

By Sakthi Sir

Principal component analysis

19/9/22 Monday

DATE: 16/9/22

PAGE NO.:

Friday

PCA

$$f_1 = x \quad f_2 = y$$

$$X = P_1 \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

1. Make mean-subtracted data

$$x \rightarrow x', \quad x' = \begin{bmatrix} x_1 - \mu_x & y_1 - \mu_y \\ x_2 - \mu_x & y_2 - \mu_y \\ x_3 - \mu_x & y_3 - \mu_y \end{bmatrix}$$

2. calculate the dispersion Matrix $\Sigma = \frac{1}{N} x^T x$.

$$\Sigma = \begin{bmatrix} \text{var}(x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{var}(y) \end{bmatrix}$$

$$\text{var}(x) = \frac{1}{N} \sum (x_i - \mu_x)^2$$

$$\text{cov}(x, y) = \text{cov}(y, x) = \frac{1}{N} \sum (x_i - \mu_x)(y_i - \mu_y)$$

$$\Sigma = \frac{1}{N} \begin{bmatrix} x_1 - \mu_x & x_2 - \mu_x & x_3 - \mu_x \\ y_1 - \mu_y & y_2 - \mu_y & y_3 - \mu_y \end{bmatrix} \begin{bmatrix} x_1 - \mu_x & y_1 - \mu_y \\ x_2 - \mu_x & y_2 - \mu_y \\ x_3 - \mu_x & y_3 - \mu_y \end{bmatrix}$$

$$\Sigma = \frac{1}{N} \begin{bmatrix} \sum_{i=1}^3 (x_i - \mu_x)^2 & \sum_{i=1}^3 (x_i - \mu_x)(y_i - \mu_y) \\ \sum_{i=1}^3 (x_i - \mu_x)(y_i - \mu_y) & \sum_{i=1}^3 (y_i - \mu_y)^2 \end{bmatrix}$$

3. Find D no. of eigen values of eigen vectors of Σ .

$$\lambda_1 > \lambda_2 > - - - > \lambda_D \quad (D=2 \text{ in 2 attributes})$$

$\downarrow \quad \downarrow \quad \downarrow$
 $v_1 \quad v_2 \quad v_D$

$$\vec{v}_i = \begin{bmatrix} v_{i1} \\ v_{i2} \\ \vdots \\ v_{iD} \end{bmatrix}$$

$$W = \begin{bmatrix} v_1 & v_2 & - & - & - & v_D \\ v_{11} & v_{21} & & & & v_{D1} \\ v_{12} & v_{22} & & & & v_{D2} \\ | & | & & & & | \\ | & | & & & & | \\ v_{1D} & v_{2D} & & & & v_{DD} \end{bmatrix}_{B \times D}$$

$\lambda_1 > \lambda_2 > - - - > \lambda_D$

Eg:

$$y = \underbrace{w^T x}_\uparrow$$

data in another D-dimensional space

$$\text{Trace}(\Sigma) = \sum_{i=1}^D \lambda_i$$

for 'd' eigen values

$$\frac{\sum_{i=1}^d \lambda_i}{\text{Trace}(\Sigma)} \times 100 > 98$$

or

$$\sum_{i=1}^d \lambda_i$$

let $D = 100$, $d = 20$.

- \therefore From adding 1st 20 eigen values, we are covering 98% of data acc to its significance
 \therefore other 80 eigen values are ~~not~~ not significant.
 \therefore discard $D-d'$ eigen values.

$$\therefore \text{We consider } W = \begin{bmatrix} v_1 & -v_{d+1} \\ v_2 & v_{d+2} \\ \vdots & \vdots \\ v_D & -v_{D+d} \end{bmatrix}_{(D \times d)}$$

$$\cancel{W^T X} \quad Y = W^T (D \times D) \times X (D \times 1)$$

$$Y_{D \times 1} =$$

Eg:

$$X = \begin{bmatrix} 2 & 1 \\ 3 & 5 \\ 4 & 3 \\ 5 & 6 \\ 6 & 7 \\ 7 & 8 \end{bmatrix}$$

