# Gesture Recognition and its Application

Ziyang Wang
ShanghaiTech University
wangzy4@shanghaitech.edu.cn

Kefei Wu
ShanghaiTech University
wukf@shanghaitech.edu.cn

Donghao Zhao
ShanghaiTech University
zhdh@shanghaitech.edu.cn

Yiduo Gu
ShanghaiTech University
guyd@shanghaitech.edu.cn

## Abstract

*In this project we focus on how to successfully classify four kinds of SMSs including spams, notifications, verifications and normal messages with different features including sender and message information, and finally classify with five methods including Naive Bayes Classifier, Logistic Regression, Support Vector Machine, Random Forests and Gradient Boosting Decision Tree, to reach the goal to teach the program to classify automatically.*

## 1. Introduction

We receive many different kinds of text messages in our daily life, one of the simplest way to classify them of which is to divide them into spam and non-spam messages. We hope to further divide non-spam messages into more detailed categories. One of the widely used natural language processing task in different business problems is "Message Classification". The goal of Message classification is to automatically classify the text message into one or more defined categories.

In this project, the message classification is implemented by python. This project classifies messages based on python's machine learning library scikit-learn and the complete Chinese word segmentation tool jieba.

## 2. Motivation

As described above, the importance of filtering the useless messages and garbage information from the strangers, and the classifier that iOS system can only classify the messages according to the address book, which is naïve enough for classifying that may not work well to our expecting result. What we're going to do here is to do the real 'classify' that the application should understand more than what is garbage or not, but also the kind of it such as verifying code or relatives' communication, by analyzing the information given in the text and number, which will largely improve our using experience.

## 3. Table of Contents

Text Classification is an example of supervised machine learning task since a labelled dataset containing text documents and their labels is used for train a classifier. Our message classification project is composed of three main components:

1. Dataset preparation: The first step is the dataset preparation step which includes the process of loading a dataset and performing basic pre-processing. The dataset is then split into training and testing sets.

2. Feature engineering: The next step is the feature engineering in which the raw dataset is transformed into flat features which can be used in a machine learning model. This step also includes the process of creating new features from the existing data.

3. Model training: The final step is the Model Building step in which a machine learning model is trained on a labelled dataset.

4. Testing result analysis: the result of different features along with different models. We evaluate the performance and analyze the reason of it.

5. Possible improvement: In this article, we will also look at the different ways to improve the performance of message classifiers.

## 4. Dataset preparation

We've collected more than 2500 messages from our own mobile phones and manually labeled them with 4 categories: normal message, spam, verification code, and notification.

Our dataset reflects the real distribution of each categories and their features. features from our dataset.

Our dataset is stored in csv format.

# 5. Feature engineering

The next step is feature engineering. In this step, raw text data will be transformed into feature vectors and new features will be created using the existing dataset. We will implement the following different ideas in order to obtain relevant

## 5.1. Bag of Words Representation

First, we split each message into combination of words, and use each words that more than 2 characters as a local feature. We put local features of all message into a "vocabulary bag" as a dictionary.

For each message, we extract the local features in dictionary and count the number of the occurrences of each local feature to form a histogram. By $CountVectorizer$ of $sklearn$, we could store each message's histogram into a sparse matrix.



An example of bag of words representation.

## 5.2. TF-IDF vectors as features

Term Frequency(TF) is the number of times a given word appears in a message.Inverse Document Frequency(IDF) diminishes the weight of terms that occur very frequently in the document set and increases the weight of terms that occur rarely.

$$TF_t = \frac{t \ occurrences \ in \ a \ message}{Total \ number \ of \ words \ in \ a \ message} \quad (1)$$

$$IDF_t = log(\frac{Number \ of \ messages}{Messages \ include \ word \ t \ + 1}) \quad (2)$$

$$TF - IDF = TF * IDF \quad (3)$$

TF-IDF score represents the relative importance of a term in the document and the entire corpus.TF-IDF vectors tends to filter out common words and retain important words.

TF-IDF vectors can be generated at different levels of input tokens (words, characters, n-grams) by using $TfidfTransformer$ of $sklearn$ on bag of words representation matrix.

### 5.2.1 Word Level TF-IDF

Matrix representing tf-idf scores of every term in different documents

### 5.2.2 N-gram Level TF-IDF

N-grams are the combination of N terms together. This Matrix representing tf-idf scores of N-grams

### 5.2.3 Character Level TF-IDF

Matrix representing tf-idf scores of character level n-grams in the corpus

## 5.3. The length of the sender's number as features

In real life, normal messages are always sent by humans, whose phone numbers are usually 11 digits (or 13 digits when the country number is included in front), while notifications from banks or airplane companies are always from a 5-digit number.

Therefore, we use the length of sender's number as a feature outside the text.

