

CH5115 - Parameter and State Estimation  
Assignment 3

Question 1

Given,

$$f(y) = \begin{cases} e^{-(y-\theta)}, & x > \theta, -\infty < \theta < \infty \\ 0, & \text{otherwise.} \end{cases}$$

Part (a)

$$T_N = 2 \min(Y_N)$$

Let random variable  $Z = 2Y$ , and the pdf becomes

$$f_Z(z) = \begin{cases} \frac{1}{2} e^{-\left(\frac{z}{2}-\theta\right)}, & z > 2\theta \\ 0, & \text{otherwise} \end{cases}$$

And the statistic  $T_N$  becomes

$$T_N = \min(Z_N)$$

To determine the pdf of  $T_N$ , the 1<sup>st</sup> step is to find the cdf of  $Z$ .  $f_{Z_N}$

$$F_Z(z) = \int_{2\theta}^z f_Z(z) dz = \int_{2\theta}^z \frac{1}{2} e^{-\left(\frac{z}{2}-\theta\right)} dz$$

which gives,

$$P_{r_Z}(Z < z) = F_Z(z) = \begin{cases} 1 - e^{-(\frac{z}{2} - \theta)} & , \text{ when } z > 2\theta \\ 0 & , \text{ otherwise} \end{cases}$$

Considering  $z$  as a fixed point and given that we have  $N$  observations of  $Z_N$ , the cdf of  $T_N$  will be the Probability that at least one observation is less than  $z$ .

$$\begin{aligned} F_{T_N}(z) &= 1 - \Pr_z(Z_i > z \ \forall i = 1, 2, \dots, N) \\ &= 1 - (1 - F_Z(z))^N \\ &= 1 - e^{-N(\frac{z}{2} - \theta)}, \quad \text{when } z > 2\theta \end{aligned}$$

Pdf can be obtained by differentiating  $F_{T_N}$ .

$$f_{T_N}(z) = \frac{N}{2} e^{-N(\frac{z}{2} - \theta)}, \quad z > 2\theta$$
$$0, \quad \text{otherwise.}$$

where  $z = 2y$

### Part 1b)

We have derived that , for  $z \geq 2\theta$  , the pdf of  $T_N$  is

$$f_{T_N}(z) = \frac{N}{2} e^{-N(\frac{z}{2} - \theta)}$$

$$E(T_N) = \int_{2\theta}^{\infty} z f(z; \theta) dz$$

$$= \frac{N}{2} \int_{2\theta}^{\infty} z e^{-N(z/2 - \theta)} dz$$

$$= \frac{N}{2} \left[ z \int e^{-N(z/2 - \theta)} dz - \int \left( \int e^{-N(z/2 - \theta)} dz \right) dz \right]_{2\theta}^{\infty}$$

$$= \frac{N}{2} \left[ -\frac{z}{N} e^{-N(z/2 - \theta)} + \frac{2}{N} \int e^{-N(z/2 - \theta)} dz \right]_{2\theta}^{\infty}$$

$$= \left. -\frac{z}{N} e^{-N(z/2 - \theta)} \right|_{2\theta}^{\infty} - \left. \frac{2}{N} e^{-N(z/2 - \theta)} \right|_{2\theta}^{\infty}$$

$$= 2\theta + \frac{2}{N}$$

$$\text{Therefore bias } \Delta T_N = E(T_N) - \theta = \theta + \frac{2}{N}$$

To correct the bias, we obtain statistic such that

$$T_N' = \frac{T_N}{2} - \frac{1}{N} = \min(y_N) - \frac{1}{N}, \text{ where } E(T_N') = \theta$$

### Part (c)

In order to prove  $T_N'$  converges to  $\theta$  in probability, we have to show,

$$\lim_{N \rightarrow \infty} \Pr(|T_N' - \theta| > \epsilon) = 0$$

Consider  $T_N' > \theta$

$$\begin{aligned}\Pr(T_N' - \theta > \epsilon) &= \Pr(T_N' > \theta + \epsilon) \\ &= 1 - \Pr(T_N' < \theta + \epsilon) \\ &= 1 - F_{T_N'}(\theta + \epsilon) = 1 - F_{T_N'}(\theta + \epsilon)\end{aligned}$$

To derive the pdf and cdf of  $T_N'$ ,

$$f_{T_N'}(z; \theta) = f_{T_N}(z, \theta - 1/N) = N e^{-N(z - (\theta - 1/N))}, z$$

when  $z > \theta - 1/N$   
0, otherwise

$$\begin{aligned}F_{T_N'}(z; \theta) &= \int_{\theta - 1/N}^z e^{-N(z - (\theta - 1/N))} dz \\ &= -e^{-N(z - (\theta - 1/N))} \Big|_{\theta - 1/N}^z \\ &= 1 - e^{-N(z - (\theta - 1/N))} \\ F_{T_N'}(\theta + \epsilon) &= 1 - e^{-N(\epsilon + 1/N)}\end{aligned}$$

$$\Pr(T_N - \theta > \epsilon) = 1 - F_{T_N}(\theta + \epsilon)$$

$$= e^{-N(\epsilon + 1/N)}.$$

$$\lim_{N \rightarrow \infty} \Pr(T_N' - \theta > \epsilon) = \lim_{N \rightarrow \infty} e^{-N(\epsilon + 1/N)} = 0$$

when  $T_N'$

When  $\theta - \frac{1}{N} < T_N' < \theta$

$$\Pr(\theta - T_N' > \epsilon) = \Pr(T_N' < \theta - \epsilon) = F_{T_N'}(\theta - \epsilon)$$

$$F_{T_N'}(\theta - \epsilon) = \begin{cases} 1 - e^{-N(-\epsilon + 1/N)}, & \epsilon < 1/N \\ 0, & \epsilon > 1/N. \end{cases}$$

As  $N$  increases to infinity,  $1/N \rightarrow 0$  and  $\epsilon > 1/N$ ,

$$\lim_{N \rightarrow \infty} \Pr(\theta - T_N' > \epsilon) = 0, \text{ when } \theta - \frac{1}{N} < T_N' < \theta$$

Hence,  $\lim_{N \rightarrow \infty} \Pr(|T_N' - \theta| > \epsilon) = 0$ , this implies the statistic  $T_N'$  converges to the true value as  $N$  tends to infinity.

## Question 2

Given

$$V[n] \triangleq V(f_n) = \sum_{k=0}^{N-1} v[k] \exp(-j2\pi f_n k), \quad f_n = \frac{n}{N}$$

$$\left| \frac{V(f_n)}{N} \right|^2 = \frac{a_n^2 + b_n^2}{N}$$

$$a_n = \operatorname{Re}(V[n]), \quad b_n = \operatorname{Im}(V[n])$$

This gives

$$a_n = \sum_{k=0}^{N-1} v[k] \cos(2\pi f_n k)$$

$$b_n = - \sum_{k=0}^{N-1} v[k] \sin(2\pi f_n k)$$

Part (a)

We invoke 2 properties of Gaussian distributions

$$\textcircled{1} \quad \text{If } X \sim N(\mu, \sigma^2) \Rightarrow (\alpha X) \sim N(\mu, (\alpha\sigma)^2)$$

$$\textcircled{2} \quad \text{If } X \sim (\mu_x, \sigma_x^2) \text{ and } Y \sim (\mu_y, \sigma_y^2) \text{ and } X, Y \text{ are independent, } (X+Y) \sim N(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$$

When

$$(a_n = \sum_{k=0}^{N-1} v[k] \cos(2\pi f_n k))$$

Given  $v[k] \sim \text{GWN}(0, 1)$ , let  $u[k] = v[k] \cos(2\pi f_n k)$   
where  $u[k] \sim N(0, \cos^2(2\pi f_n k))$

$$\Rightarrow a_n = \sum_{k=0}^{N-1} u[k]$$

$$\text{III}^{\text{by}} \quad w[k] = u[k] \sin(2\pi f_n k), \quad w[k] \sim N(0, \sin^2(2\pi f_n k))$$

$$b_n = \sum_{k=0}^{N-1} w[k]$$

AS the processes  $u[k]$  and  $w[k]$  are Gaussian distributed with 0 mean, and each of the process has ~~an~~ 0 autocovariance for non zero lags\* (even though it is variance non stationary), we can say that  ~~$a_n, b_n$~~  both  $a_n$  and  $b_n$  are Gaussian distributed, as they are the sum of uncorrelated Gaussian random variables.

B. The mean and variance of  $a_n$  is given by

$$\mu_{a_n} = \sum_{k=0}^{N-1} \mu_{u[k]} = 0$$

$$\sigma_{a_n}^2 = \sum_{k=0}^{N-1} \text{var}(u[k]) = \sum_{k=0}^{N-1} \cos^2 2\pi f_n k$$

