

UVA CS 6316/4501

– Fall 2016

Machine Learning

Lecture 12: Generative Bayes Classifiers

Dr. Yanjun Qi

University of Virginia

Department of
Computer Science

Where are we ? →

Five major sections of this course

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

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

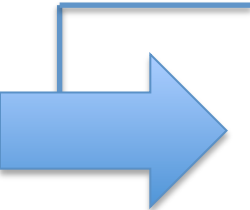
3. Instance based classifiers

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

Last Lecture Recap: Probability Review

- The big picture : data \leftrightarrow probabilistic model
- Sample space, Events and Event spaces
- Random variables
- Joint probability, Marginal probability, conditional probability,
- Chain rule, Bayes Rule, Law of total probability, etc.
- Structural properties
 - Independence, conditional independence

Today : Generative Bayes Classifiers

- 
- ✓ Bayes Classifier
 - MAP classification rule
 - Generative Bayes Classifier
 - ✓ Naïve Bayes Classifier
 - ✓ Gaussian Bayes Classifiers
 - Gaussian distribution
 - Gaussian NBC
 - LDA, QDA

X_1	X_2	X_3	C

A Dataset for classification

$$f : X \longrightarrow C$$

Output as Discrete Class Label
 C_1, C_2, \dots, C_L

$P(C | X)$

(X_1, X_2, \dots, X_p)

- **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]

Bayes classifiers

- Treat each feature attribute and the class label as random variables.

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 value of C_i that maximizes $p(C_i | x_1, x_2, \dots, x_p)$.

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 . $\rightarrow \{C_1, C_2, \dots, C_L\}$
 - Specifically, we want to find the value of C_i that maximizes $p(C_i | x_1, x_2, \dots, x_p)$.
- Can we estimate $p(C_i | \mathbf{x}) = p(C_i | x_1, x_2, \dots, x_p)$ directly from data?

Bayes classifiers

→ MAP classification rule

- Establishing a probabilistic model for classification
- **MAP** classification rule
 - **MAP**: **M**aximum **A** **P**osterior

Bayes classifiers

→ MAP classification rule

- Establishing a probabilistic model for classification

→ **MAP** classification rule

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

$$c^* = \operatorname{argmax}_{c_i \in \{c_1, c_2, \dots, c_L\}} P(c_i | x)$$

$$P(c^* | x)$$

$$P(C = c^* | \mathbf{X} = \mathbf{x}) > P(C = c | \mathbf{X} = \mathbf{x})$$

$$\text{for } c \neq c^*, c = c_1, \dots, c_L$$

Bayes classifiers

→ MAP classification rule

- Establishing a probabilistic model for classification
- **MAP** classification rule
 - **MAP**: **M**aximum **A** **P**osterior
 - Assign x to c^* if

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

$$\left. \begin{array}{l} P(C = c_1 | \mathbf{x}) \\ P(C = c_2 | \mathbf{x}) \\ P(C = c_3 | \mathbf{x}) \end{array} \right\} \max \Rightarrow c_i$$

Bayes classifiers

→ MAP classification rule

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

(1) Discriminative

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

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

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

**Discriminative
Probabilistic Classifier**

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

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

*logistic
regression*

(2) Generative

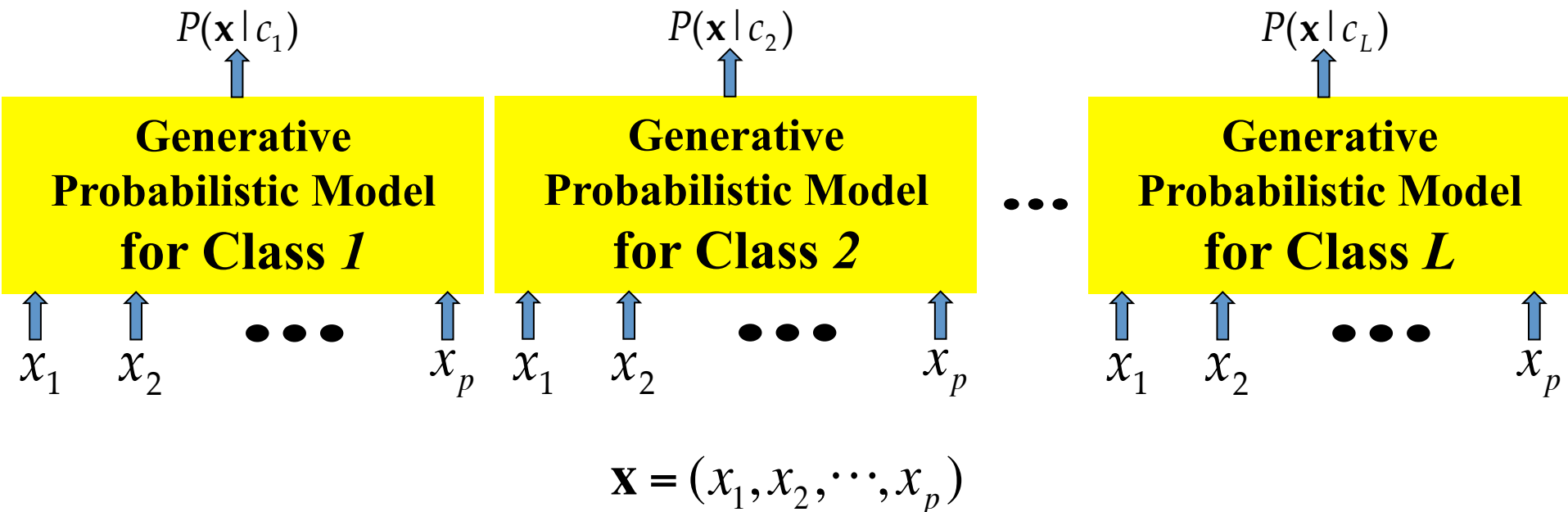
$$P(x|c_1), P(x|c_2), \dots, P(x|c_L)$$

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

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

$$\Rightarrow P(C|X)$$

MAP rule



(2) Generative BayesRule

generative model

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

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

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

$$\phi(X=x_i | C=c_i)$$

$$(x_1, x_2, \dots, x_p)$$

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

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

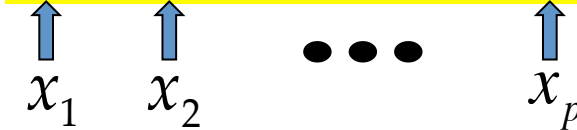
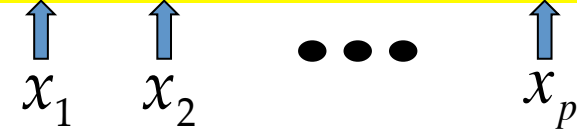
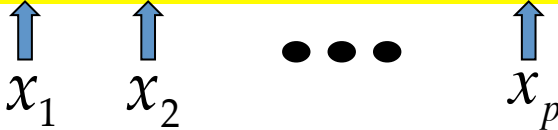
$$P(\mathbf{x}|c_L)$$

Generative Probabilistic Model for Class 1

Generative Probabilistic Model for Class 2

...

Generative Probabilistic Model for Class L



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

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 : Bayes' Rule

– for Generative Bayes Classifiers

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

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

max $\Rightarrow C^$*

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

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

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

Review : Bayes' Rule

– for Generative Bayes Classifiers

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

[MAP rule]

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

Posterior

Prior

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

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

max $\Rightarrow C^*$

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

Summary:

Generative classification with the MAP rule

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

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

- Generative classification with the MAP rule

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

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

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

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

- Then apply the MAP rule

$$\begin{aligned} & \arg \max_{i=1,2,\dots,L} \left\{ p(c_i | \mathbf{x}) = \frac{P(\mathbf{x} | c_i) P(c_i)}{P(\mathbf{x})} \right\} \\ & = \arg \max_{i=1,2,\dots,L} \left\{ P(\mathbf{x} | c_i) P(c_i) \right\} \end{aligned}$$

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

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

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

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

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

MAP = Maximum A Posteriori

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

\mathcal{X}_1 ↓ \mathcal{X}_2 ↑ \mathcal{X}_3 ↑ \mathcal{X}_4 ↑ C
PlayTennis: training examples

$k_2 = 3$

\mathcal{X}_2 :
 {Hot, Mild, Cool}

$\mathcal{X}_3 = \{ \text{High, Normal} \}$

$k_3 = 2$

$\mathcal{X}_4 = (W, S)$
 $k_4 =$

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

$\mathcal{X}_1: \{ \text{sunny, overc, rain} \}$
 $(k=3)$

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

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

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

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

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

$3 \times 3 \times 2 \times 2 \times 2 \Rightarrow 72$

parameter from train

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

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

Example

- Example: Play Tennis $\rightarrow 36 \times 2$

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

$P(C=C_i)$

{Yes}
{No}

$P(C=Yes)$
= 9/14

$P(C=No)$
= 5/14

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

- maximum likelihood estimates (explain later)
 - 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.
 $p(\text{overcast, hot, high, weak} \mid \text{No}) = \frac{0}{9}$

Generative Bayes Classifier:

- Learning Phase

	X_1	X_2	X_3	C
S_1				
S_2				
S_3				
S_4				
S_5				
S_6				

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

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

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

a look up table of cond. prob

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

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


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

- Given a new instance,

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

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

Today : Generative Bayes Classifiers

- ✓ Bayes Classifier
 - MAP classification rule
 - Generative Bayes Classifier
-  ✓ Naïve Bayes Classifier
- ✓ Gaussian Bayes Classifiers
 - Gaussian distribution
 - Gaussian NBC
 - LDA, QDA

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!**
given C variable

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

$$P(x, y | z) = P(x | z) P(y | z)$$

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

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

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

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

Naïve Bayes Classifier

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

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

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

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

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

$$\left\{ \begin{array}{l} j=1,2,\dots,p \\ i=1,2,\dots,L \end{array} \right\} P(X_j | c_i)$$

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 ;

⇒ L

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 ($c_i = c_1, \dots, c_L$)

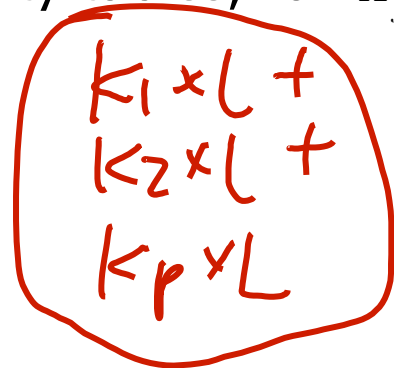
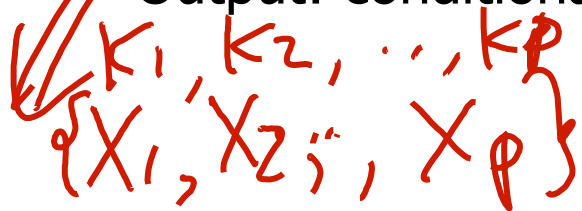
→ $\hat{P}(C = c_i) \leftarrow$ estimate $P(C = c_i)$ with examples in S ;



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



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

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

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

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

An Example

- Example: Play Tennis

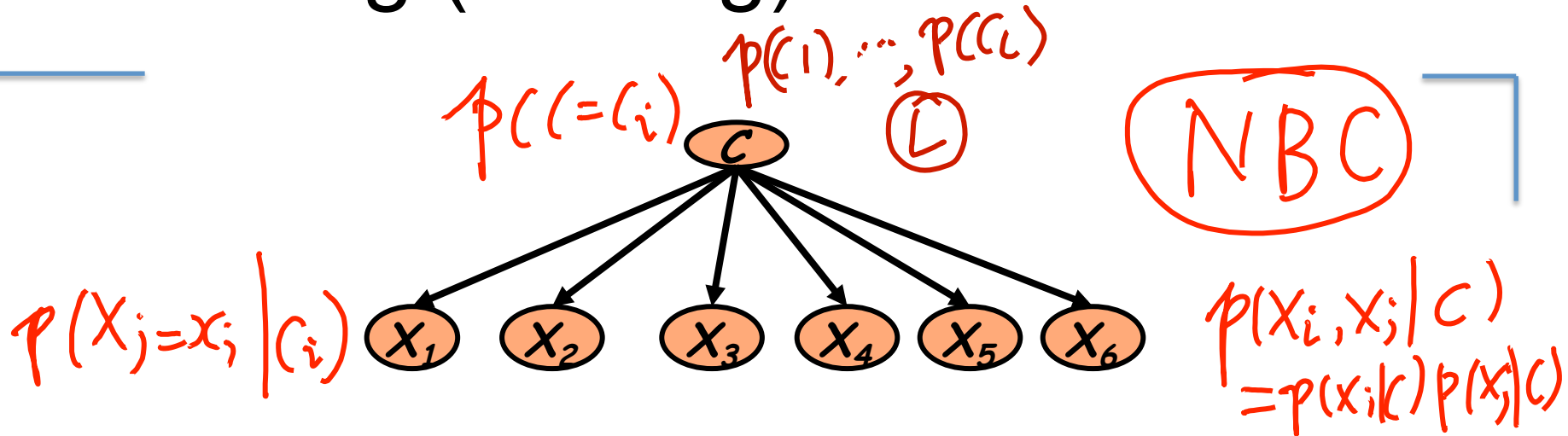
Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal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

X_1 X_2 X_3 X_4
 X_1 ↓ X_2 ↑ X_3 ↑ X_4 ↑
 PlayTennis: training examples C

$k_2 = 3$
 $X_2 = \{ \text{Hot, Mild, Cool} \}$
 $X_3 = \{ \text{High, Normal} \}$
 $k_3 = 2$
 $X_4 = (W, S)$
 $k_4 =$

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

Learning (training) the NBC Model



- maximum likelihood estimates (explain later)
 - simply use the frequencies in the data

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

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

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

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

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

$$P(x_1, x_2, x_3, x_4 | C_i)$$

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

Counting
↑

Learning Phase

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

X_1
3

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

2

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

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

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

$P(C_i)$

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

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,

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

$$\begin{aligned} \rightarrow & P(c_1) P(\textit{Sunny} | c_1) P(\textit{Cool} | c_1) P(\textit{High} | c_1) P(\textit{Strong} | c_1) \\ & = \frac{9}{14} \times \frac{2}{9} \cdots = \end{aligned}$$

$$\begin{aligned} \rightarrow & P(c_2) P(\textit{Su} | c_2) P(\textit{Co} | c_2) P(\textit{hi} | c_2) P(\textit{St} | c_2) \\ & = \frac{5}{14} \times \frac{3}{5} \times \cdots = \end{aligned}$$

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

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

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

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

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

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

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

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

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

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

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

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

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

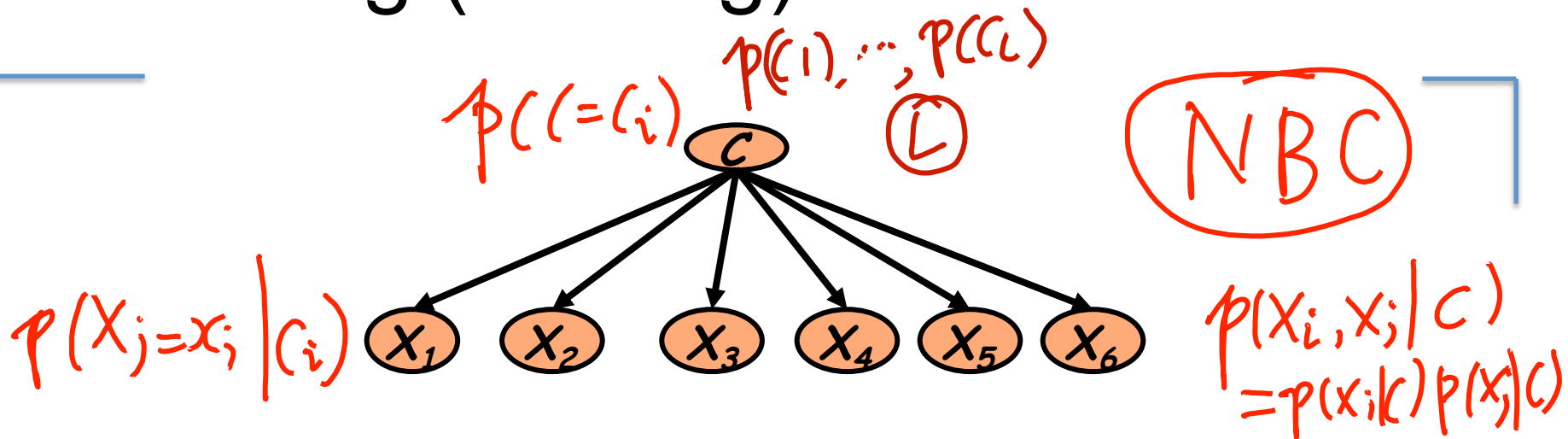
Not
Naïve

Naïve

DETOUR: Course Schedule

- Midterm @ WED / In CLASS / 70mins
- Open to your notes + (printed) lecture + Four HWs we had so far
 - Nothing else is allowed
 - Please turn off your phone at the beginning
 - No Electronic Devices (other than basic calculator)
- Final Exam
 - Will be close-note !

Learning (training) the NBC Model

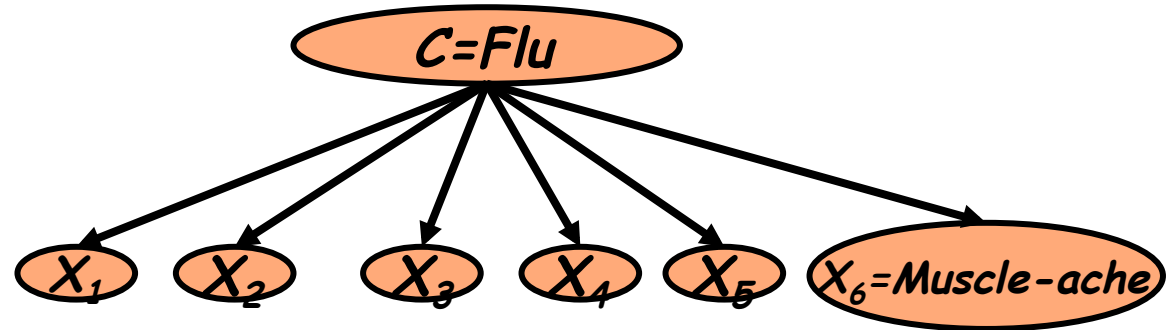


- maximum likelihood estimates (explain later)
 - simply use the **frequencies** in the data

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

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

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 $X_i = x_{i,k}$

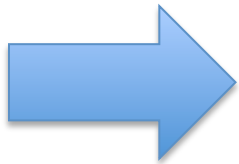
$$\hat{P}(x_{i,k} | c_j) = \frac{N(X_i = x_{i,k}, C = c_j) + mp_{i,k}}{N(C = c_j) + m}$$

→ $K \in \{1, 2, \dots, k_i\}$

extent of "smoothing"

Today : Generative Bayes Classifiers

- ✓ Bayes Classifier
 - MAP classification rule
 - Generative Bayes Classifier
- ✓ Naïve Bayes Classifier
- ✓ Gaussian Bayes Classifiers
 - Gaussian distribution
 - Gaussian NBC
 - LDA, QDA



Review: Continuous Random Variables

- Probability density function (pdf) instead of probability mass function (pmf)
 - For discrete RV: Probability mass function (pmf):
 $P(X = x_i)$
- A pdf (prob. Density func.) is any function $f(x)$ that describes the probability density in terms of the input variable x .

