```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn import metrics
%matplotlib inline
import warnings
import os
from pandas import read_csv
warnings.filterwarnings('ignore')
print(os.listdir("C:/Users/anishm/OneDrive - Adobe/Documents/testdata"))
```

['housing.csv']

```
column_names = ['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX', 'PTRATIO', 'B', 'LSTAT'
data = pd.read_csv('C:/Users/anishm/OneDrive - Adobe/Documents/testdata/housing.csv', header=None,delimiter=r"\
data.head(5)
```

Out[114]:

|   | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT | MEDV |
|---|------|-----|-------|------|-------|-------|------|--------|-----|-------|---------|--------|-------|------|
| 0 | 0.00632 | 18.0 | 2.31 | 0 | 0.538 | 6.575 | 65.2 | 4.0900 | 1 | 296.0 | 15.3 | 396.90 | 4.98 | 24.0 |
| 1 | 0.02731 | 0.0 | 7.07 | 0 | 0.469 | 6.421 | 78.9 | 4.9671 | 2 | 242.0 | 17.8 | 396.90 | 9.14 | 21.6 |
| 2 | 0.02729 | 0.0 | 7.07 | 0 | 0.469 | 7.185 | 61.1 | 4.9671 | 2 | 242.0 | 17.8 | 392.83 | 4.03 | 34.7 |
| 3 | 0.03237 | 0.0 | 2.18 | 0 | 0.458 | 6.998 | 45.8 | 6.0622 | 3 | 222.0 | 18.7 | 394.63 | 2.94 | 33.4 |
| 4 | 0.06905 | 0.0 | 2.18 | 0 | 0.458 | 7.147 | 54.2 | 6.0622 | 3 | 222.0 | 18.7 | 396.90 | 5.33 | 36.2 |

```
print(data.describe())
x = data.drop(["MEDV"],axis =1)
y = data.filter(["MEDV"],axis = 1)
```

```
              CRIM          ZN       INDUS        CHAS         NOX          RM  \
count  506.000000  506.000000  506.000000  506.000000  506.000000  506.000000
mean     3.613524   11.363636   11.136779    0.069170    0.554695    6.284634
std      8.601545   23.322453    6.860353    0.253994    0.115878    0.702617
min      0.006320    0.000000    0.460000    0.000000    0.385000    3.561000
25%      0.082045    0.000000    5.190000    0.000000    0.449000    5.885500
50%      0.256510    0.000000    9.690000    0.000000    0.538000    6.208500
75%      3.677083   12.500000   18.100000    0.000000    0.624000    6.623500
max     88.976200  100.000000   27.740000    1.000000    0.871000    8.780000

              AGE         DIS         RAD         TAX     PTRATIO           B  \
count  506.000000  506.000000  506.000000  506.000000  506.000000  506.000000
mean    68.574901    3.795043    9.549407  408.237154   18.455534  356.674032
std     28.148861    2.105710    8.707259  168.537116    2.164946   91.294864
min      2.900000    1.129600    1.000000  187.000000   12.600000    0.320000
25%     45.025000    2.100175    4.000000  279.000000   17.400000  375.377500
50%     77.500000    3.207450    5.000000  330.000000   19.050000  391.440000
75%     94.075000    5.188425   24.000000  666.000000   20.200000  396.225000
max    100.000000   12.126500   24.000000  711.000000   22.000000  396.900000

            LSTAT        MEDV
count  506.000000  506.000000
mean    12.653063   22.532806
std      7.141062    9.197104
min      1.730000    5.000000
25%      6.950000   17.025000
50%     11.360000   21.200000
75%     16.955000   25.000000
max     37.970000   50.000000
```

```
x.head(5)
```

Out[86]:

|   | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT |
|---|------|-----|-------|------|-------|-------|------|--------|-----|-------|---------|--------|-------|
| 0 | 0.00632 | 18.0 | 2.31 | 0 | 0.538 | 6.575 | 65.2 | 4.0900 | 1 | 296.0 | 15.3 | 396.90 | 4.98 |
| 1 | 0.02731 | 0.0 | 7.07 | 0 | 0.469 | 6.421 | 78.9 | 4.9671 | 2 | 242.0 | 17.8 | 396.90 | 9.14 |
| 2 | 0.02729 | 0.0 | 7.07 | 0 | 0.469 | 7.185 | 61.1 | 4.9671 | 2 | 242.0 | 17.8 | 392.83 | 4.03 |
| 3 | 0.03237 | 0.0 | 2.18 | 0 | 0.458 | 6.998 | 45.8 | 6.0622 | 3 | 222.0 | 18.7 | 394.63 | 2.94 |
| 4 | 0.06905 | 0.0 | 2.18 | 0 | 0.458 | 7.147 | 54.2 | 6.0622 | 3 | 222.0 | 18.7 | 396.90 | 5.33 |

```
y.head(5)
```

