



Machine Learning

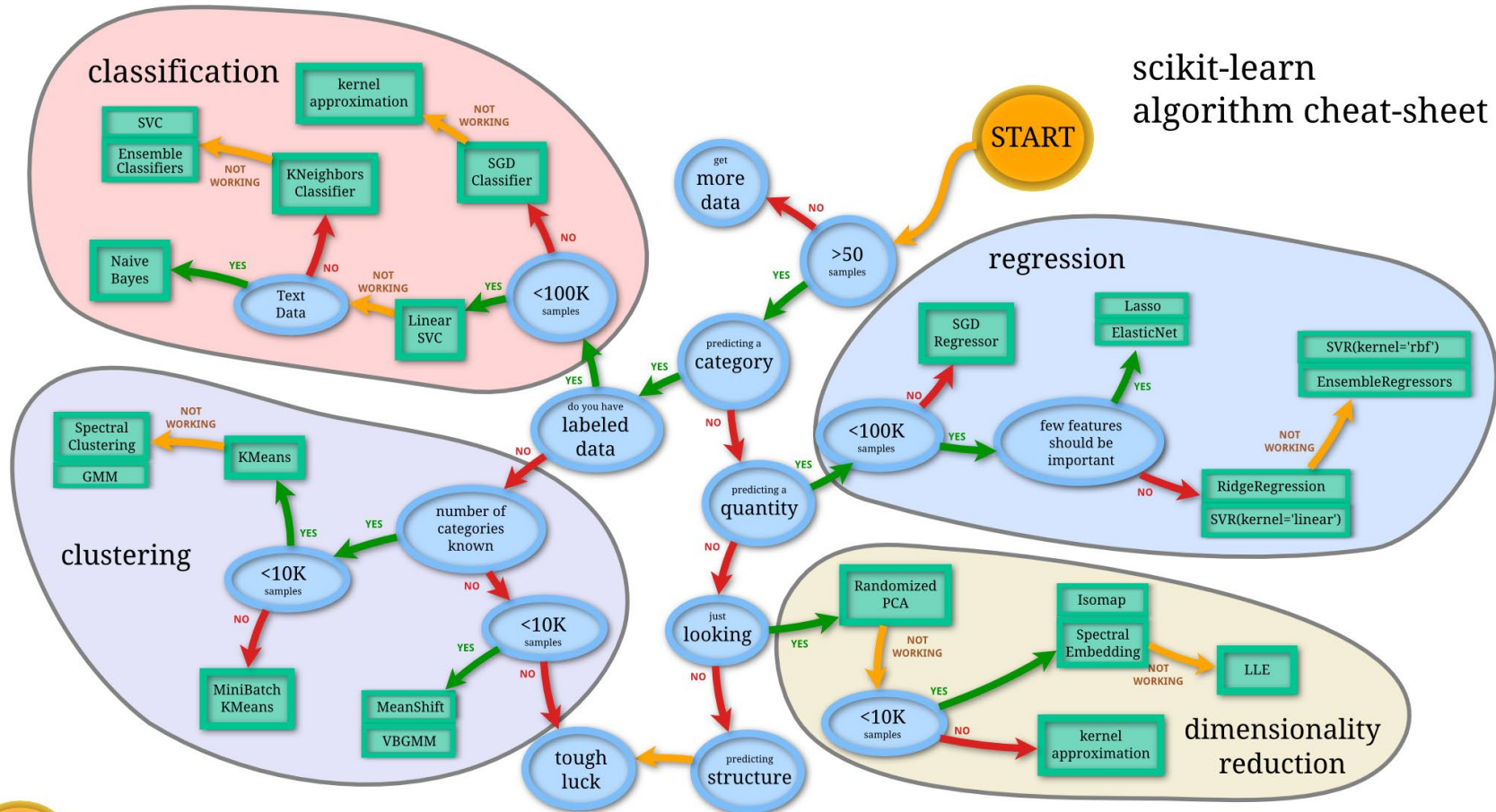
Lecture 17a: Generative Bayes Classifiers

