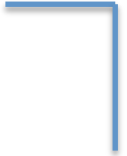


UVA CS 4501: Machine Learning



Lecture 18: Generative Bayes Classifiers

Dr. Yanjun Qi

University of Virginia

Department of
Computer Science

Where are we ? ➔

Major sections of this course

- ❑ Regression (supervised)
- ➔ ❑ Classification (supervised)
- ❑ Unsupervised models
- ❑ Learning theory

Where are we ? ➔

Three major sections for classification

- We can divide the large variety of classification approaches into **roughly three major types**

1. Discriminative

- directly estimate a decision rule/boundary
- e.g., **support vector machine**, decision tree



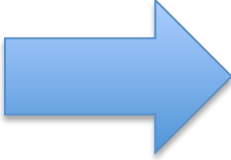
2. Generative:

- build a generative statistical model
- e.g., **naïve bayes classifier**, Bayesian networks

3. Instance based classifiers

- Use observation directly (no models)
- e.g. **K nearest neighbors**

Today : Generative Bayes Classifiers

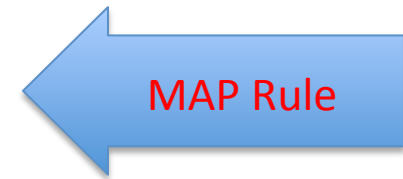
- 
- ✓ Bayes Classifier
 - MAP classification rule
 - Generative Bayes Classifier
 - ✓ Naïve Bayes Classifier

Review: Notations

- Inputs
 - X, X_j (jth element of vector X) : random variables written in capital letter
 - p #input features, n #observations
 - X : matrix written in bold capital
 - Vectors are assumed to be column vectors
- Outputs
 - quantitative Y
 - qualitative C (for categorical)

Review: Bayes classifiers

- Treat each feature attribute and the class label as random variables.
- Given a sample \mathbf{x} with attributes (x_1, x_2, \dots, x_p) :
 - Goal is to predict its class c .
 - Specifically, we want to find the class that maximizes $p(c | x_1, x_2, \dots, x_p)$.



Review: Bayes Classifiers – MAP Rule

Task: Classify a new instance X based on a tuple of attribute values $X = \langle X_1, X_2, \dots, X_p \rangle$ into one of the classes

$$c_{MAP} = \operatorname{argmax}_{c_j \in C} P(c_j | x_1, x_2, \dots, x_p)$$



MAP Rule

MAP = Maximum A posteriori Probability

Please read the L13-Logistic for details

Review: Establishing a probabilistic model for classification

– (1) Discriminative model

$$\operatorname{argmax}_{c \in \mathcal{C}} P(c | \mathbf{X}), \quad \mathcal{C} = \{c_1, \dots, c_L\}$$

$$P(c_1 | \mathbf{x}) \quad P(c_2 | \mathbf{x}) \quad \dots \quad P(c_L | \mathbf{x})$$

**Discriminative
Probabilistic Classifier**

$$x_1 \quad x_2 \quad \dots \quad x_p$$

$$\mathbf{x} = (x_1, x_2, \dots, x_p)$$

*logistic
regression*

Bayes classifiers

→ MAP classification rule

- Establishing a probabilistic model for classification
 - **(1) Discriminative**
 - **(2) Generative**

X_1	X_2	X_3	C

A Dataset for classification

$$f : X \rightarrow C$$

Output as Discrete
Class Label

C_1, C_2, \dots, C_L

Discriminative

$$\operatorname{argmax}_{c \in C} P(c | \mathbf{X}) \quad C = \{c_1, \dots, c_L\}$$

Generative

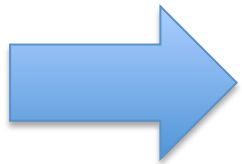
$$\operatorname{argmax}_{c \in C} P(c | X) = \operatorname{argmax}_{c \in C} P(X, c) = \operatorname{argmax}_{c \in C} P(X | c) P(c)$$

- **Data/points/instances/examples/samples/records:** [rows]
- **Features/attributes/dimensions/independent variables/covariates/predictors/regressors:** [columns, except the last]
- **Target/outcome/response/label/dependent variable:** special column to be predicted [last column]

this lecture!

Today : Generative Bayes Classifiers

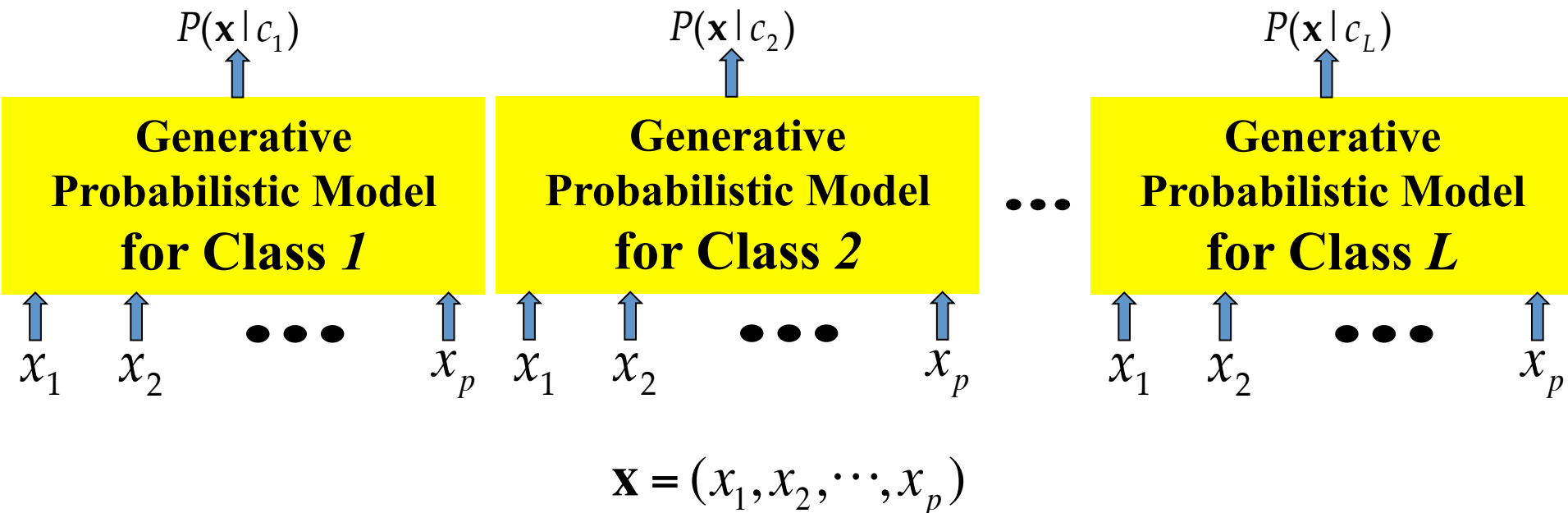
- ✓ Bayes Classifier
 - MAP classification rule
 - Generative Bayes Classifier
- ✓ Naïve Bayes Classifier



