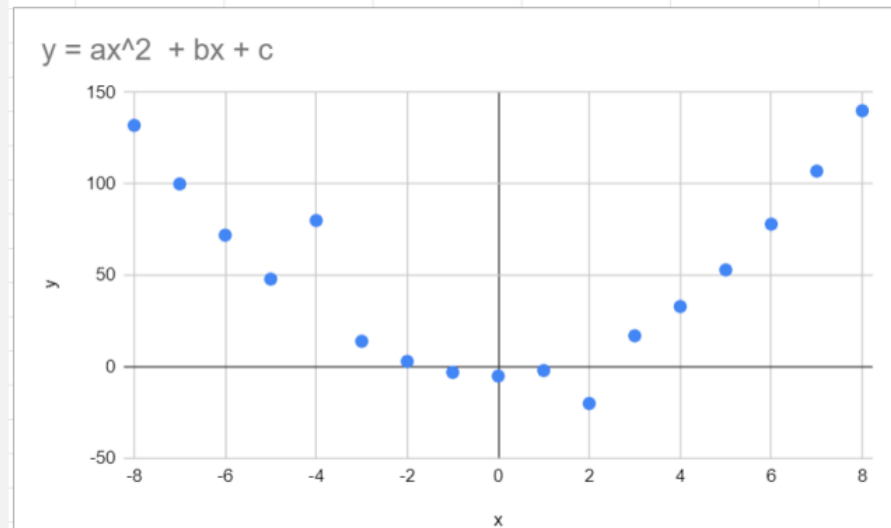


Quiz 9

1. **20 points. 2 hours. RANSAC Regression.** Use RANSAC to find a , b , c for the following dataset where points (x_i, y_i) are discrete samples from a function $f(x) = ax^2 + bx + c$ with 2 outliers. *Hint:* You should get a , b , and c close to 2.2, 0.5, -4.5, respectively.

x_i	y_i
-8	132
-7	100
-6	72
-5	48
-4	80
-3	14
-2	3
-1	-3
0	-5
1	-2
2	-20
3	17
4	33
5	53
6	78
7	107
8	140



Link: https://colab.research.google.com/drive/19yYWmq_iwGvI_EV9CoikddZyTMfYwbNp?usp=sharing

```

1 from sklearn import datasets
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LinearRegression, RANSACRegressor
5 from sklearn.metrics import r2_score, mean_squared_error
6 import numpy as np
7 import matplotlib.pyplot as plt
8 import pandas as pd
9 import math
10 import random

```

```

1 # Create dataset with some random perturbation
2 x_list = []
3 y_list = []
4 a=2.2
5 b=0.5
6 c=-4.5
7 for x in range (-8, 9): #x = -8, -7, ..., 7, 8
8     x_list.append(x)
9     y_hat = int(a*x**2 + b*x + c)
10    y_list.append(y_hat)
11
12 X = np.array(x_list)
13 y = np.array(y_list)

```

```

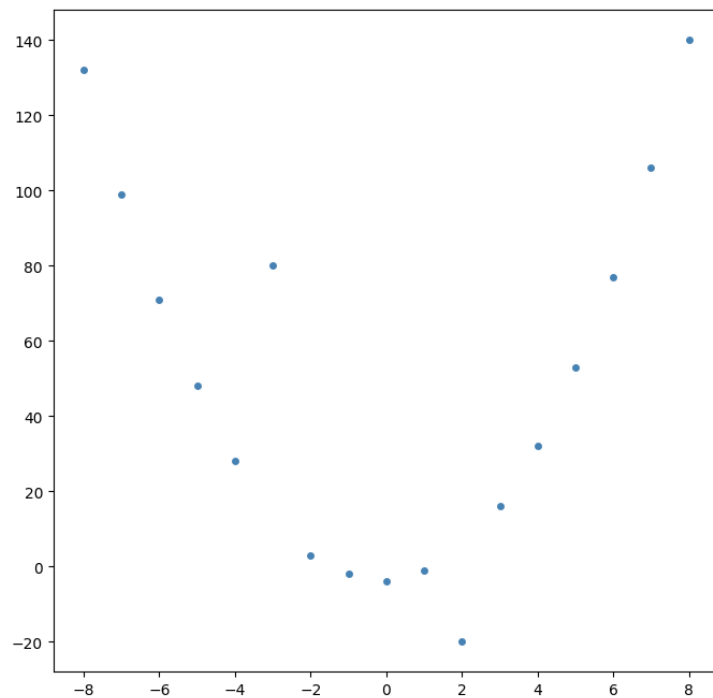
1 # Add 2 outlier points
2 y[2-9] = -20
3 y[-4+9] = 80

```

```

1 # Create scatter plot for Input Data
2 plt.figure(figsize=(8, 8))
3 plt.scatter(X, y,c='steelblue', edgecolor='white',marker='o', label='Input')
4 plt.show()

```



```

1 # Select Avg. No of rooms per dwelling as feature
2 # and fit the model
3 #
4 X = X.reshape(-1, 1)
5 print("X = ", X)
6 print("y = ", y)

```

```

X = [[-8]
      [-7]
      [-6]
      [-5]
      [-4]
      [-3]
      [-2]
      [-1]
      [0]
      [1]
      [2]
      [3]
      [4]
      [5]
      [6]
      [7]
      [8]]
y = [132 99 71 48 28 80 3 -2 -4 -1 -20 16 32 53 77 106 140]

```

```

1 import warnings
2 import numpy as np
3 from sklearn.linear_model import RANSACRegressor
4 from sklearn.metrics import mean_squared_error
5 from sklearn.datasets import make_regression
6 warnings.filterwarnings('ignore')
7
8 class PolynomialRegression(object):
9     #See https://scikit-learn.org/stable/developers/develop.html for Sklearn estimator attributes and methods
10    #Attributes: degree, coeffs
11    #Key methods: fit, predict, and score.
12    def __init__(self, degree=2):
13        print(f"Degree: {degree}")
14        self.degree = degree
15
16    def fit(self, X, y):
17        self.coefs = np.polyfit(X.ravel(), y, self.degree)
18
19    def get_params(self, deep=False):
20        return {'degree': self.degree}
21
22    def set_params(self, **parameters):
23        for parameter, value in parameters.items():
24            setattr(self, parameter, value)
25        return self
26
27    def predict(self, X):
28        poly_eqn = np.poly1d(self.coefs)
29        y_hat = poly_eqn(X.ravel())
30        return y_hat
31
32    def score(self, X, y):
33        return mean_squared_error(y, self.predict(X))

```

```

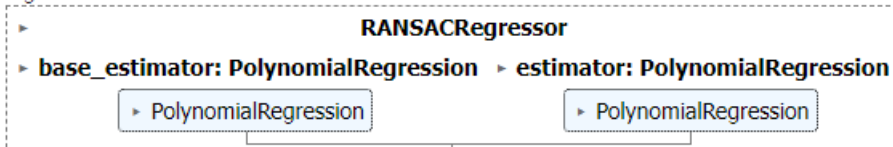
1 # Also try residual_threshold = 20
2 ransac = RANSACRegressor(
3     base_estimator=PolynomialRegression(degree=2),
4     residual_threshold= 20,
5     random_state=0,
6     min_samples=5
7 )

```

Degree: 2

```
1 ransac.fit(X, y)
```

Degree: 2



```

1 print("Estimated a b c: ", ransac.estimator_.coefs)
2 print("Actual a b c: ", np.array([a, b, c]))

```

Estimated a b c: [2.1870699 0.4958774 -4.42589906]
 Actual a b c: [2.2 0.5 -4.5]

```
1 X[inlier_mask]
```

```
array([[ -8],  
       [ -7],  
       [ -6],  
       [ -5],  
       [ -4],  
       [ -2],  
       [ -1],  
       [  0],  
       [  1],  
       [  3],  
       [  4],  
       [  5],  
       [  6],  
       [  7],  
       [  8]])
```

```
1 y[inlier_mask]
```

```
array([132, 99, 71, 48, 28,  3, -2, -4, -1, 16, 32, 53, 77,  
       106, 140])
```

```
1 X[outlier_mask]
```

```
array([[ -3],  
       [  2]])
```

```
1 y[outlier_mask]
```

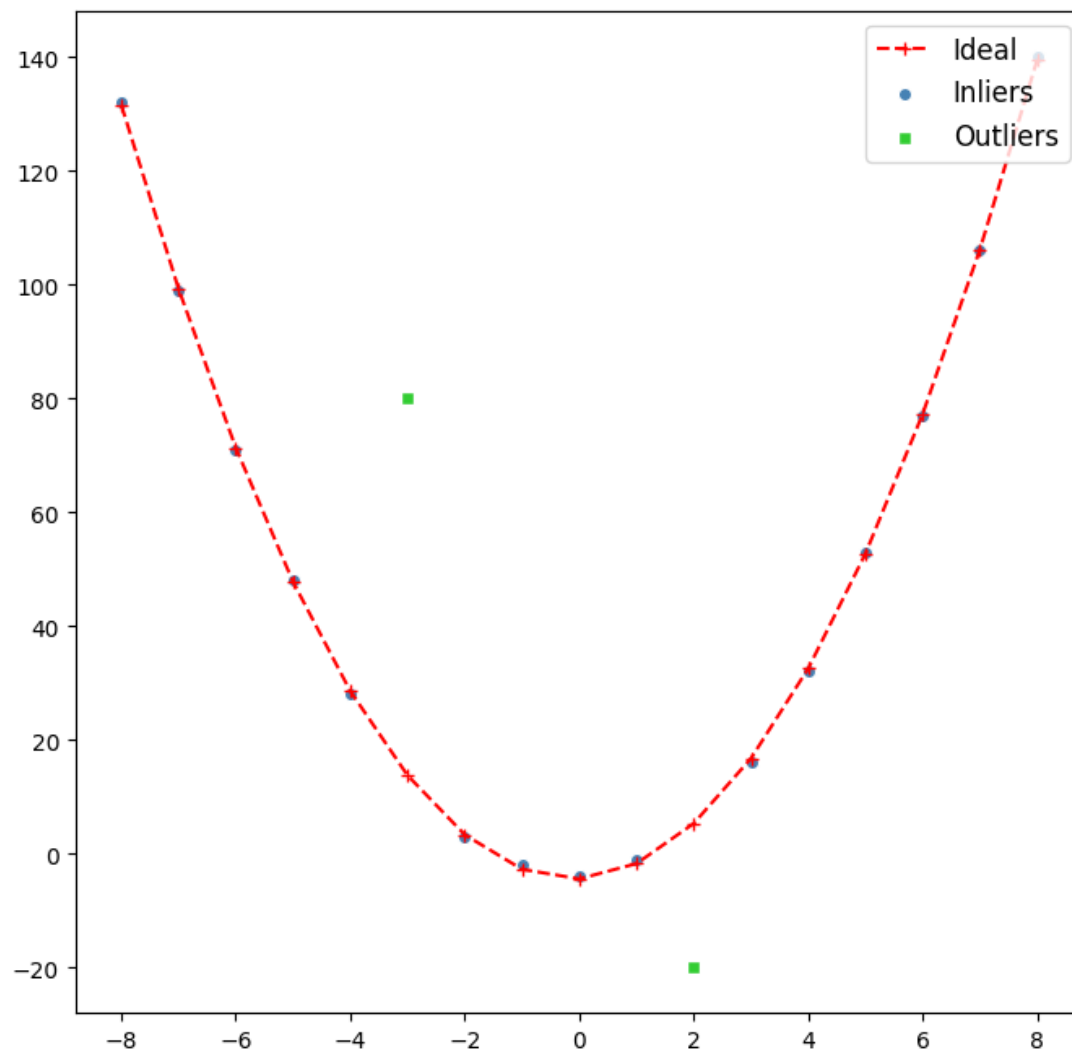
```
array([ 80, -20])
```

```
1 x_list = []  
2 y_list = []  
3 co = ransac.estimator_.coeffs  
4 for i in X:  
5     y_hat = co[0]*i**2 + co[1]*i + co[2]  
6     x_list.append(i)  
7     y_list.append(y_hat)  
8 X_perfect = np.array(x_list)  
9 y_perfect = np.array(y_list)
```

```

1 # Draw
2 #
3 # Create scatter plot for inlier dataset
4 plt.figure(figsize=(8, 8))
5
6 # perfect curve
7 plt.plot(X_perfect, y_perfect,
8 | | | c='red', marker='+', linestyle= 'dashed', label='Ideal')
9
10 # Inlier scatter plot
11 plt.scatter(X[inlier_mask], y[inlier_mask],
12 | | | c='steelblue', edgecolor='white',
13 | | | marker='o', label='Inliers')
14
15 # Create scatter plot for outlier dataset
16 plt.scatter(X[outlier_mask], y[outlier_mask],
17 | | | c='limegreen', edgecolor='white',
18 | | | marker='s', label='Outliers')
19 plt.legend(loc='upper right', fontsize=12)
20 plt.show()

```



2. Use K Means clustering on the IRIS dataset.

Link: <https://colab.research.google.com/drive/16psaSiAJXaYogzPRx2BoujsN7awBB4uh?usp=sharing>

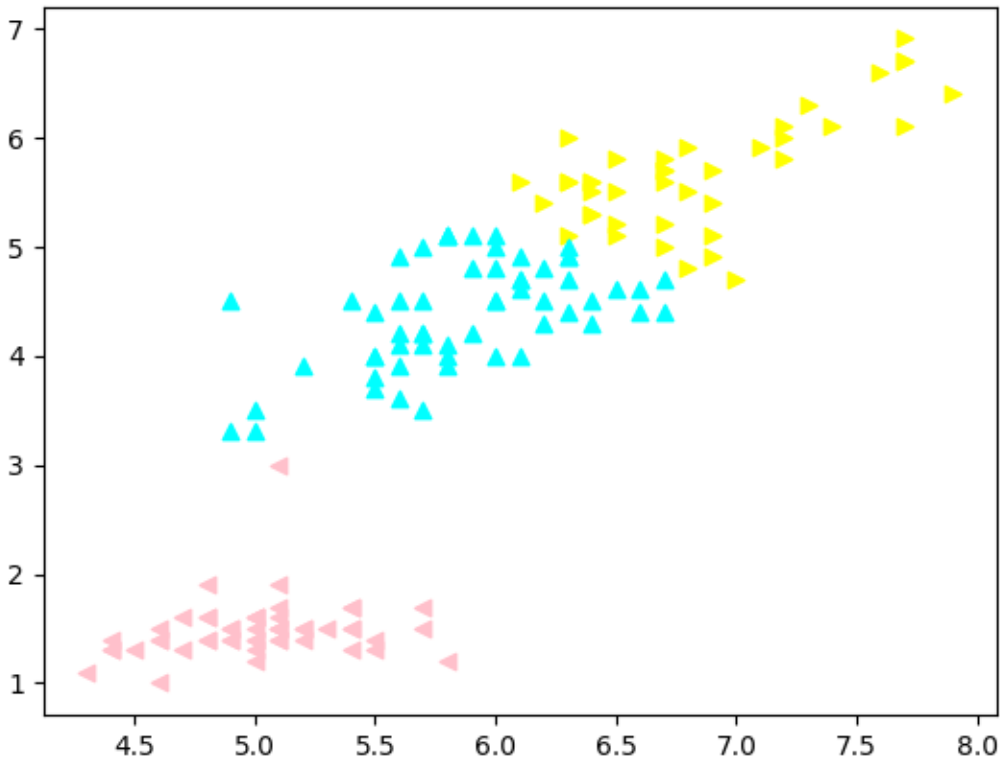
```
1 #copied from https://towardsdatascience.com/k-means-clustering-algorithm-applications-evaluation-methods-and-drawbacks-aa03e644b48a
2 import matplotlib.pyplot as plt
3 from matplotlib.image import imread
4 from sklearn.cluster import KMeans, SpectralClustering
5 from sklearn.preprocessing import StandardScaler
6 from sklearn.metrics import silhouette_samples, silhouette_score
7 from sklearn.datasets import load_iris
8 import numpy as np
```

2.1 *10 points*. 0.5 hour. Using $K = 3$, cluster the entire dataset into 3 labels using only features 1 & 3; namely, sepal length and petal length (Note: the example in class used all 4 features for clustering). Show a scatter plot based on these 2 features using known training 3 classes using markers "<" for class 1 (Setosa), ">" for class 2 (Versicolor), and "^" for class 3 (Virginica) while also using colors based on the 3 computed clusters using colors of "pink" for cluster 1, "yellow" for cluster 2, and "cyan" for cluster 3.

```
1 # Load the IRIS dataset
2 iris = load_iris()
3 X = iris.data[:, [0, 2]] # Select features 1 and 3
4 y = iris.target
5
6 # Run the Kmeans algorithm using all 1 3 features
7 km = KMeans(n_clusters=3, n_init= "auto")
8 labels = km.fit_predict(X)
9 print(labels)
```

```
[00000000000000000000000000000000000000000000000000000  
000000000000000002121111111111111111111111111111111111  
1122111111111111111111111111111110121222212222  
2211222212121221122222222222212221222122212  
21]
```

```
1 # Plotting the results of clustering
2 # Plot of actual class labels using <, >, ^.
3 plt.scatter(X[labels == 0][:,0], X[labels == 0][:,1], marker = "<", color = 'pink')
4 plt.scatter(X[labels == 1][:,0], X[labels == 1][:,1], marker = ">", color = 'yellow')
5 plt.scatter(X[labels == 2][:,0], X[labels == 2][:,1], marker = "^", color = 'cyan')
6 plt.show()
```



2.2 5 points. Report based on known labels what percent is misclassified when using 2 features.

```
1 # Calculate misclassification percentage
2 misclassified = np.sum(labels != y)
3 misclassification_percentage = (misclassified / len(y)) * 100
4
5 print(f"Misclassification percentage: {misclassification_percentage:.2f}%")
```

Misclassification percentage: 55.33%

2.3 10 points. 0.5 hour. Plot the result of K Means clustering using all 4 features with $K = 4$.

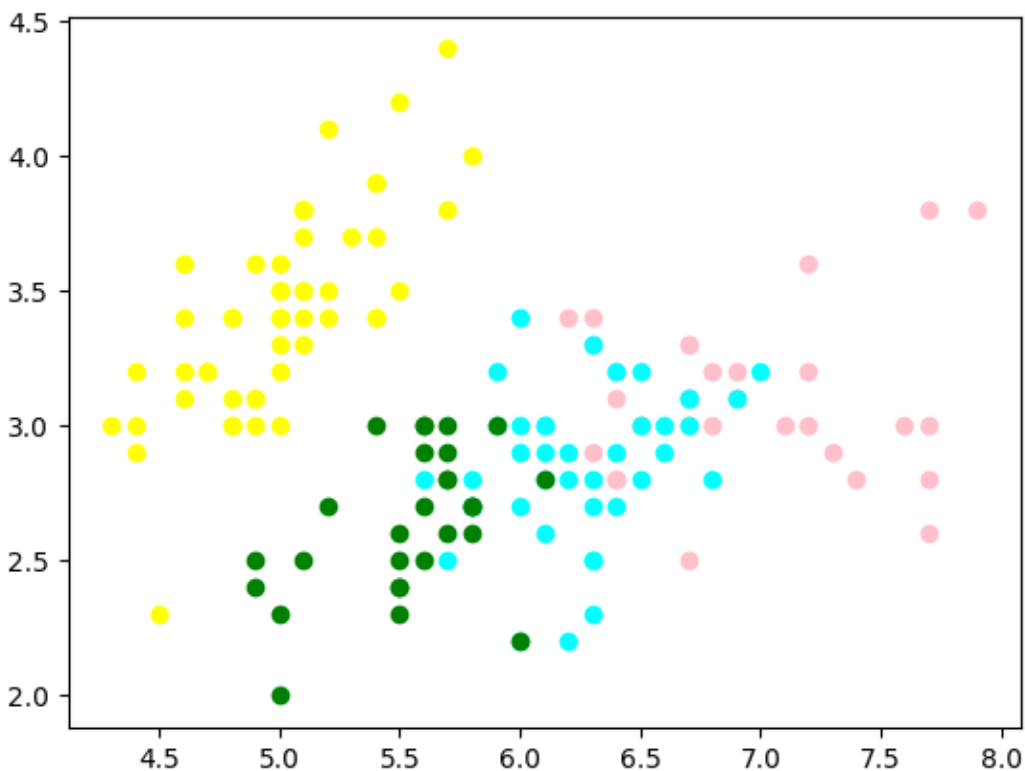
```
1 # Run the Kmeans algorithm using all 4 features
2 X = iris.data[:,:]
3 km = KMeans(n_clusters=4, n_init= "auto")
4 label = km.fit_predict(X)
5 print(label)
```

```
[1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 3 2 3 2 3 2 3 3 3 2 3 2 3 2 3 2 2 2  
2 2 2 2 2 3 3 3 3 2 3 2 2 2 3 3 3 2 3 3 3 3 3 2 3 3 0 2 0 0 0 0 3 0 0 0 2  
2 0 2 2 0 0 0 0 2 0 2 0 2 0 0 2 2 0 0 0 0 0 2 2 0 0 0 2 2 0 0 0 2 0 0 0 2 2  
0 2]
```

```

1 #Plotting the results of clustering
2 plt.scatter(X[label == 0][:,0] , X[label == 0][:,1] , color = 'pink')
3 plt.scatter(X[label == 1][:,0] , X[label == 1][:,1] , color = 'yellow')
4 plt.scatter(X[label == 2][:,0] , X[label == 2][:,1] , color = 'cyan')
5 plt.scatter(X[label == 3][:,0] , X[label == 3][:,1] , color = 'green')
6 plt.show()

```



2.4 15 points. 1 hour. Reduce the 4 features (sepal length, sepal width, petal length, petal width) into 2 PCA features (an example is also provided in class). Use $K = 3$ to cluster the entire dataset using these 2 PCA features. Show a scatter plot like in problem 2.1 along with percent misclassified as in problem 2.2.

```

1 # Apply PCA for dimensionality reduction
2 X = iris.data[:, :]
3 pca = PCA(n_components=2)
4 X_pca = pca.fit_transform(X)
5
6 k = 3
7
8 # Perform K-means clustering on PCA-transformed data
9 km = KMeans(n_clusters=k, n_init= "auto")
10 labels = km.fit_predict(X_pca)

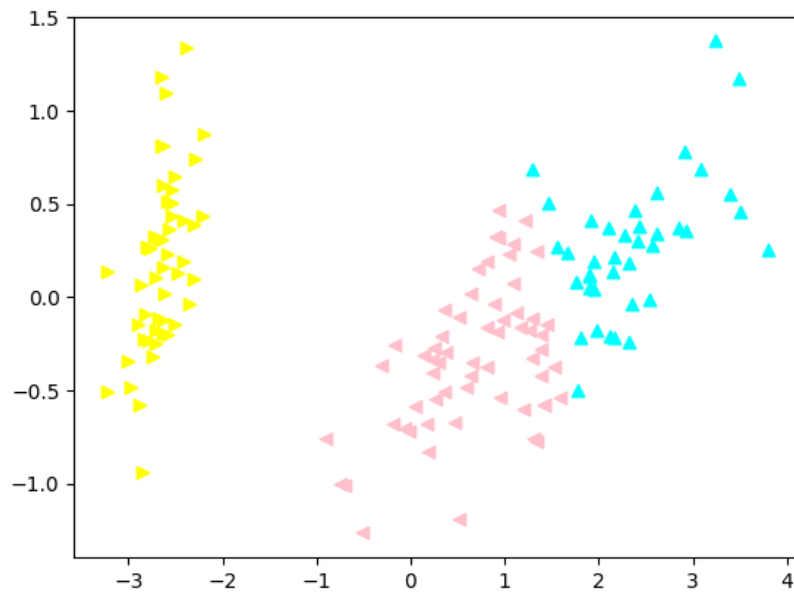
```



```

1 # Plotting the results of clustering
2 # Plot of actual class labels using <, >, ^.
3 plt.scatter(X_pca[labels == 0][:,0], X_pca[labels == 0][:,1], marker = "<", color = 'pink')
4 plt.scatter(X_pca[labels == 1][:,0], X_pca[labels == 1][:,1], marker = ">", color = 'yellow')
5 plt.scatter(X_pca[labels == 2][:,0], X_pca[labels == 2][:,1], marker = "^", color = 'cyan')
6 plt.show()

