



**TÉCNICO**  
LISBOA

# **Multimodal and Longitudinal Approaches to Alzheimer's Disease Classification**

**Beatriz Oliveira Ferreira**

Thesis to obtain the Master of Science Degree in

## **Electrical and Computer Engineering**

Supervisor(s): Prof. Maria Margarida Campos da Silveira  
Prof. Jaime dos Santos Cardoso

### **Examination Committee**

Chairperson: Prof. : João Fernando Cardoso Silva Sequeira

Supervisor: Prof. Maria Margarida Campos da Silveira

Member of the Committee: Prof. Ana Luísa Nobre Fred

**16 November 2020**



Declaro que o presente documento é um trabalho original da minha autoria e que cumpre todos os requisitos do Código de Conduta e Boas Práticas da Universidade de Lisboa.

I declare that this document is an original work of my own authorship and that it fulfills all the requirements of the Code of Conduct and Good Practices of the Universidade de Lisboa.



## Acknowledgments

Gostaria de começar por agradecer aos meus orientadores, a professora Margarida Silveira e o professor Jaime Cardoso, pelo acompanhamento, ideias, sugestões e paciência. Foram formidáveis e fundamentais neste fim de ciclo.

Tenho também de agradecer ao Instituto Superior Técnico. Passei aqui os anos mais duros e desafiadores da minha formação, mas também os mais recompensadores. Aqui cresci muito, tanto a nível académico como a nível pessoal.

Neste percurso conheci pessoas que com certeza serão companheiros para a vida: quero agradecer ao João, à Mariana e ao Duarte. Pelos dias passados a passear e pelos dias passados a trabalhar. Pelos momentos de descontração e por me perguntarem "Como está a tese?". E agradecer também de modo muito especial ao Guilherme. Sabe sempre o que dizer, para além de dar muito bons abraços.

Por fim, tenho de agradecer à minha família o apoio, amor e carinho: mãe, pai, Teté, Luís Pedro e Sofia.

Obrigada a todos.



## Resumo

A doença de Alzheimer é uma doença neurodegenerativa, para a qual não se conhece a cura. No entanto, o seu progresso pode ser atrasado, quando é precocemente diagnosticada. MRI e PET são ferramentas fulcrais para este propósito, uma vez que podem ajudar a identificar a doença anos antes de aparecerem sintomas. Este trabalho analisou como é que vários classificadores se comportavam em diferentes situações. Estas incluíram classificação binária e multiclasse (até 4 classes), prever o diagnóstico futuro, prever o diagnóstico atual com base num histórico, e prever se o diagnóstico de um sujeito iria piorar. Este último cenário teve uma dificuldade acrescida, já que a representatividade de uma das classes era muito superior à da outra. As técnicas utilizadas para lidar com este efeito foram SMOTE, "random undersampling", e uma combinação das duas. Os classificadores usados foram SVM, RF, NN e RNN. Em ocasiões onde foram considerados vários intervalos de tempo, os dados de entrada tiveram de ser adaptados para SVM e RF. Para isto as "features" de todos os intervalos de tempo foram concatenadas. Para além dos classificadores já conhecidos, o kSVM foi proposto e testado. O objectivo deste classificador era melhorar o resultado do SVM para situações em que há vários intervalos de tempo, ao receber as "features" de cada intervalo de tempo separadas, em vez de concatenadas. No entanto a sua performance foi pior que a do SVM, o que significa que este classificador não é muito promissor. Foi também proposto um classificador de mudanças de diagnóstico.

**Palavras-chave:** AD, MCI, MRI, RF, RNN, SVM.





## Abstract

Alzheimer's disease is a degenerative brain disease, without a known cure. However, its progress can be delayed, when diagnosed in early stages. MRI and PET imaging are precious tools for this purpose, as they can help to identify this disease years before symptoms appear.

The present work analyzed how several classifiers performed in different scenarios. These scenarios included binary and multiclass classification (up to 4 classes), predicting the future diagnosis, predicting the current diagnosis having access to data from multiple previous timesteps and predicting whether a subject's diagnosis would get worst. In this last scenario there was an extra obstacle, since the data was highly imbalanced. SMOTE, random undersampling, and a combination of both were the approaches used to deal with the imbalance.

The classifiers used were SVM, RF, NN and RNN. In scenarios where more than one timestep was being considered, the input data was adapted to SVM and RF, by concatenating the features of all timesteps. In addition to the standard classifiers, kSVM was proposed and tested. This classifier's objective was to improve the SVM score when there were several timesteps, by breaking up the features by timestep. However, its performance proved to be worse than the one of SVM, which makes this classifier not very promising. A classifier to predict changes of diagnosis was also proposed.

**Keywords:** AD, MCI, MRI, RF, RNN, SVM.



# Contents

- Acknowledgments . . . . . v
- Resumo . . . . . vii
- Abstract . . . . . ix
- List of Tables . . . . . xiii
- List of Figures . . . . . xv
- Acronyms . . . . . xix
  
- 1 Introduction . . . . . 1**
- 1.1 Motivation . . . . . 1
- 1.1.1 Alzheimer’s Disease . . . . . 1
- 1.1.2 Diagnosis . . . . . 2
- 1.2 Objectives . . . . . 2
- 1.3 Thesis Outline . . . . . 3
  
- 2 Theoretical Background . . . . . 5**
- 2.1 Classification methods . . . . . 5
- 2.1.1 Support Vector Machines . . . . . 5
- 2.1.2 Random forests . . . . . 8
- 2.1.3 Neural networks . . . . . 9
- 2.1.4 kSVM . . . . . 13
- 2.1.5 Training . . . . . 13
- 2.2 Data properties . . . . . 14
- 2.2.1 Ordinal data . . . . . 14
- 2.2.2 Imbalanced data . . . . . 16
- 2.3 Performance metrics . . . . . 18
- 2.3.1 Confusion matrix . . . . . 18
- 2.3.2 Accuracy . . . . . 18
- 2.3.3 F1-score . . . . . 19
  
- 3 State-of-the-Art . . . . . 21**

<b>4 Experiments and Results</b>	<b>25</b>
4.1 Data description . . . . .	25
4.2 Experimental design . . . . .	28
4.2.1 Data preparation . . . . .	28
4.2.2 Performance score . . . . .	28
4.2.3 SVM and kSVM . . . . .	28
4.2.4 RF . . . . .	28
4.2.5 NN . . . . .	29
4.3 Predictions for a timestep, given data from that timestep . . . . .	29
4.4 Predictions for a timestep, given data from the previous timestep . . . . .	31
4.5 Predictions for a timestep, given data from more than one previous timestep . . . . .	32
4.6 Predicting changes of diagnosis . . . . .	33
<b>5 Conclusions and Future Work</b>	<b>37</b>
5.1 Conclusions . . . . .	37
5.2 Future Work . . . . .	38
<b>Bibliography</b>	<b>41</b>

# List of Tables

3.1	Summary of the performance of different proposed methods. Metrics considered in multiclass problems: MAE - Mean Absolute Error. Metrics considered in binary problems: BACC - Balanced Accuracy. SENS - Sensivity. SPEC - Specificity. AUC - Area Under Curve. . . . .	23
4.1	Number of participants in each ADNI phase, according to their clinical status. . . . .	25
4.2	Average f1-score for classifying subjects into CN or AD, with data relative to month 0. . . . .	29
4.3	Average f1-score for classifying subjects into CN, MCI or AD, with data relative to month 0. . . . .	30
4.4	Average f1-score for classifying subjects into CN, EMCI, LMCI or AD, with data relative to month 0. . . . .	30
4.5	Average f1-score for predicting the subjects' class on month 24, with baseline information from month 24, 12, 6 or 0. . . . .	31
4.6	Average f1-score for predicting the subjects' class on month 12, with information from month 12, 6 or 0. . . . .	32
4.7	Average f1-score obtained with different methods for classifying subjects into 3 classes, in function of the number of timesteps. . . . .	33
4.8	Average f1-score obtained with different strategies to deal with class imbalance and different classifiers, for classifying subjects as having changed their diagnosis between month 0 and month 24. . . . .	34



# List of Figures

1.1	Age of people with AD in the US. Image from 2018 Alzheimer’s disease: facts & figures [1].	1
2.1	Simple neural network structure. Image from Artificial Neural Networks Tutorial [27].	9
2.2	Simple RNN structure. Image from Deep Learning with Python [21].	12
2.3	Simple RNN structure unfolded in time. Image from Deep Learning with Python [21].	13
2.4	Structure of the classifier composed by 3SVM. The three sets of numbers represent hypothetical outputs for three input patterns.	13
2.5	3-fold cross validation. Image from Deep Learning with Python [21].	14
2.6	Example of the creation of a new sample using the SMOTE algorithm. Image from Learning from Imbalanced Data [36].	17
2.7	Structure of the confusion matrix for n classes. Cells shown in blue correspond to correct classifications whereas cells in white correspond to incorrect ones.	18
4.1	MR images of the same section of the brain, using different weightings. Image from Introduction to the Basics of Magnetic Resonance Imaging [46].	26
4.2	Number of subjects by visit.	26
4.3	Data available in the ADNI database, broken down by diagnosis and ADNI phases, and by timesteps. Figures from [43].	27
4.4	Confusion matrix obtained for classification of subjects in 3 classes, using MRI and PET imaging. The classifier used was a neural network.	30
4.5	Confusion matrix obtained for classification of subjects in 4 classes, using MRI and PET imaging. The classifier used was a SVM.	31
4.6	Average f1-score obtained with different methods for classifying subjects into three classes, using information from months 0, 12, 24.	33
4.7	Comparison of average f1-score obtained with different methods for classifying subjects into three classes, with different number of time steps.	34
4.8	Confusion matrices obtained for classifying subjects into change/no change, using SVM. The time window considered was from month 0 to month 24.	35
4.9	Confusion matrices obtained for classifying subjects into change/no change, using RF. The time window considered was from month 0 to month 24.	35

4.10 Confusion matrices obtained for classifying subjects into change/no change, using NN.

The time window considered was from month 0 to month 24. . . . . 36





# Acronyms

**AD** Alzheimer's disease.

**CSF** Cerebrospinal fluid.

**DAGSVM** Directed acyclic graph support vector machine.

**EMCI** Early mild cognitive impairment.

**GRU** Gated recurrent unit.

**LMCI** Late mild cognitive impairment.

**LSTM** Long short term memory.

**PCA** Principal component analysis.

**RBF** Radial basis function.

**CAD** Computer aided diagnosis.

**CART** Classification and regression trees.

**CM** Confusion matrix.

**CN** Cognitive normal.

**CNN** Convolutional neural network.

**ID3** Iterative Dichotomiser 3.

**MCI** Mild cognitive impairment.

**MLP** Multilayer perceptron.

**MRI** Magnetic resonance imaging.

**NN** Neural network.

**PET** Positron emission tomography.

**RELU** Rectified linear unit.

**RF** Random forest.

**RNN** Recurrent neural network.

**RU** Random undersampling.

**SDPSO** Switching delayed particle swarm optimization.

**SMOTE** Synthetic minority oversampling technique.

**SVM** Support vector machine.



# Chapter 1

## Introduction

### 1.1 Motivation

#### 1.1.1 Alzheimer’s Disease

Alzheimer’s disease (AD) is a degenerative brain disease, and the most common type of dementia. In the early stages it causes damage to brain cells in parts of the brain involved in cognitive functions, which translates in difficulties with memory, language, problem-solving and other cognitive skills. In later stages, neurons in other parts of the brain are also damaged or destroyed, affecting a person’s ability to perform basic body functions like walking or swallowing, which means that people in these stages are bed-bound and need around-the-clock-care. Alzheimer’s disease is ultimately fatal [1].

Alzheimer’s disease prevalence increases with age, as shown in figure 1.1. Due to the fact that the population in developed countries is aging, the worldwide prevalence of AD is increasing: in 2010 there were 35.6 million people living with the disease, while in 2015 there were 46.8 million, which represents an increase of 31%. Therefore, the costs associated with the disease have increased 35% from US \$604 bn in 2010 to \$818 bn in 2015. These costs can be divided into three main categories: direct medical costs, social care, and informal care. Their distribution has not changed substantially: in 2015 their relative weight was, respectively, 19.5%, 40.1% and 40.4% [2].

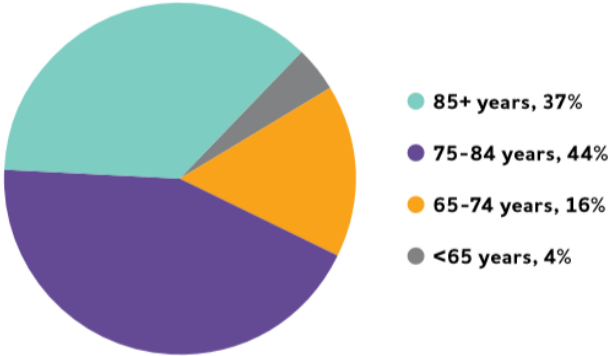


Figure 1.1: Age of people with AD in the US. Image from 2018 Alzheimer’s disease: facts & figures [1].

