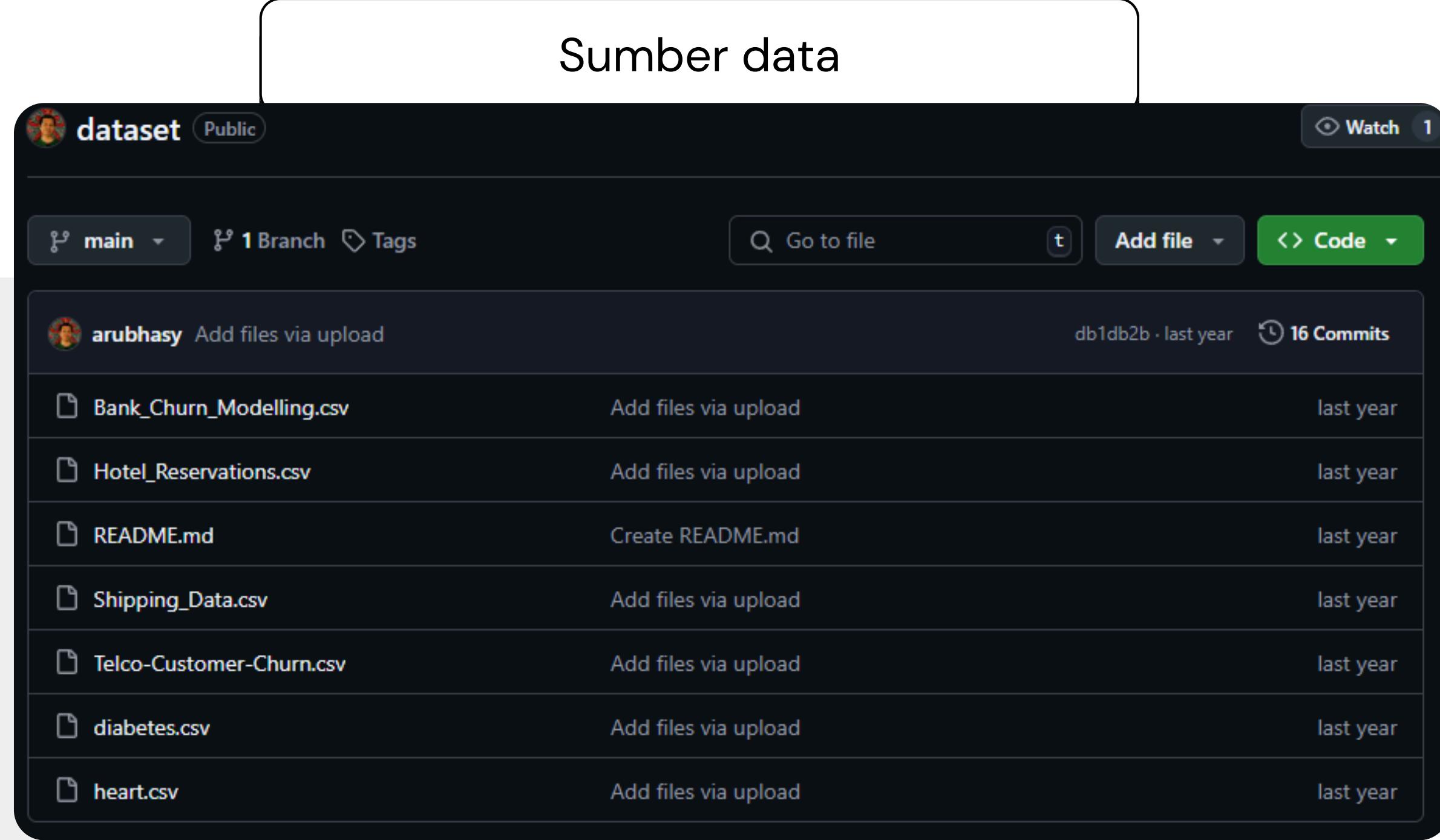


# ASSESSMENT PRESENTATION

Presented by  
**Sunday B. Putera Mandiri**

# MENGUMPULKAN DATA

Sumber data



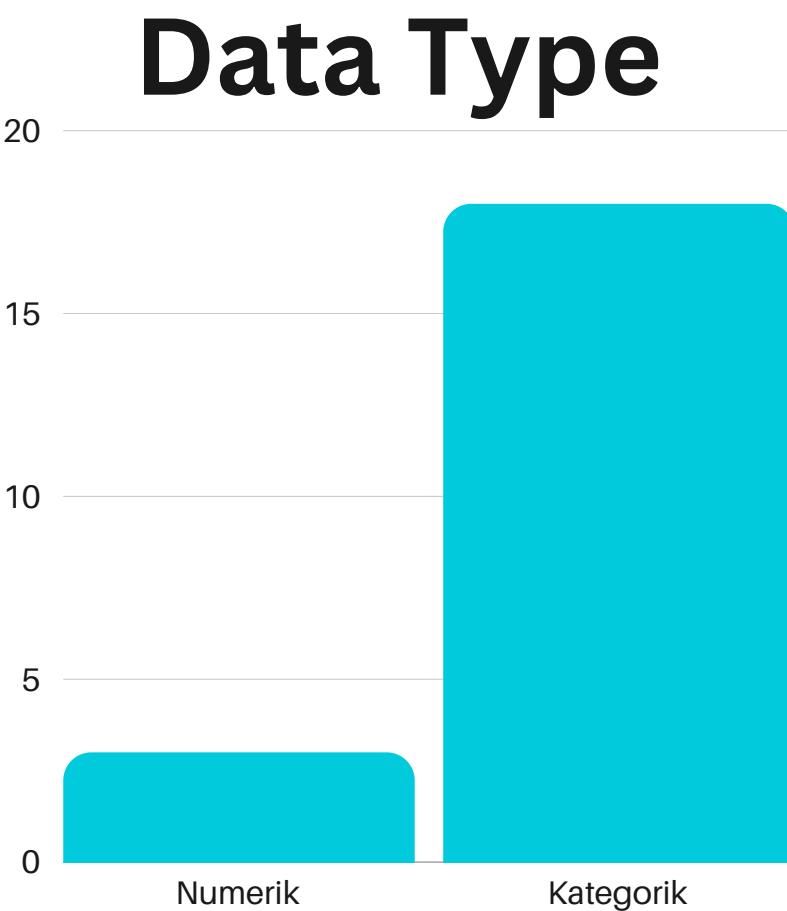
The screenshot shows a GitHub repository interface for a public dataset. The repository has 1 branch and 16 commits. It contains several CSV files: Bank\_Churn\_Modelling.csv, Hotel\_Reservations.csv, diabetes.csv, heart.csv, and Telco-Customer-Churn.csv. There is also a README.md file. The repository was last updated last year.

File	Action	Last Updated
arubhasy Add files via upload		last year
Bank_Churn_Modelling.csv	Add files via upload	last year
Hotel_Reservations.csv	Add files via upload	last year
README.md	Create README.md	last year
Shipping_Data.csv	Add files via upload	last year
Telco-Customer-Churn.csv	Add files via upload	last year
diabetes.csv	Add files via upload	last year
heart.csv	Add files via upload	last year



# MENELAAH DATA

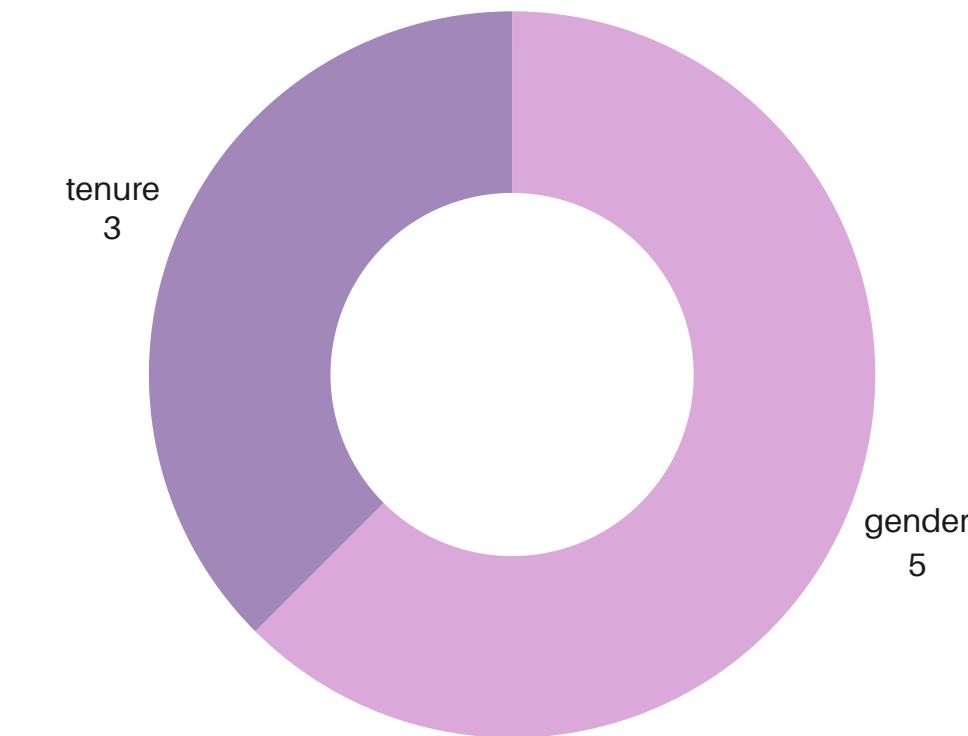
Sindhya Central Bank



## Zero Duplicate Data

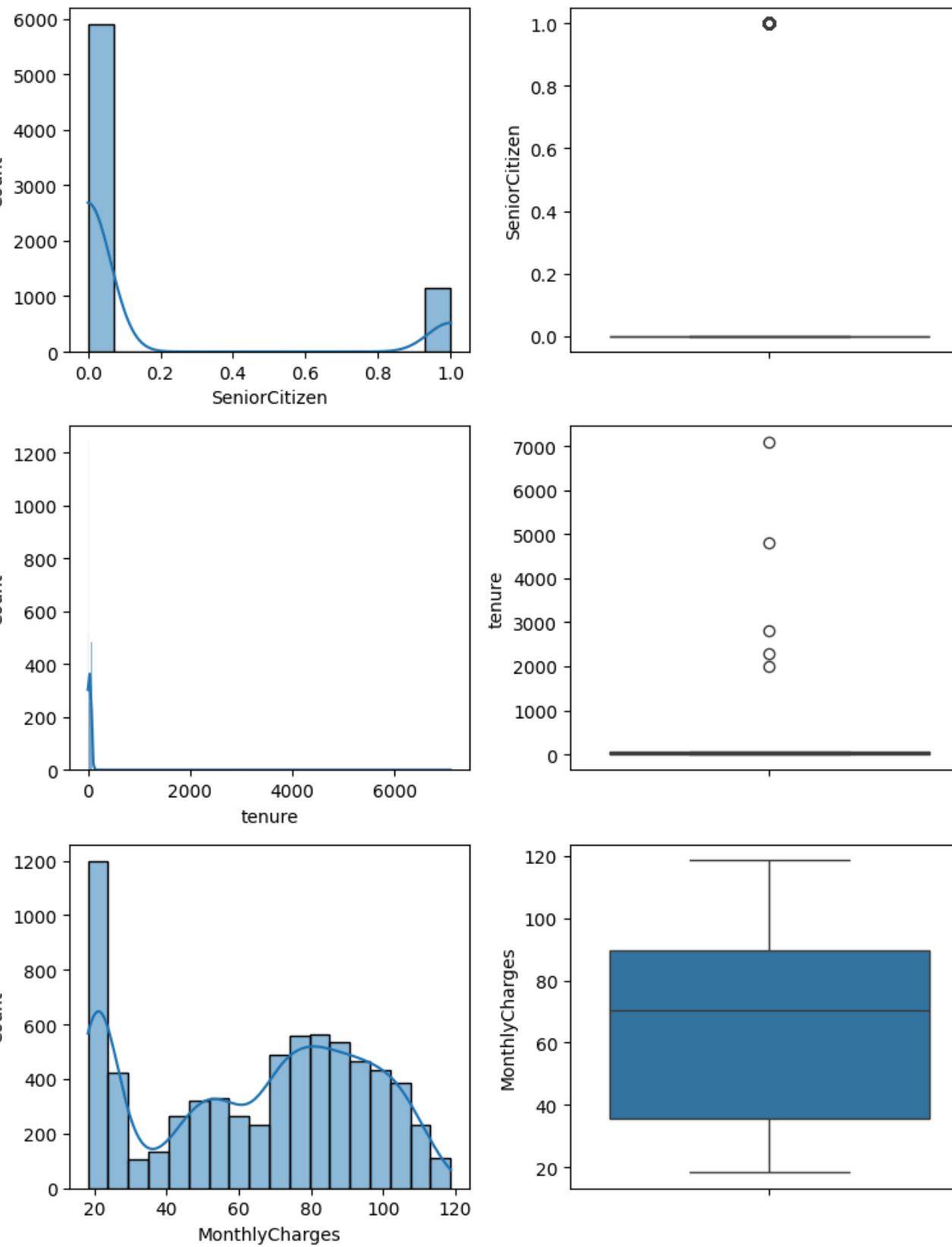


## Missing Value



# MENELAAH DATA

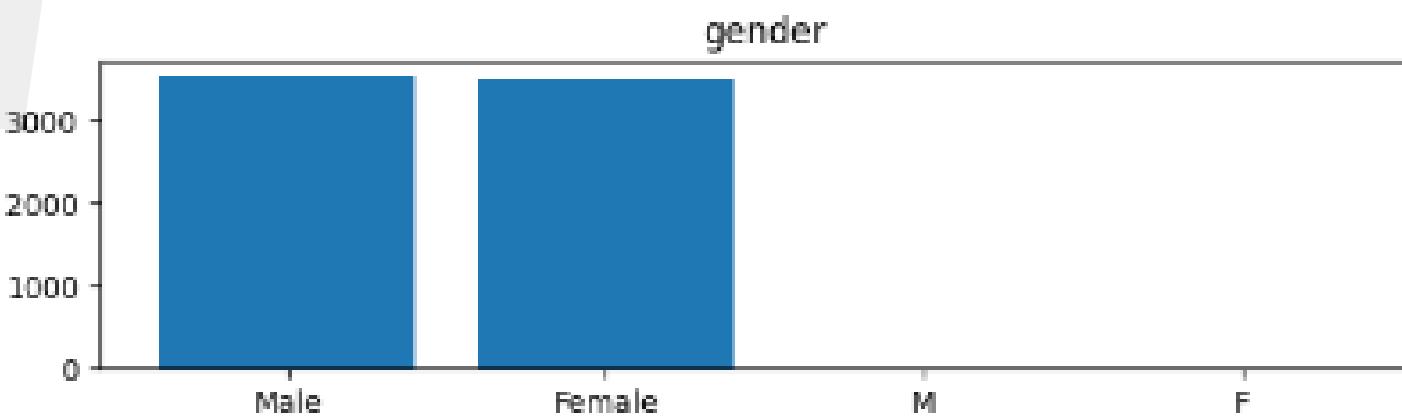
## Numerik



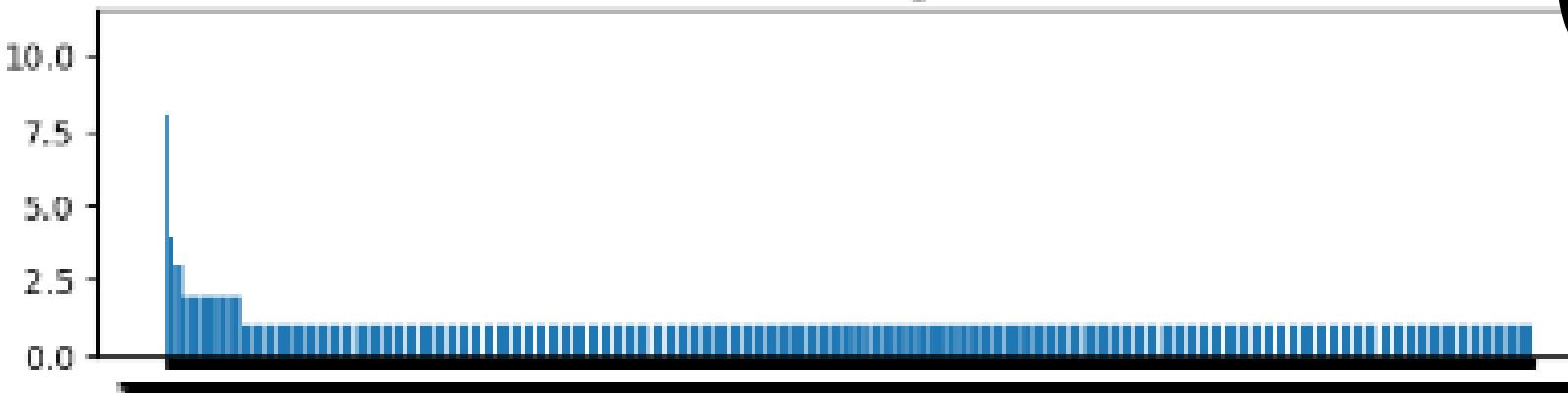
# MENELAAH DATA



Kategorik



TotalCharges



Kok  
bentuknya  
gitu sih?



# MENELAAH DATA

Kesimpulan



Terdapat Missing Value

Outliers

Terdapat variable dengan tipe data yang tidak  
sesuai

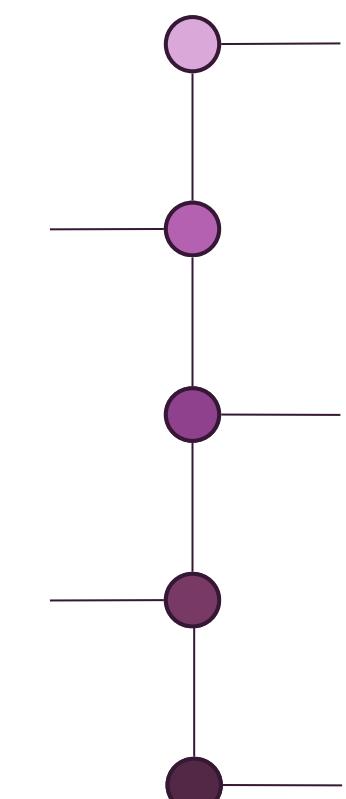
Inkonsisten Data

# Agenda



**Stage 2**  
mengatasi inkonsistensi data dengan menyamakan format pada gender

**Stage 4**  
transformasi atau rekonstruksi data pada MonthlyCharges untuk distribusi yang lebih seimbang (Optional)



**Stage 1**  
mengatasi missing value

**Stage 3**  
mengubah tipe data TotalCharges menjadi numerik

**Stage 5**  
Mengatasih string kosong pada TotalCharges

# DATA CLEANING

## Gender



### Menyamakan value

```
method replcae  
df['gender'] = df['gender'].replace('F',  
'Female')  
df['gender'] = df['gender'].replace('M',  
'Male')
```

### Mengisih missing value berdasarkan proporsi

```
fill_values = random.choices(['Male',  
'Female'], weights=[male_proportion,  
female_proportion], k=missing_count)  
df.loc[df['gender'].isnull(), 'gender'] =  
fill_values
```

### Output

	count
gender	
Male	3557
Female	3486
dtype: int64	

# DATA CLEANING

tenure

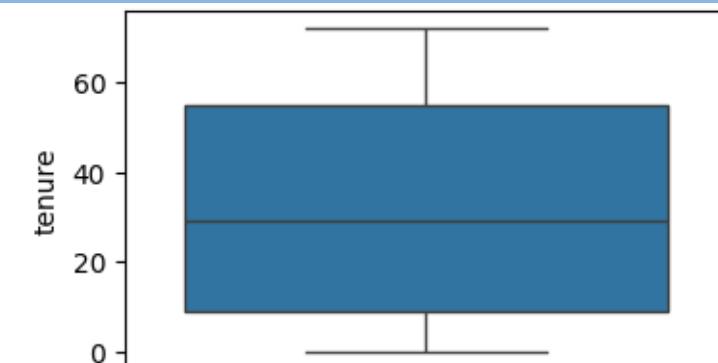
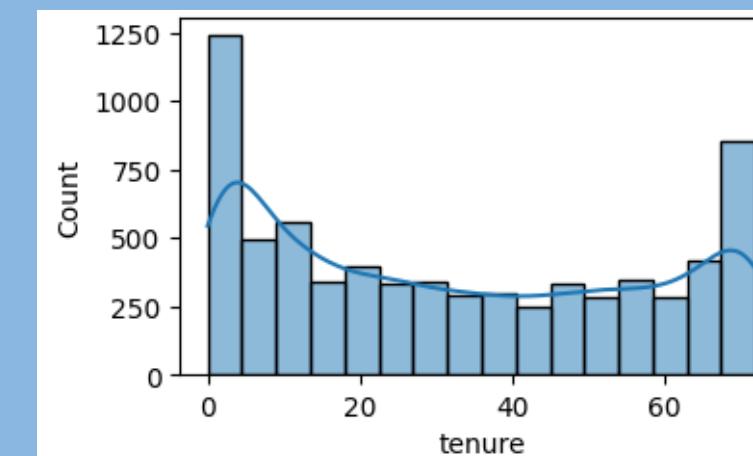
Mengisih missing value

```
df['tenure'] =  
df['tenure'].fillna(df['tenure'].mean())
```

Mengatasih outliers

```
# Menghitung percentile ke-95  
percintile_95 = df['tenure'].quantile(0.95)  
# Mengganti outliers dengan nilai percentile  
ke-95  
df['tenure'] = np.where(df['tenure'] >  
percintile_95, percintile_95, df['tenure'])
```

Output



# DATA CLEANING

TotalCharges



Mengubah tipe data

```
df['TotalCharges'] =  
pd.to_numeric(df['TotalCharges'],error  
s='coerce')
```

Mengisih missing value

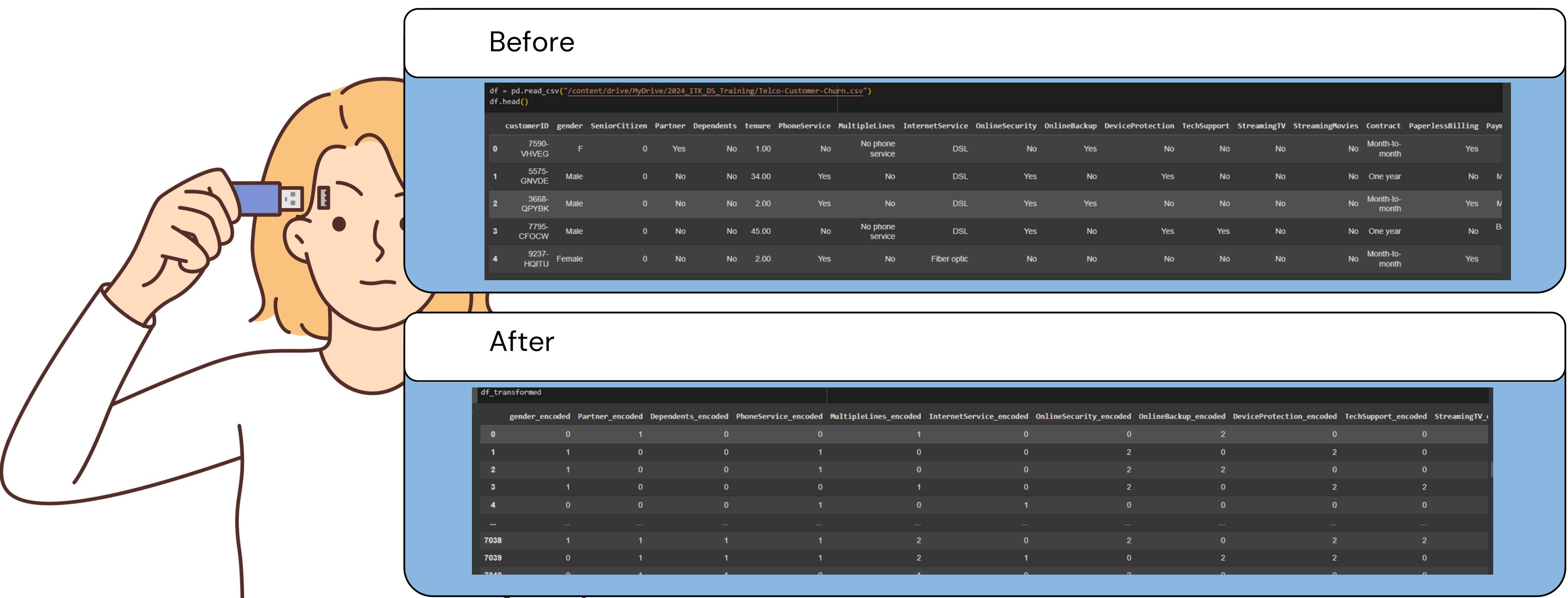
```
df['TotalCharges'] =  
df['TotalCharges'].fillna(df['TotalCharges'].me  
dian())
```

Output

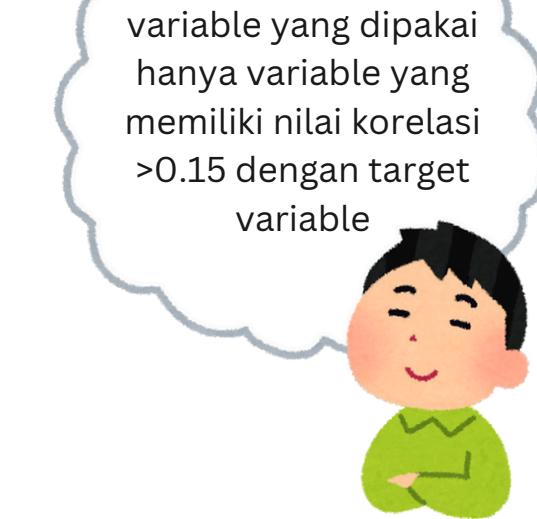
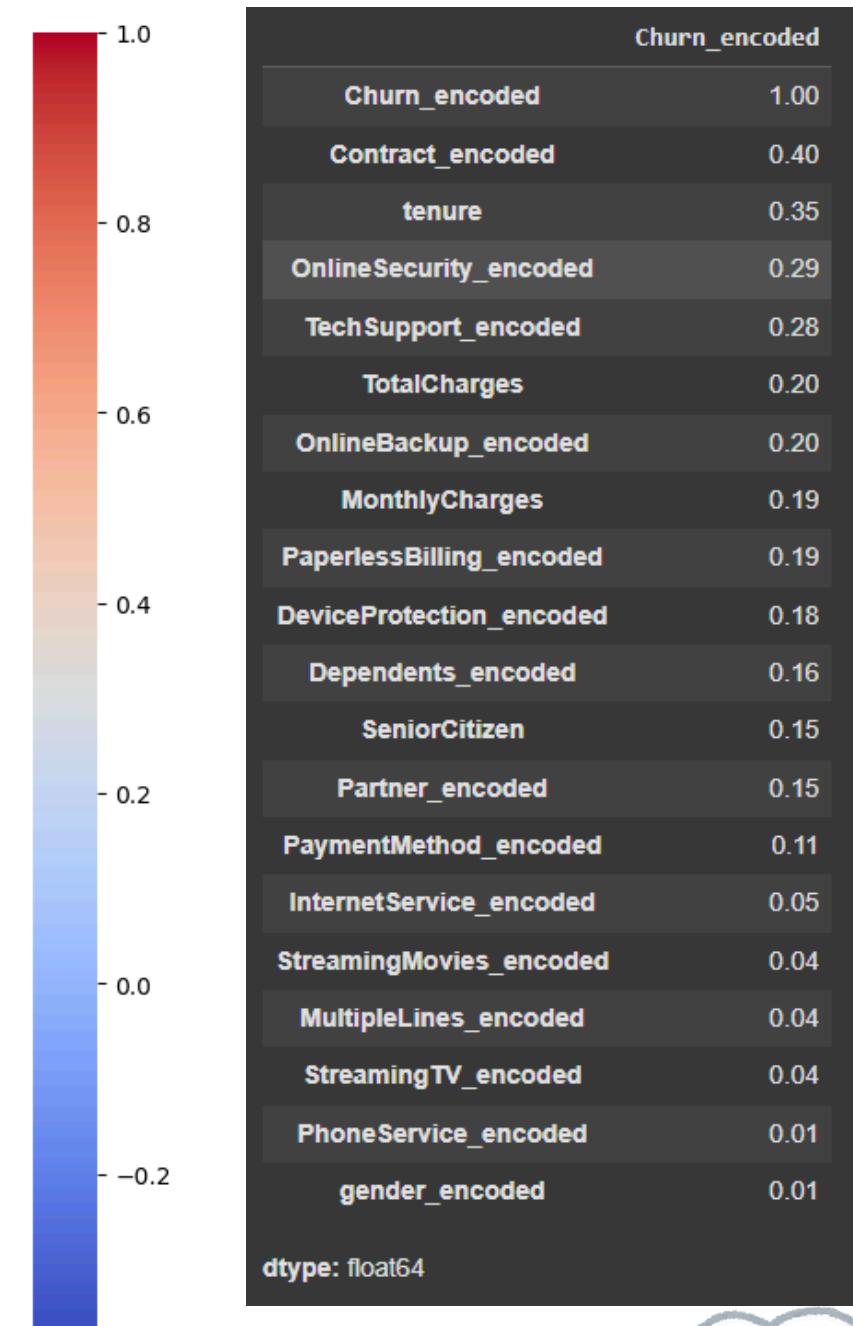
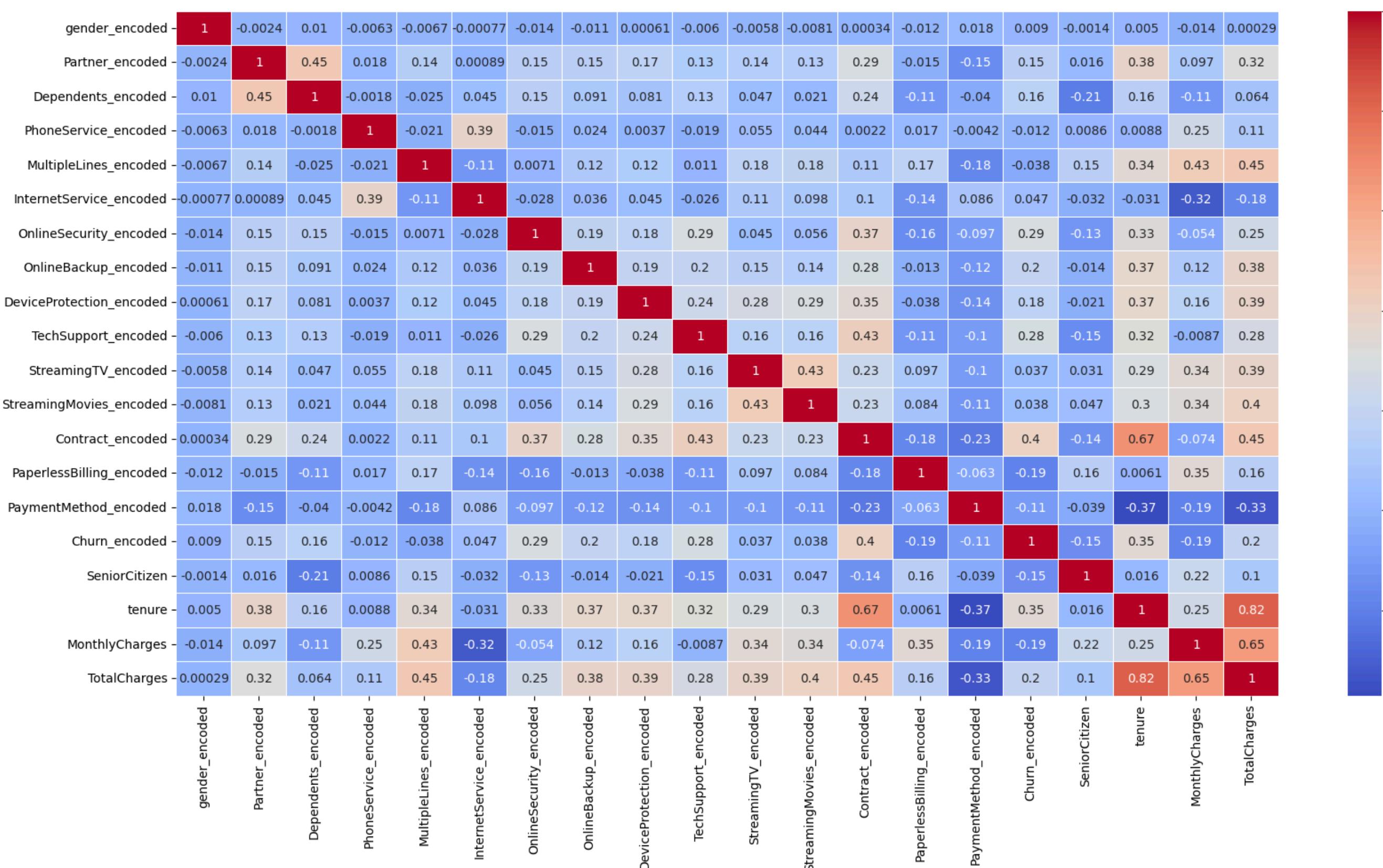
```
df['TotalCharges'].isnull().sum()
```

# TRANSFORMAS

Transformasi dilakukan menggunakan method LabelEncoder()

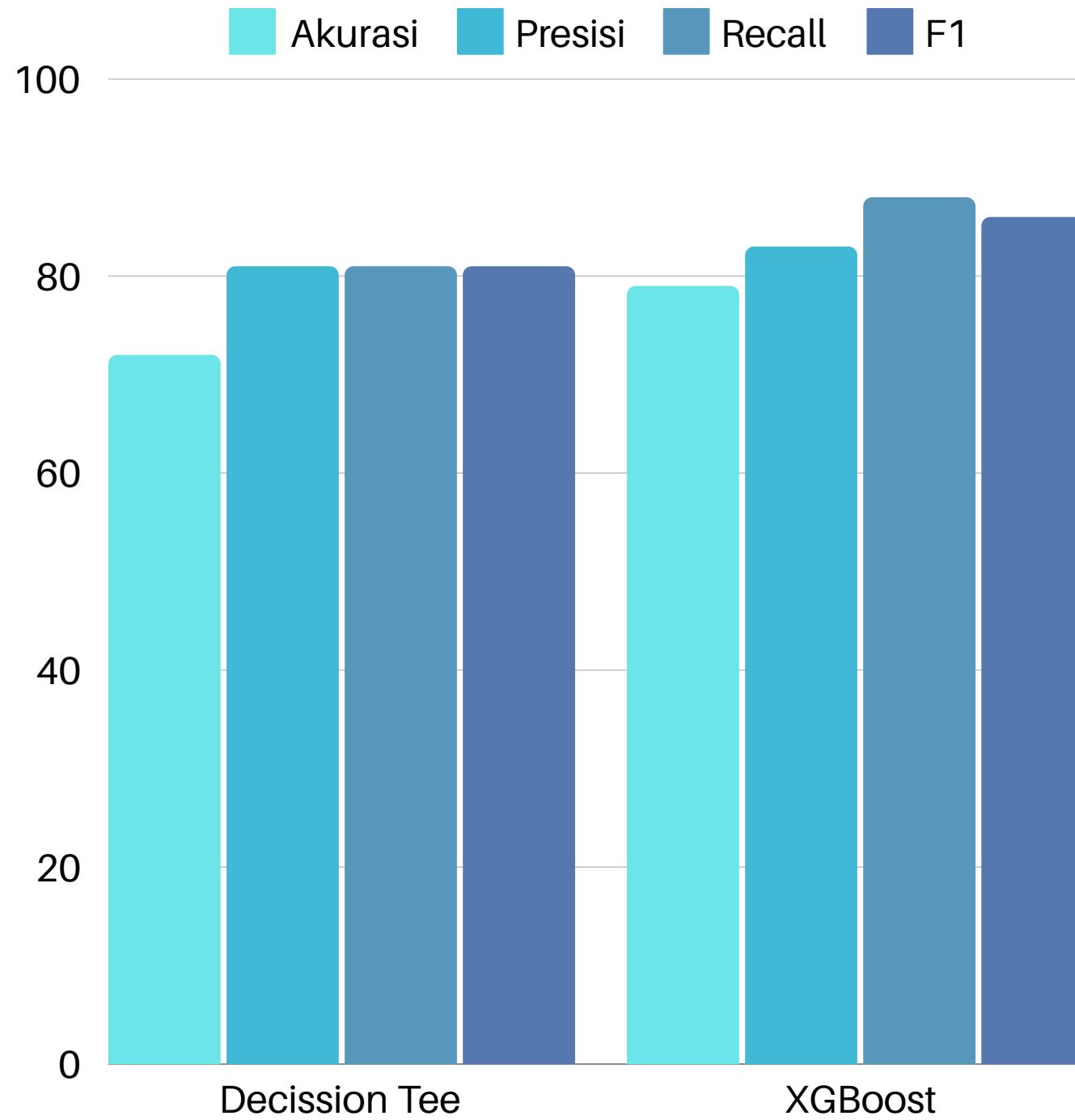


# KORELASI

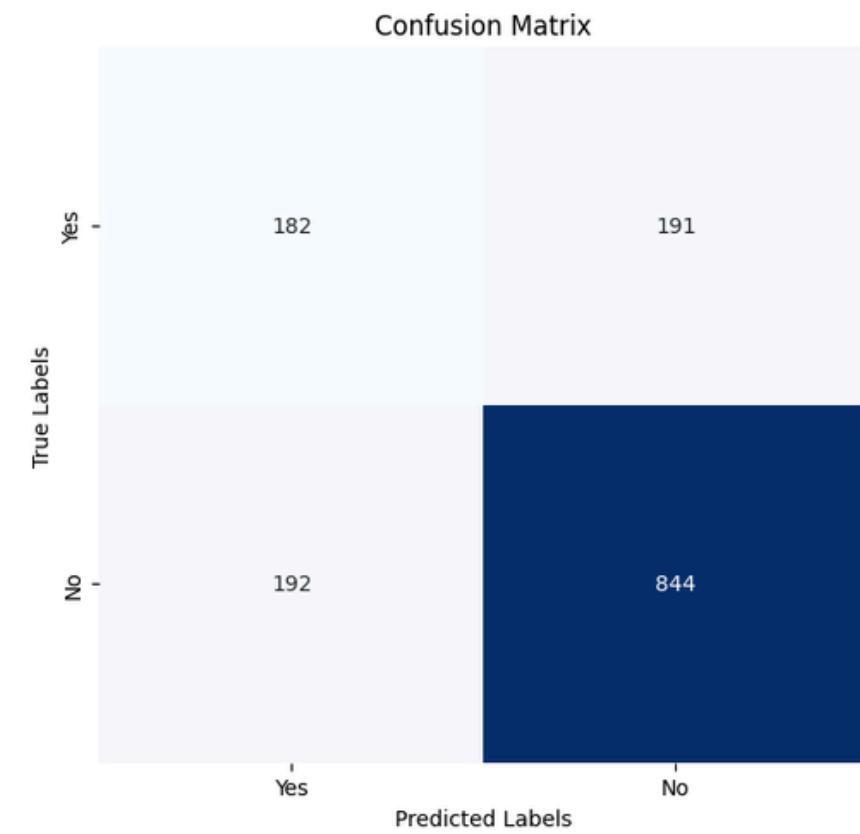


# MODEL

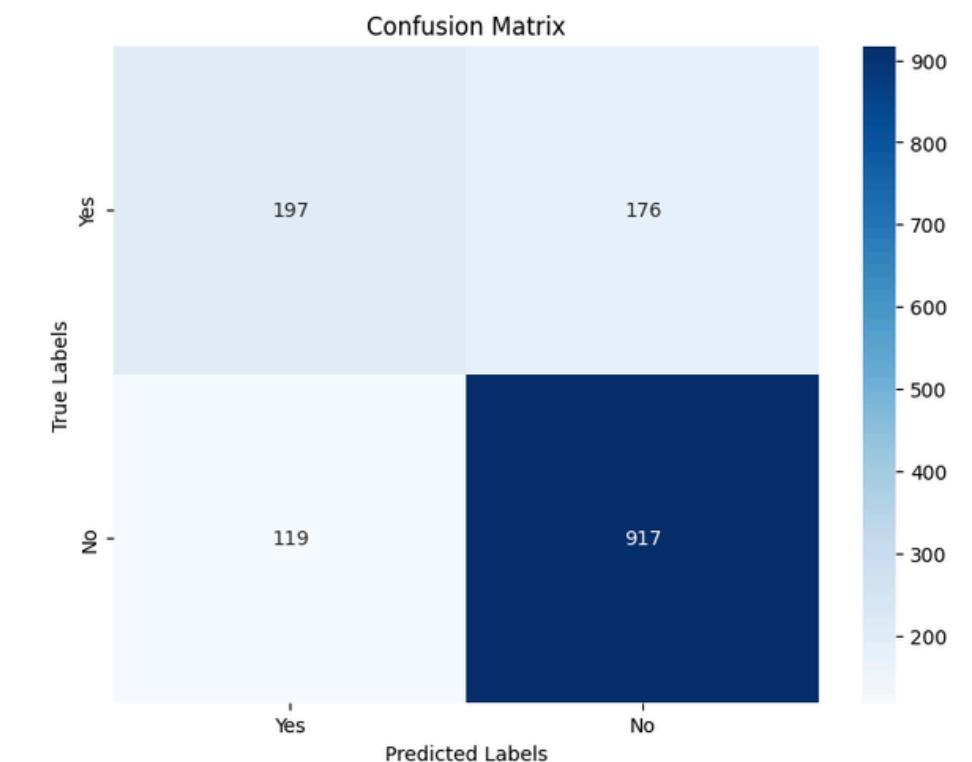
## Metrik Evaluasi



## Decission Tree



## XGBoost

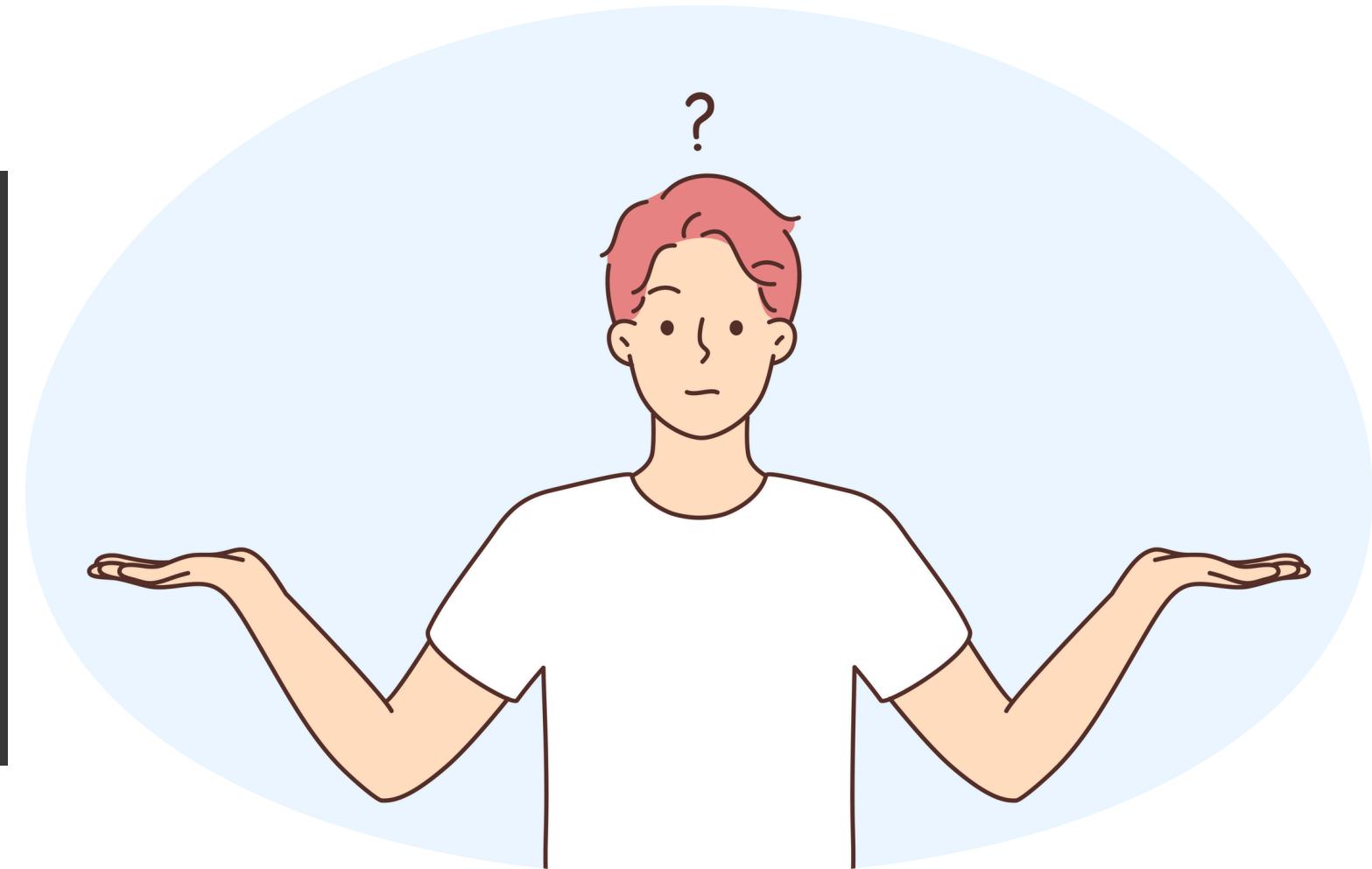


# MODEL

Features Importance

## Decission Tree

	Feature	Importance
6	MonthlyCharges	0.27
4	TotalCharges	0.23
0	Contract_encoded	0.17
1	tenure	0.12
2	OnlineSecurity_encoded	0.05
11	Partner_encoded	0.03
7	PaperlessBilling_encoded	0.03
9	Dependents_encoded	0.03
10	SeniorCitizen	0.02
3	TechSupport_encoded	0.02
5	OnlineBackup_encoded	0.02
8	DeviceProtection_encoded	0.02



## XGBoost

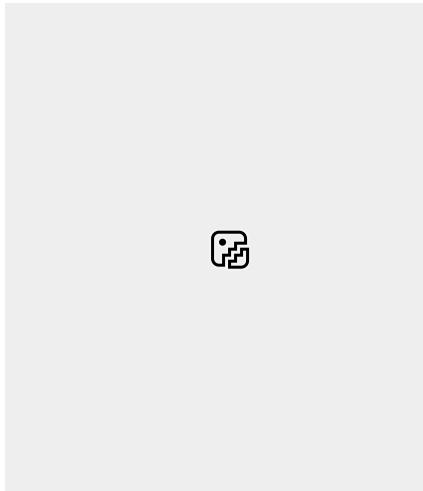
	Feature	Importances
0	Contract_encoded	0.52
2	OnlineSecurity_encoded	0.12
3	TechSupport_encoded	0.07
1	tenure	0.04
6	MonthlyCharges	0.04
7	PaperlessBilling_encoded	0.04
10	SeniorCitizen	0.03
5	OnlineBackup_encoded	0.03
4	TotalCharges	0.03
8	DeviceProtection_encoded	0.03
9	Dependents_encoded	0.03
11	Partner_encoded	0.02

Dari kedua model tersebut, bisa disimpulkan bahwa variable Contract memiliki pengaruh paling besar dan disusul oleh variabel tenure yang memiliki pengaruh sedang.

# TESTING

Diketahui data customer baru sebut saja Asep

```
Contract_encoded = 'Two year'  
tenure = 10  
OnlineSecurity_encoded = 'Yes'  
TechSupport_encoded = 'Yes'  
TotalCharges = 100  
OnlineBackup_encoded = "Yes"  
MonthlyCharges = 80  
PaperlessBilling_encoded = 'Yes'  
DeviceProtection_encoded = "Yes"  
Dependents_encoded = "Yes"  
SeniorCitizen = 1  
Partner_encoded = "Yes"
```



Hasil test pada rill menggunakan model Decission Tree

```
# Lakukan prediksi pada data baru  
new_predictions = model.predict(new_df)  
  
# Jika Anda ingin melihat hasil prediksi  
print("Hasil Prediksi untuk data baru:")  
if new_predictions[0] == 1:  
    print("Pelanggan akan churn")  
else:  
    print("Pelanggan tidak akan churn")
```

Hasil Prediksi untuk data baru:  
Pelanggan akan churn

Hasil test pada rill menggunakan model XGBoost

```
# Lakukan prediksi pada data baru  
new_predictions = model.predict(new_df)  
  
# Jika Anda ingin melihat hasil prediksi  
print("Hasil Prediksi untuk data baru:")  
if new_predictions[0] == 1:  
    print("Pelanggan akan churn")  
else:  
    print("Pelanggan tidak akan churn")
```

Hasil Prediksi untuk data baru:  
Pelanggan akan churn



**Thank you  
very much!**

**[www.colab.research.google.com](https://www.colab.research.google.com)**