

# Bayes Decision Theory

Minimum-Error-Rate Classification

Classifiers, Discriminant Functions and Decision Surfaces

The Normal Density

# Minimum-Error-Rate Classification

- Actions are decisions on classes

If action  $\alpha_i$  is taken and the true state of nature is  $\omega_j$   
then: decision is correct if  $i = j$  and in error if  $i \neq j$

Seek a decision rule that minimizes the  
*probability of error* which is the *error rate*

# Minimum Error Rate Classifier Derivation

- zero-one loss function:

$$\lambda(\alpha_i, \omega_j) = \begin{cases} 0 & i = j \\ 1 & i \neq j \end{cases} \quad i, j = 1, \dots, c$$

- Therefore, the conditional risk is:

$$\begin{aligned} R(\alpha_i / x) &= \sum_{j=1}^{j=c} \lambda(\alpha_i / \omega_j) P(\omega_j / x) \\ &= \sum_{j \neq i} P(\omega_j / x) = 1 - P(\omega_i / x) \end{aligned}$$

*The risk corresponding to this loss function is the average probability error"*

- Minimize the risk requires maximize  $P(\omega_i / x)$   
(since  $R(\alpha_i / x) = 1 - P(\omega_i / x)$ )  
For Minimum error rate  
Decide  $\omega_i$  if  $P(\omega_i / x) > P(\omega_j / x) \quad \forall j \neq i$

# Likelihood Ratio Classification

- Regions of decision and zero-one loss function,

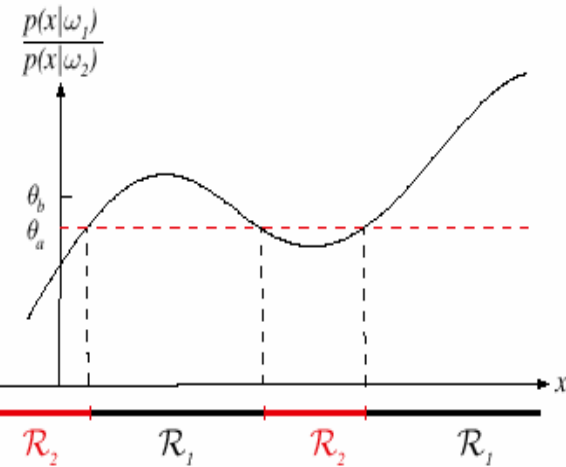
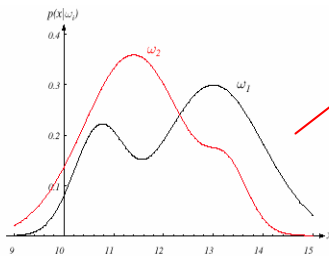
$$\text{Let } \frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}} \cdot \frac{P(\omega_2)}{P(\omega_1)} = \theta_\lambda \text{ then decide } \omega_1 \text{ if : } \frac{P(x/\omega_1)}{P(x/\omega_2)} > \theta_\lambda$$

- If  $\lambda$  is the zero-one loss

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

$$\text{then } \theta_\lambda = \frac{P(\omega_2)}{P(\omega_1)} = \theta_a$$

$$\text{if } \lambda = \begin{pmatrix} 0 & 2 \\ 1 & 0 \end{pmatrix} \text{ then } \theta_\lambda = \frac{2P(\omega_2)}{P(\omega_1)} = \theta_b$$