Range of 1st dimension
 2 to 7

$$\mu = (4.5) \ (5)$$

$$\therefore X' = \begin{bmatrix} -2.5 & -4 \\ -1.5 & 0 \\ -0.5 & -2 \\ 0.5 & 1 \\ 1.5 & 2 \\ 2.5 & 3 \end{bmatrix}$$

$$\Sigma = \frac{1}{6} \begin{bmatrix} -2.5 & -1.5 & -0.5 & 0.5 & 1.5 & 2.5 \\ -4 & 0 & -2 & 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} -2.5 \\ -1.5 \\ -0.5 \\ 0.5 \\ 1.5 \\ 2.5 \end{bmatrix}$$

$$= \frac{1}{6} \begin{bmatrix} 17.5 & 22 \\ 22 & 34 \end{bmatrix} = \begin{bmatrix} 2.92 & 3.67 \\ 3.67 & 5.67 \end{bmatrix}$$

Now eigen values.

$$\begin{bmatrix} 2.92 & 3.67 \\ 3.67 & 5.67 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = 0$$

$$\Rightarrow \lambda_1 = 8.22 \text{ and } \lambda_2 = 0.38$$

} for eigen vectors:
 $\sum \bar{v}_i = \lambda_i \bar{v}_i$

$$\text{Trace}(\Sigma) = 2.92 + 5.67 = 8.59 \\ = \lambda_1 + \lambda_2$$

we see $\lambda_1 \gg \lambda_2$

$$\frac{\lambda_1}{8.59} \times 100 = 95.6$$

$$\Sigma \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 8.22 \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

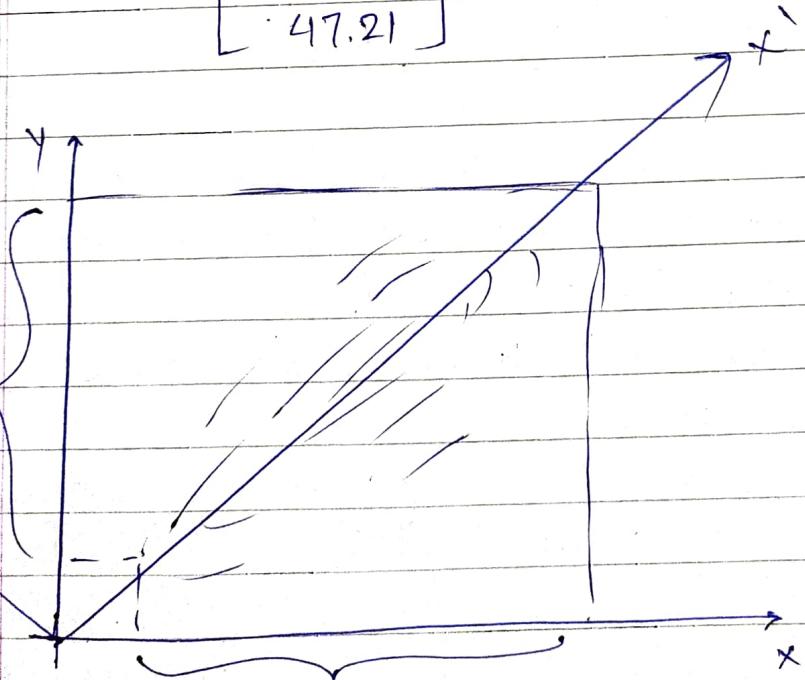
$v = \begin{bmatrix} 2.55 \\ 3.67 \end{bmatrix}$ eigen vector
 $= w$

$$y = w^T x$$

$$y_1 = [2.55 \quad 3.67] \begin{bmatrix} 2 \\ 1 \end{bmatrix} = [8.77]$$

(1st data entry)

$$\therefore y = \begin{bmatrix} 8.77 \\ 26 \\ 21.21 \\ 34.77 \\ 40.92 \\ 47.21 \end{bmatrix} \rightarrow \text{Range: } 8.77 \text{ to } 47.21$$



We are rotating the axes such that the variance of 1st dimension is increased.

Advantages

1. computational complexity
→ By changing dimension space, we compute less differences.
2. In dataset, don't know the meaning of the attributes, we also don't know if any attribute is redundant or if $attr_1 + attr_2 = attr_3$.
∴ by this process, we can remove redundancy.
 - For 8.77, we have considered 2 & 1 both (both features),
 - In transformed space, each value has considered every attribute's value.

Disadvantages

1. In OG dataset, all stored data is transformed in Y, although all values of OG set has contributions, the new values obtained, we don't know what this attribute denotes.

$$X = \begin{bmatrix} x & y \\ 2 & 1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix} \rightarrow Y = \begin{bmatrix} f' \\ 8.77 \end{bmatrix}$$

← don't know
this
attribute's
meaning

→ meaning of dataset is lost in transformed space.

2. If the weightage of attributes is given, no transformation required.

→ When machine improves performance, minimizes error, we call it Machine Learning.

Supervised

	f_1	f_2	Class Label
P ₁	3	4	1
P ₂	4	8	2
P ₃	9	10	1

Unsupervised

→ Data w/o class label info.
 → We 1st have to group that data, build an algo for it, based on similarity in patterns.

→ This is unsupervised learning

→ Done in clustering techniques.

- Known: Class label with data, we call it Training dataset
- we know data pattern belongs to what class.

- Then we can apply class labels to unknown data based on learning from training dataset.
- This is supervised learning
- Lesser chance of error
- Error minimized.

- Supervised: Prior info given
- Unsupervised: No prior info given, clustering on basis of similarity.

Classification

- Supervised learning
- Given a new or unknown pattern and assign it into a predefined classes.
- Generate and train a classification model with the help of training data (or other form of data).

2

clustering

- Unsupervised learning
- Based on similarity or dissimilarity measures among patterns, develop an algorithm to find the natural grouping of the data.
- Look for intraclass (within) distance minimized, inter-class distⁿ max (between)

Classification Techniques

1. Probabilistic: Bayes Classifier

2. Similarity-based: MDC, KNN, LDC, QDC
(minⁿ distance)

3. SVM

Probabilistic or Bayes' Rule of classification

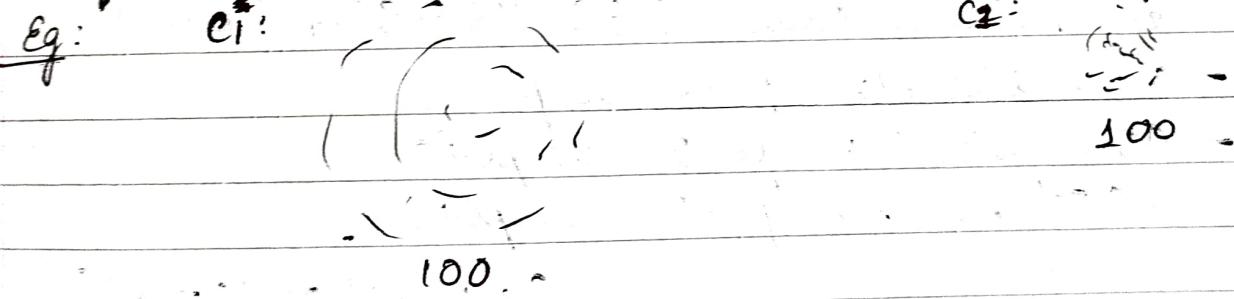
Prior info could be in form of level data or any other form.

conditional Prob. Density funcⁿ given.

