## Machine Learning 10-601

Tom M. Mitchell
Machine Learning Department
Carnegie Mellon University

September 13, 2012

#### Today:

- Bayes Classifiers
- · Naïve Bayes
- Gaussian Naïve Bayes

#### Readings:

Mitchell:

"Naïve Bayes and Logistic Regression" (available on class website)

#### **Estimating Parameters**

• Maximum Likelihood Estimate (MLE): choose  $\theta$  that maximizes probability of observed data  $\mathcal D$ 

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

• Maximum a Posteriori (MAP) estimate: choose  $\theta$  that is most probable given prior probability and the data

$$\begin{split} \widehat{\theta} &= \arg\max_{\theta} \ P(\theta \mid \mathcal{D}) \\ &= \arg\max_{\theta} \ = \ \frac{P(\mathcal{D} \mid \theta)P(\theta)}{P(\mathcal{D})} \end{split}$$

## Conjugate priors

XE {0,1}

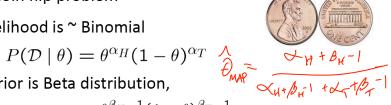
•  $P(\theta)$  and  $P(\theta \mid D)$  have the same form

Eg. 1 Coin flip problem

Likelihood is ~ Binomial

$$P(\mathcal{D} \mid \theta) = \theta^{\alpha_H} (1 - \theta)^{\alpha_T}$$

If prior is Beta distribution,



$$P(\theta) = \frac{\theta^{\beta_H - 1} (1 - \theta)^{\beta_T - 1}}{B(\beta_H, \beta_T)} \sim Beta(\beta_H, \beta_T)$$

Then posterior is Beta distribution

$$P(\theta|D) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

For Binomial, conjugate prior is Beta distribution.

[A. Singh]

## Conjugate priors



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

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

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

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

If prior is Dirichlet distribution,

$$P(\theta) = \frac{\theta_1^{\beta_1 - 1} \theta_2^{\beta_2 - 1} \dots \theta_k^{\beta_k - 1}}{B(\beta_1, \beta_2, \dots \beta_K)} \sim \text{Dirichlet}(\beta_1 \dots \beta_k)$$

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.

[A. Singh]

## Conjugate priors

- $P(\theta)$  and  $P(\theta \mid D)$  have the same form
- Eg. 2 Dice roll problem (6 outcomes instead of 2

Likelihood is ~ Multinomial( $\theta = \{\theta_1, \theta_2, \dots, \theta_n\}$ 

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

If prior is Dirichlet distribution,

$$P(\theta) = \frac{\theta_1^{\beta_1 - 1} \theta_2^{\beta_2 - 1} \dots \theta_k^{\beta_k - 1}}{B(\beta_1, \beta_2, \dots \beta_K)} \sim \text{Dirichlet}$$

Then posterior is Dirichlet distribution

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



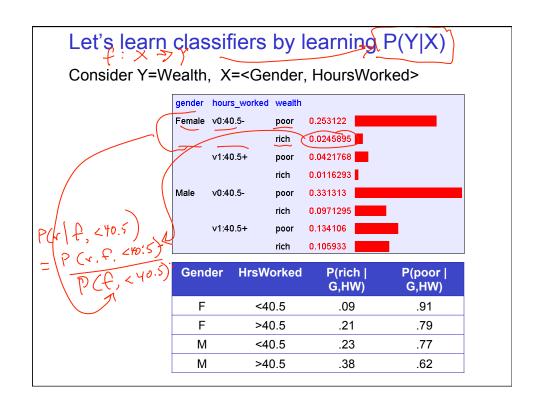
a mater University of Bon
toral advisor Simeon Poisson
Joseph Fourier
toral students Ferdinand Eisens
Leopold Krones
Rudolf Lipschitz

Ferdinand Eisenste
Leopold Kronecker
Rudolf Lipschitz
Carl Wilhelm Borch
Dirichlet function

Dirichlet function
Dirichlet eta function

For Multinomial, conjugate prior is Dirichlet distribution.

[A. Singh]



# How many parameters must we estimate?

Suppose  $X = < X_1,... X_n >$ where  $X_i$  and Y are boolean RV's

Gender	HrsWorked	P(rich   S.HW		W) I	P(poor   G,HW)
F	<40.5		.09	,	.91
F	>40.5	-	.21		.79 •
М	<40.5		.23		.77 •
М	>40.5	4	.38	1	.62-5 =
X,	×2		4		

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

If we have 30 boolean  $X_i$ 's:  $P(Y | X_1, X_2, ... X_{30})$ 



# **Bayes Rule**

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

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 params using Bayes Rule?

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

Suppose X =<X<sub>1</sub>,... X<sub>n</sub>> where X<sub>i</sub> and Y are boolean RV's P(Y|X) = P(X|Y)P(Y) P(X)  $P(X|Y) \Rightarrow P(X|Y) \Rightarrow P(X|Y)$   $P(X|Y) \Rightarrow P(X|Y)$  P

## Naïve Bayes

Naïve Bayes assumes

$$P(X_1 ... X_n | Y) = \prod_i P(X_i | Y)$$

i.e., that X<sub>i</sub> and X<sub>j</sub> are conditionally independent given Y, for all i≠j

### 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<sub>i</sub> are conditionally independent, given Y

X, LX2/Y

Given this assumption, then:

Given this assumption, then:
$$P(X_1, X_2|Y) = P(X_1|X_2, Y)P(X_2|Y)$$

$$= P(X_1|Y)P(X_2|Y)$$

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

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

- Without conditional indep assumption? 2(2\*-1)
- With conditional indep assumption? 24 Params + 1

## Naïve Bayes in a Nutshell

Bayes rule:

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

Assuming conditional independence among 
$$X_i$$
's: 
$$P(Y = y_k | X_1 ... X_n) = \frac{P(Y = y_k) \prod_i P(X_i | Y = y_k)}{\sum_j P(Y = y_j) \prod_i P(X_i | Y = y_j)}$$

So, to pick most probable Y for  $X^{new} = \langle X_1, ..., X_n \rangle$ 

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

### Naïve Bayes Algorithm – discrete X<sub>i</sub>

- 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<sup>new</sup>) 
  $$\begin{split} Y^{new} &\leftarrow \arg\max_{y_k} \ P(Y = y_k) \prod_i P(X_i^{new} | Y = y_k) \\ Y^{new} &\leftarrow \arg\max_{y_k} \ \pi_k \prod_i \theta_{ijk} \end{split}$$

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):

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

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

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

#### Example: Live in Sq Hill? P(S|G,D,E)

- S=1 iff live in Squirrel Hill
- D=1 iff Drive to CMU
- G=1 iff shop at SH Giant Eagle
  - E=1 iff even # of letters in last name

What probability parameters must we estimate?

```
Example: Live in Sq Hill? P(S|G,D,E)

• S=1 iff live in Squirrel Hill

• G=1 iff shop at SH Giant Eagle

• E=1 iff Even # letters last name

P(S=1): 26

P(D=1|S=1): 26

P(D=1|S=1): 26

P(D=0|S=1): 26
```

#### Naïve Bayes: Subtlety #1

If unlucky, our MLE estimate for  $P(X_i | Y)$  might be zero. (e.g.,  $X_i$  = 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  $\theta$  that maximizes probability of observed data  $\mathcal{D}$ 

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

 Maximum a Posteriori (MAP) estimate: choose θ that is most probable given prior probability and the data

#### Estimating Parameters: Y, X<sub>i</sub> discrete-valued

Maximum likelihood estimates:

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

$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_j | Y = y_k) = \frac{\#D\{X_i = x_j \land Y = y_k\}}{\#D\{Y = y_k\}}$$

MAP estimates (Beta, Dirichlet priors):

$$\hat{\pi}_k = \hat{P}(Y = y_k) = \frac{\#D\{Y = y_k\} + (\beta_k - 1)}{|D| + \sum_m (\beta_m - 1)}$$
 "imaginary" examples 
$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_j | Y = y_k) = \frac{\#D\{X_i = x_j \land Y = y_k\} + (\beta_k - 1)}{\#D\{Y = y_k\} + \sum_m (\beta_m - 1)}$$

Only difference:

## 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:  $X_i = X_k$