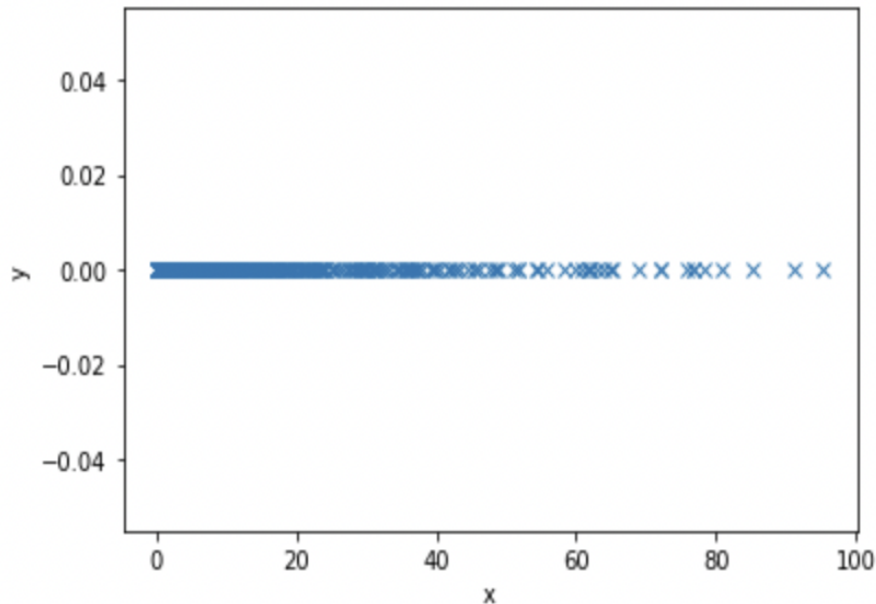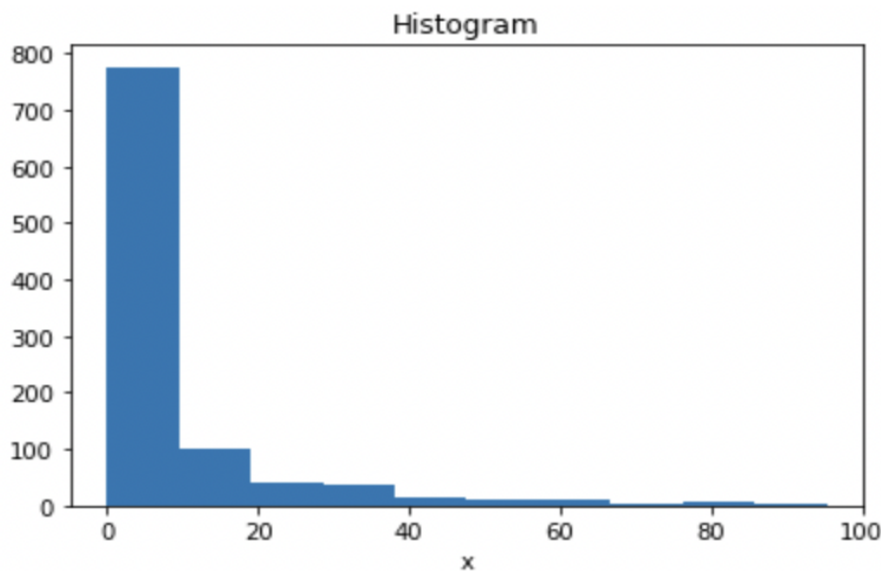## Question 1

(i)

The following image displays all the 1000 data points. It is clear from the image that most of the data points lie within (0, 20).



The histogram also suggests that most of the data points lie within the interval (0, 20).



If we sample a large number of datapoints from an exponential distribution, it is very likely that we will see a similar pattern. Hence I am assuming that the 1000 data points are generated from a mixture of exponentials.

# Mixture of exponentials :—

mixture of exponentials →

**Step1:** Pick which mixture a datapoint comes from

**Step2:** Generate datapoint from that mixture.

parameters to estimate?

$$\lambda_1, \lambda_2, \dots \lambda_K \text{ for the 'K' exponentials.}$$

$$\pi_1, \pi_2, \pi_3 \dots \pi_K$$

## Maximum Likelihood :—

$$L(\lambda_1, \lambda_2 \dots, \lambda_K, \pi_1, \pi_2, \dots \pi_K) = P(x_1, x_2, \dots x_n ; \lambda_1, \lambda_2 \dots, \lambda_K, \pi_1, \pi_2, \dots \pi_K)$$

$$= \prod_{i=1}^{n} P(x_i ; \lambda_1, \lambda_2, \dots \lambda_K, \pi_1, \pi_2, \dots \pi_K)$$

$$= \prod_{i=1}^{n} f_{mix}(x_i ; \lambda_1, \lambda_2, \dots \lambda_K, \pi_1, \pi_2, \dots \pi_K)$$

$$= \prod_{i=1}^{n} \left[ \sum_{k=1}^{K} \pi_k \, f(x_i ; \lambda_k) \right]$$

$\longrightarrow$ exponential distribution.

$$L(\theta) = \prod_{i=1}^{n} \left[ \sum_{k=1}^{K} \pi_k \, \lambda e^{-\lambda_k x_i} \right]$$

$\downarrow$
all parameters

$$\log L(\theta) = \sum_{i=1}^{n} \log \left( \sum_{k=1}^{K} \pi_k \, \lambda_k \, e^{-\lambda_k x_i} \right)$$

For every data point $i$, we have parameters

$$\{ \gamma_1^i, \gamma_2^i, \dots, \gamma_k^i \} \text{ such that } \forall i \ \sum_{k=1}^{K} \gamma_k^i = 1, \ 0 \le \gamma_k^i \le 1 \quad \forall i, k.$$

$$\log L(\theta) = \sum_{i=1}^{n} \log \left\{ \sum_{k=1}^{K} \gamma_k^i \left( \frac{\pi_k \lambda_k e^{-\lambda_k x_i}}{\gamma_k^i} \right) \right\}$$

By Jensen's inequality,

$$\log L(\theta) \ge \text{modified\_log} L(\theta, \gamma)$$

$$= \sum_{i=1}^{n} \sum_{k=1}^{K} \gamma_k^i \log \left( \frac{\pi_k \lambda_k e^{-\lambda_k x_i}}{\gamma_k^i} \right)$$

## Fix $\gamma$ and maximize over $\theta$

$$\max_{\theta} \sum_{i=1}^{n} \sum_{k=1}^{K} \gamma_k^i \log\left(\frac{\pi_k \lambda_k e^{-\lambda_k x_i}}{\gamma_k^i}\right)$$

$$\{\lambda_1, \lambda_2, \ldots \lambda_K \atop \pi_1, \pi_2, \ldots \pi_K\} \quad \sum_{i=1}^{n} \sum_{k=1}^{K} \gamma_k^i \left\{ \log\left(\pi_k \lambda_k e^{-\lambda_k x_i}\right) - \log \gamma_k^i \right\}$$

$$\sum_{i=1}^{n} \sum_{k=1}^{K} \gamma_k^i \left\{ \log \pi_k + \log \lambda_k + \log e^{-\lambda_k x_i} - \log \gamma_k^i \right\}$$

$$\sum_{i=1}^{n} \sum_{k=1}^{K} \left( \gamma_k^i \log \pi_k + \gamma_k^i \log \lambda_k - \gamma_k^i \lambda_k x_i - \underbrace{\gamma_k^i \log \gamma_k^i}_{\text{constant}} \right)$$

Taking derivative w.r.t $\lambda$, we get

$$\boxed{\frac{1}{\lambda_k} = \frac{\sum_{i=1}^{n} \gamma_k^i x_i}{\sum_{i=1}^{n} \gamma_k^i}}$$

similarly,

$$\boxed{\pi_k = \frac{\sum_{i=1}^{n} \gamma_k^i}{n}}$$

Thus when $\gamma$ is fixed, we get

$$\frac{1}{\hat{\lambda}_k^{MML}} = \frac{\sum_{i=1}^{n} \gamma_k^i x_i}{\sum_{i=1}^{n} \gamma_k^i} \quad\quad\quad (i)$$

$$\hat{\pi}_k^{MML} = \frac{\sum_{i=1}^{n} \gamma_k^i}{n} \quad\quad\quad (ii)$$

$$\sum_{i=1}^{n} \sum_{k=1}^{K} \gamma_k^i \log\left(\frac{\pi_k \lambda_k e^{-\lambda_k x_i}}{\gamma_k^i}\right)$$

