

Fall 2020 CX464/CS7641 A Homework 2

Instructor: Dr. Mahdi Roozbahani

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

- No unapproved extension of the deadline is allowed. Late submission will lead to 0 credit.
- Discussion is 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 answers but will not be accepted.
- If a question requires a picture, you could use this syntax `` to include them within your python 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 you are satisfied with your code. 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 a problem.
- 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 3.7](#), 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]: #####
#####
from __future__ import absolute_import
from __future__ import print_function
from __future__ import division

%matplotlib inline

import sys
import matplotlib
import numpy as np
import matplotlib.pyplot as plt
from matplotlib import image
from mpl_toolkits.mplot3d import axes3d
from tqdm import tqdm
import math

print('Version information')

print('python: {}'.format(sys.version))
print('matplotlib: {}'.format(matplotlib.__version__))
print('numpy: {}'.format(np.__version__))

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

# Load images
import imageio

Version information
python: 3.8.3 (default, Jul 2 2020, 11:26:31)
[Clang 10.0.0 ]
matplotlib: 3.2.2
numpy: 1.18.5

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

KMeans is trying to solve the following optimization problem:

arg min_K \sum_{i=1}^N ||x_i - \mu_k||^2

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

1.1 pairwise distance [5pts]

In this section, you are asked to implement pairwise_distance function.

Given X \in R^{N \times D} and Y \in R^{M \times D}, obtain the pairwise distance dist \in R^{N \times M} using the euclidean distance metric, where dist_{ij} = ||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]
"""
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, :]
"""
# raise NotImplementedError
N, D = np.shape(x)
reshaped_x = np.reshape(x, (N, 1, D))
euclidean = np.sqrt(np.sum(np.square(reshaped_x - y), axis=2))
return euclidean

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.
"""
# raise NotImplementedError
_, D = np.shape(points)
options = points.flatten()
centersArray = np.random.choice(options, (K, D))
return centersArray

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.
# raise NotImplementedError
_, D = np.shape(points)
distanceArray = self.pairwise_dist(points, centers)
cluster_idx = np.argmax(distanceArray, axis=1)
return cluster_idx

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, _ = np.shape(old_centers)
centers = np.copy(old_centers)
for k in range(K):
idxes = np.where(cluster_idx == k)[0]
if idxes.size > 0:
members = points[idxes]
centers[k, :] = np.mean(points[idxes], axis=0)
return centers

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.
"""
K, _ = np.shape(centers)
loss = 0.0
for k in range(K):
idxes = np.where(cluster_idx == k)[0]
if idxes.size > 0:
members = points[idxes]
center = centers[k]
loss += np.sum(np.square(members - center))
return loss

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_iter: 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: final loss value of the objective function of KMeans
"""
centers = self._init_centers(points, K, **kwargs)
for it 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 it:
diff = np.abs(prev_loss - loss)
if diff <= abs_tol and diff <= rel_tol:
break
prev_loss = loss
if verbose:
print('Iter %d, loss: %.4e' % (it, 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: number of clusters
Return:
losses: an array of loss denoting the loss of each number of clusters
"""
# raise NotImplementedError
x_values = []
y_values = []
for k in range(1, max_K + 1):
x_values.append(k)
_, loss = self._call(data, k)
y_values.append(loss)
plt.plot(x_values, y_values)
return np.array(y_values)
```

```
In [3]: # Helper function for checking the implementation of pairwise_distance function. Please DO NOT change this function
% TEST CASE
np.random.seed(0)
x = np.random.randn(2, 2)
y = np.random.randn(3, 2)

print('*** Expected Answer ***')
print('==x==')
[[ 0.62434536 -0.61175641]
 [ 0.52817175 -1.07296862]]
print('==y==')
[[ 0.86540763 -2.3015387 ]
 [ 1.74481176 -0.7612069 ]
 [ 0.3190391 -0.24937038]]
==dist==
[[ 1.85239052 0.19195729 1.35467638]
 [ 1.85780729 2.29426447 1.18155842]]

print('*** My Answer ***')
print('==x==')
[[ 0.62434536 -0.61175641]
 [ 0.52817175 -1.07296862]]
print('==y==')
[[ 0.86540763 -2.3015387 ]
 [ 1.74481176 -0.7612069 ]
 [ 0.3190391 -0.24937038]]
==dist==
[[ 1.85239052 0.19195729 1.35467638]
 [ 1.85780729 2.29426447 1.18155842]]

*** My Answer ***
==x==
[[ 0.62434536 -0.61175641]
 [ 0.52817175 -1.07296862]]
print('==y==')
[[ 0.86540763 -2.3015387 ]
 [ 1.74481176 -0.7612069 ]
 [ 0.3190391 -0.24937038]]
==dist==
[[ 1.85239052 0.19195729 1.35467638]
 [ 1.85780729 2.29426447 1.18155842]]

In [4]: # Test kmeans
np.random.seed(1)
points = np.random.randn(100, 2)

cluster_idx2, centers2, loss2 = KMeans()(points, 2)
cluster_idx5, centers5, loss5 = KMeans()(points, 5)

print('*** Expected Answer ***')
print('==centers2==')
[[ -0.23265213 0.66957783]
 [ 0.61791745 -0.59496962]]
==centers5==
[[ 0.94945532 -1.42382563]
 [ 0.64137918 0.09830081]
 [ -0.51672295 -0.35410285]
 [ -0.07747868 1.08964491]
 [ 1.93010934 0.48561944]]
==loss2==
105.0662377653986
==loss5==
53.0865571656247

*** My Answer ***
==centers2==
[[ 0.35647907 -0.77810103]
 [ 0.39337497 -0.72135245]]
==centers5==
[[ 0.83959261 -1.65541692]
 [ -0.87953914 0.59594955]
 [ -0.24569191 -0.46721266]
 [ 0.39043328 -0.1441221 ]
 [ 0.61753564 0.98757263]]
==loss2==
104.93750103699625
==loss5==
44.980763304237215
```

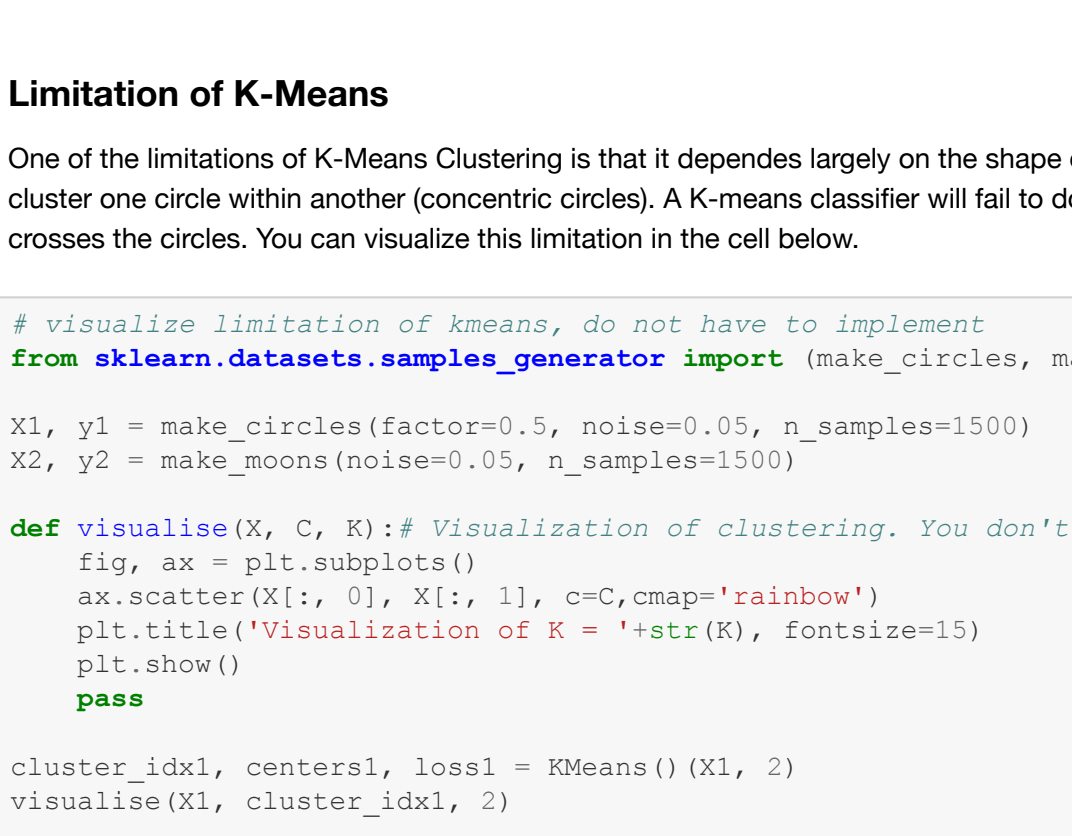
```
In [5]: def image_to_matrix(image_file, grays=False):
"""
Convert .png image to matrix
of values.
params:
image_file = str
grays = Boolean
returns:
img = (color) np.ndarray[np.ndarray[np.ndarray[...]]]
grayscale np.ndarray[np.ndarray[np.ndarray[...]]]
"""
img = plt.imread(image_file)
# in case of 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 [6]: image_values = image_to_matrix('./images/01rd_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 = 5 # feel free to change this value
cluster_idx, centers, loss = KMeans()(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]

