```
1  import pandas as pd
2  import math as m
3  import numpy as np
4  import scipy.stats as stats
5  import matplotlib.pyplot as plt
6  plt.style.use('ggplot')
7  import statsmodels.api as sm
8
9
10 from sklearn.cross_validation import KFold
11 from sklearn import cross_validation
12 from sklearn.preprocessing import Imputer
13 from sklearn.linear_model import LinearRegression
14 from sklearn.ensemble import RandomForestRegressor
15 from sklearn.ensemble import AdaBoostRegressor
16 from sklearn.decomposition import PCA
17 from sklearn.linear_model import Ridge
18
19
20
21 df_train_rdkit = pd.read_csv("train_rdkit.csv")
22 df_train_SM = pd.read_csv("df_train_withSMFeatures.csv")
23 # df_train = pd.read_csv("train.csv")
24 # dfs = [df_train_rdkit,df_train_SM,df_train]
25 # df_train_all = reduce(lambda left,right: pd.merge(left,right,on='smiles'),
26 df_train_all = pd.merge(df_train_SM,df_train_rdkit,on="smiles")
27
28
29 ############ Data preprocessing ############
30
31 ## Clean the data a little ##
32 df_train_X = df_train_all.drop(['smiles'],axis=1) ##Use all the columns exce
33 df_train_X=df_train_X.dropna(axis=1,how='all') ## Drop all the columns that
34 df_train_X =df_train_X.loc[:, (df_train_X != 0).any(axis=0)] ## Drop the all
35
36 for column in df_train_X:
37     df_train_X[column].replace([np.inf, -np.inf], 0)  ## Replace all the inf
38
39 ## Remove NaNs and center them at the mean
40 imp = Imputer(missing_values='NaN', strategy='mean', axis=0)
41 imp.fit(df_train_X)
42 df_train_X = imp.transform(df_train_X)
43
44 ## Create the Ys ##
45 df_train_Y =  pd.read_csv("train.csv",usecols =['gap']).values
46
```

```
47
48  print "Shape of X train is ",df_train_X.shape
49  print "Shape of Y train is ", df_train_Y.shape
50
51
52
53  degree = 1 ## Set the degree you want your basis function
54  ### Sinsusoidal basis transformation
55  def makeMatrixsinusoid(yourMatrix,deg):
56      f= np.vectorize((lambda x: m.sin(x/deg)))
57      holder = yourMatrix
58      for i in xrange(1,deg+1):
59          holder= np.append(holder,f(yourMatrix),1)
60      return holder
61
62
63  # def makeMatrixsinusoid(yourMatrix,deg,numtimes):
64  #    holder = yourMatrix
65  #    for i in xrange(1,numtimes+1):
66  #        for j in xrange(yourMatrix.shape[1]):
67  #            print holder.shape
68  #            holder = np.append(holder,makeArraysinusoid(yourMatrix[:,j],i).r
69  #        return holder
70
71
72  def runLR(X_tr,y_tr,X_ts,y_ts):
73      LR = LinearRegression()
74      LR.fit(X_train, y_train)
75
76      RMSE_LR = m.sqrt(np.mean((LR.predict(X_test)-y_test)**2))
77      print "RMSE for Baseline Linear Regression is %s"%RMSE_LR
78
79
80  ## Use PCA for dimensionality reduction
81  def PCAreduction(matrix,dims):
82
83      pca = PCA(n_components=dims, whiten=True).fit(matrix)
84      print "Explaing variance for each PC is %s"%pca.explained_variance_ratio
85      print "Total variance explained is %s"%pca.explained_variance_ratio_.sum
86      return pca.transform(matrix)
87
88  def main():
89      df_train_X_sine = makeArraysinusoid(df_train_X,2)
90      # Create simple train/validate set from the training set
91      X_train, X_test, y_train, y_test = cross_validation.train_test_split(df_
92
```

```
 93
 94
 95
 96
 97  if __name__ == "__main__":
 98      # execute only if run as a script
 99      main()
100
101
```