

Post Graduation in Data Science & Business  
Analytics - 11<sup>th</sup> Edition  
Data Science Project

**Decision Support Model for Determining  
Appropriate Loan Interest Rates for Customers**



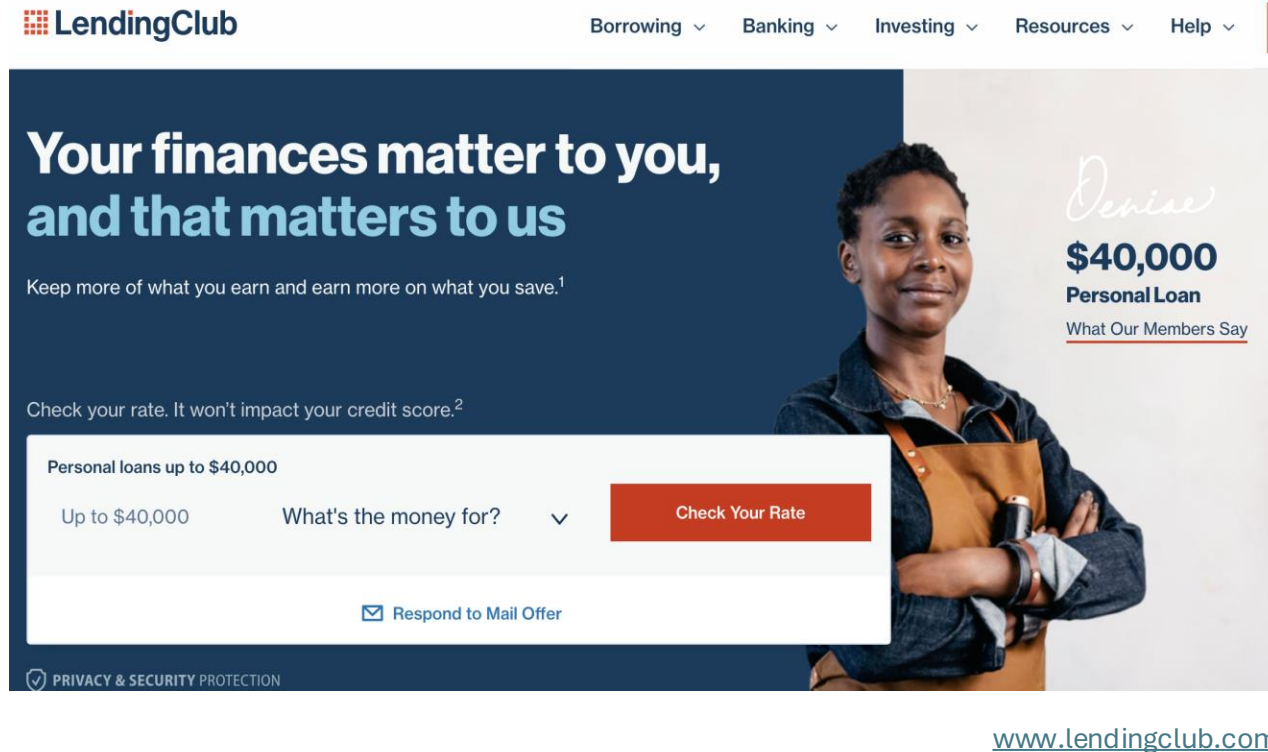
Carlos Ramagem | i30386  
Luiza Campos | i30387  
Thaís Almeida | i30394

## INDEX

- Introduction to the use case and project objectives
- Presentation of the main technical issues implemented
  - 2.1 Dimensional Model
  - 2.2 ETL
  - 2.3 Semantic Model
- Presentation of conclusions obtained through report navigation (Power BI)
- Presentation of conclusions obtained and explanation of the implemented Machine Learning model (ML Studio)
- Conclusions

# 1. Introduction to the use case and project objectives

## Use case



- Online Crowdfunding Loan Platform
- Provides accessible credit to individuals/entities by establishing a "bridge" between investors (who provide the loan) and borrowers.
- Transactions are carried out through direct interactions between the different agents on the platform (borrowers and investors/subscribers of bonds).

### Objective:

To provide recommendations for more competitive interest rates in the market, aiming to offer better profit margins for companies while reducing the risk of default.

