# Skoltech
Skolkovo Institute of Science and Technology

# Multilabel text classification

Sofia Medvedeva
Tatiana Smirnova
Almir Dzhumaev
Adel Galimov
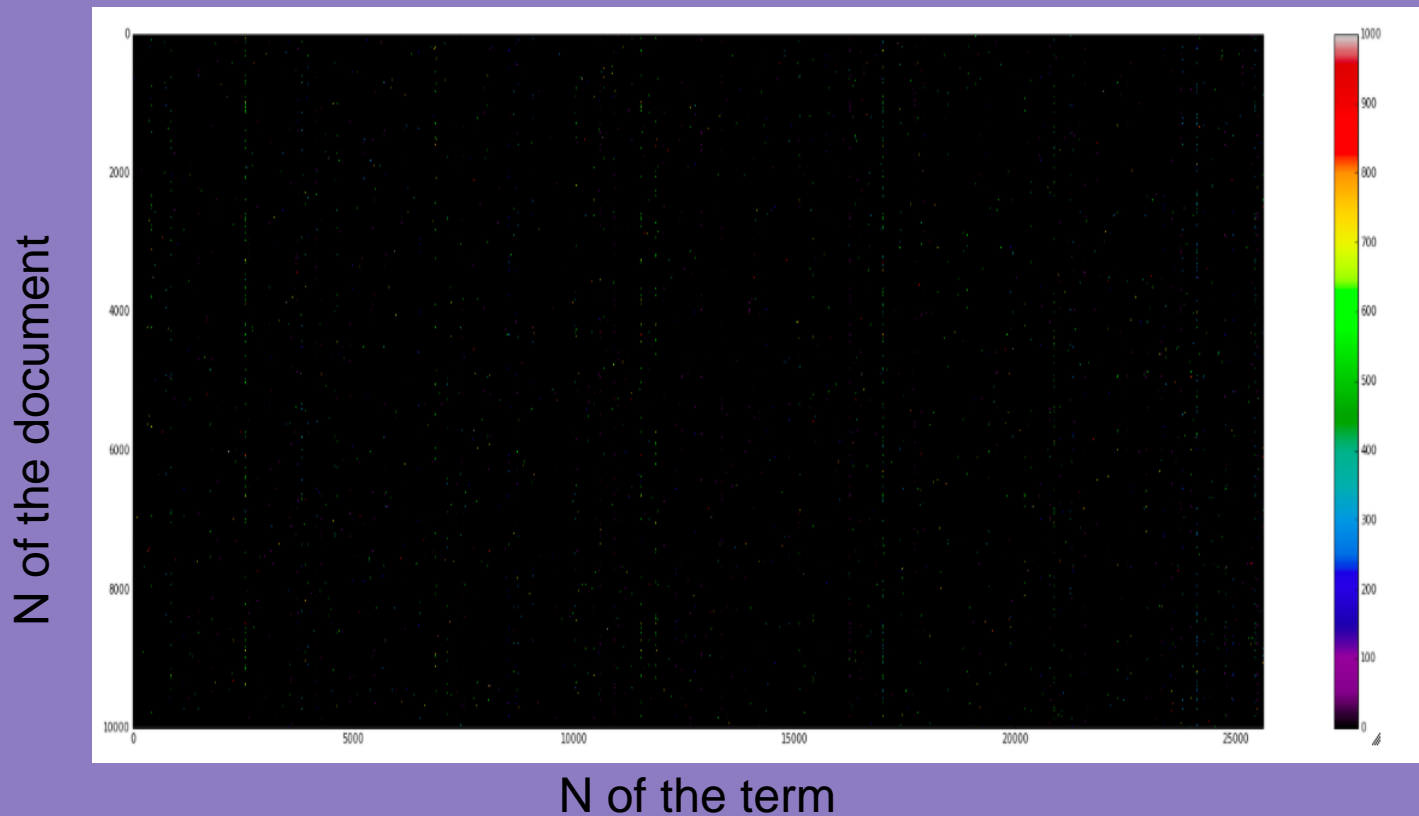
Moscow 2014

# Task and goal

- We have data about scientific terms importances in different articles

- Each article may belong to several topics

- Goal: to build a classifier for topic prediction (and win Kaggle competition!)

