

# Development of text classification system

**A Sabyrkulova\***

*Faculty of Information Technology, Kazakh-British Technical University, Tole bi Str.59, Almaty, Kazakhstan*

*\*Corresponding author's e-mail: sabyraliya@gmail.com*

## Abstract

In this article we will consider the approaches applied to classification of texts in a natural language processing by their topics. The problem of automatic classification of texts is considered by sentiment analysis, there are described methods of machine learning for the solution of this problem. The description of an algorithm of classification of the naive Bayesian classifier is provided. Utilization of classifiers, allows to limit search of necessary information to rather small subset of documents.

*Keywords:* bag of words, machine learning, Bayesian classifier

## 1 Introduction

The Internet is full of information and every day we deal with constantly increasing volume of the processed and accumulated information that does a problem of text classification more and more urgent. For example, such necessity at automatic processing of a news stream and distribution of news texts by catalogs, Email Classification and Spam Filtering, recognition of emotional coloring of texts in documents. In all cases, the main idea is to assign the suitable category or label to each document that needs to be classified.

## 2 Goals

The main purpose of this work is creation of an automatic texts classification system. The idea is to build web-based platform by using well-known NLP tools and libraries for the reading and processing widely distributed and dynamic data for classification into different categories. For the construction of a system it is necessary to know the Naive Bayesian Classifier, sentiment analysis, moreover system will build on high-level web framework Django, and written in python language with different useful library such as OpenNLP, NLTK.

## 3 Feature Representation

Document representation and feature selection, it's a one of the most fundamental tasks which needs to be realized, preceding any classification problem. To understand better we use 'bag-of-words' (bag-of-features or bag-of-keypoints) [1] word-based document representation method, where we ignore the structure and sequence of words in the document. It is regularly used in methods of text classification where the frequency of each word is used as a feature for training a classifier. The feature vectors describe the words noticed in the documents. By this representation a document is considered to be simply a group of words which occur in it at least once. With this technique, it is possible to have thousands of words occurring in a small set of texts.

## 4 Classification Algorithm. Naive Bayes classifier

Nowadays we have different well-known often used machine learning algorithms for text classification such as decision trees, support - vector machines, naive Bayes classifier, nearest neighbor algorithms.

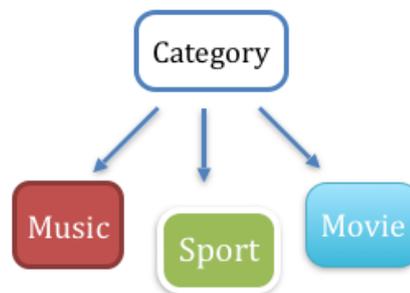


FIGURE 1 Classification by category

The naive Bayesian classifier [4] (NBC) is one of examples of use of methods of the vector analysis. This model of classification is based on a concept of conditional probability of the document  $d$  to a class  $c$ . NBC – one of the most often used classifiers, because of comparative simplicity in implementation and testing.

Basic formula of the naive Bayesian classifier (1):

$$P(A \setminus B) = \frac{P(B \setminus A)P(A)}{P(B)}. \quad (1)$$

For this model, the document is a vector:  $d = \{ w_1, w_2, \dots, w_n \}$ , where  $w_i$  - weight of  $i$ -th term, and  $n$  – size of the dictionary of selection. Thus, according to Bayes's theorem, probability of a  $c$  class for the document  $d$  will be (2):

$$P(c \setminus d) = \frac{P(d \setminus c)P(c)}{P(d)}. \quad (2)$$

Thus, for finding of the most probable class for the document  $d = \{ w_1, w_2, \dots, w_n \}$  by using naive Bayesian classifier, it is necessary to count conditional probabilities of accessory of the document  $d$  to each of the provided

classes separately and to choose the class having the maximum probability:

$$c_{NB} = \operatorname{argmax}_c [P(c_j) * \prod_i P(\omega_i | c_j)]. \quad (3)$$

```

classfier.py
1 import nltk
2 from nltk import word_tokenize, WordNetLemmatizer
3 from nltk.corpus import stopwords
4 from nltk import NaiveBayesClassifier, classify
5 from collections import Counter
6 import os
7 import random
8 import time
9 import io
10
11
12 lemmatizer = WordNetLemmatizer()
13
14 def getdataList(folder):
15     out_list = []
16     file_list = os.listdir(folder)
17     for fl in file_list:
18         try:
19             curr_file = io.open(folder + fl, "r", encoding="utf8")
20             out_list.append(curr_file.read())
21         except UnicodeDecodeError:
22             pass
23     return out_list
24
25
26 def preprocess(sentence):
27     def getFeatures(text, model=""):
28         tokens = word_tokenize(text)
29         lemmatized_tokens = [lemmatizer.lemmatize(token) for token in tokens]
30         return lemmatized_tokens
31     def train(features, samples, proportion):
32         train_data, test_data = train_test_split(features, samples, test_size=proportion, random_state=42)
33         stop_words = stopwords.words("english")
34         spam_emails = getdataList("data/aliya/spam/")
35         ham_emails = getdataList("data/aliya/ham/")
36         all_emails = [(email, "spam") for email in spam_emails] + [(email, "ham") for email in ham_emails]
37         assert(len(all_emails) == len(spam_emails) + len(ham_emails))
38         random.shuffle(all_emails)
    
```

There is a small problem connected to this formula. If a word never appears in the given training data, its relative frequency estimate will be zero. For the solution of this problem we applied the Laplace law of succession [2] to estimate  $P(\omega_i | c_j)$ . The estimate of the probability  $P(\omega_i | c_j)$  is given as:

$$P(\omega_i | c_j) = \frac{n_{ij} + 1}{n_j + k_j}, \quad (4)$$

where  $n_j$  is the total number of words in class  $c_j$ ,  $n_{ij}$  is the number of occurrences of word  $\omega_i$  in class  $c_j$  and  $k_j$  is the vocabulary size of class  $c_j$ . This is the result of the Bayesian estimation with a uniform prior assumption [3]. Bayesian method has the high speed of work and simplicity of mathematical models. This method is often used as a basic method when comparing various methods of machine learning.

### 5 Example

In order to better understand how naive Bayesian classifier works, let's assume that we are faced with a spam problem, and we need to create simple Spam filter which will attempt to classify incoming email messages in two different categories as 'spam', 'ham' (good, non-spam email) or 'unsure'. To solve this problem we will manually create data set with «emails», and determine bag of words which usually come in spam such as "sale", "buy now", "lowest prices", "click here» and etc. By counting probability of how many times spam words appears in email, system will identify, is it spam or ham. The program is written in Python with using NLTK libraries.

In result, we retrieve information about email in the table:

email	class
Sales! come and take your favorite boots	spam
Hi Mariya! Whats going on?	ham
Your application was declined. Please come tomorrow	ham
Viagra is available in our site. Order now	spam
Sale Sale Sale	spam
Buy now! only one day lowest prices	spam
Can you bring my book Mariya	ham
Hey hi! i'm going to celebrate my birthday, please come	ham

### 6 Conclusion

Nowadays the classification problem is one of the most fundamental problems in the machine learning and data mining. The amount of data that needs to be processed is very large, consequently, text analysis techniques must be

designed effectively to manage large numbers of elements with varying frequencies. In this research paper we mentioned work with classification algorithms, moreover have got acquainted with the principle of Bayesian classifier work and reviewed an example of simple spam filter where we categorized e-mails in two groups (spam and not spam).

### References

[1] Barakhnin V, Lukpanova L, Solovyov A 2014 An algorithm for constructing word forms using inflected classes for systems of morphological analysis of the Kazakh language *Bulletin of NGU, Series: Information technology* **12**(2) 25-31 (in Russian)

[2] Belonogov G, Zelenkov Y 1985 Algorithm for automatic analysis of Russian word *Questions of information theory and practice* **53** 62-93 (in Russian)

[3] Valiayeva T *The grammar of the Kazakh language* <http://kaz-tili.kz> (in Russian)

[4] Fedotov A, Tusupov D, Sambetbayeva M, Yerimbetova A, Bakiyeva A, Idrisova I 2015 The model determine the normal form of the word for the Kazakh language *Bulletin of NGU, Series: Information technology* **13**(1) 107-116 (in Russian)

[5] Porter M F 1980 *An algorithm for suffix stripping* *Program* **14**(3) 130-7

[6] Bakiyeva A *Program generation of word forms of the Kazakh language* [http://poem.ict.nsc.ru/~bakiyeva\\_aigerim/kazGen/](http://poem.ict.nsc.ru/~bakiyeva_aigerim/kazGen/) (in Russian)