

# ECE421 Assignment 3 Report

## 1.1 Learning K-means

1. Implementing `distance_func()` and running the K-means Algorithm for  $K=3$  on `data2D.npy`. Some of the parameters used: `iterations=200`, use `stddev = 0.5`. K-means using the Adam Optimizer and Tensorflow 1 (`assignments(X, mu)` is a helper function)

```
# Distance function for K-means
def distance_func(X, mu):
    """ Inputs:
        X: is an NxD matrix (N observations and D dimensions)
        mu: is an KxD matrix (K means and D dimensions)

        Output:
        pair_dist: is the squared pairwise distance matrix (NxK)
    """
    X = tf.expand_dims(X, 0)
    mu = tf.expand_dims(mu, 1)
    pair_dist = tf.reduce_sum(tf.square(tf.subtract(X, mu)), 2)
    # returns KxN so transpose it because we need NxK
    return tf.transpose(pair_dist)
```

Figure 1: code for `distance_func`

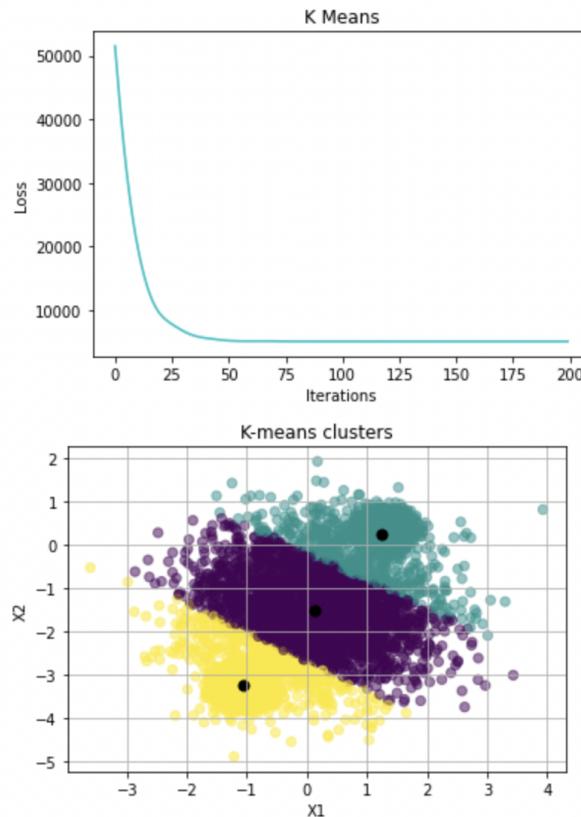


Figure 2: Loss and clusters for  $K=3$

```

def k_means(K, is_valid):
    # Loading data
    data = np.load('data2D.npy')
    #data = np.load('data100D.npy')
    [num_pts, dim] = np.shape(data)

    if is_valid:
        valid_batch = int(num_pts / 3.0)
        np.random.seed(45689)
        rnd_idx = np.arange(num_pts)
        np.random.shuffle(rnd_idx)
        val_data = data[rnd_idx[:valid_batch]]
        data = data[rnd_idx[valid_batch:]]

    loss_array = []
    loss_array_val = []
    N = num_pts
    D = dim
    iterations = 200

    # create x as a placeholder (nonexD)
    X = tf.placeholder("float", shape=[None, D])
    #X = tf.constant(data, dtype=tf.float32)
    # need standard distributions values (KxD)
    mu_std = tf.truncated_normal([K, D], stddev=0.5)
    mu = tf.Variable(mu_std)

    # calculate distance between cluster centers and points
    # https://docs.w3cub.com/tensorflow-python/tf/reduce_sum
    # loss is minimizing the distance
    pair_dist = distance_func(X, mu)
    loss = tf.reduce_sum(tf.reduce_min(pair_dist, axis=1))

    # using adamOptimizer as stated in handout
    optimizer = tf.train.AdamOptimizer(learning_rate = 0.1, beta1 = 0.9, beta2 = 0.99, epsilon=1e-5).minimize(loss)

    # init globals and start the session
    globals_init = tf.global_variables_initializer()
    sess = tf.Session()
    sess.run(globals_init)

    # start training
    for i in range(iterations):
        # for validation data
        # set X=val_data and run session
        if is_valid:
            center_val, loss_curr_val, opt_val = sess.run([mu, loss, optimizer], feed_dict={X:val_data})
            loss_array_val.append(loss_curr_val)
            # print info every 10 iterations
            if i%10 == 0:
                print("Iteration: ", i, " Loss: ", loss_curr_val)

        else:
            # for training data
            # set X=data and run session
            center, loss_curr, opt = sess.run([mu, loss, optimizer], feed_dict={X:data})
            loss_array.append(loss_curr)

            # print info every 10 iterations
            if i%10 == 0:
                print("Iteration: ", i, " Loss: ", loss_curr)

    assign_clusters = sess.run(assignments(X, mu), feed_dict={X:data, mu:center})

# this function finds the closest cluster center to a point and extracts it
def assignments(X, mu):
    dist = distance_func(X, mu)
    # extract the most small value from dist
    return tf.argmin(dist, 1)

