# INF554 - Team TVRPZ

Romain Fouilland - Philémon Gamet - Jacques Song

École polytechnique

January 11, 2019

# Plan

Basic features, citations graph and first predictive model
Solving the overfitting issue
Final model and parameters tuning

First features
Graph-based features
First predictive model
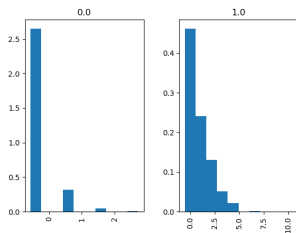
## Raw data & data exploration
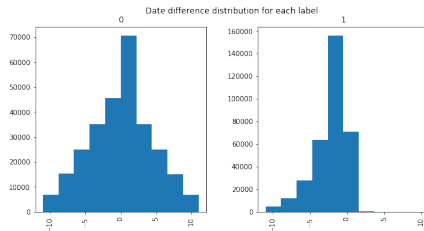


Figure 1: Overlap between titles



Figure 2: Temporal difference

- Raw data: date, title, authors, journal, abstract
- Features: difference between dates, overlap between titles, authors in common...

Basic features, citations graph and first predictive model
Solving the overfitting issue
Final model and parameters tuning

First features
Graph-based features
First predictive model

## Abstract embedding

How to use the abstract?

- Overlapping words between articles
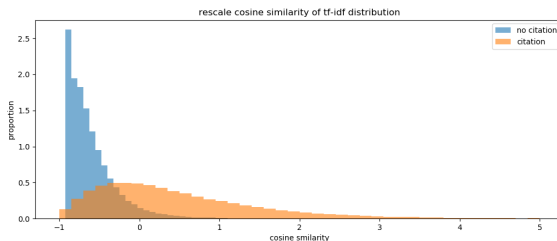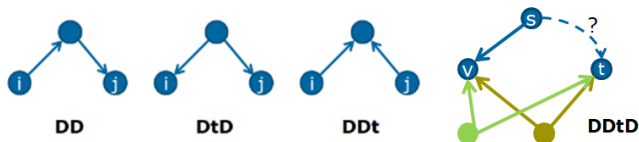- Embedding: a representation of text in high dimensional vector space



Figure 3: Tf-idf cosine similarity

Basic features, citations graph and first predictive model
Solving the overfitting issue
Final model and parameters tuning

First features
**Graph-based features**
First predictive model

# The graph of documents

Nodes: documents, edges: citations



$D$ the adjacency matrix of the graph:

- $DD^T$: number of common neighbors
- $DD^T D$: quoted simultaneously

Compute the degrees: $d_{in}(s), d_{in}(t), d_{out}(s), d_{out}(t)$

Basic features, citations graph and first predictive model
Solving the overfitting issue
Final model and parameters tuning

First features
Graph-based features
**First predictive model**

# First neural network

Score with SVM classifier trained on 5% of the dataset: **96,52 %**
$\rightarrow$ Computed features are meaningful

Multi Layer Perceptron

- 1 hidden layer of 16 perceptrons
- Trained on 5% of the dataset
- Accuracy: **97,01%** on kaggle

Basic features, citations graph and first predictive model
Solving the overfitting issue
Final model and parameters tuning

First features
Graph-based features
First predictive model

## Bad results with all the data

To improve predictions: train over a larger part of the dataset

- Good results on our validation dataset
- Very poor results on kaggle (around 85%)...

# Plan

Basic features, citations graph and first predictive model
**Solving the overfitting issue**
Final model and parameters tuning
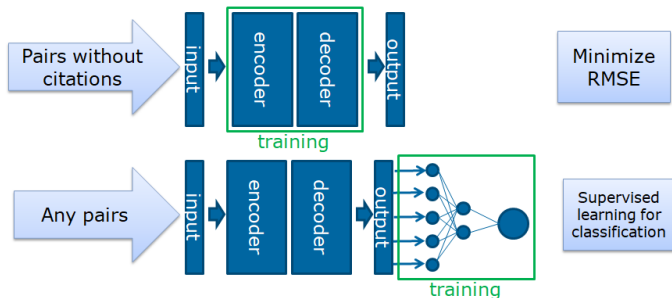
Change the model
Deeper data exploration
Random forests

# Autoencoder



Figure 4: Autoencoder pipeline

- No overfitting thanks to compressed information
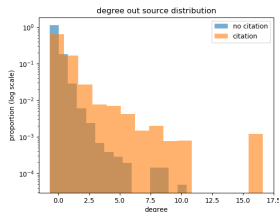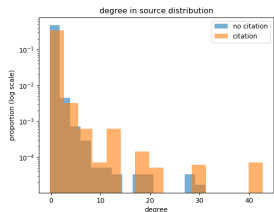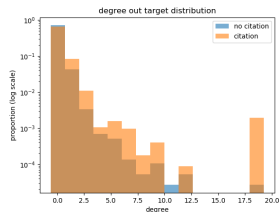- Accuracy: **95,24%** on kaggle

