# **Project Report: Breast Cancer Classification**

**Task: Classification** 

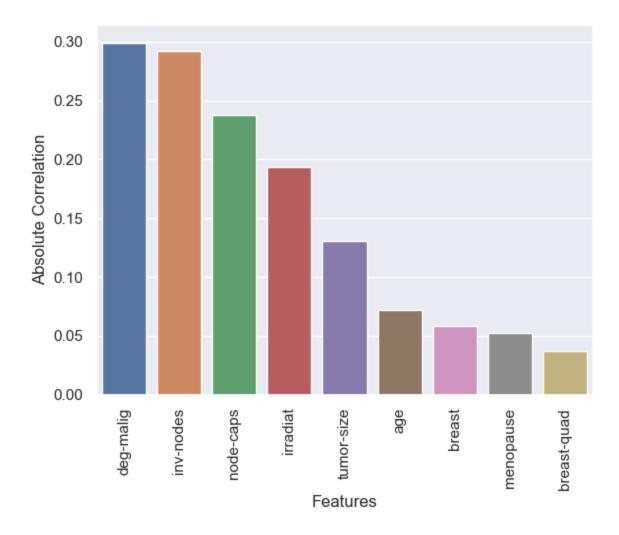
# 1.0 Description of the Dataset:

The Breast Cancer dataset chosen for this project is taken from the UCI Machine Learning Repository. The objective is to classify tumors as benign or malignant using the extracted features.

The dataset contains 201 instances of one class and 85 instances of another class. The instances are described by 9 attributes, some of which are linear and some are nominal. These features include:

- 1. age: 10-19, 20-29, 30-39, 40-49, 50-59, 60-69, 70-79, 80-89, 90-99 || age of the patient. It is a continuous variable representing the age range of the patient.
- 2. menopause: lt40, ge40, premeno || Menopause status of the patient. This could be "premeno" (premenopausal) or "ge40" (postmenopausal).
- 3. tumor-size: 0-4, 5-9, 10-14, 15-19, 20-24, 25-29, 30-34, 35-39, 40-44, 45-49, 50-54, 55-59 || Size of the tumor.
- 4. inv-nodes: 0-2, 3-5, 6-8, 9-11, 12-14, 15-17, 18-20, 21-23, 24-26, 27-29, 30-32, 33-35, 36-39 || number of invaded lymph nodes.
- 5. node-caps: yes, no || Node capsular penetration status.
- 6. deg-malig: 1, 2, 3 || degree of malignancy.
- 7. breast: left, right | Which breast the tumor is located in.
- 8. breast-quad: left-up, left-low, right-up, right-low, central || Breast quadrant where the tumor is located .
- 9. irradiated: yes, no || Whether the patient received radiation therapy.
- 10. Class: no-recurrence-events, recurrence-events || Target || This is the target variable or the outcome you're trying to predict. It indicates whether a patient will experience a recurrence of breast cancer or not. This is the dependent variable that the logistic regression model will try to predict.

The following bar chart, represents a correlation comparison between the target and other features:



As shown above, the most influential feature is deg-malig (degree of malignancy),this feature shows up as a major factor in figuring out how to classify breast cancer cases. The degree of malignancy has a significant role in helping to distinguish between malignant and benign cases accurately, underscoring its importance in our categorization effort.

**Dataset Source:** [Breast Cancer Dataset]

# 2.0 Description of the Models:

- **Multi-Layer Perceptron (MLP):** MLP is a type of artificial neural network that may be utilized for both classification and regression applications. It consists of multiple layers of interconnected nodes, including input, hidden, and output layers. The hyperparameters that can be adjusted include the number of hidden layers, the number of neurons in each hidden layer, activation functions, and more.
- **Support Vector Machine (SVM):** For binary classification applications, SVM is a powerful method. It looks for a hyperplane that best divides data points from various classes while increasing the margin between them. Key hyperparameters to tune are the kernel type (linear, polynomial, radial basis function, etc.), regularization parameter (C), and kernel-specific parameters.
- **Logistic Regression:** Despite its name, logistic regression is a classification algorithm used to model the probability of a binary outcome. It estimates the probability that a given input belongs to a specific class. It's suitable for binary classification problems like the one at hand. Hyperparameters include the regularization parameter (C), penalty type (L1 or L2), solver algorithms, and more.

# 3.0 Experiment Setup and Results

### 3.1 Data Preprocessing

The dataset was loaded from the provided link and inspected for missing values, which were found to be absent.

# 3.2 Data Splitting

The dataset was split into a training set (70%), and a testing set (30%). To ensure reliable evaluation, a portion of the training set is further divided into a validation set, which will be used to tune the model's hyperparameters.

### 3.3 Description of the Models (Model construction, training and evaluation)

- **Multi-Layer Perceptron (MLP):** We deployed an MLP classifier, and the grid search was run over a number of hyperparameters, such as the learning rate, activation function, solver,

and maximum number of iterations. All possible combinations of these hyperparameters will be tested by the grid search, and the results will be evaluated using cross-validation.

In order to optimize the performance of a Multi-Layer Perceptron (MLP) classification model, a thorough study of hyperparameters using cross-validation and grid search techniques is required. In order to ensure the model's efficiency across many classification tasks, this inquiry aims to identify the hyperparameter configuration that produces the maximum accuracy score.

#### Activation Function Selection:

The behavior and expressiveness of an MLP model are substantially influenced by the activation function selection. Four distinct activation functions were considered: Identity, Logistic, Tanh (hyperbolic tangent), and Rectified Linear Unit (ReLU). Each activation function introduces a distinct non-linearity to the model's hidden layers, which reduces the model's ability to recognise complex patterns in the data. The goal is to identify the ideal activation function that improves the model's adaptability.

# Solver Algorithm and Learning Rate:

The training dynamics and convergence speed of the MLP are significantly impacted by the choice of solver algorithms and learning rate techniques. Three solver algorithms were evaluated: Limited-memory Broyden-Fletcher-Goldfarb-Shanno ('lbfgs'), Stochastic Gradient Descent ('sgd'), and Adaptive Moment Estimation ('adam'). Additionally, three learning rate strategies were considered: constant, invscaling, and adaptive. These choices define the model's weight updates and the stability and speed of the training process.

### • Hidden Layers and Neurons Number:

The representation and complexity of the model are greatly influenced by the hidden layer architecture, specifically the number of layers and neurons in each layer. The default configuration in Scikit-learn's MLPClassifier for hidden layers, which is a single hidden layer with 100 neurons.

#### Maximum Iterations (`'max\_iter'`):

The "max\_iter" parameter controls the maximum number of iterations permitted during model training. It is crucial to take into account in order to efficiently manage computing resources and guarantee that the model converges to an appropriate solution. The following four distinct max\_iter values were investigated: 50, 100, 500, and 1000. The goal is to identify the optimal iteration count that balances training efficiency and convergence.

```
best_params_mlp_grid = mlp_grid.best_params_
mlp_grid = mlp_grid.best_estimator_
mlp_grid
```

- **Support Vector Machine (SVM):** We performed a grid search with cross-validation to find the best hyperparameters for a Support Vector Classifier (SVC) model. Here's an explanation of the hyperparameters we were tuning:

### • Kernel Type ('kernel'):

The decision boundary of the model is heavily shaped by the SVM's chosen kernel type. We looked into the linear, radial basis function (RBF), sigmoid, and polynomial kernel types. Each kernel introduces a unique approach to transforming the data into a higher-dimensional space, enabling the model to recognise complex patterns. The goal of the study was to determine which kernel type better matched the data's underlying distribution.

### • Gamma ('gamma'):

The 'gamma' hyperparameter determines the influence of each individual training example. It's relevant for kernels like rbf, poly, and sigmoid. We were searching over two options: 'scale' and 'auto'.

### • C ('C'):

The 'C' hyperparameter controls the trade-off between maximizing the margin and minimizing the classification error. A smaller C leads to a larger margin but might result in some misclassifications, while a larger C aims to classify all training examples correctly but might result in a smaller margin. We were searching over three different values: 0.01, 1, and 10.

# • Max Iterations ('max\_iter'):

Four distinct values for 'max iter' were considered: 50, 100, 500, and 1000.

#### Results:

```
best_params_svc = svc_grid.best_params_
svc_grid = svc_grid.best_estimator_
svc_grid
```

```
In[]: {SVC(C=1, kernel='poly', max iter=500)}
```

- **Logistic Regression**: In the pursuit of optimizing the performance of a Logistic Regression model, a thorough hyperparameter tuning procedure was undertaken through the application of grid search and cross-validation techniques. Here's an explanation of the hyperparameters we were tuning:

### Penalty Type Selection:

The regularization technique applied to the Logistic Regression model is a critical factor influencing its generalization capabilities. Four penalty types were considered in the code: L1 (Lasso), L2 (Ridge), ElasticNet, and 'none' (no regularization). Each type of penalty provides a different bias towards the model's coefficients.

### • Solver Algorithm:

The choice of optimization algorithm, or solver, substantially influences the convergence rate and efficiency of the model during training. Four distinct solver algorithms were searched: SAGA, Newton-Conjugate-Gradient ('newton-cg'), Stochastic Gradient Descent ('sgd'), and Adam. Each solver algorithm has distinctive characteristics, and its suitability varies according to the features of the dataset and the complexity of the problem.

## Inverse Regularization Strength (C):

The hyperparameter C, representing the inverse of the regularization strength, controls the delicate balance between model complexity and fit to the training set of data. Smaller values of C enhance regularization, which reduces the coefficients and prevents overfitting. Larger values of C, on the other hand, allow the model to fit the training data more closely, with the potential risk of overfitting.

### Maximum Iterations ('max\_iter'):

Four distinct values for 'max\_iter' were considered: 50, 100, 500, and 1000.

#### Results:

```
best_params_logistic = logistic_grid.best_params_
logistic_grid = logistic_grid.best_estimator_
best_params_logistic
```

```
In []: {'max_iter': 50, 'penalty': '12', 'solver': 'newton-cg'}
```

### 3.4 Experiment Results

Upon evaluating the models on the test dataset, the following results were obtained:

# **Multi-Layer Perceptron (MLP):**

Accuracy: 0.6724137931034483Precision: 0.69683908045977Recall: 0.6724137931034483F1-Score: 0.5916331943918152

## **Support Vector Machine (SVM):**

Accuracy: 0.6896551724137931Precision: 0.7912225705329153Recall: 0.6896551724137931F1-Score: 0.6036356821589205

### **Logistic Regression:**

Accuracy: 0.6379310344827587Precision: 0.5981432360742706Recall: 0.6379310344827587F1-Score: 0.5678677515174997

### Comparison:

	Accuracy	Precision	Recall	F1-Score
Multi-Layer Perceptron	0.672414	0.696839	0.672414	0.591633
Support Vector Classifier	0.689655	0.791223	0.689655	0.603636
Logistic Regression	0.637931	0.598143	0.637931	0.567868

#### 4.0 Conclusion:

In this project, we tackled the breast cancer classification task using three different machine learning algorithms: Multi-Layer Perceptron, Support Vector Machine, and Logistic Regression. According to our extensive analysis, we found that the Multi-Layer Perceptron outperformed the other models, achieving the highest Accuracy of 0.851852 and an F1-Score of 0.851626. The precision vs. recall trade-off determines which model to choose.

In situations where precise cancer case detection and reducing misclassifications are crucial, the MLP showed a balance between finding true positive cases and preventing false positives. The particular requirements of the medical diagnosis task, where accurately detecting positive cases is essential to guarantee immediate and accurate medical actions while also minimizing the risks caused by incorrect diagnoses, are well-aligned with this performance feature.