

HW-1

1) Provide an intuitive example to show that $P(A|B)$ and $P(B|A)$ are in general not the same.

→ Here to prove $P(A|B) \neq P(B|A)$ in general lets take an example of someone having COVID-19 and Testing positive.

So, Here our two events are:

A: Testing V+ve

B: Having COVID-19.

Suppose that COVID-19 test is 80% accurate. Which means that probability testing positive while having COVID-19 is 80%.

$$\text{i.e } P(A|B) = 0.8$$

But, alternatively there might be only 15% chance of having being tested positive and having COVID-19 as well. This happens due to "false positive".

$$\text{i.e } P(B|A) = 0.15$$

Hence it is proved that in general $P(A|B) \neq P(B|A)$.

(2) (i) Here X and Y are two Continuous Random Variables
And it's also given that ~~X and Y are independent~~

~~$P(X=x, Y=y) = P(X=x) \cdot P(Y=y)$~~

$$\text{So, } P(X=x, Y=y) = P(X=x) \cdot P(Y=y)$$

$$\left. \begin{aligned} E(X, Y) &= \iint_{\mathbb{R}^2} x \cdot y \cdot P(X=x, Y=y) dx dy \\ &= \int_{\mathbb{R}} x \cdot P_x(x) dx \cdot \int_{\mathbb{R}} y \cdot P_y(y) dy \\ &= \int x P_x(x) dx \cdot \int y P_y(y) dy \end{aligned} \right\}$$

$$\begin{aligned} \text{Now } E(X, Y) &= \iint_{\mathbb{R}^2} x \cdot y \cdot P(X=x, Y=y) dx dy \\ &= \int_{\mathbb{R}} x \cdot P(X=x) dx \cdot \int_{\mathbb{R}} y \cdot P(Y=y) dy \\ &= \cancel{\int x P(X=x) dx \cdot \int y P(Y=y) dy} \\ &= \int x P(X=x) dx \cdot \int y P(Y=y) dy \end{aligned}$$

$$E(X, Y) = E[X] \cdot E[Y] \text{ by def.}$$

~~$E(X, Y) - E[X] \cdot E[Y] = 0$~~

$$\text{Now coefficient constant} = \frac{\text{cov}[X, Y]}{\sqrt{\text{var}[X] \cdot \text{var}[Y]}}$$

$$= \frac{E(X, Y) - E[X] \cdot E[Y]}{\sqrt{\text{var}[X] \cdot \text{var}[Y]}}$$

$$= 0 \quad (\because E(X, Y) - E[X] \cdot E[Y] = 0)$$

Hence proved.

(2) Suppose X and Y are uncorrelated, can we conclude X and Y are independent? If so, prove it or give one counter example.

No to prove that $\uparrow X \sim \text{Normal}(0, 1)$

assume

$$\text{and } Y = X^2$$

$$\text{and } E(X) = 0 \text{ as } u=0 \text{ and } E(X^3) = 0$$

i.e. Y is dependent upon X . (third moment mean)

And if we still be able to prove $\text{cov}(X, Y) = 0$ then it will be clear that X and Y are uncorrelated but not independent.

$$\text{cov}(X, Y) = \cancel{\text{def}} - E[(X - \bar{x})(Y - \bar{y})]$$

from question (1) simplification we can write

$$\begin{aligned}\text{cov}(X, Y) &= E(XY) - E(X)E(Y) \\ &= E(X^3) - E(X)E(Y) \quad \text{i.e. } Y = X^2 \\ &= 0 - 0 \cdot E(Y) \\ &= 0\end{aligned}$$

Hence X and Y are uncorrelated

But from $Y = X^2$ we can say that they are dependent

Hence proved that if X and Y are uncorrelated then they may not be independent.

3)

[Minimum Probability of Error, Discriminant function]

Let the components of the vector $x = [x_1, \dots, x_d]^T$ be binary valued (0,1), and let $P(\omega_j)$ be the prior probability for the state of nature ω_j and $j = 1, \dots, c$. We define

$$P_{ij} = P(x_i=1 | \omega_j), i=1, \dots, d, j=1, \dots, c$$

With the components x_i being statistical independent for all x in ω_j . Show the minimum probability of error is achieved by the following decision rule.

