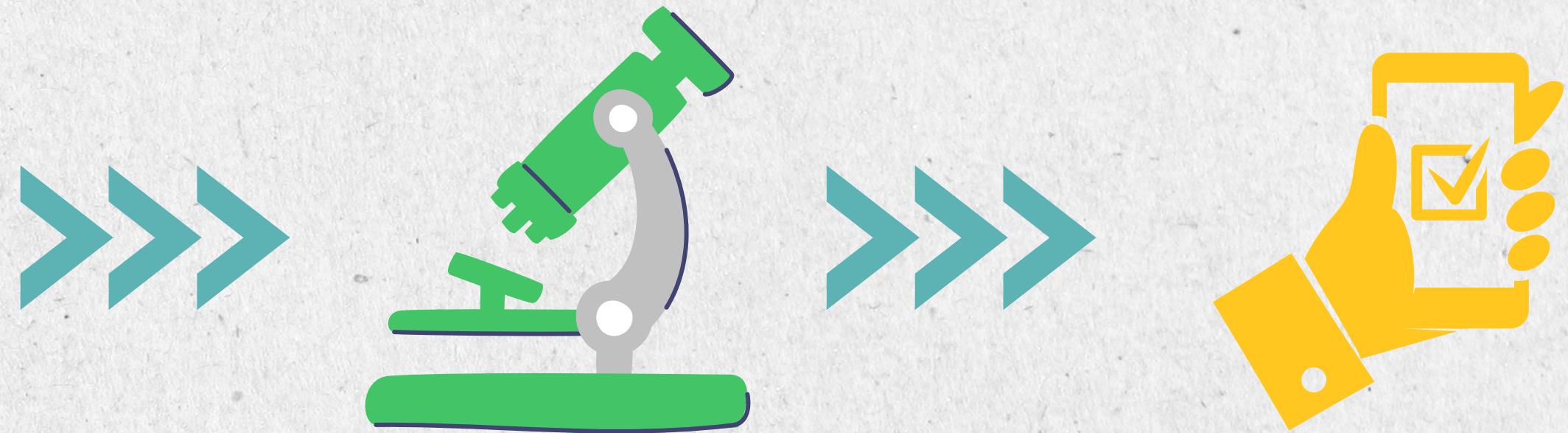


Diabetes Predictive Model



Feras Rafeh – Ianel González – Bruna Santos



Agenda

1

Problem 



2

Data Introduction



3

Data Journey

4

Getting into the insights

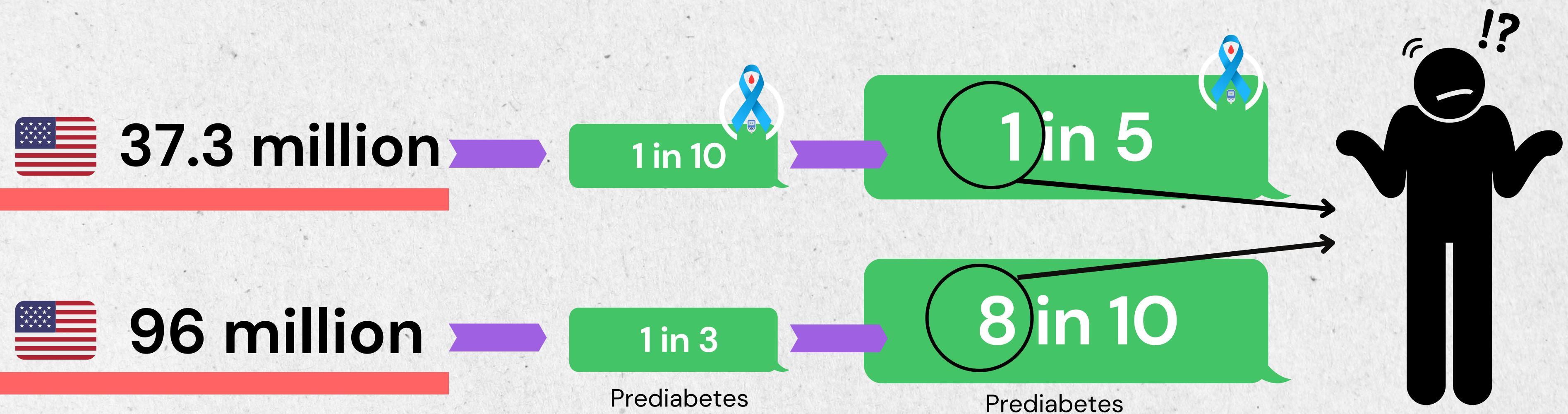


5

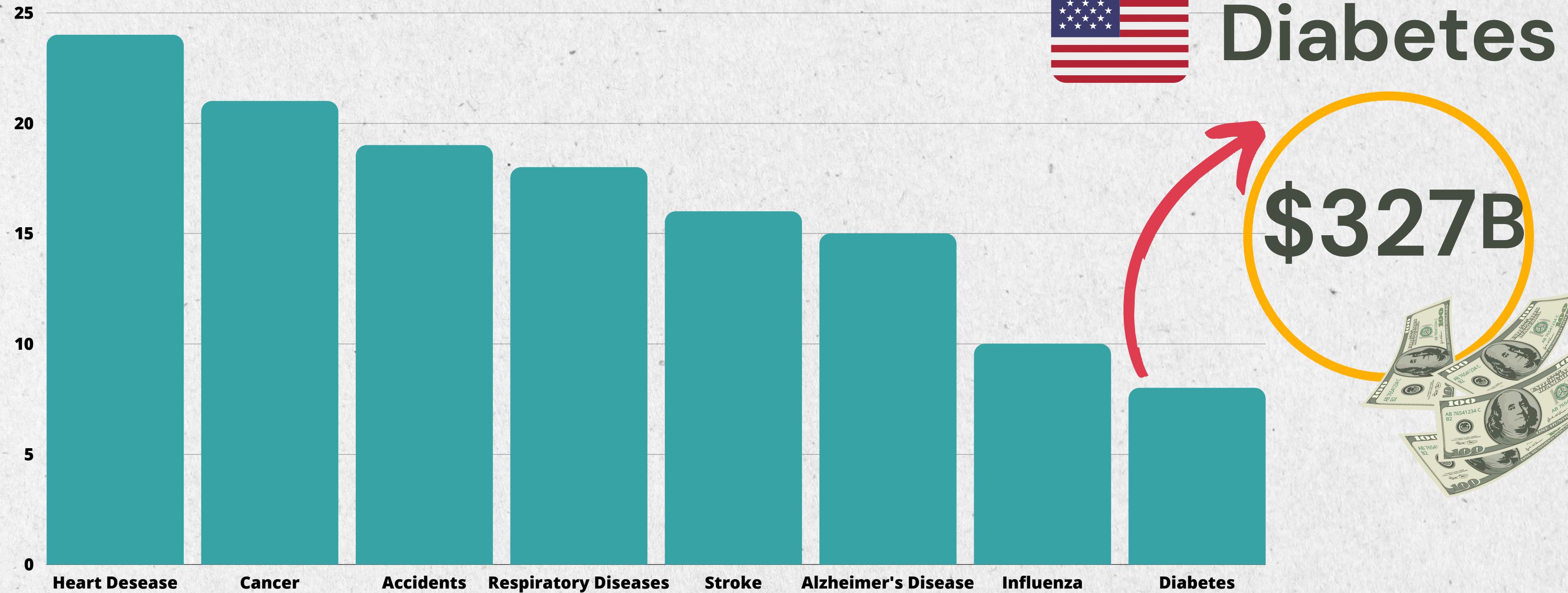
Conclusions



The problem



The problem



Data Introduction



1 Data



The Diabetes prediction dataset is a collection of medical and demographic data from patients, along with their diabetes status (positive or negative).

Data Shape: (100.000, 9) [kaggle](#)



2 Features



Medical Data

- Hypertension
- Heart disease
- Body mass index (BMI)
- Glycosylated hemoglobin
- Blood glucose level



Demographic Data

- Age
- Gender
- Smoking history

3 Objectives

🎯 Developing a predictive model to identify patients at risk of developing Diabetes based on their medical history and demographic information.

🎯 Analyze the impact and importance of the features (independent variables) used to identify the main factors that increase the risk to develop Diabetes.



