# DAI assignment-1

**23B0924-B.somanath vasanth**
**23B1032-V.Srikar nihal**
**23B1055-S.Anirudh**

# Files included in the zip file

For q6 code is in the file q6.py

For q7 code for plots is in q7violinplot.py,q7paretoplot.py,q7coxcomb.py,q7waterfallplot.py

For q8 code for q8monalisa.py

1) Let's Gamble.

Given,

A has $n+1$ dice (fair)

B has $n$ dice (fair)

Probability of getting a prime number on top $\sim \not{p}$

$$\not{p} = \frac{3}{6} = \frac{1}{2}.$$  ($\because$ Each dice has 3 prime numbers i.e.,

$$\{2,3,5\})$$

let $p$ be the probability that A will have more wins than B

let $x$ be the ~~no.of~~ no. of wins of A and Y be the no. of wins of B

$$P = P(X=1)P(Y<1) + P(X=2)P(Y<2) + P(x=3)P(Y<3) \cdots \underset{\text{n+1}}{P(x=n+1)P(Y<n+1)}$$

$$\cdots \cdots P(X=n+1)P(Y<n+1)$$

$$P(X=\gamma) = (n+1)_{c\gamma}\left(\frac{1}{2}\right)^{\gamma}\left(\frac{1}{2}\right)^{n+1-\gamma}.$$

$$= {}^{n+1}_{c\gamma}\left(\frac{1}{2}\right)^{n+1}$$

$$P(Y \leq \gamma) = \sum_{\alpha=0}^{\gamma-1} P(Y=\alpha) = \sum_{\alpha=0}^{\gamma-1} n_{c\alpha}\left(\frac{1}{2}\right)^{\alpha}\left(\frac{1}{2}\right)^{n-\alpha}.$$

$$P = {}^{n+1}c_1\left({}^{n}c_0\left(\frac{1}{2}\right)^{n}\right)\left(\frac{1}{2}\right)^{n+1} + {}^{n+1}c_2\cdot\left(\frac{1}{2}\right)^{n+1}\left[n_{c_0}+n_{c_1}\right]\left(\frac{1}{2}\right)^{n} - - - -$$

$$\cdots {}^{n+1}c_{n+1}\left(n_{c_0}+n_{c_1}\cdots n_{cn}\right)\left(\frac{1}{2}\right)^{2n+1}$$

$$= \left(\frac{1}{2}\right)^{2n+1}\left[\left({}^{n+1}c_1 n_{c_0} + {}^{n+1}c_2 n_{c_1} + \cdots {}^{n+1}c_{n+1}n_{cn}\right) + \left({}^{n+1}c_2 n_{c_0} + {}^{n+1}c_3 n_{c_1}\right.\right.$$

$$\cdots {}^{n+1}c_{n+1}n_{cn}\right)$$

$$+ \cdots \cdots \left({}^{n+1}c_{n+1}\right)(n_{c_0})\Big]$$

$$= \left(\frac{1}{2}\right)^{2n+1}\left[\left({}^{n+1}c_1 n_{cn} + {}^{n+1}c_2 n_{cn-1} \cdots {}^{n+1}c_{n+1}n_{c_0}\right) + \left({}^{n+1}c_2 n_{cn} + {}^{n+1}c_3 n_{cn-1}\right.\right.$$

$$\cdots {}^{n+1}c_{n+1}n_{c_1}\Big)$$

$$+ \cdots \left({}^{n+1}c_{n+1}n_{cn}\right)\Big]$$

$$(\because n_{c\gamma} = n_{c(n-\gamma)})$$

$$P = \left(\frac{1}{2}\right)^{2n+1} \left( {}^{2n+1}C_{n+1} + {}^{2n+1}C_{n+2} \cdots \cdots {}^{2n+1}C_{2n+1} \right)$$

$$\left( \because \sum {}^{a}C_r \, {}^{b}C_{n-r} = {}^{a+b}C_n \right)$$

$$= \left(\frac{1}{2}\right)^{2n+1} \left( 2^{2n} \right)$$

$$\left( \because \sum_{r=0}^{2n+1} {}^{2n+1}C_r = 2^{2n+1} \quad \text{and} \right.$$

$${}^{2n+1}C_r = {}^{2n+1}C_{2n+1-r}.$$

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

$$\left. \Rightarrow \sum_{r=0}^{n} {}^{2n+1}C_q = \sum_{r=0}^{n} {}^{2n+1}C_{2n+1-r} \right)$$

$\therefore$ The probability that A will have more wins than B $= \underline{\frac{1}{2}}$

## 2) Two Trading Teams

Given, Team B is better at trading than Team A.

let $P(A)$ be the probability of we win against A. and $P(B)$ be the probability of we win against B.

