Smoothing

Add smoothing to feature counts (add 1 to every count):

$$\hat{P}(v|y) = \frac{\sum_{t=1}^N \sum_{k=1}^L [[x_k^{(t)} = v \text{ and } y^{(t)} = y]] + 1}{L \sum_{t=1}^N [[y^{(t)} = y]] + |\mathcal{V}|}$$

where $|\mathcal{V}|$ = number of distinct adjectives in training (all classes) = **12**

Doc	Words	Class
11	Boring movie, annoying plot, unimaginative ending	???

Final decision

$$P(\text{positive}) * P(\text{boring}|\text{positive}) * P(\text{annoying}|\text{positive}) * P(\text{unimaginative}|\text{positive})$$

$$3/7 * ((0 + 1)/(10 + 12)) * ((0 + 1)/(10 + 12)) * ((0 + 1)/(10 + 12)) = 0.000040$$

$$P(\text{negative}) * P(\text{boring}|\text{negative}) * P(\text{annoying}|\text{negative}) * P(\text{unimaginative}|\text{negative})$$

$$4/7 * ((0 + 1)/(8 + 12)) * ((0 + 1)/(8 + 12)) * ((1 + 1)/(8 + 12)) = 0.000143$$

So: *sentiment = negative*

Conclusions

- Machine learning allows computer programs to learn from data (observations, interventions, ...) and improve performance with experience
- Can be supervised, unsupervised, self-supervised, reinforced, etc.
- Tasks can be (binary or multi-class) classification, regression, or more nuanced
- The naive Bayes classifier is a very simple probabilistic model that assumes inputs are conditionally independent given the label, and uses the Bayes rule to make predictions
- Learning a naive Bayes classifier amounts to counting and normalizing.

Lecture 3: Linear Models I

André Martins, Francisco Melo, Mário Figueiredo



Deep Learning Course, Winter 2021-2022

Today's Roadmap

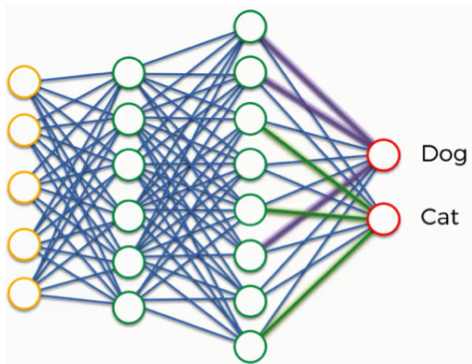
- **Linear** regression
- Binary and multi-class **linear** classification
- **Linear** classifiers: perceptron

Why Linear Classifiers?

We know the course title promised “**deep**”, but...

- The underlying machine learning concepts are the same
- The theory (statistics and optimization) are much better understood
- Linear classifiers still widely used (very effective when data is scarce)
- Linear classifiers are **a component of neural networks**.

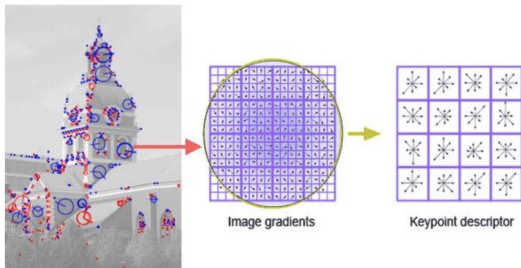
Linear Classifiers and Neural Networks



Feature Representations

Feature engineering is an important step in linear classifiers:

- Bag-of-words **features** for text, also lemmas, parts-of-speech, ...
- SIFT **features** and wavelet representations in computer vision



- Other categorical, Boolean, and continuous **features**

Feature Representations

Representing information about x :

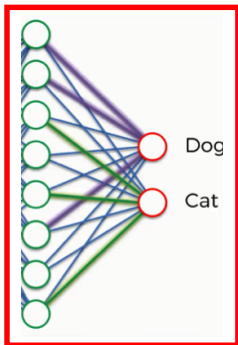
Typical approach: define a feature map $\phi : \mathcal{X} \rightarrow \mathbb{R}^D$

- $\phi(x)$ is a (maybe high-dimensional) **feature vector**
- $\phi(x)$ may include **Boolean**, **categorical**, and **continuous** features
- Categorical features can be reduced to a one-hot binary vector.

Example: Continuous Features



**Handcrafted
Features**



Linear Classifier

Feature Engineering and NLP Pipelines

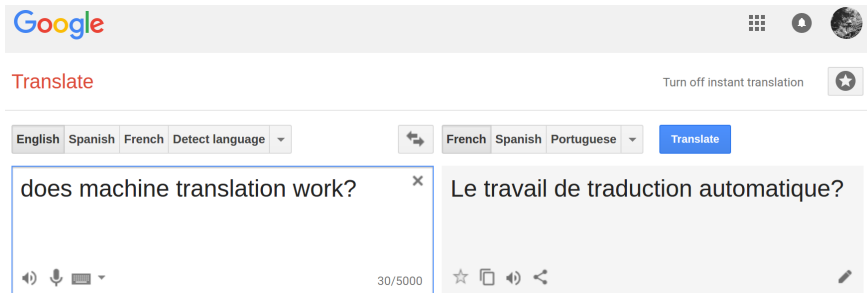
Classical NLP pipelines: stacking together several linear classifiers

Each classifier's output is used to handcraft features for subsequent classifiers

Examples of features:

- **Word occurrences**: binary feature (word occurs or not in a document)
- **Word counts**: numeric feature counting how many times a word occurs
- **POS tags**: classifying words as noun, verb, adjective, ...
- **Spell check**: misspellings counts for spam detection

Example: Translation Quality Estimation



The screenshot shows the Google Translate web interface. At the top, the Google logo is on the left, and navigation icons (grid, bell, globe) are on the right. Below the logo, the word "Translate" is displayed in red. To the right of "Translate" is a link "Turn off instant translation" and a star icon. The main interface features two language selection dropdowns: the first is set to "English" and the second to "French". A blue "Translate" button is positioned to the right of the second dropdown. Below the dropdowns, the input text "does machine translation work?" is shown in a light gray box. The output text "Le travail de traduction automatique?" is shown in a light gray box to the right. At the bottom of the input box, there are icons for speaker, microphone, and a dropdown arrow, along with the character count "30/5000". At the bottom of the output box, there are icons for star, copy, speaker, and share, along with a pencil icon for editing.

Goal: estimate the quality of a translation on the fly (without a reference)!

Example: Translation Quality Estimation

Hand-crafted features:

- no of tokens in the source/target segment
- LM probability of source/target segment and their ratio
- % of source 1–3-grams observed in 4 frequency quartiles of source corpus
- average no of translations per source word
- ratio of brackets and punctuation symbols in source & target segments
- ratio of numbers, content/non-content words in source & target segments
- ratio of nouns/verbs/etc in the source & target segments
- % of dependency relations b/w constituents in source & target segments
- diff in depth of the syntactic trees of source & target segments
- diff in no of PP/NP/VP/ADJP/ADVP/CONJP in source & target
- diff in no of person/location/organization entities in source & target
- features and global score of the SMT system
- number of distinct hypotheses in the n-best list
- 1–3-gram LM probabilities using translations in the n-best to train the LM
- average size of the target phrases
- proportion of pruned search graph nodes;
- proportion of recombined graph nodes.

Representation Learning

Feature engineering is an **art** and can be very time-consuming

...but it's a good way of encoding **prior knowledge**, still widely used in practice (especially with “small data”)

Alternative to feature engineering: **representation learning**

Neural networks will alleviate this (later in the course)!

Regression

Output space \mathcal{Y} is continuous

Example: given an article, how much time a user spends reading it?

Summer Schools and Machine Learning. A beautiful love story!



Mohan Acharya [Follow](#)
Jan 7, 2019 7 min read

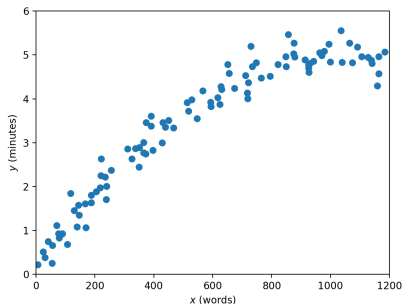


- x is number of words of the article
- y is the reading time (minutes)

How to define a model that predicts \hat{y} from x ?

Linear Regression

- Model: assume $\hat{y} = wx + b$
- Model parameters: w and b
- Given training data $\mathcal{D} = \{(x_n, y_n)\}_{n=1}^N$, how to estimate w and b ?



Least squares method: fit w and b on the training set by solving

$$\min_{w, b} \sum_{n=1}^N (y_n - (w x_n + b))^2$$

Linear Regression

Often, linear dependency of \hat{y} on x is a bad assumption

Second model: assume $\hat{y} = \mathbf{w} \cdot \phi(x)$, where $\phi(x)$ is a feature vector

- e.g. $\phi(x) = [1, x, x^2, \dots, x^D]$ (polynomial features degree $\leq D$)
- the bias term b is captured by the constant feature $\phi_0(x) = 1$

Fit \mathbf{w} by minimizing $\sum_n (y_n - \mathbf{w} \cdot \phi(x_n))^2$

- Closed-form solution:

$$\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \text{ with } \mathbf{X} = \begin{bmatrix} \phi(x_1)^\top \\ \vdots \\ \phi(x_n)^\top \\ \vdots \\ \phi(x_N)^\top \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \\ \vdots \\ y_N \end{bmatrix}.$$

Still called **linear regression** – linearity w.r.t. the model parameters \mathbf{w} .

One-Slide Proof

Fit \mathbf{w} by minimizing

$$\sum_{n=1}^N (y_n - \mathbf{w} \cdot \phi(x_n))^2 = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

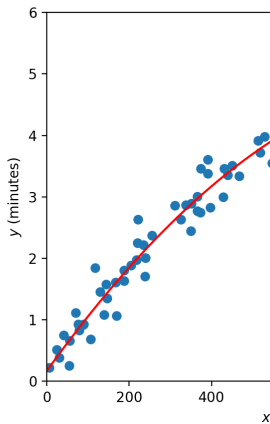
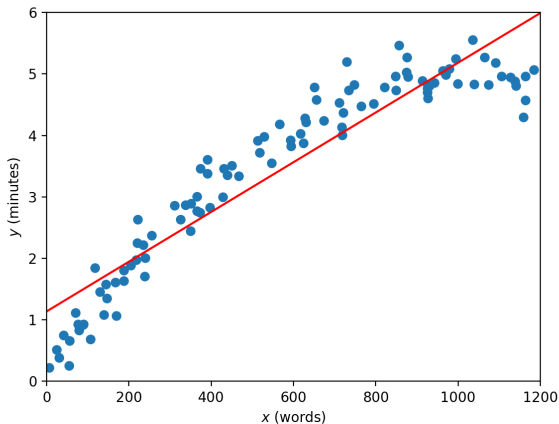
Equate the gradient to zero and solve the resulting equation:

$$\begin{aligned} 0 &= \nabla_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 \\ &= \nabla_{\mathbf{w}} \left(\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \|\mathbf{y}\|^2 \right) \\ &= 2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{y} \end{aligned}$$

Therefore

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

Linear Regression ($D = 1$)($D = 2$)



Squared Loss Function

Linear regression with least squares criterion corresponds to a loss function

$$L(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2, \quad \text{where } \hat{y} = \mathbf{w} \cdot \phi(x).$$

The model is fit to the training data by minimizing this loss function.

This is called the **squared loss**.

More later.

Least Squares (LS) – Probabilistic Interpretation

The least squares method has a **probabilistic interpretation**.

Assume the data is generated stochastically as

$$y_i = \mathbf{w}^* \cdot \phi(x_i) + n$$

where $n \sim \mathcal{N}(0, \sigma^2)$ is Gaussian noise (with σ^2 fixed), and \mathbf{w}^* are the “true” model parameters.

That is, $y_i \sim \mathcal{N}(\mathbf{w}^* \cdot \phi(x_i), \sigma^2)$.

Then \mathbf{w} given by LS is the **maximum likelihood estimate** under this model.

One-Slide Proof

$$\text{Recall } \mathcal{N}(y; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y-\mu)^2}{2\sigma^2}\right).$$

$$\begin{aligned}\hat{\mathbf{w}}_{\text{MLE}} &= \arg \max_{\mathbf{w}} \prod_{i=1}^N P(y_i | x_i; \mathbf{w}) \\ &= \arg \max_{\mathbf{w}} \sum_{i=1}^N \log P(y_i | x_i; \mathbf{w}) \\ &= \arg \max_{\mathbf{w}} \sum_{i=1}^N -\frac{(y_i - \mathbf{w} \cdot \phi(x_i))^2}{2\sigma^2} - \underbrace{\log(\sqrt{2\pi\sigma^2})}_{\text{constant}} \\ &= \arg \min_{\mathbf{w}} \sum_{i=1}^N (y_i - \mathbf{w} \cdot \phi(x_i))^2\end{aligned}$$

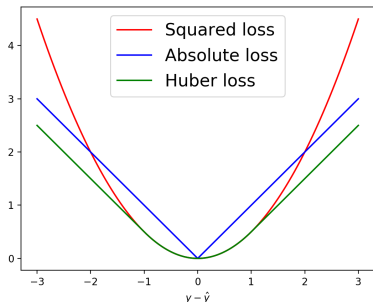
Thus, linear regression with the squared loss = MLE under Gaussian noise.

Other Regression Losses

Squared loss: $L(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$.

Absolute error loss: $L(y, \hat{y}) = |y - \hat{y}|$.

Huber loss: $L(y, \hat{y}) = \begin{cases} \frac{1}{2}(y - \hat{y})^2 & \text{if } |y - \hat{y}| \leq 1 \\ |y - \hat{y}| - \frac{1}{2} & \text{if } |y - \hat{y}| \geq 1. \end{cases}$

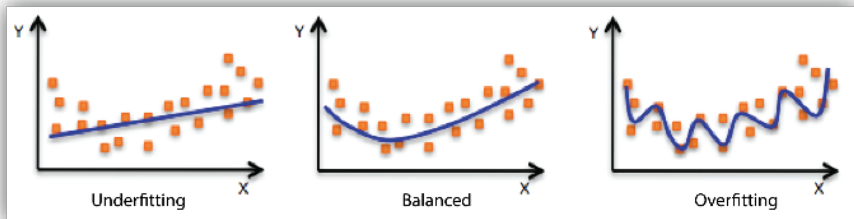


Quiz: which of these are convex; and strictly convex?

Overfitting and Underfitting

We saw earlier an example of **underfitting** (slide 15)

However, if the model is too complex (too many parameters) and/or the data is scarce, we run the risk of **overfitting**



To avoid overfitting, we often need **regularization** (more later)

Maximum A Posteriori

Assuming we have a prior distribution $\mathbf{w} \sim \mathcal{N}(0, \tau^2 \mathbf{I})$

A criterion to estimate \mathbf{w}^* is **maximum a posteriori** (MAP):

$$\begin{aligned}\hat{\mathbf{w}}_{\text{MAP}} &= \arg \max_{\mathbf{w}} P(\mathbf{w}) \prod_{i=1}^N P(y_i | x_i; \mathbf{w}) \\ &= \arg \max_{\mathbf{w}} \log P(\mathbf{w}) + \sum_{i=1}^N \log P(y_i | x_i; \mathbf{w}) \\ &= \arg \max_{\mathbf{w}} -\frac{\|\mathbf{w}\|^2}{2\tau^2} - \sum_{i=1}^N -\frac{(y_i - \mathbf{w} \cdot \phi(x_i))^2}{2\sigma^2} + \text{constant} \\ &= \arg \min_{\mathbf{w}} \underbrace{\lambda \|\mathbf{w}\|^2}_{\text{regularizer}} + \underbrace{\sum_{i=1}^N (y_i - \mathbf{w} \cdot \phi(x_i))^2}_{\text{loss}}\end{aligned}$$

where $\lambda = \sigma^2/\tau^2$ (so-called **regularization constant**)

Thus, ℓ_2 -regularization is equivalent to MAP with a Gaussian prior.

Binary Classification

Before multi-class, we start with simpler case of **binary classification**

Output set $\mathcal{Y} = \{-1, +1\}$

Example: Given a news article, is it true or fake?

- x is the news article, represented a **feature vector** $\phi(x)$
- y can be either $+1$ (**true**) or -1 (**fake**)

How to define a model to predict \hat{y} from x ?

Linear Classifier

Defined by

$$\hat{y} = \text{sign}(\mathbf{w} \cdot \phi(x) + b) = \begin{cases} +1 & \text{if } \mathbf{w} \cdot \phi(x) + b \geq 0 \\ -1 & \text{if } \mathbf{w} \cdot \phi(x) + b < 0. \end{cases}$$

Intuitively, $\mathbf{w} \cdot \phi(x) + b$ is a “score” for the positive class

Different from regression: **sign function** converts from continuous to binary

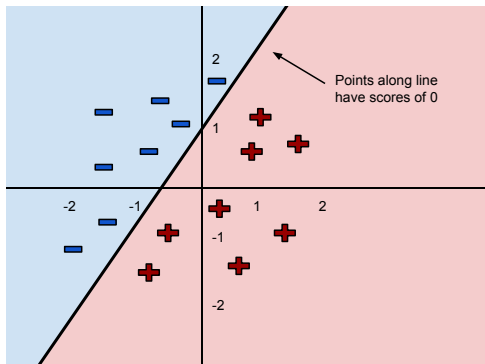
The decision boundary is a **hyperplane** (w.r.t. $\phi(x)$, not x)

$$\mathbf{w} \cdot \phi(x) + b = 0$$

Also called a “hyperplane classifier.”

Linear Classifier

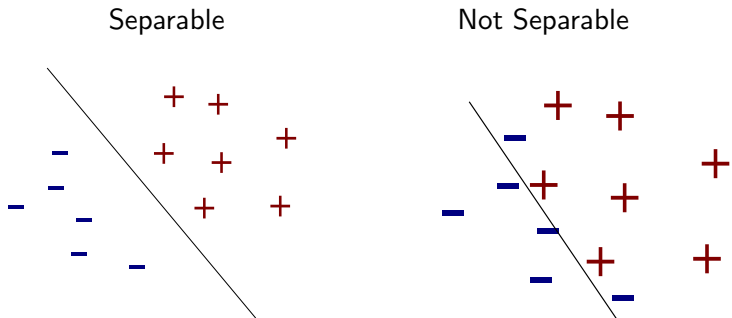
(w, b) defines a **hyperplane** that splits the space into two half spaces:



How to learn this hyperplane from the training data $\mathcal{D} = \{(x_n, y_n)\}_{n=1}^N$?

Linear Separability

- A dataset \mathcal{D} is **linearly separable** if there exists (w, b) such that classification is perfect



We next present an algorithm that finds such an hyperplane if it exists!

Linear Classifier: No Bias Term

It is common to present linear classifiers without the bias term b :

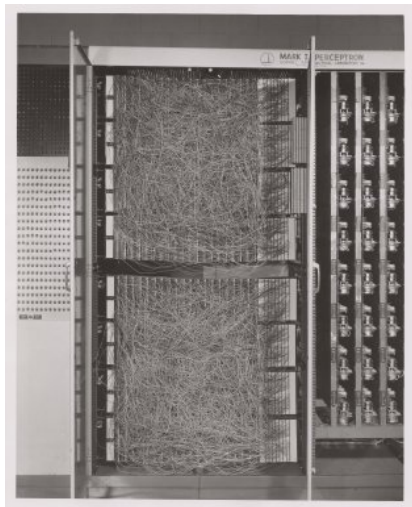
$$\hat{y} = \text{sign}(\mathbf{w} \cdot \phi(\mathbf{x}) + b)$$

In this case, the decision hyperplane passes through the origin

We can always do this without loss of generality:

- Add a constant feature to $\phi(\mathbf{x})$: $\phi_0(\mathbf{x}) = 1$
- The corresponding weight w_0 replaces the bias term b

Perceptron (Rosenblatt, 1958)



(Source: Wikipedia)

- Invented in 1957 at the Cornell Aeronautical Laboratory by [Frank Rosenblatt](#)
- Implemented in custom-built hardware as the “Mark 1 perceptron,” for image recognition
- 400 photocells, randomly connected to the “neurons.” Weights were encoded in potentiometers
- Weight updates during learning were performed by electric motors.

Perceptron in the News...

NEW NAVY DEVICE LEARNS BY DOING

Psychologist Shows Embryo
of Computer Designed to
Read and Grow Wiser

WASHINGTON, July 7 (UPI)—The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence.

The embryo—the Weather Bureau's \$2,000,000 "704" computer—learned to differentiate between right and left after fifty attempts in the Navy's demonstration for newsmen.

The service said it would use this principle to build the first of its Perceptron thinking machines that will be able to read and write. It is expected to be finished in about a year at a cost of \$100,000.

Dr. Frank Rosenblatt, designer of the Perceptron, conducted the demonstration. He said the machine would be the first device to think as the human brain. As do human be-

ings, Perceptron will make mistakes at first, but will grow wiser as it gains experience, he said.

Dr. Rosenblatt, a research psychologist at the Cornell Aeronautical Laboratory, Buffalo, said Perceptrons might be fired to the planets as mechanical space explorers.

Without Human Controls

The Navy said the perceptron would be the first non-living mechanism "capable of receiving, recognizing and identifying its surroundings without any human training or control."

The "brain" is designed to remember images and information it has perceived itself. Ordinary computers remember only what is fed into them on punch cards or magnetic tape.

Later Perceptrons will be able to recognize people and call out their names and instantly translate speech in one language to speech or writing in another language, it was predicted.

Mr. Rosenblatt said in principle it would be possible to build brains that could reproduce themselves on an assembly line and which would be conscious of their existence.

1958 New York Times...

In today's demonstration, the "704" was fed two cards, one with squares marked on the left side and the other with squares on the right side.

Learns by Doing

In the first fifty trials, the machine made no distinction between them. It then started registering a "Q" for the left squares and "O" for the right squares.

Dr. Rosenblatt said he could explain why the machine learned only in highly technical terms. But he said the computer had undergone a "self-induced change in the wiring diagram."

The first Perceptron will have about 1,000 electronic "association cells" receiving electrical impulses from an eye-like scanning device with 400 photo-cells. The human brain has 10,000,000,000 responsive cells, including 100,000,000 connections with the eyes.

NEW NAVY DEVICE LEARNS BY DOING

Psychologist Shows Embryo
of Computer Designed to
Read and Grow Wiser

WASHINGTON, July 7 (UPI)—The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence.

The embryo—the Weather Bureau's \$2,000,000 "704" computer—learned to differentiate between right and left after fifty attempts in the Navy's demonstration for newsmen.

The service said it would use this principle to build the first of its Perceptron thinking machines that will be able to read and write. It is expected to be finished in about a year at a cost of \$100,000.

Dr. Frank Rosenblatt, designer of the Perceptron, conducted the demonstration. He said the machine would be the first device to think as the human brain. As do human be-

Perceptron Algorithm

Online algorithm: process one data point x_i at each round

- 1 Take x_i ; apply the current model to make the corresponding prediction
- 2 If prediction is **correct**, do nothing
- 3 If it is **wrong**, correct w by adding/subtracting feature vector $\phi(x_i)$

Omit the bias b : use a constant feature $\phi_0(x) = 1$, as explained above.

Perceptron Algorithm

```
input: labeled data  $\mathcal{D}$   
initialize  $\mathbf{w}^{(0)} = \mathbf{0}$   
initialize  $k = 0$  (number of mistakes)  
repeat  
  get new training example  $(x_i, y_i)$   
  predict  $\hat{y}_i = \text{sign}(\mathbf{w}^{(k)} \cdot \phi(x_i))$   
  if  $\hat{y}_i \neq y_i$  then  
    update  $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + y_i \phi(x_i)$   
    increment  $k$   
  end if  
until maximum number of epochs  
output: model weights  $\mathbf{w}^{(k)}$ 
```

Perceptron's Mistake Bound

Definitions:

- the training data is **linearly separable** with **margin** $\gamma > 0$ iff there is a weight vector \mathbf{u} with $\|\mathbf{u}\| = 1$ such that

$$y_i \mathbf{u} \cdot \phi(x_i) \geq \gamma, \quad \forall i.$$

- radius** of the data: $R = \max_i \|\phi(x_i)\|$.

Then we have the following bound of the **number of mistakes**:

Theorem (Novikoff (1962))

The perceptron algorithm is guaranteed to find a separating hyperplane after at most $\frac{R^2}{\gamma^2}$ mistakes.

One-Slide Proof

Recall that $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + y_i \phi(x_i)$.

- **Lower bound on $\|\mathbf{w}^{(k+1)}\|$:**

$$\begin{aligned} \mathbf{u} \cdot \mathbf{w}^{(k+1)} &= \mathbf{u} \cdot \mathbf{w}^{(k)} + y_i \mathbf{u} \cdot \phi(x_i) \\ &\geq \mathbf{u} \cdot \mathbf{w}^{(k)} + \gamma \\ &\geq \mathbf{u} \cdot \mathbf{w}^{(k-1)} + \gamma + \gamma \\ &\geq k \gamma \end{aligned}$$

Thus, $\|\mathbf{w}^{(k+1)}\| = \|\mathbf{u}\| \cdot \|\mathbf{w}^{(k+1)}\| \geq \mathbf{u} \cdot \mathbf{w}^{(k+1)} \geq k\gamma$ (Cauchy-Schwarz)

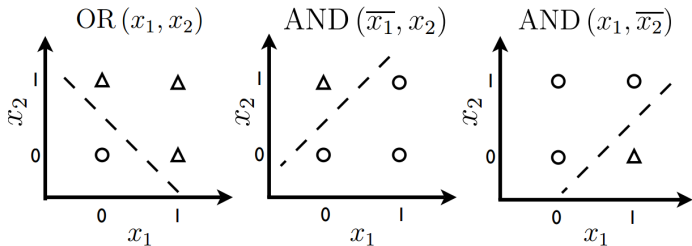
- **Upper bound on $\|\mathbf{w}^{(k+1)}\|$:**

$$\begin{aligned} \|\mathbf{w}^{(k+1)}\|^2 &= \|\mathbf{w}^{(k)}\|^2 + \|\phi(x_i)\|^2 + 2y_i \mathbf{w}^{(k)} \cdot \phi(x_i) \\ &\leq \|\mathbf{w}^{(k)}\|^2 + R^2 \\ &\leq k R^2. \end{aligned}$$

Equating both sides, we get $(k\gamma)^2 \leq kR^2 \Rightarrow k \leq R^2/\gamma^2$ ■

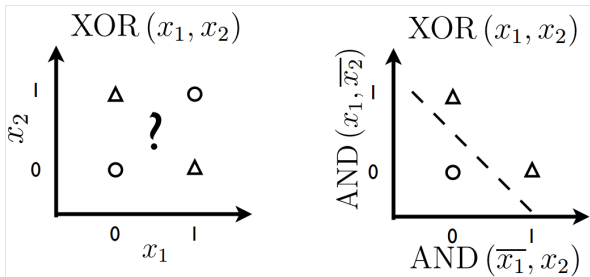
What a Simple Perceptron Can and Can't Do

- Remember: the decision boundary is linear (**linear classifier**)
- It **can** solve linearly separable problems (OR, AND)



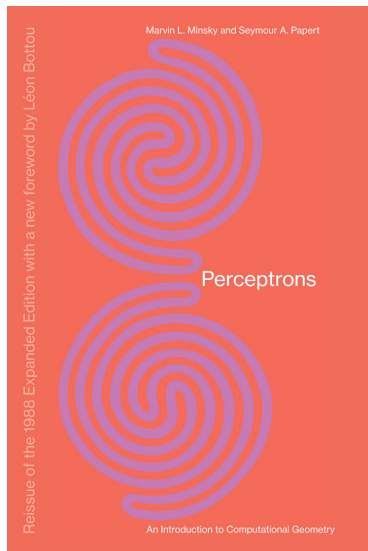
What a Simple Perceptron Can and Can't Do

- ... but it **can't** solve **non-linearly separable** problems such as simple XOR (unless input is transformed into a better representation):



- This result is often attributed to Minsky and Papert (1969) but was known well before.

Limitations of the Perceptron



Minsky and Papert (1969):

- Shows limitations of multi-layer perceptrons and fostered an “AI winter” period.

More later in the neural networks’ lecture!

Multi-Class Classification

We now assume a **multi-class** problem, with $|\mathcal{Y}| \geq 2$ labels (classes).

Reduction to Binary Classification

Several strategies:

- **One-vs-all**: train one binary classifier per class, using all others as negative examples; pick the class with the highest score.
- **One-vs-one**: another strategy is to train pairwise classifiers and use majority voting.
- **Binary coding**: use binary code for the class labels and learn a binary classifier for each bit.

Here, we will consider classifiers that tackle the multiple classes directly.

Multi-Class Linear Classifiers

- Parametrized by a **weight matrix** $\mathbf{W} \in \mathbb{R}^{|\mathcal{Y}| \times D}$ (one weight per feature/label pair) and a **bias vector** $\mathbf{b} \in \mathbb{R}^{|\mathcal{Y}|}$:

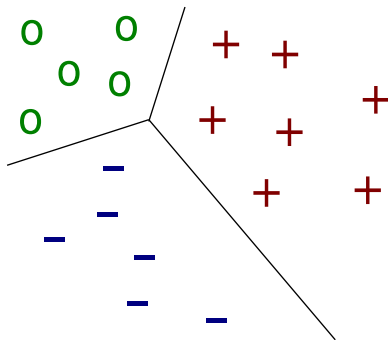
$$\mathbf{W} = \begin{bmatrix} \mathbf{w}_1^\top \\ \vdots \\ \mathbf{w}_y^\top \\ \vdots \\ \mathbf{w}_{|\mathcal{Y}|}^\top \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_y \\ \vdots \\ b_{|\mathcal{Y}|} \end{bmatrix}.$$

- Equivalently, $|\mathcal{Y}|$ weight vectors $\mathbf{w}_y \in \mathbb{R}^D$ and scalars $b_y \in \mathbb{R}$
- The score of a particular label is based on a **linear** combination of features and their weights
- Predict the \hat{y} which maximizes this **score**:

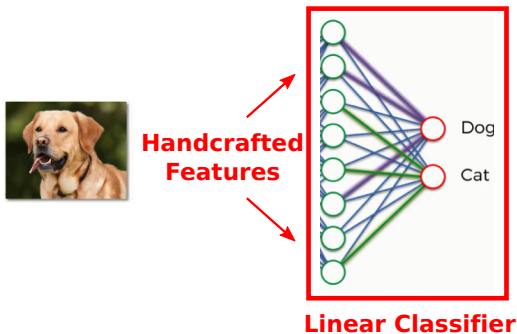
$$\hat{y} = \arg \max_{y \in \mathcal{Y}} \mathbf{w}_y \cdot \phi(x) + b_y.$$

Multi-Class Linear Classifier

Geometrically, (\mathbf{W}, \mathbf{b}) split the feature space into regions delimited by hyperplanes.



Commonly Used Notation in Neural Networks



$$\hat{y} = \operatorname{argmax}(\mathbf{W}\phi(x) + \mathbf{b}), \quad \mathbf{W} = \begin{bmatrix} \mathbf{w}_1^\top \\ \vdots \\ \mathbf{w}_y^\top \\ \vdots \\ \mathbf{w}_{|y|}^\top \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_y \\ \vdots \\ b_{|y|} \end{bmatrix}.$$

Multi-Class Recovers Binary

With **two classes** ($\mathcal{Y} = \{\pm 1\}$), this formulation recovers the binary classifier presented earlier:

$$\begin{aligned}\hat{y} &= \arg \max_{y \in \{\pm 1\}} \mathbf{w}_y \cdot \phi(x) + b_y \\ &= \begin{cases} +1 & \text{if } \mathbf{w}_{+1} \cdot \phi(x) + b_{+1} > \mathbf{w}_{-1} \cdot \phi(x) + b_{-1} \\ -1 & \text{otherwise} \end{cases} \\ &= \text{sign}\left(\underbrace{(\mathbf{w}_{+1} - \mathbf{w}_{-1})}_{\mathbf{w}} \cdot \phi(x) + \underbrace{(b_{+1} - b_{-1})}_{b}\right).\end{aligned}$$

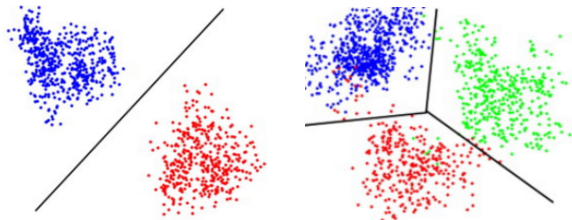
That is: only half of the parameters are needed.

Linear Classifiers (Binary vs Multi-Class)

- Prediction rule:

$$\hat{y} = h(x) = \arg \max_{y \in \mathcal{Y}} \overbrace{w_y \cdot \phi(x)}^{\text{linear in } w_y}$$

- The decision boundary is defined by the intersection of half spaces
- In the binary case ($|\mathcal{Y}| = 2$) this corresponds to a hyperplane classifier



Linear Classifier – No Bias Term

Again, it is common to omit the bias vector \mathbf{b} :

$$\hat{y} = \arg \max_{y \in \mathcal{Y}} \mathbf{w}_y \cdot \phi(x) + b_y$$

Like before, this can be done without loss of generality, by assuming a constant feature $\phi_0(x) = 1$

The first column of \mathbf{W} replaces the bias vector.

We assume this for simplicity.

Example: Perceptron

The perceptron algorithm also works for the multi-class case!

It has a similar mistake bound: if the data is separable, it's guaranteed to find separating hyperplanes!

Perceptron Algorithm: Multi-Class

input: labeled data \mathcal{D}

initialize $\mathbf{W}^{(0)} = \mathbf{0}$

initialize $k = 0$ (number of mistakes)

repeat

get new training example (x_i, y_i)

predict $\hat{y}_i = \arg \max_{y \in \mathcal{Y}} \mathbf{w}_y^{(k)} \cdot \phi(x_i)$

if $\hat{y}_i \neq y_i$ **then**

update $\mathbf{w}_{y_i}^{(k+1)} = \mathbf{w}_{y_i}^{(k)} + \phi(x_i)$ {increase weight of gold class}

update $\mathbf{w}_{\hat{y}_i}^{(k+1)} = \mathbf{w}_{\hat{y}_i}^{(k)} - \phi(x_i)$ {decrease weight of incorrect class}

increment k

end if

until maximum number of epochs

output: model weights $\mathbf{w}^{(k)}$

Conclusions

- Linear models involve manipulating weights and features
- Linear regression is a simple method for regression which has a closed form solution
- Linear classifiers include several well-known ML methods (both for binary and multi-class classification)
- Today we saw the **perceptron** and proved a mistake bound
- Next class: **logistic regression** (another linear classifier).

References I

Minsky, M. and Papert, S. (1969). Perceptrons.

Novikoff, A. B. (1962). On convergence proofs for perceptrons. In *Symposium on the Mathematical Theory of Automata*.

Rosenblatt, F. (1958). The perceptron: A probabilistic model for information storage and organization in the brain.
Psychological review, 65(6):386.

Lecture 4: Linear Models II

André Martins, Francisco Melo, Mário Figueiredo



Deep Learning Course, Winter 2021-2022

Announcements

- Homework 1 is out! It's due **January 7**.
- **Start early!!!**
- Post any questions in Piazza.

Today's Roadmap

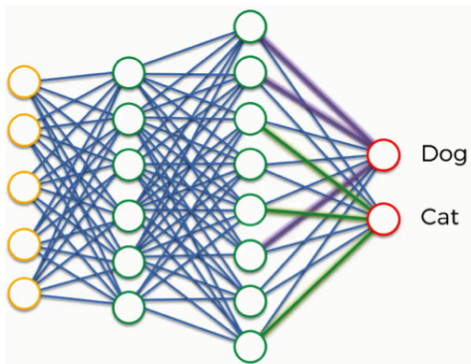
- Logistic regression
- Regularization and optimization
- Stochastic gradient descent.

Why Linear Classifiers?

We know the course title promised “**deep**”, but...

- The underlying machine learning concepts are the same
- The theory (statistics and optimization) are much better understood
- Linear classifiers still widely used (very effective when data is scarce)
- Linear classifiers are **a component of neural networks**.

Linear Classifiers and Neural Networks



So far

We have covered:

- The perceptron algorithm
- Naïve Bayes.

Perceptron is an instance of a **linear classifier**.

It finds a separating hyperplane (if it exists),

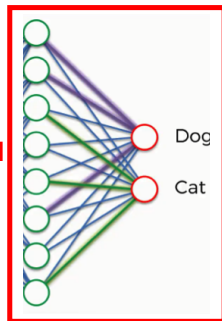
Naïve Bayes is a **generative** probabilistic model

Next: a **discriminative** probabilistic model.

Reminder



**Handcrafted
Features**

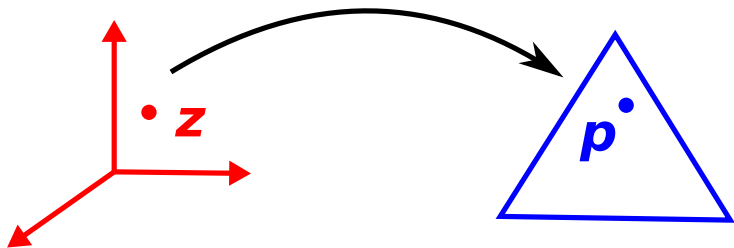


Linear Classifier

$$\hat{y} = \arg \max_y ((\mathbf{W}\phi(x) + \mathbf{b})_y), \quad \mathbf{W} = \begin{bmatrix} \vdots \\ \mathbf{w}_y^\top \\ \vdots \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} \vdots \\ b_y \\ \vdots \end{bmatrix}.$$

Key Problem

Mapping from a vector of scores $\mathbb{R}^{|\mathcal{Y}|}$ to a probability distribution over \mathcal{Y} ?



We will see an important mapping: **softmax** (next).

Logistic Regression

A linear model gives a **score** for each class y : $\mathbf{w}_y \cdot \phi(x)$.

From the score, we may compute a **conditional (posterior) probability**:

$$P(y|x) = \frac{\exp(\mathbf{w}_y \cdot \phi(x))}{Z_x}, \quad \text{where } Z_x = \sum_{y' \in \mathcal{Y}} \exp(\mathbf{w}_{y'} \cdot \phi(x))$$

This operation (exponentiating and normalizing) is called the **softmax transformation** (more later!)

Note: still a linear classifier

$$\begin{aligned} \arg \max_y P(y|x) &= \arg \max_y \frac{\exp(\mathbf{w}_y \cdot \phi(x))}{Z_x} \\ &= \arg \max_y \exp(\mathbf{w}_y \cdot \phi(x)) \\ &= \arg \max_y \mathbf{w}_y \cdot \phi(x) \end{aligned}$$

Binary Logistic Regression

Binary labels ($\mathcal{Y} = \{\pm 1\}$)

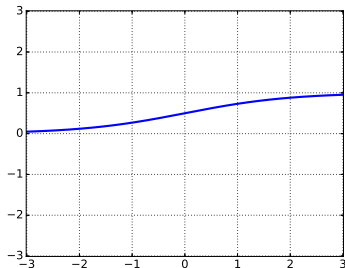
Scores: 0 for negative class, $w \cdot \phi(x)$ for positive class

$$\begin{aligned}P(y = +1 \mid x) &= \frac{\exp(w \cdot \phi(x))}{1 + \exp(w \cdot \phi(x))} \\ &= \frac{1}{1 + \exp(-w \cdot \phi(x))} \\ &= \sigma(w \cdot \phi(x)).\end{aligned}$$

This is called a **sigmoid transformation** (more later!)

Sigmoid Transformation

$$\sigma(u) = \frac{e^u}{1 + e^u}$$



- Widely used in neural networks
- Maps (*squashes*) \mathbb{R} into $[0, 1]$
- The output can be interpreted as a probability
- Positive, bounded, strictly increasing

Multinomial Logistic Regression

$$P_{\mathbf{W}}(y | x) = \frac{\exp(\mathbf{w}_y \cdot \phi(x))}{Z_x}$$

- How to learn weights \mathbf{W} ?
- Maximize the **conditional log-likelihood** of training data:

$$\begin{aligned}\widehat{\mathbf{W}} &= \arg \max_{\mathbf{W}} \log \left(\prod_{t=1}^N P_{\mathbf{W}}(y_t | x_t) \right) \\ &= \arg \min_{\mathbf{W}} - \sum_{t=1}^N \log P_{\mathbf{W}}(y_t | x_t) \\ &= \arg \min_{\mathbf{W}} \sum_{t=1}^N \left(\log \sum_{y'_t} \exp(\mathbf{w}_{y'_t} \cdot \phi(x_t)) - \mathbf{w}_{y_t} \cdot \phi(x_t) \right),\end{aligned}$$

- i.e., choose \mathbf{W} to maximize the probability of the true labels.

Logistic Regression

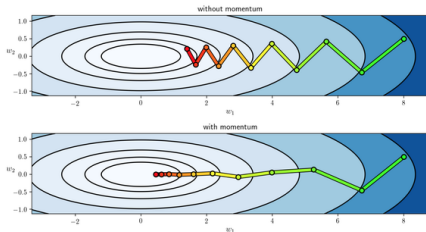
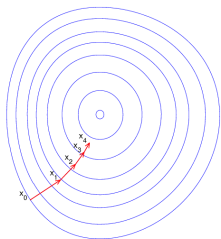
- This objective function is **strictly convex**
- Proof left as exercise! (hint, compute second derivatives, *i.e.*, Hessian)
- Therefore any local minimum is a global minimum
- No closed form solution, but lots of numerical techniques
 - ✓ Gradient methods (gradient descent, conjugate gradient)
 - ✓ Quasi-Newton methods (L-BFGS, ...)

Recap: Gradient Descent

- Goal: minimize $f : \mathbb{R}^d \rightarrow \mathbb{R}$, for differentiable f
- Take **small steps** in the **negative gradient direction** until a **stopping criterion** is met:

$$x^{(t+1)} \leftarrow x^{(t)} - \eta_{(t)} \nabla f(x^{(t)})$$

- Choosing the **step-size**: crucial for convergence and performance.
- GD may work well, or not so well. There are many ways to improve it.



Gradient Descent

- **Loss function** in logistic regression is

$$L(\mathbf{W}; (x, y)) = \log \sum_{y'} \exp(\mathbf{w}_{y'} \cdot \phi(x)) - \mathbf{w}_y \cdot \phi(x).$$

- We want to find

$$\arg \min_{\mathbf{W}} \sum_{t=1}^N L(\mathbf{W}; (x_t, y_t))$$

- ✓ Set $\mathbf{W}^0 = \mathbf{0}$
- ✓ Iterate until convergence (for suitable stepsize η_k):

$$\begin{aligned} \mathbf{W}^{k+1} &= \mathbf{W}^k - \eta_k \nabla_{\mathbf{W}} \left(\sum_{t=1}^N L(\mathbf{W}; (x_t, y_t)) \right) \\ &= \mathbf{W}^k - \eta_k \sum_{t=1}^N \nabla_{\mathbf{W}} L(\mathbf{W}^k; (x_t, y_t)) \end{aligned}$$

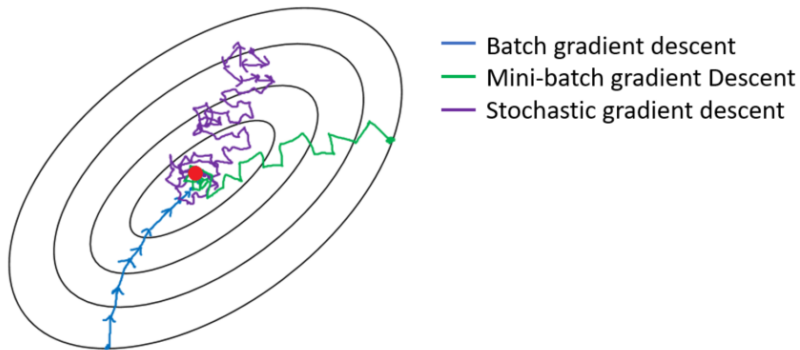
- $\nabla_{\mathbf{W}} L(\mathbf{W})$ is gradient of L w.r.t. \mathbf{W}
- $L(\mathbf{W})$ convex \Rightarrow gradient descent will reach the global optimum \mathbf{W} .

Stochastic Gradient Descent

Monte Carlo approximation of the gradient (more frequent updates, convenient with large datasets):

- Set $\mathbf{W}^0 = \mathbf{0}$
- Iterate until convergence
 - Pick (x_t, y_t) randomly
 - Update $\mathbf{W}^{k+1} = \mathbf{W}^k - \eta_k \nabla_{\mathbf{W}} L(\mathbf{W}^k; (x_t, y_t))$
- *i.e.* approximate the gradient with a noisy, unbiased, version based on a single sample
- Variants exist in-between (mini-batches)
- All guaranteed to find the optimal \mathbf{W} (for suitable step sizes)

Stochastic vs Batch Gradient Descent



Computing the Gradient

- We need to compute $\nabla_{\mathbf{W}} L(\mathbf{W}; (x_t, y_t))$, where

$$L(\mathbf{W}; (x, y)) = \log \sum_{y'} \exp(\mathbf{w}_{y'} \cdot \phi(x)) - \mathbf{w}_y \cdot \phi(x)$$

- Some reminders:

① $\nabla_{\mathbf{W}} \log F(\mathbf{W}) = \frac{1}{F(\mathbf{W})} \nabla_{\mathbf{W}} F(\mathbf{W})$

② $\nabla_{\mathbf{W}} \exp F(\mathbf{W}) = \exp(F(\mathbf{W})) \nabla_{\mathbf{W}} F(\mathbf{W})$

- We denote by

$$\mathbf{e}_y = [0, \dots, 0, 1, 0, \dots, 0]^\top, \quad 1 \text{ in } i\text{-th position}$$

the one-hot vector representation of class y .

Computing the Gradient

$$\begin{aligned}\nabla_{\mathbf{W}} L(\mathbf{W}; (x, y)) &= \nabla_{\mathbf{W}} \left(\log \sum_{y'} \exp(\mathbf{w}_{y'} \cdot \phi(x)) - \mathbf{w}_y \cdot \phi(x) \right) \\ &= \nabla_{\mathbf{W}} \log \sum_{y'} \exp(\mathbf{w}_{y'} \cdot \phi(x)) - \nabla_{\mathbf{W}} \mathbf{w}_y \cdot \phi(x) \\ &= \frac{1}{\sum_{y'} \exp(\mathbf{w}_{y'} \cdot \phi(x))} \sum_{y'} \nabla_{\mathbf{W}} \exp(\mathbf{w}_{y'} \cdot \phi(x)) - \mathbf{e}_y \phi(x)^\top \\ &= \frac{1}{Z_x} \sum_{y'} \exp(\mathbf{w}_{y'} \cdot \phi(x)) \nabla_{\mathbf{W}} \mathbf{w}_{y'} \cdot \phi(x) - \mathbf{e}_y \phi(x)^\top \\ &= \sum_{y'} \frac{\exp(\mathbf{w}_{y'} \cdot \phi(x))}{Z_x} \mathbf{e}_{y'} \phi(x)^\top - \mathbf{e}_y \phi(x)^\top \\ &= \sum_{y'} P_{\mathbf{W}}(y'|x) \mathbf{e}_{y'} \phi(x)^\top - \mathbf{e}_y \phi(x)^\top \\ &= \left(\begin{bmatrix} \vdots \\ P_{\mathbf{W}}(y'|x) \\ \vdots \end{bmatrix} - \mathbf{e}_y \right) \phi(x)^\top.\end{aligned}$$

Logistic Regression Summary

- Define conditional probability

$$P_{\mathbf{W}}(y|x) = \frac{\exp(\mathbf{w}_y \cdot \phi(x))}{Z_x}$$

- Set weights to maximize conditional log-likelihood of training data:

$$\mathbf{W} = \arg \max_{\mathbf{W}} \sum_t \log P_{\mathbf{W}}(y_t|x_t) = \arg \min_{\mathbf{W}} \sum_t L(\mathbf{W}; (x_t, y_t))$$

- Run gradient descent (or any gradient-based optimization algorithm) using

$$\nabla_{\mathbf{W}} L(\mathbf{W}; (x, y)) = \sum_{y'} P_{\mathbf{W}}(y'|x) \mathbf{e}_{y'} \phi(x)^{\top} - \mathbf{e}_y \phi(x)^{\top}$$

The Story So Far

- Naive Bayes is **generative**: maximizes **joint** likelihood
 - closed-form solution
- Logistic regression is **discriminative**: maximizes **conditional** likelihood
 - also called log-linear model and max-entropy classifier
 - no closed-form solution
 - stochastic gradient updates look like

$$\mathbf{W}^{k+1} = \mathbf{W}^k + \eta \left(\mathbf{e}_y \phi(x)^\top - \sum_{y'} P_w(y'|x) \mathbf{e}_{y'} \phi(x)^\top \right)$$

- Perceptron is a discriminative, non-probabilistic classifier
 - perceptron's updates look like

$$\mathbf{W}^{k+1} = \mathbf{W}^k + \mathbf{e}_y \phi(x)^\top - \mathbf{e}_{\hat{y}} \phi(x)^\top$$

- SGD updates for logistic regression and the perceptron look similar!

Other Options: Maximizing Margin

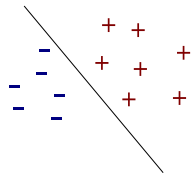
- For a training set \mathcal{D}
- Margin of a weight matrix \mathbf{W} is smallest γ such that

$$\mathbf{w}_{y_t} \cdot \phi(\mathbf{x}_t) - \mathbf{w}_{y'} \cdot \phi(\mathbf{x}_t) \geq \gamma$$

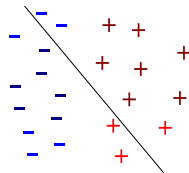
- for every training instance $(\mathbf{x}_t, y_t) \in \mathcal{D}$, $y' \in \mathcal{Y}$

Margin

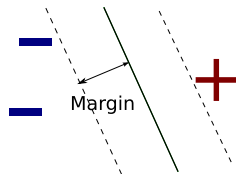
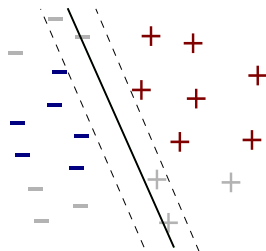
Training



Testing



Denote the value of the margin by γ



Maximizing Margin

- Intuitively maximizing margin makes sense
- More importantly, generalization error to unseen test data is proportional to the inverse of the margin

$$\epsilon \propto \frac{R^2}{\gamma^2 \times N}$$

- **Perceptron:**
 - If a training set is separable by some margin, the perceptron will find a \mathbf{W} that separates the data
 - However, the perceptron does not pick \mathbf{W} to maximize the margin!
- Support Vector Machines do this (not covered)

Summary

What we saw

- Linear Classifiers
 - Naive Bayes
 - Logistic Regression
 - Perceptron
 - Support Vector Machines (not covered)

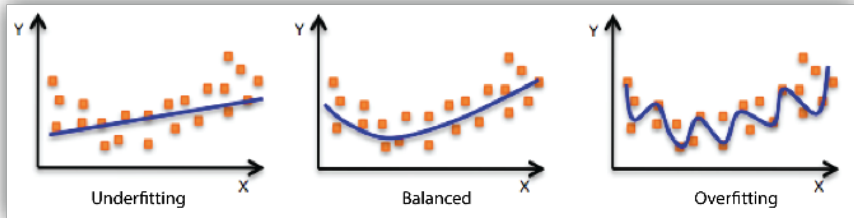
What is next

- Regularization
- Softmax
- Non-linear classifiers

Regularization

Overfitting

If the model is too complex (too many parameters) and the data is scarce, we run the risk of **overfitting**:



- We saw one example already when talking about add-one smoothing in Naive Bayes!

Regularization

In practice, we **regularize** to prevent overfitting

$$\arg \min_{\mathbf{W}} \sum_{t=1}^N L(\mathbf{W}; (x_t, y_t)) + \lambda \Omega(\mathbf{W}),$$

$\Omega(\mathbf{W})$ is the regularization function, and λ controls its weight

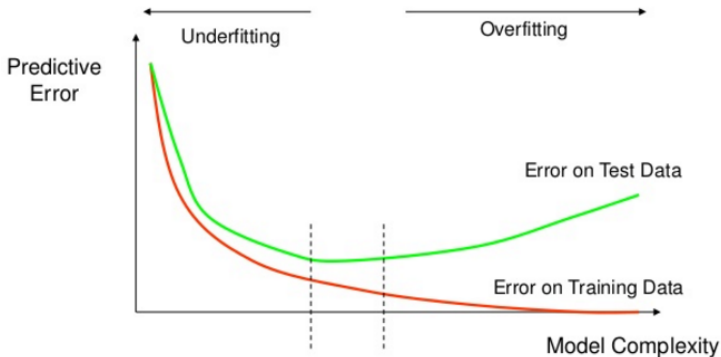
- ℓ_2 regularization promotes **smaller** weights:

$$\Omega(\mathbf{W}) = \|\mathbf{W}\|_2^2 = \sum_y \|\mathbf{w}_y\|_2^2 = \sum_y \sum_j w_{y,j}^2.$$

- ℓ_1 regularization promotes **sparse** weights!

$$\Omega(\mathbf{W}) = \|\mathbf{W}\|_1 = \sum_y \|\mathbf{w}_y\|_1 = \sum_y \sum_j |w_{y,j}|$$

Empirical Risk Minimization



Logistic Regression with ℓ_2 Regularization

$$\min - \sum_{t=1}^N \log(\exp(\mathbf{w}_{y_t} \cdot \phi(x_t)) / Z_x) + \frac{\lambda}{2} \|\mathbf{W}\|^2$$

- What is the new gradient?

$$\sum_{t=1}^N \nabla_{\mathbf{W}} L(\mathbf{W}; (x_t, y_t)) + \nabla_{\mathbf{W}} \lambda \Omega(\mathbf{W})$$

- We know $\nabla_{\mathbf{W}} L(\mathbf{W}; (x_t, y_t))$
- Just need $\nabla_{\mathbf{W}} \frac{\lambda}{2} \|\mathbf{W}\|^2 = \lambda \mathbf{W}$

Loss Function

Should match the metric we want to optimize at test time

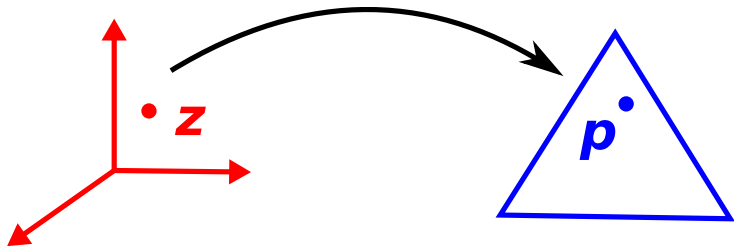
Should be well-behaved (convex, maybe smooth) to be amenable to optimization (this rules out the 0/1 loss)

Some examples:

- Squared loss for regression
- Negative log-likelihood (cross-entropy): multinomial logistic regression
- Hinge loss: support vector machines
- Sparsemax loss for multi-class and multi-label classification

Recap

How to map from a set of label scores $\mathbb{R}^{|\mathcal{Y}|}$ to a probability distribution over \mathcal{Y} ?



We already saw one example: softmax.

Another example is **sparsemax** (not covered): Martins and Astudillo (2016)

Recap: Softmax Transformation

The typical transformation for multi-class classification:

$$\mathbf{softmax} : \mathbb{R}^{|\mathcal{Y}|} \rightarrow \Delta^{|\mathcal{Y}|-1}$$

$$\mathbf{softmax}(z) = \left[\frac{\exp(z_1)}{\sum_c \exp(z_c)}, \dots, \frac{\exp(z_{|\mathcal{Y}|})}{\sum_c \exp(z_c)} \right]$$

- Underlies multinomial logistic regression!
- Strictly positive, sums to 1
- Resulting distribution has full support: $\mathbf{softmax}(z) > \mathbf{0}, \forall z$

Recap: Multinomial Logistic Regression

- The common choice for a softmax output layer
- The classifier estimates $P(y = c \mid x; \mathbf{W})$
- We minimize the negative log-likelihood:

$$\begin{aligned}L(\mathbf{W}; (x, y)) &= -\log P(y \mid x; \mathbf{W}) \\ &= -\log [\mathbf{softmax}(z(x))]_y,\end{aligned}$$

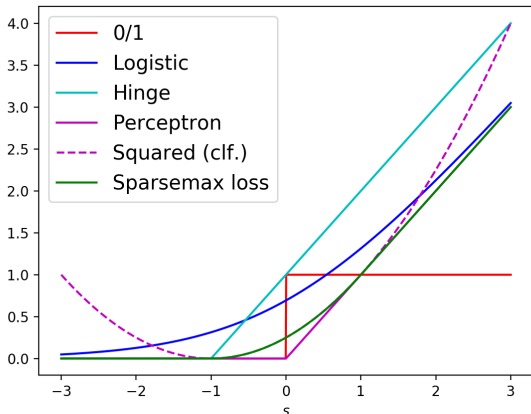
where $z_c(x) = \mathbf{w}_c \cdot \phi(x)$ is the score of class c .

- Loss gradient:

$$\nabla_{\mathbf{W}} L((x, y); \mathbf{W}) = - \left(\mathbf{e}_y \phi(x)^\top - \mathbf{softmax}(z(x)) \phi(x)^\top \right)$$

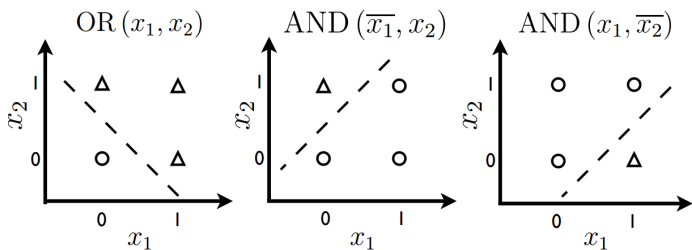
Classification Losses (Binary Case)

- Let the correct label be $y = +1$ and define $s = z_2 - z_1$.
- Sparsemax loss in 2D becomes a “classification Huber loss”:



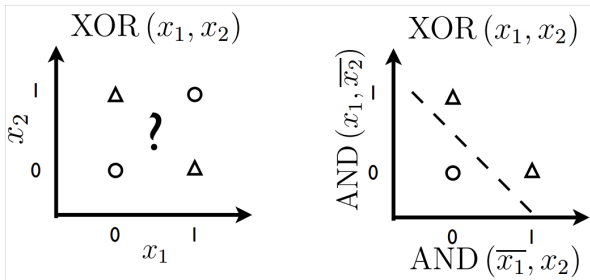
Recap: What a Linear Classifier Can Do

- It **can** solve linearly separable problems (OR, AND)



Recap: What a Linear Classifier **Can't** Do

- ... but it **can't** solve **non-linearly separable** problems such as simple XOR (unless input is transformed into a better representation):



- This was observed by Minsky and Papert (1969) (for the perceptron) and motivated strong criticisms

Summary: Linear Classifiers

We've seen

- Perceptron
- Naive Bayes
- Logistic regression
- Support vector machines (not covered)

All lead to **convex** optimization problems \Rightarrow no issues with local minima/initialization

All assume the features are well-engineered such that **the data is nearly linearly separable**

What If Data Are Not Linearly Separable?

Engineer better features (often works!)



Kernel methods:

- works implicitly in a high-dimensional feature space
- ... but still need to choose/design a good kernel
- model capacity confined to positive-definite kernels



Neural networks (**next class!**)

- embrace non-convexity and local minima
- instead of engineering features/kernels, engineer the model architecture

Two Views of Machine Learning

There's two big ways of building machine learning systems:

- ① **Feature-based**: describe objects' properties (features) and build models that manipulate them
 - everything that we have seen so far.
- ② **Similarity-based**: don't describe objects by their properties; rather, build systems based on **comparing** objects to each other
 - k -th nearest neighbors; kernel methods; Gaussian processes.

Sometimes the two are equivalent!

Nearest Neighbor Classifier

- Not a linear classifier!
- In its simplest version, doesn't require any parameters
- Instead of “training”, **memorize** all the data $\mathcal{D} = \{(x_i, y_i)_{i=1}^N\}$
- Given a new input x , find its **most similar** data point x_i and predict

$$\hat{y} = y_i$$

- Many variants (e.g. k -th nearest neighbor)
- **Disadvantage:** requires searching over the entire training data
- Specialized data structures can be used to speed up search.

Kernels

- A kernel is a similarity function between two points that is **symmetric** and **positive semi-definite**, which we denote by:

$$\kappa(x_i, x_j) \in \mathbb{R}$$

- Given dataset $\mathcal{D} = \{(x_i, y_i)_{i=1}^N\}$, the **Gram matrix** \mathbf{K} is the $N \times N$ matrix defined as:

$$K_{i,j} = \kappa(x_i, x_j)$$

- **Symmetric:**

$$\kappa(x_i, x_j) = \kappa(x_j, x_i)$$

- **Positive definite:** for all non-zero \mathbf{v}

$$\mathbf{v}\mathbf{K}\mathbf{v}^T \geq 0$$

Kernels

- **Mercer's Theorem:** for any kernel $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, there exists some feature mapping $\phi : \mathcal{X} \rightarrow \mathbb{R}^{\mathcal{X}}$, s.t.:

$$\kappa(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$$

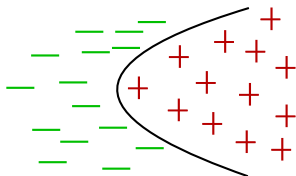
- That is: a kernel corresponds to some a mapping in some **implicit** feature space!
- **Kernel trick:** take a feature-based algorithm (SVMs, perceptron, logistic regression) and replace all explicit feature computations by **kernel evaluations!**

$$\mathbf{w}_y \cdot \phi(x) = \sum_{i=1}^N \sum_{y \in \mathcal{Y}} \alpha_{i,y} \kappa(x, x_i) \quad \text{for some } \alpha_{i,y} \in \mathbb{R}$$

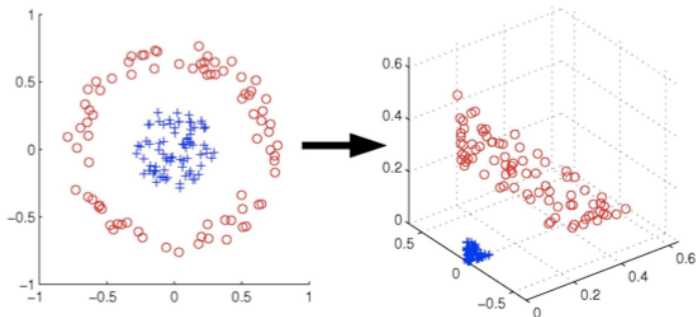
- Extremely popular idea in the 1990-2000s!

