

CS178 Homework 1

Due: Monday, October 7th 2024 (11:59 PM)

Instructions

Welcome to CS 178!

This homework (and many subsequent ones) will involve data analysis and reporting on methods and results using Python code. You will submit a **single PDF file** that contains everything to Gradescope. This includes any text you wish to include to describe your results, the complete code snippets of how you attempted each problem, any figures that were generated, and scans of any work on paper that you wish to include. It is important that you include enough detail that we know how you solved the problem, since otherwise we will be unable to grade it.

Your homeworks will be given to you as Jupyter notebooks containing the problem descriptions and some template code that will help you get started. You are encouraged to modify these starter Jupyter notebooks to complete your assignment and to write your report. You may add additional cells (containing either code or text) as needed. This will help you not only ensure that all of the code for the solutions is included, but also will provide an easy way to export your results to a PDF file (for example, doing *print preview* and *printing to pdf*). I recommend liberal use of Markdown cells to create headers for each problem and sub-problem, explaining your implementation/answers, and including any mathematical equations. For parts of the homework you do on paper, scan it in such that it is legible (there are a number of free Android/iOS scanning apps, if you do not have access to a scanner), and include it as an image in the Jupyter notebook.

Several problems in this assignment require you to create plots. Use `matplotlib.pyplot` to do this, which is already imported for you as `plt`. Do not use any other plotting libraries, such as `pandas` or `seaborn`. Unless you are told otherwise, you should call `pyplot` plotting functions with their default arguments.

If you have any questions/concerns about the homework problems or using Jupyter notebooks, ask us on EdD. If you decide not to use Jupyter notebooks, but go with Microsoft Word or Latex to create your PDF file, make sure that all of the answers

can be generated from the code snippets included in the document.

Summary of Assignment: 100 total points

- Problem 1: Exploring a NYC Housing Dataset (30 points)
 - Problem 1.1: Numpy Arrays (5 points)
 - Problem 1.2: Feature Statistics (5 points)
 - Problem 1.3: Logical Indexing (5 points)
 - Problem 1.4: Histograms (5 points)
 - Problem 1.5: Scatter Plots (10 points)
- Problem 2: Building a Nearest Centroid Classifier (50 points)
 - Problem 2.1: Implementing Nearest Centroids (30 points)
 - Problem 2.2: Evaluating Nearest Centroids (20 points)
- Problem 3: Decision Boundaries (15 points)
 - Problem 3.1: Visualize 2D Centroid Classifier (5 points)
 - Problem 3.2: Visualize 2D Gaussian Bayes Classifier (5 points)
 - Problem 3.3: Analysis (5 points)
- Statement of Collaboration (5 points)



Before we get started, let's import some libraries that you will make use of in this assignment. Make sure that you run the code cell below in order to import these libraries.

Important: In the code block below, we set `seed=123` . This is to ensure your code has reproducible results and is important for grading. Do not change this. If you are not using the provided Jupyter notebook, make sure to also set the random seed as below.

```
In [42]: # If you haven't installed numpy, pyplot, scikit, etc., do so:
!pip3 install -U scikit-learn
```

Requirement already satisfied: scikit-learn in /Library/Frameworks/Python.framework/Versions/3.10/lib/python3.10/site-packages (1.5.2)
Requirement already satisfied: numpy>=1.19.5 in /Users/adityasingh/Library/Python/3.10/lib/python/site-packages (from scikit-learn) (1.24.1)
Requirement already satisfied: scipy>=1.6.0 in /Library/Frameworks/Python.framework/Versions/3.10/lib/python3.10/site-packages (from scikit-learn) (1.12.0)
Requirement already satisfied: joblib>=1.2.0 in /Library/Frameworks/Python.framework/Versions/3.10/lib/python3.10/site-packages (from scikit-learn) (1.2.0)
Requirement already satisfied: threadpoolctl>=3.1.0 in /Library/Frameworks/Python.framework/Versions/3.10/lib/python3.10/site-packages (from scikit-learn) (3.3.0)

[notice] A new release of pip is available: 24.0 -> 24.2

[notice] To update, run: `python3 -m pip install --upgrade pip`

```
In [43]: import numpy as np
import matplotlib.pyplot as plt

from sklearn.datasets import fetch_openml
from sklearn.neighbors import KNeighborsClassifier, NearestCentroid
from sklearn.metrics import accuracy_score, zero_one_loss, confusion_matrix, ConfusionMatrixDisplay
from sklearn.model_selection import train_test_split
from sklearn.inspection import DecisionBoundaryDisplay

import requests          # we'll use these for reading data from a url
from io import StringIO

# Fix the random seed for reproducibility
# !! Important !! : do not change this
seed = 123
np.random.seed(seed)
```

Problem 1: Exploring a NYC Housing Dataset

In this problem, you will explore some basic data manipulation and visualizations with a small dataset of real estate prices from NYC. For every datapoint, we are given several real-valued features which will be used to predict the target variable, y , representing in which borough the property is located. Let's first load in the dataset by running the code cell below:

```
In [44]: # Load the features and labels from an online text file
url = 'https://ics.uci.edu/~ihler/classes/cs178/data/nyc_housing.txt'
with requests.get(url) as link:
    datafile = StringIO(link.text)
    nych = np.genfromtxt(datafile, delimiter=',')
    nych_X, nych_y = nych[:, :-1], nych[:, -1]
```

These data correspond to (a small subset of) property sales in New York in 2014. The target, y , represents the borough in which the property was located (0: Manhattan; 1: Bronx; 2: Staten Island). The observed features correspond to the property size (square feet), price (USD), and year built; the first two features have been log2-transformed (e.g., $x_1 = \log_2(\text{size})$) for convenience.

