

03_k_fold

December 21, 2019

1 Kaggle Titanic survival - K-fold stratification

In our previous example* using logistic regression to classify passengers as likely to survive the Titanic, we used a random split for training and test data. But doing a single assessment like this may lead to an inaccurate assesment of the accuracy.

We could use repeated random splits, but a more robust method is to use ‘stratified k-fold validation’. In this method the model is repeated k times, so that all the data is used once, but only once, as part of the test set. This, alone, is k-fold validation. Stratified k-fold validation adds an extra level of robustness by ensuring that in each of the k training/test splits, the balance of outcomes represents the balance of outcomes (between survivors and non-survivors)in the overall data set. Most commonly 5 or 10 different splits of the data are used.

In a full project it is common to also have some hold-back test data that is used only at the end of model development (with k-fold validation used during model development).

*https://github.com/MichaelAllen1966/1804_python_healthcare/blob/master/titanic/02_logistic_regression.ipynb

In this notebook we assume that you have run through the basic logistic regression example in the previous example. We will not explain all setps fully.

1.1 Load modules

```
[1]: import numpy as np
import pandas as pd
from sklearn.linear_model import LogisticRegression
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import StratifiedKFold
```

1.2 Load data

The loading of data assumes that data has been downloaded and saved. If this has not been done yet, please refer to the loading of data in the previous example (see link in introduction).

```
[2]: data = pd.read_csv('data/processed_data.csv')
```

```
[3]: # Drop Passengerid (axis=1 indicates we are removing a column rather than a row)
# We drop passenger ID as it is not original data

data.drop('PassengerId', inplace=True, axis=1)
```

1.3 Divide into X (features) and y (lables)

We will split into features (X) and label (y) and convert from a Pandas DataFrame to NumPy arrays. NumPy arrays are simpler to refer to by row/column index numbers, and sklearn's k-fold method provides row indices for each set.

```
[4]: # Split data into two DataFrames
X_df = data.drop('Survived',axis=1)
y_df = data['Survived']

# Convert DataFrames to NumPy arrays
X = X_df.values
y = y_df.values
```

1.4 Define function to standardise data

Standardisation subtracts the mean and divides by the standard deviation, for each feature. Here we use the sklearn built-in method for standardisation.

```
[5]: def standardise_data(X_train, X_test):
    """
    Converts all data to a similar scale.
    Standardisation subtracts mean and divides by standard deviation
    for each feature.
    Standardised data will have a mean of 0 and standard deviation of 1.
    The training data mean and standard deviation is used to standardise both
    training and test set data.
    """

    # Initialise a new scaling object for normalising input data
    sc = StandardScaler()

    # Set up the scaler just on the training set
    sc.fit(X_train)

    # Apply the scaler to the training and test sets
    train_std=sc.transform(X_train)
    test_std=sc.transform(X_test)

    return train_std, test_std
```

1.5 Training and testing the model for all k-fold splits

The following code: * Sets up lists to hold results for each k-fold split * Sets up the k-fold splits using sklearn's StratifiedKFold method * Trains a logistic regression model, and test it, for each k-fold split * Adds each k-fold training/test accuracy to the lists

Here we are using a simple accuracy score, the proportion of predicts that are correct. K-fold validation may also be used for more complicated accuracy assessment.

```
[6]: # Set up lists to hold results for each k-fold run
training_acc_results = []
test_acc_results = []

# Set up splits
number_of_splits = 10
skf = StratifiedKFold(n_splits = number_of_splits)
skf.get_n_splits(X, y)

# Loop through the k-fold splits
for train_index, test_index in skf.split(X, y):

    # Get X and Y train/test
    X_train, X_test = X[train_index], X[test_index]
    y_train, y_test = y[train_index], y[test_index]

    # Standardise X data
    X_train_std, X_test_std = standardise_data(X_train, X_test)

    # Set up and fit model
    model = LogisticRegression(solver='lbfgs')
    model.fit(X_train_std, y_train)

    # Predict training and test set labels
    y_pred_train = model.predict(X_train_std)
    y_pred_test = model.predict(X_test_std)

    # Calculate accuracy of training and test sets
    accuracy_train = np.mean(y_pred_train == y_train)
    accuracy_test = np.mean(y_pred_test == y_test)

    # Add accuracy to lists
    training_acc_results.append(accuracy_train)
    test_acc_results.append(accuracy_test)
```

1.6 Show training and test results

```
[7]: # Show individual accuracies on training data  
training_acc_results
```

```
[7]: [0.8139825218476904,  
      0.8102372034956304,  
      0.816708229426434,  
      0.8017456359102244,  
      0.816708229426434,  
      0.8154613466334164,  
      0.814214463840399,  
      0.8117206982543641,  
      0.8104738154613467,  
      0.8144458281444583]
```

```
[8]: # Show individual accuracies on test data  
test_acc_results
```

```
[8]: [0.7666666666666667,  
      0.7888888888888889,  
      0.7752808988764045,  
      0.8651685393258427,  
      0.7640449438202247,  
      0.7640449438202247,  
      0.7865168539325843,  
      0.797752808988764,  
      0.8202247191011236,  
      0.8068181818181818]
```

```
[9]: # Get mean results  
mean_training = np.mean(training_acc_results)  
mean_test = np.mean(test_acc_results)  
  
# Display each to three decimal places  
print ('{0:.3f}, {1:.3f}'.format(mean_training, mean_test))
```

0.813, 0.794

As expected, the average accuracy is better for the training than the test sets. This is due to the model being slightly ‘over-fitted’ to the training data, a topic we shall return to in a later notebook.

1.7 Plot results: Box Plot

Box plots show median (orange line), the seconds and third quartiles (the box), the range (excluding outliers), and any outliers as ‘whisker’ points. Outliers, by convention, are considered to be any

points outside of the quartiles ± 1.5 times the interquartile range. The limit for outliers may be changed using the optional `whis` argument in the boxplot.

Medians tend to be an easy reliable guide to the centre of a distribution (i.e. look at the medians to see whether a fit is improving or not, but also look at the box plot to see how much variability there is).

Test sets tend to be more variable in their accuracy measures. Can you think why?

```
[10]: import matplotlib.pyplot as plt
      %matplotlib inline

      # Set up X data
      x_for_box = [training_acc_results, test_acc_results]

      # Set up X labels
      labels = ['Training', 'Test']

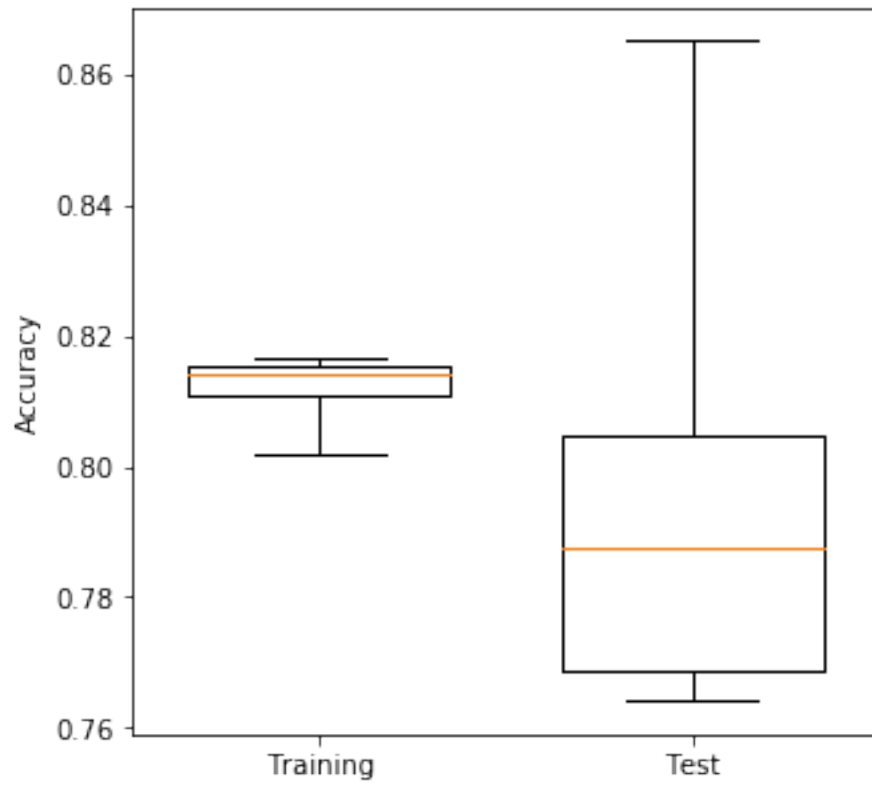
      # Set up figure
      fig = plt.figure(figsize=(5,5))

      # Add subplot (can be used to define multiple plots in same figure)
      ax1 = fig.add_subplot(1,1,1)

      # Define Box Plot (`widths` is optional)
      ax1.boxplot(x_for_box,
                  widths=0.7,
                  whis=10)

      # Set X and Y labels
      ax1.set_xticklabels(labels)
      ax1.set_ylabel('Accuracy')

      # Show plot
      plt.show()
```



[]: