

Assignment 3 - Supervised Learning: model training and evaluation

DANIELA JIMENEZ LARA

Netid: dj216

Names of students you worked with on this assignment: Bárbara Flores

Note: this assignment falls under collaboration Mode 2: Individual Assignment – Collaboration Permitted. Please refer to the syllabus for additional information.

Instructions for all assignments can be found [here](#).

Total points in the assignment add up to 90; an additional 10 points are allocated to presentation quality.

Learning Objectives:

This assignment will provide structured practice to help enable you to...

1. Understand the primary workflow in machine learning: (1) identifying a hypothesis function set of models, (2) determining a loss/cost/error/objective function to minimize, and (3) minimizing that function through gradient descent
2. Understand the inner workings of logistic regression and how linear models for classification can be developed.
3. Gain practice in implementing machine learning algorithms from the most basic building blocks to understand the math and programming behind them to achieve practical proficiency with the techniques
4. Implement batch gradient descent and become familiar with how that technique is used and its dependence on the choice of learning rate
5. Evaluate supervised learning algorithm performance through ROC curves and using cross validation
6. Apply regularization to linear models to improve model generalization performance

1

Classification using logistic regression: build it from the ground up

[60 points]

This exercise will walk you through the full life-cycle of a supervised machine learning classification problem. Classification problem consists of two features/predictors (e.g. petal width and petal length) and your goal is to predict one of two possible classes (class 0 or class 1). You will build, train, and evaluate the performance of a logistic regression classifier on the data provided. Before you begin any modeling, you'll load and explore your data in Part I to familiarize yourself with it – and check for any missing or erroneous data. Then, in Part II, we will review an appropriate hypothesis set of functions to fit to the data: in this case, logistic regression. In Part III, we will derive an appropriate cost function for the data (spoiler alert: it's cross-entropy) as well as the gradient descent update equation that will allow you to optimize that cost function to identify the parameters that minimize the cost for the training data. In Part IV, all the pieces come together and you will implement your logistic regression model class including methods for fitting the data using gradient descent. Using that model you'll test it out and plot learning curves to verify the model learns as you train it and to identify an appropriate learning rate hyperparameter. Lastly, in Part V you will apply the model you designed, implemented, and verified to your actual data and evaluate and visualize its generalization performance as compared to a KNN algorithm. **When complete, you will have accomplished learning objectives 1-5 above!**

I. Load, prepare, and plot your data

You are given some data for which you are tasked with constructing a classifier. The first step when facing any machine learning project: look at your data!

(a) Load the data.

- In the data folder in the same directory of this notebook, you'll find the data in `A3_Q1_data.csv`. This file contains the binary class labels, y , and the features x_1 and x_2 .
- Divide your data into a training and testing set where the test set accounts for 30 percent of the data and the training set the remaining 70 percent.
- Plot the training data by class.
- Comment on the data: do the data appear separable? May logistic regression be a good choice for these data? Why or why not?

(b) Do the data require any preprocessing due to missing values, scale differences (e.g. different ranges of values), etc.? If so, how did you handle these issues?

Next, we walk through our key steps for model fitting: choose a hypothesis set of models to train (in this case, logistic regression); identify a cost function to measure the model fit to our training data; optimize model parameters to minimize cost (in this case using gradient descent). Once we've completed model fitting, we will evaluate the performance of our model and compare performance to another approach (a KNN classifier).

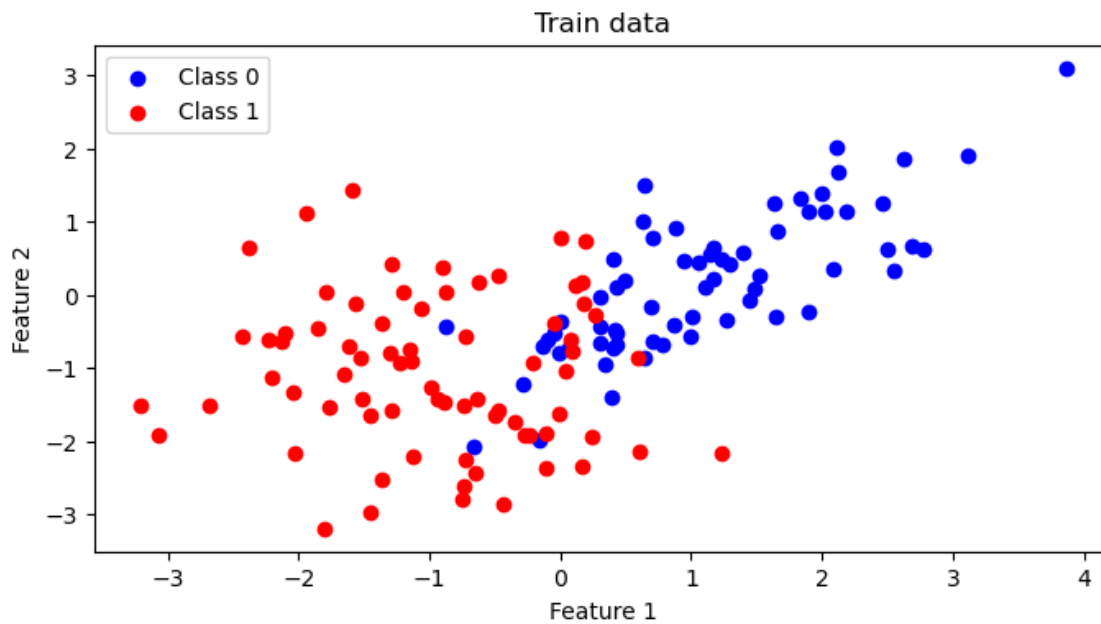
```
In [ ]: import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt
import warnings
import sklearn

# import dataset and split in test and train
db = pd.read_csv(
    "https://github.com/kylebradbury/ids705/raw/main/assignments/data/A3_Q1_data.csv"
)

X = db[["x1", "x2"]].values
y = db["y"].values
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.30, random_state=42
)

# plot train data

plt.figure(figsize=(8, 4))
plt.scatter(
    X_train[y_train == 0][:, 0],
    X_train[y_train == 0][:, 1],
    color="blue",
    label="Class 0",
)
plt.scatter(
    X_train[y_train == 1][:, 0],
    X_train[y_train == 1][:, 1],
    color="red",
    label="Class 1",
)
plt.title("Train data")
plt.xlabel("Feature 1")
plt.ylabel("Feature 2")
plt.legend()
plt.show()
```



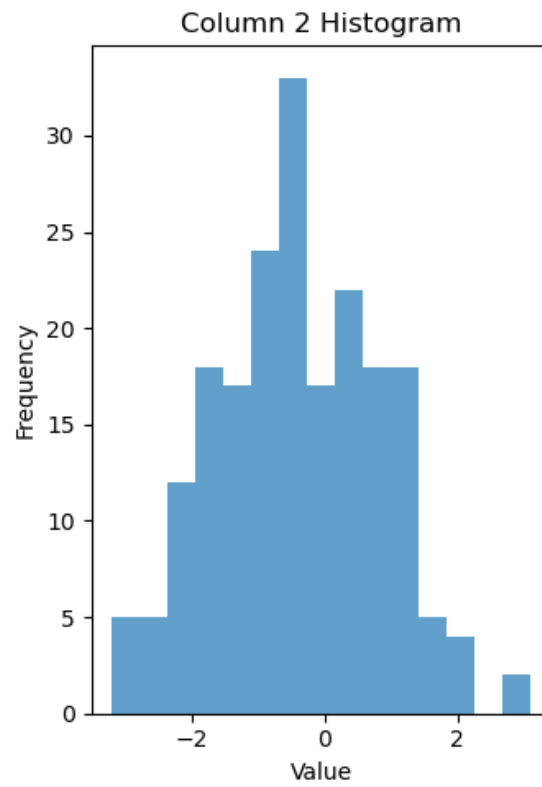
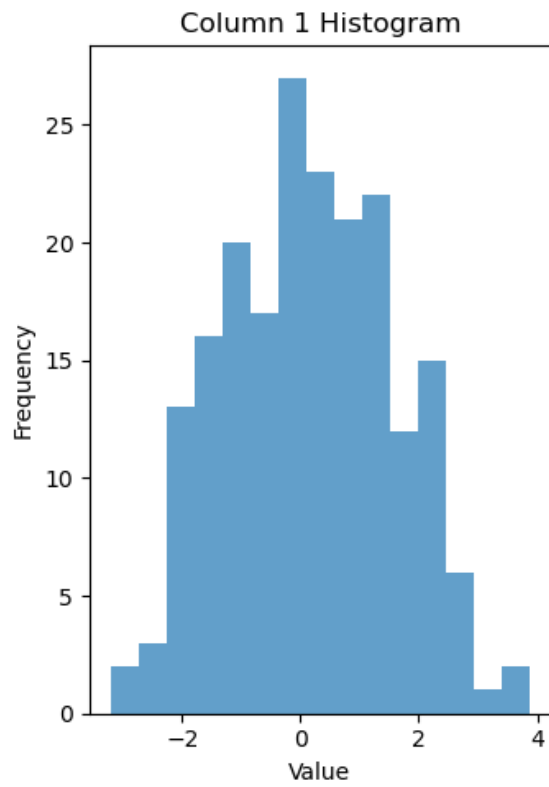
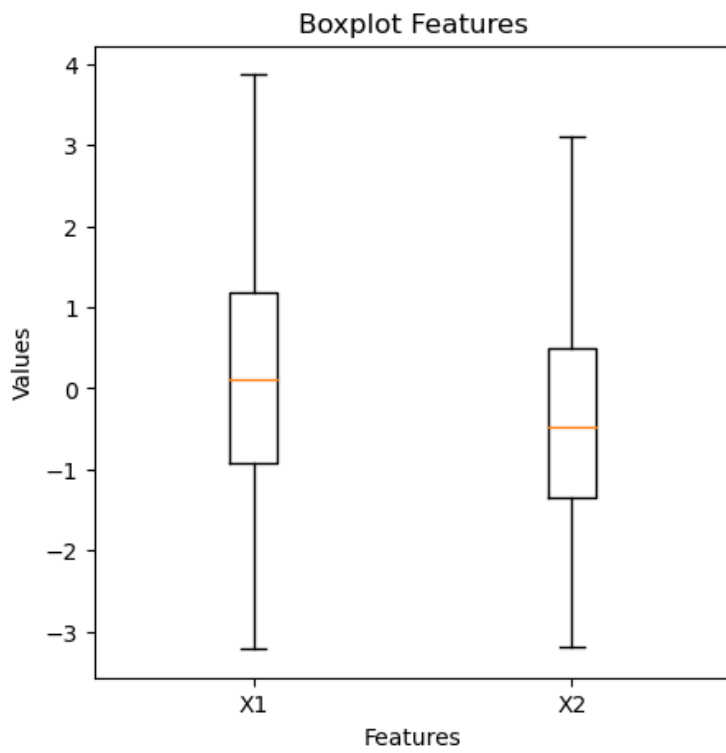
The data appears to be separable, a logistic regression might be the best choice as a sigmoid function could fit the data best. However such function might have trouble classifying the values shown at the center of the scatter plot.

```
In [ ]: plt.figure(figsize=(5, 5))
plt.boxplot(X, labels=["X1", "X2"])
plt.title("Boxplot Features")
plt.xlabel("Features")
plt.ylabel("Values")
plt.show()

num_bins = 10

# Create histograms
plt.figure(figsize=(7, 5))
for i in range(X.shape[1]):
    plt.subplot(1, X.shape[1], i + 1)
    plt.hist(X[:, i], bins=15, alpha=0.7)
    plt.title(f"Column {i+1} Histogram")
    plt.xlabel("Value")
    plt.ylabel("Frequency")

plt.tight_layout()
plt.show()
```



```
In [ ]: db.describe()
```

