

# Probability and Naïve Bayes

# What is Probability?

- “The probability the coin will land heads is 0.5”
  - Q: what does this mean?
- 2 Interpretations:
  - Frequentist (Repeated trials)
    - If we flip the coin many times...
  - Bayesian
    - We believe there is equal chance of heads/tails
    - Advantage: events that do not have long term frequencies

Q: What is the probability the polar ice caps will melt by 2050?

# Probability Review

$$\sum_x P(X = x) = 1$$

Conditional  
Probability

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

Chain Rule

$$P(A|B)P(B) = P(A, B)$$

# Probability Review

Disjunction / Union:  $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$

Negation:  $P(\neg A) = 1 - P(A)$

$$\sum_x P(X = x, Y) = ??$$

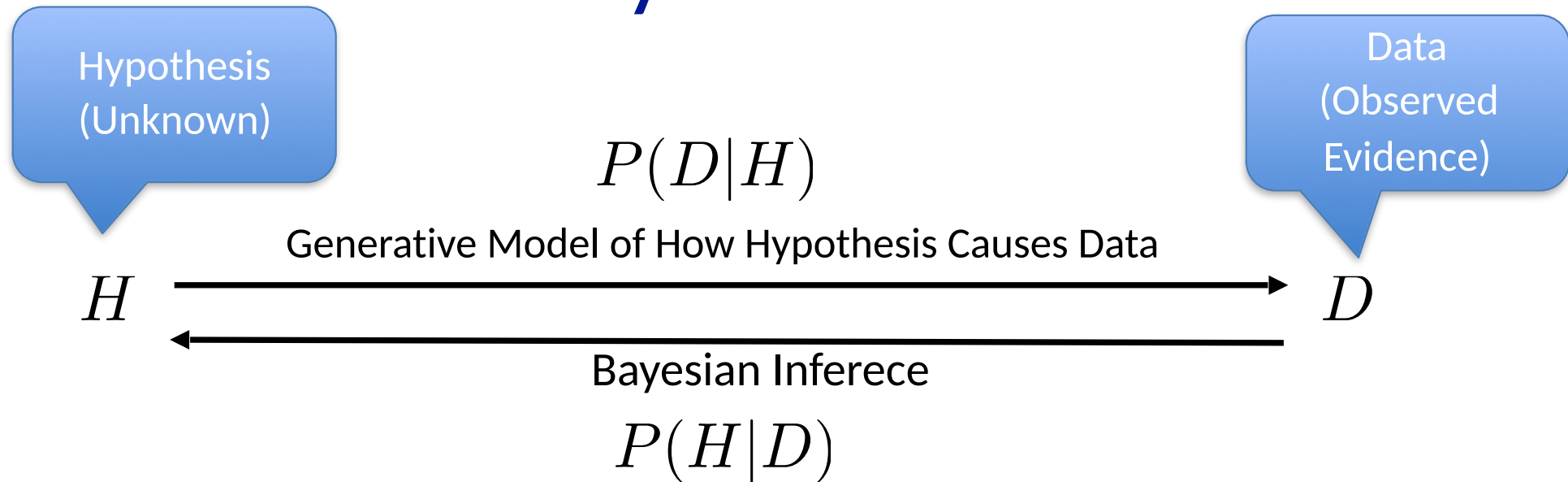
# Probability Review

Disjunction / Union:  $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$

Negation:  $P(\neg A) = 1 - P(A)$

$$\sum_x P(X = x, Y) = P(Y)$$

# Bayes Rule



Bayes Rule tells us how to flip the conditional  
Reason about effects to causes  
Useful if you assume a generative model for your data

$$P(H|D) = \frac{P(D|H)P(H)}{P(D)}$$

# Bayes Rule

Bayes Rule tells us how to flip the conditional  
Reason about effects to causes

Useful if you assume a generative model for your data

The diagram illustrates the components of Bayes' Rule. The formula is  $P(H|D) = \frac{P(D|H)P(H)}{P(D)}$ . Arrows point from labels to parts of the formula: 'Likelihood' points to  $P(D|H)$ , 'Prior' points to  $P(H)$ , 'Posterior' points to  $P(H|D)$ , and 'Normalizer' points to  $P(D)$ .

$$P(H|D) = \frac{P(D|H)P(H)}{P(D)}$$

Labels and arrows:

- Likelihood points to  $P(D|H)$
- Prior points to  $P(H)$
- Posterior points to  $P(H|D)$
- Normalizer points to  $P(D)$

# Bayes Rule

Bayes Rule tells us how to flip the conditional  
Reason about effects to causes

Useful if you assume a generative model for your data

The diagram illustrates the components of Bayes' Rule. The formula is  $P(H|D) = \frac{P(D|H)P(H)}{\sum_h P(D|H)P(H)}$ . Arrows point from labels to parts of the formula: 'Likelihood' points to  $P(D|H)$  in the numerator; 'Prior' points to  $P(H)$  in the numerator; 'Posterior' points to  $P(H|D)$  on the left; and 'Normalizer' points to the denominator  $\sum_h P(D|H)P(H)$ .

