

Bayesian Learning Method

Nguyen Van Vinh
UET-Hanoi VNU

Content

- Bayesian Learning (NB) method
- Examples for NB
- Application (Text Classification, Spam Mail)

Introduction

Thomas Bayes (c. 1702 – 17 April 1761)
was a British mathematician and
Presbyterian minister. (wikipedia)

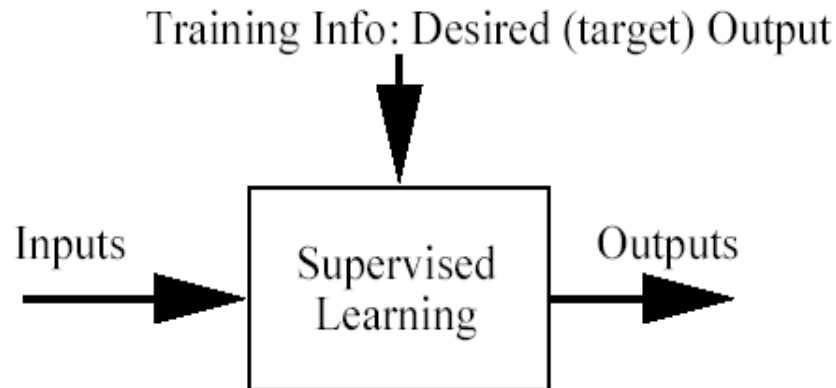


Today we learn:

- Bayesian classification
 - E.g. How to decide if a patient is ill or healthy, based on
 - A probabilistic model of the observed data
 - Prior knowledge

Classification problem

- Training data: examples of the form $(d, h(d))$
 - where d are the data objects to classify (inputs)
 - and $h(d)$ are the correct class info for d , $h(d) \in \{1, \dots, K\}$
- Goal: given d_{new} , provide $h(d_{\text{new}})$



Why Bayesian?

- Provides practical learning algorithms
 - E.g. Naïve Bayes
- Prior knowledge and observed data can be combined
- It is a generative (model based) approach, which offers a useful conceptual framework
 - E.g. sequences could also be classified, based on a probabilistic model specification
 - Any kind of objects can be classified, based on a probabilistic model specification

Bayes' Rule

$$p(h | d) = \frac{P(d | h)P(h)}{P(d)}$$

Understanding Bayes' rule

d = data

h = hypothesis (model)

- rearranging

$$p(h | d)P(d) = P(d | h)P(h)$$

$$P(d, h) = P(d, h)$$

the same joint probability

on both sides

Who is who in Bayes' rule

$P(h)$: prior belief (probability of hypothesis h before seeing any data)

$P(d | h)$: likelihood (probability of the data if the hypothesis h is true)

$P(d) = \sum_h P(d | h)P(h)$: data evidence (marginal probability of the data)

$P(h | d)$: posterior (probability of hypothesis h after having seen the data d)

Does patient have cancer or not?

- A patient takes a lab test and the result comes back positive. It is known that the test returns a correct positive result in only 98% of the cases and a correct negative result in only 97% of the cases. Furthermore, only 0.008 of the entire population has this disease.
 1. What is the probability that this patient has cancer?
 2. What is the probability that he does not have cancer?
 3. What is the diagnosis?

$hypothesis1: 'cancer'$
 $hypothesis2: '\neg cancer'$ } *hypothesis space H*
 $- data: '+'$

$$1. P(cancer | +) = \frac{P(+ | cancer)P(cancer)}{P(+)} = \frac{\dots\dots\dots}{\dots\dots\dots} = \dots\dots\dots$$

$$P(+ | cancer) = 0.98$$

$$P(cancer) = 0.008$$

$$P(+)= P(+ | cancer)P(cancer) + P(+ | \neg cancer)P(\neg cancer)$$

$$= \dots\dots\dots$$

$$P(+ | \neg cancer) = 0.03$$

$$P(\neg cancer) = \dots\dots\dots$$

$$2. P(\neg cancer | +) = \dots\dots\dots$$

3. *Diagnosis??*

Choosing Hypotheses

- *Maximum Likelihood* hypothesis:

$$h_{ML} = \arg \max_{h \in H} P(d | h)$$

- Generally we want the most probable hypothesis given training data. This is the *maximum a posteriori* hypothesis:
 - Useful observation: it does not depend on the denominator $P(d)$

$$h_{MAP} = \arg \max_{h \in H} P(h | d)$$

Now we compute the diagnosis

- To find the Maximum Likelihood hypothesis, we evaluate $P(d|h)$ for the data d , which is the positive lab test and chose the hypothesis (diagnosis) that maximises it:

$$P(+ | cancer) = \dots\dots\dots$$

$$P(+ | \neg cancer) = \dots\dots\dots$$

$$\Rightarrow \text{Diagnosis} : h_{ML} = \dots\dots\dots$$

- To find the Maximum A Posteriori (**MAP**) hypothesis, we evaluate $P(d|h)P(h)$ for the data d , which is the positive lab test and chose the hypothesis (diagnosis) that maximises it. This is the same as choosing the hypotheses gives the higher posterior probability.

$$P(+ | cancer)P(cancer) = \dots\dots\dots$$

$$P(+ | \neg cancer)P(\neg cancer) = \dots\dots\dots$$

$$\Rightarrow \text{Diagnosis} : h_{MAP} = \dots\dots\dots$$

Bayesian decision theory

- Let x be the value predicted by the agent and x^* be the true value of X .
- The agent has a **loss function**, which is 0 if $x = x^*$ and 1 otherwise
- Expected loss for predicting x :

$$\sum_{x^*} L(x, x^*) P(x^* | e)$$

- What is the estimate of X that minimizes the expected loss?
 - The one that has the greatest posterior probability $P(x|e)$
 - This is called the **Maximum a Posteriori (MAP)** decision

MAP decision

- Value x of X that has the highest posterior probability given the evidence $E = e$:

$$x^* = \arg \max_x P(X = x \mid E = e) = \frac{P(E = e \mid X = x)P(X = x)}{P(E = e)}$$

$$\propto \arg \max_x P(E = e \mid X = x)P(X = x)$$

$$\underbrace{P(x \mid e)}_{\text{posterior}} \propto \underbrace{P(e \mid x)}_{\text{likelihood}} \underbrace{P(x)}_{\text{prior}}$$

- Maximum likelihood (ML) decision:

$$x^* = \arg \max_x P(e \mid x)$$

Naïve Bayes Classifier

- What can we do if our data d has several attributes?
- Naïve Bayes assumption: Attributes that describe data instances are conditionally independent given the classification hypothesis

$$P(\mathbf{d} | h) = P(a_1, \dots, a_T | h) = \prod_t P(a_t | h)$$

- it is a simplifying assumption, obviously it may be violated in reality
 - in spite of that, it works well in practice
- The Bayesian classifier that uses the Naïve Bayes assumption and computes the MAP hypothesis is called Naïve Bayes classifier
- One of the most practical learning methods
- Successful applications:
 - Medical Diagnosis
 - Text classification

Example. 'Play Tennis' data

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes
Day11	Sunny	Mild	Normal	Strong	Yes
Day12	Overcast	Mild	High	Strong	Yes
Day13	Overcast	Hot	Normal	Weak	Yes
Day14	Rain	Mild	High	Strong	No

Naïve Bayes solution

Classify any new datum instance $\mathbf{x}=(a_1, \dots a_T)$ as:

$$h_{Naive\ Bayes} = \arg \max_h P(h)P(\mathbf{x} | h) = \arg \max_h P(h) \prod_t P(a_t | h)$$

- To do this based on training examples, we need to estimate the parameters from the training examples:

- For each target value (hypothesis) h

$$\hat{P}(h) := \text{estimate } P(h)$$

- For each attribute value a_t of each datum instance

$$\hat{P}(a_t | h) := \text{estimate } P(a_t | h)$$

Based on the examples in the table, classify the following datum \mathbf{x} :