## 5.4. Length Encoding

Inspired by label encoder, we use label 1, 2, 0 to label 5-digit number, 11(13)-digit number and other numbers.

```
def lengthencode(x):
    if x == 5:
        return 1
    elif x == 11:
        return 2
    elif x == 13:
        return 2
    return 0
```

2

## 6. Model building

The final step in the text classification framework is to train a classifier using the features created in the previous step. There are many different choices of machine learning models which can be used to train a final model. We will implement following different classifiers for this purpose:

### 6.1. Naive Bayes Classifier

Naive Bayes is a classification technique based on Bayes' Theorem, which is a simple technique for constructing classifiers: models that assign class labels to problem instances, represented as vectors of feature values, where the class labels are drawn from some finite set. There is not a single algorithm for training such classifiers, but a family of algorithms based on a common principle: all naive Bayes classifiers assume that the value of a particular feature is independent of the value of any other feature, given the class variable.

The multinomial Naive Bayes classifier is suitable for classification with discrete features (e.g., word counts for text classification). The multinomial distribution normally requires integer feature counts. However, in practice, fractional counts such as tf-idf may also work.

We've implemented a Naive Bayes classifier by ourselves, based on this fundamental principle.

The table below shows the relationship between arguments in Laplace smoothing equation and final results.

| Weight | Assumed Probability | Accuracy |
|--------|--------------------|----------|
| 1 | 0.5 | 88.12% |
| 1 | 0.25 | 88.51% |
| 0.5 | 0.25 | 89.27% |
| 0.5 | 0.5 | 88.51% |
| 2 | 0.5 | 87.35% |
| 2 | 0.25 | 87.36% |

**Bayes' theorem**:

$$P(A|B) = \frac{P(A) \times P(B|A)}{P(B)} \tag{4}$$

### 6.2. Logistic Regression

As the assumption of Naive Bayes classifier, Logistic regression also assume that the value of a particular feature is independent of the value of any other feature, given the class variable.

Given the sigmoid function:

$$g(z) = \frac{1}{1 + e^{-z}} \tag{5}$$

The probability of a message belong to a class given it's

feature $x$ and a parameter $\theta$ is:

$$
\begin{aligned}
P(Y = 1|x; \theta) &= g(\theta^T x) \\
&= \frac{1}{1 + e^{-\theta^T x}} \\
&= (g(\theta^T x_i))_i^y * (1 - g(\theta^T x_i))^{1-y_i}
\end{aligned} \tag{6}
$$

Joint distribution of all $n$ samples is :

$$
\begin{aligned}
L(\theta) &= \prod_{i=1}^{n} P(y_i|x_i, \theta) \\
&= \prod_{i=1}^{n} (g(\theta^T x_i))_i^y * (1 - g(\theta^T x_i))^{1-y_i}
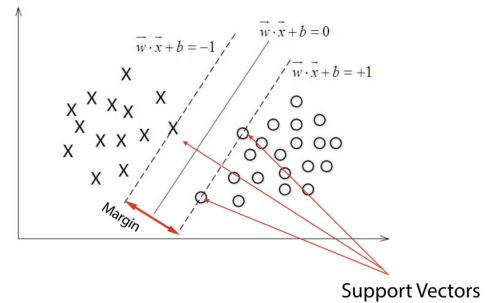\end{aligned} \tag{7}
$$

Take a log-likelihood function:

$$
\begin{aligned}
l(\theta) &= log(L(\theta)) \\
&= \sum_{i=1}^{n} y_i log(g(\theta^T x_i)) + (1 - y_i)log(1 - g(\theta^T x_i))
\end{aligned} \tag{8}
$$

The maximum likelihood estimation is the $\theta$ required to take $l(\theta)$ to the maximum value, which can be solved by the gradient ascending method.

### 6.3. Support Vector Machine

Support Vector Machine (SVM) is a supervised machine learning algorithm which can be used for both classification and regression challenges. The model extracts a best possible hyper-plane / line that maximizes the margin between the positive and negative examples.
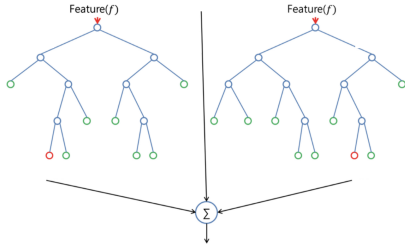


Positive: $\vec{w} \cdot \vec{x} + b \geq 1$
Negative: $\vec{w} \cdot \vec{x} + b \leq -1$
Support vector: $\vec{w} \cdot \vec{x} + b = \pm 1$

### 6.4. Random Forests

Random forests construct a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees to get a more accurate and stable prediction.

Random Forest Algorithm

## 6.5. Gradient Boosting Decision Tree

Boosting algorithm iteratively learning weak classifiers with respect to a distribution and adding them to a final strong classifier.

In gradient boosting, comparing to boosting algorithm, each new model is established to reduce the residual of the previous model to the gradient direction instead of adding with previous models.



