

DATA SCIENCE

FRAMEWORK REPORT

MAURICIO CARVAJAL

AGENDA

The background of the slide is a dark, textured surface. On the right side, there is a close-up, slightly out-of-focus image of a green ballpoint pen. On the left side, there is a faint, light-colored notepad with a checklist. The checklist has four items, each in a square box. The first two boxes are checked with a green checkmark, and the last two are empty.

- PROBLEM STATEMENT
- PROPOSAL DATA SCIENCE PROCESS FRAMEWORK
- DATA MANAGEMENT
- RISK MANAGEMENT
- INSIGHTS

PROBLEM STATEMENT

- OVER THE PAST YEAR OR SO CREDIT ONE HAS SEEN AN INCREASE IN THE NUMBER OF CUSTOMERS WHO HAVE DEFAULTED ON LOANS THEY HAVE SECURED FROM VARIOUS PARTNERS, AND CREDIT ONE, AS THEIR CREDIT SCORING SERVICE, COULD RISK LOSING BUSINESS IF THE PROBLEM IS NOT SOLVED RIGHT AWAY.





Provide a robust model to predict the default payment on clients



Analyzed and understand what are the current situation of the bank

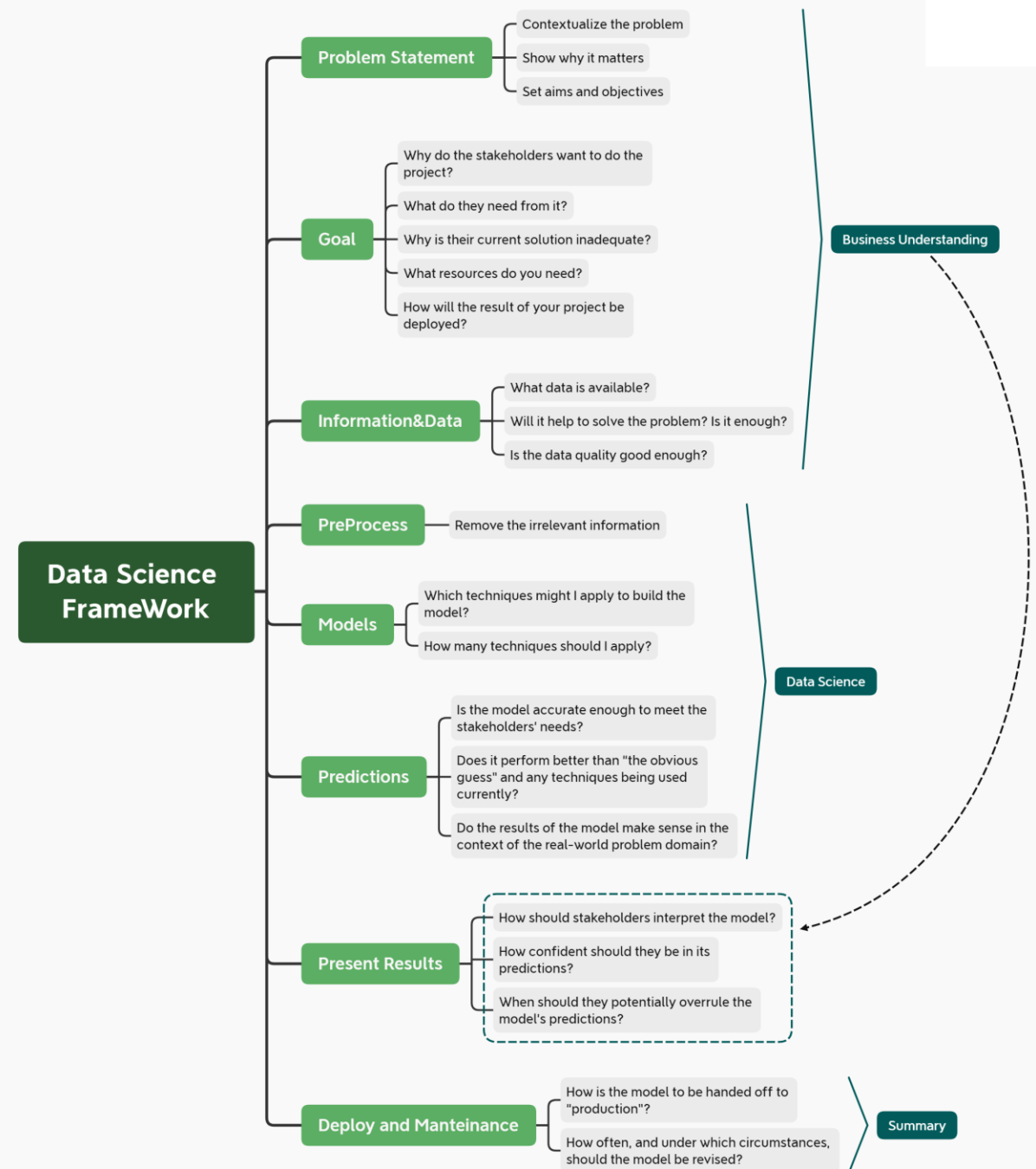


Base on data, recommendations on how improve the current situation

GOALS

DATA SCIENCE FRAMEWORK

- BASE ON THE MODEL :
 - ZUMEL AND MOUNT
- DIVIDED IN 3 GROUPS:
 - BUSINESS UNDERSTANDING
 - DATA SCIENCE (DATA MANIPULATION)
 - DEPLOY AND MAINTENANCE



DATA SCIENCE FRAMEWORK (BUSINESS UNDERSTANDING)