$\mathbf{x}=(\text{Outl}=\text{Sunny}, \text{Temp}=\text{Cool}, \text{Hum}=\text{High}, \text{Wind}=\text{strong})$

- That means: Play tennis or not?

$$h_{NB} = \arg \max_{h \in [\text{yes}, \text{no}]} P(h)P(\mathbf{x} | h) = \arg \max_{h \in [\text{yes}, \text{no}]} P(h) \prod_t P(a_t | h)$$

$$= \arg \max_{h \in [\text{yes}, \text{no}]} P(h)P(\text{Outlook} = \text{sunny} | h)P(\text{Temp} = \text{cool} | h)P(\text{Humidity} = \text{high} | h)P(\text{Wind} = \text{strong} | h)$$

- Working:

$$P(\text{PlayTennis} = \text{yes}) = 9 / 14 = 0.64$$

$$P(\text{PlayTennis} = \text{no}) = 5 / 14 = 0.36$$

$$P(\text{Wind} = \text{strong} | \text{PlayTennis} = \text{yes}) = 3 / 9 = 0.33$$

$$P(\text{Wind} = \text{strong} | \text{PlayTennis} = \text{no}) = 3 / 5 = 0.60$$

etc.

$$P(\text{yes})P(\text{sunny} | \text{yes})P(\text{cool} | \text{yes})P(\text{high} | \text{yes})P(\text{strong} | \text{yes}) = 0.0053$$

$$P(\text{no})P(\text{sunny} | \text{no})P(\text{cool} | \text{no})P(\text{high} | \text{no})P(\text{strong} | \text{no}) = \mathbf{0.0206}$$

$$\Rightarrow \text{answer} : \text{PlayTennis}(x) = \text{no}$$

Example: Training Dataset

age	income	student	credit_rating	comp
<=30	high	no	fair	no
<=30	high	no	excellent	no
31...40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31...40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes
>40	medium	no	excellent	no

Class:

C1:buys_computer = 'yes'

C2:buys_computer = 'no'

