

Homework 2

DUE NOV 1th at 11:59 PM

You are expected to turn in a pdf version of this notebook with all your codes, results, and figures (Use the print option). Make sure the figures and results are visible as you want them to appear in the pdf before turning it in. Please do not modify the instructions as doing so will limit our ability to follow and grade your answers.

Problem 1 (HW1 Problem 4)

In this problem, you will work on the clustering problem using Bottom-up Agglomerative clustering and K-mean clustering.

a) A 4-D dataset is given in 'iris.csv' with the last column being the ground truth label. Load the file. Store the data in a variable **X** and store the label in a variable **y**. Because clustering is an unsupervised task, there is no need for the labels during training.

```
In [191]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

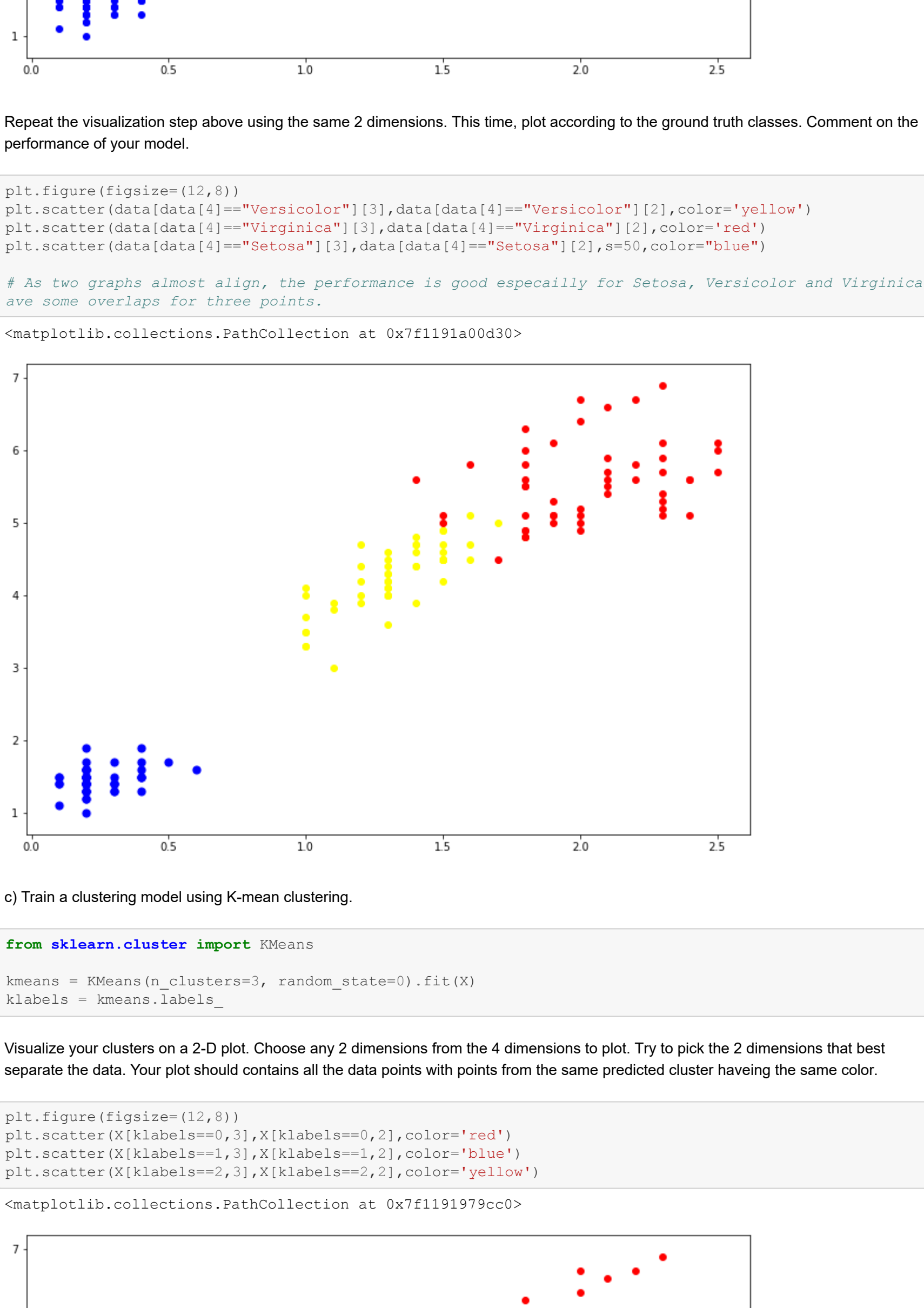
data = pd.read_csv('iris.csv',header = None)
X = data[range(4)].values
y = data[4].values

b) Train a clustering model using Bottom-up Agglomerative clustering.
```

Visualize your clusters on a 2-D plot. Choose any 2 dimensions from the 4 dimensions to plot. Try to pick the 2 dimensions that best separate the data. Your plot should contains all the data points with points from the same predicted cluster having the same color.

```
In [195]: from sklearn.cluster import AgglomerativeClustering
clustering = AgglomerativeClustering(n_clusters=3,affinity="euclidean",linkage="ward")
clustering.fit(X)
labels=clustering.labels_
```

