# Statistical Data Analysis Problem sheet 3

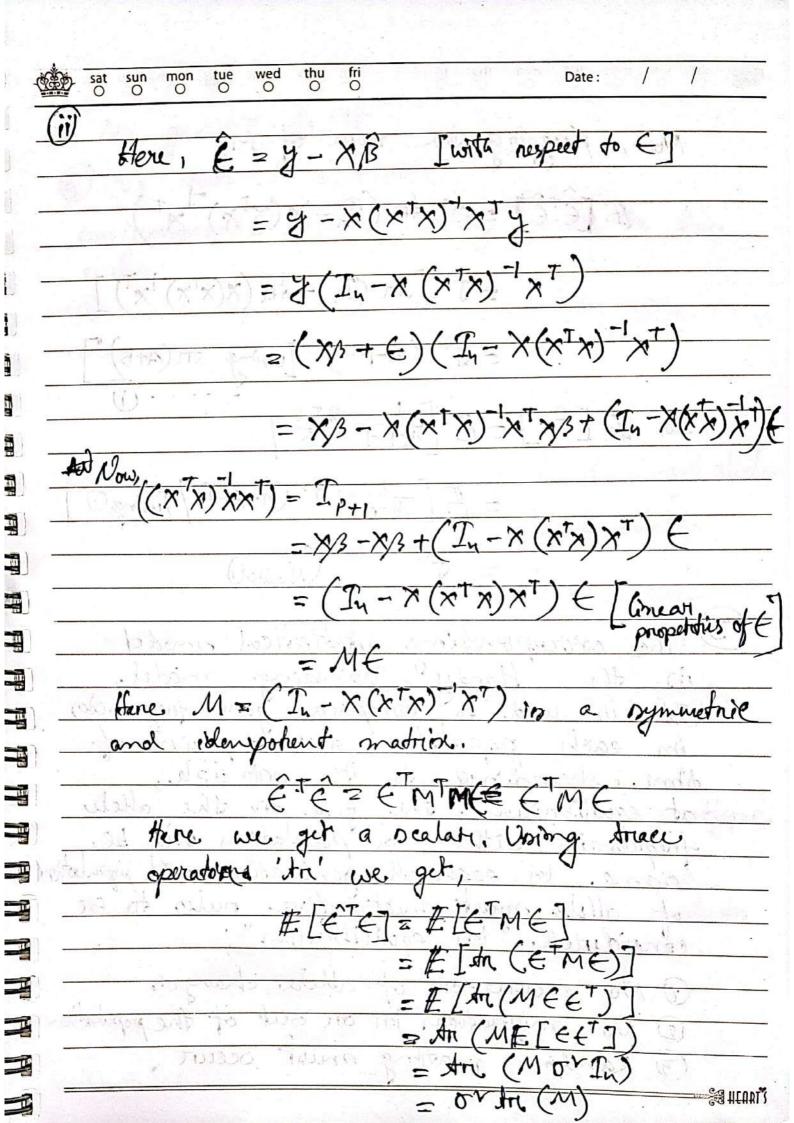## 1. Exercise 1

```
[37]    1 #importing all the libraries
        2 import pandas as pd
        3 import numpy as np
```

```
[38]    1 #getting the X values as a dataframe
        2 dfx = pd.read_csv('/content/drive/MyDrive/1-DS/X.txt', sep=",", header=None)
```

```
[39]    1 #getting the Y values as a dataframe
        2 dfy = pd.read_csv('/content/drive/MyDrive/1-DS/Y.txt', header= None)
```

Now we need to convert them into matrix

```
[40]    1 #converting to numpy array to get matrix
        2 x = dfx.to_numpy()
        3 y = dfy.to_numpy()
```

```
[41]    1 #adding an identity column to the x to equalize the full rank p + 1
        2 x = np.insert(x, 0, 1.0, axis=1)
```

