

Name: Rohit Grarla
NET ID: rg3365

classmate

Date _____

Page _____

Homework 3.

Part I.

(a)

In Dataset I,

$$P(+)=\frac{4}{9} \quad P(-)=\frac{5}{9}$$

In Dataset II,

$$P(+)=\frac{3}{9}=\frac{1}{3} \quad P(-)=\frac{6}{9}=\frac{2}{3}$$

As we can see, the data set I is rather close to a 50-50 split, so the uncertainty to choose a class is higher.

(b)

A)

	x_1	x_2	x_3	y
$x(1)$	F	F	F	+
$x(2)$	T	F	T	+
$x(3)$	F	F	T	-
$x(4)$	F	T	F	+
$x(5)$	T	T	F	-
$x(6)$	T	T	T	-
$x(7)$	F	F	F	-

v - possible values for ~~one~~ attribute (T, F)
 S_v - subset of eq S satisfying $x_i = v$

$$\text{Infogain} = \text{Entropy}(S) - \sum_{v \in V} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\text{For } x_1, S_T = 3; S_F = 4$$

$$\text{Entropy}(S) = - \sum_{l \in L} \frac{N_l}{N} \log_2 \frac{N_l}{N}$$

l - class labels

$N_l = \text{no. of obj in } S \text{ that have label } l$

$$N_+ = 3, N_- = 4$$

$$\text{Entropy} = - \left[\frac{3}{7} \log_2 \frac{3}{7} + \frac{4}{7} \log_2 \frac{4}{7} \right]$$

$$= -\frac{1}{7} [3 \log_2 3 - 3 \log_2 7 + 4 \log_2 4 - 4 \log_2 7]$$

$$= -\frac{1}{7} \left[3 \log_2 3 + 3 \log_2 4 - 7 \log_2 7 \right]$$

$$= \log_2 7 - \frac{3}{7} \log_2 3 - \frac{4}{7} \log_2 4$$

$$= \boxed{-0.1577 - 0.1388} - \boxed{(0.5238) - (0.4613)}$$

$$\approx \boxed{0.9851}$$

$$\sum_{v \in V} \frac{|S_v|}{|S|} \text{ Entropy}(S_v) = \frac{|S_F|}{|S|} \text{ Entropy } S_F + \frac{|S_T|}{|S|} \text{ Entropy } S_T$$

$$= \frac{4}{7} \text{ Entropy}(S_F) + \frac{3}{7} \text{ Entropy}(S_T)$$

$$= \frac{4}{7} \left[-\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right] + \frac{3}{7} \left[-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right]$$

$$= \frac{4}{7} \left[-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right] + \frac{3}{7} \left[-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right]$$

$$= \frac{4}{7} [1] + \frac{3}{7} [0.9182] = \frac{4}{7} + 0.3935$$

$$= 0.9649$$

$$\text{Info Gain} = 0.9851 - 0.9649 = \boxed{0.0202}$$

(C)

We know from b that Information gain of X_1 is 0.0202

Similarly info gain of X_2 is

$$\text{Entropy}(S) = 0.9851$$

$$\sum_{v \in V} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$= \frac{4}{7} \text{Entropy}(S_F) + \frac{3}{7} \text{Entropy}(S_T)$$

$$= \frac{4}{7} \left[-2 \log \frac{2}{4} - \frac{2}{4} \log \frac{2}{4} \right] + \frac{3}{7} \left[-\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3} \right]$$

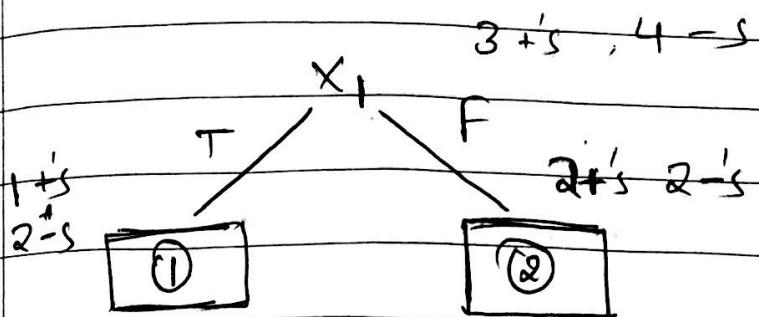
$$= 0.9649$$

Info Gain of X_2 is = 0.0202.