Problem 1.1 (5 points): Numpy Arrays

The variable `nych_X` is a numpy array containing the feature vectors in our dataset, and `nych_y` is a numpy array containing the corresponding labels.

- What is the shape of `nych_X` and `nych_y`? (Hint)
- How many datapoints are in our dataset, and how many features does each datapoint have?
- How many different classes (i.e. labels) are there?
- Print rows 3, 4, 5, and 6 of the feature matrix and their corresponding labels. Since Python is zero-indexed, we will count our rows starting at zero -- for example, by "row 0" we mean `nych_X[0, :]`, and "row 1" means `nych_X[1, :]`, etc. (Hint: you can do this in two lines of code with slicing).

```
In [45]: #Answer1
print("Answer 1:")
nychX_shape = nych_X.shape
nychY_shape = nych_y.shape
print(f"Shape of nych_X is: {nychX_shape}")
print(f"Shape of nych_y is: {nychY_shape}")

#Answer2
print("Answer 2:")
datapoints_nychX, features_nychX = nychX_shape
datapoints_nychY = nychY_shape
print(f"nych_X feature vector has {datapoints_nychX} datapoints and {features_nychX} features each")
```

```
print(f"nych_Y label array has {datapoints_nychY} datapoints and 0 features each")

#Answer3
num_diff_classes = np.unique(nych_y).size
print(f"The label array has {num_diff_classes} unique labels")

#Answer4
print("Rows 3, 4, 5, and 6 of feature matrix:")
print(nych_X[3:7,:])
print("Corresponding labels:")
print(nych_y[3:7])
```

Answer 1:

Shape of nych_X is: (300, 3)

Shape of nych_y is: (300,)

Answer 2:

nych_X feature vector has 300 datapoints and 3 features each

nych_Y label array has (300,) datapoints and 0 features each

The label array has 3 unique labels

Rows 3, 4, 5, and 6 of feature matrix:

```
[ [ 11.839204  19.416995 1980.    ]
  [ 18.517396  25.357833 1973.    ]
  [ 11.050529  19.041723 2014.    ]
  [ 17.255803  26.280297 1917.    ]]
```

Corresponding labels:

```
[2. 1. 2. 0.]
```

Problem 1.2 (5 points): Feature Statistics

Let's compute some statistics about our features. You are allowed to use `numpy` to help you with this problem -- for example, you might find some of the `numpy` functions listed [here](#) or [here](#) useful.

- Compute the mean, variance, and standard deviation of each feature.
- Compute the minimum and maximum value for each feature.

Make sure to print out each of these values, and indicate clearly which value corresponds to which computation.

```
In [46]: #Mean:
property_size_mean = np.mean(nych_X[:, 0])
property_price_mean = np.mean(nych_X[:, 1])
```

```
property_year_mean = np.mean(nych_X[:, 2])
print(f"Mean of property sizes: {property_size_mean}")
print(f"Mean of property prices: {property_price_mean}")
print(f"Mean of property built year: {property_year_mean}")

#Variance:
property_size_variance = np.var(nych_X[:, 0])
property_price_variance = np.var(nych_X[:, 1])
property_year_variance = np.var(nych_X[:, 2])
print(f"Variance of property sizes: {property_size_variance}")
print(f"Variance of property prices: {property_price_variance}")
print(f"Variance of property built year: {property_year_variance}")

#STD:
property_size_std = np.std(nych_X[:, 0])
property_price_std = np.std(nych_X[:, 1])
property_year_std = np.std(nych_X[:, 2])
print(f"Standard deviation of property sizes: {property_size_std}")
print(f"Standard deviation of property prices: {property_price_std}")
print(f"Standard deviation of property built year: {property_year_std}")

#Minimum:
property_size_min = np.min(nych_X[:, 0])
property_price_min = np.min(nych_X[:, 1])
property_year_min = np.min(nych_X[:, 2])
print(f"Minimum of property sizes: {property_size_min}")
print(f"Minimum of property prices: {property_price_min}")
print(f"Minimum of property built year: {property_year_min}")

#Maximum:
property_size_max = np.max(nych_X[:, 0])
property_price_max = np.max(nych_X[:, 1])
property_year_max = np.max(nych_X[:, 2])
print(f"Maximum of property sizes: {property_size_max}")
print(f"Maximum of property prices: {property_price_max}")
print(f"Maximum of property built year: {property_year_max}")
```

```

Mean of property sizes: 14.118392473333333
Mean of property prices: 21.907116153333334
Mean of property built year: 1946.3533333333332
Variance of property sizes: 6.60022491794569
Variance of property prices: 8.871930118164771
Variance of property built year: 1253.0818222222222
Standard deviation of property sizes: 2.5690902899559
Standard deviation of property prices: 2.9785785398684337
Standard deviation of property built year: 35.39889577687731
Minimum of property sizes: 10.366322
Minimum of property prices: 16.872675
Minimum of property built year: 1893.0
Maximum of property sizes: 20.152714
Maximum of property prices: 29.123861
Maximum of property built year: 2014.0

```

Problem 1.3 (5 points): Logical Indexing

Use numpy's logical (boolean) indexing to extract only those data corresponding to $y = 0$ (Manhattan). Then, compute the mean and standard deviation of *only these* data points. Then, do the same for $y = 1$ (Bronx).

Again, print out each of these vectors and indicate clearly which value corresponds to which computation.

```

In [47]: # Manhattan(y = 0):

manhattan_rows = nych_X[nych_y == 0]
mean_manhattan = np.mean(manhattan_rows, axis=0)
std_manhattan = np.std(manhattan_rows, axis=0)
print(f"mean of all manhattan columns: {mean_manhattan} and std of all columns: {std_manhattan}")

# Bronx(y = 1):

bronx_rows = nych_X[nych_y == 1]
mean_bronx = np.mean(bronx_rows, axis=0)
std_bronx = np.std(bronx_rows, axis=0)
print(f"mean of all bronx columns: {mean_bronx} and std of all columns: {std_bronx}")

```

mean of all manhattan columns: [16.1489863 25.07251963 1926.94] and std of all columns: [2.1941 6051 2.09812353 28.14562843]
mean of all bronx columns: [14.60837771 21.4446885 1935.29] and std of all columns: [1.89678446 1.99063026 22.96619037]

Problem 1.4 (5 points): Feature Histograms

Now, you will visualize the distribution of each feature with histograms. Use `matplotlib.pyplot` to do this, which is already imported for you as `plt`. Do not use any other plotting libraries, such as `pandas` or `seaborn`.

- For every feature in `nych_X`, plot a histogram of the values of the feature. Your plot should consist of a grid of subplots with 1 row and 3 columns.
- Include a title above each subplot to indicate which feature we are plotting. For example, you can call the first feature "Feature 0", the second feature "Feature 1", etc.

Some starter code is provided for you below. (Hint: `axes[0].hist(...)` will create a histogram in the first subplot.)

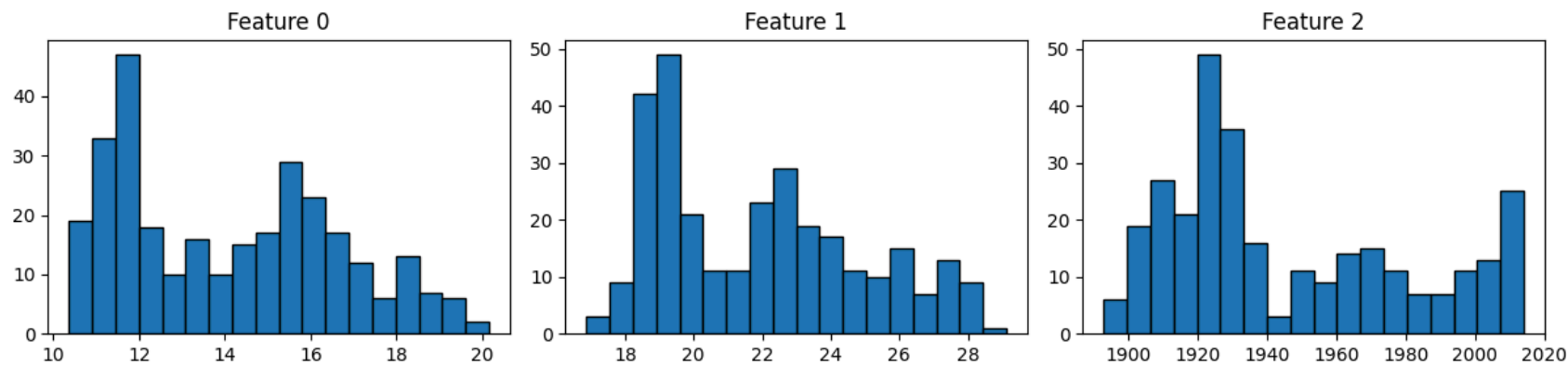
```
In [48]: # Create a figure with 1 row and 3 columns
fig, axes = plt.subplots(1, 3, figsize=(12, 3))

### YOUR CODE STARTS HERE ###

for i in range(nych_X.shape[1]):
    axes[i].hist(nych_X[:, i], bins = 18, edgecolor='black')
    axes[i].set_title(f'Feature {i}')

### YOUR CODE ENDS HERE ###

fig.tight_layout()
```

Problem 1.5 (10 points): Feature Scatter Plots

To help further visualize the NYC-Housing dataset, you will now create several scatter plots of the features. Use `matplotlib.pyplot` to do this, which is already imported for you as `plt`. Do not use any other plotting libraries, such as `pandas` or `seaborn`.

- For every pair of features in `nyc_h_X`, plot a scatter plot of the feature values, colored according to their labels. For example, plot all data points with $y = 0$ as blue, $y = 1$ as green, etc. Your plot should be a grid of subplots with 3 rows and 3 columns. (Hint: `axes[0, 0].scatter(...)` will create a scatter plot in the first column and first row).
- Include an x-label and a y-label on each subplot to indicate which features we are plotting. For example, you can call the first feature "Feature 0", the second feature "Feature 1", etc. (Hint: `axes[0, 0].set_xlabel(...)` might help you with the first subplot.)

Some starter code is provided for you below.

```
In [49]: # Create a figure with 3 rows and 3 columns
fig, axes = plt.subplots(3, 3, figsize=(8, 8))

### YOUR CODE STARTS HERE ###

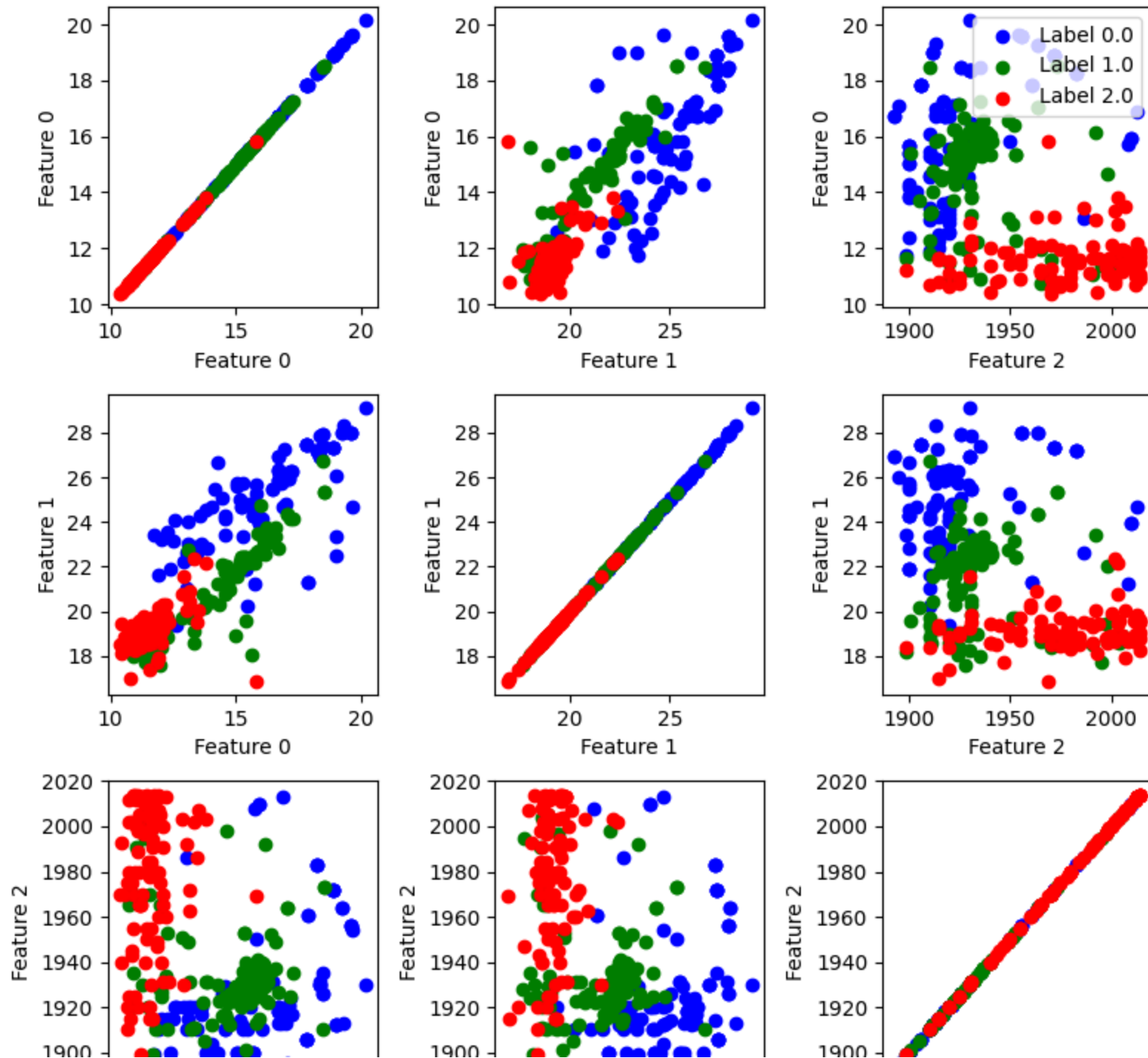
colors = {0: 'blue', 1: 'green', 2: 'red'}

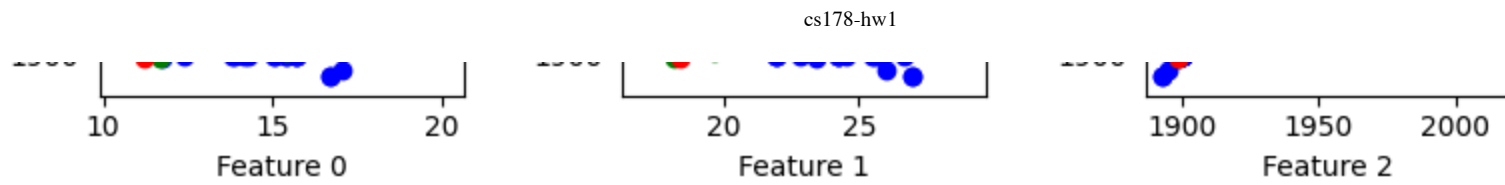
for i in range(3):
    for j in range(3):
        for label in np.unique(nyc_h_y):
```

```
        axes[i, j].scatter(nych_X[nych_y == label, j], nych_X[nych_y == label, i],
                           color=colors[label], label=f'Label {label}')
    axes[i, j].set_xlabel(f'Feature {j}')
    axes[i, j].set_ylabel(f'Feature {i}')
    if i == 0 and j == 2:
        axes[i, j].legend()

### YOUR CODE ENDS HERE ###

fig.tight_layout()
```





Problem 2: Nearest Centroid Classifiers

In this problem, you will implement a nearest centroid classifier and train it on the NYC data.

Problem 2.1 (30 points): Implementing a Nearest Centroid Classifier

In the code given below, we define the class `NearestCentroidClassifier` which has an unfinished implementation of a nearest centroid classifier. For this problem, you will complete this implementation. Your nearest centroid classifier will use the Euclidean distance, which is defined for two feature vectors x and x' as

$$d_E(x, x') = \sqrt{\sum_{j=1}^d (x_j - x'_j)^2}.$$

- Implement the method `fit`, which takes in an array of features `X` and an array of labels `y` and trains our classifier. You should store your computed centroids in the list `self.centroids`, and their y values in `self.classes_` (whose name is chosen to match `sklearn` conventions).
- Test your implementation of `fit` by training a `NearestCentroidClassifier` on the NYC data, and printing out the list of centroids. (These should match the means in Problem 1.3.)
- Implement the method `predict`, which takes in an array of feature vectors `X` and predicts their class labels based on the centroids you computed in the method `fit`.
- Print the predicted labels (using your `predict` function) and the true labels for the first ten data points in the NYCH dataset. Make sure to indicate which are the predicted labels and which are the true labels.