Basic features, citations graph and first predictive model
Solving the overfitting issue
Final model and parameters tuning

Change the model
Deeper data exploration
Random forests

## Distributions



Figure 5: Source degrees



Figure 6: Target degrees

Basic features, citations graph and first predictive model
**Solving the overfitting issue**
Final model and parameters tuning

Change the model
**Deeper data exploration**
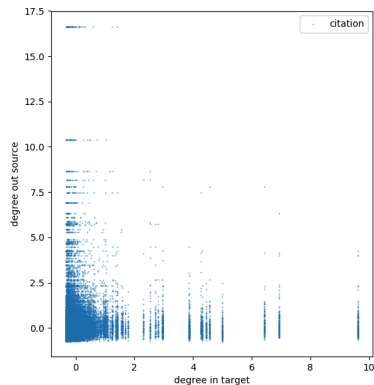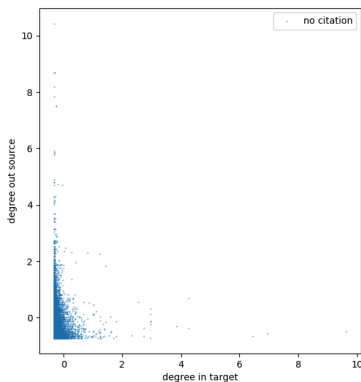Random forests

# Correlations



Figure 7: Correlation between source degree out and target degree in

Basic features, citations graph and first predictive model
**Solving the overfitting issue**
Final model and parameters tuning

Change the model
Deeper data exploration
**Random forests**

# Decision Tree analysis



Figure 8: Decision Tree understanding of the data

3 features: DDtD, year, tf-idf and 95.62% accuracy

Basic features, citations graph and first predictive model
Solving the overfitting issue
Final model and parameters tuning

Change the model
Deeper data exploration
Random forests

# Random forest results

Random forest on basic features: 77.9%
Random forest on all features:

- 99.7% on validation
- 80.1% on Kaggle

$\Rightarrow$ Is our validation set correct ?

Basic features, citations graph and first predictive model
**Solving the overfitting issue**
Final model and parameters tuning

Change the model
Deeper data exploration
**Random forests**

# Data splitting issue



Figure 9: Initial data split

The results are already in the training data !

Romain Fouilland - Philémon Gamet - Jacques Song    INF554 - Team TVRPZ

Basic features, citations graph and first predictive model
**Solving the overfitting issue**
Final model and parameters tuning

Change the model
Deeper data exploration
**Random forests**

# Data splitting issue



Figure 10: No overfit on few data (5%)



Figure 11: Overfit with more data (70%)

Romain Fouilland - Philémon Gamet - Jacques Song    INF554 - Team TVRPZ

Basic features, citations graph and first predictive model
Solving the overfitting issue
Final model and parameters tuning

Change the model
Deeper data exploration
Random forests

# Data splitting issue

The results are already in the training data !

Figure 12: Initial data split

It's now safe to learn.

Figure 13: Final data split

Romain Fouilland - Philémon Gamet - Jacques Song     INF554 - Team TVRPZ

Basic features, citations graph and first predictive model
**Solving the overfitting issue**
Final model and parameters tuning

Change the model
Deeper data exploration
**Random forests**

## Random Forest analysis



Figure 14: Features importance in a random forest

Well performing forest (96.59% accuracy)
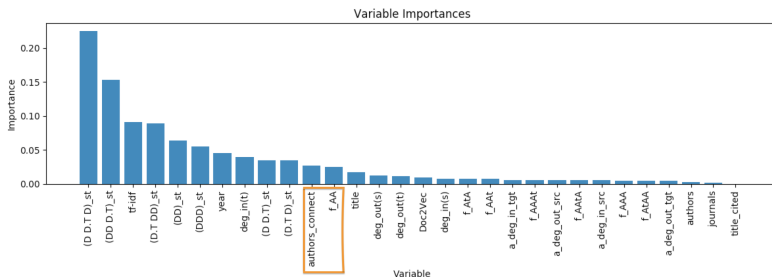$\Rightarrow$ Add new features: abstract and authors graph

Basic features, citations graph and first predictive model
**Solving the overfitting issue**
Final model and parameters tuning

Change the model
Deeper data exploration
**Random forests**

# Random Forest analysis



Figure 15: Features importance in our best random forest
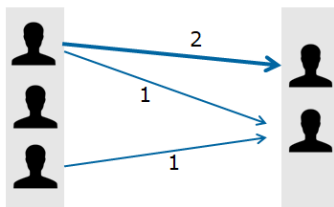
High performing forest (97,11% accuracy)

# Plan

Basic features, citations graph and first predictive model
Solving the overfitting issue
Final model and parameters tuning

Enriching the data
Parameters tuning
Final results

# The graph of authors

Build a citation network for the authors

- Nodes: authors, edges: number of citations
- Multiple authors for each paper



$$author\_connection = \sum_{s \in S} \sum_{t \in T} A_{st}$$

Basic features, citations graph and first predictive model
Solving the overfitting issue
Final model and parameters tuning

**Enriching the data**
Parameters tuning
Final results

## Convolution Network for abstract embedding

**Idea**: convert abstract into features maps, train CNN on these maps and collect intermediate layer output to reuse as features in the final model.

**Input data**:

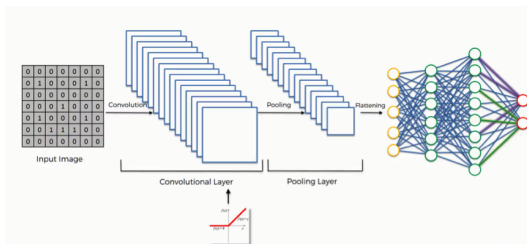- Word embedding: word2vec
- Abstract embedding: array of word2vec vectors



Figure 16: Structure of a CNN

Basic features, citations graph and first predictive model
Solving the overfitting issue
**Final model and parameters tuning**

**Enriching the data**
Parameters tuning
Final results

# Convolution Network for abstract embedding

Results of the CNN: too much overfitting
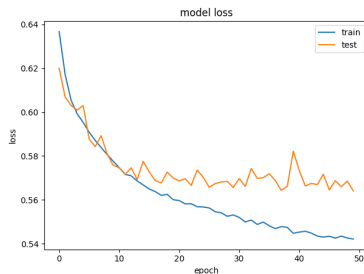


Figure 17: CNN model accuracy showing overfitting



Figure 18: CNN model loss showing overfitting

Romain Fouilland - Philémon Gamet - Jacques Song          INF554 - Team TVRPZ

Basic features, citations graph and first predictive model
Solving the overfitting issue
Final model and parameters tuning

Enriching the data
Parameters tuning
Final results

# Overfit vs. Regularization

How to reduce overfitting directly in the neural network ?

- Dropout
- Regularization
- Batch size



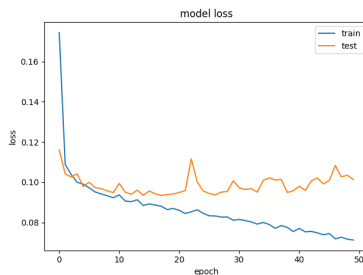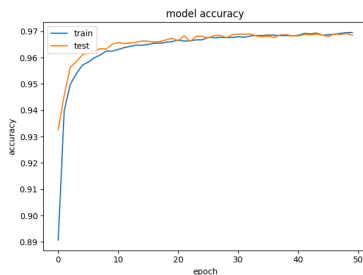Figure 19: Model accuracy without regularization



Figure 20: Model loss without regularization

Basic features, citations graph and first predictive model
Solving the overfitting issue
Final model and parameters tuning

Enriching the data
Parameters tuning
Final results

# Overfit vs. Regularization

- Dropout: add dropout layers between each dense layer
- Regularization: $l_2$-regularization
- Batch size: increase size



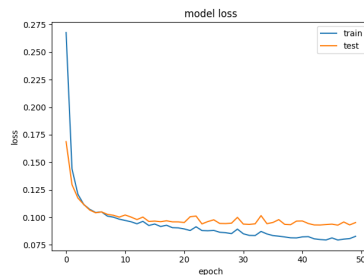Figure 21: Model accuracy with regularization

Figure 22: Model loss with regularization

Basic features, citations graph and first predictive model
Solving the overfitting issue
Final model and parameters tuning

Enriching the data
Parameters tuning
Final results

# Best results

| 38 | — | ugfdd908 | | 0.97035 | 7 | 6d |
|---|---|---|---|---|---|---|
| 39 | ▼ 16 | Pumpkin | | 0.97035 | 22 | 2d |
| 40 | ▼ 14 | TVRPZ | | 0.97010 | 28 | 3d |
| 41 | ▲ 1 | MBS | | 0.96986 | 23 | 2d |
| 42 | ▼ 2 | BDC | | 0.96961 | 17 | 4d |
| 43 | ▲ 1 | Navy | | 0.96955 | 10 | 2d |



LIKE A BOSS