

Fall 2020 CX4641/CS7641 A Homework 2

Instructor: Dr. Mahdi Roozbahani

Deadline: Oct 6th, Tuesday, 11:59 pm AOE

- No discussion/extension of the deadline is allowed. Late submission will lead to 0 credit.
- Discussions are encouraged on Piazza as part of the Q/A. However, all assignments should be done individually.

Instructions for the assignment

- This assignment consists of both programming and theory questions.
- Q4 is bonus for both undergraduate and graduate students.
- To switch between cells for code and for markdown, see the menu -> Cell -> Cell Type
- You can directly type LaTeX equations into markdown cells.
- Typing with LaTeX/markdown is required for all the written questions. Handwritten answers will not be accepted.
- If a question requires a picture, you could use this syntax `<img src="" style="width: 300px; height: 300px;"/>` to include them within your jupyter notebook.

Using the autograder

- You will find two assignments on Gradescope that correspond to HW2: 'HW2 - Programming' and 'HW2 - Non-programming'.
- You will submit your code for the autograder on 'HW2 - Programming' in the following format:
  - kmeans.py
  - gmm.py
  - semisupervised.py
- All you will have to do is to copy your implementations of the classes "Kmeans", "GMM", "CleanData", "SemiSupervised" onto the respective files. We provided you different .py files and we added libraries in those files please DO NOT remove those lines and add your code after those lines. Note that these are the only allowed libraries that you can use for the homework.
- You are allowed to make as many submissions until the deadline as you like. Additionally, note that the autograder tests each function separately, therefore it can serve as a useful tool to help you debug your code if you are not sure of what part of your implementation might have an issue.
- For the 'HW2 - Non-programming' part, you will download your jupyter notebook as html and submit it as a PDF on Gradescope. To download the notebook as PDF click on "File" on the top left corner of this page and select "Download as > PDF". The non-programming part corresponds to Q2, Q3.3 (both your response and the generated images with your implementation) and Q4.2
- When submitting to Gradescope, please make sure to mark the page(s) corresponding to each problem/sub-problem.

0 Set up

This notebook is tested under [python 3.1.1](#), and the corresponding packages can be downloaded from [miniconda](#). You may also want to get yourself familiar with several packages:

- [jupyter notebook](#)
- [numpy](#)
- [matplotlib](#)

Please implement the functions that have "raise NotImplementedError", and after you finish the coding, please delete or comment "raise NotImplementedError".

```
In [1]: #####
# Do NOT change this cell. ##
#####

from __future__ import absolute_import
from __future__ import print_function
from __future__ import division

import matplotlib inline

import sys
import matplotlib.pyplot as plt
import numpy as np
import matplotlib.pyplot as plt
from matplotlib.pyplot import *
from mpl_toolkits.mplot3d import Axes3D
from tqdm import tqdm

print("Version information")
print("python: {}".format(sys.version))
print("matplotlib: {}".format(matplotlib.__version__))
print("tqdm: {}".format(tqdm.__version__))

# Set random seed so output is all same
np.random.seed(1)

# Load image
import imageio

Version information
python: 3.7.3 (default, Mar 27 2019, 16:54:48)
[Clang 4.0.1 (tags/RELEASE_401/final)]
matplotlib: 3.1.0
numpy: 1.16.4
```

1. KMeans Clustering [5 + 30 + 10 + 5 + 10 pts]

KMeans is trying to solve the following optimization problem:

$$\arg \min_S \sum_{j=1}^K \sum_{x_i \in S_j} \|x_i - \mu_j\|^2$$

where one needs to partition the N observations into K clusters:  $S = \{S_1, S_2, \dots, S_K\}$  and each cluster has  $\mu_j$  as its center.

1.1 pairwise distance [5pts]

In this section, you are asked to implement pairwise\_dist function.

Given  $X \in \mathbb{R}^{N \times D}$  and  $Y \in \mathbb{R}^{M \times D}$ , obtain the pairwise distance matrix  $dist \in \mathbb{R}^{N \times M}$  using the euclidean distance metric, where  $dist_{i,j} = \|X_i - Y_j\|_2$ .

DO NOT USE FOR LOOP in your implementation -- they are slow and will make your code too slow to pass our grader. Use array broadcasting instead.

1.2 KMeans Implementation [30pts]

In this section, you are asked to implement \_init\_centers [5pts], \_update\_assignment [10pts], \_update\_centers [10pts] and \_get\_loss function [5pts].

For the function signature, please see the corresponding doc strings.

1.3 Find the optimal number of clusters [10 pts]

In this section, you are asked to implement find\_optimal\_num\_clusters function.

You will now use the elbow method to find the optimal number of clusters.

