

# UVA CS 4501: Machine Learning

## Lecture 17b: Gaussian BC and Generative vs. Discriminative Classifier

Dr. Yanjun Qi

University of Virginia  
Department of Computer Science

# Course Content Plan →

Six major sections of this course

~~Regression (supervised)~~

Y is a continuous

Classification (supervised)

Y is a discrete

Unsupervised models

NO Y

Learning theory

About  $f()$

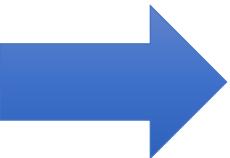
Graphical models

About interactions among  $X_1, \dots, X_p$

Reinforcement Learning

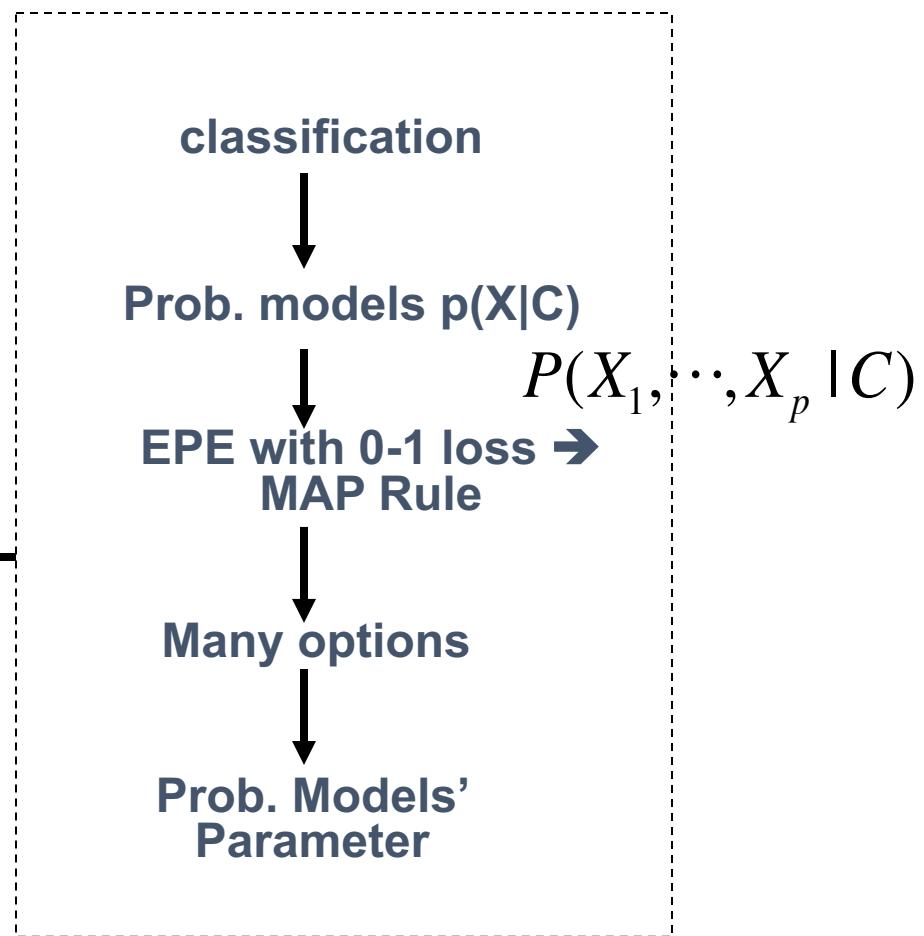
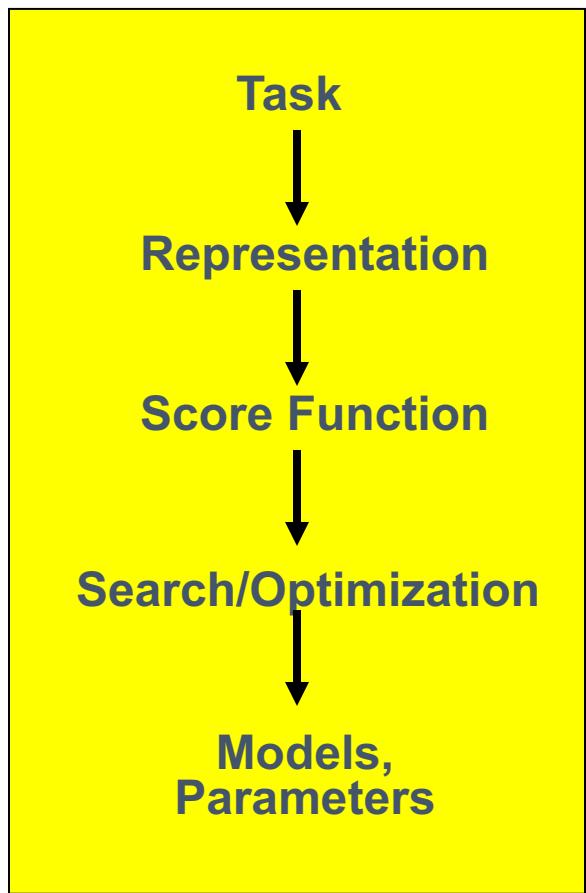
Learn program to Interact with its environment

# Today: More Generative Bayes Classifiers

- 
- ✓ Generative Bayes Classifier
  - ✓ Naïve Bayes Classifier
  - ✓ Gaussian Bayes Classifiers
    - Gaussian distribution
    - Naïve Gaussian BC
    - Not-naïve Gaussian BC → LDA, QDA
  - ✓ Discriminative vs. Generative classifier

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

## Generative Bayes Classifier



*Bernoulli  
Naïve*

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

*Gaussian  
Naïve*

*Multinomial*

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

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

# 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  $x$ .

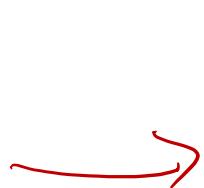
# Review: Probability of Continuous RV

- Properties of pdf

- 
- 

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

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



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

- 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:

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

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

# Review: Mean and Variance of RV

- Variance:  $Var(X) = E((X - \mu)^2)$

- Discrete RVs:

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

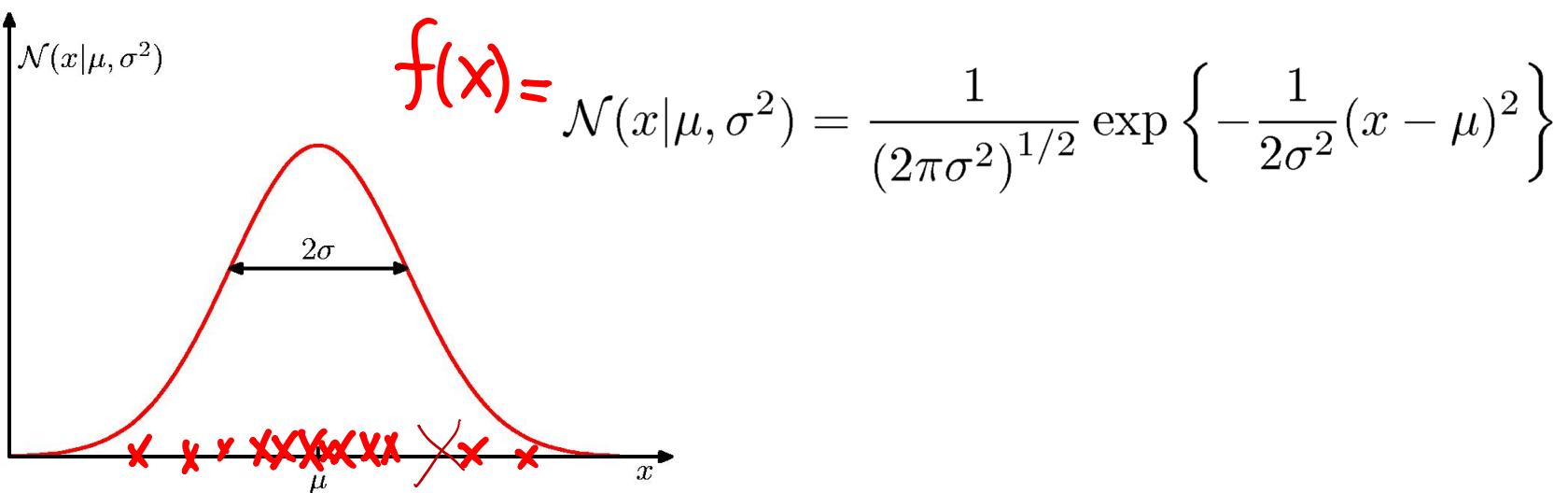
- Continuous RVs:

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

- Covariance:

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

# Single-Variate Gaussian Distribution



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

# Multivariate Normal (Gaussian) PDFs

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

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

Where  $|*$  represents determinant

Mean                          Covariance Matrix

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

#para :  $O(P + P^2)$

$\overset{\rightarrow}{\boldsymbol{\mu}}_{P \times 1}$  : mean vector

$\boldsymbol{\Sigma}_{P \times P}$  : covariance matrix

$$\begin{bmatrix} \sigma_1^2 & \text{---} \\ \sigma_2^2 & \text{---} \\ \vdots & \ddots \\ \sigma_P^2 & \text{---} \end{bmatrix} \xrightarrow{i} \text{cov}(X_i, X_j) \xrightarrow{j}$$

Review: Discrete RV  
 $p(x_1, x_2, \dots, x_p)$   
→ Nonnaive:  $2^P$   
Naive:

# 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  $|*$  represents determinant

Mean                          Covariance Matrix

- 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

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

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

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

$[\sigma_{12} = \rho\sigma_1\sigma_2]$

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

# Example: the Bivariate Normal distribution

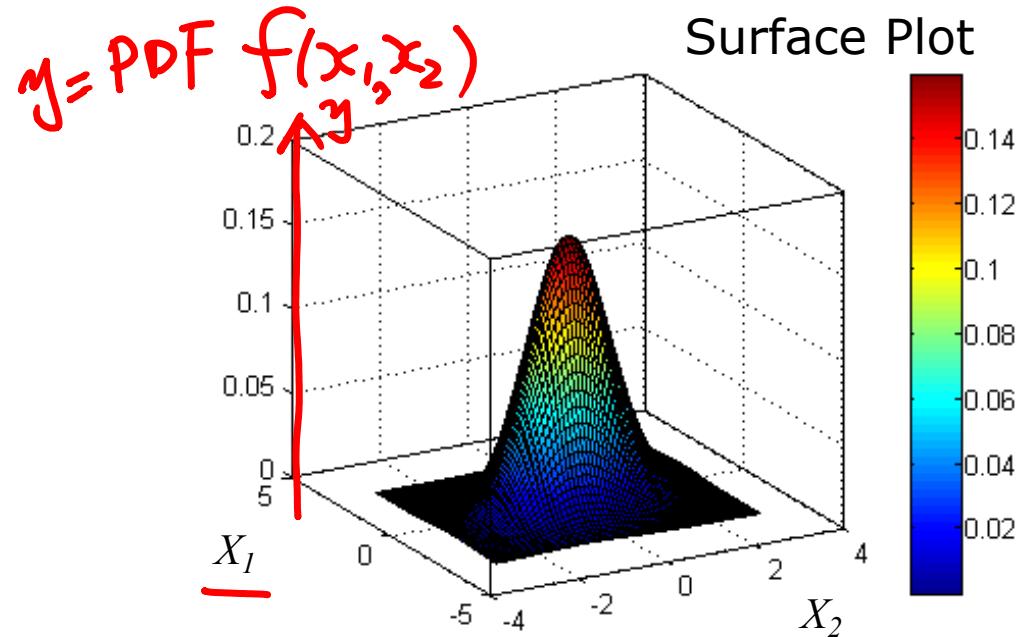
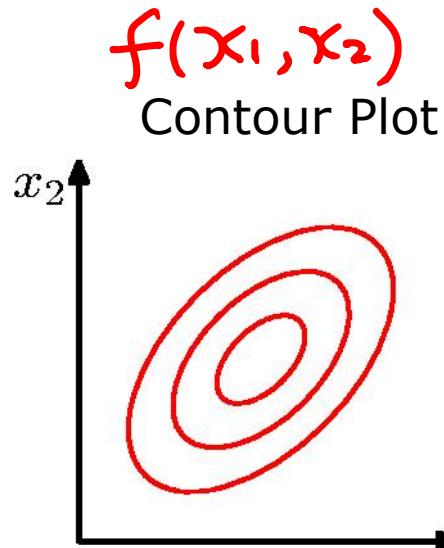
$$f(x_1, x_2) = \frac{1}{(2\pi)^{1/2} |\Sigma|} e^{-\frac{1}{2} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})}$$

with  $\vec{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}_{2 \times 1}$  and

$$\Sigma_{2 \times 2} = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \text{Cov}(x_1, x_2) \\ \underbrace{\rho \sigma_1 \sigma_2}_{\text{Cov}(x_1, x_2)} & \sigma_2^2 \end{bmatrix}_{2 \times 2}$$

