# CoGrammar

# Data Preprocessing & Exploratory Data Analysis

# Data Science Lecture Housekeeping

- The use of disrespectful language is prohibited in the questions, this is a supportive, learning environment for all - please engage accordingly. **(FBV: Mutual Respect.)**

- No question is daft or silly - **ask them!**

- There are **Q&A sessions** midway and at the end of the session, should you wish to ask any follow-up questions. Moderators are going to be answering questions as the session progresses as well.

- If you have any questions outside of this lecture, or that are not answered during this lecture, please do submit these for upcoming Open Classes. You can submit these questions here: **Open Class Questions**

CoGrammar

# Data Science Lecture Housekeeping cont.

- For all **non-academic questions**, please submit a query:
  **www.hyperiondev.com/support**

- Report a **safeguarding** incident:
  **www.hyperiondev.com/safeguardreporting**

- We would love your **feedback** on lectures: **Feedback on Lectures**

CoGrammar

# Lecture Objectives

- **Recapping data cleaning & preprocessing.**

- **Introducing and explaining normalisation and standardisation.**

- **Introduce the exploratory data analysis.**

# Recap : Data Analysis

★ **When trying to understand a dataset, we need to think analytically.**

★ **However, it might not be enough and we would need the right tools to analyse our dataset.**

★ **With the main goal being that we gain an in-depth understanding of our dataset.**

# Step One : Clean our Data

★ **Data that has inconsistencies / missing values, introduces noise, which affects our analysis.**

★ **There are techniques we can implement to clean our data:**
  - **Dropping columns**
  - **Filling and replacing**
  - **Indexing**
  - **Creating sub-datasets**

# Dropping Columns

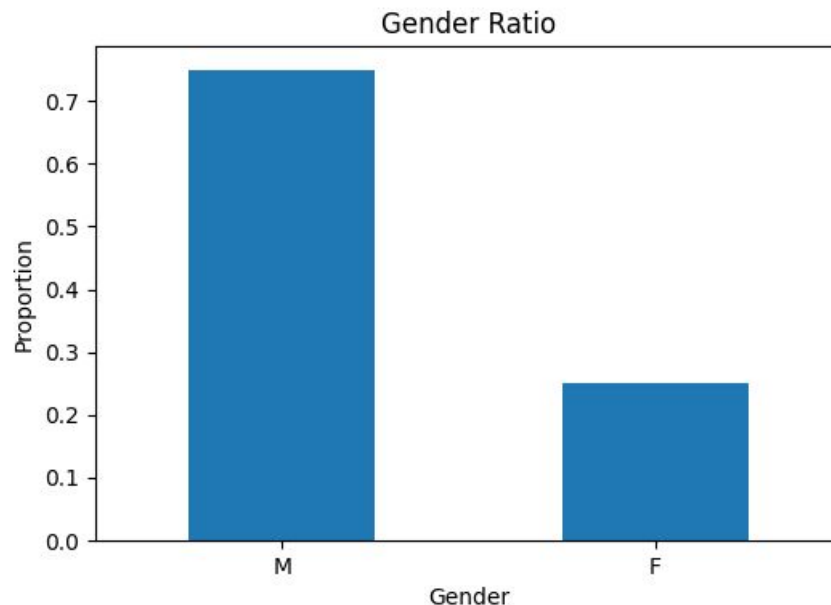| first_name | last_name | student_id | gender | math_score | english_score | science_score | Mother | Father | Address |
|---|---|---|---|---|---|---|---|---|---|
| Theresa | Imore | 1 | F | 89 | 78 | 50 | Jane Claremint | Michael Scofield | 90 East Buckingham |
| Brook | Gratrex | 2 | M | 67 | 56 | 65 | Joesephine Brow | Lincon Burrows | 60 Fairway Ave. |
| Jerrold | Isenor | 3 | M | 98 | 67 | 42 | Claire Mint | Ballack Benet | 749 West Kingston Street |
| Jo | Pretsel | 4 | M | 56 | 72 | 87 | Alexis Mahone | Duke Harry | Lawrence Township, NJ |

★ **Dropping columns depends on our goal of our analysis. For instance, we want to analyse the proportion of Male / Female, location data. We wouldn't need grades in this case.**

```python
data.drop(['math_score', 'english_score', 'science_score'],
          axis=1, inplace=True)
```

**CoGrammar**

★ **Here is the result of us using the dataset after dropping the redundant columns.**

★ **As you can see, we only drop data that is not useful.**

**Typically we try not to drop data, so this method comes in when it becomes very obvious the data isn't needed.**

Gender Ratio

# Filling and Replacing

★ **Often, datasets are filled inconsistently, missing values and inconsistencies in data entry exist. We want to ensure data consistency and keep entries uniform.**

★ **For example, an individual's nationality could be entered as "UK", while others are entered as "United Kingdom". In this case we opt to stick with "United Kingdom" in the data. If they are not from the UK, we label it "Other"**

```python
data.replace('UK', 'United Kingdom', inplace=True)
data.fillna('Other', inplace=True)
```

# Indexing & Creating sub-datasets

★ **Often, we want to be able to find out if a specific rule works generally. If we can make a prediction on seen data, how well would it hold up to unseen data?**
  ○ **This concept will be looked at further in Machine Learning.**

★ **It would be necessary to split up our dataset:**

```
train_data = data.loc[1:80]
test_data = data.loc[81:100]
```

★ **Common splits are 80:20 for train:test in industry.**

# Step Two : Missing Data

★ **The method .fillna() can be dangerous tool. If done incorrectly, it is possible to introduce bias to our dataset.**

★ **How do we use this responsibly? That would depend on the type of missingness.**
  - **Missing Completely at Random (MCAR)**
  - **Missing at Random (MAR)**
  - **Missing Not at Random (MNAR)**

# Missing Completely at Random (MCAR)

★ **Missing data is not in any way affected by other variables in the dataset.**

    ○ **E.g. A test paper was lost, or a blood sample is missing.**

    ○ **There is no reasonable way to guess the missing value.**

# Missing at Random (MAR)

★ **Missing data is in some way related to other variables in the dataset.**

★ **Arises as a result of implicit biases in the data itself.**
  ○ **E.g. Consider a survey on depression. Statistically, males are less likely to fill in a survey on the severity of their depression.**
  ○ **As a result, missing data can be guessed from the proportion of males in the survey.**

# Missing Not at Random (MNAR)

★ **Similar to MAR, it arises as a result of implicit biases.**

★ **However, these biases haven't been measured on the data itself.**
   ○ **E.g. Fewer reported COVID cases once restrictions have been lifted.**
   ○ **It does not mean COVID has disappeared, simply means that less people are getting tested.**

★ **Different to MCAR, as it is not random event that caused the data to go missing.**
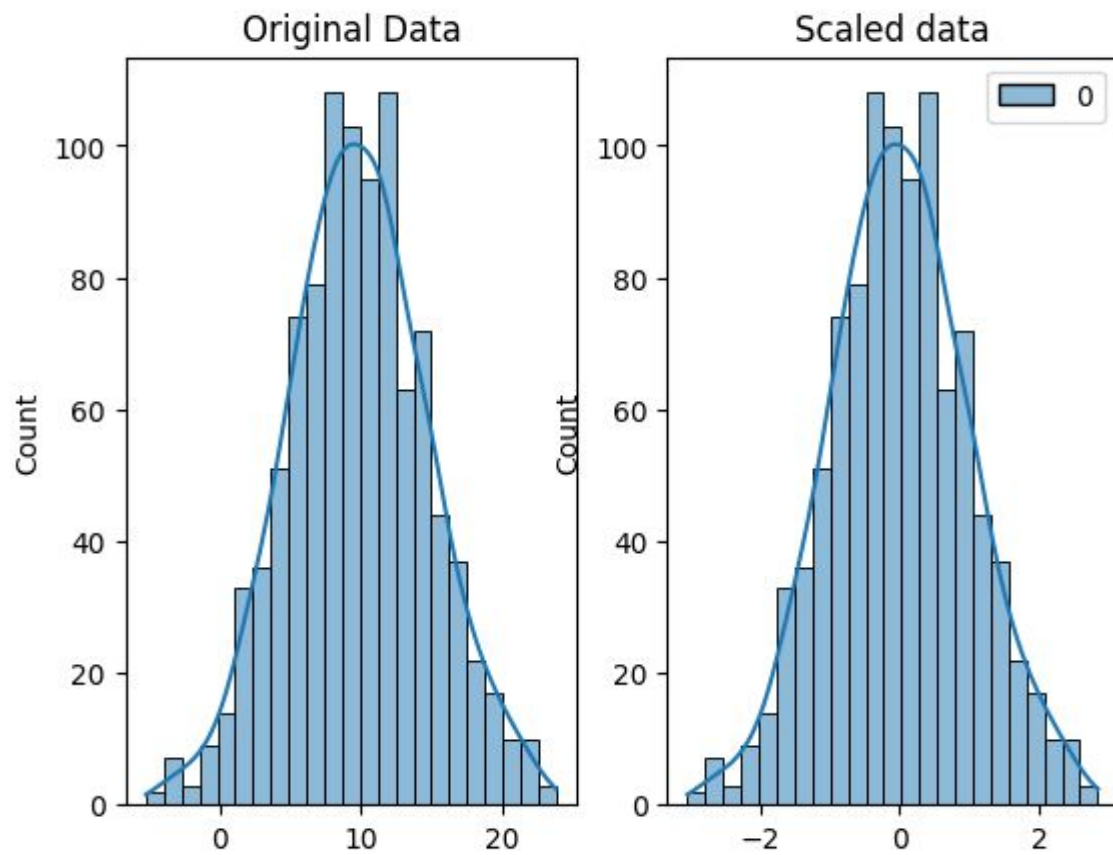
# Let's Breathe

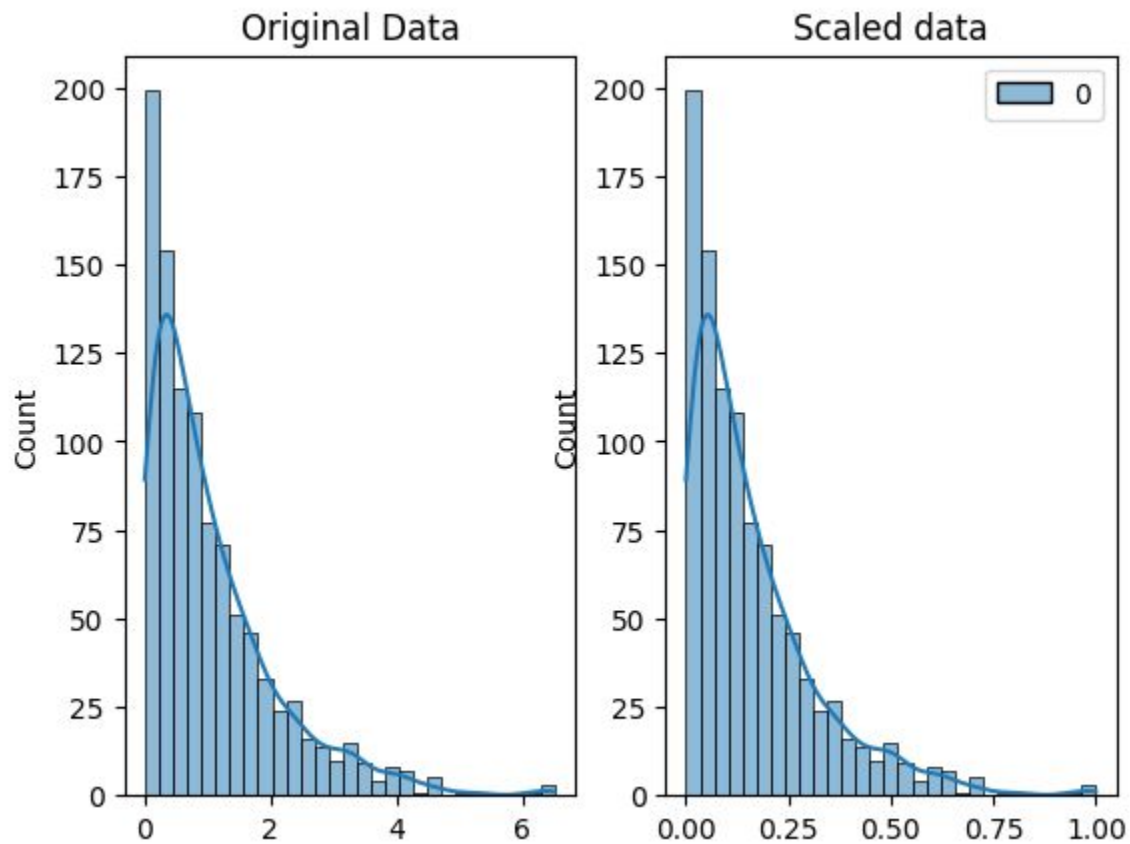Let's take a small break before moving on to the next topic.

