

Math 7243

Machine Learning and Statistical Learning Theory

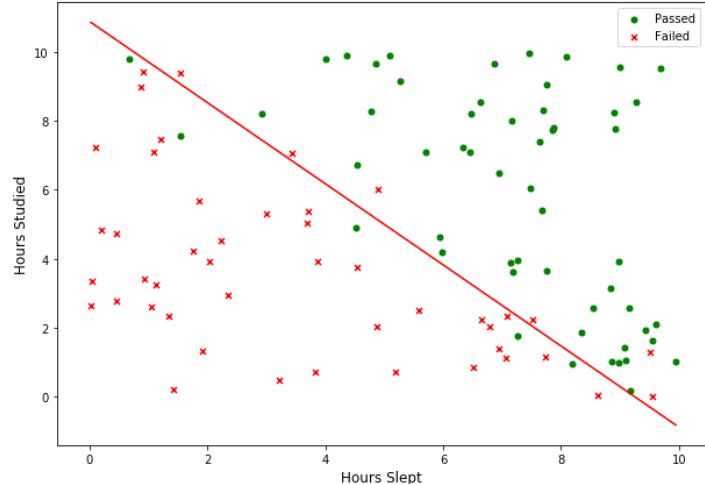
## Section 7 Gaussian Discriminant Analysis

Instructor: He Wang

Department of Mathematics

Northeastern University

## Review :



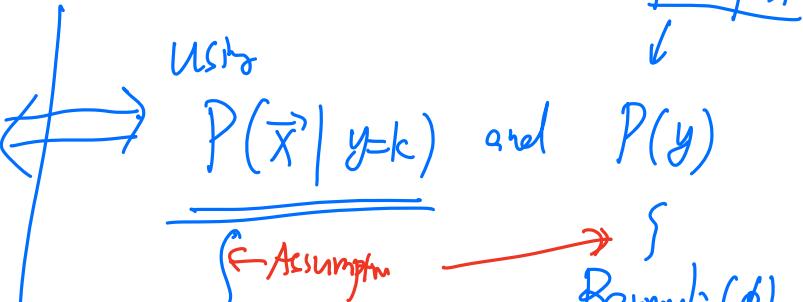
## Today:

1. Gaussian Discriminant Analysis
2. Review Gaussian distribution
3. Linear Discriminant Analysis (LDA)  $\approx$  QDA
4. LDA v.s. Logistics regression

➤ Discriminative learning algorithms.

Data  $(\vec{x}^{(i)}, y^{(i)}) \rightarrow \text{classification}$

Model 1:  $P(y|\vec{x}; \theta)$



➤ Generative learning algorithms.

$$P(y=k | \vec{x}) = \frac{P(y=k, \vec{x})}{P(\vec{x})}$$

Posterior prob.

$$= \frac{P(\vec{x}|y=k)P(y=k)}{\sum_i P(\vec{x}|y=i)P(y=i)}$$

Bayes' Thm

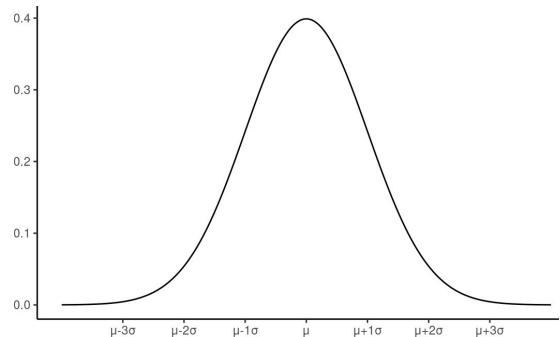
$$\arg \max_{\vec{y}} P(\vec{y} | \vec{X}) = \arg \max_{\vec{y}} (P(\vec{x} | \vec{y}) P(\vec{y}))$$

➤ Normal distribution.

Polyt  $f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$  for  $-\infty < x < \infty$ .

The mean is  $E(X) = \mu$ .

The variance is  $\text{Var}(X) = \sigma^2$ .



➤ Multivariate normal distribution.

