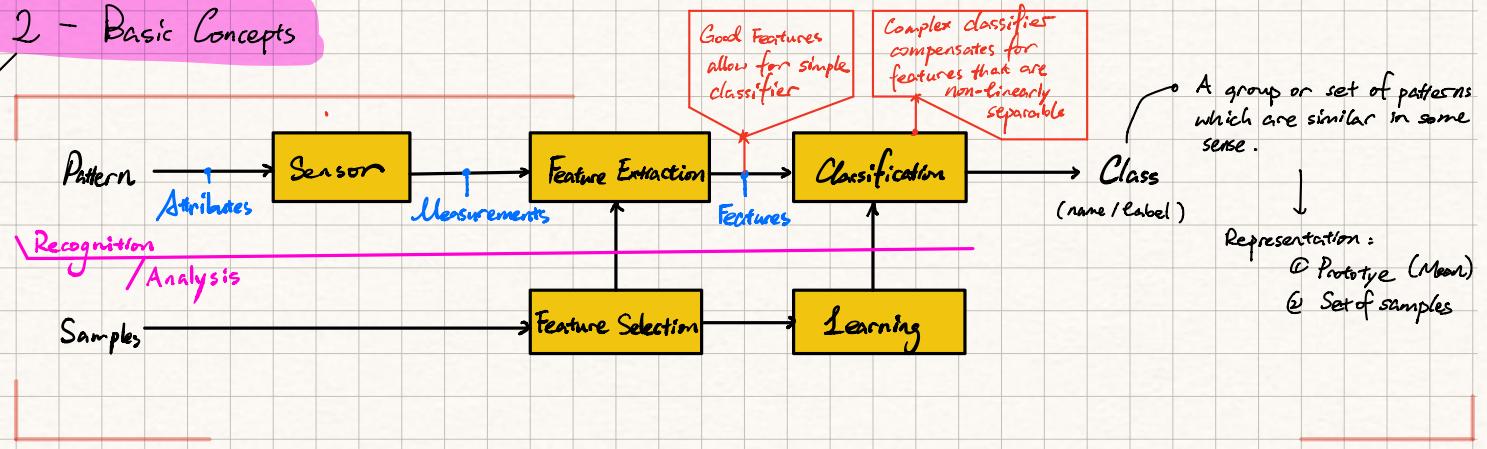


§ 2 - Basic Concepts



Classification Types (Three)

① Known probability model for each class. → Approach: Statistical decision theory ⇒ find optimal classifier by minimizing prob. of error.

② Known class labels for the samples. → Approach: ① Learn empirical probability model based on samples
② Derive classifiers directly from distribution of samples in feature space.

③ No known label info for the sample → Need to determine # of classes & class definition.

↳ Refers as Clustering problem.

Approach: look for naturally occurring order, groupings of clusters in data.

Summary:

Two Patterns are Similar \Leftrightarrow sharing common properties

(vector space rep):

Closeness in feature space

measured by distance metric $d(\vec{y}_1, \vec{y}_2)$.

§ 3 - Pattern Representation Random Vectors

- In Statistical Pattern Rec,

Pattern \Leftrightarrow Random Vectors : Vector of random variables

* Random Variables

- a single scalar that's random (x)
- Discrete : random integers
- Continuous : random real #.
- Characterized by:

CDF: $F_x(\tau) = \Pr(x < \tau)$

PDF: $p(x)$, where:

$$F_x(\tau) = \Pr(x < \tau) = \int_{-\infty}^{\tau} p(x) dx.$$

with

$$\int_{-\infty}^{\infty} p(x) dx = 1$$

Expectation:

$$E[f(x)] = \int_{-\infty}^{\infty} f(x) p(x) dx$$

Choice of
 $f(x)$

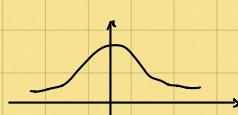
$$\text{Mean} : \mu_x = E[x]$$

$$\text{Variance} : \sigma_x^2 = E[x^2] - E[x]^2 = E[(x - E[x])^2]$$

* Common Distribution Models

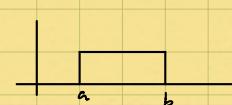
Gaussian:

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma_x} e^{-\frac{1}{2}(\frac{x-\mu_x}{\sigma_x})^2}$$



Uniform:

$$p(x) = \begin{cases} 0 & x < a \\ \frac{1}{b-a} & a \leq x \leq b \\ 0 & b < x \end{cases}$$



Exponential:

$$p(x) = \begin{cases} 0 & x < 0 \\ \lambda e^{-\lambda x} & x \geq 0 \end{cases}$$



* Joint Statistics

- relationship of 2 random variables x & y

Joint Probability Distribution $p(x,y)$:

$$\Pr(x < \alpha, y < \beta) = \int_{-\infty}^{\alpha} \int_{-\infty}^{\beta} p(x,y) dx dy$$

Expectation:

$$E[f(x,y)] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} p(x,y) f(x,y) dx dy$$

$$E[x+y] = E[x] + E[y]$$

Marginal Distribution $p(x)$:

$$p(x) = \int_{-\infty}^{\infty} p(x,y) dy.$$

Correlation:

$$\text{Variance: } E[(x - \mu_x)^2]$$

Normalized by std. deviation.

$$E[(x - \mu_x)(y - \mu_y)]$$

* Independence

$$\text{if } p(x,y) = p(x)p(y)$$

\Rightarrow knowing x tells nothing about y .

Independence \Rightarrow Uncorrelated

Uncorrelated

$$\text{if } E[xy] = E[x]E[y]$$

Uncorrelated \neq Independence

Correlation Coefficient:

$$\rho_{x,y} = \frac{E[(x - \mu_x)(y - \mu_y)]}{\sigma_x \sigma_y}$$

\downarrow measures ability to predict ' y ' as linear func of ' x '.

- $\rho_{x,y} = 0 \Rightarrow$ No predictability (uncorrelated)
- $\rho_{x,y} = \pm 1 \Rightarrow$ Perfect predictability.
(x & y deterministically correlated)

* Conditional Statistics. \rightarrow assessing unknown conditioned on known (identity of pattern) (measurement)

$p(x|y)$: PDF for x conditional on y (measurement)

$p(x|A)$: PDF for x given event A (class).

[distribution of measurements given a class
[Likelihood Distribution]]

$p(A|x)$: [Posterior Distribution]

[distribution of class given measurements.]

* Baye's Rule

$$p(x|y) = \frac{p(x,y)}{p(y)} = \frac{p(y|x)p(x)}{p(y)}$$

Joint
conditional prob. marginal

$$\Rightarrow p(x|A) = \frac{p(A|x)p(x)}{p(A)}$$

Likelihood Posterior (Hard to get)
prior/marginal of A

$$p(A|x) = \frac{p(x|A)p(A)}{p(x)}$$

* Random Vectors

$$x \vec{x} = [x_1, \dots, x_n]^T$$

* PDF: (Joint density function) of random variables in vector:

$$Pr(x_1 < \tau_1, \dots, x_n < \tau_n) = \int_{-\infty}^{\tau_1} \dots \int_{-\infty}^{\tau_n} p(\vec{x}) d\vec{x}$$

Marginal Distribution of a subset:

$$P(x_1, \dots, x_{l-1}, x_{l+1}, x_n) = \int_{-\infty}^{\infty} p(\vec{x}) dx_l$$

Expectation:

$$E[f(\vec{x})] = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} p(\vec{x}) f(\vec{x}) d\vec{x}$$

p(x_1, \dots, x_n) (Joint prob.)

Mean = (vector)ⁿ

$$\vec{\mu}_x = E[\vec{x}]$$

Covariance: (matrix)^{nxn} \rightarrow correlation between x_i & x_j .

$$\Sigma_x = E[(\vec{x} - \vec{\mu}_x)(\vec{x} - \vec{\mu}_x)^T]$$

Diagonal(Σ_x) = variance of random variables
(+ve)

$$\Sigma_x = \Sigma_x^T \text{ (Diag Symmetry)}$$

Independent (\vec{x} & \vec{y})

$$\text{if } p(\vec{x}, \vec{y}) = p(\vec{x})p(\vec{y})$$

Bayes Rule:

$$p(\vec{x}|A) = \frac{p(A|\vec{x})p(\vec{x})}{p(A)}$$

Uncorrelated

$$\text{if } E[\vec{x}\vec{y}^T] = E[\vec{x}]E[\vec{y}^T]$$

* Sample Statistics

In reality, we don't know prob-density or mean or covariance of a random vector.
we infer actual observations: $\{x_1, x_2, \dots, x_n\}$

Sample Statistics

More samples $n \uparrow \Rightarrow$ Better Distribution

Sample Mean:

$$\vec{m}_x = \frac{1}{N} \sum_{i=1}^N \vec{x}_i$$

Sample Covariance:

$$\Sigma_x = E[\vec{x}\vec{x}^T] - E[\vec{x}]E[\vec{x}]^T$$

$$S_x = \frac{1}{N} \sum_{i=1}^N (\vec{x}_i - \vec{m}_x)(\vec{x}_i - \vec{m}_x)^T = \frac{1}{N} \sum_{i=1}^N (\vec{x}_i \vec{x}_i^T) - \vec{m}_x \vec{m}_x^T$$

Multivariable Gaussian Distribution

$$p(\vec{x}) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} e^{\left\{ -\frac{1}{2} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu}) \right\}}$$

Scalar

determinants of covariance
=(Magnitude of matrix)

$$\mathbb{R}^1 \quad \Sigma = \sigma^2 \quad |\Sigma| = \sigma^2$$

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

\mathbb{R}^n , Diagonal Only

$$\Sigma = \text{diag}(\frac{1}{\sigma_i^2}) \quad (\Sigma) = \prod_{i=1}^n \sigma_i^2$$

$$(\Sigma) = \prod_{i=1}^n \sigma_i^2$$

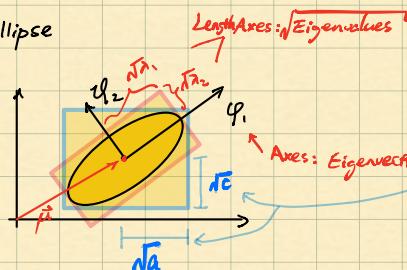
$$p(\vec{x}) = \prod_{i=1}^n \frac{1}{(2\pi)^{\frac{1}{2}} \sigma_i} e^{\left\{ -\frac{1}{2} [(x_i - \mu_i)^2 / \sigma_i^2] \right\}}$$

$\mathbb{R}^2 \neq$

$$\Sigma = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$$

$$\text{Equiprobability Contour: } p(\vec{x}) = (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu}) = C$$

