```
import pandas as pd
import numpy as np
df=pd.read_csv('/content/boston.csv')
```

```
df.head()
```

|   | 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 |

```
df.columns
```

```
Index(['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX',
       'PTRATIO', 'B', 'LSTAT', 'MEDV'],
      dtype='object')
```

```
df.isnull().sum()
```

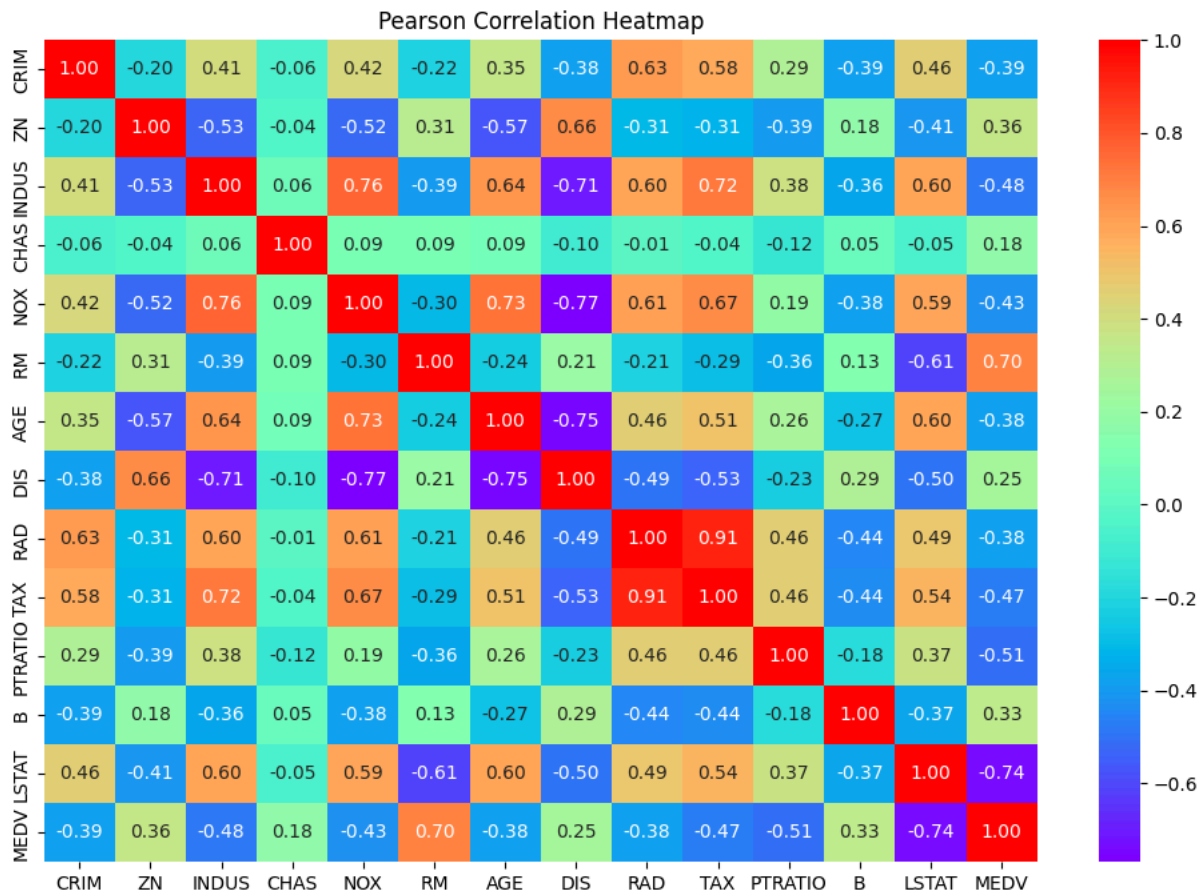|         | 0 |
|---------|---|
| CRIM    | 0 |
| ZN      | 0 |
| INDUS   | 0 |
| CHAS    | 0 |
| NOX     | 0 |
| RM      | 0 |
| AGE     | 0 |
| DIS     | 0 |
| RAD     | 0 |
| TAX     | 0 |
| PTRATIO | 0 |
| B       | 0 |
| LSTAT   | 0 |
| MEDV    | 0 |

**dtype:** int64

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 506 entries, 0 to 505
Data columns (total 14 columns):
 #   Column   Non-Null Count  Dtype
---  ------   --------------  -----
 0   CRIM     506 non-null    float64
 1   ZN       506 non-null    float64
 2   INDUS    506 non-null    float64
 3   CHAS     506 non-null    int64
 4   NOX      506 non-null    float64
 5   RM       506 non-null    float64
 6   AGE      506 non-null    float64
 7   DIS      506 non-null    float64
 8   RAD      506 non-null    int64
 9   TAX      506 non-null    float64
 10  PTRATIO  506 non-null    float64
 11  B        506 non-null    float64
 12  LSTAT    506 non-null    float64
 13  MEDV     506 non-null    float64
dtypes: float64(12), int64(2)
memory usage: 55.5 KB
```

```
import seaborn as sns
import matplotlib.pyplot as plt
```

```
corr_matrix = df.corr() #correlation matrix

plt.figure(figsize=(12,8)) #heatmap
sns.heatmap(corr_matrix, annot=True, cmap='rainbow', fmt=".2f")
plt.title("Pearson Correlation Heatmap")
plt.show()
```

## Pearson Correlation Heatmap

|        | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT | MEDV |
|--------|------|------|-------|------|------|------|------|------|------|------|---------|------|-------|------|
| CRIM   | 1.00 | -0.20 | 0.41 | -0.06 | 0.42 | -0.22 | 0.35 | -0.38 | 0.63 | 0.58 | 0.29 | -0.39 | 0.46 | -0.39 |
| ZN     | -0.20 | 1.00 | -0.53 | -0.04 | -0.52 | 0.31 | -0.57 | 0.66 | -0.31 | -0.31 | -0.39 | 0.18 | -0.41 | 0.36 |
| INDUS  | 0.41 | -0.53 | 1.00 | 0.06 | 0.76 | -0.39 | 0.64 | -0.71 | 0.60 | 0.72 | 0.38 | -0.36 | 0.60 | -0.48 |
| CHAS   | -0.06 | -0.04 | 0.06 | 1.00 | 0.09 | 0.09 | 0.09 | -0.10 | -0.01 | -0.04 | -0.12 | 0.05 | -0.05 | 0.18 |
| NOX    | 0.42 | -0.52 | 0.76 | 0.09 | 1.00 | -0.30 | 0.73 | -0.77 | 0.61 | 0.67 | 0.19 | -0.38 | 0.59 | -0.43 |
| RM     | -0.22 | 0.31 | -0.39 | 0.09 | -0.30 | 1.00 | -0.24 | 0.21 | -0.21 | -0.29 | -0.36 | 0.13 | -0.61 | 0.70 |
| AGE    | 0.35 | -0.57 | 0.64 | 0.09 | 0.73 | -0.24 | 1.00 | -0.75 | 0.46 | 0.51 | 0.26 | -0.27 | 0.60 | -0.38 |
| DIS    | -0.38 | 0.66 | -0.71 | -0.10 | -0.77 | 0.21 | -0.75 | 1.00 | -0.49 | -0.53 | -0.23 | 0.29 | -0.50 | 0.25 |
| RAD    | 0.63 | -0.31 | 0.60 | -0.01 | 0.61 | -0.21 | 0.46 | -0.49 | 1.00 | 0.91 | 0.46 | -0.44 | 0.49 | -0.38 |
| TAX    | 0.58 | -0.31 | 0.72 | -0.04 | 0.67 | -0.29 | 0.51 | -0.53 | 0.91 | 1.00 | 0.46 | -0.44 | 0.54 | -0.47 |
| PTRATIO | 0.29 | -0.39 | 0.38 | -0.12 | 0.19 | -0.36 | 0.26 | -0.23 | 0.46 | 0.46 | 1.00 | -0.18 | 0.37 | -0.51 |
| B      | -0.39 | 0.18 | -0.36 | 0.05 | -0.38 | 0.13 | -0.27 | 0.29 | -0.44 | -0.44 | -0.18 | 1.00 | -0.37 | 0.33 |
| LSTAT  | 0.46 | -0.41 | 0.60 | -0.05 | 0.59 | -0.61 | 0.60 | -0.50 | 0.49 | 0.54 | 0.37 | -0.37 | 1.00 | -0.74 |
| MEDV   | -0.39 | 0.36 | -0.48 | 0.18 | -0.43 | 0.70 | -0.38 | 0.25 | -0.38 | -0.47 | -0.51 | 0.33 | -0.74 | 1.00 |

```
corr_with_target = corr_matrix['MEDV'].sort_values(ascending=False)
print(corr_with_target)
```

```
MEDV       1.000000
RM         0.695360
ZN         0.360445
B          0.333461
DIS        0.249929
CHAS       0.175260
AGE       -0.376955
RAD       -0.381626
CRIM      -0.388305
NOX       -0.427321
TAX       -0.468536
INDUS     -0.483725
PTRATIO   -0.507787
LSTAT     -0.737663
Name: MEDV, dtype: float64
```

```
high_corr_features=corr_with_target[abs(corr_with_target) > 0.5]
print(high_corr_features)
```

```
MEDV       1.000000
RM         0.695360
PTRATIO   -0.507787
LSTAT     -0.737663
Name: MEDV, dtype: float64
```

```
def remove_outiers(df, columns):
  for col in columns:
    Q1 = df[col].quantile(0.25)
    Q3 = df[col].quantile(0.75)
    IQR = Q3 - Q1
    lower = Q1 - 1.5 * IQR
    upper = Q3 + 1.5 * IQR
    df = df[(df[col] >= lower) & (df[col] <= upper)]
  return df
```
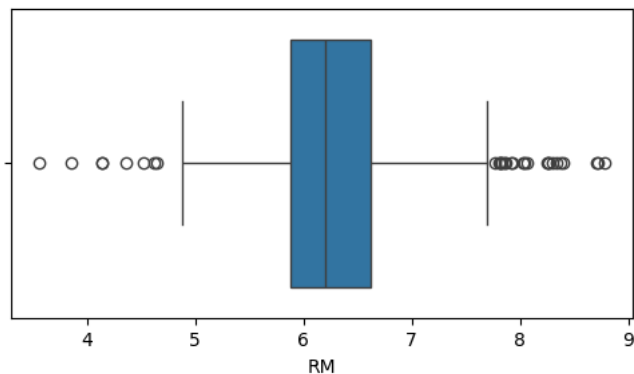
```
important_cols = high_corr_features.index.tolist()
important_cols.remove('MEDV')

df_clean = remove_outiers(df, important_cols)
```
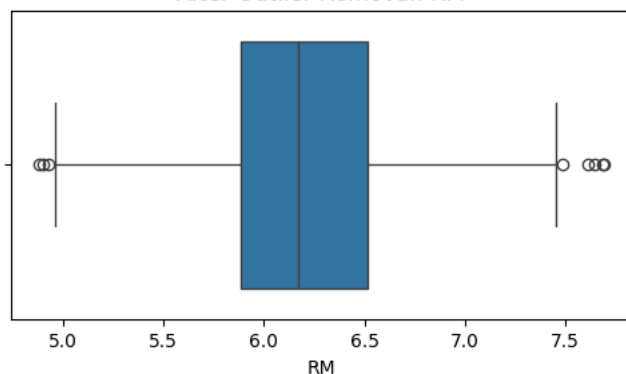
```
for col in important_cols:
  plt.figure(figsize=(6,3))
  sns.boxplot(x=df[col])
  plt.title(f"Before Outlier Removal: {col}")
  plt.show()

  plt.figure(figsize=(6,3))
  sns.boxplot(x=df_clean[col])
  plt.title(f"After Outlier Removal: {col}")
  plt.show()
```
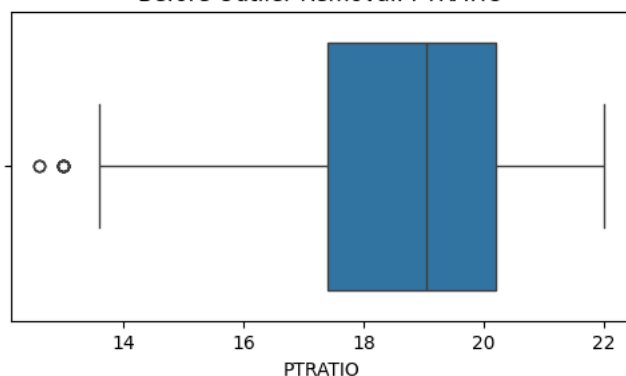
### Before Outlier Removal: RM

### After Outlier Removal: RM

### Before Outlier Removal: PTRATIO

### After Outlier Removal: PTRATIO

```
x = df_clean('MEDV', axis=1)
y = df_clean['MEDV']
```

```
---------------------------------------------------------------------------
NameError                                 Traceback (most recent call last)
/tmp/ipython-input-460421389.py in <cell line: 0>()
----> 1 x = df_clean('MEDV', axis=1)
      2 y = df_clean['MEDV']

NameError: name 'df_clean' is not defined
```

Start coding or generate with AI.



After Outlier Removal: LSTAT