# PHASE - 3

# Development Part 1

| DATE | 23 - 10 - 2023 |
|---|---|
| TEAM ID | 8939 |
| PROJECT NAME | 8301 – CUSTOMER CHURN PREDICTION |
| TEAM NAME | Proj_207142_Team_1 |

# ANALYTICS OBJECTIVES

## Data Preprocessing :

1. Data Inspection
2. Data Cleaning
3. Data Transformation
4. Data Splitting
5. Data Normalization
6. Data Validation
7. Data Visualization
8. Data Collection

# 1.Loading Data :

Use pandas.read_csv() to load data from a CSV file.

Use pandas.read_excel() for Excel files.

```python
In [1]:

import pandas as pd
df=pd.read_csv("E:/Churn.csv")
```

# 2. Exploring Data:

Use df.head() to view the first few rows of the dataset.

In [2]: df.head()

Out[2]:

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

5 rows × 21 columns

Use df.info() to get information about data types and missing values.

Use df.describe() for summary statistics.python



```
In [3]: df.info()

<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
 17  PaymentMethod     7043 non-null   object
 18  MonthlyCharges    7043 non-null   float64
 19  TotalCharges      7043 non-null   object
 20  Churn             7043 non-null   object
dtypes: float64(1), int64(2), object(18)
memory usage: 1.1+ MB
```

```
In [4]: df.describe()
```
Out[4]:

|       | SeniorCitizen | tenure       | MonthlyCharges |
|-------|---------------|--------------|----------------|
| count | 7043.000000   | 7043.000000  | 7043.000000    |
| mean  | 0.162147      | 32.371149    | 64.761692      |
| std   | 0.368612      | 24.559481    | 30.090047      |
| min   | 0.000000      | 0.000000     | 18.250000      |
| 25%   | 0.000000      | 9.000000     | 35.500000      |
| 50%   | 0.000000      | 29.000000    | 70.350000      |
| 75%   | 0.000000      | 55.000000    | 89.850000      |
| max   | 1.000000      | 72.000000    | 118.750000     |

# 3.Handling Missing Values:

Use df.isnull() to identify missing values.

Use df.fillna() or df.dropna() to handle missing values.

```
In [7]: df.isnull()
```

Out[7]:

| | customerID | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | ... | DeviceProtection | TechS |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | False | False | False | False | False | False | False | False | False | False | ... | False | |
| 1 | False | False | False | False | False | False | False | False | False | False | ... | False | |
| 2 | False | False | False | False | False | False | False | False | False | False | ... | False | |
| 3 | False | False | False | False | False | False | False | False | False | False | ... | False | |
| 4 | False | False | False | False | False | False | False | False | False | False | ... | False | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 7038 | False | False | False | False | False | False | False | False | False | False | ... | False | |
| 7039 | False | False | False | False | False | False | False | False | False | False | ... | False | |
| 7040 | False | False | False | False | False | False | False | False | False | False | ... | False | |
| 7041 | False | False | False | False | False | False | False | False | False | False | ... | False | |
| 7042 | False | False | False | False | False | False | False | False | False | False | ... | False | |

7043 rows × 21 columns

# 4.Data Cleaning:

Remove duplicate rows with df.drop_duplicates().

Rename columns using df.rename() if necessary.Convert data types with df.astype().
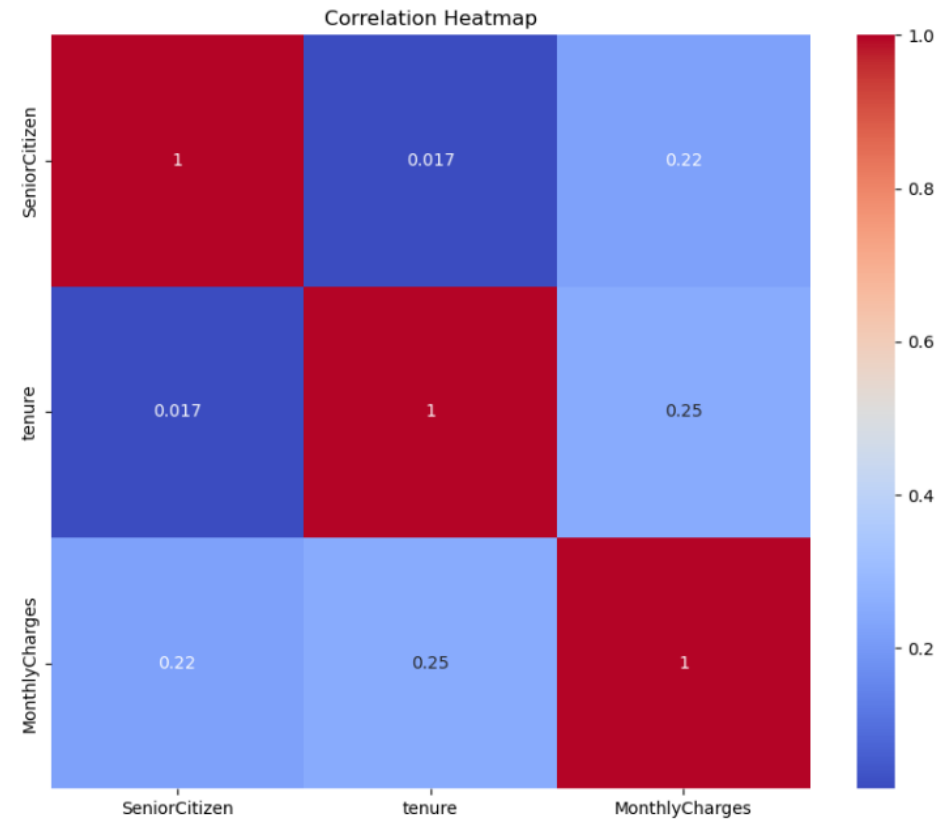
# 5.Handling Outliers:

Detect and deal with outliers using statistical methods or visualization.

You can use techniques like z-scores or IQR (Interquartile Range).



```
In [11]:  import seaborn as sns
          import matplotlib.pyplot as plt
          corr_matrix = df.corr()
          plt.figure(figsize=(10, 8))
          sns.heatmap(corr_matrix, annot=True, cmap="coolwarm")
          plt.title("Correlation Heatmap")
          plt.show()
```

# 6.Saving Data:

Save the preprocessed data back to a file if needed.

```
In [13]: df.to_csv('Downloads/preprocessed_churn.csv', index=False)
```

# IBM COGNOS ANALYTICS

# IN COGNOS - DATA LOADING

# IN COGNOS - DATA RELATIONSHIP
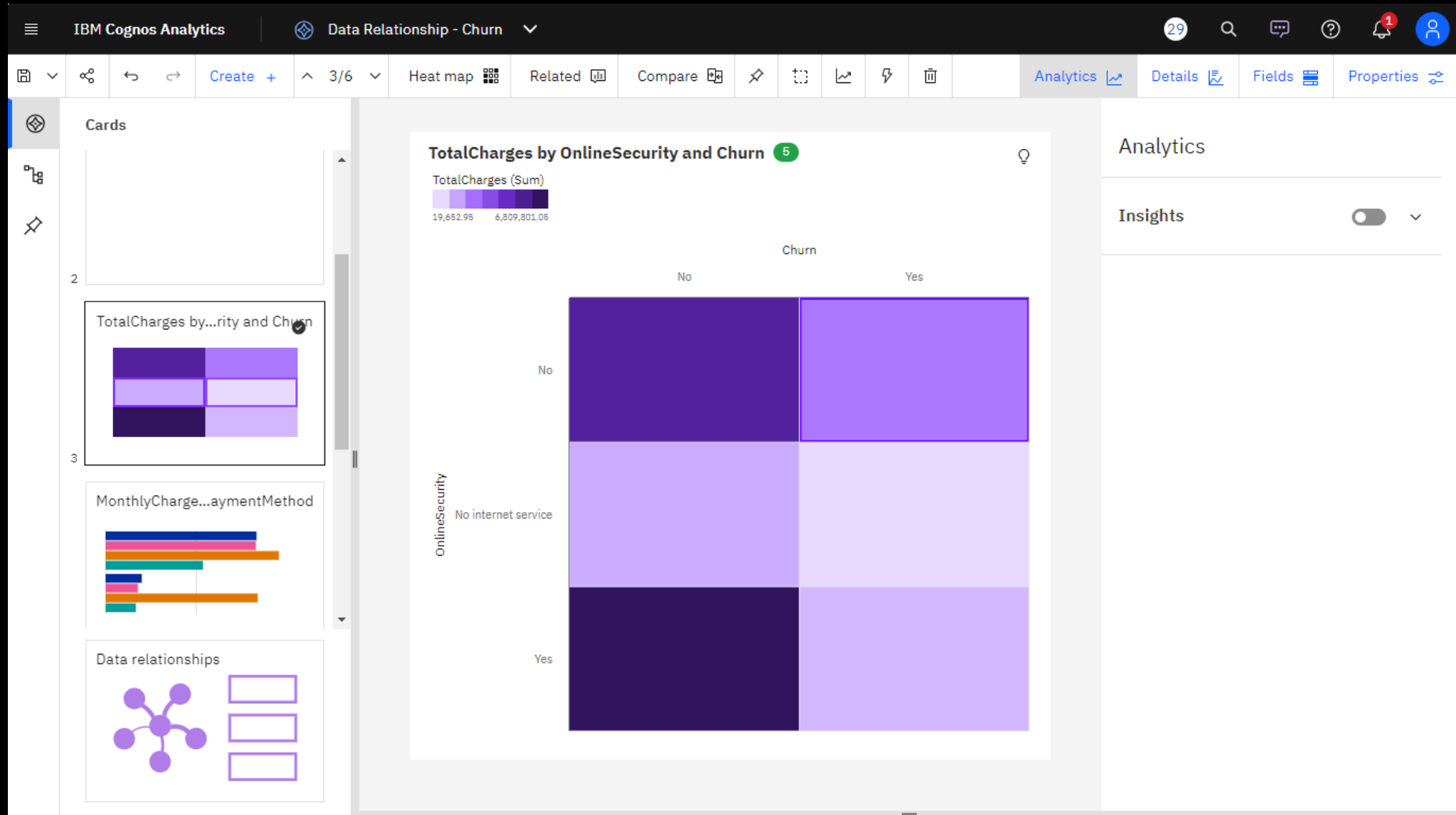
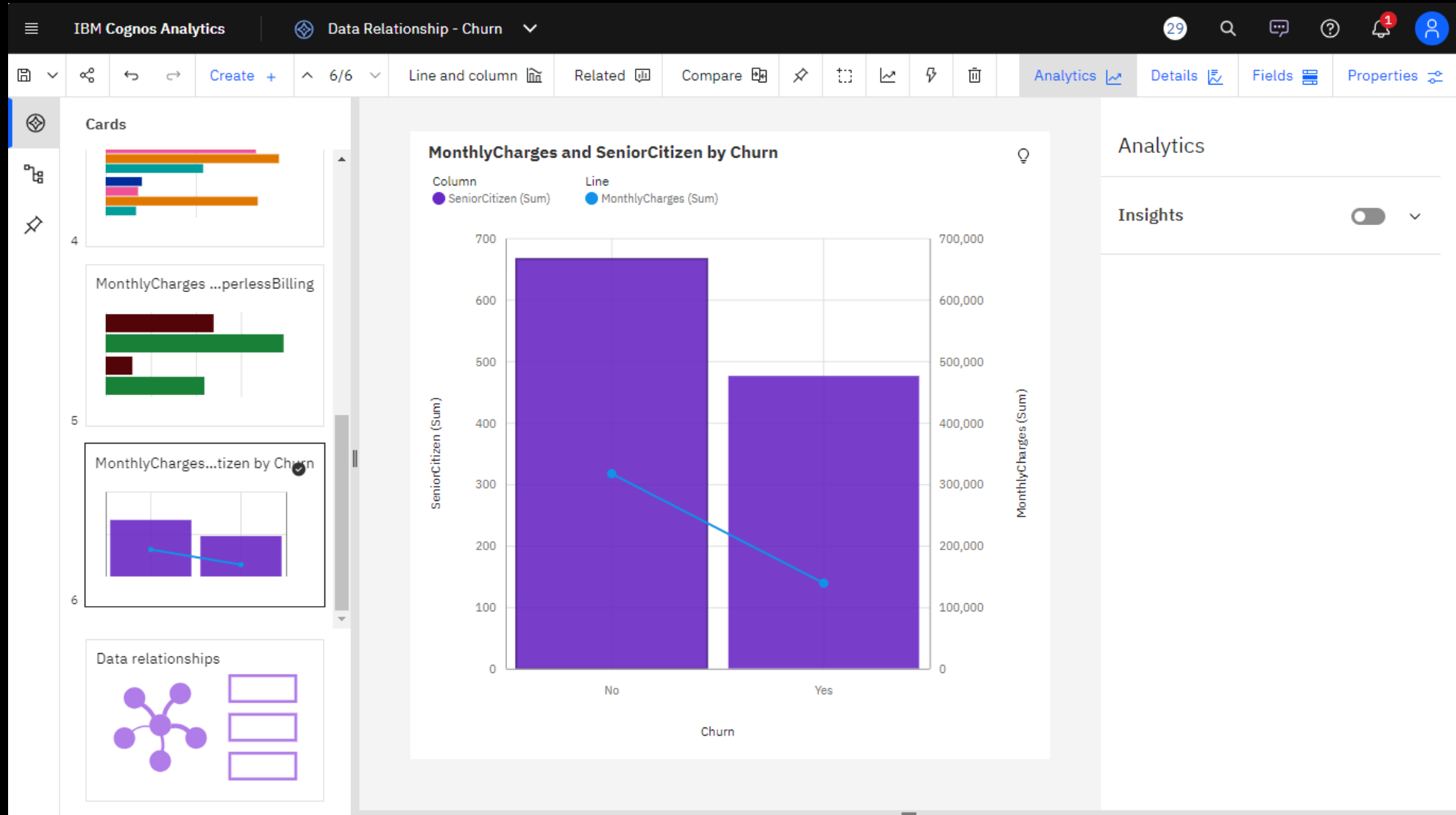# IN COGNOS - CHURN DASHBOARD

# IN COGNOS - TENURE AND MONTHLY CHARGES

# MONTHLY CHARGES BY CHURN COLORED BY PAYMENT METHOD

# TOTAL CHARGES BY ONLINE SECURITY

# SENIOR CITIZER BY CHURN

# VIEW OF TABLES

# CONCLUSION

Data preprocessing is a crucial step in preparing data for analysis and machine learning. It involves collecting, inspecting, cleaning, transforming, and organizing data. The main steps include data collection, inspection, cleaning, transformation, splitting, normalization, and validation. It ensures data is ready for analysis and model training, improving the success of data-related projects.