$\rightarrow \mathbb{R}^2$: Ellipse



$$\left\{ -\frac{1}{2} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu}) \right\}$$

$$(A - \lambda I) \vec{v} = 0$$

Eigenvalue & Eigenvectors

$$\textcircled{1} \quad \det(S - \lambda I) = 0$$

$$\hookrightarrow |\Sigma - \lambda I| = 0$$

$$\Downarrow$$

$$\lambda_1, \lambda_2$$

$$\Downarrow$$

$$\textcircled{2} \quad S \vec{v} = \lambda_i \vec{v}, \forall i = \{1, 2\}$$

$$\Downarrow$$

$$\varphi_1, \varphi_2$$

§ 4/5 - Classification

Patterns

$$\vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \text{ vectors in feature space } \mathbb{R}^n$$

Classification : • Compare two patterns \vec{x} & \vec{z} \Rightarrow how similar their attributes are
 pattern & prototype $\vec{z}_k \Rightarrow$ how similar it is with a class

① MED (Minimum Euclidean Distance)

$$\vec{x} \in C_k, \text{ if } d_E(\vec{x}, \vec{z}_k) < d_E(\vec{x}, \vec{z}_l), \forall l \neq k$$

$$\text{where } d_E(\vec{x}, \vec{z}_k) = \|\vec{x} - \vec{z}_k\|_2 = \sqrt{(\vec{x} - \vec{z}_k)^T (\vec{x} - \vec{z}_k)}$$

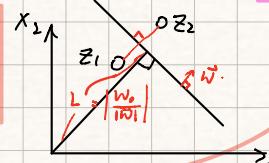
$$\text{• Criteria: } -\vec{z}_1^T \vec{x} + \frac{1}{2} \vec{z}_1^T \vec{z}_1 < -\vec{z}_2^T \vec{x} + \frac{1}{2} \vec{z}_2^T \vec{z}_2$$

$$\text{• Discriminant func: } g_k(\vec{x}) = -\vec{z}_k^T \vec{x} + \frac{1}{2} \vec{z}_k^T \vec{z}_k$$

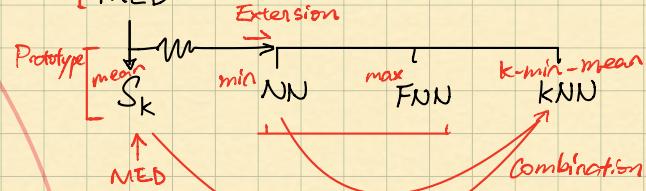
$$\text{• MED: } \vec{x} \in C_k, \text{ if } g_k(\vec{x}) < g_l(\vec{x}), \forall l \neq k$$

• Decision Boundary:

$$g(\vec{x}) = \langle \vec{w}, \vec{x} \rangle + w_0 = 0.$$



Bases [-MED]



Prototype Selection:

$$\text{• Sample Mean: } \vec{z}_k(x) = \frac{1}{N_k} \sum_{i=1}^{N_k} \vec{x}_i$$

$$\text{• Nearest Neighbors: } \vec{z}_k(x) = \vec{x}_k$$

$$\text{s.t. } d_E(\vec{x}, \vec{x}_k) = \min_i d_E(\vec{x}, \vec{x}_i) \quad \forall \vec{x}_i \in C_k$$

• Less sensitive to Noise & Outliers 😊

• Poor @ handling long, thin, dendrite-like clusters 😞 (non-linearly separable clusters)

• Better @

(non-linearly separable) 😊

• More sensitive to noise & outliers 😞

• Computationally complex

• Hard to store, Privacy/Priority

• FNN (Farthest Neighbor): $\vec{z}_k(x) = \vec{x}_k$ choose max dist in the class 'k' as prototype ' \vec{x}_k ' for class 'k'

$$\text{s.t. } d_E(\vec{x}, \vec{x}_k) = \max_i d_E(\vec{x}, \vec{x}_i) \quad \forall \vec{x}_i \in C_k$$

⇒ More Compact Clusters 😊

😊

② kNN (K-Nearest Neighbor)

• For given \vec{x} to classify,

compute $d_E(\vec{x}, \vec{x}_k)$, where \vec{x}_k is the mean of 'k' nearest points in that class.

• Less sensitive to noise & outliers. 😊

• Better at long, thin, dendrite-like, non-linearly separable clusters.

😊 — Computationally Complex

(5 Part 2) →

Weighted Euclidean Distance Metric (WED)

$$d_{W_0}(\vec{x}, \vec{z}) = \left[(\vec{x} - \vec{z})^T W^T W (\vec{x} - \vec{z}) \right]^{\frac{1}{2}}$$

$$\text{where: } W = \begin{bmatrix} w_{11} & \cdots & w_1 \\ \vdots & \ddots & \vdots \\ w_{n1} & \cdots & w_{nn} \end{bmatrix}$$

Scaling & rotation

Assume $W_0 = \text{diag. } \mathbf{I}$ (No rotation)

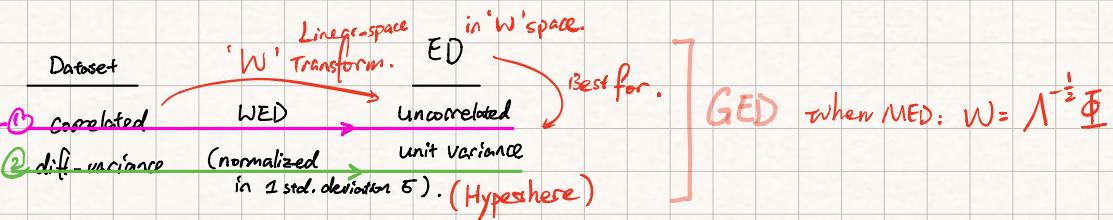
$$d_{W_0}(\vec{x}, \vec{z}) = d_E(\vec{x}', \vec{z}')$$

Transformed Space!!

$$\text{, where } \vec{x}' = W_0 \vec{x}, \vec{z}' = W_0 \vec{z}$$

$$W_0 = \text{diag}([w_{11}, \dots, w_{nn}]) = \begin{bmatrix} w_{11} & \cdots & 0 & 0 \\ \vdots & \ddots & \vdots & 0 \\ 0 & \cdots & 0 & w_{nn} \end{bmatrix}$$

Linear space ⇒ Quadratic Space



→ Orthonormal Covariance Transforms. ⇒ ① Uncorrelated Features. (Rotate Coordinate Axes using 'Eigen.')

find transform A ,

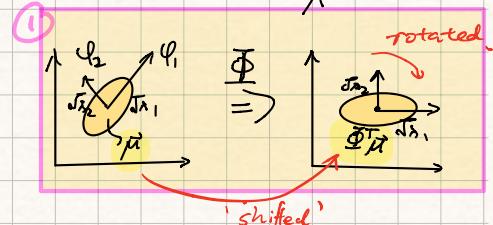
$$\text{s.t. } A \Sigma A^T = \Lambda \quad (\text{that diagonalize } \Sigma)$$

↑ diag-covariance matrix.

where:

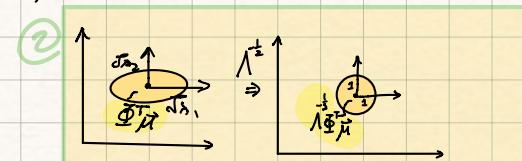
$$A = [\vec{\Phi}_1^T \dots \vec{\Phi}_n^T] = \Phi^T \quad \Lambda = [\lambda_1 \dots 0 \quad 0 \dots \lambda_n] = \text{diag}(\vec{\lambda}_1, \dots, \vec{\lambda}_n)$$

↑ to get
diagonal covariance matrix
 Λ'



→ Whitening Transformation ⇒ ② Unity Variance (Normalize)

$$\Lambda^{-\frac{1}{2}} = \begin{bmatrix} \lambda_1^{-\frac{1}{2}} & \dots & 0 \\ 0 & \ddots & \vdots \\ 0 & \dots & \lambda_n^{-\frac{1}{2}} \end{bmatrix}$$



→ Generalized Euclidean Distance Metrics (GED)

$$W = \Lambda^{-\frac{1}{2}} \Phi^T, W \in \mathbb{R}^{n \times n} \text{ in MED.}$$

GED

$$d_G(\vec{x}, \vec{z}) = \left[(\vec{x} - \vec{z})^T (\Sigma^{-1}) (\vec{x} - \vec{z}) \right]^{\frac{1}{2}}$$

$$\text{where } \Sigma = \Phi \Lambda^{-1} \Phi^T$$

→ Minimum Intra-Class Distance Classifier (MICD)

$$\vec{x} \in A \text{ if } (\vec{x} - \vec{m}_A)^T S_A^{-1} (\vec{x} - \vec{m}_A) < (\vec{x} - \vec{m}_B)^T S_B^{-1} (\vec{x} - \vec{m}_B)$$

MICD Metric: GED metric

intra distance measurement ⇒ compare.

(normalize space into std. deviation)

Decision Boundaries: Non-linear (Quadratic Surface)

- Lies on intersection of equidistance contours around classes.

$$(\vec{x} - \vec{m}_A)^T S_A^{-1} (\vec{x} - \vec{m}_A) = C$$

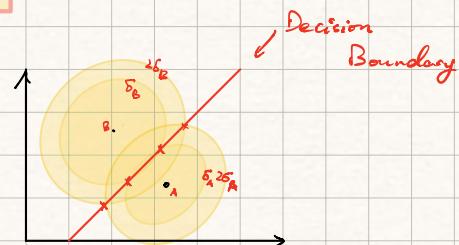
- Quadratic Surface in Hyperspace:

$$\vec{x}^T Q_0 \vec{x} + Q_1 \vec{x} + Q_2 = 0$$

$$\text{where } Q_0 = S_A^{-1} - S_B^{-1}$$

$$Q_1 = 2[\vec{m}_B^T S_B^{-1} - \vec{m}_A^T S_A^{-1}]$$

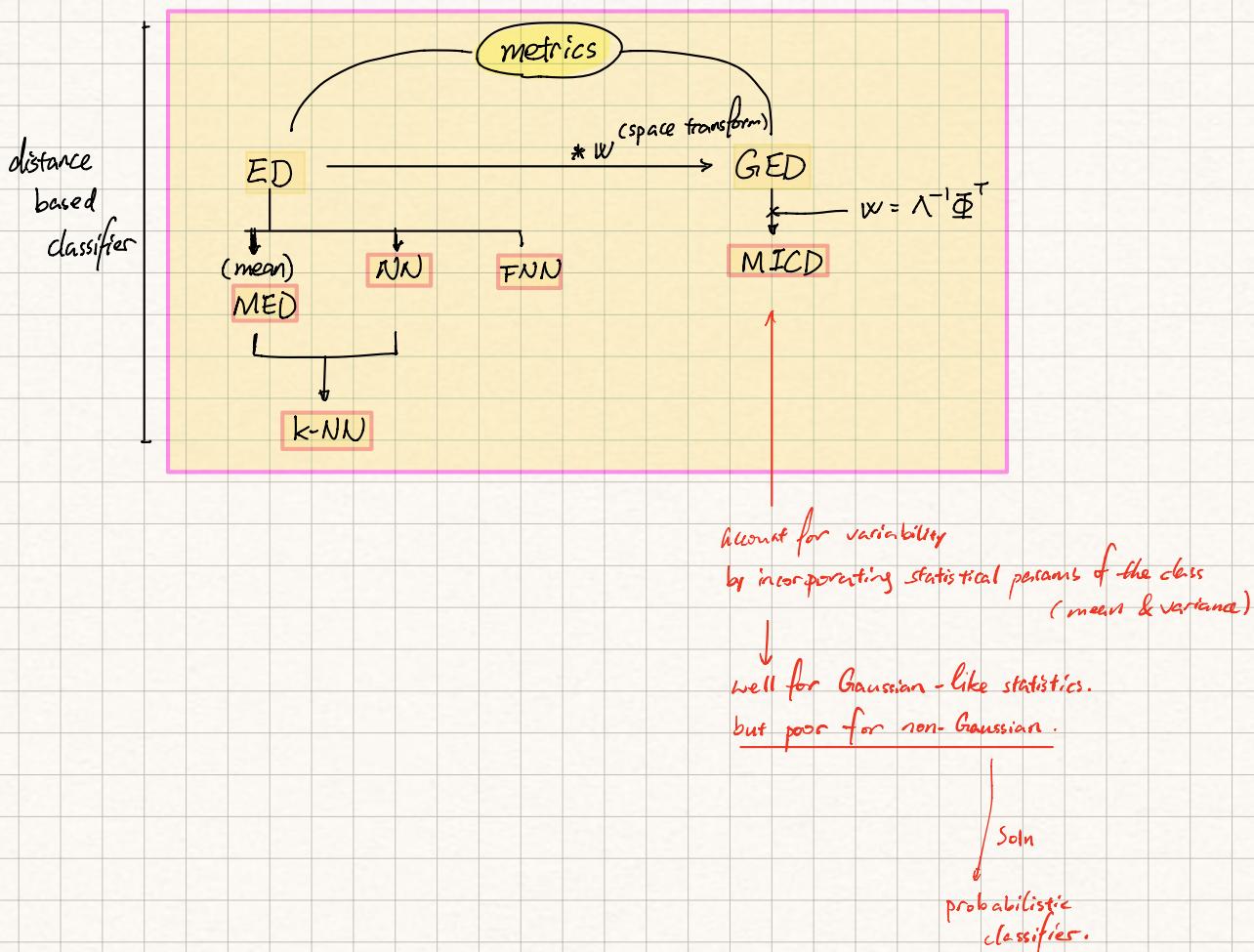
$$Q_2 = \vec{m}_A^T S_A^{-1} \vec{m}_A - \vec{m}_B^T S_B^{-1} \vec{m}_B.$$



- Issue: in favor of the class with largest variance, regardless of x !

- Pro: Lower sensitivity to noise & outliers
 Great handling class distribution & Gaussian models. (Gaussian-like)
- Con: Poor @ handling more complex class distribution.


↳ (Does not take into account of probabilistic info of classes.)



§ 6/7 - Probabilistic Classification

Idea: Given known "class conditional probability density distribution", we can create powerful similarity measures that tell us the **Likelihood / Probability** of each class given an observed pattern.

Basis: Optimal in the minimum "**probability of error**" sense.

Bayesian Classifiers

- Let L_{ij} : cost of deciding on class ' c_j ' when the true class is ' c_i '

- Total Risk with deciding \vec{x} belongs to c_j :

$$\text{The expected cost: } \hat{r}_j(\vec{x}) = \sum_{i=1}^K L_{ij} P(c_i | \vec{x})$$

Total Risk # of classes
 Probability I'm wrong
 (Posterior distribution of class ' c_j ' given pattern ' \vec{x} ')

- \mathbb{R}^2 case:

$$\begin{aligned} \text{Total Cost: } r_1(\vec{x}) &= L_{11} P(c_1 | \vec{x}) + L_{21} P(c_2 | \vec{x}) = \frac{L_{11} P(\vec{x} | c_1) P(c_1)}{p(\vec{x})} + \frac{L_{21} P(\vec{x} | c_2) P(c_2)}{p(\vec{x})} \\ r_2(\vec{x}) &= L_{12} P(c_1 | \vec{x}) + L_{22} P(c_2 | \vec{x}) = \frac{L_{12} P(\vec{x} | c_1) P(c_1)}{p(\vec{x})} + \frac{L_{22} P(\vec{x} | c_2) P(c_2)}{p(\vec{x})} \end{aligned}$$

Bayes' rule

K-classes Bayesian Classifiers:

$$\vec{x} \in C_i \text{ if } \hat{r}_i(\vec{x}) < \hat{r}_j(\vec{x}) \quad \forall j \neq i$$

- \mathbb{R}^2 case: $(L_{12} - L_{11}) P(\vec{x} | c_1) P(c_1) \stackrel{1}{>} (L_{21} - L_{22}) P(\vec{x} | c_2) P(c_2)$

Cost Function

- "zero-one" loss (most common):

$$L_{ij} = \begin{cases} 0 & i=j \\ 1 & i \neq j \end{cases}$$

⇒ Total Risk Func:

$$\hat{r}_j(\vec{x}) = \sum_{\substack{i=1 \\ i \neq j}}^K P(c_i | \vec{x}) = P(E | \vec{x})$$

cond. prob.

↳ Meaning: All errors have equal costs

↳ $\min \{\hat{r}_j(\vec{x})\} \equiv \min \{\text{prob. of errors}\}$ (total risk)

- Types:
 - MAP (Maximum A Posterior)
 - ML (Maximum Likelihood)

⇒ choose most probable class

probability classifiers.

⇒ choose the class that makes observed pattern \vec{x} most probable.

MAP Classifiers

$$P(A | \vec{x}) \stackrel{A}{>} P(B | \vec{x}) \Rightarrow \text{All patterns with higher probability for 'A' than for 'B' will be classified as 'A'}$$

$$\therefore \frac{P(\vec{x}|A)}{P(\vec{x}|B)} \stackrel{A}{>} \frac{P(B)}{P(A)} \quad \text{Bayes.}$$

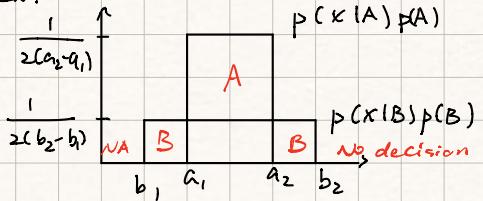
$$\therefore l(x) \stackrel{A}{>} \theta \quad \Rightarrow \quad \log(l(x)) \stackrel{A}{>} \log \theta \quad (\text{log-likelihood})$$

↑ likelihood Ratio ↑ threshold.

when dealing with PDF with exponential dependence (Gamma/Beta)

reduce computation 'e'

* Ex:



Often, $P(A)$ & $P(B)$ priors are unknown.
hence impossible to use the posterior $p(A|\vec{x})$ & $p(B|\vec{x})$.

ML Classifies

Special case of MAP when $P(A) = P(B)$

$$p(\vec{x}|A) \stackrel{A}{>} p(\vec{x}|B) \Rightarrow \frac{p(\vec{x}|A)}{p(\vec{x}|B)} \stackrel{A}{>} 1$$

(§ 7 Part II)

MAP for Normal Distribution

Gaussian: $\begin{cases} p(x|A) = N(\mu_A, \sigma_A^2) = \frac{1}{\sqrt{2\pi}\sigma_A} e^{-\frac{1}{2}(\frac{x-\mu_A}{\sigma_A})^2} \\ p(x|B) = N(\mu_B, \sigma_B^2) \end{cases}$

$$\text{MAP: } \frac{N(\mu_A, \sigma_A^2)}{N(\mu_B, \sigma_B^2)} \stackrel{A}{<} \frac{p(B)}{p(A)}$$

$$\Rightarrow \frac{\exp\{-\frac{1}{2}(\frac{x-\mu_A}{\sigma_A})^2\}}{\exp\{-\frac{1}{2}(\frac{x-\mu_B}{\sigma_B})^2\}} \stackrel{A}{<} \frac{p(B)\sigma_A}{p(A)\sigma_B}$$

$$\text{Log-likelihood} \Rightarrow \left\{ -\frac{1}{2} \left(\frac{x-\mu_A}{\sigma_A} \right)^2 \right\} - \left\{ -\frac{1}{2} \left(\frac{x-\mu_B}{\sigma_B} \right)^2 \right\} \stackrel{A}{<} \ln[\sigma_A p(B)] - \ln[\sigma_B p(A)]$$

∴

$$\text{MAP (Normal)} \quad \left[\left(\frac{x-\mu_B}{\sigma_B} \right)^2 \right] - \left[\left(\frac{x-\mu_A}{\sigma_A} \right)^2 \right] \stackrel{A}{<} 2 \{ \ln[\sigma_A p(B)] - \ln[\sigma_B p(A)] \}$$

' = '

solve for decision boundary

* Case: $\sigma_A = \sigma_B$, $P(A) = P(B) = \frac{1}{2}$.

$$x = \frac{(\mu_B^2 \sigma_A^2 - \mu_A^2 \sigma_A^2)}{2(\mu_B \sigma_A^2 - \mu_A \sigma_A^2)} = \frac{(\mu_B^2 - \mu_A^2)}{2(\mu_B - \mu_A)} = \frac{(\mu_B + \mu_A)}{2}$$

For equal likely, equal-variance classes,

MAP \rightarrow (threshold midway between the means)

$\rightarrow \mathbb{R}^n$ case: $p(\vec{x}|A) = \mathcal{N}(\vec{\mu}_A, \Sigma_A^2)$ $\Rightarrow \frac{\exp[-\frac{1}{2}(\vec{x}-\vec{\mu}_A)^T \Sigma_A^{-1} (\vec{x}-\vec{\mu}_A)]}{\exp[-\frac{1}{2}(\vec{x}-\vec{\mu}_B)^T \Sigma_B^{-1} (\vec{x}-\vec{\mu}_B)]} \stackrel{A}{>} \frac{|\Sigma_A|^{\frac{1}{2}} P(A)}{|\Sigma_B|^{\frac{1}{2}} P(B)}$

$p(\vec{x}|B) = \mathcal{N}(\vec{\mu}_B, \Sigma_B^2)$

$\left(\begin{array}{c} \text{Log} \\ \hline \end{array} \right)$

Bias decision in favor based on prior probabilities

Bias on smaller volume

MAP $(\vec{x}-\vec{\mu}_B)^T \Sigma_B^{-1} (\vec{x}-\vec{\mu}_B) - (\vec{x}-\vec{\mu}_A)^T \Sigma_A^{-1} (\vec{x}-\vec{\mu}_A) \stackrel{A}{<} \stackrel{B}{<} 2 \ln \left[\frac{P(B)}{P(A)} \right] + \ln \left[\frac{|\Sigma_A|}{|\Sigma_B|} \right]$

Decision Boundary: $\vec{x}^T Q_0 \vec{x} + Q_1 \vec{x} + Q_2 + 2Q_3 + Q_4 = 0$

$Q_0 = \Sigma_A^{-1} - \Sigma_B^{-1}$
 $Q_1 = 2 [\vec{\mu}_B^T \Sigma_B^{-1} - \vec{\mu}_A^T \Sigma_A^{-1}]$
 $Q_2 = \vec{\mu}_A^T \Sigma_A^{-1} \vec{\mu}_A - \vec{\mu}_B^T \Sigma_B^{-1} \vec{\mu}_B$
 $Q_3 = \ln \left[\frac{P(B)}{P(A)} \right]$
 $Q_4 = \ln \left[\frac{|\Sigma_A|}{|\Sigma_B|} \right]$

Same as MICD RH.

