APPLIED DATA SCIENCE PROJECT

A.A 2021/2022

# PPHI PROJECT

PRESENT BY: TENKAMTE KENGNE ARSÈNE BOLIVAR AND MAMMADLI FIDAN

COMPAGNY SUPERVISORS:

◇  LUCA SCHIATTI

◇  ANH-DUNG LE

◇  GIUSEPPE GIORDANO

SUPERVISE BY:

GIUSSEPPE RIZZO

# PLAN

INTRODUCTION

DESIGN

DEVELOPMENT

MANAGEMENT

CONCLUSION
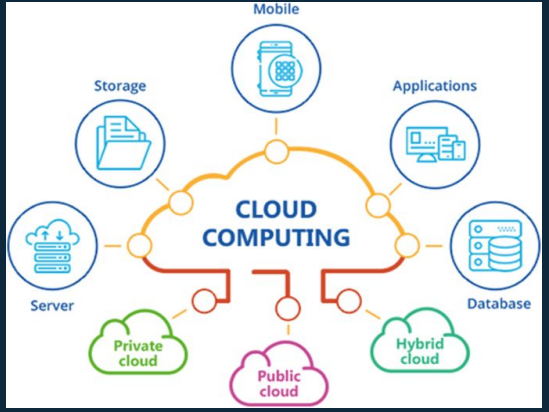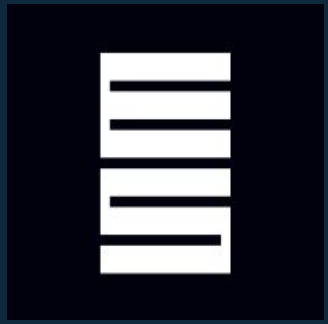
# INTRODUCTION
## OF THE PROJECT

**1**

## ML Model
predictive model to determine if their Health insurance policyholders would be also interested in a Vehicle Insurance

## TEE
A creation of a secure area that guarantees optimal protection for highly sensitive data in all states, with respect to the confidentiality and integrity

# DESIGN

2

Presentation of the general aspect with also the ongoing scenerio of the different point of the projet

PERSONA CANVAS: CROSS INSURANCE

PERSONAS CANVAS: TEE

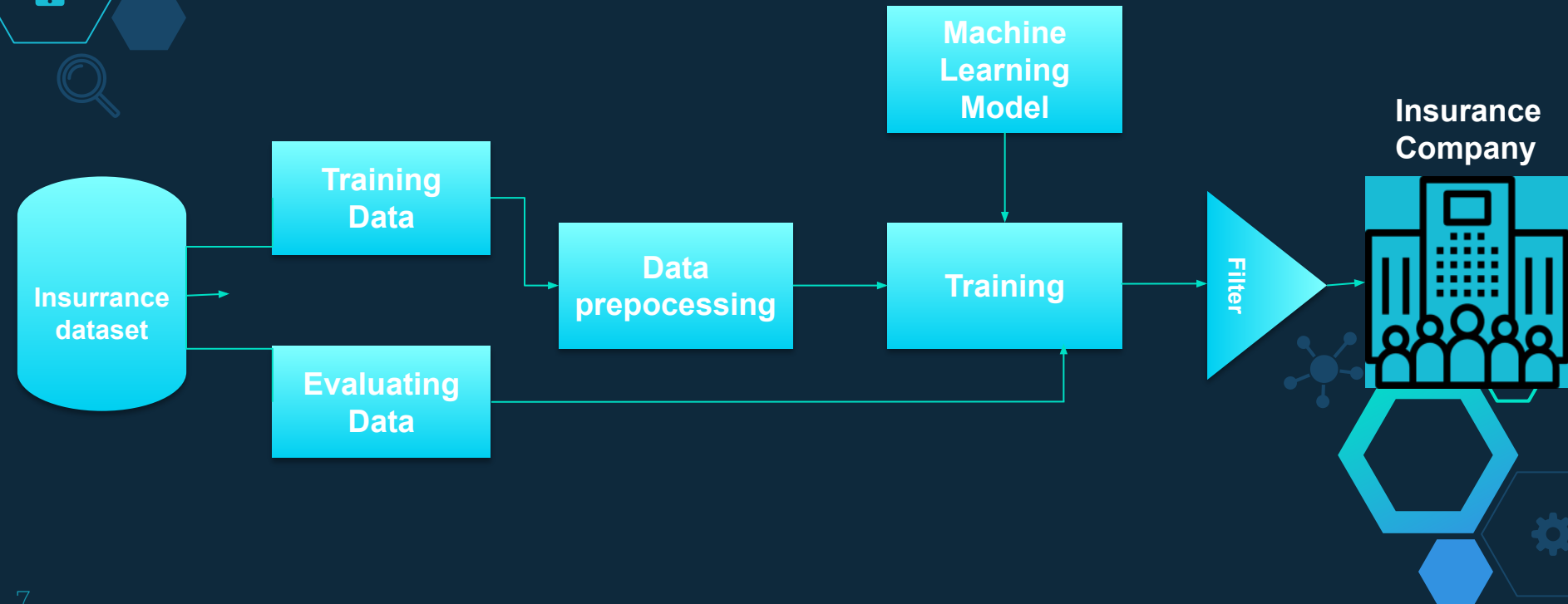FUNCTIONAL DIAGRAM: CROSS INSURANCE

FUNCTIONAL DIAGRAM: TEE

# Persona Canvas: Cross Insurance Model

| Frustrations | Needs |
|---|---|
| The company don't know who want to apply to a vehicle insurance | The companies need to optimize it's business model and revenue. |
| | A model to predict which customer will apply to the vehicle insurance |

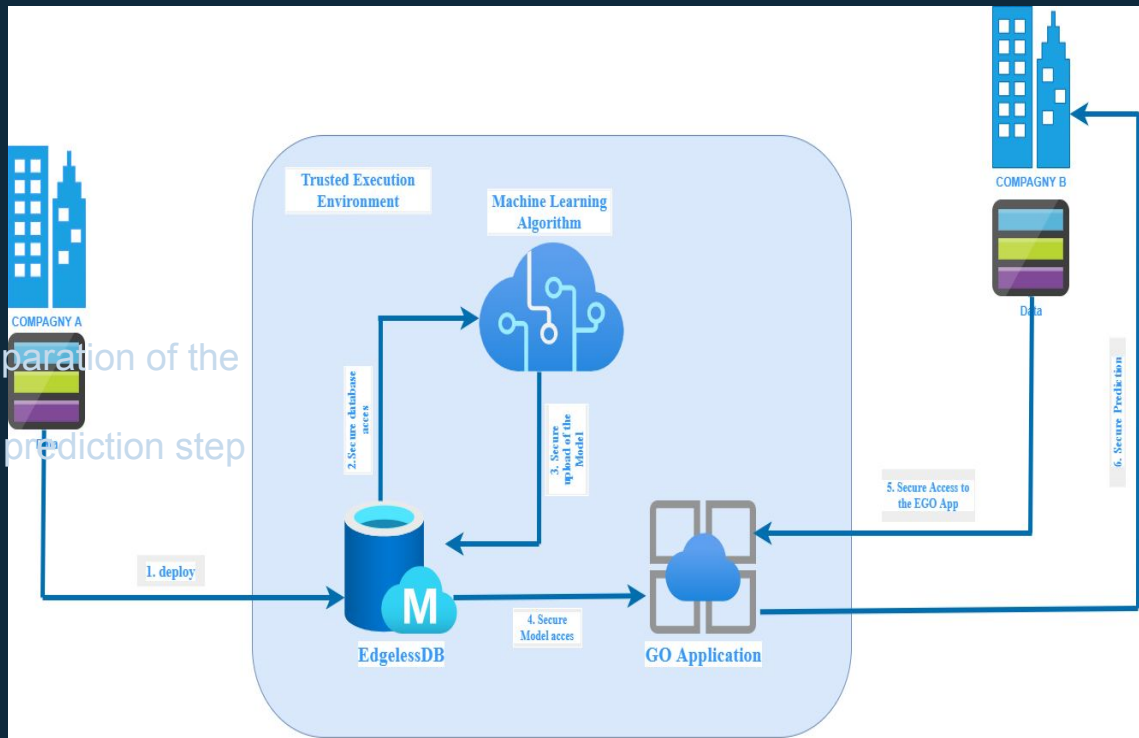**Insurance companies**

# Persona Canvas: TEE

| Frustrations | Needs |
| --- | --- |
| The absence of confidence between the different companies involve on the project.<br><br>The lack of privacy regarding the sharing of the data through the cloud. | Creation of a secure enclave environment.<br><br>Storage of the dataset on the database |
| The lack of security throw the download and the processing of the model.<br><br>The lack of privacy through the evaluation of the data. | Execution of the model in an enclave environment.<br><br>Creation of a app that will be launch into the enclave |

**Companies**

# Functional Diagram: Cross Insurance Model

# Functional Diagram: TEE



**First step**: the preparation of the model

**Second step:** the prediction step

# Development

3

DATA ANALYSIS

DATA PREPROCESSING

IMPLEMENTATION TUNNING

INTRODUCTION: TEE

CONFIGURATION AND DEPLOYMENT

PRATICAL SIMULATION

# Data Analysis

## Numerical attribute

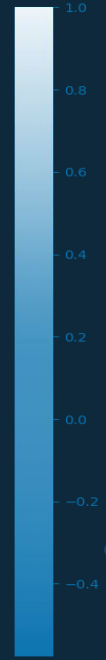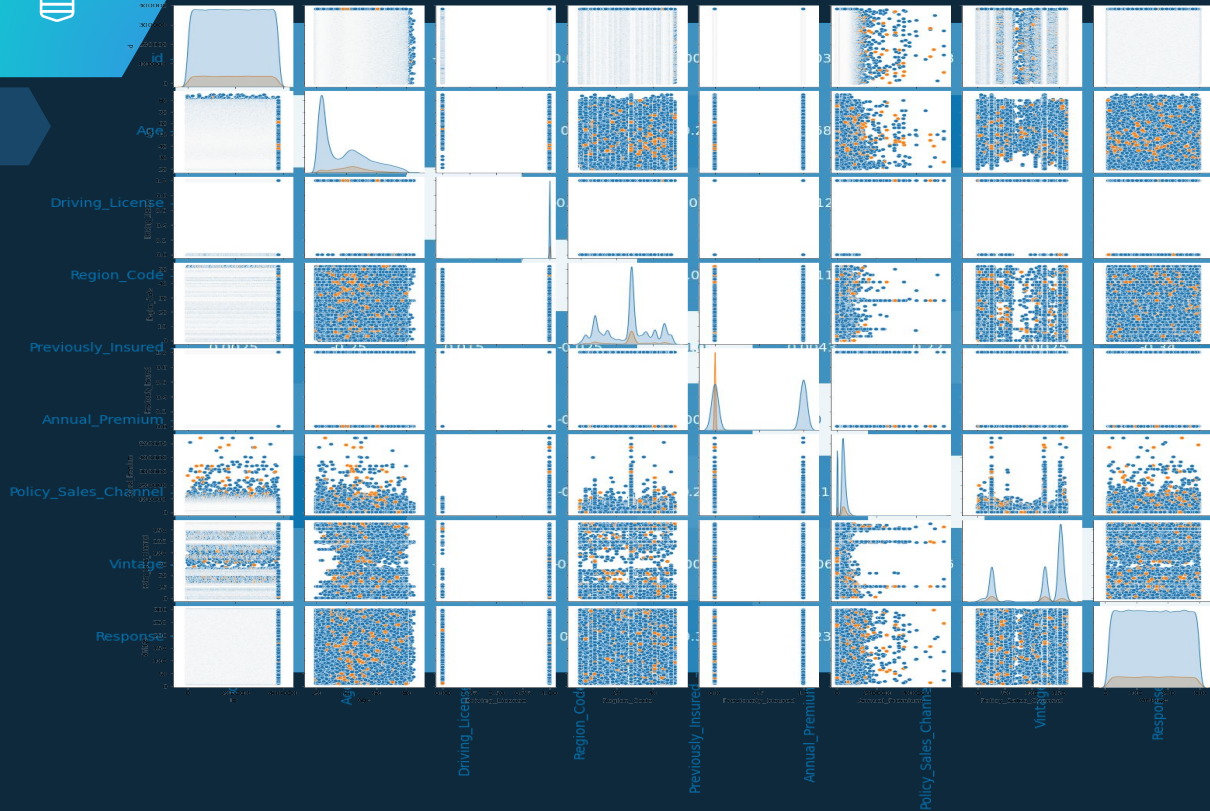◇ No missing values on the different features

◇ Few amount of uniques

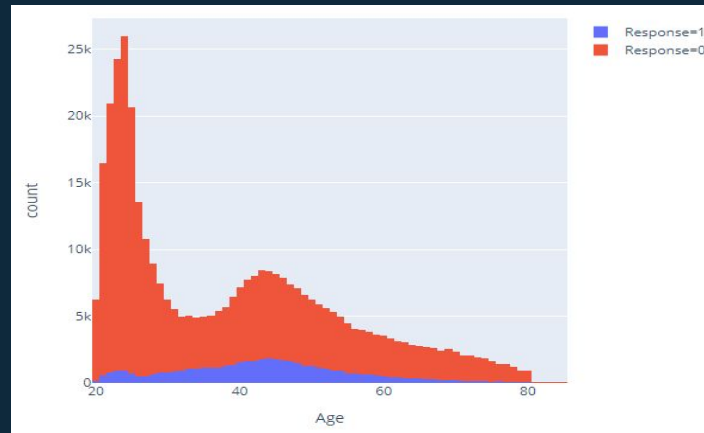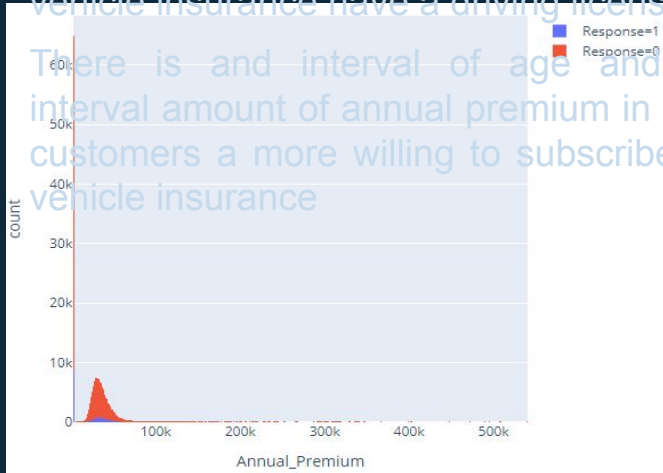| | Missing value, % | N unique value | dtype |
|---|---|---|---|
| id | 0.0 | 381109 | int64 |
| Gender | 0.0 | 2 | object |
| Age | 0.0 | 66 | int64 |
| Driving_License | 0.0 | 2 | int64 |
| Region_Code | 0.0 | 53 | float64 |
| Previously_Insured | 0.0 | 2 | int64 |
| Vehicle_Age | 0.0 | 3 | object |
| Vehicle_Damage | 0.0 | 2 | object |
| Annual_Premium | 0.0 | 48838 | float64 |
| Policy_Sales_Channel | 0.0 | 155 | float64 |
| Vintage | 0.0 | 290 | int64 |
| Response | 0.0 | 2 | int64 |
| Gender_Code | 0.0 | 2 | int8 |
| Vehicle_Age_code | 0.0 | 3 | int8 |
| Vehicle_Damage_code | 0.0 | 2 | int8 |

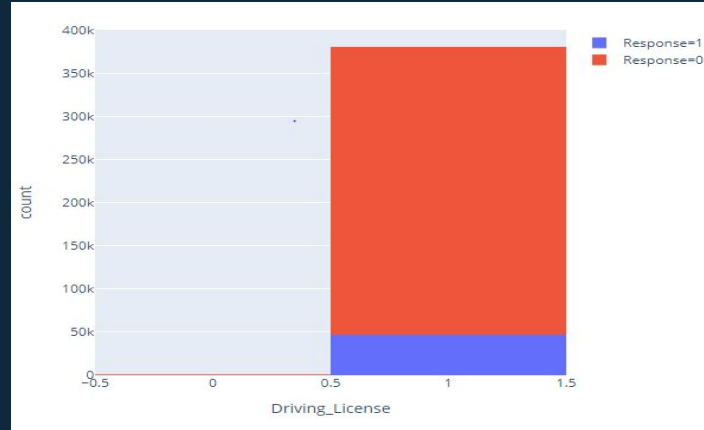# Data Analysis

## Numerical attribute

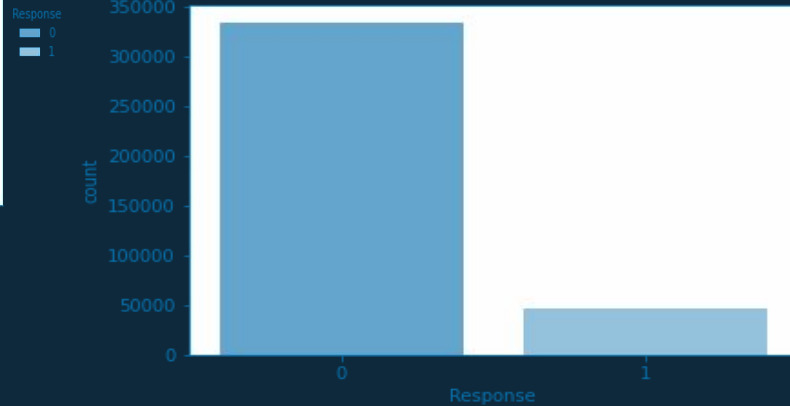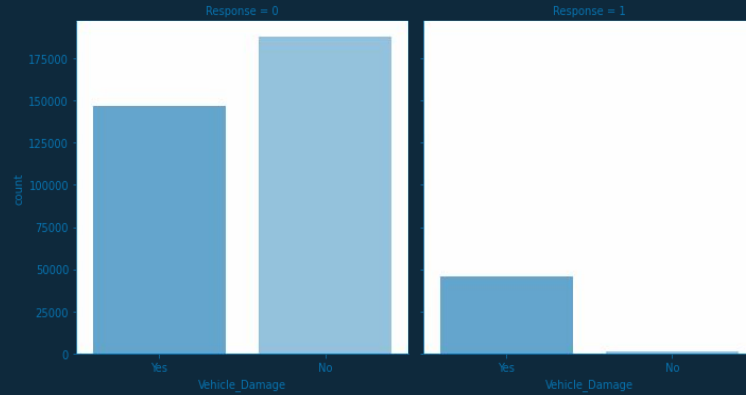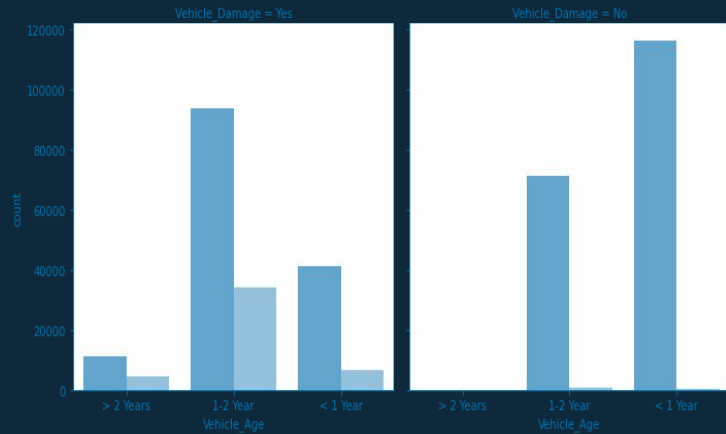# Data Analysis

## Numerical attribute

The majority of the customers that apply to vehicle insurance have a driving license

There is and interval of age and also interval amount of annual premium in which customers a more willing to subscribe to a vehicle insurance

# Data Analysis

## Categorical attribute

# Data Preprocessing

- **Dropping** of the ID feature

- Use of the **One-hot encoding** on the feature Vehicule_Age

- **Label Encoding**: with the **OrdinalEncoder** for the feature Vehicle_name and Gender

- **Spliting** of the dataset in trainning and test data

    - Unbalanced data

    - Balanced data (SMOTE)

- Four data files to utilize in the construction and the analysis of our models

**Steps involved in Data Preprocessing**

- Data Cleaning
- Data Integration
- Data Trans formation
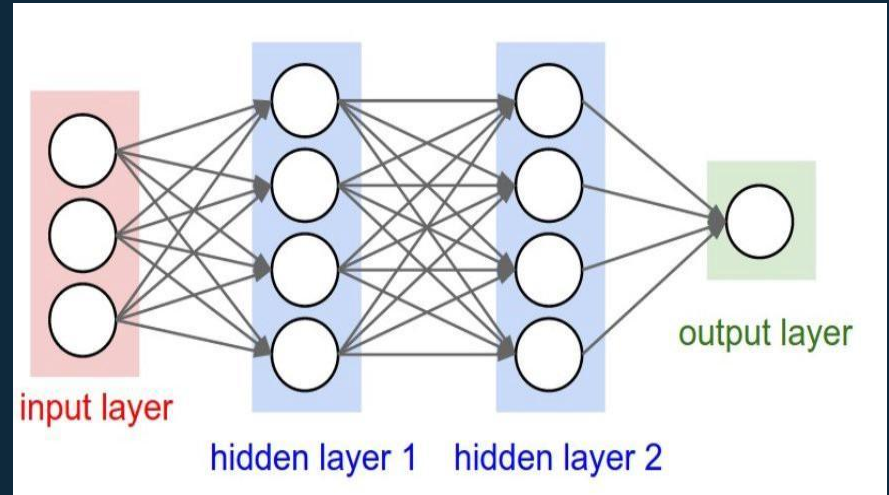- Data Reduction
- Data Discretization
- Data Sampling

14

# Normalisation

| | Gender | Age | Driving_License | Region_Code | Previously_Insured | Vehicle_Damage | Annual_Premium | Policy_Sales_Channel | Vintage |
|---|---|---|---|---|---|---|---|---|---|
| count | 3.811090e+05 | 3.811090e+05 | 3.811090e+05 | 3.811090e+05 | 3.811090e+05 | 3.811090e+05 | 3.811090e+05 | 3.811090e+05 | 3.811090e+05 |
| mean | 8.502237e-16 | -9.258629e-16 | -2.971048e-15 | -6.254366e-16 | -2.781409e-15 | 1.470378e-15 | -5.018194e-16 | -4.104990e-15 | -6.921441e-17 |
| std | 1.000001e+00 | 1.000001e+00 | 1.000001e+00 | 1.000001e+00 | 1.000001e+00 | 1.000001e+00 | 1.000001e+00 | 1.000001e+00 | 1.000001e+00 |
| min | -1.085134e+00 | -1.213453e+00 | -2.164130e+01 | -1.994638e+00 | -9.196380e-01 | -1.009801e+00 | -1.622853e+00 | -2.048455e+00 | -1.725174e+00 |
| 25% | -1.085134e+00 | -8.911132e-01 | 4.620794e-02 | -8.608404e-01 | -9.196380e-01 | -1.009801e+00 | -3.578308e-01 | -1.531887e+00 | -8.646631e-01 |
| 50% | 9.215448e-01 | -1.819661e-01 | 4.620794e-02 | 1.217845e-01 | -9.196380e-01 | 9.902940e-01 | 6.417254e-02 | 3.867931e-01 | -4.151927e-03 |
| 75% | 9.215448e-01 | 6.561169e-01 | 4.620794e-02 | 6.508902e-01 | 1.087384e+00 | 9.902940e-01 | 5.133064e-01 | 7.373213e-01 | 8.683108e-01 |
| max | 9.215448e-01 | 2.976962e+00 | 4.620794e-02 | 1.935861e+00 | 1.087384e+00 | 9.902940e-01 | 2.960534e+01 | 9.402586e-01 | 1.728822e+00 |

After the normalisation

◇ Decision tree

◇ Catboost

◇ Random Forest

◇ Keras Sequential Model

◇ Multilayer Perceptron

# Tuning

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.90 | 0.95 | 0.92 | 100320 |
| 1 | 0.38 | 0.24 | 0.30 | 14013 |
| | | | | |
| accuracy | | | 0.86 | 114333 |
| macro avg | 0.64 | 0.59 | 0.61 | 114333 |
| weighted avg | 0.84 | 0.86 | 0.85 | 114333 |

| | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.94 | 0.84 | 0.89 | 62157 |
| 1 | 0.80 | 0.93 | 0.86 | 43313 |
| | | | | |
| accuracy | | | 0.88 | 105470 |
| macro avg | 0.87 | 0.89 | 0.88 | 105470 |
| weighted avg | 0.89 | 0.88 | 0.88 | 105470 |

To find the best parameters of the different models we run a GridsearchCV with a number of cv=3 and which give us the opportunity to bypass the training and check the score on the testing data (validation data)
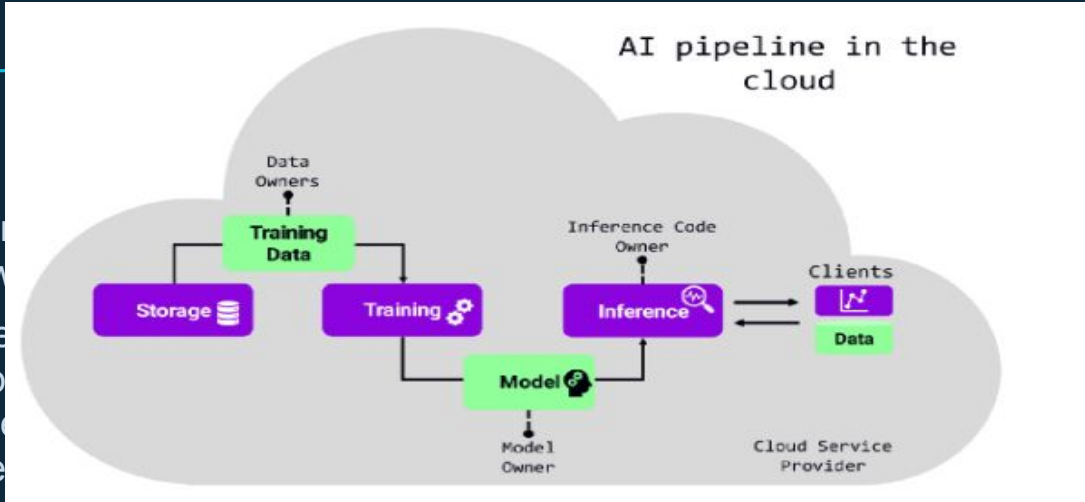
# Implementation Tuning

| | | |
|---|---|---|
| **Random Forest** | *max_features* = ['auto']<br>*criterion* = ['mse']<br>*bootstrap* = [False]<br>*n_estimators* = [ 200] | **0.97712452288843219** |
| **CatBoost** | *n_estimators =[300]* = [100, 200, 500, 1000]<br>*depth =[5]*estimators = [100, 200, 500, 1000]<br>*learning_rate =[0.01]*,6,4,5,7,8,9,10]<br>*border_count =[50]*e =[0.03,0.001,0.01,0.1,0.2,0.3]<br>*ctr_border_count=[100]*32,5,10,20,50,100,200]<br>*thread_count = 4_count=[50,5,10,20,100,200]*<br>*thread_count = 4* | **0.981938** |
| **Keras Sequential model** | *epoch=[100]*<br>*InputLayer=[64,32,16,8]* | **0.9575998937222543** |
| **Multilayer Perceptron** | *Hidden_layer_sizes = (50,50,50)*<br>*Activation = ['relu']*Activation = ['relu']<br>*Solver =['adam']*Solver =['adam']<br>*Alpha = [0.0001]*Alpha = [0.0001]<br>*learning_rate = ['adaptive']*learning_rate = ['adaptive'] | **0.9634193509552655** |

# W

◇ Ther
betw

◇ Ente

◇ cryp
prot
type
and exploits.



AI pipeline in the cloud

Data Owners

Training Data

Storage

Training

Inference Code Owner

Inference

Clients

Data

Model

Model Owner

Cloud Service Provider

20

# Configuration

## FSGSBASE

which is a feature in recent processors which allows direct access to the FS and GS segment base addresses
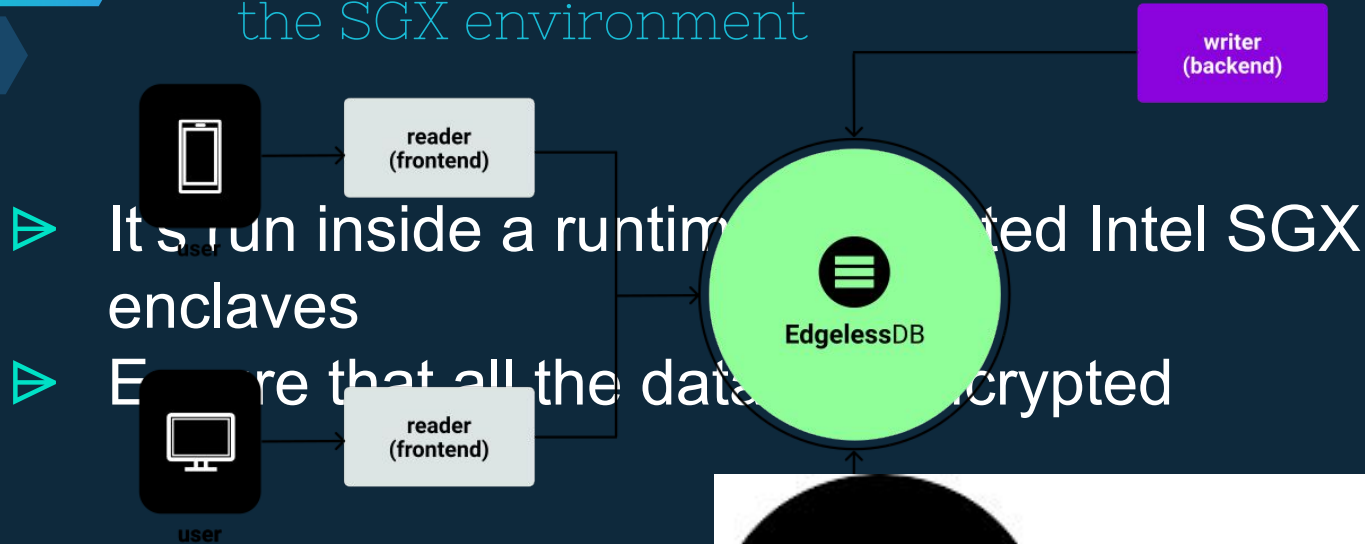
## Intel SGX driver

To improve data protection and enhance application code security

## EdgelessRT

An SDK for Trusted Execution Environments (TEE)

# EdgelessDB

EdgelessDB is a SQL database architected for the SGX environment



- ▷ It's run inside a runtime rated Intel SGX enclaves
- ▷ Ensure that all the data encrypted

# EGO

**Ego** is a framework for
building confidential apps in Go

▷ Go apps always-encrypted run in
encrypted enclaves on Intel
SGX-enabled hardware
▷ It's simplified enclave deployment

  ▷ Ego-go

  ▷ Ego

# Practical simulation

Run the system

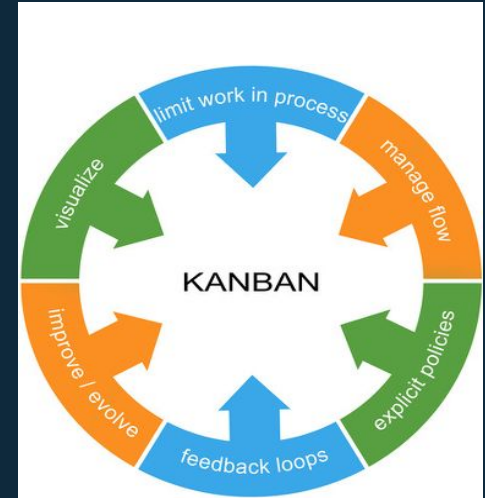Test the Database

Test the model

# MANAGEMENT

**4**

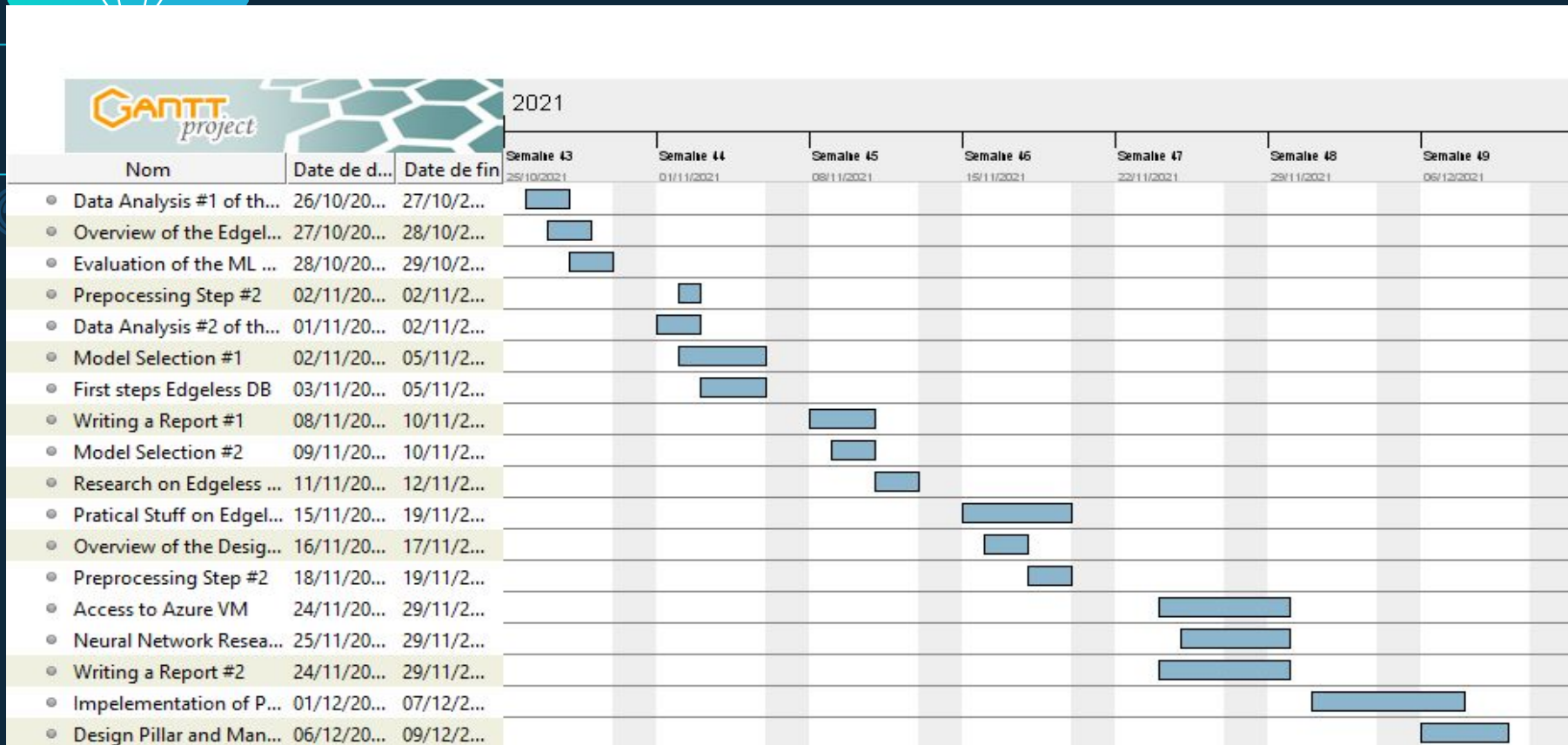How do we solve the project ? Which has been utilize ?

COLLABORATION TOOLS

GANTT DIAGRAMM

# COLLABORATION TOOLS



◇ **Kanbanflow** : for the follow up of the team work collaboration.
◇ **Scrum poker**: to estimate the complexity and effort of the different task.
◇ **Google Collab:** online cloud-based Jupyter notebook environment that allows us to train our machine learning and deep learning models on CPUs, GPUs, and TPUs.

# GANTT DIAGRAMM

# 1

## CONC

Just one last thing...

# Thanks!

◇ For your time
◇ And for your attention

## Any questions?