2.

a. We are given a covariance matrix;

$$C = \begin{bmatrix} 1.6250 & -1.9486 \\ -1.9486 & 3.8750 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

We already know,

if $X = AZ + \mu$

where $X$ has $\mu$ as mean and $AA^T$ as covariance

$$\Rightarrow \quad C = AA^T$$

we have to find solutions for $A$.

Many people are attempting to find a triangular matrix $A$. But there is more better method !!! using the fact that $C$ is POSITIVE DEFINITE

⟨eigen values must be non-negative⟩.

Let

$$AA^T = C$$

Let $C = V Z V^T$

Then $A = V\sqrt{Z}$ is a solution where

$\sqrt{Z}$ is diagonal matrix with entries as square root of eigen values.

$$\boxed{A^T = (\sqrt{Z})^T V^T = \sqrt{Z} V^T}$$

We used eig method in numpy to get $V$ and $Z$.

Then used general math to find $A$.

After finding $A$ computed $X$ using $A, M$ and the random value array of length $N$ we got using random function.

After finding the $X$ array we have computed the mean and covariance for $X$.

It can be observed that as $N$ value increases the mean and covariance will be getting closer to the true mean and covariance we got which makes sense.

b.

For each N, we have to repeat the experiment 100 times draw a boxplot of the error between the true mean $\mu$ and ML estimate $\hat{\mu}_N$, where the error is $\|\mu - \hat{\mu}_N\|_L / \|\mu\|_L$ as a function of $\log_{10} N$.

For each $n = 10^l$   err $[l][k]$ ↳ This one is would
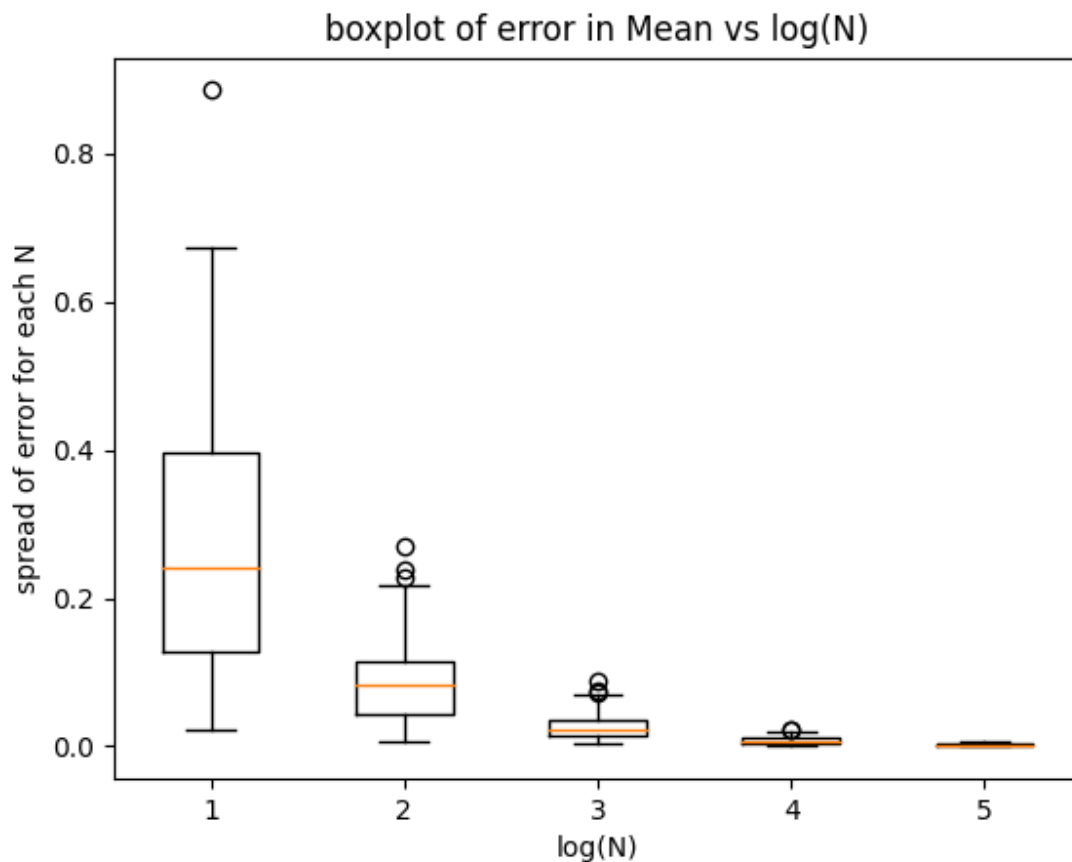signify 100 values
[errors in each case].

$\|\mu - \hat{\mu}_N\|_L / \|\mu\|_L$   is the measure of error.

→ Gen-data function is called 100-times to generate data 100 times.

OBSERVATIONS:

1. Exactly like law of large numbers for some estimate we got error becoming smaller as more data is taken, similarly here also the error for the mean estimate goes on decreasing as N is increasing exponentially.

- We can see from the above plot that as $\log_{10} N$ increases the spread of error from true mean decreases. This can be clearly seen from the box-plot graphs.