Early diagnosis of the disease is very important, since there is still no cure for AD. There are some drugs that can slow or even stop the progression of the disease, like rivastigmine, galantamine, donepezil, memantine, memantine combined with donepezil, and tacrine. These drugs work by increasing the amount of neurotransmitters in the brain. Their effectiveness varies from person to person and is limited in duration [1].

Part of identifying dementia in early stages is to identify Mild Cognitive Impairment, or MCI, which is a condition in which an individual has mild but measurable changes in thinking abilities, that affects 15-20% of people aged 65 or more. Individuals with MCI can usually carry everyday tasks, but the cognitive changes are noticeable to friends and family [1]. MCI is a major indicator that a subject may develop Alzheimer or other dementia: 32% of individuals with MCI developed AD in the following 5 years [1].

### **1.1.2 Diagnosis**

It is not always simple to diagnose Alzheimer's disease: doctors can very accurately determine that a subject has a type of dementia, but to specify which one is harder. During diagnosis they must look for specific signs that distinguish AD from other memory-affecting dementias. That is why the diagnosis is usually a combination of behavioral changes and biomarkers [3]. These biomarkers are part of computer aided diagnosis (CAD), which aims to complement the information that doctors have, so their diagnosis can improve, rather than replace the current medical knowledge [4].

The behavioral changes are identified by obtaining a medical and family history from the individual, asking a family member to provide input about changes in thinking skills and behavior, and conducting cognitive tests and physical/neurological examinations are some of the tools to provide AD diagnosis [1].

The biomarkers are signs of physical changes in the brain, like the accumulation of beta-amyloid plaques outside neurons and of protein tau inside neurons. The first change contributes to the decay of number of cells by interfering with synapses, and the second one causes the block of transportation of nutrients and other essential molecules to neurons. Other brain changes include inflammation and atrophy [1, 5]. AD (and MCI, with less severity) is also associated with reduced cortical thickness and surface area in a few brain regions associated with cognitive impairment. These changes are identified with the use of PET, MR imaging and CSF biomarkers [1, 5–7]. These methods can hint that a subject will convert to AD up to 15 years before symptoms appear [1, 6].

## **1.2 Objectives**

The problem of diagnosing AD that is being tackled here is multi-class and ordinal, since it consists of 3 classes and there is an implicit chronological order: subjects usually evolve from CN (cognitive normal) to MCI, and later to AD.

This thesis will study how different classifiers behave in different setups. It also briefly looks at how these classifiers are influenced by different imaging modalities. The three main classifiers are SVM, RF and NN (or RNN when 2 or more timesteps are being considered, except in section 4.6). A new classifier,

kSVM, is briefly introduced and tested. The setups analyzed are: the "classical" one, where the model predicts labels given training data from the same point in time, the one that looks at the future, predicting labels for a point in time given training data from a previous one, the one that takes into account the past and the present, by predicting labels given training data from the current time and previous ones and changes of diagnosis, where instead of predicting the class, the model tries to predict if subjects will keep the same diagnosis or get worse.

### **1.3 Thesis Outline**

The outline of the remaining document is as follows: in chapter 2 all relevant concepts for this work, like classification methods, data properties and performance metrics, are introduced and explained. Then, in chapter 3, there is a review of the latest methods used to classify AD using MRI and PET imaging and their results. In chapter 4, the ADNI database and the experimental results are presented. At last, in chapter 5, the conclusions of this work are presented, and some suggestions for future work are discussed.





# Chapter 2

## Theoretical Background

### 2.1 Classification methods

#### 2.1.1 Support Vector Machines

A support vector machine (SVM) is a supervised machine learning method for binary classification problems, that can be applied to both linearly separable and non-linearly separable data. It can also be extended to multi-class problems.

SVM was first proposed by Boser, Guyon and Vapnik, in 1992 [8]. In 1995, Cortes and Vapnik published the notion of soft margin, making SVM applicable to non-linearly separable training data [9].

The main idea behind SVM is to find the optimal hyperplane to separate the two classes in the training data. This hyperplane is the decision boundary. New data is classified simply by checking in which side of the boundary it falls.

#### Linear classification

##### Linearly separable data

The simplest implementation of SVM is for linearly separable data. Although this case is not useful for many real world applications, it is a good starting point to a more complex model.

Consider a set of  $L$  training points,  $x_i, i = 1, \dots, L$ , where each input  $x_i$  belongs to  $\mathbb{R}^D$  and to a class  $y_i \in \{-1, 1\}$ . To separate the training points from both classes, a hyperplane is defined:

$$w \cdot x + b = 0, \tag{2.1}$$

where  $w$  is the normal to the hyperplane and  $\frac{b}{\|w\|}$  is the perpendicular distance from the hyperplane to the origin. This hyperplane should be as far away as possible from the closest members of both classes, which are called support vectors. This means that the hyperplane is equidistant to both classes [10].

The constraint for the positive class, where  $y_i$  is equal to 1, is:

$$\mathbf{x}_i \cdot \mathbf{w} + b \geq +1 \Leftrightarrow y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0 \quad (2.2)$$

The constraint for the negative class, where  $y_i$  is equal to -1, is:

$$\mathbf{x}_i \cdot \mathbf{w} + b \leq -1 \Leftrightarrow \mathbf{x}_i \cdot \mathbf{w} + b + 1 \leq 0 \Leftrightarrow y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0 \quad (2.3)$$

Considering equations 2.2 and 2.3, it is possible to write both constraints as:

$$y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0, \forall i \quad (2.4)$$

To maximize the margin (distance from the hyperplane to the support vectors), which is equal to  $\frac{1}{\|\mathbf{w}\|}$ , is equivalent to minimizing  $\|\mathbf{w}\|$ , which in turn is equivalent to minimizing  $\frac{1}{2}\|\mathbf{w}\|^2$ . Therefore the problem can be formulated as:

$$\begin{aligned} & \text{minimize} && \frac{1}{2}\|\mathbf{w}\|^2 \\ & \text{subject to} && y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0, \forall i \end{aligned} \quad (2.5)$$

It is possible to adopt the Lagrangian formulation in order to deal with the constraints:

$$L_P \equiv \frac{1}{2}\|\mathbf{w}\|^2 - \boldsymbol{\alpha}[y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1, \forall i] \quad (2.6)$$

$$\equiv \frac{1}{2}\|\mathbf{w}\|^2 - \sum_{i=1}^L \alpha_i y_i(\mathbf{x}_i \cdot \mathbf{w} + b) + \sum_{i=1}^L \alpha_i \quad (2.7)$$

By differentiating  $L_P$  with respect to  $\mathbf{w}$  and  $b$  and reintroducing the results in 2.7, the dual form of the problem is found:

$$\begin{aligned} L_D & \equiv \sum_{i=1}^L \alpha_i - \frac{1}{2}\boldsymbol{\alpha}^T \mathbf{H} \boldsymbol{\alpha} \\ & \text{s. t. } \alpha_i \geq 0 \forall i, \sum_{i=1}^L \alpha_i y_i = 0 \end{aligned} \quad (2.8)$$

where  $H_{ij} \equiv y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$ .

The whole demonstration, without omitting steps, is shown in [10].

### Not fully linearly separable data

In order to extend SVM to data that is not fully linearly separable it is necessary to allow misclassified points, that is, points on the wrong side of the decision boundary. This is a soft margin SVM [9]. The soft margin is achieved by relaxing the constraint in equation 2.4, with the introduction of a positive slack variable  $\xi_i, i = 1, \dots, L$ . The new constraint is given by:

$$y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 + \xi_i \geq 0, \text{ and } \xi_i \geq 0 \forall i \quad (2.9)$$

Besides minimizing  $\mathbf{w}$ , now it is also necessary to minimize the number of misclassified points, so

the new optimization problem is given by:

$$\begin{aligned} &\text{minimize} && \frac{1}{2}\|w\|^2 + C \sum_{i=1}^n \xi_i \\ &\text{subject to} && y_i(x_i \cdot w + b) - 1 + \xi_i \geq 0, \forall i \end{aligned} \tag{2.10}$$

where C is an adjustable parameter.

## Non-linear classification

Using SVM to perform non-linear classification could be achieved by simply pre-processing the training patterns  $x_i$  by a map  $\phi : X \rightarrow F$  into a feature space F, followed by the standard support vector regression algorithm. However, this approach would easily become computationally infeasible, for both polynomial features of higher order and higher dimensionality [11].

This is where the *Kernel trick* comes in: the SVM algorithm depends only on the inner products between pairs of points  $x_i x_j$ , and not on each input  $x_i$ . It is enough to know  $k(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$ , and therefore there is no need to explicitly compute  $\phi$ .

The most common kernels are the polynomial (equation 2.11), the sigmoid (equation 2.12) and the radial basis function (RBF) (equation 2.13).

$$k(x_i, x_j) = (x_i \cdot x_j + a)^b \tag{2.11}$$

$$k(x_i, x_j) = \tanh(ax_i \cdot x_j - b) \tag{2.12}$$

$$k(x_i, x_j) = e^{-\frac{1}{2\sigma^2}\|x_i - x_j\|^2} = e^{-\gamma\|x_i - x_j\|^2} \tag{2.13}$$

## Multi-class classification

SVM can also be used for multi-class classification, however the best way to do this extension is still an on-going research issue. The most common methods are based in binary classification: one-against-all [12], one-against-one [12] and directed acyclic graph SVM (DAGSVM) [12].

**One-against-all** To classify data that belongs to  $k$  different classes, this method starts by building  $k$  different SVM models. The  $m^{\text{th}}$  SVM is trained by labeling all patterns of the  $m^{\text{th}}$  class in the training dataset with  $+1$  and all of the other patterns with  $-1$ . This leads to having  $k$  different decision functions,  $f_1(x), f_2(x), \dots, f_k(x)$ . A new data point  $x$  is labeled with the class that has the largest value of the decision function:  $x \equiv \operatorname{argmax}_{m=1, \dots, k} (f_m(x))$ .

**One-against-one** This method trains a model for every pair of classes, so if there are  $k$  classes it will train  $k(k-1)/2$  classifiers. The decision for a new data point  $x$  is made by majority voting: if the SVM for the pair of classes  $i, j$  says  $x$  is in the  $i^{\text{th}}$  class, the vote for this class is incremented. Otherwise is the  $j^{\text{th}}$  class that is incremented. This is also called the "Max Wins" strategy.

**DAGSVM** The training phase for this method is the same as for the one-against-one method (training  $k(k - 1)/2$  classifiers). To classify a new data point, rather than voting, it uses a rooted binary directed acyclic graph, with  $k(k - 1)/2$  nodes and  $k$  leaves. Each node corresponds to a binary SVM, and decides between the  $i^{\text{th}}$  and the  $j^{\text{th}}$  classes. Each leaf corresponds to a class. To classify a new sample the algorithm starts at the root node and decides left or right, based on the corresponding SVM. It goes like this through the graph, until it reaches a leaf node. This leaf node corresponds to the predicted class.

## 2.1.2 Random forests

Random forest is an ensemble learning technique, that can be used for both classification and regression. Although it is very simple, it is very powerful, achieving state-of-the-art accuracy in many problems [13].

The base classifier of random forest is the decision tree. A decision tree is a classification algorithm that recursively partitions a data set until all data belongs to a single class. It can be adapted to be used for regression. It is attractive due to being intuitive, with straight-forward training and the classification is extremely fast [14].

Decision trees' structure is made of a root and nodes, which can be internal or leaf nodes. The former correspond to data being internally split and the later correspond to single classes [15]. There are several algorithms to build decision trees; the most used ones are ID3, C4.5 and CART [15]. They all have the same basis: starting with a training set  $T$  and a root node, find the best attribute to splitting the node, using a measure of node impurity (can be misclassification error, entropy or Gini index) [16].

Besides decision trees, the random forest algorithm uses the bagging algorithm and the random subspace method. The idea of sampling from a random subspace was proposed by Tin Kam Ho in 1995 [14]. Later, in 2001, Breiman combined this idea with the idea of bagging [17], which resulted in the random forests that are currently used.

The random forest algorithm benefits from having multiple tree classifiers, because this compensates for the bias of a single classifier. The strength of random forest classifiers increases with the strength of individual tree classifiers, but decreases with their correlation. So, to improve accuracy, the "randomness" has to minimize the correlation between different trees, while maintaining strength [17]. This is achieved using sampling from a random feature subspace and bagging.

One of the methods used to ensure that the different trees are independent is randomly sampling from the feature subspace [14]. This means that each tree is constructed with only some features, rather than all of the available ones, that are randomly selected. The number of this feature subset is a hyper-parameter of the random forest algorithm.

The other method, bagging or bootstrap aggregation, consists in creating multiple training sets from a single one, by drawing samples with replacement from the original training set [17]. The number of training sets, and therefore of trained trees, is another hyper-parameter of random forests.

At last, the random forest algorithm needs to combine the decisions of the individual trees. It does

this by averaging the conditional probability of each class at the leaves [18].

