**TEAM MEMBERS**

María José García
Daniel Murillo
Paul Garay
Roberto Barrón

# PREDICTING BREAST CANCER

## with Machine Learning

https://github.com/MajoGarciaMontes/FINAL-PROJECT

# CONTENT

**01**
Objective / main questions

**02**
Architecture of our ML solution

**03**
Data gathering and cleansing

**04**
ML model design

**05**
API integration

**06**
Visualization / Dashboard

**07**
Conclusions / Next steps

# OBJECTIVE / MAIN QUESTIONS

## CONTEXT

- *'A woman born today has about a 1 in 8 chance of being diagnosed with breast cancer at some time during her life'*
  *National Cancer Institute*

- *The costs of having more advanced tests done to determine whether a tumor is benign or malignant can go up exponentially, so not every patient will be able to have them done.*

### OBJECTIVE

**Develop a Machine learning (ML) model and train it with relevant breast cancer variables to predict the probability of developing breast cancer.**

### MAIN QUESTIONS TO SOLVE

- **What classification fits the patient? (benign / malignant tumor).**
- **Which is the best model to use?**
- **What is the accuracy of the classification?**
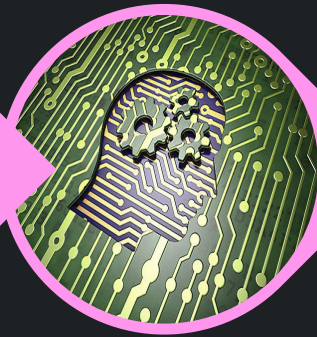
# ARCHITECTURE OF OUR ML SOLUTION

During the last 6 months we have learned to create data solutions:
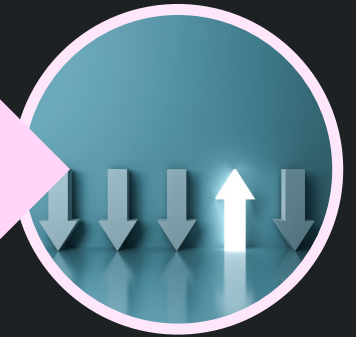


| PROJECT PROPOSAL | ETL | ML DESIGN | API | DASHBOARD CODING | USER TESTING | PROJECT EVALUATION |

# DATA GATHERING / CLEANSING

## GATHERING

- We were looking for a dataset with relevant information and with a significant size to train the data with ML.

### WHERE?
University of Wisconsin (DS) repository

### FORMAT
CSV file

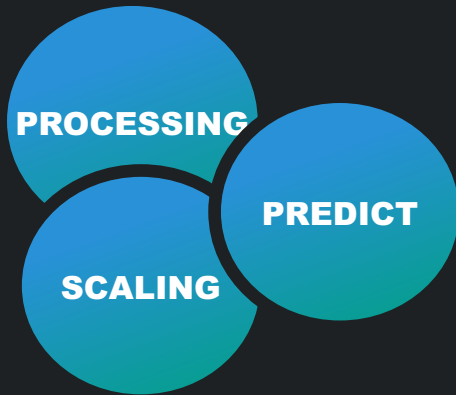### SIZE
569 datapoints

### CONTENT
30 variables of characteristics of the tumor and the final diagnosis

## CLEANSING

1) Cleansing for database creation
   - Drop of columns without data
   - Check for missing values

2) Cleansing for ML model
   - Pre-process of categorical data
   - Drop of unnecessary columns

# ML MODEL DESING

## ● LOGISTIC REGRESSION

|  | Precision | Recall | F1-score |
|---|---|---|---|
| **BENIGN** | 0.97 | 0.97 | 0.97 |
| **MALIGNANT** | 0.96 | 0.94 | 0.95 |
| **Accuracy** | 0.97 | | |

## ● RANDOM FOREST

|  | Precision | Recall | F1-score |
|---|---|---|---|
| **BENIGN** | 0.99 | 0.98 | 0.98 |
| **MALIGNANT** | 0.96 | 0.98 | 0.97 |
| **Accuracy** | 0.98 | | |

