# Outline of the course

- **Intro on AI, ML and NLP**

- **Text Processing**

- **Words and Corpora**

- **Lexical similarity**

- **Language Modeling**

- **Text Classification**

- **Semantic similarity**

- **Knowledge Graphs**

- **Intro to Large Language Models**

# The task

## Text Classification

# 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.

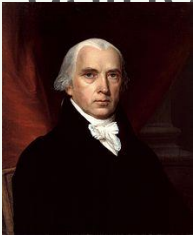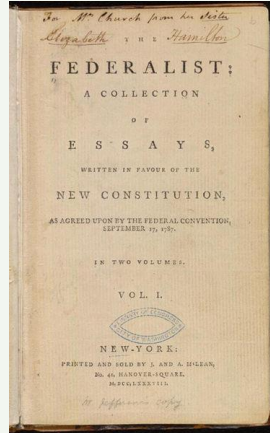http://www.123contactform.com/contact-form-StanfordNew1-236335.html

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.
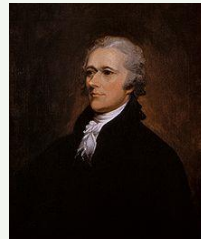
© Stanford University. All Rights Reserved.

# 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
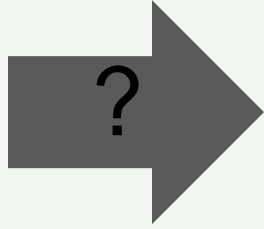
James Madison

Alexander Hamilton

# What is the subject of this medical article?

MEDLINE Article



**MeSH Subject Category Hierar**

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

# Positive or negative movie review?

**+** *...zany characters and richly applied satire, and some great plot twists*

**−** *It was pathetic. The worst part about it was the boxing scenes...*

**+** *...awesome caramel sauce and sweet toasty almonds. I love this place!*

**−** *...awful pizza and ridiculously overpriced...*

# Positive or negative movie review?

**+** *...zany characters and **richly** applied satire, and some **great** plot twists*

**−** *It was **pathetic**. The **worst** part about it was the boxing scenes...*

**+** *...**awesome** caramel sauce and sweet toasty almonds. I **love** this place!*

**−** *...**awful** pizza and **ridiculously** overpriced...*

# Why sentiment analysis?

- *Movie*:  is this review positive or negative?
- *Products*: what do people think about the new iPhone?
- *Public sentiment*: how is consumer confidence?
- *Politics*: what do people think about this candidate or issue?
- *Prediction*: predict election outcomes or market trends from sentiment

# Scherer Typology of Affective States

- **Emotion**: brief organically synchronized … evaluation of a major event
  - *angry, sad, joyful, fearful, ashamed, proud, elated*
- **Mood**: diffuse non-caused low-intensity long-duration change in subjective feeling
  - *cheerful, gloomy, irritable, listless, depressed, buoyant*
- **Interpersonal stances**: affective stance toward another person in a specific interaction
  - *friendly, flirtatious, distant, cold, warm, supportive, contemptuous*
- **Attitudes**: enduring, affectively colored beliefs, dispositions towards objects or persons
  - *liking, loving, hating, valuing, desiring*
- **Personality traits**: stable personality dispositions and typical behavior tendencies
  - *nervous, anxious, reckless, morose, hostile, jealous*

# Scherer Typology of Affective States

- **Emotion**: brief organically synchronized … evaluation of a major event
  - *angry, sad, joyful, fearful, ashamed, proud, elated*
- **Mood**: diffuse non-caused low-intensity long-duration change in subjective feeling
  - *cheerful, gloomy, irritable, listless, depressed, buoyant*
- **Interpersonal stances**: affective stance toward another person in a specific interaction
  - *friendly, flirtatious, distant, cold, warm, supportive, contemptuous*
- **Attitudes: enduring, affectively colored beliefs, dispositions towards objects or persons**
  - *liking, loving, hating, valuing, desiring*
- **Personality traits**: stable personality dispositions and typical behavior tendencies
  - *nervous, anxious, reckless, morose, hostile, jealous*

# Basic Sentiment Classification

- Sentiment analysis is the detection of **attitudes**
- Simple task we focus on in this chapter
  - Is the attitude of this text positive or negative?
- We return to affect classification in later chapters

# Summary: Text Classification

- Sentiment analysis

- Spam detection

- Authorship identification

- Language Identification

- Assigning subject categories, topics, or genres

- ...

# Text Classification: definition

- *Input*:

  ○ a document $d$

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

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

# Classaification Methods:  Hand-coded rules

- Rules based on combinations of words or other features
  - spam: black-list-address OR ("dollars" AND "you 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
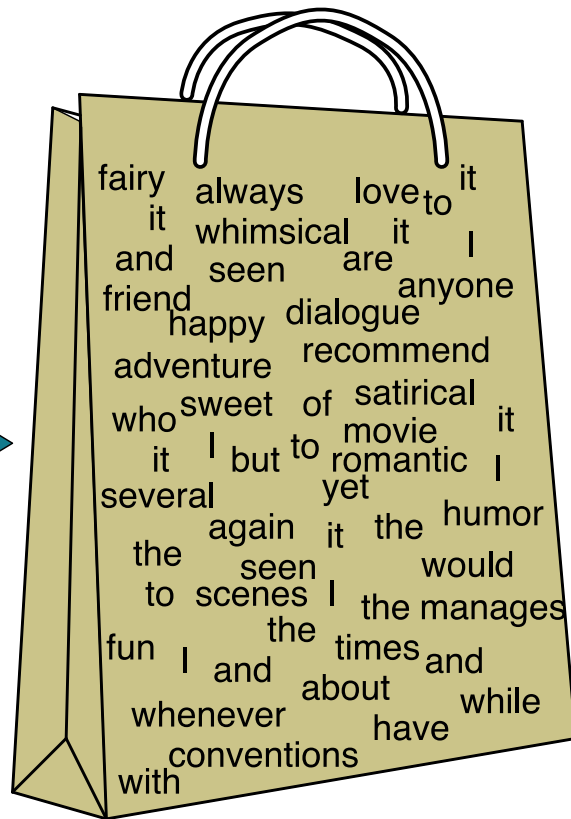  - Neural networks
  - k-Nearest Neighbors
  - …

# The task

## Text Classification

# Naive Bayes Intuition

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

# The Bag of Words Representation

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



| it | 6 |
|---|---|
| I | 5 |
| the | 4 |
| to | 3 |
| and | 3 |
| seen | 2 |
| yet | 1 |
| would | 1 |
| whimsical | 1 |
| times | 1 |
| sweet | 1 |
| satirical | 1 |
| adventure | 1 |
| genre | 1 |
| fairy | 1 |
| humor | 1 |
| have | 1 |
| great | 1 |
| … | … |

# The bag of words representation

$$\gamma \left( \begin{array}{|l|l|} \hline \text{seen} & 2 \\ \hline \text{sweet} & 1 \\ \hline \text{whimsical} & 1 \\ \hline \text{recommend} & 1 \\ \hline \text{happy} & 1 \\ \hline \ldots & \ldots \\ \hline \end{array} \right) = c$$

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

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

# Naive Bayes Classifier (I)

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

MAP is "maximum a posteriori" = most likely class

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

Bayes Rule

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

Dropping the denominator

# Naive Bayes Classifier (II)

"Likelihood"    "Prior"

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

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

Document d represented as features x1..xn

# Naïve Bayes Classifier (IV)

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

$O(|X|^n \bullet |C|)$ 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

# Multinomial Naive Bayes Independence Assumptions

$$P(x_1, x_2, \square \;, x_n \mid c)$$

- **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 $c$.

$$P(x_1, \square \;, x_n \mid c) = P(x_1 \mid c) \cdot P(x_2 \mid c) \cdot P(x_3 \mid c) \cdot \ldots \cdot P(x_n \mid c)$$

# Multinomial Naive Bayes Classifier

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

$$c_{NB} = \underset{c \in C}{\mathrm{argmax}}\, P(c_j) \prod_{x \in X} P(x \mid c)$$

# Applying Multinomial Naive Bayes Classifiers to Text Classification

positions $\leftarrow$ all word positions in test document

$$c_{NB} = \underset{c_j \hat{\in} \, C}{\mathrm{argmax}} \, P(c_j) \underset{i \hat{\in} \, positions}{\widetilde{O}} P(x_i \mid c_j)$$

# Problems with multiplying lots of probs

- There's a problem with this:

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

- Multiplying lots of probabilities can result in floating-point underflow!
- .0006 * .0007 * .0009 * .01 * .5 * .000008….
- Idea:  Use logs, because  $\log(ab) = \log(a) + \log(b)$
- We'll sum logs of probabilities instead of multiplying probabilities!

# We actually do everything in log space

Instead of this:
$$c_{NB} = \operatorname*{argmax}_{c_j \in C} P(c_j) \widetilde{\bigcirc}_{i \in positions} P(x_i \mid c_j)$$

This:
$$c_{\mathrm{NB}} = \operatorname*{argmax}_{c_j \in C} \left[ \log P(c_j) + \sum_{i \in \mathrm{positions}} \log P(x_i | c_j) \right]$$

Notes:

1) Taking log doesn't change the ranking of classes!

The class with highest probability also has highest log probability!

2) It's a linear model:

Just a max of a sum of weights: a **linear** function of the inputs

So naive bayes is a **linear classifier**

# Naive Bayes: Learning

**Text Classification**

# Learning the Multinomial Naive Bayes Model

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

$$\hat{P}(c_j) = \frac{N_{c_j}}{N_{total}}$$

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

# Parameter estimation

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

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)*?

$$\hat{P}(\text{"fantastic"} \mid \text{positive}) = \frac{count(\text{"fantastic"}, \text{positive})}{\overset{\circ}{a} \, count(w, \text{positive})} = 0$$

$$w \hat{I} \, V$$

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

$$c_{MAP} = \text{argmax}_c \, \hat{P}(c) \tilde{O}_i \, \hat{P}(x_i \mid c)$$

# Laplace (add-1) smoothing for Naïve Bayes

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

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

# Multinomial Naïve Bayes: Learning

- From training corpus, extract *Vocabulary*

- Calculate $P(c_j)$ terms
  - For each $c_j$ in $C$ do
    $docs_j \leftarrow$ all docs with class $=c_j$

$$P(c_j) \neg \frac{|docs_j|}{|\text{total \# documents}|}$$

- Calculate $P(w_k \mid c_j)$ terms
  - $Text_j \leftarrow$ single doc containing all $docs_j$
  - For each word $w_k$ in *Vocabulary*
    $n_k \leftarrow$ \# of occurrences of $w_k$ in $Text_j$

$$P(w_k \mid c_j) \neg \frac{n_k + a}{n + a|Vocabulary|}$$

# Unknown words

- What about unknown words
  - that appear in our test data
  - but not in our training data or vocabulary?
- We **ignore** them
  - Remove them from the test document!
  - Pretend they weren't there!
  - Don't include any probability for them at all!
- Why don't we build an unknown word model?
  - It doesn't help: knowing which class has more unknown words is not generally helpful!

# Stop words

- Some systems ignore stop words
  - **Stop words:** very frequent words like *the* and *a*.
    - Sort the vocabulary by word frequency in training set
    - Call the top 10 or 50 words the **stopword list**.
    - Remove all stop words from both training and test sets
      - As if they were never there!
- But removing stop words doesn't usually help
  - So in practice most NB algorithms use **all** words and **don't** use stopword lists

# Naive Bayes: Learning

**Text Classification**

# Naive Bayes: Sentiment and Binary NB

**Text Classification**

# Let's do a worked sentiment example!

| Cat | Documents |
|---|---|
| Training - | just plain boring |
| - | entirely predictable and lacks energy |
| - | no surprises and very few laughs |
| + | very powerful |
| + | the most fun film of the summer |
| Test ? | predictable with no fun |

# A worked sentiment example with add-1 smoothing

| Cat | Documents |
|---|---|
| Training - | just plain boring |
| - | entirely predictable and lacks energy |
| - | no surprises and very few laughs |
| + | very powerful |
| + | the most fun film of the summer |
| Test ? | predictable ~~with~~ no fun |