# Visualization of image values
plt.figure(figsize=(9,12))
plt.imshow(updated_image_values)
plt.show()
```



```
In [7]: KMeans().find_optimal_num_clusters(image_values)
Out[7]: array([25575.97020067, 14136.69482422, 10225.6109557 , 7720.83544922,
2585.13922119, 4221.15609741, 4432.25927734, 3511.813797 ,
2320.43060703, 2368.43972655, 2106.70523071, 1352.51284002,
1851.87907028, 1707.80400848, 1552.42409134])
```

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 an object has 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 [8]: 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
"""
# raise NotImplementedError
members_idxes = np.where(labels == cluster_idx)
members = data[members_idxes]
n, d = np.shape(members)
reshaped_members = np.reshape(members, (n, 1, d))
intra_dist = np.sqrt(np.sum(np.square(reshaped_members - members), axis=2))
intra_dist_cluster = np.sum(intra_dist, axis=0) / (n - 1)
return intra_dist_cluster

def inter_cluster_dist(cluster_idx, data, labels): # [4 pts]
"""
Calculates the average distance from one cluster 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
"""
# raise NotImplementedError
members_idxes = np.where(labels == cluster_idx)
members = data[members_idxes]
n, d = np.shape(members)
reshaped_members = np.reshape(members, (n, 1, d))
clusters = np.unique(labels)
inter_dist_cluster = np.full(n, math.inf)
for cluster in clusters:
neigh_idxes = np.where(labels == cluster)
neigh = data[neigh_idxes]
inter_dist = np.sqrt(np.sum(np.square(reshaped_members - neigh), axis=2))
inter_dist_cluster_holder = np.mean(inter_dist, axis=0)
inter_dist_cluster = np.minimum(inter_dist_cluster_holder, inter_dist_cluster)
return inter_dist_cluster

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
"""
# raise NotImplementedError
s_clusters = np.unique(labels)
W_ = np.shape(data)
sum = 0.0
for s_cluster in s_clusters:
intra_dist = intra_cluster_dist(s_cluster, data, labels)
inter_dist = inter_cluster_dist(s_cluster, data, labels)
max_denominator = np.maximum(intra_dist, inter_dist)
coeff_array = (inter_dist - intra_dist) / max_denominator
sum += np.sum(coeff_array)
return sum / N
```

```
In [9]: 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()(data, k)
silhouette_coefficients.append(silhouette_coefficient(data, labels))
plt.plot(clusters, silhouette_coefficients)
return silhouette_coefficients
```

```
Out[9]: [0.2517641125992784,
0.2664612313313146,
0.2715465952747097,
0.2795505088461313,
0.2780854520051484,
0.293374974263646,
0.293142930312764013,
0.28481059954311105,
0.2780854520051484,
0.2764539194534896,
0.2788395422930645,
0.2899429875436111]
```



Limitation of K-Means

One of the limitations of K-Means Clustering is that it depends largely on the shape of the dataset. A common example of this is trying to cluster one circle within another (concentric circles). A K-means classifier will fail to do this and will end up effectively drawing a line which crosses the circles. You can visualize this limitation in the cell below.

```
In [10]: # visualize limitation of k-means, do not have to implement
from sklearn.datasets.samples_generator import make_circles, make_moons

X1, y1 = make_circles(factor=0.5, noise=0.05, n_samples=1500)
X2, y2 = make_moons(noise=0.05, n_samples=1500)

def visualize(X, C, R): # Visualization of clustering. You don't need to change this function
fig, ax = plt.subplots()
ax.scatter(X[:, 0], X[:, 1], c=C, cmap='rainbow')
plt.title('Visualization of K = {}'.format(R), fontsize=15)
plt.show()

pass

cluster_idx1, centers1, loss1 = KMeans()(X1, X2)
visualize(X1, cluster_idx1, 2)

cluster_idx2, centers2, loss2 = KMeans()(X2, X2)
visualize(X2, cluster_idx2, 2)

/Users/vicentelaf@anacoda3/lib/python3.8/site-packages/sklearn/utils/deprecation.py:143: FutureWarning
ing: The sklearn.datasets.samples_generator module is deprecated in version 0.22 and will be removed
in version 0.24. The corresponding classes / functions should instead be imported from sklearn.datasets. Anything that cannot be imported from sklearn.datasets is now part of the private API.
warnings.warn(message, FutureWarning)
```

Visualization of K = 2

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:

$$\begin{aligned} \mathbf{x} &\sim \text{Bernoulli}(\theta) \\ \mathbf{z} & \sim \mathcal{N}(\mu, \sigma^2) \\ \mathbf{p}(\mathbf{x}|\mathbf{z}) &\sim \mathcal{N}(\mu, 3\sigma^2) \end{aligned}$$

