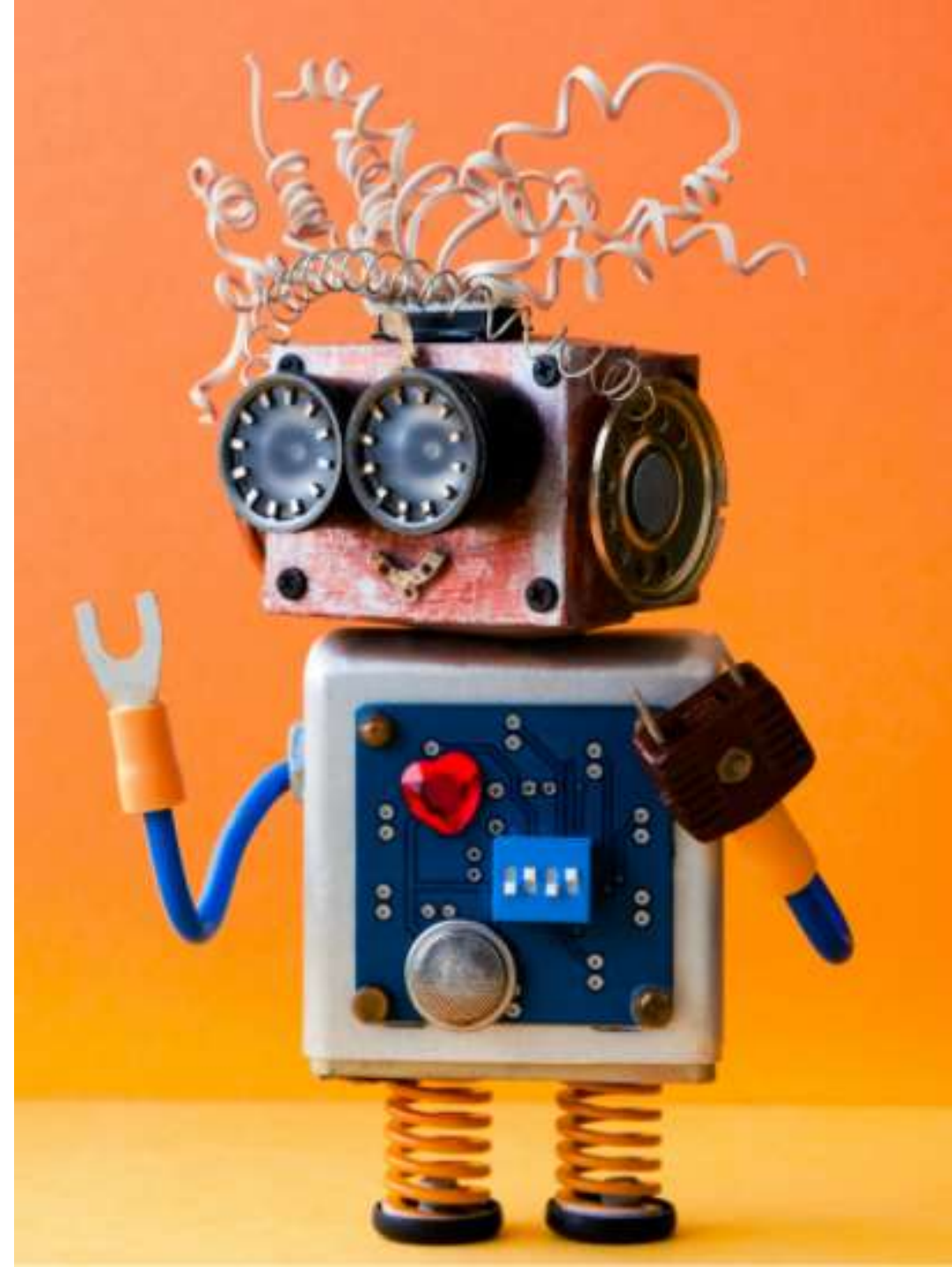# PROJECT OVERVIEW AND KEY LEARNING OUTCOMES



