




Lecture 9: 3 May, 2021

Madhavan Mukund

<https://www.cmi.ac.in/~madhavan>

Data Mining and Machine Learning
April–July 2021

Bayesian classifiers

- As before
 - Attributes $\{A_1, A_2, \dots, A_k\}$ and
 - Classes $C = \{c_1, c_2, \dots, c_\ell\}$
- Each class c_i defines a probabilistic model for attributes
 - $Pr(A_1 = a_1, \dots, A_k = a_k \mid C = c_i)$ 
- Given a data item $d = (a_1, a_2, \dots, a_k)$, identify the best class c for d

- Maximize $Pr(C = c_i \mid A_1 = a_1, \dots, A_k = a_k)$


Bayesian classification

- Maximize $Pr(C = c_i \mid A_1 = a_1, \dots, A_k = a_k)$
- By Bayes' rule,

$$\begin{aligned} & Pr(C = c_i \mid A_1 = a_1, \dots, A_k = a_k) \\ &= \frac{Pr(A_1 = a_1, \dots, A_k = a_k \mid C = c_i) \cdot Pr(C = c_i)}{\underbrace{Pr(A_1 = a_1, \dots, A_k = a_k)}} \\ &= \frac{Pr(A_1 = a_1, \dots, A_k = a_k \mid C = c_i) \cdot Pr(C = c_i)}{\underbrace{\sum_{j=1}^{\ell} Pr(A_1 = a_1, \dots, A_k = a_k \mid C = c_j) \cdot Pr(C = c_j)}} \end{aligned}$$

- Denominator is the same for all c_i , so sufficient to maximize

$$Pr(A_1 = a_1, \dots, A_k = a_k \mid C = c_i) \cdot Pr(C = c_i)$$

Example

- To classify $A = g, B = q$

- $Pr(C = t) = 5/10 = 1/2$

- $Pr(A = g, B = q \mid C = t) = 2/5$

- $Pr(A = g, B = q \mid C = t) \cdot Pr(C = t) = 1/5$

- $Pr(C = f) = 5/10 = 1/2$

- $Pr(A = g, B = q \mid C = f) = \underline{1/5}$

- $Pr(A = g, B = q \mid C = f) \cdot Pr(C = f) = 1/10$

- Hence, predict $C = t$

$\rightarrow P(C=t \mid A=g, B=q)$

A	B	C
m	b	t
m	s	t
g	q	t
h	s	t
g	q	t
g	q	f
g	s	f
h	b	f
h	q	f
m	b	f

Example ...

- What if we want to classify $A = m, B = q$?
- $Pr(A = m, B = q \mid C = t) = 0$
- Also $Pr(A = m, B = q \mid C = f) = 0$!

A	B	C
m	b	t
m	s	t
g	q	t
h	s	t
g	q	t
g	q	f
g	s	f
h	b	f
h	q	f
m	b	f

Example ...

- What if we want to classify $A = m, B = q$?
- $Pr(A = m, B = q \mid C = t) = 0$
- Also $Pr(A = m, B = q \mid C = f) = 0!$
- To estimate joint probabilities across all combinations of attributes, we need a much larger set of training data

$A \times B$

A	B	C
m	b	t
m	s	t
g	q	t
h	s	t
g	q	t
g	q	f
g	s	f
h	b	f
h	q	f
m	b	f

Naïve Bayes classifier

- Strong simplifying assumption: attributes are pairwise independent

$$\Pr(\underline{A_1 = a_1, \dots, A_k = a_k} \mid C = c_i) = \prod_{j=1}^k \Pr(\underline{A_j = a_j} \mid \underline{C = c_i})$$

- $\Pr(C = c_i)$ is fraction of training data with class c_i
- $\Pr(\underline{A_j = a_j} \mid C = c_i)$ is fraction of training data labelled c_i for which $A_j = a_j$

Naïve Bayes classifier

- Strong simplifying assumption: attributes are pairwise independent

$$Pr(A_1 = a_1, \dots, A_k = a_k \mid C = c_i) = \prod_{j=1}^k Pr(A_j = a_j \mid C = c_i) \quad \text{a}$$