III<sup>by</sup>

$$\mu_{b_n} = 0$$

$$\sigma_{b_n}^2 = \sum_{k=0}^{N-1} \sin^2(2\pi f_n k)$$

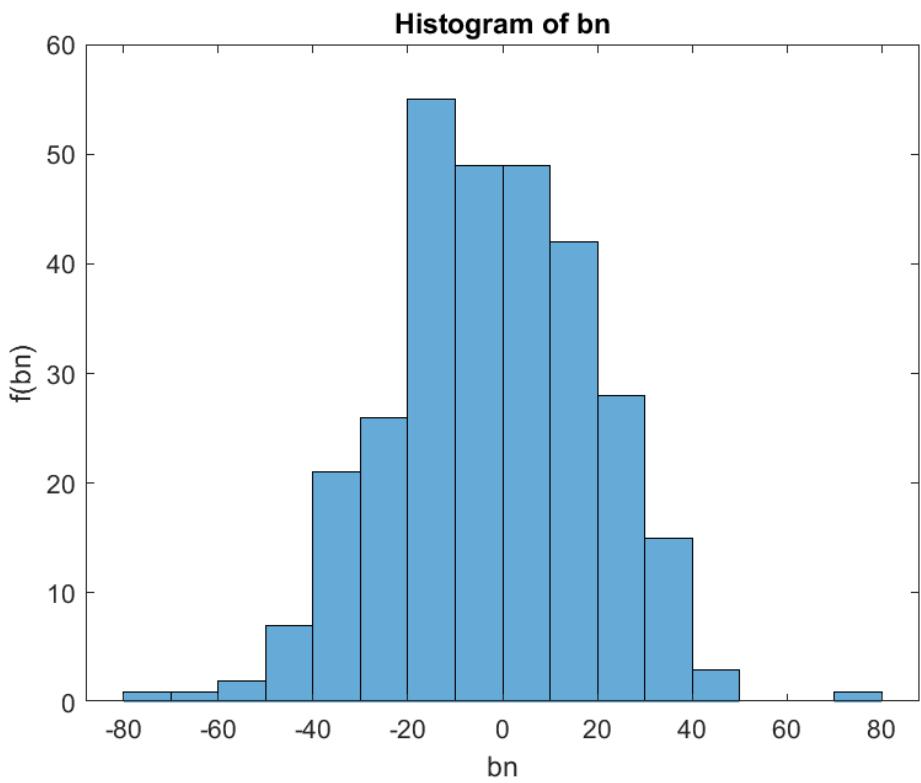
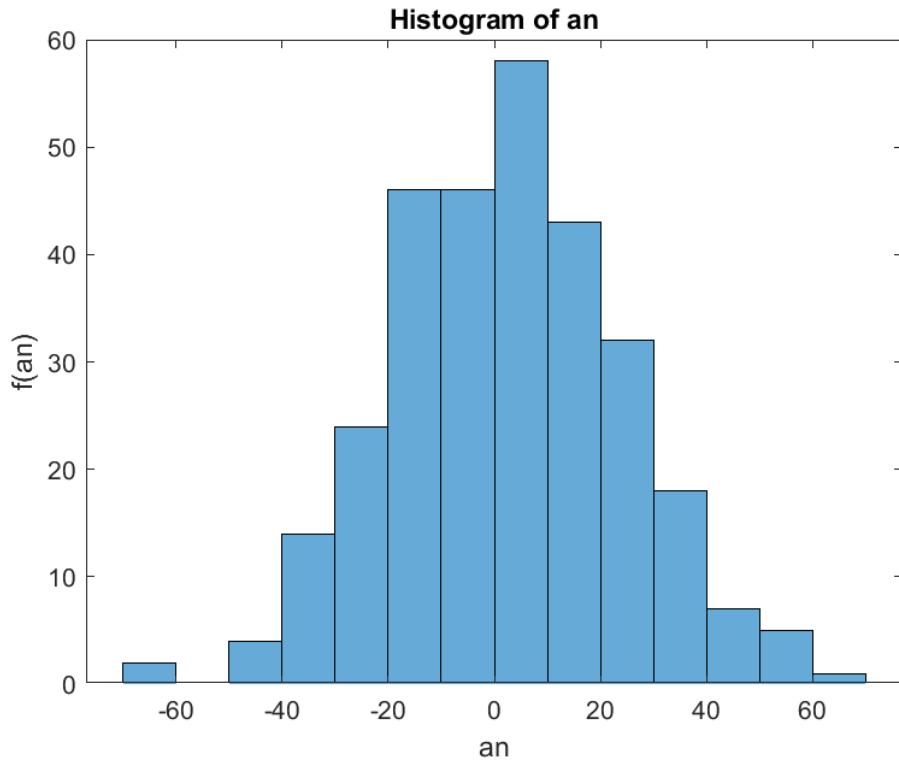
\*Autocovariance of  $u[k]$

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} u[k] u[k+l] f(u[k], u[k+l]) du[k] du[k+l] = 0$$

~~III<sup>by</sup>~~ III<sup>by</sup> for  $w[k]$

### Question 2(a)

The histograms of  $a_n$  and  $b_n$  is shown below, where  $n$  was chosen to be 12



These visualizations suggest that  $a_n$  and  $b_n$  are Gaussian distributed. This was further verified by performing Kolmogorov-Smirnov Test on the samples, where the null hypothesis (that the standardized version follows a standard normal distribution) was not rejected in both cases.

The details of the sample statistics are given below:

$$\bar{a}_n = 2.2657$$

$$\bar{b}_n = -2.0661$$

$$s_a^2 = 473.1549$$

$$s_b^2 = 464.7418$$

$$\text{cov}(a_n, b_n) = 4.9168$$

These estimates might not match with the theoretical values because of small number of samples and realizations.

### Part (b)

$$\rho_{anbn} = \frac{\text{cov}(an, bn)}{\sigma_{an} \sigma_{bn}} = \frac{E(anbn) - E(an)E(bn)}{\sigma_{an} \sigma_{bn}}$$

We have seen,  $E(an) = E(bn) = 0$ .

$\Rightarrow \cancel{\rho_{anb}}$

$$\rho_{anbn} = \frac{E(anbn)}{\sigma_{an} \sigma_{bn}}$$

$$anbn = \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} v[i]v[j] \cos(2\pi f_n i) \sin(2\pi f_n j)$$

Again,  $E(v[i]v[j] \cos(2\pi f_n i) \sin(2\pi f_n j)) = 0 \quad \forall i, j = 0 \text{ to } N-1$

$$\begin{aligned} \Rightarrow E(anbn) &= \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} E(v[i]v[j] \cos(2\pi f_n i) \sin(2\pi f_n j)) \\ &= 0 \end{aligned}$$

Hence, the correlation between  $a_n$  and  $b_n$

$$\rho_{anbn} = 0$$

### Part (c)

We are given

$$P(f_n) = \frac{a_n^2 + b_n^2}{N}$$

Since,  $a_n, b_n$  are normally distributed,  
then  $a_n^2, b_n^2$  are  $\chi^2$  distributed (when standardized)

As,  $\& C_n = \frac{2P_{vv}(f_n)}{\gamma_{vv}(f_n)}$

$\& C_n$  is also  $\chi^2$  distributed ~~non-standard~~

$$\gamma_{vv}(f_n) = \sum_{l=-\infty}^{\infty} \sigma_{vv}[l] \exp(-j2\pi f l)$$

Since  $\sigma_{vv}[l] = 0 \quad \forall l \neq 0$  (V[K] is GWN)

$$\gamma_{vv}(f_n) = 1$$

$$\Rightarrow C_n = 2P_{vv}(f_n)$$

$$\begin{aligned} E(C_n) &= 2E(P_{vv}(f_n)) = 2E\left(\frac{a_n^2 + b_n^2}{N}\right) \\ &= \frac{2}{N} (E(a_n^2) + E(b_n^2)) \\ &= \frac{2}{N} (\text{var}(a_n) + \text{var}(b_n)) \end{aligned}$$

$$= \frac{2}{N} \left( \sum_{k=0}^{N-1} \cos^2(2\pi f_n k) + \sum_{k=0}^{N-1} \sin^2(2\pi f_n k) \right)$$