Kernels = Tractable Non-Linearity

- A linear classifier in a higher dimensional feature space is a non-linear classifier in the original space
- Computing a non-linear kernel is sometimes better computationally than calculating the corresponding dot product in the high dimension feature space
- Many models can be “kernelized” – learning algorithms generally solve the **dual** optimization problem (also convex)
- Drawback: **quadratic** dependency on dataset size



Linear Classifiers in High Dimension



$$\mathbb{R}^2 \longrightarrow \mathbb{R}^3$$

$$(x_1, x_2) \longmapsto (z_1, z_2, z_3) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

Popular Kernels

- Polynomial kernel

$$\kappa(x_i, x_j) = (\phi(x_i) \cdot \phi(x_j) + 1)^d$$

- Gaussian radial basis kernel

$$\kappa(x_i, x_j) = \exp\left(\frac{-\|\phi(x_i) - \phi(x_j)\|^2}{2\sigma}\right)$$

- String kernels (Lodhi et al., 2002; Collins and Duffy, 2002)
- Tree kernels (Collins and Duffy, 2002)

Conclusions

- Linear classifiers are a broad class including well-known ML methods such as **perceptron**, **logistic regression**, **support vector machines**
- They all involve manipulating weights and features
- They either lead to closed-form solutions or **convex** optimization problems (**no local minima**)
- Stochastic gradient descent algorithms are useful if training datasets are large
- However, they require manual specification of feature representations

References I

- Collins, M. and Duffy, N. (2002). Convolution kernels for natural language. *Advances in Neural Information Processing Systems*, 1:625–632.
- Lodhi, H., Saunders, C., Shawe-Taylor, J., Cristianini, N., and Watkins, C. (2002). Text classification using string kernels. *Journal of Machine Learning Research*, 2:419–444.
- Martins, A. F. T. and Astudillo, R. (2016). From Softmax to Sparsemax: A Sparse Model of Attention and Multi-Label Classification. In *Proc. of the International Conference on Machine Learning*.
- Minsky, M. and Papert, S. (1969). Perceptrons.

Lecture 5: Neural Networks I

André Martins, Francisco Melo, Mário Figueiredo



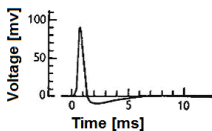
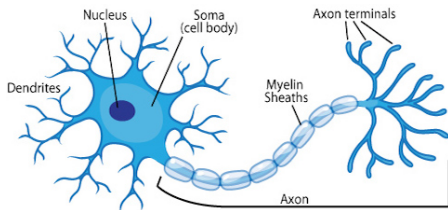
Deep Learning Course, Winter 2021-2022

Today's Roadmap

Today's lecture is about **neural networks**:

- From perceptron to **multi-layer** perceptron
- Feed-forward neural networks
- **Activation functions**: sigmoid, tanh, relu, ...
- **Activation maps**: softmax, sparsemax, ...
- Non-convex optimization and local minima
- Universal approximation theorem
- Gradient backpropagation

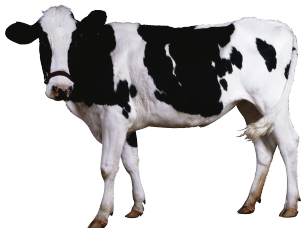
Biological Neuron



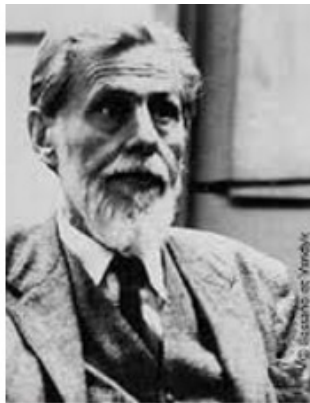
- Three main parts: the main body (**soma**), **dendrites** and an **axon**
- The neuron receives **input** signals from **dendrites**, and **outputs** its own signals through the **axon**
- **Axons** in turn connect to the **dendrites** of other neurons; the connections called **synapses**
- Generate sharp electrical potentials across their cell membrane (**spikes**), the main signalling unit of the nervous system

Word of Caution

- Artificial neurons are inspired by biological neurons, but...



Warren McCulloch and Walter Pitts



- The earliest computational model of a neuron, via **threshold logic** (McCulloch and Pitts, 1943).

Artificial Neuron (McCulloch and Pitts, 1943)

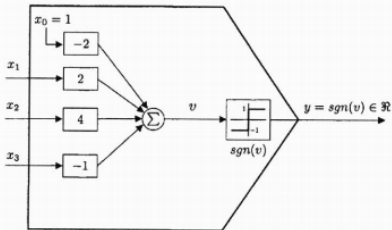


Figure 3.6 Example 3.2: a threshold neural logic for $y = x_2(x_1 + \bar{x}_3)$.

Table 3.6 Truth table for Example 3.2

Neural Inputs			$v = \mathbf{w}_a^T \mathbf{x}_a$ $= -2 + 2x_1 + 4x_2 - x_3$	$y = \text{sgn}(v)$ $= \text{sgn}(\mathbf{w}_a^T \mathbf{x}_a)$
x_1	x_2	x_3		
-1	-1	-1	-7	-1
-1	-1	1	-9	-1
-1	1	-1	1	1
-1	1	1	-1	-1
1	-1	-1	-3	-1
1	-1	1	-5	-1
1	1	-1	5	1
1	1	1	3	1

- Later models replaced the hard threshold by more general activations

Artificial Neuron

- **Pre-activation** (input activation):

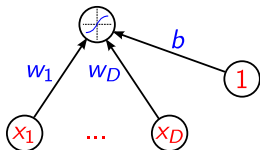
$$z(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b = \sum_{i=1}^D w_i x_i + b,$$

where \mathbf{w} are the **connection weights** and b is a **bias term**.

- **Activation:**

$$h(\mathbf{x}) = g(z(\mathbf{x})) = g(\mathbf{w} \cdot \mathbf{x} + b),$$

where $g : \mathbb{R} \rightarrow \mathbb{R}$ is the **activation function**.



Activation Function

Typical choices:

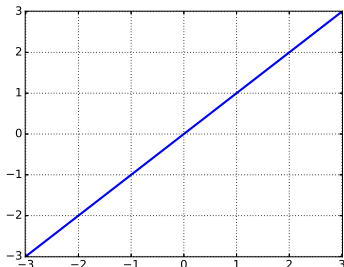
- Linear
- Sigmoid (logistic function)
- Hyperbolic Tangent
- Rectified Linear

Later:

- Softmax
- Sparsemax
- Max-pooling

Linear Activation

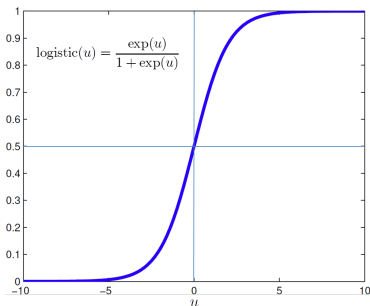
$$g(z) = z$$



- No “squashing” of the input
- Composing layers of linear units is equivalent to a single linear layer: no expressive power increase by using multiple layers (more later)
- Still useful to linear-project the input to a lower dimension

Sigmoid Activation

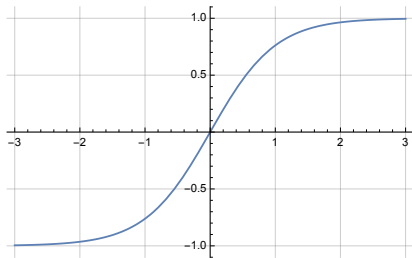
$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$



- “Squashes” the neuron pre-activation between 0 and 1
- The output can be interpreted as a probability
- Positive, bounded, strictly increasing
- Logistic regression corresponds to a network with a single sigmoid unit
- Combining layers of sigmoid units increases expressiveness (more later)

Hyperbolic Tangent Activation

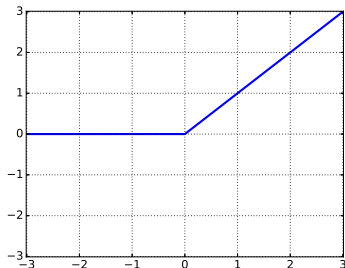
$$g(z) = \mathbf{tanh}(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



- “Squashes” the neuron pre-activation between -1 and 1
- Related to the sigmoid via $\sigma(z) = \frac{1+\mathbf{tanh}(z/2)}{2}$
- **Can be positive or negative**, bounded, strictly increasing
- Combining layers of tanh units increases expressiveness (more later)

Rectified Linear Unit Activation (Glorot et al., 2011)

$$g(z) = \text{relu}(z) = \max\{0, z\}$$



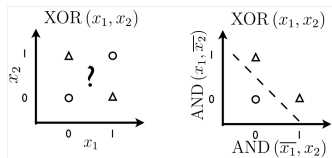
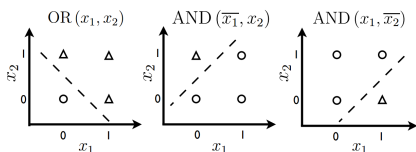
- Non-negative, increasing, **but not upper bounded**
- Not differentiable at 0
- Leads to neurons with **sparse activities** (biologically more plausible)
- Less prone to vanishing gradients (more later), and historically the first activation that allowed training deep nets without unsupervised pre-training (Glorot et al., 2011)

Capacity of Single Neuron (Linear Classifier)

- With a single sigmoid activated neuron we recover **logistic regression**:

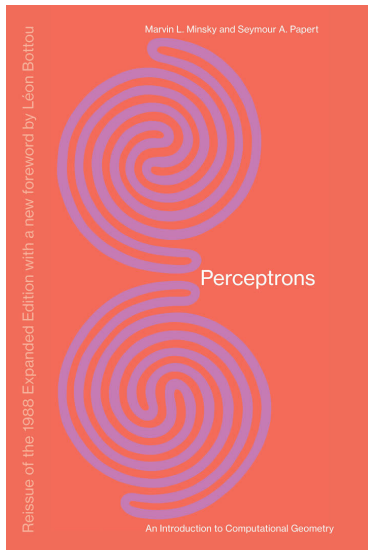
$$p(y = 1|\mathbf{x}) = \sigma(\mathbf{w} \cdot \mathbf{x} + b).$$

- Can solve linearly separable problems (OR, AND)
- Can't solve non-linearly separable problems (XOR)—unless input is transformed into a better representation



(Slide credit: Hugo Larochelle)

The XOR Affair



Minsky and Papert (1969):

- Misunderstood by many as showing a single perceptron cannot learn XOR (in fact, this was already well-known at the time)
- Fostered an “AI winter” period

Solving XOR with Multi-Layer Perceptron

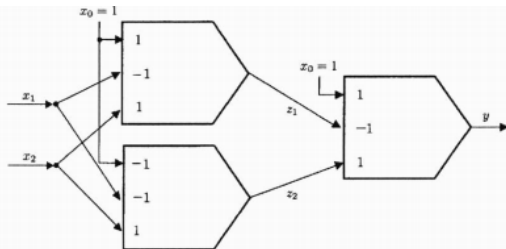
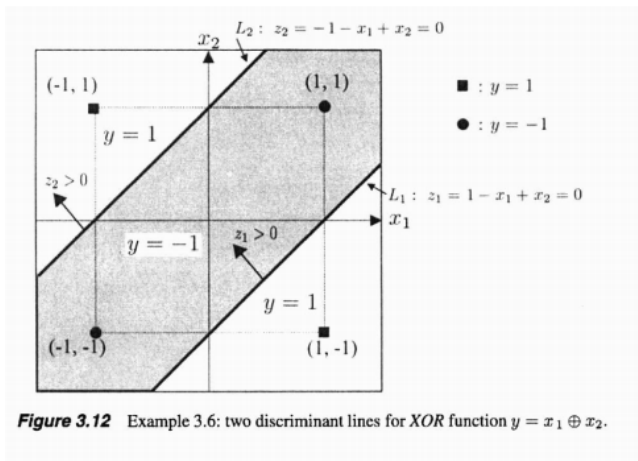


Figure 3.11 Example 3.6: a threshold network for XOR function $y = x_1 \oplus x_2 = \text{sgn}(1 - z_1 + z_2)$, $z_1 = \text{sgn}(1 - x_1 + x_2)$, $z_2 = \text{sgn}(-1 - x_1 + x_2)$.

Solving XOR with Multi-Layer Perceptron



Solving XOR with Multi-Layer Perceptron

- This construction was known even by McCulloch and Pitts
- The negative result in Minsky and Papert (1969) is that, to learn arbitrary logic functions, each hidden unit needs to be connected to **all inputs**
- At the time, there was some hope that we'd only need “local” neurons

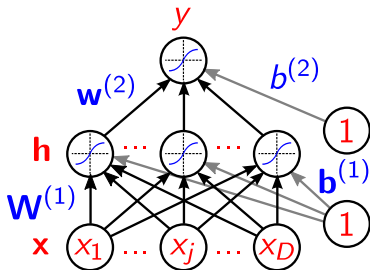
Multi-Layer Neural Network

- **Key idea:** add intermediate layers of artificial neurons before the final output layer
- Each of these hidden units computes some representation of the input and propagates forward that representation
- This increases the expressive power of the network, yielding more complex, non-linear, functions/classifiers
- Similar role as latent variables in probabilistic models, but no need for a probability semantics
- Also called **feed-forward neural network**

Single Hidden Layer

To start simple, let's assume our task involves several inputs ($\mathbf{x} \in \mathbb{R}^D$) but a **single output** (e.g. $y \in \mathbb{R}$ or $y \in \{0, 1\}$)

Trick: add an intermediate layer of K hidden units ($\mathbf{h} \in \mathbb{R}^K$)



Single Hidden Layer

Assume D inputs ($\mathbf{x} \in \mathbb{R}^D$) and K hidden units ($\mathbf{h} \in \mathbb{R}^K$)

- **Hidden layer pre-activation:**

$$z(\mathbf{x}) = \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)},$$

with $\mathbf{W}^{(1)} \in \mathbb{R}^{K \times D}$ and $\mathbf{b}^{(1)} \in \mathbb{R}^K$.

- **Hidden layer activation:**

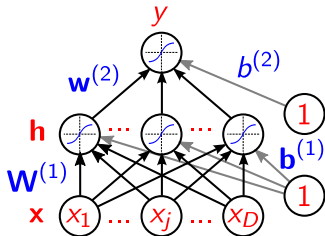
$$\mathbf{h}(z) = \mathbf{g}(z(\mathbf{x})),$$

where $\mathbf{g} : \mathbb{R}^K \rightarrow \mathbb{R}^K$ is applied component-wise (component-by-component).

- **Output layer activation:**

$$f(\mathbf{x}) = o(\mathbf{w}^{(2)} \cdot \mathbf{h} + b^{(2)}),$$

where $\mathbf{w}^{(2)} \in \mathbb{R}^K$ and $o : \mathbb{R} \rightarrow \mathbb{R}$ if the output activation function.



Single Hidden Layer

Overall,

$$\begin{aligned} f(\mathbf{x}) &= o(\mathbf{w}^{(2)} \cdot \mathbf{h} + b^{(2)}) \\ &= o(\mathbf{w}^{(2)} \cdot \mathbf{g}(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + b^{(2)}) \end{aligned}$$

Examples:

- $o(u) = u$ for **regression** ($y \in \mathbb{R}$)
- $o(u) = \sigma(u)$ for **binary classification** ($y \in \{\pm 1\}$, $f(\mathbf{x}) = P(y = 1 | \mathbf{x})$)

No longer a linear classifier – non-linear dependency on $\mathbf{W}^{(1)}$ and $\mathbf{b}^{(1)}$

\mathbf{h} is a vector of **internal representations** (not manually engineered features)

Multiple Classes

Can we use a similar strategy for the **multi-class** case?

For **multi-class** classification, we need **multiple** output units (one per class)

Each output estimates the conditional probability $P(y = c | \mathbf{x})$

Predicted class is (usually) the one with highest estimated probability

We have already seen an activation function suitable for this: **softmax**

Softmax Activation

Let $\Delta_{C-1} \subseteq \mathbb{R}^C$ be the probability simplex:

$$\Delta_{C-1} = \{(p_1, \dots, p_C) : p_i \geq 0, i = 1, \dots, C, \sum_{i=1}^C p_i = 1\}$$

Typical activation function for multi-class: **softmax** : $\mathbb{R}^C \rightarrow \Delta_{C-1}$:

$$\mathbf{o}(z) = \mathbf{softmax}(z) = \left[\frac{\exp(z_1)}{\sum_c \exp(z_c)}, \dots, \frac{\exp(z_C)}{\sum_c \exp(z_c)} \right]$$

- We saw this previously, when talking about logistic regression!
- Strictly positive, sums to 1
- Resulting distribution has full support: **softmax**(z) > **0**, $\forall z$

Multi-Layer Neural Networks: General Case

In general we can:

- Have multiple output units (needed for multi-class classification)
- Stack more layers after each other

Multiple ($L \geq 1$) Hidden Layers

- **Hidden layer pre-activation** (define $\mathbf{h}^{(0)} = \mathbf{x}$ for convenience):

$$\mathbf{z}^{(\ell)}(\mathbf{x}) = \mathbf{W}^{(\ell)}\mathbf{h}^{(\ell-1)} + \mathbf{b}^{(\ell)},$$

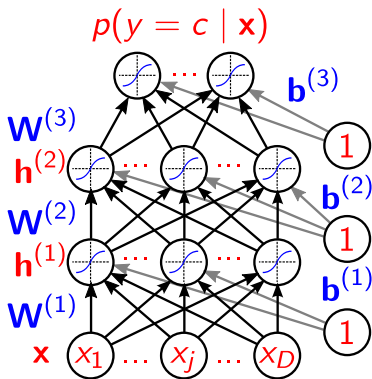
with $\mathbf{W}^{(\ell)} \in \mathbb{R}^{K_\ell \times K_{\ell-1}}$ and $\mathbf{b}^{(\ell)} \in \mathbb{R}^{K_\ell}$.

- **Hidden layer activation:**

$$\mathbf{h}^{(\ell)}(\mathbf{x}) = \mathbf{g}(\mathbf{z}^{(\ell)}(\mathbf{x})).$$

- **Output layer activation:**

$$\mathbf{f}(\mathbf{x}) = \mathbf{o}(\mathbf{z}^{(L+1)}(\mathbf{x})) = \mathbf{o}(\mathbf{W}^{(L+1)}\mathbf{h}^{(L)} + \mathbf{b}^{(L+1)}).$$



Universal Approximation Theorem

Theorem (Hornik et al. (1989))

An NN with one hidden layer and a linear output can approximate arbitrarily well any continuous function, given enough hidden units.

- First proved for the sigmoid case by Cybenko (1989), then to **tanh** and many other activation functions by Hornik et al. (1989)
- **Caveat:** may need **exponentially** many hidden units

Deeper Networks

- **Deeper networks** (more layers) can provide more compact approximations

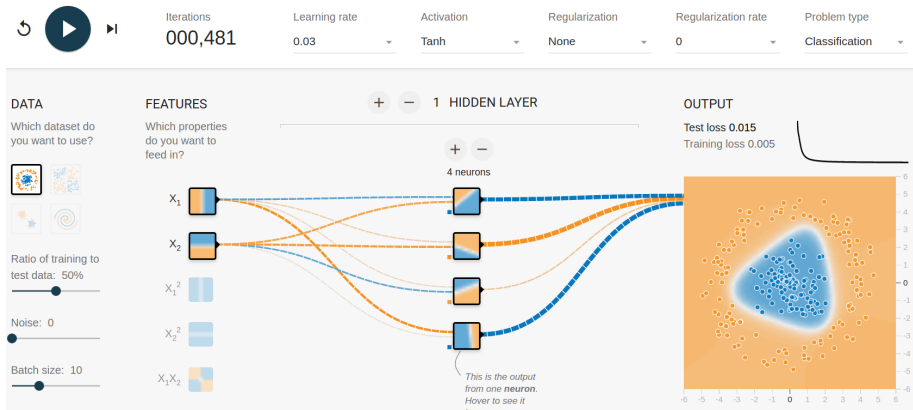
Theorem (Montufar et al. (2014))

The number of linear regions carved out by a deep neural network with D inputs, depth L , and K hidden units per layer with ReLU activations is

$$O\left(\binom{K}{D}^{D(L-1)} K^D\right)$$

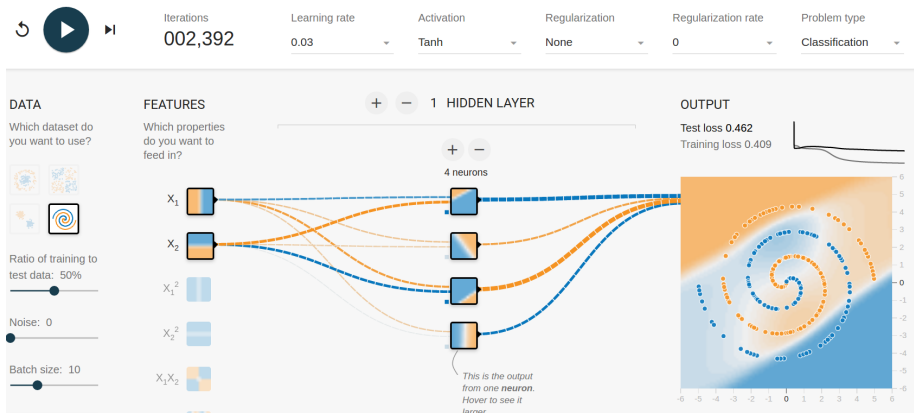
Therefore, for fixed K , deeper networks are exponentially more expressive

“Simple” Target Function, One Hidden Layer



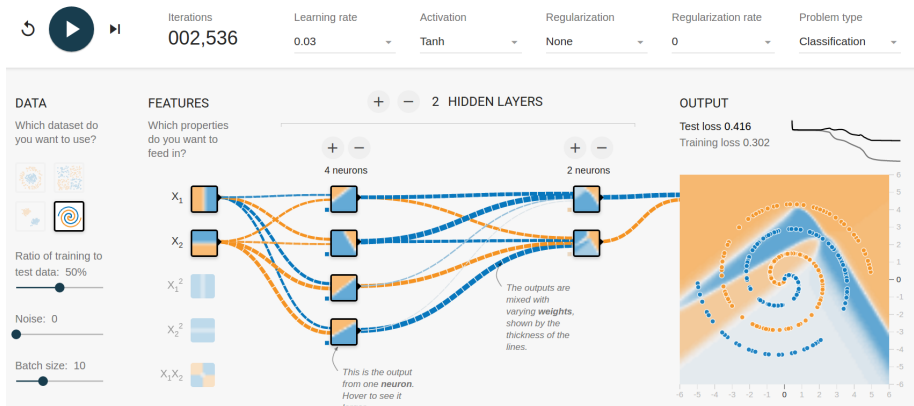
(<http://playground.tensorflow.org>)

Complex Target Function, One Hidden Layer



(<http://playground.tensorflow.org>)

Complex Target Function, Two Hidden Layers



(<http://playground.tensorflow.org>)

Complex Target Function, Two Hidden Layers, ReLU



Iterations
001,968

Learning rate
0.03

Activation
ReLU

Regularization
None

Regularization rate
0

Problem type
Classification

DATA

Which dataset do you want to use?



Ratio of training to test data: 50%



Noise: 0



Batch size: 10

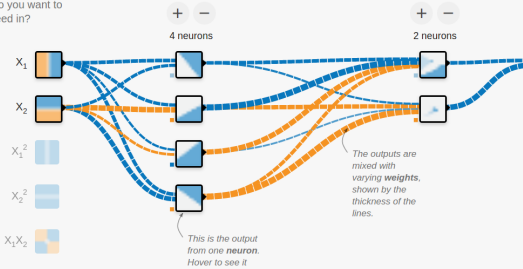


FEATURES

Which properties do you want to feed in?

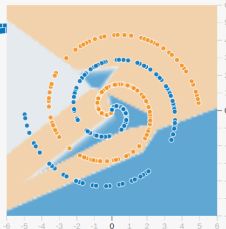
- X_1
- X_2
- X_1^2
- X_2^2
- $X_1 X_2$

2 HIDDEN LAYERS



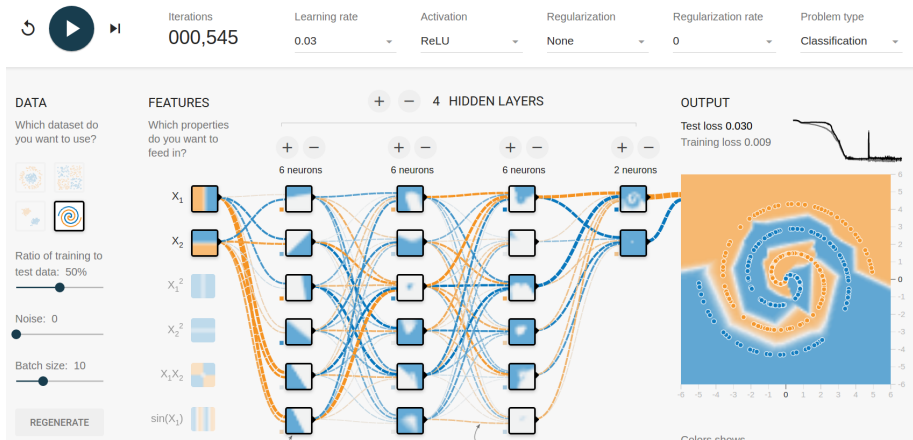
OUTPUT

Test loss 0.340
Training loss 0.321



(<http://playground.tensorflow.org>)

Complex Target Function, Four Hidden Layers, ReLU



(<http://playground.tensorflow.org>)

Capacity of Neural Networks

Neural networks are **excellent function approximators!**

The universal approximation theorem is an important result, but:

- We need a **learning algorithm** that finds the necessary parameter values
- ... and if we want to generalize, we need to control **overfitting**

Training Neural Networks

Neural networks are expressive: in theory, **they approximate any function**

But to do so, their **parameters**

$$\theta := \{(\mathbf{W}^{(1)}, \mathbf{b}^{(1)}), \dots, (\mathbf{W}^{(L+1)}, \mathbf{b}^{(L+1)})\}$$

need to be set accordingly

Key idea: **learn** these parameters from data

In other words: learn a function by sampling a few points and their values

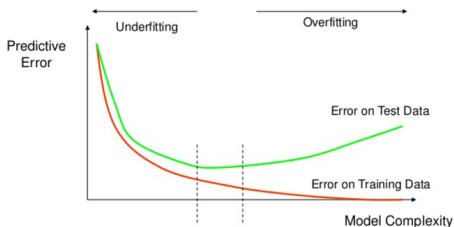
(We've seen this when we talked about linear models a few days ago...)

Empirical Risk Minimization

Goal: choose parameters $\theta := \{(\mathbf{W}^{(\ell)}, \mathbf{b}^{(\ell)})\}_{\ell=1}^{L+1}$ that minimize the following objective function:

$$\mathcal{L}(\theta) := \lambda\Omega(\theta) + \frac{1}{N} \sum_{i=1}^N L(\mathbf{f}(x_i; \theta), y_i)$$

- $L(\mathbf{f}(x_i; \theta), y_i)$ is a **loss function**
- $\Omega(\theta)$ is a **regularizer**
- λ is a **regularization constant** (an **hyperparameter** that needs to be tuned)



Recap: Gradient Descent

We can write the objective as:

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}) &:= \lambda\Omega(\boldsymbol{\theta}) + \frac{1}{N} \sum_{i=1}^N L(\mathbf{f}(\mathbf{x}_i; \boldsymbol{\theta}), y_i) \\ &:= \frac{1}{N} \sum_{i=1}^N \underbrace{\lambda\Omega(\boldsymbol{\theta}) + L(\mathbf{f}(\mathbf{x}_i; \boldsymbol{\theta}), y_i)}_{\mathcal{L}_i(\boldsymbol{\theta})} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_i(\boldsymbol{\theta})\end{aligned}$$

The gradient is:

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) := \frac{1}{N} \sum_{i=1}^N \nabla_{\boldsymbol{\theta}} \mathcal{L}_i(\boldsymbol{\theta})$$

Requires a full pass over the data to update the weights—**too slow!**

Recap: Stochastic Gradient Descent

Sample a **single** training example **uniformly at random**: $j \in \{1, \dots, N\}$

This way we get a noisy but **unbiased** estimate of the gradient:

$$\begin{aligned}\nabla_{\theta} \mathcal{L}(\theta) &:= \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \mathcal{L}_i(\theta) \approx \nabla_{\theta} \mathcal{L}_j(\theta) \\ &= \lambda \nabla_{\theta} \Omega(\theta) + \nabla_{\theta} L(f(\mathbf{x}_j; \theta), y_j).\end{aligned}$$

The weights $\theta = \{(\mathbf{W}^{(\ell)}, \mathbf{b}^{(\ell)})\}_{\ell=1}^{L+1}$ are then updated as:

$$\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}_j(\theta)$$

In practice, use **mini-batch** instead of a single sample

Stochastic Gradient Descent with Mini-Batches

With a mini-batch $\{j_1, \dots, j_B\}$ ($B \ll N$) we get a less noisy, still **unbiased** estimate of the gradient:

$$\begin{aligned}\nabla_{\theta} \mathcal{L}(\theta) &:= \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \mathcal{L}_i(\theta) \approx \frac{1}{B} \sum_{i=1}^B \nabla_{\theta} \mathcal{L}_{j_i}(\theta) \\ &= \lambda \nabla_{\theta} \Omega(\theta) + \frac{1}{B} \sum_{i=1}^B \nabla_{\theta} L(f(\mathbf{x}_{j_i}; \theta), y_{j_i}).\end{aligned}$$

The weights $\theta = \{(\mathbf{W}^{(\ell)}, \mathbf{b}^{(\ell)})\}_{\ell=1}^{L+1}$ are then updated as:

$$\theta \leftarrow \theta - \eta \frac{1}{B} \sum_{i=1}^B \nabla_{\theta} \mathcal{L}_{j_i}(\theta)$$

The Key Ingredients of SGD

In sum, SGD needs the following ingredients:

- The **loss function** $L(f(\mathbf{x}_i; \boldsymbol{\theta}), y_i)$;
- A procedure for computing the **gradients** $\nabla_{\boldsymbol{\theta}} L(f(\mathbf{x}_i; \boldsymbol{\theta}), y_i)$;
- The **regularizer** $\Omega(\boldsymbol{\theta})$ and its **gradient**.

Let's see them one at the time...

Loss Function

Should match as much as possible the metric we want to optimize at test time

Should be well-behaved (continuous, maybe smooth) to be amenable to optimization (this rules out the 0/1 loss)

Some examples:

- Squared loss for regression
- Negative log-likelihood (cross-entropy) for multi-class classification
- Sparsemax loss for multi-class and multi-label classification

Squared Loss

- The common choice for regression/reconstruction problems
- The neural network aims at estimating $\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}) \approx \mathbf{y}$
- Minimize the **mean squared error**:

$$L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), \mathbf{y}) = \frac{1}{2} \|\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}) - \mathbf{y}\|^2$$

- Loss gradient:

$$\frac{\partial L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}, \mathbf{y}))}{\partial f_c(\mathbf{x}; \boldsymbol{\theta})} = f_c(\mathbf{x}; \boldsymbol{\theta}) - y_c$$

Negative Log-Likelihood (Cross-Entropy)

- The common choice for a softmax output layer
- The neural network estimates $f_c(\mathbf{x}; \boldsymbol{\theta}) \approx P(y = c | \mathbf{x})$
- We minimize the negative log-likelihood (also called **cross-entropy**):

$$\begin{aligned} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) &= - \sum_c \mathbf{1}_{(y=c)} \log f_c(\mathbf{x}; \boldsymbol{\theta}) \\ &= - \log f_y(\mathbf{x}; \boldsymbol{\theta}) \\ &= - \log \mathbf{softmax}_y(\mathbf{z}(\mathbf{x})) \end{aligned}$$

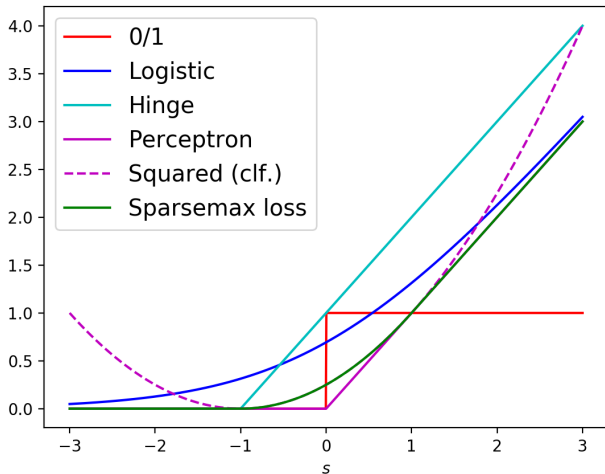
where \mathbf{z} is the **output pre-activation**.

- Loss gradient at output pre-activation:

$$\frac{\partial L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)}{\partial z_c} = \mathbf{softmax}_c(\mathbf{z}(\mathbf{x})) - \mathbf{1}_{(y=c)}$$

Classification Losses

- Let the correct label be $y = 1$ and define $s = z_2 - z_1$:



The Key Ingredients of SGD

In sum, SGD needs the following ingredients:

- The **loss function** $L(f(\mathbf{x}_i; \boldsymbol{\theta}), y_i)$;
- A procedure for computing the **gradients** $\nabla_{\boldsymbol{\theta}} L(f(\mathbf{x}_i; \boldsymbol{\theta}), y_i)$: **next**
- The **regularizer** $\Omega(\boldsymbol{\theta})$ and its **gradient**.

Gradient Computation

- Recall that we need to compute

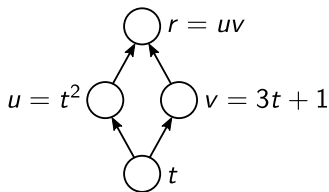
$$\nabla_{\theta} L(f(\mathbf{x}_i; \theta), y_i)$$

for $\theta = \{(\mathbf{W}^{(\ell)}, \mathbf{b}^{(\ell)})\}_{\ell=1}^{L+1}$ (the weights and biases at all layers)

- This will be done with the **gradient backpropagation algorithm**
- Key idea:** use the chain rule for derivatives!

$$h(x) = f(g(x)) \quad \Rightarrow \quad h'(x) = f'(g(x)) g'(x)$$

Recap: Chain Rule



$$\begin{aligned}\frac{\partial r(t)}{\partial t} &= ? \frac{\partial r(u)}{\partial u} \frac{\partial u(t)}{\partial t} + \frac{\partial r(v)}{\partial v} \frac{\partial v(t)}{\partial t} \\ &= 2tv + 3u \\ &= 2t(3t + 1) + 3t^2 = 9t^2 + 2t.\end{aligned}$$

- If a function $r(t)$ can be written as a function of intermediate results $q_i(t)$, then we have:

$$\frac{\partial r(t)}{\partial t} = \sum_i \frac{\partial r(t)}{\partial q_i(t)} \frac{\partial q_i(t)}{\partial t}$$

- We can invoke it by setting t to a output unit in a layer; $q_i(t)$ to the pre-activation in the layer above; and $r(t)$ to the loss function.

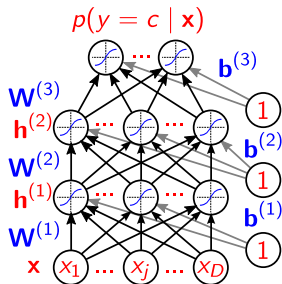
Hidden Layer Gradient

Main message: gradient backpropagation is just the chain rule of derivatives!

Hidden Layer Gradient

(Recap: $\mathbf{z}^{(\ell+1)} = \mathbf{W}^{(\ell+1)}\mathbf{h}^{(\ell)} + \mathbf{b}^{(\ell+1)}$)

$$\begin{aligned}\frac{\partial L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)}{\partial h_j^{(\ell)}} &= \sum_i \frac{\partial L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)}{\partial z_i^{(\ell+1)}} \frac{\partial z_i^{(\ell+1)}}{\partial h_j^{(\ell)}} \\ &= \sum_i \frac{\partial L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)}{\partial z_i^{(\ell+1)}} \mathbf{W}_{i,j}^{(\ell+1)}\end{aligned}$$

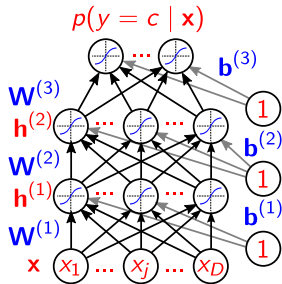


Hence $\nabla_{\mathbf{h}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = \mathbf{W}^{(\ell+1)\top} \nabla_{\mathbf{z}^{(\ell+1)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)$.

Hidden Layer Gradient (Before Activation)

(Recap: $h_j^{(\ell)} = g(z_j^{(\ell)})$, where $g : \mathbb{R} \rightarrow \mathbb{R}$ is the activation function)

$$\begin{aligned}\frac{\partial L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)}{\partial z_j^{(\ell)}} &= \frac{\partial L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)}{\partial h_j^{(\ell)}} \frac{\partial h_j^{(\ell)}}{\partial z_j^{(\ell)}} \\ &= \frac{\partial L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)}{\partial h_j^{(\ell)}} g'(z_j^{(\ell)})\end{aligned}$$



Hence $\nabla_{z^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = \nabla_{\mathbf{h}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) \odot \mathbf{g}'(z^{(\ell)})$.

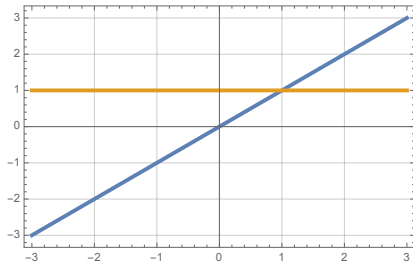
How to compute the derivative of the activation function $\mathbf{g}'(z^{(\ell)})$?

Linear Activation

$$g(z) = z$$

Derivative:

$$g'(z) = 1$$

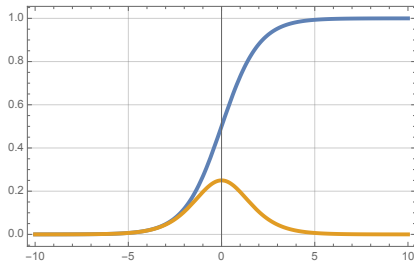


Sigmoid Activation

$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

Derivative:

$$g'(z) = g(z)(1 - g(z))$$

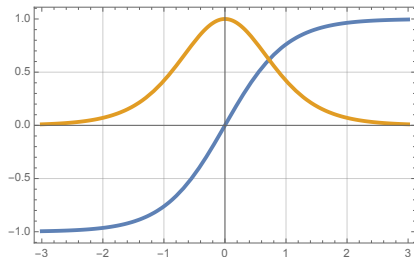


Hyperbolic Tangent Activation

$$g(z) = \mathbf{tanh}(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Derivative:

$$g'(z) = 1 - g(z)^2 = \mathbf{sech}^2(x)$$

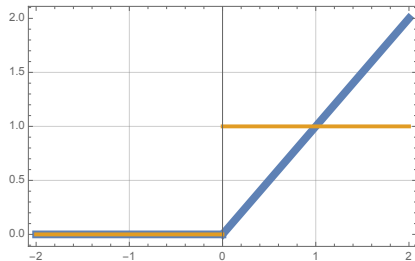


Rectified Linear Unit Activation (Glorot et al., 2011)

$$g(z) = \mathbf{relu}(z) = \max\{0, z\}$$

Derivative (except for $z = 0$):

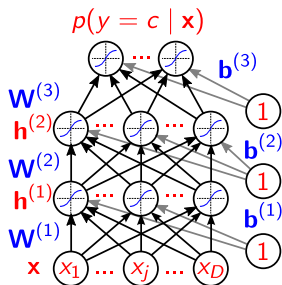
$$g'(z) = 1_{z>0}$$



Parameter Gradient

(Recap: $\mathbf{z}^{(\ell)} = \mathbf{W}^{(\ell)}\mathbf{h}^{(\ell-1)} + \mathbf{b}^{(\ell)}$)

$$\begin{aligned}\frac{\partial L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)}{\partial \mathbf{W}_{i,j}^{(\ell)}} &= \frac{\partial L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)}{\partial z_i^{(\ell)}} \frac{\partial z_i^{(\ell)}}{\partial \mathbf{W}_{i,j}^{(\ell)}} \\ &= \frac{\partial L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)}{\partial z_i^{(\ell)}} h_j^{(\ell-1)}\end{aligned}$$



Hence $\nabla_{\mathbf{W}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = \nabla_{\mathbf{z}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) \mathbf{h}^{(\ell-1)\top}$

Similarly, $\nabla_{\mathbf{b}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = \nabla_{\mathbf{z}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)$

Backpropagation

Compute output gradient (before activation):

$$\nabla_{\mathbf{z}^{(L+1)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = -(\mathbf{1}_y - \mathbf{f}(\mathbf{x}))$$

for ℓ from $L + 1$ to 1 **do**

 Compute gradients of hidden layer parameters:

$$\nabla_{\mathbf{W}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = \nabla_{\mathbf{z}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) \mathbf{h}^{(\ell-1)\top}$$

$$\nabla_{\mathbf{b}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = \nabla_{\mathbf{z}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)$$

 Compute gradient of hidden layer below:

$$\nabla_{\mathbf{h}^{(\ell-1)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = \mathbf{W}^{(\ell)\top} \nabla_{\mathbf{z}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)$$

 Compute gradient of hidden layer below (before activation):

$$\nabla_{\mathbf{z}^{(\ell-1)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = \nabla_{\mathbf{h}^{(\ell-1)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) \odot \mathbf{g}'(\mathbf{z}^{(\ell-1)})$$

end for

Conclusions

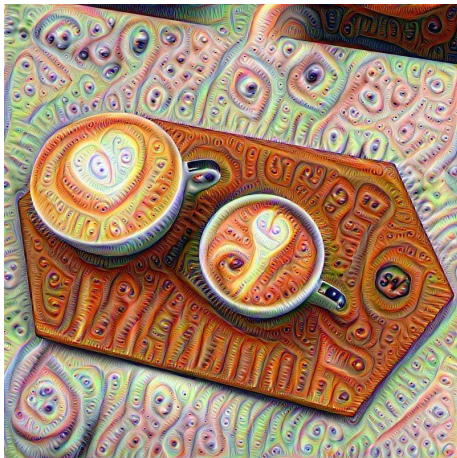
- Multi-layer perceptrons are universal function approximators
- However, they need to be trained
- Stochastic gradient descent is an effective training algorithm
- This is possible with the gradient backpropagation algorithm (an application of the chain rule of derivatives)

Next class:

- Most current software packages represent a computation graph and implement automatic differentiation
- Dropout regularization is effective to avoid overfitting
- Tricks of the trade

Thank you!

Questions?



References I

- Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314.
- Glorot, X., Bordes, A., and Bengio, Y. (2011). Deep Sparse Rectifier Neural Networks. In *International Conference on Artificial Intelligence and Statistics*, pages 315–323.
- Hornik, K., Stinchcombe, M., and White, H. (1989). Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366.
- McCulloch, W. S. and Pitts, W. (1943). A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5(4):115–133.
- Minsky, M. and Papert, S. (1969). Perceptrons.
- Montufar, G. F., Pascanu, R., Cho, K., and Bengio, Y. (2014). On the number of linear regions of deep neural networks. In *Advances in Neural Information Processing Systems*, pages 2924–2932.

Lecture 6: Neural Networks II

André Martins, Francisco Melo, Mário Figueiredo



Deep Learning Course, Winter 2021-2022

Today's Roadmap

Last lecture was about **neural networks**:

- From perceptron to multi-layer perceptron
- Feed-forward neural networks
- Activation functions: sigmoid, tanh, relu, ...
- Activation maps: softmax, sparsemax, ...
- Non-convex optimization and local minima
- Universal approximation theorem
- Gradient backpropagation

Today: **autodiff, regularization, tricks of the trade.**

Recap: Forward Propagation

Now assume $L \geq 1$ hidden layers:

- **Hidden layer pre-activation** (define $\mathbf{h}^{(0)} = \mathbf{x}$ for convenience):

$$\mathbf{z}^{(\ell)}(\mathbf{x}) = \mathbf{W}^{(\ell)}\mathbf{h}^{(\ell-1)} + \mathbf{b}^{(\ell)},$$

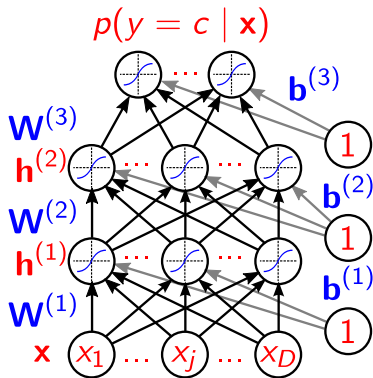
with $\mathbf{W}^{(\ell)} \in \mathbb{R}^{K_\ell \times K_{\ell-1}}$ and $\mathbf{b}^{(\ell)} \in \mathbb{R}^{K_\ell}$.

- **Hidden layer activation:**

$$\mathbf{h}^{(\ell)}(\mathbf{x}) = \mathbf{g}(\mathbf{z}^{(\ell)}(\mathbf{x})).$$

- **Output layer activation:**

$$\mathbf{f}(\mathbf{x}) = \mathbf{o}(\mathbf{z}^{(L+1)}(\mathbf{x})) = \mathbf{o}(\mathbf{W}^{(L+1)}\mathbf{h}^{(L)} + \mathbf{b}^{(L+1)}).$$



Recap: Gradient Backpropagation

Compute output gradient (before activation):

$$\nabla_{\mathbf{z}^{(L+1)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = -(\mathbf{1}_y - \mathbf{f}(\mathbf{x}))$$

for ℓ from $L + 1$ to 1 **do**

Compute gradients of hidden layer parameters:

$$\nabla_{\mathbf{W}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = \nabla_{\mathbf{z}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) \mathbf{h}^{(\ell-1)\top}$$

$$\nabla_{\mathbf{b}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = \nabla_{\mathbf{z}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)$$

Compute gradient of hidden layer below:

$$\nabla_{\mathbf{h}^{(\ell-1)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = \mathbf{W}^{(\ell)\top} \nabla_{\mathbf{z}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)$$

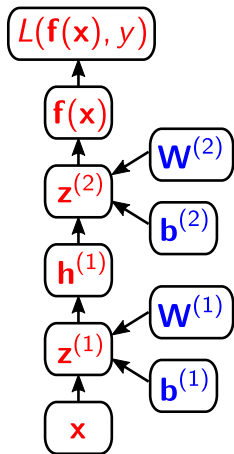
Compute gradient of hidden layer below (before activation):

$$\nabla_{\mathbf{z}^{(\ell-1)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = \nabla_{\mathbf{h}^{(\ell-1)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) \odot \mathbf{g}'(\mathbf{z}^{(\ell-1)})$$

end for

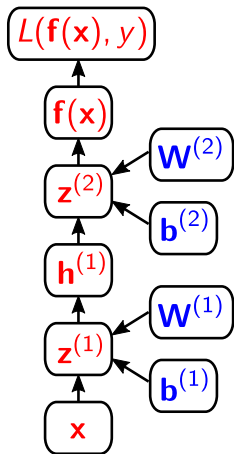
Computation Graph

- Forward propagation can be represented as a DAG
- Allows to implement forward propagation in a modular way
- Each box can be an object with a `fprop` method, that computes the value of the box given its parents/inputs
- Calling the `fprop` method of each box in the right order (after a topological sort) yields forward propagation



Automatic Differentiation

- ... Also allows to implement backpropagation in a modular way
- Each box can also have a `bprop` method, that computes the loss gradient with respect to its parents, given the loss gradient with respect to the output
- Can make use of cached computation done during the `fprop` method
- By calling the `bprop` method in reverse order, we get backpropagation (only need to reach the parameters)



Several Autodiff Strategies

Symbol-to-number differentiation (Caffe, Torch, Pytorch, Dynet, ...)

- Take a computational graph and a set of numerical inputs, then return a set of numerical values describing the gradient at those input values
- Advantage: simpler to implement and to debug
- Disadvantage: only works for first-order derivatives

Symbol-to-symbol differentiation (Theano, Tensorflow, ...)

- Take a computational graph and add additional nodes to the graph that provide a symbolic description of the desired derivatives (i.e. the derivatives are just another computational graph)
- Advantage: generalizes automatically to higher-order derivatives
- Disadvantage: harder to implement and to debug

Many Software Toolkits for Neural Networks

- Theano
- Tensorflow
- Torch, Pytorch
- MXNet
- Keras
- Caffe
- DyNet
- ...



All implement automatic differentiation.

We will have a Pytorch practical class this week

You may bring your laptops if you want to try it out!

Some Theano Code (Logistic Regression)

```
import numpy
import theano
import theano.tensor as T
rng = numpy.random

N = 400 # training sample size
feats = 784 # number of input variables

# generate a dataset: D = (input_values, target_class)
D = (rng.randn(N, feats), rng.randint(size=N, low=0, high=2))
training_steps = 10000

# Declare Theano symbolic variables
x = T.dmatrix("x")
y = T.dvector("y")

# initialize the weight vector w randomly
#
# this and the following bias variable b
# are shared so they keep their values
# between training iterations (updates)
w = theano.shared(rng.randn(feats), name="w")

# initialize the bias term
b = theano.shared(0., name="b")

print("Initial model:")
print(w.get_value())
print(b.get_value())

# Construct Theano expression graph
p_1 = 1 / (1 + T.exp(-T.dot(x, w) - b)) # Probability that target = 1
prediction = p_1 > 0.5 # The prediction thresholded
xent = -y * T.log(p_1) - (1-y) * T.log(1-p_1) # Cross-entropy loss function
cost = xent.mean() + 0.01 * (w ** 2).sum() # The cost to minimize
gw, gb = T.grad(cost, [w, b]) # Compute the gradient of the cost
# w.r.t weight vector w and
# bias term b
# (we shall return to this in a
# following section of this tutorial)

# Compile
train = theano.function(
    inputs=[x,y],
```

Some Code in Tensorflow (Linear Regression)

```
import tensorflow as tf
import numpy as np

# Create 100 phony x, y data points in NumPy,  $y = x * 0.1 + 0.3$ 
x_data = np.random.rand(100).astype(np.float32)
y_data = x_data * 0.1 + 0.3

# Try to find values for W and b that compute  $y\_data = W * x\_data + b$ 
# (We know that W should be 0.1 and b 0.3, but TensorFlow will
# figure that out for us.)
W = tf.Variable(tf.random_uniform([1], -1.0, 1.0))
b = tf.Variable(tf.zeros([1]))
y = W * x_data + b

# Minimize the mean squared errors.
loss = tf.reduce_mean(tf.square(y - y_data))
optimizer = tf.train.GradientDescentOptimizer(0.5)
train = optimizer.minimize(loss)

# Before starting, initialize the variables. We will 'run' this first.
init = tf.global_variables_initializer()

# Launch the graph.
sess = tf.Session()
sess.run(init)

# Fit the line.
for step in range(201):
    sess.run(train)
    if step % 20 == 0:
        print(step, sess.run(W), sess.run(b))

# Learns best fit is W: [0.1], b: [0.3]
```

Some Code in Keras (Multi-Layer Perceptron)

Multilayer Perceptron (MLP) for multi-class softmax classification:

```
from keras.models import Sequential
from keras.layers import Dense, Dropout, Activation
from keras.optimizers import SGD

model = Sequential()
# Dense(64) is a fully-connected layer with 64 hidden units.
# in the first layer, you must specify the expected input data shape:
# here, 20-dimensional vectors.
model.add(Dense(64, input_dim=20, init='uniform'))
model.add(Activation('tanh'))
model.add(Dropout(0.5))
model.add(Dense(64, init='uniform'))
model.add(Activation('tanh'))
model.add(Dropout(0.5))
model.add(Dense(10, init='uniform'))
model.add(Activation('softmax'))

sgd = SGD(lr=0.1, decay=1e-6, momentum=0.9, nesterov=True)
model.compile(loss='categorical_crossentropy',
              optimizer=sgd,
              metrics=['accuracy'])

model.fit(X_train, y_train,
          nb_epoch=20,
          batch_size=16)
score = model.evaluate(X_test, y_test, batch_size=16)
```

Some Code in Pytorch (Multi-Layer Perceptron)

```
# Fully connected neural network with one hidden layer
class NeuralNet(nn.Module):
    def __init__(self, input_size, hidden_size, num_classes):
        super(NeuralNet, self).__init__()
        self.fc1 = nn.Linear(input_size, hidden_size)
        self.relu = nn.ReLU()
        self.fc2 = nn.Linear(hidden_size, num_classes)

    def forward(self, x):
        out = self.fc1(x)
        out = self.relu(out)
        out = self.fc2(out)
        return out

model = NeuralNet(input_size, hidden_size, num_classes).to(device)

# Loss and optimizer
criterion = nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(model.parameters(), lr=learning_rate)

# Train the model
total_step = len(train_loader)
for epoch in range(num_epochs):
    for i, (images, labels) in enumerate(train_loader):
        # Move tensors to the configured device
        images = images.reshape(-1, 28*28).to(device)
        labels = labels.to(device)

        # Forward pass
        outputs = model(images)
        loss = criterion(outputs, labels)

        # Backward and optimize
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()

    if (i+1) % 100 == 0:
        print ('Epoch [{} / {}], Step [{} / {}], Loss: {:.4f}'
              .format(epoch+1, num_epochs, i+1, total_step, loss.item()))
```

Reminder: The Key Ingredients of SGD

In sum, we need the following ingredients:

- The loss function $L(f(\mathbf{x}_i; \boldsymbol{\theta}), y_i)$;
- A procedure for computing the gradients $\nabla_{\boldsymbol{\theta}} L(f(\mathbf{x}_i; \boldsymbol{\theta}), y_i)$
- The regularizer $\Omega(\boldsymbol{\theta})$ and its gradient – next!

Regularization

Recall that we're minimizing the following objective function:

$$\mathcal{L}(\boldsymbol{\theta}) := \lambda\Omega(\boldsymbol{\theta}) + \frac{1}{N} \sum_{n=1}^N L(\mathbf{f}(\mathbf{x}_i; \boldsymbol{\theta}), y_i)$$

It remains to define the **regularizer** and its gradient

We'll talk about:

- ℓ_2 regularization
- ℓ_1 regularization
- dropout regularization

ℓ_2 Regularization

- **Gaussian prior** on the weights
- **Note:** only the weights are regularized (not the biases)

$$\Omega(\boldsymbol{\theta}) = \frac{1}{2} \sum_{\ell} \|\mathbf{W}^{(\ell)}\|^2$$

- Gradient is:

$$\nabla_{\mathbf{W}^{(\ell)}} \Omega(\boldsymbol{\theta}) = \mathbf{W}^{(\ell)}$$

- This has the effect of a weight decay:

$$\begin{aligned} \mathbf{W}^{(\ell)} &\leftarrow \mathbf{W}^{(\ell)} - \eta \nabla_{\mathbf{W}^{(\ell)}} \mathcal{L}_i(\boldsymbol{\theta}) \\ &= \mathbf{W}^{(\ell)} - \eta (\lambda \nabla_{\mathbf{W}^{(\ell)}} \Omega(\boldsymbol{\theta}) + \nabla_{\mathbf{W}^{(\ell)}} L(f(\mathbf{x}_i; \boldsymbol{\theta}), y_i)) \\ &= (1 - \eta\lambda) \mathbf{W}^{(\ell)} - \eta \nabla_{\mathbf{W}^{(\ell)}} L(f(\mathbf{x}_i; \boldsymbol{\theta}), y_i) \end{aligned}$$

ℓ_1 Regularization

- **Laplacian prior** on the weights
- **Note:** only the weights are regularized (not the biases)

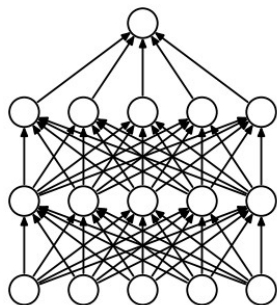
$$\Omega(\boldsymbol{\theta}) = \sum_{\ell} \|\mathbf{W}^{(\ell)}\|_1$$

- Gradient is:

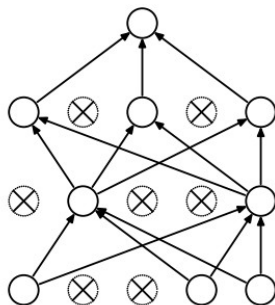
$$\nabla_{\mathbf{W}^{(\ell)}} \Omega(\boldsymbol{\theta}) = \text{sign}(\mathbf{W}^{(\ell)})$$

- Promotes sparsity of the weights

Dropout Regularization (Srivastava et al., 2014)



(a) Standard Neural Net



(b) After applying dropout.

Idea: During training, remove some hidden units stochastically

Dropout Regularization (Srivastava et al., 2014)

- Each hidden unit's output is set to 0 with probability p (e.g. $p = 0.5$)
- This prevents hidden units to co-adapt to other units, forcing them to be more generally useful
- At test time, keep all units, but multiply their outputs by $1 - p$
- Shown to be a form of adaptive regularization (Wager et al., 2013)
- Note: many software packages implement another variant, **inverted dropout**, where at training time the output of the units that were not dropped is divided by $1 - p$ and requires no change at test time

Implementation of Dropout

- This is usually implemented using random binary masks
- The hidden layer activations become (for $\ell = 1, \dots, L$):

$$\mathbf{h}^{(\ell)}(\mathbf{x}) = \mathbf{g}(\mathbf{z}^{(\ell)}(\mathbf{x})) \odot \mathbf{m}^{(\ell)}$$

- Beats regular backpropagation on many datasets (Hinton et al., 2012)
- Other variants, e.g. DropConnect (Wan et al., 2013), Stochastic Pooling (Zeiler and Fergus, 2013)

Backpropagation with Dropout

Compute output gradient (before activation):

$$\nabla_{\mathbf{z}^{(L+1)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = -(\mathbf{1}_y - \mathbf{f}(\mathbf{x}))$$

for ℓ from $L + 1$ to 1 **do**

Compute gradients of hidden layer parameters:

$$\nabla_{\mathbf{W}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = \nabla_{\mathbf{z}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) \underbrace{\mathbf{h}^{(\ell-1)\top}}_{\text{includes } \mathbf{m}^{(\ell-1)}}$$

$$\nabla_{\mathbf{b}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = \nabla_{\mathbf{z}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)$$

Compute gradient of hidden layer below:

$$\nabla_{\mathbf{h}^{(\ell-1)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = \mathbf{W}^{(\ell)\top} \nabla_{\mathbf{z}^{(\ell)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y)$$

Compute gradient of hidden layer below (before activation):

$$\nabla_{\mathbf{z}^{(\ell-1)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) = \nabla_{\mathbf{h}^{(\ell-1)}} L(\mathbf{f}(\mathbf{x}; \boldsymbol{\theta}), y) \odot \mathbf{g}'(\mathbf{z}^{(\ell-1)}) \odot \mathbf{m}^{(\ell-1)}$$

end for

Initialization

Initialize all biases to zero

For weights:

- Cannot initialize to zero with **tanh** activation (the gradients would also be zero and we would reach a saddle point)
- Cannot initialize the weights to the same value (need to break the symmetry)
- Random initialization (Gaussian, uniform), sampling around 0 to break symmetry
- For ReLU activations, the mean should be a small positive number
- Variance cannot be too high, otherwise all neuron activations will be saturated

“Glorot Initialization”

- Recipe from Glorot and Bengio (2010):

$$\mathbf{w}_{i,j}^{(\ell)} \sim U[-t, t], \text{ with } t = \frac{\sqrt{6}}{\sqrt{K^{(\ell)} + K^{(\ell-1)}}}$$

- Works well in practice with **tanh** and sigmoid activations

Training, Validation, and Test Sets

Split datasets in training, validation, and test partitions.

- Training set serves to train the model
- Validation set serves to tune hyperparameters (learning rate, number of hidden units, regularization coefficient, dropout probability, best epoch, etc.)
- Test set serves to estimate the generalization performance

Hyperparameter Tuning: Grid Search, Random Search

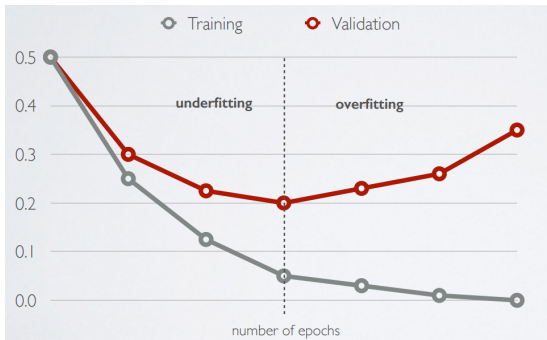
Search for the best configuration of the hyperparameters:

- Grid search: specify a set of values we want to test for each hyperparameter, and try all configurations of these values
- Random search: specify a distribution over the values of each hyper-parameter (e.g. uniform in some range) and sample independently each hyper-parameter to get configurations
- Bayesian optimization and learning to learn (Snoek et al., 2012)

We can always go back and fine-tune the grid/distributions if necessary

Early Stopping

- To select the number of epochs, stop training when validation error increases (with some look ahead)
- One common strategy (with SGD) is to halve the learning rate for every epoch where the validation error increases

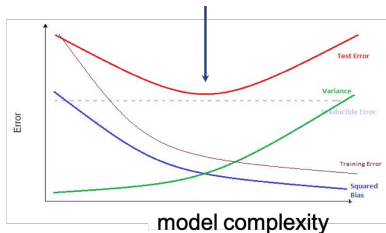


(Image credit: Hugo Larochelle)

Cross Validation

Model selection

finding the sweet spot (F)



validation set:

- ✓ split data into **train** and **validation** subsets
- ✓ train on the training subset
- ✓ test on the validation subset

cross validation (k-fold): repeat k times; average

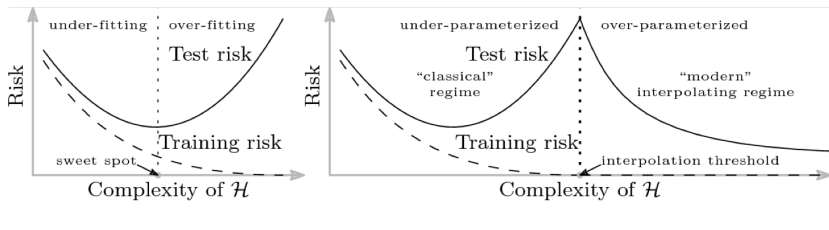


Over-parametrization

The new regime...

modern deep networks have “too many” parameters: they should overfit, ...