- $Pr(C = c_i)$ is fraction of training data with class c_i
- $Pr(A_j = a_j \mid C = c_i)$ is fraction of training data labelled c_i for which $A_j = a_j$
- Final classification is

$$\arg \max_{c_i} Pr(C = c_i) \prod_{j=1}^k Pr(A_j = a_j \mid C = c_i)$$

$$P(C=c_i) \cdot P(A_1=a_1 \rightarrow A_k=a_k \mid C=c_i)$$

Naïve Bayes classifier ...

- Conditional independence is not theoretically justified

Naïve Bayes classifier ...

- Conditional independence is not theoretically justified
- For instance, text classification
 - Items are documents, attributes are words (absent or present)
 - Classes are topics
 - Conditional independence says that a document is a set of words: ignores sequence of words
 - Meaning of words is clearly affected by relative position, ordering

Naïve Bayes classifier ...

- Conditional independence is not theoretically justified
- For instance, text classification
 - Items are documents, attributes are words (absent or present)
 - Classes are topics
 - Conditional independence says that a document is a set of words: ignores sequence of words
 - Meaning of words is clearly affected by relative position, ordering
- However, naive Bayes classifiers work well in practice, even for text classification!
 - Many spam filters are built using this model

Example revisited

- Want to classify $A = m, B = q$
- $Pr(A = m, B = q \mid C = t) = Pr(A = m, B = q \mid C = f) = 0$

A	B	C
m	b	t
m	s	t
g	q	t
h	s	t
g	q	t
g	q	f
g	s	f
h	b	f
h	q	f
m	b	f

Example revisited

- Want to classify $A = m, B = q$
- $Pr(A = m, B = q \mid C = t) = Pr(A = m, B = q \mid C = f) = 0$
- $Pr(A = m \mid C = t) = 2/5$
- $Pr(B = q \mid C = t) = 2/5$

A	B	C
m	b	t
m	s	t
g	q	t
h	s	t
g	q	t
g	q	f
g	s	f
h	b	f
h	q	f
m	b	f

Example revisited

- Want to classify $A = m, B = q$
- $Pr(A = m, B = q \mid C = t) = Pr(A = m, B = q \mid C = f) = 0$
- $Pr(A = m \mid C = t) = 2/5$
- $Pr(B = q \mid C = t) = 2/5$
- $Pr(A = m \mid C = f) = 1/5$ —
- $Pr(B = q \mid C = f) = 2/5$ —

A	B	C
m	b	t
m	s	t
g	q	t
h	s	t
g	q	t
g	q	f
g	s	f
h	b	f
h	q	f
m	b	f

Example revisited

- Want to classify $A = m, B = q$
- $\Pr(A = m, B = q \mid C = t) = \Pr(A = m, B = q \mid C = f) = 0$

- $\Pr(A = m \mid C = t) = 2/5$
- $\Pr(B = q \mid C = t) = 2/5$

$$\frac{2}{5} \cdot \frac{2}{5} \cdot \frac{1}{2} \rightarrow \frac{2}{25}$$

- $\Pr(A = m \mid C = f) = 1/5$
- $\Pr(B = q \mid C = f) = 2/5$
- $\Pr(A = m \mid C = t) \cdot \Pr(B = q \mid C = t) \cdot \Pr(C = t) = 2/25$

$$P(A=m, B=q \mid C=t)$$

A	B	C
m	b	t
m	s	t
g	q	t
h	s	t
g	q	t
g	q	f
g	s	f
h	b	f
h	q	f
m	b	f

Example revisited

- Want to classify $A = m, B = q$
- $Pr(A = m, B = q \mid C = t) = Pr(A = m, B = q \mid C = f) = 0$
- $Pr(A = m \mid C = t) = 2/5$
- $Pr(B = q \mid C = t) = 2/5$
- $Pr(A = m \mid C = f) = 1/5$
- $Pr(B = q \mid C = f) = 2/5$
- $Pr(A = m \mid C = t) \cdot Pr(B = q \mid C = t) \cdot Pr(C = t) = 2/25$
- $Pr(A = m \mid C = f) \cdot Pr(B = q \mid C = f) \cdot Pr(C = f) = 1/25$