1.4 Autograder test to find centers for data points [5 pts]

To obtain these 5 points, you need to be pass the tests set up in the autograder. These will test the centers created by your implementation. Be sure to upload the correct files to obtain these points.

```
In [2]: class KMeans(object):
    def __init__(self): #No need to implement pass
    def pairwise_dist(self, x, y): # [5 pts]
        Args:
            x: N x D numpy array
            y: M x D numpy array
        Return:
            dist: N x M array, where dist[i, j] is the euclidean distance between x[i, : ] and y[j, : ]
            #12 norm, ord=2
            dist = np.linalg.norm(x[:, np.newaxis, :] - y, ord=2, axis=2)
        return dist
        raise NotImplementedError
    def _init_centers(self, points, K, **kwargs): # [5 pts]
        Args:
            points: NxD numpy array, where N is # points and D is the dimensionality
            K: number of clusters
            kwargs: any additional arguments you want
        Return:
            centers: K x D numpy array, the centers.
        """
        #K random indices from points
        indices = np.random.choice(points.shape[0], size=K, replace=False)
        #from the numpy array with these indices
        centers = points[indices, :]
        return centers
        raise NotImplementedError
    def _update_assignment(self, centers, points): # [10 pts]
        Args:
            centers: KxD numpy array, where K is the number of clusters, and D is the dimension
            points: NxD numpy array, the observations
        Return:
            cluster_idx: numpy array of length N, the cluster assignment for each point
        Hint: You could call pairwise_dist() function.
        """
        #distances between points and all cluster centers
        distances = self.pairwise_dist(points, centers)
        #index of minimum distance from each row
        cluster_idx = np.argmin(distances, axis=1)
        return cluster_idx
        raise NotImplementedError
    def _update_centers(self, old_centers, cluster_idx, points): # [10 pts]
        Args:
            old_centers: old centers KxD numpy array, where K is the number of clusters, and D is the dimension
            cluster_idx: numpy array of length N, the cluster assignment for each point
            points: NxD numpy array, the observations
        Return:
            centers: new centers, K x D numpy array, where K is the number of clusters, and D is the dimension.
        """
        K,D = old_centers.shape[0], old_centers.shape[1]
        #initialize centers as zero array
        centers = np.zeros((K,D))
        for i in range(K):
            #find mean of all points having i as cluster id
            centers[i] = np.mean(points[cluster_idx==i, :], axis=0)
        return centers
        raise NotImplementedError
    def _get_loss(self, centers, cluster_idx, points): # [5 pts]
        Args:
            centers: KxD numpy array, where K is the number of clusters, and D is the dimension
            cluster_idx: numpy array of length N, the cluster assignment for each point
            points: NxD numpy array, the observations
        Return:
            loss: a single float number, which is the objective function of KMeans.
        """
        #find squared distances between all points and cluster centers
        distances = np.linalg.norm(points[:, np.newaxis, :] - centers, ord=2, axis=2) ** 2
        #element distance from cluster center
        distance_from_cluster_center = distances[np.arange(len(distances), cluster_idx)]
        #loss is sum of all these distances
        loss = np.sum(distance_from_cluster_center)
        return loss
        raise NotImplementedError
    def _call(self, points, K, max_iters=100, abs_tol=1e-16, rel_tol=1e-16, verbose=False, **kwargs):
        Args:
            points: NxD numpy array, where N is # points and D is the dimensionality
            K: number of clusters
            max_iters: maximum number of iterations (Hint: You could change it when debugging)
            abs_tol: convergence criteria w.r.t. absolute change of loss
            rel_tol: convergence criteria w.r.t. relative change of loss
            verbose: boolean to set whether method should print loss (Hint: helpful for debugging)
            kwargs: any additional arguments you want
        Return:
            cluster assignments: Nx1 int numpy array
            cluster centers: K x D numpy array, the centers
        """
        loss = self._init_centers(points, K, **kwargs)
        for i in range(max_iters):
            cluster_idx = self._update_assignment(centers, points)
            centers = self._update_centers(centers, cluster_idx, points)
            loss = self._get_loss(centers, cluster_idx, points)
            K = centers.shape[0]
            if i:
                diff = np.abs(prev_loss - loss)
                if diff < abs_tol and diff / prev_loss < rel_tol:
                    break
            prev_loss = loss
            if verbose:
                print('iter %d, loss: %.4f' % (i, loss))
        return cluster_idx, centers, loss
    def find_optimal_num_clusters(self, data, max_K=15): # [10 pts]
        """Plots loss values for different number of clusters in K-Means"""
        Args:
            image: input image of shape (H, W, 3)
            max_K: the number of clusters
        Return:
            losses: an array of loss denoting the loss of each number of clusters
        """
        losses = []
        #for each value of K, find the loss and append to array
        for i in range(1, max_K+1):
            cluster_idx, centers, loss = self._call(data, i)
            losses.append(loss)
        return losses
        raise NotImplementedError
    # Helper function for checking the implementation of pairwise_distance function. Please DO NOT change this function
    # TEST CASE
    x = np.random.randn(2, 2)
    y = np.random.randn(3, 2)

    print("Expected Answer ***")
    print("=====")
    print([1.62434536 -0.61175641]
          [-0.52817175 -1.07296862])
    print("=====")
    print([0.86540763 -2.3015387 ]
          [1.74481176 -0.7612069 ]
          [0.3190391 -0.24937038])
    print("=====")
    print([1.85239052 0.19195729 1.35467638]
          [1.85780729 2.29426447 1.18155842])
    print("My Answer ****")
    print("=====")
    print([1.62434536 -0.61175641]
          [-0.52817175 -1.07296862])
    print("=====")
    print([0.86540763 -2.3015387 ]
          [1.74481176 -0.7612069 ]
          [0.3190391 -0.24937038])
    print("=====")
    print([1.85239052 0.19195729 1.35467638]
          [1.85780729 2.29426447 1.18155842])

In [3]: # Helper function for checking the implementation of pairwise_distance function. Please DO NOT change this function
# TEST CASE
x = np.random.randn(2, 2)
y = np.random.randn(3, 2)

print("Expected Answer ***")
print("=====")
print([1.62434536 -0.61175641]
      [-0.52817175 -1.07296862])
print("=====")
print([0.86540763 -2.3015387 ]
      [1.74481176 -0.7612069 ]
      [0.3190391 -0.24937038])
print("=====")
print([1.85239052 0.19195729 1.35467638]
      [1.85780729 2.29426447 1.18155842])
print("My Answer ****")
print("=====")
print([1.62434536 -0.61175641]
      [-0.52817175 -1.07296862])
print("=====")
print([0.86540763 -2.3015387 ]
      [1.74481176 -0.7612069 ]
      [0.3190391 -0.24937038])
print("=====")
print([1.85239052 0.19195729 1.35467638]
      [1.85780729 2.29426447 1.18155842])

In [4]: def image_to_matrix(image_file, grays=False):
    """
    Convert .png image to matrix
    of values.
    params:
    image_file = str
    grays = Boolean
    returns:
    img = (color) np.ndarray[np.ndarray[np.ndarray[...]]]
    or (grayscale) np.ndarray[np.ndarray[np.ndarray[...]]]
    """
    img = plt.imread(image_file)
    # is case of transparency values
    if len(img.shape) == 3 and img.shape[2] > 3:
        height, width, depth = img.shape
        new_img = np.zeros((height, width, 3))
        for r in range(height):
            for c in range(width):
                new_img[r, c, :] = img[r, c, 0:3]
    img = np.copy(new_img)
    if grays and len(img.shape) == 3:
        height, width = img.shape[0:2]
        new_img = np.zeros((height, width))
        for r in range(height):
            for c in range(width):
                new_img[r, c] = img[r, c, 0]
    img = new_img
    return img

In [5]: image_values = image_to_matrix('./images/bird_color_24.png')

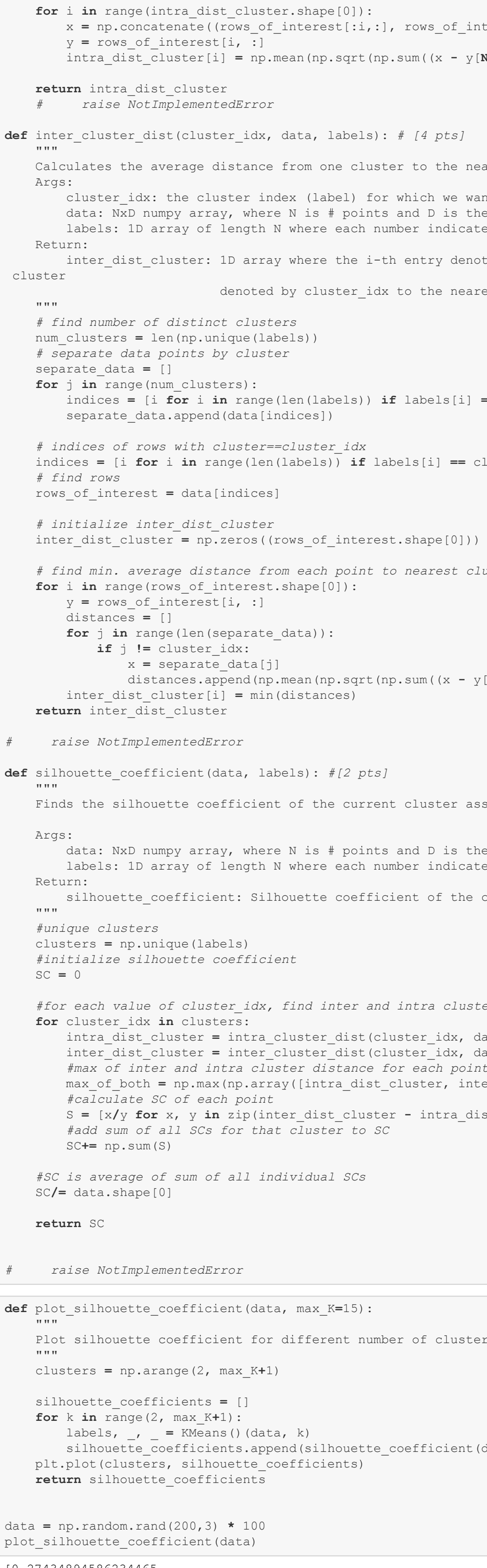
r = image_values.shape[0]
c = image_values.shape[1]
ch = image_values.shape[2]
# Flatten the image values
image_values = image_values.reshape(r*c,ch)

k = 6 # feel free to change this value
cluster_idx, centers, loss = KMeans().fit(image_values, k)
updated_image_values = np.copy(image_values)

# assign each pixel to cluster mean
for i in range(0,k):
    indices_current_cluster = np.where(cluster_idx == i)[0]
    updated_image_values[indices_current_cluster] = centers[i]

updated_image_values = updated_image_values.reshape(r,c,ch)

plt.figure(figsize=(9,12))
plt.imshow(updated_image_values)
plt.show()
```



```
In [6]: KMeans().find_optimal_num_clusters(image_values)
Out [6]: [25575.970929107656,
14136.693932942087,
10225.6098515669,
6993.510281152104,
2738.2662257049277,
2554.129469081243,
2065.325028335666,
1929.4098314821665,
1965.2534688501296,
1706.0740171535149,
1580.6903867544215]
```

Silhouette Coefficient Evaluation [10 pts]

The average silhouette of the data is another useful criterion for assessing the natural number of clusters. The silhouette of a data instance is a measure of how closely it is matched to data within its cluster and how loosely it is matched to data of the neighbouring cluster.

The silhouette value is a measure of how similar an object is to its own cluster (cohesion) compared to other clusters (separation). The silhouette ranges from -1 to 1, where a high value indicates that the object is well matched to its own cluster and poorly matched to neighboring clusters. If most objects have a high value, then the clustering configuration is appropriate. If many points have a low or negative value, then the clustering configuration may have too many or too few clusters.

```
In [7]: def intra_cluster_dist(cluster_idx, data, labels): # [4 pts]
    """
    Calculates the average distance from a point to other points within the same cluster
    Args:
        cluster_idx: the cluster index (label) for which we want to find the intra cluster distance
        data: NxD numpy array, where N is # points and D is the dimensionality
        labels: 1D array of length N where each number indicates of cluster assignment for that point
    Return:
        intra_dist_cluster: 1D array where the i-th entry denotes the average distance from point i
        in cluster denoted by cluster_idx to other points within the same cluster
    """
    # indices of rows with cluster==cluster_idx
    indices = [i for i in range(len(labels)) if labels[i]==cluster_idx]
    # find rows
    rows_of_interest = data[indices]
    # initialize intra_dist_cluster
    intra_dist_cluster = np.zeros((rows_of_interest.shape[0]))

    for i in range(intra_dist_cluster.shape[0]):
        x = np.concatenate((rows_of_interest[i,:], rows_of_interest[i+1,:]))
        y = rows_of_interest[i, :]
        intra_dist_cluster[i] = np.mean(np.sqrt(np.sum((x - y)**2, axis=1)))

    return intra_dist_cluster
    raise NotImplementedError

def inter_cluster_dist(cluster_idx, data, labels): # [4 pts]
    """
    Calculates the average distance from a point to the nearest cluster
    Args:
        cluster_idx: the cluster index (label) for which we want to find the intra cluster distance
        data: NxD numpy array, where N is # points and D is the dimensionality
        labels: 1D array of length N where each number indicates of cluster assignment for that point
    Return:
        inter_dist_cluster: 1D array where the i-th entry denotes the average distance from point i in
        cluster denoted by cluster_idx to the nearest neighboring cluster
    """
    # find number of distinct clusters
    num_clusters = len(np.unique(labels))
    # separate data points by cluster
    separate_data = {}
    for j in range(num_clusters):
        indices = [i for i in range(len(labels)) if labels[i] == j]
        separate_data.append(data[indices])

    # indices of rows with cluster==cluster_idx
    indices = [i for i in range(len(labels)) if labels[i] == cluster_idx]
    # find rows
    rows_of_interest = data[indices]
    # initialize inter_dist_cluster
    inter_dist_cluster = np.zeros((rows_of_interest.shape[0]))

    # find min. average distance from each point to nearest cluster
    for i in range(rows_of_interest.shape[0]):
        y = rows_of_interest[i, :]
        distances = []
        for j in range(len(separate_data)):
            if j != cluster_idx:
                x = separate_data[j]
                distances.append(np.mean(np.sqrt(np.sum((x - y)**2, axis=1)))
            inter_dist_cluster[i] = min(distances)

    return inter_dist_cluster
    raise NotImplementedError

def silhouette_coefficient(data, labels): # [2 pts]
    """
    Finds the silhouette coefficient of the current cluster assignment
    Args:
        data: NxD numpy array, where N is # points and D is the dimensionality
        labels: 1D array of length N where each number indicates of cluster assignment for that point
    Return:
        silhouette_coefficient: silhouette coefficient of the current cluster assignment
    """
    #unique clusters
    clusters = np.unique(labels)
    #initialize silhouette coefficient
    SC = 0

    for cluster_idx in clusters:
        intra_dist_cluster = intra_cluster_dist(cluster_idx, data, labels)
        inter_dist_cluster = inter_cluster_dist(cluster_idx, data, labels)
        #max of intra and inter cluster distance for each point
        max_of_both = np.max(np.array([intra_dist_cluster, inter_dist_cluster]), axis=0)
        #calculate SC of each point
        s = [s for s, y in zip(inter_dist_cluster - intra_dist_cluster, max_of_both)]
        #add sum of all SCs for that cluster to SC
        SC += np.sum(s)

    #SC is average of sum of all individual SCs
    SC /= data.shape[0]

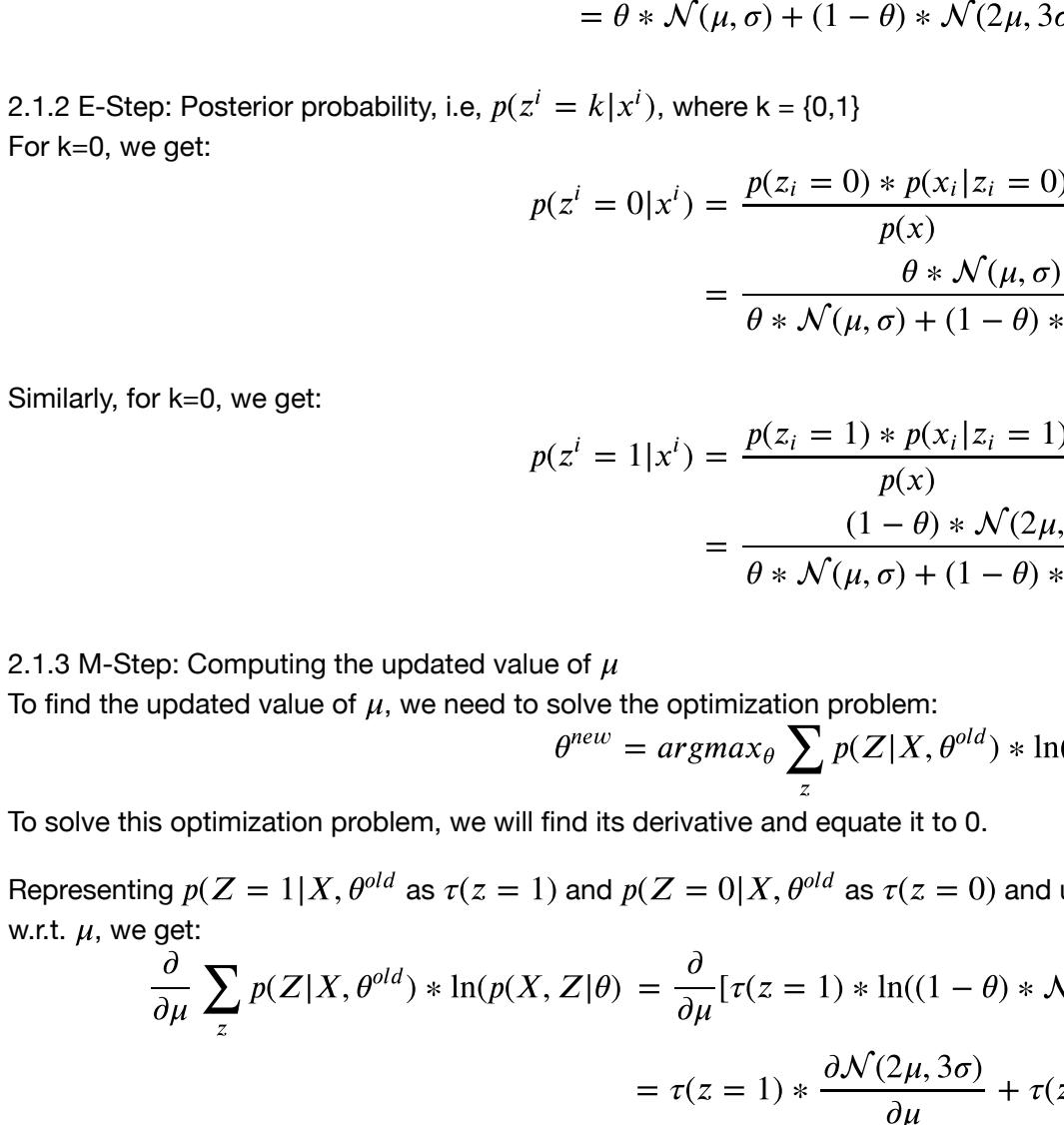
    return SC

# raise NotImplementedError

In [8]: def plot_silhouette_coefficient(data, max_K=15):
    """
    Plot silhouette coefficient for different number of clusters, no need to implement
    """
    clusters = np.arange(2, max_K+1)

    silhouette_coefficients = []
    for k in range(2, max_K+1):
        labels, _ = KMeans().fit(data, k)
        silhouette_coefficients.append(silhouette_coefficient(data, labels))
    plt.plot(clusters, silhouette_coefficients)
    return silhouette_coefficients

data = np.random.rand(200,3) * 100
plot_silhouette_coefficient(data)
```



2. EM algorithm [20 pts]

2.1 Performing EM Update [10 pts]

A univariate Gaussian Mixture Model (GMM) has two components, both of which have their own mean and standard deviation. The model is defined by the following parameters:

$$z \sim \text{Bernoulli}(\theta)$$
$$p(x|z = 0) = \mathcal{N}(\mu, \sigma)$$
$$p(x|z = 1) = \mathcal{N}(\mu, 3\sigma)$$

For a dataset of N datapoints, find the value of  $\theta$ .

2.1.1 Write the marginal probability of  $x$ , i.e.  $p(x)$  [2pts]

2.1.2 E-Step: Compute the posterior probability, i.e.  $p(z' = k|x')$ , where  $k \in \{0, 1\}$  [2pts]

2.1.3 M-Step: Compute the updated value of  $\theta$  (You can keep  $\theta$  fixed for this) [3pts]

2.1.4 M-Step: Compute the updated value for  $\sigma$  (You can keep  $\mu$  fixed for this) [3pts]

2.1.1 Marginal probability of  $x$ , i.e.  $p(x)$

Given: A univariate Gaussian mixture model with two components, defined by their parameters as:

$$z \sim \text{Bernoulli}(\theta)$$
$$p(x|z = 0) = \mathcal{N}(\mu, \sigma)$$
$$p(x|z = 1) = \mathcal{N}(\mu, 3\sigma)$$

Note: Considering  $p(z = 0) = \theta$ , and  $\sigma$  and  $3\sigma$  as standard deviation as suggested by TAs on piazza:

Marginal probability of  $x$  can be calculated as:

$$p(x) = \sum_{z \in \{0, 1\}} p(z) * p(x|z)$$
$$= p(z = 0) * p(x|z = 0) + p(z = 1) * p(x|z = 1)$$
$$= \theta * \mathcal{N}(\mu, \sigma) + (1 - \theta) * \mathcal{N}(\mu, 3\sigma)$$

2.1.2 E-Step: Posterior probability, i.e.  $p(z' = k|x')$ , where  $k \in \{0, 1\}$

For  $k=0$ , we get:

$$p(z' = 0|x') = \frac{p(z = 0) * p(x|z = 0)}{p(x)}$$
$$= \frac{\theta * \mathcal{N}(\mu, \sigma)}{\theta * \mathcal{N}(\mu, \sigma) + (1 - \theta) * \mathcal{N}(\mu, 3\sigma)}$$

Similarly, for  $k=0$ , we get:

$$p(z' = 1|x') = \frac{p(z = 1) * p(x|z = 1)}{p(x)}$$
$$= \frac{(1 - \theta) * \mathcal{N}(\mu, 3\sigma)}{\theta * \mathcal{N}(\mu, \sigma) + (1 - \theta) * \mathcal{N}(\mu, 3\sigma)}$$

2.1.3 M-Step: Computing the updated value of  $\mu$

To find the updated value of  $\mu$ , we need to solve the optimization problem:

$$\theta^{\text{EM}} = \arg \max_{\theta} \sum_{i=1}^N p(X_i, \theta^{\text{EM}}) \approx \ln p(X, \theta^{\text{EM}}) \approx \ln p(X, Z, \theta^{\text{EM}})$$

To solve this optimization problem, we will find its derivative and equate it to 0.

Representing  $\theta$  as  $\theta = P(A) = P(A, z=0) + P(A, z=1)$  and using the value of the derivative of the Gaussian pdf w.r.t.  $\mu$ , we get:

$$\frac{\partial}{\partial \mu} \sum_{i=1}^N p(X_i, \theta^{\text{EM}}) \approx \ln p(X, Z, \theta) = \frac{\partial}{\partial \mu} [\tau(z = 1) * \ln((1 - \theta) * \mathcal{N}(\mu, 3\sigma)) + \tau(z = 0) * \ln(\theta * \mathcal{N}(\mu, \sigma))]$$
$$= \tau(z = 1) * \frac{\partial \mathcal{N}(\mu, 3\sigma)}{\partial \mu} + \tau(z = 0) * \frac{\partial \mathcal{N}(\mu, \sigma)}{\partial \mu}$$
$$= \tau(z = 1) * \frac{2(x - 2\mu) - 9\sigma^2}{9\sigma^3} + \tau(z = 0) * \frac{(x - \mu) - \sigma^2}{\sigma^3}$$

Equating the above value of derivative to 0, we get:

$$\tau(z = 1) * \frac{2(x - 2\mu) - 9\sigma^2}{9\sigma^3} + \tau(z = 0) * \frac{(x - \mu) - \sigma^2}{\sigma^3} = 0$$
$$2x * \tau(z = 1) - 4\mu * \tau(z = 1) + 9x * \tau(z = 0) - 9\mu * \tau(z = 0) = 0$$
$$x(2 * \tau(z = 1) + 9 * \tau(z = 0)) - \mu(2 * \tau(z = 1) + 9 * \tau(z = 0)) = 0$$

This gives us the updated value of  $\mu$  as:

$$\mu = \frac{x(2\tau(z = 1) + 9\tau(z = 0))}{2\tau(z = 1) + 9\tau(z = 0)}$$

2.1.4 M-Step: Computing the updated value for  $\sigma$ :

For this part, we are solving the same optimization problem, but this time, we will find the partial derivative w.r.t.  $\sigma$ .

Representing  $\theta$  as  $\theta = P(A) = P(A, z=0) + P(A, z=1)$  and using the value of the derivative of the Gaussian pdf w.r.t.  $\sigma$ , we get:

$$\frac{\partial}{\partial \sigma} \sum_{i=1}^N p(X_i, \theta^{\text{EM}}) \approx \ln p(X, Z, \theta) = \frac{\partial}{\partial \sigma} [\tau(z = 1) * \ln((1 - \theta) * \mathcal{N}(\mu, 3\sigma)) + \tau(z = 0) * \ln(\theta * \mathcal{N}(\mu, \sigma))]$$
$$= \tau(z = 1) * \frac{\partial \mathcal{N}(\mu, 3\sigma)}{\partial \sigma} + \tau(z = 0) * \frac{\partial \mathcal{N}(\mu, \sigma)}{\partial \sigma}$$
$$= \tau(z = 1) * \frac{(x - 2\mu) - 9\sigma^2}{9\sigma^3} + \tau(z = 0) * \frac{(x - \mu) - \sigma^2}{\sigma^3}$$

Equating the above value of derivative to 0, we get:

$$\tau(z = 1) * \frac{(x - 2\mu) - 9\sigma^2}{9\sigma^3} + \tau(z = 0) * \frac{(x - \mu) - \sigma^2}{\sigma^3} = 0$$
$$\tau(z = 1) * ((x - 2\mu) - 9\sigma^2) + 9\tau(z = 0) * ((x - \mu) - \sigma^2) = 0$$
$$\tau(z = 1)(x - 2\mu) - 9\tau(z = 1)\sigma^2 + 9\tau(z = 0)(x - \mu) - 9\tau(z = 0)\sigma^2 = 0$$
$$\sigma^2 = \frac{\tau(z = 1)(x - 2\mu) + 9\tau(z = 0)(x - \mu)}{9(\tau(z = 1) + \tau(z = 0))}$$

This gives us the updated value of  $\sigma$  as:

$$\sigma = \frac{1}{3} \sqrt{\frac{\tau(z = 1)(x - 2\mu) + 9\tau(z = 0)(x - \mu)}{\tau(z = 1) + \tau(z = 0)}}$$

2.2 EM Algorithm in ABO Blood Groups [10 pts]

In the ABO blood group system, each individual has a phenotype and a genotype as shown below. The genotype is made of underlying alleles (A, B, O).

Phenotype	Genotype
A	AA
A	AO
B	BB
B	BO
O	OO
AB	AB

In a research experiment, scientists wanted to model the distribution of the population. They collected the phenotype information from the participants as this could be directly observed from the individual's blood group. The scientists, however want to use this data to model the underlying genotype information. In order to help them obtain an understanding, you suggest using the EM algorithm to find out the genotype distribution.

You know that the probability of an allele is present in an individual is independent of the probability of any other allele, i.e.  $P(A, O) = P(A) * P(O)$  and so on. Also note that the genotype pairs: (AO, OA) and (BO, OB) are identical and can be treated as AO, BO respectively. You also know that the alleles follow a multinomial distribution.

$$p(O) = 1 - p(A) - p(B)$$

Let  $n_A, n_B, n_{AB}$  be the number of individuals with the phenotypes A, B, O and AB respectively. Let  $n_{AA}, n_{AO}, n_{BB}, n_{BO}, n_{AB}$  be the numbers of individuals with genotypes AA, AO, BB, BO and AB respectively. The satisfy the following conditions:

$$n_A = n_{AA} + n_{AO}$$
$$n_B = n_{BB} + n_{BO}$$
$$n_{AB} = n_{AB} + n_{BA} + n_{AB} = n$$

Given:

$$p_A = p_B = p_O = \frac{1}{3}$$
$$n_A = 186, n_B = 38, n_O = 284, n_{AB} = 13$$

2.2.1 In the E step, compute the value of  $n_{AA}, n_{AO}, n_{BB}, n_{BO}$  [5pts]

2.2.2 In the M step, find the new value of  $p_A, p_B$  given the updated values from E-step above. (Round off the answer to 3 decimal places) [5pts]



2.2.1 We have:

$$\begin{aligned}p_{AA} &= p_A = \frac{1}{3} = \frac{2}{9} \\p_{AO} &= 2p_Ap_O = 2\left(\frac{1}{3}\right)^2 = \frac{2}{9} \\p_{BB} &= p_B^2 = \left(\frac{1}{3}\right)^2 = \frac{1}{9} \\p_{BO} &= 2p_Bp_O = 2\left(\frac{1}{3}\right)^2 = \frac{2}{9} \\p_{OO} &= p_O^2 = \left(\frac{1}{3}\right)^2 = \frac{1}{9} \\p_{AB} &= 2p_Ap_B = \left(\frac{1}{3}\right)^2 = \frac{2}{9}\end{aligned}$$

We can calculate:

$$\begin{aligned}n_{AA} &= n_A + \frac{p_{AA} + p_{AO}}{p_A} \\&= 186 + \frac{\frac{2}{9} + \frac{1}{9}}{\frac{1}{9}} \\&= 186 + \frac{\frac{3}{9}}{\frac{1}{9}} \\&= 186 \\&= 62 \\n_{AO} &= n_A + \frac{p_{AA} + p_{AO}}{p_A} \\&= 186 + \frac{\frac{2}{9} + \frac{1}{9}}{\frac{1}{9}} \\&= 186 + \frac{\frac{3}{9}}{\frac{1}{9}} \\&= 186 + \frac{2}{3} \\&= \frac{186 \times 3}{3} + \frac{2}{3} \\&= \frac{558 + 2}{3} \\&= \frac{560}{3} \\n_{BB} &= n_B + \frac{p_{BB}}{p_B} \\&= 38 + \frac{\frac{1}{9}}{\frac{1}{9}} \\&= 38 + \frac{1}{9} \\&= \frac{38 \times 9}{9} + \frac{1}{9} \\&= \frac{342 + 1}{9} \\&= \frac{343}{9} \\n_{BO} &= n_B + \frac{p_{BB} + p_{BO}}{p_B} \\&= 38 + \frac{\frac{1}{9} + \frac{2}{9}}{\frac{1}{9}} \\&= 38 + \frac{\frac{3}{9}}{\frac{1}{9}} \\&= 38 + \frac{2}{3} \\&= \frac{114 + 4}{3} \\&= \frac{118}{3} \\&= 39\end{aligned}$$

2.2.2 For the M-step, we will solve an optimization problem, given the constraint  $p_A + p_B + p_O = 1$ .

We have, the log-likelihood function:

$$\begin{aligned}LL(p_A, p_B, p_C) &= \log(p_A^{n_{AA}} + p_{AO}^{n_{AO}} + p_B^{n_{BB}} + p_{BO}^{n_{BO}} + p_O^{n_{OO}} + p_{AB}^{n_{AB}}) \\&= n_{AA} \log p_A + n_{AO} \log p_{AO} + n_{BB} \log p_B + n_{BO} \log p_{BO} + n_{OO} \log p_O + n_{AB} \log p_{AB}\end{aligned}$$

Thus, the Lagrange function would be:

$$L(p_A, p_B, p_C) = LL(p_A, p_B, p_C) - \lambda(p_A + p_B + p_O - 1)$$

To solve the optimization problems, we will find the partial derivatives and equate it to 0.

$$\begin{aligned}\frac{\partial L(p_A, p_B, p_C)}{\partial p_A} &= \frac{n_{AA}}{p_A} - \lambda \\&= \frac{2n_{AA}}{p_A} + \frac{2n_{AO}p_O}{2p_Ap_O} + \frac{2n_{AB}p_B}{2p_Ap_B} - \lambda \\&= \frac{2n_{AA} + n_{AO} + n_{AB}}{p_A} - \lambda = 0\end{aligned}$$

$$\begin{aligned}\frac{\partial L(p_A, p_B, p_C)}{\partial p_B} &= \frac{n_{BB}}{p_B} - \lambda \\&= \frac{2n_{BB}}{p_B} + \frac{2n_{BO}p_O}{2p_Bp_O} + \frac{2n_{AB}p_A}{2p_Bp_A} - \lambda \\&= \frac{2n_{BB} + n_{BO} + n_{AB}}{p_B} - \lambda = 0\end{aligned}$$

$$\begin{aligned}\frac{\partial L(p_A, p_B, p_C)}{\partial p_O} &= \frac{n_{OO}}{p_O} - \lambda \\&= \frac{2n_{AO}}{2p_Ap_O} + \frac{2n_{BO}p_B}{2p_Bp_O} + \frac{2n_{AB}p_A}{2p_Op_A} - \lambda \\&= \frac{2n_{AO} + n_{BO} + n_{AB}}{p_O} - \lambda = 0\end{aligned}$$

$$\begin{aligned}\frac{\partial L(p_A, p_B, p_C)}{\partial \lambda} &= \frac{\partial}{\partial \lambda} [-\lambda(p_A + p_B + p_O - 1)] \\&= (p_A + p_B + p_O - 1) = 0\end{aligned}$$

$$2(n_{AA} + n_{AO} + n_{BB} + n_{BO} + n_{OO} + n_{AB}) - \lambda(p_A + p_B + p_O) = 0$$

Since the sum of all frequencies is  $n$  and the sum of all probabilities is 1, we get:

$$\begin{aligned}\lambda &= 2n = 2 * (186 + 38 + 244 + 13) = 2 * 521 = 1042 \\&= \frac{2n_{AA} + n_{AO} + n_{AB}}{p_A} \\&= \frac{2n_{AA} + n_{AO} + n_{AB}}{\frac{1}{3}} \\&= \frac{2 * 62 + 124 + 13}{\frac{1}{3}} \\&= \frac{1042}{\frac{1}{3}} \\&= 3126\end{aligned}$$

Now, for  $p_B$ ,

$$\begin{aligned}p_B &= \frac{2n_{BB} + n_{BO} + n_{AB}}{p_B} \\&= \frac{2 * 39 + 25 + 13}{\frac{1}{3}} \\&= \frac{142}{\frac{1}{3}} \\&= 426\end{aligned}$$

3. GMM implementation [40 + 10 + 5(bonus) pts]

A Gaussian Mixture Model(GMM) is a probabilistic model that assumes all the data points are generated from a mixture of a finite number of Gaussian Distributions with unknown parameters. Both clustering algorithm in a sense that each data point is assigned to a cluster with a probability. In order to do that, we need to convert our clustering problem into an inference problem.

Given  $N$  samples  $X = [x_1, x_2, \dots, x_N]^T$ , where  $x_i \in \mathbb{R}^D$ . Let  $x$  be a  $K$ -dimensional probability distribution and  $(\mu_i, \Sigma_i)$  be the mean and covariance matrix of the  $k^{th}$  Gaussian distribution in  $\mathbb{R}^D$ .

The GMM object implements EM algorithms for fitting the model and MLE for optimizing its parameters. It also has some particular hypothesis on how the data was generated:

- Each data point  $x_i$  is assigned to a cluster  $k$  with probability of  $\pi_k$  where  $\sum_{k=1}^K \pi_k = 1$
- Each data point  $x_i$  is generated from Multivariate Normal Distribution  $N(\mu_i, \Sigma_i)$  where  $\mu_i \in \mathbb{R}^D$  and  $\Sigma_i \in \mathbb{R}^{D \times D}$

Our goal is to find  $K$ -dimension Gaussian distributions to model our data  $X$ . This can be done by learning the parameters  $\pi, \mu$  and  $\Sigma$  through likelihood function. Detailed derivation can be found in our slide of GMM. The log-likelihood function now becomes:

$$\ln p(x_1, \dots, x_N | \pi, \mu, \Sigma) = \sum_{i=1}^N \ln \left( \sum_{k=1}^K \pi(k) N(x_i | \mu_i, \Sigma_i) \right)$$

From the lecture we know that MLEs for GMM all depend on each other and the responsibility  $\tau$ . Thus, we need to use an iterative algorithm (the EM algorithm) to find the estimate of parameters that maximize our likelihood function. **All detailed derivations can be found in the lecture slide of GMM.**

- E-step: Evaluate the responsibilities

In this step, we need to calculate the responsibility  $\tau$ , which is the conditional probability that a data point belongs to a specific cluster  $k$  if we are given the datapoint, i.e.  $P(z_k | x_i)$ . The formula for  $\tau$  is given below:

$$\tau(z_k) = \frac{\pi_k N(x_i | \mu_i, \Sigma_i)}{\sum_{j=1}^K \pi_j N(x_i | \mu_j, \Sigma_j)}, \quad \text{for } k = 1, \dots, K$$

Note that each data point should have one probability for each component/cluster. For this homework, you will work with  $\tau(z_k)$  which has a size of  $N \times K$  and you should have all the responsibility values in one matrix. **We use gamma as  $\tau$  in this homework.**

- M-step: Re-estimate Parameters

After we obtained the responsibility, we can find the update of parameters, which are given below:

$$\begin{aligned}\mu_k^{new} &= \frac{\sum_{i=1}^N \tau(z_k) x_i}{N_k} \\ \Sigma_k^{new} &= \frac{1}{N_k} \sum_{i=1}^N \tau(z_k) (x_i - \mu_k^{new})^T (x_i - \mu_k^{new}) \\ \pi_k^{new} &= \frac{N_k}{N}\end{aligned}$$

where  $N_k = \sum_{i=1}^N \tau(z_k)$ . Note that the updated value for  $\mu_k$  is used when updating  $\Sigma_k$ . The multiplication of  $\tau(z_k)^T (x_i - \mu_k^{new})^T$  is element-wise so it will preserve the dimensions of  $(x_i - \mu_k^{new})^T$ .

- We repeat E and M steps until the incremental improvement to the likelihood function is small.

### Special Notes

- For undergraduate students: you may assume that the covariance matrix  $\Sigma_i$  is a diagonal matrix, which means the features are independent. (i.e. the red intensity of a pixel is independent of its blue intensity, etc).
- For graduate students: please assume a full covariance matrix.
- The class notes assume that your dataset  $X$  is  $(N, D, N)$ . However, the homework dataset is  $(N, D)$  as mentioned on the instructions, so the formula is a little different from the lecture note in order to obtain the right dimensions of parameters.

### Hints

- DO NOT USE FOR LOOPS OVER N.** You can always find a way to avoid looping over the observation data points in our homework problem. If you have to loop over D or K, that would be fine.
- In part you can initiate  $\pi(k)$  for each of the  $k$ , i.e.  $\pi(k) = \frac{1}{K}, \forall k = 1, 2, \dots, K$ .
- In part 3 you are asked to generate the model for pixel clustering of image. We will need to use a multivariate Gaussian because each image will have  $N$  pixels and  $D = 3$  features, which correspond to red, green, and blue color intensities. It means that each image is a  $(N \times 3)$  dataset matrix. In the following parts, remember  $D = 3$  in this problem.
- To avoid using for loops in numpy, GMMs use the concept of broadcasting. You take a look at the concept [Broadcasting in Numpy](#). Also, some calculations that required different shapes of arrays can be achieved by broadcasting.
- Be careful of the dimensions of your parameters. Before you test anything on the autograder, please look at the instructions below on the shapes of the variables you need to output. This could enhance the functionality of your code and help you debug. Also notice that a **numpy array in shape  $(N, 1)$  is NOT the same as that in shape  $(N, )$** , so be careful and consistent on what you are using. You can see the detailed explanation here: [Difference between numpy.array.shape \[\(R, 1\)\] and \[\(R,\)\]](#)

- The dataset  $X: (N, D)$
- $\mu: (K, D)$
- $\Sigma: (K, D, D)$
- $\pi: (N, K)$
- $\pi$ : array of length  $K$
- $\Pi$ : joint  $(N, K)$

### 3.1 Helper functions [15 pts]

To facilitate some of the operations in the GMM implementation, we would like you to implement the following three helper functions. In these functions, "logit" refers to an input array of size  $(N, D)$ . Remember the goal of helper functions is to facilitate our calculation so **DO NOT USE FOR LOOP ON N**.

#### 3.1.1. softmax [5 pts]

Given  $\text{logit} \in \mathbb{R}^{N \times D}$ , calculate  $\text{prob} \in \mathbb{R}^{N \times D}$ , where  $\text{prob}_{ij} = \frac{\exp(\text{logit}_{ij})}{\sum_{k=1}^D \exp(\text{logit}_{ik})}$ .

Note: it is possible that  $\text{logit}_{ij}$  is very large, making  $\exp(\cdot)$  of it to explode. To make sure it is numerically stable, you need to subtract the maximum for each row of  $\text{logit}$ , and then add it back in your result.

#### 3.1.2. logsumexp [5 pts]

Given  $\text{logit} \in \mathbb{R}^{N \times D}$ , calculate  $s \in \mathbb{R}^N$ , where  $s_i = \log \left( \sum_{j=1}^D \exp(\text{logit}_{ij}) \right)$ . Again, pay attention to the numerical problem. You may want to use similar trick as in the softmax function. Note: This function is used in the call() function which is given, so you will not need in your own implementation. It helps calculate the loss of log-likelihood.

#### 3.1.3. Multivariate Gaussian PDF [5 pts]

You should be able to write your own function based on the following formula, and you are NOT allowed to use outside resource packages other than those we provided.

(for undergrads only) **normaPDF**

Using the covariance matrix as a diagonal matrix with variances of the individual variables appearing on the main diagonal of the matrix and zeros everywhere else means that we assume the features are independent. In this case, the multivariate normal density function simplifies to the expression below:

$$N(x; \mu, \Sigma) = \prod_{i=1}^D \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{1}{2\sigma_i^2}(x_i - \mu_i)^2\right)$$

where  $\sigma_i^2$  is the variance for the  $i^{th}$  feature, which is the diagonal element of the covariance matrix.

(for grads only) **multinormalPDF**

Given the dataset  $X \in \mathbb{R}^{N \times D}$ , the mean vector  $\mu \in \mathbb{R}^D$  and covariance matrix  $\Sigma \in \mathbb{R}^{D \times D}$  for a multivariate Gaussian distribution, calculate the probability  $p \in \mathbb{R}^N$  of each data. The PDF is given by

$$N(X; \mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(X - \mu)^T \Sigma^{-1} (X - \mu)\right)$$

where  $|\Sigma|$  is the determinant of the covariance matrix.

### Hints

- If you encounter "LinAlgError", you can mitigate your number/array by summing a small value before taking the operation, e.g. `np.linalg.inv(Sigma_n + 1e-32)`. You can arrest and handle such error by using [Try and Exception Block](#) in Python.
- In the above calculation, you must avoid computing a  $(N, N)$  matrix. Using the above equation for large  $N$  will crash your kernel and/or give you a memory error on Gradescope. Instead, you can do this same operation by calculating  $(X - \mu)^T \Sigma^{-1} (X - \mu)$  in  $(N, D)$  matrix, transpose it to be a  $(D, N)$  matrix and do an element-wise multiplication with  $(X - \mu)^T$ , which is also a  $(D, N)$  matrix. Lastly, you will need to sum over the 0 axis to get a  $(1, N)$  matrix before proceeding with the rest of the calculation. This uses the fact that doing an element-wise multiplication and summing over the 0 axis is the same as taking the diagonal of the  $(N, N)$  matrix from the matrix multiplication.
- In Numpy implementation for  $\mu$ , you can either use a 2-D array with dimension  $(1, D)$  for each Gaussian Distribution, or a 1-D array with length  $D$ . Same goes for grads, you can also initialize the  $\Sigma$  by  $K$  diagonal matrices. It will become a full matrix after one iteration, as long as you adopt the correct computation.
- Other ways of initialization are acceptable and welcome.

### 3.2 GMM Implementation [25 pts]

Things to do in this problem:

#### 3.2.1. Initialize parameters in `_init_components()` [5 pts]

Examples of how you can initialize the parameters.

- Set the prior probability  $\pi$  the same for each class.
- Initialize  $\mu$  by randomly selecting  $K$  numbers of observations as the initial mean vectors, and initialize the covariance matrix with `np.eye()` for each  $K$ . For grads, you can also initialize the  $\Sigma$  by  $K$  diagonal matrices. It will become a full matrix after one iteration, as long as you adopt the correct computation.
- Other ways of initialization are acceptable and welcome.

#### 3.2.2. Formulate the log-likelihood function `_ll_joint()` [5 pts]

The log-likelihood function is given by:

$$\ell(\theta) = \ln \left( \prod_{i=1}^N \left( \sum_{k=1}^K \pi(k) N(x_i | \mu_i, \Sigma_i) \right) \right)$$

In this part, we will generate a  $(N, K)$  matrix where each datapoint  $x_i, \forall i = 1, \dots, N$  has  $K$  log-likelihood numbers. Thus, for each  $i = 1, \dots, N$  and  $k = 1, \dots, K$ :

$$\text{log\_likelihood}[i, k] = \log \pi_k + \log N(x_i | \mu_k, \Sigma_k)$$

### Hints:

- If you encounter "ZeroDivisionError" or "RuntimeWarning: divide by zero encountered in log", you can mitigate your number/array by summing a small value before taking the operation, e.g. `np.log(Sigma_k + 1e-32)`.
- You need to use the Multivariate Normal PDF function you created in the last part. Remember the PDF function is for each Gaussian Distribution (i.e. for each  $k$ ) so you need to use a for loop over  $K$ .

#### 3.2.3. Setup Iterative steps for EM Algorithm [5+10 pts]

You can find the detail instruction in the above description box.

### Hints:

- For E steps, we already get the log-likelihood `_ll_joint()` function. This is not the same as responsibilities ( $\tau$ ), but you should be able to finish this part with just a few lines of code by using `_ll_joint()` and `softmax` as defined above.
- For Undergrads: Try to simplify your calculation for  $\Sigma$  in M steps as you assumed independent components. Make sure you are only taking the diagonal terms of your calculated covariance matrix.

```
In [13]: class GMM(object):
def __init__(self, X, K, max_iters = 100): # No need to change
    """
    Args:
        X: the observations/datapoints, N x D numpy array
        K: number of clusters/components
        max_iters: maximum number of iterations (used in EM implementation)
    """
    self.points = X
    self.max_iters = max_iters

    self.N = self.points.shape[0] #number of observations
    self.D = self.points.shape[1] #number of features
    self.K = K #number of components/clusters

    #Helper function for you to implement
    def softmax(self, logit): # [5pts]
    """
    Args:
        logit: N x D numpy array
    Return:
        prob: N x D numpy array. See the above function.
    """
    #max of each row in logit
    max_logit = np.max(logit, axis=1)
    #subtract max of each row from logit
    logit = max_logit[:, None]
    #initialize prob
    prob = np.exp(logit)
    prob /= np.sum(np.exp(logit), axis=1)[:, None]
    return prob
    raise NotImplementedError

def logsumexp(self, logit): # [5pts]
    """
    Args:
        logit: N x D numpy array
    Return:
        s: N x 1 array where s[i,0] = logsumexp(logit[i,:]). See the above function
    """
    #max of each row in logit
    max_logit = np.max(logit, axis=1)
    #subtract max of each row from logit
    logit = max_logit[:, None]
    #initialize s
    s = np.log(np.sum(np.exp(logit), axis=1)[:, None] + 1e-32)
    #add max back to s
    s = max_logit[:, None] + s
    return s
    raise NotImplementedError

# for undergrad student
def normalPDF(self, logit, mu_i, sigma_i): # [5pts]
    """
    Args:
        logit: N x D numpy array
        mu_i: 1xD numpy array (or array of length D), the center for the ith gaussian.
        sigma_i: 1xDxD 3-D numpy array (or 2-D 2-D numpy array), the covariance matrix of the ith gaussian.
    Return:
        pdf: 1xN numpy array (or array of length N), the probability distribution of N data for the ith gaussian.
    Hint:
        np.diagonal() should be handy.
    """
    raise NotImplementedError

# for grad students
def multinormalPDF(self, logits, mu_i, sigma_i): # [5pts]
    """
    Args:
        logit: N x D numpy array
        mu_i: 1xD numpy array (or array of length D), the center for the ith gaussian.
        sigma_i: 1xDxD 3-D numpy array (or 2-D 2-D numpy array), the covariance matrix of the ith gaussian.
    Return:
        pdf: 1xN numpy array (or array of length N), the probability distribution of N data for the ith gaussian.
    Hint:
        np.linalg.det() and np.linalg.inv() should be handy.
    """
    N, D = logits.shape[0], logits.shape[1]
    mu_i = logits.reshape((D,1))
    #term1
    t1 = np.dot(logits - mu_i.T[None,:], np.linalg.inv(sigma_i*0.00001)).T
    #term2
    t2 = (logits - mu_i.T[None,:]).T
    #element-wise multiplication
    t3 = np.exp(np.sum(np.multiply(t1,t2), axis=0)*(0.5))
    #term4
    t4 = np.power(np.linalg.det(sigma_i), -0.5)
    #term5
    t5 = np.power(np.power((2*np.pi), D/2), -1)
    #putting it all together
    pdf = t5*t4*t3
    pdf.reshape((1,N))
    return pdf
    raise NotImplementedError

#
def _init_components(self, **kwargs): # [5pts]
    """
    Args:
        kwargs: any other arguments you want
    Return:
        pi: numpy array of length K, prior
        mu: KxD numpy array, the center for each gaussian.
        sigma: KxDxD numpy array, the diagonal standard deviation of each gaussian.
    """
    N, D, K = self.N, self.D, self.K
    #initialize pi as same probability for each class
    pi = np.array([1/K] * K)

    #Randomly select K observations as mean vectors
    indices = np.random.choice(points.shape[0], size=K, replace=False)
    #from the numpy array with these indices
    mu = points[indices,:]

    #sigma is np.eye for each K
    sigma = np.array([np.eye(D)*K])
    return pi, mu, sigma
    raise NotImplementedError

def _ll_joint(self, pi, mu, sigma, **kwargs): # [10 pts]
    """
    Args:
        pi: np array of length K, the prior of each component
        mu: KxD numpy array, the center for each gaussian.
        sigma: KxDxD numpy array, the diagonal standard deviation of each gaussian. You will have K xxDxD numpy array for full covariance matrix case
    Return:
        ll(log-likelihood): NxK array, where ll(i, k) = log pi(k) + log NormalPDF(points_i | mu[k], sigma[k])
    """
    N, D, K = self.N, self.D, self.K
    logits = self.points
    # initialize ll
    ll = np.zeros((N, K))

    # term1
    t1 = np.log(pi + 1e-32)
    # add ll to ll
    ll = ll + t1

    # term2
    for i in range(K):
        mu_i = mu[i,:]
        sigma_i = sigma[i,:,:]
        pdf = self.multinormalPDF(logits, mu_i, sigma_i)
        ll[:, i] += np.log(pdf + 1e-32)
    #
    np.nan_to_num(ll)

    return ll

def _E_step(self, pi, mu, sigma, **kwargs): # [5pts]
    """
    Args:
        pi: np array of length K, the prior of each component
        mu: KxD numpy array, the center for each gaussian.
        sigma: KxDxD numpy array, the diagonal standard deviation of each gaussian. You will have Kx DxD numpy array for full covariance matrix case
    Return:
        gamma(tau): NxK array, the posterior distribution (a.k.a, the soft cluster assignment) for each observation.
    Hint:
        You should be able to do this with just a few lines of code by using _ll_joint() and softmax() defined above.
    """
    N, D, K = self.N, self.D, self.K

    #initialize tau
    tau = np.zeros((N, K))

    for i in range(K):
        mu_i = mu[i,:]
        sigma_i = sigma[i,:,:]
        pdf = self.multinormalPDF(logits, mu_i, sigma_i)
        if np.isnan(pdf).any():
            print('mu has 0')
            tau[:, i] = np.dot(pi[i], pdf).T

    #term2, summation
    t2 = np.sum(tau, axis=1)
    tau = tau/t2[:, None]

    return tau
    raise NotImplementedError

def _M_step(self, gamma, **kwargs): # [10pts]
    """
    Args:
        gamma(tau): NxK array, the posterior distribution (a.k.a, the soft cluster assignment) for each observation.
    Return:
        pi: np array of length K, the prior of each component
        mu: KxD numpy array, the center for each gaussian.
        sigma: KxDxD numpy array, the diagonal standard deviation of each gaussian. You will have K xxDxD numpy array for full covariance matrix case
    Hint:
        There are formulas in the slide and in the above description box.
    """
    N, D, K = self.N, self.D, self.K
    pi, mu, sigma = self._init_components(**kwargs)
    logits = self.points
    # initialize ll
    ll = np.zeros((N, K))

    # term1
    t1 = (logits - mu_new[i,:]).T
    #term2, element-wise multiplication
    t2 = np.multiply(gamma[:, i], t1)
    #term3, dot product
    t3 = np.dot(t2, (logits - mu_new[i,:]))
    sigma_new[i] = t3/N_k[i]

    #pi_k
    pi_new = np.zeros(K)
    pi_new = N_k/N
    pi_new.reshape((K))

    return pi_new, mu_new, sigma_new
    raise NotImplementedError

def _call(self, abs_tol=1e-16, rel_tol=1e-16, **kwargs): # No need to change
    """
    Args:
        abs_tol: convergence criteria w.r.t absolute change of loss
        rel_tol: convergence criteria w.r.t relative change of loss
        kwargs: any additional arguments you want
    Return:
        gamma(tau): NxK array, the posterior distribution (a.k.a, the soft cluster assignment) for each observation
        (pi, mu, sigma): (1xK) np array, KxD numpy array, KxDxD numpy array
    Hint:
        You do not need to change it. For each iteration, we process E and M steps, then update the parameters.
    """
    pi, mu, sigma = self._init_components(**kwargs)
    if np.isnan(pi).any():
        print('pi has 0')
    if np.isnan(mu).any():
        print('mu has 0')
    if np.isnan(sigma).any():
        print('sigma has 0')
    pbar = tqdm(range(self.max_iters))
    for it in pbar:
        # E-step
        gamma = self._E_step(pi, mu, sigma)
        # M-step
        pi, mu, sigma = self._M_step(gamma)
        pi = np.nan_to_num(pi)
        mu = np.nan_to_num(mu)
        sigma = np.nan_to_num(sigma)
        if np.isnan(pi).any():
            print('pi_new has 0')
        if np.isnan(mu).any():
            print('mu_new has 0')
        if np.isnan(sigma).any():
            print('sigma_new has 0')
        # calculate the negative log-likelihood of observation
        joint_ll = self._ll_joint(pi, mu, sigma)
        joint_ll = np.nan_to_num(joint_ll)
        if np.isnan(joint_ll).any():
            print('joint_ll has 0')
        loss = -np.sum(self.logsumexp(joint_ll))
        if it:
            diff = np.abs(prev_loss - loss)
            if diff < abs_tol and diff / prev_loss < rel_tol:
                break
        prev_loss = loss
        pbar.set_description(f'Iter {it}, loss: {4*diff} % (it, loss)')
    return gamma, (pi, mu, sigma)
```

### 3.3 Japanese art and pixel clustering [10pts + 5pts]

Ukiyo-e is a Japanese art genre predominant from the 17th through 19th centuries. In order to produce the intricate patterns that came to represent the genre, artists carved wood blocks with the patterns for each color in a design. Paint would be applied to the block and later transferred to the print to form the image. In this section, you will use your GMM algorithm implementation to do pixel clustering and estimate how many wood blocks were likely used to produce a single print. That is to say, how many wood blocks would be used to produce the original print. (Hint: you can justify your answer based on visual inspection of the resulting images or on a different metric of your choosing)

You do NOT need to submit your code for this question to the autograder. Instead you should include whatever images/information you find relevant in the report.

```
In [14]: # helper function for performing pixel clustering. You don't have to modify it
def cluster_pixels_gmm(image, K):
    """Clusters pixels in the input image
    Args:
        image: input image of shape(H, W, 3)
        K: number of components
    Return:
        clustered_img: image of shape(H, W, 3) after pixel clustering
    """
    im_height, im_width, im_channel = image.shape
    flat_img = np.reshape(image, [-1, im_channel]).astype(np.float32)
    gamma, (pi, mu, sigma) = GMM(flat_img, K = K, max_iters = 100)()
    cluster_img = np.argmax(gamma, axis=1)
    centers = mu

    gmm_img = np.reshape(centers[cluster_ids], (im_height, im_width, im_channel))

    return gmm_img

# helper function for plotting images. You don't have to modify it
def plot_images(img_list, title_list, figsize=(20, 10)):
    assert len(img_list) == len(title_list)
    fig, axes = plt.subplots(1, len(title_list), figsize=figsize)
    for i, ax in enumerate(axes):
        ax.imshow(img_list[i] / 255.0)
        ax.set_title(title_list[i])
        ax.axis('off')
```

```
In [15]: # pick 2 of the images in this list:
url10 = "https://upload.wikimedia.org/wikipedia/commons/b/1/Utagawa_Kunisada_I_328c_1832k29_Dawn_at_Fu_tami-ga-ura.jpg"
url11 = "https://upload.wikimedia.org/wikipedia/commons/9/95/Hokusai_328182k29_Cuckoo_and_Azaleas.jpg"
url12 = "https://upload.wikimedia.org/wikipedia/commons/7/74/Kitao-Shigemasa_328177k29_Geisha_and_a_serveant_carrying_her_shamisen_box.jpg"
url13 = "https://upload.wikimedia.org/wikipedia/commons/4/10/Kuniyoshi_Utagawa32C_Saikuden_Series_4.jpg"

# example of loading image from url10
image10 = imageio.imread(image10_core.urlopen(url10).read())
image11 = imageio.imread(image11_core.urlopen(url11).read())

# this is for you to implement
def find_n_woodblocks(image, min_clusters=5, max_clusters=15):
    """
    Using the helper function above to find the optimal number of woodblocks that can appropriately produce a single image.
    You can simply examine the answer based on your visual inspection (i.e. looking at the resulting images) or provide any metrics you prefer.
    Args:
        image: input image of shape(H, W, 3)
        min_clusters, max_clusters: the minimum and maximum number of clusters you should test with. Default are 5 <= n <= 15. (Usually the maximum number of clusters would not exceed 15)
    Return:
        plot: comparison between original image and image pixel clustering.
        optional: any other information/metric/plot you think is necessary.
    """
    img_list = [image]
    for num_clusters in range(min_clusters, max_clusters+1):
        np.seterr(divide='ignore', invalid='ignore')
        gmm_img = cluster_pixels_gmm(image, num_clusters)
        img_list.append(gmm_img)
        title_list.append(f"GMM {num_clusters} = " + str(num_clusters))

    plot_images(img_list, title_list)

    find_n_woodblocks(image0, 5, 8)
    find_n_woodblocks(image1, 5, 8)
    # raise NotImplementedError
```

For my experiments, I used the images 2 and 4, and ran the GMM clustering 6 for various number of clusters. As we can see from the images above, using 8 clusters almost generates the original image. Hence, I estimate that the number of woodblocks used to produce each print were 8.

### (Bonus for All) [5 pts]

Compare the full covariance matrix with the diagonal covariance matrix in GMM. Can you explain why the images are different with the same clusters? Note: You will have to implement both multinormalPDF and normalPDF, and add a few arguments in the original `_ll_joint` and `_M_step` function to indicate which matrix you are using. You will earn full credit only if you implement both functions AND explain the reason.

```
In [ ]: def compare_matrix(image, K):
    """
    Args:
        image: input image of shape(H, W, 3)
        K: number of components
    Return:
        plot: comparison between full covariance matrix and diagonal covariance matrix.
    """
    raise NotImplementedError
```

```
In [ ]: compare_matrix(image1, 5)
```

### 4. (Bonus for Grad and Undergrad) A Wrench in the Machine [30pts]

Learning to work with messy data is a hallmark of a well-rounded data scientist. In most real-world settings the data given



#### 4.1.a Data Cleaning

The first step is to break up the whole dataset into clear parts. All the data is randomly shuffled in one csv file. In order to move forward, the data needs to be split into three separate arrays:

- labeled\_complete: containing the complete characterization data and corresponding labels (broken = 1 and OK = 0)
- labeled\_incomplete: containing partial characterization data and corresponding labels (broken = 1 and OK = 0)
- unlabeled\_complete: containing only complete material characterization results

```
In [ ]: def complete_data():
    """
    Args:
        data: N x D numpy array
    Return:
        labeled_complete: n x D array where values contain both complete features and labels
    """
    raise NotImplementedError

def incomplete_data():
    """
    Args:
        data: N x D numpy array
    Return:
        labeled_incomplete: n x D array where values contain incomplete features but complete labels
    """
    raise NotImplementedError

def unlabeled_data():
    """
    Args:
        data: N x D numpy array
    Return:
        unlabeled_complete: n x D array where values contain complete features but incomplete labels
    """
    raise NotImplementedError
```

#### 4.1.b KNN [10pts]

The second step in this task is to clean the Labeled\_incomplete dataset by filling in the missing values with probable ones derived from complete data. A useful approach to this type of problem is using a k-nearest neighbors (k-NN) algorithm. For this application, the method consists of replacing the missing value of a given point with the mean of the closest k-neighbors to that point.

```
In [ ]: class CleanData(object):
    def __init__(self): # No need to implement
        pass

    def pairwise_dist(self, x, y): # [0pts] - copy from kmeans
        """
        Args:
            x: N x D numpy array
            y: M x D numpy array
        Return:
            dist: N x M array, where dist2[i, j] is the euclidean distance between
                x[i, :] and y[j, :]
        """
        #12 norm, ord=2
        dist = np.linalg.norm(x[:, np.newaxis, :] - y, ord=2, axis=2)
        return dist
    # raise NotImplementedError

    def __call__(self, incomplete_points, complete_points, K, **kwargs): # [10pts]
        """
        Args:
            incomplete_points: N_incomplete x (D+1) numpy array, the incomplete labeled observations
            complete_points: N_complete x (D+1) numpy array, the complete labeled observations
            K: integer, corresponding to the number of nearest neighbors you want to base your calculation on
            kwargs: any other args you want
        Return:
            clean_points: (N_incomplete + N_complete) x (D+1) X D numpy array of length K, containing both complete points and recently filled points

        Hints:
            (1) You want to find the k-nearest neighbors within each class separately;
            (2) There are missing values in all of the features. It might be more convenient to address each feature at a time.
        """
        raise NotImplementedError
```

Below is a good expectation of what the process should look like on a toy dataset. If your output matches the answer below, you are on the right track.

```
In [ ]: complete_data = np.array([[1.,2.,3.,1],[7.,8.,9.,0],[16.,17.,18.,1],[22.,23.,24.,0]])
incomplete_data = np.array([[11.,np.nan,3.,1],[7.,np.nan,9.,0],[np.nan,17.,18.,1],[np.nan,23.,24.,0]])

clean_data = CleanData().complete_data(complete_data, 2)
print("""Expected Answer - k = 2 """)
print("""complete data==
[[ 1.  5.  3.  1.]
 [ 7.  8.  9.  0.]
 [16. 17. 18.  1.]
 [22. 23. 24.  0.]]
incomplete data==
[[ 1. nan  3.  1.]
 [ 7. nan  9.  0.]
 [nan 17. 18.  1.]
 [nan 23. 24.  0.]]
==clean_data==
[[ 1.  2.  3.  1.]
 [ 7.  8.  9.  0.]
 [16. 17. 18.  1.]
 [22. 23. 24.  0.]
 [14.9 23.  24.  0.]
 [ 7. 15.5  9.  0.]
 [ 8.5 17. 18.  1.]
 [ 1.  9.5  3.  1.]]""")

print("""\n*** My Answer - k = 2 ***\n""")
print(clean_data)
```

#### 4.2 Getting acquainted with semi-supervised learning approaches. [5pts]

You will implement a version of the algorithm presented in Table 1 of the paper "[Text Classification from Labeled and Unlabeled Documents](#)" by Nigam et al. (2000). While you are recommended to read the whole paper this assignment focuses on items 1–5.2 and 6.1. Write a brief summary of three interesting highlights of the paper (50-word maximum).

#### 4.3 Implementing the EM algorithm. [10 pts]

In your implementation of the EM algorithm proposed by Nigam et al. (2000) on Table 1, you will use a Gaussian Naive Bayes (GNB) classifier as opposed to a naive Bayes (NB) classifier. (Hint: Using a GNB in place of an NB will enable you to reuse most of the implementation you developed for GMM in this assignment. In fact, you can successfully solve the problem by simply modifying the call method)

```
In [ ]: class SemiSupervised(object):
    def __init__(self): # No need to implement
        pass

    def softmax(self, logits): # [0 pts] - can use same as for GMM
        """
        Args:
            logits: N x D numpy array
        """
        raise NotImplementedError

    def logsumexp(self, logits): # [0 pts] - can use same as for GMM
        """
        Args:
            logits: N x D numpy array
        Return:
            s: N x 1 array where s[i,0] = logsumexp(logits[i,:])
        """
        raise NotImplementedError

    def _init_components(self, points, K, **kwargs): # [5 pts] - modify from GMM
        """
        Args:
            points: Nx(D+1) numpy array, the observations
            K: number of components
            kwargs: any other args you want
        Return:
            pi: numpy array of length K, prior
            mu: KxD numpy array, the center for each gaussian.
            sigma: KxDxD numpy array, the diagonal standard deviation of each gaussian.

        Hint: The paper describes how you should initialize your algorithm.
        """
        raise NotImplementedError

    def _ll_joint(self, points, pi, mu, sigma, **kwargs): # [0 pts] - can use same as for GMM
        """
        Args:
            points: Nx(D) numpy array, the observations
            pi: np array of length K, the prior of each component
            mu: KxD numpy array, the center for each gaussian.
            sigma: KxDxD numpy array, the diagonal standard deviation of each gaussian.
        Return:
            ll(log-likelihood): NKx array, where ll(i, j) = log pi(j) + log NormalPDF(points_i | mu[j], sigma[j])

        Hint: Assume that the three properties of the lithium-ion batteries (multivariate gaussian) are independent.
            This allows you to treat it as a product of univariate gaussians.
        """
        raise NotImplementedError

    def _E_step(self, points, pi, mu, sigma, **kwargs): # [0 pts] - can use same as for GMM
        """
        Args:
            points: Nx(D) numpy array, the observations
            gamma: NKx array, the posterior distribution (a.k.a, the soft cluster assignment) for each observation.
        Return:
            pi: np array of length K, the prior of each component
            mu: KxD numpy array, the center for each gaussian.
            sigma: KxDxD numpy array, the diagonal standard deviation of each gaussian.

        Hint: There are formulas in the slide.
        """
        raise NotImplementedError

    def _call__(self, points, K, max_iters=100, abs_tol=1e-16, rel_tol=1e-16, **kwargs): # [5 pts] - m
from GMM
        """
        Args:
            points: Nx(D) numpy array, where N is # points and D is the dimensionality
            K: number of clusters
            max_iters: maximum number of iterations
            abs_tol: convergence criteria w.r.t absolute change of loss
            rel_tol: convergence criteria w.r.t relative change of loss
            kwargs: any additional arguments you want
        Return:
            gamma: NKx array, the posterior distribution (a.k.a, the soft cluster assignment) for each observation.
            (pi, mu, sigma): (1xK np array, KxD numpy array, KxD numpy array), mu and sigma.
        """
        raise NotImplementedError
```

#### 4.4 Demonstrating the performance of the algorithm. [5pts]

Compare the classification error based on the Gaussian Naive Bayes (GNB) classifier you implemented following the Nigam et al. (2000) approach to the performance of a GNB classifier trained using only labeled data. Since you have not covered supervised learning in class, you are allowed to use the scikit learn library for training the GNB classifier based only on labeled data: [https://scikit-learn.org/stable/modules/generated/sklearn.naive\\_bayes.GaussianNB.html](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.GaussianNB.html).

```
In [ ]: from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score

class ComparePerformance(object):
    def __init__(self): #No need to implement
        pass

    def accuracy_semi_supervised(self, points, independent, n=8):
        """
        Args:
            points: Nx(D+1) numpy array, where N is the number of points in the training set, D is the dimensionality, the last column represents the labels (when available) or a flag that allows you to separate the unlabeled data.
            independent: Nx(D+1) numpy array, where N is # points and D is the dimensionality and the last column are the correct labels
        Return:
            accuracy: floating number
        """
        raise NotImplementedError

    def accuracy_GNB_onlycomplete(self, points, independent, n=8):
        """
        Args:
            points: Nx(D+1) numpy array, where N is the number of only initially complete labeled points in the training set, D is the dimensionality, the last column represents the labels.
            independent: Nx(D+1) numpy array, where N is # points and D is the dimensionality and the last column are the correct labels
        Return:
            accuracy: floating number
        """
        raise NotImplementedError

    def accuracy_GNB_cleandata(self, points, independent, n=8):
        """
        Args:
            points: Nx(D+1) numpy array, where N is the number of clean labeled points in the training set, D is the dimensionality, the last column represents the labels.
            independent: Nx(D+1) numpy array, where N is # points and D is the dimensionality and the last column are the correct labels
        Return:
            accuracy: floating number
        """
        raise NotImplementedError
```

```
In [ ]: from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score
# load and clean data for the next section
telemetry = np.loadtxt('data/telemetry.csv', delimiter=',')

labeled_complete = complete(telemetry)
labeled_incomplete = incomplete(telemetry)
unlabeled = unlabeled(telemetry)

clean_data = CleanData().labeled_incomplete, labeled_complete, 7)
# load unlabeled set
# append unlabeled flag
unlabeled_flag = ~np.ones((unlabeled.shape[0],1))
unlabeled = np.concatenate((unlabeled, unlabeled_flag), 1)
unlabeled = np.delete(unlabeled, -1, axis=1)

# =====
# SEMI SUPERVISED

# format training data
points = np.concatenate((clean_data, unlabeled),0)

# train model
(pi, mu, sigma) = SemiSupervised()._E_step(points, 7)

# COMPARISON

# load test data
independent = np.loadtxt('data/validation.csv', delimiter=',')

# classify test data
classification = SemiSupervised()._E_step(independent[:,1:8], pi, mu, sigma)
classification = np.argmax(classification,axis=1)

# =====
print("""=====COMPARISON=====""")
print("""SemiSupervised Accuracy: """, ComparePerformance().accuracy_semi_supervised(classification, independent))
print("""Supervised with clean data: GNB Accuracy: """, ComparePerformance().accuracy_GNB_onlycomplete(labeled_complete, independent))
print("""Supervised with only complete data: GNB Accuracy: """, ComparePerformance().accuracy_GNB_cleandata(clean_data, independent))
```