

Fall 2021 CS4641/CS7641 A Homework 2

Instructor: Dr. Mahdi Roozbahani

Deadline: October 19th, Tuesday, 11:59 pm AOE

- No unapproved extension of the deadline is allowed. Late submission will lead to 0 credit.
- Discussion is encouraged on Ed as part of the Q/A. However, all assignments should be done individually.
- Plagiarism is a **serious offense**. You are responsible for completing your own work. You are not allowed to copy and paste, or paraphrase, or submit materials created or published by others, as if you created the materials. All materials submitted must be your own.
- All incidents of suspected dishonesty, plagiarism, or violations of the Georgia Tech Honor Code will be subject to the institute's Academic Integrity procedures (e.g., reported to and directly handled by the Office of Student Integrity (OSI)). **Consequences can be severe, e.g., academic probation or dismissal, grade penalties, a 0 grade for assignments concerned, and prohibition from withdrawing from the class.**

Instructions for the assignment

- This assignment consists of both programming and theory questions.
- To switch between cell for code and for markdown, see the menu -> Cell -> Cell Type
- You can directly type Latex equations into markdown cells.
- If a question requires a picture, you could use this syntax
"`<imgsrc ="" style ="" width : 300px; />`" to include them within your ipython notebook.
- Your write up must be submitted in PDF form. You may use either Latex, markdown, or any word processing software. **We will **NOT** accept handwritten work.** Make sure that your work is formatted correctly, for example submit $\sum_{i=0} x_i$ instead of `\text{sum}_{i=0} x_i`
- When submitting the non-programming part of your assignment, you must correctly map pages of your PDF to each question/subquestion to reflect where they appear. Improperly mapped questions may not be graded correctly.
- Discussion is encouraged on Edstem as part of the Q/A. You may discuss high-level ideas with other students at the "whiteboard" level (e.g. how cross validation works, using matmul instead of dot) and review any relevant materials online. However, all assignments should be done individually, each student must write up and submit their own answers.
- **Graduate Students:** You are required to complete any sections marked as Bonus for Undergrads

Using the autograder

- You will find two assignments on Gradescope that correspond to HW2: "Assignment 2 Programming" and "Assignment 2 Programming - Bonus for all".
- You will submit your code for the autograder in the Assignment 2 Programming sections. Please refer to the Deliverables and Point Distribution section for what parts are considered required, bonus for undergrads, and bonus for all.
- 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 "Assignment 2 - Non-programming" part, you will download your Jupyter Notebook as html and submit it as a PDF on Gradescope. To download the notebook as html, click on "File" on the top left corner of this page and select "Download as > html". Then, open the html file and print to PDF.** Please refer to the Deliverables and Point Distribution section for an outline of the non-programming questions.
- **When submitting to Gradescope, please make sure to mark the page(s) corresponding to each problem/sub-problem.**

Deliverables and Points Distribution

Q1: KMeans Clustering & DBScan [65 pts]

Deliverables: **kmeans.py, dbscan.py, and Written Report**

- **pairwise_dist** [5 pts] - *programming*
- **KMeans Implementation** [35pts] - *programming*
 - `_init_centers` [5pts]
 - `_update_assignment` [10pts]
 - `_update_centers` [10pts]
 - `_get_loss` function [5pts]
 - `_test_centers` (Additional Autograder Tests - no need to implement) [5pts]
- **Elbow Method** [15 pts]
 - `find_num_optimal_clusters` [10pts] - *programming*

- written questions [5pts] - *non-programming*
- **DBScan** [10 pts] - *programming*
 - regionQuery [2pts]
 - expandClusters [4pts]
 - fit [4pts]

Q2: EM Algorithm [10pts; 20pts Bonus for Undergrad]

Deliverables: **Written Report**

- **2.1 Performing EM Update** [10 pts] - *non-programming*
 - 2.1.1 [2pts] - *non-programming*
 - 2.1.2 [2pts] - *non-programming*
 - 2.1.3 [6pts] - *non-programming*
- **2.2 EM Algorithm in Coin Toss Problem** [20 pts Undergrad Bonus] - *non-programming*
 - 2.2.1 [10pts] - *non-programming*
 - 2.2.2 [10pts] - *non-programming*

Q3: GMM implementation [55pts; 5pts Bonus for All]

Deliverables: **gmm.py and Written Report**

- 3.1 Helper Functions [15pts] - *programming & non-programming*
 - 3.1.1. softmax [5pts]
 - 3.1.2. logsumexp [3pts + 2pts] - *programming & non-programming*
 - 3.1.3. normalPDF [5pts] - *for CS4641 students only*
 - 3.1.3. multinormalPDF [5pts] - *for CS7641 students only*
- 3.2 GMM Implementation [30pts] - *programming*
 - 3.2.1. init_components [5pts]
 - 3.2.2. _ll_joint [10pts]
 - 3.2.3. Setup iterative steps for EM algorithm [15pts]
- 3.3 Image Compression and Pixel clustering [10pts] - *non-programming*
- 3.4 Compare Full Covariance Matrix with Diagonal Covariance Matrix [5pts Bonus for All]

Q4: Cleaning Super Duper Messy data with semi-supervised learning [30pts Bonus for All]

Deliverables: [semisupervised.py](#) and Written Report

- 4.1: KNN [10pts] - *programming*
- 4.2: Getting acquainted with semi-supervised learning approaches [5pts] - *non-programming*
- 4.3: Implementing the EM algorithm [10pts] - *programming*
- 4.4: *Demonstrating the performance of the algorithm* [5pts] - *programming*

0 Set up

This notebook is tested under [python 3.*.*](#), 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 [3]: #####  
### DO NOT CHANGE THIS CELL ###  
#####  
  
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 mpl_toolkits.mplot3d import axes3d  
from tqdm import tqdm  
  
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 image  
import imageio
```

```
%load_ext autoreload
%autoreload 2

Version information
python: 3.8.5 (default, Sep 3 2020, 21:29:08) [MSC v.1916 64 bit (AMD64)]
matplotlib: 3.3.2
numpy: 1.19.2
```

1. Rob's Balloon Commission - KMeans Clustering [55 pts + 10 pts]

Rob Boss, after being stuck in his dorm for multiple nights finishing Homework 1 for his Machine Learning class, decides to get some well deserved fresh air. He plans on visiting the hot air balloon festival downtown this weekend with a few friends. Once there, he takes many pictures and ultimately decides that his favorite is the picture attached below (please compliment Rob on his photography and machine learning skills the next time you see him) and wants to hang it up in his dorm. He decides to commission his artist friend to convert his picture into a painting but finds out his friend charges him for each color used when painting. Being the cash-strapped college student he is, help Rob build a KMeans clustering algorithm that can help him decide how many colors he wants to commission for his painting so that he can see the pattern of the hot air balloons without spending too much money.



KMeans is trying to solve the following optimization problem:

$$\arg \min_S \sum_{i=1}^K \sum_{x_j \in S_i} \|x_j - \mu_i\|^2 \quad (1)$$

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

1.1 pairwise distance [5pts]

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

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

DO NOT USE FOR LOOPS in your implementation -- they are slow and will make your code too slow to pass our grader. [Use array broadcasting instead.](#)

Test Case:

```
In [4]: #####
### DO NOT CHANGE THIS CELL ####
#####

from kmeans import pairwise_dist

x = np.random.randn(2, 2)
y = np.random.randn(3, 2)
```

```

print("/** Expected Answer **")
print("==x=="
[[ 1.62434536 -0.61175641]
 [-0.52817175 -1.07296862]]
==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("\n** My Answer **")
print("==x==")
print(x)
print("==y==")
print(y)
print("==dist==")
print(pairwise_dist(x, y))

```

```

*** Expected Answer ***
==x==
[[ 1.62434536 -0.61175641]
 [-0.52817175 -1.07296862]]
==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==
[[ 1.62434536 -0.61175641]
 [-0.52817175 -1.07296862]]
==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]]

```

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.

```

In [5]: #####
### NO NEED TO CHANGE THIS CELL ###
#####

from kmeans import KMeans

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[float]]]
or (grayscale) np.ndarray[np.ndarray[float]]
"""
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

def update_image_values(k):
    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]

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

def plot_image(img_list, title_list, figsize=(9, 12)):
    fig, axes = plt.subplots(1, len(img_list), figsize=figsize)
    for i, ax in enumerate(axes):
        ax.imshow(img_list[i])
        ax.set_title(title_list[i])
        ax.axis('off')

```

In [159...]

```

#####
### NO NEED TO CHANGE THIS CELL #####
#####

#Note that because of a different file structure, students' paths will be different

image_values = image_to_matrix('./images/balloon.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)

print('Loading...')

```

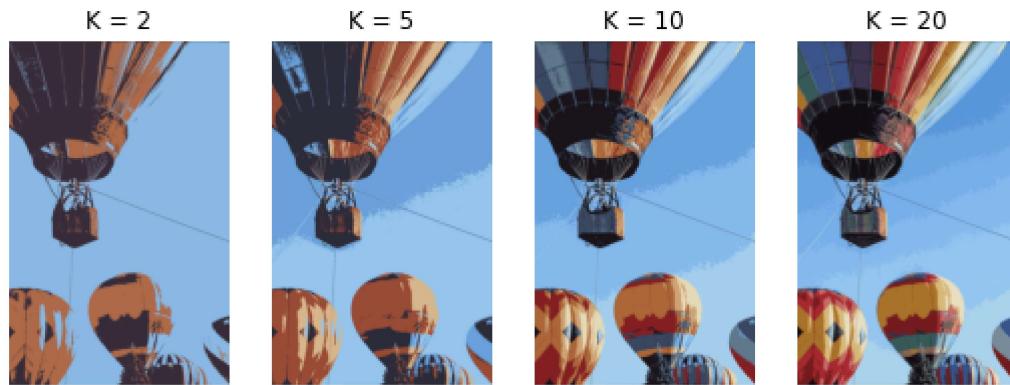
```

image_2 = update_image_values(2).reshape(r, c, ch)
image_2 = update_image_values(3).reshape(r, c, ch)
image_5 = update_image_values(5).reshape(r, c, ch)
image_10 = update_image_values(10).reshape(r, c, ch)
image_20 = update_image_values(20).reshape(r, c, ch)

plot_image([image_2, image_5, image_10, image_20], ['K = 2', 'K = 5', 'K = 10', 'K = 20'])

```

Loading...

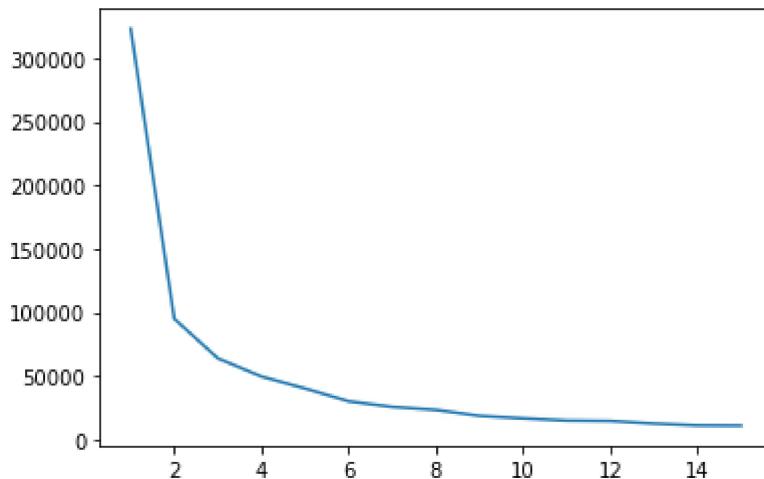


1.3 Elbow Method [6 pts Programming + 4 pts Written Questions]

One of the biggest drawbacks of KMeans is that we need to know the number of clusters beforehand. Let's see if we can help Rob's paint optimization problem by implementing the elbow method to find the optimal number of clusters in the function `find_optimal_num_clusters` below.

```
In [62]: #####
### NO NEED TO CHANGE THIS CELL ###
#####

from kmeans import find_optimal_num_clusters
find_optimal_num_clusters(image_values)
```



Written Questions:

- Approximately what value does the elbow method give? Was the elbow method roughly accurate when compared to the images displayed?

2) Rob is not quite sure to interpret his loss graph using the elbow method, so he decides to optimize the number of clusters by choosing k to have the loss be as close to 0 as possible. Would this idea be a good solution for making sure he gets the most bang for his buck when it comes to displaying balloon patterns? Why or why not?

Solutions:

1) Elbow value is between 2 and 3. I think the elbow is pretty accurate compared to the images displayed. I added in the case for running K means for K = 3, and I think that gave me the best results.

2) To have the lowest possible loss, each point shall have basically its own cluster, which defeats the point of clustering and further will lead to Rob running out of his budget. For effective clustering to occur there must be some loss.

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 [8]: #####
### NO NEED TO CHANGE THIS CELL #####
#####

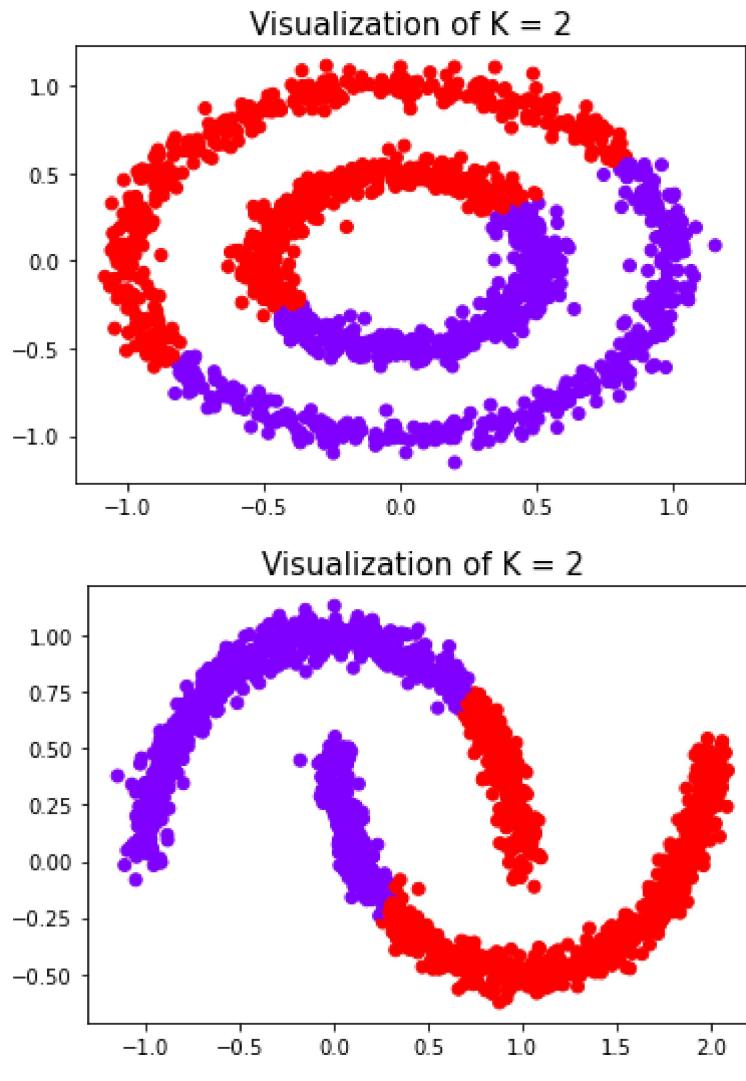
# visualize limitation of kmeans
from sklearn.datasets import make_circles, make_moons

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

def visualise(X, C, K=None):# Visualization of clustering. You don't need to change this
    fig, ax = plt.subplots()
    ax.scatter(X[:, 0], X[:, 1], c=C, cmap='rainbow')
    if K:
        plt.title('Visualization of K = '+str(K), fontsize=15)
    plt.show()
    pass

cluster_idx1, centers1, loss1 = KMeans()(X1, 2)
visualise(X1, cluster_idx1, 2)

cluster_idx2, centers2, loss2 = KMeans()(X2, 2)
visualise(X2, cluster_idx2, 2)
```



