# ▾ Problem 5: Simple Linear Regression

ESE402 HW4

Sheil Sarda <sheils@seas.upenn.edu>

In this question, you will implement simple linear regression from scratch. The dataset you will work with is called the Boston data set. You can find more information about the data set here:
https://www.cs.toronto.edu/~delve/data/boston/bostonDetail.html

You will use the pandas library to load the csv file into a dataframe:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 %matplotlib inline
```

```
1 # read the csv file and load into a pandas dataframe
2 # make sure Boston.csv is in the same file path as this notebook
3 boston = pd.read_csv('Boston.csv')
```

```
1 # read the above link to learn more about what each of the columns indicate
2 boston.head()
```

|   | crim | zn | indus | chas | nox | rm | age | dis | rad | tax | ptratio | lstat |
|---|------|-----|-------|------|-------|-------|------|--------|-----|-----|---------|-------|
| 0 | 0.00632 | 18.0 | 2.31 | 0 | 0.538 | 6.575 | 65.2 | 4.0900 | 1 | 296 | 15.3 | 4.98 |
| 1 | 0.02731 | 0.0 | 7.07 | 0 | 0.469 | 6.421 | 78.9 | 4.9671 | 2 | 242 | 17.8 | 9.14 |
| 2 | 0.02729 | 0.0 | 7.07 | 0 | 0.469 | 7.185 | 61.1 | 4.9671 | 2 | 242 | 17.8 | 4.03 |
| 3 | 0.03237 | 0.0 | 2.18 | 0 | 0.458 | 6.998 | 45.8 | 6.0622 | 3 | 222 | 18.7 | 2.94 |
| 4 | 0.06905 | 0.0 | 2.18 | 0 | 0.458 | 7.147 | 54.2 | 6.0622 | 3 | 222 | 18.7 | 5.33 |

Simple linear regression builds a linear relationship between an input variable $X$ and an output variable $Y$. We can define this linear relationship as follows:

$$Y = \beta_0 + \beta_1 X$$

Objective: find the linear relationship between the proportion of non-retail business acres per town (indus) and the full-value property-tax rate per 10,000 dollars (tax)

So our equation will look like:

$$TAX = \beta_0 + \beta_1 INDUS$$

Here, the coefficient $\beta_0$ is the intercept, and $\beta_1$ is the scale factor or slope. How do we determine the values of these coefficients?

There are several different methods to do so, but we will focus on the Ordinary Least Squares (OLS) method. This method minimizes the sum of the squares of the differences between the observed dependent variable and those predicted by the linear function.

Recall that a residual is the difference between any data point and the line of regression. When we develop a regression model, we want the sum of the residuals squared to be minimized, indicating that the model is a close fit to the data.
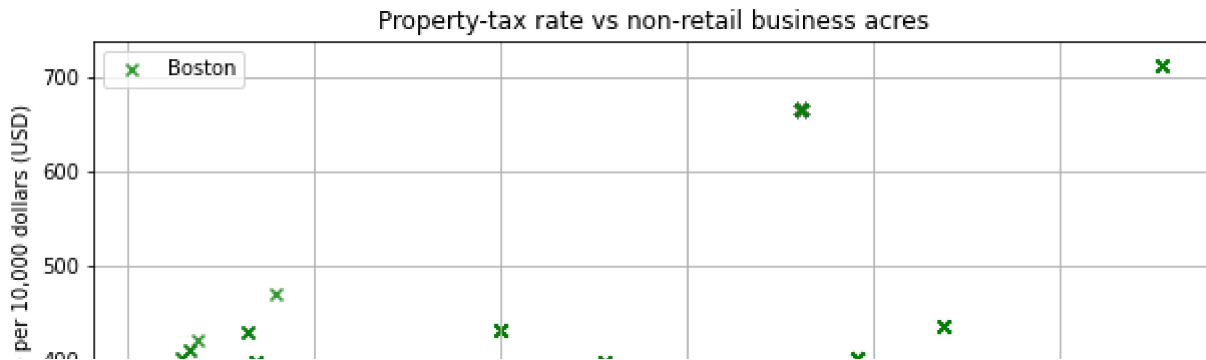
$$RSS = \sum_{i=1}^{n}(y_i - f(x_i))^2$$
$$= \sum_{i=1}^{n}(y_i - (\beta_0 + \beta_1 x_i))^2$$

This is the objective function we minimize to find $\beta_0$ and $\beta_1$.

```
1 # set X to 'indus' and y to 'tax'
2 X = boston['indus']
3 y = boston['tax']
```

First, visualize the data by plotting X and y using matplotlib. Be sure to include a title and axis labels.

```
1 # TODO: display plot
2 plt.rcParams["figure.figsize"] = (10,5)
3 plt.scatter(X, y, label="Boston", marker='x', alpha=0.8, c='Green')
4 # TODO: labels and title
5 plt.title("Property-tax rate vs non-retail business acres")
6 plt.xlabel("Proportion of non-retail business acres per town (%)")
7 plt.ylabel("Property-tax rate per 10,000 dollars (USD)")
8 plt.grid()
9 plt.legend()
10 plt.savefig("boston_indus_tax.png")
11 plt.show()
```

## Property-tax rate vs non-retail business acres



▼ TODO: What do you notice about the relationship between the variables?

A: In general, there is a positive correlation between the proportion of non-retail business acres per town and business taxes. A potential driver for this relationship could be that as the number of businesses in a town increase, the town's need to provide lower taxes as an incentive for companies to operate there and provide employment to the region declines.

Next, find the coefficients. The values for $\beta_0$ and $\beta_1$ are given by the following equations, where $n$ is the total number of values. This derivation was done in class.

$$\beta_1 = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2}$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

```
1 # TODO: implement function
2 def get_coeffs(X, y):
3     """
4     Params:
5         X: the X vector
6         y: the y vector
7     Returns:
8         a tuple (b1, b0)
9     """
10     X_mean = sum(X) / float(len(X))
11     y_mean = sum(y) / float(len(y))
12     X_var = sum((X - X_mean)**2)/float(len(X) - 1)
13
14     covariance = 0.0
15     for i in range(len(X)):
16         covariance += (X[i] - X_mean) * (y[i] - y_mean)
17     covariance = covariance / float(len(X) - 1)
18
19     b1 = covariance / X_var
20     b0 = y_mean - b1 * X_mean
21     return (b1, b0)
```

```
1 # run cell to call function and display the regression line
2 # the values are rounded for display convenience
3 b1, b0 = get_coeffs(X, y)
4 print("Regression line: TAX = " + str(round(b0)) + " + " + str(round(b1)) +"*INDUS")
```
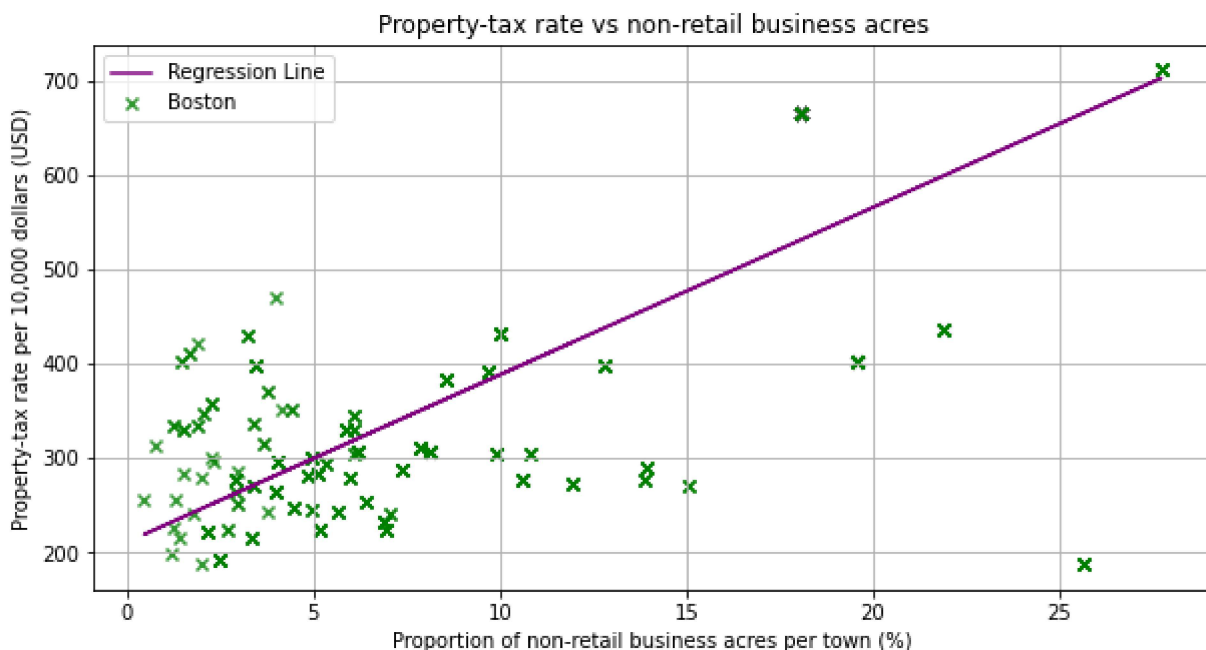
```
Regression line: TAX = 211 + 18*INDUS
```

Plot the regression line overlayed on the real y-values.

