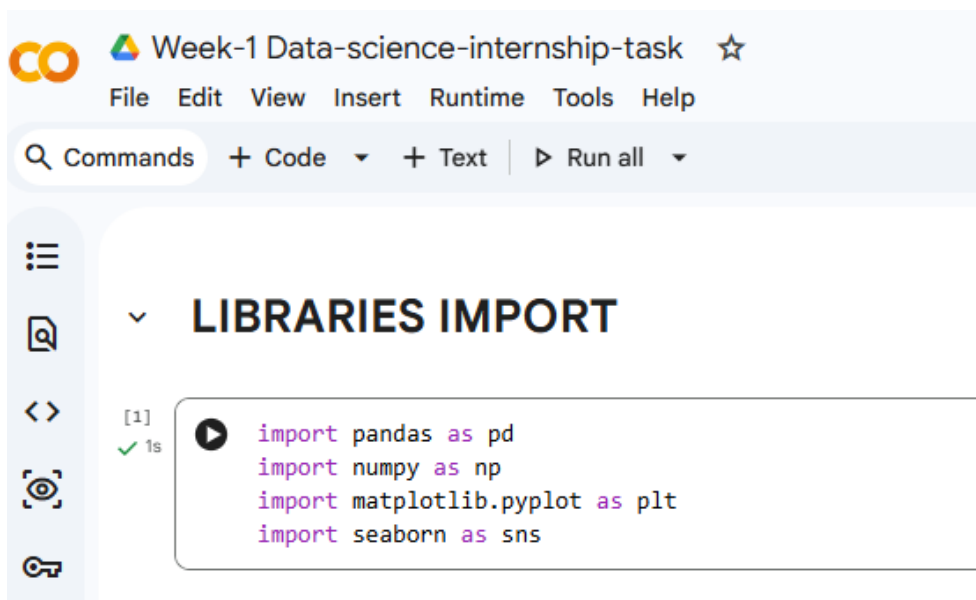
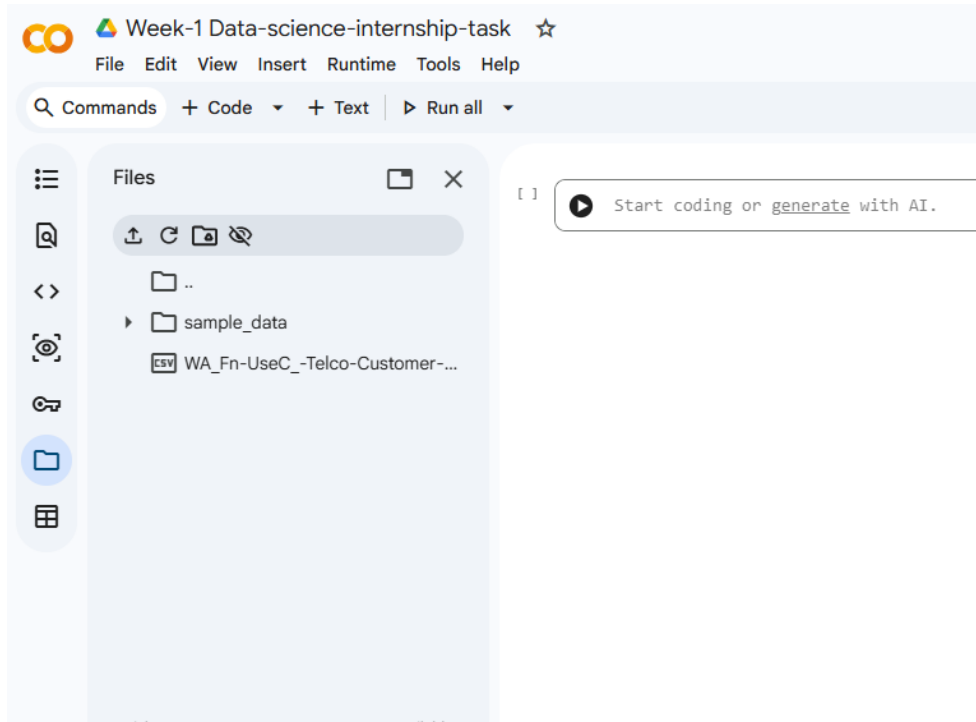


**Introduction:** This report documents the work completed for assignment in the Data Science Internship. The objective of this task was to perform advanced data cleaning and preprocessing on a real-world business dataset. The focus was on identifying data quality issues and applying professional techniques to make the dataset ready for analysis.



**DATASET LOAD**

```
[2] df = pd.read_csv("WA_Fn-UseC_-Telco-Customer-Churn.csv")
df.head()
```

	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetService	OnlineSecurity	...	DeviceProtection	T...
0	7590-VHVEG	Female	0	Yes	No	1	No	No phone service	DSL	No	...	No	
1	5575-GNVDE	Male	0	No	No	34	Yes	No	DSL	Yes	...	Yes	
2	3668-QPYBK	Male	0	No	No	2	Yes	No	DSL	Yes	...	No	
3	7795-CFOCW	Male	0	No	No	45	No	No phone service	DSL	Yes	...	Yes	
4	9237-HQITU	Female	0	No	No	2	Yes	No	Fiber optic	No	...	No	

Variables Terminal 23:21 Python 3

**Dataset Description:** The Telco Customer Churn dataset was selected for this assignment. It contains detailed information about telecom customers including demographic details, billing information, and service usage. The dataset included several missing values, inconsistent formats, and potential outliers, making it suitable for this cleaning task.

Week-1 Data-science-internship-task ☆ Unsaved changes since 17:35

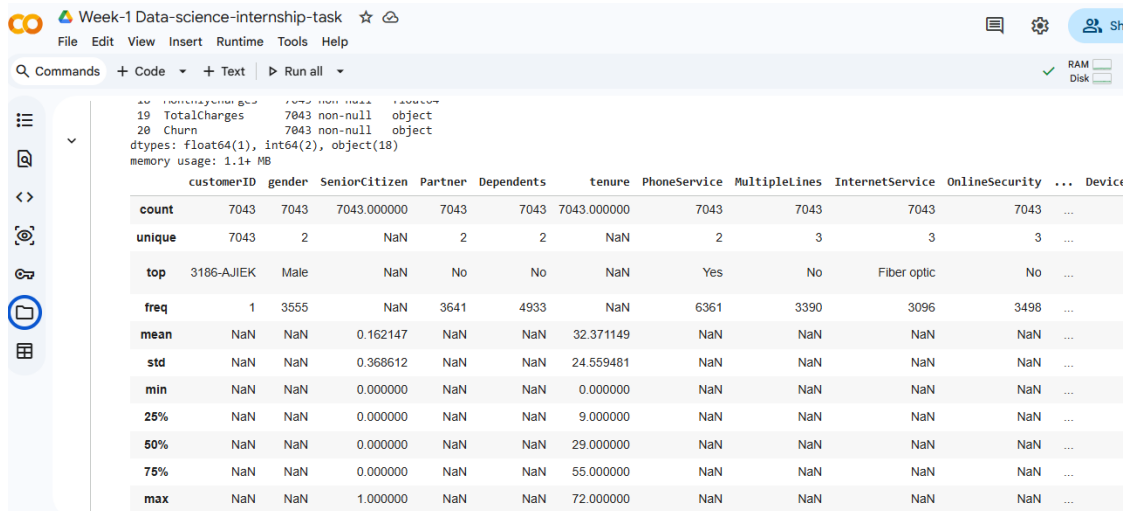
File Edit View Insert Runtime Tools Help

Commands + Code + Text ▶ Run all ▼

```
[6] df.shape
df.info()
df.describe(include='all')
```

<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 7043 entries, 0 to 7042  
Data columns (total 21 columns):

#	Column	Non-Null Count	Dtype
0	customerID	7043 non-null	object
1	gender	7043 non-null	object
2	SeniorCitizen	7043 non-null	int64
3	Partner	7043 non-null	object
4	Dependents	7043 non-null	object
5	tenure	7043 non-null	int64
6	PhoneService	7043 non-null	object
7	MultipleLines	7043 non-null	object
8	InternetService	7043 non-null	object
9	OnlineSecurity	7043 non-null	object
10	OnlineBackup	7043 non-null	object
11	DeviceProtection	7043 non-null	object
12	TechSupport	7043 non-null	object
13	StreamingTV	7043 non-null	object
14	StreamingMovies	7043 non-null	object
15	Contract	7043 non-null	object
16	PaperlessBilling	7043 non-null	object



