



# Machine Learning

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

Dr. Beilun Wang

Southeast University  
School of Computer Science  
and Engineering

# 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  $X_1, \dots, X_p$

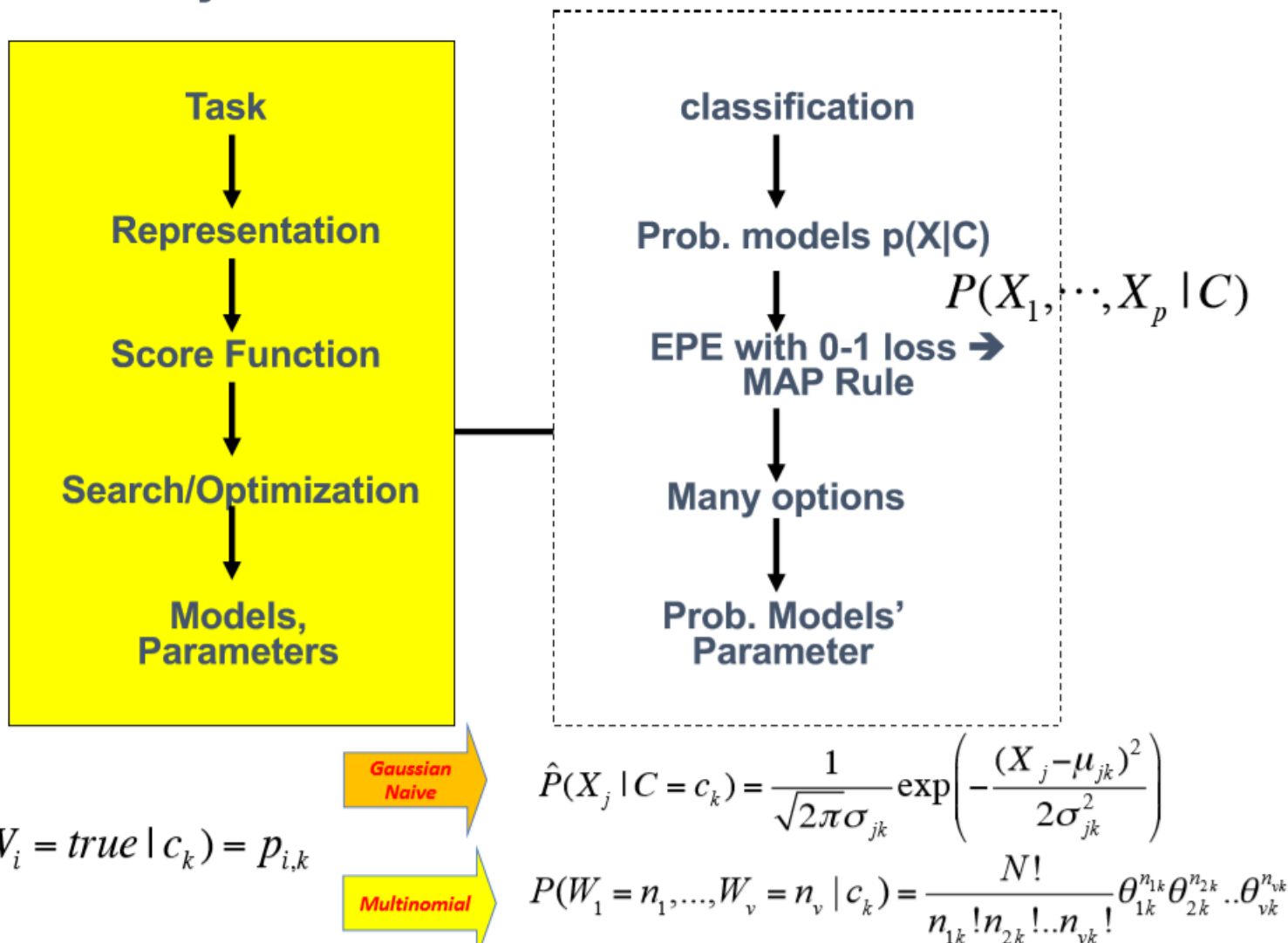
Learn program to Interact with its environment

# Today: More Generative Bayes Classifiers

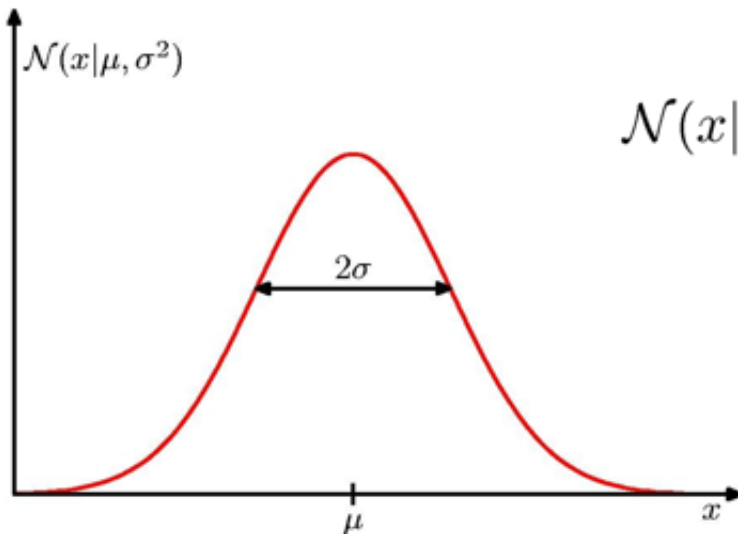
- Generative Bayes Classifier
- Naïve Bayes Classifier
- • Gaussian Bayes Classifiers
  - Gaussian distribution
  - Naïve Gaussian BC
  - Not-naïve Gaussian BC → LDA, QDA
    - LDA: Linear Discriminant Analysis
    - QDA: Quadratic Discriminant Analysis
- Extra: Discriminative vs. Generative classifier

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

## Generative Bayes Classifier



# Review: Single-Variate Gaussian Distribution



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

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

# Multivariate Normal (Gaussian) PDFs

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

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

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

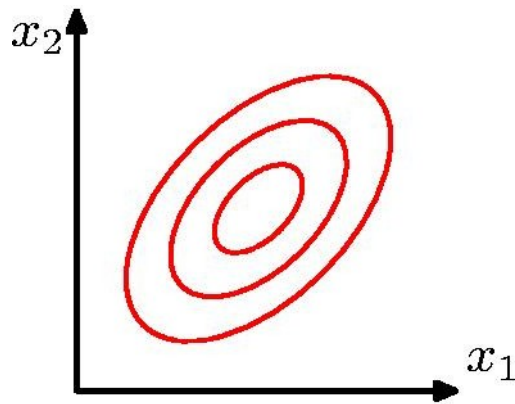
$$f(x_1, x_2) = \frac{1}{(2\pi)|\Sigma|^{1/2}} e^{-\frac{1}{2}(\vec{x}-\vec{\mu})^T \Sigma^{-1}(\vec{x}-\vec{\mu})}$$

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

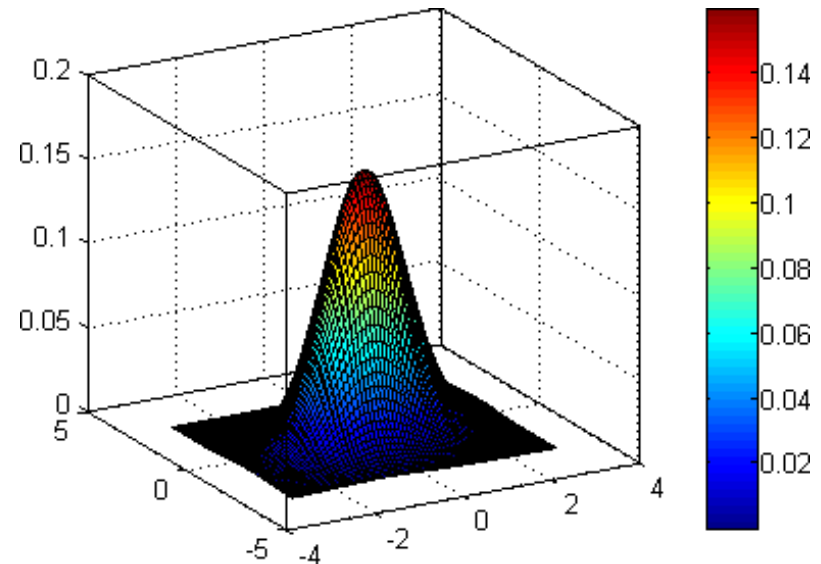
$$\Sigma_{2 \times 2} = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}_{2 \times 2}$$

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

# Bi-Variate Gaussian Distribution



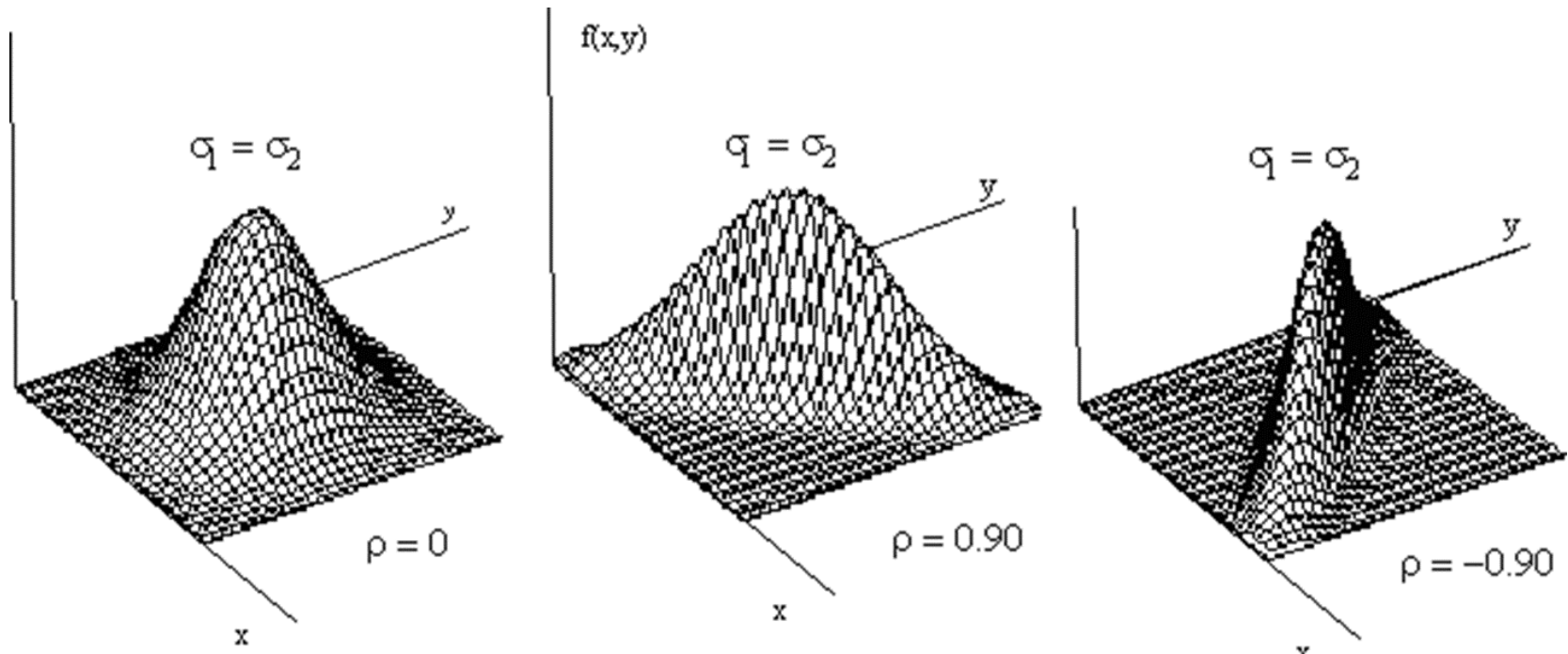
Bivariate normal PDF



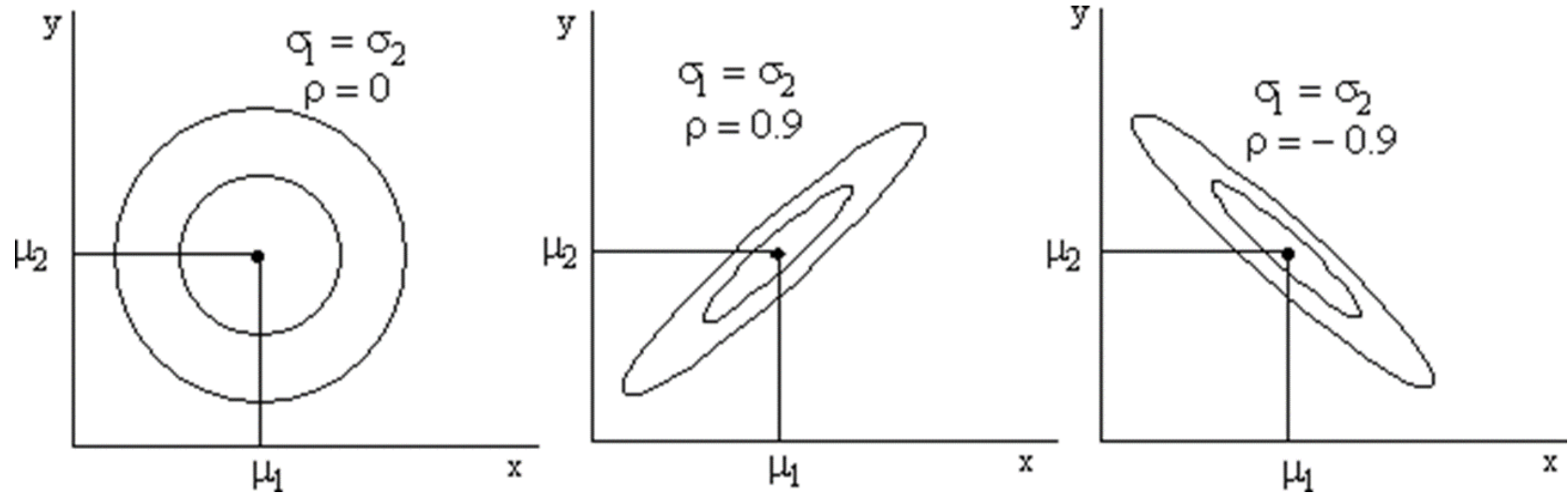
- Mean of normal PDF is at peak value. Contours of equal PDF form ellipses.
- The covariance matrix captures linear dependencies among the variables