Data sample

X = (**age <=30,**

Income = medium,

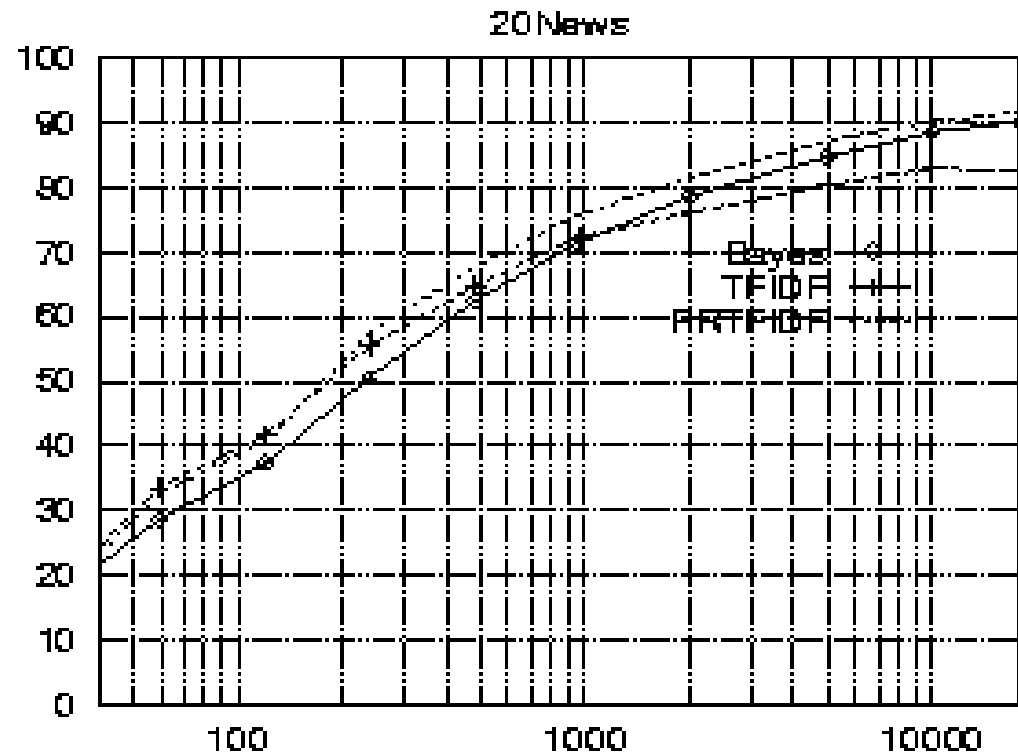
Student = yes

Credit_rating = Fair)

Learning to classify text

- Learn from examples which articles are of interest
- The attributes are the words
- Observe the Naïve Bayes assumption just means that we have a random sequence model within each class!
- NB classifiers are one of the most effective for this task
- Resources for those interested:
 - Tom Mitchell: Machine Learning (book) Chapter 6.

Results on a benchmark text corpus

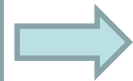


Accuracy vs. Training set size (1/3 withheld for test)

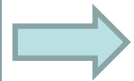
Learning and inference pipeline

Learning

Training Samples



Features



Training Labels



Training



Learned model

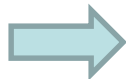
Learned model



Inference



Test Sample



Features



Prediction

Features

- A measurable variable that is (rather, should be) distinctive of something we want to model.
- We usually choose features that are useful to **identify** something, i.e., to do **classification**
 - **Ex:** Cô gái đó **rất đẹp** trong bữa tiệc hôm đó.
- We often need **several** features to adequately model something – *but not too many!*

Feature vectors

- **Values for several features of an observation can be put into a single vector**



Rush Limbaugh looks like if someone put a normal human being in landscape mode.

Reply Retweet Favorite More



BREAKING: Apple Maps projecting Barack Obama to win Brazil.

Reply Retweet Favorite More



If there was an award for most pessimistic, I probably wouldn't even be nominated.

Reply Retweet Favorite More

# proper nouns	# 1st person pronouns	# commas
2	0	0
5	0	0
0	1	1

Feature vectors

- Features should be useful in **discriminating** between categories.

Table 3: Features to be computed for each text

- Counts:
 - First person pronouns
 - Second person pronouns
 - Third person pronouns
 - Coordinating conjunctions
 - Past-tense verbs
 - Future-tense verbs
 - Commas
 - Colons and semi-colons
 - Dashes
 - Parentheses
 - Ellipses
 - Common nouns
 - Proper nouns
 - Adverbs
 - *wh*-words
 - Modern slang acronyms
 - Words all in upper case (at least 2 letters long)
- Average length of sentences (in tokens)
- Average length of tokens, excluding punctuation tokens (in characters)
- Number of sentences

Higher values → this person is referring to themselves (to their opinion, too?)

Higher values → looking forward to (or dreading) some future event?

Lower values → this tweet is more formal. Perhaps not overly sentimental?

Feature Representation

this movie was great! would watch again Positive

- Convert this example to a vector using *bag-of-words features*

[contains <i>the</i>] position 0	[contains <i>a</i>] position 1	[contains <i>was</i>] position 2	[contains <i>movie</i>] position 3	[contains <i>film</i>] position 4	...
$f(x) = [0$	0	1	1	0	$...$

- Very large vector space (size of vocabulary), sparse features
- Requires *indexing* the features (mapping them to axes)
- More sophisticated feature mappings possible (m-idf), as well as lots of other features: character n-grams, parts of speech, ...

Case study:

Text document classification

- **MAP decision:** assign a document to the class with the highest posterior $P(\text{class} \mid \text{document})$
- Example: spam classification
 - Classify a message as spam if $P(\text{spam} \mid \text{message}) > P(\neg \text{spam} \mid \text{message})$



Dear Sir.

First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret. ...



TO BE REMOVED FROM FUTURE MAILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.

99 MILLION EMAIL ADDRESSES
FOR ONLY \$99



Ok, I know this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

Case study:

Text document classification

- **MAP decision:** assign a document to the class with the highest posterior $P(\text{class} \mid \text{document})$
- We have $P(\text{class} \mid \text{document}) \propto P(\text{document} \mid \text{class})P(\text{class})$
- To enable classification, we need to be able to estimate the **likelihoods** $P(\text{document} \mid \text{class})$ for all classes and **priors** $P(\text{class})$

Naïve Bayes Representation

- Goal: estimate likelihoods $P(\text{document} \mid \text{class})$ and priors $P(\text{class})$
- Likelihood: **bag of words** representation
 - The document is a sequence of words (w_1, \dots, w_n)
 - The order of the words in the document is not important
 - Each word is conditionally independent of the others given document class



Dear Sir.

First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret. ...



TO BE REMOVED FROM FUTURE MAILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.

99 MILLION EMAIL ADDRESSES
FOR ONLY \$99



Ok, I know this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

Naïve Bayes Representation

- Goal: estimate likelihoods $P(\text{document} \mid \text{class})$ and priors $P(\text{class})$
- Likelihood: **bag of words** representation
 - The document is a sequence of words (w_1, \dots, w_n)
 - The order of the words in the document is not important
 - Each word is conditionally independent of the others given document class

$$P(\text{document} \mid \text{class}) = P(w_1, \dots, w_n \mid \text{class}) = \prod_{i=1}^n P(w_i \mid \text{class})$$

Bag of words illustration

2007-01-23: State of the Union Address

George W. Bush (2001-)

abandon accountable affordable afghanistan africa aided ally anbar armed army baghdad bless challenges chamber chaos
choices civilians coalition commanders commitment confident confront congressman constitution corps debates deduction
deficit deliver democratic deploy dikembe diplomacy disruptions earmarks economy einstein elections eliminates
expand extremists failing faithful families freedom fuel funding god haven ideology immigration impose
insurgents iran **iraq** islam julie lebanon love madam marine math medicare moderation neighborhoods nuclear offensive
palestinian payroll province pursuing **qaeda** radical regimes resolve retreat rieman sacrifices science sectarian senate
september shia stays strength students succeed sunni tax territories **terrorists** threats uphold victory
violence violent **war** washington weapons wesley

US Presidential Speeches Tag Cloud

<http://chir.ag/projects/preztags/>

Bag of words illustration

2007-01-23: State of the Union Address

George W. Bush (2001-)

1962-10-22: Soviet Missiles in Cuba

John F. Kennedy (1961-63)

abandon achieving adversaries aggression agricultural appropriate armaments **arms** assessments atlantic ballistic berlin
buildup burdens cargo college commitment communist constitution consumers cooperation crisis **cuba** dangers
declined **defensive** deficit depended disarmament divisions domination doubled **economic** education
elimination emergence endangered equals **europa** expand exports fact false family forum **freedom** fulfill gromyko
halt hazards **hemisphere** hospitals ideals **independent** industries inflation labor latin limiting minister **missiles**
modernization neglect **nuclear** oas obligation observer **offensive** peril pledged predicted purchasing quarantine **quote**
recession rejection republics retaliatory safeguard sites solution **soviet** space spur stability standby **strength**
surveillance **tax** territory treaty undertakings unemployment **war** warhead **weapons** welfare western widen withdraw

US Presidential Speeches Tag Cloud

<http://chir.ag/projects/preztags/>

Bag of words illustration

2007-01-23: State of the Union Address

George W. Bush (2001-)

1962-10-22: Soviet Missiles in Cuba

John F. Kennedy (1961-63)

1941-12-08: Request for a Declaration of War

Franklin D. Roosevelt (1933-45)

abandoning acknowledge aggression aggressors airplanes armaments **armed** army assault assembly authorizations bombing
britain british cheerfully claiming constitution curtail december defeats defending delays democratic dictators disclose
economic empire endanger **facts** false forgotten fortunes france **freedom** fulfilled fullness fundamental gangsters
german germany **god** guam harbor hawaii **hemisphere** hint **hitler** hostilities immune improving indies innumerable
japanese
invasion **islands** isolate labor metals midst midway **navy** nazis obligation offensive
officially **pacific** partisanship patriotism pearl peril perpetrated perpetual philippine preservation privilege reject
repaired **resisting** retain revealing rumors seas soldiers speaks speedy **stamina** **strength** sunday sunk supremacy tanks taxes
treachery true tyranny undertaken victory **war** wartime washington

