### CSE4334/5334 Data Mining Classification: Bayesian Classifiers

Chengkai Li

Department of Computer Science and Engineering University of Texas at Arlington Fall 2020 (Slides courtesy of Pang-Ning Tan, Michael Steinbach and Vipin Kumar)



A probabilistic framework for solving classification problems

Conditional Probability:

$$P(Y \mid X) = \frac{P(X, Y)}{P(X)}$$
$$P(X \mid Y) = \frac{P(X, Y)}{P(Y)}$$

Bayes theorem:

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



# Example of Bayes Theorem

### Given:

A doctor knows that meningitis causes stiff neck 50% of the time
Prior probability of any patient having meningitis is 1/50,000
Prior probability of any patient having stiff neck is 1/20
If a patient has stiff neck, what's the probability he/she has meningitis?

$$P(M \mid S) = \frac{P(S \mid M)P(M)}{P(S)} = \frac{0.5 \times 1/50000}{1/20} = 0.0002$$

## Using Bayes Theorem for Classification

Consider each attribute and class label as random variables

Given a record with attributes  $(X_1, X_2, ..., X_d)$ 

• Goal is to predict class Y

• Specifically, we want to find the value of Y that maximizes  $P(Y | X_1, X_2, ..., X_d)$ 

Can we estimate  $P(Y | X_1, X_2, ..., X_d)$  directly from data?

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

## Using Bayes Theorem for Classification

#### Approach:

 $\circ\,$  compute posterior probability P(Y  $\mid\, X_1, X_2, \, ..., \, X_d)$  using the Bayes theorem

$$P(Y \mid X_1 X_2 \dots X_d) = \frac{P(X_1 X_2 \dots X_d \mid Y) P(Y)}{P(X_1 X_2 \dots X_d)}$$

- $\circ$  Maximum a-posteriori: Choose Y that maximizes P(Y | X<sub>1</sub>, X<sub>2</sub>, ..., X<sub>d</sub>)
- Equivalent to choosing value of Y that maximizes P(X<sub>1</sub>, X<sub>2</sub>, ..., X<sub>d</sub> | Y) P(Y)

How to estimate  $P(X_1, X_2, ..., X_d | Y)$ ?

## Example Data

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

#### Given a Test Record:

X = (Refund = No, Divorced, Income = 120K)

• Can we estimate

P(Evade = Yes | X) and P(Evade = No | X)?

In the following we will replace Evade = Yes by Yes, and Evade = No by No



## Example Data

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married 100K		No
3	No	Single	Single 70K	
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

#### Given a Test Record:

X = (Refund = No, Divorced, Income = 120K)

Using Bayes Theorem:  $P(Yes | X) = \frac{P(X | Yes)P(Yes)}{P(X)}$   $P(No | X) = \frac{P(X | No)P(No)}{P(X)}$ 

• How to estimate

P(X | Yes) and P(X | No)?







### Conditional Independence

**X** and **Y** are independent if  $P(\mathbf{X} | \mathbf{Y}) = P(\mathbf{X})$  and  $P(\mathbf{Y} | \mathbf{X}) = P(\mathbf{Y})$ **X** and **Y** are conditionally independent given **Z** if  $P(\mathbf{X} | \mathbf{YZ}) = P(\mathbf{X} | \mathbf{Z})$ and  $P(\mathbf{Y} | \mathbf{XZ}) = P(\mathbf{Y} | \mathbf{Z})$ 

#### Example: Arm length and reading skills

- Young child has shorter arm length and limited reading skills, compared to adults
- If age is fixed, no apparent relationship between arm length and reading skills
- Arm length and reading skills are conditionally independent given age

### Naïve Bayes Classifier



Assume independence among attributes  $X_i$  when class is given:  $\circ P(X_1, X_2, ..., X_d | Y) = P(X_1 | Y) P(X_2 | Y) ... P(X_d | Y)$ 

- $\circ~$  Now we can estimate  $P(X_i \mid~Y)$  for all value combinations of  $X_i$  and Y from the training data
- New point is classified to y if  $P(y) \prod P(X_i | y)$  is maximal.

