

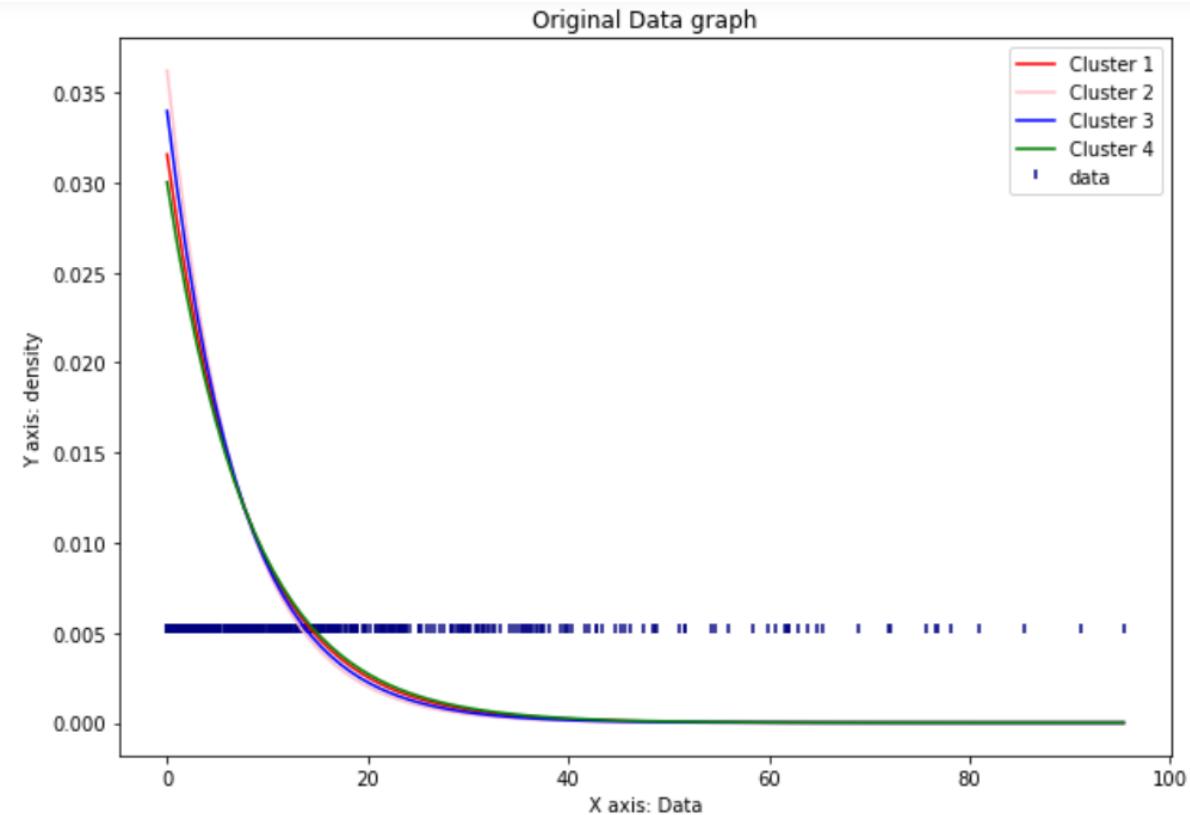
NAME : SNEHA MAGANAHALLI
RAJENDRANATH

ROLL NO: CS21M522

Machine Learning, Assignment: 2

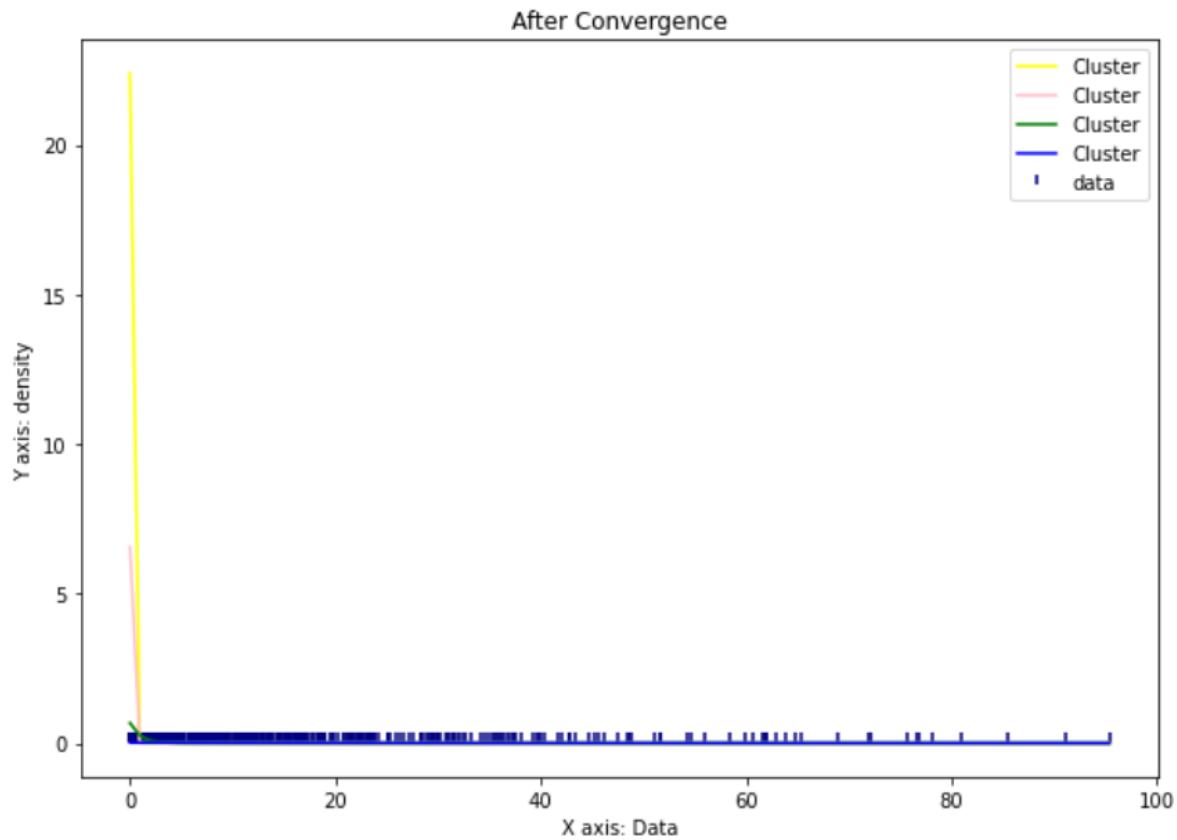
1) i)

Fig 1:



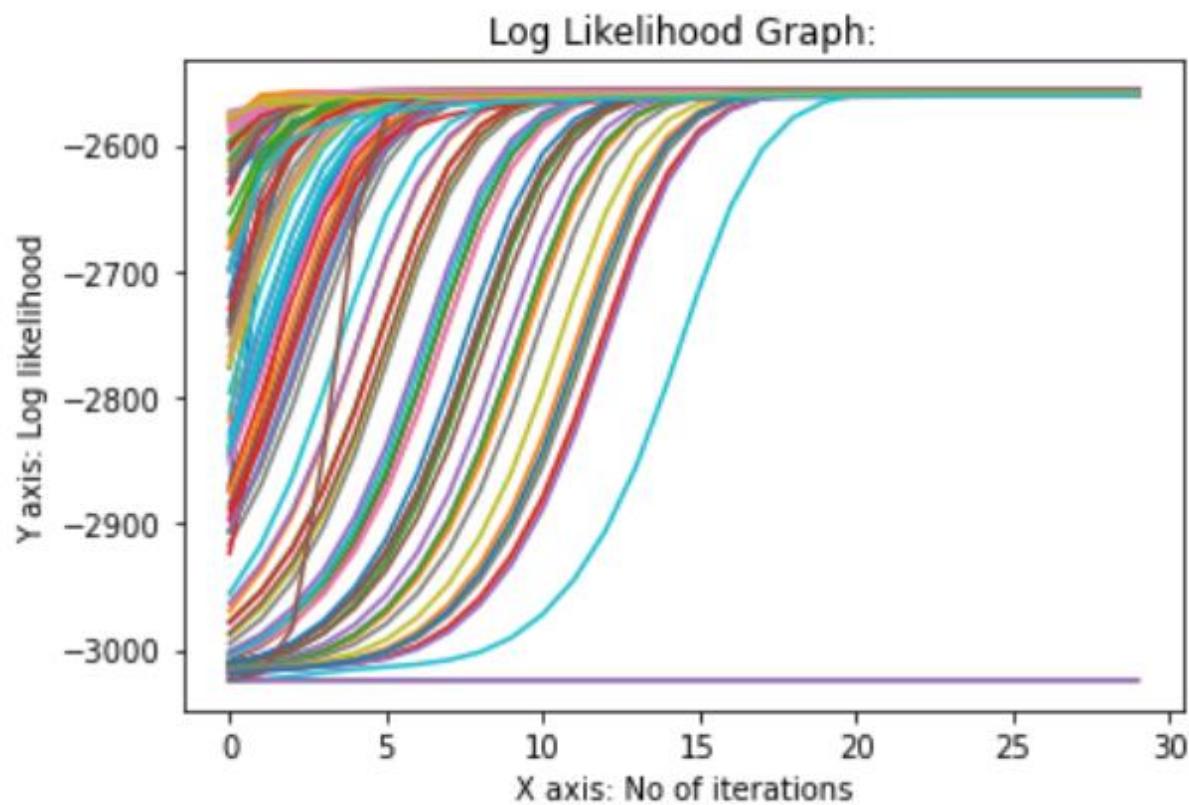
1) i)