In conclusion, the random forest algorithm starts by generating a predetermined number of training sets using bootstrap aggregation, then trains a decision tree for each training set, using a subset of features randomly selected from the features subset, and finally it aggregates all distributions.

### 2.1.3 Neural networks

The core ideas of neural networks started to be investigated in the 1950s, by Rosenblatt, who created the single layer perceptron [19], using as basis the mathematical model of an artificial neuron proposed by McCulloch and Pitts [20]. The next breakthrough happened in the mid-1980s, when multiple people independently rediscovered the backpropagation algorithm [21]. This led to the first successful practical application of neural nets, in 1989, when Yann LeCun combined the earlier ideas of convolutional neural networks and backpropagation, and applied them to the problem of classifying handwritten digits [22].

A generalization of the backpropagation algorithm, where it is used to forecast over time, proposed by Werbos in 1988 [23], led to the creation of RNN. However RNNs do not work well for long sequences of data: they tend to forget the distant past, due to the vanishing gradient problem. LSTM (Long Short Term Memory) [24] and GRU (Gated Recurrent Unit) [25] layers were invented to solve this problem, by adding a way to carry information across many timesteps. RNNs are discussed in more detail in the following sections.

Another popular type of neural networks is the convolutional neural network (CNN). This type of network is primarily used to recognize patterns in images, because there are layers which encode image-specific features in the architecture. However, CNN are still similar to "regular" NN: they are composed by layers, which in turn are composed by neurons. These have a score function, or weight, which is optimized through learning [21, 26].

Neural networks are made of several units (also called neurons), organized in different layers. The structure of these networks consists of an input layer, one or more hidden layers, and an output layer, with as many units as network outputs [27]. An example of a neural network is shown in figure 2.1.

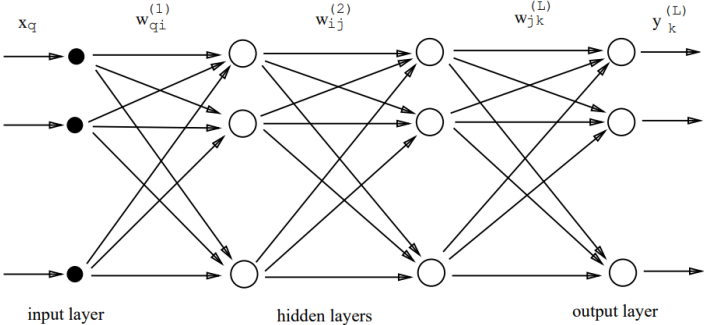


Figure 2.1: Simple neural network structure. Image from Artificial Neural Networks Tutorial [27].

The network in figure 2.1 is an example of a feed-forward architecture, because signals only go to the next layer's neurons. This is opposed to a feedback architecture, where there are connections between

neurons of the same or previous layer [27].

Each node receives  $n$  input signals  $x_i$  (with connection weights  $w_i$ ) which sum to an activation value  $y$  (that is,  $y = w_i \cdot x_i + b$ , where  $b$  is the bias of the node). A suitable transfer (or activation) function  $F$  transforms it into the output  $F(y)$ . The knowledge of a network is contained in the connection weights, which assume their values in the training phase. This training is often done using the backpropagation learning algorithm [28].

## Activation functions

Activation functions are needed to get access to a rich hypothesis space. Without them, a neural network dense layer would consist simply of two linear operations, and therefore the net could only learn linear transformations [21].

The most popular activation function in deep learning, for both input and hidden layers, is RELU (equation 2.14).

$$f(x) = \begin{cases} 0, & \text{for } x \leq 0 \\ x, & \text{for } x > 0 \end{cases} \quad (2.14)$$

The final layer can have different activation functions, depending on the purpose of the network [21]:

- for regression there is no need for an activation function at all, because the output of the layer is already linear;
- for binary classification or multiclass and multilabel classification a sigmoid activation function (like arctangent or the logistic function) is used. These functions returns a value between 0 and 1, that can be interpreted as a probability;
- for multiclass and single-label classification the softmax function is used. This function outputs a probability distribution over all the output classes.

## Backpropagation algorithm

The backpropagation algorithm is used to modify the connection weights  $w_i$ , in such a way that they minimize the value of the loss function. To do so, the gradient-descent algorithm is used. It starts from a generic data point  $p_0$  and calculates the gradient  $\nabla$ , which gives the direction in which to move in order to reach the maximum increase (or decrease, if  $-\nabla$  is considered). A new point  $p_1$  is found by moving from  $p_0$  a distance of  $\eta\nabla$  in the found direction, where  $\eta$  is the learning rate. The gradient is recalculated in the new point  $p_1$ , and algorithm continues iteratively until the gradient is smaller than a given threshold. The backpropagation algorithm can be divided into two steps: forward step, where the input data to the network is propagated through all the layers, and then the loss function is computed; and backward step, where the loss is propagated backwards, and the weights are updated appropriately [28]. The loss function that yields better results depends on the problem's characteristics (for example,

whether it is regression, binary classification or multiclass classification). Usually, binary or categorical cross-entropy are used.

### **Cross-entropy**

In machine learning, cross-entropy is used to define a loss function. When using NN, it is usually the better option to train a model for classification. For regression problems, the mean squared error loss function is the usual choice [21].

Cross-entropy was first introduced in the field of Information Theory, as a measure of how different the original and the received messages are. In machine learning, it can be interpreted as the distance between the real probability distribution and the predicted one. For binary classification problems, binary cross-entropy is used, and in multiclass problems, where there is only a correct label for each observation, categorical cross-entropy is used [29].

### **Batch-size**

It was stated in the previous section that the backpropagation algorithm updates the network weights. The amount of training patterns used in each iteration depends on the training mode: either all training patterns are used, which is called batch mode, or a single training pattern is selected, which is called on-line mode. The third option, which also the most used one, is using a subset of training patterns, which is called mini-batch mode. To choose the best mini-batch size, a compromise has to be reached: if it is too small, it is not possible to use optimized matrix libraries, saving computation time. On the other hand, if it is too large the weights are not updated often enough.

The optimal mini-batch size does not depend much on the other hyper-parameters, so by using some acceptable values for them and trying different mini-batch sizes, and then plotting the accuracy vs time, it is possible to find a mini-batch size that yields fast learning [30].

### **Overfitting**

In general, when training a neural network the performance of the model starts by improving rapidly: both the test and the validation losses decrease substantially. However, after a certain number of iterations over the training data, the model starts to overfit: the loss on the training data may continue to drop, but it is because the model has lost generalization ability, and is now learning patterns specific to the data at hand. Meanwhile, the loss on the validation data starts increasing [21].

The most straightforward solution to overfitting is to get more training data, although that is not always possible. Other solutions include reducing the network's size, weight regularization and dropout.

Reducing the network size is the simplest way of preventing overfitting. Smaller networks usually start overfitting later and, more importantly, their performance degrades more slowly once it starts. By reducing the number of layer or units per layer, the number of parameters to be learned are also reduced. This leads to the loss of memorization capacity, so the network will have to learn more general

representations rather than a mapping from the training input to the training output. There has to be some caution when reducing the network size, so that the network wont start to underfit [21].

Weight regularization is a way of forcing the network weights to be smaller, because larger weights means that the network is more complex. Like it was seen previously, with the network size, a simple explanation leads to better generalization. To apply weight regularization, a cost is added to the loss function of the network. The cost added can be either proportional to the L1 or L2 norm of the weights, which leads to L1 or L2 regularization [21].

Dropout is one of the most effective techniques to deal with overfitting. It is applied to layers, rather than to the whole network. It consists in randomly dropping a fraction of the weights (usually 0.2 to 0.5) of the layer during training. The idea behind this technique is that by adding some noise to the output of the layers the patterns with less significance, which are specific to the training data rather than a general trend, are broken [21].

## Recurrent neural networks

A recurrent neural network (RNN) is a type of neural network commonly used with sequential data that has an internal loop, that is, it processes sequences by iterating through the sequence elements and maintaining a state containing information relative to what it has seen so far [21].

Recurrent neural networks have a feedback architecture, presented in figure 2.2. This means that they have memory: they processes sequences by iterating through its elements and maintaining a state containing information relative to what it has seen so far [21].

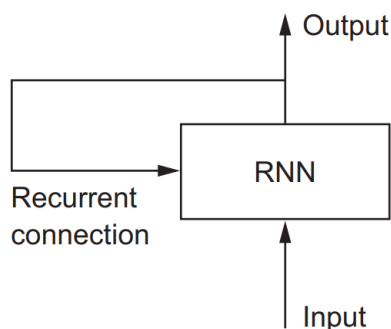


Figure 2.2: Simple RNN structure. Image from Deep Learning with Python [21].

It is also possible to visualize a neural network having the time explicitly represented, like in figure 2.3.

When training the RNN, the weights of the different layers are adjusted according to the backpropagation algorithm, which can be easily automated, using proper software like *Keras* and *TensorFlow*. However, it is not as simple to find the optimal hyper-parameters, like the architecture of the network (number of layers and how many units per layer), batch size, dropout rate and regularization. To tune hyper-parameters, the data is split into test, train and validation sets, and then cross-validation is used.



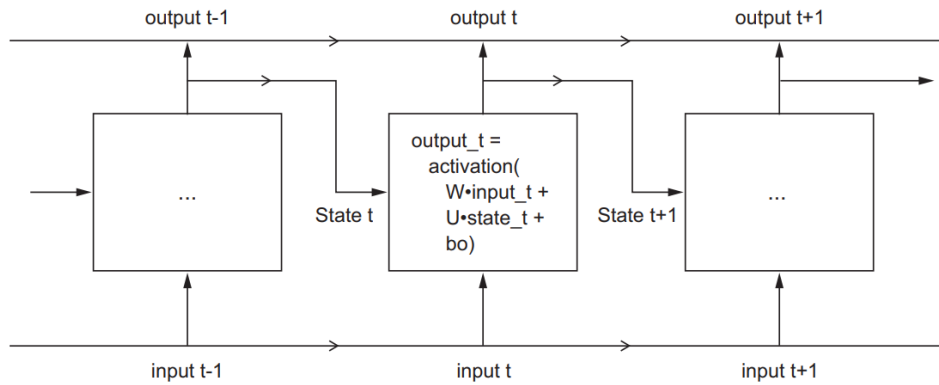


Figure 2.3: Simple RNN structure unfolded in time. Image from Deep Learning with Python [21].

### 2.1.4 kSVM

A fourth classification method was considered, with the objective of improving the SVM classification for multiple timesteps, as it is a simplified version of what RNN are to NN. The first SVM receives features from the first timestep, and proceeds as described above, except that it outputs the probability of each class for each subject in the dataset, rather than a single label.

These probabilities were an input to the following SVM, along with the features from the next timestep. The second SVM is similar to the first one. The  $k$  SVM, rather than outputting the probability for each class, outputs the final classification for each subject. A schematic for  $k = 3$  is presented in figure 2.4.

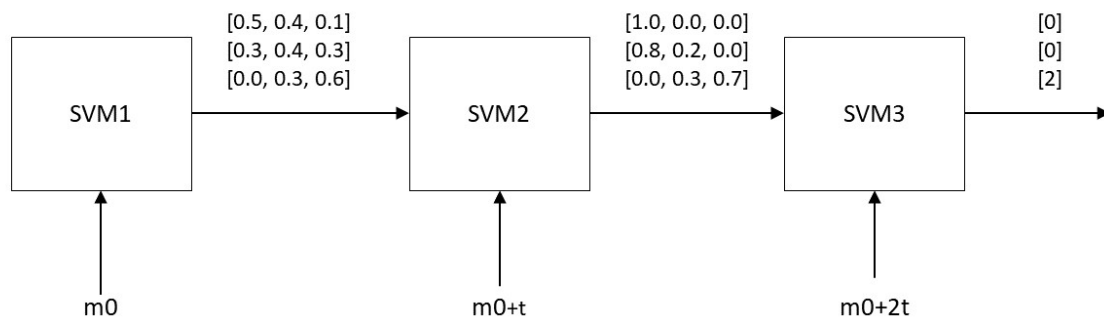


Figure 2.4: Structure of the classifier composed by 3SVM. The three sets of numbers represent hypothetical outputs for three input patterns.

### 2.1.5 Training

After building a model, and having the tools to train it, it is necessary to evaluate how well it behaves with new data. In order to do so, rather than training the model with all of the available data, the data is split into three independent sets, called training, test and validation.

The training set is used to train the model, using predetermined hyper-parameters. Then the model is evaluated using the validation set. The hyper-parameters are changed, the model is trained again on the training set and evaluated on the validation set, and so on, until the model has a good performance on the validation set. The hyper-parameters chosen in each iteration can be found either by grid search or by trial and error (which is not as computationally heavy, although it is less effective).

After the model is trained, it is evaluated on data it has never seen before, the testing set. It is very important to make no changes to the model based on information from the testing set, because it would insert a bias into the model and defeat the purpose of having an independent evaluation of the model [21].

## Cross-validation

If there are not many data points for training, the validation score may have a high variance depending on the validation split. Cross-validation is a tool to have a more robust validation.

There are several data resampling strategies, with different degrees of complexity. In practice,  $k$ -fold stratified cross-validation is often applied, with  $k$  equal to 5 or 10 [31].  $K$ -fold cross-validation consists in splitting the data in  $k$  partitions, and train  $k$  different models. Each model is trained into  $k - 1$  partitions and evaluated on the remaining one. The validation score of the model is the average of the  $k$  scores obtained from each of the  $k$  models. This method is represented in figure 2.5. The stratified part means that the proportions of classes in the different partitions reflect the proportions in the learning set.

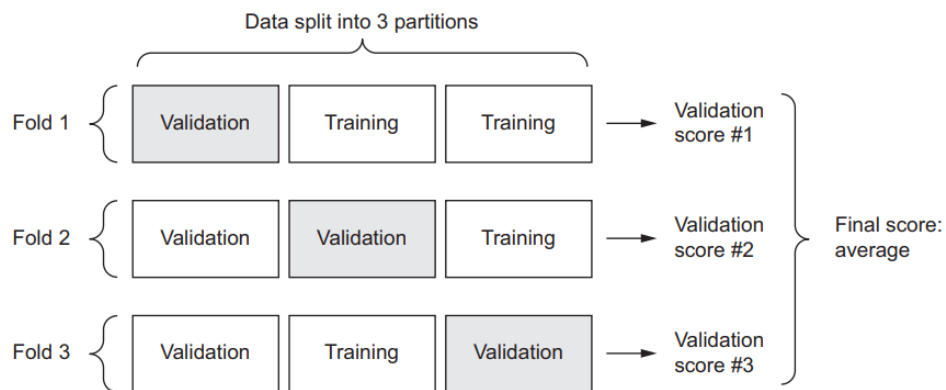


Figure 2.5: 3-fold cross validation. Image from Deep Learning with Python [21].

## 2.2 Data properties

### 2.2.1 Ordinal data

Measurements can be classified according to four different levels of measurement: nominal, ordinal, interval, and ratio quantities. Ordinal quantities, unlike nominal ones, have an intrinsic order in the different values they can assume. Interval quantities, in addition to this property, have values that can be measured in fixed and equal units, so the difference between two values is simply their subtraction. Ratio quantities have a zero point, which means that values can be multiplied [32].

Problems that have ordinal data can be solved using algorithms optimized for nominal data, however this throws away information. So, using algorithms adapted to ordinal data usually improves the results [32, 33]. Some of these algorithms are described in the next paragraphs.

## Adaptation of nominal algorithms to ordinal data

The method proposed in [32] to deal with ordinal data is a general adaptation of nominal algorithms to ordinal data. It starts with transforming the  $k$ -class ordinal problem into  $k - 1$  binary class problems. Consider an attribute  $X$ , with an ordinal label  $y \in \{V_1, V_2, \dots, V_k\}$ . The  $i^{th}$  binary problem corresponds to "is  $y$  larger than  $V_i$ ?", or simply  $y > V_i$ . Then, a new training set is derived for each of the  $k - 1$  binary problems, according to the following rule: for each attribute in the  $i^{th}$  problem, if in the original set the label was larger than  $V_i$  then the new label is 1, otherwise the new label is 0. In the next step, any suitable binary classification algorithm can be applied, to generate a model. For the  $i^{th}$  problem this is the probability that the class of an observation is larger than  $V_i$ , knowing the observations features, i.e.  $\Pr(y > V_i|X)$  for each of the  $k - 1$  datasets.

At last, it is necessary to merge the  $k - 1$  models to predict the probability of the original  $k$  classes, as described in equation 2.15.

$$\begin{cases} \Pr(V_1) = 1 - \Pr(y > V_1) \\ \Pr(V_i) = \Pr(y > V_{i-1}) - \Pr(y > V_i), 1 < i < k \\ \Pr(V_k) = \Pr(y > V_{k-1}) \end{cases} \quad (2.15)$$

To predict the class of a new observation, the probability of each  $k$  class is computed, using the model in equation 2.15. The class with maximum probability is assigned to the observation.

This method has been applied to the problem of detecting MCI and AD based on multimodal neuroimages (PET and MRI) and CSF biomarkers in [33]. The classification method used was a SVM, with an extra step of using a cost matrix to also tackle imbalance (imbalanced data is discussed in detail in 2.2.2).

It has also been used in [34] to estimate age based on pictures of subjects' faces, using a CNN as the classification method.

An alternative method, proposed in [35], is balanced ranking. It takes advantage of the ordinal nature of data, while also tackling imbalanced data. It has been used to classify subjects into AD, MCI and NC, according to multimodal neuroimages.

It uses pairwise ranking, which means that observations are trained in pairs. To create these pairs, it starts by transforming the original data domain  $(x_i, k_i)$  into a domain of binary differences  $(x_i - x_j, +1), (x_j - x_i, -1)$ , when  $x_i \succ x_j$ , i.e., the rank of  $x_i$  is greater than the rank of  $x_j$ . The distance of  $x_i$  to the hyperplane is the ranking score  $s_i$ . The definition of  $\succ$  only makes sense for ordinal data, so this algorithm cannot be applied to nominal data.

The next step is to train a SVM in this space. Then, it is necessary to convert the predictions from  $+1, -1$  back to the original classes, using the score  $s_i$  to train thresholds. These are defined by breaking points,  $\hat{k}$ , that are found according to equation 2.16, where  $\varepsilon = [\varepsilon_{k_i k}]$  is a cost matrix.

$$f(s_i, k_i, \hat{k}) = \min[\varepsilon_{k_i k} + f(s_{i+1, i+1}, \hat{k}), f(s_i, k_i, \hat{k} + 1)] \quad (2.16)$$

## 2.2.2 Imbalanced data

Although any data set that has unequal distribution between classes can be considered imbalanced, the term imbalanced data is more commonly used when the data set has a significant imbalance, in the order of 100:1, 1,000:1 or even 10,000:1. The main problem when using imbalanced data sets is that they significantly compromise the performance of most standard learning algorithms: they often fail to properly represent the distributive characteristics of the data, and in result provide unfavorable accuracy across the classes of the data. To deal with this problem there are some solutions: resampling, cost-sensitive methods, kernel-based methods, active learning methods, among others [36]. The first two are more common, and achieve state-of-the-art performance in many problems, so they will be explained in more detail. Despite not actually improving the classification, changing the performance metrics can also be a helpful tool when dealing with imbalanced data [36]. Performance metrics are discussed in more detail in Section 2.3.

### Resampling

Learning algorithms yield better overall classification performance when dealing with balanced sets rather than imbalanced. The main idea of resampling is to balance the number of samples between classes [36].

#### Random oversampling and undersampling

Random oversampling consists in appending to the original data set,  $S$ , an extra data set,  $E$ , gathered by sampling from the minority class. This means that the number of examples of the minority class in  $S$  will increase by  $|E|$ , which allows to adjust the degree of imbalance. However, if the number of samples in  $E$  is too large, the examples of the minority class are repeated many times, which leads to very specific rules, which in turn leads to overfitting [36].

Random undersampling is very similar to random oversampling: rather than appending data from the minority class, it removes data that belongs to the majority class from the original data set,  $S$ . The problem with this technique is that removing examples from the majority class may cause the classifier to miss important concepts regarding this class [36].

#### Informed undersampling

The objective of informed undersampling is to overcome the problem of information loss that occurs with random undersampling. The main idea is to use ensemble learning algorithms to learn which samples can be discarded without losing information. Two algorithms that have shown good results are EasyEnsemble and BalanceCascade. Another idea is to use a clustering technique or k-nearest neighbor to obtain a smaller dataset, so that the smaller dataset includes observations of every data group from the majority class [36].

## Synthetic samples

The SMOTE (synthetic minority oversampling technique) algorithm [37] creates artificial data points for the minority class based on the similarities between minority examples. To create a synthetic sample,  $x_{new}$ , it starts by considering an element of the minority class,  $x_i$  and its  $K$  nearest neighbors that belong to the same class. Then it randomly chooses one of the  $K$  neighbors,  $\hat{x}_i$ , and a  $\delta$  between 0 and 1. At last this vector is added to the original sample:  $x_{new} = x_i + (\hat{x}_i - x_i) \times \delta$ . Figure 2.6 presents an example of the generation of a synthetic sample using SMOTE, for  $K = 6$ .

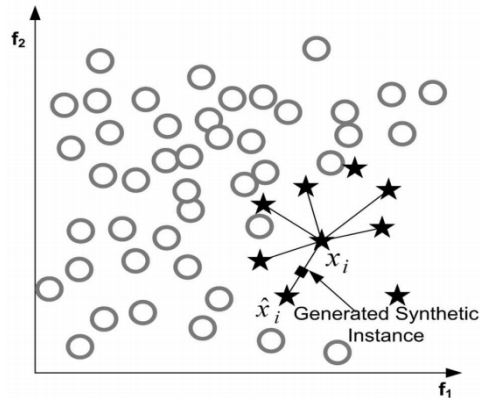


Figure 2.6: Example of the creation of a new sample using the SMOTE algorithm. Image from Learning from Imbalanced Data [36].

This method has the same benefits as random oversampling, without the drawbacks. However, the SMOTE algorithm still has some problems, including over generalization and variance [36].

## Cost-sensitive methods

Cost-sensitive methods, rather than trying to balance the data set, tackle the imbalanced data problem by using different cost matrices to describe misclassification. These square matrices are composed by as many rows as classes and each cell represents the cost of misclassifying a sample of the class  $j$  into the class  $i$ ,  $C(i, j)$  [36].

There are several ways of implementing these methods. The most general ones include dataspace weighting, where the misclassification costs are used to select the best training distribution, and applying cost-minimizing techniques to the combination schemes of ensemble methods. These two methods are commonly used together. However, if there is not an explicit cost matrix available, they cannot be applied.

The alternative is to incorporate cost-sensitive functions directly into classification paradigms, to fit the cost-sensitive framework into these classifiers, but these techniques are often specific to each particular paradigm. For example, in decision trees, cost-sensitive fitting can be applied in the adjustment of the decision threshold, at the split criteria at each node, and in pruning schemes. In neural networks, cost sensitivity can be introduced in four ways: first, cost sensitive modifications can be applied to the probabilistic estimate; second, the neural network outputs can be made cost-sensitive; third, cost-sensitive modifications can be applied to the learning rate; and fourth, the error minimization function can be adapted to account for expected costs [36].

### 2.3 Performance metrics

Performance methods are key to evaluating a classifier’s performance, and to compare different methods. While changing the performance metrics doesn’t actually improve the classification, it may provide a better understanding of the classifier. The most frequently used metrics, like accuracy and error rate, provide a simple, and often useful, way of describing a classifier’s performance. However, these metrics can be deceiving when dealing with imbalanced data sets [36, 38], which is why in these situations more complex methods, like f1-score, are preferred.

#### 2.3.1 Confusion matrix

Although it is not a performance metric, it is an useful indicator of the classifier’s behavior: it presents the number (or percentage) of hits/misses discriminated by class. If the data set is imbalanced, a naive classifier (or even a more complex one that has not been properly trained) that always choose the most frequent class may achieve a high accuracy, but by analyzing the confusion matrix is possible to detect this behavior. The structure of a confusion matrix is presented in figure 2.7. It consists in a square matrix, with as many rows as class. The cell (i,i) contains the number of elements that belong to class i and were correctly classified into class i. The cell (i,j) contains the number of elements that belong to class j, but were incorrectly classified into class i.

		True Class		
		Class 1	...	Class n
Predicted Class	Class 1	Number/% of elements from class 1 correctly classified into class 1		Number/% of elements from class n incorrectly classified into class 1
	...			
	Class n	Number/% of elements from class 1 incorrectly classified into class n		Number/% of elements from class n correctly classified into class n

Figure 2.7: Structure of the confusion matrix for n classes. Cells shown in blue correspond to correct classifications whereas cells in white correspond to incorrect ones.

#### 2.3.2 Accuracy

Accuracy is the most widely used metric for evaluating performances. It is defined simply as the percentage of correctly classified samples [38]. However, there are shortcomings when dealing with imbalanced data, as it was already mentioned. Accuracy also has problems when dealing with ordinal data: since only correct hits are considered, it does not care about how close/far the misses were.

### 2.3.3 F1-score

F1-score, or f-measure, is a performance metric designed for binary classification, however it can be extended for multiclass classification. It can be interpreted as a weighted average of precision and recall, as presented in equation 2.17, where  $\beta$  is usually 1. Precision and recall are defined as in equations 2.18 and 2.20, respectively, where TP is the number of true positives, FP is the number of false positives, and FN is the number of false negatives [38].

$$f1\text{-score} = \frac{(1 + \beta)^2 \cdot Recall \cdot Precision}{\beta^2 \cdot Recall + Precision} \quad (2.17)$$

$$Precision = \frac{TP}{TP + FP} \quad (2.18)$$

$$Specificity = \frac{TN}{TN + FP} \quad (2.19)$$

$$Recall = Sensitivity = \frac{TP}{TP + FN} \quad (2.20)$$

To apply this method to multiclass classification, the final f1-score equals the weighted average f1-score of each class.





## Chapter 3

# State-of-the-Art

The investigation of applying CAD to AD is very active at the moment. However, the large majority of studies considers only binary classification and a single time point. In this chapter we review approaches that tackle multiclass classification and/or multiple timesteps.

The following studies tackle multiclass while considering the ordinal nature of data, using an adaptation of SVM.

In [33], MRI and PET imaging and CSF biomarkers are used to classify subjects into CN, MCI-NC (MCI non-converter), MCI-C (MCI converter) or AD. The classifier used, SVM adapted to ordinal data, has already been described in section 2.2.1, subsection "Adaptation of nominal algorithms to ordinal data".

In [35], 2 experiments were done. In the first, the goal was to classify subjects into CN, s-MCI (stable MCI), p-MCI (progressive MCI) and AD. In the second, s-MCI and p-MCI were fused into a single class, MCI. It used FDG-PET and 1.5T MR image from the ADNI database. The proposed approach, balanced ranking, was described in section 2.2.1, subsection "Adaptation of nominal algorithms to ordinal data".

In [39], 3T-MR images from the ADNI database was used to classify subjects into CN, s-MCI, p-MCI and AD. It uses a SDPSO-SVM-PCA (SDPSO: Switching delayed particle swarm optimization, PCA: principal component analysis) approach. Their approach proposed a new SDPSO algorithm to optimize the SVM parameters. The PCA is used to improve the computational efficiency, by transforming the 32490 features into 38 linearly uncorrelated variables.

The following studies tackle binary problems, using sequential data to feed RNN or its variants. The ordinal nature of data is not taken into consideration.

In [40], a RNN was used to classify subjects into AD, CN or AUD (Alcohol Use Disorder). Regarding only the AD vs. CN part, the dataset used 404 patients with AD, from the ADNI database, and 274 healthy controls. It used CNNs to extract features from MRI images, that were then injected at each timestep in the RNN cells. Besides, it has a longitudinal pooling layer to fuse the features of a visit with the features of the previous ones, rather than feature concatenation or fully connected networks.

In [41], an ensemble of Bagged models was created, by combining three different models: CNN,

RNN and LSTM. It used MRI scan data from the OASIS Brain dataset to classify subjects into CN or AD.

In [42], subjects are classified into AD or CN using MRI data from the ADNI database. The method used consists in a MLP followed by a BGRU (bidirectional gated recurrent units) network (composed by the following layers: layer, masking, 2 BGRU, full-connected, activation, dropout, softmax). The spatial and longitudinal features are learned automatically from the imaging data of multiple time points.

Table 3.1 presents a summary of these studies, which includes the classification method, the number of participants and the main results.

Each of these studies presents a strong solution to classify AD, but there is not a global view. The present work tackles the classification of Alzheimer's disease by applying several classifiers, in several scenarios. By having the same dataset as starting point across all the experiments, it allows a better comparison between the different methods.

Table 3.1: Summary of the performance of different proposed methods. Metrics considered in multiclass problems: MAE - Mean Absolute Error. Metrics considered in binary problems: BACC - Balanced Accuracy. SENS - Sensivity. SPEC - Specificity. AUC - Area Under Curve.

Authors	Classification method	Participants	Results (%)
Y. Fan et al. [33]	SVM adapted to ordinal data	51 AD	$\rho$ 49.9
		44 MCI-C	Acc 39.6
R. Cruz, M. Silveira and J. Cardoso [35]	SVM adapted into a pairwise scoring ranker	57 MCI-NC	MAE 80.6
		55 CN	
R. Cruz, M. Silveira and J. Cardoso [35]	SVM adapted into a pairwise scoring ranker	58 AD	$\rho$
		51 p-MCI	3 classes 49.2
		84 s-MCI	50.7
		75 CN	55.2
			4 classes 56.8
		45.1	
		72.0	
N. Zeng et al. [39]	SDPSO-SVM-PCA	92 AD	Acc
		95 p-MCI	sMCI vs. pMCI 69.2
		82 s-MCI	NC vs. AD 81.3
		92 CN	NC vs. sMCI 76.9
			NC vs. pMCI 85.7
J. Ouyang et al. [40]	RNN	404 AD	sMCI vs. AD 71.2
		274 CN	pMCI vs. AD 57.1
M. Dua et al. [41]	Bagging ensemble of CNN, RNN and LSTM	BACC	90.4
		SENS	88.9
R. Cui & M. Liu & G. Li [42]	MLP + BGRU	SPEC	91.9
		AUC	91.5
M. Dua et al. [41]	Bagging ensemble of CNN, RNN and LSTM	ACC	92.2
		SENS	92.6
R. Cui & M. Liu & G. Li [42]	MLP + BGRU	SPEC	91.9
		# time points	ACC
		1	SENS 87.1
		2	SPEC 85.4
		3	88.6
R. Cui & M. Liu & G. Li [42]	MLP + BGRU	198 AD	2
		229 CN	3
			4
			5
			88.3
	85.9		
	90.4		
	86.4		
	90.8		
	86.9		
	92.1		
	86.9		
	92.6		



# Chapter 4

## Experiments and Results

### 4.1 Data description

The data used for the experiments comes from the ADNI database [43]. ADNI has data from 3 different phases: ADNI 1, from 2004 to 2009, and ADNI GO/ADNI 2, from 2010 to 2016 (ADNI3 data collection is currently in progress). This work uses data from the 3 available phases (ADNI 1, ADNI 2 and ADNI GO). The number of participants in each one, as well as their clinical status, are presented in table 4.1, where CN means Control Normal, EMCI is Early Mild Cognitive Impairment, LMCI is Late Mild Cognitive Impairment and AD is Alzheimer’s Disease. IN ADNI1 there was no distinction between EMCI and LMCI, so there are 400 subjects labeled simply as MCI. ADNI GO only collected data on EMCI subjects.

Table 4.1: Number of participants in each ADNI phase, according to their clinical status.

	ADNI1	ADNI GO	ADNI2	Total
CN	200	-	150	350
EMCI	400	200	150	900
LMCI		-	150	
AD	200	-	200	400

In the ADNI database there are different types of data: CSF markers of amyloid-beta and tau deposition, MRI, PET with three tracers (FDG, AV45 and AV1451), cognitive assessments, genetic information and general demographic information [44]. This work focuses on MRI and PET data.

MRI (magnetic resonance imaging) is a noninvasive imaging technique. Its developers, Lauterbur and Mansfield, were awarded the Medicine Nobel Prize in 2003 [45].

The main advantage of MRI (besides being noninvasive) is that its images have rich and versatile tissue contrast, which can be tuned by different parameters. Some of these parameters are proton density,  $\rho$ , longitudinal relaxation time,  $T_1$ , and transverse relaxation times,  $T_2$  and  $T_2^*$  [46].

Figure 4.1 illustrates how different weightings affect the contrast. These weightings are obtained by varying the echo time (TE) and repetition time (TR).

PET, or positron emission tomography, is an imaging technique that uses radioactive substances, or

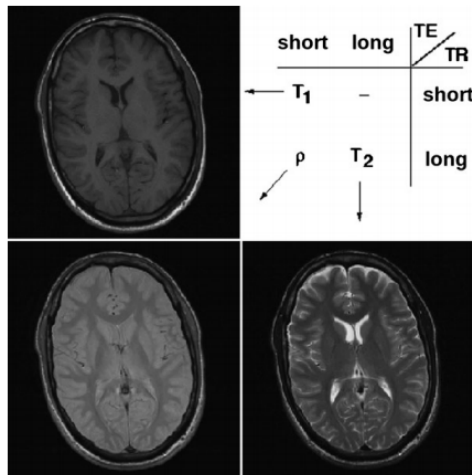


Figure 4.1: MR images of the same section of the brain, using different weightings. Image from Introduction to the Basics of Magnetic Resonance Imaging [46].

tracers, to detect functional changes. It is used both in the field of oncology and neuroimaging.

One of the most common tracers is 2-Deoxy-2-[ $^{18}\text{F}$ ]fluorodeoxyglucose (or [ $^{18}\text{F}$ ]FDG). This substance, which is a glucose analog, is consumed rapidly by healthy brain cells and tends to accumulate in the temporo-parietal cortical regions affected by AD. Although FDG PET is a well established method to diagnose AD, it has shown lower accuracy in older patients [7]. This has led to the development of more specific biomarkers, the amyloid precursors. The detection of amyloid precursors is able to predict the conversion of a subject at risk into AD 10 years prior to the onset of AD symptoms [7]. One of these substances is [ $^{18}\text{F}$ ]Florbetapir, or  $^{18}\text{F}$ -AV-45). This tracer is very specific to beta amyloid plaques ( $A\beta$ ) detection, and is quickly cleared from blood circulation [7].

MRI and PET data was collected for every subject at the beginning of the study, the baseline (or month 0). Most subjects have measurements for several timesteps: figure 4.2 presents the number of subjects by visit (every subject that went to visit  $n$  went also to visit  $n - 1$ ). The time interval between subjects' visits depends on the diagnosis and ADNI phase. A summary of the data collected in function of time and diagnosis is displayed in figure 4.3.

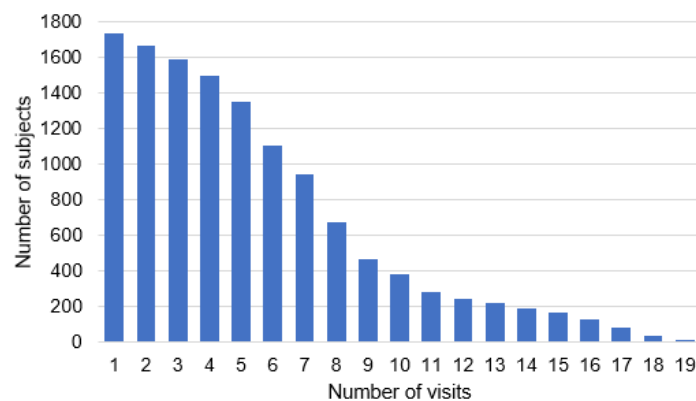
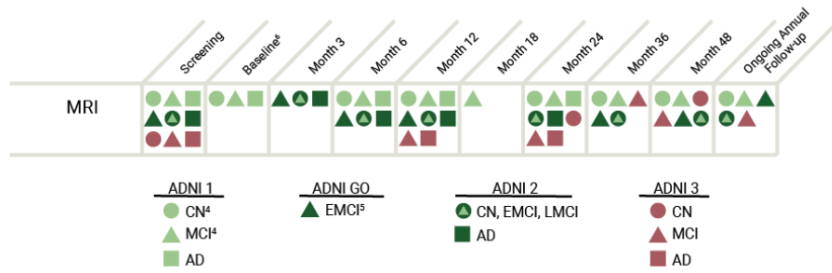
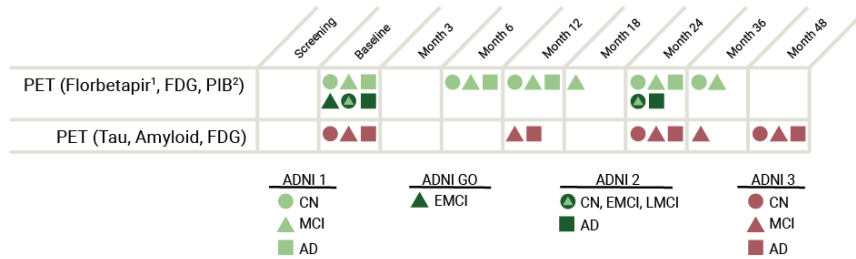


Figure 4.2: Number of subjects by visit.



(a) MRI



(b) PET

Figure 4.3: Data available in the ADNI database, broken down by diagnosis and ADNI phases, and by timesteps. Figures from [43].

The imaging data used was made available by the TADPOLE challenge [44]. Since it was pre-processed, the data used for classification was already in the form of quantitative numeric measurements.

The pre-processing of MRI scans included correction of gradient non-linearity, B1 non-uniformity scans and peak scans. The relevant regional features, like volume and cortical thickness, were extracted using Freesurfer cross-sectional and longitudinal pipelines. For PET scans, each image is composed by a series of dynamic frames. Its frames were co-registered, averaged across the dynamic range, standardized with respect to the orientation and voxel size and smoothed to produce a uniform resolution of 8mm full-width/half-max [44].

The start point for the remainder of this work was the resulting database.

The features used in the model, rather than selected through feature engineering, were features known to be relevant, based on previous research. More specifically, the features used were the surface area and volume of several regions of interest: ventricles, hippocampus, whole brain, fusiform, middle temporal gyrus and entorhinal cortex. The volumes were normalized with respect to the intracranial volume. These regions are known to show significant differences between NC, MCI and AD subjects [5], and are the same that were used successfully in [47].

After handling missing data, the MRI dataset had 83 features and the PET dataset had 54 features. The number of subjects in each dataset depended on the timestep being considered. Since PET data from ADNI2 was available only for month 0 and 24, in the following sections the dataset consists only of MRI data, unless stated otherwise.

The maximum number of timesteps considered was 5. There were more timesteps available (figure

4.3), but they were not evenly distributed.

## 4.2 Experimental design

This section describes the specific parameters and optimization process of each classifier and the methods used in the following experiments.

### 4.2.1 Data preparation

Often, subjects don't have measurements for all timesteps and for all features. RNN can deal with missing data, so it is as simple as using a masking layer, which is a type of layer used to signal the downstream layers that a timestep should be skipped [48]. Although RF can also deal with missing data, SVM cannot. To have a more direct comparison between SVM and RF, the database used is the same. It results from the following process: the features that are blank for more than a predetermined threshold of subjects are dropped. Then, all the subjects that are missing information for the remainder features are also dropped. The threshold is chosen so that the number of subjects dropped is minimized, with the condition of keeping at least 10 features.

Before every classification, the features are standardized (for every sample the mean is subtracted and it is divided by the standard deviation).

### 4.2.2 Performance score

All experiments were evaluated using a 5-fold cross-validation process. The final result was the average of the f1-score of each fold. For multiclass problems, the f1-score was computed using the weighted average of the f1-score of each label. For binary problems, the f1-score was simply the one of the positive class.

### 4.2.3 SVM and kSVM

For each fold, the train data is divided into train and validation sets. The train set is used to find the optimal hyper-parameters (kernel,  $C$  and  $\gamma$ ), through a grid-search (the kernel can be either linear or RBF,  $C$  is chosen from  $[10^{-4}, 10^{-3}, 10^{-2}, 0.1, 1, 10, 100]$ , and for a RBF kernel  $\gamma$  is chosen from  $[10^{-4}, 10^{-3}, 10^{-2}, 0.1, 1, 10, 100]$ ). The model is then trained using these parameters and the training set. The f1-score is obtained with the predictions from the test set.

### 4.2.4 RF

For each fold, the train set is split, to find the optimal hyper-parameters. However, since it would be computationally heavy to optimize a large number of parameters, only the number of trees and the max features are regarded in the optimization process. The number of trees is chosen from  $[100, 200, 500,$



750, 1000, 1500, 2000] and the maximum number of features is chosen from [None,  $\sqrt{n}$ ,  $\log_2(n)$ ], where  $n$  is the number of features.

Other parameters, like maximum depth, minimum samples before splitting and minimum samples per leaf, are not optimized. The model is then trained in the training set and evaluated in the test set.

#### 4.2.5 NN

The process starts by converting integer labels into categorical ones, using one hot encoding.

For each fold, two callbacks are used: early stopping and model checkpoint. The early stop one monitors the validation loss. When it has not decreased for a predetermined number of epochs the model stops its training. The model checkpoint monitors the validation accuracy, and saves the model weights that performed better. The model was built using keras library. After some trial and error, the network that showed better results was a sequential network, composed of three layers. The input layer had 25 neurons and the second one had 20. Both have RELU as the activation function. The dropout between these layers was set to 0.2. The number of neurons in the output layer is equal to the number of classes being considered. The output layer is activated by softmax.

### 4.3 Predictions for a timestep, given data from that timestep

The first classification problem tackled was to classify subjects into only two classes, AD and CN. It is a simple problem, to establish if there are significant differences in using different imaging modalities (MRI, PET, and both) and in using different classifiers (SVM, RF, and NN).

For data corresponding to month 0, there are 1737 subjects. Of those, 865 belong to either AD or CN. After dealing with missing data, there were 852 subjects and 83 features left for MRI data, and 431 subjects and 54 features for PET data. A dataset containing all the common subjects in the MRI and PET database was also considered. This dataset had 424 subjects and 137 features. The labels were converted to numeric values (0 for CN and 1 for AD). F1-scores in percentage and their respective standard deviation are presented in table 4.2.

A second set of results was obtained for month 0, using 3 classes (CN, MCI and AD). The number of features remained equal to what was described in the previous paragraph. The number of subjects in the MRI dataset changed to 1714, in the PET dataset to 890 and in the combined dataset to 877. The average f1-scores in percentage are presented in table 4.3.

A third set of results was obtained for month 0, using 4 classes: MCI was split into EMCI and LMCI. This information is available only for month 0. Table 4.4 presents the average f1-scores.

Table 4.2: Average f1-score for classifying subjects into CN or AD, with data relative to month 0.

	PET	MRI	MRI & PET
SVM	83.1 ± 2.8	84.4 ± 3.9	<b>86.4 ± 5.5</b>
RF	81.9 ± 2.4	83.8 ± 1.4	84.8 ± 5.7
NN	82.8 ± 3.4	83.4 ± 2.6	84.9 ± 5.6

Table 4.3: Average f1-score for classifying subjects into CN, MCI or AD, with data relative to month 0.

	PET	MRI	MRI & PET
SVM	51.9 ± 5.4	50.9 ± 1.7	53.7 ± 2.3
RF	51.2 ± 3.4	52.2 ± 2.0	53.1 ± 1.8
NN	53.4 ± 3.8	53.7 ± 2.6	<b>55.7 ± 4.1</b>

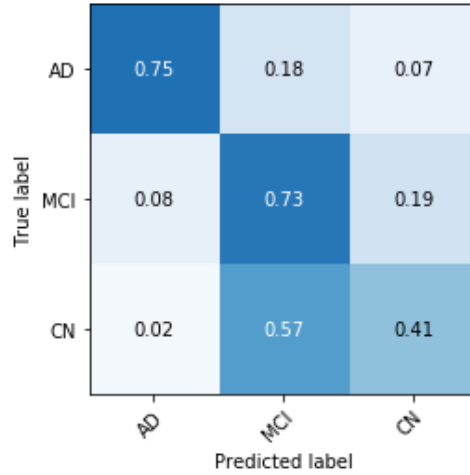


Figure 4.4: Confusion matrix obtained for classification of subjects in 3 classes, using MRI and PET imaging. The classifier used was a neural network.

Table 4.4: Average f1-score for classifying subjects into CN, EMCI, LMCI or AD, with data relative to month 0.

	PET	MRI	MRI & PET
SVM	41.9 ± 1.9	41.4 ± 7.1	47.5 ± 1.9
RF	41.1 ± 4.0	43.1 ± 3.4	<b>48.0 ± 3.8</b>
NN	42.3 ± 5.7	43.0 ± 2.8	47.6 ± 2.1

RF score could be improved by optimizing more parameters. NN could also be improved, by having a systematic optimization of parameters, rather than trial and error. However, both classifiers are already quite slow. SVM is by far the fastest classifier, and at the same time is the most optimized method. Its results could still be improved, by using a narrower grid to find the optimal parameters.

None of the classifiers under performed or outperformed significantly the others, so all of them were considered in the following sections. Having a combination of MRI and PET measurements always yielded better results than considering a single imaging modality. However, PET data from ADNI2 was available only for month 0 and 24, so in the following sections the dataset consists only of MRI data.

Adding a class between AD and CN lowers the scores significantly. The confusion matrices show that most of the misclassifications are between MCI and CN, as shown in figure 4.4. This CM was obtained for a combination of MRI and PET features, using a NN as the classifier. It is worth mentioning that approximately half of the subjects in the database belong to the MCI class, so it is expected that this class is favored. When MCI is split into EMCI and LMCI (this data is only available for month 0) there is no longer a bias towards this class. It becomes clear CN and EMCI are often mixed-up, and the same happens for AD and LMCI. This was verified for all the classifiers. As an example, the confusion matrix obtained using SVM, for the database composed of both PET and MRI features, is presented in figure

4.5.

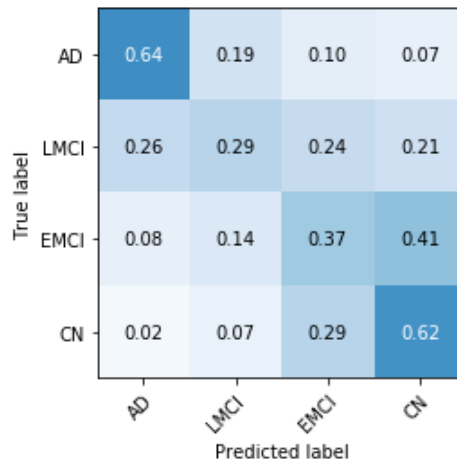


Figure 4.5: Confusion matrix obtained for classification of subjects in 4 classes, using MRI and PET imaging. The classifier used was a SVM.

It is likely that classification scores would be improved by filling in missing data rather than just not considering features/subjects with missing data.

#### 4.4 Predictions for a timestep, given data from the previous timestep

The classification methods used in this section were the same as the ones used in the previous one. The difference was in the training: rather than using data from a given month to obtain a prediction for that same month, the model is trained with features from a previous one and given labels of the given month. Therefore the model has data from the current month, and is trying to predict the future diagnosis.

Below are the results of trying to predict the labels of month 24, using features from month 24 (as a baseline), 12, 6 and 0. The classes considered were CN, MCI and AD. The subjects that converted from one class to another were considered into the class they were converting into (for example, if a subject is classified as "CN to MCI", it is considered MCI). The results are presented in table 4.5. The experiment was repeated using month 12 as reference. The results are presented in table 4.6.

Table 4.5: Average f1-score for predicting the subjects' class on month 24, with baseline information from month 24, 12, 6 or 0.

	Month 0	Month 6	Month 12	Month 24
SVM	52.0 ± 2.6	56.4 ± 4.1	58.2 ± 2.9	57.2 ± 2.6
RF	54.5 ± 2.5	58.3 ± 3.6	<b>59.8 ± 2.3</b>	58.5 ± 2.6
NN	53.8 ± 1.0	55.3 ± 2.5	56.4 ± 3.2	55.7 ± 3.4

The results show that the predictions get better when using more recent data, which was expected, since the test data used is more similar to the training data. However, unlike what was expected, the results are worse when using features and labels from the same month. For example, in table 4.5 the prediction using data from month 0 is worse than the prediction using data from month 6, which in turn

Table 4.6: Average f1-score for predicting the subjects' class on month 12, with information from month 12, 6 or 0.

	Month 0	Month 6	Month 12
SVM	51.3 ± 1.9	51.8 ± 3.9	52.2 ± 1.8
RF	50.2 ± 3.4	<b>55.7 ± 2.7</b>	54.4 ± 0.8
NN	53.7 ± 2.3	54.0 ± 3.2	53.8 ± 5.0

is worse than the prediction using data from month 12. However, the prediction using data from month 24 is worse than the one using data from month 12. The results in table 4.6 are similar, which discards the possibility of a statistical outlier. This behavior could be explained by some kind of systematic error, like the time between visitations not being exactly the one considered.

## 4.5 Predictions for a timestep, given data from more than one previous timestep

First, subjects were classified into three different classes (AD, MCI and CN) using data from 3 points in time: month 0, 12 and 24.

Two of the classifiers used were the usual SVM and RF. The "memory", or the ability to use information from more than 1 period of time, was created by concatenating features from several timesteps. In this particular case, the information was from month 0, 12 and 24. So, rather than having the usual 83 features, these classifiers had to deal with 249 features (83 features times 3 timesteps). This means that these methods had to look at all the information at once, without having the advantage of "broken-down" sequences.

The third classifier was a RNN. Since memory is its biggest advantage, there was no need to adapt it to receive information from several timesteps. Another upside is that it can easily deal with missing data, by using a masking layer. However, to allow for a more direct comparison, the database used consisted of the same subjects and features as the one used for other classifying methods. The network used consisted of 2 simple RNN layers, each with 32 neurons, dropout equal to 0.1, l1 regularization equal to 0.0001, l2 regularization equal to 0.001 and relu as activation function. The final layer was a dense layer with 3 units, with softmax as activation function. The callbacks Early Stop and Model Checkpoint were used to monitor the validation loss and the validation accuracy, respectively.

The fourth classifier considered was kSVM, with  $k = 3$ , which means that 3 SVM's were used.

For RF and SVM, the dataset had 249 features and 1016 subjects. The dataset used for both RNN and 3SVM had the same 1016 subjects, with 83 features for each timestep. The subjects were classified into AD, MCI and CN. The results obtained for each method are presented in figure 4.6.

RF achieved the best score, but the 5-fold validation took 44 minutes to run. SVM took only 89 seconds. RNN and 3SVM took 26 minutes and 6 minutes, respectively.

