

```

1  # -*- coding: utf-8 -*-
2  """
3  Created on Mon Mar 12 20:21:42 2018
4
5  @author: cml
6  """
7
8  import numpy as np
9  from sklearn import linear_model
10 from sklearn.metrics import mean_squared_error, r2_score
11 import pandas as pd
12 from sklearn.model_selection import KFold # import KFold
13
14 df = pd.read_csv('Auto.csv', usecols=range(1,10))
15
16 X = df["horsepower"].values.reshape(-1,1) # our independent variable
17 y = df["mpg"].values.reshape(-1,1) # our dependent variable
18
19 kf = KFold(n_splits=10) # Define the split into 2 folds
20 print('Splits: ', kf.get_n_splits(X))
21
22 #Arrays to store test data and predictions for each run
23 ytests = []
24 ypreds = []
25
26 #for each KFold split in X, fit a model using current x_train
27 #and y_train value. Save the array
28 for train_index, test_index in kf.split(X):
29     print("TRAIN: ", train_index, "TEST: ", test_index)
30     X_train, X_test = X[train_index], X[test_index]
31     y_train, y_test = y[train_index], y[test_index]
32
33     model = linear_model.LinearRegression()
34     model.fit(X = X_train, y = y_train)
35     y_pred = model.predict(X_test)
36
37     ytests += list(y_test)
38     ypreds += list(y_pred)
39
40 rr = r2_score(ytests, ypreds)
41 ms_error = mean_squared_error(ytests, ypreds)
42
43 print("KFOLD results:")
44 print("R^2: {:.5f}%, MSE: {:.5f}".format(rr*100, ms_error))
45
46 #ms_error -> ~24 when n_splits increases.
47 #We use k-1 subsets to train our data.
48 #Large k means less bias towards overestimating the true expected error.
49

```