For a dataset of N datapoints, find the following:

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

$$\mathbf{p}(\mathbf{x}) = \sum_{\mathbf{z}} \mathbf{p}(\mathbf{z}) \mathbf{p}(\mathbf{x}|\mathbf{z}) = \sum_{\mathbf{z}} \mathbf{p}(\mathbf{z}_k) \mathcal{N}(\mathbf{x}|\mu_k, \sigma_k^2)$$

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

$$\mathbf{p}(\mathbf{z}^i = k|\mathbf{x}^i) = \frac{\mathbf{p}(\mathbf{z}^i = k) \mathbf{p}(\mathbf{x}^i|\mathbf{z}^i = k)}{\sum_{\mathbf{z}^i} \mathbf{p}(\mathbf{z}^i = k) \mathbf{p}(\mathbf{x}^i|\mathbf{z}^i = k)} = \frac{\mathbf{p}(\mathbf{z}^i = k) \mathcal{N}(\mathbf{x}^i|\mu_k, \sigma_k^2)}{\sum_{\mathbf{z}^i} \mathbf{p}(\mathbf{z}^i = k) \mathcal{N}(\mathbf{x}^i|\mu_k, \sigma_k^2)}$$

for $k = 0$:

$$\mathbf{p}(\mathbf{z}^i = 0|\mathbf{x}^i) = \frac{\mathbf{p}(\mathbf{z}^i = 0) \mathcal{N}(\mathbf{x}^i|\mu_0, \sigma_0^2)}{\sum_{\mathbf{z}^i} \mathbf{p}(\mathbf{z}^i = k) \mathcal{N}(\mathbf{x}^i|\mu_k, \sigma_k^2)} = \frac{\theta \mathcal{N}(\mathbf{x}^i|\mu_0, \sigma_0^2)}{\theta \mathcal{N}(\mathbf{x}^i|\mu_0, \sigma_0^2) + (1-\theta) \mathcal{N}(\mathbf{x}^i|\mu_1, \sigma_1^2)}$$

for $k = 1$:

$$\mathbf{p}(\mathbf{z}^i = 1|\mathbf{x}^i) = \frac{\mathbf{p}(\mathbf{z}^i = 1) \mathcal{N}(\mathbf{x}^i|\mu_1, \sigma_1^2)}{\sum_{\mathbf{z}^i} \mathbf{p}(\mathbf{z}^i = k) \mathcal{N}(\mathbf{x}^i|\mu_k, \sigma_k^2)} = \frac{(1-\theta) \mathcal{N}(\mathbf{x}^i|\mu_1, \sigma_1^2)}{\theta \mathcal{N}(\mathbf{x}^i|\mu_0, \sigma_0^2) + (1-\theta) \mathcal{N}(\mathbf{x}^i|\mu_1, \sigma_1^2)}$$

2.1.3. M-Step: Compute the updated value of μ (You can keep σ fixed for this) [3pts]

$$l(\mathbf{x}|\theta) = \sum_{\mathbf{x}} \sum_{\mathbf{z}} \mathbf{p}(\mathbf{z}_k|\mathbf{x}_n, \theta_{old}) \ln [\mathbf{p}(\mathbf{z}_k) \mathcal{N}(\mathbf{x}_n|\mu_k, \sigma_k^2)] = \sum_{\mathbf{x}} \sum_{\mathbf{z}} \mathbf{p}(\mathbf{z}_k|\mathbf{x}_n, \theta_{old}) \ln [\mathbf{p}(\mathbf{z}_k) \mathcal{N}(\mathbf{x}_n|\mu_k, \sigma_k^2)]$$

$$l(\mathbf{x}|\theta) = \sum_{\mathbf{x}} \sum_{\mathbf{z}} \mathbf{p}(\mathbf{z}_k|\mathbf{x}_n, \theta_{old}) \ln [\mathbf{p}(\mathbf{z}_k)] - \frac{1}{2\sigma_k^2} \sum_{\mathbf{x}} \frac{(\mathbf{x}_n - \mu_k)^2}{\sigma_k^2}$$

$$l(\mathbf{x}|\theta) = \sum_{\mathbf{x}} \sum_{\mathbf{z}} \mathbf{p}(\mathbf{z}_k|\mathbf{x}_n, \theta_{old}) \ln [\mathbf{p}(\mathbf{z}_k)] + \ln \frac{1}{\sqrt{2\pi\sigma_k^2}} - \frac{(\mathbf{x}_n - \mu_k)^2}{2\sigma_k^2}$$

