# 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 INSTRUCIONS.txt file.

## Question 1 :

**(a)** $X_1$ = no. of times additionaly a book has to be picked so that we move from 0 distinct colors to 1 distinct color

$= 1$, (as the 1st book itself serves as the 1st distinct color)

$$\Rightarrow \boxed{X_1 = 1}$$ (Ans.)

**(b)** $P(X_i = k) = P($ (k-1) books picked belonged to the (i-1) colors & the kth book belongs to the (n-i+1) other colors)

$$= \left(\frac{i-1}{n}\right)^{k-1}\left(1 - \frac{i-1}{n}\right)$$

$$= \left(1 - \frac{n-i+1}{n}\right)^{k-1}\left(\frac{n-i+1}{n}\right)$$

$$= (1-p)^{k-1}\, p .$$

$\Rightarrow X_i \sim$ Geometric random variable with

$$\text{parameter} = \boxed{\dfrac{n-i+1}{n}}$$ (Ans.)

(C) $P(Z=K) = (1-p)^{k-1} p$           To prove : $\underline{E(Z) = 1/p}$.

$$E(Z) = \sum_{K=1}^{\infty} K p (1-p)^{k-1}$$

Now g consider $g(x) = \sum_{K=1}^{\infty} x(1-x)^k$

$$g(x) = \frac{x(1-x)}{1-(1-x)} = (1-x) \qquad \{ x \in (0,1) \}$$

$$g'(x) = \sum_{K=1}^{\infty} (1-x)^k - \sum_{K=1}^{\infty} Kx(1-x)^{k-1}$$

$$\Rightarrow \frac{d}{dx}\left((1-x)\right) = \frac{(1-x)}{x} - \sum_{1}^{\infty} Kx(1-x)^{k-1}$$

$$\Rightarrow -1 = \frac{1}{x} - 1 - \sum_{1}^{\infty} Kx(1-x)^{k-1}$$

$$\Rightarrow \boxed{\sum_{1}^{\infty} Kx(1-x)^{k-1} = \frac{1}{x}} \qquad \{ x \in (0,1) \}$$

Put $x=p$       ; valid as $0 < p \leq 1$

$$\Rightarrow \sum_{1}^{\infty} K p (1-p)^{k-1} = \boxed{\frac{1}{p} = E(Z)}$$

Hence proved..

$$Var(Z) = E(Z^2) - E(Z)^2.$$

$$E(Z) = 1/p.$$

$$E(Z^2) = E(Z^2) = \sum_{k=1}^{\infty} k^2 \, p \, (1-p)^{k-1}$$

$$= \sum_{1}^{\infty} (k(k-1) + k)(1-p)^{k-1} \, p$$

$$= \sum_{1}^{\infty} p(1-p) \, k(k-1)(1-p)^{k-2} + \overbrace{\left( \sum_{1}^{\infty} p k (1-p)^{k-1} \right)}^{E(Z)}$$

$$= p(1-p) \sum_{k=1}^{\infty} k(k-1)(1-p)^{k-2} + E(Z)$$

Now consider again $g(x) = \sum_{k=1}^{\infty} (1-x)^{k}$

$$g'(x) = \sum_{1}^{\infty} -k(1-x)^{k-1}$$

$$g''(x) = \sum_{1}^{\infty} k(k-1)(1-x)^{k-2}$$

$$g(x) = \frac{(1-x)}{1-(1-x)} = \frac{1-x}{x} \cdot = \frac{1}{x} - 1 = g(x).$$

$$g'(x) = -\frac{1}{x^2} \qquad \Rightarrow \boxed{g''(x) = \frac{2}{x^3}}$$

$$\Rightarrow \boxed{\frac{2}{p^3} = \sum_{1}^{\infty} k(k-1)(1-p)^{k-2}}$$

$$\Rightarrow E(Z^2) = \frac{2p(1-p)}{p^3} + \frac{1}{p} = \frac{2}{p^2} - \frac{2}{p} + \frac{1}{p} = \boxed{\frac{2}{p^2} - \frac{1}{p}}$$

$$\Rightarrow Var(z) = E(z^2) - E(z)^2$$

$$= \frac{2}{p^2} - \frac{1}{p} - \frac{1}{p^2} = \frac{1}{p^2} - \frac{1}{p} = \boxed{\frac{(1-p)}{p^2}}$$

$$\Rightarrow \boxed{Var(z) = \frac{1-p}{p^2}} \quad \text{Variance of Geometric}$$
$$\text{(Ans)} \qquad \text{random}$$
$$\text{variable.}$$