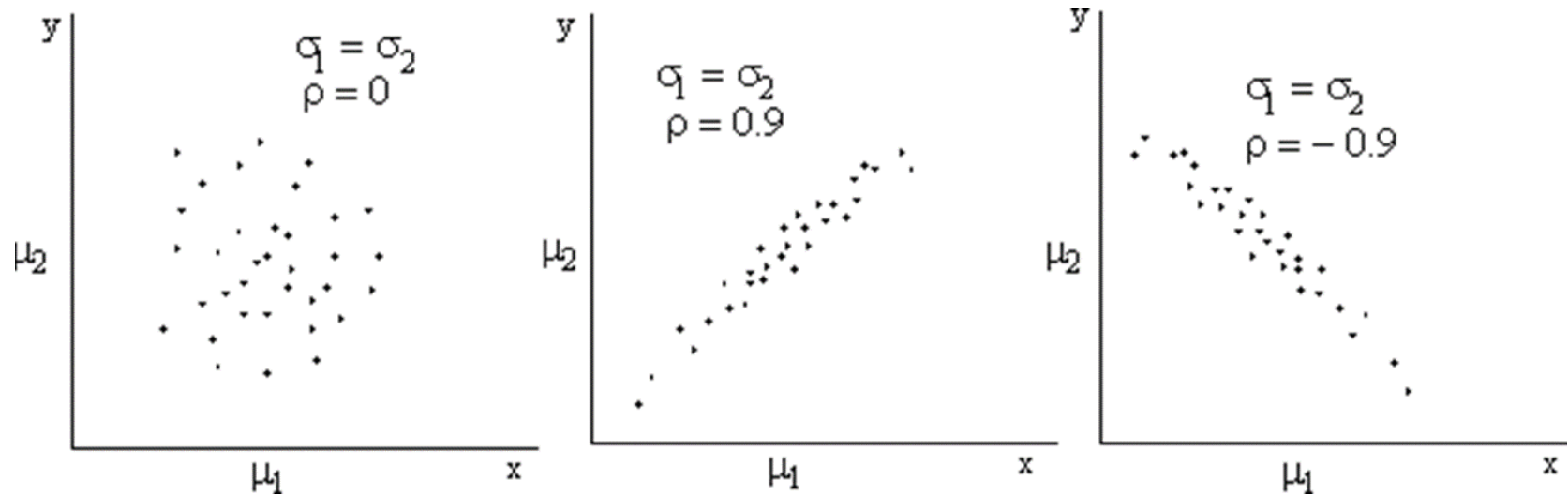
# Surface Plots of the bivariate Normal distribution



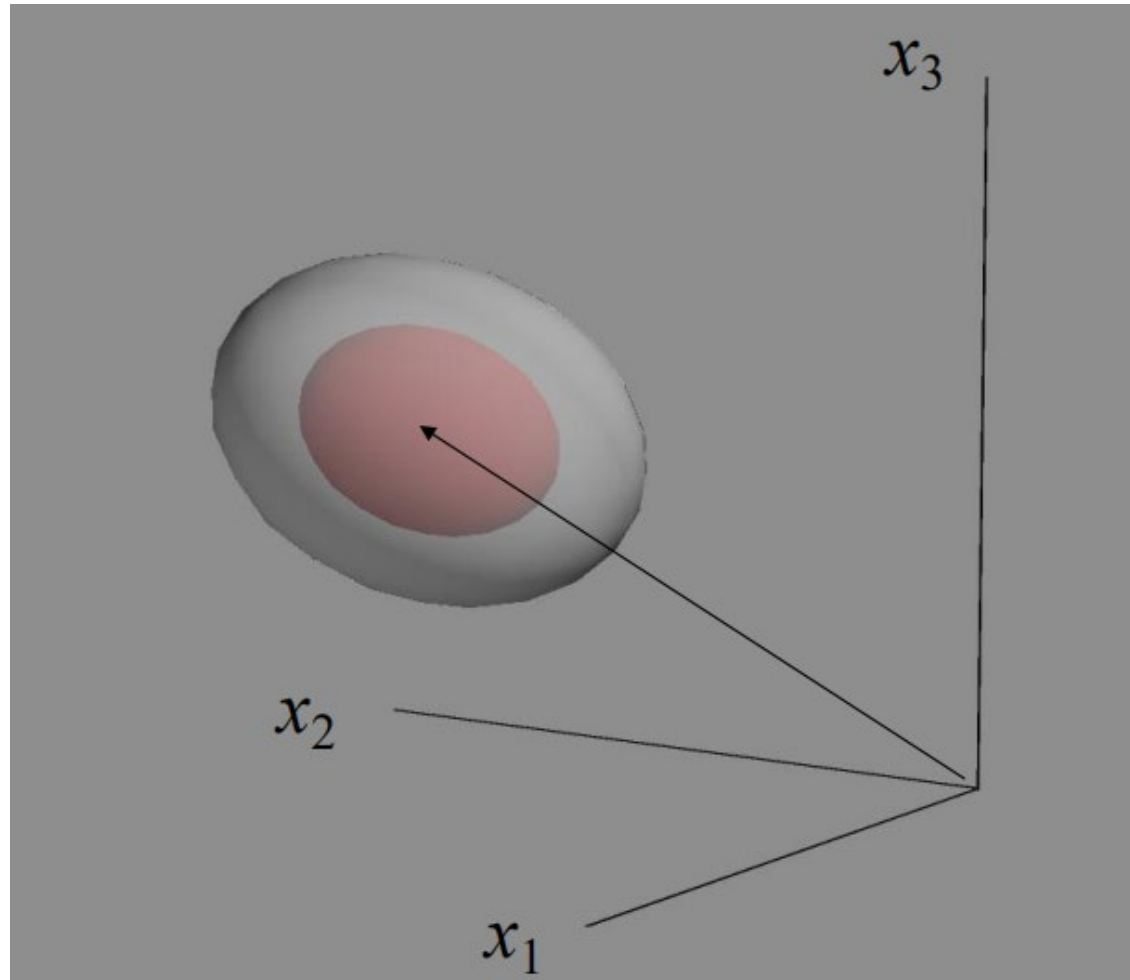
# Contour Plots of the bivariate Normal distribution



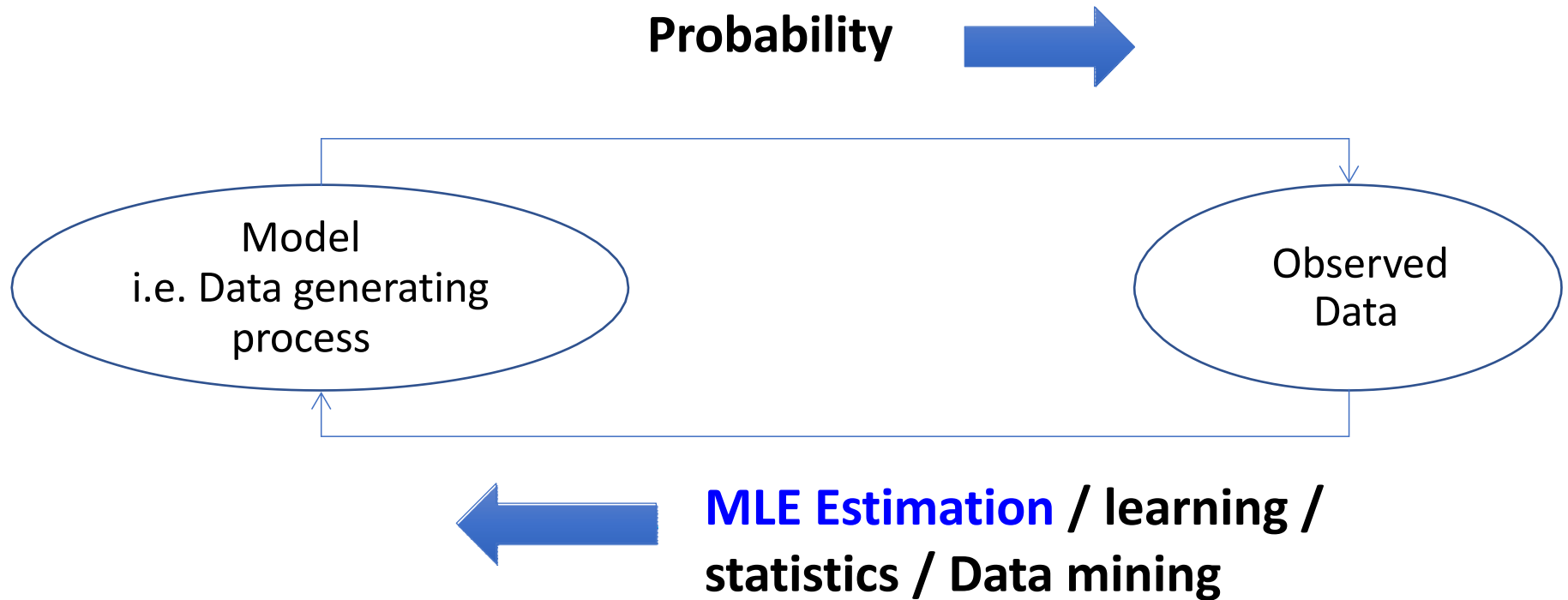
# Scatter Plots of the bivariate Normal distribution



# Trivariate Normal distribution

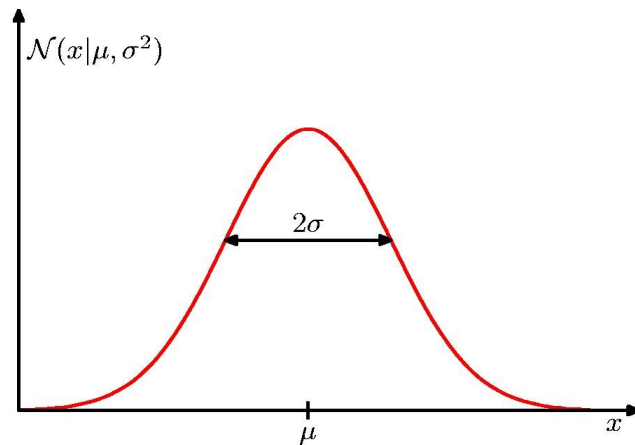


# The Big Picture



# 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

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

$$\vec{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_p \end{bmatrix} \quad \Sigma_{p \times p} = \begin{bmatrix} \text{var}(X_1) & \dots & \text{cov}(X_1, X_p) \\ \vdots & \ddots & \vdots \\ \text{cov}(X_p, X_1) & \dots & \text{var}(X_p) \end{bmatrix}$$

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i \quad \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})(x_i - \hat{\mu})^T$$

# 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: 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
    - LDA: Linear Discriminant Analysis
    - QDA: Quadratic Discriminant Analysis
- Extra: Discriminative vs. Generative classifier

# 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 (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  $C = c_i$