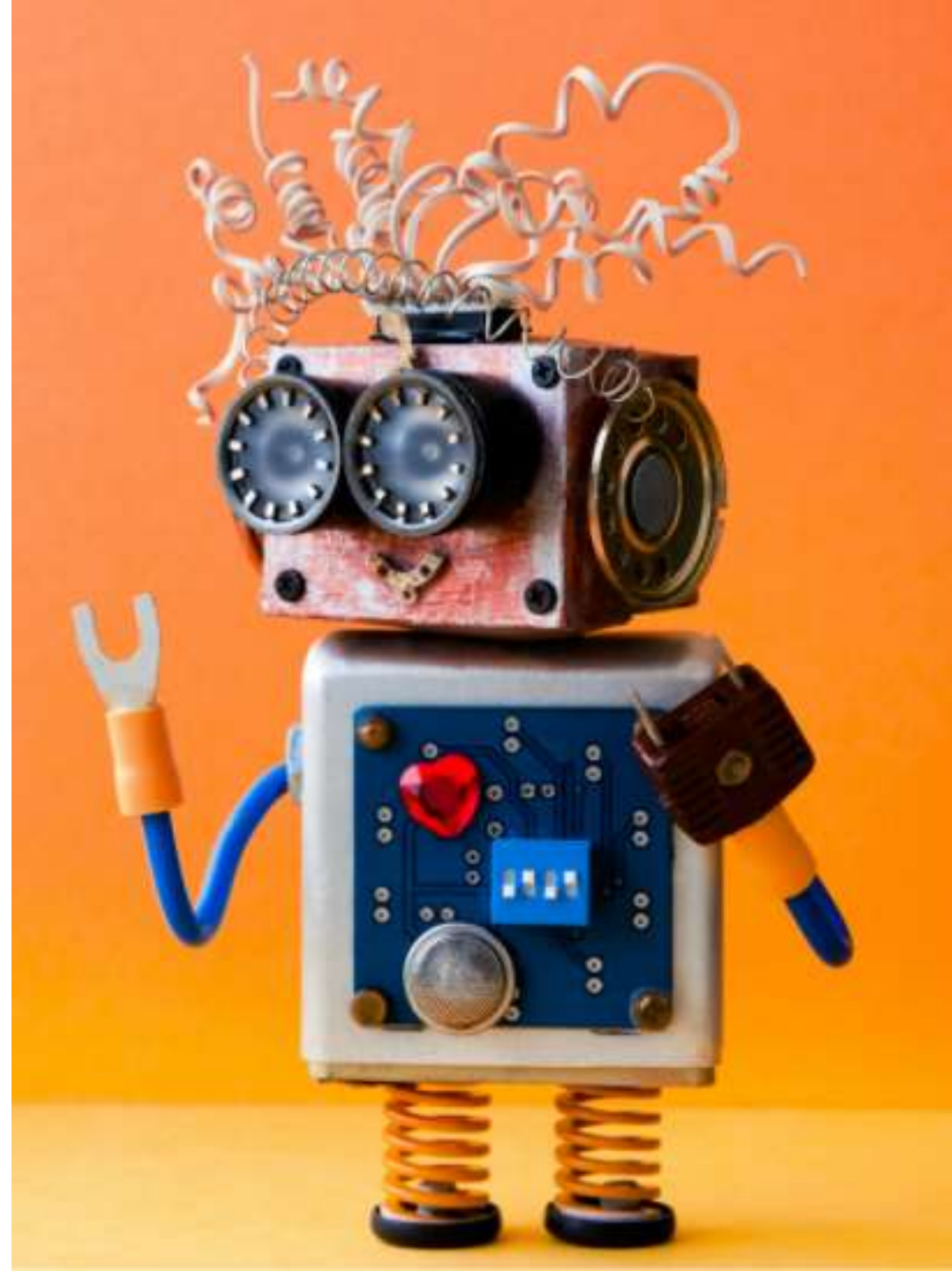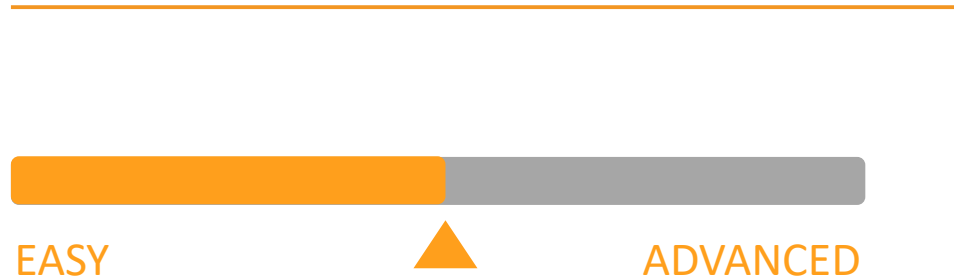EASY ▲                    ADVANCED