Decide ω_k if $g_k(x) \geq g_j(x)$ for all j and k where

$$g_j(x) = \sum_{i=1}^d x_i \ln \frac{P_{ij}}{1-P_{ij}} + \sum_{i=1}^d \ln(1-P_{ij}) + \ln P(\omega_j)$$

→ Here we can use formula of the posterior probability to derive the $g_j(x)$

Let's consider it as following function

$$\begin{aligned} g_j(x) &= p(x/\omega_j) \cdot P(\omega_j) \\ &= \ln[p(x/\omega_j) \cdot P(\omega_j)] \text{ taking } \ln. \\ &= \ln p(x/\omega_j) + \ln P(\omega_j) \quad \cdots [1] \end{aligned}$$

\hookrightarrow density \hookrightarrow prior

~~Here we write it is given that x are in ω_j~~

Here it is given that x are statistically independent for all x in ω_j

So we can write the the density as product over $i=1$ to d .

$$\textcircled{1} \quad p(x|w_j) = \prod_{i=1}^d p(x_i|w_j)$$

Now we can convert above equation in terms of p_{ij}

$$= \prod_{i=1}^d p_{ij}^{x_i} (1-p_{ij})^{1-x_i} \quad \text{--- [2]}$$

Plug in [2] in [1]

$$g_j(x) = \sum_{i=1}^d \ln [p_{ij}^{x_i} (1-p_{ij})^{1-x_i}] + \ln p(w_j)$$

$$= \sum_{i=1}^d [x_i \ln p_{ij} + (1-x_i) \ln (1-p_{ij})] + \ln p(w_j)$$

$$= \sum_{i=1}^d x_i \ln \left(\frac{p_{ij}}{1-p_{ij}} \right) + \sum_{i=1}^d \ln (1-p_{ij}) + \ln p(w_j)$$

Hence derived.

4) [Likelihood Ratio] Suppose we consider two category classification, the class condition 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 $p(\omega_2) = 1/4$. We do not penalize ~~for~~ for correct classification, while for misclassification, we put 1 unit penalty for misclassifying ω_1 ~~and~~ ω_2 and put 3 unit for misclassifying ω_2 to ω_1 . Derive bayesian decision rule using likely hood ratio.

→ Here we are given following informations

$$\begin{aligned} p(x|\omega_1) &= N(4, 1) & p(\omega_2) &= 1/4 \\ p(x|\omega_2) &= N(8, 1) & p(\omega_1) &= 3/4 \end{aligned}$$

$$\text{and penalty } \lambda = \begin{bmatrix} \lambda_{11} & \lambda_{12} \\ \lambda_{21} & \lambda_{22} \end{bmatrix}$$

Now we know that

~~$$N = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$~~

So, from above equation we can write

$$p(x|\omega_1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-4)^2}{2}} \quad \dots [1]$$

$$p(x|\omega_2) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-8)^2}{2}} \quad \dots [2]$$

From Likelihood ratio test we also know following is true

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

Now plug in value of λ_{12} & λ_{22} in above formula

$$\frac{1/\sqrt{\pi} e^{-(x-4)^2/2}}{1/\sqrt{\pi} e^{-(x-8)^2/2}} > \frac{3-0}{1-0} \times \frac{1}{4} \times \frac{1}{3}$$

$$e^{-\frac{(x-4)^2}{2} + \frac{(x-8)^2}{2}} > 1$$

$$\log \left(e^{-\frac{(x-4)^2}{2} + \frac{(x-8)^2}{2}} \right) < \log 11$$

$$-\frac{(x-4)^2}{2} + \frac{(x-8)^2}{2} < 0$$

$$-(x-4)^2 + (x-8)^2 < 0$$

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

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

$$-8x + 48 < 0$$

$$48 < 8x$$

$$x > 6$$

So if $x > 6$ then we can conclude following

if $x < 6$ then we go for ω_1 .
otherwise we go for ω_2 .

5) [Minimum Risk, Reject option] In many machine learning applications, one has either option to assign the pattern to one of c classes, or to reject it as being unrecognizable. If cost of reject is λ_s not too high rejection may be desirable action. Let

$$\lambda(\alpha_i|\omega_i) = \begin{cases} 0, & i=j \text{ and } i,j=1,2,\dots,c \\ \lambda_r, & i=c+1 \\ \lambda_s, & \text{otherwise.} \end{cases}$$

When λ_r is loss function incurred for choosing the $(c+1)$ th action, rejection, and λ_s is the loss incurred for making any substitution error.

(1) Derive the decision rule with minimum risk.

~~Minimum Risk~~

$$\text{Risk is defined as } R(\alpha_i|x) = \sum_{j=1}^c \lambda(\alpha_i|\omega_j) P(\omega_j|x)$$

Now for $i = 1, 2, \dots, c$

$$\begin{aligned} \text{So } R(\alpha_i|x) &= \sum_{j=1}^c \lambda(\alpha_i|\omega_j) P(\omega_j|x) \\ &= \sum_{\substack{j=1 \\ j \neq i}}^c 0 + \sum_{\substack{j=1 \\ j \neq i}}^c \lambda_s P(\omega_j|x) \\ &= 0 + \lambda_s \sum_{\substack{j=1 \\ j \neq i}}^c P(\omega_j|x) \end{aligned}$$

$$= \lambda_s [1 - p(\omega_i|x)]$$

$$\boxed{= \lambda_s [1 - p(\omega_i|x)]}$$

Now for $i=c+1$

$$R(\alpha_{c+1} | x) = \lambda_s \quad (\because \text{given})$$

So, we can achieve min risk if we decide ω_i if

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

Here we know $R(\omega_i | x)$ and $R(\alpha_{c+1} | x)$

$$\text{So } \lambda_s [1 - p(\omega_i | x)] \leq \lambda_s$$

$$1 - p(\omega_i | x) \leq \lambda_s / \lambda_s$$

$$p(\omega_i | x) \geq 1 - \lambda_s / \lambda_s$$

(2) What happens if $\lambda_s = 0$

if $\lambda_s = 0$ then $p(\omega_i | x) \geq 1$

And we always REJECT.

(3) What happens if $\lambda_s > \lambda_s$?

then $p(\omega_i | x) > 1$

And we will ~~NEVER~~ NEVER REJECT.

6) [Maximum Likelihood Estimation (MLE)] A general representation of a exponential family given by the following probability function:

$$p(x|\eta) = h(x) \exp \{ \eta^T T(x) - A(\eta) \}$$

η is natural parameter.

$h(x)$ is the base density which ensure x in right space.

$T(x)$ is the sufficient statistics.

$A(\eta)$ is the log normalizer which is determined by $T(x) \cdot h(x)$
 $\exp(\cdot)$ represent the exponential function.

(1) Write down the expression of $A(\eta)$ in terms of $T(x)$ & $h(x)$

$$p(x|\eta) = h(x) \cdot \exp \{ \eta^T T(x) - A(\eta) \}$$

Taking limit on the both sides. in term of x .

$$\int_x p(x|\eta) = \int h(x) \cdot \exp \{ \eta^T T(x) - A(\eta) \}$$

Now we can assign ~~any~~ integration of probability to 1.

$$1 = \int_x h(x) \cdot \exp \{ \eta^T T(x) - A(\eta) \}$$

Taking log on both side

$$\log(1) = \log \left[\int_x h(x) \cdot \exp \{ \eta^T T(x) - A(\eta) \} \right]$$

$$O = \log \left[\int_{\mathcal{X}} h(x) \cdot \exp(\eta^T T(x)) \right] - A(n)$$

$$O = \log \int_{\mathcal{X}} h(x) \cdot \exp(\eta^T T(x)) + \log (\exp(-A(n)))$$

$\because A(n)$ constant to \int

$$O = \log \int_{\mathcal{X}} h(x) \cdot \exp(\eta^T T(x)) - A(n)$$

$$A(n) = \log \int_{\mathcal{X}} h(x) \cdot \exp(\eta^T T(x))$$

(2) Show that $\frac{\partial}{\partial n} A(n) = E_n T(x)$ where $E_n(\cdot)$ is the expectation w.r.t $p(x|n)$

$$A_n = \log \int_{\mathcal{X}} (h(x) \cdot \exp(\eta^T T(x)))$$

$$= e \log \cancel{\int_{\mathcal{X}} h(x)} \cdot \cancel{T(x)}$$

$$\frac{\partial A(n)}{\partial n} = \left[\frac{\partial}{\partial n} \int_{\mathcal{X}} h(x) \exp(\eta^T T(x)) dx \right] \cdot \frac{1}{\int_{\mathcal{X}} h(x) \exp(\eta^T T(x))}$$

$$= \left[\frac{\partial}{\partial n} \int_{\mathcal{X}} h(x) \cdot \exp(\eta^T T(x)) \cdot \overset{(T(x))}{dx} \right] \cdot \frac{1}{A(n)}$$

$$= \int_x h(x) \cdot \exp \{ n^T T(x) \} T(x) dx$$

$\because 1/e^{A(n)}$ is independent \int

$$= \int_x h(x) \cdot \exp \{ n^T T(x) - A(n) \}^2 T(x) dx$$

$$= \int_x p(x/n) \cdot T(x) dx$$

$$= E_n T(x)$$

(3) Suppose we have n i.i.d samples x_1, x_2, \dots, x_n derive the maximum likelihood estimator for η .

$$p(x_i | n) = \prod p(x_i | n) \quad [\because \text{prob (3)}]$$

$$p(x_i | n) = h(x_i) \exp \{ n^T T(x_i) - A(n) \}^2$$

$$\text{Hence } p(x | n) = \prod_i h(x_i) \exp \{ n^T \sum_i T(x_i) - n A(n) \}^2$$

Now lets take Log of Likelihood.

$$\log (n: x_1, x_2, \dots, x_n) = \log (p(x_1, x_2, \dots, x_n | \eta))$$

$$= \log (h(x_1, \dots, x_n) + \eta^T \sum_i T(x_i) - n A(n))$$

Taking gradients on both side

$$\frac{\partial}{\partial n} \log(n: x_1, x_2, \dots, x_n) = \frac{\partial}{\partial n} \left[h(x_1, x_2, \dots, x_n) + \eta^T \sum T(x_i) \right] - \eta A(n)$$

$$0 = 0 + \sum_i T(x_i) - n \frac{\partial A(n)}{\partial n}$$

$$0 = \sum_i T(x_i) - n \frac{\partial A(n)}{\partial n}$$

$$E_n(T(x)) = \frac{\sum_i T(x_i)}{n}$$

$$\begin{aligned} y_n &= \sum_i \frac{T(x_i)}{n} \\ &= \boxed{\frac{\partial A(n)}{\partial n}} \end{aligned}$$

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

$$y = \mathbf{x}^T \boldsymbol{\alpha}$$

$$\text{So, } \sigma(y) = \frac{1}{1 + e^{-\mathbf{x}^T \boldsymbol{\alpha}}} \quad [\text{Logistic function}]$$

Now let's assume following:

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

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

Then likely hood fun will be

$$\text{Log Likelihood} = \prod_{i=1}^N \sigma(\mathbf{x}_i^T \boldsymbol{\alpha})^{y_i} (1 - \sigma(\mathbf{x}_i^T \boldsymbol{\alpha}))^{1-y_i}$$

$$p(y|\mathbf{x}; \boldsymbol{\alpha}) = \prod_{i=1}^N y_i^{y_i} (1 - \sigma(\mathbf{x}_i^T \boldsymbol{\alpha}))^{1-y_i}$$

$$= \prod_{i=1}^N [\sigma(\mathbf{x}_i^T \boldsymbol{\alpha})]^{y_i} [1 - \sigma(\mathbf{x}_i^T \boldsymbol{\alpha})]^{1-y_i}$$

Log likely hood

$$LL(\alpha) = \sum_{i=1}^N y_i \log \sigma(x^T \alpha) + (1-y_i) \log (1-\sigma(x^T \alpha))$$

Taking gradient on both side

$$\frac{\partial LL(\alpha)}{\partial \alpha} = \frac{-y_i}{\sigma(x^T \alpha)} \frac{\partial}{\partial \alpha} \sigma(x^T \alpha) - (1-y_i) \frac{1}{(1-\sigma(x^T \alpha))} \cdot \frac{\partial}{\partial \alpha} (1-\sigma(x^T \alpha))$$

$$= \frac{-y_i}{\sigma(x^T \alpha)} \cdot \sigma(x^T \alpha) (1-\sigma(x^T \alpha)) \cdot x^T$$

$$- \frac{(1-y_i)}{(1-\sigma(x^T \alpha))} \cdot (\sigma(x^T \alpha)(1-\sigma(x^T \alpha))) x^T$$

• Taking x^T common

$$= -x^T [y_i - y_i \sigma(x^T \alpha) - \sigma(x^T \alpha) + \sigma(x^T \alpha)]$$

$$= -x^T [y_i - \sigma(x^T \alpha)]$$

Finally $\frac{\partial LL(\alpha)}{\partial \alpha} = \sum_{i=1}^N [y_i - \sigma(x^T \alpha)] x_i^T$