Fig 2:



1) i)

Fig 3:



① ②
Determine which probabilistic mixture could have generated this data.
Exponential distribution has generated this data.

Fig 1: (attached)
It shows the original data (it is one-dimensional).
It shows the original data (it is one-dimensional)
It shows the original data (it is one-dimensional)

It shows before convergence (initial state)

Fig 2:
It shows after convergence.
Here we see that all the clusters are overlapped on each other at the beginning.

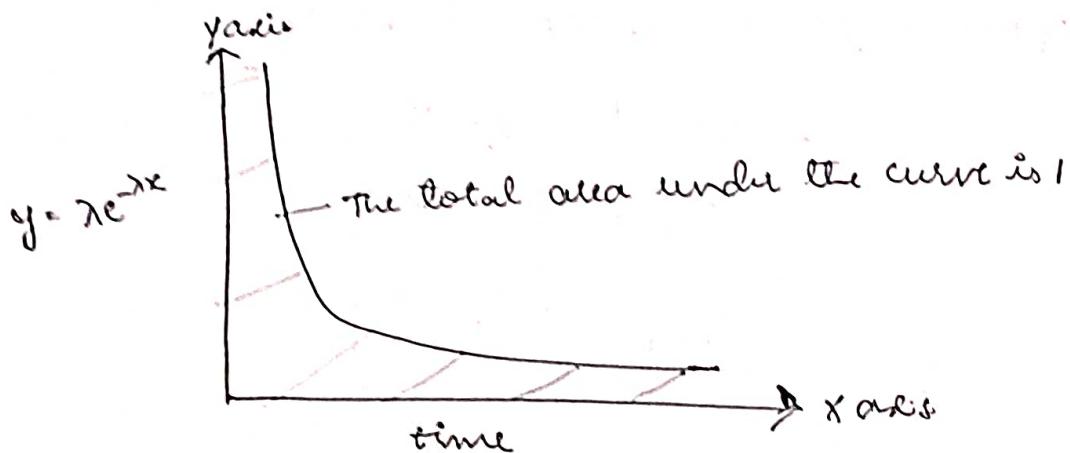
Fig 3:
It shows the log-likelihood graph for 100-different

random initialization.
we see that likelihood is 0 initially & raises slowly as the number of iterations increases. And reaches maximum to give the maximum likelihood.
Each color in the graph shows one random-initialization of data.

What is exponential distribution?

It is a statistical distribution that models the time between events.

- ↳ How long you wait before you get another text message



(amount of time between events)

probability of an event, (like person waiting for text message) happening within 0 to 5 seconds \Rightarrow check the area under the curve from $x=0$ to $x=5$ seconds

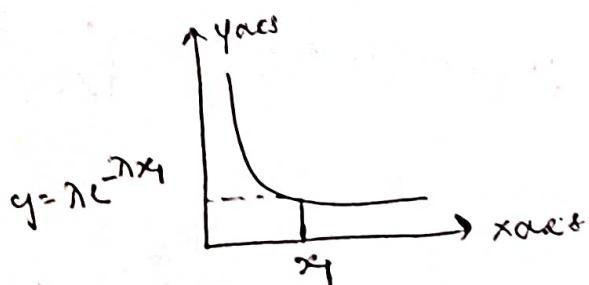
formula : $y = \lambda e^{-\lambda x}$

λ : lambda is called a rate parameter, & it is proportional to how quickly things happen

Goal of Maximum likelihood :

Given a set of measurements, to find an optimal value for λ

$$L(\lambda | x_1) = \lambda e^{-\lambda x_1} \quad \text{--- What is the likelihood of } \lambda \text{ given our first measurement is } x_1$$



$$III^4 \quad L(\lambda | x_1) = \lambda e^{-\lambda x_1}$$

What is the likelihood of λ given x_1 and x_2

$$L(\lambda | x_1) = \lambda e^{-\lambda x_1}$$

$$L(\lambda | x_2) = \lambda e^{-\lambda x_2}$$

$$L(\lambda | x_1 \text{ and } x_2) = L(\lambda | x_1) L(\lambda | x_2)$$

$$= \lambda e^{-\lambda x_1} \lambda e^{-\lambda x_2}$$

$$= \lambda^2 [e^{-\lambda x_1} e^{-\lambda x_2}]$$

$$= \underline{\lambda^2 [e^{-\lambda(x_1+x_2)}]}$$

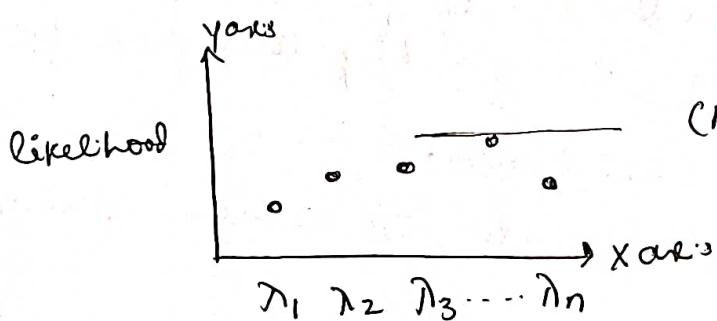
III⁴

$$L(\lambda | x_1, x_2, x_3, \dots, x_n) = \lambda^n [e^{-\lambda(x_1+x_2+\dots+x_n)}] \quad \star$$

To find the Maximum likelihood

steps:

1. Take derivative of \star
2. solve for λ when the derivative is set to be equal to 0



(At maximum likelihood, the slope (and thus the derivative) will be 0)

Step 1: $\frac{\partial}{\partial \lambda} L(\lambda | x_1, x_2, \dots, x_n) = \frac{\partial}{\partial \lambda} \lambda^n [e^{-\lambda(x_1+x_2+\dots+x_n)}]$

$$= \frac{\partial}{\partial \lambda} \log [\lambda^n [e^{-\lambda(x_1+x_2+\dots+x_n)}]]$$

{^o derivative of a function & the derivative of the log of a function is equal to 0 at same place}

$$= \frac{\partial}{\partial \lambda} \log(\lambda^n) + \log [e^{-\lambda(x_1+x_2+\dots+x_n)}]$$

$$= \frac{\partial}{\partial \lambda} n \log \lambda - \lambda(x_1+x_2+\dots+x_n) = \boxed{n \frac{1}{\lambda} - (x_1+x_2+\dots+x_n)}$$

Step 2

Set the derivative to be 0 & solve for λ .

$$0 = n \frac{1}{\lambda} - (x_1+x_2+\dots+x_n)$$

$$\frac{n}{\lambda} = x_1+x_2+\dots+x_n$$

$$\therefore \lambda = \frac{x_1+x_2+\dots+x_n}{n}$$

$$\boxed{0: \lambda = \frac{n}{x_1+x_2+\dots+x_n}}$$

... Max likelihood estimate
for λ

EM Algorithm for exponential distribution:

2 steps

1. find expected log likelihood
2. Adjust the parameter in the likelihood to maximize this.
Iterate using new parameter values

Algorithm:

Detailed steps

i) Initialize the initial probability.

