# TP: Data Wrangling with Pandas

## Your Name

### December 18, 2024

## Introduction

In this practical session, we will explore various data wrangling techniques using the pandas library in Python. Each mini-project will focus on a specific task, and solutions will be provided for each.

# 1 Mini-Project 1: Titanic Survivor Analysis with Visualization

## Description

You have a dataset of Titanic passengers. Your task is to create a pandas DataFrame from the data, then calculate the number of survivors and non-survivors by passenger class. Next, visualize this data using a bar chart.

## Solution

```python
import pandas as pd
import matplotlib.pyplot as plt

# Load the data
url = 'https://raw.githubusercontent.com/chrisalbon/sim_data/
    master/titanic.csv'
dataframe = pd.read_csv(url)

# Calculate the number of survivors and non-survivors by
    class
survivors_by_class = dataframe[dataframe['Survived'] == 1].
    groupby('PClass').size()
non_survivors_by_class = dataframe[dataframe['Survived'] ==
    0].groupby('PClass').size()
```

```
11
12 # Visualize the data
13 classes = survivors_by_class.index
14 survivors_counts = survivors_by_class.values
15 non_survivors_counts = non_survivors_by_class.values
16
17 plt.bar(classes, survivors_counts, label='Survivors')
18 plt.bar(classes, non_survivors_counts, bottom=
      survivors_counts, label='Non-Survivors')
19 plt.xlabel('Class')
20 plt.ylabel('Number of Passengers')
21 plt.title('Survivors and Non-Survivors by Class')
22 plt.legend()
23 plt.show()
```

# 2 Mini-Project 2: Handling Missing Values and Age Analysis

## Description

You have a DataFrame containing data on Titanic passengers, but some age values are missing. Your task is to replace these missing values with the mean age. Then, analyze the age distribution by creating a histogram.

## Solution

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 # Load the data
6 url = 'https://raw.githubusercontent.com/chrisalbon/sim_data/
      master/titanic.csv'
7 dataframe = pd.read_csv(url)
8
9 # Replace missing values with the mean age
10 mean_age = dataframe['Age'].mean()
11 dataframe['Age'].fillna(mean_age, inplace=True)
12
13 # Display a histogram of the ages
14 plt.hist(dataframe['Age'], bins=20, color='blue', edgecolor='
      black')
15 plt.xlabel('Age')
16 plt.ylabel('Number of Passengers')
17 plt.title('Age Distribution of Passengers')
```

```
18 plt.show()
```

# 3 Mini-Project 3: Renaming Columns and Adding a New Feature

## Description

You have a DataFrame with unclear column names. Your task is to rename the columns to make them more descriptive. Then, add a new column that indicates whether a passenger is a child (age < 18) or an adult.

## Solution

```python
1 import pandas as pd
2
3 # Load the data
4 url = 'https://raw.githubusercontent.com/chrisalbon/sim_data/
    master/titanic.csv'
5 dataframe = pd.read_csv(url)
6
7 # Rename the columns
8 dataframe.rename(columns={'PClass': 'PassengerClass', 'Sex':
    'Gender'}, inplace=True)
9
10 # Add a new column to indicate if the passenger is a child or
    an adult
11 dataframe['IsChild'] = dataframe['Age'] < 18
12
13 # Display the first few rows to verify
14 print(dataframe.head())
```

# 4 Mini-Project 4: Grouping by Class and Calculating Detailed Statistics

## Description

You have a DataFrame containing information about Titanic passengers, grouped by class (1st, 2nd, 3rd). Your task is to calculate the mean, median, minimum, and maximum age for each class. Then, display these statistics in a DataFrame.

## Solution

```python
import pandas as pd

# Load the data
url = 'https://raw.githubusercontent.com/chrisalbon/sim_data/
    master/titanic.csv'
dataframe = pd.read_csv(url)

# Group by class and calculate detailed statistics
stats_by_class = dataframe.groupby('PClass')['Age'].agg(['
    mean', 'median', 'min', 'max'])

print(stats_by_class)
```

# 5 Mini-Project 5: Merging DataFrames and Sales Analysis

## Description

You have two DataFrames, one containing information about employees and the other containing their sales. Your task is to merge these two DataFrames on a common column (e.g., employee ID). Then, calculate the total sales by employee and display the top 5 sellers.

## Solution

```python
import pandas as pd

# Create the DataFrames
employee_data = {'employee_id': ['1', '2', '3', '4'], 'name':
     ['Amy Jones', 'Allen Keys', 'Alice Bees', 'Tim Horton']}
sales_data = {'employee_id': ['3', '4', '5', '6'], '
    total_sales': [23456, 2512, 2345, 1455]}

dataframe_employees = pd.DataFrame(employee_data)
dataframe_sales = pd.DataFrame(sales_data)

# Merge the DataFrames
merged_dataframe = pd.merge(dataframe_employees,
    dataframe_sales, on='employee_id', how='inner')

# Calculate the total sales by employee
```

```
14 total_sales_by_employee = merged_dataframe.groupby('name')['
       total_sales'].sum().reset_index()
15
16 # Display the top 5 sellers
17 top_5_sellers = total_sales_by_employee.sort_values(by='
       total_sales', ascending=False).head(5)
18
19 print(top_5_sellers)
```