Dataset : c_1	Data	Patterns : 100
: c_2	class	500
: c_3		1000

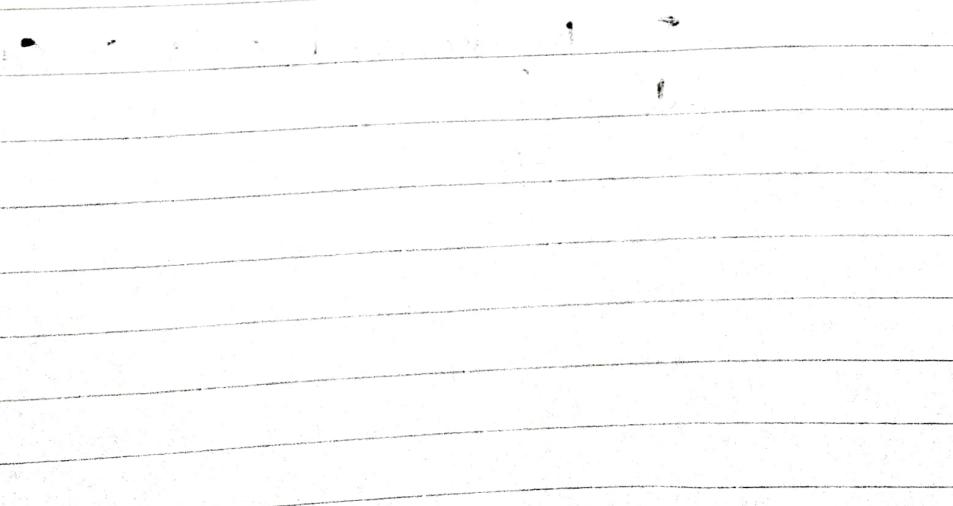
Now for a new pattern, $P(\text{class})$ to be in class c_1, c_2, c_3 ,
 $c_3 \Rightarrow \frac{1}{16}, \frac{5}{16}, \frac{10}{16}$ resp \rightarrow Prior probability

Unique class

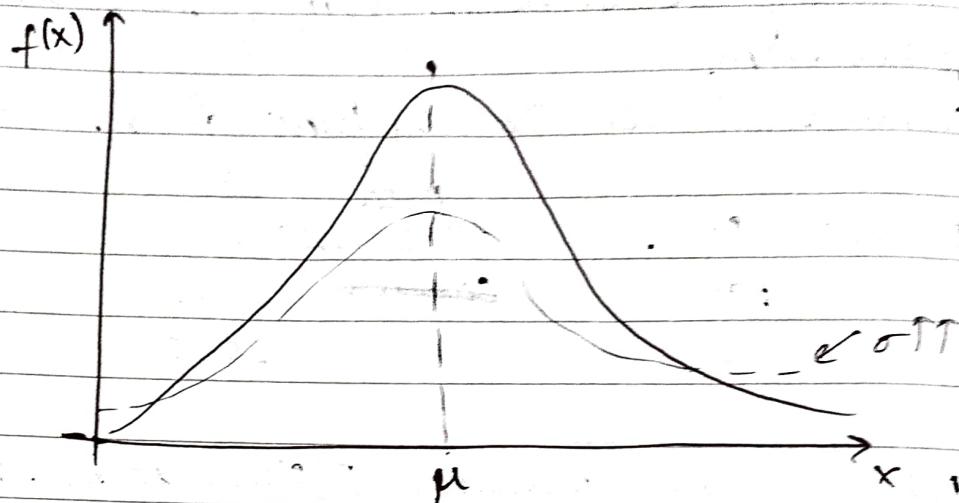


We also need to know how data is distributed within a class.

Given by conditional Prob. Density function.



$$f(x) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$



- If $\sigma \gg \mu$, peak is lower.

\Rightarrow Empirical Rules (68-95-99.7% Rule)

- 68% lies b/w $\mu - \sigma$ to $\mu + \sigma$.
- 95% lies b/w $\mu - 2\sigma$ to $\mu + 2\sigma$
- 99.7% lies b/w $\mu - 3\sigma$ to $\mu + 3\sigma$.

$$f(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right\}}$$

d : No. of dimensions

Σ : Dispersion Matrix or ~~sd~~ variance

$$\begin{cases} \sqrt{\Sigma} = \sigma \\ \uparrow \\ SD. \end{cases}$$

Bayes' Rule

→ Let there be M classes ($M \geq 2$),
Prior probability be p_1, p_2, \dots, p_M ,
where $p_i \in [0, 1]$ and $\sum_{i=1}^M p_i = 1$.