Dr. Beilun Wang

Southeast University
School of Computer Science
and Engineering

Roadmap

scikit-learn
algorithm cheat-sheet



Course Content Plan

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

- ❑ Graphical models

- ❑ Reinforcement Learning

Y is a continuous

Y is a discrete

NO Y

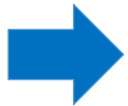
About $f()$

About interactions among X_1, \dots, X_p

Learn program to Interact with its environment

Three major sections for classification

- We can divide the large variety of classification approaches into **roughly three major types**
 - Discriminative
 - directly estimate a decision rule/boundary
 - e.g., support vector machine, decision tree, logistic regression,
 - e.g. neural networks (NN), deep NN
 - Generative:
 - build a generative statistical model
 - e.g., Bayesian networks, **Naïve Bayes classifier**
 - Instance based classifiers
 - Use observation directly (no models)
 - e.g. K nearest neighbors



Today: Generative Bayes Classifiers

- ➔ • Bayes Classifier (BC)
 - Generative Bayes Classifier
- Naïve Bayes Classifier
- Gaussian Bayes Classifiers
 - Gaussian distribution
 - Naïve Gaussian BC
 - Not-naïve Gaussian BC ➔ LDA, QDA

Review: Bayes classifiers (BC)

- Treat each feature attribute and the class label as random variables.
- Testing: 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 $p(c|x_1, x_2, \dots, x_p)$
- Training: can we estimate $p(c_i|\mathbf{x}) = p(c_i|x_1, x_2, \dots, x_p)$ directly from data?

$$C_{MAP} = \operatorname{argmax}_{c_i \in C} p(c_i|x_1, x_2, \dots, x_p)$$

MAP Rule

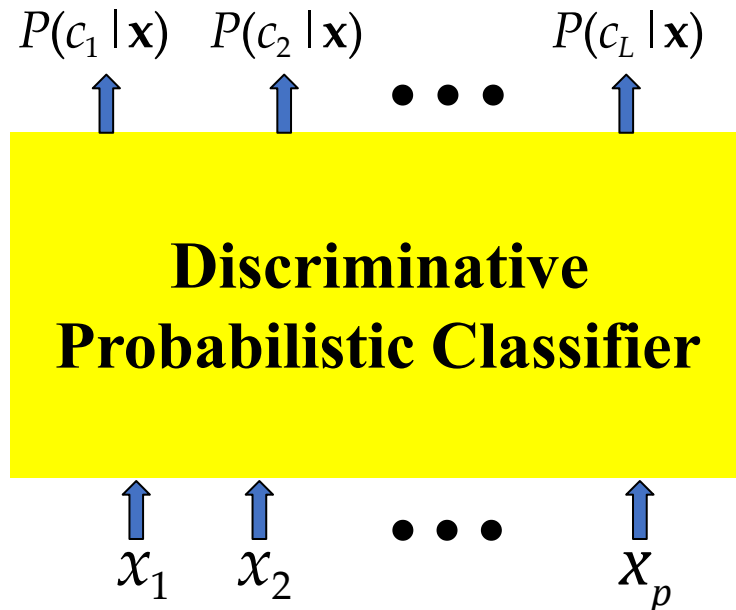
Review: Two kinds of Bayes classifiers via MAP classification rule



- Establishing a probabilistic model for classification
 - Discriminative
 - Generative

Review: Discriminative BC

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



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

Review: Generative BC

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

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

$$P(\mathbf{x} | c_1)$$



**Generative
Probabilistic Model
for Class 1**

x_1

x_2

...

x_p

$$P(\mathbf{x} | c_2)$$



**Generative
Probabilistic Model
for Class 2**

x_1

x_2

...

x_p

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

...

$$P(\mathbf{x} | c_L)$$



**Generative
Probabilistic Model
for Class L**

x_1

x_2

...

