

# ASSIGNMENT 1

CS 215

Details of students in the group for the assignment :

- 1.) Name : TATHAGAT VERMA      Roll Number : 180050111.
- 2.) Name : NEEL ARYAN GUPTA      Roll Number : 180050067.

ANSWERS : Follow from next page.

- Instructions for running the MATLAB code is specified in the INSTRUCTIONS.txt file.

### Question - 1

Given  $n$  distinct values  $\{x_i\}_{i=1}^n$

$$\therefore (x_1 - \mu)^2 + (x_2 - \mu)^2 + \dots + (x_{k-1} - \mu)^2 + (x_{k+1} - \mu)^2 + \dots + (x_n - \mu)^2 \geq 0$$

$$\Rightarrow \sum (x_i - \mu)^2 \geq (x_k - \mu)^2$$

$$\therefore \sigma^2 (\text{variance}) = \frac{\sum (x_i - \mu)^2}{n-1}$$

$$\Rightarrow \sigma^2 (n-1) = \sum (x_i - \mu)^2 \geq (x_k - \mu)^2$$

$$\Rightarrow \sigma^2 (n-1) \geq (x_k - \mu)^2$$

$$\Rightarrow \sigma \sqrt{n-1} \geq |x_k - \mu|$$

Chebyshev's inequality:  $P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}$

$$\Rightarrow P(|X - \mu| \leq k\sigma) \geq 1 - \frac{1}{k^2}$$

$$\text{Put } k = \sqrt{n-1}$$

$$\Rightarrow P(|X - \mu| \leq \sigma \sqrt{n-1}) \geq 1 - \frac{1}{n-1}$$

As  $n$  increases, the probability will approach the value 1 (the probability should be 1 but isn't).

$\Rightarrow$  Chebyshev's inequality agrees with the given inequality only for large  $n$ .

This means that Chebyshev's inequality gives correct bounds but the bounds are very loose.

## Question - 2

$\mu$  = mean

$T$  = median

$\sigma$  = standard deviation

To prove:  $|\mu - T| \leq \sigma$

By Chebyshev - Cantelli inequality,

$$P(X - \mu \geq k\sigma) \leq \frac{1}{1+k^2}$$

$$\text{Also, } P(X - \mu \leq -k\sigma) \leq \frac{1}{1+k^2}$$

Put  $k=1$

$$\Rightarrow P(X - \mu \geq \sigma) \leq \frac{1}{2} \quad \& \quad P(X - \mu \leq -\sigma) \leq \frac{1}{2}$$

$$\therefore \frac{1}{2} = P(X \geq T) = P(X \leq T)$$

$$\Rightarrow P(X \geq \sigma + \mu) \leq P(X \geq T) \quad \& \quad P(X \leq \mu - \sigma) \leq P(X \leq T)$$

Since  $P(X \geq T)$  is greater than  $P(X \geq \sigma + \mu)$ , it means that ~~it must be a greater bound than~~  $\sigma + \mu$  values greater than  $T$  are more ~~than~~ in number than values greater than  $\sigma + \mu$ .

$$\Rightarrow T \leq \sigma + \mu \Rightarrow T - \mu \leq \sigma$$

Similarly from 2<sup>nd</sup> inequality, we have

$$T \geq -\sigma + \mu$$

$$\Rightarrow T - \mu \geq -\sigma$$

$$\text{So, } -\sigma \leq T - \mu \leq \sigma$$

$$\Rightarrow |T - \mu| \leq \sigma$$

Hence, proved.



3.] Given: There exist 100 rickshaws of which 1 is red and 99 are blue.