→ Let $p_1(x), p_2(x), \dots, p_n(x)$ be the conditional probability density functions

→ Given a noisy unknown pattern x_0 ,
 x_0 be in class i if

$$p_i p_i(x_0) \geq p_j p_j(x_0); \quad \forall i \neq j.$$

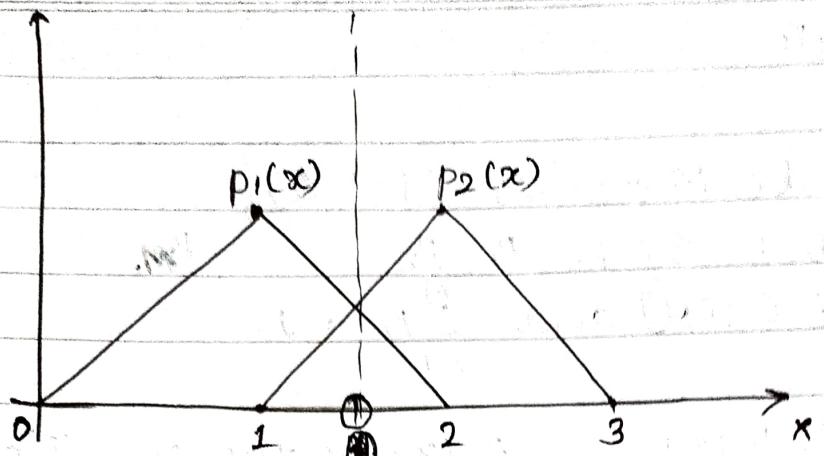
Eg: Let there be 2 classes ($M=2$)

Prior probability = $P, (1-P)$

$$p_1(x) = \begin{cases} x & 0 \leq x \leq 1 \\ 2-x & 1 \leq x \leq 2 \\ 0 & \text{otherwise} \end{cases}$$

$$p_2(x) = \begin{cases} x-1 & 1 \leq x \leq 2 \\ 3-x & 2 \leq x \leq 3 \\ 0 & \text{otherwise} \end{cases}$$

Find the classification model



$\leftarrow \underline{x = 1.5}$
 \leftarrow line of classification

\therefore if $x < 1.5 \rightarrow$ class 1
 $\quad \quad \quad x > 1.5 \rightarrow$ class 2.

But for Bayes' Rule, we need to calculate & find the line of classification.

Domain of space : $x \in [0, 3]$

Break it : $(0-1), (1-$

$$0 \text{ to } 1 : P_1 = P \quad P_2 = 1 - P$$

$$p_1(x) = x \quad p_2(x) = 0$$

$$P_1 p_1(x) = (\cancel{Px}) \quad P_2 p_2(x) = (0)$$

$$\therefore P_1 p_1(x) \geq P_2 p_2(x)$$

$\therefore x$ should be in class 1

if $x \in [0, 1]$.

Case 2:

2 to 3:

$$P_1 = P$$

$$P_2 = 1 - P$$

$$p_1(x) = 0$$

$$p_2(x) = 3 - x.$$

$$P_1 p_1(x) = 0 < P_2 p_2(x) \quad (1-P)(3-x)$$

\therefore if $x \in [2, 3]$, x should be in class 2.

Case 3:

1 to 2:

$$P_1 = P$$

$$P_2 = 1 - P$$

$$p_1(x) = 2 - x$$

$$p_2(x) = x - 1.$$

If $P_1 p_1(x) > P_2 p_2(x) \rightarrow x$ in class 1

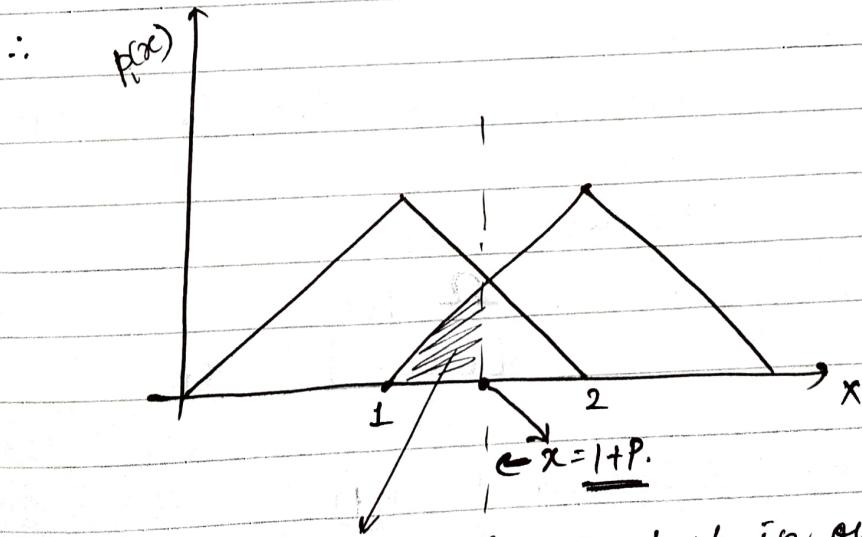
$$\Rightarrow P(2-x) \geq (1-P)(x-1)$$

$$2P - Px \geq x - 1 - P + P$$

$$x \leq P+1.$$

$\therefore x$ will be in class 1 if $x \leq P+1$

and x will be in class 2 if $x > P+1$



- belong to class 2, but in our rule, they ~~are~~ are counted in class 1

- \therefore Bayes Rule has ~~#~~ some error.

Reason: ~~Dataset has data i~~

- Data is not well-separated

- There is overlapping.

$P_i p_i(x) \geq P_j p_j(x)$, then x should be in class i

- Overlapping of pattern from class 1 onto class 2 due to classification = ERROR.

$$\text{ERROR: } P_1 \times \int_{\Omega} p_1(x) dx$$

$\underbrace{\quad}_{1+P}$

$\left\{ \begin{array}{l} (1+P \text{ to } 3) \in \text{outside} \\ \text{class 1} \end{array} \right.$

No. of patterns that lie outside the class 1 and overlap in class 2

Let $\Omega \Rightarrow \text{Domain} \Rightarrow [0 \text{ to } 3]$

$\Omega_1 \Rightarrow \text{Domain of class 1 in } \Omega$,

$$\Rightarrow [0 \text{ to } 1.5]$$

$\Omega_2 \Rightarrow \text{Domain of class 2 in } \Omega$,

$$\Rightarrow [1.5 \text{ to } 3]$$

$\therefore \Omega_1^c \Rightarrow \Omega_1 \text{ complement}$

$\Rightarrow \text{Range in Domain that is outside class 1}$

$$\Rightarrow \Omega - \Omega_1$$

$$\Rightarrow [1.5 \text{ to } 3] = \Omega_2$$

Similarly $\Rightarrow \Omega_2^c = \Omega - \Omega_2$

$$[0 \text{ to } 1.5] = \Omega_1$$

error due to class 1 = \textcircled{e}_1

$$e_1 = P_1 \int_{\Omega_1^c} p_1(x) dx$$

~~Total Error~~ (E).

$$E = \sum_{i=1}^c e_i = \sum_{i=1}^c \int_{\Omega_i^c} p_i(x) dx.$$

This is called misclassification Probability or Error Rate.

[In the given example

$$\begin{aligned}
 E &= P_1 \int_{\Omega_1^c} p_1(x) dx + P_2 \int_{\Omega_2^c} p_2(x) dx \\
 &= P \int_{1-P}^3 p_1(x) dx + (1-P) \int_0^{1+P} p_2(x) dx \\
 &= P \left[\int_{1-P}^2 (2-x) dx + \int_2^3 0 dx \right] + (1-P) \left[\int_0^1 0 dx + \int_1^{1+P} (x-1) dx \right] \\
 &= P \left[2x - \frac{x^2}{2} \Big|_{1-P}^2 \right] + (1-P) \left[\frac{x^2}{2} - x \Big|_1^{1+P} \right] \\
 &= P \left[2(1-P) + \frac{(P^2 + 2P - 3)}{2} \right] + (1-P) \left[\frac{P^2 + 2P}{2} - P \right] \\
 &= P \left[\cancel{P} \frac{P^2 - 2P + 1}{2} \right] + (1-P) \left[\frac{P^2}{2} \right] \\
 &\quad \circ \frac{P(P-1)^2}{2} + \frac{P^2(1-P)}{2} \\
 &= \frac{P(1-P)}{2} (1-P+P) = \boxed{\frac{P(1-P)}{2} = E}
 \end{aligned}$$

Classifiers

→ Different types of classifiers divide in different categories

Eg: classifier 2.

$$\text{class 1} \rightarrow x \leq 1.5$$

$$\text{class 2} \rightarrow \text{o/w}$$

classifier 3

$$\text{class 1} \rightarrow x \leq 1.4$$

$$\text{class 2} \rightarrow \text{o/w.}$$

Now for classifier 2, $E = ?$.

$$E = P_1 \int_{-1.5}^2 p_1(x) dx + P_2 \int_{-1.5}^{1.5} p_2(x) dx.$$

$\left(1-P \right)$ $\xrightarrow{\quad}$

$$= P \left[\frac{2x - x^2}{2} \right]_{-1.5}^{1.5} + (1-P) \left[\frac{x^2 - x}{2} \right]_{-1.5}^{1.5}$$

$$= P \left[1 - \frac{7}{8} \right] + (1-P) \left(\frac{5 - 1}{8} \right)$$

$$= P \left(\frac{1}{8} \right) + (1-P) \left(\frac{1}{8} \right)$$

$$= \boxed{\frac{1}{8}}$$

$$= \boxed{0.125}$$

- Bayes' classifier is the best classifier.
- No other classifier works better.
- E_{Bayes} classifier \Rightarrow is minimum.

i.e. $\frac{P(1-P)}{2} \leq (\text{Any other error})$ (in the given example)

$$\therefore \frac{P(1-P)}{2} \leq 0.125 \quad (E_{1.5} = 0.125)$$

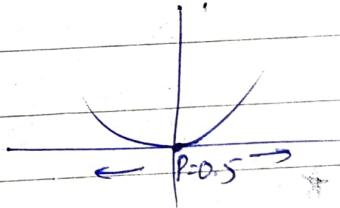
Let $P = 0.4$

$$\therefore E_{\text{Bayes}} = \frac{\frac{4}{10} \times \frac{6}{10} \times 1}{2} \Rightarrow 0.12$$

Let $P = 0.6$

$$\therefore E_{\text{Bayes}} = 0.12$$

$$P = 0.5 \rightarrow E_{\text{Bayes}} = 0.125 = E_{1.5}$$



If given data is in form of training data, like

Eg: $f_1 f_2 f_3 \rightarrow \text{Classes}$

$c_1 \downarrow c_2 \downarrow c_3 \downarrow$ } can estimate prior probs.

No. of patterns: 100 200 300

$$\Rightarrow p_1 = \frac{1}{6}, p_2 = \frac{2}{6}, p_3 = \frac{3}{6}$$

for c_1

$p_1(x) = \text{Assume } \text{false} \text{ dataset forms Gaussian/Normal distribution.}$

$$p_i(x) = \frac{\{ -\frac{1}{2} (x - \mu_i)' \Sigma_i^{-1} (x - \mu_i) \}}{(2\pi)^{d/2} |\Sigma_i|^{1/2}}$$

$$\begin{aligned} c_1 &\rightarrow 100 \rightarrow \mu_1, \Sigma_1 \\ c_2 &\rightarrow 200 \rightarrow \mu_2, \Sigma_2 \\ c_3 &\rightarrow 300 \rightarrow \mu_3, \Sigma_3 \end{aligned} \quad \left. \begin{array}{l} \text{calculate} \\ \hline \end{array} \right.$$

- Then apply Bayes' Rule for classification
 - Baye's classifier used very less in Data science
Reason:
→ To solve real time prob, we get training dataset and we assume Gaussian Distribution but in real world probs, the dist' is not normal.
- $p_i(x) = ?$
- * ① Estimation of ^{conditional} prob. density func' is very difficult.
- ② It is computationally very complex.