$$\frac{\partial l(\theta)}{\partial \mu_k} = \sum_{\mathbf{x}} \mathbf{p}(\mathbf{z}_k|\mathbf{x}_n, \theta_{old}) \left[\frac{(\mathbf{x}_n - \mu_k)}{\sigma_k^2} \right] = 0$$

$$\sum_{\mathbf{x}} \mathbf{p}(\mathbf{z}_k|\mathbf{x}_n, \theta_{old}) (\mathbf{x}_n - \mu_k) = 0$$

$$\sum_{\mathbf{x}} \mathbf{p}(\mathbf{z}_k|\mathbf{x}_n, \theta_{old}) \mathbf{x}_k^2 = \sum_{\mathbf{x}} \mathbf{p}(\mathbf{z}_k|\mathbf{x}_n, \theta_{old}) (\mathbf{x}_n - \mu_k)^2$$

$$\sigma_k^{new} = \sqrt{\frac{\sum_{\mathbf{x}} \mathbf{p}(\mathbf{z}_k|\mathbf{x}_n, \theta_{old}) (\mathbf{x}_n - \mu_k)^2}{\sum_{\mathbf{x}} \mathbf{p}(\mathbf{z}_k|\mathbf{x}_n, \theta_{old})}}$$

2.2 EM Algorithm in AB Blood Groups [10 pts]

In the AB 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
A	OA
B	BB
B	BO
B	OB
O	OO
AB	AB

In a research experiment, scientists wanted to model the distribution of the genotypes of the population. They collected the phenotype information from the participants as it 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 that 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 $n_{AA}, n_{AB}, n_{AO}, n_{BB}, n_{BO}, n_{OB}$ 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:

$$\begin{aligned}n_{AA} &= n_{AA} + n_{AO} \\n_{BB} &= n_{BB} + n_{BO} \\n_{AB} &= n_{AB} + n_{AO} + n_{BO} + n_{AB} = n\end{aligned}$$

Given:

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

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

$$p(A|A) = \frac{p(A,A)}{p(A)} = \frac{p(A)p(A)}{p(A,A)+2p(A,O)} = \frac{p(A)p(A)}{p(A)p(A)+2p(A)p(O)} = \frac{p(A)}{p(A)+2p(O)}$$

$$p(A|A) = \frac{p(A,A)}{p(A,A)+2p(A,O)} = \frac{\frac{1}{3}}{\frac{1}{3}+\frac{2}{3}} = \frac{1}{3}$$

$$n_{AA} = p(A|A)n_{AA} = \frac{1}{3}186 = 62$$

$$p(A|O) = \frac{p(A,O)}{p(A,A)+2p(A,O)} = \frac{\frac{1}{3}}{\frac{1}{3}+\frac{2}{3}} = \frac{1}{3}$$

$$n_{AO} = 2p(A|O)n_{AO} = \frac{2}{3}186 = 124$$

$$p(B|B) = \frac{p(B,B)}{p(B,B)+2p(B,O)} = \frac{\frac{1}{3}}{\frac{1}{3}+\frac{2}{3}} = \frac{1}{3}$$

$$n_{BB} = p(B|B)n_{BB} = \frac{1}{3}38 = 12.667 \approx 13$$

$$p(B|O) = \frac{p(B,O)}{p(B,B)+2p(B,O)} = \frac{\frac{1}{3}}{\frac{1}{3}+\frac{2}{3}} = \frac{1}{3}$$

$$n_{BO} = 2p(B|O)n_{BO} = \frac{2}{3}38 \approx 25.333 \approx 25$$

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

$$l(p|\theta) = \sum_k n_k \ln p_k$$

$$l(p|\theta) = n_{AA} \ln [p_{AA}] + n_{AO} \ln [p_{AO}] + n_{BB} \ln [p_{BB}] + n_{BO} \ln [p_{BO}] + n_{AB} \ln [p_{AB}] + n_{OO} \ln [p_{OO}]$$

$$l(p|\theta) = n_{AA} \ln [p_A] + n_{AO} \ln [2p_A p_B] + n_{BB} \ln [p_B] + n_{BO} \ln [2p_B p_O] + n_{AB} \ln [p_A p_B] + n_{OO} \ln [p_O]$$

$$\mathcal{L}(p, \lambda) = l(p|\theta) - \lambda(p_A + p_B + p_O - 1)$$

$$\frac{\partial \mathcal{L}(p, \lambda)}{\partial p_A} = \frac{2n_{AA} + n_{AO} + n_{AB}}{p_A} + \lambda = 0$$

$$\frac{\partial \mathcal{L}(p, \lambda)}{\partial p_O} = \frac{n_{OO} + n_{BO} + n_{AB}}{p_O} + \lambda = 0$$

$$\frac{\partial \mathcal{L}(p, \lambda)}{\partial \lambda} = p_A + p_B + p_O - 1 = 0$$