Info Gain of X_3 is = 0.0202.

Comparing all we can see that they are of the same information so we can choose either variable

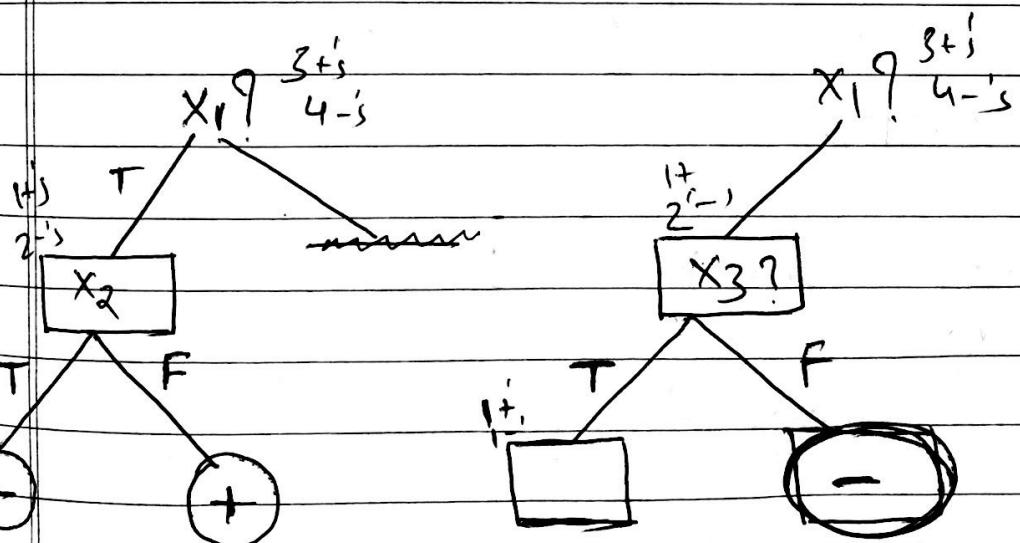
choosing X_1 as root we get.



lets measure the left subtree.

$$\text{if } X_1 = T,$$

x_1	x_2	x_3	π
T	F	T	+
T	T	F	-
T	T	T	-



Now we calculate ~~for~~ info gain for x_2 and x_3 at position ①

Info gain for x_2 at ①

$$= \text{Entropy}(S) - \sum_{v \in V} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$= \left[-\frac{2}{3} \log_2 \left(\frac{2}{3} \right) - \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right] -$$

$$\cancel{-\frac{2}{3} \log_2 \left(\frac{2}{3} \right)} - \cancel{\frac{1}{3} \log_2 \left(\frac{1}{3} \right)}$$

$$= \left[\frac{2}{3} \left[-\frac{2}{3} \log_2 \left(\frac{2}{3} \right) - 0 \log_2(0) \right] + \frac{1}{3} \left[-\frac{1}{3} \log_2 1 - 0 \log_2 0 \right] \right]$$

$$= 0.9182 - [0] + [0]$$

$$= 0.9182$$

Info gain for x_3 at ①

$$= \text{Entropy}(S) - \sum_{v \in V} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$= \left[-\frac{2}{3} \log_2 \left(\frac{2}{3} \right) - \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right] -$$

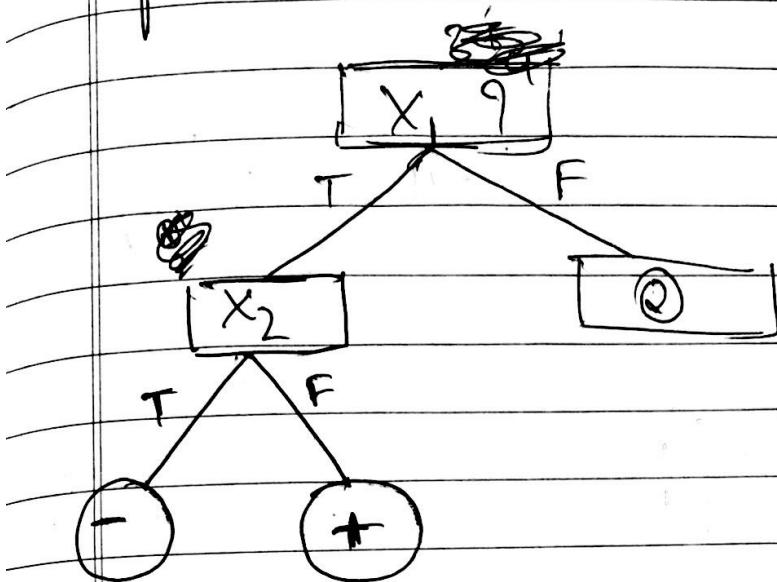
$$\left[+\frac{2}{3} \left[-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right] + \frac{1}{3} \left[-\frac{1}{1} \log_2 1 - 0 \log_2 0 \right] \right]$$

