# CLASSIFICATION:

# BREAST CANCER DETECTION

**REPRESENTED BY: SAOUABEDDINE YOUNES**

**GUIDED BY BY: KHAMJANE AZIZ**

# CONTENT

# INTORDUCTION

Breast cancer is a significant public health concern worldwide, with a high incidence rate among women. Early detection of breast cancer plays a crucial role in improving patient outcomes and survival rates.

Manual detection methods, such as physical examination and palpation, may have limited sensitivity (ability to detect true positive cases) and specificity (ability to avoid false positive cases).

"and that is what we are looking for in this project is to increase sensitivity and low rate of FP cases"

k

(569, 32)

# ARCHITECTURE



DataSet → **Data Processing / Cleaning / Splitting / Standardizing** → **Training**

- Logistic Regression
- SVC
- SGB
- KNN
- Neural Network
- GBC

→ **Testing** → **Benchmarking**

- Accuracy
- Confusion Matrix

→ **The Best Model** → Deploy the model in webPage

# DATA PROCESSING

```
1  # Importing the required libraries
2
3  import pandas as pd
4  import numpy as np
5  import matplotlib.pyplot as plt
6  import seaborn as sns
```

```
1  data = pd.read_csv("breast-cancer.csv")
```

```
1  data.head()
```

| | id int64 | diagnosis object | radius_mean float64 | texture_mean float... | perimeter_mean fl... | area_mean float64 | smoothness_mean f. | c |
|---|---|---|---|---|---|---|---|---|
| 0 | 842302 | M | 17.99 | 10.38 | 122.8 | 1001 | 0.1184 | |
| 1 | 842517 | M | 20.57 | 17.77 | 132.9 | 1326 | 0.08474 | |
| 2 | 84300903 | M | 19.69 | 21.25 | 130 | 1203 | 0.1096 | |
| 3 | 84348301 | M | 11.42 | 20.38 | 77.58 | 386.1 | 0.1425 | |
| 4 | 84358402 | M | 20.29 | 14.34 | 135.1 | 1297 | 0.1003 | |

5 rows, showing [10 ▾] per page          ≪ < Page [1] of 1 > ≫

## Data features

```
1  data.columns
```

```
Index(['id', 'diagnosis', 'radius_mean', 'texture_mean', 'perimeter_mean',
       'area_mean', 'smoothness_mean', 'compactness_mean', 'concavity_mean',
       'concave points_mean', 'symmetry_mean', 'fractal_dimension_mean',
       'radius_se', 'texture_se', 'perimeter_se', 'area_se', 'smoothness_se',
       'compactness_se', 'concavity_se', 'concave points_se', 'symmetry_se',
       'fractal_dimension_se', 'radius_worst', 'texture_worst',
       'perimeter_worst', 'area_worst', 'smoothness_worst',
       'compactness_worst', 'concavity_worst', 'concave points_worst',
       'symmetry_worst', 'fractal_dimension_worst'],
      dtype='object')
```

- We can notice that we have a feature called id, which is not important in our study so we will remove it soon.
- Diagnosis is a categorical variable.

```
1  data.drop('id', axis=1, inplace= True)
2
```
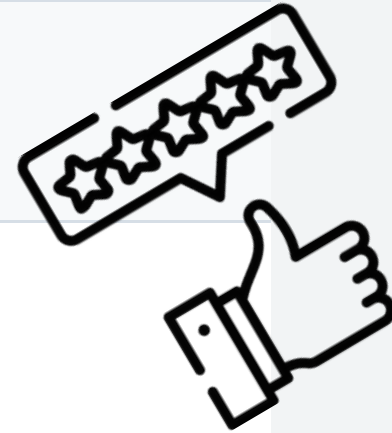
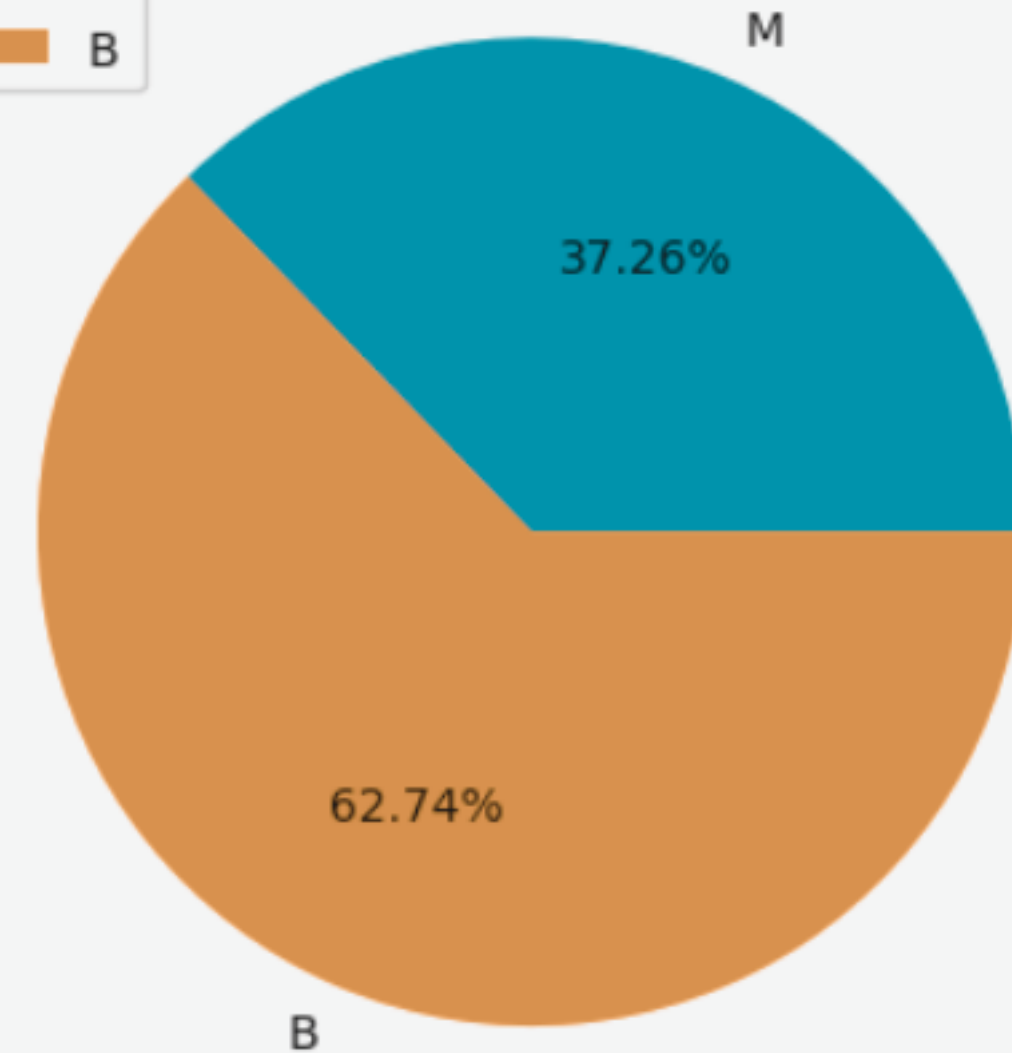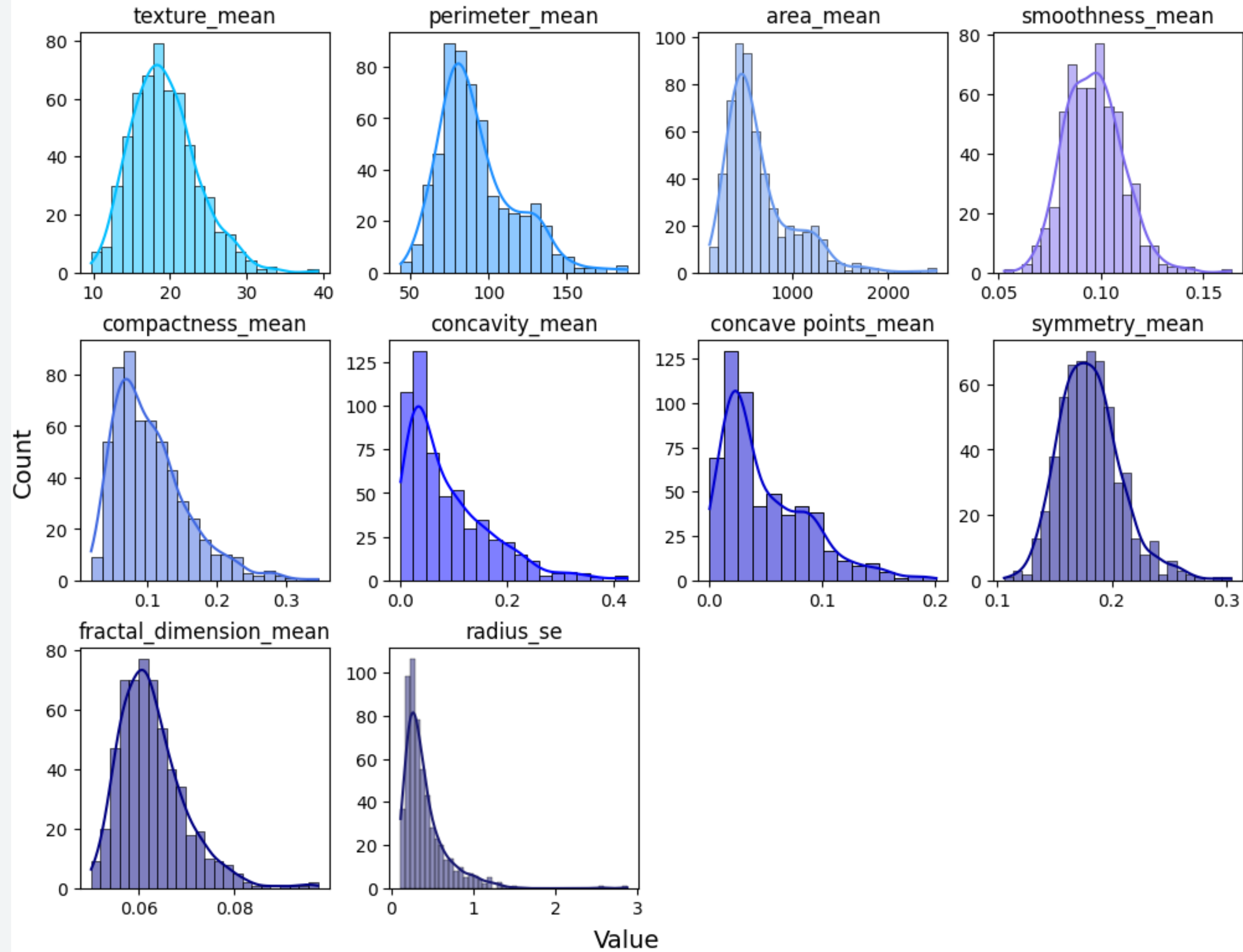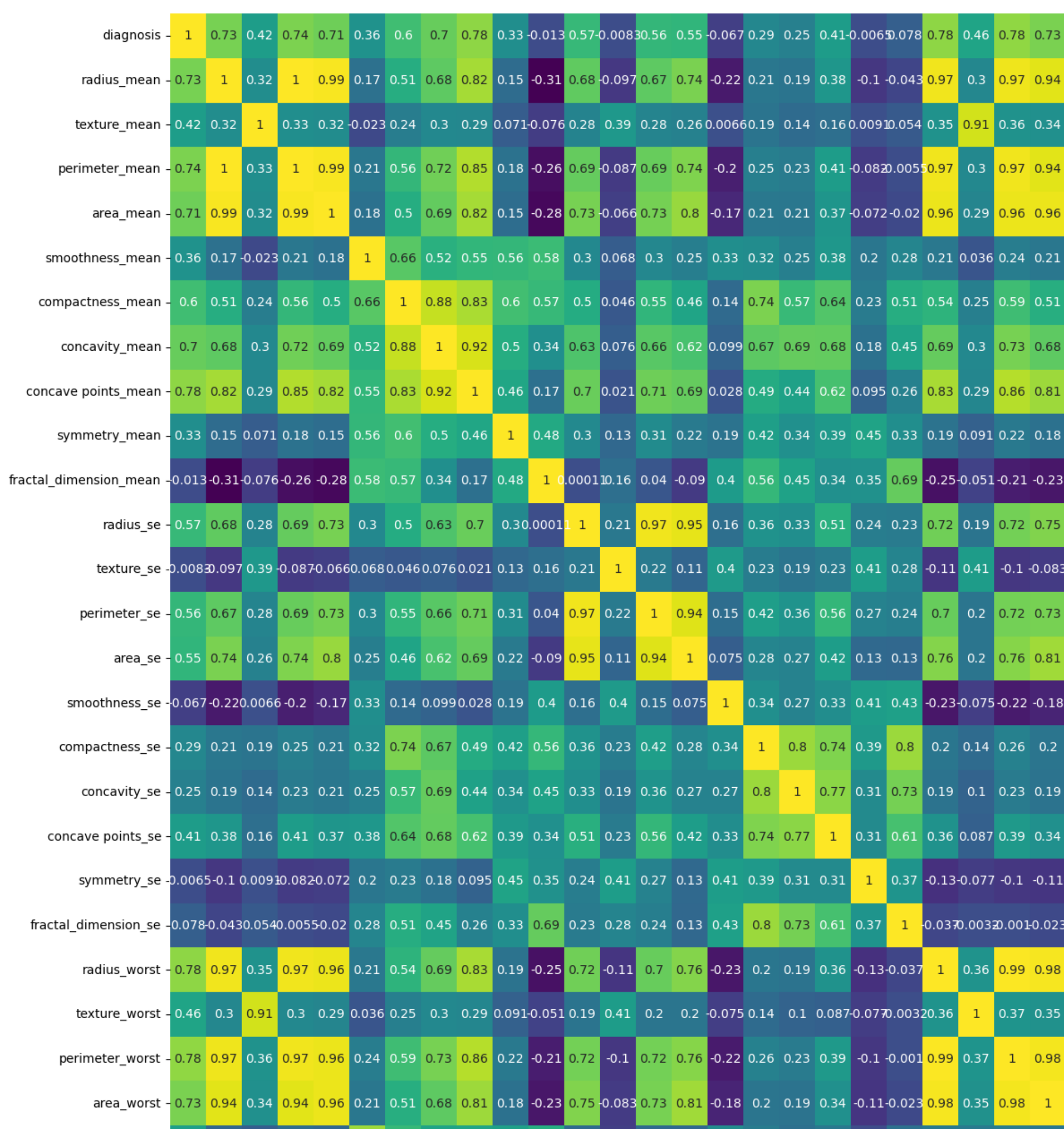# DATA VISUALIZATION

```
1  data.isnull().sum().sum()

0
```

