

$$\lambda_{ij} - \lambda_{jj} = \lambda_{ij}$$

$$R(i|x) = \sum_{j=1}^c \lambda_{ij} P(y=j|x)$$

$$\frac{P(x|y=1)}{P(x|y=2)} \gtrsim \frac{\frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}}}{\frac{P(y=2)}{P(y=1)}}$$

1

$$\lambda_{11} = \lambda_{22} = 0$$

$$\lambda_{12} = \lambda_{21} \approx 1$$

Multi-variate Normal (Gaussian) Distribution:

$$x \in \mathbb{R}^d$$
$$P(x) = N(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{d}{2}} \underbrace{|\Sigma|^{\frac{1}{2}}}_{< 0}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

$\Sigma$  is symmetric and positive.

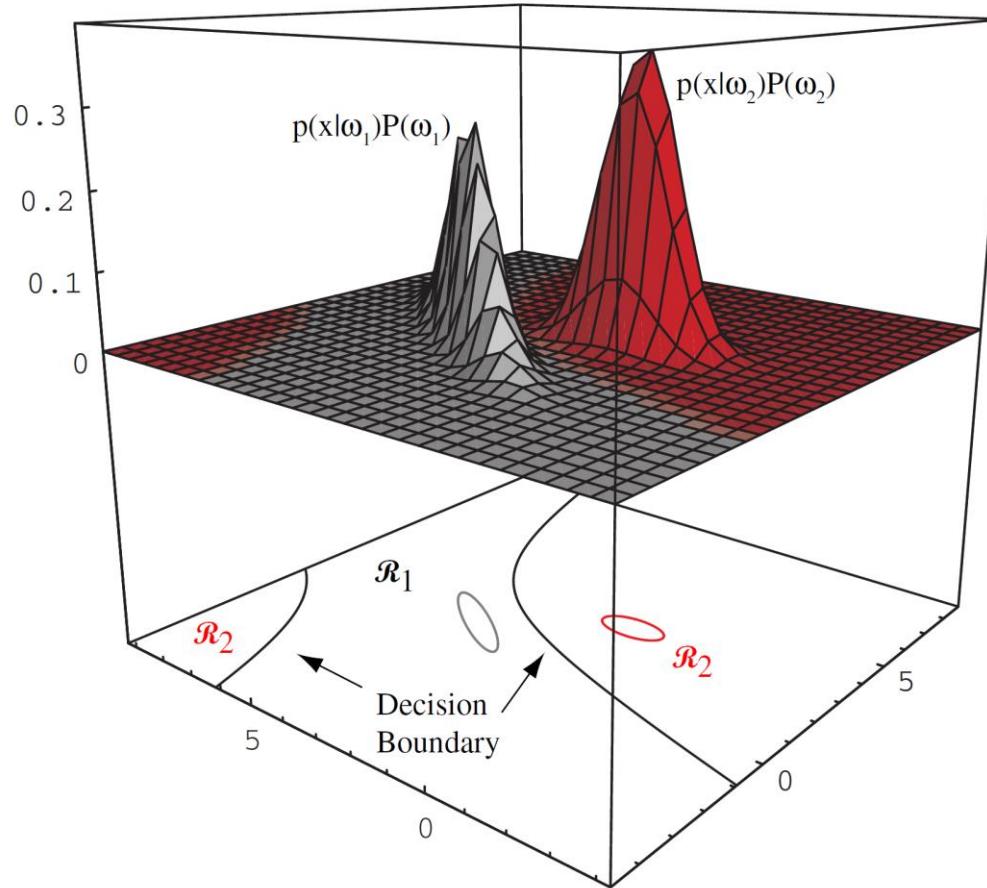
$$x^T A x > 0$$

$$\Sigma = B^T B$$

$$\Sigma^{-1} = B^{-1} B^{-T} = (B^{-1}) (B^{-1})^T$$

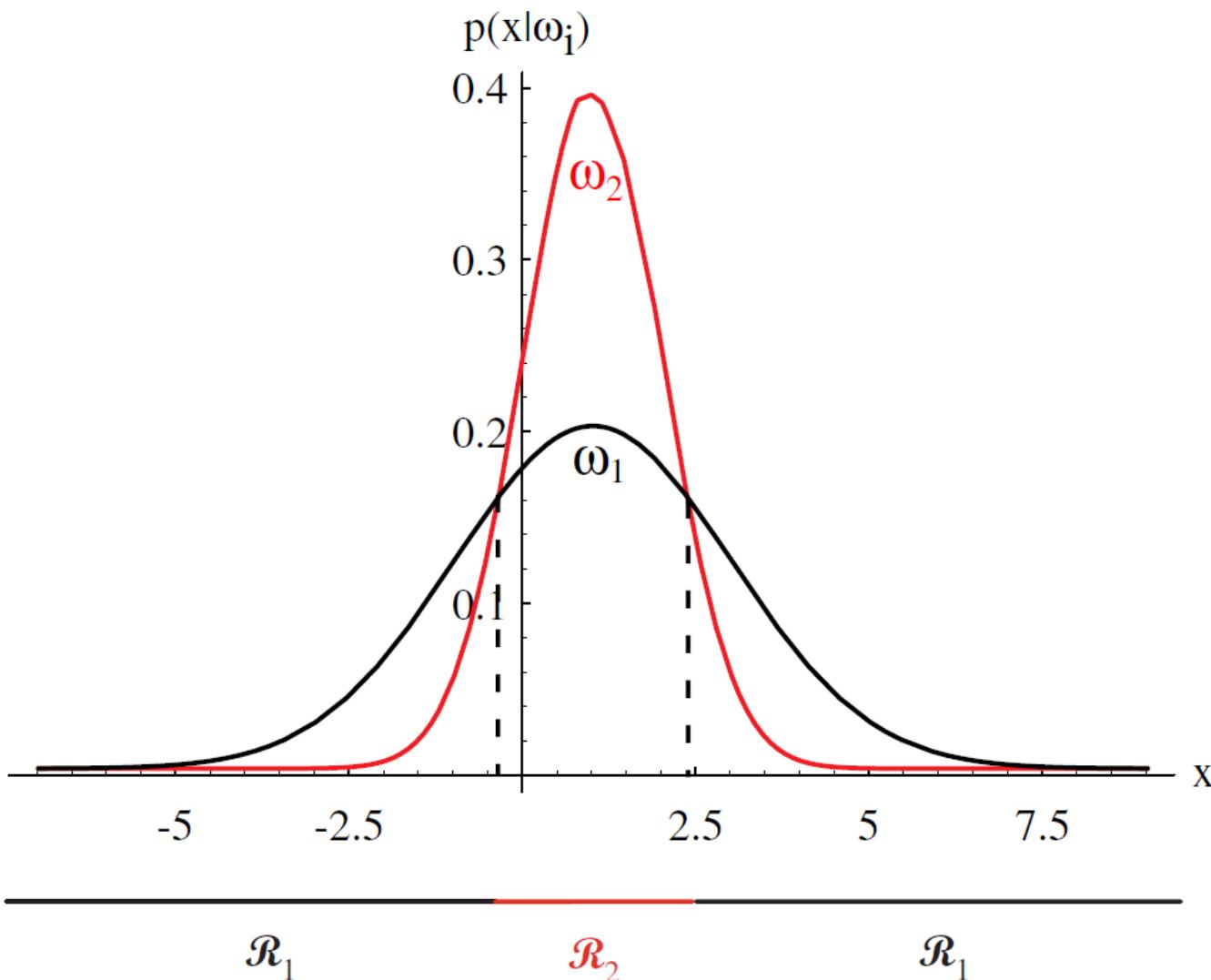
1

# Two-dimensional two-category classifier



the probability densities are Gaussian (with  $1/e$  ellipses shown), the **decision boundary consists of two hyperbolas**, and thus the decision region  $R_2$  is not simply connected.