(2) Generative

$$P(\mathbf{X} | C),$$

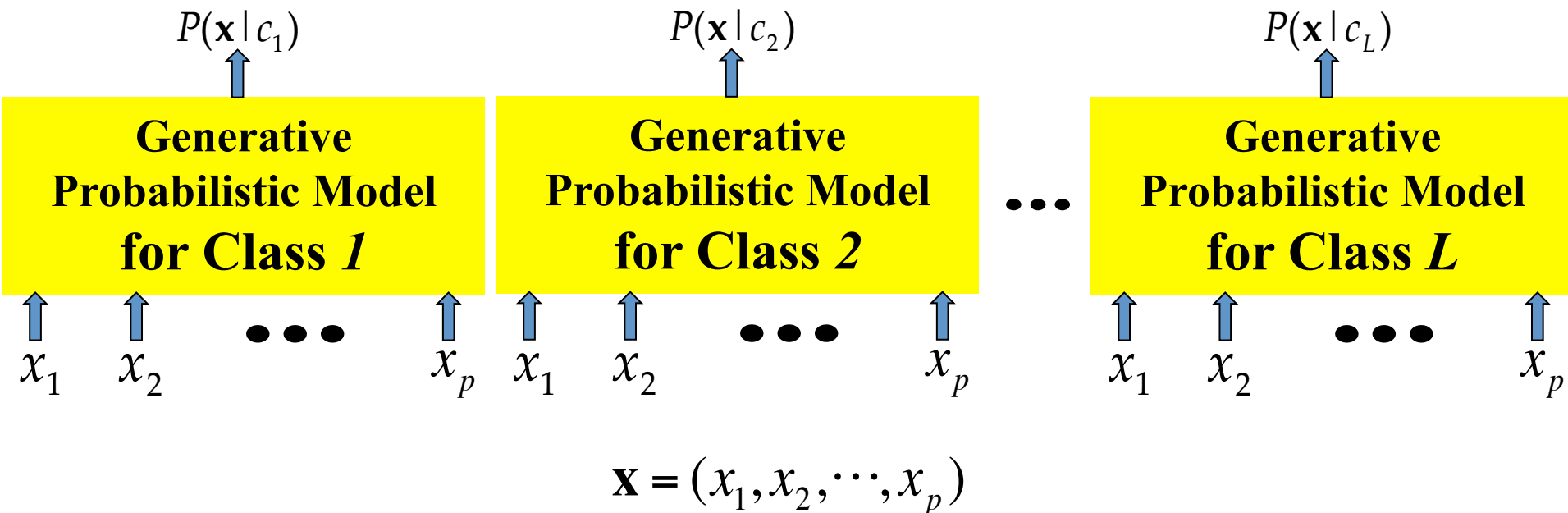
$$C = c_1, \dots, c_L, \mathbf{X} = (X_1, \dots, X_p)$$



Generative BC

$$P(\mathbf{X} | C),$$

$$C = c_1, \dots, c_L, \mathbf{X} = (X_1, \dots, X_p)$$



Generative BC

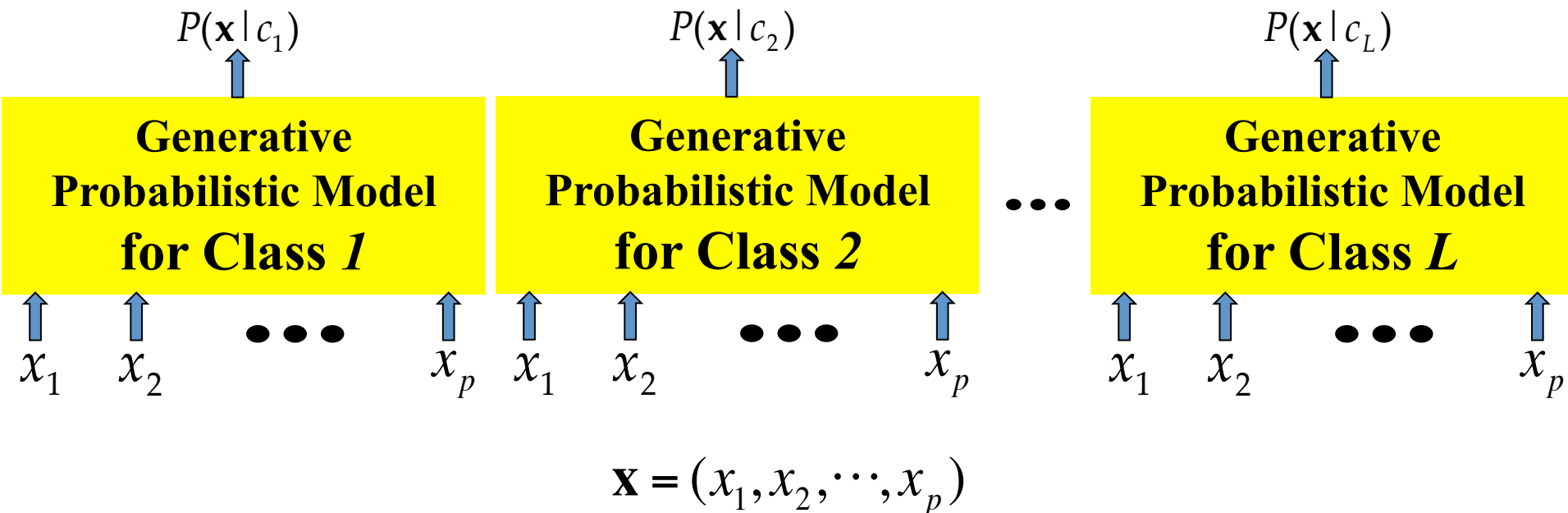
$$p(x|c_1), p(x|c_2), \dots, p(x|c_L)$$

$$P(\mathbf{X} | C),$$

$$C = c_1, \dots, c_L, \mathbf{X} = (X_1, \dots, X_p)$$

$$\Rightarrow p(c|x)$$

MAP rule



Review : Bayes' Rule

– for Generative Bayes Classifiers

$$P(C, X) = P(C | X)P(X) = P(X | C)P(C)$$



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

Review Probability:

If hard to directly estimate from data, most likely we can estimate

- 1. Joint probability
 - Use Chain Rule
- 2. Marginal probability
 - Use the total law of probability
- 3. Conditional probability
 - Use the Bayes Rule

Review : Bayes' Rule

– for Generative Bayes Classifiers

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

$$P(C_1|x), P(C_2|x), \dots, P(C_L|x)$$

$$P(C_1), P(C_2), \dots, P(C_L)$$

$$P(C_i | \mathbf{X}) = \frac{P(\mathbf{X} | C_i)P(C_i)}{P(\mathbf{X})}$$

Review : Bayes' Rule

– for Generative Bayes Classifiers

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

Posterior

$P(C_1|x), P(C_2|x), \dots, P(C_L|x)$

Prior

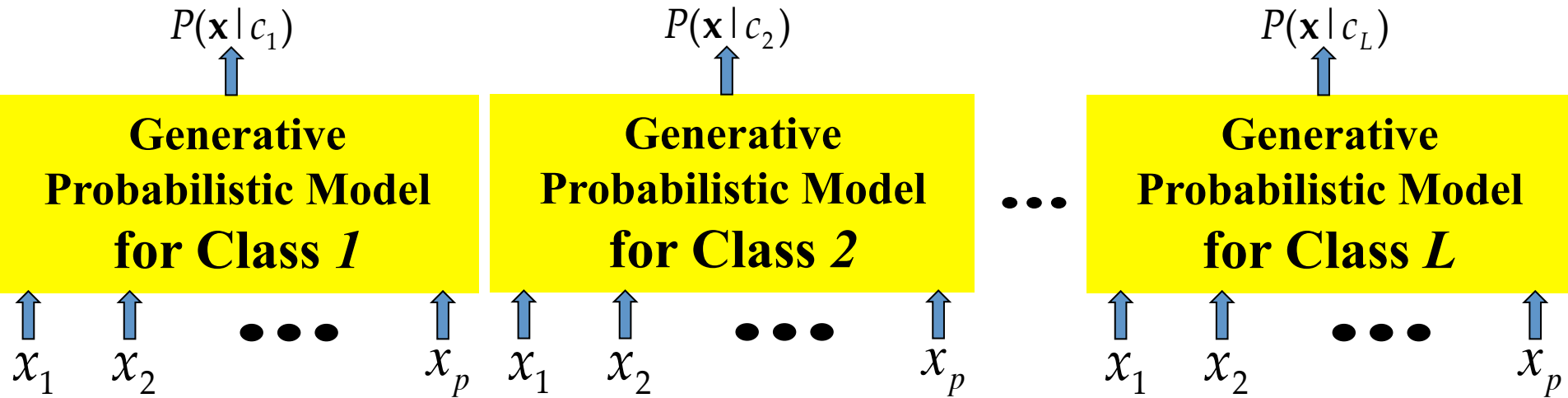
$P(C_1), P(C_2), \dots, P(C_L)$

$$P(C | \mathbf{X}) = \frac{P(\mathbf{X} | C)P(C)}{P(\mathbf{X})}$$

Establishing a probabilistic model for classification through generative modeling

