

# Fathers and Sons

Sebastiano Caccaro  
`sebastiano.caccaro@studenti.unimi.it`

Università degli Studi di Milano

## 1 Introduction

The aim of this project is compare the speeches of political leaders in the 1<sup>st</sup>, 12<sup>th</sup>, 17<sup>th</sup> and 18<sup>th</sup> legislature, in order to find out which political parties from different legislature speak in a more similar way.

The proposed solution consists in creating a model of speech for every political party by extracting the main topics from the speeches. The comparison between political party speeches is therefore performed using their topics. As such, two models are considered similar if their topics are similar.

Such topic extraction is performed with LDA (Latent Dirichlet Allocation), proposed in [1]. LDA is a generative statistical model that can be used on text corpora. In LDA, each document is described by a mixture of topic probabilities; in turn, each topic can be represented as set of word probabilities.

## 2 Research question and methodology

The goal of the project is to analyze the transcriptions of parliamentary debates of the 1<sup>st</sup>, 12<sup>th</sup>, 17<sup>th</sup> and 18<sup>th</sup> legislature in order to find out which present political party's speeches resemble the most the speeches of ones in the past.

The problem can be deconstructed in the following steps:

- Text preprocessing
- Language modeling
- Models comparison

### 2.1 Text preprocessing

For each debate transcription, text needs to be divided into speeches. Parts of the document which are not speeches are discarded. Every speech gets tagged with its speaker name. This process results in every speaker having a collection of speeches. Every speech is then divided into tokens.

Data in the 1<sup>st</sup> and 12<sup>th</sup> legislature is particularly noisy (see subsection 3.1). Therefore an additional correction phase has been deemed necessary.

For every legislature, every token (that is, a word) is checked against a vocabulary. This process eliminates words that cannot be corrected and weird OCR artifacts, especially in the 1<sup>st</sup> and 12<sup>th</sup> legislature; stopwords are also filtered out.

In the end, every token is stemmed. Therefore words that have the same root and thus have same meaning are effectively mapped into the same word. This is beneficial both in modeling topics and in classifying documents. Moreover, this approach will reduce the effective number of words, leading to a smaller document-word frequency matrix.

## 2.2 Language modeling

Each deputy's set of speeches is then vectorized in its own sparse document-word frequency matrix. The matrix is used to create a language model using LDA (Latent Dirichlet Allocation). For each model a perplexity score is calculated. The output of this phase, for each model, is:

- A list of topics: each topic is a set of the words in paired with their probability of being in the topic.
- For every topic, the probability of any document to be about a certain topic.

## 2.3 Models comparison

Model similarity is estimated taking into account the topic-word probability part of the previous output. Model similarity between two models  $M_1$  and  $M_2$  is calculated as follows:

- For every couple of topics in  $M_1$  and  $M_2$  a similarity score is calculated. Each score is put in a difference matrix where rows represent topics in  $M_1$  and columns topics in  $M_2$ .
- Recursively the highest value of the matrix is summed to the final result, and the corresponding row and column are removed from the matrix.

The similarity score between two topics  $T_1$  and  $T_2$  is computed as follows:

- Each topic  $T_x$  gets vectorized: the  $i^{th}$  cell in the vector contains the probability that the  $i^{th}$  word in the set of words of  $T_1 \cup T_2$  belongs to  $T_x$ .
- The similarity score corresponds to the cosine distance between the two topic vectors.

Results plotted in a table show three different "blocks" inside which similarity is greater than outside; these blocks are formed as follows:

1. 1<sup>st</sup> legislature
2. 12<sup>th</sup> legislature
3. 17<sup>th</sup> and 18<sup>th</sup> legislature

This is expected and it is almost certainly caused by the years separating the legislatures. Therefore, for every couple of blocks  $B_1$  and  $B_2$ , the following normalization steps are applied:

- The combined average  $avg_b$  of all the cells in  $B_1$  and  $B_2$  is computed. The cells on the diagonal (the ones in which the same party models is confronted with itself) are not taken into account.

- The section of the table  $S_1$  and  $S_2$  are defined as the group of cells sitting in the same columns as  $B_1$  and same rows as  $B_2$  and vice versa. Since the table is symmetrical the average  $avg_s$  of their cells is the same.
- The absolute value of  $(avg_b - avg_s)$  is summed to all cells in  $S_1$  and  $S_2$ .

### 3 Experimental results

#### 3.1 Dataset

The dataset consists in the full transcriptions of parliamentary debates in the Italian House Of Representative of four different legislatures:

- 1<sup>st</sup> legislature
- 12<sup>th</sup> legislature
- 17<sup>th</sup> legislature
- 18<sup>th</sup> legislature

For every legislature the text is extracted from the PDF of the transcripts themselves. In the case of the 17<sup>th</sup> and 18<sup>th</sup> legislature, the PDF transcripts are native digital files. Therefore characters are encoded in UTF8 and the extracted text has been basically copy-pasted into the dataset.

In the case of the 1<sup>st</sup> and 12<sup>th</sup> legislature the extraction has been performed through OCR (Optical Character Recognition). This leads to many error and artifacts being present in the extracted text, such as misspelled, missing or fragmented words, and wrong or missing punctuation. As such, on this two legislatures a correction step is performed to try to correct some of the errors. Errors are detected by checking a word against a preexistent dictionary.

Since OCR errors are quite particular in the types of additions or deletions generated, standard distance-based algorithms for finding correction candidates do not work well. Therefore, an OCR-specific algorithm called TICCL [2] is used. To choose between generated candidates, a simple heuristic is used: the most frequent candidate in the corpus is the one that get chosen. Of course, a word-frequency list for the corpus needs to be calculated beforehand.

On the last two legislatures such correction step was deemed not necessary, as the amount of errors in the text it minimal.

In every legislature the extracted text also contains unwanted lines such as page headings and footings, titles, results of parliamentary votes and so on. Moreover, the text formatting follows different rules in every legislature: this requires a different parser to be written for every legislature.

#### 3.2 Experimental methodology

In order to evaluate the models performances, the mandated metrics is perplexity. Therefore an attempt has been made to optimize the number of topics parameter through a grid search. For each party, speeches are shuffled and divided in train set and test set in a 85/15 split. The expected behaviour for train

perplexity would be to decrease as the number of topic increases, as stated in [1]. Unfortunately, the grid search showed increasing values of perplexity in both training and test sets. This could be likely caused by bug<sup>1</sup> in the LDA implementation of the used library. Because of this, an arbitrary number of topics of 40 is set.

### 3.3 Results

The models produced by LDA are obtained and put in table as described in previous sections. For each couple of models the score is a float number between 0 (totally different models) and 40 (the models are the same). Parties with less than 500 documents are discarded as results showed to be unreliable and differed too much from others.

