

Student name:

Student ID:

SIT225: Data Capture Technologies

Activity 7.1: Data analysis and interpretation

Data analysis is a broad term that covers a wide range of techniques that enable you to reveal any insights and relationships that may exist within raw data. As you might expect, Python lends itself readily to data analysis. Once Python has analyzed your data, you can then use your findings to make good business decisions, improve procedures, and even make informed predictions based on what you've discovered.

You have done data wrangling using Python Pandas module already in activity 5.2. In this activity, you will learn Data science statistics and linear regression models.

Hardware Required

No hardware is required.

Software Required

Python 3

Python packages including Pandas, Numpy, Scikit-learn, seaborn, plotly

Steps:

Step	Action
1	A Jupyter Notebook is provided for Data Science exploration here (https://github.com/deakin-deep-dreamer/sit225/tree/main/week_7). You will need to fill in your student ID and name and run all the cells to observe the output. Convert the Notebook into PDF and merge with this activity sheet which needs to be combined with this week's task for OnTrack submission.

	<p>Question: There are sections in the Notebook. After running the cells and observing the outputs, provide your reflection in brief on the topic items for each section of the Notebook.</p> <p>Answer: <Your answer></p>
2	<p>Question: In the 1.1 Percentile subsection of Descriptive statistics section in the Notebook, you have calculated 10%, 25%, 50% and 75% percentiles for <i>Max_Pulse</i>. Compare these percentiles with <i>Average_Pulse</i> percentiles for any trend, if exists.</p> <p>Answer: <Your answer></p>
3	<p>Question: In the “Correlation Does not imply Causality” section answer the question regarding the increase of ice cream sale in your own understanding.</p> <p>Answer: <Your answer></p>
4	<p>Question: In the 1.7 Linear Regression section in the Notebook, a linear regression model was used to predict Calorie_Burnage from attributes such as Average_Pulse. The Duration value was predicted from the model for all the value range of Average_Pulse and a regression line was drawn. You will need to answer the follow up question next to 1.7 section where it is required to generate a linear regression model for Duration instead of Average_Pulse to predict the Calorie_Burnage. Take a screenshot of the regression line and paste it here. Also, comment on both the regression lines.</p> <p>Answer: <Your answer></p>