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

# Bi-Variate Gaussian (normal) PDF

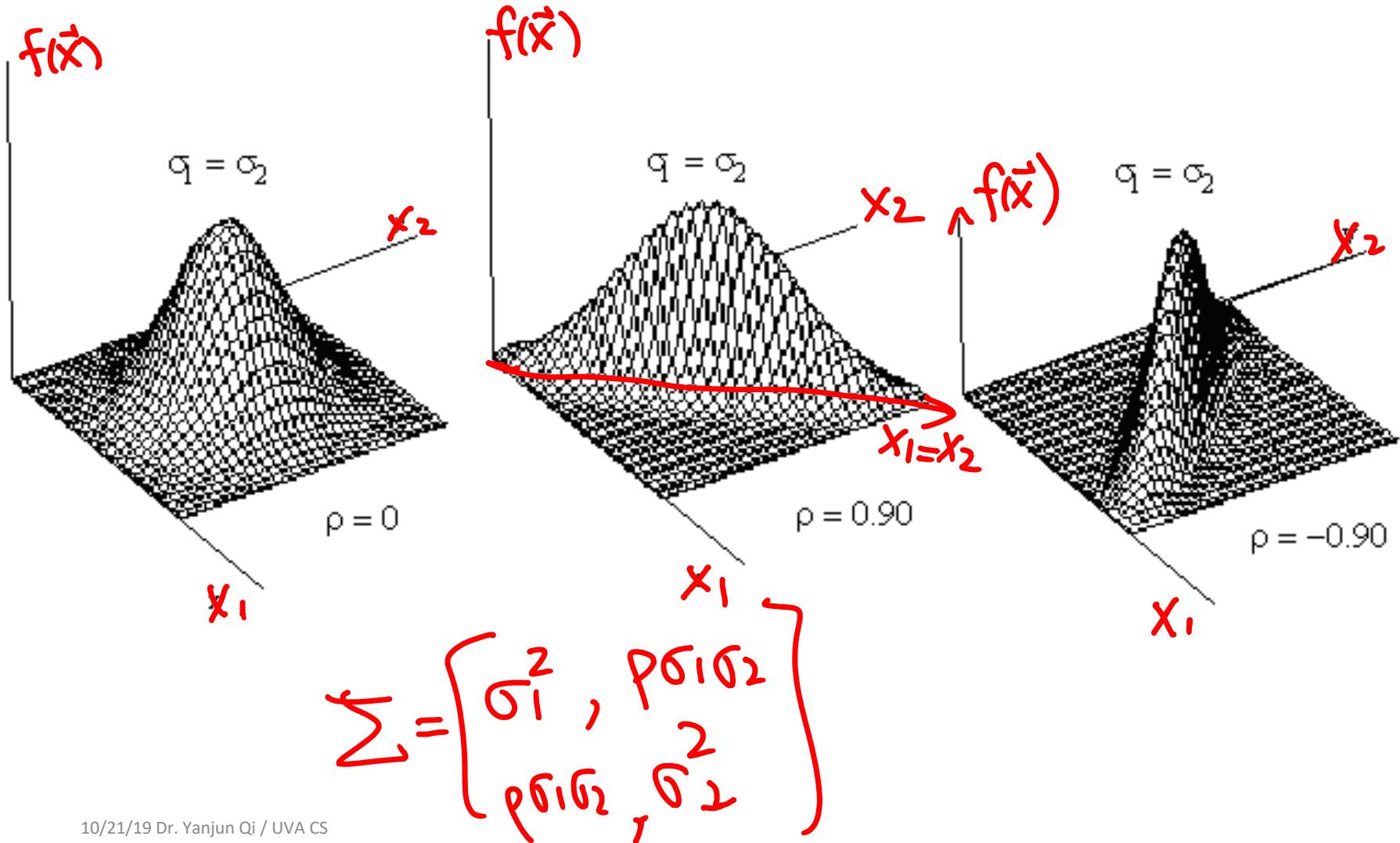


- Mean of a normal PDF is at peak value. Contours of equal PDF form ellipses.

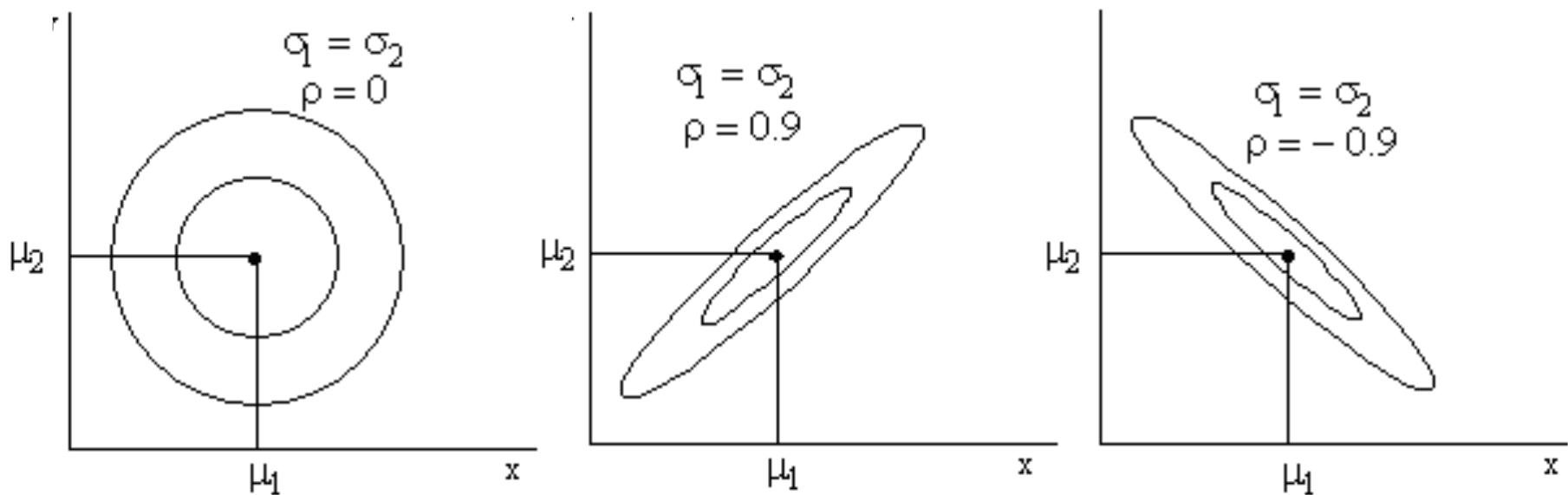
$$\vec{\mu} = \begin{bmatrix} \mu_{x_1} \\ \mu_{x_2} \end{bmatrix}$$

- The covariance matrix captures linear dependencies among the variables

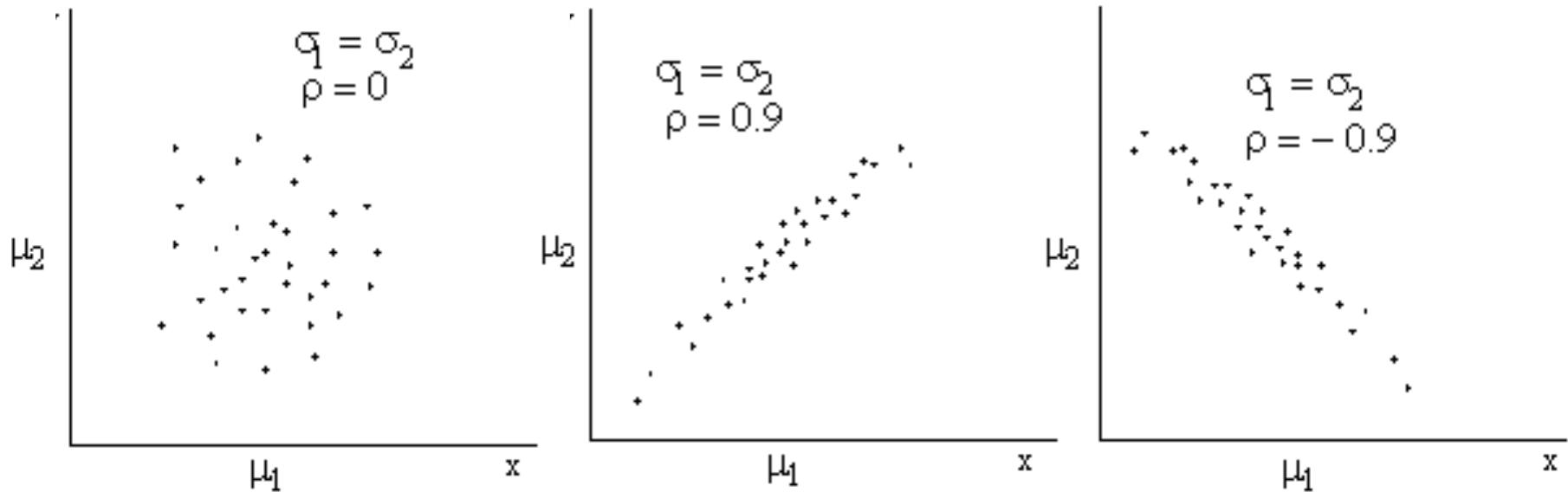
# Surface Plots of the bivariate Normal distribution