	1_DC	1_PC1	1_PSI	1_PBI	12_PD	12_F1	12_Ra	12_AN	12_LE	12_PP	12_Ri	17_FD	17_AI	17_F1	17_LE	17_MS	17_Ar	17_PD	17_Si	17_Je	18_F1	18_FD	18_Le	18_MS	18_PD	18_R	18_LE	18_Mi
1_DC	0.0	13.4	12.57	12.04	15.08	15.32	13.69	15.02	14.62	14.14	15.53	15.86	16.08	15.0	15.96	15.86	16.33	16.02	15.31	15.9	16.51	16.37	15.35	16.38	15.45	15.13	16.26	
1_PC1	13.4	0.0	13.21	16.72	16.49	16.2	13.23	14.45	14.19	14.36	16.03	15.08	16.18	15.12	15.51	14.62	16.47	15.71	16.01	15.07	15.24	15.77	15.88	14.25	16.11	15.4	15.53	15.98
1_PSI	12.57	13.21	0.0	10.96	15.3	14.59	13.84	14.01	13.72	14.51	15.08	15.01	14.51	14.1	14.45	14.18	15.44	14.25	15.15	14.63	15.29	14.61	14.32	12.63	15.33	14.5	14.24	14.92
1_PBI	12.04	16.72	10.96	0.0	13.88	13.9	12.55	13.05	12.7	13.8	14.17	13.76	14.7	13.66	13.84	13.74	14.19	13.6	13.77	13.28	14.05	14.41	14.28	12.5	14.34	13.96	13.53	14.03
12_PD	15.08	16.49	15.3	13.88	0.0	17.54	13.61	17.6	15.61	16.08	17.68	15.36	15.32	15.87	16.32	14.45	16.4	15.22	16.83	14.5	16.01	15.86	15.84	13.65	16.16	15.97	14.78	15.95
12_F1	15.32	16.2	14.59	13.9	17.54	0.0	13.51	17.13	17.15	15.38	15.99	16.36	15.82	16.73	16.44	15.41	16.48	15.97	17.07	14.79	16.39	16.33	16.71	15.82	17.1	16.08	14.88	17.36
12_Ra	13.69	13.23	13.84	12.55	13.61	13.51	0.0	13.2	11.95	12.83	13.19	14.13	13.25	13.63	13.72	12.91	14.27	12.45	14.45	13.12	13.86	14.14	14.34	12.4	14.11	13.95	13.12	14.94
12_AN	15.02	14.65	14.01	13.05	17.6	17.13	13.2	0.0	14.85	13.74	16.47	14.75	15.53	15.48	14.75	14.64	15.2	13.61	14.84	13.3	15.48	15.07	14.81	14.06	14.78	13.73	14.61	14.85
12_LE	14.62	14.19	13.72	12.7	15.61	17.15	11.93	14.85	0.0	14.82	15.23	14.95	14.66	15.06	15.62	14.28	15.67	15.38	15.47	13.66	14.47	14.33	15.26	14.66	13.62	14.47	13.48	15.46
12_PP	14.14	14.36	14.51	13.8	16.08	15.38	12.83	13.74	14.02	0.0	14.34	14.96	14.33	15.35	15.15	14.09	15.78	14.03	16.3	14.38	16.03	14.74	15.83	13.35	15.47	15.3	14.13	15.26
12_Ri	15.53	16.03	15.08	14.17	17.68	15.99	13.19	16.47	15.23	14.34	0.0	15.4	15.47	16.05	16.55	15.01	16.84	15.62	16.62	15.1	15.99	15.63	15.58	14.59	16.45	15.73	15.9	16.1
17_FD	15.86	15.08	15.01	13.76	15.39	16.36	14.13	14.75	14.95	14.96	15.4	0.0	14.38	15.2	16.34	15.94	16.38	15.26	15.51	14.86	16.32	17.27	14.77	12.56	15.33	15.13	13.59	16.47
17_AI	15.96	16.18	14.51	14.7	15.32	15.82	13.25	15.53	14.66	14.33	15.47	14.28	0.0	15.44	13.89	15.39	15.62	16.72	15.22	13.48	15.11	14.33	15.83	14.96	14.41	13.37	14.05	14.8
17_F1	16.08	15.12	14.1	13.66	15.87	16.73	13.63	15.48	15.06	15.35	16.05	15.2	15.44	0.0	14.36	17.59	15.81	16.0	16.3	13.2	15.44	14.8	14.71	14.86	14.57	13.02	13.16	15.18
17_LE	15.0	15.51	14.45	13.84	16.32	16.44	13.72	14.75	15.62	15.15	16.55	16.24	13.89	14.36	0.0	16.88	17.13	15.13	15.38	12.56	15.85	15.28	15.82	12.36	15.85	15.24	13.54	15.5
17_MS	15.96	14.62	14.18	13.74	14.45	15.41	12.91	14.64	14.29	14.09	15.01	15.94	15.39	17.55	16.96	0.0	16.95	14.21	15.71	13.66	16.19	14.68	15.4	15.35	15.75	14.5	13.28	15.31
17_Ar	15.86	16.47	15.44	14.19	16.4	16.48	14.27	15.2	15.87	15.78	16.84	16.38	15.62	15.81	17.12	16.95	0.0	16.54	16.59	15.44	16.88	16.4	16.56	13.13	16.74	16.82	14.83	16.87
17_PD	16.33	15.71	14.25	13.6	15.22	15.97	12.45	13.61	15.38	14.93	15.62	15.26	16.72	16.0	15.13	16.11	16.54	0.0	15.91	14.22	15.12	14.67	17.36	15.84	16.13	14.34	12.78	16.01
17_Si	16.62	16.81	15.35	13.77	16.83	17.07	14.45	14.84	15.47	16.3	16.82	15.51	15.22	16.3	15.38	15.71	16.58	15.81	0.0	14.84	18.05	15.14	15.34	13.47	16.27	15.75	15.29	16.08
17_Je	15.31	15.07	14.63	13.28	14.5	14.79	13.12	13.3	13.66	14.38	15.1	14.88	13.48	13.2	12.56	13.66	14.44	14.22	14.84	0.0	14.93	15.36	13.8	11.22	14.34	14.09	13.4	14.16
18_F1	15.9	15.24	15.29	14.05	16.01	16.39	13.86	15.48	14.47	16.83	15.99	16.32	15.11	15.44	15.85	16.15	16.88	15.12	16.05	14.93	0.0	16.64	16.48	13.54	17.46	16.13	14.2	16.7
18_FD	16.51	15.77	14.61	14.41	15.66	16.33	14.14	15.07	14.35	14.74	15.63	17.23	14.53	14.8	15.26	14.68	16.4	14.67	15.14	15.36	16.64	0.0	14.85	12.49	15.75	16.05	13.93	15.53
18_Le	16.37	15.98	14.32	14.28	15.64	16.71	14.34	14.81	15.26	15.93	15.58	14.77	15.93	14.71	15.02	15.4	16.56	17.36	15.34	13.0	16.49	14.85	0.0	15.11	15.77	14.64	14.16	15.58
18_MS	15.35	14.25	12.63	12.5	13.65	15.82	12.4	14.06	14.66	13.35	14.59	12.56	14.86	14.86	12.36	15.35	13.13	15.94	13.47	11.22	13.54	12.68	15.11	0.0	13.49	13.39	12.5	14.42
18_PD	16.38	16.11	15.33	14.54	16.16	17.1	14.11	14.78	15.62	15.47	16.85	15.33	14.41	14.57	15.87	15.75	16.74	16.13	16.37	14.34	17.48	15.75	15.77	13.49	0.0	16.63	14.5	16.32
18_R	15.45	15.4	14.5	13.86	15.87	16.08	13.95	13.73	14.47	15.3	15.73	15.13	13.37	13.02	15.24	14.5	16.82	14.34	15.75	14.09	16.33	16.05	14.64	11.39	16.43	0.0	13.56	16.4
18_LE	15.13	15.53	14.24	13.53	14.78	14.88	13.12	14.81	13.48	14.13	15.9	13.59	14.05	13.16	13.54	13.28	14.83	12.78	15.29	13.4	14.2	13.93	14.16	12.5	14.5	13.56	0.0	14.33
18_Mi	16.26	15.99	14.82	14.03	15.85	17.36	14.94	14.85	15.46	15.26	16.1	16.47	14.8	15.18	15.5	15.31	16.87	16.01	16.09	14.16	16.7	15.53	15.58	14.42	16.32	16.4	14.33	0.0

