



# **Abstractive Multi-document Summarization using Topical Simplicial Curves**

**João Rui Serote Nunes Martins Cruz**

Thesis to obtain the Master of Science Degree in  
**Computer Science and Engineering**

Supervisors: Prof. David Manuel Martins de Matos  
Prof. Ricardo Daniel Santos Faros Marques Ribeiro

## **Examination Committee**

Chairperson: Prof. José Luís Brinquete Borbinha  
Supervisor: Prof. David Manuel Martins de Matos  
Member of the Committee: Prof. Alexandre Paulo Lourenço Francisco

**January 2021**



# Agradecimentos

Queria começar por agradecer aos meus orientadores, David Matos e Ricardo Ribeiro, por todo o apoio e ajuda que me têm dado ao longo destes dois anos. Influenciaram mais a minha vida e a minha maneira de pensar do que pensam; terei sempre em mim uma parte deles.

Um grande abraço a todos os meus primeiro colegas, depois amigos do laboratório do HLT, onde tivemos muitas conversas muito profundas (e muitas conversas nada profundas) no meio de muito (talvez demasiado) riso. Grande dia que foi aquele em que decidi perseguir um estágio no então L<sup>2</sup>F – não sabia bem no que me estava a meter, mas adorei todos os momentos.

Queria também agradecer aos Bengas de Setúbal, muitos para nomear aqui, e pedir desculpa por estar tão ausente durante este tempo todo. Não obstante, o seu apoio nos dias piores foi completamente imprescindível para me manter acima de água. Em particular um abraço ao Francisco, por me ter acompanhado mais perto que os outros durante esta viagem.

Um abraço a todos os meus amigos do NII com quem passei por grandes aventuras e com quem ainda mantenho contacto. Em particular ao Ben, ao Padipat (e um tributo especial a todo o dinheiro que perdemos nas arcades), e à Pla – ainda estou à espera do meu guarda-chuva!

Aos meus pais, claro, por todo o apoio que me têm dado, não só nesta fase da minha vida mas em todas as outras.

Por fim, uma palavra de agradecimento muito especial ao Atsuhiko por toda a paciência e tempo que teve para mim; à Takenaka por me ter tratado como se fosse da família dela: estou longe, mas parece que estou presente. Por último, um sentimento especial à Fujimoto, por me manter grounded durante todo este tempo.

Lisbon, February 9, 2021

João Cruz



# Resumo

Explorámos a eficácia das curvas simpliciais, um método de representação de palavras sensível ao contexto, motivados pelas suas propriedades matemáticas intrínsecas (e.g., diferenciação e facilidade de combinação de representações), na tarefa de sumarização multi-documento. Para este efeito, adaptámos o framework das curvas simpliciais para uma nova representação matricial com base em representações densas de palavras e desenvolvemos uma álgebra sobre objetos no simplex. Utilizamos os corpora de sumarização multi-documento DUC 2006 e DUC 2007. Os sumários gerados são comparados com os sumários de referência utilizando as métricas de avaliação ROUGE-1, ROUGE-2 e ROUGE-L. Comparado com a pontuação ROUGE-1 de 0.29 da baseline mais simples escolhida, o nosso método obtém uma pontuação ROUGE-1 de 0.04, ficando assim aquém das expectativas. Concluimos com uma exploração dos resultados obtidos e sugerimos outras aplicações do método das curvas simpliciais.



# Abstract

We explore the effectiveness of simplicial curves, a word-representation method that is context-sensitive, motivated by its intrinsic mathematical properties (e.g., differentiation and ease of combining representations), in the multi-document summarization task. To this effect, we adapt the simplicial curves framework to use dense word-representations as its basis matrix representation, and we develop an algebra over objects in the simplex. We use the DUC 2006 and DUC 2007 multi-document summarization corpora. The generated summaries are compared with the reference summaries using the ROUGE-1, ROUGE-2 and ROUGE-L evaluation metrics. Compared to the ROUGE-1 score of 0.29 of the simplest chosen baseline, our method achieves a ROUGE-1 score of 0.04, falling below our expectations. We conclude with an exploration of the obtained results and suggest other applications of the simplicial curves method.





# Palavras Chave Keywords

## *Palavras Chave*

Sumarização multi-documento

Curvas simpliciais

Representação de palavras

## *Keywords*

Multi-document summarization

Simplicial curves

Word representations



# Index

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Introduction</b>                          | <b>1</b>  |
| <b>2</b> | <b>Simplicial Curves</b>                     | <b>3</b>  |
| 2.1      | Background on Word Representations . . . . . | 3         |
| 2.1.1    | Traditional Representations . . . . .        | 3         |
| 2.1.2    | Neural Representations . . . . .             | 4         |
| 2.2      | Objects in the Simplex . . . . .             | 6         |
| 2.2.1    | Vocabulary Simplex . . . . .                 | 6         |
| 2.2.2    | Curve Overview . . . . .                     | 7         |
| 2.2.3    | Curve Construction . . . . .                 | 7         |
| 2.2.3.1  | One-hot Bag-of-Words-based . . . . .         | 10        |
| 2.2.3.2  | Dense Representation-based . . . . .         | 10        |
| 2.3      | Curve Algebra in the Simplex . . . . .       | 12        |
| 2.4      | Summary . . . . .                            | 14        |
| <b>3</b> | <b>Multi-document Summarization</b>          | <b>15</b> |
| 3.1      | Corpora . . . . .                            | 15        |
| 3.2      | Evaluation Metrics . . . . .                 | 16        |
| 3.3      | Related Work . . . . .                       | 18        |
| 3.4      | Summary . . . . .                            | 24        |

|          |                                   |           |
|----------|-----------------------------------|-----------|
| <b>4</b> | <b>Summarization Experiments</b>  | <b>25</b> |
| 4.1      | Datasets . . . . .                | 25        |
| 4.2      | Curve Construction . . . . .      | 25        |
| 4.3      | Summary Construction . . . . .    | 26        |
| 4.4      | Evaluation . . . . .              | 27        |
| 4.4.1    | Baselines . . . . .               | 27        |
| 4.5      | Results and Discussion . . . . .  | 28        |
| 4.6      | Summary . . . . .                 | 31        |
| <b>5</b> | <b>Conclusion and Future Work</b> | <b>33</b> |

# List of Figures

|     |   |    |
|-----|---|----|
| 2.1 | An example of $\Sigma_3$ , where each vertex corresponds to a word. The middle point $v$ is a distribution over each of the words. Being equidistant from all vertices, $v$ is a uniform distribution, that is, $v = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ . . . . . | 6  |
| 2.2 | LDA plate model. The outer plate represents documents, while the inner plate represents the repeated choice of topics and words within a document. In our case, $T$ is $\theta$ in the image. Image adapted from <a href="#">Blei et al. (2003)</a> . . . . .           | 11 |
| 3.1 | The four red dots represent the selected sentences since the polytope spanned by them covers all the other sentences in the document. Image adapted from <a href="#">Yogatama et al. (2015)</a> . . . . .   | 21 |
| 3.2 | Example of an AMR graph for the sentence “The Japanese Government stated on April 8, 2002 its policy of holding no nuclear warheads”. Image taken from <a href="#">Liao et al. (2018)</a> . . . . .   | 22 |
| 4.1 | Shortcoming of the curve-averaging method. $\gamma'$ has no clear relation with $\gamma_1$ or $\gamma_2$ . . . . .  | 29 |



# List of Tables

|     |   |    |
|-----|---|----|
| 4.1 | Baseline and curve results on the DUC 2006 dataset. . . . . | 29 |
| 4.2 | Baseline and curve results on the DUC 2007 dataset. . . . . | 29 |





# 1 Introduction

Automatic text summarization is a task where the goal is to, given one or more documents, produce a small text that accurately captures the information contained in the documents being summarized. This can be mainly done in two ways: a) by selecting important words or passages of the original text to preserve in the summary — extractive summarization, or b) by generating new words (or new sentences with the original words) that better synthesize what was in the original text (i.e., rephrase the text) — abstractive summarization. Despite working towards the same goal, these two approaches have differences in metrics and corpora used, and implementation methods.

For this purpose, is it fundamental to have a foundational framework that transforms the words of a text (such as the ones in this sentence) into objects that can be manipulated using mathematical operations, and, as such, can be used in computer applications. Different transformation methods encode different aspects of language (such as syntactic properties), and using one over the other is usually a matter of what trade-off is acceptable for a particular application. For example, if we are interested in automatically tagging the parts-of-speech of a text, we should choose a representation that enhances its syntactical aspects.

In this work, we are concerned with the re-exploration of a method to represent text in a manner that deviates from the current, well-established representation methods. This method — simplicial curves ([Lebanon et al., 2007](#)) — was chosen for its rich mathematical properties and the potential for finding parallels between fundamental analytical operations (integrals, derivatives) and results in the textual domain. Also, simplicial curves inherently encode the sequencing of text and, since we can define an algebra over this representation, lend themselves to composition. Simplicial curves have an intuitive sense of a document traversing in the space composed by its parts (often words but, as we will see, we can admit other definitions for lexical units), which can help with the explainability of the obtained results.

Our task is to explore the effectiveness of the simplicial curves approach in the field of

abstractive multi-document summarization (MDS — produce a textual summary from multiple documents, possibly from different sources, and talking about the same thing, while focusing on slightly different aspects), as opposed to single-document summarization (SDS — produce a textual summary from a single document), using MDS datasets. We also explore different methods of combining documents in a single representation and extracting text from it, evaluating our results using standard summarization literature metrics (ROUGE (Lin, 2004); see section 3.2).

This document is structured as follows: a) chapter 2 describes the simplicial curves method, with section 2.1 introducing historical context for different word representations in linguistic tasks; section 2.2 introducing the notion of simplex and objects embedded in it; and section 2.3 building an algebra of curves, showing how we can combine two curves into a third one or transform a curve into a numeric value; b) chapter 3 introduces the most used corpora, evaluation metrics, and addresses the related work done in the MDS task; c) chapter 4 delimits the corpora, evaluation metrics and baselines that we have used in the MDS task, with section 4.5 presenting the results from our experiments; and d) chapter 5 concludes the document, also pointing out some future work to be done.

# Simplicial Curves

In this chapter we introduce different word representation methods, introduce the concept of the simplex, and then develop the theory of simplicial curves with an accompanying algebra.

## 2.1 *Background on Word Representations*

In the following sections we will provide a description of different dimensional word representations. We are interested in different types of word representations because they will serve as the base representation upon which the simplex (see section 2.2.1) will be built.