CoGrammar

# Standardisation & Normalisation

★ Useful when studying relationships between variables in a dataset.

★ Units of measurement can sometimes affect your observations.

★ This also affects Machine Learning applications.

★ Also known as Feature Scaling.
  ○ Feature scaling is the process of restricting the values in a particular feature to ensure that all features lie on the same scale.

# What's the Difference?

One of the reasons that it's easy to get confused between standardisation and normalisation is because the terms are sometimes used interchangeably and, to make it even more confusing, they are very similar. In both cases, you're transforming the values of numeric variables so that the transformed data points have specific helpful properties. The difference is that, in standardisation, you're changing the range of your data to have a mean of 0 and a standard deviation of 1, while in normalisation you're changing the range of your data to have a maximum of 1 and a minimum of 0.

Original Data / Scaled data

# Exploratory Data Analysis

★ We explore data via analysis. Machine learning works by making assumptions of the properties of the data.

★ Exploratory Data Analysis (EDA) is how we learn about those properties.

★ We explore the data, find patterns within the data and extract insights for the EDA.
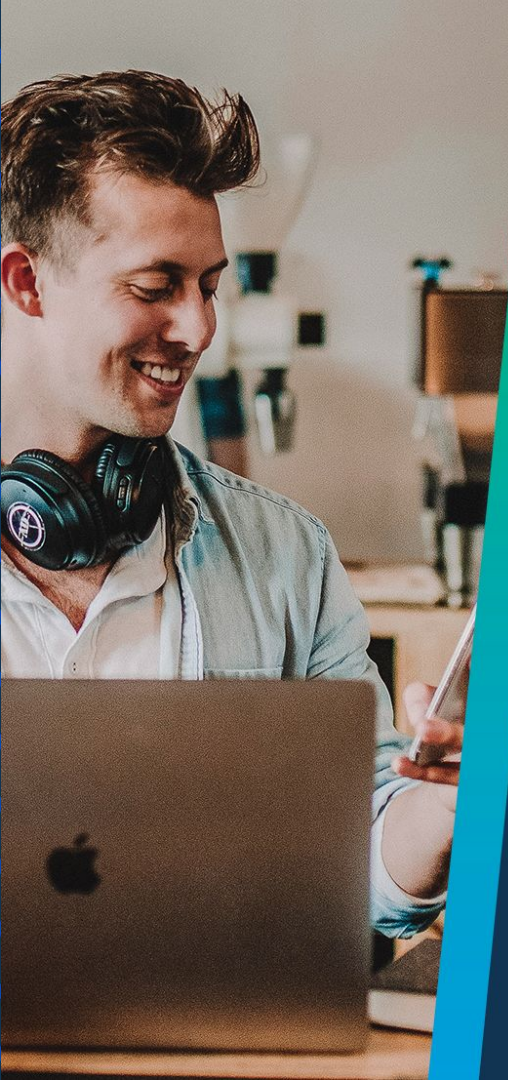
# How to Explore your Dataset

★ **There is no prescribed method, do what feels right, and benefits most in terms of understanding your dataset.**

★ **It's vital to ask questions as well :**
  ○ **What kind of data do I have?**
  ○ **What and where is the missing data? How do I deal with the missing values.**
  ○ **Where are the outliers? Should I remove them?**
  ○ **How do I add, change or remove features to get the most of the data?**

# CoGrammar

## Q & A SECTION

**Please use this time to ask any questions relating to the topic, should you have any.**

**CoGrammar**

# Thank you for joining us

1. Take regular breaks
2. Stay hydrated
3. Avoid prolonged screen time
4. Practice good posture
5. Get regular exercise

*"With great power comes great responsibility"*