- This is like the what we saw in the case of law of large numbers in uni-variate case.

boxplot of error in Mean vs log(N)



BOX-PLOT GRAPH

For each $N$, we have to repeat the experiment 100 times draw a boxplot of the error between the true variance $C$ and ML estimate $\hat{C}_N$, where the error is $\lVert C - \hat{C}_N \rVert_{Fro} / \lVert C \rVert_{Fro}$ as a function of

$$\log_{10} N.$$

Definition:- $\lVert X \rVert_{Fro} = \sum (a_{ij})^2$

For each $n = 10^l$ err $[C][k] \hookrightarrow$ This axis would signify 100 values [errors in each case].

$\lVert C - \hat{C}_N \rVert_{Fro} / \lVert C \rVert_{Fro}$ is the measure of error.

$\rightarrow$ Gen-data function is called 100-times to generate data 100 times.

OBSERVATIONS:

boxplot of error in C vs log(N)

1. Exactly like law of large numbers for some estimate we got error becoming smaller as more data is taken, similarly here also the error for the mean estimate goes on decreasing as N is increasing exponentially.

- We can see from the above plot that as $\log_{10} N$ increases the spread of error from true covariance decreases. This can be clearly seen from the box-plot graphs.

- This is like the what we saw in the case of law of <u>large numbers</u> in uni-variate case.

---

d.

for each N

Steps :

1. Generate a data sample [In the code all the required data is actually generated all at once.]

2. Plot the 2D-scatter plot using the data generated form the above step.

3. Using eig function in numpy find the eigen values and eigen vectors for the covariance matrix C.

4. Using the above eigenvalue and eigenvector draw the line showing the principle mode of variation.

→ Principle modes of variation for drawing this here we have to pick the largest eigenvalue and the corresponding eigenvector.
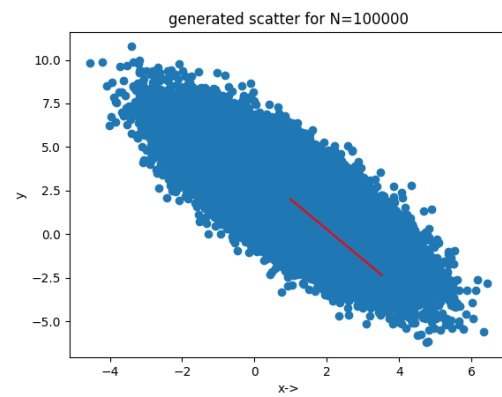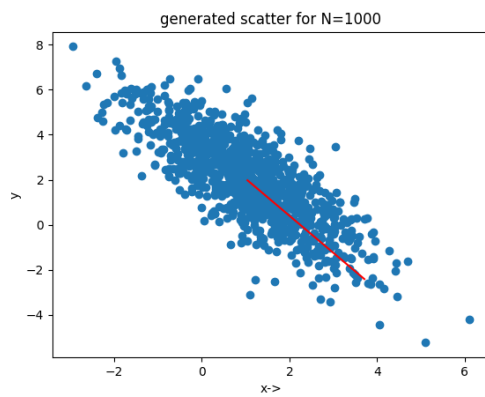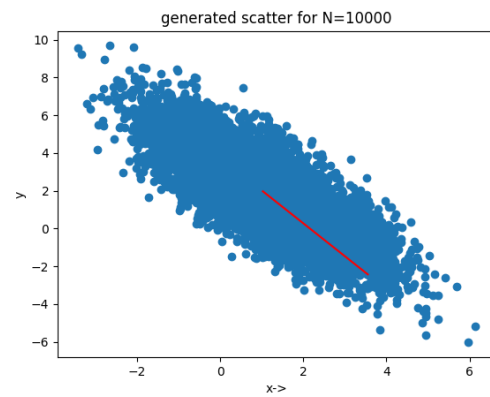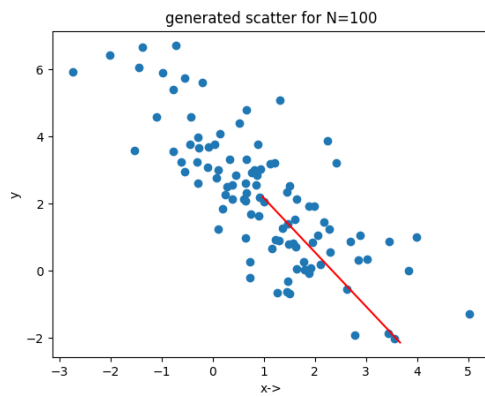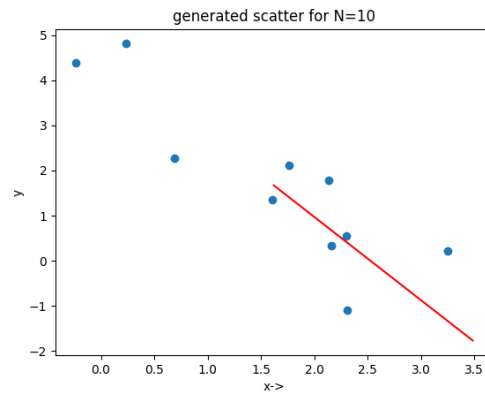
We do this because eigenvalue is directly related to variance hence large eigen value ⟹ large variance ⟹ Most of the data will be spread farthest along that direction.

∴ From the plot's we can also observe that most the data is maximum spread along that direction of line we have drawn.

Also as N is increasing we can clearly see the <u>hyper-elliptic</u> shape of the distribution and also our line staying in the direction of major axis of the ellipse (2D-case).

Which is what we expected theoritically also.

→ The following are the scatter-plots for all N ranging from 10 to $10^5$.

Scatter Plots for N = 10,100,1000,10000,100000