# Contour Plots of the bivariate Normal distribution

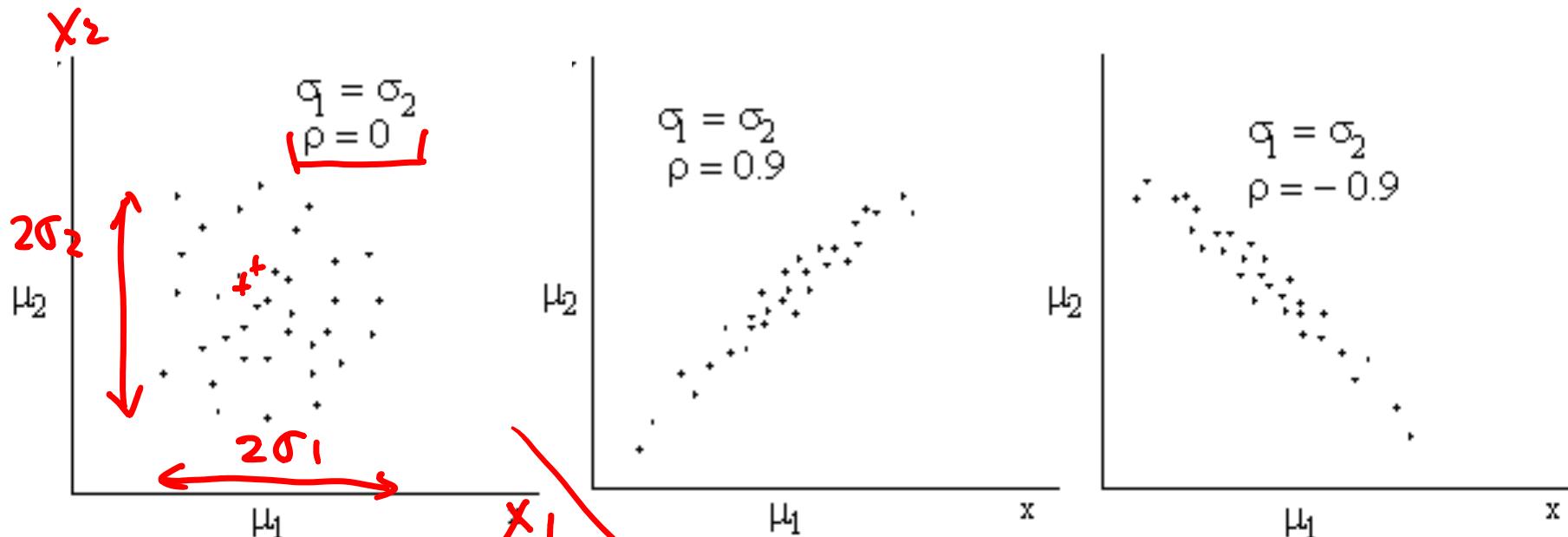


# Scatter Plots of samples from the three bivariate Normal distributions



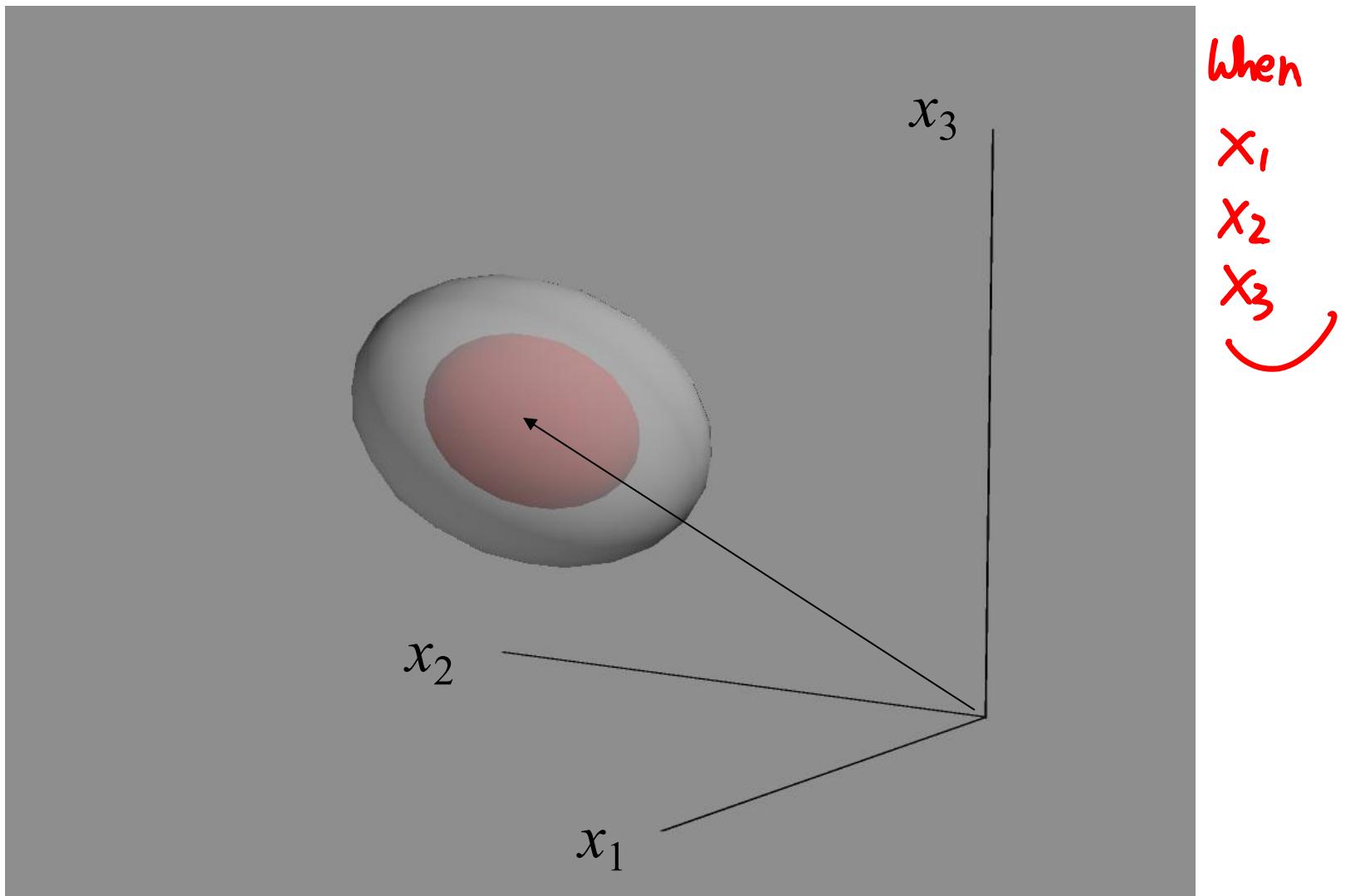
$N(\vec{\mu}, \Sigma)$

## Scatter Plots of samples from the three bivariate Normal distributions

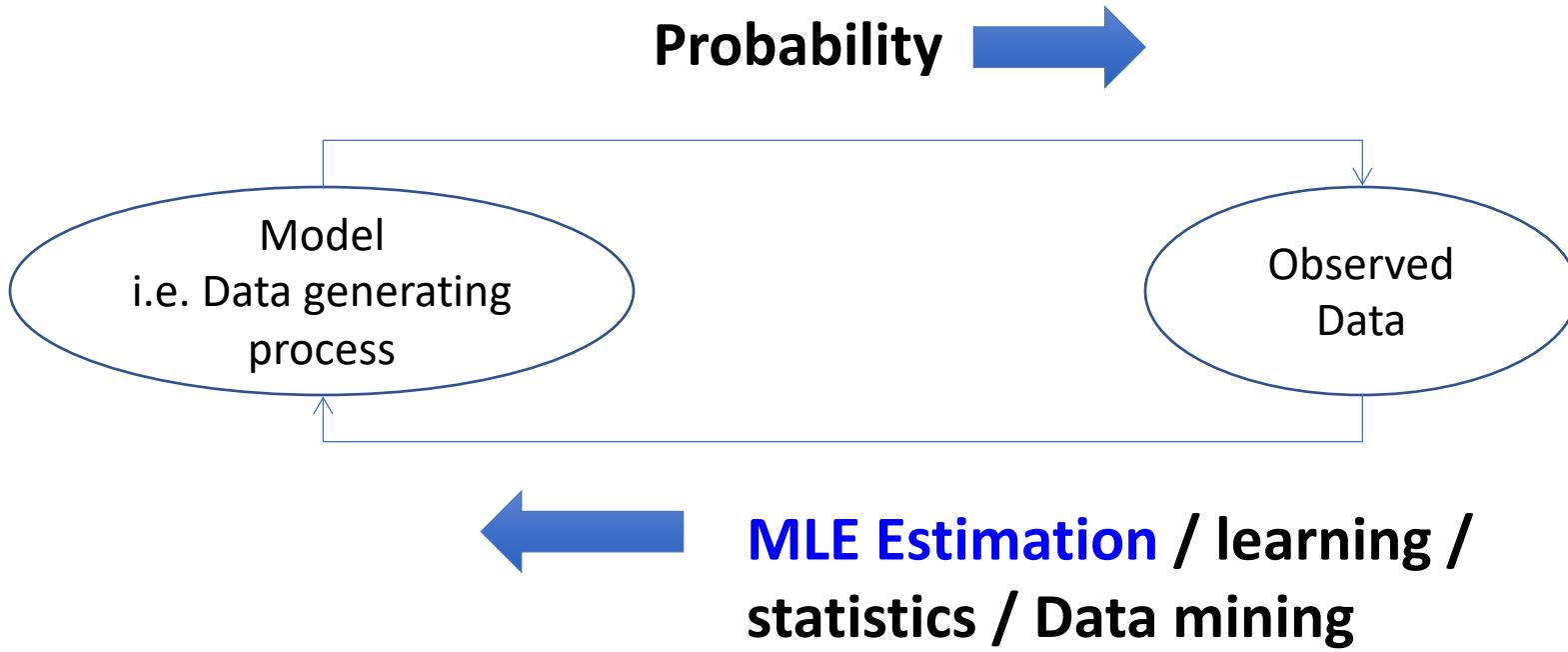


where  $\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$   $\Rightarrow f(x_1, x_2) = f(x_1)f(x_2) \xrightarrow{\text{data}} \begin{cases} \mu_1, \dots, \mu_p \\ \sigma_1, \dots, \sigma_p \\ 0(2\rho) \end{cases}$

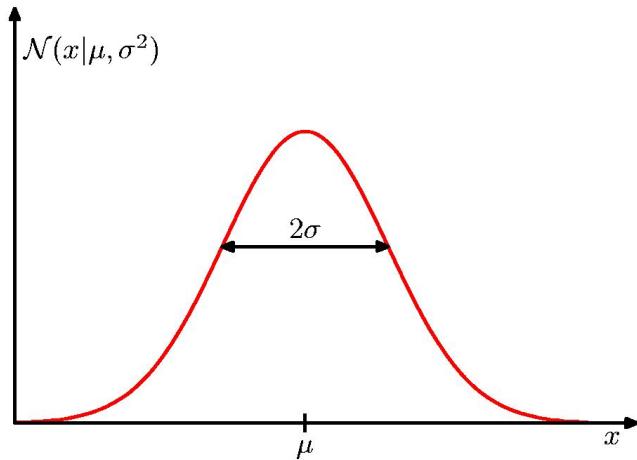
# Trivariate Normal distribution (Contour plot)



# The Big Picture



# How to Estimate 1D Gaussian: MLE



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

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

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

# How to Estimate p-D Gaussian: MLE

$\in \{1, 2, \dots, p\}$

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

# How to Estimate p-D Gaussian: MLE

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

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

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

$\in \{1, 2, \dots, p\}$

$i$ -th feature

$j$ -th sample

$\in \{1, 2, \dots, N\}$

$$\Sigma = \begin{bmatrix} \text{Var}(X_1) & & \\ & \ddots & \\ & & \text{Var}(X_p) \end{bmatrix} \quad p \times p$$

$i$

$j$

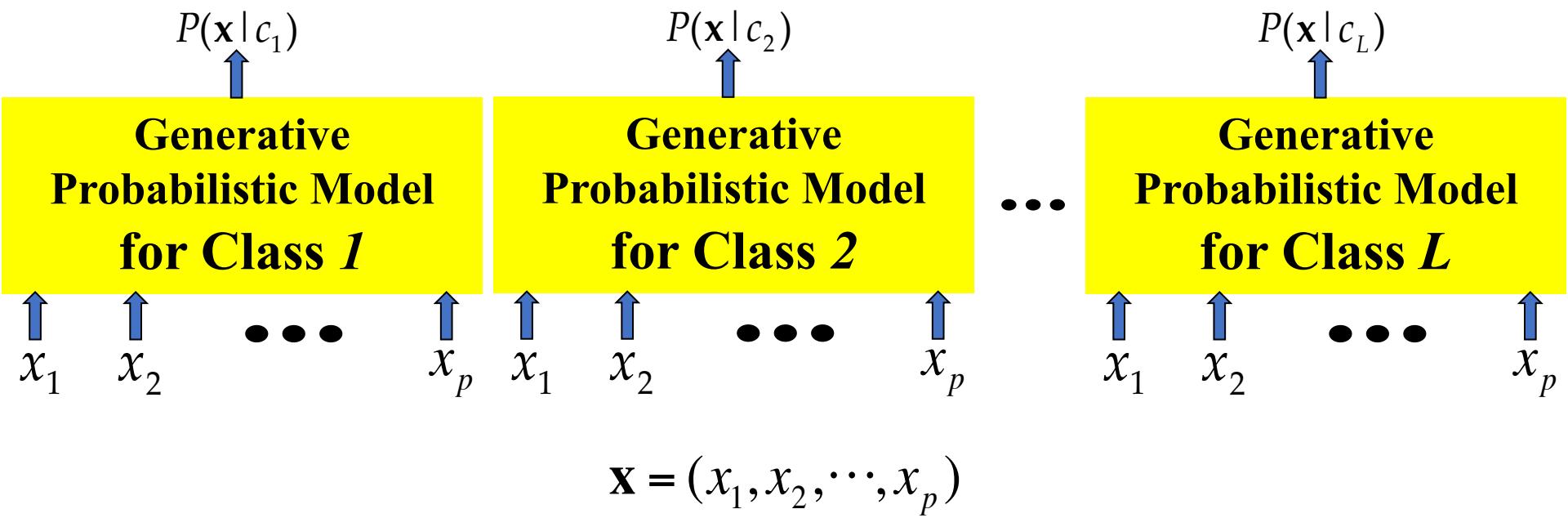
$\text{Cov}(X_i, X_j)$

#  $O(p + p^2)$

# Review: Generative BC

$$\begin{aligned} c^* &= \operatorname{argmax} P(C = c_i | \mathbf{X} = \mathbf{x}) \\ &\propto P(\mathbf{X} = \mathbf{x} | C = c_i) P(C = c_i) \\ &\text{for } i = 1, 2, \dots, L \end{aligned}$$

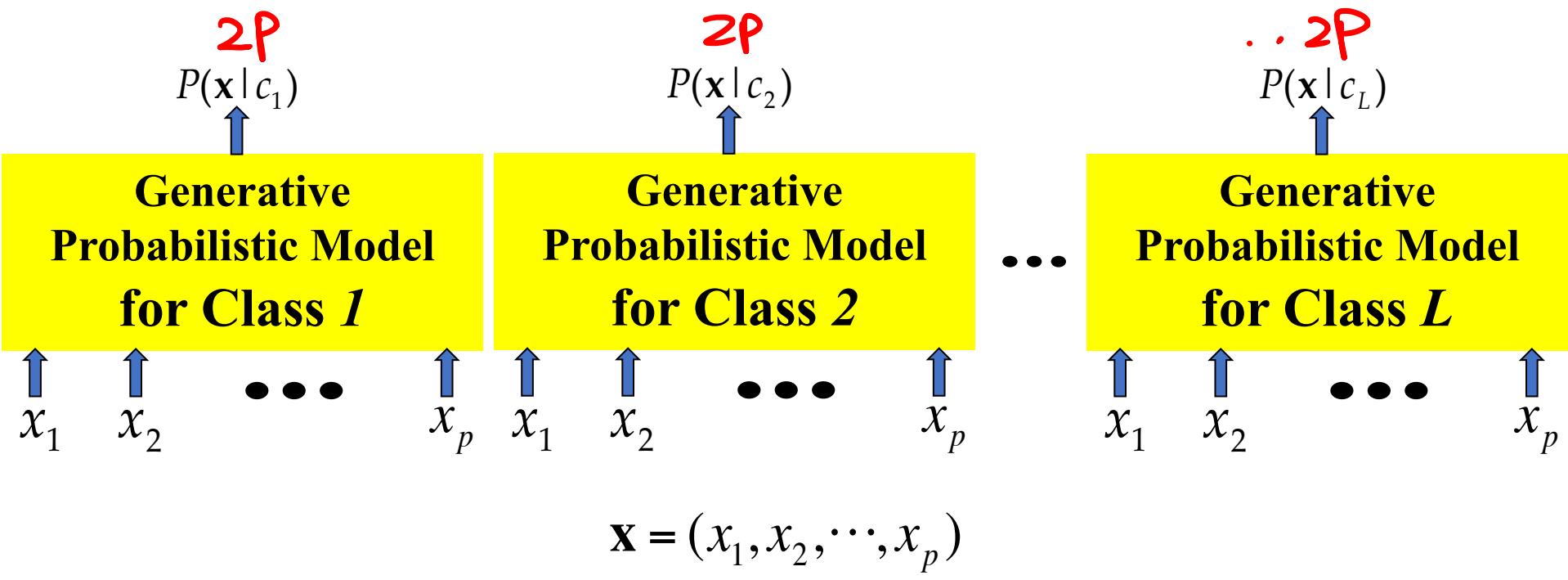
$$\begin{aligned} &P(\mathbf{X} | C), \\ &C = c_1, \dots, c_L, \mathbf{X} = (X_1, \dots, X_p) \end{aligned}$$



# Review: Generative BC

$$c^* = \operatorname{argmax} P(C = c_i | \mathbf{X} = \mathbf{x})$$
$$\propto P(\mathbf{X} = \mathbf{x} | C = c_i)P(C = c_i)$$
$$\text{for } i = 1, 2, \dots, L$$

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



# Review: 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

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

# Today: More Generative Bayes Classifiers

- ✓ Generative Bayes Classifier
  - ✓ Naïve Bayes Classifier
  - ✓ Gaussian Bayes Classifiers
    - Gaussian distribution
    - Naïve Gaussian BC
    - Not-naïve Gaussian BC → LDA, QDA
  - ✓ Discriminative vs. Generative
- 

# 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

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

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

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

$\sigma_{ji}$  : 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

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

$\mu_{ji}$  : mean (avearage) of attribute values  $X_j$  of examples for which  $C = c_i$

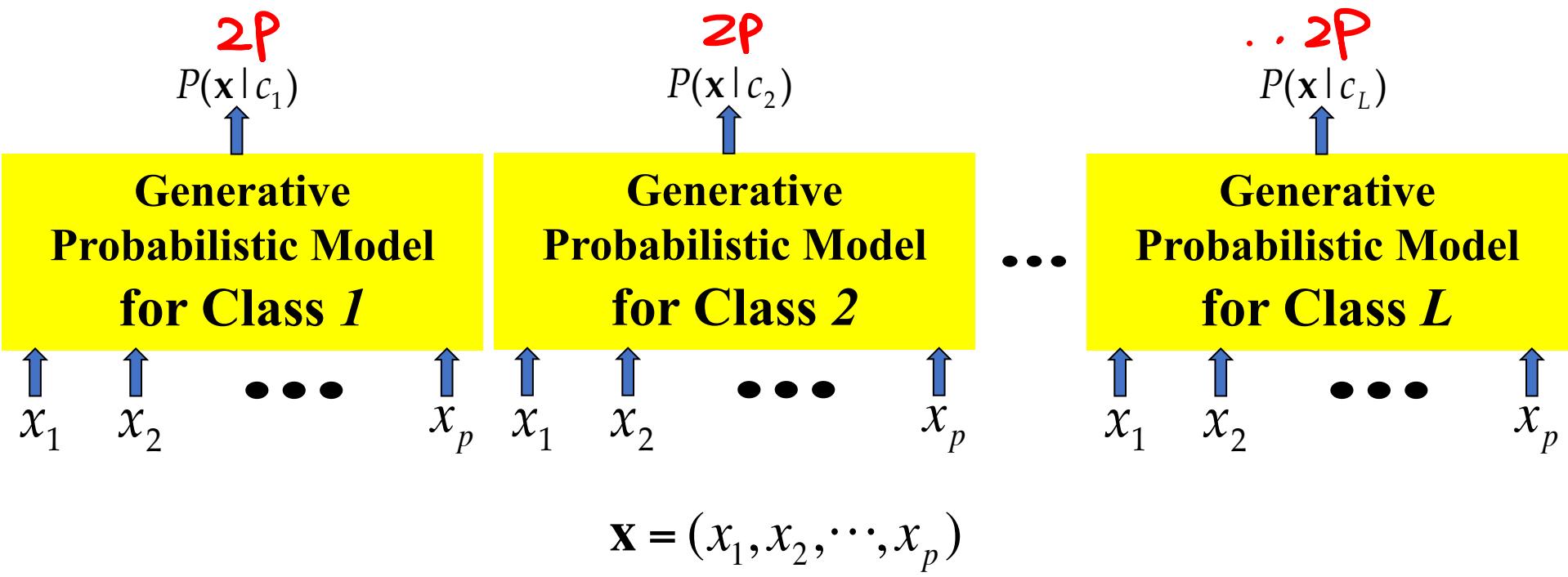
$\sigma_{ji}$  : standard deviation of attribute values  $X_j$  of examples for which  $C = c_i$

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

# Review: Generative BC

$$\begin{aligned} c^* &= \operatorname{argmax} P(C = c_i | \mathbf{X} = \mathbf{x}) \\ &\propto P(\mathbf{X} = \mathbf{x} | C = c_i) P(C = c_i) \\ &\text{for } i = 1, 2, \dots, L \end{aligned}$$

$$\begin{aligned} &P(\mathbf{X} | C), \\ &C = c_1, \dots, c_L, \mathbf{X} = (X_1, \dots, X_p) \end{aligned}$$



# Gaussian Naïve Bayes Classifier

- Continuous-valued Input Attributes
  - Conditional probability modeled with the normal distribution

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

$\mu_{ji}$  : mean (avearage) of attribute values  $X_j$  of examples for which  $C = c_i$

$\sigma_{ji}$  : 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: L different p-normal distributions and  $P(C = c_i)$   $i = 1, \dots, L$

# Gaussian Naïve Bayes Classifier

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

Naïve  
Bayes  
Classifier

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

O(L \times 2P + L)

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

$\mu_{ji}$  : mean (avearage) of attribute values  $X_j$  of examples for which  $C = c_i$

$\sigma_{ji}$  : 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

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

$\mu_{ji}$  : mean (avearage) of attribute values  $X_j$  of examples for which  $C = c_i$

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

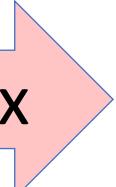
- **Learning Phase:** for  $\mathbf{X} = (X_1, \dots, X_p)$ ,  $C = c_1, \dots, c_L$   
Output: L different p-normal distributions and  $P(C = c_i) \ i = 1, \dots, L$

- **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  $\arg\max_i P(C=c_i) P(X'_1|c_i) \dots P(X'_p|c_i)$

when  $\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$    $f(x_1, x_2) = f(x_1)f(x_2)$    $\xrightarrow{\text{data}} \begin{cases} M_1, \dots, M_p \\ \sigma_1, \dots, \sigma_p \end{cases}$   
 $O(2P)$

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

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


Diagonal Matrix 

$$\sum_{-c_k} = \Lambda_{-c_k}$$

Each class' covariance matrix is diagonal

when  $\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \Rightarrow f(x_1, x_2) = f(x_1)f(x_2)$   $\xrightarrow{\text{data}} \begin{cases} M_1, \dots, M_p \\ \sigma_1, \dots, \sigma_p \end{cases}$   
 $O(2P)$



Total #param  $\xrightarrow{} L \times (P + P_f)$

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

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

$\sum |C_i| = \begin{bmatrix} \sigma_{11} & & & \\ & \sigma_{22} & & \\ & & \ddots & \\ & & & \sigma_{pp} \end{bmatrix}$

Diagonal Matrix

$$\Sigma - c_k = \Lambda - c_k$$

Each class' covariance matrix is diagonal

# Today: More Generative Bayes Classifiers

- ✓ Generative Bayes Classifier
- ✓ Naïve Bayes Classifier
- ✓ Gaussian Bayes Classifiers
  - Gaussian distribution
  - Naïve Gaussian BC
  - Not-naïve Gaussian BC → LDA, QDA
    - LDA: Linear Discriminant Analysis
    - QDA: Quadratic Discriminant Analysis
- ✓ Discriminative vs. Generative

# Not Naïve Gaussian means ?

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

$$\begin{aligned} P(X_1, X_2, \dots, X_p | C = c_j) &= P(X_1 | C)P(X_2 | C) \cdots P(X_p | C) \\ &= \prod_i \frac{1}{\sqrt{2\pi}\sigma_{ji}} \exp \left( -\frac{(X_j - \mu_{ji})^2}{2\sigma_{ji}^2} \right) \end{aligned}$$

Diagonal Matrix

$$\sum_c c_k = \Lambda_c c_k$$

Each class' covariance matrix is diagonal

# Not Naïve Gaussian means ?

$$P=28 \times 28, L \sim 10^3, \Rightarrow O(10^7)$$

$$\vec{\Sigma}_c, \vec{\mu}_c \Rightarrow O(LP + L \cdot P^2)$$

Not  
Naïve

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

$$\Rightarrow O(2PL)$$

Naïve

$$P(X_1, X_2, \dots, X_p | C = c_j) = P(X_1 | C)P(X_2 | C) \cdots P(X_p | C) \\ = \prod_i \frac{1}{\sqrt{2\pi}\sigma_{ji}} \exp \left( -\frac{(X_j - \mu_{ji})^2}{2\sigma_{ji}^2} \right)$$

Diagonal Matrix

11/11/19

$$\sum_c c_k = \Lambda c_k$$

Each class' covariance matrix is diagonal

38

# Not Naïve Gaussian means ?

Total # param  $\Rightarrow L \times \{P + P \times P\}$

$\mu/C$   
 $\Sigma/C$

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

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

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

$\sum |C_i| = \begin{bmatrix} \sigma_{11} & & & \\ & \sigma_{22} & & \\ & & \ddots & \\ & & & \sigma_{pp} \end{bmatrix}$

Diagonal Matrix

$$\sum c_k = \Lambda c_k$$

Each class' covariance matrix is diagonal

# Not-naïve Gaussian BC

- LDA: Linear Discriminant Analysis
- QDA: Quadratic Discriminant Analysis

$$\Sigma_1 = \dots = \Sigma_L = \Sigma$$

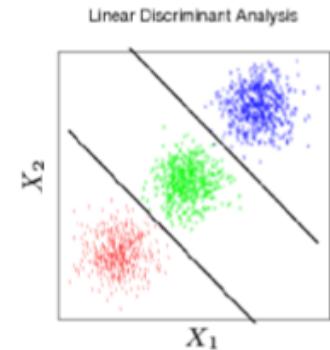
$\Sigma \Rightarrow P^2$ ,  $P \sim 100$ ,  $L \sim 10$

$O(n) < \underbrace{10k}_{10^4}$   $\cancel{\Rightarrow} O(LP^2 + LP) \sim \underline{10^5} \xrightarrow{\text{LDA}} O(P^2 + \overbrace{LP}^{\overline{m}_c})$

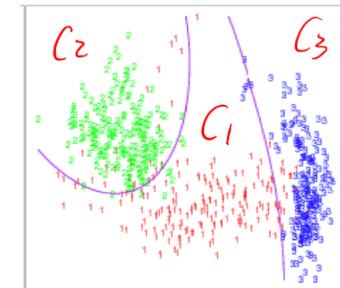
# Not-naïve Gaussian BC



- LDA: Linear Discriminant Analysis



- QDA: Quadratic Discriminant Analysis

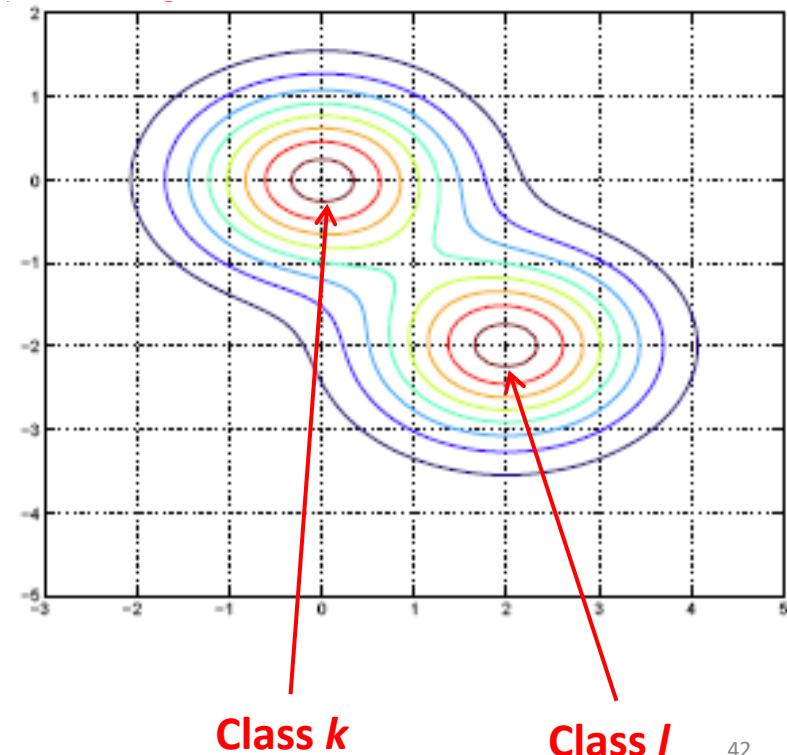
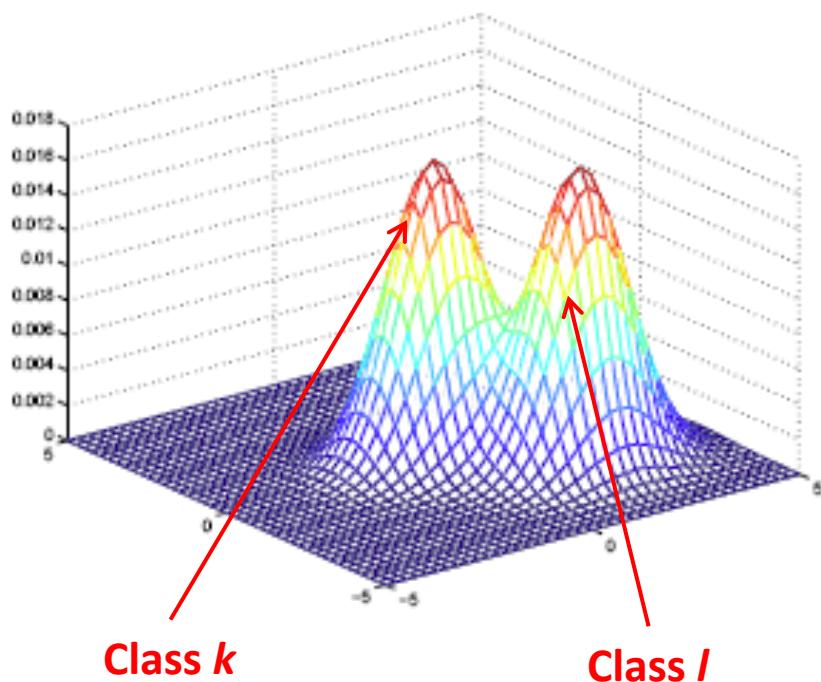


(1) covariance matrix are the same across classes  
→ LDA (Linear Discriminant Analysis)

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

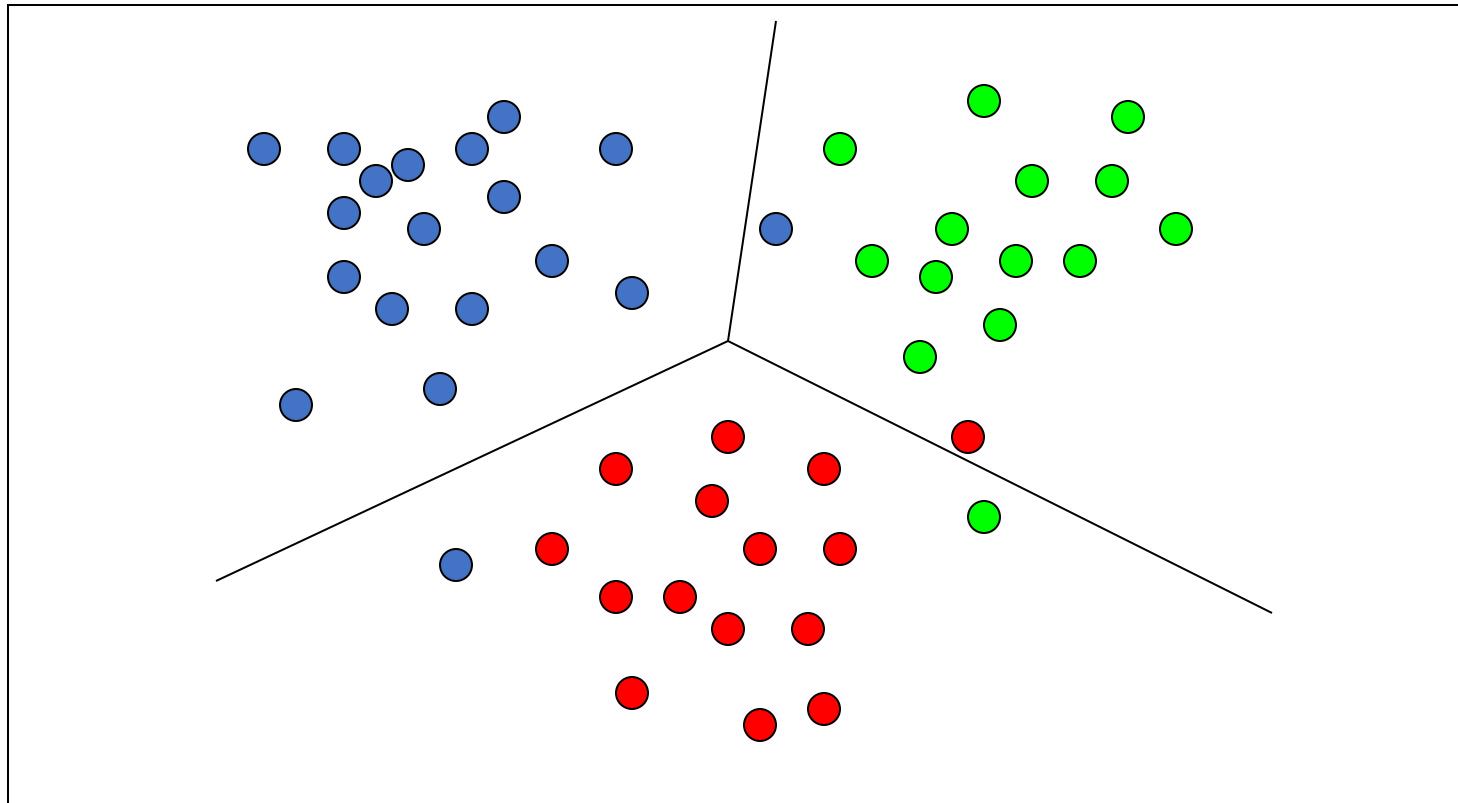
Each class' covariance matrix is the same

