

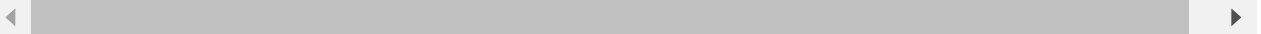
```
In [1]: import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns
%matplotlib inline
import warnings
warnings.filterwarnings('ignore')
```

```
In [2]: from sklearn.datasets import load_boston
boston_dataset = load_boston()
```

```
In [3]: boston = pd.DataFrame(boston_dataset.data, columns=boston_dataset.feature_names)
boston['MEDV'] = boston_dataset.target
boston.head()
```

```
Out[3]:
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	MED
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	396.90	4.98	24.
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14	21.
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03	34.
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94	33.
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33	36.



```
In [4]: boston.shape
```

```
Out[4]: (506, 14)
```

```
In [5]: boston.info() # All variables are numerical, there's no need for dummy encoding
```

```
<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    float64
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    float64
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(14)
memory usage: 55.5 KB

In [6]: `boston.describe()`

Out[6]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS
count	506.000000	506.000000	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	68.574901	3.795043
std	8.601545	23.322453	6.860353	0.253994	0.115878	0.702617	28.148861	2.105710
min	0.006320	0.000000	0.460000	0.000000	0.385000	3.561000	2.900000	1.129600
25%	0.082045	0.000000	5.190000	0.000000	0.449000	5.885500	45.025000	2.100175
50%	0.256510	0.000000	9.690000	0.000000	0.538000	6.208500	77.500000	3.207450
75%	3.677083	12.500000	18.100000	0.000000	0.624000	6.623500	94.075000	5.188425
max	88.976200	100.000000	27.740000	1.000000	0.871000	8.780000	100.000000	12.126500

In [8]: `boston.MEDV.describe()`

Out[8]:

```
count    506.000000
mean      22.532806
std        9.197104
min         5.000000
25%       17.025000
50%       21.200000
75%       25.000000
max       50.000000
Name: MEDV, dtype: float64
```

In [9]: `boston.isnull().sum()` *# ALL data are complete, there's no need for removing or filling*

Out[9]:

```
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
```

In [10]:

```
sns.set(rc={'figure.figsize':(11.7,8.27)})
correlation_matrix = boston.corr().round(2)
# annot = True to print the values inside the square
sns.heatmap(data=correlation_matrix, annot=True)
# There's indeed obvious correlation between the data
```

Out[10]: <AxesSubplot:>



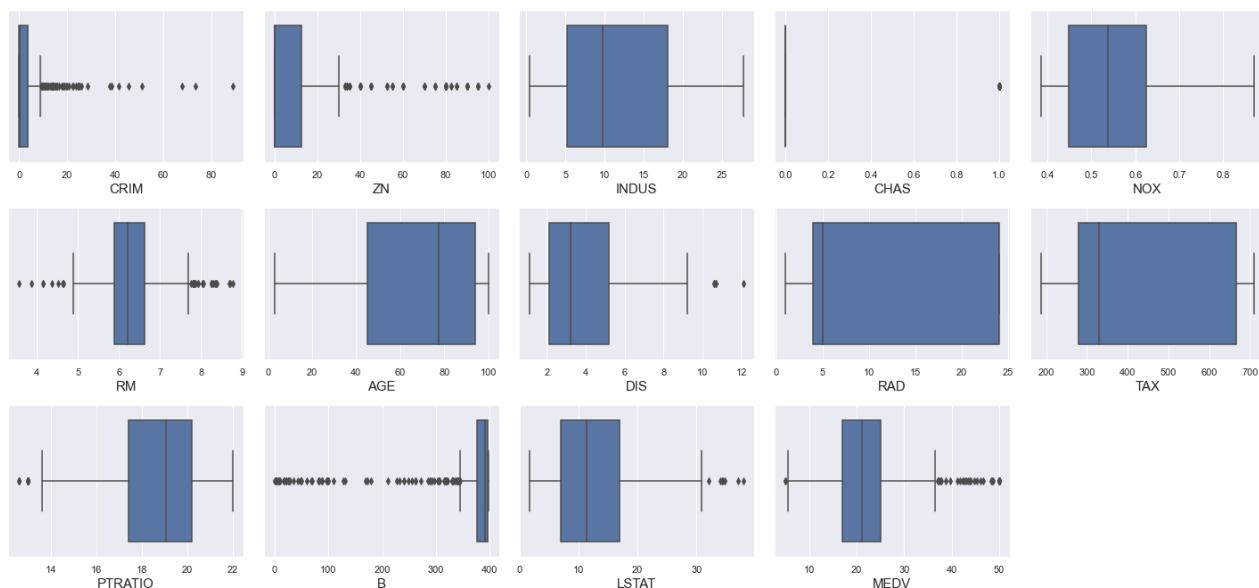
```

In [11]: plt.figure(figsize = (20, 15))
plotnumber = 1

for col in boston.columns:
    if plotnumber <= 30:
        ax = plt.subplot(5, 5, plotnumber)
        sns.boxplot(boston[col])
        plt.xlabel(col, fontsize = 15)

        plotnumber += 1
plt.tight_layout()
plt.show()
#There's indeed outlier in the dataset

```

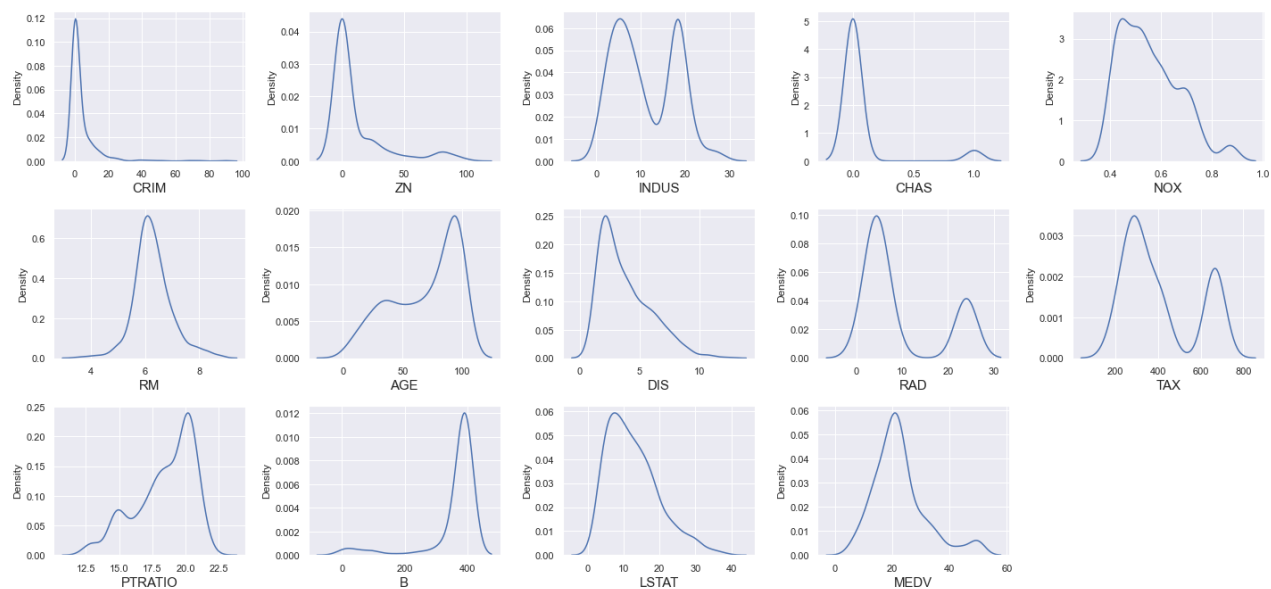


In [12]:

```
plt.figure(figsize = (20, 15))
plotnumber = 1

for col in boston.columns:
    if plotnumber <= 30:
        ax = plt.subplot(5, 5, plotnumber)
        sns.distplot(boston[col], rug=False, hist=False)
        plt.xlabel(col, fontsize = 15)

        plotnumber += 1
plt.tight_layout()
plt.show()
#There's indeed some sparsity in the data
```



In [13]:

```
boston.to_csv('Boston.csv')
```