# Computational Biology

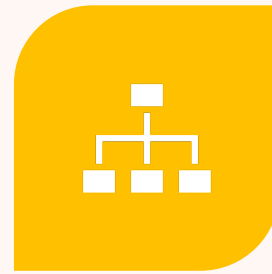Breast Cancer classification

# Outline

INTRODUCTION

DATASET DESCRIPTION
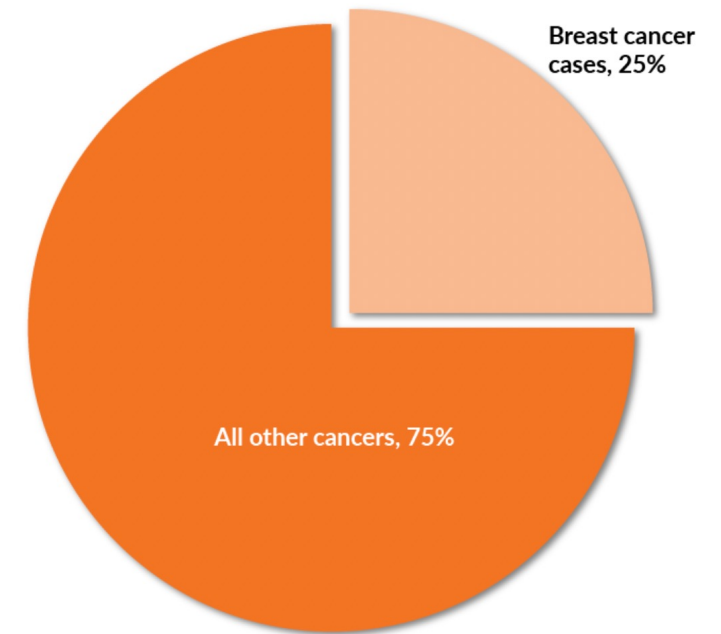
PROJECT STRUCTURE

IMPORTANT DATES

# Breast Cancer

It is estimated that in 2022:

- 28,600 Canadian women will be diagnosed with breast cancer. This represents 25% of all new cancer cases in women in 2022.

- 5,500 Canadian women will die from breast cancer. This represents 14% of all cancer deaths in women in 2022.

- On average, 78 Canadian women will be diagnosed with breast cancer every day.

- On average, 15 Canadian women will die from breast cancer every day.

- 270 Canadian men will be diagnosed with breast cancer and 55 will die from breast cancer.

**Percentage of All Estimated New Cancer Cases in Women in 2022**

Breast cancer cases, 25%

All other cancers, 75%
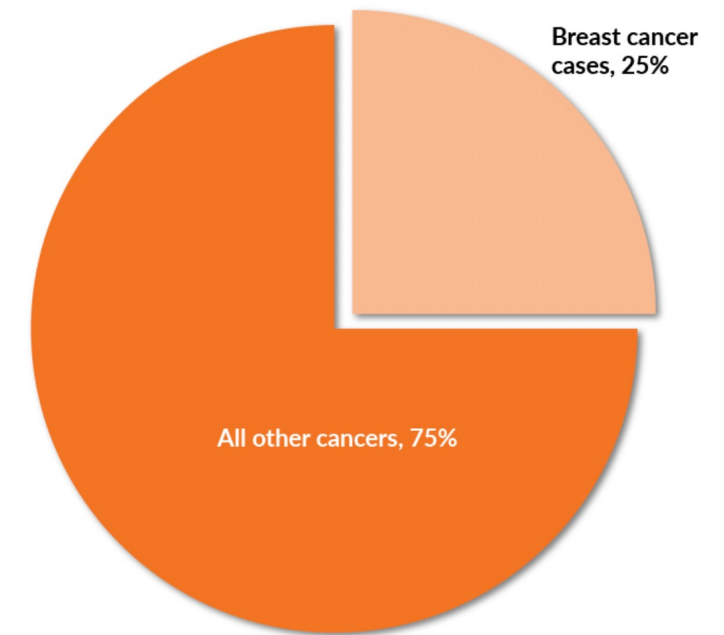
© Canadian Cancer Society

# Breast Cancer

It is estimated that in 2022:

- 28,600 Canadian women will be diagnosed with breast cancer. This represents 25% of all new cancer cases in women in 2022.

- 5,500 Canadian women will die from breast cancer. This represents 14% of all cancer deaths in women in 2022.

- On average, 78 Canadian women will be diagnosed with breast cancer every day.

- On average, 15 Canadian women will die from breast cancer every day.