- PROBLEM STATEMENT:
 - INCREASE IN THE NUMBER OF CUSTOMERS WHO HAVE DEFAULTED ON LOANS THAT WILL POTENTIAL CAUSES SERIOUS PROBLEM TO THE BANK.
- GOALS:
 - ANALYZED WHAT FEATURES ARE PRODUCING THE CLIENTS DEFAULTED
 - CREATE A ROBUST MODEL TO PREDICT POSSIBLE DEFAULT CLIENTS
 - PROVIDE RECOMMENDATIONS ON HOW TO DECREASE THE CLIENT DEFAULT PAUMENT

DATA SCIENCE FRAMEWORK (DATA SCIENCE)

- **DATA UNDERSTANDING:**

- FIRST STEP IS TO UNDERSTAND THE DATA, FIND IF THERE ANY MISSING VALUE, WHAT REPRESENT THE VALUES, NUMERICAL, CATEGORICAL, THE INFORMATION IS RELEVANT TO THE PROBLEM STATEMENT.

- **PRE-PROCESS**

- TRANSFORM THE DATA THAT IS NEEDED, IS THE DATA NEED TO BE ESCALATED, OR NEED TO BE REDUCED, MAYBE NEED TO BE BINNED

- **FEATURE ENGINEERING**

- SPLIT THE DATA SETS, GROUP BY CATEGORY, PRE PROCESS IT, MAYBE APPLY PCA FOR NUMERICAL VALUES I.E AMOUNT OF BILL STATEMENT.

DATA SCIENCE FRAMEWORK (DATA SCIENCE)

- MODELS:
 - IDENTIFY IF IS REGRESSION, OR CLASSIFICATION PROBLEM. I.E IN THIS CASE IS CLASSIFICATION
 - FIND A GROUP OF MACHINE LEARNING ALGORITHMS THAT PERFORM WELL ACCORDING WITH THE SIZE OF THE DATA SET.
- EVALUATION
 - USE ACCURACY, KAPPA, CONFUSION MATRIX TO EVALUATE HOW WELL THE IS THE MODEL.
- PREDICTIONS
 - BASE ON THE MODEL AND THE RESULTS OF THE EVALUATION, MAKE AND EVALUATE THE RESULTS OF THE PERDITIONS. THIS PROCESS CAN BE ITERATIVE TOO.

DATA SCIENCE FRAMEWORK (PRESENT RESULTS)

- ONE THE MODEL IS BEEN PROVEN THE RESULTS IS CRITICAL PART FOR THE SUCCESSFUL OF THE PROJECT.
 - USE THE LEARNINGS OF THE BUSINESS UNDERSTANDING
 - USE A LOT VISUALIZATION RATHER TO GO TECHNICAL
 - TRY TO USE THE BEST PRACTICES ON “DATA TELLING”

INFORMATION OF THE DATA



Amount of the given credit

This is a numerical value,
can be binned



Gender

Needs to convert to
categorical value



Education

Categorical value,
important to see if there
any relationship between
educational levels.



Marital status

Categorical



Age (year)

This is the age of the
customer

INFORMATION OF THE DATA



History of past payment

Categorical value



Amount of bill statement

Numerical value, this can be reduced using PCA



Amount of previous payment

Numerical value, this can be reduced using PCA



client's behavior

This is the dependent value

DATA MANAGEMENT



1- Understanding Data

Need to convert the data, numerical to categorical



2-Preprocess the information

Find a method to discard irrelevant information

- Correlation, PCA, Feature variance, recursive feature elimination



3- Find the best model

Use different machine learnings models: RF,SVM,KNN C5.0

Evaluate the best performance

- Kappa, accuracy



4-Use the best model to predict the results define in our problem statement

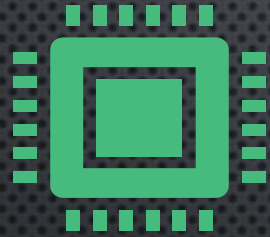
RISK DATA MANAGEMENT

- CHANGE THE CUSTOMER SENDING HABITS IS NOT POSSIBLE
- NEED TO HAVE CLEAR UNDERSTANDING OF THE TYPE OF THE DATA.
 - CONVERSION FROM NUMERICAL TO CATEGORIAL AS NEEDED
- UNDERSTAND IF THE PROBLEM IS CLASSIFICATION OR REGRESSION PROBLEM
- THE TIME THAT IS NEEDED TO BUILD THE MODELS
- HUMAN BEHAVIOR IS DIFFICULT TO MODEL

INITIAL INSIGHTS

- HISTORY OF PREVIOUS BEHAVIOR IS KEY TO CREATE A MODEL THAT CAN PREDICT THE POSSIBLE BEHAVIOR OF THE CUSTOMER.
- THE AGE OF THE CUSTOMER CAN ANOTHER FACTOR TO TAKE INTO CONSIDERATION
- GENDER, CAN BE RISKY TO ANALYZE THIS FEATURE DUE TO ANY DISCRIMINATIONS IN THE FINDINGS.
- THE AMOUNT OF GIVEN CREDIT CAN BE BINNED OR GROUPED TO HAVE BETTER UNDERSTANDING.
- THE AGE CAN BE GROUPED AS: YOUTH, YOUTH ADULT, ADULT, SENIORS.

ENVIRONMENTAL CHANGES



Python

For this project the idea is to migrate from r studio to Python.

Python is more widely known

Python have similar library “sklearn” as R
Caret”.



Start using git

Starting implemented the best practices of the industry using Git.

Git: Control version system



**THANK
YOU**

LMAURICIOCARVAJAL@GMAIL.COM