Likelihood

Prior

Posterior

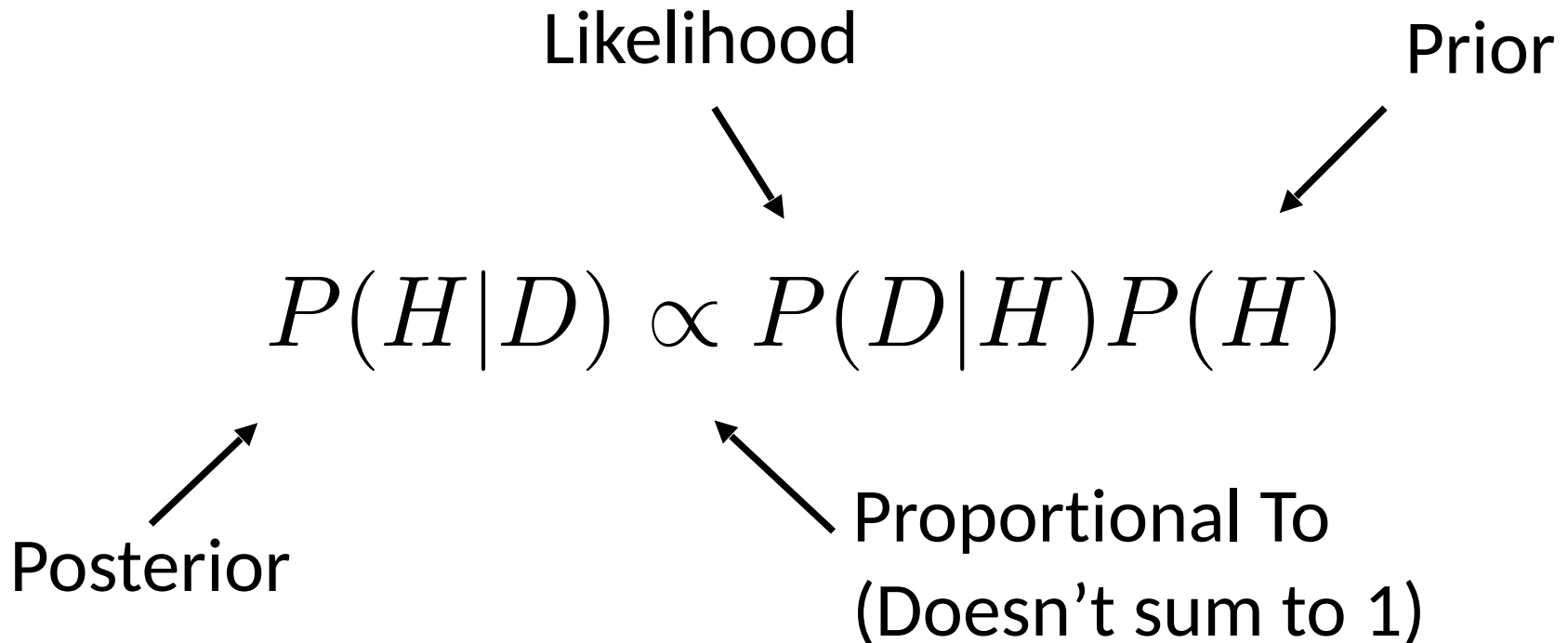
Normalizer

$$P(H|D) = \frac{P(D|H)P(H)}{\sum_h P(D|H)P(H)}$$



# Bayes Rule

Bayes Rule tells us how to flip the conditional  
Reason about effects to causes  
Useful if you assume a generative model for your data