Data Journey



Cleaning

- Changing Data type in "Age" column
- Dropped ages from 0 to 3
- Dropped "No info" values in "smoking history" (31964 rows)
- Merge some data from "smoking history"

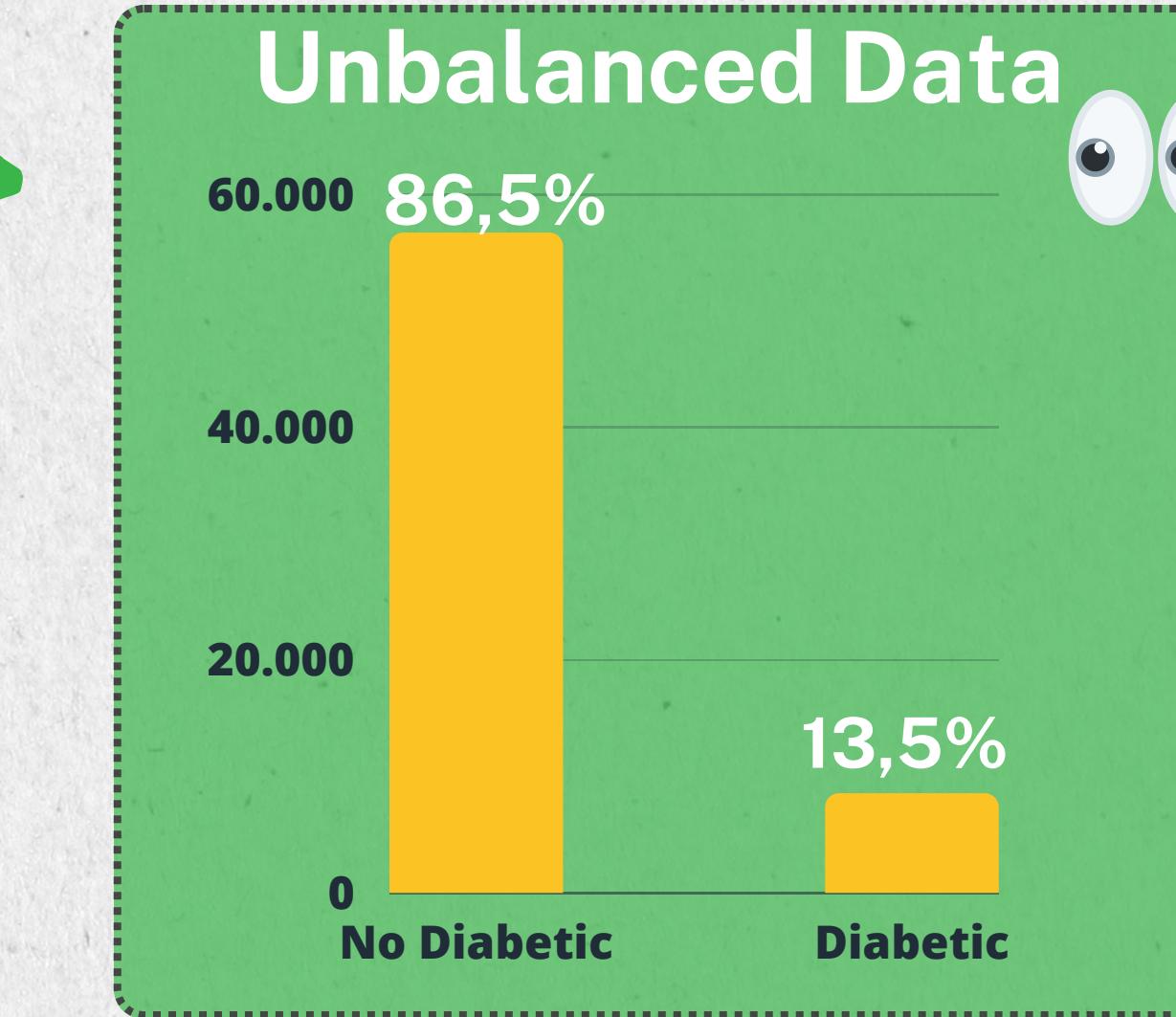


Cleaning

- Changing Data type in "Age" column
- Dropped ages from 0 to 3
- Dropped "No info" values in "smoking history" (31964 rows)
- Merge some data from "smoking history"

Explore

- The data was unbalanced: 56,665 negative results (no diabetes, 86.5% of our data) and 8,500 positive results, 13.5% of our data



Data Journey



Cleaning

- Changing Data type in "Age" column
- Dropped "No info" values in "smoking history"
- Merge some data from "smoking history"
- Dropped ages from 0 to 3



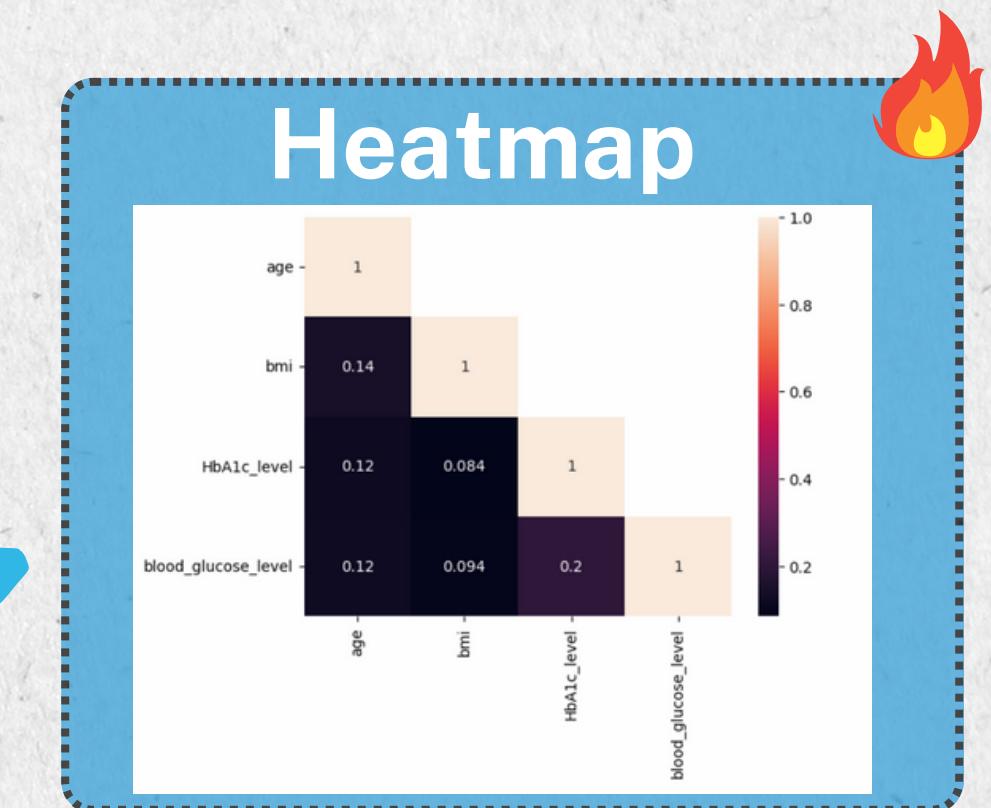
Explore

- The data was unbalanced: 56,665 negative results (no diabetes) and 8,500 positive results)

