They are identical if $\hat{P}_j = $ constant and $|\hat{\underline{\underline{C}}}_j| = $ constant

Important: $N_j$ (no of training examples for class $\omega_j$) large enough for

- $\hat{\underline{\underline{C}}}_j$ invertible : $N_j \gg d$   else $\hat{\underline{\underline{C}}}_j$ is rank deficient.

   Hence NECESSARY CONDITION

- good estimate for $\hat{\underline{\mu}}_j$, $\hat{\underline{\underline{C}}}_j$ : $N_j \gg d$

## 4.3.3   Naive Bayes plug-in

A simplified version of Bayes plugin

- Naive assumption : all features $x_1 \cdots x_d$ are independent

   i.e. $p(\underline{x} | \omega_j) = \prod_{i=1}^{d} p(x_i | \omega_j, \underline{\nu}_{ij})$    $\forall \omega_j$

- Bayes plugin : $p(x_i | \omega_j; \underline{\nu}_{ij})$ known except for $\underline{\nu}_{ij}$

$\Rightarrow$ Estimate $\underline{\nu}_{ij}$ from $x_i$ of $N_j$ training samples.

Special case : Naive Gaussian classifier

- Gaussian likelihood : $\underline{x} | \omega_j \sim N(\underline{\mu}_j, \underline{\underline{C}}_j)$

- "naive" : $\underline{\underline{C}}_j = \begin{bmatrix} \sigma_{1,j}^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_{d,j}^2 \end{bmatrix}$

- ML parameter estimator:

  $\hat{\mu}_j$ : as in E 4.4

  $\delta_{ij}^2$ : diagonal elements of $\hat{\underline{\underline{C}}}_{\theta j}$ in E 4.4

## Benefit:

- Less parameters to estimate
- Less training samples as we have less parameters to estimate

  → lower complexity

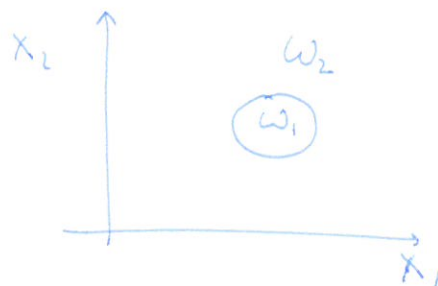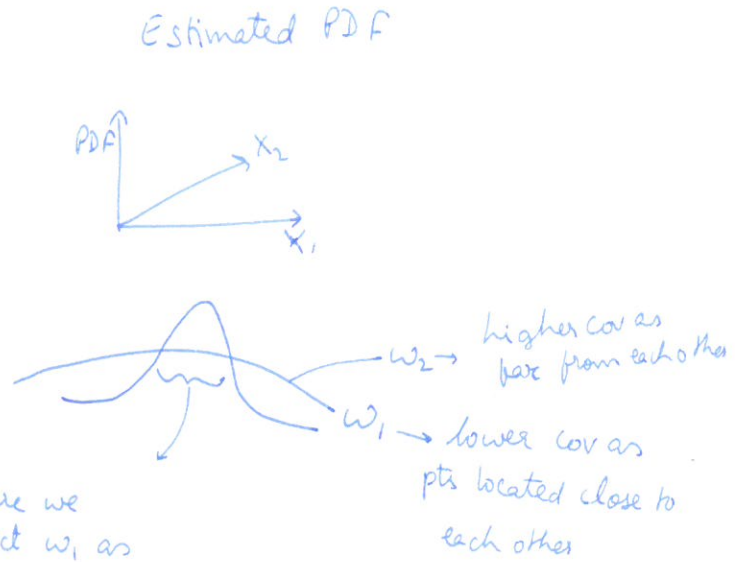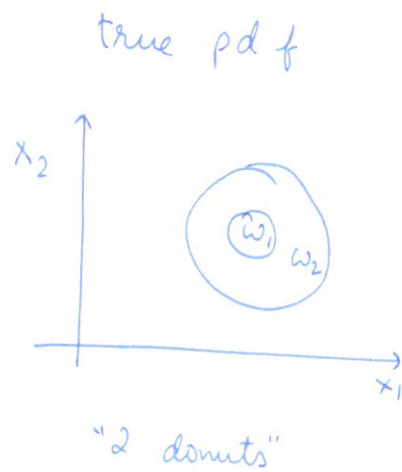| # parameters | Gaussian | naive Gaussian |
|---|---|---|
| $\underline{\mu}_j$ ▯ | $d$ | $d$ |
| $\underline{\underline{C}}_j$ ▽ | $\dfrac{d(d+1)}{2}$ | $\setminus d$ |
| $\Sigma$ | $\dfrac{d(d+3)}{2}$ | $2d$ |

eg $d = 100$    5150    200    per class

## Drawbacks

- rarely independent features in real-life, naive approx
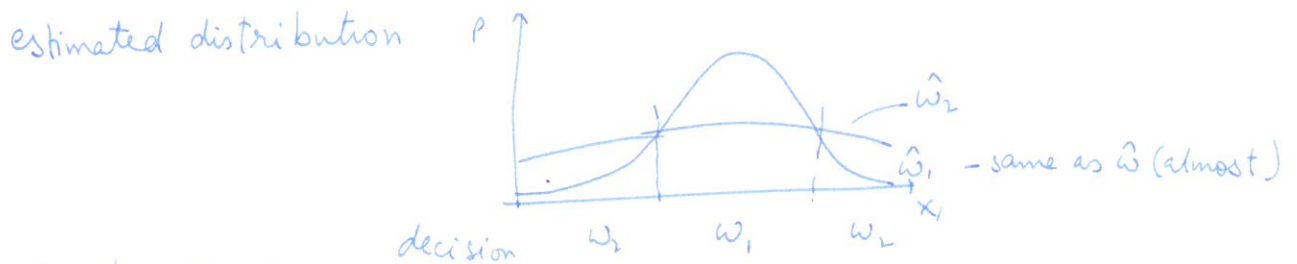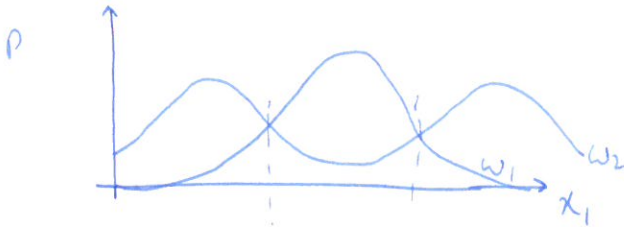
Explaination to data set 2 :

true pdf



"2 donuts"

Estimated PDF



$w_2 \to$ higher cov as far from each other

$w_1 \to$ lower cov as pts located close to each other

here we predict $w_1$ as
$$P(w_1) > P(w_2)$$
outside we have $w_2$ as
$$P(w_2) > P(w_1)$$
Hence    we get,



bad estimate of $P(\underline{x}/w_1)$, $P(\underline{x}/w_2)$, but we have still satisfied decision boundary. $\to$ we have luck.

Explaination to dataset 3:

true PDF :    $x_1, x_2 \rightarrow$ independent

$x_2 \rightarrow$ uniform in $[0,1]$



estimated distribution



– same as $\hat{\omega}$ (almost)

good estimate for $P(\underline{x}|\omega_1)$

bad estimate for $P(\underline{x}|\omega_2)$     } we have luck

satisfied decision boundary

Explaination to dataset 4:

true PDF



we need 3 lines/boundaries

estimated PDF



only 1 decision boundaries

bad estimate of $P(\underline{x}|\omega_1), P(\underline{x}|\omega_2)$

bad decision boundary ; hence high error rate

# 4.3.4  Gaussian Mixture Models

Upto now: One Gaussian likelihood $N(\underline{\mu}, \underline{\underline{\Sigma}})$ per class

$$\underline{\mathcal{V}} = \{\underline{\mu}, \underline{\underline{\Sigma}}\} \text{ contains } \frac{d(d+3)}{2} \text{ parameters}$$

Limitations → only one mode/component/maximum
  → unimodal pdf

Wish → a multimodal pdf

E4.7: Multimodal PDF

(a) Fish:
    multiple species of fish:



$x = $ fish length

(b) Emotion recognition:
    from speech



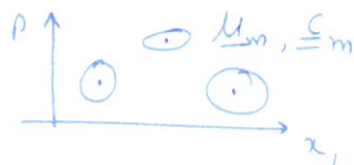$x = $ pitch freq

male     female   kid
         male

Idea: Gaussian Mixture Model (GMM), a mixture of
      Gaussian PDFs per class

GMM for one class:

· M Gaussian modes/sub-classes:

$$P_m(\underline{x}; \underline{\nu}_m) \sim N(\underline{\mu}_m, \underline{\underline{C}}_m)$$

$$\underline{\nu}_m = \{\underline{\mu}_m, \underline{\underline{C}}_m\} \quad 1 \leq m \leq M$$



- Random mode variable $Z \in \{1 \dots M\}$
  (sub class label)

  $Z = m$ : $\underline{x}$ generated by mode $m \;\hat{=}\;$ subclass label $m$.

- law of total probability:

  PDF $\quad p(\underline{x}; \underline{\nu}_{\emptyset}) = \sum\limits_{m=1}^{M} \underbrace{p(\underline{x}|z=m)}_{P_m(\underline{x}; \underline{\nu}_m)} \underbrace{P(z=m)}_{\alpha_m}$

  $$= \sum\limits_{m=1}^{M} \alpha_m P_m(\underline{x}; \underline{\nu}_m)$$

- $\alpha_m = P(z=m) \geqslant 0 \quad \sum\limits_{m=1}^{M} \alpha_m = 1$ : prior of $Z$, mode weight,
  
  $\hspace{8cm} M-1$ parameters

- $\underline{\nu} = \{\underline{\nu}_1 \dots \underline{\nu}_M^{\bullet}; \alpha_1 \dots \alpha_{M-1}\}$ contains $M \underbrace{\dfrac{d(d+3)}{2}}_{\text{for each } \underline{\nu}_i} + M-1$ parameters

  $\hspace{7cm}$ if $\underline{\underline{C}}_m$ is non-diagonal

  If naive then $\quad 2Md + M-1$ if $\underline{\underline{C}}_m$ is diagonal.

i.e. 2 levels of class label:
- class label $\omega$: $\hspace{3cm} \omega_1 \hspace{3cm} \omega_2$
  $\hspace{2cm}$ to be classified with labelled samples

- mode $m$: $\hspace{2cm} 1 \dots \quad M_1 \hspace{3cm} 1 \dots M_2$
  $\hspace{2cm}$ not relevant for classification, no labels
  $\hspace{2cm}$ but relevant for GMM parameter estimation

$\underline{Q}$ How to estimate $\underline{\nu}$ ?

<u>A</u>   Direct ML parameter estimation

$$\max_{\underline{v}} \; p(\underline{x}; \underline{v})$$

. $M = 1$ : closed form solution see 4.4

. $M > 1$ : no closed form sol$^n$ , hard to solve.

slide 4.34

<u>B</u>   Assumption: $N$ iids measurements of $\underline{x}_n \in \mathbb{R}^d$

and $z_n \in \{1, \cdots M\}$ $1 \leq n \leq N$

$z_n$ : additional mode measurements

Let $\underline{X} = \begin{bmatrix} \underline{x}_1 \\ \vdots \\ \underline{x}_N \end{bmatrix}$ $\underline{Z} = \begin{bmatrix} z_1 \\ \vdots \\ z_N \end{bmatrix}$

continuous          discrete
valued               valued

joint PDF/PMF model for $(\underline{X}, \underline{Z})$:

$$p(\underline{X}, \underline{Z}; \underline{v}) = \prod_{n=1}^{N} p(\underline{x}_n, z_n; \underline{v})$$

$$p(\underline{x}_n, z_n; \underline{v}) = \underbrace{p(\underline{x}_n | z_n)}_{P_{z_n}(\underline{x}_n \text{ ; } \underline{v})} \underbrace{p(z_n)}_{\alpha_{z_n}}$$

$$= \alpha_{z_n} \, P_{z_n}(\underline{x}_n; \underline{v})$$

$$\max_{\underline{v}} \; p(\underline{X}, \underline{Z}; \underline{v})$$

unrealistic due to missing $\underline{Z}$

c) Consider $\underline{Z}$ as unknown: ("missing variable)

$$\max_{\underline{\nu}} \ P \ (\underline{X}, \underline{z}; \underline{\nu})$$

How? Expectation Maximization Algorithm (EM Algorithm)