
Analyzing Air Quality and Liver Disease Using Binary Logistic Regression

Max Ardito

maxwell.ardito@mail.mcgill.ca

Edwin Meriaux

edwin.meriaux@mail.mcgill.ca

Ohood Sabr

ohood.sabr@mail.mcgill.ca

Abstract

This work focuses on using a binary logistic regression model to try and solve two different problems related to health and medicine. The first problem attempts to classify polluted air samples provided a dataset containing information about airborne chemicals, and the second problem attempts to classify instances of liver disease in patients provided similar biochemical data. A common gradient descent approach was used to train a binary classifier, and k-fold cross validation was used to choose the final model. After implementing the base regression model, various techniques were used in an attempt to improve accuracy, including brute-force feature reduction, regularization, feature expansion, and principal component analysis (PCA). It was observed that most of these techniques made very some difference in improving the mean accuracy of the training and validation set. However, these differences were very subtle, even when using optimal hyperparameters. We concluded that the largest qualitative improvement in the model was found in using a combination of feature reduction, feature expansion, and PCA.

1 Introduction

There are a number of approaches for forecasting data of various forms in the field of machine learning. Classification is one of these approaches, which includes many types of algorithms such as Naive Bayes, Decision Trees, and Logistic Regression. The type of data used, the dimensions of the dataset, and the aptitude of the dataset's features have a significant impact on the model or classifier that is chosen [8, 10].

This work focuses on using binary logistic regression to try and predict the likelihood of a patient testing positive for liver disease and the likelihood of an air sample being polluted, given a number of respective chemical features. Logistic regression is a supervised learning algorithm that is used to solve binary classification problems. It calculates the likelihood of an event occurring based on a collection of variables that are assumed to be independent from one another, thus drawing decision boundaries between data points in N-dimensional feature space in an attempt to classify future data [3]. This type of classification is widely applied in studies about the health sciences because it is especially suitable for models involving illness condition (diseased or healthy) and decision making (yes or no) [4].

In order to build a logistic regression model from scratch, we first must define an error function such as the negative log likelihood or "cross-entropy" function (1) that we attempt to minimize by taking its derivative and using gradient descent (2) to slowly and iteratively modify the weights and the bias of the linear model. [7]

$$Err(w) = - \sum_{i=1}^N y_i \ln(\mu_i) + (1 - y_i) \ln(1 - \mu_i) \quad (1)$$

$$\Delta = \frac{\partial Err(w_k)}{\partial(w_k)} = - \sum_{i=1}^N x_i(y_i - \mu_i) \quad (2)$$

In (1) and (2), μ is equal to the sigmoid function $\sigma(w^T x_i) = P(y_i|x_i, w) = \frac{1}{1+e^{-w^T x}}$ which is used to transform the continuous output into a discrete binary measure (0 or 1) during classification, depending on whether its output is greater than or less than 0.5.[\[11\]](#) The weights in this model are updated following the rule $w_{k+1} = w_k + \alpha_k \Delta$, where α_k is a hyper-parameter representing the step-size [\[2\]](#).

The overall success of logistic regression depends heavily on finding ideal hyper-parameters, variations on the given dataset, and expansions on the model. In Section 2, we describe the Datasets and give an overview of the features analysis. In Section 3, the various results are presented which include a model selection and feature selection and expansion. In Section 4, we the discuss our findings in conclusion and propose some future work. Section 5 gives a statement of contribution, and finally an Appendix is given containing the complete complementary code, plots, and results.

2 Datasets

2.1 Overview

In this study, the aforementioned logistic regression algorithm is employed to carry out classification on two different datasets. The first dataset represents information on air quality to estimate the probability of air pollution. The second dataset contains similar information to determine the probability of a patient having liver disease. In order to classify future data with known feature information and unknown class information, we have built and tested a binary logistic regression model from scratch by using the Python programming language and the Google Collab environment.

The original air quality dataset contains 1599 samples and 11 features each, with a class value of 1 representing a sample of polluted air and a class value of 0 representing normal air. The dataset for liver disease contains 330 samples of 8 features each, with class 1 assigned to a positive diagnosis for liver disease and class 0 assigned to a negative diagnosis. Furthermore, all of the features in each dataset are continuous.

2.2 Feature Analysis

In Figure 1, we can see that most of the features in the liver disease dataset form Gaussian distributions. This is the same case for the air quality dataset (see Appendix A for air quality histograms). In a hypothetical feature selection stage in which features deemed unnecessary or noisy are removed, these distributions can serve as qualitative tools to both improve and justify the accuracy of our model. For instance, the variance between the two binary classes of each respective feature can be calculated, but also judged qualitatively by looking at the histogram plots. Features that hold the least amount of variance—in other words, features with the most overlap in their distributions—might serve as unnecessary in the training process due to their indistinguishable nature between class 0 and class 1. These methods will be discussed further in Section 3.2, as well as methods for improving the model by expanding the feature space.