So, $P(A) > P(B)$

Let $P(A-B-A)$ (case 1: A-B-A

$P_1$ be the probability to win in two sets in a row.

$$P_1 = P(A\,win)\,P(B\,win) + P(A)\,P(A\,lose)\,P(B\,win)\,P(A\,win)$$

$$\left( \begin{array}{l} \because P(A\,win) \\ \text{is probability we} \\ \text{winning if we play} \\ \text{with A similar for B)} \end{array} \right.$$

$$= P(A)P(B) + \{1 - P(A)\}\, P(B)\, P(A)$$

Case 2: B-A-B

$$P_2 = P(B)P(A) + (1-P(B))(P(A))P(B)$$

$$1 - P(A) < 1 - P(B), \qquad (\because P(A) > P(B))$$

$$P(A)P(B) + (1-P(A))P(B)P(A) \;<\; P(A)P(B) + (1-P(B))\,P(A)P(B)$$

$$P_1 \qquad\qquad\qquad < \qquad P_2.$$

Hence the probability of we win is more in the case -2

i.e., B-A-B order

3)

3.1) Given that,

$\&$ $Q_1, Q_2, q_1, q_2$ are non-negative

$$P(Q_1 < q_1) \geq 1 - P_1$$

$$P(Q_2 < q_2) \geq 1 - P_2$$

We need to prove prove $P(Q_1 Q_2 < q_1 q_2) \geq 1 - P_1 P_2$

$$P(Q_1 < q_1) \geq 1 - P_1 \implies P(Q_1 \geq q_1) \leq P_1$$

Similarly $P(Q_2 \geq q_2) \leq P_2$

We know that for $Q_1 Q_2 \geq q_1 q_2$ we need to have $\&$ $Q_1 \geq q_1$
or $Q_2 \geq q_2$. (if $Q_1 < q_1$ and $q_2 < q_2$ then $Q_1 Q_2 < q_1 q_2$)

So let $A$ be the event that $Q_1 \geq q_1$ and $B$ be the event that $Q_2 \geq q_2$ and $C$ be the event that $Q_1 Q_2 \geq q_1 q_2$

We know that $P(A \cup B) \leq P(A) + P(B)$

$$\implies P(Q_1 Q_2 \geq q_1 q_2) \leq P(Q_1 \geq q_1) + P(Q_2 \geq q_2)$$

$$\implies P(Q_1 Q_2 \geq q_1 q_2) \leq P_1 + P_2$$

$$\implies P(Q_1 Q_2 < q_1 q_2) > 1 - (P_1 + P_2) \quad (\because P(A) < k \text{ then } P(\bar{A}) \geq 1 - k)$$

Hence Proved.

3.2) Given, for data values $\{x_i\}_{i=1}^{n}$, $\mu$ is mean and $a$ is standard deviation.

We need to prove that $|x - \mu_i| \leq a\sqrt{n-1}$

Now let us consider $a = \sqrt{\dfrac{(x_1 - \mu)^2 + (x_2 - \mu)^2 + \cdots (x_n - \mu)^2}{n-1}}$

from the formula we can say that

$$a\sqrt{n-1} = \sqrt{\sum_{i=1}^{n}(x_i-\mu)^2}$$

$$\Rightarrow a\sqrt{n-1} = \sqrt{\sum_{i=1}^{n}(x_i-\mu)^2} \geq \sqrt{(x_i-\mu)^2}$$

$$\Rightarrow \sigma\sqrt{n-1} \geq \sqrt{(x_i-\mu)^2}$$

$$\Rightarrow \sqrt{(x_i-\mu)^2} \leq \sigma\sqrt{n-1}$$

$$\Rightarrow |x_i-\mu| \leq \sigma\sqrt{n-1}$$

chebyshev's inequality states that

The proportion of sample points $k$ or more than $k$ $(k>0)$ standard deviation away from the sample mean is less than or equal to $1/k^2$.

$$S_K = \{x_i : |x_i - \bar{x}_\sigma| \geq k\sigma\}$$

$$\frac{S_K}{N} \leq \frac{1}{k^2}$$

if we substitute $k = \sqrt{n-1}$ we have

$$S_K = \{x_i : |x_i-\mu| \geq \sigma\sqrt{n-1}\}$$

$$\frac{S_K}{N} \leq \frac{1}{n-1}$$

$$\Rightarrow S'_K = \{x_i : |x_i-\mu| \leq \sigma\sqrt{n-1}\}$$

$$\frac{S'_K}{N} \geq \frac{n-2}{n-1}$$

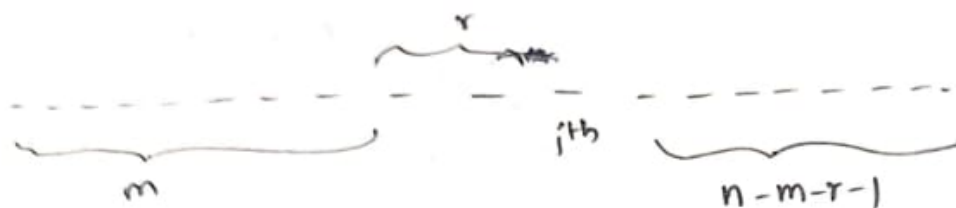as $n$ increases the chebyshev's inequality tends to the given inequality.

# 4) Staff Assitant

(a) Given,

Event $E_i$ be that $i^{th}$ candidate is best and we hire him.

