

Computational Social Science Methods in R

Text Classification

Taehee Kim

Summer Semester 2022

University of Oldenburg

Tokenization

- Segmenting an input stream into an ordered sequence of tokens is called tokenization.
- A system which splits texts into word tokens is called a tokenizer.
- Example:

Input text: "Tony likes pizza and Mike likes pasta."

output: {"Tony", "likes", "pizza", "and", "Mike", "likes", "pasta", "."}

Text Normalization: Stemming and Lemmatization

- different forms of a word: e.g., organize, organizes, and organizing
- The goal of both stemming and lemmatization is to reduce inflectional and derivationally related forms of a word to a common base form
- Example: 'the boy's cars are different colors' → the boy car be differ color

¹Tokenization:

<https://nlp.stanford.edu/IR-book/html/htmledition/tokenization-1.html>, Text

Normalization: <https://nlp.stanford.edu/IR-book/html/htmledition/stemming-and-lemmatization-1.html>

<https://nlp.stanford.edu/IR-book/html/htmledition/stemming-and-lemmatization-1.html>

Stemming

- Stemming: usually refers to **chop off the ends of words** (often includes the removal of derivational affixes)
- beautifully → beauti

Lemmatization

- Lemmatization: usually refers to doing things properly with the use of a vocabulary and morphological analysis of words, normally aiming to remove inflectional endings only and to **return the base or dictionary form** of a word, which is known as the lemma.
- beautifully → beautifully

Stop words

- refers common words in a language so that does not contain important significance to be used in Search Queries.
- in English: as, the, is, are, with etc..

German

- Python NLTK package → **install**
- **punkt** in modules install

Goal

- Assign 'texts' from a universe to two or more classes or categories

Applications

- Spam vs. Non-spam email
- Email filtering
- News events tracked and filtered by topics
- Journal articles indexed by subject categories

Who wrote which Federalist papers?

- 1987: anonymous essays try to convince New York to ratify U.S. Constitution
- Jay, Madison, and Hamilton
- Authorship of 12 of the letters in dispute
- 1963: solved by Mosteller and Wallace using Bayesian methods

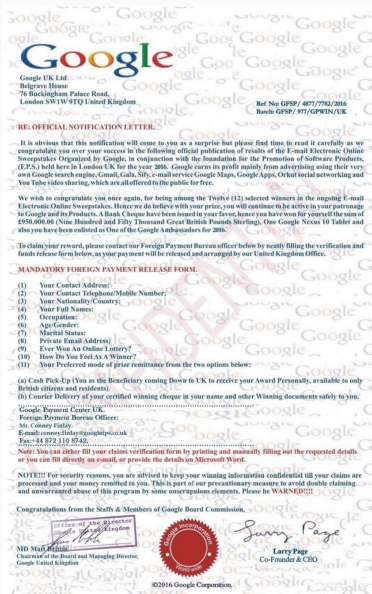


Figure 1: Example of a spam email

Hand-coded rules

- Rules based on combinations of words or other features
 - Spam mail: black-list-address OR ("winning" AND "payment")
- Accuracy can be high if rules are carefully refined by experts
- However, building and maintaining these rules are expensive

Input

- Document D
- A fixed set of classes C
- A training set (hand-labeled documents)

Output

- A learned classifier

- One of the most important classifiers
- Relies on **bag of words**

Bag of words approach

- does not consider the order of words
- but **the counts of individual words**
- works for multiple classification

Example

I will vote for XXX party! I agree with their ideas and support their policies. Not only they are sophisticated and smart but also sincere PATRIOT!

I will **vote** for XXX party! I **agree** with their ideas and **support** their policies. Not only they are **sophisticated** and **smart** but also sincere PATRIOT!

Bayes Rule

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

Text Classification Naive Bayes (Multinomial)

$P(c|d) = \frac{P(d|c)P(c)}{P(d)}$, where **d**: a document and **c**: a class

We compute $P(c_i|d)$, where $i=1,\dots,n$ (number of class)
and adopt c_i which maximize $P(c_i|d)$

²Recommended reading: Raschka, Sebastian. "Naive Bayes and Text Classification I– Introduction and Theory." arXiv preprint arXiv:1410.5329 (2014)

$$P(c_i|d) = \frac{P(d|c_i)P(c_i)}{P(d)}$$

$$P(d|c_i)P(c_i)$$

$$= P(w_1, w_2, \dots, w_n|c_i)P(c_i)$$

$$= P(w_1|c_i)P(w_2|c_i)\dots P(w_n|c_i)P(c_i), \text{ where } w: \text{ word, } n: \text{ number of words in } d$$

- $P(c_i) = \frac{(d, c_i)}{\sum (d, c)} = \frac{\text{number of doc. in class } i}{\text{number of all doc.}}$

- $P(w_j|c_i) = \frac{(w_j, c_i)}{\sum (w, c_i)} = \frac{\text{number of word } j \text{ in class } i}{\text{number of all words in class } i}$

Laplace Smoothing

$$P(w_j|c_i) = \frac{(w_j, c_i) + 1}{\sum (w, c_i) + |V|}$$

V : number of all unique words appeared in training documents

Training:

- dog, cat, cat \rightarrow class A
- dog, dog, bird \rightarrow class A
- dog, fox \rightarrow class A
- fish, dolphin, dog \rightarrow class B

Test

- dog, dog, dog, fish, dolphin \rightarrow ???

Support Vector Machines (SVM)

- SVM is one of the popular machine learning algorithms since it shows good performance in a diverse types of data, including text
- it determines the best decision boundary between vectors that belong to a given class
- SVM shows good performance with data that have lots of features → such as text data
- most text classification problems are linearly separable, and linear Kernel shows good performance for data with a lot of features, such as text data

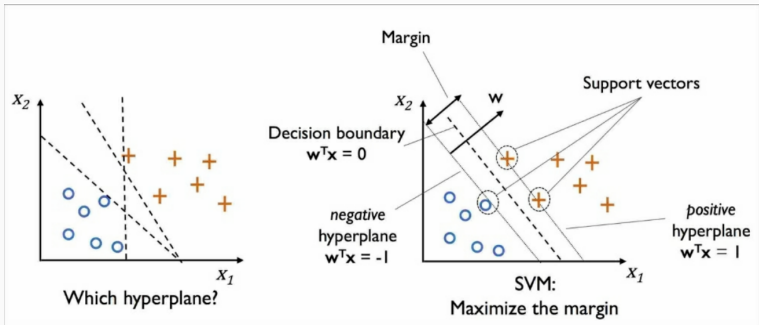
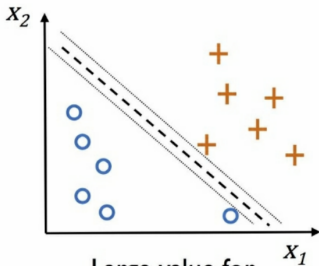
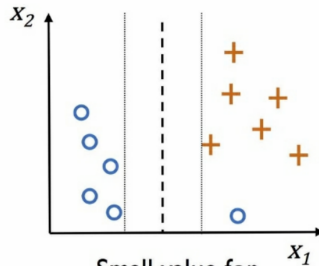


Figure 2: Support Vector Machine



Large value for
parameter C



Small value for
parameter C

Figure 2: Support Vector Machine

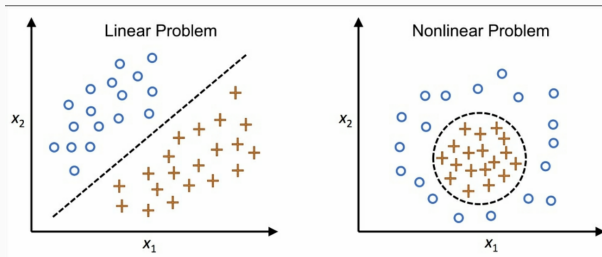


Figure 3: Linear and Nonlinear Problem

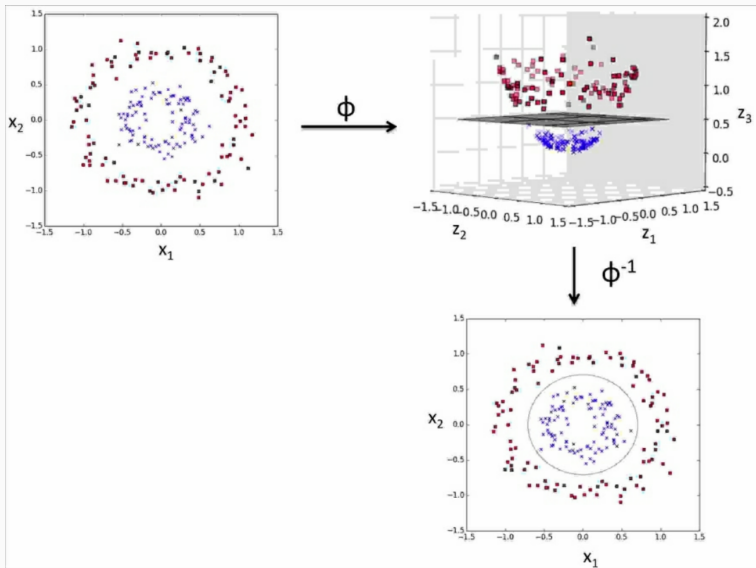


Figure 3: Linear and Nonlinear Problem

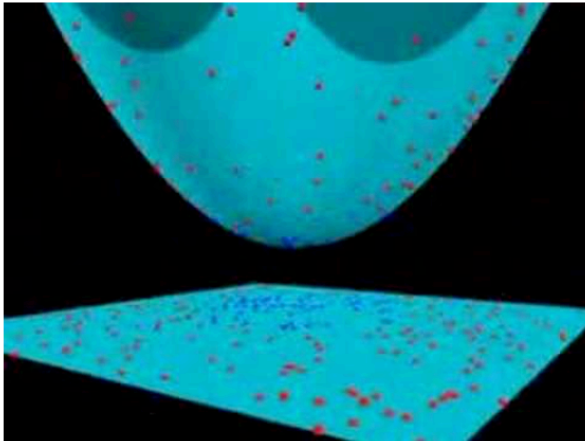


Figure 4: <https://www.youtube.com/watch?v=3liCbRZPrZA>

Text vectorization

- transform text into a vector
- define a fixed length vector where each entry corresponds to a word appeared in training texts.
- vector entries: count a word appeared in the text

Example

- Suppose that we have following words in training texts
(I, do, not, am, like, happy, pizza)
- We vectorize the text "I like pizza"
(1, 0, 0, 0, 1, 0, 1)

Term frequency–inverse document frequency(tf-idf)

- numerical statistic which represent the importance of a word
- it is often used as a weight

Term frequency

- The term frequency $tf_{t,d}$ of term t in document d is defined as the number of times that t occurs in d
- A document with 10 occurrences of the term is more relevant than a document with one occurrence of the term (but the relevance does not increase proportionally).

Document frequency

- Are rare terms or frequent terms more informative? → rare terms
- we want to a high weight for rare terms
- df_t is the document frequency, the number of documents that t occurs in, N is the total number of documents

Inverse document frequency

- a measure of the informativeness of the term.
- $idf_t = \log_{10} \frac{N}{df_t}$

TF-IDF

- We assign a tf-idf weight for each term t in each document d :
- $w_{t,d} = (1 + \log tf_{t,d}) \cdot \log \frac{N}{df_t}$

N : total number of documents

Increases with the number of occurrences within a document

Increases with the rarity of the term in the collection