

Course	DS 203: Programming for Data Science
Activity	Type: Exercise Title: E03 – Simple Linear Regression using the Python Programming Language
Background	<ul style="list-style-type: none"> <li>• The previous exercises have introduced various SLR concepts using Excel as a tool</li> <li>• The Python programming environment provides rich features to do all that spread-sheets enable – and much more.</li> <li>• This exercise is part of acquiring Python programming skills in the context of solving Data Science problems</li> </ul>
Expected outcomes of this exercise	<ul style="list-style-type: none"> <li>• Get familiar with the Python packages Numpy and Matplotlib</li> <li>• Start solving data problems using arrays and their compact representations, backed by powerful and efficient operations.</li> <li>• Start using Matplotlib functions for data visualization</li> </ul>
Tools	<ul style="list-style-type: none"> <li>• Python Notebook using either Jupyter or VSC</li> </ul>
Effort estimate	<ul style="list-style-type: none"> <li>• <b>3 hours</b></li> </ul>
Submission type	<ul style="list-style-type: none"> <li>• <b>Mandatory submission</b></li> </ul>
Due date and time	<ul style="list-style-type: none"> <li>• <b>24<sup>th</sup> August, 23:55 Hrs</b></li> </ul>
Submission instruction	<ul style="list-style-type: none"> <li>• Submit to the appropriate Moodle submission point</li> <li>• Your final submission should be in the form of a single PDF file.</li> <li>• Your solutions and answers should include explanations / graphs / charts / Tables that are necessary to fully explain the solution(s).</li> <li>• <b>Complete this exercise and submit well within time. No extensions will be granted, no email submissions will be accepted.</b></li> </ul>
Marks for the exercise	<ul style="list-style-type: none"> <li>• Credit will be given for <b>complete and timely</b> submission to Moodle</li> <li>• The exercise itself will not carry marks.</li> <li>• <b>Your understanding and skill – expected to be gained by completing this exercise – will be gauged in a quiz, test, or viva that will be conducted subsequently.</b></li> </ul>
References	<ul style="list-style-type: none"> <li>• Python Notebooks uploaded to course page on Moodle. Especially units 5 and 6.</li> <li>• Lecture notes</li> <li>• Excel spread-sheets uploaded to Moodle</li> <li>• Help documentation for Python, Numpy and Matplotlib.</li> <li>• Articles, blogs and other sources</li> </ul>

Note:

- This exercise requires skills related to Numpy and Matplotlib libraries. Peruse the Python Notebooks 5 (Numpy) and 6 (Matplotlib) uploaded to the course page on Moodle for an introduction.

Create a Python Notebook to program and complete the following tasks (estimated effort: 3 hours):

1. Load the data file **linear-data-set-for-regression.csv** into a numpy array
2. Create a scatter plot of the data using Python matplotlib functions
3. Using numpy functions calculate the regression coefficients a, b for this data set
4. Create an array of predicted values based on the regression model
5. Create an array of prediction errors
  - a) There exists a function '**scipy.stats.normaltest**'. Understand this function and use the p-value returned by it to check if the prediction errors follow normal distribution. Print the results.
6. Calculate SST, SSR, SSE and R2 values
7. Validate the expression:  $SST = SSR + SSE$
8. Create a plot with two sub-plots that have aligned horizontal scales:
  - a) Scatter plot of original data + predicted points
  - b) Error plot
9. Create a nicely formatted output text file with all the metrics you have created. Open this file in a text editor, capture a screen-shot and add it to the Notebook.
10. Compare all the numbers with those in the file **linear-data-set-for-regression-soln-aug-11.xlsx** and ensure that the outputs of this exercise match with them! Mention the result of this check in the Notebook.

Convert the Notebook into PDF (your program segments and all the generated / added outputs) and upload it as your submission for this assignment.

\*\*\*\*\*