# Text Classification and Naïve Bayes

## The Task of Text Classification

Many slides are adapted from slides by Dan Jurafsky

# Is this spam?

**Subject:** **Important notice!**

**From:** Stanford University <newsforum@stanford.edu>

**Date:** October 28, 2011 12:34:16 PM PDT

**To:** undisclosed-recipients:;

---

Greats News!

You can now access the latest news by using the link below to login to Stanford University News Forum.

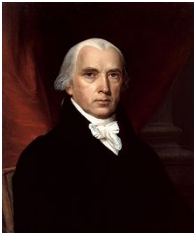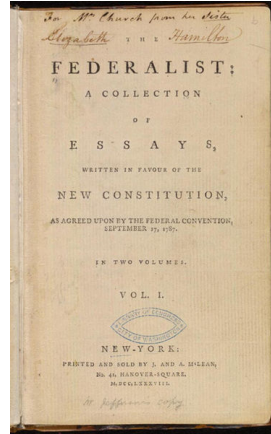http://www.123contactform.com/contact-form-StanfordNew1-236335.html

Click on the above link to login for more information about this new exciting forum. You can also copy the above link to your browser bar and login for more information about the new services.
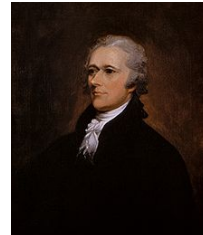
© Stanford University. All Rights Reserved.

# Who wrote which Federalist papers?

- 1787-8: anonymous essays try to convince New York to ratify U.S Constitution:  Jay, Madison, Hamilton.

- Authorship of 12 of the letters in dispute

- 1963: solved by Mosteller and Wallace using Bayesian methods

James Madison

Alexander Hamilton

# Male or female author?

1. By 1925 present-day Vietnam was divided into three parts under French colonial rule. The southern region embracing Saigon and the Mekong delta was the colony of Cochin-China; the central area with its imperial capital at Hue was the protectorate of Annam...

2. Clara never failed to be astonished by the extraordinary felicity of her own name. She found it hard to trust herself to the mercy of fate, which had managed over the years to convert her greatest shame into one of her greatest assets...

S. Argamon, M. Koppel, J. Fine, A. R. Shimoni, 2003. "Gender, Genre, and Writing Style in Formal Written Texts," Text, volume 23, number 3, pp. 321–346

# Positive or negative movie review?

- 👎 unbelievably disappointing
- 👍 Full of zany characters and richly applied satire, and some great plot twists
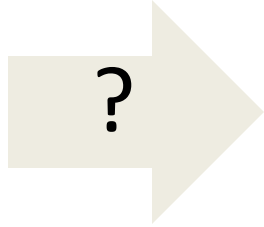- 👍 this is the greatest screwball comedy ever filmed
- 👎 It was pathetic. The worst part about it was the boxing scenes.

# What is the subject of this article?

## MEDLINE Article



?

## MeSH Subject Category Hierarchy

- Antogonists and Inhibitors
- Blood Supply
- Chemistry
- Drug Therapy
- Embryology
- Epidemiology
- …

6

# Text Classification

- Assigning subject categories, topics, or genres
- Spam detection
- Authorship identification
- Age/gender identification
- Language Identification
- Sentiment analysis
- …

# Text Classification: definition

- *Input*:
  - a document $d$
  - a fixed set of classes $C = \{c_1, c_2, ..., c_J\}$


- *Output*: a predicted class $c \in C$

# Classification Methods:
# Hand-coded rules

- Rules based on combinations of words or other features
  - spam: black-list-address OR ("dollars" AND"have been selected")
- Accuracy can be high
  - If rules carefully refined by expert
- But building and maintaining these rules is expensive

# Classification Methods:
# Supervised Machine Learning

- *Input:*

  - a document $d$

  - a fixed set of classes $C = \{c_1, c_2, ..., c_J\}$

  - A training set of $m$ hand-labeled documents $(d_1, c_1), ...., (d_m, c_m)$

- *Output:*

  - a learned classifier $\gamma: d \rightarrow c$

# Classification Methods:
# Supervised Machine Learning

- Any kind of classifier
  - Naïve Bayes
  - Logistic regression, maxent
  - Support-vector machines
  - k-Nearest Neighbors