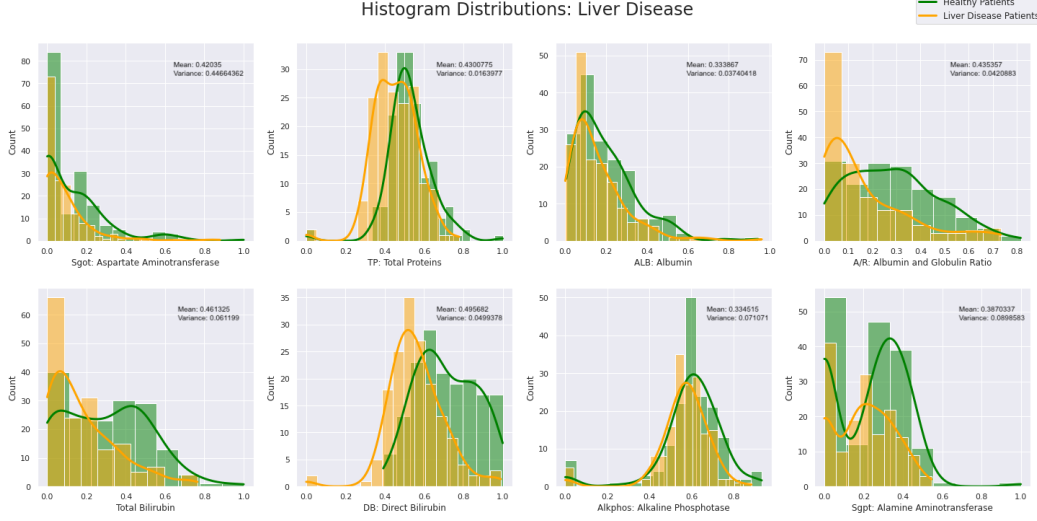


Figure 1: Histogram plot of each liver disease feature comparing data points from people who tested positive (yellow) and tested negative (green). See Appendix A for a similar histogram of the air quality dataset.

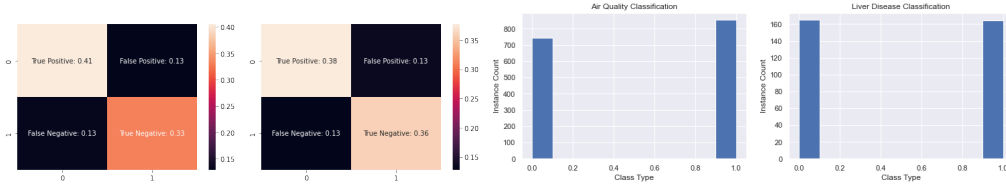


Figure 2: (Left) Confusion matrices for base logistic regression model \mathcal{M} run on air quality dataset and liver disease dataset. (Right) Histogram representing the total count of data instances of each class.

3 Results

3.1 Model Selection

In experimenting with different variations of the logistic regression model, we started with a base model $\mathcal{M} = \{X, y, \alpha, k\}$ where X and y are the training set's feature matrix and class vector respectively, α represents the step size, and k represents the number of steps taken during gradient descent.

The base model \mathcal{M} trains the binary logistic regression model on each dataset given strictly as-is with no additional feature reduction or feature expansion methods. However, one row-reduction step was implemented prior to training that removes identical rows. The logic behind this is to remove biases in the data. The odds that two points are identical over 8-11 features are low. In this case there were only a few duplicates over each individual dataset, which means duplicates don't currently affect the model too much, but very well could affect performance if more data was added.

It was observed that the performance of each model discussed in this section demonstrated similar trends in accuracy between both the liver disease dataset and the air quality dataset. But a few variations of the base model displayed some key differences in the performance of each dataset. These differences will be discussed further while each variation of the base model is discussed. A table containing the differences in performance between the two datasets is provided in Appendix A, and can be consulted throughout this section.

The base model \mathcal{M} was found to yield the highest accuracy when setting $\alpha = 0.001$ and $k = 100000$. When trained on a 10-fold cross validation, the model yields a mean accuracy score of 73.41% on the training data and 71.88% on the validation set data for liver disease, as well as a mean accuracy of 74.42% on the air quality training set and 73.63% on the air quality validation set. In the event that these two datasets were not initially normalized, we added a normalization algorithm prior to the model training. This was found to increase the accuracy of the model by anywhere from 1-3% on the air quality dataset, but made little to no improvement on the liver disease dataset.

For both datasets, the fluctuation in accuracy on each fold of the training set was found to range anywhere from 72% to 76%, while the validation set for each fold was observed to range anywhere between 59.38% and 81.25%, even despite the model’s mean accuracy on the training and validation sets being nearly equal. It is possible to conclude from this information that model \mathcal{M} might produce an unpredictably high variance, given a larger version of the dataset. An L1 regularization was thus implemented on the model to penalize any overfitting that might be occurring.

$$\hat{w} = \operatorname{argmin}_w \left\{ \sum_{i=1}^N y_i - (w^T x_i)^2 + \Lambda \right\} \begin{cases} \Lambda \triangleq \lambda \sum_{j=0}^M |w_j| & \text{for L1} \\ \Lambda \triangleq \lambda \sum_{j=0}^M w_j^2 & \text{for L2} \end{cases} \quad (3)$$

By implementing the L1 cost function in (3), we derived a model $\hat{\mathcal{M}} = \{X, y, \alpha = 0.001, k = 100000, \lambda = 0.001\}$ containing the added λ hyper-parameter. We concluded that L1 regularization tends to improve this fluctuation in the validation set’s results of both datasets, but makes little to no difference in the mean accuracy of each model. Additionally an L2 regularization was also implemented but made no quantitative difference to the any of the base model’s accuracy metrics.

3.2 Feature Selection and Expansion

Feature selection techniques proved to be more successful in decreasing the base model’s variance. Our approach combined implementing a brute force technique to calculate which features to remove and qualitatively analyzing the aforementioned histograms.

$$\mathcal{P}(y) = 2^N, |y| = N \quad (4)$$

In order to figure out how many features to remove from the model, we first derived the power set $\mathcal{P}(y)$ containing all N subsets of our feature vector y and calculated the mean accuracy for each subset $\{s\}^y$. [6] We then compared each of these accuracy measures for each subset using cross-validation and calculated which one had the least error (See Appendix A - Figure 4).

When manually removing features our findings suggested that the accuracy of the model decreases significantly after removing more than one feature, but increases slightly when removing a single feature just like from the air quality dataset. As discussed in Section 2.2, analyzing the histograms in Figure 1 led us to believe that removing a feature whose binary class distributions overlapped the most might be advantageous. These can be quantitatively decided by using the variance values captured in Figure 1, and were decidedly TP: Total Proteins in the liver disease dataset and Temperature in the air quality dataset. Indeed, dropping these two features and retraining resulted in a consistent 1% increase in both the validation set accuracy and the total accuracy of the model trained on each dataset. But removing multiple features can increase the accuracy when all combos are checked. The top non brute force accuracy in the air quality dataset is with PCA and Expansion with a 75.07 percent accuracy but this is increased to 76.7 percent with brute force reduction where features 4, 6, 8, and 9 were removed (from the original dataset). A similar effect is seen in the liver disease dataset where the top accuracy goes from 74.6 percent to 74.74 percent. The only features removed were

1, 4, 5, and 8 (from the original dataset). This method of course is very time intensive so its utility decrease as the dataset grows more and more at an exponential rate.

After focusing on variance, we decided to try and construct a more complex feature space, as this is a common methodology for lowering a model's bias [8]. Starting from the feature-reduced model from above, a feature expansion on the model was implemented, thus adding powers 2, 3, 1/2, 1/3, and the natural log into the data's feature space. Of course, there are also other possible values that the data can be expanded to like exponential but for the sake of this experiment, this proved the point of feature expansion. By itself feature expansion reduced validation accuracy, but in combination with a Principal Component Analysis (PCA) it helped increase accuracy by about 1-1.5%.

Deriving a dataset from performing Principal Component Analysis (PCA) reduces the dimensionality of the feature space into the first N principal components, privileging new columns that describe the most variance in the data over lower-ranking components that contain unnecessary information.

$$\operatorname{argmin}_{P,U} \left\{ \sum_{i=1}^N \|x_i - UP^T x_i\|^2 \right\} \quad (5)$$

Using eigenvalue decomposition, a matrix of eigenvectors P is derived and sorted, representing N vectors ordered by the amount of variance captured in the original data. This decomposition method is optimized by (5) to ensure that the minimum amount of information is lost from compression of the original data x_i to decompression of the principal components P^T through multiplying by U [1].

We cross validated the model for training sets that used the first 1 - 10 principal components and compared the accuracy with one another. Our findings showed that by simply using the first 4 principal components, our liver disease model yielded both the highest accuracy and lowest variance yet scoring 74.44% on the training set and 74.06% on the validation set. Similarly, the air quality dataset scored an accuracy of 75.07% trained on the first 9 principal components. These models also used the feature reduction, expansion, and normalization steps mentioned earlier.

4 Discussion and Conclusion

We have seen that using a number of variations on the binary logistic regression model have proven to marginally increase the accuracy of classifying unknown air quality and medical patient data. The extension on the base model that yielded the highest accuracy was a model that used a number of techniques prior to the training process including data normalization, feature reduction, feature expansion, and PCA. These techniques provided an improved dataset for which we trained the logistic regression model.

Possible future work could include investigating augmented gradient descent techniques such as Stochastic Gradient Descent (SGD) and the Stochastic Average Gradient (SAG) [9]. These techniques are used heavily in the SciKit Learn implementation of the binary logistic regression model [5]. An additional step for improving the model could be to collect more data, which might improve the accuracy of our model significantly.

5 Statement of Contribution

1. Max Ardito: Logistic regression, K-Fold, evaluation methods, writing report, principal component analysis, plots, regularization
2. Ohood Sabr: Logistic regression, K-Fold, writing report
3. Edwin Meriaux: Logistic regression, principal component analysis, brute force feature selection, feature expansion, average feature decrease

References

- [1] Narges Armanfard. Ecse 551 - machine learning for engineers lecture 13 — dimension reduction and knn, October 2022.
- [2] Narges Armanfard. Ecse 551 - machine learning for engineers lecture 5—linear regression, overfitting, logistic regression, September 2022.
- [3] Hoss Belyadi and Alireza Haghighat. Chapter 5 - supervised learning. In Hoss Belyadi and Alireza Haghighat, editors, *Machine Learning Guide for Oil and Gas Using Python*, pages 169–295. Gulf Professional Publishing, 2021.
- [4] Ernest Yeboah Boateng and Daniel A Abaye. A review of the logistic regression model with emphasis on medical research. *Journal of data analysis and information processing*, 7(4):190–207, 2019.
- [5] https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html. `sklearn.linear_model.logisticregression`, Oct 2022.
- [6] <https://www.geeksforgeeks.org/power-set/>. Power set, Aug 2022.
- [7] Kevin P. Murphy. *Machine learning: A probabilistic perspective*, page 247–249. MIT Press, 2021.
- [8] omdena.com/blog/machine-learning-classification-algorithms. Best machine learning classification algorithms + real-world projects, 2022.
- [9] Mark Schmidt, Nicolas Le Roux, and Francis Bach. Minimizing finite sums with the stochastic average gradient. *Mathematical Programming*, 162(1):83–112, 2017.
- [10] Pratap Chandra Sen, Mahimarnab Hajra, and Mitadru Ghosh. Supervised classification algorithms in machine learning: A survey and review. In *Emerging technology in modelling and graphics*, pages 99–111. Springer, 2020.
- [11] Xiaonan Zou, Yong Hu, Zhewen Tian, and Kaiyuan Shen. Logistic regression model optimization and case analysis. In *2019 IEEE 7th International Conference on Computer Science and Network Technology (ICCSNT)*, pages 135–139, 2019.