- As we can see, there is no NULL values in our data.

**Pourcentage of Malignant/Begnin in Dataset**

M
B

37.26%

62.74%

B

Distribution of data via Histograms

# CORRELATION MATRIX

| | diagnosis | radius_mean | texture_mean | perimeter_mean | area_mean | smoothness_mean | compactness_mean | concavity_mean | concave points_mean | symmetry_mean | fractal_dimension_mean | radius_se | texture_se | perimeter_se | area_se | smoothness_se | compactness_se | concavity_se | concave points_se | symmetry_se | fractal_dimension_se | radius_worst | texture_worst | perimeter_worst | area_worst |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| diagnosis | 1 | 0.73 | 0.42 | 0.74 | 0.71 | 0.36 | 0.6 | 0.7 | 0.78 | 0.33 | -0.013 | 0.57 | -0.0083 | 0.56 | 0.55 | -0.067 | 0.29 | 0.25 | 0.41 | -0.0065 | -0.078 | 0.78 | 0.46 | 0.78 | 0.73 |
| radius_mean | 0.73 | 1 | 0.32 | 1 | 0.99 | 0.17 | 0.51 | 0.68 | 0.82 | 0.15 | -0.31 | 0.68 | -0.097 | 0.67 | 0.74 | -0.22 | 0.21 | 0.19 | 0.38 | -0.1 | -0.043 | 0.97 | 0.3 | 0.97 | 0.94 |
| texture_mean | 0.42 | 0.32 | 1 | 0.33 | 0.32 | -0.023 | 0.24 | 0.3 | 0.29 | 0.071 | -0.076 | 0.28 | 0.39 | 0.28 | 0.26 | 0.0066 | 0.19 | 0.14 | 0.16 | 0.0091 | 0.054 | 0.35 | 0.91 | 0.36 | 0.34 |
| perimeter_mean | 0.74 | 1 | 0.33 | 1 | 0.99 | 0.21 | 0.56 | 0.72 | 0.85 | 0.18 | -0.26 | 0.69 | -0.087 | 0.69 | 0.74 | -0.2 | 0.25 | 0.23 | 0.41 | -0.082 | -0.0055 | 0.97 | 0.3 | 0.97 | 0.94 |
| area_mean | 0.71 | 0.99 | 0.32 | 0.99 | 1 | 0.18 | 0.5 | 0.69 | 0.82 | 0.15 | -0.28 | 0.73 | -0.066 | 0.73 | 0.8 | -0.17 | 0.21 | 0.21 | 0.37 | -0.072 | -0.02 | 0.96 | 0.29 | 0.96 | 0.96 |
| smoothness_mean | 0.36 | 0.17 | -0.023 | 0.21 | 0.18 | 1 | 0.66 | 0.52 | 0.55 | 0.56 | 0.58 | 0.3 | 0.068 | 0.3 | 0.25 | 0.33 | 0.32 | 0.25 | 0.38 | 0.2 | 0.28 | 0.21 | 0.036 | 0.24 | 0.21 |
| compactness_mean | 0.6 | 0.51 | 0.24 | 0.56 | 0.5 | 0.66 | 1 | 0.88 | 0.83 | 0.6 | 0.57 | 0.5 | 0.046 | 0.55 | 0.46 | 0.14 | 0.74 | 0.57 | 0.64 | 0.23 | 0.51 | 0.54 | 0.25 | 0.59 | 0.51 |
| concavity_mean | 0.7 | 0.68 | 0.3 | 0.72 | 0.69 | 0.52 | 0.88 | 1 | 0.92 | 0.5 | 0.34 | 0.63 | 0.076 | 0.66 | 0.62 | 0.099 | 0.67 | 0.69 | 0.68 | 0.18 | 0.45 | 0.69 | 0.3 | 0.73 | 0.68 |
| concave points_mean | 0.78 | 0.82 | 0.29 | 0.85 | 0.82 | 0.55 | 0.83 | 0.92 | 1 | 0.46 | 0.17 | 0.7 | 0.021 | 0.71 | 0.69 | 0.028 | 0.49 | 0.44 | 0.62 | 0.095 | 0.26 | 0.83 | 0.29 | 0.86 | 0.81 |
| symmetry_mean | 0.33 | 0.15 | 0.071 | 0.18 | 0.15 | 0.56 | 0.6 | 0.5 | 0.46 | 1 | 0.48 | 0.3 | 0.13 | 0.31 | 0.22 | 0.19 | 0.42 | 0.34 | 0.39 | 0.45 | 0.33 | 0.19 | 0.091 | 0.22 | 0.18 |
| fractal_dimension_mean | -0.013 | -0.31 | -0.076 | -0.26 | -0.28 | 0.58 | 0.57 | 0.34 | 0.17 | 0.48 | 1 | 0.00011 | 0.16 | 0.04 | -0.09 | 0.4 | 0.56 | 0.45 | 0.34 | 0.35 | 0.69 | -0.25 | -0.051 | -0.21 | -0.23 |
| radius_se | 0.57 | 0.68 | 0.28 | 0.69 | 0.73 | 0.3 | 0.5 | 0.63 | 0.7 | 0.3 | 0.00011 | 1 | 0.21 | 0.97 | 0.95 | 0.16 | 0.36 | 0.33 | 0.51 | 0.24 | 0.23 | 0.72 | 0.19 | 0.72 | 0.75 |
| texture_se | -0.0083 | -0.097 | 0.39 | -0.087 | -0.066 | 0.068 | 0.046 | 0.076 | 0.021 | 0.13 | 0.16 | 0.21 | 1 | 0.22 | 0.11 | 0.4 | 0.23 | 0.19 | 0.23 | 0.41 | 0.28 | -0.11 | 0.41 | -0.1 | -0.083 |
| perimeter_se | 0.56 | 0.67 | 0.28 | 0.69 | 0.73 | 0.3 | 0.55 | 0.66 | 0.71 | 0.31 | 0.04 | 0.97 | 0.22 | 1 | 0.94 | 0.15 | 0.42 | 0.36 | 0.56 | 0.27 | 0.24 | 0.7 | 0.2 | 0.72 | 0.73 |
| area_se | 0.55 | 0.74 | 0.26 | 0.74 | 0.8 | 0.25 | 0.46 | 0.62 | 0.69 | 0.22 | -0.09 | 0.95 | 0.11 | 0.94 | 1 | 0.075 | 0.28 | 0.27 | 0.42 | 0.13 | 0.13 | 0.76 | 0.2 | 0.76 | 0.81 |
| smoothness_se | -0.067 | -0.22 | 0.0066 | -0.2 | -0.17 | 0.33 | 0.14 | 0.099 | 0.028 | 0.19 | 0.4 | 0.16 | 0.4 | 0.15 | 0.075 | 1 | 0.34 | 0.27 | 0.33 | 0.41 | 0.43 | -0.23 | -0.075 | -0.22 | -0.18 |
| compactness_se | 0.29 | 0.21 | 0.19 | 0.25 | 0.21 | 0.32 | 0.74 | 0.67 | 0.49 | 0.42 | 0.56 | 0.36 | 0.23 | 0.42 | 0.28 | 0.34 | 1 | 0.8 | 0.74 | 0.39 | 0.8 | 0.2 | 0.14 | 0.26 | 0.2 |
| concavity_se | 0.25 | 0.19 | 0.14 | 0.23 | 0.21 | 0.25 | 0.57 | 0.69 | 0.44 | 0.34 | 0.45 | 0.33 | 0.19 | 0.36 | 0.27 | 0.27 | 0.8 | 1 | 0.77 | 0.31 | 0.73 | 0.19 | 0.1 | 0.23 | 0.19 |
| concave points_se | 0.41 | 0.38 | 0.16 | 0.41 | 0.37 | 0.38 | 0.64 | 0.68 | 0.62 | 0.39 | 0.34 | 0.51 | 0.23 | 0.56 | 0.42 | 0.33 | 0.74 | 0.77 | 1 | 0.31 | 0.61 | 0.36 | 0.087 | 0.39 | 0.34 |
| symmetry_se | -0.0065 | -0.1 | 0.0091 | -0.082 | -0.072 | 0.2 | 0.23 | 0.18 | 0.095 | 0.45 | 0.35 | 0.24 | 0.41 | 0.27 | 0.13 | 0.41 | 0.39 | 0.31 | 0.31 | 1 | 0.37 | -0.13 | -0.077 | -0.1 | -0.11 |
| fractal_dimension_se | -0.078 | -0.043 | 0.054 | -0.0055 | -0.02 | 0.28 | 0.51 | 0.45 | 0.26 | 0.33 | 0.69 | 0.23 | 0.28 | 0.24 | 0.13 | 0.43 | 0.8 | 0.73 | 0.61 | 0.37 | 1 | -0.037 | -0.0032 | -0.001 | -0.023 |
| radius_worst | 0.78 | 0.97 | 0.35 | 0.97 | 0.96 | 0.21 | 0.54 | 0.69 | 0.83 | 0.19 | -0.25 | 0.72 | -0.11 | 0.7 | 0.76 | -0.23 | 0.2 | 0.19 | 0.36 | -0.13 | -0.037 | 1 | 0.36 | 0.99 | 0.98 |
| texture_worst | 0.46 | 0.3 | 0.91 | 0.3 | 0.29 | 0.036 | 0.25 | 0.3 | 0.29 | 0.091 | -0.051 | 0.19 | 0.41 | 0.2 | 0.2 | -0.075 | 0.14 | 0.1 | 0.087 | -0.077 | -0.0032 | 0.36 | 1 | 0.37 | 0.35 |
| perimeter_worst | 0.78 | 0.97 | 0.36 | 0.97 | 0.96 | 0.24 | 0.59 | 0.73 | 0.86 | 0.22 | -0.21 | 0.72 | -0.1 | 0.72 | 0.76 | -0.22 | 0.26 | 0.23 | 0.39 | -0.1 | -0.001 | 0.99 | 0.37 | 1 | 0.98 |
| area_worst | 0.73 | 0.94 | 0.34 | 0.94 | 0.96 | 0.21 | 0.51 | 0.68 | 0.81 | 0.18 | -0.23 | 0.75 | -0.083 | 0.73 | 0.81 | -0.18 | 0.2 | 0.19 | 0.34 | -0.11 | -0.023 | 0.98 | 0.35 | 0.98 | 1 |

'A lot of projects preserve only the features with correlation with >0.6 or >0.5 ,but it's not worth it'

# DATA SPLITTING

```python
1  X = data.drop(['diagnosis'], axis=1)
2  y = data['diagnosis']
```

```python
1  from sklearn.model_selection import train_test_split
2  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30, random_state=42)
```

# DATA STANDARDIZATION

$$x' = \frac{x - \bar{x}}{\sigma}$$

```
1  # scaling data
2  from sklearn.preprocessing import StandardScaler
3  s = StandardScaler()
4
5  X_train = s.fit_transform(X_train)
6  X_test  = s.fit_transform(X_test)
```

- Scaling (standardization or normalization) is required when we use any machine learning algorithm that require **gradient calculation.**

# TRAINING & TESTING MODELS

# TRAINING & TESTING MODELS



## Logistic regression

```
1  from sklearn.linear_model import LogisticRegression
2
3  # training the model
4  logModel = LogisticRegression()
5  logModel.fit(X_train, y_train)
6  pred = logModel.predict(X_test)
```

0.9766 accuracy

# TRAINING & TESTING MODELS

## Neural Network

```
1   # importing tensorflow and Keras
2   import tensorflow as tf
3   tf.random.set_seed(3)
4   from tensorflow import keras
```

```
1   # setting up the layers of Neural Network
2
3   NN_model = keras.Sequential([
4                           keras.layers.Flatten(input_shape=(30,)),
5                           keras.layers.Dense(60, activation='relu'),
6                           keras.layers.Dense(2, activation='sigmoid')
7   ])
```

```
1   # compiling the Neural Network
2
3   NN_model.compile(optimizer='adam',
4                    loss='sparse_categorical_crossentropy',
5                    metrics=['accuracy'])
```