$$p_i = 0.25, 0.25, 0.25, 0.25$$

ii) Compute the probability (posterior) that a point in feature space comes from a given class $r_k(x) = P(T=t | y=k)$ and equal to the

$$r_k(x) = \frac{\pi_k f_{T|Y=k}(x; \theta_k)}{\sum_{k \in K} \pi_k f_{T|Y=k}(x; \theta_k)}$$

Bayes rule as

$$r_k(x) = \frac{\pi_k f_{T|Y=k}(x; \theta_k)}{\sum_{k \in K} \pi_k f_{T|Y=k}(x; \theta_k)}$$

→ Bayes classifier

E-step:

Responsibilities at iteration m for data point n in class k is

$$\gamma_{n,k}^m = \frac{\pi_k^m f_{T_n | Y_n=k}(t_n; \gamma_k^m)}{\sum_{k \in K} \pi_k^m f_{T_n | Y_n=k}(t_n; \gamma_k^m)}$$

Where the probability that data point n occurs given that it is in class k is

$$f_{T_n | Y_n=k}(t_n; \gamma_k^m) = \gamma_{k,1}^m e^{-t_n,1} \gamma_{k,1}^m \gamma_{k,2}^m e^{-t_n,2} \gamma_{k,2}^m$$

M-step:

m -step update rules for class k

$$\pi_k^{m+1} = \frac{\sum_{n \in N} \gamma_{n,k}^m}{\sum_{n \in N, k \in K} \gamma_{n,k}^m}$$

$$\gamma_{k,i}^{m+1} = \frac{\sum_{n \in N} \gamma_{n,k}^m}{\sum_{n \in N} t_{n,i} \gamma_{n,k}^m}$$

Iterate, until the algorithm converges.

We used the new π and γ and iterate again in every step to check whether the algorithm has converged. After few iterations, maximum likelihood is found with new π, γ . After few iterations, when the algorithm converges. The likelihood will be Maximum!

The log of the likelihood is taken to plot the likelihood graph for each iteration.

Calculation:

Initial π & μ i.e. (Initial probabilities & centroid)

π , μ

0.25	7.7	→ cluster 1
0.25	6.9	→ cluster 2
0.25	7.3	→ cluster 3
0.25	8.3	→ cluster 4

Run Expectation Maximization

Run 30 iterations [It can be more or less than 30] → E step (result from M step)

- for 30 iterations [calculate exponential Pdt] → E step (result from M step)
- E calculate new π
 - calculate new μ
 - calculate log-likelihood
 - plot log-likelihood graph
- M step (result from E step)
- new π & new μ

3 stop once it converges.

E

E step:

$$\text{log gamma} = \begin{bmatrix} 2.7 & \dots & \log \\ \vdots & \ddots & 7.9 \\ 1.1 & \dots & 7.9 \end{bmatrix}_{1000 \times 4} \text{ matrix}$$

result from exponential distribution

formula:

$$f(x, \lambda) = \begin{cases} \lambda e^{-\lambda x} - x > 0 \\ 0 & x \leq 0 \end{cases}$$

M step:

$$\text{new } \pi = \begin{bmatrix} 1.1 & 1.2 & 1.3 & -1.9 \\ 3.0 & 0.25 & 0.8 & -1.9 \end{bmatrix}$$

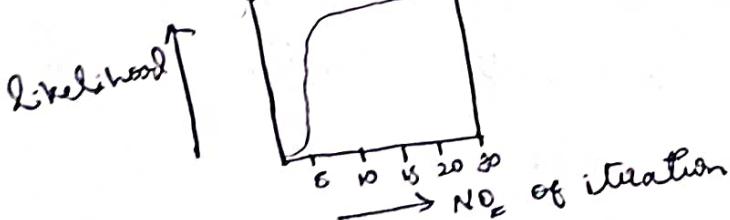
$$\text{new } \mu = \begin{bmatrix} 3.0 & 0.25 & 0.8 & -1.9 \end{bmatrix}$$

$$\text{log-likelihood } [-2985, 2870, \dots, -2556]$$

plot the graph of likelihood

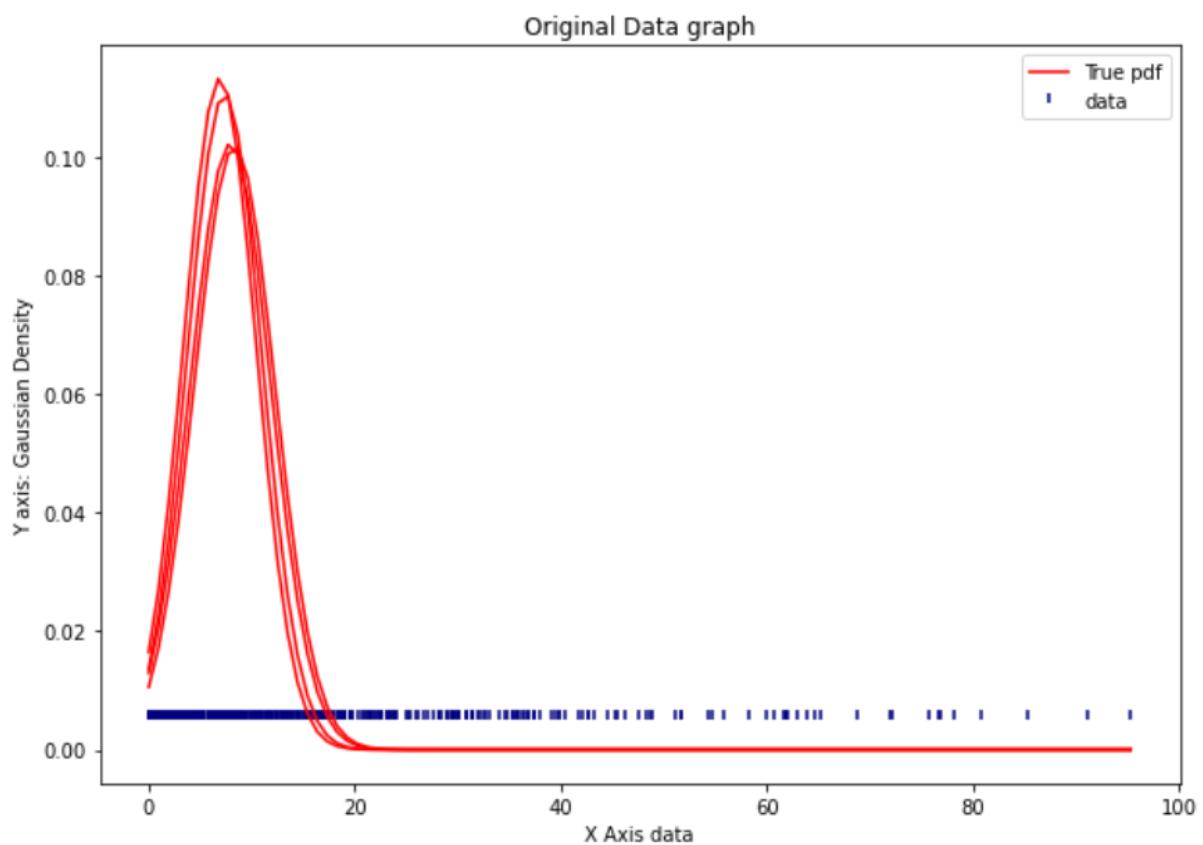
stop once it converges.