Week-1 Data-science-internship-task

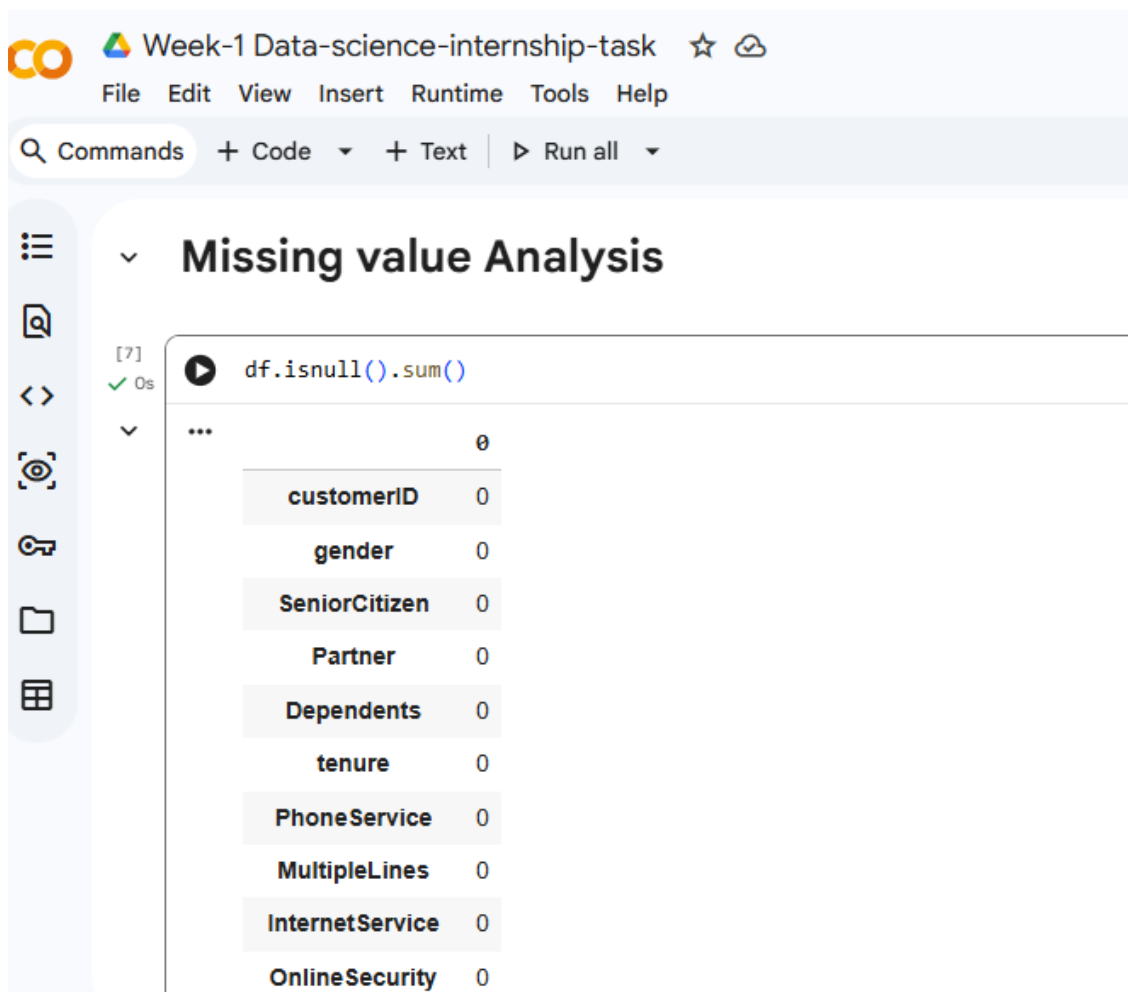
File Edit View Insert Runtime Tools Help