### 2.1.1 Traditional Representations

The most basic transformation from words to a mathematical object that can be manipulated is called the one-hot approach. In this approach, words in a vocabulary  $V$  (with  $\#V = n$ ) correspond to dimensions in some space  $\mathbb{N}^n$ , and a word vector is represented by a 1 in the position for the word and 0 everywhere else. For example, if we have a text with two words, “test” and “red”, the vector for “test” is  $(1, 0)$  and for “red” is  $(0, 1)$ .

This basic model was further improved with the use of term-frequency (TF) (Luhn, 1957), where a word is represented by a one-hot vector multiplied by the word’s frequency in some document, which then is composed of stacks of word-vectors (a matrix). Another possible definition is that a document is the sum of the vectors of the words that are in the document. Under this model, a word is considered important in a document if it appears many times in it. This approach does not perform well in downstream applications when applied to texts where something as simple as function words (such as “the” — the most frequent word in English) are abundant, which are most texts. To counter this, the TF approach was enhanced by the addition of an “inverse-document-frequency” (IDF) term (Jones, 1972), where the multiplicative component of TF is weighted down by a function of the number of documents in which the

word occurs. The combination of both methods is known as TF-IDF:

$$\text{tfidf}(w, d) := \text{tf}(w, d) \times \text{idf}(w) \quad (2.1)$$

$$\text{tf}(w, d) := \frac{\text{times } w \text{ appears in } d}{\max\{\text{times } t \text{ appears in } d \mid \forall t \in d\}} \quad (2.2)$$

$$\text{idf}(w) := \log \frac{\#D}{\#\{d \in D \mid w \in d\}} \quad (2.3)$$

where  $D$  is a corpus,  $d \in D$  is a document, and  $w, t \in d$  are words in that document. Note that there are other sensible definitions for the TF function; we chose to define it as the normalized frequency.

Representing words as TF-IDF vectors accurately modeled text for many applications, but this method for word representation results in vectors that are too sparse for applications such as text classification and others. To counteract this, Latent Semantic Indexing (Deerwester et al., 1990) took the resulting sparse matrix and decomposed it using Singular Value Decomposition, allowing for the extraction of dense representations for words and documents that allowed the use of other similarity notions, e.g., two words/documents are similar if the cosine of the angle between their vectors is close to 1.

## 2.1.2 Neural Representations

Although representations derived from neural networks were already being studied (Collobert & Weston, 2008; Turian et al., 2010; Mnih & Hinton, 2007) since 2000 with their introduction to NLP in Bengio et al. (2000), it was in 2013 that the revolution of dense vector space representations of words derived from a neural network was kick-started by Mikolov et al. (2013) (here called SGNS, after the main method: “skip-gram with negative-sampling”). In this approach, word vectors are extracted from the inner state of a neural network after training on some proxy task (in their case, word similarity of a random word with the rest of the words in the enclosing sentence). The advantages of SGNS-based representations are that a) the resulting word vectors capture  $A:X::B:Y$  ( $A$  is to  $X$  as  $B$  is to  $Y$ ) analogies by means of simple arithmetic: if we want to solve  $A:X::B:?$  using SGNS embeddings, we can find the vector of  $?$  by  $v_? = v_X - v_A + v_B$ ; b) they are easy to incorporate in downstream applications, seeing as to get a word-embedding for a word  $w$ , all one has to do is lookup the row of  $w$  in the embedding matrix  $W$  (serving as a lookup table), after the neural network has been trained; and c) they indiscriminately improved

task scores on many different NLP benchmarks (Baroni et al., 2014).

Subsequent research expanded on this trend, with various extensions and modifications appearing over the decade. Of note are fastText in 2017 (Bojanowski et al., 2017) (rather than considering words to be the smallest textual unit, build instead vectors for each character in the text; a word vector is then the sum of the vectors of its characters), ELMo in 2018 (Peters et al., 2018) (words have different vector representations, depending on the context in which they appear), and more recently BERT in 2019 (Devlin et al., 2019) (a generalization of ELMo).

In 2014, a count-based approach (like TF-IDF) called GloVe was presented by Pennington et al. (2014), which aimed to compete with SGNS-embeddings in their representational power. Although their results in the word analogy, word similarity and named entity recognition tasks indeed showed better results than SGNS-based solutions, Levy et al. (2015) later showed empirically that GloVe is worse than SGNS for the word analogy and word similarity tasks on various datasets (they do not conclude that representations derived from counting approaches are worse than those derived from neural approaches, however), which reveals the variance of the impact of the representation for a given task, illustrating the importance of choosing the appropriate word representation for the task at hand. Moreover, further research (Levy & Goldberg, 2014) showed that SGNS are doing nothing more than factoring a Pointwise Mutual Information matrix derived from the text. In fact, how effective a method is in solving A:X::B:? analogies was discovered by Ethayarajh et al. (2019) to be caused by variance shifts in a modified formulation of the Pointwise Mutual Information between any two words.

$$\text{PMI}(x, y) := \log \frac{\Pr(x, y)}{\Pr(x) \Pr(y)} \quad (2.4)$$

Also, one major drawback of any of the above representations (apart from ELMo and BERT, indirectly) is that word order in a text is not preserved in the vector representation. These so-called bag-of-words methods construct their vectors as if any two words in a text are interchangeable, which is a fundamental oversight: text is sequentially structured, and we can obtain much information by order alone (e.g., in SOV languages we know that the last word in a sentence is most likely the verb). The relevance of this order information is, nevertheless, dependent on the task.

## 2.2 Objects in the Simplex

In the following sections we describe the notion of a simplex over an object space, and of curves in this simplex.

### 2.2.1 Vocabulary Simplex

A simplex  $\Sigma_n$  over some set  $C$  of size  $n$  is a subset of an  $\mathbb{R}^n$ -dimensional space where each dimension represents an object in  $C$ , and the coordinates of each point in this subset are non-negative and sum to one, i.e.,  $\Sigma_n := \{\mathbf{x} \in \mathbb{R}^n \mid x_i \geq 0, \sum_{i=1}^n x_i = 1\}$ . We can thus understand points in the simplex as probability distributions over the objects in  $C$ . We call  $C$  the collection and its objects the items. These can be concrete (such as real words in a vocabulary) or abstract (such as the topics of a document).

Geometrically, the simplex can be thought of as an  $(n - 1)$ -dimensional triangle, e.g.,  $\Sigma_3$  is the surface  $x + y + z = 1$ . Each dimension is thus a vertex in this triangle, and the notion of probability in this space is how close (under the  $L_2$  metric) a given point is to a vertex (see Fig. 2.1).

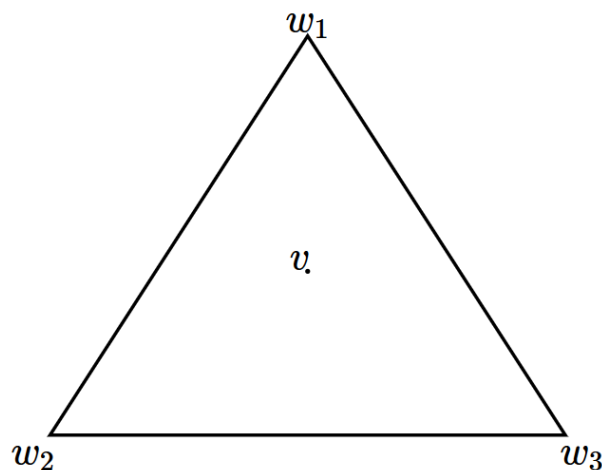


Figure 2.1: An example of  $\Sigma_3$ , where each vertex corresponds to a word. The middle point  $v$  is a distribution over each of the words. Being equidistant from all vertices,  $v$  is a uniform distribution, that is,  $v = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ .

### 2.2.2 Curve Overview

Simplicial curves were first introduced by [Lebanon et al. \(2007\)](#), motivated by the, at the time, lack of representation methods to model sequential content in textual documents. As a generalization of the bag-of-words representation, simplicial curves aim to preserve the same vector-space analogy as traditional representation methods, with the added benefit of modeling words locally — i.e., the document is a time-dependent histogram of words rather than a plain vector of all the words in the document.

The main idea is to model a text document  $y$  (or any sequence of discrete objects) as a continuous, sequential mathematical object, i.e., a parametric curve, where one can use standard calculus tools (e.g. derivatives, integrals, metrics) to model properties of the sequence. This also introduces the concept of time (represented by  $\mu$ ) in a document, taken to be between 0 and 1.  $\mu = 0$  represents the beginning of the curve, which maps to the beginning of the sequence, and  $\mu = 1$  represents the end of the curve, mapping to the end of the sequence.

The way the curve is built is flexible in that it allows one to model the original sequence at different levels of detail. If we choose a lower level of detail then the curve will focus more on the individual objects of the sequence (in the case of text, words); if we choose a higher level of detail then the curve will tend to model the sequence as a whole. Different representations of the same curve (or even of different curves) can be combined to produce a single curve that models the sequences at a varying level of detail.

A curve can be seen as a function  $\gamma_y$  from  $\mu \in [0, 1]$  (the point of the curve we want to query) to a member of  $\Sigma_n \subset \mathbb{R}^n$ , a point in the simplex (i.e., a distribution) over the objects in the original space.

Intuitively, the curve is a mapping from the normalized position of the document to a histogram of the words in that position in the document. All curves are the same length to allow comparisons between curves built from different-length sequences.

### 2.2.3 Curve Construction

Given an item sequence  $y$  of length  $N$ , the method starts by building an initial  $N \times C$  matrix,  $M_y$ , where the columns are the features of the objects (the vocabulary space) and the rows are

the vector representations of each of them. The matrix indices are then made continuous in time ( $t \in [0, 1]$ ) by making  $M_y$  be indexed not by its row numbers, but by a function  $\varphi$  defined as

$$\varphi_{M_y}(t, w) := [M_y]_{\lceil tN \rceil, w} \quad (2.5)$$

where  $[M]_{i,j}$  indicates the  $(i, j)$ 'th element of the matrix.  $M_y$  can be built in a variety of manners, which allows us to leverage several base-representations and the added benefits they contain. More details will be shown in the following sections.