x_p

Review: Bayes Rule for Generative BC

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

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

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

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

Summary of Generative BC

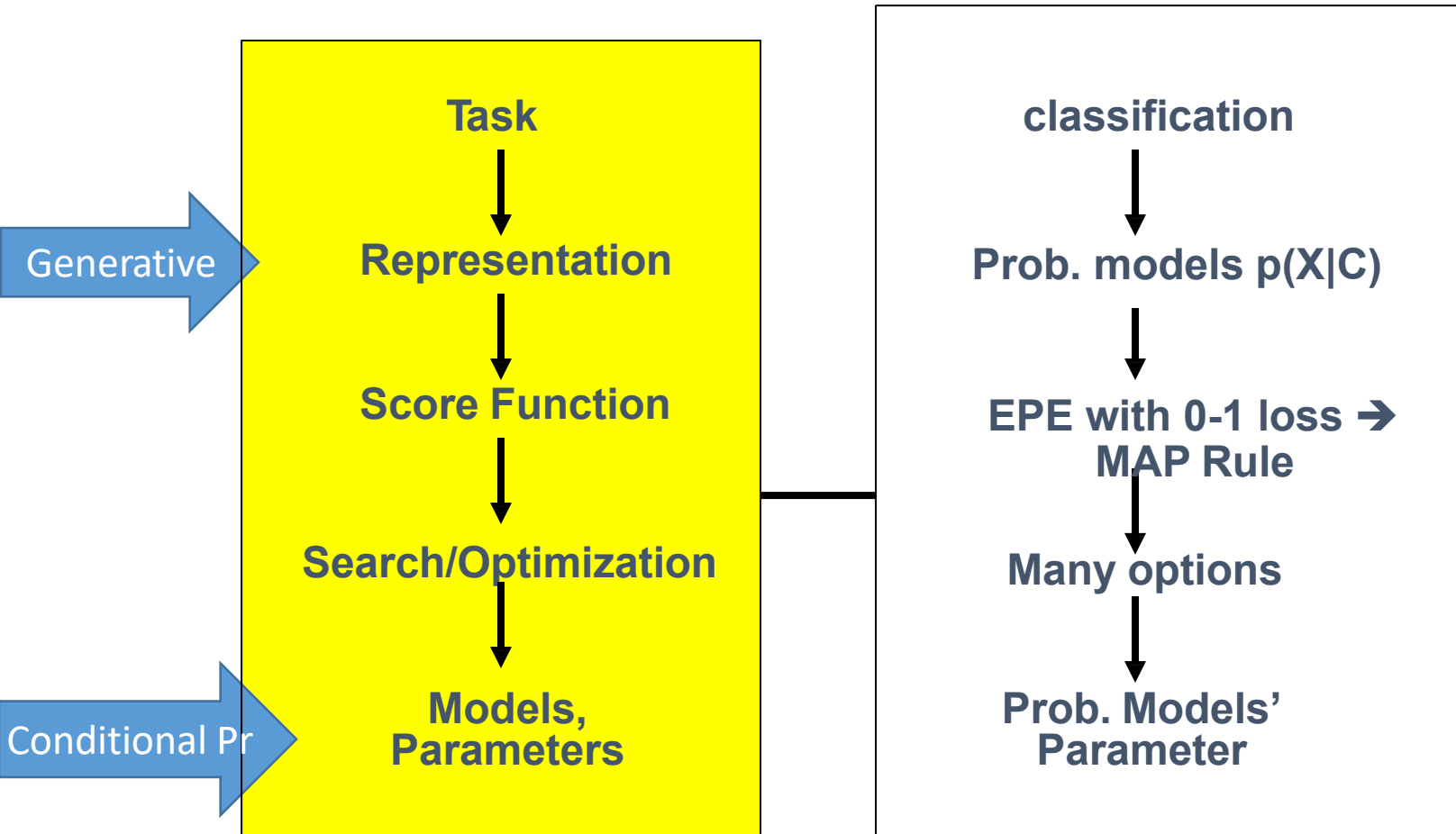
- Apply Bayes rule to get posterior probabilities