Q Commands + Code + Text ▶ Run all

memory usage: 1.1+ MB

	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetService	OnlineSecurity	...	Device
count	7043	7043	7043.000000	7043	7043	7043.000000	7043	7043	7043	7043	...	...
unique	7043	2	NaN	2	2	NaN	2	3	3	3	...	...
top	3186-AJIEK	Male	NaN	No	No	NaN	Yes	No	Fiber optic	No	...	...
freq	1	3555	NaN	3641	4933	NaN	6361	3390	3096	3498	...	...
mean	NaN	NaN	0.162147	NaN	NaN	32.371149	NaN	NaN	NaN	NaN	...	...
std	NaN	NaN	0.368612	NaN	NaN	24.559481	NaN	NaN	NaN	NaN	...	...
min	NaN	NaN	0.000000	NaN	NaN	0.000000	NaN	NaN	NaN	NaN	...	...
25%	NaN	NaN	0.000000	NaN	NaN	9.000000	NaN	NaN	NaN	NaN	...	...
50%	NaN	NaN	0.000000	NaN	NaN	29.000000	NaN	NaN	NaN	NaN	...	...
75%	NaN	NaN	0.000000	NaN	NaN	55.000000	NaN	NaN	NaN	NaN	...	...
max	NaN	NaN	1.000000	NaN	NaN	72.000000	NaN	NaN	NaN	NaN	...	...

**Missing Value Analysis:** An in-depth missing value analysis was conducted to understand which columns required attention. Both tabular summaries and visualization techniques were used to analyze the distribution of missing data.



Week-1 Data-science-internship-task

File Edit View Insert Runtime Tools Help

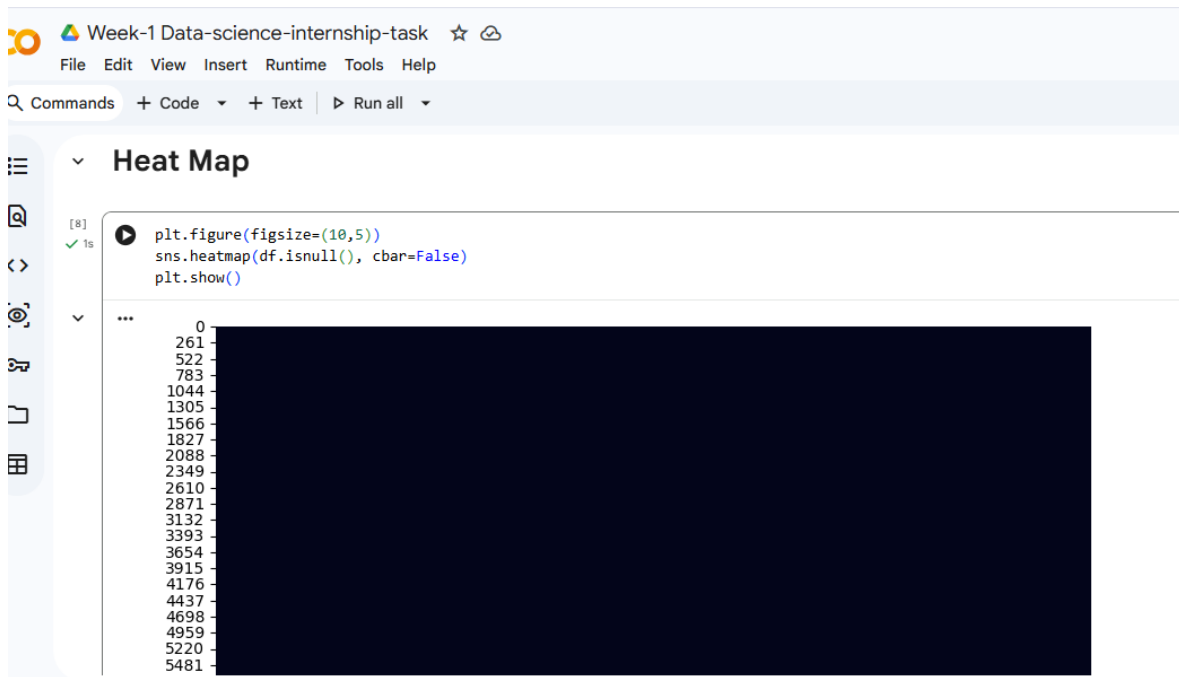
Q Commands + Code + Text ▶ Run all

## Missing value Analysis

[7] ✓ 0s

df.isnull().sum()