(d) $E(x^{(n)}) = \sum_{1}^{n} E(x_i)$

$$= \sum_{i=1}^{n} \frac{1}{P_i} \qquad\qquad P_i = \frac{n-i+1}{n}$$

$$= \sum_{i=1}^{n} \frac{n}{n-i+1}$$

$$\boxed{E(x^{(n)}) = n\left(\frac{1}{n} + \frac{1}{n-1} + \cdots + \frac{1}{2} + \frac{1}{1}\right)} \quad \underline{\text{(Ans)}}$$

$$= \sum_{i=1}^{n} \left(\frac{n}{n-i+1}\right)$$

(e.) $Var(x^{(n)}) = Var\left(\sum_i^n x_i\right)$

$(x_i, x_j)$ are independent.

$$= \sum_i^n Var(x_i) + 2\sum\sum_{i<j} Cov(x_i, x_j)^{\nearrow 0}$$

$$= \sum Var(x_i)$$

$$= \sum_i^n \frac{1-p_i}{p_i^2} \qquad p_i = \frac{n-i+1}{n}$$

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

$$Var(x^{(n)}) = n\left(\sum_i^n \frac{(i-1)}{(n-i+1)^2}\right)$$

$$< n \times n \left(\sum_i^n \frac{1}{(n-i+1)^2}\right)$$

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

$$< n^2 \frac{\pi^2}{6}$$

$$\Rightarrow \boxed{Var(x^{(n)}) < \frac{n^2 \pi^2}{6}}$$

$\Rightarrow$ The required upper bound is $\boxed{\dfrac{n^2 \pi^2}{6}}$ (Ans)

f) $E(x^{(n)}) = n\left(\frac{1}{1} + \frac{1}{2} + \cdots + \frac{1}{n}\right)$ . $= n \cdot S_n$

$S_n = \frac{1}{1} + \frac{1}{2} + \cdots + \frac{1}{n}$ .

for large $n$;

$$S_n \approx \int_{1/n}^{1} \frac{1}{x} \, dx \qquad \left\{ \sum \frac{1}{i/n} \left(\frac{1}{n}\right) \right.$$

$$= \sum \frac{1}{i}$$

$$= \ln x \Big|_{1/n}^{1} \qquad = \int_{\alpha}^{1} \frac{dx}{x}$$

$$= \log 1 - \log(1/n) \qquad \left. \alpha = 1/n \right\}$$

$$= \log n$$

$$\Rightarrow \lim_{n \to \infty} S_n = \log n.$$

$$\Rightarrow E(x^{(n)}) = \Theta(n S_n)$$

$$= \Theta(n \log(n))$$

$$\Rightarrow \boxed{f(n) = n \log(n)}$$

(Ans.)

Plot of $E(x^{(n)})$ vs $n$ is attached below.

E(X^{(n)}) vs n

## Question - 2

(a)  Given $u_i \sim$ uniform $(0, 1)$
  $\Rightarrow P(u_i \leq k) = k$ for $k \in (0, 1)$
  Since $F$ is a cdf and is given invertible
  $\Rightarrow F$ is strictly increasing (since $F$ is always
  non-decreasing as it is a cdf)
  $\therefore v_i = F^{-1}(u_i)$ for $i = \{1, 2 \dots n\}$
  $\therefore F_v(x) = P(v_i \leq x)$
    $= P(F^{-1}(u_i) \leq x)$
  Since $F$ is increasing
  $F^{-1}(u_i) \leq x \Rightarrow u_i \leq F(x)$
  $\therefore F_v(x) = P(F^{-1}(u_i) \leq x)$
    $= P(u_i \leq F(x))$
    $= F(x)$
  $\Rightarrow v_i$ has the distribution defined by $F$.
  Hence, proved.


(b)  For simplicity, we assume that $F$ is strictly
  increasing instead of being non-decreasing
  ie. we assume that $F^{-1}$ exists.
  $F^{-1}$ is also strictly increasing
  We have $D = \max_x |F_e(x) - F(x)|$
  Since we assumed $F^{-1}$ exists, there exists
  $y$ such that ~~$\text{\#\#\#\#\#\#}$~~ $y = F(x)$ which is unique
  $\Rightarrow D = \max_x |F_e(x) - F(x)|$
    $= \max_{-\infty < x < \infty} |F_e(F^{-1}(y)) - F(F^{-1}(y))|$
    $= \max_{0 < y < 1} |F_e(F^{-1}(y)) - y|$

The last equality comes from the fact that for every $\infty < x < \infty$, there is a $0 \leq y \leq 1$ such that $y = F(x)$. So we can always establish a map from $x$ to $y$ such that we can index them.

Note that,
$$F_e(F^{-1}(y)) = \frac{1}{n} \sum 1(Y_i \leq F^{-1}(y))$$

$$= \frac{1}{n} \sum 1(F(Y_i) \leq y)$$

Now $\{F(Y_i)\}$ are also i.i.d random variables.

Note that,
$$P(F(Y_i) \leq y) = P(Y_i \leq F^{-1}(y))$$
$$\therefore \quad P(Y_i \leq k) = F(k)$$
$$\Rightarrow P(F(Y_i) \leq y) = P(Y_i \leq F^{-1}(y))$$
$$= F(F^{-1}(y))$$
$$= y$$
$$\Rightarrow F(Y_i) \sim \text{Uniform } (0, 1) \quad [\text{since } P(U_i \leq k) = k]$$
$$\Rightarrow F(Y_i) = U_j \quad (\text{for some } j)$$

$$\therefore \quad F_e(F^{-1}(y)) = \frac{1}{n} \sum 1(F(Y_i) \leq y)$$

$$= \frac{1}{n} \sum 1(U_j \leq y)$$

$$\Rightarrow D = \max_{0 \leq y \leq 1} \left| F_e(F^{-1}(y)) - y \right|$$

$$= \max_{0 \leq y \leq 1} \left| \frac{1}{n} \sum 1(U_j \leq y) - y \right|$$

$$\Rightarrow D = \max_{0 \le y \le 1} \left| \frac{1}{n} \sum 1(V_i \le y) - y \right|$$

$$= E$$

$\Rightarrow$ So $D$ & $E$ have the same distribution

i.e. $P(D \ge x) = P(E \ge x)$

In the general case, where $F^{-1}$ does not exist
we can use another function $G$ which
is defined as follows-
$$G(x) = \min \{ y : F(x) > y \}$$
Note that $G(x)$ is a one-one mapping
& we can replace $F^{-1}(x)$ everywhere by
$G(x)$ to get the same result.

The practical significance for this result is
given on next page.

The interpretation of the preceding result can be as follows -

Any arbitrary empirical distribution will converge to the actual distribution to the same extent as the uniform distribution w.r.t. $n$, the number of samples.

Since this relationship is transitive in nature, we can extend this relation to the empirical distribution of any two random variables provided their distribution functions are continuous.

The most remarkable property of this result is that the random variable D is independent of the underlying function F i.e. it is independent of the distribution of the original random variables.

## Question -3

Given N points $(x_i, y_i)$, we have
$$z_i = ax_i + by_i + c + \varepsilon_i$$
where $\varepsilon_i \sim N(0, \sigma^2)$