Review: Probability of Continuous RV

- Properties of pdf

- $f(x) \geq 0, \forall x$

- $\int_{-\infty}^{+\infty} f(x) = 1$

$$\longrightarrow \sum_{i=1}^{k_i} P(X = x_i) = 1$$

- Actual probability can be obtained by taking the integral of pdf

- E.g. the probability of X being between 5 and 6 is

$$P(5 \leq X \leq 6) = \int_5^6 f(x) dx$$

Review: Mean and Variance of RV

- Mean (Expectation): $\mu = E(X)$

- Discrete RVs:
$$E(X) = \sum_{v_i} v_i P(X = v_i)$$

$$E(g(X)) = \sum_{v_i} g(v_i) P(X = v_i)$$

- Continuous RVs:
$$E(X) = \int_{-\infty}^{+\infty} xf(x)dx$$

$$E(g(X)) = \int_{-\infty}^{+\infty} g(x)f(x)dx$$

Review: Mean and Variance of RV

- Variance: $Var(X) = E((X - \mu)^2)$ $\sigma_x = \sqrt{V(x)}$

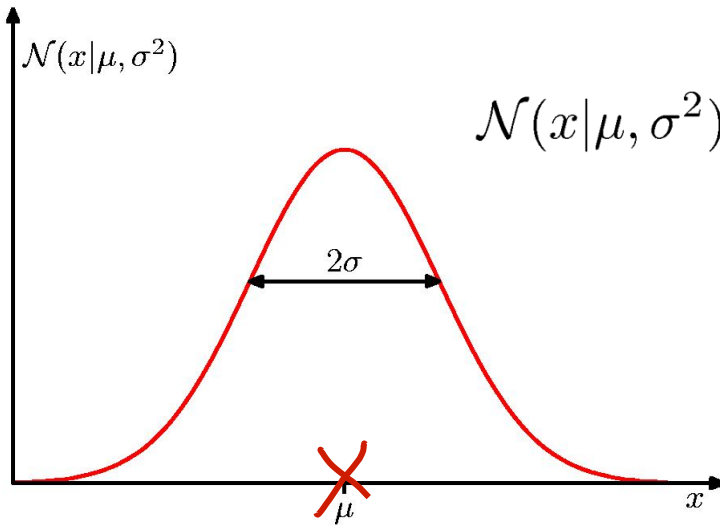
- Discrete RVs: $V(X) = \sum_{v_i} (v_i - \mu)^2 P(X = v_i)$

- Continuous RVs: $V(X) = \int_{-\infty}^{+\infty} (x - \mu)^2 f(x) dx$

- Covariance: Correlation
 $\rho_{X,Y} = Corr(X, Y) = \frac{Cov(X, Y)}{\sigma_x \sigma_y}$

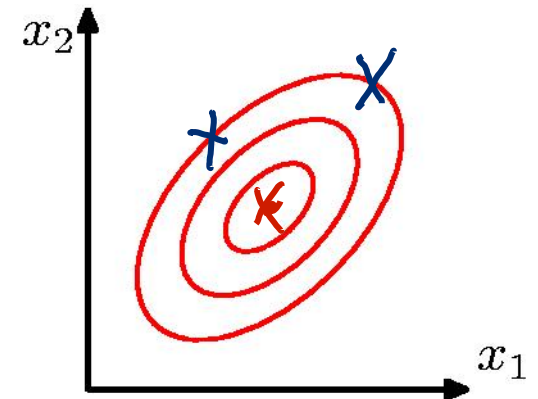
$$Cov(X, Y) = E((X - \mu_x)(Y - \mu_y)) = E(XY) - \mu_x \mu_y$$

Gaussian Distribution



$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\}$$

$$X \sim \mathcal{N}(\mu, \sigma^2)$$



$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{\mathcal{P}/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Mean
Covariance Matrix

Multivariate Normal (Gaussian) PDFs

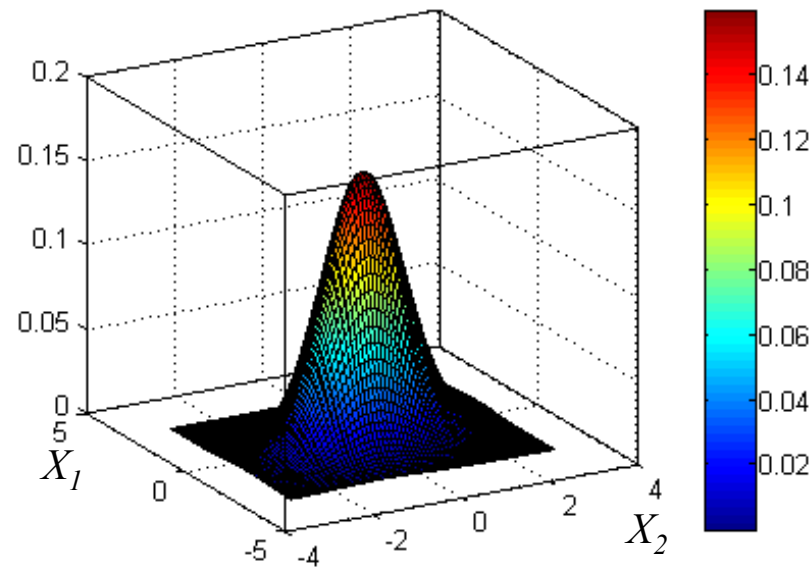
The only widely used continuous joint PDF is the multivariate normal (or Gaussian):

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{p/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Where $|\ast|$ represents **determinant**

Bivariate
normal PDF:

- Mean of normal PDF is at peak value. Contours of equal PDF form ellipses.



- The covariance matrix captures linear dependencies among the variables

Example: the Bivariate Normal distribution

$$f(x_1, x_2) = \frac{1}{(2\pi) |\Sigma|^{1/2}} e^{-\frac{1}{2}(\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})}$$

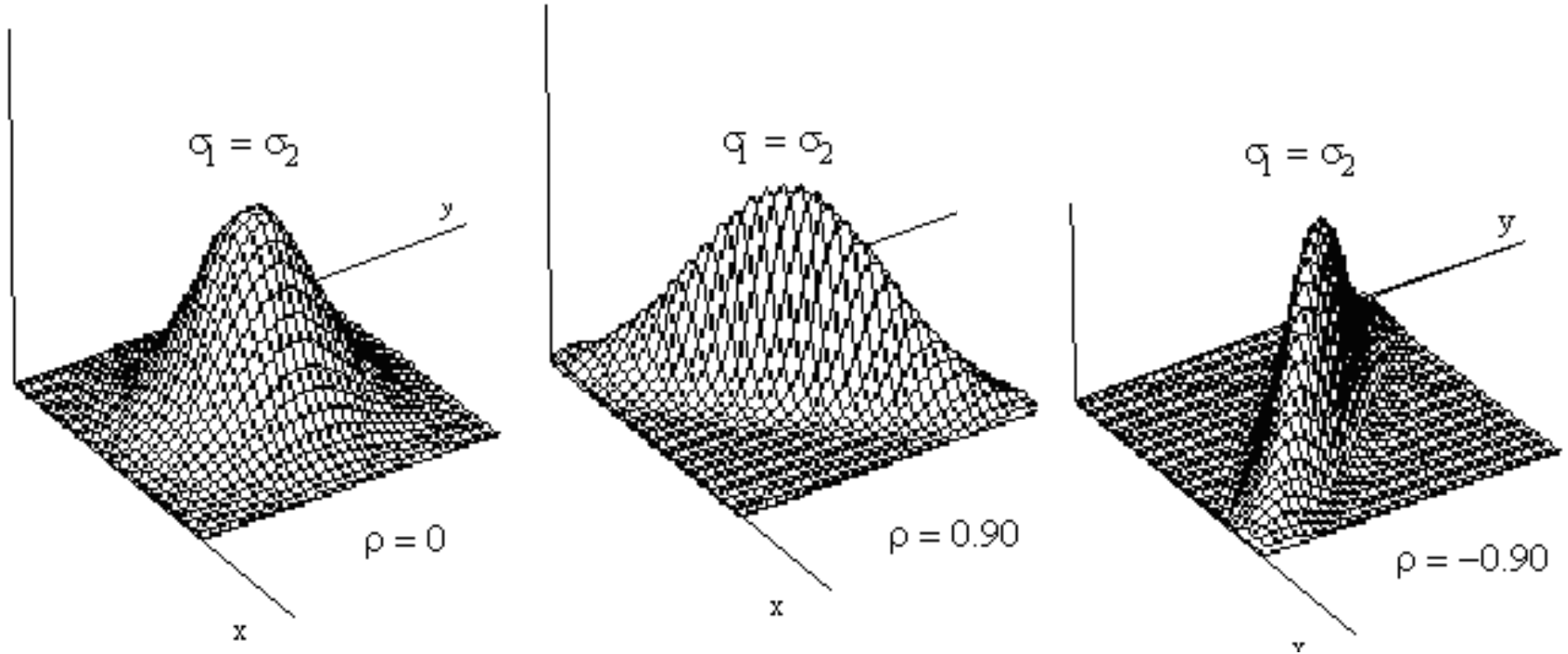
with $\vec{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$ and

$$\Sigma_{2 \times 2} = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{bmatrix}$$