- **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

- 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 (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  $C = c_i$

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

# Gaussian Naïve Bayes Classifier

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

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

Each class' covariance matrix is diagonal

11/6/19 Dr. Yanjun Qi / UVA CS

# 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
- Extra: Discriminative vs. Generative classifier

# 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

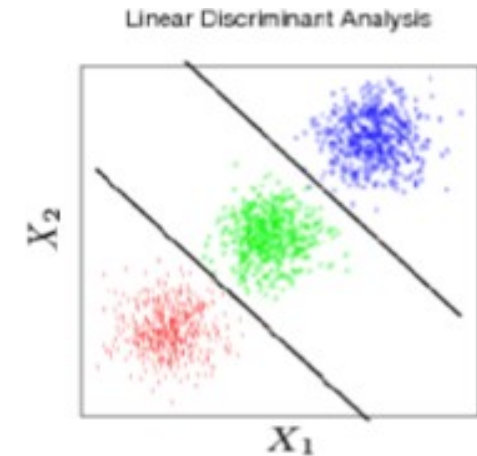
$$\boldsymbol{\Sigma} \mathbf{c}_k = \boldsymbol{\Lambda} \mathbf{c}_k$$

Each class' covariance matrix is diagonal

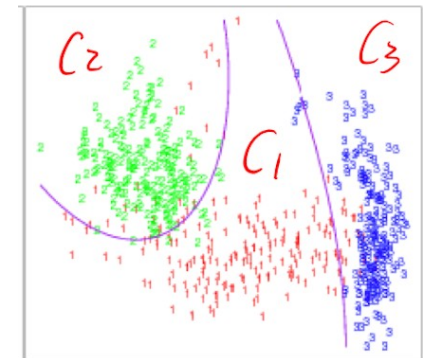


# Not-naïve Gaussian BC

- LDA: Linear Discriminant Analysis



- QDA: Quadratic Discriminant Analysis



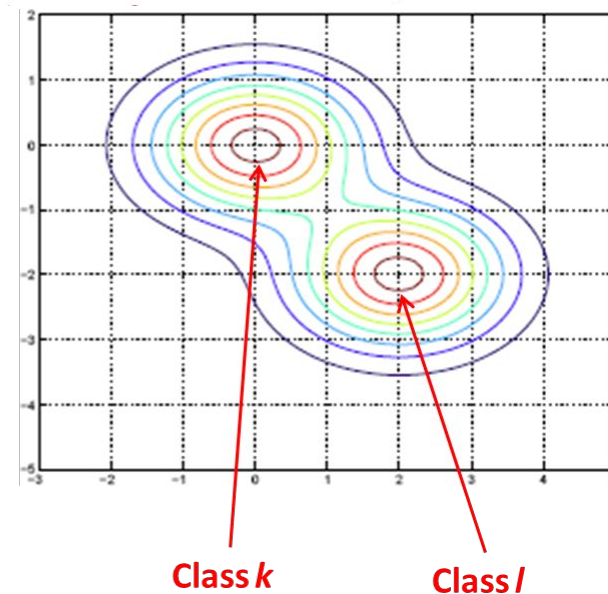
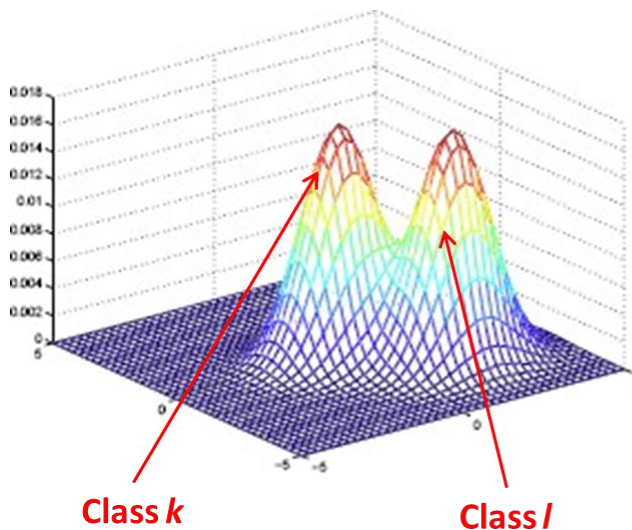
# covariance matrix are the same across classes

- 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



$$\begin{aligned}\arg\max_k P(C_k | X) &= \arg\max_k P(X, C_k) = \arg\max_k P(X | C_k) P(C_k) \\ &= \arg\max_k \log\{P(X | C_k) P(C_k)\}\end{aligned}$$

Decision Boundary Points  
satisfying:

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

$$\begin{aligned}\arg\max_k P(C_k | X) &= \arg\max_k P(X, C_k) = \arg\max_k P(X | C_k) P(C_k) \\ &= \arg\max_k \log\{P(X | C_k) P(C_k)\}\end{aligned}$$

$$= \arg\max_k \log P(x | C_k) + \log P(C_k) \longrightarrow \pi_k$$

### Decision Boundary Points

$$\begin{aligned}\log \frac{P(C_k | X)}{P(C_l | X)} = 0 &= \log \frac{P(X | C_k)}{P(X | C_l)} + \log \frac{\pi_k}{\pi_l} \\ &= \log P(X | C_k) - \log P(X | C_l) + \log \frac{\pi_k}{\pi_l}\end{aligned}$$