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

To obtain these 5 points, you need to 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.

1.6 DBSCAN [10 pts]

Let us try to solve these limitations using another clustering algorithm: DBSCAN. As mentioned in lecture, DBSCAN tries to find dense regions in the data space, separated by regions of lower density. DBSCAN is parameterized by two parameters (ϵ and minPts):

- ϵ : Maximum radius of neighborhood
- MinPts : Maximum number of points in ϵ -neighborhood of a point

The algorithm fits on a dataset as follows:

For each unvisited point in dataset:

1. Set point as visited
2. Extract the ϵ -neighborhood *Neighbors* of this point using `regionQuery()` (all points which are within ϵ distance of this neighborhood)

3. If the size of this neighborhood is greater than minPts , then we run `expandCluster()`, which does the following for cluster C :

- A. Set point P as visited
- B. Add P to cluster C
- C. For each point P' in P 's neighborhood Neighbors :
 - a. Set point P' as visited
 - b. Extract the Eps-neighborhood $\text{Neighbors}'$ of this point using `regionQuery()`
 - c. If the size of $\text{Neighbors}'$ is greater than minPts , we add $\text{Neighbors}'$ to our current neighborhood Neighbors :
 - d. If P' has not been assigned a cluster yet, we add it to cluster C
- D. Move onto next cluster $C + 1$

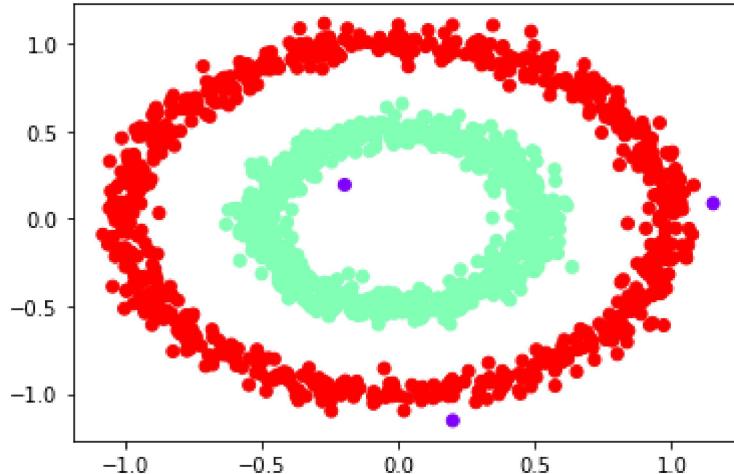
Using the above description, complete `fit()`, `expandCluster()`, and `regionQuery()` in `dbSCAN.py`.

Then, test your code by running the cell below. You should be able to get a perfect clustering for the two circles dataset, which you can observe quantitatively by checking whether the clusters returned by `cluster_idx` and the ground truth clusters are the same and qualitatively by visualizing the clusters.

```
In [160]: #####
### DO NOT CHANGE THIS CELL #####
#####

BEST_EPS = 0.11
BEST_POINTS = 3
# from dbSCAN import DBSCAN
dbSCAN = DBSCAN(BEST_EPS, BEST_POINTS, X1)
cluster_idx = dbSCAN.fit()
## Note that one of the two cells should print True for a correct implementation
# print(y1.tolist())
# print(cluster_idx.tolist())
# print(cluster_idx[760])
print(np.array_equal(y1, cluster_idx)) #Checks if y1 == cluster_idx
print(np.array_equal(y1, 1-cluster_idx)) ## Checks if y1 is the exact opposite of cluster_idx
visualise(X1, cluster_idx)
```

False
False



2. EM algorithm [30 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:

$$\mathbf{z} \sim Bernoulli(\theta)$$

$$\mathbf{p}(\mathbf{x}|\mathbf{z} = \mathbf{0}) \sim \mathcal{N}(2\mu, 5\sigma^2)$$

$$\mathbf{p}(\mathbf{x}|\mathbf{z} = \mathbf{1}) \sim \mathcal{N}(\mu, \sigma^2)$$

For a dataset of N datapoints, find the following:

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

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

2.1.3. M-Step: Compute the updated value of σ^2 (You can keep μ fixed for this) [6pts]

Answers:

2.1.1

$P(X)$ can be written as a sum of a partition of the sample, using the law of total probability,

$$p(x) = p(x|z = 0)(\theta) + p(x|z = 1)(1 - \theta) = N(2\mu, 5\sigma^2)(\theta) + N(\mu, \sigma^2)(1 - \theta)$$

2.1.2

$$p(z^i = 0|x^i) = \frac{p(x^i|z^i = 0)p(z^i = 0)}{p(x^i)} = \frac{(N(2\mu, 5\sigma^2)(\theta))^i}{(N(2\mu, 5\sigma^2)(\theta) + N(\mu, \sigma^2)(1 - \theta))^i}$$

$$p(z^i = 1|x^i) = \frac{p(x^i|z^i = 1)p(z^i = 1)}{p(x^i)} = \frac{(N(\mu, \sigma^2)(1 - \theta))^i}{(N(2\mu, 5\sigma^2)(\theta) + N(\mu, \sigma^2)(1 - \theta))^i}$$

2.1.3

$$\sigma_k^2 = \frac{1}{N_k} \sum_{n=1}^N \tau(z_{nk})(x_n - \mu_k)(x_n - \mu_k)^\tau$$

$$\tau_k = \frac{\pi_k N(x_n | \mu_k, \sigma_k^2)}{N(2\mu, 5\sigma^2)\theta + N(\mu, \sigma^2)(1-\theta)}$$

$$\mu_k = \sum_{n=1}^N \frac{\tau(z_{nk})x_n}{N_k}$$

$$\sigma_0^2 = \frac{1}{N_0} \sum_{n=1}^N \tau(z_{n0})(x_n - \mu)(x_n - \mu)^\tau$$

$$= \frac{1}{N_0} \sum_{n=1}^N \frac{\pi_0 N(x_n | \mu, \sigma_0^2)(x_n - \mu)(x_n - \mu)^\tau}{N(2\mu, 5\sigma^2)\theta + N(\mu, \sigma^2)(1-\theta)}$$

$$= \frac{1}{N_0} \sum_{n=1}^N \frac{N(2\mu, 5\sigma^2)\theta (x_n - \mu)(x_n - \mu)^\tau}{N(2\mu, 5\sigma^2)\theta + N(\mu, \sigma^2)(1-\theta)}$$

$$= \frac{1}{\sum_{n=1}^N p(z_n = 0 | x_n)} \sum_{n=1}^N p(z_n = 0 | x_n)(x_n - \mu)(x_n - \mu)^\tau$$

$$\sigma_1^2 = \frac{1}{N_1} \sum_{n=1}^N \tau(z_{n1})(x_n - \mu)(x_n - \mu)^\tau$$

$$= \frac{1}{N_1} \sum_{n=1}^N \frac{\pi_1 N(x_n | \mu, \sigma_1^2)(x_n - \mu)(x_n - \mu)^\tau}{N(2\mu, 5\sigma^2)\theta + N(\mu, \sigma^2)(1-\theta)}$$

$$= \frac{1}{N_1} \sum_{n=1}^N \frac{N(\mu, \sigma^2)(1-\theta) (x_n - \mu)(x_n - \mu)^\tau}{N(2\mu, 5\sigma^2)\theta + N(\mu, \sigma^2)(1-\theta)}$$

$$= \frac{1}{\sum_{n=1}^N p(z_n = 1 | x_n)} \sum_{n=1}^N p(z_n = 1 | x_n)(x_n - \mu)(x_n - \mu)^\tau$$

2.2 EM Algorithm in Coin Toss Problem [20 pts Undergrad Bonus]

Suppose we have 2 coins, C_1 and C_2 . We have the following table that shows information related to 6 trials along with their head and tail results. In each trial, one of the coins is flipped exactly 8 times, meaning every row of the table represents a single trial performed by a single coin. For example, let's say we knew that the first trial with results TTHTHHHT was performed by C_1 . Then we know that for that trial, C_1 obtained 4 heads in total and C_2 obtained 0 heads in total since it was not the one flipped on this specific trial.

<i>Coin</i>	<i>FlipResults</i>	<i>Coin1Heads</i>	<i>Coin2Heads</i>
	<i>TTHTHHHT</i>		
	<i>HHTTHHHT</i>		
	<i>TTTTTHHT</i>		
	<i>THTHTTHT</i>		
	<i>HHTHHTHH</i>		
	<i>THHTTTTT</i>		

As you can see, we know the results of the coin flips for each trial, but we don't know which coin was flipped on each trial. (Example: we know that the first trial resulted in TTHTHHHT, but we don't know which coin out of the two coins available was used to produce such results). Therefore, we also can't count how many heads a coin produced in a trial since we don't know which coin was flipped on such trial.

Using the information we have available, we want to take advantage of the EM Algorithm to figure out the following:

- Make an initial guess of the biases of the 2 coins to estimate the order of coins for the 6 trials performed
- Estimate the probability that each coin could have been thrown at a certain trial
- Recalculate to make our biases guess a better one in the end of the EM process

To start, it's important to note that both coins have equal chance of being picked for each trial:

- $P(C_1) = 0.5$
- $P(C_2) = 0.5$

Regarding the biases of the coins, we will start by assuming that they have the following:

- $C_1 = 0.4$ heads (and 0.6 tails)
- $C_2 = 0.7$ heads (and 0.3 tails)

2.2.1 You will fill in the table below, calculating the probability that C_x was picked for each trial, as well as the amount of heads flipped by C_x for each trial. (You can round your answers to 3 decimal places). Some entries in the table are already filled as hints to your calculations. To be specific, each column of the table represents: [10pts]

- $P(C_1|T)$ is the probability that coin 1 (C_1) was picked for the trial
- $P(C_2|T)$ is the probability that coin 2 (C_2) was picked for the trial

- Coin 1 Heads is the amount of heads flipped in a certain trial times the probability that Coin 1 was picked for the trial
- Coin 2 Heads is the amount of heads flipped in a certain trial times the probability that Coin 2 was picked for the trial

(Hint: you will find Bayes' Theorem useful for this question

$$P(A|B) = \frac{P(B|A)*P(A)}{P(B)} = \frac{P(B|A)*P(A)}{P(B|A)*P(A) + P(B|-A)*P(-A)}$$

(Hint: $P(T_1|C_1) = P(TTHTHHHT|\text{picked } C_1) = \frac{8!}{4!4!} * 0.4^4 * 0.6^4$)

FlipResults	$P(C1 T)$	$P(C2 T)$	Coin1Heads	Coin2Heads
TTHTHHHT	0.630			
HHTTHHHT				3.360
TTTTTHHT			1.908	
THTHTTHT		0.878		
HHTHHTHH				
THTTTTTT				

2.2.2 Now you want to apply the M-Step of the algorithm to improve our first estimate on the coin biases, calculating θ^{new_1} and θ^{new_2} . You would be maximizing the likelihood of the estimated count of flips we previously calculated in the E-Step [10pts]

(Hint: $\theta_{new_x} = \frac{\sum c_x \text{heads}}{\text{number of flips in a trial} * \sum_{i=1}^n P(C_x|T_i)}$, where $n = \text{number of trials}$)

Answer:

2.2.1



FlipResults	$P(C1 T)$	$P(C2 T)$	Coin1Heads	Coin2Heads
TTHTHHHT	0.630	0.370	2.52	1.48
HHTTHHHT	0.328	0.672	1.64	3.360
TTTTTHHT	0.954	0.046	1.908	0.092
THTHTTHT	0.857	0.143	2.571	0.429
HHTHHTHH	0.122	0.878	0.732	5.268
THTTTTTT	0.987	0.013	0.987	0.013



2.2.2

$$\theta_{new_1} = \frac{2.52 + 1.64 + 1.908 + 2.571 + 0.732 + 0.987}{8 * (0.63 + 0.328 + 0.954 + 0.857 + 0.122 + 0.987)} = 0.333871$$

$$\theta_{new_2} = \frac{1.48 + 3.360 + 0.092 + 0.429 + 5.268 + 0.013}{8 * (0.37 + 0.672 + 0.046 + 0.143 + 0.878 + 0.013)} = 0.627$$

3. GMM implementation [43pts Programming + 12pts Written Questions + 5pts Bonus for all 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 a 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 our slide of GMM. The log-likelihood function now becomes:

$$\ln p(x_1, \dots, x_N | \pi, \mu, \Sigma) = \sum_{i=1}^N \ln \left(\sum_{k=1}^K \pi(k) \mathcal{N}(x_i | \mu_k, \Sigma_k) \right) \quad (2)$$