30 values
↓
No of iteration



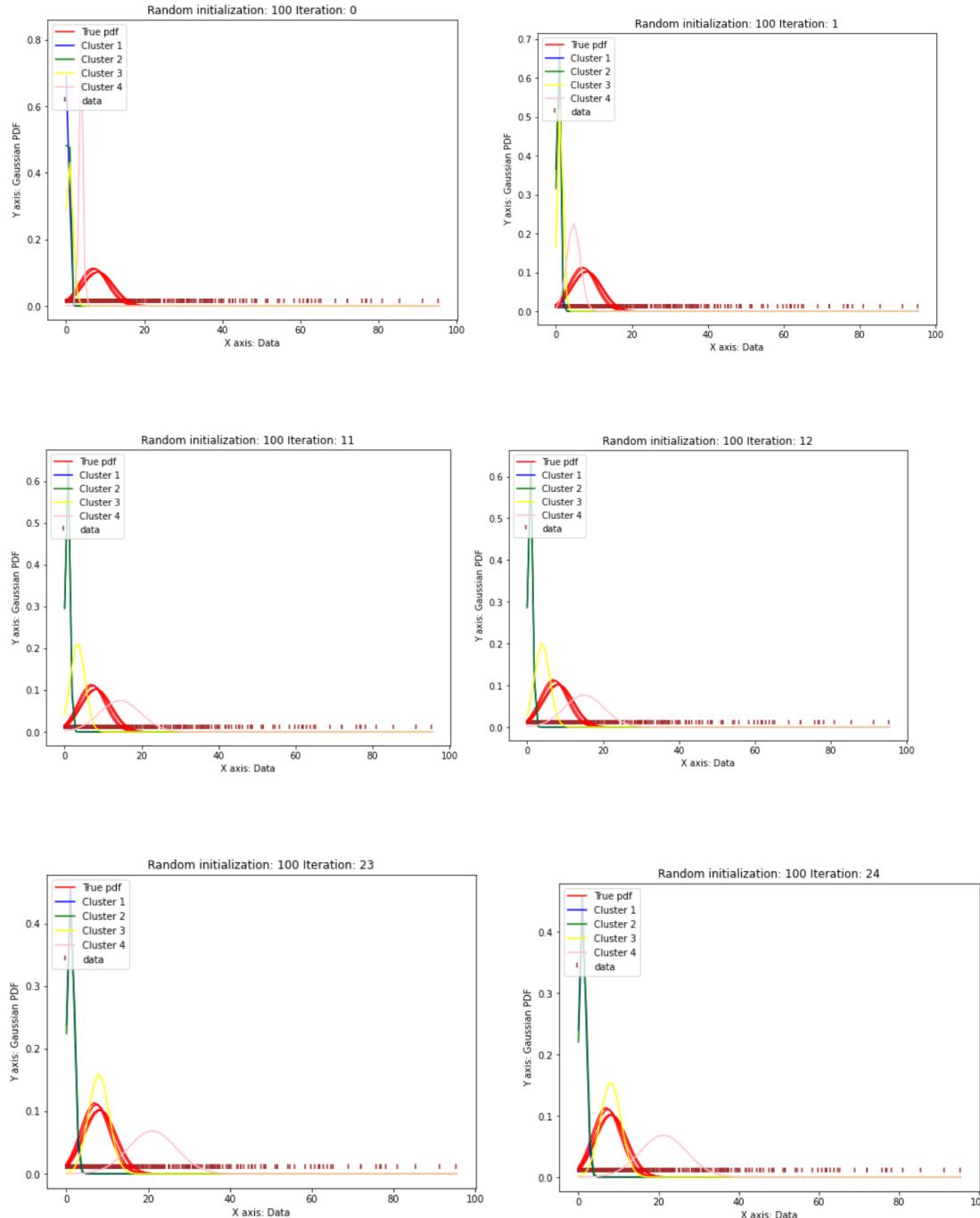
1) ii)

Fig a:



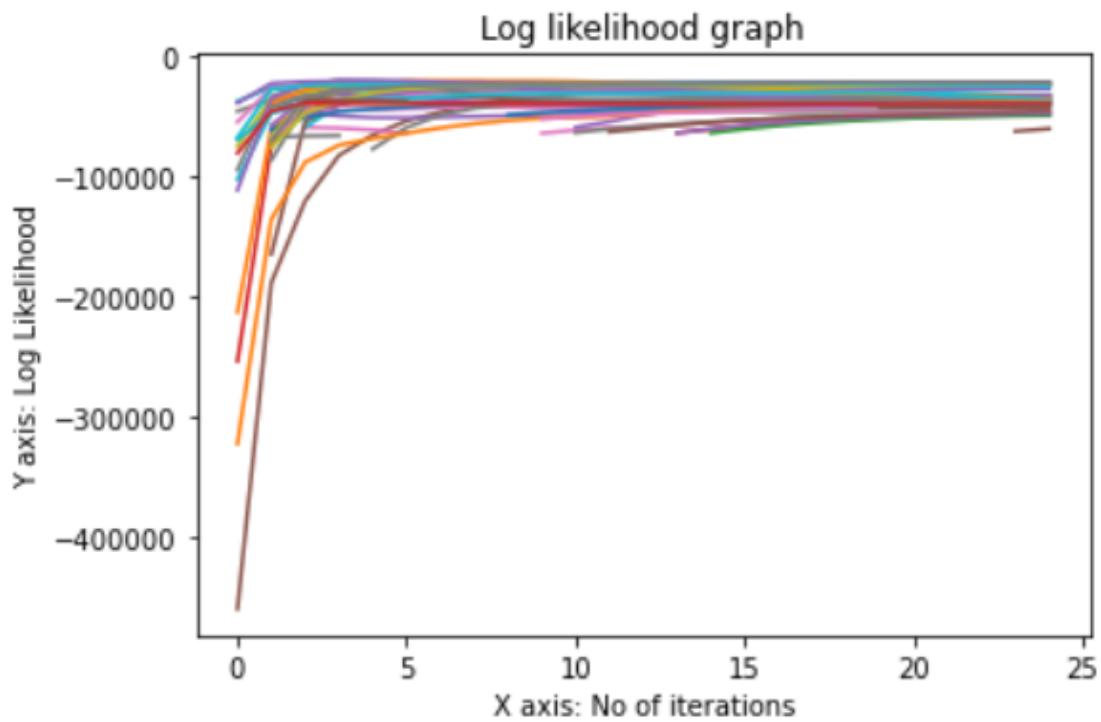
1) ii)

Fig 2:



1) ii)

Fig 3:



(1)
(ii)

Apply Gaussian distribution with 4 mixtures

Fig 1:

- Shows the original data [blue color, 1 dimension]

- Shows 4-Gaussian bell curve [Red color].
Each as initial probability as 0.25

$$0.25 + 0.25 + 0.25 + 0.25 = 1$$

$$[\text{cluster 1 probability}] + [\text{cluster 2 probability}] + [\text{cluster 3 probability}] + [\text{cluster 4 probability}] = 1$$

Fig 2:

- It shows one of the random initialization of the data [here 100th random initialization]

- Iteration 0 & 1
we see that all the four clusters are very near to y-axis (overlapping on each other)

- Iteration 11 & 12
we see that clusters are moving to their appropriate positions

- Iteration 23 & 24
we see that cluster 1 (green), cluster 2 (blue) overlap on each other. cluster 3, cluster 4 are spread on the x axis to cover all the data points.

we see that at the end of the x axis, there is no cluster, since density of the data is less.
we see more clusters at the beginning since the density is higher at the beginning.

Fig 3

It shows log-likelihood graph for 100 random initialization. Each color shows, one random initialization. we see that likelihood almost starts from a lower-value & goes on increasing as the no of iterations increase. Finally we get the maximum likelihood.

Difference between Exponential distribution mixture & Gaussian distribution mixture.

Exponential distribution mixture

Gaussian distribution mixture

Compare:

Fig ① (i) Fig ②

All the clusters are touching y axis after convergence since the data-density is higher at the beginning of the x axis.

If data density is too high at one point & too low at other point, exponential distribution mixture suits.

- ⇒ Good & better than Gaussian if data density is high at one point & decreases drastically at other point.
- ⇒ For current data, exponential distribution is better than Gaussian

① (ii) Fig ②

The clusters are spread spread along the x-axis. It is good if data density is spread-out evenly.

Good than exponential if data density is spread out evenly

⇒ For current data fairly does the job! not better than exponential

Fig ① (i) Fig 3

⇒ we see that log likelihood starts from minimum & rises slowly for each iteration.

This is true for all random initialization of data.

Better than Gaussian

Infact best!!

