

Section 7 Gaussian Discriminant Analysis

1. Gaussian Discriminant Analysis (GDA)
2. Review Gaussian distribution
3. Linear Discriminant Analysis (LDA)
(QDA)
quadratic
4. LDA v.s. Logistics regression

Question:

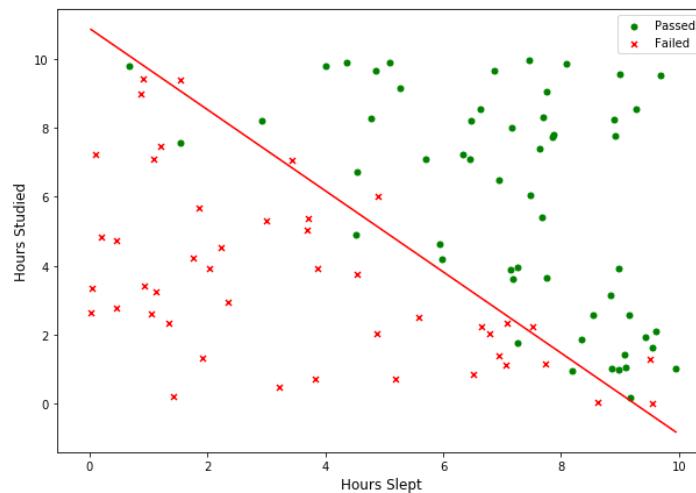
Classification Data:

$$\text{or } y^{(i)} \in \{0, 1\}$$

$$D = \{(\vec{x}^{(i)}, y^{(i)}), i = 1, \dots, n\} \quad \underline{y^{(i)} \in \{1, 2, \dots, K\}},$$

Goal: Find conditional (posterior) probability

$$P(\underline{Y = k} | \vec{X} = \vec{x}) \quad \text{for } k = 1, 2, \dots, K$$

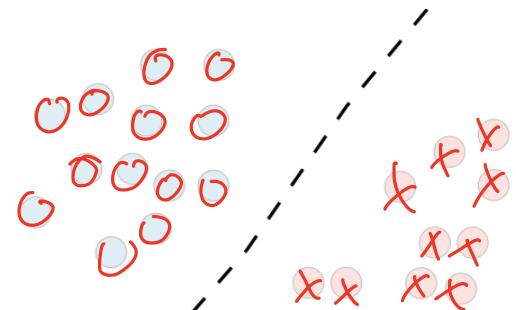


➤ Discriminative learning algorithms.

- Based on a model of the conditional probability.

$$P(Y = k | \vec{X} = \vec{x}) = \text{some model}$$

Examples: logistics/softmax, SVM,



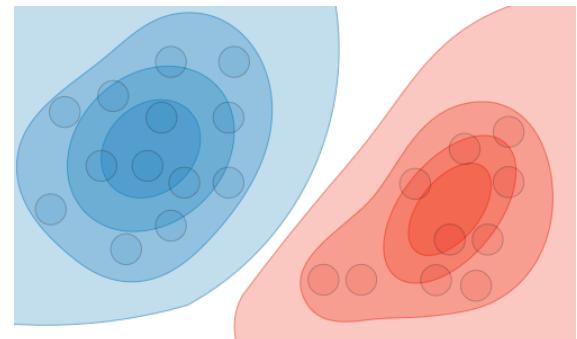
➤ Generative learning algorithms.

$$P(Y=k | \vec{X}=\vec{x}) = \frac{P(X|Y=k) \cdot P(Y=k)}{P(X)}$$

- Based on models of the distributions of the dataset:

- Prior probability $P(Y)$
- Likelihood probability $P(X|Y = k)$

Examples: Gauss discriminant analysis (GDA):
LDA/QDA), Navis Bayes,



➤ Gauss discriminant analysis

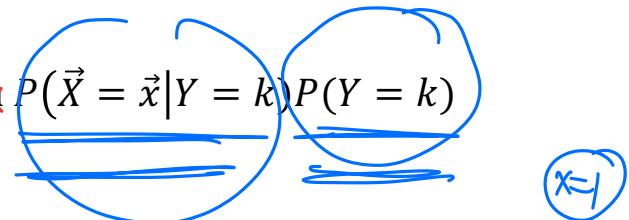
By Bayes Rule:

$$P(Y = k | \vec{X} = \vec{x}) = \frac{P(\vec{X} = \vec{x} | Y = k)P(Y = k)}{\sum_{all\ i} P(\vec{X} = \vec{x} | Y = i)P(Y = i)}$$

$$= \frac{P(\vec{X} = \vec{x} | Y = k)P(Y = k)}{P(\vec{X} = \vec{x})}$$

$$\underline{Posterior} = \frac{\underline{Likelihood} \times \underline{Prior}}{\underline{Evidence}}$$

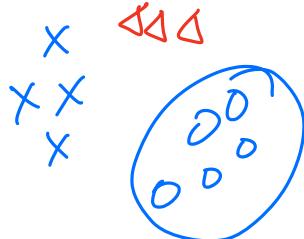
$$\arg\max_y P(Y = k | \vec{X} = \vec{x}) = \arg\max_y P(\vec{X} = \vec{x} | Y = k)P(Y = k)$$



➤ Gaussian Discriminant Analysis (GDA) assumptions

- Prior probability $P(Y)$

$$\begin{array}{c|c} Y=0 & Y=1 \\ \hline \phi & \phi \end{array}$$



Assume $Y \sim \text{Bernoulli}(\phi)$ or $\text{Categorical}(\phi_1, \dots, \phi_K)$

- Likelihood probability $P(X|Y = k)$

$$k=0, 1, \dots, K$$

Assume $X|Y = k$ is a normal distribution for each k .

- **Binary Classification** $y \in \{0,1\}$

Assume $Y \sim \text{Bernoulli}(\phi)$

$$\text{Pdf function } p(y) = \boxed{\phi^y(1-\phi)^{1-y}} = \boxed{\begin{cases} \phi & \text{if } y=1 \\ 1-\phi & \text{if } y=0 \end{cases}} = \phi^{\mathbb{I}(y=1)} (1-\phi)^{\mathbb{I}(y=0)}$$

$$\phi_1 + \phi_0 = 1$$

- **Multiclass Classification** $y \in \{1, 2, \dots, K\}$

Assume $Y \sim \text{Categorical}(\phi_1, \dots, \phi_K)$ such that $\phi_1 + \dots + \phi_K = 1$

$$\text{Pdf function } p(y) = \phi_1^{\mathbb{I}(y=1)} \phi_2^{\mathbb{I}(y=2)} \dots \phi_K^{\mathbb{I}(y=K)}$$

$\mathbb{I}(\cdot)$ indicator

$$= \begin{cases} \phi_1 & \text{if } y=1 \\ \vdots & \vdots \\ \phi_K & \text{if } y=K \end{cases}$$

Normal (Gaussian) distribution (single variable).

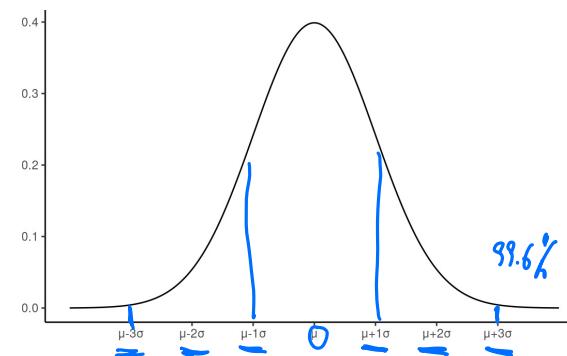
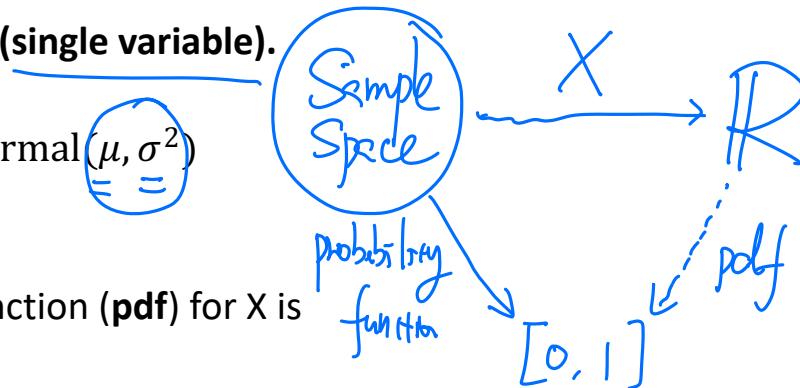
Random variable $X \sim \text{Normal}(\mu, \sigma^2)$

- The probability density function (**pdf**) for X is

$$f_X(x) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \quad -\infty < x < \infty$$

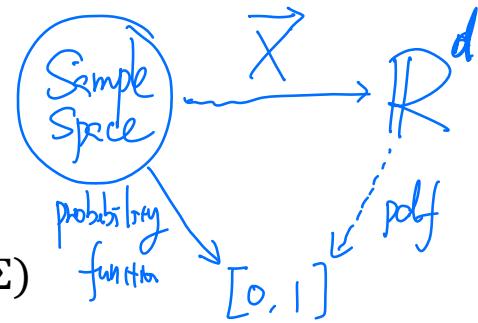
- The **mean** of X is $E(X) = \mu$
- The **variance** of X is $\text{Var}(X) = \sigma^2$
- Probability

$P(a < X < b) = \int_a^b f_X(x)dx$



➤ Multivariate normal distribution.

Vector random variable $\vec{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_d \end{bmatrix} \sim \text{Normal}(\vec{\mu}, \Sigma)$



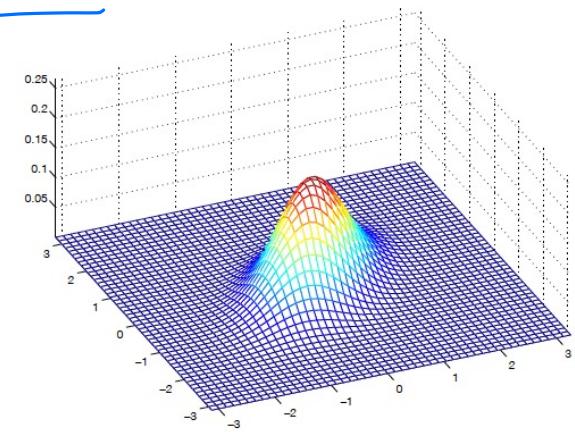
Here $\vec{\mu} \in \mathbb{R}^d$ and Σ is an $d \times d$ symmetric, positive definite matrix.
 $\Sigma = \Sigma^T$ eigenvalues > 0 invertible.

- The joint probability density function (**pdf**) for \vec{X} is

$$f_{\vec{X}}(\vec{x}) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})\right)$$

$$\Sigma = \begin{bmatrix} \cdot & \cdots & \cdot \\ \vdots & \ddots & \vdots \\ \cdot & \cdots & \cdot \end{bmatrix}$$

$$\int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f_{\vec{X}}(\vec{x}) \, d\vec{x} = 1$$



$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \vec{x} f(\vec{x}) d\vec{x}_1 \dots d\vec{x}_n$$

- The **mean vector** of \vec{X} is $E(\vec{X}) = \vec{\mu}$

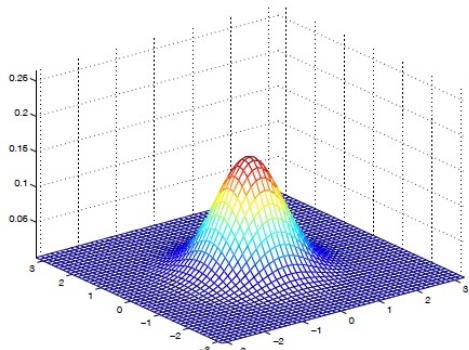
- The **(co)variance matrix** is $\text{Cov}(\vec{X}) = \Sigma$

$$E(\vec{X}\vec{X}^T) - E(\vec{X})E(\vec{X})^T$$

$$\vec{\mu} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

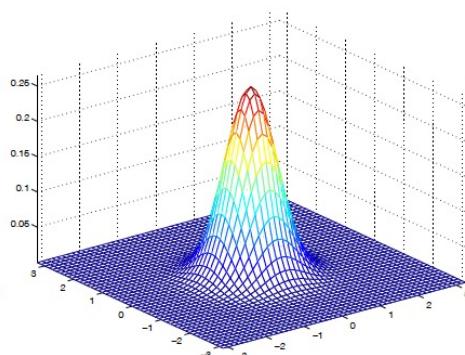
$$\Sigma = I_2$$

Standard normal



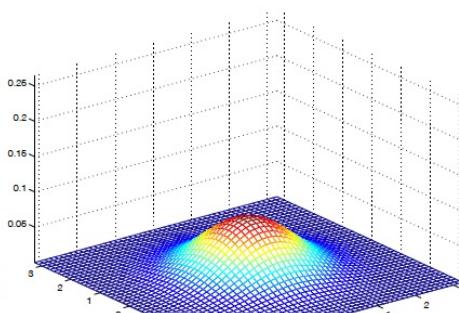
$$\Sigma = 0.6 I_2$$

Compressed

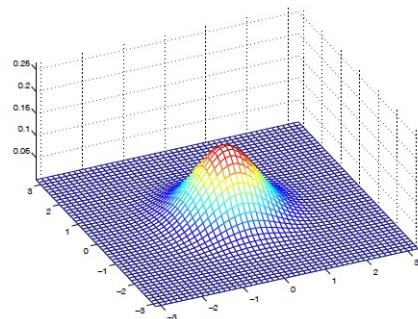


$$\Sigma = 2 I_2$$

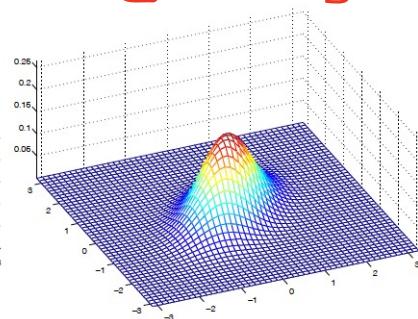
Spread-out



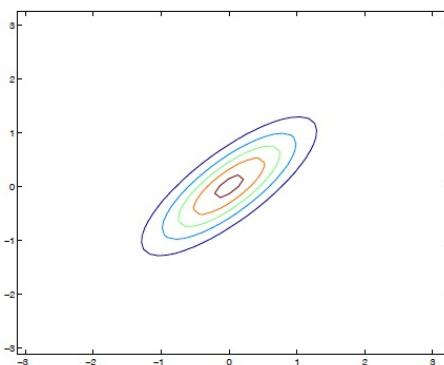
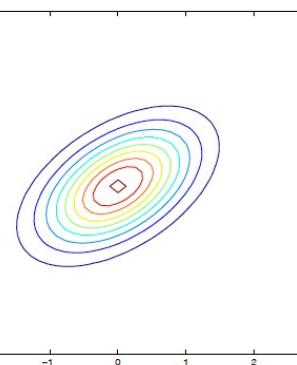
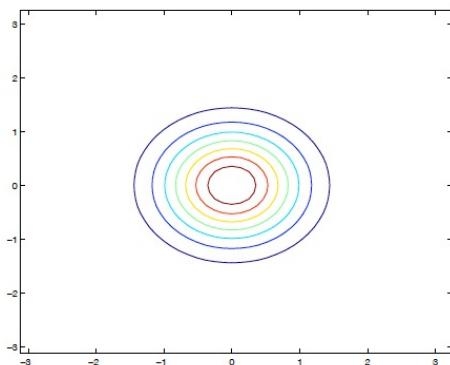
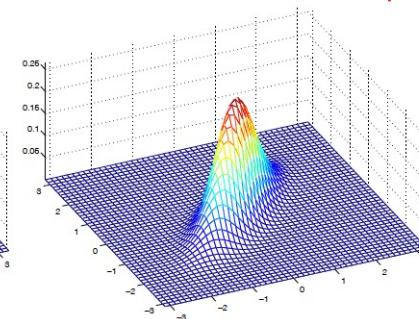
$$\Sigma_1 = I_2$$



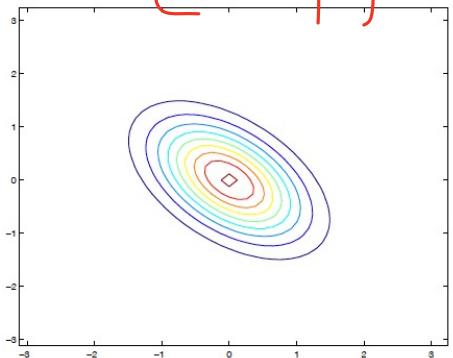
$$\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$



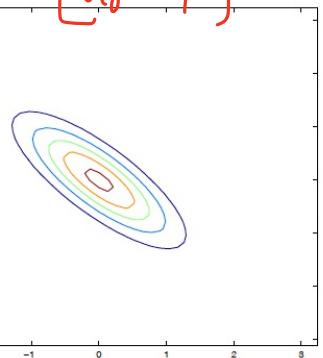
$$\Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$



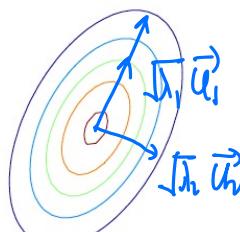
$$\begin{pmatrix} 1 & -0.5 \\ -0.5 & 1 \end{pmatrix}$$



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



$$\begin{pmatrix} 3 & 0.8 \\ 0.8 & 1 \end{pmatrix}$$



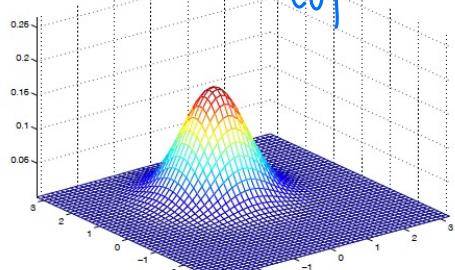
$$\Sigma = P D P^T$$

$$P = [\vec{u}_1 \ \vec{u}_2]$$

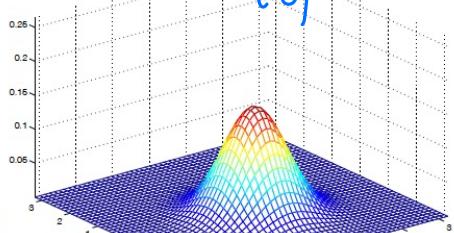
$$D = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \quad \lambda_1 > \lambda_2 > 0$$

$$P^T = P^{-1}$$

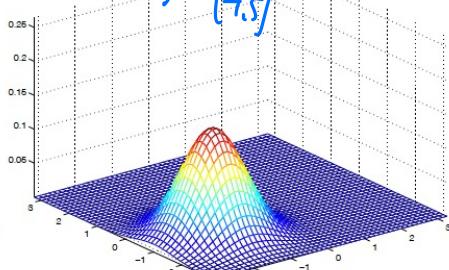
$$\mu = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$



$$\mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$



$$\mu = \begin{pmatrix} -1 \\ -0.5 \end{pmatrix}$$



See more the probability review or Chapter 2 "Pattern Recognition and Machine Learning" -Chris Bishop

➤ Gaussian (Linear) Discriminant Analysis (LDA).

Binary Classification Data: $D = \{(\vec{x}^{(i)}, y^{(i)}), i = 1, \dots, n\}$ $y^{(i)} \in \{0, 1\}$,

Goal: Find conditional (posterior) probability

$$P(Y = k | \vec{X} = \vec{x}) \quad \text{for } k = 0, 1$$

GDA Method: We need to find

$$\left. \begin{array}{l} P(\vec{X} = \vec{x} | Y = k) \\ P(Y = k) \end{array} \right\}$$

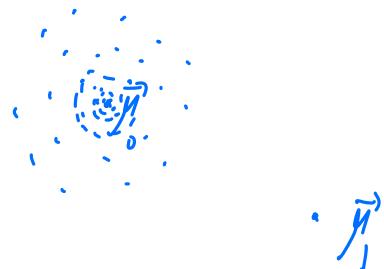
• **Assume** $Y \sim \text{Bernouli}(\phi)$

• **Assume** $\vec{X} | Y = 0 \sim \text{Normal}(\vec{\mu}_0, \Sigma_0)$

$\vec{X} | Y = 1 \sim \text{Normal}(\vec{\mu}_1, \Sigma_1)$

• **LDA Assume:** $\Sigma_0 = \Sigma_1 = \Sigma$

• QDA $\Sigma_0 \neq \Sigma_1$



pdf functions:

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

$$p(\vec{X}|Y=0) = \frac{1}{\sqrt{(2\pi)^d |\Sigma_0|}} \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu}_0)^T \Sigma_0^{-1} (\vec{x} - \vec{\mu}_0)\right)$$

$$p(\vec{X}|Y=1) = \frac{1}{\sqrt{(2\pi)^d |\Sigma_1|}} \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu}_1)^T \Sigma_1^{-1} (\vec{x} - \vec{\mu}_1)\right)$$

Given data $D = \{(\vec{x}^{(i)}, y^{(i)}), i = 1, \dots, n\}$, we want to **maximize likelihood**

$$P(\text{data}) = P(X, \vec{y}) = \prod_{i=1}^n P(\vec{X} = \vec{x}^{(i)}, Y = y^{(i)})$$

independent

Equivalently, we **maximize likelihood**

$$L(\phi, \vec{\mu}_0, \vec{\mu}_1, \Sigma) = \prod_{i=1}^n p(\vec{X} = \vec{x}^{(i)} \mid Y = y^{(i)}) p_Y(y^{(i)})$$

Equivalently, we **maximize log likelihood**

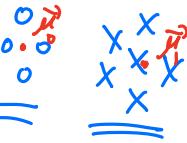
$$l(\phi, \vec{\mu}_0, \vec{\mu}_1, \Sigma) = \log L(\phi, \vec{\mu}_0, \vec{\mu}_1, \Sigma)$$

$$= \sum_{i=1}^n \left(\log p(\vec{X} = \vec{x}^{(i)} \mid Y = y^{(i)}) + \log p_Y(y^{(i)}) \right)$$

Calculate $\nabla l(\phi, \vec{\mu}_0, \vec{\mu}_1, \Sigma) = 0$ and find critical points. (Practice.)

~~Practice.~~
Given

Maximum Likelihood estimates: $(\vec{x}^{(i)}, y^{(i)})$ $i=1, \dots, n$
 $y^{(i)} \in \{0, 1\}$



We obtain formulas for the parameters maximizing the likelihood:

$$\hat{\phi} := \frac{1}{n} \sum_{i=1}^n \mathbb{I}(y^{(i)} = 1) = \frac{\#(\text{label 1's})}{\#(\text{total data})}$$

$$\hat{\vec{\mu}}_0 := \frac{\sum_{i=1}^n \mathbb{I}(y^{(i)} = 0) \vec{x}^{(i)}}{\sum_{i=1}^n \mathbb{I}(y^{(i)} = 0)} = \text{Sample mean for } (\underline{\vec{x}^{(i)}}, y^{(i)}=0)$$

$$\hat{\vec{\mu}}_1 := \frac{\sum_{i=1}^n \mathbb{I}(y^{(i)} = 1) \vec{x}^{(i)}}{\sum_{i=1}^n \mathbb{I}(y^{(i)} = 1)}$$

$$\hat{\Sigma} = \frac{1}{n-2} \sum_{i=1}^n \left(\vec{x}^{(i)} - \hat{\vec{\mu}}_{y^{(i)}} \right) \left(\vec{x}^{(i)} - \hat{\vec{\mu}}_{y^{(i)}} \right)^T$$

Sample covariance.

↑
 Assume $\Sigma_0 = \Sigma_1 = \hat{\Sigma}$

A $\Sigma_0 = \Sigma_1 = \hat{\Sigma}$

~~XXX XXX - 00000~~

We have the optimal distribution models with pdf functions:

$$\Sigma_1 = \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (\vec{x}^{(i)} - \vec{\mu}_1)(\vec{x}^{(i)} - \vec{\mu}_1)^T$$

$$\Sigma^{(i)} = \text{...}$$

$$p_Y(k) = \phi^k (1 - \phi)^{1-k}$$

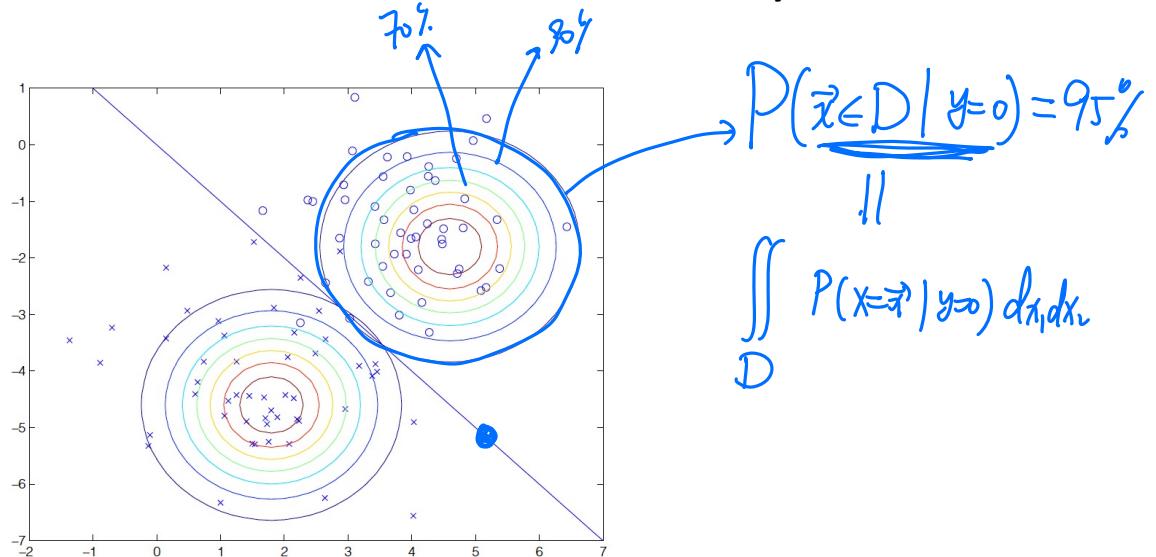
$$p(\vec{X} = \vec{x} | Y = 0) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2} (\vec{x} - \vec{\mu}_0)^T \Sigma^{-1} (\vec{x} - \vec{\mu}_0)\right)$$

$$p(\vec{X} = \vec{x} | Y = 1) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2} (\vec{x} - \vec{\mu}_1)^T \Sigma^{-1} (\vec{x} - \vec{\mu}_1)\right)$$

We have the formulas for the conditional pdf of Y given $\vec{X} = \vec{x}$

$$p(Y = k | \vec{X} = \vec{x}) = \frac{p(\vec{X} = \vec{x} | Y = k) p_Y(k)}{\sum_{all \ i} p(\vec{X} = \vec{x} | Y = i) p_Y(i)}$$

We can find the level curves of the distributions and the **boundary**:



On **boundary** points \vec{x} , we have

$$\log p(Y = 0 | \vec{X} = \vec{x}) = p(Y = 1 | \vec{X} = \vec{x})$$

So,

$$\log \frac{p(Y = 0 | \vec{X} = \vec{x})}{p(Y = 1 | \vec{X} = \vec{x})} = 0$$

$$p(Y = k | \vec{X} = \vec{x}) = \frac{p(\vec{X} = \vec{x} | Y = k) p_Y(k)}{p(\vec{X} = \vec{x})}$$

$$\log p(Y = k | \vec{X} = \vec{x}) = \log p(\vec{X} = \vec{x} | Y = k) + \log p_Y(k) - \log p(\vec{X} = \vec{x})$$

$$= -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (\vec{x} - \vec{\mu}_k)^T \Sigma_k^{-1} (\vec{x} - \vec{\mu}_k) + \log \phi_k + \text{constant}$$

Quadratic discriminant function $\delta_k^Q(\vec{x})$

LDA assumption:

$$\Sigma_0 = \Sigma_1 = \Sigma$$

$$= \vec{x}^T \Sigma^{-1} \vec{\mu}_k - \frac{1}{2} \vec{\mu}_k^T \Sigma^{-1} \vec{\mu}_k + \log \phi_k + \text{constant}$$

Linear discriminant function $\delta_k(\vec{x})$

Boundary formula

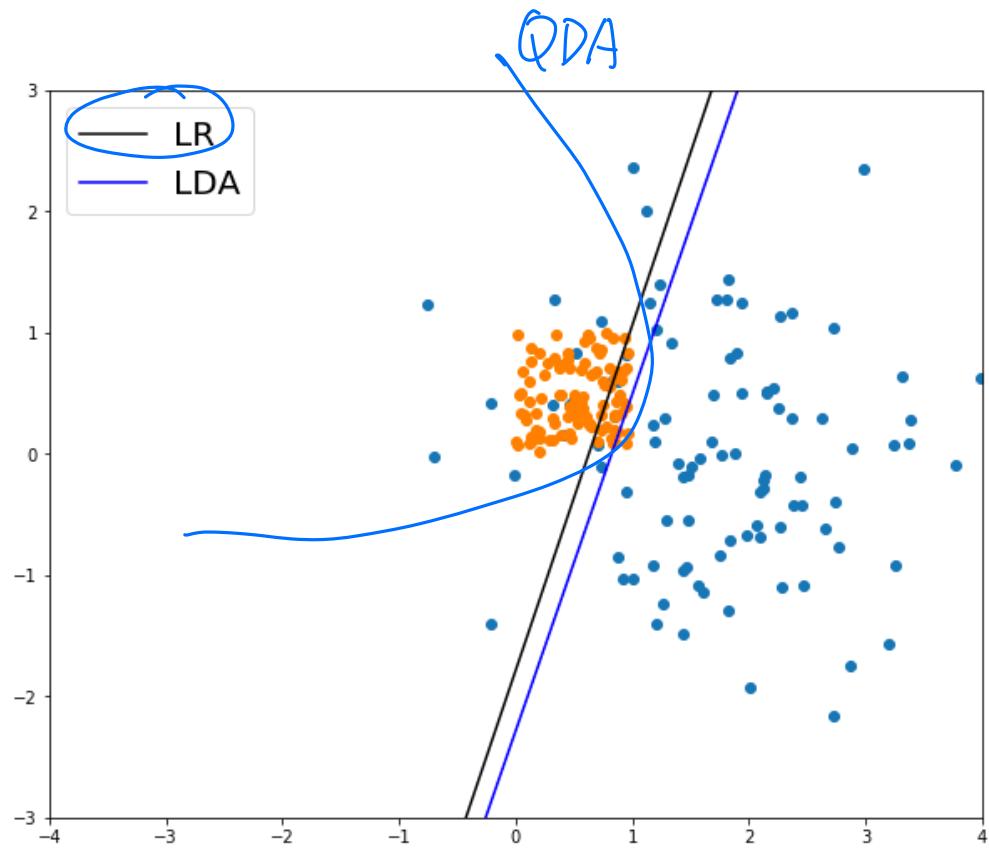
$$\log p(Y = 0 | \vec{X} = \vec{x}) = \log p(Y = 1 | \vec{X} = \vec{x})$$

- LDA boundary Equivalent to $\delta_0(\vec{x}) = \delta_1(\vec{x}) + \text{constant}$

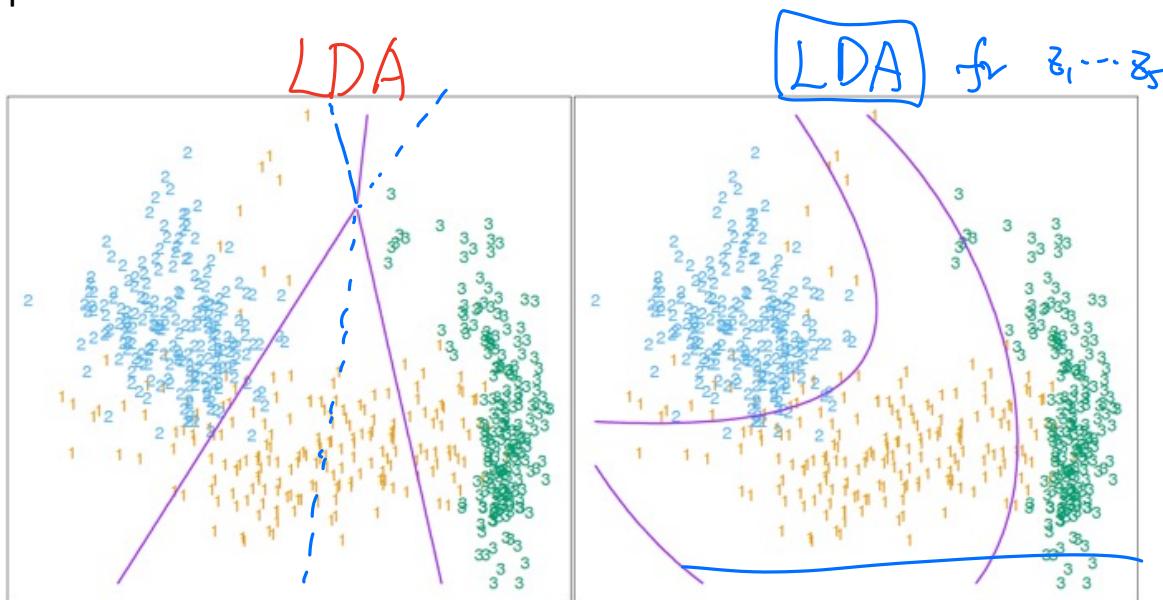
$$\vec{x}^T \Sigma^{-1} \vec{\mu}_0 - \frac{1}{2} \vec{\mu}_0^T \Sigma^{-1} \vec{\mu}_0 + \log \phi_0 = \vec{x}^T \Sigma^{-1} \vec{\mu}_1 - \frac{1}{2} \vec{\mu}_1^T \Sigma^{-1} \vec{\mu}_1 + \log \phi_1$$

- QDA boundary Equivalent to $\delta_0^Q(\vec{x}) = \delta_1^Q(\vec{x})$

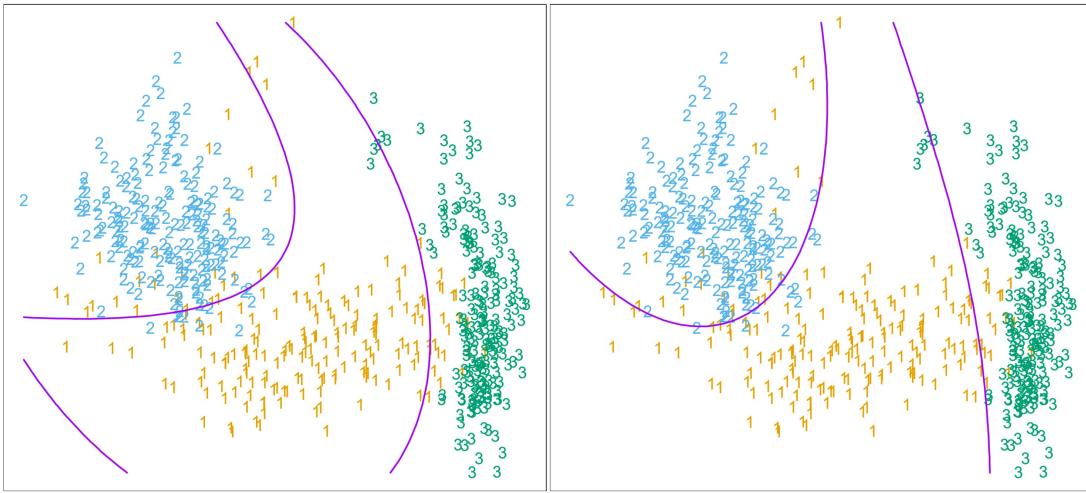
$$-\frac{1}{2} \log |\Sigma_0| - \frac{1}{2} (\vec{x} - \vec{\mu}_0)^T \Sigma_0^{-1} (\vec{x} - \vec{\mu}_0) + \log \phi_0 =$$
$$-\frac{1}{2} \log |\Sigma_1| - \frac{1}{2} (\vec{x} - \vec{\mu}_1)^T \Sigma_1^{-1} (\vec{x} - \vec{\mu}_1) + \log \phi_1$$



Example from the book.



- The left plot shows some data from three classes, with linear decision boundaries found by **linear discriminant analysis**.
- The right plot shows quadratic decision boundaries. These were obtained by finding linear boundaries in the five-dimensional space $(x_1, x_2, x_1x_2, x_1^2, x_2^2)$. Linear inequalities in this space are quadratic inequalities in the original space.



LDA with

QDA

- Two methods for fitting quadratic boundaries.
- The left plot shows the quadratic decision boundaries for the data (obtained using LDA in the five-dimensional space $x_1, x_2, x_1x_2, x_1^2, x_2^2$).
- The right plot shows the quadratic decision boundaries found by QDA. The differences are small, as is usually the case.

LDA/QDA

Compare to Logistic regression

$$P(y=1|\vec{x}) = \frac{1}{1+e^{-\theta^T \vec{x}}} \quad \leftarrow \# \text{ (parameters)} = d+1$$

Goal $P(y=1|\vec{x})$

~~logistic~~ \uparrow ~~not~~

\uparrow

~~QDA~~ \downarrow $\sim \text{Bernoulli}(\phi)$

$\Sigma_1 = \Sigma_0 = \Sigma$

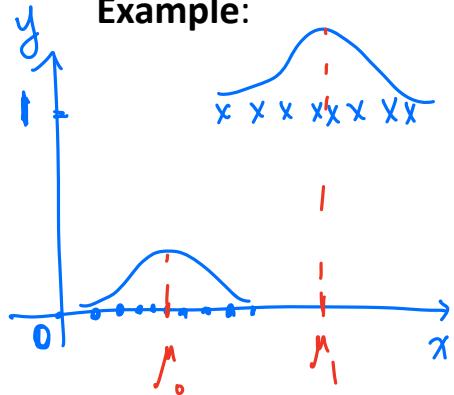
$\vec{x} | Y=1 \sim \text{Normal}(\vec{\mu}_1, \Sigma_i)$

$\leftarrow \# \text{ (parameters)} = 1 + d + d + d^2$

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

$\vec{x} \in \mathbb{R}^1$

Example:



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

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

$i=0, 1$

$\sigma_0 = \sigma_1 = \sigma$

$$P(y=1 | \vec{x}) \stackrel{\text{LDA}}{=} \frac{P(\vec{x} | y_1) P(y=1)}{P(\vec{x} | y=1) P(y=1) + P(\vec{x} | y=0) P(y=0)}$$

Logistic //

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

$$= \frac{\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu_1}{\sigma}\right)^2} \cdot \phi}{\left(\dots \right) + \left(\dots \right) (1-\phi)}$$

$$= \frac{1}{1 + e^{-(\frac{1}{2}\left(\frac{x-\mu_1}{\sigma}\right)^2 - \frac{1}{2}\left(\frac{x-\mu_0}{\sigma}\right)^2 + \ln(\frac{\phi}{1-\phi}))}}$$

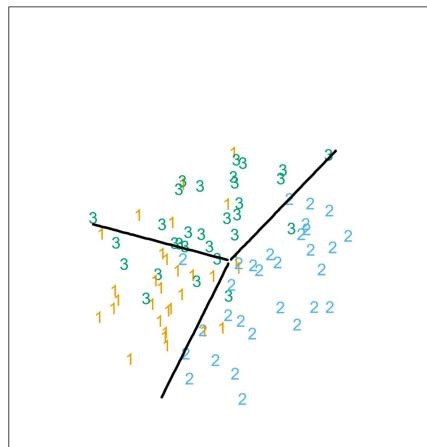
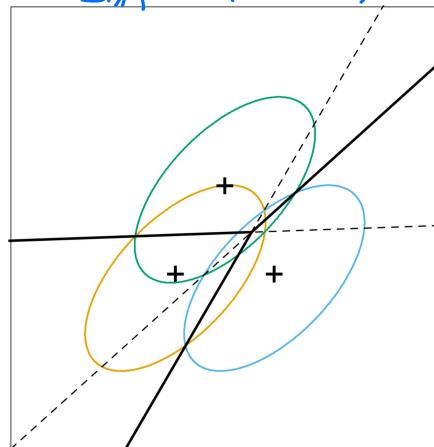
➤ Multiclass Classification

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

Assume $Y \sim \text{Categorical}(\phi_1, \dots, \phi_K)$ such that $\phi_1 + \dots + \phi_K = 1$

Pdf function $p(y) = \phi_1^{\mathbb{I}(y=1)} \phi_2^{\mathbb{I}(y=2)} \dots \phi_K^{\mathbb{I}(y=K)}$

LDA $\Sigma_1 = \Sigma_2 = \Sigma = \bar{\Sigma}$



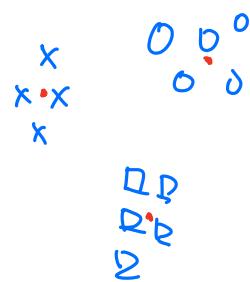
Maximum Likelihood estimates:

We obtain formulas for the parameters maximizing the likelihood:

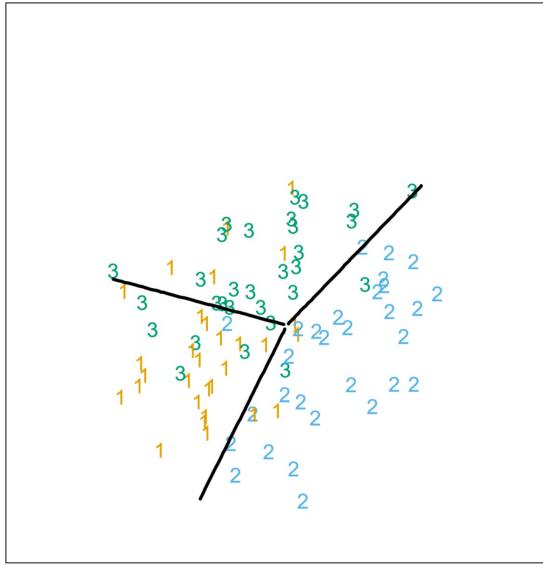
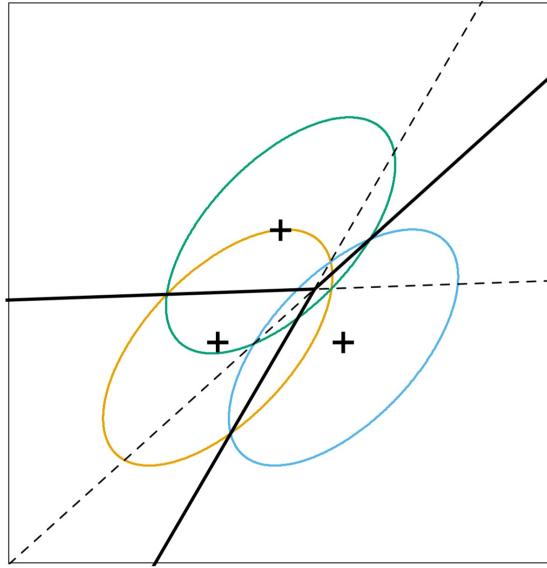
$$\phi_j = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(y^{(i)} = j) = \frac{\#(\text{blue } | \text{ } j)}{n}$$

$$\vec{\mu}_j = \frac{\sum_{i=1}^n \mathbb{I}(y^{(i)} = j) x^{(i)}}{\sum_{i=1}^n \mathbb{I}(y^{(i)} = j)}$$

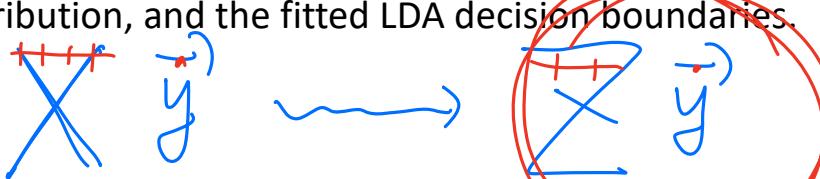
$$\Sigma = \frac{1}{n - K} \sum_{i=1}^n \left(\vec{x}^{(i)} - \vec{\mu}_{y^{(i)}} \right) \left(\vec{x}^{(i)} - \vec{\mu}_{y^{(i)}} \right)^T$$



LDA assumption: $\Sigma_1 = \dots = \Sigma_K = \Sigma$



The left panel shows three Gaussian distributions, with the same covariance and different means. Included are the contours of constant density enclosing 95% of the probability in each case. The Bayes decision boundaries between each pair of classes are shown (broken straight lines), and the Bayes decision boundaries separating all three classes are the thicker solid lines (a subset of the former). On the right we see a sample of 30 drawn from each Gaussian distribution, and the fitted LDA decision boundaries.

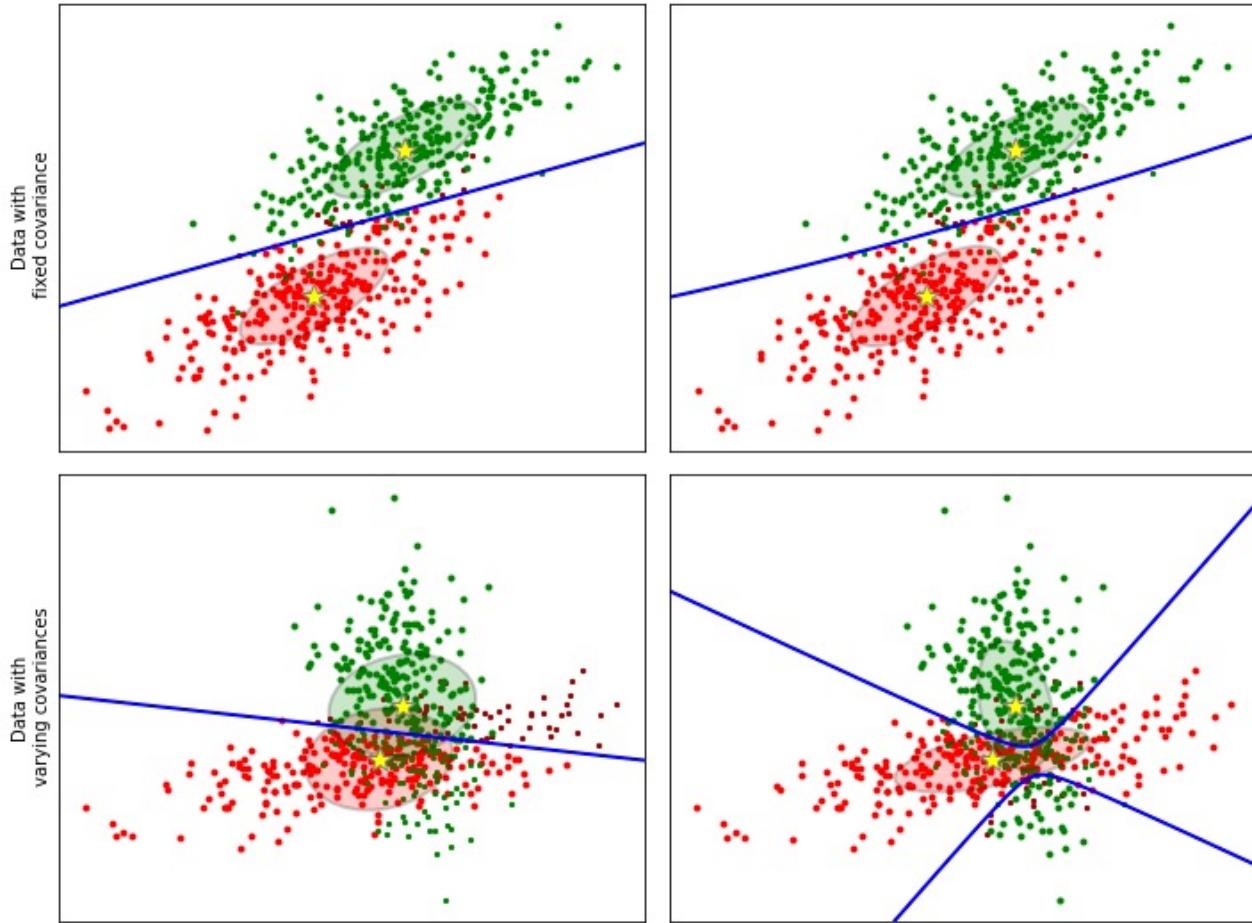


Iris dataset



Linear Discriminant Analysis vs Quadratic Discriminant Analysis

Linear Discriminant Analysis Quadratic Discriminant Analysis



https://scikit-learn.org/stable/modules/lda_qda.html

$$\text{Core } \underline{\underline{J(\vec{\theta})}} =$$

[Hw.] 2.1]

Gradient descent :

$$\begin{bmatrix} \vec{\theta}_0 \\ \vec{\theta}_1 \\ \vdots \\ \vec{\theta}_d \end{bmatrix} = \vec{\theta}^{\text{next}} = \vec{\theta} \boxed{||} - \alpha \nabla J(\vec{\theta})$$

$$\nabla J = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_d} \end{bmatrix}$$

- batch G.D.

$$[X \quad \vec{y}]$$

$$\begin{bmatrix} / / / / / / / / \\ X \end{bmatrix} \quad | \quad \vec{y}$$

- minibatch G.D

$$\begin{bmatrix} / / / / / / / / \\ X \end{bmatrix} \quad | \quad \vec{y}$$

- SGD

$$J^{\text{Ridge}} = \text{RSS}(\vec{\theta}) + \lambda \sum_{i=1}^d \vec{\theta}_i^2$$

assume $\vec{\theta}_0 = 0$

$$= \|X\vec{\theta} - \vec{y}\|^2 + \lambda \|\vec{\theta}\|^2$$

$$\vec{\theta} = \begin{bmatrix} \vec{\theta}_0 \\ \vec{\theta}_1 \\ \vdots \\ \vec{\theta}_d \end{bmatrix}$$

$$\arg \min J^{\text{Ridge}} = (X^T X + \lambda I)^{-1} X^T \vec{y}$$

$$(X, \vec{y})$$

$$\begin{aligned} \tilde{x} &= \vec{x} - \bar{x} \\ \tilde{y} &= y - \bar{y} \end{aligned}$$

$$\underline{\underline{X}}, \bar{y}$$

ridge

$$\hat{y} = \theta_i \tilde{x}$$

$$y - \bar{y} = \theta_i (x - \bar{x})$$

$$\tilde{X} = \begin{bmatrix} (\vec{x}^{(1)} - \bar{\vec{x}})^T \\ \vdots \\ (\vec{x}^{(n)} - \bar{\vec{x}})^T \end{bmatrix}$$

$$\tilde{y} = \begin{bmatrix} y^{(1)} - \bar{y} \\ \vdots \\ y^{(n)} - \bar{y} \end{bmatrix}$$