

Object Recognition

Javier González Jiménez

Reference Books:

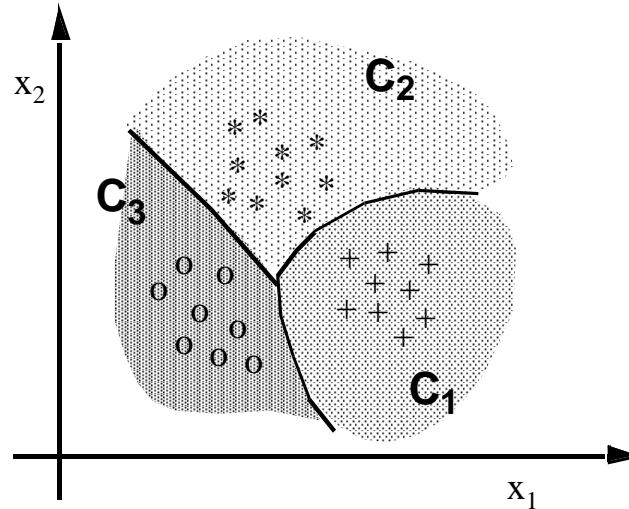
- *Computer Vision: Algorithms and Applications*. Richard Szeliski. Springer. 2010.
<http://szeliski.org/Book>
- *Pattern Recognition and Machine Learning*. Christopher Bishop. Springer-Verlag New York,.2006.

Content

- Introduction
- Discriminant functions
 - Linear discriminant functions
 - Linear Basis Function Models
- Naive Bayes Classifiers
 - Basis
 - Binomial distribution
 - Gaussian distribution
- Support Vector Machine (not included)

1. Introduction

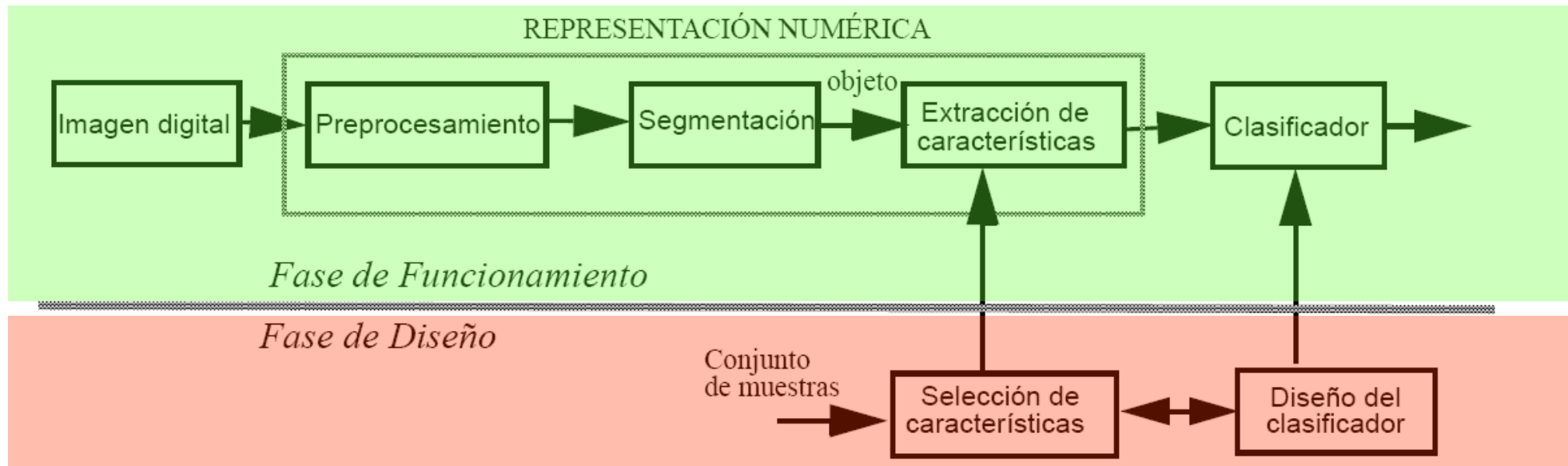
- In previous lecture we learned how to extract a feature vector \mathbf{x} to describe an image region
- Now, we want to **classify the region**, based on that vector \mathbf{x} , as one of M possible object classes (or categories)
- For that, we **divide the feature space** into a number of prediction subspaces C_i : if a feature \mathbf{x} lies in C_i it is assigned to the object class C_i represents



1. Introduction

The recognition process comprises two steps:

- a **training (design) phase**, where sample vectors of known objects are used to learn the classifier (supervised learning)
- a **prediction (online) phase**, where the image objects are classified to one of the classes based on the learned prediction model



1. Introduction

Approaches:

Statistical classifiers:




- We assume the feature vectors \mathbf{x} of the classes C follow a **statistical distribution**
- The parameters of such distribution need to be **learned from known objects**
- Two statistical models can be considered:
 - **Generative models**: the parameters of the joint $p(C, \mathbf{x}) = p(\mathbf{x}|C)p(C)$ are learned
Example: *Naïve Bayesian Classifier*
 - **Discriminative models**: the parameter of the posterior $p(C|\mathbf{x})$ are learned
Example: *Logistic regression, conditional random fields*

Non-statistical classifiers:

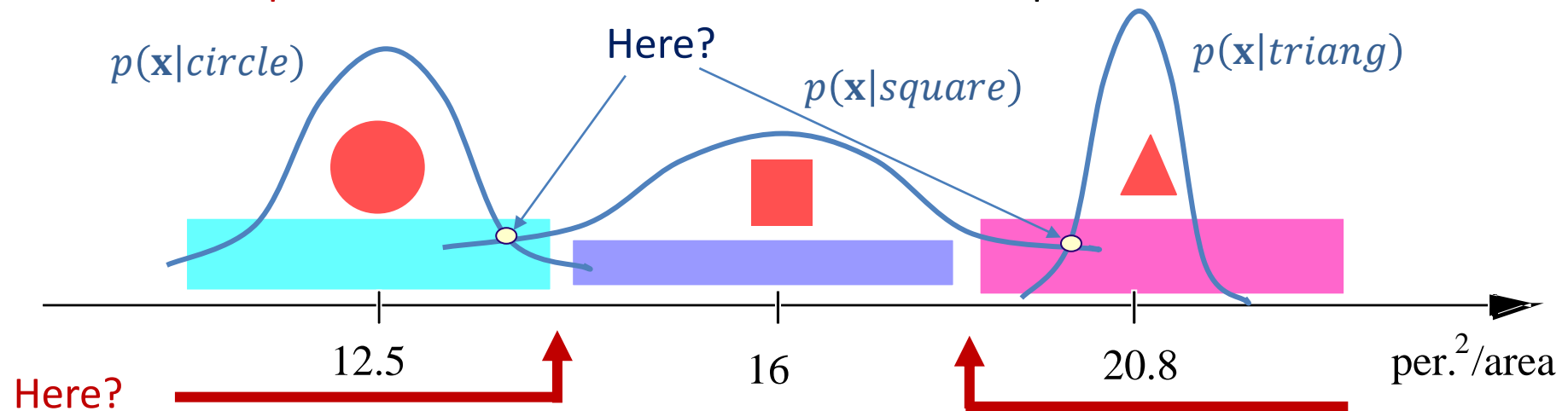
- No assumption is made on the statistical distribution of the feature vector
- The coefficient of deterministic discriminant functions are learned
Example: *Support Vector Machine, Perceptron, AdaBoost*

Example:

Classification of objects, based on compactness (\mathbf{x}), into 3 classes (C)

		AREA	PERIMETER	$\mathbf{x} = \text{PER.}^2/\text{AREA}$ (COMPACTNESS ⁻¹)
circle		πR^2	$2\pi R$	12.5
square		L^2	$4L$	16
equilateral triangle		$(\sqrt{3}/4)L^2$	$3L$	20.8

Partition of the \mathbf{x} space: Where to set the threshold for an optimal classification?



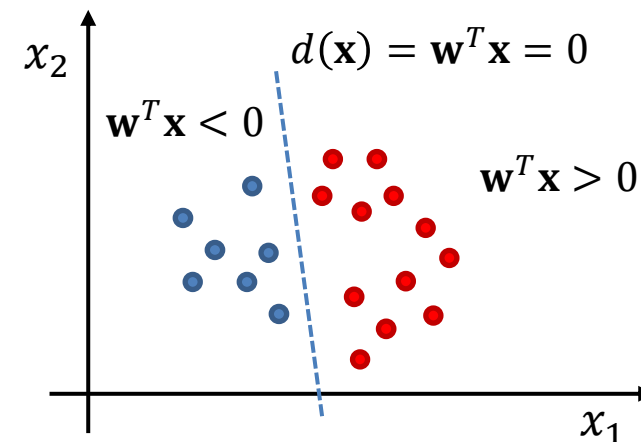
2. Discriminant functions

Linear discriminant functions

For two classes C_1, C_2 :

Hyper-planes in the n-dimensional space

$$d(\mathbf{x}) = w_1 \cdot x_1 + w_2 \cdot x_2 + \dots + w_n \cdot x_n + w_{n+1} = \mathbf{w}_0^T \mathbf{x} + w_{n+1}$$



$\mathbf{w} = [w_1 \quad w_2 \quad \dots \quad w_n]^T$ **Weight Vector** \rightarrow To be learned in the training phase

$\mathbf{x} = [x_1 \quad x_2 \quad \dots \quad x_n]^T$ **Feature Vector**

More convenient as an **augmented form**:

$$\left. \begin{array}{l} \mathbf{w} = [w_1 \quad w_2 \quad \dots \quad w_n \quad w_{n+1}]^T \\ \mathbf{x} = [x_1 \quad x_2 \quad \dots \quad x_n \quad \underbrace{1}_{\text{augmentation}}]^T \end{array} \right\} \Rightarrow d(\mathbf{x}) = \mathbf{w}^T \mathbf{x} \quad \text{Dot (scalar) product of the two vectors}$$

2. Discriminant functions

Linear discriminant functions

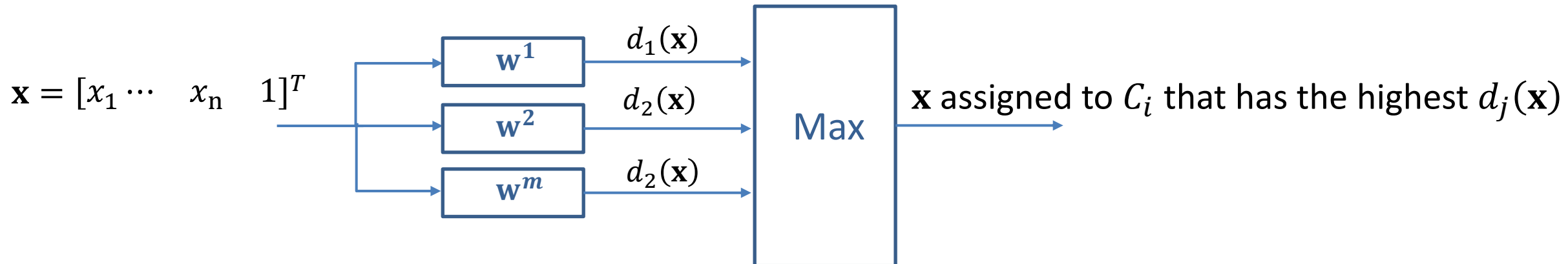
More than two classes (m classes): A linear function for each class

$$d_i(\mathbf{x}) = w_1^i \cdot x_1 + w_2^i \cdot x_2 + \dots + w_n^i \cdot x_n + w_{n+1}^i = (\mathbf{w}^i)^T \cdot \mathbf{x}$$

Classification criterion:

$$\text{IF } d_i(\mathbf{x}) > d_j(\mathbf{x}) \quad \forall i \neq j \quad \text{THEN } \mathbf{x} \in C_i$$

Weights to be
computed during the
training phase

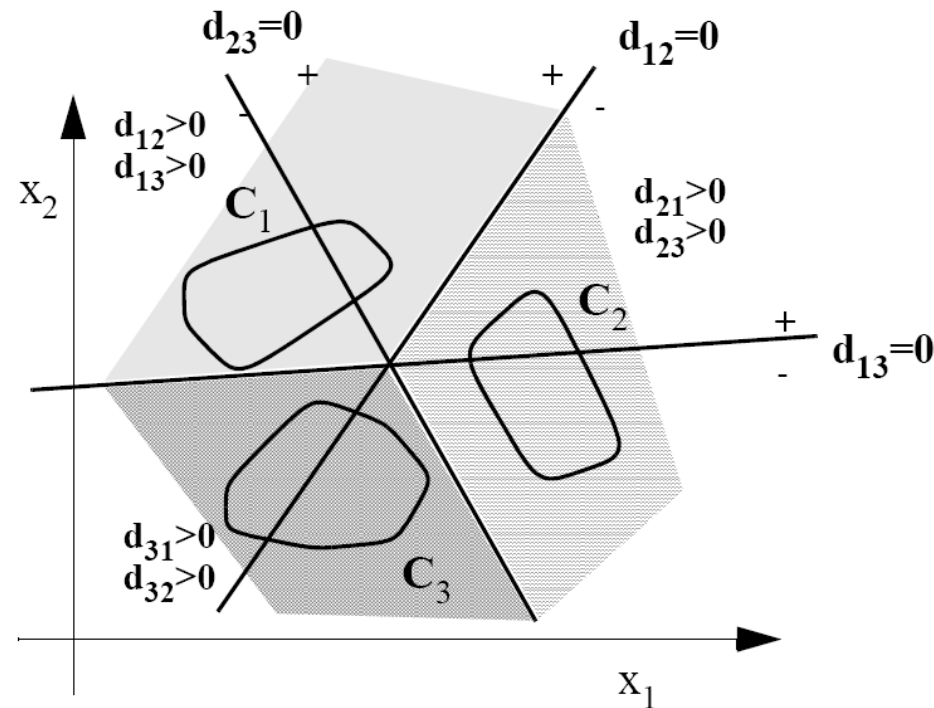


Linear discriminant functions

More than two classes (m classes):

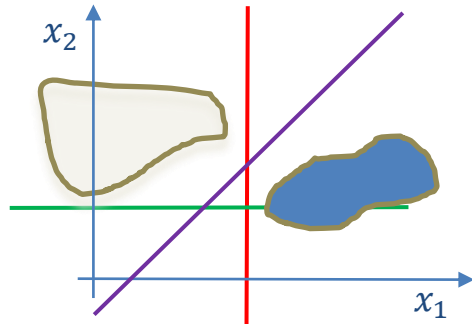
The **border (linear) function** separating the classes C_i y C_j is computed as:

$$d_{ij}(\mathbf{x}) = d_i(\mathbf{x}) - d_j(\mathbf{x}) = (\mathbf{w}_i^T - \mathbf{w}_j^T) \cdot \mathbf{x} = \mathbf{w}_{ij}^T \cdot \mathbf{x} \quad \left\{ \begin{array}{ll} > 0 & d_i(\mathbf{x}) > d_j(\mathbf{x}), \mathbf{x} \notin C_j \\ = 0 & \text{the frontier} \\ < 0 & d_j(\mathbf{x}) > d_i(\mathbf{x}), \mathbf{x} \notin C_i \end{array} \right.$$

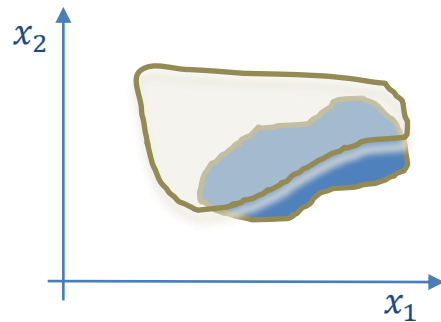


2. Discriminant functions

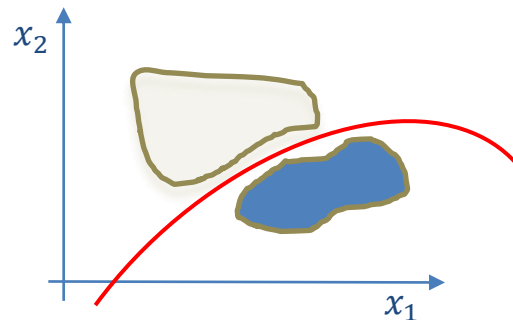
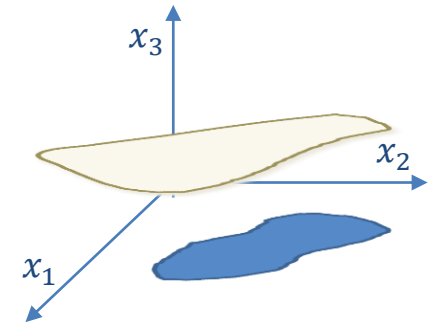
SEPARABILITY:



- Classes separable just with x_1 —
- Classes DO NOT separable with x_2 —
- x_1 more discriminative than x_2
- In the x_1 - x_2 space, classes are more clearly separables (x_2 helps) —



- Classes are not separable in x_1 - x_2
- The feature x_1 provides nothing! (No discriminative at all)
- Can be separated with an additional feature x_3

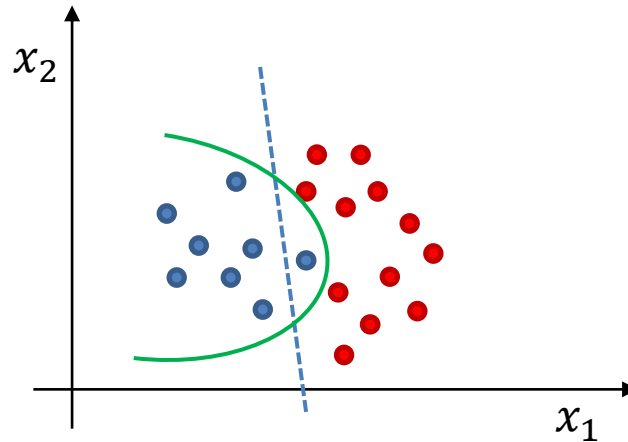


- Classes NOT-linearly separables in x_1 - x_2
- Can be separated with a **Not-linear** discriminant function

2. Discriminant functions

Linear Basis Function Models (also called *generalized discriminant functions*)

Needed when the classes are better separated with not non-linear functions:



How: Through a transformation from the \mathbf{x} space to a $\mathbf{x}' = f(\mathbf{x})$ space ($\dim(\mathbf{x}) < \dim(\mathbf{x}')$):


$$\begin{aligned}
 d(\mathbf{x}) &= w_1 \cdot f_1(\mathbf{x}) + w_2 \cdot f_2(\mathbf{x}) + \dots + w_k \cdot f_k(\mathbf{x}) + w_{k+1} = \sum_{i=1}^{k+1} w_i \cdot \boxed{f_i(\mathbf{x})} = \mathbf{w}^T \cdot \mathbf{f}(\mathbf{x}) \\
 f_{k+1}(\mathbf{x}) &= 1 \\
 \text{New space} \\
 &= w_1 \cdot x'_1 + w_2 \cdot x'_2 + \dots + w_k \cdot x'_k + w_{k+1} = \sum_{i=1}^{k+1} w_i \cdot x'_i = \mathbf{w}^T \cdot \boxed{\mathbf{x}'} = d'(\mathbf{x}')
 \end{aligned}$$

Arrows indicate the mapping from $f_i(\mathbf{x})$ to x'_i for $i=1, 2, k$.

2. Discriminant functions

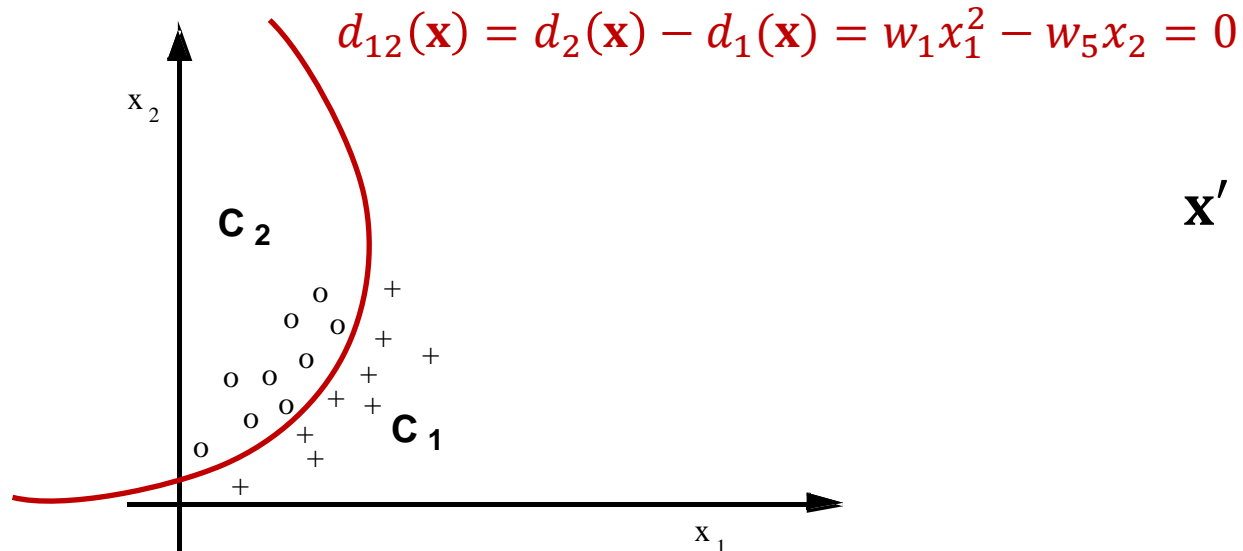
Linear Basis Function Models

$$\left. \begin{aligned} \mathbf{x}' = \mathbf{f}(\mathbf{x}) &= [f_1(\mathbf{x}) \quad f_2(\mathbf{x}) \quad \cdots \quad f_k(\mathbf{x}) \quad 1]^T \\ \mathbf{w}^i &= [w_1^i \quad w_2^i \quad \cdots \quad w_k^i \quad w_{k+1}^i]^T \end{aligned} \right\} d_i(\mathbf{x}) = \mathbf{w}^{iT} \mathbf{x}' = \mathbf{w}^{iT} \mathbf{f}(\mathbf{x})$$




Basis functions

EXAMPLE : Quadratic function ($n=2, k=5$)




$$\mathbf{x}' = [x_1^2 \quad x_1 x_2 \quad x_2^2 \quad x_1 \quad x_2 \quad 1]^T$$


$w_2, w_3, w_4 = 0$




$f_1(\mathbf{x}) = x'_1$




$f_2(\mathbf{x}) = x'_2$



$f_3(\mathbf{x}) = x'_3$



$f_2(\mathbf{x}) = x'_2$



$f_5(\mathbf{x}) = x'_5$

RECALL: Linear and quadratic functions

- **Linear function:**

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = (\mathbf{w}^T \mathbf{x})^T = \mathbf{x}^T \mathbf{w} \quad (\text{Dot product of vectors})$$

Special case: *Square Euclidean distance* (o square 2-norm)

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{x} = \sum_i x_i^2 = \|\mathbf{x}\|_2^2$$

- **Quadratic function:**

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{Q} \mathbf{x} = \sum_i x_i x_j q_{ij} \quad \text{with } \mathbf{Q} = [q_{ij}]_{n \times n}$$

Special cases:

- \mathbf{Q} symmetric ($\mathbf{Q}^T = \mathbf{Q}$) and positive-semidefinite ($\mathbf{x}^T \mathbf{Q} \mathbf{x} \geq 0, \forall \mathbf{x}$)

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{Q} \mathbf{x} = 0 \quad (\text{equation of an ellipse})$$

- Besides, if \mathbf{Q} diagonal ($q_{ij} = 0$ if $i \neq j$)

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{Q} \mathbf{x} = \sum_i x_i^2 q_{ii} = 0 \quad (\text{equation of an ellipse aligned with the axes } x_i)$$

3. Bayesian Classifier

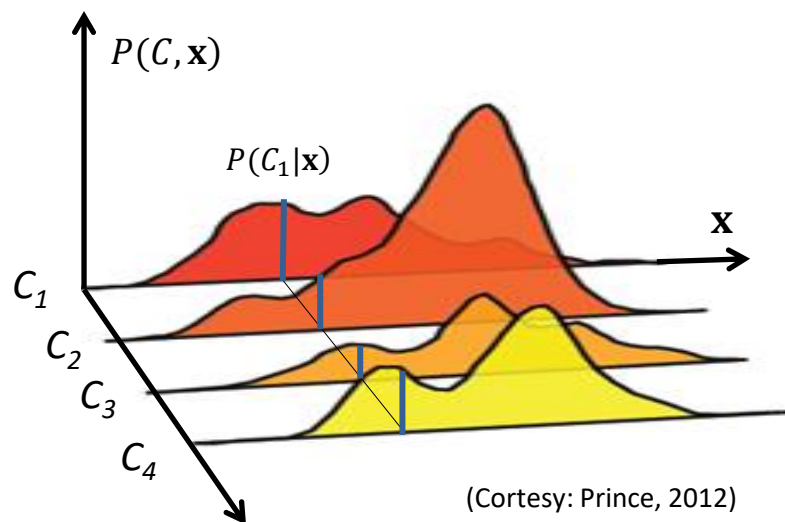
Principle: assign \mathbf{x} to the class C_i that has the highest **posterior probability**

Why: the more probable the less probability of making a mistake

Example: 4 classes, \mathbf{x} belongs to class C_2 . Making an error means assigning it to C_1 or C_3

$$P(\text{error}|\mathbf{x}) = P(C_1|\mathbf{x}) + P(C_3|\mathbf{x}) + P(C_4|\mathbf{x}) = 1 - P(C_2|\mathbf{x}) \quad \text{Recall: } \sum_{k=1}^{M \text{ classes}} P(C_k|\mathbf{x}) = 1$$

Minimize $P(\text{error}|\mathbf{x}) \rightarrow$ assign \mathbf{x} to the class with the highest probability

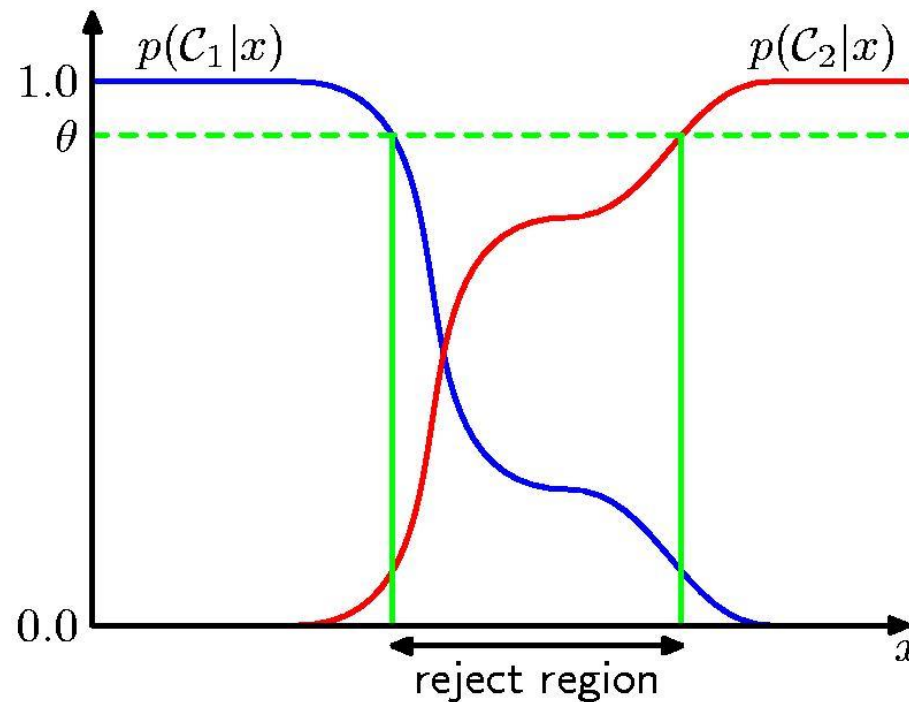


In this example, given a feature vector $\mathbf{x} = \mathbf{x}_1$, it must be assigned to C_1 , since $P(C_1|\mathbf{x})$ is the highest value.

This is called: a **MAP** prediction
(**M**aximum **A** Posteriori)

3. Bayesian Classifier

Sometimes, it is convenient to have a *reject region*, where **no decision is made**

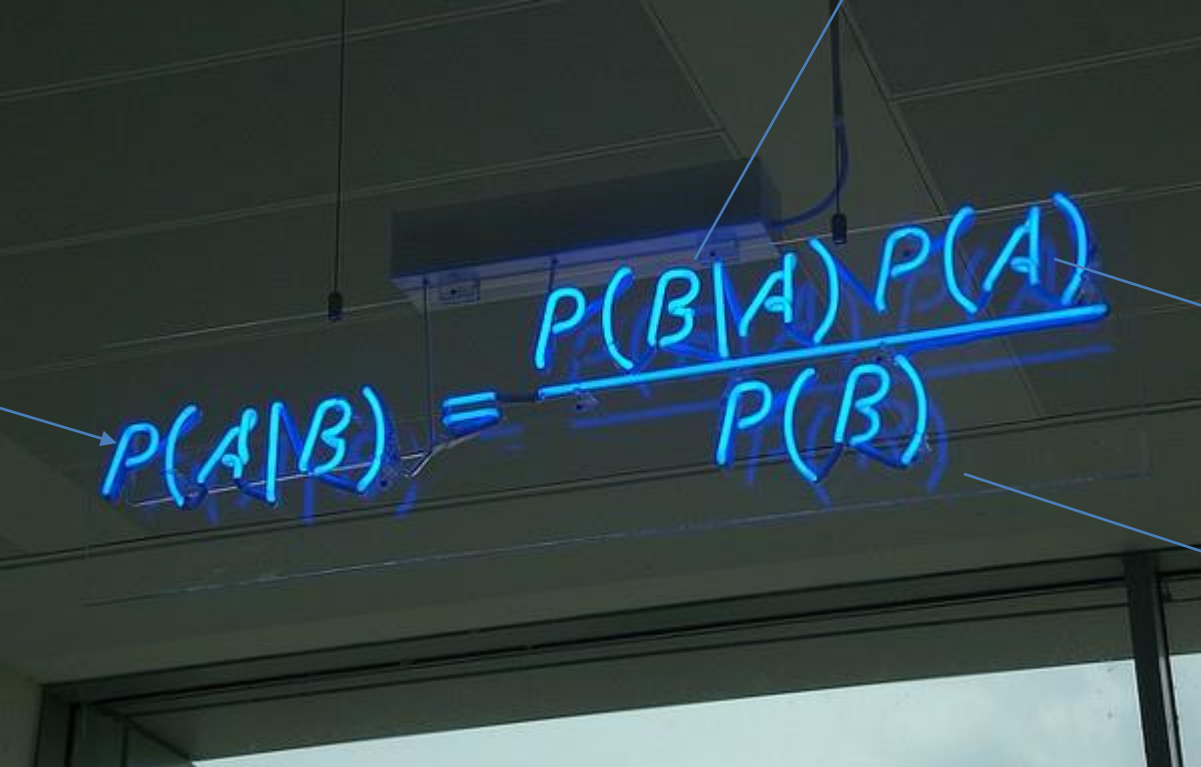


(Courtesy of Bishop, 2006)

Reject region: None probability is high enough (above a certain posterior probability θ)

3. Bayesian Classifier

RECALL: Bayes theorem



Likelihood (function of A)

A posteriori
(B given)

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

A priori

Total probability

3. Bayesian Classifier

We can build a classifier based on discriminant functions for a Bayesian classifier


How: Create a discriminant function $d_i(\mathbf{x})$ for each class C_i , such that $d_i(\mathbf{x}) > d_j(\mathbf{x})$ whenever $P(C_i|\mathbf{x}) > P(C_j|\mathbf{x})$

Then:


$P(\mathbf{x})$ is a scale factor
for all the classes C_i

$\text{Max } \ln(f(\mathbf{x})) \rightarrow \text{Max } f(\mathbf{x})$

If $P(C_i) = P(C_j) \quad \forall i, j$


$$d_i(\mathbf{x}) = P(C_i|\mathbf{x}) = \frac{p(\mathbf{x}/C_i)P(C_i)}{P(\mathbf{x})}$$

MAP estimation
(Maximum **A** Posteriori)


$$d_i(\mathbf{x}) = p(\mathbf{x}/C_i)P(C_i)$$


$$d_i(\mathbf{x}) = \ln p(\mathbf{x}/C_i) + \ln P(C_i)$$

$$d_i(\mathbf{x}) = \ln p(\mathbf{x}/C_i)$$

Maximum Log-Likelihood estimation

3. Bayesian Classifier

BINOMIAL DISTRIBUTION:

$\mathbf{x} = [x_1 \quad x_2 \quad \cdots \quad x_f \quad \cdots \quad x_n]^T$ Features are either **0** or **1** (Bernoulli trial)

$$p(x_f/C_i) = p_f^i x_f \cdot (1 - p_f^i)^{(1-x_f)}$$

p_f^i probability that $x_f = 1$ if $\mathbf{x} \in C_i$
 $1 - p_f^i$ probability that $x_f = 0$ if $\mathbf{x} \in C_i$

If features independent
(Naïve Bayes Classifier)

$$p(\mathbf{x}/C_i) = \prod_{f=1}^n p(x_f/C_i) = \prod_{f=1}^n p_f^i x_f \cdot (1 - p_f^i)^{(1-x_f)}$$

$$\ln p(\mathbf{x}/C_i) = \ln \prod_{f=1}^n p_f^i x_f \cdot (1 - p_f^i)^{(1-x_f)} = \sum_{f=1}^n [x_f \cdot \ln p_f^i + (1 - x_f) \cdot \ln(1 - p_f^i)] = \sum_{f=1}^n x_f \cdot \ln \frac{p_f^i}{1 - p_f^i} + \sum_{f=1}^n \ln(1 - p_f^i)$$

$$d_i(\mathbf{x}) = \ln P(C_i) + \ln p(\mathbf{x}/C_i) = \underbrace{\ln P(C_i) + \sum_{f=1}^n \ln(1 - p_f^i)}_{w_{n+1}^i} + \sum_{f=1}^n x_f \cdot \underbrace{\ln \frac{p_f^i}{1 - p_f^i}}_{w_f^i} = w_{n+1}^i + \sum_{f=1}^n w_f^i \cdot x_f$$

Linear function!

WATCH OUT: we can not use p_f^i exactly equals to 0 or 1 because we would have numerical problems. Instead, take values close to 0 or 1.

3. Bayesian Classifier

EXAMPLE BINOMIAL DISTRIBUTION . Number recognition

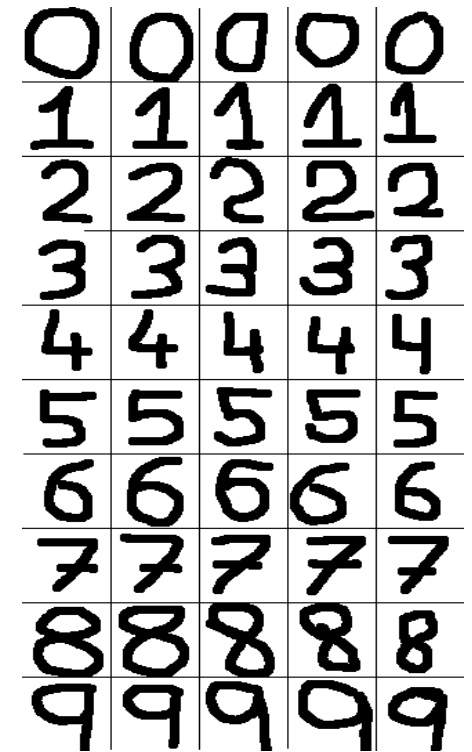
Learning the classifier: Estimate the probabilities p_f^i

Input: many binary images with the numbers handwritten

- The numbers are segmented and binarized.
- The bounding box around the segmented number is resampled to have a fixed size of 16 x 16 pixels
- The bounding box image is rearranged in a vector of 16x16= 256 elements

Output: 10 discriminant functions

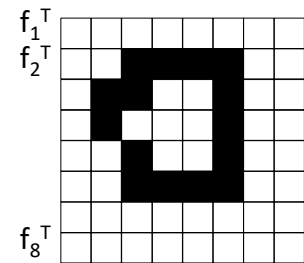
[More detail next]



EXAMPLE BINOMIAL DISTRIBUTION . Number recognition

Learning the classifier

Reshape the segmented image to be a vector:

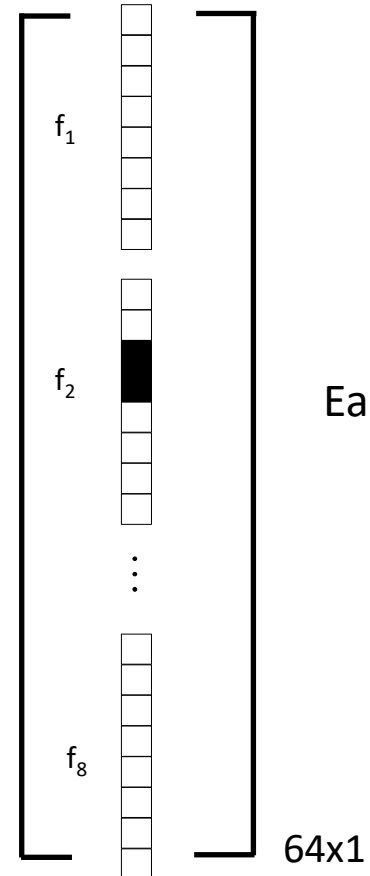


8x8

segmented image

→ $X =$

Feature vector

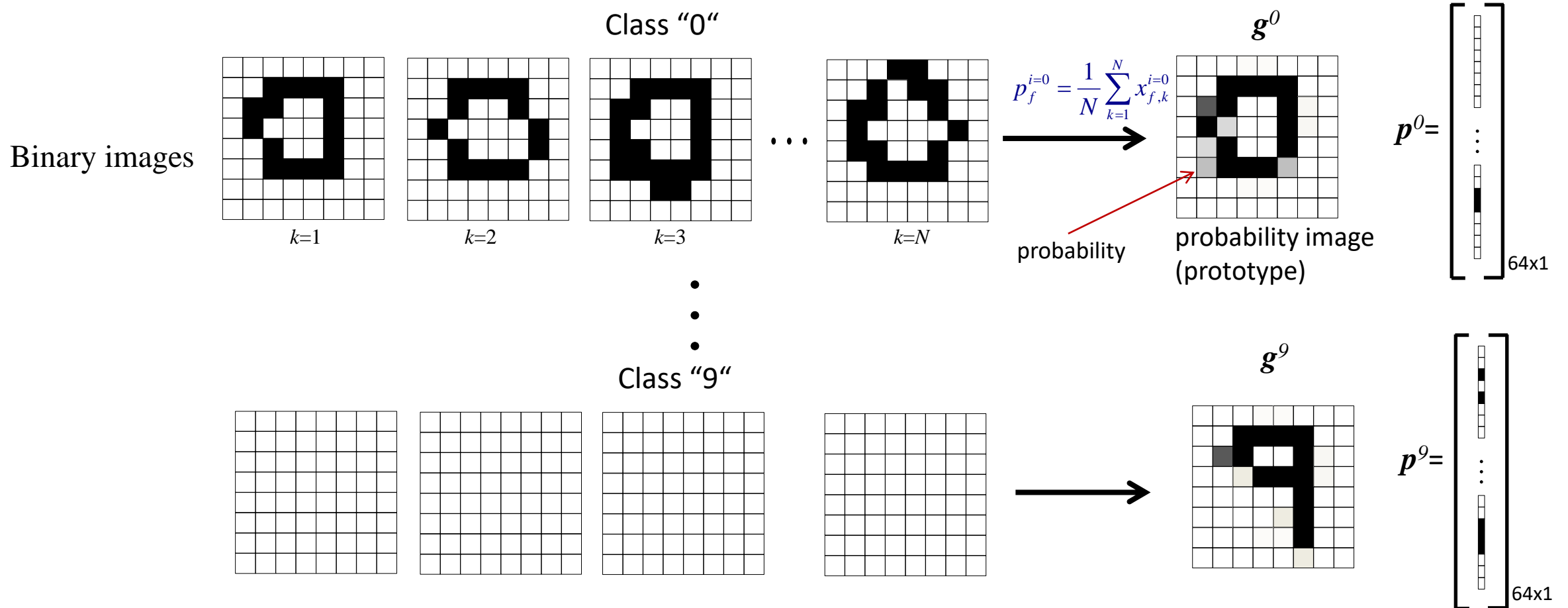


Each element equals 0 or 1

EXAMPLE BINOMIAL DISTRIBUTION . Number recognition

Learning the classifier

Computation of the p_f for each class $i=0, 9$



EXAMPLE BINOMIAL DISTRIBUTION . Number recognition

Learning the classifier

Prototype images \mathbf{g}^i : The value of a pixel $g^i(i, j)$ represents the probability of pixel (i, j) of being “1” (black on a white paper sheet)

$$g^i(i, j) = p_f = \frac{1}{N} \sum_{k=1}^N x_{f,k}^i$$

N is the number of training samples
(given to the recognizer)



EXAMPLE BINOMIAL DISTRIBUTION . Number recognition

Learning the classifier

Computation of the discriminant functions for each class $i = 0, \dots, 9$

Compute the 256 weights (f) for each class $i = 0 \dots, 9$

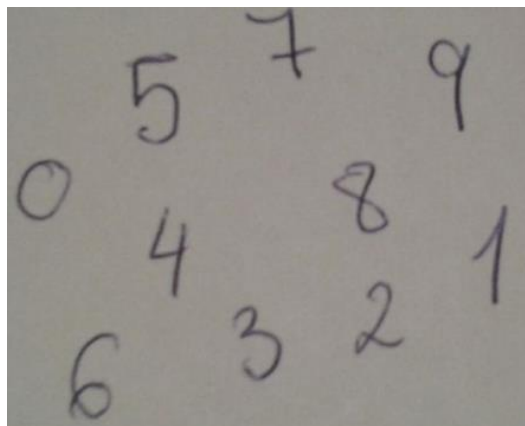
$$= w_f^i \ln \frac{p_f^i}{1 - p_f^i}$$

$$w_{n+1}^i = \ln P(C_k) + \sum_{f=1}^n \ln(1 - p_f^i)$$

EXAMPLE BINOMIAL DISTRIBUTION . Number recognition

Prediction phase:

Input: Image of handwritten numbers (not rotated)



- The numbers are segmented and binarized.
- The bounding box around the segmented number is resampled to have a fixed size of 16 x 16 pixels
- The bounding box image is rearranged in a vector of $f = 16 \times 16 = 256$ elements

EXAMPLE BINOMIAL DISTRIBUTION . Number recognition

Prediction phase:

Evaluate the 10 discriminant functions $d_i(\mathbf{x})$ ($i = 0, \dots, 9$) for each vector \mathbf{x}

$$d_i(\mathbf{x}) = w_{n+1}^i + \sum_{f=1}^{256} w_f^i \cdot x_f$$

Output:

- class C_i with the highest $d_i(\mathbf{x})$
- Probability of \mathbf{x} to belong to any class C_i

$$p(C_i/\mathbf{x}) = \eta p(\mathbf{x}/C_i)P(C_i)$$

$$\text{with } \eta \sum_{i=0}^9 p(\mathbf{x}/C_i)P(C_i) = 1 \quad \Rightarrow \quad \eta = 1 / \sum_{i=0}^9 p(\mathbf{x}/C_i)P(C_i)$$

3. Bayesian Classifier

GAUSSIAN DISTRIBUTION:

Feature vector of dimension n : $\mathbf{x} = [x_1 \quad x_2 \quad \cdots \quad x_f \quad \cdots \quad x_n]^T$

Features can take continuous values following the probability density function (pdf):

$$p(\mathbf{x}/C_i) = \frac{1}{(2\pi)^{n/2} |\Sigma^i|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}^i)^T \Sigma^{i-1} (\mathbf{x}-\boldsymbol{\mu}^i)}$$

Given by two set of parameters:

- **Mean vector:** $\boldsymbol{\mu} = [\mu_1 \quad \mu_2 \quad \cdots \quad \mu_f \quad \cdots \quad \mu_n]^T$

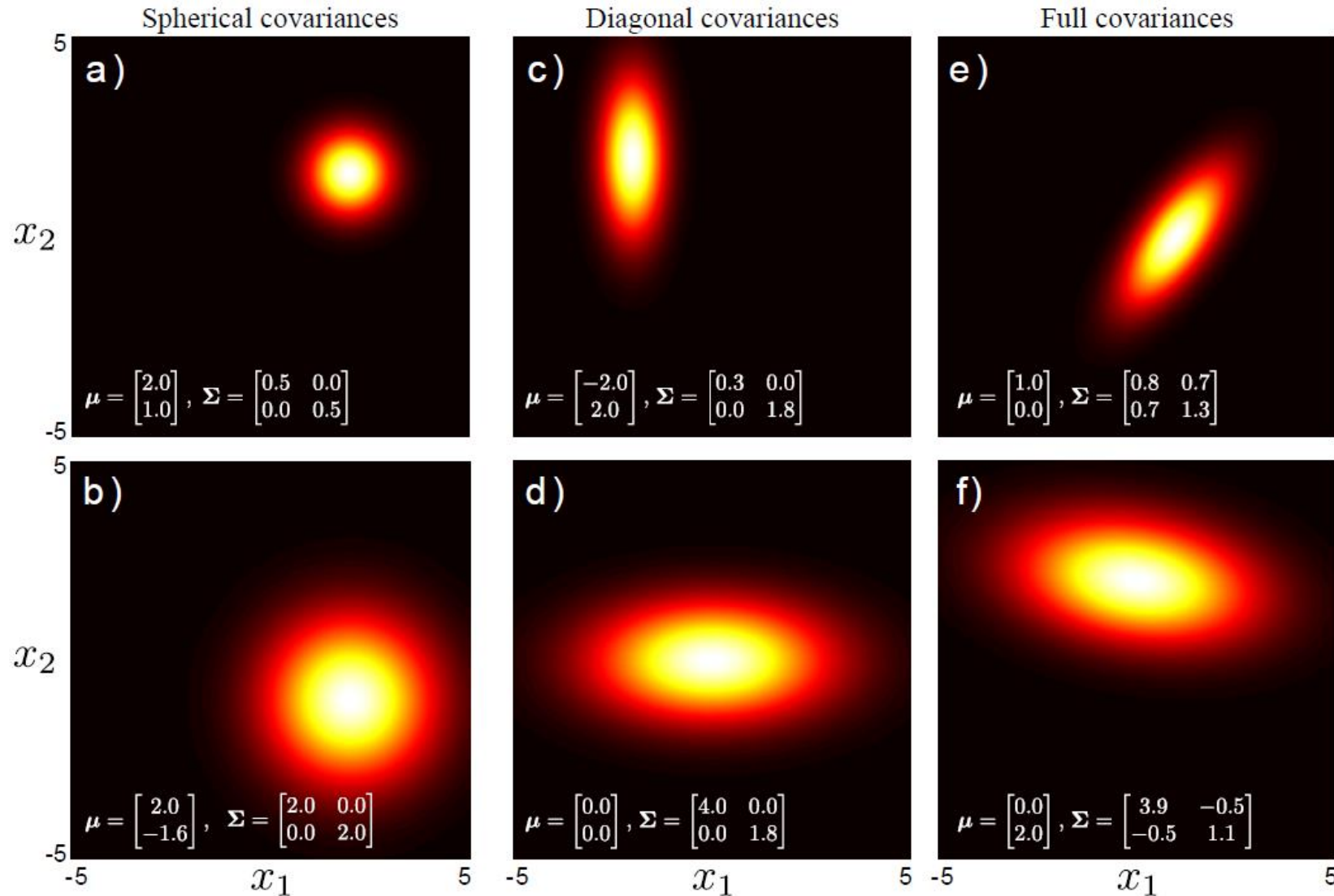
- **Covariance matrix:**

$$\Sigma = E[(\mathbf{x} - \boldsymbol{\mu}) \cdot (\mathbf{x} - \boldsymbol{\mu})^T] = \begin{bmatrix} \sigma_{11} & \cdots & \sigma_{1f} & \cdots & \sigma_{1n} \\ \vdots & & \vdots & & \vdots \\ \sigma_{n1} & \cdots & \sigma_{nf} & \cdots & \sigma_{nn} \end{bmatrix} = \begin{bmatrix} \sigma_{11} & \cdots & 0 & \cdots & 0 \\ \vdots & & \vdots & & \vdots \\ 0 & \cdots & 0 & \cdots & \sigma_{nn} \end{bmatrix}$$

Features are independent (no correlated)
(Naïve Bayes Classifier)

GAUSSIAN DISTRIBUTION:

RECALL: How does a Gaussian depend on its two parameters μ y Σ



(Cortesía: Prince, 2012)

3. Bayesian Classifier

GAUSSIAN DISTRIBUTION:

$$d_k(\mathbf{x}) = \ln P(C_k) + \ln p(\mathbf{x}/C_k) = \ln P(C_k) + \ln \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma^k|^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}^k)^T (\Sigma^k)^{-1} (\mathbf{x}-\boldsymbol{\mu}^k)} =$$

Squared Mahalanobis distance $D_k^2(\mathbf{x})$

$$= \ln P(C_k) + \ln \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma^k|^{\frac{1}{2}}} - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}^k)^T (\Sigma^k)^{-1} (\mathbf{x} - \boldsymbol{\mu}^k)$$

$$= \ln P(C_k) - \frac{1}{2} [n \ln(2\pi) + \ln |\Sigma^k| + D_k^2(\mathbf{x})]$$

Constant → can be removed

$$d_k(\mathbf{x}) = \underbrace{\ln P(C_k) - \frac{1}{2} [\ln |\Sigma^k| + \boldsymbol{\mu}^{kT} (\Sigma^k)^{-1} \boldsymbol{\mu}^k]}_{\text{Independent term}} + \underbrace{\mathbf{x}^T (\Sigma^k)^{-1} \boldsymbol{\mu}^k}_{\text{Linear weights}} - \frac{1}{2} \underbrace{\mathbf{x}^T (\Sigma^k)^{-1} \mathbf{x}}_{\text{Quadratic weights}}$$

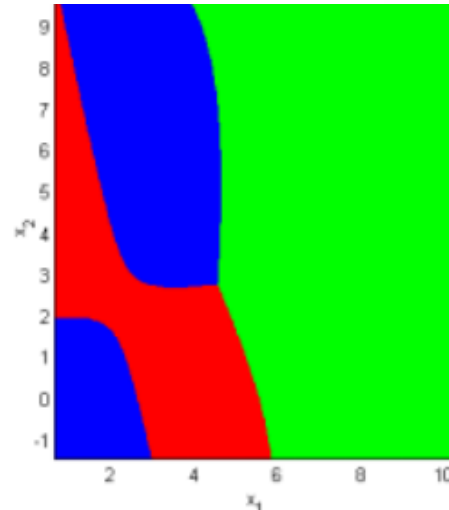
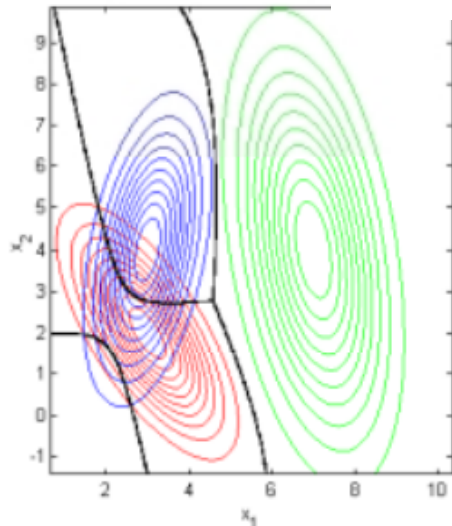
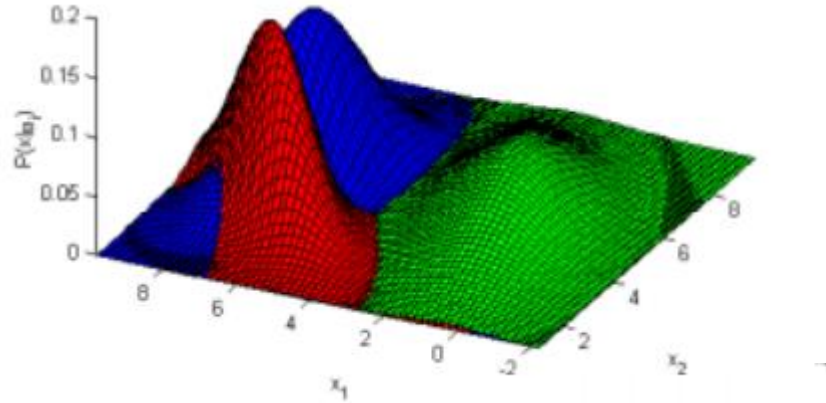
The discriminant function is a quadratic polynomial! $d_k(\mathbf{x}) = w_{n+1} + \mathbf{x}^T \mathbf{w} + \mathbf{x}^T \mathbf{Q} \mathbf{x}$

GAUSSIAN DISTRIBUTION :

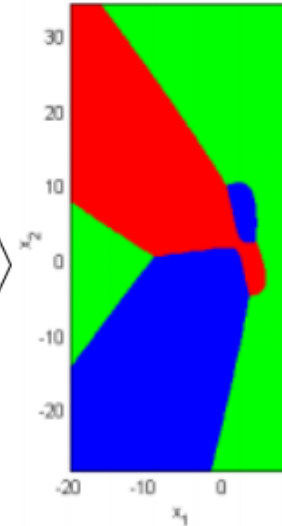
Visually:

$$d_k(\mathbf{x}) = \ln P(C_k) - \frac{1}{2} \ln |\Sigma^k| - \left[(\mathbf{x} - \boldsymbol{\mu}^k)^T (\Sigma^k)^{-1} (\mathbf{x} - \boldsymbol{\mu}^k) \right]$$

$D_k^2(x)$ Squared Mahalanobis distance



Zoom out



NORMAL DISTRIBUTION:

SIMPLIFICATIONS

Equal priors: $P(C_1) = P(C_2) = \dots = P(C_m) = P(C)$

Same covariance matrices: $\Sigma^1 = \Sigma^2 = \dots = \Sigma^m = \Sigma$

$$d_k(\mathbf{x}) = \ln P(C_k) - \frac{1}{2} \left[n \ln 2\pi + \ln |\Sigma^k| + \mathbf{x}^T \Sigma^{k-1} \mathbf{x} - 2\mathbf{x}^T \Sigma^{k-1} \boldsymbol{\mu}^k + \boldsymbol{\mu}^{kT} \Sigma^{k-1} \boldsymbol{\mu}^k \right]$$

$$= \ln P(C) - \frac{1}{2} \left[n \ln 2\pi + \ln |\Sigma| + \mathbf{x}^T \Sigma^{-1} \mathbf{x} - 2\mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}^k + \boldsymbol{\mu}^{kT} \Sigma^{-1} \boldsymbol{\mu}^k \right]$$

Constant → can be removed

$$d_k(\mathbf{x}) = -(\mathbf{x} - \boldsymbol{\mu}^k)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}^k) = -D_k^2(\mathbf{x})$$

Classifier based on Square Mahalanobis distance

Notice: $\mathbf{x}^T \Sigma^{-1} \mathbf{x}$ is a constant term for all the classes → can be removed:

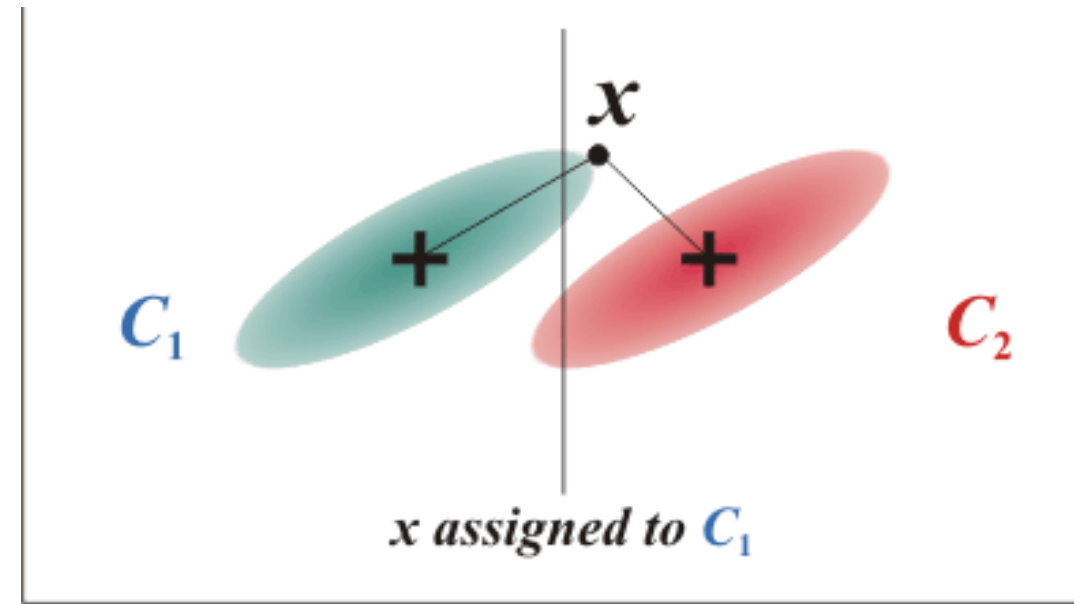
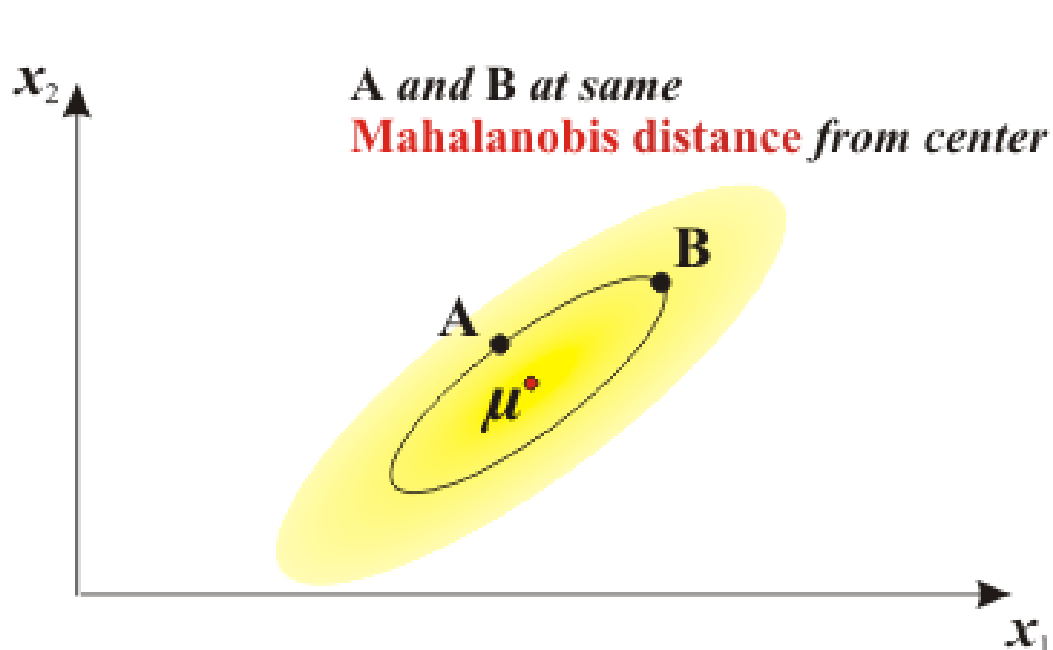
$$d_k(\mathbf{x}) = -2\mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}^k + \boldsymbol{\mu}^{kT} \Sigma^{-1} \boldsymbol{\mu}^k = w_{n+1} + \mathbf{x}^T \mathbf{w} \quad \text{The discriminant functions are Linear!}$$

Classifier based on Square Mahalanobis distance

$$d_k(\mathbf{x}) = -(\mathbf{x} - \boldsymbol{\mu}^k)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}^k) = -D_k^2(\mathbf{x})$$

Recall: Quadratic function:

$$f(\mathbf{x}) = \mathbf{x}^T \cdot \mathbf{Q} \cdot \mathbf{x} = \sum_i x_i x_j q_{ij} \quad \text{with } \mathbf{Q} = [q_{ij}]_{n \times n}$$



x is assigned to C_1 though the closest centroid is C_2

NORMAL DISTRIBUTION:

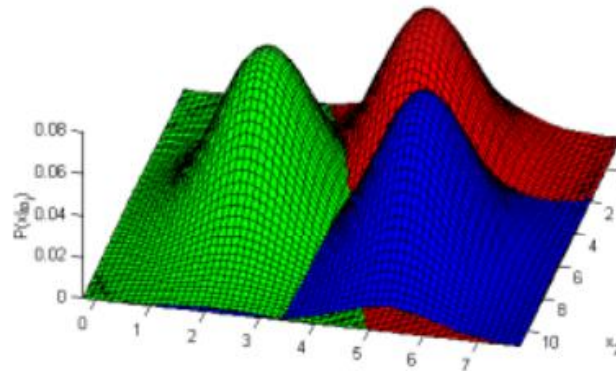
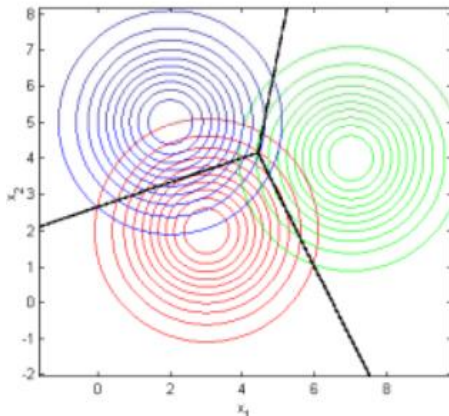
MORE SIMPLIFICATIONS

Isotropic covariance matrix: $\Sigma^k = \Sigma = \sigma^2 \cdot \mathbf{I} = \sigma^2 \begin{bmatrix} 1 & \dots & 0 \\ \vdots & \dots & \vdots \\ 0 & \dots & 1 \end{bmatrix}$

$$d_k(\mathbf{x}) = -(\mathbf{x} - \boldsymbol{\mu}^k)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}^k) = -\frac{1}{\sigma^2} (\mathbf{x} - \boldsymbol{\mu}^k)^T \mathbf{I} (\mathbf{x} - \boldsymbol{\mu}^k) = -\underbrace{(\mathbf{x} - \boldsymbol{\mu}^k)^T (\mathbf{x} - \boldsymbol{\mu}^k)}_{-\|\mathbf{x} - \boldsymbol{\mu}^k\|^2} = -D_k^2(\mathbf{x})$$

Linear classifier!

EUCLIDEAN distance ←

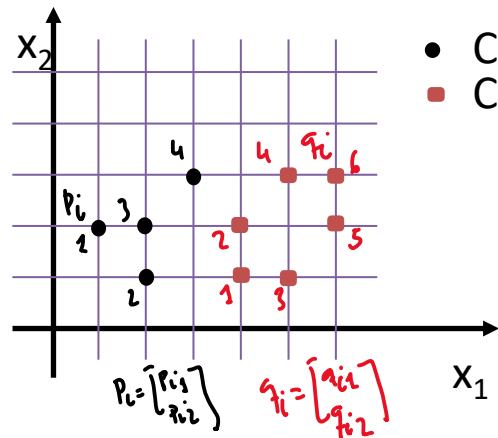


DECISION RULE:

assign \mathbf{x} to the class C_k such that:
 $d_k(\mathbf{x}) > d_j(\mathbf{x}) \quad \forall k \neq j$

It's called the NATURAL classifier

EXAMPLE: Design a classifier based on minimum Mahalanobis distance for the following example.



- Class 1
- Class 2

- Draw approximately the two ellipses representing the Covariance matrices
- What class do $p_1 = [3, 2]^T$ and $p_2 = [4, 3]^T$ belong to?

Classifier based on Square Mahalanobis distance $d_k(\mathbf{x}) = -(\mathbf{x} - \boldsymbol{\mu}^k)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}^k) = -D_k^2(\mathbf{x})$
 we need the two means μ^k and a single Σ for the two set of data points.

$$\mu^1 = \frac{1}{4} \sum_i p_i = \begin{bmatrix} 2 \\ 2 \end{bmatrix} \quad \mu^2 = \frac{1}{6} \sum_i q_i = \begin{bmatrix} 3 \\ 3 \end{bmatrix}$$

$$\Sigma = \frac{1}{4} \sum (p_i - \mu^1)(p_i - \mu^1)^T + \frac{1}{6} \sum (q_i - \mu^2)(q_i - \mu^2)^T$$

$$\Sigma = \frac{1}{4} \left\{ \begin{pmatrix} -1 \\ 0 \end{pmatrix} \begin{pmatrix} -1 & 0 \end{pmatrix} + \begin{pmatrix} 0 \\ -1 \end{pmatrix} \begin{pmatrix} 0 & -1 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \end{pmatrix} \begin{pmatrix} 0 & 0 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \end{pmatrix} \right\} + \frac{1}{6} \left\{ \begin{pmatrix} -1 \\ -1 \end{pmatrix} \begin{pmatrix} -1 & -1 \end{pmatrix} + \begin{pmatrix} -1 \\ 0 \end{pmatrix} \begin{pmatrix} -1 & 0 \end{pmatrix} + \begin{pmatrix} 0 \\ -1 \end{pmatrix} \begin{pmatrix} 0 & -1 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \end{pmatrix} \begin{pmatrix} 0 & 0 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 1 & 0 \end{pmatrix} + \begin{pmatrix} 0 \\ 1 \end{pmatrix} \begin{pmatrix} 0 & 1 \end{pmatrix} \right\}$$

$$\Sigma = \frac{1}{4} \begin{bmatrix} 2 & 2 \\ 1 & 2 \end{bmatrix} + \frac{1}{3} \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} = \frac{1}{12} \begin{bmatrix} 6 & 3 \\ 3 & 6 \end{bmatrix} + \frac{1}{3} \begin{bmatrix} 4 & 4 \\ 4 & 8 \end{bmatrix} = \frac{1}{12} \begin{bmatrix} 14 & 7 \\ 7 & 14 \end{bmatrix} = \frac{7}{12} \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \rightarrow \Sigma^{-1} = \frac{4}{7} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$$

$$\Sigma \Sigma^{-1} = \frac{7 \cdot 4}{12 \cdot 7} \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \text{check}$$

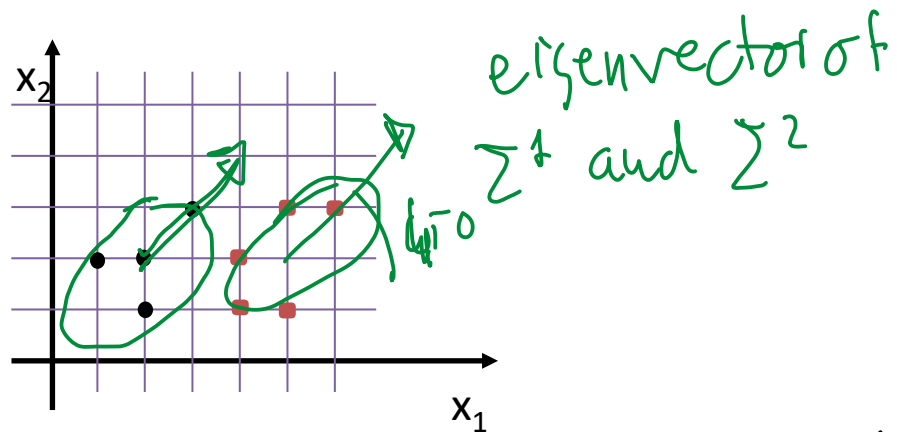
Both gaussians have the same orientation (45°) why? lets compute the eigenvectors

$$|\Sigma - \lambda I| = 0 \rightarrow \begin{vmatrix} 2-\lambda & 1 \\ 1 & 2-\lambda \end{vmatrix} = 0 \rightarrow (2-\lambda)^2 - 1 = 0$$

$$\lambda^2 - 4\lambda + 4 - 1 = 0 \quad \lambda = \frac{4 \pm \sqrt{16-4}}{2} = 2 \pm \sqrt{3}$$

$$\lambda_1 = 3 \rightarrow (\Sigma - \lambda_1 I)u = 0 \quad \begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} u_x \\ u_y \end{bmatrix} = 0 \quad u = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\lambda_2 = 1 \rightarrow \theta = \arctan \frac{1}{1} = 45^\circ$$



Recall: $(x_1, x_2) \begin{pmatrix} a & b \\ b & c \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = ax_1^2 + bx_2^2 + 2bx_1x_2$

Assuming the same Σ for both classes

$$\Sigma^{-1} = \frac{1}{7} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$$

$$d^1(x) = - \begin{bmatrix} (x_1 - 2) & (x_2 - 2) \end{bmatrix} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 - 2 \\ x_2 - 2 \end{bmatrix} = - \left[2(x_1 - 2)^2 + 2(x_2 - 2)^2 - 2(x_1 - 2)(x_2 - 2) \right]$$

$$d^2(x) = - \begin{bmatrix} (x_1 - 5) & (x_2 - 2) \end{bmatrix} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} x_1 - 5 \\ x_2 - 2 \end{bmatrix} = - \left[2(x_1 - 5)^2 + 2(x_2 - 2)^2 - 2(x_1 - 5)(x_2 - 2) \right]$$

$$d^1 \left(\begin{pmatrix} 3 \\ 2 \end{pmatrix} \right) = - (2 + 0 - 2(1)(0)) = -2$$

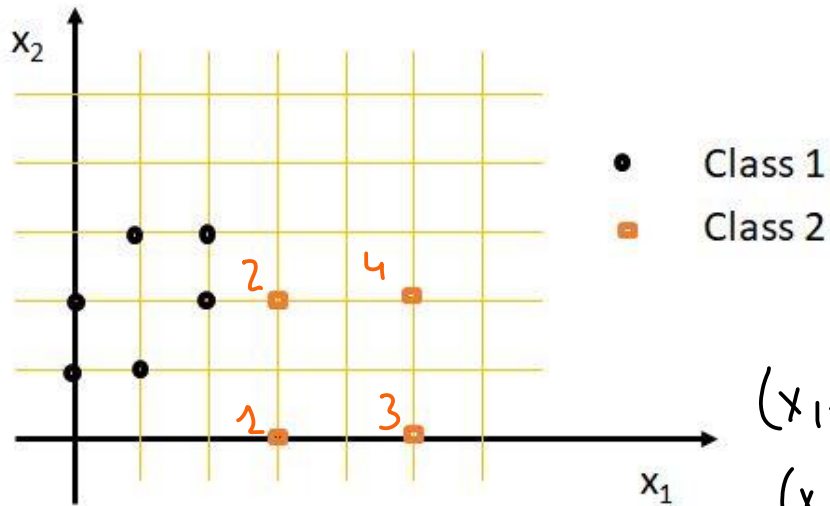
$$d^2 \left(\begin{pmatrix} 3 \\ 2 \end{pmatrix} \right) = - (2 \cdot 4 + 0 + 0) = -8$$

$$d^1 > d^2 \Rightarrow \begin{bmatrix} 3 \\ 2 \end{bmatrix} \in C_1$$

Frontier between classes: $d^1(x) = d^2(x) \Rightarrow 0$ The quadratic terms (x_1^2, x_2^2, x_1x_2) cancel out $\Rightarrow d^1(x) = 0$ at $(4, 0)$

Pregunta examen:

1. Compute the value of the feature vector $\mathbf{x}=[2,0]^T$ for the discriminant function for Class 2 (d_2). Assume $\ln P(C_k) = -0,5$
2. Which is the orientation of the Gaussian of the Class 1 .



$$\mu_2 = \frac{1}{4} \left[\begin{pmatrix} 3 \\ 0 \end{pmatrix} + \begin{pmatrix} 3 \\ 2 \end{pmatrix} + \begin{pmatrix} 5 \\ 0 \end{pmatrix} + \begin{pmatrix} 5 \\ 2 \end{pmatrix} \right] = \begin{bmatrix} 4 \\ 1 \end{bmatrix}$$

$$\Sigma_2 = \frac{1}{N} \sum_i (\mathbf{x}_i - \mu) (\mathbf{x}_i - \mu)^T = \frac{1}{4} \begin{bmatrix} 4 & 0 \\ 0 & 4 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \rightarrow \Sigma^{-1} = \frac{1}{1} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\begin{aligned} (\mathbf{x}_1 - \mu) (\mathbf{x}_1 - \mu)^T &= \begin{pmatrix} -1 \\ -1 \end{pmatrix} \begin{pmatrix} -1 & -1 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} & (\mathbf{x}_3 - \mu) (\mathbf{x}_3 - \mu)^T &= \begin{pmatrix} -1 \\ -1 \end{pmatrix} \begin{pmatrix} -1 & -1 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \\ (\mathbf{x}_2 - \mu) (\mathbf{x}_2 - \mu)^T &= \begin{pmatrix} -1 \\ 1 \end{pmatrix} \begin{pmatrix} -1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} & (\mathbf{x}_4 - \mu) (\mathbf{x}_4 - \mu)^T &= \begin{pmatrix} 1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \end{aligned}$$

$$d_k(\mathbf{x}) = \ln P(C_k) - \frac{1}{2} \left[\ln |\Sigma^k| + \mu^{kT} (\Sigma^k)^{-1} \mu^k \right] + \mathbf{x}^T (\Sigma^k)^{-1} \mu^k - \frac{1}{2} \mathbf{x}^T (\Sigma^k)^{-1} \mathbf{x}$$

$$\mu^T \Sigma^{-1} \mu = (4 \ 1) \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 4 \\ 1 \end{pmatrix} = 16 + 1 = 17$$

$$\mathbf{x}^T \Sigma^{-1} \mu = (x_1 \ x_2) \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 4 \\ 1 \end{pmatrix} = 4x_1 + x_2$$

$$\mathbf{x}^T \Sigma \mathbf{x} = \mathbf{x}^T \mathbf{x} = x_1^2 + x_2^2$$

$$d_2(\mathbf{x}) = -0,5 - \frac{1}{2} \left[0 + 17 \right] + 4x_1 + x_2 + x_1^2 + x_2^2 = x_1^2 + x_2^2 + 4x_1 + x_2 - 9$$

$$d \begin{bmatrix} 2 \\ 0 \end{bmatrix} = 4 + 0 + 8 + 0 - 9 = 3$$