



Pattern Recognition

Approximating class densities, Bayesian classifier



Bayesian classification

- Suppose we have 2 classes and we know probability density functions of their feature vectors. How some new pattern should be classified?
- Bayes classification rule: classify x to the class w_i which has biggest posterior probability $P(w_i | x)$

$$P(w_1 | x) > P(w_2 | x) ? \quad w_1 : w_2$$

posterior

Using Bayes formula, we can rewrite classification rule:

$$p(x | w_1)P(w_1) > p(x | w_2)P(w_2) ? \quad w_1 : w_2$$

likelihood *prior*



Estimating probability density function.

- Parametric pdf estimation: model unknown probability density function $p(x | w_i)$ of class w_i by some parametric function $p_i(x; \theta)$ and determine parameters based on training samples.

Example: Gaussian function

$$p(x; \mu) = \frac{1}{(2\pi)^{1/2}} e^{-\frac{1}{2}(x-\mu)^2}$$

- Non-parametric pdf estimation:
 1. Histogram
 2. K nearest neighbor
 3. Kernel methods (Parzen kernels or windows)

$$\hat{p}(x) = \frac{1}{N} \sum_{i=1}^N \left(\frac{1}{h} \varphi \left(\frac{x_i - x}{h} \right) \right) \quad N \text{ is the number of training samples}$$

4. Other methods (estimating cumulative distribution function first, SVM density estimation, etc.)



Estimating kernel width

- Non-parametric pdf estimation:
 - Fixed kernels:

$$\hat{p}(x) = \frac{1}{N} \sum_{i=1}^N \left(\frac{1}{h} \varphi \left(\frac{x_i - x}{h} \right) \right)$$

- Adaptive kernels:

$$\hat{p}(x) = \frac{1}{N} \sum_{i=1}^N \left(\frac{1}{h_i} \varphi \left(\frac{x - x_i}{h_i} \right) \right)$$

or

$$\hat{p}(x) = \frac{1}{N} \sum_{i=1}^N \left(\frac{1}{h_i} \varphi_i \left(\frac{x - x_i}{h_i} \right) \right)$$

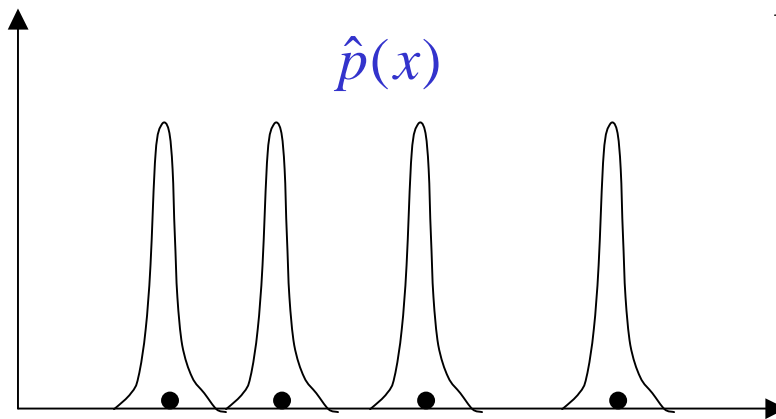


Estimating kernel width

Recall, we used maximum likelihood method for parametric pdf estimation:

$$\max_{\theta} \hat{p}(X; \theta) = \max_{\theta} \hat{p}(x_1, x_2, \dots, x_N | \theta) = \max_{\theta} \prod_{k=1}^N \hat{p}(x_k; \theta)$$

Can we use same method for estimating the kernel width h ?



No, the max is not achievable:

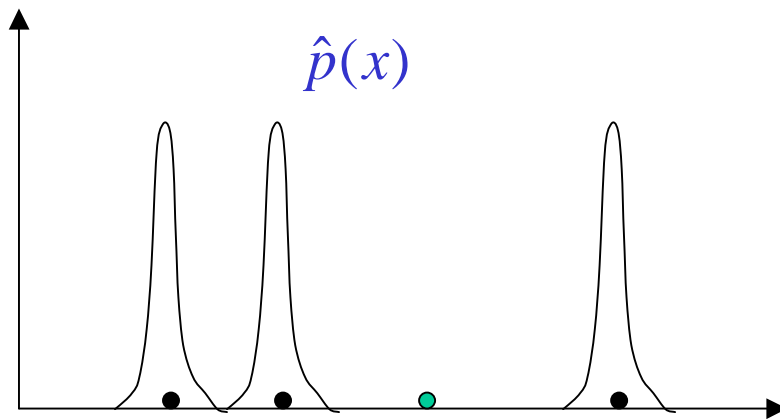
$$\begin{aligned} \max_h \prod_{k=1}^N \hat{p}(x_k; h) &= \\ \max_h \prod_{k=1}^N \left(\frac{1}{N} \sum_{i=1}^N \frac{1}{h} \varphi\left(\frac{x_i - x_k}{h}\right) \right) &\rightarrow \infty \\ \text{if } h &\rightarrow 0 \end{aligned}$$