Known: $\{x_i, y_i\}$ accurately

The $\{z_i\}$ have been corrupted by noise from
$N(0, \sigma^2)$

Now, $z_i \sim N(ax_i + by_i + c, \sigma^2)$

$$P(z_i | x_i, y_i, a, b, c) = \frac{e^{-\frac{(z_i - (ax_i + by_i + c))^2}{2\sigma^2}}}{\sigma \sqrt{2\pi}}$$

$$= \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-(z_i - (ax_i + by_i + c))^2 / 2\sigma^2\right)$$

$$\log(P(z_i | x_i, y_i, a, b, c)) = -\sum_{i=1}^{N} \left[\frac{(z_i - (ax_i + by_i + c))^2}{2\sigma^2}\right]$$
$(\forall i)$
$$-n \log \sqrt{2\pi} - n \log \sigma$$

$$\Rightarrow L(x_i, y_i, z_i, a, b, c) = -\sum_{i=1}^{N} \left[\frac{(z_i - (ax_i + by_i + c))^2}{2\sigma^2}\right] - n \log \sqrt{2\pi} - n \log \sigma$$

$$\therefore \frac{\partial L}{\partial a} = +\sum_{i=1}^{N} \frac{2(z_i - ax_i - by_i - c) \cdot x_i}{2\sigma^2} = 0$$

$$\Rightarrow \sum^{N} x_i z_i = a \sum^{N} x_i^2 + b \sum^{N} x_i y_i + c \sum^{N} x_i \quad \text{———①}$$

$$\therefore \frac{\partial L}{\partial b} = +\sum_{i=1}^{N} \frac{2(z_i - ax_i - by_i - c) \cdot y_i}{2\sigma^2} = 0$$

$$\Rightarrow \sum^{N} y_i z_i = a \sum^{N} x_i y_i + b \sum^{N} y_i^2 + c \sum^{N} y_i \quad \text{———②}$$

$$\therefore \frac{\partial L}{\partial c} = +\sum_{i=1}^{N} \frac{2(z_i - ax_i - by_i - c) \cdot 1}{2\sigma^2} = 0$$

$$\Rightarrow \sum^{N} z_i = a \sum^{N} x_i + b \sum^{N} y_i + Nc \quad \text{———③}$$

## In matrix/vector form.

$$\begin{bmatrix} \sum x_i^2 & \sum x_i y_i & \sum x_i \\ \sum x_i y_i & \sum y_i^2 & \sum y_i \\ \sum x_i & \sum y_i & \sum 1 \end{bmatrix} \begin{bmatrix} \hat{a} \\ \hat{b} \\ \hat{c} \end{bmatrix} = \begin{bmatrix} \sum x_i z_i \\ \sum y_i z_i \\ \sum z_i \end{bmatrix}$$

Let $I$ denote the identity vector: $(1, 1, \underbrace{\qquad}_{n \text{ times}}, 1)$

& $\bar{X}, \bar{Y}, \bar{Z}$ denote the $\{x_i\}, \{y_i\}, \{z_i\}$ vectors

So, $$\begin{bmatrix} \bar{X}.\bar{X} & \bar{X}.\bar{Y} & \bar{X}.\bar{I} \\ \bar{X}.\bar{Y} & \bar{Y}.\bar{Y} & \bar{Y}.\bar{I} \\ \bar{X}.\bar{I} & \bar{Y}.\bar{I} & \bar{I}.\bar{I} \end{bmatrix} \begin{bmatrix} \hat{a} \\ \hat{b} \\ \hat{c} \end{bmatrix} = \begin{bmatrix} \bar{X}.\bar{Z} \\ \bar{Y}.\bar{Z} \\ \bar{Z}.\bar{I} \end{bmatrix}$$

Using Matlab for solving the above matrix.
We get $a = 10.002208$
$$b = 19.998022$$
$$c = 29.951579$$

Expected noise variance = 23.068503

Eqn.: $z = 10.002208 x + 19.998022 y + 29.951579$

Matrix form: $AV = B$

Note that here $A$ is fixed ie. known with certainty

$\Rightarrow \quad E(AV) = E(B)$

$\Rightarrow \quad A \, E(V) = E(B)$

$$\Rightarrow \begin{bmatrix} \sum x_i^2 & \sum x_i y_i & \sum x_i \\ \sum x_i y_i & \sum y_i^2 & \sum y_i \\ \sum x_i & \sum y_i & \sum 1 \end{bmatrix} \begin{bmatrix} E(\hat{a}) \\ E(\hat{b}) \\ E(\hat{c}) \end{bmatrix} = \begin{bmatrix} \sum x_i E(z_i) \\ \sum y_i E(z_i) \\ \sum E(z_i) \end{bmatrix}$$