[MAP rule]

$$\operatorname{argmax}_{C_i} P(C_i | X) = \operatorname{argmax}_{C_i} P(X, C_i) = \operatorname{argmax}_{C_i} P(X | C_i) P(C_i)$$



$$\mathbf{X} = (x_1, x_2, \dots, x_p)$$

Summary:

Generative classification with the MAP rule

- MAP classification rule
 - **MAP**: **M**aximum **A** **P**osterior
 - Assign x to c^* if

$$P(C = c^* | \mathbf{X} = \mathbf{x}) > P(C = c | \mathbf{X} = \mathbf{x}) \quad c \neq c^*, c = c_1, \dots, c_L$$

Summary:

Generative classification with the MAP rule

$$P(C = c^* | \mathbf{X} = \mathbf{x}) > P(C = c | \mathbf{X} = \mathbf{x}) \quad c \neq c^*, c = c_1, \dots, c_L$$

- Generative classification with the MAP rule
 - Apply Bayes rule to convert them into posterior probabilities

$$\begin{aligned} P(C = c_i | \mathbf{X} = \mathbf{x}) &= \frac{P(\mathbf{X} = \mathbf{x} | C = c_i)P(C = c_i)}{P(\mathbf{X} = \mathbf{x})} \\ &\propto P(\mathbf{X} = \mathbf{x} | C = c_i)P(C = c_i) \\ &\quad \text{for } i = 1, 2, \dots, L \end{aligned}$$

- Then apply the MAP rule

An Example

- Example: Play Tennis

PlayTennis: training examples

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

An Example

- Example: Play Tennis

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

$k_2=3$

X_2 :
{Hot, Mild, Cool}

$X_3 = \{ \text{High, Normal} \}$

$k_3=2$

$X_4 = (W, S)$
 $k_4=$

PlayTennis: training examples

C

$C: \{ \text{Yes, No} \}$
 $(L=2)$

$X_1: \{ \text{sunny, overc, rain} \}$
 $(k_1=3)$

Example

$P(X_1, X_2, X_3, X_4 | \text{Yes})$
 $P(X_1, X_2, X_3, X_4 | \text{No})$

- Example: Play Tennis

3 PlayTennis: training examples

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

$P(C=C_i)$

$\begin{cases} \text{Yes} \\ \text{No} \end{cases}$

$P(C=\text{Yes}) = 9/14$

$P(C=\text{No}) = 5/14$

$3 \times 3 \times 2 \times 2 = 36$

$$P(C = \text{Yes} \mid X_1, X_2, X_3, X_4)$$

$$P(C = \text{No} \mid X_1, X_2, X_3, X_4)$$

$$\rightarrow P(C_1 = \text{Yes}) = 9/14$$

$$P(C_2 = \text{No}) = 5/14$$

$$\rightarrow P(X_1, X_2, X_3, X_4 \mid C_i)$$

3 \times 3 \times 2 \times 2 \times 2 \Rightarrow 72 parameters from train

$$\rightarrow \underset{\hat{i}=1,2}{\operatorname{argmax}} P(\bar{X}_{ts} \mid C_i) P(C_i) \quad \text{Generative BC}$$

- maximum likelihood estimates
 - simply use the frequencies in the data

e.g. $p(\text{overcast, hot, high, weak} \mid \text{Yes}) = \frac{1}{9}$

PlayTennis: training examples

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

e.g.
 $p(\text{overcast, hot, high, weak} \mid \text{No})$

$$= \frac{0}{9}$$

Check L12-MLE
Lecture for
Why

Generative Bayes Classifier:

- Learning Phase

$$P(C_1), P(C_2), \dots, P(C_L)$$

$$P(\text{Play}=\text{Yes}) = 9/14 \quad P(\text{Play}=\text{No}) = 5/14$$

$$P(X_1, X_2, \dots, X_p | C_1), P(X_1, X_2, \dots, X_p | C_2)$$

a look up table of cond. prob

	X_1	X_2	X_3	C
S_1				
S_2				
S_3				
S_4				
S_5				
S_6				

Outlook (3 values)	Temperature (3 values)	Humidity (2 values)	Wind (2 values)	Play=Yes	Play=No
<i>sunny</i>	<i>hot</i>	<i>high</i>	<i>weak</i>	0/9	1/5
<i>sunny</i>	<i>hot</i>	<i>high</i>	<i>strong</i>	.../9	.../5
<i>sunny</i>	<i>hot</i>	<i>normal</i>	<i>weak</i>	.../9	.../5
<i>sunny</i>	<i>hot</i>	<i>normal</i>	<i>strong</i>	.../9	.../5
....
....
....
....

3*3*2*2 [conjunctions of attributes] * 2 [two classes]=72 parameters

Generative Bayes Classifier:

$$[\hat{P}(a'_1 | c^*) \cdots \hat{P}(a'_p | c^*)] \hat{P}(c^*) > [\hat{P}(a'_1 | c) \cdots \hat{P}(a'_p | c)] \hat{P}(c)$$

- Test Phase

- Given an unknown instance $\mathbf{X}'_{ts} = (a'_1, \dots, a'_p)$
- Look up tables to assign the label c^* to \mathbf{X}_{ts} if

$$\hat{P}(a'_1, \dots, a'_p | c^*) \hat{P}(c^*) > \hat{P}(a'_1, \dots, a'_p | c) \hat{P}(c),$$

$$c \neq c^*, c = c_1, \dots, c_L$$

Last
Page

- Given a new instance,

$\mathbf{x}' = (\text{Outlook}=\text{Sunny}, \text{Temperature}=\text{Cool}, \text{Humidity}=\text{High}, \text{Wind}=\text{Strong})$

$$\left\{ \begin{array}{l} p(\mathbf{x}' | \text{Yes}) p(c = \text{Yes}) \\ p(\mathbf{x}' | \text{No}) p(c = \text{No}) \end{array} \right\} \Rightarrow \arg \max_c \Rightarrow \text{predicted } c^*$$

Today : Generative Bayes Classifiers

- ✓ Bayes Classifier
 - MAP classification rule
 - Generative Bayes Classifier
- ✓ Naïve Bayes Classifier

Naïve Bayes Classifier

- Bayes classification

$$\operatorname{argmax}_{c_j \in \mathcal{C}} P(x_1, x_2, \dots, x_p | c_j) P(c_j)$$

Difficulty: learning the joint probability

- Naïve Bayes classification
 - Assumption that **all input attributes are conditionally independent!**

Naïve Bayes Classifier

- Bayes classification

$$\operatorname{argmax}_{c_j \in \mathcal{C}} P(x_1, x_2, \dots, x_p | c_j) P(c_j)$$

$$= p(x_1 | c_j) p(x_2 | c_j) \dots p(x_p | c_j)$$

Difficulty: learning the joint probability

- Naïve Bayes classification
 - Assumption that **all input attributes are conditionally independent!**
given C variable