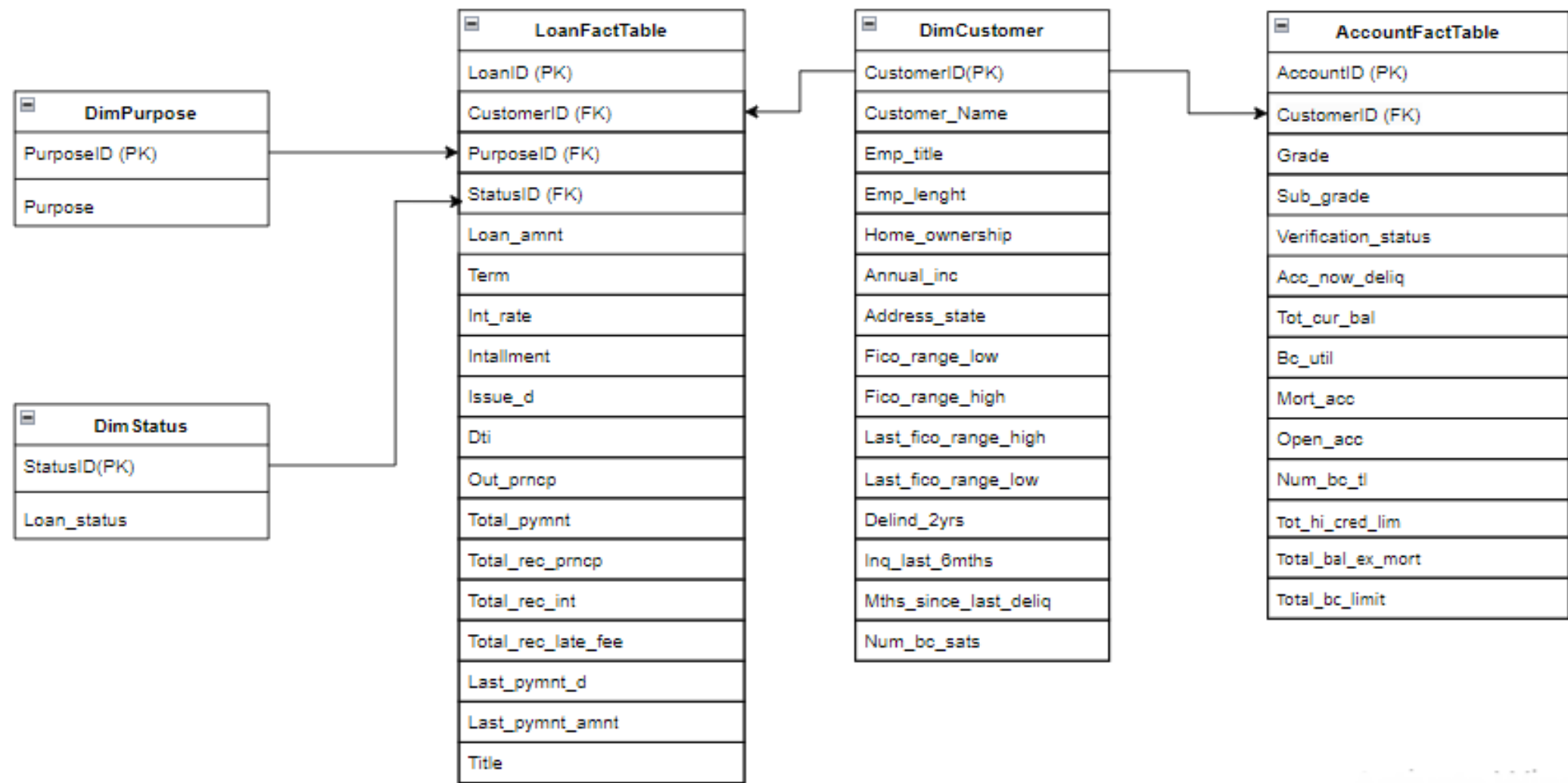
# Script to create the star schema

SQL script defining all variables contained in the star schema model

```
1 /*Criando a tabela DSP11_GR10.DimCustomer */
2 Create table DSP11_GR10.DimCustomer(
3 [CustomerID] [BIGINT] PRIMARY KEY,
4 [Customer_name] [varchar](100) NULL,
5 [Emp_title] [varchar](100) NULL,
6 [Emp_length] [int] NULL,
7 [Home_ownership] [varchar](100) NULL,
8 [Annual_inc] [float] NULL,
9 [Addr_state] [varchar](100) NULL,
10 [Fico_range_low] [int] NULL,
11 [Fico_range_high] [int] NULL,
12 [Last_fico_range_high] [int] NULL,
13 [Last_fico_range_low] [int] NULL,
14 [Delinq_2yrs] [int] NULL,
15 [Inq_last_6mths] [int] NULL,
16 [Mths_since_last_delinq] [int] NULL,
17 [Num_bc_sats] [int] NULL,
18 );
19
20 DROP TABLE DSP11_GR10.DimCustomer
21
22 select * from DSP11_GR10.DimCustomer
23
24 GO
25
26 /*Criando a tabela DSP11_GR10.DimPurpose */
27 Create table DSP11_GR10.DimPurpose(
28 [PurposeID] [BIGINT] PRIMARY KEY,
29 [Purpose] [varchar](100) NULL,
30 );
31
32 GO
33
34 DROP TABLE DSP11_GR10.DimPurpose
35
36 select * from DSP11_GR10.DimPurpose
37
38 /*Criando a tabela DSP11_GR10.DimStatus */
39 Create table DSP11_GR10.DimStatus(
40 [StatusID] [BIGINT] PRIMARY KEY,
41 [Loan_status] [varchar](100) NULL,
42 );
43
44 select * from DSP11_GR10.DimStatus
```

```
1 /*Criando a tabela DSP11_GR10.LoanFactTable */
2 CREATE TABLE DSP11_GR10.LoanFactTable(
3 [LoanID] [BIGINT] PRIMARY KEY,
4 [CustomerID] [int] NULL,
5 [PurposeID] [int] NULL,
6 [StatusID] [int] NULL,
7 [Loan_amnt] [int] NULL,
8 [Term] [int] NULL,
9 [Int_rate] [float] NULL,
10 [Installment] [float] NULL,
11 [Issue_d] [date],
12 [Dti] [float] NULL,
13 [Out_prncp] [float] NULL,
14 [Total_pymnt] [float] NULL,
15 [Total_rec_prncp] [float] NULL,
16 [Total_rec_int] [float] NULL,
17 [Total_rec_late_fee] [float] NULL,
18 [Last_pymnt_d] [date],
19 [Last_pymnt_amnt] [float] NULL,
20 [Title] [varchar](100) NULL);
21
22 select * from DSP11_GR10.LoanFactTable
23
24 DROP TABLE DSP11_GR10.LoanFactTable
25
26 /*Criando a tabela DSP11_GR10.AccountFactTable */
27 CREATE TABLE DSP11_GR10.AccountFactTable(
28 [AccountID] [BIGINT] PRIMARY KEY,
29 [CustomerID] [int] NULL,
30 [Grade] [varchar](50) NULL,
31 [Sub_grade] [varchar](50) NULL,
32 [Verification_status] [varchar](100) NULL,
33 [Acc_now_delinq] [int] NULL,
34 [Tot_cur_bal] [int] NULL,
35 [Bc_util] [float] NULL,
36 [Mort_acc] [int] NULL,
37 [Open_acc] [int] NULL,
38 [Num_bc_tl] [int] NULL,
39 [Tot_hi_cred_lim] [int] NULL,
40 [Total_bal_ex_mort] [int] NULL,
41 [Total_bc_limit] [int] NULL);
42
43 DROP TABLE DSP11_GR10.AccountFactTable
```

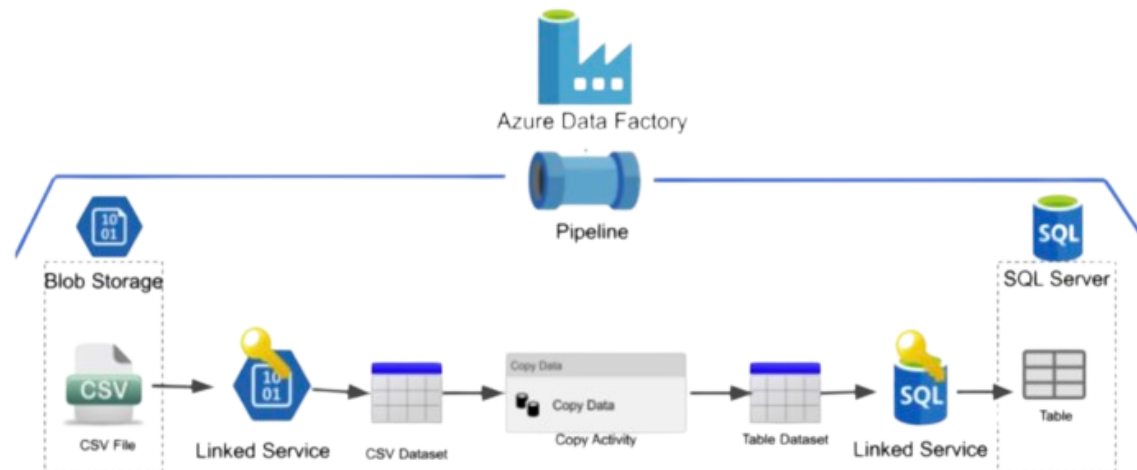
Star schema



**Objective:** Automate the process of data ingestion, transformation, and importation using Azure Data Factory and Azure Data Studio.

**Strategy Used:**

- Ingestion: Extraction of data from a blob container in Azure.
- Transformation: Use of the "FlowData" artifact to process and prepare the data.
- Importation: Loading the transformed data into Azure SQL for analysis and centralized storage

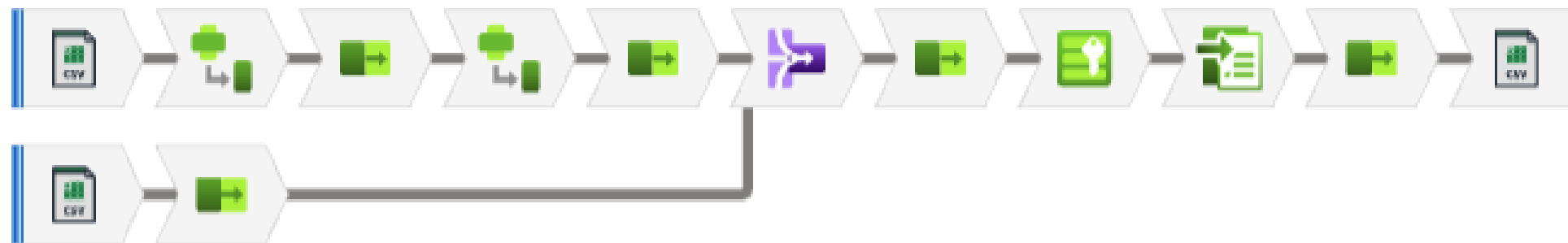


## Transformation

- DimCustomer

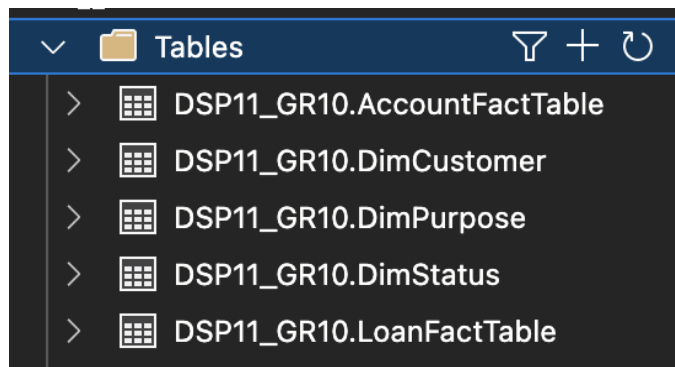


- FactAccount



### Importation

- Transformed data from each DataFlow in Azure Data Factory, related to each entity, was loaded into the Azure SQL account created for this project.
- Pipeline was used to load the data into Azure SQL. For each table, SQL code was executed to set up the corresponding structure according to the data schema.

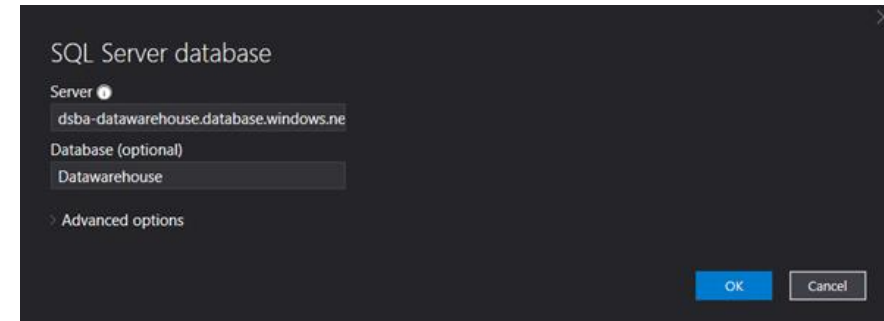


```
/*Criando a tabela DSP11_GR10.DimCustomer */
Create table DSP11_GR10.DimCustomer(
[CustomerId] int identity(1000,1),
[emp_title] [varchar](100) NULL,
[emp_length] [varchar](100) NULL,
[home_ownership] [varchar](100) NULL,
[annual_inc] [float] NULL,
[addr_state] [varchar](100) NULL,
[fico_range_low] [int] NULL,
[fico_range_high] [int] NULL,
[last_fico_range_high][int] NULL,
[last_fico_range_low][int] NULL,
[delinq_2yrs] [int] NULL,
[inq_last_6mths] [int] NULL,
[mths_since_last_delinq] [int] NULL,
[num_bc_sats][int] NULL,
);
GO
```



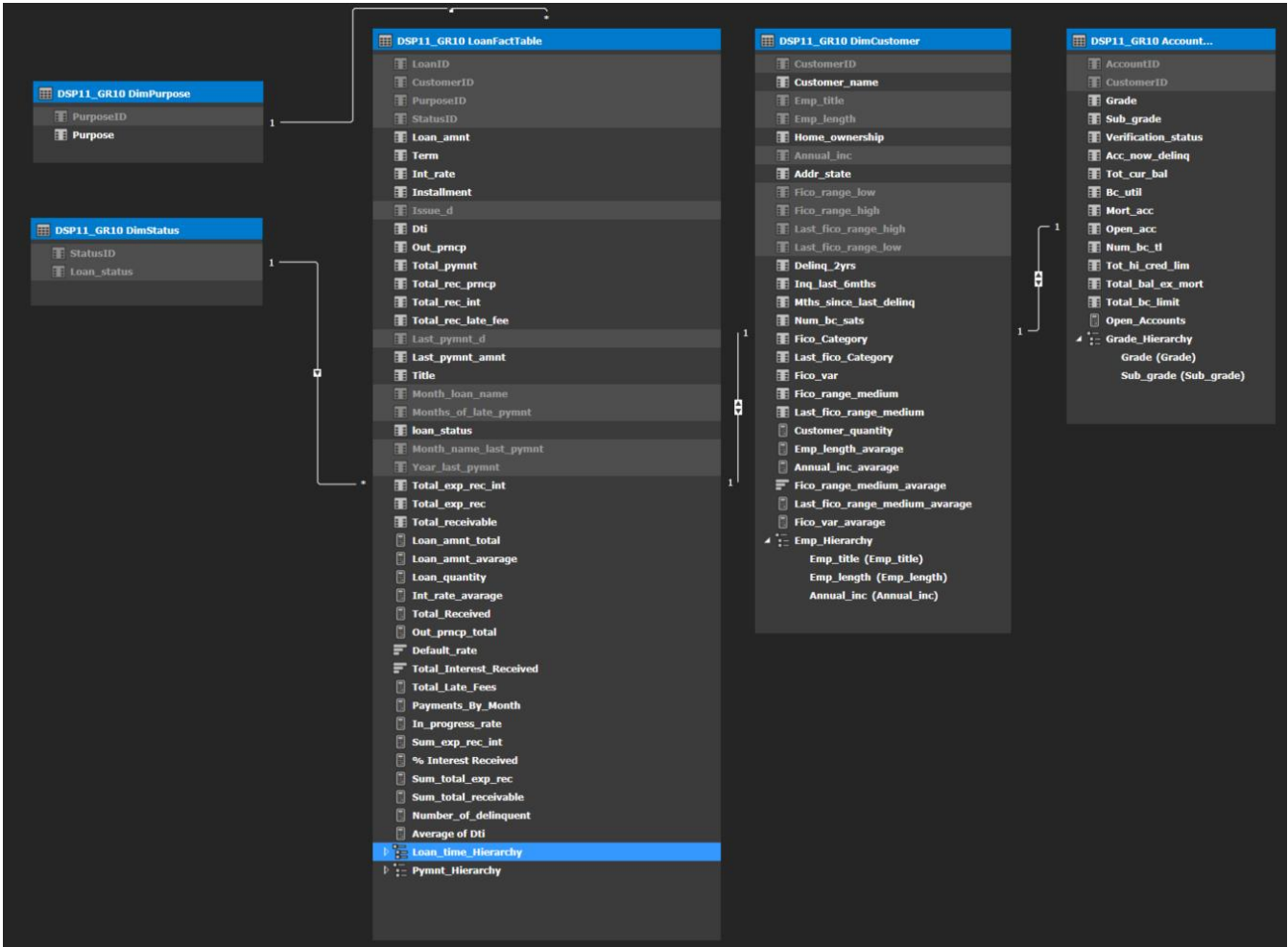
## 2.3 Modelo Semântico

- The transformed and stored data tables in Azure were connected to Visual Studio using the Microsoft Analysis Services template.
- This integration allowed for the import and organization of data tables within the Visual Studio development environment, preparing them for the construction of the semantic model.



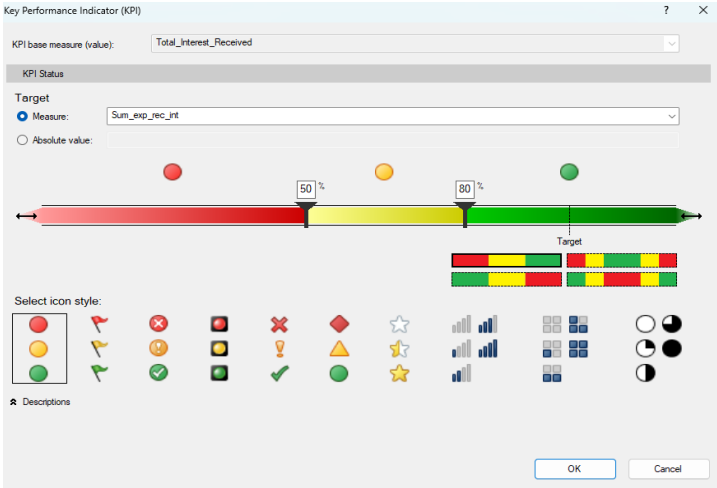
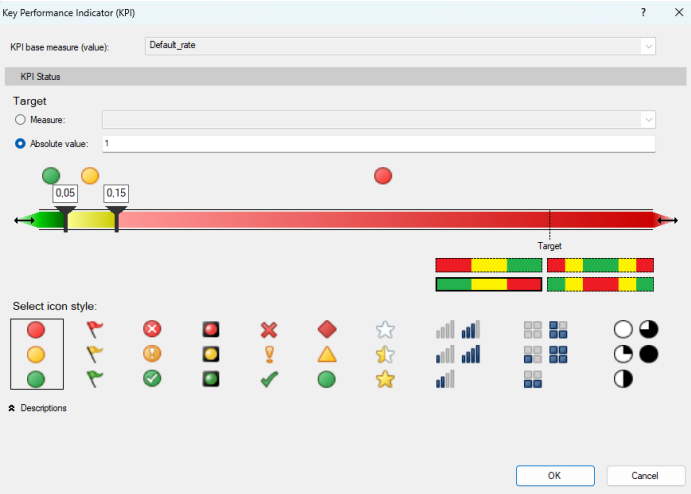
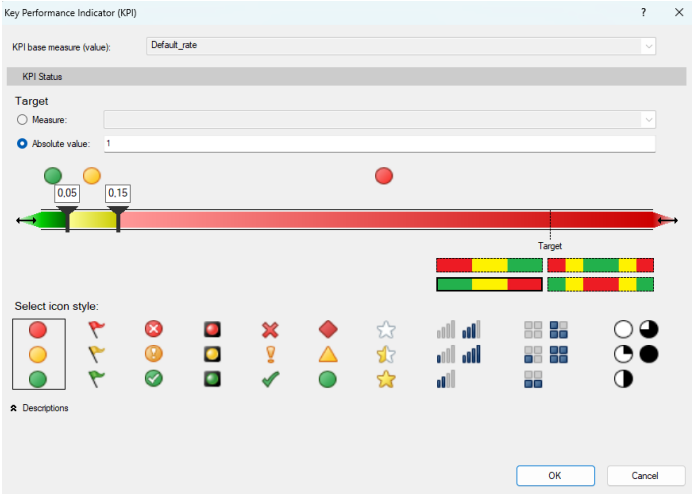
## 2.3 Semantic Model

- It was necessary to define the relationships between the tables to facilitate data analysis, including the addition of calculated columns and the definition of measures for business evaluation.

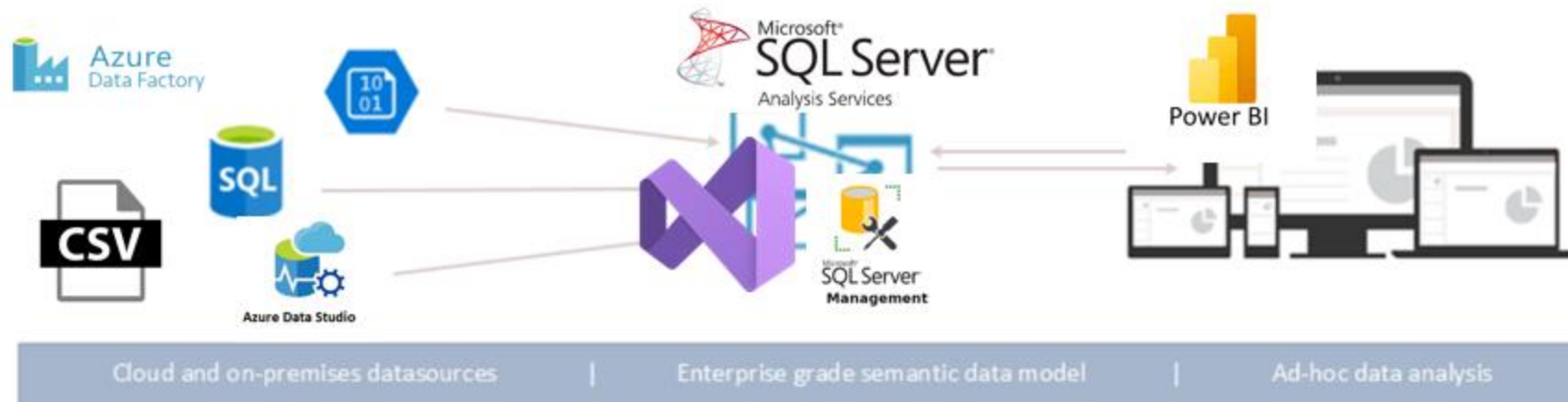


## 2.3 Semantic Model

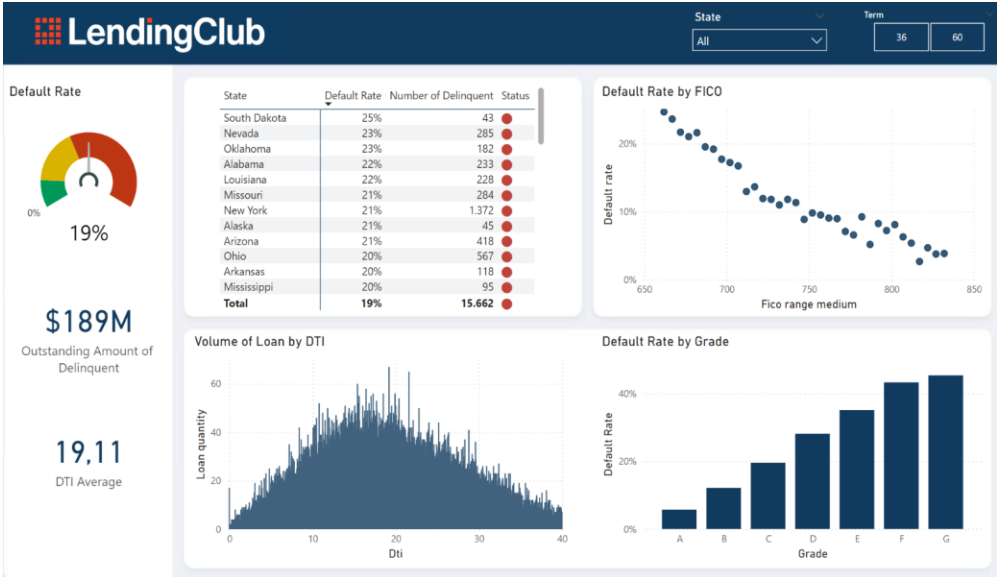
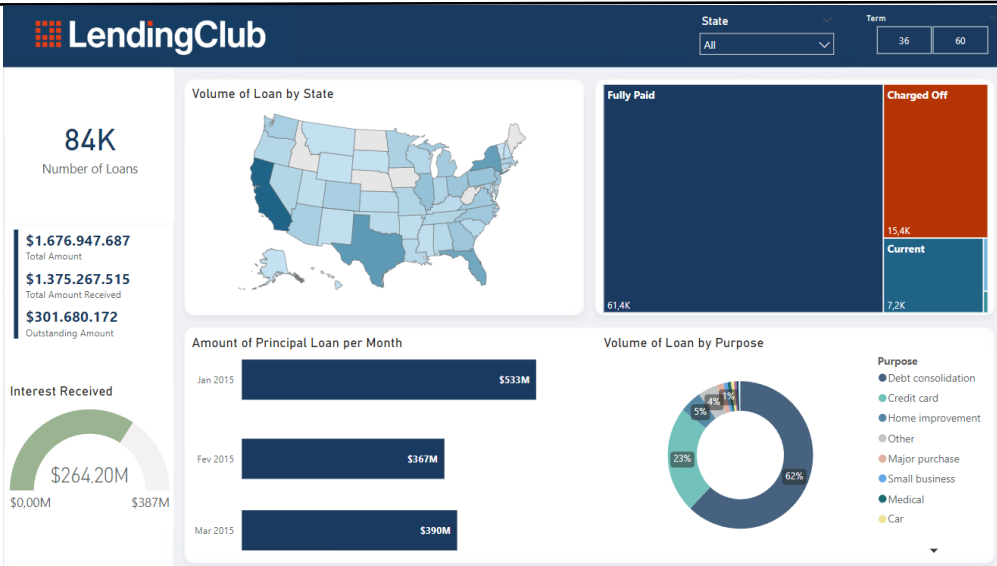
### Creation of KPIs for Performance Monitoring



### PoweBI's deployment

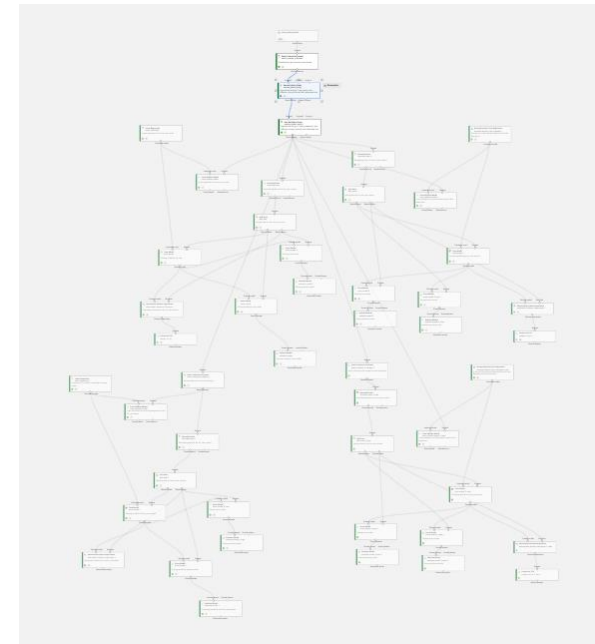


3. Power BI Dashboards



## 4. Implemented Machine Learning Model (ML Studio)

- **Objective:** Predict the appropriate interest rate for a loan, specifically tailored to each client's profile and history.
- **Strategy Used:** From Simple to Complex  
An iterative approach to developing the Machine Learning model.
- **Chosen Feature:** Loan Interest Rate (Int\_rate)
- The methods used were **Linear Regression** and **Boosted Decision Tree Regression**.
- For Linear Regression, **L2 regularization** and feature selection strategies were applied as model tuning processes to find the best fit.
- Still not satisfied with the results, **Boosted Decision Tree Regression** was tested with hyperparameter tuning and feature selection to find the optimal model.



## 4. Implemented Machine Learning Model (ML Studio)

MAE (Mean Absolute Error):  
Lowest MAE of  $0.90 \pm 0.01$ , indicating higher model accuracy.

RMSE (Root Mean Square Error):  
Lowest RMSE of  $1.29 \pm 0.03$ , suggesting a better model fit.

RSE (Relative Squared Error):  
Lowest RSE of  $0.09 \pm 0.00$ , indicating that 9% of the variance in interest rates is explained, highlighting a superior model fit compared to Linear Regression.

RAE (Relative Absolute Error):  
Lowest RAE of  $0.26 \pm 0.00$ , indicating minimal deviation from actual interest rates compared to other models.

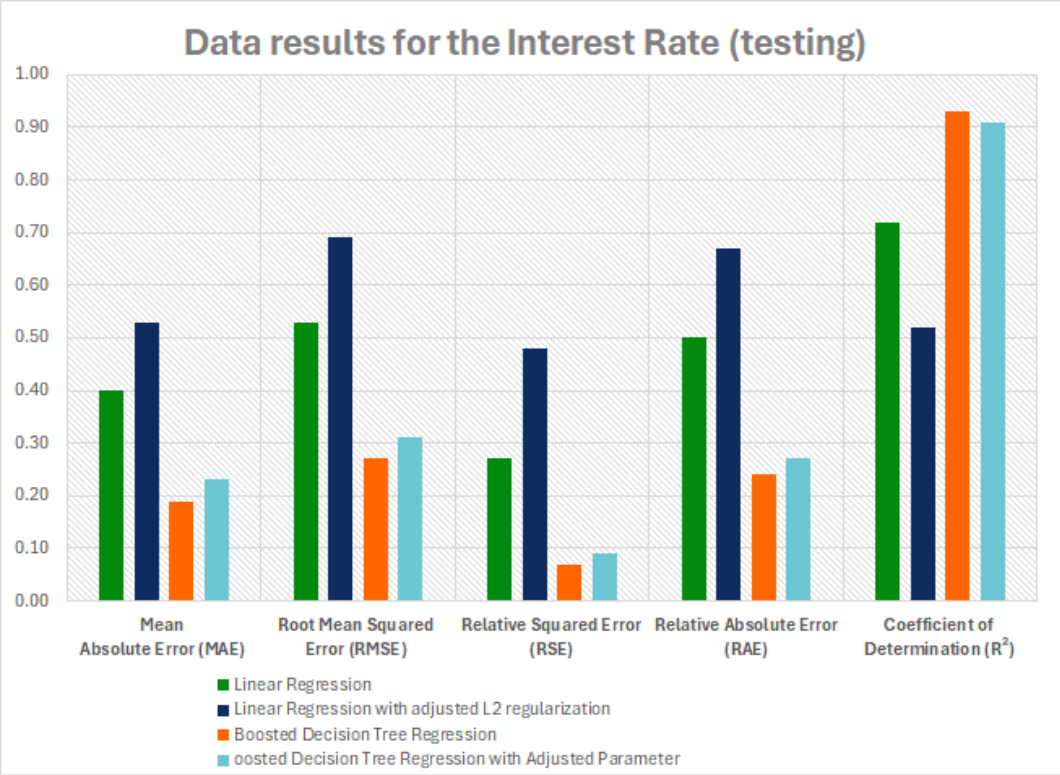
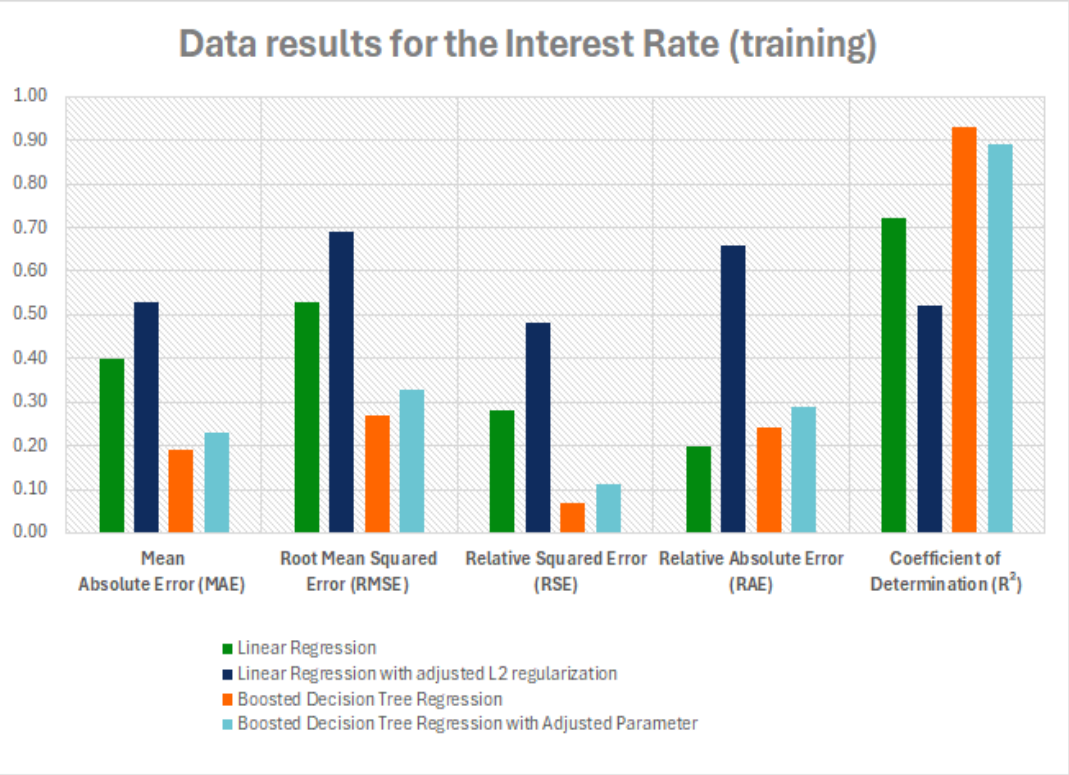
$R^2$  (Coefficient of Determination):  
Highest  $R^2$  of  $0.91 \pm 0.00$ , indicating 91% accuracy in predicting interest rates, meaning the model is highly capable of explaining and predicting variations in interest rates based on the analyzed data.

Table 1 – Analysis of Cross validation Results for the Interested Rate (mean across 10 folds)

Model Version (Cross validation for the <u>Int_rate</u> column)	Mean Absolute Error (MAE) with $\sigma$	Root Mean Squared Error (RMSE) with $\sigma$	Relative Squared Error (RSE) with $\sigma$	Relative Absolute Error (RAE) with $\sigma$	Coefficient of Determinatio n ( $R^2$ ) with $\sigma$
Linear Regression	$1.74 \pm 0.02$	$2.26 \pm 0.03$	$0.28 \pm 0.01$	$0.51 \pm 0.00$	$0.72 \pm 0.01$
Linear Regression with adjusted L2 regularization	$2.36 \pm 0.02$	$3.06 \pm 0.03$	$0.51 \pm 0.01$	$0.69 \pm 0.01$	$0.49 \pm 0.01$
Boosted Decision Tree Regression	$0.90 \pm 0.01$	$1.29 \pm 0.03$	$0.09 \pm 0.00$	$0.26 \pm 0.00$	$0.91 \pm 0.00$
Boosted Decision Tree Regression with Adjusted Parameters	$0.99 \pm 0.01$	$1.40 \pm 0.02$	$0.11 \pm 0.00$	$0.29 \pm 0.00$	$0.89 \pm 0.00$



## 4. Modelo de Machine Learning implementado (ML Studio)



The similarity between the training model result metrics suggests that the model is well-fitted, generalizes well, and is robust.



## 5. Conclusões

---

- The company faces challenges with customer profiles and a high default rate of 19%, highlighting the need to improve risk management.
- The analysis underscores the importance of rigorous evaluation criteria, given the negative correlation between FICO scores and defaults.
- Implementing the ML techniques proposed in this report for credit approval analysis can reduce loan risks and optimize credit management, critically managing the outstanding amount of \$189 million.



THANK YOU