# Project Case Scenario

**Project Scenario:** You are a Data Scientist with a housing agency in Boston MA, you have been given access to a previous dataset on housing prices derived from the U.S. Census Service to present insights to higher management. Based on your experience in Statistics, what information can you provide them to help with making an informed decision? Upper management will like to get some insight into the following.

- Is there a significant difference in the median value of houses bounded by the Charles river or not?
- Is there a difference in median values of houses of each proportion of owner-occupied units built before 1940?
- Can we conclude that there is no relationship between Nitric oxide concentrations and the proportion of non-retail business acres per town?
- What is the impact of an additional weighted distance to the five Boston employment centres on the median value of owner-occupied homes?

Using the appropriate graphs and charts, generate basic statistics and visualizations that you think will be useful for the upper management to give them important insight given the question they are asking, in your graphs, include an explanation of each statistic.

Details regarding the project will be broken down in the next reading sections.

# Overview of Project Tasks

**Final Project:** For the project scenario, here is an overview of your tasks. The following reading sections will provide you with detailed instructions for each task.

**Project Tasks:**

**Task 1:** Familiarize yourself with the dataset

**Task 2:** (Optional) If you do not already have an instance of Watson Studio, create an IBM Cloud Lite account and provision an instance of Waston Studio.

**Task 3:** Load the dataset in a Jupyter Notebook using Watson Studio.

**Task 4:** Generate basic statistics and visualizations for upper management.

**Task 5:** Use the appropriate tests to answer the questions provided.

**Task 6**: Share your Jupyter Notebook.

 This project is worth 15% of your final grade. Detailed instructions for each of these tasks follow.

总结:

# Task 1: Become familiar with the dataset

The following describes the dataset variables:

· CRIM - per capita crime rate by town

· ZN - proportion of residential land zoned for lots over 25,000 sq.ft.

· INDUS - proportion of non-retail business acres per town.

· CHAS - Charles River dummy variable (1 if tract bounds river; 0 otherwise)

· NOX - nitric oxides concentration (parts per 10 million)

· RM - average number of rooms per dwelling

· AGE - proportion of owner-occupied units built prior to 1940

· DIS - weighted distances to five Boston employment centres

· RAD - index of accessibility to radial highways

· TAX - full-value property-tax rate per $10,000

· PTRATIO - pupil-teacher ratio by town

· LSTAT - % lower status of the population

· MEDV - Median value of owner-occupied homes in $1000's

*This data was modified for this course and the link to the complete dataset can be found in the sklearn.datasets library*

## sklearn.datasets.load_boston

sklearn.datasets. **load_boston**(*, *return_X_y=False*)

Load and return the boston house-prices dataset (regression).

| Samples total | 506 |
|---|---|
| Dimensionality | 13 |
| Features | real, positive |
| Targets | real 5. - 50. |

Read more in the User Guide.

总结:

# Task 2: Create or Login into IBM cloud to use Watson Studio.

For this project, you will be required to use Watson studio and at the end of the project, you will share a link to your Jupyter notebook in your Watson Studio project for a peer review.

If you already have an IBM Cloud account and an instance of Watson Studio with a Jupyter Notebook Project, please skip ahead to the next task.

Step1: If you don't have IBM cloud, please use the instructions found here to create an account and use Watson Studio. If you already have an account, sign-in here

Step 2: Use the instructions found here to create a Project in Watson studio, and create a Jupyter Notebook that you will utilize in the following tasks.

# Task 3: Load in the Dataset in your Jupyter Notebook

In the Jupyter notebook you created in the previous task, add a code cell, and copy the contents below to load the Boston housing dataset in the notebook:

```
1    boston_url = 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsN
2    boston_df=pd.read_csv(boston_url)
```

boston_url = 'https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-ST0151EN-SkillsNetwork/labs/boston_housing.csv'

boston_df=pd.read_csv(boston_url)

# Task 4: Generate Descriptive Statistics and Visualizations

For all visualizations, please include a title in each graph and appropriate labels

Generate the following and explain your findings:

- For the "Median value of owner-occupied homes" provide a **boxplot**

- Provide a  **histogram** for the Charles river variable

- Provide a **boxplot** for the MEDV variable vs the AGE variable. (Discretize the age variable into three groups of 35 years and younger, between 35 and 70 years and 70 years and older)

- Provide a **scatter plot** to show the relationship between Nitric oxide concentrations and the proportion of non-retail business acres per town. What can you say about the relationship?

- Create a **histogram** for the pupil to teacher ratio variable

**总结:**

# Task 5: Use the appropriate tests to answer the questions provided.

For each of the following questions;

- Is there a significant difference in median value of houses bounded by the Charles river or not? (T-test for independent samples)

- Is there a difference in Median values of houses (MEDV) for each proportion of owner occupied units built prior to 1940 (AGE)? (ANOVA)

- Can we conclude that there is no relationship between Nitric oxide concentrations and proportion of non-retail business acres per town? (Pearson Correlation)

- What is the impact of an additional weighted distance to the five Boston employment centres on the median value of owner occupied homes? (Regression analysis)

Be sure to:

1. State your hypothesis.

2. Use $\alpha = 0.05$

3. Perform the test Statistics.

4. State the conclusion from the test.

# Task 6: Share your Jupyter Notebook.

Follow the instructions here to share your notebook.