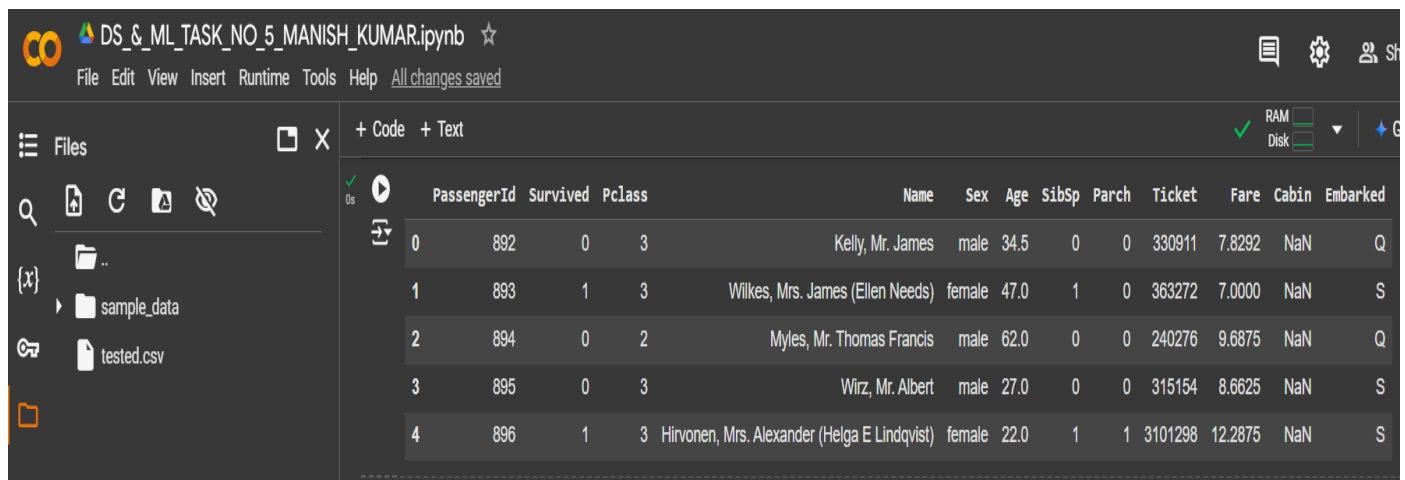# Task 5 Report
# Time Series Analysis with Pandas

## 1. Task Description

This task focuses on performing **time series analysis** using Pandas. The key objectives include:

- ❖ Simulating a datetime column for time-based indexing and operations.
- ❖ Handling missing values in numeric columns.
- ❖ Resampling the dataset to calculate weekly averages.
- ❖ Computing rolling statistics like rolling mean and rolling standard deviation.
- ❖ Creating lagged features to identify time-based dependencies.
- ❖ Visualizing the time series using Matplotlib to understand trends, patterns, and variability.
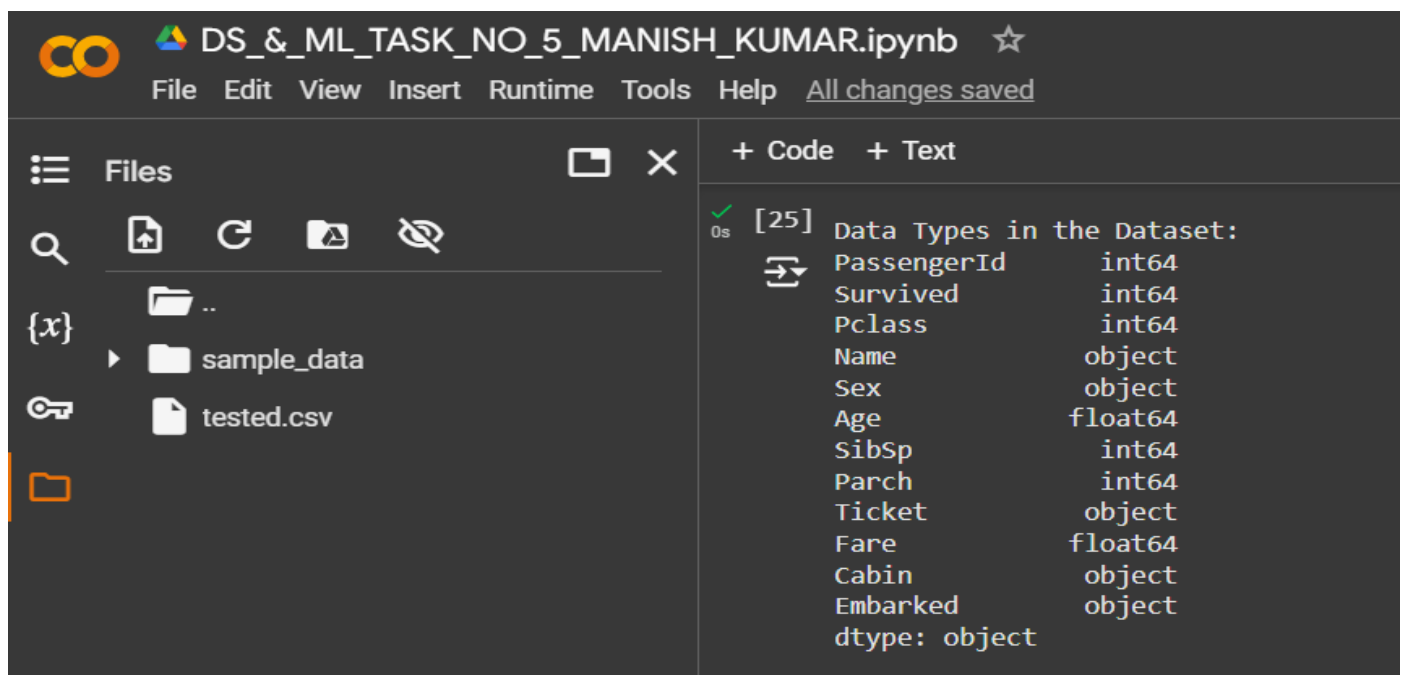
## 2. Attach Screenshot of Output

❖ **Dataset Preview:**



❖ **Data Types in the Dataset:**

# Task 5 Report
## Time Series Analysis with Pandas

❖ **Missing values in the Dataset:**



```
✓ [26]  Missing Values in the Dataset:
0s        PassengerId     0
⮕        Survived        0
          Pclass          0
          Name            0
          Sex             0
          Age             0
          SibSp           0
          Parch           0
          Ticket          0
          Fare            0
          Cabin         327
          Embarked        0
          dtype: int64
```

❖ **Dataset with Synthetic Datetime Column:**



```
✓        Dataset with Synthetic Datetime Column:
0s                     PassengerId  Survived  Pclass  \
⮕        date
          2023-01-01         892        0        3
          2023-01-02         893        1        3
          2023-01-03         894        0        2
          2023-01-04         895        0        3
          2023-01-05         896        1        3

                                               Name     Sex   Age  SibSp  \
          date
          2023-01-01                Kelly, Mr. James    male  34.5      0
          2023-01-02  Wilkes, Mrs. James (Ellen Needs)  female  47.0      1
```
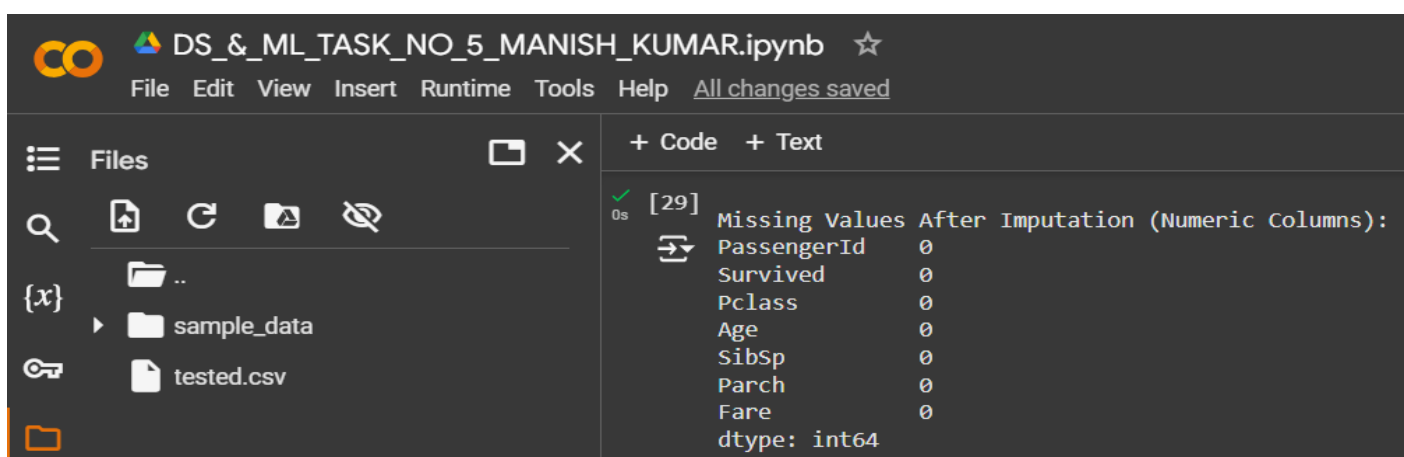
❖ **Missing values after Imputation:**



```
✓ [29]  Missing Values After Imputation (Numeric Columns):
0s        PassengerId     0
⮕        Survived        0
          Pclass          0
          Age             0
          SibSp           0
          Parch           0
          Fare            0
          dtype: int64
```

# Task 5 Report
## Time Series Analysis with Pandas

❖ **Weekly Resampled Data:**



❖ **Time Series Analysis:**



❖ **Processed Dataset Exported as 'processed_time_series.csv':**

# Task 5 Report
## Time Series Analysis with Pandas

### 3. Describe Widget/Algorithm Used in Task

**Algorithm Used:**

The task leverages the following algorithms and techniques:

❖ **Data Preprocessing**:

- **Datetime Simulation**: A synthetic datetime column is added to mimic time-series data, and the column is set as the index to enable time-based operations.
- **Handling Missing Values**: Missing values in numeric columns are replaced with the column mean for smoother analysis.

❖ **Resampling**:

- The numeric data is resampled to a weekly frequency using Pandas' resample method, calculating weekly averages to understand long-term trends.

❖ **Rolling Statistics**:

- **Rolling Mean**: A 7-day rolling window is applied to calculate the average value over the last 7 days, smoothing short-term fluctuations.
- **Rolling Standard Deviation**: A 7-day rolling window is applied to measure variability in the data.

❖ **Lagged Features**:

- **Lagging**: Shifted versions of the data are created (e.g., lagged by 1 day or 7 days) to capture time-dependent patterns.

❖ **Visualization**:

- **Line Plots**: Rolling mean, rolling standard deviation, and the original data are plotted to visually compare trends and variability over time.

**Tools Used:**

❖ **Pandas**:
\
- For data manipulation, resampling, rolling statistics, and handling missing values.
- Used the resample, rolling, and shift functions to perform time-series-specific operations.

❖ **NumPy**:

- Assisted in handling numerical operations while imputing missing values.

❖ **Matplotlib**:

- For creating visualizations such as line plots to showcase trends, rolling statistics, and lagged features.

❖ **Jupyter/Colab Environment**:

- Provided an interactive Python environment for implementing the time-series analysis and visualizing results.

**\*\*\* The End \*\*\***