$z_i = \hat{a} x_i + \hat{b} y_i + \hat{c} + \varepsilon \qquad$ where $\varepsilon \sim N(0, \sigma^2)$

$\Rightarrow E(z_i) = a x_i + b y_i + c + 0 = a x_i + b y_i + c$

$$\Rightarrow \begin{bmatrix} \sum x_i^2 & \sum x_i y_i & \sum x_i \\ \sum x_i y_i & \sum y_i^2 & \sum y_i \\ \sum x_i & \sum y_i & \sum 1 \end{bmatrix} \begin{bmatrix} E(\hat{a}) \\ E(\hat{b}) \\ E(\hat{c}) \end{bmatrix} = \begin{bmatrix} \sum x_i(ax_i + by_i + c) \\ \sum y_i(ax_i + by_i + c) \\ \sum (ax_i + by_i + c) \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} \sum x_i^2 & \sum x_i y_i & \sum x_i \\ \sum x_i y_i & \sum y_i^2 & \sum y_i \\ \sum x_i & \sum y_i & \sum 1 \end{bmatrix} \begin{bmatrix} E(\hat{a}) \\ E(\hat{b}) \\ E(\hat{c}) \end{bmatrix} = \begin{bmatrix} a\sum x_i^2 + b\sum x_i y_i + c\sum x_i \\ a\sum x_i y_i + b\sum y_i^2 + c\sum y_i \\ a\sum x_i + b\sum y_i + c\sum 1 \end{bmatrix}$$

$$= \begin{bmatrix} \sum x_i^2 & \sum x_i y_i & \sum x_i \\ \sum x_i y_i & \sum y_i^2 & \sum y_i \\ \sum x_i & \sum y_i & \sum 1 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \end{bmatrix}$$

$$\Rightarrow A \begin{bmatrix} E(\hat{a}) \\ E(\hat{b}) \\ E(\hat{c}) \end{bmatrix} = A \begin{bmatrix} a \\ b \\ c \end{bmatrix}$$

Since $A$ is invertible

$$\Rightarrow E(\hat{a}) = a$$
$$E(\hat{b}) = b$$
$$E(\hat{c}) = c$$

For variance,

$$\text{Var}(z_i) = \text{Var}(ax_i + by_i + c + N(0, \sigma^2))$$
$$= \sigma^2$$

$$\begin{bmatrix} (\sum x_i^2)^2 & (\sum x_i y_i)^2 & (\sum x_i)^2 \\ (\sum x_i y_i)^2 & (\sum y_i^2)^2 & (\sum y_i)^2 \\ (\sum x_i)^2 & (\sum y_i)^2 & (n)^2 \end{bmatrix} \begin{bmatrix} \text{Var}(\hat{a}) \\ \text{Var}(\hat{b}) \\ \text{Var}(\hat{c}) \end{bmatrix} = \begin{bmatrix} \sum \text{Var}(z_i) x_i^2 \\ \sum \text{Var}(z_i) y_i^2 \\ \sum \text{Var}(z_i) \end{bmatrix}$$

$$= \begin{bmatrix} \sigma^2 \sum x_i^2 \\ \sigma^2 \sum y_i^2 \\ \sigma^2 \sum 1 \end{bmatrix}$$

Question 4 :