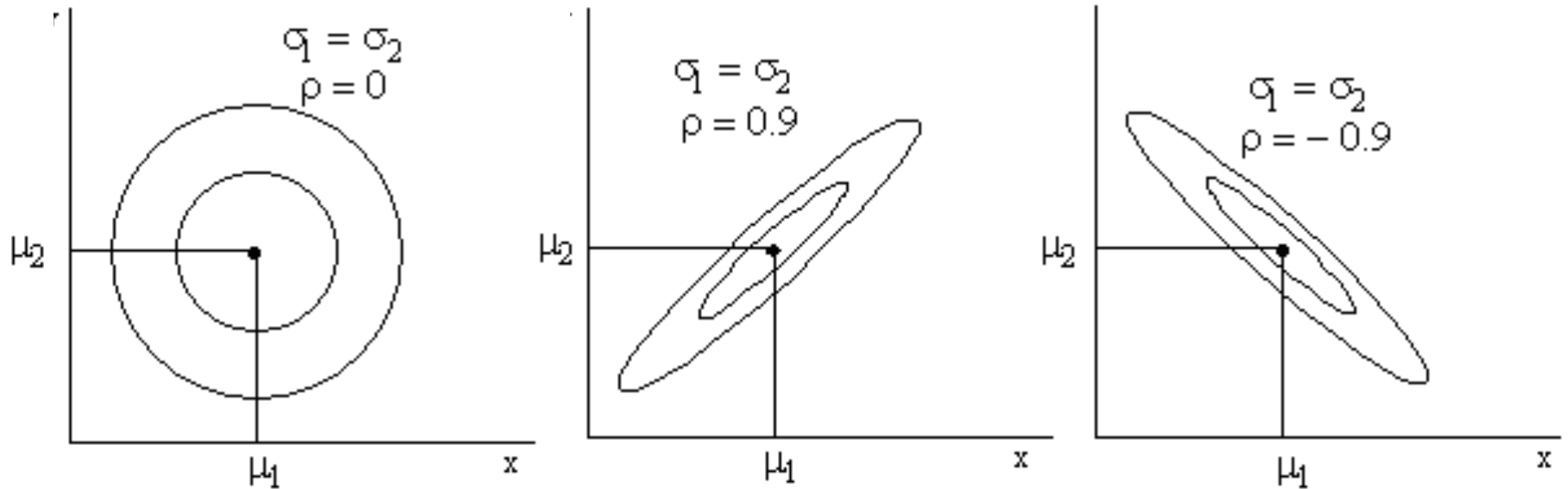
$V(X_1)$
 $Cov(X_1, X_2)$
 $V(X_2)$

$$|\Sigma| = \sigma_{11} \sigma_{22} - \sigma_{12}^2 = \sigma_1^2 \sigma_2^2 (1 - \rho^2)$$

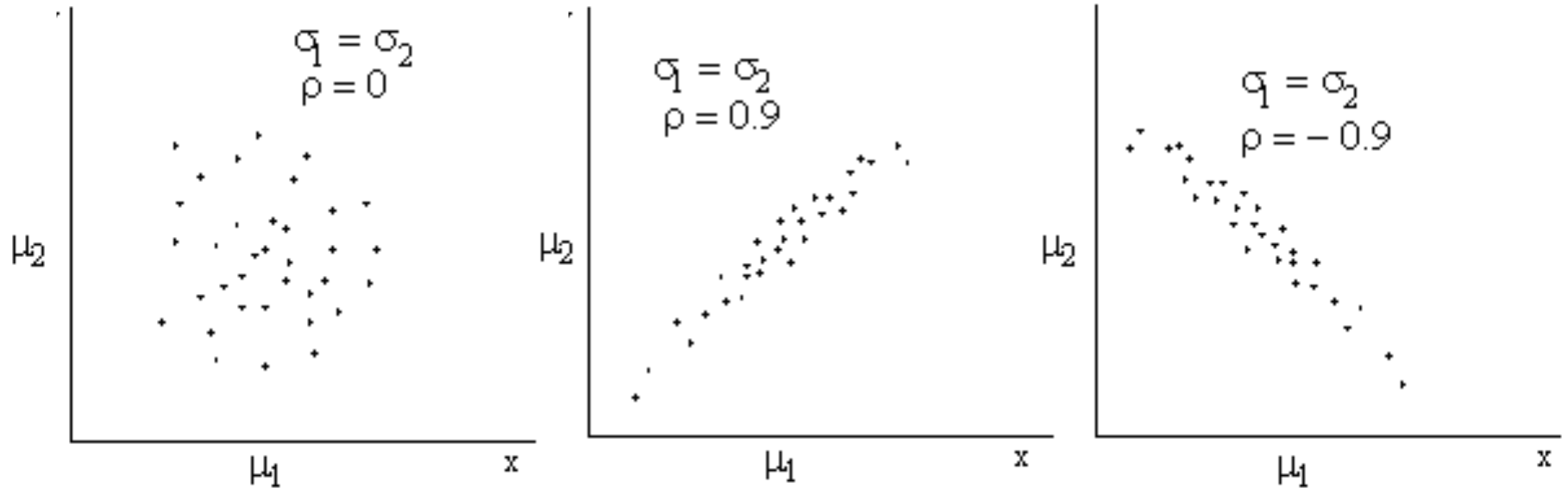
Surface Plots of the bivariate Normal distribution



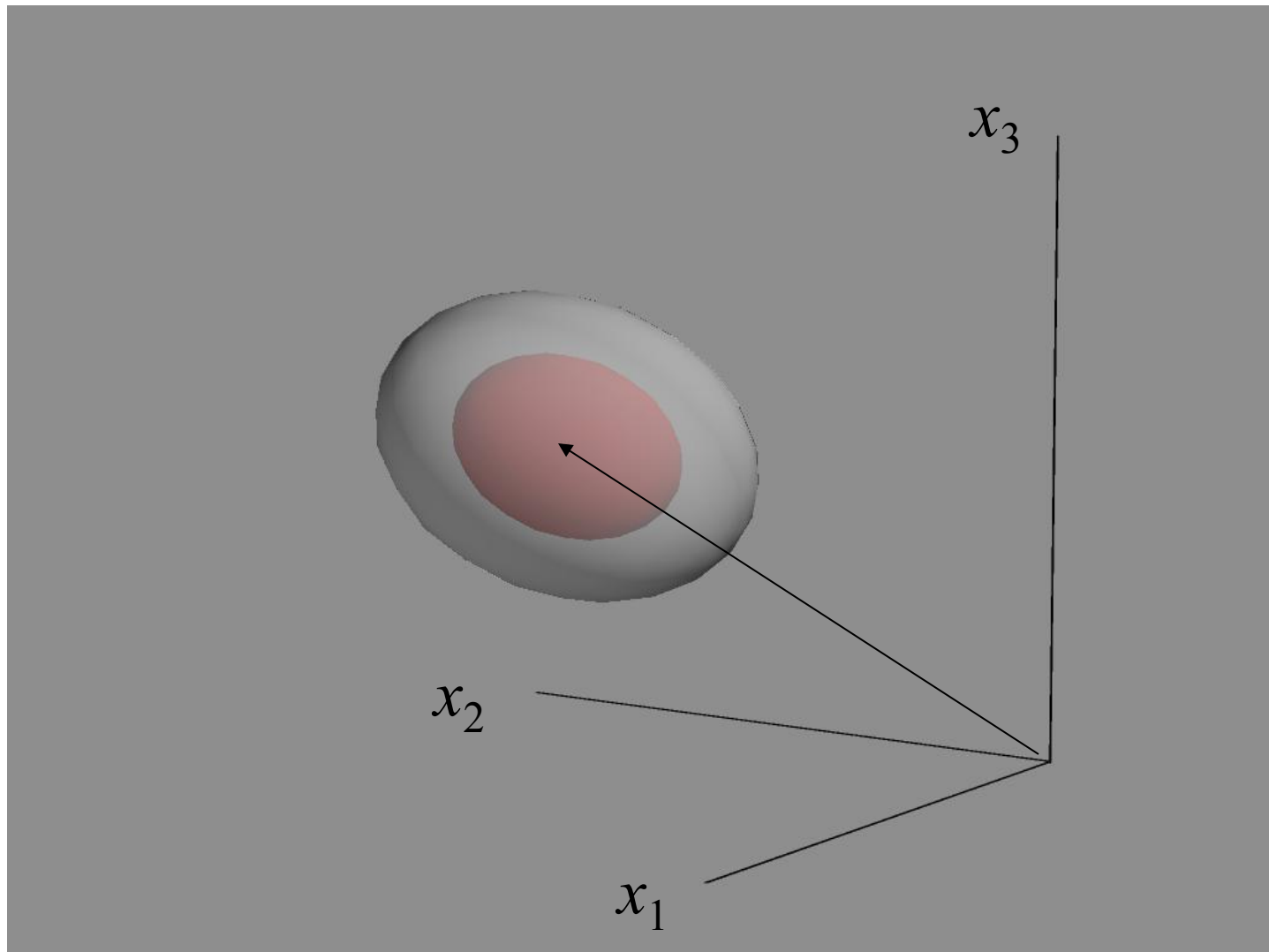
Contour Plots of the bivariate Normal distribution



Scatter Plots of data from the bivariate Normal distribution

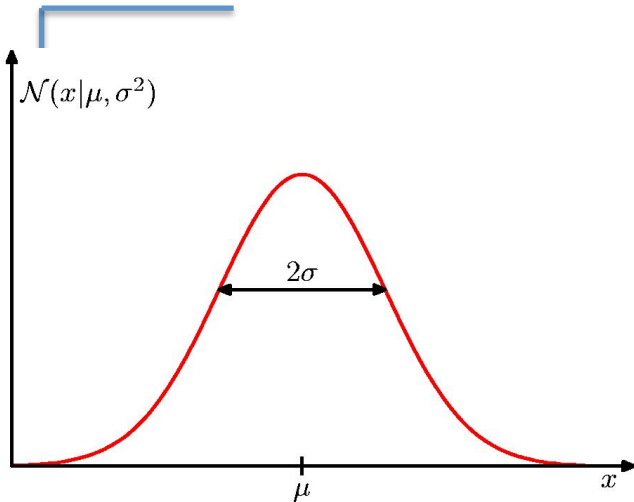


Trivariate Normal distribution



How to Estimate Gaussian:

MLE (Later)



- We can fit statistical models by maximizing the probability / likelihood of generating the observed samples:

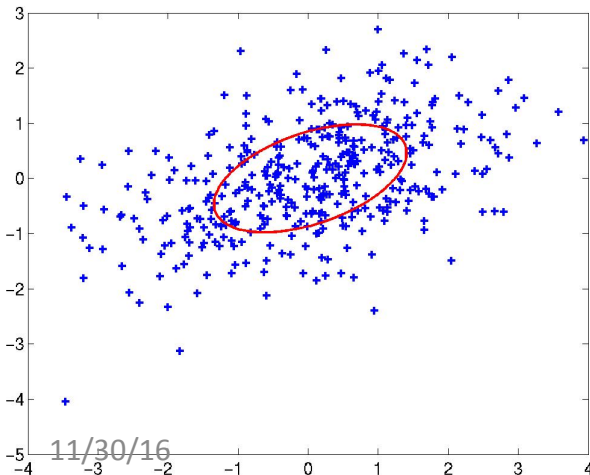
$$L(x_1, \dots, x_n | \theta) = p(x_1 | \theta) \dots p(x_n | \theta)$$

(the samples are assumed to be IID)

- In the 1D Gaussian case, we simply set the mean and the variance to the **sample mean** and the **sample variance**:

$$\bar{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\bar{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{\mu})^2$$



The p-multivariate Normal distribution

$$\langle X_1, X_2, \dots, X_p \rangle \sim N(\vec{\mu}, \Sigma)$$

$$\vec{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix} \quad p \times 1$$

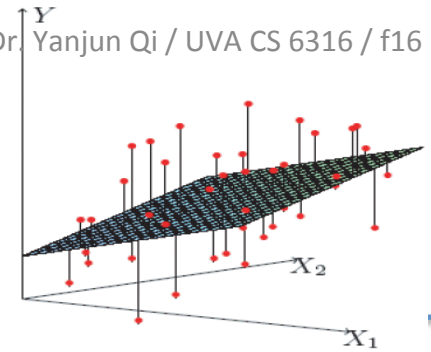
$$\mu_i = \frac{1}{n} \sum_{j=1}^N X_j^{(i)}$$

$\in \{1, 2, \dots, p\}$
i-th feature
 $\in \{1, 2, \dots, N\}$
j-th sample

$$\Sigma = \begin{bmatrix} \text{Var}(X_1) & & & \\ & \ddots & & \\ & & \text{Cov}(X_i, X_j) & \\ & & & \ddots \\ & & & & \text{Var}(X_p) \end{bmatrix} \quad \begin{matrix} | \\ | \\ | \\ | \\ | \end{matrix}$$

DETOUR: Probabilistic

Interpretation of Linear Regression



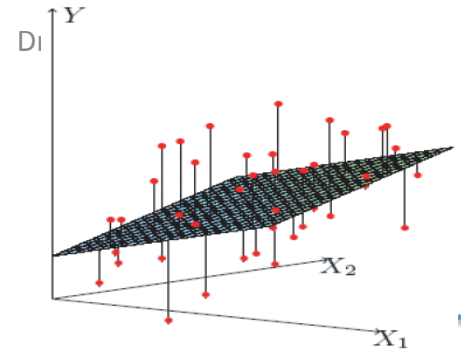
- Let us assume that the target variable and the inputs are related by the equation:

$$y_i = \theta^T \mathbf{x}_i + \varepsilon_i$$

where ε is an error term of unmodeled effects or random noise

DETOUR: Probabilistic

Interpretation of Linear Regression



- Let us assume that the target variable and the inputs are related by the equation:

$$y_i = \theta^T \mathbf{x}_i + \varepsilon_i$$

$$\text{RV } \varepsilon \sim N(0, \sigma^2)$$

where ε is an error term of unmodeled effects or random noise

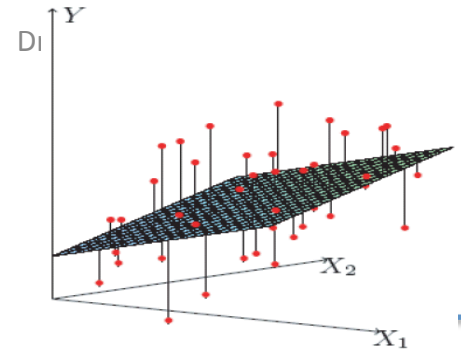
- Now assume that ε follows a Gaussian $N(0, \sigma^2)$, then we have:

$$p(y_i | x_i; \theta) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(y_i - \theta^T \mathbf{x}_i)^2}{2\sigma^2}\right)$$

$$\text{RV } y | x; \theta \sim N(\theta^T x, \sigma)$$

DETOUR: Probabilistic

Interpretation of Linear Regression



- Let us assume that the target variable and the inputs are related by the equation:

$$y_i = \theta^T \mathbf{x}_i + \varepsilon_i$$

where ε is an error term of unmodeled effects or random noise

- Now assume that ε follows a Gaussian $N(0, \sigma^2)$, then we have:

$$p(y_i | x_i; \theta) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(y_i - \theta^T \mathbf{x}_i)^2}{2\sigma^2}\right)$$

- By IID (independent and identically distributed) assumption:

$$L(\theta) = \prod_{i=1}^n p(y_i | x_i; \theta) = \left(\frac{1}{\sqrt{2\pi\sigma}}\right)^n \exp\left(-\frac{\sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2}{2\sigma^2}\right)$$

$$L(\theta) = \prod_{i=1}^n p(y_i | x_i; \theta) = \left(\frac{1}{\sqrt{2\pi\sigma}} \right)^n \exp \left(-\frac{\sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2}{2\sigma^2} \right)$$

We can learn θ by maximizing the probability / likelihood of generating the observed samples:

$$\begin{aligned}
 & P \left\{ (\vec{x}_1, y_1) \wedge (\vec{x}_2, y_2) \wedge \dots \wedge (\vec{x}_N, y_N) \right\} \\
 &= \prod_{i=1}^N p(y_i, \vec{x}_i) \stackrel{\text{IID}}{=} \prod_{i=1}^N p(y_i | \vec{x}_i; \theta) p(\vec{x}_i) \\
 &\theta^* = \operatorname{argmax}_{\theta} \prod_{i=1}^N p(y_i | \vec{x}_i; \theta)
 \end{aligned}$$

$$L(\theta) = \prod_{i=1}^n p(y_i | x_i; \theta) = \left(\frac{1}{\sqrt{2\pi\sigma}} \right)^n \exp \left(-\frac{\sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2}{2\sigma^2} \right)$$

We can learn θ by maximizing the probability / likelihood of generating the observed samples:

$$E \sim N(0, \sigma^2)$$

$$l(\theta) = \log(L(\theta)) = n \log \frac{1}{\sqrt{2\pi\sigma}} \left\{ -\frac{1}{\sigma^2} \frac{1}{2} \sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2 \right\}$$

$$\underset{\theta}{\operatorname{argmax}} l(\theta) \Rightarrow \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2$$

SSE

$$L(\theta) = \prod_{i=1}^n p(y_i | x_i; \theta) = \left(\frac{1}{\sqrt{2\pi\sigma}} \right)^n \exp \left(-\frac{\sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2}{2\sigma^2} \right)$$

We can learn θ by maximizing the probability / likelihood of generating the observed samples:



$$l(\theta) = \log(L(\theta)) = n \log \frac{1}{\sqrt{2\pi\sigma}} - \frac{1}{\sigma^2} \frac{1}{2} \sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2$$



$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \theta - y_i)^2$$

Maximum Likelihood Estimation

A general Statement

Consider a sample set $T=(X_1\dots X_n)$ which is drawn from a probability distribution $P(X|\theta)$ where θ are parameters.

If the X s are independent with probability density function $P(X_i|\theta)$, the joint probability of the whole set is

$$P(X_1\dots X_n|\theta) = \prod_{i=1}^n P(X_i|\theta)$$

Maximum Likelihood Estimation

A general Statement

Consider a sample set $T=(X_1\dots X_n)$ which is drawn from a probability distribution $P(X|\theta)$ where θ are parameters.

If the X s are independent with probability density function $P(X_i|\theta)$, the joint probability of the whole set is

$L(\theta)$ likelihood

$$P(X_1\dots X_n|\theta) = \prod_{i=1}^n P(X_i|\theta)$$

$LL(\theta) = \log(L(\theta))$

→ This may be maximised with respect to θ to give the maximum likelihood estimates (MLE) of θ :

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} P(X_1\dots X_n|\theta)$$

The idea is to

- ✓ assume a particular **model with unknown parameters**: θ
- ✓ we can then define the probability of observing a given event conditional on a particular set of parameters. $P(X_i | \theta)$
- ✓ We have observed **a set of outcomes** in the real world. x_1, x_2, \dots, x_n

The idea is to

- ✓ assume a particular **model with unknown parameters**, θ
- ✓ we can then define the probability of observing a given event conditional on a particular set of parameters. $P(X_i | \theta)$
- ✓ We have observed **a set of outcomes** in the real world. x_1, x_2, \dots, x_n
- ✓ It is then possible to choose a set of parameters which are most likely to have produced the observed results.

$$\hat{\theta} = \operatorname{argmax}_{\theta} \underbrace{P(X_1 \dots X_n | \theta)}_{\text{likelihood}}$$

The idea is to

- ✓ assume a particular **model with unknown parameters**, θ
- ✓ we can then define the probability of observing a given event conditional on a particular set of parameters. $P(X_i | \theta)$
- ✓ We have observed **a set of outcomes** in the real world. x_1, x_2, \dots, x_n
- ✓ It is then possible to choose a set of parameters which are most likely to have produced the observed results.

$$\hat{\theta} = \operatorname{argmax}_{\theta} \left[P(X_1 \dots X_n | \theta) \right] \rightarrow \text{likelihood}$$

This is maximum likelihood. In most cases it is **both consistent and efficient**. It provides a standard to compare other estimation techniques.

$$\log(L(\theta)) = \sum_{i=1}^n \left[\log(P(X_i | \theta)) \right] \rightarrow \text{log likelihood}$$

It is often convenient to work with the Log of the likelihood function.

DETOUR: Probabilistic

Interpretation of Linear Regression

- Hence the log-likelihood is:

$$l(\theta) = \log(L(\theta)) = n \log \frac{1}{\sqrt{2\pi\sigma}} - \frac{1}{\sigma^2} \frac{1}{2} \sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2$$

- Recognize the last term?

Yes it is:
$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \theta - y_i)^2$$

- Thus under independence Gaussian residual assumption, residual square error is equivalent to **MLE** of ϑ !

$$y_i \sim N(\exp(wx_i), 1)$$

(b) (6 points) (no explanation required) Suppose you decide to do a maximum likelihood estimation of w . You do the math and figure out that you need w to satisfy one of the following equations. Which one?

- A. $\sum_i x_i \exp(wx_i) = \sum_i x_i y_i \exp(wx_i)$
- B. $\sum_i x_i \exp(2wx_i) = \sum_i x_i y_i \exp(wx_i)$
- C. $\sum_i x_i^2 \exp(wx_i) = \sum_i x_i y_i \exp(wx_i)$
- D. $\sum_i x_i^2 \exp(wx_i) = \sum_i x_i y_i \exp(wx_i/2)$
- E. $\sum_i \exp(wx_i) = \sum_i y_i \exp(wx_i)$

$$y_i \sim N(\exp(wx_i), 1)$$

Answer: B (this is an extra credit question.)

$$L(\theta)$$

$$\downarrow$$

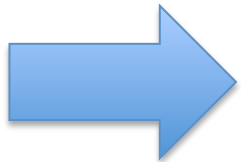
$$L(\theta)$$

$$\downarrow$$

$$\frac{\partial L(\theta)}{\partial \theta} = 0 \Rightarrow (B)$$

Today : Generative Bayes Classifiers

- ✓ Bayes Classifier
 - MAP classification rule
 - Generative Bayes Classifier
- ✓ Naïve Bayes Classifier
- ✓ Gaussian Bayes Classifiers
 - Gaussian distribution
 - Gaussian NBC
 - Not-naïve Gaussian BC → LDA, QDA



Gaussian Naïve Bayes Classifier

$$\operatorname{argmax}_C P(C | X) = \operatorname{argmax}_C P(X, C) = \operatorname{argmax}_C P(X | C)P(C)$$

Naïve
Bayes
Classifier

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

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

$N(\mu_{ji}, \sigma_{ji}^2)$

μ_{ji} : mean (average) of attribute values X_j of examples for which $C = c_i$

σ_{ji} : standard deviation of attribute values X_j of examples for which $C = c_i$

Gaussian Naïve Bayes Classifier

- Continuous-valued Input Attributes

- Conditional probability modeled with the normal distribution

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

μ_{ji} : mean (average) of attribute values X_j of examples for which $C = c_i$

σ_{ji} : standard deviation of attribute values X_j of examples for which $C = c_i$

- **Learning Phase:** for $\mathbf{X} = (X_1, \dots, X_p)$, $C = c_1, \dots, c_L$
Output: $p \times L$ normal distributions and $P(C = c_i) \quad i = 1, \dots, L$

$\left\{ \begin{array}{l} \mu_{ji} \\ \sigma_{ji} \end{array} \right.$

\rightarrow
 $j \in \{1, 2, \dots, p\}$
 $i \in \{1, 2, \dots, L\}$

MLE

$\left\{ \begin{array}{l} \text{sample mean} \\ \text{sample variance} \end{array} \right.$

Gaussian Naïve Bayes Classifier

- Continuous-valued Input Attributes

- Conditional probability modeled with the normal distribution

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

μ_{ji} : mean (average) of attribute values X_j of examples for which $C = c_i$

σ_{ji} : standard deviation of attribute values X_j of examples for which $C = c_i$

- **Learning Phase:** for $\mathbf{X} = (X_1, \dots, X_p)$, $C = c_1, \dots, c_L$
Output: $p \times L$ normal distributions and $P(C = c_i) \quad i = 1, \dots, L$

- **Test Phase:** for $\mathbf{X}' = (X'_1, \dots, X'_p)$

- Calculate conditional probabilities with all the normal distributions
- Apply the MAP rule to make a decision

$$\operatorname{argmax}_i P(C=c_i) P(X_1|c_i) \dots P(X_p|c_i)$$

Naïve Gaussian means ?

Total # para $\Rightarrow L \times \{p + p \times p\}$ $\xrightarrow{\mu/c}$ $\xrightarrow{\Sigma/c}$

Not Naïve

$$P(X_1, X_2, \dots, X_p | C) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Naïve

Total # para $\Rightarrow L \times (p + p)$

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

$$= \prod_i \frac{1}{\sqrt{2\pi\sigma_{ji}}} \exp \left(-\frac{(X_j - \mu_{ji})^2}{2\sigma_{ji}^2} \right)$$

$\Sigma | c_i = \begin{bmatrix} \sigma_{1i} & & \\ & \sigma_{2i} & \\ & & \ddots \\ & & & \sigma_{pi} \end{bmatrix}$

Diagonal Matrix

$$\boldsymbol{\Sigma} _ c_k = \boldsymbol{\Lambda} _ c_k$$

Each class' covariance matrix is diagonal

Today : Generative Bayes Classifiers

- ✓ Bayes Classifier
 - MAP classification rule
 - Generative Bayes Classifier
- ✓ Naïve Bayes Classifier
- ✓ Gaussian Bayes Classifiers
 - Gaussian distribution
 - Gaussian NBC
 - Not-naïve Gaussian BC → LDA, QDA

(1) covariance matrix are the same across classes

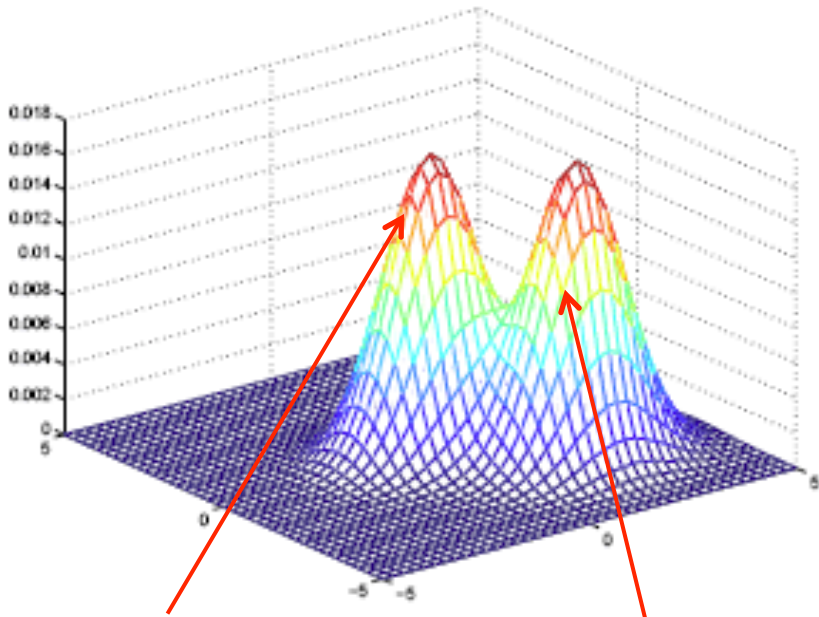
→ LDA (Linear Discriminant Analysis)

Linear Discriminant Analysis : $\Sigma_k = \Sigma, \forall k$

Each class' covariance matrix is the same

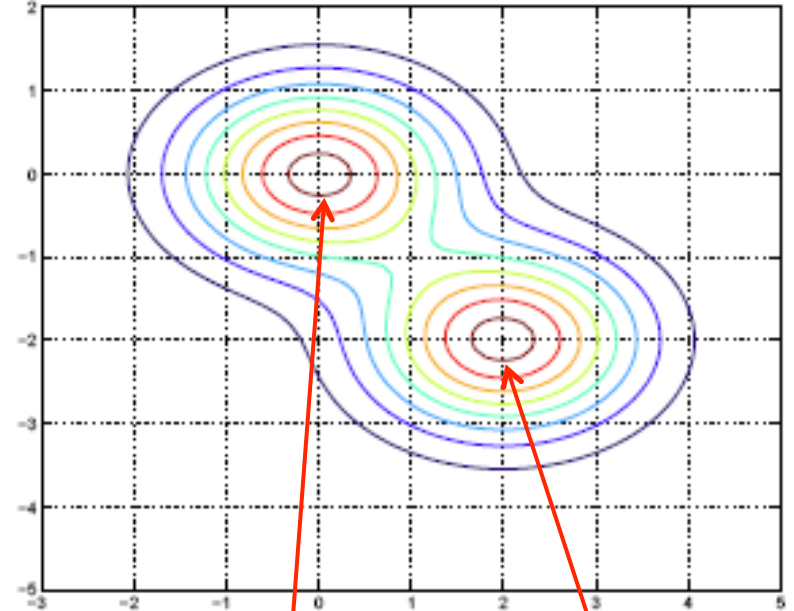
The Gaussian Distribution are shifted versions of each other