As an concrete example, consider the following word sequence  $y = w_1 w_3 w_2 w_1$ , whose matrix is represented in a one-hot fashion:

$$M_y = \begin{array}{l} t = 0 \\ t = 1 \\ t = 2 \\ t = 3 \end{array} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

Surrounding the matrix are the indices we would normally use to index the rows. Using  $\varphi_{M_y}(t, w)$ , we get the following continuous access pattern:

$$\varphi_{M_y}(t, w) = \begin{array}{l} 0 \leq t \leq 1/4 \\ 1/4 < t \leq 2/4 \\ 2/4 < t \leq 3/4 \\ 3/4 < t \leq 1 \end{array} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

where  $t$  is used to index the rows and  $w$ , as before, the columns.

Once the continuous-access matrix representation is obtained, we smooth the entire matrix by multiplying in time (so only the  $t$  variable is involved) the access function  $\varphi_{M_y}$  with some smoothing kernel  $K_{\mu, \sigma}$ , where  $\mu \in [0, 1]$  is the center and  $\sigma > 0$  is the scale of the kernel.



For simplicity, a convenient choice for a smoothing kernel is a restricted Gaussian in  $[0, 1]$ :

$$K_{\mu,\sigma}(x) := \begin{cases} \frac{\mathcal{N}(x;\mu,\sigma)}{\Phi(\frac{1-\mu}{\sigma}) - \Phi(-\frac{\mu}{\sigma})}, & \text{if } x \in [0, 1] \\ 0, & \text{otherwise} \end{cases} \quad (2.6)$$

where  $\mathcal{N}$  is the probability density function of the Normal distribution and  $\Phi$  is its cumulative density function.

However, given that we only care about the shape of the kernel (it should have a bell-shape to emphasize points near its center, and have support only on  $[0, 1]$ ), we can choose to use other distributions as well. Another good candidate would be the Beta distribution  $B$  with appropriate parameters.

$$B(x; \alpha, \beta) := \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1} \quad (2.7)$$

The center parameter enhances the part of the curve we wish to focus on. Since the curve access is normalized to be in  $[0, 1]$ , we inspect a neighborhood of  $\mu$  by centering the smoothing kernel there. The scale parameter controls the shape of the kernel, which is fundamental to control the level of detail we want the curve to encode. If the kernel is too wide (i.e., the smoothing is too strong) then the multiplication will yield a not-too-detailed view of the document (this is, a higher-level representation of the document as a whole). If the kernel is too narrow (i.e., the smoothing is too weak) then the multiplication will yield a very peaky view of the document, precisely focusing on the original words. Formally:

$$\gamma_y^\sigma(\mu)_w := \int_0^1 \varphi_{M_y}(t, w) \times K_{\mu,\sigma}(t) dt \quad (2.8)$$

where  $\gamma_y^\sigma$  is a distribution over the original collection  $C$  built using the items in  $y$  and  $\gamma_y^\sigma(\mu)_w$  is the probability of item  $w$  from  $C$  at time  $\mu \in [0, 1]$  in the sequence. For convenience, we will generally drop the  $y$  and  $\sigma$  indices, only using them when we need to refer to curves built from different documents or using different kernel scales.

Next, we explore different ways of obtaining  $M_y$ .

### 2.2.3.1 One-hot Bag-of-Words-based

In the bag-of-words approach,  $M_y$  is built by stacking one-hot vector representations of the words in the order that they appear in the sequence. To avoid the sparse representation pitfall (see section 2.1.1), we apply some form of smoothing. Any classical form of smoothing can be used so, for simplicity, we use Laplace Smoothing by  $c \in \mathbb{R}$ .

Note that the larger  $c$  is, the closer all the previously-zero indices of the one-hot vectors will be, which means that the corresponding point will tend towards the center of the vocabulary simplex. Since all words will be smoothed the same way, the resulting curve will predominantly stay near the center of the simplex.

This approach offers two advantages: a) it is easier to grasp its intuition since every word is its own dimension, and b) it is easy to implement. However, preliminary experimental results show that the resulting representation may not capture important information that should be in the summary. A possible explanation for this is that, since a one-hot vector space has no intrinsic meaning, we are not leveraging information contained in the vector spaces of other forms of base representations (LDA (Latent Dirichlet Allocation (Blei et al., 2003)) with topics, SGNS with semantic similarity — see below).

### 2.2.3.2 Dense Representation-based

Instead of stacking the one-hot representation for words into a matrix, we can instead use the word  $\times$  topic weighting matrix given by LDA. In this case, the simplex dimensions will be the topics selected by LDA as being prevalent in the input text.

Formally, Latent Dirichlet Allocation (Blei et al., 2003) outputs a matrix  $T \in C \times K \subseteq \mathbb{R}^2$ , where  $C$  is as before and  $K$  is an LDA hyper-parameter that specifies the number of topics that should be extracted from the text. Each word  $w_i$  in  $C$  corresponds to a row in  $T$ , and it is this row, after normalization, that we take to be the  $i$ 'th row in  $M_y$  when we see  $w_i$  in  $y$ . Rows in  $M_y$  can and will be repeated for each repeated word that appears in  $y$ .

As an example, for a document with four words and three topics, a possible topic weighing

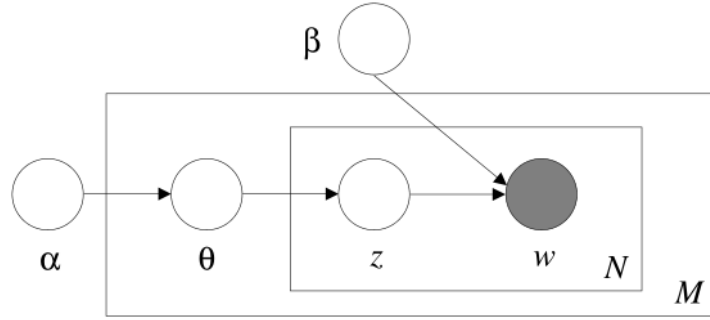


Figure 2.2: LDA plate model. The outer plate represents documents, while the inner plate represents the repeated choice of topics and words within a document. In our case,  $T$  is  $\theta$  in the image. Image adapted from [Blei et al. \(2003\)](#).

of those words could be:

$$T = \begin{bmatrix} 0.1 & 34 & 24.1 \\ 23.4 & 4.3 & 1.4 \\ 19.6 & 8.3 & 11.2 \\ 4.6 & 6.6 & 50.4 \end{bmatrix}$$

Note that each weight-vector of  $w_i$  given by LDA does not sum to one (does not form a probability distribution), so to map them onto the simplex we need to normalize them using a non-linear function (since if we used a linear function then vectors that were linearly dependent would collapse to the same point when mapped to the simplex). We choose to use the softmax function:

$$\text{softmax}(\mathbf{x}) := \frac{\exp(\mathbf{x})}{\sum_i \exp(x_i)} \quad (2.9)$$

where the exponential in the numerator is applied to the vector  $\mathbf{x}$  (a row-vector of the matrix) element-wise.  $M_y$  is then a dense matrix by construction, so no additional smoothing is necessary.

More generally,  $M_y$  can be built using pre-built vector representations such as those given by methods like word2Vec (SGNS) ([Mikolov et al., 2013](#)) or ELMo ([Peters et al., 2018](#)). In this way, we can adapt the benefits that these representations provide (analogy resolution, dense base-representations) with what simplicial curves gives us (ordering, algebraic operations). The same normalization precautions described above also apply in this case.

## 2.3 Curve Algebra in the Simplex

The main motivating idea is doing importance selection in the curve space and then mapping the results back to text. To achieve this, we need to have a way of combining curves in various ways, emphasizing different points.

With that in mind, we can define a basic algebra in the simplex, using familiar concepts such as addition and concatenation, allowing us to combine different points (or sets of points), indirectly combining the original items.

Since a curve is just a function  $\gamma : [0, 1] \rightarrow \Sigma^n \subset \mathbb{R}^n$ , we can consider the algebra of vector valued functions, with some additional operations to best allow the combination of curves, as long as they remain closed in the simplex.

Let  $\gamma_1$  and  $\gamma_2$  be two curves and  $\gamma'$  the result of combining them in some way. We define the following:

- Curve addition

$$\gamma'(\mu) = \frac{1}{2} (\gamma_1(\mu) + \gamma_2(\mu)) \quad (2.10)$$

- Curve subtraction

$$\gamma'(\mu) = \text{softmax}(\gamma_1(\mu) - \gamma_2(\mu)) \quad (2.11)$$

- Curve concatenation

$$\gamma'(\mu) = \text{if } \mu < \frac{1}{2} \text{ then } \gamma_1(2\mu) \text{ else } \gamma_2(2\mu) \quad (2.12)$$

- Curve conflation

$$\gamma'(\mu) = \frac{\gamma_1(\mu) \otimes \gamma_2(\mu)}{\sum_w \gamma_1(\mu)_w \otimes \gamma_2(\mu)_w} \quad (2.13)$$

All of these generalize to a higher number of curves.

Curve addition has an intuitive motivation: just return the curve in the geometric space between  $\gamma_1$  and  $\gamma_2$ .

Conflation (Hill, 2011) is a method used to compose different probability distributions over the same underlying objects whilst ensuring important statistical properties (e.g., conflation

minimizes the loss of Shannon Information, and yields a maximum likelihood estimator, among others). The vector multiplications are done element-wise.

Curve subtraction runs the risk of yielding a negative probability distribution, so we need to normalize it to positive by applying the softmax function.

Concatenation also has an intuitive meaning — take the beginning of the second curve and attach it to the end of the first, correcting the access argument accordingly. This can be useful for, e.g., forming a curve for a document by sequentially composing the curves for its sentences.

Do note that it does not make sense to consider curve scaling by some scalar in  $\mathbb{R}$  since we would immediately have to re-normalize, losing the scaling operation.

We can also define the curve inner-product, allowing us to see how much two given curves “agree” with each other, i.e., how similarly they travel in the simplex space.

- Curve inner-product

$$\gamma_1 \cdot \gamma_2 = \int_0^1 \gamma_1(\mu) \cdot \gamma_2(\mu) \, d\mu$$

Since the result of evaluating a curve at a given point is a distribution, we can also generalize common probabilistic descriptors such as entropy or the Fisher information:

- Curve Entropy

$$H(\gamma) := \int_0^1 H(\gamma(\mu)) \, d\mu \quad (2.14)$$

- Curve Fisher Information

$$\mathcal{I}(\gamma^\sigma) := \int_0^1 \mathbb{E}_{w \sim \gamma^\sigma(\mu)} \left[ \left( \frac{\partial}{\partial \sigma} \log(\gamma^\sigma(\mu)_w) \right)^2 \middle| \sigma \right] \, d\mu \quad (2.15)$$

We can also compare two different curves by finding their distance  $d$  under some metric  $\mathcal{M}$ , e.g.,  $L_2$  distance (yielding the mean euclidean distance of one curve to another in space), or the Jensen-Shannon metric, defined as  $JS(P \parallel Q) := \frac{1}{2} (\text{KL}(P \parallel A) + \text{KL}(Q \parallel A))$ , with  $A = \frac{1}{2} (P + Q)$ , where KL is the Kullback-Leibler divergence, yielding how similar two curves are from one another, in terms of their probability distributions.

- Curve difference under metric  $\mathcal{M}$

$$d_{\mathcal{M}}(\gamma_1, \gamma_2) := \int_0^1 \mathcal{M}(\gamma_1(\mu), \gamma_2(\mu)) \, d\mu \quad (2.16)$$

## 2.4 Summary

In this chapter we have given an overview of the different methods for word representations, traditional and neural-based, that have been developed. We then introduced the idea of a simplex and objects on that simplex. We detailed how to construct curves on simplices using different word representations as a basis, and developed an original algebra for these objects.

# Multi-document Summarization

In this chapter we are going to give an overview of some foundational and state-of-the-art methods in multi-document summarization, as well as detail some of the most used evaluation metrics and corpora.

## 3.1 Corpora

The Document Understanding Conference<sup>1</sup> (DUC) were a series of challenges running from 2001 to 2007 whose aim was to evolve the state-of-the-art in text summarization. To this end, a corpus for MDS was published each year, which invited competing implementations. The top-ranked systems became good baselines for MDS. Of note are the DUC 2006 and DUC 2007 corpora, which, for our purposes, are comprised of, respectively, 50 and 45 document clusters of English news from the Associated Press and the New York Times. Each document cluster has, on average, 25 documents.

Since 2008, DUC became the summarization track of the Text Analysis Conference<sup>2</sup> (TAC), where the goal was the same. The track ran from 2008 to 2011 (and uniquely in 2014) but, in recent years, TAC has grown to focus on knowledge-based systems. TAC challenges were more diverse, ranging from MDS to just summary evaluation, opinion summarization or even multilingual summarization. Of note is the TAC 2009 corpus, which is a dataset of 44 topics and 20 documents clusters per topic. The dataset is a subset of AQUAINT-2 (Vorhees & Graff, 2008), a collection of 907k documents in English, comprised from articles from October 2004 to March 2006 from Agence France-Presse, Central News Agency (Taiwan), Xinhua News Agency, Los Angeles Times, Washington Post News Service, New York Times and Associated Press.

Recently, Fabbri et al. (2019) introduced a new dataset for MDS along with some baselines on that dataset using MMR and end-to-end methods. It consists of 56216 documents taken from

---

<sup>1</sup><https://duc.nist.gov/> (Accessed January 20, 2021)

<sup>2</sup><https://tac.nist.gov/> (Accessed January 20, 2021)

scrapped full news articles and summaries from [newser.com](http://newser.com). Each summary is handmade and has at least two or more sources from where it was obtained. The baseline methods used were also tested in DUC and reported worse performance than known values for DUC. This highlights the fact that the effectiveness of different summarization techniques is also highly dependent on the type of text we want to summarize.

## 3.2 Evaluation Metrics

The classical metric used to automatically measure summary quality is Recall-Oriented Understudy for Gisting Evaluation (ROUGE) (Lin, 2004), which measures lexical overlap between the produced summary and some reference summary in various ways. The overlaps can be computed at the word level (ROUGE-1), bi-gram level (ROUGE-2), bi-gram with  $n$  words in between (ROUGE- $S_n$ ), bi-gram with  $n$  words in between and uni-gram overlap (ROUGE- $SU_n$ ), longest common subsequence (ROUGE-L), and some subsequent extensions considering dense vectors built from  $n$ -grams (ROUGE- $n$ -WE) (Ng & Abrecht, 2015) and co-occurrence statistics (Lin & Och, 2004). While ROUGE correlates well with human judgments for extractive summarization, it does not perform as well for abstractive summarization since the chosen new words may not overlap with the reference summary, although possibly preserving the general meaning of the text. Some work has been done in trying to take advantage of dense representations to measure similarity rather than semantic overlap (Ng & Abrecht, 2015), which also generalizes the ROUGE framework for abstractive summarization settings.

If ROUGE essentially measures the recall of the generated sentences (how much of the candidate sentence is in the reference summary), Bilingual Evaluation Understudy (BLEU) (Papineni et al., 2002) measures the accuracy (how much of the reference summary is in the candidate sentence). Although it originated in machine translation (ranking possible translations), BLEU was applied to summarization in the very first DUC challenge, but subsequent challenges evaluated performance using ROUGE.

In domains where we know *a priori* that a “good” summary of a document will necessarily paraphrase the original text, Summarization Evaluation by Relevance Analysis (SERA) (Cohan & Goharian, 2016) was developed to be a better (has higher correlation with human scores) metric than ROUGE. By measuring the summaries’ content relevance rather than lexical overlap,



SERA displays a perfect rank correlation ( $\rho = 1$ ) in the domain of scientific articles.

More recently, [W. Zhao et al. \(2019\)](#) advanced a metric — MoverScore — that measures the semantics of a candidate summary sentence in comparison to a reference summary, by formulating the text generation task as an optimal transport problem (i.e., find the cheapest (under some cost metric) deformation of a probability distribution into another). They tested this metric in a series of tasks such as machine translation, image captioning and, importantly, MDS, and found that MoverScore generally correlates better with human judgments than ROUGE.

For manual evaluation, [Nenkova et al. \(2007\)](#) proposed the Pyramid method: a framework for combining hand-crafted summaries from different human annotators, all the while accounting for the fact that there is high variability between the content different annotators produce (from the focus of the summary to the way it is produced — if in an extractive or abstractive manner). At the basis of the Pyramid method are Summary Content Units (SCU), passages that appear repeated throughout summaries (not necessarily lexically equal) weighted by how many times they appeared in a summary. SCUs with equal weight are then partitioned into layers (i.e., a pyramid) where the layers at the bottom (less weight) are informationally less important because they came from fewer summaries. After this arrangement, the pyramid is transformed into a final score (in  $[0, 1]$ ), where higher scores indicate that more of the content is as highly weighted as possible, via the following formula, where there are  $n$  layers  $T_i$  (each layer then has  $i$  weight; the bottom-most layer is  $T_1$ ) in the pyramid,  $D_i$  is the number of SCUs in the summary that appear in  $T_i$ , and  $X$  is the summary size in SCUs:

$$M := \left( \sum_{i=j+1}^n i \times \#T_i \right) + j \times \left( X - \sum_{i=j+1}^n \#T_i \right) \quad (3.1)$$

$$\text{Score} := \left( \sum_{i=1}^n i \times D_i \right) / M \quad (3.2)$$

where  $j = \max_{1 \leq i \leq n} \{ \sum_{t=i}^n |T_t| \geq X \}$ . Reference summaries are used to build the pyramid and a system summary can be assigned a number to see how informative it is.

### 3.3 Related Work

At its core, extractive summarization is a text ranking problem, where we have to choose the most important words to preserve in a final, shorter text. This formulation has a simple translation to SDS: rank and select the parts of the original text that should appear in the summary. However, this simplicity breaks down when we pass to the multi-document world: the set of documents to summarize may not talk about the same thing at the same level of detail, so we must identify and eliminate some redundancy. Also, we need to ensure that the produced summary is coherent with respect to the different source texts (Radev et al., 2002). A common way to solve this problem is to collapse the task into SDS: just concatenate all the texts. This, however, also creates new problems. News articles, for example, make an effort to have their first sentence be the most prominent (or “summarizable”), so the concatenated document of news articles would have multiple “first sentences” throughout. This also destroys the text’s narrative: the content no longer begins in the introduction and ends in the conclusion, it has now instead various phases where it begins and ends anew.

Multi-document summarization is a well-studied field, featuring many different approaches with various degrees of success, most of which fall into the pitfall of modeling it as an SDS problem. Nonetheless, we highlight below some recent or seminal work done in the area.

Goldstein et al. (2000) generalized the Maximum Marginal Relevance (MMR) (Carbonell & Goldstein, 1998) approach to extractive summarization to the MDS setting. MMR is a method for selecting sentences to include in a summary that a) provide new information, and b) are not similar to the already included sentences. Let  $s \in D$  be a sentence in a document  $D$ ,  $R$  be the set of sentences already chosen as relevant and that have been selected to appear in the summary, and  $Q$  be some user query. MMR is then defined as

$$\text{MMR}(D, Q, R) := \operatorname{argmax}_{s \in D \setminus R} \left\{ \lambda \operatorname{sim}_1(s, Q) - (1 - \lambda) \max_{s' \in R} \{\operatorname{sim}_2(s, s')\} \right\} \quad (3.3)$$

where  $\lambda \in [0, 1]$  is a parameter controlling if we want the selected sentence to be more relevant towards the query (as measured by some similarity metric  $\operatorname{sim}_1$ ) or more diverse (by metric  $\operatorname{sim}_2$ ) towards the already selected sentences. In Goldstein et al. (2000), the authors applied MMR to MDS by incorporating a series of document-independent statistical heuristics, such as number of documents that contain the query, the document where a selected sentence comes

from, the timestamps between document publications, among others. Testing was done using the TIPSTER corpus (Harman & Liberman, 1993), evaluated using both compression ratio and cosine similarity to reference human summaries.

Erkan & Radev (2004) introduced LexRank, a graph-based method for ranking TF-IDF (see section 2.1.1) represented sentences in a text document. LexRank first constructs an undirected graph with nodes representing sentences and edges representing the cosine similarity between sentences (where edges are only present if this similarity is above some threshold). It then applies the concept of eigenvalue centrality in graph-theory, achieved by multiplying the adjacency matrix by some initial distribution over all the vertices in the graph until convergence. A summary is then built by selecting the top  $n$ -th sentences, and evaluated using ROUGE-1 on the DUC 2003 and DUC 2004 corpora.

L. Zhao et al. (2009) did query-focused graph-based MDS extraction by selecting the top sentences that are closest to the query using LexRank. Afterward, these sentences are added to the user query, after which all the sentences are re-ranked according to this new query, paying attention to redundancy. This is done to reduce the information noise in the documents, and, as such, should generate better summaries. Testing was done in DUC 2005 and DUC 2006 using ROUGE-1, ROUGE-2, ROUGE-S, and ROUGE-SU4, showing that this method is comparable to the top performing systems in the DUC challenge for those corpora.

