# ML Mid-2 :-

* **Bayes Decision Theory** :- (Continuous Features).

→ Allowing actions (decisions) other than classification.

→ Loss function ⇒ States ~~too costly~~ cost of each action taken.

**Loss function :-**

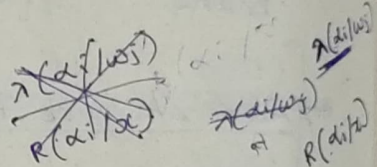$\lambda(\alpha_i/\omega_j)$ ⇒ loss of taking @ $\alpha_i$ action when the state of nature is $\omega_j$

$\Omega = \{\omega_1, \omega_2, \ldots \omega_c\}$ ⇒ set of states of nature (or) classes (or) Categories

$A = \{\alpha_1, \alpha_2, \ldots \alpha_a\}$ ⇒ set of possible actions.

Given a pattern $x$, function $\alpha(x)$ → action taken.

$$\alpha : x \to \alpha(x).$$

**How to find the best action?**

$R(\alpha_i/x)$ ⇒ Conditional risk.

⇒ risk of taking action $\alpha_i$ when given pattern is $x$.

Best action $\alpha_k$ ⇒ min $\{R(\alpha_1/x), R(\alpha_2/x), \ldots R(\alpha_a/x)\}$.
and min. Risk ⇒ Bayes risk.

$$R(\alpha_i/x) = \sum_{j=1}^{c} \lambda(\alpha_i/\omega_j) \cdot P(\omega_j/x).$$

$$R = \int R(\alpha(x)/x) \cdot P(x) \cdot dx. \Rightarrow \left(\begin{array}{c}\text{Should be}\\ \text{Minimum}\end{array}\right)$$

**Two-category classification :-**

$\alpha_1 \to$ decide '$\omega_1$'
$\alpha_2 \to$ decide '$\omega_2$'.

$$\lambda_{ij} = \lambda(\alpha_i/\omega_j)$$

$$R(\alpha_i/x) = \sum_{j=1}^{c} \lambda(\alpha_i/\omega_j) \cdot P(\omega_j/x)$$

$\alpha_1 \Rightarrow R(\alpha_1/x) = \cancel{\lambda\ell\alpha\cancel{/}} \lambda_{11} \cdot P\left(\frac{\omega_1}{x}\right) + \lambda_{12} \cdot P\left(\frac{\omega_2}{x}\right).$

$\alpha_2 \Rightarrow R(\alpha_2/x) = \lambda_{21} \cdot P\left(\frac{\omega_1}{x}\right) + \lambda_{22} P\left(\frac{\omega_2}{x}\right).$

if $R(\alpha_1/x) < R(\alpha_2/x)$.

then    action   $\alpha_1$ : "decide   $\omega_1$"   is taken

Otherwise     action   $\alpha_2$ : "decide   $\omega_2$" is taken.

**Eg :-**

$\alpha_1 \longrightarrow$ decide '$\omega_1$' (criminal).

$\alpha_2 \rightarrow$ decide '$\omega_2$' (innocent)

$R(\alpha_i/x) = \sum\limits_{j=1}^{c} \lambda(\alpha_i/\omega_j) \cdot P(\omega_j/x).$

$\alpha_1 \Rightarrow R(\alpha_1/x) = \lambda_{11} \cdot P\left(\frac{\omega_1}{x}\right) + \lambda_{12} \cdot P\left(\frac{\omega_2}{x}\right).$

$\qquad = 0 + 10 \times \dfrac{P(\omega_2) \cdot P(x/\omega_2)}{P(x)}$

$\qquad = 10 \times \dfrac{9.5 \times 0.86}{0.7}$

$\qquad = 3\cancel{4}00$ over $7$

$\alpha_2 \Rightarrow R(\alpha_2/x) = \lambda_{21} \cdot P\left(\frac{\omega_1}{x}\right) + \lambda_{22} \cdot P\left(\frac{\omega_2}{x}\right).$