\downarrow
 $k \in \{1, 2, \dots, L\}$



Class k

Class l

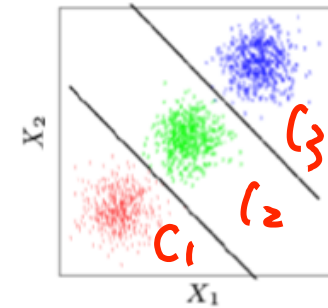


Class k

Class l

Optimal Classification

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



$$\operatorname{argmax}_k P(C_k | X) = \operatorname{argmax}_k P(X, C_k) = \operatorname{argmax}_k P(X | C_k) P(C_k)$$

$$= \operatorname{argmax}_k \left[-\log((2\pi)^{p/2} |\Sigma|^{1/2}) - \frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log(\pi_k) \right]$$

Short for $\pi_k = P(C_k)$
 $k \in \{1, 2, \dots, L\}$

$$= \operatorname{argmax}_k \left[-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log(\pi_k) \right]$$

- Note

Linear Discriminant Function for LDA

$$-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k - \frac{1}{2} \underline{x^T \Sigma^{-1} x}$$

$$\begin{aligned} \operatorname{argmax}_k P(C_k | X) &= \operatorname{argmax}_k P(X, C_k) = \operatorname{argmax}_k P(X | C_k) P(C_k) \\ &= \operatorname{argmax}_k \log\{P(X | C_k) P(C_k)\} \end{aligned}$$

Decision Boundary means those points

satisfying: $P(C_i | X) = P(C_j | X)$

$$\begin{aligned} \frac{P(C_i | X)}{P(C_j | X)} &= 1 \\ \Rightarrow \log \frac{P(C_i | X)}{P(C_j | X)} &= 0 \end{aligned}$$

$$\operatorname{argmax}_k P(C_k | X) = \operatorname{argmax}_k P(X, C_k) = \operatorname{argmax}_k P(X | C_k) P(C_k)$$

$$= \operatorname{argmax}_k \log \{ P(X | C_k) P(C_k) \}$$

$$= \operatorname{argmax}_k \log P(X | C_k) + \log P(C_k) \Rightarrow \pi_k$$

Decision Boundary points,

$$\log \frac{P(C_k | X)}{P(C_l | X)} = 0 = \log \frac{P(X | C_k)}{P(X | C_l)} + \log \frac{\pi_k}{\pi_l}$$

$$= \log P(X | C_k) - \log P(X | C_l) + \log \frac{\pi_k}{\pi_l}$$

$$\log \frac{P(C_k | X)}{P(C_l | X)} = \log \frac{P(X | C_k)}{P(X | C_l)} + \log \frac{P(C_k)}{P(C_l)}$$

$$= \log \frac{\pi_k}{\pi_l} - \frac{1}{2} (\mu_k + \mu_l)^T \Sigma^{-1} (\mu_k - \mu_l) \quad (4.9)$$

$$+ \underbrace{x^T \Sigma^{-1} (\mu_k - \mu_l)}_a = 0 \quad b$$

$\Rightarrow x^T a + b = 0 \Rightarrow$ a linear line decision boundary

Define **Linear Discriminant Function**

$$\delta(x) = -\frac{1}{2}(x - \mu_k)^T \Sigma^{-1}(x - \mu_k) + \log \pi_k$$

→ The **Decision Boundary Between class k and l** , $\{x : \delta_k(x) = \delta_l(x)\}$, **is linear**

$$= \log \frac{\pi_k}{\pi_l} - \frac{1}{2}(\mu_k + \mu_l)^T \Sigma^{-1}(\mu_k - \mu_l) \quad (4.9)$$
$$+ x^T \Sigma^{-1}(\mu_k - \mu_l),$$

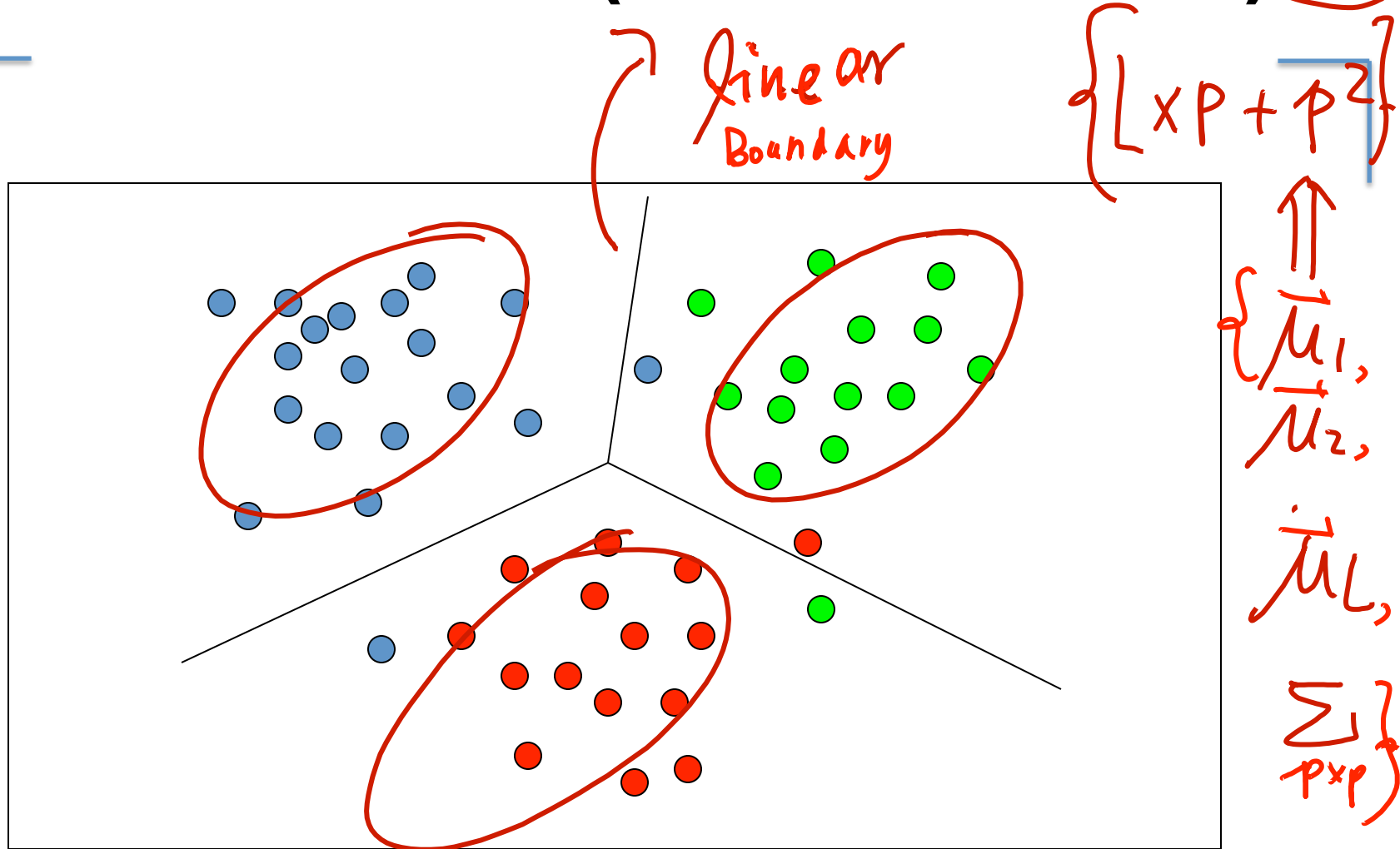
Equals to
zero

Boundary points X :

when $P(c_k|X) == P(c_l|X)$,

the left linear equation $==0$, a linear
line / plane

Visualization (three classes) LDA



(2) If covariance matrix are not same

e.g. → **QDA (Quadratic Discriminant Analysis)**

- ▶ Estimate the covariance matrix Σ_k separately for each class k , $k = 1, 2, \dots, K$.

- ▶ Quadratic discriminant function:

$$\delta_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log \pi_k .$$

$\{\Sigma_1, \Sigma_2, \dots, \Sigma_K, \mu_1, \mu_2, \dots, \mu_K\}$
 ↓
 Total # para

- ▶ Classification rule:

$$\hat{G}(x) = \arg \max_k \delta_k(x) .$$

$\delta_k(x) - \delta_l(x) = 0$
 ↑

Total # para
 $K \cdot (p + p^2)$
 $\{\mu_k, \Sigma_k\}$

- ▶ [Decision boundaries] are quadratic equations in x .