$$= \sum_{i=1}^{n} \sum_{k=1}^{K} \gamma_k^i \log\left(\frac{a_{ik}}{\gamma_k^i}\right) \longrightarrow \text{constant equal to } \pi_k \lambda_k e^{-\lambda_k x_i}$$

$$= \sum_{i=1}^{n} \sum_{k=1}^{K} \left[\gamma_k^i \log a_{ik} - \gamma_k^i \log \gamma_k^i\right]$$

There is a set of $\gamma_k^i$'s for every $i$, so can maximize separately for every $i$
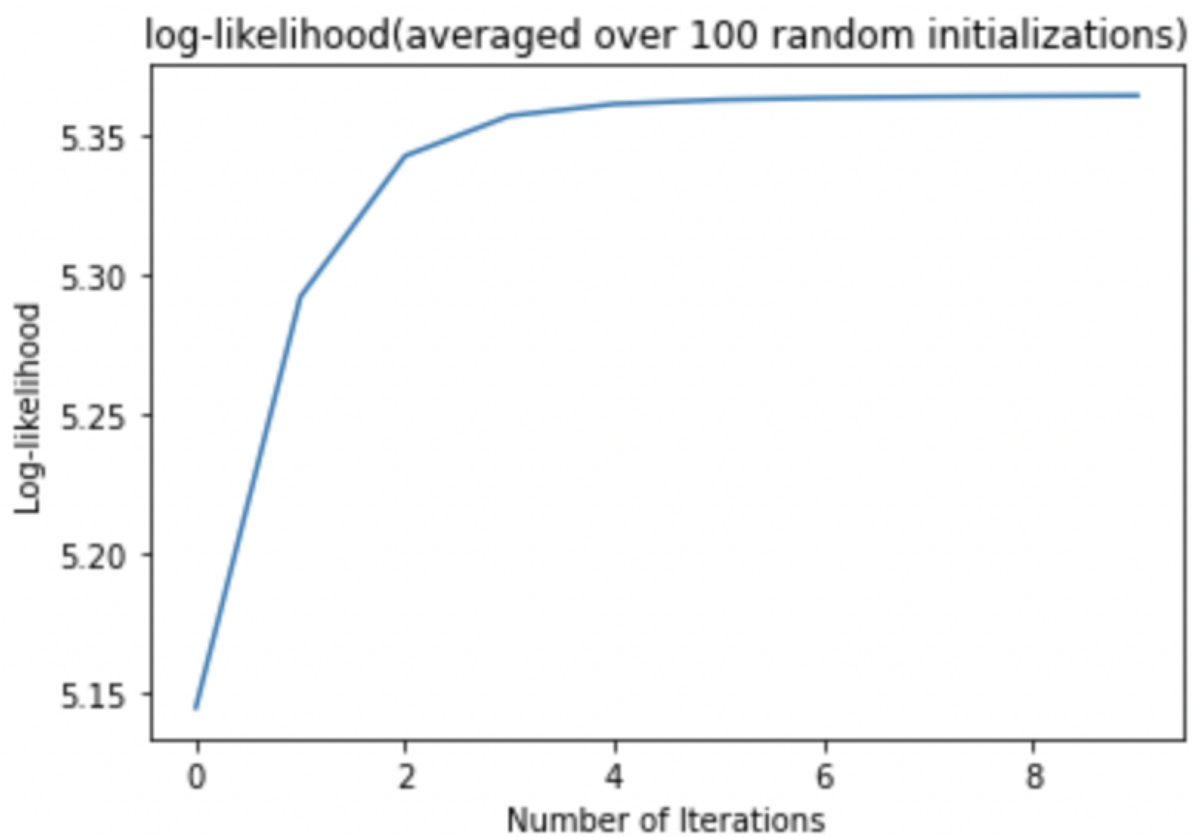
For a fixed $i$, we want

$$\max_{\gamma_1^i, \gamma_2^i, \ldots \gamma_K^i} \sum_{k=1}^{K} \gamma_k^i \log(a_{ik}) - \gamma_k^i \log \gamma_k^i \quad \text{such that } \sum_{k=1}^{K} \gamma_k^i = 1$$

$$0 \leq \gamma_k^i \leq 1$$

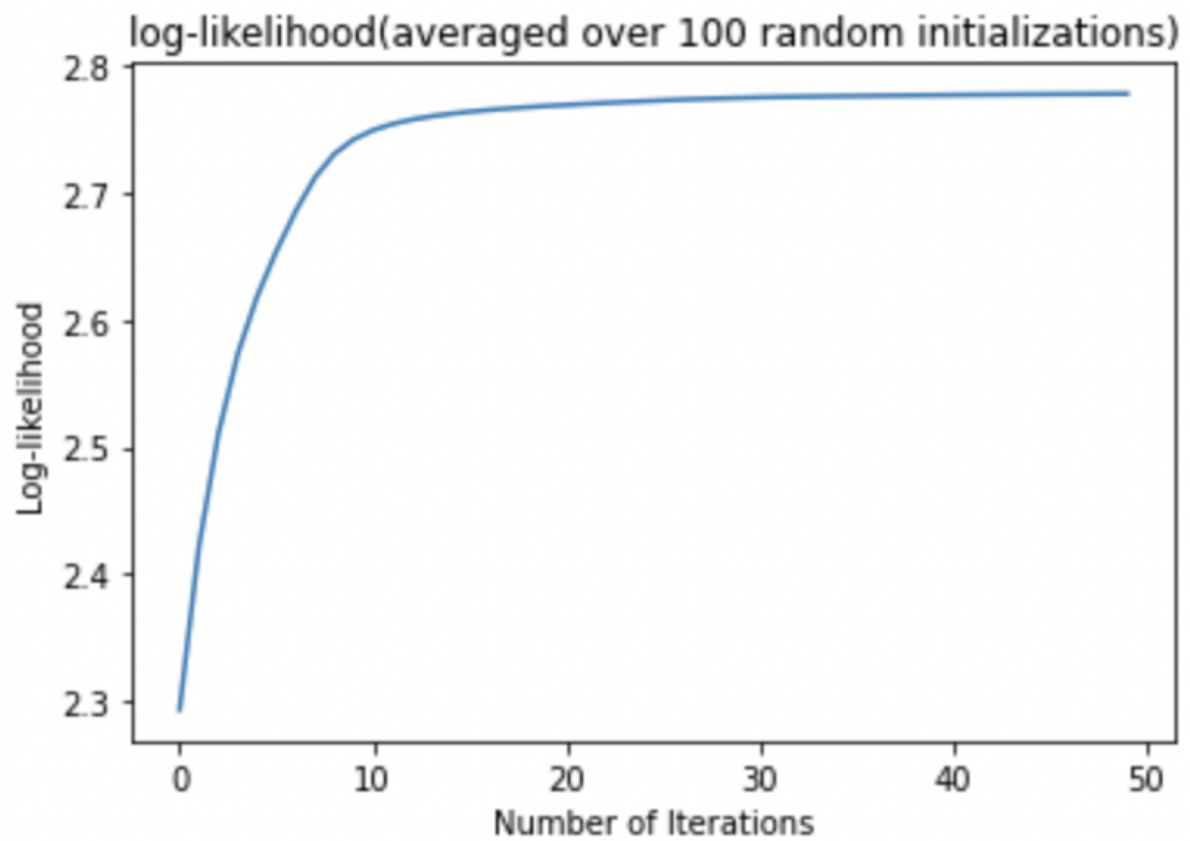this has an analytical solution that can be obtained using method of Lagrange multipliers.

$$\gamma_k^{i \cdot MML} = \frac{\lambda_k e^{-\lambda_k x_i} \cdot \pi_k}{\sum_{l=1}^{K}\left(\lambda_l e^{-\lambda_l x_i} \cdot \pi_l\right)} \quad\quad —— (III)$$

Implemented the EM algorithm for mixture of exponentials by setting the number of mixtures K = 4. Please refer to the "Mixture of Exponentials.ipynb" file for implementation.



log-likelihood(averaged over 100 random initializations)

(ii)

Implemented the EM algorithm for mixture of gaussians by setting the number of mixtures K
= 4. Please refer to the "Mixture of Gaussians.ipynb" file for implementation.
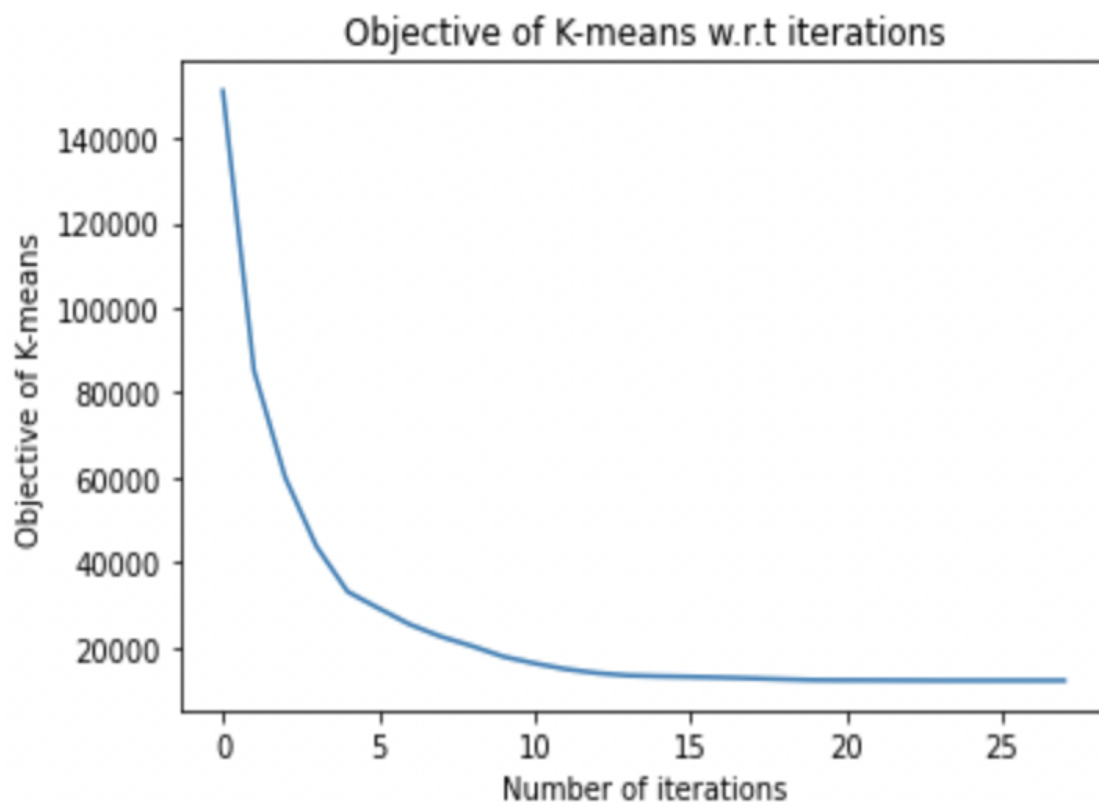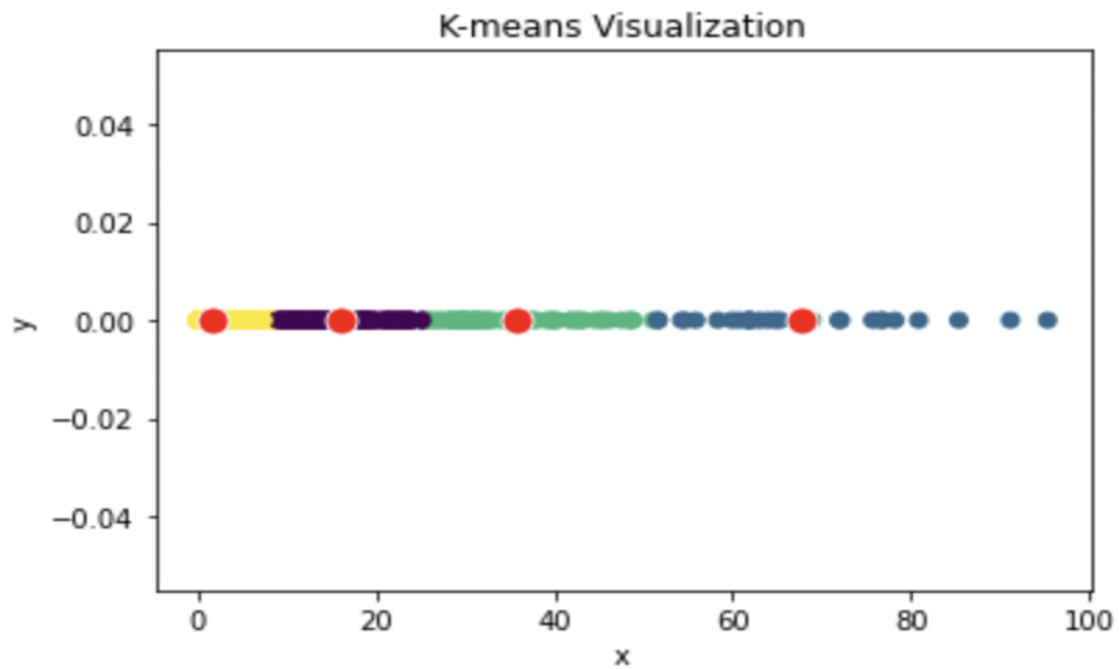
In part (i) I assumed that the data points were generated from a mixture of exponentials. I implemented the EM algorithm for a mixture of exponentials in order to estimate the parameters of the model. Our goal is to find the set of parameter values that maximize the log likelihood function. As we know after every iteration of the EM algorithm the value of the log likelihood function increases until we converge to a local maxima. The log-likelihood vs number of iterations plot for the part (i) suggests the same. For a mixture of exponentials the EM algorithm converged after 9 or 10 iterations. The average log likelihood value obtained was approximately 5.36.

