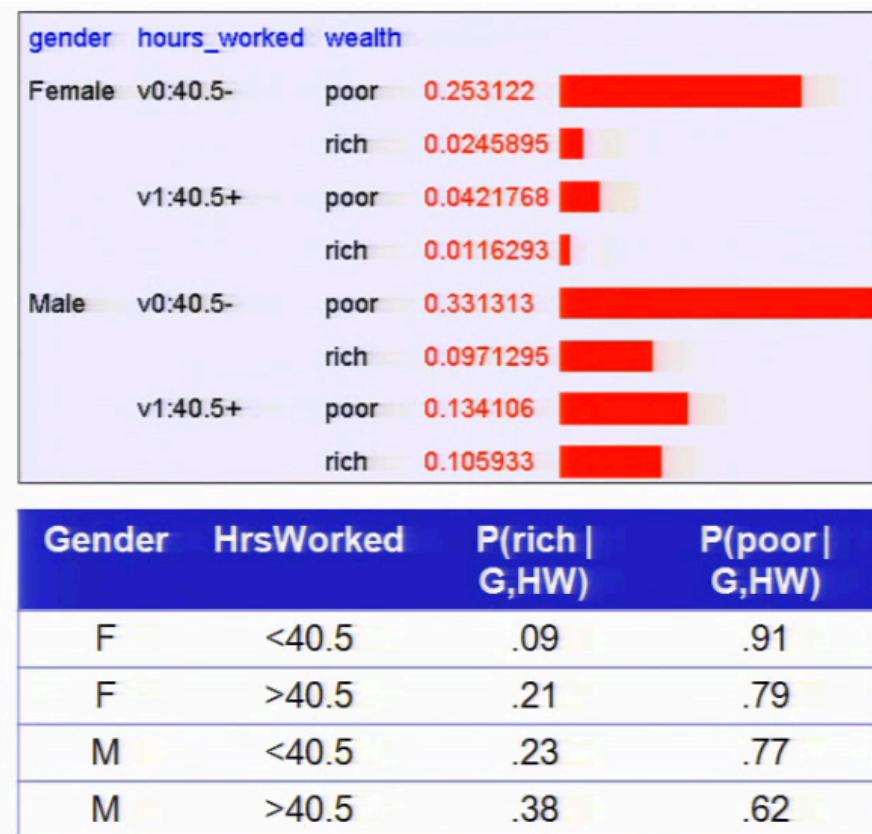


Bayesian Learning

Adapted from Tom Michell's Lecture

Learn Classifiers by Learning $P(X|Y)$

- Consider $Y=\text{Wealth}$, $X=\langle \text{Gender}, \text{HoursWork} \rangle$



How many Parameters must we estimate?

Suppose $X = \langle X_1, \dots, X_n \rangle$

where X_i and Y are boolean RV's

Gender	HrsWorked	$P(\text{rich} G, HW)$	$P(\text{poor} G, HW)$
F	<40.5	.09	.91
F	>40.5	.21	.79
M	<40.5	.23	.77
M	>40.5	.38	.62

To estimate $P(Y | X_1, X_2, \dots, X_n)$

If we have 30 X_i 's instead of 2?



Bayes Rules

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

Which is shorthand for:

$$(\forall i, j) P(Y = y_i | X = x_j) = \frac{P(X = x_j | Y = y_i)P(Y = y_i)}{P(X = x_j)}$$

Equivalently:

$$(\forall i, j) P(Y = y_i | X = x_j) = \frac{P(X = x_j | Y = y_i)P(Y = y_i)}{\sum_k P(X = x_j | Y = y_k)P(Y = y_k)}$$



Can we reduce parameters using Bayes Rule?

Suppose $X = \langle X_1, \dots, X_n \rangle$
where X_i and Y are boolean RV's

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$



Bayesian Learning



Conditional Independence

Definition: X is conditionally independent of Y given Z, if the probability distribution governing X is independent of the value of Y, given the value of Z

$$(\forall i, j, k) P(X = x_i | Y = y_j, Z = z_k) = P(X = x_i | Z = z_k)$$

Which we often write

$$P(X|Y, Z) = P(X|Z)$$

E.g.,

$$P(Thunder|Rain, Lightning) = P(Thunder|Lightning)$$



Naïve Bayes uses assumption that the X_i are conditionally independent, given Y

Given this assumption, then:

$$\begin{aligned} P(X_1, X_2|Y) &= P(X_1|X_2, Y)P(X_2|Y) \\ &= P(X_1|Y)P(X_2|Y) \end{aligned}$$

in general: $P(X_1\dots X_n|Y) = \prod_i P(X_i|Y)$

How many parameters to describe $P(X_1\dots X_n|Y)$? $P(Y)$?

- Without conditional indep assumption?
- With conditional indep assumption?

An Overview Of Naïve Bayes

- **Bayes rule :**

$$P(Y = y_k | X_1 \dots X_n) = \frac{P(Y = y_k)P(X_1 \dots X_n | Y = y_k)}{\sum_j P(Y = y_j)P(X_1 \dots X_n | Y = y_j)}$$

- **Assuming conditional independence among X_i 's:**

$$P(Y = y_k | X_1 \dots X_n) = \frac{P(Y = y_k) \prod_i P(X_i | Y = y_k)}{\sum_j P(Y = y_j) \prod_i P(X_i | Y = y_j)}$$

- **So, classification rule for $X^{new} = \langle X_1, \dots, X_n \rangle$ is :**

$$Y^{new} \leftarrow \arg \max_{y_k} P(Y = y_k) \prod_i P(X_i^{new} | Y = y_k)$$

Naïve Bayes Algorithm – discrete X_i

- **Train Naïve Bayes (examples)**

for each* value y_k

estimate $\pi_k \equiv P(Y = y_k)$

for each* value x_{ij} of each attribute X_i

estimate $\theta_{ijk} \equiv P(X_i = x_{ij} | Y = y_k)$

- **Classify (X^{new})**

$$Y^{new} \leftarrow \arg \max_{y_k} P(Y = y_k) \prod_i P(X_i^{new} | Y = y_k)$$

$$Y^{new} \leftarrow \arg \max_{y_k} \pi_k \prod_i \theta_{ijk}$$

*probabilities must sum to 1, so need estimate only n-1 of these ...

Estimating Parameters : Y, X_i discrete-valued

- Maximum likelihood estimates (MLE's)

$$\hat{\pi}_k = \hat{P}(Y = y_k) = \frac{\# D\{Y = y_k\}}{|D|}$$

$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_{ij} \mid Y = y_k) = \frac{\# D\{X_i = x_{ij} \wedge Y = y_k\}}{\# D\{Y = y_k\}}$$

Number of items in
dataset D for which $Y = y_k$

Example: Living in campus vicinity? $P(C|T,B)$

- $C=1$ iff currently living in Cisitu, Sangkuring, Pelesiran, Kebon Binatang
- $T=1$ iff go to campus on foot
- $B=1$ iff home town is Bandung

What parameters must we estimate?

Example: Living in campus vicinity? $P(C|T,B)$

- $C=1$ iff currently living in Cisitu, Sangkuring, Pelesiran, Kebon Binatang, or in another area within 2 km.
- $T=1$ iff go to campus on foot
- $B=1$ iff home town is Bandung

- | | |
|--------------------|--------------------|
| • $P(C=1) :$ | • $P(C=0) :$ |
| • $P(T=1 C=1) :$ | • $P(T=0 C=1) :$ |
| • $P(T=1 C=0) :$ | • $P(T=0 C=0) :$ |
| • $P(B=1 C=1) :$ | • $P(B=0 C=1) :$ |
| • $P(B=1 C=0) :$ | • $P(B=0 C=0) :$ |

Naïve Bayes : Sublety #1

If unlucky, our MLE estimate for $P(X_i | Y)$ might be zero (e.g., $X_i = \text{Birthday_Is_January_30_1990}$)

- Why worry about just one parameter out of many?
- What can be done to avoid this ?

Estimating Parameters

- Maximum Likelihood Estimate (MLE) : choose θ that maximize probability of observed data D

$$\hat{\theta} = \arg \max_{\theta} P(D | \theta)$$

- Maximum a Posteriori (MAP)P estimatee: choose that is most probable given prior probability and the data

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} P(\theta | D) \\ &= \arg \max_{\theta} = \frac{P(D | \theta)P(\theta)}{P(D)}\end{aligned}$$

Estimating Parameters : Y, X_i discrete-valued

- Maximum likelihood estimates :

$$\hat{\pi}_k = \hat{P}(Y = y_k) = \frac{\# D\{Y = y_k\}}{|D|}$$

$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_{ij} \mid Y = y_k) = \frac{\# D\{X_i = x_{ij} \wedge Y = y_k\}}{\# D\{Y = y_k\}}$$

- MAP estimates (Beta, Dirichlet priors)

$$\hat{\pi}_k = \hat{P}(Y = y_k) = \frac{\# D\{Y = y_k\} + \alpha_k}{|D| + \sum_m \alpha_m}$$

Only difference :
“imaginary” examples

$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_{ij} \mid Y = y_k) = \frac{\# D\{X_i = x_{ij} \wedge Y = y_k\} + \alpha'_k}{\# D\{Y = y_k\} + \sum_m \alpha'_m}$$

Another way to view Naïve Bayes (Boolean Y) :
Decision rule : is this quantity greater or less than 1 ?

$$\frac{P(Y=1 | X_1 \dots X_n)}{P(Y=0 | X_1 \dots X_n)} = \frac{P(Y=1) \prod_i P(X_i | Y=1)}{P(Y=0) \prod_i P(X_i | Y=0)}$$

Naïve Bayes : classifying text documents

- Classify which email are spam ?
- Classify which emails promise an attachment ?

I am pleased to announce that Bob Frederking of the Language Technologies Institute is our new Associate Dean for Graduate Programs. In this role, he oversees the many issues that arise with our multiple masters and PhD programs. Bob brings to this position considerable experience with the masters and PhD programs in the LTI.

I would like to thank Frank Pfenning, who has served ably in this role for the past two years.

Randal E. Bryant
Dean and University Professor

How shall we represent text documents for Naïve Bayes ?

Learning to classify documents : $P(Y|X)$

- **Y discrete valued.**
 - e.g., Spam or not
- **$X = \langle X_1, X_2, \dots, X_n \rangle = \text{document}$**

I am pleased to announce that Bob Frederking of the Language Technologies Institute is our new Associate Dean for Graduate Programs. In this role, he oversees the many issues that arise with our multiple masters and PhD programs. Bob brings to this position considerable experience with the masters and PhD programs in the LTI.

I would like to thank Frank Pfenning, who has served ably in this role for the past two years.

Randal E. Bryant
Dean and University Professor

- **X_i is a random variable describing ...**

Answer 1 : X_i is boolean, 1 if word i is in document, else 0

- e.g., $X_{\text{pleased}} = 1$

Issues ?

Learning to classify documents : $P(Y|X)$

- Y discrete valued.
 - e.g., Spam or not
- $X = \langle X_1, X_2, \dots, X_n \rangle = \text{document}$

I am pleased to announce that Bob Frederking of the Language Technologies Institute is our new Associate Dean for Graduate Programs. In this role, he oversees the many issues that arise with our multiple masters and PhD programs. Bob brings to this position considerable experience with the masters and PhD programs in the LTI.

I would like to thank Frank Pfenning, who has served ably in this role for the past two years.

Randal E. Bryant
Dean and University Professor

- X_i is a random variable describing ...

Answer 2 :

- X_i represents the i^{th} word position in document
- $X_1 = \text{"I"}, X_2 = \text{"am"}, X_3 = \text{"pleased"}$
- and, let's assume the X_i are iid (indep, identically distributed)

$$P(X_i | Y) = P(X_j | Y) (\forall_i, j)$$

Learning to classify documents : $P(Y|X)$

the “Bag of Words” model

- **Y discrete valued.** e.g., Spam or not
- **$X = \langle X_1, X_2, \dots, X_n \rangle = \text{document}$**
- **X_i is a random variables. Each represents the word at its position i in the document**
- Generating a document according to this distribution = rolling a 50,000 sided die, once for each word position in the document
- The observed counts for each word follow a ??? distribution

Multinomial Distribution

- $P(\theta)$ and $P(\theta|D)$ have the same form

Eg. 2 Dice roll problem (6 outcomes instead of 2)

Likelihood is \sim Multinomial ($\theta = \{\theta_1, \theta_2, \dots, \theta_k\}$)

$$P(D | \theta) = \theta_1^{\alpha_1} \theta_2^{\alpha_2} \dots \theta_k^{\alpha_k}$$

If prior is Dirichlet distribution,

$$P(\theta) = \frac{\prod_{i=1}^k \theta_i^{\beta_i - 1}}{B(\beta_1, \dots, \beta_k)} \sim \text{Dirichlet}(\beta_1, \dots, \beta_k)$$

Then posterior is Dirichlet distribution

$$P(\theta | D) \sim \text{Dirichlet}(\beta_1 + \alpha_1, \dots, \beta_k + \alpha_k)$$

For Multinomial, conjugate prior is Dirichlet distribution.



Multinomial Bag of Words

the world of **TOTAL**



all about the company

Our energy exploration, production, and distribution operations span the globe, with activities in more than 100 countries.

At TOTAL, we draw our greatest strength from our fast-growing oil and gas reserves. Our strategic emphasis on natural gas provides a strong position in a rapidly expanding market.

Our expanding refining and marketing operations in Asia and the Mediterranean Rim complement already solid positions in Europe, Africa, and the U.S.

Our growing specialty chemicals sector adds balance and profit to the core energy business.

All About The Company

- Global Activities
- Corporate Structure
- TOTAL's Story
- Upstream Strategy
- Downstream Strategy
- Chemicals Strategy
- TOTAL Foundation
- Homepage



aardvark	0
about	2
all	2
Africa	1
apple	0
anxious	0
...	
gas	1
...	
oil	1
...	
Zaire	0

Naïve Bayes Algorithm – discrete X_i

- Train Naïve Bayes (examples)

for each* value y_k

estimate $\pi_k \equiv P(Y = y_k)$

for each* value x_{ij} of each attribute X_i

estimate $\theta_{ijk} \equiv P(X_i = x_{ij} | Y = y_k)$

prob that word x_{ij} appears
in position i, given $Y = y_k$

- Classify (X^{new})

$$Y^{new} \leftarrow \arg \max_{y_k} P(Y = y_k) \prod_i P(X_i^{new} | Y = y_k)$$

$$Y^{new} \leftarrow \arg \max_{y_k} \pi_k \prod_i \theta_{ijk}$$

*Additional assumption : word probability are position independent

$$\theta_{ijk} = \theta_{mjk} \quad \text{for } i \neq m$$

MAP estimates for bag of words

- MAP estimate for multinomial

$$\theta_i = \frac{\alpha_i + \beta_i - 1}{\sum_{m=1}^k \alpha_m + \sum_{m=1}^k \beta_m - 1}$$

$$P(X_i = \text{aardvark}) = \frac{\# \text{observed 'aardvark'} + \# \text{hallucinated 'aardvark} - 1}{\# \text{observed words} + \# \text{hallucinated words} - k}$$

- What β 's should we choose ?

Twenty Newsgroups

Given 1000 training documents from each group
Learn to classify new documents according to
which newsgroup it came from

comp.graphics	misc.forsale
comp.os.ms-windows.misc	rec.autos
comp.sys.ibm.pc.hardware	rec.motorcycles
comp.sys.mac.hardware	rec.sport.baseball
comp.windows.x	rec.sport.hockey
alt.atheism	sci.space
soc.religion.christian	sci.crypt
talk.religion.misc	sci.electronics
talk.politics.mideast	sci.med
talk.politics.misc	
talk.politics.guns	

Naive Bayes: 89% classification accuracy

LEARN_NAIVE_BAYES_TEXT(*Examples*, V)

1. collect all words and other tokens that occur in *Examples*
- $Vocabulary \leftarrow$ all distinct words and other tokens in *Examples*
2. calculate the required $P(v_j)$ and $P(w_k|v_j)$ probability terms
- For each target value v_j in V do
 - $docs_j \leftarrow$ subset of *Examples* for which the target value is v_j
 - $P(v_j) \leftarrow \frac{|docs_j|}{|Examples|}$
 - $Text_j \leftarrow$ a single document created by concatenating all members of $docs_j$
 - $n \leftarrow$ total number of words in $Text_j$ (counting duplicate words multiple times)
 - for each word w_k in $Vocabulary$
 - * $n_k \leftarrow$ number of times word w_k occurs in $Text_j$
 - * $P(w_k|v_j) \leftarrow \frac{n_k+1}{n+|Vocabulary|}$

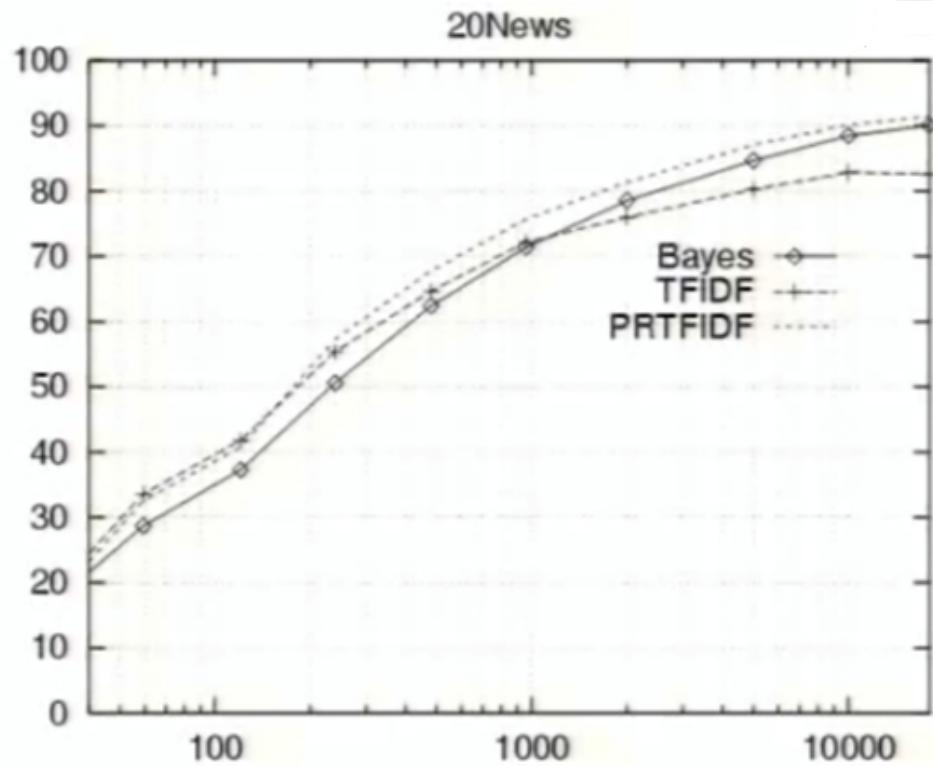
`CLASSIFY_NAIVE_BAYES_TEXT(Doc)`

- $positions \leftarrow$ all word positions in *Doc* that contain tokens found in *Vocabulary*
- Return v_{NB} , where

$$v_{NB} = \operatorname{argmax}_{v_j \in V} P(v_j) \prod_{i \in positions} P(a_i | v_j)$$

Learning Curve for 20 Newsgroups

For code and data, see
www.cs.cmu.edu/~tom/mlbook.html
click on "Software and Data"



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

What if we have continuous X_i ?

- Eg., image classification : X_i is real-valued i^{th} pixel
- Naïve Bayes requires $P(X_i \mid Y=y_k)$, but X_i is real (continuous)

$$P(Y = y_k \mid X_1 \dots X_n) = \frac{P(Y = y_k) \prod_i P(X_i \mid Y = y_k)}{\sum_j P(Y = y_j) \prod_i P(X_i \mid Y = y_j)}$$

- Common approach : assume $P(X_i \mid Y=y_k)$ follows a Normal (Gaussian) distribution

What if we have continuous X_i ?

- Eg., image classification : X_i is ith pixel
- Gaussian Naïve Bayes : assume

$$P(X_i = x \mid Y = y_k) = \frac{1}{\sigma_{ik} \sqrt{2\pi}} e^{\frac{-(x - \mu_{ik})^2}{2\sigma_{ik}^2}}$$

- Sometimes assume σ_{ik}
 - is independent of Y (i.e., σ_i)
 - or independent of X_i (i.e., σ_k)
 - Or both (i.e., σ)

Gaussian Naïve Bayes Algorithm – continuous X_i (but still discrete Y)

- **Train Naïve Bayes (examples)**

for each* value y_k

estimate $\pi_k \equiv P(Y = y_k)$

for each* value x_{ij} of each attribute X_i

estimate $\theta_{ijk} \equiv P(X_i = x_{ij} | Y = y_k)$

- **Classify (X^{new})**

$$Y^{new} \leftarrow \arg \max_{y_k} P(Y = y_k) \prod_i P(X_i^{new} | Y = y_k)$$

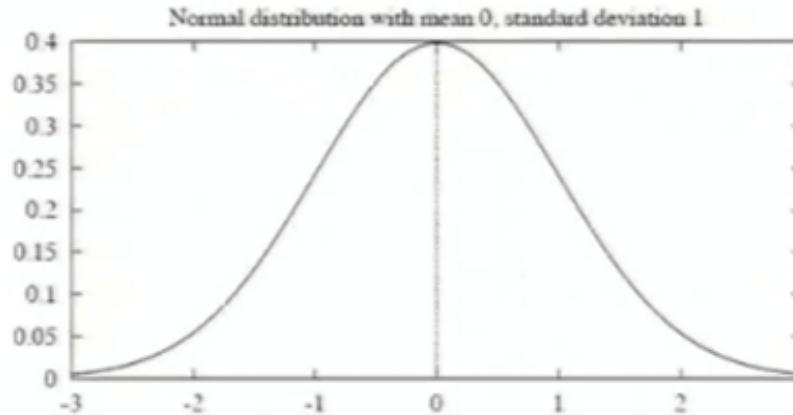
$$Y^{new} \leftarrow \arg \max_{y_k} \pi_k \prod_i Normal(X_i^{new}, \mu_{ik}, \sigma_{ik})$$

*probabilities must sum to 1, so need estimate only n-1 parameters ...

Gaussian Distribution

(also called “Normal”)

$p(x)$ is a *probability density function*, whose integral (not sum) is 1



$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$$

The probability that X will fall into the interval (a, b) is given by

$$\int_a^b p(x)dx$$

- Expected, or menu value of X , $E[X]$, is

$$E[X] = \mu$$

- Variance of X is

$$Var(X) = \sigma^2$$

- Standard deviation of X , σ_X , is

$$\sigma_X = \sigma$$

Gaussian Naïve Bayes – Big Picture

Consider boolean Y , continuous X_i , Assume $P(Y=1)=0.5$

$$Y^{new} \leftarrow \arg \max_{y \in \{0,1\}} P(Y = y) \prod_i P(X_i^{new} | Y = y)$$

What is the minimum possible error ?

Best case :

- conditional independence assumption is satisfied
- we know $P(Y)$, $P(X|Y)$ perfectly (e.g., infinite training data)