

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')

In [2]: df = pd.read_csv(r'C:\Users\sutharsan\Downloads\Inbytes airpollution T2\PSRA_data_2010.1.1-2014.12.31.csv')

In [3]: df.head(5)

Out[3]:
   No  year  month  day  hour  pm2.5  DEWP  TEMP  PRES  cbwd  lws  ls  lr
0  1  2010    1    1    0  NaN      -21  -11.0  1021.0  NW  1.79  0  0
1  2  2010    1    1    1  NaN      -21  -12.0  1020.0  NW  4.92  0  0
2  3  2010    1    1    2  NaN      -21  -11.0  1019.0  NW  6.71  0  0
3  4  2010    1    1    3  NaN      -21  -14.0  1019.0  NW  9.94  0  0
4  5  2010    1    1    4  NaN      -20  -12.0  1018.0  NW  12.97  0  0

In [4]: df.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 43824 entries, 0 to 43823
Data columns (total 13 columns):
 #   Column  Non-Null Count  Dtype
--  --
0  No      43824 non-null      int64
1  year    43824 non-null      int64
2  month   43824 non-null      int64
3  day     43824 non-null      int64
4  hour    43824 non-null      int64
5  pm2.5   41757 non-null      float64
6  DEWP    43824 non-null      int64
7  TEMP    43824 non-null      float64
8  PRES    43824 non-null      float64
9  cbwd    43824 non-null      object
10 lws     43824 non-null      float64
11 ls     43824 non-null      int64
12 lr     43824 non-null      int64
dtypes: float64(4), int64(8), object(1)
memory usage: 4.2+ MB

In [5]: df.shape # rows and columns

Out[5]:
(43824, 13)

In [6]: df.describe().transpose()

Out[6]:
      count      mean      std      min      25%      50%      75%      max
No  43824.0  21912.500000  12691.043435    1.00  10956.75  21912.50  32868.25  43824.0
year  43824.0  2012.000000    1.413842  2010.00  2011.00  2012.00  2013.00  2014.0
month  43824.0  15.527549    3.448573    1.00    4.00    7.00   10.00   12.0
day  43824.0  15.727820    8.799425    1.00    8.00   16.00   23.00   31.0
hour  43824.0  11.500000    6.922266    0.00    5.75   11.50   17.25   23.0
pm2.5  41757.0  98.613215    92.050387    0.00   29.00   72.00  137.00  994.0
DEWP  43824.0  1.817261    14.433440   -40.00  -10.00    2.00   15.00   28.0
TEMP  43824.0  12.448521    12.198613   -19.00    2.00   14.00   23.00   42.0
PRES  43824.0  1016.447654    10.268698   991.00  1008.00  1016.00  1025.00  1046.0
lws  43824.0  23.889140    50.010635    0.45    1.79    5.37   21.91  585.6
ls  43824.0  0.052734    0.760375    0.00    0.00    0.00    0.00   27.0
lr  43824.0  0.194916    1.415967    0.00    0.00    0.00    0.00   36.0

In [7]: # check missing values
df.isnull().sum() # in target we have null values.

Out[7]:
No      0
year     0
month    0
day       0
hour      0
pm2.5  2667
DEWP      0
TEMP      0
PRES      0
cbwd      0
lws       0
ls        0
ls        0
lr        0
dtype: int64

In [8]: # to deal with null value we can dropna,ffill, bfill,3m.
df.dropna(axis=0,inplace=True)

In [9]: df.isnull().sum() # drop the missing rows-->axis=0

Out[9]:
No      0
year     0
month    0
day       0
hour      0
pm2.5    0
DEWP      0
TEMP      0
PRES      0
cbwd      0
lws       0
ls        0
ls        0
lr        0
dtype: int64

In [10]: # categorical variable
df.dtypes=="object"

Out[10]:
No      False
year    False
month   False
day     False
hour    False
pm2.5   False
DEWP    False
PRES    False
cbwd    True
lws     False
ls      False
ls      False
lr      False
dtype: bool

In [11]: df = df[df['DEWP'] > 0]
df = df[df['TEMP'] > 0]

In [12]: df.shape

Out[12]:
(21758, 13)

In [13]: df['cbwd'].nunique() # 4 unique value

Out[13]:
4

In [14]: df1 = pd.get_dummies(df['cbwd'])

In [15]: df = pd.concat([df,df1,axis=1)

In [16]: df

Out[16]:
      No  year  month  day  hour  pm2.5  DEWP  TEMP  PRES  cbwd  lws  ls  lr  NE  NW  SE  cv
1315  1316  2010    2    24  19  368.0    2    6.0  1006.0  cv  0.89  0  0  0  0  1  0
1316  1317  2010    2    24  20  309.0    3    7.0  1007.0  SE  1.79  0  0  0  0  1  0
1317  1318  2010    2    24  21  284.0    3    6.0  1008.0  SE  3.58  0  0  0  0  1  0
1318  1319  2010    2    24  22  244.0    2    6.0  1009.0  SE  6.71  0  0  0  0  1  0
1319  1320  2010    2    24  23  218.0    3    5.0  1010.0  cv  0.89  0  0  0  0  0  1
...
43059 43060 2014    11    30  3  194.0    3    4.0  1021.0  NW  3.58  0  0  0  1  0  0
43060 43061 2014    11    30  4  227.0    2    4.0  1021.0  NW  6.71  0  0  0  1  0  0
43061 43062 2014    11    30  5  242.0    2    3.0  1021.0  NW  9.94  0  0  0  1  0  0
43062 43063 2014    11    30  6  235.0    2    3.0  1021.0  NW  11.63  0  0  0  1  1  0
43063 43064 2014    11    30  7  168.0    2    4.0  1022.0  NW  14.76  0  0  0  1  0  0

21758 rows x 17 columns

In [17]: df.drop('cbwd',axis=1,inplace=True)

In [18]: df.head(2) # cbwd is removed and got dummies.

Out[18]:
      No  year  month  day  hour  pm2.5  DEWP  TEMP  PRES  lws  ls  lr  NE  NW  SE  cv
1315  1316  2010    2    24  19  368.0    2    6.0  1006.0  0.89  0  0  0  0  0  1
1316  1317  2010    2    24  20  309.0    3    7.0  1007.0  1.79  0  0  0  0  1  0

In [19]: # plt.figure(figsize=(18,10))
sns.distplot(df['pm2.5'],bins=50)

<AxesSubplot: xlabel='pm2.5', ylabel='Density'>

Density
0.006
0.005
0.004
0.003
0.002
0.001
0.000
0 100 200 300 400 500 600 700
pm2.5

In [20]: # df['PM2.5'] = np.log(df['pm2.5'])

In [21]: # visual
plt.figure(figsize=(18,10))
sns.heatmap(df.corr(),annot=True)

<AxesSubplot: >

No 1 0.99 0.081 0.024 -0.0032 -0.1 -0.027 0.011 0.031 -0.084 -0.0022 -0.037 0.02 -0.022 -0.027 0.037
year -0.99 1 0.033 0.027 -0.005 -0.11 -0.033 0.029 -0.02 -0.068 -0.00031 -0.04 0.015 -0.032 -0.0072 0.027
month 0.081 -0.033 1 0.18 0.01 0.047 0.049 -0.16 0.43 -0.14 -0.015 0.026 0.042 0.094 -0.18 0.087
day 0.024 0.027 -0.18 1 0.0017 0.093 0.0067 3.9e-06 0.01 0.015 -0.011 0.0023 -0.009 -0.04 0.039 -0.002
hour -0.032 -0.005 0.01 -0.0017 1 0.018 0.0014 0.25 -0.033 0.11 -0.0035 0.0094 -0.078 -0.14 0.24 -0.093
pm2.5 -0.1 -0.11 0.047 0.093 -0.018 1 0.06 -0.22 0.11 -0.11 0.00076 -0.096 -0.03 -0.18 0.081 0.091
DEWP 0.027 0.033 0.049 0.0067 0.0014 0.06 1 0.62 -0.55 0.08 0.013 0.038 0.014 -0.14 0.14 -0.019
TEMP -0.011 0.029 -0.16 3.9e-06 0.25 -0.22 0.62 1 0.66 0.076 -0.019 -0.11 -0.073 -0.17 0.28 -0.12
PRES 0.031 -0.02 0.43 0.01 -0.053 0.11 -0.55 0.66 1 0.068 0.016 0.014 0.033 0.076 -0.16 0.09
lws -0.084 0.068 0.14 0.015 0.11 -0.11 -0.08 0.076 -0.068 1 0.0029 0.057 -0.14 0.063 0.28 -0.28
ls -0.0022 -0.00031 0.015 -0.011 0.0035 0.00076 -0.013 -0.019 0.016 -0.0028 1 0.0013 -0.0024 -0.0034 0.0075 0.0038
lr -0.037 -0.04 0.026 -0.0023 -0.0094 -0.096 0.038 -0.11 0.014 0.057 -0.0013 1 0.053 0.12 -0.091 -0.04
NE 0.02 0.015 0.042 -0.009 -0.078 -0.03 -0.014 -0.073 0.033 -0.14 -0.0024 0.053 1 -0.17 -0.32 -0.2
NW -0.022 -0.032 0.094 -0.04 -0.14 -0.18 -0.14 -0.17 0.076 0.063 -0.0034 0.12 -0.17 1 -0.45 -0.28
SE -0.027 -0.0072 -0.18 0.039 0.24 0.081 0.14 0.28 -0.16 0.28 0.0075 -0.091 -0.32 -0.45 1 -0.51
cv 0.037 0.027 0.087 -0.002 0.093 0.091 -0.019 -0.12 0.09 -0.28 -0.0038 -0.04 -0.2 -0.28 -0.51 1

No year month day hour pm2.5 DEWP TEMP PRES lws ls lr NE NW SE cv

In [22]: df['year'].unique()

Out[22]:
array([2010, 2011, 2012, 2013, 2014], dtype=int64)

In [23]: plt.figure(figsize=(18,10))
sns.histplot(df['TEMP'],hue='year',data=df)

<AxesSubplot: xlabel='TEMP', ylabel='pm2.5'>

pm2.5
700
600
500
400
300
200
100
0
0 10 20 30 40
TEMP
year
2010
2011
2012
2013
2014

In [24]: df['pm2.5'][100:1000].plot()

<AxesSubplot: >

400
350
300
250
200
150
100
50
0
2400 2600 2800 3000 3200 3400 3600

In [25]: df['TEMP'][100:1000].plot()

<AxesSubplot: >

35
30
25
20
15
10
5
0
2400 2600 2800 3000 3200 3400 3600

In [26]: df['DEWP'][100:1000].plot()

<AxesSubplot: >

17.5
15.0
12.5
10.0
7.5
5.0
2.5
0
2400 2600 2800 3000 3200 3400 3600

In [27]: df.columns

Index(['No', 'year', 'month', 'day', 'hour', 'pm2.5', 'DEWP', 'TEMP', 'PRES', 'lws', 'ls', 'lr', 'NE', 'NW', 'SE', 'cv'],
      dtype='object')

In [28]: from sklearn.stats.outliers_influence import variance_inflation_factor

df1 = df[['DEWP', 'TEMP', 'PRES', 'lws', 'ls', 'lr', 'NE', 'NW', 'SE', 'cv']]
import statsmodels as sm
data = sm.tools.add_constant(df1)
series = pd.Series([variance_inflation_factor(df1.values,i) for i in range(df1.shape[1])],index=df1.columns)

Out[29]:
DEWP      1.817447
TEMP      2.324708
PRES      1.896835
lws       1.174493
ls        1.080595
lr        1.057735
NE       4187.294182
NW       7581.569333
SE       17215.065842
cv        8298.934220
dtype: float64

In [30]: df2 = df[['DEWP', 'TEMP', 'PRES', 'lws', 'ls', 'lr']]
import statsmodels as sm
data = sm.tools.add_constant(df2)
series = pd.Series([variance_inflation_factor(df2.values,i) for i in range(df2.shape[1])],index=df2.columns)

Out[30]:
DEWP      9.888935
TEMP     16.565522
PRES      9.901542
lws       1.418709
ls        1.000427
lr        1.077952
dtype: float64

In [31]: y = df['pm2.5']
x = df[['DEWP', 'TEMP', 'PRES', 'lws', 'ls', 'lr', 'NE', 'NW', 'SE', 'cv']]

In [32]: y.shape, x.shape

Out[32]:
((21758,), (21758, 10))

In [58]: # model implementation - splitting, scaling and model building
from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.3,random_state=42)

In [59]: from sklearn.preprocessing import StandardScaler
ss = StandardScaler()

In [60]: x_train = ss.fit_transform(x_train)
x_test = ss.transform(x_test)
print(x_train.shape)
print(x_test.shape)

(15230, 10)
(6528, 10)

In [61]: from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_squared_error,mean_absolute_error,explained_variance_score
mse = mean_squared_error
mae = mean_absolute_error

In [62]: RF = RandomForestRegressor()
RF.fit(x_train,y_train)
RF_score = RF.score(x_train,y_train)
RF_pred = RF.predict(x_test)
print('RF_Regressor',RF_score) # 88 % on training

RF_Regressor 0.881766924606274

In [63]: print('RF on testing',explained_variance_score(RF_pred,y_test)) # 66 % on testing like overfit so hyperparameter tuning
RF on testing -0.6691269415904197

In [64]: RF_mae = mae(y_test,RF_pred)
RF_rmse = mse(y_test,RF_pred)
RF_rmse = np.sqrt(RF_rmse)

In [65]: print('Random Forest Regressor')
RF_mae,RF_rmse,RF_rmse

Random Forest Regressor
(47.636537969136554, 4570.119042484793, 67.60265558692723)

Overfitting - no cross validation

In [86]: from sklearn.model_selection import RandomizedSearchCV
param_grid = {'n_estimators': [100, 200, 300],
              'max_depth': [None, 5, 10],
              'max_features': ['sqrt', 'log2']}

rf = RandomizedSearchCV(RF, param_distributions=param_grid,cv=10,n_iter=10)
rf.fit(x_train,y_train)
print('best_param',rf.best_params_)
print('score',rf.best_score_)

best_param {'max_features': 'log2', 'max_depth': 10}
score 0.296333787569263

In [89]: RF_model = rf.best_estimator_
test_score_rf = RF_model.score(x_test,y_test)
print('Random Forest score',test_score_rf) # 30 % on testing

Random Forest score 0.2915229793893898

Linear Regression

In [68]: from sklearn.decomposition import PCA
pca=PCA()
pca.fit(x_train)
print(pca.components_)
print(pca.explained_variance_)

[[[-6.2783922e-01 -5.47218003e-01 4.91659383e-01 -1.33807426e-01
-3.46944695e-16 6.95181227e-02 8.86380327e-02 2.00110377e-01
-3.74495206e-01 1.83769345e-01]
[-1.42548847e-02 -4.73640799e-02 6.47258021e-02 1.96384141e-01 5.05873782e-01
0.00808080e+00 4.53692858e-02 -1.81280346e-01 1.46512706e-01
4.71243226e-01 -5.51211112e-01]
[ 8.14637197e-01 9.78616253e-02 2.37547957e-01 8.25353658e-02
-2.22844695e-16 4.11809781e-01 1.68911923e-01 6.74907988e-01
-3.81667476e-01 -3.08461694e-01]
[ -1.42548847e-02 -4.73640799e-02 6.47258021e-02 1.23901835e-01
4.16333634e-17 5.0895543e-02 8.46335180e-01 -3.34408374e-01
6.23551222e-02 -3.78123229e-01]
[ 7.82526236e-02 9.27838304e-02 4.62046063e-02 1.67873438e-01
3.33866907e-16 6.68542064e-01 7.37822389e-02 -3.48629189e-01
1.11520255e-01 2.48913623e-01]
[-1.31804890e-01 8.17454439e-02 1.54676591e-01 7.92859554e-01
0.00808080e+00 -1.06308359e-01 2.57713778e-01 -1.19204457e-01
-3.3183704e-01 3.07488719e-01]
[-7.11848114e-01 2.68488954e-02 -6.57483830e-01 -2.09189156e-01
1.11022302e-16 1.07738202e-01 -1.75346377e-02 -5.07297861e-02
4.46939508e-02 8.13457324e-03]
[-3.53882140e-01 8.04317004e-01 4.48620957e-01 -6.12952115e-02
-1.11022302e-16 1.26432901e-01 1.93276361e-02 3.79676757e-02
-6.9635297e-02 2.67460384e-02]
[ 1.14062185e-16 -8.47323084e-16 3.56555229e-16 3.61793537e-16
-1.03974380e-04 -1.64840256e-16 3.76731293e-01 4.81033381e-01
5.90749031e-01 5.16781589e-01]
[ 6.01952919e-18 -1.14979795e-16 -1.14125537e-16 9.06440140e-17
-9.99999999e-01 1.51272569e-16 -3.91704829e-05 -5.00151459e-05
-6.23586278e-05 -5.37237194e-05]
[2.45432469e+00 1.62625214e+00 1.33726217e+00 1.16499753e+00
9.49404622e-01 7.42438330e-01 4.25624683e-01 3.00287417e-01
5.13286289e-30 3.78521440e-35]

x_train = pca.transform(x_train)
x_train.shape

(15230, 10)

In [70]: x_test = pca.transform(x_test)
x_test.shape

(6528, 10)

In [71]: from sklearn.linear_model import LinearRegression, Ridge

In [72]: lr = LinearRegression()
lr.fit(x_train,y_train)
train_score = lr.score(x_train,y_train)
lr_pred = lr.predict(x_test)
print('Linear Regression training score',train_score)

Linear Regression training score 0.1847777576132398

In [73]: lr_mae = mae(y_test,lr_pred)

In [74]: print('Linear Regression testing score',explained_variance_score(y_test,lr_pred))

Linear Regression testing score 0.178211899108883

In [79]: lr_mae = mae(y_test,lr_pred)
lr_rmse = mse(y_test,lr_pred)
lr_rmse = np.sqrt(lr_rmse)

In [80]: print('Ridge metrics')
lr_mae,lr_rmse,lr_rmse

Ridge metrics
(54.68661616561064, 5486.986978853476, 74.07419913879242)

In [81]: from sklearn.model_selection import RandomizedSearchCV
param_grid = {'solver':['auto','sqr','cholesky','lsqr','sparse_cg','slsq','saga','lbfgs']

ridge_model = RandomizedSearchCV(r,param_grid,cv=5)
ridge_train = fit(x_train,y_train)
print('best_param',ridge_train.best_params_)
print('score',ridge_train.best_score_)

best_param {'solver': 'sag', 'alpha': 0.01}
score 0.18351749966488105

In [82]: ridge_model = ridge_train.best_estimator_
test_score = ridge_model.score(x_test,y_test)
print('ridge test score',test_score)

ridge test score 0.1778910564257666

In [ ]:

In [ ]:
```