A	B	C
m	b	t
m	s	t
g	q	t
h	s	t
g	q	t
g	q	f
g	s	f
h	b	f
h	q	f
m	b	f

$$\frac{1}{5} \times \frac{2}{5} \times \frac{1}{2}$$

Example revisited

- Want to classify $A = m, B = q$
- $Pr(A = m, B = q \mid C = t) = Pr(A = m, B = q \mid C = f) = 0$
- $Pr(A = m \mid C = t) = 2/5$
- $Pr(B = q \mid C = t) = 2/5$
- $Pr(A = m \mid C = f) = 1/5$
- $Pr(B = q \mid C = f) = 2/5$
- $Pr(A = m \mid C = t) \cdot Pr(B = q \mid C = t) \cdot Pr(C = t) = 2/25$
- $Pr(A = m \mid C = f) \cdot Pr(B = q \mid C = f) \cdot Pr(C = f) = 1/25$
- Hence predict $C = t$

A	B	C
m	b	t
m	s	t
g	q	t
h	s	t
g	q	t
g	q	f
g	s	f
h	b	f
h	q	f
m	b	f

Zero counts

- Suppose $A = a$ never occurs in the test set with $C = c$

$$P(A=a|C=c)$$

Zero counts

- Suppose $A = a$ never occurs in the test set with $C = c$
- Setting $Pr(A = a \mid C = c) = 0$ wipes out any product $\prod_{i=1}^k Pr(A_i = a_i \mid C = c)$ in which this term appears

Zero counts

- Suppose $A = a$ never occurs in the ^{training} test set with $C = c$

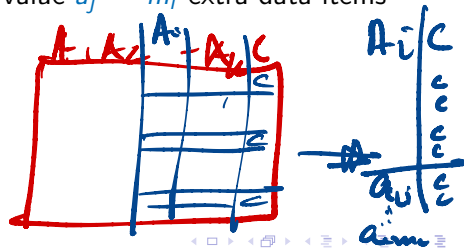
- Setting $Pr(A = a \mid C = c) = 0$ wipes out any product $\prod_{i=1}^k Pr(A_i = a_i \mid C = c)$ in which this term appears

- Assume A_i takes m_i values $\{a_{i1}, \dots, a_{im_i}\}$

$A_i = a_{ik}$ does not occur for $C = c$

Zero counts

- Suppose $A = a$ never occurs in the test set with $C = c$
- Setting $Pr(A = a \mid C = c) = 0$ wipes out any product $\prod_{i=1}^k Pr(A_i = a_i \mid C = c)$ in which this term appears
- Assume A_i takes m_i values $\{a_{i1}, \dots, a_{im_i}\}$
- “Pad” training data with one sample for each value a_j — m_i extra data items



Zero counts

- Suppose $A = a$ never occurs in the test set with $C = c$

- Setting $Pr(A = a \mid C = c) = 0$ wipes out any product $\prod_{i=1}^k Pr(A_i = a_i \mid C = c)$ in which this term appears

- Assume A_i takes m_i values $\{a_{i1}, \dots, a_{im_i}\}$

- “Pad” training data with one sample for each value a_j — m_i extra data items

- Adjust $Pr(A_i = a_i \mid C = c_j)$ to $\frac{n_{ij} + 1}{n_j + m_i}$

where

- n_{ij} is number of samples with $A_i = a_i, C = c_j$
- n_j is number of samples with $C = c_j$

$$P(A = a_j \mid C = c_j) = \frac{n_{ij} + 1}{n_j + m_i}$$

A_i	C
a_i	c_j
\vdots	c_j
\vdots	c_j

$n_j + m_i$

- Laplace's law of succession

$$Pr(A_i = a_i \mid C = c_j) = \frac{n_{ij} + 1}{n_j + m_i}$$

Smoothing

- Laplace's law of succession

$$Pr(A_i = a_i \mid C = c_j) = \frac{n_{ij} + 1}{n_j + m_i}$$

+1 pertains
count

- More generally, Lidstone's law of succession, or smoothing

$$Pr(A_i = a_i \mid C = c_j) = \frac{n_{ij} + \lambda}{n_j + \lambda m_i}$$

$\lambda \leq 1$

- Laplace's law of succession

$$Pr(A_i = a_i \mid C = c_j) = \frac{n_{ij} + 1}{n_j + m_i}$$

$$A=m, B=q$$

- More generally, Lidstone's law of succession, or smoothing

$$m, q, t$$

$$Pr(A_i = a_i \mid C = c_j) = \frac{n_{ij} + \lambda}{n_j + \lambda m_i}$$

- $\lambda = 1$ is Laplace's law of succession

Text classification

- Classify text documents using topics

Text classification

- Classify text documents using topics
- Useful for automatic segregation of newsfeeds, other internet content

Text classification

- Classify text documents using topics
- Useful for automatic segregation of newsfeeds, other internet content
- Training data has a unique topic label per document — e.g., Sports, Politics, Entertainment

Text classification

- Classify text documents using topics
- Useful for automatic segregation of newsfeeds, other internet content
- Training data has a unique topic label per document — e.g., Sports, Politics, Entertainment
- Want to use a naïve Bayes classifier

Text classification

- Classify text documents using topics
- Useful for automatic segregation of newsfeeds, other internet content
- Training data has a unique topic label per document — e.g., Sports, Politics, Entertainment
- Want to use a naïve Bayes classifier
- Need to define a generative model

Text classification

- Classify text documents using topics
- Useful for automatic segregation of newsfeeds, other internet content
- Training data has a unique topic label per document — e.g., Sports, Politics, Entertainment
- Want to use a naïve Bayes classifier
- Need to define a generative model
- How do we represent documents?

Set of words model

- Each document is a **set** of words over a vocabulary $V = \{w_1, w_2, \dots, w_m\}$

ignore
multiplicity

ignore
order

A hit B
B hit A

Set of words model

- Each document is a **set** of words over a vocabulary $V = \{w_1, w_2, \dots, w_m\}$
- Topics come from a set $C = \{c_1, c_2, \dots, c_k\}$

↓
Focus on words of interest

x the, a, in ...

✓ HP Pavilion,



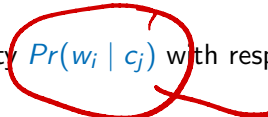
IBM ThinkPad

CZ22X7

Set of words model

- Each document is a **set** of words over a vocabulary $V = \{w_1, w_2, \dots, w_m\}$
- Topics come from a set $C = \{c_1, c_2, \dots, c_k\}$
- Each topic c has probability $Pr(c)$ **4**

Set of words model

- Each document is a **set** of words over a vocabulary $V = \{w_1, w_2, \dots, w_m\}$
- Topics come from a set $C = \{c_1, c_2, \dots, c_k\}$ 
- Each topic c has probability $Pr(c)$ 
- Each word $w_i \in V$ has conditional probability $Pr(w_i | c_j)$ with respect to each $c_j \in C$ 

Set of words model

- Each document is a **set** of words over a vocabulary $V = \{w_1, w_2, \dots, w_m\}$
- Topics come from a set $C = \{c_1, c_2, \dots, c_k\}$
- Each topic c has probability $Pr(c)$
- Each word $w_i \in V$ has conditional probability $Pr(w_i | c_j)$ with respect to each $c_j \in C$
- Generating a random document d
 - Choose a topic c with probability $Pr(c)$
 - For each $w \in V$, toss a coin, include w in d with probability $Pr(w | c)$

Set of words model

- Each document is a **set** of words over a vocabulary $V = \{w_1, w_2, \dots, w_m\}$
- Topics come from a set $C = \{c_1, c_2, \dots, c_k\}$
- Each topic c has probability $Pr(c)$
- Each word $w_i \in V$ has conditional probability $Pr(w_i | c_j)$ with respect to each $c_j \in C$
- Generating a random document d
 - Choose a topic c with probability $Pr(c)$
 - For each $w \in V$, toss a coin, include w in d with probability $Pr(w | c)$

$$Pr(d | c) = \prod_{w_i \in d} Pr(w_i | c) \prod_{w_i \notin d} (1 - Pr(w_i | c))$$

$d \subseteq V$

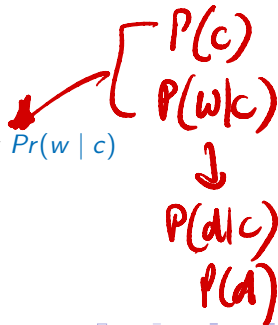
$$p^k (1-p)^{n-k}$$

Set of words model

- Each document is a **set** of words over a vocabulary $V = \{w_1, w_2, \dots, w_m\}$
- Topics come from a set $C = \{c_1, c_2, \dots, c_k\}$
- Each topic c has probability $Pr(c)$
- Each word $w_i \in V$ has conditional probability $Pr(w_i | c_j)$ with respect to each $c_j \in C$
- Generating a random document d
 - Choose a topic c with probability $Pr(c)$
 - For each $w \in V$, toss a coin, include w in d with probability $Pr(w | c)$

$$\blacksquare Pr(d | c) = \prod_{w_i \in d} Pr(w_i | c) \prod_{w_i \notin d} (1 - Pr(w_i | c))$$

$$\blacksquare Pr(d) = \sum_{c \in C} Pr(d | c) \cdot Pr(c)$$



Naïve Bayes classifier

- Training set $D = \{d_1, d_2, \dots, d_n\}$
 - Each $d_i \subseteq V$ is assigned a unique label from C

	w_1	w_2	w_m	c		
d_i	0	1	...	0	0	...	1	7

$$\text{Set } X \subseteq U \quad \equiv \quad f_x: U \rightarrow \{0, 1\}$$

$$f_x(x) = 1 \quad \forall \quad x \in X$$

Naïve Bayes classifier

- Training set $D = \{d_1, d_2, \dots, d_n\}$
 - Each $d_i \subseteq V$ is assigned a unique label from C
- $Pr(c_j)$ is fraction of D labelled c_j

Naïve Bayes classifier

- Training set $D = \{d_1, d_2, \dots, d_n\}$
 - Each $d_i \subseteq V$ is assigned a unique label from C
- $Pr(c_j)$ is fraction of D labelled c_j
- $Pr(w_i | c_j)$ is fraction of documents labelled c_j in which w_i appears

$w_i = 1$

Naïve Bayes classifier

- Training set $D = \{d_1, d_2, \dots, d_n\}$
 - Each $d_i \subseteq V$ is assigned a unique label from C
- $Pr(c_j)$ is fraction of D labelled c_j
- $Pr(w_i \mid c_j)$ is fraction of documents labelled c_j in which w_i appears
- Given a new document $d \subseteq V$, we want to compute $\arg \max_c Pr(c \mid d)$

Naïve Bayes classifier

- Training set $D = \{d_1, d_2, \dots, d_n\}$
 - Each $d_i \subseteq V$ is assigned a unique label from C
- $Pr(c_j)$ is fraction of D labelled c_j
- $Pr(w_i | c_j)$ is fraction of documents labelled c_j in which w_i appears
- Given a new document $d \subseteq V$, we want to compute $\arg \max_c Pr(c | d)$
- By Bayes' rule, $\underline{Pr(c | d)} = \frac{Pr(d | c)Pr(c)}{Pr(d)}$ $\rightarrow = \sum_c P(d|c) \cdot P(c)$
 - As usual, discard the common denominator and compute $\arg \max_c Pr(d | c)Pr(c)$

Naïve Bayes classifier

- Training set $D = \{d_1, d_2, \dots, d_n\}$
 - Each $d_i \subseteq V$ is assigned a unique label from C
- $Pr(c_j)$ is fraction of D labelled c_j
- $Pr(w_i | c_j)$ is fraction of documents labelled c_j in which w_i appears
- Given a new document $d \subseteq V$, we want to compute $\arg \max_c Pr(c | d)$
- By Bayes' rule, $Pr(c | d) = \frac{Pr(d | c)Pr(c)}{Pr(d)}$
 - As usual, discard the common denominator and compute $\arg \max_c Pr(d | c)Pr(c)$
- Recall $Pr(d | c) = \prod_{w_i \in d} Pr(w_i | c) \prod_{w_i \notin d} (1 - Pr(w_i | c))$



Bag of words model

- Each document is a **multiset** or **bag** of words over a vocabulary

$$V = \{w_1, w_2, \dots, w_m\}$$

- Count multiplicities of each word

Bag of words model

- Each document is a **multiset** or **bag** of words over a vocabulary

$$V = \{w_1, w_2, \dots, w_m\}$$

- Count multiplicities of each word
- As before
 - Each topic c has probability $Pr(c)$
 - Each word $w_i \in V$ has conditional probability $Pr(w_i | c_j)$ with respect to each $c_j \in C$
 - Note that $\sum_{i=1}^m Pr(w_i | c_j) = 1$
 - Assume document length is independent of the class

Bag of words model

- Generating a random document d

- Choose a document length ℓ with $Pr(\ell)$ ←

- Choose a topic c with probability $Pr(c)$

- Recall $|V| = m$.

- To generate a single word, throw an m -sided die that displays w with probability

$Pr(w | c)$

- Repeat ℓ times

$$V = \{w_1, \dots, w_m\}$$
$$\downarrow$$
$$P(w_i | c_j)$$



ℓ -words

u_1

u_2

...

u_e

$w_i?$



Bag of words model

- Generating a random document d
 - Choose a document length ℓ with $Pr(\ell)$
 - Choose a topic c with probability $Pr(c)$
 - Recall $|V| = m$.
 - To generate a single word, throw an m -sided die that displays w with probability $Pr(w | c)$
 - Repeat ℓ times
- Let n_j be the number of occurrences of w_j in d

Bag of words model

- Generating a random document d
 - Choose a document length ℓ with $Pr(\ell)$
 - Choose a topic c with probability $Pr(c)$
 - Recall $|V| = m$.
 - To generate a single word, throw an m -sided die that displays w with probability $Pr(w | c)$
 - Repeat ℓ times
- Let n_j be the number of occurrences of w_j in d

$$\text{--- } \underline{\underline{Pr(d | c)}} = \text{Pr}(\ell) \ell! \prod_{j=1}^m \frac{Pr(w_j | c)^{n_j}}{n_j!}$$

7h 3+
hhhhh e h
e h —

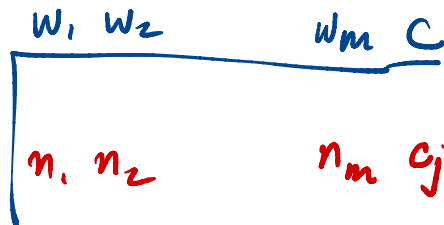
$$P^k$$
$$P(w_j | c)^{n_j}$$

Parameter estimation

- Training set $D = \{d_1, d_2, \dots, d_n\}$
 - Each d_i is a multiset over V of size ℓ_i

Parameter estimation

- Training set $D = \{d_1, d_2, \dots, d_n\}$
 - Each d_i is a multiset over V of size ℓ_i
- As before, $Pr(c_j)$ is fraction of D labelled c_j



Parameter estimation

- Training set $D = \{d_1, d_2, \dots, d_n\}$
 - Each d_i is a multiset over V of size ℓ_i
- As before, $\Pr(c_j)$ is fraction of D labelled c_j
- $\Pr(w_i | c_j)$ — fraction of occurrences of w_i over documents $D_j \subseteq D$ labelled c_j
 - n_{id} — occurrences of w_i in d

$$\Pr(w_i | c_j) = \frac{\sum_{d \in D_j} n_{id}}{\sum_{t=1}^m \sum_{d \in D_j} n_{td}}$$

D_j = documents
labelled c_j

Parameter estimation

- Training set $D = \{d_1, d_2, \dots, d_n\}$
 - Each d_i is a multiset over V of size ℓ_i
- As before, $Pr(c_j)$ is fraction of D labelled c_j
- $Pr(w_i | c_j)$ — fraction of occurrences of w_i over documents $D_j \subseteq D$ labelled c_j
 - n_{id} — occurrences of w_i in d

$$Pr(w_i | c_j) = \frac{\sum_{d \in D_j} n_{id}}{m \sum_{t=1} \sum_{d \in D_j} n_{td}} = \frac{\sum_{d \in D} n_{id} Pr(c_j | d)}{m \sum_{t=1} \sum_{d \in D} n_{td} Pr(c_j | d)}$$

$$Pr(c_j | d)$$

$$d \in D_j \Rightarrow c_j = 1$$

$$d \notin D_j \Rightarrow c_j = 0$$

$$\text{since } Pr(c_j | d) = \begin{cases} 1 & \text{if } d \in D_j, \\ 0 & \text{otherwise} \end{cases}$$

$$d \in D \cap D_j$$

$$Pr(c_j | d) = 0$$

$$d \in D_j$$

Classification

$$\blacksquare \Pr(c \mid d) = \frac{\Pr(d \mid c) \Pr(c)}{\Pr(d)}$$

↑ ↑

Classification

- $Pr(c | d) = \frac{Pr(d | c) Pr(c)}{Pr(d)}$
- Want $\arg \max_c Pr(c | d)$

Classification

- $Pr(c | d) = \frac{Pr(d | c) Pr(c)}{Pr(d)}$
- Want $\arg \max_c Pr(c | d)$
- As before, discard the denominator $Pr(d)$

Classification

- $Pr(c | d) = \frac{Pr(d | c) Pr(c)}{Pr(d)}$
- Want $\arg \max_c Pr(c | d)$
- As before, discard the denominator $Pr(d)$
- Recall, $Pr(d | c) = \cancel{Pr(\ell)} \ell! \prod_{j=1}^m \frac{Pr(w_j | c)^{n_j}}{\underline{n_j!}}$, where $|d| = \ell$

Classification

- $Pr(c | d) = \frac{Pr(d | c) Pr(c)}{Pr(d)}$
- Want $\arg \max_c Pr(c | d)$
- As before, discard the denominator $Pr(d)$
- Recall, $Pr(d | c) = Pr(\ell) \ell! \prod_{j=1}^m \frac{Pr(w_j | c)^{n_j}}{n_j!}$, where $|d| = \ell$
- Discard $Pr(\ell), \ell!$ since they do not depend on c

Classification

- $Pr(c | d) = \frac{Pr(d | c) Pr(c)}{Pr(d)}$

- Want $\arg \max_c Pr(c | d)$

- As before, discard the denominator $Pr(d)$

- Recall, $Pr(d | c) = Pr(\ell) \ell! \prod_{j=1}^m \frac{Pr(w_j | c)^{n_j}}{n_j!}$, where $|d| = \ell$

- Discard $Pr(\ell), \ell!$ since they do not depend on c

- Compute $\arg \max_c Pr(c) \prod_{j=1}^m \frac{Pr(w_j | c)^{n_j}}{n_j!}$

d_1 dog¹ dog² cat³ $c=0$
 d_2 dog⁴ cat⁵ wolf⁶ $c=0$

$$P(\text{dog} | c=0) = \frac{2}{2} = 1$$

↳ Set of words

$$P(\text{dog} | c=0) = \frac{3}{6} \begin{matrix} \text{dog} \\ \text{total} \end{matrix}$$