Gradient Boosting Decision Tree

## 7. Result analysis

| Feature | Classifier | Accuracy |
|---|---|---|
| Bag of Words | Naive Bayes | 90.80% |
| Bag of Words | Logistic Regression | 92.72% |
| Bag of Words | SVM | 51.72% |
| Bag of Words | Random Forest | 89.27% |
| Bag of Words | Gradient Boosting Decision Tree | 90.04% |

| Feature | Classifier | Accuracy |
|---|---|---|
| TF-IDF | Naive Bayes | 89.27% |
| TF-IDF | Logistic Regression | 90.42% |
| TF-IDF | SVM | 51.72% |
| TF-IDF | Random Forest | 91.18% |
| TF-IDF | Gradient Boosting Decision Tree | 90.80% |

From the result above, we can know that TF-IDF and Bag of Words basically has no difference, and TF-IDF' s results is worse than Bag of Words feature.

| Feature | Classifier | Accuracy |
|---|---|---|
| Bag of Words+Sender | Naive Bayes | 91.18% |
| Bag of Words+Sender | Logistic Regression | 93.10% |
| Bag of Words+Sender | SVM | 51.72% |
| Bag of Words+Sender | Random Forest | 88.89% |
| Bag of Words+Sender | Gradient Boosting Decision Tree | 90.80% |

| Feature | Classifier | Accuracy |
|---|---|---|
| Bag of Words+Sender with encoding | Naive Bayes | 90.80% |
| Bag of Words+Sender with encoding | Logistic Regression | 93.10% |
| Bag of Words+Sender with encoding | SVM | 51.72% |
| Bag of Words+Sender with encoding | Random Forest | 89.66% |
| Bag of Words+Sender with encoding | Gradient Boosting Decision Tree | 90.42% |

We found that use sender's number as feature will help to improve training results, but encode sender's number doesn't help at all. We guess maybe the following reasons: 1. Our dataset is lack of normal message (which meets real condition). 2. Banks sends too much spam messages.

| Feature | Classifier | Dataset | Accuracy |
|---|---|---|---|
| Bag of Words | Naive Bayes | 1600 | 87.89% |
| Bag of Words | Logistic Regression | 1600 | 92.99% |
| Bag of Words | Naive Bayes | 2600 | 90.80% |
| Bag of Words | Logistic Regression | 2600 | 92.72% |

After training with a smaller dataset, we found that even if logistic regression performs better in these two dataset, Naive Bayes performs better with the dataset increasing. So we think Naive Bayes might have better performance than logistic regression when the dataset is larger.

We also realize our own Naive Bayes function and Logistic Regression classifier by ourselves, due to some reasons relating to the smoothing function arguments, but which mostly has reached our expectations.

| Feature | Classifier | Dataset | Accuracy |
|---|---|---|---|
| Bag of Words | Naive Bayes (by ourselves) | 1600 | 85.35% |
| Bag of Words | Naive Bayes (by ourselves) | 1600 | 88.12% |
| Bag of Words | Logistic Regression (by ourselves) | 1600 | 33.14% |

## 8. Possible improvement

While the above framework can be applied to message classification problems, but to achieve a good accuracy some improvements can be done in the overall framework.

### 8.1. Text Cleaning

text cleaning can help to reduce the noise present in text data in the form of stopwords, punctuations marks, suffix variations etc.

### 8.2. NLP features with text feature vectors

In the feature engineering section, we generated a number of different feature vectors, combining them together can help to improve the accuracy of the classifier.

### 8.3. Ensemble Models

Stacking different models and blending their outputs can help to further improve the results.

We have envisioned it, but it still needs further verification.

# References

1.Unsupervised learning of finite mixture models Figueiredo, MAT (Figueiredo, MAT); Jain, AK (Jain, AK) IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE 24 3 381-396

2.Short Text Classification in Twitter to Improve Information Filtering Sriram, B (Sriram, Bharath)[ 1 ] ; Fuhry, D (Fuhry, David)[ 1 ] ; Demir, E (Demir, Engin)[ 1 ] ; Ferhatosmanoglu, H (Ferhatosmanoglu, Hakan)[ 1 ] ; Demirbas, M (Demirbas, Murat) SIGIR 2010: PROCEEDINGS OF THE 33RD ANNUAL INTERNATIONAL ACM SIGIR CONFERENCE ON RESEARCH DEVELOPMENT IN INFORMATION RETRIEVAL 841-842

3.Sentiment Analysis of Short Informal Text Kiritchenko, S (Kiritchenko, Svetlana)[ 1 ] ; Zhu, XD (Zhu, Xiaodan)[ 1 ] ; Mohammad, SM (Mohammad, Saif M.)[ 1 ] JOURNAL OF ARTIFICIAL INTELLIGENCE RESEARCH 50 723-762

4.Rough set-aided keyword reduction for text categorization Chouchoulas, A (Chouchoulas, A); Shen, Q (Shen, Q) APPLIED ARTIFICIAL INTELLIGENCE 15 9 843-873

5.Gradient Boost Decision Tree OR Treelink

6.Baidu Grammar analysis     7. Jie Ba analysis