(b)  $$JL = \prod_{j=1}^{250} \left( \frac{\sum_{i=1}^{750} \exp\left(-(v_j - t_i)^2 / 2\sigma^2\right)}{750 \, \sigma \sqrt{2\pi}} \right)$$
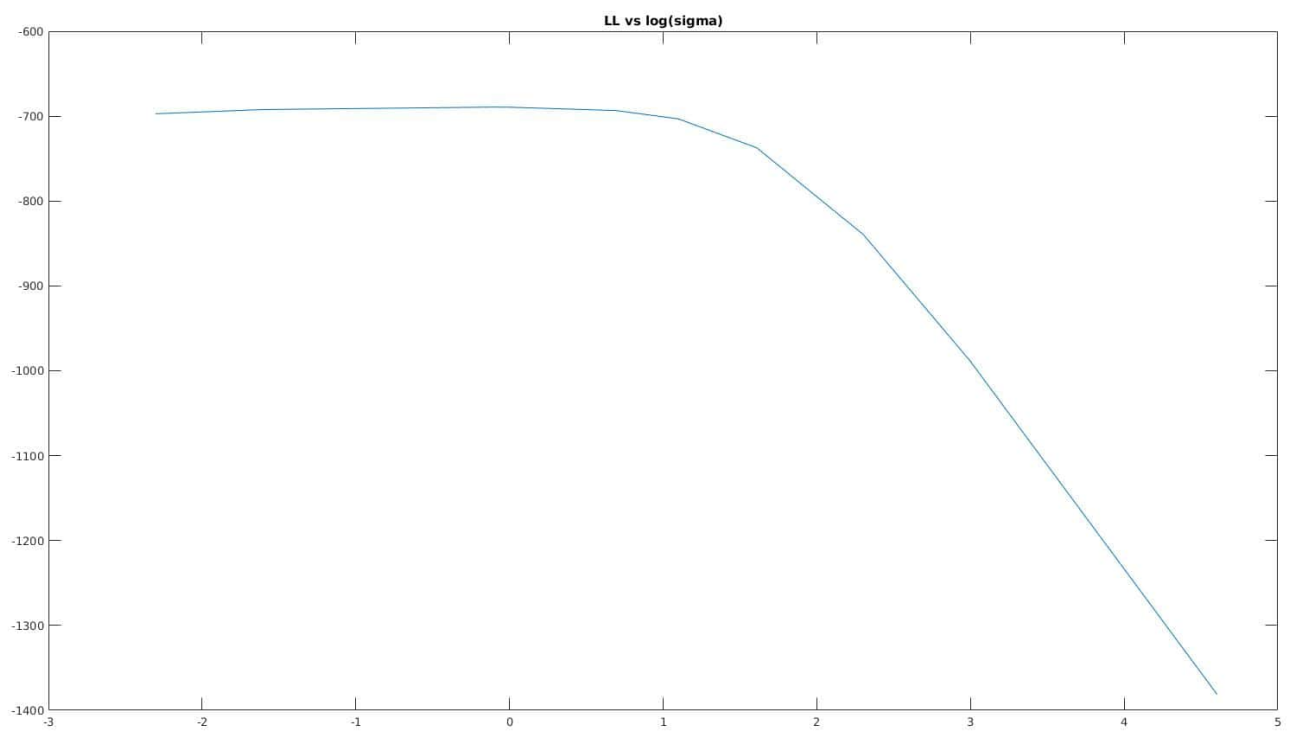
where $\{v_1, v_2, \ldots, v_{250}\}$ is set $V$.