US Presidential Speeches Tag Cloud

<http://chir.ag/projects/preztags/>

Naïve Bayes Representation

- Goal: estimate likelihoods $P(\text{document} \mid \text{class})$ and $P(\text{class})$
- Likelihood: **bag of words** representation
 - The document is a sequence of words (w_1, \dots, w_n)
 - The order of the words in the document is not important
 - Each word is conditionally independent of the others given document class

$$P(\text{document} \mid \text{class}) = P(w_1, \dots, w_n \mid \text{class}) = \prod_{i=1}^n P(w_i \mid \text{class})$$

- Thus, the problem is reduced to estimating marginal likelihoods of individual words $P(w_i \mid \text{class})$

Parameter estimation

- Model parameters: feature likelihoods $P(\text{word} \mid \text{class})$ and priors $P(\text{class})$
 - How do we obtain the values of these parameters?

prior

spam:	0.33
\neg spam:	0.67

$P(\text{word} \mid \text{spam})$

the :	0.0156
to :	0.0153
and :	0.0115
of :	0.0095
you :	0.0093
a :	0.0086
with:	0.0080
from:	0.0075
...	

$P(\text{word} \mid \neg \text{spam})$

the :	0.0210
to :	0.0133
of :	0.0119
2002:	0.0110
with:	0.0108
from:	0.0107
and :	0.0105
a :	0.0100
...	

Parameter estimation

- Model parameters: feature likelihoods $P(\text{word} \mid \text{class})$ and priors $P(\text{class})$
 - How do we obtain the values of these parameters?
 - Need *training set* of labeled samples from both classes

$$P(\text{word} \mid \text{class}) = \frac{\text{\# of occurrences of this word in docs from this class}}{\text{total \# of words in docs from this class}}$$

- This is the *maximum likelihood* (ML) estimate, or estimate that maximizes the likelihood of the training data:

$$\prod_{d=1}^D \prod_{i=1}^{n_d} P(w_{d,i} \mid \text{class}_{d,i})$$

d : index of training document, i : index of a word

Parameter estimation

- Parameter estimate:

$$P(\text{word} \mid \text{class}) = \frac{\text{\# of occurrences of this word in docs from this class}}{\text{total \# of words in docs from this class}}$$

- Parameter smoothing: dealing with words that were never seen or seen too few times
 - **Laplacian smoothing:** pretend you have seen every vocabulary word one more time than you actually did

$$P(\text{word} \mid \text{class}) = \frac{\text{\# of occurrences of this word in docs from this class} + 1}{\text{total \# of words in docs from this class} + V}$$