# Roles in the team

- Tatiana, Adel - data visualizations, preprocessing
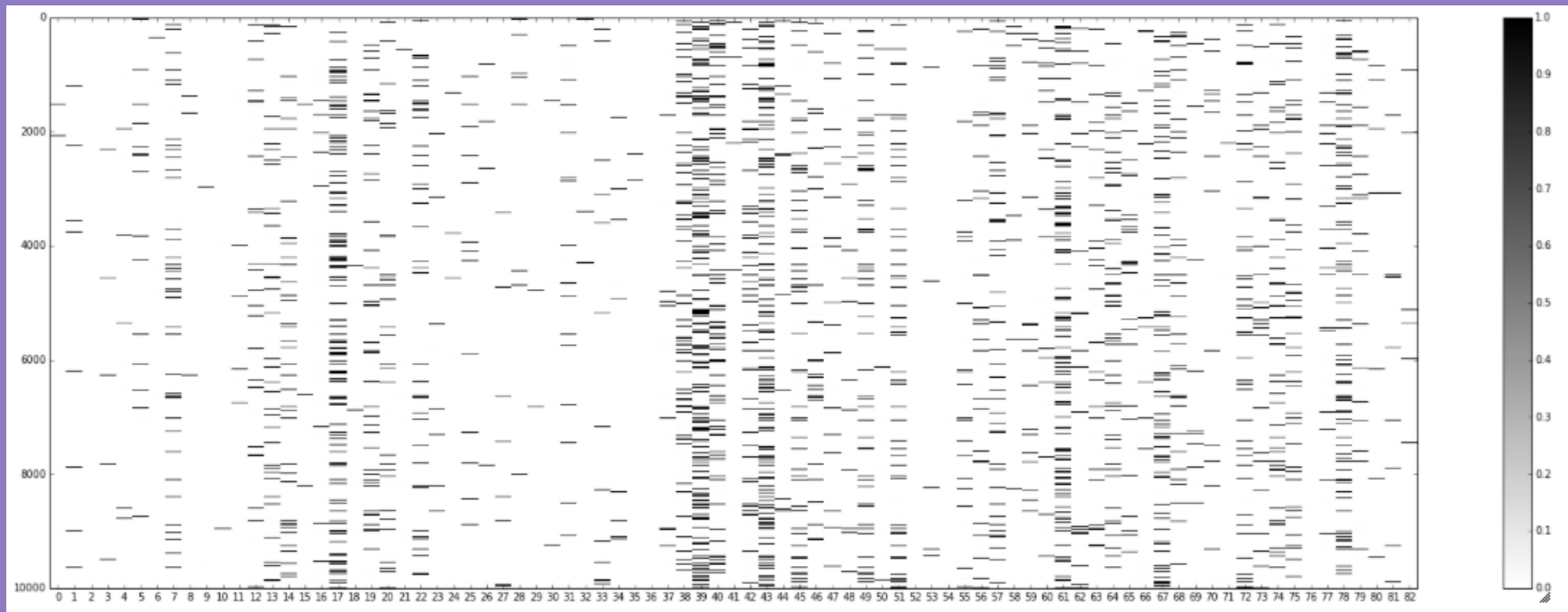
- Sofia, Almir - classificators' training

# Data Structure. X_train matrix.
# 10,000x25,000 sparse matrix

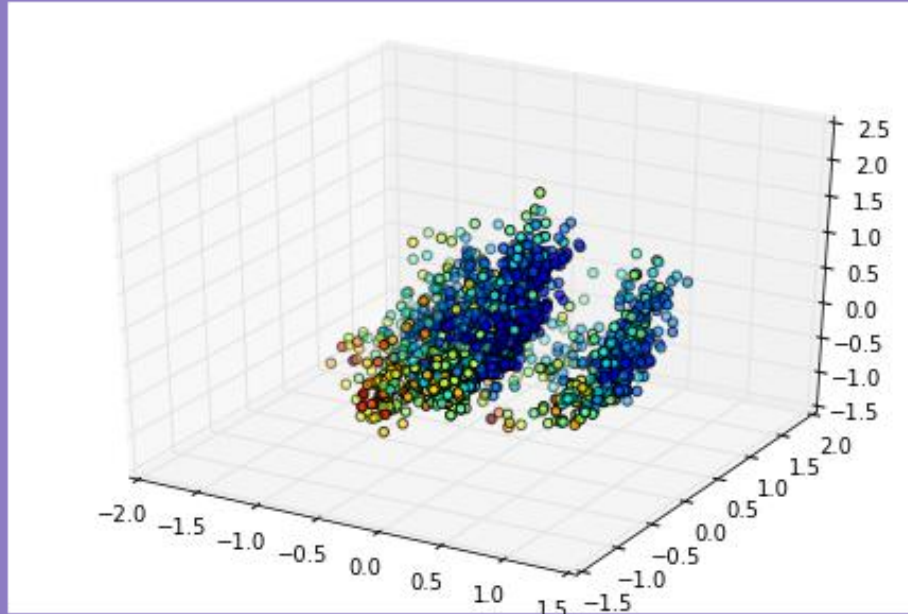# Labels Structure. y_train. 10,000x83 sparse matrix
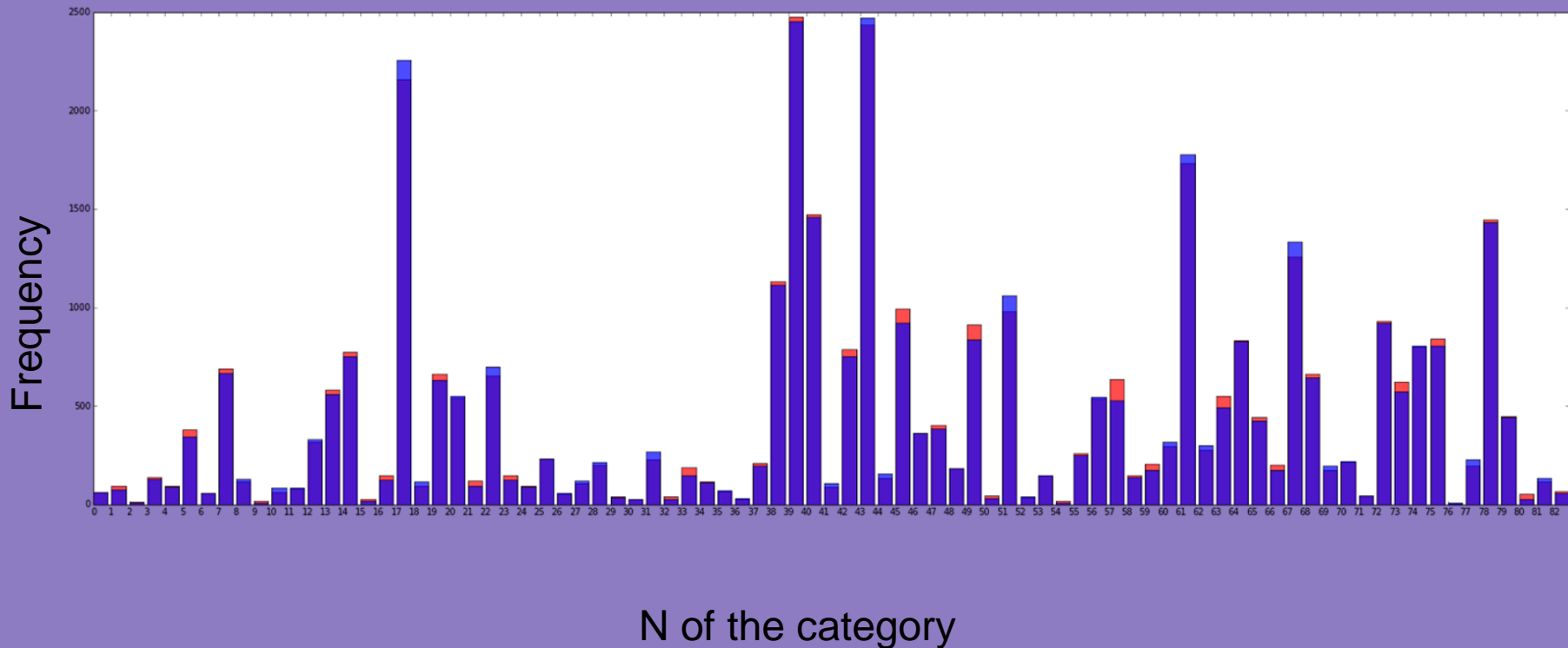