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

Decision Boundary Points of LDA classifier →

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

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 →

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

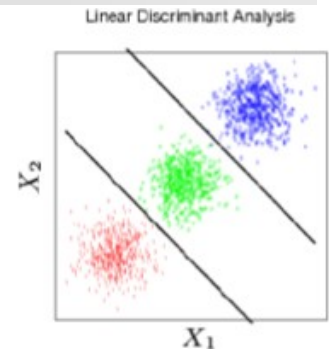
*b*

*a*

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

# LDA Classification Rule

- Also called as Linear discriminant function



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

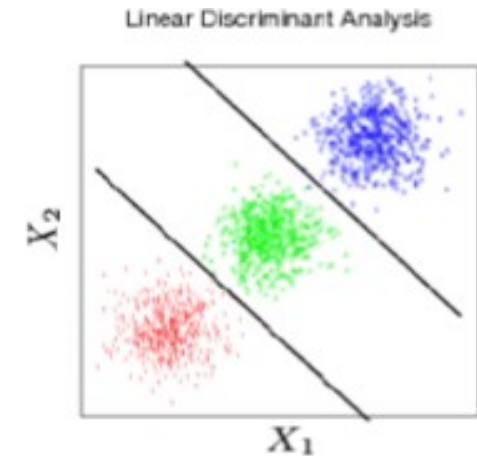
$$= \operatorname{argmax}_k \left[ -\log((2\pi)^{p/2} |\Sigma|^{1/2}) - \frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log(\pi_k) \right]$$

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

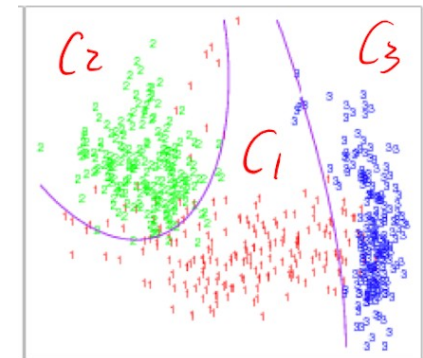
**Linear Discriminant Function for LDA**

# Not-naïve Gaussian BC

- LDA: Linear Discriminant Analysis



- ➔ • QDA: Quadratic Discriminant Analysis





# If covariance matrix are not the same

- QDA (Quadratic Discriminant Analysis)

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

- ▶ Quadratic discriminant function:

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

- ▶ Classification rule:

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

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

# 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:

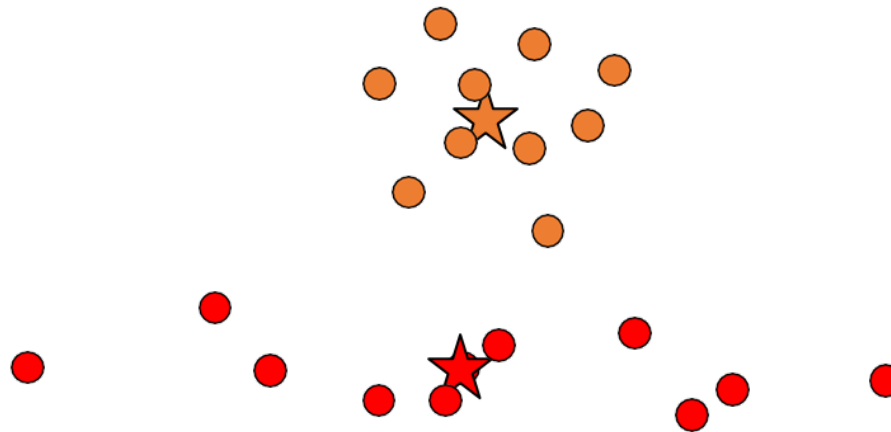
$$\hat{\Sigma}_k(\alpha) = \alpha \hat{\Sigma}_k + (1 - \alpha) \hat{\Sigma}.$$

- ▶ The quadratic discriminant function  $\delta_k(x)$  is defined using the shrunk 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



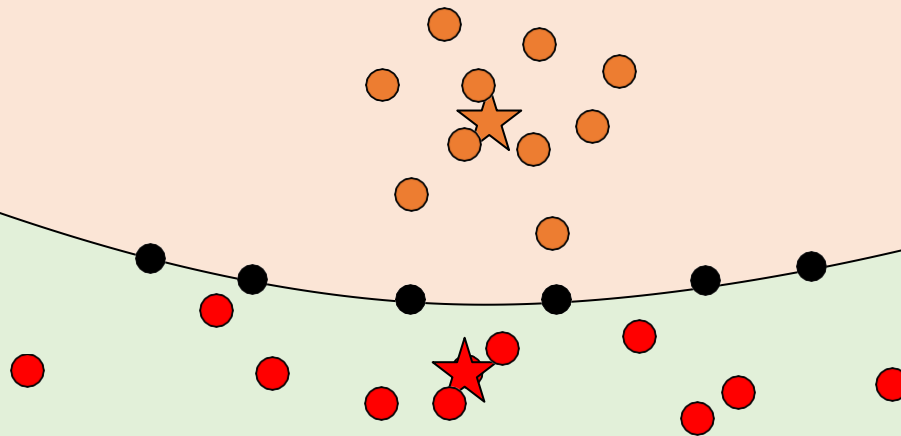
Red Team

Naïve Gaussian Bayes Classifier is not a linear classifier!

# Gaussian Naïve Bayes Classifier

Orange Team

Red Team



Naïve Gaussian Bayes Classifier is not a linear classifier!

# 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
- ➔ • Extra: Discriminative vs. Generative classifier

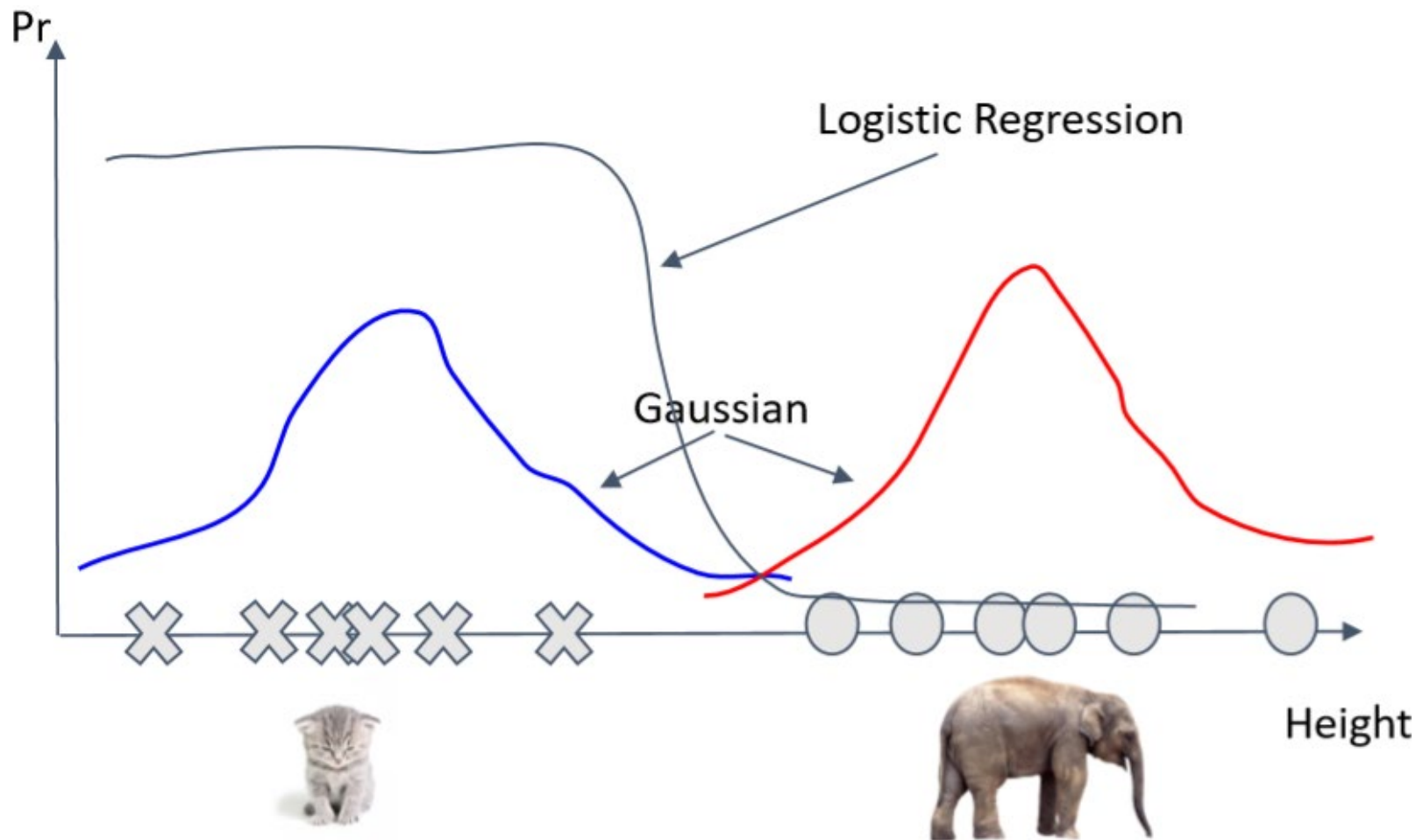
# Discriminative vs. Generative

- Generative approach
  - Model the joint distribution  $p(X, C)$  using  $p(X | C = c_k)$  and  $p(C = c_k)$
- 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)

- 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,  
 $Kp + \frac{p(p+1)}{2} + (K - 1)$  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)

- 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

- 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.



# 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

# References

- <https://qiyanjun.github.io/2019f-UVA-CS6316-MachineLearning/>
- Prof. Tan, Steinbach, Kumar's "Introduction to Data Mining" slide
- Prof. Andrew Moore's slides
- Prof. Eric Xing's slides
- Prof. KeChen NB slides qHastie, Trevor, et al. The elements of statistical learning. Vol. 2. No. 1. New York: Springer, 2009.



*Thanks for listening*