```

Figure 3: Code for K-means for Training data and K=3

## 2. K-means for Validation Data for K={1, 2, 3, 4, 5}

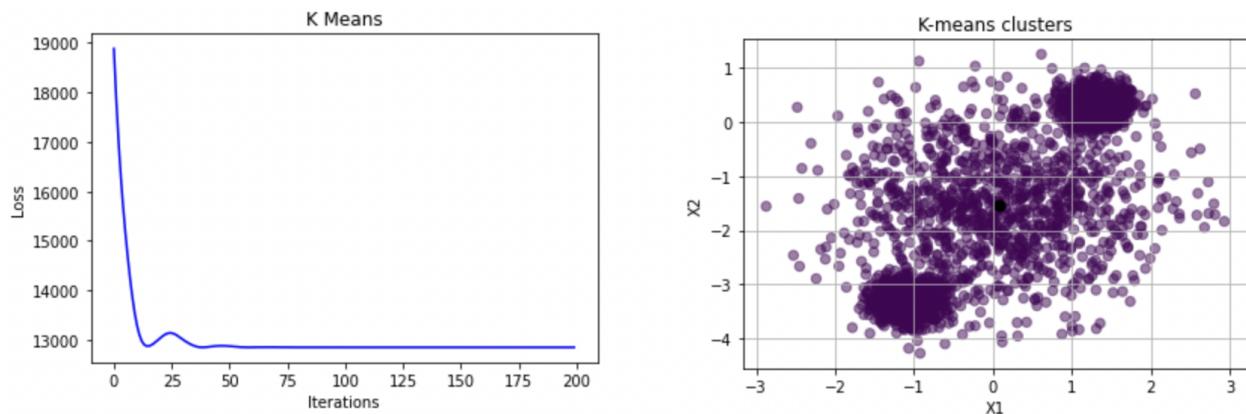


Figure 4: Clusters and Loss for K=1

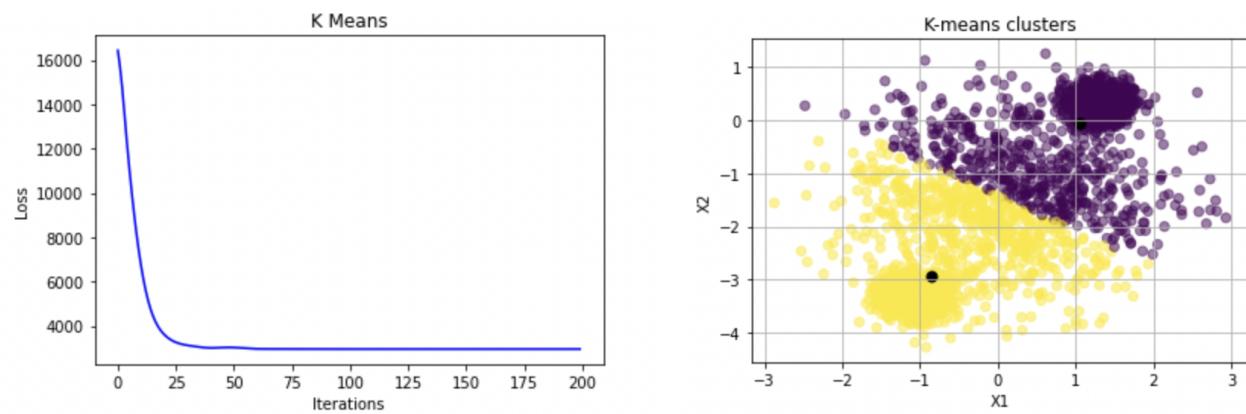


Figure 5: Clusters and Loss for K=2

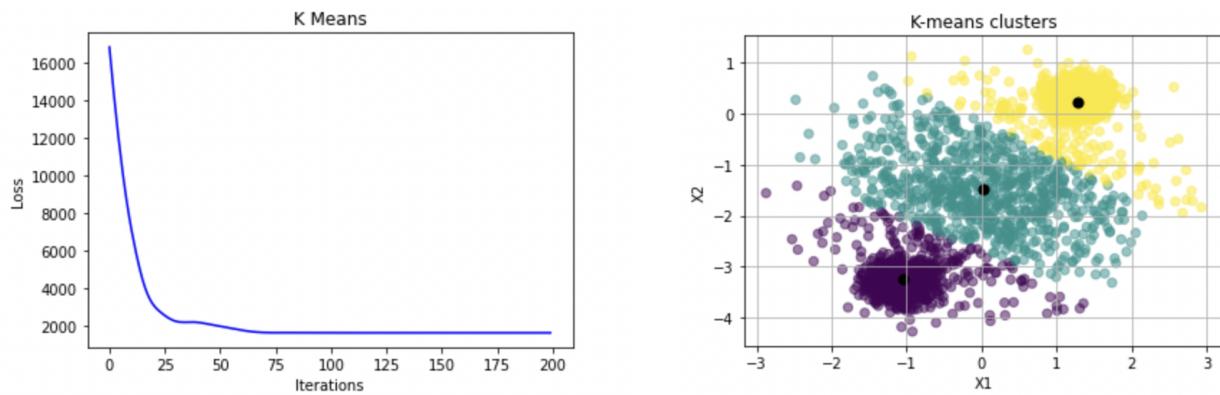


Figure 6: Clusters and Loss for K=3

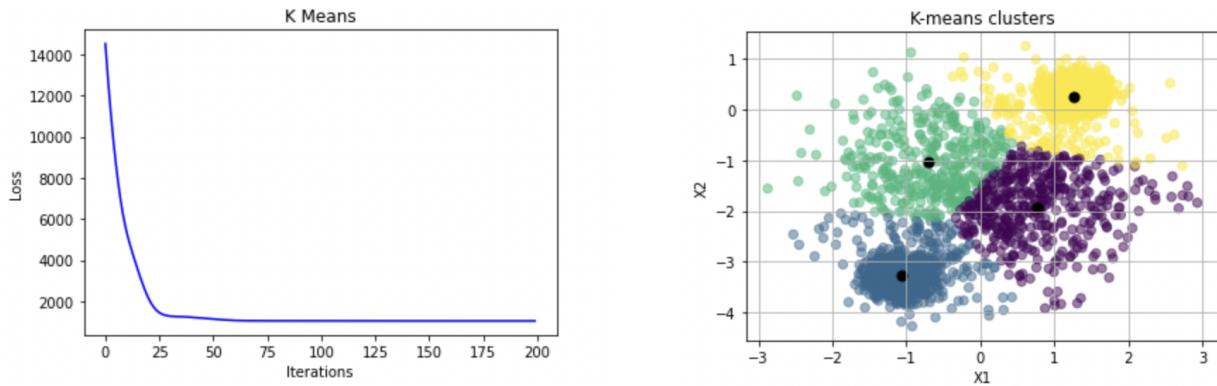


Figure 7: Clusters and Loss for K=4

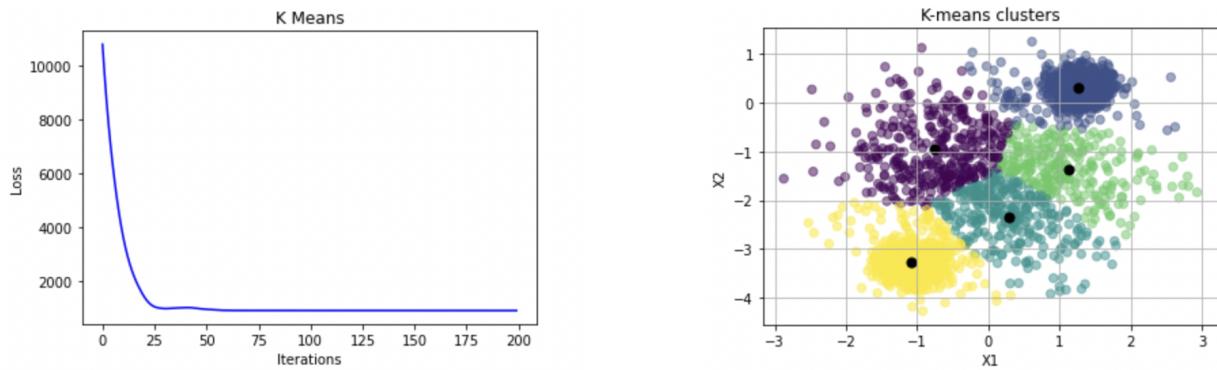


Figure 8: Clusters and Loss for K=5

Table 1: Shows the percentage of points belonging to a cluster with different K's

K	% Class 1	% Class 2	% Class 3	% Class 4	% Class 5
1	100				
2	51.82	48.18			
3	37.50	39.54	22.95		
4	13.14	38.82	11.40	36.63	
5	10.62	7.71	38.34	6.96	36.36

Table 2: Shows the final loss values depending on number of clusters

K (# clusters)	Validation Loss
1	12853.26
2	2958.65
3	1608.84
4	1052.84
5	884.77