We can find λ by taking the sum of the first three partials, which that $p_A = p_B = p_O = \frac{1}{3}$ and

$$p_A = n_{AA} + n_{AO} + n_{AB} + n_{BB} + n_{BO} + n_{OO} = 186 + 124 + 13 + 12.667 + 25.333 + 284 = 665.667 \approx 666$$

$$2n = -\lambda$$

$$\lambda = -2n$$

To find p_A^{new} we solve for p_A in the partial with respect to p_A :

$$\frac{2n_{AA} + n_{AO} + n_{AB}}{p_A^{new}} - 2n = 0$$

$$p_A^{new} = \frac{2n_{AA} + n_{AO} + n_{AB}}{2n} = \frac{186 + 124 + 13}{1042} = 0.2505$$

To find p_B^{new} we solve for p_B in the partial with respect to p_B :

$$\frac{2n_{BB} + n_{BO} + n_{AB}}{p_B^{new}} - 2n = 0$$

$$p_B^{new} = \frac{2n_{BB} + n_{BO} + n_{AB}}{2n} = \frac{12.667 + 25.333 + 13}{1042} = 0.0614$$

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 Distribution. In a nutshell, GMM is a soft 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 π be a K -dimensional probability distribution and (μ_k, Σ_k) 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 π_k where $\sum_{k=1}^K \pi_k = 1$
- Each data point x_i is generated from Multivariate Normal Distribution $\mathcal{N}(\mu_k, \Sigma_k)$ where $\mu_k \in \mathbb{R}^D$ and $\Sigma_k \in \mathbb{R}^{D \times D}$

Our goal is to find K -dimension Gaussian distributions to model our data X . This can be done by learning the parameters π, μ and Σ through likelihood function. Detailed derivation can be found in the sum of $\ln \left(\sum_{k=1}^K \pi(k) \mathcal{N}(x_i | \mu_k, \Sigma_k) \right)$

From the lecture we know that MLEs for GMM all depend on each other and the responsibility r . 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 r , which is the conditional probability that a data point belongs to a specific cluster k if we are given the datapoint. $P(z_k | x)$. The formula for r is given below:

$$r(z_k) = \frac{\pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_i | \mu_j, \Sigma_j)}$$

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

- M-step: Re-estimate Parameters

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

$$\mu_k^{new} = \frac{\sum_{i=1}^N r(z_k) x_i}{N_k}$$

$$\Sigma_k^{new} = \frac{1}{N_k} \sum_{i=1}^N r(z_k) \bar{x}_i (\bar{x}_i - \mu_k^{new})^T (\bar{x}_i - \mu_k^{new})$$

$$\pi_k^{new} = \frac{N_k}{N}$$

where $N_k = \sum_{i=1}^N r(z_k)$. Note that the updated value for μ_k is used when updating Σ_k . The multiplication of $r(z_k)^T (x_i - \mu_k^{new})^T$ 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 undergrads/students: you may assume that the covariance matrix Σ_k is a diagonal matrix, which means the features are independent, i.e. the 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 (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.
- You can initiate $\pi(k)$ the same for each 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 your code, we recommend you take a look at the concept [Array 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 \(N, 1\) and \(N,\)](#)

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

3.1 Helper functions [15 pts]

To facilitate some of the operations in the GMM implementation, we will give 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** OVER 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 $\text{prob}_{ij} = \frac{\exp(\text{logit}_{ij})}{\sum_{k=1}^K \exp(\text{logit}_{ik})}$.

Note: It is possible that 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 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 $s_i = \log \left(\sum_{k=1}^K \exp(\text{logit}_{ik}) \right)$. 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 it 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) normalPDF

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

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

where σ_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

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

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

Hints

- If you encounter "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.linalg.inv(Sigma_k + 1e-32)`. You can arrest and handle such error by using [Try and Except/Block](#) in Python.
- In `np.linalg.solve`, 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 $(X - \mu)\Sigma^{-1}$, a (N, D) matrix, and 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 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 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 π the same for each cluster.
- Initialize μ 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 also can initialize the Σ 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:

$$\ell(\theta) = \sum_{i=1}^N \ln \left(\sum_{k=1}^K \pi(k) \mathcal{N}(x_i | \mu_k, \Sigma_k) \right)$$