XYZ sees red objects as red 99% of the time  
& blue objects as red 2% of the time.

$$\therefore P(\text{rickshaw was really red} \mid \text{XYZ observed it to be red}) \quad \{\text{using ①}\}$$

$$= \frac{P(\text{rickshaw was red, XYZ observed it to be red})}{P(\text{XYZ observed it to be red})} \quad \{\text{Refer NOTE}\}$$

$$= \frac{P(\text{rickshaw was red}) \times P(\text{XYZ observed red rickshaw as red})}{P(\text{XYZ observed a rickshaw to be red})}$$

$$= \frac{1/100 \times 99/100}{P(\text{rickshaw was blue \& XYZ observed red}) + P(\text{rickshaw was red \& XYZ observed red})}$$

$$= \frac{1/100 \times 99/100}{\dots} \quad \{\text{Refer NOTE}\}$$

$$P(\text{rickshaw was blue}) \times P(\text{XYZ observed blue object as red}) + P(\text{rickshaw was red}) \times P(\text{XYZ observed red as red})$$

$$= \frac{1/100 \times 99/100}{99/100 \times 2/100 + 1/100 \times 99/100} = \frac{99}{99 \times 2 + 99} = \boxed{1/3}$$

Answer: probability that the rickshaw was really a red when XYZ observed it to be a red one is  $\frac{1}{3}$ .

used in solution

$$\textcircled{1} P(X|Y) = \frac{P(X, Y)}{P(Y)}$$

NOTE: We also have considered that the color of the rickshaw is independent of the probability of XYZ seeing red objects as red or blue objects as red.  
Hence we used  $P(X, Y) = P(X)P(Y)$ .

$$4.] (a) P(C_i | Z_1) = \frac{P(C_i, Z_1)}{P(Z_1)}$$

$$= \frac{P(C_i) P(Z_1)}{P(Z_1)} = \frac{\frac{1}{3} \times \cancel{\frac{1}{3}}}{\cancel{\frac{1}{3}}} \text{ for all } i=1,2,3$$

~~for~~ ( $P(C_i, Z_1) = P(C_i) \times P(Z_1)$  as the car being behind a particular door is independent of the door being chosen first)

$$P(C_i | Z_1) = \boxed{\frac{1}{3}} \text{ for } i=1,2,3. \quad [\text{ANS}]$$

$$(b) P(H_3 | C_i, Z_1) = \frac{P(H_3, C_i, Z_1)}{P(C_i, Z_1)}$$

$$i=1; P(H_3 | C_1, Z_1) = \frac{P(H_3, C_1, Z_1)}{P(C_1, Z_1)} = \frac{\cancel{1/3} \times \cancel{1/3} \times 1/2}{\cancel{1/3} \times \cancel{1/3}}$$

(as the host can open door 2 or 3 with equal probability)  
&  $P(C_1) = 1/3$   
 $P(Z_1) = 1/3$

$$i=1 \Rightarrow \boxed{P(H_3 | C_1, Z_1) = 1/2}$$

$$i=2; P(H_3 | C_2, Z_1) = \frac{P(H_3, C_2, Z_1)}{P(C_2, Z_1)} = \frac{1/3 \times 1/3 \times 1}{1/3 \times 1/3} = 1$$

(as when car is behind door 2 & contestant chose door 1, the host will open door 3 with probability 1).  
&  $P(C_2) = P(Z_1) = 1/3$ .

$$i=2 \Rightarrow \boxed{P(H_3 | C_2, Z_1) = 1}$$

$$i=3; \boxed{P(H_3 | C_3, Z_1) = 0}$$

(As when the car is behind door 3, the host will NOT open door 3).

$$(c) P(C_2 | H_3, Z_1) = \frac{P(H_3 | C_2, Z_1) P(C_2, Z_1)}{P(H_3, Z_1)}$$

$$= \frac{1 \times 1/9}{P(H_3, C_1, Z_1) + P(H_3, C_2, Z_1) + P(H_3, C_3, Z_1)}$$

$$= \frac{1/9}{\sum_i P(H_3 | C_i, Z_1) \times P(C_i, Z_1)}$$

$$= \frac{1/9}{\sum_i P(H_3 | C_i, Z_1) \times P(C_i, Z_1)}$$

[as  $C_1, C_2, C_3$  are mutually exclusive & exhaustive]

$$= \frac{1/9}{1/9 \times 1/2 + 1/9 \times 1 + 1/9 \times 0} = \frac{\cancel{1/9}}{\cancel{1/9} \times 3/2} = \frac{2}{3}$$

$$\Rightarrow \boxed{P(C_2 | H_3, Z_1) = 2/3}$$

$$(d) P(C_1 | H_3, Z_1) = \frac{P(H_3 | C_1, Z_1) P(C_1, Z_1)}{P(H_3, Z_1)}$$

$$= \frac{(1/2) \times (1/9)}{1/9 \times 1/2 + 1/9 \times 1 + 1/9 \times 0}$$

{ same explanation as that of part (c) }

$$= \frac{1/2 \times \cancel{1/9}}{3/2 \times \cancel{1/9}} = \frac{1}{3}$$

$$\Rightarrow \boxed{P(C_1 | H_3, Z_1) = 1/3}$$

(e) Since probability of winning by switching i.e.  $P(C_2 | H_3, Z_1) = 2/3$  > probability of winning by not switching i.e.  $P(C_1 | H_3, Z_1) = 1/3$ , we conclude that switching is indeed beneficial.

(f.) Repeating calculations when the host chooses to open doors with equal probability.



$$(i) P(C_i | Z_1) = \underline{1/3} \quad \text{for } i=1,2,3.$$

remains unchanged as this doesn't depend on which door the host opens.

$$(ii) P(H_3 | C_i, Z_1) = 1/2 \quad \text{for } i=1,2,3$$

as no-matter where the car is, the host will ~~choose~~ <sup>open</sup> any one of doors 2 & 3, with the equal probability i.e. 1/2.

$$(iii) P(C_2 | H_3, Z_1) = \frac{P(H_3 | C_2, Z_1) P(C_2, Z_1)}{P(H_3, Z_1)}$$

$$= \frac{1/2 \times 1/3 \times 1/3}{\sum_1^3 P(H_3 | C_i, Z_1) P(C_i, Z_1)}$$

{ ~~sim~~ same logic as part (c) }

$$= \frac{1/2 \times 1/9}{1/2 \times 1/9 \times 3} = \underline{1/3}$$

$$\Rightarrow \boxed{P(C_2 | H_3, Z_1) = 1/3}$$

$$(iv) P(C_1 | H_3, Z_1) = \frac{P(H_3 | C_1, Z_1) P(C_1, Z_1)}{P(H_3, Z_1)}$$

$$= \frac{1/2 \times 1/9}{1/2 \times 1/9 \times 3}$$

{ similar calculations as that of previous part }

$$= \underline{1/3}$$

$$\Rightarrow \boxed{P(C_1 | H_3, Z_1) = 1/3}$$

Since probability of winning with & without switching is the same i.e.  $P(C_2 | H_3, Z_1) = P(C_1 | H_3, Z_1) = 1/3$ ,



It is NOT Beneficial to switch. (as chances to win remains unchanged).

5.] For  $f = 30\%$ , Relative mean squared error with:

$$\text{Median} = 2.99 \times 10^1.$$

$$\text{Mean} = 9.22 \times 10^1$$

$$\text{1st Quartile} = 1.177 \times 10^{-2}.$$

For  $f = 60\%$ , Relative mean squared error with:

$$\text{Median} = 6.38 \times 10^2$$

$$\text{Mean} = 3.54 \times 10^2$$

$$\text{1st Quartile} = 1.05 \times 10^2.$$

NOTE: These error values are for one iteration (one time), with every iteration, different random values and positions will be generated leading to different error values. However the range would be around these values.

and  $f = 60\%$   
Now; For  $f = 30\%$ , We can observe Relative Mean Squared Error to be much less for 1st Quartile than for Mean & Median. Hence quartile filtering produces best relative mean squared error.

This can be justified as;  
The error caused is mainly due to addition of random values in the 100-120 range. (Taking mean/median/quartile also generates error but this is magnified only due to the large random values).

- In mean filtering, every ~~random~~ random value added is ~~to~~ taken into account,
- While for median, the random value would be accounted for only if number of random values added in a particular interval of  $i-8$  to  $i+8$  was  $> 8$ . (This has less probability)
- Eg. original values of an interval are (we consider length 5 subarray)  
 $y = 0.10 \quad 0.15 \quad 0.2 \quad 0.17 \quad 0.12$

If only 1 random value was added;  
 mean would change, median & 25% quantile would remain unchanged.

- Median would ~~can~~ change only if number of random values added was  $> 2$ . (As values added are greater than original values of  $y$ .)
- 25% quantile would change only if number of random values added was  $> [75\% \times (\text{length of subarray})]$   

$$= \left[ \frac{3}{4} \times 5 \right] = [3.75] = 3$$

- Due to this reason, 1st quartile produces least mean squared error. ~~£~~

Relative

- For  $f=60\%$ , we observe <sup>Relative</sup> Mean squared error to be least for 25% quantile, then Mean, then Median.

In this case, since 60% (many) random values are added, in many intervals of  $i-8$  to  $i+8$ , there would be  $> 8$  random values added due to which median went in the range of 100-120 causing lot of error, ~~is~~

while the mean in this case uses the un-corrupted values as well to reduce the filtered value to some extent. Due to this relative mean squared error is less for mean than for median.

In this case, still 1st quartile performs best due to the same reason stated before<sup>(i.e.  $f=30\%$ )</sup>, but we can observe the extent to which 1st quartile performs better than mean & Median has reduced. The reason for this is same as the reason why median has started performing poorer than the mean.

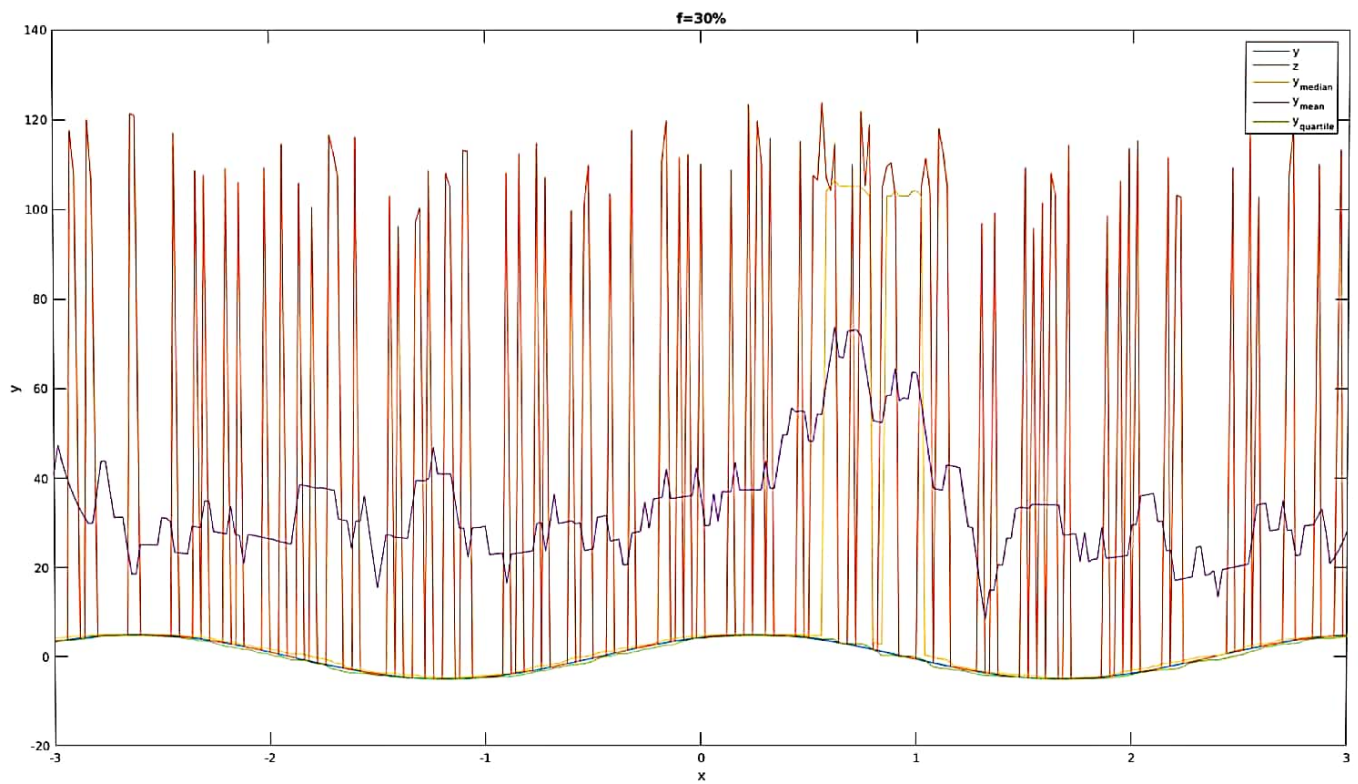
- For  $f=30\%$ , what we can see is the robustness of median, and more robustness of 1st quartile.

Instructions to execute: Go in the directory of the Zip folder & execute :-

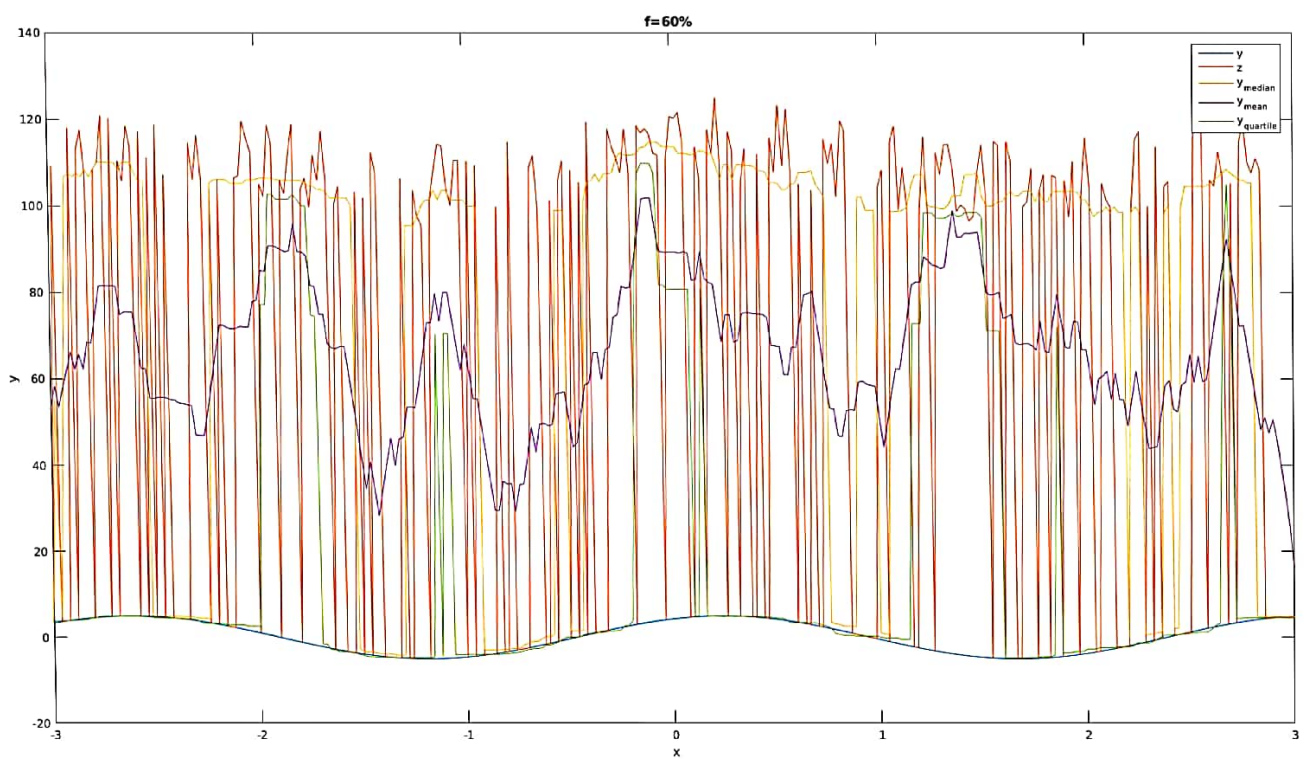
`run('a1q5.m');`

in the command window.

- Graphs are supplied on the next pages.







## Question - 6 REPORT

Let  $\{x_i\}_{i=1}^n$  denote the old array elements and  $x$  be the new element

### → Updating the mean

$$\mu_{old} = \frac{\sum x_i}{n} \Rightarrow \sum x_i = n\mu_{old}$$

$$\mu_{new} = \frac{\sum x_i + x}{n+1} = \frac{n\mu_{old} + x}{n+1}$$

### → Updating the median

Case 1:  $n$  is odd

$x_1, x_2, \dots, x_{k-1}, x_k, x_{k+1}, \dots, x_n$   
be the sorted array, without loss of generality  
 $\therefore k = (n+1)/2$ ,  $\mu_{old} = x_k$ . Since  $n+1$  is even,  
If  $x \geq x_{k+1}$  then  $\mu_{new} = (x_k + x_{k+1})/2$   
If  $x < x_{k-1}$  then  $\mu_{new} = (x_k + x_{k-1})/2$   
else  $\mu_{new} = (x_k + x)/2$  [ $x_{k-1} \leq x \leq x_{k+1}$ ]

Case 2:  $n$  is even

$x_1, x_2, \dots, x_k, x_{k+1}, \dots, x_n$   
be the sorted array, WLOG  
 $\therefore k = n/2$ ,  $\mu_{old} = (x_k + x_{k+1})/2$ , since  $n+1$  is odd  
If  $x \geq x_{k+1}$  then  $\mu_{new} = x_{k+1}$   
If  $x \leq x_k$  then  $\mu_{new} = x_k$   
else  $\mu_{new} = x$  [ $x_k < x < x_{k+1}$ ]

### → Updating the standard deviation

$$\text{Let } \mu_{new} - \mu_{old} = d$$

$$\text{i.e. } \mu_{new} = \mu_{old} + d$$

$$\sigma_{old}^2 = \frac{\sum (x_i - \mu_{old})^2}{n-1}$$

$$\Rightarrow \sum (x_i - \mu_{old})^2 = (n-1)\sigma_{old}^2$$

$$\therefore (x_i - \mu_{\text{new}})^2 = (x_i - \mu_{\text{old}} - d)^2$$

$$\sum (x_i - \mu_{\text{new}})^2 = \sum (x_i - \mu_{\text{old}} - d)^2$$

$$= \sum (x_i - \mu_{\text{old}})^2 + \sum d^2 - 2d \sum (x_i - \mu_{\text{old}})$$

$$\text{Also } \sum x_i = n \mu_{\text{old}} = \sum \mu_{\text{old}}$$

$$\Rightarrow \sum (x_i - \mu_{\text{old}}) = 0$$

$$\Rightarrow \sum (x_i - \mu_{\text{new}})^2 = \sum (x_i - \mu_{\text{old}})^2 + nd^2 - 0$$

$$= (n-1)\sigma_{\text{old}}^2 + nd^2$$

$$\sigma_{\text{new}}^2 = \frac{\sum (x_i - \mu_{\text{new}})^2 + (x - \mu_{\text{new}})^2}{n}$$

$$= \frac{(n-1)\sigma_{\text{old}}^2 + nd^2 + (x - \mu_{\text{new}})^2}{n}$$

$$\Rightarrow \sigma_{\text{new}} = \sqrt{\frac{(n-1)\sigma_{\text{old}}^2 + n(\mu_{\text{new}} - \mu_{\text{old}})^2 + (x - \mu_{\text{new}})^2}{n}}$$

→ Updating the histogram

We will just increment the value of the bin by 1 in which the new value/element lies.