| | MEDV |
|---|---|
| 0 | 24.0 |
| 1 | 21.6 |
| 2 | 34.7 |
| 3 | 33.4 |
| 4 | 36.2 |

```python
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=42)
house_predictor = LinearRegression()
house_predictor.fit(x_train,y_train)
y_pred=house_predictor.predict(x_test)
```

```python
print('Mean Absolute Error:' ,metrics.mean_absolute_error(y_test,y_pred))
print('Mean Squared Error:' ,metrics.mean_squared_error(y_test,y_pred))
print('Root Mean Squared Error:' , np.sqrt(metrics.mean_squared_error(y_test,y_pred)))
```

```
Mean Absolute Error: 3.1890919658878745
Mean Squared Error: 24.29111947497371
Root Mean Squared Error: 4.928602182665355
```

```python
comparison_df = pd.DataFrame({'Actual' : y_test.values.tolist(), 'Predicted': y_pred.tolist()})
comparison_df.head(5)
```

| | Actual | Predicted |
|---|---|---|
| 0 | [23.6] | [28.99672361982493] |
| 1 | [32.4] | [36.02556533567232] |
| 2 | [13.6] | [14.816944045388338] |
| 3 | [22.8] | [25.031979150399636] |
| 4 | [16.1] | [18.76987991524812] |

```python
print(house_predictor.coef_)
```

```
[[-1.13055924e-01  3.01104641e-02  4.03807204e-02  2.78443820e+00
  -1.72026334e+01  4.43883520e+00 -6.29636221e-03 -1.44786537e+00
   2.62429736e-01 -1.06467863e-02 -9.15456240e-01  1.23513347e-02
  -5.08571424e-01]]
```

```python
single_point = x_test.values[1].reshape(1,-1)
house_predictor.predict(x_test.values[1].reshape(1,-1))
```
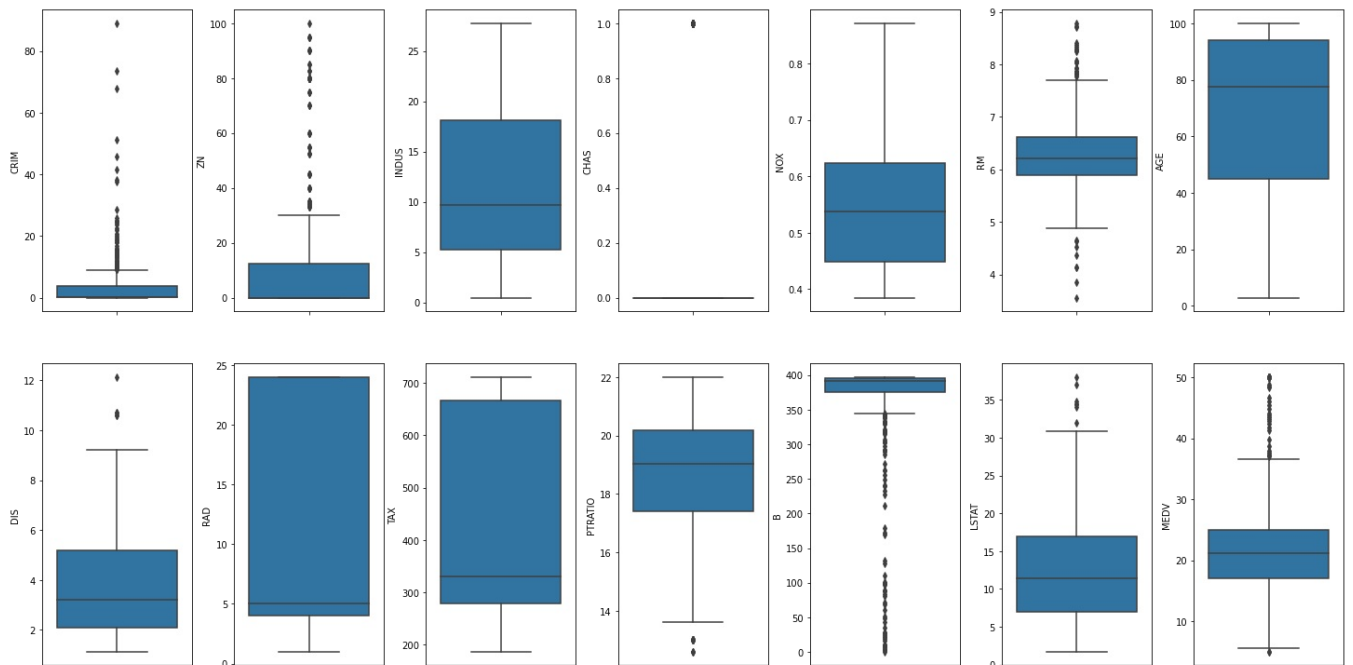
array([[36.02556534]])

```python
y_test.values[1]
```