Naïve Bayes Classifier

- Naïve Bayes classification
 - Assumption that **all input attributes are conditionally independent!**

$$\begin{aligned}P(X_1, X_2, \dots, X_p | C) &= P(X_1 | X_2, \dots, X_p, C) P(X_2, \dots, X_p | C) \\&= P(X_1 | C) P(X_2, \dots, X_p | C) \\&= P(X_1 | C) P(X_2 | C) \cdots P(X_p | C)\end{aligned}$$

Naïve Bayes Classifier

- Naïve Bayes classification
 - Assumption that **all input attributes are conditionally independent!**

$$P(X_1, X_2, \dots, X_p | C) = P(X_1 | C)P(X_2 | C) \cdots P(X_p | C)$$

- MAP classification rule: for a sample $\mathbf{x} = (x_1, x_2, \dots, x_p)$

$$[P(x_1 | c^*) \cdots P(x_p | c^*)]P(c^*) > [P(x_1 | c) \cdots P(x_p | c)]P(c),$$

$$c \neq c^*, c = c_1, \dots, c_L$$

Naïve Bayes Classifier

- Naïve Bayes classification
 - Assumption that **all input attributes are conditionally independent!**

$$P(X_1, X_2, \dots, X_p | C) = P(X_1 | C) P(X_2 | C) \cdots P(X_p | C)$$

- MAP classification rule: for a sample $\mathbf{x} = (x_1, x_2, \dots, x_p)$

$$\underbrace{[P(x_1 | c^*) \cdots P(x_p | c^*)] P(c^*)}_{> [P(x_1 | c) \cdots P(x_p | c)] P(c)},$$

$$c \neq c^*, c = c_1, \dots, c_L$$

$$\Rightarrow \arg \max_{i=1, \dots, L} P(c_i) P(x_1 | c_i) P(x_2 | c_i) \cdots P(x_p | c_i)$$

Naïve Bayes Classifier

- Naïve Bayes classification
 - Assumption that **all input attributes are conditionally independent!**

$$P(X_1, X_2, \dots, X_p | C) = P(X_1 | C) P(X_2 | C) \dots P(X_p | C)$$

- MAP classification rule: for a sample $\mathbf{x} = (x_1, x_2, \dots, x_p)$

$$[P(x_1 | c^*) \dots P(x_p | c^*)]P(c^*) > [P(x_1 | c) \dots P(x_p | c)]P(c),$$

$$c \neq c^*, c = c_1, \dots, c_L$$

$\{i=1,2,\dots,p\} \quad \{i=1,2,\dots,L\}$

$$P(X_i | c_i)$$

Bernoulli (P)
 Binomial (K,P)
 Multinomial
 Gaussian

Naïve Bayes Classifier (for discrete input attributes) - training

- Naïve Bayes Algorithm (for discrete input attributes)
 - **Learning Phase:** Given a training set S ,

For each target value of c_i ($c_i = c_1, \dots, c_L$)

$\hat{P}(C = c_i) \leftarrow$ estimate $P(C = c_i)$ with examples in S ;

Naïve Bayes Classifier

(for discrete input attributes) - training

- Naïve Bayes Algorithm (for discrete input attributes)
 - **Learning Phase:** Given a training set S ,

For each target value of c_i ($c_i = c_1, \dots, c_L$)

$\hat{P}(C = c_i) \leftarrow$ estimate $P(C = c_i)$ with examples in S ;

$p(c_1) p(c_2) \dots p(c_L)$
 L parameters

For every attribute value x_{jk} of each attribute X_j ($j = 1, \dots, p$; $k = 1, \dots, K_j$)

$\hat{P}(X_j = x_{jk} | C = c_i) \leftarrow$ estimate $P(X_j = x_{jk} | C = c_i)$ with examples in S ;

Output: conditional probability tables; for $X_j, K_j \times L$ elements

Naïve Bayes Classifier

(for discrete input attributes) - training

- Naïve Bayes Algorithm (for discrete input attributes)
 - **Learning Phase:** Given a training set S ,

For each target value of c_i ($c_i = c_1, \dots, c_L$)

$\hat{P}(C = c_i) \leftarrow$ estimate $P(C = c_i)$ with examples in S ;

For every attribute value x_{jk} of each attribute X_j ($j = 1, \dots, p$; $k = 1, \dots, K_j$)

$\hat{P}(X_j = x_{jk} | C = c_i) \leftarrow$ estimate $P(X_j = x_{jk} | C = c_i)$ with examples in S ;

Output: conditional probability tables; for $X_j, K_j \times L$ elements

K_1, K_2, \dots, K_p
 $\{X_1, X_2, \dots, X_p\}$

$K_1 \times L +$
 $K_2 \times L +$
 $K_p \times L$

Naïve Bayes

(for discrete input attributes) - testing

- Naïve Bayes Algorithm (for discrete input attributes)

– **Test Phase:** Given an unknown instance $\mathbf{X}' = (a'_1, \dots, a'_p)$

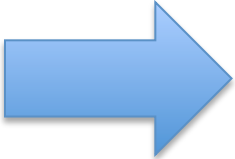
Look up tables to assign the label c^* to \mathbf{X}' if

$$[\hat{P}(a'_1 | c^*) \cdots \hat{P}(a'_p | c^*)] \hat{P}(c^*) > \underbrace{[\hat{P}(a'_1 | c) \cdots \hat{P}(a'_p | c)] \hat{P}(c)},$$

$$c \neq c^*, c = c_1, \dots, c_L$$

$$\begin{aligned} & P(\mathbf{X}' | c_i) P(c_i) \\ &= P(a'_1 | c_i) P(a'_2 | c_i) \cdots P(a'_p | c_i) P(c_i) \\ & \quad i=1, 2, \dots, L \end{aligned}$$

Today : Generative Bayes Classifiers

- ✓ Bayes Classifier
 - MAP classification rule
 - Generative Bayes Classifier
- ✓ Naïve Bayes Classifier
-  ✓ NBC for discrete input variables

An Example

- Example: Play Tennis

PlayTennis: training examples

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

X_1	X_2	X_3	C

An Example

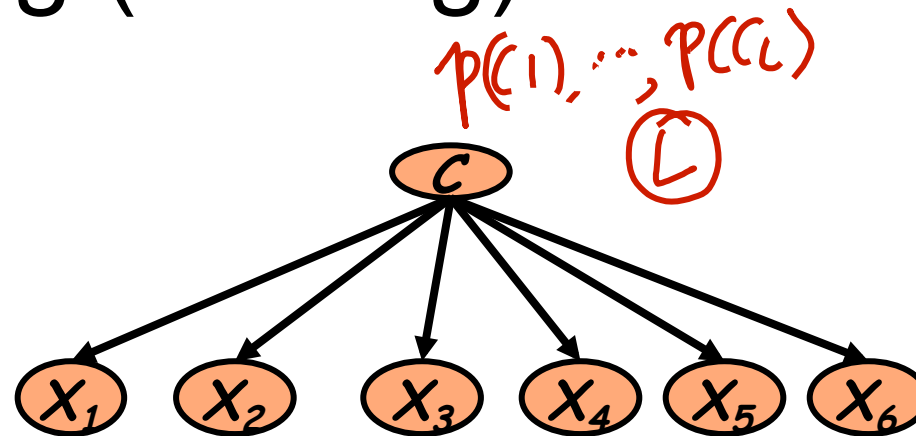
- Example: Play Tennis

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

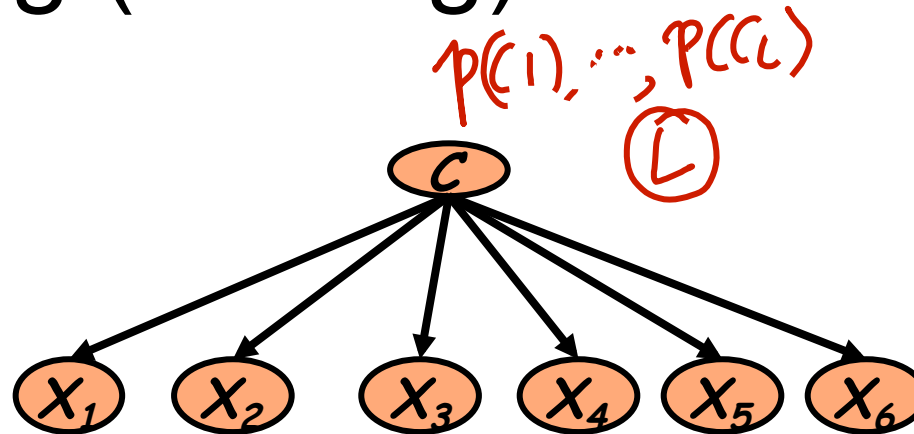
[illegible]

