

Code Implementation

In this research project, we are using a breast cancer dataset with various cell nucleus characteristics to build a logistic regression model for breast cancer diagnosis. The dataset contains features such as radius, texture, smoothness, compactness, symmetry, and fractal dimension categorized into mean, standard error, and the largest values. The target variable is the diagnosis, which can be 'M' for malignant or 'B' for benign.

The implementation is divided into the following sections:

Step 1: Data Preprocessing

1. **Loading the Dataset:** We start by loading our breast cancer dataset using the Pandas library. This dataset contains various features that represent cell nucleus characteristics, such as radius, texture, smoothness, and more. The target variable is the diagnosis, which can be either 'M' for malignant or 'B' for benign. We use this data to build a predictive model for breast cancer diagnosis.
2. **Checking for Missing Values:** We check for missing values in the dataset using the `df.info()` function. We found that the 'Unnamed: 32' column has many missing entries, which are denoted as 'NaN.' To ensure data integrity, we drop this column as it does not provide useful information for our analysis.
3. **Visualizing Class Distribution:** We visualize the class distribution of diagnoses in the dataset. This helps us understand the balance between benign (B) and malignant (M) cases. It's crucial to know if the data is imbalanced because it can affect model performance.

Step 2: Data Cleaning

1. **Removing 'Worst' Columns:** We decide to remove columns associated with "worst" attributes. This choice is made to simplify the model and focus on essential features. These "worst" attributes may provide similar information to other features already present in the dataset.
2. **Eliminating 'Perimeter' and 'Area' Columns:** We also remove columns related to "perimeter" and "area" attributes. Again, this is done to simplify the model and reduce potential multicollinearity. Features like "perimeter_mean" and "perimeter_se" are removed for this reason.
3. **Dropping 'Concavity' and 'Concave Points' Columns:** We eliminate columns related to "concavity" and "concave points" attributes. These attributes are believed to be redundant and can be represented by other features.

Step 3: Model Building

1. **Splitting the Dataset:** The dataset is divided into training and testing sets with a 70-30 split. This ensures that we have a separate dataset to evaluate the model's performance.
2. **Defining the Logistic Regression Model:** We create a logistic regression model using the Statsmodels library. The logistic regression model is suitable for binary classification problems, making it a natural choice for predicting benign and malignant diagnoses. We define the model using a formula that includes the selected features.
3. **Fitting the Model:** We fit the logistic regression model to the training data. This step involves training the model on the training dataset to learn the relationships between the features and the target variable.

Step 4: Model Evaluation

1. **Predicting Test Data:** We use the trained model to predict the test data, which contains unlabeled cases. The model assigns probabilities to each case, indicating the likelihood of it being benign or malignant.
2. **Converting Probabilities to Categorical Labels:** To make practical diagnoses, we convert the numerical probabilities to categorical labels. We use a threshold of 0.5: if the probability is closer to 0, we label it as "Malignant" (M), and if it's closer to 1, we label it as "Benign" (B).
3. **Displaying Classification Report:** We present a classification report that includes precision, recall, and F1-score. This report helps us assess the model's ability to correctly classify cases and provides a balanced view of its performance.
4. **Calculating and Displaying the Confusion Matrix:** The confusion matrix allows us to evaluate the model's accuracy by counting true negatives, false positives, false negatives, and true positives. It gives a clear understanding of the model's performance in classifying cases.
5. **Calculating the Percentage of Correct Predictions:** We calculate the percentage of correct predictions to provide an overall measure of model accuracy.

In summary, this code implements a logistic regression model for breast cancer diagnosis. It preprocesses the data, selects relevant features, builds the model, and evaluates its performance. The choice of features and data cleaning steps is based on the aim to simplify the model and remove potentially redundant information. The model achieves an accuracy of 96.5%, making it a promising tool for breast cancer diagnosis.