$Pr(E_i) = 0 \qquad 1 \le i \le m$ (∵ We reject the first $m$ candidates)

$Pr(E_i)$ if $i > m$

let $i = m + r + 1$        $i^{th}$ is the best candidate



$$Pr(E_i) = \frac{{}^{n-1}C_{m+r} \times \left({}^{m}C_1 \times (m+r-1)!\right) \times (n-m-r-1)!}{n!}$$

- $\rightarrow$ Total no of arrangements

selecting $m+r$. candidates from $n$ excluding the best ones.

candidate

As the best of $m+r$ should be in first $m$ if not then he will be selected

arranging of remaining. $m+r-1$ candidates

arranging $(n-m-r-1)$ candidate.

$$= \frac{(n-1)!}{(m+r)!\,(n-m-r-1)!} \times \frac{m \times (m+r-1)! \times (n-m-r-1)!}{n!}$$
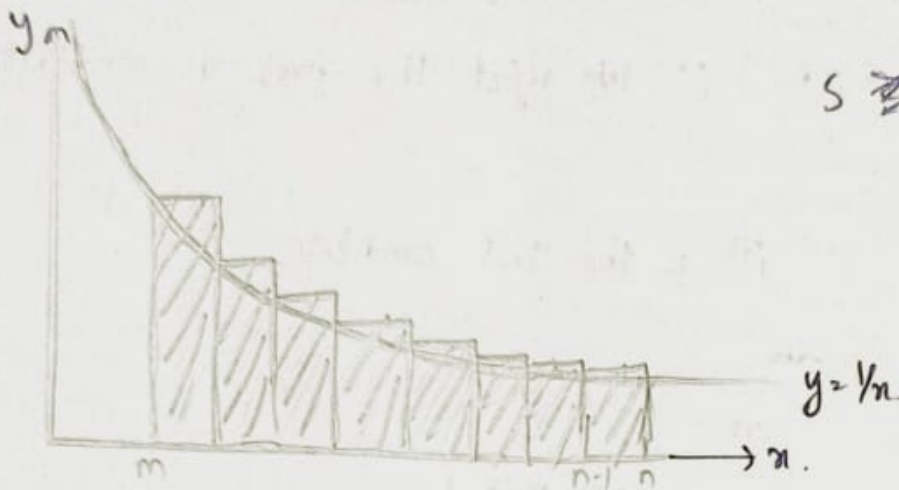
$$= \frac{(n-1)! \times m \times (m+r-1)!}{n\,(n-1)!\,(m+r)\,(m+r-1)!}$$

$$= \frac{m}{n} \times \frac{1}{m+r} \qquad = \frac{m}{n} \times \frac{1}{\cdot i-1}$$

$$Pr(E_i) = \begin{cases} 0 & 1 \le i \le m \\[2mm] \dfrac{m}{n} \times \dfrac{1}{i-1} & m+1 \le i \le n \end{cases}$$
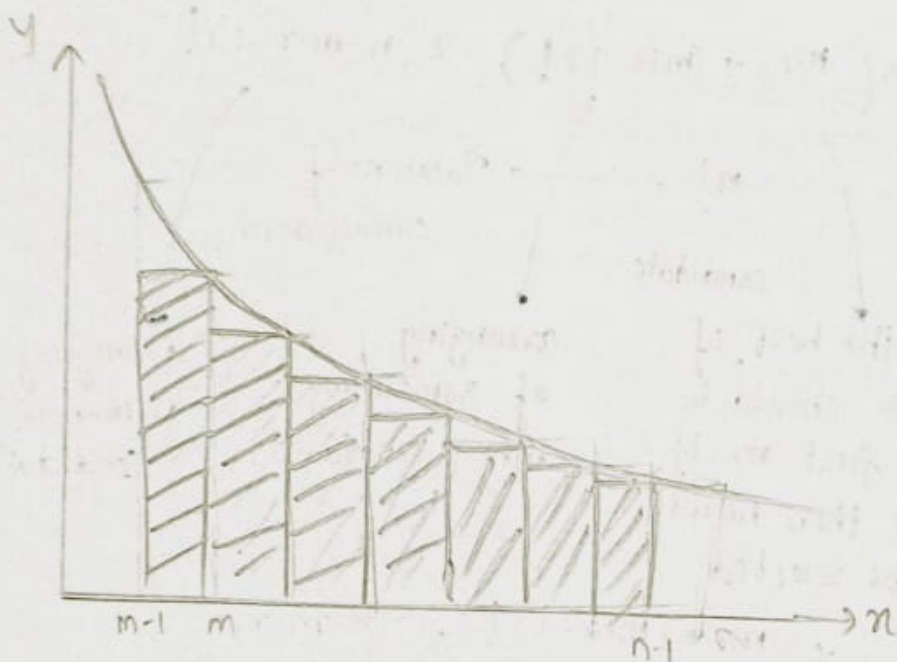
(b) $S = \sum\limits_{j=m+1}^{n} \dfrac{1}{j-1} = \sum\limits_{i=m}^{n-1} \dfrac{1}{i} = \dfrac{1}{m} + \dfrac{1}{m+1} \cdots \dfrac{1}{n-1}$



$y = 1/x$

$S \geq \displaystyle\int_{m}^{n} \dfrac{1}{x}\, dx \quad \left(\because \int \frac{1}{x}dx = \ln x\right)$

$S \geq \left[\ln x\right]_{m}^{n}$

$S \geq \ln n - \ln m$



$S \leq \displaystyle\int_{m-1}^{n-1} \dfrac{1}{x}\, dx$

$S \leq \left[\ln x\right]_{m-1}^{n-1}$

$S \leq \ln(n-1) - \ln(m-1)$

$Pr(t) = m/n \; S$

$\therefore \dfrac{m}{n}\left(\ln(n) - \ln(m)\right) \leq Pr(t) \leq \dfrac{m}{n}\left(\ln(n-1) - \ln(m-1)\right)$

(c) Given equation is $\dfrac{m}{n}\left[\ln(n) - \ln(m)\right]$ now differentiate it with respective to m and equate it to zero

$\dfrac{1}{n}\left(\ln(n) - \ln(m)\right) + \dfrac{m}{n}\left(0 - \dfrac{1}{m}\right) = 0$

$\Rightarrow \dfrac{\ln(n)}{n} - \dfrac{\ln(m)}{n} - \dfrac{1}{n} = 0$

$\Rightarrow \quad \dfrac{\ln(n) - (1 + \ln(m))}{n} = 0$

$\Rightarrow \quad (\ln(n) - 1) = \ln(m)$

$\Rightarrow \quad m = n/e$

Now if we differentiate it again with respective to $m$ we get

$$\frac{1}{n}\left(-1/m\right) + 0$$

after substituting $m = \dfrac{n}{e}$ we get $-e/n^2$ which is negative so $n/e$ is where $m$ is maximising the curve.

now if we substitute $m = n/e$ in

$\dfrac{m}{n} \ln(n/m) \leq P_r(E)$ (from previous part)

we get

$$P_r(E) \geq 1/e$$

7) **Uses of Violin-Plot**

① Violet plots allow for a direct comparison of distributions across multiple catagories. For instance

② Violet plot can reveal if a dataset has multiple peaks or modes which might indicate the presence of subgroups within the data

③ By visualizing density, violen plots help asses whether the data is symmetrically distributed or skewed

## ② Uses of Pareto Plots

① They are used to identify the most frequent defects or issues in manfacturing and service industries

② They help in prioritizing resources by highlighting the few Critical areas that cause the most significant impact

③ These plots are commonly used to determine the primary Causes of Problems

## ③ Uses of Coxcomb Plots

① This plot is particularly effective for visualizing cyclic or seasonal data, where the cyclical nature of the data can be naturally represented.

② The Coxcomb plot was famously used by Florence Nightingale to present mortality data, making it easier to understand the impact of different causes of death.

## ④ Uses of Waterfall plots:

① Waterfall plots are commonly used in finance to break down changes in finance metrics like profit or revenue

