

### 3. Diabetes Prediction

In this section, we will train a logistic regression model using stochastic gradient descent on the diabetes dataset.

The example assumes that a CSV copy of the dataset is in the current working directory with the filename **pima-indians-diabetes.csv**.

The dataset is first loaded, the string values converted to numeric and each column is normalized to values in the range of 0 to 1. This is achieved with the helper functions **load\_csv()** and **str\_column\_to\_float()** to load and prepare the dataset and **dataset\_minmax()** and **normalize\_dataset()** to normalize it.

We will use k-fold cross validation to estimate the performance of the learned model on unseen data. This means that we will construct and evaluate k models and estimate the performance as the mean model performance. Classification accuracy will be used to evaluate each model. These behaviors are provided in the **cross\_validation\_split()**, **accuracy\_metric()** and **evaluate\_algorithm()** helper functions.

We will use the **predict()**, **coefficients\_sgd()** functions created above and a new **logistic\_regression()** function to train the model.

Below is the complete example.



```
1 # Logistic Regression on Diabetes Dataset
2 from random import seed
3 from random import randrange
4 from csv import reader
5 from math import exp
6
7 # Load a CSV file
8 def load_csv(filename):
9     dataset = list()
10    with open(filename, 'r') as file:
11        csv_reader = reader(file)
12        for row in csv_reader:
13            if not row:
14                continue
15            dataset.append(row)
16    return dataset
17
18 # Convert string column to float
19 def str_column_to_float(dataset, column):
20     for row in dataset:
21         row[column] = float(row[column].strip())
22
23 # Find the min and max values for each column
24 def dataset_minmax(dataset):
25     minmax = list()
26     for i in range(len(dataset[0])):
27         col_values = [row[i] for row in dataset]
28         value_min = min(col_values)
29         value_max = max(col_values)
30         minmax.append([value_min, value_max])
31     return minmax
32
33 # Rescale dataset columns to the range 0-1
34 def normalize_dataset(dataset, minmax):
35     for row in dataset:
36         for i in range(len(row)):
37             row[i] = (row[i] - minmax[i][0]) / (minmax[i][1] - minmax[i][0])
```

```

38
39 # Split a dataset into k folds
40 def cross_validation_split(dataset, n_folds):
41     dataset_split = list()
42     dataset_copy = list(dataset)
43     fold_size = int(len(dataset) / n_folds)
44     for i in range(n_folds):
45         fold = list()
46         while len(fold) < fold_size:
47             index = randrange(len(dataset_copy))
48             fold.append(dataset_copy.pop(index))
49         dataset_split.append(fold)
50     return dataset_split
51
52 # Calculate accuracy percentage
53 def accuracy_metric(actual, predicted):
54     correct = 0
55     for i in range(len(actual)):
56         if actual[i] == predicted[i]:
57             correct += 1
58     return correct / float(len(actual)) * 100.0
59
60 # Evaluate an algorithm using a cross validation split
61 def evaluate_algorithm(dataset, algorithm, n_folds, *args):
62     folds = cross_validation_split(dataset, n_folds)
63     scores = list()
64     for fold in folds:
65         train_set = list(folds)
66         train_set.remove(fold)
67         train_set = sum(train_set, [])
68         test_set = list()
69         for row in fold:
70             row_copy = list(row)
71             test_set.append(row_copy)
72             row_copy[-1] = None
73         predicted = algorithm(train_set, test_set, *args)
74         actual = [row[-1] for row in fold]
75         accuracy = accuracy_metric(actual, predicted)
76         scores.append(accuracy)
77     return scores
78
79 # Make a prediction with coefficients
80 def predict(row, coefficients):
81     yhat = coefficients[0]
82     for i in range(len(row)-1):
83         yhat += coefficients[i + 1] * row[i]
84     return 1.0 / (1.0 + exp(-yhat))
85
86 # Estimate logistic regression coefficients using stochastic gradient descent
87 def coefficients_sgd(train, l_rate, n_epoch):
88     coef = [0.0 for i in range(len(train[0]))]
89     for epoch in range(n_epoch):
90         for row in train:
91             yhat = predict(row, coef)
92             error = row[-1] - yhat
93             coef[0] = coef[0] + l_rate * error * yhat * (1.0 - yhat)
94             for i in range(len(row)-1):
95                 coef[i + 1] = coef[i + 1] + l_rate * error * yhat * (1.0 - yhat) * row[i]
96     return coef
97
98 # Linear Regression Algorithm With Stochastic Gradient Descent
99 def logistic_regression(train, test, l_rate, n_epoch):
100     predictions = list()
101     coef = coefficients_sgd(train, l_rate, n_epoch)
102     for row in test:
103         yhat = predict(row, coef)
104         yhat = round(yhat)
105         predictions.append(yhat)
106     return(predictions)
107
108 # Test the logistic regression algorithm on the diabetes dataset

```

```

109 seed(1)
110 # load and prepare data
111 filename = 'pima-indians-diabetes.csv'
112 dataset = load_csv(filename)
113 for i in range(len(dataset[0])):
114     str_column_to_float(dataset, i)
115 # normalize
116 minmax = dataset_minmax(dataset)
117 normalize_dataset(dataset, minmax)
118 # evaluate algorithm
119 n_folds = 5
120 l_rate = 0.1
121 n_epoch = 100
122 scores = evaluate_algorithm(dataset, logistic_regression, n_folds, l_rate, n_epoch)
123 print('Scores: %s' % scores)
124 print('Mean Accuracy: %.3f%%' % (sum(scores)/float(len(scores))))

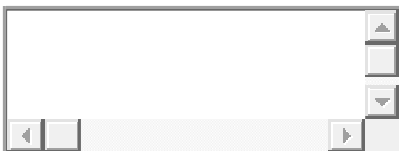
```

A k value of 5 was used for cross-validation, giving each fold  $768/5 = 153.6$  or just over 150 records to be evaluated upon each iteration. A learning rate of 0.1 and 100 training epochs were chosen with a little experimentation.

You can try your own configurations and see if you can beat my score.

Running this example prints the scores for each of the 5 cross-validation folds, then prints the mean classification accuracy.

We can see that the accuracy is about 77%, higher than the baseline value of 65% if we just predicted the majority class using the Zero Rule Algorithm.



```

1 Scores: [73.8562091503268, 78.43137254901961, 81.69934640522875, 75.81699346405229, 75.81699346405229]
2 Mean Accuracy: 77.124%

```

## Extensions

This section lists a number of extensions to this tutorial that you may wish to consider exploring.

- **Tune The Example.** Tune the learning rate, number of epochs and even data preparation method to get an improved score on the dataset.
- **Batch Stochastic Gradient Descent.** Change the stochastic gradient descent algorithm to accumulate updates across each epoch and only update the coefficients in a batch at the end of the epoch.
- **Additional Classification Problems.** Apply the technique to other binary (2 class) classification problems on the UCI machine learning repository.

**Did you explore any of these extensions?**

Let me know about it in the comments below.

## Review

In this tutorial, you discovered how to implement logistic regression using stochastic gradient descent from scratch with Python.

You learned.

- How to make predictions for a multivariate classification problem.
- How to optimize a set of coefficients using stochastic gradient descent.
- How to apply the technique to a real classification predictive modeling problem.