In part (ii) I implemented the EM algorithm for mixture of gaussians with the same goal of estimating the parameters. In this case also the log likelihood value increased after every iteration. The log likelihood vs number of iterations plot for part(ii) suggests the same as well. For a mixture of gaussians the EM algorithm took 50 iterations to converge. The average log likelihood value obtained was approximately 2.78.

The log-likelihood value of a probabilistic model is a way to measure the goodness of the model. The higher the value of the log-likelihood, the better a model describes the data. In this case if we assume the data points were generated using a mixture of exponentials then we have an average log-likelihood value of 5.36. On the other hand if we assume a mixture of gaussians as the underlying probabilistic model, then we have an average log-likelihood value of 2.78. So a mixture of exponentials is a better probabilistic model for describing the observed data points.

(iii)

Implemented the K-means algorithm with K = 4. Please refer to the "K Means.ipynb" file for implementation.

(iv)

In part (i) I assumed that the data points were generated from a mixture of four exponentials. I estimated the parameters of the model using the EM algorithm.

In part (ii) I considered a mixture of four gaussians as the data generation story and estimated the parameters using the EM algorithm.

And in part (iii) I implemented the K means algorithm with K = 4.

There are some fundamental differences between what the EM algorithm does and what the K means algorithm does. For every datapoint K means performs hard clustering in the sense that it assigns a datapoint to a particular cluster with certainty. On the other hand the EM algorithm produces a softer version of clustering. EM gives us a probability distribution for every datapoint, which allows us to know what is the probability that a particular data point belongs to a particular cluster.

K means only considers the means in the sense that a mean describes a cluster. In EM mean and variance describes a cluster/mixture. When the clusters are mixed(which is the case in this problem), we can use the variance information to separate one cluster from another. EM has more parameters than K means, hence EM allows us to understand the data better. EM has more directional information to separate out the clusters as opposed to K means.

For the above reasons we could say that EM is a better choice for this particular problem.

We implemented the EM algorithm for two mixture models - Exponential Mixture and Gaussian Mixture. The log-likelihood value of a probabilistic model is a way to measure the goodness of the model. The higher the value of the log-likelihood, the better a model describes the data. In this case if we assume the data points were generated using a mixture of exponentials then we have an average log-likelihood value of 5.36. On the other hand if we assume a mixture of gaussians as the underlying probabilistic model, then we have an average log-likelihood value of 2.78. So a mixture of exponentials is a better probabilistic model for describing the observed data points.

Hence for this dataset, I would choose a mixture of exponentials as the underlying probabilistic model and I will run the EM algorithm in order to estimate the parameters of the model.