$\qquad = 1 \times \dfrac{P(\omega_1) \cdot P(x/\omega_1)}{P(x)} + 0$

$\qquad = \dfrac{0.5 \times 0.8}{0.7} = \dfrac{4}{7}$

$P(x) = P(\omega_1) \cdot P(x/\omega_1)$
$\qquad\quad +$
$\qquad P(\omega_2) \cdot P(x/\omega_2)$

$P(x) = 0.5 \times 0.8 +$
$\qquad\quad 0.5 \times 0.6$
$\qquad = 0.5 \times 1.4$
$\qquad = 0.7$

$\therefore$   '$\alpha_2$' action taken $\Rightarrow$ innocent

**\* Likelihood ratio :-**

$\qquad$ Likelihood ratio $\Rightarrow \dfrac{P(x/\omega_1)}{P(x/\omega_2)}$

**Eg:-**

$R(\alpha_1/x) = \lambda_{11} \cdot P\left(\frac{\omega_1}{x}\right) + \lambda_{12} \cdot P\left(\frac{\omega_2}{x}\right) + \lambda_{13} P\left(\frac{\omega_3}{x}\right)$ ; $= 0 + 1 \times 0.4 + 2 \times 0.5$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad = 1.4$

$R(\alpha_2/x) = \lambda_{21} \cdot P\left(\frac{\omega_1}{x}\right) + \lambda_{22} \cdot P\left(\frac{\omega_2}{x}\right) + \lambda_{23} \cdot P\left(\frac{\omega_3}{x}\right) = 1 \times 0.1 + 0$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad + 2 \times 0.5$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad = 1.1$

$R(\alpha_3/x) = \lambda_{31} \cdot P\left(\frac{\omega_1}{x}\right) + \lambda_{32} \cdot P\left(\frac{\omega_2}{x}\right) + \lambda_{33} \cdot P\left(\frac{\omega_3}{x}\right) = 3 \times 0.1$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad + 10 \times 0.4$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad + 0$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad = 4.3$

$\therefore$   '$\alpha_2$' action is taken.

Reject :-

→ Actions :- assign a class label or reject.

→ When misclassification is costly, we can reject to classify.

* Naive Bayes classification :-

(See slides)

* Maximum likelihood parameter estimation :-

→ Here, we assume parameters are unknown but fixed.

Let the parameters we are trying to estimate be

$\theta = (\mu, \Sigma)^t$.

$D \to$ training set. Contains 'n' samples $x_1, x_2, \ldots, x_n$.

Likelihood of '$\theta$', $P(D/\theta) = \prod_{i=1}^{n} P(x_i/\theta)$.

Maximum - likelihood estimate $\Rightarrow \hat{\theta}$, that maximizes $P(D/\theta)$.

log - likelihood function,

$$l(\theta) = \ln \cdot P(D/\theta).$$

$$= \ln \prod_{i=1}^{n} \cdot P(x_i/\theta).$$

$$= \sum_{i=1}^{n} \cdot \ln(P(x_i/\theta))$$

Let parameter vector $\theta = (\theta_1, \theta_2, \ldots, \theta_p)^t$.
(P-parameter).

$\nabla_\theta \Rightarrow$ gradient operator.

$$\nabla_\theta = \begin{bmatrix} \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_p} \end{bmatrix}$$

t to classify.

$$\Rightarrow \nabla_\theta l = \begin{bmatrix} \dfrac{\partial l}{\partial \theta_1} \\ \vdots \\ \dfrac{\partial l}{\partial \theta_p} \end{bmatrix}$$

For maximum - likelihood . estimate, $\nabla_\theta l = 0$.

Estimate the parameters in '$\theta$'.

* ## Maximum - likelihood estimation of Gaussian Distribution :-

be

Let `D` training set contains 'n' samples.

$$D = \{x_1, x_2, \cdots x_n\}$$

···· $x_n$.

$$P(D/\theta) = \prod_{i=1}^{n} \cdot P(x_i/\theta).$$

).

Gaussian Distribution,    $\theta = (\mu, \sigma^2)$

imizes $P(D/\theta)$.

$$L\left(\dfrac{x}{\theta}\right) = P\left(\dfrac{x}{\theta}\right) = \dfrac{1}{\sigma\sqrt{2\pi}} \propto e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

Log - likelihood,

$$\ln P\left(\dfrac{x}{\theta}\right) = \ln\left(\dfrac{1}{\sigma\sqrt{2\pi}} \times e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}\right).$$

$$LL\cdot(x/\theta) = - \ln\sigma - \dfrac{1}{2}\ln(2\pi) - \dfrac{1}{2\sigma^2}\sum_{i=1}^{n}(x_i-\mu)^2$$

Differentiate it with $\dfrac{\partial}{\partial\mu}$ and $\dfrac{\partial}{\partial\sigma}$ to estimate '$\theta$'. 

maximise

)t.

$$\dfrac{\partial}{\partial\mu}\left(LL(x/\theta)\right) = 0 + 0 - \dfrac{1}{2\sigma^2}\sum_{i=1}^{n}2(x_i-\mu) = 0$$

$$\Rightarrow \dfrac{\sum_{i=1}^{n}(x_i-\mu)}{\sigma^2} = 0.$$

$$\sum_{i=1}^{n}(x_i-\mu) = 0.$$

$$\sum_{i=1}^{n}x_i - \sum_{i=1}^{n}\mu = 0.$$

$$\Rightarrow \boxed{\mu = \dfrac{\sum_{i=1}^{n}x_i}{n}}$$

$$\frac{\partial}{\partial \sigma} \left( LL(x/\theta) \right) = \sum_{i=1}^{n} -\frac{1}{\sigma} + \sum_{i=1}^{n} \frac{2}{2\sigma^3} (x_i - \mu)^2 = 0$$

$$\sum_{i=1}^{n} \frac{1}{\sigma^2} \cdot (x_i - \mu)^2 = \sum_{i=1}^{n} \frac{1}{\sigma}$$

$$\frac{1}{\sigma^2} \sum_{i=1}^{n} (x_i - \mu)^2 = n$$

$$\Rightarrow \boxed{\sigma^2 = \frac{\sum_{i=1}^{n} (x_i - \mu)^2}{n}}$$

**Univariate :-**

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

**Multivariate :-**

$$P(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \, e^{-\frac{1}{2} \left( (x-\mu)^t \cdot \Sigma^{-1} (x-\mu) \right)}$$

$$\mu = \frac{\sum_{i=1}^{n} x_i}{n}$$

$$\Sigma = \frac{\sum_{i=1}^{n} (x_i - \mu)(x_i - \mu)^t}{n}$$

**\* Regression :-**

→ Regression analysis investigates the relationship b/w two (or) more variables in non-deterministic fashion.

→ Linear regression attempts to model the relationship b/w two variables by fitting a linear equation to observed data.

→ Scatter plot helps to determine the strength of relationship b/w two variables.

$$y = b_0 + b_1 x \quad \Rightarrow \text{Linear regression model}$$

$$\boxed{b_0 = \bar{y} - b_1 \bar{x}} \quad \Rightarrow Y\text{-intercept}$$

$$\boxed{b_1 = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2}} \quad \Rightarrow \text{slope}$$

If $b_1 > 0 \Rightarrow x$ (predictor) and $y$ (target) have positive relationship. $(x\uparrow \quad y\uparrow)$.

If $b_1 < 0 \Rightarrow x$ (predictor) and $y$ (target) have negative relationship $(x\downarrow \quad y\downarrow)$.

$$y = b_0 + b_1 x. \longrightarrow \text{predicted output}$$

Sum of squared error $\Rightarrow \boxed{\sum_{i=1}^{n}(\text{actual output} - \text{predicted output})^2}$

Eg:-

$$\begin{bmatrix} 1 & 6 & 4 & 11 \\ 1 & 8 & 5 & 15 \\ 1 & 12 & 9 & 25 \\ 1 & 2 & 1 & 3 \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} = \begin{bmatrix} 20 \\ 30 \\ 50 \\ 7 \end{bmatrix}$$

$$\theta = (z^T.z)^{-1}.z^T y$$

$$\begin{bmatrix} 1 \\ 6 \\ 4 \\ 11 \end{bmatrix}$$