- 270 Canadian men will be diagnosed with breast cancer and 55 will die from breast cancer.
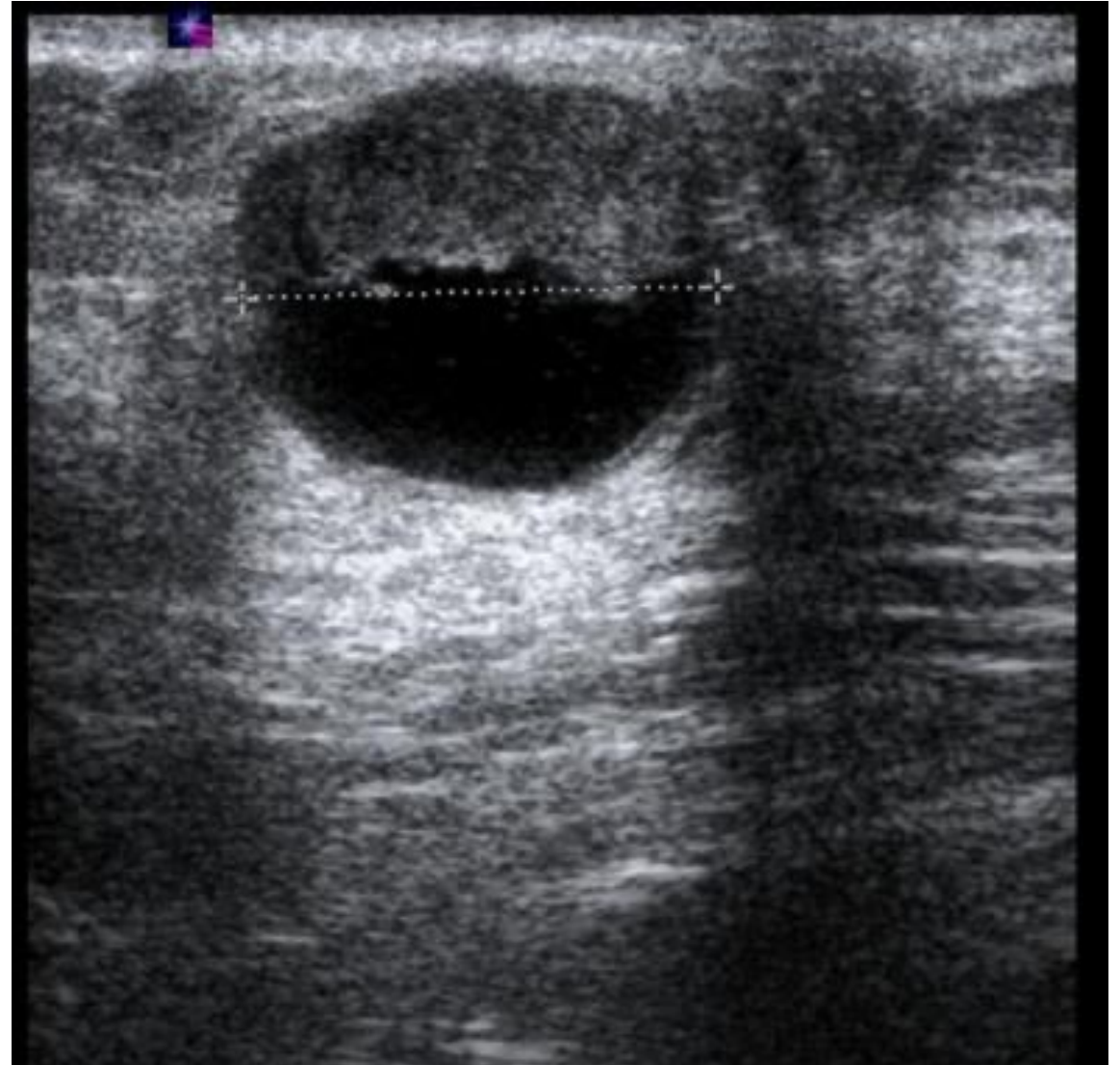
Motivation: the early diagnosis of breast cancer can improve the prognosis and chance of survival significantly, as it can promote timely clinical treatment to patients.

