

## Multicollinearity of Features - Lab

### Introduction

In this lab you'll identify multicollinearity in the Boston Housing Data set.

### Objectives

You will be able to:

- Plot heatmaps for the predictors of the Boston dataset
- Understand and calculate correlation matrices

### Correlation matrix for the Boston Housing data

Let's reimport the Boston Housing data and let's use the data with the categorical variables for `tax_dummy` and `rad_dummy`

```
In [5]: import pandas as pd
from sklearn.datasets import load_boston
boston = load_boston()

boston_features = pd.DataFrame(boston.data, columns = boston.feature_names)

# first, create bins for based on the values observed. 5 values will result in 4 bins
bins = [0, 3, 4, 5, 24]
bins_rad = pd.cut(boston_features['RAD'], bins)
bins_rad = bins_rad.cat.as_unordered()

# first, create bins for based on the values observed. 5 values will result in 4 bins
bins = [0, 250, 300, 360, 460, 712]
bins_tax = pd.cut(boston_features['TAX'], bins)
bins_tax = bins_tax.cat.as_unordered()

tax_dummy = pd.get_dummies(bins_tax, prefix="TAX")
rad_dummy = pd.get_dummies(bins_rad, prefix="RAD")
boston_features = boston_features.drop(["RAD", "TAX"], axis=1)
boston_features = pd.concat([boston_features, rad_dummy, tax_dummy], axis=1)
```

```
In [18]: boston_features.head()
```

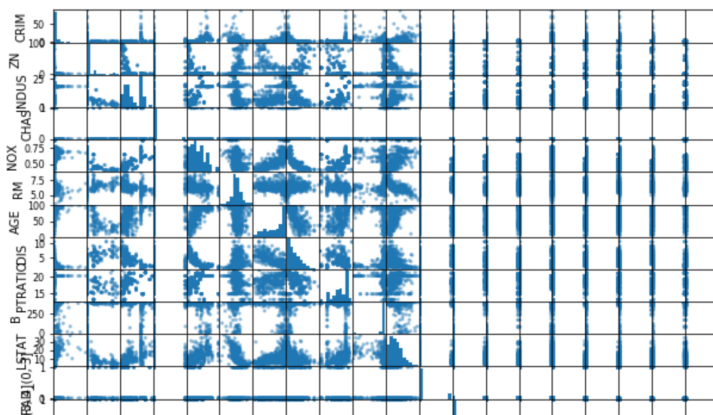
```
Out[18]:
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	PTRATIO	B	LSTAT	RAD_(0, 3]	RAD_(3, 4]	RAD_(4, 5]	RAD_(5, 24]	TAX_(0, 250]	TAX_(250, 300]	TAX_(300, 360]
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	15.3	396.90	4.98	1	0	0	0	0	1	0
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	17.8	396.90	9.14	1	0	0	0	1	0	0
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	17.8	392.83	4.03	1	0	0	0	1	0	0
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	18.7	394.63	2.94	1	0	0	0	1	0	0
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	18.7	396.90	5.33	1	0	0	0	1	0	0

### Scatter matrix

Create the scatter matrix for the Boston Housing data.

```
In [19]: pd.plotting.scatter_matrix(boston_features, figsize = [10, 10]);
```





<b>ZN</b>	False	True	False	False	False	False	False	False	False	False	False
<b>INDUS</b>	False	False	True	False	True	False	False	False	False	False	False
<b>CHAS</b>	False	False	False	True	False	False	False	False	False	False	False
<b>NOX</b>	False	False	True	False	True	False	False	True	False	False	False
<b>RM</b>	False	False	False	False	False	True	False	False	False	False	False
<b>AGE</b>	False	False	False	False	False	False	True	False	False	False	False
<b>DIS</b>	False	False	False	False	True	False	False	True	False	False	False
<b>PTRATIO</b>	False	False	False	False	False	False	False	False	True	False	False
<b>B</b>	False	False	False	False	False	False	False	False	False	True	False
<b>LSTAT</b>	False	False	False	False	False	False	False	False	False	False	True

Remove the most problematic feature from the data.

```
In [31]: boston_features = boston_features.drop("NOX",axis=1)
```

### Summary

Good job! You've now edited the Boston Housing Data so highly correlated variables are removed.