array([32.4])

```python
import seaborn as sns
import matplotlib.pyplot as plt
from scipy import stats

fig, axs = plt.subplots(ncols=7, nrows=2, figsize=(20, 10))
index = 0
axs = axs.flatten()
for k,v in data.items():
    sns.boxplot(y=k, data=data, ax=axs[index])
    index += 1
plt.tight_layout(pad=0.4, w_pad=0.5, h_pad=5.0)
```

```
In [95]: for k, v in data.items():
             q1 = v.quantile(0.25)
             q3 = v.quantile(0.75)
             irq = q3 - q1
             v_col = v[(v <= q1 - 1.5 * irq) | (v >= q3 + 1.5 * irq)]
             perc = np.shape(v_col)[0] * 100.0 / np.shape(data)[0]
             print("Column %s outliers = %.2f%%" % (k, perc))
```
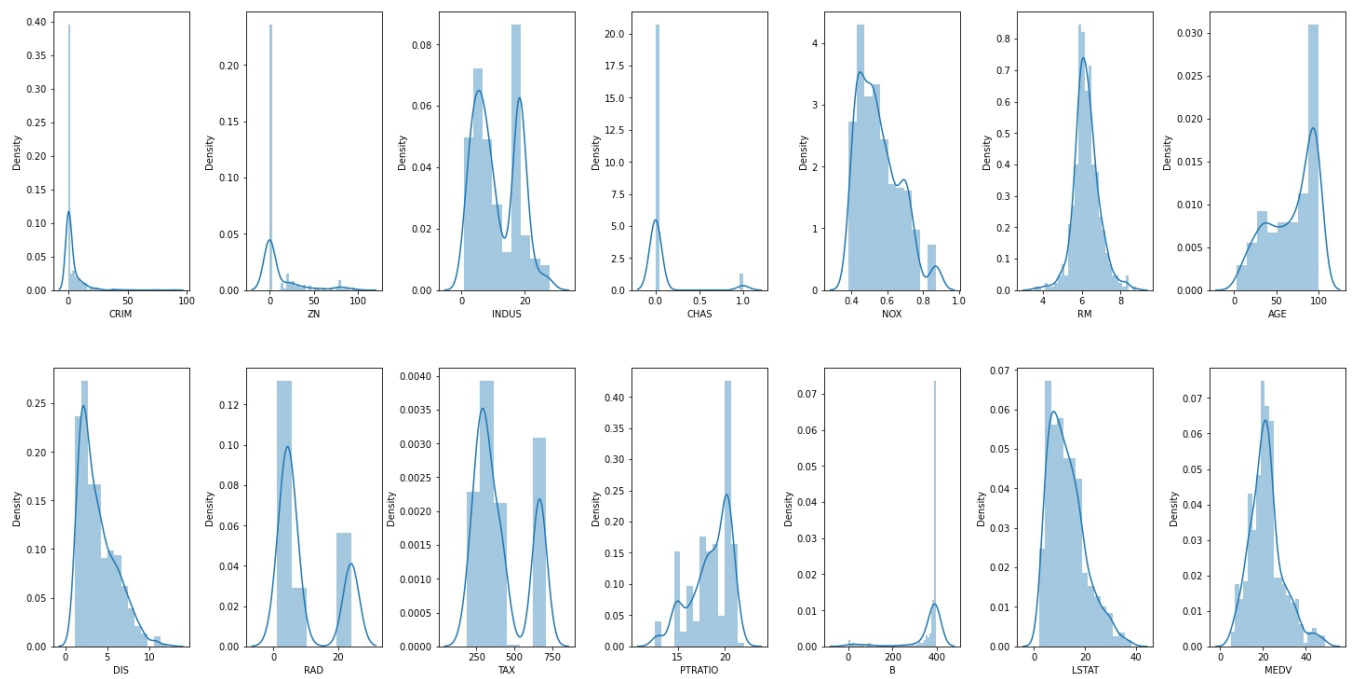
```
Column CRIM outliers = 13.04%
Column ZN outliers = 13.44%
Column INDUS outliers = 0.00%
Column CHAS outliers = 100.00%