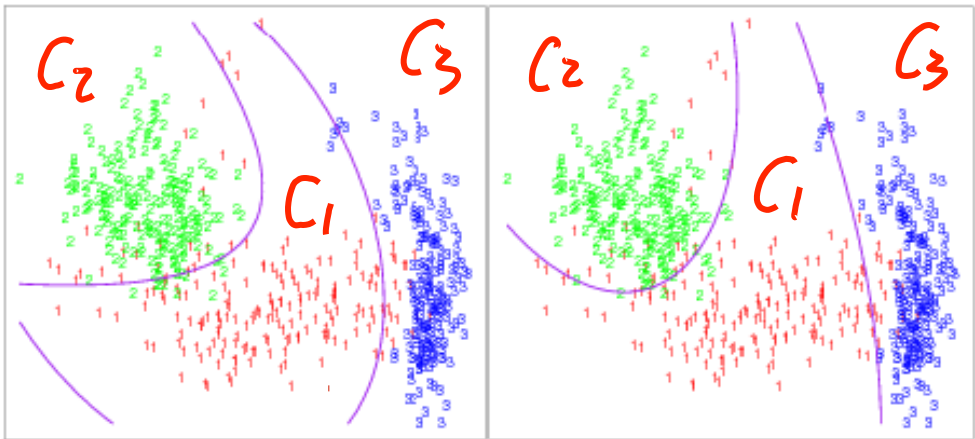
- ▶ QDA fits the data better than LDA, but has [more parameters] to estimate.

LDA on Expanded Basis

- ▶ Expand input space to include X_1X_2 , X_1^2 , and X_2^2 .
- ▶ Input is five dimensional: $X = (X_1, X_2, X_1X_2, X_1^2, X_2^2)$.

$\Phi(x) \leftarrow \text{LDA}$

LDA
with
 $\Phi(x)$



QDA

LDA with quadratic basis
Versus
QDA

Figure 4.6: Two methods for fitting quadratic boundaries. The left plot shows the quadratic decision boundaries for the data in Figure 4.1 (obtained using LDA in the five-dimensional space $x_1, x_2, x_{12}, x_1^2, x_2^2$). The right plot shows the quadratic decision boundaries found by QDA. The differences are small, as is usually the case.

Both with
Quadratic
decision
Boundary

(3) Regularized Discriminant Analysis

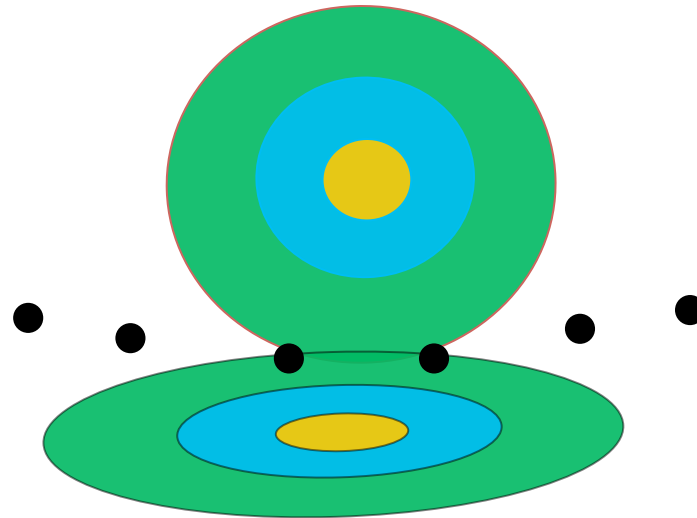
- ▶ A compromise between LDA and QDA.
- ▶ Shrink the separate covariances of QDA toward a common covariance as in LDA.
- ▶ Regularized covariance matrices:

$$\hat{\Sigma}_k(\alpha) = \alpha \hat{\Sigma}_k + (1 - \alpha) \hat{\Sigma} .$$

- ▶ The quadratic discriminant function $\delta_k(x)$ is defined using the shrunken covariance matrices $\hat{\Sigma}_k(\alpha)$.
- ▶ The parameter α controls the complexity of the model.

An example: Gaussian Bayes Classifier

Naive BC



$$\Sigma_1 \neq \Sigma_2$$

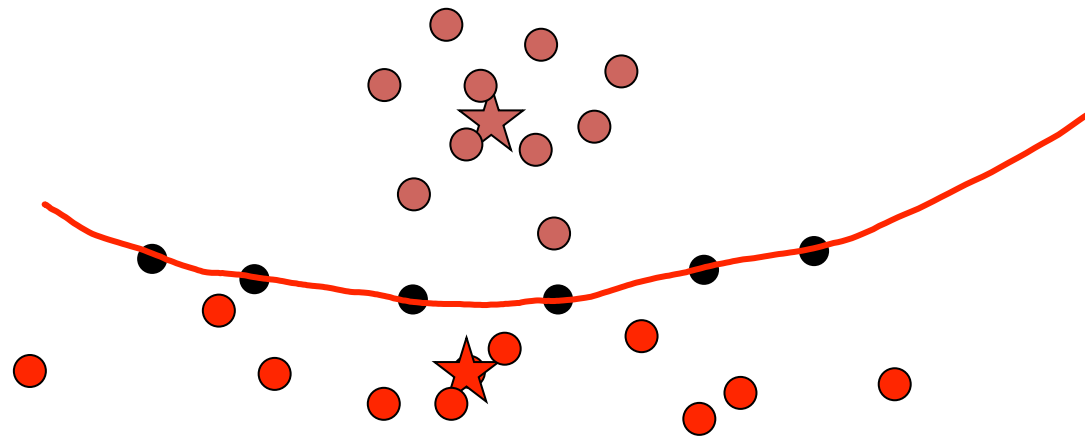
diagonal

$$\Sigma_1 = \Lambda_1$$

$$\Sigma_2 = \Lambda_2$$

diagonal

Gaussian Bayes Classifier



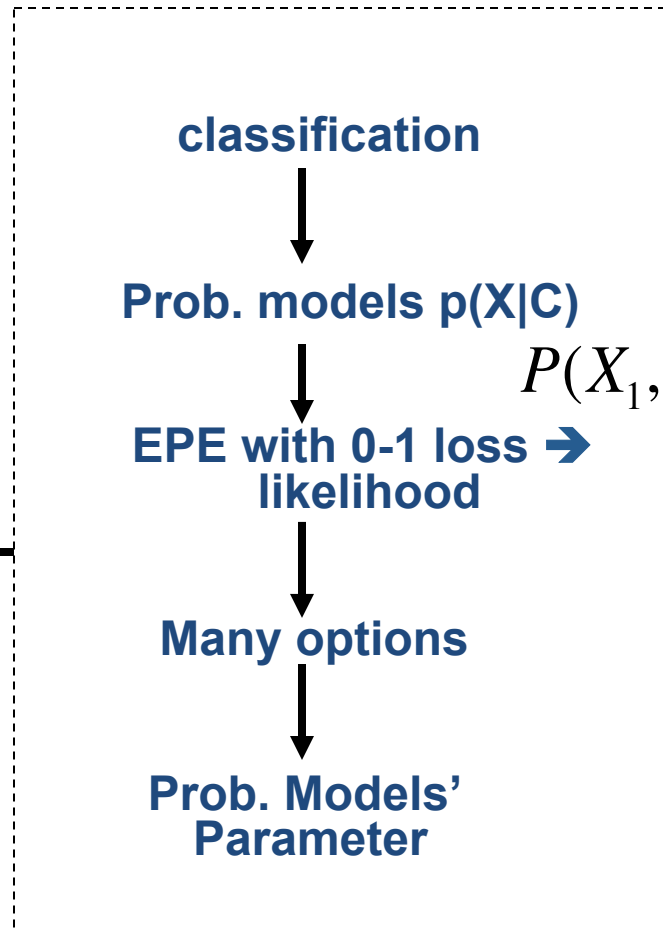
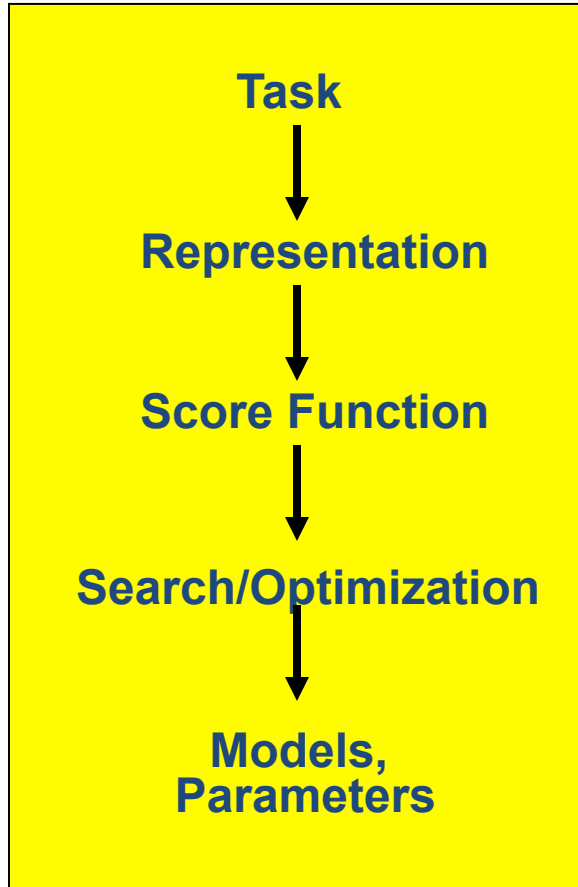
Quadratic decision Boundary

Today Recap : Generative Bayes Classifiers

- ✓ Bayes Classifier
 - MAP classification rule
 - Generative Bayes Classifier
- ✓ Naïve Bayes Classifier
- ✓ Gaussian Naïve Bayes Classifiers
 - Gaussian distribution
 - Gaussian NBC
 - Not-naïve Gaussian BC → LDA, QDA

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

Generative Bayes Classifier



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

Bernoulli Naïve

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

Gaussian Naïve

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

Multinomial

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

References

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