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

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

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 τ in this homework.**

- **M-step:** Re-estimate Parameters

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

$$\mu_k^{new} = \frac{\sum_{n=1}^N \tau(z_k) x_n}{N_k} \quad (3)$$

$$\Sigma_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \tau(z_k)^T (x_n - \mu_k^{new})^T (x_n - \mu_k^{new}) \quad (4)$$

$$\pi_k^{new} = \frac{N_k}{N} \quad (5)$$

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

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

Special Notes

- For undergraduate student: you may assume that the covariance matrix Σ is diagonal matrix, which means the features are independent. (i.e. the red intensity of a pixel is independent from its blue intensity, etc). Make sure you set **Full_matrix = False** before you submit your code to Gradscope.
- For graduate student: please assume full covariance matrix. Make sure you set **Full_matrix = True** before you submit your code to Gradscope
- The class notes assume that your dataset X is (D, N) but **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

1. **DO NOT USE FOR LOOPS OVER N.** You can always find a way to avoid looping over the observation datapoints in our homework problem. If you have to loop over D or K, that would be fine.
2. You can initiate $\pi(k)$ the same for each k , i.e. $\pi(k) = \frac{1}{K}, \forall k = 1, 2, \dots, K$.
3. 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 corresponds 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.
4. To avoid using for loops in your code, we recommend you take a look at the concept [Array Broadcasting in Numpy](#). Also, certain calculations that required different shapes of arrays can

also be achieved by broadcasting.

5. 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 and how to format your return statement. Print the shape of an array by `print(array.shape)` could enhance the functionality of your code and help you debugging. 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 $X: (N, D)$
- $\mu: (K, D)$.
- $\Sigma: (K, D, D)$
- $\tau: (N, K)$
- $\pi: \text{array of length } K$
- $\Pi_{\text{joint}}: (N, K)$

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 $\text{prob}_{i,j} = \frac{\exp(\text{logit}_{i,j})}{\sum_{d=1}^D \exp(\text{logit}_{i,d})}$.

Note: it is possible that $\text{logit}_{i,j}$ 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 *logits*.

Special Notes

- Do not add back the maximum for each row.
- Add **keepdims=True** in your `np.sum()` function to avoid broadcast error.

3.1.2. logsumexp [3 pts Programming + 2 pts Written Question]

Given $\text{logit} \in \mathbb{R}^{N \times D}$, calculate $s \in \mathbb{R}^N$, where $s_i = \log \left(\sum_{j=1}^D \exp(\text{logit}_{i,j}) \right)$. Again, pay attention to the numerical problem. You may face similar condition as in the softmax function. In this case, add the maximum for each row of *logit* back for your functions 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.

Written Questions [2 pts]:

1) Why **logsumexp()** function add the maximum for each row of *logit* back?

Hint: start with a simple example like $\text{logit} \in \mathbb{R}^{1 \times D}$

Answer:

Let x_1, x_2, \dots, x_d be the elements in a row and let m be the max of all those elements. Then,

$$\begin{aligned}s_i &= \log(e^{x_1-m} + \dots + e^{x_d-m}) = \log\left(\frac{1}{e^m}(e^{x_1} + \dots + e^{x_d})\right) = \log(e^{x_1} + \dots + e^{x_d}) - \log(e^m) \\ &= \log(e^{x_1} + \dots + e^{x_d}) - m = s_i - m\end{aligned}$$

Therefore by re adding the max value, m , we get the result we desire.

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 everywhere else 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

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

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_k + SIGMA_CONST)`. 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 $(X - \mu)\Sigma^{-1}$, a (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 each individual μ , 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 to other array parameters. Both ways should be acceptable but pay attention to the shape mismatch problem and be **consistent all the time** when you implement such arrays.

3.2 GMM Implementation [30 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.

1. Set the prior probability π the same for each class.
2. Initialize μ by randomly selecting K numbers of observations as the initial mean vectors. You can use `int(np.random.uniform())` to get the row index number of the datapoints randomly.
3. Initialize the covariance matrix with `np.eye()` for each k . For grads, you can also initialize the Σ by K diagonal matrices. It will become a full matrix after one iteration, as long as you adopt the correct computation.
4. Other ways of initialization are acceptable and welcome. The autograder will only test the shape of your π, μ, σ . Make sure you pass other evaluations in the autograder.

3.2.2. Formulate the log-likelihood function `_ll_joint()` [10 pts]

The log-likelihood function is given by:

$$\ell(\theta) = \sum_{i=1}^N \ln \left(\sum_{k=1}^K \pi(k) \mathcal{N}(x_i | \mu_k, \Sigma_k) \right) \quad (6)$$

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 $k = 1, \dots, K$,

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

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(pi_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 (τ), but you should be able to finish this part with just a few lines of code by using `_ll_joint()` and `softmax()` 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.

Function Tests

Use these to test if your implementation of functions in GMM work as expected

```
In [63]: #####  
### DO NOT CHANGE THIS CELL ###  
#####  
  
from gmm import GMM
```

```
In [64]: #####  
## NO NEED TO CHANGE THIS CELL ##  
#####  
  
np.random.seed(1)  
  
data = np.random.randn(4, 3)  
  
# test softmax utility  
my_softmax = GMM(data, 3).softmax(data)  
expected_softmax = np.array([[0.81761761, 0.08738232, 0.09500007],  
    [0.12135669, 0.84312089, 0.03552242],  
    [0.75647821, 0.0617229, 0.18179889],  
    [0.14923883, 0.82635643, 0.02440474]])  
  
print("Your softmax works within the expected range: ", np.allclose(expected_softmax, my_softmax))  
  
# test logsumexp utility  
my_logsumexp = GMM(data, 3).logsumexp(data)  
expected_logsumexp = np.array([[1.82570589],  
    [1.03605256],  
    [2.02389331],  
    [1.65283702]])  
  
print("Your logsumexp works within the expected range: ", np.allclose(expected_logsumexp, my_logsumexp))
```

```
# init random
points = np.random.randn(12, 3)

mu = np.array([[-0.74715829,  1.6924546,   0.05080775],
               [-1.09989127, -0.17242821, -0.87785842],
               [-0.3224172,  -0.38405435,  1.13376944]])

sigma = np.array([[1., 0., 0.],
                  [0., 1., 0.],
                  [0., 0., 1.]],
                 [[1., 0., 0.],
                  [0., 1., 0.],
                  [0., 0., 1.]],
                 [[1., 0., 0.],
                  [0., 1., 0.],
                  [0., 0., 1.]]])
pi = np.ones(3)/3
```

Your softmax works within the expected range: True
 Your logsumexp works within the expected range: True