5.1 Appendix A. - Plots

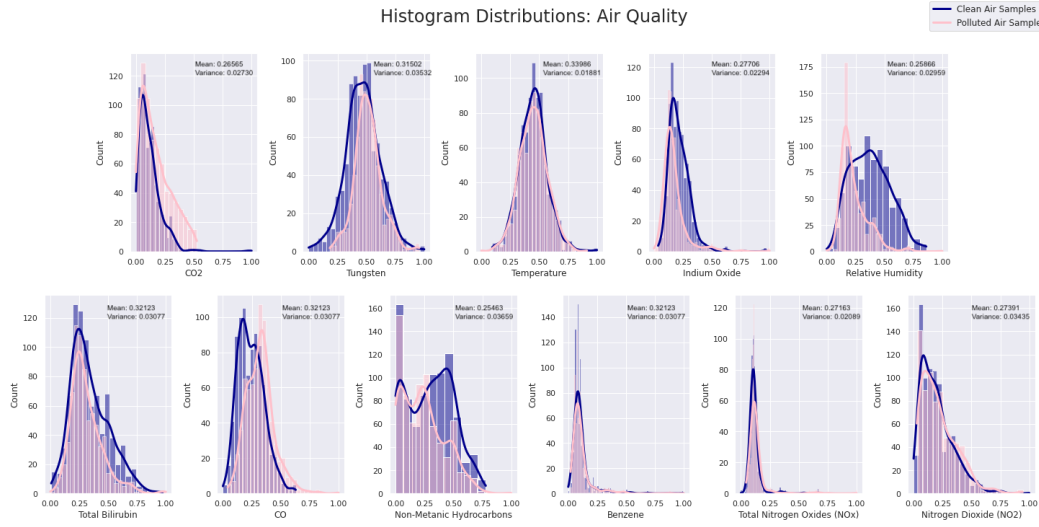


Figure 3: Histogram plot of each air quality feature comparing data points from clean air samples (purple) and polluted air samples (purple).

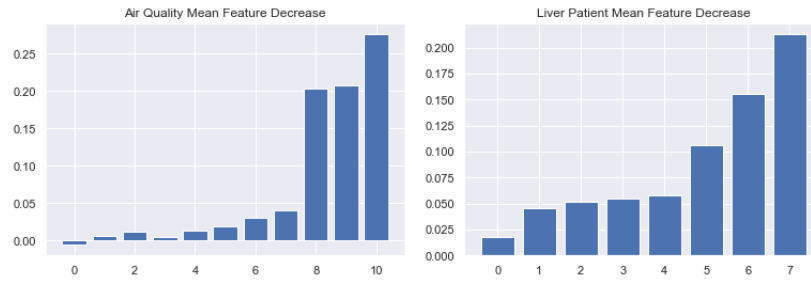


Figure 4: Brute force tests for the mean error of the model when dropping the number of features displayed on the X axis.

5.2 Appendix B. - Accuracy Table

Table 1: Mean accuracy metrics between air quality data and liver disease data for various logistic regression models trained using 10-fold cross-validation

Model	Air Quality	Liver Disease
Base Logistic Regression	72.3%	71.56%
Normalized	73.78%	70.95%
Regularized (L1) + Normalized	74.07%	70.93%
Regularized (L2) + Normalized	73.4%	71.56%
Feature Reduction + Normalized	74.5%	72.5%
Feature Reduction + Feature Expansion + Normalized	74.9%	71.56%
Feature Reduction + Feature Expansion + PCA + Normalized	75.07% [PC = 9]	74.68% [PC = 4]

5.3 Appendix C. - Code

```
1  # -*- coding: utf-8 -*-
2  # -*- coding: utf-8 -*-
3  """Logistic_Regression_10-10.ipynb
4
5  Automatically generated by Colaboratory.
6
7  Original file is located at
8      https://colab.research.google.com/drive/1vFsscCXpbgXur3w80n\_tpæEa\_feU0pYj
9
10 # Logistic Regression
11
12 # Regression Class
13
14 Class containing tools for a binary logistic regression
15 """
16
17 import numpy as np
18 from numpy import log as ln
19 from enum import Enum
20 import math
21
22 """
23 Class for implementing a binary logistic regression. Includes additional
24 functionality for penalizing weights using L1 and L2 regularization.
25
26 Code citations: https://github.com/python-engineer/MLfromscratch
27 """
28 class LogisticRegression:
29
30     def __sigmoid(self, a: float):
31         return (1 / (1 + np.exp(-a)))
32
33     class Regularization(Enum):
34         NONE = 0
35         L1 = 1
36         L2 = 2
37
38     def __init__(self, step_size: float, iterations: int):
39         self.step_size = step_size
40         self.iterations = iterations
41
42     def fit(self,
43            X: np.ndarray,
44            Y: np.ndarray,
45            regularization: Regularization = Regularization.NONE,
46            lmda: float = 0):
47         # Initialize random weights and bias
48         self.w = np.random.rand(X.shape[1])
49         self.bias = np.random.rand()
50         y_hat = np.ndarray(Y.shape)
51
52         # Perform gradient descent
53         for _ in range(self.iterations):
54             wTx = np.dot(X, self.w) + self.bias
55             y_hat = self.__sigmoid(wTx)
56
57         # Implement a Lasso or Ridge regularization if specified
58         if (regularization == self.Regularization.L1):
59             dW = (1 / X.shape[0]) * np.add(np.dot(X.transpose(),
60                                                    (y_hat - Y.flatten())),
```



```

61         elif (regularization == self.Regularization.L2):
62             dW = (1 / X.shape[0]) * np.add(np.dot(X.transpose(),
63                 (y_hat - Y.flatten()))
64                 ↪ np.dot(self.w, np.sign(self.w)
65                 ↪ * lmda))
66
67         else:
68             dW = (1 / X.shape[0]) * np.dot(X.transpose(), (y_hat -
69                 ↪ Y.flatten()))
70
71         dB = (1 / X.shape[0]) * np.sum(y_hat - Y.flatten())
72
73         self.w -= self.step_size * dW
74         self.bias -= self.step_size * dB
75
76     def predict(self, X_new):
77         #prediction of the LR model
78         wTx = np.dot(X_new, self.w) + self.bias
79         y_hat = self.__sigmoid(wTx)
80
81         for i in range(y_hat.shape[0]):
82             if y_hat[i] > 0.5:
83                 y_hat[i] = 1
84             else:
85                 y_hat[i] = 0
86
87         return y_hat
88
89     """# Evaluation
90
91     Class containing various evaluation tools
92     """
93
94     import seaborn as sns
95     import numpy as np
96     import matplotlib.pyplot as plt
97     import pandas as pd
98
99     Class for evaluating the performance of a logistic regression model.
100     Includes tools for calculating confusion matrix, accuracy, precision,
101     recall, specificity, and false positive rate.
102
103     Assumes that the input data is the result of a binary logistic regression
104     (e.g. y && y_hat = {0, 1})
105     """
106
107     class LogisticEvaluation:
108         """
109         Initializes the evaluation from a vector of predicted binary values (y)
110         and a vector of actual values (y_hat). Stores these values in a
111         confusion matrix variable (cm) as well as individual cell
112         values (tp, tn, fp, fn)
113         """
114
115         def __init__(self, y_hat: np.ndarray, y: np.ndarray):
116             size = y_hat.size
117             self.cm = np.array([[0, 0], [0, 0]])
118             for i in range(size):
119                 if (y_hat[i] == 1) and (y[i] == 1):
120                     self.cm += np.array([[1, 0], [0, 0]])
121                 elif (y_hat[i] == 0) and (y[i] == 0):
122                     self.cm += np.array([[0, 0], [0, 1]])
123                 elif (y_hat[i] == 1) and (y[i] == 0):
124                     self.cm += np.array([[0, 1], [0, 0]])

```

```

122         else:
123             self.cm += np.array([[0, 0], [1, 0]])
124         # Normalize confusion matrix
125         self.cm = np.divide(self.cm, np.array(size))
126
127         # Store individual table values
128         self.tp = self.cm[[0], [0]].item()
129         self.tn = self.cm[[1], [1]].item()
130         self.fp = self.cm[[1], [0]].item()
131         self.fn = self.cm[[0], [1]].item()
132
133     """
134     Prints a heatmap of the confusion matrix
135     """
136
137     def confusion_matrix(self):
138         df = pd.DataFrame(self.cm)
139         fig = plt.figure()
140         cell_labels = np.array(
141             [[
142                 "True Positive: " + str(round(self.tp, 2)),
143                 "False Positive: " + str(round(self.fp, 2))
144             ],
145             [
146                 "False Negative: " + str(round(self.fn, 2)),
147                 "True Negative: " + str(round(self.tn, 2))
148             ]])
149         print(round(self.tp, 2), round(self.fp, 2), round(self.fn,
150             ↪ 2), round(self.tn, 2))
151
152         sns.heatmap(df, fmt='', annot=cell_labels)
153         plt.savefig("Heatmap.png")
154         return self.cm
155
156     """
157     Returns the model's accuracy
158     """
159
160     def accu_eval(self):
161         accuracy = (self.tp + self.tn) / (self.cm.sum())
162         return accuracy
163
164     """
165     Returns the model's precision
166     """
167
168     def precision_eval(self):
169         precision = self.tp / (self.tp + self.fp)
170         return precision
171
172     """
173     Returns the model's recall
174     """
175
176     def recall_eval(self):
177         recall = self.tp / (self.tp + self.fn)
178         return recall
179
180     """
181     Returns the model's specificity
182     """
183
184     def spec_eval(self):
185         specificity = self.tn / (self.fp + self.tn)

```

```

185         return specificity
186
187     """
188     Returns the model's false positive rate
189     """
190
191     def fpr_eval(self):
192         fp_rate = self.fp / (self.fp + self.tn)
193         return fp_rate
194
195     """# K-Fold Cross Validation Class
196
197     Class that uses the logistic regression class in the context of a K-fold cross
198     ↪ validation. Optional flags can be set for normalizing, expanding, and/or
199     ↪ using the first N principal components to train instead of the features
200     ↪ themselves
201     """
202
203     import pandas as pd, numpy as np
204     from sklearn.preprocessing import StandardScaler
205     """
206     Class for performing K-fold cross validation and
207     returning its mean error. Ideally used for comparing
208     the performance of multiple logistic regression models.
209     """
210
211     class KFold:
212         """
213         Initializes class by shuffling the input data and obtaining
214         a validation set size based on the input dimensions and specified
215         K value.
216         """
217
218         def __init__(self,
219                     X: np.ndarray,
220                     k: int,
221                     step_size: float,
222                     iterations: int,
223                     Normalization_toggle: int,
224                     Regularization_toggle: int,
225                     PCA_toggle: int,
226                     Expansion_toggle: int,
227                     pca_size: int,
228                     broute_force_toggle: int):
229             # Remove dupliacate data points
230             self.X = np.unique(X, axis=0)
231             self.k = k
232             self.step_size = step_size
233             self.iterations = iterations
234             self.PCA_toggle = PCA_toggle
235             self.Normalization_toggle = Normalization_toggle
236             self.Regularization_toggle = Regularization_toggle
237             self.Expansion_toggle = Expansion_toggle
238             self.pca_size = pca_size
239             self.brout_force_toggle = broute_force_toggle
240             self.best_accuracy = 0
241             self.best_combo = []
242
243             # Validation size set = number of rows / k
244             self.validation_set_size = int(self.X.shape[0] / k)
245
246             # Shuffle the rows of the input data

```

```

246         np.random.shuffle(self.X)
247
248     """
249     Performs the cross validation on K iterations of the input data.
250     The cross validation is performed by taking the first validation
251     set from the top of the input data and then subsequently shifting
252     (rolling) the input data N elements, where N = validation set size.
253
254     Note that this method of cycling and partitioning will automatically
255     throw any remainder of input data into the validation set if the data
256     set cannot be evenly divided into K sets.
257     """
258     #calculates variance and expectation based on np libraries
259     def variance_expectation(self,Xdata):
260         print(Xdata.shape[1])
261         for i in range(Xdata.shape[1]):
262             print("mean: ",np.mean(Xdata[i]))
263             print("variance: ",np.var(Xdata[i]))
264
265     #power set calculation for the brute for feature reduction
266     #https://www.geeksforgeeks.org/power-set/ source used to help with the bit
    ↪ shift if statement
267     def powerset(self,fullset):
268         listsub = list(fullset)
269         subsets = []
270         for i in range(2**len(listsub)):
271             subset = []
272             for k in range(len(listsub)):
273                 if i & 1<=k:
274                     subset.append(int(listsub[k]))
275             subsets.append(subset)
276         return subsets
277
278     #normalization calculation based on expectation (mu) and standard
    ↪ deviation (std)
279     def normalization(self, Xtrain):
280         normalized = Xtrain
281         #operation on each feature in dataset
282         #in theory this can be accelerated if a custom np function was used
283         #would remove the need for the for loop
284         #the function would effectively be everything in the current for loop
285         for i in range(Xtrain.shape[1]):
286             #mu calculation
287             mu = np.sum(Xtrain[:,i])/Xtrain.shape[0]
288             #std calculation
289             std_subtraction = np.subtract(Xtrain[:,i],mu)
290             std_square = np.square(std_subtraction)
291             std_sum = np.sum (std_square)/Xtrain.shape[0]
292             std = std_sum**(1/2)
293             #normalization calculation
294             normalized[:,i] = np.subtract(Xtrain[:,i],mu)/std
295
296         return normalized
297
298     # Principal Component Analysis function for dimensionality reduction
299     # Code citation:
    ↪ https://www.askpython.com/python/examples/principal-component-analysis
300     def PCA(self, X, num_components):
301         # Center the data set by subtracting the mean
302         X_mean = X - np.mean(X , axis = 0)
303
304         # Get the covariance matrix
305         covariance_matrix = np.cov(X_mean , rowvar = False)
306

```

```

307     # Calculate the eigenvalues and eigenvectors of the covariance matrix
308     eigen_values , eigen_vectors = np.linalg.eigh(covariance_matrix)
309
310     # Sort the eigenvalues and eigenvectors in descending order
311     sorted_index = np.argsort(eigen_values)[::-1]
312     sorted_eigenvalue = eigen_values[sorted_index]
313     sorted_eigenvectors = eigen_vectors[:,sorted_index]
314
315     # Take a subset of the sorted vectors to represent the first N principal
316     ↪ components
317     eigenvector_subset = sorted_eigenvectors[:,0:num_components]
318
319     # Reduce dimensionality by taking the dot product of the transposed
320     # principal components mean of the original data set
321     X_reduced = np.dot(eigenvector_subset.transpose(),
322     ↪ X_mean.transpose()).transpose()
323
324     return X_reduced
325
326 #power calculation for feature expansion for which ever power needed ~2
327 ↪ ~3...
328 def to_power(self, data, power):
329     return data ** power
330
331 #power log calculation for the features expansion
332 def natural_log(self, data):
333     math.log1p(data)
334     return math.log1p(data)
335
336 #mean feature decrease calculation to find the importance of each feature
337 def feature_importance(self,X_train,Y_train):
338
339     iterations = 10000
340     learning_rate = 0.0015
341
342     #setting up the baseline for the test results assuming all features
343     ↪ tested on
344     W, B, true_cost_list, true_accuracy, true_confusion = model(X_train,
345     ↪ Y_train, learning_rate = learning_rate, iterations = iterations)
346     feature_number_list = [i for i in range(X_train.shape[1])]
347
348     #data storage for feature drop comparison
349     importance_accuracy = []
350     importance_confusion = []
351     importance_cost = []
352     for i in range(X_train.shape[1]):
353         #runs test on each set of features. Set of features being n-1
354         ↪ features of n combinations
355         tmp_list = feature_number_list
356         tmp_list.remove(i)
357         tmp_array = X_train[:,tmp_list]
358         #current Logistic Regression test with modified dataset (see line
359         ↪ above)
360         W, B, cost_list, accuracy, confusion = model(tmp_array, Y_train,
361         ↪ learning_rate = learning_rate, iterations = iterations)
362
363         #data storage for future plots
364         importance_accuracy.append(true_accuracy - accuracy)
365         importance_confusion.append(confusion)
366         importance_cost.append(cost_list)
367         print(true_accuracy - accuracy)
368     print(importance_accuracy)
369     return importance_accuracy #data return

```

```

363     #plotting mean feature decrease data
364     def plot_results(self,data,val_list):
365         plt.title("Air Quality Data Importance") #line changed depending on the
366         ↪ data being run for feature decrease
367         plt.bar(val_list,data)
368         plt.show()
369
370     #feature expansion for the data set to certain powers
371     # powers: 2,3,1/2,1/3 and natural log
372     def Feature_Increase(self, xdata):
373         new_xdata = xdata
374         for i in range(xdata.shape[1]):
375             #higher exponential powers
376             for j in range(2,4):
377                 tmp = xdata[:,i]
378                 tmp_vector = np.vectorize(self.to_power)
379                 tmp = np.asmatrix(tmp_vector(tmp,j))
380                 new_xdata = np.hstack((new_xdata, tmp.T))
381             #1/n power of exponential powers
382             for j in range(2,4):
383                 tmp = xdata[:,i]
384                 tmp_vector = np.vectorize(self.to_power)
385                 tmp = np.asmatrix(tmp_vector(tmp,1/j))
386                 new_xdata = np.hstack((new_xdata, tmp.T))
387             #natural log power
388             tmp = xdata[:,i]
389             tmp_vector = np.vectorize(self.natural_log)
390             tmp = np.asarray(tmp_vector(tmp))
391             tmp = np.asmatrix(tmp)
392
393             new_xdata = np.hstack((new_xdata, tmp.T))
394         print(new_xdata)
395         print(new_xdata.shape)
396         return new_xdata.real
397
398     def cross_validation(self):
399         # Initialize error count
400         self.train_err = 0
401         self.test_err = 0
402
403         # Separate the data's classes (Y) from the features (X)
404         # and use expansion, normalization, and/or PCA if flags are set
405         self.Y = np.array([self.X[:, -1]]).transpose()
406         self.X = self.X[:, :-1]
407         if self.Expansion_toggle == 1:
408             self.X = self.Feature_Increase(self.X)
409             self.X = np.asarray(self.X)
410         if self.Normalization_toggle == 1:
411             self.X = self.normalization(self.X)
412         if self.PCA_toggle == 1:
413             self.X = self.PCA(self.X,self.pca_size)
414         if self.brout_force_toggle == 1:
415             print(self.X.shape[1])
416             print(list(range(self.X.shape[1])))
417             test_set = self.powerset(list(range(self.X.shape[1])))
418
419             for combo in test_set:
420                 print("-----")
421                 print("current combo: ",combo)
422                 print("best combo:", self.best_combo)
423                 print("best acc:", self.best_accuracy)
424                 print("-----")
425                 for i in range(self.k):
426                     tmpX = self.X[:,combo]

```

```

426         # Training set = all rows where index is larger than
427         ↪ validation set size
428     X_train = tmpX[self.validation_set_size:, :]
429     Y_train = self.Y[self.validation_set_size:, :]
430
431     # Test set = all rows where index is below validation set size
432     X_test = tmpX[:self.validation_set_size, :]
433     Y_test = self.Y[:self.validation_set_size, :]
434
435     # Initialize logistic regression class for training data
436     logistic = LogisticRegression(self.step_size, self.iterations)
437
438     # Fit the model (get Ws)
439     if (self.Regularization_toggle == 2):
440         logistic.fit(X_train, Y_train, logistic.Regularization.L2,
441             ↪ 0.001)
442     elif (self.Regularization_toggle == 1):
443         logistic.fit(X_train, Y_train, logistic.Regularization.L2,
444             ↪ 0.001)
445     else:
446         logistic.fit(X_train, Y_train)
447
448     # Predict Y output on X training data to measure accuracy
449     Y_hat_train = logistic.predict(X_train)
450
451     # Predict Y output on X test data
452     Y_hat_test = logistic.predict(X_test)
453
454     # Evaluate the predicted Y with the actual Y from the training
455     ↪ data
456     self.train_eval = LogisticEvaluation(
457         Y_hat_train, Y_train)
458
459     # Evaluate the predicted Y with the actual Y from the test data
460     self.test_eval = LogisticEvaluation(Y_hat_test, Y_test)
461
462     # Accumulate the error
463     self.train_err += self.train_eval.accu_eval()
464     self.test_err += self.test_eval.accu_eval()
465
466     print("Fold - ", i + 1, " / ", self.k)
467
468     acc_train = self.train_eval.accu_eval() * 100
469     acc_test = self.test_eval.accu_eval() * 100
470
471     print("Accuracy (Training Data): ", float(f'{acc_train:.2f}'),
472         ↪ "%")
473     print("Accuracy (Validation Data): ", float(f'{acc_test:.2f}'),
474         ↪ "%")
475     if acc_train > self.best_accuracy:
476         self.best_accuracy = acc_train
477         self.best_combo = combo
478
479     # Shift the X data over a validation set size to ensure a new
480     # validation set data for the next training iteration
481     tmpX = np.roll(tmpX, -self.validation_set_size, axis=0)
482     self.Y = np.roll(self.Y, -self.validation_set_size, axis=0)
483
484 if self.brout_force_toggle == 0:
485
486     # Train on training set, test on validation set
487     for i in range(self.k):

```

```

484         # Training set = all rows where index is larger than validation
         ↪ set size
485     X_train = self.X[self.validation_set_size:, :]
486     Y_train = self.Y[self.validation_set_size:, :]
487
488     # Test set = all rows where index is below validation set size
489     X_test = self.X[:self.validation_set_size, :]
490     Y_test = self.Y[:self.validation_set_size, :]
491
492     # Initialize logistic regression class for training data
493     logistic = LogisticRegression(self.step_size, self.iterations)
494
495     # Fit the model (get Ws)
496     if(self.Regularization_toggle):
497         logistic.fit(X_train, Y_train, logistic.Regularization.L2,
         ↪ 0.001)
498     else:
499         logistic.fit(X_train, Y_train)
500
501     # Predict Y output on X training data to measure accuracy
502     Y_hat_train = logistic.predict(X_train)
503
504     # Predict Y output on X test data
505     Y_hat_test = logistic.predict(X_test)
506
507     # Evaluate the predicted Y with the actual Y from the training
     ↪ data
508     self.train_eval = LogisticEvaluation(
509         Y_hat_train, Y_train)
510
511     # Evaluate the predicted Y with the actual Y from the test data
512     self.test_eval = LogisticEvaluation(Y_hat_test, Y_test)
513
514     # Accumulate the error
515     self.train_err += self.train_eval.accu_eval()
516     self.test_err += self.test_eval.accu_eval()
517
518     print("Fold - ", i + 1, " / ", self.k)
519
520     acc_train = self.train_eval.accu_eval() * 100
521     acc_test = self.test_eval.accu_eval() * 100
522
523     print("Accuracy (Training Data): ", float(f'{acc_train:.2f}'),
     ↪ "%")
524     print("Accuracy (Validation Data): ", float(f'{acc_test:.2f}'),
     ↪ "%")
525
526
527     # Shift the X data over a validation set size to ensure a new
528     # validation set data for the next training iteration
529     self.X = np.roll(self.X, -self.validation_set_size, axis=0)
530     self.Y = np.roll(self.Y, -self.validation_set_size, axis=0)
531
532     # Normalize error of all iterations of the validation set
533     self.train_err /= self.k
534     self.test_err /= self.k
535
536     print("-----")
537
538     total_acc_train = self.train_err * 100
539     total_acc_test = self.test_err * 100
540
541     print("Total Accuracy (Training Data): ",
542         float(f'{total_acc_train:.2f}'), "%")
543     print("Total Accuracy (Validation Data): ",

```



```

544         float(f'{total_acc_test:.2f}'), "%")
545
546     print("Confusion matrix for final fold: ")
547     self.train_eval.confusion_matrix()
548     return(total_acc_test)
549
550
551
552     """# File Upload
553
554     Here you can upload the CSV you intend to use for the rest of the notebook.
555
556
557     NOTE: Running these cells prompts you to upload a CSV files. Choose the
558     ↪ original `air_quality.csv` or `liver_disease.csv` files. Be sure to
559     ↪ replace the filename in the `io.BytesIO` function with the name of the
560     ↪ file you intend to upload
561
562     """
563
564     #importing data from either google collab or locally
565     import io
566     from google.colab import files, drive
567
568     uploaded = files.upload()
569     df = pd.read_csv(io.BytesIO(uploaded['air_quality.csv']))
570     data = df.to_numpy()
571
572     # Main Program - Base Model
573     """# Main Program - Optimal Normalized Model
574
575     Below is the optimal model discussed in the report: a variation on the base
576     ↪ model using a normalization step implemented on the dataset before
577     ↪ training, a feature reduction, a feature expansion, and 11 different PCA
578     ↪ sets cross-validated with each other using 10-fold.
579
580     """
581
582     import pandas as pd, numpy as np
583
584     def main():
585         # Hyperparameters
586         step_size = 0.001 #step sizes on the gradient descent
587         iterations = 100000 #max epoch count of the test
588         folds = 10 #number of folds on the kfolds
589         normalization = 1 #normalization enabling bit (1 = active, 0 = inactive)
590         regularization = 0 #regularization enabling bit (2 = active L2
591         ↪ regularization, 1 = active L1 regularization, 0 = inactive)
592         pca = 1 #pca enabling bit (1 = active, 0 = inactive)
593         pca_size = 11 #initiates the size of the pca set
594         expansion = 1 #expansion enabling bit (1 = active, 0 = inactive)
595         brute_force = 0 #brute_force enabling bit (1 = active, 0 = inactive)
596
597         # Trains the model with no additional techniques using 10-fold cross
598         ↪ validation
599         k = KFold(data,
600                 folds,
601                 step_size,
602                 iterations,
603                 normalization,
604                 regularization,
605                 pca,
606                 expansion,
607                 pca_size,

```

```
600         brute_force)
601
602     acc = k.cross_validation()
603
604
605     print("-----")
606     print("Total Accuracy: ", acc)
607     print("-----")
608
609
610 if __name__ == "__main__":
611     main()
612
```