# Data and Artificial Intelligence
# Cyber Shujaa Program

## Week 2 Assignment
## Data Wrangling in Python

**Student Name:** Faith Jeptoo

**Student ID:** CS-DA02-25005

# Table of Contents

# Introduction

This week's assignment was to practice data wrangling concepts using the Netflix dataset from Kaggle.

The objectives were:
- Load the dataset and explore its structure.
- Discover data types, missing values, and quality issues.
- Clean the dataset by handling duplicates, missing values, and formatting inconsistencies.
- Transform and enrich the dataset.
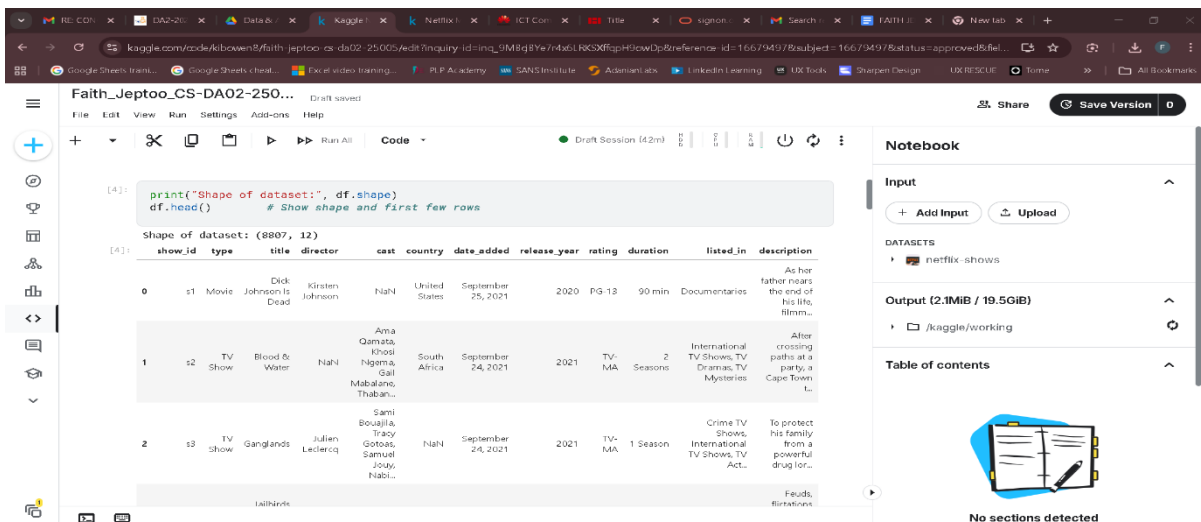- Validate and export the final dataset.

# Tasks Completed

## Step 1: Import Libraries and Load Dataset

```python
import pandas as pd
import os

# Check current working directory
print(os.getcwd())

# Load Netflix dataset
filepath = '/kaggle/input/netflix-shows/netflix_titles.csv'
df = pd.read_csv(filepath)

# Show shape and first few rows
print("Shape of dataset:", df.shape)
df.head()
```