Estimating kernel width

Solution: separate model data (kernel centers) from testing data
 - cross-validation technique

$$\max_h \prod_{k=1}^N \left(\frac{1}{N} \sum_{i \neq k} \frac{1}{h} \varphi \left(\frac{x_i - x_k}{h} \right) \right)$$



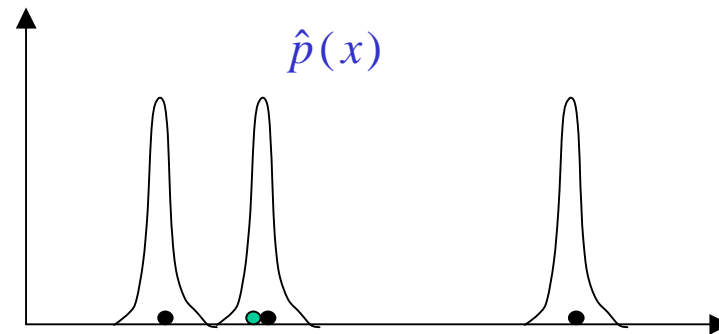


Estimating kernel width

Tried maximum likelihood cross-validation and still diverges?

$$\max_h \prod_{k=1}^N \left(\frac{1}{N} \sum_{i \neq k} \frac{1}{h} \varphi \left(\frac{x_i - x_k}{h} \right) \right) \rightarrow \infty$$

This might happen if data is somewhat discrete:

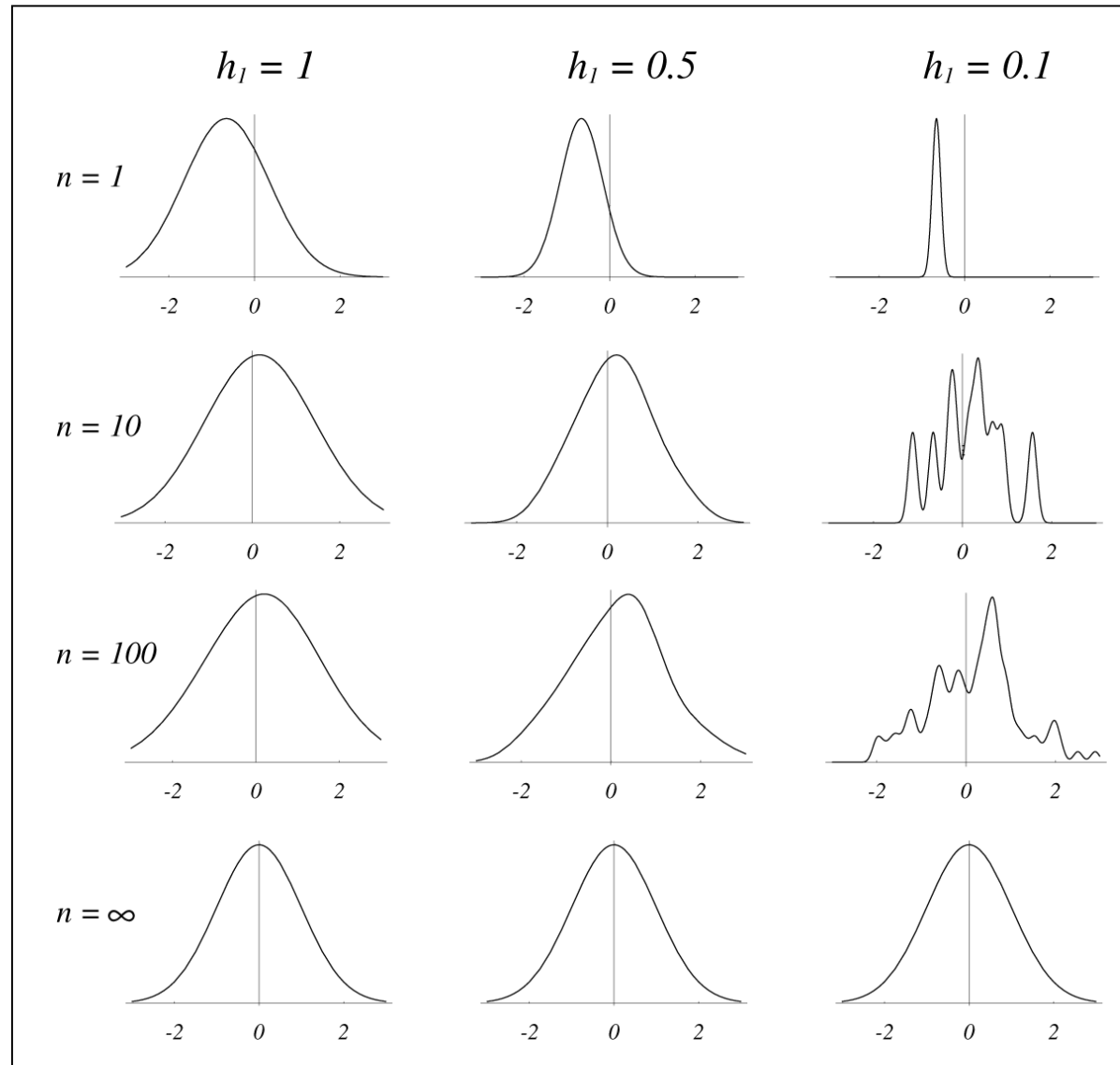


Solution - truly separate model data from testing data:

$$\max_h \prod_{k=1}^N \left(\frac{1}{N} \sum_{x_i \neq x_k} \frac{1}{h} \varphi \left(\frac{x_i - x_k}{h} \right) \right)$$



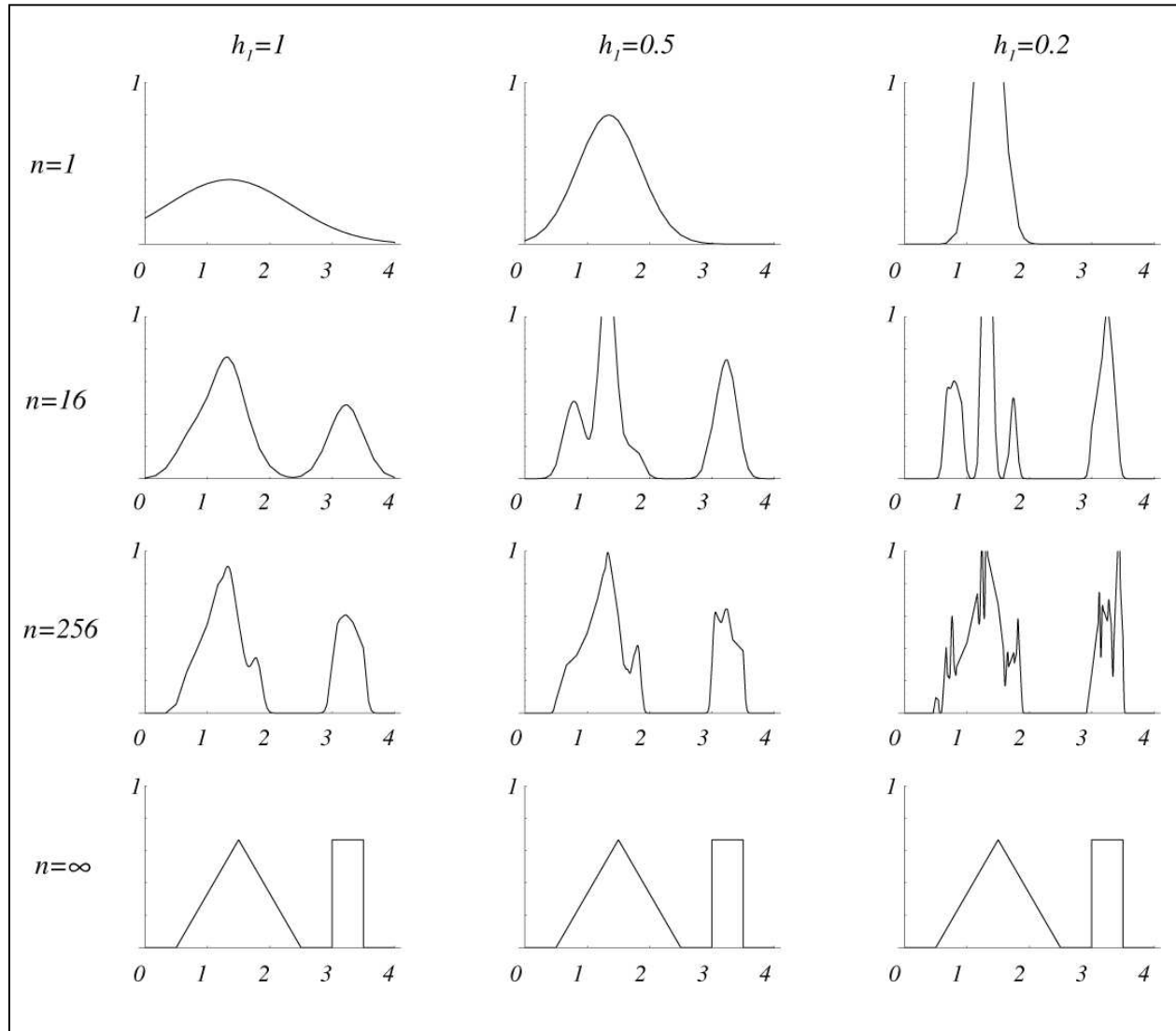
Examples of pdf estimation



Parzen-window (kernel) estimates of a univariate normal density using different window widths and numbers of samples.