$d_{\text{MICD}}^2(\vec{x}, \vec{\mu}_B, \Sigma_B) - d_{\text{MICD}}^2(\vec{x}, \vec{\mu}_A, \Sigma_A)$

$\stackrel{A}{<} \stackrel{B}{<} 2 \ln \left[\frac{P(B)}{P(A)} \right] + \ln \left[\frac{|\Sigma_A|}{|\Sigma_B|} \right]$

\rightarrow Relationships

when normal distribution

MAP: $(\vec{x}-\vec{\mu}_B)^T \Sigma_B^{-1} (\vec{x}-\vec{\mu}_B) - (\vec{x}-\vec{\mu}_A)^T \Sigma_A^{-1} (\vec{x}-\vec{\mu}_A) \stackrel{A}{>} \stackrel{B}{<} 2 \ln \left[\frac{P(B)}{P(A)} \right] + \ln \left[\frac{|\Sigma_A|}{|\Sigma_B|} \right]$

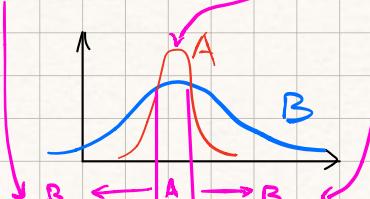
ML: $(\vec{x}-\vec{\mu}_B)^T \Sigma_B^{-1} (\vec{x}-\vec{\mu}_B) - (\vec{x}-\vec{\mu}_A)^T \Sigma_A^{-1} (\vec{x}-\vec{\mu}_A) \stackrel{A}{>} \stackrel{B}{<} \ln \left[\frac{|\Sigma_A|}{|\Sigma_B|} \right]$

MICD: $(\vec{x}-\vec{\mu}_B)^T \Sigma_B^{-1} (\vec{x}-\vec{\mu}_B) - (\vec{x}-\vec{\mu}_A)^T \Sigma_A^{-1} (\vec{x}-\vec{\mu}_A) \stackrel{A}{>} \stackrel{B}{<} 0$

MED: $(\vec{x}-\vec{\mu}_B)^T (\vec{x}-\vec{\mu}_B) - (\vec{x}-\vec{\mu}_A)^T (\vec{x}-\vec{\mu}_A) \stackrel{A}{>} \stackrel{B}{<} 0$

$\left\{ \begin{array}{l} \text{MAP} = \text{MICD} , \text{ iff } \left\{ \begin{array}{l} \text{normal distribution: } \mathcal{N}(\vec{\mu}, \Sigma) \\ \text{equal prior prob. } (P(A) = P(B)) \\ \text{equal volume } (|\Sigma_A| = |\Sigma_B|) \end{array} \right\} \end{array} \right. \quad \ln(1) = 0$

$\left\{ \begin{array}{l} \text{MICD in favour of class with largest variance regardless of } \vec{x} . \\ \text{MED in favor of class with lowest variance close to mean, highest variance beyond a certain point in both class.} \end{array} \right. \quad \begin{array}{l} \text{Fix this issue.} \\ \text{by taking account of the probability.} \end{array}$



Performance of Bayes Classifier

Can compute
Theoretical Limit
without data.

Probability of Error: $P(E|x) = \min [P(A|x), P(B|x)]$

We choose maximum posterior as our class.
 \therefore error would be the minimum posterior.

Expected probability of error: $P(\varepsilon) = \int P(\varepsilon|x) p(x) dx$

$$= \int \min [P(x|A) P_A, P(x|B) P_B] dx.$$

Hard

define decision Regions R_A & R_B :

$$R_A = \{x \text{ s.t. } P(A|x) > P(B|x)\}$$

$$R_B = \{x \text{ s.t. } P(B|x) > P(A|x)\}$$

$$\therefore P(\varepsilon) = \int_{R_A} P(x|B) P_B dx + \int_{R_B} P(x|A) P_A dx$$

Ex. univariate Normal, equal variance, equally likely, two class

- $\times (1D)$
- $n=1, P(A)=P(B)=0.5, \mu_A=\mu_B=6, \sigma^2 = \sigma_A^2 = \sigma_B^2$
 - Likelihood:

$$\begin{cases} p(x|A) = N(\mu_A, \sigma_A^2) \\ p(x|B) = N(\mu_B, \sigma_B^2) \end{cases}$$

• find $p(\varepsilon)$

• Recall this case: MAP $\xrightarrow[\text{Boundary}]{\text{reduce}} x = \frac{\mu_B + \mu_A}{2}$

$$\cdot \mu_A < \mu_B \Rightarrow R_A = \{x \text{ s.t. } x < \frac{\mu_B + \mu_A}{2}\}$$

$$R_B = \{x \text{ s.t. } x > \frac{\mu_B + \mu_A}{2}\}$$

$$\therefore P(\varepsilon) = \int_{R_A} P(x|B) P(x|B) dx + \int_{R_B} P(x|A) P(x|A) dx$$

$$= \frac{1}{2} \int_{-\infty}^{\frac{\mu_B + \mu_A}{2}} N(\mu_B, \sigma^2) dx + \frac{1}{2} \int_{\frac{\mu_B + \mu_A}{2}}^{\infty} N(\mu_A, \sigma^2) dx$$

\therefore symmetric class. $\Rightarrow P(\varepsilon|A) = P(\varepsilon|B)$

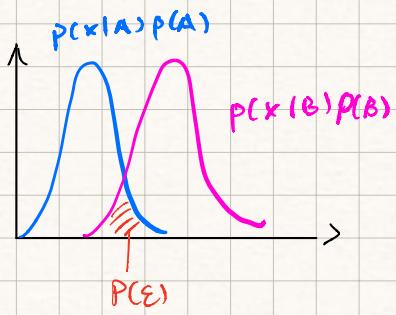
$$\therefore P(\varepsilon) = \int_{\frac{\mu_B + \mu_A}{2}}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2}\left(\frac{x-\mu_A}{\sigma}\right)^2\right\} dx.$$

\Rightarrow Let $y = \frac{x-\mu_A}{\sigma}$, $dx = \sigma dy$ \int Normalize to $N(0, 1) \Rightarrow$ for lookup table. (Q func.)

$$P(\varepsilon) = \int_{\frac{(\mu_B - \mu_A)}{\sigma}}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}y^2\right] dy$$

$$\hookrightarrow Q(Q) = \int_{\alpha}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}y^2\right] dy$$

$$\therefore P(\varepsilon) = Q\left(\frac{\mu_B - \mu_A}{\sigma}\right) \text{ (lookup table)}$$

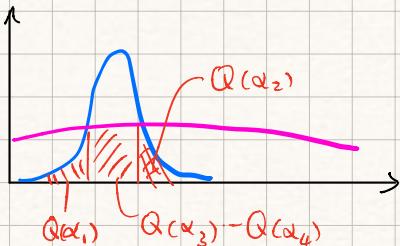


Observation:

- as dist. btwn means \uparrow
shaded area($P(E)$) monotonically \downarrow
- At $\alpha = 0$, $\mu_A = \mu_B = 0$ & $P(E) = \frac{1}{2}$
- $\lim_{\alpha \rightarrow \infty} P(E) = 0$.

For case $p(A) \neq p(B) / \sigma_A \neq \sigma_B$,

(1D) Decision boundary change AND additional boundary is introduced.



$$P(E) = p(A)Q(\alpha_1) + p(B)[Q(\alpha_3) - Q(\alpha_4)] + p(A)Q(\alpha_2)$$

n-D case (Multivariate, $n \geq 1$)

$$\begin{cases} p(\vec{x}|A) = \mathcal{N}(\vec{\mu}_A, \Sigma) \\ p(\vec{x}|B) = \mathcal{N}(\vec{\mu}_B, \Sigma) \\ p(A) = p(B) \end{cases}$$

$$\hookrightarrow P(E) = Q(d_M(\vec{\mu}_A, \vec{\mu}_B)/2)$$

where $d_M(\vec{\mu}_A, \vec{\mu}_B)$ is Mahalanobis Dist. btwn classes.

$$d_M(\vec{\mu}_A, \vec{\mu}_B) = [(\vec{\mu}_A - \vec{\mu}_B)^T \Sigma^{-1} (\vec{\mu}_A - \vec{\mu}_B)]^{1/2}$$

Explanation:

• If cases where $\Sigma_A = \Sigma_B$, $p(A) = p(B)$,

the decision boundary btwn classes is ALWAYS **straight line** in hyperspace.

GED

- slope based on Σ (since orthonormal whitening ($\Lambda^{-1}\Phi^T$) is identical)
- intersects with the midpoint of the line segment btwn ($\vec{\mu}_A$ & $\vec{\mu}_B$)

- The probability of err. \equiv the area under $p(\vec{x}|A)p(A)$ on the class B side of this DBnd.
+ the area under $p(\vec{x}|B)p(B)$ on the class A side of the DBnd.