Nayem et al. (2018) did sentence fusion via walks on a graph. First, it constructs a graph where vertices are words and directed edges mean that the source word appeared before the target word in a sentence. Compression can be given by the shortest path between two words. This is done as a way to merge sentences from different documents, hence achieving higher coverage. Then it transforms the vertices into vectors via representations taken from the hidden state in a Recurrent Neural Network (the state itself is the word vector) and passed to TextRank (Mihalcea & Tarau, 2004) (similar to LexRank but with different criteria to build the graph) to capture semantic meaning in the absence of lexical overlap. The final sentences are retrieved from the TextRank output. For every word, they also decide if they want to include it in the final summary or if they want to replace it by a possibly more informative word for that context. The optimality of the length of the generated summary was validated on DUC 2004 using ROUGE-1, ROUGE-2, ROUGE-1-WE, and ROUGE-2-WE, beating the state-of-the-art in all the metrics of the dataset.

Tohalino & Amancio (2018) explored a multi-level graph approach where documents are modeled as layers and sentences are vertices in the layer. Inter/intra-layer edges are the cosine similarity of sentences (modeled as TF-IDF vectors; see section 2.1.1), where inter-layer edges are reinforced to give more weight to similar sentences in different documents. By using different layers for the documents, the method can naturally discriminate from which documents do sentences come from, helping with information redundancy, i.e., do not assign much weight to vertices (sentences) that are in non-informative layers (documents). Summaries are generated by picking the top vertices according to some measured classical graph-theoretic criterion such as degree, average shortest path, accessibility, and absorption time, achieving competitive results in ROUGE-1 on DUC 2002 and DUC 2004 and state-of-the-art in CSTNews (Cardoso et al., 2011), a corpus of Brazilian Portuguese journalistic texts, when compared against other graph-based systems (Ribaldo et al., 2012).

Kågeback et al. (2014) explored the viability for summarization of modeling sentences with semantically-aware representations (such as SGNS vectors, see section 2.1.2). To achieve this, they represented sentences using a simple (sentence vectors are given by the sum of the word vectors) and a complex method (sentence vectors are given by a recursive auto-encoder (Socher et al., 2011) that explicitly models the word order in the sentence and the grammar used). The dataset used is Opinosis (Ganesan et al., 2010) (short user reviews on different topics), evaluated using ROUGE-1, ROUGE-2, and ROUGE-SU4; they concluded that simpler methods (i.e., sentences are the sum of their words) outperform more complex ones.

Yogatama et al. (2015) represented each sentence by a vector given by Latent Semantic Indexing (see section 2.1.1). Given that a set of vectors (points) forms a polytope, a summary is built by selecting the sentences corresponding to the set of points that form the widest polytope over all the other sentences in the document cluster, the assumption being that the generated summary will be the one with the highest coverage (maximum volume) and least redundancy (points chosen are, by construction, at the boundary – see Figure 3.1). Their method was evaluated on TAC2008 and TAC2009 using ROUGE-1, ROUGE-2, and ROUGE-SU4, and compared against MMR and Coverage-Based Summarization (Gillick et al., 2008), where the generated summaries that cover more diverse bi-grams scored higher.

One should note that this is the same as finding the set of points that span the polytope with the largest area. This is similar to an approach that we explore to construct summaries, with the

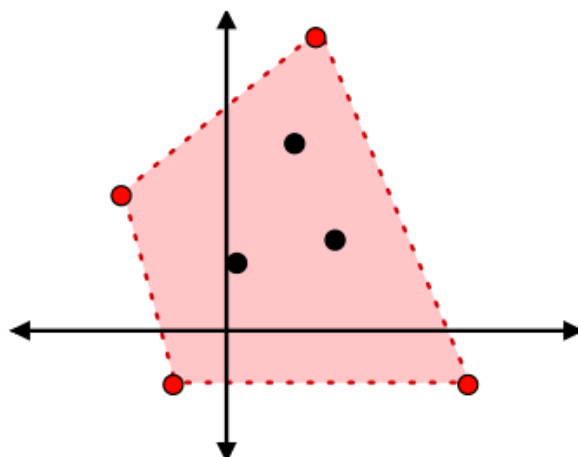


Figure 3.1: The four red dots represent the selected sentences since the polytope spanned by them covers all the other sentences in the document. Image adapted from [Yogatama et al. \(2015\)](#).

difference being that we find the set of points that span the smallest volume polytope.

[Mani et al. \(2018\)](#) represented a document by a paragraph vector ([Le & Mikolov, 2014](#)) and, for a document cluster, its centroid is seen as the “content average” of all the documents in the cluster. Summarization is done by selecting sentences that minimize the euclidean distance to the centroid of the cluster. They evaluate on DUC 2006 and DUC 2007 using ROUGE-1, ROUGE-2, and ROUGE-SU4. The idea of averaging representations of different objects into a final statistic is close to one of our proposed approaches, the main drawback in the case of [Mani et al. \(2018\)](#) is that their representation model needs to be pre-trained on some other dataset (they chose Thomson Reuters Text Research Collection ([Lewis et al., 2004](#)) and CNN/Dailymail ([Hermann et al., 2015](#)) corpora).

[Peyrard et al. \(2017\)](#) chose to focus on a different aspect of MDS: the metric. They combine a series of metrics such as ROUGE-N, ROUGE-L, ROUGE-WE, Jensen-Shannon Divergence, cosine similarity, and  $n$ -gram coverage in one final combination that optimizes correlation with human judgments. They trained a model based on a Support Vector Machine to optimize the combination of metrics and tested on the TAC 2008 and TAC 2009 corpora, concluding that this combination of metrics has a high correlation ( $r \approx 0.77$ ,  $\rho \approx 0.70$ ) with pyramid scores, higher than the same model trained with each individual metric.

[Rioux et al. \(2014\)](#) created a Reinforcement Learning agent that selects sentences from a document cluster to include in the summary. The sentence features are based on bigrams and

heuristics like “longest common subsequence length”. Delaying the reward the agent receives for a well-selected sentence is found to have much better performance than not delaying. Summaries were produced for DUC 2004, evaluated using ROUGE-1, ROUGE-2, and ROUGE-L, and compared mainly against Automatic Summarization using Reinforcement Learning (Ryang & Abekawa, 2012) (another RL agent), achieving higher performance. By incorporating these metrics in the reward function for the agent, the generated summaries were found to be much more grammatical (e.g., correct determiner use). This is due to the system extracting passages that are more than one word long (i.e., increasing ROUGE-2 and ROUGE-L).

Liao et al. (2018) used Abstract Meaning Representation (Banarescu et al., 2013) for abstractive summarization. Under the AMR framework, documents are represented by graphs where concepts in the text are vertices and semantic relations are edges. After parsing a document to obtain an AMR graph, sentences are generated which preserve the core semantics of the text. Here, this approach is applied in the MDS setting by considering a document cluster to be just a bag of documents, where each document is a bag of sentences. Documents were taken from DUC 2004 and TAC 2011, and compared against a series (see section 4.4.1 — we also test our method against the systems they chose) of both extractive and abstractive baselines using ROUGE-1 and ROUGE-2.

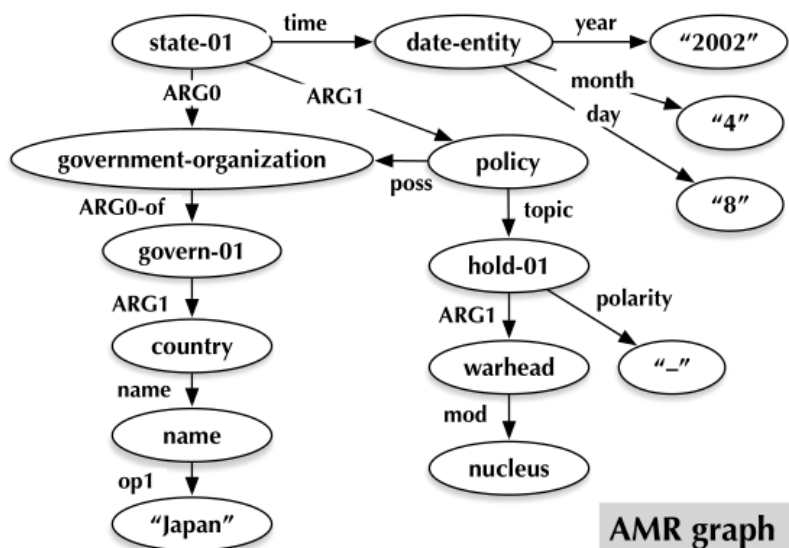


Figure 3.2: Example of an AMR graph for the sentence “The Japanese Government stated on April 8, 2002 its policy of holding no nuclear warheads”. Image taken from Liao et al. (2018).

Lebanoff, Muchovej, et al. (2019) proposed testing the effectiveness of sentence fusion in abstractive summarization settings. To this end they found that current state-of-the-art systems are not doing (either implicitly or explicitly) sentence fusion, ending with the note that sentence fusion is not as effective as previously thought. Evaluation was done using the CNN/Dailymail corpus by human evaluators, where the assessed metrics were a) faithfulness: if the summary remains true to the original text, b) grammaticality: if the summary is grammatically acceptable, and c) coverage: if the summary has information pertaining to selected article highlights. They concluded that the systems that perform fusion rank the lowest on faithfulness and that higher ROUGE scores do not necessarily lead to more faithful summaries. Within the systems that perform fusion, they found that the systems that fuse sentences by simple concatenation are the ones that have the highest faithfulness. In the other metrics, fusion systems are found to generate readable, grammatically correct summaries, but not as much as the state-of-the-art encoder-decoder systems.

Lebanoff et al. (2018) trained an encoder-decoder model to learn how to fuse disparate sentences to generate the summary in an abstractive manner, with an attention mechanism to regulate which sentences to fuse. MMR is also used to calibrate the selection, to account for redundancy in the summary. They evaluate their performance using ROUGE-1, ROUGE-2, ROUGE-SU4 and human judgments on documents from DUC 2004 and TAC 2011, comparing it with the baselines used by Liao et al. (2018) plus an integer linear-programming model for summarization (Gillick et al., 2009). The implications of the mixture of extractive and abstractive summarization is seen as a point to develop further, since, despite not performing as well as some extractive baselines, the summaries generated were more highly ranked by human annotators in faithfulness and coverage. One advanced possibility is that they are not optimizing the extractive part of the method separately from the abstractive one.