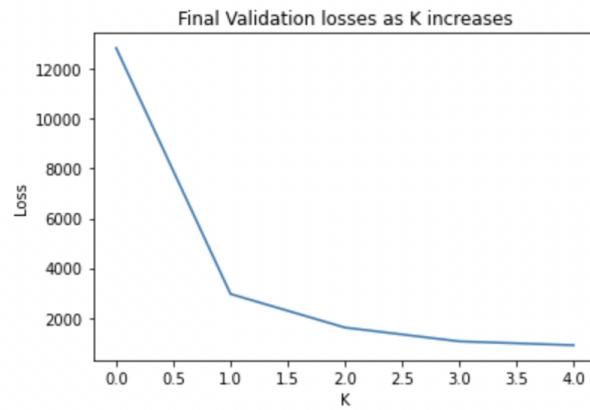


Figure 9: Validation losses as K changes

By observing the percentages, we can see that for  $K=3$ , the distribution of points in the 3 clusters is more even compared to  $K=4, 5$  where many data points are concentrated in one cluster. We can also see the loss starts to sort of level off at  $K=3$  so the best choice is  $K=3$ .

## 2.1 The Gaussian Cluster Mode

1. Derived  $\log P(x|\mu_k, \sigma_k^2)$  using the pdf equation for Gaussians. Derivation of all equations is attached in Appendix at the end of the report.

Log PDF :

$$\begin{aligned}
 & N(x|\mu_k, \sigma_k^2) \\
 &= P(x|\mu_k, \sigma_k^2) \\
 &= \log P(x|\mu_k, \sigma_k^2) \rightarrow \text{bcz given as logP} \\
 &= \log \left( \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(\frac{-(x-\mu_k)^2}{2\sigma_k^2}\right) \right) \rightarrow \text{stated in Tutorial} \\
 &= \log\left(\frac{1}{\sqrt{2\pi}\sigma_k}\right) + \left(\frac{-(x-\mu_k)^2}{2\sigma_k^2}\right)
 \end{aligned}$$

Figure 10: Derivation of  $\log P(x|\mu_k, \sigma_k^2)$

```

def log_gauss_pdf(X, mu, sigma):
    """ Inputs:
        X: N X D
        mu: K X D
        sigma: K X 1

        Outputs:
            log Gaussian PDF (N X K)
    """
    x_mu = distance_func(X, mu)
    x_mu = tf.multiply(x_mu, x_mu)
    sigma = tf.squeeze(sigma)
    second_term = x_mu/(2*sigma)
    first_term = tf.log(2*np.pi*sigma) * tf.to_float(tf.rank(X))
    pdf = -0.5*first_term - second_term
    return pdf

```

Figure 11: Code for `log_gauss_pdf()`

2. Derived  $\log P(z|x)$  and wrote `log_posterior()`. It is important to use `reduce_logsumexp()` because it is more stable than using `reduce_sum()`. Taking  $\exp$  of large values can cause overflows and  $\exp$  of small numbers can cause underflows and this is handled well by the function.

$P(z=k|x)$  : (tutorial)

$$\begin{aligned}
 P(z=k|x) &= \frac{P(x|z)P(z)}{P(x)} \rightarrow \text{Bayes rule} \\
 &= \frac{\prod_i P(x|z_i)P(z_i)}{\sum_i P(x|z_i)P(z_i)}
 \end{aligned}$$

$$\begin{aligned}
 P(z_i|x) &= \log P(z_i|x) \xrightarrow{\text{from ln a logPi}} \\
 &= \log P(x_i|z_i) + \log P(z_i) - \log \underbrace{\sum_j e^{\log P(x_j|z_j) + \log P(z_j)}}_{LSE}
 \end{aligned}$$

Figure 12: Derivation of  $\log P(z|x)$

```

def log_posterior(log_PDF, log_pi):
    """ Inputs:
        log_PDF: log Gaussian PDF N X K
        log_pi: K X 1

        Outputs
        log_post: N X K
    """
    logPI = tf.squeeze(log_pi)
    logP = tf.add(logPI, log_PDF)
    logSummation = hlp.reduce_logsumexp(log_PDF + logPI, keep_dims=True)
    final = logP - logSummation
    return final

```

Figure 13: Code for log\_posterior()

## 2.2 Learning the MoG

1. *The model parameters that the algorithm has learnt are pi, mu, sigma which maximize  $P(x_i)$  and minimize the log Loss.*

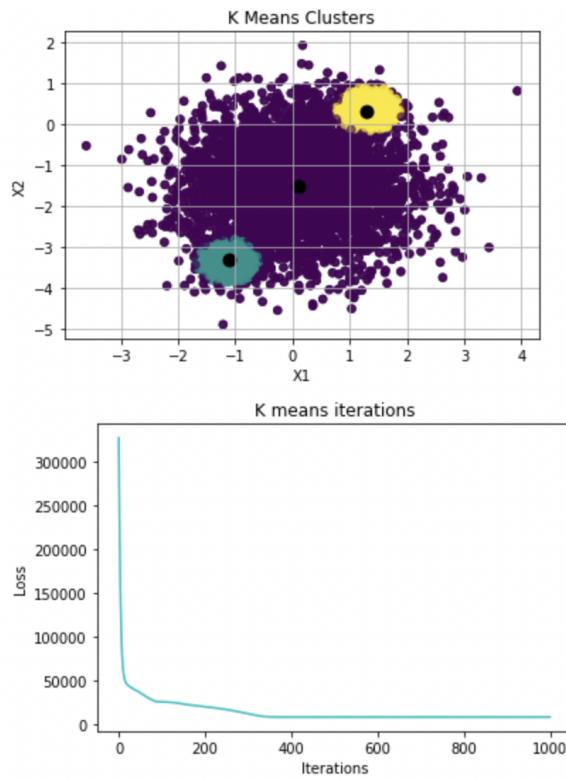


Figure 14: Clusters and Training loss for K=3

2. Algorithm was run for  $K=\{1,2,3,4, 5\}$ . The most reasonable choice for  $K$  is  $k=3$ . There was very little change in loss for  $K=4,5$  and also for  $K=3$  the data was more evenly distributed.