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

for $i = 1, 2, \dots, L$

- Then apply the MAP rule

Generative Bayes Classifier



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
 X_4
 C


X_1		X_2		X_3		X_4		C	
S (Sunny) O R		H (Hot) M C		H (High) N		W (Weak) S		$C_1 = Yes$ $C_2 = No$	
3	×	3	×	2	×	2	×	2	—

→ $P(X_1, X_2, X_3, X_4 | C), P(C)$

$$\textcircled{1} \begin{cases} P(C = Yes) = \frac{N(Yes)}{N(train)} = \frac{9}{14} \\ P(C = No) = 1 - P(C = Yes) = \frac{5}{14} \end{cases}$$

$$\textcircled{2} \begin{cases} P(X_1 = S, X_2 = H, X_3 = H, X_4 = W | C = No) = \frac{1}{5} \\ P(X_1 = S, X_2 = H, X_3 = H, X_4 = W | C = Yes) = 0 \end{cases}$$

Learning: maximum likelihood estimates

- simply use the frequencies in the data

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. $P(O, H, H, W | Yes) = \frac{1}{9}$

Directly
estimate from
data

e.g. $P(O, H, H, W | No) = \frac{0}{5} = 0$

Generative BC: Learning Phase

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

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

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

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

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

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

Generative BC: Testing 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



$$\hat{P}(a'_1, \dots, a'_p | c^*) \hat{P}(c^*) > \hat{P}(a'_1, \dots, a'_p | c) \hat{P}(c),$$
$$c \neq c^*, c = c_1, \dots, c_L$$

- Given a new instance:

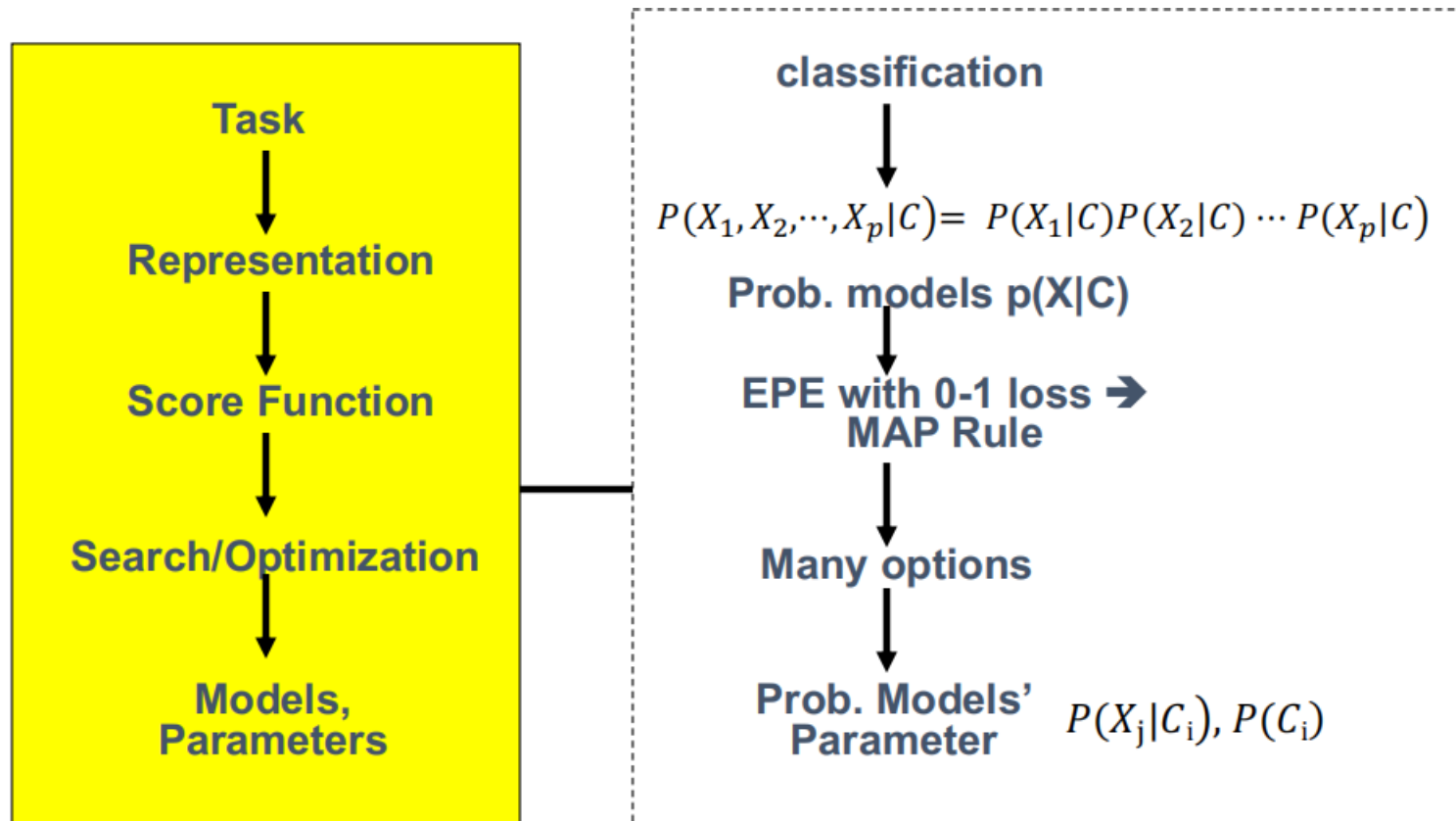
$x' = (\text{Outlook} = \text{Sunny}, \text{Temperature} = \text{Cool}, \text{Humidity} = \text{High}, \text{Wind} = \text{Strong})$

$$\begin{cases} P(x'|Yes) P(C = Yes) \\ P(x'|No) P(C = No) \end{cases} \Rightarrow \underset{C}{argmax} \Rightarrow \text{predicted } C^*$$

Today: Generative Bayes Classifiers

- Bayes Classifier (BC)
 - Generative Bayes Classifier
-  • Naïve Bayes Classifier
- Gaussian Bayes Classifiers
 - Gaussian distribution
 - Naïve Gaussian BC
 - Not-naïve Gaussian BC  LDA, QDA

Naïve Bayes Classifier



Naïve Bayes Classifier

- Bayes classification

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

- Difficulty: learning the joint probability
- Naïve Bayes classification
 - Assume that **all input attributes are conditionally independent**

Naïve Bayes Classifier

- Naïve Bayes classification
 - Assume 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 (for discrete input attributes) – Training / Learning

- 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

Naïve Bayes Classifier (for discrete input attributes) – Testing

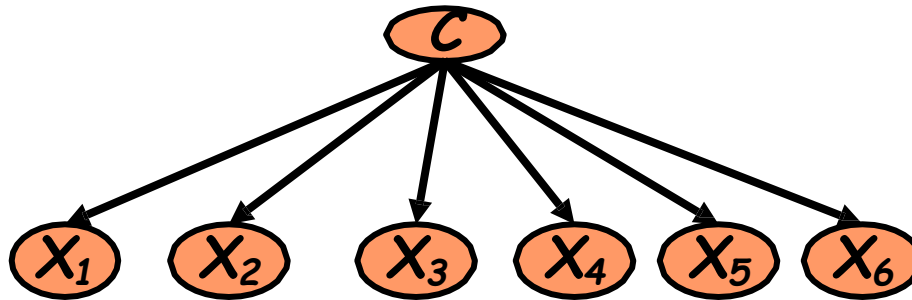
Test Phase: Given an unknown instance

Look up tables to assign the label c^* to \mathbf{X}' if $\mathbf{X}' = (a'_1, \dots, a'_p)$

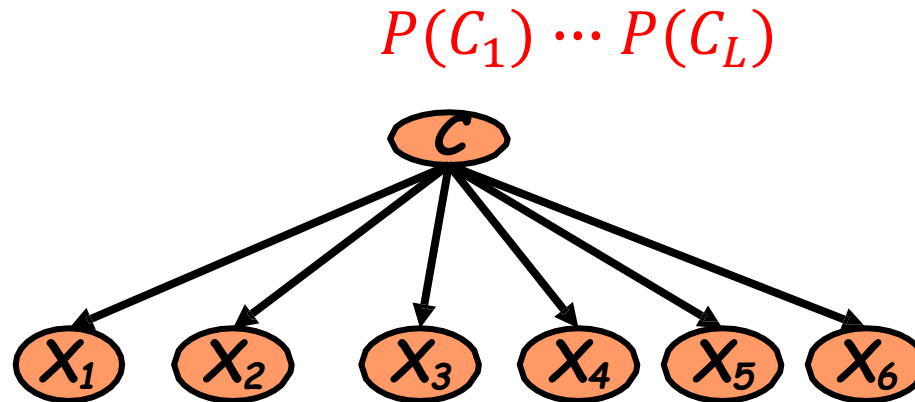
$$[\hat{P}(a'_1 | c^*) \cdots \hat{P}(a'_p | c^*)] \hat{P}(c^*) > [\hat{P}(a'_1 | c) \cdots \hat{P}(a'_p | c)] \hat{P}(c), \\ c \neq c^*, c = c_1, \dots, c_L$$

Learning (training) the NBC Model

$$P(C_1) \cdots P(C_L)$$



Learning (training) the NBC Model



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

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

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

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

$$P(X_1 = \text{Rain} | C = \text{Yes}) = \frac{3}{9}$$

$$P(X_1 = \text{Rain} | C = \text{No}) = \frac{2}{5}$$

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

- Learning Phase

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

Outlook	Play=Yes	Play=No
Sunny		
Overcast		
Rain		

Temperature	Play=Yes	Play=No
Hot		
Mild		
Cool		

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

Humidity	Play=Yes	Play=No
High		
Normal		

Wind	Play=Yes	Play=No
Strong		
Weak		

$P(\text{Play=Yes}) = ??$

$P(\text{Play=No}) = ??$

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

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

- Learning Phase

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

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

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

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

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

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

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

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

Example: Play tennis

- Test Phase

Given a new instance,

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

Example: Play tennis

- 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

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

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

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

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

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

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

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

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

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

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

Testing the NBC Model

$P(\text{Outlook}=\text{Sunny} | \text{Play}=\text{Yes}) = 2/9$
 $P(\text{Temperature}=\text{Cool} | \text{Play}=\text{Yes}) = 3/9$
 $P(\text{Humidity}=\text{High} | \text{Play}=\text{Yes}) = 3/9$
 $P(\text{Wind}=\text{Strong} | \text{Play}=\text{Yes}) = 3/9$
 $P(\text{Play}=\text{Yes}) = 9/14$

$P(\text{Outlook}=\text{Sunny} | \text{Play}=\text{No}) = 3/5$
 $P(\text{Temperature}=\text{Cool} | \text{Play}=\text{No}) = 1/5$
 $P(\text{Humidity}=\text{High} | \text{Play}=\text{No}) = 4/5$
 $P(\text{Wind}=\text{Strong} | \text{Play}=\text{No}) = 3/5$
 $P(\text{Play}=\text{No}) = 5/14$

– MAP rule

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

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



Given the fact $P(\text{Yes} | \mathbf{x}') < P(\text{No} | \mathbf{x}')$, we label \mathbf{x}' to be “No”.

Why Naïve Bayes Assumption

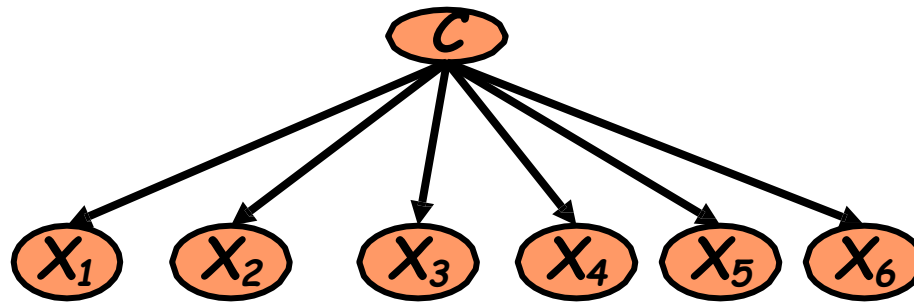
- $P(c_j)$
 - Can be estimated from the frequency of classes in the training examples.
- $P(x_1, x_2, \dots, x_p | c_j)$
 - $O(|X_1| \cdot |X_2| \cdot |X_3| \dots |X_p| \cdot |C|)$ parameters
 - Could only be estimated if a very, very large number of training examples was available.
- $P(x_k | c_j)$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

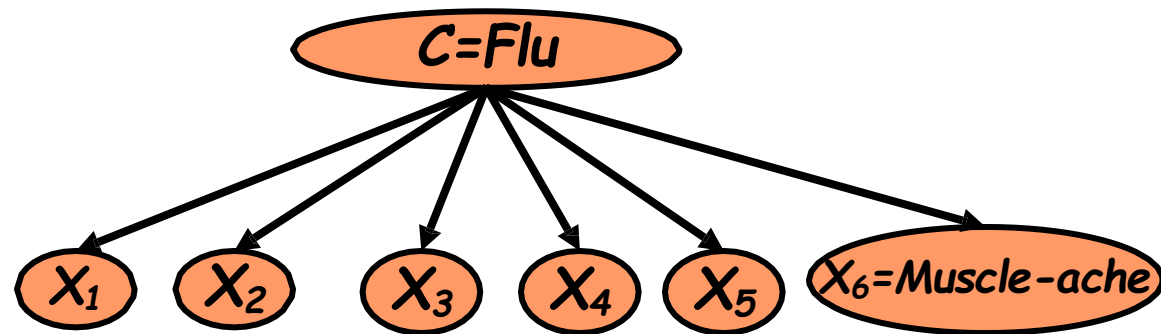
Naïve

Challenges during learning the NBC Model

$$P(C_1) \cdots P(C_L)$$



For instance:



Challenges during learning the NBC Model

- For instance:
 - What if we have seen no training cases where patient had no flu and muscle aches?

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

- Zero probabilities cannot be conditioned away, no matter the other evidence!

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

Smoothing to Avoid Overfitting

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

number of values of feature X_i

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

Smoothing to Avoid Overfitting

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

number of values of X_i

- Somewhat more subtle version

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

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

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

extent of
"smoothing"



Summary: Generative Bayes Classifier

- 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

$$\begin{aligned} c_{MAP} &= \operatorname{argmax}_{c_j \in C} P(c_j \mid x_1, x_2, \dots, x_p) \\ &= \operatorname{argmax}_{c_j \in C} \frac{P(x_1, x_2, \dots, x_p \mid c_j) P(c_j)}{P(x_1, x_2, \dots, x_p)} \\ &= \operatorname{argmax}_{c_j \in C} P(x_1, x_2, \dots, x_p \mid c_j) P(c_j) \end{aligned}$$

MAP = Maximum A Posteriori

Next: Generative Bayes Classifiers

- Bayes Classifier (BC)
 - Generative Bayes Classifier
- Naïve Bayes Classifier
-  • Gaussian Bayes Classifiers
 - Gaussian distribution
 - Naïve Gaussian BC
 - Not-naïve Gaussian BC  LDA, QDA

References

- <https://qiyanjun.github.io/2019f-UVA-CS6316-MachineLearning/>
- Prof. Andrew Moore's review tutorial
- Prof. Ke Chen NB slides
- Prof. Carlos Guestrin recitation slides



Thanks for listening