# Non-simply connected decision regions





# Univariate Gaussian

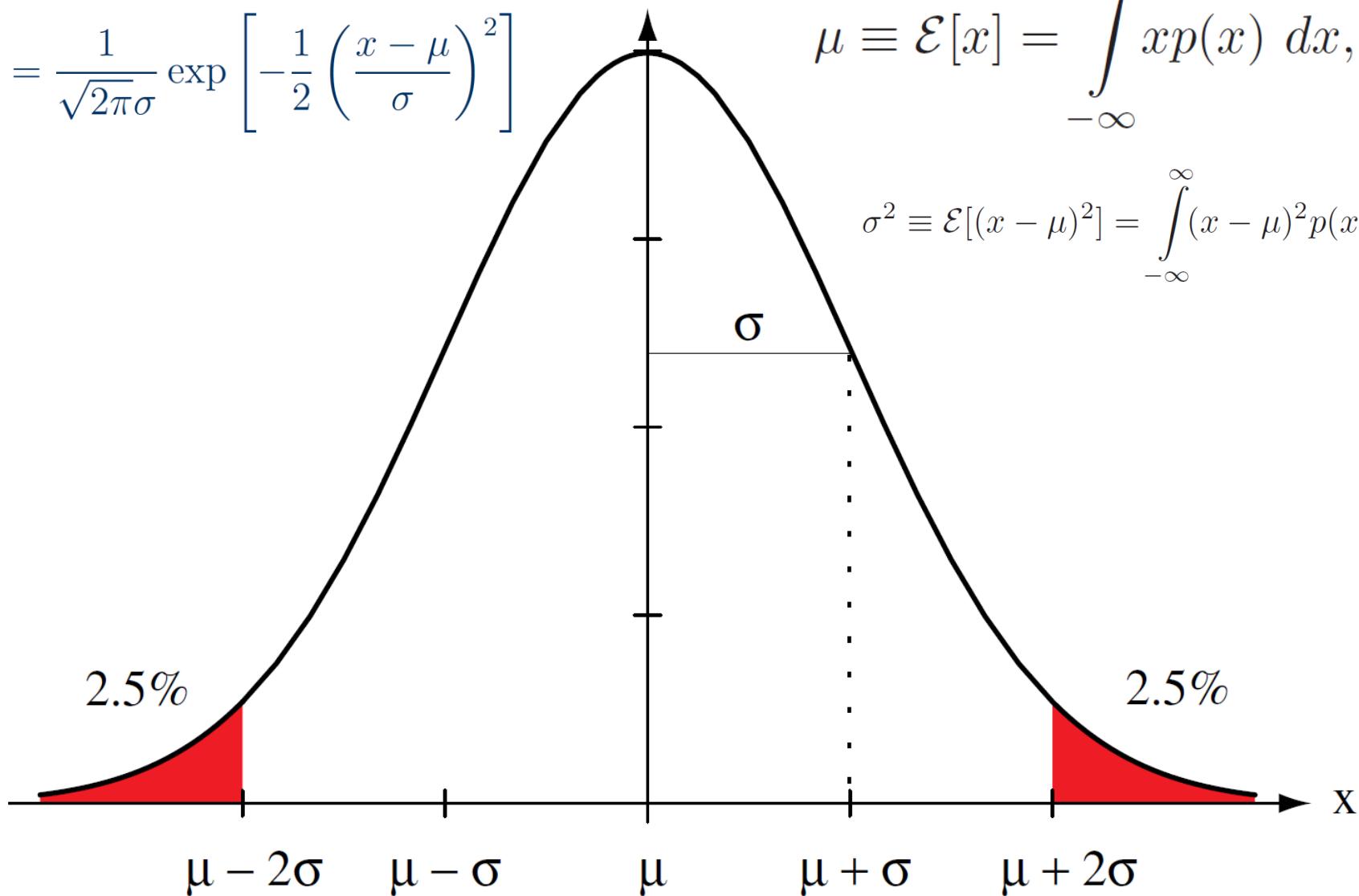
$$p(x) = N(\mu, \sigma^2)$$

$$= \frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{1}{2} \left( \frac{x-\mu}{\sigma} \right)^2 \right]$$

$$p(x)$$

$$\mu \equiv \mathcal{E}[x] = \int_{-\infty}^{\infty} xp(x) \, dx,$$

$$\sigma^2 \equiv \mathcal{E}[(x-\mu)^2] = \int_{-\infty}^{\infty} (x-\mu)^2 p(x) \, dx.$$





# Multivariate Gaussian

- General multivariate normal density in  $d$  dimensions

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^t \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right],$$



- Where

$$\boldsymbol{\mu} \equiv \mathcal{E}[\mathbf{x}] = \int \mathbf{x} p(\mathbf{x}) d\mathbf{x}$$

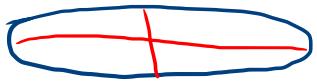
$$\Sigma \equiv \mathcal{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^t] = \int (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^t p(\mathbf{x}) d\mathbf{x},$$

$$\mu_i = \mathcal{E}[x_i] \quad \sigma_{ij} = \mathcal{E}[(x_i - \mu_i)(x_j - \mu_j)].$$

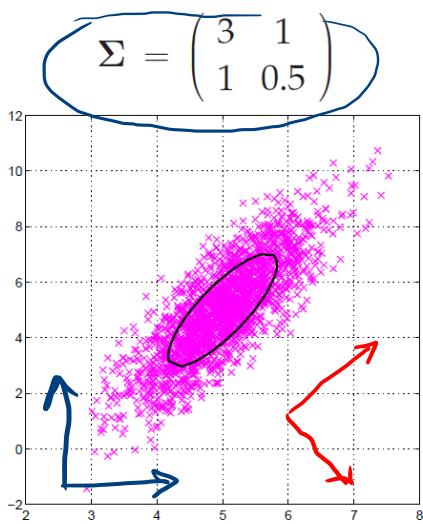
- $\Sigma$  is **positive definite**, so that the determinant of  $\Sigma$  is strictly **positive**

# Example

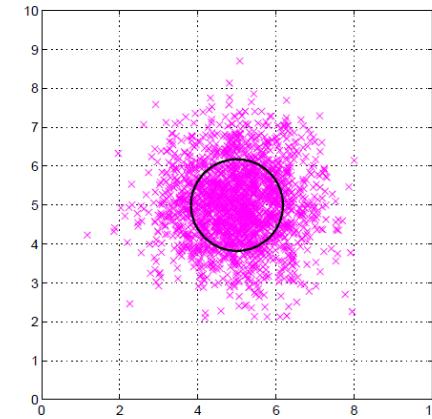
$$\Sigma \approx \begin{bmatrix} 5 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 5 \end{bmatrix}$$

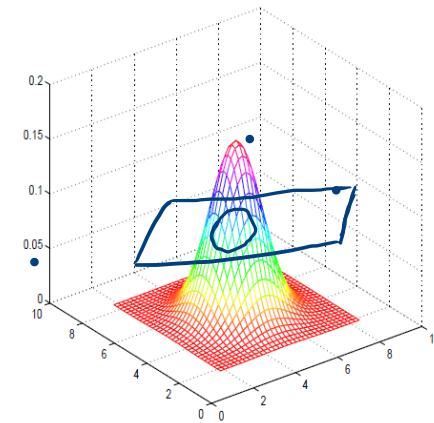
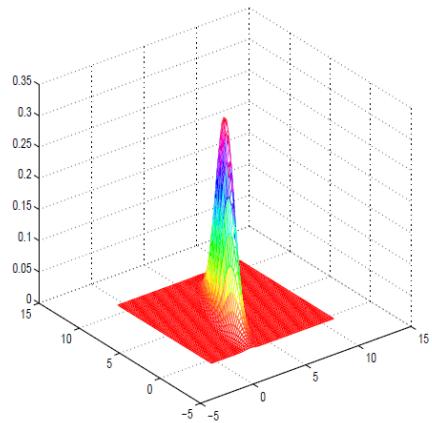


$$\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$



•  $M = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

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



A. Dehaqani, UT

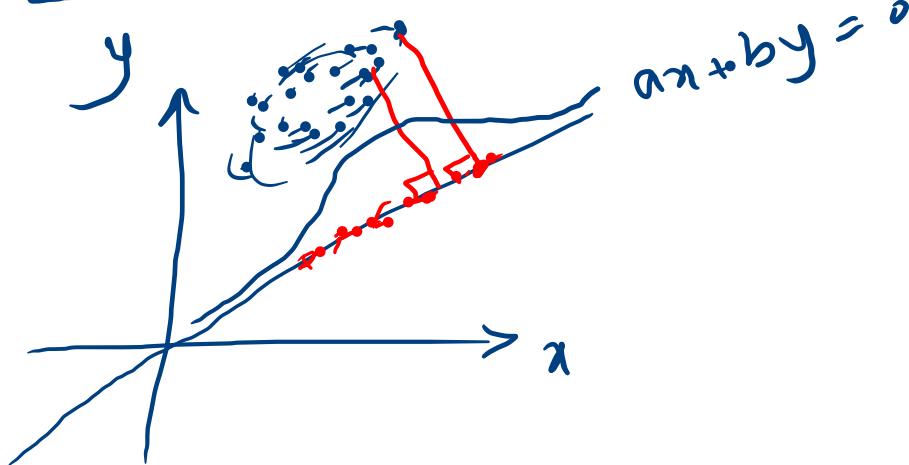
$$P(x, y) = \mathcal{N}\left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \Sigma_{2 \times 2}\right)$$

$$\begin{array}{ll} P(x) & P(x|y) \\ \backslash & / \\ P(y) & P(y|x) \end{array}$$