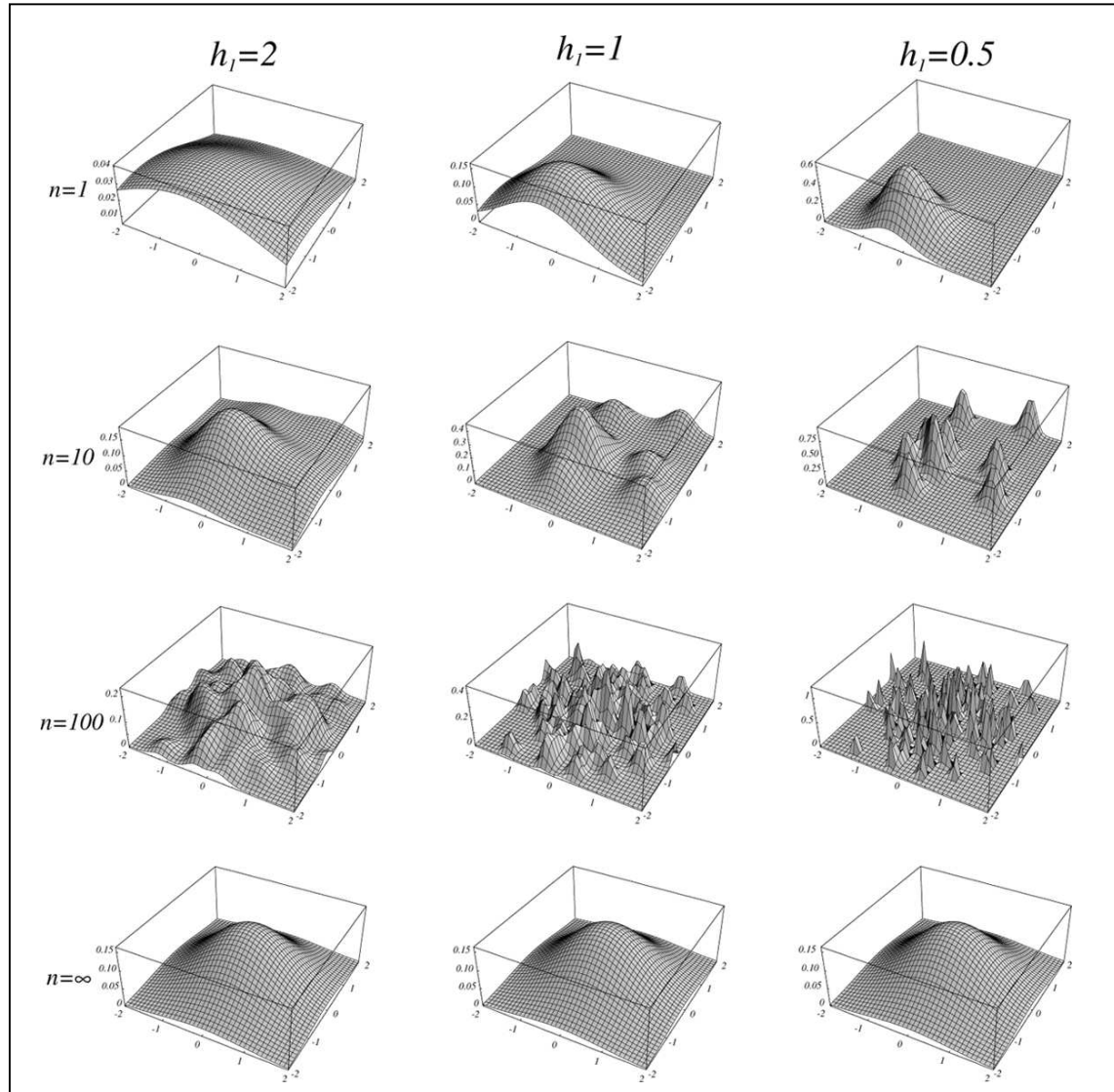
Examples of pdf estimation



Parzen-window (kernel) estimates of a bimodal density using different window widths and numbers of samples.



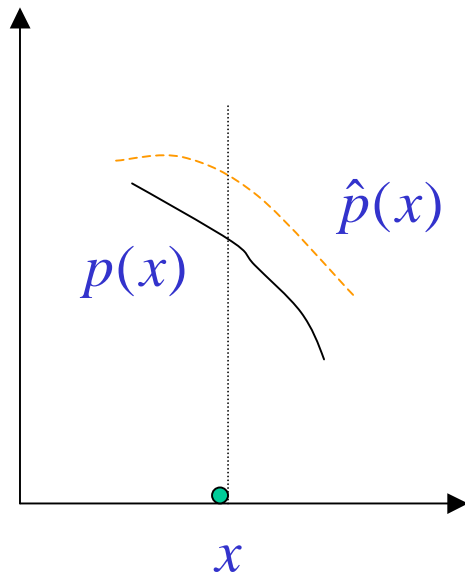
Examples of pdf estimation



Parzen-window (kernel) estimates of a bivariate normal density using different window widths and numbers of samples.



Error in pdf estimation



Discrepancy between true density $p(x)$ and its estimation $\hat{p}(x)$:

$$MSE_x(\hat{p}) = E\{\hat{p}(x) - p(x)\}^2$$

- Mean Square Error

$$MISE(\hat{p}) = \int E\{\hat{p}(x) - p(x)\}^2 dx$$

- Mean Integrated Square Error

$$MSE_x(\hat{p}) = E\{\hat{p} - p\}^2 = E\{\hat{p}^2 - 2\hat{p}p + p^2\}$$

$$= E\{\hat{p}^2\} - 2E\{\hat{p}\}p + p^2$$

$$= \{E\hat{p}\}^2 - 2\{E\hat{p}\}p + p^2 + [E\{\hat{p}^2\} - \{E\hat{p}\}^2]$$

$$= [E\hat{p} - p]^2 + [E\{E\hat{p} - \hat{p}\}^2]$$



Bias and variance of estimation error

$$MSE_x(\hat{p}) = \underbrace{[E\hat{p} - p]^2}_{\text{Bias}} + \underbrace{[E\{E\hat{p} - \hat{p}\}^2]}_{\text{Variance}}$$

Bias

Variance

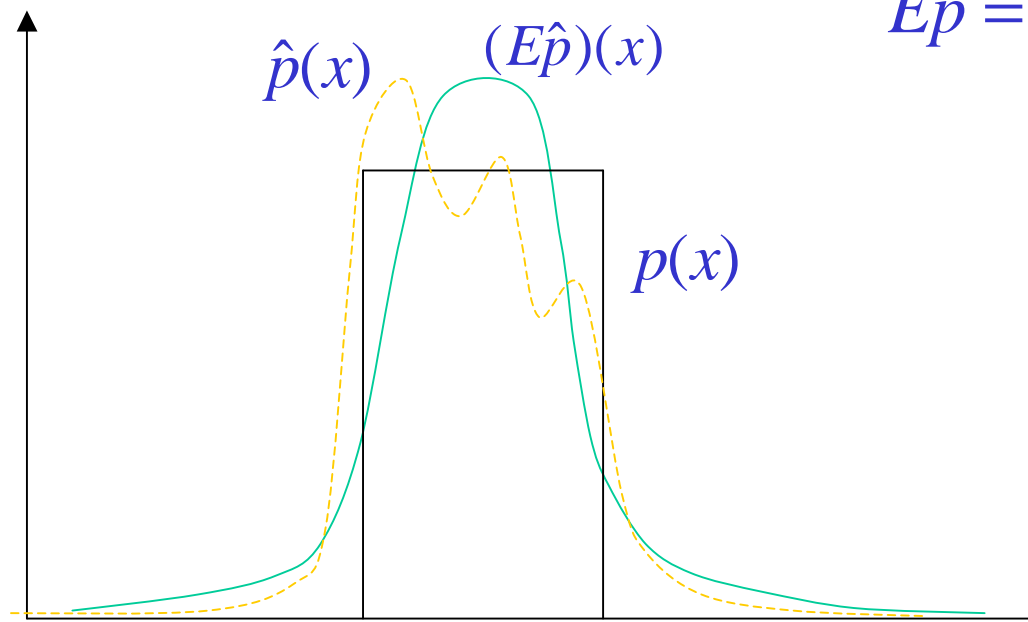
$E\hat{p}$ - Average approximation

$$E\hat{p} = \int \frac{1}{h} K\left(\frac{x-y}{h}\right) p(y) dy$$

Bias is the difference between true density and average approximation

Variance is the difference between average approximation and individual approximations

Smaller kernel width reduces bias, but increases variance.





Bias and variance of estimation error

If some assumptions on the true density are made (e.g. $\int (p''(x))^2 dx < \infty$) then it is possible to analytically find the kernel width which gives smallest $MISE(\hat{p})$

Silverman (Parzen):

$$h_{opt} = k_2^{-2/5} \left\{ \int \varphi(t)^2 dt \right\}^{1/5} \left\{ \int p''(x)^2 dx \right\}^{-1/5} n^{-1/5}$$

Optimal kernel width gets smaller when the number of training samples n increases. For optimal kernel width $MISE(\hat{p})$ also decreases:

$$MISE \sim C(\varphi) \left\{ \int p''(x)^2 dx \right\}^{1/5} n^{-4/5}$$

Note, that $p(x)$ is unknown. Above formulas are useful for theory, but not for practical applications.

For multivariate pdf approximation: $MISE \sim n^{-4/(4+d)}$

The performance decreases exponentially when the number of dimensions increases



Bayesian classification

- Bayes classification rule: classify x to the class w_i which has biggest posterior probability $P(w_i | x)$

$$P(w_1 | x) > P(w_2 | x) ? \quad w_1 : w_2$$

- Bayes classification rule minimizes the total probability of misclassification.

Cost of errors.

- Errors happen when samples of class 1 are incorrectly classified to belong to class 2, and samples of class 2 are classified to belong to class 1.
- The cost of making these errors can be different :

λ_1 - the cost of misclassifying samples of class 1

λ_2 - the cost of misclassifying samples of class 2



Total cost (or risk) of classification

Classification algorithm splits feature space into two decision regions:

R_1 - samples in this region are classified as being in class 1

R_2 - samples in this region are classified as being in class 2

$\int_{R_2} p(x | w_1) dx$ - the proportion of samples of class 1 being classified as class 2

$\int_{R_1} p(x | w_2) dx$ - the proportion of samples of class 2 being classified as class 1

$P(w_1) \int_{R_2} p(x | w_1) dx$ - the proportion of all input samples being class 1 but classified as being in class 2

$P(w_2) \int_{R_1} p(x | w_2) dx$ - the proportion of all input samples being class 2 but classified as being in class 1

$$\text{Cost} = \lambda_1 P(w_1) \int_{R_2} p(x | w_1) dx + \lambda_2 P(w_2) \int_{R_1} p(x | w_2) dx \quad \text{- total cost}$$



Minimizing total cost of classification

Since R_1 and R_2 cover whole feature space

$$\int_{R_1} p(x | w_1) dx + \int_{R_2} p(x | w_1) dx = 1$$

Thus

$$\begin{aligned}
 \text{Cost} &= \lambda_1 P(w_1) \left\{ 1 - \int_{R_1} p(x | w_1) dx \right\} + \lambda_2 P(w_2) \int_{R_1} p(x | w_2) dx \\
 &= \lambda_1 P(w_1) + \int_{R_1} (\lambda_2 P(w_2) p(x | w_2) - \lambda_1 P(w_1) p(x | w_1)) dx
 \end{aligned}$$

Cost is minimized if R_1 includes only points where

$$\lambda_2 P(w_2) p(x | w_2) - \lambda_1 P(w_1) p(x | w_1) < 0$$



Bayesian classification

Bayesian classifier is an optimal classifier minimizing total classification cost. Such classifier is possible only if we have full knowledge about class distributions.

If $\lambda_1 P(w_1) p(x | w_1) > \lambda_2 P(w_2) p(x | w_2)$ then classify x as class 1.

If $\lambda_1 P(w_1) p(x | w_1) \leq \lambda_2 P(w_2) p(x | w_2)$ then classify x as class 2.

Alternatively, assuming non-zero terms, the class assignment is based on

testing whether $\frac{p(x | w_1)}{p(x | w_2)} > \frac{\lambda_2 P(w_2)}{\lambda_1 P(w_1)}$ or $\frac{p(x | w_1)}{p(x | w_2)} \leq \frac{\lambda_2 P(w_2)}{\lambda_1 P(w_1)}$

Decision surface $\frac{p(x | w_1)}{p(x | w_2)} = \frac{\lambda_2 P(w_2)}{\lambda_1 P(w_1)}$ separates two decision regions.

$\frac{p(x | w_1)}{p(x | w_2)}$ - likelihood ratio

$\frac{p(x | w_1)}{p(x | w_2)} > (<) \frac{\lambda_2 P(w_2)}{\lambda_1 P(w_1)}$ - likelihood ratio test



Performance of Bayesian classification

Denote:

$$t = \frac{\lambda_2 P(w_2)}{\lambda_1 P(w_1)} \quad \text{- decision threshold}$$

$$R_1(t) = \left\{ x \mid \frac{p(x | w_1)}{p(x | w_2)} > t \right\} \quad \text{- decision region of class 1 for threshold } t$$

$$R_2(t) = \left\{ x \mid \frac{p(x | w_1)}{p(x | w_2)} \leq t \right\} \quad \text{- decision region of class 2 for threshold } t$$

$$MR_1(t) = \int_{R_2(t)} p(x | w_1) dx \quad \text{- misclassification rate for class 1 and threshold } t$$

$$MR_2(t) = \int_{R_1(t)} p(x | w_2) dx \quad \text{- misclassification rate for class 2 and threshold } t$$



Performance of Bayesian classification

$MR_1(t)$ and $MR_2(t)$ completely characterize the performance of a Bayesian classifier

For a given misclassification costs λ_1, λ_2 and prior class probabilities $P(w_1), P(w_2)$ we find $t = \frac{\lambda_2 P(w_2)}{\lambda_1 P(w_1)}$

Then the (mis)classification cost is

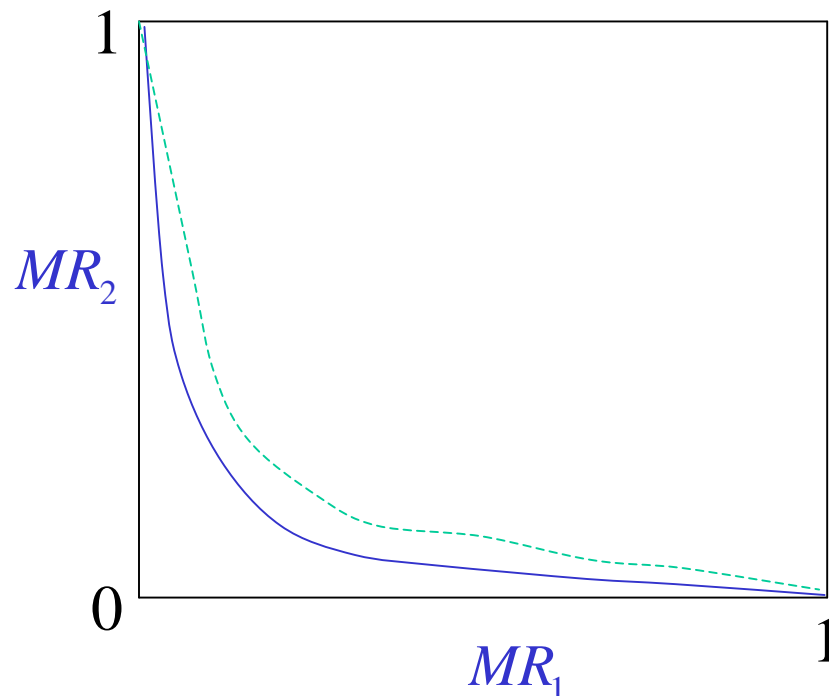
$$\text{Cost} = \lambda_1 P(w_1) MR_1(t) + \lambda_2 P(w_2) MR_2(t)$$



ROC of a Bayesian classification

$MR_1(t)$ and $MR_2(t)$ are used only with the the same t .

Thus the parameter t is not important and the performance of a Bayesian classifier can be characterized only by the relationships between $MR_1(t)$ and $MR_2(t)$.



Example of an optimal Bayesian ROC curve (—) and some non-optimal classifier's ROC curve (- - -).

For a given MR_1 the MR_2 of a non optimal classifier should be bigger; otherwise non-optimal classifier would outperform optimal.



Performance of Biometric Matchers

For biometric matchers (person identity verification) we distinguish two classes:

- Genuine – person's claimed identity is correct
- Impostor - person's claimed identity is incorrect

The decision for genuine class is to accept, and the decision for the impostor class is to reject. The decision is usually done based on a single matching score of input biometric with the enrolled biometric template of claimed identity person.

Instead of optimal $\frac{p(x | w_1)}{p(x | w_2)} > (<) \theta$ use $x > (<) \theta$

If $\frac{p(x | w_1)}{p(x | w_2)}$ is monotonous, these decisions are equivalent.

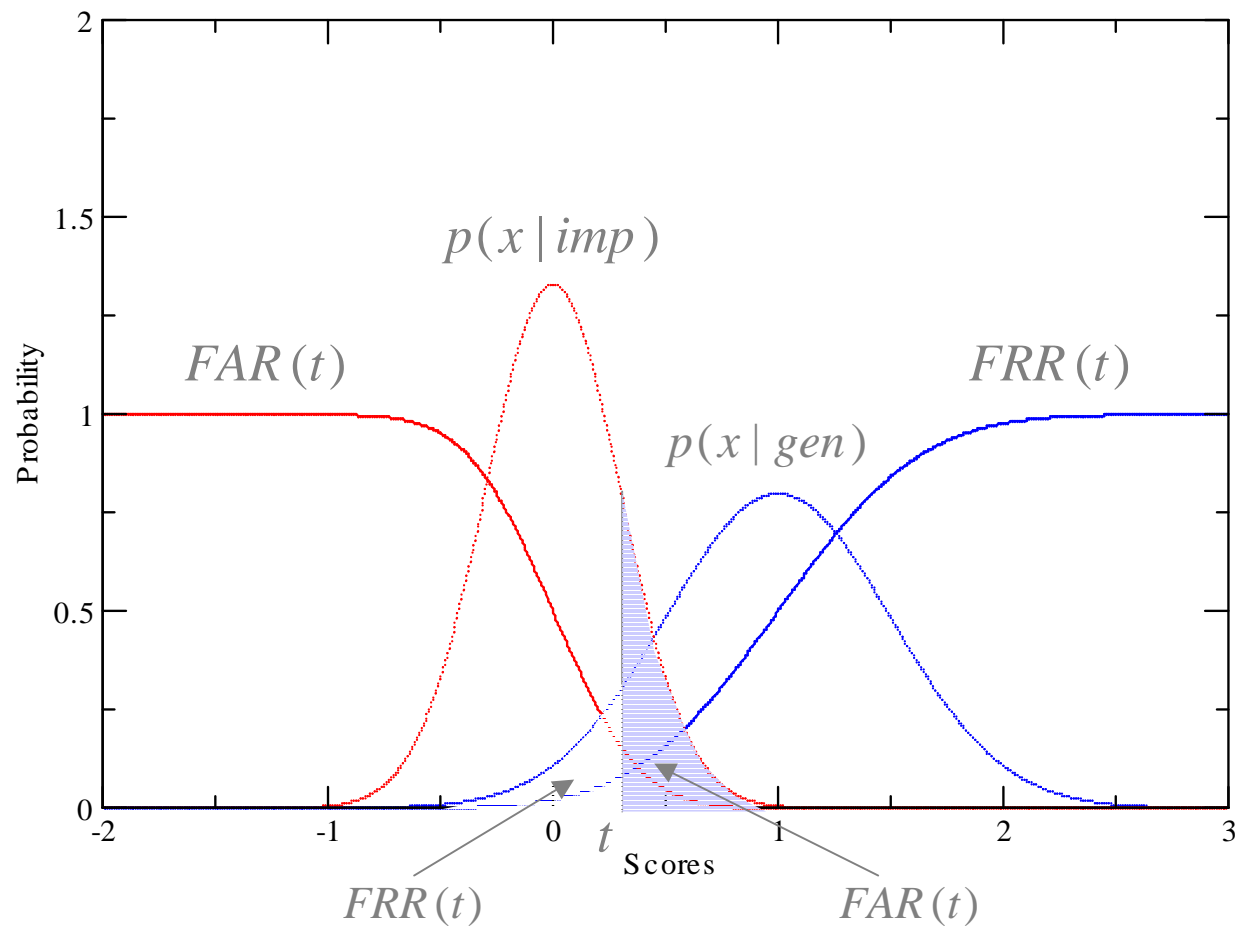
Instead of $MR_1(t)$ and $MR_2(t)$ use

$$FAR(t) = \int_{x>t} p(x | imp) dx \quad \text{- false accept rate for threshold } t$$

$$FRR(t) = \int_{x<t} p(x | gen) dx \quad \text{- false reject rate for threshold } t$$



Performance of Biometric Matchers





Estimating FAR and FRR

In contrast to estimating pdf, FAR and FRR are easily estimated:

$$FAR(t) = \int_{x>t} p(x | imp) dx \approx \frac{|\{x_i | x_i > t, x_i \text{ is impostor}\}|}{|\{x_i | x_i \text{ is impostor}\}|}$$

$$FRR(t) = \int_{x<t} p(x | gen) dx \approx \frac{|\{x_i | x_i < t, x_i \text{ is genuine}\}|}{|\{x_i | x_i \text{ is genuine}\}|}$$

Types of ROC curves:

$$\{FRR(t), FAR(t)\}_{-\infty < t < \infty}$$

$$\{FAR(t), P(gen)(1 - FRR(t)) + P(imp)FAR(t)\}_{-\infty < t < \infty}$$

$$\{\log FRR(t), \log FAR(t)\}_{-\infty < t < \infty}$$