



# Predictive Diagnostic Assistant for Breast Cancer Screening

---

A use case for  
machine learning  
tools in the health  
care setting by:

Anna Remler

Antonette Goroch

Daniel Rodriguez

Grace Yoo

Luis Hernandez

Pooja Rajesekharan



# About the Project



# What need does this tool address?

- Breast cancer affects about 13% (about **1 in 8**) of U.S. women. 8\*
- Misdiagnosis by medical professionals is a serious problem (46% of 2155 case were misdiagnosed in one often quoted study)
- Providing an **easy-to-access diagnostic tool** for health care professionals, which could utilize machine learning tools and existing research data to isolate key diagnostic features and direct to additional resources, could help reduce misdiagnosis.
- Since early detection is one of the best tools in reducing negative outcomes, such a tool could truly help save lives.

*\*\*Source: American Cancer Society, based on 2022 data.*





# Who & How

## Target Audience:

A wide spectrum of healthcare professionals who represent the first line of diagnostic defense in treating breast cancer. Such as:

- Lab Techs
- Nurses
- Physician Assistants



## How the tool works:

- User inputs five cell nuclei features into a web interface, based on a patient's lab data.
- System returns a suggested course of action based on the probability of malignancy based on the machine learning model.

**Input Feature Data**

concave points\_mean  
(0, 0.20)

radius\_worst  
(7.93, 36.04)

perimeter\_worst  
(50.41, 251.2)

area\_worst  
(185.2, 4254)

concavity\_worst  
(0, 1.25)

Predict



# Our Process



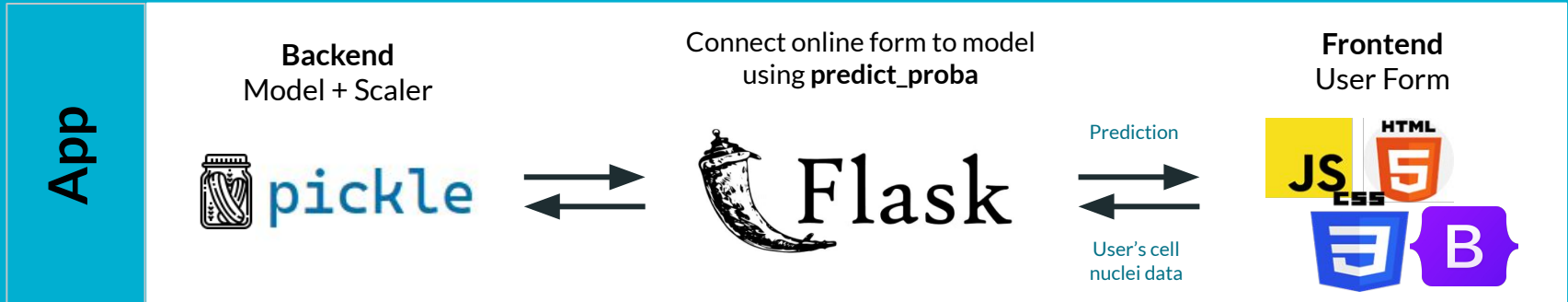
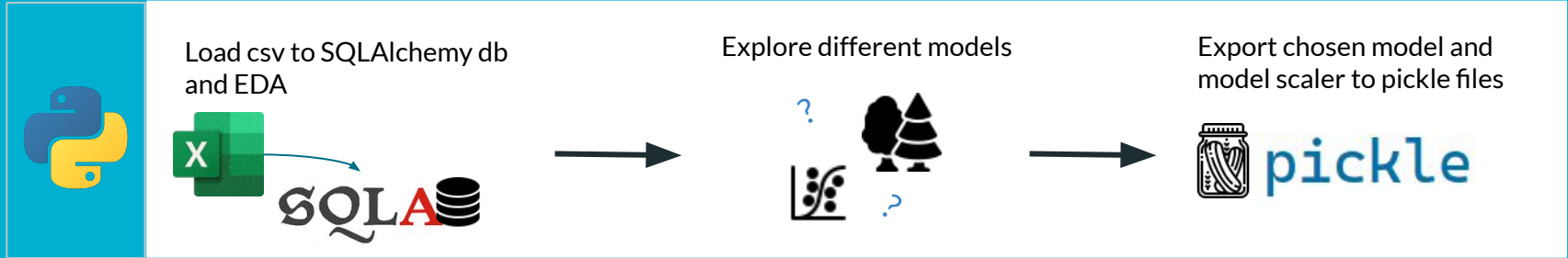
# Steps Taken + Who Did What

- Dataset Discovery (Pooja)
- Initial scoping & planning
- Data cleaning & pre-processing
  - Daniel converted the initial dataset to an SQL database.
- Exploratory analysis
  - Initial research on the dataset .
  - Histograms using Numpy (Grace)
  - Logistical Regression model (Luis)
  - Random Forest model (Grace)
- Deploy Model & Dashboard
  - Pooja, Grace & Luis worked on the Flask app and HTML to display the dashboard using Bootstraps.
- Final Presentation & Repository
  - Antonette created the final presentation deck for the repository.
  - Anna created the final ReadMe file & supporting image files.
  - Luis maintained our central Github repository.

- Data Cleaning & Preprocessing
  - PANDAS
  - SQL
- Exploratory Data Analysis
  - PANDAS
  - Numpy
  - SciKitLearn
  - Matplotlib
- Dashboard Creation
  - Pickle
  - Flask
  - HTML
  - Bootstrap



# Data Processing & Pipeline



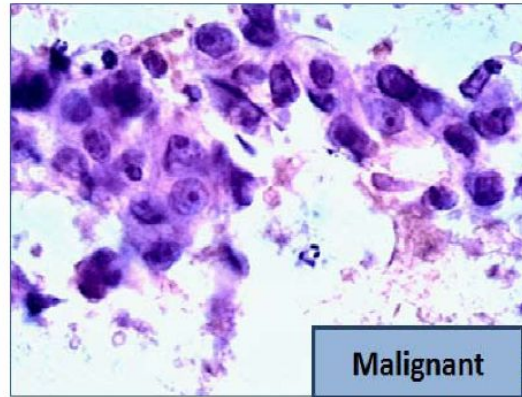
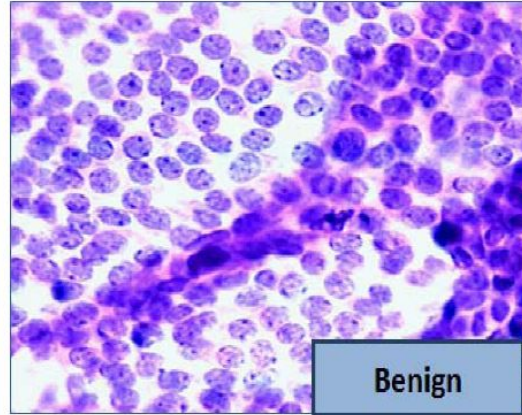


# About the Data

The data was obtained on Kaggle and uploaded by UC Irvine Machine Learning Repository but the results were found by University of Wisconsin.

The dataset was developed in 1995.

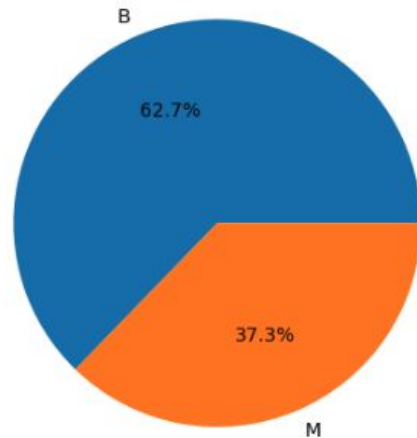
The dataset consists of 30 features—derived from 10 features by getting the mean, the standard error and the mean of the three largest values.







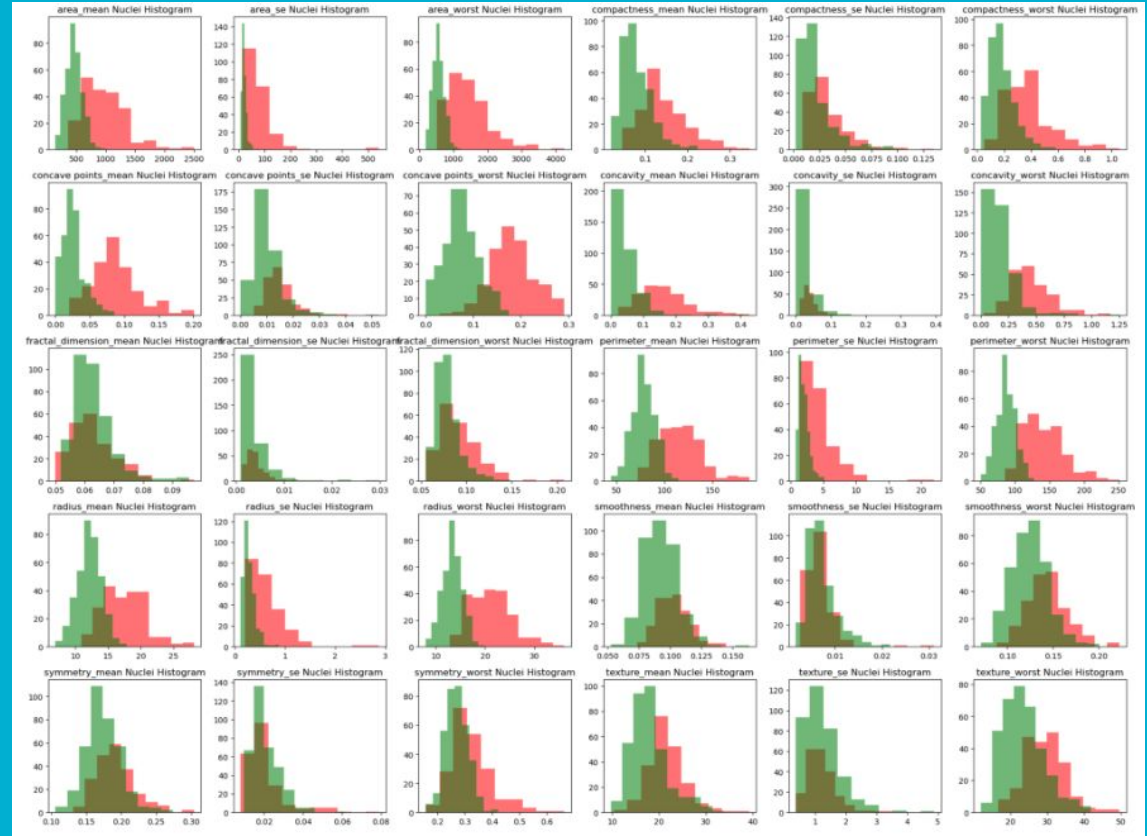
## diagnosis





# Pre-Processing & EDA

- From initial CSV, created a SQL DB
- Used numpy to run histogram analysis, showing initial favorable variables.
- Reducing to most responsive variables.





# Model Training, Tuning & Evaluation

- Goal: Find a model that reduces the number of user inputs with an acceptable level of accuracy.
- When including all 30 features, the Logistic Regression performed 1% better than Random Forests.



## Logistic Regression with All Features

	precision	recall	f1-score	support
benign	0.98	0.99	0.98	88
malignant	0.98	0.96	0.97	55
accuracy			0.98	143
macro avg	0.98	0.98	0.98	143
weighted avg	0.98	0.98	0.98	143



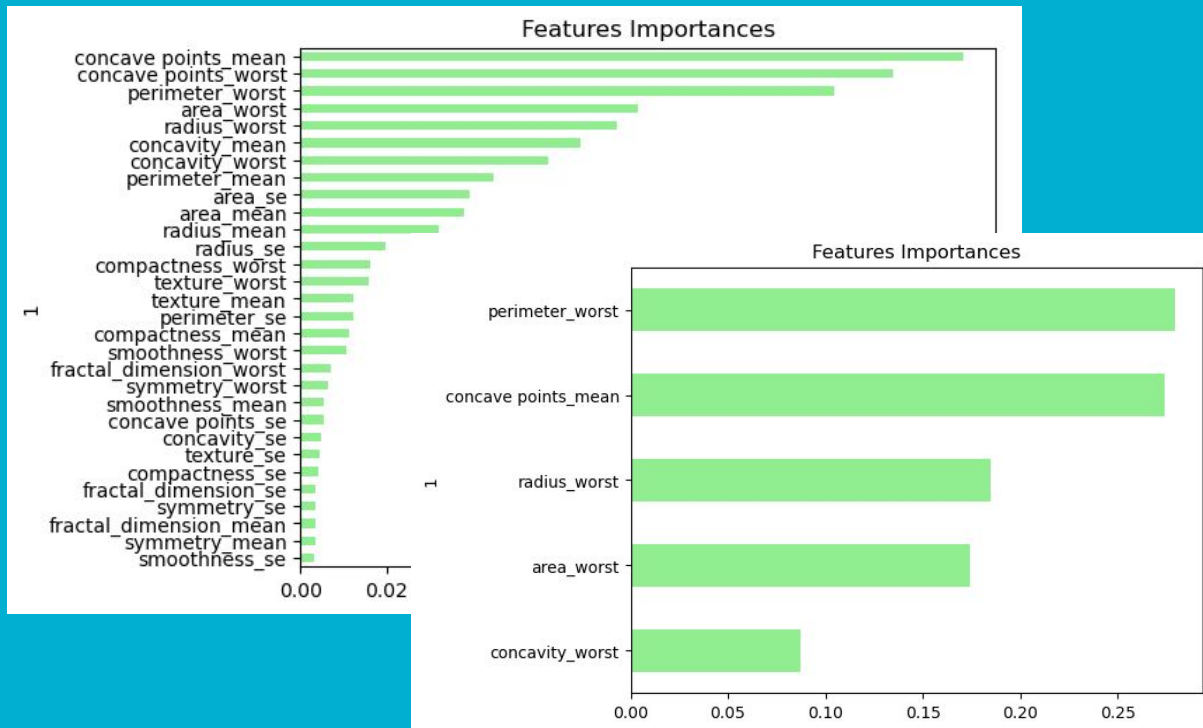
## Random Forests with All Features

	precision	recall	f1-score	support
B	0.98	0.96	0.97	84
M	0.95	0.97	0.96	59
accuracy			0.97	143
macro avg	0.96	0.97	0.96	143
weighted avg	0.97	0.97	0.97	143



# Reducing the Number of Features

- By analyzing the feature importance through sklearn and the histograms of the initial data analysis, we found the top five most important features.
- We experimented with removing features from the models without sacrificing too much accuracy.





# Model Training with Fewer Features

- Reducing to the top five features reduced the overall accuracy by 2-3% on both models.
- The models had equal accuracy at 94%.
- However, Random Forests performed 4% better than Logistic Regression at predicting malignant diagnosis.



## Logistic Regression with Fewer Features

	precision	recall	f1-score	support
benign	0.94	0.97	0.96	88
malignant	0.94	0.91	0.93	55
accuracy			0.94	143
macro avg	0.94	0.94	0.94	143
weighted avg	0.94	0.94	0.94	143



## Random Forests with Fewer Features

	precision	recall	f1-score	support
B	0.96	0.94	0.95	84
M	0.92	0.95	0.93	59
accuracy			0.94	143
macro avg	0.94	0.94	0.94	143
weighted avg	0.94	0.94	0.94	143



# Creating the Dashboard

Pickle

Flask

Java

## Predictive Diagnostic Assistant for Breast Cancer Screening

### About the Tool

Breast cancer is the leading type of cancer globally, accounting for 12.5% of all new annual cancer cases worldwide. In the United States, it is the most commonly diagnosed cancer among women. Our tool is designed to be used by a lab technician to predict whether a patient's tumor cell nuclei can be classified as benign or malignant based on specific feature data inputted. Below is a list of 5 features that were identified through analysis as the most predictive from a diagnostic standpoint.

#### Input Feature Data

concave\_points\_mean  
(0, 0.20)

radius\_worst  
(7.93, 36.04)

perimeter\_worst  
(50.41, 251.2)

area\_worst  
(185.2, 4254)

concavity\_worst  
(0, 1.25)

Predict

#### Features Description

Source: [https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+\(diagnostic\)](https://archive.ics.uci.edu/ml/datasets/breast+cancer+wisconsin+(diagnostic))

- **concave\_points\_mean**: the average count of the number of points on the nuclear border that lie on an indentation for all nuclei
- **radius\_worst**: the mean averaged length of radial line segments from the center of the nuclear mass to each of the points of the nuclear border for the three largest nuclei
- **perimeter\_worst**: mean distance around the nuclear border for the largest three nuclei
- **area\_worst**: mean area within the outlined nuclear border for the three largest nuclei
- **concavity\_worst**: mean severity of concave portions on the contour for the three largest nuclei



# Challenges

---

- The main limitation we encountered was in the dataset itself, since it was somewhat limited in both size and scope (limited to one state, Wisconsin).



## Next Steps

---

Several other areas of development would add significant value to the app:

- Refine the tool with additional/better data (i.e. data from more states, more current data).
- Further tools which build on the given outcomes (i.e. 30% chance of malignancy leads to a prompt for follow up with a specialist, or 3% might suggest retesting)
- Further research about the best type of health personnel to utilize the tool.





# Thank you!

**Special Thanks to Kevin and Mounika for the help!**  
**Congratulations to everyone!**