② Waterfall plots are used for analizing the step-by-step impact of changes in business strategies, cost reductions, or sales initiatives, helping to identify which factors had most significant effect.
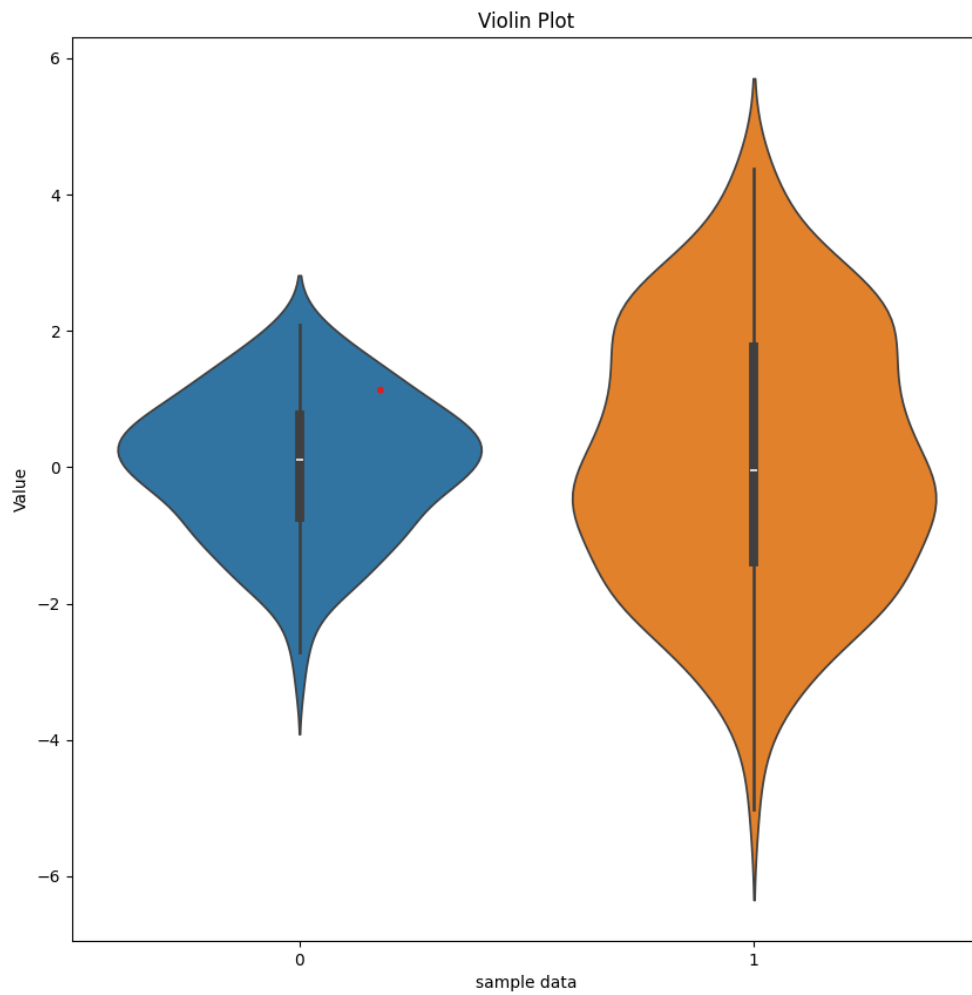
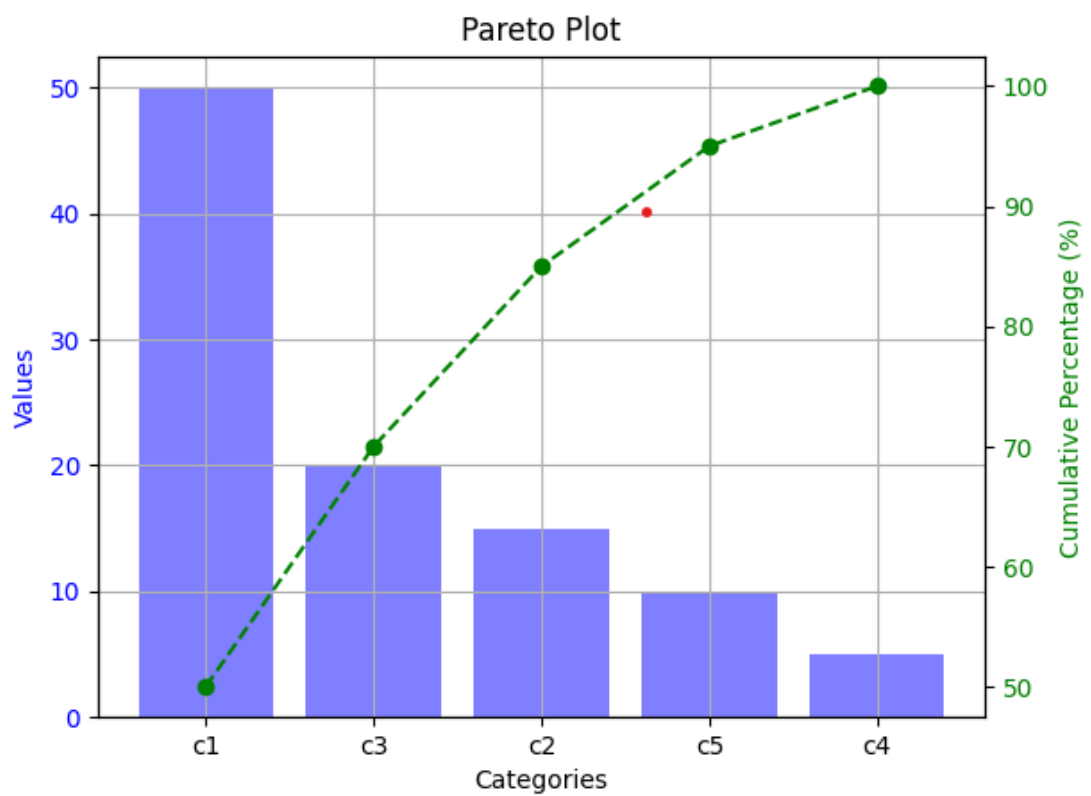Figure 1: violin plot generated for sample data