$$= \frac{2}{N} (N) = 2$$

Hence, mean of  $e_{jn} = E(e_{jn}) = 2$

For a  $\chi^2$  distribution, since  $E(a_n) = E(b_n) = 3$

$$\text{var}(e_{jn}) = 2E(e_{jn}) = 4$$

Hence, the variance of  $e_{jn} = E(e_{jn}^2) - E(e_{jn})^2 = 4$

### Part (d)

We know that, the true value of PSD

$$\gamma_{rr}(f) = 1$$

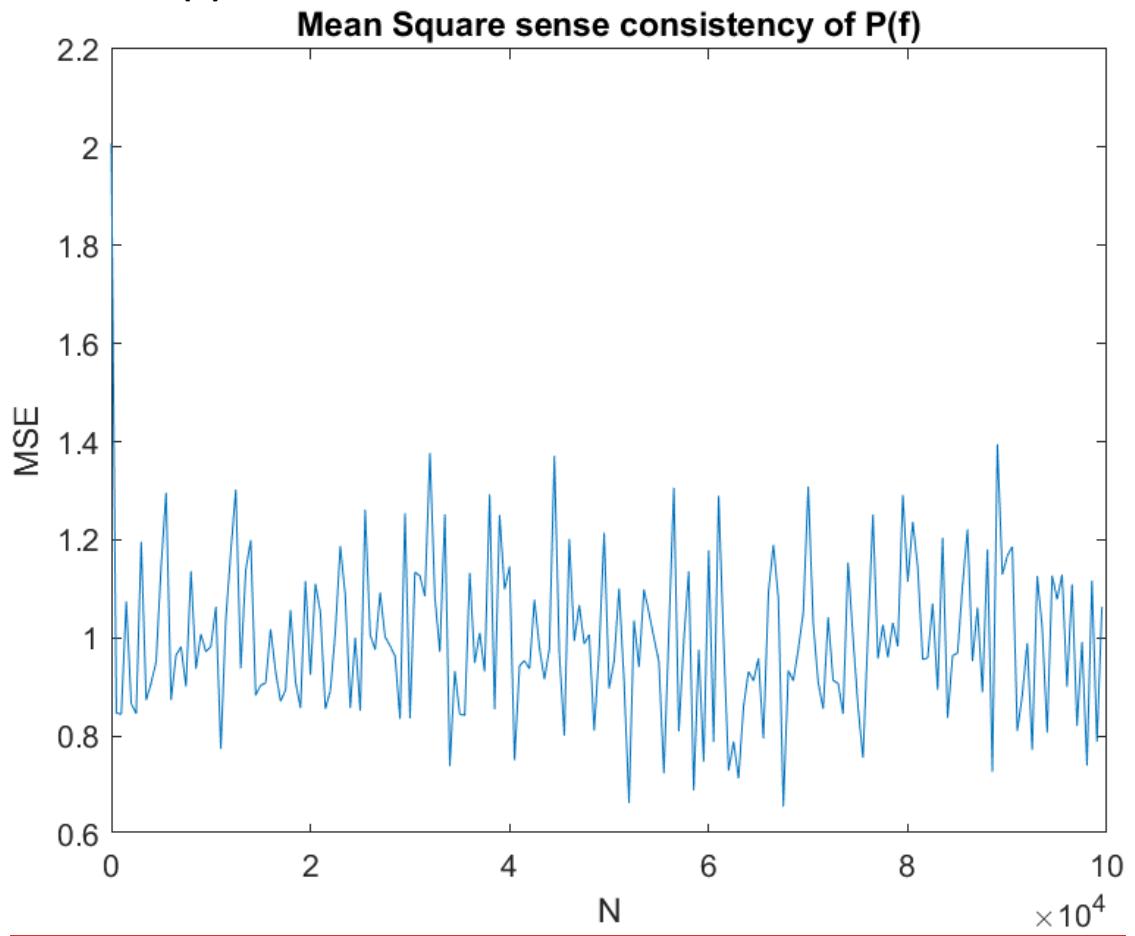
$P(f_n)$  is mean square consistent iff

$$\lim_{N \rightarrow \infty} E((P(f_n) - 1)^2) = 0$$

$$E((P(f_n) - 1)^2) = E(P_m^2) - 1$$

$$= \text{var}(P_m) = \frac{\text{var}(e_{jn})}{4} = \frac{4}{4} = 1 \neq 0$$

Therefore, since the MSE is 1 irrespective of  $f$ ,  $P(f_n)$  is not a consistent estimator of  $\gamma_{rr}(f_n)$  in mean square sense.

**Question 2(d)**

It can be observed that the MSE of estimator does not converge to 0, and hence the estimator  $P(f_n)$  is not consistent for estimating  $\gamma_{vv}(f_n)$ . This is in accordance with what was proved theoretically.

### Question 3

#### Part (a)

The random process is given by,

$$y[k] = A \sin(2\pi f_0 k) + e[k], \quad e[k] \sim N(0, \sigma_e^2)$$

where both  $A$  (amplitude) and frequency ( $f_0$ ) are unknown.

First, we determine Fisher information of  $A$  and  $f_0$ .

Pdf of  $y_N \rightarrow N$  observations

$$f(y_N; A, f_0) = \frac{1}{(2\pi\sigma_e^2)^{N/2}} \exp\left(-\frac{1}{2} \sum_{k=0}^{N-1} \frac{(y[k] - A \sin(2\pi f_0 k))^2}{\sigma_e^2}\right)$$
$$= l(A, f_0; y_N) \rightarrow \text{likelihood}$$

Loglikelihood

$$L(A, f_0; y_N) = -\frac{N}{2} \ln(2\pi\sigma_e^2) - \frac{1}{2} \sum_{k=0}^{N-1} \frac{(y[k] - A \sin(2\pi f_0 k))^2}{\sigma_e^2}$$

Score of  $A$

$$S_A = \frac{\partial L}{\partial A} = \cancel{2\pi f_0} \frac{\sum_{k=0}^{N-1} (y[k] - A \sin(2\pi f_0 k)) \sin(2\pi f_0 k)}{\sigma_e^2}$$

$$\frac{\partial S_A}{\partial A} = - \frac{\sum_{k=0}^{N-1} \sin^2(2\pi f_0 k)}{\sigma_e^2}$$

$$\Rightarrow I(A) = -E\left(\frac{\partial S_A}{\partial A}\right) = -E\left(-\sum_{k=0}^{N-1} \frac{\sin^2(2\pi f_0 k)}{\sigma_e^2}\right)$$

$$I(A) = \frac{\sum_{k=0}^{N-1} \sin^2(2\pi f_0 k)}{\sigma_e^2}$$

Score of  $B f_0$

$$S_{f_0} = \frac{\partial L}{\partial B} \frac{\partial L}{\partial f_0} = \frac{2\pi \sum_{k=0}^{N-1} A_k \cos(2\pi f_0 k) (y[k] - A \sin(2\pi f_0 k))}{\sigma_e^2}$$

$$\frac{\partial S_{f_0}}{\partial f_0} = -\frac{4\pi^2 A}{\sigma_e^2} \sum_{k=0}^{N-1} k^2 (y[k] \sin(2\pi f_0 k) + A \cos^2(2\pi f_0 k) - A \sin^2(2\pi f_0 k))$$

$$I(f_0) = -E\left(\frac{\partial S_{f_0}}{\partial f_0}\right) = \frac{4\pi^2 A^2}{\sigma_e^2} \sum_{k=0}^{N-1} k^2 \cos^2(2\pi f_0 k)$$