The same experiment was repeated, but using only the subjects who had measurements for all timesteps (that is, months 0, 12, 24, 36 and 48). The final database had only 196 subjects and 83

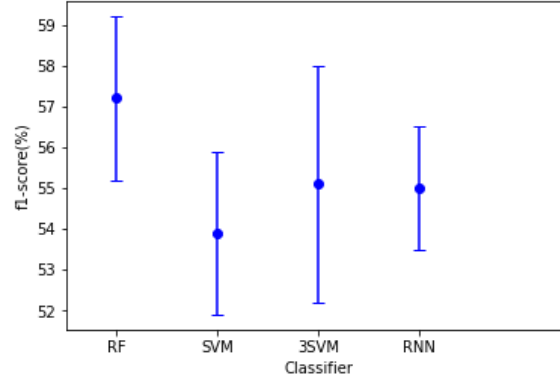


Figure 4.6: Average f1-score obtained with different methods for classifying subjects into three classes, using information from months 0, 12, 24.

features. Table 4.7 shows the average f1-score and its respective error for each method, obtained through a k-fold cross-validation with  $k = 5$ . The results of all methods plotted together, so they can be more easily compared, are in figure 4.7.

Apart from the result of RF for 3 timesteps, all methods have a similar behavior. Figure 4.7 shows that RNN is the worst method for only two timesteps, but that it improves with each new timestep. Unlike SVM and RF, that can perform well with few training data, RNNs tend to work better with more data. It would be interesting to see if, either with more subjects or with more timesteps, RNN would become the best classifier.

kSVM and SVM have a very similar behavior, despite kSVM being more complex and computationally heavier. Therefore there was no significant advantage in breaking up the features by timesteps before feeding it to kSVM, rather than just using the concatenated version.

Table 4.7: Average f1-score obtained with different methods for classifying subjects into 3 classes, in function of the number of timesteps.

# of Timesteps	2	3	4	5
SVM	$52.1 \pm 4.4$	$52.0 \pm 6.3$	$52.7 \pm 4.6$	$55.4 \pm 7.8$
RF	$52.3 \pm 5.8$	$56.6 \pm 5.2$	$54.2 \pm 5.1$	<b><math>56.7 \pm 7.1</math></b>
RNN	$50.4 \pm 8.2$	$51.2 \pm 6.4$	$52.9 \pm 5.5$	$55.9 \pm 4.6$
kSVM	$52.7 \pm 6.7$	$51.7 \pm 2.0$	$52.2 \pm 5.5$	$54.8 \pm 4.1$

## 4.6 Predicting changes of diagnosis

In this experiment, rather than predicting if a subject belongs to CN, MCI or AD, the focus was in predicting if a subject's condition got more severe. Using month 0 as starting point and month 24 as the ending point, there were 1020 subjects. In these 2 years, 30 subjects showed improvement (transitioned from MCI to CN or from AD to MCI), so they were removed from the database. 172 subjects transitioned to a worse state, so they were labeled as 1. 818 subjects kept their diagnosis, and were labeled as 0. The features used for SVM, RF and NN were the concatenation of the usual 83 features for months 0 and 24.

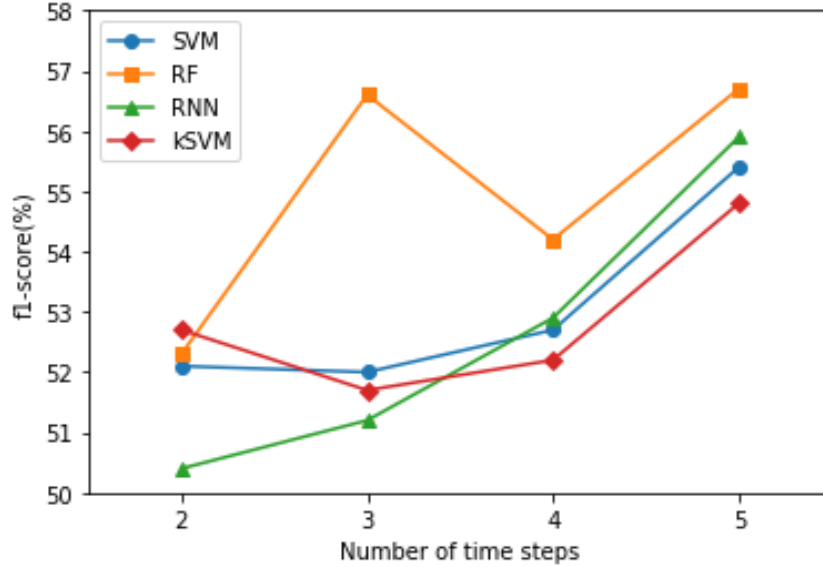


Figure 4.7: Comparison of average f1-score obtained with different methods for classifying subjects into three classes, with different number of time steps.

This dataset is much more imbalanced than the previous ones, so it was necessary to implement strategies to deal with it: random undersampling (RU), SMOTE and a combination of both. The results are presented in table 4.8.

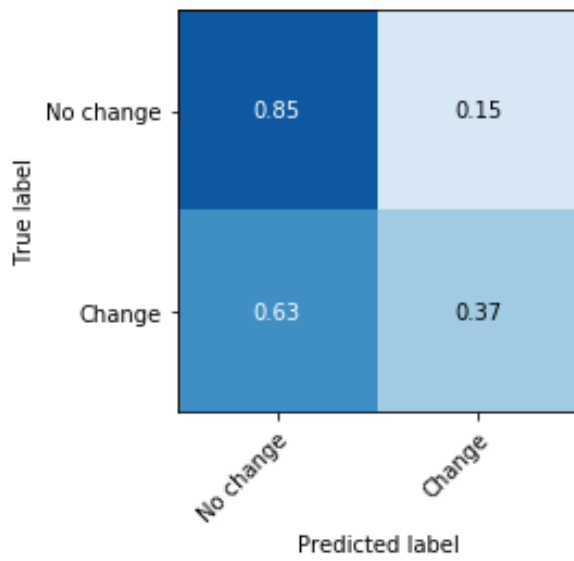
Table 4.8: Average f1-score obtained with different strategies to deal with class imbalance and different classifiers, for classifying subjects as having changed their diagnosis between month 0 and month 24.

	Imbalanced	SMOTE	RU	RU & SMOTE
SVM	$33.1 \pm 5.7$	$16.4 \pm 4.4$	$40.8 \pm 6.3$	<b><math>40.9 \pm 5.2</math></b>
RF	$12.7 \pm 9.7$	$16.6 \pm 8.7$	$27.1 \pm 6.8$	$35.4 \pm 2.1$
NN	$10.8 \pm 5.6$	$30.7 \pm 8.2$	$24.8 \pm 3.7$	$35.3 \pm 4.9$

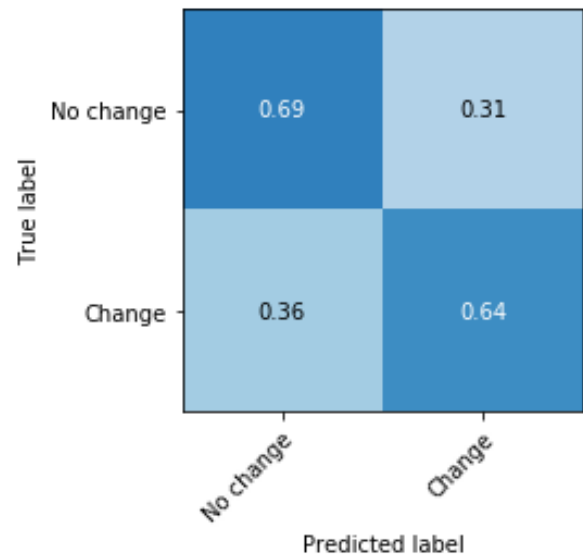
Figure 4.8(a) presents the confusion matrix obtained for SVM without using any approach to deal with the imbalance and figure 4.8(b) presents the confusion matrix obtained when using SMOTE & RU. Figures 4.9 and 4.10 present the same information, for RF and NN, respectively.

The combination of SMOTE and RU was the best approach for every classifier. As the confusion matrices in figures 4.8, 4.9 and 4.10 show, the accuracy of the minority class as always improved, at the cost of decreasing the accuracy of the majority class. Without dealing with the imbalance, both RF and NN classified almost all test patterns as "no change", which means they were initially as good as choosing always the majority class. Using SMOTE & RU, this behavior was altered.

Despite none of the classifiers having achieved a very good performance, SVM and NN reached an accuracy greater than 50% for both classes. To improve the overall performance, tuning the parameters of SMOTE and RU could be attempted. Another strategy could be choosing 2 points in time further apart, which would lead to more subjects converting to a more serious diagnosis, which in turn would result in a more balanced dataset.

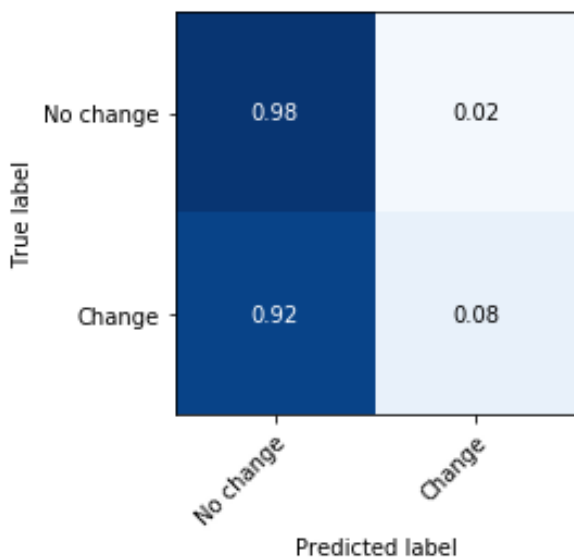


(a) Classification without mitigation of the imbalance

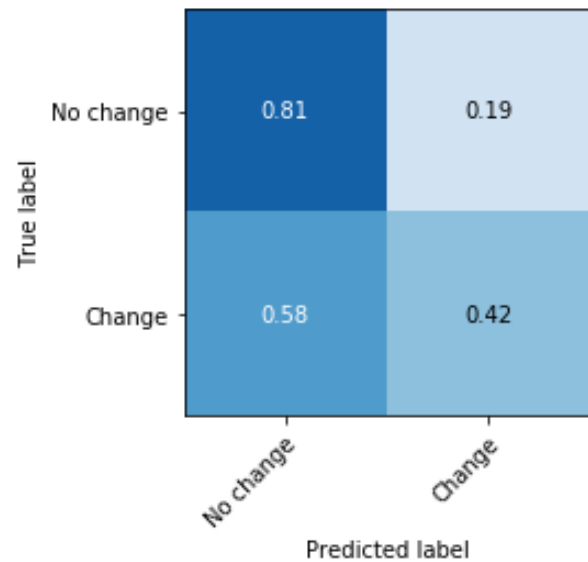


(b) Classification using SMOTE & RU to mitigate the imbalance

Figure 4.8: Confusion matrices obtained for classifying subjects into change/no change, using SVM. The time window considered was from month 0 to month 24.

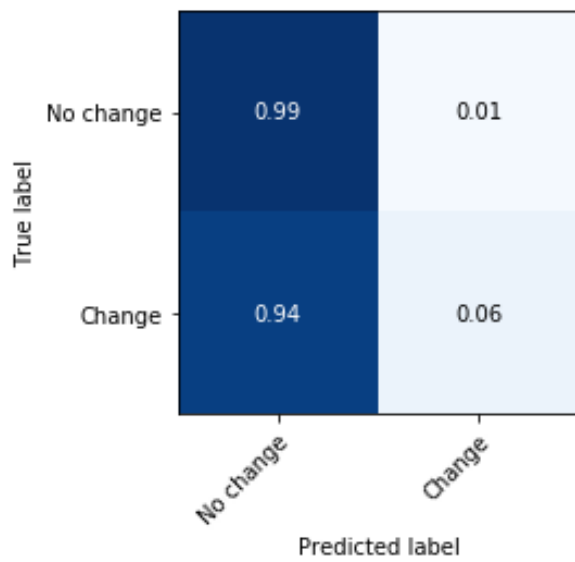


(a) Classification without mitigation of the imbalance

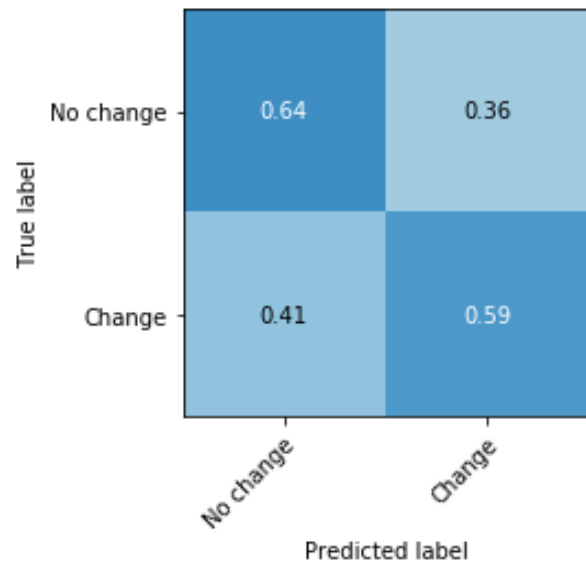


(b) Classification using SMOTE & RU to mitigate the imbalance

Figure 4.9: Confusion matrices obtained for classifying subjects into change/no change, using RF. The time window considered was from month 0 to month 24.