Lebanoff, Song, et al. (2019) used vector representations for words that change depending on the context the word is inserted in (see section 2.1.2). With this representation, they created a system that compresses or joins two sentences – hence it is an abstractive procedure – selected by an attention mechanism that is sensible to the fact that the sentences may have come from different documents. Tests were done both in SDS and MDS settings, and, for the used MDS corpus (DUC 2004), the baseline extractive methods (see section 4.4.1, plus  $N$ -LEAD, where the first  $N$  sentences are taken from each document to form the summary — in their case,

$N$  is the average number of sentences in the reference summaries) fared better in ROUGE-1, ROUGE-2, and ROUGE-L scores, indicating that just having the awareness that sentences come from different documents is not enough to properly select important content. Also, for the MDS setting, they find that simpler representations that leverage word frequency across documents (TF-IDF, see section 2.1.1) outperform more complex ones.

### 3.4 Summary

In this chapter we have given an overview of the related work on multi-document summarization. We started by introducing the main corpora used in the MDS task, DUC 2006 and DUC 2007, and the main metrics used for automatic or manual evaluation of the generated summaries, namely ROUGE and Pyramid. We then talked about some elementary or state-of-the-art methods in the MDS task, providing a short summary of each method, as well as the corpora and metrics used.



# 4 Summarization Experiments

This chapter details the datasets, metrics, and steps taken for the summarization algorithm. We also present and analyse the results of our experiments.

## 4.1 *Datasets*

We chose to test simplicial curves in the MDS problem using the DUC 2006 and DUC 2007 datasets (c.f. section 3.1). These datasets were chosen since they are the ones most suitable for MDS, as well as being the ones where most MDS articles have focused on, facilitating comparison of results. Our objective is, given a document cluster, to produce a summary (where summary size varies with the chosen dataset) that is acceptable (measured by some metric — see below) when compared to some human-made reference summary for that document cluster.

## 4.2 *Curve Construction*

Every DUC document is first converted from XML to plain text. We chose to construct the curves at three levels: a) the sentence level (each curve represents a sentence); b) the document level (each curve represents a document); and c) the mixed level (each curve represents a document, built by concatenating smaller curves that represent the sentences of that document). Sentences are extracted from every document using Apache OpenNLP<sup>1</sup>. No stopwords were removed and no stemming was done, in order to preserve function words in the generated summaries.

The base matrices were created in two ways: a) 100-dimension vectors from word2Vec (viz. 2.2.3.2); and b) 10-dimension vectors, obtained by applying UMAP (McInnes et al., 2018) to the 100-dimension vectors. UMAP is a dimensionality reduction technique that maintains positional relativity: objects close in the high-dimensional space are mapped to close objects in

---

<sup>1</sup><https://opennlp.apache.org/> (Accessed January 20, 2021)

the low-dimensional space and objects further apart in the high-dimensional space are mapped to distant objects in the low-dimensional space..

Curves were built with a smoothing value for the restricted Gaussian kernel of  $\sigma \in \{0.05; 0.005; 0.003; 0.001\}$ . These values were chosen in-line with the original article. Finding a clear relation between smoothing value and curve performance can be a future area of enquiry. UMAP was run with the default parameters from the authors' implementation.

### 4.3 Summary Construction

The most straightforward approach to building summaries using curves is the average curve approach. The resulting curve should intuitively model both documents: to get a summary, it suffices to synthesize words from it.

The success of this method is highly dependent on the underlying representation used to build  $M_y$ . To summarize a single document, we can also consider curves built with different scales for the kernel, i.e., synthesizing words from the curve  $\gamma'_y = \frac{1}{2}(\gamma_y^{\sigma_1} + \gamma_y^{\sigma_2})$  for document  $y$ .

We construct the summary curve by a) averaging and b) conflating the curves for all the documents. Reconstructing the text is done by sampling uniform-spaced points from the summary curve and retrieving the word associated with the dimension with the highest probability (in the case of one-hot built curves), or by training an encoder-decoder model to map between curves to sentences.

An LSTM (Hochreiter & Schmidhuber, 1997) was used to create a mapping from curves to sentences in the cases where the curve dimensions do not have any extrinsic meaning. This was done by training a neural-network to match curve representations (dense matrices) of sentences to a vector representation of those sentences. This vector representation has length equal to the length of the sentence, and the entries of the vector are the index positions in the vocabulary of the word in the sentence.

As an example, consider the sentence "How are you, you villain?". If we create a curve from a base matrix representation with 10-dimension vectors, and we sample 5 points of that curve to

create a summary, the network will have to map:

$$\mathbb{R}^{5 \times 10} \ni \begin{bmatrix} 0.1 & \cdots & 0.15 \\ \vdots & \ddots & \vdots \\ 0.45 & \cdots & 0.03 \end{bmatrix} \mapsto [0 \ 1 \ 2 \ 2 \ 3]$$

Since any unlabelled collection of texts can be used for this purpose, we trained the model in the DUC 2006 dataset, using only the original documents as the source for our training sentences. The model was constructed using the Keras Framework (Chollet et al., 2015) with the Tensorflow backend. It was trained for 50 iterations using the Sparse Categorical Cross-entropy loss and the RMSprop optimizer. The hidden layer of the LSTM had dimension 100 and its activation function was softmax.

## 4.4 Evaluation

We compare the generated summaries with the reference summaries using ROUGE-1, ROUGE-2, and ROUGE-L, since these are the most widely used automatic metrics in the MDS task. Although it was also presented in section 3.2, we do not evaluate our summaries using SERA because this metric has relevance chiefly in the field of summarizing scientific articles, and not general news articles.

### 4.4.1 Baselines

We compare simplicial curves with some strong extractive (*ext*) and abstractive (*abs*) baselines that have been applied successfully in multi-document summarization:

- SumBasic (*ext*) (Vanderwende et al., 2007) is a greedy algorithm for sentence selection that chooses to include in the summary the sentence with the highest probability, as given by  $P(S) = \frac{1}{\#S} \sum_{w \in S} P(w)$ , where  $P(w)$  is a unigram distribution of all the words in the corpus. This process is repeated until the desired length of the summary is reached.
- KLSum (*ext*) (Haghighi & Vanderwende, 2009) builds upon the above idea but, instead of including in the summary the sentence with the highest probability, it selects the sentence

that, when added to the summary, most reduces the KL-divergence between the unigram distribution of the sentence and the unigram distribution of the corpus.

- TextRank (*ext*) (Mihalcea & Tarau, 2004) is very similar to LexRank (viz. section 3.3) but, instead of finding the eigenvalue centrality of the sentence graph, sentences are ranked for extraction by the PageRank algorithm. Further, the original TextRank algorithm builds the weighted sentence graph by considering the weight of an edge to be the amount of lexical overlap between two sentences.
- Pointer-Generator (PG) networks (*abs*) (See et al., 2017) is a mixture of extractive and abstractive approaches: when constructing the summary, the model decides (based on an attention mechanism) if, for the current position of the under-construction summary, it is better to generate a new word or to copy a word from the source text. Both the attention mechanism and the underlying model need to be trained using a different corpus from the one we want to summarize.
- PG-MMR (*abs*) (Lebanoff et al., 2018) builds upon the above idea, only that the summary construction step is interleaved with MMR (section 3.3), where it is used to pick  $K$  sentences to pass on to PG. After each summary construction round, the sentences are re-ranked and the process repeats until the summary has the desired length.

## 4.5 Results and Discussion

The ROUGE-1, ROUGE-2 and ROUGE-L scores for each baseline and the proposed method are presented in Tables 4.1 and 4.2, for the DUC 2006 and DUC 2007 datasets. In the simplicial curves entry is the dimension of the word-embeddings used (10 or 100). “Cat” means that the curve was done at the document level by concatenating curves at the sentence level. The results presented are for curves built with smoothing kernel  $\sigma = 0.003$ .

All presented curve results were obtained by combining curves by average. The conflation method, as discussed in section 2.3, is not shown as it was not successful: the resulting curve would always concentrate all of its mass around one point, making the curve generate only one word.

We can see that the ROUGE scores achieved for the simplicial curves are not competitive

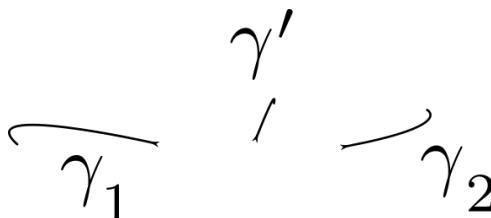
| DUC 2006                  | R-1         | R-2         | R-L         |
|---------------------------|-------------|-------------|-------------|
| SumBasic                  | 0.27        | 0.03        | 0.12        |
| KLSum                     | 0.27        | 0.03        | 0.13        |
| TextRank                  | <b>0.33</b> | <b>0.06</b> | <b>0.16</b> |
| PG                        | 0.24        | 0.04        | 0.13        |
| PG-MMR                    | 0.30        | <b>0.06</b> | 0.15        |
| Simplicial Curves 10      | 0.02        | 0.001       | 0.02        |
| Simplicial Curves 100     | 0.04        | 0.001       | 0.04        |
| Simplicial Curves 100 Cat | 0.05        | 0.004       | 0.04        |

Table 4.1: Baseline and curve results on the DUC 2006 dataset.

| DUC 2007                  | R-1         | R-2         | R-L         |
|---------------------------|-------------|-------------|-------------|
| SumBasic                  | 0.29        | 0.04        | 0.14        |
| KLSum                     | 0.28        | 0.04        | 0.13        |
| TextRank                  | <b>0.36</b> | <b>0.07</b> | <b>0.17</b> |
| PG                        | 0.26        | 0.05        | 0.14        |
| PG-MMR                    | 0.32        | <b>0.07</b> | <b>0.17</b> |
| Simplicial Curves 10      | 0.02        | 0.001       | 0.01        |
| Simplicial Curves 100     | 0.03        | 0.001       | 0.03        |
| Simplicial Curves 100 Cat | 0.04        | 0.003       | 0.03        |

Table 4.2: Baseline and curve results on the DUC 2007 dataset.

with the baselines. One possible explanation for this is that the curve-averaging method generates curves that are poor information-wise, due to the original curves' distance apart in the 100-dimension word-embedding space (as illustrated in Figure 4.1, with  $\gamma' = \frac{1}{2}(\gamma_1 + \gamma_2)$ ). This, combined with the fact that the curves pass through the same region many times (because of the stopwords), leads to the resulting average curve being condensed in some specific areas and every so often shooting off into regions with content.

Figure 4.1: Shortcoming of the curve-averaging method.  $\gamma'$  has no clear relation with  $\gamma_1$  or  $\gamma_2$ .