① (ii) Fig 3

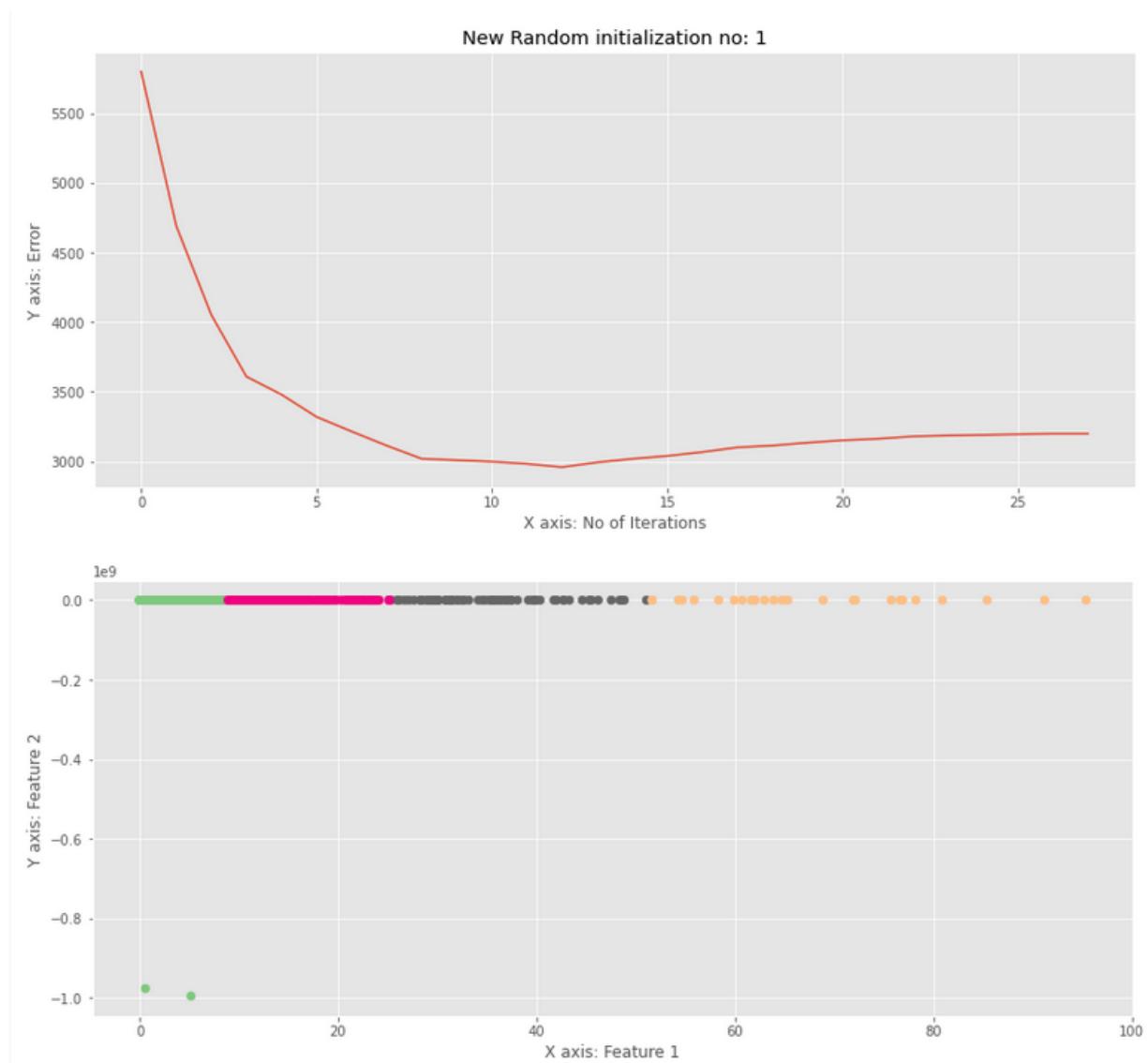
⇒ For some random initializations of data, the likelihood starts at the maximum itself or even crosses maximum point!!

Many random initialization should be considered & averaged to get the best maximum likelihood

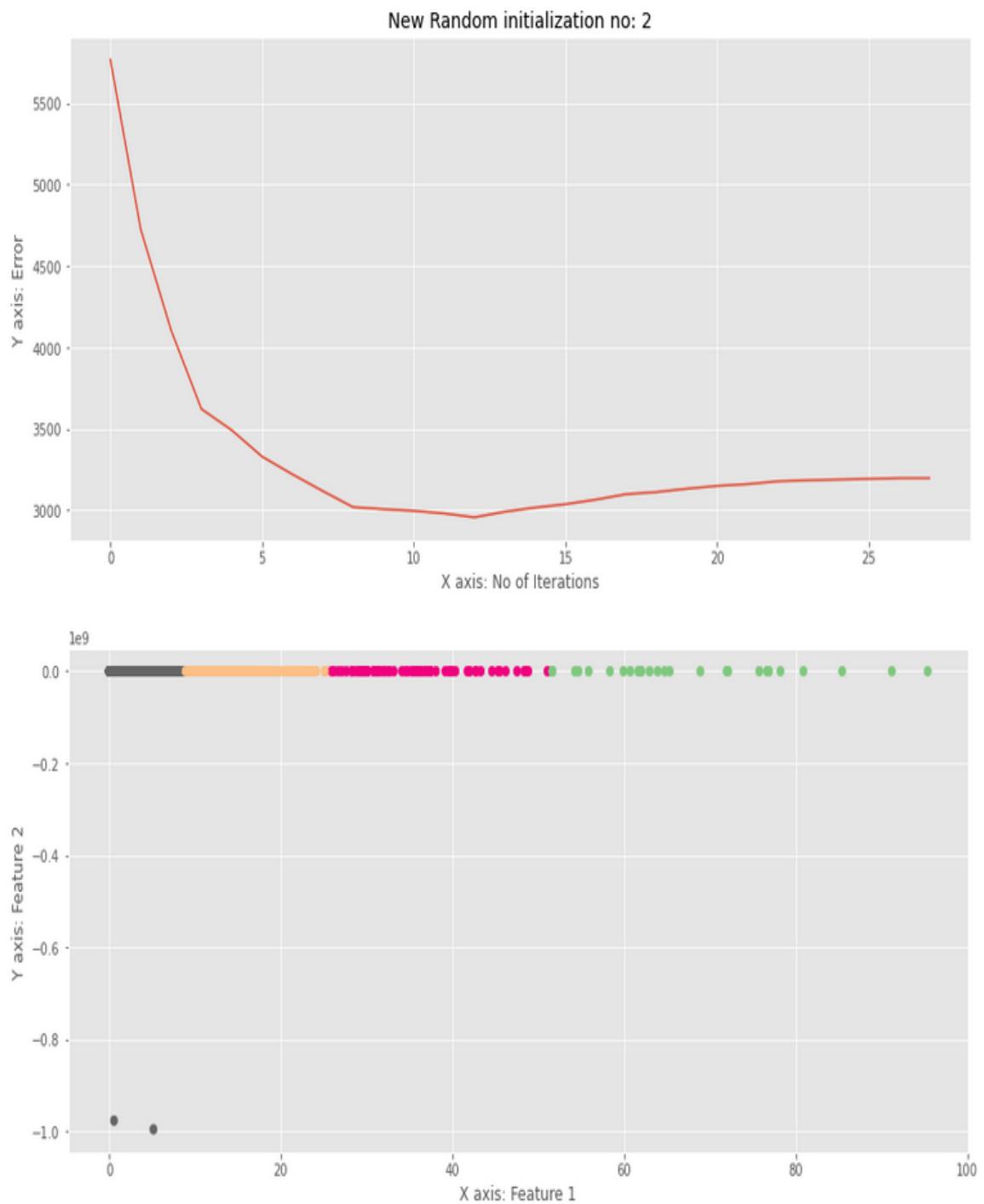
Not very good like Exponential (iC)

1) iii)