Percentage of All Estimated New Cancer Cases in Women in 2022

Breast cancer cases, 25%

All other cancers, 75%
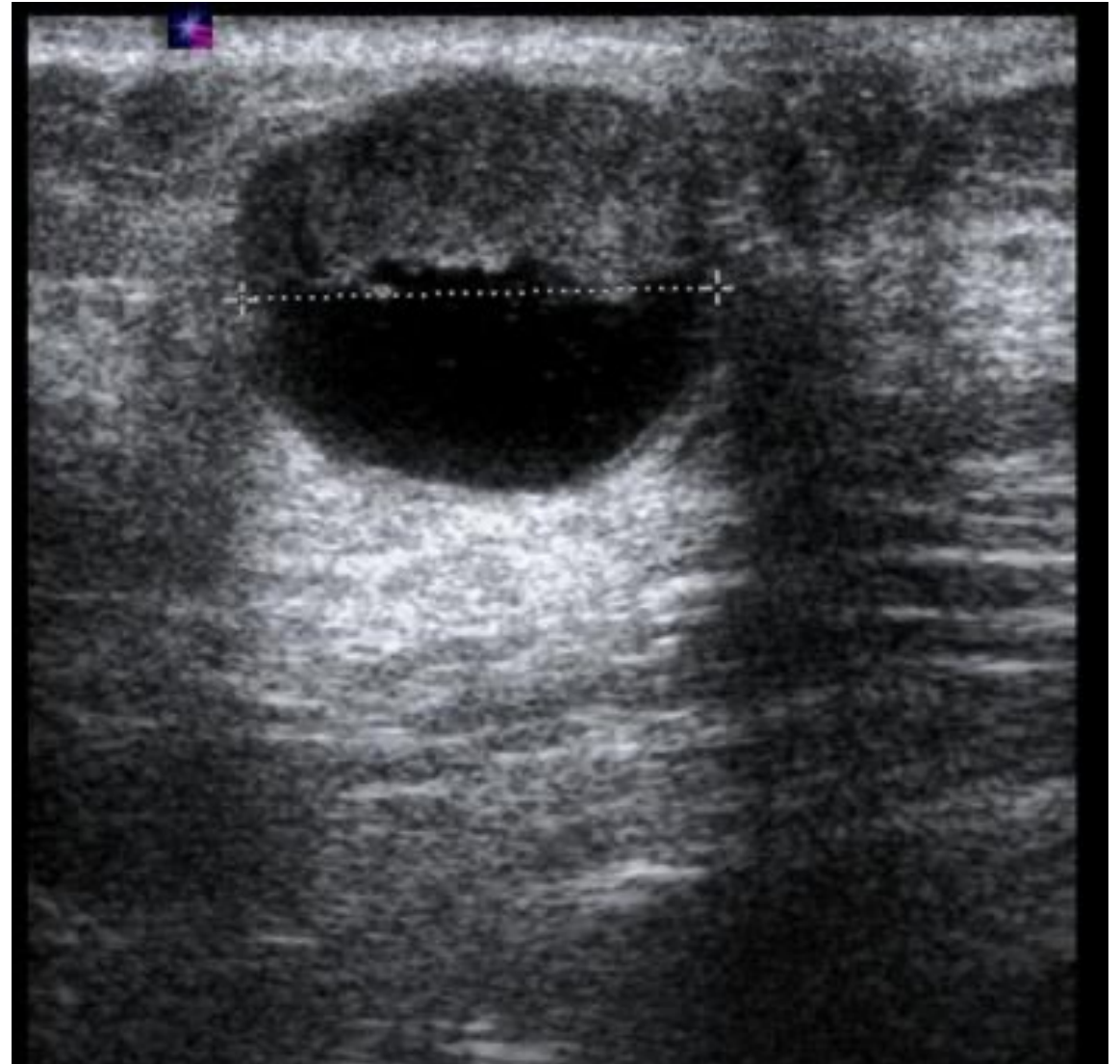
© Canadian Cancer Society

# Breast Cancer Wisconsin (Diagnostic) Data Set

- The dataset is publicly available and was created by Dr. William H. Wolberg, physician at the University Of Wisconsin Hospital at Madison, Wisconsin, USA.