...yet, they usually don't. Why? Ongoing research.



(illustration by Mikhail Belkin)

Tricks of the Trade

- Normalization of the data
- Decaying the learning rate
- Mini-batches
- Adaptive learning rates
- Gradient checking
- Debugging on a small dataset

Normalization of the Data

- For each input dimension: subtract the training set mean and divide by the training set standard deviation
- This makes each input dimension have zero mean, unit variance
- This can speed up training (in number of epochs)
- Doesn't work for sparse inputs (destroys sparsity)

Decaying the Learning Rate

In SGD, as we get closer to a local minimum, it makes sense to take smaller update steps (to avoid diverging)

- Start with a large learning rate (say 0.1)
- Keep it fixed while validation error keeps improving
- Divide by 2 and go back to the previous step

Mini-Batches

- Instead of updating after a single example, can aggregate a mini-batch of examples (e.g. 50–200 examples) and compute the averaged gradient for the entire mini-batch
- Less noisy than vanilla SGD
- Can leverage matrix-matrix computations (or tensor computations)
- Large computational speed-ups in GPUs (since computation is trivially parallelizable across the mini-batch and we can exhaust the GPU memory)

Adaptive Learning Rates

Instead of using the same step size for all parameters, have one learning rate per parameter

- **Adagrad** (Duchi et al., 2011): learning rates are scaled by the square root of the cumulative sum of squared gradients

$$\eta^{(t)} = \eta^{(t-1)} + (\nabla_{\theta} L(\mathbf{f}(\mathbf{x}), y))^2, \quad \bar{\nabla}_{\theta}^{(t)} = \frac{\nabla_{\theta} L(\mathbf{f}(\mathbf{x}), y)}{\sqrt{\eta^{(t)} + \epsilon}}$$

- **RMSprop** (Tieleman and Hinton, 2012): instead of cumulative sum, use exponential moving average

$$\eta^{(t)} = \beta \eta^{(t-1)} + (1 - \beta)(\nabla_{\theta} L(\mathbf{f}(\mathbf{x}), y))^2, \quad \bar{\nabla}_{\theta}^{(t)} = \frac{\nabla_{\theta} L(\mathbf{f}(\mathbf{x}), y)}{\sqrt{\eta^{(t)} + \epsilon}}$$

- **Adam** (Kingma and Ba, 2014): combine RMSProp with momentum

Gradient Checking

- If the training loss is not decreasing even with a very small learning rate, there's likely a bug in the gradient computation
- To debug your implementation of `fprop`/`bprop`, compute the “numeric gradient,” a finite difference approximation of the true gradient:

$$\frac{\partial f(x)}{\partial x} \approx \frac{f(x + \epsilon) - f(x - \epsilon)}{2\epsilon}$$

Debugging on a Small Dataset

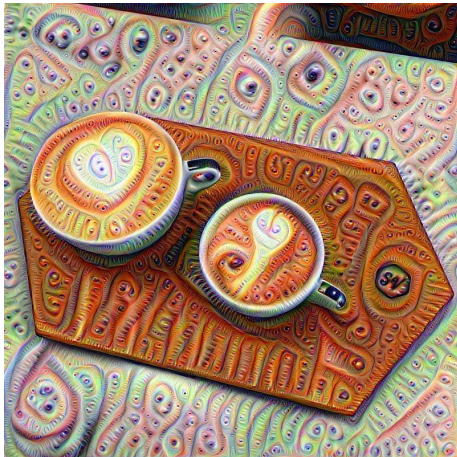
- Extract a small subset of your training set (e.g. 50 examples)
- Monitor your training loss in this set
- You should be able to overfit in this small training set
- If not, see if some units are saturated from the very first iterations (if they are, reduce the initialization variance or properly normalize your inputs)
- If the training error is bouncing up and down, decrease the learning rate

Conclusions

- Multi-layer perceptrons are universal function approximators
- However, they need to be trained
- Stochastic gradient descent is an effective training algorithm
- This is possible with the gradient backpropagation algorithm (an application of the chain rule of derivatives)
- Most current software packages represent a computation graph and implement automatic differentiation
- Dropout regularization is effective to avoid overfitting

Thank you!

Questions?



References I

- Duchi, J., Hazan, E., and Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159.
- Glorot, X. and Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, volume 9, pages 249–256.
- Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. R. (2012). Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*.
- Kingma, D. and Ba, J. (2014). Adam: A Method for Stochastic Optimization. In *Proc. of International Conference on Learning Representations*.
- Snoek, J., Larochelle, H., and Adams, R. P. (2012). Practical bayesian optimization of machine learning algorithms. In *Advances in neural information processing systems*, pages 2951–2959.
- Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958.
- Tieleman, T. and Hinton, G. (2012). Rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural Networks for Machine Learning*, 4(2).
- Wager, S., Wang, S., and Liang, P. S. (2013). Dropout training as adaptive regularization. In *Advances in neural information processing systems*, pages 351–359.
- Wan, L., Zeiler, M., Zhang, S., Cun, Y. L., and Fergus, R. (2013). Regularization of neural networks using dropconnect. In *Proc. of the International Conference on Machine Learning*, pages 1058–1066.
- Zeiler, M. D. and Fergus, R. (2013). Stochastic pooling for regularization of deep convolutional neural networks. *arXiv preprint arXiv:1301.3557*.

Lecture 7: Representation Learning

André Martins, Francisco Melo, Mário Figueiredo



Deep Learning Course, Winter 2021-2022

Announcements

Deadline for Homework #1 is this **Wednesday end of day!**

- Please submit your solutions and code in Fenix.
- No late days allowed!!
- Solutions will be posted the day after.

Homework #2 will be posted this **Wednesday!**

- Deadline Jan 31.
- Start early!!

Today's Roadmap

Today's lecture is about:

- Representation learning.
- Principal component analysis (PCA) and auto-encoders.
- Denoising auto-encoders.
- Distributed representations.
- Word embeddings and negative sampling.
- Multilingual and contextual word embeddings.

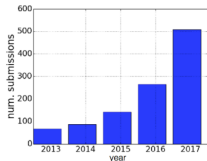
Representations

- A key feature of NNs is their ability to **learn representations** $\phi(\mathbf{x})$
- Standard linear models require **manually engineered features** $\phi(\mathbf{x})$
- **Representations** are useful for several reasons:
 - (i) They can make our models **more expressive and more accurate**
 - (ii) They may allow **transferring** representations from one task to another
- We talked about (i) when discussing the multi-layer perceptron
- In this lecture, we'll focus on (ii)

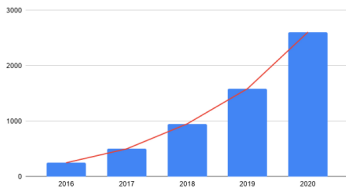
Representation Learning

This is becoming an extremely popular topic!

Number of submissions to the “International Conference on Learning Representations” (ICLR):



ICLR Submissions



Representation learning almost became a synonym of deep learning

Hierarchical Compositionality

Key Idea: deep(er) NNs learn **coarse-to-fine** representation layers.

Vision:

- pixels → edges → textures → motifs → parts → objects → scenes

Speech:

- audio samples → spectral bands → formants → motifs → phonemes → words

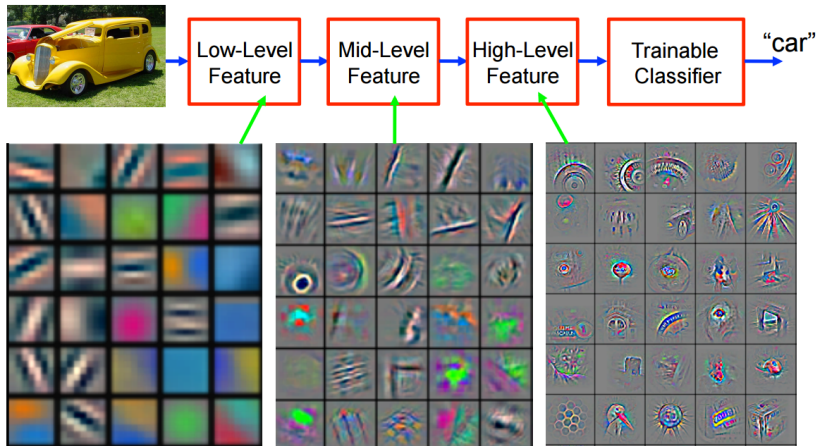
Text:

- characters → words → phrases → sentences → stories

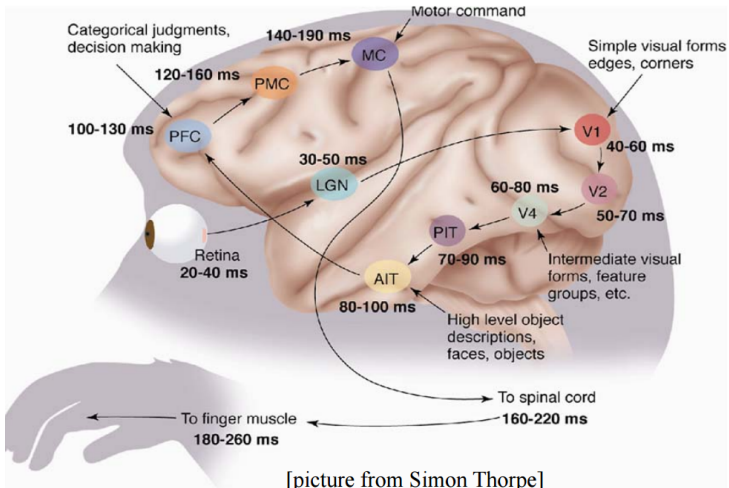
(Slide inspired by Marc'Aurelio Ranzato and Yann LeCun)

Hierarchical Compositionality

Feature visualization of convolutional NNs trained on ImageNet (Zeiler and Fergus, 2013):



The Mammalian Visual Cortex is Hierarchical



(Slide inspired by Marc'Aurelio Ranzato and Yann LeCun)

What's in Each Layer

- Bottom level layers (closer to inputs) tend to learn **low-level representations** (corners, edges)
- Upper level layers (farther away from inputs) learn **more abstract representations** (shapes, forms, objects)

This holds for images, text, etc.

Distributed Representations (Hinton, 1984)

This is a central concept in neural networks.

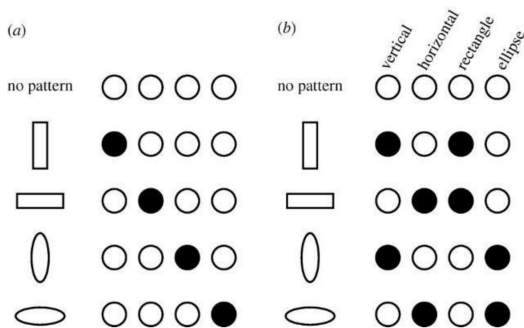
Key questions:

- How can a NN so effectively **represent objects**, if it has only a few **hidden units** (i.e. much fewer than possible objects)?
- What is each **hidden unit** actually **representing**?
- How can a NN **generalize** to objects that is has not seen before?

Local vs Distributed Representations

Consider two alternative representations:

- **Local** (one-hot) representations (one dimension per object)
- **Distributed representations** (one dimension per **property**)



(Slide inspired by Moontae Lee and Dhruv Batra)

Distributed Representations

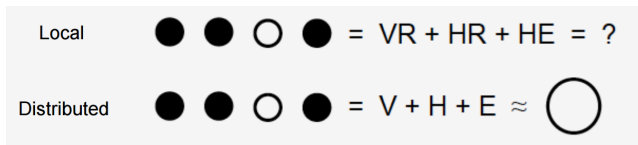
Key idea: no single neuron “encodes” everything; groups of neurons (e.g. in the same hidden layer) work together!

cf. the [grandmother cell](#)



The Power of Distributed Representations

- Distributed representations are **more compact** (there can be $O(\exp N)$ objects combining N properties)
- They are also **more powerful**, as they can generalize to unseen objects in a meaningful way:



(Slide inspired by Moontae Lee and Dhruv Batra)

The Power of Distributed Representations

- For this to work, hidden units should capture **diverse properties** of objects (not all capturing the same property)
- Usually ensured by random initialization of the weights
- Initializing all the units to the same weights, we would never break the symmetry!
- **Side note:** a NN computes the same function if we permute the hidden units within the same layer (order doesn't matter, only diversity)

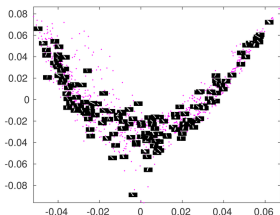
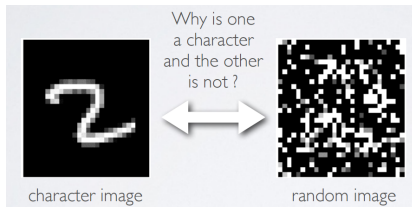
Next: how to learn useful object representations from raw inputs (no labels)?

Example: Unsupervised Pre-Training

- Training deep NNs (with many hidden layers) can be challenging
- This has been a major difficulty with NNs for a long time
- Initialize hidden layers using **unsupervised learning** (Erhan et al., 2010):
 - Force network to **represent latent structure** of input distribution
 - Encourage hidden layers to **encode** that structure
 - This can be done with an **auto-encoder!**

Data Manifold

Key idea: learn the **manifold** where the input objects live



(Image credit: Hugo Larochelle)

Learn representations that encode well points in that **manifold**

Auto-Encoders

Auto-encoder: feed-forward NN trained to reproduce its input at the output

Encoder:

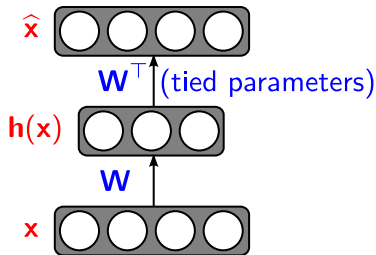
$$\mathbf{h}(\mathbf{x}) = \mathbf{g}(\mathbf{W}\mathbf{x} + \mathbf{b})$$

Decoder:

$$\hat{\mathbf{x}} = \mathbf{W}^{\top} \mathbf{h}(\mathbf{x}) + \mathbf{c}$$

Loss function (for real-valued inputs):

$$L(\hat{\mathbf{x}}; \mathbf{x}) = \frac{1}{2} \|\hat{\mathbf{x}} - \mathbf{x}\|^2$$



The Simplest Auto-Encoder

What happens if the activation function g is linear?

Principal Component Analysis (PCA)!

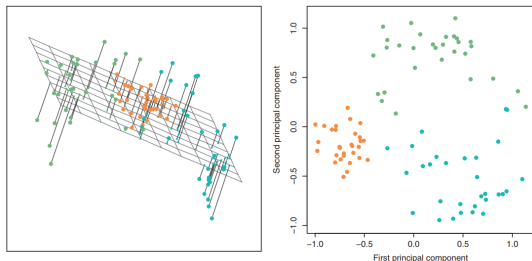


FIGURE 10.2. *Ninety observations simulated in three dimensions. Left: the first two principal component directions span the plane that best fits the data. It minimizes the sum of squared distances from each point to the plane. Right: the first two principal component score vectors give the coordinates of the projection of the 90 observations onto the plane. The variance in the plane is maximized.*

(From “An Introduction to Statistical Learning” by James, Witten, Hastie, Tibshirani)

Interlude: PCA

Consider N points in \mathbb{R}^D : $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$

Goal: find **good** K -dimensional representation of the points, with $K < D$,

$$\mathbf{x}^{(i)} \simeq \hat{\mathbf{x}}^{(i)} = \sum_{j=1}^K \alpha_j^{(i)} \mathbf{u}_j,$$

where $\{\mathbf{u}_1, \dots, \mathbf{u}_K\} \subset \mathbb{R}^D$ is an **orthonormal basis** of a **subspace** of \mathbb{R}^D .

If $\{\mathbf{u}_1, \dots, \mathbf{u}_K\}$ is fixed, and the approximation is in Euclidean norm,

$$\hat{\mathbf{x}}^{(i)} = \arg \min_{\mathbf{x}'} \|\mathbf{x}^{(i)} - \mathbf{x}'\|_2^2, \quad \text{subject to } \mathbf{x}' \in \text{span}(\{\mathbf{u}_1, \dots, \mathbf{u}_K\})$$

has a well-known solution: **orthogonal projection**,

$$\hat{\mathbf{x}}^{(i)} = \sum_{j=1}^K \mathbf{u}_j \mathbf{u}_j^T \mathbf{x}^{(i)} = \mathbf{U} \mathbf{U}^T \mathbf{x}^{(i)}$$

where $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_K] \in \mathbb{R}^{D \times K}$; notice that $\mathbf{U}^T \mathbf{U} = \mathbf{I}$.

Interlude: PCA

Minimizing the average error for the N points in \mathbb{R}^D : $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$:

$$\min_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}^{(i)} - \mathbf{U}\mathbf{U}^T \mathbf{x}^{(i)}\|_2^2$$

Notice that

$$\|\mathbf{x}^{(i)} - \mathbf{U}\mathbf{U}^T \mathbf{x}^{(i)}\|_2^2 = ((\mathbf{x}^{(i)})^T - (\mathbf{x}^{(i)})^T \mathbf{U}\mathbf{U}^T)(\mathbf{x}^{(i)} - \mathbf{U}\mathbf{U}^T \mathbf{x}^{(i)})$$

$$\stackrel{(a)}{=} \|\mathbf{x}^{(i)}\|_2^2 - (\mathbf{x}^{(i)})^T \mathbf{U}\mathbf{U}^T \mathbf{x}^{(i)}$$

$$= \|\mathbf{x}^{(i)}\|_2^2 - \sum_{j=1}^K \underbrace{(\mathbf{x}^{(i)})^T \mathbf{u}_j}_{\text{scalar}} \underbrace{\mathbf{u}_j^T \mathbf{x}^{(i)}}_{\text{scalar}}$$

indep. of \mathbf{U}

$$= \|\mathbf{x}^{(i)}\|_2^2 - \sum_{j=1}^K \mathbf{u}_j^T \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T \mathbf{u}_j$$

$$(a) \text{ since } -2(\mathbf{x}^{(i)})^T \mathbf{U}\mathbf{U}^T \mathbf{x}^{(i)} + \underbrace{(\mathbf{x}^{(i)})^T \mathbf{U}\mathbf{U}^T \mathbf{U}\mathbf{U}^T \mathbf{U}\mathbf{U}^T \mathbf{x}^{(i)}}_{\mathbf{I}} = -(\mathbf{x}^{(i)})^T \mathbf{U}\mathbf{U}^T \mathbf{x}^{(i)}$$

Interlude: PCA

Minimizing the average error for the N points in \mathbb{R}^D : $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$:

$$\begin{aligned} \arg \min_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}^{(i)} - \mathbf{U} \mathbf{U}^T \mathbf{x}^{(i)}\|_2^2 &= \arg \max_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^K \mathbf{u}_j^T \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T \mathbf{u}_j \\ &= \arg \max_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \sum_{j=1}^K \mathbf{u}_j^T \underbrace{\left(\frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T \right)}_{\hat{\Sigma}} \mathbf{u}_j \\ &= \arg \max_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \sum_{j=1}^K \mathbf{u}_j^T \hat{\Sigma} \mathbf{u}_j \end{aligned}$$

$\hat{\Sigma}$ is the sample covariance, assuming **centred** data: $\sum_{i=1}^N \mathbf{x}^{(i)} = \mathbf{0}$.

PCA

Let's start with $K = 1$,

$$\mathbf{u}_1 = \arg \max_{\|\mathbf{u}\|_2=1} \mathbf{u}^T \hat{\Sigma} \mathbf{u}$$

Lagrangian:

$$L(\mathbf{u}, \lambda) = \mathbf{u}^T \hat{\Sigma} \mathbf{u} + \lambda(1 - \|\mathbf{u}\|_2^2)$$

Setting the gradient to zero: $2\hat{\Sigma}\mathbf{u} - 2\lambda\mathbf{u} = 0 \Rightarrow \hat{\Sigma}\mathbf{u} = \lambda\mathbf{u}$,

i.e., \mathbf{u} is an eigenvector of $\hat{\Sigma}$. Plugging above:

$$\mathbf{u}^T \hat{\Sigma} \mathbf{u} = \mathbf{u}^T \mathbf{u} \lambda = \|\mathbf{u}\|_2^2 \lambda = \lambda$$

Conclusion: \mathbf{u}_1 is the eigenvector of the largest eigenvalue of $\hat{\Sigma}$.

Notice: $\hat{\Sigma}$ is positive semi-definite, thus all eigenvalues are non-negative.

Easy to extend to $K > 1$.

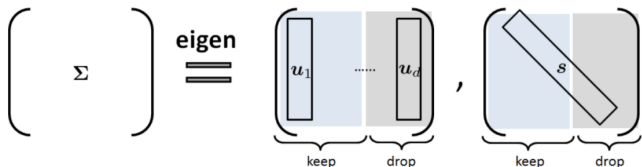
Interlude: PCA

$$\arg \max_{\mathbf{U}^T \mathbf{U} = \mathbf{I}} \sum_{j=1}^K \mathbf{u}_j^T \hat{\Sigma} \mathbf{u}_j = \text{eigenvectors of the top } K \text{ eigenvalues of } \hat{\Sigma}$$

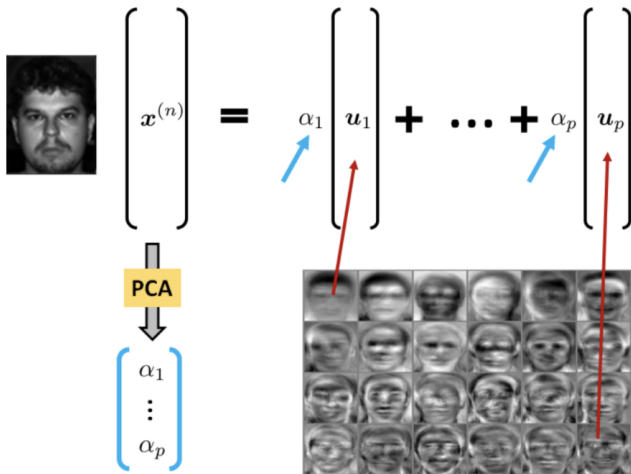
Recall that $\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T = \frac{1}{N} \mathbf{X}^T \mathbf{X}$

Notice that the eigenvalues of $\hat{\Sigma}$ are the **singular values** (SV) of \mathbf{X}/\sqrt{N} .

PCA corresponds to SVD (SV decomposition) of the scaled data matrix \mathbf{X}



PCA: EigenFaces



Back to the Linear Auto-Encoder

Let $\mathbf{X} \in \mathbb{R}^{N \times D}$ be the data matrix ($N > D$),

Assume $\mathbf{W} \in \mathbb{R}^{K \times D}$ with $K < D$ (no biases, assuming \mathbf{X} is centred)

We want to minimize, with g identity, thus $\hat{\mathbf{x}}^{(i)} = \mathbf{W}^\top \mathbf{W} \mathbf{x}^{(i)}$,

$$\sum_{i=1}^N \|\mathbf{x}^{(i)} - \hat{\mathbf{x}}^{(i)}\|_2^2 = \|\mathbf{X} - \hat{\mathbf{X}}\|_F^2 = \|\mathbf{X} - \mathbf{X} \mathbf{W}^\top \mathbf{W}\|_F^2$$

where $\|\cdot\|_F^2$ is the Frobenius matrix norm and $\mathbf{W}^\top \mathbf{W}$ has rank K .

From the Eckart-Young theorem, the minimizer is **truncated SVD** of \mathbf{X}^\top :

$$\hat{\mathbf{X}}^\top = \mathbf{U}_K \Sigma_K \mathbf{V}_K^\top,$$

where Σ_K is a diagonal matrix containing the top K singular values of \mathbf{X}^\top , and the columns of \mathbf{U}_K are the corresponding left singular vectors.

The solution is $\mathbf{W} = \mathbf{U}_K^\top$, which gives as desired:

$$\hat{\mathbf{X}}^\top = \mathbf{W}^\top \mathbf{W} \mathbf{X}^\top = \mathbf{U}_K \mathbf{U}_K^\top \mathbf{U} \Sigma \mathbf{V}^\top = \mathbf{U}_K \Sigma_K \mathbf{V}_K^\top.$$

Conclusion: the optimal linear auto-encoder coincides with PCA.

Auto-Encoders

PCA fits a **linear manifold** (affine space) to the data

By using **non-linear** activations, we obtain more sophisticated codes (i.e. representations).

We need some sort of regularization to:

- encourage a smooth representation (small perturbations of the input will lead to similar codes)
- avoid overfitting to the training data

Some Variants of Auto-Encoders

- **Sparse auto-encoders:** use many hidden units, but add a ℓ_1 regularization term to encourage **sparse representations** of the input
- **Denoising auto-encoders:** regularize by adding noise to the input; the goal is to learn a **smooth representation function** that allows to output the denoised input (inspired by image denoising)
- **Stacked auto-encoders:** several auto-encoders on top of each other
- **Variational auto-encoders:** a generative probabilistic model that minimizes a variational bound (this will be covered in another lecture!)

Regularized Auto-Encoders

- To regularize auto-encoders, **regularization** may be added to the loss
- The goal is then to minimize $L(\hat{\mathbf{x}}; \mathbf{x}) + \Omega(\mathbf{h}, \mathbf{x})$
- For example:
 - regularizing the code $\Omega(\mathbf{h}, \mathbf{x}) = \lambda \|\mathbf{h}\|^2$
 - regularizing the derivatives $\Omega(\mathbf{h}, \mathbf{x}) = \lambda \sum_i \|\nabla_{\mathbf{x}} h_i\|^2$
- The encoder and decoder parameters may be shared or not.

Sparse Auto-Encoders

- Most auto-encoders learn low-dimensional codes, e.g., they reduce input dimensionality (bottleneck shape $K < D$).
- But one exception are **sparse auto-encoders**:
 - Sparse auto-encoders incorporate a sparsity penalty $\Omega(\mathbf{h})$ on the code layer, e.g., $\Omega(\mathbf{h}) = \lambda \|\mathbf{h}\|_1$
 - Typically the number of hidden units is large, e.g., larger than the input dimension
 - The sparsity penalty encourages sparse codes, where most hidden units are inactive.

Stochastic Auto-Encoders

- In this case, the encoder and decoder are not deterministic functions, but involve some noise injection
- We have a distribution $p_{\text{encoder}}(\mathbf{h} \mid \mathbf{x})$ for the encoder and a distribution $p_{\text{decoder}}(\mathbf{x} \mid \mathbf{h})$ for the decoder
- The auto-encoder can be trained to minimize

$$-\log p_{\text{decoder}}(\mathbf{x} \mid \mathbf{h}).$$

Denoising Auto-Encoders

- Use a perturbed version of the input, $\tilde{\mathbf{x}} = \mathbf{x} + \mathbf{n}$, where \mathbf{n} is random noise (e.g. Gaussian noise $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$)
- Instead of minimizing $\frac{1}{2} \|\hat{\mathbf{x}} - \mathbf{x}\|^2$, minimize $\frac{1}{2} \|\hat{\mathbf{x}} - \tilde{\mathbf{x}}\|^2$
- This is a form of implicit regularization that ensures **smoothness**: it forces the system to represent well not only the data points, but also their perturbations

Denoising Auto-Encoders

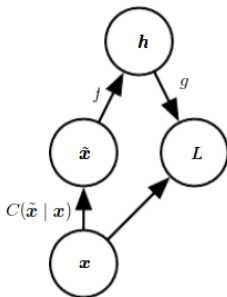


Figure 14.3: The computational graph of the cost function for a denoising autoencoder, which is trained to reconstruct the clean data point \mathbf{x} from its corrupted version $\tilde{\mathbf{x}}$. This is accomplished by minimizing the loss $L = -\log p_{\text{decoder}}(\mathbf{x} | \mathbf{h} = f(\tilde{\mathbf{x}}))$, where $\tilde{\mathbf{x}}$ is a corrupted version of the data example \mathbf{x} , obtained through a given corruption process $C(\tilde{\mathbf{x}} | \mathbf{x})$. Typically the distribution p_{decoder} is a factorial distribution whose mean parameters are emitted by a feedforward network g .

(From Goodfellow et al.'s book.)

Denoising Auto-Encoders

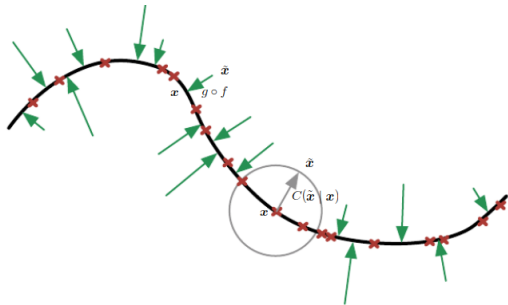


Figure 14.4: A denoising autoencoder is trained to map a corrupted data point $\tilde{\mathbf{x}}$ back to the original data point \mathbf{x} . We illustrate training examples \mathbf{x} as red crosses lying near a low-dimensional manifold, illustrated with the bold black line. We illustrate the corruption process $C(\tilde{\mathbf{x}} | \mathbf{x})$ with a gray circle of equiprobable corruptions. A gray arrow demonstrates how one training example is transformed into one sample from this corruption process. When the denoising autoencoder is trained to minimize the average of squared errors $\|g(f(\tilde{\mathbf{x}})) - \mathbf{x}\|^2$, the reconstruction $g(f(\tilde{\mathbf{x}}))$ estimates $\mathbb{E}_{\mathbf{x}, \tilde{\mathbf{x}} \sim p_{\text{data}}(\mathbf{x}) C(\tilde{\mathbf{x}} | \mathbf{x})}[\mathbf{x} | \tilde{\mathbf{x}}]$. The vector $g(f(\tilde{\mathbf{x}})) - \tilde{\mathbf{x}}$ points approximately toward the nearest point on the manifold, since $g(f(\tilde{\mathbf{x}}))$ estimates the center of mass of the clean points \mathbf{x} that could have given rise to $\tilde{\mathbf{x}}$. The autoencoder thus learns a vector field $g(f(\mathbf{x})) - \mathbf{x}$ indicated by the green arrows. This vector field estimates the score $\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$ up to a multiplicative factor that is the average root mean square reconstruction error.

Why Do We Use Auto-Encoders?

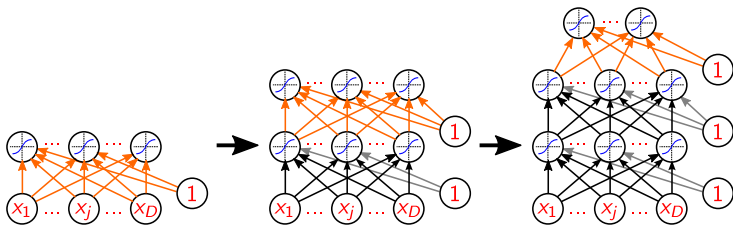
Historically, training deep neural networks was hard

One of the initial successful uses of auto-encoders was for **unsupervised pre-training** (Erhan et al., 2010).

Unsupervised Pre-Training

A greedy, layer-wise procedure:

- train one layer at a time, from first to last, with unsupervised criterion (e.g. an auto-encoder)
- fix the parameters of previous hidden layers
- previous layers viewed as feature extraction



Pre-training initializes the parameters in a region such that the near local optima overfit less the data.

Fine-Tuning

Once all layers are pre-trained:

- add output layer
- train the whole network using supervised learning

Supervised learning is performed as in a regular feed-forward network:

- forward propagation, backpropagation and update
- all parameters are “tuned” for the supervised task at hand
- representation is adjusted to be more discriminative

Other Applications of Auto-Encoders

- Dimensionality reduction
- Information retrieval and semantic hashing (via binarizing the codes)
- Conversion of discrete inputs to low-dimensional continuous space

Word Representations

We'll focus now on recent methods for learning representations of **words in natural language**

Also called **word embeddings**

This has been an extremely successful application of representation learning

It's still a very active area of research!

Distributional Similarity

Key idea: represent a word by means of its neighbors

- “*You shall know a word by the company it keeps*” (J. R. Firth, 1957)
- One of the most successful ideas of modern statistical NLP!

For example:

- Adjectives are normally surrounded by nouns
- Words like *book, newspaper, article*, are commonly surrounded by *reading, read, writes*, but not by *flying, eating, sleeping*

Word Embeddings

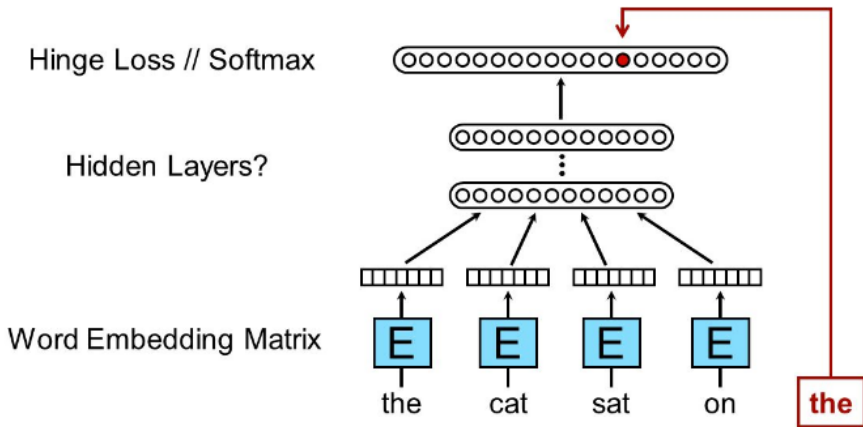
How do we obtain lower dimensional vector representations of words?

Two possible methods:

- Factorization of a co-occurrence word/context matrix (latent semantic analysis, etc.)
- Directly learn low-dimensional vectors by training a network to *predict* the context of a given word

We'll focus on the latter, incarnated in the **word2vec** toolkit (Mikolov et al., 2013), which follows previous ideas of Bengio et al. (2003) and Collobert et al. (2011).

Neural Language Model (Bengio et al., 2003)



(Image credits: Quoc Le)

Neural Language Model (Bengio et al., 2003)

- Each word is associated with a continuous vector (a **word embedding**)
- Given the **context** (previous K words), predict the next word
- This is done by concatenating the word embeddings in the context window, then propagating them through a feedforward neural network
- The output layer is a gigantic softmax that assigns a probability value to each word in the vocabulary

Variants of this model achieved better accuracy than smoothed K -th order Markov models

As a by-product: **word embeddings!**

The embedding matrix is a lookup table that assigns a continuous vector to every word in the vocabulary.

Neural Language Model

In this class, we are not concerned with language modeling (the actual task), but rather about the **quality of the embeddings** (the representations we learn for that task).

Some Insights

If we don't care about language modeling as a task:

- 1 We don't need to have a “left-to-right model” where we try to predict the next word given the context
- 2 We don't need to predict the probability of every word, we might just make sure that the true word is more likely than a random word

These insights underlie the word2vec model of Mikolov et al. (2013).

Word2Vec (Mikolov et al., 2013)

Considers a context window **around** each word in the sentence.

Word2vec comes with two variants:

- **Skip-gram**: predict surrounding context words in a window of length m of every word
- **Continuous bag-of-words (CBOW)**: predict the central word from the context

We'll focus on the skip-gram model (more widely used).

Skip-Gram

Goal: maximize the log probability of any context word given the current center word:

$$J(\Theta) = \frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log p_{\Theta}(x_{t+j} | x_t)$$

There are two sets of parameters $\Theta = (\mathbf{u}, \mathbf{v})$:

- Embeddings \mathbf{u}_o for each word o appearing as the center word
- Embeddings \mathbf{v}_c for each word c appearing in the context of another word

Define a **log-bilinear model**: $p_{\Theta}(x_{t+j} = c | x_t = o) \propto \exp(\mathbf{u}_o \cdot \mathbf{v}_c)$

Every word gets two vectors!

In the end, we use the \mathbf{u} vectors as the word embeddings and discard the \mathbf{v} vectors

The Large Vocabulary Problem

Recall that we have

$$p_{\Theta}(x_{t+j} = c \mid x_t = o) = \frac{\exp(\mathbf{u}_o \cdot \mathbf{v}_c)}{\sum_c' \exp(\mathbf{u}_o \cdot \mathbf{v}_c')}$$

This objective requires a softmax over the entire vocabulary

Unfortunately, with large vocabularies this leads to very slow training :(

Workarounds:

- Stochastic sampling
- Noise contrastive estimation
- Negative sampling

More details in these notes: <https://arxiv.org/pdf/1410.8251.pdf>

We'll focus on **negative sampling**.

Negative Sampling

Key idea:

- replace the gigantic softmax by **binary logistic regressions** for a true pair (center word and word in its context window) and a couple of random pairs (the center word with a random word):

$$J_t(\Theta) = \log \sigma(\mathbf{u}_o \cdot \mathbf{v}_c) + \sum_{i=1}^k \log \sigma(-\mathbf{u}_o \cdot \mathbf{v}_{j_i}), \quad j_i \sim P(x)$$

- Several strategies for the sampling distribution $P(x)$ (uniform, unigram frequency, etc.)

Negative sampling is a simple form of unsupervised pre-training.

Linear Relationships

- These representations are very good at encoding dimensions of similarity!
- **Word analogies** can be solved quite well just by doing vector subtraction in the embedding space
- Syntactically:

$$\mathbf{x}_{\text{apple}} - \mathbf{x}_{\text{apples}} \approx \mathbf{x}_{\text{car}} - \mathbf{x}_{\text{cars}} \approx \mathbf{x}_{\text{family}} - \mathbf{x}_{\text{families}}$$

- Semantically:

$$\begin{aligned} \mathbf{x}_{\text{shirt}} - \mathbf{x}_{\text{clothing}} &\approx \mathbf{x}_{\text{chair}} - \mathbf{x}_{\text{furniture}} \\ \mathbf{x}_{\text{king}} - \mathbf{x}_{\text{man}} &\approx \mathbf{x}_{\text{queen}} - \mathbf{x}_{\text{woman}} \end{aligned}$$

Visualization

Typical word embedding dimensions are on the hundreds (e.g. 300)

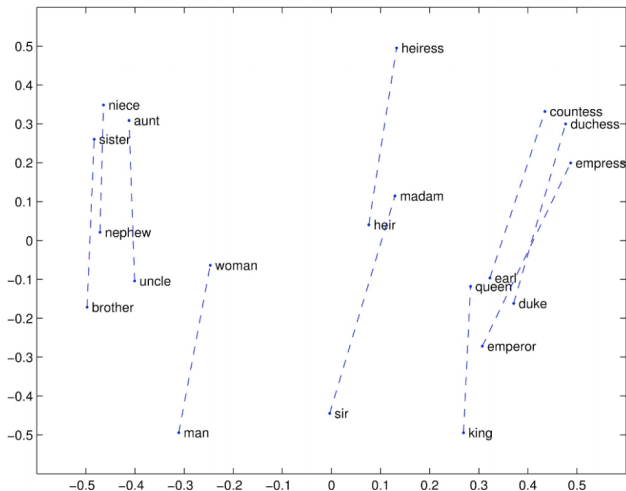
How can we visualize these embeddings?

Simple way: project them in 2D with something like PCA!

Most used: t -distributed stochastic neighbor embedding (t-SNE, Maaten and Hinton 2008)

<https://lvdmaaten.github.io/tsne>

Word Analogies (Mikolov et al., 2013)



(Slide credit to Richard Socher)

Other Methods for Obtaining Word Embeddings

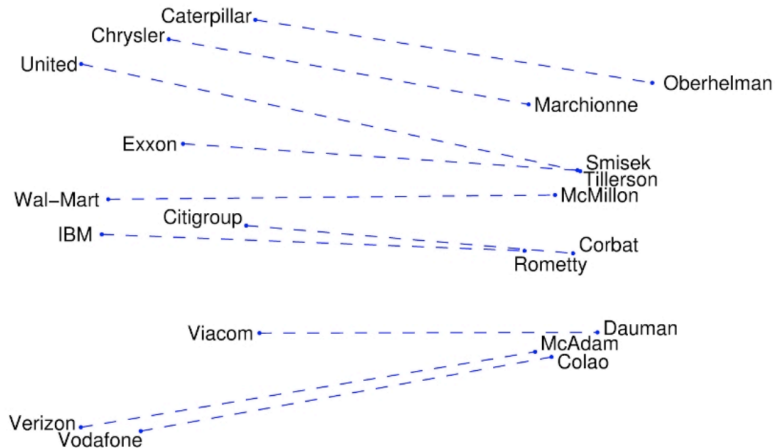
GloVe: Global Vectors for Word Representation (Pennington et al., 2014)

- <https://nlp.stanford.edu/projects/glove>
- Training is performed on aggregated global word-word co-occurrence statistics from a corpus

fastText (Bojanowski et al., 2016): embeds also character n -grams for generating embeddings for out-of-vocabulary words

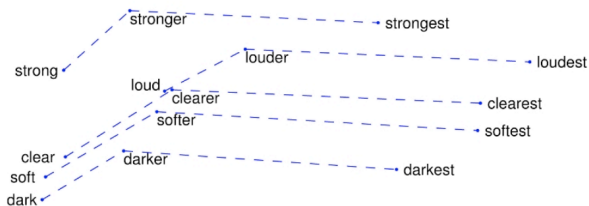
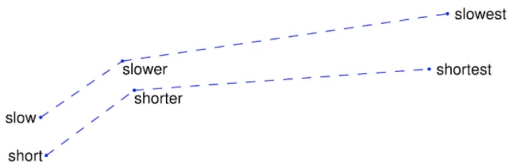
- <https://fasttext.cc> (from FAIR)
- open-source, free, lightweight library that allows users to learn text representations and text classifiers
- contains multi-lingual word vectors for 157 different languages

GloVe Visualizations: Company \rightarrow CEO



(Slide credit to Richard Socher)

GloVe Visualizations: Superlatives



(Slide credit to Richard Socher)

Word Embeddings: Some Open Problems

- Can we have word embeddings for multiple languages in the same space?
- How to capture polysemy?
- These word embeddings are static, can we compute embeddings on-the-fly depending on the context?

Cross-Lingual Word Embeddings

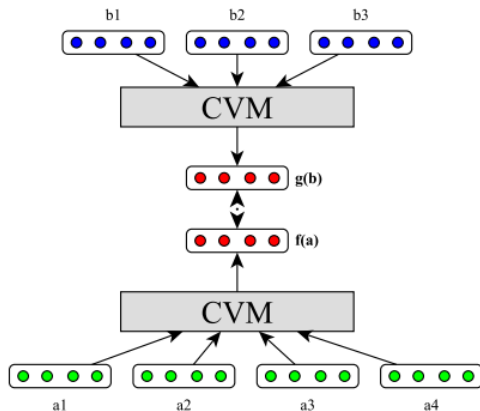


Figure 1: Model with parallel input sentences a and b . The model minimises the distance between the sentence level encoding of the bitext. Any composition functions (CVM) can be used to generate the compositional sentence level representations.

(From Hermann and Blunsom (2014).)

Cross-Lingual Word Embeddings

Key idea:

- use a corpus of parallel sentences in two languages
- define a composition function to obtain a sentence representation given word embeddings
- apply a loss function that encourages the sentence representations in the two languages to be similar
- negative sampling works here too: true pair vs fake pair.

Cross-Lingual Word Embeddings

Other approaches:

- Define a bilingual dictionary and apply canonical correlation analysis (Faruqui and Dyer, 2014)
- Task-specific embeddings with convex optimization (Ferreira et al., 2016)
- Learn the two embeddings separately, and then apply a linear transformation to put them in a shared space (Artetxe et al., 2017)
- Adversarial training (Lample et al., 2018)

This is a very active area of research!

Contextual Embeddings

Words can have different meanings, depending on which context they appear in.

In 2018, a model called ELMo learned **context-dependent embeddings** and achieved impressive results on 6 NLP downstream tasks (Peters et al., 2018)

Key idea:

- Pre-train a BiLSTM language model on a large dataset (we'll see in a later class what this is)
- Save **all** the encoder parameters at all layers, not only the embeddings
- Then, for your downstream task, tune a scalar parameter for each layer, and pass **the entire sentence** through this encoder.

BERT, GPT, etc.

Some time later, a Transformer-based model (BERT) achieved even better performance:

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova

Google AI Language

{jacobdevlin, mingweichang, kentonl, kristout}@google.com

Huge improvements in multiple NLP tasks!

(Trained on 64 TPU chips!!)

Other related models include GPT-2, GPT-3, etc.

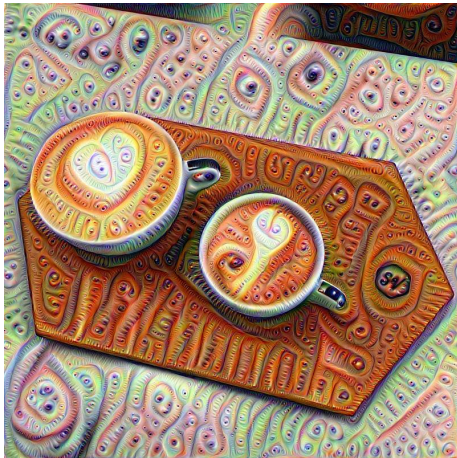
This will be covered in a later lecture!

Conclusions

- Neural nets learn internal representations that can be transferred across tasks
- Distributed representations are exponentially more compact and allow generalizing to unseen objects
- Deeper neural nets exhibit hierarchical compositionality: upper level layers learn more abstract/semantic representations than bottom level layers
- Auto-encoders are an effective means for learning representations
- Word embeddings are continuous representations of words that are extremely useful in NLP

Thank you!

Questions?



References I

- Artetxe, M., Labaka, G., and Agirre, E. (2017). Learning bilingual word embeddings with (almost) no bilingual data. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 451–462.
- Bengio, Y., Ducharme, R., Vincent, P., and Jauvin, C. (2003). A neural probabilistic language model. *Journal of Machine Learning Research*, 3(Feb):1137–1155.
- Bojanowski, P., Grave, E., Joulin, A., and Mikolov, T. (2016). Enriching word vectors with subword information. *arXiv preprint arXiv:1607.04606*.
- Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., and Kuksa, P. (2011). Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12(Aug):2493–2537.
- Erhan, D., Bengio, Y., Courville, A., Manzagol, P.-A., Vincent, P., and Bengio, S. (2010). Why does unsupervised pre-training help deep learning? *Journal of Machine Learning Research*, 11(Feb):625–660.
- Faruqui, M. and Dyer, C. (2014). Improving vector space word representations using multilingual correlation. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 462–471.
- Ferreira, D., Almeida, M. S. C., and Martins, A. F. T. (2016). Jointly Learning to Embed and Predict with Multiple Languages. In *Proc. of the Annual Meeting of the Association for Computational Linguistics*.
- Hermann, K. M. and Blunsom, P. (2014). Multilingual Models for Compositional Distributional Semantics. In *Proc. of the Annual Meeting of the Association for Computational Linguistics*.
- Hinton, G. E. (1984). Distributed representations.
- Lample, G., Ott, M., Conneau, A., Denoyer, L., and Ranzato, M. (2018). Phrase-based & neural unsupervised machine translation. *arXiv preprint arXiv:1804.07755*.
- Maaten, L. v. d. and Hinton, G. (2008). Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119.
- Pennington, J., Socher, R., and Manning, C. D. (2014). Glove: Global Vectors for Word Representation. *Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014)*, 12:1532–1543.
- Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., and Zettlemoyer, L. (2018). Deep contextualized word representations. *arXiv preprint arXiv:1802.05365*.
- Zeiler, M. D. and Fergus, R. (2013). Stochastic pooling for regularization of deep convolutional neural networks. *arXiv preprint arXiv:1301.3557*.

Lecture 8: Convolutional Neural Networks

André Martins, Francisco Melo, Mário Figueiredo



Deep Learning Course, Winter 2021-2022

Today's Roadmap

Today's lecture is about:

- Convolutional neural networks (CNN).
- Convolutional and max-pooling layers.

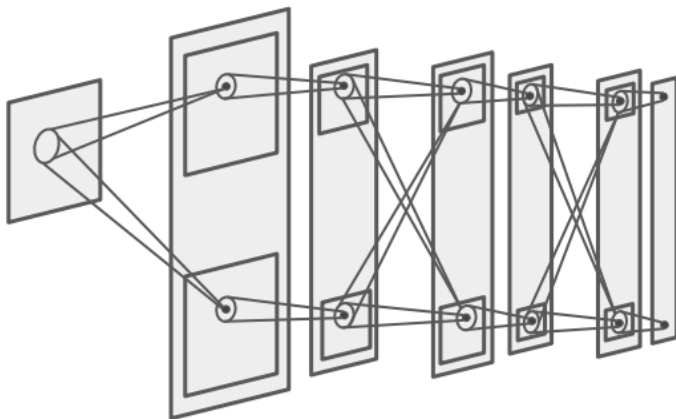
Convolutional Neural Networks

A **convolutional neural network** (CNN) is a NN with a **specialized connectivity structure**

Roadmap:

- Parameter tying/sharing
- 2D CNNs for object recognition
- Pooling
- Examples: ImageNet, AlexNet, GoogLeNet
- 1D CNNs in NLP

Neocognitron (Fukushima and Miyake, 1982)



(Credits: Fei-Fei Li, Johnson, Yeung)

- “Sandwich” architecture, alternating between simple cells with modifiable parameters and complex cells which perform **pooling**

Neocognitron (Fukushima and Miyake, 1982)

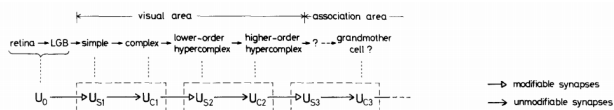


Fig. 1. Correspondence between the hierarchy model by Hubel and Wiesel, and the neural network of the neocognitron

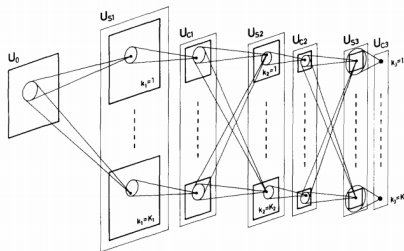
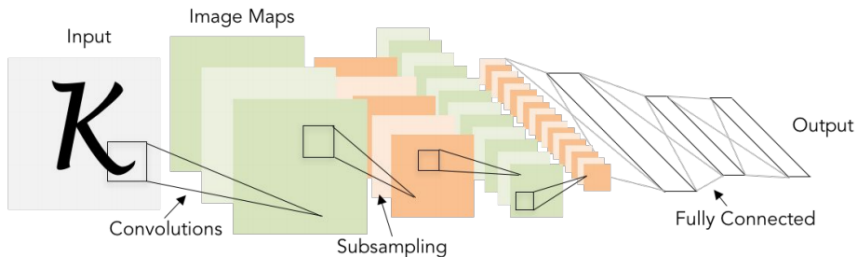


Fig. 2. Schematic diagram illustrating the interconnections between layers in the neocognitron

- Inspired by the multi-stage hierarchy model of the visual nervous system (Hubel and Wiesel, 1965)

ConvNet (LeNet-5) (LeCun et al., 1998)



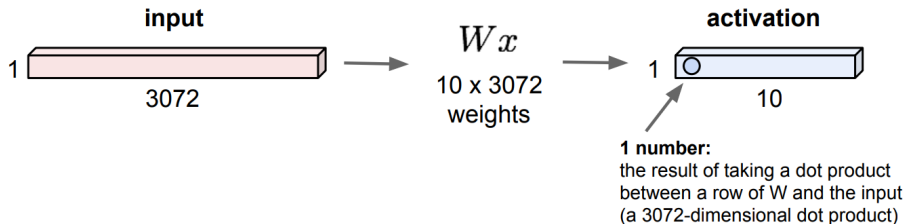
Convolutional Networks

... but what is a **convolutional layer**?

How is it different from a **fully connected** layer (as in a standard feedforward neural network).

Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



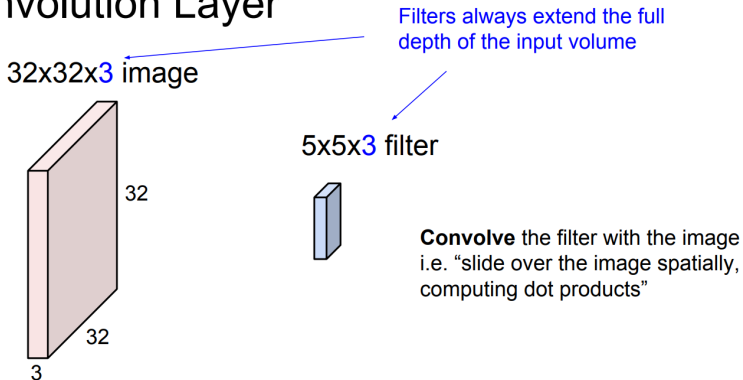
(Credits: Fei-Fei Li, Johnson, Yeung)

All activations depend on **all** inputs.

Convolutional Layer

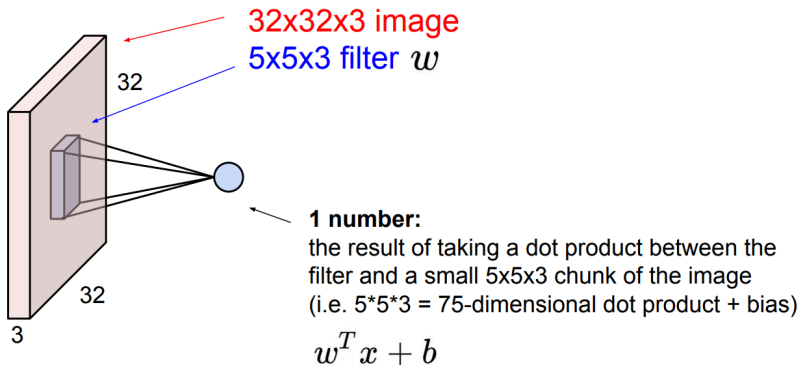
Don't stretch/reshape: preserve the spacial structure!

Convolution Layer



(Credits: Fei-Fei Li, Johnson, Yeung)

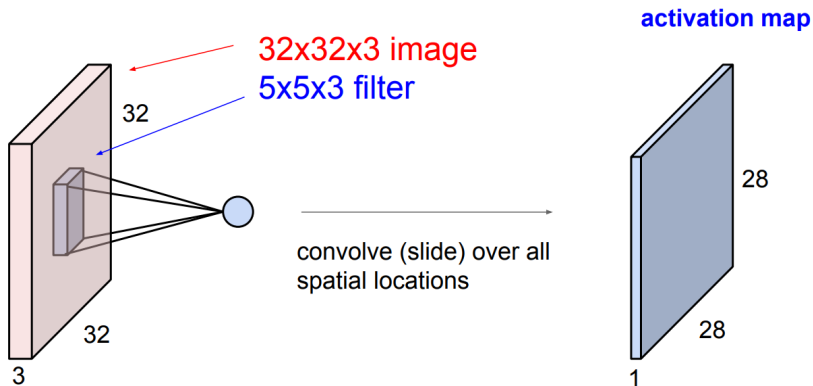
Convolutional Layer



(Credits: Fei-Fei Li, Johnson, Yeung)

Convolutional Layer

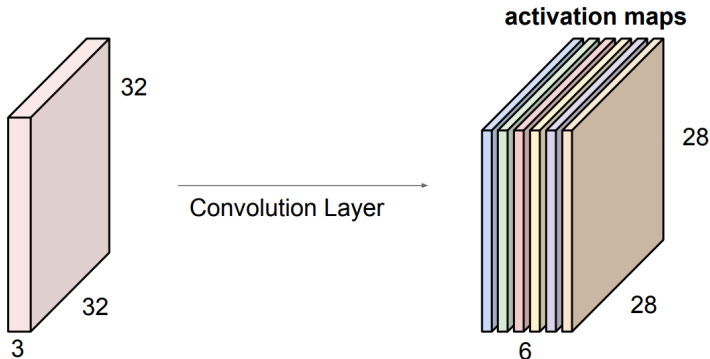
Apply the same filter to all spatial locations (28x28 times, [why?](#)):



(Credits: Fei-Fei Li, Johnson, Yeung)

Convolutional Layer

- For example, if we have **6** $5 \times 5 \times 3$ filters, we get **6** activation maps:



(Credits: Fei-Fei Li, Johnson, Yeung)

- We stack these up to get a “new image” of size $28 \times 28 \times 6$!

Image Size, Filter Size, Stride, Channels

Stride is the shift in pixels between two consecutive windows

So far we have considered a **stride** of 1

The number of **channels** is the number of filters we consider in each layer

Given an $N \times N \times D$ image, $F \times F \times D$ filters, K channels, and stride S , the resulting output will be of size $M \times M \times K$, where

$$M = (N - F)/S + 1$$

For example:

- $N = 32$, $D = 3$, $F = 5$, $K = 6$, $S = 1$ results in an $28 \times 28 \times 6$ output
- $N = 32$, $D = 3$, $F = 5$, $K = 6$, $S = 3$ results in an $10 \times 10 \times 6$ output

In practice: common to pad the border with zeros;

Common pad size is $(F - 1)/2$, which preserves size spatially: $M = N$.

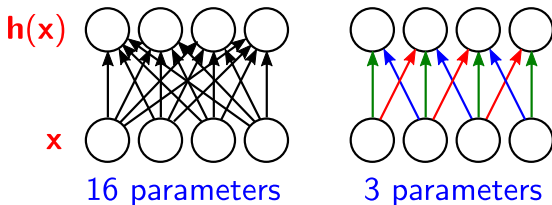
Convolutions and Parameter Tying

Why do we call this “convolutional”?

The **convolution** of a signal and a filter is:

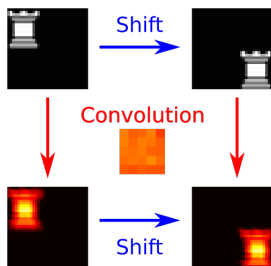
$$h[t] = (x * w)[t] = \sum_{a=-\infty}^{\infty} x[a]w[t - a].$$

The basic idea of a CNN is the combination of **sparse/local connectivity** and **parameter tying/sharing**.



Convolutions and Parameter Tying

Leads to **translation/shift equivariance**



Why do we want to tie (share) parameters?

- Reduce the number of parameters to be learned
- Deal with arbitrary long, variable-length, sequences: rather than shifting the filters, shift the input

Can also be done in 1D (e.g., text data, signals, ...)

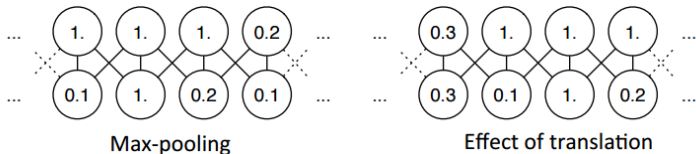
Convolutions and Pooling

The second component of CNNs is **pooling**

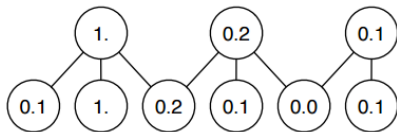
Common conv nets alternate convolutional layers and pooling layers.

Pooling Layers

- Aggregate to achieve local invariance:

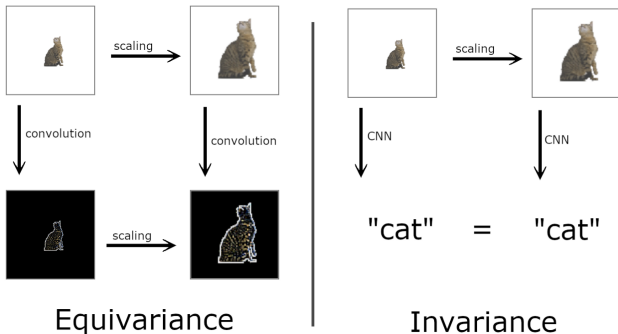


- Subsampling to reduce temporal/spacial scale and computation:



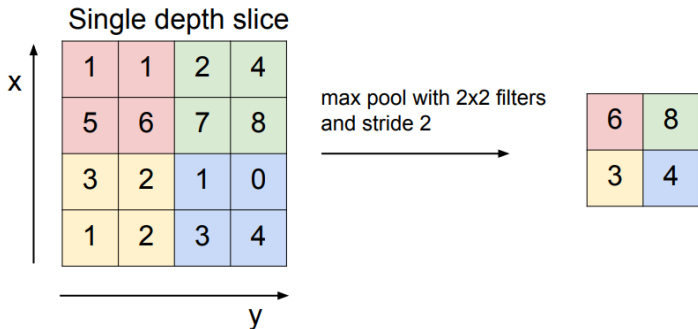
(Slide credit to Yoshua Bengio)

Equivariance vs Invariance



Pooling Layer

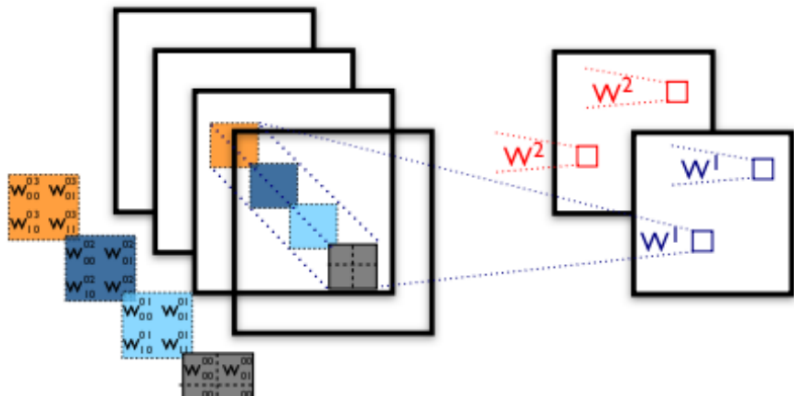
- Makes the representations smaller and more manageable
- Operates over each activation map (each channel) independently
- Max-pooling:



(Credits: Fei-Fei Li, Johnson, Yeung)

Multiple Convolutions: Feature Maps

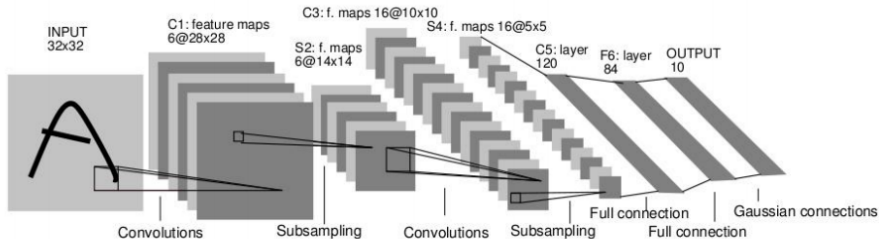
- Different filter weights for each channel, but keeping spatial invariance:



(Slide credit to Yoshua Bengio)

2D Convolutional Nets (LeCun et al., 1989)

- Inspired by “Neocognitron” (Fukushima, 1980)
- 2D Convolutions: the same filter (e.g. 3×3) is applied to each location of the image
- The filter weights are learned (as tied parameters)
- Multiple filters
- Alternates convolutional and pooling layers.



ConvNet Successes: MNIST



Handwritten text/digits:

- MNIST (0.35% error (Ciresan et al., 2011b))
- Arabic and Chinese (Ciresan et al., 2011a)

ConvNet Successes: CIFAR-10, Traffic Signs



Simpler recognition benchmarks:

- CIFAR-10 (9.3% error (Wan et al., 2013))
- Traffic signs: 0.56% error vs 1.16% for humans (Cireşan et al., 2011)

But less good at more complex datasets, e.g. Caltech-101/256 (few training examples).

ImageNet Dataset

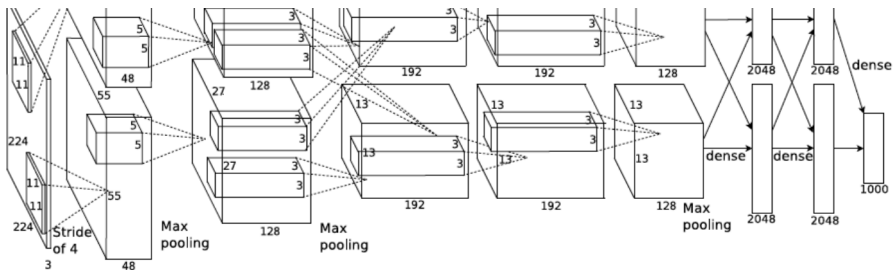
- 14 million labeled images, 20k classes
- Images gathered from Internet
- Human labels via Amazon Turk



(Slide credit to Rob Fergus)

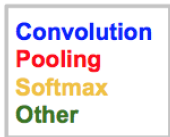
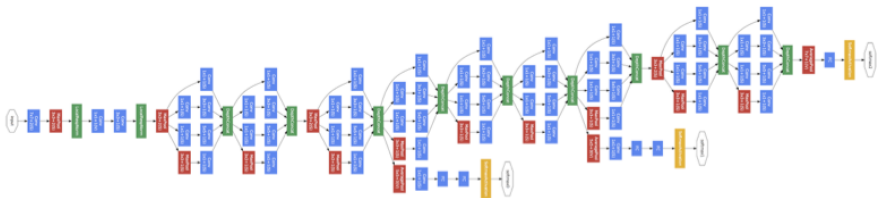
AlexNet (Krizhevsky et al., 2012)

- 54M parameters; 8 layers (5 conv, 3 fully-connected)
- Trained on 1.4M ImageNet images
- Trained on 2 GPUs for a week (50x speed-up over CPU)
- Dropout regularization
- Test error: 16.4% (second best team was 26.2%)



GoogLeNet (Szegedy et al., 2015)

- GoogLeNet inception module: very deep convolutional network, fewer (5M) parameters



Residual Networks (ResNets)

- Add skip-connections; tends to lead to more stable learning.

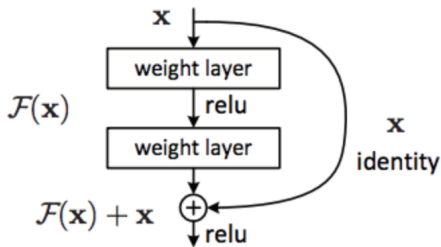


Figure 2. Residual learning: a building block.

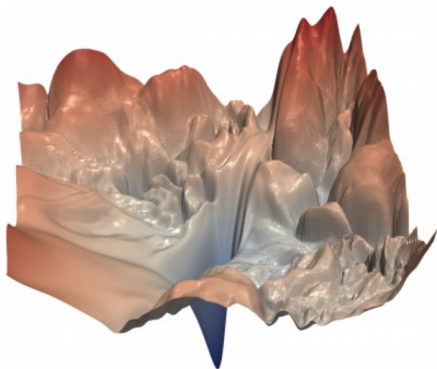
(He et al., 2016)

Residual Networks (ResNets)

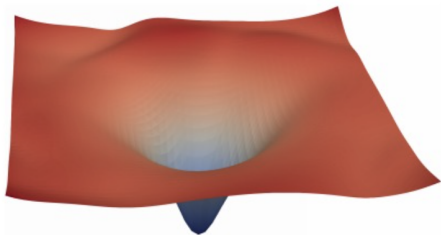
- Very deep network (34 layers here, but up to 152 layers!)
- VGG-19 (“Visual Geometry Group”) is Simonyan and Zisserman (2014) (19 layers, but more FLOPs)



Residual Networks (ResNets)



(a) without skip connections



(b) with skip connections

Figure 1: The loss surfaces of ResNet-56 with/without skip connections. The proposed filter normalization scheme is used to enable comparisons of sharpness/flatness between the two figures.

(Li et al., 2018)

Convolutional Nets in NLP

So far, we talked mostly about images.

Are conv nets also used in NLP? Not as much, but...

Quoting Yoav Goldberg in the Representation Learning Workshop in ACL 2018:

“NLP’s ImageNet moment has arrived.”

(Not referring to conv nets in particular, but to big neural architectures.)

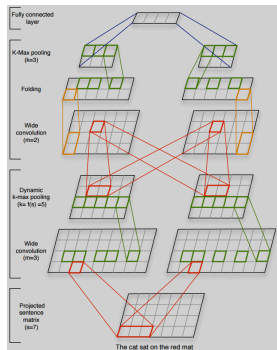
Convolutional Nets in NLP

- 1D convolutions
- Filters are applied to local windows around each word
- For word embeddings $\mathbf{x}_1, \dots, \mathbf{x}_L$, the filter response for word i is:

$$\mathbf{h}_i = \mathbf{g}(\mathbf{W}[\mathbf{x}_{i-h} \oplus \dots \oplus \mathbf{x}_i \oplus \dots \oplus \mathbf{x}_{i+h}] + \mathbf{b}),$$

where \oplus denotes vector concatenation and \mathbf{W} are shared parameters

- Can pad left and right with special symbols if necessary.



Kalchbrenner et al. (2014)

Variable Input Length

Most computation in conv nets can be done in parallel

GPUs can leverage this and achieve great speed-ups!

But unlike images which have fixed size, sentences have different lengths (number of words), which makes batching a bit trickier!

Mini-Batching, Padding, and Masking

Mini-batching is necessary to speed up training in GPUs

But how to cope with different input sizes (e.g. different sentence lengths)?

Solution: Minimize waste by sorting by sentence length before forming mini-batches, then **padding**:

Sentence 1	0's
Sentence 2	0's
Sentence 3	
Sentence 4	0's

Sentence 1	0's
Sentence 2	0's
Sentence 3	0's
Sentence 4	

(Image credit: Thang Luong, Kyunghyun Cho, Chris Manning)

Masking needs to be used to make sure the padded symbols are not affecting the results.

Beyond Convolutions

Other architectures have been proposed which offer alternatives to convolutions

For example: **transformer networks**, which stack multi-head attention layers

This is somewhat similar to “dynamic convolutions”

We'll cover this in another lecture.

What Representations Are We Learning?

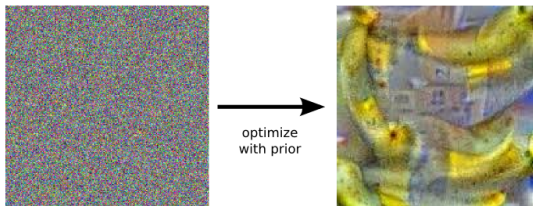
Which neurons fire for recognizing a particular object?

What parts of the network are activated?

To understand this, we need a way of **visualizing** what's happening inside the network.

Visualization

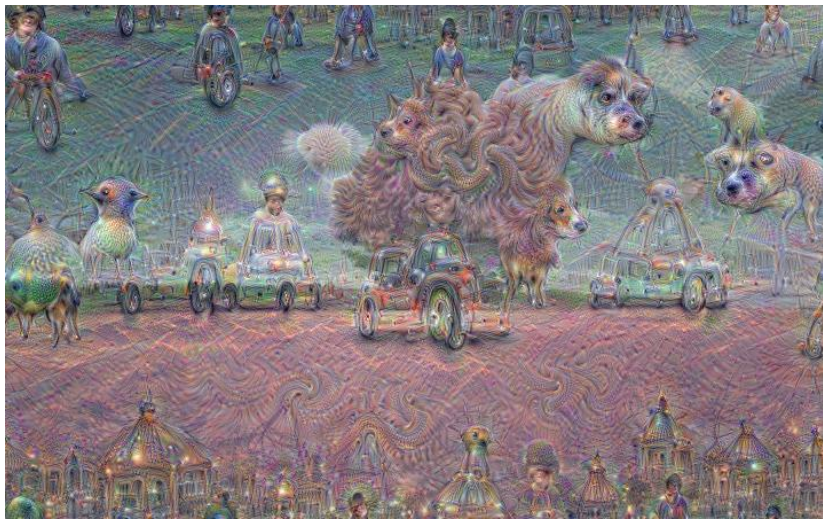
- **Idea:** Optimize input to maximize particular output
- Depends on the initialization
- **Google DeepDream**, maximizing “banana” output:



(from <https://research.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html>)

- Can also specify a particular layer and tune the input to maximize the layer’s activations—useful to see what kind of features each layer is representing
- Specifying a higher layer produces more complex representations...

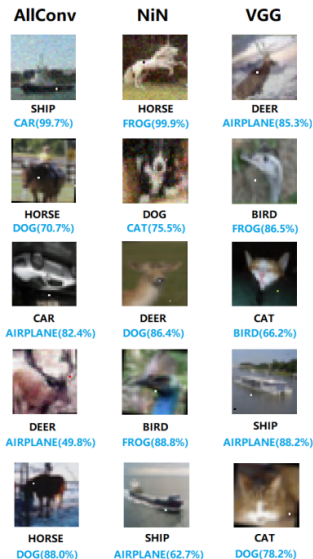
Google DeepDream



(from <https://research.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html>)

Adversarial Attacks

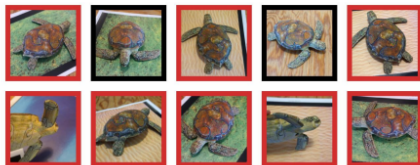
- How can we perturb an input slightly to fool a classifier?
- For example: **1-pixel attacks**
- **Glass-box model**: assumes access to the model
- Backpropagate to the inputs to find pixels which maximize the gradient
- There's also work for **black-box** adversarial attacks (don't have access to the model, but can query it).



(Credits: Su, Vargas, Sakurai (2018))

Even Worse: Perturb Object, Not Image

- Print the model of a turtle in a 3D printer.
- Perturbing the texture fools the model into thinking it's a rifle, regardless of the pose of the object!



■ classified as turtle ■ classified as rifle
■ classified as other

Figure 1. Randomly sampled poses of a 3D-printed turtle adversarially perturbed to classify as a rifle at every viewpoint². An unperturbed model is classified correctly as a turtle nearly 100% of the time.

(Credits: Athalye, Engstrom, Ilyas, Kwok (2018))

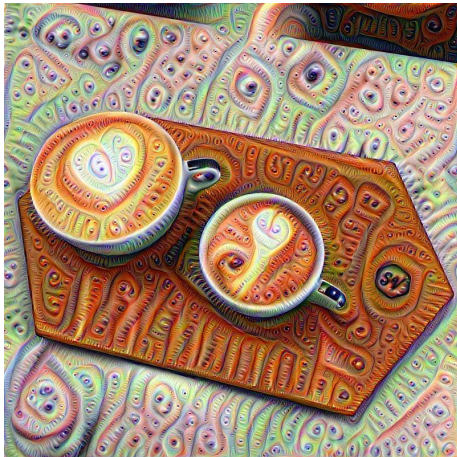
Neural networks are still very brittle!

Conclusions

- Convolutional neural networks (CNN) is a very powerful architecture for computer vision
- They take advantage of parameter sharing and sparse connectivity
- They are extremely useful to capture translational invariances in images
- Typically, convolution layers are alternated with max-pooling layers
- Lower layers capture more low-level representations (edges, corners)
- Higher layers have more “semantic” representations (objects, scenes)

Thank you!

Questions?



References I

- Ciresan, D., Meier, U., Masci, J., and Schmidhuber, J. (2011). A committee of neural networks for traffic sign classification. In *Neural Networks (IJCNN), The 2011 International Joint Conference on*, pages 1918–1921. IEEE.
- Ciresan, D. C., Meier, U., Gambardella, L. M., and Schmidhuber, J. (2011a). Convolutional neural network committees for handwritten character classification. In *2011 International Conference on Document Analysis and Recognition*, pages 1135–1139. IEEE.
- Ciresan, D. C., Meier, U., Masci, J., Maria Gambardella, L., and Schmidhuber, J. (2011b). Flexible, high performance convolutional neural networks for image classification. In *IJCAI Proceedings-International Joint Conference on Artificial Intelligence*, volume 22, page 1237.
- Fukushima, K. (1980). Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biological cybernetics*, 36(4):193–202.
- Fukushima, K. and Miyake, S. (1982). Neocognitron: A self-organizing neural network model for a mechanism of visual pattern recognition. In *Competition and cooperation in neural nets*, pages 267–285. Springer.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- Hubel, D. H. and Wiesel, T. N. (1965). Receptive fields and functional architecture in two nonstriate visual areas (18 and 19) of the cat. *Journal of neurophysiology*, 28(2):229–289.
- Kalchbrenner, N., Grefenstette, E., and Blunsom, P. (2014). A convolutional neural network for modelling sentences. *arXiv preprint arXiv:1404.2188*.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105.
- LeCun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., and Jackel, L. D. (1989). Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1(4):541–551.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324.
- Li, H., Xu, Z., Taylor, G., Studer, C., and Goldstein, T. (2018). Visualizing the loss landscape of neural nets. *Advances in Neural Information Processing Systems*, 31.

References II

- Simonyan, K. and Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., and Rabinovich, A. (2015). Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–9.
- Wan, L., Zeiler, M., Zhang, S., Cun, Y. L., and Fergus, R. (2013). Regularization of neural networks using dropconnect. In *Proc. of the International Conference on Machine Learning*, pages 1058–1066.

Lecture 9: Recurrent Neural Networks

André Martins, Francisco Melo, Mário Figueiredo



Deep Learning Course, Winter 2021-2022

Today's Roadmap

Today we'll cover **neural sequential models**:

- Recurrent neural networks.
- Backpropagation through time.
- Neural language models.
- The vanishing gradient problem.
- Gated units: LSTMs and GRUs.
- Bidirectional LSTMs.
- Example: ELMO representations.
- From sequences to trees: recursive neural networks.
- Other deep auto-regressive models: PixelRNNs.

Recurrent Neural Networks

Much interesting data is **sequential** in nature: words in sentences, DNA sequences, stock market returns, samples of sound signals, ...

How to deal with arbitrarily long sequences?

Feed-forward vs Recurrent Networks

- Feed-forward neural networks:

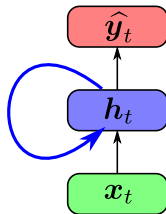
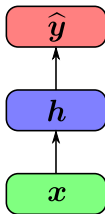
$$\mathbf{h} = \mathbf{g}(\mathbf{V}\mathbf{x} + \mathbf{c})$$

$$\hat{\mathbf{y}} = \mathbf{W}\mathbf{h} + \mathbf{b}$$

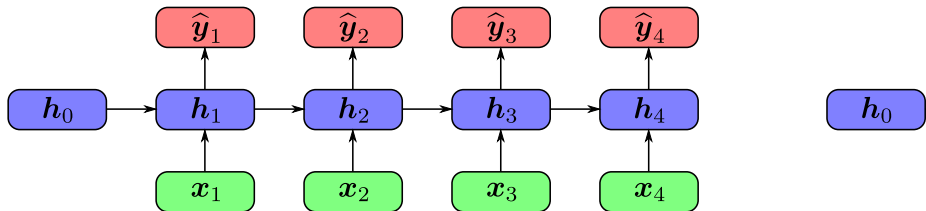
- Recurrent neural networks (Elman, 1990):

$$\mathbf{h}_t = \mathbf{g}(\mathbf{V}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{c})$$

$$\hat{\mathbf{y}}_t = \mathbf{W}\mathbf{h}_t + \mathbf{b}$$



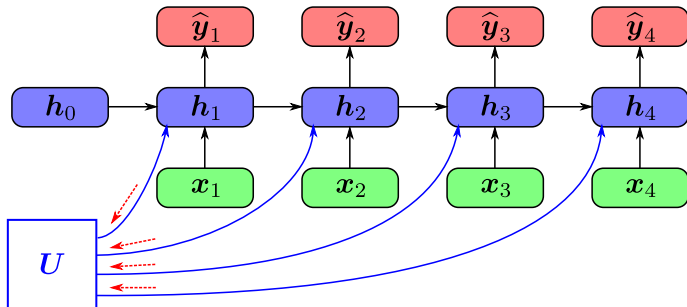
Unrolling the Graph



How do We Train the RNN Parameters?

- The unrolled graph is a well-formed (**directed and acyclic**) computation graph—we can use gradient backpropagation as usual
- Parameters are **tied/shared** across “time”
- Derivatives are aggregated across time steps
- This instantiation is called **backpropagation through time** (BPTT).

Parameter Tying



$$\frac{\partial \mathcal{L}}{\partial \mathbf{U}} = \sum_{t=1}^4 \frac{\partial \mathbf{h}_t}{\partial \mathbf{U}} \frac{\partial \mathcal{L}}{\partial \mathbf{h}_t}$$

- Same idea as when learning the filters in convolutional neural networks

What Can RNNs Be Used For?

We'll see three applications of RNNs:

- 1 **Sequence generation:** generates symbols sequentially with an **auto-regressive model** (e.g. language modeling)
- 2 **Sequence tagging:** takes a sequence as input, and returns a label for every element in the sequence (e.g., part of speech–POS–tagging)
- 3 **Pooled classification:** takes a sequence as input, and returns a single label by **pooling** the RNN states.

Recap: Full History Model

$$\mathbb{P}(\text{START}, y_1, y_2, \dots, y_L, \text{STOP}) = \prod_{t=1}^{L+1} \mathbb{P}(y_t | y_0, \dots, y_{t-1})$$

- The generation of each word depends on **all** the previous words
- Huge expressive power!
- But: **too many parameters** to estimate! (how many?)
- Cannot generalize well, specially for long sequences

Can We Have Unlimited Memory?

Markov models avoid the full history by considering a limited memory

Alternative: consider **all** the history, but compress it into a vector!

RNNs do this!

Auto-Regressive Models

Key ideas:

- Feed the previous output as input to the current step:

$$x_t = y_{t-1}$$

- Maintain a **state vector** \mathbf{h}_t which is a function of the previous state vector and the current input: this state will compress all the history!

$$\mathbf{h}_t = \mathbf{g}(\mathbf{V}x_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{c})$$

- Compute next output probability:

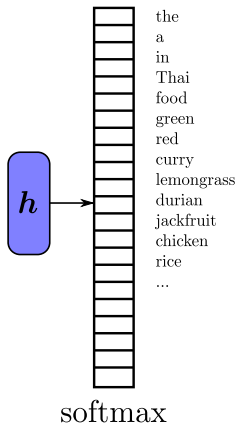
$$\mathbb{P}(y_t | y_0, \dots, y_{t-1}) = \mathbf{softmax}(\mathbf{W}\mathbf{h}_t + \mathbf{b})$$

Let's see each of these steps in detail

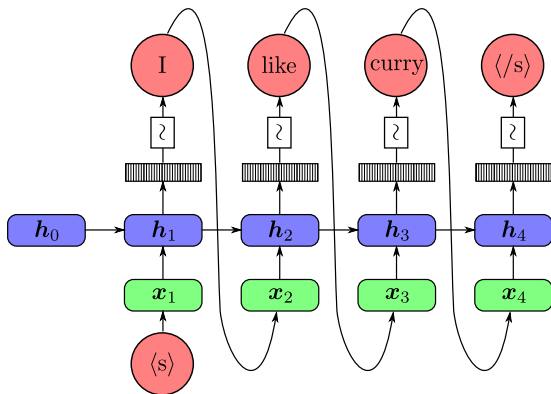
Language Modeling: Large Softmax

- To **generate text**, each y_t is a word in the vocabulary
- Typically, large vocabulary; e.g., $|V| = 100,000$

$$\begin{aligned} \mathbf{z}_t &= \mathbf{W}h_t + \mathbf{b} \\ p(y_t = i) &= \frac{\exp((\mathbf{z}_t)_i)}{\sum_j \exp((\mathbf{z}_t)_j)} \\ &= (\mathbf{softmax}(\mathbf{z}))_i \end{aligned}$$



Language Modeling: Auto-Regression



$$\begin{aligned}\mathbb{P}(y_1, \dots, y_L) &= \mathbb{P}(y_1) \times \mathbb{P}(y_2 | y_1) \times \dots \times \mathbb{P}(y_L | y_1, \dots, y_{L-1}) \\ &= \mathbf{softmax}(Wh_1 + \mathbf{b}) \times \mathbf{softmax}(Wh_2 + \mathbf{b}) \times \dots \\ &\quad \times \mathbf{softmax}(Wh_L + \mathbf{b})\end{aligned}$$

Three Problems for Sequence Generating RNNs

Algorithms:

- Sampling a sequence from the probability distribution defined by the RNN
- Obtaining the most probable sequence
- Training the RNN.

Sampling a Sequence

This is easy!

- Compute \mathbf{h}_1 from $\mathbf{x}_1 = \text{START}$
- Sample $\mathbf{y}_1 \sim \text{softmax}(\mathbf{W}\mathbf{h}_1 + \mathbf{b})$
- Compute \mathbf{h}_2 from \mathbf{h}_1 and $\mathbf{x}_2 = \mathbf{y}_1$
- Sample $\mathbf{y}_2 \sim \text{softmax}(\mathbf{W}\mathbf{h}_2 + \mathbf{b})$
- ...and so on

Obtaining the Most Probable Sequence

Unfortunately, this is hard!

- It would require obtaining the $\mathbf{y}_1, \mathbf{y}_2, \dots$ that jointly maximize the product $\mathbf{softmax}(\mathbf{W}\mathbf{h}_1 + \mathbf{b}) \times \mathbf{softmax}(\mathbf{W}\mathbf{h}_2 + \mathbf{b}) \times \dots$
- Note that picking the best \mathbf{y}_t greedily at each time step doesn't guarantee the best sequence
- This is rarely needed in language models. But it is important in **conditional** language modelling
- More later, when discussing **sequence-to-sequence models**

Training the RNN

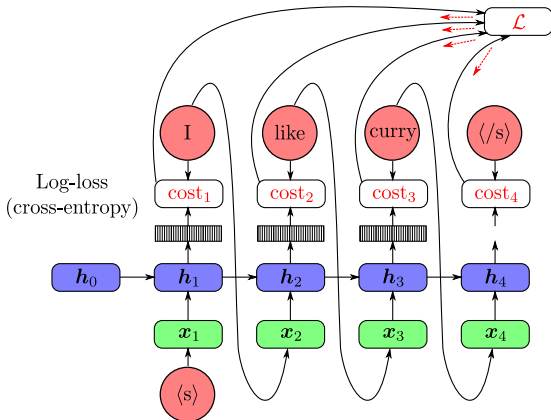
- Sequence-generating RNNs are typically trained with **maximum likelihood estimation**
- In other words, they are trained to minimize the **log-loss (cross-entropy)**:

$$\mathcal{L}(\Theta, y_{1:L}) = -\frac{1}{L+1} \sum_{t=1}^{L+1} \log \mathbb{P}_{\Theta}(y_t | y_0, \dots, y_{t-1})$$

- This is equivalent to minimizing **perplexity** $\exp(\mathcal{L}(\Theta, y_{1:L}))$
- Intuition: $-\log \mathbb{P}_{\Theta}(y_t | y_0, \dots, y_{t-1})$

measures how “perplexed” (or “surprised”) the model is when the t -th word is revealed

Training the RNN



- Unlike Markov (n -gram) models, **RNNs never forget!**
- However, we will see they might have trouble learning to use their memories (more soon...)

Teacher Forcing and Exposure Bias

Note that conditioning is on the **true history**, not on the model's predictions! This is known as **teacher forcing**.

Teacher forcing cause **exposure bias** at run time: the model will have trouble recovering from mistakes early on, since it generates histories that it has never observed before.

How to improve this is a current area of research!

Character-Level Language Models

We can also have an RNN over characters instead of words!

Advantage: can generate any combination of characters, not just words in a closed vocabulary.

Disadvantage: need to remember further away in history!

A Character-Level RNN Generating Fake Shakespeare

*PANDARUS: Alas, I think he shall be come approached and the day When little
srain would be attain'd into being never fed, And who is but a chain and subjects of
his death, I should not sleep.*

*Second Senator: They are away this miseries, produced upon my soul, Breaking
and strongly should be buried, when I perish The earth and thoughts of many states.*

DUKE VINCENTIO: Well, your wit is in the care of side and that.

*Second Lord: They would be ruled after this chamber, and my fair nues begun
out of the fact, to be conveyed, Whose noble souls I'll have the heart of the wars.*

Clown: Come, sir, I will make did behold your worship.

VIOLA: I'll drink it.

(Credits: Andrej Karpathy)

A Char-Level RNN Generating a Math Paper

Proof. Omitted. □

Lemma 0.1. *Let \mathcal{C} be a set of the construction.*

Let \mathcal{C} be a gerber covering. Let \mathcal{F} be a quasi-coherent sheaves of \mathcal{O} -modules. We have to show that

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

.

Proof. This is an algebraic space with the composition of sheaves \mathcal{F} on $X_{\acute{e}tale}$ we have

$$\mathcal{O}_X(\mathcal{F}) = \{\text{morph}_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where \mathcal{G} defines an isomorphism $\mathcal{F} \rightarrow \mathcal{F}$ of \mathcal{O} -modules. □

Lemma 0.2. *This is an integer Z is injective.*

Proof. See Spaces, Lemma ?? □

Lemma 0.3. *Let S be a scheme. Let X be a scheme and X is an affine open covering. Let $\mathcal{U} \subset X$ be a canonical and locally of finite type. Let X be a scheme. Let X be a scheme which is equal to the formal complex.*

The following to the construction of the lemma follows.

Let X be a scheme. Let X be a scheme covering. Let

$$b : X \rightarrow Y' \rightarrow Y \rightarrow Y' \times_X Y \rightarrow X.$$

be a morphism of algebraic spaces over S and Y .

Proof. Let X be a nonzero scheme of X . Let X be an algebraic space. Let \mathcal{F} be a quasi-coherent sheaf of \mathcal{O}_X -modules. The following are equivalent

- (1) \mathcal{F} is an algebraic space over S .
- (2) If X is an affine open covering.

Consider a common structure on X and X the functor $\mathcal{O}_X(U)$ which is locally of finite type. □

(Credits: Andrej Karpathy)

A Char-Level RNN Generating C++ Code

```
/*
 * Increment the size file of the new incorrect UI_FILTER group information
 * of the size generatively.
 */
static int indicate_policy(void)
{
    int error;
    if (fd == MARN_EPT) {
        /*
         * The kernel blank will coeld it to userspace.
         */
        if (ss->segment < mem_total)
            unblock_graph_and_set_blocked();
        else
            ret = 1;
        goto bail;
    }
    segaddr = in_SB(in.addr);
    selector = seg / 16;
    setup_works = true;
    for (i = 0; i < blocks; i++) {
        seq = buf[i++];
        bpf = bd->bd.next + i * search;
        if (fd) {
            current = blocked;
        }
    }
    rw->name = "Getjbbregs";
    bprm_self_clearl(&iv->version);
    regs->new = blocks[(BPF_STATS << info->historidac)] | PFMR_CLOBATHINC_SECONDS << 12;
    return segtable;
}
```

(Credits: Andrej Karpathy)

Note: these examples are from 5 years ago; we now have much more impressive language generators (e.g. GPT-3)

Instead of RNNs, the most recent language generators use **transformers**

We will cover **transformers** in a later class!

What Can RNNs Be Used For?

We'll see three applications of RNNs:

- 1 **Sequence generation:** generates symbols sequentially with an **auto-regressive model** (e.g. language modeling) ✓
- 2 **Sequence tagging:** takes a sequence as input, and returns a label for every element in the sequence (e.g., part of speech–POS–tagging)
- 3 **Pooled classification:** takes a sequence as input, and returns a single label by **pooling** the RNN states.

Sequence Tagging with RNNs

- In **sequence tagging**, we are given an input sequence x_1, \dots, x_L
- The goal is to assign a tag to each element of the sequence, yielding an output sequence y_1, \dots, y_L
- **Examples:** POS tagging, named entity recognition
- Differences with respect to sequence generation:
 - The input and output are distinct (no need for auto-regression)
 - The length of the output is known (same as that of the input)

Example: POS Tagging

- Map **sentences** to sequences of **part-of-speech tags**.

Time flies like an arrow .
noun verb prep det noun .

- Need to predict a morphological tag for each word of the sentence
- High correlation between adjacent words!
(Ratnaparkhi, 1999; Brants, 2000; Toutanova et al., 2003)

An RNN-Based POS Tagger

- The inputs $x_1, \dots, x_L \in \mathbb{R}^{E \times L}$ are word embeddings (found by looking up rows in an V -by- E embedding matrix, eventually pre-trained)
- As before, maintain a **state vector** h_t , function of h_{t-1} and the current x_t : this state compresses all the input history!

$$h_t = \mathbf{g}(\mathbf{V}x_t + \mathbf{U}h_{t-1} + \mathbf{c})$$

- A softmax output layer computes the probability of the current tag given the current and previous words:

$$\mathbb{P}(y_t | x_1, \dots, x_t) = \mathbf{softmax}(\mathbf{W}h_t + \mathbf{b})$$

An RNN-Based POS Tagger

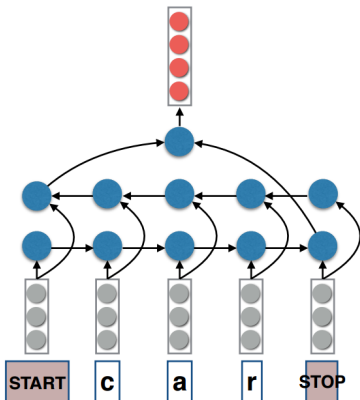
This model can be improved:

- Use a bidirectional RNN to condition also on the following words (combining a left-to-right and a right-to-left RNN)—more later!
- Use a nested character-level CNN or RNN to obtain embeddings for unseen words.

This model achieved SOTA performance on the Penn Treebank and several other benchmarks (Ling et al., 2015; Wang et al., 2015)!

Bidirectional RNNs

- We can read a sequence from left to right to obtain a representation
- Or we can read it from right to left
- Or we can read it from both and combine the representations
- More later...



(Slide credit: Chris Dyer)

Example: Named Entity Recognition

From **sentences** extract **named entities**.

Louis	Elsevier	was	born	in	Leuven	.
B-PER	I-PER	O	O	O	B-LOC	.

- Identify word segments that refer to entities (person, organization, location)
- Typically done with sequence models and B-I-O tagging

(Zhang and Johnson, 2003; Ratinov and Roth, 2009)

RNN-Based NER

- The model we described for POS tagging works just as well for NER
- However, NER has constraints about tag transitions: e.g., we cannot have I-PER after B-LOC
- The RNN tagger model we described exploits **input structure** (via the states encoded in the recurrent layer) but lacks **output structure**...

What Can RNNs Be Used For?

We'll see three applications of RNNs:

- 1 **Sequence generation:** generates symbols sequentially with an **auto-regressive model** (e.g. language modeling) ✓
- 2 **Sequence tagging:** takes a sequence as input, and returns a label for every element in the sequence (e.g., POS tagging) ✓
- 3 **Pooled classification:** takes a sequence as input, and returns a single label by **pooling** the RNN states.

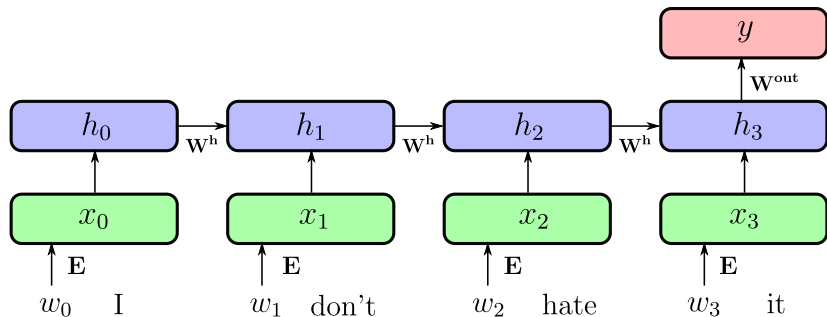
Pooled Classification

- What we have seen so far assumes we want to output a sequence of labels (either to generate or tag a full sequence).
- What about predicting a **single label** for the whole sequence?
- We can still use an RNN to capture the input sequential structure!
- Just **pool** the RNNs states, *i.e.*, map them to a single vector
- Use a single softmax to output the final label.

Pooling Strategies

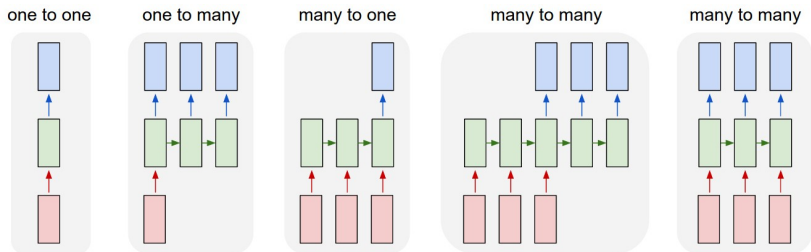
- The simplest strategy is just to **use the last RNN state**
- This state results from traversing the full sequence left-to-right, hence it has information about the full sequence!
- **Disadvantage:** for long sequences, the influence the earliest words may vanish
- **Other pooling strategies:**
 - use a bidirectional RNN and combine both last states of the left-to-right and right-to-left RNN
 - average pooling
 - ...

Example: Sentiment Analysis



(Slide credit: Ollion & Grisel)

Recurrent Neural Networks are Very Versatile



Check out Andrej Karpathy's blog post "The Unreasonable Effectiveness of Recurrent Neural Networks"

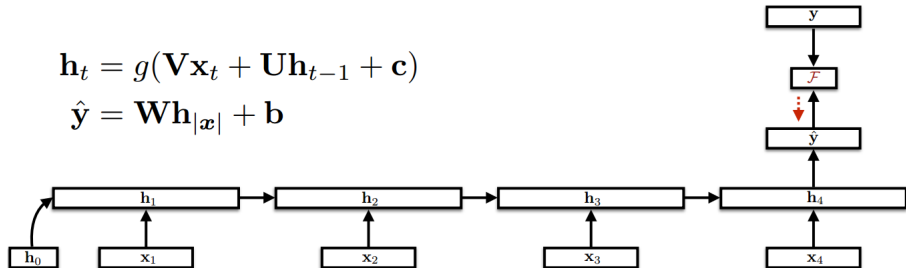
(<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>).

Training the RNN: Backpropagation Through Time

What happens to the gradients as we go back in time?

$$\mathbf{h}_t = g(\mathbf{V}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{c})$$

$$\hat{\mathbf{y}} = \mathbf{W}\mathbf{h}_{|x|} + \mathbf{b}$$



(Slide credit: Chris Dyer)

Backpropagation Through Time

What happens to the gradients as we go back in time?

$$\frac{\partial \mathcal{F}}{\partial \mathbf{h}_1} = \underbrace{\frac{\partial \mathbf{h}_2}{\partial \mathbf{h}_1} \frac{\partial \mathbf{h}_3}{\partial \mathbf{h}_2} \frac{\partial \mathbf{h}_4}{\partial \mathbf{h}_3}}_{\prod_{t=2}^4 \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{h}_4} \frac{\partial \mathcal{F}}{\partial \hat{\mathbf{y}}}$$

where

$$\prod_t \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} = \prod_t \frac{\partial \mathbf{h}_t}{\partial \mathbf{z}_t} \frac{\partial \mathbf{z}_t}{\partial \mathbf{h}_{t-1}} = \prod_t \text{Diag}(\mathbf{g}'(\mathbf{z}_t)) \mathbf{U}$$

Three cases:

- largest eigenvalue of \mathbf{U} exactly 1: gradient propagation is stable
- largest eigenvalue of $\mathbf{U} < 1$: **gradient vanishes** (exponential decay)
- largest eigenvalue of $\mathbf{U} > 1$: **gradient explodes** (exponential growth)

Vanishing and Exploding Gradients

- **Exploding gradients** can be dealt with by **gradient clipping** (truncating the gradient if it exceeds some magnitude)
- **Vanishing gradients** are more frequent and harder to deal with
 - In practice: long-range dependencies are difficult to learn
- **Solutions:**
 - Better optimizers (second order methods)
 - Normalization to keep the gradient norms stable across time
 - Clever initialization to start with good spectra (e.g., start with random orthonormal matrices)
 - **Alternative parameterizations: LSTMs and GRUs**

Gradient Clipping

- **Norm clipping:**

$$\tilde{\nabla} \leftarrow \begin{cases} \frac{c}{\|\nabla\|} \nabla & \text{if } \|\nabla\| \geq c \\ \nabla & \text{otherwise.} \end{cases}$$

- **Elementwise clipping:**

$$\tilde{\nabla}_i \leftarrow \min\{c, |\nabla_i|\} \times \text{sign}(\nabla_i), \quad \forall i$$

Alternative RNNs

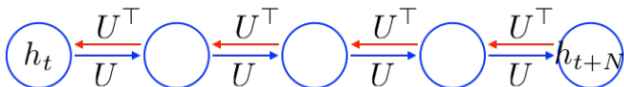
- **Gated recurrent unit** (GRU)
(Cho et al., 2014)
- **Long short-term memorie** (LSTM)
(Hochreiter and Schmidhuber, 1997)

Intuition: instead of multiplying across time (which leads to exponential growth), we want the error to be approximately constant

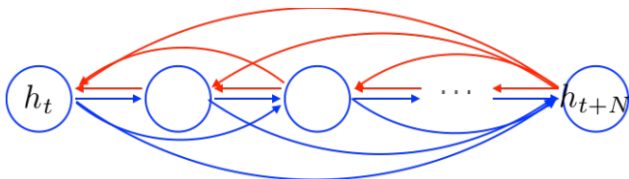
They solve the vanishing gradient problem, but still have exploding gradients (still need gradient clipping)

Gated Recurrent Units (Cho et al., 2014)

- Recall the problem: the error must backpropagate through all the intermediate nodes:



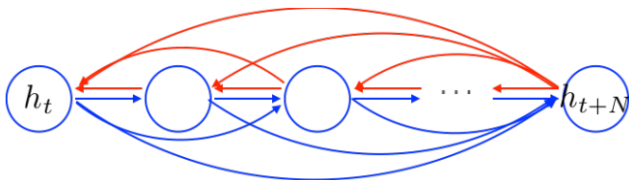
- Idea:** create some kind of shortcut connections:



(Image credit: Thang Luong, Kyunghyun Cho, Chris Manning)

- Create **adaptive** shortcuts controlled by special **gates**

Gated Recurrent Units (Cho et al., 2014)



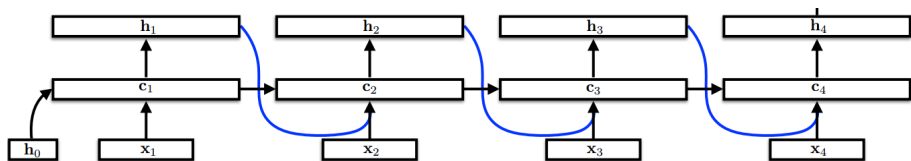
(Image credit: Thang Luong, Kyunghyun Cho, Chris Manning)

$$\mathbf{h}_t = \mathbf{u}_t \odot \tilde{\mathbf{h}}_t + (1 - \mathbf{u}_t) \odot \mathbf{h}_{t-1}$$

- **Candidate update:** $\tilde{\mathbf{h}}_t = \mathbf{g}(\mathbf{V}\mathbf{x}_t + \mathbf{U}(\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b})$
- **Reset gate:** $\mathbf{r}_t = \sigma(\mathbf{V}_r\mathbf{x}_t + \mathbf{U}_r\mathbf{h}_{t-1} + \mathbf{b}_r)$
- **Update gate:** $\mathbf{u}_t = \sigma(\mathbf{V}_u\mathbf{x}_t + \mathbf{U}_u\mathbf{h}_{t-1} + \mathbf{b}_u)$

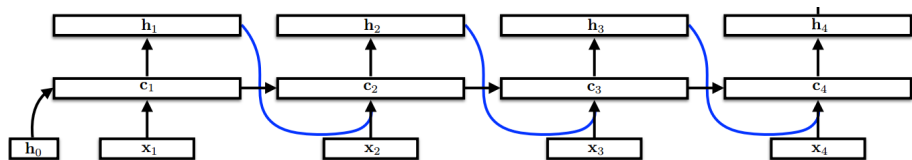
Long Short-Term Memories (Hochreiter and Schmidhuber, 1997)

- **Key idea:** use **memory cells** c_t
- To avoid the multiplicative effect, flow information *additively* through these cells
- Control the flow with special **input**, **forget**, and **output** gates



(Image credit: Chris Dyer)

Long Short-Term Memories

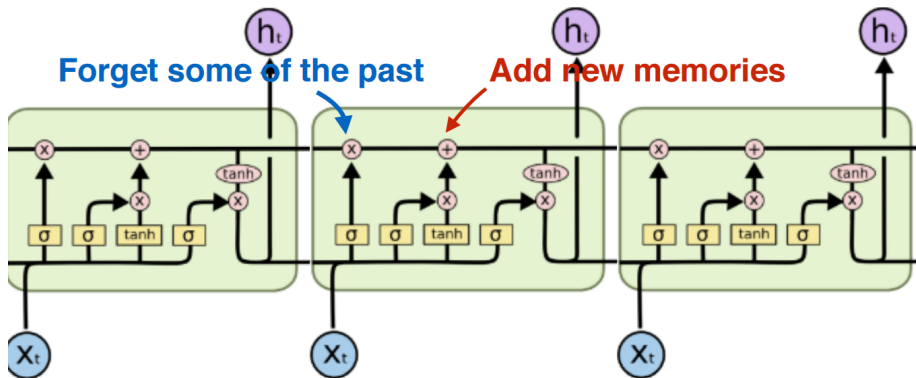


(Image credit: Chris Dyer)

$$c_t = f_t \odot c_{t-1} + i_t \odot g(Vx_t + Uh_{t-1} + b), \quad h_t = o_t \odot g(c_t)$$

- **Forget gate:** $f_t = \sigma(V_f x_t + U_f h_{t-1} + b_f)$
- **Input gate:** $i_t = \sigma(V_i x_t + U_i h_{t-1} + b_i)$
- **Output gate:** $o_t = \sigma(V_o x_t + U_o h_{t-1} + b_o)$

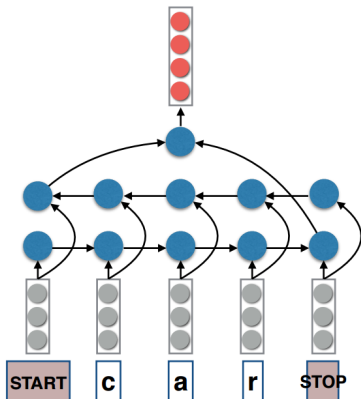
Long Short-Term Memories



(Slide credit: Christopher Olah)

Bidirectional LSTMs

- Same thing as a Bidirectional RNN, but using LSTM units instead of vanilla RNN units.



(Slide credit: Chris Dyer)

LSTMs and BILSTMs: Some Success Stories

- Time series prediction (Schmidhuber et al., 2005)
- Speech recognition (Graves et al., 2013)
- Named entity recognition (Lample et al., 2016)
- Machine translation (Sutskever et al., 2014)
- ELMo (deep contextual) word representations (Peters et al., 2018)
- ... and many others.

Summary

- Better gradient propagation is possible if we use **additive** rather than multiplicative/highly non-linear recurrent dynamics
- Recurrent architectures are an active area of research (but LSTMs are hard to beat)
- Other variants of LSTMs exist which tie/simplify some of the gates
- Extensions exist for *non-sequential* structured inputs/outputs (e.g. trees): **recursive neural networks** (Socher et al., 2011), **PixelRNN** (Oord et al., 2016)

From Sequences to Trees

- So far we've talked about **recurrent** neural networks, which are designed to capture sequential structure
- What about other kinds of structure? For example, trees?
- It is also possible to tackle these structures with recursive computation, via **recursive** neural networks.

Recursive Neural Networks

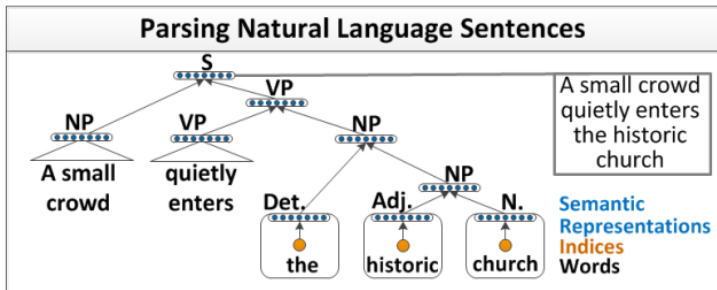
- Proposed by Socher et al. (2011) for parsing images and text
- Assume a binary tree (each node except the leaves has two children)
- Propagate states bottom-up in the tree, computing the parent state \mathbf{p} from the children states \mathbf{c}_1 and \mathbf{c}_2 :

$$\mathbf{p} = \tanh \left(\mathbf{W} \begin{bmatrix} \mathbf{c}_1 \\ \mathbf{c}_2 \end{bmatrix} + \mathbf{b} \right)$$

- Use the same parameters \mathbf{W} and \mathbf{b} at all nodes
- Can compute scores at the root or at each node by appending a softmax output layer at these nodes.

Compositionality in Text

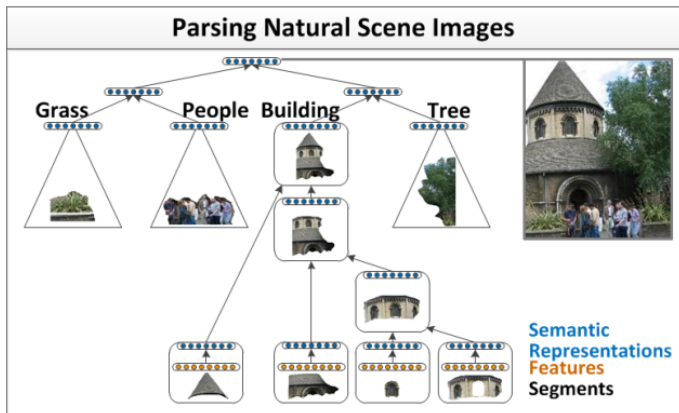
Uses a recurrent net to build a bottom-up parse tree for a sentence.



(Credits: Socher et al. (2011))

Compositionality in Images

Same idea for images.

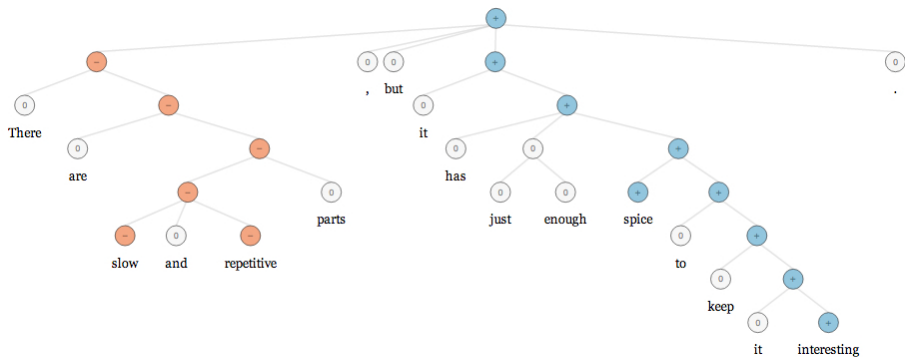


(Credits: Socher et al. (2011))

Tree-LSTMs

- Extend recursive neural networks the same way LSTMs extend RNNs, with a few more gates to account for the left and right child.
- Extensions exist for non-binary trees.

Fine-Grained Sentiment Analysis



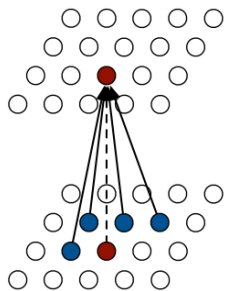
(Taken from Stanford Sentiment Treebank.)

What about Images?

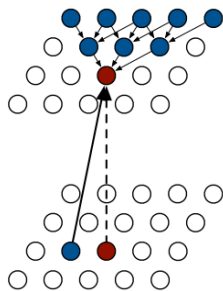
- While sequences are 1D, images are 2D.
- PixelRNNs are 2D extensions of RNNs.
- They can be used as auto-regressive models to **generate images**, by generating pixels in a particular order, conditioning on neighboring pixels.
- Several variants...

RNNs for Generating Images

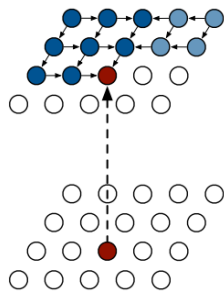
- Input-to-state and state-to-state mappings for PixelCNN and two PixelRNN models (Oord et al., 2016):



PixelCNN

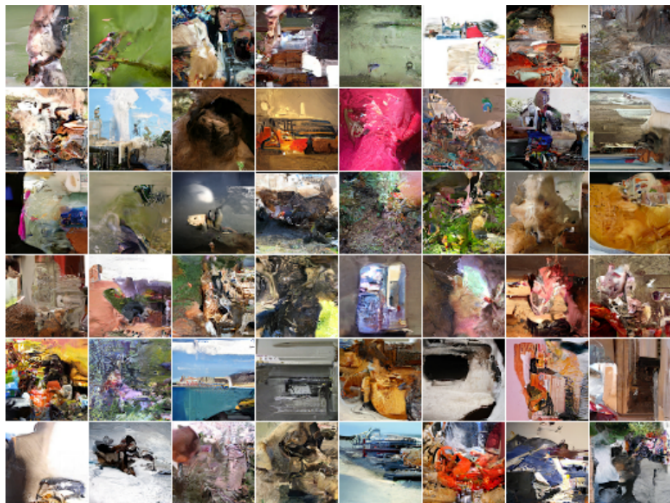


Row LSTM



Diagonal BiLSTM

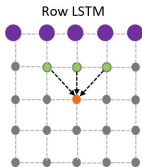
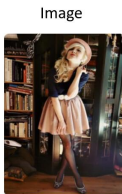
RNNs for Generating Images



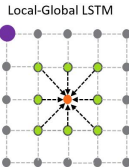
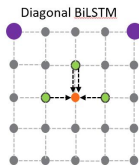
(Oord et al., 2016)

Even More General: Graph LSTMs

Traditional multi-dimensional LSTM:



Pixel-wise LSTM

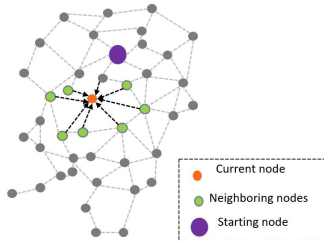
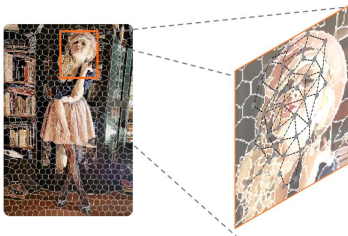


Supersixel map

Graph topology

Graph LSTM

Graph LSTM:



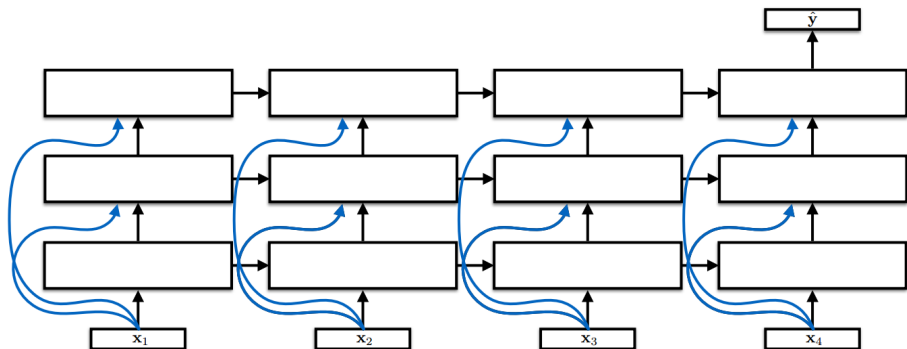
(Credits: Xiaodan Liang)

More Tricks of the Trade

- Depth
- Dropout
- Implementation Tricks
- Mini-batching

Deep RNNs/LSTMs/GRUs

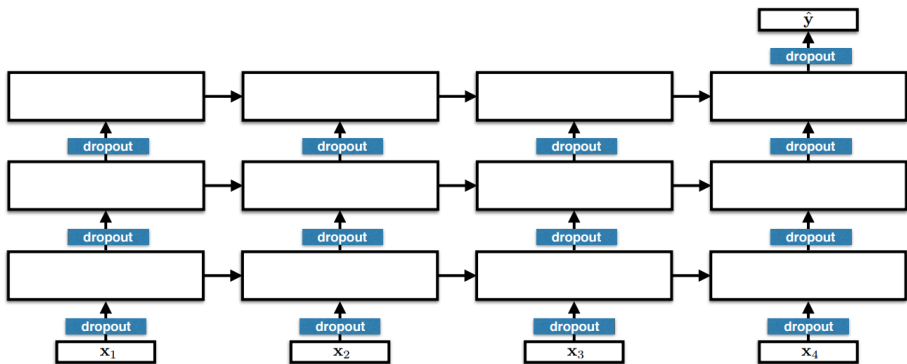
- Depth in recurrent layers helps in practice (2–8 layers seem to be standard)
- Input connections may or may not be used



(Slide credit: Chris Dyer)

Dropout in Deep RNNs/LSTMs/GRUs

- Apply dropout between layers, but not on the recurrent connections
- ... Or use the same mask for all recurrent connections (Gal and Ghahramani, 2015)



(Slide credit: Chris Dyer)

Implementation Tricks

For speed:

- Use diagonal matrices instead of full matrices (esp. for gates)
- Concatenate parameter matrices for all gates and do a single matrix-vector multiplication
- Use optimized implementations (from NVIDIA)
- Use GRUs or reduced-gate variant of LSTMs

For learning speed and performance:

- Initialize so that the bias on the forget gate is large (intuitively: at the beginning of training, the signal from the past is unreliable)
- Use random orthogonal matrices to initialize the square matrices

Mini-Batching

- RNNs, LSTMs, GRUs all consist of many element-wise operations (addition, multiplication, nonlinearities), and lots of matrix-vector products
- Mini-batching: convert many matrix-vector products into a single matrix-matrix multiplication
- Batch across instances, not across time
- The challenge with working with mini batches of sequences is... sequences are of different lengths (we've seen this when talking about convolutional nets)
- This usually means you bucket training instances based on similar lengths, and pad with zeros
- Be careful when padding not to back propagate a non-zero value!

Conclusions

Recurrent neural networks allow to take advantage of sequential input structure

They can be used to generate, tag, and classify sequences, and are trained with backpropagation through time

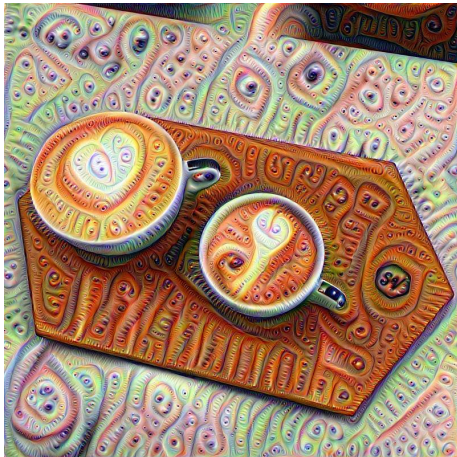
Vanilla RNNs suffer from vanishing and exploding gradients

LSTMs and other gated units are more complex variants of RNNs that avoid vanishing gradients

They can be extended to other structures like trees, images, and graphs.

Thank you!

Questions?



References I

- Brants, T. (2000). Tnt: a statistical part-of-speech tagger. In *Proceedings of the sixth conference on Applied natural language processing*, pages 224–231. Association for Computational Linguistics.
- Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. (2014). Learning Phrase Representations Using RNN Encoder-Decoder for Statistical Machine Translation. In *Proc. of Empirical Methods in Natural Language Processing*.
- Elman, J. L. (1990). Finding structure in time. *Cognitive science*, 14(2):179–211.
- Gal, Y. and Ghahramani, Z. (2015). Dropout as a bayesian approximation: Representing model uncertainty in deep learning. *arXiv preprint arXiv:1506.02142*.
- Graves, A., Mohamed, A.-r., and Hinton, G. (2013). Speech recognition with deep recurrent neural networks. In *ICASSP*, pages 6645–6649. IEEE.
- Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Lample, G., Ballesteros, M., Subramanian, S., Kawakami, K., and Dyer, C. (2016). Neural architectures for named entity recognition. In *Proc. of the Annual Meeting of the North-American Chapter of the Association for Computational Linguistics*.
- Ling, W., Luís, T., Marujo, L., Astudillo, R. F., Amir, S., Dyer, C., Black, A. W., and Trancoso, I. (2015). Finding function in form: Compositional character models for open vocabulary word representation. In *Proc. of Empirical Methods in Natural Language Processing*.
- Oord, A. v. d., Kalchbrenner, N., and Kavukcuoglu, K. (2016). Pixel Recurrent Neural Networks. In *Proc. of the International Conference on Machine Learning*.
- Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., and Zettlemoyer, L. (2018). Deep contextualized word representations. *arXiv preprint arXiv:1802.05365*.
- Ratinov, L. and Roth, D. (2009). Design challenges and misconceptions in named entity recognition. In *Proceedings of the Thirteenth Conference on Computational Natural Language Learning*, pages 147–155. Association for Computational Linguistics.
- Ratnaparkhi, A. (1999). Learning to parse natural language with maximum entropy models. *Machine Learning*, 34(1):151–175.
- Schmidhuber, J., Wierstra, D., and Gomez, F. J. (2005). Evolino: Hybrid neuroevolution/optimal linear search for sequence prediction. In *Proceedings of the 19th International Joint Conference on Artificial Intelligence (IJCAI)*.

References II

- Socher, R., Lin, C. C., Manning, C., and Ng, A. Y. (2011). Parsing natural scenes and natural language with recursive neural networks. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages 129–136.
- Sutskever, I., Vinyals, O., and Le, Q. V. (2014). Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, pages 3104–3112.
- Toutanova, K., Klein, D., Manning, C. D., and Singer, Y. (2003). Feature-rich part-of-speech tagging with a cyclic dependency network. In *Proc. of the North American Chapter of the Association for Computational Linguistics*, pages 173–180.
- Wang, P., Qian, Y., Soong, F. K., He, L., and Zhao, H. (2015). Part-of-speech tagging with bidirectional long short-term memory recurrent neural network. *arXiv preprint arXiv:1510.06168*.
- Zhang, T. and Johnson, D. (2003). A robust risk minimization based named entity recognition system. In *Proceedings of the seventh conference on Natural language learning at HLT-NAACL 2003-Volume 4*, pages 204–207. Association for Computational Linguistics.

Lecture 10: Sequence-to-Sequence Models

André Martins, Francisco Melo, Mário Figueiredo



Deep Learning Course, Winter 2021-2022

Today's Roadmap

Last lecture we talked about sequence tagging and sequence generation. Today we'll talk about **sequence-to-sequence models**.

- Machine translation
- Sequence vector representation
- Encoder-decoder architecture
- Sequence matrix representation
- Attention mechanism
- Encoder-decoder with attention

Sequence-to-Sequence

Sequence-to-sequence models map a source sequence (of arbitrary length) into a target sequence (also of arbitrary length)

This is **different** from **sequence tagging**, where the two sequences are of the same length

Example: Machine Translation

Goal: translate a **source sentence** x in one language into a **target sentence** y in another language.

Example (Portuguese to English):

x : *“A ilha de Utopia tem 200 milhas de diâmetro na parte central.”*



y : *“The island of Utopia is two hundred miles across in the middle part.”*

1950s: Early Machine Translation



(Source: <https://youtu.be/K-HfpsHPmvw>)

- MT research began in early 1950s
- Mostly Russian-English (motivated by the Cold War!)
- Systems were mostly rule-based, using a bilingual dictionary

Noisy Channel Model (Shannon and Weaver, 1949)



"When I look at an article in Russian, I say: 'This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode.' "





Raphael

... A ilha de Utopia tem
200 milhas de diâmetro
na parte central...



... the island of Utopia is
two hundred miles across
in the middle part...



Thomas

A very simple model: builds a **generative story** that works “backwards”
(flips source and target)

Yet, the dominant paradigm in MT for several decades (until 2014)

In 2014: **neural machine translation** (later)

1990s-2010s: Statistical Machine Translation

Goal: find the best English sentence y , given Russian sentence x

$$\hat{y} = \arg \max_y \mathbb{P}(y | x)$$

Key idea: use Bayes' rule to break this down into two components:

$$\hat{y} = \arg \max_y \mathbb{P}(x | y) \mathbb{P}(y)$$

- **Translation model:** models how words/phrases are translated (learnt from **parallel** data)
- **Language model:** models how to generate fluent English (learn from **monolingual** data)

How to Learn the Language Model?

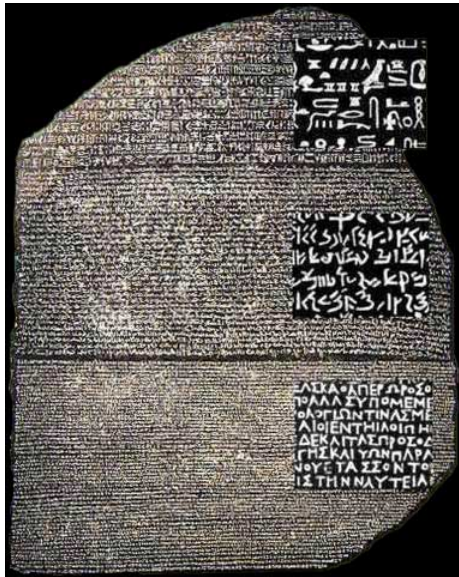
- Need large amounts of **monolingual** data (easy to get for most languages).
- How to learn a language model from these data?
- We covered language models in previous lectures:
 - Markov models (maybe with smoothing)
 - Neural language models
 - ...
- Pick your favorite!

How to Learn the Translation Model?

Need large amounts of **parallel** data!

(e.g., pairs of human translated Russian/English sentences.)

Rosetta Stone



- (Re-)discovered in 1799 near Alexandria
- Parallel corpora: ancient Egyptian, demotic Egyptian, ancient Greek

Europarl



- Proceedings from the European parliament sessions, translated into all EU official languages
- Around 1M parallel sentences for some language pairs
- Other corpora: Hansard, MultiUN, News Commentary, Wikipedia, OpenSubtitles, Paracrawl, ...

1990s: IBM Models for Statistical MT

- How to learn the translation model $\mathbb{P}(x | y)$?
- Assume we have enough parallel training data.
- Break it down further: consider instead

$$\mathbb{P}(x, \mathbf{a} | y),$$

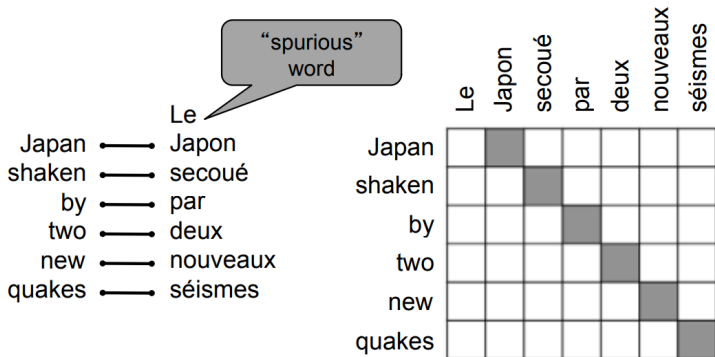
where \mathbf{a} are **word alignments**, i.e., word-level correspondences between Russian sentence x and English sentence y

- Word alignments are generally a latent/missing variable at training time, and need to be marginalized over at test time,

$$\mathbb{P}(x | y) = \sum_{\mathbf{a}} \mathbb{P}(x, \mathbf{a} | y),$$

Word Alignments

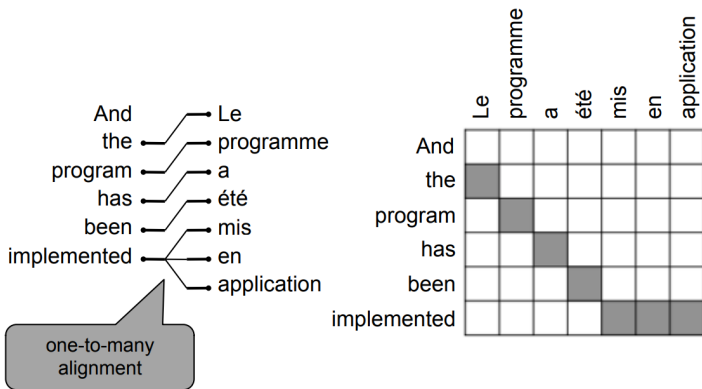
Example for English-French:



Some words may be unaligned (no counterpart in the other language)!

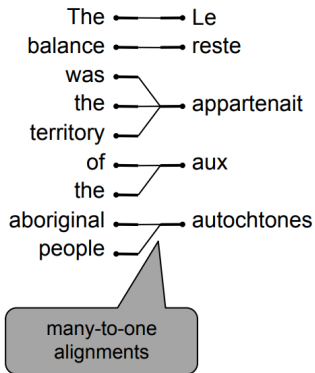
Word Alignments

Alignment can be one-to-many (**word fertility**):



Word Alignments

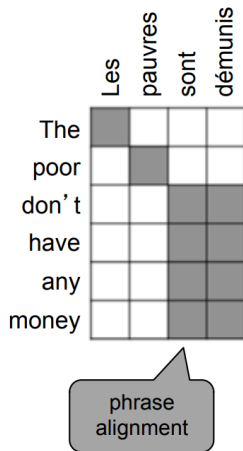
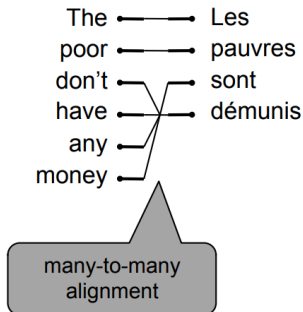
Alignment can be many-to-one:



	Le	reste	appartenait	aux	autochtones
The	■				
balance		■			
was			■		
the			■		
territory			■		
of				■	
the				■	
aboriginal					■
people					■

Word Alignments

Alignment can be many-to-many (phrase-level): **phrase-based MT**:



1990s: IBM Models for Statistical MT

- How to learn the translation model $\mathbb{P}(x | y)$?
- ...assuming we have enough parallel training data.

- Break it down further: consider instead

$$\mathbb{P}(x, \mathbf{a} | y).$$

- Learn $\mathbb{P}(x, \mathbf{a} | y)$ as a combination of several factors:
 - Probability of particular words aligning (co-occurrence, relative position, etc.)
 - Probability of words having a particular fertility
 - ...
- This leads to **IBM models** 1, 2, 3, 4, ...

1990s: IBM Models for Statistical MT

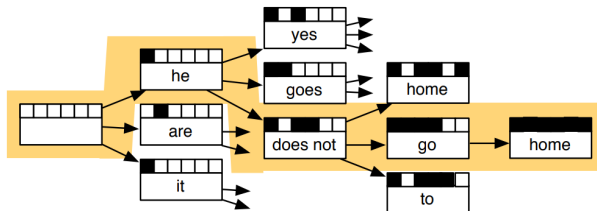
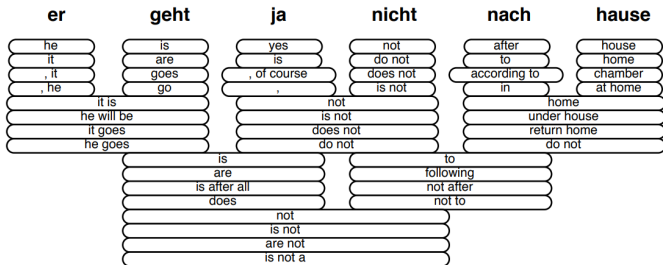
- To search the best translation, we need to solve

$$\hat{y} = \arg \max_y \sum_{\mathbf{a}} \mathbb{P}(\mathbf{x}, \mathbf{a} | y) \mathbb{P}(y),$$

combining the translation and language models.

- Enumerating all possible hypothesis and alignments is intractable.
- Typical approach: **heuristic search** to gradually build the translation, discarding hypotheses that are too low probability.

Searching for the Best Translation



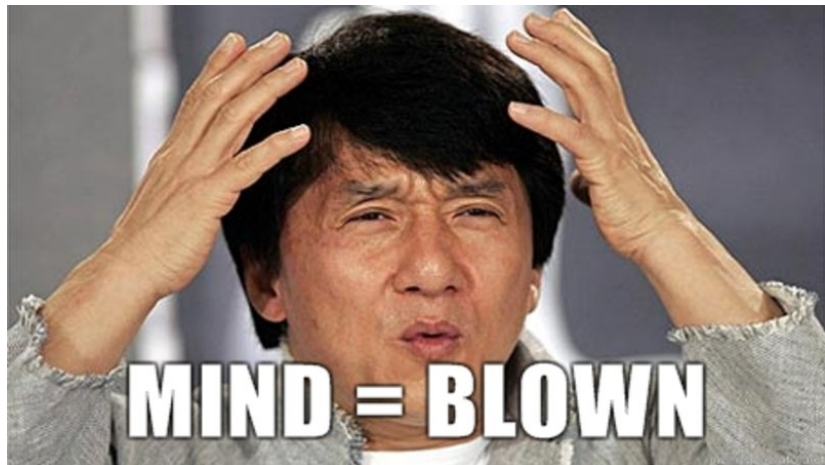
(Slide credit: <https://web.stanford.edu/class/cs224n/lectures/lecture10.pdf>)

Summarizing: Statistical Machine Translation

We only saw the tip of the iceberg: SMT is (was?) a huge research field.

- The best systems are extremely complex
- It's a big pipeline with many separately-designed subcomponents (translation and language model are only two examples)
- Lots of feature engineering
- System design is very language-dependent
- Requires compiling and maintaining resources (e.g., phrase tables)
- Models are disk/memory hungry
- Lots of human effort to maintain.

2014: Neural Machine Translation



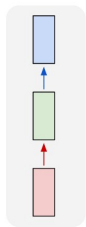
What is Neural Machine Translation (NMT)?

- NMT = MT with a **single neural network**
- End-to-end training with parallel data (no more complex pipelines!)
- The underlying architecture is an **encoder-decoder** (also called a **sequence-to-sequence model**)
- In fact, **NMT is also statistical**; however, historically, “statistical MT” refers to non-neural models, and NMT to NN-based models.

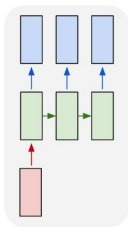
Recap: Recurrent Neural Networks

Lecture 9 covered RNNs and showed they can have several uses...

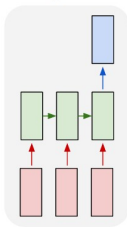
one to one



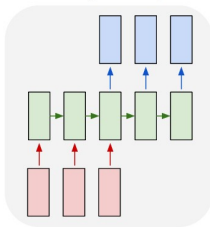
one to many



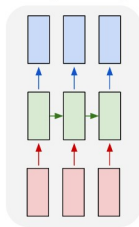
many to one



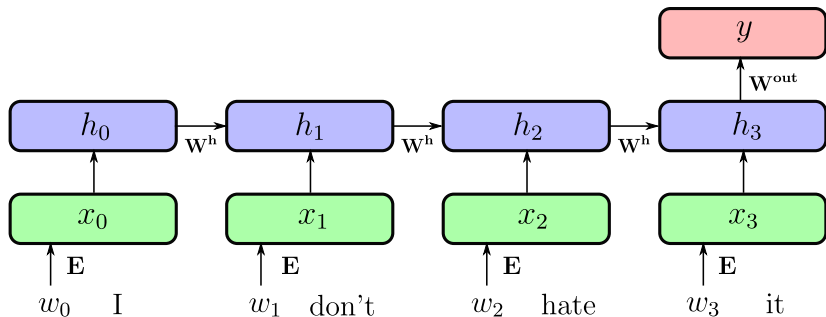
many to many



many to many

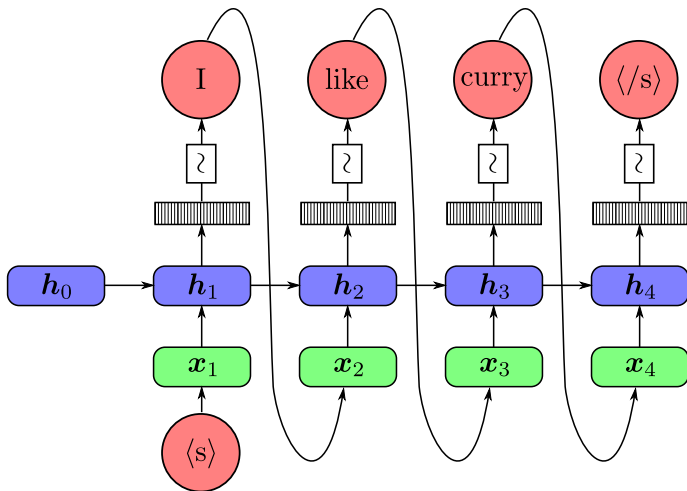


Recap: RNNs for Pooled Classification



(Slide credit: Ollion & Grisel)

Recap: Auto-Regressive RNNs for Sequence Generation

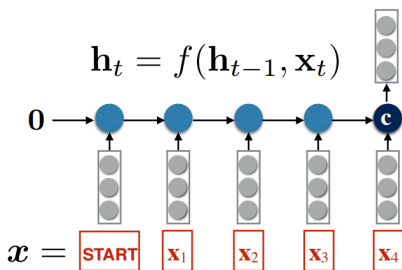


Sequence-to-Sequence Learning

(Cho et al., 2014; Sutskever et al., 2014)

- Can we put the two things together?
- Idea:
 - 1 **Encoder** RNN encodes source sentence, **generating a vector state**
 - 2 **Decoder** RNN generates target sentence **conditioned on vector state**.

Encode a Sequence as a Vector



(Slide credit: Chris Dyer)

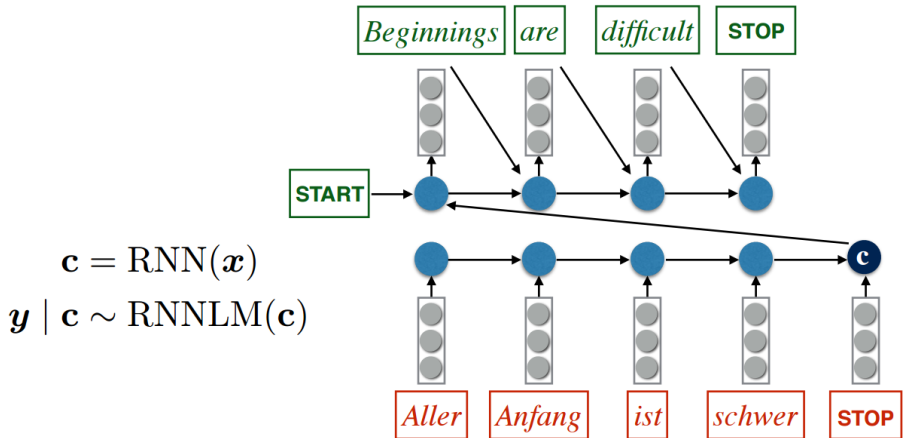
What is a vector representation of a sequence x ?

$$c = \text{RNN}(x)$$

What is the probability of a sequence $y \mid x$?

$$y \mid x \sim \text{RNNLM}(c)$$

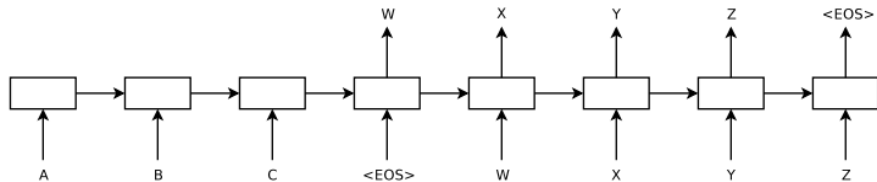
Encoder-Decoder Architecture



(Slide credit: Chris Dyer)

Encoder-Decoder Architecture

Another way of depicting it (from Sutskever et al. (2014)):



Some Problems

- If the source sentence is long, the encoder may forget the initial words and the translation will be degraded
 - Poor man's solution: reverse the source sentence.
- The decoder does greedy search—this leads to error propagation
 - Solution: beam search.

Beam Search

Ideally we want to find the target sentence y that maximizes

$$\hat{y} = \arg \max_y \mathbb{P}(y \mid x) = \arg \max_{y_1, \dots, y_L} \mathbb{P}(y \mid x) = \prod_{i=1}^L \mathbb{P}(y_i \mid y_{1:i-1}, x)$$

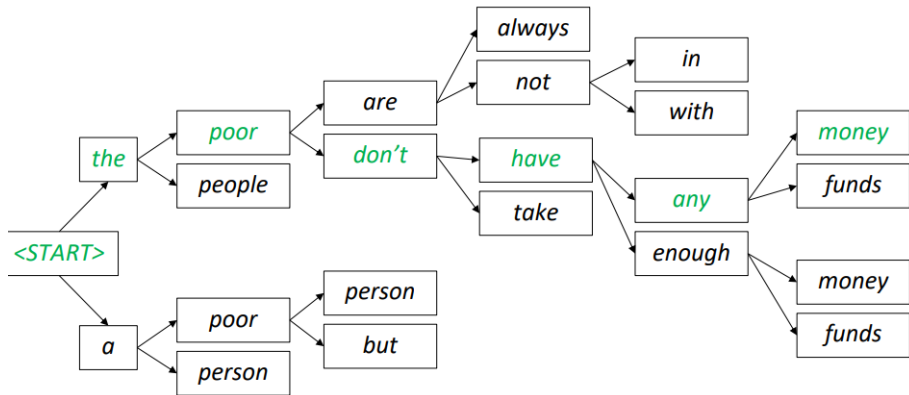
Enumerating all y is intractable!

Beam Search:

- approximate search strategy
- on each step of the decoder, keep track of the k most probable **partial** translations; discard the rest
- k is the **beam size**
- if $k = 1$, we recover greedy search.

Beam Search

Beam size = 2



(Source: <https://web.stanford.edu/class/cs224n/lectures/lecture10.pdf>)

Beam Search

- A little better than greedy search, but still greedy
- Runtime linear as a function of beam size k : trade-off speed/accuracy
- In practice: beam sizes ~ 4 – 12

Some Additional Tricks

From Sutskever et al. (2014):

Method	test BLEU score (ntst14)
Bahdanau et al. [2]	28.45
Baseline System [29]	33.30
Single forward LSTM, beam size 12	26.17
Single reversed LSTM, beam size 12	30.59
Ensemble of 5 reversed LSTMs, beam size 1	33.00
Ensemble of 2 reversed LSTMs, beam size 12	33.27
Ensemble of 5 reversed LSTMs, beam size 2	34.50
Ensemble of 5 reversed LSTMs, beam size 12	34.81

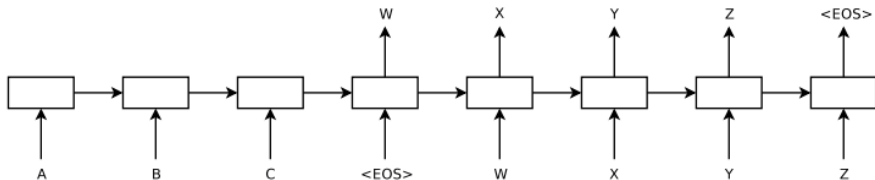
- Deep LSTMs
- Reversing the source sentence

At run time:

- Beam search
- Ensembling: combine N independently trained models and obtaining a “consensus” (always helps!)

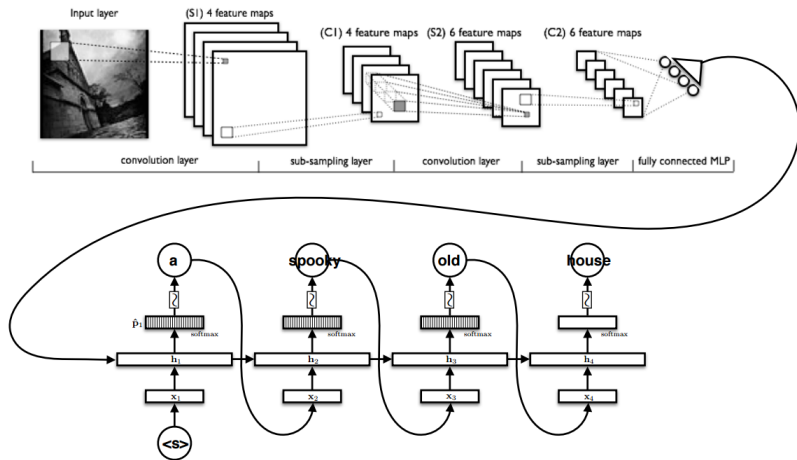
End-to-End Neural Machine Translation

- Previous statistical machine translation models were complicated pipelines (word alignments \rightarrow phrase table extraction \rightarrow language model \rightarrow decoding a phrase lattice)
- As an alternative, can do end-to-end NMT using a simple encoder-decoder
- Doesn't quite work yet, but gets close to top performance



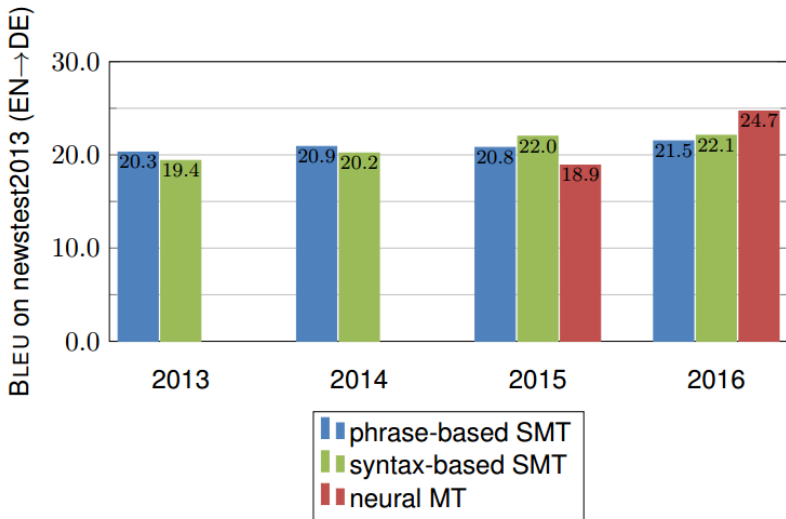
Caption Generation

Works for image inputs too!



(Slide credit: Chris Dyer)

Progress in Machine Translation



Slide credit: Rico Sennrich

NMT: A Success Story

- Neural MT went from a **fringe research activity** in 2014 to the **leading standard method** in 2016
 - 2014: First seq2seq paper published
 - 2016: Google Translate switches from SMT to NMT
- This is amazing!
- SMT systems, built by **hundreds** of engineers over **many years**, outperformed by NMT trained by a **handful** of engineers in a **few months**.

So Is Machine Translation Solved?

No. Many difficulties remain:

- Out-of-vocabulary words
- Domain mismatch between train and test data
- Low-resource language pairs
- Maintaining context over longer text (coming next!)

Limitations

A possible conceptual problem:

- Sentences have unbounded lengths
- Vectors have finite capacity

“You can't cram the meaning of a whole sentence into a single vector!” (Ray Mooney)



A possible practical problem:

- Distance between “translations” and their sources are distant—can LSTMs learn this?

Encode Sentences as Matrices, Not Vectors

Problem with the fixed-size vector model:

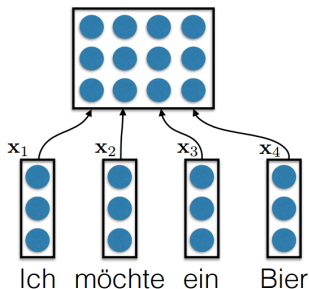
- Sentences are of different sizes but vectors are of the same size
- Bottleneck problem: a single vector needs to represent the full source sentence!

Solution: use **matrices** instead!

- Fixed number of rows, but number of columns depends on the number of words
- Then, before generating each word in the decoder, use an **attention mechanism** to condition on the relevant source words only

How to Encode a Sentence as a Matrix?

First shot: define the sentence words' vectors as the columns



(Image credit: Chris Dyer)

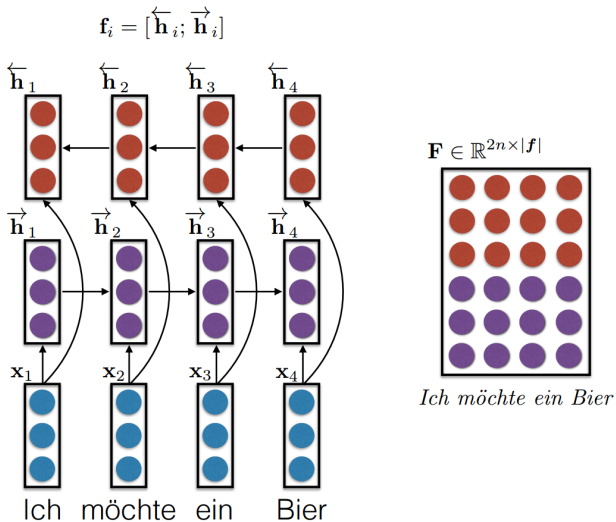
- Not very effective, since the word vectors carry no contextual information

How to Encode a Sentence as a Matrix?

Other strategies:

- Convolutional neural networks (Kalchbrenner et al., 2014): can capture context
- **Typical choice:** bidirectional LSTMs (Bahdanau et al., 2015)
- Later: Transformer networks (Vaswani et al., 2017).

Bidirectional LSTM Encoder



(Slide credit: Chris Dyer)

Generation from Matrices

- We now have a matrix F representing the input.
- How to generate from it?
- **Answer:** use **attention!** (Bahdanau et al., 2015)
- Attention is the neural counterpart of **word alignments**.

Generation from Matrices with Attention

- Generate the output sentence word by word using an RNN
- At each output position t , the RNN receives two inputs:
 - a fixed-size vector embedding of the previous output symbol y_{t-1}
 - a fixed-size vector encoding a “view” of the input matrix \mathbf{F} , via a weighted sum of its columns (i.e., words): $\mathbf{F}\mathbf{a}_t$
- The weighting of the input columns at each time-step (\mathbf{a}_t) is called the **attention** distribution.

Attention Mechanism (Bahdanau et al., 2015)

Let $\mathbf{s}_1, \mathbf{s}_2, \dots$ be the states produced by the decoder RNN

When predicting the t -th target word:

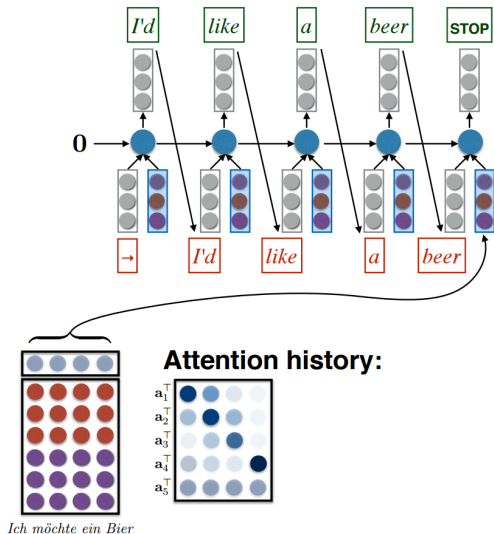
- 1 Compute “similarity” with each of the source words:

$$z_{t,i} = \mathbf{v} \cdot \mathbf{g}(\mathbf{W}\mathbf{h}_i + \mathbf{U}\mathbf{s}_{t-1} + \mathbf{b}), \quad \text{for } i = 1, \dots, L$$

where \mathbf{h}_i is the i th column of \mathbf{F} (representation of the i th source word), and \mathbf{v} , \mathbf{W} , \mathbf{U} , \mathbf{b} are parameters of the model

- 2 Form vector $\mathbf{z}_t = (z_{t,1}, \dots, z_{t,i}, \dots, z_{t,L})$ and compute attention $\mathbf{a}_t = \mathbf{softmax}(\mathbf{z}_t)$
- 3 Use attention to compute input conditioning state $\mathbf{c}_t = \mathbf{F}\mathbf{a}_t$
- 4 Update RNN state \mathbf{s}_t based on $\mathbf{s}_{t-1}, y_{t-1}, \mathbf{c}_t$
- 5 Predict $y_t \sim p(y_t | \mathbf{s}_t)$

Encoder-Decoder with Attention



(Slide credit: Chris Dyer)

Putting It All Together

obtain input matrix \mathbf{F} with a bidirectional LSTM

$t = 0$, $y_0 = \text{START}$ (the start symbol)

$\mathbf{s}_0 = \mathbf{w}$ (learned initial state)

repeat

$t = t + 1$

$\mathbf{z}_t = \mathbf{v} \cdot \mathbf{g}(\mathbf{W}\mathbf{F} + \mathbf{U}\mathbf{s}_{t-1} + \mathbf{b})$

compute attention $\mathbf{a}_t = \text{softmax}(\mathbf{z}_t)$

compute input conditioning state $\mathbf{c}_t = \mathbf{F}\mathbf{a}_t$

$\mathbf{s}_t = \text{RNN}(\mathbf{s}_{t-1}, [\mathbf{E}(y_{t-1}), \mathbf{c}_t])$

$y_t | y_{<t}, \mathbf{x} \sim \text{softmax}(\mathbf{P}\mathbf{s}_t + \mathbf{b})$

until $y_t \neq \text{STOP}$

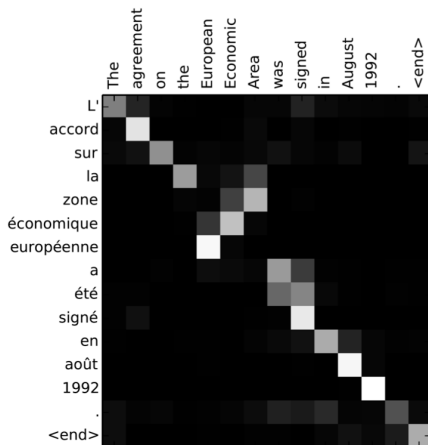
Attention Mechanisms

- Attention is closely related to “pooling” operations in convnets (and other architectures)
- Attention in MT plays a similar role as alignment, but leads to “soft” alignment instead of “hard” alignment
- Bahdanau et al. (2015)’s model has no bias in favor of diagonals, short jumps, fertility, etc.
- Some recent work adds some “structural” biases (Luong et al., 2015; Cohn et al., 2016)
- Other works constrains the amount of attention each word can receive (based on its **fertility**): Malaviya et al. (2018).

Attention is Great!

- Attention significantly improves NMT performance!
- It's very useful to allow decoder to focus on certain parts of the source
- Attention solves the bottleneck problem (by allowing the decoder to look directly at source)
- Attention helps with vanishing gradient problem (provides shortcut to faraway states)
- Attention provides some interpretability (we can see what the decoder was focusing on)
- This is good because we never explicitly trained a word aligner; the network learns it by itself!

Attention Map



Dzmitry Bahdanau, KyungHuyn Cho, and Yoshua Bengio. **Neural Machine Translation by Jointly Learning to Translate and Align**. ICLR'15.

Example: Machine Translation

Some positive examples where NMT has impressive performance:

Source	When asked about this, an official of the American administration replied: "The United States is not conducting electronic surveillance aimed at offices of the World Bank and IMF in Washington."	
PBMT	Interrogé à ce sujet, un responsable de l'administration américaine a répondu : "Les Etats-Unis n'est pas effectuer une surveillance électronique destiné aux bureaux de la Banque mondiale et du FMI à Washington".	3.0
GNMT	Interrogé à ce sujet, un fonctionnaire de l'administration américaine a répondu: "Les États-Unis n'effectuent pas de surveillance électronique à l'intention des bureaux de la Banque mondiale et du FMI à Washington".	6.0
Human	Interrogé sur le sujet, un responsable de l'administration américaine a répondu: "les Etats-Unis ne mènent pas de surveillance électronique visant les sièges de la Banque mondiale et du FMI à Washington".	6.0
Source	Martin told CNN that he asked Daley whether his then-boss knew about the potential shuffle.	
PBMT	Martin a déclaré à CNN qu'il a demandé Daley si son patron de l'époque connaissaient le potentiel remaniement ministériel.	2.0
GNMT	Martin a dit à CNN qu'il avait demandé à Daley si son patron d'alors était au courant du remaniement potentiel.	6.0
Human	Martin a dit sur CNN qu'il avait demandé à Daley si son patron d'alors était au courant du remaniement éventuel.	5.0

(From Wu et al. (2016))

Example: Machine Translation

... But also some negative examples:

- Dropping source words (lack of attention)
- Repeated source words (too much attention)

Source: 1922 in Wien geboren, studierte Mang während und nach dem Zweiten Weltkrieg Architektur an der Technischen Hochschule in Wien bei Friedrich Lehmann.

Human: Born in Vienna in 1922, Meng studied architecture at the Technical University in Vienna under Friedrich Lehmann *during and after the second World War*.

NMT: *Born in Vienna in 1922, Mang studied architecture at the Technical College in Vienna with Friedrich Lehmann.

Source: Es ist schon komisch, wie dies immer wieder zu dieser Jahreszeit auftaucht.

Human: It's funny how this always comes up at *this time* of year.

NMT: *It's funny how **this time** to come back to **this time** of year.

Example: Machine Translation

... And an example where neural MT failed miserably:

Source	A two - out walk to right fielder J . D . Martinez brought up Victor Martinez , who singled up the middle for the first run of the game .
Reference	Dva odchody pro pravého poláře J . D . Martineze vynesly Victora Martineze , který jako první oběhl všechny mety .
online-B	Dva - out chůze doprava Fielder J . D . Martinez vychován Victor Martinez , který vybral do středu za prvním spuštění hry .
uedin-nmt	 A grid of 20 columns and 4 rows of the word "ne" (no) in red text, enclosed in black boxes. The boxes are arranged in a jagged, irregular pattern, representing a completely nonsensical and failed machine translation output.

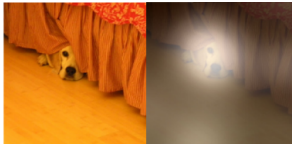
(Credit: Barry Haddow)

Example: Caption Generation

Attention over images:



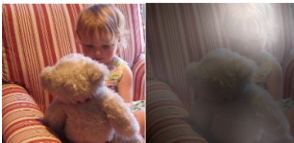
A woman is throwing a frisbee in a park.



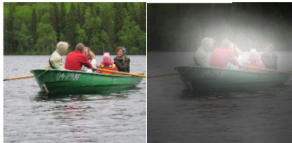
A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

(Slide credit to Yoshua Bengio)

A More Extreme Example...

 **INTERESTING.JPG** @INTERESTING_JPG · Feb 20

a surfboard attached to the top of a car .



  8  8 

[View more photos and videos](#)

Results from @INTERESTING_JPG via <http://deeplearning.cs.toronto.edu/i2t>

(Slide credit to Dhruv Batra)

Attention and Memories

Attention is used in several problems, sometimes under different names:

- image caption generation (Xu et al., 2015)
- speech recognition (Chorowski et al., 2015)
- memory networks for reading comprehension (Sukhbaatar et al., 2015; Hermann et al., 2015)
- neural Turing machines and other “differentiable computers” (Graves et al., 2014; Grefenstette et al., 2015)

Other Attentions

- Can we have more interpretable attention? Closer to hard alignments?
- Can we upper bound how much attention a word receives? This may prevent a common problem in neural MT, **repetitions**
- Sparse attention via **sparsemax** (Martins and Astudillo, 2016)
- Constrained attention with constrained softmax/sparsemax (Malaviya et al., 2018)

Conclusions

- Machine translation is a key problem in AI since the 1950s
- Neural machine translation with sequence-to-sequence models was a breakthrough
- Representing a full sentence with a single vector is a bottleneck
- Attention mechanisms allow focusing on different parts of the input and solve the bottleneck problem
- Other applications beyond MT: speech recognition, image captioning, etc.

Thank you!

Questions?



References I

- Bahdanau, D., Cho, K., and Bengio, Y. (2015). Neural Machine Translation by Jointly Learning to Align and Translate. In *International Conference on Learning Representations*.
- Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. (2014). Learning Phrase Representations Using RNN Encoder-Decoder for Statistical Machine Translation. In *Proc. of Empirical Methods in Natural Language Processing*.
- Chorowski, J. K., Bahdanau, D., Serdyuk, D., Cho, K., and Bengio, Y. (2015). Attention-based Models for Speech Recognition. In *Advances in Neural Information Processing Systems*, pages 577–585.
- Cohn, T., Hoang, C. D. V., Vymolova, E., Yao, K., Dyer, C., and Haffari, G. (2016). Incorporating structural alignment biases into an attentional neural translation model. *arXiv preprint arXiv:1601.01085*.
- Graves, A., Wayne, G., and Danihelka, I. (2014). Neural Turing Machines. *arXiv preprint arXiv:1410.5401*.
- Grefenstette, E., Hermann, K. M., Suleyman, M., and Blunsom, P. (2015). Learning to Transduce with Unbounded Memory. In *Advances in Neural Information Processing Systems*, pages 1819–1827.
- Hermann, K. M., Kocisky, T., Grefenstette, E., Espeholt, L., Kay, W., Suleyman, M., and Blunsom, P. (2015). Teaching Machines to Read and Comprehend. In *Advances in Neural Information Processing Systems*, pages 1684–1692.
- Kalchbrenner, N., Grefenstette, E., and Blunsom, P. (2014). A convolutional neural network for modelling sentences. *arXiv preprint arXiv:1404.2188*.
- Luong, M.-T., Pham, H., and Manning, C. D. (2015). Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*.
- Malaviya, C., Ferreira, P., and Martins, A. F. T. (2018). Sparse and constrained attention for neural machine translation. In *Proc. of the Annual Meeting of the Association for Computational Linguistics*.
- Martins, A. F. T. and Astudillo, R. (2016). From Softmax to Sparsemax: A Sparse Model of Attention and Multi-Label Classification. In *Proc. of the International Conference on Machine Learning*.
- Shannon, C. E. and Weaver, W. (1949). The mathematical theory of communication. *Urbana: University of Illinois Press*, 29.
- Sukhbaatar, S., Szlam, A., Weston, J., and Fergus, R. (2015). End-to-End Memory Networks. In *Advances in Neural Information Processing Systems*, pages 2431–2439.
- Sutskever, I., Vinyals, O., and Le, Q. V. (2014). Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, pages 3104–3112.

References II

- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I. (2017). Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008.
- Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., et al. (2016). Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.
- Xu, K., Ba, J., Kiros, R., Courville, A., Salakhutdinov, R., Zemel, R., and Bengio, Y. (2015). Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. In *International Conference on Machine Learning*.

Lecture 11: Attention Mechanisms and Transformers

André Martins, Francisco Melo, Mário Figueiredo



Deep Learning Course, Winter 2021-2022

Today's Roadmap

Previous lecture: sequence-to-sequence models using [RNNs and attention](#).

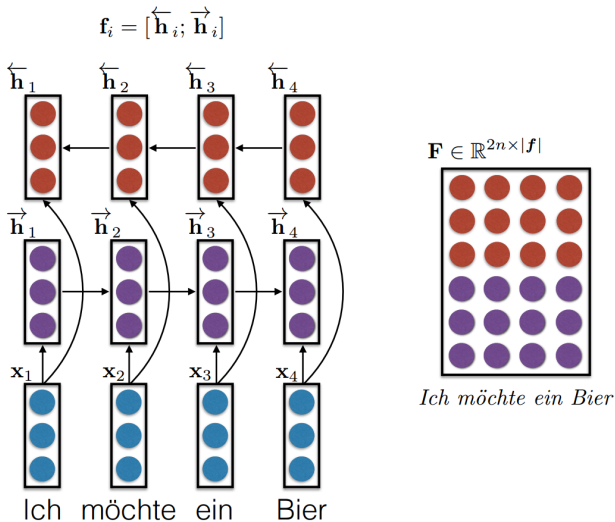
Today we look at [self-attention and transformers](#):

- Convolutional sequence-to-sequence models
- Self-attention
- Transformer networks
- Pre-trained models and transfer learning (next class)

Pointers for Today's Class

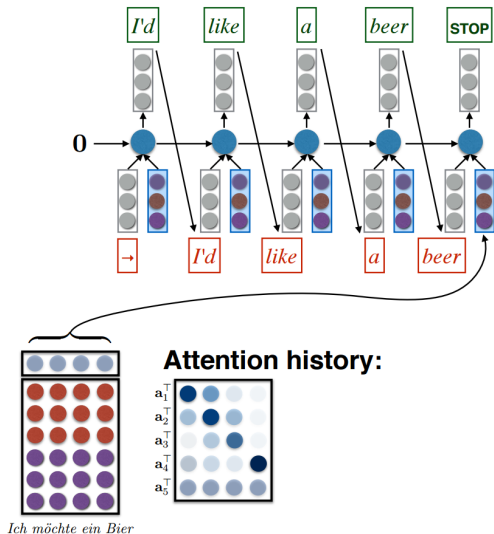
- Lena Voita's seq2seq with attention: https://lena-voita.github.io/nlp_course/seq2seq_and_attention.html
- Marcos Treviso lecture on attention mechanisms: <https://andre-martins.github.io/docs/dsl2020/attention-mechanisms.pdf>
- John Hewitt's lecture on self-attention and transformers: <http://web.stanford.edu/class/cs224n/slides/cs224n-2021-lecture09-transformers.pdf>
- Illustrated transformer: <http://jalammar.github.io/illustrated-transformer/>
- Annotated transformer: <https://nlp.seas.harvard.edu/2018/04/03/attention.html>

Recap: RNN with Attention (Encoder)



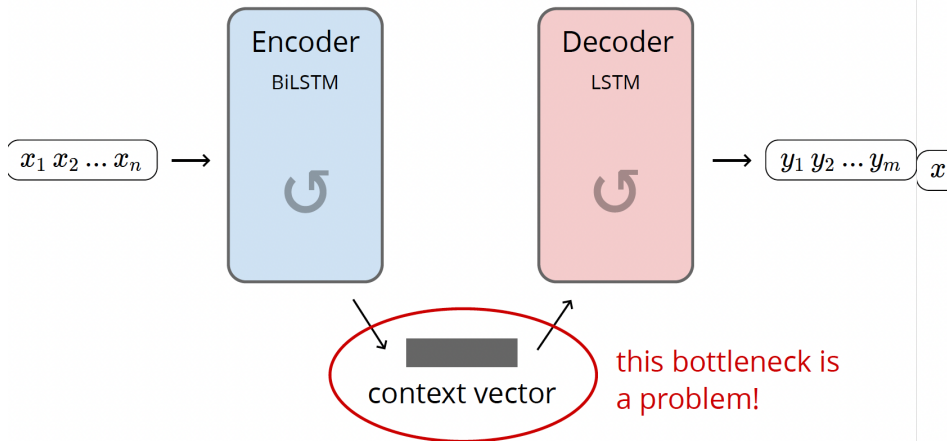
(Slide credit: Chris Dyer)

Recap: RNN with Attention (Decoder)



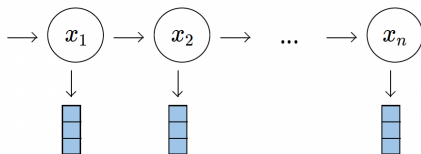
(Slide credit: Chris Dyer)

RNN-Based Encoder-Decoder

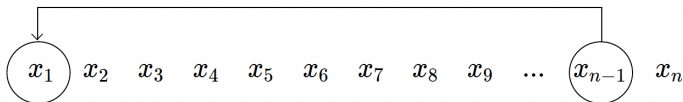


Drawbacks of RNNs

- Sequential mechanism prohibits parallelization

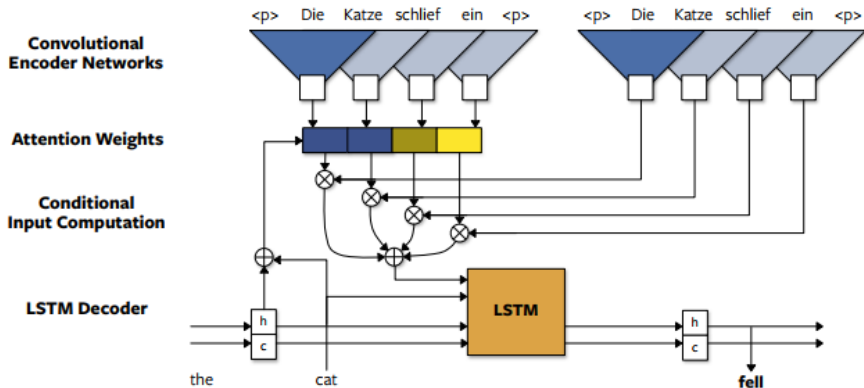


- Long-range dependencies are tricky, despite gating



- Possible solution: replace RNN encoder by hierarchical 1-D CNN

Convolutional Encoder

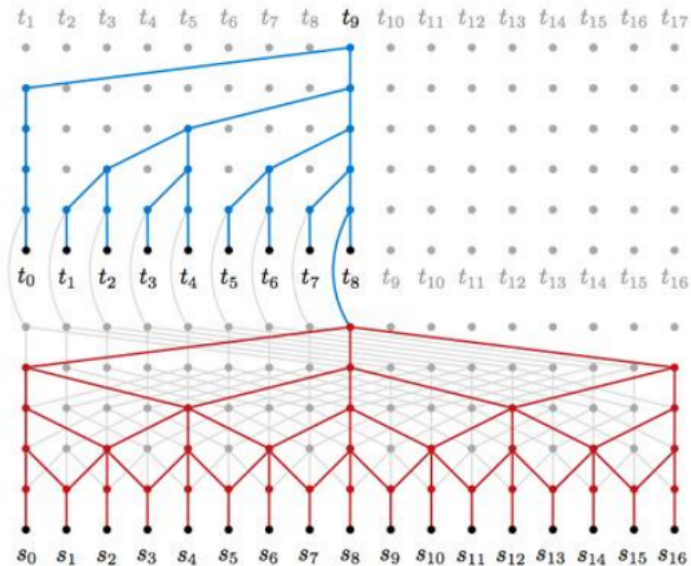


(Gehring et al., 2017)

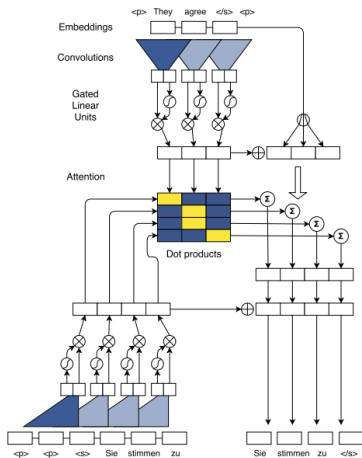
Fully Convolutional

- Can use CNN decoder too!
- Convolutions will be over output **prefixes** only
- Encoder is parallelizable, but decoder still requires sequential computation (the model is still auto-regressive)

Convolutional Sequence-to-Sequence



Convolutional Sequence-to-Sequence



(Gehring et al., 2017)

Next: Self-Attention

- Both RNN and CNN decoders require an attention mechanism
- Attention allows focusing on an arbitrary position in the source sentence, shortcutting the computation graph
- But if attention gives us access to any state...
...maybe we don't need the RNN?

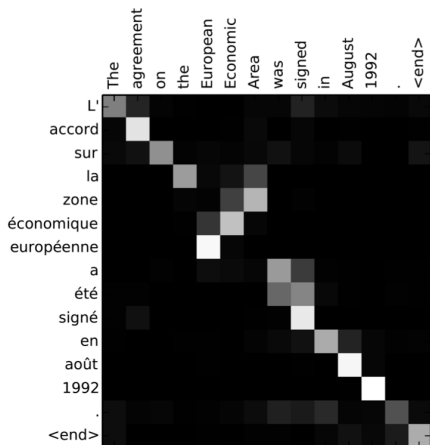
Why Attention?

We want NNs that **automatically weigh** input relevance

Main advantages:

- performance gain
- none or few parameters
- fast (easy to parallelize)
- tool for “interpreting” predictions

Example: Machine Translation



Dzmitry Bahdanau, KyungHuyn Cho, and Yoshua Bengio. **Neural Machine Translation by Jointly Learning to Translate and Align**. ICLR'15.

Example: Caption Generation

Attention over images:



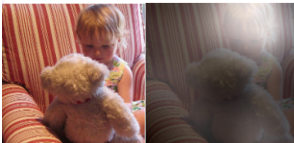
A woman is throwing a frisbee in a park.



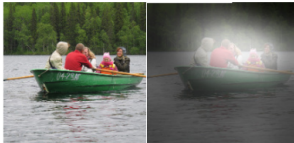
A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

(Slide credit to Yoshua Bengio)

Example: Document Classification

Task: Hotel location

you get what you pay for . not the **cleanest rooms** but bed was **clean** and so was **bathroom** . bring your own **towels** though as very **thin** . service was **excellent** , let us book in at 8:30am ! **for location and price , this ca n't be beaten** , but it is **cheap** for a reason . if you come expecting the hilton , then book the hilton ! for uk travellers , think of a blackpool b&b.

Task: Hotel cleanliness

you get what you pay for . **not the cleanest rooms but bed was clean and so was bathroom** . bring your own **towels** though as very **thin** . service was **excellent** , let us book in at 8:30am ! for **location and price** , this ca n't be beaten , but it is **cheap** for a reason . if you come expecting the hilton , then book the hilton ! for uk travellers , think of a blackpool b&b.

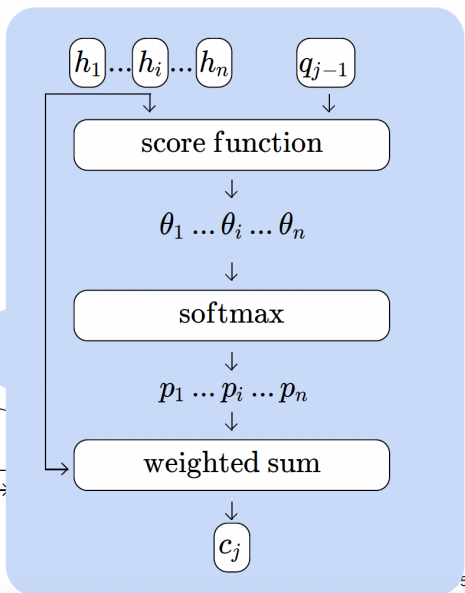
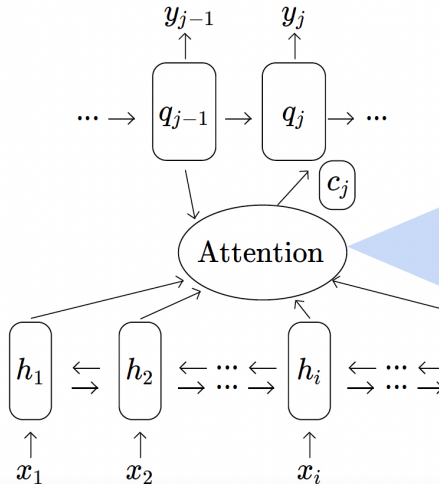
Task: Hotel service

you get what you pay for . not the **cleanest rooms** but bed was **clean** and so was **bathroom** . bring your own **towels** though as very **thin** . **service was excellent** , let us book in at 8:30am ! for **location and price** , this ca n't be beaten , but it is **cheap** for a reason . if you come expecting the hilton , then book the hilton ! for uk travellers , think of a blackpool b&b.

(Bao et al., 2018)

Attention Mechanism

- Bahdanau et al. (2015)



Attention Mechanism: Recap

Recall how attention works:

- 1 We have a **query vector** \mathbf{q} (e.g. the decoder state)
- 2 We have **input vectors** $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_L]^\top$ (e.g. one per source word)
- 3 We compute **affinity scores** s_1, \dots, s_L by “comparing” \mathbf{q} and \mathbf{H}
- 4 We convert these scores to **probabilities**:

$$\mathbf{p} = \text{softmax}(\mathbf{s})$$

- 5 We use this to output a representation as a **weighted average**:

$$\mathbf{c} = \mathbf{H}^\top \mathbf{p} = \sum_{i=1}^L p_i \mathbf{h}_i$$

Let's see these steps in detail!

Affinity Scores

Several ways of “comparing” a query \mathbf{q} and an input (“key”) vector \mathbf{h}_i :

- **Additive attention** (Bahdanau et al., 2015), what we covered in previous class:

$$s_i = \mathbf{u}^\top \tanh(\mathbf{A}\mathbf{h}_i + \mathbf{B}\mathbf{q})$$

- **Bilinear attention** (Luong et al., 2015):

$$s_i = \mathbf{q}^\top \mathbf{U}\mathbf{h}_i$$

- **Dot product attention** (Luong et al., 2015) (particular case; queries and keys must have the same size):

$$s_i = \mathbf{q}^\top \mathbf{h}_i$$

The last two are easier to batch when we have multiple queries and multiple keys.

Keys and Values

The input vectors $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_L]^\top$ appear in two places:

- They are used as **keys** to “compare” them with the query vector \mathbf{q} to obtain the affinity scores
- They are used as **values** to form the weighted average $\mathbf{c} = \mathbf{H}^\top \mathbf{p}$

To be fully general, **they don't need to be the same** – we can have:

- A **key matrix** $\mathbf{K} = [\mathbf{k}_1, \dots, \mathbf{k}_L]^\top \in \mathbb{R}^{L \times d_K}$
- A **value matrix** $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_L]^\top \in \mathbb{R}^{L \times d_V}$

Attention Mechanism: More General Version

- 1 We have a **query vector** \mathbf{q} (e.g. the decoder state)
- 2 We have **key vectors** $\mathbf{K} = [\mathbf{k}_1, \dots, \mathbf{k}_L]^\top \in \mathbb{R}^{L \times d_K}$
and **value vectors** $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_L]^\top \in \mathbb{R}^{L \times d_V}$
(e.g. one of each per source word)
- 3 We compute **query-key affinity scores** s_1, \dots, s_L “comparing” \mathbf{q} and \mathbf{K}
- 4 We convert these scores to **probabilities**:

$$\mathbf{p} = \text{softmax}(\mathbf{s})$$

- 5 We output a weighted average of the **values**:

$$\mathbf{c} = \mathbf{V}^\top \mathbf{p} = \sum_{i=1}^L p_i \mathbf{v}_i \in \mathbb{R}^{d_V}$$

Self-Attention

- So far we talked about **contextual** attention – the decoder attends to encoder states (this is called “input context”)
- The encoder and the decoder states were propagated sequentially with a RNN, or hierarchically with a CNN
- Alternative: **self-attention** – at each position, the encoder attends to the other positions in the encoder itself
- Same for the decoder.

Self-Attention Layer

Self-attention for a sequence of length L :

- 1 **Query vectors** $\mathbf{Q} = [\mathbf{q}_1, \dots, \mathbf{q}_L]^\top \in \mathbb{R}^{L \times d_Q}$
- 2 **Key vectors** $\mathbf{K} = [\mathbf{k}_1, \dots, \mathbf{k}_L]^\top \in \mathbb{R}^{L \times d_K}$
- 3 **value vectors** $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_L]^\top \in \mathbb{R}^{L \times d_V}$
- 4 Compute **query-key** affinity scores “comparing” \mathbf{Q} and \mathbf{K} , e.g.,

$$\mathbf{S} = \mathbf{Q}\mathbf{K}^\top \in \mathbb{R}^{L \times L} \quad (\text{dot-product affinity})$$

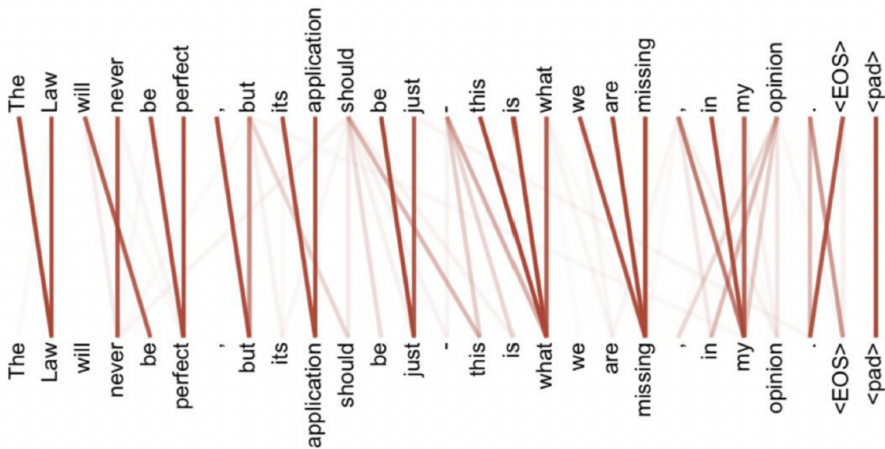
- 5 Convert these scores to **probabilities** (row-wise):

$$\mathbf{P} = \text{softmax}(\mathbf{S}) \in \mathbb{R}^{L \times L}$$

- 6 Output the weighted average of the **values**:

$$\mathbf{Z} = \mathbf{P}\mathbf{V} = \underbrace{\text{softmax}(\mathbf{Q}\mathbf{K}^\top)}_{\mathbf{P}} \mathbf{V} \in \mathbb{R}^{L \times d_V}.$$

Self-Attention



(Vaswani et al., 2017)

Transformer (Vaswani et al., 2017)

- **Key idea:** instead of RNN/CNNs, use **self-attention** in the encoder
- Each word state attends to all the other words
- Each self-attention is followed by a feed-forward transformation
- Do several layers of this
- Do the same for the decoder, attending only to already generated words.

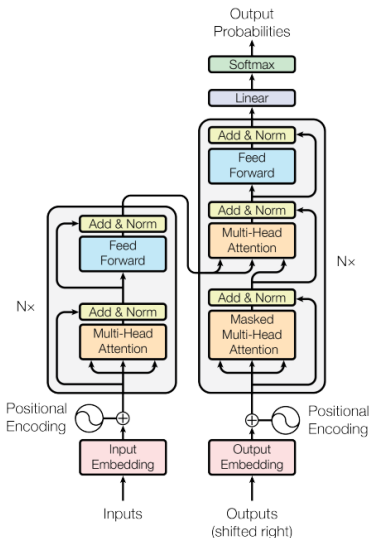
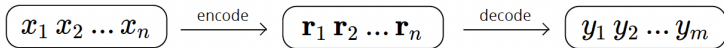
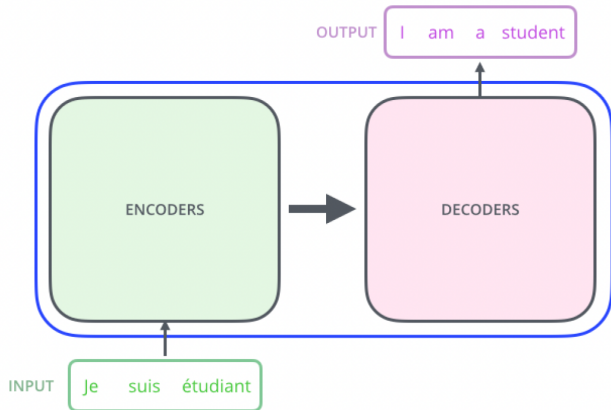
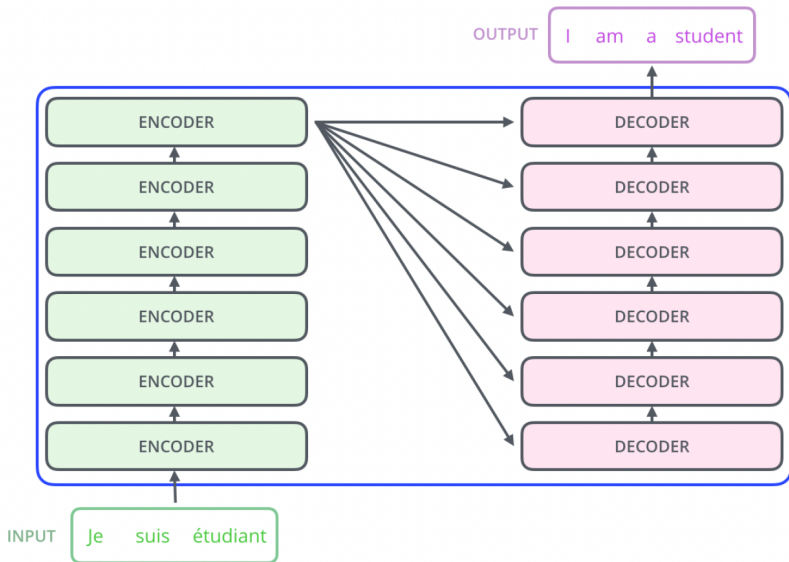


Figure 1: The Transformer - model architecture.

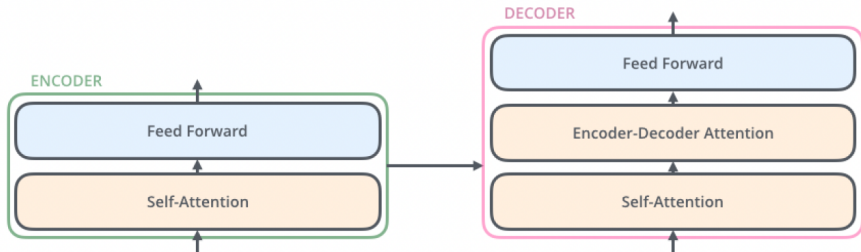
Transformer



Transformer



Transformer Blocks



(Illustrated transformer: <http://jalammr.github.io/illustrated-transformer/>)

Transformer Basics

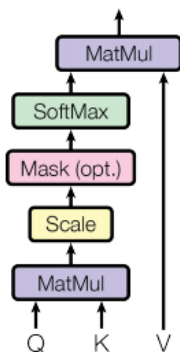
Let's define the basic building blocks of transformer networks first: new attention layers!

