

Logistic Regression

Definition

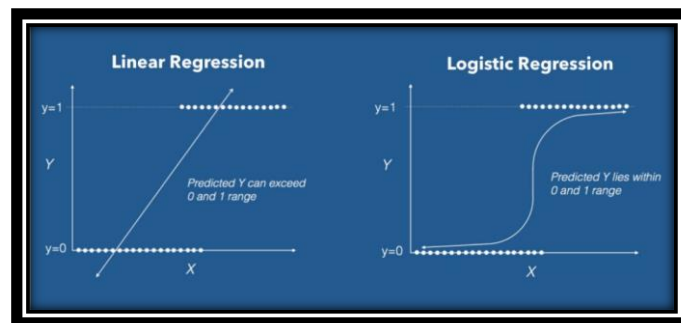
Logistic regression is a supervised classification algorithm used for categorical variables which can take only discrete values as the output such as 0/1, True/False, Yes/No, etc.

Use of logistic regression for classification

Logistic regression creates a regression model by computing the probability of the given observation belonging to a particular category. It becomes a classification algorithm when a decision threshold is set making threshold value an important element which in turn is affected by the values of precision and recall.

Linear Regression vs Logistic Regression

Linear regression and logistic regression are similar in the essence that both are parametric regressions but that is it. Linear Regression computes continuous values for the target/output variable whereas Logistic Regression computes discrete values for the target/output variable. Also, Linear Regression uses a linear function which results in a straight line graph whereas Logistic Regression uses a sigmoid function which results in a curved line graph.



_(Graph)

Types of logistic regression

There are 3 types of Logistic regression based on the number of categories which are as follows:

- Binomial: The target variable has only two possible output values. For eg: Yes/No, 0/1, Win/Loss
 - Multinomial: The target variable has three possible unordered output values. For eg: “Disease A” vs “Disease B” vs “Disease C”
 - Ordinal: The target variable has ordered output values. For eg: “Good”, “Better”, “Best” or “Excellent”
- (Types of Logistic Regression)

Which type are we using for our purpose?

We need binomial logistic regression since there are only two classes/categories (Fire: Yes/No)

Core of logistic regression

Logistic regression has three main functions:

- Hypothesis
- Cost
- Gradient Descent

The notation that is going to be used for the algorithm explanation is as follows:

X	input data matrix of shape m x n
m	number of observations
n	number of features

y	output/target value
x(i), y(i)	i th training example
θ	regression coefficient
α	Learning rate

Hypothesis

To start building the hypothesis for the logistic regression algorithm, regression coefficients need to be introduced.

Regression coefficients are associated with the features which tell how important a feature is for the classification problem.

The hypothesis for linear regression is:

$$h_{\theta}(x) = \theta^T X$$

Logistic regression also uses the hypothesis function which is slightly different because the hypothesis for linear regression is a linear function that can have values greater than 1 or less than 0. Therefore, the sigmoid function is used which maps the predicted real values to other values between 0 and 1. In other words, the sigmoid function maps predictions to probabilities.

The formula for **sigmoid function** (σ) is:

$$\sigma(z) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + \exp(-z)}$$

For logistic regression, $\theta = \theta_0 + \theta_1 X_1 + \theta_2 X_2 + \dots + \theta_n X_n$

$$h_{\theta}(x) = \sigma(\theta) = \sigma(\theta_0 + \theta_1 X_1 + \theta_2 X_2 + \dots + \theta_n X_n)$$

i.e.

$$h_{\theta}(x) = \frac{1}{1 + e^{\theta^T x}}$$

Cost Function

The objective of the cost function is optimization i.e., evaluate the errors the logistic model is going to make so that we can develop a model with minimum error.

The cost function for Linear regression is:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

However, the cost function for linear regression can't be used for logistic regression as it gives multiple local minimums whereas we need to minimize the cost value and find the global minimum.

For logistic regression, the cost function is as follows:

$$Cost(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases}$$

The graph for the cost function is shown below:



For a binary classification problem, the cost function can be simplified as:

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right]$$

Gradient Descent

Gradient descent is an optimization technique that finds the optimal coefficients and minimizes the cost function. In this, we iterate through the training set repetitively to find the gradient of the cost function at the current point and move in the opposite direction.

$$\theta_j = \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) * x_j^{(i)}$$

Design Decision

- a) **Threshold Value:** The selection of threshold value depends on error preferences which are: false positive and false negative. High threshold value results in a more false negative error and a low threshold value results in a more false positive error.

In the wildfires.txt data, we need to predict fire where we can't afford to have more false negative but also, it would be an overhead to have more false positives as it would take up more resources. To reach a balance, we have chosen the threshold value to be 0.5.

$$P(Y=1|X) \geq 0.5, \text{ class} = \text{Yes, and if } P(Y=1|X) < 0.5 \text{ class} = \text{No}$$

- b) **Type of gradient descent:** Batch guide gradient descent is used for this data. The reason behind choosing this is as follows:
- Coefficients are updated after computing the error gradient for the entire training set.
 - It makes the least noisy updates to the coefficients when compared to other gradient descents.
 - Since for batch guide gradient descent, we need to loop through the entire data which makes the calculation slow and expensive and takes time to reach convergence but, in our case, there are just 136 observations in the training set which is not huge.
- c) **Learning rate:** It determines the magnitude of the amount to move in gradient descent. The change in coefficient is learning rates times the gradient (Martin, 2021). This we need to specify explicitly. If it is too small, the algorithm will converge slowly but if it is large, it can lead to non-convergence i.e., bouncing back and forth between the convex function of gradient descent but never reaching the local minimum. The learning rate used in our code is **0.01**.
- d) **Convergence:** Gradient descent minimizes the cost and when the minimum cost is achieved, it means it has converged. In the code, convergence is used to find the minimum cost instead of iterating it a specific number of times. The value of convergence used in our code is **0.01**.
- e) **Normalization:** It helps to model the data correctly by creating new values. It maintains the general distribution and ratios in the source data. It keeps values within a scale that is applied to all numeric columns in the model. The normalization used in our code is Z-normalization. It normalizes every value of the dataset in a way that mean of all the values is 0 and the standard deviation is 1.

$$Z = \frac{x - \mu}{\sigma}$$

Z = standard score

x = observed value

μ = mean of the sample

σ = standard deviation of the sample

- f) Define $x_0 = 1$ in the input matrix as we need to perform matrix product between feature vector and coefficients vector and there was no value to be multiplied with θ_0 .

Comparison between Sklearn Implementation and Custom Implementation

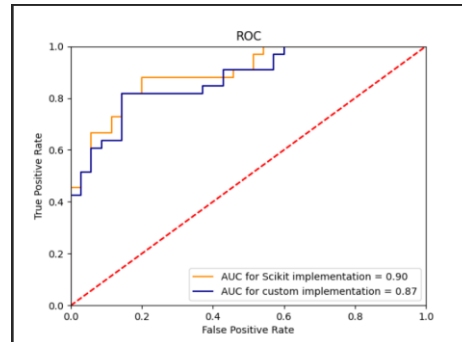
Scikit package has been used as a reference package.

As instructed, the data is split randomly into training data and test data in the ratio 2:1, and the algorithms are executed 10 times.

The difference between the accuracies achieved by the two implementations for different learning rates and convergence is shown in the table below.

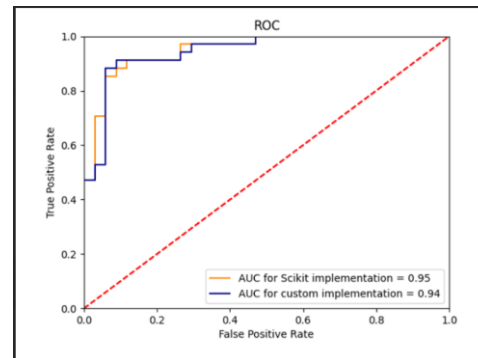
accuracy_sk	accuracy_lr
83.82352941	75
83.82352941	83.82352941
76.47058824	79.41176471
94.11764706	82.35294118
79.41176471	79.41176471
77.94117647	70.58823529
83.82352941	72.05882353
89.70588235	88.23529412
79.41176471	76.47058824
82.35294118	72.05882353
Average	
83.08823529	77.94117647

$\alpha = 0.01$, convergence = 0.01



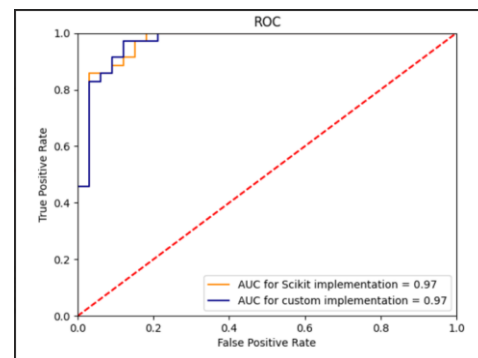
accuracy_sk	accuracy_lr
88.23529412	88.23529412
88.23529412	88.23529412
82.35294118	82.35294118
79.41176471	79.41176471
86.76470588	86.76470588
86.76470588	80.88235294
85.29411765	80.88235294
85.29411765	76.47058824
83.82352941	83.82352941
88.23529412	89.70588235
Average	
85.44117647	83.67647059

$\alpha = 0.01$, convergence = 0.001



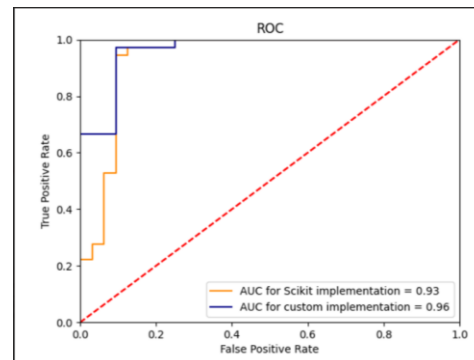
accuracy_sk	accuracy_lr
82.35294118	89.70588235
85.29411765	79.41176471
83.82352941	89.70588235
82.35294118	82.35294118
89.70588235	80.88235294
82.35294118	86.76470588
86.76470588	88.23529412
85.29411765	80.88235294
77.94117647	88.23529412
91.17647059	91.17647059
Average	
84.70588235	85.73529412

$\alpha = 0.1$, convergence = 0.001



accuracy_sk	accuracy_lr
91.17647059	85.29411765
82.35294118	63.23529412
83.82352941	86.76470588
91.17647059	88.23529412
88.23529412	86.76470588
88.23529412	82.35294118
88.23529412	89.70588235
85.29411765	85.29411765
91.17647059	85.29411765
91.17647059	85.29411765
Average	
88.08823529	83.82352941

$\alpha = 0.1$, convergence = 0.00001



Note: The excel sheets containing the detailed output for the four test cases shown above are available in the .zip folder which can be used for reference.

Observations

- a) The accuracy achieved with custom implementation is quite close to the accuracy achieved with the sklearn implementation.
- b) The difference between accuracies stems from the reason that sklearn implementation uses regularization which prevents the algorithm from overfitting the training dataset. The regularization is not used in the custom implementation because we have a small dataset that is being split into training and test data and the chances of overfitting are less.
- c) Sklearn implementation uses a number of iterations = 100, as well as convergence (it is called tolerance in sklearn implementation), and the default value for this is 0.0001. As per our understanding, the algorithm runs to find the minimum cost but if it is not achieved in 100 iterations, it considers the cost computed at the 100th execution as the minimum cost.
- d) With different values of learning rate and convergence, we achieved different accuracies but the difference is not huge.
- e) Decreasing the value of convergence results in more time being taken to find the minimum cost.
- f) Accuracy can be used as an analysis tool in our model since the data is not unbalanced.

Conclusion

Logistic regression works well with linearly separable data. It is easier to implement, interpret, and efficient to train. Logistic regression assigns weights to the features of the data which is used to determine the impact of the feature along with the direction of association (positive or negative) and can be used to eliminate features that are less important in the case of a very large dataset. This makes logistic regression a suitable classification algorithm.

Division of Work

Section	Task	Contributor
Code	sigmoid, cost, predict and normalize	Saumya Goel
Code	gradient_descent, logistic_gradient, fit, accuracy	Diksha Srivastava
Report	Introduction to logistic regression (definition, use of logistic regression for classification, Linear vs Logistic Regression, Types of Logistic Regression, which type are we using for our purpose?) Core of logistic regression (Hypothesis)	Diksha Srivastava
Report	Core of logistic regression (Cost and gradient descent), design decisions (Threshold value, type of gradient descent, learning rate, convergence and normalization)	Saumya Goel
Report	Comparison, Observation and Conclusion	Worked together

References

Graph, L. v. (n.d.). *Introduction to Logistic Regression*. Retrieved from towardsdatascience: towardsdatascience.com/introduction-to-logistic-regression-66248243c148

Martin, D. J. (2021, September 21). *Logistic Regression*. Retrieved from Stanford: www.web.stanford.edu/~jurafsky/slp3/5.pdf

Understanding Logistic Regression. (n.d.). Retrieved from Geeksfor geeks: www.geeksforgeeks.org/understanding-logistic-regression/

[Notes Logistic Regression - Stanford ML Andrew Ng \(joparga3.github.io\)](https://joparga3.github.io/Notes-Logistic-Regression-Stanford-ML-Andrew-Ng/)

APPENDIX

```
# Custom Implementation
import numpy as np

"""
Assignment 2
Name: Diksha Srivastava (21235117)
      Saumya Goel (21238562)
Course: Msc Data Analytics (CSD1)
"""

class Logistic:
    def normalize(self, X):
        """
        This method applies Z normalization over the input array X and returns the
        normalized array.
        @author: Saumya Goel
        :param X: input training observations
        :return: normalized training observations
        """
        n = X.shape[1]
        for i in range(n):
            X = (X - X.mean(axis=0))/X.std(axis=0)
        return X

    def sigmoid(self, theta, X):
        """
        This method returns the hypothesis of logistic regression which is the
        probability
        ranging between 0 and 1. This probability can be further used for classification
        task based on threshold value.
        @author: Saumya Goel
        :param theta: regression coefficient
        :param X: input training observations
        :return: hypothesis of logistic regression i.e probabilities ranging between 0
        and 1
        """
        return 1.0 / (1 + np.exp(-np.dot(X, theta.T)))

    def cost(self, theta, X, y):
        """
        This method calculates the cost i.e. evaluate the errors a logistic model is
        going to make.
        @author: Saumya Goel
        :param theta: regression coefficient
        :param X: input training observations
        :param y: labels/target value
        :return: cost value
        """
        y_hat = self.sigmoid(theta, X)
        y = np.squeeze(y)
        prob_y_1 = y * np.log(y_hat)
        prob_y_0 = (1 - y) * np.log(1 - y_hat)
        cost = -prob_y_1 - prob_y_0
        return np.mean(cost)

    def logistic_gradient(self, theta, X, y):
        """
        This method calculates the gradient using formula  $(h(x) - y) \cdot X$ 
        @author: Diksha Srivastava
        :param theta: regression coefficient
        :param X: input training observations
        :param y: labels/target value
        :return: the gradient
        """
```

```

64     first_gradient = self.sigmoid(theta, X) - y.reshape(X.shape[0], -1)
65     final_gradient = np.dot(first_gradient.T, X)
66     return final_gradient
67
68 def gradient_descent(self, X, y, theta, learning_rate, converge):
69     """
70     This method is used to find the optimal coefficient i.e., by minimizing the
71     cost function.
72     @author: Diksha Srivastava
73     :param X: input training observations
74     :param y: labels/target value
75     :param theta: regression coefficient
76     :param learning_rate: magnitude of the amount to move in gradient descent
77     :param converge: the rate through which convergence can be achieved
78     :return: optimized regression coefficient
79     """
80     cost = self.cost(theta, X, y)
81     change_in_cost = 1
82     m = X.shape[0]
83     while change_in_cost > converge:
84         prev_cost = cost
85         derivative = self.logistic_gradient(theta, X, y)/m
86         theta = theta - (learning_rate * derivative)
87         cost = self.cost(theta, X, y)
88         change_in_cost = prev_cost - cost
89     return theta
90
91 def fit(self, X, y):
92     """
93     This method calculates the theta coefficient at 0.01 learning rate and
94     0.01 convergence.
95     @author: Diksha Srivastava
96     :param X: input training observations
97     :param y: labels/target value
98     :return: regression coefficient
99     """
100     # stacking columns with 1's in feature matrix
101     X = np.hstack((np.matrix(np.ones(X.shape[0])).T, X))
102     # initializing theta values
103     theta = np.matrix(np.zeros(X.shape[1]))
104     theta = self.gradient_descent(X, y, theta, 0.01, 0.01)
105     return theta
106
107 def predict(self, theta, X):
108     """
109     This methods makes prediction on the trained model and return the predicted
110     value.
111     @author: Saumya Goel
112     :param theta: regression coefficient
113     :param X: input training observations
114     :return: predictions and predicted probabilities
115     """
116     # stacking columns with 1's in matrix
117     X = np.hstack((np.matrix(np.ones(X.shape[0])).T, X))
118     y_hat = self.sigmoid(theta, X)
119     predictions = np.where(y_hat >= 0.5, 1, 0)
120     return y_hat, np.squeeze(predictions)
121
122 def accuracy(self, y_test, y_pred):
123     """
124     This method counts the number of correct predictions and calculates the
125     accuracy score.
126     @author: Diksha Srivastava
127     :param y_test: actual label
128     :param y_pred: predicted label
129     :return: accuracy
130     """

```

```

129         # counter
130         correctly_classified = 0
131         for count in range(np.size(y_pred)):
132             if y_test[count] == y_pred[count]:
133                 correctly_classified = correctly_classified + 1
134         return (correctly_classified / len(y_test)) * 100
135
136     # Sklearn Implementation
137     from sklearn.preprocessing import StandardScaler
138     from sklearn.linear_model import LogisticRegression
139     from sklearn.metrics import accuracy_score
140
141     """
142     Assignment 2
143     Name: Diksha Srivastava (21235117)
144           Saumya Goel (21238562)
145     Course: Msc Data Analytics (CSD1)
146     """
147     class LogisticSklearn:
148         def normalize(self, X_train, X_test):
149             """
150             This method performs standard scaling and returns the normalized value.
151             @author: Diksha Srivastava
152             :param X_train: input training observations
153             :param X_test: input test observations
154             :return: normalized training and test set
155             """
156             sc = StandardScaler()
157             X_train = sc.fit_transform(X_train)
158             X_test = sc.transform(X_test)
159             return X_train, X_test
160
161         def train_and_predict(self, X_train, X_test, y_train, y_test):
162             """
163             This method trains the classifier with X_train and y_train and predicts on the
164             y_test value.
165             @author: Saumya Goel
166             :param X_train: input training observations
167             :param X_test: input test observations
168             :param y_train: input training labels
169             :param y_test: test labels
170             :return: accuracy and predicted labels/probabilities
171             """
172             # logistic regression
173             classifier = LogisticRegression(tol=0.01)
174             classifier.fit(X_train, y_train.ravel())
175             y_pred = classifier.predict(X_test)
176             probs = classifier.predict_proba(X_test)
177             prob_sk = probs[:, 1]
178             accuracy = accuracy_score(y_test, y_pred) * 100
179             return accuracy, y_pred, prob_sk
180
181     # Calling Class
182     import pandas as pd
183     from sklearn.model_selection import train_test_split
184     from Logistic import Logistic
185     from LogisticSklearn import LogisticSklearn
186     import sklearn.metrics as metrics
187     import matplotlib.pyplot as plt
188     import warnings
189
190     """
191     Assignment 2
192     Name: Diksha Srivastava (21235117)
193           Saumya Goel (21238562)
194     Course: Msc Data Analytics (CSD1)

```



```

195 """
196 class LogisticRegressionTest:
197     warnings.simplefilter(action="ignore", category=UserWarning)
198
199     def read_data(self, file):
200         """
201         This method reads wildfires.csv file and creates independent variable X
202         containing features and dependent variable y containing target value.
203         @author: Diksha Srivastava
204         :param file: input wildfires.csv file
205         :return: X, y which is an array of features and labels
206         """
207         # Data pre-processing.
208         data = pd.read_csv(file)
209         # Converting dependent variable 'fire' in binary form where 'yes' is 1 and 'no'
210         # is 0.
211         data['fire'] = data['fire'].str.strip().map({'yes': 1, 'no': 0})
212         X = data.iloc[:, 1:].values
213         y = data.iloc[:, :1].values
214         return X, y
215
216     def save_predictions(self, writer, y_test, y_pred_sk, y_pred_lr, i):
217         """
218         This method writes the actual and predicted label achieved from sklearn and
219         custom model
220         to Excel sheet for 10 iterations.
221         @author: Diksha Srivastava
222         :param writer: it is the excel writer object of pandas
223         :param y_test: test set of target values
224         :param y_pred_sk: prediction array of sklearn model
225         :param y_pred_lr: prediction array of custom model
226         :param i: the number of iterations
227         """
228         y_test_flat = y_test.ravel()
229         df = pd.DataFrame(columns=['y_test', 'y_pred_sk', 'y_pred_lr'])
230         df['y_test'] = pd.Series(y_test_flat)
231         df['y_pred_sk'] = pd.Series(y_pred_sk)
232         df['y_pred_lr'] = pd.Series(y_pred_lr)
233         df = df.replace({1: 'yes', 0: 'no'})
234         df.to_excel(writer, sheet_name='Exec ' + str(i+1), index=False)
235
236     def plot_roc(self, y_test, preds, color, label):
237         """
238         This method plots the ROC curve for the given model.
239         @author: Saumya Goel
240         :param y_test: test set of target values
241         :param preds: predicted probability of y_test values
242         :param color: provides different color to each model
243         :param label: provides different label to each model
244         """
245         fpr, tpr, threshold = metrics.roc_curve(y_test, preds)
246         roc_auc = metrics.auc(fpr, tpr)
247         plt.title('ROC')
248         plt.plot(fpr, tpr, 'b', label=label+' = %.2f' % roc_auc, color=color)
249         plt.legend(loc='lower right')
250         plt.plot([0, 1], [0, 1], 'r--')
251         plt.xlim([0, 1])
252         plt.ylim([0, 1])
253         plt.ylabel('True Positive Rate')
254         plt.xlabel('False Positive Rate')
255
256     def main(self):
257         """
258         This is the main method which creates the object of sklearn and custom
259         implementation
260         and drives the whole program.
261         @author: Diksha Srivastava

```

```

259 @author: Saumya Goel
260 """
261 global X_test_sk, y_test, y_hat, prob_sk
262 # A dataframe to store the accuracies over 10 iterations and the average
    accuracy as well.
263 df = pd.DataFrame(columns=['accuracy_sk', 'accuracy_lr'])
264 writer = pd.ExcelWriter('Output.xlsx', engine='xlsxwriter')
265 for i in range(0, 10):
266     # Reading and splitting data
267     X, y = self.read_data('wildfires.csv')
268     X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=1 / 3)
269
270     # Training model through sklearn logistic regression
271     sklearn_model = LogisticSklearn()
272     X_train_sk, X_test_sk = sklearn_model.normalize(X_train, X_test)
273     accuracy_sklearn, y_pred_sk, prob_sk = sklearn_model.train_and_predict(
274         X_train_sk, X_test_sk, y_train,
275         y_test)
276
277     # Training model through custom logistic regression
278     lr_model = Logistic()
279     X_train_lr = lr_model.normalize(X_train)
280     X_test_lr = lr_model.normalize(X_test)
281     theta = lr_model.fit(X_train_lr, y_train)
282     y_hat, y_pred_lr = lr_model.predict(theta, X_test_lr)
283     accuracy_lr = lr_model.accuracy(y_test, y_pred_lr)
284     self.save_predictions(writer, y_test, y_pred_sk, y_pred_lr, i)
285
286     # Storing accuracy in dataframe
287     df.loc[i] = [accuracy_sklearn, accuracy_lr]
288
289 print(df)
290 print('Average Accuracy of sklearn model: ', df['accuracy_sk'].mean())
291 print('Average Accuracy of custom model: ', df['accuracy_lr'].mean())
292
293 # Writing the average accuracy to excel.
294 df.to_excel(writer, sheet_name='Accuracy', index=False)
295 writer.sheets['Accuracy'].activate()
296 writer.sheets['Accuracy'].write(11, 0, 'Average')
297 writer.sheets['Accuracy'].write(12, 0, df['accuracy_sk'].mean())
298 writer.sheets['Accuracy'].write(12, 1, df['accuracy_lr'].mean())
299 writer.save()
300
301 # Plotting ROC curve
302 self.plot_roc(y_test, prob_sk, 'darkorange', 'AUC for Scikit implementation')
303 self.plot_roc(y_test, y_hat, 'darkblue', 'AUC for custom implementation')
304 plt.show()
305
306 if __name__ == "__main__":
307     main_obj = LogisticRegressionTest()
308     main_obj.main()
309
310
311
312
313

```