

Assignment_3

February 18, 2022

1 Assignment 3 - Supervised Learning: model training and evaluation

1.1 *Satvik Kishore*

Netid: sk741

Names of students you worked with on this assignment: LIST HERE IF APPLICABLE (delete if not)

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](#), and is also linked to from the [course syllabus](#).

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

2 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

3 1

3.1 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 and 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!**

3.1.1 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).

Answer

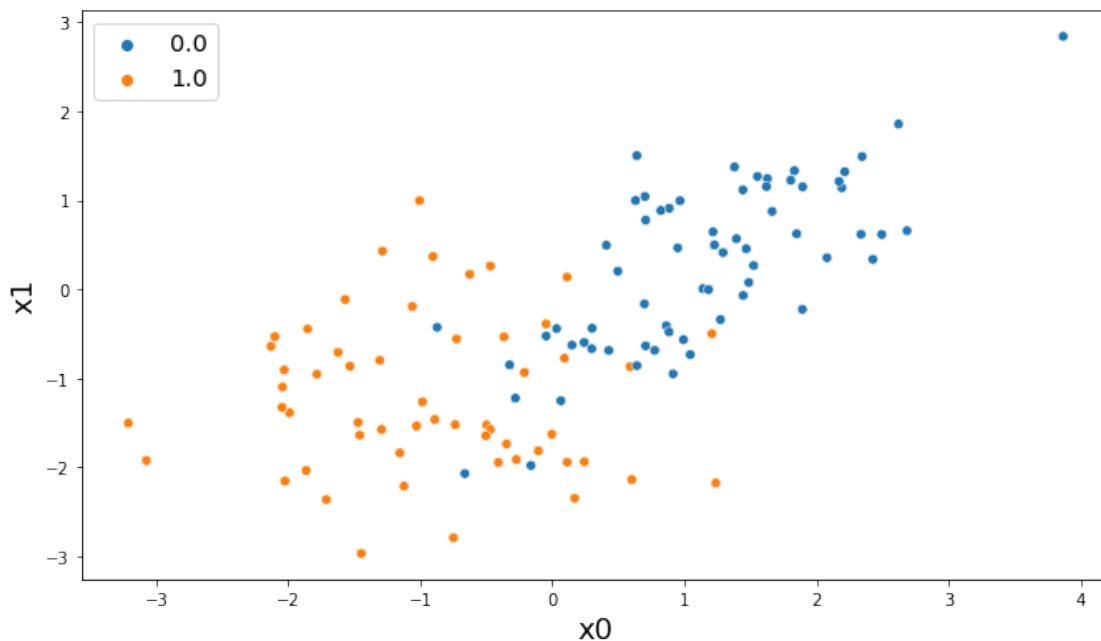
(a)

```
[1]: import numpy as np
from IPython.display import Markdown as md

dat = np.genfromtxt("./data/A3_Q1_data.csv", skip_header=1, dtype=float,
    ↪delimiter=",")
np.random.seed(12345)
np.random.shuffle(dat)
train_x = dat[0:120,0:2]
train_y = dat[0:120,2]
test_x = dat[120:,0:2]
test_y = dat[120:,2]
```

```
[2]: import seaborn as sb
import matplotlib.pyplot as plt

fig_scale = 0.6
fig, axs = plt.subplots(1)
fs = 18
fig.set_size_inches(18.5*fig_scale, 10.5*fig_scale)
sb.scatterplot(x=train_x[:,0], y=train_x[:,1], hue=train_y, ax=axs)
axs.set_xlabel("x0", fontsize=fs)
axs.set_ylabel("x1", fontsize=fs)
plt.legend(fontsize='x-large', title_fontsize=fs)
None
```



The classes do not appear to be completely separable linearly, but a linear separator with a negative slope could separate the classes with a reasonable accuracy. Thus logistic regression is a good choice as probability of observation being 1 (vs 0) increases as we move towards the upper right.

(b)

The data does not have any missing values. The variables x_1 and x_2 have similar ranges, and y is binary (0,1). This data does not need any preprocessing.

3.1.2 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 $\mathbf{w}^\top \mathbf{x}_i = w_0x_{i,0} + w_1x_{i,1} + w_2x_{i,2}$ 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.

3.1.3 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.

(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 $w_0x_{i,0} + w_1x_{i,1} + w_2x_{i,2}$ when it appears each time (where $x_{i,0} = 1$ for all i). You are also welcome to use $\sigma()$ 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.

Answer

(c)

$$\begin{aligned} L(\mathbf{w}|\mathbf{y}, \mathbf{X}) &= \prod_{i=1}^N P(y_i|\mathbf{x}_i) \\ &= \prod_{i=1}^N [1 - \sigma(\mathbf{w}^\top \mathbf{x}_i)]^{1-y_i} [\sigma(\mathbf{w}^\top \mathbf{x}_i)]^{y_i} \end{aligned}$$

(d)

$$\begin{aligned} C(\mathbf{w}) &= -\log[L(\mathbf{w}|\mathbf{y}, \mathbf{X})] \\ &= -\log\left[\prod_{i=1}^N [1 - \sigma(\mathbf{w}^\top \mathbf{x}_i)]^{1-y_i} [\sigma(\mathbf{w}^\top \mathbf{x}_i)]^{y_i}\right] \\ &= -\log \prod_{i=1}^N [1 - \sigma(\mathbf{w}^\top \mathbf{x}_i)]^{1-y_i} - \sum_{i=1}^N \log [\sigma(\mathbf{w}^\top \mathbf{x}_i)]^{y_i} \\ &= -\log \prod_{i=1}^N [1 - \sigma(\mathbf{w}^\top \mathbf{x}_i)]^{1-y_i} - \sum_{i=1}^N \log [\sigma(\mathbf{w}^\top \mathbf{x}_i)]^{y_i} \\ &= -\sum_{i=1}^N y_i \log(\sigma(\mathbf{w}^\top \mathbf{x}_i)) + (1 - y_i) \log(1 - \sigma(\mathbf{w}^\top \mathbf{x}_i)) \end{aligned}$$

(e)

$$\begin{aligned} \nabla_{\mathbf{w}} C(\mathbf{w}) &= \nabla_{\mathbf{w}} \sum_{i=1}^N y_i \log(1 + e^{-\mathbf{w}^\top \mathbf{x}_i}) + (1 - y_i) \log(1 + e^{\mathbf{w}^\top \mathbf{x}_i}) \\ &= \sum_{i=1}^N \frac{y_i e^{-\mathbf{w}^\top \mathbf{x}_i}}{1 + e^{-\mathbf{w}^\top \mathbf{x}_i}} (-\mathbf{x}_i) + \sum_{i=1}^N \frac{y_i e^{\mathbf{w}^\top \mathbf{x}_i}}{1 + e^{\mathbf{w}^\top \mathbf{x}_i}} (\mathbf{x}_i) \\ &= \sum_{i=1}^N (\mathbf{x}_i) [(1 - y_i) \sigma(\mathbf{w}^\top \mathbf{x}_i) - y_i \sigma(-\mathbf{w}^\top \mathbf{x}_i)] \\ &= \sum_{i=1}^N (\mathbf{x}_i) [(1 - y_i) \sigma(\mathbf{w}^\top \mathbf{x}_i) - y_i (1 - \sigma(\mathbf{w}^\top \mathbf{x}_i))] \\ &= \sum_{i=1}^N (\mathbf{x}_i) [\sigma(\mathbf{w}^\top \mathbf{x}_i) - y_i] \\ &= \begin{bmatrix} \sum_{i=1}^N x_{i,0} [\sigma(\mathbf{w}^\top \mathbf{x}_i) - y_i] \\ \sum_{i=1}^N x_{i,1} [\sigma(\mathbf{w}^\top \mathbf{x}_i) - y_i] \\ \sum_{i=1}^N x_{i,2} [\sigma(\mathbf{w}^\top \mathbf{x}_i) - y_i] \end{bmatrix} \end{aligned}$$

(f)

$$\begin{aligned} \mathbf{w}^{(k+1)} &= \mathbf{w}^{(k)} - \eta \nabla_{\mathbf{w}} C(\mathbf{w}^{(k)}) \\ &= \mathbf{w}^{(k)} - \eta \sum_{i=1}^N (\mathbf{x}_i) [\sigma(\mathbf{w}^{(k)\top} \mathbf{x}_i) - y_i] \end{aligned}$$

Therefore,

$$\begin{bmatrix} w_0^{(k+1)} \\ w_1^{(k+1)} \\ w_2^{(k+1)} \end{bmatrix} = \begin{bmatrix} w_0^{(k)} \\ w_1^{(k)} \\ w_2^{(k)} \end{bmatrix} - \eta \begin{bmatrix} \sum_{i=1}^N x_{i,0} [\sigma(\mathbf{w}^{(k)\top} \mathbf{x}_i) - y_i] \\ \sum_{i=1}^N x_{i,1} [\sigma(\mathbf{w}^{(k)\top} \mathbf{x}_i) - y_i] \\ \sum_{i=1}^N x_{i,2} [\sigma(\mathbf{w}^{(k)\top} \mathbf{x}_i) - y_i] \end{bmatrix}$$

3.1.4 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.).

```

[3]: # Logistic regression class
class Logistic_regression:
    # Class constructor
    def __init__(self, n_params=3):
        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

        pass

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

    # Cost function for an input set of weights
    def cost(self, X, y, w):
        # returns the average cross entropy cost
        X1 = X[y==1,:]
        X0 = X[y==0,:]
        out = np.log(1 - self.sigmoid(X0,w)).sum() + np.log(self.sigmoid(X1,w)).
        ↪ sum()
        return - out/X.shape[0]

    # 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
        grad = (X.T @ (self.sigmoid(X, self.w) - y))
        self.w = self.w - lr*grad
        self.saved_w.append(self.w)
        return np.linalg.norm(self.w)

    # Fit the logistic regression model to the data through gradient descent
    def fit(self, X, y, lr, w_init=None, 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)
        X = self.prepare_x(X)
        if w_init is not None:
            self.w = w_init
            self.saved_w.append(w_init)
        else:

```

```

        self.w = np.random.rand(3)
        self.saved_w.append(self.w)
    for i in range(max_iter):
        grad_mag = self.gradient_descent(X,y,lr)
        if verbose is True:
            print(self.w)
            pass
        if grad_mag < delta_thresh:
            break
        pass
    pass

    # Use the trained model to predict the confidence scores (prob of positive
    →class in this case)
    def predict_proba(self, X):
        # returns the confidence score for the each sample
        X = self.prepare_x(X)
        return self.sigmoid(X, self.w)

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

    # 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
        X = self.prepare_x(X)
        curve = []
        for w in self.saved_w:
            curve.append(self.cost(X, y, w))
        return curve

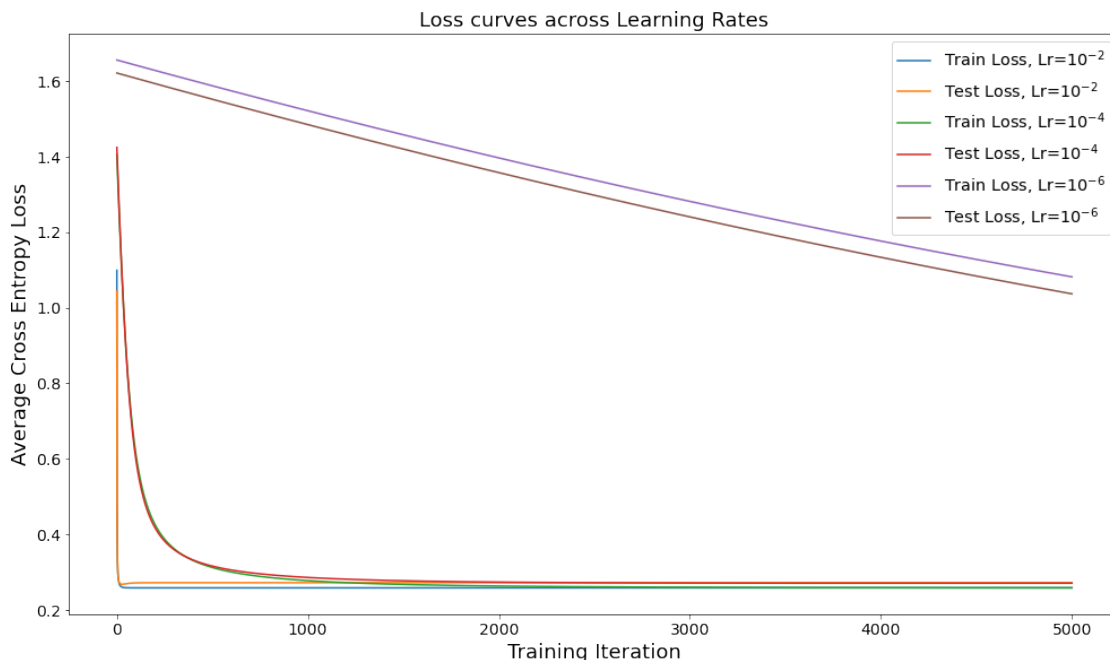
    # Appends a column of ones as the first feature to account for the bias term
    def prepare_x(self, X):
        # returns the X with a new feature of all ones (a column that is the
        →new column 0)
        return np.c_[np.ones(X.shape[0]), X]

```

(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 insight about how model training is proceeding. Show the learning curves for the gradient descent process for learning rates of $\{10^{-2}, 10^{-4}, 10^{-6}\}$. 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). Each of

the 6 resulting curves (train and test average cost for each learning rate) should be plotted on the same set of axes for direct comparison. *Note: make sure you're using average cost per sample, not total cost* - Try running this process for a really big learning rate for this problem: 10^0 . Look at the weights that the fitting process generates over the first 50 iterations and how they change. You may simply print these first 50 iterations as output or plot them. What happens and why? - What is the impact that the different values of learning has 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.

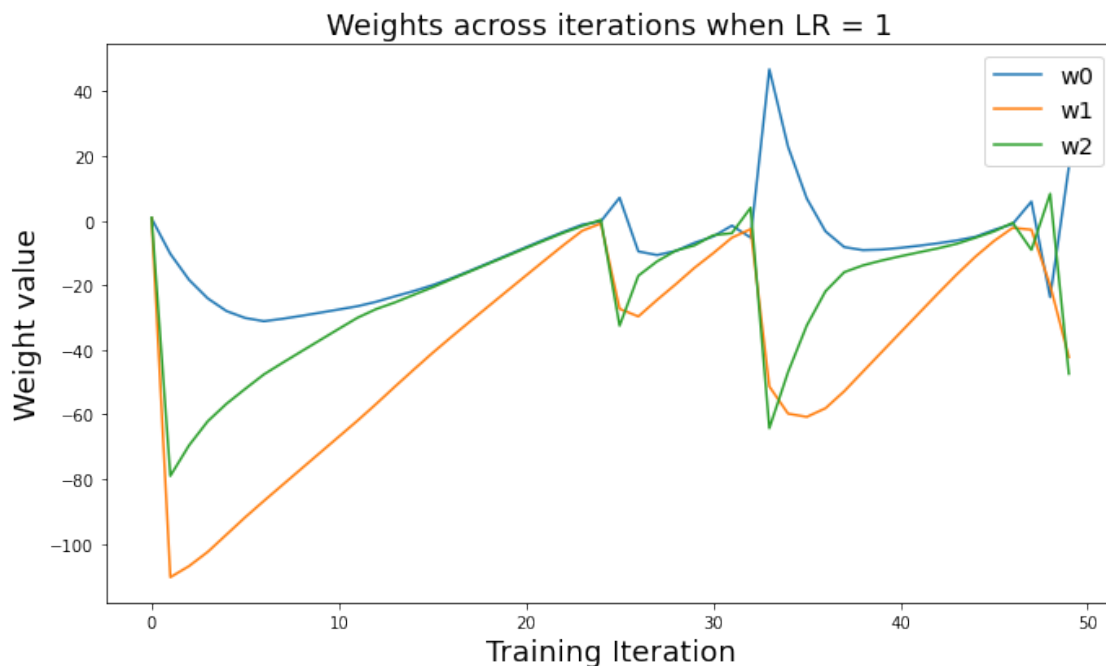
```
[4]: np.random.seed(12345)
fs = 18
fig,axs = plt.subplots(1)
axs.tick_params(axis='both', which='major', labelsize=14)
fig_scale = 0.9
fig.set_size_inches(18.5*fig_scale, 10.5*fig_scale)
for lr_s,lr in [("$10^{-2}$",1e-2), ("10^{-4}$",1e-4),("$10^{-6}$",1e-6)]:
    m = Logistic_regression()
    m.fit(train_x, train_y, lr=lr, verbose=False)
    axs.plot(range(len(m.saved_w)), m.learning_curve(train_x, train_y),
    ↪label="Train Loss, Lr="+lr_s)
    axs.plot(range(len(m.saved_w)), m.learning_curve(test_x, test_y),
    ↪label="Test Loss, Lr="+lr_s)
    pass
axs.set_xlabel("Training Iteration", fontsize=fs)
axs.set_ylabel("Average Cross Entropy Loss", fontsize=fs)
axs.set_title("Loss curves across Learning Rates", fontsize=fs)
axs.legend(fontsize='x-large', title_fontsize=fs)
None
```



We observe that when learning rate = 10^{-6} , the change in gradient is too small and the algorithm is not able to converge. The 10^{-4} learning rate is small enough to complete learning in the 5000 iterations limit but is not fast enough. The 10^{-2} manages to converge very early and is clearly the best hyperparameter choice. In these three choices, 10^{-2} would be the best.

```
[5]: fig,axs = plt.subplots(1)
fig_scale = 0.6
fig.set_size_inches(18.5*fig_scale, 10.5*fig_scale)
m = Logistic_regression()
m.fit(train_x, train_y, lr=1, verbose=False)
axs.plot(range(min(50,len(m.saved_w))), [i[0] for i in m.saved_w][0:50],
        ↪label="w0")
axs.plot(range(min(50,len(m.saved_w))), [i[1] for i in m.saved_w][0:50],
        ↪label="w1")
axs.plot(range(min(50,len(m.saved_w))), [i[2] for i in m.saved_w][0:50],
        ↪label="w2")
axs.set_xlabel("Training Iteration", fontsize=fs)
axs.set_ylabel("Weight value", fontsize=fs)
axs.legend(fontsize='x-large', title_fontsize=fs)
axs.set_title("Weights across iterations when LR = 1",fontsize=fs)
None

opt_lr = 0.01
```



For $lr = 10^0$, we observe weights to be changing very haphazardly, to extreme values. This indicates that this learning rate is too high, and that the weights overshoot the target minima at every iteration and will probably fail at learning the best parameters.

The best learning rates we observed were for 10^{-2} , as it seems to have converged the fastest at the optimal solution.

3.1.5 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 two plots corresponding to each model (logistic regression and KNN): one with the training data, and one for 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 - 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

```
[6]: from sklearn.neighbors import KNeighborsClassifier
      from matplotlib.colors import ListedColormap

      knn = KNeighborsClassifier(n_neighbors=7)
      knn.fit(train_x, train_y)
      logistic = LogisticRegression()
      logistic.fit(train_x, train_y, lr=opt_lr)

[7]: # construct cmap
      bg_cmap = ListedColormap(sb.color_palette("pastel",N=2))
      dot_cmap = sb.color_palette("dark",2)

      flex = 0.5
      h=0.02
      xx, yy = np.meshgrid(np.arange(dat[:,0].min() - flex,
                                     dat[:,0].max() + flex, h),
                           np.arange(dat[:,1].min() - flex,
                                     dat[:,1].max() + flex, h))
      grid = np.c_[xx.ravel(), yy.ravel()]
      gridpred_logistic = logistic.predict(grid)
      gridpred_logistic = gridpred_logistic.reshape(xx.shape)
```

```

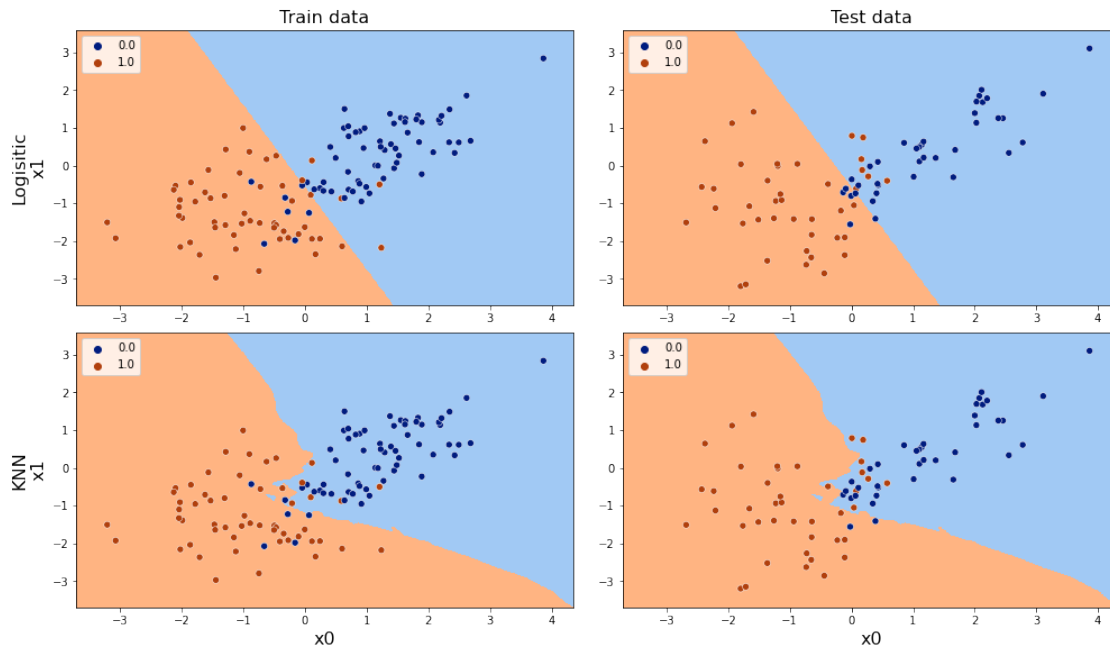
gridpred_knn = knn.predict(grid)
gridpred_knn = gridpred_knn.reshape(xx.shape)

fig, ax = plt.subplots(2,2)
fig.subplots_adjust(hspace=0.1, wspace=0.1)
fs = 16
fig_scale = 0.9
fig.set_size_inches(18.5*fig_scale, 10.5*fig_scale)
ax[0,0].contourf(xx, yy, gridpred_logistic,cmap=bg_cmap)
ax[0,1].contourf(xx, yy, gridpred_logistic,cmap=bg_cmap)
ax[1,0].contourf(xx, yy, gridpred_knn,cmap=bg_cmap)
ax[1,1].contourf(xx, yy, gridpred_knn,cmap=bg_cmap)

sb.scatterplot(x=train_x[:,0], y=train_x[:,1], hue=train_y, ax=ax[0,0],
    ↪palette=dot_cmap)
sb.scatterplot(x=test_x[:,0], y=test_x[:,1], hue=test_y, ax=ax[0,1],
    ↪palette=dot_cmap)
sb.scatterplot(x=train_x[:,0], y=train_x[:,1], hue=train_y, ax=ax[1,0],
    ↪palette=dot_cmap)
sb.scatterplot(x=test_x[:,0], y=test_x[:,1], hue=test_y, ax=ax[1,1],
    ↪palette=dot_cmap)

ax[0,0].set_ylabel("Logisitic\nx1", fontsize=fs)
ax[1,0].set_ylabel("KNN\nx1", fontsize=fs)
ax[1,0].set_xlabel("x0", fontsize=fs)
ax[1,1].set_xlabel("x0", fontsize=fs)
ax[0,0].set_title("Train data", fontsize=fs)
ax[0,1].set_title("Test data", fontsize=fs)
None

```



Here, Decision Boundary Threshold = 0.5

```
[8]: from sklearn.model_selection import StratifiedKFold

logistic_prob = np.zeros(dat.shape[0])
knn_prob = np.zeros(dat.shape[0])
skf = StratifiedKFold(n_splits=10)
for train_index, test_index in skf.split(dat[:,0:2], dat[:,2]):
    xtrain = dat[train_index, 0:2]
    ytrain = dat[train_index, 2]
    xtest = dat[test_index, 0:2]

    knn = KNeighborsClassifier(n_neighbors=7)
    knn.fit(xtrain, ytrain)
    logistic = LogisticRegression()
    logistic.fit(xtrain, ytrain, lr=opt_lr)

    logistic_prob[test_index] = logistic.predict_proba(xtest)
    knn_prob[test_index] = knn.predict_proba(xtest)[:,1]
pass
```

```
[9]: logistic_tpr = [(dat[:,2] * (logistic_prob >= th/100)).sum() / dat[:,2].sum()
    ↪ for th in range(0,101)]
logistic_fpr = [((1-dat[:,2]) * (logistic_prob >= th/100)).sum() / (dat.
    ↪ shape[0] - dat[:,2].sum()) for th in range(0,101)]
```

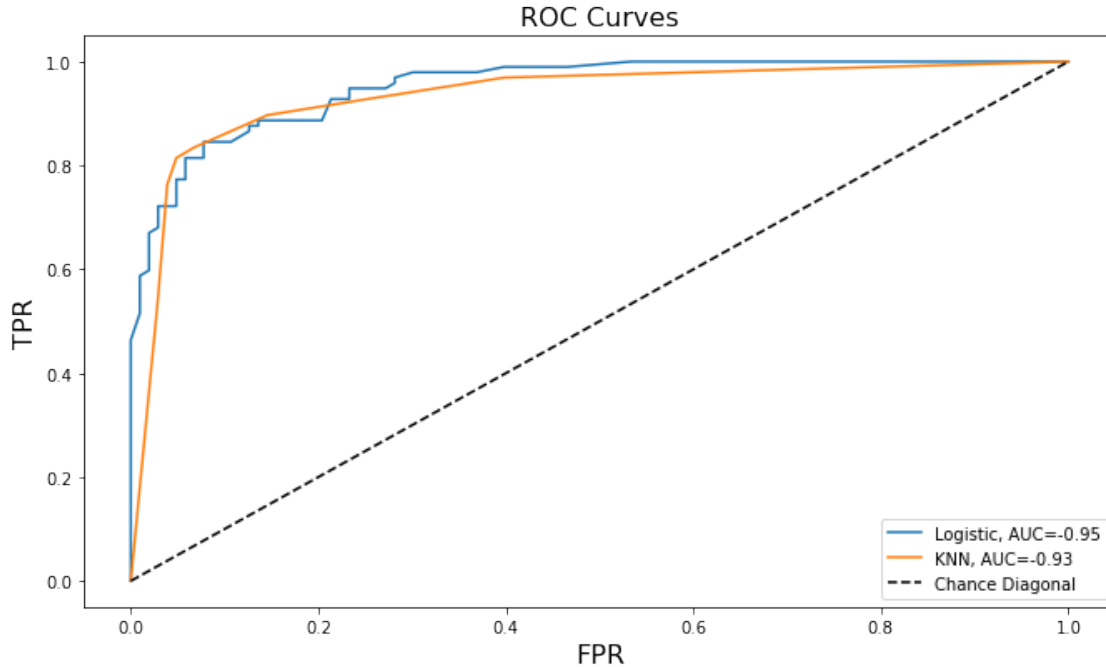
```

knn_tpr = [(dat[:,2] * (knn_prob >= th/100)).sum() / dat[:,2].sum() for th in
    range(0,101)]
knn_fpr = [(1-dat[:,2]) * (knn_prob >= th/100)).sum() / (dat.shape[0] - dat[:,
    2].sum()) for th in range(0,101)]
knn_tpr.append(0.0)
knn_fpr.append(0.0)

logistic_auc = 0
knn_auc = 0
for i in range(1,len(logistic_tpr)):
    logistic_auc += 0.5*(logistic_fpr[i] - logistic_fpr[i-1])*(logistic_tpr[i]
    + logistic_tpr[i-1])
    pass
for i in range(1,len(knn_tpr)):
    knn_auc += 0.5*(knn_fpr[i] - knn_fpr[i-1])*(knn_tpr[i] + knn_tpr[i-1])
    pass

fig,axs = plt.subplots(1)
fs = 16
fig_scale = 0.6
fig.set_size_inches(18.5*fig_scale, 10.5*fig_scale)
axs.plot(logistic_fpr, logistic_tpr, label="Logistic,
    AUC="+str(round(logistic_auc,2)))
axs.plot(knn_fpr, knn_tpr, label="KNN, AUC="+str(round(knn_auc,2)))
axs.plot([0,1], [0,1], "--",label="Chance Diagonal", color="black")
axs.set_title("ROC Curves", fontsize=fs)
axs.set_xlabel("FPR", fontsize=fs)
axs.set_ylabel("TPR", fontsize=fs)
axs.legend()
None

```



Cross Validation helps us evaluate model performance on unseen data, giving estimates of generalizability of the model. In this example, our sample size is small, and K Fold validation cross validation allows us to test our metrics on all data points.

In the decision boundary plots, KNN appears to mimic the data distribution better, but due to overlap between the distributions of the two categories, there is always some misclassification, which both models appear to suffer from equally.

The AUC of the ROC curve suggests that Logistic regression is marginally better, and we would choose logistic regression when it is not known if preference is to be given to False Negatives or False Positives. If costs associated with False Negatives/Positives were known, our answer could change as there are FPR/TPR values where KNN appears to perform better.

4 2

4.1 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) Using a logistic regression classifier, apply lasso regularization and retrain the model and evaluate its performance 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 20 different values of C . As you vary the regularization coefficient, Plot - 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. - 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.


```
[10]: # 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 = False

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()
    pass
#X_train /= 100
#X_test /= 100
```

ANSWER

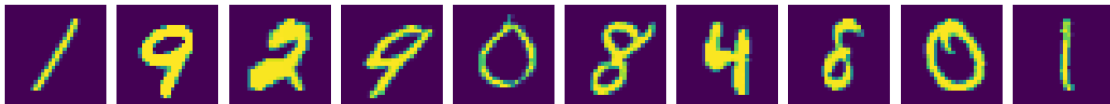
(a)

```
[11]: np.random.seed(1234)
samples0, samples1 = (np.random.choice(X_train[y_train==i].shape[0], 10) for i in
→[0,1])
samples0, samples1 = (X_train[y_train==i][s,:].reshape((-1,28,28)) for i,s in
→enumerate((samples0, samples1)))
fig, axs = plt.subplots(2,10, gridspec_kw={'hspace': 0.0, 'wspace': 0.1})
fig_scale = 1
```

```

fig.set_size_inches(18.5*fig_scale, 10.5*fig_scale)
#fig.subplots_adjust(hspace=0)
#fig.tight_layout()
for ii in range(10):
    axs[0,ii].imshow(samples0[ii,:,:])
    axs[1,ii].imshow(samples1[ii,:,:])
    axs[0,ii].set_axis_off()
    axs[1,ii].set_axis_off()
    pass
None

```



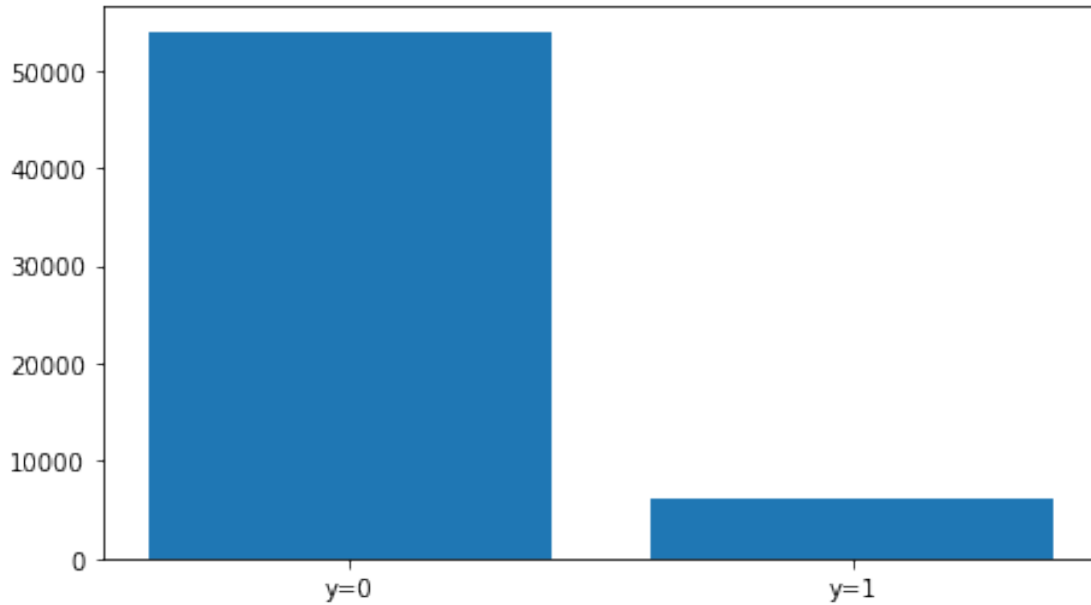
(b)

```

[12]: hh = [(y_train==d).sum() for d in [0,1]]
fig, ax = plt.subplots(1)
fig_scale = 0.4
fig.set_size_inches(18.5*fig_scale, 10.5*fig_scale)
ax.bar(x=["y=0", "y=1"], height = hh)
md(f"{100*hh[1]/X_train.shape[0]}% of the samples are positive")

```

[12]: 10.215% of the samples are positive



The low number of positive cases might cause imbalance issues and the classifier might not learn enough from the positive cases, that is give equal weightage to False Positives and False Negatives, which is not desirable

(c)

```
[13]: from sklearn.linear_model import LogisticRegression
      from sklearn.metrics import log_loss
      from sklearn.metrics import roc_auc_score as auc
      from sklearn.metrics import f1_score
      import sklearn
      import matplotlib.ticker as ticker

      lrms = []
      penalizers = np.logspace(-4,4,20)
      for l1c in penalizers:
          lrms.append(LogisticRegression(penalty = "l1", C = l1c, solver="liblinear",
          ↪random_state=42))
          lrms[-1].fit(X_train, y_train)
```

```
[14]: nparams = [(m.coef_!=0).sum() for m in lrms]
      log_losses = [log_loss(y_test, m.predict_proba(X_test)) for m in lrms]
      aucs = [auc(y_test, m.predict_proba(X_test)[:,-1]) for m in lrms]
      f1s = [f1_score(y_test, m.predict_proba(X_test)[:,-1]>=0.5) for m in lrms]

      fig_scale = 1
      fig, axs = plt.subplots(2,2)
```

```

fig.subplots_adjust(hspace=0.3, wspace=0.2)
fig.set_size_inches(18.5*fig_scale, 10.5*fig_scale)
for axi in axs:
    for ax in axi:
        ax.tick_params(axis='both', which='major', labelsize=15)
fs = 16

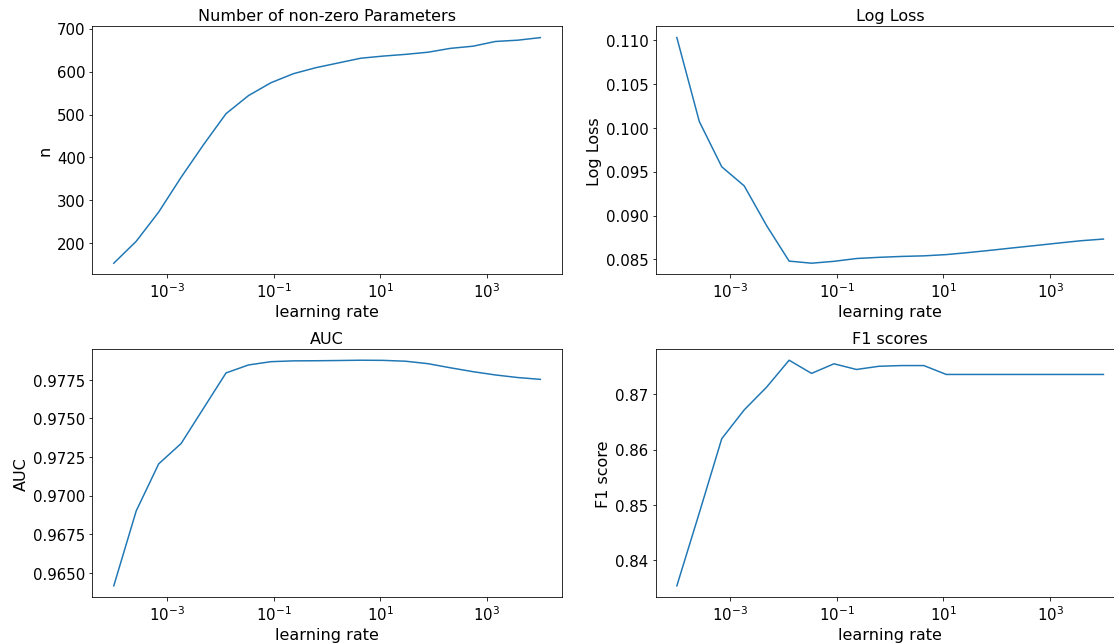
axs[0,0].plot(penalizers, nparams)
axs[0,0].set_title("Number of non-zero Parameters", fontsize=fs)
axs[0,1].plot(penalizers, log_losses)
axs[0,1].set_title("Log Loss", fontsize=fs)
axs[1,0].plot(penalizers, aucs)
axs[1,0].set_title("AUC", fontsize=fs)
axs[1,1].plot(penalizers, f1s)
axs[1,1].set_title("F1 scores", fontsize=fs)

axs[0,0].set_xscale('log')
axs[0,1].set_xscale('log')
axs[1,0].set_xscale('log')
axs[1,1].set_xscale('log')

axs[0,0].set_xlabel("learning rate", fontsize=fs)
axs[1,0].set_xlabel("learning rate", fontsize=fs)
axs[0,1].set_xlabel("learning rate", fontsize=fs)
axs[1,1].set_xlabel("learning rate", fontsize=fs)

axs[0,0].set_ylabel("n", fontsize=fs)
axs[1,0].set_ylabel("AUC", fontsize=fs)
axs[1,1].set_ylabel("F1 score", fontsize=fs)
axs[0,1].set_ylabel("Log Loss", fontsize=fs)
None

```



```
[15]: best = np.argmax(np.array(aucs))
      best_c = penalizers[best]
      md(f"I chose to optimize on the AUC scores. The best AUC score was found to be_
      ↪{best}, with the best value of C being {round(best_c,2)}")
```

[15]: I chose to optimize on the AUC scores. The best AUC score was found to be 11, with the best value of C being 4.28

```
[16]: best_c = 10
```

```
[17]: from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
      from sklearn.ensemble import RandomForestClassifier
      models = {}
      models['logist100'] = LogisticRegression(penalty = "l1", C = 1e100,
      ↪solver="liblinear")
      models['logist100'].fit(X_train, y_train)

      models['logistoptim'] = LogisticRegression(penalty = "l1", C = best_c,
      ↪solver="liblinear")
      models['logistoptim'].fit(X_train, y_train)

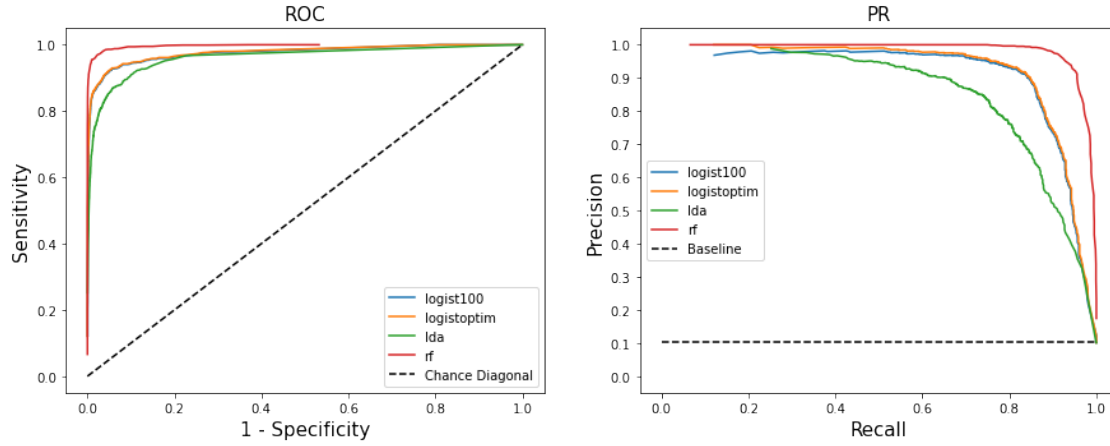
      models['lda'] = LinearDiscriminantAnalysis()
      models['lda'].fit(X_train, y_train)

      models['rf'] = RandomForestClassifier()
      models['rf'].fit(X_train, y_train)
```

None

```
[18]: tprs = {}
fprs = {}
precisions = {}
pr_baseline = y_train.sum()/y_train.shape[0]
probs_record = {}
ths = np.linspace(0,1,1000)
ths[1] = 1e-6
ths = ths[:-1]
for mk in models.keys():
    probs = models[mk].predict_proba(X_test)[:,-1]
    tprs[mk] = [(y_test * (probs > th)).sum() / y_test.sum() for th in ths]
    fprs[mk] = [((1-y_test) * (probs > th)).sum() / (y_test.shape[0] - y_test.
    ↳sum()) for th in ths]
    precisions[mk] = [(y_test * (probs > th)).sum() / (probs > th).sum() for th
    ↳in ths]
    probs_record[mk] = probs
    pass

fig, axs = plt.subplots(1,2)
fig_scale = 0.8
fs = 15
fig.set_size_inches(18.5*fig_scale, 10.5*fig_scale/1.6)
axs[0].set_title("ROC", fontsize=fs)
axs[1].set_title("PR", fontsize=fs)
for mk in models.keys():
    axs[0].plot(fprs[mk], tprs[mk], label=mk)
    axs[1].plot(tprs[mk], precisions[mk], label=mk)
axs[0].plot([0,1], [0,1], '--', label="Chance Diagonal", color="black")
axs[1].plot([0,1], [pr_baseline,pr_baseline], '--', label="Baseline",
    ↳color="black")
axs[0].legend()
axs[1].legend()
axs[1].yaxis.set_major_locator(ticker.MultipleLocator(0.1))
axs[1].set_xlim([-0.05, 1.05])
axs[1].set_ylim([-0.05, 1.05])
axs[0].set_xlabel("1 - Specificity", fontsize=fs)
axs[1].set_xlabel("Recall", fontsize=fs)
axs[0].set_ylabel("Sensitivity", fontsize=fs)
axs[1].set_ylabel("Precision", fontsize=fs)
None
```



In these plots, we can see that our hyperparameter optimized logistic regression is only barely better than the model with arbitrarily high C value. However, The Randomforest model performs the best on the holdout data and we would expect it to generalize the best. The LDA model was worse than our logistic regression models.