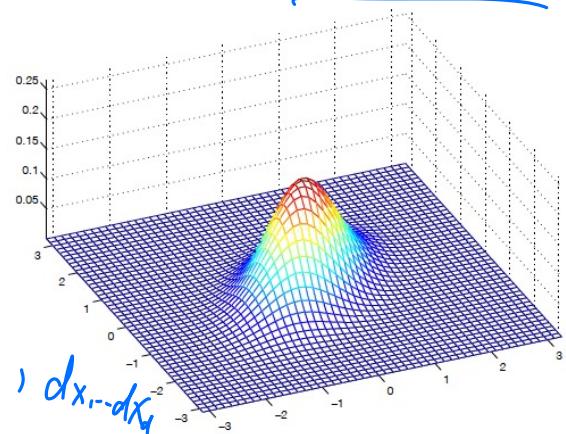
$$\vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

Polyt  $f(\vec{x}; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left( -\frac{1}{2} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu}) \right)$

$$\vec{x} \in \mathbb{R}^d \quad \vec{\mu} \in \mathbb{R}^d$$

$$\cdot \int_{-\infty}^{\infty} \int f \, dx_1 \cdots dx_d = 1$$

$\Sigma \in \mathbb{R}^{d \times d}$   
• symmetric  
• positive definite



①  $\vec{\mu} \in \mathbb{R}^d \quad \vec{\mu} = E(\vec{X}) = \int_{-\infty}^{\infty} \int \vec{x} f(\vec{x}) \, d\vec{x}_1 \cdots d\vec{x}_d$

$$\textcircled{2} \quad \Sigma = \text{cov}(\vec{X}) \in \mathbb{R}^{d \times d} \quad \text{Covariance matrix of } \vec{X}$$

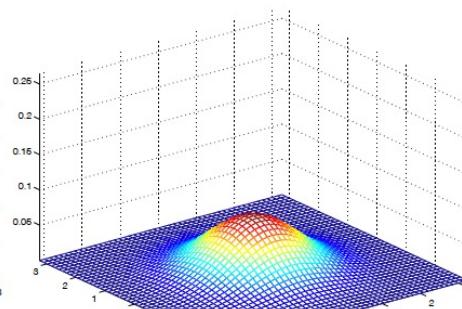
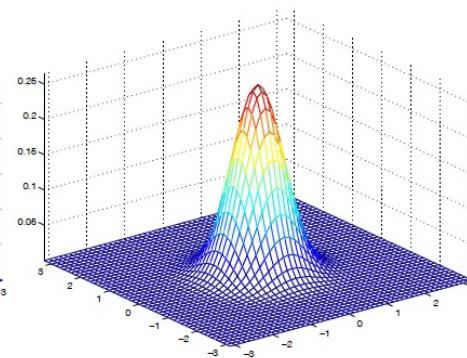
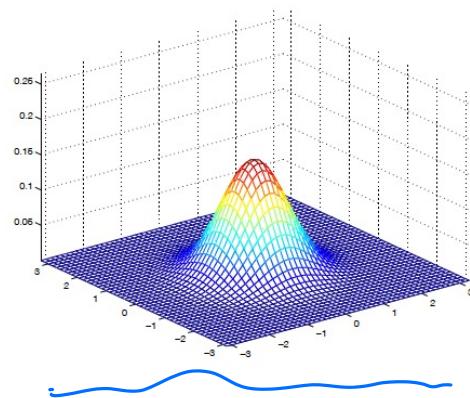
$$X \sim \text{Normal} (\vec{\mu}, \Sigma)$$