Score function corresponding to cross-derivative of the parameter.

$$S_{Af_0} = \frac{\partial L}{\partial A \partial f_0} = \frac{2\pi}{\sigma_e^2} \sum_{k=0}^{N-1} (y[k] \cos(2\pi f_0 k) - 2A \sin(2\pi f_0 k) \cos(2\pi f_0 k))$$

$$I(A, f_0) = -E(S_{Af_0}) = \frac{2\pi A}{\sigma_e^2} \sum_{k=0}^{N-1} k \sin(2\pi f_0 k) \cos(2\pi f_0 k)$$

Fisher Information matrix

$$I(\theta) = \begin{bmatrix} I(A) & I(A, f_0) \\ I(A, f_0) & I(f_0) \end{bmatrix}$$

$$\therefore (\mathbb{I}(\theta))^{-1} = \frac{1}{\mathbb{I}(A) \mathbb{I}(f_0) - \mathbb{I}(A, f_0)^2} \begin{bmatrix} \mathbb{I}(f_0) & -\mathbb{I}(A, f_0) \\ -\mathbb{I}(A, f_0) & \mathbb{I}(A) \end{bmatrix}$$

$$CRLB(A) = (\mathbb{I}(\theta))_{1,1}^{-1}$$

$$= \frac{\mathbb{I}(f_0)}{\mathbb{I}(A) \mathbb{I}(f_0) - \mathbb{I}(A, f_0)^2}$$

$$= \frac{\sigma e^2 \sum_{k=0}^{N-1} k^2 \cos^2(2\pi f_0 k)}{\sum_{k=0}^{N-1} \sin^2(2\pi f_0 k) \sum_{k=0}^{N-1} k^2 \cos^2(2\pi f_0 k) - \left( \sum_{k=0}^{N-1} k \sin(2\pi f_0 k) \cos(2\pi f_0 k) \right)^2}$$

Similarly,

$$CRLB(f_0) = (\mathbb{I}(\theta))_{2,2}^{-1}$$

$$= \frac{\mathbb{I}(A)}{\mathbb{I}(A) \mathbb{I}(f_0) - \mathbb{I}(A, f_0)^2}$$

$$= \frac{\sigma e^2}{4\pi^2 A^2} \frac{\sum_{k=0}^{N-1} \sin^2(2\pi f_0 k)}{\sum_{k=0}^{N-1} \sin^2(2\pi f_0 k) \sum_{k=0}^{N-1} k^2 \cos^2(2\pi f_0 k) - \left( \sum_{k=0}^{N-1} k \sin(2\pi f_0 k) \cos(2\pi f_0 k) \right)^2}$$

However, when ~~only~~ one of the parameters is known and the other is unknown, the CLR of the estimator of the unknown parameter is only dependent on the Fisher Information of that parameter. Here,

$$CRLB(A) = (\mathbb{I}(A))^{-1} \text{ and } CRLB(f_0) = (\mathbb{I}(f_0))^{-1}$$

In the case when  $A$  is unknown and  $f_0$  is known,

$$CRLB(A) = I(A)^{-1} = \frac{\sigma_e^2}{\sum_{k=0}^{N-1} \sin^2(2\pi f_0 k)}$$

When  $f_0$  is unknown &  $A$  and  $A$  is known,

$$CRLB(f_0) = I(f_0)^{-1} = \frac{\sigma_e^2}{4\pi A^2 \sum_{k=0}^{N-1} k^2 \cos^2(2\pi f_0 k)}$$

It can be noticed that the CRLBs for amplitude and frequency estimates in the case when both the parameters are unknown, are different from that when derived in the case when one of the parameters is known.

As ~~the~~ both parameters are unknown, the corresponding variance lower bound for each parameter is highly dependent on the other parameter as well.

### Part (b)

Given GWN process

$$y[k] = e[k], \quad e[k] \sim N(0, \sigma_e^2)$$

we can rewrite it as

$$y[k] = \sigma u[k], \quad u[k] \sim N(0, 1)$$

In such a setting, the parameter  $\sigma$  which is to be estimated is associated with the noise and is hence not linear, i.e., it cannot be expressed as  $y = \sigma + v$ .

This suggests that BLUE estimator cannot be designed directly for this process. However, we can log-transform the data to obtain a linear relationship in  $\sigma$ .

Let

$$z[k] = \log(|y[k]|) = \log(\sigma u[k]) \\ = \log \sigma + \log(|u[k]|)$$

Here, we take the absolute value and log transform the data.

$$z[k] = \sigma' + v[k]$$

$$\text{where } \sigma' = \log(\sigma), \quad v[k] = \log(|u[k]|)$$

Also,

$$E(v[k]) = E(\log(|u[k]|)) = \int_0^{\infty} \log(u) (2\pi N(0, 1)) du \\ = \int_0^{\infty} \log u \left( \frac{2}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right) \right) du \\ = -0.6352 \neq 0$$

(Integration performed using MATLAB).

To correct the nonzero expectation of noise, we define another noise process,  $w[k] = v[k] + 0.6352$ , such that  $E(w[k]) = 0$ . Therefore,

$$z[k] = \sigma'' + w[k]$$

$$\text{where } \sigma'' = \sigma' - 0.6352$$

Rewriting the process in vector notation,

$$\mathbf{z}_{N \times 1} = \mathbf{l}_{N \times 1}^T \boldsymbol{\sigma}'' + \mathbf{w}_{N \times 1}$$

where  $\mathbf{l} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}_{N \times 1}$

constraint :  $\mathbf{a}^T \mathbf{l} = 1$

objective function:  $\min_{\mathbf{a}} \mathbf{a}^T \sum \mathbf{w} \mathbf{a}$ ,

$$= \min_{\mathbf{a}} \mathbf{a}^T \mathbf{Q} \mathbf{a}, \quad \text{as } \mathbf{Q} \text{ is symmetric}$$

Solving, we get.

$$\hat{\mathbf{a}}^* = (\mathbf{l}^T \sum \mathbf{w} \mathbf{l})^{-1} \mathbf{l}^T \sum \mathbf{w}$$

and  $\hat{\sigma}''_{\text{BLUE}} = \hat{\mathbf{a}}^{*\top} \mathbf{Z} = \left( \frac{N}{\sum \mathbf{w}^2} \right)^{-1} (\sum \mathbf{w}^2) \sum_{k=0}^{N-1} Z[k]$

$$\hat{\sigma}''_{\text{BLUE}} = \frac{\sum_{k=0}^{N-1} Z[k]}{N}$$

$$\hat{\sigma}'_{\text{BLUE}} = \hat{\sigma}''_{\text{BLUE}} + 0.6352$$

$$= \frac{\sum_{k=0}^{N-1} Z[k]}{N} + 0.6352$$

$$\hat{\sigma}_{\text{BLUE}} = e^{\hat{\sigma}'_{\text{BLUE}}} = \exp\left(\frac{\sum_{k=0}^{N-1} \log(Y[k])}{N} + 0.6352\right)$$

## Question 4

### Part (a)

Given  $N = 100$  (number of samples)

$m = 14578$  (sample mean)

$s = 1845$  (sample standard deviation)

Let the true mean be  $\mu$

To determine  $Pr(12000 < \mu < 16000)$ , we can determine  $Pr(\mu < 16000) - Pr(\mu < 12000)$ .

To

We construct two t statistics  $t_1$  and  $t_2$ , where

$$t_1 = \frac{12000 - m}{s/\sqrt{N}} = \frac{12000 - 14578}{1845/10} = -13.9729$$

$$t_2 = \frac{16000 - m}{s/\sqrt{N}} = \frac{16000 - 14578}{1845/10} = 7.7073$$

and we assume

$T_t \sim t_{\text{dist}}(v)$ , where  $v = N-1 = 99$   
(degrees of freedom)

$$\begin{aligned}
 P_r(12000 > M) &= tCDF(t_1, \nu) \\
 &= tCDF(-13.9729, 99) \\
 &= 1.87 \times 10^{-25} \\
 &\approx 0.
 \end{aligned}$$

$$\begin{aligned}
 P_r(16000 > M) &= tCDF(t_2, \nu) \\
 &= tCDF(7.7073, 99) \\
 &= 99.999\dots \\
 &\approx 1
 \end{aligned}$$

$$\text{Hence, } P_r(12000 < M < 16000) = 1$$

Therefore, with 100% confidence, we can say that the average molecular weight of the polymer lies between 12000 and 16000.

### Part (b)

Sample mean of institution A =  $\bar{x}_1 = 85.2$

$$n_1 = 60, S_1 = 6.8$$

Sample mean of institution B =  $\bar{x}_2 = 87.2$

$$n_2 = 55, S_2 = 8.8$$

One tailed 2-sample hypothesis test with unknown standard deviation.

We choose  $\alpha = 0.05$ , (95% confidence).

Let  $M_1$  and  $M_2$  be the true means of institutions A and B respectively.

$$H_0 : M_1 = M_2 \quad (\text{Precisely, } M_1 \leq M_2)$$

$$H_a : M_1 > M_2$$

$$t_s = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}} = -1.355$$

We reject null hypothesis when

$$t_s > t_{1-\alpha, v}$$

$$\text{where } v = \frac{\left( \frac{S_1^2}{n_1} + \frac{S_2^2}{n_2} \right)^2}{\left( \frac{S_1^2}{n_1} \right)^2 / (n_1 - 1) + \left( \frac{S_2^2}{n_2} \right)^2 / (n_2 - 1)} = 101.469$$

$$t_{1-\alpha, v} = 1.66$$

As  $t_{\text{obs}} < t_{1-\alpha/2, v}$ , we do not reject the null hypothesis, and thereby with 95% confidence, we reject the claim that students from elite institution A perform better than students from elite institution B. In other words, we fail to reject the null hypothesis with 95% confidence.

## Question 5

### Part (a)

From Q1, we have

$$f(y) = \begin{cases} e^{-(y-\theta)}, & y > \theta \\ 0, & \text{otherwise} \end{cases}$$

For  $N$  observations  $y_N$

$$f(y_N) = \begin{cases} e^{-\sum_{k=0}^{N-1} (y[k] - \theta)}, & y[k] > \theta \ \forall k \\ 0, & \text{otherwise} \end{cases}$$

Consider step function,

$$H(x) = \begin{cases} 1, & x \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

We can rewrite the pdf of  $y_N$  as

$$\begin{aligned} f(y_N; \theta) &= H(\min(y_N) - \theta) e^{-\sum_{k=0}^{N-1} (y[k] - \theta)} \\ &= \Phi(H(\frac{y_N}{2} - \theta)) e^{-\sum_{k=0}^{N-1} (y[k] - \theta)} \\ &= \Phi(\frac{y_N}{2} - \theta) \\ &= H(\frac{y_N}{2} - \theta) e^{N\theta} \cdot e^{-\sum_{k=0}^{N-1} y[k]} \end{aligned}$$

$$\phi(T_N; \theta) = H\left(\frac{T_N}{2} - \theta\right) e^{N\theta}$$

$$K(y_N) = e^{-\sum_{k=0}^{N-1} y_k}$$

$$\Rightarrow f(y_N; \theta) = \phi(T_N; \theta) K(y_N)$$

Hence, by Neyman-Fisher factorization theorem, we can say  $T_N$  is a sufficient statistic for estimating  $\theta$ .

However, as we saw in Q1.b), the estimator  $T_N$  is biased with  $E(T_N) = 2\theta + \frac{2}{N}$ .

To correct the bias, we derived  $T_N'$  as

$$T_N' = \frac{T_N}{2} - \frac{1}{N} = \min(y_N) - \frac{1}{N}.$$

Therefore, we can say that the estimator  $T_N'$  is sufficient, unbiased and hence complete.

By Rao-Blackwell theorem,  $T_N'$  is the MVUE estimator of  $\theta$ ,

$$\therefore \hat{\theta}_{MVUE} = T_N' = \min(y_N) - \frac{1}{N}$$

### Question 5 (b)

#### Part (b)

We derived that  $T_N' = \min(y_N) - \frac{1}{N}$  is the minimum unbiased estimator for  $\theta$ ,

$$\begin{aligned} \text{Var}(\hat{\theta}^*) &= \text{Var}(T_N') = E(T_N'^2) - E(T_N')^2 \\ &= \left[ \left( \theta - \frac{1}{N} \right)^2 + \frac{2\theta}{N} \right] - \theta^2 \\ &= \theta^2 + \frac{1}{N^2} \cancel{+ 2\theta} - \theta^2 \\ &= \cancel{\frac{1}{N}} \left( 2\theta + \frac{1}{N} \right) = \frac{1}{N^2} \\ &= \left( \frac{1}{N} + 1 \right)^2 - 1 \end{aligned}$$

$$\Rightarrow \text{var}(\hat{\theta}^*) = 1/N^2$$

For the purpose,  $\text{var}(\hat{\theta}^*)$  was assigned to be the variance of  $T_N'$  across R realizations.

It was therefore shown that efficiency

$$n_{\hat{\theta}_1} = \frac{\text{var}(\hat{\theta}^*)}{\text{var}(\hat{\theta}_1)} = 9.2604 \times 10^{-4}$$

$$n_{\hat{\theta}_2} = \frac{\text{var}(\hat{\theta}^*)}{\text{var}(\hat{\theta}_2)} = 9.3133 \times 10^{-4}$$

where  $\hat{\theta}_1 = \bar{y} - 1$ ,  $\hat{\theta}_2 = \tilde{y} - \ln 2$

Since  $\eta_{\hat{\theta}_2} > \eta_{\hat{\theta}_1}$ ,  $\hat{\theta}_2$  is more efficient than  $\hat{\theta}_1$ .

NOTE: This result differed each time the simulation was performed. A seed value of 42 was therefore set for this simulation.

**Question 5(b)**

$$\widehat{\theta}_1 = \bar{y} - 1$$

$$\widehat{\theta}_2 = \tilde{y} - \ln(2)$$

With a random state of 42, the efficiencies of both estimators were calculated as:

$$\eta_{\theta_1} = \frac{var(T'_N)}{var(\widehat{\theta}_1)} = 9.2604 \times 10^{-4}$$

$$\eta_{\theta_2} = \frac{var(T'_N)}{var(\widehat{\theta}_2)} = 9.3133 \times 10^{-4}$$

It can be observed that  $\widehat{\theta}_2$  is a comparatively more efficient estimator than  $\widehat{\theta}_1$