$$= 0.9182 - \left[+\frac{2}{3} [1] - \frac{1}{3} [0] \right] = 0.9182 - \frac{2}{3}$$

Infogain for x_3 at N = ~~0.2843~~ 0.2515

So, we can see that infogain is more

for x_2 at N.



Now, let's measure the right subtree

if $x_1 \neq F$

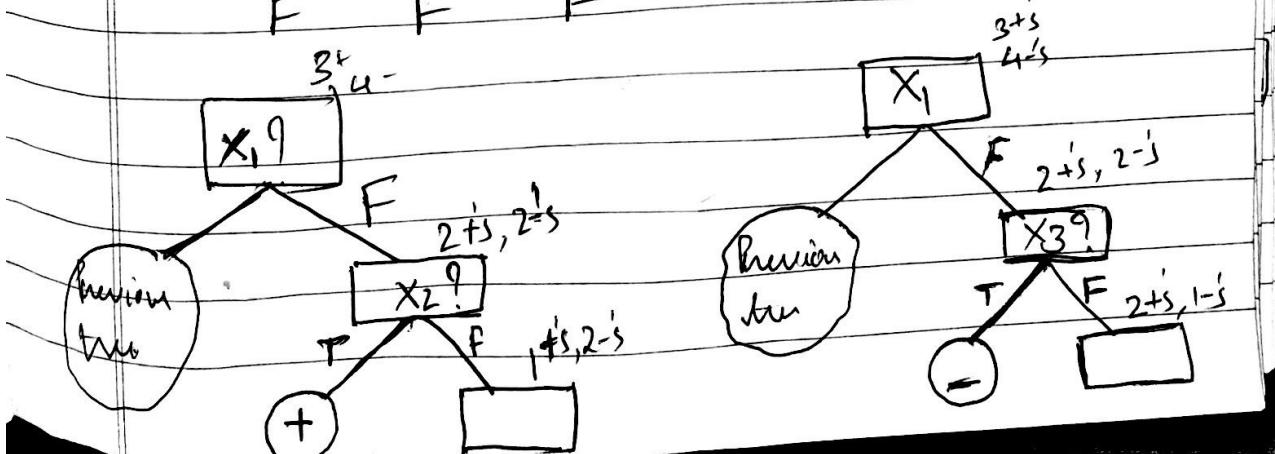
$x_1 \cdot x_2 \cdot x_3$ or

F F F +

F F T -

F T F +

F F F -



Info gain with x_2 at ②

$$= \text{Entropy}(S) - \sum_{v \in V} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$= \left[-\frac{2}{4} \log\left(\frac{2}{4}\right) - \frac{2}{4} \log\left(\frac{2}{4}\right) \right]$$

$$- \left[\frac{1}{4} (\text{Entropy}(S_T)) + \frac{3}{4} (\text{Entropy}(S_F)) \right]$$

$$= [1] - \left[\frac{1}{4} \left[\frac{1}{2} \log\frac{1}{2} - 0 \log 0 \right] + \frac{3}{4} \left[-\frac{2}{3} \log\frac{2}{3} - \frac{1}{3} \log\frac{1}{3} \right] \right]$$

$$= 1 - \left[0 + \frac{3}{4} [0.9182] \right]$$

$$1 - \frac{3}{4}(0.9182) = 0.31135$$

Info gain with $x_2 = 0.31135$

Ran

Infogain with X_3 at ②

$$= \text{Entropy}(S) - \sum_{v \in V} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