**As the LS estimator is equivalent to the ML estimator based on the maximization of the log likelihood, we can estimate the beta hat from it.**

```
[42]    1 #beta hat estimating
        2 xt = np.transpose(x)
        3 xtx = xt.dot(x)
        4
        5 beta = np.linalg.inv(xtx).dot(xt).dot(y)
        6 beta
```

```
[42]    1 #beta hat estimating
        2 xt = np.transpose(x)
        3 xtx = xt.dot(x)
        4
        5 beta = np.linalg.inv(xtx).dot(xt).dot(y)
        6 beta
```

```
array([[-0.00800698],
       [ 0.88161162],
       [-2.45938171],
       [-0.97715699]])
```

- *Here we get B0, B1, B2 and B3 for the multiple linear regression model*

- As the mean is 0 and the constant variance is 1 which gives us that N is a standard normal distrbution.

```
        1 #now we can get the sigma hat square in the multiple linear regression model
        2 sgms =  (np.transpose(y-x.dot(beta))*(y-x.dot(beta))) / 201
        3 sgms
```

```
array([[ 1.43756944e-04, -2.60179376e-04,  4.88288715e-05, ...,
        -8.27929177e-04, -5.93215518e-04, -6.80968604e-04],
       [-2.60179376e-04,  4.70887219e-04, -8.83732288e-05, ...,
         1.49843263e-03,  1.07363470e-03,  1.23245515e-03],
       [ 4.88288715e-05, -8.83732288e-05,  1.65853462e-05, ...,
        -2.81216658e-04, -2.01493183e-04, -2.31299633e-04],
       ...,
       [-8.27929177e-04,  1.49843263e-03, -2.81216658e-04, ...,
         4.76823381e-03,  3.41646407e-03,  3.92185420e-03],
       [-5.93215518e-04,  1.07363470e-03, -2.01493183e-04, ...,
         3.41646407e-03,  2.44791410e-03,  2.81002872e-03],
       [-6.80968604e-04,  1.23245515e-03, -2.31299633e-04, ...,
         3.92185420e-03,  2.81002872e-03,  3.22571018e-03]])
```

```
[44]    1 #now we can also get the adjusted estimator of the variance hat
        2 sgad = ((np.transpose(y).dot(y))- (np.transpose(beta).dot(xt).dot(y))) / (201 - 3 - 1)
        3 sgad
```
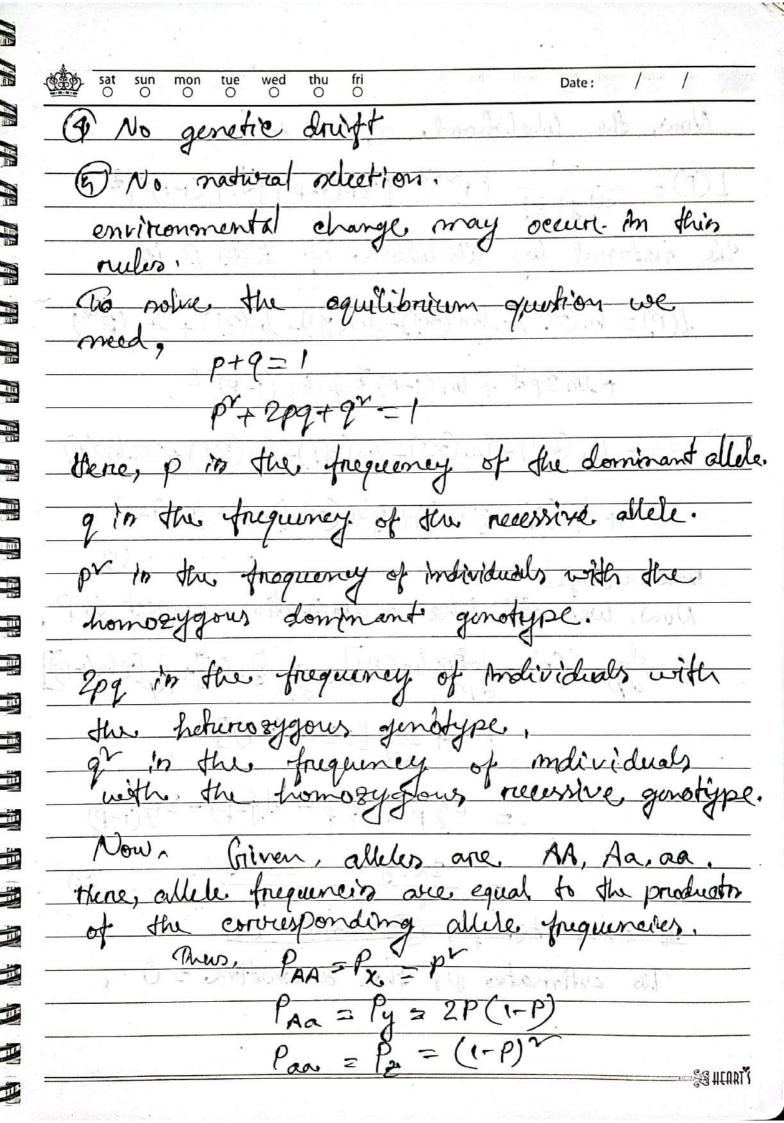
```
array([[0.97422819]])
```

② ⓘ Here,

Least squares estimator $\hat{\beta} = (X^TX)^{-1}X^Ty$

and REML estimator $\hat{\sigma}^2_{ad} = \frac{1}{n-p-1}\hat{\epsilon}^T\hat{\epsilon}$

Now, $Cov(\hat{\beta}) = Cov((X^TX)^{-1}X^Ty)$ [Plugging $\hat{\beta}$ value] ⓘ

We know, $Cov(Xy) = E[(Xy - E[Xy])(Xy - E[Xy])^T]$

$$= E[X(y - E[y]) - (y - E[y])X^T]$$

$$= X E[(y - E[y])(y - E[y])^T]X^T$$

$$= X Cov(y)X^T \text{ [putting value of } Cov(y)]$$ ⓘⓘ

By using ⓘⓘ in ⓘ, we get,

$$Cov(\hat{\beta}) = (X^TX)^{-1}X^T Cov(y)[(X^TX)^{-1})X^T]^T$$

$$= (X^TX)^{-1}X^T Cov(y)X[(X^TX)^{-1}]^T$$

$$= (X^TX)^{-1}X^TX[(X^TX)^T]^{-1} Cov(y)$$

$$= (X^TX)^{-1}(X^TX)(X^TX)^{-1} Cov(y)$$

$$= (X^TX)^{-1} Cov(y)$$

$$\therefore Cov(\hat{\beta}) = \sigma^2(X^TX^{-1}) \quad [Cov(y) = \sigma^2 I]$$

(showed)

HEART

(ii)

Here, $\hat{\epsilon} = y - X\hat{\beta}$    [with respect to $\epsilon$]

$$= y - X(X^TX)^{-1}X^Ty$$

$$= y(I_n - X(X^TX)^{-1}X^T)$$

$$= (X\beta + \epsilon)(I_n - X(X^TX)^{-1}X^T)$$

$$= X\beta - X(X^TX)^{-1}X^TX\beta + (I_n - X(X^TX)^{-1}X^T)\epsilon$$

Now, $((X^TX)^{-1}XX^T) = I_{p+1}$

$$= X\beta - X\beta + (I_n - X(X^TX)X^T)\epsilon$$

$$= (I_n - X(X^TX)X^T)\epsilon \quad \left[\begin{array}{l}\text{linear} \\ \text{properties of }\epsilon\end{array}\right]$$

$$= M\epsilon$$

Here $M = (I_n - X(X^TX)^{-1}X^T)$ is a symmetric and idempotent matrix.

$$\hat{\epsilon}^T\hat{\epsilon} = \epsilon^TM^TM\epsilon = \epsilon^TM\epsilon$$

Here we get a scalar. Using trace operators 'tr' we get,

$$\mathbb{E}[\hat{\epsilon}^T\epsilon] = \mathbb{E}[\epsilon^TM\epsilon]$$

$$= \mathbb{E}[tr(\epsilon^TM\epsilon)]$$

$$= \mathbb{E}[tr(M\epsilon\epsilon^T)]$$

$$= tr(M\mathbb{E}[\epsilon\epsilon^T])$$

$$= tr(M\sigma^2I_n)$$

$$= \sigma^2 tr(M)$$

Now, plugging the value of M,

$$\mathbb{E}[\hat{\epsilon}^{\top}\hat{\epsilon}] = \sigma^{2} \, \text{tr}\left(\mathbb{I}_{n} - X(X^{\top}X)^{-1}X^{\top}\right)$$

$$= \sigma^{2}\left[\text{tr}(\mathbb{I}_{n}) - \text{tr}(X(X^{\top}X)^{-1}X^{\top})\right]$$

$$= \sigma^{2}(n-p-1) \quad \left[\text{using } \text{tr}(A+B)\right] \cdots \cdots (1)$$

So, $\mathbb{E}[\hat{\sigma}^{2}] = \mathbb{E}\left[\frac{1}{n-p-1}\hat{\epsilon}^{\top}\hat{\epsilon}\right]$

$$= \mathbb{E}\left[\frac{1}{n-p-1}\cdot \sigma^{2}(n-p-1)\right] \left[\text{using } (1)\right]$$

$$= \sigma^{2} \quad \text{(showed)}$$

(3) The corresponding statistical model is the Hardy - Weinberg model. It is used to compare gene frequencies in each population over a period of time. According to it's principle, at equilibrium, the gene or the allele frequencies within a population will be same, in ~~each~~ all generations. A population of alleles must meet five rules to be considered "in equilibrium":

① No occurances of allele changes.
② No migration in or out of the population.
③ Random mating must occur.

④ No genetic drift

⑤ No natural selection.

environmental change, may occur. In this rules.

To solve the equilibrium question we need,

$$p + q = 1$$

$$p^2 + 2pq + q^2 = 1$$

Here, p is the frequency of the dominant allele.

q is the frequency of the recessive allele.

$p^2$ is the frequency of individuals with the homozygous dominant genotype.

$2pq$ is the frequency of individuals with the heterozygous genotype,

$q^2$ is the frequency of individuals with the homozygous recessive genotype.

Now, Given, alleles are AA, Aa, aa.

Here, allele frequencies are equal to the products of the corresponding allele frequencies.

Thus, $P_{AA} = P_x = p^2$

$$P_{Aa} = P_y = 2P(1-P)$$

$$P_{aa} = P_z = (1-P)^2$$

HEART

Now, the likelihood of P is:

$$L(P) = \frac{n!}{x! \cdot y! \cdot z!} (p^2)^x \left[2P(1-P)\right]^y \cdot \left[(1-P)^2\right]^z$$

The natural log likelihood of $\cancel{\mathcal{L}(P)}$ P is

$$\ell(P) = \ln(n!) - \ln(x!) - \ln(y!) - \ln(z!) + \ln(p^{2x})$$

$$+ \ln 2p^y + \ln(1-P)^y + \ln(1-P)^{2z}$$

$$= \ln(n!) - \ln(x!) - \ln(y!) - \ln(z!) + 2x\ln(P)$$

$$+ \ln(2) + y\ln(P) + y\ln(1-P) + 2z\ln(1-P) \quad \cdots\cdots (1)$$

Now, we will take a derivative respect to P.

$$\therefore \frac{d}{dp}\ell(P) = \frac{d}{dp}\left[2x\ln(P)\right] + \frac{d}{dp}\left[y\ln(P)\right] + \frac{d}{dp}\left[y\ln(1-P)\right]$$

$$+ \frac{d}{dp}\left[2z\ln(1-P)\right]$$

$$= 2x/p + y/p - y/(1-P) - 2z/(1-P)$$

$$= \frac{2x+y}{P} - \frac{2z+y}{1-P} \quad \cdots\cdots\cdots (ii)$$

To estimate p, the derivative = 0.

$$\therefore \quad \frac{2x+y}{P} - \frac{2z+y}{1-P} = 0$$

or, $\dfrac{2x+y}{P} = \dfrac{2z+y}{1-P}$

or, $(1-P)(2x+y) = 2zP + Py$

or, $2x+y - 2xp - py - py = 2zp$

or, $2p(x+y+z) = 2x+y$

or, $P = \dfrac{2x+y}{2(x+y+z)} \qquad [\because x+y+z = n]$

$\therefore$ Maximum likelihood estimator for $p = \dfrac{2x+y}{2n}$.