# Machine Learning (IV): Conditional Independence: Naive Bayes Classifiers and Other Examples

Le Wang

## Conditional Distribution and Independence

$$p(x_i, y_j) = p(x_i)p(y_j)$$

## Conditional Distribution and Independence

$$p(x_i, y_j) = p(x_i)p(y_j)$$
$$\frac{p(x_i, y_j)}{p(x_i)} = p(y_j)$$

### Conditional Distribution and Independence

$$p(x_i, y_j) = p(x_i)p(y_j)$$
$$\frac{p(x_i, y_j)}{p(x_i)} = p(y_j)$$

In other words

$$p(y_j \mid x_i) = p(y_j)$$

Independence implies that conditional distribution is marginal distribution!

**Intuitively**, no predictive power at all! as it should be for independent variables!

### Unconditional Independence to Conditional Independence

To identify the causal effects of x on y, we will later rely on variants of the assumption of **conditional independence** (it is not the same as **unconditional independence**).

$$p(x, y|z) = p(x|z)p(y|z)$$
 for all  $x, y, z$   
 $p(y|x, z) = p(y|z)$ 

**Interpretation**: The conditional distribution of Y, given X, Z is in fact completely determined by the value of Z alone, Y being superfluous once Z is given.

### Visualization of Conditional Independence

**Question:** Can you think of a way to visualize the conditional independence for all possible combinations?

### Visualization of Conditional Independence

**Question:** Can you think of a way to visualize the conditional independence for all possible combinations?

A scatter plot of conditional probability and the product of two marginal probabilities with a 45 degree line. Why?

#### Further Applications

As we will see later, conditional independence plays an important role in causal inference, designing different experiments (e.g., conditional randomization), and resolving many of the empirical paradoxes (e.g., Simpson's paradox).

Conditional independence is, sometimes, a direct implication of economic theory. For example, in the literature of insurance, the presence of positive conditional dependence between coverage and risk is known to be a direct consequence of adverse selection under information asymmetry [e.g., Chiappori and Salanie (2000)].

The literature of testing conditional independence for continuous variables appears rather recent and includes relatively few researches [Song, 2009, Testing Conditional Independence via Transforms, Annals of Statistics, Vol. 37, No. 6B, 4011–4045.]

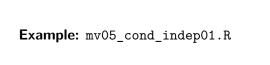
## Application (I): The Problem of Operating a "Fair" Admission Policy

An important issue: Operating a "fair" procedure for the selection of minority group members for universeity admission. One solution is to require that the probability of such selection should depend only on the academic promise of the candidates, and not on race, sex and so on. How to examine whether this criterion is indeed met?

Let Y=1 (selection) Y=0 (rejection), let X denote sex or race, and let Z be a test-score or other measures of academic promise. Our goal is then to test

$$Y \perp X \mid Z$$

Given Z (academic performance), there should not be any predictive power of X (sex or race) for Y (admission or not).



## Application (II): Markov Chain Assumption

Recall that a **Markov Chain** is a stochastic process for whic, given the current state, future states of the random variable Y are independent of past states.

$$\Pr[Y_{t+1} = y \mid y_t, y_{t-1}, \dots] = \Pr[Y_{t+1} = y \mid y_t]$$

### **Applications**

#### Combining conditional independence and Bayes' Rule

- 1. Statistical Language Model
- 2. Naive Bayes Classifier
- 3. Predicting Unobservable Ethnicity in Political Science.

## Application (III): Statistical Language Model

$$\mathsf{Sender} \overset{s_1, s_2, s_3, \dots}{\Longrightarrow} \mathsf{Chanel} \overset{o_1, o_2, o_3, \dots}{\Longrightarrow} \mathsf{Receiver}$$

You send some signals to the receiver. The task of communication is to decode whatever signals the receiver receives. In other words, you would like to back out the original meaning  $(s_1, s_2, s_3, \ldots)$  of the received cell-phone signals  $(o_1, o_2, o_3, \ldots)$ .

Nearly all natural language processing problems can be thought of as decoding the communication process.

This model is called **Acoustic Model** in voice recognition, **Translation Model** in machine learning, **Correction Model** in autocorrection.