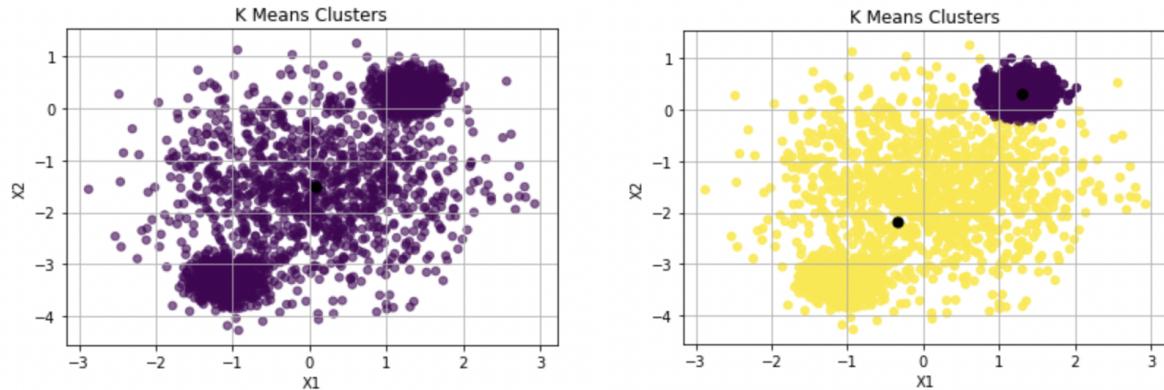


Figure 15: Clusters for  $K=1, 2$

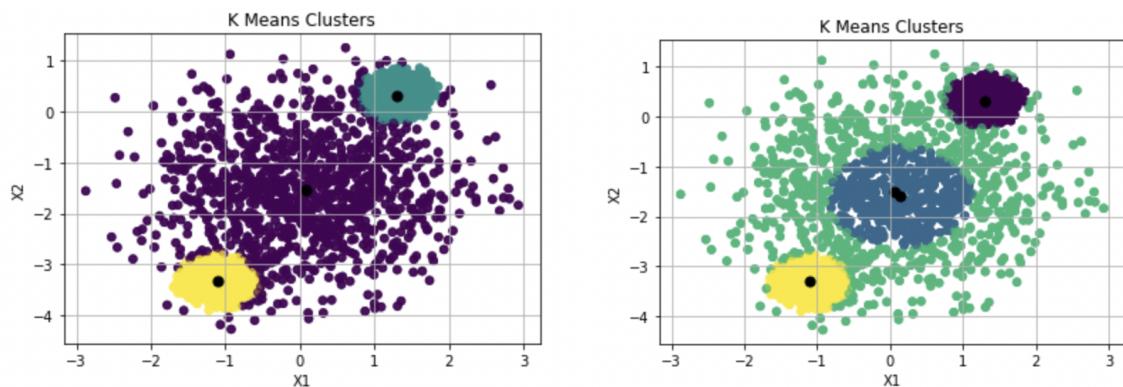


Figure 16: Clusters for  $K=3, 4$

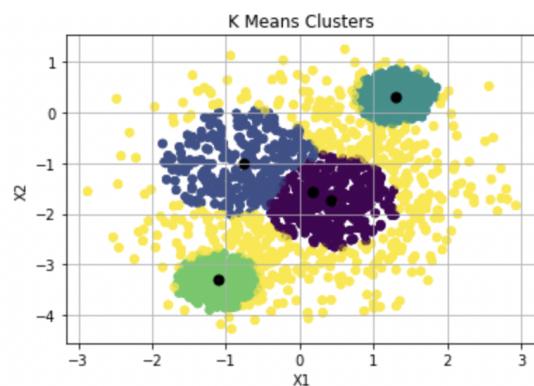


Figure 17: Clusters for  $K=5$

Table 3: Shows the percentage of points belonging to a cluster with different K's

K	% Class 1	% Class 2	% Class 3	% Class 4	% Class 5
1	100				
2	33.60	66.40			
3	32.94	32.79	34.26		
4	32.82	12.66	20.10	34.41	
5	9.84	8.07	32.97	34.47	14.64

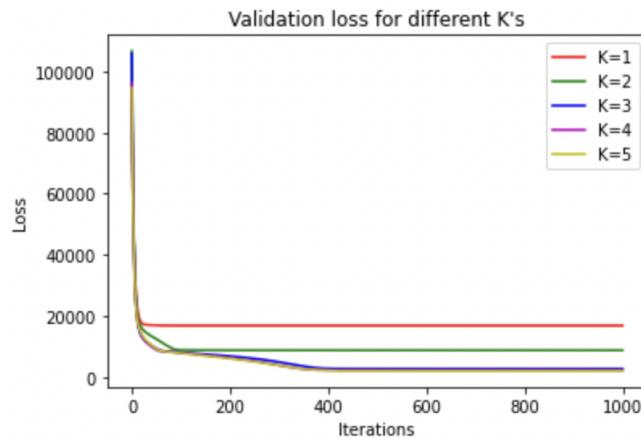


Figure 18: Validation Loss for different K values

Table 4: Shows the final loss values depending on number of clusters

K (# clusters)	Validation Loss
1	16815.32
2	8732.02
3	2678.17
4	2257.59
5	2009.33

3. Many clusters past  $K=10$ , had 0% of points in some clusters for  $K=15,20,30$ .  
 For the K-means, the good choice is around 10 clusters because there is very little change in loss for 15,20,30 and half the clusters are empty in those.  
 However, we get lower overall losses using GMM in higher dimensions as we can see in Table 5. For GMM, around 15 clusters would be a good choice but the loss seemed to be fluctuating quite a bit, however it did have much better performance than K-means in higher dimensional data (data100D.py).

Table 5: Shows the final loss values depending on number of clusters

<b>K</b>	<b>5</b>	<b>10</b>	<b>15</b>	<b>20</b>	<b>30</b>
<b>K-means</b>	71653.69	69918.68	68553.46	67972.11	67119.99
<b>GMM</b>	36884.13	37925.97	30505.32	34864.34	30573.30

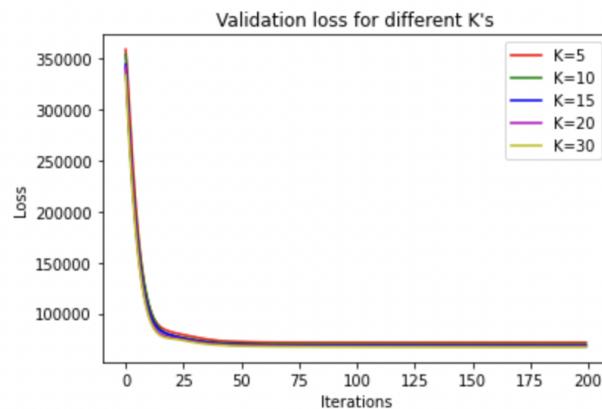


Figure 18: K-Means

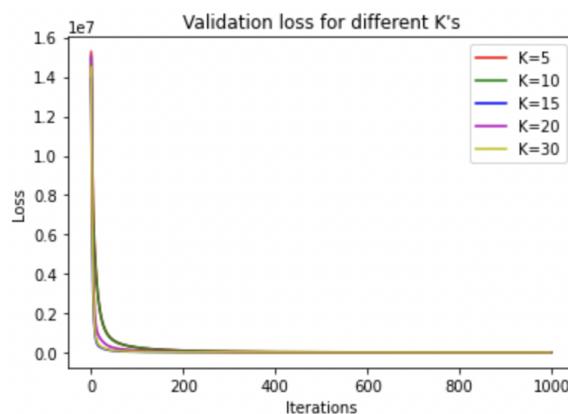


Figure 19: GMM

### 3.0 Appendix

Log PDF :

$$\begin{aligned}
 & N(x|N_K, \sigma_K^2) \\
 & = P(x|N_K, \sigma_K^2) \\
 & = \log P(x|N_K, \sigma_K^2) \rightarrow \text{bcz given as log p} \\
 & = \log \left( \frac{1}{\sqrt{2\pi}\sigma_K} e^{-\frac{(x-N_K)^2}{2\sigma_K^2}} \right) \rightarrow \text{stated in Tutorial} \\
 & = \log \left( \frac{1}{\sqrt{2\pi}\sigma_K} \right) + \left( -\frac{(x-N_K)^2}{2\sigma_K^2} \right)
 \end{aligned}$$

$P(z=k|x)$  : (tutorial)

$$\begin{aligned}
 P(z=k|x) &= \frac{P(x|z) P(z)}{P(x)} \rightarrow \text{Bayes rule} \\
 &= \frac{P(x|z) P(z)}{\sum_{i=1}^k P(x|z_i) P(z_i)}
 \end{aligned}$$

$$\begin{aligned}
 P(z_i|x) &= \log P(z_i|x) \xrightarrow{\text{from in A}} \log \pi_i \\
 &= \log P(x_i|z_i) + \log P(z_i) - \underbrace{\log \sum_j e^{\log P(x_j|z_j) + \log P(z_j)}}_{LSE}
 \end{aligned}$$

$$\left. \begin{array}{l} K=3 \\ \Pi_1 \left\{ \begin{array}{l} N_1 \\ \sigma_1 \end{array} \right. \\ \Pi_2 \left\{ \begin{array}{l} N_2 \\ \sigma_2 \end{array} \right. \\ \Pi_3 \left\{ \begin{array}{l} N_3 \\ \sigma_3 \end{array} \right. \end{array} \right\} P(x) = P(x|z_1) P(z_1) + P(x|z_2) P(z_2) + P(x|z_3) P(z_3) = \sum_{i=1}^3 P(x|z_i)$$

prob of x given we selected 1st gaussian  
 x prob we select 1st gaussian  
 $\rightarrow x$  either generated from  $z_1, z_2, z_3, \dots$

$\downarrow$

$$\Pi_1 + \Pi_2 + \Pi_3 = 1$$

Want to maximize this  $P(\underline{x})$

$$\theta = \{\pi_k, \mu_k, \sigma_k^2\}_{k=1}^K \text{ (parameters)}$$

↪ want setting that maximizes  $P$ .

$$\hookrightarrow \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^N p_{\theta}(x_i)$$

↑ training - ML

↪  $\underline{x}$  given an  $\underline{x}$ , want to know prob of it generated from 1, 2, 3

$$P(z|x) = \begin{bmatrix} P(z=1|x) \\ P(z=2|x) \\ P(z=3|x) \end{bmatrix}$$

$$P(z_i|x) = \frac{P(x|z_i)P(z_i)}{\sum_{j=1}^k P(x|z_j)P(z_j)} \quad \text{Bayes}$$

$$\text{g: } P(z=0|x) = \frac{P(x|z=0)P(z=0)}{P(x|z=0)P(z_0) + P(x|z=1)P(z_1)}.$$

↑ inference

$$\underset{\theta}{\operatorname{argmax}} \prod_{i=1}^N p_{\theta}(x_i)$$

$$\begin{aligned} \underset{\theta}{\operatorname{argmax}} & \sum_i \log p_{\theta}(x_i) && \text{we have log of this not only } \\ & = -\sum_{i=1}^N \log \left( \sum_{k=1}^K P(x_i|z=k) \underbrace{P(z=k)}_{\substack{\text{pdf of} \\ \text{gaussian}}} \right) && \pi^k \end{aligned}$$

$$\begin{aligned}
 &= -\sum_{i=1}^n \log \sum_{k=1}^K e^{\log P(x_i | z=k) + \log \pi(z=k)} \\
 &= -\sum_{i=1}^n \log \sum_{k=1}^K e^{\log P(x_i | z=k) + \log \pi_k} \quad \textcircled{1} \quad \textcircled{II} \\
 &\quad \text{L S E} \quad \text{sum[LSE(1, II)]} \\
 &\quad \text{by sum_exp.} \quad \log \sum e^0 \\
 &\hookrightarrow \text{similar to softmax.} \\
 &\hookrightarrow \text{in helper.py.}
 \end{aligned}$$

$$\begin{aligned}
 P(z_i | x) &= \log P(z_i | x) \xrightarrow{\text{from ln a}} \log \pi_i \\
 &= \log P(x_i | z_i) + \log \pi(z_i) + \underbrace{\log \sum e^{\log P(x_i | z_j) + \log \pi_j}}_{\text{LSE}}
 \end{aligned}$$

$$\theta = [\pi_k, \mu_k, \sigma_k^2]_{k=1}^K \quad (\text{parameters})$$

$$\begin{aligned}
 \psi_i &\rightarrow \text{use log parameter.} & \rightarrow \text{unconstraint } \pi \\
 \pi &= \underbrace{\text{Softmax}(\psi)}_{\text{use this in grad descent.}}
 \end{aligned}$$