

# Data Engineer – Data Fellowship 12 IYKRA

## Assignment 4

Alya Mutiara Firdausyi

---

### Instruction

#### Scenario

You work as a data analytics engineer for a Ritz Jager Bank, analyzing customer behavior and predicting churn. The dataset provides information on existing credit card customers, including demographic data, financial behavior, and interaction history. Your goal is to build a predictive model to identify customers likely to churn (Attrition\_Flag = "Attrited Customer") based on their attributes and behavior. This model will help the bank retain valuable customers and reduce customer churn.

#### Dataset Link

<https://drive.google.com/file/d/1uD1MQWducAAu1-gEJhIRZ0wV-zd6-fC/view?usp=sharing>

#### Answer the following questions

1. Create the data quality report of the dataset you received. What is the issue? How to solve the issue?
  2. Solve the data quality issue on the dataset.
  3. Create the exploratory data analysis of the dataset.
  4. Explain 5 univariate and bivariate analyses that you found interesting.
  5. Your company requests that you create a predictive analytics model based on the above problem. Demonstrate the process and show the model performance result
- 

### Data Quality Report: Ritz Jager Bank Customer Churn Dataset

#### Introduction

This report assesses the quality of the Ritz Jager Bank customer dataset delivered to the data analytics engineer team. The goal is to identify potential data quality issues that would hinder the model's accuracy and effectiveness in predicting customer churn.

#### Data Quality Evaluation

The following dimensions were evaluated to assess the data quality within the dataset.

1. **Completeness:** All attributes appear to have values and there is no missing value from the dataset.

Column	Type	% valid	% invalid	% non empty	% empty
Attrition_Flag	Text	100.00%	0.00%	100.00%	0.00%
Customer_Age	Integer	100.00%	0.00%	100.00%	0.00%
Gender	Text	100.00%	0.00%	100.00%	0.00%
Dependent_count	Integer	100.00%	0.00%	100.00%	0.00%
Education_Level	Text	100.00%	0.00%	100.00%	0.00%
Marital_Status	Text	100.00%	0.00%	100.00%	0.00%
Income_Category	Text	100.00%	0.00%	100.00%	0.00%
Card_Category	Text	100.00%	0.00%	100.00%	0.00%
Months_on_book	Integer	100.00%	0.00%	100.00%	0.00%
Total_Relationship_Count	Integer	100.00%	0.00%	100.00%	0.00%
Months_Inactive_12_mon	Integer	100.00%	0.00%	100.00%	0.00%
Contacts_Count_12_mon	Integer	100.00%	0.00%	100.00%	0.00%
Credit_Limit	Decimal	100.00%	0.00%	100.00%	0.00%
Total_Revolving_Bal	Integer	100.00%	0.00%	100.00%	0.00%
Avg_Open_To_Buy	Decimal	100.00%	0.00%	100.00%	0.00%
Total_Amt_Chng_Q4_Q1	Decimal	100.00%	0.00%	100.00%	0.00%
Total_Trans_Amt	Integer	100.00%	0.00%	100.00%	0.00%
Total_Trans_Ct	Integer	100.00%	0.00%	100.00%	0.00%
Total_Ct_Chng_Q4_Q1	Decimal	100.00%	0.00%	100.00%	0.00%
Avg_Utilization_Ratio	Decimal	100.00%	0.00%	100.00%	0.00%
Naive_Bayes_Classifier_Attrition_Flag_Card_Category_Contacts_Count_12_mon_Dependent_count_Education_Level_Months_Inactive_12_mon_1	Decimal	100.00%	0.00%	100.00%	0.00%
Naive_Bayes_Classifier_Attrition_Flag_Card_Category_Contacts_Count_12_mon_Dependent_count_Education_Level_Months_Inactive_12_mon_2	Decimal	100.00%	0.00%	100.00%	0.00%

2. **Consistency:** It seems we need to justify some column types to better fit the table schema.

CLIENTNUM	Attrition_Flag	Customer_Age	Gender	Dependent_count	Education_Level	Marital_Status	Income_Category	Card_Category	Months_on_book	Total_Relationship_Count
768805383	Existing Customer	45	M	3	High School	Married	\$60K - \$80K	Blue	39	
818770008	Existing Customer	49	F	5	Graduate	Single	Less than \$40K	Blue	44	
713982108	Existing Customer	51	M	3	Graduate	Married	\$80K - \$120K	Blue	36	

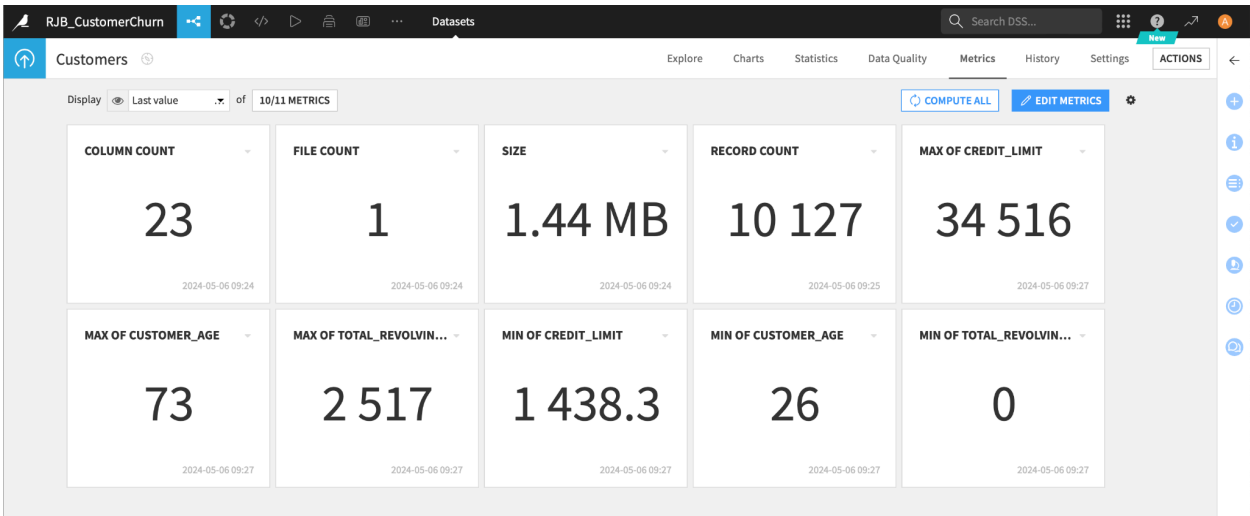
Months_Inactive_12_mon	Contacts_Count_12_mon	Credit_Limit	Total_Revolving_Bal	Avg_Open_To_Buy	Total_Amt_Chng_Q4_Q1	Total_Trans_Amt	Total_Trans_Ct	Total_Ct_Chng_Q4_Q1	Avg_Util
1	3	12691	777	11914	1.335	1144	42	1.625	
1	2	8256	864	7392	1.541	1291	33	3.714	
1	0	3418	0	3418	2.594	1887	20	2.333	
4	1	3313	2517	796	1.405	1171	20	2.333	

The screenshot shows a data table with 10,127 rows and 23 columns. The columns include: **\_Buy** (string), **Total\_Amt\_Chng\_Q4\_Q1** (Decimal), **Total\_Trans\_Amt** (Integer), **Total\_Trans\_Ct** (Integer), **Total\_Ct\_Chng\_Q4\_Q1** (Decimal), **Avg\_Utilization\_Ratio** (Decimal), **Naive\_Bayes\_Classifier\_Attrition\_Flag\_Card\_Ca...** (string), and **Naive\_Bayes\_Classifier\_Attrition\_Flag\_Card\_Ca...** (string). The first few rows of data are visible.

_Buy	Total_Amt_Chng_Q4_Q1	Total_Trans_Amt	Total_Trans_Ct	Total_Ct_Chng_Q4_Q1	Avg_Utilization_Ratio	Naive_Bayes_Classifier_Attrition_Flag_Card_Ca...	Naive_Bayes_Classifier_Attrition_Flag_Card_Ca...
11914	1.335	1144	42	1.625	0.061	0.000093448	0.99991
7392	1.541	1291	33	3.714	0.105	0.000056861	0.99994
3418	2.594	1887	20	2.333	0	0.000021081	0.99998
796	1.405	1171	20	2.333	0.76	0.00013366	0.99987

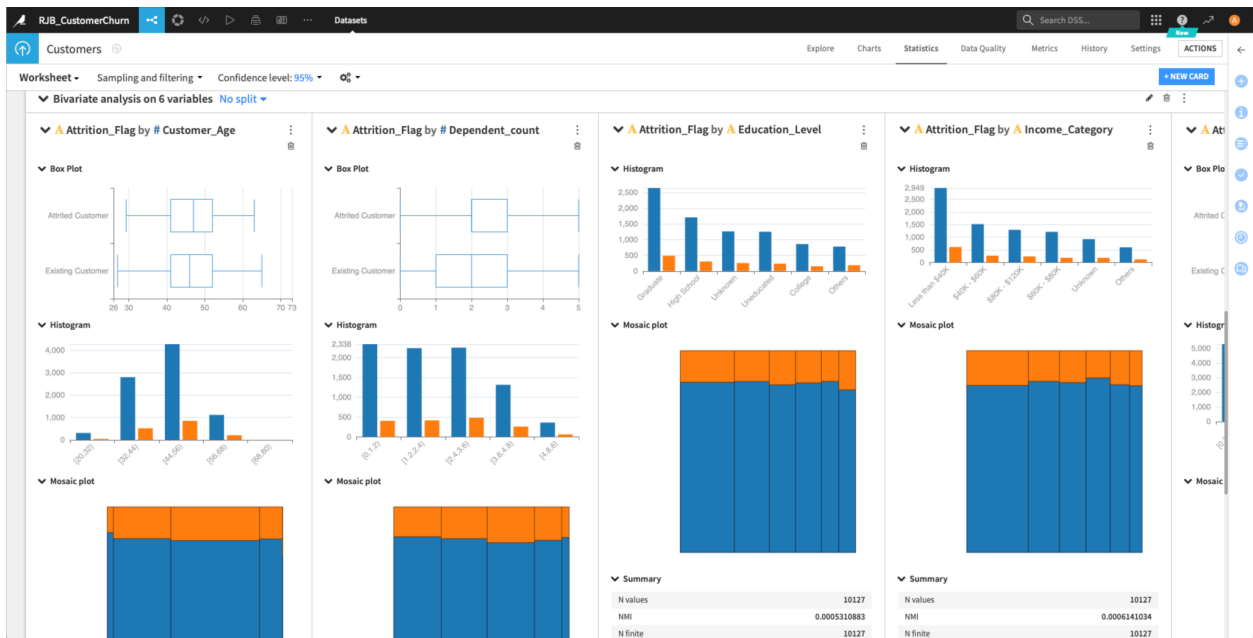
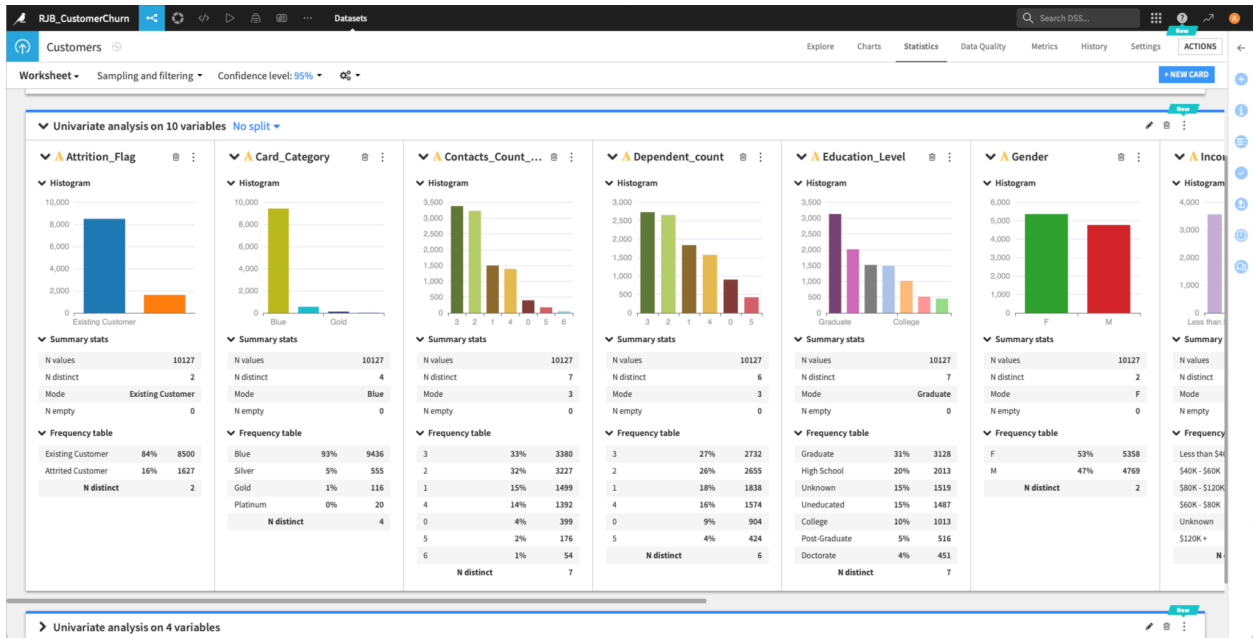
3. Validity:

Metadata



Exploratory Data Analysis

Attritied Customer Analysis



# Predictive Analytics Model