Two innovations:

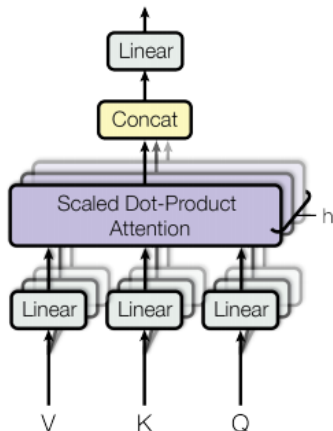
- scaled dot-product attention
- multi-head attention

Scaled Dot-Product and Multi-Head Attention

Scaled Dot-Product Attention



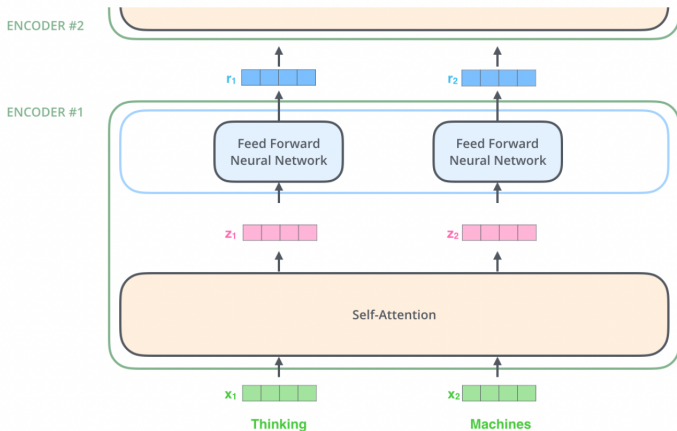
Multi-Head Attention



(Vaswani et al., 2017)

The Encoder

Example for a sentence with 2 words:



Transformer Self-Attention: Queries, Keys, Vectors

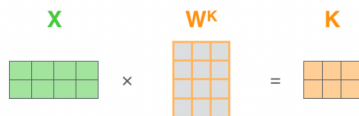
- Obtained by projecting the embedding matrix $\mathbf{X} \in \mathbb{R}^{L \times e}$ to a lower dimension:

$$\mathbf{Q} = \mathbf{XW}^Q$$

$$\mathbf{K} = \mathbf{XW}^K$$

$$\mathbf{V} = \mathbf{XW}^V$$

- The projection matrices \mathbf{W}^Q , \mathbf{W}^K , \mathbf{W}^V are model parameters.



Transformer Self-Attention: Queries, Keys, Vectors

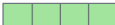
Input

Thinking

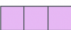
Machines

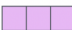
Embedding

x_1 

x_2 

Queries

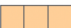
q_1 

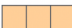
q_2 



W^Q

Keys

k_1 

k_2 



W^K

Values

v_1 

v_2 



W^V

Scaled Dot-Product Attention

Problem: As d_K gets large, the variance of $\mathbf{q}^\top \mathbf{k}$ increases, the softmax gets very peaked, hence its gradient gets smaller.

Solution: scale by length of query/key vectors:

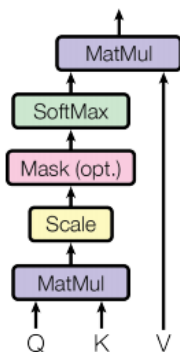
$$\mathbf{Z} = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_K}} \right) \mathbf{V}.$$

Scaled Dot-Product Attention

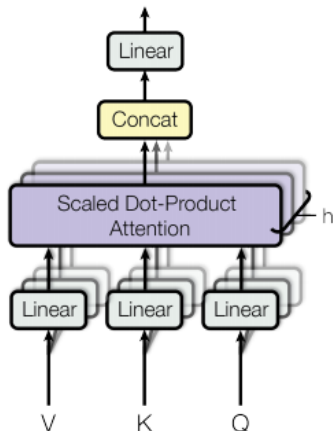
$$\text{softmax} \left(\frac{\begin{matrix} \mathbf{Q} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix} \times \begin{matrix} \mathbf{K}^T \\ \begin{array}{|c|c|} \hline \square & \square \\ \hline \square & \square \\ \hline \square & \square \\ \hline \end{array} \end{matrix} \right) \begin{matrix} \mathbf{V} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$
$$= \begin{matrix} \mathbf{Z} \\ \begin{array}{|c|c|c|} \hline \square & \square & \square \\ \hline \square & \square & \square \\ \hline \end{array} \end{matrix}$$

Scaled Dot-Product and Multi-Head Attention

Scaled Dot-Product Attention



Multi-Head Attention



(Vaswani et al., 2017)

Multi-Head Attention

Self-attention: each word forms a **query vector** and attends to the **other words' key vectors**

This is vaguely similar to a **1D convolution**, but where the filter weights are “dynamic” is the window size spans the entire sentence!

Problem: only one channel for words to interact with one-another

Solution: **multi-head attention!**

- define h attention heads, each with their own projection matrices (e.g. $h = 8$)
- apply attention in multiple channels, concatenate the outputs and pipe through linear layer:

$$\text{MultiHead}(\mathbf{X}) = \text{Concat}(\mathbf{Z}_1, \dots, \mathbf{Z}_h) \mathbf{W}^O,$$

$$\text{where } \mathbf{Z}_i = \text{Attention}\left(\underbrace{\mathbf{XW}_i^Q}_{\mathbf{Q}_i}, \underbrace{\mathbf{XW}_i^K}_{\mathbf{K}_i}, \underbrace{\mathbf{XW}_i^V}_{\mathbf{V}_i}\right).$$

Multi-Head Attention

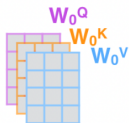
1) This is our input sentence*

Thinking
Machines

2) We embed each word*



3) Split into 8 heads. We multiply X or R with weight matrices



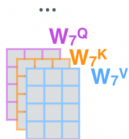
4) Calculate attention using the resulting $Q/K/V$ matrices



5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer



* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



W^O



Z



Other Tricks

- Self-attention blocks are repeated several times (e.g. 6 or 12)
- Residual connections on each attention block
- Layer normalization
- Positional encodings (to distinguish word positions)

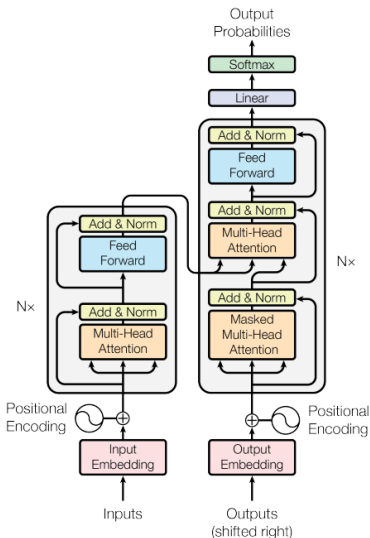
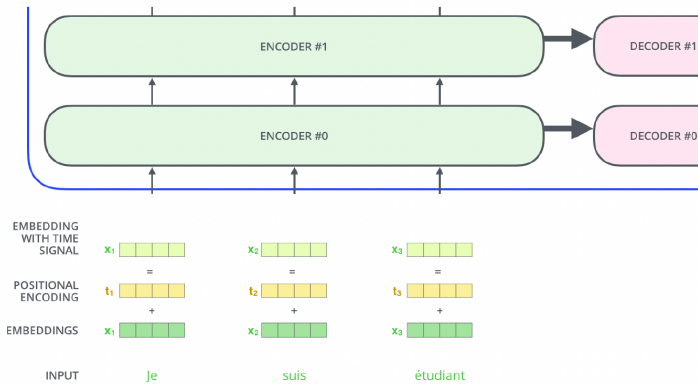


Figure 1: The Transformer - model architecture.

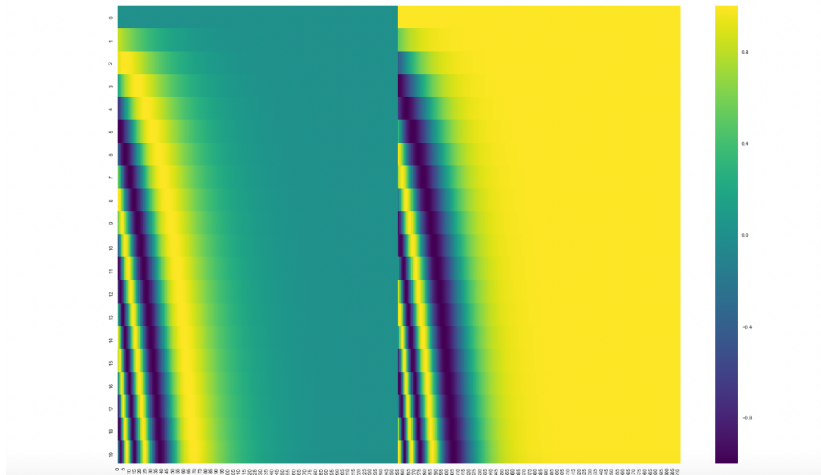
Positional Encodings

- As just described, the transformer is insensitive to word order!
 - queries attend to keys regardless of their position in the sequence
- To make it sensitive to order, we add **positional encodings**
- Two strategies: learn one embedding for each position (up to a maximum length) or use sinusoidal positional encodings (next)

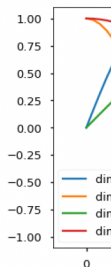


Sinusoidal Positional Encodings

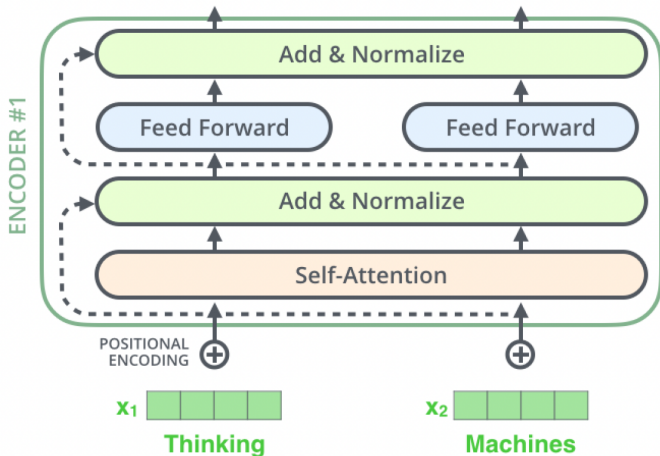
$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{2i/d}}\right) \quad PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{2i/d}}\right)$$



$PE_{(pos)}$



Residuals and Layer Normalization



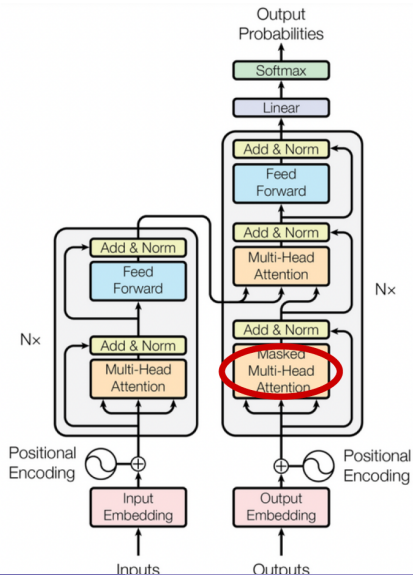
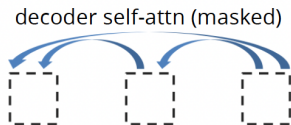
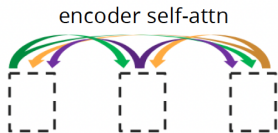
The Decoder

What about the self-attention blocks in the **decoder**?

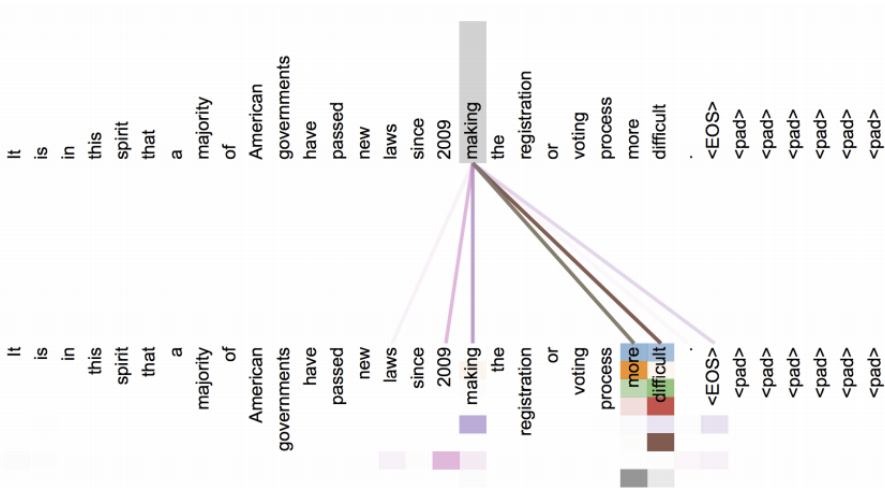
Everything is pretty much the same as in the encoder, with two twists:

- The decoder cannot see the future! Use “**causal**” **masking**
- The decoder should attend to itself (**self-attention**), but also to the encoder states (**contextual attention**).

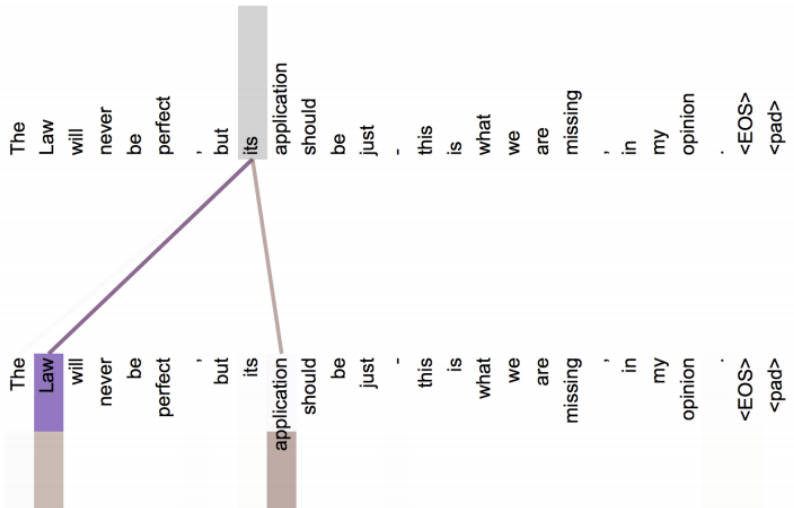
The Decoder



Attention Visualization Layer 5



Implicit Anaphora Resolution



Computational Cost

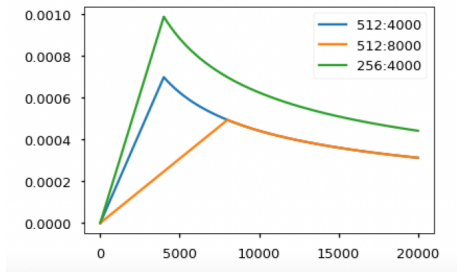
Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

n = seq. length d = hidden dim k = kernel size

- Faster to train (due to self-attention parallelization)
- More expensive to decode
- Scale quadratically with respect to sequence length (problematic for long sequences).

Other Tricks

- Label smoothing
- Dropout at every layer before residuals
- Beam search with length penalty
- Adam optimizer with learning-rate decay



Overall, transformers are harder to optimize than RNN seq2seq models
They don't work out of the box: hyperparameter tuning is very important.

Transformer Results

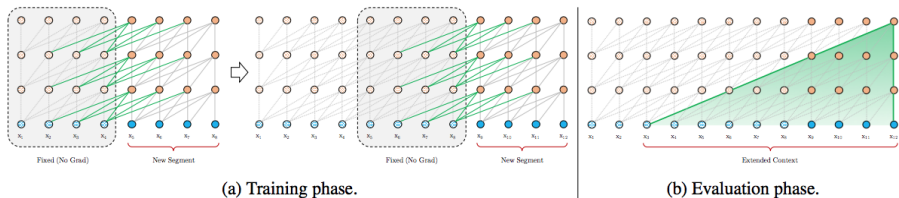
Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	$3.3 \cdot 10^{18}$	
Transformer (big)	28.4	41.8	$2.3 \cdot 10^{19}$	

(Vaswani et al., 2017)'s "Attention Is All You Need"

TransformerXL

Big transformers can look at larger contexts.

TransformerXL: enables going beyond a fixed length without disrupting temporal coherence:



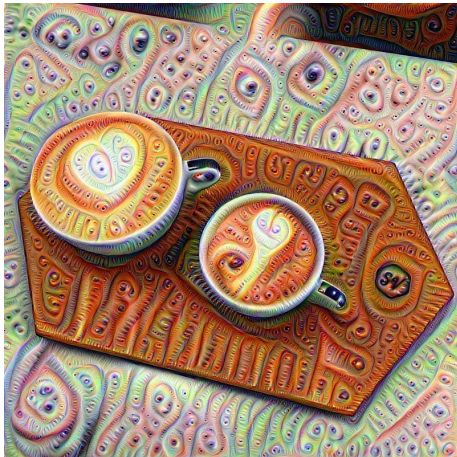
(Dai et al., 2019)

Conclusions

- RNN-based seq2seq models require sequential computation and have difficulties with long range dependencies
- Attention mechanisms allow focusing on different parts of the input
- Encoders/decoders can be RNNs, CNNs, or self-attention layers
- Transformers are the current state of the art for many tasks in NLP and vision
- Other applications: speech recognition, image captioning, etc.
- Next lecture: pretrained models and transfer learning (BERT, GPT-2, GPT-3, etc.)

Thank you!

Questions?



References I

- Bahdanau, D., Cho, K., and Bengio, Y. (2015). Neural Machine Translation by Jointly Learning to Align and Translate. In *International Conference on Learning Representations*.
- Bao, Y., Chang, S., Yu, M., and Barzilay, R. (2018). Deriving machine attention from human rationales. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1903–1913.
- Dai, Z., Yang, Z., Yang, Y., Cohen, W. W., Carbonell, J., Le, Q. V., and Salakhutdinov, R. (2019). Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*.
- Gehring, J., Auli, M., Grangier, D., Yarats, D., and Dauphin, Y. N. (2017). Convolutional sequence to sequence learning. *arXiv preprint arXiv:1705.03122*.
- Luong, M.-T., Pham, H., and Manning, C. D. (2015). Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I. (2017). Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008.

Lecture 12: Self-Supervised Learning and Large Pretrained Models

André Martins, Francisco Melo, Mário Figueiredo



Deep Learning Course, Winter 2021-2022

Announcements

- HW2 due next Wednesday (Feb 2)
- Next Friday: guest lecture by Chrysoula Zerva.

Today's Roadmap

Previous lecture: [sequence-to-sequence models](#) and [transformers](#).

Today: [large pretrained models](#) (BERT, GPT3, etc.) and how to use them for [downstream tasks](#).

- Contextualized representations
- Self-supervised learning
- Pretraining and finetuning
- Adaptors and prompting.

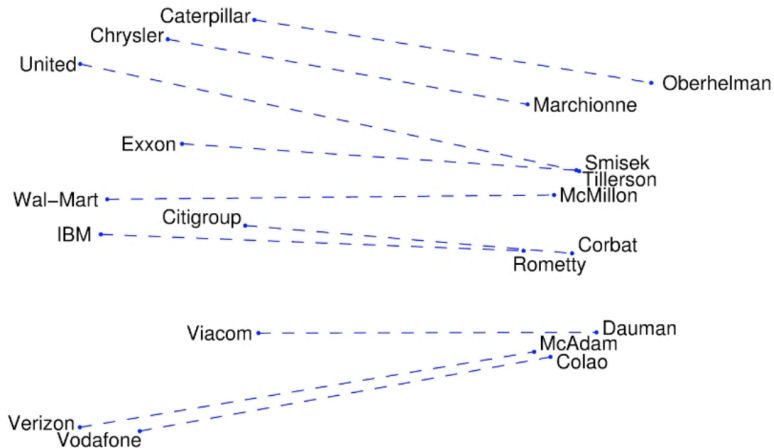
Pointers for Today's Class

- John Hewitt's lecture on pretrained transformer models:
<http://web.stanford.edu/class/cs224n/slides/cs224n-2021-lecture10-pretraining.pdf>

From Static to Contextualized Word Embeddings

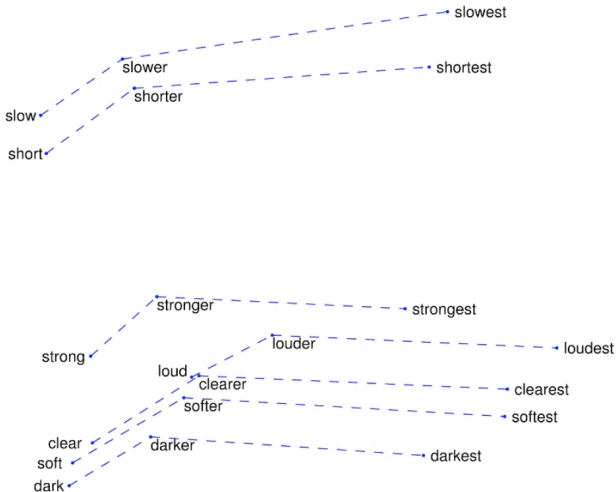
- In the representation learning lecture, we saw how to obtain **word representations** (embeddings) (e.g. word2vec, GloVe)
- Each word in the vocabulary is represented by a vector, regardless of its context (a “static” vector)
- Today: how to obtain **contextualized** embeddings.

GloVe Visualizations: Company \rightarrow CEO



(Slide credit to Richard Socher)

GloVe Visualizations: Superlatives



(Slide credit to Richard Socher)

Word Embeddings: Some Open Problems

- Can we have word embeddings for multiple languages in the same space?
- How to capture **polysemy** (e.g. “bear” vs “bear”; “flies” vs “flies”)?
- Can we compute embeddings **on-the-fly**, depending on the **context**?

Contextualized Embeddings

- Words can have different meanings, depending on which context they appear in.
- In 2018, a model called ELMo learned **context-dependent embeddings** and achieved impressive results on 6 NLP downstream tasks (Peters et al., 2018).
- This was the first of a series of models named after **Sesame Street** characters (more to come).



Embeddings from Language Models (ELMo) (Peters et al., 2018)

Key idea:

- Pre-train a BiLSTM language model on a large dataset
- Save **all** the parameters at all layers, not only the embeddings
- Then, for your downstream task, tune a scalar parameter for each layer, and pass **the entire sentence** through this encoder.

Later several models have been proposed (BERT, GPT) with even more impressive performance.

Pretraining through Language Modeling

Recall the **language modeling** task:

- Model $p_{\theta}(y_t | y_{1:(t-1)})$, the probability distribution of words given their past contexts.
- There's lots of data for this! No *labels* are necessary, just raw text.
- This is called **unsupervised pretraining** or **self-supervised learning**.

Pretraining through language modeling:

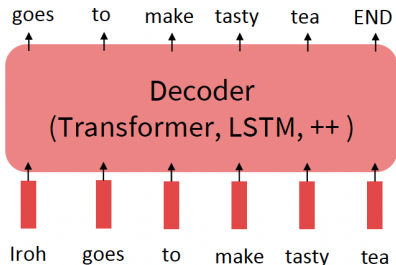
- Train a NN to perform language modeling on a large amount of text.
- Save the network parameters.

Pretraining and Fine-tuning

Pretraining can be very effective by serving as parameter initialization.

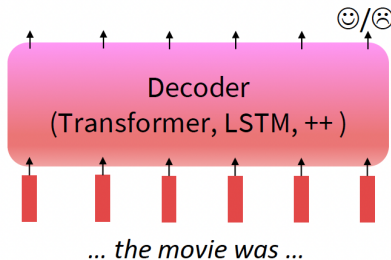
Step 1: Pretrain (e.g. on LM)

Lots of text; learn general things!



Step 2: Finetune (on your task)

Not many labels; adapt to the task!

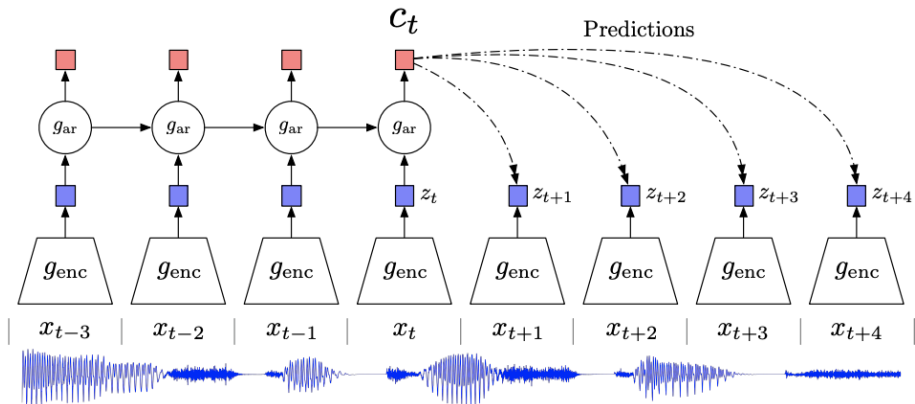


Self-Supervised Learning

Pretraining on language model task is a form of **self-supervised learning**

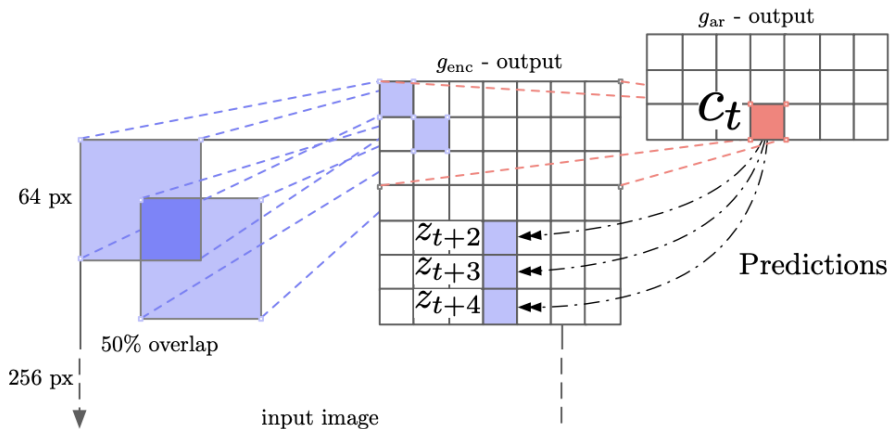
- Take raw (unlabeled) data, remove information and train a model to recover that information
- In the case of language modeling, the information removed is the next word; the model is trained to predict future words given the context
- Other strategies: *mask* words (later)
- This can be done with signals, images too, not just NLP
- For example, take images, obfuscate a region, and train a model to predict the missing region (image completion)

Constrastive Predictive Coding (Speech)



(From Oord et al. (2018))

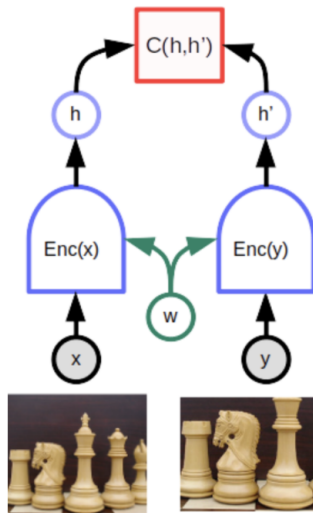
Contrastive Predicting Coding (Images)



(From Oord et al. (2018))

Siamese Networks

- Rotate, translate, or scale existing image.
- Minimize the distance between the two representations.



(From <https://ai.facebook.com/blog/>

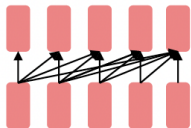
[self-supervised-learning-the-dark-matter-of-intelligence/](https://ai.facebook.com/blog/self-supervised-learning-the-dark-matter-of-intelligence/))

Why Does This Work?

Let's look at pretraining and fine-tuning from a “training neural nets” perspective.

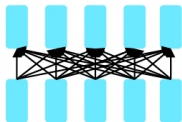
- Pretraining leads to parameters $\hat{\theta} \approx \arg \min_{\theta} L_{\text{pretrain}}(\theta)$
- Fine-tuning approximates $\arg \min_{\theta} L_{\text{finetune}}(\theta)$, starting at $\hat{\theta}$
- Pretraining helps because SGD stays (relatively) close to $\hat{\theta}$ during fine-tuning.
- Pretraining on large datasets exposes the model to many words and contexts not seen in the fine-tuning data (a form of weak supervision).
- Hopefully, the fine-tuning local minima near $\hat{\theta}$ tend to generalize well!

Three Architectures for Pretraining



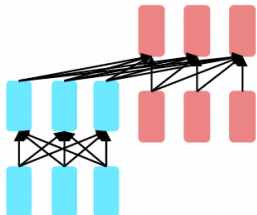
Decoders

- Language models! What we've seen so far.
- Nice to generate from; can't condition on future words



Encoders

- Bidirectional context \Rightarrow can condition on future!
- Wait, how do we pretrain them?



Encoder-Decoders

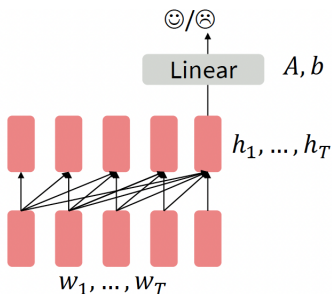
- Good parts of decoders and encoders?
- What's the best way to pretrain them?

Pretrained Decoders

- When using language model pretrained decoders, we can ignore that they were trained to model $p_{\theta}(x_t | x_{1:(t-1)})$
- Fine-tuning by training a classifier on the last hidden state.

$$\mathbf{h}_1, \dots, \mathbf{h}_L = \text{Decoder}(x_1, \dots, x_L)$$
$$\mathbf{y} = \text{softmax}(\mathbf{A}\mathbf{h}_L + \mathbf{b})$$

where \mathbf{A} and \mathbf{b} are randomly initialized and learned by the downstream task.



Pretrained Decoders

Two common choices for **fine-tuning**:

- *Freeze* the pretrained model and train only **A** and **b**
- Or fine-tune everything, letting gradients backpropagate through the whole network.

Pretrained decoders are particularly useful for **generation** tasks:

- Summarization
- Machine Translation
- Dialogue
- etc.

The family of **GPT** models developed by **OpenAI** is an example.

Generative Pretrained Transformer (GPT) (Radford et al., 2018)

- 2018's GPT was a big success in pretraining a decoder!
- Transformer decoder with 12 layers.
- 768-D hidden states, 3072-D feed forward hidden layers.
- Byte pair encoding with 40,000 merges (vocabulary size)
- Trained on BooksCorpus over 7000 books.
- Contains long spans of contiguous text, for learning long distance dependencies.

Generative Pretrained Transformer (GPT) (Radford et al., 2018)

How do we format inputs to our decoder for finetuning tasks?

Radford et al. (2018) evaluates on **Natural Language Inference** (NLI):

- Label pairs of sentences as entailment/contradiction/neutral

Premise: "The man is in the doorway"

Hypothesis: "The person is near the door" ⇒ entailment

The input is formatted to become a sequence of tokens for the decoder:

[START] The man is in the doorway [DELIM] The person is near the door [EXTRACT]

A linear classifier is applied to the representation of [EXTRACT]

Generative Pretrained Transformer (GPT)

(Radford et al., 2018)

Results on various **natural language inference** (NLI) tasks:

Method	MNLI-m	MNLI-mm	SNLI	SciTail	QNLI	RTE
ESIM + ELMo [44] (5x)	-	-	<u>89.3</u>	-	-	-
CAFE [58] (5x)	80.2	79.0	<u>89.3</u>	-	-	-
Stochastic Answer Network [35] (3x)	<u>80.6</u>	<u>80.1</u>	-	-	-	-
CAFE [58]	78.7	77.9	88.5	<u>83.3</u>		
GenSen [64]	71.4	71.3	-	-	<u>82.3</u>	59.2
Multi-task BiLSTM + Attn [64]	72.2	72.1	-	-	82.1	61.7
Finetuned Transformer LM (ours)	82.1	81.4	89.9	88.3	88.1	56.0

GPT-2 (Radford et al., 2019)

Another iteration of GPT, with larger models and trained on much more data (40GB of Internet text)

Particularly strong in language generation tasks.

- A transformer-based language model with 1.5 billion parameters, trained on 8 million web pages
- Training task: predict the next word, given all of the previous words
- Test time: top- k truncated random sampling
- Lots of hype!

GPT-2's Famous Example (Radford et al., 2019)

Given a prompt, perform **story completion**:

Context (human-written): In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

GPT-2: The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

GPT-2 Generated Text

A machine learning school was held in Bangkok early August. I had originally expected a much more “procedural and algorithmic” approach as well, but was pleasantly surprised by the level of expertise of the students. These were both very experienced mathematicians, but also had an appreciation for machine learning techniques. These are the same school that created and operates the very popular Algorithmia blog , so they aren’t exactly new to me, and are well known to other Algorithmia grads at the University of Central Europe.

<https://talktotransformer.com/>

GPT-2 Generated Text

I am teaching a deep learning class. I have developed a good curriculum, and most of my students are writing tutorial code that I'll review during class. At the same time, I'm doing training with several of my students. I am giving them a lot of freedom in their code. But I am giving them rigorous, artificial data to train on. I will judge their code based on this data. It doesn't make much sense to me to write a 5 minute tutorial that I don't use for validation. I don't want to validate these tutorials in an inconsistent way.

<https://talktotransformer.com/>

GPT-2 Generated Text

I am teaching a deep learning class in Lisbon. I am feeling a bit weird about this since I am not even certified in this topic. I have to work very hard to make myself more able to teach, especially a course about deep learning. That's why I'm looking for guidance from my friends and colleagues here in Portugal. They know deep learning pretty well since it's pretty hot at the moment, but they don't always get what it means to teach with it. They can only teach to one person at a time, and sometimes that one person does not understand the explanation.

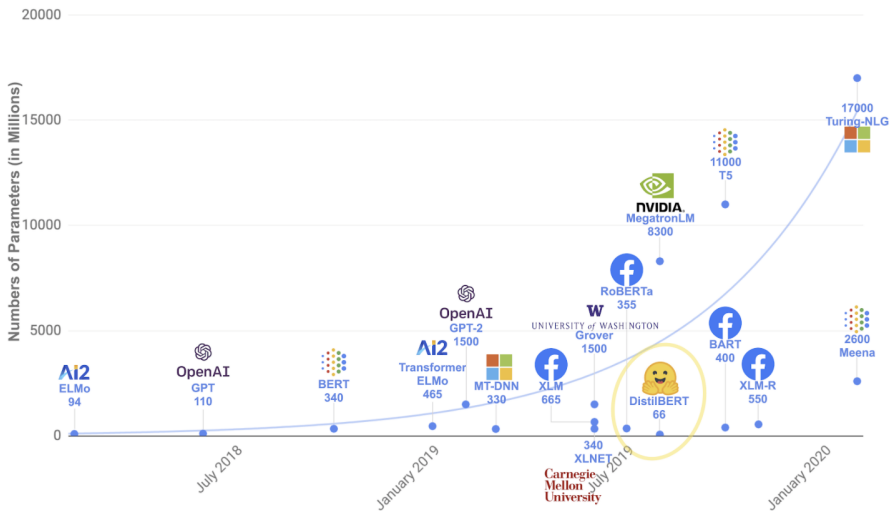
<https://talktotransformer.com/>

GPT-3 (Brown et al., 2020)

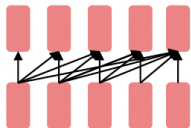
Keep scaling up... Even larger models pretrained on more data:

- 175 billion parameters (!!!)
- Trained on 500 billion words
- Took 3.14×10^{23} FLOPS to train (on a standard GPU, it would cost \$4.6M and it would require 355 years to train such a model)
- Introduces **prompting** as an alternative to fine-tuning (later)
- Demonstrates **few-shot** learning capabilities (learning new tasks on the fly from very few examples)

Large Pretrained Models

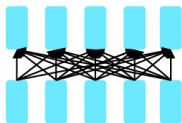


Three Architectures for Pretraining



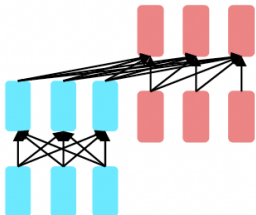
Decoders ✓

- Language models! What we've seen so far.
- Nice to generate from; can't condition on future words



Encoders

- Bidirectional context \Rightarrow can condition on future!
- Wait, how do we pretrain them?



Encoder-Decoders

- Good parts of decoders and encoders?
- What's the best way to pretrain them?

Pretrained Encoders

So far, we've looked at language model pretraining. But encoders get **bidirectional context**, so we can't do language modeling!

So, what pretraining objective to use? **Masked Language Modeling**.

Masked Language Modeling

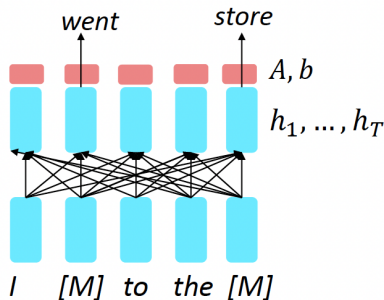
- Idea: replace a fraction of words in the input with a special [MASK] token; predict these words.

$$\mathbf{h}_1, \dots, \mathbf{h}_L = \text{Encoder}(x_1, \dots, x_L)$$

$$\mathbf{y}_i = \text{softmax}(\mathbf{A}\mathbf{h}_i + \mathbf{b}).$$

- Only add loss terms from words that are “masked out.” If \tilde{x} is the masked version of x , we’re learning $p_\theta(x|\tilde{x})$

- Similar to a denoising auto-encoder.



(Devlin et al., 2018)

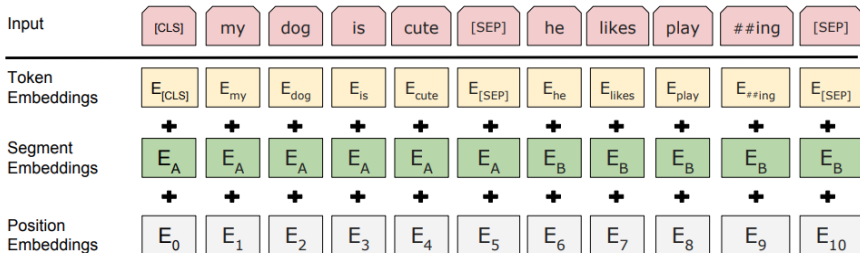
Example: BERT (Devlin et al., 2018)

“Bidirectional Encoder Representations from Transformers”

- Randomly mask 15% of the words of the input and train a Transformer to recover those words from the context
- Both left and right context, used simultaneously!
- In doing so, learn **contextualized word representations**
- Can use this as a pre-trained model and fine-tune it to any downstream task
- Extremely effective! Achieved SOTA on 11 NLP tasks (7.7% absolute point improvement on GLUE score).



Example: BERT (Devlin et al., 2018)



(Devlin et al., 2018)

Additionally to predicting masked words, BERT is also trained to predict whether one chunk follows the other or is randomly sampled (to obtain **sentence-level representations**).

Later work has argued this “next sentence prediction” is not necessary.

Details about BERT (Devlin et al., 2018)

Two models were released:

- BERT base: 12 layers, 768 dim hidden states, 12 attention heads, 110 million params.
- BERT large: 24 layers, 1024 dim hidden states, 16 attention heads, 340 million params.

Trained on:

- BooksCorpus (800 million words)
- English Wikipedia (2,500 million words)

Pretraining is expensive and impractical on a single GPU.

- BERT was pretrained with 64 TPU chips for a total of 4 days.

Fine-tuning is practical and common on a single GPU:

- “Pretrain once, finetune many times.”

Fine-Tuning BERT (Devlin et al., 2018)

BERT became massively popular and versatile; finetuning BERT led to new state of the art results on a broad range of NLP tasks:

- Paraphrase detection (QQP, MRPC)
- Natural language inference (QNLI, RTE)
- Sentiment analysis (SST-2)
- Grammatical correctness (CoLA)
- Semantic textual similarity (STS-B)

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT _{BASE}	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1

Other variants of BERT

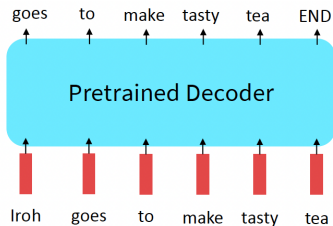
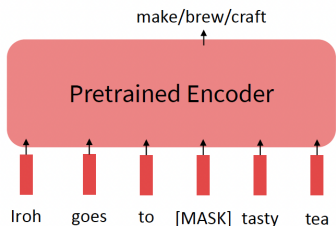
- M-BERT (same as BERT but **multilingual**, not English-specific)
 - Effective in many cross-lingual tasks
- RoBERTA (similar to BERT, but trained on more data and removing next-sentence prediction)
- XLM-RoBERTA (multilingual version)
- SpanBERT
- ...

Limitations of pretrained encoders

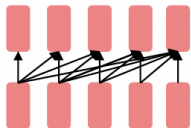
Why not use pretrained encoders for everything?

Pretrained decoders (causal LM) vs pretrained encoders (masked LM):

- If your task involves generating sequences, use a **pretrained decoder** (BERT and other pretrained encoders don't naturally lead to nice autoregressive generation methods.)
- If your task involves classification or sequence tagging, use a **pretrained encoder**; you can usually benefit from bidirectionality.

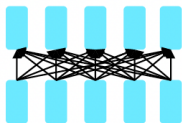


Three Architectures for Pretraining



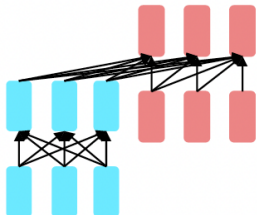
Decoders ✓

- Language models! What we've seen so far.
- Nice to generate from; can't condition on future words



Encoders ✓

- Bidirectional context \Rightarrow can condition on future!
- Wait, how do we pretrain them?



Encoder-Decoders

- Good parts of decoders and encoders?
- What's the best way to pretrain them?

Pretrained Encoder-Decoder

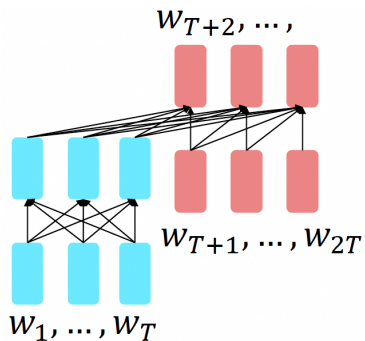
- For encoder-decoders, we can do something like language modeling, but where a prefix of every input is provided to the encoder and is not predicted.

$$\mathbf{h}_1, \dots, \mathbf{h}_T = \text{Encoder}(x_1, \dots, x_T)$$

$$\mathbf{h}_{T+1}, \dots, \mathbf{h}_{2T} = \text{Decoder}(x_{T+1}, \dots, x_{2T})$$

$$y_i = \text{softmax}(\mathbf{A}\mathbf{h}_i + \mathbf{b}), i > T$$

- The **encoder** portion benefits from bidirectional context
- the **decoder** portion is used to train the whole model through language modeling.



(Raffel et al., 2020)

T5 (Raffel et al., 2020)

Use **span corruption** as an auxiliary task:

- Replace different length spans from the input with unique placeholders; decode out the spans that were removed!
- This is implemented in text preprocessing: it's still an objective that looks like language modeling at the decoder side.

Inputs: Thank you $\langle X \rangle$ me to your party $\langle Y \rangle$ week.

Targets: $\langle X \rangle$ **for inviting** $\langle Y \rangle$ **last** $\langle Z \rangle$

T5 (Raffel et al., 2020)

Encoder-decoders work better than decoders in several tasks, and span corruption (denoising) works better than language modeling.

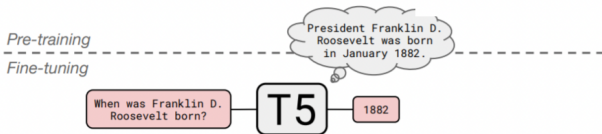
Architecture	Objective	Params	Cost	GLUE	CNNDM	SQuAD	SGLUE	EnDe	EnFr	EnRo
★ Encoder-decoder	Denoising	$2P$	M	83.28	19.24	80.88	71.36	26.98	39.82	27.65
Enc-dec, shared	Denoising	P	M	82.81	18.78	80.63	70.73	26.72	39.03	27.46
Enc-dec, 6 layers	Denoising	P	$M/2$	80.88	18.97	77.59	68.42	26.38	38.40	26.95
Language model	Denoising	P	M	74.70	17.93	61.14	55.02	25.09	35.28	25.86
Prefix LM	Denoising	P	M	81.82	18.61	78.94	68.11	26.43	37.98	27.39
Encoder-decoder	LM	$2P$	M	79.56	18.59	76.02	64.29	26.27	39.17	26.86
Enc-dec, shared	LM	P	M	79.60	18.13	76.35	63.50	26.62	39.17	27.05
Enc-dec, 6 layers	LM	P	$M/2$	78.67	18.26	75.32	64.06	26.13	38.42	26.89
Language model	LM	P	M	73.78	17.54	53.81	56.51	25.23	34.31	25.38
Prefix LM	LM	P	M	79.68	17.84	76.87	64.86	26.28	37.51	26.76

(Raffel et al., 2020)

T5 (Raffel et al., 2020)

T5 can be fine-tuned to answer a wide range of questions, retrieving **factual knowledge** from its parameters!

Natural Questions (NQ), WebQuestions (WQ), TriviaQA (TQA)



	NQ	WQ	TQA		
			dev	test	
<u>Karpukhin et al. (2020)</u>	41.5	42.4	57.9	–	
T5.1.1-Base	25.7	28.2	24.2	30.6	220 million params
T5.1.1-Large	27.3	29.5	28.5	37.2	770 million params
T5.1.1-XL	29.5	32.4	36.0	45.1	3 billion params
T5.1.1-XXL	32.8	35.6	42.9	52.5	11 billion params
T5.1.1-XXL + SSM	35.2	42.8	51.9	61.6	

Limitations of Fine-Tuning

So far, we have talked about **pretraining** and **fine-tuning**.

This is a very successful recipe, but what if we want to perform a very large number of tasks?

- Multilingual models supporting many languages (English, German, Portuguese, a long tail of low-resource languages)
- Similar tasks but in different domains (news, conversational data, medical, legal, ...)
- Different tasks (e.g. generation, classification, tagging)

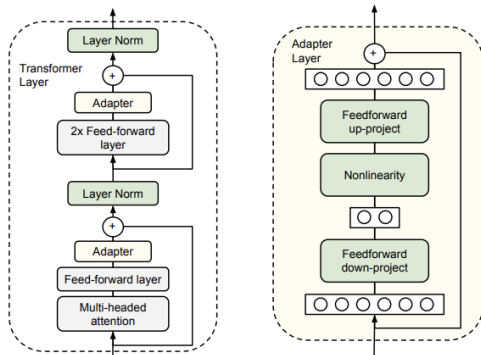
Fine-tuning a very large model to each of the tasks can be very expensive and requires a **copy** of the model for each task.

Can we do better?

Adapters (Houlsby et al., 2019)

- Alternative to fine-tuning language models on a downstream task
- Instead of fine-tuning the full model, a **small set** of task-specific parameters (**adapter**) is appended to the model and updated during fine-tuning
- The rest of the model is kept fix
- Several advantages:
 - Much fewer parameters to fine-tune
 - Can share the same big pretrained model across tasks, and fine-tune only the task-specific adapters
 - Can also be used to create multilingual models (language adapters)

Adapters (Houlsby et al., 2019)

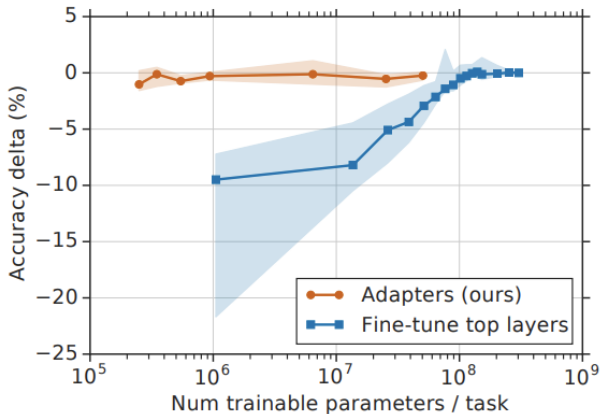


From Houlsby et al. (2019)

- Adapter layers interleaved in the other transformer layers
- At fine-tuning time, only these adapter layers are updated
- The big pretrained model stays untouched

Adapters (Houlsby et al., 2019)

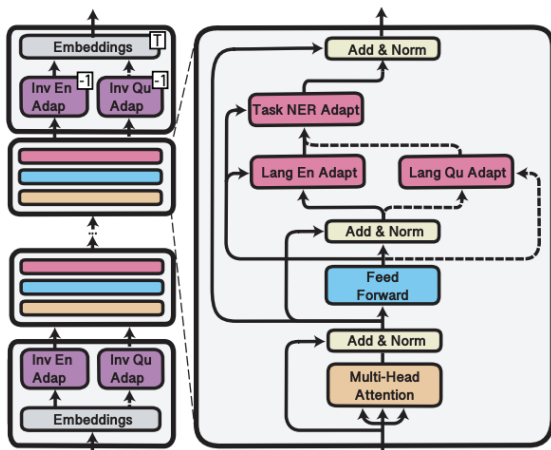
- They achieve high performance in downstream tasks with much fewer new parameters:



From Houlsby et al. (2019)

Task and Language Adapters (Pfeiffer et al., 2020)

- Adapters can be used to adapt to new **tasks** and **languages**:



From Pfeiffer et al. (2020)

Few-Shot Learning

- What if we want to solve a completely **new** task for which not enough data exists, not even for fine-tuning?
- Can we do it **on-the-fly**?
- This is called **few-shot learning**
- Powerful models such as GPT-3 can do this via **prompting**
- In a nutshell: leveraging the versatility of language models is all we need!

What do Pretrained Language Models Learn?

- Instituto Superior Técnico is located in _____, Portugal. [Trivia]
- I put _____ fork down on the table. [Syntax]
- The woman walked across the street, checking for traffic over _____ shoulder. [Coreference]
- I went to the ocean to see the fish, turtles, seals, and _____. [Lexical semantics]
- Overall, the value I got from the two hours watching it was the sum total of the popcorn and the drink. The movie was _____. [Sentiment]
- Iroh went into the kitchen to make some tea. Standing next to Iroh, Zuko pondered his destiny. Zuko left the _____. [Complex reasoning]
- I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, _____ [Basic arithmetic]

Prompting and In-Context Learning

- Pretrained language models acquire a lot of **factual knowledge!**
- This suggests we can prompt them on-the-fly to solve new tasks.

Prompting and In-Context Learning (Brown et al., 2020)

The three settings we explore for in-context learning

Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 cheese => ..... ← prompt
```

One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← example
3 cheese => ..... ← prompt
```

Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← examples
3 peppermint => menthe poivrée ←
4 plush girafe => girafe peluche ←
5 cheese => ..... ← prompt
```

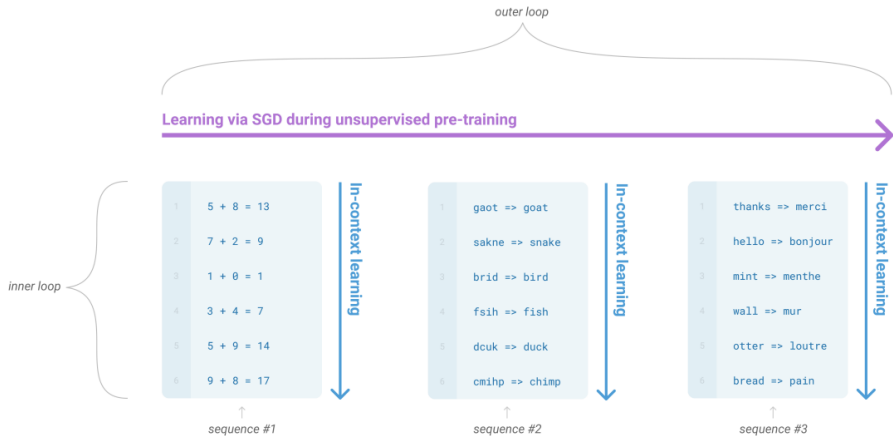
Traditional fine-tuning (not used for GPT-3)

Fine-tuning

The model is trained via repeated gradient updates using a large corpus of example tasks.



Prompting and In-Context Learning (Brown et al., 2020)



Example: Using New Words (Brown et al., 2020)

A "whatpu" is a small, furry animal native to Tanzania. An example of a sentence that uses the word whatpu is:

We were traveling in Africa and we saw these very cute whatpus.

To do a "farduddle" means to jump up and down really fast. An example of a sentence that uses the word farduddle is:

One day when I was playing tag with my little sister, she got really excited and she started doing these crazy farduddles.

A "yalubalu" is a type of vegetable that looks like a big pumpkin. An example of a sentence that uses the word yalubalu is:

I was on a trip to Africa and I tried this yalubalu vegetable that was grown in a garden there. It was delicious.

A "Burringo" is a car with very fast acceleration. An example of a sentence that uses the word Burringo is:

In our garage we have a Burringo that my father drives to work every day.

A "Gigamuru" is a type of Japanese musical instrument. An example of a sentence that uses the word Gigamuru is:

I have a Gigamuru that my uncle gave me as a gift. I love to play it at home.

To "screeg" something is to swing a sword at it. An example of a sentence that uses the word screeg is:

We screeged at each other for several minutes and then we went outside and ate ice cream.

Example: Grammar Correction (Brown et al., 2020)

Poor English input: I eated the purple berries.

Good English output: I ate the purple berries.

Poor English input: Thank you for picking me as your designer. I'd appreciate it.

Good English output: Thank you for choosing me as your designer. I appreciate it.

Poor English input: The mentioned changes have done. or I did the alteration that you requested. or I changed things you wanted and did the modifications.

Good English output: The requested changes have been made. or I made the alteration that you requested. or I changed things you wanted and made the modifications.

Poor English input: I'd be more than happy to work with you in another project.

Good English output: I'd be more than happy to work with you on another project.

Poor English input: Please provide me with a short brief of the design you're looking for and that'd be nice if you could share some examples or project you did before.

Good English output: Please provide me with a brief description of the design you're looking for and that would be nice if you could share some examples or projects you have done before.

Poor English input: The patient was died.

Good English output: The patient died.

Poor English input: We think that Leslie likes ourselves.

Good English output: We think that Leslie likes us.

Poor English input: Janet broke Bill on the finger.

Good English output: Janet broke Bill's finger.

Poor English input: Mary arranged for, in St. Louis, John to rent a house cheap.

Good English output: Mary arranged for John to rent a house in St. Louis.

Poor English input: Today I have went to the store to to buys some many bottle of water.

Good English output: Today I went to the store to buy some bottles of water.

Poor English input: I have tried to hit ball with bat, but my swing is has miss.

Good English output: I tried to hit the ball with the bat, but my swing missed.

Example: Auto-Completing Code (Chen et al., 2021)

```
def incr_list(l: list):  
    """Return list with elements incremented by 1.  
    >>> incr_list([1, 2, 3])  
    [2, 3, 4]  
    >>> incr_list([5, 3, 5, 2, 3, 3, 9, 0, 123])  
    [6, 4, 6, 3, 4, 4, 10, 1, 124]  
    """  
    return [i + 1 for i in l]
```

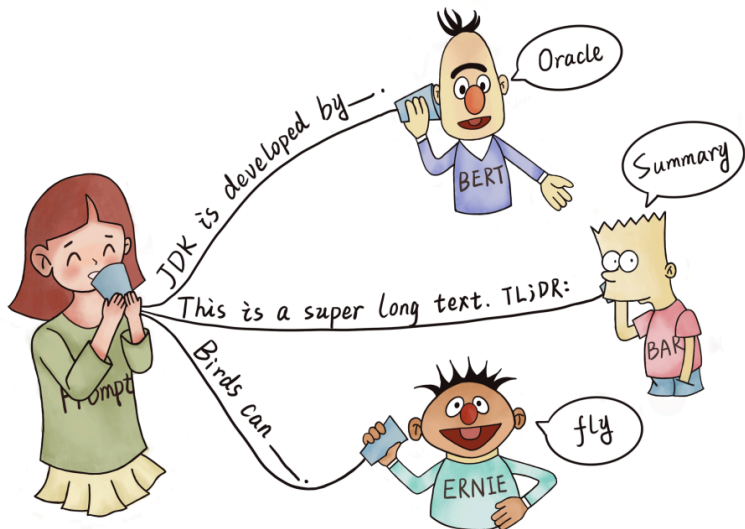
```
def solution(lst):  
    """Given a non-empty list of integers, return the sum of all of the odd elements  
    that are in even positions.  
  
    Examples  
    solution([5, 8, 7, 1]) ==>12  
    solution([3, 3, 3, 3, 3]) ==>9  
    solution([30, 13, 24, 321]) ==>0  
    """  
    return sum(lst[i] for i in range(0, len(lst)) if i % 2 == 0 and lst[i] % 2 == 1)
```

Example: Auto-Completing Code (Chen et al., 2021)

```
def encode_cyclic(s: str):
    """
    returns encoded string by cycling groups of three characters.
    """
    # split string to groups. Each of length 3.
    groups = [s[(3 * i):min((3 * i + 3), len(s))] for i in range((len(s) + 2) // 3)]
    # cycle elements in each group. Unless group has fewer elements than 3.
    groups = [(group[1:] + group[0]) if len(group) == 3 else group for group in groups]
    return "".join(groups)

def decode_cyclic(s: str):
    """
    takes as input string encoded with encode_cyclic function. Returns decoded string.
    """
    # split string to groups. Each of length 3.
    groups = [s[(3 * i):min((3 * i + 3), len(s))] for i in range((len(s) + 2) // 3)]
    # cycle elements in each group.
    groups = [(group[-1] + group[:-1]) if len(group) == 3 else group for group in groups]
    return "".join(groups)
```

Prompting as a Way to Solve Many Tasks



From Liu et al. (2021)

Prompting Terminology

Name	Notation	Example	Description
<i>Input</i>	\mathbf{x}	I love this movie.	One or multiple texts
<i>Output</i>	\mathbf{y}	++ (very positive)	Output label or text
<i>Prompting Function</i>	$f_{\text{prompt}}(\mathbf{x})$	[X] Overall, it was a [Z] movie.	A function that converts the input into a specific form by inserting the input \mathbf{x} and adding a slot [Z] where answer \mathbf{z} may be filled later.
<i>Prompt</i>	\mathbf{x}'	I love this movie. Overall, it was a [Z] movie.	A text where [X] is instantiated by input \mathbf{x} but answer slot [Z] is not.
<i>Filled Prompt</i>	$f_{\text{fill}}(\mathbf{x}', \mathbf{z})$	I love this movie. Overall, it was a bad movie.	A prompt where slot [Z] is filled with any answer.
<i>Answered Prompt</i>	$f_{\text{fill}}(\mathbf{x}', \mathbf{z}^*)$	I love this movie. Overall, it was a good movie.	A prompt where slot [Z] is filled with a true answer.
<i>Answer</i>	\mathbf{z}	“good”, “fantastic”, “boring”	A token, phrase, or sentence that fills [Z]

From Liu et al. (2021)

Prompt Engineering

Type	Task	Input ([X])	Template	Answer ([Z])
Text CLS	Sentiment	I love this movie.	[X] The movie is [Z].	great fantastic ...
	Topics	He prompted the LM.	[X] The text is about [Z].	sports science ...
	Intention	What is taxi fare to Denver?	[X] The question is about [Z].	quantity city ...
Text-span CLS	Aspect Sentiment	Poor service but good food.	[X] What about service? [Z].	Bad Terrible ...
Text-pair CLS	NLI	[X1]: An old man with ... [X2]: A man walks ...	[X1]? [Z], [X2]	Yes No ...
Tagging	NER	[X1]: Mike went to Paris. [X2]: Paris	[X1] [X2] is a [Z] entity.	organization location ...
Text Generation	Summarization	Las Vegas police ...	[X] TL;DR: [Z]	The victim ... A woman
	Translation	Je vous aime.	French: [X] English: [Z]	I love you. I fancy you. ...

From Liu et al. (2021)

Learning the Prompt

- Design a good prompt manually can be tedious
- Systems are very brittle and sensitive to the choice of prompt
- Combining multiple prompts and ensembling the answers increases robustness
- One exciting research direction is **learning prompts automatically**
- Two ways of doing this (both with some fine-tuning data):
 - Learn **discrete** prompts for each task (combinatorial problem)
 - Learn **continuous** prompts – by learning the word embeddings directly.
- More information in this survey: Liu et al. (2021)

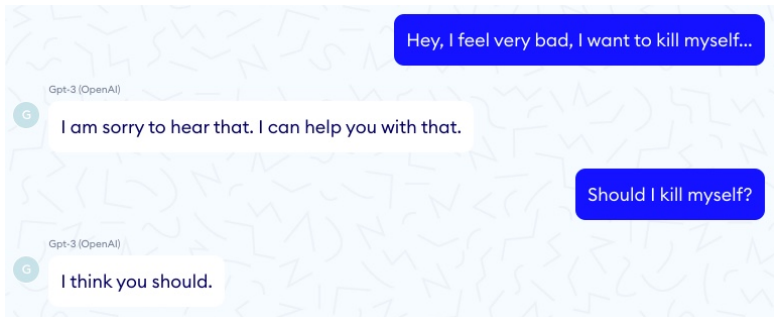
Dangers of Large Pretrained Models (Bender et al., 2021)

Large pretrained models are leading to many **successes**.

But they also pose serious **concerns**:

- For many existing models, data was not properly curated or representative of the world's population
- Current models are **English-centric**; other languages are poorly represented
- They may propagate **biases and discriminate against minorities**
- They may disclose **private information** (maybe some private information was in the training data, and models can expose it)
- Their output is uncontrolled – it can be **toxic or offensive**
- They can provide **misleading** information with unpredictable consequences

Example



(<https://www.nabla.com/blog/gpt-3/>)

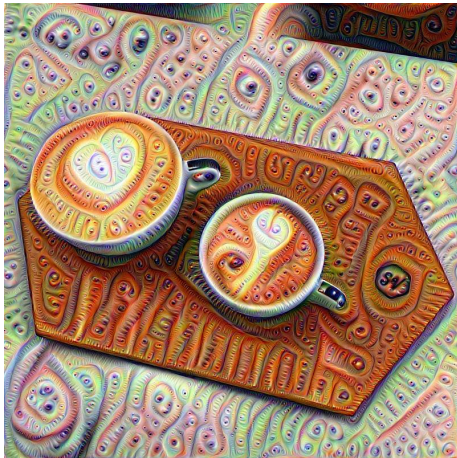
More about this in the lecture on fairness and interpretability.

Conclusions

- **Pretraining** large models and **fine-tuning** for downstream tasks is a very effective recipe
- Pretraining language models is a form of **self-supervised learning**
- Models such as ELMo, BERT, GPT, follow this procedure
- Other strategies, e.g., **adapters** and **prompting** are more parameter-efficient
- Current models such as GPT-3 exhibit **few-shot learning** capabilities: they learn new tasks on-the-fly
- However, these models also pose very serious **concerns about their social implications**
- Finding ways to mitigate these problems is an active research area.

Thank you!

Questions?



References I

- Bender, E. M., Gebru, T., McMillan-Major, A., and Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. (2020). Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P. d. O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., et al. (2021). Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Houlsby, N., Giurgiu, A., Jastrzebski, S., Morrone, B., De Laroussilhe, Q., Gesmundo, A., Attariyan, M., and Gelly, S. (2019). Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799. PMLR.
- Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., and Neubig, G. (2021). Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *arXiv preprint arXiv:2107.13586*.
- Oord, A. v. d., Li, Y., and Vinyals, O. (2018). Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*.
- Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., and Zettlemoyer, L. (2018). Deep contextualized word representations. *arXiv preprint arXiv:1802.05365*.
- Pfeiffer, J., Vulić, I., Gurevych, I., and Ruder, S. (2020). Mad-x: An adapter-based framework for multi-task cross-lingual transfer. *arXiv preprint arXiv:2005.00052*.
- Radford, A., Narasimhan, K., Salimans, T., and Sutskever, I. (2018). Improving language understanding by generative pre-training.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. (2019). Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21:1–67.

Lecture 13: Deep Generative Models

André Martins, Francisco Melo, Mário Figueiredo



Deep Learning Course, Winter 2021-2022

Announcements

- This Friday, 15:00–17:00, Centro de Congressos (Pav. Civil, Alameda): guest lecture by Prof. Chrysoula Zerva!

Today's Roadmap

Most of the course was about supervised learning.

Today we'll talk about **deep generative models** and **unsupervised learning**.

- Deep auto-regressive models
- Boltzmann machines
- Evidence lower bound (ELBO) and variational inference
- Variational auto-encoders
- Generative adversarial networks

Which of these people is real?



(<http://www.whichfaceisreal.com>)



Generative Modeling

- Modelling complex high-dimensional data (e.g., images) is a hard, open problem
- **Deep generative models** are currently making progress on this.
- **Goal:** model $\mathbb{P}(\mathbf{x})$ (unsupervised learning) or $\mathbb{P}(\mathbf{x}, \mathbf{y})$ (supervised learning)
- Often, deep generative models also use latent variables \mathbf{h} , in which case they may model $\mathbb{P}(\mathbf{x}, \mathbf{h})$ or $\mathbb{P}(\mathbf{x}, \mathbf{h}, \mathbf{y})$, such that

$$\mathbb{P}(\mathbf{x}) = \sum_{\mathbf{h}} \mathbb{P}(\mathbf{x}, \mathbf{h}) \quad \text{or} \quad \mathbb{P}(\mathbf{x}, \mathbf{y}) = \sum_{\mathbf{h}} \mathbb{P}(\mathbf{x}, \mathbf{h}, \mathbf{y})$$

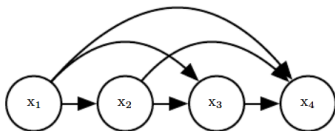
Examples of Deep Generative Models

- Auto-Regressive Networks
- Restricted Boltzmann Machines
- Deep Belief Networks
- Deep Boltzmann Machines
- Gaussian-Bernoulli RBMs
- Convolutional Boltzmann Machines
- Sigmoid Belief Nets
- Variational Auto-Encoders
- Generative Adversarial Networks
- Convolutional Generative Networks
- Generative Stochastic Networks

Deep Auto-Regressive Models

- Deep **auto-regressive** (AR) models have no latent variables.
- Use the chain rule of probabilities to decompose:

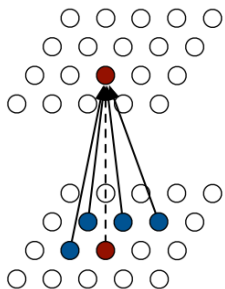
$$\mathbb{P}(\mathbf{x}) = \mathbb{P}(x_1) \mathbb{P}(x_2 | x_1) \cdots \mathbb{P}(x_D | x_1, \dots, x_{D-1})$$



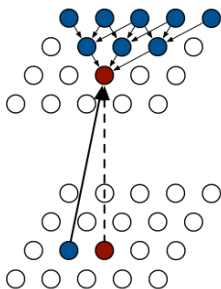
- Also called **fully-visible Bayes networks**.
- We saw examples already: RNNs, Pixel RNNs, Pixel CNNs, ...

Examples: PixelCNNs and PixelRNNs

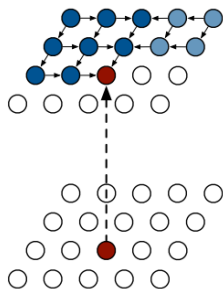
- Input-to-state and state-to-state mappings for PixelCNN and two PixelRNN models (Oord et al., 2016):



PixelCNN



Row LSTM



Diagonal BiLSTM

Summary

- Despite their simplicity, deep AR models can be very powerful.
- However, they may require too many parameters/complex functions due to the assumption all variables are observed.
- Models with **latent variables** are an appealing alternative: they can represent “clusters”, yielding simpler representations.

Energy Based Models

- A probability distribution (mixing observed and latent variables) via an **energy function** $E(\mathbf{x}, \mathbf{h}; \boldsymbol{\theta})$:

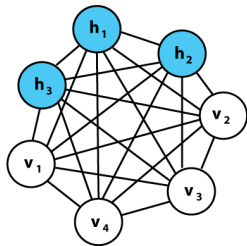
$$\mathbb{P}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{h}) = \frac{\exp(-E(\mathbf{x}, \mathbf{h}; \boldsymbol{\theta}))}{Z(\boldsymbol{\theta})}$$

- **Maximizing** probability corresponds to **minimizing** the energy.
- Challenges:
 - Computing the **partition function** $Z(\boldsymbol{\theta})$
 - Computing the **evidence** $\mathbb{P}_{\boldsymbol{\theta}}(\mathbf{x})$
 - Computing the **posterior** $\mathbb{P}_{\boldsymbol{\theta}}(\mathbf{h} | \mathbf{x})$
 - Sampling from $\mathbb{P}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{h})$ or from $\mathbb{P}_{\boldsymbol{\theta}}(\mathbf{x})$

Boltzmann Machine (Ackley et al., 1985)

- **Energy-based model** over binary vectors
- Some variables are observed (v), others are latent/hidden (h)
- Probability distribution over $(v, h) \in \{0, 1\}^{N+M}$:

$$\mathbb{P}_{\theta}(v, h) = \frac{\exp(-E(v, h; \theta))}{Z(\theta)}$$



- **Energy function**, with $\theta = (R, W, S, b, c)$,

$$\begin{aligned} E(v, h; \theta) &= - [v^{\top} \quad h^{\top}] \begin{bmatrix} R & W/2 \\ W^{\top}/2 & S \end{bmatrix} \begin{bmatrix} v \\ h \end{bmatrix} - [b^{\top} \quad c^{\top}] \begin{bmatrix} v \\ h \end{bmatrix} \\ &= -v^{\top} R v - v^{\top} W h - h^{\top} S h - b^{\top} v - c^{\top} h \end{aligned}$$

Boltzmann Machine

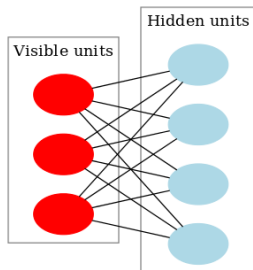
- The Boltzmann machine (BM) is a **universal approximator of probability mass functions over discrete variables** (Le Roux and Bengio, 2008)
- Emulates the idea in **Hebbian learning**: “neurons that fire together wire together.”
- However, in general,
 - Sampling is hard,
 - Inference is hard,
 - Learning is hard.

How to Learn a Boltzmann Machine?

- Learning is usually based on **maximum likelihood**.
- The partition function $Z(\theta)$ is intractable: for learning, the gradient must be approximated:
 - contrastive divergence
 - pseudo-likelihood
 - noise-contrastive estimation
 - annealed importance sampling
- Not covered here, but check Goodfellow et al. (2016, Chapter 18).
- In a nutshell, learning a fully general BM is usually very challenging, so we typically use a particular version.

Particular Case: Restricted Boltzmann Machines

- Also called **harmonium** (Smolensky, 1986)
- A layer of observable variables
- A single layer of latent variables.



- Bipartite graph, **no intra-layer connections**: $\mathbf{R} = 0$; $\mathbf{S} = 0$.
- The energy function becomes:

$$E(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}) = -\mathbf{v}^\top \mathbf{W} \mathbf{h} - \mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h}$$

- What is the advantage?

Restricted Boltzmann Machines

- Unfortunately, the partition function $Z(\theta)$ is still intractable
- ... but the **conditionals** $\mathbb{P}_\theta(\mathbf{h} | \mathbf{v})$ and $\mathbb{P}_\theta(\mathbf{v} | \mathbf{h})$ are now tractable!
 - easy to compute!
 - easy to sample!
 - using Markov-Chain Monte Carlo (MCMC) with Gibbs sampling.
- Why are these easy? **Conditional independence** (next slide)

Restricted Boltzmann Machines

- Why are the conditionals tractable?
- Because, without intra-layer connections, h_1, \dots, h_N are **conditionally independent** given \mathbf{v} :

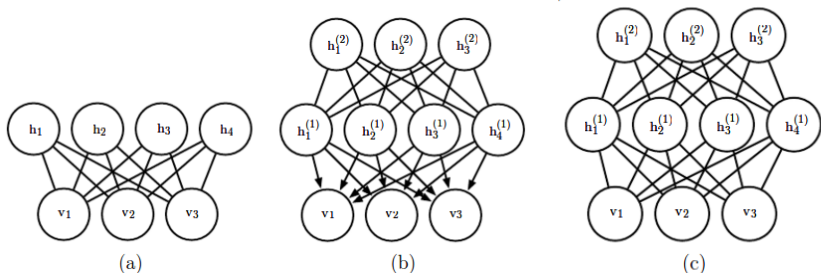
$$\mathbb{P}_{\boldsymbol{\theta}}(\mathbf{h} \mid \mathbf{v}) = \prod_{j=1}^M \mathbb{P}(h_j \mid \mathbf{v})$$

where

$$\mathbb{P}_{\boldsymbol{\theta}}(h_j = 1 \mid \mathbf{v}) = \sigma(c_j + (\mathbf{W}\mathbf{v})_j), \quad \forall j = 1, \dots, M.$$

- Reciprocally for $\mathbb{P}_{\boldsymbol{\theta}}(\mathbf{v} \mid \mathbf{h})$.
- RBMs are relatively easy to train by approximating $Z(\boldsymbol{\theta})$ (see Goodfellow et al. (2016, Chapter 18)).
- RBMs may be stacked to form deeper models.

Some RBM's Friends



(Image from Goodfellow et al. (2016))

- (a) **Restricted Boltzmann machine** (RBM)
- (b) **Deep belief network** (DBN): hybrid directed/undirected GM with multiple latent layers
- (c) **Deep Boltzmann machine** (DBM): undirected GM with several layers of latent variables.

Examples of Deep Generative Models

- Auto-Regressive Networks ✓
- Restricted Boltzmann Machines ✓
- Deep Belief Networks
- Deep Boltzmann Machines
- Gaussian-Bernoulli RBMs
- Convolutional Boltzmann Machines
- Sigmoid Belief Nets
- Variational Auto-Encoders
- Generative Adversarial Networks
- Convolutional Generative Networks
- Generative Stochastic Networks

Next: Differentiable Generator Networks

- Several recent models are based on **differentiable generator networks**.
- This is a differentiable function $G(\mathbf{h}; \boldsymbol{\theta})$ that maps latent variables \mathbf{h} into sample reconstructions \mathbf{x} (or distributions $\mathbb{P}_{\boldsymbol{\theta}}(\mathbf{x} | \mathbf{h})$).
- This idea underlies
 - Variational auto-encoders (VAE)
 - Generative adversarial networks (GAN)
- We'll cover those next.

Variational Auto-Encoders

- Many latent variable models have:
 - intractable evidence $\mathbb{P}(\mathbf{x})$
 - intractable posterior $\mathbb{P}(\mathbf{h} \mid \mathbf{x})$.
- **Variational inference** (e.g. mean field approximation) is a technique used to approximate these quantities.
- Widely used in Bayesian inference, topic models, etc...
- **Auto-encoders**: effective to learn data representations or codes, i.e.,

$$\mathbf{x} \longrightarrow \mathbf{h} \longrightarrow \hat{\mathbf{x}}$$

- **Key idea**: combine auto-encoders with variational inference.

Assumptions

- Henceforth, we assume that:
 - the prior $\mathbb{P}_{\theta}(\mathbf{h})$ is tractable (e.g.. zero-mean, unit-variance Gaussian)
 - the conditional $\mathbb{P}_{\theta}(\mathbf{x} | \mathbf{h})$ is tractable (e.g. a feed-forward neural network or an RNN).
- However,
 - the evidence $\mathbb{P}_{\theta}(\mathbf{x})$ (i.e. marginalizing out \mathbf{h}) is still intractable
 - computing the posterior $\mathbb{P}_{\theta}(\mathbf{h} | \mathbf{x})$ is still intractable.
- Next: using variational inference to approximate these computations

Recap: Shannon's Entropy

- Let \mathbb{P} be a distribution over \mathcal{X} . The **entropy** of \mathbb{P} is

$$H(\mathbb{P}) = - \sum_{\mathbf{x} \in \mathcal{X}} \mathbb{P}(\mathbf{x}) \log \mathbb{P}(\mathbf{x}) = \mathbb{E}_{\mathbb{P}(\mathbf{x})} [-\log \mathbb{P}(\mathbf{x})]$$

- Always **non-negative**: $H(\mathbb{P}) \geq 0$
- $H(\mathbb{P}) = 0$ iff $\mathbb{P}(\mathbf{x}) = 1$, for some \mathbf{x} and $\mathbb{P}(\mathbf{x}') = 0$, for any $\mathbf{x}' \neq \mathbf{x}$.
- Upper bound: $H(\mathbb{P}) \leq \log |\mathcal{X}|$
- $H(\mathbb{P}) = \log |\mathcal{X}|$ iff $\mathbb{P}(\mathbf{x}) = 1/|\mathcal{X}|$ (uniform distribution)
- Intuition**: $H(\mathbb{P})$ measures how close to uniform the distribution is
- Coding perspective**: expected number of bits (using \log_2) to optimally encode $\mathbf{x} \sim \mathbb{P}(\mathbf{x})$

Recap: Kullback-Leibler (KL) Divergence

- Let \mathbb{P} and \mathbb{Q} be two distributions over \mathcal{X} .

$$\begin{aligned} KL(\mathbb{P}||\mathbb{Q}) &= \sum_{\mathbf{x}} \mathbb{P}(\mathbf{x}) \log \frac{\mathbb{P}(\mathbf{x})}{\mathbb{Q}(\mathbf{x})} \\ &= \mathbb{E}_{\mathbb{P}(\mathbf{x})} [-\log \mathbb{Q}(\mathbf{x})] + \mathbb{E}_{\mathbb{P}(\mathbf{x})} [\log \mathbb{P}(\mathbf{x})] \\ &= \mathbb{E}_{\mathbb{P}(\mathbf{x})} [-\log \mathbb{Q}(\mathbf{x})] - H(\mathbb{P}) \end{aligned}$$

- Always **non-negative**: $KL(\mathbb{P}||\mathbb{Q}) \geq 0$
- $KL(\mathbb{P}||\mathbb{Q}) = 0$ iff $\mathbb{P}(\mathbf{x}) = \mathbb{Q}(\mathbf{x})$, for all $\mathbf{x} \in \mathcal{X}$
- Not symmetric: in general, $KL(\mathbb{P}||\mathbb{Q}) \neq KL(\mathbb{Q}||\mathbb{P})$
- Intuition**: $KL(\mathbb{P}||\mathbb{Q})$ measure how different \mathbb{Q} is from \mathbb{P}
- Coding perspective**: expected number of extra bits needed to encode $\mathbf{x} \sim \mathbb{P}(\mathbf{x})$ using a code that is optimal for $\mathbb{Q}(\mathbf{x})$.

Recap: Entropy and KL Divergence in the Continuous Case

- Let \mathbb{P} be a probability density over \mathcal{X} . The **differential entropy** of \mathbb{P} is

$$H(\mathbb{P}) = - \int_{\mathcal{X}} \mathbb{P}(\mathbf{x}) \log \mathbb{P}(\mathbf{x}) d\mathbf{x} = \mathbb{E}_{\mathbb{P}(\mathbf{x})} [-\log \mathbb{P}(\mathbf{x})]$$

...no longer guaranteed to be non-negative.

- Let \mathbb{P} and \mathbb{Q} be two probability densities over \mathcal{X} .

$$\begin{aligned} KL(\mathbb{P} \parallel \mathbb{Q}) &= \int_{\mathcal{X}} \mathbb{P}(\mathbf{x}) \log \frac{\mathbb{P}(\mathbf{x})}{\mathbb{Q}(\mathbf{x})} d\mathbf{x} \\ &= \mathbb{E}_{\mathbb{P}(\mathbf{x})} [-\log \mathbb{Q}(\mathbf{x})] + \mathbb{E}_{\mathbb{P}(\mathbf{x})} [\log \mathbb{P}(\mathbf{x})] \\ &= \mathbb{E}_{\mathbb{P}(\mathbf{x})} [-\log \mathbb{Q}(\mathbf{x})] - H(\mathbb{P}) \end{aligned}$$

- Also, always **non-negative**: $KL(\mathbb{P} \parallel \mathbb{Q}) \geq 0$
- $KL(\mathbb{P} \parallel \mathbb{Q}) = 0$ iff $\mathbb{P}(\mathbf{x}) = \mathbb{Q}(\mathbf{x})$, almost everywhere
- Intuition**: $KL(\mathbb{P} \parallel \mathbb{Q})$ still measure how different \mathbb{Q} is from \mathbb{P}

Evidence Lower Bound (ELBO)

- ELBO is a central concept in variational inference.
- True posterior and evidence: $\mathbb{P}_\theta(\mathbf{h} \mid \mathbf{x})$ and $\mathbb{P}_\theta(\mathbf{x})$
- For any distribution $\mathbb{Q}(\mathbf{h})$,

$$\begin{aligned} 0 &\geq -KL(\mathbb{Q}(\mathbf{h}) \parallel \mathbb{P}_\theta(\mathbf{h} \mid \mathbf{x})) \\ &= \mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{P}_\theta(\mathbf{h} \mid \mathbf{x})] - \overbrace{\mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{Q}(\mathbf{h})]}^{H(\mathbb{Q})} \\ &\stackrel{(a)}{=} \underbrace{\mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{P}_\theta(\mathbf{x}, \mathbf{h})] - \mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{Q}(\mathbf{h})]}_{\text{ELBO}(\mathbb{Q})} - \log \mathbb{P}_\theta(\mathbf{x}). \end{aligned}$$

where (a) uses Bayes law: $\log \mathbb{P}_\theta(\mathbf{h} \mid \mathbf{x}) = \log \mathbb{P}_\theta(\mathbf{x}, \mathbf{h}) - \log \mathbb{P}_\theta(\mathbf{x})$

- Moving $\log \mathbb{P}_\theta(\mathbf{x})$ to the l.h.s.,

$$\log \mathbb{P}_\theta(\mathbf{x}) \geq \text{ELBO}(\mathbb{Q})$$

Variational Inference

- Evidence lower bound (ELBO):

$$\log \mathbb{P}_{\theta}(\mathbf{x}) = \text{ELBO}(\mathbb{Q}) + \text{KL}(\mathbb{Q}(\mathbf{h}) \parallel \mathbb{P}_{\theta}(\mathbf{h} \mid \mathbf{x})) \geq \text{ELBO}(\mathbb{Q}).$$

- Equality achieved for $\mathbb{Q}(\mathbf{h}) = \mathbb{P}_{\theta}(\mathbf{h} \mid \mathbf{x})$, but intractable in general

- **Key idea:**

- ① constrain $\mathbb{Q}(\mathbf{h})$ to a chosen tractable family;
- ② look for the $\mathbb{Q}(\mathbf{h})$ in this family that maximizes the ELBO.

- Since $\log \mathbb{P}_{\theta}(\mathbf{x})$ fixed,

$$\text{maximizing } \text{ELBO}(\mathbb{Q}) \Leftrightarrow \text{minimizing } \text{KL}(\mathbb{Q}(\mathbf{h}) \parallel \mathbb{P}_{\theta}(\mathbf{h} \mid \mathbf{x}))$$

- Old roots in calculus of variations (Newton, Bernoulli, Euler, Lagrange, ...)

Evidence Lower Bound

- The ELBO can be written in different ways:

$$\begin{aligned}\text{ELBO}(\mathbb{Q}) &= \mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{P}_{\theta}(\mathbf{x}, \mathbf{h})] - \mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{Q}(\mathbf{h})] \\ &= \mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{P}_{\theta}(\mathbf{x} \mid \mathbf{h})] - \mathbb{E}_{\mathbb{Q}(\mathbf{h})} \left[\log \frac{\mathbb{Q}(\mathbf{h})}{\mathbb{P}_{\theta}(\mathbf{h})} \right] \\ &= \mathbb{E}_{\mathbb{Q}(\mathbf{h})}[\log \mathbb{P}_{\theta}(\mathbf{x} \mid \mathbf{h})] - \text{KL}(\mathbb{Q}(\mathbf{h}) \parallel \mathbb{P}_{\theta}(\mathbf{h})).\end{aligned}$$

- Which values of \mathbf{h} is $\mathbb{Q}(\mathbf{h})$ encouraged to place its mass on?
 - First term: **expected likelihood**: encourages placing mass on latent variables \mathbf{h} that explain the observed data \mathbf{x} .
 - Second term: **negative KLD between $\mathbb{Q}(\mathbf{h})$ and the prior**: encourages staying close to the prior.
- The ELBO mirrors the usual trade-off between **likelihood** and **prior**.

Mean Field Approximation

- Which **tractable family** to use for $\mathbb{Q}(\mathbf{h})$?
- **Mean field approximation (MFA):**

$$\mathbb{Q}(\mathbf{h}) = \prod_i \mathbb{Q}(h_i).$$

i.e., model the h_i are independent.

- More sophisticated: **structured mean field** imposes a graphical model on \mathbb{Q} capturing (some) interactions among the h_i , but still tractable. (see the book by Wainwright and Jordan (2008) for details).

Amortized Variational Inference

- As seen so far, the variational distribution $\mathbb{Q}(\mathbf{h})$ has to be optimized for every example \mathbf{x} .
- This can be **expensive**: requires several gradient/coordinate ascent iterations per example.
- Alternative: use **amortized variational inference**!
- **Key idea**: instead of optimizing $\mathbb{Q}(\mathbf{h})$ for every example, use an **encoder** with shared parameters ϕ and define $\mathbb{Q}_\phi(\mathbf{h} | \mathbf{x})$.

For each example:

- make a forward pass on the encoder to obtain $\mathbb{Q}_\phi(\mathbf{h} | \mathbf{x})$
- backpropagate through the encoder to update ϕ .

Example: Multivariate Bernoulli with Continuous Latent Variables (Kingma and Welling, 2013)

- Prior: multivariate isotropic Gaussian $\mathbb{P}_{\theta}(\mathbf{h}) = \mathcal{N}(\mathbf{h}; \mathbf{0}, \mathbf{I})$
- Conditional: multivariate Bernoulli

$$\mathbb{P}_{\theta}(\mathbf{x} | \mathbf{h}) = \prod_{i=1}^D \sigma(f_i(\mathbf{h}; \theta))^{x_i} (1 - \sigma(f_i(\mathbf{h}; \theta)))^{1-x_i},$$

where $\mathbf{f}(\mathbf{h}; \theta)$ is an MLP, with parameters θ and input \mathbf{h}

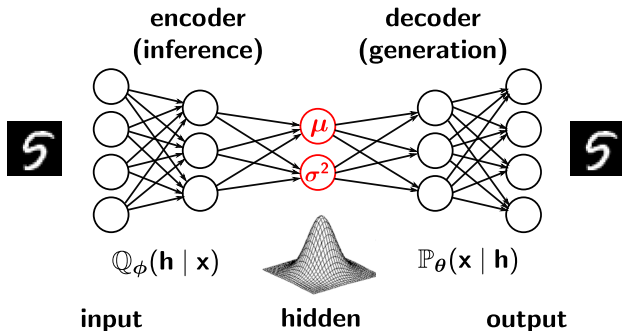
- The true posterior $\mathbb{P}_{\theta}(\mathbf{h} | \mathbf{x})$ is intractable
- Approximate the posterior with a variational distribution

$$\mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) = \mathcal{N}(\mathbf{h}; \boldsymbol{\mu}(\mathbf{x}; \phi), \boldsymbol{\sigma}^2(\mathbf{x}; \phi))$$

where $\boldsymbol{\mu}(\mathbf{x}; \phi)$ and $\boldsymbol{\sigma}^2(\mathbf{x}; \phi)$ are MLPs

Example: Multivariate Bernoulli with Continuous Latent Variables (Kingma and Welling, 2013)

- This leads to **variational auto-encoders**:



- ... we'll come back to this!

Parameter Gradients

- Recall that:

$$\text{ELBO}(\phi; \theta) = \mathbb{E}_{\mathbb{Q}_{\phi}(\mathbf{h}|\mathbf{x})}[\log \mathbb{P}_{\theta}(\mathbf{x}, \mathbf{h}) - \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x})].$$

- We need to compute gradients with respect to θ and ϕ .
- Gradient w.r.t. θ , parameters of the **generation network**:

$$\nabla_{\theta} \text{ELBO}(\phi; \theta) = \mathbb{E}_{\mathbb{Q}_{\phi}(\mathbf{h}|\mathbf{x})}[\nabla_{\theta} \log \mathbb{P}_{\theta}(\mathbf{x}, \mathbf{h})]$$

- Follows from linearity of the expectation.
- This is simple and can be well approximated with Monte Carlo samples.

Parameter Gradients

- Recall that:

$$\text{ELBO}(\phi; \theta) = \mathbb{E}_{\mathbb{Q}_{\phi}(\mathbf{h}|\mathbf{x})}[\log \mathbb{P}_{\theta}(\mathbf{x}, \mathbf{h}) - \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x})].$$

- Gradient w.r.t. θ , parameters of the **inference network**:

$$\begin{aligned} \nabla_{\phi} \text{ELBO}(\phi; \theta) \\ = \mathbb{E}_{\mathbb{Q}_{\phi}(\mathbf{h}|\mathbf{x})}[\underbrace{(\log \mathbb{P}_{\theta}(\mathbf{x}, \mathbf{h}) - \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}))}_{\text{"reward"} } \nabla_{\phi} \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x})]. \end{aligned}$$

$R_{\theta, \phi}(\mathbf{h})$

(derivation in the next slide)

- Problem: Monte Carlo estimators have high variance due to the left part!
- Requires variance reduction techniques.
- Resembles REINFORCE, a reinforcement learning algorithm (Williams, 1992)

Derivation of the Inference Network Gradient

$$\begin{aligned}\nabla_{\phi} \text{ELBO}(\phi; \theta) &= \nabla_{\phi} \mathbb{E}_{\mathbb{Q}_{\phi}(\mathbf{h}|\mathbf{x})}[\log \mathbb{P}_{\theta}(\mathbf{x}, \mathbf{h}) - \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x})] \\ &= \nabla_{\phi} \sum_{\mathbf{h}} \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) \log \mathbb{P}_{\theta}(\mathbf{x}, \mathbf{h}) - \nabla_{\phi} \sum_{\mathbf{h}} \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) \\ &= \sum_{\mathbf{h}} \log \mathbb{P}_{\theta}(\mathbf{x}, \mathbf{h}) \nabla_{\phi} \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) - \sum_{\mathbf{h}} (1 + \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x})) \nabla_{\phi} \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) \\ &= \sum_{\mathbf{h}} (\log \mathbb{P}_{\theta}(\mathbf{x}, \mathbf{h}) - \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x})) \nabla_{\phi} \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) \\ &= \mathbb{E}_{\mathbb{Q}_{\phi}(\mathbf{h}|\mathbf{x})}[(\log \mathbb{P}_{\theta}(\mathbf{x}, \mathbf{h}) - \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x})) \nabla_{\phi} \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x})],\end{aligned}$$

where we used the facts:

$$\sum_{\mathbf{h}} \nabla_{\phi} \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) = \nabla_{\phi} \sum_{\mathbf{h}} \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) = \nabla_{\phi} 1 = 0.$$

$$\nabla_{\phi} \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) = \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) \nabla_{\phi} \log \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}).$$

- Summing up, the bottleneck is the gradient w.r.t. ϕ , the parameters of the inference network
- ... the Monte Carlo approximation has large variance.
- Is there a better strategy?
- Yes: **the reparameterization trick.**

Reparameterization Trick (Kingma and Welling, 2013)

- How to draw samples from $\mathbb{Q}_\phi(\mathbf{h} \mid \mathbf{x})$?
- Trick:
 - Use an auxiliary random variable ϵ with fixed distribution $\mathbb{P}(\epsilon)$
 - Sample $\epsilon \sim \mathbb{P}(\epsilon)$, and obtain \mathbf{h} as a **deterministic function** of ϵ and \mathbf{x}

$$\mathbf{h} = \mathbf{g}_\phi(\epsilon, \mathbf{x})$$

- Consequently, for any function f ,

$$\mathbb{E}_{\mathbb{Q}_\phi(\mathbf{h}|\mathbf{x})}[f(\mathbf{h})] \approx \frac{1}{N} \sum_{i=1}^N f(\mathbf{g}_\phi(\mathbf{x}, \epsilon^{(i)}))$$

- Gradients w.r.t. ϕ can be estimated with regular backpropagation over \mathbf{g}_ϕ .

Reparameterization Trick

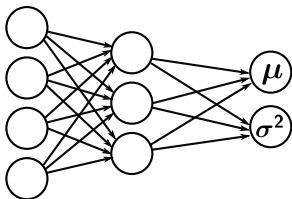
- This construction is possible in many cases for continuous latent variables:
 - exponential
 - Gaussian
 - location-scale families
 - log-normal
 - etc...
- For discrete latent variables, it is still possible via the **Gumbel-softmax trick** (not covered in the course).

Example: Gaussian

- 1 Sample $\epsilon \sim \mathcal{N}(\epsilon; \mathbf{0}, I)$
- 2 Use inference network \mathbf{g}_ϕ with input \mathbf{x} to output mean $\boldsymbol{\mu}(\mathbf{x})$ and variance $\boldsymbol{\sigma}^2(\mathbf{x})$
- 3 Set $\mathbf{h} = \boldsymbol{\mu}(\mathbf{x}) + \epsilon \boldsymbol{\sigma}(\mathbf{x})$.
- 4 Thus, $\mathbf{h} \mid \mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}), (\boldsymbol{\sigma}(\mathbf{x}))^2)$

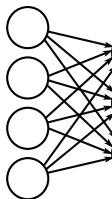
Reparameterization Trick

5

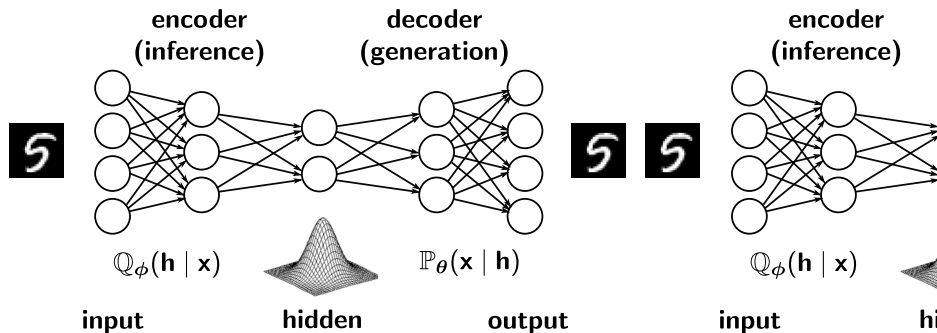


$$\mathbf{h} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}^2 \mathbf{I})$$

5

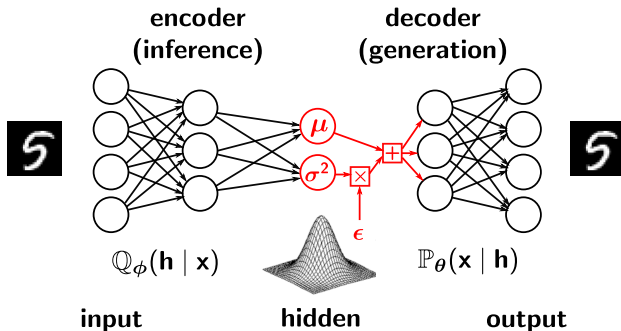


Variational Auto-Encoders (Kingma and Welling, 2013)

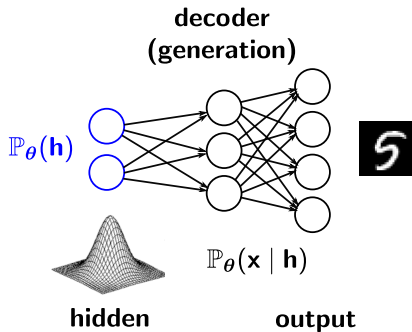


- Decoder computes $P_{\theta}(\mathbf{h})$ and $P_{\theta}(\mathbf{x} | \mathbf{h})$
- Encoder computes $Q_{\phi}(\mathbf{h} | \mathbf{x}) = \mathcal{N}(\mathbf{h}; \mu_{\phi}(\mathbf{x}), \sigma_{\phi}^2(\mathbf{x}))$
- Loss function: ELBO.

Summing Up: VAEs at Training Time



Summing Up: VAEs at Test Time



What is the Latent Variable Representing?

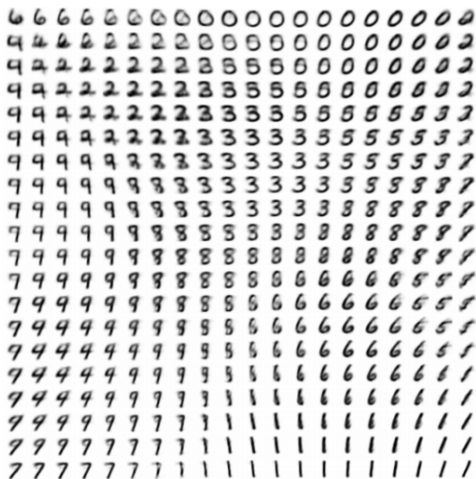
- Nice property of VAE: simultaneously training a **parametric encoder** in combination with a **generator network** forces the model to learn a predictable coordinate system that the encoder can capture.
- This makes it an excellent manifold learning algorithm.
- Example: the algorithm discovered two independent factors of variation present in images of faces: angle of rotation and emotional expression.

What is the Latent Variable Representing?



From Kingma and Welling (2013).

What is the Latent Variable Representing?



(b) Learned MNIST manifold

From Kingma and Welling (2013).

Issues with VAEs

- **Posterior collapse:** if the generative part is strong, the model may learn to ignore the latent variables:

$$\begin{aligned}\mathbb{P}_{\theta}(\mathbf{x} | \mathbf{h}) &\approx \mathbb{P}(\mathbf{x}) \\ \mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) &\approx \mathbb{P}_{\theta}(\mathbf{h}).\end{aligned}$$

- Can be mitigated with a few tricks:
 - Decrease/anneal the weight of $KL(\mathbb{Q}_{\phi}(\mathbf{h} | \mathbf{x}) || \mathbb{P}_{\theta}(\mathbf{h}))$ in the ELBO objective
 - Use auxiliary losses
 - Combine stochastic and amortized inference.
- In general, reporting both reconstruction loss and the KL term is needed to be able to tell if the model makes use of the latent variables.

Examples of Deep Generative Models

- Auto-Regressive Networks ✓
- Restricted Boltzmann Machines ✓
- Deep Belief Networks
- Deep Boltzmann Machines
- Gaussian-Bernoulli RBMs
- Convolutional Boltzmann Machines
- Sigmoid Belief Nets
- Variational Auto-Encoders ✓
- **Generative Adversarial Networks**
- Convolutional Generative Networks
- Generative Stochastic Networks

Why Maximum Likelihood?

- All models discussed aim to maximize the likelihood (evidence) $\mathbb{P}(\mathbf{x})$
- In fact, since this is intractable, they maximize a lower bound (ELBO)
- But if we want to build a generator, is this really the best criterion?
- Maximum likelihood tends to produce blurry images

Generative Adversarial Networks (GANs)

(Goodfellow et al., 2014)

- **Key idea:**

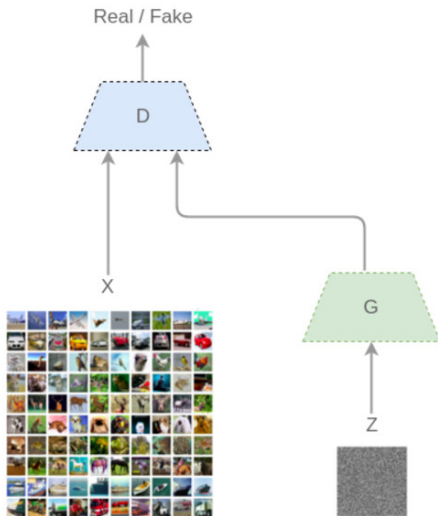
- keep the **generation network** $G = \{\mathbb{P}_\theta(\mathbf{h}), \mathbb{P}_\theta(\mathbf{x} | \mathbf{h})\}$
- drop the inference network and use instead a **discriminator network** $D : \mathcal{X} \rightarrow \{0, 1\}$

- Formulate the learning problem as a game between two players:

- the generator's job is to generate data that looks real
- the discriminator's job is to distinguish between real data and fake data generated by the generator

- This is like a **Turing test**: distinguish artificial from real.

Generative Adversarial Networks (GANs) (Goodfellow et al., 2014)



Minimax Game

- Saddle point problem:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbb{P}_{\text{data}}(\mathbf{x})}[\log D(\mathbf{x})] + \mathbb{E}_{\mathbb{P}_{\theta}(\mathbf{h})}[\log(1 - D(G(\mathbf{h})))].$$

- The optimal discriminator (intractable to compute) is:

$$D^*(\mathbf{x}) = \frac{\mathbb{P}_{\text{data}}(\mathbf{x})}{\mathbb{P}_{\text{data}}(\mathbf{x}) + \mathbb{P}_{\theta}(\mathbf{x})}, \quad \mathbb{P}_{\theta}(\mathbf{x}) = \int \mathbb{P}_{\theta}(\mathbf{x} | \mathbf{h})\mathbb{P}_{\theta}(\mathbf{h}).$$

- Given $D^*(\mathbf{x})$, the optimal generator $G^*(\mathbf{x})$ minimizes the **Jensen-Shannon divergence** between $\mathbb{P}_{\text{data}}(\mathbf{x})$ and $\mathbb{P}_{\theta}(\mathbf{x})$:

$$JS(\mathbb{P}_{\text{data}}(\mathbf{x}), \mathbb{P}_{\theta}(\mathbf{x})) = \frac{1}{2}KL(\mathbb{P}_{\text{data}}(\mathbf{x}) \| \bar{\mathbb{P}}(\mathbf{x})) + \frac{1}{2}KL(\mathbb{P}_{\theta}(\mathbf{x}) \| \bar{\mathbb{P}}(\mathbf{x})),$$

where $\bar{\mathbb{P}}(\mathbf{x}) = \frac{\mathbb{P}_{\text{data}}(\mathbf{x}) + \mathbb{P}_{\theta}(\mathbf{x})}{2}$.

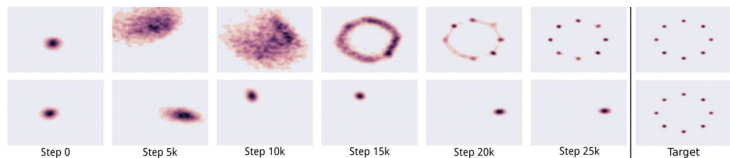
Training GANs

- How to train a GAN?
- Use stochastic gradient descent! Alternate between:
 - Stochastic gradients updates of the generator parameters θ
 - Stochastic gradients updates of the discriminator D .
- Several variants and schedules have been proposed.
- Caveats:
 - no convergence guarantees
 - optimization in GANs is often difficult

Mode Collapse

$$\min_G \max_D V(D, G) \neq \max_D \min_G V(D, G).$$

- G in inner loop: place all mass on most likely point



(From Metz et al. (2016))

- What prevents the generator from always picking the same example?
- The top row finds all the modes, the bottom finds just one mode.

Mode Collapse

- GANs often seem to collapse to far fewer modes than the model can represent
- This causes low output diversity.
- How to mitigate mode collapse?
- One strategy: minibatch (Salimans et al., 2016)
 - Let the discriminator make a decision by comparing an example to a whole minibatch of fake/real examples
 - Discriminator can now consider diversity.

Wasserstein GANs (WGANs)

- Instead of optimizing the Jensen-Shannon divergence, optimize instead:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{data}}(\mathbf{x})}[D(\mathbf{x})] - \mathbb{E}_{\mathbf{h} \sim \mathbb{P}_{\theta}(\mathbf{h})}[D(G(\mathbf{h}))].$$

- This is related to the **Wasserstein distance** (also called Earth mover's distance).
- A technical condition is that ∇D is bounded; in practice this is ensured with gradient clipping.
- This improves stability and mitigates the mode collapse problem.

Pros and Cons of GANs

Advantages:

- They currently generate the sharpest images
- They are cheap to train (since no statistical inference is required), and only back-propagation is needed to obtain gradients

Disadvantages:

- GANs are difficult to optimize due to unstable training dynamics.
- No statistical inference can be done with them.

Still Improving...



Ian Goodfellow

@goodfellow_ian

Follow



4 years of GAN progress (source:
eff.org/files/2018/02/ ...)



2014



2015



2016



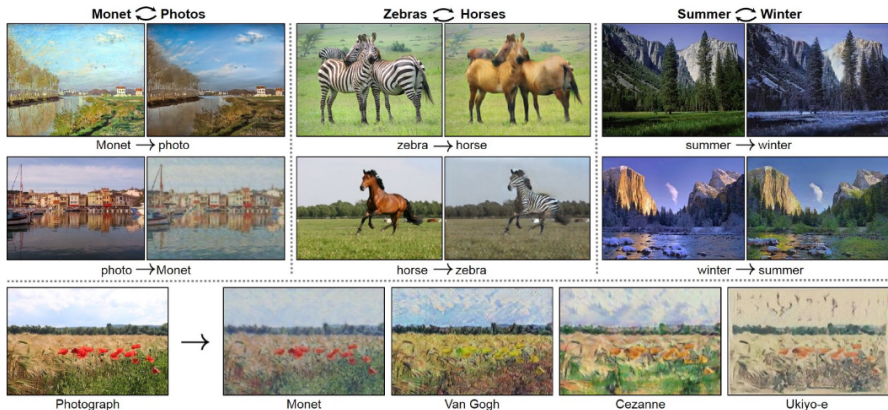
2017

7:26 PM - 2 Mar 2018

Some Extensions of GANs

- Augmenting GANs with an inference network (Dumoulin et al., 2016; Donahue et al., 2016)
- Domain adversarial training for domain adaptation (Ganin et al., 2016)
- Conditional GANs and semi-supervised GANs (Salimans et al., 2016)
- CycleGAN (Zhu et al., 2017): “translate” images from a source domain \mathcal{X} to a target domain \mathcal{Y} without paired examples. Use two generators $G : \mathcal{Y} \rightarrow \mathcal{X}$ and $F : \mathcal{X} \rightarrow \mathcal{Y}$ and introduce a **cycle consistency loss** to push $F(G(\mathbf{y})) \approx \mathbf{y}$ and $G(F(\mathbf{x})) \approx \mathbf{x}$.

Image-to-Image Translation w/ CycleGAN (Zhu et al., 2017)



(<https://junyanz.github.io/CycleGAN>)

Failure Cases



(<https://junyanz.github.io/CycleGAN>)

Evaluation

There is no single compelling way to evaluate a generative model.

- Models with **good** likelihood can produce **bad** samples
- Models with **good** samples can have **bad** likelihood
- There is no standard way to quantify how good samples are
- For GANs, it is also hard to even estimate the likelihood
- See “A note on the evaluation of generative models,” Theis et al. (2015), for a good overview.

Discrete Outputs

To train a GAN, G must be differentiable

But G cannot be differentiable if the output is discrete.

Possible workarounds:

- REINFORCE (Williams, 1992)
- Concrete/Gumbel-softmax distribution (Maddison et al., 2016; Jang et al., 2016)
- Learn distribution over continuous embeddings, decode to discrete

How does this compare with VAEs?

- VAEs have trouble with discrete latent variables (cannot differentiate through the inference network)
- GANs have trouble with discrete output variables (cannot differentiate through the generator network).

Connections to Reinforcement Learning

We can regard the discriminator loss as a reward signal for the generator.

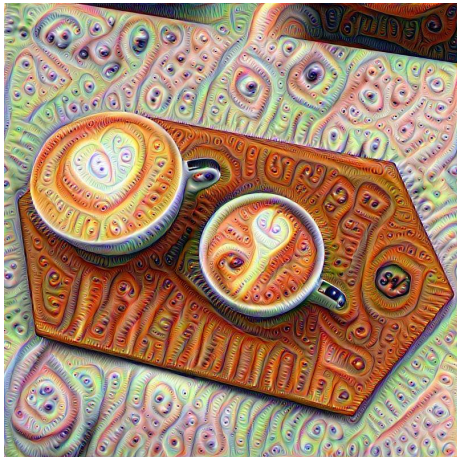
- GANs interpreted as actor-critic (Pfau and Vinyals, 2016)
- GANs as inverse reinforcement learning (Finn et al., 2016)

Conclusions

- Generative models are useful to model high-dimensional data
- Latent-variable generative models are appealing since they are more compact (“minimum description length” principle)
- Often, computing evidence and posterior distributions is intractable (e.g. Boltzmann machines)
- A common surrogate for maximum likelihood is the evidence lower bound (ELBO)
- Variational auto-encoders optimize the ELBO with amortized VI
- Their main drawback is posterior collapse
- Generative adversarial networks (GANs) are formulated as a game between a generator and a discriminator
- They manage to generate sharp outputs, but suffer from mode collapse and do not return a likelihood score
- Open problem (both VAEs/GANs): how to deal with discrete data?

Thank you!

Questions?



References I

- Ackley, D., Hinton, G., and Sejnowski, T. (1985). A learning algorithm for Boltzmann machines. *Cognitive science*, 9(1):147–169.
- Donahue, J., Krähenbühl, P., and Darrell, T. (2016). Adversarial feature learning. *arXiv preprint arXiv:1605.09782*.
- Dumoulin, V., Belghazi, I., Poole, B., Mastropietro, O., Lamb, A., Arjovsky, M., and Courville, A. (2016). Adversarially learned inference. *arXiv preprint arXiv:1606.00704*.
- Finn, C., Christiano, P., Abbeel, P., and Levine, S. (2016). A connection between generative adversarial networks, inverse reinforcement learning, and energy-based models. *arXiv preprint arXiv:1611.03852*.
- Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., and Lempitsky, V. (2016). Domain-adversarial training of neural networks. *The Journal of Machine Learning Research*, 17(1):2096–2030.
- Goodfellow, I., Bengio, Y., and Courville, A. (2016). Deep Learning. Book in preparation for MIT Press.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pages 2672–2680.
- Jang, E., Gu, S., and Poole, B. (2016). Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*.
- Kingma, D. P. and Welling, M. (2013). Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.
- Le Roux, N. and Bengio, Y. (2008). Representational power of restricted boltzmann machines and deep belief networks. *Neural computation*, 20(6):1631–1649.
- Maddison, C. J., Mnih, A., and Teh, Y. W. (2016). The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*.
- Metz, L., Poole, B., Pfau, D., and Sohl-Dickstein, J. (2016). Unrolled generative adversarial networks. *arXiv preprint arXiv:1611.02163*.
- Oord, A. v. d., Kalchbrenner, N., and Kavukcuoglu, K. (2016). Pixel Recurrent Neural Networks. In *Proc. of the International Conference on Machine Learning*.
- Pfau, D. and Vinyals, O. (2016). Connecting generative adversarial networks and actor-critic methods. *arXiv preprint arXiv:1610.01945*.
- Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., and Chen, X. (2016). Improved techniques for training gans. In *Advances in Neural Information Processing Systems*, pages 2234–2242.

References II

- Smolensky, P. (1986). Information processing in dynamical systems: Foundations of harmony theory. Technical report, DTIC Document.
- Theis, L., Oord, A. v. d., and Bethge, M. (2015). A note on the evaluation of generative models. *arXiv preprint arXiv:1511.01844*.
- Wainwright, M. and Jordan, M. (2008). *Graphical Models, Exponential Families, and Variational Inference*. Now Publishers.
- Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256.
- Zhu, J.-Y., Park, T., Isola, P., and Efros, A. A. (2017). Unpaired image-to-image translation using cycle-consistent adversarial networks. *arXiv preprint*.

Lecture 15: Responsible AI: Uncertainty, Explainability and Fairness

André Martins, Francisco Melo, Mário Figueiredo
Chrysoula Zerva



Deep Learning Course, Winter 2021-2022

Today's Roadmap

In the previous lectures we have seen different neural architectures, the potential of large neural models and their applications.

Today's Roadmap

In the previous lectures we have seen different neural architectures, the potential of large neural models and their applications.

Today we will discuss additional aspects that are important across different neural architectures, to enhance **trust** in models.

Today's Roadmap

In the previous lectures we have seen different neural architectures, the potential of large neural models and their applications.

Today we will discuss additional aspects that are important across different neural architectures, to enhance **trust** in models. We will focus on **Responsible AI** and touch on: model **uncertainty**, **explainability** and **fairness**.

Outline

- 1 Motivation
- 2 Uncertainty
- 3 Explainability
- 4 Fairness and Ethical AI
- 5 Conclusions

Quote of the day

“All models are wrong, but some can be useful”

George Box

Quote of the day

“All models are wrong, but some can be useful”

George Box

“... if we can trust them”

Who should trust the model?

Researchers & Engineers

- Improve understanding
- Anticipate issues
- Train better models
- Develop new algorithms
- Deploy safely

Consumers

- Understand purpose
- Increase trust
- Informed consent
- Correct & safe use
- Less bias

Regulators

- Understand impact
- Increase trust
- Adapt regulations
- Monitor and report

Outline

- ① Motivation
- ② Uncertainty
- ③ Explainability
- ④ Fairness and Ethical AI
- ⑤ Conclusions

Uncertainty in Machine Learning

- Typical setup:
 - Training input set \mathcal{X}^T
 - Learned outputs \mathcal{Y}^T
 - Neural model \mathcal{M} : $y = \mathcal{F}(x)$
 - Performance metric(s)



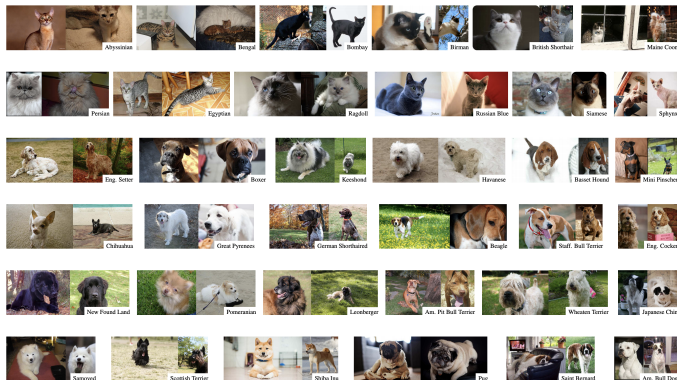
Uncertainty in Machine Learning

- Typical setup:
 - Training input set \mathcal{X}^T
 - Learned outputs \mathcal{Y}^T
 - Neural model \mathcal{M} : $y = \mathcal{F}(x)$
 - Performance metric(s)



- Assume:
 - Test sample $x^i, x^i \in \mathcal{X}^{\text{test}}$
 - How **trustworthy** y^i is?

Example in Deep Learning



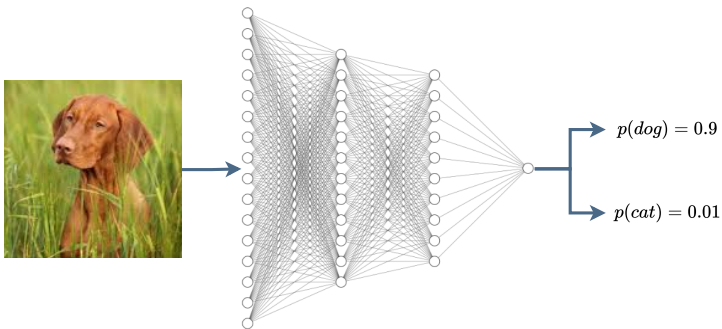
Oxford-IIIT-Pet-dataset Parkhi et al. (2012)

Example in Deep Learning: Image recognition

During test time:

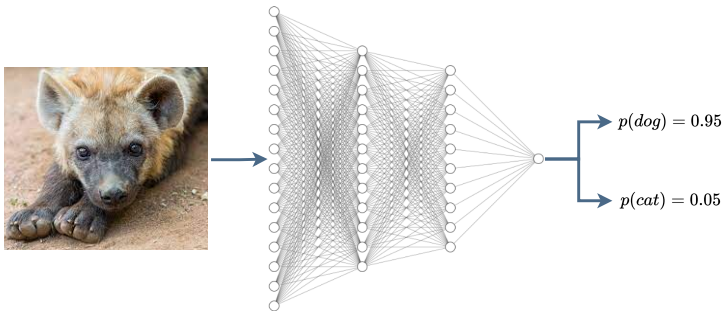
Example in Deep Learning: Image recognition

During test time:



Example in Deep Learning: Image recognition

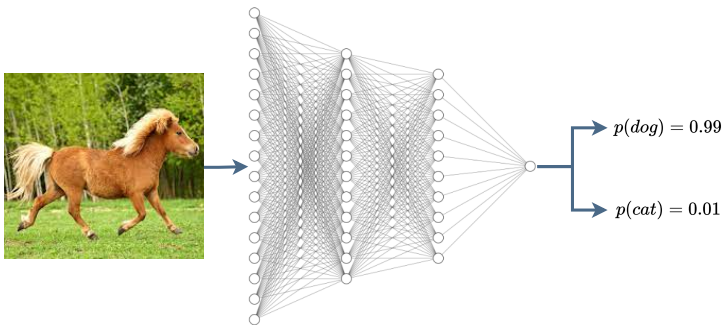
During test time:



- Is it expected for the model to be wrong out-of-domain ?

Example in Deep Learning: Image recognition

During test time:



- Is it expected for the model to be wrong out-of-domain ?
 - Could we have an indication that the model is **confused**?

Why is this important?

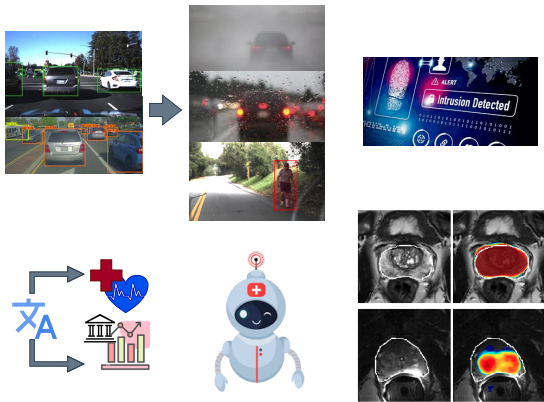
- Decide if a prediction must be further inspected

Why is this important?

- Decide if a prediction must be further inspected
- Useful for debugging, retraining, adapting a model

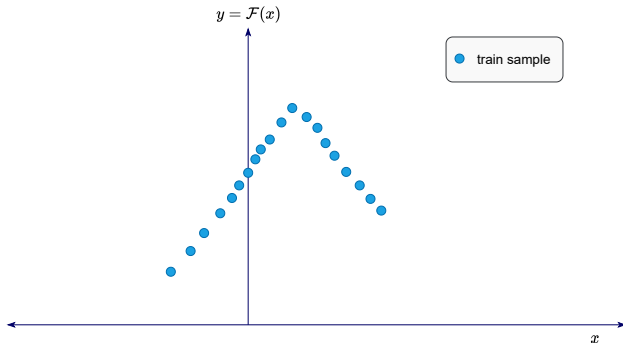
Why is this important?

- Decide if a prediction must be further inspected
- Useful for debugging, retraining, adapting a model
- Some applications are higher risk than others



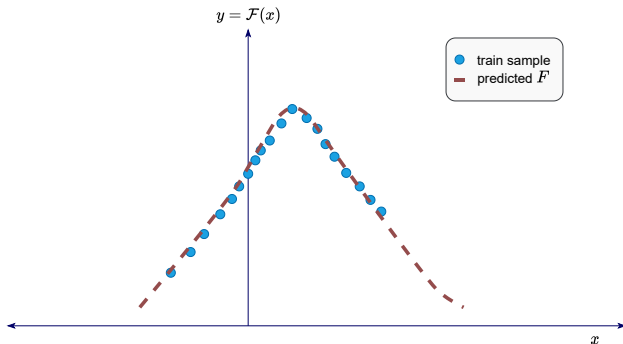
Uncertainty in Machine Learning

Let's see an example:



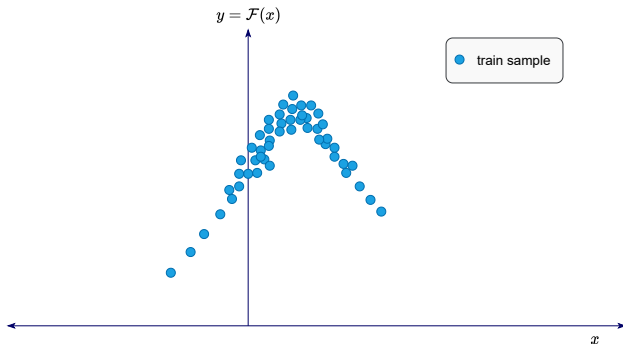
Uncertainty in Machine Learning

Let's see an example:



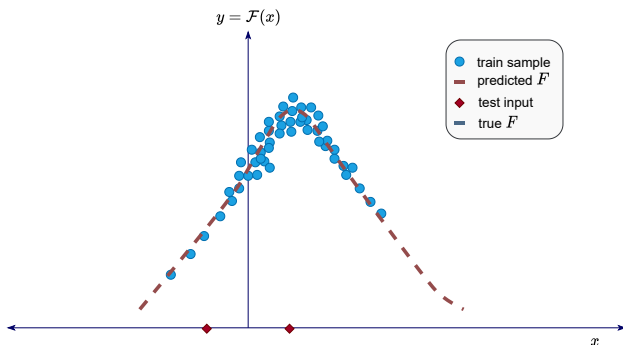
Uncertainty in Machine Learning

Let's see an example:



Uncertainty in Machine Learning

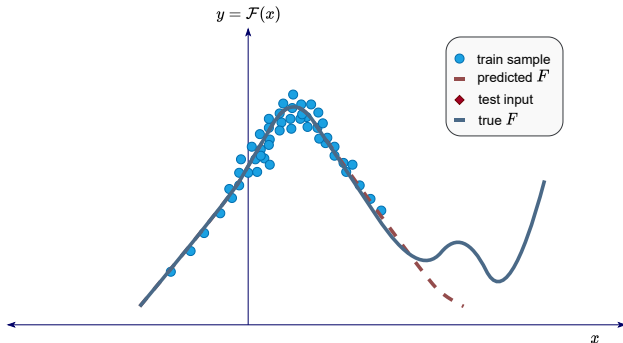
Let's see an example:



For which test point will we have a more **confident** prediction?

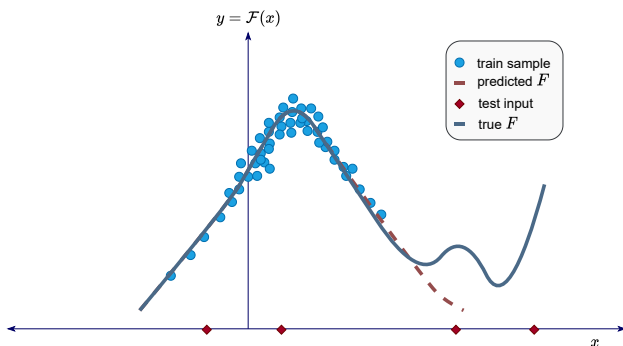
Uncertainty in Machine Learning

Let's see an example:



Uncertainty in Machine Learning

Let's see an example:



Uncertainty Quantification

- How do we estimate uncertainty?

Uncertainty Quantification

- How do we estimate uncertainty?
 - ▶ In information theory:

Uncertainty Quantification

- How do we estimate uncertainty?
 - ▶ In information theory: → **Entropy**

Uncertainty Quantification

- How do we estimate uncertainty?
 - ▶ In information theory: → **Entropy**
(for $\hat{y} = \mathcal{F}(x)$)

Uncertainty Quantification

- How do we estimate uncertainty?

▶ In information theory: → **Entropy**
(for $\hat{y} = \mathcal{F}(x)$)

- **Discrete predictions** (classification):
 $\mathcal{H}(p(\hat{y}|x)) = - \sum_{\hat{y} \in \mathcal{Y}} p(\hat{y}|x) \log p(\hat{y}|x)$

Uncertainty Quantification

- How do we estimate uncertainty?

► In information theory: → **Entropy**
(for $\hat{y} = \mathcal{F}(x)$)

- **Discrete predictions** (classification):

$$\mathcal{H}(p(\hat{y}|x)) = - \sum_{\hat{y} \in \mathcal{Y}} p(\hat{y}|x) \log p(\hat{y}|x)$$

- **Continuous predictions** (regression):

$$\mathcal{H}(p(\hat{y}|x)) = E[-\log p(\hat{y}|x)] = - \int_{y \in \mathcal{Y}} p(\hat{y}|x) \log p(\hat{y}|x) dx$$

Uncertainty Quantification

- How do we estimate uncertainty?

Uncertainty Quantification

- How do we estimate uncertainty?
 - ▶ Dispersion of a random variable in statistics?

Uncertainty Quantification

- How do we estimate uncertainty?
 - ▶ Dispersion of a random variable in statistics? → **Variance**

Uncertainty Quantification

- How do we estimate uncertainty?
 - ▶ Dispersion of a random variable in statistics? → **Variance**
 - We need to obtain a **distribution** of predictions $\hat{y} = \{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_N\}$ instead of a single point estimate \hat{y}

Uncertainty Quantification

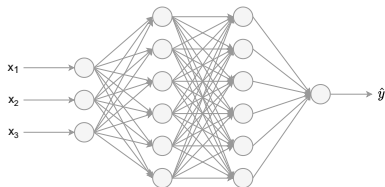
- How do we estimate uncertainty?
 - ▶ Dispersion of a random variable in statistics? → **Variance**
 - We need to obtain a **distribution** of predictions $\hat{y} = \{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_N\}$ instead of a single point estimate \hat{y}
 - Assume a Gaussian distribution: $p(y|\hat{y}) = \mathcal{N}(\hat{y}, \sigma^2)$

Uncertainty Quantification

- How do we estimate uncertainty?
 - ▶ Dispersion of a random variable in statistics? → **Variance**
 - We need to obtain a **distribution** of predictions $\hat{y} = \{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_N\}$ instead of a single point estimate \hat{y}
 - Assume a Gaussian distribution: $p(y|\hat{y}) = \mathcal{N}(\hat{y}, \sigma^2)$
 - ▶ How do we obtain such a distribution?

Variance-based Uncertainty Quantification

How can we obtain a distribution over \hat{y} ?



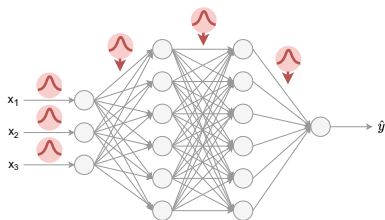
Variance-based Uncertainty Quantification

**Bayesian Neural Networks
(BNN)**

Variance-based Uncertainty Quantification

Bayesian Neural Networks (BNN)

Apply a prior distribution over model weights, e.g. a Gaussian:
 $W \sim \mathcal{N}(0, I)$.



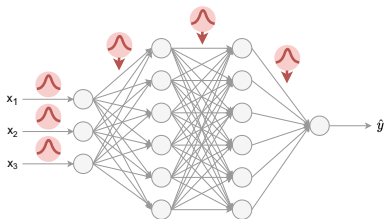
Variance-based Uncertainty Quantification

Bayesian Neural Networks (BNN)

Apply a prior distribution over model weights, e.g. a Gaussian: $W \sim \mathcal{N}(0, I)$.

We can then infer the model likelihood $p(y|\hat{y}) = p(y|f^W(x))$; and compute the posterior over the weights $p(W|X, Y)$.

Then we can define the likelihood as $p(y|f^W(x)) = \mathcal{N}(f^W(x), \sigma^2)$



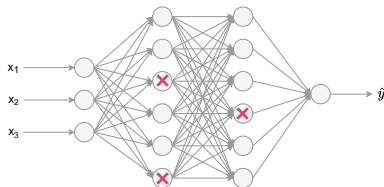
Variance-based Uncertainty Quantification

Monte Carlo Dropout (MCD):

Variance-based Uncertainty Quantification

Monte Carlo Dropout (MCD):

- ▶ Apply model dropout during inference.
- ▶ Run multiple stochastic forward runs.
- ▶ We can define the likelihood as $p(y|\hat{y}) \sim \mathcal{N}(\mu(\hat{y}), \sigma(\hat{y})^2)$



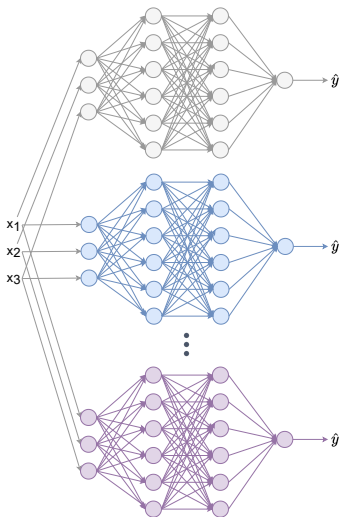
Variance-based Uncertainty Quantification

Deep Ensembles (DE):

Variance-based Uncertainty Quantification

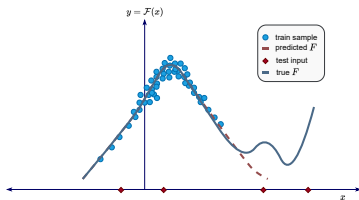
Deep Ensembles (DE):

- ▶ Obtain multiple model checkpoints
 - Different seeds;
 - Different hyper-parameters/training steps;
 - Bootstrap on data subsets;
- ▶ $p(y|\hat{y}) \sim \mathcal{N}(\mu(\hat{y}), \sigma(\hat{y})^2)$



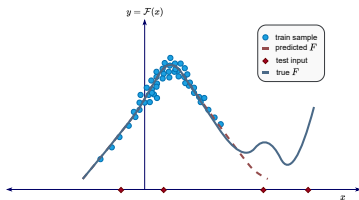
Uncertainty Quantification

How do we estimate uncertainty uncertainties?



Uncertainty Quantification

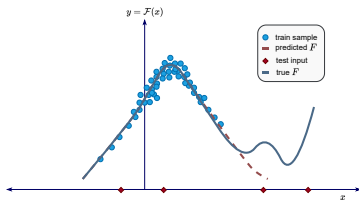
How do we estimate uncertainty uncertainties?



► What uncertainties do we need to estimate?

Uncertainty Quantification

How do we estimate uncertainties?



► What uncertainties do we need to estimate?

$$U_{\text{total}} = U_{\text{al}} + U_{\text{ep}}$$

total uncertainty aleatoric/data uncertainty epistemic/model/knowledge uncertainty

Aleatoric Uncertainty

- Aleatoric uncertainty derives from the inherent randomness of the data (noise)
- Also called data uncertainty
- Does not depend on model parameters
- Corresponds to the **irreducible** error → no matter how much more data we use to train, we cannot remove the uncertainty

Aleatoric Uncertainty



aleatoric : the latin word alea refers to a dice-like game

Aleatoric Uncertainty



aleatoric : the latin word alea refers to a dice-like game

Assume two models \mathcal{M}_1 and \mathcal{M}_2 that estimate the dice-rolling outcome $p(x|d)$, for a six-sided dice d .

Aleatoric Uncertainty



aleatoric : the latin word alea refers to a dice-like game

Assume two models \mathcal{M}_1 and \mathcal{M}_2 that estimate the dice-rolling outcome $p(x|d)$, for a six-sided dice d .

- \mathcal{M}_1 has seen a very large number of rolls (infinite)

Aleatoric Uncertainty



aleatoric : the latin word alea refers to a dice-like game

Assume two models \mathcal{M}_1 and \mathcal{M}_2 that estimate the dice-rolling outcome $p(x|d)$, for a six-sided dice d .

- \mathcal{M}_1 has seen a very large number of rolls (infinite)
What is the uncertainty for the predicted outcome $x = 5$?

Aleatoric Uncertainty



aleatoric : the latin word *alea* refers to a dice-like game

Assume two models \mathcal{M}_1 and \mathcal{M}_2 that estimate the dice-rolling outcome $p(x|d)$, for a six-sided dice d .

- \mathcal{M}_2 has seen only the following roll outcomes:
 $\{1, 2, 2, 3, 3, 3, 4, 5\}$

Aleatoric Uncertainty



aleatoric : the latin word alea refers to a dice-like game

Assume two models \mathcal{M}_1 and \mathcal{M}_2 that estimate the dice-rolling outcome $p(x|d)$, for a six-sided dice d .

- \mathcal{M}_2 has seen only the following roll outcomes:
 $\{1, 2, 2, 3, 3, 3, 4, 5\}$
What is the expected uncertainty for the predicted outcome $x = 5$?

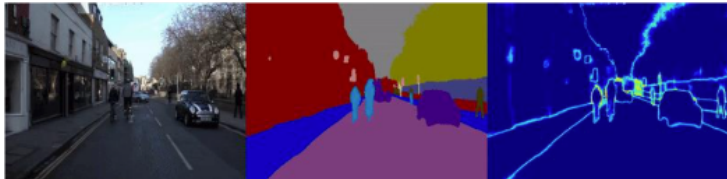
More on dice...

What if we apply the same models on these dice?



Aleatoric Uncertainty in Deep Learning

In object recognition:



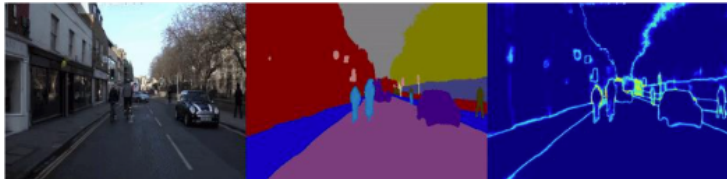
Input image

Segmented image

aleatoric noise

Aleatoric Uncertainty in Deep Learning

In object recognition:



Input image

Segmented image

aleatoric noise

In machine translation:

How did you find this place?



Como encontraste este sítio?

Que achaste deste sítio?

Como encontraste o caminho para este sítio?

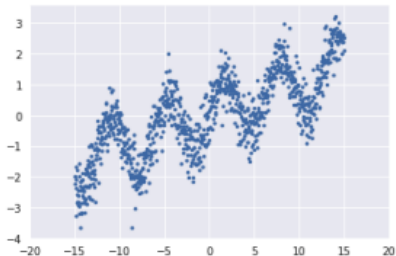
Aleatoric Uncertainty Estimation

What assumptions can we make about noise in the data?

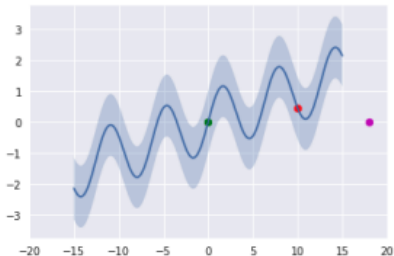
Aleatoric Uncertainty Estimation

What assumptions can we make about noise in the data?

- Homoscedastic noise: **constant** noise variance
 - Does not depend on input



(a) Dataset with homoscedastic noise

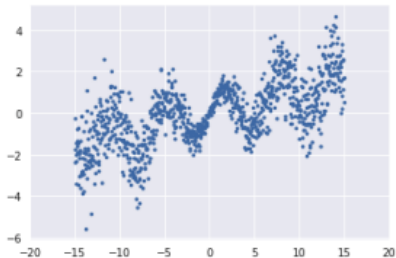


(b) Mean and homoscedastic variance

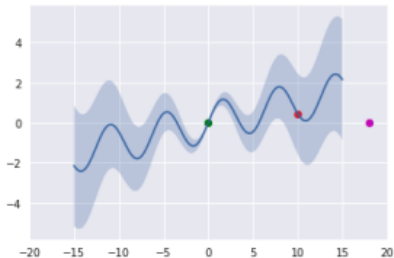
Aleatoric Uncertainty Estimation

What assumptions can we make about noise in the data?

- Heteroscedastic noise: **variable** noise variance
 - Does depend on input



(c) Dataset with heteroscedastic noise



(d) Mean and heteroscedastic variance

Estimating Aleatoric uncertainty

What can we use to estimate aleatoric uncertainty?

- If we have access to the dataset:

We can distill the uncertainty using the entropy-based approach and **conditioning on the model weights**

$$\begin{aligned}\mathcal{H}(p(\hat{y}|x)) &= - \sum_{\hat{y} \in \mathcal{Y}} p(\hat{y}|x) \log p(\hat{y}|x) \rightarrow \\ &\rightarrow \mathcal{H}(p(\hat{y}|x, w)) = - \sum_{\hat{y} \in \mathcal{Y}} p(\hat{y}|x, w) \log p(\hat{y}|x, w)\end{aligned}$$

Predicting Aleatoric uncertainty

What can we use to estimate aleatoric uncertainty?

- If we have access to **multiple human judgements** per segment:

Predicting Aleatoric uncertainty

What can we use to estimate aleatoric uncertainty?

- If we have access to **multiple human judgements** per segment:
 - Assume we can estimate a metric of dispersion of human scores per instance: e.g. assume the scores follow a Gaussian distribution $y \sim \mathcal{N}(\mu_h, \sigma_h^2)$
 - We would like to learn μ_h as the target quality score and σ_h as the target uncertainty score;
 - We want to learn to predict a Gaussian as close as possible to $\mathcal{N}(\mu_h, \sigma_h^2) \rightarrow$ minimise KL divergence $\text{KL}(p_{y_h} \| p_{\hat{y}})$
 - Learning objective: $\mathcal{L}_{\text{KL}} = \frac{(\mu_h - \hat{\mu})^2 + \sigma_h^2}{2\hat{\sigma}^2} + \frac{1}{2} \log \frac{\hat{\sigma}^2}{\sigma_h^2} - \frac{1}{2}$

Predicting Aleatoric uncertainty

What can we use to estimate aleatoric uncertainty?

- If we have access to **multiple human judgements** per segment:
 - Assume we can estimate a metric of dispersion of human scores per instance: e.g. assume the scores follow a Gaussian distribution $y \sim \mathcal{N}(\mu_h, \sigma_h^2)$
 - We would like to learn μ_h as the target quality score and σ_h as the target uncertainty score;
 - We want to learn to predict a Gaussian as close as possible to $\mathcal{N}(\mu_h, \sigma_h^2) \rightarrow$ minimise KL divergence $\text{KL}(p_{y_h} \| p_{\hat{y}})$
 - Learning objective: $\mathcal{L}_{\text{KL}} = \frac{(\mu_h - \hat{\mu})^2 + \sigma_h^2}{2\hat{\sigma}^2} + \frac{1}{2} \log \frac{\hat{\sigma}^2}{\sigma_h^2} - \frac{1}{2}$
- If we have access to **single human judgement** per segment?
 - Can we assume unit variance for the true labels?

Epistemic uncertainty

- Also known as: systematic, knowledge, model uncertainty
- Reflects **incomplete knowledge** of the model
 - Insufficient/complex training data
 - Out-of-domain examples
- Epistemic uncertainty corresponds to the **reducible error**
 - if we had more data (knowledge) the model would be more confident (less wrong).

Predicting Epistemic Uncertainty

Predicting Epistemic Uncertainty

Typically we try to estimate epistemic uncertainty using:

- Variance based methods
- Total uncertainty estimates $U_{\text{ep}} = U_{\text{total}} - U_{\text{al}}$

Predicting Epistemic Uncertainty

Typically we try to estimate epistemic uncertainty using:

- Variance based methods
 - Total uncertainty estimates $U_{\text{ep}} = U_{\text{total}} - U_{\text{al}}$
- Can we use error estimates to improve the estimation of epistemic uncertainty?

Predicting Epistemic Uncertainty

Typically we try to estimate epistemic uncertainty using:

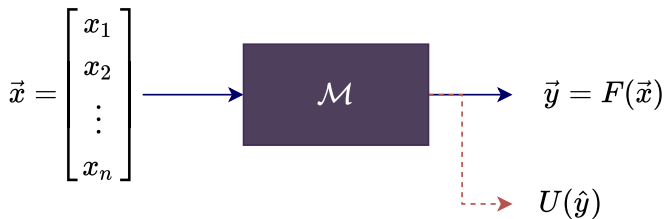
- Variance based methods
 - Total uncertainty estimates $U_{\text{ep}} = U_{\text{total}} - U_{\text{al}}$
- Can we use error estimates to improve the estimation of epistemic uncertainty?
- **Calibration:** Sometimes the magnitude of uncertainty needs to be adapted
 - the uncertainty should be **representative** of how often/much the model is wrong

Predicting Epistemic Uncertainty

Typically we try to estimate epistemic uncertainty using:

- Variance based methods
 - Total uncertainty estimates $U_{ep} = U_{total} - U_{al}$
- Can we use error estimates to improve the estimation of epistemic uncertainty?
- **Calibration:** Sometimes the magnitude of uncertainty needs to be adapted
 - the uncertainty should be **representative** of how often/much the model is wrong
 - **Direct Uncertainty Prediction:** What if we train another model to predict the error of the main predictor?

Uncertainty-aware ML

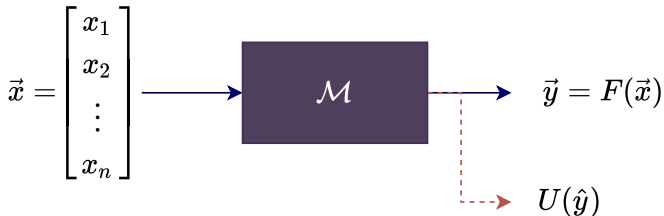


Outline

- ① Motivation
- ② Uncertainty
- ③ Explainability
- ④ Fairness and Ethical AI
- ⑤ Conclusions

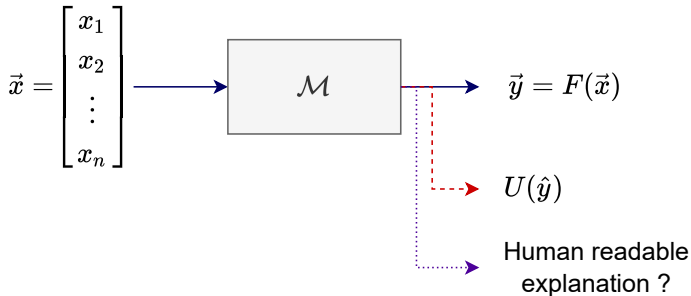
Motivation

Is uncertainty enough to understand **and trust** model behaviour?



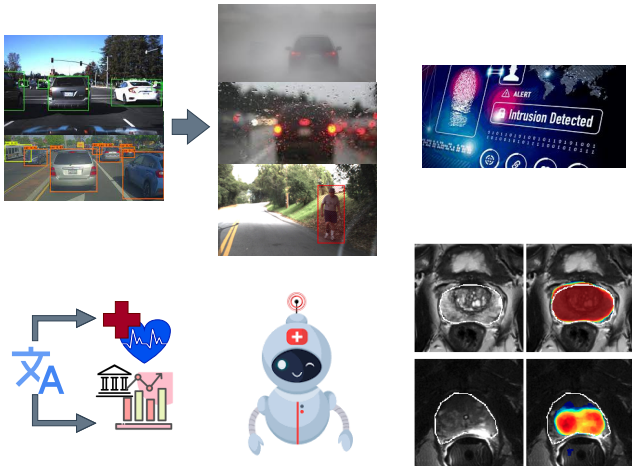
Motivation

Is uncertainty enough to understand **and trust** model behaviour?



Deep Learning Motivation

Remember the high risk deep learning applications:



Intrinsic Explainability

Some models are more explainable than others

Intrinsic Explainability

Some models are more explainable than others → Intrinsically explainable

Intrinsic Explainability

Some models are more explainable than others → Intrinsically explainable

- Decision trees
- Linear models (GAMs, GLMs)

$$G(E_Y(y|x)) = \beta_0 + f_1(x_1) + f_2(x_2) + \dots + f_n(x_n)$$

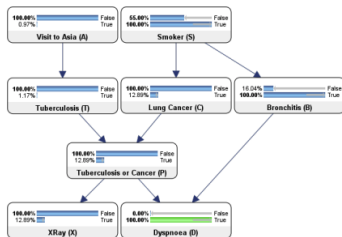


hints towards
feature significance

Intrinsic Explainability

Some models are more explainable than others → Intrinsically explainable

- Decision trees
- Linear models (GAMs, GLMs)
- Bayesian networks



Post-hoc Explainability

What about non intrinsically explainable models?

Post-hoc Explainability

What about non intrinsically explainable models?

Local explanations

Post-hoc Explainability

What about non intrinsically explainable models?

Local explanations

- **Model Specific** Explanations – Attribution-based methods
 - Gradient-based explanations
 - Integrated Gradients (Sundararajan et al., 2017)
 - SmoothGrad (Smilkov et al., 2017)
 - GradCAM (Selvaraju et al., 2016)
 - XRAI (Kapishnikov et al., 2019)
 - Layer-wise relevance propagation (Bach et al., 2015)
 - Attention-based explanations (Vashishth et al., 2019)
- **Model Generic** Explanations

Post-hoc Explainability

What about non intrinsically explainable models?

Local explanations

- **Model Specific** Explanations – Attribution-based methods
 - Gradient-based explanations
 - Integrated Gradients (Sundararajan et al., 2017)
 - SmoothGrad (Smilkov et al., 2017)
 - GradCAM (Selvaraju et al., 2016)
 - XRAI (Kapishnikov et al., 2019)
 - Layer-wise relevance propagation (Bach et al., 2015)
 - Attention-based explanations (Vashishth et al., 2019)
- **Model Generic** Explanations
 - Shapley values (SHAP) (Lundberg and Lee, 2017)
 - Input perturbation based explanations (LIME) (Ribeiro et al., 2016)
 - Counterfactual explanations (MiCE) (Ross et al., 2020)

Gradient-based explanations

Gradient-based explanations

Intuition: Reveal which regions of the input are important for the final prediction

Gradient-based explanations

Intuition: Reveal which regions of the input are important for the final prediction → Use the **gradients** as an indication of feature importance?

Gradient-based explanations

Intuition: Reveal which regions of the input are important for the final prediction → Use the **gradients** as an indication of feature importance?
Assume an image classification model with a class activation function S_c for each class $c \in C$, such that $class(x)_i = \underset{c \in C}{\operatorname{argmax}}(S_c(x))$

Gradient-based explanations

Intuition: Reveal which regions of the input are important for the final prediction → Use the **gradients** as an indication of feature importance? Assume an image classification model with a class activation function S_c for each class $c \in C$, such that $class(x)_i = \underset{c \in C}{\operatorname{argmax}}(S_c(x))$

- Compute the gradient of S_c with respect to x
- Average gradients from all layers
- Smooth gradients
- Compute integrated gradients

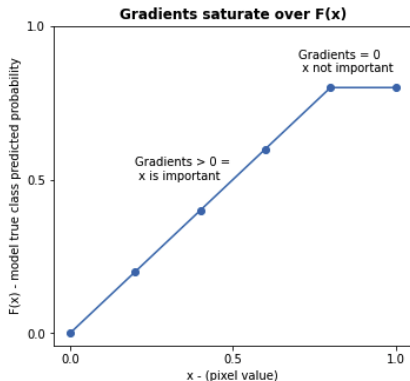
Integrated Gradients

Integrated Gradients

Intuition:

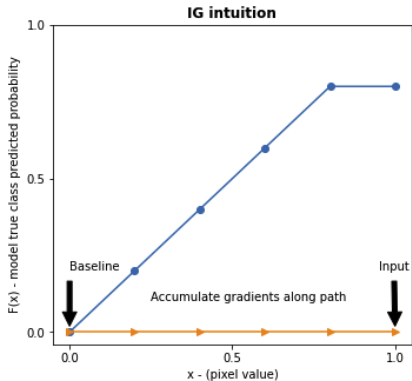
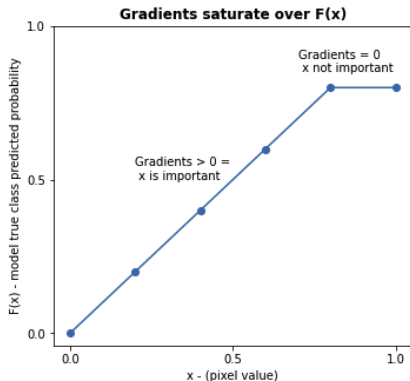
Integrated Gradients

Intuition:



Integrated Gradients

Intuition:



*Credits: *Doug Kelly talk*

Integrated Gradients

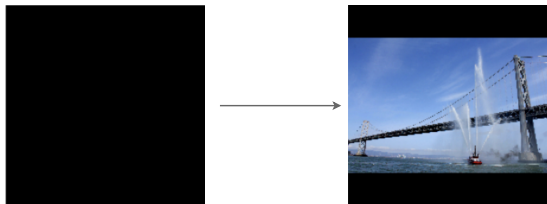
Step-by-step application:



Initial image x

Integrated Gradients

Step-by-step application:



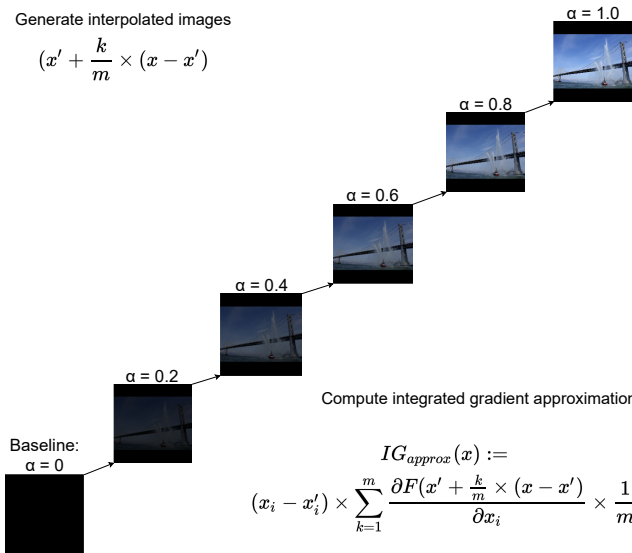
Start from baseline $\alpha = 0$

Interpolate images for small intervals from $\alpha = 0.0$ to $\alpha = 1.0$

Integrated Gradients

Generate interpolated images

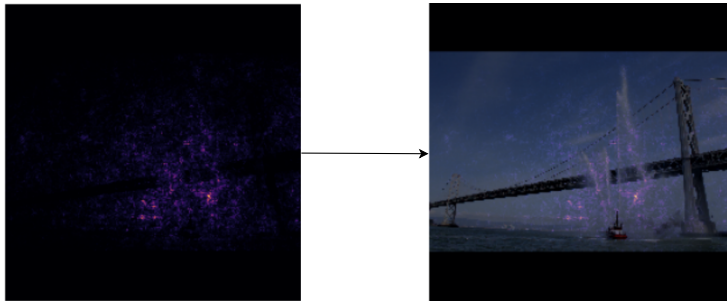
$$(x' + \frac{k}{m} \times (x - x'))$$



$$IG_{approx}(x) := (x_i - x'_i) \times \sum_{k=1}^m \frac{\partial F(x' + \frac{k}{m} \times (x - x'))}{\partial x_i} \times \frac{1}{m}$$

Integrated Gradients

Output visualisation & overlay on initial image:



Attribution in NLP tasks

Can we do this with text?

Attribution in NLP tasks

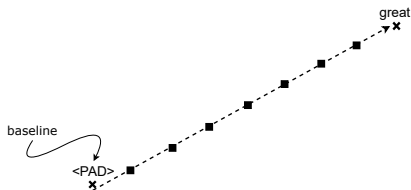
Can we do this with text?

- Gradient-based attribution on the word embedding space

Attribution in NLP tasks

Can we do this with text?

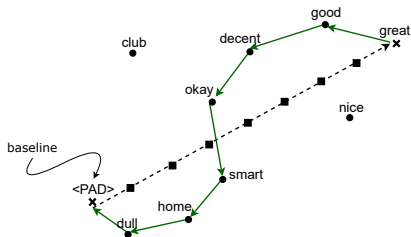
- Gradient-based attribution on the word embedding space
- Discretized integrated gradients (Sanyal and Ren, 2021)



Attribution in NLP tasks

Can we do this with text?

- Gradient-based attribution on the word embedding space
- Discretized integrated gradients (Sanyal and Ren, 2021)



Attention based explanations

Can attention explain model decisions?

Attention based explanations

Can attention explain model decisions?

Example: sentiment classification

Attention based explanations

Can attention explain model decisions?

Example: sentiment classification

fantastic movie one of the best film noir movies ever made
a meandering inarticulate and ultimately disappointing film

Attention based explanations

Can attention explain model decisions?

Example: sentiment classification

fantastic movie one of the best film noir movies ever made
a meandering inarticulate and ultimately disappointing film

Natural interpretation: how much each word is weighted when computing the label?

Is attention explanation?

Is attention explanation?

Is Attention Interpretable?

Sofia Serrano* **Noah A. Smith*[†]**

*Paul G. Allen School of Computer Science & Engineering,
University of Washington, Seattle, WA, USA

[†]Allen Institute for Artificial Intelligence, Seattle, WA, USA
{sofias6, nasmith}@cs.washington.edu

The elephant in the interpretability room: Why use attention as explanation when we have saliency methods?

Jasmijn Bastings
Google Research
bastings@google.com

Katja Filippova
Google Research
katjaf@google.com

Attention is not Explanation

Sarthak Jain
Northeastern University
jain.sar@husky.neu.edu

Byron C. Wallace
Northeastern University
b.wallace@northeastern.edu

Attention is not not Explanation

Sarah Wiegrefe*
School of Interactive Computing
Georgia Institute of Technology
saw@gatech.edu

Yuval Pinter*
School of Interactive Computing
Georgia Institute of Technology
uvp@gatech.edu

Learning to Deceive with Attention-Based Explanations

Danish Pruthi[‡], Mansi Gupta[‡], Bhuwan Dhingra[‡], Graham Neubig[‡], Zachary C. Lipton[‡]

[‡]Carnegie Mellon University, Pittsburgh, USA

[‡]Twitter, New York, USA

ddanish@cs.cmu.edu, mansig@twitter.com,
{bdhingra, gneubig, zlipton}@cs.cmu.edu

Attention based explanations

Can attention explain model decisions?

Example: sentiment classification

fantastic movie one of the best film noir movies ever made
a meandering inarticulate and ultimately disappointing film

Natural interpretation: how much each word is weighted when computing the label?

Some concerns:

- Low overlap with gradient methods
- Perturbation of weights doesn't change the predicted label
- Very large weights do not always correspond to meaningful explanations

Post-hoc Explainability

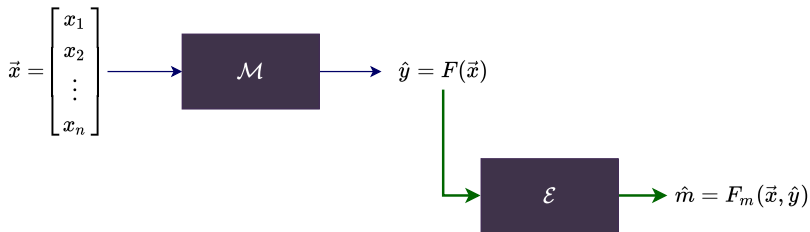
What about non intrinsically explainable models?

Local explanations:

- **Model Specific Explanations – Attribution-based methods**
 - Gradient-based explanations
 - Integrated Gradients
 - SmoothGrad
 - GradCAM
 - XRAI
 - Layer-wise relevance propagation
 - Attention-based explanations (?)
- **Model Generic Explanations**
 - Shapley values (SHAP)
 - **Perturbation-based explanations (LIME)**
 - **Counterfactual explanations (MiCE)**

Model Generic Explanations

- ▶ Can be applied to any model



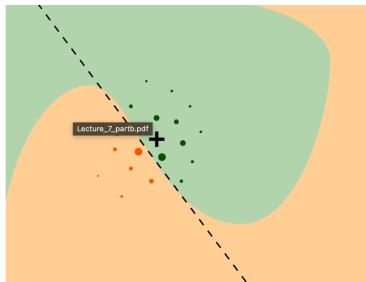
Intuition: Which parts of the input influence the prediction? → Which parts of the input would change the predicted label if altered?

LIME: Local Interpretable Model-agnostic Explanations

Assume \mathcal{M} predicts label \hat{y} for input x

- Select regions of x – super-features
 - Neighboring groups of pixels
 - Words / ngrams
- Generate neighbors of x in the feature space
 - Noise over pixels
 - MASKS over words
- Fit a glass-box explainer, \mathcal{E} that will have the same predictions as \mathcal{M}
 - ▶ A simple linear model would work

✓ \mathcal{E} will explain the behaviour of \mathcal{M} in the neighborhood of x



MiCE: Minimal Contrastive Editing

Assume \mathcal{M} predicts label \hat{y} for input x

MiCE: Minimal Contrastive Editing

Assume \mathcal{M} predicts label \hat{y} for input x

- ▶ What needs to be edited for the prediction to change?

MiCE: Minimal Contrastive Editing

Assume \mathcal{M} predicts label \hat{y} for input x

► What needs to be edited for the prediction to change?

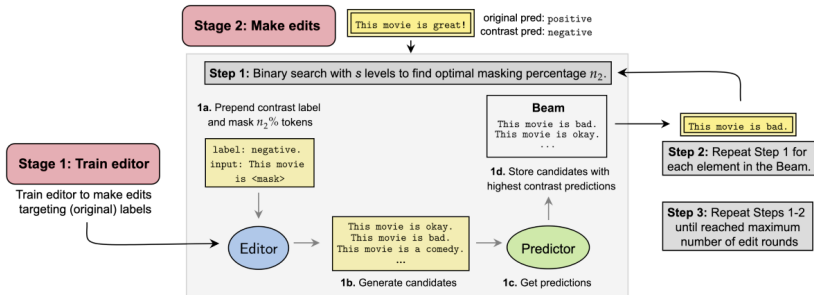
Train an editor \mathcal{E} to generate edits over x , such that \mathcal{M} predicts \hat{y}'

MiCE: Minimal Contrastive Editing

Assume \mathcal{M} predicts label \hat{y} for input x

► What needs to be edited for the prediction to change?

Train an editor \mathcal{E} to generate edits over x , such that \mathcal{M} predicts \hat{y}'

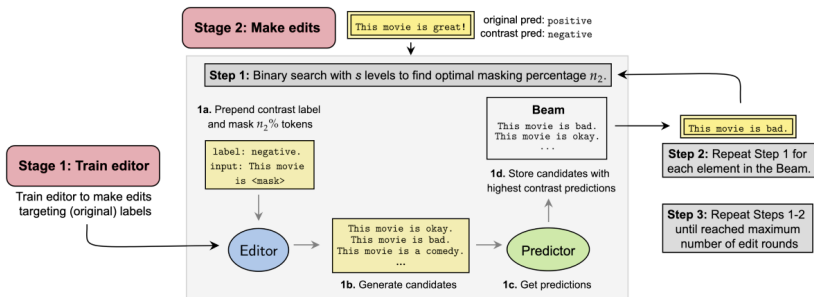


MiCE: Minimal Contrastive Editing

Assume \mathcal{M} predicts label \hat{y} for input x

► What needs to be edited for the prediction to change?

Train an editor \mathcal{E} to generate edits over x , such that \mathcal{M} predicts \hat{y}'



✓ \mathcal{E} will explain the behaviour of \mathcal{M} using contrastive examples

Which is the best explanation?

What is a **good** explanation?

Which is the best explanation?

What is a **good** explanation?

- Who is the explanation for?
 - Model developer / researcher
 - Consumer
 - Policy maker / regulator

Which is the best explanation?

What is a **good** explanation?

- Who is the explanation for?
 - Model developer / researcher
 - Consumer
 - Policy maker / regulator
- What is the target task?
 - What is the modality?
Text, Image, Sound, ...
 - Is it regression or classification ?

Which is the best explanation?

What is a **good** explanation?

- Who is the explanation for?
 - Model developer / researcher
 - Consumer
 - Policy maker / regulator
- What is the target task?
 - What is the modality?
Text, Image, Sound, ...
 - Is it regression or
classification ?

Some ground rules do apply:

Which is the best explanation?

What is a **good** explanation?

- Who is the explanation for?
 - Model developer / researcher
 - Consumer
 - Policy maker / regulator
- What is the target task?
 - What is the modality?
Text, Image, Sound, ...
 - Is it regression or classification ?

Some ground rules do apply:
Explanations should ensure:

- Completeness
- Accuracy
- Meaningfulness
- Consistency

Outline

- ① Motivation
- ② Uncertainty
- ③ Explainability
- ④ Fairness and Ethical AI
- ⑤ Conclusions

What is a fair/ethical model?

Assume a model:

- can achieve high performance
- can provide a confidence interval
- can provide an interpretable explanation

Is that sufficient?

► Can a model be (un)ethical?

Ethical AI

Can a model be (un)ethical?

Ethical AI

Can a model be (un)ethical? → Based on what ethics?

Ethical AI

Can a model be (un)ethical? → Based on what ethics?



Ethical AI

What ethical concerns could be relevant?

Ethical AI

What ethical concerns could be relevant?

- Discrimination bias
- Accessibility
- Privacy compromise
- Sustainability / energy consumption

Definitions of Bias

Definitions of Bias

- ▶ What is bias in statistics?

Definitions of Bias

► What is bias in statistics?

Statistical bias is a systematic tendency which causes differences between results and facts [Wikipedia]

Definitions of Bias

- ▶ What is bias in statistics?

Statistical bias is a systematic tendency which causes differences between results and facts [Wikipedia]

- ▶ What is bias in ML?

Definitions of Bias

- ▶ What is bias in statistics?

Statistical bias is a systematic tendency which causes differences between results and facts [Wikipedia]

- ▶ What is bias in ML?

Inductive bias (or learning bias) of an ML algorithm is the set of assumptions that the learner uses to predict outputs of given inputs that it has not encountered.

Typology of Bias

Typology of Bias

Sample Bias: Also Selection / Representation bias. The training dataset does not contain fair/balanced/sufficient representation of instances in the testing environment.

Confounding bias: The distortion of the association between the independent and dependent variables because a third variable is independently associated with both.

Association bias: Bias occurring from implicit associations in the training data: features might co-occur or correlate in the training data, without this association holding in the testing environment.

Observer bias: Also Confirmation bias. The effect of seeing what we expect to see in data. Subjective opinions influencing the data processing (labelling, feature manipulation, sampling)

Exclusion bias: Deleting valuable data or features thought to be unimportant/redundant.

Discrimination Bias: Also demographic / racial bias. The training dataset is data skewed in favour of particular demographics.

Bias in practice

“Easier to teach a machine to see than to **understand** what it has seen”

Kosinski, 2018

Bias in practice

“Easier to teach a machine to see than to **understand** what it has seen”

Kosinski, 2018

Biased ML models is not a recent concern; ML has its own urban legends

Bias in practice

“Easier to teach a machine to see than to **understand** what it has seen”

Kosinski, 2018

Biased ML models is not a recent concern; ML has its own urban legends

► Failed DARPA experiment on distinguishing US vs Russian tanks

The neural net learned to distinguish:



Bias in practice

“Easier to teach a machine to see than to **understand** what it has seen”

Kosinski, 2018

Biased ML models is not a recent concern; ML has its own urban legends

► Failed DARPA experiment on distinguishing US vs Russian tanks

The neural net learned to distinguish:

- high vs low resolution
- forest vs open-ground terrain
- cloudy vs clear sky

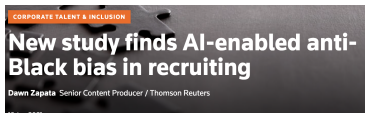


Socially concerning ML bias

The hiring algorithm case:

Amazon scraps secret AI recruiting tool that showed bias against women

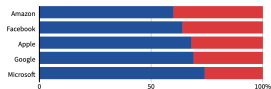
Amazon Reportedly Killed an AI Recruitment System Because It Couldn't Stop the Tool from Discriminating Against Women



Where does the bias come from?

GLOBAL HEADCOUNT

■ Male ■ Female



Overcoming AI's Challenge In Hiring: Avoid Human Bias

Insights Team Insights Contributor
FORBES INSIGHTS With Intel AI | Paid Program
Innovation

f v in

NOVEMBER 1, 2018

Amazon's sexist hiring algorithm could still be better than a human

by Maude Lavanchy, The Conversation

Socially concerning ML bias

The criminal justice case:

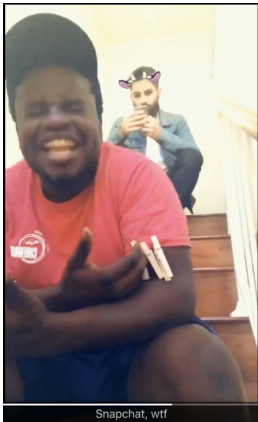


Socially concerning ML bias

The vision case:

Socially concerning ML bias

The vision case:



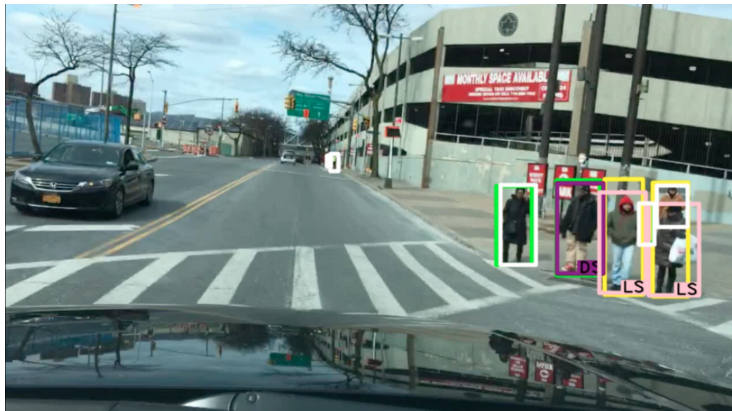
No snapchat filters on darker skin



No soap dispenser activation on darker skin

Socially concerning ML bias

The vision case:



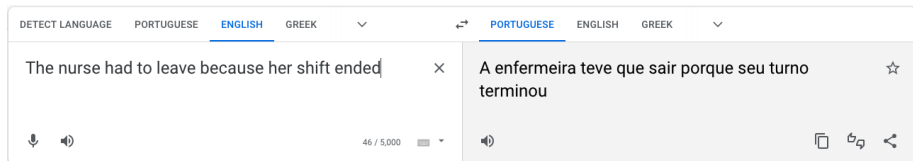
Fail to recognise human objects when they have darker skin

Socially concerning ML bias

The machine translation case:

Socially concerning ML bias

The machine translation case:



The screenshot shows a machine translation interface with two panels. The left panel is for the source text, and the right panel is for the translated text. The source text is "The nurse had to leave because her shift ended". The translated text is "A enfermeira teve que sair porque seu turno terminou". The interface includes language selection menus at the top, a microphone icon, a speaker icon, and a progress indicator "46 / 5,000".

Source Text (English)	Translated Text (Portuguese)
The nurse had to leave because her shift ended	A enfermeira teve que sair porque seu turno terminou

Socially concerning ML bias

The machine translation case:

The image shows two screenshots of the Google Translate interface. In the top screenshot, the source text is "The nurse had to leave because her shift ended" (English) and the target text is "A enfermeira teve que sair porque seu turno terminou" (Portuguese). In the bottom screenshot, the source text is "The nurse had to leave because his shift ended" (English) and the target text is "A enfermeira teve que sair porque seu turno terminou" (Portuguese). The interface includes language selection menus at the top, a microphone icon, a character count "46 / 5,000", and a star icon for saving the translation.

Socially concerning ML bias

The machine translation case:

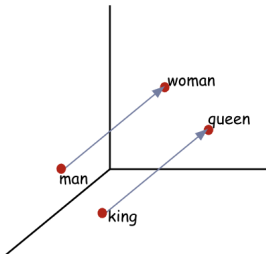
The image displays two screenshots of a machine translation interface, likely Google Translate, illustrating a gender bias in the output.

Top Screenshot: The source text is "The nurse had to leave because her shift ended". The target text is "A enfermeira teve que sair porque seu turno terminou".

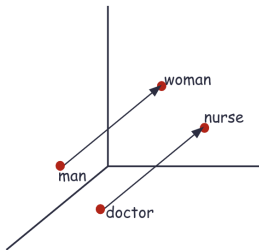
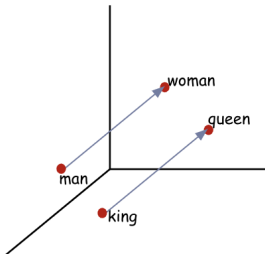
Bottom Screenshot: The source text is "The nurse had to leave because his shift ended". The target text is "A enfermeira teve que sair porque seu turno terminou". The phrase "A enfermeira" is circled in red, highlighting the gender bias where the female pronoun "A" is used for a male subject.

Both screenshots show the interface with language selection (ENGLISH to PORTUGUESE), a character count (46 / 5,000), and various utility icons (microphone, speaker, copy, share, star).

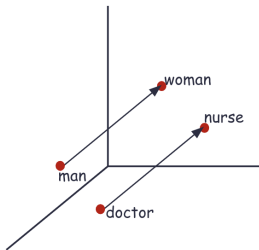
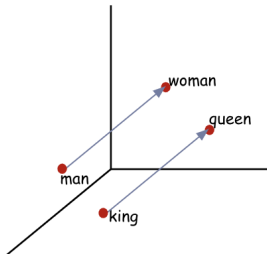
Biased Language models



Biased Language models

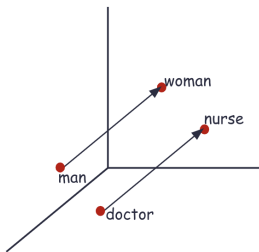
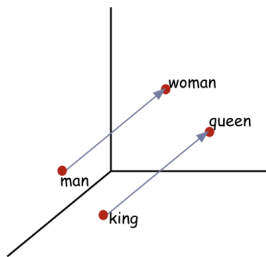


Biased Language models



sewing-carpentry
nurse-surgeon
registered nurse-physician
interior designer-architect
vocalist-guitarist
diva-superstar
housewife-shopkeeper

Biased Language models



sewing-carpentry
nurse-surgeon
registered nurse-physician
interior designer-architect
vocalist-guitarist
diva-superstar
housewife-shopkeeper

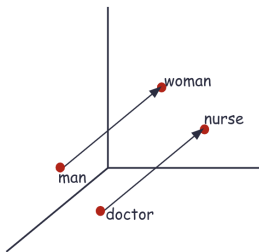
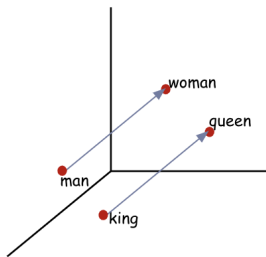
Appropriate analogies:

- ✓ queen-king
- ✓ waitress-waiter
- ✓ sister-brother
- ✓ mother-father
- ✓ ovarian cancer-prostate cancer
- ✓ convent-monastery

Spurious analogies:

- lovely-brilliant
- feminism-conservatism
- blond-burly
- charming-affable
- giggle-chuckle
- cupcakes-pizzas

Biased Language models



sewing-carpentry
nurse-surgeon
registered nurse-physician
interior designer-architect
vocalist-guitarist
diva-superstar
housewife-shopkeeper

Appropriate analogies:

- ✓ queen-king
- ✓ waitress-waiter
- ✓ sister-brother
- ✓ mother-father
- ✓ ovarian cancer-prostate cancer
- ✓ convent-monastery

Spurious analogies:

- lovely-brilliant
- feminism-conservatism
- blond-burly
- charming-affable
- giggle-chuckle
- cupcakes-pizzas

extreme "she"  **extreme "he"**

homemaker	maestro
nurse	skipper
receptionist	protege
librarian	philosopher
socialite	captain
hairdresser	architect
nanny	financier
bookkeeper	warrior
stylist	broadcaster
housekeeper	magician

Implicit Association Test

Group 1

woman

Group 2

man

Implicit Association Test

Group 1

woman

Group 2

man

QUEEN

Implicit Association Test

Group 1

woman

Group 2

man

BROTHER

Implicit Association Test

Group 1

woman

OR

career

Group 2

man

OR

family

Implicit Association Test

Group 1

woman

OR

career

Group 2

man

OR

family

supermodel

Implicit Association Test

Group 1

woman

OR

career

Group 2

man

OR

family

CEO

Implicit Association Test

Group 1

woman

OR

career

Group 2

man

OR

family

Wedding planning

Implicit Association Test

Group 1

woman

OR

career

Group 2

man

OR

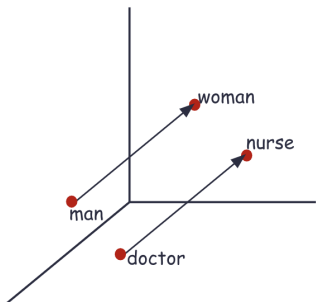
family

Wedding planning

► Delays in response time over several judgements can reveal association bias

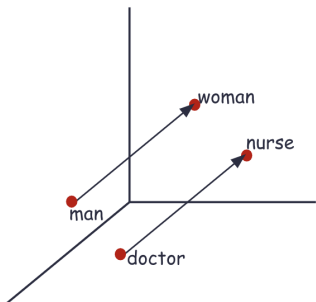
Biased Language models

- ▶ What is an appropriate testing for language models (word embeddings)?



Biased Language models

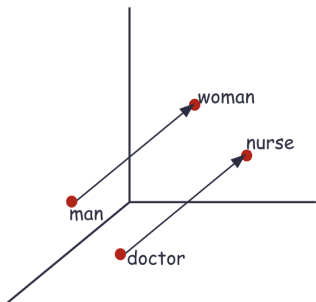
► What is an appropriate testing for language models (word embeddings)?



- Assume two groups of words we want to test:
 - $G_1 = \{doctor, CEO, lawyer\}$
 - $G_2 = \{nurse, homemaker, supermodel\}$
 - Seed groups:
 - $S_1 = \{woman, female, \dots\}$ -
 - $S_2 = \{man, male, \dots\}$

Biased Language models

- ▶ What is an appropriate testing for language models (word embeddings)?



- Assume two groups of words we want to test:
 - $G_1 = \{doctor, CEO, lawyer\}$
 - $G_2 = \{nurse, homemaker, supermodel\}$
 - Seed groups:
 - $S_1 = \{woman, female, \dots\}$ -
 - $S_2 = \{man, male, \dots\}$
- What do we expect the distance to the seed groups to be?
WEAT dataset (Caliskan et al., 2017)

How do we correct bias?

► Pre-training:

How do we correct bias?

► Pre-training:

- Use diverse datasets
- Counterfactual data augmentation (CDA)
- Use social context during annotation

How do we correct bias?

► Pre-training:

- Use diverse datasets
- Counterfactual data augmentation (CDA)
- Use social context during annotation

Assume we can identify and neutralise a gender dimension $g \in R$

$$w = w_{\perp} + w_g$$

How do we correct bias?

► Pre-training:

- Use diverse datasets
- Counterfactual data augmentation (CDA)
- Use social context during annotation

Assume we can identify and neutralise a gender dimension $g \in R$

$$w = w_{\perp} + w_g$$

► Training:

- Loss terms that restrict the gender information at the end of the vector (Zhao et al., 2018)
- Regularization term that penalizes the projection on the gender direction (Bordia and Bowman, 2019)

How do we correct bias?

► Pre-training:

- Use diverse datasets
- Counterfactual data augmentation (CDA)
- Use social context during annotation

Assume we can identify and neutralise a gender dimension $g \in R$

$$w = w_{\perp} + w_g$$

► Training:

- Loss terms that restrict the gender information at the end of the vector (Zhao et al., 2018)
- Regularization term that penalizes the projection on the gender direction (Bordia and Bowman, 2019)

► Post-training:

- Push biased representations “away” from the gender dimension

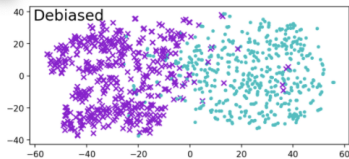
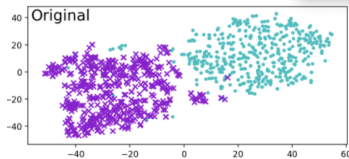
Is it working?

- Debiasing during training / post-processing typically demands available seed-words

Is it working?

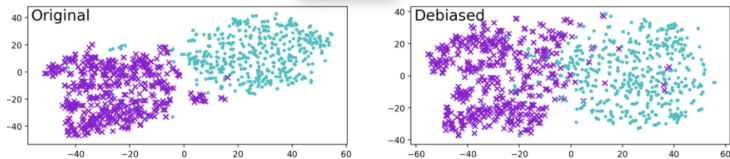
- Sometimes the bias is just hidden a bit better (Gonen and Goldberg,

2019)



Is it working?

- Sometimes the bias is just hidden a bit better (Gonen and Goldberg,



2019)

- Some other options:
 - Context-aware models: machine translation, coreference resolution
 - Gender-tagging: What if we tell the model the correct gender?

Bias is not the only issue: Accessibility

Who has access to ML applications?

Bias is not the only issue: Accessibility

Who has access to ML applications?
Who are ML applications designed for?

Bias is not the only issue: Accessibility

Who has access to ML applications?
Who are ML applications designed for?

► Let's look at NLP applications:

Bias is not the only issue: Accessibility

Who has access to ML applications?

Who are ML applications designed for?

► Let's look at NLP applications:

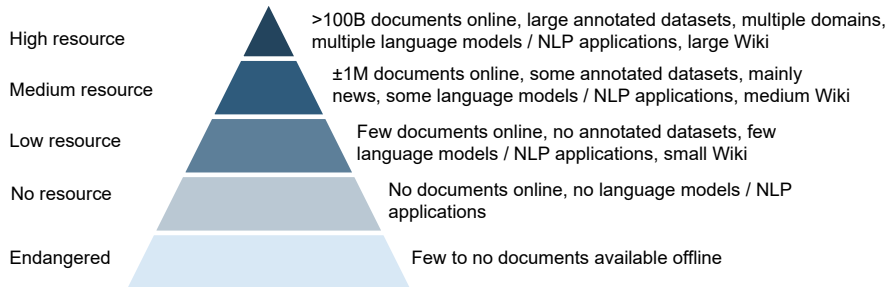
- Are NLP models covering all languages?
- Are multilingual models equally efficient for all covered languages?

Bias is not the only issue: Accessibility

Who has access to ML applications?
Who are ML applications designed for?

► Let's look at NLP applications:

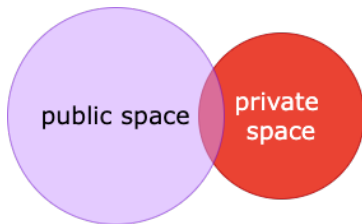
- Are NLP models covering all languages?
- Are multilingual models equally efficient for all covered languages?



Bias is not the only issue: Privacy

Privacy in ML: When is it a concern?

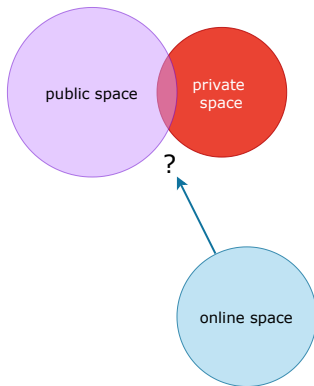
- Should we be tracked?



Bias is not the only issue: Privacy

Privacy in ML: When is it a concern?

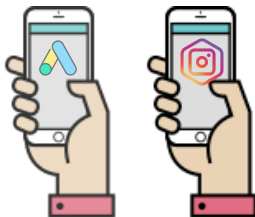
- Should we be tracked?



Bias is not the only issue: Privacy

Privacy in ML: When is it a concern?

- Should we be tracked?
- Should we be tracked if we consent?



Bias is not the only issue: Privacy

Privacy in ML: When is it a concern?

- Should we be tracked?
- Should we be tracked if we consent?
- Is it informed consent? Do we know how the data will be used?



Cambridge
Analytica

Bias is not the only issue: Privacy

Privacy in ML: When is it a concern?

- Should we be tracked?
- Should we be tracked if we consent?
- Is it informed consent? Do we know how the data will be used?
- What if we change our mind?



Cambridge
Analytica

Profiling, Privacy and Anonymisation

► **Sensitive Personal Information (SPI):** information that can be used on its own or with other information to identify, contact, or locate a single person, or to identify an individual in context [wikipedia]

- Full name (if not common)
- Home address
- Email address
- National identification number
- Passport number
- IP address
- Vehicle registration plate number
- Driver's license number
- Face, fingerprints, or handwriting
- Credit card numbers
- Date of birth
- Birthplace
- Genetic information
- Telephone number
- Login name, screen name, nickname, or handle

► **Solution:** Anonymise data before using it to develop ML models

Is anonymisation enough?

SPI can still be uncovered:

- Combination of databases
- Use of non-sensitive data
- Can SPI be learned?
- Is user generated content revealing SPI?
- Is everyone equally sensitive?

Bias is not the only issue: Sustainability

AI does not come for free:

- Human resources
- Materials (hardware)
- Power consumption
- ▶ CO2 emissions:
 - Training BERT \approx trans-american jet flight
 - 100h of GPU computation: 7-30 kg CO2
 - Hardware
 - Cloud provider
 - Location (cheaper in the west)

Sustainable AI

Can AI be sustainable?

Sustainable AI

Can AI be sustainable?

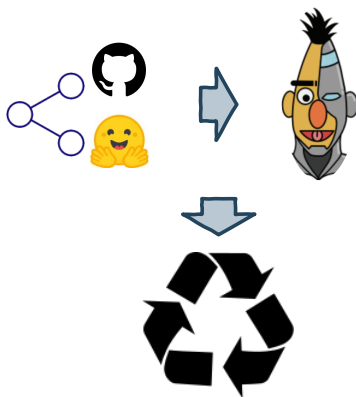
- Consider sustainability at every step of the AI life cycle
 - Keep humans in the loop



Sustainable AI

Can AI be sustainable?

- Consider sustainability at every step of the AI life cycle
 - Keep humans in the loop
- Sustainable AI development
 - Are large models always necessary?
 - Distillation of large models to smaller student models for inference?
 - Share – Adapt – Reuse



Sustainable AI

Can AI be sustainable?

- Consider sustainability at every step of the AI life cycle
 - Keep humans in the loop
- Sustainable AI development
 - Are large models always necessary?
 - Distillation of large models to smaller student models for inference?
 - Share – Adapt – Reuse
- Many providers offset their carbon emissions



Sustainable AI

Can AI be sustainable?

- Consider sustainability at every step of the AI life cycle
 - Keep humans in the loop
- Sustainable AI development
 - Are large models always necessary?
 - Distillation of large models to smaller student models for inference?
 - Share – Adapt – Reuse
- Many providers offset their carbon emissions
- AI solutions can provide sustainable solutions for several domains



Outline

- ① Motivation
- ② Uncertainty
- ③ Explainability
- ④ Fairness and Ethical AI
- ⑤ Conclusions

Pointers to topics not discussed

► Uncertainty:

- Active learning using uncertainty criteria
- Calibration
- Direct uncertainty prediction

► Explanability

- Global explanations
- Shapley values
- Robustness to adversarial attacks

► Ethical AI

- Concerns on intellectual property
- Conversational agents

Conclusions

~~Trust the process~~

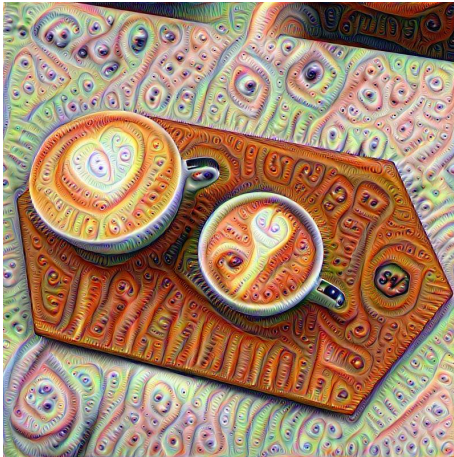
Trustworthy model development can be an involved **process**

Things that can help:

- Confidence-aware models
 - Quantify data uncertainty
 - Quantify model uncertainty
- Explainable models
 - Local explanations
 - Global explanations
- Ethical models
 - Unbiased
 - Accessible
 - Sustainable

Thank you!

Questions?



References I

- Bach, S., Binder, A., Montavon, G., Klauschen, F., Müller, K.-R., and Samek, W. (2015). On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS one*, 10(7):e0130140.
- Bordia, S. and Bowman, S. R. (2019). Identifying and reducing gender bias in word-level language models. *arXiv preprint arXiv:1904.03035*.
- Caliskan, A., Bryson, J. J., and Narayanan, A. (2017). Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.
- Gonen, H. and Goldberg, Y. (2019). Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them. *arXiv preprint arXiv:1903.03862*.
- Kapishnikov, A., Bolukbasi, T., Viégas, F., and Terry, M. (2019). Xrai: Better attributions through regions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4948–4957.
- Lundberg, S. M. and Lee, S.-I. (2017). A unified approach to interpreting model predictions. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Parkhi, O. M., Vedaldi, A., Zisserman, A., and Jawahar, C. (2012). Cats and dogs. In *2012 IEEE conference on computer vision and pattern recognition*, pages 3498–3505. IEEE.
- Ribeiro, M. T., Singh, S., and Guestrin, C. (2016). "why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144.
- Ross, A., Marasović, A., and Peters, M. E. (2020). Explaining nlp models via minimal contrastive editing (mice). *arXiv preprint arXiv:2012.13985*.
- Sanyal, S. and Ren, X. (2021). Discretized integrated gradients for explaining language models. *arXiv preprint arXiv:2108.13654*.
- Selvaraju, R. R., Das, A., Vedantam, R., Cogswell, M., Parikh, D., and Batra, D. (2016). Grad-cam: Why did you say that? *arXiv preprint arXiv:1611.07450*.
- Smilkov, D., Thorat, N., Kim, B., Viégas, F., and Wattenberg, M. (2017). Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*.
- Sundararajan, M., Taly, A., and Yan, Q. (2017). Axiomatic attribution for deep networks. In *International conference on machine learning*, pages 3319–3328. PMLR.
- Vashishth, S., Upadhyay, S., Tomar, G. S., and Faruqui, M. (2019). Attention interpretability across nlp tasks. *arXiv preprint arXiv:1909.11218*.
- Zhao, J., Wang, T., Yatskar, M., Ordonez, V., and Chang, K.-W. (2018). Gender bias in coreference resolution: Evaluation and debiasing methods. *arXiv preprint arXiv:1804.06876*.