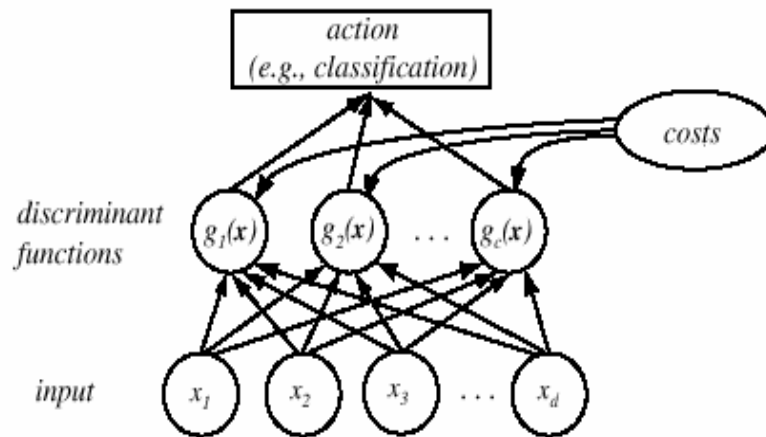


- Likelihood Ratio  $p(x/\omega_1)/p(x/\omega_2)$ .
- If we use a zero-one loss function decision boundaries are determined by threshold  $\theta_a$ .
- If loss function penalizes miscategorizing  $\omega_2$  as  $\omega_1$  more than converse we get larger threshold  $\theta_b$  and hence  $\mathbf{R}_1$  becomes smaller

Class-conditional pdfs

# Classifiers, Discriminant Functions and Decision Surfaces

- Many methods of representing pattern classifiers  
Set of discriminant functions  $g_i(x), i = 1, \dots, c$   
Classifier assigns feature  $x$  to class  $\omega_i$  if  $g_i(x) > g_j(x) \forall j \neq i$



Classifier is a machine  
that computes  $c$   
discriminant functions

Functional structure of a general statistical pattern  
Classifier with  $d$  inputs and  $c$  discriminant functions  
 $g_i(x)$

# Forms of Discriminant Functions

- Let  $g_i(x) = -R(\omega_i | x)$   
(max. discriminant corresponds to min. risk!)

- For the minimum error rate, we take

$$g_i(x) = P(\omega_i | x)$$

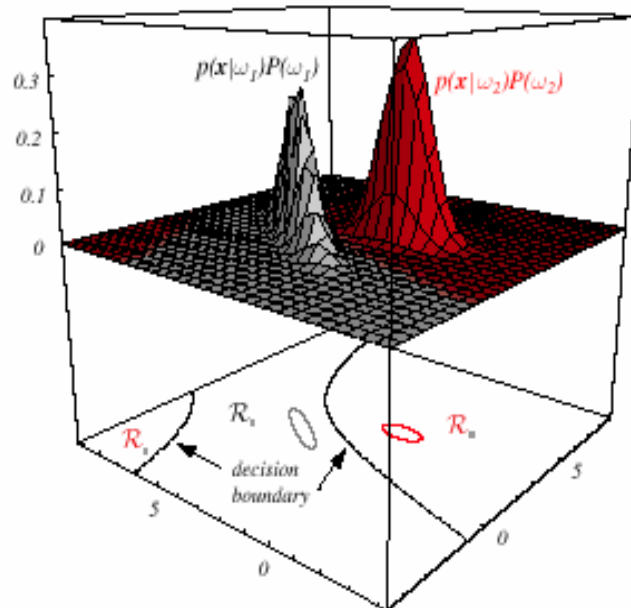
(max. discrimination corresponds to max. posterior!)

$$g_i(x) \equiv P(x | \omega_i) P(\omega_i)$$

$$g_i(x) = \ln P(x | \omega_i) + \ln P(\omega_i)$$

# Decision Region

- Feature space divided into  $c$  decision regions  
*if  $g_i(x) > g_j(x) \forall j \neq i$  then  $x$  is in  $R_i$*



2-D, two-category classifier with Gaussian pdfs

Decision Boundary = two hyperbolas

Hence decision region  $R_2$  is not simply connected

Ellipses mark where density is  $1/e$  times that of peak distribution

# The Two-Category case

A classifier is a *dichotomizer*

that has two discriminant functions  $g_1$  and  $g_2$

Let  $g(x) \equiv g_1(x) - g_2(x)$

Decide  $\omega_1$  if  $g(x) > 0$ ; Otherwise decide  $\omega_2$

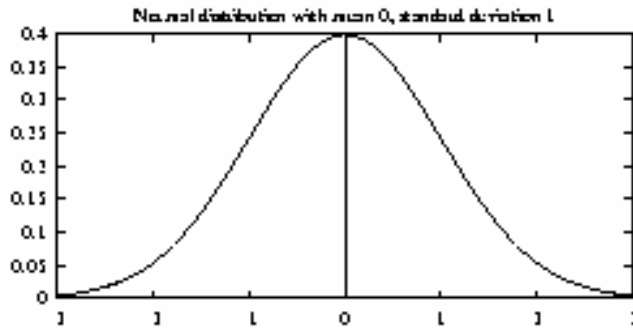
The computation of  $g(x)$

$$\begin{aligned} g(x) &= P(\omega_1 | x) - P(\omega_2 | x) \\ &= \ln \frac{P(x | \omega_1)}{P(x | \omega_2)} + \ln \frac{P(\omega_1)}{P(\omega_2)} \end{aligned}$$



# The Normal Distribution

A bell-shaped distribution defined by the probability density function



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

If the random variable  $X$  follows a normal distribution, then

- The probability that  $X$  will fall into the interval  $(a, b)$  is given by

$$\int_a^b p(x) dx$$

- Expected, or mean, value of  $X$  is  $E[X] = \int_{-\infty}^{\infty} xp(x) dx = \mu$

- Variance of  $X$  is

$$\text{Var}(x) = E[(x - \mu)^2] = \int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx = \sigma^2$$

- Standard deviation of  $X$ ,  $\sigma$ , is

$$\sigma_x = \sigma$$

# Relationship between Entropy and Normal Density

Entropy of a distribution

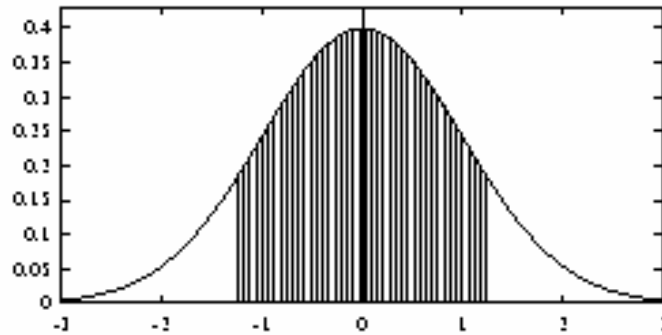
$$H(p(x)) = \int_{-\infty}^{\infty} p(x) \ln p(x) dx$$

Measured in nats. If  $\log_2$  is used the unit is bits

Entropy measures uncertainty in the values of points selected randomly from a distribution

Normal distribution has maximum entropy over all distributions having a given mean and variance

# Normal Distribution, Mean 0, Standard Deviation 1



80% of area (probability) lies in  $\mu \pm 1.28\sigma$

N% of area (probability) lies in  $\mu \pm z_N\sigma$

N%:	50%	68%	80%	90%	95%	98%	99%
$z_N$ :	0.67	1.00	1.28	1.64	1.96	2.33	2.58

With 80% confidence the r.v. will lie in the two-sided interval[-1.28,1.28]

# The Normal Density in Pattern Recognition

- Univariate density

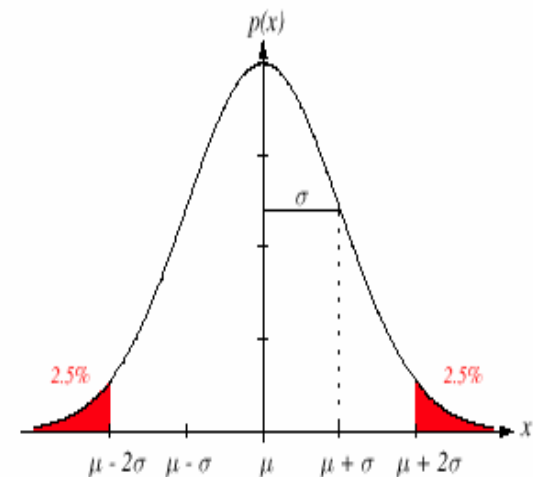
- Analytically tractable, continuous
- A lot of processes are asymptotically Gaussian
- Central Limit Theorem: aggregate effect of a sum of a large number of small, independent random disturbances will lead to a Gaussian distribution
- Handwritten characters, speech sounds are ideal or prototype corrupted by random process

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

Where:

$\mu$  = mean (or expected value) of  $x$

$\sigma^2$  = expected squared deviation  
or variance



CSE 555: Srihar

Univariate normal distribution has roughly 95% of its area in the range

$|x-\mu| < 2\sigma$ .

The peak of the distribution has value  $p(\mu) = 1/\sqrt{2\pi}\sigma$

# Multivariate density

Multivariate normal density in  $d$  dimensions is:

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

*abbreviated as*

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

where:

$x = (x_1, x_2, \dots, x_d)^t$  (t stands for the transpose vector form)

$\mu = (\mu_1, \mu_2, \dots, \mu_d)^t$  mean vector

$\Sigma = d \times d$  covariance matrix

$|\Sigma|$  and  $\Sigma^{-1}$  are determinant and inverse respectively

# Mean and Covariance Matrix

- Formal Definitions

$$\mu = E[x] = \int xp(x)dx$$

$$\Sigma = E[(x - \mu)(x - \mu)^t] = \int (x - \mu)(x - \mu)^t p(x)dx$$

Mean vector has its components which are means of variables

Covariance :

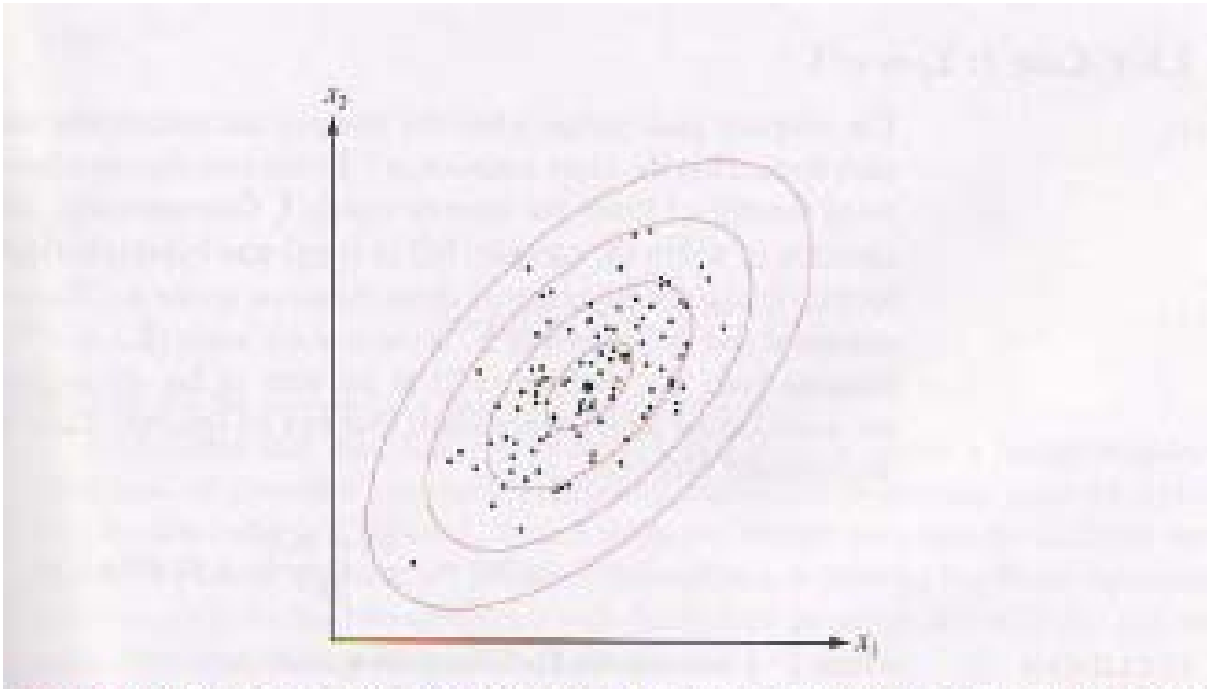
Diagonal elements are variances of variables

Cross-diagonal elements are covariances of pairs of variables

Statistical independence means off-diagonal elements are zero

# Multivariate Normal Density

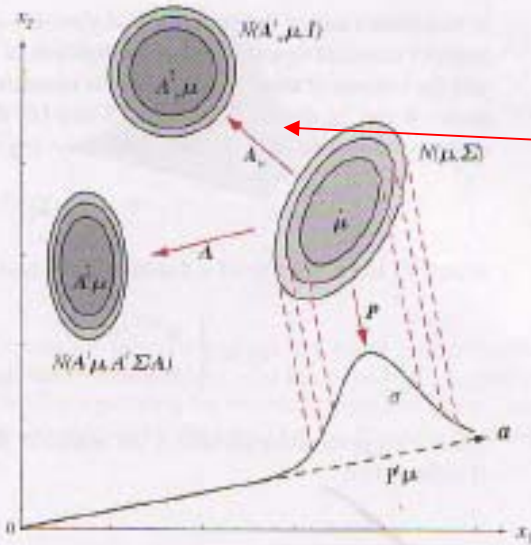
- Specified by  $d + d(d+1)/2$  parameters: mean and independent elements of covariance matrix



Locii of points of constant density are hyperellipsoids

Samples drawn from a 2-D Gaussian lie in a cloud centered at the mean  $\mu$ . Ellipses show lines of equal probability density of the Gaussian

# Linear Combinations of Normally distributed variables are normally distributed



$\Phi$  is matrix whose columns are the orthonormal Eigen vectors of  $\Sigma$  and  $A$  is diagonal matrix of corresponding Eigen values then the transformation  $A_w = \Phi A^{-1/2}$  applied to the coordinates ensures that transformed distribution has covariance matrix equal to the identity matrix

Action of a linear transformation on the feature space will convert an arbitrary normal distribution into another normal distribution. One transformation  $A$ , takes the source distribution into distribution  $N(A^t\mu, A^t\Sigma A)$ . Another linear transformation- a projection  $P$  onto a Line defined by vector  $a$ - leads to  $N(\mu, \sigma^2)$  measured along that line. While the transforms yield distributions in a different space they are shown superimposed on the original  $x_1$ - $x_2$  space. A whitening transform  $A_w$  leads to a circularly symmetric Gaussian, here shown displaced.



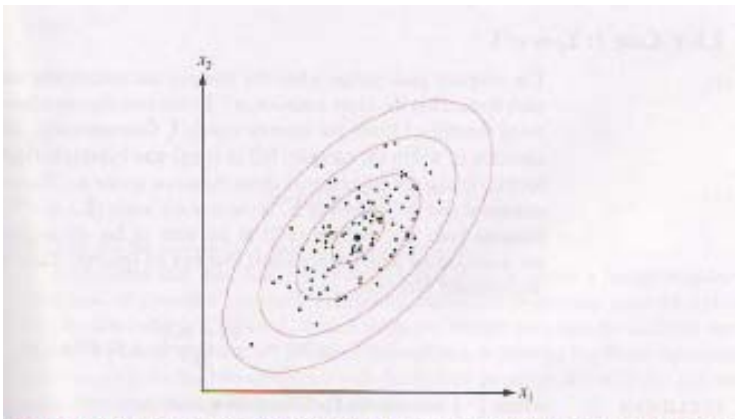
# Mahalanobis Distance

$$r^2 = (x - \mu)^t \Sigma^{-1} (x - \mu)$$

is the Mahalanobis distance from  $x$  to  $\mu$

Contours of constant Density are hyperellipsoids of constant Mahalanobis Distance

For a given dimensionality Scatter of samples varies directly with  $|E|^{1/2}$



Samples drawn from a 2-D Gaussian lie in a cloud centered at the mean  $\mu$ . Ellipses show lines of equal probability density of the Gaussian