# Press-Release Classification && Sorting







**Tatiana Fomicheva MSc in Big Data** 

**Scraping** 

Identification of the common link part to the documents; PDF/PPT/Word automated download;

Conversion to .TXT

**Tools:** BeautifulSoup, Mechanize

Model **Evaluation&&Extension** 

Accuracy testing; Generalize to be used for Press Releases from different companies

**Tools:** Sklearn metrics

#### **Tokenization**

Separate the PR's in words/per document **Tools:** NLTK

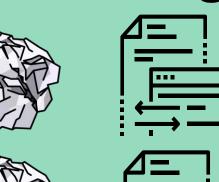
#### **Pre-Processing**

specific common news identifiers and custom stop-words; Stemming and Lemmatization

**Tools:** NLTK

Data Acquisition











Modelling



## **Feature Engineering**

Bag of Word OR

Word Embedding OR

Combination of BoW & Word Embedding

Tools: TF-IDF, Word2Vec, GloVe,

FastText

Goal

Create a generalized framework which automates the text acquisition from different sources, pre-processing and modelling to attribute labels per text file accurately

### **Vectorizer&&Model Tuning**

Hyperparameters Tuning;

Weight adjustment, reducing false negatives

Tools: GridSearchCV, Pipeline, Matplotlib, ...



**Model (Classifier) Building** 

Supervised Models testing

Tools: Pipeline, RandomForest, GradientBoosting, ExtraTrees, Bagging, Kneighbours, MLP, LogReg, ...