## Step 2: Data Discovery

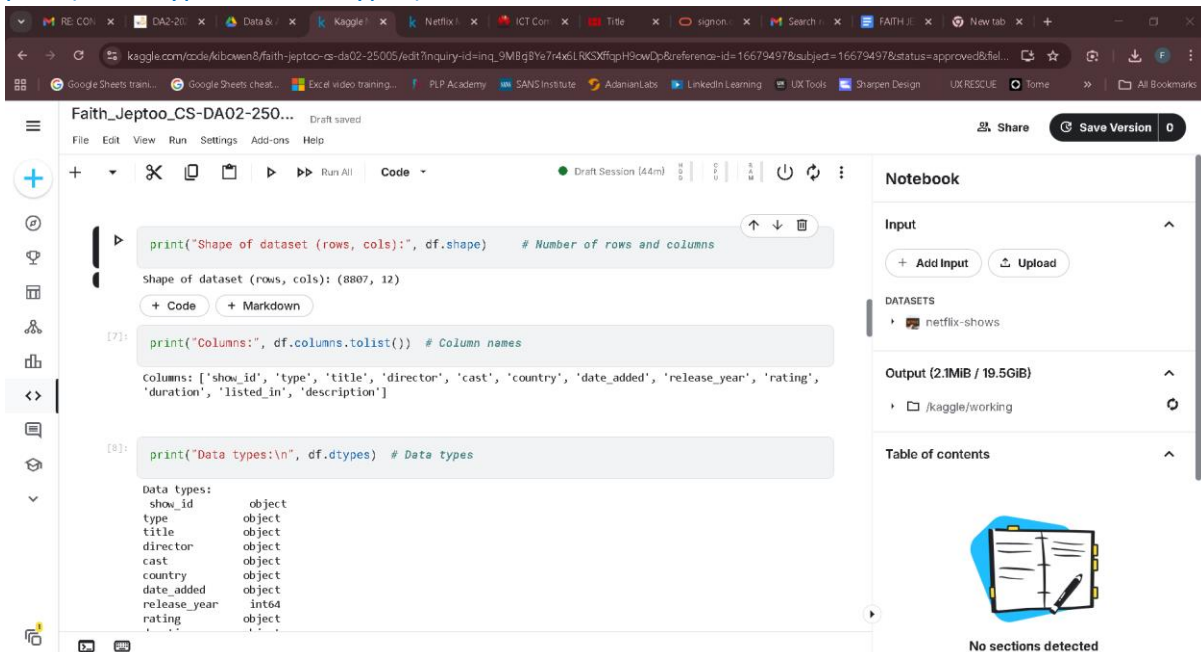# Overview of the dataset
df.info()



# Number of rows and columns
print("Shape of dataset (rows, cols):", df.shape)

# Column names
print("Columns:", df.columns.tolist())

# Data types
print("Data types:\n", df.dtypes)

# Missing values (counts)
print("Missing values per column:\n", df.isnull().sum())



# Missing values in percentage (overall and per column)
print("Average missing % across dataset:", df.isnull().sum().mean() * 100)
print("Missing % per column:\n", df.isnull().mean() * 100)



# Duplicates
print("Duplicate rows:", df.duplicated().sum())

## Step 3: Structuring

# Standardize column names: lowercase + replace spaces with underscores
df.columns = df.columns.str.lower().str.replace(' ', '_')

# Convert 'date_added' to datetime
df['date_added'] = pd.to_datetime(df['date_added'], errors='coerce')

# Separate duration into numeric value and unit
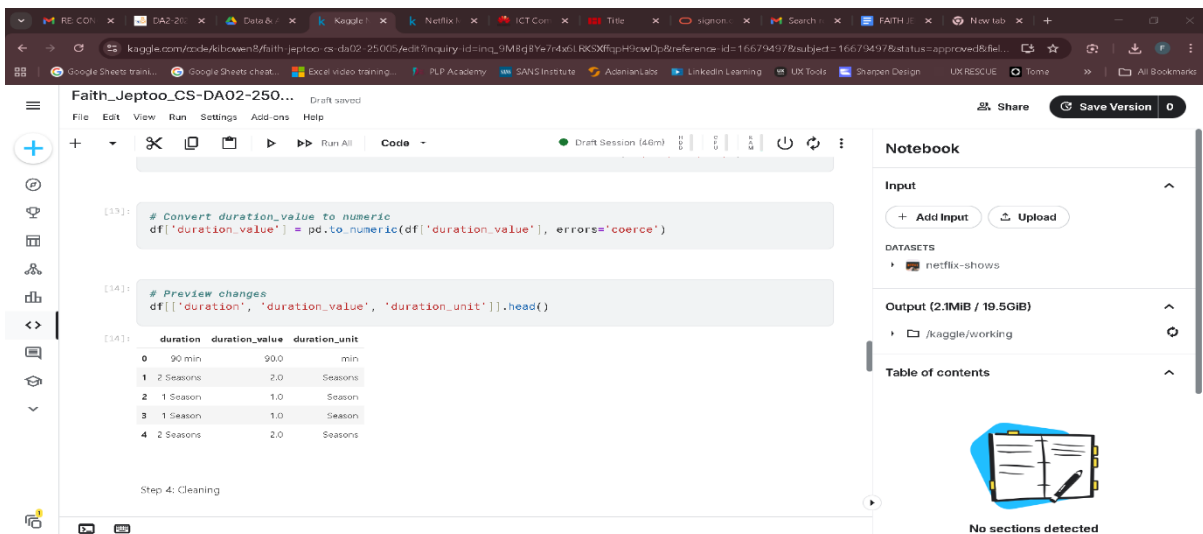df[['duration_value', 'duration_unit']] = df['duration'].str.extract(r'(\d+)\s*(\w+)')