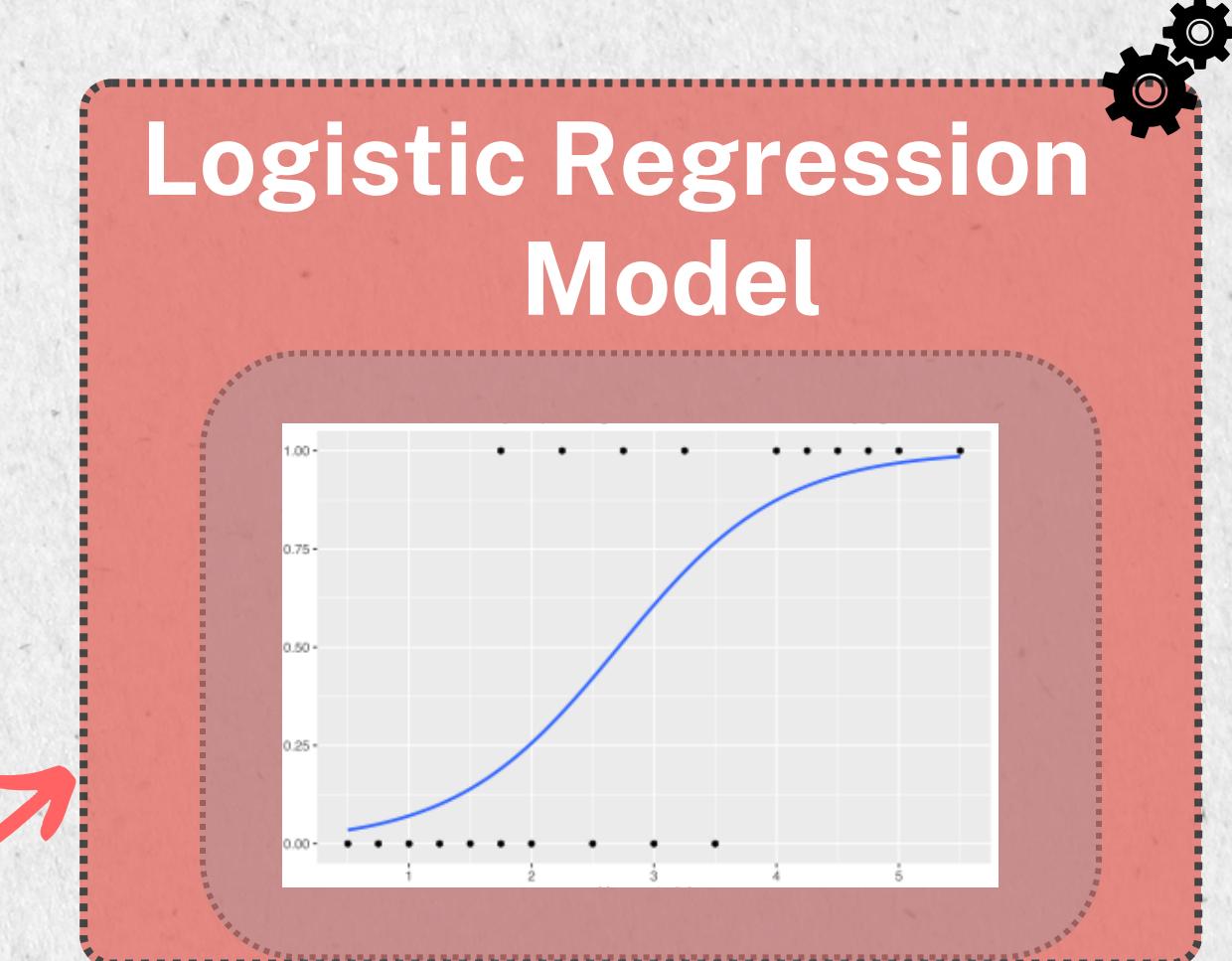
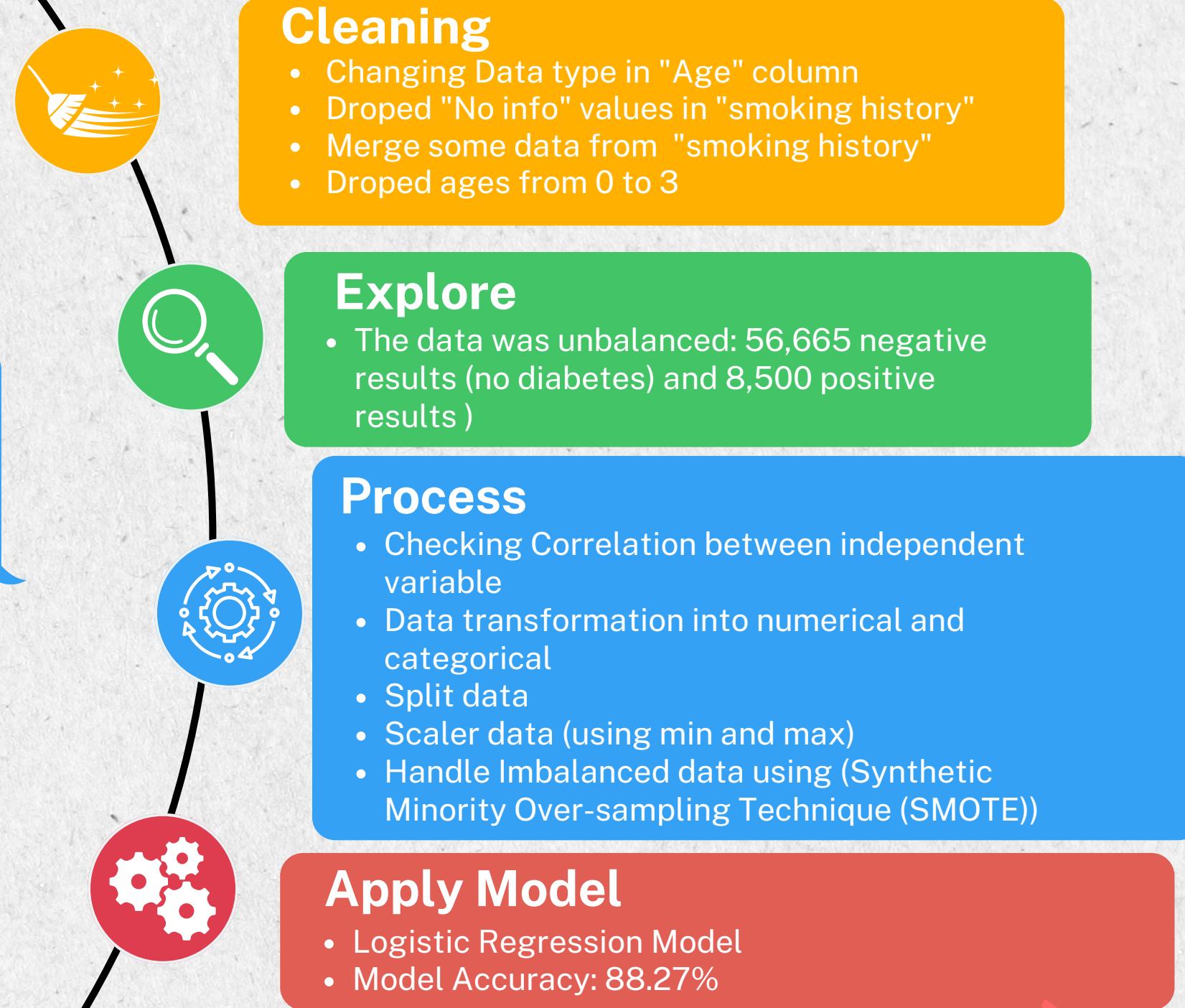


Process

- Checking Correlation between independent variable
- Data transformation into numerical and categorical
- Split data into train and test data
- Scaler data (using min and max)
- Handle Imbalanced data using (Synthetic Minority Over-sampling Technique (SMOTE))



Data Journey





Data Journey



Data Journey

Confusion Matrix

Confusion Matrix

Y-Test	Y-Pred	0	1
0		15.109	1.926
1		316	1.763

- confusion
- False Negative
 - False Positive
 - True Negative
 - True Positive

Count of Y-Test broken down by Y-Pred vs. Y-Test
Colour shows details about confusion.

- categorical
- Split data
 - Scaler data (using min and max)
 - Handle Imbalanced data
- Minority Over-sampling Techniques

Apply Model

- Logistic Regression Model
- Model Accuracy: 88.27%

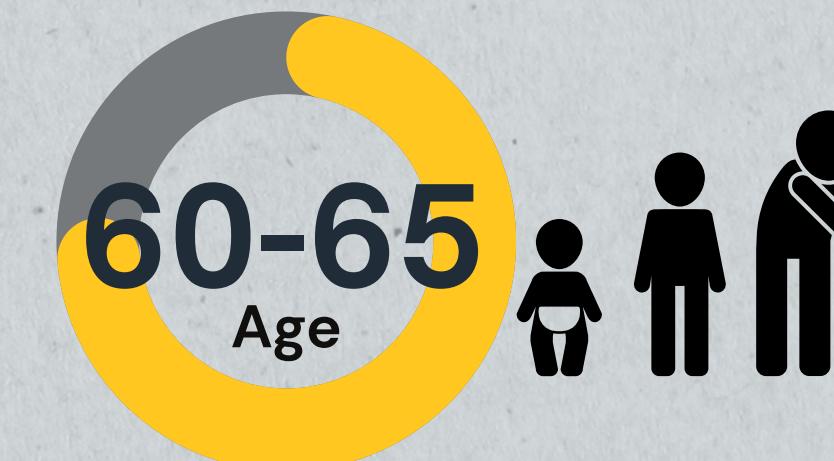
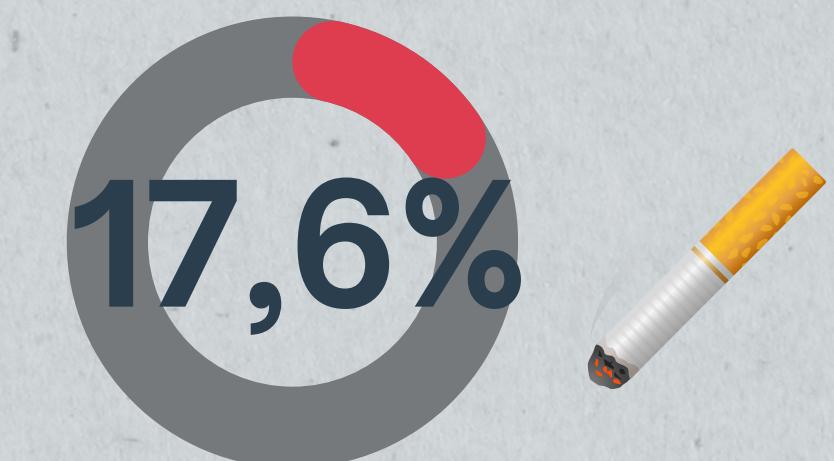
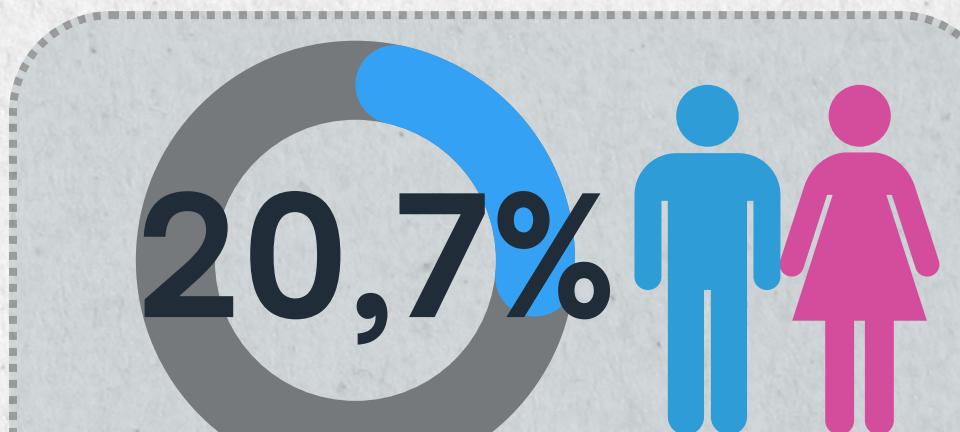
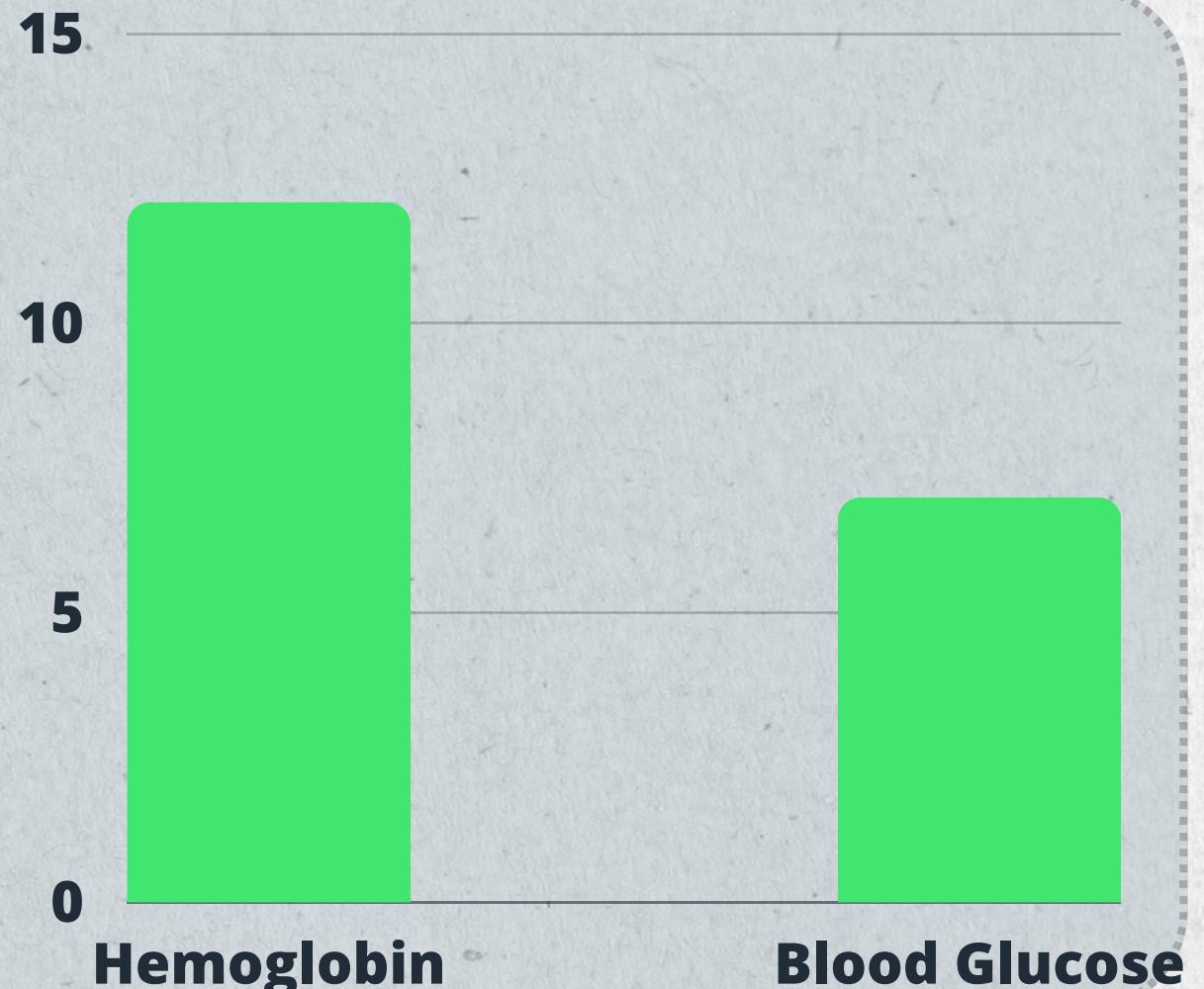
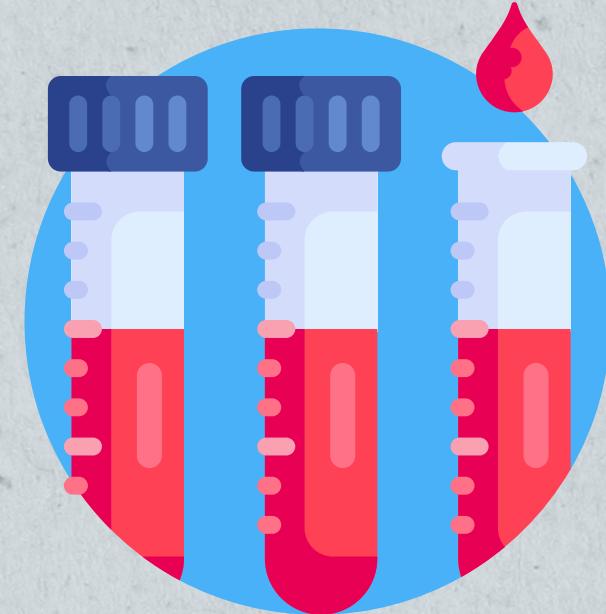
Classification Report

No/Yes	Precision	Recall	F1-score	Support
0	0,98	0,89	0,93	17035
1	0,48	0,85	0,61	2079

Validate

- Classification Report
- Confusion Matrix

Getting into the insights



- Overweight 25.0 - 29.9
- Obesity 1st Class 30.0 - 34.9
- Obesity 2nd Class 35.0 - 39.9

Coefficient: 8.39

Let's see our model in action



Conclusions

✓ We developed a predictive model that identifies patients at risk of developing diabetes based on their medical history and demographic information.



✓ We analyzed the impact and importance that our independent variables have in detecting patients at risk of developing diabetes.



- ✓ Glycosylated hemoglobin
- ✓ BMI
- ✓ Blood glucose

Model Accuracy:

88.27%

Performance Validation:

316 FN

- Age

60–65

- Gender



- Smoking Story



Thank you for your attention!

Find us on
Slack

Ferass Rafeh – Irael González – Bruna Santos