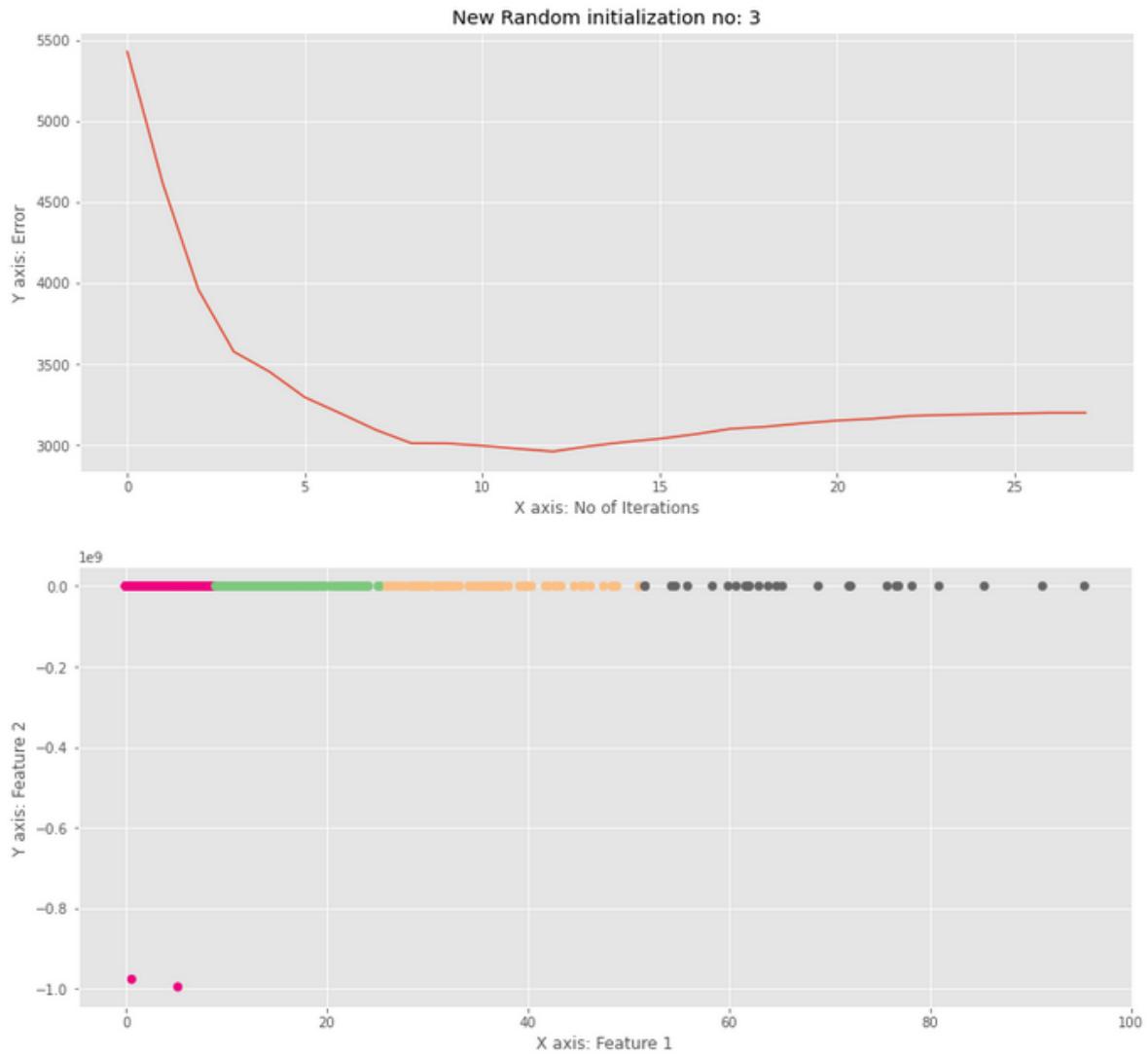
Fig 1:



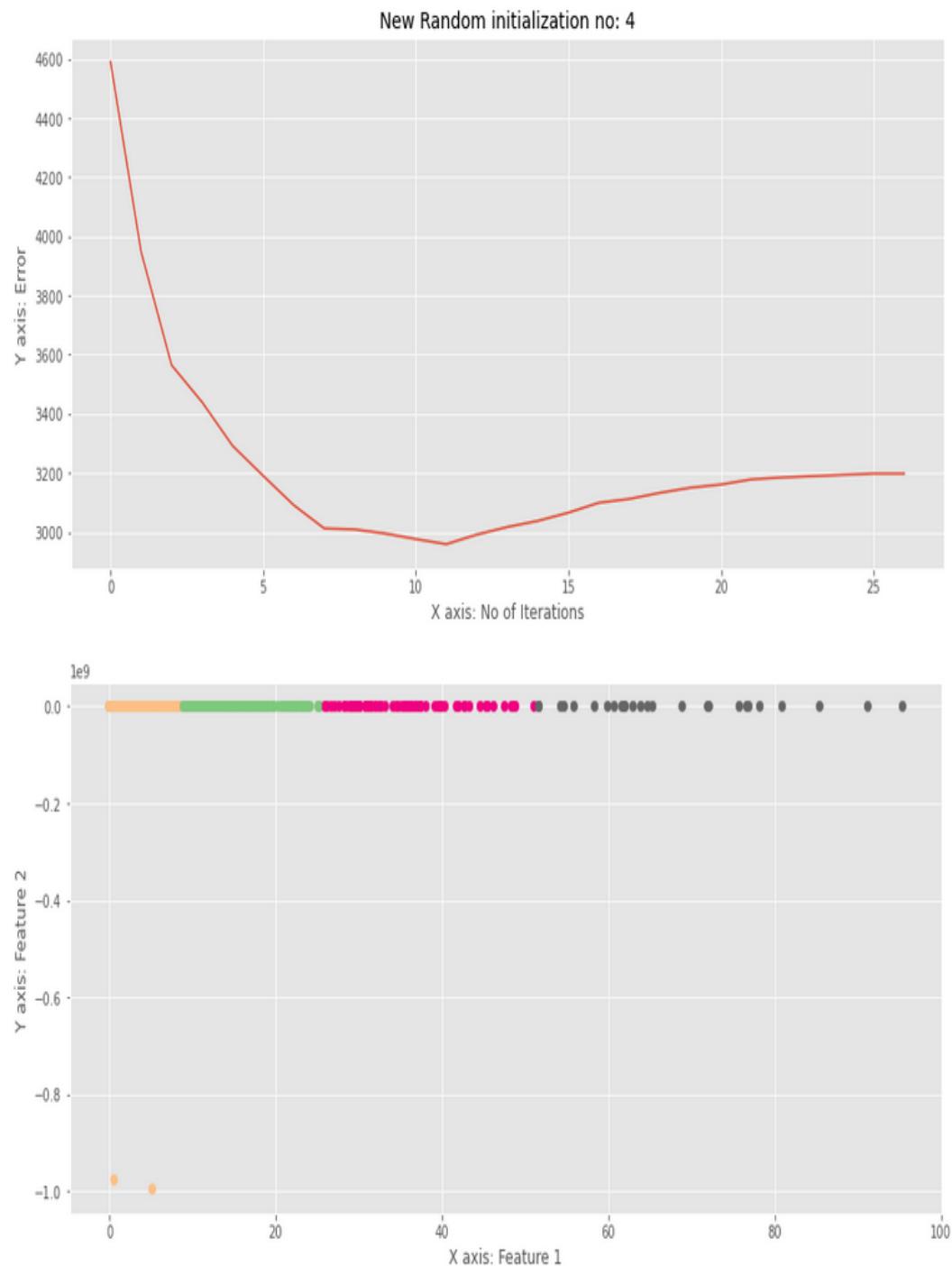
1) iii)



