

### 9.1 Independence and Uncorrelation.

(1) Suppose  $x$  and  $y$  are two continuous variables, show that if  $x$  and  $y$  are independent, then they are uncorrelated.

$\Rightarrow$  Proof:

Two Random variables  $x, y$

Two random variable ' $x$ ' and ' $y$ ' are independent when their joint probability distribution is a product of their marginal probability distribution (for  $\forall x$  and  $\forall y$ )

$$\text{i.e. } P_{xy}(x,y) = P_x(x) P_y(y) \quad \text{--- (i)}$$

where  $x$  and  $y$  are elements or ~~variables~~ variables  
( $x$  and  $y$  respectively)

Two random variables  $x$  and  $y$  are uncorrelated when their correlation co-efficient is zero.

$$\text{i.e. } \rho_{xy} = \frac{\text{Cov}[x, y]}{\sqrt{\text{Var}[x] \text{Var}[y]}}$$

but even  $\text{Cov}[x, y]$  is zero.

$$\therefore \text{Cov}[x, y] = E[x, y] - E[x]E[y]$$

$$\therefore 0 = E[x, y] - E[x]E[y]$$

$$\therefore E[x, y] = E[x]E[y] \quad \text{--- (ii)}$$

We know that,

$$E[x,y] = \iint x \cdot y \cdot P_{x,y}(u,y) du dy$$

but from equation (i),

$$E[x,y] = \iint x \cdot y \cdot (P_x[x] \cdot P_y[y]) du dy$$

$$= \iint x P_x[x] \cdot y P_y[y] du dy$$

$$= \int x P_x[x] dx \cdot \int y P_y[y] dy$$

$$E[x,y] = E[x] \cdot E[y] \quad \text{--- (iii)}$$

So from (iii) we can conclude that

when  $E[x,y] = E[x] \cdot E[y]$  the two random variables are uncorrelated. And from (i) we

can conclude they are also independent.

Hence proved.

81

- (iii) Suppose  $x$  and  $y$  are uncorrelated, can we conclude  $x$  and  $y$  are independent? If so, prove it, otherwise, give one counter example.

$\Rightarrow$  Let's take an example for two random variables ' $x$ ' and ' $y$ ' with 0 correlation and check if they satisfy independency condition.

10

$$x \sim \text{Uniform } [-1, 1]$$

$$y = x^2$$

Let probability of  $x = -1$  and  $x = 1$  be  $\frac{1}{2}$  and  $\frac{1}{2}$

15

$$\therefore y = x^2$$

$$y = [1, 1]$$

Let probability of  $y = 1$ ,  $y = 1$  be  $\frac{1}{2}$  and  $\frac{1}{2}$

20

$$E[x] = -1 \times \frac{1}{2} + 1 \times \frac{1}{2}$$

$$= 0$$

25

$$E[y] = 1 \times \frac{1}{2} + 1 \times \frac{1}{2}$$

$$= 1$$

$$E[x, y] = 2 \times \frac{1}{2} \times -1 \times \frac{1}{2} + 1 \times \frac{1}{2} = 0$$

$$\therefore E[x, y] = E[x] \cdot E[y] = 0$$

thus  $x$  and  $y$  are uncorrelated.

Now let's check the marginal distributions

$$P[X] = \text{for } x=1 = \frac{1}{2} \quad P[Y] = \text{for } y=1 = \frac{1}{2}$$

$$P[X] = \text{for } x=-1 = \frac{1}{2} \quad "$$

$$P[X, Y] \text{ for } [-1, 1] = \frac{1}{2} \quad P[X, Y] \text{ for } [1, 1] = \frac{1}{2}$$

So for  $P[X, Y] = P[X] \cdot P[Y]$  (if independent)

$$\frac{1}{2} = \frac{1}{2} \times \frac{1}{2}$$

$$\frac{1}{2} \neq \frac{1}{4}$$

As we can conclude from the above example  
not all uncorrelated random variables are necessarily  
independent.

8.2 [Minimum Error Rate Decision] Let  $w_{\max}(u)$  be state of nature for which  $P(w_{\max}|u) \geq P(w_i|u)$  for all  $i = 1, \dots, c$

(i) Show that  $P(w_{\max}|u) \geq \frac{1}{c}$

⇒ AS,

Given  $P(w_{\max}|u) \geq P(w_i|u)$ ,

$$\sum_{i=1}^c P(w_{\max}|u) \geq \sum_{i=1}^c P(w_i|u)$$

but  $\sum_{i=1}^c P(w_i|u) = 1$

$$\therefore c P(w_{\max}|u) \geq 1$$

$$\therefore P(w_{\max}|u) \geq \frac{1}{c} \quad \text{--- (i)}$$

Hence proved.

(ii) Show that for minimum-error-rate decision rule, the average probability of error is given by,

$$P(\text{error}) = 1 - \int P(w_{\max}|u) p(u) du$$

$$P(\text{error}) = \int_{-\infty}^{\omega_{\max}} P(\text{error}|u) p(u) du$$

$$= \int_1^c (1 - P(w_{\max}|u)) p(u) du$$

$$= 1 - \int_1^c P(w_{\max}|u) p(u) du, \quad \text{--- (ii)}$$

Hence proved

(iii) Show that  $P(\text{error}) \leq \frac{c-1}{c}$

From (i) and (ii)

$$\int p(\omega_{\max}(x)) p(x) dx \geq \int \frac{1}{c} p(x) dx$$

$$\therefore 1 - \int p(\omega_{\max}(x)) p(x) dx \leq 1 - \int \frac{1}{c} p(x) dx$$

$$\begin{aligned} P(\text{error}) &\leq 1 - \int \frac{1}{c} p(x) dx \\ &\leq 1 - \frac{1}{c} \end{aligned}$$

$$\therefore P(\text{error}) \leq \frac{c-1}{c}$$

Hence proved.

Q3 [Likelihood Ratio] Suppose we consider two category classification, the class conditionals are assumed to be Gaussian i.e.  $P(x|\omega_1) = N(4,1)$  and  $P(x|\omega_2) = N(8,1)$ , based on prior knowledge, we have  $P(\omega_2) = 1$ . We do not penalize for correct classification, while for misclassification, we put 1 unit penalty for misclassifying  $\omega_1$  to  $\omega_2$  and put 3 unit for misclassifying  $\omega_2$  to  $\omega_1$ . Derive the Bayesian decision rule using likelihood ratio

⇒

$$P(x|\omega_1) = N(4,1)$$

$$P(x|\omega_2) = N(8,1)$$

$$P(\omega_2) = \frac{1}{2}$$

$$\lambda = \begin{bmatrix} \lambda_{11} & \lambda_{12} \\ 0 & 3 \\ -1 & 0 & 0 \\ \lambda_{21} & \lambda_{22} \end{bmatrix}$$

by definition,

$$\frac{P(x|\omega_1)}{P(x|\omega_2)} > \frac{\lambda_{12} - \lambda_{21}}{\lambda_{22} - \lambda_{11}} \quad \frac{P(\omega_2)}{P(\omega_1)}$$

The Posterior Probability =  $P(x|\omega_1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-4)^2}{2}}$  - (ii)  
 (In Gaussian form)

$$= P(x|\omega_2) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-8)^2}{2}}$$
 - (iii)

from (ii) and (iii),

$$\frac{P(x|\omega_1)}{P(x|\omega_2)} = \frac{\frac{1}{\sqrt{2\pi}} e^{-\frac{(x-4)^2}{2}}}{\frac{1}{\sqrt{2\pi}} e^{-\frac{(x-8)^2}{2}}}$$

$$= e^{\frac{-(x-h)^2}{2}} + \frac{(x-8)^2}{2}$$

$$\alpha = \frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}} \frac{P(w_2)}{P(w_1)}$$

$$= \frac{3-0}{13-0} \times \frac{3\pi}{\sqrt{314}} \frac{1/\pi}{314}$$

$$= \frac{1}{13} \times \frac{3\pi}{\sqrt{314}}$$

$$\therefore e^{\frac{(x-h)^2}{2} + \frac{(x-8)^2}{2}} > 1$$

Take logs,

$$\frac{-(x-h)^2}{2} + \frac{(x-8)^2}{2} < 0$$

$$-(x^2 - 8x + 16) + (x^2 - 16x + 64) < 0$$

~~$$-8x + 48 < 0$$~~

~~$$8x - 48 < 0$$~~

$$8x < 48$$

~~$$x < 6$$~~

When  $x$  is smaller than 6 we choose  $w_1$ .

when  $x$  is larger than 6 we choose  $w_2$

(8.4) [Minimum Risk, Reject Option] In many machine learning applications, one has the option either to assign the pattern to one of  $c$  class, or to reject it as being unrecognizable. If the cost for reject is not too high, rejection maybe a desirable action. Let

$$\lambda(\alpha_i | w_j) = \begin{cases} 0, & i=j \text{ and } i, j = 1, \dots, c \\ \lambda_r, & i=c+1 \\ \lambda_s, & \text{otherwise} \end{cases}$$

where  $\lambda_r$  is the loss incurred for choosing the  $(c+1)$ -th action, rejection, and  $\lambda_s$  is the loss incurred for making any substitution ever.

(1) Derive the decision rule with minimum risk.

(2) what happens if  $\lambda_r = 0$ ?

(3) what happens if  $\lambda_r > \lambda_s$ ?

(1) For  $i=1, 2, 3, \dots, c$

$$\begin{aligned} R(\alpha_i | x) &= \sum_{j=1}^c \lambda(\alpha_i | w_j) P(w_j | x) \\ &= \lambda_s \sum_{j=1, j \neq i}^c P(w_j | x) \\ &= \lambda_s [1 - P(w_i | x)] \end{aligned}$$

for  $i=c+1$

$$R(\alpha_{c+1} | x) = \lambda_r$$

$\therefore$  the minimum risk is achieved if we decide  $w_i$  if

$$R(\alpha_i | x) \leq R(\alpha_{c+1} | x)$$

i.e.  $P(w_i | x) \geq 1 - \frac{\lambda_r}{\lambda_s}$  and reject otherwise.

2. If  $\lambda_r = 0$ , we always reject ~~the~~ the ~~class~~ class

3. If  $\lambda_r > \lambda_v$ , we always accept the class.

Q.5 [Maximum Likelihood Estimation] Suppose we have training samples  $x_1, x_2, \dots, x_n$ . Consider the following distributions.

5 (ii) Exponential density:  $f(x_i; \alpha) = \alpha e^{-\alpha x_i}$ ,  $x_i \geq 0$ ,  $\alpha > 0$   
Find MLE for  $\alpha$

⇒

10 Likelihood function

$$L(\alpha, x) = \prod_{i=1}^n f(x_i; \alpha)$$

Take log likelihood to maximize,

$$\ln L(\alpha, x) = \sum_{i=1}^n \ln f(x_i; \alpha)$$

$$\ln L(\alpha, x) = \sum_{i=1}^n \ln (\alpha e^{-\alpha x_i})$$

$$\ln L(\alpha, x) = \sum_{i=1}^n -(\alpha x_i) + n \ln \alpha$$

$$\ln L(\alpha, x) = - \left( \sum_{i=1}^n x_i \right) \alpha + n \ln \alpha$$

$$\frac{d}{d\alpha} \ln L(\alpha, x) = - \left( \sum_{i=1}^n x_i \right) + \frac{n}{\alpha}$$

..... - (Taking derivatives)

But,

$$\frac{d}{d\alpha} \ln L(\alpha, x) = 0$$

$$\therefore 0 = - \left( \sum_{i=1}^n x_i \right) + \frac{n}{\alpha}$$

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

$$\therefore \bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

5

$$\therefore \text{MLE of } \alpha = \hat{\alpha} = \frac{n}{\sum_{i=1}^n x_i} //$$

10 (2) Uniform distribution  $\text{Unif}(\alpha_1, \alpha_2)$ , find MLE for  $\alpha_1, \alpha_2$

$$f(x_i | \alpha_1, \alpha_2) = \begin{cases} \frac{1}{\alpha_2 - \alpha_1} & \text{for } \alpha_1 \leq x_i \leq \alpha_2 \\ 0 & \text{otherwise} \end{cases}$$

∴ Likelihood function has the form.

$$20 L(\alpha_1, \alpha_2) = \begin{cases} \frac{1}{(\alpha_2 - \alpha_1)^n} & \text{for } \alpha_1 \leq x_i \leq \alpha_2 \quad (i=1, \dots, n) \\ 0 & \text{otherwise.} \end{cases}$$

25 We can see that for the given function we need to maximize  $(\alpha_2 - \alpha_1)^n$ , which is a decreasing function.

30 ∴  $(\alpha_2 - \alpha_1)^n$  should be minimum.  
i.e.  $\alpha_2$  should be ~~maximum~~ such that  $\alpha_2 \geq x_i$   
 $\alpha_1$  should be maximum such that  $x_i \leq \alpha_1$

so it can be concluded that

MLE of  $\alpha_2$  is  $\hat{\alpha}_2 = \max [u_1, \dots, x_n]$

5 MLE of  $\alpha_1$  is  $\hat{\alpha}_1 = \min [u_1, \dots, x_n]$

10

15

20

25

30

Q.6 [Logistic Regression, MLE] In ~~this~~ this problem, you need to use MLE to derive and build a logistic regression classifier (suppose the target/response  $y \in \{0, 1\}$ ):

- (1) Suppose the classifier is  $y = \mathbf{x}^T \boldsymbol{\alpha}$ , where  $\boldsymbol{\alpha}$  contains the weight as well as bias parameter  
 - The log-likelihood function is  $LL(\boldsymbol{\alpha})$ , what is  $\frac{\partial LL(\boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}}$

$\Rightarrow$  Logistic regression is actually a classification method, LR ~~introduces~~ introduces an extra non-linearity over a linear classifier.

$f(u) = \mathbf{x}^T \boldsymbol{\alpha}$  by using a logistic function

$$\mathbf{x}^T \boldsymbol{\alpha} = m \sum_{i=1}^m x_i \alpha_i = x_1 \alpha_1 + \dots + x_m \alpha_m$$

$$\sigma(f(u)) = 1/(1+e^{-f(u)}) \quad \text{i.e. } \sigma(z) = 1/(1+e^{-z})$$

This is called the logistic function

Assume that,

$$p(y=1|\mathbf{x}; \boldsymbol{\alpha}) = \sigma(\mathbf{x}^T \boldsymbol{\alpha})$$

$$p(y=0|\mathbf{x}; \boldsymbol{\alpha}) = 1 - \sigma(\mathbf{x}^T \boldsymbol{\alpha})$$

Then the likelihood (assuming data independence) is,

$$p(Y|\mathbf{x}; \boldsymbol{\alpha}) \sim \prod_{i=1}^n [\sigma(\mathbf{x}^T \boldsymbol{\alpha})]^{y_i} [1 - \sigma(\mathbf{x}^T \boldsymbol{\alpha})]^{1-y_i}$$

The log likelihood equation is,

$$LL(\boldsymbol{\alpha}) = \sum_{i=1}^n y_i \log \sigma(\mathbf{x}^T \boldsymbol{\alpha}) + (1-y_i) \log [1 - \sigma(\mathbf{x}^T \boldsymbol{\alpha})]$$

$$\frac{\partial LL(\alpha)}{\partial \alpha_j} = \frac{\partial LL(\alpha)}{\partial p} \cdot \frac{\partial p}{\partial \alpha_j} = [y_{1p} - (1-y_{11}-p)] \cdot \sigma'(z)$$

$$= [y_{1p} - y_{11}] \cdot p \frac{1-p}{1-p} \cdot x_j$$

$$= \left[ \frac{y(1-p) - y(p-1)}{(p-1)(1-p)} \times (1-p) \right] x_j$$

$$= [y - p] x_j$$

$$\therefore \frac{\partial LL(\alpha)}{\partial \alpha_j} = n \sum_{i=1}^n [y_i - \sigma(x^T \alpha)] x_{ij}$$

$$\frac{\partial LL(\alpha)}{\partial \alpha_3} = \sum_{i=1}^n [y_i - \sigma(x^T \alpha)] x_{ij}$$