Out[]:	x1	x2	y
count	200.000000	200.000000	200.000000
mean	0.151376	-0.385426	0.485000
std	1.411722	1.217490	0.501029
min	-3.210005	-3.193456	0.000000
25%	-0.912029	-1.341047	0.000000
50%	0.112286	-0.479684	0.000000
75%	1.174400	0.495114	1.000000
max	3.867647	3.103541	1.000000

b)

The data does not require any preprocessing as the range of values does not show outliers and missing values are not present. All of the former can be determined by visualizing the distribution via boxplots, histograms and using tools to describe the data.

II. Stating the hypothesis set of models to evaluate (we'll use logistic regression)

Given that our data consists of two features, our logistic regression problem will be applied to a two-dimensional feature space. Recall that our logistic regression model is:

$$f(\mathbf{x}_i, \mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x}_i)$$

where the sigmoid function is defined as $\sigma(x) = \frac{e^x}{1 + e^x} = \frac{1}{1 + e^{-x}}$. Also, since this is a two-dimensional problem, we define $\mathbf{w}^\top \mathbf{x}_i = w_0 x_{i,0} + w_1 x_{i,1} + w_2 x_{i,2}$ and here, $\mathbf{x}_i = [x_{i,0}, x_{i,1}, x_{i,2}]^\top$, and $x_{i,0} \triangleq 1$

Remember from class that we interpret our logistic regression classifier output (or confidence score) as the conditional probability that the target variable for a given sample y_i is from class "1", given the observed features, \mathbf{x}_i . For one sample, (y_i, \mathbf{x}_i) , this is given as:

$$P(Y = 1 | X = \mathbf{x}_i) = f(\mathbf{x}_i, \mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x}_i)$$

In the context of maximizing the likelihood of our parameters given the data, we define this to be the likelihood function $L(\mathbf{w} | y_i, \mathbf{x}_i)$, corresponding to one sample observation from the training dataset.

*Aside: the careful reader will recognize this expression looks different from when we talk about the likelihood of our data given the true class label, typically expressed as $P(x|y)$, or the posterior probability of a class label given our data, typically expressed as $P(y|x)$. In the context of training a logistic regression model, the likelihood we are interested in is the likelihood function of our logistic regression **parameters**, \mathbf{w} . It's our goal to use this to choose the parameters to maximize the likelihood function.*

No output is required for this section - just read and use this information in the later sections.

III. Find the cost function that we can use to choose the model parameters, \mathbf{w} , that best fit the training data.

(c) What is the likelihood function that corresponds to all the N samples in our training dataset that we will wish to maximize? Unlike the likelihood function written above which gives the likelihood function for a *single training data pair* (y_i, \mathbf{x}_i) , this question asks for the likelihood function for the *entire training dataset* $\{(y_1, \mathbf{x}_1), (y_2, \mathbf{x}_2), \dots, (y_N, \mathbf{x}_N)\}$.

(d) Since a logarithm is a monotonic function, maximizing the $f(x)$ is equivalent to maximizing $\ln[f(x)]$. Express the likelihood from the last question as a cost function of the model parameters, $C(\mathbf{w})$; that is the negative of the logarithm of the likelihood. Express this cost as an average cost per sample (i.e. divide your final value by N), and use this quantity going forward as the cost function to optimize.

(e) Calculate the gradient of the cost function with respect to the model parameters $\nabla_{\mathbf{w}} C(\mathbf{w})$. Express this in terms of the partial derivatives of the cost function with respect to each of the parameters, e.g. $\nabla_{\mathbf{w}} C(\mathbf{w}) = \left[\frac{\partial C}{\partial w_0}, \frac{\partial C}{\partial w_1}, \frac{\partial C}{\partial w_2} \right]$.

To simplify notation, please use $\mathbf{w}^\top \mathbf{x}$ instead of writing out $w_0 x_{i,0} + w_1 x_{i,1} + w_2 x_{i,2}$ when it appears each time (where $x_{i,0} = 1$ for all i). You are also welcome to use $\sigma(\cdot)$ to represent the sigmoid function. Lastly, this will be a function the features, $x_{i,j}$ (with the first index in the subscript representing the observation and the second the feature; targets, y_i ; and the logistic regression model parameters, w_j .

(f) Write out the gradient descent update equation. This should clearly express how to update each weight from one step in gradient descent $w_j^{(k)}$ to the next $w_j^{(k+1)}$. There should be one equation for each model logistic regression model parameter (or you can represent it in vectorized form). Assume that η represents the learning rate.

c)

$$L(\omega) = \prod_{i=1}^N P(\omega, x_1)^{y_i} [1 - P(\omega x_1)]^{(1-y_i)}$$

\$\$

$$\begin{aligned} L(\omega \mid y, x) &= \prod_{i=1}^N \sigma(\omega^\top x_i)^{y_i} [1 - \sigma(\omega^\top x_i)]^{1-y_i} \\ &= \prod_{i=1}^N \hat{y}_1^{y_i} [1 - \hat{y}_i]^{1-y_i} \leftarrow \text{assuming } \hat{y} \triangleq \sigma(\omega^\top x_i) \end{aligned}$$

d)

$$\begin{aligned} c(\omega) &= -\log L(\omega \mid y_i x) \\ c(\omega) &= -\frac{1}{N} \left[\sum_{i=1}^N y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \right] \leftarrow \text{assuming } \hat{y} \triangleq \sigma(\omega^\top x_i) \end{aligned}$$

e)

$$c(\omega) = -\frac{1}{N} \sum_{i=1}^N \underbrace{y_i \log(\sigma z) + (1 - y_i) \log(1 - \sigma)}_L, \text{ where } \hat{y} = \underbrace{\sigma(\omega^\top x)}_Z$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \sigma} \frac{\partial \sigma}{\partial z} \frac{\partial z}{\partial \omega}$$

$$\begin{aligned} \frac{\partial L}{\partial \sigma} &= \frac{Y_i}{\sigma z_i} + \frac{(-1 + y_i)}{1 - \sigma z_i} = \frac{Y_i - \sigma z_i Y_i - \sigma z_i + y \sigma z_i}{\sigma z_i (1 - \sigma z_i)} = \\ &= \frac{Y_i - \sigma z_i}{\sigma(z_i) (1 - \sigma(z_i))} \end{aligned}$$

$$\begin{aligned} \frac{\partial \sigma}{\partial z} &= \frac{e^{-z}}{(1 + e^{-z})^2} = \frac{1}{1 + e^{-2}} - \frac{1}{(1 + e^{-2})^2} = \\ &= \sigma(z_i) - \sigma^2(z_i) = \\ &= \sigma(z_i) [1 - \sigma(z_i)] \end{aligned}$$

$$\frac{\partial z_i}{\partial \omega} = x_i$$

$$\begin{aligned} \frac{\partial L}{\partial w} &= \left[\frac{y_1 - \sigma(z)}{\sigma(z)(1 - \sigma(z))} \right] \sigma(z_i) [1 - \sigma(z_i)] x_i = \\ &= [y_i - \sigma(z)] x_i \end{aligned}$$

$$\nabla_{\mathbf{w}} C(\mathbf{w}) = -\frac{1}{N} \sum_{i=1}^N [(y_i - \sigma(z)) x_i]$$

f)

$$\omega^{(n+1)} = \omega^{(n)} - \lambda \nabla C(\omega), \text{ where } \nabla C(\omega) = -\frac{1}{N} \sum_{i=1}^N [y_i - \sigma(z)] x_i, z = \omega^T x_i$$

$$\begin{aligned} \begin{bmatrix} \omega_0^{k+1} \\ \omega_1^{k+1} \\ \omega_2^{k+1} \end{bmatrix} &= \begin{bmatrix} \omega_0^k \\ \omega_1^k \\ \omega_2^k \end{bmatrix} - \lambda \begin{bmatrix} \frac{\partial C}{\partial \omega_0} \\ \frac{\partial C}{\partial \omega_1} \\ \frac{\partial C}{\partial \omega_2} \end{bmatrix} \\ &= \begin{bmatrix} \omega_0^k \\ \omega_1^k \\ \omega_2^k \end{bmatrix} - \lambda \begin{pmatrix} \frac{1}{N} \sum_{i=1}^N [\sigma z - y_i] x_{i,0} \\ \frac{1}{N} \sum_{i=1}^N [\sigma z - y_i] x_{i,1} \\ \frac{1}{N} \sum_{i=1}^N [\sigma z - y_i] x_{i,2} \end{pmatrix} \end{aligned}$$

IV. Implement gradient descent and your logistic regression algorithm

(g) Implement your logistic regression model.