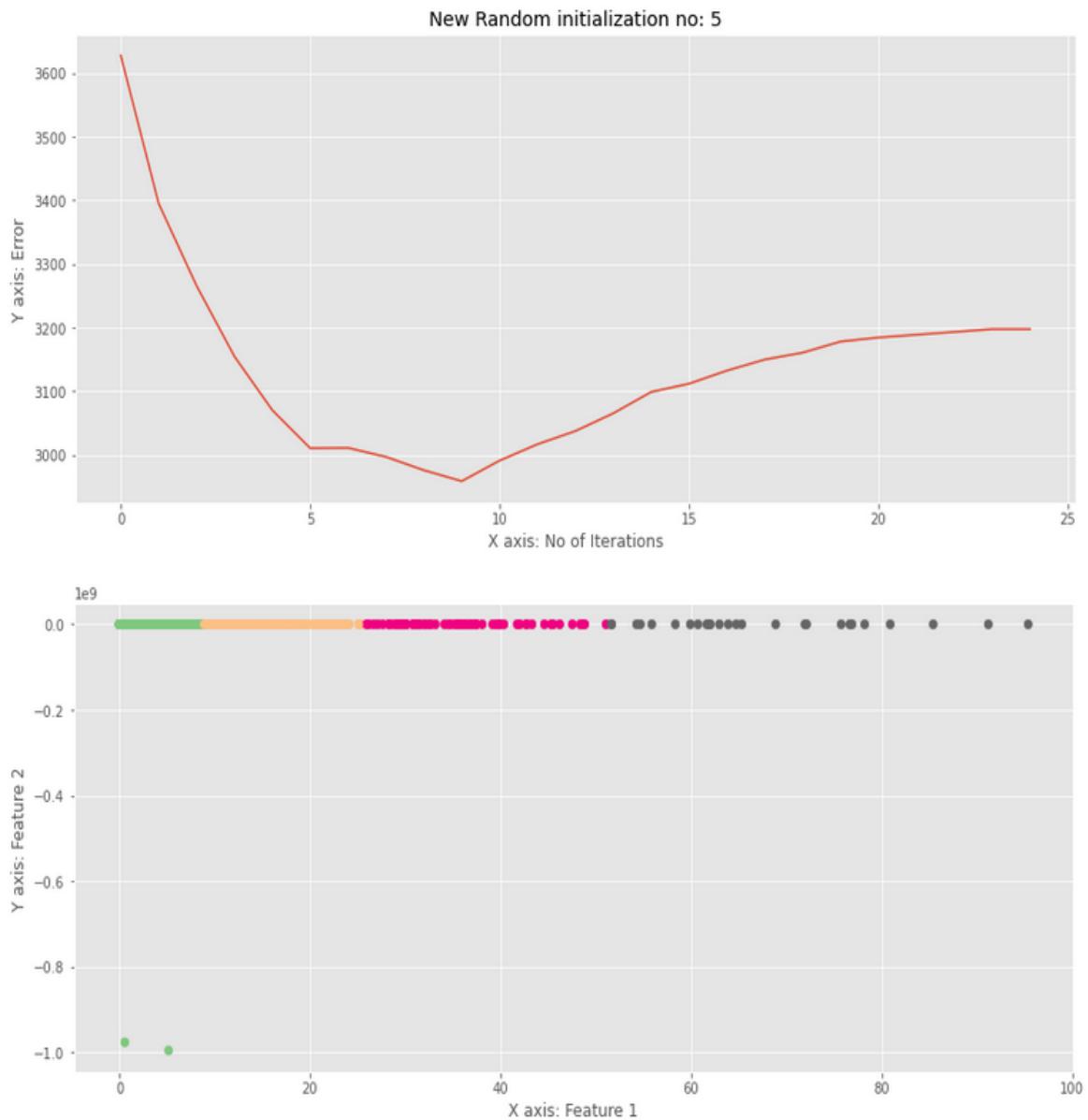
1) iii)



1) iii)



1) iii)



①

iii

K-means clustering :

K-means clustering does a hard clustering. i.e; it assigns every point to a specific cluster (& we color it) unlike Gaussian or other mixture models.

The mixture models does soft clustering. i.e; at each iteration, it will calculate the probability of the data-point belonging to a particular cluster.

Fig ①: (error function vs no of iterations)

It shows, as the number of iteration increase, the error decreases. And finally the algorithm converges. Each data-point will be assigned a particular best (nearest) cluster.

We can also see the Graph, to see that how each data-point is grouped into the best clusters.

W

Comparison between Exponential / Gaussian / K-Means for given data set.

Exponential distribution	Gaussian distribution	K-means
<ul style="list-style-type: none"> • Best 	<p>Better:</p> <ul style="list-style-type: none"> a) find k to minimize $\sum_{i=1}^n \frac{(x_i - \mu_k)^2}{\sigma^2}$ ⇒ denominator is different. It will consider the variance also while clustering. b) c) Not sensitive to initial guess of centroids d) It will analyze more complex data. e) can handle outliers * f) Since in current example data is not in circular shape, hence k-means is not good. If Gaussian is good. g) Handles any shaped data. 	<p>Good.</p> <p>find k to minimize $\sum_{i=1}^n (x_i - \mu_k)^2$</p> <p>does not consider variance for clustering</p> <ul style="list-style-type: none"> b) No mechanism to handle uncertainty when a data-point is close to more than 1 cluster centroid. c) Sensitive to initial guess of centroids. Different initialization will give different clusters. d) cannot analyse more complex data e) outliers misplace the centroid f) Data is in \mathbb{C} lar shape ⇒ k-means suite. g) Cluster created are in \mathbb{C} lar shape (∴ centroid is updated using mean value)

Exponential distribution

- i) Since the original data has a good exponential curve, exponential mixture model suits.
- ii) For each random-initialization, the likelihood starts from minimum & increases as the number of iterations increases.
- iii) No of random-initialization needed will be less when calculating maximum likelihood. (look at graph attached)
- j) converges early
- k) the initial plot of the data & after the convergence are almost similar. It says exponential distribution has done good job!

Gaussian distribution

All the cluster come at one location itself. Hence not suited if the data density statis exponentially.

Some random-initializa¹ has the higher likelihood at the beginning itself or may also cross the maximum likelihood point.

⇒ No of random initialization needed is more when calculating maximum likelihood.

j) converges little late

The initial plot of the data & after the convergence differ (cluster position). Gaussian distribution has done not pretty good job!

- o% I Rank : Exponential distribution mixture model
 II Rank : Gaussian distribution mixture model
III Rank : k-means model.

K-means

(2)

② Fitting a line to data / least squares / Linear regression.

Mean squared error on training data is 0.03968

Mean squared error on testing set is 0.370

$\begin{bmatrix} \text{Mean squared error} \\ \text{on training data} \end{bmatrix}$ is approximately equal $\begin{bmatrix} \text{Mean squared error} \\ \text{on testing data} \end{bmatrix}$

Therefore we can say that a line has ~~been~~ been fitted to the data well. MSE should be less \Rightarrow line fitted to the data is having less error.

Let us see what is MSE?

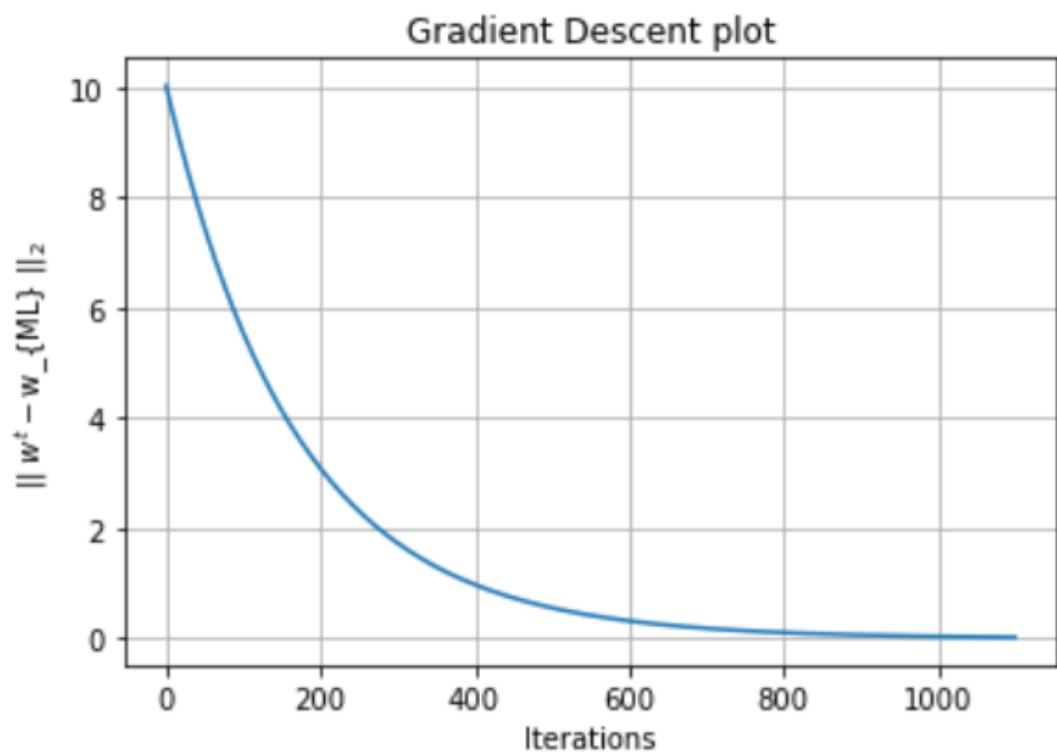
$$\begin{bmatrix} \text{Mean squared} \\ \text{error} \end{bmatrix} \text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