- To create the dataset Dr. Wolberg used fluid samples, taken from patients with solid breast masses and an easy-to-use graphical computer program called Xcyt
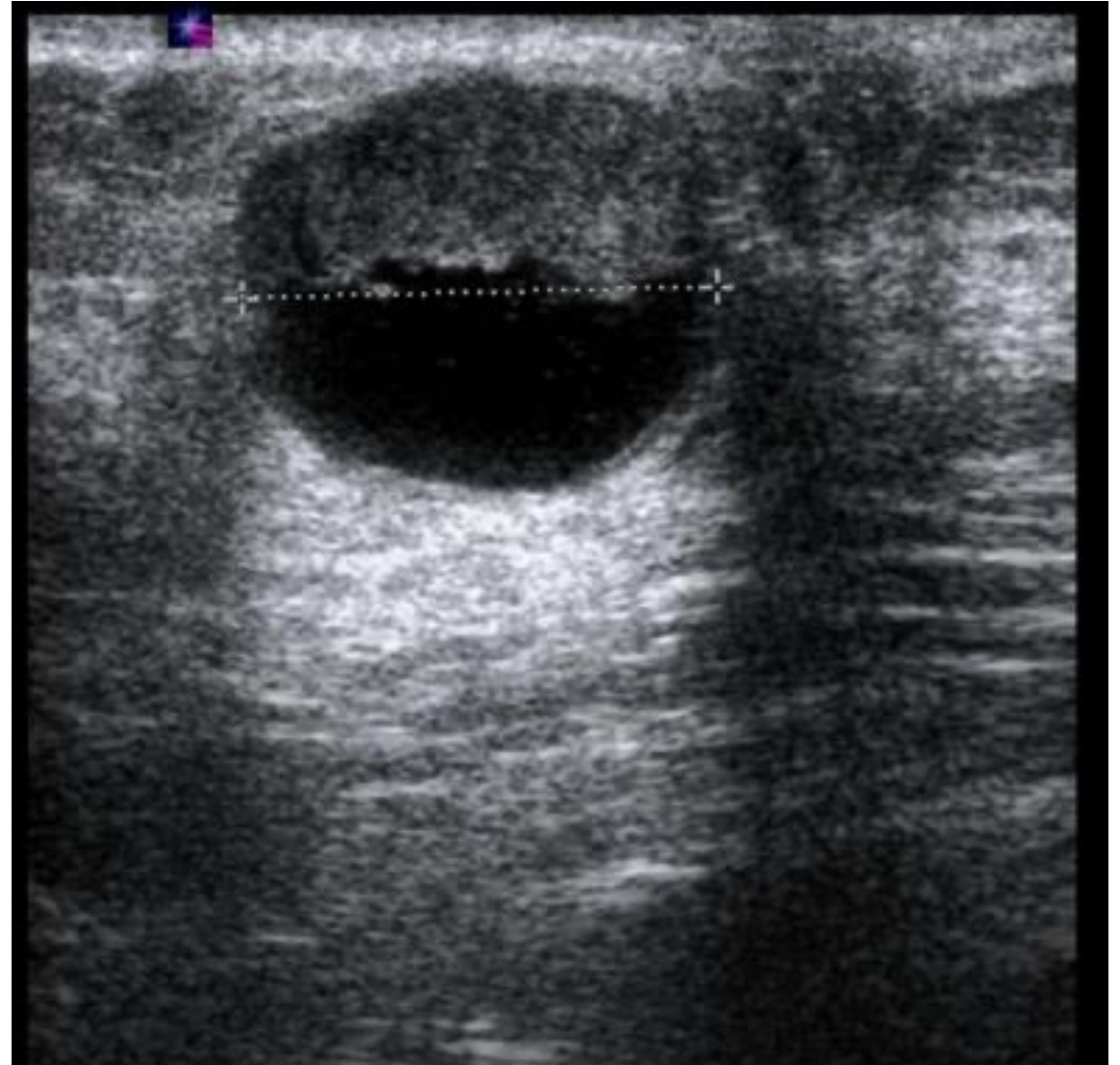
# Breast Cancer Wisconsin (Diagnostic) Data Set

- ID number
- Diagnosis (M = malignant, B = benign)
- radius
- texture
- perimeter
- area
- smoothness (local variation in radius lengths)
- compactness (perimeter² / area — 1.0)
- concavity (severity of concave portions of the contour)
- concave points (number of concave portions of the contour)
- symmetry
- fractal dimension