You are allowed to modify the given code as necessary to complete the problem, e.g. you may create helper functions.

```
In [50]: class NearestCentroidClassifier:
    def __init__(self):
        # A list containing the centroids; to be filled in with the fit method.
        self.centroids = []

    def fit(self, X, y):
        """ Fits the nearest centroid classifier with training features X and training labels y.

        X: array of training features; shape (m,n), where m is the number of datapoints,
           and n is the number of features.
        y: array training labels; shape (m, ), where m is the number of datapoints.

        """
        # First, identify what possible classes exist in the training data set:
```

```
self.classes_ = np.unique(y)

### YOUR CODE STARTS HERE ###
# Hint: you should append to self.centroids with the corresponding centroid for each class.
# The centroid (mean vector) can be computed in a similar way to P1.2, for example.

for i in self.classes_:
    X_temp = X[y == i]
    self.centroids.append(np.mean(X_temp, axis=0))

### YOUR CODE ENDS HERE ###

def predict(self, X):
    """ Makes predictions with the nearest centroid classifier on the features in X.

    X: array of features; shape (m,n), where m is the number of datapoints,
        and n is the number of features.

    Returns:
    y_pred: a numpy array of predicted labels; shape (m, ), where m is the number of datapoints.
    """
    ### YOUR CODE STARTS HERE ###
    # Hint: find the distance from each x[i] to the centroids, and predict the closest.
    y_pred = []
    for x in X:
        distanceToCentroid = [np.linalg.norm(x - centroid) for centroid in self.centroids]
        closest = self.classes_[np.argmin(distanceToCentroid)]
        y_pred.append(closest)

    ### YOUR CODE ENDS HERE ###
```

```
return y_pred
```

Here is some code illustrating how to use your `NearestCentroidClassifier`. You can run this code to fit your classifier and to plot the centroids. You should write your implementation above such that you don't need to modify the code in the next cell. As a sanity check, you should find that the 3rd centroid (for Staten Island) has a "year build" coordinate value of around 1976.8 (i.e., the rightmost column).

```
In [51]: nc_classifier = NearestCentroidClassifier() # Create a NearestCentroidClassifier object
nc_classifier.fit(nych_X, nych_y)                  # Fit to the NYC training data

print(nc_classifier.centroids)

[array([ 16.1489863 ,  25.07251963, 1926.94      ]), array([ 14.60837771,  21.4446885 , 1935.29
]), array([ 11.59781341,  19.20414033, 1976.83      ])]
```

```
In [52]: # Print the predicted and true labels for the first ten data points in the NYCH testing set
        ### YOUR CODE STARTS HERE ###
```

```
predicted = nc_classifier.predict(nych_X[:10])
trueLabels = nych_y[:10]

for i in range(10):
    print(f"predicted label : {predicted[i]} -> true label : {trueLabels[i]}")
```

```
### YOUR CODE ENDS HERE ###
```

```

predicted label : 0.0 -> true label : 1.0
predicted label : 2.0 -> true label : 2.0
predicted label : 0.0 -> true label : 0.0
predicted label : 2.0 -> true label : 2.0
predicted label : 2.0 -> true label : 1.0
predicted label : 2.0 -> true label : 2.0
predicted label : 0.0 -> true label : 0.0
predicted label : 0.0 -> true label : 0.0
predicted label : 2.0 -> true label : 1.0
predicted label : 1.0 -> true label : 1.0

```

Problem 2.2 (20 points): Evaluating the Nearest Centroids Classifier

Now that you've implemented the nearest centroid classifier, it is time to evaluate its performance.

- Write a function `compute_error_rate` that computes the error rate (fraction of misclassifications) of a model's predictions. That is, your function should take in an array of true labels `y` and an array of predicted labels `y_pred`, and return the error rate of the predictions. You may use `numpy` to help you do this, but do not use `sklearn` or any other machine learning libraries.
- Write a function `compute_confusion_matrix` that computes the confusion matrix of a model's predictions. That is, your function should take in an array of true labels `y` and an array of predicted labels `y_pred`, and return the corresponding $C \times C$ confusion matrix as a numpy array, where C is the number of classes. You may use `numpy` to help you do this, but do not use `sklearn` or any other machine learning libraries.
- Verify that your implementations of `NearestCentroidClassifier`, `compute_error_rate`, and `compute_confusion_matrix` are correct. To help you do this, you are given the functions `eval_sklearn_implementation` and `eval_my_implementation`. The function `eval_sklearn_implementation` will use the relevant `sklearn` implementations to compute the error rate and confusion matrix of a nearest centroid classifier. The function `eval_my_implementation` will do the same, but using your implementations. If your code is correct, the outputs of the two functions should be the same.

```

In [53]: def compute_error_rate(y, y_pred):
          """ Computes the error rate of an array of predictions.

          y: true labels; shape (n, ), where n is the number of datapoints.
          y_pred: predicted labels; shape (n, ), where n is the number of datapoints.

```



```

Returns:
error_rate: the error rate of y_pred compared to y; scalar expressed as a decimal (e.g. 0.5)
"""
### YOUR CODE STARTS HERE ###

error_rate = np.mean(y != y_pred)

### YOUR CODE ENDS HERE ###

return error_rate

```

```

In [54]: def compute_confusion_matrix(y, y_pred):
        """ Computes the confusion matrix of an array of predictions.

        y: true labels; shape (n, ), where n is the number of datapoints.
        y_pred: predicted labels; shape (n, ), where n is the number of datapoints.

        Returns:
        confusion_matrix: a numpy array corresponding to the confusion matrix from y and y_pred; shape (C, C),
            where C is the number of unique classes. The (i,j)th entry is the number of examples of class i
            that are classified as being from class j.
        """

        ### YOUR CODE STARTS HERE ###
        classes = np.unique(np.concatenate((y, y_pred)))
        C = len(classes)
        class_index = {label : index for index, label in enumerate(classes)}

        confusion_matrix = np.zeros((C, C), dtype=int)

        for trueLabel, predLabel in zip(y, y_pred):
            trueidx = class_index[trueLabel]
            predidx = class_index[predLabel]

            confusion_matrix[trueidx, predidx] += 1

        ### YOUR CODE ENDS HERE ###

```

```
return confusion_matrix
```

You can run the two code cells below to compare your answers to the implementations in `sklearn`. If your answers are correct, the outputs of these two functions should be the same. Do not modify the functions

`eval_sklearn_implementation` and `eval_my_implementation`, but make sure that you read and understand this code.

```
In [55]: #####
### Results with the sklearn implementation ###
#####

def eval_sklearn_implementation(X, y):
    # Nearest centroid classifier implemented in sklearn
    sklearn_nearest_centroid = NearestCentroid()

    # Fit on training dataset
    sklearn_nearest_centroid.fit(X, y)

    # Make predictions on training and testing data
    sklearn_y_pred = sklearn_nearest_centroid.predict(X)

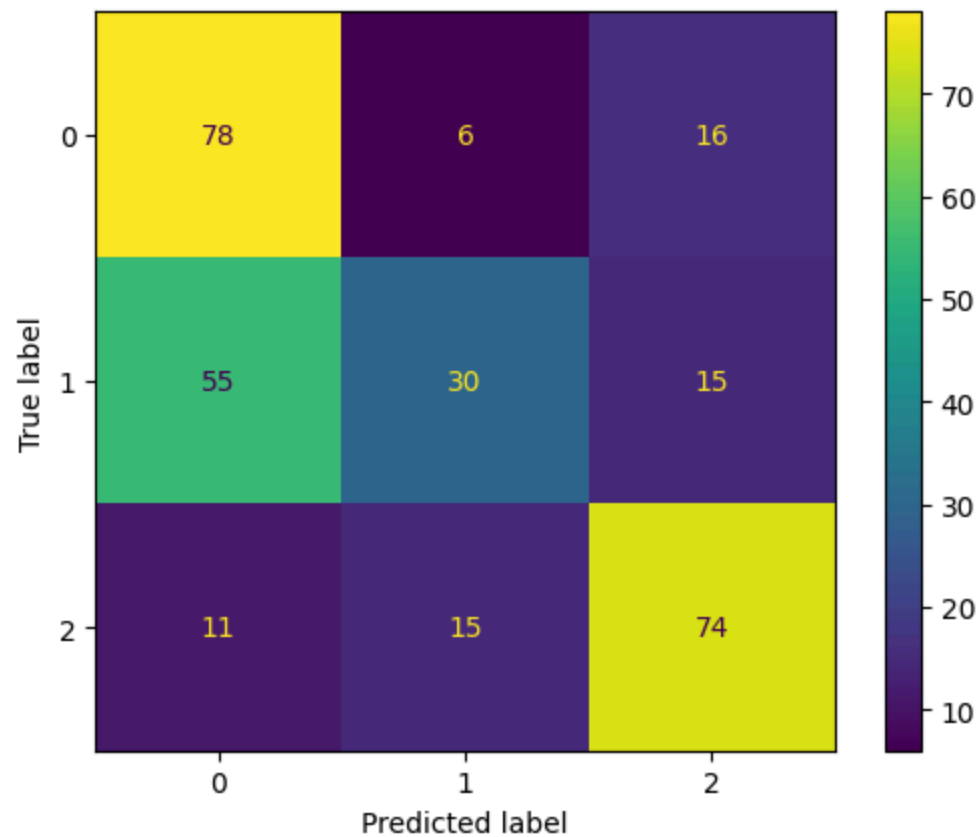
    # Evaluate accuracies using the sklearn function accuracy_score
    sklearn_err = zero_one_loss(y, sklearn_y_pred)

    print(f'Sklearn Results:')
    print(f'--- Error Rate (0/1): {sklearn_err}')

    # Evaluate confusion matrix using the sklearn function confusion_matrix
    sklearn_cm = confusion_matrix(y, sklearn_y_pred)
    sklearn_disp = ConfusionMatrixDisplay(confusion_matrix = sklearn_cm)
    sklearn_disp.plot();

# Call the function
eval_sklearn_implementation(nych_X, nych_y)
```

```
Sklearn Results:
--- Error Rate (0/1): 0.3933333333333333
```



```
In [56]: #####
### Results with your implementation ###
#####

def eval_my_implementation(X, y):
    # Now test your implementation of NearestCentroidClassifier
    nearest_centroid = NearestCentroidClassifier()

    # Fit on training dataset
    nearest_centroid.fit(X, y)

    # Make predictions on training and testing data
    y_pred = nearest_centroid.predict(X)

    # Evaluate accuracies using your function compute_accuracy
```

```
err = zero_one_loss(y, y_pred)

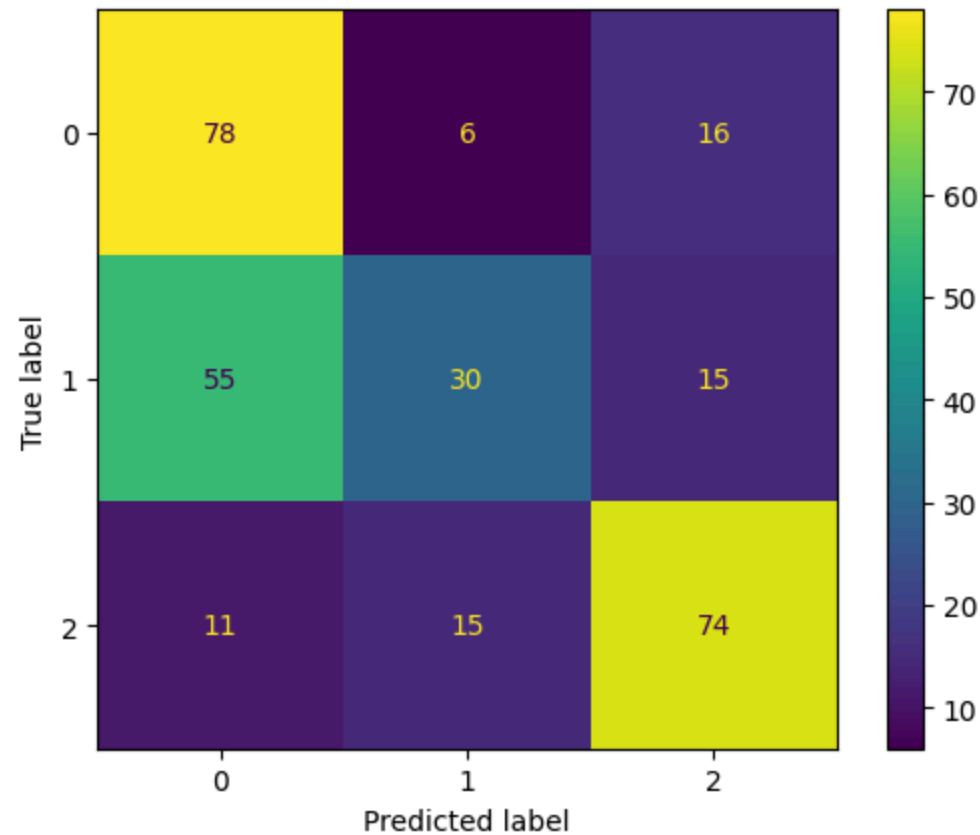
print(f'Your Results:')
print(f'--- Error Rate (0/1): {err}')

# Evaluate confusion matrix using your function compute_confusion_matrix
cm = compute_confusion_matrix(y, y_pred)
disp = ConfusionMatrixDisplay(confusion_matrix = cm)
disp.plot()

# Call the function
eval_my_implementation(nych_X, nych_y)
```

Your Results:

--- Error Rate (0/1): 0.3933333333333333





Problem 3: Decision Boundaries

For the final problem of this homework, you will visualize the decision function and decision boundary of your nearest centroid classifier on 2D data, and compare it to the similar but more flexible Gaussian Bayes classifier discussed in class. Code for drawing the decision function (which simply evaluates the prediction on a grid) and superimposing the data points is provided.

Problem 3.1 (5 points): Visualize 2D Centroid Classifier

We will use only the first two features of the NYCH data set, to facilitate visualization.

```
In [57]: # Plot the decision boundary for your classifier

# Some keyword arguments for making nice looking plots.
plot_kwargs = {'cmap': 'jet',      # another option: viridis
               'response_method': 'predict',
               'plot_method': 'pcolormesh',
               'shading': 'auto',
               'alpha': 0.5,
               'grid_resolution': 100}

figure, axes = plt.subplots(1, 1, figsize=(4,4))

learner = NearestCentroidClassifier()

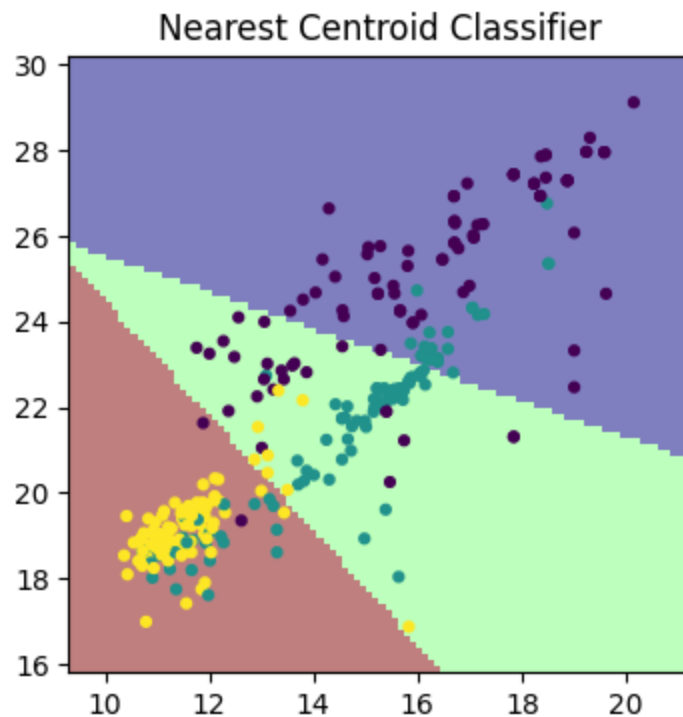
### YOUR CODE STARTS HERE ###
nych_first_two = nych_X[:, :2]

# get just the first two features of X
nych_X2 = learner.fit(nych_first_two, nych_y)
    # Fit "learner" to nych 2-feature data

### YOUR CODE ENDS HERE ###
```

```
DecisionBoundaryDisplay.from_estimator(learner, nych_first_two, ax=axes, **plot_kwargs)
axes.scatter(nych_first_two[:, 0], nych_first_two[:, 1], c=nych_y, edgecolor=None, s=12)
axes.set_title(f'Nearest Centroid Classifier')
```

Out[57]: Text(0.5, 1.0, 'Nearest Centroid Classifier')



Problem 3.2 (5 points): Visualize a 2D Gaussian Bayes Classifier

In class, we discussed building a Bayes classifier using an estimate of the class-conditional probabilities $p(X|Y = y)$, for example, a Gaussian distribution. It turns out this is relatively easy to implement and fairly similar to your Nearest Centroid classifier (in fact, Nearest Centroid is a special case of this model).

An implementation of a Gaussian Bayes classifier is provided:

```
In [58]: class GaussianBayesClassifier:
         def __init__(self):
```

```

"""Initialize the Gaussian Bayes Classifier"""
self.pY = []          # class prior probabilities, p(Y=c)
self.pXgY = []        # class-conditional probabilities, p(X|Y=c)
self.classes_ = []    # list of possible class values

def fit(self, X, y):
    """ Fits a Gaussian Bayes classifier with training features X and training labels y.
        X, y : (m,n) and (m,) arrays of training features and target class values
    """
    from sklearn.mixture import GaussianMixture
    self.classes_ = np.unique(y)      # Identify the class labels; then
    for c in self.classes_:          # for each class:
        self.pY.append(np.mean(y==c)) # estimate p(Y=c) (a float)
        model_c = GaussianMixture(1)  #
        model_c.fit(X[y==c,:])        # and a Gaussian for p(X|Y=c)
        self.pXgY.append(model_c)     #

def predict(self, X):
    """ Makes predictions with the nearest centroid classifier on the features in X.
        X : (m,n) array of features for prediction
        Returns: y : (m,) numpy array of predicted labels
    """
    pXY = np.stack(tuple(np.exp(p.score_samples(X)) for p in self.pXgY)).T
    pXY *= np.array(self.pY).reshape(1,-1) # evaluate p(X=x|Y=c) * p(Y=c)
    pYgX = pXY/pXY.sum(1,keepdims=True)    # normalize to p(Y=c|X=x) (not required)
    return self.classes_[np.argmax(pYgX, axis=1)] # find the max index & return its class ID

```

Using this learner, evaluate the predictions and error rate on the training data, and plot the decision boundary. The code should be the same as your Nearest Centroid, but using the new learner object.

```

In [59]: # Plot the decision boundary for your classifier

# Some keyword arguments for making nice looking plots.
plot_kwargs = {'cmap': 'jet',      # another option: viridis
               'response_method': 'predict',
               'plot_method': 'pcolormesh',
               'shading': 'auto',
               'alpha': 0.5,
               'grid_resolution': 100}

figure, axes = plt.subplots(1, 1, figsize=(4,4))

```

```
learner = GaussianBayesClassifier()

### YOUR CODE STARTS HERE ###
nych_first_two = nych_X[:, :2]

# get just the first two features of X
nych_X2 = learner.fit(nych_first_two, nych_y)
# Fit "learner" to nych 2-feature data

gbc_y_pred = learner.predict(nych_first_two)
# Use "learner" to predict on same data used in training

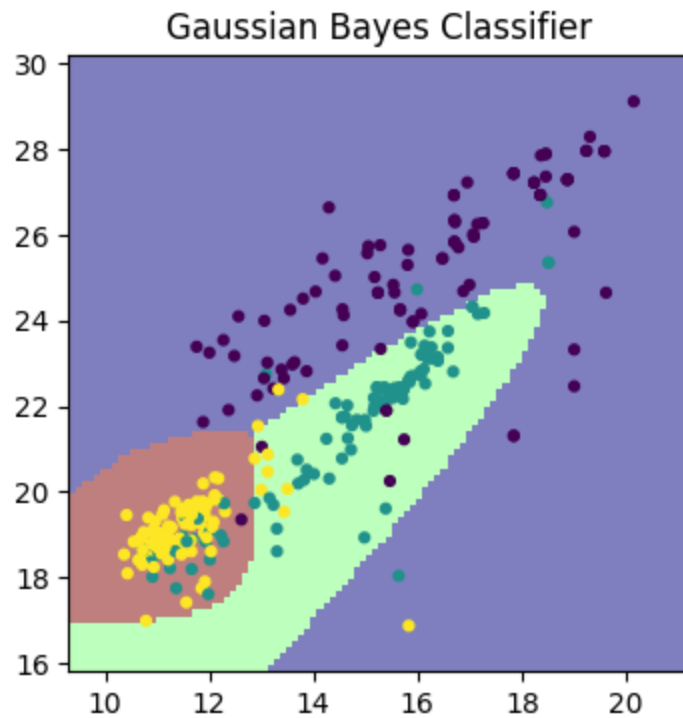
### YOUR CODE ENDS HERE ###

err = zero_one_loss(nych_y, gbc_y_pred)
print(f'Gaussian Bayes Error Rate (0/1): {err}')

DecisionBoundaryDisplay.from_estimator(learner, nych_first_two, ax=axes, **plot_kwargs)
axes.scatter(nych_first_two[:, 0], nych_first_two[:, 1], c=nych_y, edgecolor=None, s=12)
axes.set_title(f'Gaussian Bayes Classifier')
```

Gaussian Bayes Error Rate (0/1): 0.15000000000000002

Out[59]: Text(0.5, 1.0, 'Gaussian Bayes Classifier')



Problem 3.3 (5 points): Analysis

Did the error increase or decrease? Why do you think this is?

The error decreased, Gaussian Bayes reduces the error rate compared to Nearest Centroid because it models feature distributions, capturing non-linear decision boundaries, while Nearest Centroid assumes linear separability and ignores variance in the data.



Statement of Collaboration (5 points)

It is **mandatory** to include a Statement of Collaboration in each submission, with respect to the guidelines below. Include the names of everyone involved in the discussions (especially in-person ones), and what was discussed.

All students are required to follow the academic honesty guidelines posted on the course website. For programming assignments, in particular, I encourage the students to organize (perhaps using EdD) to discuss the task descriptions, requirements, bugs in my code, and the relevant technical content before they start working on it. However, you should not discuss the specific solutions, and, as a guiding principle, you are not allowed to take anything written or drawn away from these discussions (i.e. no photographs of the blackboard, written notes, referring to EdD, etc.). Especially after you have started working on the assignment, try to restrict the discussion to EdD as much as possible, so that there is no doubt as to the extent of your collaboration.

Self-completed