

Etude comparative de vecteurs pour l'identification du parti politique d'interventions parlementaires

Pauline Degez and Florian Philippe and Valentine Fleith

Address line

...

43017467@parisnanterre.fr

Abstract

This document is a supplement to the general instructions for *ACL authors. It contains instructions for using the \LaTeX style files for ACL conferences. The document itself conforms to its own specifications, and is therefore an example of what your manuscript should look like. These instructions should be used both for papers submitted for review and for final versions of accepted papers.

1 Introduction

La cinquième édition du défi Fouille de Textes (DEFT) porte sur la fouille d'opinions sur des corpus multilingues. Trois tâches ont été proposées, dans trois langues : le français, l'anglais et l'italien. Cet article se concentre sur la 3ème tâche, dont l'objet est l'identification automatique du parti politique d'appartenance de chacun des intervenants dans un corpus de débats parlementaires européens. Il s'agit d'une tâche de classification à 5 classes: Verts-ALE, GUE-NGL, PSE, ELDR et PPE-DE.

Le but de nos expériences sera ainsi de trouver un/des classifieur(s) permettant de réaliser cette tâche. Pour ce faire, nous utiliserons les algorithmes de Machine Learning implementés dans la bibliothèque Python `scikit-learn`.

1.1 Travaux présentés en 2009

Parti	ELDR	GUE-NGL	PPE-DE	PSE	Verts/ALE
F-mesure	0.21	0.37	0.47	0.37	0.25

Table 1: Moyennes des F-mesures par parti politique.

En 2009, un seul participant a soumis un travail pour la tâche 3 ; la Présentation de l'édition 2009¹ évoque, pour expliquer cela, les faibles résultats des logiciels sur cette tâche, bien que conformes à ceux que des humains obtiendraient manuellement.

¹Actes du cinquième défi fouille de texte, DEFT2009, Paris, France, 22 juin 2009

L'équipe de l'Université de Montréal (D. Forest and al.) a obtenu en moyenne les f-mesures présentées dans la Table 1. En moyenne, cela donne donc une f-mesure 0.331.

1.2 Notre approche et travaux antérieurs

Pour ce travail, notre approche a été comparative sur plusieurs niveaux. Tout d'abord, nous comparons différents classifieurs : Random Forest, Régression logistique, Perceptron et Support Vector Machine. De plus, nous testons aussi différentes vectorisations du corpus sur l'ensemble de ces modèles: TF-IDF, Doc2Vec, et des Bert embeddings.

Plusieurs travaux de recherches explorent les comparaisons entre performances des modèles selon les techniques de vectorisation utilisées. Nous pouvons par exemples évoquer ceux de P. Joseph et S. Y. Yerima² en 2022, qui compare les performances des N-grams, TF-IDF, Sac de mots, Word2Vec, Doc2Vec, etc. Leur objectif est de comparer l'impact de la vectorisation sur la précision des modèles. Dans leur article, les modèles doc2vec et TF-IDF démontrent de bons résultats, nous allons ainsi les tester dans notre expérience. Nous décidons d'ajouter à ces deux dernier les embeddings de Bert afin d'avoir trois techniques variées : une méthode statistique, une méthode fondée sur un ANN classique et une sur un Transformer.

2 Méthode

2.1 Dataset

Ici dire que le dataset est nul a chier

²P. Joseph and S. Y. Yerima, "A comparative study of word embedding techniques for SMS spam detection," 2022 14th International Conference on Computational Intelligence and Communication Networks (CICN), Al-Khobar, Saudi Arabia, 2022, pp. 149-155,

Command	Output	Command	Output
<code>{\"a}</code>	ä	<code>{\c c}</code>	ç
<code>{^e}</code>	ê	<code>{\u g}</code>	ğ
<code>{`i}</code>	ì	<code>{\l}</code>	ł
<code>{\ .I}</code>	İ	<code>{\~n}</code>	ñ
<code>{\o}</code>	ø	<code>{\H o}</code>	ő
<code>{\'u}</code>	ú	<code>{\v r}</code>	ř
<code>{\aa}</code>	å	<code>{\ss}</code>	ß

Table 2: Example commands for accented characters, to be used in, e.g., Bib_{TEX} entries.

3 Résultats

3.1 Vecteurs TF-IDF

ca marche trop bien wtf

4 Preamble

The first line of the file must be

```
\documentclass[11pt]{article}
```

To load the style file in the review version:

```
\usepackage[review]{acl}
```

For the final version, omit the review option:

```
\usepackage{acl}
```

To use Times Roman, put the following in the preamble:

```
\usepackage{times}
```

Please see the \LaTeX source of this document for comments on other packages that may be useful.

By default, the box containing the title and author names is set to the minimum of 5 cm. If you need more space, include the following in the preamble:

```
\setlength\titlebox{<dim>}
```

where <dim> is replaced with a length. Do not set this length smaller than 5 cm.

5 Document Body

5.1 Footnotes

5.2 Tables and figures

See Table ?? for an example of a table and its caption. **Do not override the default caption sizes.**

As much as possible, fonts in figures should conform to the document fonts. See Figure ?? for an example of a figure and its caption.



Figure 1: A figure with a caption that runs for more than one line. Example image is usually available through the mwe package without even mentioning it in the preamble.

environment at an appropriate point within the text. The `graphicx` package supports various optional arguments to control the appearance of the figure. You must include it explicitly in the \LaTeX preamble (after the `\documentclass` declaration and before `\begin{document}`) using `\usepackage{graphicx}`.