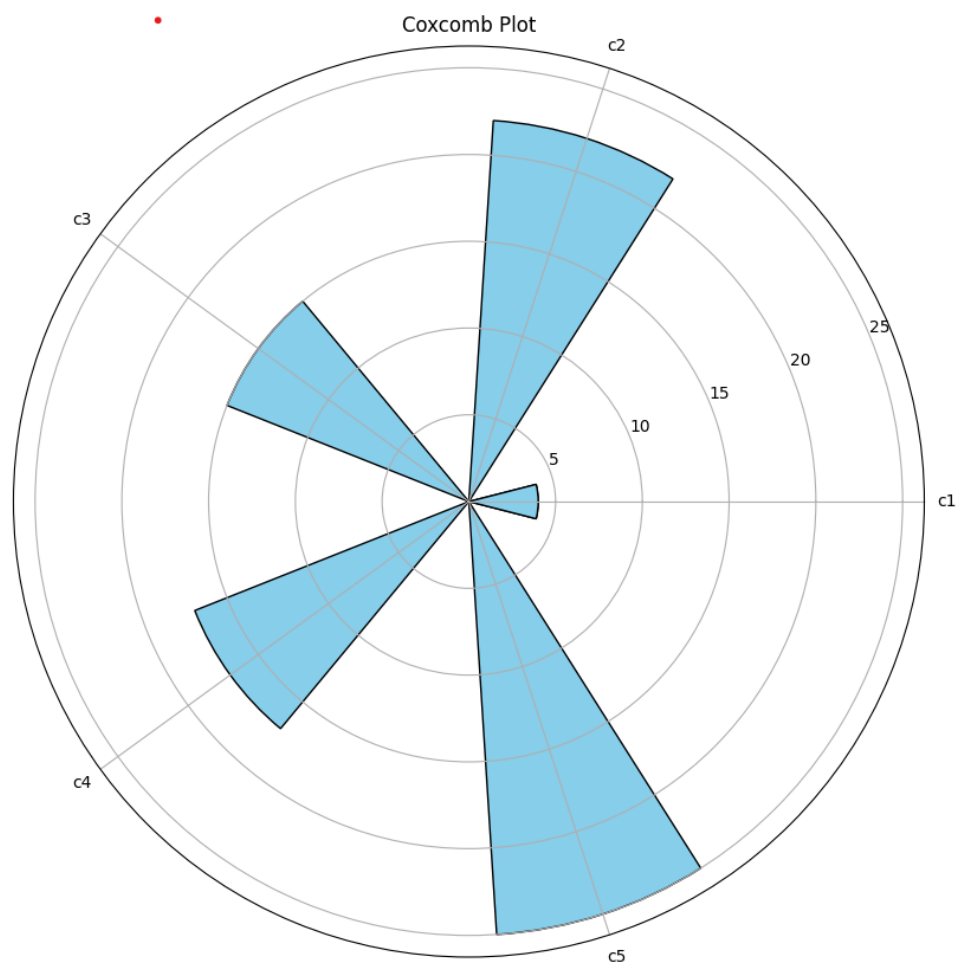Figure 2: pareto plot generated for sample data
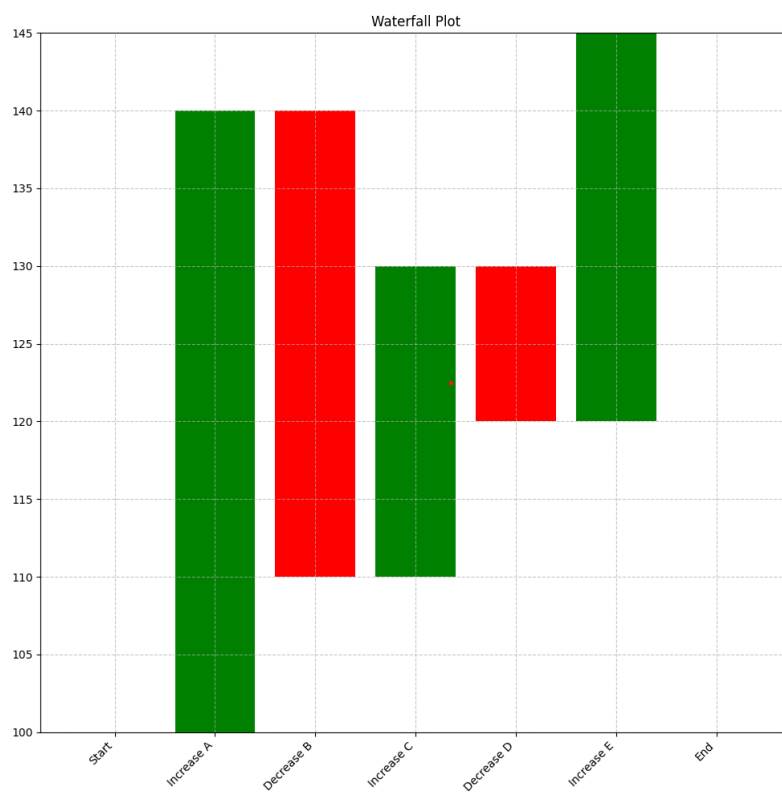
Figure 3: coxcomplot generated for sample data



Figure 4: pareto plot generated for sample data

**5)** given    200 tokens.

we  have to find the  pasn  in  which  probability  of  getting
a free tradeoff  is  maximum.

let  the  pasn  be  $k^{th}$  pasn.



$\underbrace{\qquad\qquad\qquad}_{k}$

now  let  the  token  at this $k^{th}$  person  be  $i$

the  tokens  from  1st to  $k-1^{th}$  pasn  should  have  1 token of $i$

&  $k-2$ tokens  of  other  199 tokens ( no  token should be  with
                                        more than 1 person )

selecting  $k-2$ tokens  from   199 $= 199 C_{k-2}$

now  arranging  these  token  in  different  orders

no. of orders $= (k-1)!$

The  probability  of  $k^{th}$ person getting ~~trade~~ a free trade

$= \dfrac{\text{no. of orders satisfying the given condition}}{\text{Total no. of orders}}$

no. of  order satisfying  given condition $= 199 C_{k-2} \times (k-1)!$

Total  no. of orders $= 200^{k-1}$ ( each  pasn can get
                                      any token)

Probability $= \dfrac{199 C_{k-2} \times (k-1)!}{200^{k-1}} = \dfrac{\dfrac{199!}{(k-2)!\,(201-k)!} \times (k-1)!}{200^{k-1}}$

$= \dfrac{(199!) \times (k-1)}{(201-k)!\,(200^{k-1})}$

now  ~~we~~  have to  maximize this  probability

$$P(K) = \frac{199! \, (k-1)}{(201-k)! \, (200)^{k-1}}$$

• calculating $P(i)/P(i+1)$

$$\frac{P(i)}{P(i+1)} = \frac{\dfrac{i-1}{(201-i)! \, 200^{i-1}}}{\dfrac{i}{(200-i)! \, 200^{i}}} = \left(\frac{i-1}{i}\right)\left(\frac{200}{\cdot}\right)\left(\frac{(200-i)!}{(201-i)!}\right)$$

$$= \left(\frac{i-1}{i}\right)\left(200\right)\frac{1}{201-i}$$

$$\frac{P(i)}{P(i+1)} = \frac{200\,(i-1)}{(i)(201-i)}$$

now    calculate. when

$$\frac{P(i)}{P(i+1)} \leq 1 \qquad \frac{200\,(i-1)}{i\,(201-i)} \leq 1$$

$$200i - 200 \leq 201i - i^2 \qquad\qquad\quad \frac{P(14)}{P(15)} \leq 1$$

$$i^2 - i - 200 \leq 0 \qquad\qquad\qquad\quad P(14) \leq P(15)$$

$$\frac{1 - \sqrt{1+4(200)}}{2} < i < \frac{1 + \sqrt{1+4(200)}}{2} \qquad \frac{P(15)}{P(16)} \geq 1 \quad P(15) \geq P(16)$$

$$-13.6525 \;\leq\; i \;\leq\; 14.6525 \qquad\qquad\qquad\qquad \text{15 is maximum}$$

by this $\frac{P(i)}{P(i+1)} \leq 1$ for 0 to 14 & $\frac{P(i)}{P(i+1)} \geq 1$ for greater than 15
    by this    calculation    we    can    say.

    P(i) increases    from    0 - 14 & reaches maximum

    at i = 15

So   we   should choose   the   15th posn of queue to get
maximum probability.

6) given

mean, median, standard deviation. of a set of n numbers.

### MEAN

$$mean = \frac{\sum\limits_{n=1}^{n} x_i}{n} \rightarrow ①$$

$$new\ mean = \frac{\sum\limits_{n=1}^{n} x_i + new\ element}{n+1}$$

$$from ① \Rightarrow \sum\limits_{n=1}^{n} x_i = n\,(mean)$$

$$\boxed{new\ mean = \frac{n\,(mean) + new\ element}{n+1}}$$

without calculating mean from individual elements we can use this formula.

### STANDARD DEVIATION:

$$St\ dev = \sqrt{\left(\frac{1}{n-1}\right) \sum (x_i - \mu)^2}$$

$$\sum (x_i - \mu)^2 = \sum x_i^2 + \mu^2 - 2\mu x_i$$

$$= \sum x_i^2 + n\mu^2 - 2\mu(n)$$

$$= \sum x_i^2 - \mu^2 n$$

$$St\ dev = \sqrt{\left(\frac{1}{n-1}\right)\left(\sum\limits_{1}^{n} x_i^2 - \mu^2 n\right)} \rightarrow ①$$

$$new\ st\ dev = \sqrt{\left(\frac{1}{n}\right)\left(\sum\limits_{1}^{n} x_i^2 + (new\ elem)^2 - (new\ mean)^2(n+1)\right)}$$

by squaring ①

$$(n-1)(St\ dev)^2 = \sum\limits_{1}^{n} x_i^2 - \mu^2 n \qquad \sum\limits_{1}^{n} x_i^2 = \mu^2 n + (n-1)(St\ dev)^2$$

by substituting this in formula

new std dev =

$$\sqrt{\left(\frac{1}{n}\right)\left((\text{newelem})^2 + (\text{oldmean})^2 n + (n-1)(\text{oldstddev})^2 - (n+1)(\text{newmean})^2\right)}$$

↓

by using this formula we can get new stddev.

## Median

Assumption: Assuming a A to be sorted

if n is even: $(n/2)$ $(n/2)$ . Indexing is done from 0.

— — — — — ... ↓ ↓ . — — — — — } n number
0  1  2  3       ⌣

→ if the new element is less than $(n/2-1)^{th}$ element
  it will be inserted somewhere ^before it & $\left(\frac{n}{2}-1\right)^{th}$ element
  will become median / middle element

→ like this if new element is greater tha $\frac{n}{2}^{th}$ element
  it will be inserted after $n/2$ th element & $\left(\frac{n}{2}\right)^{th}$ element
  will become median.
                                                both including
→ if new element is b/w $n/2$ & $\frac{n}{2}-1$ elements ^ then
  the new element will become median.

If n is odd

$\left(\frac{n-3}{2}\right)^{th}$ $\left(\frac{n+1}{2}\right)^{th}$ [Indexing done from 0]

$$\bar{0} \quad \bar{1} \quad \bar{2} \quad \cdots \quad \underset{\uparrow}{-} \quad \underset{\downarrow}{-} \quad \underset{\uparrow}{-} \quad - \quad - \quad - \quad -$$

$\left(\frac{n-1}{2}\right)^{th}$

if new element is greater than $\left(\frac{n+1}{2}\right)^{th}$ element.

then $\left(\frac{n-1}{2}\right)^{th}$ & $\left(\frac{n+1}{2}\right)^{th}$ elements will be come middle elements

and median is $\dfrac{\left(\frac{n-1}{2}\right)^{th} \text{element} + \left(\frac{n+1}{2}\right)^{th} \text{element}}{2}$

if new element is lesser than $\left(\frac{n-3}{2}\right)^{th}$ element.

then $\left(\frac{n-3}{2}\right)^{th}$ element & $\left(\frac{n-1}{2}\right)^{th}$ element will become middle elements

and median is $\dfrac{\left(\frac{n-3}{2}\right)^{th} \text{elem} + \left(\frac{n-1}{2}\right)^{th} \text{elem}}{2}$

if new element is b/w $\left(\frac{n-3}{2}\right)^{th}$ & $\left(\frac{n+1}{2}\right)^{th}$ element & [both including]

then $\left(\frac{n+1}{2}\right)^{th}$ & new element will become middle elements.

and median is $\dfrac{\text{new element} + \left(\frac{n+1}{2}\right)^{th} \text{element}}{2}$

Like this we find median.

How to update histogram of A:

↳ we can check for the bin which contains the range in which new element lies & increase its count of that bin ~~frequency~~ by 1. every thing else remains the same

If the new value doesn't fall in any bin we should create a new bin for this ~~new element~~ new element that is added
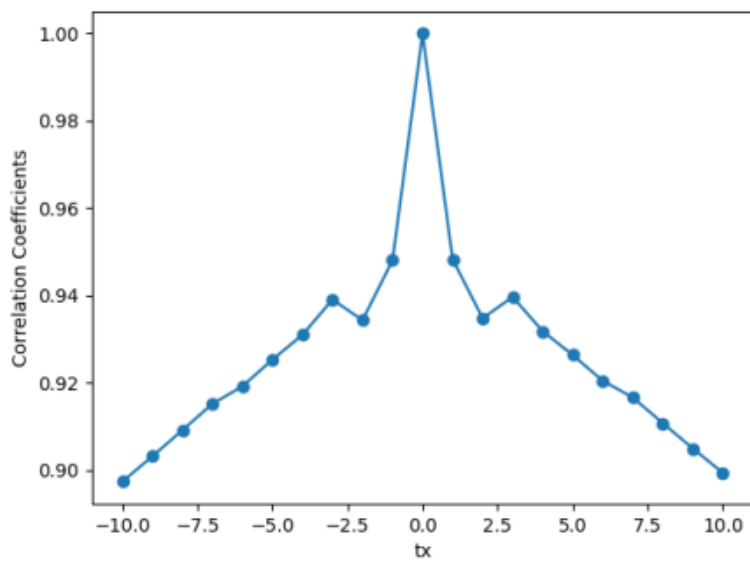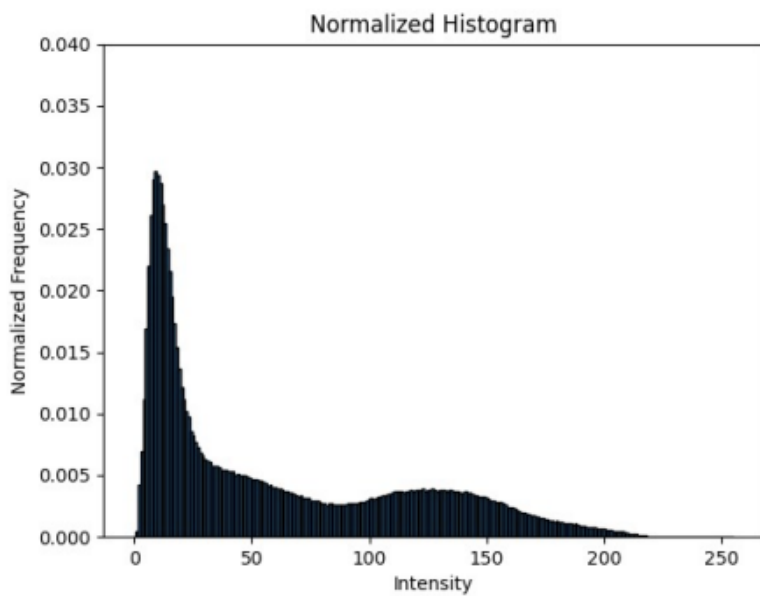
# Question 8



Figure 1: correlation coefficient vs tx



Figure 2: Normalized frequency vs intensity