The Gaussian Distribution are shifted versions of each other

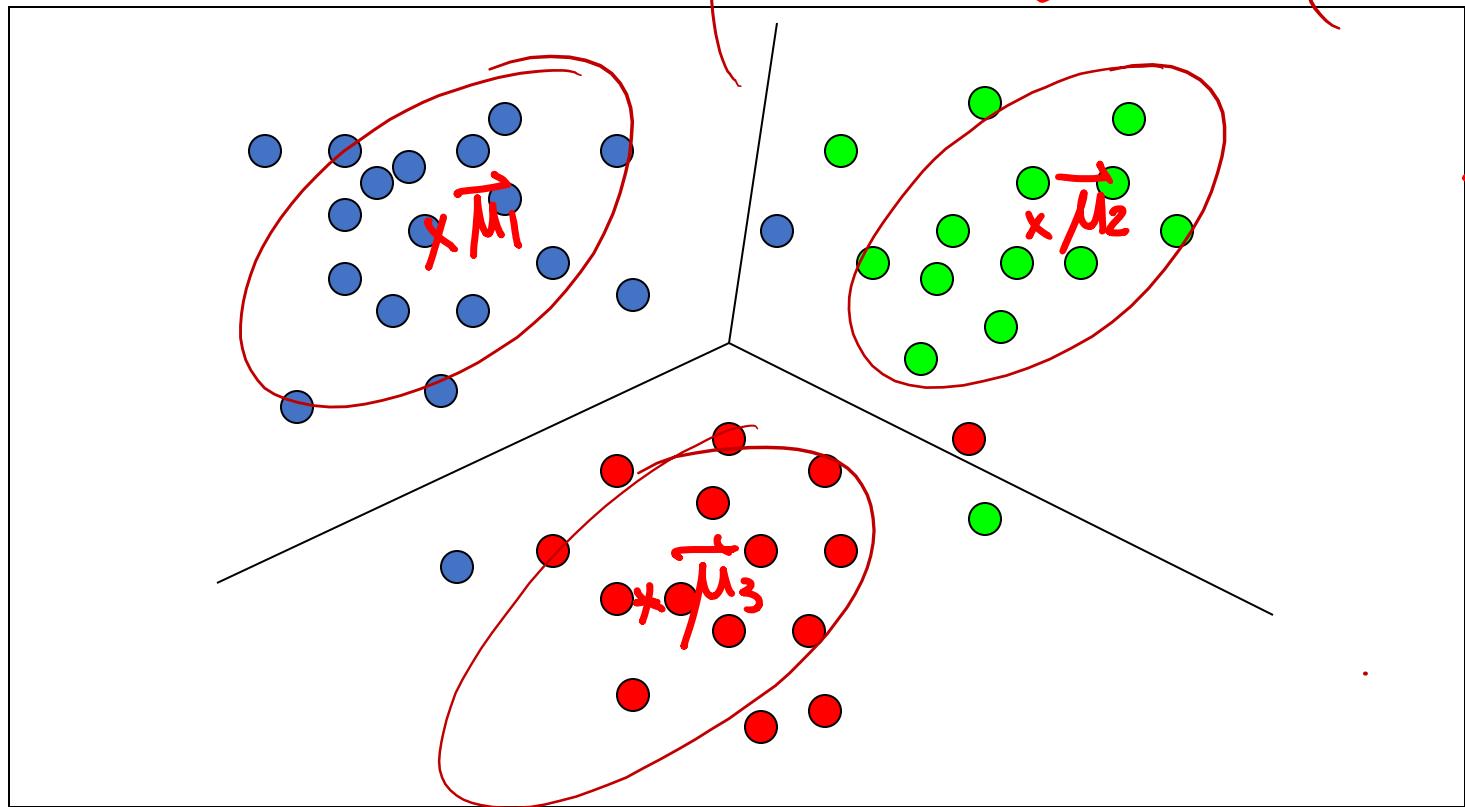


# Visualization (three classes)

$$\Sigma_1 = \Sigma_2 = \dots = \Sigma_L \Rightarrow \text{linear}$$



# Visualization (three classes)



LDA

$$\left\{ \begin{matrix} Xp + p_2 \\ \bar{\mu}_1, \\ \bar{\mu}_2, \\ \bar{\mu}_3, \\ \sum p_x p \end{matrix} \right\}$$

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

Decision Boundary Points →

Satisfying:  $\hat{P}(C_i | X) = \hat{P}(C_j | X)$

$$\frac{\hat{P}(C_i | X)}{\hat{P}(C_j | X)} = 1$$

$$\Rightarrow \log \frac{\hat{P}(C_i | X)}{\hat{P}(C_j | X)} = 0$$

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

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

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

Decision Boundary points

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

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

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

Decision Boundary Points of LDA classifier →

$$\begin{aligned} &= \log \frac{\pi_k}{\pi_\ell} - \frac{1}{2}(\mu_k + \mu_\ell)^T \Sigma^{-1} (\mu_k - \mu_\ell) \\ &\quad + x^T \Sigma^{-1} (\mu_k - \mu_\ell), \end{aligned} \tag{4.9}$$

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

Decision Boundary Points of LDA classifier →

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

The above is derived from the following :

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

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

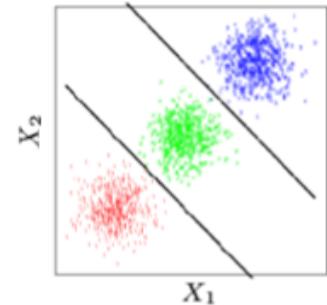
Decision Boundary Points of LDA classifier →

$$= \underbrace{\log \frac{\pi_k}{\pi_\ell} - \frac{1}{2}(\mu_k + \mu_\ell)^T \Sigma^{-1}(\mu_k - \mu_\ell)}_{+ \underbrace{x^T \Sigma^{-1}(\mu_k - \mu_\ell)}_a, = 0} \quad (4.9)$$

b

$$\Rightarrow x^T a + b = 0 \Rightarrow \text{a linear line}$$

decision boundary



## LDA Classification Rule (also called as Linear discriminant function:)

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

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

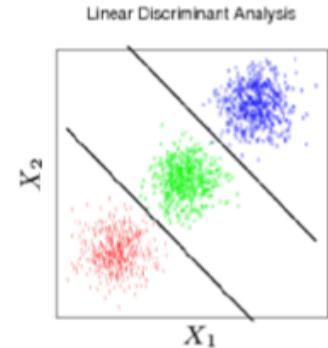
$$= \operatorname{argmax}_k \boxed{-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log(\pi_k)}$$

**Linear Discriminant Function for LDA**

- Note

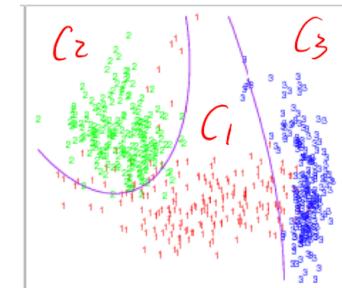
# Not-naïve Gaussian BC

- LDA: Linear Discriminant Analysis



- ■ QDA: Quadratic Discriminant Analysis

Quadratic decision Boundary



(2) If covariance matrix are not the same  
e.g. → QDA (Quadratic Discriminant Analysis)

- ▶ Estimate the covariance matrix  $\Sigma_k$  separately for each class  $k$ ,  
 $k = 1, 2, \dots, K$ .
- ▶ *Quadratic discriminant function:*

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

- ▶ Classification rule:

$$\log p(x|c_k) p(c_k)$$

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

- ▶ Decision boundaries are quadratic equations in  $x$ .
- ▶ QDA fits the data better than LDA, but has more parameters to estimate.

(2) If covariance matrix are not the same

e.g. → QDA (Quadratic Discriminant Analysis)

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

- ▶ Quadratic discriminant function:

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

$\{\Sigma_1, \Sigma_2, \dots, \Sigma_K, \mu_1, \mu_2, \dots, \mu_K\}$

- ▶ Classification rule:

$$\delta_1(x) - \delta_2(x) = 0$$

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

Total # para

$$K \times (P + P^2)$$

$\{\mu_k, \Sigma_k\}$

- ▶ Decision boundaries are quadratic equations in  $x$ .

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

### (3) Regularized Discriminant Analysis

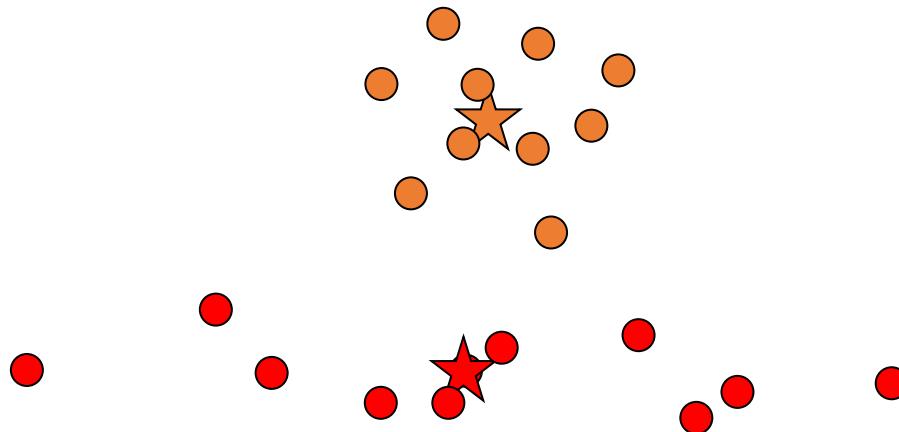
- ▶ A compromise between LDA and QDA.
- ▶ Shrink the separate covariances of QDA toward a common covariance as in LDA.
- ▶ Regularized covariance matrices:

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

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

# More: Decision Boundary of Gaussian naïve Bayes Classifiers ???

Orange Team

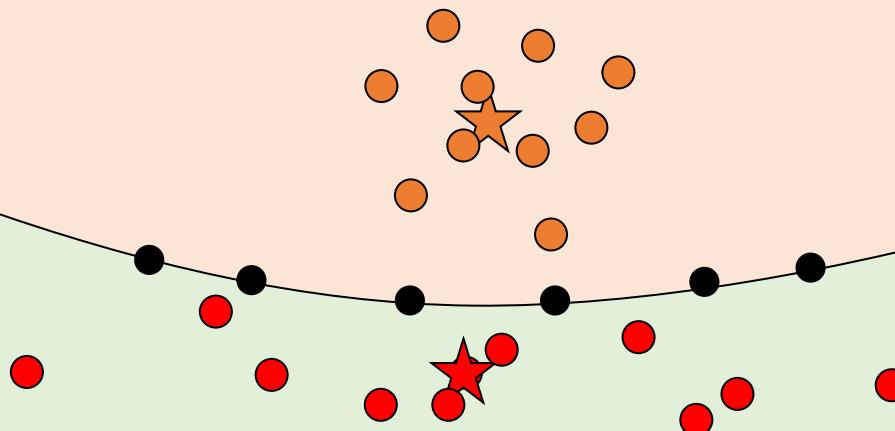


Green Team

Naïve Gaussian Bayes Classifier is  
not a linear classifier!

# Gaussian Naïve Bayes Classifier

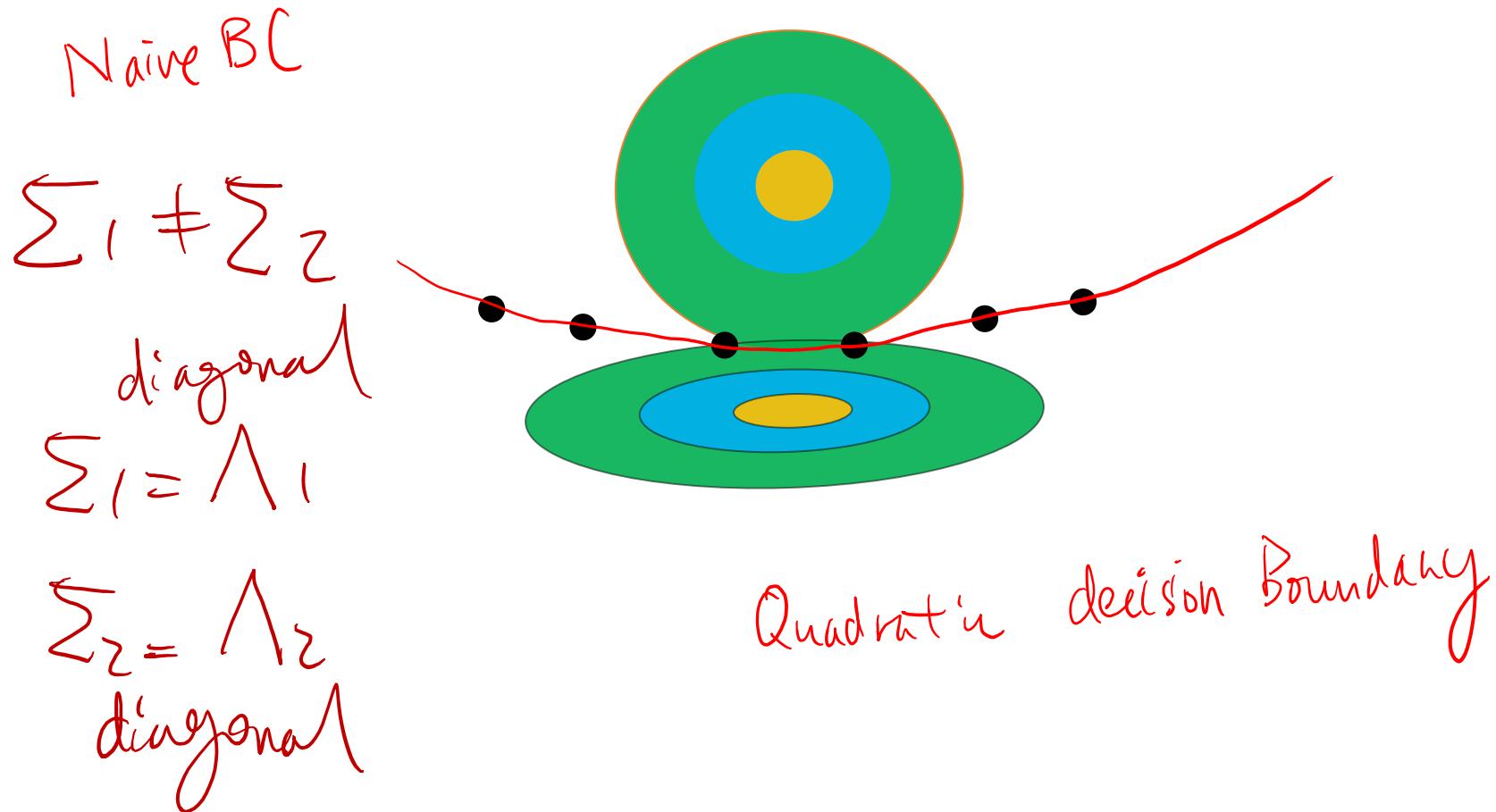
Orange Team



Green Team

Naïve Gaussian Bayes Classifier is  
not a linear classifier!

# Decision Boundary of Gaussian naïve Bayes Classifiers ???



# Today: More Generative Bayes Classifiers

- ✓ Generative Bayes Classifier
- ✓ Naïve Bayes Classifier
- ✓ Gaussian Bayes Classifiers
  - Gaussian distribution
  - Naïve Gaussian BC
  - Not-naïve Gaussian BC → LDA, QDA
- ✓ Discriminative vs. Generative classifier

# Discriminative vs. Generative

Generative approach

- Model the joint distribution  $p(X, C)$  using  
 $\underline{p(X | C = c_k)}$  and  $\underline{p(C = c_k)}$

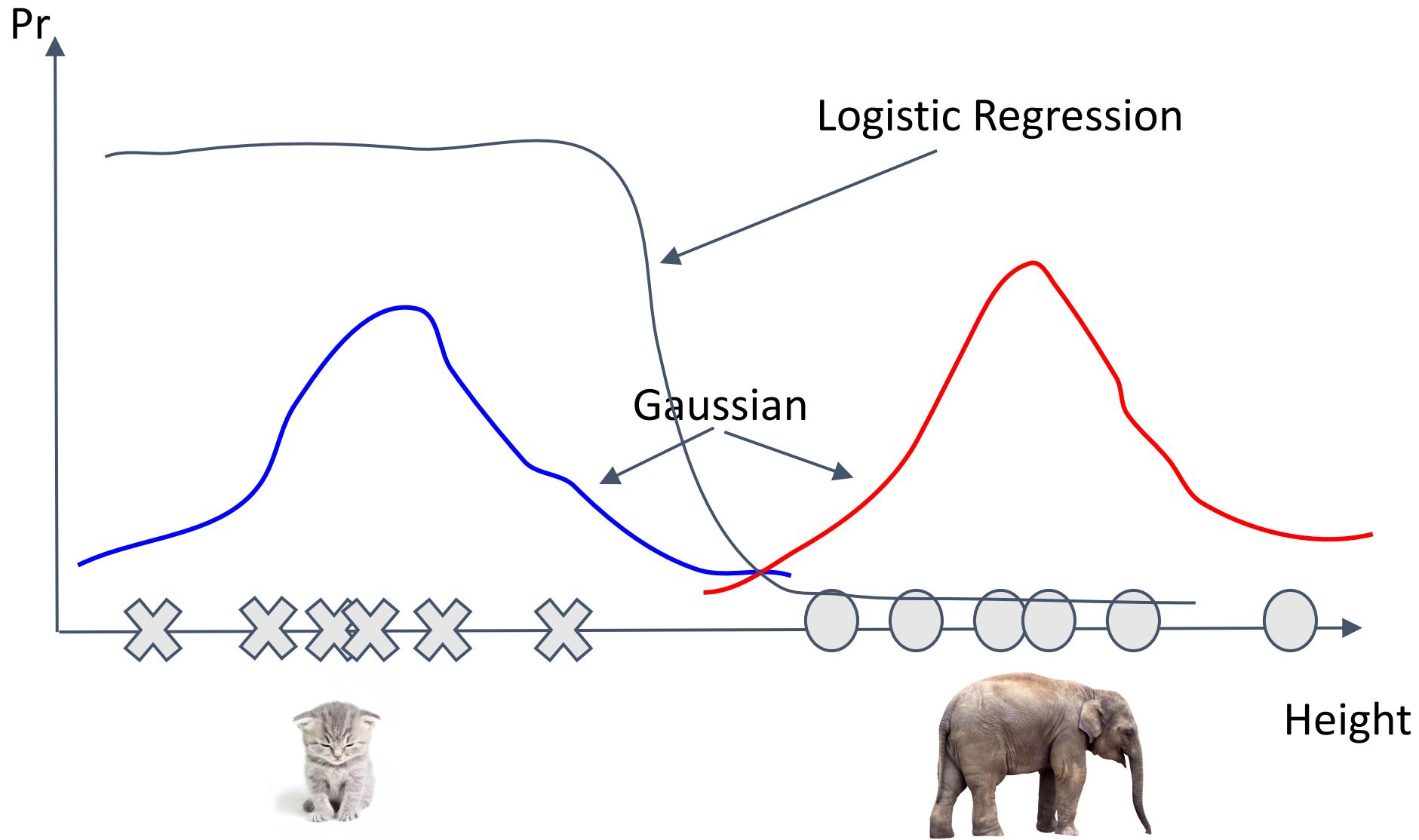
Discriminative approach

- Model the conditional distribution  $p(c | X)$  directly

e.g.,

$$P(C=1 | X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 * X)}}$$

# Discriminative vs. Generative



# LDA vs. Logistic Regression

- **LDA (Generative model)**

*linear*

- Assumes Gaussian class-conditional densities and a common covariance
- Model parameters are estimated by maximizing the full log likelihood, parameters for each class are estimated independently of other classes,  $K p + \frac{p(p+1)}{2} + (K - 1)$  parameters
- Makes use of marginal density information  $\Pr(x)$
- Easier to train, low variance, more efficient if model is correct
- Higher asymptotic error, but converges faster

- **Logistic Regression (Discriminative model)**

*linear*

- Assumes class-conditional densities are members of the (same) exponential family distribution
- Model parameters are estimated by maximizing the conditional log likelihood, simultaneous consideration of all other classes,  $(K - 1)(p + 1)$  parameters
- Ignores marginal density information  $\Pr(x)$
- Harder to train, robust to uncertainty about the data generation process
- Lower asymptotic error, but converges more slowly

# LDA vs. Logistic Regression

## • LDA (Generative model)

- Assumes Gaussian class-conditional densities and a common covariance
- Model parameters are estimated by maximizing the [full log likelihood,]  
parameters for each class are estimated independently of other classes,  
 $K p + \frac{p(p+1)}{2} + (K - 1)$  parameters
- Makes use of marginal density information  $\Pr(x)$
- Easier to train, low variance, more efficient if model is correct
- Higher asymptotic error, but converges faster

$$p(x_{p+1} | c_i)$$

$$\Rightarrow \text{mean } Kp + p^2_{\text{Conv}}$$

## • Logistic Regression (Discriminative model)

- Assumes class-conditional densities are members of the (same) exponential family distribution  $p(c_i|x)$
- Model parameters are estimated by maximizing the [conditional log likelihood]  
simultaneous consideration of all other classes,  $(K - 1)(p + 1)$  parameters
- Ignores marginal density information  $\Pr(x)$
- Harder to train, robust to uncertainty about the data generation process
- Lower asymptotic error, but converges more slowly

$$\Rightarrow (K-1)(p+1)$$

# asymptotic classifiers

- Definitions
  - $h_{gen}$  and  $h_{dis}$ : generative and discriminative classifiers
  - $h_{gen, inf}$  and  $h_{dis, inf}$ : same classifiers but trained on the entire population (asymptotic classifiers)
  - $n \rightarrow \text{infinity}$ ,  $h_{gen} \rightarrow h_{gen, inf}$  and  $h_{dis} \rightarrow h_{dis, inf}$

Ng, Jordan,. "On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes." *Advances in neural information processing systems* 14 (2002): 841.

# Discriminative vs. Generative

Proposition 1:

$$\epsilon(h_{dis,\inf}) \leq \epsilon(h_{gen,\inf})$$

- p : number of dimensions
- n : number of observations
- $\epsilon$  : asymptotic generalization error

Proposition 1 states that asymptotically, the error of the discriminative logistic regression is smaller than that of the generative naive Bayes. This is easily shown

# Logistic Regression vs. Naïve BC

Discriminative classifier (Logistic Regression)