```



```

1 # Calculate misclassification percentage
2 misclassified = np.sum(labels != y)
3 misclassification_percentage = (misclassified / len(y)) * 100
4
5 print(f"Misclassification percentage: {misclassification_percentage:.2f}%")

```

Misclassification percentage: 76.00%

2.5 20 points. Redo the example in class with all 4 features and $K = 3$, but using your own class or function **my_k_means** in Python that has initialization parameters: K , X , max_iterations , $\text{centroid_move_epsilon}$ and returns y as a 1-D array of integer labels of 1, 2, ..., K .. Each input N -dimensional data $X[i]$ will have a 1-dimensional output label $y[i]$ for $i = 1..M$ where M is the number of data points. The algorithm should start by assigning K cluster centers based on random values from the (min, max) range of each dimension in the N -dimensional data X . It should stop when all centers have moved by less than the $\text{centroid_move_epsilon}$ or when the max_iterations is reached. Make sure your results are similar to the K Means library class.

```

1 def my_k_means(K, X, max_iterations, centroid_move_epsilon):
2     # Random initialization of cluster centers within the range of each dimension
3     min_vals = np.min(X, axis=0)
4     max_vals = np.max(X, axis=0)
5     centroids = np.random.uniform(min_vals, max_vals, size=(K, X.shape[1]))
6
7     # Iterative K-means algorithm
8     for _ in range(max_iterations):
9         # Assign each data point to the nearest centroid
10        distances = np.linalg.norm(X[:, np.newaxis] - centroids, axis=-1)
11        labels = np.argmin(distances, axis=-1)
12
13        # Move centroids to the mean of their assigned data points
14        new_centroids = np.array([np.mean(X[labels == k], axis=0) for k in range(K)])
15
16        # Check if centroids have moved by less than the threshold
17        centroid_moves = np.linalg.norm(new_centroids - centroids, axis=-1)
18        if np.all(centroid_moves < centroid_move_epsilon):
19            break
20
21        centroids = new_centroids
22
23    return labels

```

```

1 # Run my_k_means
2 K = 3
3 max_iterations = 100
4 centroid_move_epsilon = 1e-4
5 labels = my_k_means(K, X, max_iterations, centroid_move_epsilon)
6
7 print(labels)

```

```

[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 0 0 0 0 0 0 0 0 0 0 0 0 1 2 1 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
 2 2 2 1 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 1 2 1 1 1 2 1 1 1
 1 1 2 2 1 1 1 1 2 1 2 1 2 1 1 2 2 1 1 1 1 2 1 1 1 1 2 1 1 1 2 1
 1 2]

```

```

1 km = KMeans(n_clusters=3, n_init="auto")
2 labels = km.fit_predict(X)
3 print(labels)

```

```

[1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 1 1 1 1 1 1 1 1 1 1 1 2 2 0 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
 2 2 2 0 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 0 2 0 0 0 2 0 0 0
 0 0 2 2 0 0 0 2 0 2 0 2 0 0 2 2 0 0 0 0 2 0 0 0 0 2 0 0 0 2 0 0 0
 0 2]

```

3. 15 points. 1 hour. Decision Trees. Change the “IRIS Decision Tree.ipynb” shown in class, to use SKlearn’s Wine Recognition Dataset instead. Report the classification accuracy % for a single tree using 70% training samples and for a random forest with 100 estimators.

Link: https://colab.research.google.com/drive/1bpDNMyA3W_AyvNaVzaQLeBz8USJzR33N?usp=sharing

```
1 # Import scikit-learn dataset library
2 from sklearn import datasets
3
4 # Load the Wine dataset
5 wine = datasets.load_wine()
```

```
1 # Print the label species
2 print(wine.target_names)
3
4 # Print the names of the features
5 print(wine.feature_names)
```

```
['class_0', 'class_1', 'class_2']
['alcohol', 'malic_acid', 'ash', 'alkalinity_of_ash', 'magnesium', 'total_phenols', 'flavanoids', 'nonflavanoid_phenols', 'proanthocyanins', 'color_intensity', 'hue', 'od280/od315_of_diluted_wines', 'proline']
```

```
1 # Print the wine data (top 5 records)
2 print(wine.data[0:5])
3
4 # Print the wine labels
5 print(wine.target)
```

[illegible]

```

1 # Creating a DataFrame of the wine dataset
2 import pandas as pd
3 data = pd.DataFrame({
4     'alcohol': wine.data[:, 0],
5     'malic_acid': wine.data[:, 1],
6     'ash': wine.data[:, 2],
7     'alcalinity_of_ash': wine.data[:, 3],
8     'magnesium': wine.data[:, 4],
9     'total_phenols': wine.data[:, 5],
10    'flavanoids': wine.data[:, 6],
11    'nonflavanoid_phenols': wine.data[:, 7],
12    'proanthocyanins': wine.data[:, 8],
13    'color_intensity': wine.data[:, 9],
14    'hue': wine.data[:, 10],
15    'od280/od315_of_diluted_wines': wine.data[:, 11],
16    'proline': wine.data[:, 12],
17    'species': wine.target
18 })
19 data.head()

```

	alcohol	malic_acid	ash	alcalinity_of_ash	magnesium	total_phenols	flavanoids	nonflavanoid_phenols	proanthocyanins	color_intensity	hue	od280/od315_of_diluted_wines	proline	species
0	14.23	1.71	2.43	15.6	127.0	2.80	3.06	0.28	2.29	5.64	1.04	3.92	1065.0	0
1	13.20	1.78	2.14	11.2	100.0	2.65	2.76	0.26	1.28	4.38	1.05	3.40	1050.0	0
2	13.16	2.36	2.67	18.6	101.0	2.80	3.24	0.30	2.81	5.68	1.03	3.17	1185.0	0
3	14.37	1.95	2.50	16.8	113.0	3.85	3.49	0.24	2.18	7.80	0.86	3.45	1480.0	0
4	13.24	2.59	2.87	21.0	118.0	2.80	2.69	0.39	1.82	4.32	1.04	2.93	735.0	0

```

1 # Import train_test_split function
2 from sklearn.model_selection import train_test_split
3
4 X=data[['alcohol', 'malic_acid', 'ash', 'alcalinity_of_ash', 'magnesium',
5         'total_phenols', 'flavanoids', 'nonflavanoid_phenols', 'proanthocyanins',
6         'color_intensity', 'hue', 'od280/od315_of_diluted_wines', 'proline']] # Features
7 y=data['species'] # Labels
8
9 # Split the dataset into training set and test set
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3) # 70% training and 30% test

```

```

1 # Import the DecisionTreeClassifier and RandomForestClassifier Model
2 from sklearn.tree import DecisionTreeClassifier
3 from sklearn.ensemble import RandomForestClassifier
4
5 # Create a Decision Tree Classifier
6 clf_tree = DecisionTreeClassifier()
7
8 # Train the Decision Tree Classifier using the training sets
9 clf_tree.fit(X_train, y_train)
10
11 # Make predictions on the test set
12 y_pred_tree = clf_tree.predict(X_test)

```

```
1 #Import scikit-learn metrics module for accuracy calculation
2 from sklearn import metrics
3
4 # Calculate the accuracy of the Decision Tree Classifier
5 accuracy_tree = metrics.accuracy_score(y_test, y_pred_tree)
6 print("Decision Tree Accuracy:", accuracy_tree)
7
8 # Create a Random Forest Classifier with 100 estimators
9 clf_forest = RandomForestClassifier(n_estimators=100)
10
11 # Train the Random Forest Classifier using the training sets
12 clf_forest.fit(X_train, y_train)
13
14 # Make predictions on the test set
15 y_pred_forest = clf_forest.predict(X_test)
16
17 # Calculate the accuracy of the Random Forest Classifier
18 accuracy_forest = metrics.accuracy_score(y_test, y_pred_forest)
19 print("Random Forest Accuracy:", accuracy_forest)
```

Decision Tree Accuracy: 0.8888888888888888

Random Forest Accuracy: 0.9629629629629629