

```
#####
```

```
''' Project No. 1: Mercedes-Benz Greener Manufacturing '''
```

```
#####
```

```
# Step1: Import the required libraries
```

```
# linear algebra
```

```
import numpy as np
```

```
# data processing, CSV file I/O (e.g. pd.read_csv)
```

```
import pandas as pd
```

```
# for dimensionality reduction
```

```
from sklearn.decomposition import PCA
```

```
# Step2: Read the data from train.csv
```

```
df_train = pd.read_csv('train.csv')
```

```
# let us understand the data
```

```
print('Size of training set: {} rows and {} columns'
```

```
      .format(*df_train.shape))
```

```
# print few rows and see how the data looks like
```

```
df_train.head()
```

```
# Step3: Collect the Y values into an array
```

```
# seperate the y from the data as we will use this to learn as
```

```
# the prediction output
```

```
y_train = df_train['y'].values
```

Step4: Understand the data types we have

iterate through all the columns which has X in the name of the column

```
cols = [c for c in df_train.columns if 'X' in c]
```

```
print('Number of features: {}'.format(len(cols)))
```

```
print('Feature types:')
```

```
df_train[cols].dtypes.value_counts()
```

Step5: Count the data in each of the columns

```
counts = [], [], []
```

```
for c in cols:
```

```
    typ = df_train[c].dtype
```

```
    uniq = len(np.unique(df_train[c]))
```

```
    if uniq == 1:
```

```
        counts[0].append(c)
```

```
    elif uniq == 2 and typ == np.int64:
```

```
        counts[1].append(c)
```

```
    else:
```

```
        counts[2].append(c)
```

```
print('Constant features: {} Binary features: {} Categorical features: {}'.format(*[len(c) for c in counts]))
```

```
print('Constant features:', counts[0])
```

```
print('Categorical features:', counts[2])
```

Step6: Read the test.csv data

```
df_test = pd.read_csv('test.csv')
```

remove columns ID and Y from the data as they are not used for learning

```
usable_columns = list(set(df_train.columns) - set(['ID', 'y']))
```

```
y_train = df_train['y'].values
```

```
id_test = df_test['ID'].values
```

```
x_train = df_train[usable_columns]
```

```
x_test = df_test[usable_columns]
```

Step7: Check for null and unique values for test and train sets

```
def check_missing_values(df):
```

```
    if df.isnull().any().any():
```

```
        print("There are missing values in the dataframe")
```

```
    else:
```

```
        print("There are no missing values in the dataframe")
```

```
check_missing_values(x_train)
```

```
check_missing_values(x_test)
```

Step8: If for any column(s), the variance is equal to zero,

then you need to remove those variable(s).

Apply label encoder

```
for column in usable_columns:

    cardinality = len(np.unique(x_train[column]))

    if cardinality == 1:

        x_train.drop(column, axis=1) # Column with only one
        # value is useless so we drop it
        x_test.drop(column, axis=1)

    if cardinality > 2: # Column is categorical

        mapper = lambda x: sum([ord(digit) for digit in x])

        x_train[column] = x_train[column].apply(mapper)

        x_test[column] = x_test[column].apply(mapper)

x_train.head()
```

Step9: Make sure the data is now changed into numericals

```
print('Feature types:')

x_train[cols].dtypes.value_counts()
```

Step10: Perform dimensionality reduction

Linear dimensionality reduction using Singular Value Decomposition of
the data to project it to a lower dimensional space.

```
n_comp = 12

pca = PCA(n_components=n_comp, random_state=420)

pca2_results_train = pca.fit_transform(x_train)

pca2_results_test = pca.transform(x_test)
```

Step11: Training using xgboost

```
import xgboost as xgb

from sklearn.metrics import r2_score

from sklearn.model_selection import train_test_split


x_train, x_valid, y_train, y_valid = train_test_split(
    pca2_results_train,
    y_train, test_size=0.2,
    random_state=4242)


d_train = xgb.DMatrix(x_train, label=y_train)
d_valid = xgb.DMatrix(x_valid, label=y_valid)
#d_test = xgb.DMatrix(x_test)
d_test = xgb.DMatrix(pca2_results_test)


params = {}

params['objective'] = 'reg:linear'

params['eta'] = 0.02

params['max_depth'] = 4


def xgb_r2_score(preds, dtrain):
    labels = dtrain.get_label()
    return 'r2', r2_score(labels, preds)


watchlist = [(d_train, 'train'), (d_valid, 'valid')]


clf = xgb.train(params, d_train,
    1000, watchlist, early_stopping_rounds=50,
```

```
feval=xgb_r2_score, maximize=True, verbose_eval=10)
```

```
# Step12: Predict your test_df values using xgboost
```

```
p_test = clf.predict(d_test)
```

```
sub = pd.DataFrame()
```

```
sub['ID'] = id_test
```

```
sub['y'] = p_test
```

```
sub.to_csv('xgb.csv', index=False)
```

```
sub.head()
```

```
#####
```

```
'''          End          '''
```

```
#####
```