- You are provided with a template, below, for a class with key methods to help with your model development. It is modeled on the Scikit-Learn convention. For this, you only need to create a version of logistic regression for the case of two feature variables (i.e. two predictors).
- Create a method called `sigmoid` that calculates the sigmoid function
- Create a method called `cost` that computes the cost function $C(\mathbf{w})$ for a given dataset and corresponding class labels. This should be the **average cost** (make sure your total cost is divided by your number of samples in the dataset).
- Create a method called `gradient_descent` to run **one step** of gradient descent on your training data. We'll refer to this as "batch" gradient descent since it takes into account the gradient based on all our data at each iteration of the algorithm.
- Create a method called `fit` that fits the model to the data (i.e. sets the model parameters to minimize cost) using your `gradient_descent` method. In doing this we'll need to make some assumptions about the following:
 - Weight initialization. What should you initialize the model parameters to? For this, randomly initialize the weights to a different values between 0 and 1.
 - Learning rate. How slow/fast should the algorithm step towards the minimum? This you will vary in a later part of this problem.
 - Stopping criteria. When should the algorithm be finished searching for the optimum? There are two stopping criteria: small changes in the gradient descent step size and a maximum number of iterations. The first is whether there was a sufficiently small change in the gradient; this is evaluated as whether the magnitude of the step that the gradient descent algorithm takes changes by less than 10^{-6} between iterations. Since we have a weight vector, we can compute the change in the weight by evaluating the L_2 norm (Euclidean norm) of the change in the vector between iterations. From our gradient descent update equation we know that mathematically this is $\|\eta \nabla_{\mathbf{w}} C(\mathbf{w})\|$. The second criterion is met if a maximum number of iterations has been reached (5,000 in this case, to prevent infinite loops from poor choices of learning rates).
 - Design your approach so that at each step in the gradient descent algorithm you evaluate the cost function for both the training and the test data for each new value for the model weights. You should be able to plot cost vs gradient descent iteration for both the training and the test data. This will allow you to plot "learning curves" that can be informative for how the model training process is proceeding.
- Create a method called `predict_proba` that predicts confidence scores (that can be thresholded into the predictions of the `predict` method).
- Create a method called `predict` that makes predictions based on the trained model, selecting the most probable class, given the data, as the prediction, that is class that yields the larger $P(y|\mathbf{x})$.
- (Optional, but recommended) Create a method called `learning_curve` that produces the cost function values that correspond to each step from a previously run gradient descent operation.
- (Optional, but recommended) Create a method called `prepare_x` which appends a column of ones as the first feature of the dataset \mathbf{X} to account for the bias term ($x_{i,1} = 1$).

This structure is strongly encouraged; however, you're welcome to adjust this to your needs (adding helper methods, modifying parameters, etc.).

g)

```
In [ ]: # Logistic regression class
class Logistic_regression:
    # Class constructor
    def __init__(self):
        self.w = None # logistic regression weights
        self.saved_w = [] # Since this is a small problem, we can save the weights
        # at each iteration of gradient descent to build our
        # learning curves
        # returns nothing
        pass

    # Method for calculating the sigmoid function of  $w^T X$  for an input set of weights
    def sigmoid(self, X, w):
        z = w.T @ X
        sigmoid = 1 / (1 + np.exp(-z))
        return sigmoid # returns the value of the sigmoid
        # pass

    # Cost function for an input set of weights
    def cost(self, X, y, w):
        N = X.shape[1]
        L = y * np.log(self.sigmoid(X, w)) + (1 - y) * np.log(1 - self.sigmoid(X, w))
        ce = -(1 / N) * np.sum(L)
        return ce

    # Update the weights in an iteration of gradient descent
    def gradient_descent(self, X, y, lr):
        # returns a scalar of the magnitude of the Euclidean norm
        # of the change in the weights during one gradient descent step
        N = X.shape[1]
        gd = (-1 / N) * ((y - self.sigmoid(X, self.w)) @ X.T) # cambiar my_w por self.w
        me = np.linalg.norm(gd)
        self.w = self.w - (lr * gd).T
        return me

    # returns a scalar of the magnitude of the Euclidean norm
    # of the change in the weights during one gradient descent step

    # Fit the logistic regression model to the data through gradient descent
    def fit(self, X, y, w_init, lr, delta_thresh=1e-6, max_iter=5000, verbose=False):
        # Note the verbose flag enables you to print out the weights at each iteration
        # (optional - but may help with one of the questions)
        # returns nothing
        self.w = w_init

        for i in range(max_iter):
            # Update weights using gradient descent
            self.gradient_descent(X, y, lr)
            self.saved_w.append(self.w)
            if self.gradient_descent(X, y, lr) < delta_thresh:
                break
        if verbose:
            print(self.saved_w)

    # Use the trained model to predict the confidence scores (prob of positive class in this case)
    def predict_proba(self, X):
        proba = self.sigmoid(X, self.w)
        return proba

    # Use the trained model to make binary predictions
    def predict(self, X, thresh=0.5):
        # returns a binary prediction for each sample
        predited = (self.predict_proba(X) > thresh).astype(int)
        return predited

    # Stores the learning curves from saved weights from gradient descent
```



```

def learning_curve(self, X, y):
    # returns the value of the cost function from each step in gradient descent
    # from the last model fitting process
    costs = []
    for myw in self.saved_w:
        cost_value = self.cost(X, y, myw)
        costs.append(cost_value)
    return costs

# Appends a column of ones as the first feature to account for the bias term
def prepare_x(self, X):
    # returns the X with a new feature of all ones (a column that is the new column 0)
    xplus1 = np.vstack((np.ones((1, X.shape[1])), X))
    return xplus1

```

(h) Choose a learning rate and fit your model. Learning curves are a plot of metrics of model performance evaluated through the process of model training to provide insights about how model training is proceeding. Show the learning curves for the gradient descent process for learning rates of $\{10^{-0}, 10^{-2}, 10^{-4}\}$. For each learning rate plot the learning curves by plotting **both the training and test data average cost** as a function of each iteration of gradient descent. You should run the model fitting process until it completes (up to 5,000 iterations of gradient descent). All of the 6 resulting curves (train and test average cost for each learning rate) should be plotted on the **same set of axes** to enable direct comparison. *Note: make sure you're using average cost per sample, not the total cost.*

- Try running this process for a really big learning rate for this problem: 10^2 . Look at the weights that the fitting process generates over the first 50 iterations and how they change. Either print these first 50 iterations as console output or plot them. What happens? How does the output compare to that corresponding to a learning rate of 10^0 and why?
- What is the impact that the different values of learning have on the speed of the process and the results?
- Of the options explored, what learning rate do you prefer and why?
- Use your chosen learning rate for the remainder of this problem.

h)

```

In [ ]: import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split

db = pd.read_csv(
    "https://github.com/kylebradbury/ids705/raw/main/assignments/data/A3_Q1_data.csv"
)

X = db[["x1", "x2"]].values
y = db["y"].values
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.30, random_state=42
)

X_train = X_train.T
X_test = X_test.T

```

```

In [ ]: clf = LogisticRegression()

w_init = np.random.rand(3, 1)

X_train = clf.prepare_x(X_train)
X_test = clf.prepare_x(X_test)

model_1_t = LogisticRegression()
model_2_t = LogisticRegression()
model_3_t = LogisticRegression()

model_1_t.fit(
    X_train,
    y_train,
    w_init,

```

```

    10 ** (-0),
    delta_thresh=1e-6,
    max_iter=5000,
    verbose=False,
)
model_2_t.fit(
    X_train,
    y_train,
    w_init,
    10 ** (-2),
    delta_thresh=1e-6,
    max_iter=5000,
    verbose=False,
)
model_3_t.fit(
    X_train,
    y_train,
    w_init,
    10 ** (-4),
    delta_thresh=1e-6,
    max_iter=5000,
    verbose=False,
)

```

```

In [ ]: m_1_test_c = model_1_t.learning_curve(X_train, y_train)
m_1_train_c = model_1_t.learning_curve(X_test, y_test)
m_2_test_c = model_2_t.learning_curve(X_train, y_train)
m_2_train_c = model_2_t.learning_curve(X_test, y_test)
m_3_test_c = model_3_t.learning_curve(X_train, y_train)
m_3_train_c = model_3_t.learning_curve(X_test, y_test)

```

```

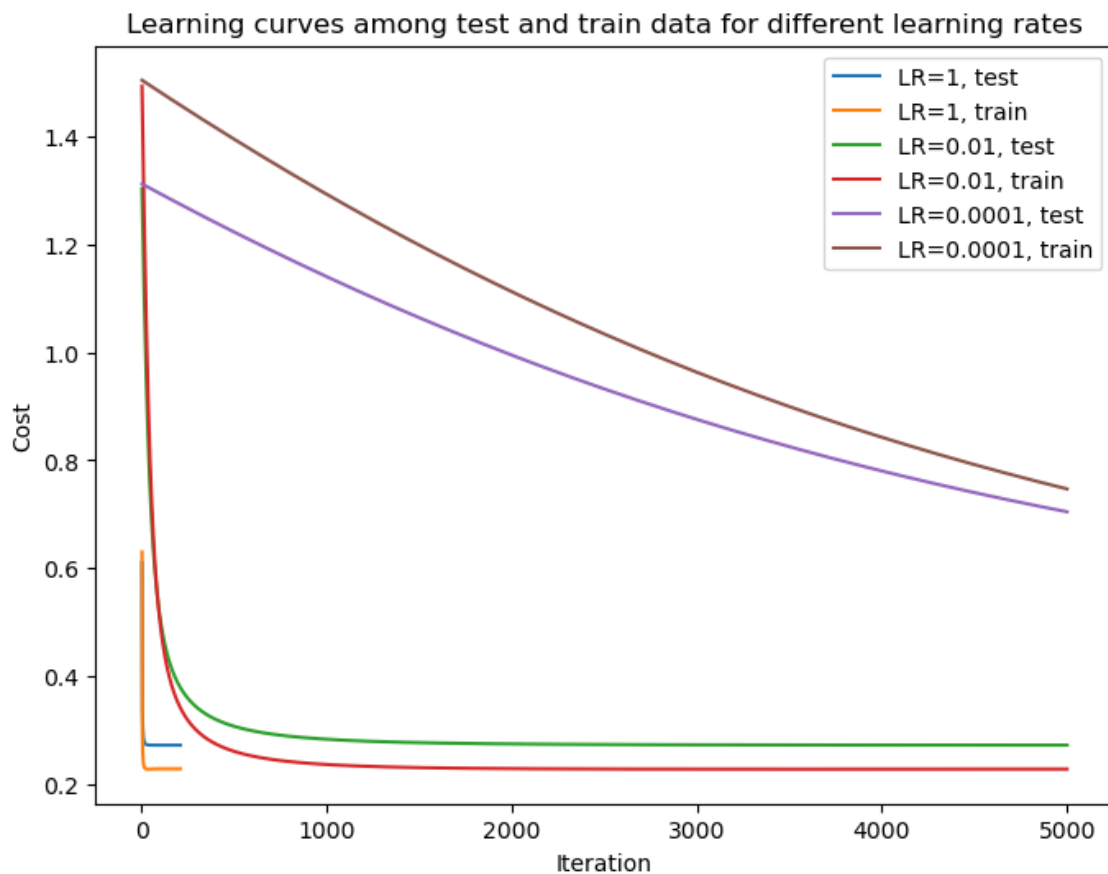
In [ ]: import matplotlib.pyplot as plt

x_values_1 = range(1, len(m_1_test_c) + 1)
x_values_2 = range(1, len(m_2_test_c) + 1)
x_values_3 = range(1, len(m_3_test_c) + 1)

plt.figure(figsize=(8, 6))
# Plot the curves
plt.plot(x_values_1, m_1_test_c, label="LR=1, test")
plt.plot(x_values_1, m_1_train_c, label="LR=1, train")
plt.plot(x_values_2, m_2_test_c, label="LR=0.01, test")
plt.plot(x_values_2, m_2_train_c, label="LR=0.01, train")
plt.plot(x_values_3, m_3_test_c, label="LR=0.0001, test")
plt.plot(x_values_3, m_3_train_c, label="LR=0.0001, train")

plt.title("Learning curves among test and train data for different learning rates")
plt.xlabel("Iteration")
plt.ylabel("Cost")
plt.legend()
plt.show()

```



```
In [ ]: import warnings

warnings.simplefilter(action="ignore")

model_4_t = Logistic_regression()
model_4_t.fit(
    X_train,
    y_train,
    w_init,
    10 ** (2),
    delta_thresh=1e-6,
    max_iter=5000,
    verbose=False,
)
m_4_test_c = model_4_t.learning_curve(X_train, y_train)
m_4_train_c = model_4_t.learning_curve(X_test, y_test)
wights = model_4_t.saved_w

we1 = [matrix[0, 0] for matrix in wights]
we2 = [matrix[1, 0] for matrix in wights]
we3 = [matrix[2, 0] for matrix in wights]

we3 = we3[:50]
we1 = we1[:50]
we2 = we2[:50]

x_ax = range(1, len(we2) + 1)

plt.figure(figsize=(5, 3))
# Plot the curves
plt.plot(x_ax, we1, label="weight 1")
plt.plot(x_ax, we2, label="weight 2")
plt.plot(x_ax, we3, label="weight 3")

plt.title("Weights for first 50 iterations, learning rate 100")
plt.xlabel("Iteration")
plt.ylabel("Weights")
plt.legend()
```

```

plt.show()

model_5_t = Logistic_regression()
model_5_t.fit(
    X_train,
    y_train,
    w_init,
    10 ** (0),
    delta_thresh=1e-6,
    max_iter=5000,
    verbose=False,
)
# m_5_test_c = model_5_t.learning_curve(X_train, y_train)
# m_5_train_c = model_5_t.learning_curve(X_test, y_test)
wights = model_5_t.saved_w

we1 = [matrix[0, 0] for matrix in wights]
we2 = [matrix[1, 0] for matrix in wights]
we3 = [matrix[2, 0] for matrix in wights]

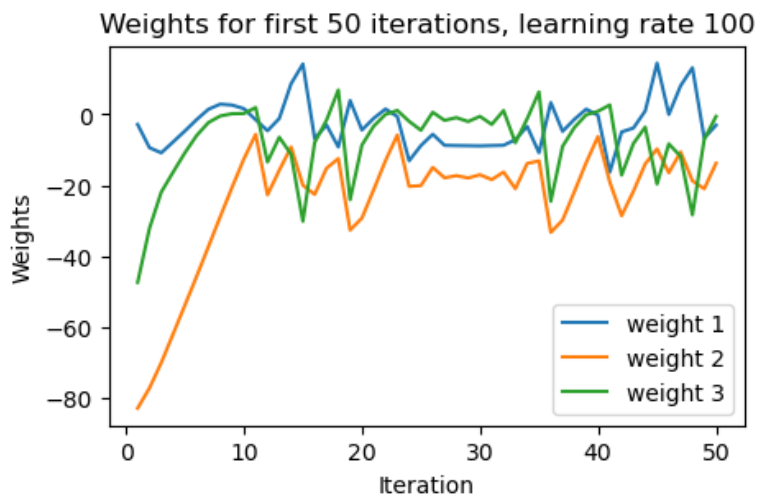
we1 = we1[:50]
we2 = we2[:50]
we3 = we3[:50]

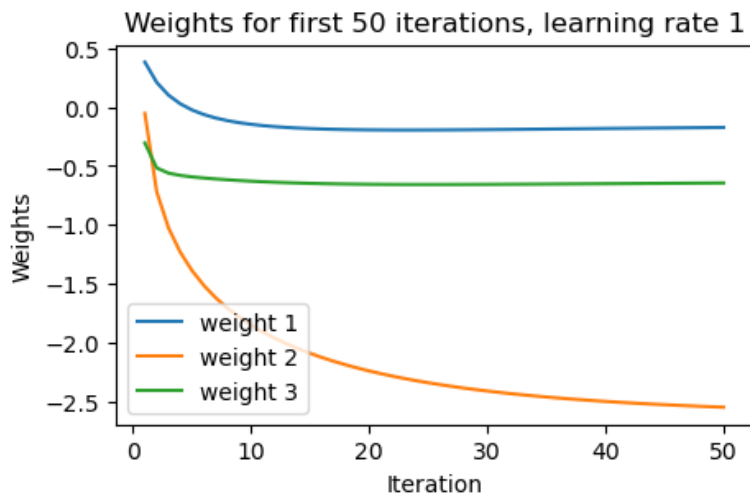
x_ax = range(1, len(we2) + 1)

plt.figure(figsize=(5, 3))
# Plot the curves
plt.plot(x_ax, we1, label="weight 1")
plt.plot(x_ax, we2, label="weight 2")
plt.plot(x_ax, we3, label="weight 3")

plt.title("Weights for first 50 iterations, learning rate 1")
plt.xlabel("Iteration")
plt.ylabel("Weights")
plt.legend()
plt.show()

```





- Try running this process for a really big learning rate for this problem: 10^2 . Look at the weights that the fitting process generates over the first 50 iterations and how they change. Either print these first 50 iterations as console output or plot them. What happens? How does the output compare to that corresponding to a learning rate of 10^0 and why?
- What is the impact that the different values of learning have on the speed of the process and the results?
- Of the options explored, what learning rate do you prefer and why?

The impact of different values of learning have on the speed of the process is the amount of time (iterations) it takes to reach the convergence point.

When the learning patterns changes from 100 to 1, each weight smoothes down.

The preferred learning rate is 1 as its the fastest to reach the convergence point.

V. Evaluate your model performance through cross validation

(i) Test the performance of your trained classifier using K-folds cross validation resampling technique. The scikit-learn package [StratifiedKFolds](#) may be helpful.

- Train your logistic regression model and a K-Nearest Neighbor classification model with $k = 7$ nearest neighbors.
- Using the trained models, make four plots: two for logistic regression and two for KNN. For each model have one plot showing the training data used for fitting the model, and the other showing the test data. On each plot, include the decision boundary resulting from your trained classifier.
- Produce a Receiver Operating Characteristic curve (ROC curve) that represents the performance from cross validated performance evaluation for each classifier (your logistic regression model and the KNN model, with $k = 7$ nearest neighbors). For the cross validation, use $k = 10$ folds.
 - Plot these curves on the same set of axes to compare them. You should not plot one curve for each fold of k-folds; instead, you should plot one ROC curve for Logistic Regression and one for KNN (each should incorporate all 10 folds of validation). Also, don't forget to plot the "chance" line.
 - On the ROC curve plot, also include the chance diagonal for reference (this represents the performance of the worst possible classifier). This is represented as a line from (0, 0) to (1, 1).
 - Calculate the Area Under the Curve for each model and include this measure in the legend of the ROC plot.
- Comment on the following:
 - What is the purpose of using cross validation for this problem?
 - How do the models compare in terms of performance (both ROC curves and decision boundaries) and which model (logistic regression or KNN) would you select to use on previously unseen data for this problem and why?

ANSWER

i)

```

In [ ]: from sklearn.neighbors import KNeighborsClassifier
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split

db = pd.read_csv(
    "https://github.com/kylebradbury/ids705/raw/main/assignments/data/A3_Q1_data.csv"
)

# KNN
X = db[["x1", "x2"]].values
y = db["y"].values
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.30, random_state=42
)

model_knn_1 = KNeighborsClassifier(n_neighbors=7)
model_knn_1.fit(X_train, y_train)

# Logistic reg

X = db[["x1", "x2"]].values
y = db["y"].values
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.30, random_state=42
)
model_h_t = LogisticRegression()

X_train2 = X_train.T
X_test2 = X_test.T
X_train2 = model_h_t.prepare_x(X_train2)
X_test2 = model_h_t.prepare_x(X_test2)
w_init = np.random.rand(3, 1)

model_h_t.fit(
    X_train2,
    y_train,
    w_init,
    1,
    delta_thresh=1e-6,
    max_iter=5000,
    verbose=False,
)

```

```

In [ ]: import matplotlib.pyplot as plt

def plot_decision_boundary(X, y, model, title, ax=None):
    if ax is None:
        ax = plt.gca() # If ax is not provided, use current axis

    # Plot data points
    ax.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.coolwarm, s=20, edgecolors="k")

    # Create meshgrid for decision boundary
    x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
    y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
    xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.1), np.arange(y_min, y_max, 0.1))
    Z = model.predict(np.c_[xx.ravel(), yy.ravel()])

    # Plot decision boundary
    Z = Z.reshape(xx.shape)
    ax.contourf(xx, yy, Z, alpha=0.3, cmap=plt.cm.coolwarm)

    ax.set_title(title)
    ax.set_xlabel("Feature 1")
    ax.set_ylabel("Feature 2")

```

```

In [ ]: from matplotlib.lines import Line2D

legend_elements = [
    Line2D(
        [0],
        [0],
        marker="o",
        color="w",
        label="Class 0",
        markerfacecolor="blue",
        markersize=10,
    ),
    Line2D(
        [0],
        [0],
        marker="o",
        color="w",
        label="Class 1",
        markerfacecolor="red",
        markersize=10,
    ),
]

# Plot LGR
# TRAIN
x_min, x_max = X_train2[1, :].min() - 1, X_train2[1, :].max() + 1
y_min, y_max = X_train2[2, :].min() - 1, X_train2[2, :].max() + 1

xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.1), np.arange(y_min, y_max, 0.1))

meshgrid_data = np.c_[xx.ravel(), yy.ravel()].T
meshgrid_data = model_h_t.prepare_x(meshgrid_data)

Z = model_h_t.predict(meshgrid_data)

Z = Z.reshape(xx.shape)

fig, axi = plt.subplots(2, 2, figsize=(9, 7))

axi[0][0].contourf(xx, yy, Z, alpha=0.3, cmap=plt.cm.coolwarm)
axi[0][0].scatter(
    X_train[:, 0], X_train[:, 1], c=y_train, cmap=plt.cm.coolwarm, s=20, edgecolors="k"
)
axi[0][0].set_xlabel("Feature 1")
axi[0][0].set_ylabel("Feature 2")
axi[0][0].set_title("Logistic Regression Decision Boundary --Train")

axi[0][0].legend(handles=legend_elements, loc="upper right")

# TEST
x_min, x_max = X_test2[1, :].min() - 1, X_test2[1, :].max() + 1
y_min, y_max = X_test2[2, :].min() - 1, X_test2[2, :].max() + 1

xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.1), np.arange(y_min, y_max, 0.1))

meshgrid_data = np.c_[xx.ravel(), yy.ravel()].T
meshgrid_data = model_h_t.prepare_x(meshgrid_data)

Z = model_h_t.predict(meshgrid_data)

Z = Z.reshape(xx.shape)

axi[0][1].contourf(xx, yy, Z, alpha=0.3, cmap=plt.cm.coolwarm)
axi[0][1].scatter(
    X_test[:, 0], X_test[:, 1], c=y_test, cmap=plt.cm.coolwarm, s=20, edgecolors="k"
)
axi[0][1].set_xlabel("Feature 1")
axi[0][1].set_ylabel("Feature 2")

```

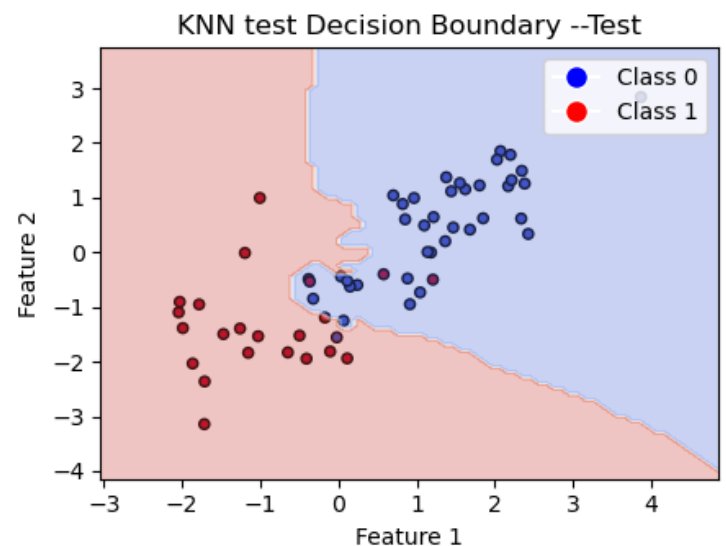
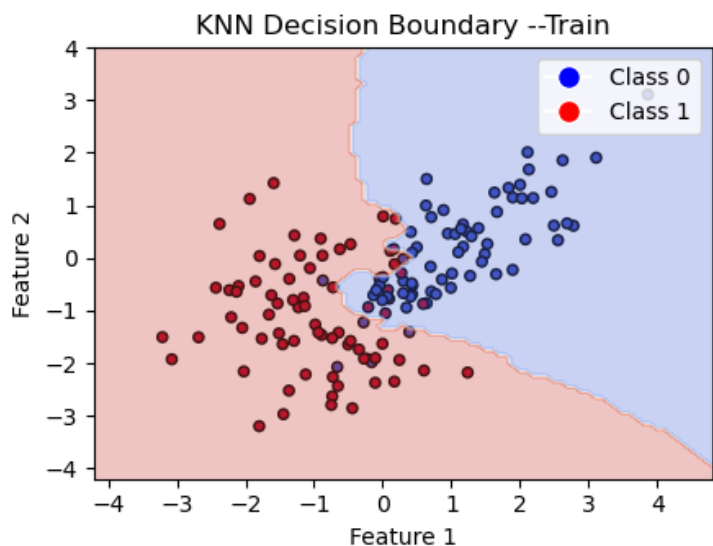
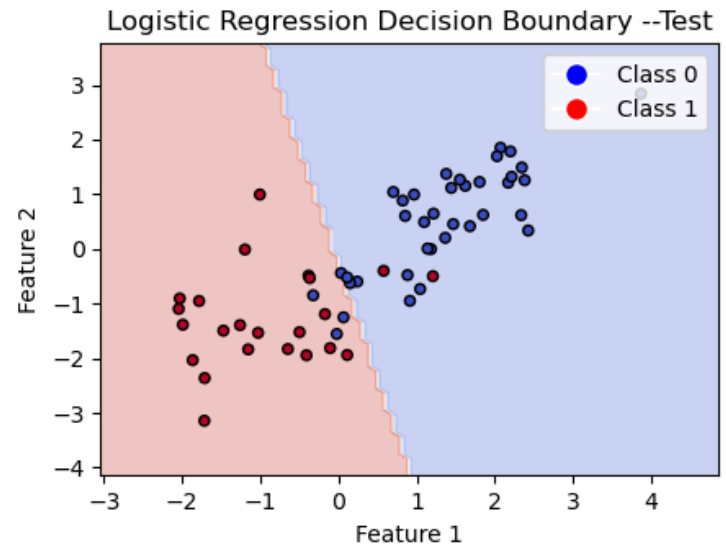
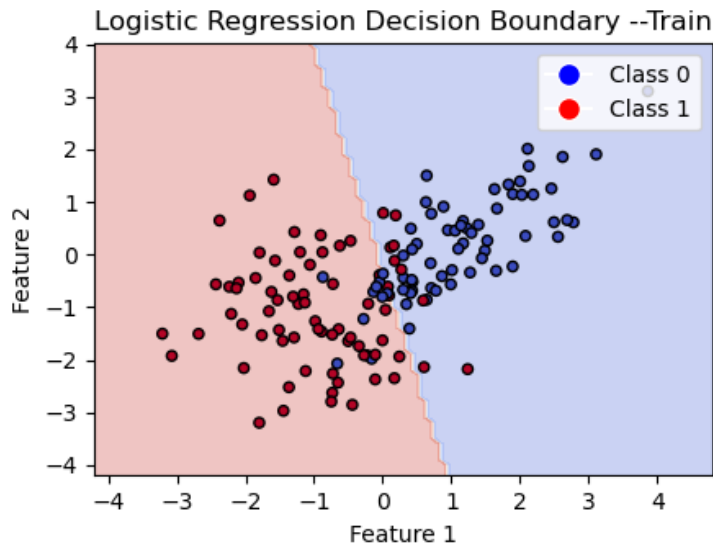
```

axi[0][1].set_title("Logistic Regression Decision Boundary --Test")
axi[0][1].legend(handles=legend_elements, loc="upper right")
axi[1][0].legend(handles=legend_elements, loc="upper right")
axi[1][1].legend(handles=legend_elements, loc="upper right")

# Plot KNN
plot_decision_boundary(
    X_train, y_train, model_knn_1, "KNN Decision Boundary --Train", ax=axi[1][0]
)
plot_decision_boundary(
    X_test, y_test, model_knn_1, "KNN test Decision Boundary --Test", ax=axi[1][1]
)

# plot all
plt.tight_layout()
plt.show()

```



```

In [ ]: from sklearn.model_selection import StratifiedKFold, cross_val_score
        from sklearn.neighbors import KNeighborsClassifier
        from sklearn.metrics import roc_curve, auc
        import matplotlib.pyplot as plt

        X = db[["x1", "x2"]].values
        y = db["y"].values

        model_lr = LogisticRegression()
        model_knn = KNeighborsClassifier(n_neighbors=7)
        cv = StratifiedKFold(n_splits=10)

        # lists for tpr and fdr

```



```

mean_tpr_knn = 0.0
mean_fpr_knn = np.linspace(0, 1, 100)
mean_tpr_lgr = 0.0
mean_fpr_lgr = np.linspace(0, 1, 100)

# 10-fold cv for knn
scores_knn = cross_val_score(model_knn, X, y, cv=10, scoring="roc_auc")

# StratifiedKFold cv object:
cv = StratifiedKFold(n_splits=10)

# Iterate over each fold
for train_idx, test_idx in cv.split(X, y):
    X_train_fold, X_test_fold = X[train_idx], X[test_idx]
    y_train_fold, y_test_fold = y[train_idx], y[test_idx]

    """knn"""
    # Train KNN model
    model_knn.fit(X_train_fold, y_train_fold)

    # predict
    y_hat_knn = model_knn.predict_proba(X_test_fold[:, 1])

    # ROC
    fpr, tpr, _ = roc_curve(y_test_fold, y_hat_knn)
    mean_tpr_knn += np.interp(mean_fpr_knn, fpr, tpr)
    mean_tpr_knn[0] = 0.0

    """lgr"""
    # Train lgr, necessary to transpose X for class to work
    X_train_fold_p = model_lr.prepare_x(X_train_fold.T)
    w_init = np.zeros(X_train_fold_p.shape[0])
    model_lr.fit(X_train_fold_p, y_train_fold, w_init, 1, max_iter=1000)

    # prepare X test and transpose X
    X_test_fold_p = model_lr.prepare_x(X_test_fold.T)

    # predict
    y_hat_lr = model_lr.predict_proba(X_test_fold_p)

    # ROC
    fpr, tpr, _ = roc_curve(y_test_fold, y_hat_lr)
    mean_tpr_lgr += np.interp(mean_fpr_lgr, fpr, tpr)
    mean_tpr_lgr[0] = 0.0

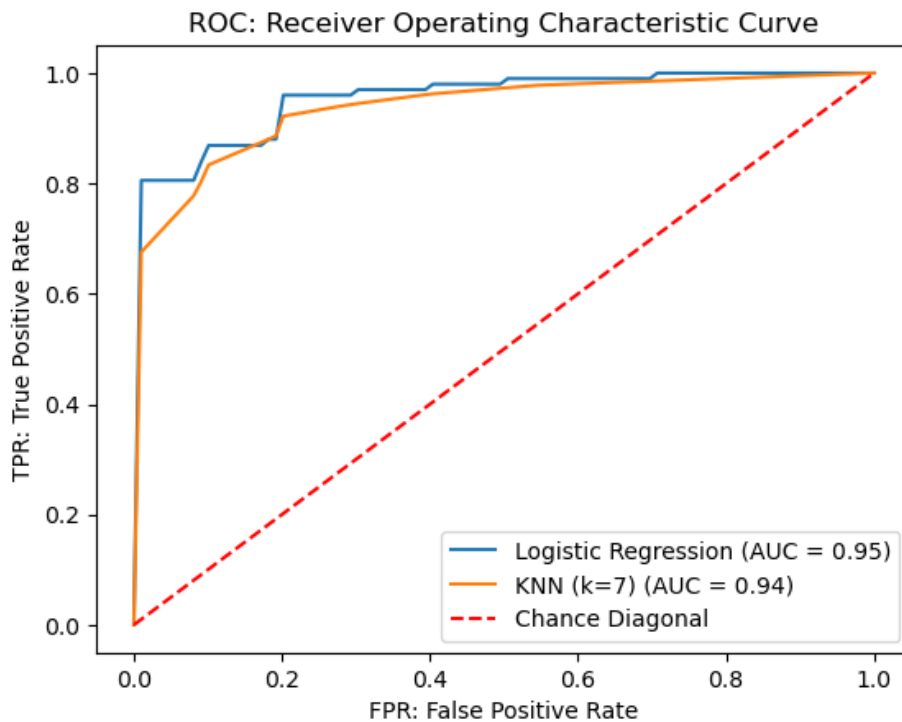
# Average tprs
mean_tpr_lgr /= cv.get_n_splits()
mean_tpr_knn /= cv.get_n_splits()

# Average AUC
mean_auc_lgr = auc(mean_fpr_lgr, mean_tpr_lgr)
mean_auc_knn = auc(mean_fpr_knn, mean_tpr_knn)

# Plot ROC curves for logistic regression and KNN
plt.plot(
    mean_fpr_lgr, mean_tpr_lgr, label=f"Logistic Regression (AUC = {mean_auc_lgr:.2f})"
)
plt.plot(mean_fpr_knn, mean_tpr_knn, label=f"KNN (k=7) (AUC = {mean_auc_knn:.2f})")
plt.plot([0, 1], [0, 1], color="red", linestyle="--", label="Chance Diagonal")
# Add labels and legend
plt.xlabel("FPR: False Positive Rate")
plt.ylabel("TPR: True Positive Rate")
plt.title("ROC: Receiver Operating Characteristic Curve")
plt.legend()

# Show plot
plt.show()

```



2

Digits classification

[30 points]

An exploration of regularization, imbalanced classes, ROC and PR curves

The goal of this exercise is to apply your supervised learning skills on a very different dataset: in this case, image data; MNIST: a collection of images of handwritten digits. Your goal is to train a classifier that is able to distinguish the number "3" from all possible numbers and to do so as accurately as possible. You will first explore your data (this should always be your starting point to gain domain knowledge about the problem.). Since the feature space in this problem is 784-dimensional, overfitting is possible. To avoid overfitting you will investigate the impact of regularization on generalization performance (test accuracy) and compare regularized and unregularized logistic regression model test error against other classification techniques such as linear discriminant analysis and random forests and draw conclusions about the best-performing model.

Start by loading your dataset from the [MNIST dataset](#) of handwritten digits, using the code provided below. MNIST has a training set of 60,000 examples, and a test set of 10,000 examples. The digits have been size-normalized and centered in a fixed-size image.

Your goal is to classify whether or not an example digit is a 3. Your binary classifier should predict $y = 1$ if the digit is a 3, and $y = 0$ otherwise. Create your dataset by transforming your labels into a binary format (3's are class 1, and all other digits are class 0).

(a) Plot 10 examples of each class (i.e. class $y = 0$, which are not 3's and class $y = 1$ which are 3's), from the training dataset.

- Note that the data are composed of samples of length 784. These represent 28 x 28 images, but have been reshaped for storage convenience. To plot digit examples, you'll need to reshape the data to be 28 x 28 (which can be done with numpy `reshape`).

(b) How many examples are present in each class? Show a plot of samples by class (bar plot). What fraction of samples are positive? What issues might this cause?

(c) Identify the value of the regularization parameter that optimizes model performance on out-of-sample data. Using a logistic regression classifier, apply lasso regularization and retrain the model and evaluate its performance on the test set over a range of values on the regularization coefficient. You can implement this using the [LogisticRegression](#) module and activating the 'l1' penalty;

the parameter C is the inverse of the regularization strength. Vary the value of C logarithmically from 10^{-4} to 10^4 (and make your x-axes logarithmic in scale) and evaluate it at least 20 different values of C . As you vary the regularization coefficient, Plot the following four quantities (this should result in 4 separate plots)...

- The number of model parameters that are estimated to be nonzero (in the logistic regression model, one attribute is `coef_`, which gives you access to the model parameters for a trained model)
- The cross entropy loss (which can be evaluated with the Scikit Learn `log_loss` function)
- Area under the ROC curve (AUC)
- The F_1 -score (assuming a threshold of 0.5 on the predicted confidence scores, that is, scores above 0.5 are predicted as Class 1, otherwise Class 0). Scikit Learn also has a `f1_score` function which may be useful. -Which value of C seems best for this problem? Please select the closest power of 10. You will use this in the next part of this exercise.

(d) Train and test a (1) logistic regression classifier with minimal regularization (using the Scikit Learn package, set `penalty='l1'`, `C=1e100` to approximate this), (2) a logistic regression classifier with the best value of the regularization parameter from the last section, (3) a Linear Discriminant Analysis (LDA) Classifier, and (4) a Random Forest (RF) classifier (using default parameters for the LDA and RF classifiers).

- Compare your classifiers' performance using ROC and Precision Recall (PR) curves. For the ROC curves, all your curves should be plotted on the same set of axes so that you can directly compare them. Please do the same with the PR curves.
- Plot the line that represents randomly guessing the class (50% of the time a "3", 50% not a "3"). You SHOULD NOT actually create random guesses. Instead, you should think through the theory behind how ROC and PR curves work and plot the appropriate lines. It's a good practice to include these in ROC and PR curve plots as a reference point.
- For PR curves, an excellent resource on how to correctly plot them can be found [here](#) (ignore the section on "non-linear interpolation between two points"). This describes how a random classifier is represented in PR curves and demonstrates that it should provide a lower bound on performance.
- When training your logistic regression model, it's recommended that you use `solver='liblinear'`; otherwise, your results may not converge.
- Describe the performance of the classifiers you compared. Did the regularization of the logistic regression model make much difference here? Which classifier you would select for application to unseen data.

ANSWER

a)

```
In [ ]: # Load the MNIST Data
from sklearn.datasets import fetch_openml
from sklearn.model_selection import train_test_split
import numpy as np
import matplotlib.pyplot as plt
import pickle

# Set this to True to download the data for the first time and False after the first time
# so that you just load the data locally instead
download_data = True

if download_data:
    # Load data from https://www.openml.org/d/554
    X, y = fetch_openml("mnist_784", return_X_y=True, as_frame=False)

    # Adjust the labels to be '1' if y==3, and '0' otherwise
    y[y != "3"] = 0
    y[y == "3"] = 1
    y = y.astype("int")

    # Divide the data into a training and test split
    X_train, X_test, y_train, y_test = train_test_split(
        X, y, test_size=1 / 7, random_state=88
    )

    file = open("tmpdata", "wb")
    pickle.dump((X_train, X_test, y_train, y_test), file)
    file.close()
else:
    file = open("tmpdata", "rb")
```

```
X_train, X_test, y_train, y_test = pickle.load(file)
file.close()
```

```
In [ ]: import matplotlib.pyplot as plt
```

```
num_row = 4
num_col = 5

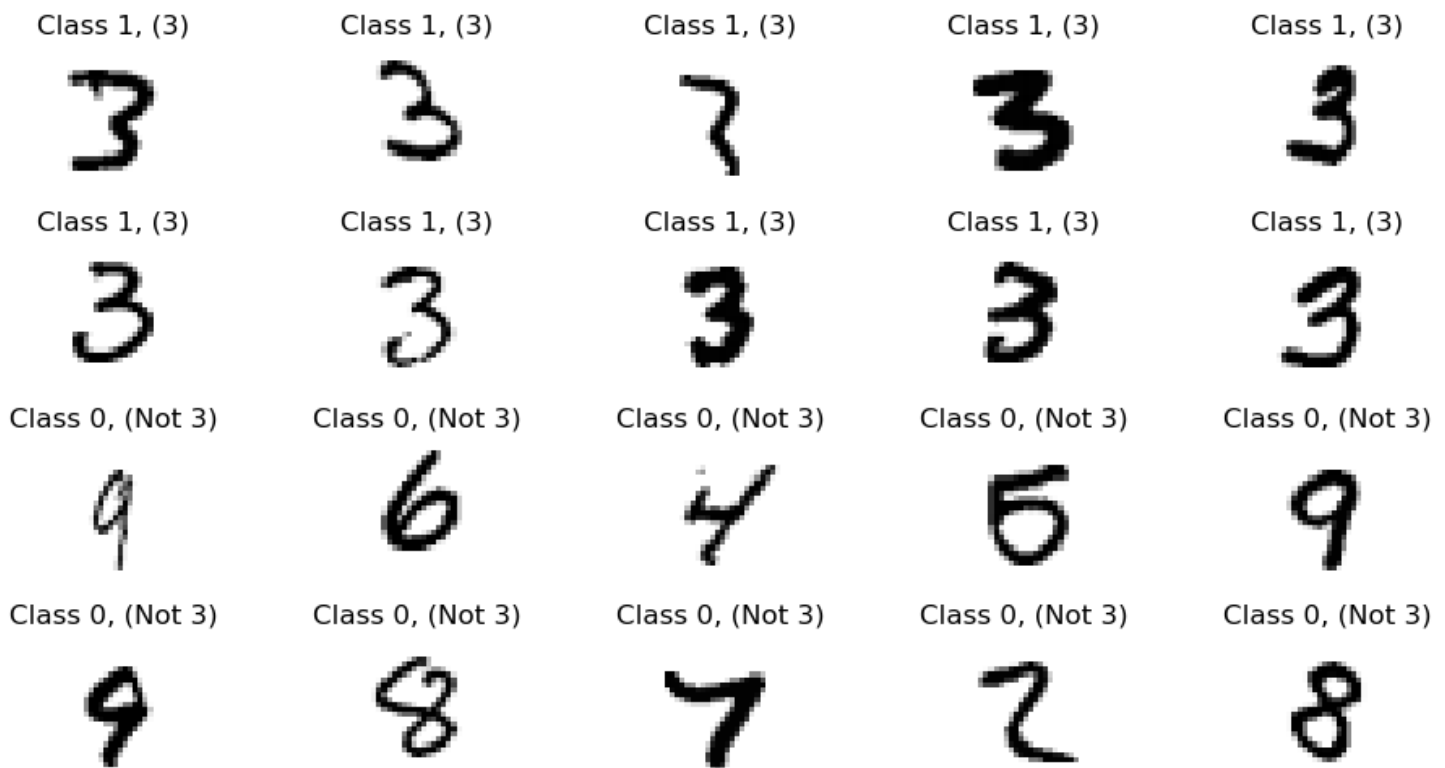
np.random.rand(268)
random_0 = np.where(y_train == 0)[0]
random_1 = np.where(y_train == 1)[0]
list_not_3 = np.random.choice(random_0, size=10, replace=False).tolist()
list_3 = np.random.choice(random_1, size=10, replace=False).tolist()

fig, axes = plt.subplots(num_row, num_col, figsize=(10, 5))

# 3
for i, idx in enumerate(list_3):
    ax = axes[i // num_col, i % num_col]
    ax.imshow(np.reshape(X_train[idx], (28, 28)), cmap="gray_r")
    ax.set_title("Class 1, (3)")
    ax.axis("off")

# not 3
for i, idx in enumerate(list_not_3):
    ax = axes[(i + len(list_3)) // num_col, (i + len(list_3)) % num_col]
    ax.imshow(np.reshape(X_train[idx], (28, 28)), cmap="gray_r")
    ax.set_title("Class 0, (Not 3)")
    ax.axis("off")

plt.tight_layout()
plt.show()
```



(b)

- For the train data, 53871 examples are present for class 0 (not 3) and 6129 are present for class 1 (3).
- The fraction of samples that are positive is $\sim 1/10$ or 10%
- The issues that this might is a biased model that is overfit and not generalizable.

```

In [ ]: u_v, counts = np.unique(y_train, return_counts=True)

colors = ["gold", "purple"]
plt.figure(figsize=(5, 5))

for i, (value, count) in enumerate(zip(u_v, counts)):
    plt.bar(i, count, color=colors[value])

plt.xticks([])

# Add count value on top of each bar
for i, count in enumerate(counts):
    plt.text(i, count, count, ha="center", va="bottom")

plt.legend(["Not 3", "3"], title="Unique Values", loc="upper right")

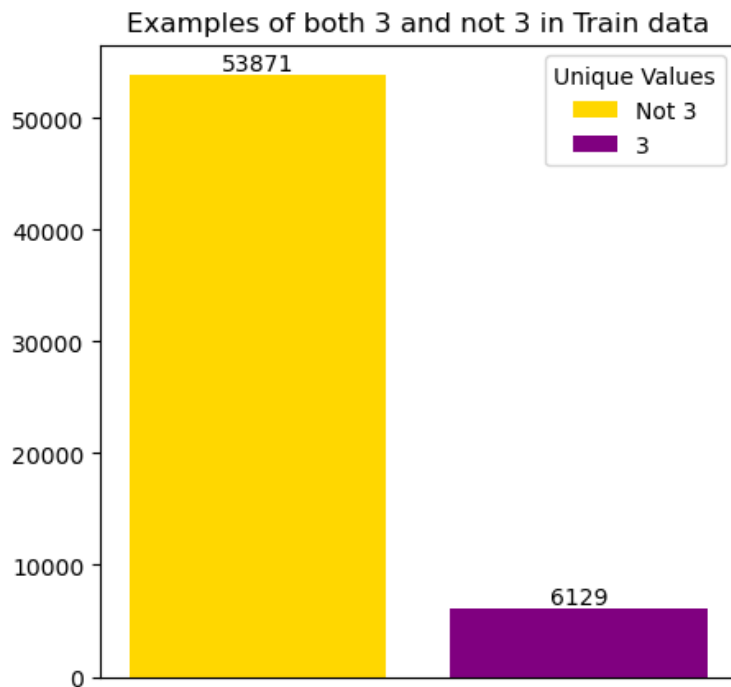
plt.title("Examples of both 3 and not 3 in Train data")

```

```

Out[ ]: Text(0.5, 1.0, 'Examples of both 3 and not 3 in Train data')

```



```

In [ ]: 6129 / (6129 + 53871)

```

```

Out[ ]: 0.10215

```

(c)

The value of C that seems best for this problem is 10^{-2}

```

In [ ]: from sklearn.metrics import auc
        from sklearn.metrics import roc_curve
        from sklearn.metrics import f1_score
        from sklearn.linear_model import LogisticRegression
        from sklearn.metrics import log_loss

        c_list = np.logspace(-4, 4, 20)
        param_list = []
        logloss_list = []
        auc_list = []
        f1_list = []

        for c in c_list:
            model = LogisticRegression(penalty="l1", C=c, solver="liblinear", random_state=42)
            model.fit(X_train, y_train)
            y_hat = model.predict(X_test)

```

```

# print('C:', c)

# params
non_zero_parms = np.count_nonzero(model.coef_)
param_list.append(non_zero_parms)

# log loss
logloss = log_loss(y_test, y_hat)
# print('logloss', logloss)
logloss_list.append(logloss)

# AUC
fpr, tpr, thresholds = roc_curve(y_test, y_hat)
my_auc = auc(fpr, tpr)
# print('auc', my_auc)
auc_list.append(my_auc)

# f1
my_f = f1_score(y_test, y_hat, zero_division=1.0)
# print('f1 score:', my_f)
f1_list.append(my_f)

```

```

In [ ]: import matplotlib.pyplot as plt

fig, axs = plt.subplots(2, 2, figsize=(8, 8))

# Parameters
axs[0, 0].plot(c_list, param_list)
axs[0, 0].set_xscale("log")
axs[0, 0].set_xlabel("Regularization parameter (C)")
axs[0, 0].set_ylabel("Number of non-zero parameters")
axs[0, 0].set_title("Non Zero Parameters")

# cross en
axs[0, 1].plot(c_list, logloss_list)
axs[0, 1].set_xscale("log")
axs[0, 1].set_xlabel("Regularization parameter (C)")
axs[0, 1].set_ylabel("Cross Entropy")
axs[0, 1].set_title("Cross Entropy")

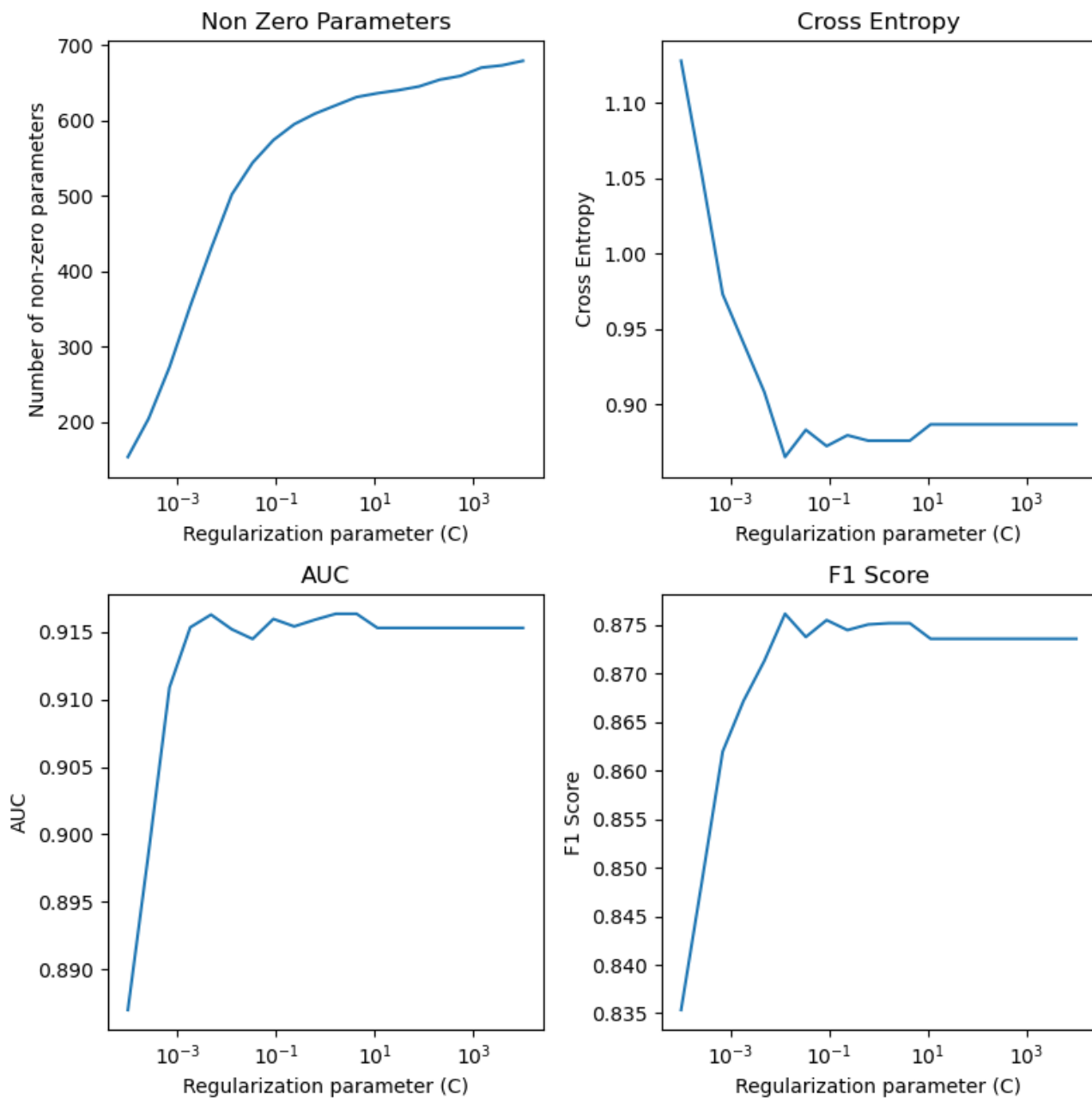
# AUC
axs[1, 0].plot(c_list, auc_list)
axs[1, 0].set_xscale("log")
axs[1, 0].set_xlabel("Regularization parameter (C)")
axs[1, 0].set_ylabel("AUC")
axs[1, 0].set_title("AUC")

# F1
axs[1, 1].plot(c_list, f1_list)
axs[1, 1].set_xscale("log")
axs[1, 1].set_xlabel("Regularization parameter (C)")
axs[1, 1].set_ylabel("F1 Score")
axs[1, 1].set_title("F1 Score")

# Adjust layout
plt.tight_layout()

# Show plot
plt.show()

```



```
In [ ]: f1_list.index(max(f1_list))
        c_list[5]
```

```
Out[ ]: 0.012742749857031334
```

(d)

```
In [ ]: from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
        from sklearn.ensemble import RandomForestClassifier
        from sklearn.metrics import precision_recall_curve

        """Model 1 Logistic Regression C=1e100"""
        lg_1 = LogisticRegression(penalty="l1", C=1e100, solver="liblinear", random_state=42)
        lg_1.fit(X_train, y_train)
        y_hat_lg1 = lg_1.predict_proba(X_test)[: , 1]

        fpr_lg1, tpr_lg1, _ = roc_curve(y_test, y_hat_lg1)
        roc_auc_lg1 = auc(fpr_lg1, tpr_lg1)

        p_lg1, r_lg1, _ = precision_recall_curve(y_test, y_hat_lg1)

        """Model 2 Logistic Regression C=1e-2"""
```

```

lg_2 = LogisticRegression(penalty="l1", C=1e-2, solver="liblinear", random_state=42)
lg_2.fit(X_train, y_train)
y_hat_lg2 = lg_2.predict_proba(X_test)[: , 1]

fpr_lg2, tpr_lg2, _ = roc_curve(y_test, y_hat_lg2)
roc_auc_lg2 = auc(fpr_lg2, tpr_lg2)

p_lg2, r_lg2, _ = precision_recall_curve(y_test, y_hat_lg2)

"""Model 3 Linear Discriminant Analysis (LDA) """
lda = LinearDiscriminantAnalysis()
lda.fit(X_train, y_train)
y_hat_lda = lda.predict_proba(X_test)[: , 1]

fpr_lda, tpr_lda, _ = roc_curve(y_test, y_hat_lda)
roc_auc_lda = auc(fpr_lda, tpr_lda)

p_lda, r_lda, _ = precision_recall_curve(y_test, y_hat_lda)

"""Model 4 Random Forest (RF) classifier"""
rf = RandomForestClassifier()
rf.fit(X_train, y_train)
y_hat_rf = rf.predict_proba(X_test)[: , 1]

fpr_rf, tpr_rf, _ = roc_curve(y_test, y_hat_rf)
roc_auc_rf = auc(fpr_rf, tpr_rf)

p_rf, r_rf, _ = precision_recall_curve(y_test, y_hat_rf)

```

```

In [ ]: plt.figure(figsize=(6.5, 6.5))

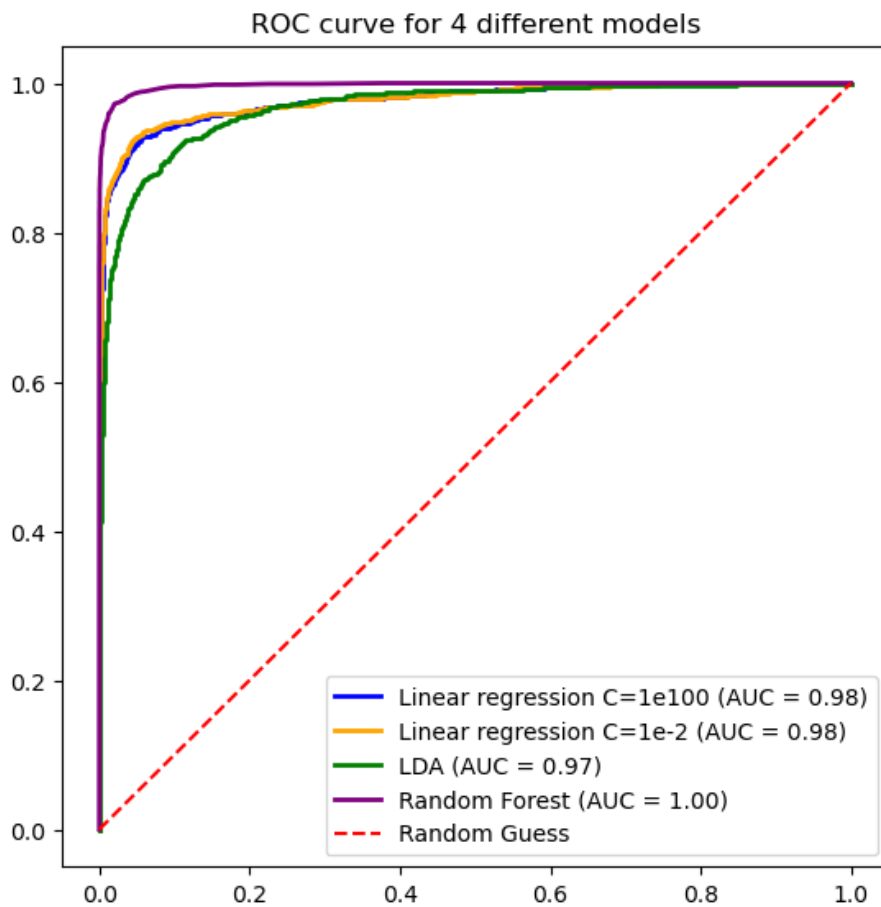
plt.plot(
    fpr_lg1,
    tpr_lg1,
    color="blue",
    lw=2,
    label="Linear regression C=1e100 (AUC = %0.2f)" % roc_auc_lg1,
)
plt.plot(
    fpr_lg2,
    tpr_lg2,
    color="orange",
    lw=2,
    label="Linear regression C=1e-2 (AUC = %0.2f)" % roc_auc_lg2,
)
plt.plot(fpr_lda, tpr_lda, color="green", lw=2, label="LDA (AUC = %0.2f)" % roc_auc_lda)
plt.plot(
    fpr_rf,
    tpr_rf,
    color="purple",
    lw=2,
    label="Random Forest (AUC = %0.2f)" % roc_auc_rf,
)

plt.plot([0, 1], [0, 1], color="red", linestyle="--", label="Random Guess")

plt.title("ROC curve for 4 different models")

plt.legend(loc="lower right")
plt.show()

```

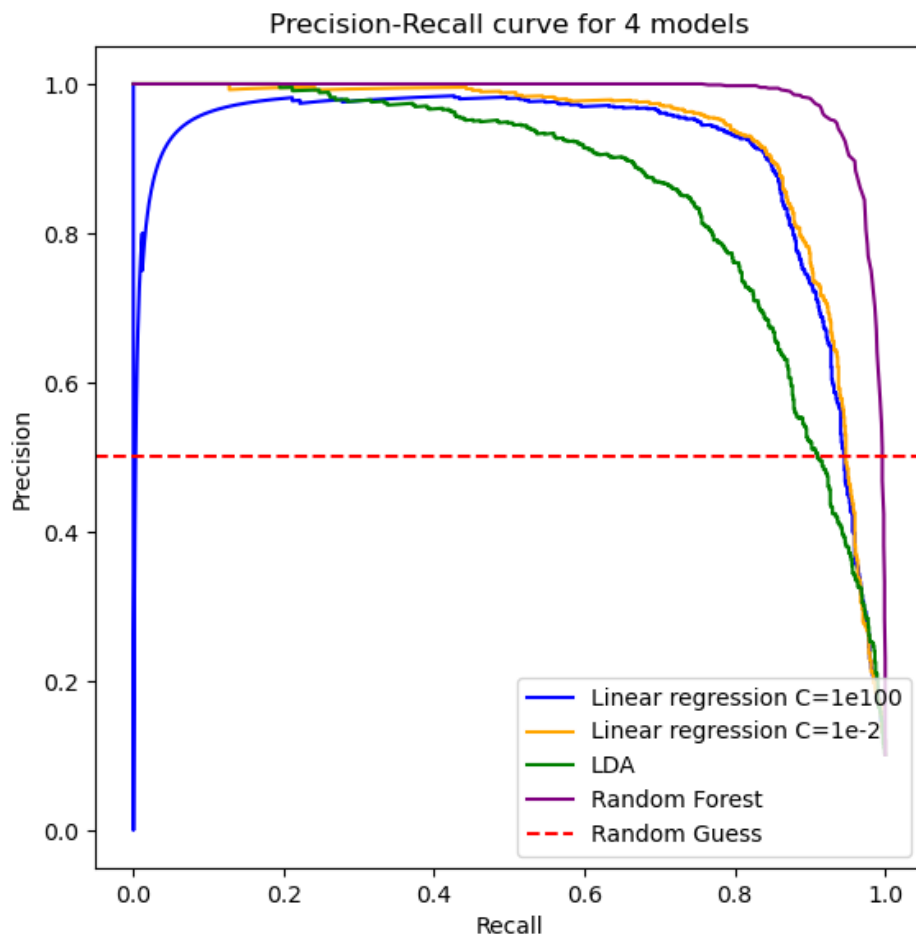
```
In [ ]: plt.figure(figsize=(6.5, 6.5))

plt.plot(r_lg1, p_lg1, label="Linear regression C=1e100", color="blue")
plt.plot(r_lg2, p_lg2, label="Linear regression C=1e-2", color="orange")
plt.plot(r_lda, p_lda, label="LDA", color="green")
plt.plot(r_rf, p_rf, label="Random Forest", color="purple")

plt.axhline(y=0.5, color="red", linestyle="--", label="Random Guess")

plt.legend(loc="lower right")
plt.xlabel("Recall")
plt.ylabel("Precision")
plt.title("Precision-Recall curve for 4 models")

plt.show()
```



- The four classifiers have a high performance per the AUC value, showing high discrimination between classes. LDA performs the worst but still with a 98% value. The regularization does not make a huge difference in the AUC for the logistic regression classifiers.
- The regularization of the logistic regression does not seem to carry much effect when comparing to a Random Forest or LDA classifier. However comparing between the two logistic regression models that have different regularization parameters, the regularization does have an effect between them.
- I would select the Random Forest classifier as its precision curve shows less sensitivity and it would work best for the unbalanced dataset.