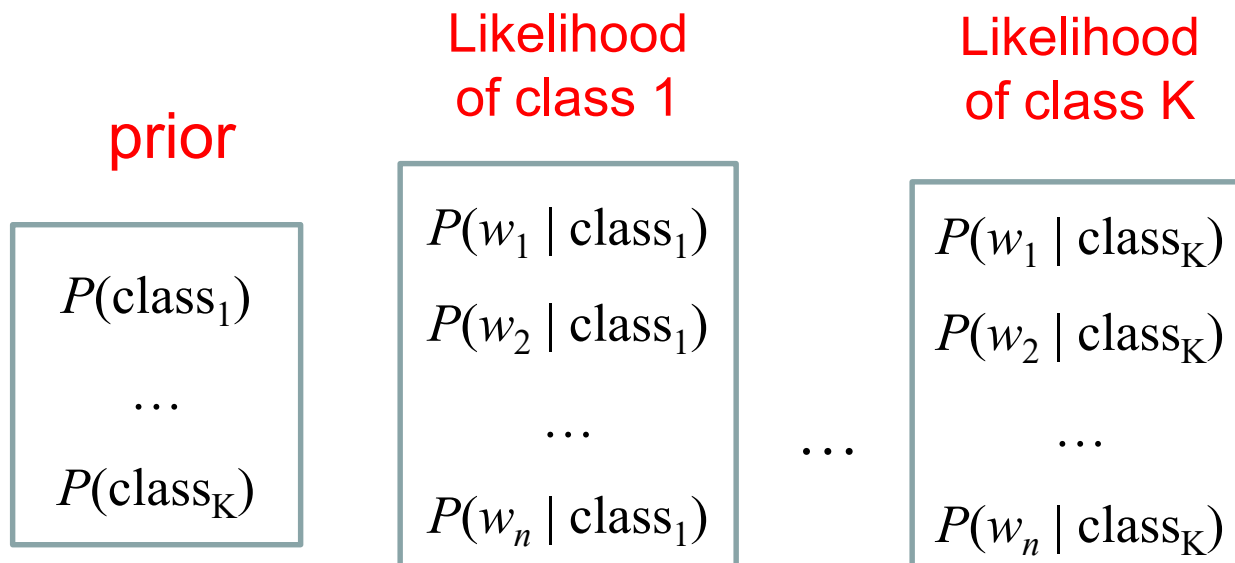
(V: total number of unique words)

Summary: Naïve Bayes for Document Classification

- Naïve Bayes model: assign the document to the class with the highest posterior

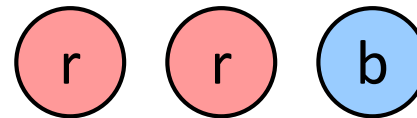
$$P(\text{class} \mid \text{document}) \propto P(\text{class}) \prod_{i=1}^n P(w_i \mid \text{class})$$

- Model parameters:



Laplace Smoothing

- Laplace's estimate:
 - Pretend you saw every outcome once more than you actually did



$$P_{ML}(X) =$$

$$\begin{aligned} P_{LAP}(x) &= \frac{c(x) + 1}{\sum_x [c(x) + 1]} \\ &= \frac{c(x) + 1}{N + |X|} \end{aligned}$$

$$P_{LAP}(X) =$$

Laplace Smoothing

- Laplace's estimate (extended):

- Pretend you saw every outcome k extra times

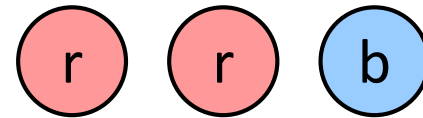
$$P_{LAP,k}(x) = \frac{c(x) + k}{N + k|X|}$$

- What's Laplace with $k = 0$?
- k is the **strength** of the prior

- Laplace for conditionals:

- Smooth each condition independently:

$$P_{LAP,k}(x|y) = \frac{c(x, y) + k}{c(y) + k|X|}$$



$$P_{LAP,0}(X) =$$

$$P_{LAP,1}(X) =$$

$$P_{LAP,100}(X) =$$

Summarization

- Bayes' rule can be turned into a classifier
- Maximum A Posteriori (MAP) hypothesis estimation incorporates prior knowledge; Max Likelihood doesn't
- Naive Bayes Classifier is a simple but effective Bayesian classifier for vector data (i.e. data with several attributes) that assumes that attributes are independent given the class.
- Bayesian classification is a generative approach to classification

Reference

- Slides of ML Course, University of Birmingham
- Textbook reading (contains details about using Naïve Bayes for text classification):
Tom Mitchell, Machine Learning (book), Chapter 6.
- Software: NB for classifying text:
<http://www-2.cs.cmu.edu/afs/cs/project/theo-11/www/naive-bayes.html>
<https://julianstier.com/posts/2021/01/text-classification-with-naive-bayes-in-numpy/>
- Useful reading for those interested to learn more about NB classification, beyond the scope of this module:
<http://www-2.cs.cmu.edu/~tom/NewChapters.html>