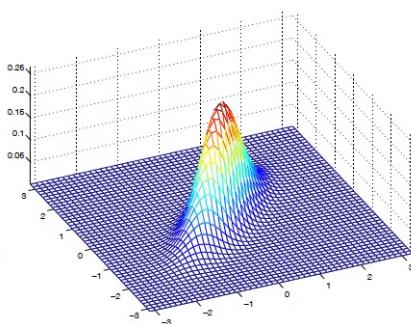
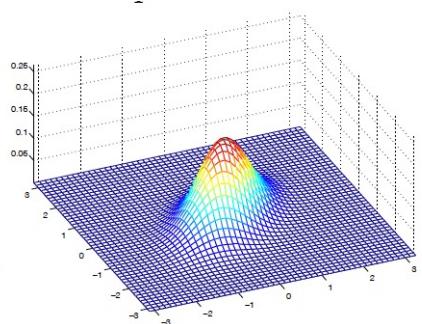
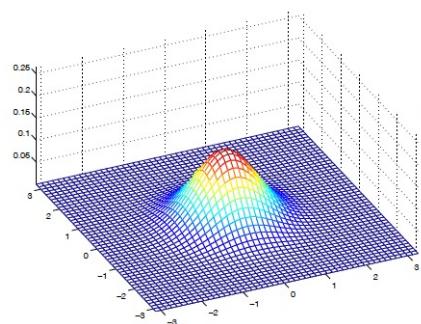
$$\text{Cov}(X) = E(X X^T - E(X)E(X)^T)$$

$$\vec{\mu} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \alpha^6 & 0 \\ 0 & \alpha^6 \end{bmatrix}$$

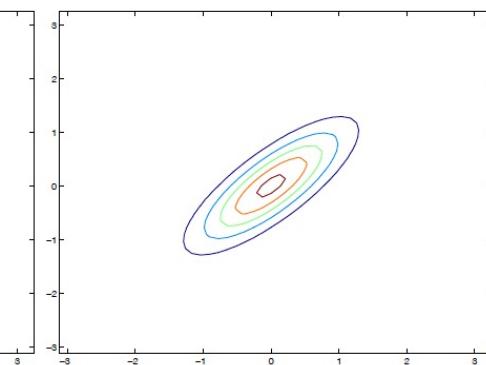
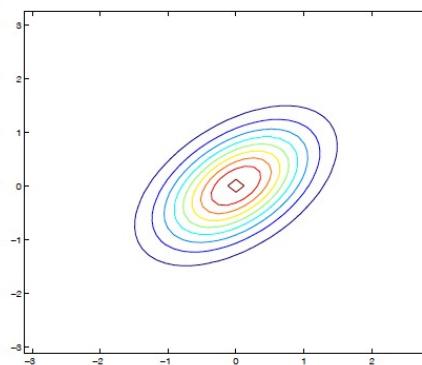
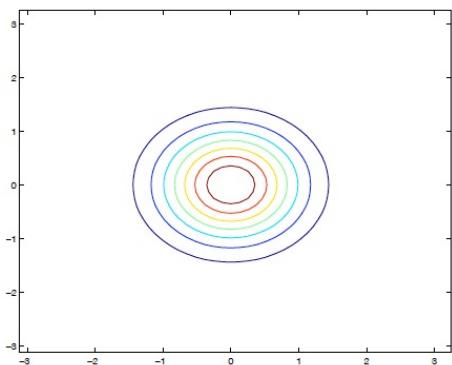
$$\Sigma = 2 I$$

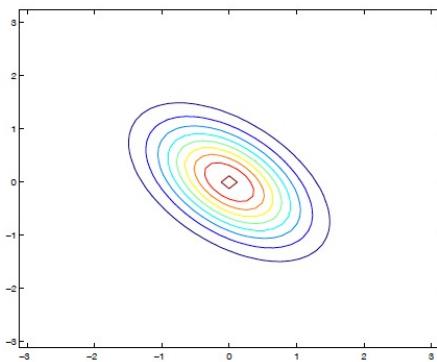




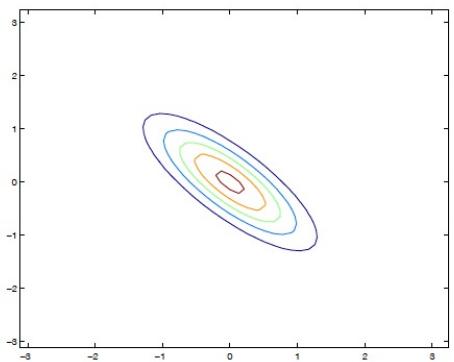
$$\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

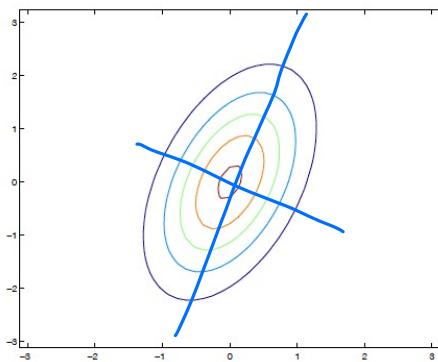




$$\Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$$

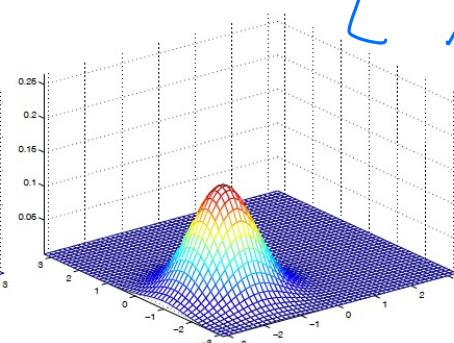
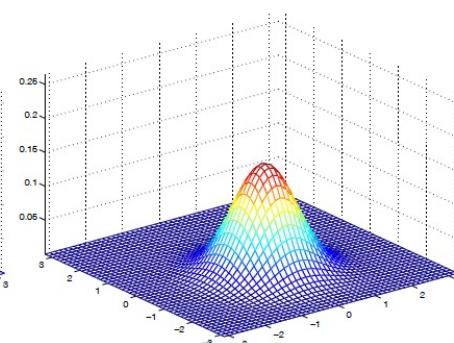
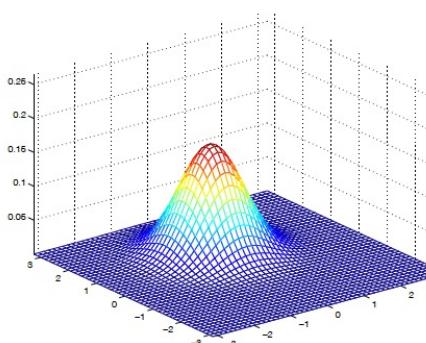


$$\begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}$$



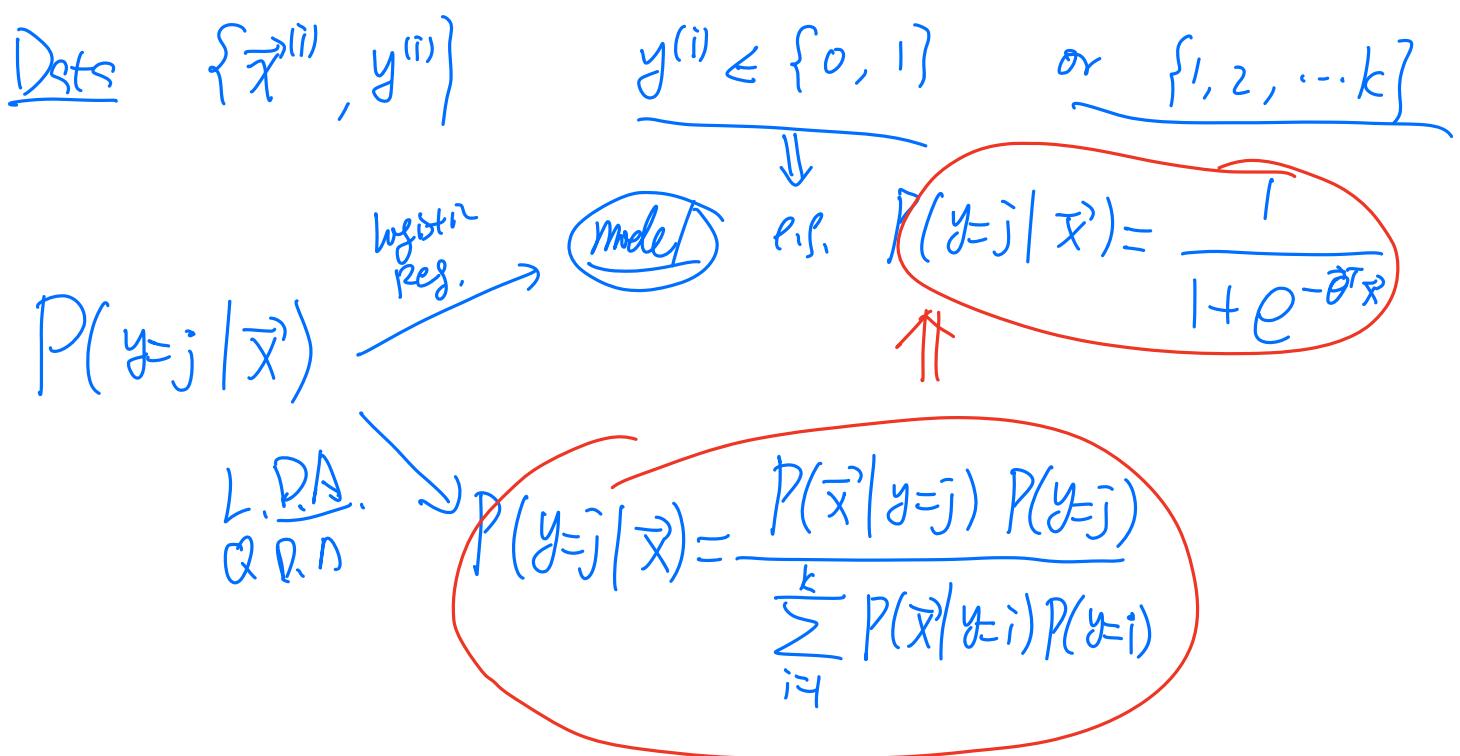
$$\Sigma = P D P^T$$

$P \in \mathbb{R}^{n \times n}$



$$D \in \mathbb{R}^{n \times n}$$

$$\begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$$



Assumption :  $\vec{x} | y=j \sim \text{Normal}(\vec{\mu}_j, \Sigma_j)$

②  $y \sim \text{Bernoulli}(\phi)$

or Categorical  $(\phi_1, \dots, \phi_k)$

(LDA)  
 Assumptions:  $\Sigma_j = \Sigma$  for all  $j$

$$\phi_1 + \dots + \phi_k = 1$$

➤ Gaussian Linear Discriminant Analysis(LDA).

$$\cdot \quad y \sim \text{Bernoulli}(\phi)$$

$$P(y) = \begin{cases} \phi & \text{if } y=1 \\ 1-\phi & \text{if } y=0 \end{cases}$$

$$\mathbb{1} (\text{true}) = 1$$

$$\mathbb{1} (\text{false}) = 0$$

$$= \phi^y (1-\phi)^{1-y}$$

$$= \phi^{\mathbb{1}(y=1)} (1-\phi)^{\mathbb{1}(y=0)}$$

$$P(\vec{x} | y=0) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left( -\frac{1}{2} (x - \mu_0)^T \underline{\Sigma^{-1}} (x - \mu_0) \right)$$

$$P(\vec{x} | y=1) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \left( -\frac{1}{2} (x - \mu_1)^T \underline{\Sigma^{-1}} (x - \mu_1) \right)$$

Maximize the likelihood function

$$L(\phi, \vec{\mu}_0, \vec{\mu}_1, \Sigma) \stackrel{\text{I.i.d.}}{=} \prod_{i=1}^m P(\vec{x} = \vec{x}^{(i)}, y = y^{(i)}) = \underbrace{\prod_{i=1}^m P(\vec{x} = \vec{x}^{(i)} | y = y^{(i)})}_{\text{Marginalization}} P(y = y^{(i)})$$

$$\begin{aligned}\ell(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^m p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \\ &= \log \prod_{i=1}^m p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi)\end{aligned}$$

$\frac{\partial \ell}{\partial \phi} = 0 \quad \frac{\partial \ell}{\partial \mu_i} = 0 \quad \frac{\partial \ell}{\partial \Sigma} = 0 \quad \Rightarrow \quad y \in \{0, 1\}$

$$= \sum_{i=1}^m \left( \log P(\vec{x}^{(i)} | y^{(i)}) + \log P(y^{(i)}) \right)$$

$P(\vec{x} | y=0) = 90\%$

$$\phi = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{y^{(i)} = 1\} \quad \text{mean of } y$$

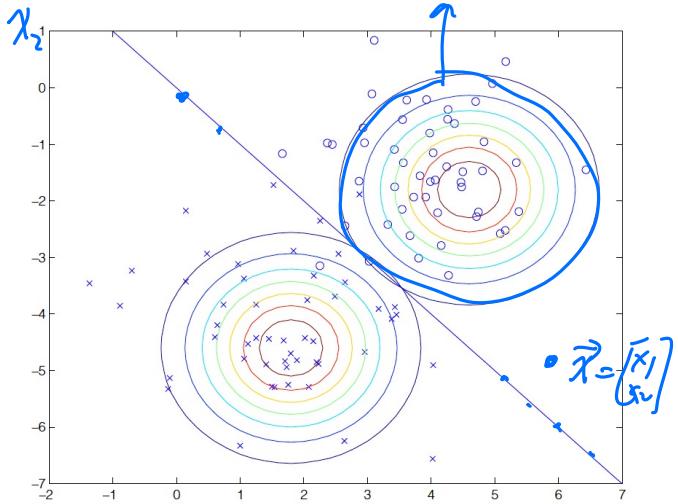
$$\mu_0 = \frac{\sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\} \vec{x}^{(i)}}{\sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\}} \quad \text{mean of } \vec{x}^{(i)}, 0$$

$$\mu_1 = \frac{\sum_{i=1}^m \mathbb{1}\{y^{(i)} = 1\} \vec{x}^{(i)}}{\sum_{i=1}^m \mathbb{1}\{y^{(i)} = 1\}} \quad \text{mean of } \vec{x}^{(i)}, 1$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (\vec{x}^{(i)} - \mu_{y^{(i)}})(\vec{x}^{(i)} - \mu_{y^{(i)}})^T$$

$m \leftarrow$  Sample cov.

Determine  $\vec{x}^{(i)} \quad y^{(i)} \quad i=1, \dots, m$

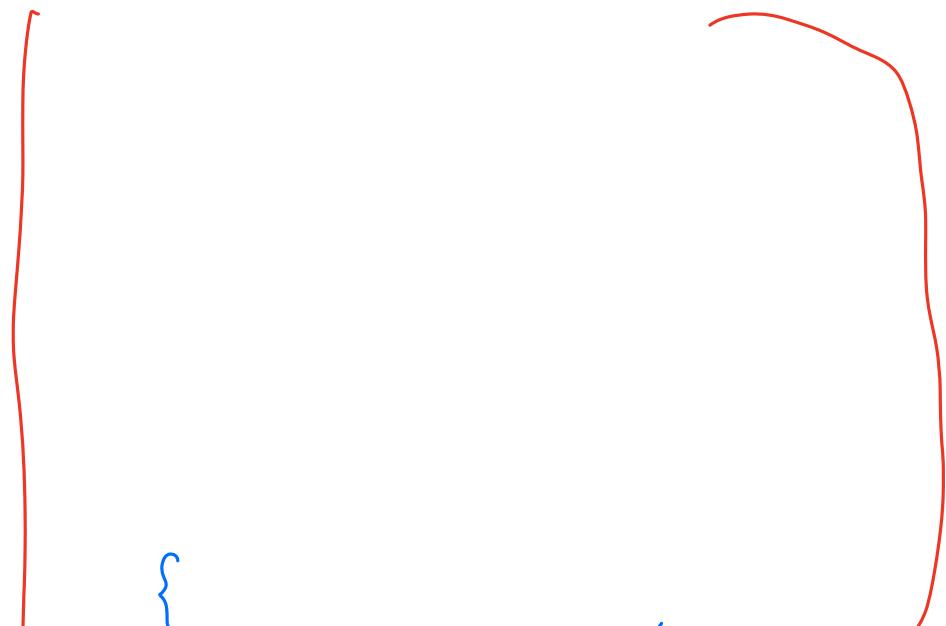


$$\textcircled{0.5} = P(y=j | \vec{x}) = \frac{P(\vec{x} | y=j) P(y=j)}{\sum_{i=1}^k P(\vec{x} | y=i) P(y=i)}$$

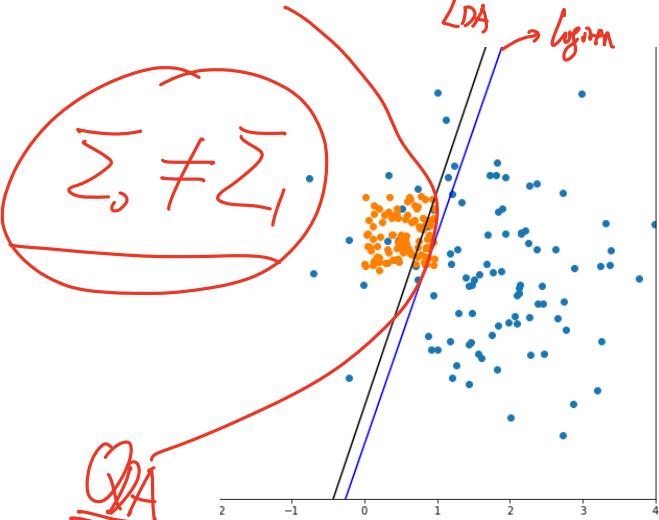
$$\Leftrightarrow P(y_0 | \vec{x}) = P(y_1 | \vec{x})$$

$$\Leftrightarrow \log \frac{P(y_0 | \vec{x})}{P(y_1 | \vec{x})} = 0$$

$\Rightarrow$  boundary  $f(\vec{x}) = 0$  linear



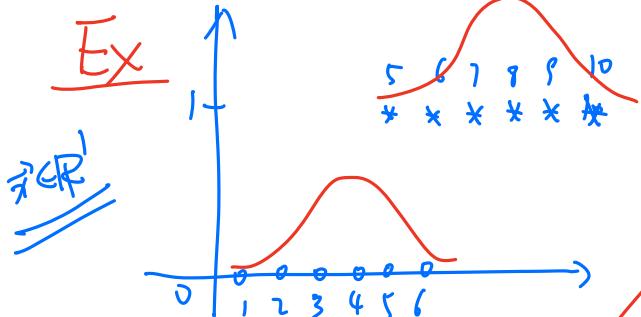
- Compare to Logistic regression



When these modeling assumptions are correct, then GDA will be better fits to the data.

(GDA) will be a better algorithm than logistic regression for small training set sizes. LDA  $\sigma_1 = \sigma_2 = \sigma$

logistic regression makes significantly weaker assumptions. So it is more robust and less sensitive to incorrect modeling assumptions. GDA



Assume

$$P(x|y=i) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}(\frac{x-\mu_i}{\sigma})^2}$$

$i=0, 1$

$$P(y=i) = \phi^i (1-\phi)^{1-i}$$

$$P(y=1|x) = \frac{P(x|y=1) P(y=1)}{P(x|y=0) P(y=0) + P(x|y=1) P(y=1)} = \underline{\hspace{10cm}}$$

Remark 1

$$y \in \{-1, 1\}$$

$$P(y=j) = \phi_j \begin{cases} 1 & (j=1) \\ 1-\phi_j & (j=2) \end{cases}$$

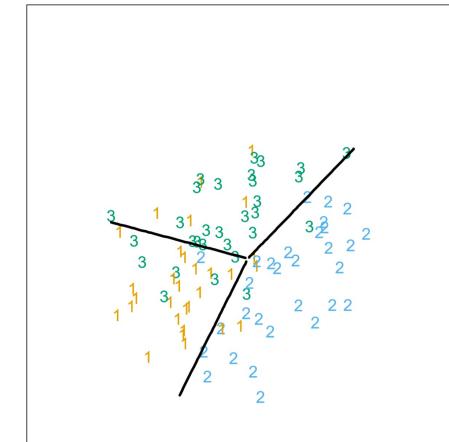
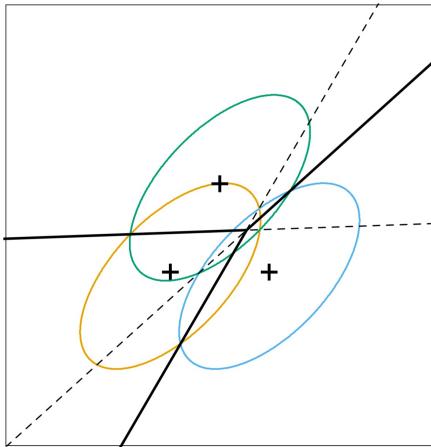
$$\phi_1 + \phi_2 = 1$$

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

$$y \sim Gx(\phi_j)$$

$$= \frac{1}{1 + e^{\frac{1}{2} \left( \frac{x-\mu_1}{\sigma_1} \right)^2 - \frac{1}{2} \left( \frac{x-\mu_2}{\sigma_2} \right)^2 + \log \left( \frac{1-\phi}{\phi} \right)}}$$

$$= \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}}$$



$$P(y) = \phi_1^{\mathbb{1}(y=1)} \phi_2^{\mathbb{1}(y=2)} \cdots \phi_k^{\mathbb{1}(y=k)} \quad \phi_1 + \cdots + \phi_k = 1$$

$$\phi_j = \frac{\sum_{i=1}^m \mathbb{1}(y^{(i)}=j)}{m} = \frac{\#(\text{class-}j \text{ observation})}{\#(\text{observations})} =: \frac{N_j}{m}$$

$$\vec{\mu}_j = \frac{\sum_{i=1}^m \mathbb{1}(y^{(i)}=j) \cdot \vec{x}^{(i)}}{N_j} \quad j=1, 2, \dots, k \quad \text{Sample mean}$$

$$\Sigma = \frac{1}{m-k} \sum_{i=1}^m (\vec{x}^{(i)} - \vec{\mu}_{y^{(i)}})(\vec{x}^{(i)} - \vec{\mu}_{y^{(i)}})^T \quad \text{Sample Cov.}$$

Remark    boundaries:

$$l_{ij} = \left\{ \vec{x} \in \mathbb{R}^d \middle| P(y=i|\vec{x}) = P(y=j|\vec{x}) \right\}$$

$$= \left\{ \vec{x} \in \mathbb{R}^d \middle| \log \frac{P(y=i|\vec{x})}{P(y=j|\vec{x})} = 0 \right\}$$

Remark: Maximize

$$P(y=j|\vec{x}) = \frac{P(\vec{x}|y=j) P(y=j)}{P(\vec{x})}$$

$$\log(P(y=j | \vec{x})) = \underbrace{\log P(\vec{x} | y=j)}_{\text{constant}} + \underbrace{\log P(y=j)}_{\text{constant}} - \underbrace{\log(P(\vec{x}))}_{\text{constant}}$$

$$= \boxed{-\frac{1}{2} \log |\Sigma_j| - \frac{1}{2} (\vec{x} - \mu_j)^T \Sigma_j^{-1} (\vec{x} - \mu_j) + \log \phi_j} + \text{constant.}$$

Assume  $\Sigma_1 = \dots = \Sigma_k = \Sigma$

$$= \boxed{\vec{x}^T \Sigma_j^{-1} \mu_j - \frac{1}{2} \mu_j^T \Sigma_j^{-1} \mu_j + \log \phi_j} - \frac{1}{2} \log |\Sigma| + \text{constant} \quad \text{Gaussian discriminant func.}$$

$\delta_j(\vec{x})$

$$\frac{\partial \delta_j(\vec{x})}{\partial \phi_j} = 0$$

• Class  $G(\vec{x}) = \arg \max_j \delta_j(\vec{x})$

Remark:  $x \in \mathbb{Z}^+$ .

$$P(x | y=i) = \frac{\lambda_i^x e^{-\lambda_i}}{x!}$$

Assumption:

$$\begin{aligned} x | y=1 &\sim \text{Poisson}(\lambda_1) \\ x | y=0 &\sim \text{Poisson}(\lambda_0) \end{aligned}$$

$y \sim \text{Bernoulli}(\phi)$

$$z = f(x, \lambda_i, \phi)$$

$$P(y=1 | \vec{x}) = \dots = \frac{1}{1 + e^{-z}}$$

E.x.

## Linear Discriminant Analysis vs Quadratic Discriminant Analysis

Linear Discriminant Analysis      Quadratic Discriminant Analysis