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

$$P(x, z)$$

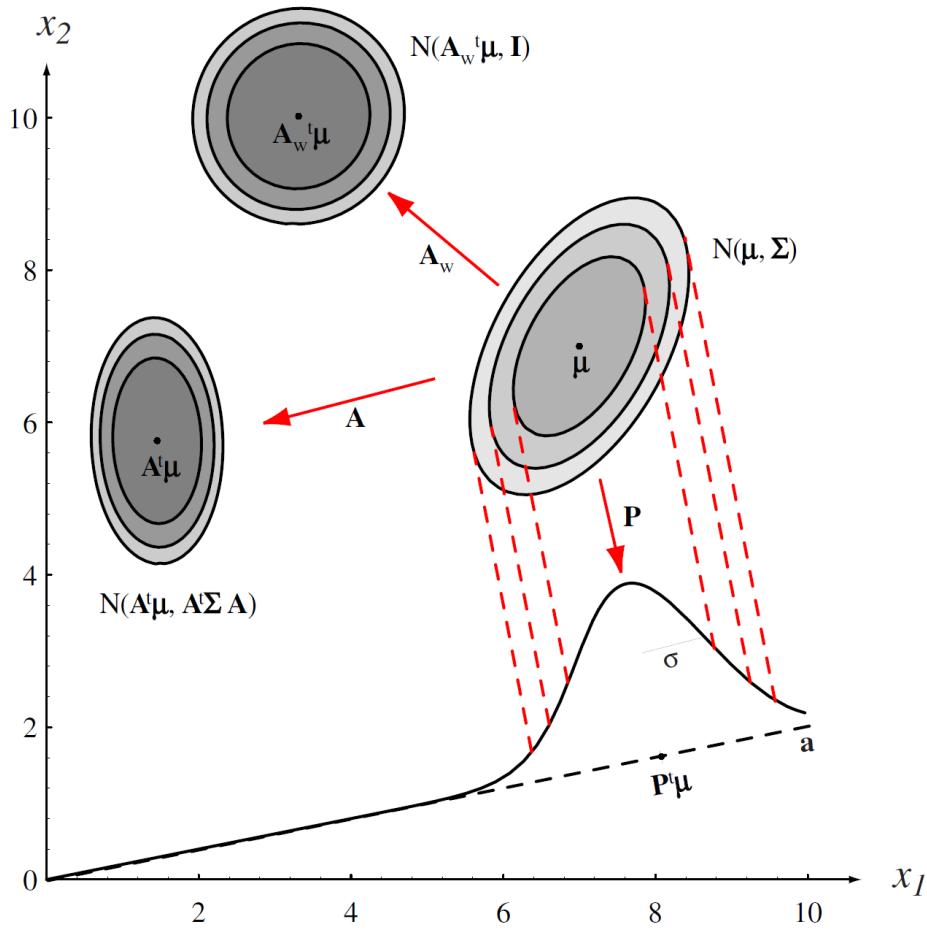
$$\boxed{ax + by = z}$$





# Linear Transformations

- $\mathbf{A}$  is a  $d$ -by- $k$  matrix and  $\mathbf{y} = \mathbf{A}^t \mathbf{x}$  is a  $k$ -component vector,
- If  $k = 1$  and  $\mathbf{A}$  is a unit-length vector  $\mathbf{a}$ ,  $y = \mathbf{a}^t \mathbf{x}$  is a scalar that represents the **projection of  $x$  onto a line in the direction of  $a$** ;
- In general then, **knowledge of the covariance matrix** allows us to calculate the dispersion of the data in **any direction**, or in any **subspace**.



$$x \in \mathbb{R}^d$$

$$\underline{y = Ax}$$

$$A \in \mathbb{R}^{m \times d}$$

$$y \in \mathbb{R}^m$$

$$E[a+b]$$

$$x \sim N(\mu, \Sigma)$$

$$y \sim N(\mu_y, \Sigma_y)$$

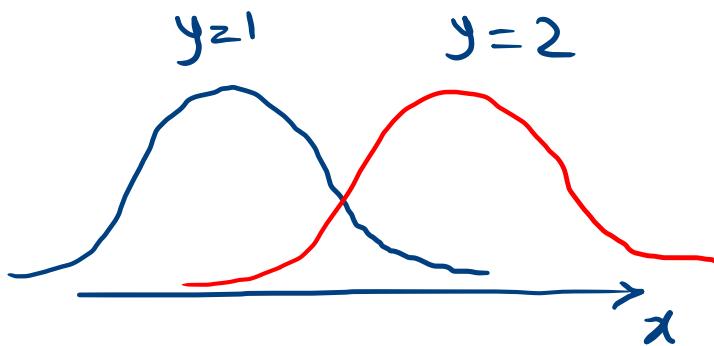
$$\mu_y = E[y] = E[Ax] = A E[x] = \underline{\underline{A\mu}}$$

$$\Sigma_y = E[(y - \mu_y)(y - \mu_y)^T] = \underbrace{E[yy^T]}_{E[Ax x^T A^T]} - \mu_y \mu_y^T = \underline{\underline{A \Sigma_x A^T}}$$
$$E[Ax x^T A^T] = A E[x x^T] A^T$$
$$A E[x x^T] A^T - A \mu \mu^T A^T$$

$$\frac{P(y=1|x)}{P(y=2|x)} \stackrel{1}{\underset{2}{\gtrless}} 1 \Rightarrow \underbrace{\ln P(y=1|x) - \ln P(y=2|x)}_{\underbrace{P(x|y=1)P(y=1)}_{\text{خوب نیم زمان است}}} \stackrel{g(n)}{\underset{2}{\gtrless}} 0$$

$$P(y|x) =$$

$$\left\{ \begin{array}{l} P(x|y=1) = N(x; \mu_1, \Sigma_1) \\ P(x|y=2) = N(x; \mu_2, \Sigma_2) \end{array} \right.$$



!

$$g(x) = \ln \underbrace{P(x|y=1)}_{\cancel{x^T \Sigma_1^{-1} x}} P(y=1) - \ln \underbrace{P(x|y=2)}_{\cancel{x^T \Sigma_2^{-1} x}} P(y=2)$$

$$\ln P(x|y=1) = -\frac{d}{2} \cancel{\ln(2\pi)} - \frac{1}{2} \cancel{\ln |\Sigma_1|} - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1)$$

$$\ln P(x|y=2) = -\frac{d}{2} \cancel{\ln 2\pi} - \frac{1}{2} \cancel{\ln |\Sigma_2|} - \frac{1}{2} (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2)$$

$$\Sigma_1 = \Sigma_2 = \Sigma$$

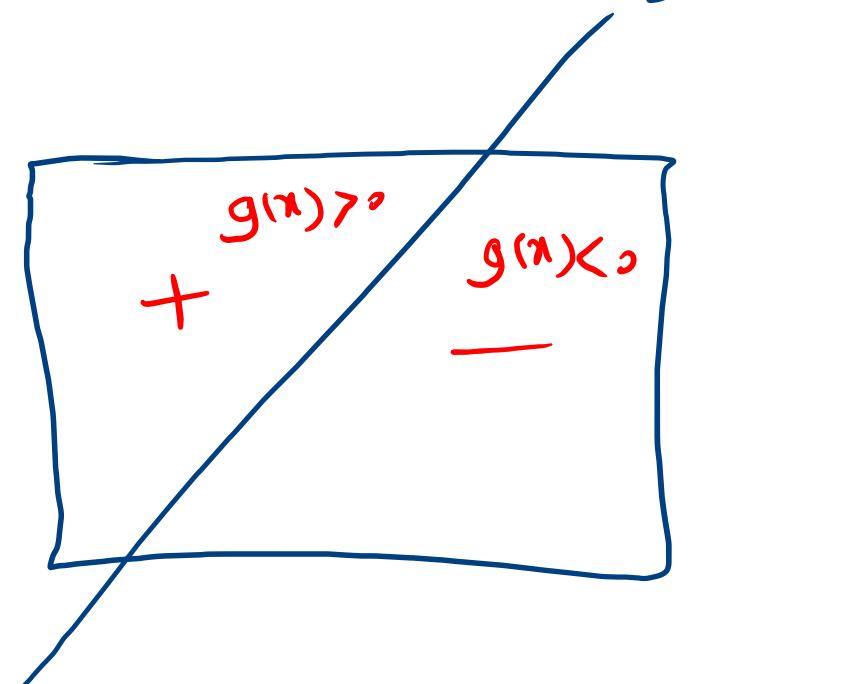
$$g(x) = \underline{\underline{\omega}}^T \underline{\underline{x}} + \underline{\underline{\omega}_0}$$

$$\underbrace{\omega_1 x_1 + \omega_2 x_2 + \omega_0 = 0}_{\text{---}}$$

$$\omega = \begin{bmatrix} \omega_1 \\ \omega_2 \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$g(x) = \omega^T x + \omega_0$$

$$\omega^T x + \omega_0 = 0$$





# Case 1: $\Sigma_i = \sigma^2 \mathbf{I}$

- Features are **statistically independent**, and each feature has the **same variance**,  $\sigma^2$
- Samples fall in equal-size **hyperspherical clusters**
- $|\Sigma_i| = \sigma^{2d}$  and  $\Sigma_i^{-1} = (1/\sigma^2)\mathbf{I}$
- $|\Sigma_i|$  and the  $(d/2) \ln 2\pi$  term are **independent of  $i$**

$$g_i(\mathbf{x}) = -\frac{\|\mathbf{x} - \boldsymbol{\mu}_i\|^2}{2\sigma^2} + \ln P(\omega_i),$$