5.3 Hyperlinks

Users of older versions of \LaTeX may encounter the following error during compilation:

This happens when pdf \LaTeX is used and a citation splits across a page boundary. The best way to fix this is to upgrade \LaTeX to 2018-12-01 or later.

5.4 Citations

Table ?? shows the syntax supported by the style files. We encourage you to use the natbib styles. You can use the command `\citet` (cite in text) to get “author (year)” citations, like this citation to a paper by ?. You can use the command `\citep` (cite in parentheses) to get “(author, year)” citations (?). You can use the command `\citealp` (alternative cite without parentheses) to get “author, year” citations, which is useful for using citations within parentheses (e.g. ?).

A possessive citation can be made with the command `\citepos`. This is not a standard natbib command, so it is generally not compatible with other style files.

5.5 References

The \LaTeX and Bib_{TEX} style files provided roughly follow the American Psychological Association format. If your own bib file is named `custom.bib`,

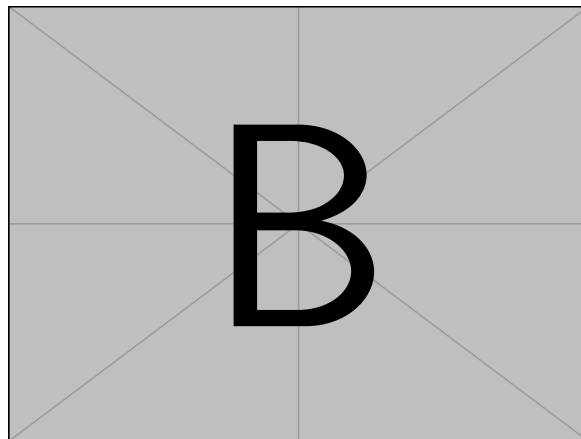
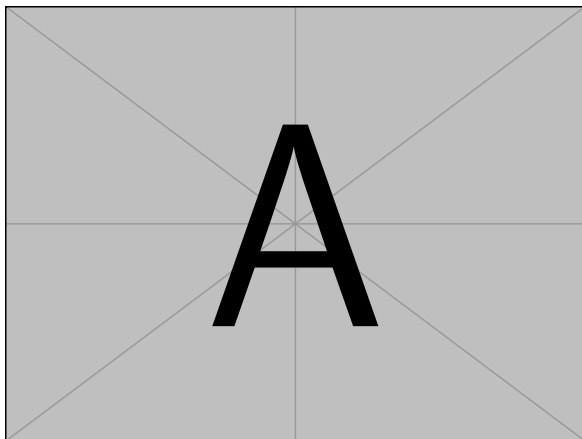


Figure 2: A minimal working example to demonstrate how to place two images side-by-side.

then placing the following before any appendices in your \LaTeX file will generate the references section for you:

```
\bibliography{custom}
```

You can obtain the complete ACL Anthology as a Bib \TeX file from <https://aclweb.org/anthology/anthology.bib.gz>. To include both the Anthology and your own .bib file, use the following instead of the above.

```
\bibliography{anthology,custom}
```

Please see Section ?? for information on preparing Bib \TeX files.

5.6 Equations

An example equation is shown below:

$$A = \pi r^2 \quad (1)$$

Labels for equation numbers, sections, subsections, figures and tables are all defined with the $\text{\label{label}}$ command and cross references to them are made with the $\text{\ref{label}}$ command.

This an example cross-reference to Equation ??.

5.7 Appendices

Use \appendix before any appendix section to switch the section numbering over to letters. See Appendix ?? for an example.

6 Bib \TeX Files

Unicode cannot be used in Bib \TeX entries, and some ways of typing special characters can disrupt Bib \TeX 's alphabetization. The recommended way of typing special characters is shown in Table ??.

Please ensure that Bib \TeX records contain DOIs or URLs when possible, and for all the ACL materials that you reference. Use the doi field for DOIs and the url field for URLs. If a Bib \TeX entry has a URL or DOI field, the paper title in the references section will appear as a hyperlink to the paper, using the hyperref \LaTeX package.

Acknowledgments

This document has been adapted by Steven Bethard, Ryan Cotterell and Rui Yan from the instructions for earlier ACL and NAACL proceedings, including those for ACL 2019 by Douwe Kiela and Ivan Vulić, NAACL 2019 by Stephanie Lukin and Alla Roskovskaya, ACL 2018 by Shay Cohen, Kevin Gimpel, and Wei Lu, NAACL 2018 by Margaret Mitchell and Stephanie Lukin, Bib \TeX suggestions for (NA)ACL 2017/2018 from Jason Eisner, ACL 2017 by Dan Gildea and Min-Yen Kan, NAACL 2017 by Margaret Mitchell, ACL 2012 by Maggie Li and Michael White, ACL 2010 by Jing-Shin Chang and Philipp Koehn, ACL 2008 by Johanna D. Moore, Simone Teufel, James Allan, and Sadaoki Furui, ACL 2005 by Hwee Tou Ng and Kemal Oflazer, ACL 2002 by Eugene Charniak and Dekang Lin, and earlier ACL and EACL formats written by several people, including John Chen, Henry S. Thompson and Donald Walker. Additional elements were taken from the formatting instructions of the *International Joint Conference on Artificial Intelligence* and the *Conference on Computer Vision and Pattern Recognition*.

A Example Appendix

This is an appendix.