N : The number of samples we are testing against

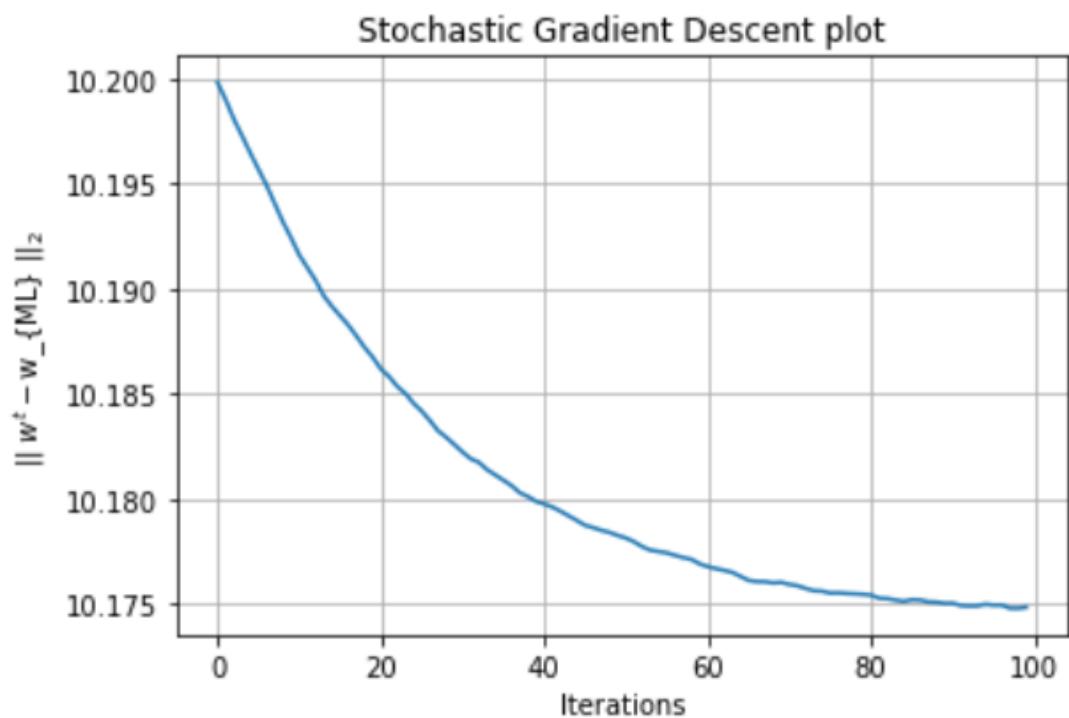
y_i : predicted Y value

y_i : Actual Y value

2) ii)



2) iii)



(2)
(ii)

By looking at the gradient descent plot, we can say that the sum of squared residuals decreases as the number of iterations increases. At a particular point [sweet spot] we see that the squared residuals is very less (almost equal to 0). At this point we should stop calculating the sum of squared residuals.

$$\left\{ \begin{array}{l} \text{MSE on training} \\ \text{data set} \end{array} \right\}_{\text{gradient descent.}} = 0.370$$

Steps:

1. Define the step size.
2. calculate the MSE. If it is not equal to 0.
calculate the step size. If MSE is approaching near to 0 then reduce the step size & calculate MSE again [iterate]
3. If MSE is almost equal to 0 \Rightarrow stop

(2)
(iii)

By looking at the stochastic Gradient Descent plot, we can say that the sum of squared residuals decreases as the number of iteration increases.

$$\left\{ \begin{array}{l} \text{MSE on training data-set} \\ \text{for Gradient descent} \end{array} \right\}_{= 0.370} < \left\{ \begin{array}{l} \text{MSE on training data-set} \\ \text{for stochastic gradient} \\ \text{descent} = 6.688 \end{array} \right\}$$

- o Gradient descent does a good job! than stochastic Gradient descent. It does not do very bad also (not too much difference) $6.6 - 0.3 = 6.3$
- o Stochastic Gradient descent algorithm should be used when there is lot of redundancy in the data.
- o Instead of calculating squared residuals for each data point. It calculates on batches (one 100 data points = 1 batch)

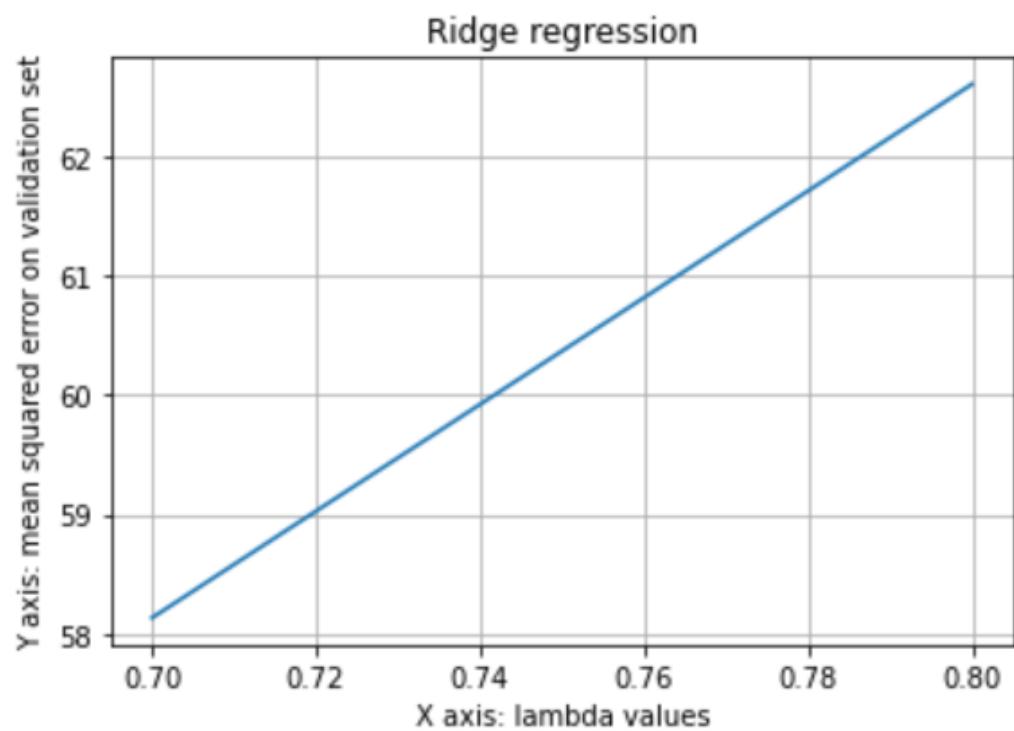
Steps for stochastic gradient descent:

- i) Make the data into batches.
- ii) Define the step size.
- iii) Calculate MSE on 1 batch.
- iv) Check the MSE.
if MSE is approaching zero, then decrease the
step size. i.e., update step size appropriately
- v) Go to step ii) & iterate.
- vi) $\{MSE\} \underset{\substack{\text{approximately} \\ \text{equal to 0}}}{\approx} \{0\}$ \Rightarrow Stop.

Advantages:

- i) stochastic gradient descent algorithm reduces lot of computation.
- ii) Very good when there is lot of data redundancy.
Or when data-points are very near to each other

2) iv)



(2)
iv

Ridge Regression :-

By looking at the ridge regression graph, as the λ -values increase, the MSE (mean squared error) also increases.

So Ridge regression for the current data set is of no-use. Without applying ridge regression algorithm, the MSE was good.

$$\begin{aligned} \left[\begin{array}{l} \text{MSE without} \\ \text{Ridge regression} \end{array} \right] &< \left[\begin{array}{l} \text{MSE with ridge} \\ \text{regression} \end{array} \right] \\ [0.3] &< [8.86] \end{aligned}$$

When Ridge Regression should be used???

- When the MSE is higher ^{for the} after data set, then Ridge regression should be used.

When the MSE will be higher??

- Sometimes, we may adjust & fit a line properly & thus decrease the MSE on the training data-set. \rightarrow [over fitting]
- But after using the testing data-set, the MSE may increase abruptly. To reduce this abrupt increase & fit the line properly for both training & testing data, we use ridge regression.
- we try to tune the λ -value so that $(\text{MSE} \approx 0)$ MSE is almost equal to 0.

P.T.O

Steps for Ridge regression:

i) Define step size initial λ value.

ii) calculate MSE

$$[MSE]_{Ridge} = MSE + [\text{Ridge penalty}] \\ (\text{includes } \lambda \text{ parameter})$$

If MSE_{Ridge} is decreasing as the lambda (λ) values increase, then ridge-regression is working bbb

if $[MSE_{Ridge}] > 0$ continue calculating MSE_{Ridge} for different λ values. (Ridge regression penalty)

Once the $[MSE_{Ridge}] \approx 0$ i.e., MSE_{Ridge} is approximately equal to zero. \Rightarrow Stop.

iii) plot the Graph.