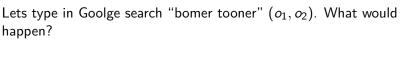
- 1. Speech recognition: You say something, and the computer receives the digits and cracks out what you say.
- 2. Translations between languages

How is this related to what we have learned? Let me translate for you.

Just look for the most likely  $s_1, s_2, s_3 \dots$ !

$$\Pr[s_1, s_2, s_3, \dots | o_1, o_2, o_3, \dots]$$

With Conditional Independence, Bayes' Rule and Hidden Markov Model, we can solve this.



happen?

 $\Pr[s_1, s_2, s_3, \dots | o_1, o_2, o_3, \dots] =$ 

$$\frac{\Pr[s_1, s_2, s_3, \dots | o_1, o_2, o_3, \dots]}{\Pr[o_1, o_2, o_3, \dots | s_1, s_2, s_3, \dots] \cdot \Pr[s_1, s_2, s_3, \dots]}}{\Pr[o_1, o_2, o_3, \dots]}$$

When I recevie the signal,  $o_1, o_2, o_3 \dots$  is already known, we can then ignore it. Because it is a constant, our goal is to maximize the quantity above. Then, the question becomes maximizing the following

$$Pr[o_1, o_2, o_3, \dots | s_1, s_2, s_3, \dots] \cdot Pr[s_1, s_2, s_3, \dots]$$

How can we solve this?

Solution: Hidden Markov Model

Hidden Markov Model:

Leonard E. Baum and others developed Hidden Markov Model in

60s and 70s. The idea is simple:

#### Hidden Markov Model:

Leonard E. Baum and others developed Hidden Markov Model in 60s and 70s. The idea is simple:

We do not observe the original messages (hence **Hidden**), but I assume that it follows a Markov Model.

#### Hidden Markov Model:

Leonard E. Baum and others developed Hidden Markov Model in 60s and 70s. The idea is simple:

We do not observe the original messages (hence **Hidden**), but I assume that it follows a Markov Model.

In the meantime, even though I do not observe the original message, I assume that for every message  $s_t$ , I could observe a signal  $o_t$ . This signal  $o_t$  depends on  $s_t$  and **only**  $s_t$ !

$$Pr[o_1, o_2, o_3, \dots | s_1, s_2, s_3, \dots] \cdot Pr[s_1, s_2, s_3, \dots]$$

1. Conditional Independence:

$$Pr[o_1, o_2, o_3, \dots | s_1, s_2, s_3, \dots] = Pr[o_1 | s_1] \cdot Pr[o_2 | s_2] \cdots$$

2. (Hidden) Markov Model:

$$\Pr[s_1, s_2, s_3, \dots] = \Pr[s_2 | s_1] \cdot \Pr[s_3 | s_2] \cdots$$

We can solve this problem using, e.g., Viterbi Algorithm

In the 70s, James and Janet Baker (co-founders of Dragon) propos
using Hidden Markov model for voice recognition. The

misclassification rates reduced from 30% to 10%.

Kai-fu Lee used this type of model to develop Sphinx.

### **Applications**

#### Combining conditional independence and Bayes' Rule

- 1. Statistical Language Model
- 2. Naive Bayes Classifier
- 3. Predicting Unobservable Ethnicity in Political Science.

## Application (IV): Naive Bayes Classifier

Suppose that we now have K predictors  $(X_1, X_2, ..., X_K)$  and the outcome has J classes  $y_1, y_2, ..., y_J$ .

## Application (IV): Naive Bayes Classifier

Suppose that we now have K predictors  $(X_1, X_2, ..., X_K)$  and the outcome has J classes  $y_1, y_2, ..., y_J$ .