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^t \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i).$$

where  $\|\cdot\|$  is the *Euclidean norm*, that is,

$$\|\mathbf{x} - \boldsymbol{\mu}_i\|^2 = (\mathbf{x} - \boldsymbol{\mu}_i)^t (\mathbf{x} - \boldsymbol{\mu}_i).$$



# Discriminant Functions for the Gaussian Density

- Discriminant functions for **minimum-error-rate** classification can be written as

$$g_i(\mathbf{x}) = \ln p(\mathbf{x}|\omega_i) + \ln P(\omega_i).$$

- For  $p(\mathbf{x}|\omega_i) \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ .

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^t \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}_i| + \ln P(\omega_i).$$



# Case 1: $\Sigma_i = \sigma^2 \mathbf{I}$

- Features are **statistically independent**, and each feature has the **same variance**,  $\sigma^2$
- Samples fall in equal-size **hyperspherical clusters**
- $|\Sigma_i| = \sigma^{2d}$  and  $\Sigma_i^{-1} = (1/\sigma^2)\mathbf{I}$
- $|\Sigma_i|$  and the  $(d/2) \ln 2\pi$  term are **independent of  $i$**

$$g_i(\mathbf{x}) = -\frac{\|\mathbf{x} - \boldsymbol{\mu}_i\|^2}{2\sigma^2} + \ln P(\omega_i),$$

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^t \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i).$$

where  $\|\cdot\|$  is the *Euclidean norm*, that is,

$$\|\mathbf{x} - \boldsymbol{\mu}_i\|^2 = (\mathbf{x} - \boldsymbol{\mu}_i)^t (\mathbf{x} - \boldsymbol{\mu}_i).$$



# linear discriminant functions

- Expansion of the quadratic form  $(\mathbf{x} - \boldsymbol{\mu})^T(\mathbf{x} - \boldsymbol{\mu})$  yields

$$g_i(\mathbf{x}) = -\frac{1}{2\sigma^2} [\mathbf{x}^t \mathbf{x} - 2\boldsymbol{\mu}_i^t \mathbf{x} + \boldsymbol{\mu}_i^t \boldsymbol{\mu}_i] + \ln P(\omega_i),$$

- It seems quadratic, but term  $\mathbf{x}^t \mathbf{x}$  is the same for all  $i$ , making it an **ignorable additive constant**
- Thus, we obtain the **linear discriminant functions**

$$g_i(\mathbf{x}) = \mathbf{w}_i^t \mathbf{x} + w_{i0},$$

- Where

$$\mathbf{w}_i = \frac{1}{\sigma^2} \boldsymbol{\mu}_i$$

$$w_{i0} = \frac{-1}{2\sigma^2} \boldsymbol{\mu}_i^t \boldsymbol{\mu}_i + \ln P(\omega_i).$$

We call  $w_{i0}$  the **threshold** or **bias** in the  $i^{th}$  category



# Decision Boundaries

- Decision boundaries are the hyperplanes  $g_i(x)=g_j(x)$ , and can be written as

$$\mathbf{w}^t(\mathbf{x} - \mathbf{x}_0) = 0,$$

- Where

$$\mathbf{w} = \boldsymbol{\mu}_i - \boldsymbol{\mu}_j$$

$$\mathbf{x}_0 = \frac{1}{2}(\boldsymbol{\mu}_i + \boldsymbol{\mu}_j) - \frac{\sigma^2}{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2} \ln \frac{P(\omega_i)}{P(\omega_j)} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j).$$

- Hyperplane separating  $R_i$  and  $R_j$  passes through the point  $\mathbf{x}_0$  and is **orthogonal** to the **vector  $\mathbf{w}$** .



# Minimum-distance classifier

- Special case when  $P(w_i)$  are the same for  $c$  classes the  $\ln P(w_i)$  term becomes unimportant additive **constant** that can be ignored and it is **minimum-distance classifier**

$$g_i(\mathbf{x}) = -\frac{\|\mathbf{x} - \boldsymbol{\mu}_i\|^2}{2\sigma^2} + \ln P(\omega_i) \longrightarrow g_i(\mathbf{x}) = -\|\mathbf{x} - \boldsymbol{\mu}_i\|^2$$

- The minimum-distance classifier that uses the decision rule

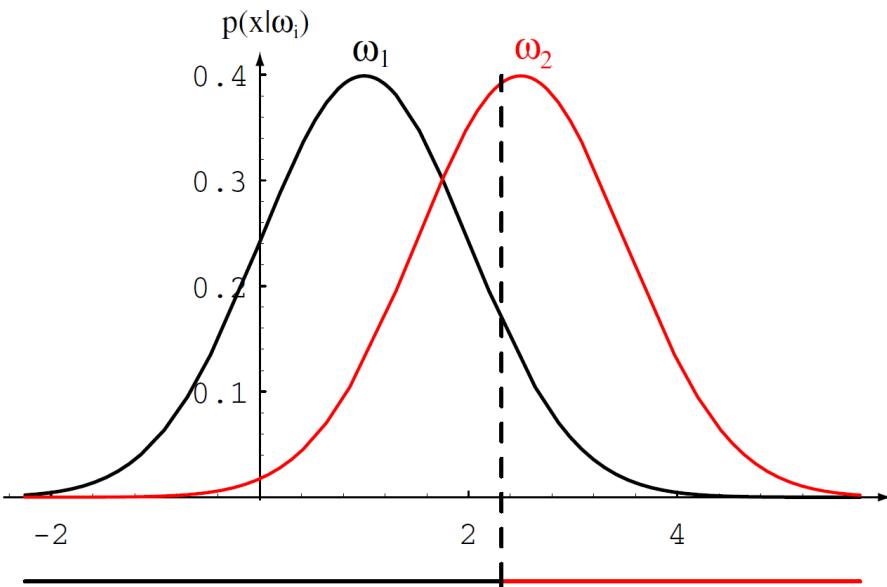
assign  $\mathbf{x}$  to  $w_{i^*}$  where  $i^* = \arg \min_{i=1,\dots,c} \|\mathbf{x} - \boldsymbol{\mu}_i\|$ .



# Shift of decision boundary by changing priors

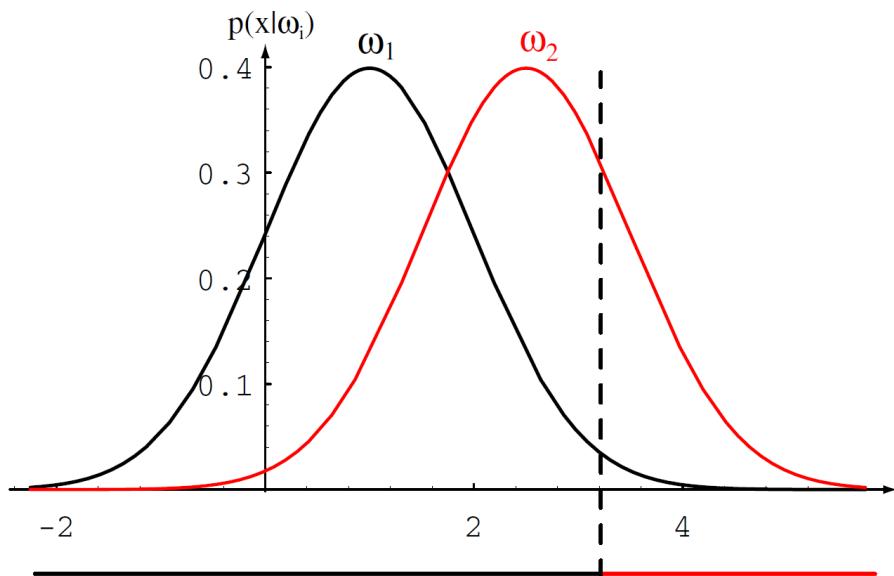
- If  $P(\omega_i) \neq P(\omega_j)$ , then  $x_0$  shifts away from the most likely category.
- If  $\sigma$  is very small, the position of the boundary is insensitive to  $P(\omega_i)$  and  $P(\omega_j)$

$$w_{i0} = \frac{-1}{2\sigma^2} \boldsymbol{\mu}_i^t \boldsymbol{\mu}_i + \ln P(\omega_i).$$



$\mathcal{R}_1$   
 $P(\omega_1) = .7$

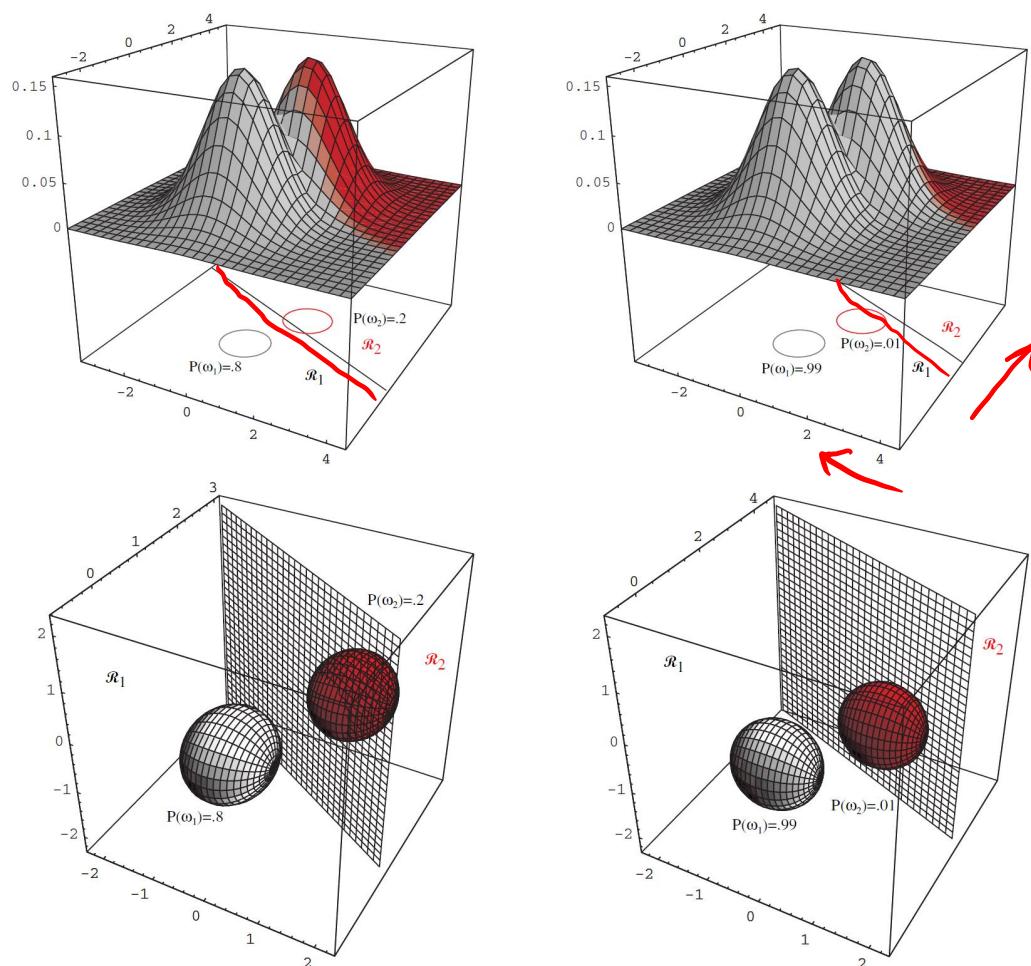
$\mathcal{R}_2$   
 $P(\omega_2) = .3$



$\mathcal{R}_1$   
 $P(\omega_1) = .9$

$\mathcal{R}_2$   
 $P(\omega_2) = .1$

$$\Sigma_i = \underline{\sigma^2 \mathbf{I}}$$



The distributions are **spherical in d dimensions**, and the boundary is a **generalized hyperplane of  $d-1$  dimensions, perpendicular to the line separating the means**.

The decision boundary **shifts** as the priors are changed.



## Case 2: $\Sigma_i = \Sigma$

- When the **covariance matrices for all of the classes are identical** but otherwise arbitrary
- Samples fall in **hyperellipsoidal** clusters of **equal size and shape**
- $|\Sigma_i|$  and the  $(d/2) \ln 2\pi$  term are **independent of  $i$** 
$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^t \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) + \ln P(\omega_i).$$
- To classify a feature vector  $\mathbf{x}$ , measure the squared **Mahalanobis distance from  $\mathbf{x}$**  to each of the  $c$  mean vectors, and assign  $\mathbf{x}$  to the category of the **nearest mean**



# Discriminant functions

- In the expansion of quadratic form  $(\mathbf{x} - \mathbf{m}_i)^T \mathbf{S}^{-1} (\mathbf{x} - \mathbf{m}_i)$ , the **quadratic** term  $\mathbf{x}^T \Sigma^{-1} \mathbf{x}$  is independent of  $i$ ; so

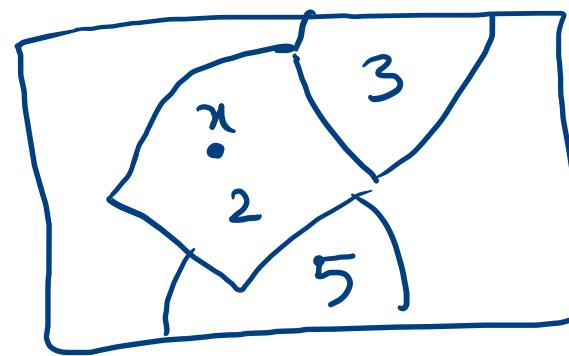
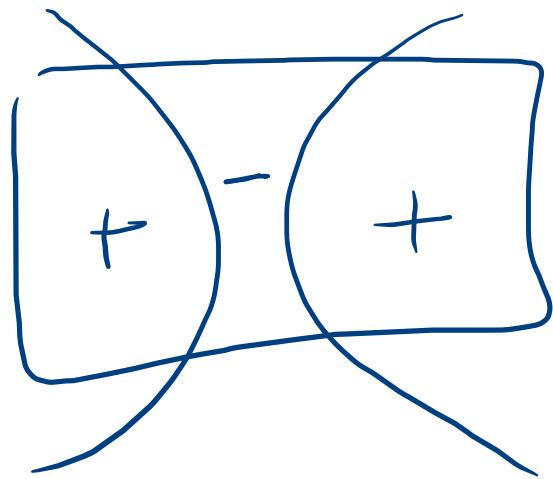
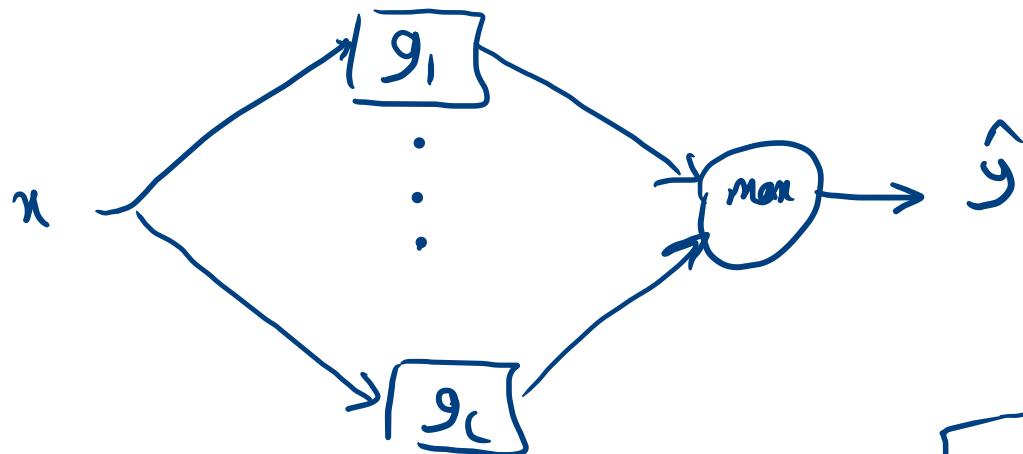
$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0},$$