Decision Boundary

Non-Gaussian PDF:

Ex! Suppose two classes have density functions

& a priori probabilities:

$$p(x|C_1) = \begin{cases} ce^{-\lambda x}, & 0 \leq x \leq 1 \\ 0, & \text{else} \end{cases}$$

$$p(x|C_2) = \begin{cases} ce^{-\lambda(1-x)}, & 0 \leq x \leq 1 \\ 0, & \text{else} \end{cases}$$

where $c = \frac{\lambda}{1-e^{-\lambda}}$ = normalizing the PDF.

$$P(C_1) = P(C_2) = \frac{1}{2}$$

∴ Expected Prob of error:

$$\begin{aligned} P(E) &= \int \min [P(\bar{x}|C_1), P(\bar{x}|C_2) P(C_2)] dx \\ &= \int_{R_{C_1}} P(\bar{x}|C_1) P(C_1) dx + \int_{R_{C_2}} P(\bar{x}|C_2) P(C_2) dx \\ &= \int_0^{0.5} 0.5 P(\bar{x}|C_1) dx + \int_{0.5}^1 0.5 P(\bar{x}|C_2) dx \\ &\because \text{Symmetry} \\ &\quad \hookrightarrow = \int_{0.5}^{1.0} ce^{-\lambda x} dx \\ P(E|C_1) &= P(E|C_2) \\ &= \frac{c}{\lambda} [e^{-\lambda/2} - e^{-\lambda}] \end{aligned}$$

$$\therefore P(E|x) = \begin{cases} P(C_2|\bar{x}) & 0 \leq x \leq \frac{1}{2} \\ P(C_1|\bar{x}) & \frac{1}{2} \leq x \leq 1 \end{cases} \Rightarrow \begin{cases} \frac{P(x|C_2) P(C_2)}{P(x)} & 0 \leq x \leq \frac{1}{2} \\ \frac{P(x|C_1) P(C_1)}{P(x)} & \frac{1}{2} \leq x \leq 1 \end{cases} \Rightarrow \begin{cases} \frac{e^{-\lambda x} 0.5}{e^{-\lambda x} + e^{\lambda(1-x)}} & 0 \leq x \leq \frac{1}{2} \\ \frac{e^{-\lambda(1-x)} 0.5}{e^{-\lambda x} + e^{\lambda(1-x)}} & \frac{1}{2} \leq x \leq 1 \end{cases}$$

* Error Bound

- * In practice, the exact $P(E)$ is only easy to compute for simple cases as shown before
- * We may determine Bounds on $P(E)$ instead of exact $P(E)$
 - Easier to compute
 - Leads to estimates of classifier performance.

* Bhattacharyya Bound

Derivation:

$$\text{Approximate } \min[a, b] \leq \sqrt{ab}$$

∴

$$P(E) = \int \min [P(\bar{x}|A) P(A), P(\bar{x}|B) P(B)] dx$$

∴

$$P(E) \leq \sqrt{P(A)P(B)} \int \sqrt{P(\bar{x}|A)P(\bar{x}|B)} dx$$

$$\therefore P(A) + P(B) = 1$$

Let define Bhattacharyya Coeff as

$$\rho = \int \sqrt{P(\bar{x}|A)P(\bar{x}|B)} dx$$

∴ Upper Bound

(Bhattacharyya Bound) of $P(E)$: $P(E) \leq \frac{1}{2}\rho$

max possible error:

$$\text{when } P(A) = P(B) = \frac{1}{2}$$

$$\Rightarrow \sqrt{P(A)P(B)} = \frac{1}{2}$$

§ 8/9 - Estimation & Learning

Motivation

Recall: C Bayesian Classifiers achieves Minimum Probability of Error

\Leftrightarrow Better than Naïve Bayes classifier for situations,

where the class conditional PDFs ' $p(\vec{x} | C_i)$ ' are known!

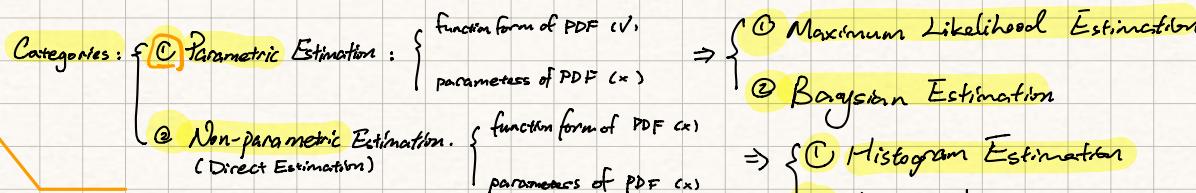
Idea: Learn Empirical PDFs \Rightarrow to apply Bayesian Classifiers.

But normally unknown

Remark: ① Bayesian Classification is optimal

② Empirical PDFs result in sub-optimal classifiers.

③ Performance to the theoretical limit (minimum P(E)) depends on accuracy of estimated PDF.



① Parametric Learning

$$\left\{ \begin{array}{l} \text{function form of PDF } p(\vec{x}) \\ \text{parameters of PDF } (\mu, \Sigma) \end{array} \right. \quad \text{Ex: Gaussian: } N(\vec{\mu}_A, \Sigma_A) \quad \left\{ \begin{array}{l} \vec{\mu}_A = ? \quad \Sigma_A = ? \quad (\mu, \Sigma) \end{array} \right.$$

Two Approaches

- ① Maximum Likelihood Estimation
- ② Bayesian (Maximum Posterior) Estimation

Params fixed (\checkmark)
 Params unknown quantity (x) \Rightarrow Goal: find estimate values
 s.t. max probability that given samples
 & the resulting PDF: $p(\vec{x}|\theta)$

L: Params as random variables
 L: known prior distribution (\checkmark)

Goal: obtain a posterior distribution $p(\theta|\vec{x})$
 which indicates the estimate value based on the given samples.

Maximum Likelihood Estimation

Given: ① set of samples: $\{\vec{x}_i\}$

② PDF form: $p(\vec{x})$ Ex: Gaussian.

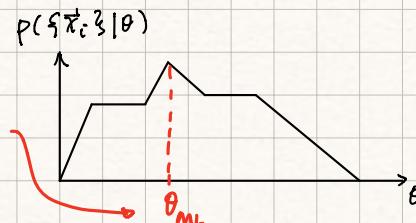
Goal: obtain parameters $\theta \Rightarrow$ PDF

Ex: Given @ PDF: $p(\vec{x}|A) = N(\vec{\mu}_A, \Sigma_A)$

\Rightarrow Goal: $\theta = (\vec{\mu}_A, \Sigma_A)$

Formulation:

$$\hat{\theta}_{ML} = \underset{\theta}{\operatorname{argmax}} [p(\{\vec{x}_i\} | \theta)]$$



Verbal: Given \vec{x}_i (observation),

the maximum estimate of parameter θ

is chosen to be that value which maximizes the PDF: $p(\vec{x}_i | \theta)$

Derivation: Assume I.I.D. condition. (samples are independent of each other) \rightarrow

$$\hookrightarrow p(\{\vec{x}_i\} | \theta) = p(\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N | \theta) = \prod_{i=1}^N p(\vec{x}_i | \theta)$$

Joint Distribution
Likelihood prob.

To: $\max \{ p(\{\vec{x}_i\} | \theta) \}$, take derivative & set to 0

$$\hookrightarrow \frac{\partial}{\partial \theta} p(\{\vec{x}_i\} | \theta) \Big|_{\theta = \hat{\theta}_{ML}} = 0$$

product of individual sample prob.

• Log Likelihood:

$$l(\theta) = \log \{ p(x_i | \theta) \} = \sum_{i=1}^N \log p(x_i | \theta)$$

∴ Maximum Likelihood Condition:

$$\frac{\partial l(\theta)}{\partial \theta} \Big|_{\theta=\hat{\theta}_{ML}} = 0$$

Estimation Bias

* Definition:

An estimate $\hat{\theta}$ is unbiased

if its expected value is equal to the true value: $E[\hat{\theta}] = \theta$

* Ex] ML estimate of the mean is Unbiased:

$$\begin{aligned} E[\hat{\mu}_{ML}] &= E\left[\frac{1}{N} \sum_{i=1}^N \vec{x}_i\right] \\ &= \frac{1}{N} \sum_{i=1}^N E[\vec{x}_i] \end{aligned}$$

$$\because \vec{\mu} = E[\vec{x}]$$

$$\therefore = \frac{1}{N} \sum_{i=1}^N \vec{\mu}$$

$$= \vec{\mu}$$

$$\therefore E[\hat{\mu}_{ML}] = \vec{\mu} \Rightarrow \text{Unbiased } \circlearrowleft$$

* Ex] ML estimate of covariance matrix is Biased:

$$\begin{aligned} E[\hat{\Sigma}_{ML}] &= E\left[\frac{1}{N} \sum_{i=1}^N (\vec{x}_i - \hat{\mu}_{ML})(\vec{x}_i - \hat{\mu}_{ML})^T\right] \\ &= \frac{1}{N} \sum_{i=1}^N E[(\vec{x}_i - \hat{\mu}_{ML})(\vec{x}_i - \hat{\mu}_{ML})^T] \\ \hat{\mu}_{ML} &= \frac{1}{N} \sum_{i=1}^N \vec{x}_i \quad \left(\begin{aligned} &= \frac{1}{N} \sum_{i=1}^N E[(\vec{x}_i - \vec{\mu}) - (\hat{\mu}_{ML} - \vec{\mu})] \underbrace{E[(\vec{x}_i - \vec{\mu})]}_{\Sigma} \underbrace{E[(\hat{\mu}_{ML} - \vec{\mu})]}_{\hat{\mu}_{ML} - \vec{\mu}}^T \end{aligned} \right) \\ &= E\left[\frac{1}{N} \sum_{i=1}^N (\vec{x}_i - \vec{\mu})(\vec{x}_i - \vec{\mu})^T\right] - E[(\hat{\mu}_{ML} - \vec{\mu})(\hat{\mu}_{ML} - \vec{\mu})^T] \\ &= \Sigma - E[(\hat{\mu}_{ML} - \vec{\mu})(\hat{\mu}_{ML} - \vec{\mu})^T] \\ &= \Sigma - E\left[(\frac{1}{N} \sum_{i=1}^N \vec{x}_i - \vec{\mu})(\frac{1}{N} \sum_{i=1}^N \vec{x}_i - \vec{\mu})^T\right] \\ &= \Sigma - \frac{1}{N^2} E\left[\sum_{i=1}^N \sum_{j=1}^N (\vec{x}_i - \vec{\mu})(\vec{x}_j - \vec{\mu})^T\right] \end{aligned}$$

$$\because \text{IID} \Rightarrow E[(\vec{x}_i - \vec{\mu})(\vec{x}_j - \vec{\mu})^T] = 0 \quad \forall i \neq j.$$

$$\therefore = \Sigma - \frac{1}{N^2} \sum_{i=1}^N E[(\vec{x}_i - \vec{\mu})(\vec{x}_i - \vec{\mu})^T]$$

$$\therefore \boxed{E[\hat{\Sigma}_{ML}] = \Sigma - \frac{1}{N^2} N \Sigma = \frac{N-1}{N} \Sigma} \Rightarrow \text{Biased } \circlearrowleft$$

↳ As $N \rightarrow \infty$, Bias negligible.

↳ or $(\times \frac{N}{N-1})$ for ML Estimate: $E\left[\frac{N}{N-1} \hat{\Sigma}_{ML}\right] = \Sigma$!

* Bayesian Estimation

- Idea:
 - Instead of MLE treating params fixed finding param θ \Rightarrow maximize samples come from estimated PDF
 - We do BE of treating params random var. with an assumed a priori distribution use the observed samples \Rightarrow obtain posterior distribution which indicates params.

- Procedure:
 - Let $p(\theta)$ be a priori distribution $\{x_i\}$ be the set of samples.

• Posteriori Distribution: $p(\theta | \{x_i\}) = \frac{p(\{x_i\} | \theta) p(\theta)}{p(\{x_i\})}$ from the requirement of PDF:
 scale factor $\Leftarrow \int p(\theta | \{x_i\}) d\theta = 1$.

→ EX / Suppose: $\begin{cases} \text{PDF: } p(\theta) = N(\mu, \sigma^2) & (\checkmark) \\ \mu: (\times) \\ \sigma^2: (\checkmark) = \sigma_0^2 \end{cases}$

↳ Bayesian Estimate of μ ?

- ① Assume a priori PDF for param $\theta = \mu$.

↳ $p(\mu) = N(\mu | \mu_0, \sigma_0^2)$
 ↓ initial guess of the guess.

- ② Given Samples, compute $p(\mu | \{x_i\})$

$$p(\mu | \{x_i\}) = \alpha p(\{x_i\} | \mu) p(\mu)$$

- Assume samples IID.

$$= \alpha \prod_{i=1}^N p(x_i | \mu) p(\mu)$$

where $\alpha = \frac{1}{p(\{x_i\})}$ as scale factor independent of μ .

- o Sub:

$$\begin{aligned} &= \alpha \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2 \right\} \cdot \frac{1}{\sqrt{2\pi}\sigma_0} \exp \left\{ -\frac{1}{2} \left(\frac{\mu - \mu_0}{\sigma_0} \right)^2 \right\} \\ &= \alpha' \exp \left\{ -\frac{1}{2} \left[\sum_{i=1}^N \frac{x_i^2 - 2x_i\mu + \mu^2}{\sigma^2} + \frac{\mu^2 - 2\mu\mu_0 + \mu_0^2}{\sigma_0^2} \right] \right\} \\ &= \alpha'' \exp \left\{ -\frac{1}{2} \left[\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2} \right] \mu^2 - 2 \left[\frac{Nm_N}{\sigma^2} + \frac{\mu_0}{\sigma_0^2} \right] \mu \right\} \end{aligned}$$

where m_N is sample mean \Rightarrow quadratic in μ . \Rightarrow Gaussian Form!

- In short:

As posteriori Density:

$$p(\mu | \{x_i\}) = N(\mu | \mu_N, \sigma_N^2)$$

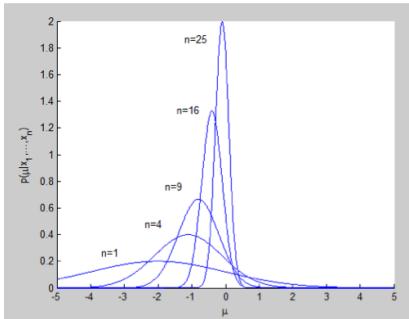
$$\text{where } \mu_N = \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} m_N + \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0$$

∴ Peak Density @ μ_N , with variance of σ_N^2
 ∴ Bayesian estimate of μ is $\hat{\mu}_B = \mu_N$.

* Observations:

- Bayesian estimate can be interpreted as weighted average of initial guess μ_0 & sample mean m_N .

- If $\sigma_0 = 0$, \Rightarrow sure of initial guess that we ignore the samples.
- $\sigma_0 > 0$, \Rightarrow some uncertainty, sample mean μ_N has greater dominance
- If $\sigma_0 \gg \sigma$, \Rightarrow initial uncertainty is relatively large, samples weighted heavily.
- As $N \rightarrow \infty$, $\sigma_N \rightarrow 0$ & $\mu_N \rightarrow \mu_0$
- ↳ # of samples \nearrow ,
 the density narrows & peaks @ true mean during BL !



As measures arrive, the PDF of the estimate narrows, implying that the estimation error is decreasing.

* (Slide 9) Non-parametric

Non-parametric Learning

- Given N samples $\{x_i\}$ with labels $\{y_i\}$.
- estimate distribution directly

- 3 Main categories
- ① Histogram Estimation —> Grp labeled samples into discrete regions to approx. $p(x)$
 - ② Parzen window Estimation —> Approx. $p(x)$ in continuous manner based on the local distribution of each sample, thus controlling resolution along x-axis explicitly, with resolution along PDF axis data dependent.
 - ③ kNN Estimation.
- ↳ Approx. $p(x)$ in continuous manner based on contribution of nearest neighbour samples, \Rightarrow thus controlling res along PDF axis explicitly, with resolution along x-axis data dependent.

Histogram Estimation.

- simplest approach.
- construct normalized histogram.

Derivation / Proof:

- Consider some interval $R = [a, b]$
 - If $p(x)$ is const. over this region,
- $$\Rightarrow P_{R^c}(x \in R) = P_R = \int_a^b p(x) dx = p(a) \cdot (b-a)$$

- Suppose N samples $\{x_1, \dots, x_N\}$ from PDF $p(x)$
- \Rightarrow # of samples M that fall within region R
must obey a Binomial Distribution:

$$P(M) = \binom{N}{M} P_R^M (1-P_R)^{N-M}$$

- Log Form:

$$\log p(M) = \log \left[\frac{N}{M} \right] + M \log [P_R] + (N-M) \log [1-P_R]$$

- Derivative = 0 :

$$\frac{\partial \log p(M)}{\partial P_R} = \frac{M}{P_R} - \frac{N-M}{1-P_R} = 0$$

$$\Rightarrow (1-P_R)M = P_R(N-M)$$

- ML estimate of P_R as:

$$\hat{P}_R = \frac{M}{N}$$

- Recall:

$$P_R = P(a) \cdot c(b-a)$$

- ML estimate:

$$\hat{P}_R(x) = \hat{P}_{(x)} \cdot (b-a)$$

$$\hat{P}_{(x)} = \frac{M}{N|R_i|}$$

\approx size of region R_i

Steps :

- Given a set of bins R_i :

- ① Count # of samples M_i that falls into each bin 'i'
- ② Count # total of samples 'N'
- ③ For a particular pattern x ,

$$\hat{p}(x) = \frac{M_i}{N|R_i|} \quad \forall x \in R_i$$

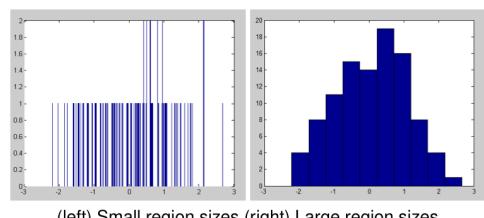
TRADE OFF:

- For good res along x , \Rightarrow small-sized regions \Rightarrow small $M_i \Rightarrow$ poor PDF $p(x)$ resolution.
- For good res along $p(x)$, \Rightarrow large $M_i \Rightarrow$ large-sized region \Rightarrow poor ' x ' resolution.

Disadvantages:

- Estimated PDF is discontinuous.
- Shifting origin changes the shape of the estimated PDF.
- $p(x)$ const. over each region.

Example: 100 samples from a Gaussian distribution



Parzen Windows Estimation

- Advantages:
 - ① No need to predefine region sizes
 - ② Estimated PDF is **ALWAYS** continuous
 - ③ Origin-Independent

- * Idea:
- Every sample x_i locally influences the estimated PDF in the vicinity of x_i .
 - So if observed x_i , PDF can't be too small
if see lots of x_i , PDF in that region larger
 - Up estimated PDF = sum of the contributions.

$$\hat{p}(x) \propto \sum_i \Phi(x - x_i)$$

~ Linear combination of window function (local distributions)

window function. \Rightarrow controls how observed sample x_i influences the PDF.

* Window Function Φ

- ① Must be normalized:

$$\int_{-\infty}^{\infty} \Phi(x) dx = 1$$

- ② Stretch / Compress the window function \Rightarrow change the locality of influence of a sample.

$$\Phi\left(\frac{x - x_i}{h}\right)$$

\ddot{x} scaling factor

- ③ Keep things normalized:

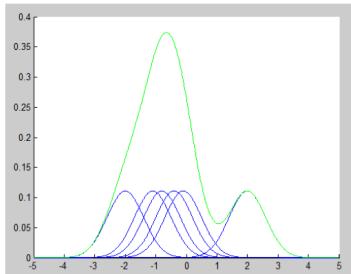
Paren Window Est. $\int_{-\infty}^{\infty} \frac{1}{h} \Phi\left(\frac{x - x_i}{h}\right) dx = 1$

\Rightarrow Given ' N ' samples $\{x_1, \dots, x_N\}$,

$$\hat{p}(x) = \frac{1}{N} \sum_i \frac{1}{h} \Phi\left(\frac{x - x_i}{h}\right)$$

- Common window func: rect, triangular, gaussian, & exponential.

Example: 6 samples, Gaussian window function



* Disadvantages:

- ① Need to choose a window function & scale factor h .

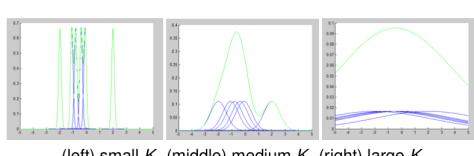
Rule of Thumb:

• One option, try $h = K/\sqrt{N}$ for some const. K .

- ② For small N , const K is important

- if too small, the estimate is noisy with sharp peaks at the samples
- if too large, the estimate is smeared with low resolution.

Example: 6 samples, Gaussian window function, different K



* K-NN Est.

Recall:

Histogram : region size
Parzen Window : window function width } \Rightarrow explicitly ctrl res along x-axis,
res along PDF is data independent.

KNN

- Determine the size of region required at each point to enclose this many samples.

↳

explicitly ctrl's the res along PDF axis.

↳

x-axis res is now data dependent.

Procedure:

- create an interval $[x-\alpha, x+\alpha]$, centred around 'x'
- Increase α till it contains a suitable # of observations M
- Compute estimate of $p(x)$ as:

$$\hat{p}(x) = \frac{M}{N|R(x)|} = \frac{M}{N \cdot 2\alpha}$$

where $R(x)$ is the smallest region,

centred around x , which encloses M sample points.

Remark:

- Freqly, we set $M=\sqrt{N}$, \Rightarrow no free params

↳ sample density high $\Rightarrow |R(x)| \downarrow$, small \Rightarrow high res where it's needed

↳ low $\Rightarrow |R(x)| \uparrow$, P large \Rightarrow low res where probably acceptable in sparse regions.

Pro: ① Avoid setting $p(x)$ identically to zero for regions have no samples.