$C: \{ \text{Yes,} \}$
 $(L=2) \text{ No} \}$
 $X_1: \{ \text{sunny,} \}$
 overc,
 $(k=3) \text{ rain} \}$

Learning (training) the NBC Model



Learning (training) the NBC Model

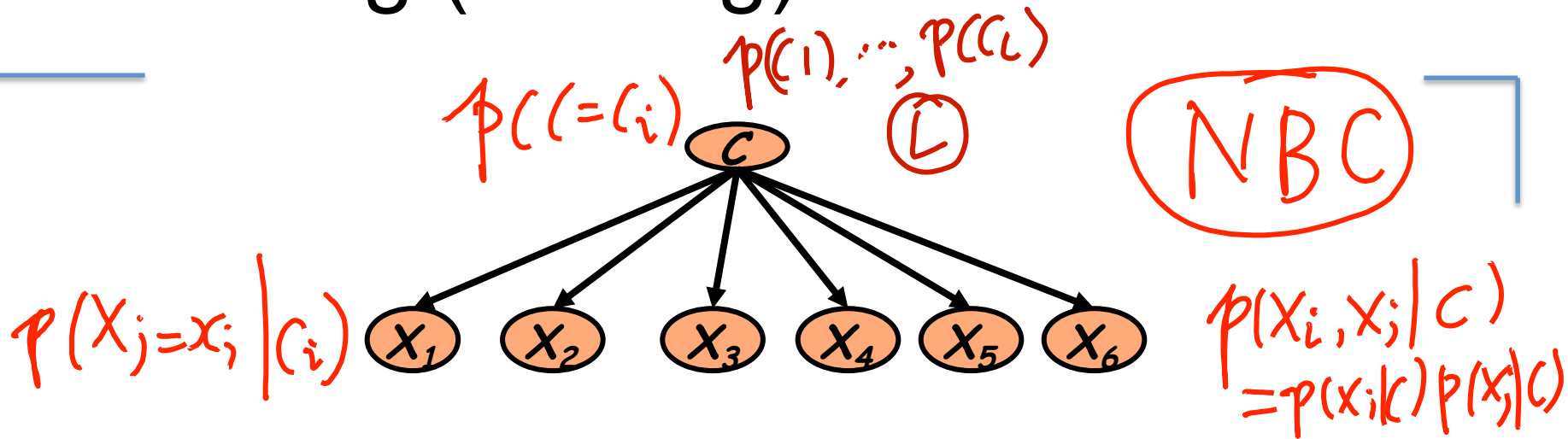


- maximum likelihood estimates:
 - simply use the frequencies in the data

$$\hat{P}(c_j) = \frac{N(C = c_j)}{N}$$

$$\hat{P}(x_i | c_j) = \frac{N(X_i = x_i, C = c_j)}{N(C = c_j)}$$

Learning (training) the NBC Model



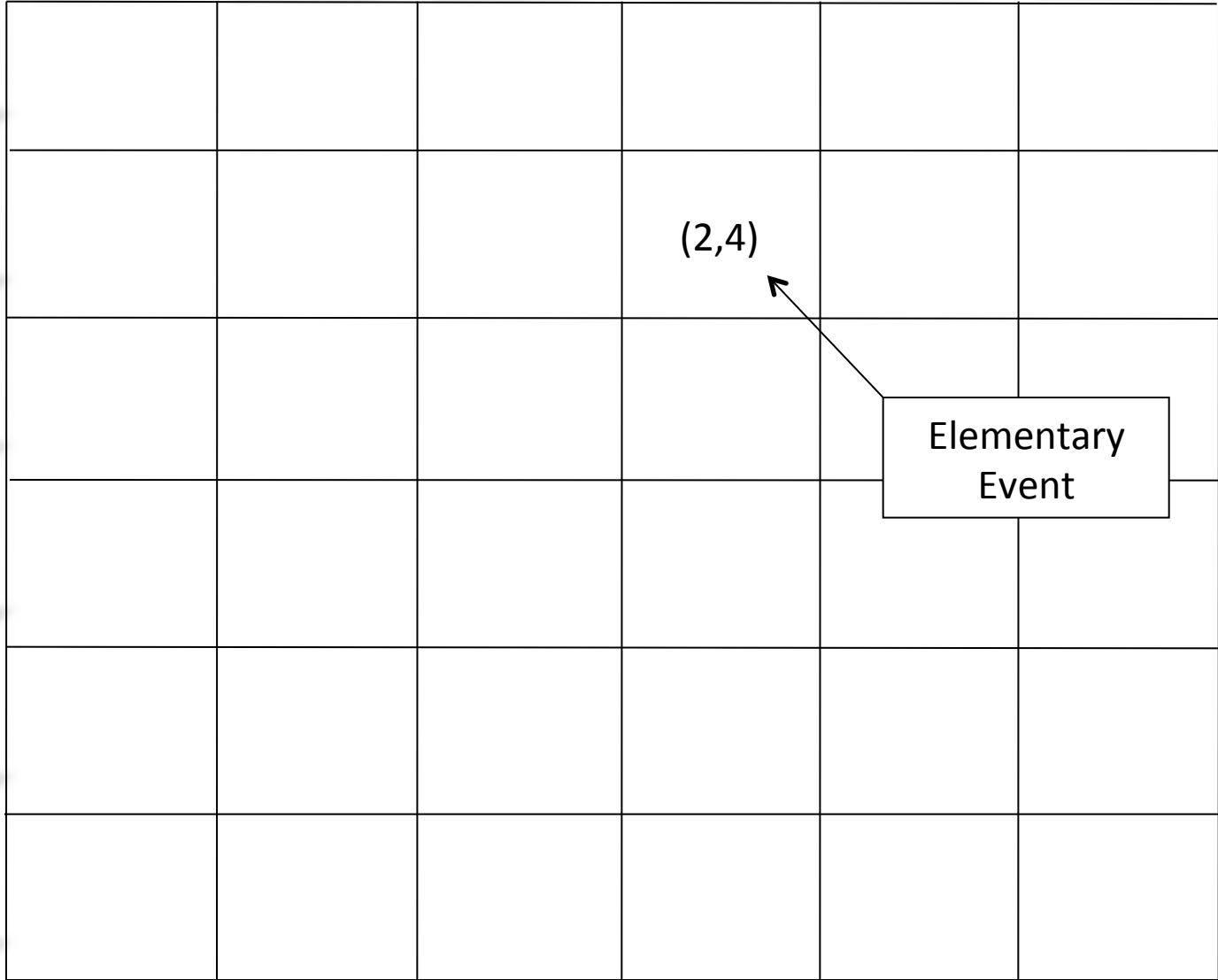
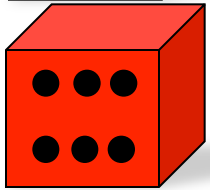
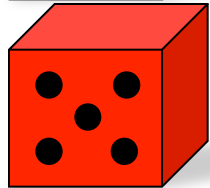
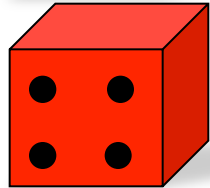
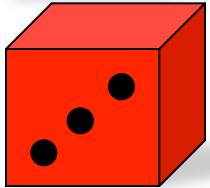
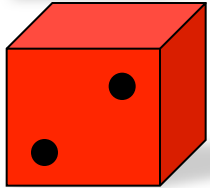
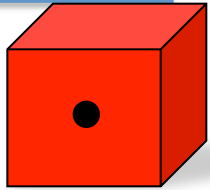
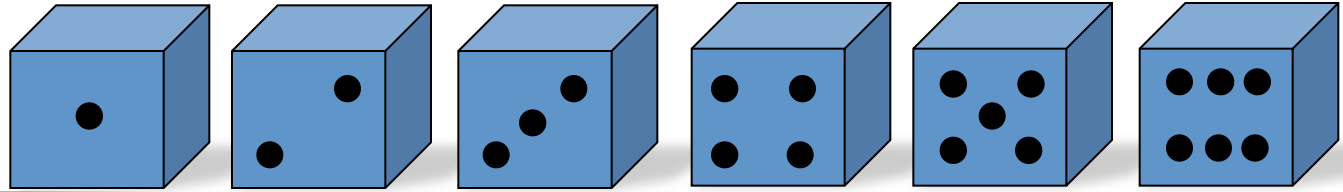
- maximum likelihood estimates
 - simply use the frequencies in the data