0.9942 accuracy

# TRAINING & TESTING MODELS

## Gradient Boosting Classifier

```python
1  from sklearn.ensemble import GradientBoostingClassifier
2  from sklearn.model_selection import GridSearchCV
3
4  gbc = GradientBoostingClassifier()
5
6  parameters = {
7      'loss': ['deviance', 'exponential'],
8      'learning_rate': [0.001, 0.1, 1, 10],
9      'n_estimators': [100, 150, 180, 200]
10 }
11
12 grid_search_gbc = GridSearchCV(gbc, parameters, cv = 5, n_jobs = -1, verbose = 1)
13 grid_search_gbc.fit(X_train, y_train)
```

```python
1  # best parameters
2  grid_search_gbc.best_params_
```

```
{'learning_rate': 1, 'loss': 'exponential', 'n_estimators': 100}
```

```python
1  gbc = GradientBoostingClassifier(learning_rate = 0.1, loss = 'exponential', n_estimators = 180)
2  gbc.fit(X_train, y_train)
```

```
                    GradientBoostingClassifier
GradientBoostingClassifier(loss='exponential', n_estimators=180)
```



0.97 accuracy

# TRAINING & TESTING MODELS

## Stochastic Gradient Boosting (SGB)

```
1  sgbc = GradientBoostingClassifier(max_depth=4, subsample=0.9, max_features=0.75, n_estimators=200, random_state=0)
2  sgbc.fit(X_train, y_train)
```

## K Neighbors Classifier (KNN)

```
1  from sklearn.neighbors import KNeighborsClassifier
2
3  knn = KNeighborsClassifier()
4  knn.fit(X_train, y_train)
```

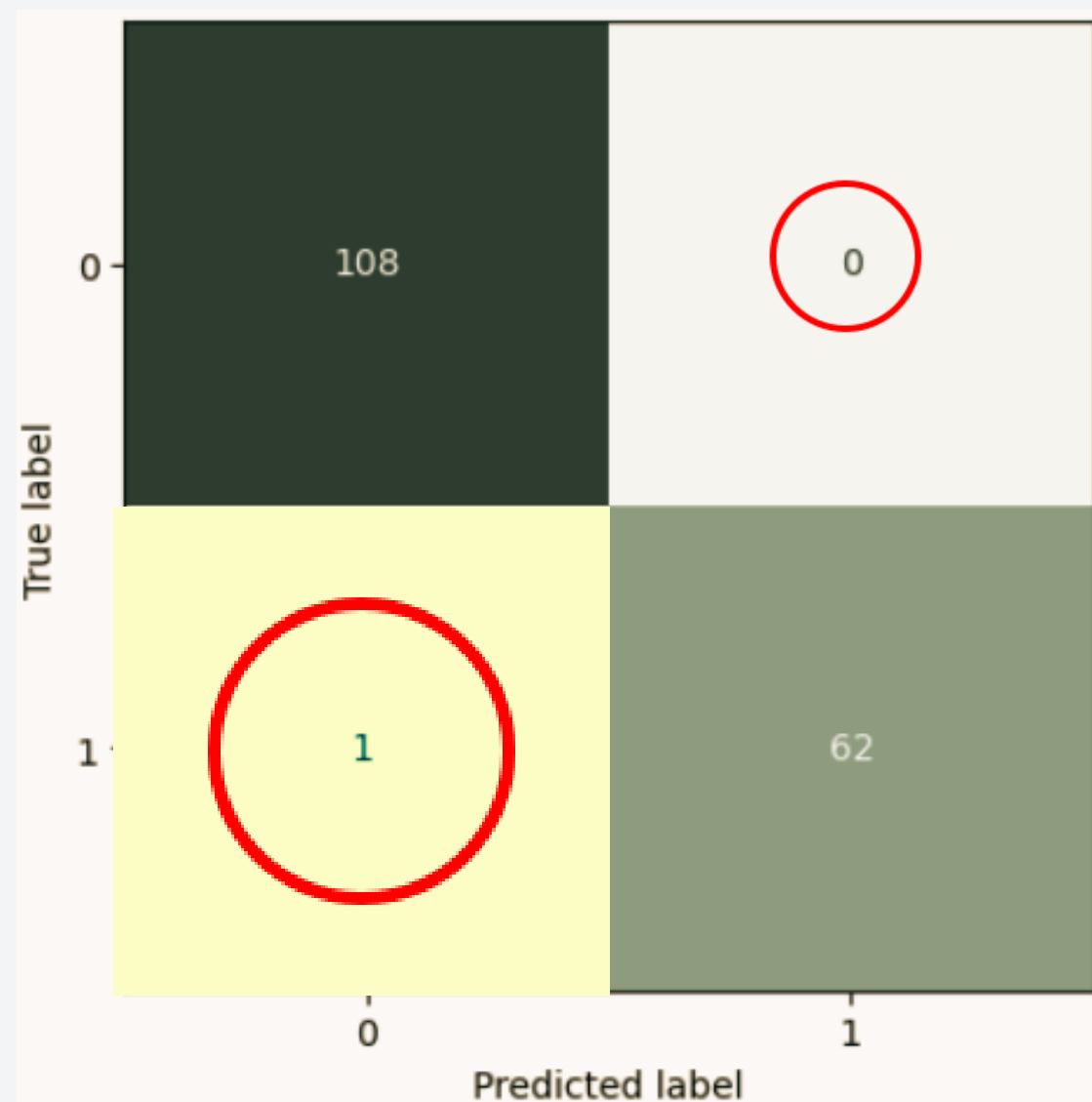the same accuracy & and confusion Matrix as GBC

# TRAINING & TESTING MODELS

**Support Vector Classifier (SVC)**

`': 20, 'gamma': 0.001}`



**0.988 accuracy**

# Neural Network



**0.9942 accuracy**

# Model Deployment using Flask FW

**Inputed Data:**

| 412.3 | 0.1001 | 0.07348 | 0 | 0 | 0.2458 | 0.06592 |

**Class:**
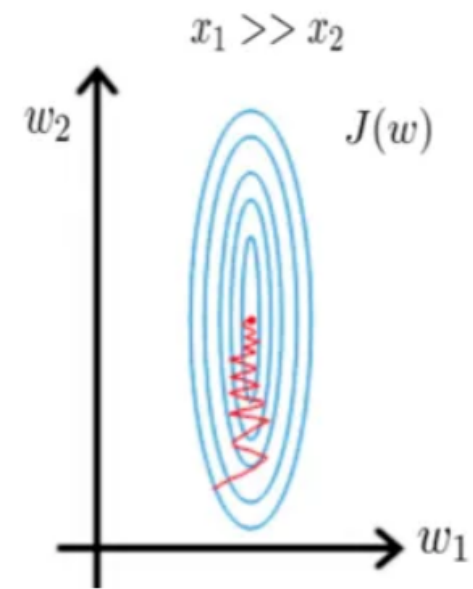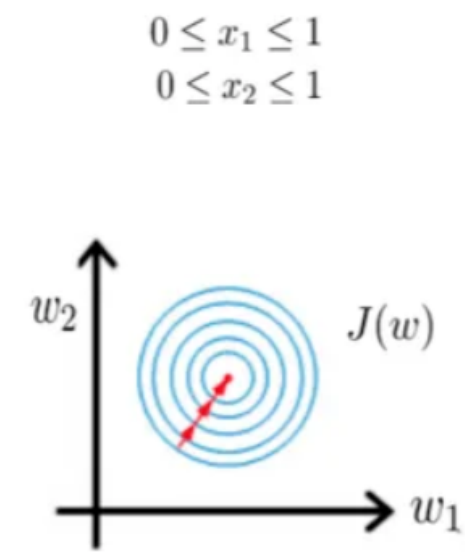
**The tumor is Benin**

Predict

# CHALLENGES

# CONCLUSION

**FUTURE WORK !**

# CNN implementation from scratch

1. Mean Centering: Subtracting the mean from each feature ensures that the mean of the feature becomes zero. This is important because it removes any bias from the features. When features have different means, they can bias the learning algorithm, especially in algorithms that are sensitive to feature scales, like gradient descent-based algorithms.
2. Variance Scaling: Scaling the features to unit variance means that the variance of each feature becomes 1. This is important because it gives equal importance to all features, preventing features with larger variances from dominating the learning process. It also helps algorithms converge faster, especially for algorithms that are sensitive to feature scales.

Gradient descent without scaling

Gradient descent after scaling variables

$x_1 \gg x_2$

$w_2$    $J(w)$

$0 \leq x_1 \leq 1$
$0 \leq x_2 \leq 1$

$w_2$    $J(w)$

$w_1$

$w_1$

Reference



With Scaling

Without Scaling