N of the document

N of the category

# Data Structure. Labels visualization

PCA with 4 dimensions

# Train and test labels structure. Categories frequencies.
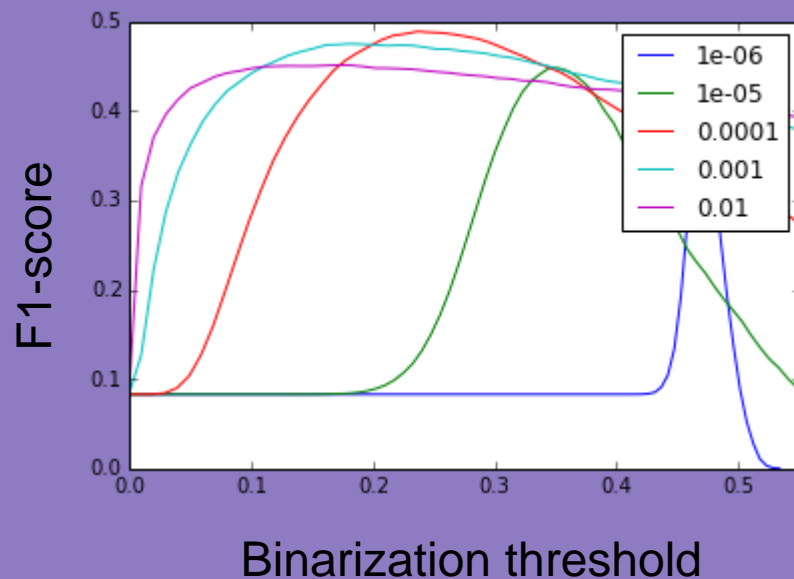
# ML methods used

- Random Forest

- SVM

- Naive Bayes

- Logistic regression

# Strategies used to improve the result

- Combination of methods (blending, pipeline)
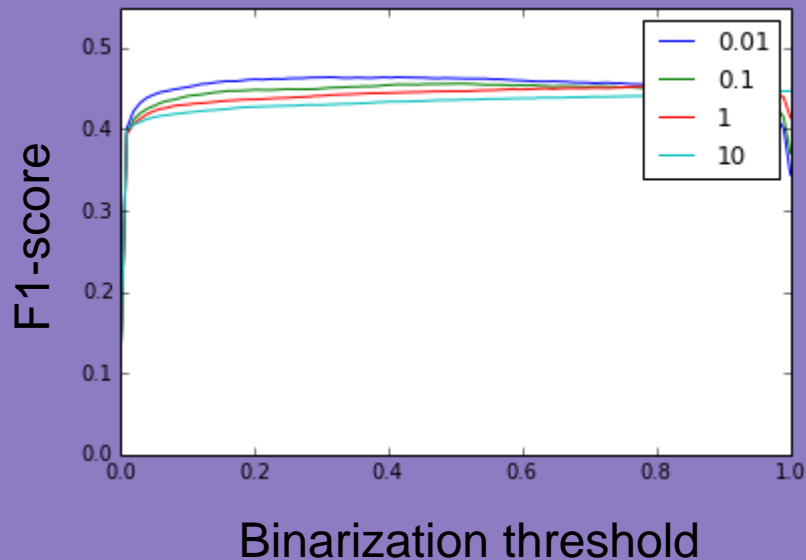
- Binarization

- Feature selection

# Binarization

Binarization can help!

# Feature selection

- l1 regularization

- K-best

# Final result - 1st place

Best result achieved by logistic regression with optimized parameters' values - 0.51 f1-score