## Naïve Bayes

Dr. Xiaowei Huang https://cgi.csc.liv.ac.uk/~xiaowei/

#### Up to now,

- Three machine learning algorithms:
  - decision tree learning
  - k-nn
  - linear regression + gradient descent
    - linear regression
    - linear classification
    - logistic regression
    - gradient descent
    - gradient descent on Linear Regression
    - Linear Regression: Analytical Solution

#### Topics

- MLE (maximum Likelihood Estimation) and MAP
- Naïve Bayes

# Recall: MAP Queries (Most Probable Explanation)

- Finding a high probability assignment to some subset of variables
- Most likely assignment to all non-evidence variables W

$$MAP(W | e) = \arg\max_{w} P(w, e) \qquad P(w, e) = P(w|e) P(e)$$

i.e., value of w for which P(w,e) is maximum

#### **Estimating Parameters**

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

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

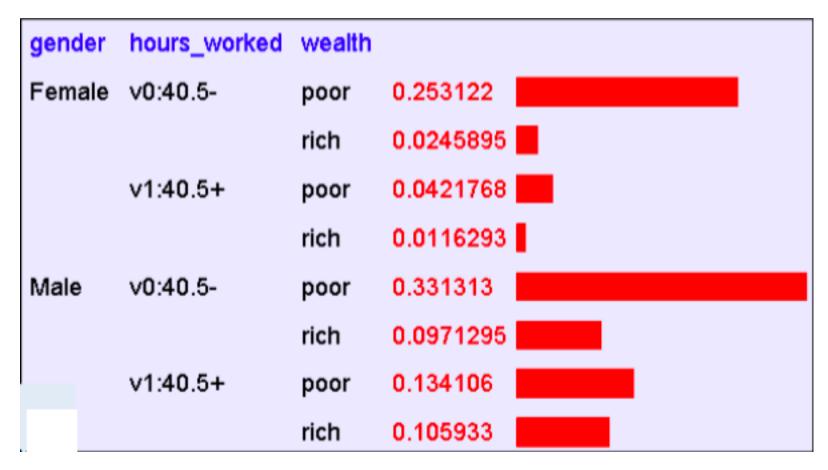
 Maximum a Posteriori (MAP) estimate: choose θ that is most probable given prior probability and the data

$$\hat{ heta} = rg\max_{ heta} P( heta|D)$$
 — posterior  
 $= rg\max_{ heta} rac{P(D| heta)P( heta)}{P(D)} = rg\max_{ heta} P(D| heta)P( heta)$ 

Reducing the number of parameters to estimate

### Let's learn classifiers by learning P(Y|X)

• Consider Y=Wealth, X=<Gender, HoursWorked>



## Let's learn classifiers by learning P(Y|X)

• P(gender, hoursWorked, wealth) => P(wealth|gender, hoursWorked)

Gender	HrsWorked	P(rich   G,HW)	P(poor   G,HW)
F	<40.5	.09	.91
F	>40.5	.21	.79
М	<40.5	.23	.77
М	>40.5	.38	.62

#### How many parameters must we estimate?

feature vector

- Suppose  $X = \langle X_1, ..., X_n \rangle$  where  $X_i$  and Y are Boolean real variables
- To estimate P(Y|X<sub>1</sub>, X<sub>2</sub>, ... X<sub>n</sub>)
  2<sup>n</sup> quantities need to be estimated or collected!

Gender	HrsWorked	P(rich   G,HW)	P(poor   G,HW)
F	<40.5	.09	.91
F	>40.5	.21	.79
М	<40.5	.23	.77
М	>40.5	.38	.62

- If we have 30 Boolean X<sub>i</sub>'s: P(Y | X<sub>1</sub>, X<sub>2</sub>, ... X<sub>30</sub>)
  2<sup>30</sup> ~ 1 billion!
- You need lots of data or a very small *n*

### Can we reduce parameters using Bayes Rule?

- Suppose  $X = \langle X_1, ..., X_n \rangle$  where  $X_i$  and Y are Boolean real variables
- By Bayes rule:

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

• How many parameters for  $P(X|Y) = P(X_1, ..., X_n|Y)$ ? (2<sup>n</sup>-1)x2

How many parameters for P(	Y)?
1	

For example, P(Gender,HrsWorked|Wealth)

Gender	HrsWorked	P(rich   G,HW)	P(poor   G,HW)
F	<40.5	.09	.91
F	>40.5	.21	.79
М	<40.5	.23	.77
М	>40.5	.38	.62

For example, P(Wealth)

#### Naïve Bayes

• Naïve Bayes assumes

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

i.e., that  $X_i$  and  $X_i$  are conditionally independent given Y, for all  $i \neq j$ 

For example, P(Gender,HrsWorked|Wealth) = P(Gender|Wealth) \* P(HrsWorked|Wealth)

#### Recap: Conditional independence

• Two variables A, B are *independent* if

 $P(A \land B) = P(A)P(B)$  $\forall a, b : P(A = a \land B = b) = P(A = a)P(B = b)$ 

