# Information Retrieval
## Text classification: Naive Bayes

Hamid Beigy

Sharif university of technology

November 5, 2018

# Table of contents

## Introduction

1 How would you write a program that would automatically detect and delete this type of message?

```
From: ''' <takworlld@hotmail.com>
Subject: real estate is the only way... gem  oalvgkay
Anyone can buy real estate with no money down
Stop paying rent TODAY !
There is no need to spend hundreds or even thousands for similar courses
I am 22 years old and I have already purchased 6 properties using the
methods outlined in this truly INCREDIBLE ebook.
Change your life NOW !
=================================================
Click Below to order:
http://www.wholesaledaily.com/sales/nmd.htm
=================================================
```

# Examples of how search engines use classification

1. Query classification (types of queries)
2. Spelling correction
3. Document/webpage classification
4. Automatic detection of spam pages (spam vs. non-spam)
5. Topic classification (relevant to topic vs. not)
6. Language identification (classes: English vs. French etc.)
7. User classification (personalised search)

# Classification terminology

1. A document space $\mathbb{X}$
   Documents are represented in this space (typically some type of high-dimensional space).

2. A fixed set of classes $\mathbb{C} = \{c_1, c_2, \ldots, c_J\}$
   The classes are human-defined for the needs of an application (e.g., spam vs. nonspam).

3. A training set $\mathbb{D}$ of labeled documents.
   Each labeled document $\langle d, c \rangle \in \mathbb{X} \times \mathbb{C}$

4. Using a learning method or learning algorithm, we then wish to learn a classifier $\gamma$ that maps documents to classes:

$$\gamma : \mathbb{X} \to \mathbb{C}$$

5. Classification task: Given: a description $d \in \mathbb{X}$ of a document
   Determine: $\gamma(d) \in \mathbb{C}$, that is, the class that is most appropriate for $d$.

# Classification terminology

1. Features: measurable properties of the data.
2. Classes: labels associated with the data.
3. Consider the following example
   - Sentiment classification: automatically classify text based on the sentiment it contains (e.g., movie reviews).
   - Features: the words the text contains, parts of speech, grammatical constructions etc.
   - Classes: positive or negative sentiment (binary classification).
4. Classification is the function that maps input features to a class.

# Classification terminology

1. Consider a text classification with six classes {UK, China, poultry, coffee, elections, sports}
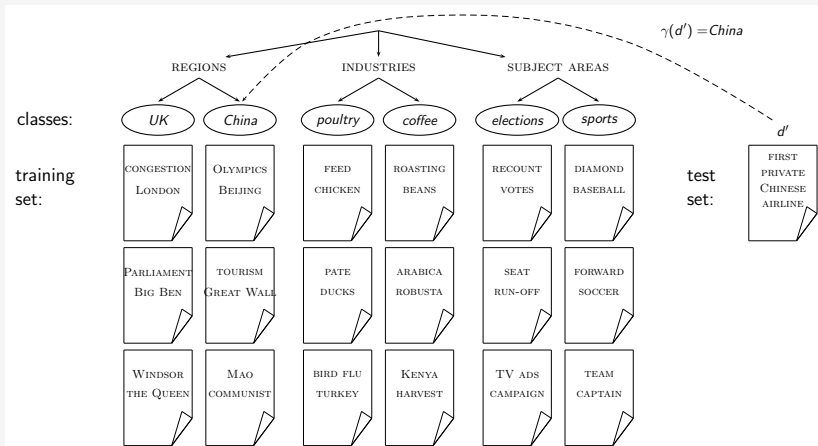
# Table of contents

# Classification methods: Manual

1. Manual classification was used by Yahoo in the beginning of the web. Also: PubMed, ODP
2. Very accurate if job is done by experts
3. Consistent when the problem size and team is small
4. Scaling manual classification is difficult and expensive.
5. Hence, we need automatic methods for classification.

# Classification methods: Rule-based

1. E.g., Google Alerts is rule-based classification.
2. There are IDE-type development environments for writing very complex rules efficiently. (e.g., Verity)
3. Often: Boolean combinations (as in Google Alerts)
4. Accuracy is very high if a rule has been carefully refined over time by a subject expert.
5. Building and maintaining rule-based classification systems is cumbersome and expensive.
6. Email classification in email clients such as outlook.

# Classification methods: Statistical/Probabilistic

1. This was our definition of the classification problem – text classification as a learning problem

2. (i) Supervised learning of a the classification function $\gamma$ and (ii) application of $\gamma$ to classifying new documents

3. We will look at two methods for doing this: Naive Bayes and SVMs

4. No free lunch: requires hand-classified training data

5. But this manual classification can be done by non-experts.

# Table of contents

## Bayes classifier

1. We compute the probability of a document $d$ being in a class $c$ as follows:

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

$$P(c|d) \propto P(c)P(d|c)$$

2. $P(d)$ is constant during a given classification and wont affect the result.

# Naive Bayes classifier

1. The Naive Bayes classifier is a probabilistic classifier.

$$P(c|d) \propto P(c) \prod_{1 \le k \le n_d} P(t_k|c)$$

2. $n_d$ is the length of the document. (number of tokens)
3. $P(t_k|c)$ is the conditional probability of term $t_k$ occurring in a document of class $c$
4. $P(t_k|c)$ as a measure of how much evidence $t_k$ contributes that $c$ is the correct class.
5. $P(c)$ is the prior probability of $c$.
6. If a document's terms do not provide clear evidence for one class vs. another, we choose the $c$ with highest $P(c)$.

# Maximum a posteriori class

1. Our goal in Naive Bayes classification is to find the "best" class.
2. The best class is the most likely or maximum a posteriori (MAP) class $c_{\mathrm{map}}$:

$$c_{\mathrm{map}} = \arg\max_{c \in \mathbb{C}} \hat{P}(c|d) = \arg\max_{c \in \mathbb{C}} \hat{P}(c) \prod_{1 \leq k \leq n_d} \hat{P}(t_k|c)$$

## Maximum a posteriori class

1. Multiplying lots of small probabilities can result in floating point underflow.
2. Since $\log(xy) = \log(x) + \log(y)$, we can sum log probabilities instead of multiplying probabilities.
3. Since log is a monotonic function, the class with the highest score does not change.
4. So what we usually compute in practice is:

$$c_{\mathrm{map}} = \arg\max_{c \in \mathbb{C}} [\log \hat{P}(c) + \sum_{1 \leq k \leq n_d} \log \hat{P}(t_k|c)]$$

# Naive Bayes classifier

1 Classification rule:

$$c_{\mathsf{map}} = \underset{c \in \mathbb{C}}{\arg\max} \left[ \log \hat{P}(c) + \sum_{1 \le k \le n_d} \log \hat{P}(t_k|c) \right]$$

2 Simple interpretation:

- Each conditional parameter $\log \hat{P}(t_k|c)$ is a weight that indicates how good an indicator $t_k$ is for $c$.
- The prior $\log \hat{P}(c)$ is a weight that indicates the relative frequency of $c$.
- The sum of log prior and term weights is then a measure of how much evidence there is for the document being in the class.
- We select the class with the most evidence.

## Parameter estimation

1. Estimate parameters $\hat{P}(c)$ and $\hat{P}(t_k|c)$ from train data: How?
2. Prior:

$$\hat{P}(c) = \frac{N_c}{N}$$

3. $N_c$: number of docs in class $c$; $N$: total number of docs
4. Conditional probabilities:

$$\hat{P}(t|c) = \frac{T_{ct}}{\sum_{t' \in V} T_{ct'}}$$

5. $T_{ct}$ is the number of tokens of $t$ in training documents from class $c$ (includes multiple occurrences)
6. We've made a Naive Bayes independence assumption here.

# Maximum likelihood

1. Is it possible to compute $\hat{P}(c)$? (why?)
2. In some cases, we consider $\hat{P}(c)$ to be equal for all documents.
3. Hence, we only compute $\hat{P}(t|c)$
4. Conditional probabilities:

$$\hat{P}(t|c) = \frac{T_{ct}}{\sum_{t' \in V} T_{ct'}}$$

5. $T_{ct}$ is the number of tokens of $t$ in training documents from class $c$ (includes multiple occurrences)
6. We've made a Naive Bayes independence assumption here.
7. Classification rule:

$$c_{\text{map}} = \arg\max_{c \in \mathbb{C}} \sum_{1 \le k \le n_d} \log \hat{P}(t_k|c)$$

## Add-one smoothing

**1** Before:

$$\hat{P}(t|c) = \frac{T_{ct}}{\sum_{t' \in V} T_{ct'}}$$

**2** Now: Add one to each count to avoid zeros:

$$\hat{P}(t|c) = \frac{T_{ct} + 1}{\sum_{t' \in V}(T_{ct'} + 1)} = \frac{T_{ct} + 1}{(\sum_{t' \in V} T_{ct'}) + B}$$

**3** $B$ is the number of bins – in this case the number of different words or the size of the vocabulary $|V| = M$

# Naive Bayes classifier: Interpretation

**1** Classification rule:

$$c_{\mathrm{map}} = \arg\max_{c \in \mathbb{C}} \left[ \log \hat{P}(c) + \sum_{1 \leq k \leq n_d} \log \hat{P}(t_k | c) \right]$$

**2** Simple interpretation:

- Each conditional parameter $\log \hat{P}(t_k | c)$ is a weight that indicates how good an indicator term $t_k$ is for class $c$.
- The prior $\log \hat{P}(c)$ is a weight that indicates how likely we are to see class $c$.
- The sum of log prior and term weights is then a measure of how much evidence there is for the document being in the class.
- We select the class with the most evidence.

# Time complexity of Naive Bayes

1. Time complexity is

   | mode | time complexity |
   |------|-----------------|
   | training | $\Theta(|\mathbb{D}|L_{\text{ave}} + |\mathbb{C}||V|)$ |
   | testing | $\Theta(L_{\text{a}} + |\mathbb{C}|M_{\text{a}}) = \Theta(|\mathbb{C}|M_{\text{a}})$ |

2. $L_{\text{ave}}$: average length of a training doc, $L_{\text{a}}$: length of the test doc, $M_{\text{a}}$: number of distinct terms in the test doc, $\mathbb{D}$: training set, $V$: vocabulary, $\mathbb{C}$: set of classes

3. $\Theta(|\mathbb{D}|L_{\text{ave}})$ is the time it takes to compute all counts.

4. $\Theta(|\mathbb{C}||V|)$ is the time it takes to compute the parameters from the counts.

5. Generally: $|\mathbb{C}||V| < |\mathbb{D}|L_{\text{ave}}$

6. Test time is also linear (in the length of the test document).

7. Thus: Naive Bayes is linear in the size of the training set (training) and the test document (testing). This is optimal.

## An example

1. Consider the following dataset.

|  | docID | words in document | in $c = China$? |
|---|---|---|---|
| training set | 1 | Chinese Beijing Chinese | yes |
|  | 2 | Chinese Chinese Shanghai | yes |
|  | 3 | Chinese Macao | yes |
|  | 4 | Tokyo Japan Chinese | no |
| test set | 5 | Chinese Chinese Chinese Tokyo Japan | ? |

2. we have

$$\hat{P}(c) = \frac{N_c}{N}$$

$$\hat{P}(t|c) = \frac{T_{ct} + 1}{\sum_{t' \in V}(T_{ct'} + 1)} = \frac{T_{ct} + 1}{(\sum_{t' \in V} T_{ct'}) + B}$$

$B$ is the size of the vocabulary $|V| = M$.

$$c_{\text{map}} = \arg\max_{c \in \mathbb{C}} \left[ \hat{P}(c) \cdot \prod_{1 \leq k \leq n_d} \hat{P}(t_k|c) \right]$$

## An example: training

1. Priors: $\hat{P}(c) = \frac{3}{4}$ and $\hat{P}(\overline{c}) = \frac{3}{4}$ Conditional probabilities:

$$
\begin{aligned}
\hat{P}(\text{CHINESE}|c) &= (5+1)/(8+6) = 6/14 = 3/7 \\
\hat{P}(\text{TOKYO}|c) = \hat{P}(\text{JAPAN}|c) &= (0+1)/(8+6) = 1/14 \\
\hat{P}(\text{CHINESE}|\overline{c}) &= (1+1)/(3+6) = 2/9 \\
\hat{P}(\text{TOKYO}|\overline{c}) = \hat{P}(\text{JAPAN}|\overline{c}) &= (1+1)/(3+6) = 2/9
\end{aligned}
$$

2. The denominators are $(8+6)$ and $(3+6)$ because the lengths of $text_c$ and $text_{\overline{c}}$ are 8 and 3, respectively, and because the constant $B$ is 6 as the vocabulary consists of six terms.

## An example: classification

1. For classification, we have

$$\hat{P}(c|d_5) \propto \frac{3}{4} \cdot (3/7)^3 \cdot 1/14 \cdot 1/14 \approx 0.0003$$
$$\hat{P}(\bar{c}|d_5) \propto 1/4 \cdot (2/9)^3 \cdot 2/9 \cdot 2/9 \approx 0.0001$$

2. Thus, the classifier assigns the test document to $c = China$.

3. The reason for this classification decision is that the three occurrences of the positive indicator CHINESE in $d_5$ outweigh the occurrences of the two negative indicators JAPAN and TOKYO.

# Evaluating classification

1. Evaluation must be done on test data that are independent of the training data, i.e., training and test sets are disjoint.
2. It's easy to get good performance on a test set that was available to the learner during training (e.g., just memorize the test set).
3. Measures: Precision, recall, $F_1$, classification accuracy

## Evaluating classification

1. Precision, recall, and $F_1$ can be calculated using

|  | in the class | not in the class |
|---|---|---|
| predicted to be in the class | true positives (TP) | false positives (FP) |
| predicted to not be in the class | false negatives (FN) | true negatives (TN) |

2. TP, FP, FN, TN are counts of documents. The sum of these four counts is the total number of documents.

$$\begin{aligned} \text{precision:} P &= TP/(TP + FP) \\ \text{recall:} R &= TP/(TP + FN) \end{aligned}$$

3. $F_1$ allows us to trade off precision against recall.

$$F_1 = \frac{1}{\frac{1}{2}\frac{1}{P} + \frac{1}{2}\frac{1}{R}} = \frac{2PR}{P + R}$$

4. This is the harmonic mean of $P$ and $R$: $\frac{1}{F} = \frac{1}{2}(\frac{1}{P} + \frac{1}{R})$

# Table of contents

# Reading

Please read chapter 13.1 – 13.4 of Information Retrieval Book.