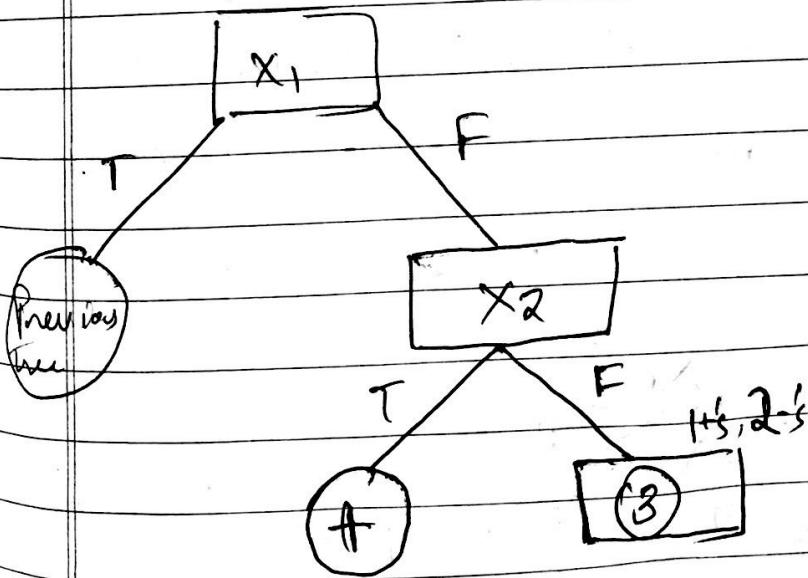
$$= 1 - \left[\frac{1}{4} \left(-\frac{1}{1} \log_2 \frac{1}{1} - \frac{0}{0} \log_2 0 \right) + \frac{3}{4} \left(\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right) \right]$$

$$= 1 - \frac{3}{4} (0.9182) = 0.31135$$

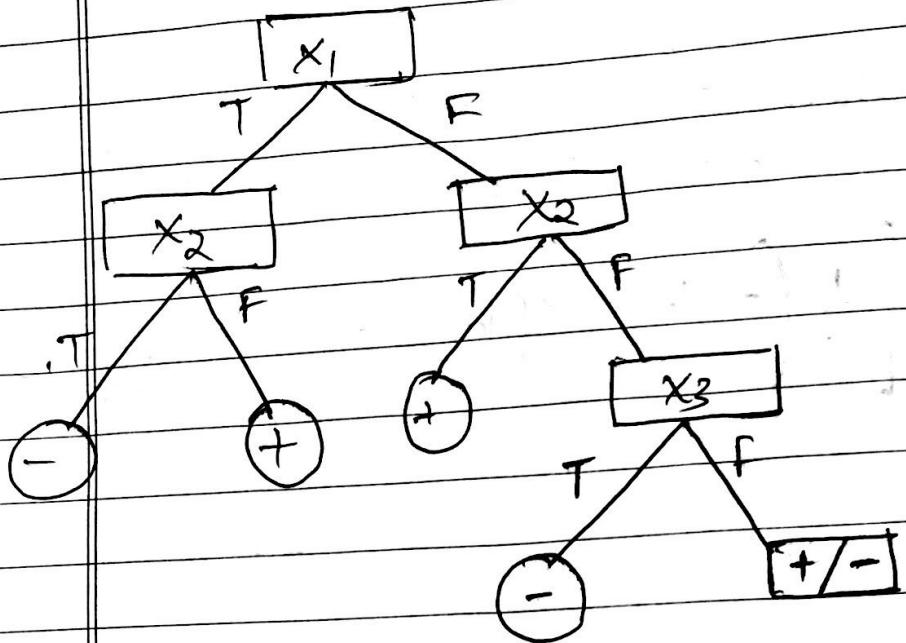
Now, we can see that both infogain are equal

so we can choose either

If we choose ~~X_2~~



Now for position ③ we only have attribute X_3 left so, the resulting tree would be



This will be the resulting tree.

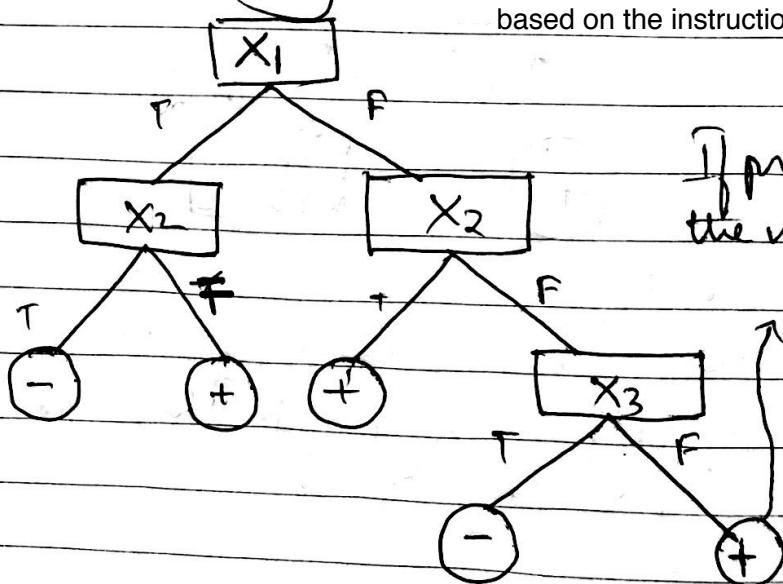
The final value in $[+/-]$ since there are two labels assigned to examples with similar attribute , so, the tree can never be 100% accurate

if we choose only one value .

Answer

There is a tie with 1+ and 1-. We choose + based on the instructions

If predicting the majority label



A) From the given information we can say that

$$H(X) = - \sum_{i=1}^n P[x=i] \log_2 P[x=i]$$

and conditional entropy is

$$H(Y/X) = \sum_x P[X=x] * \left(\sum_y -P[Y=y|X=x] * \log_2 P[Y=y|X=x] \right)$$

Now, we calculate $H(X)$

$$H(X) = - \sum_{i=1}^n P[X=i] \log_2 P[X=i]$$

$i \in \{T, F\}$

$$= \cancel{\sum_i} = - \left[P[X=T] \log_2 P[X=T] + P[X=F] \log_2 P[X=F] \right]$$
$$= - \left[\frac{3}{7} \log_2 \frac{3}{7} + \frac{4}{7} \log_2 \frac{4}{7} \right]$$

$$= 0.9851$$

$$H(Y/X) = \sum_x P[X=x] * \left(\sum_y -P[Y=y|X=x] * \log_2 P[Y=y|X=x] \right)$$

$\therefore x \in \{T, F\}, y \in \{+, -\}$

$$= \cancel{\sum_x} \left[P[X=T] * \left(- \left[P[Y=+/X=T] * \log_2 P[Y=+/X=T] \right. \right. \right. \\ \left. \left. \left. + \left[P[Y=-/X=T] * \log_2 P[Y=-/X=T] \right] \right) \right]$$

$$P[X=F] * \left(- \left[P[Y=+/X=F] * \log_2 P[Y=+/X=F] \right. \right. \\ \left. \left. + \left[P[Y=-/X=F] * \log_2 P[Y=-/X=F] \right] \right) \right)$$

$$= \frac{3}{7} * \left[- \left[\frac{1}{3} \log\left(\frac{1}{3}\right) + \frac{2}{3} \log\left(\frac{2}{3}\right) \right] \right] + \\ \frac{4}{7} * \left[- \left[\frac{2}{3} \log\left(\frac{2}{3}\right) + \frac{1}{2} \log\left(\frac{1}{2}\right) \right] \right]$$

~~$$= \frac{3}{7} [1.0223 + 0.8899] \rightarrow 0.3935 + 0.5714$$~~
~~$$= 0.96492$$~~

$$H(X) = 0.985$$

$$H(Y|X) = 0.96492$$

$$H(Y) - H(Y|X)$$

$H(Y)$ ~~means~~ \rightarrow x takes values (+, -)

$$H(Y) = - \left[P[x=+] \log_2 P[x=+] + P[x=-] \log_2 P[x=-] \right]$$

$$= - \left[\frac{3}{7} \log_2 \frac{3}{7} + \frac{4}{7} \log_2 \frac{4}{7} \right]$$

$$= 0.9851$$

0.02018

$$H(Y) - H(Y|X) = 0.9851 - 0.96492 = 0.02018$$

~~$$= 0.02018$$~~

$$= 0.0202$$

(e)

Let no. of examples be n

$$\text{Entropy} = -\sum_{i=1}^z P(x=i) \log_2 P(x=i)$$

$$= - \left[\frac{n}{z} \log_2 \frac{n}{z} \right]$$

$$= - \left[\frac{1}{z} \log_2 \frac{1}{z} \right]$$

$$= - \left[z \times \left(\frac{1}{z} \log_2 \frac{1}{z} \right) \right]$$

$$= - \left[\log_2 \frac{1}{z} \right]$$

$$\Rightarrow - \left[\log_2 (1 - \log_2 \frac{1}{z}) \right]$$

$$= \log_2 z$$

So entropy is $\log_2 z$.

Name: Rishit Gaurav

Page No. _____
Date _____

Part II : Programming exercise

Q)
a)
x)

$$A = \begin{bmatrix} 0 & 14 \\ 6 & q \end{bmatrix}$$

~~→~~ the characteristic polynomial of
matrix A

$$\det [A - kI]$$

$$\epsilon \quad \det \begin{bmatrix} 0-k & 14 \\ 6 & q-k \end{bmatrix}$$

$$\begin{aligned} &= (0-k)(q-k) \\ &= qk + k^2 - 84 \\ &= k^2 - qk - 84 \end{aligned}$$

$$-k(q-k) = 84$$

$$= k^2 - qk - 84$$

(b) The eigen values are the solutions to
the equation $\lambda^2 - 9\lambda - 84 = 0$

$$\frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{-9 \pm \sqrt{81 - 4(-84)}}{2}$$

$$\lambda_1 = 14.71, \quad \lambda_2 = -5.71$$

(c) Eigen vectors:

$$\lambda_1 = 14.71$$

$$\lambda_2 = -5.71$$

$$6x + (9 - 14.71)y = 0$$

$$6x + (9 + 5.71)y = 0$$

$$6x = 5.71y$$

$$x = -2.451y$$

$$x = 0.9516y$$

$$\sqrt{x^2 + y^2} = 1$$

$$\Rightarrow \sqrt{(0.9516y)^2 + y^2} = 1$$

$$y = 0.724$$

$$\sqrt{(-2.451y)^2 + y^2} = 1$$

$$2.6497y = 1$$

$$y = 0.3776$$

$$x = 0.9516 \times 0.724$$

$$= 0.6893$$

$$x = -2.451 \times 0.3776$$

$$= -0.9256$$

$$\Rightarrow [0.6893, 0.724]$$

$$[-0.9256, 0.3776]$$

(d) In the notebook :

(3)

a) Matrix X

$$= \begin{bmatrix} 5 & 2 & 4 \\ 9 & 6 & 4 \\ 7 & 1 & 0 \\ 2 & 5 & 6 \end{bmatrix}$$

Calculating Sample mean :

$$S_1 = \frac{\bar{5+9+7+2}}{4} = 5.75$$

$$S_2 = \frac{\bar{2+6+1+5}}{4} = 3.5$$

$$S_3 = \frac{\bar{4+6+0+6}}{4} = 3.5$$

$$B = \begin{bmatrix} -0.75 & -1.5 & 0.5 \\ 3.25 & 2.5 & 0.5 \\ 1.25 & -2.5 & -3.5 \\ -3.75 & 1.5 & 2.5 \end{bmatrix}$$

(b) Sample Covariance : } $S_{1,3}$:

$$\begin{aligned}
 &= \sum_{i=1}^{N-1} (x_i^t - S_i) * (x_j^t - S_j) \\
 &= \frac{(0.75 \times 0.5) + (3.25 \times 0.5) + ((1.25 \times 3.5) + (-3.75 \times 2.5))}{4-1}
 \end{aligned}$$

~~$= -4 \times 1.667$~~

$$= -4.1667 \dots$$

(c) The value is equal to the 13 entry which is $-4.1667 \cancel{7}$

Largest value is 12.97881

(d)

First two columns are:

0.26018674	-1.41900435
-0.8735472	4.03721245
-4.04749635	-1.8486773
-4.66084433	-0.7695308