Column NOX outliers = 0.00%
Column RM outliers = 5.93%
Column AGE outliers = 0.00%
Column DIS outliers = 0.99%
Column RAD outliers = 0.00%
Column TAX outliers = 0.00%
Column PTRATIO outliers = 2.96%
Column B outliers = 15.22%
Column LSTAT outliers = 1.38%
Column MEDV outliers = 7.91%
```

```
In [96]: data = data[~(data['MEDV'] >= 50.0)]
         print(np.shape(data))


         fig, axs = plt.subplots(ncols=7, nrows=2, figsize=(20, 10))
         index = 0
         axs = axs.flatten()
         for k,v in data.items():
             sns.distplot(v, ax=axs[index])
             index += 1
         plt.tight_layout(pad=0.4, w_pad=0.5, h_pad=5.0)
```
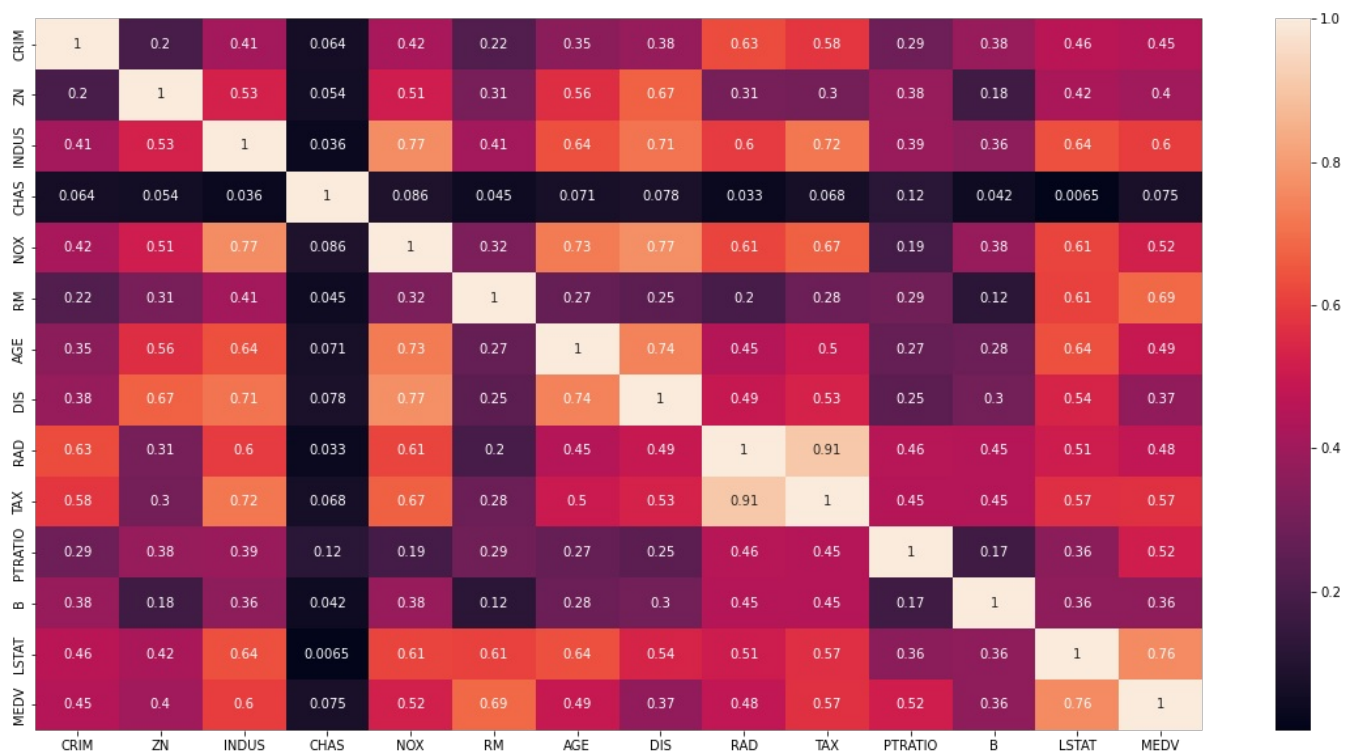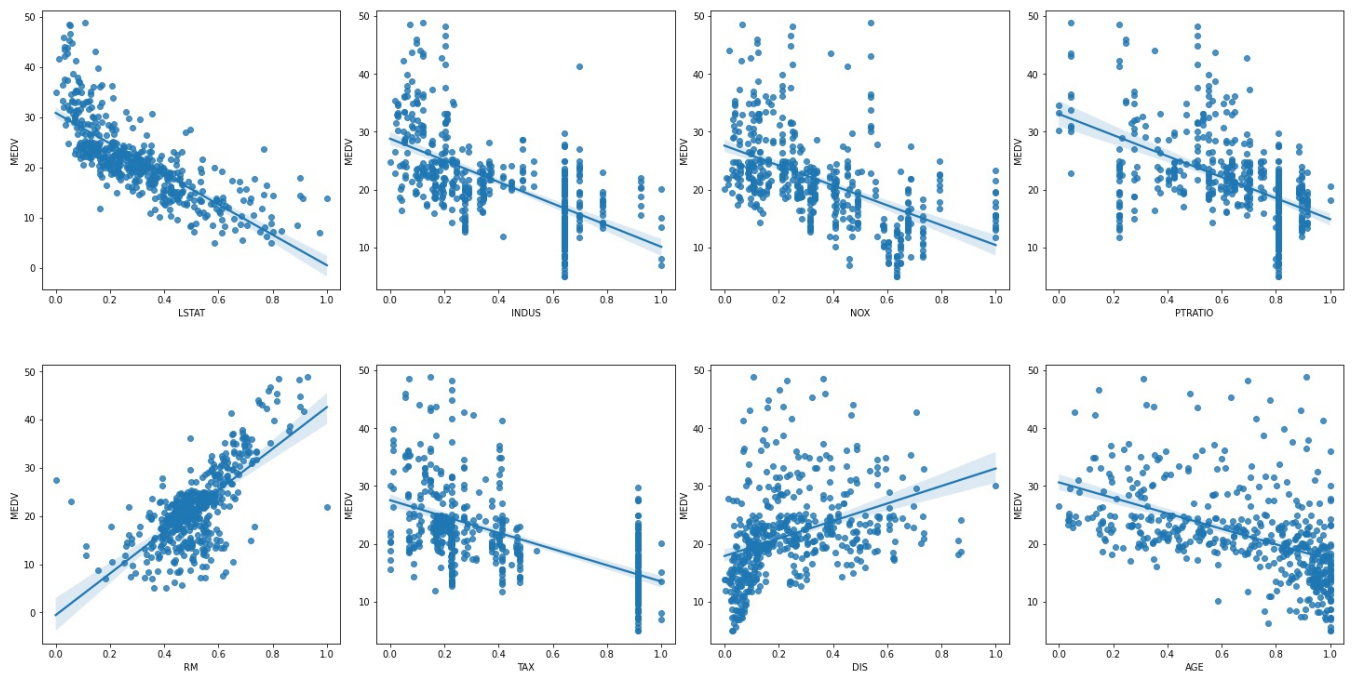
```
(490, 14)
```

```
In [97]: plt.figure(figsize=(20, 10))
         sns.heatmap(data.corr().abs(), annot=True)
```

