

Instructions for *ACL Proceedings

Anonymous ACL submission

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

Blabla introduction The templates include the \LaTeX source of this document (`acl_latex.tex`), the \LaTeX style file used to format it (`acl.sty`), an ACL bibliography style (`acl_natbib.bst`), an example bibliography (`custom.bib`), and the bibliography for the ACL Anthology (`anthology.bib`).

2 Méthode

2.1 Dataset

Nous utilisons le dataset fourni pour la tâche 3 de l'édition 2009 de DEFT. Au moment de faire des statistiques descriptives, nous nous sommes rendus comptes que le corpus présentait des doublons, et ce majoritairement dans la partition test.

Après suppression des doublons, la partition prévue (40/60) est changée : elle est maintenant de 20/80¹. Nous avons envisagé de réimplémenter le partitionnement prévu, mais avons renoncé pour deux raisons : refaire le partitionnement nous éloigne, encore, du corpus initial, et les résultats de quelques modèles sur un corpus repartitionné étaient proches des résultats sur cette partition 20/80. Par ailleurs, la répartition des classes est déséquilibrée : les classes PPE-DE et PSE sont plus grandes et forment à elles deux 63,5 % du corpus. Ceci devra être pris en compte dans le prétraitement.²

¹0.79 pour le train et 0.21 pour le test

²la figure correspond au train, mais la répartition est sensi-

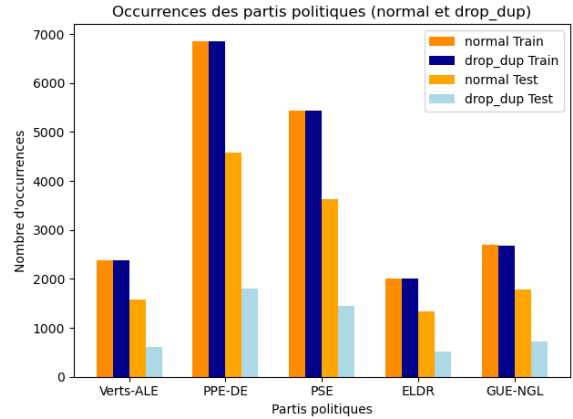


Figure 1: Nombre d'interventions par parti par partition test/train, dans le corpus original et dans la version sans doublons

013			
014			
015			
016			
017	Statistique	Train	Test
018	Moyenne	3871.4	1021.2
019	STD	2149.1	569.7
020	Min	2005.0	525.0
021	1er quartile	2376.0	615.0
022	Médiane	2687.0	715.0
023	3eme quartile	5431.0	1448.0
024	Max	6858.0	1803.0

Table 1: Nombre d'intervention des partis par partition

2.2 Prétraitements

Le texte des interventions a été soumis à un pré-traitement simple :

- (1) Suppression de la ponctuation
- (2) Unification de la casse en minuscules
- (3) Tokenisation³

Pour résoudre le problème de déséquilibre des classes, nous avons opté pour le *downsampling*

blement la même dans le test

³Une lemmatisation avec la bibliothèque SpaCy a été envisagée, mais ce corpus multilingue aurait nécessité le chargement de 3 modèles linguistiques différents et ralenti le temps de traitement

Répartition des partis politiques dans le corpus train

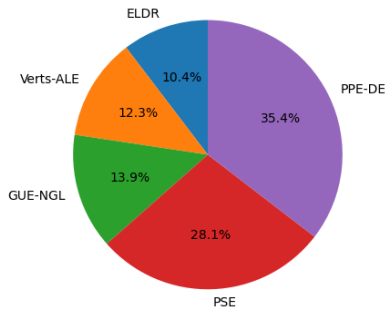


Figure 2: Répartition des interventions par parti dans la partition train sans doublons

afin d’obtenir des classes relativement équilibrées, en utilisant la fonction *resample* de la bibliothèque *scikit-learn*. Après *downsampling*, les classes PPE-DE et PSE ne représentent plus que 43% du corpus, ayant chacune été ramenée autour de 21,5% du corpus

Parti	Train	Test
ELDR	2531 (16,1%)	525 (16,2%)
PPE-DE	3402 (21,63%)	715 (22%)
GUE-NGL	3402 (21,63%)	687 (21,2%)
PSE	3402 (21,63%)	693 (21,4%)
Verts-ALE	2990 (19%)	614 (19%)
Total	15727	3235

Table 2: Nombre d’intervention par parti par partition pour une langue après *downsampling*, exemple de l’italien

2.3 Les différents embeddings testés

Nous avons choisi dans notre étude de comparer les résultats obtenus sur une tâche de classification en utilisant 3 techniques de vectorisation différentes.

La vectorisation TF-IDF⁴ (*Term Frequency-Inverse Document Frequency*) est la méthode la plus "ancienne" que nous présentons. Elle se base sur une mécanique de comptage des mots: la fréquence d’apparition de chaque mot dans un document est divisée par sa fréquence d’apparition dans le corpus, permettant de donner plus d’importance aux mots significatifs dans leurs documents d’apparition.

La vectorisation Doc2Vec⁵ qui génère des

⁴Implémentée avec la fonction *tfidfvectorizer* de *scikit-learn*

⁵Implémentée à l’aide de la bibliothèque *gensim*

vecteurs de document plutôt que de mot. Ces vecteurs sont l’output d’un réseau de neurones et nécessitent donc une phrase d’entraînement. Pour cette étude, nous avons choisi de générer des vecteur Doc2Vec à 100 dimensions⁶ avec une fenêtre glissante de 5 mots et d’ignorer les mots n’apparaissant pas au moins 3 fois.

La vectorisation avec BERT multilingue⁷ qui se base sur les réseaux de neurones mais s’appuie sur une architecture transformer et nécessite également une phase d’entraînement. Contrairement à Doc2Vec, elle renvoie des vecteurs de mots, et utilise un mécanisme d’attention lors de la génération des vecteurs, lui permettant de prendre en compte l’importance d’un mot en fonction du contexte local dans lequel il apparaît. Nous avons généré des vecteurs à 768 dimensions (valeur de base).⁰⁴⁷

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 L^AT_EX 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.

⁶il serait possible de faire plus, mais nous essayons de ne pas saturer nos machines

⁷bert-base-multilingual-uncased

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



Figure 3: 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.

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.

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 `\citeposs`. 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`, 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 `\label{label}` command and cross references to them are made with the `\ref{label}` command.

This is 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.

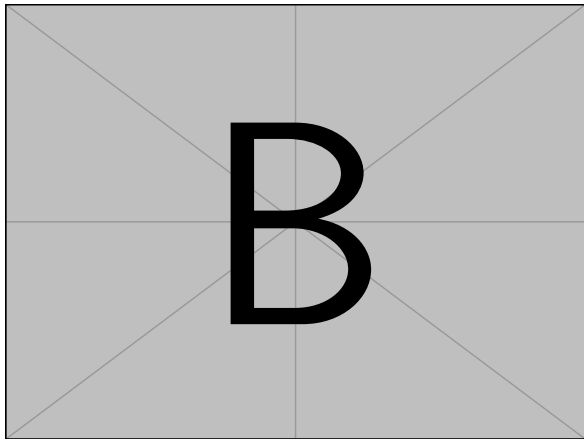
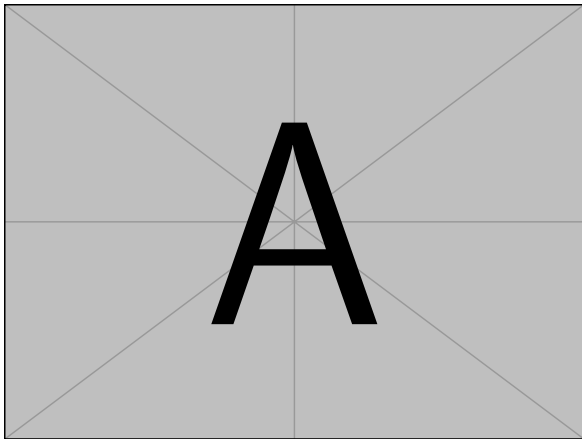


Figure 4: A minimal working example to demonstrate how to place two images side-by-side.

6 BibT_EX Files

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

Please ensure that BibT_EX 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 BibT_EX 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 L^AT_EX 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, BibT_EX 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.

206
207
208