# Convert duration_value to numeric
df['duration_value'] = pd.to_numeric(df['duration_value'], errors='coerce')

# Preview changes
df[['duration', 'duration_value', 'duration_unit']].head()

## Step 4: Cleaning

```python
# Remove duplicates
print("Duplicates before:", df.duplicated().sum())
df = df.drop_duplicates()
print("Duplicates after:", df.duplicated().sum())

# Drop unnecessary column
df = df.drop(columns=['description'])

# Fill missing values
df['director'] = df['director'].fillna('Not Given')
df['cast'] = df['cast'].fillna('Not Given')
df['country'] = df['country'].fillna('Not Given')

# Drop rows where critical fields are missing
df = df.dropna(subset=['date_added', 'rating', 'duration'])
```
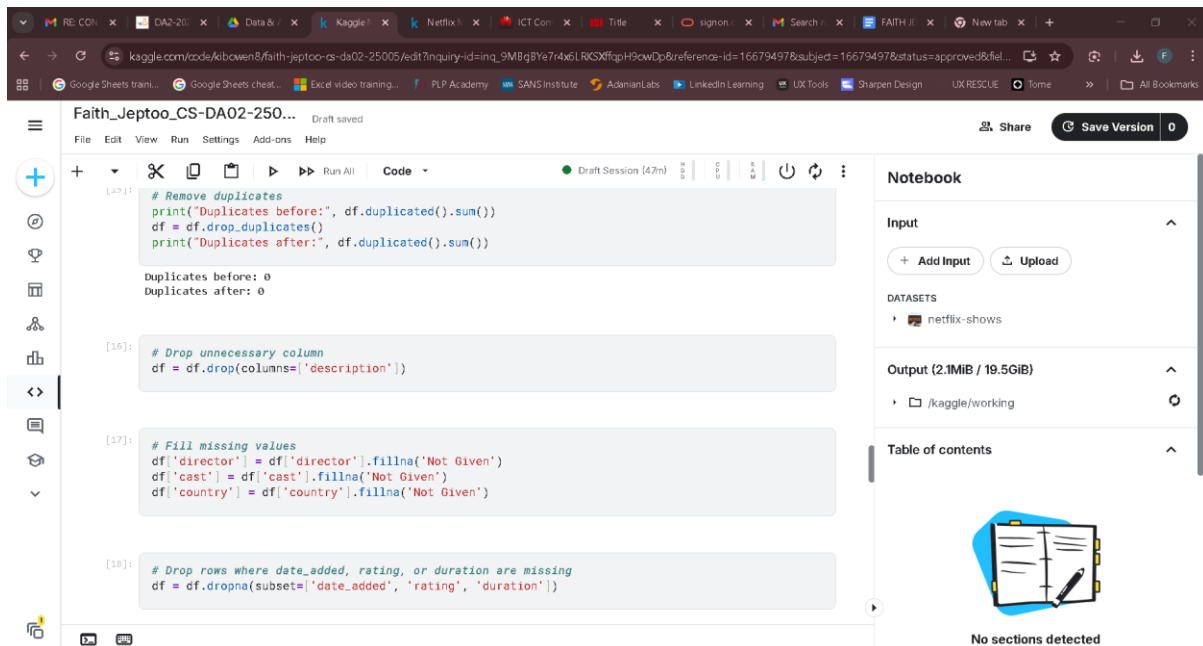


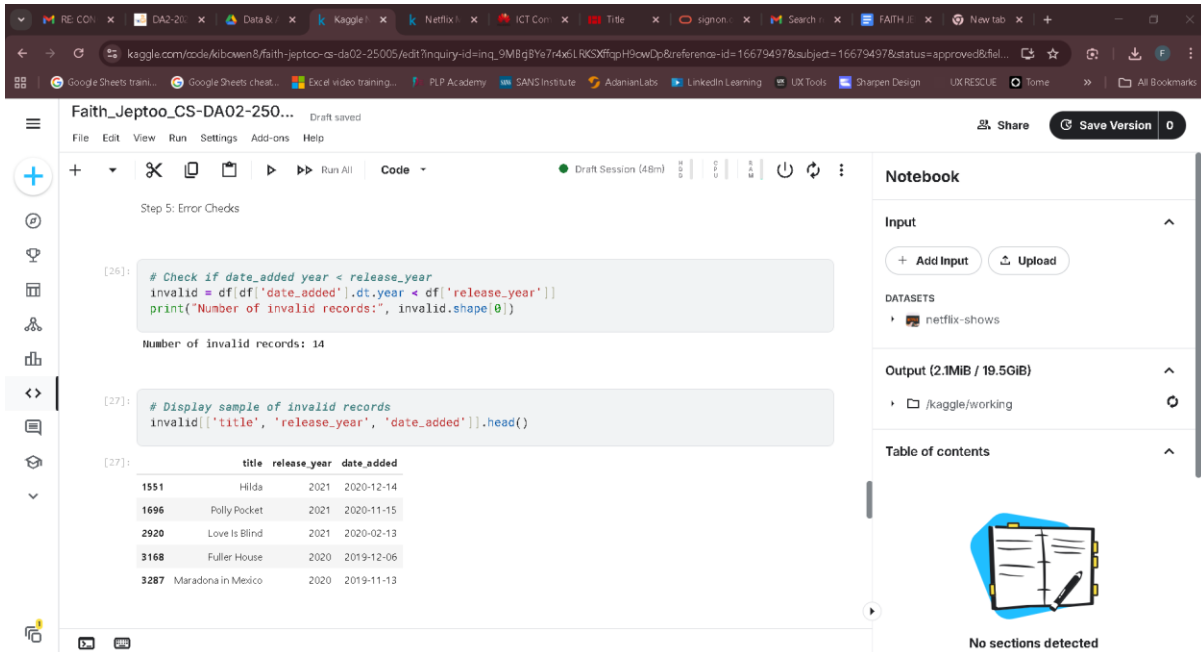## Step 5: Error Checks

```python
# Check if date_added year < release_year
invalid = df[df['date_added'].dt.year < df['release_year']]
print("Number of invalid records:", invalid.shape[0])

# Display sample of invalid records
invalid[['title', 'release_year', 'date_added']].head()
```
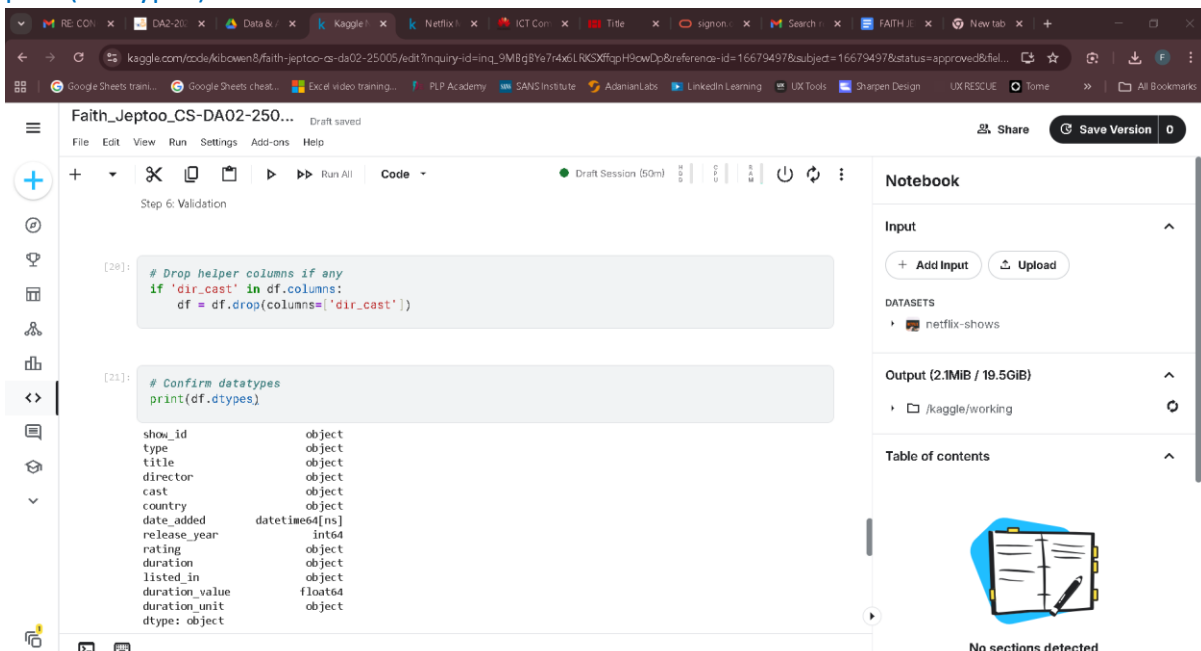
## Step 6: Validation

```python
# Drop helper columns if any
if 'dir_cast' in df.columns:
    df = df.drop(columns=['dir_cast'])
```

```python
# Confirm datatypes
print(df.dtypes)
```



```python
# Check for missing values again
print("Missing values after cleaning:\n", df.isnull().sum())
```

# Sample few rows
df.sample(5)
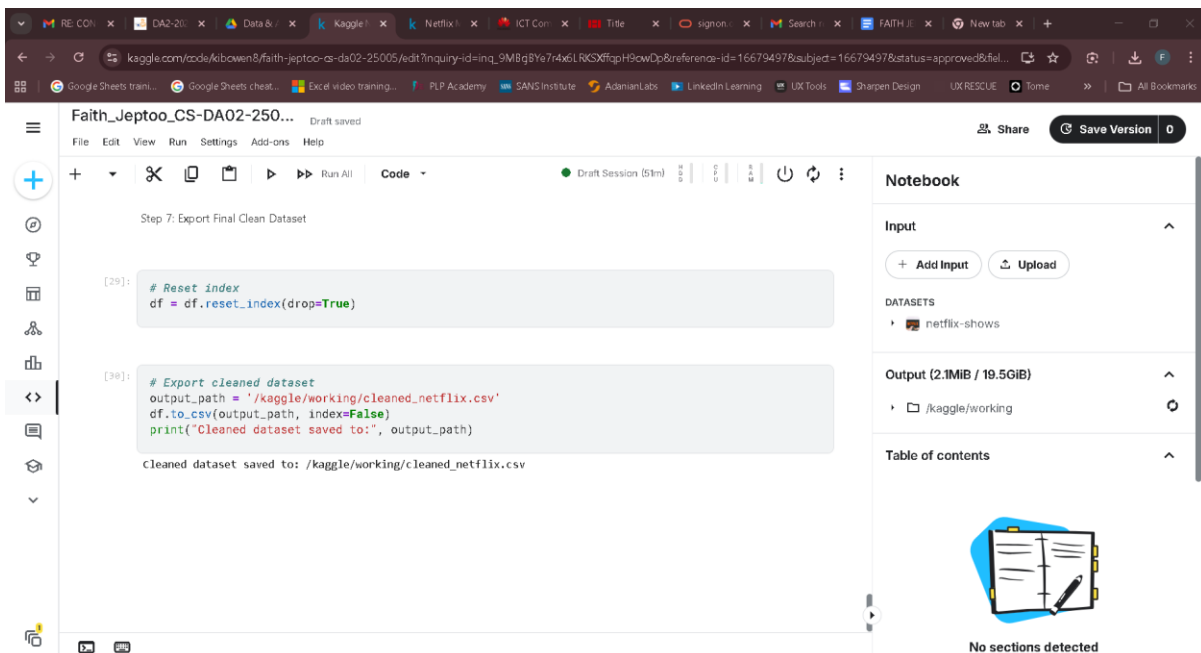


## Step 7: Export Final Dataset

# Reset index
df = df.reset_index(drop=True)

# Export cleaned dataset
output_path = '/kaggle/working/cleaned_netflix.csv'
df.to_csv(output_path, index=False)
print("Cleaned dataset saved to:", output_path)

# Conclusion

This assignment provided hands-on experience in data wrangling. I learned how to explore, clean, structure, and validate real-world datasets. The final Netflix dataset is now ready for analysis and visualization.

# Link to Notebook

https://www.kaggle.com/code/kibowen8/faith-jeptoo-cs-da02-25005-week-2