Out[97]: <AxesSubplot:>

```
In [98]:  from sklearn import preprocessing
          min_max_scaler = preprocessing.MinMaxScaler()
          column_sels = ['LSTAT', 'INDUS', 'NOX', 'PTRATIO', 'RM', 'TAX', 'DIS', 'AGE']
          x = data.loc[:,column_sels]
          y = data['MEDV']
          x = pd.DataFrame(data=min_max_scaler.fit_transform(x), columns=column_sels)
          fig, axs = plt.subplots(ncols=4, nrows=2, figsize=(20, 10))
          index = 0
          axs = axs.flatten()
          for i, k in enumerate(column_sels):
              sns.regplot(y=y, x=x[k], ax=axs[i])
          plt.tight_layout(pad=0.4, w_pad=0.5, h_pad=5.0)
```

```
In [99]:  y = np.log1p(y)
          for col in x.columns:
              if np.abs(x[col].skew()) > 0.3:
                  x[col] = np.log1p(x[col])
```

```
In [100… from sklearn import datasets, linear_model
         from sklearn.model_selection import cross_val_score
         from sklearn.model_selection import KFold
         import numpy as np

         l_regression = linear_model.LinearRegression()
         kf = KFold(n_splits=10)
         min_max_scaler = preprocessing.MinMaxScaler()
         x_scaled = min_max_scaler.fit_transform(x)
         scores = cross_val_score(l_regression, x_scaled, y, cv=kf, scoring='neg_mean_squared_error')
         print("MSE: %0.2f (+/- %0.2f)" % (scores.mean(), scores.std()))

         scores_map = {}
         scores_map['LinearRegression'] = scores
         l_ridge = linear_model.Ridge()
         scores = cross_val_score(l_ridge, x_scaled, y, cv=kf, scoring='neg_mean_squared_error')
         scores_map['Ridge'] = scores
         print("MSE: %0.2f (+/- %0.2f)" % (scores.mean(), scores.std()))

         # Lets try polinomial regression with L2 with degree for the best fit
         from sklearn.pipeline import make_pipeline
         from sklearn.preprocessing import PolynomialFeatures
         #for degree in range(2, 6):
         #    model = make_pipeline(PolynomialFeatures(degree=degree), linear_model.Ridge())
         #    scores = cross_val_score(model, x_scaled, y, cv=kf, scoring='neg_mean_squared_error')
         #    print("MSE: %0.2f (+/- %0.2f)" % (scores.mean(), scores.std()))
         model = make_pipeline(PolynomialFeatures(degree=3), linear_model.Ridge())
         scores = cross_val_score(model, x_scaled, y, cv=kf, scoring='neg_mean_squared_error')
         scores_map['PolyRidge'] = scores
         print("MSE: %0.2f (+/- %0.2f)" % (scores.mean(), scores.std()))
```

```
MSE: -0.04 (+/- 0.04)
MSE: -0.04 (+/- 0.04)
MSE: -0.03 (+/- 0.03)
```

```
In [101… from sklearn.svm import SVR
         from sklearn.model_selection import GridSearchCV

         svr_rbf = SVR(kernel='rbf', C=1e3, gamma=0.1)
         #grid_sv = GridSearchCV(svr_rbf, cv=kf, param_grid={"C": [1e0, 1e1, 1e2, 1e3], "gamma": np.logspace(-2, 2, 5)},
         #grid_sv.fit(x_scaled, y)
         #print("Best classifier :", grid_sv.best_estimator_)
         scores = cross_val_score(svr_rbf, x_scaled, y, cv=kf, scoring='neg_mean_squared_error')
         scores_map['SVR'] = scores
         print("MSE: %0.2f (+/- %0.2f)" % (scores.mean(), scores.std()))
```

```
MSE: -0.04 (+/- 0.03)
```

```
In [102… from sklearn.tree import DecisionTreeRegressor

         desc_tr = DecisionTreeRegressor(max_depth=5)
         #grid_sv = GridSearchCV(desc_tr, cv=kf, param_grid={"max_depth" : [1, 2, 3, 4, 5, 6, 7]}, scoring='neg_mean_squ
         #grid_sv.fit(x_scaled, y)
         #print("Best classifier :", grid_sv.best_estimator_)
         scores = cross_val_score(desc_tr, x_scaled, y, cv=kf, scoring='neg_mean_squared_error')
```

```
scores_map['DecisionTreeRegressor'] = scores
print("MSE: %0.2f (+/- %0.2f)" % (scores.mean(), scores.std()))
```