...	0
customerID	0
gender	0
SeniorCitizen	0
Partner	0
Dependents	0
tenure	0
PhoneService	0
MultipleLines	0
InternetService	0
OnlineSecurity	0



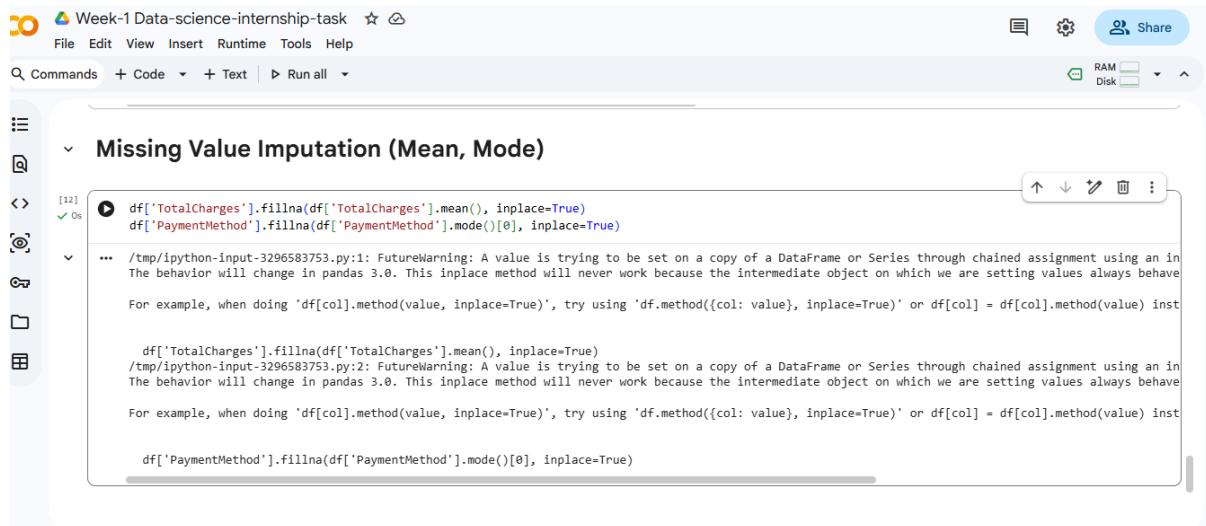
**Data Type Correction:** During preprocessing, the TotalCharges column was found to contain non-numeric values stored as text. This issue was corrected by converting the column into proper numeric format to ensure consistency and accuracy.

Fix Data Type

```
[11] df['TotalCharges'] = pd.to_numeric(df['TotalCharges'], errors='coerce')
```

	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetService	OnlineSecurity	...	DeviceProtection
0	7590-VHVEG	Female	0	Yes	No	1	No	No phone service	DSL	No	...	No
1	5575-GNVDE	Male	0	No	No	34	Yes	No	DSL	Yes	...	Yes
2	3668-QPYBK	Male	0	No	No	2	Yes	No	DSL	Yes	...	No
3	7795-CFOCW	Male	0	No	No	45	No	No phone service	DSL	Yes	...	Yes
4	9237-HQITU	Female	0	No	No	2	Yes	No	Fiber optic	No	...	No
...	...	...	...	...	...	...	...	...	...	...	...	...
7590	6840-	...	...	...	...	...	...	...	DSL	Yes	...	Yes

**Missing Value Imputation:** Two different imputation techniques were applied. First, statistical imputation using mean and mode was used for appropriate columns. Then, KNN Imputation was implemented as an advanced technique. After imputation, the dataset was verified to ensure no missing values remained.



The screenshot shows a Jupyter Notebook titled "Week-1 Data-science-internship-task". The code cell [12] contains the following Python code:

```
df['TotalCharges'].fillna(df['TotalCharges'].mean(), inplace=True)
df['PaymentMethod'].fillna(df['PaymentMethod'].mode()[0], inplace=True)
```

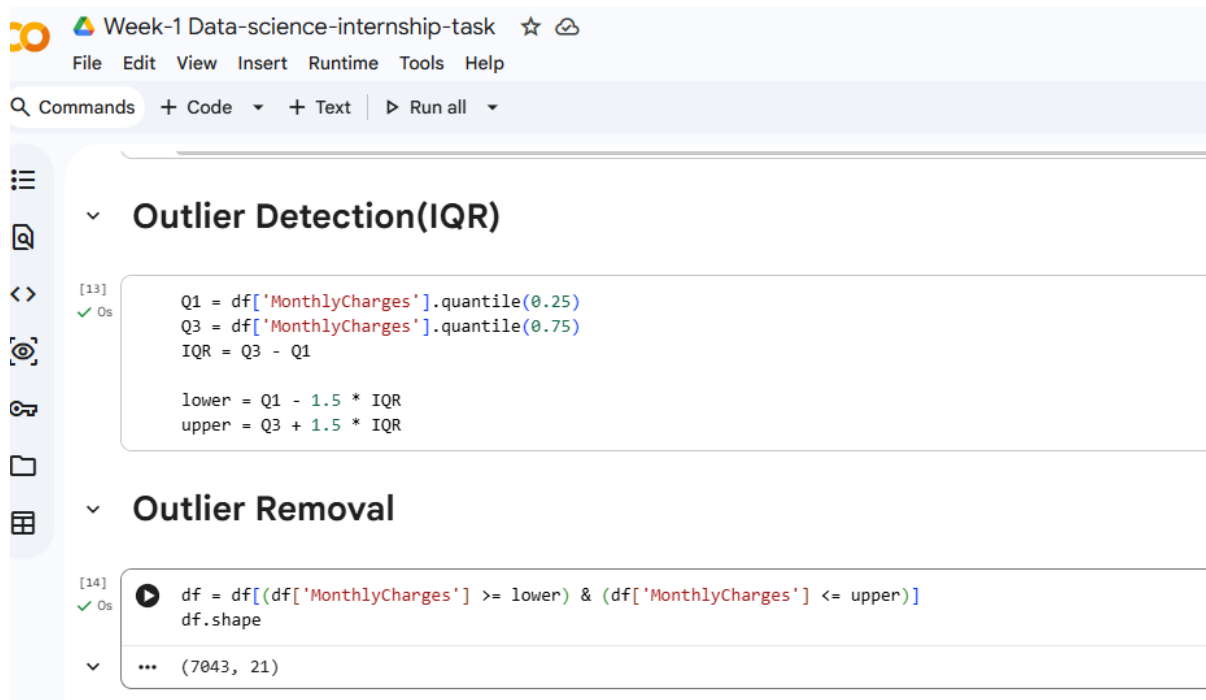
Below the code, there are two warning messages from pandas:

```
/tmp/ipython-input-3296583753.py:1: FutureWarning: A value is trying to be set on a copy of a DataFrame or Series through chained assignment using an inplace method. The behavior will change in pandas 3.0. This inplace method will never work because the intermediate object on which we are setting values always behaves like a copy. For example, when doing 'df[col].method(value, inplace=True)', try using 'df.method({col: value}, inplace=True)' or df[col] = df[col].method(value) instead.
```

The code cell [13] contains the following Python code:

```
df['TotalCharges'].fillna(df['TotalCharges'].mean(), inplace=True)
df['PaymentMethod'].fillna(df['PaymentMethod'].mode()[0], inplace=True)
```

**Outlier Detection & Treatment:** Outliers in the MonthlyCharges column were detected using the IQR method. Identified extreme values were treated through documented filtering to improve the reliability of the dataset.



The screenshot shows a Jupyter Notebook titled "Week-1 Data-science-internship-task". The code cell [13] contains the following Python code:

```
Q1 = df['MonthlyCharges'].quantile(0.25)
Q3 = df['MonthlyCharges'].quantile(0.75)
IQR = Q3 - Q1

lower = Q1 - 1.5 * IQR
upper = Q3 + 1.5 * IQR
```