\Rightarrow results in a more realistic non-zero probability.

Con: ② Estimated PDF is highly "peaked" & non-normalized

§ 10/11 - Discriminant Function

* Motivations:

Recall

- Let $g(\vec{x})$ be discriminant function

s.t. two classes classified by.

$$g(\vec{x}) \begin{cases} \geq 0 \\ < 0 \end{cases}$$

Ideas

Alternative Approach:

- Assume a particular form for the discriminant functions (e.g. hyperplane)
- Use given samples to directly estimate the params of the discriminant functions.
- Given discriminant functions, decision rules & surfaces are defined

Goal: learn discriminant functions directly from the samples.

* Linear Discriminant

$$g(\vec{x}) = \vec{w}^T \vec{x} + w_0 \Rightarrow \text{hyperplane.}$$

$$\begin{aligned} g(\vec{x}) &= w_1 x_1 + w_2 x_2 + \dots + w_n x_n + w_0 \\ &= \sum_{i=1}^n w_i x_i + w_0 \end{aligned}$$

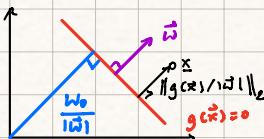
Ex:

Two-class case $\vec{x} = (x_1, x_2)$

$$\Rightarrow g(x_1, x_2) \Rightarrow x_2 = -\frac{w_1}{w_2} x_1 - \frac{w_0}{w_2}$$

Properties:

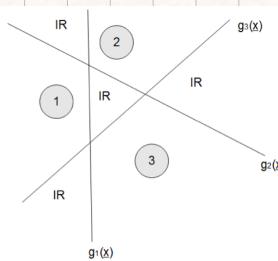
- $\vec{w} \perp$ hyperplane
- $\frac{\vec{w}}{\|\vec{w}\|}$: unit normal of hyperplane.



* Multi-class with Linear Discriminant

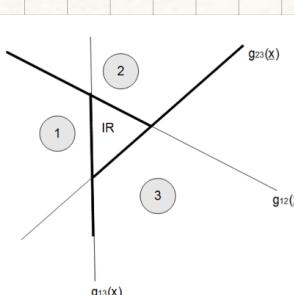
3 Strategies:

- ① one vs. all $g_i(\vec{x}) \forall i \Rightarrow \begin{cases} g_i(\vec{x}) > 0 \text{ if } \vec{x} \in C_i \\ g_i(\vec{x}) < 0 \text{ else.} \end{cases}$



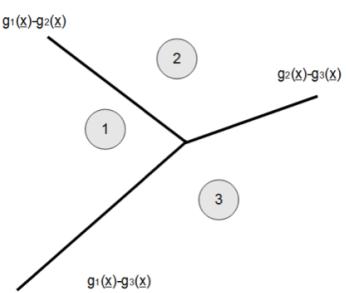
- ② one vs. one (in pair)

$$g_{ij}(\vec{x}) > 0 \quad \forall j \neq i \text{ if } \vec{x} \in C_i$$



- ③ Each class has discriminant func.

$$g_i(\vec{x}) > g_j(\vec{x}) \quad \forall j \neq i \text{ if } \vec{x} \in C_i$$



only this one avoids producing indeterminate regions

3

Generalized:

$$\vec{x} \in C_i \text{ if } g_i(\vec{x}) > g_j(\vec{x}) \quad \forall j \neq i$$

$$\text{with } g_i(\vec{x}) = \vec{w}_i^T \vec{x} + w_{i0}, \quad i=1 \dots k.$$

$$\Downarrow \text{Decision Boundary b/w } C_i \text{ & } C_j : g_i(\vec{x}) = g_j(\vec{x}) =$$

$$g_i(\vec{x}) - g_j(\vec{x}) = (\vec{w}_i - \vec{w}_j)^T \vec{x} + w_{j0} - w_{i0} = 0$$

a hyperplane with normal vector: $(\vec{w}_i - \vec{w}_j)$!

Learning Discriminant

- Assume classes are linearly separable in the original feature space (\vec{x})

$$\Rightarrow \vec{X} = \begin{bmatrix} \vec{x}_1 \\ \vdots \\ \vec{x}_n \end{bmatrix} \quad \& \quad \vec{a} = \begin{bmatrix} w_1 \\ \vdots \\ w_n \\ w_0 \end{bmatrix}$$

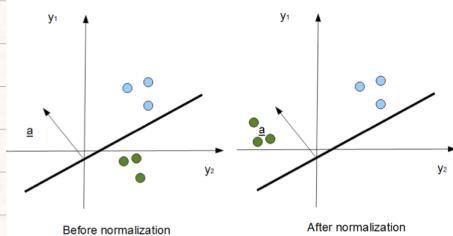
- In \vec{X} space, the decision surface is a hyperplane which contains the original & has normal vector: $\frac{\vec{a}}{\|\vec{a}\|}$

- Given labels $\{\vec{x}_1, \dots, \vec{x}_n\}$,
find \vec{a} s.t. $\vec{a}^T \vec{x}_i > 0 \quad \forall \vec{x}_i \in C_1$
 $\vec{a}^T \vec{x}_i < 0 \quad \forall \vec{x}_i \in C_2$

Simplifications:

perform normalization by replacing all \vec{x}_i by $-\vec{x}_i \quad \forall \vec{x}_i \in C_2$.

$$\Rightarrow \text{find } \vec{a}^T \vec{x}_i > 0 \quad \forall i.$$



$$\times \text{ Gradient Descent} \quad \vec{a}_{k+1} = \vec{a}_k - \rho \nabla J(\vec{a}_k)$$

{ Perception: convergence based on sum of dist. of misclassified samples to DB.

MSE : converges based on sum of squared error.

- { Non-segmental: update based on all sample @ the same time
(offline)
- | Segmental: based on one sample at time.
(online)

Perceptron approach

- Step 1: Set an initial guess for the weight vector (\underline{a}_0) and let $k = 0$
- Step 2: Based on \underline{a}_k , construct the classifier and determine the set of misclassified samples $Y(\underline{a})$. If there are no misclassified samples, stop here since we have arrived at the solution. Otherwise, continue to Step 3.
- Step 3: Compute a scalar multiple of the sum of misclassified samples $\rho_k \sum_{y \in Y(\underline{a})} (y)$
- Step 4: Determine \underline{a}_{k+1} as

$$\underline{a}_{k+1} = \underline{a}_k + \rho_k \sum_{y \in Y(\underline{a})} (y) \quad (30)$$

- Step 5: Go to Step 2.

↳ Variations:

- Fixed-Increment: $\rho_k = 1$
- Variable-Increment: $\rho_k \propto \frac{1}{k} \rightarrow$ avoid over-shooting

Single sample correction: treat samples sequentially,
change weight vector with each misclassification.

Sequential Perceptron approach

- Step 1: Set an initial guess for the weight vector (\underline{a}_0) and let $k = 0$
- Step 2: Based on \underline{a}_k , construct the classifier and determine the set of misclassified samples $Y(\underline{a})$. If there are no misclassified samples, stop here since we have arrived at the solution. Otherwise, continue to Step 3.
- Step 3: Compute a scalar multiple of the k^{th} misclassified sample $\rho_k y^k$
- Step 4: Determine \underline{a}_{k+1} as

$$\underline{a}_{k+1} = \underline{a}_k + \rho_k y^k \quad (31)$$

- Step 5: Go to Step 2.

Minimum squared error approach

- One issue with the perceptron approach is that if the classes are not linearly separable, the learning procedure will never stop since there will always be misclassified samples!
- One way around this is to terminate after a fixed number of iterations, but the resulting weight vector may or may not be appropriate for classification.
- Solution: What if we use a different criterion that will converge even if there are misclassified samples?
- The minimum squared error criterion provides a good compromise in performance for both separable and non-separable problems.

Minimum squared error approach

- Instead of solving a set of inequalities:

$$\underline{a}^T \underline{y}_i > 0, i = 1, \dots, N \quad (43)$$
- we can obtain a solution vector for a set of equations:

$$\underline{a}^T \underline{y}_i = b_i, i = 1, \dots, N \quad (44)$$
- Let the error vector \underline{e} be defined as:

$$\underline{e} = \begin{bmatrix} \underline{y}_1^T \\ \vdots \\ \underline{y}_i^T \\ \vdots \\ \underline{y}_N^T \end{bmatrix} \begin{bmatrix} a_1 \\ \vdots \\ a_i \\ \vdots \\ a_N \end{bmatrix} - \begin{bmatrix} b_1 \\ \vdots \\ b_i \\ \vdots \\ b_N \end{bmatrix} = Y\underline{a} - \underline{b} \quad (45)$$

Minimum squared error approach

- Instead of finding a solution \underline{a} that gives no misclassifications, which could be impossible if it is not a linearly separable problem, we want to find a solution \underline{a} that minimizes $|\underline{e}|^2$.
- This gives us the following sum of squared error criterion function:

$$\nabla J_s(\underline{a}) = |\underline{e}|^2 = |Y\underline{a} - \underline{b}|^2 = \sum_{i=1}^N (\underline{a}^T \underline{y}_i - b_i)^2 \quad (46)$$

Minimum squared error approach

- The gradient of $J_s(\underline{a})$ can be written as:

$$\nabla J_s(\underline{a}) = Y^T(Y\underline{a} - \underline{b}) \quad (47)$$

- This gives us the weight update formula as:

$$\underline{a}_{k+1} = \underline{a}_k + \rho_k \nabla J_p(\underline{a}) \quad (48)$$

$$\underline{a}_{k+1} = \underline{a}_k - \rho_k Y^T(Y\underline{a}_k - \underline{b}) \quad (49)$$

Minimum squared error approach

- Step 1:** Set an initial guess for the weight vector (\underline{a}_0) and let $k = 0$
- Step 2:** Determine \underline{a}_{k+1} as

$$\underline{a}_{k+1} = \underline{a}_k - \rho_k Y^T(Y\underline{a}_k - \underline{b}) \quad (50)$$

- Step 3:** If convergence reached, stop. Otherwise, go to Step 2.

Sequential variant of minimum squared error approach

- Sequential variant of minimum squared error approach:
- Step 1:** Set an initial guess for the weight vector (\underline{a}_0) and let $k = 0$
- Step 2:** Determine \underline{a}_{k+1} as

$$\underline{a}_{k+1} = \underline{a}_k - \rho_k (b_k - \underline{a}_k^T \underline{y}^k) \underline{y}^k \quad (51)$$

- Step 3:** If convergence reached, stop. Otherwise, go to Step 2.

Minimum squared error approach: parameter setting

- How do we set up parameters (i.e., ρ_k , b) for the MSE approach?
- Typically, ρ_k decreases with k (e.g., ρ_k/k) to obtain convergence
- In terms of b , useful settings include:
 - Setting \underline{b} as a vector of ones
 - Setting the first N_1 of the N components to N/N_1 and the rest to N/N_2 , where N_1 and N_2 are the number of samples in each class (e.g., if there are 10 samples in class 1 and 3 samples in class 2, then the first 10 components of \underline{b} are set to 13/10 and the rest are set to 13/3).

§ 11 - Clustering

* Definitions.

- A set of samples that are similar → Diff similarity criterions → diff clustering results
- A region in feature space containing a high density of samples → peaks in sample density func. → peaks are associated with clusters

↓

Clusters def ⇒ influences the strategies that we use to perform clustering.

* Intuitive Approaches.

① Mixture Density Strategy

- Def: high sample density region as cluster

- Give a set of N samples

↳ can estimate combined PDF ($p(\vec{x})$)
using density estimation (ex: Parzen Window)

↳ $p(\vec{x})$ refer to as a mixture density,
since it's a mixture of the K class densities:

$$p(\vec{x}) = \sum_{i=1}^K p(c_i) p(\vec{x} | c_i)$$

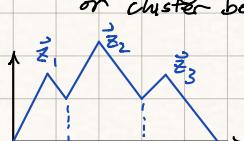
- Intuition:

- if K classes reasonably distinct & nearly equally likely,
⇒ $p(\vec{x})$ should have K peaks / local maxima,
one per class.

⇒ Define cluster prototype @ each local maxima:

$\vec{z}_i = \vec{x}$ s.t. \vec{x} is the i th local max of $p(\vec{x})$.

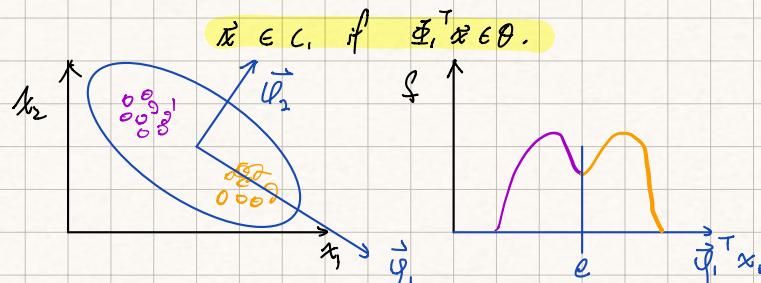
⇒ Based on the prototypes, we can determine our clusters using dist. metrics (ex: Euclidean Dist.) or cluster boundaries based on local minima.



- Cons:
 - In practice, too few samples will result in large # of spurious local maxima
 - Individual sample densities of a particular class may be skewed (ex. Gamma) and so a simple Euclidean distance metric may not define the clusters boundaries well.
 - Peaks of clusters may be too close (or even overlap) to distinguish properly

② PCA : Principal Component Analysis.

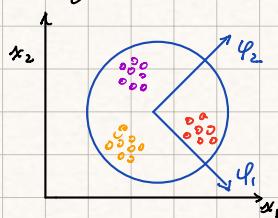
- Assume we have two fairly compact and distinct clusters.
- The combined covariance matrix describes the shape of the total sample distribution.
- If we take the maximum eigenvalue eigenvector \vec{q}_1 , as the maximum separation direction, we may do the following.
 - Project samples onto \vec{q}_1 .
 - Construct histogram.
 - Choose threshold θ @ minimum of distribution, and specify clusters based on the following rule:



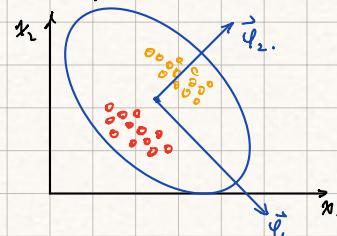
• Cons:

- While this clustering scheme is simple & intuitive, it cannot handle more complicated but typical problems such as:

• Having $K > 2$ clusters: \times



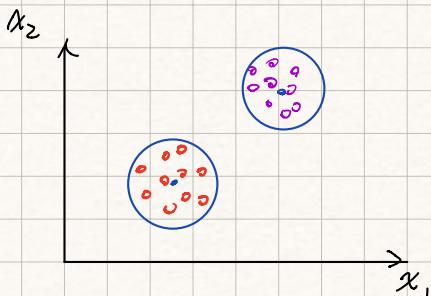
• Non-spherical Distributions \times



③ Euclidean Dist. Threshold Strategy

- Let's use the definition of clustering based on similarity with each other.
- Based on this notion, what we can do is assign a pattern to a cluster if the pattern to a cluster if the dist. between the cluster & the prototype is less than some threshold T .

- Step 1: Let $\vec{z}_1 = \vec{x}_1$, 1st sample
- Step 2: For each sample, assign it to closest prototype for which $|\vec{x}_i - \vec{z}_j| < T$, If no such prototype exists, let \vec{z}_i be a new prototype.
- Step 3: Repeat step 2 till all samples are assigned to clusters.



◦ Cons:

- T threshold arbitrary.
 - T small \Rightarrow many clusters will be found.
 - T large \Rightarrow few clusters will be found.
- Sensitivity to the order in which samples are considered
(affects specific prototypes)
- As such, most of these intuitive approaches seem rather arbitrary!
- To help remedy such issues,
we want to develop formal criteria for clustering that are less likely to impose some artificial structure on data, and more sensitive to any natural structural.

* Criterion for Clustering

- ↳ serve quantitative measure of the clusters.
- ↳ A given partitioning of the sample set is optimum, when criterion function is maximized/minimized.
- ↳ commonly used criterion functions:
 - Sum of squared errors
 - Scatter volume.

* Sum of Squared Error Criteria.

- For K clusters, cluster C_i has N_i samples.

$$J_e \rightarrow \min \text{tr}(S_w)$$

$$\Rightarrow J_i = \sum_{\vec{x} \in C_i} |\vec{x} - \vec{\bar{x}}_i|^2$$

↳ The prototype $\vec{\bar{x}}_i$ which minimize the single class error is just the class mean $\vec{\bar{m}}_i$

⇒ Total:

$$J_e = \sum_{i=1}^K J_i = \sum_{i=1}^K \sum_{\vec{x} \in C_i} |\vec{x} - \vec{\bar{m}}_i|^2$$

- A clustering / partitioning that minimizes this total error J_e is called Minimum Variance Partition.

* Scatter Volume Criteria

- For class C_i , scatter matrix: $S_i = \sum_{\vec{x} \in C_i} (\vec{x} - \vec{\bar{m}}_i)(\vec{x} - \vec{\bar{m}}_i)^T$

$$J_v \rightarrow \min |S_w|.$$

- S_i is just N_i times the covariance matrix of class i .

- Summing over K classes, we have a measure of total within class scatter:

$$S_w = \sum_{i=1}^K S_i$$

- Note: trace of S_w = sum of squared error criterion

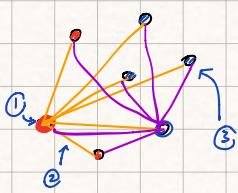
$$\text{tr}(S_w) = \sum_{i=1}^k \text{tr}(S_i) = \sum_{i=1}^k \sum_{x \in C_i} |\vec{x} - \vec{m}_i|^2 = J_e$$

- Scatter Volume: (the determinant of S_w)

$$J_v = |S_w|$$

* Minimum Variance Algo \Rightarrow minimize sum of sq. error criterion $\sim \min \text{tr}(S_w)$

* K-means Algo



Step 1: Choose prototypes $\{\vec{z}_1, \dots, \vec{z}_k\}$ arbitrarily.

Step 2: Euclidean Dist. btwn N samples to each cluster

Step 3: Assign N samples to the K clusters based on min Euclidean Dist.

$\vec{x} \in C_i$ iff $|\vec{x} - \vec{z}_i| < |\vec{x} - \vec{z}_j|, j \neq i$

Step 4: Compute new cluster prototypes as the cluster means:

$$\vec{z}_{i,\text{new}} = \frac{1}{N_i} \sum_{i=1}^{N_i} \vec{x}_i, \forall x \in C_i$$

Step 5: If any cluster prototypes change, go back to Step 2.

Cons:

- { ① Sensitivity to initialization. \Rightarrow can be reduced by running random initializations multiple times & choosing the best results.
- ② Assumes # of cluster is known \Rightarrow can be reduced by starting with large # of clusters and then consolidating them.
- ③ Sensitive to geometric properties of clusters.

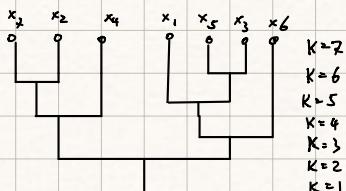
Variants:

- Complex: ISODATA : (Iterative Self Organizing Data Analysis Technique A)
 - if # of samples in any cluster is less than some threshold, the cluster may be eliminated.
 - if max variance feature for the cluster is larger than some threshold and there are sufficient samples in the cluster, split the cluster.
 - if the pairwise dist. btwn clusters is below some threshold, the clusters are combined.
- ISODATA is quite flexible, but needs # of thresholds to be defined.
- ISODATA works best when used in an interactive, empirical env.

Hierarchical Clustering: Simple Example.

- N samples, treat each as a cluster prototype.
 - We now find the two most similar clusters & merge them, $\Rightarrow (N-1)$ clusters.
 - Continue merging till some criterion is satisfied.
 - Scheme :-

Denogram



252

Measures of similarity for Hierarchical Clustering

- Popular similarity measures based on Euclidean distance.
(Each measure gives diff. clustering results)

↳ Minimum distance (nearest neighbor)

$$d_{\min}(c_i, c_j) = \min_{\tilde{x} \in c_i, \tilde{x}' \in c_j} |\tilde{x} - \tilde{x}'|$$

- Nearest Neighbor measure well-suited for string-like clusters but highly sensitive to noise and outliers.

L Maximum dist (furthest neighbor)

$$d_{\max}(c_i, c_j) = \max_{\vec{x} \in c_i, \vec{x}' \in c_j} |\vec{x} - \vec{x}'|$$

- tend to discourage growth of elongated clusters.
since two sub-clusters are only merged if
the least similar pair in the resulting cluster is sufficiently similar.
 - This makes it less sensitive to outliers & noise
 \Rightarrow well suited for compact, spherical clusters.

↳ Average list:

$$d_{avg}(c_i, g_j) = \frac{1}{N_i N_j} \sum_{x \in c_i} \sum_{x' \in g_j} |\vec{x} - \vec{x}'|$$

- less sensitive to noise & outliers.

↳ Dist. btwn means

$$d_{mean}(c_i, c_j) = |\vec{m}_i - \vec{m}_j|$$

- Inter-mean measure \Rightarrow less sensitive to noise & outliers
 - Most efficient to implement.

Termination Criterion

- Maximum # of levels
 - # of clusters
 - Error threshold