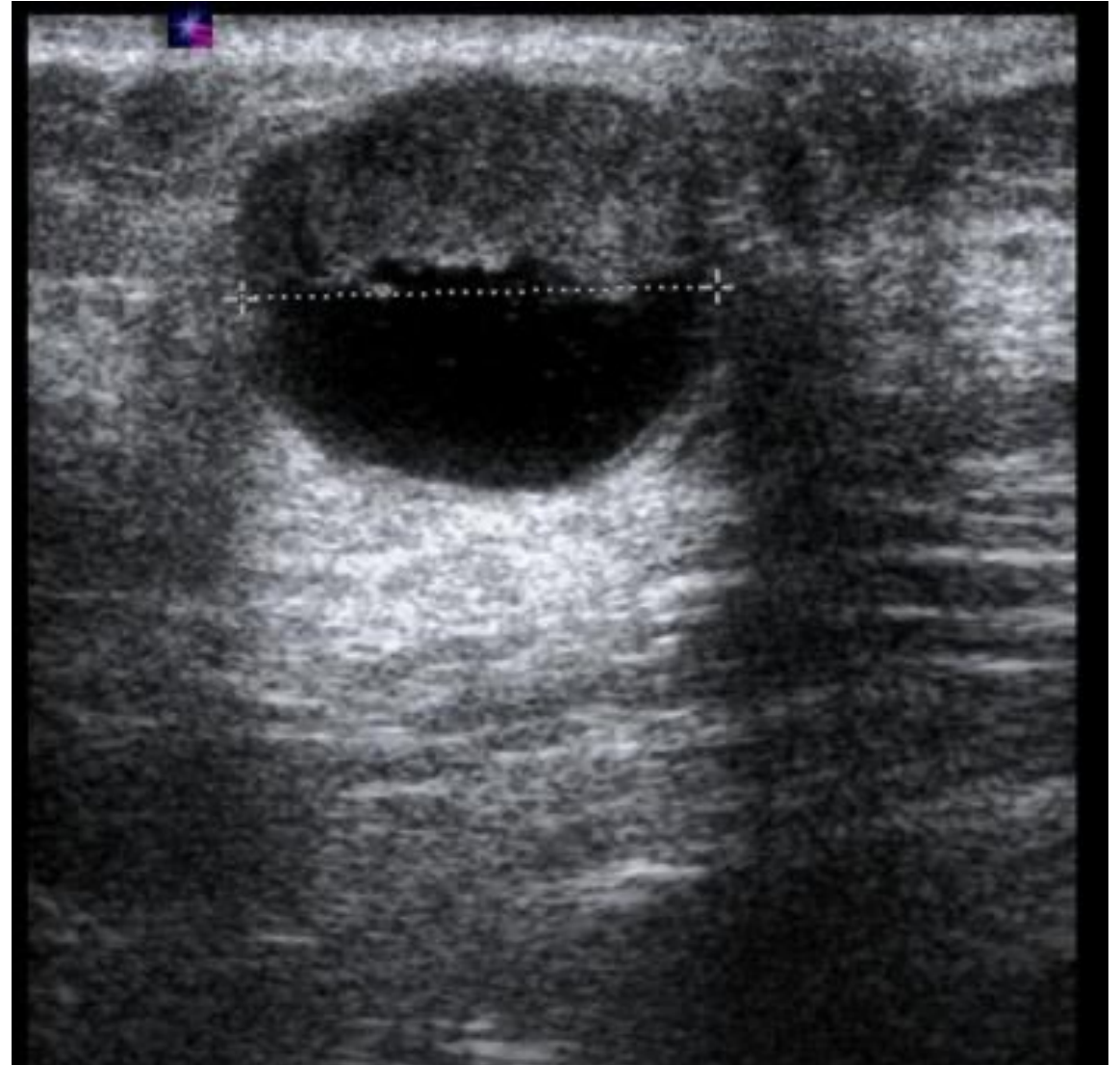# Breast Cancer Wisconsin (Diagnostic) Data Set

- 33 columns

- 569 rows (357 benign, 212 malignant)

- Numerical (int, float), categorical (b, m)

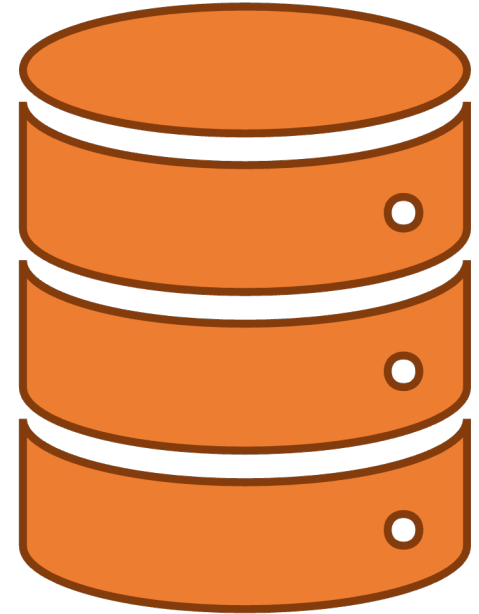# Breast Cancer Wisconsin (Diagnostic) Data Set

- 33 columns

- 569 rows (357 benign, 212 malignant)

- Numerical (int, float), categorical (b, m)

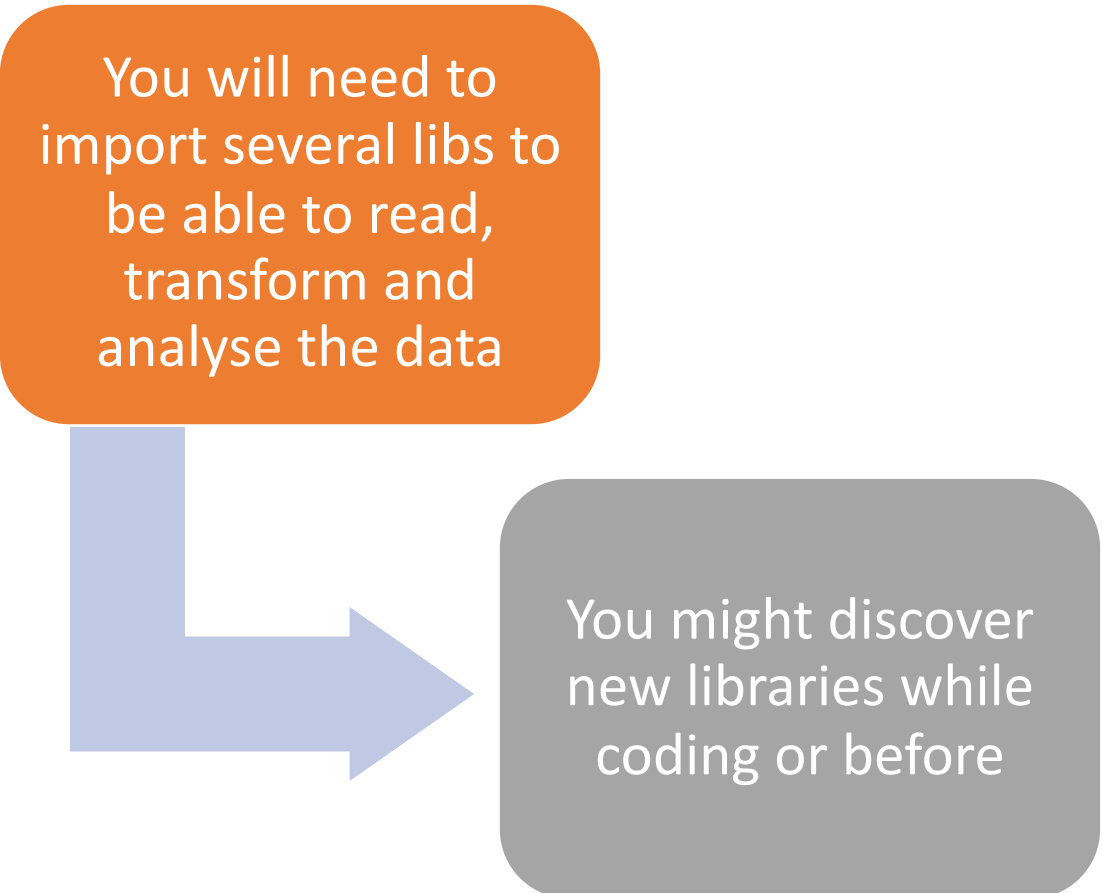Aim: we want to be able to classify the diagnosis as benign or malignant

# Project Structure

- Import libraries & dataset
- Preprocess data
- Analyze data
- Choose features
- Create and evaluate models
- Choose best model

# Import necessary libraries

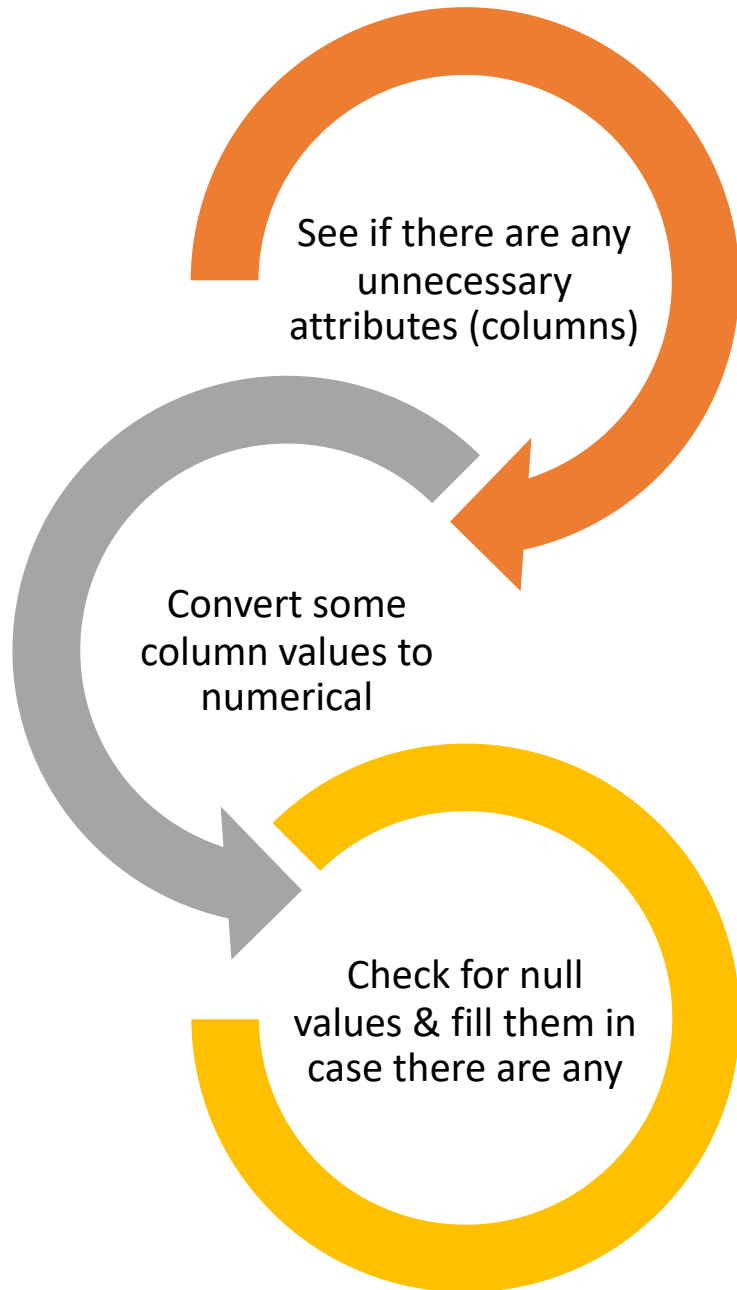You will need to import several libs to be able to read, transform and analyse the data

You might discover new libraries while coding or before

**Here are some libraries you might want to consider:**

- pandas

- numpy

- seaborn

- matplotlib.pyplot

- missingno

- sklearn

# Preprocess data

**Here are some pandas and missingno tools you might find useful:**

- df.drop(['Columns'])

- df['column'].apply(lambda val: 1 *condition* else 0)

- df.describe()

- df.info()

- df.isna().sum()

- msno.bar(df)

See if there are any unnecessary attributes (columns)

Convert some column values to numerical

Check for null values & fill them in case there are any

# Analyze data



Plot the distribution of each column

Plot the correlation matrix to see highly correlated features

# Choose features & transform data

Remove highly correlated features

Create features and labels (X and y)

Split the data into training and test set
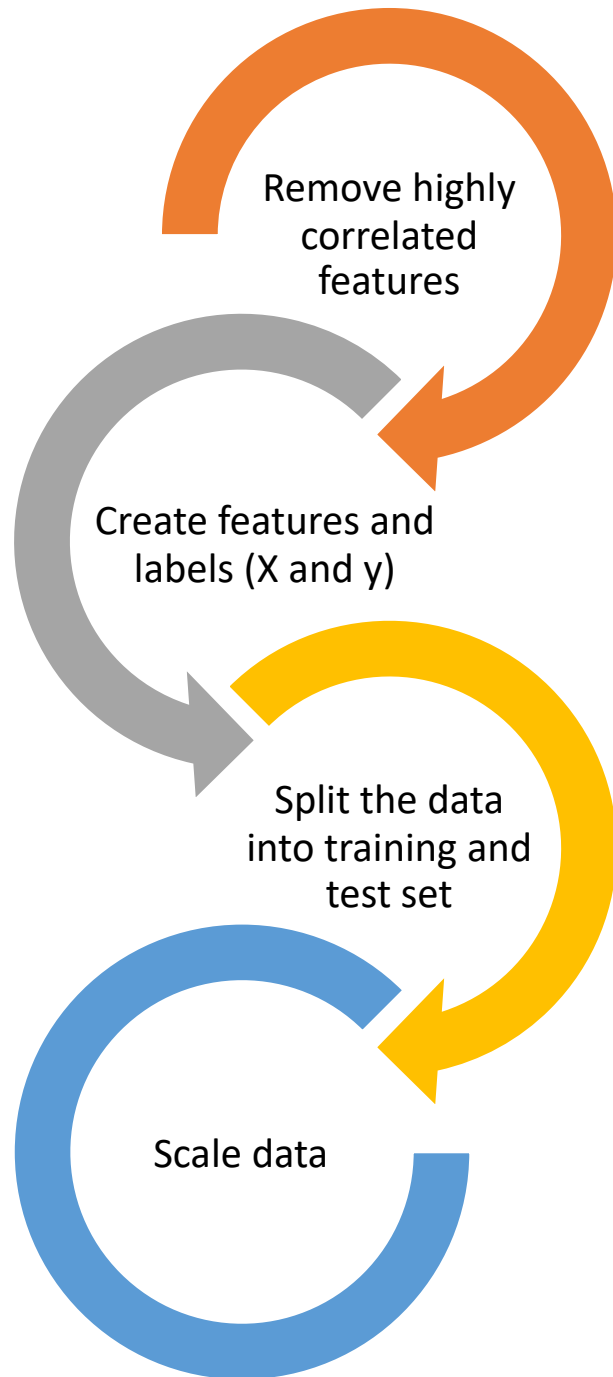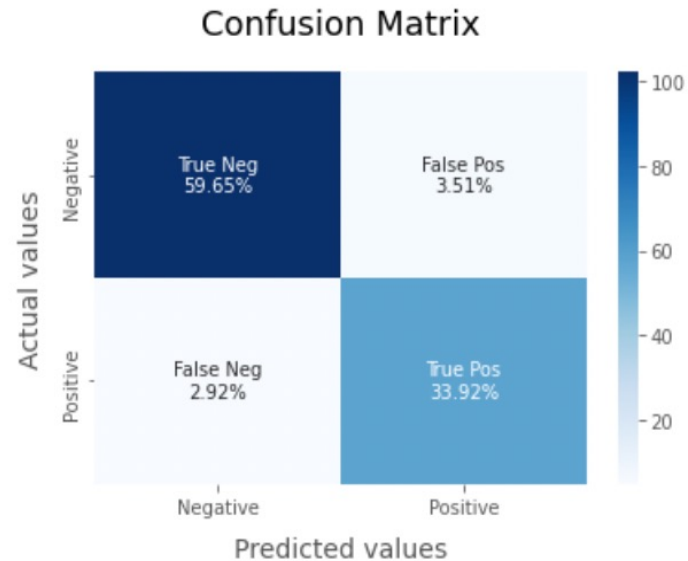
Scale data

You might find these useful:

- from sklearn.model_selection import train_test_split
- X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.30)

- from sklearn.preprocessing import StandardScaler
- scaler = StandardScaler()
- X_train = scaler.fit_transform(X_train)
- X_test = scaler.transform(X_test)

# Create and evaluate models

def model_Evaluate(model):

    # Predict values for Test dataset
    # Print the evaluation metrics for the dataset.
    # Compute and plot the Confusion matrix



Confusion Matrix

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.95 | 0.94 | 0.95 | 108 |
| 1 | 0.91 | 0.92 | 0.91 | 63 |
| accuracy |  |  | 0.94 | 171 |
| macro avg | 0.93 | 0.93 | 0.93 | 171 |
| weighted avg | 0.94 | 0.94 | 0.94 | 171 |

Create def model_Evaluate(model) function

↓

Create, predict, and evaluate the Decision Tree Classifier

↓

Create, predict, and evaluate the k-Nearest-Neighbors

↓

Create, predict, and evaluate the Support Vector Machine

# Choose best model

Sort out the accuracies of each model

Compare them

Choose the best model

| | Model | Score |
|---|---|---|
| 2 | SVC | 0.976608 |
| 0 | DT | 0.935673 |
| 1 | KNN | 0.935673 |

# Important dates

| Task | Week number | Deadline |
|---|---|---|
| Data preprocesssing | Week 1 | 7/13/2022 |
| Analyzing data | Week 1 | 7/15/2022 |
| Choosing features and transforming data | Week 2 | 7/21/2022 |
| Model creation and evaluation | Week 2 | 7/22/2022 |
| Choosing the best model | Week 3 | 7/28/2022 |
| Project presentation | Week 3 | 7/29/2022 |

# Reference

- https://radiopaedia.org/articles/complex-cystic-and-solid-breast-mass
- https://cancer.ca/en/cancer-information/cancer-types/breast/statistics