$\{t_1, t_2, \ldots, t_{750}\}$ is set $T$.
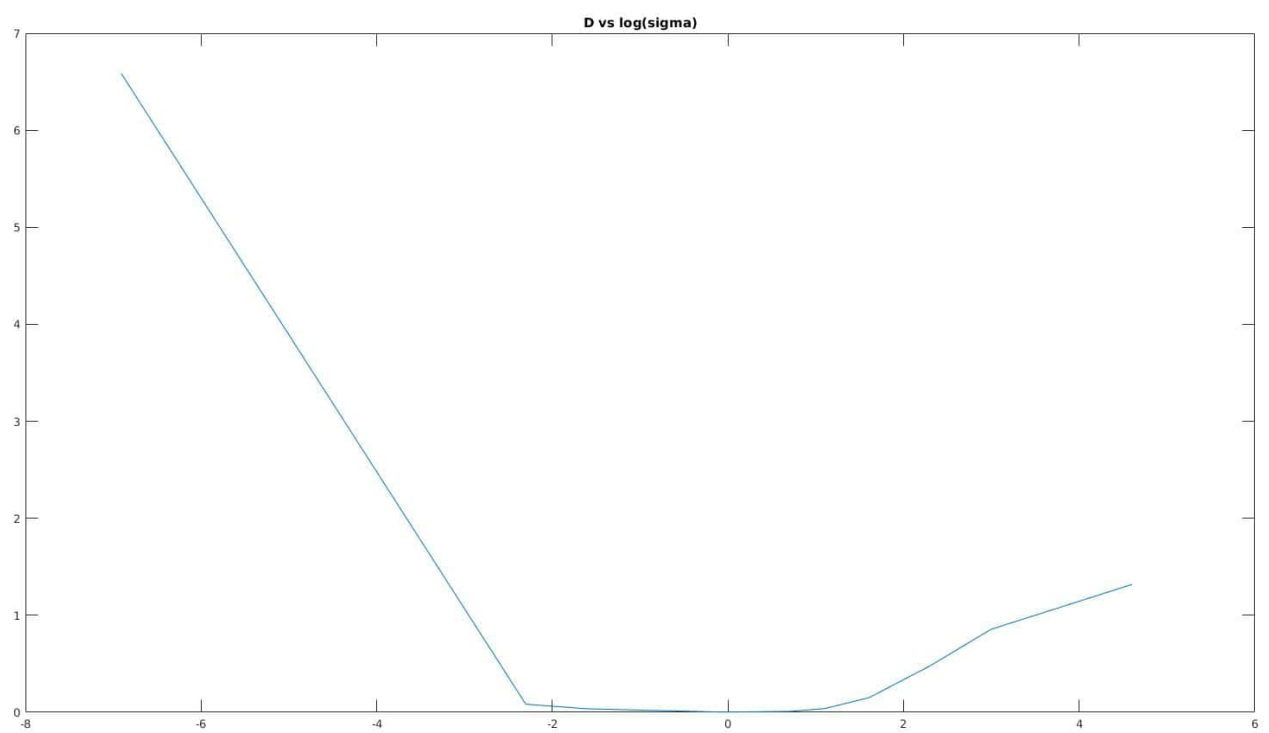
(c)  $\sigma = 0.9$  yielded the maximum value of LL (and hence the best value of LL) (since $T, V$ are randomly chosen, the best $\sigma$ value changes on repeated iterations of the same code; but it remains one of $0.9, 1$ )

Plot of graph of LL versus $\log \sigma$ is attached below.

(d)  $\sigma = 1.000$  yielded the minimum & hence best value for $D$

value of $D$ for $\sigma = 0.9$ (best LL producer) is  $0.00335$  (Ans)

LL vs log(sigma)

D vs log(sigma)

(e.) For this we first observe $\hat{p}_n$ & what $\sigma$ represents.

$$\hat{p}_n(x; \sigma) = \text{average of probabilities of occurrence of } x_i \text{ given a Gaussian centred at } x \text{ with variance } \sigma^2.$$

The graph of $\hat{p}_n(x; \sigma)$ will be <u>somewhat</u> like



$x_1 \quad x_2 \qquad x_n$

peaks like these around (not exactly) at $x_i$'s.

Now the $\sigma$ is "bandwidth", this is related to the <u>average</u> (not precisely but somewhat) <u>widths</u> of the peaks.

When cross-validation is done, if many values of set $v$ are different (or have large differences with values of $T$), then the graph of $\hat{p}_n$ will <u>widen out</u> causing larger $\sigma$. (larger sigma means that $\sigma$ which maximises LL becomes larger)

So when T & V are identical,
the graph of $\hat{p}_n$ will get more concentrated

(will have sharper peaks) at around $x_i$'s
as ~~are~~ the predictions which had before are
just reinforced.                    we

Hence for max LL, $\sigma$ will reduces as
peaks became sharper.

Thus $\sigma$ ~~red.~~ giving maximum LL
reduces (i.e. lesser than the ideal $\sigma$)

when set T & V are identical, in the
cross validation procedure.
Explanation given already!

$\int$ Also using MATLAB plots, one can confirm
that LL is ~~max~~ indeed maximised for
lower $\sigma$ }