

MLE's, Bayesian Classifiers and Naïve Bayes

Required reading:

- Mitchell draft chapter, sections 1 and 2.
(available on class website)

Machine Learning 10-601

Tom M. Mitchell
Machine Learning Department
Carnegie Mellon University

January 30, 2008

Naïve Bayes in a Nutshell

Bayes rule:

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

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})

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

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

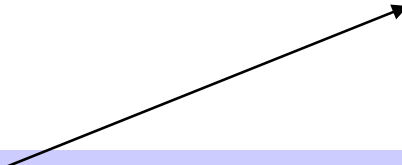
* probabilities must sum to 1, so need estimate only n-1 parameters...

Estimating Parameters: Y, X_i discrete-valued

Maximum likelihood estimates (MLE's):

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

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



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

Example: Live in Sq Hill? $P(S|G,D,M) = \frac{Q_1}{Q_1+Q_2}$

- $S=1$ iff live in Squirrel Hill
- $G=1$ iff shop at Giant Eagle
- $D=1$ iff Drive to CMU
- $M=1$ iff Dave Matthews fan

$P(S=1) = \frac{8}{48}$ $P(S=0) = 1 - P(S=1) = \frac{40}{48}$

$Q_1 = 0.0273$

$P(G=1|S=1) = \frac{6}{8}$ $P(G=0|S=1) = 1 - \frac{6}{8}$

$P(G=1|S=0) = \frac{11}{40}$:

$P(D=1|S=1) = \frac{2}{8}$:

$P(D=1|S=0) = \frac{3}{40}$:

$P(M=1|S=1) = \frac{1}{8}$ $1 - \frac{1}{8} = \frac{7}{8}$

$P(M=1|S=0) = \frac{6}{40}$:

$P(S=1|G=1, D=1, M=0) = \frac{2}{8}$
 $\propto P(S=1) P(G=1|S=1) P(D=1|S=1) P(M=0|S=1)$
 $\frac{8}{48} \cdot \frac{6}{8} \cdot \frac{7}{8}$

test	G	D	M	$P(S=1 test)$	S_{live}
1	1	1	0	$\frac{2}{3}$	0
2	0	0	0	0.02	0
3	1	1	1	0.60	1

$Q_2 = 0.0146 = P(S=0) \cdot P(G=1|S=0) \cdot P(D=1|S=0) \cdot P(M=0|S=0)$
 $\frac{40}{48} \cdot \frac{11}{40} \cdot \frac{3}{40} \cdot \frac{34}{40}$

Example: Live in Sq Hill? $P(S|G,D,M)$

- $S=1$ iff live in Squirrel Hill
- $G=1$ iff shop at Giant Eagle
- $D=1$ iff Drive to CMU
- $M=1$ iff Dave Matthews fan

Naïve Bayes: Subtlety #1

If unlucky, our MLE estimate for $P(X_i / Y)$ may be zero. (e.g., $X_{373} = \text{Birthday_Is_January30}$)

$$\hat{P}(x=1 | s=1) = 0 = \hat{P}(x=0 | s=0)$$

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

Estimating Parameters: Y, X_i discrete-valued

maximize $P(\text{data} | \theta)$

Maximum likelihood estimates:

of examples
in dataset D
that satisfy $Y = y_k$

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

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

$$\max_{\theta} P(\theta | \text{data}) = P(\text{data} | \theta) P(\theta)$$

MAP estimates (Dirichlet priors):

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

Only difference:
"imaginary" examples

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

Naïve Bayes: Subtlety #2

Often the X_i are not really conditionally independent

- We use Naïve Bayes in many cases anyway, and it often works pretty well
 - often the right classification, even when not the right probability (see [Domingos&Pazzani, 1996])
- What is effect on estimated $P(Y|X)$?
 - Special case: what if we add two copies: $X_i = X_k$

$$Q_1 = P(Y)P(X_1|Y)P(X_2|Y)P(X_3|Y)P(X_4|Y)$$

$$Q_1 + Q_2$$

Learning to classify text documents

- Classify which emails are spam
- Classify which emails are meeting invites
- Classify which web pages are student home pages

How shall we represent text documents for Naïve Bayes?

Article from rec.sport.hockey

Path: cantaloupe.srv.cs.cmu.edu!das-news.harvard.e
From: xxx@yyy.zzz.edu (John Doe)
Subject: Re: This year's biggest and worst (opinic
Date: 5 Apr 93 09:53:39 GMT

I can only comment on the Kings, but the most obvious candidate for pleasant surprise is Alex Zhitnik. He came highly touted as a defensive defenseman, but he's clearly much more than that. Great skater and hard shot (though wish he were more accurate). In fact, he pretty much allowed the Kings to trade away that huge defensive liability Paul Coffey. Kelly Hruddy is only the biggest disappointment if you thought he was any good to begin with. But, at best, he's only a mediocre goaltender. A better choice would be Tomas Sandstrom, though not through any fault of his own, but because some thugs in Toronto decided

Learning to Classify Text

Target concept *Interesting?* : *Document* $\rightarrow \{+, -\}$

1. Represent each document by vector of words
 - one attribute per word position in document
2. Learning: Use training examples to estimate
 - $P(+)$
 - $P(-)$
 - $P(doc|+)$
 - $P(doc|-)$

Naive Bayes conditional independence assumption

$$P(doc|v_j) = \prod_{i=1}^{length(doc)} P(a_i = w_k|v_j)$$

where $P(a_i = w_k|v_j)$ is probability that word in position i is w_k , given v_j

one more assumption:

$$P(a_i = w_k|v_j) = P(a_m = w_k|v_j), \forall i, m$$

Baseline: Bag of Words Approach



the world of **TOTAL**

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

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 $Text_j$

- * $P(w_k|v_j) \leftarrow \frac{n_k+1}{n+|Vocabulary|}$

For code and data, see

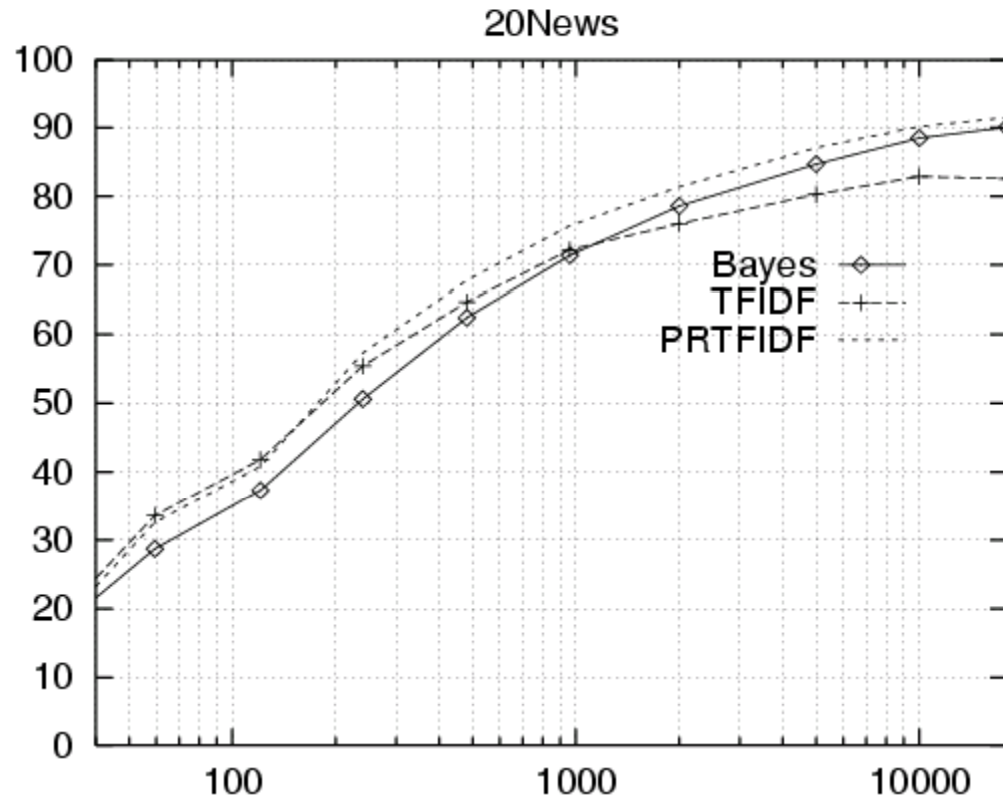
www.cs.cmu.edu/~tom/mlbook.html
click on "Software and Data"

CLASSIFY_NAIVE_BAYES_TEXT(*Doc*)

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

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

Learning Curve for 20 Newsgroups



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

What you should know:

- Training and using classifiers based on Bayes rule
- Conditional independence
 - What it is
 - Why it's important
- Naïve Bayes
 - What it is
 - Why we use it so much
 - Training using MLE, MAP estimates
 - Discrete variables (Bernoulli) and continuous (Gaussian)

$$P(x_i|Y)$$

Questions:

- Can you use Naïve Bayes for a combination of discrete and real-valued X_i ?
- How can we easily model just 2 of n attributes as dependent?
- What does the decision surface of a Naïve Bayes classifier look like?