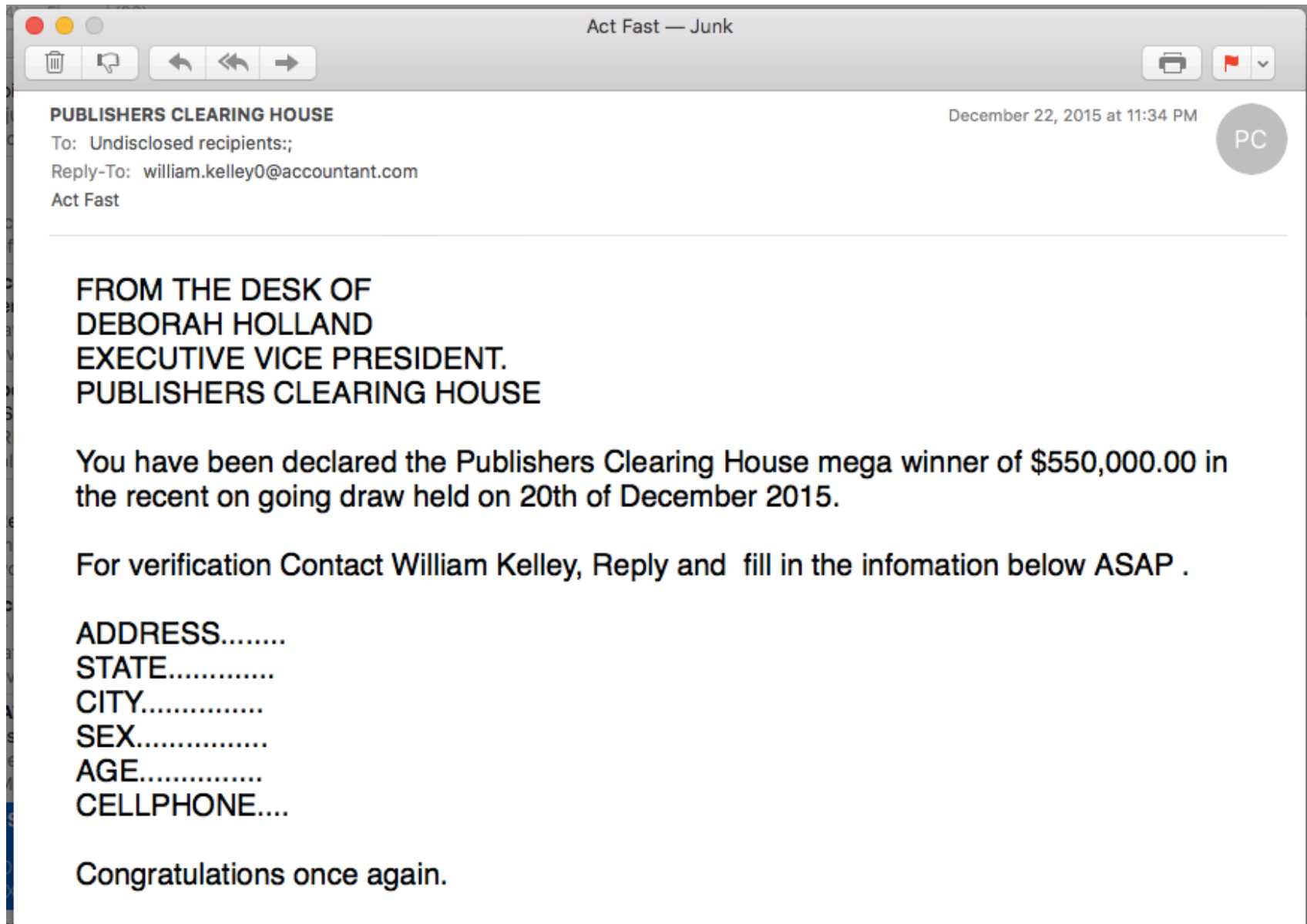
# Bayes Rule Example

- There is a disease that affects a tiny fraction of the population (0.01%)
- Symptoms include a headache and stiff neck
  - 99% of patients with the disease have these symptoms
- 1% of the general population has these symptoms

Q: assume you have the symptom, what is your probability of having the disease?

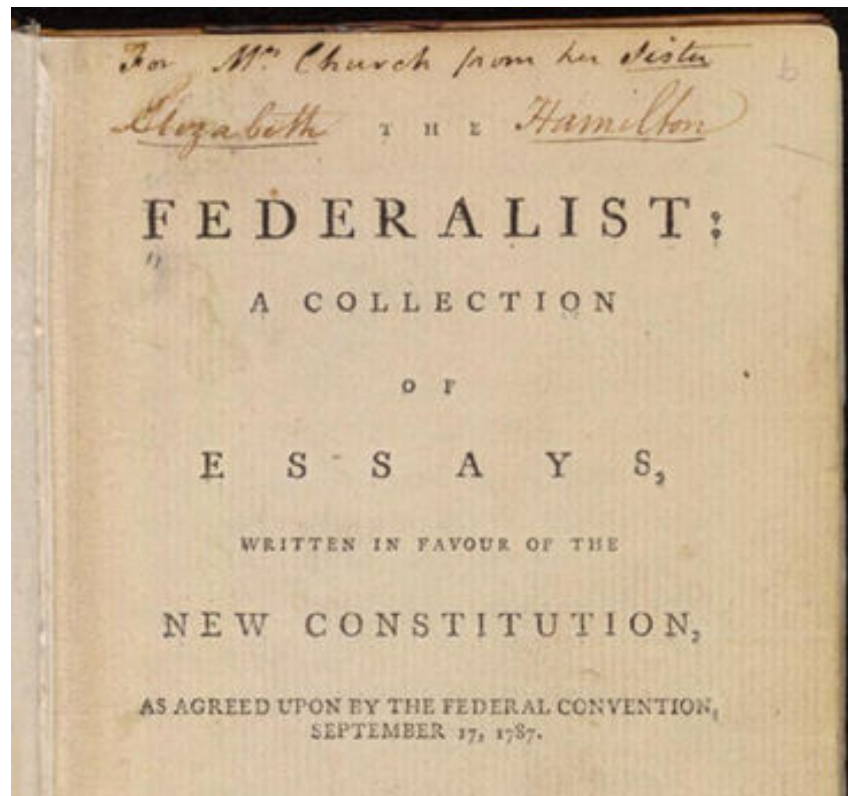
# Text Classification

# Is this Spam?



# Who wrote which Federalist papers?

- 1787-8: anonymous essays try to convince New York to ratify U.S. Constitution.
- Authorship of 12 of the letters in dispute
- 1963: solved by Mosteller and Wallace using Bayesian methods



# What is the subject of this article?





## MEDLINE Article



## MeSH Subject Category Hierarchy

- Antagonists and Inhibitors
- Blood Supply
- Chemistry
- Drug Therapy
- Embryology
- Epidemiology
- ...

# 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.

# Text Classification: definition

- *Input:*
  - a document  $d$
  - a fixed set of classes  $C = \{c_1, c_2, \dots, c_J\}$
- *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  $C = \{c_1, c_2, \dots, c_j\}$
  - A training set of  $m$  hand-labeled documents  $(d_1, c_1), \dots, (d_m, c_m)$
- Output:
  - a learned classifier  $\gamma: d \rightarrow c$

# Classification Methods: Supervised Machine Learning

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

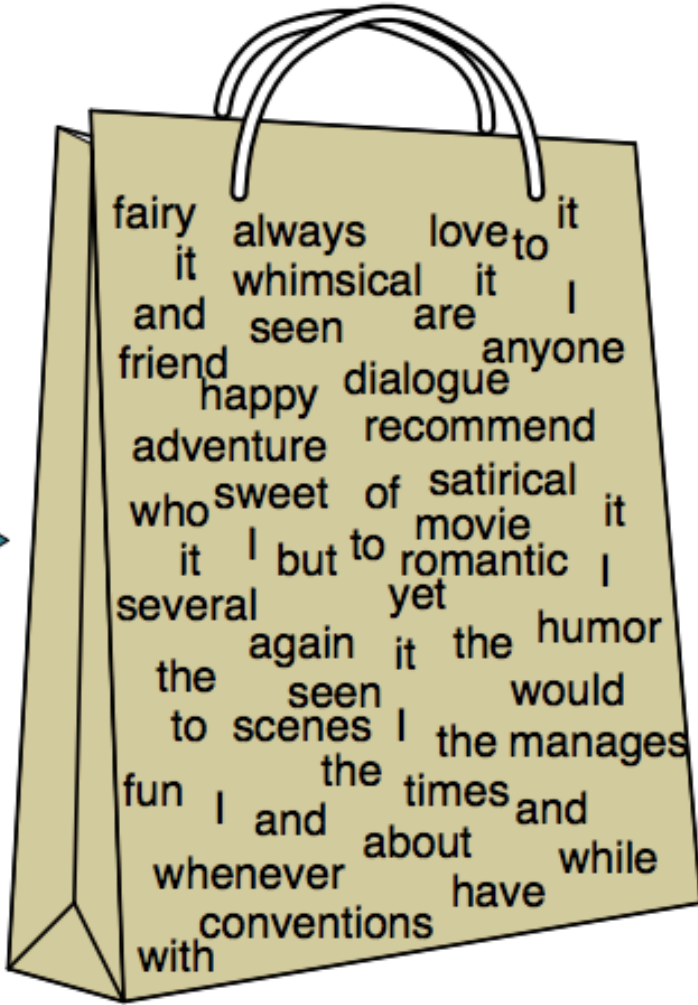
# Classification Methods: Supervised Machine Learning

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

# Naïve Bayes Intuition

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

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
...	...

# Bayes' Rule Applied to Documents and Classes

For a document  $d$  and a class  $c$

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

# Bayes' Rule Applied to Documents and Classes

For a document  $d$  and a class  $c$

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

$$C_{MAP} = \operatorname{argmax}_{c \in C} P(c|d)$$



# Naïve Bayes Classifier (I)

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

MAP is “maximum a posteriori” = most likely class

# Naïve Bayes Classifier (I)

$$\begin{aligned} C_{MAP} &= \operatorname{argmax}_{c \in C} P(c | d) \\ &= \operatorname{argmax}_{c \in C} \frac{P(d | c) P(c)}{P(d)} \end{aligned}$$

MAP is “maximum a posteriori” = most likely class

Bayes Rule

# Naïve Bayes Classifier (I)

$$\begin{aligned} C_{MAP} &= \operatorname{argmax}_{c \in C} P(c | d) \\ &= \operatorname{argmax}_{c \in C} \frac{P(d | c) P(c)}{P(d)} \\ &= \operatorname{argmax}_{c \in C} P(d | c) P(c) \end{aligned}$$

MAP is “maximum a posteriori” = most likely class

Bayes Rule

Dropping the denominator

## Naïve Bayes Classifier (II)

$$C_{MAP} = \operatorname{argmax}_{c \in C} P(d | c) P(c)$$

$$= \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c) P(c)$$

# Naïve Bayes Classifier (IV)

$$C_{MAP} = \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c) P(c)$$

How often does this class occur?

We can just count the relative frequencies in a corpus

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

# Multinomial Naïve Bayes Independence Assumptions

$$P(x_1, x_2, \dots, x_n \mid c)$$

- **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  $c$ .

# Multinomial Naïve Bayes Classifier

$$C_{MAP} = \operatorname{argmax}_{c \in C} P(x_1, x_2, \dots, x_n | c) P(c)$$

$$C_{NB} = \operatorname{argmax}_{c \in C} P(c_j) \prod_{x \in X} P(x | c)$$

# Applying Multinomial Naive Bayes Classifiers to Text Classification

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

positions  $\leftarrow$  all word positions in test document



# Learning the Multinomial Naïve Bayes Model

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

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

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

# Parameter estimation

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

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

# Parameter estimation

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

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

fraction of word in the full vocabulary that appered in topic  $c_j$

# Parameter estimation

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

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

fraction of word in the full vocabulary that appered in topic  $c_j$

# 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!

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

# Laplace (add-1) smoothing for Naïve Bayes

$$\begin{aligned}\hat{P}(w_i | c) &= \frac{\text{count}(w_i, c)}{\sum_{w \in V} (\text{count}(w, c))} \\ &= \frac{\text{count}(w_i, c) + 1}{\left( \sum_{w \in V} \text{count}(w, c) \right) + |V|}\end{aligned}$$

# Multinomial Naïve Bayes: Learning

- Calculate  $P(c_j)$  terms
  - \_ For each  $c_j$  in  $C$  do

$docs_j \leftarrow$  all docs with class  $= c_j$

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

# Multinomial Naïve Bayes: Learning

- From training corpus, extract *Vocabulary*
- Calculate  $P(w_k | c_j)$  terms
  - ▶  $Text_j \leftarrow$  single doc containing all docs<sub>j</sub>
  - ▶ For each word  $w_k$  in *Vocabulary*

$n_k \leftarrow$  # of occurrences of  $w_k$  in  $Text_j$

$$P(w_k | c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha | \text{Vocabulary} |}$$



# Exercise

# Multinomial Naïve Bayes: Learning

- Calculate  $P(c_j)$  terms
  - \_ For each  $c_j$  in  $C$  do

$docs_j \leftarrow$  all docs with class  $= c_j$

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

- From training corpus, extract *Vocabulary*
- Calculate  $P(w_k | c_j)$  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  $Text_j$

$$P(w_k | c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha |Vocabulary|}$$

Solution

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

$$P(-) = \frac{3}{5} \quad P(+) = \frac{2}{5}$$

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

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

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

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

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

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

$$P(S|-)P(-) = \frac{3}{5} \times \frac{2 \times 1 \times 2 \times 1}{34^4} = 1.8 \times 10^{-6}$$

$$P(S|+)P(+) = \frac{2}{5} \times \frac{1 \times 1 \times 1 \times 1}{29^4} = 5.7 \times 10^{-7}$$

The model thus predicts the class *negative* for the test sentence.