$$\hat{P}(c_j) = \frac{N(C = c_j)}{N}$$

$$\hat{P}(x_i | c_j) = \frac{N(X_i = x_i, C = c_j)}{N(C = c_j)}$$

$C = c_i$ 2 – Dimensional

$$X_j = x_{jk}$$



PlayTennis: training examples

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

$$P(X_1 = \text{Rain} \mid C = \text{Yes})$$

$$= \frac{3}{9}$$

$$P(X_1 = \text{Rain} \mid C = \text{No})$$

$$= \frac{2}{5}$$

Counting
↑

Learning Phase

Estimate $P(X_j = x_{jk} | C = c_i)$ with examples in training;

$P(X_2|C_1), P(X_2|C_2)$

Outlook	Play=Yes	Play=No
Sunny	2/9	3/5
Overcast	4/9	0/5
Rain	3/9	2/5

Temperature	Play=Yes	Play=No
Hot	2/9	2/5
Mild	4/9	2/5
Cool	3/9	1/5

Humidity	Play=Yes	Play=No
High	3/9	4/5
Normal	6/9	1/5

$P(X_4|C_1), P(X_4|C_2)$

Wind	Play=Yes	Play=No
Strong	3/9	3/5
Weak	6/9	2/5

3+3+2+2 [naïve assumption] * 2 [two classes]= 20 parameters

$$P(\text{Play=Yes}) = 9/14 \quad P(\text{Play=No}) = 5/14$$

$P(C_1), P(C_2), \dots, P(C_L)$

Counting
↑

Learning Phase

Estimate $P(X_j = x_{jk} | C = c_i)$ with examples in training;

$P(X_2|C_1), P(X_2|C_2)$

X_1

Outlook	Play=Yes	Play=No
Sunny	2/9	3/5
Overcast	4/9	0/5
Rain	3/9	2/5

Temperature	Play=Yes	Play=No
Hot	2/9	2/5
Mild	4/9	2/5
Cool	3/9	1/5

Humidity	Play=Yes	Play=No
High	3/9	4/5
Normal	6/9	1/5

$P(X_4|C_1), P(X_4|C_2)$

Wind	Play=Yes	Play=No
Strong	3/9	3/5
Weak	6/9	2/5

3+3+2+2 [naïve assumption] * 2 [two classes] = 20 parameters

$$P(\text{Play}=\text{Yes}) = 9/14 \quad P(\text{Play}=\text{No}) = 5/14$$

$P(C_1), P(C_2), \dots, P(C_L)$

$P(C_i)$

Testing the NBC Model

look up

$$[\hat{P}(a'_1 | c^*) \cdots \hat{P}(a'_p | c^*)] \hat{P}(c^*) > [\hat{P}(a'_1 | c) \cdots \hat{P}(a'_p | c)] \hat{P}(c)$$

- Test Phase
 - Given a new instance,
 $\mathbf{x}' = (\text{Outlook}=\textit{Sunny}, \text{Temperature}=\textit{Cool}, \text{Humidity}=\textit{High}, \text{Wind}=\textit{Strong})$

Testing the NBC Model

$$[\hat{P}(a'_1 | c^*) \cdots \hat{P}(a'_p | c^*)] \hat{P}(c^*) > [\hat{P}(a'_1 | c) \cdots \hat{P}(a'_p | c)] \hat{P}(c)$$

- Test Phase

- Given a new instance,

$\mathbf{x}' = (\text{Outlook}=\text{Sunny}, \text{Temperature}=\text{Cool}, \text{Humidity}=\text{High}, \text{Wind}=\text{Strong})$

$$\begin{aligned} &\rightarrow P(c_1) P(\text{Sunny} | c_1) P(\text{Cool} | c_1) P(\text{High} | c_1) P(\text{Strong} | c_1) \\ &= \frac{9}{14} \times \frac{2}{9} \cdots \cdots = \\ &\rightarrow P(c_2) P(\text{Su} | c_2) P(\text{Co} | c_2) P(\text{hi} | c_2) P(\text{St} | c_2) \\ &= \frac{5}{14} \times \frac{3}{5} \times \cdots = \end{aligned}$$

Testing the NBC Model

$$[\hat{P}(a'_1 | c^*) \cdots \hat{P}(a'_p | c^*)] \hat{P}(c^*) > [\hat{P}(a'_1 | c) \cdots \hat{P}(a'_p | c)] \hat{P}(c)$$

- Test Phase

- Given a new instance,

$\mathbf{x}' = (\text{Outlook}=\textit{Sunny}, \text{Temperature}=\textit{Cool}, \text{Humidity}=\textit{High}, \text{Wind}=\textit{Strong})$

- Look up in conditional-prob tables

$P(\text{Outlook}=\textit{Sunny} | \text{Play}=\textit{Yes}) = 2/9$

$P(\text{Temperature}=\textit{Cool} | \text{Play}=\textit{Yes}) = 3/9$

$P(\text{Humidity}=\textit{High} | \text{Play}=\textit{Yes}) = 3/9$

$P(\text{Wind}=\textit{Strong} | \text{Play}=\textit{Yes}) = 3/9$

$P(\text{Play}=\textit{Yes}) = 9/14$

$P(\text{Outlook}=\textit{Sunny} | \text{Play}=\textit{No}) = 3/5$

$P(\text{Temperature}=\textit{Cool} | \text{Play}=\textit{No}) = 1/5$

$P(\text{Humidity}=\textit{High} | \text{Play}=\textit{No}) = 4/5$

$P(\text{Wind}=\textit{Strong} | \text{Play}=\textit{No}) = 3/5$

$P(\text{Play}=\textit{No}) = 5/14$

Testing the NBC Model

$$[\hat{P}(a'_1 | c^*) \cdots \hat{P}(a'_p | c^*)] \hat{P}(c^*) > [\hat{P}(a'_1 | c) \cdots \hat{P}(a'_p | c)] \hat{P}(c)$$