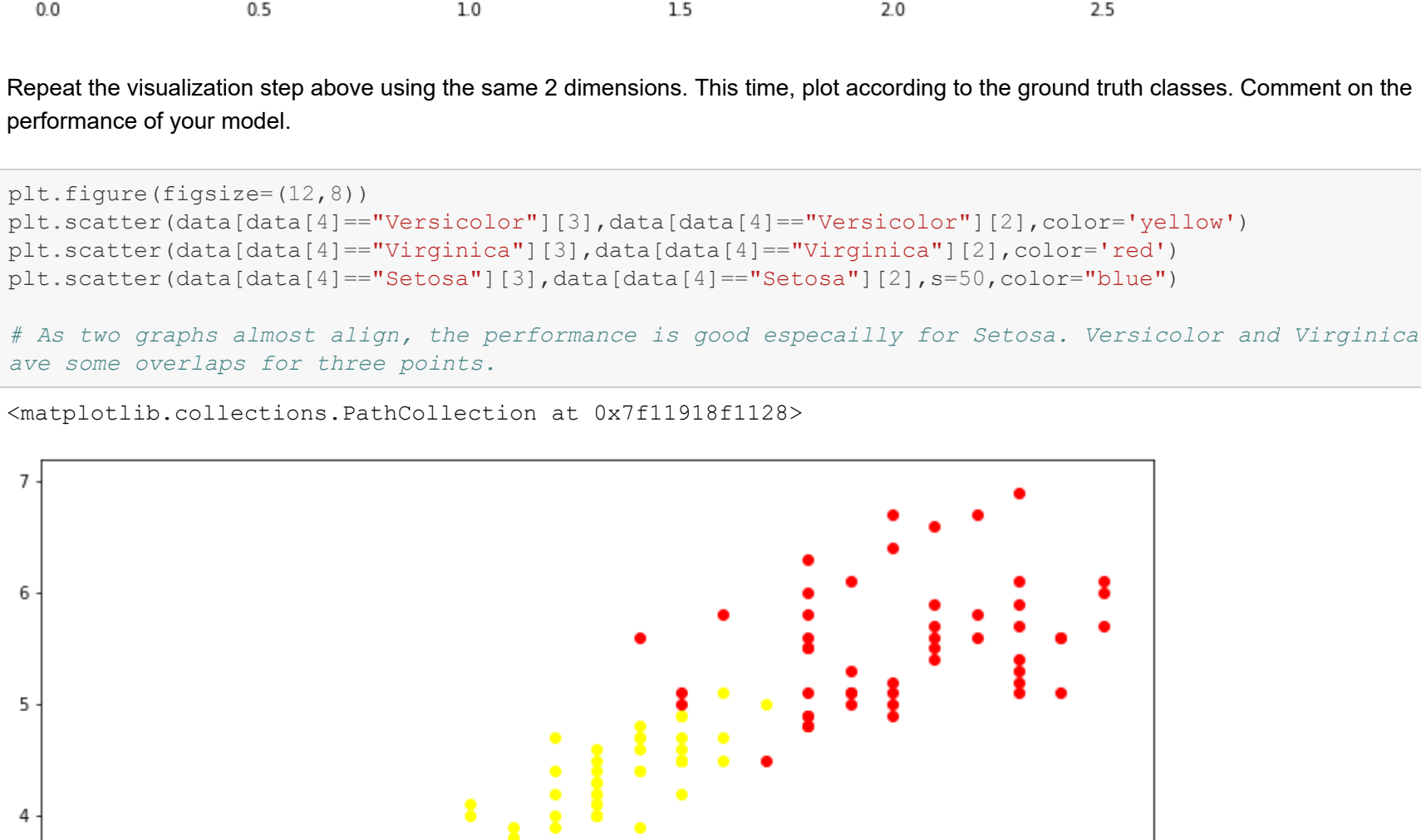
```
In [196]: plt.figure(figsize=(12,8))
plt.scatter(data[data[4]=="Versicolour"][3],data[data[4]=="Versicolour"][2],color='yellow')
plt.scatter(X[labels==1,3],X[labels==1,2],color='blue')
plt.scatter(X[labels==2,3],X[labels==2,2],color='yellow')
```



Repeat the visualization step above using the same 2 dimensions. This time, plot according to the ground truth classes. Comment on the performance of your model.

```
In [6]: plt.figure(figsize=(12,8))
plt.scatter(data[data[4]=="Versicolour"][3],data[data[4]=="Versicolour"][2],color='yellow')
plt.scatter(data[data[4]=="Virginica"][3],data[data[4]=="Virginica"][2],color='red')
plt.scatter(data[data[4]=="Setosa"][3],data[data[4]=="Setosa"][2],color='blue')

# As two graphs almost align, the performance is good especially for Setosa, Versicolour and Virginica h
# ave some overlaps for three points.
```

