

CS434a/541a: Pattern Recognition
Prof. Olga Veksler

Lecture 4

Normal Random Variable and its
discriminant functions

1

Outline

- Normal Random Variable
 - Properties
 - Discriminant functions

3

Announcement

- Assignment 1 has been posted
 - Note changes to problem 3 and 6 made today
 - Problem 3(d) corrections to 0.99 and 0.01
 - Problem 6, c = number of classes

2

Why Normal Random Variables?

- Analytically tractable
- Works well when observation comes form a corrupted single prototype (μ)
- Is an optimal distribution of data for many classifiers used in practice

4

The Univariate Normal Density

- x is a scalar (has dimension 1)

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

Where:

μ = mean (or expected value) of x

σ^2 = variance

5

Several Features

- What if we have several features x_1, x_2, \dots, x_d
 - each normally distributed
 - may have different means
 - may have different variances
 - may be dependent or independent of each other
- How do we model their joint distribution?

7

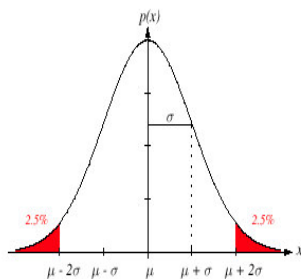


FIGURE 2.7. A univariate normal distribution has roughly 95% of its area in the range $|x - \mu| \leq 2\sigma$, as shown. The peak of the distribution has value $p(\mu) = 1/\sqrt{2\pi}\sigma$. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

6

The Multivariate Normal Density

- Multivariate normal density in d dimensions is:

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

determinant of Σ

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \dots & \sigma_{1d} \\ \vdots & \ddots & \vdots \\ \sigma_{d1} & \dots & \sigma_d^2 \end{bmatrix}$$

covariance of x_1 and x_d

$$\mathbf{x} = [x_1, x_2, \dots, x_d]^t$$

$$\boldsymbol{\mu} = [\mu_1, \mu_2, \dots, \mu_d]^t$$

- Each x_i is $N(\mu_i, \sigma_i^2)$
 - to prove this, integrate out all other features from the joint density

8

More on Σ

▪ $\Sigma = \begin{bmatrix} \sigma_1^2 & \dots & \sigma_{1d} \\ \vdots & \ddots & \vdots \\ \sigma_{d1} & \dots & \sigma_d^2 \end{bmatrix}$ plays role similar to the role that σ^2 plays in one dimension

- From Σ we can find out
 1. The individual variances of features x_1, x_2, \dots, x_d
 2. If features x_i and x_j are
 - independent $\sigma_{ij}=0$
 - have positive correlation $\sigma_{ij}>0$
 - have negative correlation $\sigma_{ij}<0$

9

The Multivariate Normal Density

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

$$p(\mathbf{x}) = c \cdot \exp\left[-\frac{1}{2} \begin{bmatrix} x_1 - \mu_1 & x_2 - \mu_2 & x_3 - \mu_3 \end{bmatrix} \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \sigma_{13} \\ \sigma_{12} & \sigma_2^2 & \sigma_{23} \\ \sigma_{13} & \sigma_{23} & \sigma_3^2 \end{bmatrix}^{-1} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \\ x_3 - \mu_3 \end{bmatrix}\right]$$

normalizing constant scalar s (single number), the closer s to 0 the larger is $p(\mathbf{x})$

- Thus $P(\mathbf{x})$ is larger for smaller $(\mathbf{x} - \boldsymbol{\mu})^t \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})$

11

The Multivariate Normal Density

- If Σ is diagonal $\begin{bmatrix} \sigma_1^2 & 0 & 0 \\ 0 & \sigma_2^2 & 0 \\ 0 & 0 & \sigma_3^2 \end{bmatrix}$ then the features x_1, \dots, x_j are independent, and

$$p(\mathbf{x}) = \prod_{i=1}^d \frac{1}{\sigma_i \sqrt{2\pi}} \exp\left[-\frac{(x_i - \mu_i)^2}{2\sigma_i^2}\right]$$

10

$(\mathbf{x} - \boldsymbol{\mu})^t \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})$

- Σ is positive semi definite ($\mathbf{x}^t \Sigma \mathbf{x} \geq 0$)
- If $\mathbf{x}^t \Sigma \mathbf{x} = 0$ for nonzero \mathbf{x} then $\det(\Sigma) = 0$. This case is not interesting, $p(\mathbf{x})$ is not defined
 1. one feature vector is a constant (has zero variance)
 2. or two components are multiples of each other
- so we will assume Σ is positive definite ($\mathbf{x}^t \Sigma \mathbf{x} > 0$)
- If Σ is positive definite then so is Σ^{-1}

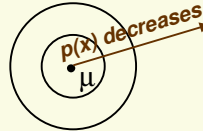
$$(\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

- Positive definite matrix of size d by d has d distinct real eigenvalues and its d eigenvectors are orthogonal
- Thus if Φ is a matrix whose columns are normalized eigenvectors of $\boldsymbol{\Sigma}$, then $\Phi^{-1} = \Phi^t$
- $\boldsymbol{\Sigma}\Phi = \Phi\boldsymbol{\Lambda}$ where $\boldsymbol{\Lambda}$ is a diagonal matrix with corresponding eigenvalues on the diagonal
- Thus $\boldsymbol{\Sigma} = \Phi\boldsymbol{\Lambda}\Phi^{-1}$ and $\boldsymbol{\Sigma}^{-1} = \Phi\boldsymbol{\Lambda}^{-1}\Phi^{-1}$
- Thus if $\boldsymbol{\Lambda}^{-1/2}$ denotes matrix s.t. $\boldsymbol{\Lambda}^{-1/2}\boldsymbol{\Lambda}^{-1/2} = \boldsymbol{\Lambda}^{-1}$

$$\boldsymbol{\Sigma}^{-1} = \left(\Phi\boldsymbol{\Lambda}^{-1/2}\right)\left(\Phi\boldsymbol{\Lambda}^{-1/2}\right)^t = \mathbf{M}\mathbf{M}^t$$

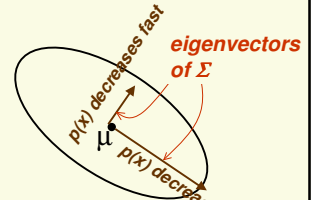
$$(\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

$(\mathbf{x} - \boldsymbol{\mu})^t (\mathbf{x} - \boldsymbol{\mu})$
usual (Euclidian)
distance between \mathbf{x} and $\boldsymbol{\mu}$



points \mathbf{x} at equal
Euclidian
distance from $\boldsymbol{\mu}$
lie on a circle

$(\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$
Mahalanobis distance
between \mathbf{x} and $\boldsymbol{\mu}$



points \mathbf{x} at equal
Mahalanobis distance from
 $\boldsymbol{\mu}$ lie on an ellipse: $\boldsymbol{\Sigma}$
stretches circles to ellipses

$$(\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

- Thus $(\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = (\mathbf{x} - \boldsymbol{\mu})^t \mathbf{M}\mathbf{M}^t (\mathbf{x} - \boldsymbol{\mu}) = (\mathbf{M}^t (\mathbf{x} - \boldsymbol{\mu}))^t (\mathbf{M}^t (\mathbf{x} - \boldsymbol{\mu})) = |\mathbf{M}^t (\mathbf{x} - \boldsymbol{\mu})|^2$

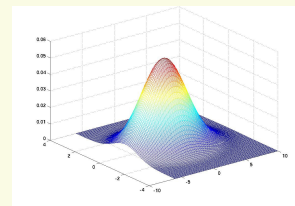
Thus $(\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = |\mathbf{M}^t (\mathbf{x} - \boldsymbol{\mu})|^2$

where $\mathbf{M}^t = \boldsymbol{\Lambda}^{-1/2} \Phi^{-1}$
scaling rotation
matrix matrix

- Points \mathbf{x} which satisfy $|\mathbf{M}^t (\mathbf{x} - \boldsymbol{\mu})|^2 = \text{const}$ lie on an ellipse

2-d Multivariate Normal Density

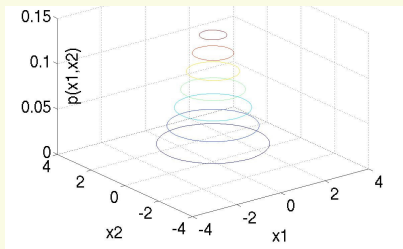
- Can you see much in this graph?



- At most you can see that the mean is around $[0,0]$, but can't really tell if \mathbf{x}_1 and \mathbf{x}_2 are correlated

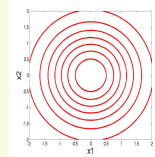
2-d Multivariate Normal Density

- How about this graph?



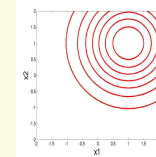
17

2-d Multivariate Normal Density



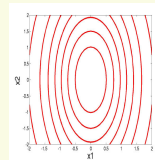
$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = [0, 0]$$



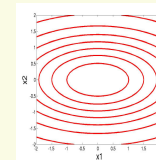
$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = [1, 1]$$



$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 4 \end{bmatrix}$$

$$\mu = [0, 0]$$



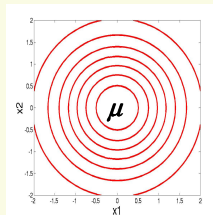
$$\Sigma = \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = [0, 0]$$

19

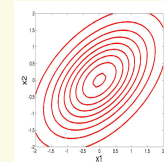
2-d Multivariate Normal Density

- Level curves graph
 - $p(x)$ is constant along each contour
 - topological map of 3-d surface
- Now we can see much more
 - x_1 and x_2 are independent
 - σ_1^2 and σ_2^2 are equal

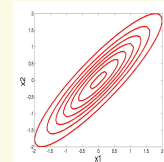


18

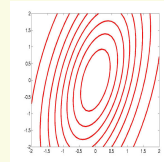
2-d Multivariate Normal Density $\mu = [0, 0]$



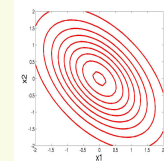
$$\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$



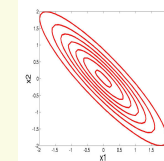
$$\Sigma = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$$



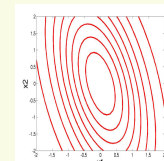
$$\Sigma = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 4 \end{bmatrix}$$



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



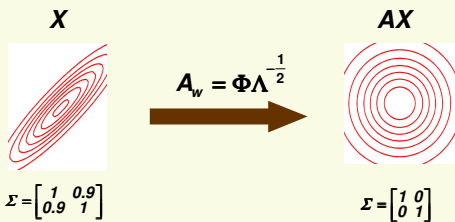
$$\Sigma = \begin{bmatrix} 1 & -0.9 \\ -0.9 & 1 \end{bmatrix}$$



$$\Sigma = \begin{bmatrix} 1 & -0.9 \\ -0.9 & 4 \end{bmatrix}$$

The Multivariate Normal Density

- If X has density $N(\mu, \Sigma)$ then AX has density $N(A\mu, A\Sigma A)$
 - Thus X can be transformed into a spherical normal variable (covariance of spherical density is the identity matrix I) with whitening transform



21

Discriminant Functions

- The minimum error-rate classification is achieved by the discriminant function

$$g_i(x) = P(c_i | x) = P(x|c_i)P(c_i)/P(x)$$

- Since the observation x is independent of the class, the equivalent discriminant function is

$$g_i(x) = P(x|c_i)P(c_i)$$

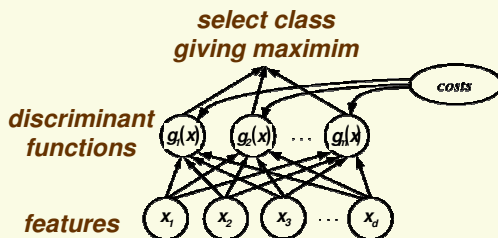
- For normal density, convenient to take logarithms. Since logarithm is a monotonically increasing function, the equivalent discriminant function is

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

23

Discriminant Functions

- Classifier can be viewed as network which computes m discriminant functions and selects category corresponding to the largest discriminant



- $g_i(x)$ can be replaced with any monotonically increasing function, the results will be unchanged

Discriminant Functions for the Normal Density

- Suppose we for class c_i its class conditional density $p(x|c_i)$ is $N(\mu_i, \Sigma_i)$

$$p(x | c_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left[-\frac{1}{2} (x - \mu_i)' \Sigma_i^{-1} (x - \mu_i) \right]$$

- Discriminant function $g_i(x) = \ln P(x|c_i) + \ln P(c_i)$

- Plug in $p(x|c_i)$ and $P(c_i)$ get

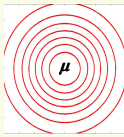
$$g_i(x) = -\frac{1}{2} (x - \mu_i)' \Sigma_i^{-1} (x - \mu_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(c_i)$$

constant for all i

$$g_i(x) = -\frac{1}{2} (x - \mu_i)' \Sigma_i^{-1} (x - \mu_i) - \frac{1}{2} \ln |\Sigma_i| + \ln P(c_i)$$

Case $\Sigma_i = \sigma^2 I$

- That is $\Sigma_i = \begin{bmatrix} \sigma^2 & 0 & 0 \\ 0 & \sigma^2 & 0 \\ 0 & 0 & \sigma^2 \end{bmatrix} = \sigma^2 \cdot \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$
- In this case, features x_1, x_2, \dots, x_d are independent with different means and equal variances σ^2

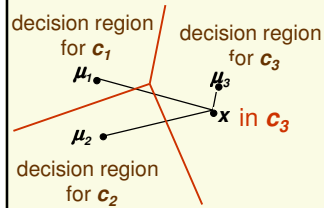


25

Case $\Sigma_i = \sigma^2 I$ Geometric Interpretation

If $\ln P(c_i) = \ln P(c_j)$, then

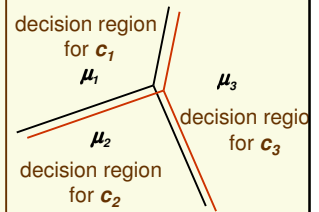
$$g_i(x) = -|x - \mu_i|^2$$



Voronoi diagram: points in each cell are closer to the mean in that cell than to any other mean

If $\ln P(c_i) \neq \ln P(c_j)$, then

$$g_i(x) = -\frac{1}{2\sigma^2}|x - \mu_i|^2 + \ln P(c_i)$$



Case $\Sigma_i = \sigma^2 I$

- Discriminant function
- $$g_i(x) = -\frac{1}{2}(x - \mu_i)^t \Sigma_i^{-1} (x - \mu_i) - \frac{1}{2} \ln |\Sigma_i| + \ln P(c_i)$$

- $\text{Det}(\Sigma_i) = \sigma^{2d}$ and $\Sigma_i^{-1} = (1/\sigma^2)I = \begin{bmatrix} \frac{1}{\sigma^2} & 0 & 0 \\ 0 & \frac{1}{\sigma^2} & 0 \\ 0 & 0 & \frac{1}{\sigma^2} \end{bmatrix}$

- Can simplify discriminant function

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^t \frac{1}{\sigma^2} (x - \mu_i) - \frac{1}{2} \ln(\sigma^{2d}) + \ln P(c_i)$$

constant for all i

$$g_i(x) = -\frac{1}{2\sigma^2}(x - \mu_i)^t (x - \mu_i) + \ln P(c_i) =$$

$$= -\frac{1}{2\sigma^2}|x - \mu_i|^2 + \ln P(c_i)$$

26

Case $\Sigma_i = \sigma^2 I$

$$g_i(x) = -\frac{1}{2\sigma^2}(x - \mu_i)^t (x - \mu_i) + \ln P(c_i) =$$

$$= -\frac{1}{2\sigma^2}(x^t x - \mu_i^t x - x^t \mu_i + \mu_i^t \mu_i) + \ln P(c_i)$$

constant for all classes

$$g_i(x) = -\frac{1}{2\sigma^2}(-2\mu_i^t x + \mu_i^t \mu_i) + \ln P(c_i) = \frac{\mu_i^t}{\sigma^2} x + \left(-\frac{\mu_i^t \mu_i}{2\sigma^2} + \ln P(c_i)\right)$$

$$g_i(x) = w_i^t x + w_{i0}$$

discriminant function is linear

28

Case $\Sigma_i = \sigma^2 I$

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

constant in \mathbf{x}

$$\mathbf{w}_i^T \mathbf{x} = \sum_{i=1}^d w_i x_i$$

linear in \mathbf{x} :

- Thus discriminant function is linear,
- Therefore the decision boundaries $g_i(\mathbf{x}) = g_j(\mathbf{x})$ are linear
 - lines if \mathbf{x} has dimension 2
 - planes if \mathbf{x} has dimension 3
 - hyper-planes if \mathbf{x} has dimension larger than 3

Case $\Sigma_i = \sigma^2 I$: Example

- Need to find out when $g_i(\mathbf{x}) < g_j(\mathbf{x})$ for $i, j = 1, 2, 3$
- Can be done by solving $g_i(\mathbf{x}) = g_j(\mathbf{x})$ for $i, j = 1, 2, 3$
- Let's take $g_1(\mathbf{x}) = g_2(\mathbf{x})$ first

$$\frac{[1 \ 2]}{3} \mathbf{x} + \left(-\frac{5}{6} - 1.38\right) = \frac{[4 \ 6]}{3} \mathbf{x} + \left(-\frac{52}{6} - 1.38\right)$$

- Simplifying, $\frac{[-3 \ -4]}{3} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = -\frac{47}{6}$

$$-x_1 - \frac{4}{3}x_2 = -\frac{47}{6}$$

line equation

31

Case $\Sigma_i = \sigma^2 I$: Example

- 3 classes, each 2-dimensional Gaussian with $\mu_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ $\mu_2 = \begin{bmatrix} 4 \\ 6 \end{bmatrix}$ $\mu_3 = \begin{bmatrix} -2 \\ 4 \end{bmatrix}$ $\Sigma_1 = \Sigma_2 = \Sigma_3 = \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}$
- Priors $P(c_1) = P(c_2) = \frac{1}{4}$ and $P(c_3) = \frac{1}{2}$
- Discriminant function is $g_i(\mathbf{x}) = \frac{\mu_i^T \mathbf{x}}{\sigma^2} + \left(-\frac{\mu_i^T \mu_i}{2\sigma^2} + \ln P(c_i)\right)$
- Plug in parameters for each class

$$g_1(\mathbf{x}) = \frac{[1 \ 2]}{3} \mathbf{x} + \left(-\frac{5}{6} - 1.38\right) \quad g_2(\mathbf{x}) = \frac{[4 \ 6]}{3} \mathbf{x} + \left(-\frac{52}{6} - 1.38\right)$$

$$g_3(\mathbf{x}) = \frac{[-2 \ 4]}{3} \mathbf{x} + \left(-\frac{20}{6} - 0.69\right)$$

30

Case $\Sigma_i = \sigma^2 I$: Example

- Next solve $g_2(\mathbf{x}) = g_3(\mathbf{x})$

$$2x_1 + \frac{2}{3}x_2 = 6.02$$
- Almost finally solve $g_1(\mathbf{x}) = g_3(\mathbf{x})$

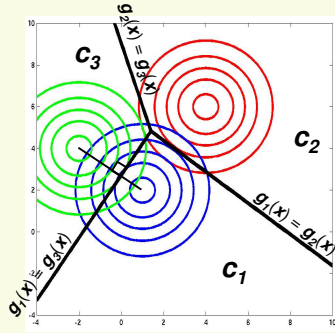
$$x_1 - \frac{2}{3}x_2 = -1.81$$
- And finally solve $g_1(\mathbf{x}) = g_2(\mathbf{x}) = g_3(\mathbf{x})$

$$x_1 = 1.4 \quad \text{and} \quad x_2 = 4.82$$

32

Case $\Sigma_i = \sigma^2 I$: Example

- Priors $P(c_1) = P(c_2) = \frac{1}{4}$ and $P(c_3) = \frac{1}{2}$



lines connecting means are perpendicular to decision boundaries

33

Case $\Sigma_i = \Sigma$

- Discriminant function

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^t \Sigma^{-1}(x - \mu_i) - \frac{1}{2} \ln |\Sigma_i| + \ln P(c_i)$$

constant for all classes

- Discriminant function becomes

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^t \Sigma^{-1}(x - \mu_i) + \ln P(c_i)$$

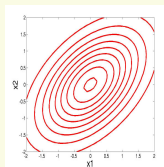
squared Mahalanobis Distance

- Mahalanobis Distance $\|x - y\|_{\Sigma^{-1}}^2 = (x - y)^t \Sigma^{-1}(x - y)$
- If $\Sigma = I$, Mahalanobis Distance becomes usual Euclidean distance

$$\|x - y\|_{I^{-1}}^2 = (x - y)^t (x - y)$$

Case $\Sigma_i = \Sigma$

- Covariance matrices are equal but arbitrary
- In this case, features x_1, x_2, \dots, x_d are not necessarily independent

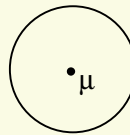


$$\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

34

Euclidean vs. Mahalanobis Distances

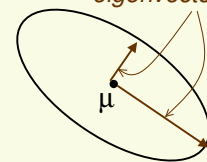
$$\|x - \mu\|^2 = (x - \mu)^t (x - \mu)$$



points x at equal Euclidean distance from μ lie on a circle

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

eigenvectors of Σ

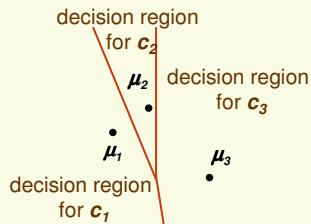


points x at equal Mahalanobis distance from μ lie on an ellipse: Σ stretches circles to ellipses

Case $\Sigma_i = \Sigma$ Geometric Interpretation

If $\ln P(c_i) = \ln P(c_j)$, then

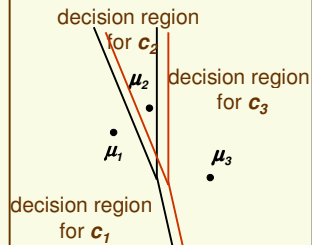
$$g_i(x) = -\|x - \mu_i\|_{\Sigma^{-1}}$$



points in each cell are closer to the mean in that cell than to any other mean under Mahalanobis distance

If $\ln P(c_i) \neq \ln P(c_j)$, then

$$g_i(x) = -\frac{1}{2}\|x - \mu_i\|_{\Sigma^{-1}} + \ln P(c_i)$$



Case $\Sigma_i = \Sigma$: Example

3 classes, each 2-dimensional Gaussian with

$$\mu_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad \mu_2 = \begin{bmatrix} -1 \\ 5 \end{bmatrix} \quad \mu_3 = \begin{bmatrix} -2 \\ 4 \end{bmatrix} \quad \Sigma_1 = \Sigma_2 = \Sigma_3 = \begin{bmatrix} 1 & -1.5 \\ -1.5 & 4 \end{bmatrix}$$

$$P(c_1) = P(c_2) = \frac{1}{4} \quad P(c_3) = \frac{1}{2}$$

Again can be done by solving $g_i(x) = g_j(x)$ for $i, j=1,2,3$

Case $\Sigma_i = \Sigma$

Can simplify discriminant function:

$$\begin{aligned} g_i(x) &= -\frac{1}{2}(x - \mu_i)' \Sigma^{-1}(x - \mu_i) + \ln P(c_i) = \\ &= -\frac{1}{2}(x' \Sigma^{-1} x - \mu_i' \Sigma^{-1} x - x' \Sigma^{-1} \mu_i + \mu_i' \Sigma^{-1} \mu_i) + \ln P(c_i) = \\ &= -\frac{1}{2}(x' \Sigma^{-1} x - 2\mu_i' \Sigma^{-1} x + \mu_i' \Sigma^{-1} \mu_i) + \ln P(c_i) = \\ &\quad \text{constant for all classes} \\ &= -\frac{1}{2}(-2\mu_i' \Sigma^{-1} x + \mu_i' \Sigma^{-1} \mu_i) + \ln P(c_i) \\ &= \mu_i' \Sigma^{-1} x + \left(\ln P(c_i) - \frac{1}{2} \mu_i' \Sigma^{-1} \mu_i \right) = w_i' x + w_{i0} \end{aligned}$$

Thus in this case discriminant is also linear

Case $\Sigma_i = \Sigma$: Example

Let's solve in general first

$$g_i(x) = g_j(x)$$

$$\mu_i' \Sigma^{-1} x + \left(\ln P(c_i) - \frac{1}{2} \mu_i' \Sigma^{-1} \mu_i \right) = \mu_j' \Sigma^{-1} x + \left(\ln P(c_j) - \frac{1}{2} \mu_j' \Sigma^{-1} \mu_j \right)$$

Let's regroup the terms

$$(\mu_i' \Sigma^{-1} - \mu_j' \Sigma^{-1}) x = \left(\ln P(c_j) - \frac{1}{2} \mu_j' \Sigma^{-1} \mu_j \right) - \left(\ln P(c_i) - \frac{1}{2} \mu_i' \Sigma^{-1} \mu_i \right)$$

We get the line where $g_i(x) = g_j(x)$

$$(\mu_j - \mu_i)' \Sigma^{-1} x = \left(\ln \frac{P(c_i)}{P(c_j)} + \frac{1}{2} \mu_j' \Sigma^{-1} \mu_j - \frac{1}{2} \mu_i' \Sigma^{-1} \mu_i \right)$$

row vector

scalar

Case $\Sigma_i = \Sigma$: Example

$$(\mu'_j - \mu'_i)\Sigma^{-1}x = \left(\ln \frac{P(c_i)}{P(c_j)} + \frac{1}{2}\mu'_j\Sigma^{-1}\mu_j - \frac{1}{2}\mu'_i\Sigma^{-1}\mu_i \right)$$

- Now substitute for $i,j=1,2$

$$\begin{bmatrix} -2 & 0 \end{bmatrix}x = 0 \\ x_1 = 0$$

- Now substitute for $i,j=2,3$

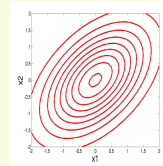
$$\begin{bmatrix} -3.14 & -1.4 \end{bmatrix}x = -2.41 \\ 3.14x_1 + 1.4x_2 = 2.41$$

- Now substitute for $i,j=1,3$

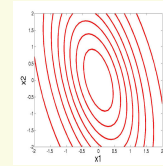
$$\begin{bmatrix} -5.14 & -1.43 \end{bmatrix}x = -2.41 \\ 5.14x_1 + 1.43x_2 = 2.41$$

General Case Σ_i are arbitrary

- Covariance matrices for each class are arbitrary
- In this case, features x_1, x_2, \dots, x_d are not necessarily independent



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

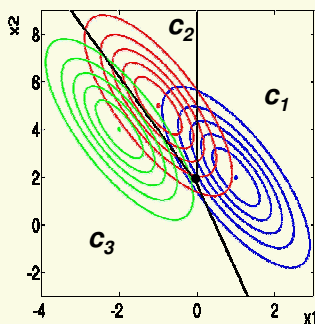


$$\Sigma_i = \begin{bmatrix} 1 & -0.9 \\ -0.9 & 4 \end{bmatrix}$$

43

Case $\Sigma_i = \Sigma$: Example

- Priors $P(c_1) = P(c_2) = \frac{1}{4}$ and $P(c_3) = \frac{1}{2}$



lines connecting means are **not** in general perpendicular to decision boundaries

42

General Case Σ_i are arbitrary

- From previous discussion,

$$g_i(x) = -\frac{1}{2}(x - \mu_i)' \Sigma_i^{-1}(x - \mu_i) - \frac{1}{2} \ln |\Sigma_i| + \ln P(c_i)$$

- This can't be simplified, but we can rearrange it:

$$g_i(x) = -\frac{1}{2}(x' \Sigma_i^{-1} x - 2\mu_i' \Sigma_i^{-1} x + \mu_i' \Sigma_i^{-1} \mu_i) - \frac{1}{2} \ln |\Sigma_i| + \ln P(c_i)$$

$$g_i(x) = x' \left(-\frac{1}{2} \Sigma_i^{-1} \right) x + \mu_i' \Sigma_i^{-1} x + \left(-\frac{1}{2} \mu_i' \Sigma_i^{-1} \mu_i - \frac{1}{2} \ln |\Sigma_i| + \ln P(c_i) \right)$$

$$g_i(x) = x' W x + w' x + w_{i0}$$

44

General Case Σ_i are arbitrary

$$g_i(x) = \underbrace{x^T W x}_{\text{quadratic in } x \text{ since}} + \underbrace{w^T x}_{\text{linear in } x} + \underbrace{w_{i0}}_{\text{constant in } x}$$

$$x^T W x = \sum_{j=1}^d \sum_{l=1}^d w_{jl} x_j x_l = \sum_{l,j=1}^d w_{lj} x_l x_j$$

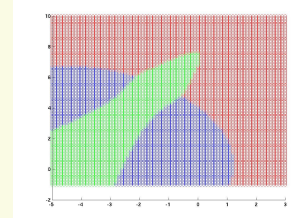
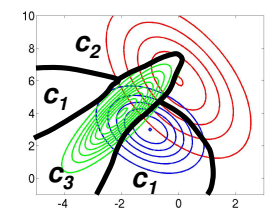
- Thus the discriminant function is quadratic
- Therefore the decision boundaries are quadratic (ellipses and paraboloids)

45

General Case Σ_i are arbitrary: Example

$$\mu_1 = \begin{bmatrix} -1 \\ 3 \end{bmatrix} \quad \mu_2 = \begin{bmatrix} 0 \\ 6 \end{bmatrix} \quad \mu_3 = \begin{bmatrix} -2 \\ 4 \end{bmatrix} \quad \Sigma_1 = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 2 \end{bmatrix} \quad \Sigma_2 = \begin{bmatrix} 2 & -2 \\ -2 & 7 \end{bmatrix} \quad \Sigma_3 = \begin{bmatrix} 1 & 1.5 \\ 1.5 & 3 \end{bmatrix}$$

$$P(c_1) = P(c_2) = \frac{1}{4} \quad P(c_3) = \frac{1}{2}$$



General Case Σ_i are arbitrary: Example

- 3 classes, each 2-dimensional Gaussian with

$$\mu_1 = \begin{bmatrix} -1 \\ 3 \end{bmatrix} \quad \mu_2 = \begin{bmatrix} 0 \\ 6 \end{bmatrix} \quad \mu_3 = \begin{bmatrix} -2 \\ 4 \end{bmatrix}$$

$$\Sigma_1 = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 2 \end{bmatrix} \quad \Sigma_2 = \begin{bmatrix} 2 & -2 \\ -2 & 7 \end{bmatrix} \quad \Sigma_3 = \begin{bmatrix} 1 & 1.5 \\ 1.5 & 3 \end{bmatrix}$$
- Priors: $P(c_1) = P(c_2) = \frac{1}{4}$ and $P(c_3) = \frac{1}{2}$
- Again can be done by solving $g_i(x) = g_j(x)$ for $i, j=1, 2, 3$

$$g_i(x) = x^T \left(-\frac{1}{2} \Sigma_i^{-1} \right) x + \mu_i^T \Sigma_i^{-1} x + \left(-\frac{1}{2} \mu_i^T \Sigma_i^{-1} \mu_i - \frac{1}{2} \ln |\Sigma_i| + \ln P(c_i) \right)$$
- Need to solve a bunch of quadratic inequalities of 2 variables

Important Points

- The Bayes classifier when classes are normally distributed is in general quadratic
 - If covariance matrices are equal and proportional to identity matrix, the Bayes classifier is linear
 - If, in addition the priors on classes are equal, the Bayes classifier is the minimum Euclidean distance classifier
 - If covariance matrices are equal, the Bayes classifier is linear
 - If, in addition the priors on classes are equal, the Bayes classifier is the minimum Mahalanobis distance classifier
- Popular classifiers (Euclidean and Mahalanobis distance) are optimal only if distribution of data is appropriate (normal)