



# T9 - Artificial Intelligence

T-AIA-901

## Travel Order Resolver

Bootstrap





You are about to build an automated tool that sorts political speeches, kind of a political compass.

More specifically it takes a bunch of speeches in french language and classifies them in two clusters, based on the vocabulary they use. The underlying question is: *“Does this clustering make sense, according to what we know about the authors?”*

## GETTING THE DATA SETS

Firstly, you need to recover the texts from *discours.zip*. You should be able to produce a **Bunch** object with several attributes, including:

- a **filenames** attribute that contains “*Alain\_Madelin\_702.txt*, *Arlette\_Laguiller\_32.txt*”, ...
- a **data** attribute that contains the texts themselves



[scikit\\_learn documentation](#) might help you with data loading and text processing

To make sure everything works all right, display a dictionary that associates all file titles with the first 40 characters of the file. Your result should look like this:

```
Terminal
~/T-AIA-901>
'Arlette_Laguiller_299.txt': 'Travailleuses, travailleurs, camarades et amis,
Pour commencer ce meeting, je', ...
```

## FROM RAW TEXT TO VECTOR

As far as pictures are concerned, the typical mapping is to code every vertex with a RGB value. When it comes to text documents, the simplest way is to count (for every possible word in the language) how many times this word appears in the text.

Next step is to apply a **tf-idf transformation**, to replace this raw count with relative frequencies.

```

. 0.      0.      0.      0.      0.      0.      0.
0.      0.      0.57501852 0.34113499 0.      0.      0.
0.      0.      0.      0.47152341 0.57501852 0.      ]
[ 0.      0.      0.      0.40824829 0.40824829 0.      ]
0.40824829 0.      0.      0.      0.      0.      0.
0.40824829 0.40824829 0.40824829 0.      0.      0.      ]
[ 0.      0.490779 0.490779 0.      0.      0.      0.
0.      0.      0.      0.      0.      0.      0.
0.      0.      0.490779 0.4024458 0.      0.3397724 ]
[ 0.26396385 0.      0.      0.      0.52792769 0.      ]
0.      0.79189154 0.      0.      0.15659897 0.      0.
0.      0.      0.      0.      0.      0.      0.      ]
[ 0.      0.      0.      0.      0.      0.      0.
0.      0.59430676 0.      0.35257792 0.59430676 0.      0.
0.      0.      0.      0.      0.      0.41144595]
[ 0.      0.      0.      0.      0.      0.      0.
0.      0.      0.      0.26990725 0.      0.9099135
0.      0.      0.      0.      0.      0.
0.31497221]]

```

Anyway, this builds a huge, sparse matrix whose number of columns is the size of the vocabulary. You're expected to produce a (9w,216zO) sparse matrix, for there are 9w texts and 216zO different words.



w and z represent two different digits



Doing this has you lose any information about positions of words in the text. This might not be the smartest thing to do for some applications.

## UNSUPERVISED LEARNING

If you don't want to feed our algorithm with *a priori* information about the political opinion of the orator, but instead, you want to associate texts freely, based on the sole vocabulary. Hence, you need to do some *unsupervised learning*.

Let's start with a simple basic classifier you may already encounter somewhere, namely *k-means*. Fix  $k = 2$  and see what comes out of the box.

Notice that most politicians have all their speeches that belong to one class, which makes sense!



Find the dates when they were candidates, then make a hypothesis on how the algorithm divides speeches in two categories here?



Be careful that k-means is based on strong hypotheses (Voronoi-based division of the space), especially when you have only two classes. You can probably do much better!

Now, play around with other algorithms, tune some parameters, investigate metrics, observe results.



Text detection can be improved a lot by removing *stopwords*. Similarly, look for *lemmatization* to help you with *inflection*.

## SUPERVISED LEARNING

Now, if you want to challenge some algorithms to identify one specific author, you can divide your sample in 3 categories:

- a folder containing a part of the texts from one politician
- a folder containing an equal part of the texts from other candidates
- a folder containing all the remaining texts from all the orators

Train a **Naive Bayes** algorithm using the first folders, then check predictions quality with the third folder.



How much learning an algorithm do you need to train? Find out an ideal training/testing ratio!



Naive Bayes is a very common, though simple, algorithm in NLP. Take time to understand the maths behind it. Other algorithms, such as Topic Modeling, are much more difficult to understand. So, best to master this one before jumping to more sophisticated ones.