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 $k = 1, \dots, K$,

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

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(np.eye(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 at `_ll_joint()` function. This is not the same as responsibilities (r), 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 Σ in M steps as you assumed independent components. Make sure you are only taking the diagonal terms of your calculated covariance matrix.

```
In [11]: 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]
    self.K = self.points.shape[1]
    self.mu = np.zeros((K, D))
    self.sigma = np.zeros((K, D, D))
    self.pi = np.ones((K,))

    # Helper function for you to implement
def softmax(self, logit): # [5pts]
    """
    Args:
        logit: N x D numpy array
    Returns:
        prob: N x D numpy array. See the above function.
    """
    # raise NotImplementedError
    max_logits = np.amax(logit, axis=1)
    max_logits = np.reshape(max_logits, (max_logits.size, 1))
    logit_norm = logit - max_logits
    logit_exp = np.exp(logit_norm)
    denominator = np.sum(logit_exp, axis=1)
    prob = logit_exp / np.reshape(denominator, (denominator.size, 1))
    return prob

def logsumexp(self, logit): # [5pts]
    """
    Args:
        logit: N x D numpy array
    Returns:
        s: N x 1 array where s[i,0] = logsumexp(logit[i,:]). See the above function
    """
    # raise NotImplementedError
    max_logits = np.amax(logit, axis=1)
    max_logits = np.reshape(max_logits, (max_logits.size, 1))
    logit_norm = logit - max_logits
    logit_exp = np.exp(logit_norm)
    logit_sum = np.sum(logit_exp, axis=1)
    s = np.log(logit_sum)
    s = s + max_logits
    return np.reshape(s, (s.size, 1))

# for undergraduate 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 1xD 2-D numpy array), the covariance matrix of the ith gaussian.
    Returns:
        pdf: 1xD numpy array (or array of length N), the probability distribution of N data for the ith gaussian
    """
    # np.diagonal() should be handy.
    # raise NotImplementedError
    N, D = np.shape(logit)
    if len(np.shape(sigma_i)) == 2:
        sigma_i = np.diagonal(sigma_i)
    else:
        sigma_i = np.diagonal(sigma_i[0])
    pdf = np.ones((1, N))
    for i in range(D):
        exponent = ((logit - sigma_i[i])**2) * np.square(logit[i, :]) - mu_i[i])
        pdf *= ((2 * np.pi * sigma_i[i])**-0.5) * np.exp(exponent)
    return np.reshape(pdf, (N,))

# 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 1xD 2-D numpy array), the covariance matrix of the ith gaussian.
    Returns:
        pdf: 1xD numpy array (or array of length N), the probability distribution of N data for the ith gaussian
    """
    # np.linalg.det() and np.linalg.inv() should be handy.
    # raise NotImplementedError

def _init_components(self, **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
    Returns:
        gamma: KxD numpy 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(x) defined above.
    """
    # raise NotImplementedError
    N, K = np.shape(gamma)
    tau_max = np.argmax(gamma, axis=1)
    pi = np.zeros(K)
    mu = np.zeros((K, D))
    sigma = np.zeros((K, D, D))
    for k in range(K):
        N_k = np.where(tau_max == k)
        X_k = self.points[N_k, :]
        tau_k = gamma[N_k, k]
        N_k = np.sum(tau_k, axis=0)
        mu_k = np.matmul(tau_k, self.points) / N_k
        mu[k] = mu_k
        pi[k] = N_k / N
        sigma_k = np.matmul(np.matmul(np.transpose(tau_k) * np.transpose(self.points - mu_k), (self.points - mu_k) / N_k)
        return pi, mu, sigma

def _call(self, self, tol=1e-6, rel_tol=1e-6, **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
    Returns:
        gamma(tau): NKX 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)
    pbar = tqdm(range(self.max_iters))
    for it in pbar:
        # E-step
        pi, mu, sigma = self._E_step(gamma)
        # M-step
        pi, mu, sigma = self._M_step(gamma)
        # calculate the negative log-likelihood of observation
        joint_ll = self._ll_joint(pi, mu, sigma)
        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*it} % (it, loss)')
    return gamma, mu, sigma
```

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

Ukiyo-e is a Japanese art genre predominant from the 17th through 19th century. In order to produce the intricate prints that came to represent the era, 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 appropriately 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 [12]: # 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
    Returns:
        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(float32)
    gamma, (pi, mu, sigma) = GMM(flat_img, K, max_iters = 100)()
    cluster_ids = 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 [13]: # pick 2 of the images in this list:
url0 = 'https://upload.wikimedia.org/wikipedia/commons/b/bl/Utagawa_Kunisada_4_k286-1832x29_Dawn_at_Fu-tami-gura.jpg'
url1 = 'https://upload.wikimedia.org/wikipedia/commons/9/95/Hokusai_k281828x29_Cuckoo_and_Azaleas.jpg'
url2 = 'https://upload.wikimedia.org/wikipedia/commons/7/74/Kitao_Shigemasa_k281777x29_Geisha_and_a_sax'
url3 = 'https://upload.wikimedia.org/wikipedia/commons/1/10/Runitoshi_Utagawa2C_Suikoden_Series_4.jpg'