In [143...]

```
#####
### NO NEED TO CHANGE THIS CELL ###
#####

# For undergrads

# test normalPDF
my_normalpdf = GMM(points, 3).normalPDF(points, mu[0], sigma[0])
expected_normal_pdf = np.array([0.0037374, 0.00681159, 0.01294674, 0.00700474, 0.00095
                                0.00813925, 0.00544499, 0.00385966, 0.00561288, 0.00228524,
                                0.06349364, 0.00250289])

print("Your normal pdf works within the expected range: ", np.allclose(expected_normal_)

# test LL-joint
my_lljoint = GMM(points, 3).ll_joint(pi, mu, sigma, False)
expected_lljoint = np.array([[-6.68797812, -6.20337699, -3.85542789],
                            [-6.08774253, -3.85542789, -6.20337698],
                            [-5.44552376, -4.81763731, -6.88557037],
                            [-6.05977986, -7.90402129, -5.95737486],
                            [-8.05161039, -6.27262476, -5.43812535],
                            [-5.90966969, -4.86498535, -4.23232988],
                            [-6.31167075, -3.98209541, -5.58159406],
                            [-6.65578947, -4.38655011, -4.69047683],
                            [-6.28130441, -7.15820124, -4.50096327],
                            [-7.17989628, -5.55202701, -6.48346667],
                            [-3.85542789, -6.08774255, -6.6879781 ],
                            [-7.08892294, -8.46315357, -4.53725014]])

print("Your lljoint works within the expected range: ", np.allclose(my_lljoint, expecte

# test E step
my_estep = GMM(points, 3).E_step(pi, mu, sigma)
expected_estep = np.array([[0.05098852, 0.08278125, 0.86623023],
                           [0.08918842, 0.83136246, 0.07944912],
                           [0.3214852 , 0.60234958, 0.07616522],
                           [0.44131069, 0.06979118, 0.48889813],
                           [0.04861361, 0.28797946, 0.66340693],
                           [0.10876892, 0.30917578, 0.5820553 ],
```

```
[0.074913 , 0.76962456, 0.15546244],
[0.0561508 , 0.54309286, 0.40075634],
[0.13609235, 0.05662422, 0.80728343],
[0.12346309, 0.62879889, 0.24773802],
[0.85752822, 0.09199548, 0.0504763 ],
[0.07101476, 0.01796916, 0.91101608]])

print("Your E step works within the expected range: ", np.allclose(my_estep, expected_e))

# test M step
my_pi, my_mu, my_sigma = GMM(points, 3).M_step(expected_estep, False)
expected_pi = np.array([0.19829313, 0.35762874, 0.44407813])
expected_mu = np.array([[[-0.20989007, 0.79579186, 0.06554929],
   [-0.35741548, -0.1535599 , -0.4876455 ],
   [-0.28772515, -0.07512445, 0.79292111]]])
expected_sigma = np.array([[[[0.64857055, 0.          , 0.          ],
   [0.          , 0.63446774, 0.          ],
   [0.          , 0.          , 0.62167826]],

   [[0.53473119, 0.          , 0.          ],
   [0.          , 0.23538075, 0.          ],
   [0.          , 0.          , 0.38671205]],

   [[0.62612107, 0.          , 0.          ],
   [0.          , 0.24611766, 0.          ],
   [0.          , 0.          , 0.88668642]]]])
print("Your M step works within the expected range: ", np.allclose(my_pi, expected_pi))
```

Your normal pdf works within the expected range: True
 Your lljoint works within the expected range: True
 Your E step works within the expected range: True
 Your M step works within the expected range: False

In [153...]

```
#####
### NO NEED TO CHANGE THIS CELL #####
#####

# For grads

## # test multinormalPDF
sigma_grad = np.array([[[[ 0.12015895, 0.61720311, 0.30017032],
   [-0.35224985, -1.1425182 , -0.34934272],
   [-0.20889423, 0.58662319, 0.83898341]],

   [[ 0.93110208, 0.28558733, 0.88514116],
   [-0.75439794, 1.25286816, 0.51292982],
   [-0.29809284, 0.48851815, -0.07557171]],

   [[ 1.13162939, 1.51981682, 2.18557541],
   [-1.39649634, -1.44411381, -0.50446586],
   [ 0.16003707, 0.87616892, 0.31563495]]]])
my_multinormalpdf = GMM(data, 3).multinormalPDF(points, mu[0], sigma_grad[0])
expected_multinormal_pdf = np.array([8.70516304e-074, 8.62201632e-001, 5.36048920e+015,
   6.91708798e+083, 9.96882978e-062, 7.03348279e-025, 2.16083146e-059,
   1.87537738e-086, 1.84295981e+075, 1.11845126e+000, 5.17746613e-097])

print("Your multinormal pdf works within the expected range: ", np.allclose(expected_mu

# test ll-joint
```

```

sigma_now = sigma * 0.5
my_lljoint = GMM(points, 3)._ll_joint(pi, mu, sigma_now, True)
expected_lljoint = np.array([[ -8.48080757, -7.51160532, -2.81570712],
                           [ -7.28033641, -2.81570712, -7.51160531],
                           [ -5.99589887, -4.74012597, -8.87599209],
                           [ -7.22441107, -10.91289393, -7.01960107],
                           [ -11.20807212, -7.65010086, -5.98110204],
                           [ -6.92419072, -4.83482203, -3.56951111],
                           [ -7.72819284, -3.06904217, -6.26803946],
                           [ -8.41643028, -3.87795155, -4.485805 ],
                           [ -7.66746017, -9.42125381, -4.10677788],
                           [ -9.4646439 , -6.20890536, -8.07178468],
                           [ -2.81570712, -7.28033643, -8.48080755],
                           [ -9.28269723, -12.03115847, -4.17935163]]))

print("Your lljoint works within the expected range: ", np.allclose(my_lljoint, expected_lljoint))

# test E step
my_estep = GMM(points, 3)._E_step(pi, mu, sigma_now)
expected_estep = np.array([[3.42169503e-03, 9.01904364e-03, 9.87559261e-01],
                           [1.12762023e-02, 9.79775837e-01, 8.94796041e-03],
                           [2.18977456e-01, 7.68731442e-01, 1.22911017e-02],
                           [4.43990237e-01, 1.11041589e-02, 5.44905604e-01],
                           [4.49802848e-03, 1.57844511e-01, 8.37657460e-01],
                           [2.65137773e-02, 2.14226353e-01, 7.59259870e-01],
                           [9.02095401e-03, 9.52129225e-01, 3.88498209e-02],
                           [6.87345697e-03, 6.43000762e-01, 3.50125781e-01],
                           [2.75025136e-02, 4.76112393e-03, 9.67736363e-01],
                           [3.22944111e-02, 8.37677122e-01, 1.30028467e-01],
                           [9.85247145e-01, 1.13391711e-02, 3.41368409e-03],
                           [6.03734929e-03, 3.86549176e-04, 9.93576102e-01]]))

print("Your E step works within the expected range: ", np.allclose(my_estep, expected_estep))

# test M step
my_pi, my_mu, my_sigma = GMM(points, 3)._M_step(expected_estep, True)
expected_pi = np.array([0.1479711 , 0.38249961, 0.46952929])
expected_mu = np.array([[ -0.15519344,  1.22500376,  0.03548931],
                        [-0.36778399, -0.18068954, -0.65203503],
                        [-0.28448252, -0.09079301,  0.92618845]])
expected_sigma = np.array([[[- 0.67247982, -0.25027742,  0.0774841 ],
                           [-0.25027742,  0.34077941,  0.04853111],
                           [ 0.0774841 ,  0.04853111,  0.2987035 ]],
                           [[ 0.49792869,  0.07842407, -0.09002534],
                           [ 0.07842407,  0.1618932 , -0.10696588],
                           [-0.09002534, -0.10696588,  0.20401203]],
                           [[ 0.65130447,  0.0049166 , -0.39258756],
                           [ 0.0049166 ,  0.22371688,  0.18769942],
                           [-0.39258756,  0.18769942,  0.70840301]]])

print("Your M step works within the expected range: ", np.allclose(my_pi, expected_pi))

```

Your multinormal pdf works within the expected range: True

Your lljoint works within the expected range: True

Your E step works within the expected range: True

Your M step works within the expected range: True

3.3 Image Compression and pixel clustering [10pts +

5pts]

Images typically need a lot of bandwidth to be transmitted over the network. In order to optimize this process, most image processors perform lossy compression of images (lossy implies some information is lost in the process of compression).

In this section, you will use your GMM algorithm implementation to do pixel clustering and compress the images. That is to say, you would develop a lossy image compression algorithm. (Hint: you can adjust the number of clusters formed and justify your answer based on visual inspection of the resulting images or on a different metric of your choosing)

Special Notes

- Try to add a small value(e.g. SIGMA_CONST and LOG_CONST) before taking the operation if the output image is solid black.
- The output images may be slightly different due to different initialization methods in GMM() function.

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 [104...]

```
#####
### NO NEED TO CHANGE THIS CELL ###
#####

# helper function for performing pixel clustering.
def cluster_pixels_gmm(image, K, full_matrix = True):
    """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 = 10)(full_matrix)
    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 [116...]

```
# the direction of two images. Both of them are from ImageNet
img1_dir ='data/images_for_GMM/bird.jpg'
img2_dir ='data/images_for_GMM/cat.jpg'

# example of loading image
# image1 = imageio.imread('../data/images_for_GMM/bird.jpg')

# this is for you to implement
def perform_compression(image, min_clusters=5, max_clusters=15):
    """
    Using the helper function above to find the optimal number of clusters that can app
    You can simply examine the answer based on your visual inspection (i.e. looking a

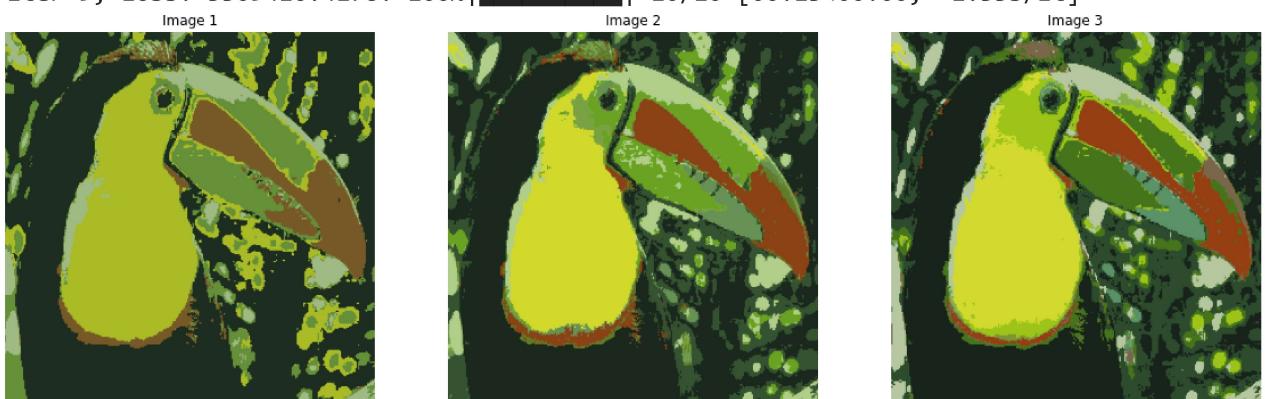
    Args:
        image: input image of shape(H, W, 3)
        min_clusters, max_clusters: the minimum and maximum number of clusters you shou
        (Usually the maximum number of clusters would not exceed 15)

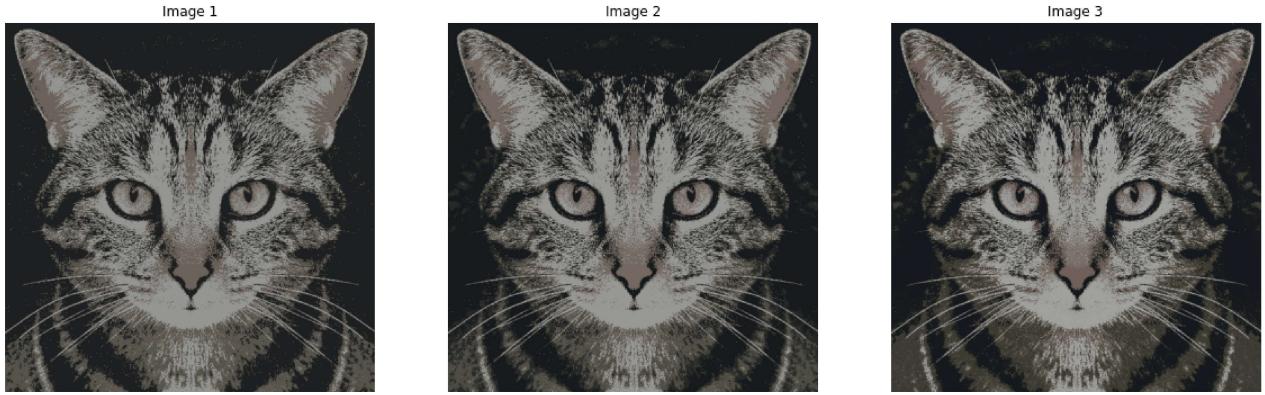
    Return:
        plot: comparison between original image and image pixel clustering.
        optional: any other information/metric/plot you think is necessary.
    """
    img_1 = cluster_pixels_gmm(image, 5)
    img_2 = cluster_pixels_gmm(image, 10)
    img_3 = cluster_pixels_gmm(image, 15)
    plot_images([img_1, img_2, img_3], ['Image 1', "Image 2", 'Image 3'])

image1 = imageio.imread(img1_dir)
perform_compression(image1, 5, 10)

image2 = imageio.imread(img2_dir)
perform_compression(image2, 5, 10)
```

iter	loss	Time	Rate
9	1140899.1259	100%	10/10 [00:01<00:00, 7.43it/s]
9	1109434.3309	100%	10/10 [00:02<00:00, 4.05it/s]
9	1101351.7215	100%	10/10 [00:04<00:00, 2.48it/s]
9	3317916.7544	100%	10/10 [00:04<00:00, 2.18it/s]
9	3311831.8969	100%	10/10 [00:08<00:00, 1.16it/s]
9	3309416.4178	100%	10/10 [00:13<00:00, 1.33s/it]





(Bonus for all) [5 pts]

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

In [157...]

```
#####
### DO NOT CHANGE THIS CELL ###
#####

def compare_matrix(image, K):
    """
    Args:
        image: input image of shape(H, W, 3)
        K: number of components

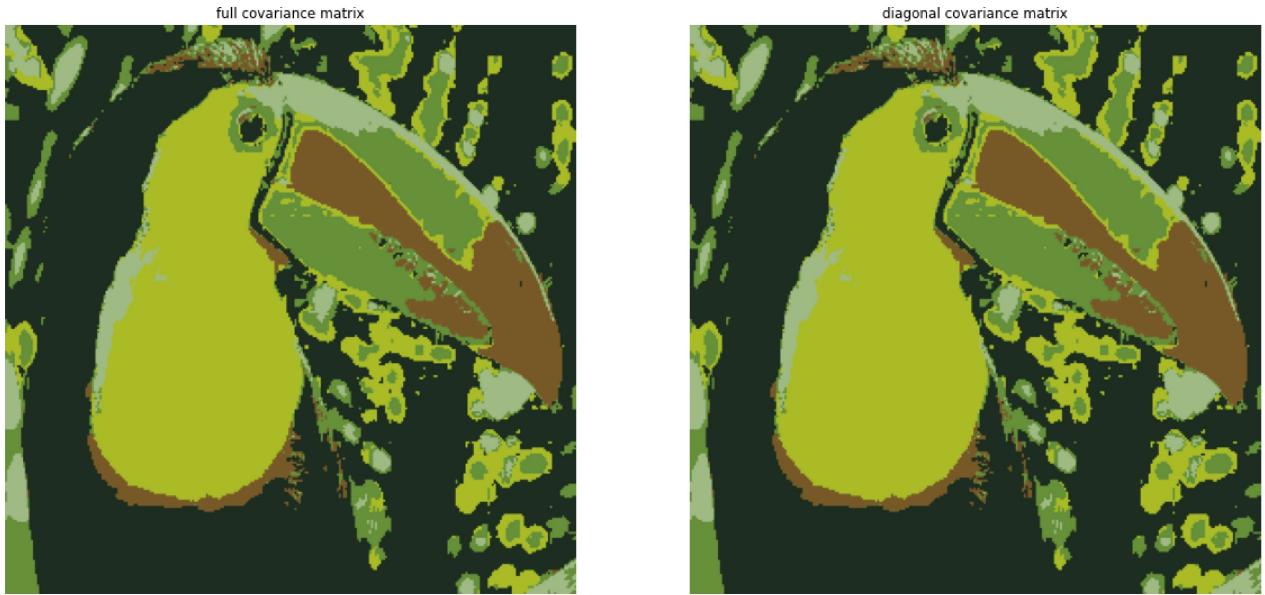
    Return:
        plot: comparison between full covariance matrix and diagonal covariance matrix.
    """
    #full covariance matrix
    gmm_image_full = cluster_pixels_gmm(image, K, full_matrix = True)
    #diagonal covariance matrix
    gmm_image_diag = cluster_pixels_gmm(image, K, full_matrix = False)

    plot_images([gmm_image_full, gmm_image_diag], ['full covariance matrix', 'diagonal
```

In [158...]

```
compare_matrix(image1, 5)
```

```
iter 9, loss: 1140899.1259: 100%|██████████| 10/10 [00:01<00:00, 7.61it/s]
iter 9, loss: 1140899.1259: 100%|██████████| 10/10 [00:01<00:00, 8.27it/s]
```



For both images analyzed above it looks like the $K = 20$ case had better clustering. This does not mean the higher K will always lead to a better clustering. As we can see in case with the cat even though $k = 20$ gave better results it had diminishing results.

4. (Bonus for All) Cleaning Messy data with semi-supervised learning [30pts]

Learning to work with messy data is a hallmark of a well-rounded data scientist. In most real-world settings the data given will usually have some issue, so it is important to learn skills to work around such impasses. This part of the assignment looks to expose you to clever ways to fix data using concepts that you have already learned in the prior questions.

Question

After graduating from Georgia Tech with your shiny new degree, you are recruited to help with safety testing for the Mars rocket at NASA. Of course NASA won't be sending rocket after rocket to stress test your fellow employees' engineering (they also graduate from Tech, so you have full confidence in them), so instead, NASA has decided to run numerous simulations on the current engineering design of the Mars rocket. The simulation collects shuttle data from its sensors, resulting in 8 features which include bypass, rad flow, etc. These features are contained within the first through eighth columns. The ninth column shows the label with 1 being a successful simulation and 0 being an unsuccessful simulation.

However, due to an intern accidentally deleting random data points, 20% of the entries are missing labels and 30% are missing characterization data. Since simply removing the corrupted entries would not reflect the true variance of the data, your job is to implement a solution to clean the data so it can be properly classified.

Your job is to assist NASA in cleaning the data and implementing a semi-supervised learning framework to help them create a general classifier for future simulations.

You are given two files for this task:

- data.csv: the entire dataset with complete and incomplete data
- validation.csv: a smaller, fully complete dataset made after after the intern deleted the datapoints

4.1.a Data Separating [3pt]

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
- labeled_incomplete: containing partial characterization data (i.e., one of the features is NaN) and corresponding labels
- unlabeled_complete: containing only complete material characterization results (i.e., the label is NaN)

```
In [ ]: #####  
### DO NOT CHANGE THIS CELL ###  
#####  
  
from semisupervised import complete_  
from semisupervised import incomplete_  
from semisupervised import unlabeled_
```

4.1.b KNN [7pts]

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 [ ]: #####  
### DO NOT CHANGE THIS CELL ###  
#####  
  
from semisupervised import CleanData
```

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 [ ]: #####  
### DO NOT CHANGE THIS CELL ###  
#####  
  
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.na  
  
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 using EM](#)" by Nigam et al. (2000). While you are recommended to read the whole paper this assignment focuses on items 5.2 and 6.1. Write a brief summary of three interesting highlights of the paper (50-word maximum).

I think the paper as a whole was a very interesting read. The thing I found most interesting was using an EM algorithm in combination with another algorithm (bayes classifier) to make a new algorithm. I liked that the approach taken by the paper was not too complicated and was pretty easy to understand, yet they were the first ones to do it. I also think the paper did some very insightful analysis on their results and was overall pretty concise and interesting.

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 [ ]: #####
### DO NOT CHANGE THIS CELL #####
#####

from semisupervised import SemiSupervised
```

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.

To achieve the full 5 points you must get these scores:

- semi_supervised_score > .87
- GNB_cleandata_score > .87
- GNB_onlycomplete_score > .87

In []:

```
#####
### DO NOT CHANGE THIS CELL ###
#####

from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score

from semisupervised import ComparePerformance
```

In []:

```
#####
### DO NOT CHANGE THIS CELL ###
#####

from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score

# Load and clean data for the next section
all_data = np.loadtxt('data.csv', delimiter=',')

labeled_complete = complete_(all_data)
labeled_incomplete = incomplete_(all_data)
unlabeled = unlabeled_(all_data)

clean_data = CleanData()(labeled_incomplete, labeled_complete, 10)
# Load unlabeled set
# append unlabeled flag
unlabeled = np.delete(unlabeled, -1, axis=1)
unlabeled_flag = -1*np.ones((unlabeled.shape[0],1))
unlabeled = np.concatenate((unlabeled, unlabeled_flag), 1)

# =====
# SEMI SUPERVISED

# format training data
points = np.concatenate((clean_data, unlabeled),0)
# train model
(pi, mu, sigma) = SemiSupervised()(points, 2)

# =====
# SUPERVISED WITH CLEAN DATA (SKLEARN)

clean_clf = GaussianNB()
clean_clf.fit(clean_data[:, :8], clean_data[:, 8])

# =====
# SUPERVISED WITH ONLY THE COMPLETE DATA (SKLEARN)

complete_clf = GaussianNB()
complete_clf.fit(labeled_complete[:, :8], labeled_complete[:, 8])

# =====
# COMPARISON
```

```
# Load test data
independent = np.loadtxt('validation.csv', delimiter=',', )

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

# =====

print("""==COMPARISON==""")
print("""SemiSupervised Accuracy:""", ComparePerformance().accuracy_semi_supervised(classification))
print("""Supervised with clean data: GNB Accuracy:""", ComparePerformance().accuracy_GNB(classification))
print("""Supervised with only complete data: GNB Accuracy:""", ComparePerformance().accuracy_GNB(classification))
```

In []: