# Text Classification and Naïve Bayes

## Formalization & Learning

Slides borrowed from Stanford with modifications

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

# 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

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

# 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

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

# 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

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

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)*?

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

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

$$c_{MAP} = \text{argmax}_c \ \hat{P}(c) \prod_i \hat{P}(x_i \mid c)$$

# Laplace (add-1) smoothing: unknown words

Add one extra word to the vocabulary, the "unknown word" w$_u$

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

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

# 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.

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

- 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

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

# 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

# Text Classification and Naïve Bayes

## Multinomial Naïve Bayes: A Worked Example