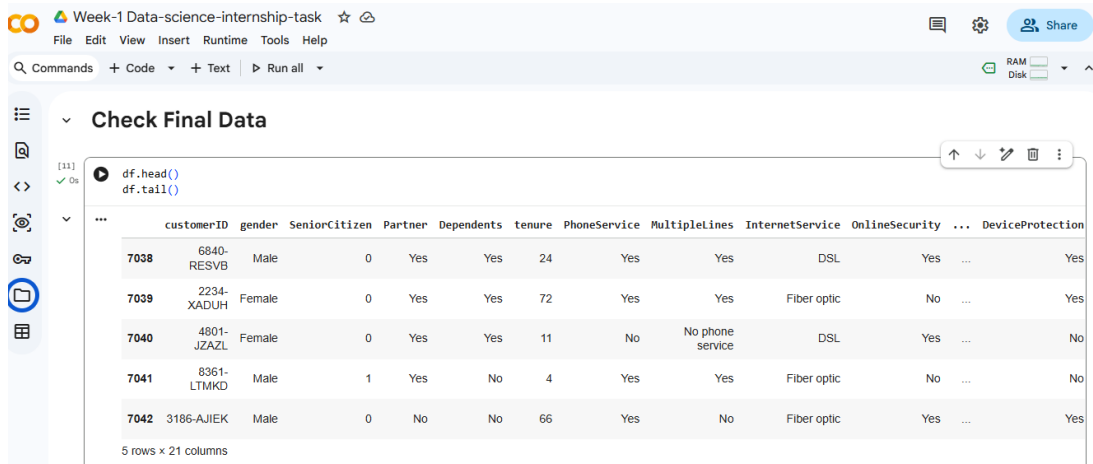
The code cell [14] contains the following Python code:

```
df = df[(df['MonthlyCharges'] >= lower) & (df['MonthlyCharges'] <= upper)]
df.shape
```

Below the code, the output of the code cell [14] is shown:

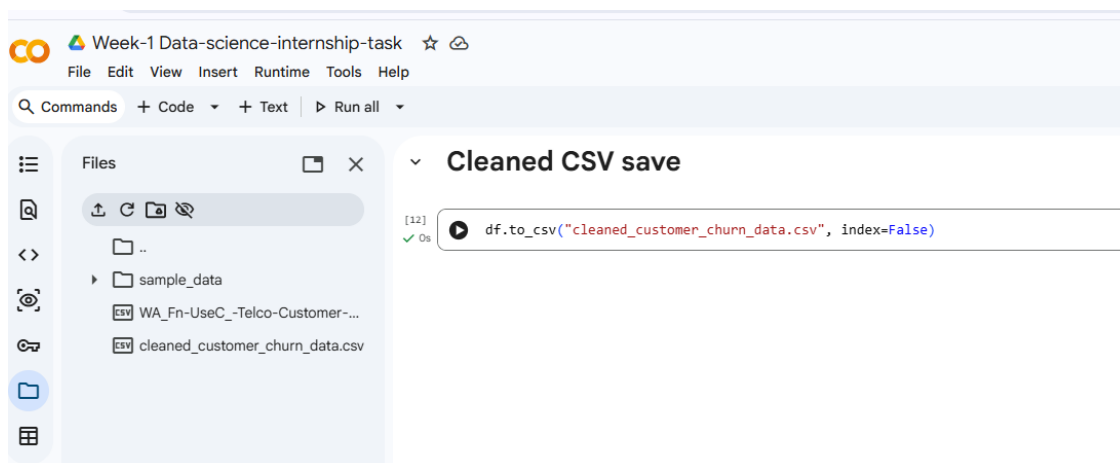
```
(7043, 21)
```

**Final Dataset Export:**



	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetService	OnlineSecurity	...	DeviceProtection
7038	6840-RESVB	Male	0	Yes	Yes	24	Yes	Yes	DSL	Yes	...	Yes
7039	2234-XADUH	Female	0	Yes	Yes	72	Yes	Yes	Fiber optic	No	...	Yes
7040	4801-JZAZL	Female	0	Yes	Yes	11	No	No phone service	DSL	Yes	...	No
7041	8361-LTMKD	Male	1	Yes	No	4	Yes	Yes	Fiber optic	No	...	No
7042	3186-AJIEK	Male	0	No	No	66	Yes	No	Fiber optic	Yes	...	Yes

5 rows x 21 columns



```
df.to_csv("cleaned_customer_churn_data.csv", index=False)
```

cleaned\_customer\_churn\_data.csv

**Conclusion:** This assignment helped me understand how real-world business data often contains many quality issues. By applying Python-based preprocessing techniques, I was able to convert messy data into a clean and structured format. The final dataset is now consistent and ready for further analytics and machine learning tasks.