Suppose that the predictors take on the following values

$$X_1 = x_1, X_2 = x_2, \dots, X_K = x_k$$

For prediction, we need to know the following values

$$Pr[Y = y_1 \mid X_1 = x_1, X_2 = x_2, ..., X_K = x_k]$$

$$Pr[Y = y_2 \mid X_1 = x_1, X_2 = x_2, ..., X_K = x_k]$$

$$Pr[Y = y_3 \mid X_1 = x_1, X_2 = x_2, ..., X_K = x_k]$$

#### Question:

What is your prediction for the employment status of a guy named Le Wang, a 38-year-old man born in Guangdong but living in Vermont?

```
\begin{split} &\text{Pr}[\text{Unemployment} = 1 \mid \text{Name} = \text{Le Wang}, \text{Age} = 38, \\ &\text{Gender} = \text{Male}, \text{ Place of Birth} = \text{Guangdong} \;, \\ &\text{Place of Residence} = \text{Vermont} \;] \end{split}
```

#### Question:

What is your prediction for the employment status of a guy named Le Wang, a 38-year-old man born in Guangdong but living in Vermont?

```
\begin{split} &\text{Pr}[\text{Unemployment} = 1 \mid \text{Name} = \text{Le Wang}, \text{Age} = 38, \\ &\text{Gender} = \text{Male}, \text{ Place of Birth} = \text{Guangdong} \;, \\ &\text{Place of Residence} = \text{Vermont} \;] \end{split}
```

$$= \frac{\mathsf{Pr}[\mathsf{Unemployment} = 1, \mathsf{Name} = \mathsf{Le}\;\mathsf{Wang}, \mathsf{Age} = 38, \dots]}{\mathsf{Pr}[\mathsf{Name} = \mathsf{Le}\;\mathsf{Wang}, \mathsf{Age} = 38, \dots]}$$

#### Naive Bayes Classififiers:

- 1. Step 1. We can simplify this using Bayes' Rule
- 2. **Step 2.** We can simplify this using **Conditional Independence Assumption**

#### Step 1. We can simplify this using Bayes' Rule

$$\Pr[Y = y_1 \mid X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]$$

$$= \frac{\Pr[Y = y_1, X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}$$

#### Step 1. We can simplify this using Bayes' Rule

$$Pr[Y = y_1 \mid X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]$$

$$= \frac{\Pr[Y = y_1, X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}$$

$$= \frac{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k \mid Y = y_1] \cdot \Pr[Y = y_1]}{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}$$

#### Step 1. We can simplify this using Bayes' Rule

$$Pr[Y = y_1 \mid X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]$$

$$= \frac{\Pr[Y = y_1, X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}$$

$$= \frac{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k \mid Y = y_1] \cdot \Pr[Y = y_1]}{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}$$

$$= \frac{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k \mid Y = y_1] \cdot \Pr[Y = y_1]}{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}$$

## **Step 2.** We can simplify this using **Conditional Independence Assumption**

$$\frac{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k \mid Y = y_1] \cdot \Pr[Y = y_1]}{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}$$

## **Step 2.** We can simplify this using **Conditional Independence Assumption**

$$\frac{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k \mid Y = y_1] \cdot \Pr[Y = y_1]}{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}$$

**Suppose** conditional on the class, every predictor is independent of each other.

## **Step 2.** We can simplify this using **Conditional Independence Assumption**

$$\frac{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k \mid Y = y_1] \cdot \Pr[Y = y_1]}{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}$$

**Suppose** conditional on the class, every predictor is independent of each other.

$$Pr[X_1 = x_1, X_2 = x_2, ..., X_K = x_k \mid Y = y_1,]$$
  
= 
$$Pr[X_1 = x_1 \mid Y = y_1] \cdot Pr[X_2 = x_2 \mid Y = y_1] \cdot,$$
  
..., 
$$Pr[X_K = x_k \mid Y = y_1]$$

#### **Our Example** (with conditional independence)

```
\begin{split} &\text{Pr}[\mathsf{Name} \ = \ \mathsf{Le} \ \mathsf{Wang}, \mathsf{Age} \ = 38, \\ &\text{Gender} \ = \ \mathsf{Male}, \ \mathsf{Place} \ \mathsf{of} \ \mathsf{Birth} \ = \ \mathsf{Guangdong} \ , \\ &\text{Place} \ \mathsf{of} \ \mathsf{Residence} \ = \ \mathsf{Vermont} \ \mid \mathsf{Unemployment} = 1] \end{split}
```

#### Our Example (with conditional independence)

```
Pr[Name = Le Wang, Age = 38,
Gender = Male, Place of Birth = Guangdong,
Place of Residence = Vermont | Unemployment = 1 |
Pr[Name = Le Wang | Unemployment = 1]
```

- $\cdot \Pr[Age = 38 \mid Unemployment = 1]$
- $\cdot \Pr[\mathsf{Gender} = \mathsf{Male} \mid \mathsf{Unemployment} = 1]$
- $\cdot$  Pr[ Place of Birth = Guangdong | Unemployment = 1]
- · Pr[ Place of Residence = Vermont | Unemployment = 1]

$$\frac{\prod_{i=1}^{K} \Pr[X_i = x_i \mid Y = y_1] \cdot \Pr[Y = y_1]}{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}$$

$$\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]$$

$$\frac{\prod_{i=1}^{K} \Pr[X_i = x_i \mid Y = y_1] \cdot \Pr[Y = y_1]}{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}$$

$$\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]$$

$$Pr[X_{1} = x_{1}, X_{2} = x_{2}, ..., X_{K} = x_{k}]$$

$$= Pr[X_{1} = x_{1}, X_{2} = x_{2}, ..., X_{K} = x_{k} \mid Y = y_{1}] \cdot Pr[Y = y_{1}] + Pr[X_{1} = x_{1}, X_{2} = x_{2}, ..., X_{K} = x_{k} \mid Y = y_{2}] \cdot Pr[Y = y_{2}] + Pr[X_{1} = x_{1}, X_{2} = x_{2}, ..., X_{K} = x_{k} \mid Y = y_{J}] \cdot Pr[Y = y_{J}]$$

$$\vdots$$

$$\frac{\prod_{i=1}^{K} \Pr[X_i = x_i \mid Y = y_1] \cdot \Pr[Y = y_1]}{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}$$

$$\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]$$

$$Pr[X_{1} = x_{1}, X_{2} = x_{2}, ..., X_{K} = x_{k}]$$

$$= \prod_{i=1}^{K} Pr[X_{i} = x_{i} \mid Y = y_{1}] \cdot Pr[Y = y_{1}] +$$

$$Pr[X_{1} = x_{1}, X_{2} = x_{2}, ..., X_{K} = x_{k} \mid Y = y_{2}] \cdot Pr[Y = y_{2}] +$$

$$\vdots$$

$$Pr[X_{1} = x_{1}, X_{2} = x_{2}, ..., X_{K} = x_{k} \mid Y = y_{J}] \cdot Pr[Y = y_{J}]$$

$$\frac{\prod_{i=1}^{K} \Pr[X_i = x_i \mid Y = y_1] \cdot \Pr[Y = y_1]}{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}$$

$$\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]$$

$$Pr[X_{1} = x_{1}, X_{2} = x_{2}, ..., X_{K} = x_{k}]$$

$$= \prod_{i=1}^{K} Pr[X_{i} = x_{i} \mid Y = y_{1}] \cdot Pr[Y = y_{1}] +$$

$$\prod_{i=1}^{K} Pr[X_{i} = x_{i} \mid Y = y_{1}] \cdot Pr[Y = y_{2}] +$$

$$\vdots$$

$$Pr[X_{1} = x_{1}, X_{2} = x_{2}, ..., X_{K} = x_{k} \mid Y = y_{I}] \cdot Pr[Y = y_{I}]$$

$$\frac{\prod_{i=1}^{K} \Pr[X_i = x_i \mid Y = y_1] \cdot \Pr[Y = y_1]}{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}$$

$$\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]$$

$$Pr[X_{1} = x_{1}, X_{2} = x_{2}, ..., X_{K} = x_{k}]$$

$$= \prod_{i=1}^{K} Pr[X_{i} = x_{i} \mid Y = y_{1}] \cdot Pr[Y = y_{1}] + \prod_{i=1}^{K} Pr[X_{i} = x_{i} \mid Y = y_{1}] \cdot Pr[Y = y_{2}] + \vdots$$

$$\vdots$$

$$\prod_{i=1}^{K} Pr[X_{i} = x_{i} \mid Y = y_{1}] \cdot Pr[Y = y_{J}]$$

 $\begin{tabular}{ll} \textbf{Step 3.} & \textbf{Apply conditional independence and law of total probability} \\ \textbf{(as we did before!)} \end{tabular}$ 

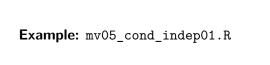
$$\frac{\prod_{i=1}^{K} \Pr[X_i = x_i \mid Y = y_1] \cdot \Pr[Y = y_1]}{\Pr[X_1 = x_1, X_2 = x_2, \dots, X_K = x_k]}$$

$$= \frac{\prod_{i=1}^{K} \Pr[X_i = x_i \mid Y = y_1] \cdot \Pr[Y = y_1]}{\sum_{i=1}^{J} \prod_{i=1}^{K} \Pr[X_i = x_i \mid Y = y_i] \cdot \Pr[Y = y_i]}$$

### Our Example (with conditional independence)

### Our Example (with conditional independence)

```
(Pr[Name = Le Wang | Unemployment = 1]
\cdot \Pr[Age = 38 \mid Unemployment = 1]
\cdot Pr[ Gender = Male | Unemployment = 1]
\cdot Pr[ Place of Birth = Guangdong | Unemployment = 1]
\cdot Pr[ Place of Residence = Vermont | Unemployment = 1])
\times Pr[Unemployment = 1] +
(Pr[Name = Le Wang | Unemployment = 0]
\cdot \Pr[Age = 38 \mid Unemployment = 0]
\cdot \Pr[\mathsf{Gender} = \mathsf{Male} \mid \mathsf{Unemployment} = 0]
\cdot Pr[ Place of Birth = Guangdong | Unemployment = 0]
\cdot Pr[ Place of Residence = Vermont | Unemployment = 0])
\times Pr[Unemployment = 0]
```



## **Applications**

### Combining conditional independence and Bayes' Rule

- 1. Statistical Language Model
- 2. Naive Bayes Classifier
- 3. Predicting Unobservable Ethnicity in Political Science.'

### Frontier Research: Predicting Unobservable Ethnicity

#### **Motivation:**

In both political behavior research and voting rights litigation, turnout and vote choice for different racial groups are often inferred using aggregate election results and racial composition. These predictions are often inaccurate.

How to reduce aggregation bias by predicting individual-level ethnicity from voter registration records.

## Frontier Research: Predicting Unobservable Ethnicity

#### **Motivation:**

In both political behavior research and voting rights litigation, turnout and vote choice for different racial groups are often inferred using aggregate election results and racial composition. These predictions are often inaccurate.

How to reduce aggregation bias by predicting individual-level ethnicity from voter registration records.

Imai and Khanna (2016), Political Analysis

- 1. Conditional Probability.
- Bayes's rule to combine the Census Bureau's Surname List with various information from geocoded voter registration records.

To classify someone's race:		
	Pr[race]	

To classify someone's race:		
	Pr[race]	
	Pr[race surname]	

To classify someone's race:

Pr[race]

Pr[race|surname]

Pr[race|surname, residence]

We will use the maximum to forecast. If no further information,

everyone is predicted to be white!

#### Purpose of this exercise

To validate the accuracy of the predictions of individual race, we will use the sample of 10,000 registered voters from Florida.

In Florida, voters are asked to self-report their race when registering.

#### **Datasets**

- 1. Florida Registered Voter Data
- 2. Census Data

# Dataset (I): Florida Registered Voter Data

surname Surname county County id of voter's residence VTD voting district id of voter's residen	Variable	Description
age age gender: $m = male$ and $f = femal$ race self-reported race	county VTD age gender	County id of voter's residence voting district id of voter's residence age gender: $m = male$ and $f = female$

# Dataset (II): Census Data

Variable	Description
surname count pctwhite	Surname number of individuals with a specific surname percentage of non-Hispanic whites among those who have a specific surname
pctblack	percentage of non-Hispanic blacks among those who have a specific surname
pctapi	percentage of non-Hispanic Asians and Pacific Islanders among those who have a specific surname
pcthispanic	percentage of Hispanic origin among those who have a specific surname