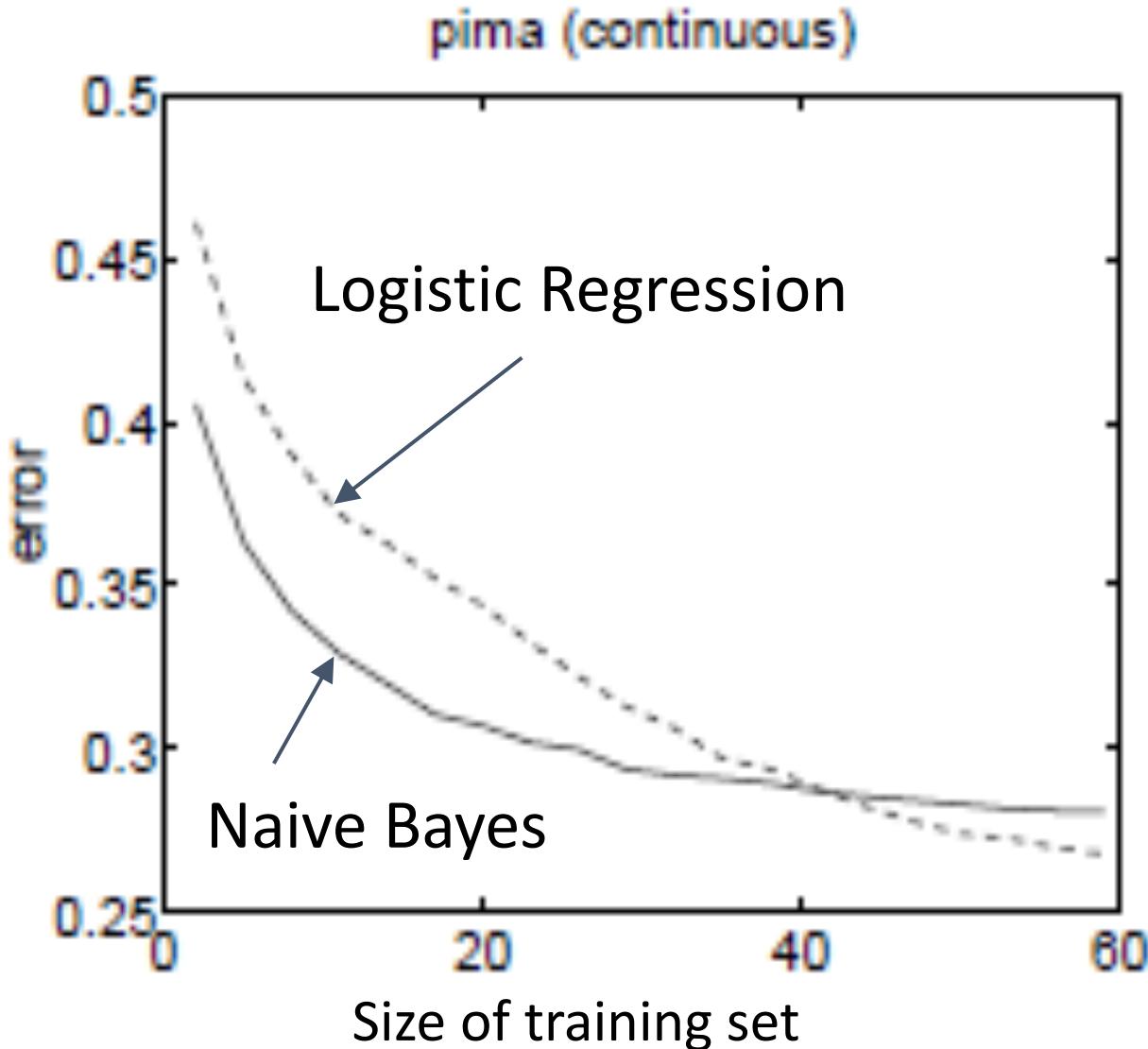
- Smaller asymptotic error
- Slow convergence  $\sim O(p)$

Generative classifier (Naive Bayes)

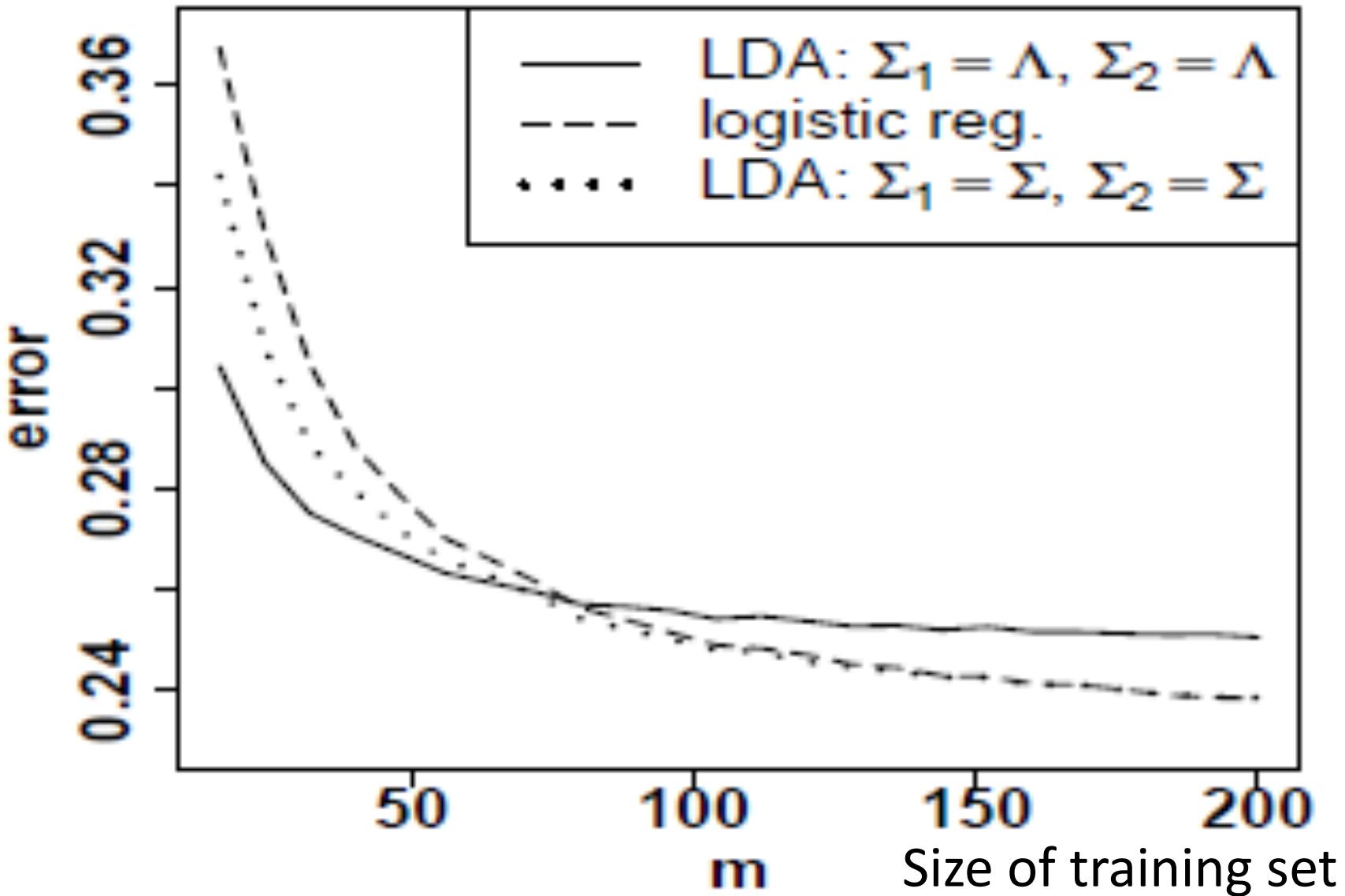
- Larger asymptotic error
- Can handle missing data (EM)
- Fast convergence  $\sim O(\lg(p))$

the speed at which a convergent sequence approaches its limit is called the rate of convergence.

Ng, Jordan,. "On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes." *Advances in neural information processing systems* 14 (2002): 841.



# Logistic regression / vs. Naïve LDA / vs. LDA



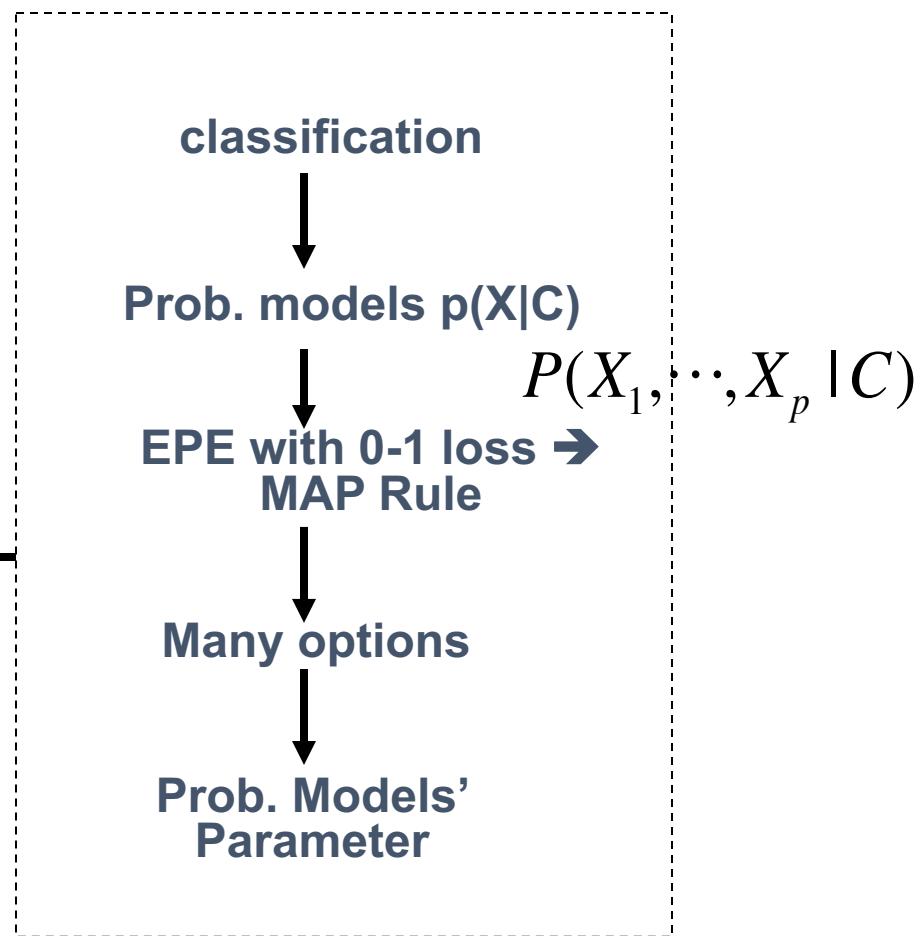
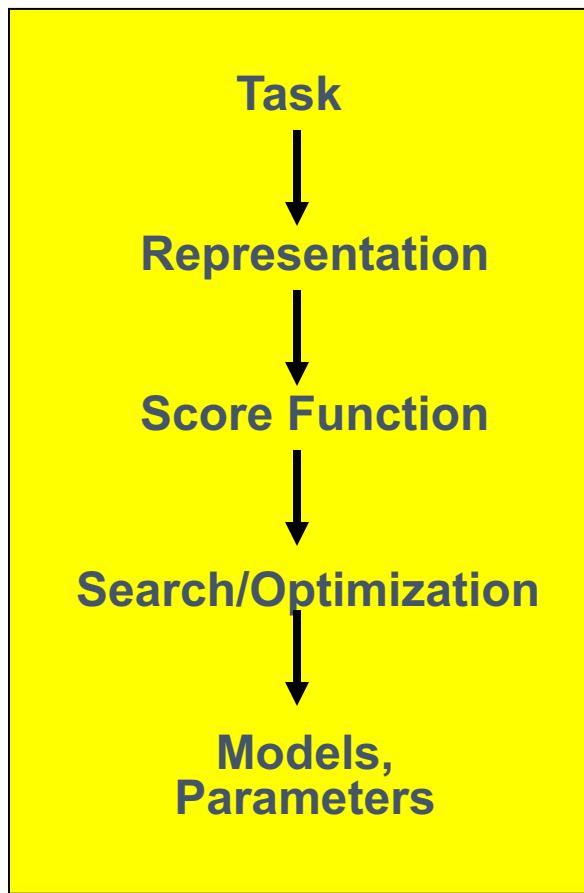
Xue, Jing-Hao, and D. Michael Titterington. "Comment on ‘On discriminative vs. generative classifiers: A comparison of logistic regression and naive Bayes’." *Neural processing letters* 28.3 (2008): 169-187.

# Summary: Discriminative vs. Generative

- Empirically, **generative** classifiers approach their asymptotic error faster than discriminative ones
  - Good for small training set
  - Handle missing data well (EM)
- Empirically, **discriminative** classifiers have lower asymptotic error than generative ones
  - Good for larger training set

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

## 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. Tan, Steinbach, Kumar's "Introduction to Data Mining" slide
- Prof. Andrew Moore's slides
- Prof. Eric Xing's slides
- Prof. Ke Chen NB slides
- Hastie, Trevor, et al. *The elements of statistical learning*. Vol. 2. No. 1. New York: Springer, 2009.