• Test Phase

- Given a new instance,

$\mathbf{x}' = (\text{Outlook}=\textit{Sunny}, \text{Temperature}=\textit{Cool}, \text{Humidity}=\textit{High}, \text{Wind}=\textit{Strong})$

- Look up in conditional-prob tables

$P(\text{Outlook}=\textit{Sunny} | \text{Play}=\textit{Yes}) = 2/9$

$P(\text{Temperature}=\textit{Cool} | \text{Play}=\textit{Yes}) = 3/9$

$P(\text{Humidity}=\textit{High} | \text{Play}=\textit{Yes}) = 3/9$

$P(\text{Wind}=\textit{Strong} | \text{Play}=\textit{Yes}) = 3/9$

$P(\text{Play}=\textit{Yes}) = 9/14$

$P(\text{Outlook}=\textit{Sunny} | \text{Play}=\textit{No}) = 3/5$

$P(\text{Temperature}=\textit{Cool} | \text{Play}=\textit{No}) = 1/5$

$P(\text{Humidity}=\textit{High} | \text{Play}=\textit{No}) = 4/5$

$P(\text{Wind}=\textit{Strong} | \text{Play}=\textit{No}) = 3/5$

$P(\text{Play}=\textit{No}) = 5/14$

- MAP rule

$P(\text{Yes} | \mathbf{x}')$: $[P(\textit{Sunny} | \textit{Yes})P(\textit{Cool} | \textit{Yes})P(\textit{High} | \textit{Yes})P(\textit{Strong} | \textit{Yes})]P(\text{Play}=\textit{Yes}) = 0.0053$

$P(\text{No} | \mathbf{x}')$: $[P(\textit{Sunny} | \textit{No})P(\textit{Cool} | \textit{No})P(\textit{High} | \textit{No})P(\textit{Strong} | \textit{No})]P(\text{Play}=\textit{No}) = 0.0206$

WHY ? Naïve Bayes Assumption

- $P(c_j)$
 - Can be estimated from the frequency of classes in the training examples.
- $P(x_1, x_2, \dots, x_p | c_j)$
 - $O(|X_1| \cdot |X_2| \cdot |X_3| \dots |X_p| \cdot |C|)$ parameters
 - Could only be estimated if a very, very large number of training examples was available.



If no naïve assumption

WHY ? Naïve Bayes Assumption

- $P(c_j)$
 - Can be estimated from the frequency of classes in the training examples.
- $P(x_1, x_2, \dots, x_p | c_j)$
 - $O(|X_1| \cdot |X_2| \cdot |X_3| \dots |X_p| \cdot |C|)$ parameters
 - Could only be estimated if a very, very large number of training examples was available.
- $P(x_k | c_j)$
 - $O([|X_1| + |X_2| + |X_3| \dots + |X_p|] \cdot |C|)$ parameters
 - Assume that the probability of observing the conjunction of attributes is equal to the product of the individual probabilities $P(x_i | c_j)$.

Not
Naïve

Naïve

WHY ? Naïve Bayes Assumption

- $P(c_j)$
 - Can be estimated from the frequency of classes in the training examples.
- $P(x_1, x_2, \dots, x_p | c_j)$
 - $O(|X_1| \cdot |X_2| \cdot |X_3| \dots |X_p| \cdot |C|)$ parameters
 - Could only be estimated if a very, very large number of training examples was available.
- $P(x_k | c_j)$
 - $O([|X_1| + |X_2| + |X_3| \dots + |X_p|] \cdot |C|)$ parameters
 - Assume that the probability of observing the conjunction of attributes is equal to the product of the individual probabilities $P(x_i | c_j)$.

Assuming $|C| = L$
num of unique values

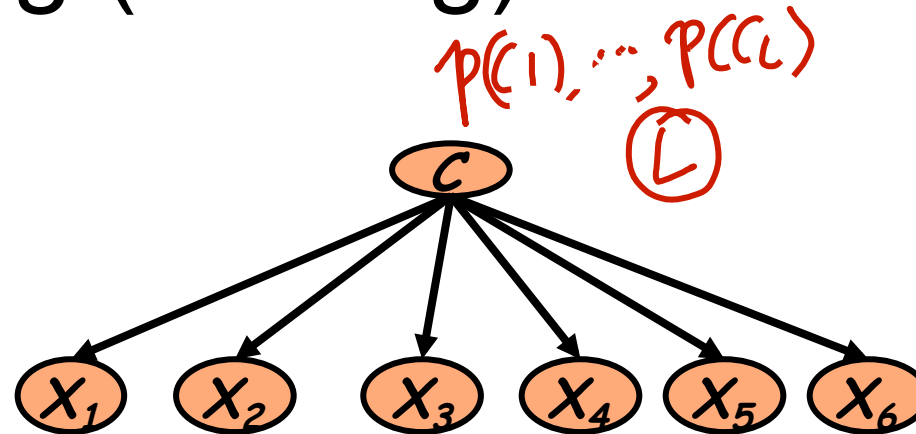
Assuming $|X_i| = 2, i=1, 2, \dots, p$
 $\Rightarrow 2^p \times L$ (Exp)

$(2+2+2+\dots+2) \times L$
 $= 2 \times p \times L$ (Linear)

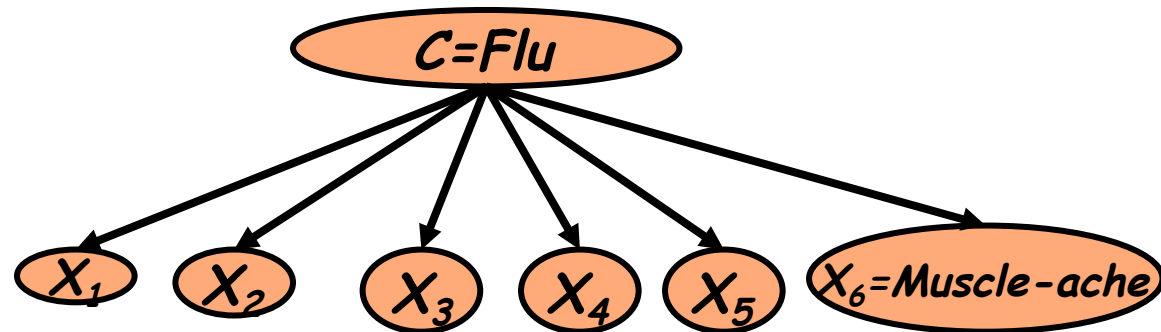
Not
Naïve

Naïve

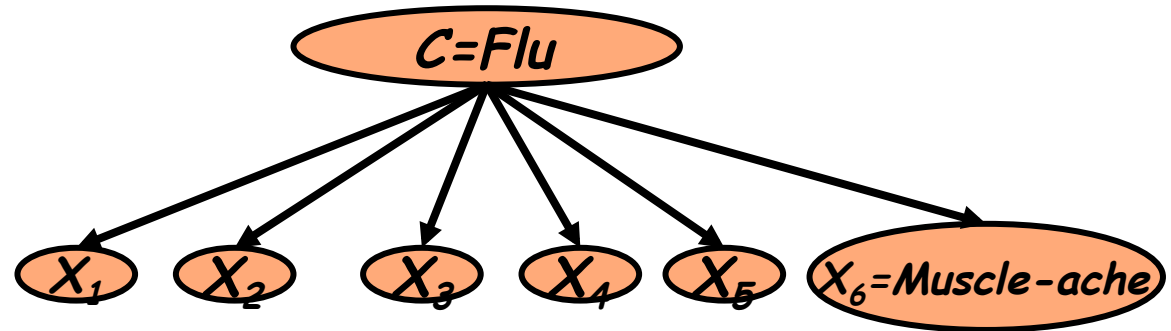
Learning (training) the NBC Model



For instance:



For instance:



- What if we have seen no training cases where patient had no flu and muscle aches?
- Zero probabilities cannot be conditioned away, no matter the other evidence!

$$\hat{P}(X_6 = t | C = \text{not_flu}) = \frac{N(X_6 = t, C = \text{nf})}{N(C = \text{nf})} = 0$$

muscle-ache-yes/no flu/nf

$$?? = \arg \max_c \hat{P}(c) \prod_i \hat{P}(x_i | c)$$

$$\delta_f = p(c=f|u) p(x_1|f) p(x_2|f) p(x_3|f) p(x_4|f) p(x_5|f) p(x_6|f)$$

$$\delta_{nf} = p(c=nf) p(x_1|nf) p(x_2|nf) p(x_3|nf) p(x_4|nf) p(x_5|nf) p(x_6|nf)$$

if any term gives 0,

$$\Rightarrow \delta_{nf} = 0$$

no matter other terms' value

Smoothing to Avoid Overfitting

Why necessary ??

$$\hat{P}(x_i | c_j) = \frac{N(X_i = x_i, C = c_j) + 1}{N(C = c_j) + k_i}$$

of values of feature X_i

To make
 $\sum_i (P(x_i | C_j)) = 1$

$$|X_i| = k_i$$

Smoothing to Avoid Overfitting

$$\hat{P}(x_i | c_j) = \frac{N(X_i = x_i, C = c_j) + 1}{N(C = c_j) + k_i}$$

of values of X_i

- Somewhat more subtle version

overall fraction in data
where $X_i = x_{i,k}$

$$\hat{P}(x_{i,k} | c_j) = \frac{N(X_i = x_{i,k}, C = c_j) + mp_{i,k}}{N(C = c_j) + m}$$

$\rightarrow k \in \{1, 2, \dots, k_i\}$

extent of
“smoothing”

Summary:

Generative Bayes Classifier with the MAP rule

Task: Classify a new instance X based on a tuple of attribute values $X = \langle X_1, X_2, \dots, X_p \rangle$ into one of the classes

$$c_{MAP} = \operatorname{argmax}_{c_j \in C} P(c_j | x_1, x_2, \dots, x_p)$$

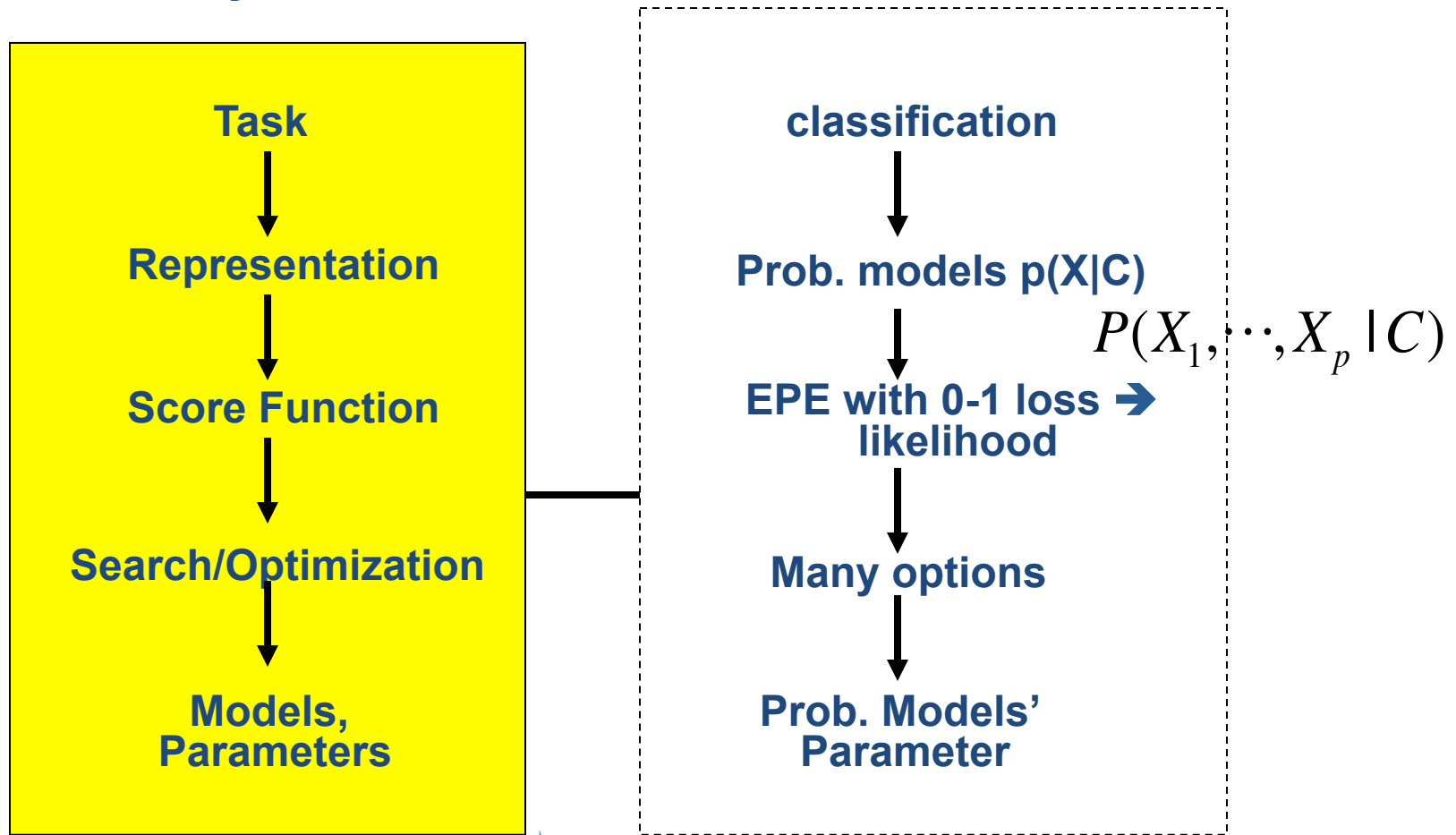
$$= \operatorname{argmax}_{c_j \in C} \frac{P(x_1, x_2, \dots, x_p | c_j) P(c_j)}{P(x_1, x_2, \dots, x_p)}$$

$$= \operatorname{argmax}_{\substack{c_j \in C \\ j=1,2,\dots,L}} \underbrace{P(x_1, x_2, \dots, x_p | c_j)} \underbrace{P(c_j)}$$

MAP = Maximum A Posteriori

$$\underset{k}{\operatorname{argmax}} P(C_k | X) = \underset{k}{\operatorname{argmax}} P(X, C) = \underset{k}{\operatorname{argmax}} P(X | C) P(C)$$

Generative Bayes Classifier



Bernoulli Naïve $p(W_i = \text{true} | c_k) = p_{i,k}$

Gaussian Naïve

Multinomial

$$\hat{P}(X_j | C = c_k) = \frac{1}{\sqrt{2\pi}\sigma_{jk}} \exp\left(-\frac{(X_j - \mu_{jk})^2}{2\sigma_{jk}^2}\right)$$

$$P(W_1 = n_1, \dots, W_v = n_v | c_k) = \frac{N!}{n_{1k}! n_{2k}! \dots n_{vk}!} \theta_{1k}^{n_{1k}} \theta_{2k}^{n_{2k}} \dots \theta_{vk}^{n_{vk}}$$

References

- Prof. Andrew Moore's review tutorial
- Prof. Ke Chen NB slides
- Prof. Carlos Guestrin recitation slides