## Putting Everything Together

Problem: Choose value of Y that maximizes  $P(Y | X_1, X_2, ..., X_d)$ 

$$\begin{split} & P(Y \mid X_{1}, X_{2}, ..., X_{d}) \\ &= \frac{P(X_{1}, X_{2}, ..., Xd \mid Y) P(Y)}{P(X_{1}, X_{2}, ..., Xd)} \text{ (Bayes Theorem)} \\ &= \frac{P(X_{1} \mid Y) P(X_{2} \mid Y) ... P(Xd \mid Y) P(Y)}{P(X_{1}, X_{2}, ..., Xd)} \text{ (Under the Attribute Independence Assumption)} \\ &= \frac{P(Y) \prod_{i=1}^{d} P(Xi \mid Y)}{P(X_{1}, X_{2}, ..., Xd)} \end{split}$$

### Naïve Bayes on Example Data

#### Given a Test Record:

X = (Refund = No, Divorced, Income = 120K)

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

P(X | Yes) =

P(Refund = No | Yes) x P(Divorced | Yes) x

P(Income = 120K | Yes)

P(X | No) = P(Refund = No | No) x P(Divorced | No) x P(Income = 120K | No)



### Estimate Probabilities from Data

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

P(y) = fraction of instances of class y- e.g., P(No) = 7/10,P(Yes) = 3/10

### For categorical attributes: $P(X_i = c | y) = n_c / n$

- where |X<sub>i</sub> =c| is number of instances having attribute value X<sub>i</sub> =c and belonging to class y
- Examples:

P(Status=Married | No) = 4/7P(Refund=Yes | Yes)=0

### Estimate Probabilities from Data



For continuous attributes:

Discretize: partition the range into bins:

• Replace continuous value with bin value (Attribute changed from continuous to ordinal)

#### Probability density estimation:

- Assume attribute follows a normal distribution
- Use data to estimate parameters of distribution (e.g., mean and standard deviation)
- $\circ~$  Once probability distribution is known, can use it to estimate the conditional probability  $P(X_i \,|\, Y)$

### How to Estimate Probabilities from Data?

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Normal distribution:  $P(X_i | Y_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} e^{-\frac{(X_i - \mu_{ij})^2}{2\sigma_{ij}^2}}$ 

- One for each (X<sub>i</sub>, Y<sub>j</sub>) pair. (X<sub>i</sub> is an attribute, Y<sub>j</sub> is a class attribute value.
   For (Income, Class=No):
- If Class=No
  - $\circ$  sample mean = 110
  - $\circ$  sample variance = 2975

$$P(Income = 120 \mid No) = \frac{1}{\sqrt{2\pi}(54.54)} e^{-\frac{(120-110)^2}{2(2975)}} = 0.0072$$



### Example of Naïve Bayes Classifier

Given a Test Record:

#### X = (Refund = No, Divorced, Income = 120K)

Naïve Bayes Classifier:

$$\begin{split} P(\text{Refund} &= \text{Yes} \mid \text{No}) = 3/7 \\ P(\text{Refund} = \text{No} \mid \text{No}) = 4/7 \\ P(\text{Refund} = \text{Yes} \mid \text{Yes}) = 0 \\ P(\text{Refund} = \text{No} \mid \text{Yes}) = 1 \\ P(\text{Marital Status} = \text{Single} \mid \text{No}) = 2/7 \\ P(\text{Marital Status} = \text{Divorced} \mid \text{No}) = 1/7 \\ P(\text{Marital Status} = \text{Married} \mid \text{No}) = 4/7 \\ P(\text{Marital Status} = \text{Single} \mid \text{Yes}) = 2/3 \\ P(\text{Marital Status} = \text{Divorced} \mid \text{Yes}) = 1/3 \\ P(\text{Marital Status} = \text{Married} \mid \text{Yes}) = 0 \end{split}$$

For Taxable Income: If class = No: sample mean = 110 sample variance = 2975 If class = Yes: sample mean = 90 sample variance = 25 •  $P(X \mid No) = P(Refund=No \mid No)$ ×  $P(Divorced \mid No)$ ×  $P(Income=120K \mid No)$ =  $4/7 \times 1/7 \times 0.0072 = 0.0006$ 

• P(X | Yes) = P(Refund=No | Yes)× P(Divorced | Yes)× P(Income=120K | Yes)=  $1 \times 1/3 \times 1.2 \times 10^{-9} = 4 \times 10^{-10}$ 

Since P(X | No)P(No) > P(X | Yes)P(Yes)Therefore P(No | X) > P(Yes | X)=> Class = No

# Naïve Bayes Classifier can make decisions with partial information about attributes in the test record

Even in absence of information about any attributes, we can use Apriori Probabilities of Class Variable:

Naïve Bayes Classifier:

$$\begin{split} P(\text{Refund} &= \text{Yes} \mid \text{No}) = 3/7 \\ P(\text{Refund} &= \text{No} \mid \text{No}) = 4/7 \\ P(\text{Refund} &= \text{Yes} \mid \text{Yes}) = 0 \\ P(\text{Refund} &= \text{No} \mid \text{Yes}) = 1 \\ P(\text{Marital Status} = \text{Single} \mid \text{No}) = 2/7 \\ P(\text{Marital Status} = \text{Divorced} \mid \text{No}) = 1/7 \\ P(\text{Marital Status} = \text{Married} \mid \text{No}) = 4/7 \\ P(\text{Marital Status} = \text{Single} \mid \text{Yes}) = 2/3 \\ P(\text{Marital Status} = \text{Divorced} \mid \text{Yes}) = 1/3 \\ P(\text{Marital Status} = \text{Married} \mid \text{Yes}) = 0 \end{split}$$

For Taxable Income: If class = No: sample mean = 110 sample variance = 2975 If class = Yes: sample mean = 90 sample variance = 25 P(Yes) = 3/10P(No) = 7/10

If we only know that marital status is Divorced, then:  $P(Yes \mid Divorced) = 1/3 \ge 3/10 / P(Divorced)$  $P(No \mid Divorced) = 1/7 \ge 7/10 / P(Divorced)$ 

If we also know that Refund = No, then

P(Yes | Refund = No, Divorced) = 1 x 1/3 x 3/10 / P(Divorced, Refund = No)

P(No | Refund = No, Divorced) = 4/7 x 1/7 x 7/10 / P(Divorced, Refund = No)

If we also know that Taxable Income = 120, then

P(Yes | Refund = No, Divorced, Income = 120) = 1.2 x10<sup>-9</sup> x 1 x 1/3 x 3/10 / P(Divorced, Refund = No, Income = 120)

P(No | Refund = No, Divorced Income = 120) = 0.0072 x 4/7 x 1/7 x 7/10 / P(Divorced, Refund = No, Income = 120)

#### Issues with Naïve Bayes Classifier

#### Given a Test Record:

X = (Married)

#### Naïve Bayes Classifier:

$$\begin{split} P(\text{Refund} &= \text{Yes} \mid \text{No}) = 3/7 \\ P(\text{Refund} = \text{No} \mid \text{No}) = 4/7 \\ P(\text{Refund} = \text{Yes} \mid \text{Yes}) = 0 \\ P(\text{Refund} = \text{No} \mid \text{Yes}) = 1 \\ P(\text{Marital Status} = \text{Single} \mid \text{No}) = 2/7 \\ P(\text{Marital Status} = \text{Divorced} \mid \text{No}) = 1/7 \\ P(\text{Marital Status} = \text{Married} \mid \text{No}) = 4/7 \\ P(\text{Marital Status} = \text{Single} \mid \text{Yes}) = 2/3 \\ P(\text{Marital Status} = \text{Divorced} \mid \text{Yes}) = 1/3 \\ P(\text{Marital Status} = \text{Married} \mid \text{Yes}) = 0 \end{split}$$

For Taxable Income: If class = No: sample mean = 110 sample variance = 2975 If class = Yes: sample mean = 90 sample variance = 25 P(Yes) = 3/10P(No) = 7/10

 $P(Yes | Married) = 0 \ge 3/10 / P(Married)$  $P(No | Married) = 4/7 \ge 7/10 / P(Married)$ 



### Issues with Naïve Bayes Classifier



Consider the table with Tid = 7 deleted

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	Single 125K I	
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married 120K		No
5	No	Divorced 95K		Yes
6	No	Married 60K		No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Given X = (Refund = Yes, Divorced, 120K)

P(X | No) = 2/6 X 0 X 0.0083 = 0P(X | Yes) = 0 X 1/3 X 1.2 X 10<sup>-9</sup> = 0 Naïve Bayes Classifier:

P(Refund = Yes | No) = 2/6P(Refund = No | No) = 4/6 $\rightarrow$  P(Refund = Yes | Yes) = 0 P(Refund = No | Yes) = 1P(Marital Status = Single | No) = 2/6 $\rightarrow$  P(Marital Status = Divorced | No) = 0 P(Marital Status = Married | No) = 4/6P(Marital Status = Single | Yes) = 2/3P(Marital Status = Divorced | Yes) = 1/3P(Marital Status = Married | Yes) = 0/3For Taxable Income: If class = No: sample mean = 91sample variance = 685If class = No: sample mean = 90sample variance = 25

> Naïve Bayes will not be able to classify X as Yes or No!

### Issues with Naïve Bayes Classifier



- If one of the conditional probability is zero, then the entire expression becomes zero.
- Need to use other estimates of conditional probabilities than simple fractions.
- Probability estimation:

original: 
$$P(X_i = c | y) = \frac{n_c}{n}$$

Laplace Estimate: 
$$P(X_i = c | y) = \frac{n_c + 1}{n + v}$$

m – estimate: 
$$P(X_i = c|y) = \frac{n_c + mp}{n + m}$$

*n*: number of training instances belonging to class *y* 

*n<sub>i</sub>*: number of instances with  $X_i = c$  and Y = y

*v*: total number of attribute values that  $X_i$  can take

*p*: initial estimate of  $P(X_i = c | y)$  known apriori, e.g., 1/v, or something else

*m*: hyper-parameter for our confidence in p

### Example of Naïve Bayes Classifier



Name	Give Birth	Can Fly	Live in Water	Have Legs	Class
human	yes	no	no	yes	mammals
python	no	no	no	no	non-mammals
salmon	no	no	yes	no	non-mammals
whale	yes	no	yes	no	mammals
frog	no	no	sometimes	yes	non-mammals
komodo	no	no	no	yes	non-mammals
bat	yes	yes	no	yes	mammals
pigeon	no	yes	no	yes	non-mammals
cat	yes	no	no	yes	mammals
leopard shark	yes	no	yes	no	non-mammals
turtle	no	no	sometimes	yes	non-mammals
penguin	no	no	sometimes	yes	non-mammals
porcupine	yes	no	no	yes	mammals
eel	no	no	yes	no	non-mammals
salamander	no	no	sometimes	yes	non-mammals
gila monster	no	no	no	yes	non-mammals
platypus	no	no	no	yes	mammals
owl	no	yes	no	yes	non-mammals
dolphin	yes	no	yes	no	mammals
eagle	no	yes	no	yes	non-mammals

Give Birth	Can Fly	Live in Water	Have Legs	Class
yes	no	yes	no	?

A: attributes M: mammals N: non-mammals

 $P(A \mid M) = \frac{6}{7} \times \frac{6}{7} \times \frac{2}{7} \times \frac{2}{7} = 0.06$  $P(A \mid N) = \frac{1}{13} \times \frac{10}{13} \times \frac{3}{13} \times \frac{4}{13} = 0.0042$  $P(A \mid M)P(M) = 0.06 \times \frac{7}{20} = 0.021$  $P(A \mid N)P(N) = 0.004 \times \frac{13}{20} = 0.0027$ 

P(A | M)P(M) > P(A | N)P(N)=> Mammals

## Naïve Bayes (Summary)



- Robust to isolated noise points
- Handle missing values by ignoring the instance during probability estimate calculations
- Robust to irrelevant attributes
- Redundant and correlated attributes will violate class conditional assumption
  - Use other techniques such as Bayesian Belief Networks (BBN)



### Naïve Bayes

How does Naïve Bayes perform on the following dataset?



Conditional independence of attributes is violated

# Bayesian Belief Networks

- Provides graphical representation of probabilistic relationships among a set of random variables
- Consists of:
  - A directed acyclic graph (dag)
    Node corresponds to a variable
    Arc corresponds to dependence relationship between a pair of variables



- A probability table associating each node to its immediate parent





### Conditional Independence



D is parent of C

A is child of C

B is descendant of D

D is ancestor of A

A node in a Bayesian network is conditionally independent of all of its nondescendants, if its parents are known



## Conditional Independence

Naïve Bayes assumption:





## Probability Tables

 If X does not have any parents, table contains prior probability P(X)

• If X has only one parent (Y), table contains conditional probability P(X | Y)

• If X has multiple parents  $(Y_1, Y_2, ..., Y_k)$ , table contains conditional probability  $P(X | Y_1, Y_2, ..., Y_k)$ 

### Example of Bayesian Belief Network





# Example of Inferencing using BBN

- Given: X = (E=No, D=Yes, CP=Yes, BP=High)
  - Compute P(HD | E,D,CP,BP)?

- P(HD=Yes | E=No,D=Yes,CP=Yes,BP=High)  $\propto 0.55 \times 0.8 \times 0.85 = 0.374$
- P(HD=No | E=No,D=Yes) = 0.45
   P(CP=Yes | HD=No) = 0.01
   P(BP=High | HD=No) = 0.2
  - P(HD=No | E=No,D=Yes,CP=Yes,BP=High)
     ∞ 0.45 × 0.01 × 0.2 = 0.0009

Classify X as Yes