

Text Mining and Text Retrieval Group Assignment

General Instructions

- Submit the Jupyter Notebook in electronic form before 23:59 on Sunday, January 26th
- The Notebook should be submitted through CANVAS
- No late homework reports will be accepted.
- The Notebook should be submitted by one individual per group.
- The group members should be identified in the Notebook.
- The Notebook text cells should be written in English.
- The Notebook can run from end to end. If necessary, have variables in the first cell that point to the directories where the files are, so I can adjust for my laptop
- It is possible to use a Google Colab Notebook
- Take into account the shell commands for downloading / unzipping the data
- Take into account the installation of libraries through pip
- Use Text Cells to discuss your choices and results
- For questions, please use Canvas.

Introduction

You will work with a corpus made of news articles. You can download the data from [Kaggle](#).

This is a collection of approximately 140.000 articles from 15 different sources, published between the years 2000 and 2017, the balance between publications is shown in Figure 1.

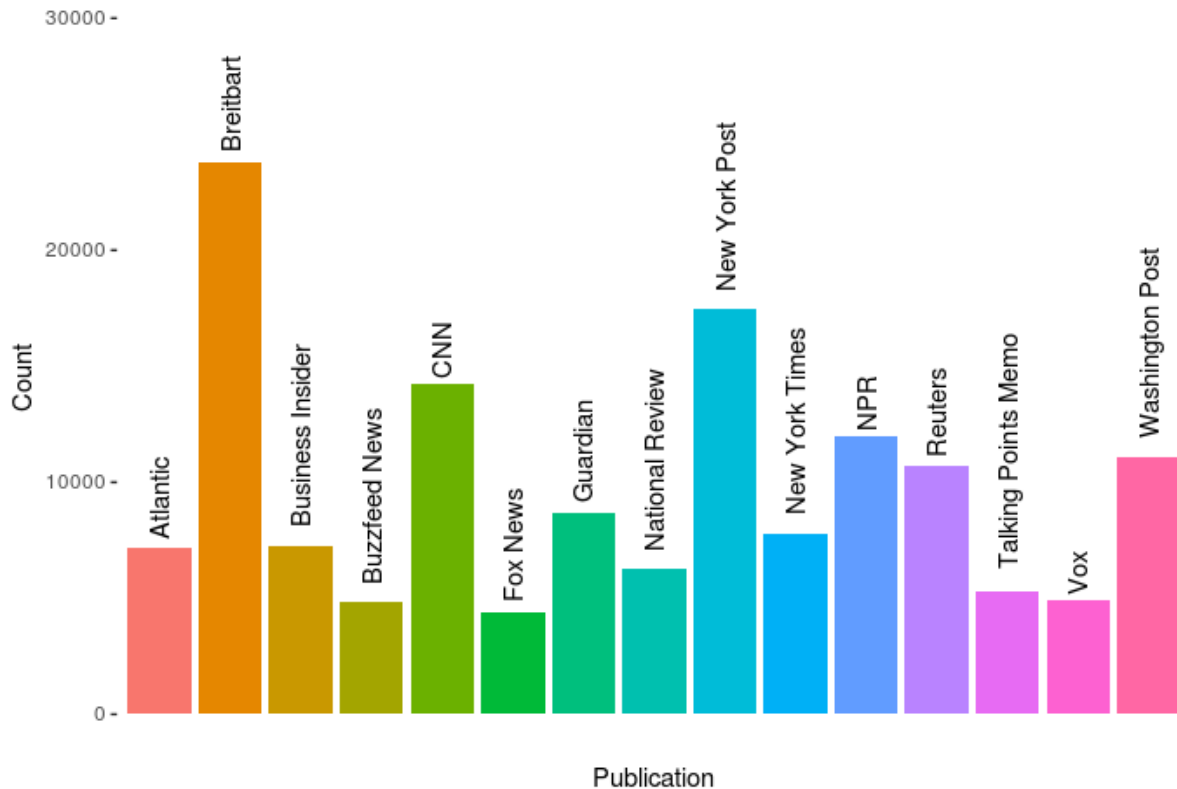


Figure 1: Balance of Publications

The data are provided as 3 CSV files, it is advised to load them as one single pandas dataframe.

The fields can be described as follows:

- Publication: the name of the media that published the article. This is one item from the list above.
- Title: The title of the article
- Date: the publication date
- Year: publication year
- Month: publication month
- Content: the full text of the article
- Author: name of the author
- ID: an id to identify each article

This homework will cover the material presented in Lectures 1 and 2: Lexical and Semantic Representation of Text. You can use any (or all) of the methods described along with anything you covered in the Machine Learning course regarding classification.

You will submit a Jupyter Notebook, in which you shall mention all the processing done to the text (stopping, stemming, ...). The Jupyter Notebook should work, it will be used as a tool for result reproducibility, but the quality of the code will not be assessed.

Regarding the text cleaning, remember you have the freedom to complete the list of stopwords, to remove words based on document frequency, etc. . . mention what you do, explain it shortly. If this is an intuition, or an assumption, state it. There are no trap, and most of the time, there are no good/bad decision.

Assignment

Question:	1	2	3	4	5	6	Total
Points:	10	10	20	20	20	20	100

You are going to analyze the content of the articles, inspired by the Vogue analysis we discussed in class.

1. (10 points) Describe the corpus
 - (a) Select a publication of your choice
 - (b) Use [AllSides](#) to identify the Political Bias (if any) of your publication of choice
 - (c) Emit an hypothesis about either the choice of topics covered by this publication, or their angle when covering topics. Differences will arise with controversial political topics, such as Diversity, Gun Control, Women's rights, Drugs Control, etc...
 - (d) Illustrate your description with statistics about number of texts, text lengths, initial dictionary size, etc...
2. (10 points) Pre-process the corpus
 - (a) Create the Dictionary, reduce it to a viable size
 - (b) Use Stopwords and Stemming
 - (c) If needed, add a cleaning step
3. (20 points) Produce a few Topic Models of your corpus
 - (a) Create the corpus BoW by incorporating from 1 to n-grams
 - (b) Use Raw Counts or TF-IDF weights for the Term-Document matrix
 - (c) Use at least 2 different values for K, the number of topics
 - (d) For each model, store the document embeddings
4. (20 points) Compare the Topics
 - (a) For a value of K, compare topics rising from SVD to topics rising from LDA
 - (b) For example, use Jaccard Index to compare the sets of Top-N words from topics
 - (c) Study similarities. Illustrate with visualizations
5. (20 points) Study your corpus
 - (a) Select a value for K

- (b) Select a topic from LDA
 - (c) Observe the saturation in your topic for each article in the corpus
 - (d) Illustrate with the Top 10 articles with highest saturation in your topic
 - (e) Discuss if this is in line with your thoughts from Question 1
6. (20 points) Predict a publication
- (a) Select a second publication, possibly from an opposite side on the political spectrum
 - (b) Learn a LDA model of the joint corpus
 - (c) Use the LDA embeddings to predict which publication originated the article
 - (d) Compare with a lexical method of your choice