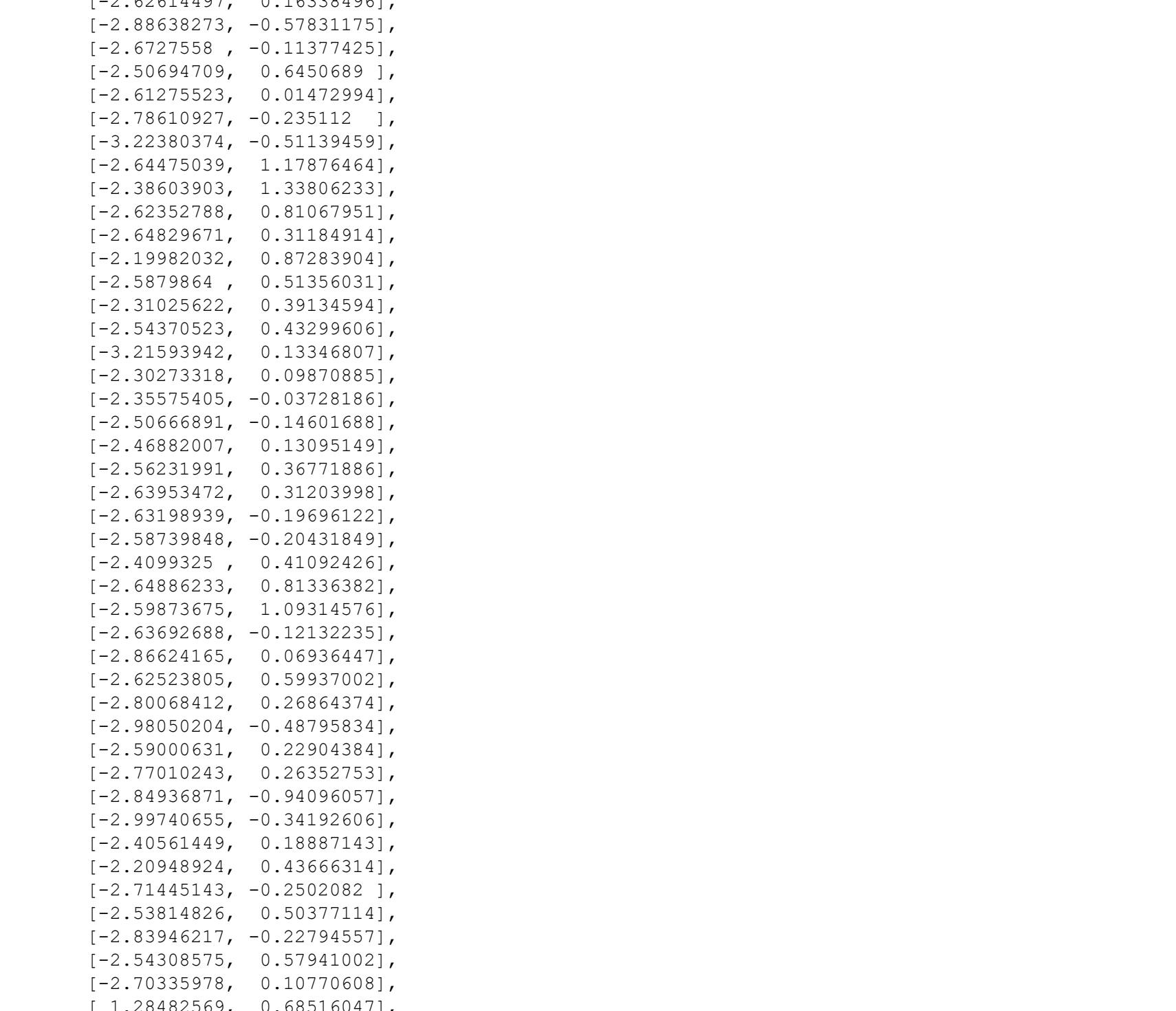


c) Train a clustering model using K-mean clustering.

```
In [7]: from sklearn.cluster import KMeans
kmeans = KMeans(n_clusters=3, random_state=0).fit(X)
klabels = kmeans.labels_
```

Visualize your clusters on a 2-D plot. Choose any 2 dimensions from the 4 dimensions to plot. Try to pick the 2 dimensions that best separate the data. Your plot should contains all the data points with points from the same predicted cluster having the same color.

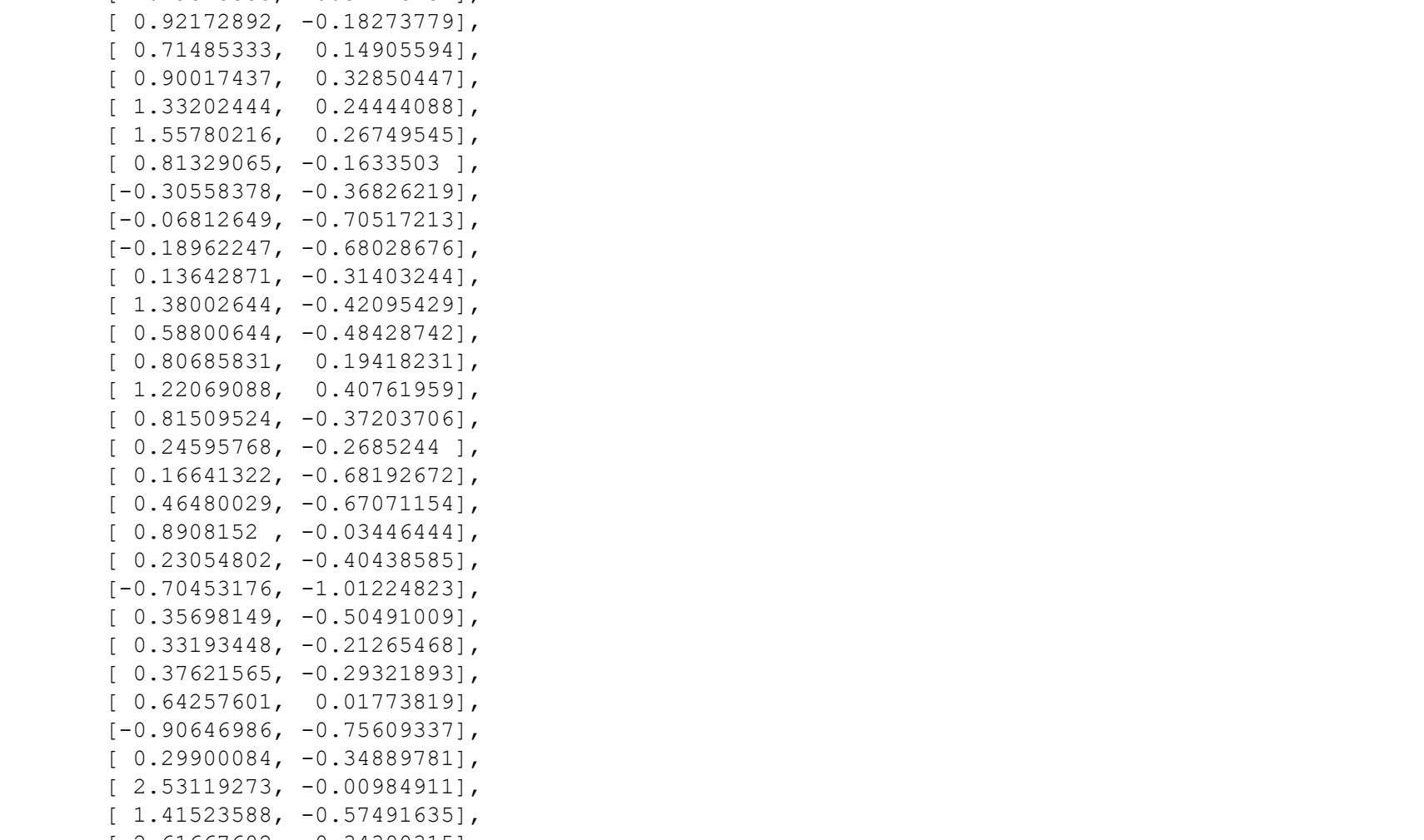
```
In [8]: plt.figure(figsize=(12,8))
plt.scatter(X[klabels==0,3],X[klabels==0,2],color='red')
plt.scatter(X[klabels==1,3],X[klabels==1,2],color='blue')
plt.scatter(X[klabels==2,3],X[klabels==2,2],color='yellow')
```



Repeat the visualization step above using the same 2 dimensions. This time, plot according to the ground truth classes. Comment on the performance of your model.

```
In [9]: plt.figure(figsize=(12,8))
plt.scatter(data[data[4]=="Versicolour"][3],data[data[4]=="Versicolour"][2],color='yellow')
plt.scatter(data[data[4]=="Virginica"][3],data[data[4]=="Virginica"][2],color='red')
plt.scatter(data[data[4]=="Setosa"][3],data[data[4]=="Setosa"][2],color='blue')

# As two graphs almost align, the performance is good especially for Setosa. Versicolour and Virginica h
# ave some overlaps for three points.
```



d) Perform Principle Component Analysis (PCA) on the data. Project the original data on the 2 largest principle components. Store this new projected 2-D data in a variable **X_projected**.

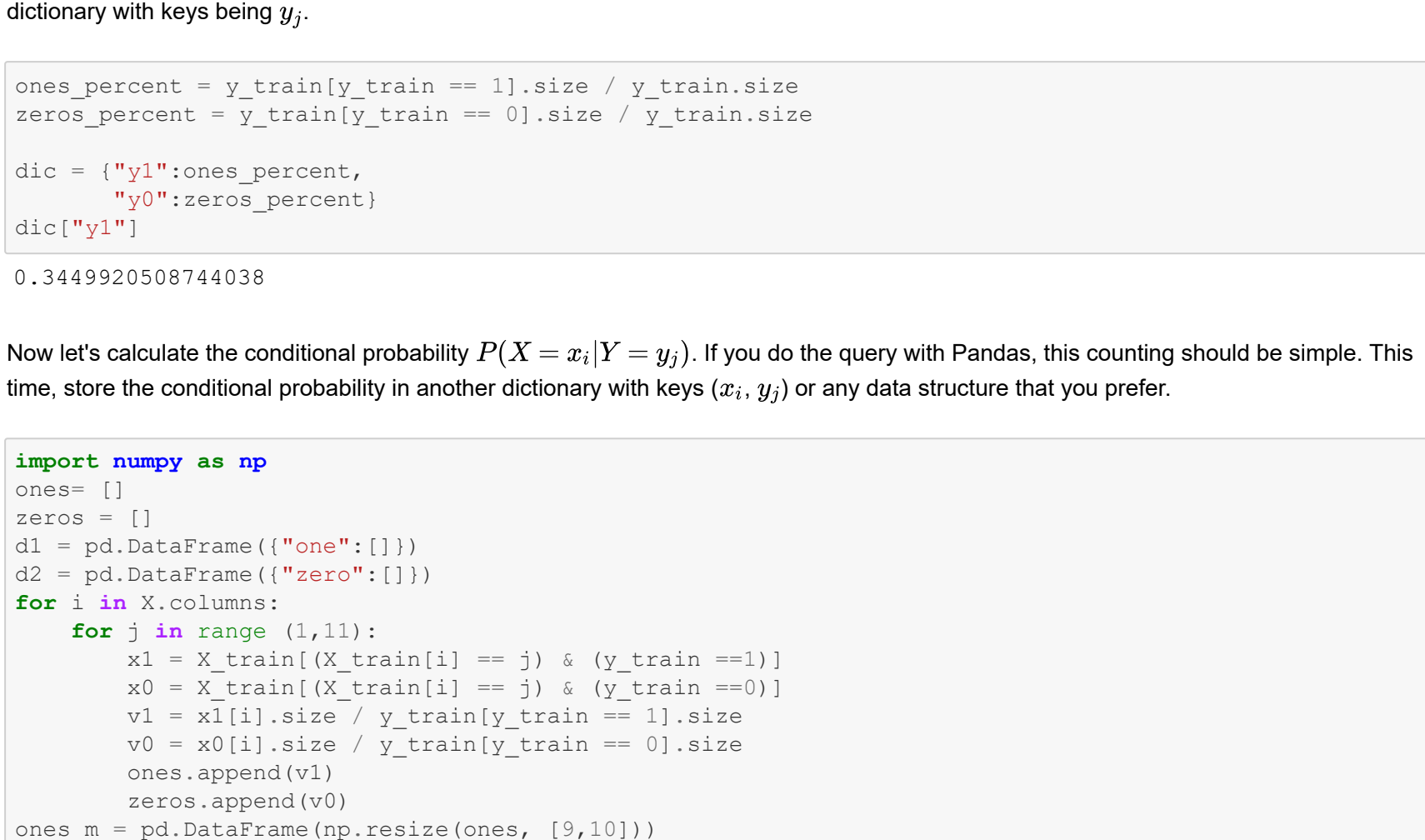
```
In [10]: from sklearn.decomposition import PCA
pca = PCA(n_components=2).fit(X)
X_pca = pca.transform(X)
```

Out[10]: array([[-2.68412563, 0.31939725],
[-2.71414169, -0.17700123],
[-2.88890957, -0.14494943],
[-2.74534286, -0.31829898],
[-2.72871654, 0.32675451],
[-2.28085963, 0.74133045],
[-2.82037375, -0.08946188],
[-2.62614497, 0.16338496],
[-2.88638273, -0.57831175],
[-2.6721758 , -0.11377429],
[-2.5684709 , 0.64506889],
[-2.61275523, 0.01472994],
[-2.78610927, -0.235112],
[-1.22380374, -0.51139456],
[-2.64475039, 1.17876464],
[-2.38603903, 1.33806233],
[-2.62332788, 0.81867951],
[-2.64829671, 0.31184914],
[-2.19982032, 0.87283904],
[-2.9379864 , 0.51250311],
[-2.8102652 , 0.39134594],
[-2.43750523, 0.43299606],
[-2.21593942, 0.13346807],
[-2.30271318, 0.09870085],
[-2.35575405, -0.03728186],
[-2.5668691 , -0.14601688],
[-2.64882007, -0.13095149],
[-2.56231991, 0.36771886],
[-2.6393472 , 0.31203998],
[-2.63188939, -0.19096122],
[-2.58789488, -0.20431693],
[-2.4099328 , 0.41092426],
[-2.6486233 , 0.81336382],
[-2.5901043 , 0.09345761],
[-2.63692688, -0.12132235],
[-2.8662165 , 0.06936447],
[-2.8523055 , 0.59937002],
[-2.73292537, -0.24864374],
[-2.98052024, -0.48759834],
[-2.5900631 , 0.22304384],
[-2.77010243, -0.24357659],
[-2.84936871, -0.94060575],
[-2.39740555, -0.34192606],
[-2.10521418, -0.18887451],
[-2.20948924, 0.43666314],
[-2.7145143 , -0.2502082],
[-2.53814826, 0.50377123],
[-2.8326217 , -0.37345591],
[-2.54308575, 0.57941002],
[-2.70359978, 0.10776089],
[-2.84823569, 0.68516047],
[0.93248853, 0.31833364],
[1.44410232, 0.50426282],
[0.18331772, -0.82759391],
[1.68810326, 0.07459068],
[0.64166908, -0.41824687],
[1.09506066, 0.28346827],
[0.47812267, -0.00489489],
[1.04413183, 0.2283619],
[-0.0087454 , -0.72308191],
[0.57894084, -0.37657119],
[0.5116956 , -0.10398124],
[0.26497651, -0.55003646],
[0.94843451, -0.12481765],
[1.94410978, -0.18732371],
[0.92786078, 0.46717949],
[0.66028376, -0.35296671],
[0.2610499 , -0.33610779],
[0.94473373, 0.34534555],
[0.04522698, -0.58383438],
[1.1623288 , -0.08461681],
[0.25788842, -0.06892503],
[1.29818388, -0.32778731],
[0.92172892, -0.18273799],
[0.7148537 , -0.14705594],
[0.90017437, 0.32850447],
[1.33202444, 0.24440889],
[1.57800216, -0.26749059],
[0.81329065, -0.1633503],
[-0.5958378 , -0.36826219],
[-0.0812649 , -0.70517213],
[-0.1892247 , -0.68028676],
[0.13642871, -0.31403244],
[1.38002644, -0.40295489],
[0.5880619 , -0.48428722],
[0.80685831, 0.19418231],
[1.20689088, 0.40761959],
[0.18502924, -0.37023011],
[0.52495768, -0.2685244],
[0.16641322, -0.68192672],
[0.64848002 , -0.67017154],
[0.48004542, -0.03466444],
[0.23054802, 0.40438585],
[-0.70453176, -0.101224823],
[0.36801419, -0.50491059],
[0.33193448, -0.21265468],
[0.37621565, -0.29321893],
[0.62437601, -0.01773931],
[-0.90646986, -0.75609337],
[0.29900084, -0.34889781],
[2.53119273, -0.00984912],
[1.4523988 , -0.57491639],
[2.61676022, 0.34390315],
[1.97153105, -0.1797279],
[2.35005092, -0.04026059],
[3.39703874, 0.55083667],
[0.52123224, -1.19275873],
[2.93287807, 0.3355],
[2.32122882, -0.2438315],
[2.91675097, 0.78279195],
[1.66177415, 0.24222841],
[0.2610499 , -0.21361762],
[2.1655918 , 0.21627559],
[1.34616358, -0.20439933],
[1.95892822, -0.53940761],
[1.9045637 , 0.11925069],
[1.94968906, 0.04194326],
[3.48705536, 1.17573154],
[1.78545442, 0.25732371],
[1.30079171, -0.76114964],
[2.42781791, 0.37819601],
[1.59001461, -0.60681053],
[3.49992004, 0.4660741],
[1.38876613, -0.20439933],
[2.7543305 , 0.33499061],
[2.61409047, 0.56090136],
[1.25850816, -0.17970479],
[1.29112056, -0.11686865],
[1.62080812, -0.20972948],
[2.38803032, 0.4646398],
[2.84167278, 0.37526917],
[3.20007366, 1.37416509],
[2.15943764, -0.21277758],
[1.44416124, -0.14341341],
[1.78123483, -0.49990168],
[3.0764999 , 0.68085568],
[2.14424331, 0.1400642],
[1.95959815, 0.04930053],
[1.69326134, -0.16459026],
[2.07611114, 0.37228787],
[2.31415471, 0.18365128],
[1.92226718, 0.40920471],
[1.41523588, -0.57491635],
[2.56301338, 0.2778626],
[2.41874618, -0.7047982],
[1.94410978, -0.18732371],
[1.52716661, -0.37531698],
[1.76434572, 0.07885885],
[1.90094161, 0.11662786],
[1.39018886, -0.28266094]]]

Repeat part b on the new 2-D data. Train the Bottom-up Agglomerative model and visualize your results.

```
In [11]: from sklearn.cluster import AgglomerativeClustering
clustering = AgglomerativeClustering(n_clusters=3,affinity="euclidean",linkage="ward")
clustering.fit(X_pca)
labels_pca=clustering.labels_
```

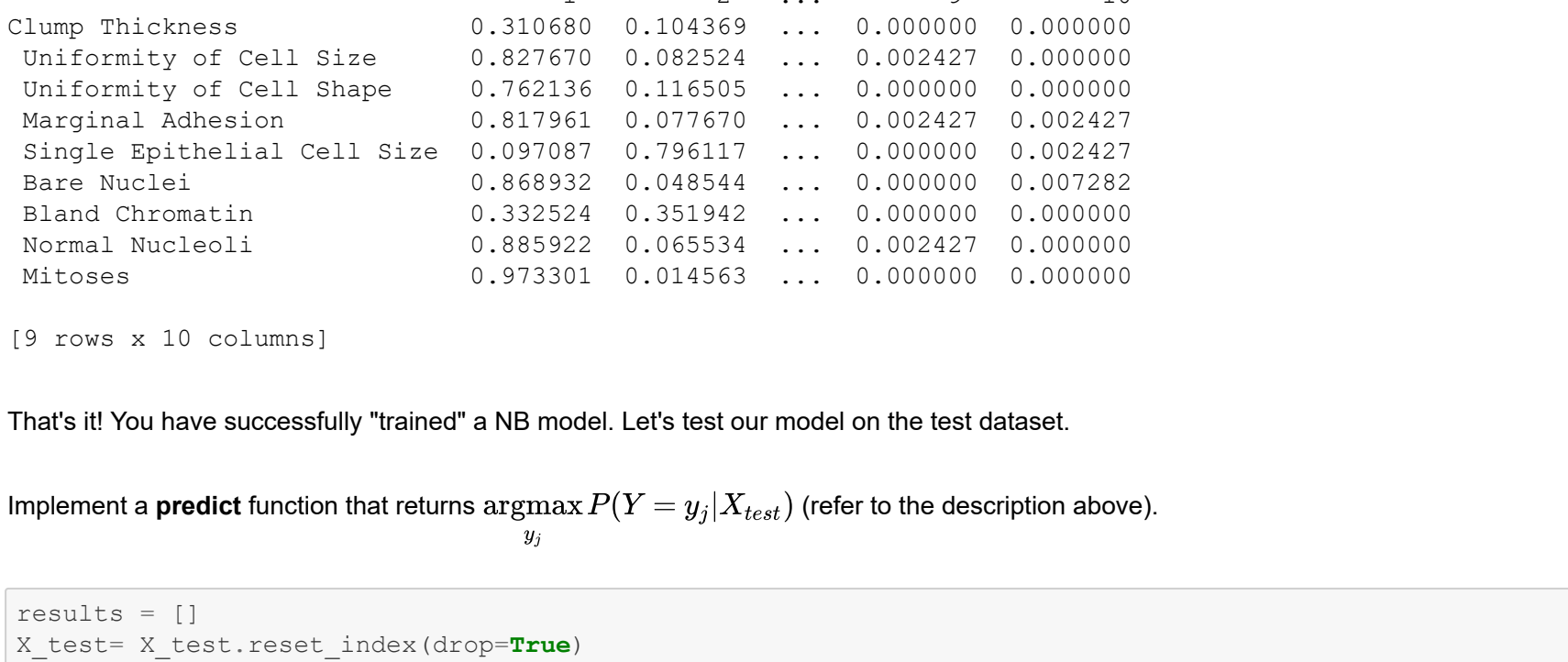
```
In [11]: plt.figure(figsize=(12,8))
plt.scatter(X_pca[labels_pca==0,1],X_pca[labels_pca==0,0],color='yellow')
plt.scatter(X_pca[labels_pca==1,1],X_pca[labels_pca==1,0],color='blue')
plt.scatter(X_pca[labels_pca==2,1],X_pca[labels_pca==2,0],color='red')
```



Repeat part c on the new 2-D data. Train the K-means model and visualize your result.

```
In [12]: from sklearn.cluster import KMeans
kmeans = KMeans(n_clusters=3, random_state=0).fit(X_pca)
klabels_pca = kmeans.labels_
```

```
In [12]: plt.figure(figsize=(12,8))
plt.scatter(X_pca[klabels_pca==0,1],X_pca[klabels_pca==0,0],color='red')
plt.scatter(X_pca[klabels_pca==1,1],X_pca[klabels_pca==1,0],color='blue')
plt.scatter(X_pca[klabels_pca==2,1],X_pca[klabels_pca==2,0],color='yellow')
```



Compare the quality of 4-D and 2-D clusterings. When would the ideas of projection and dimensionality reduction be useful?

PCA used in combination with k means clustering dimension reduction was very useful and increase accuracy, but when PCA used with Bottom-up Agglomerative clustering doesn't change much.

Overall, 2-D has better partition than 4-D. Thus PCA will be slightly better. The higher dimension, the more useful PCA will be as better clustering would be formed.

Problem 2

In this problem, you will first implement the Naive Bayes (NB) algorithm from scratch. We will use a dataset for classifying if a patient has breast cancer. Each instance (row) in the dataset is a patient described by the 9 following features (with their corresponding values).

- Clump Thickness: 1 - 10
- Uniformity of Cell Size: 1 - 10
- Uniformity of Cell Shape: 1 - 10
- Marginal Adhesion: 1 - 10
- Single Epithelial Cell Size: 1 - 10
- Bare Nuclei: 1 - 10
- Bland Chromatin: 1 - 10
- Normal Nucleoli: 1 - 10
- Mitoses: 1 - 10

Given these features, we will classify a car into one of the two classes: 0 (benign) or 1 (malignant).

NB is a very simple algorithm. Consider a feature **X** and each class label **y_{ij}**. NB calculates the value of $P(X = x_i | Y = y_j)$. For example, take the feature **Mitoses**. NB will calculate all the following values.

- $P(\text{Mitoses} = 1 | \text{Class} = 0)$, $P(\text{Mitoses} = 2 | \text{Class} = 0)$, $P(\text{Mitoses} = 3 | \text{Class} = 0)$, ... , $P(\text{Mitoses} = 10 | \text{Class} = 0)$
- $P(\text{Mitoses} = 1 | \text{Class} = 1)$, $P(\text{Mitoses} = 2 | \text{Class} = 1)$, $P(\text{Mitoses} = 3 | \text{Class} = 1)$, ... , $P(\text{Mitoses} = 10 | \text{Class} = 1)$

Repeat this calculation for all the features. In the end, NB keep a recording of all possible $P(X|Y)$. The calculation itself is intuitive:

$$P(X=x_i|Y=y_j) = \frac{\text{Number of rows with } X=x_i \text{ and } Y=y_j}{\text{Number of rows with } Y=y_j}$$

In addition, NB also calculate the priors probability $P(Y = y_j)$. Again, intuitively,

$$P(Y = y_j) = \frac{\text{Number of rows with } Y = y_j}{\text{Number of rows in the dataset}}$$

Given a test example $X_{test} = \{X_0 = x_0, X_1 = x_1, \dots, X_i = x_i\}$, for each class label **y_{ij}**, NB calculate:

$$P(Y = y_j | X_{test}) = P(X_0 = x_0 | Y = y_j) P(X_1 = x_1 | Y = y_j) \dots P(X_i = x_i | Y = y_j) P(Y = y_j)$$
$$= P(X_0 = x_0 | Y = y_j) P(X_1 = x_1 | Y = y_j) \dots P(X_i = x_i | Y = y_j) P(Y = y_j)$$

Such calculation is easy since we have bookkept all $P(X|Y)$ and all $P(Y)$ in previous steps. The output of the model is simply:

$$\text{argmax}_{y_j} P(Y = y_j | X_{test})$$

You will do each of these steps following this problem. We will use Pandas to deal with the data in this problem. Pandas can do queries like "Get all the rows in which Clump Thickness = 2 and Class = 1" with minimal syntax.

a) First let's load the dataset and store it in a Pandas dataframe. Play with the dataframe and get used to the queries (this part is not graded). This guide is a good place to start:

<https://medium.com/corbyth-in-jain-english-filtering-rows-and-columns-in-pandas-python-techniques-you-must-know-6cdf32c6b14c>

Split the dataset into a training set and a testing set. Use 10% of the data as the testing set (The splitting is graded).

```
In [211]: import pandas as pd
from sklearn.model_selection import train_test_split

data = pd.read_csv('breast-cancer-wisconsin.csv')
X = data.iloc[:,0:9]
y = data.iloc[:,9]
X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.10, random_state=42)
```

b) Now we will build up the bookkeeping. Let's start with the priors $P(Y = y_j)$, $y_j \in \{0,1\}$. For each of these labels, calculate $P(Y = y_j)$ by simply counting the number of times **y_j** appears in the dataset divided by the size of the dataset. You can bookkeep the priors in a dictionary with keys being the **y_j**.

```
In [228]: ones_percent = y_train[y_train == 1].size / y_train.size
zeros_percent = y_train[y_train == 0].size / y_train.size

dic = {}
dic["ones_percent"] = ones_percent
dic["y1"] = y_train

Out[228]: 0.344992508744038
```

Now let's calculate the conditional probability $P(X = x_i | Y = y_j)$. If you do the query with Pandas, this counting should be simple. This time, store the conditional probability in another dictionary with keys **(x_i, y_j)** or any data structure that you prefer.

```
In [229]: import numpy as np
sample = []
zeros = []
ones = pd.DataFrame({"zero":[]})
for i in range(1,10):
    x1 = X_train[X_train[i] == 1] & (y_train == 1)
    v1 = x1[i].size / y_train[y_train == 1].size
    v0 = x0[i].size / y_train[y_train == 0].size
    ones.append(v0)
    zeros.append(v0)
ones_m = pd.DataFrame(np.resize(ones, [9,10]))
ones_m.columns = ['1,2,3,4,5,6,7,8,9,10']
ones_m.index = ["Uniformity of Cell Size", "Uniformity of Cell Shape", "Marginal Adhesion", "Single Epithelial Cell Size", "Bare Nuclei", "Bland Chromatin", "Normal Nucleoli", "Mitoses"]
print(ones_m)

zeros_m = pd.DataFrame(np.resize(zeros, [9,10]))
zeros_m.columns = ['1,2,3,4,5,6,7,8,9,10']
zeros_m.index = ["Uniformity of Cell Size", "Uniformity of Cell Shape", "Marginal Adhesion", "Single Epithelial Cell Size", "Bare Nuclei", "Bland Chromatin", "Normal Nucleoli", "Mitoses"]
print(zeros_m)
```

```
Out[229]:
```

	1	2	3	4	5	6	7	8	9	10
Clump Thickness	0.13625	0.01843	...	0.05300	0.304147					
Uniformity of Cell Size	0.01843	0.03686	...	0.023041	0.262473					
Uniformity of Cell Shape	0.009217	0.032258	...	0.032258	0.239631					
Marginal Adhesion	0.129032	0.092166	...	0.018433	0.216590					
Single Epithelial Cell Size	0.004608	0.101382	...	0.009217	0.124424					
Bare Nuclei	0.069124	0.03686	...	0.03686	0.557604					
Bland Chromatin	0.009217	0.032258	...	0.050691	0.082949					
Normal Nucleoli	0.004608	0.101382	...	0.064516	0.248448					
Mitoses	0.557604	0.119816	...	0.000000	0.055300					

[9 rows x 10 columns]

	1	2	4	...	9	10
Clump Thickness	0.310680	0.104369	...	0.000000	0.000000	
Uniformity of Cell Shape	0.827670	0.082524	...	0.002427	0.000000	
Uniformity of Cell Size	0.018433	0.03686	...	0.000000	0.000000	
Marginal Adhesion	0.817961	0.077670	...	0.002427	0.002427	
Single Epithelial Cell Size	0.097087	0.077670	...	0.000000	0.002427	
Bare Nuclei	0.868932	0.048544	...	0.000000	0.007282	
Bland Chromatin	0.332524	0.351942	...	0.000000	0.000000	
Normal Nucleoli	0.885922	0.065534	...	0.002427	0.000000	
Mitoses	0.973301	0.014563	...	0.000000	0.000000	

[9 rows x 10 columns]

That's it! You have successfully "trained" a NB model. Let's test our model on the test dataset.

Implement a predict function that returns $\text{argmax}_{y_j} P(Y = y_j | X_{test})$ (refer to the description above).

```
In [230]: results = []
X_test = X_test.reset_index(drop=True)
for i in range(0, len(y_test)):
    rones = 1
    rzeros = 1
    for j in X_test.columns:
        v = X_test.at[i,j]
        rones *= ones_m.at[j,v]
        rzeros *= zeros_m.at[j,v]
    rones *= ones_percent
    rzeros *= zeros_percent
    if (rones < rzeros):
        results.append(0)
    else:
        results.append(1)
```

Predict the label of all the instances in the test dataset, calculate and print out the accuracy.

```
In [231]: accuracy_score(y_test, results)

Out[231]: 0.9857142857142858
```

Problem 3

In this problem, you will implement the Logistic Regression (LR) algorithm from scratch. Similar to NB, LR relies on $P(Y|X)$ to predict the class of an example. However, unlike NB, a generative model, LR is a discriminative model so it does not need to estimate $P(X|Y)$ and $P(Y)$. LR assumes the form of the conditional probability $P(Y|X)$ to be:

$$P(Y|X) = f(X) = \frac{1}{1 + e^{-(\alpha + \beta X)}}$$

$f(X)$ returns a value in (0,1). The model classifies **X** as 1 if $f(X)$ is closer to 1 and 0 otherwise. We have to estimate the model parameters: the vector α and from the data, which we will do via stochastic gradient descent (SGD). In SGD, a training example is shown to the model each at a time. The model makes a prediction on the training example and the error between the prediction and the ground truth label is used to update the model's parameters. We use the log-likelihood loss to estimate the error in this problem. In particular, the log-likelihood loss for classifying the $X^{(i)}$ example with the ground-truth $y^{(i)}$ is:

$$LL(y^{(i)}, f(X^{(i)})) = -(y^{(i)} \log(f(X^{(i)})) + (1 - y^{(i)}) \log(1 - f(X^{(i)})))$$

In our case, the updating is as the follows:

$$\alpha_j(t+1) = \alpha_j(t)$$

Problem 5

What is the hyperparameters, the parameters, and the objective function for finding the best hypothesis for each of the following methods:

1. Decision Tree (Entropy gain)
2. Support Vector Machine
3. K-nearest Neighbors Clustering

1.**Hyperparameters:** depth of the tree, purity, splitter, estimator, maximum leaf nodes, depth of the trees, minimumimpurity decrease and impurity split, minimum number of sample leafs, number of features.

Parameter: nodes, attribute threshold

Objective function: The Entropy Gain: maximize the information gain at each split or optimizedduring the fitting.

$$Gain(S,A) = Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

2.**Hyperparameters:** kernel, polynomial degree, C, gamma, margin, shrinking heuristic (bool)

Parameters: maximum margin of the support vectors and the w,b variables

Objective function: To find a plane that has the maximum margin, or 'optimizing' in theconstraints. kernel function:
$$g(x) = b + w^T \phi(x) = b + \sum_{i \in SY} \alpha_i y_i K(x, x_i) > 0$$

3.**Hyperparameters:** the number of clusters, complexity, weighted neighbors, precomputeddistances, algorithm to use (such as ball_tree or brute), Leaf size passed to BallTree or KDTree,Power parameter for the Minkowski metric, metric.

Parameters: the number of centroid, measure of distance, decision boundary

Objective function: the function we used to calculate and optimize the distance such as Euclidean,Minkowski or Manhattan.

$$\hat{t} = argmax_c \sum_{i: x_i \in N_c(x, \epsilon)} \delta(t_i, c)$$