$X$ — known $\underline{X}_n$

$\underline{\mathcal{V}}$ — unknown parameters

$\underline{Z}_n$ — unknown mode measurements

It is difficult to exactly or partially determined the $X$ and $\underline{\mathcal{V}}$ together i.e. coupled. However, if we decouple them then determining them becomes a relatively easy Solution. This is where EM (Expectation Maximization) Algorithm comes into play.

4-35

4-36

We have bad/missing data $\underline{z}$ and that is to be eliminated.

In the eq$^n$:

$$\underline{\varrho}(\underline{\Theta}) = \int \ln p(\underline{x}, \underline{z}, \underline{\Theta}) \, p(\underline{z} \mid \underline{x}, \underline{\Theta}) \, dz$$

we need to eliminate

Hence we do this step $p(\underline{z} \mid \underline{x}, \underline{\Theta})$

4-37

Apply the general EM algorithm to special GMM model.

4-38, 4-39

4-40

MATLAB :  gm distribution

GMM classifier :

c   GMM models

- one GMM  $p\left(\underline{x} \mid w_j ; \underline{\vartheta}_j\right)$  per class $w_j$   $1 \leq j \leq c$



No of modes $M_j$ per class $w_j$ may not be same for each class.

- run EM algorithm  c times

- then have a better idea

4-41

4-42

GMM classifier contains

    Gaussian classifier        : $M_j = 1$     $\forall j$

and naive   "      "      : $M_j = 1$ , $\underline{\underline{C}}_j$ diagonal $\forall j$

as   special cases

<u>Q</u>   How to choose the no of modes?

<u>Model order estimation</u> :

    But model orders are $M_j$ are usually unknown

    → <u>need order estimation</u> : Estimate

         $M_j$ (and $\underline{\nu}_j$) from data

         ↑         ↑

  discrete valued    continuous valued

  2 <u>popular criteria</u> :

  Akaike Information Criteria (AIC)

  Bayesian     "        "      (BIC)

<u>Q</u> How to know which one to use?

  Intuition + perform over all possibilities and then find
the one which gives the best result.
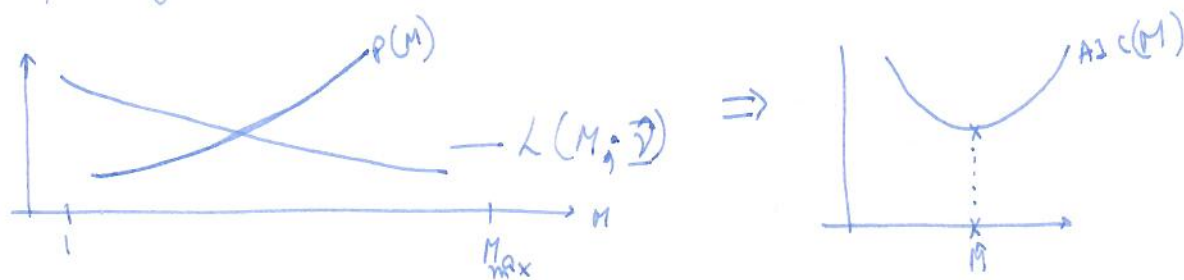
for one GMM:

- Try EM for $M = 1 \cdots M_{MAX}$

- Compute $\text{AIC}(M) = \cancel{8} - \underset{\uparrow}{\angle}(M; \hat{\underline{\vartheta}}) + \overset{\downarrow \text{penalty term}}{P(M)}$

  $\log$ likelihood $\{\ln P(\underline{X}, \underline{Z}; \underline{\vartheta})$ in EM

- $\hat{M} = \arg \underset{M}{\min} (\text{AIC}(M))$

**Q** Why penalty term



Reason for decay — if we increase M, then we have more
parameters → better fit of data to data → $\angle(M; \hat{\underline{\vartheta}}) \uparrow$
→ overfitting

Hence we add $P(M)$ to avoid overestimation of M

Difference b/w AIC and BIC:

Choice of penalty term $P(M)$

4-43, 4-42

4-1 $\longrightarrow$ Now we will talk about density estimation (Parzen window)

## 4.4 Density estimation

Idea

- Non-parametric estimate of $P(\underline{x} \mid \omega_j)$
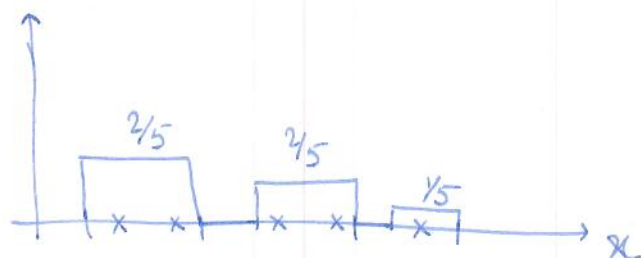
- Apply Bayesian decision theory

eg  Parzen window method

For one class:

Given: N iid sample of $\underline{x}_n$ from an unknown $P(\underline{x})$

Desired: PDF estimate $\hat{p}(\underline{x})$

Simples : Normalized histogram



Total area under curve should be 1, however this is not smooth. To smooth this, we need kernels.

Better : Kernel based PDF



smoothed version

Kernel: in $\mathbb{R}^d$ : a mathematical PDF-like function $\phi(\underline{x})$

- $\phi(\underline{x}) \geqslant 0 \quad \forall \underline{x}$

- $\int \phi(\underline{x}) \, d\underline{x} = 1$

Gaussian kernel :  $N(\underline{0}, \underline{\underline{I}})$
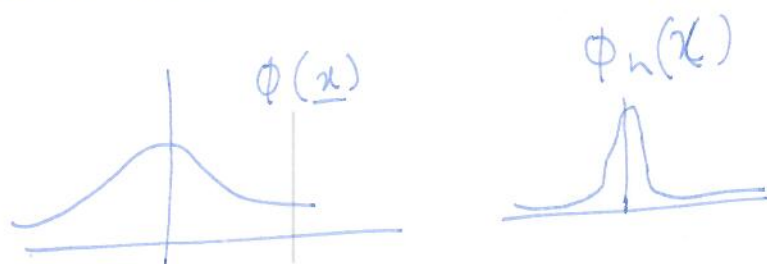
$$\phi(\underline{x}) = \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{1}{2} \|\underline{x}\|^2\right),$$

has a fixed width (SD)

Scaled kernel :

$$\phi_h(\underline{x}) = \frac{1}{h^d}\phi\left(\frac{1}{h}\underline{x}\right) \geqslant 0, \quad \int \phi_h(\underline{x})\, d\underline{x} = 1$$

$h \to$ bandwidth parameter   $h > 0$



$\phi(x)$    $\phi_h(x)$

$h \downarrow =$ narrow kernel $\to$ detailed, however sometimes
too much details in PDF

$h \uparrow =$ broad kernel $\to$ less detailes.

PDF estimate :

$$\hat{p}(\underline{x}) = \frac{1}{N} \sum_{n=1}^{N} \phi_h(\underline{x} - \underline{x}_n) \geqslant 0, \quad \int \hat{p}(\underline{x})\, d\underline{x} = 1$$

4-45
4-46