- Where

$$\mathbf{w}_i = \underline{\Sigma^{-1} \mu_i}$$

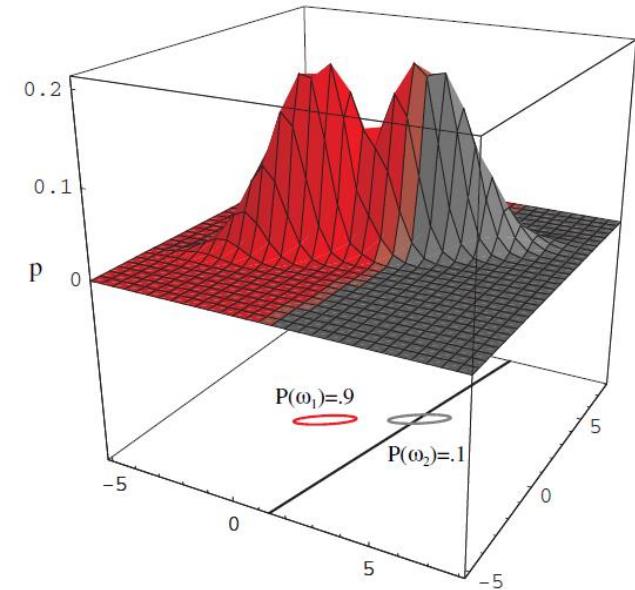
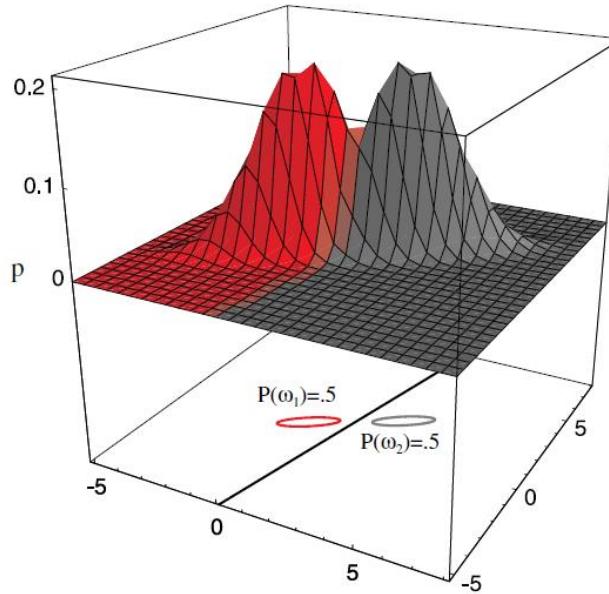
$$\omega = \Sigma^{-1} (\mu_1 - \mu_2)$$

$$w_{i0} = -\frac{1}{2} \mu_i^T \Sigma^{-1} \mu_i + \ln P(\omega_i).$$

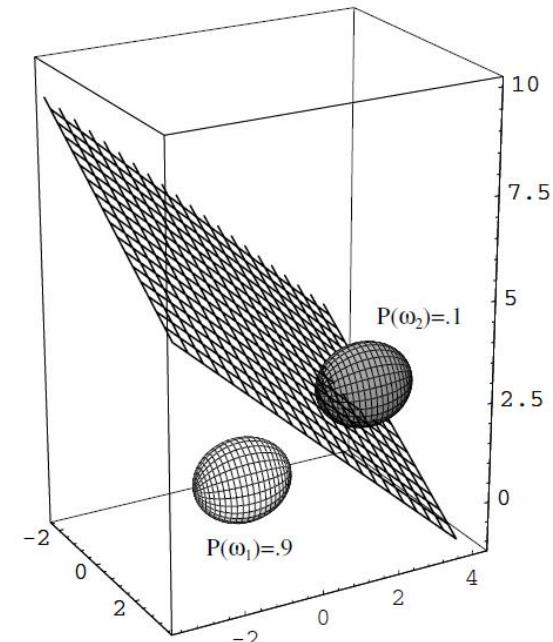
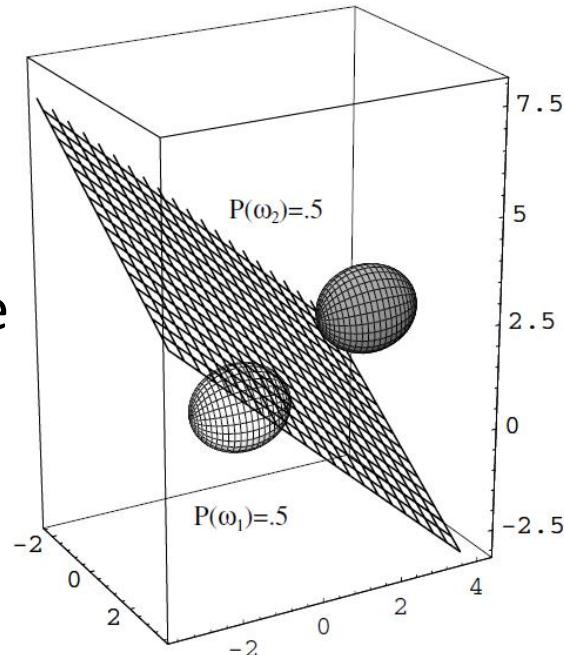


1

Probability densities  
with equal but  
**asymmetric**  
**Gaussian**  
**distributions.**



The decision  
hyperplanes **are not**  
**necessarily**  
**perpendicular** to the  
line connecting the  
means.





## Case 3: $\Sigma_i = \text{arbitrary}$

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^t \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}_i| + \ln P(\omega_i).$$

- Only  $(d/2) \ln 2\pi$  term will dropped

- Discriminant functions are:

$$g_i(\mathbf{x}) = \mathbf{x}^t \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^t \mathbf{x} + w_{i0},$$

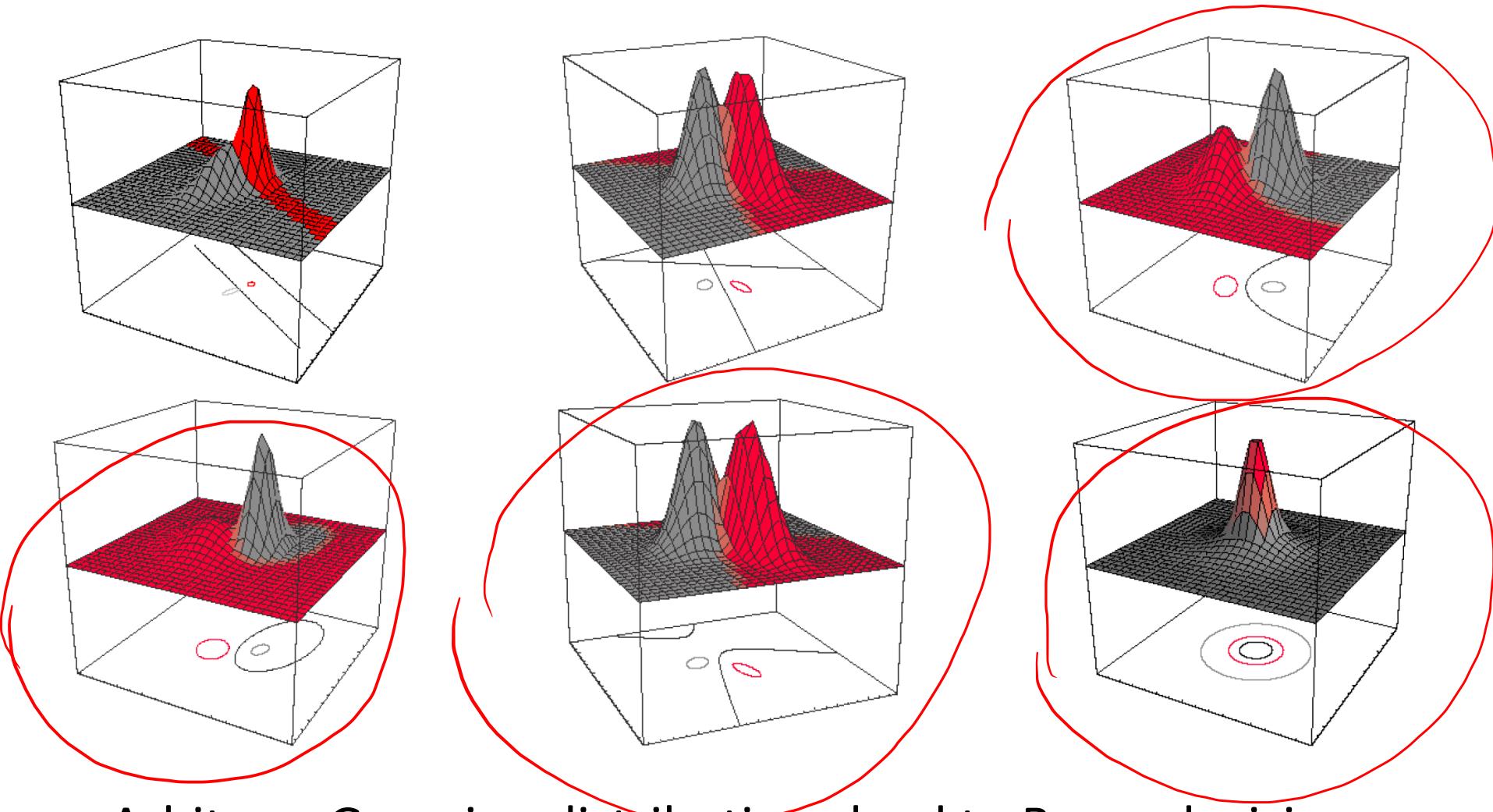
- Where

$$\mathbf{W}_i = -\frac{1}{2} \boldsymbol{\Sigma}_i^{-1},$$

$$\mathbf{w}_i = \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i$$

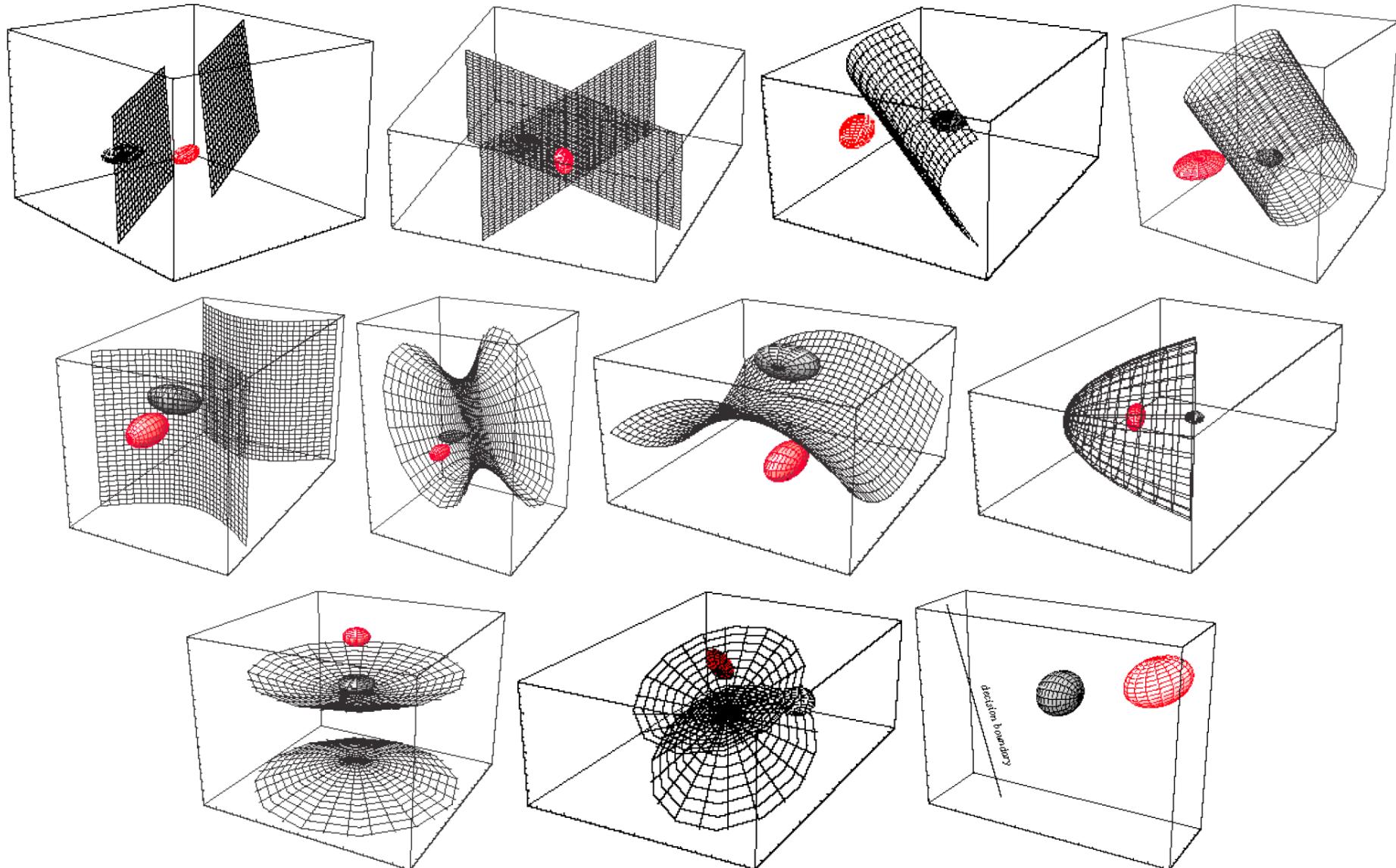
$$w_{i0} = -\frac{1}{2} \boldsymbol{\mu}_i^t \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \ln |\boldsymbol{\Sigma}_i| + \ln P(\omega_i).$$

- Decision boundaries are **hyperquadrics**



Arbitrary Gaussian distributions lead to Bayes decision  
boundaries that are **general hyperquadrics**.