PROCESSING

SCALING

PREDICT

The model is **accurate** for both benign & malignant tumor. Its predictions are nearly always correct with high **precision** scores and the model correctly finds nearly all the true 'malignant tumors' as the **recall** scores were extremely high.
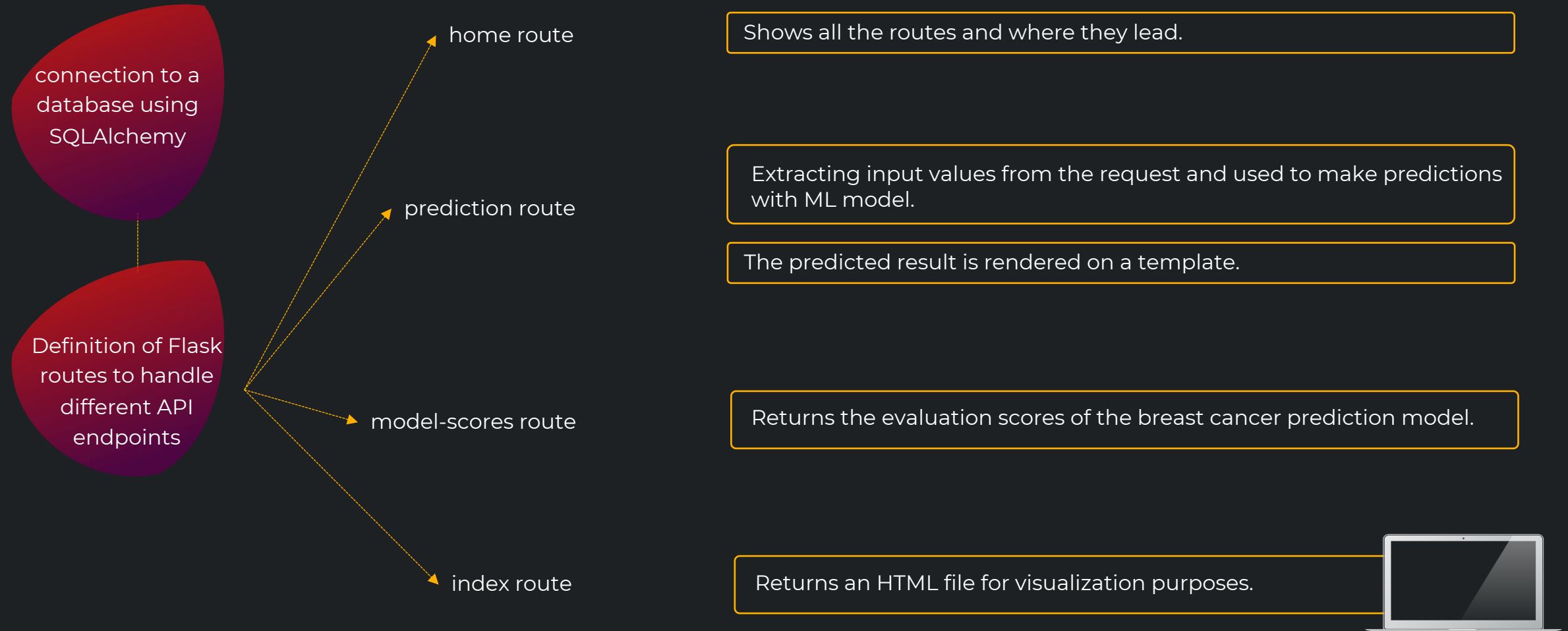
## WHY IT IS IMPORTANT TO LOOK FOR THE HIGHEST PREDICTION RECALL?
The costs of mis-classifying a 'malignant tumor' as a 'benign tumor' are extremely high. It was not acceptable that this kind of tool could misdiagnose a patient with cancer as a 'healthy patient' and send them home without treatment.
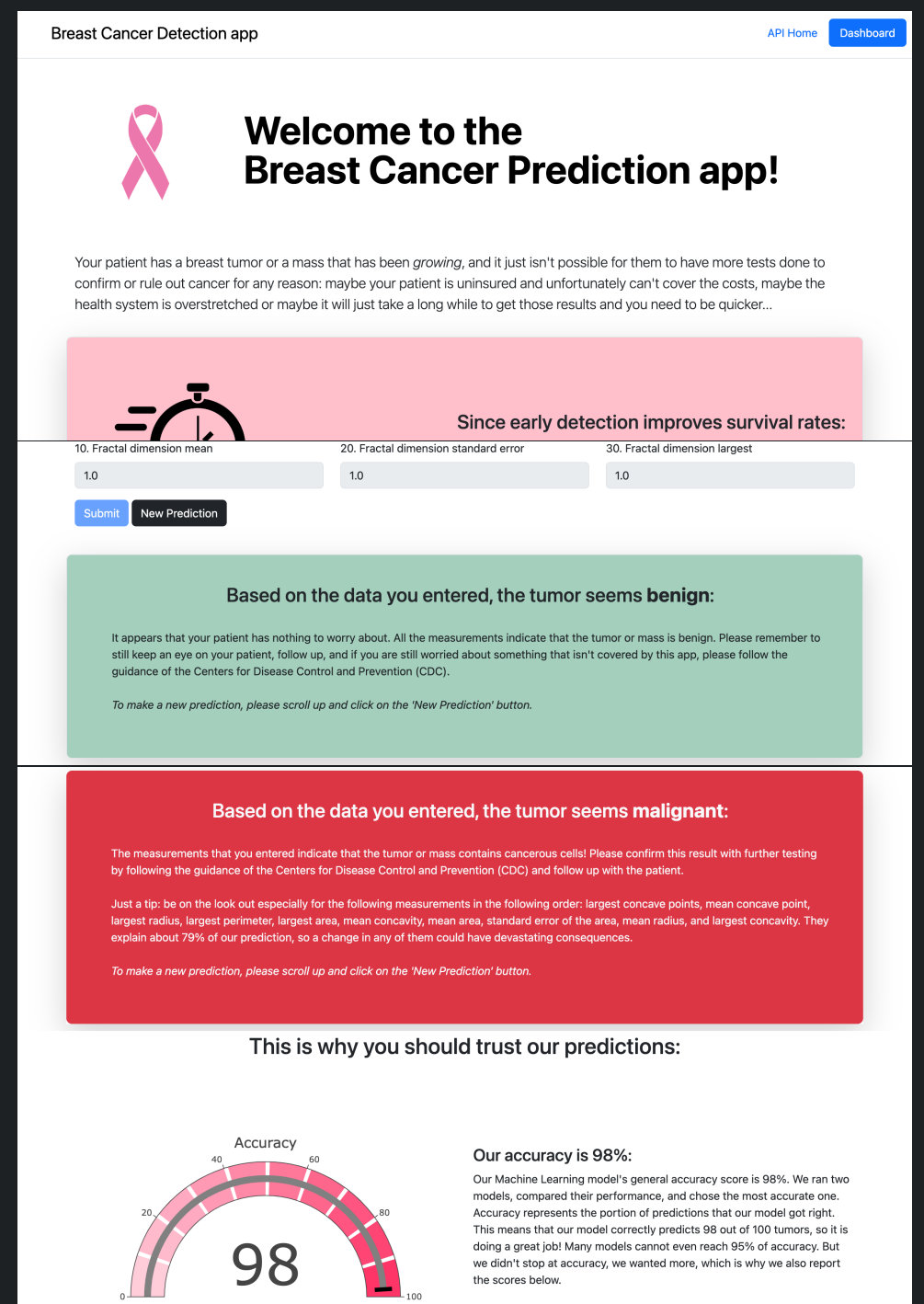
# API INTEGRATION

**WHAT?**  Flask API that allows users to submit input data for a breast cancer prediction and returns the scores of the ML model.

connection to a database using SQLAlchemy

Definition of Flask routes to handle different API endpoints

home route

Shows all the routes and where they lead.

prediction route

Extracting input values from the request and used to make predictions with ML model.

The predicted result is rendered on a template.

model-scores route

Returns the evaluation scores of the breast cancer prediction model.

index route

Returns an HTML file for visualization purposes.

# CONCLUSIONS / NEXT STEPS

- Our model was **successful** on learning with the provided dataset and developing a high level of prediction of breast cancer.

- The **random forest** model achieved the highest accuracy, precision, f1-scores and recall vs. logistic regression.

- ML has the potential to **reduce costs** while still being reliable for screening breast cancer.

## NEXT STEPS

- Having **more recent data** could improve our precision and make our model be more in touch with today.

- Having **more feature variables** could improve the model scores or even work better with other ML models, such as a Neural Network.

# ANY QUESTIONS?