(a) Classification without mitigation of the imbalance



(b) Classification using SMOTE & RU to mitigate the imbalance

Figure 4.10: Confusion matrices obtained for classifying subjects into change/no change, using NN. The time window considered was from month 0 to month 24.



## Chapter 5

# Conclusions and Future Work

### 5.1 Conclusions

The present thesis analyzed the performance of different classifiers used to assign subjects into Normal (CN), Mild Cognitive Impairment (MCI) and Alzheimer Disease (AD), by using data from the ADNI database.

In the first experiment of section 4.3, the problem of classifying subjects into AD and CN was approached. This was a simple problem, without any real challenge. Its main objective was to test if every classifier was working properly. The results were similar for all classifiers and for both imaging modalities, PET and MRI. The best scores were obtained using a combination of PET and MRI features. The maximum f1-score obtained was 86.4.

In the second experiment of the same section, 3 classes were considered. As in the previous case, the best results were obtained when using a combination of MRI and PET measurements. The maximum f1-score obtained was 55.7. CN was often classified as MCI. Splitting MCI into EMCI and LMCI, it become clear that MCI is harder to identify. LMCI was often classified as AD and EMCI was often classified as CN. The maximum f1-score obtained using 4 classes was 48.0.

In the following section, 4.4, the classification was performed having access to features from a previous timestep and labels from the current one. RF was the classifier that performed better. With the exception of the original timestep, the results are worse for timesteps further away from the original one. However, the decay is not significant for 12 or 18 months. Unlike expected, the results were better using features from the previous timestep than from the current timestep. This effect was observed for both month 12 and 24, which points to it not being a statistical outlier. There could be some kind of systematic error, like the time between visitations not being exactly the one considered.

In section 4.5, information about several points in time was used to classify subjects into 3 classes: CN, MCI and AD. 4 classifiers were used: the well-known SVM, RF and NN and a new one, kSVM. This new classifier used as many SVMs as available timesteps, and the  $i^{th}$  SVM received the output from the  $i - 1^{th}$  SVM and the features belonging to timestep  $i$ . In the first experiment, using the bigger dataset, RF obtained the best f1-result: 57.1. However, this method was very slow, and the second

experiment showed that this value could be an outlier. RNN and 3SVM obtained a f1-score of 55 and SVM of 54. All of the methods achieved a better f1-score than the ones obtained in section 4.3. In the second experiment of this section, a smaller dataset was used. RF still had the best results, but the value for 3 timesteps seems to be an outlier. All classifiers had a tendency to improve their classification with more timesteps, but RNN's performance was the one that improved the most. The new classifier, kSVM, performed worse than the simple SVM, that received all the data at once, without information regarding which features belonged to each timestep. Since kSVM is computationally heavier, more time-consuming and hard to escalate to many timesteps, it is better to use the classical SVM. Even with a smaller dataset, all of the classifiers achieved average scores significantly higher than the ones obtained in section 4.3, for 5 timesteps.

At last, in subsection 4.6, the problem tackled was different: the subjects were now classified as having changed, or not, their diagnosis between 2 points in time. This dataset proved to be very imbalanced. Using the original classifiers, the "no change" class was chosen most of the time. Using SMOTE and RU, the results improved for all the classifiers, and SVM managed to achieve a f1-score of 40.9. None of the methods proved to be very good to tackle this problem, but SVM and NN reached an accuracy greater than 50% for both classes.

## 5.2 Future Work

The results obtained in the first experiment using MRI and PET were better than the ones obtained using only MRI. Most results in the present work should improve using a combination of MRI and PET data, rather than just MRI. There could also be some gain in studying the different features, rather than using the ones already used in previous work.

The ordinal nature of the data ended up not being used.

Subjects with missing data and features not present for the majority of subjects were simply not considered, which discarded a significant number of subjects/features. Rather than "throwing away" information, these "holes" could be filled, using a number of different imputation approaches.

All of the available parameters were optimized when using SVM. Future work could test if RF results would improve by optimizing more than just 2 parameters, and if NN results would improve by systematically optimizing every parameter, or even the architecture.

In section 4.4, the reasons for the classification to be better using past features rather than the features that correspond to the labels are unclear. It could be interesting to see if these differences are due to some systematic error or something else. Using data from later months, there could also be a more thorough analysis on how the error in classification varies in function of time difference between the measurements and the given label.

In the last experiment, corresponding to section 4.6, the classifiers had to deal with an imbalanced dataset. The strategies used to deal with it were random undersampling and SMOTE. This strategies were not optimized, as they were implemented using simply the default parameters. The next step could include tuning these parameters and/or experimenting with different approaches. Another possible test

is using data from more distant timesteps, so that there are more subjects converting, and the data becomes less imbalanced naturally, rather than adding synthetic samples or removing information.



# Bibliography

- [1] Alzheimer's Association and others. 2018 Alzheimer's disease facts and figures. *Alzheimer's & Dementia*, 14(3):367–429, 2018.
- [2] A. Wimo, M. Guerchet, G.-C. Ali, Y.-T. Wu, A. M. Prina, B. Winblad, L. Jönsson, Z. Liu, and M. Prince. The worldwide costs of dementia 2015 and comparisons with 2010. *Alzheimer's & Dementia*, 13(1):1–7, 2017.
- [3] B. Dubois, H. H. Feldman, C. Jacova, H. Hampel, J. L. Molinuevo, K. Blennow, S. T. DeKosky, S. Gauthier, D. Selkoe, R. Bateman, et al. Advancing research diagnostic criteria for Alzheimer's disease: the IWG-2 criteria. *The Lancet Neurology*, 13(6):614–629, 2014.
- [4] K. Doi. Computer-aided diagnosis in medical imaging: historical review, current status and future potential. *Computerized medical imaging and graphics*, 31(4-5):198–211, 2007.
- [5] H. Yang, H. Xu, Q. Li, Y. Jin, W. Jiang, J. Wang, Y. Wu, W. Li, C. Yang, X. Li, et al. Study of brain morphology change in Alzheimer's disease and amnesic mild cognitive impairment compared with normal controls. *General psychiatry*, 32(2), 2019.
- [6] R. Bateman, C. Xiong, and T. e. a. Benzinger. Clinical and biomarker changes in dominantly inherited Alzheimer's disease. *The New England journal of medicine*, 367(9):795–804, 2012.
- [7] S. Suppiah, M.-A. Didier, and S. Vinjamuri. The who, when, why, and how of PET amyloid imaging in management of Alzheimer's disease – review of literature and interesting images. *Diagnostics*, 9(2):65, 2019.
- [8] B. E. Boser, I. M. Guyon, and V. N. Vapnik. A training algorithm for optimal margin classifiers. In *Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory*, pages 144–152. ACM Press, 1992.
- [9] C. Cortes and V. Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, Sep 1995.
- [10] T. Fletcher. Support vector machines explained. *UCL*, 2008.
- [11] A. J. Smola and B. Schölkopf. A tutorial on support vector regression. *Statistics and Computing*, 14(3):199–222, Aug 2004.

- [12] Chih-Wei Hsu and Chih-Jen Lin. A comparison of methods for multiclass support vector machines. *IEEE Transactions on Neural Networks*, 13(2):415–425, March 2002.
- [13] C. Zhang, C. Liu, X. Zhang, and G. Almpanidis. An up-to-date comparison of state-of-the-art classification algorithms. *Expert Systems with Applications*, 82, 04 2017.
- [14] Tin Kam Ho. Random decision forests. In *Proceedings of 3rd International Conference on Document Analysis and Recognition*, volume 1, pages 278–282 vol.1, Aug 1995.
- [15] M. N. Anyanwu and S. G. Shiva. Comparative analysis of serial decision tree classification algorithms. *International Journal of Computer Science and Security*, 3(3):230–240, 2009.
- [16] W.-Y. Loh. Classification and regression trees. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 1:14–23, 01 2011.
- [17] L. Breiman. Random forests. *Machine Learning*, 45(1):5–32, Oct 2001.
- [18] T. K. Ho. The random subspace method for constructing decision forests. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (8):832–844, 1998.
- [19] F. Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386, 1958.
- [20] W. S. McCulloch and W. Pitts. A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5(4):115–133, 1943.
- [21] F. Chollet. *Deep learning with Python*. Manning Publications Co., 2018.
- [22] Y. LeCun, L. Jackel, L. Bottou, C. Cortes, J. S. Denker, H. Drucker, I. Guyon, U. A. Muller, E. Sackinger, P. Simard, et al. Learning algorithms for classification: A comparison on handwritten digit recognition. *Neural networks: the statistical mechanics perspective*, 261:276, 1995.
- [23] P. Werbos. Generalization of backpropagation with application to a recurrent gas market model. *Neural Networks*, 1:339–356, 12 1988.
- [24] J. Schmidhuber and S. Hochreiter. Long short-term memory. *Neural Comput*, 9(8):1735–1780, 1997.
- [25] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 12 2014.
- [26] K. O’Shea and R. Nash. An introduction to convolutional neural networks. *arXiv preprint arXiv:1511.08458*, 11 2015.
- [27] A. K. Jain, J. Mao, and K. M. Mohiuddin. Artificial neural networks: A tutorial. *Computer*, 29(3): 31–44, 1996.
- [28] M. Bodén. A guide to recurrent neural networks and backpropagation. *The Dallas project*, 12 2001.

- [29] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [30] M. A. Nielsen. *Neural networks and deep learning*. Determination Press, 2015.
- [31] D. Berrar. Cross-validation. *Encyclopedia of Bioinformatics and Computational Biology*, pages 542–545, 2019.
- [32] E. Frank and M. Hall. A simple approach to ordinal classification. In L. De Raedt and P. Flach, editors, *Machine Learning: ECML 2001*, pages 145–156, Berlin, Heidelberg, 2001. Springer Berlin Heidelberg.
- [33] Y. Fan. Ordinal ranking for detecting mild cognitive impairment and Alzheimer’s disease based on multimodal neuroimages and CSF biomarkers. In *International Workshop on Multimodal Brain Image Analysis*, pages 44–51. Springer, 2011.
- [34] Z. Niu, M. Zhou, L. Wang, X. Gao, and G. Hua. Ordinal regression with multiple output CNN for age estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4920–4928, 2016.
- [35] R. Cruz, M. Silveira, and J. S. Cardoso. A class imbalance ordinal method for Alzheimer’s disease classification. In *2018 International Workshop on Pattern Recognition in Neuroimaging (PRNI)*, pages 1–4. IEEE, 2018.
- [36] H. He and E. A. Garcia. Learning from imbalanced data. *IEEE Transactions on knowledge and data engineering*, 21(9):1263–1284, 2009.
- [37] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer. SMOTE: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002.
- [38] L. A. Jeni, J. F. Cohn, and F. De La Torre. Facing imbalanced data—recommendations for the use of performance metrics. In *2013 Humaine association conference on affective computing and intelligent interaction*, pages 245–251. IEEE, 2013.
- [39] N. Zeng, H. Qiu, Z. Wang, W. Liu, H. Zhang, and Y. Li. A new switching-delayed-PSO-based optimized SVM algorithm for diagnosis of Alzheimer’s disease. *Neurocomputing*, 320:195 – 202, 2018.
- [40] J. Ouyang, Q. Zhao, E. V. Sullivan, A. Pfefferbaum, S. F. Tapert, E. Adeli, and K. M. Pohl. Recurrent neural networks with longitudinal pooling and consistency regularization. *arXiv preprint arXiv:2003.13958*, 2020.
- [41] M. Dua, D. Makhija, P. Manasa, and P. Mishra. A CNN–RNN–LSTM based amalgamation for Alzheimer’s disease detection. *Journal of Medical and Biological Engineering*, pages 1–19, 2020.

- [42] R. Cui, M. Liu, and G. Li. Longitudinal analysis for Alzheimer's disease diagnosis using RNN. In *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)*, pages 1398–1401, 2018.
- [43] ADNI - Alzheimer's disease neuroimaging initiative. URL <http://adni.loni.usc.edu/>.
- [44] R. V. Marinescu, N. P. Oxtoby, A. L. Young, E. E. Bron, A. W. Toga, M. W. Weiner, F. Barkhof, N. C. Fox, S. Klein, D. C. Alexander, et al. TADPOLE challenge: Prediction of longitudinal evolution in Alzheimer's disease. *arXiv preprint arXiv:1805.03909*, 2018.
- [45] The nobel prize in physiology or medicine 2003, . URL <https://www.nobelprize.org/prizes/medicine/2003/summary/>.
- [46] K. Möllenhoff, A.-M. Oros-Peusquens, and N. J. Shah. Introduction to the basics of magnetic resonance imaging. In *Molecular Imaging in the Clinical Neurosciences*, pages 75–98. Springer, 2012.
- [47] M. M. Ghazi, M. Nielsen, A. Pai, M. J. Cardoso, M. Modat, S. Ourselin, and L. Sørensen. Robust training of recurrent neural networks to handle missing data for disease progression modeling. *arXiv preprint arXiv:1808.05500*, 2018.
- [48] Keras API, . URL <https://keras.io/api/>.