1. Prior from training:

$$\hat{P}(c_j) = \frac{N_{c_j}}{N_{total}}$$

P(-) = 3/5
P(+) = 2/5

2. Drop "with"

3. Likelihoods from training:

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

$$P(\text{"predictable"}|-) = \frac{1+1}{14+20} \quad P(\text{"predictable"}|+) = \frac{0+1}{9+20}$$

$$P(\text{"no"}|-) = \frac{1+1}{14+20} \quad P(\text{"no"}|+) = \frac{0+1}{9+20}$$

$$P(\text{"fun"}|-) = \frac{0+1}{14+20} \quad P(\text{"fun"}|+) = \frac{1+1}{9+20}$$

4. Scoring the test set:

$$P(-)P(S|-) = \frac{3}{5} \times \frac{2 \times 2 \times 1}{34^3} = 6.1 \times 10^{-5}$$

$$P(+)P(S|+) = \frac{2}{5} \times \frac{1 \times 1 \times 2}{29^3} = 3.2 \times 10^{-5}$$

# Optimizing for sentiment analysis

For tasks like sentiment, word **occurrence** seems to be more important than word **frequency**.

- The occurrence of the word *fantastic* tells us a lot
- The fact that it occurs 5 times may not tell us much more.

**Binary multinominal naive bayes**, or **binary NB**

- Clip our word counts at 1
- Note: this is different than Bernoulli naive bayes; see the textbook at the end of the chapter.

# Binary Multinomial Naïve Bayes: Learning

- From training corpus, extract *Vocabulary*

- Calculate $P(c_j)$ terms
  - For each $c_j$ in $C$ do
    $docs_j \leftarrow$ all docs with class $=c_j$

$$P(c_j) \neg \frac{|docs_j|}{|\text{total \# documents}|}$$

- Calculate $P(w_k \mid c_j)$ terms
  - Remove duplicates in each doc:
  - Text$_j \leftarrow$ single doc containing all $docs_j$
  - For each word $w_k$ in *Vocabulary*
    - For each word type $w$ in doc
    - Retain only a single instance of $w$
    $n_k \leftarrow$ # of occurrences of $w_k$ in $Text_j$

$$P(w_k \mid c_j) \neg \frac{n_k + a}{n + a \mid Vocabulary \mid}$$

## Binary Maultinomial Naive Bayes
 on a test document *d*

- First remove all duplicate words from *d*
- Then compute NB using the same equation:

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

# Binary multinominal naive Bayes

**Four original documents:**

- − it was pathetic the worst part was the boxing scenes
- − no plot twists or great scenes
- + and satire and great plot twists
- + great scenes great film

# Binary multinominal naive Bayes

**Four original documents:**

- − it was pathetic the worst part was the boxing scenes
- − no plot twists or great scenes
- + and satire and great plot twists
- + great scenes great film

| | NB Counts | |
|---|---|---|
| | + | − |
| and | 2 | 0 |
| boxing | 0 | 1 |
| film | 1 | 0 |
| great | 3 | 1 |
| it | 0 | 1 |
| no | 0 | 1 |
| or | 0 | 1 |
| part | 0 | 1 |
| pathetic | 0 | 1 |
| plot | 1 | 1 |
| satire | 1 | 0 |
| scenes | 1 | 2 |
| the | 0 | 2 |
| twists | 1 | 1 |
| was | 0 | 2 |
| worst | 0 | 1 |

# Binary multinominal naive Bayes

**Four original documents:**

- − it was pathetic the worst part was the boxing scenes
- − no plot twists or great scenes
- + and satire and great plot twists
- + great scenes great film

**After per-document binarization:**

- − it was pathetic the worst part boxing scenes
- − no plot twists or great scenes
- + and satire great plot twists
- + great scenes film

| | NB Counts | |
|---|---|---|
| | + | − |
| and | 2 | 0 |
| boxing | 0 | 1 |
| film | 1 | 0 |
| great | 3 | 1 |
| it | 0 | 1 |
| no | 0 | 1 |
| or | 0 | 1 |
| part | 0 | 1 |
| pathetic | 0 | 1 |
| plot | 1 | 1 |
| satire | 1 | 0 |
| scenes | 1 | 2 |
| the | 0 | 2 |
| twists | 1 | 1 |
| was | 0 | 2 |
| worst | 0 | 1 |

# Binary multinominal naive Bayes

**Four original documents:**

- − it was pathetic the worst part was the boxing scenes
- − no plot twists or great scenes
- + and satire and great plot twists
- + great scenes great film

**After per-document binarization:**

- − it was pathetic the worst part boxing scenes
- − no plot twists or great scenes
- + and satire great plot twists
- + great scenes film

Counts can still be 2! Binarization is within-doc!

| | NB Counts | | Binary Counts | |
|---|---|---|---|---|
| | + | − | + | − |
| and | 2 | 0 | 1 | 0 |
| boxing | 0 | 1 | 0 | 1 |
| film | 1 | 0 | 1 | 0 |
| great | 3 | 1 | 2 | 1 |
| it | 0 | 1 | 0 | 1 |
| no | 0 | 1 | 0 | 1 |
| or | 0 | 1 | 0 | 1 |
| part | 0 | 1 | 0 | 1 |
| pathetic | 0 | 1 | 0 | 1 |
| plot | 1 | 1 | 1 | 1 |
| satire | 1 | 0 | 1 | 0 |
| scenes | 1 | 2 | 1 | 2 |
| the | 0 | 2 | 0 | 1 |
| twists | 1 | 1 | 1 | 1 |
| was | 0 | 2 | 0 | 1 |
| worst | 0 | 1 | 0 | 1 |

# Naive Bayes: Sentiment and Binary NB

**Text Classification**

# More on Sentiment Classification

**Text Classification**

# Sentiment Classification: Dealing with Negation

- I really like this movie

I really **don't** like this movie

Negation changes the meaning of "like" to negative.

Negation can also change negative to positive-ish
- **Don't** dismiss this film
- **Doesn't** let us get bored

# Sentiment Classification: Dealing with Negation

Das, Sanjiv and Mike Chen. 2001. Yahoo! for Amazon: Extracting market sentiment from stock message boards. In Proceedings of the Asia Pacific Finance Association Annual Conference (APFA).
Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up? Sentiment Classification using Machine Learning Techniques. EMNLP-2002, 79—86.

Simple baseline method:
Add NOT_ to every word between negation and following punctuation:

```
didn't like this movie , but I
```



```
didn't NOT_like NOT_this NOT_movie but I
```

# Sentiment Classification: Lexicons

- Sometimes we don't have enough labeled training data
- In that case, we can make use of pre-built word lists
- Called **lexicons**
- There are various publically available lexicons

# MPQA Subjectivity Cues Lexicon

Theresa Wilson, Janyce Wiebe, and Paul Hoffmann (2005). Recognizing Contextual Polarity in Phrase-Level Sentiment Analysis. Proc. of HLT-EMNLP-2005.

Riloff and Wiebe (2003). Learning extraction patterns for subjective expressions. EMNLP-2003.

- Home page: https://mpqa.cs.pitt.edu/lexicons/subj_lexicon/
- 6885 words from 8221 lemmas, annotated for intensity (strong/weak)
  - 2718 positive
  - 4912 negative
- + : *admirable, beautiful, confident, dazzling, ecstatic, favor, glee, great*
- − : *awful, bad, bias, catastrophe, cheat, deny, envious, foul, harsh, hate*

# The General Inquirer

Philip J. Stone, Dexter C Dunphy, Marshall S. Smith, Daniel M. Ogilvie. 1966. The General Inquirer: A Computer Approach to Content Analysis. MIT Press

○ Home page: http://www.wjh.harvard.edu/~inquirer
○ List of Categories:  http://www.wjh.harvard.edu/~inquirer/homecat.htm
○ Spreadsheet: http://www.wjh.harvard.edu/~inquirer/inquirerbasic.xls
● Categories:
  ○ Positiv (1915 words) and Negativ (2291 words)
  ○ Strong vs Weak, Active vs Passive, Overstated versus Understated
  ○ Pleasure, Pain, Virtue, Vice, Motivation, Cognitive Orientation, etc
● Free for Research Use

# Using Lexicons in Sentiment Classification

**Add a feature** that gets a count whenever a word from the lexicon occurs
- E.g., a feature called "**this word occurs in the positive lexicon**" or "**this word occurs in the negative lexicon**"

Now all positive words (*good, great, beautiful, wonderful*) or negative words count for that feature.

Using 1-2 features isn't as good as using all the words.
- But when training data is sparse or not representative of the test set, dense lexicon features can help

# Naive Bayes in Other tasks: Spam Filtering

- SpamAssassin Features:
  - Mentions millions of (dollar) ((dollar) NN,NNN,NNN.NN)
  - From: starts with many numbers
  - Subject is all capitals
  - HTML has a low ratio of text to image area
  - "One hundred percent guaranteed"
  - Claims you can be removed from the list

# Naive Bayes in Language ID

- Determining what language a piece of text is written in.

Features based on character n-grams do very well

- Important to train on lots of varieties of each language

    (e.g., American English varieties like African-American English, or English varieties around the world like Indian English)

# Summary: Naive Bayes is Not So Naive

- Very Fast, low storage requirements
- Work well with very small amounts of training data
- 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

# More on Sentiment Classification

**Text Classification**

# Multiclass (skip)

**Text Classification**

# Confusion Matrix for 3-class classification



|  | *gold labels* | | |
|---|---|---|---|
|  | urgent | normal | spam |
| **urgent** | 8 | 10 | 1 | $\mathbf{precision_u}=\dfrac{8}{8+10+1}$ |
| **normal** | 5 | 60 | 50 | $\mathbf{precision_n}=\dfrac{60}{5+60+50}$ |
| **spam** | 3 | 30 | 200 | $\mathbf{precision_s}=\dfrac{200}{3+30+200}$ |

*system output*

$$\mathbf{recall_u}=\frac{8}{8+5+3} \qquad \mathbf{recall_n}=\frac{60}{10+60+30} \qquad \mathbf{recall_s}=\frac{200}{1+50+200}$$

# How to combine P/R from 3 classes to get one metric

- Macroaveraging:
  - compute the performance for each class, and then average over classes
- Microaveraging:
  - collect decisions for all classes into one confusion matrix
  - compute precision and recall from that table.

# Macroaveraging and Microaveraging

| Class 1: Urgent | | |
|---|---|---|
| | true urgent | true not |
| system urgent | 8 | 11 |
| system not | 8 | 340 |

$$\text{precision} = \frac{8}{8+11} = .42$$

| Class 2: Normal | | |
|---|---|---|
| | true normal | true not |
| system normal | 60 | 55 |
| system not | 40 | 212 |

$$\text{precision} = \frac{60}{60+55} = .52$$

| Class 3: Spam | | |
|---|---|---|
| | true spam | true not |
| system spam | 200 | 33 |
| system not | 51 | 83 |

$$\text{precision} = \frac{200}{200+33} = .86$$

| Pooled | | |
|---|---|---|
| | true yes | true no |
| system yes | 268 | 99 |
| system no | 99 | 635 |

$$\text{microaverage precision} = \frac{268}{268+99} = \mathbf{.73}$$

$$\text{macroaverage precision} = \frac{.42+.52+.86}{3} = \mathbf{.60}$$

# Multiclass (skip)

**Text Classification**

# Harms in sentiment classifiers

- Kiritchenko and Mohammad (2018) found that most sentiment classifiers assign lower sentiment and more negative emotion to sentences with African American names in them.
- This perpetuates negative stereotypes that associate African Americans with negative emotions

# Harms in toxicity classification

- Toxicity detection is the task of detecting hate speech, abuse, harassment, or other kinds of toxic language
- But some toxicity classifiers incorrectly flag as being toxic sentences that are non-toxic but simply mention identities like blind people, women, or gay people.
- This could lead to censorship of discussion about these groups.

# What causes these harms?

- Can be caused by:
  - Problems in the training data; machine learning systems are known to amplify the biases in their training data.
  - Problems in the human labels
  - Problems in the resources used (like lexicons)
  - Problems in model architecture (like what the model is trained to optimized)
- Mitigation of these harms is an open research area
- Meanwhile: **model cards**

# Model Cards

(Mitchell et al., 2019)

- For each algorithm you release, document:
  - training algorithms and parameters
  - training data sources, motivation, and preprocessing
  - evaluation data sources, motivation, and preprocessing
  - intended use and users
  - model performance across different demographic or other groups and environmental situations