**Fig. 1.** Final comparison table. Some row and columns have been removed to better fit the paper. The full table can be found in the project folder.

<sup>1</sup> <https://github.com/scikit-learn/scikit-learn/issues/6777>

## 4 Concluding remarks

The results clearly show some marginally different scores for the parties. It is not sure, though, if some meaningful conclusions could be extracted from table. This considered, it is unlikely the proposed topic-extraction-based approach is the most suited to tackle the discussed problem.

However, some improvements could be made in future work to try to improve this approach's performance:

- More clear-cut distance metrics could be defined for topic and model comparison.
- Due to some problems with the scikit-learn, a proper grid search on the optimal number of topics parameter could not be conducted. Future work could use the gensim implementation of LDA to achieve some meaningful results.
- Due to some hardware and especially memory limitations, the project sticks to an unigram model. In fact, the limited amount of memory of the used PC (8GB), is not enough to work with n-grams models with  $n > 1$ . It may be interesting to try to use bigrams or even trigrams. The code already supports this option, as just a single constant would need to be changed.
- As text correction is not the main focus of this project, a simple heuristic is used to choose from some generated candidates. In future work this mechanism could be improved in order to obtain more accurate results. Moreover, the current candidate-generating algorithm could be further optimized as it is quite slow (and therefore is not used in the executable demo).

## References

1. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. *the Journal of machine Learning research* **3**, 993–1022 (2003)
2. Reynaert, M.: Non-interactive ocr post-correction for giga-scale digitization projects. In: *International Conference on Intelligent Text Processing and Computational Linguistics*. pp. 617–630. Springer (2008)