```
1 # TODO: plot y-values
2 plt.rcParams["figure.figsize"] = (10,5)
3 plt.scatter(X, y, label="Boston", marker='x', alpha=0.8, c='Green')
4
5 # TODO: plot regression line
6 y_predicted = X*b1 + b0
7 plt.plot(X, y_predicted, label="Regression Line", c='Purple')
8
9 # TODO: labels and title
10 plt.title("Property-tax rate vs non-retail business acres")
11 plt.xlabel("Proportion of non-retail business acres per town (%)")
12 plt.ylabel("Property-tax rate per 10,000 dollars (USD)")
13 plt.grid()
14 plt.legend()
15 plt.show()
```



The line appears to fit the data, but first, let us find the RSS to evaluate this model. The RSS is used to measure the amount of variance in the data set that is not explained by the regression model. Recall that

$$RSS = \sum_{i=1}^{n}(y_i - (\beta_0 + \beta_1 x_i))^2$$

```
1 # TODO: implement function
2 def get_RSS(b0, b1, X, y):
3     """
4     Params:
5         b0: beta 0
6         b1: beta 1
7         X: X vector
8         y: y vector
9     Returns:
10        residual sum of squares (RSS)
11    """
12    rss = 0.0
13    for i in range(len(y)):
14        rss += (y[i] - (b0 + b1*X[i]))**2
15    return rss
```

```
1 # run this cell to print RSS
2 print("RSS:", get_RSS(b0, b1, X, y))
```

```
RSS: 6892554.224031512
```

We can also evaluate the model through the Root Mean Squared Error (RMSE) and the Coefficient of Determination ($R^2$ score).

- The RMSE is similar to the RSS, but provides a value with more interpretable units -- in our case, tax rate per 10,000 dollars.
- The $R^2$ value represents the proportion of the variance for the dependent variable that is explained by the independent variable.

Use the following equations to find the RMSE and $R^2$ score:

$$RMSE = \sqrt{(\sum_{i=1}^{n}\frac{1}{n}(\hat{y}_i - y_i)^2)}$$

$$R^2 = 1 - \frac{SS_r}{SS_t}$$

$$SS_t = \sum_{i=1}^{n}(y_i - \bar{y})^2$$

$$SS_r = \sum_{i=1}^{n}(y_i - \hat{y}_i)^2$$

```python
1 # TODO: implement function
2 def get_RMSE(b0, b1, X, y):
3     """
4     Params:
5         b0: beta 0
6         b1: beta 1
7         X: X vector
8         y: y vector
9     Returns:
10        rmse
11    """
12    rss = get_RSS(b0, b1, X, y)
13    rmse = (rss/len(y))**0.5
14    return rmse
```

```python
1 # run cell to print RMSE
2 print("RMSE: ", get_RMSE(b0, b1, X, y))
```

```
RMSE:  116.7118188706435
```

```python
1 # TODO: implement function
2 def get_R2(b0, b1, X, y):
3     """
4     Params:
5         b0: beta 0
6         b1: beta 1
7         X: X vector
8         y: y vector
9     Returns:
10        r2 score
11    """
12    SS_t = 0.0
13    y_mean = sum(y) / float(len(y))
14    for i in range(len(y)):
15        SS_t += (y[i] - y_mean)**2
16
17    SS_r = get_RSS(b0, b1, X, y)
18    r2 = 1 - SS_r / SS_t
19    return r2
```

```python
1 # run cell to print RMSE
2 print("R2: ", get_R2(b0, b1, X, y))
```

```
R2:  0.5194952370037837
```

▼ Q: Analyze what the above $R^2$ score indicates about the model.

A: Usually, the larger the R2, the better the regression model fits your observations. In this case, an $R^2$ of 52% implies the model explains 52% of the variation in the response variable around its mean.

A caveat about this metric is that it cannot determine whether the coefficient estimates and predictions are biased, which is why we must also assess residual plots.

Now, we will compare the above results with the results from using scikit-learn, a machine learning library in Python. Read the documentation ([https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html)) to learn how to use this library. Return the $R^2$ score and RMSE.

```python
1 # TODO: scikit learn function
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import mean_squared_error
4 from sklearn.metrics import r2_score
5
6 def linear_regression_SKL(X, y):
7     """
8     Params:
9         X: X vector
10        y: y vector
11    Returns:
12        rmse and r2 as a tuple
13    """
14    linear_regression = LinearRegression()
15    X_reshaped = np.array(X).reshape(len(X),1)
16    linear_regression.fit(X_reshaped, y)
17    y_predicted = linear_regression.predict(X_reshaped)
18
19    rmse = mean_squared_error(y, y_predicted, squared=False)
20    r2 = r2_score(y, y_predicted)
21
22    return (rmse, r2)
```

```python
1 # run this cell to print results from SKL LR
2 rmse, r2 = linear_regression_SKL(X, y)
3 print("RMSE: ", rmse)
4 print("R2: ", r2)
```

```
RMSE:  116.71181887064391
R2:  0.5194952370037791
```

Q: Analyze the results and compare the RMSE and $R^2$ to the previous method.

A: Comparing $R^2$ and $RMSE$ of both implementations:

| Implementation | $R^2$ | $RMSE$ |
| --- | --- | --- |
| Sci-kit Learn | 0.5194952370037791 | 116.71181887064391 |
| Sheil's Implementation | 0.5194952370037837 | 116.7118188706435 |

This result helps validate our implementation of linear regression, since the two sets of predictions (from sci-kit learn and our own implementation) yield the same scores.