

# Author Attribution of Turkish Texts by Feature Mining

Filiz Türkoğlu, Banu Diri, and M. Fatih Amasyalı

Narrators: Mehmetcan Güleşçi, Furkan Karakoyunlu

May 13, 2017

# What is the problem?

- One of the problems in text categorization is the authorship attribution, which is used to determine the author of a text when it is not clear who wrote it
- In some occasions where two people claim to be the author of same manuscript
- or on the contrary where no one is willing to accept the authorship of a document

# The aim of the article

- They focused on author attribution of Turkish texts by extracting various feature vectors and applying different classifiers
- They studied the comparative performance of classifier algorithms using the Naive Bayes, Support Vector Machine, Random Forest, Multilayer Perceptron, and k-Nearest Neighbour
- To conclude they calculated the effectiveness of the methods by using 10-fold cross validation

# Early researches before this article

- Early researchers in authorship attribution used a variety of statistical methods but it tends to vary from author to author
  - Mosteller and Wallace - Federalist Papers, by using set of function words
  - Yule, by using complexity-based features (average sentence length, average word length, type/token ratio ..)
  - Recent technical advances in automated parsing and POS tagging, by using syntactic features such as POS n-grams
  - Peng, they modeled each author by a vector of the most frequent n-grams in the text
  - Fung - Federalist Papers, by using Support Vector Machine classifier
  - and so on (devamı çok uzun onları biz söyleriz 2 Authorship Attribution den önceki kısımda yazıyo veya 1 slayt daha devam ederiz)

# AUTHORSHIP ATTRIBUTION

- In this article, they formed feature vectors from several categories of statistics
- Stylistic analysis in order to compare the efficacy of each
- They examined features in five main categories, which are statistical, vocabulary richness, grammatical, lexical and n-grams model
- They obtained five different feature vectors from the mentioned categories
- Then, they created five feature subsets by using the feature selector to reduce the dimension of vectors

- In this article, they used a corpus, 630 documents written by a single author are obtained from 35 texts per 18 different authors
- On different subjects like sport, popular interest and economics
- From Turkish daily newspaper, “Hürriyet” and ”Vatan”
- In order to determine the authorship attribution performance in homogeneous and heterogeneous documents, and different dataset sizes, this corpus is divided into 3 parts: Dataset I, Dataset II, Dataset III.

# General Feature Vector (gfv)

- **Statistical features:** counting features in a text and applied this to word lengths and sentence lengths
- **Vocabulary richness features:** there are different statistics to determine the richness of an author's vocabulary. These features shows author's creativity.
- **Grammatical features:** based on Turkish Word Database (TWD) 35,000 words is used.
- **Function word features:** words that have little lexical meaning or have ambiguous meaning, but instead serve to express grammatical relationships with order words within a sentence.

# N-Gram Model

- They extracted bi-grams and tri-grams
- While forming the bi-grams and the tri-grams of the corpus, they observed that the number of different bi-grams and tri-grams are too much
- In order to avoid the combinatorial explosion in the feature vectors, they used a threshold value (greater than 75) to reduce the number of features.
- After removed infrequent features, the dimensions of the bi-gram *bgfv*, and tri-gram *tgfv* feature vectors are 470 and 1037 respectively.
- After that they combined *bgfv* and *tgfv*, they had a new feature vector(*btgfv*) which has a dimension of 1507.
- Finally they put together *gfv* and *btgfv* and obtained a 2148-dimensioned new vector, which is called *gbtgfv*.



# Feature Selection

- A high number of features may slow down the process while giving similar results as obtained with much smaller feature subset.
- To learn the effect of high-dimensioned feature set over success ratio, they used CfsSubsetEval function which is in WEKA.
- They reduced features of general feature vector, *gfv* and obtained a new vector, *rgfv*.
  - 24 dimension for Dataset I
  - 17 dimension for Dataset II
  - 40 dimension for Dataset III

- Same process was applied for Bi-gram feature vector, *bgrfv*, and was formed *rbgrfv*
  - 25 features for Dataset I
  - 20 features for Dataset II
  - 63 features for Dataset III
- They decreased dimension of Tri-gram feature vector, *tgrfv*, and obtained *rtgrfv*
  - It has left 60 most distinguishing features for Dataset I
  - 33 for Dataset II
  - 101 for Dataset III
- When features decreased from *btgrfv*, *rbtgrfv* is obtained. For Dataset I, vector has 61 features, for Dataset II it has 26 features, and for Dataset III it has 101 features.

- They decreased dimension of *gbtgfv*, and obtained *rgbtgfv*. It has left 69 most distinguishing features for Dataset I, 30 for Dataset II and 103 for Dataset III. All used feature vectors are shown at Table 1.