  - …

# Text Classification and Naïve Bayes

## The Task of Text Classification

# Text Classification and Naïve Bayes

Text Classification: Evaluation

# The 2-by-2 contingency table

|  | correct | not correct |
|---|---|---|
| selected | tp | fp |
| not selected | fn | tn |

# Precision and recall

- **Precision**: % of selected items that are correct
  **Recall**: % of correct items that are selected

|  | correct | not correct |
|---|---|---|
| selected | tp | fp |
| not selected | fn | tn |

# A combined measure: F

- A combined measure that assesses the P/R tradeoff is F measure (weighted harmonic mean):

$$F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha)\frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

- People usually use balanced F1 measure
  - i.e., with $\beta = 1$ (that is, $\alpha = \frac{1}{2}$):  **F = 2PR/(P+R)**

# Confusion matrix c

- For each pair of classes $<c_1, c_2>$ how many documents from $c_1$ were incorrectly assigned to $c_2$?
  - $c_{3,2}$: 90 wheat documents incorrectly assigned to poultry

| Docs in test set | Assigned UK | Assigned poultry | Assigned wheat | Assigned coffee | Assigned interest | Assigned trade |
|---|---|---|---|---|---|---|
| True UK | 95 | 1 | 13 | 0 | 1 | 0 |
| True poultry | 0 | 1 | 0 | 0 | 0 | 0 |
| True wheat | 10 | 90 | 0 | 1 | 0 | 0 |
| True coffee | 0 | 0 | 0 | 34 | 3 | 7 |
| True interest | - | 1 | 2 | 13 | 26 | 5 |
| True trade | 0 | 0 | 2 | 14 | 5 | 10 |

# Per class evaluation measures

**Recall**:

Fraction of docs in class $i$ classified correctly:

$$\frac{c_{ii}}{\sum_j c_{ij}}$$

**Precision**:

Fraction of docs assigned class $i$ that are actually about class $i$:

$$\frac{c_{ii}}{\sum_j c_{ji}}$$

**Accuracy**: (1 - error rate)

Fraction of docs classified correctly:

$$\frac{\sum_i c_{ii}}{\sum_j \sum_i c_{ij}}$$

18

# Micro- vs. Macro-Averaging

- If we have more than one class, how do we combine multiple performance measures into one quantity?

- **Macroaveraging**: Compute performance for each class, then average. Average on classes

- **Microaveraging**: Collect decisions for each instance from all classes, compute contingency table, evaluate. Average on instances

# Micro- vs. Macro-Averaging: Example

| Class 1 | | |
|---|---|---|
| | Truth: yes | Truth: no |
| Classifier: yes | 10 | 10 |
| Classifier: no | 10 | 970 |

| Class 2 | | |
|---|---|---|
| | Truth: yes | Truth: no |
| Classifier: yes | 90 | 10 |
| Classifier: no | 10 | 890 |

| Micro Ave. Table | | |
|---|---|---|
| | Truth: yes | Truth: no |
| Classifier: yes | 100 | 20 |
| Classifier: no | 20 | 1860 |

- Macroaveraged precision: (0.5 + 0.9)/2 = 0.7

- Microaveraged precision: 100/120 = .83

- Microaveraged score is dominated by score on common classes

# Development Test Sets and Cross-validation

Training set     Development Test Set     Test Set

| Training Set | Dev Test | |
|---|---|---|

| Training Set | | Dev Test |
|---|---|---|

| Dev Test | Training Set | |
|---|---|---|

- Metric: P/R/F1  or Accuracy
- Unseen test set
  - avoid overfitting ('tuning to the test set')
  - more conservative estimate of performance
- Cross-validation over multiple splits
  - Handle sampling errors from different datasets
  - Pool results over each split
  - Compute pooled dev set performance

Test Set

# Text Classification and Naïve Bayes

## Text Classification: Evaluation

# Text Classification and Naïve Bayes

Formalizing the Naïve Bayes Classifier

# Naïve Bayes Intuition

- Simple ("naïve") classification method based on Bayes rule
- Relies on very simple representation of document
  - Bag of words

# Bayes' Rule Applied to Documents and Classes

- For a document *d* and a class *c*

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

# Naïve Bayes Classifier (I)

$$c_{MAP} = \underset{c \in C}{\operatorname{argmax}} P(c \mid d)$$

MAP is "maximum a posteriori" = most likely class

$$= \underset{c \in C}{\operatorname{argmax}} \frac{P(d \mid c)P(c)}{P(d)}$$

Bayes Rule

$$= \underset{c \in C}{\operatorname{argmax}} P(d \mid c)P(c)$$

Dropping the denominator

# Naïve Bayes Classifier (II)

$$c_{MAP} = \underset{c \in C}{\operatorname{argmax}} P(d \mid c)P(c)$$

$$= \underset{c \in C}{\operatorname{argmax}} P(x_1, x_2, \ldots, x_n \mid c)P(c)$$

Document d represented as features x1..xn

# Naïve Bayes Classifier (III)

$$c_{MAP} = \underset{c \in C}{\operatorname{argmax}} P(x_1, x_2, \ldots, x_n \mid c) P(c)$$

O($|X|^n \cdot |C|$) parameters

How often does this class occur?

Could only be estimated if a very, very large number of training examples was available.

We can just count the relative frequencies in a corpus

# The bag of words representation

$$\gamma \left( \begin{array}{l} \text{I love this movie! It's sweet,} \\ \text{but with satirical humor. The} \\ \text{dialogue is great and the} \\ \text{adventure scenes are fun...  It} \\ \text{manages to be whimsical and} \\ \text{romantic while laughing at the} \\ \text{conventions of the fairy tale} \\ \text{genre. I would recommend it to} \\ \text{just about anyone. I've seen} \\ \text{it several times, and I'm} \\ \text{always happy to see it again} \\ \text{whenever I have a friend who} \\ \text{hasn't seen it yet.} \end{array} \right) = c$$

# The bag of words representation

$$\gamma \left( \begin{array}{|l|l|} \hline \text{great} & 2 \\ \hline \text{love} & 2 \\ \hline \text{recommend} & 1 \\ \hline \text{laugh} & 1 \\ \hline \text{happy} & 1 \\ \hline \dots & \dots \\ \hline \end{array} \right) = c$$

# Bag of words for document classification

?

Test document

parser
language
label
translation
…

| Machine Learning | NLP | Garbage Collection | Planning | GUI |
|---|---|---|---|---|
| learning | parser | garbage | planning | … |
| training | tag | collection | temporal | |
| algorithm | training | memory | reasoning | |
| shrinkage | translation | optimization | plan | |
| network… | language… | region… | language… | |

# Multinomial Naïve Bayes Independence Assumptions
$$P(x_1, x_2, \ldots, x_n \mid c)$$

- **Bag of Words assumption**: Assume position doesn't matter

- **Conditional Independence**: Assume the feature probabilities $P(x_i \mid c_j)$ are independent given the class $c$.

$$P(x_1, \ldots, x_n \mid c) = P(x_1 \mid c) \bullet P(x_2 \mid c) \bullet P(x_3 \mid c) \bullet \ldots \bullet P(x_n \mid c)$$

# Applying Multinomial Naive Bayes Classifiers to Text Classification

positions ← all word positions in test document

$$c_{NB} = \underset{c_j \in C}{\operatorname{argmax}} P(c_j) \prod_{i \in positions} P(x_i \mid c_j)$$

# Text Classification and Naïve Bayes

Formalizing the Naïve Bayes Classifier

# Text Classification and Naïve Bayes

## Naïve Bayes: Learning

# Learning the Multinomial Naïve Bayes Model

- First attempt: maximum likelihood estimates
  - simply use the frequencies in the data

$$\hat{P}(c_j) = \frac{doccount(C = c_j)}{N_{doc}}$$

$$\hat{P}(w_i \mid c_j) = \frac{count(w_i, c_j)}{\sum_{w \in V} count(w, c_j)}$$

# Parameter estimation

$$\hat{P}(w_i \mid c_j) = \frac{count(w_i, c_j)}{\sum_{w \in V} count(w, c_j)}$$

fraction of times word $w_i$ appears among all words in documents of topic $c_j$

- Create mega-document for topic *j* by concatenating all docs in this topic
  - Use frequency of *w* in mega-document

# Problem with Maximum Likelihood

- What if we have seen no training documents with the word *fantastic* and classified in the topic **positive** (*thumbs-up)*?

- Zero probabilities cannot be conditioned away, no matter the other evidence!

$$\hat{P}(\text{"fantastic"} \mid \text{positive}) \; = \; \frac{count(\text{"fantastic"}, \text{positive})}{\sum\limits_{w \in V} count(w, \text{positive})} \; = \; 0$$

$$c_{MAP} = \operatorname{argmax}_{c} \hat{P}(c) \prod_{i} \hat{P}(x_i \mid c)$$

# Laplace (add-1) smoothing: unknown words

Add one extra word to the vocabulary, the "unknown word" $w_u$

$$\hat{P}(w_u \mid c) = \frac{count(w_u, c) + 1}{\left( \sum_{w \in V} count(w, c) \right) + |V + 1|}$$

$$= \frac{1}{\left( \sum_{w \in V} count(w, c) \right) + |V + 1|}$$

# Underflow Prevention: log space

- Multiplying lots of probabilities can result in floating-point underflow.
- Since log($xy$) = log($x$) + log($y$)
  - Better to sum logs of probabilities instead of multiplying probabilities.
- Class with highest un-normalized log probability score is still most probable.

$$c_{NB} = \underset{c_j \in C}{\mathrm{argmax}} \log P(c_j) + \sum_{i \in positions} \log P(x_i \mid c_j)$$

- Model is now just max of sum of weights

# Text Classification and Naïve Bayes

## Naïve Bayes: Learning

# Text Classification and Naïve Bayes

## Multinomial Naïve Bayes: A Worked Example

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

$$\hat{P}(w \mid c) = \frac{count(w,c)+1}{count(c)+|V|+1}$$

| | Doc | Words | Class |
|---|---|---|---|
| Training | 1 | Chinese Beijing Chinese | c |
| | 2 | Chinese Chinese Shanghai | c |
| | 3 | Chinese Macao | c |
| | 4 | Tokyo Japan Chinese | j |
| Test | 5 | Chinese Chinese Chinese Tokyo Japan | ? |

**Priors:**
$P(c)=$  $\frac{3}{4}$
$P(j)=$  $\frac{1}{4}$

**Choosing a class:**
P(c|d5) $\propto$ 3/4 * (6/15)$^3$ * 1/15 * 1/15
                    $\approx$ 0.0002

**Conditional Probabilities:**
P(Chinese|c) =   (5+1) / (8+7) = 6/15
P(Tokyo|c)   =   (0+1) / (8+7) = 1/15
P(Japan|c)   =   (0+1) / (8+7) = 1/15
P(Chinese|j) =   (1+1) / (3+7) = 2/10
P(Tokyo|j)   =   (1+1) / (3+7) = 2/10
P(Japan|j)   =   (1+1) / (3+7) = 2/10

P(j|d5) $\propto$  1/4 * (2/10)$^3$ * 2/10 * 2/10
                    $\approx$ 0.00008

43

# Summary: Naive Bayes is Not So Naive

- Robust to Irrelevant Features

    Irrelevant Features cancel each other without affecting results

- Very good in domains with many equally important features

    Decision Trees suffer from *fragmentation* in such cases – especially if little data

- Optimal if the independence assumptions hold: If assumed independence is correct, then it is the Bayes Optimal Classifier for problem

- A good dependable baseline for text classification
    - **But we will see other classifiers that give better accuracy**

# Text Classification and Naïve Bayes

## Multinomial Naïve Bayes: A Worked Example