# Case 3: $\Sigma_i = \text{arbitrary}$

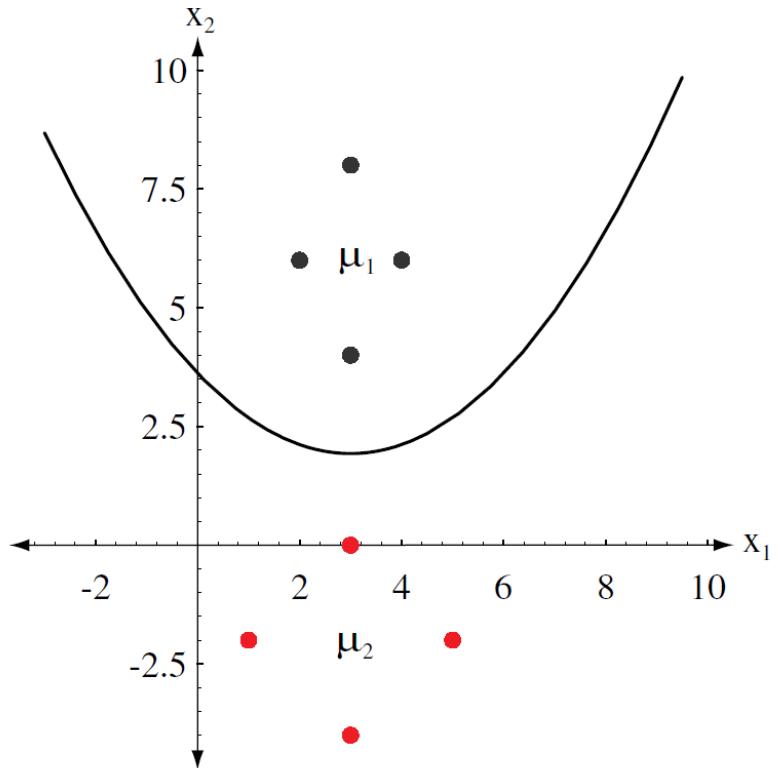




## Decision regions for two-dimensional Gaussian data

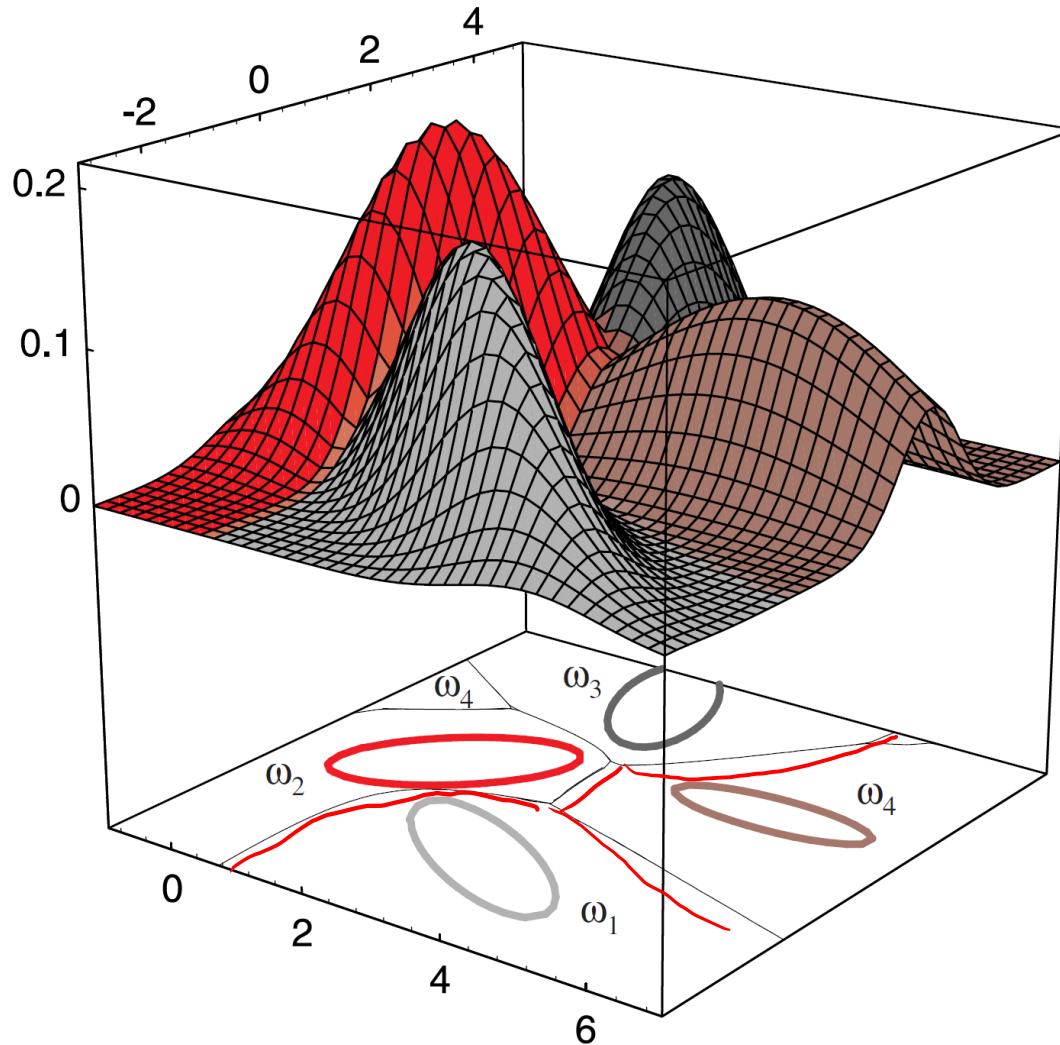
$$\mu_1 = \begin{bmatrix} 3 \\ 6 \end{bmatrix}; \quad \Sigma_1 = \begin{pmatrix} 1/2 & 0 \\ 0 & 2 \end{pmatrix} \text{ and } \mu_2 = \begin{bmatrix} 3 \\ -2 \end{bmatrix}; \quad \Sigma_2 = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}.$$

$$P(\omega_1) = P(\omega_2) = 0.5,$$



decision boundary:

$$x_2 = 3.514 - 1.125x_1 + 0.1875x_1^2.$$



$$\sum_i \cdot = \sum$$

The **decision regions** for four normal distributions. Even with such a **low number of categories**, the shapes of the boundary regions can be rather complex.

# Confusion Matrix:

$acc = 99.9$

		Predicted Class			
		Positive	Negative		
Actual Class	Positive	True Positive (TP)	False Negative (FN) <b>Type II Error</b>	Sensitivity	$= Recall = TPR$
	Negative	False Positive (FP) <b>Type I Error</b>	True Negative (TN)	Specificity	$= \frac{TP}{P}$
	Precision	$\frac{TP}{(TP + FP)}$	Negative Predictive Value	$\frac{TN}{(TN + FN)}$	$= \frac{TN}{N}$
				Accuracy	$\frac{TP + TN}{(TP + TN + FP + FN)}$

$FPR = \frac{FP}{N}$

$+ -$

$\Rightarrow$

$\Rightarrow TNR = 1 - FPR$

Actual	predict
+	
30	25 +
	5 -
20	14 -
	5 +

$$TPR = \frac{25}{30}$$

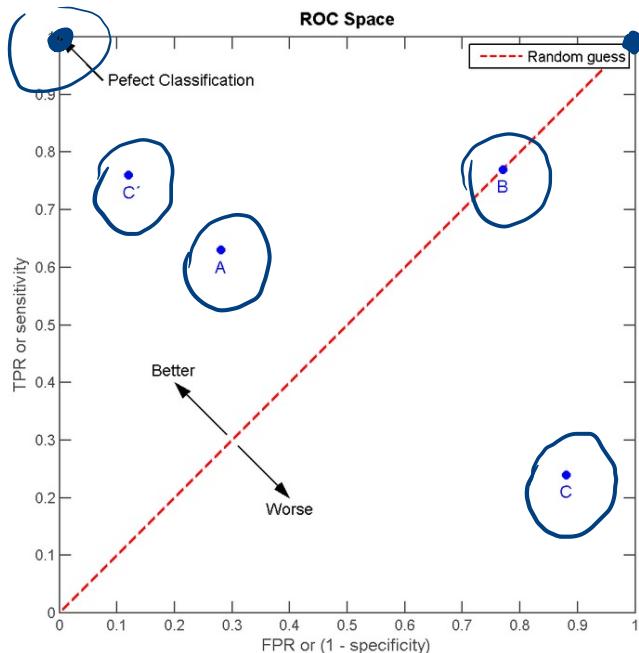
$$\text{precision} = \frac{25}{31}$$

$$TNR = \frac{14}{20}$$

Receiver operating characteristic (ROC) Curve

$A \rightarrow TPR = 1$

$A \rightarrow FPR = 0$

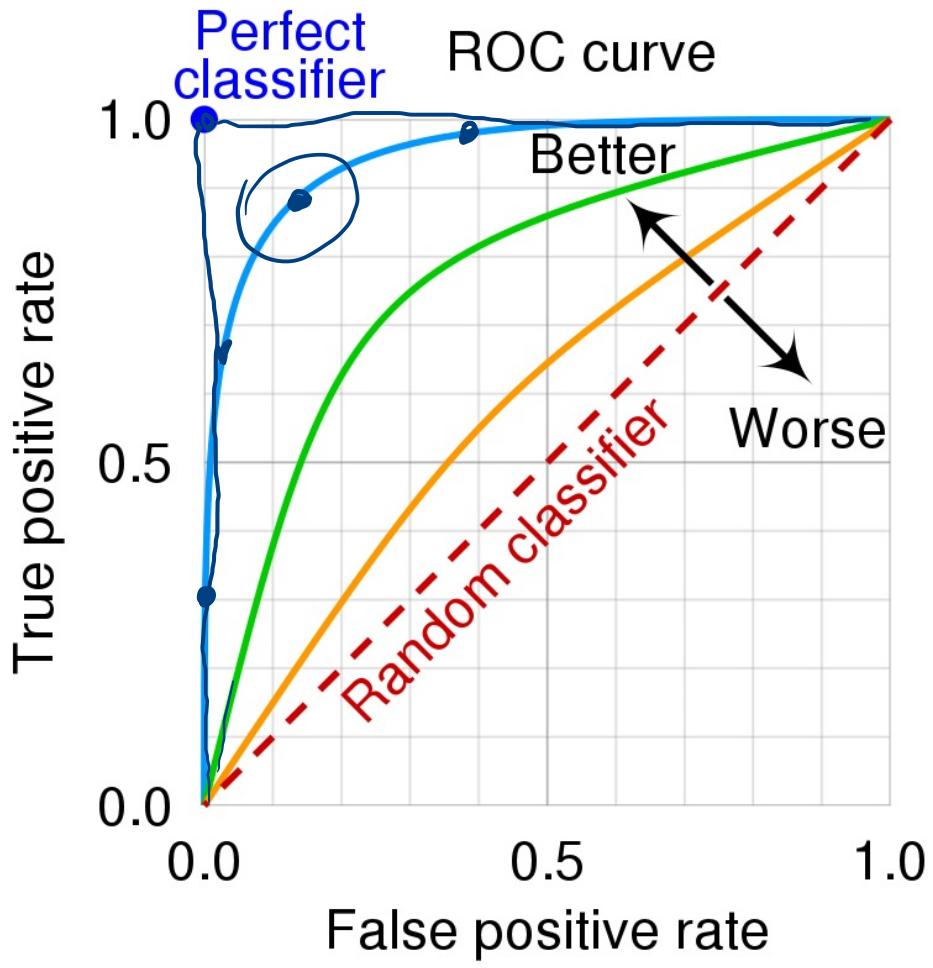


$$FPR = 1$$

$$TPR = 1$$

$FPR = 1 - \text{specificity}$

$$= \frac{FP}{N}$$



$$g(x) > \theta \uparrow$$

$$P(x|\underline{\theta}) \quad D = \{x_1, x_2, \dots, x_n\}$$

$$\text{Find } \theta = ? \quad \theta = \{ \mu, \Sigma \}$$

مین بدر توزیع احتمال

Parametric :  $P(x|\theta)$

Non-Parametric

Parameter Est.

{

frequentist :  $\theta$  is fixed, unknown

Bayesian :  $\theta$  is random var. ,  $\theta \sim P(\theta)$