• Two variables A, B are *conditionally independent given* C if

 $P(A \land B|C) = P(A|C)P(B|C)$  $\forall a, b, c : P(A = a \land B = b|C = c) = P(A = a|C = c)P(B = b|C = c)$ 

#### Recap: Conditional Independence

• A is conditionally independent of B given C, if the probability distribution governing A is independent of the value of B, given the value of C

$$\forall a, b, c : P(A = a | B = b, C = c) = P(A = a | C = c)$$

- Which we often write P(A|B,C) = P(A|C)
- Example: P(Thunder|Rain, Lightning) = P(Thunder|Lightning)

#### Assumption for Naïve Bayes

- Naïve Bayes uses assumption that the X<sub>i</sub> are conditionally independent, given Y
- Given this assumption, then:

Chain rule

$$P(X_1, X_2|Y) = P(X_1|X_2, Y)P(X_2|Y)$$
  
=  $P(X_1|Y)P(X_2|Y)$  Conditional  
Independence

• in general:  $P(X_1...X_n|Y) = \prod_i P(X_i|Y)$ (2<sup>n</sup>-1)x2 2n Why? Every P(X\_i|Y) takes a parameter, and we have n X\_i.

## Reducing the number of parameters to estimate

$$P(Y|X_1, ..., X_n) = \frac{P(X_1, ..., X_n | Y) P(Y)}{P(X_1, ..., X_n)}$$

• To make this tractable we naively assume conditional independence of the features given the class: ie

$$P(X_1, ..., X_n | Y) = P(X_1 | Y) P(X_2 | Y) ... P(X_n | Y)$$

• Now: I only need to estimate ... parameters:

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

# Reducing the number of parameters to estimate

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

- Without conditional independent assumption?
  - (2<sup>n</sup>-1)x2+1
- With conditional independent assumption?
  - 2n+1

#### Naïve Bayes Algorithm

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

- Train Naïve Bayes (given data for X and Y)
- for each value  $y_k$ 
  - Estimate  $\pi_k \equiv P(Y=y_k)$
- for each value  $x_{ij}$  of each attribute  $X_i$ 
  - estimate  $\theta_{ijk} = P(X_i = x_{ij}|Y = y_k)$

#### Training Naïve Bayes Classifier

- From the data D, estimate *class priors:* 
  - For each possible value of Y, estimate  $Pr(Y=y_1)$ ,  $Pr(Y=y_2)$ ,....  $Pr(Y=y_k)$
  - An estimate:

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

- From the data, estimate the conditional probabilities
  - If every X<sub>i</sub> has values  $x_{i1},...,x_{ik}$ 
    - for each  $y_i$  and each  $X_i$  estimate  $q(i,j,k)=Pr(X_i=x_{ij}|Y=y_k)$

• 
$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_{ij}|Y = y_k) = \frac{\#D\{X_i = x_{ij} \land Y = y_k\}}{\#D\{Y = y_k\}}$$
 Number of items in dataset D for which  $Y=y_k$ 

#### Exercise

- Consider the following dataset:
- P(Wealthy=Y) =
- P(Wealthy=N)=
- P(Gender=F | Wealthy = Y) =
- P(Gender=M | Wealthy = Y) =
- P(HrsWorked > 40.5 | Wealthy = Y) =
- P(HrsWorked < 40.5 | Wealthy = Y) =
- P(Gender=F | Wealthy = N) =
- P(Gender=M | Wealthy = N) =
- P(HrsWorked > 40.5 | Wealthy = N) =
- P(HrsWorked < 40.5 | Wealthy = N) =

Gender	HrsWorked	Wealthy?
F	39	Υ
F	45	Ν
Μ	35	N
Μ	43	N
F	32	Y
F	47	Υ
Μ	34	Υ

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

- Train Naïve Bayes (given data for X and Y)
- for each value  $y_k$ 
  - Estimate  $\pi_k \equiv P(Y=y_k)$
- for each value  $x_{ij}$  of each attribute  $X_i$ 
  - estimate  $\theta_{ijk} = P(X_i = x_{ij}|Y = y_k)$
- Classify (X<sub>new</sub>)

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

#### Exercise (Continued)

- Consider the following dataset:
- Classify a new instance
  - Gender = F / HrsWorked = 44

Gender	HrsWorked	Wealthy?
F	39	Υ
F	45	Ν
М	35	N
М	43	Ν
F	32	Υ
F	47	Υ
М	34	Υ

### Example: Live outside of Liverpool? P(L|T,D,E)

- L=1 iff live outside of Liverpool D=1 iff Drive or Carpool to Liverpool
- T=1 iff shop at Tesco E=1 iff Even # letters last name

P(L=1):	P(L=0):
P(D=1   L=1) :	P(D=0   L=1) :
P(D=1   L=0) :	P(D=0   L=0) :
P(T=1   L=1):	P(T=0   L=1):
P(T=1   L=0) :	P(T=0   L=0) :
P(E=1   L=1):	P(E=0   L=1):
P(E=1   L=0):	P(E=0   L=0):