# PROJECT OVERVIEW

- We will analyze human resources information using Pandas in AWS SageMaker Studio.
- We will learn how to:
    1. Perform statistical analysis on real world datasets.
    2. Deal with missing data using pandas
    3. Change pandas DataFrame datatypes
    4. Define a function and apply it to a Pandas DataFrame column
    5. Pandas Operations and filtering
    6. Calculate and display correlation matrix
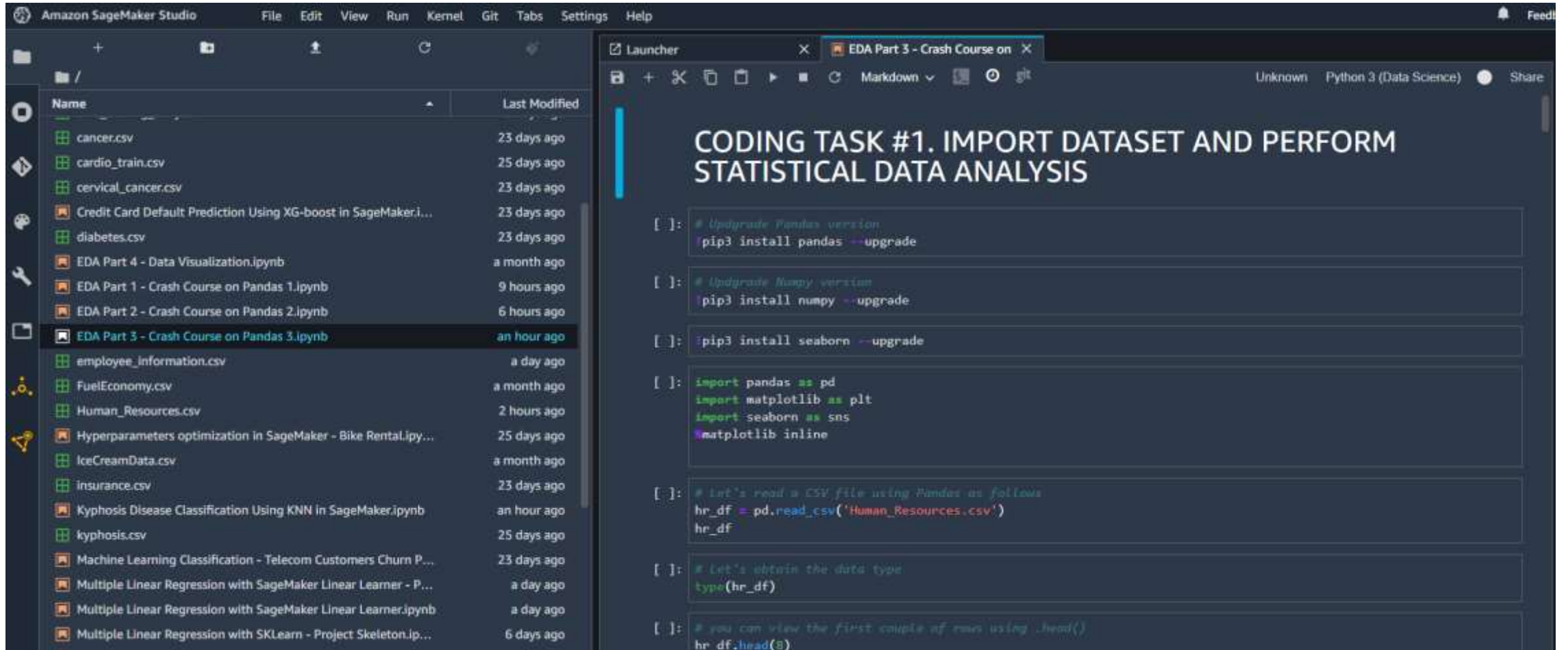    7. Use seaborn library to show heatmap

| | Age | Attrition | BusinessTravel | DailyRate | Department | DistanceFromHome | Education | EducationField | EmployeeCount | EmployeeNumber | ... | Relati |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 41 | Yes | Travel_Rarely | 1102 | Sales | 1 | 2 | Life Sciences | 1 | 1.0 | ... | |
| 1 | 49 | No | Travel_Frequently | 279 | Research & Development | 8 | 1 | Life Sciences | 1 | 2.0 | ... | |
| 2 | 37 | Yes | Travel_Rarely | 1373 | Research & Development | 2 | 2 | Other | 1 | 4.0 | ... | |
| 3 | 33 | No | Travel_Frequently | 1392 | Research & Development | 3 | 4 | Life Sciences | 1 | 5.0 | ... | |
| 4 | 27 | No | Travel_Rarely | 591 | Research & Development | 2 | 1 | Medical | 1 | 7.0 | ... | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 1465 | 36 | No | Travel_Frequently | 884 | Research & Development | 23 | 2 | Medical | 1 | 2061.0 | ... | |
| 1466 | 39 | No | Travel_Rarely | 613 | Research & Development | 6 | 1 | Medical | 1 | 2062.0 | ... | |
| 1467 | 27 | No | Travel_Rarely | 155 | Research & Development | 4 | 3 | Life Sciences | 1 | 2064.0 | ... | |
| 1468 | 49 | No | Travel_Frequently | 1023 | Sales | 2 | 3 | Medical | 1 | 2065.0 | ... | |
| 1469 | 34 | No | Travel_Rarely | 628 | Research & Development | 8 | 3 | Medical | 1 | 2068.0 | ... | |

1470 rows × 35 columns

# PROJECT DEMO

EASY       ▲       ADVANCED

# PROJECT DEMO

# FINAL CAPSTONE END-OF-DAY PROJECT

EASY ▲ ADVANCED

# FINAL PROJECT

- In this project, we will perform basic Exploratory Data Analysis (EDA) on the Kyphosis disease Dataset.
- Kyphosis is an abnormally excessive convex curvature of the spine.
- Dataset contains 81 rows and 4 columns representing data on children who have had corrective spinal surgery.
- **INPUTS:** 1. Age: in months, 2. Number: the number of vertebrae involved, 3. Start: the number of the first (topmost) vertebra operated on.
- **OUTPUTS:** Kyphosis which represents a factor with levels absent present indicating if a kyphosis (a type of deformation) was present after the operation.
- Using the "kyphosis.csv" included in the course package, write a python script to perform the following tasks:
    1. Import the "kyphosis.csv" file using Pandas
    2. Perform basic Exploratory Data Analysis (EDA) on the data
    3. List the average, minimum and maximum age (in years) considered in this study using 2 methods
    4. Plot the correlation matrix
    5. Convert the age column datatype from int64 to float64
    6. Define a function that converts age from months to years
    7. Apply the function to the "Age" column and add the results into a new column entitled "Age in Years"
    8. What are the features of the oldest and youngest child in this study?