We keep the stopwords because we need the function words to generate legible text. However,

even if we remove the stopwords, the overall result does not change: the resulting curve now focuses on superfluous words between the curves.

One thing to note is the fact that the average curve may pass through regions that do not represent the original documents in any way. To use a standard word-embedding example, if there are two documents that say “hot” and “cold”, respectively, then the average curve will pass through the “warm” region, and this would be the word that would be generated for the summary, even though the original documents have no connection to this term. This is also illustrated by Figure 4.1.

Even so, the ROUGE scores are better if the curves are built with higher dimension base matrices. This is despite the fact that, in higher dimensions, the curves have a much wider space to roam, thus aggravating the above-mentioned shortcoming of the averaging method. We can explain this dissonance in two ways: a) the 10-dimension word vectors were obtained by reducing the dimension of the 100-dimension word vectors with UMAP. This mapping may have rendered the output vectors unfit for purpose, even though UMAP retains object proximity relativity (close objects in higher dimensions remain close in lower dimensions; far away objects in higher dimensions remain far away in lower dimensions); and b) higher dimensions in the word-embedding space can model a higher number of concepts, so the generated words will be richer (this is compatible with the above shortcoming: the curves to average have more dimensions, so the average curve has more regions where it can pass through and not be informative).

It is also interesting to note that constructing document-level curves by concatenating sentence-level curves works better than constructing document-level curves or sentence-level curves by themselves. This is due to the resulting curves being bigger (more information dense) because each part of the curve explicitly models a sentence of the originating document.

Some examples of generated sentences from average curves with kernel  $\sigma$  are:

$\sigma = 0.001$  Think the other than the most common and the most common york  
 $\sigma = 0.003$  The New York Times news service the first time and the other states  
 $\sigma = 0.005$  The New York city and the New York city [*repeated*]  
 $\sigma = 0.05$  The officials said

These sentences come from curves generated from documents belonging to a document-set

where New York appears in the context of a smoking ban in restaurants, air pilot investigations and murder rates. In the generated sentences, only the main topic linking these documents is of interest, with the other words being either low content or function words. If the smoothing is too much ( $\sigma = 0.05$ ), the curve “flattens” and ceases to pass through information dense regions (explaining the small quantity of words generated), whilst if the smoothing is too little ( $\sigma = 0.001$ ) the simplicial curves method degenerates into a simple bag-of-words procedure, where words are sampled with no regard for their position in a sentence.

Given that the ROUGE results were so poor, we did not use the algebraic machinery developed in section 2.3. Our original plan was to relate an intrinsic curve feature (e.g., curve entropy) with the quality of the generated summary, as given by the ROUGE scores. However, because these were so low, any correlation score was highly probable to be just noise.

## 4.6 Summary

We have presented our proposed experiments to evaluate the simplicial curves method in the MDS task. We detailed the corpora chosen, how we evaluate the generated summaries, how we construct the curves from which we generate summaries, and a set of five baseline methods to which we compare our method to. In the end, we have analysed the results of our simplicial curves experiments. Although the results are poor metric-wise, we can find a number of explanations for this, prompting a new path forward in this line of enquiry.





# 5 Conclusion and Future Work

In this work we have expanded on the concept of simplicial curves (Lebanon et al., 2007), generalizing it to different base representations. We introduced the simplicial curves method and developed an algebra for it, which we did not use in its entirety.

We then explored the effectiveness of the method in the task of multi-document summarization, testing whether or not a representation that explicitly maintains sequencing information is useful in this task. Our experiments show that it is not, with the ROUGE-1 score of 0.04 obtained in the DUC 2007 dataset being below our simplest baseline (SumBasic) score of 0.29, where summaries are constructed simply by selecting the most important sentences in a document. It thus seems that the additional structure provided by the simplicial curves is not being used effectively in generating summaries – essentially resulting in noise.

Some of the tools we had prepared to deal with curves were left unused because of the poor results we had in the MDS task. If the results had been better – at least as good as the most basic baseline – any impact that any upstream change could have had, as guided by the intrinsic evaluation methods we developed, could at least be reliably measured.

The cause of these ROUGE results is made clear when we present the actual generated summaries: the words selected by the curves are low-content and have little relevance across documents, primarily due to the effect mentioned in Figure 4.1.

Notwithstanding the poor results in the summarization task, we still believe in the potential use for a representation that explicitly encodes sequential information (like, e.g., ELMo (Peters et al., 2018)) and that can be manipulated using standard and advanced tools of mathematics (unlike ELMo). In particular, we would like to explore ways of using curves as an intermediate representation for some domain-specialized downstream algorithms, i.e., either sample the curve for points to use as input or use the curve itself as input.

We would also like to keep exploring the effectiveness of curves in different language-related

tasks such as word segmentation or topic modeling, as well as some non-language-related tasks that would nonetheless benefit from having a method for representing objects while preserving sequential information (e.g., video processing).

## References

- Banarescu, L., Bonial, C., Cai, S., Georgescu, M., Griffitt, K., Hermjakob, U., ... Schneider, N. (2013). Abstract Meaning Representation for Sembanking. In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse* (pp. 178–186).
- Baroni, M., Dinu, G., & Kruszewski, G. (2014, 06). Don't Count, Predict! A Systematic Comparison of Context-counting vs. Context-predicting Semantic Vectors. *52nd Annual Meeting of the Association for Computational Linguistics, ACL 2014 - Proceedings of the Conference, 1*, 238-247. doi: 10.3115/v1/P14-1023
- Bengio, Y., Ducharme, R., & Vincent, P. (2000, 01). A Neural Probabilistic Language Model. *Journal of Machine Learning Research, 3*, 932-938. doi: 10.1162/153244303322533223
- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent Dirichlet Allocation. *Journal of Machine Learning Research, 3*(Jan), 993–1022.
- Bojanowski, P., Grave, E., Joulin, A., & Mikolov, T. (2017). Enriching Word Vectors with Subword Information. *Transactions of the Association for Computational Linguistics, 5*, 135-146. Retrieved from [https://doi.org/10.1162/tacl\\_a-00051](https://doi.org/10.1162/tacl_a-00051) doi: 10.1162/tacl\\_a\\_00051
- Carbonell, J. G., & Goldstein, J. (1998). The Use of MMR, Diversity-based Reranking for Reordering Documents and Producing Summaries. In *SIGIR* (Vol. 98, pp. 335–336).
- Cardoso, P. C., Maziero, E. G., Jorge, M. L., Seno, E. M., Di Felippo, A., Rino, L. H., ... Pardo, T. A. (2011). CSTnews – A Discourse-annotated Corpus for Single and Multi-document Summarization of News Texts in Brazilian Portuguese. In *Proceedings of the 3rd RST Brazilian Meeting* (pp. 88–105).
- Chollet, F., et al. (2015). *Keras*. <https://keras.io>. (Accessed November 20, 2020)
- Cohan, A., & Goharian, N. (2016, apr). Revisiting Summarization Evaluation for Scientific Articles. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation*

- (LREC 2016) (pp. 1144–1147). Portorož, Slovenia: European Language Resources Association (ELRA). Retrieved 2018-10-23, from <http://www.lrec-conf.org/proceedings/lrec2016/summaries/1144.html> (arXiv: 1604.00400)
- Collobert, R., & Weston, J. (2008). A Unified Architecture for Natural Language Processing: Deep Neural Networks with Multitask Learning. In *Proceedings of the 25th International Conference on Machine Learning* (pp. 160–167).
- Deerwester, S., Dumais, S., Furnas, G., Landauer, T., & Harshman, R. (1990, 09). Indexing by Latent Semantic Analysis. *Journal of the American Society for Information Science*, 41, 391-407. doi: 10.1002/(SICI)1097-4571(199009)41:6<391::AID-ASI1>3.0.CO;2-9
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)* (pp. 4171–4186).
- Erkan, G., & Radev, D. R. (2004). Lexrank: Graph-based Lexical Centrality as Salience in Text Summarization. *Journal of Artificial Intelligence Research*, 22, 457–479.
- Ethayarajh, K., Duvenaud, D., & Hirst, G. (2019, July). Towards Understanding Linear Word Analogies. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics* (pp. 3253–3262). Florence, Italy: Association for Computational Linguistics. Retrieved from <https://www.aclweb.org/anthology/P19-1315> doi: 10.18653/v1/P19-1315
- Fabbri, A., Li, I., She, T., Li, S., & Radev, D. (2019, July). Multi-News: A Large-Scale Multi-Document Summarization Dataset and Abstractive Hierarchical Model. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics* (pp. 1074–1084). Florence, Italy: Association for Computational Linguistics. Retrieved from <https://www.aclweb.org/anthology/P19-1102> doi: 10.18653/v1/P19-1102
- Ganesan, K., Zhai, C., & Han, J. (2010). Opinosis: A Graph-based Approach to Abstractive Summarization of Highly Redundant Opinions. In *Proceedings of the 23rd International Conference on Computational Linguistics (Coling 2010)* (pp. 340–348).
- Gillick, D., Favre, B., Hakkani-Tür, D., Bohnet, B., Liu, Y., & Xie, S. (2009). The ICSI/UTD Summarization System at TAC 2009. In *Theory and Application of Categories*.

- Gillick, D., Favre, B., & Hakkani-Tür, D. Z. (2008). The ICSI Summarization System at TAC 2008. In *Theory and Applications of Categories*.
- Goldstein, J., Mittal, V., Carbonell, J., & Kantrowitz, M. (2000). Multi-document Summarization by Sentence Extraction. In *Proceedings of the 2000 NAACL-ANLP Workshop on Automatic Summarization* (pp. 40–48).
- Haghighi, A., & Vanderwende, L. (2009, June). Exploring Content Models for Multi-Document Summarization. In *Proceedings of Human Language Technologies: The 2009 Annual Conference of the North American Chapter of the Association for Computational Linguistics* (pp. 362–370). Boulder, Colorado: Association for Computational Linguistics. Retrieved from <https://www.aclweb.org/anthology/N09-1041>
- Harman, D., & Liberman, M. (1993). TIPSTER. <https://catalog.ldc.upenn.edu/LDC93T3A>. (Accessed December 28, 2019)
- Hermann, K. M., Kocisky, T., Grefenstette, E., Espeholt, L., Kay, W., Suleyman, M., & Blunsom, P. (2015). Teaching Machines to Read and Comprehend. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, & R. Garnett (Eds.), *Advances in Neural Information Processing Systems 28* (pp. 1693–1701). Curran Associates, Inc. Retrieved from <http://papers.nips.cc/paper/5945-teaching-machines-to-read-and-comprehend.pdf>
- Hill, T. P. (2011). Conflations of Probability Distributions. *Transactions of the American Mathematical Society*, 363(6), 3351–3372. Retrieved from <http://www.jstor.org/stable/23032795>
- Hochreiter, S., & Schmidhuber, J. (1997, November). Long Short-Term Memory. *Neural Comput.*, 9(8), 1735–1780. Retrieved from <https://dl.acm.org/doi/10.1162/neco.1997.9.8.1735> doi: 10.1162/neco.1997.9.8.1735
- Jones, K. S. (1972). A Statistical Interpretation of Term Specificity and its Application in Retrieval. *Journal of Documentation*, 28, 11–21.
- Kågebäck, M., Mogren, O., Tahmasebi, N., & Dubhashi, D. (2014). Extractive Summarization Using Continuous Vector Space Models. In *Proceedings of the 2nd Workshop on Continuous Vector Space Models and their Compositionality (CVSC)* (pp. 31–39).

- Le, Q., & Mikolov, T. (2014). Distributed representations of sentences and documents. In *International Conference on Machine Learning* (pp. 1188–1196).
- Lebanoff, L., Muchovej, J., Deroncourt, F., Kim, D. S., Kim, S., Chang, W., & Liu, F. (2019). Analyzing Sentence Fusion in Abstractive Summarization. In *Proceedings of the 2nd workshop on new frontiers in summarization* (pp. 104–110).
- Lebanoff, L., Song, K., Deroncourt, F., Kim, D. S., Kim, S., Chang, W., & Liu, F. (2019, July). Scoring Sentence Singletons and Pairs for Abstractive Summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics* (pp. 2175–2189). Florence, Italy: Association for Computational Linguistics. Retrieved from <https://www.aclweb.org/anthology/P19-1209> doi: 10.18653/v1/P19-1209
- Lebanoff, L., Song, K., & Liu, F. (2018). Adapting the Neural Encoder-Decoder Framework from Single to Multi-Document Summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing* (pp. 4131–4141).
- Lebanon, G., Mao, Y., & Dillon, J. V. (2007). The Locally Weighted Bag of Words Framework for Document Representation. *Journal of Machine Learning Research*, 8, 2405-2441.
- Levy, O., & Goldberg, Y. (2014, 01). Neural Word Embedding as Implicit Matrix Factorization. *Advances in Neural Information Processing Systems*, 3, 2177-2185.
- Levy, O., Goldberg, Y., & Dagan, I. (2015, 05). Improving Distributional Similarity with Lessons Learned from Word Embeddings. *Transactions of the Association for Computational Linguistics* (2015), 3(1), 211–225. Retrieved from <https://transacl.org/ojs/index.php/tacl/article/view/570>
- Lewis, D. D., Yang, Y., Rose, T. G., & Li, F. (2004). Rcv1: A New Benchmark Collection for Text Categorization Research. *Journal of Machine Learning Research*, 5(Apr), 361–397.
- Liao, K., Lebanoff, L., & Liu, F. (2018). Abstract Meaning Representation for Multi-Document Summarization. In *Proceedings of the 27th International Conference on Computational Linguistics* (pp. 1178–1190).
- Lin, C.-Y. (2004). Rouge: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out* (pp. 74–81).

- Lin, C.-Y., & Och, F. J. (2004). Automatic Evaluation of Machine Translation Quality Using Longest Common Subsequence and Skip-bigram Statistics. In *Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics* (p. 605).
- Loper, E., & Bird, S. (2002). NLTK: The Natural Language Toolkit. In *Proceedings of the ACL-02 Workshop on Effective Tools and Methodologies for Teaching Natural Language Processing and Computational Linguistics - Volume 1* (pp. 63–70). Stroudsburg, PA, USA: Association for Computational Linguistics. Retrieved from <https://doi.org/10.3115/1118108.1118117> doi: 10.3115/1118108.1118117
- Luhn, H. P. (1957, 10). A Statistical Approach to Mechanized Encoding and Searching of Literary Information. *IBM Journal of Research and Development*, 1(4), 309-317. doi: 10.1147/rd.14.0309
- Mani, K., Verma, I., Meisheri, H., & Dey, L. (2018). Multi-document Summarization Using Distributed Bag-of-words Model. In *2018 IEEE/WIC/ACM International Conference on Web Intelligence (WI)* (pp. 672–675).
- McInnes, L., Healy, J., Saul, N., & Grossberger, L. (2018, 09). UMAP: Uniform Manifold Approximation and Projection. *Journal of Open Source Software*, 3, 861. doi: 10.21105/joss.00861
- Mihalcea, R., & Tarau, P. (2004, 07). TextRank: Bringing Order into Text. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing* (pp. 404–411). Barcelona, Spain: Association for Computational Linguistics. Retrieved from <https://www.aclweb.org/anthology/W04-3252>
- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient Estimation of Word Representations in Vector Space. *CoRR*, abs/1301.3781. Retrieved from <http://dblp.uni-trier.de/db/journals/corr/corr1301.html#abs-1301-3781>
- Mnih, A., & Hinton, G. (2007). Three New Graphical Models for Statistical Language Modelling. In *Proceedings of the 24th International Conference on Machine Learning* (pp. 641–648).
- Nayeem, M. T., Fuad, T. A., & Chali, Y. (2018). Abstractive Unsupervised Multi-document Summarization Using Paraphrastic Sentence Fusion. In *Proceedings of the 27th International Conference on Computational Linguistics* (pp. 1191–1204).

- Nenkova, A., Passonneau, R., & McKeown, K. (2007). The Pyramid Method: Incorporating Human Content Selection Variation in Summarization Evaluation. *ACM Transactions on Speech and Language Processing (TSLP)*, 4(2), 4.
- Ng, J.-P., & Abrecht, V. (2015, 09). Better Summarization Evaluation with Word Embeddings for ROUGE. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing* (pp. 1925–1930). Lisbon, Portugal: Association for Computational Linguistics. Retrieved from <https://www.aclweb.org/anthology/D15-1222> doi: 10.18653/v1/D15-1222
- Papineni, K., Roukos, S., Ward, T., & Zhu, W.-J. (2002). BLEU: A Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics* (pp. 311–318).
- Pennington, J., Socher, R., & Manning, C. (2014, 01). Glove: Global Vectors for Word Representation. *EMNLP*, 14, 1532-1543. doi: 10.3115/v1/D14-1162
- Peters, M., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., & Zettlemoyer, L. (2018). Deep Contextualized Word Representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)* (pp. 2227–2237).
- Peyrard, M., Botschen, T., & Gurevych, I. (2017). Learning to Score System Summaries for Better Content Selection Evaluation. In *Proceedings of the Workshop on New Frontiers in Summarization* (pp. 74–84).
- Porter, M. F. (1980). An Algorithm for Suffix Stripping. *Program*, 14(3), 130–137.
- Radev, D. R., Hovy, E., & McKeown, K. (2002). Introduction to the Special issue on Summarization. *Computational Linguistics*, 28(4), 399–408.
- Ribaldo, R., Akabane, A. T., Rino, L. H. M., & Pardo, T. A. S. (2012). Graph-based methods for Multi-document Summarization: Exploring Relationship Maps, Complex Networks and Discourse Information. In H. Caseli, A. Villavicencio, A. Teixeira, & F. Perdigão (Eds.), *Computational Processing of the Portuguese Language* (pp. 260–271). Berlin, Heidelberg: Springer Berlin Heidelberg.
- Rioux, C., Hasan, S. A., & Chali, Y. (2014, October). Fear the REAPER: A System for Automatic Multi-Document Summarization with Reinforcement Learning. In *Proceedings of the 2014*



- Conference on Empirical Methods in Natural Language Processing (EMNLP)* (pp. 681–690). Doha, Qatar: Association for Computational Linguistics. Retrieved from <https://www.aclweb.org/anthology/D14-1075> doi: 10.3115/v1/D14-1075
- Ryang, S., & Abekawa, T. (2012). Framework of Automatic Text Summarization using Reinforcement Learning. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning* (pp. 256–265).
- See, A., Liu, P. J., & Manning, C. D. (2017, July). Get To The Point: Summarization with Pointer-Generator Networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (pp. 1073–1083). Vancouver, Canada: Association for Computational Linguistics. Retrieved from <https://www.aclweb.org/anthology/P17-1099> doi: 10.18653/v1/P17-1099
- Socher, R., Huang, E. H., Pennin, J., Manning, C. D., & Ng, A. Y. (2011). Dynamic Pooling and Unfolding Recursive Autoencoders for Paraphrase Detection. In *Advances in Neural Information Processing Systems* (pp. 801–809).
- Tohalino, J. V., & Amancio, D. R. (2018). Extractive Multi-document Summarization Using Multilayer Networks. *Physica A: Statistical Mechanics and its Applications*, 503, 526–539.
- Turian, J., Ratinov, L., & Bengio, Y. (2010). Word Representations: A Simple and General Method for Semi-supervised Learning. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics* (pp. 384–394). Association for Computational Linguistics. Retrieved from <http://dl.acm.org/citation.cfm?id=1858681.1858721>
- Vanderwende, L., Suzuki, H., Brockett, C., & Nenkova, A. (2007). Beyond SumBasic: Task-focused Summarization with Sentence Simplification and Lexical Expansion. *Information Processing & Management*, 43(6), 1606–1618.
- Vorhees, E., & Graff, D. (2008). AQUAINT-2. <https://catalog.ldc.upenn.edu/LDC2008T25>. (Accessed December 28, 2019)
- Yogatama, D., Liu, F., & Smith, N. A. (2015). Extractive Summarization by Maximizing Semantic Volume. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing* (pp. 1961–1966). Lisbon, Portugal: Association for Computational Linguistics. Retrieved from <http://aclweb.org/anthology/D15-1228> doi: 10.18653/v1/D15-1228

- Zhao, L., Wu, L., & Huang, X. (2009). Using Query Expansion in Graph-based Approach for Query-focused Multi-document Summarization. *Information Processing & Management*, 45(1), 35 - 41. Retrieved from <http://www.sciencedirect.com/science/article/pii/S030645730800071X> doi: <https://doi.org/10.1016/j.ipm.2008.07.001>
- Zhao, W., Peyrard, M., Liu, F., Gao, Y., Meyer, C. M., & Eger, S. (2019). MoverScore: Text Generation Evaluating with Contextualized Embeddings and Earth Mover Distance. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)* (pp. 563–578).