# example of loading image from url0
image0 = imageio.imread(image0_core.urlopen(url1).read())
image3 = imageio.imread(image0_core.urlopen(url3).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 prod
    use 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: the minimum and maximum number of clusters you should test with. De
        fault are 5 and 15.
        (Usually the maximum number of clusters would not exceed 15)
    Returns:
        plot: comparison between original image and image pixel clustering.
        optional: any other information/metric/plot you think is necessary.
    """
    # raise NotImplementedError
    img_array = []
    for k in range(min_clusters, max_clusters+1):
        clustered_img = cluster_pixels_gmm(image, k)
        img_array.append(clustered_img)
    plot_images(img_array, ['image', 'k=5', 'k=6', 'k=7', 'k=8', 'k=9', 'k=10', 'k=11', 'k=12', 'k=13', 'k=14', 'k=15'])
    find_n_woodblocks(image)
    find_n_woodblocks(image3)
```

```
Iter 99, loss: 110329264.435: 100%
Iter 99, loss: 10043532.847: 100%
Iter 99, loss: 9729212.969: 100%
Iter 99, loss: 9628612.802: 100%
Iter 99, loss: 959466.594: 100%
Iter 99, loss: 9550773.437: 100%
Iter 99, loss: 9325898.917: 100%
Iter 99, loss: 9492650.511: 100%
Iter 99, loss: 9284258.090: 100%
Iter 99, loss: 9254103.152: 100%
Iter 99, loss: 9129672.664: 100%
Iter 99, loss: 1020359.412: 100%
Iter 99, loss: 10163177.072: 100%
Iter 99, loss: 1009649.772: 100%
Iter 99, loss: 1000978.613: 100%
Iter 99, loss: 9974826.263: 100%
Iter 99, loss: 9945855.219: 100%
Iter 99, loss: 9841324.084: 100%
Iter 99, loss: 9901072.97: 100%
Iter 99, loss: 985654.768: 100%
Iter 99, loss: 983562.793: 100%
Iter 99, loss: 9813929.591: 100%
```

By visually inspecting the reconstructed images using GMM with increasing clusters from 5 to 15, we can see how the prints would look like if only that number of color wood blocks was used to create them. In order to decide which was the number of wood blocks used for the original we must not inspect each reconstructed image with the original by increasing number of clusters until we reach the first one that only represents all colors in the original. For the first image selected this happens at $k=13$, meaning that for that print they probably used around 13 wood blocks. For the second one, it looks like they used 11 wood blocks (or 11 clusters in the case of the GMM algorithm).

(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


```
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, :]
"""
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: M_complete x (D+1) numpy array, the complete labeled observations
K: integer, corresponding to the number of nearest neighbors you want to base your calculat
ion on
kwargs: any other args you want
Return:
clean_points: (N_incomplete + M_complete) x (D+1) X D numpy array of length K, containing b
oth 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 addr
ess each feature at a time.
"""
raise NotImplementedError
```

Below is the 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([[1.,np.nan,3.,1],[7.,np.nan,9.,0],[np.nan,17.,18.,1],[np.nan,23.,24.,0]])
clean_data = CleanData().incomplete_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.]
 [14.5 23. 24.  0.]
 [ 7. 15.5 9.  0.]
 [ 8.5 17. 18.  1.]
 [ 1.  9.5 3.  1.]]""")
print("""\n*** My Answer - k = 2 ****""")
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. (2003). 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, the prior for each gaussian.
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: NxD 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): NxK array, where ll(i, j) = log pi(j) + log NormalPDF(points_i | mu[j],
sigma[i])
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: NxD 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:
gamma: 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 soft
max() defined above.
"""
raise NotImplementedError

def _M_step(self, points, gamma, **kwargs): # [0 pts] - can use same as for GMM
"""
Args:
points: NxD numpy array, the observations
gamma: 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.
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
odify from GMM
"""
Args:
points: NxD 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: NxK 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.

```
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 l
ast 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 point
s 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 l
ast 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 l
ast 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 = ~1*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()(points, 7)

# =====
# COMPARISON

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

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

# =====
print("""=====COMPARISON=====""")
print("""SemiSupervised Accuracy: """, ComparePerformance().accuracy_semi_supervised(classification, ind
ependent))
print("""Supervised with clean data: GNB Accuracy: """, ComparePerformance().accuracy_GNB_onlycomplete(l
abeled_complete, independent))
print("""Supervised with only complete data: GNB Accuracy: """, ComparePerformance().accuracy_GNB_cleand
ata(clean_data, independent))
```

```
In [ ]:
```