MSE: -0.05 (+/- 0.04)

In [103...
```python
from sklearn.neighbors import KNeighborsRegressor

knn = KNeighborsRegressor(n_neighbors=7)
scores = cross_val_score(knn, x_scaled, y, cv=kf, scoring='neg_mean_squared_error')
scores_map['KNeighborsRegressor'] = scores
#grid_sv = GridSearchCV(knn, cv=kf, param_grid={"n_neighbors" : [2, 3, 4, 5, 6, 7]}, scoring='neg_mean_squared_
#grid_sv.fit(x_scaled, y)
#print("Best classifier :", grid_sv.best_estimator_)
print("KNN Accuracy: %0.2f (+/- %0.2f)" % (scores.mean(), scores.std()))
```
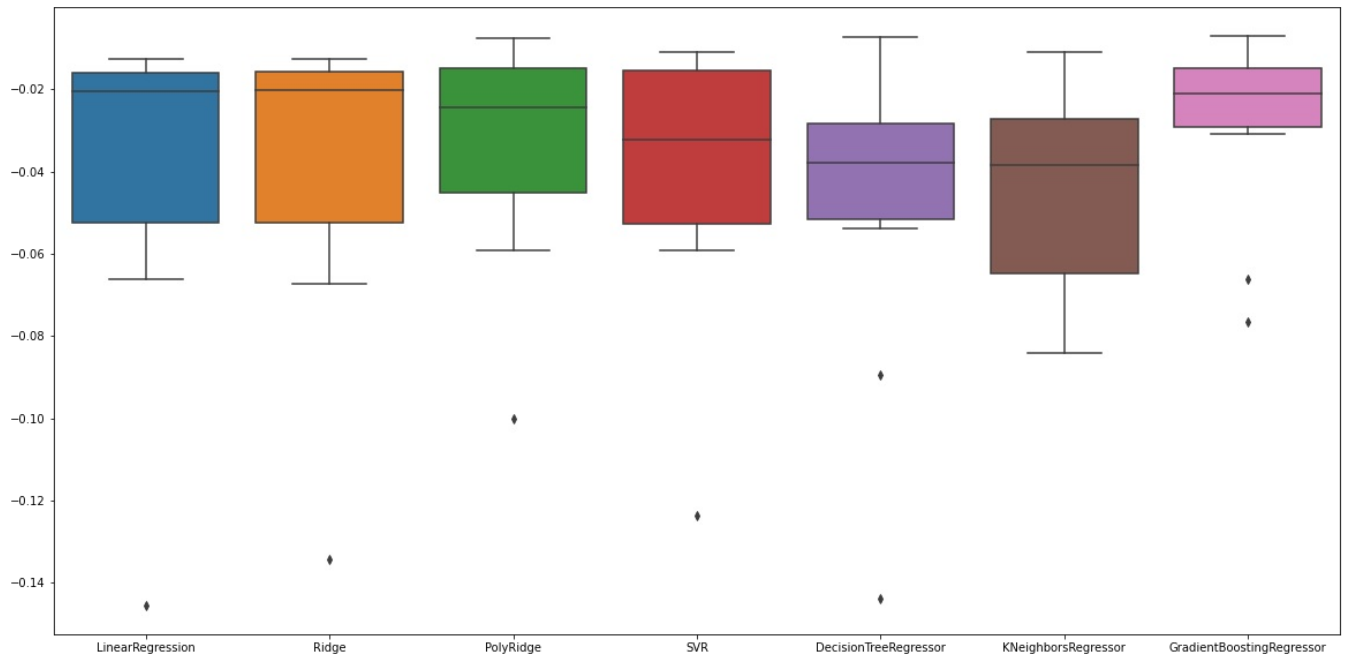
KNN Accuracy: -0.04 (+/- 0.02)

In [104...
```python
from sklearn.ensemble import GradientBoostingRegressor

gbr = GradientBoostingRegressor(alpha=0.9,learning_rate=0.05, max_depth=2, min_samples_leaf=5, min_samples_spli
#param_grid={'n_estimators':[100, 200], 'learning_rate': [0.1,0.05,0.02], 'max_depth':[2, 4,6], 'min_samples_le
#grid_sv = GridSearchCV(gbr, cv=kf, param_grid=param_grid, scoring='neg_mean_squared_error')
#grid_sv.fit(x_scaled, y)
#print("Best classifier :", grid_sv.best_estimator_)
scores = cross_val_score(gbr, x_scaled, y, cv=kf, scoring='neg_mean_squared_error')
scores_map['GradientBoostingRegressor'] = scores
print("MSE: %0.2f (+/- %0.2f)" % (scores.mean(), scores.std()))
```

MSE: -0.03 (+/- 0.02)

In [105...
```python
plt.figure(figsize=(20, 10))
scores_map = pd.DataFrame(scores_map)
sns.boxplot(data=scores_map)
```

Out[105]:    <AxesSubplot:>



In [106...
*#The models SVR and GradientBoostingRegressor show better performance with -11.62 (+/- 5.91) and -12.39 (+/- 5.*

In [ ]: