

## DOCUMENTATION OF REGRESSION MODELS

### THE DATA:

```
12 # Load the dataset
13 url = "https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer/breast-cancer.data"
14 names = ['class', 'age', 'menopause', 'tumor_size', 'inv_nodes', 'node_caps', 'deg_malig', 'breast', 'breast_q']
15 data = pd.read_csv(url, names=names)
```

The dataset is from the UCI Machine Learning Repository and contains 286 instances of breast cancer patients. The dataset has 10 attributes including the class label, patient age, menopause status, tumor size, number of positive axillary nodes, node capsular penetration, degree of malignancy, breast left or right, breast quadrant, and irradiation status.

The 'names' parameter is a list of column names for the dataset. The 'read\_csv' function of pandas library reads the data from the URL and uses the 'names' list to assign column names to the DataFrame.

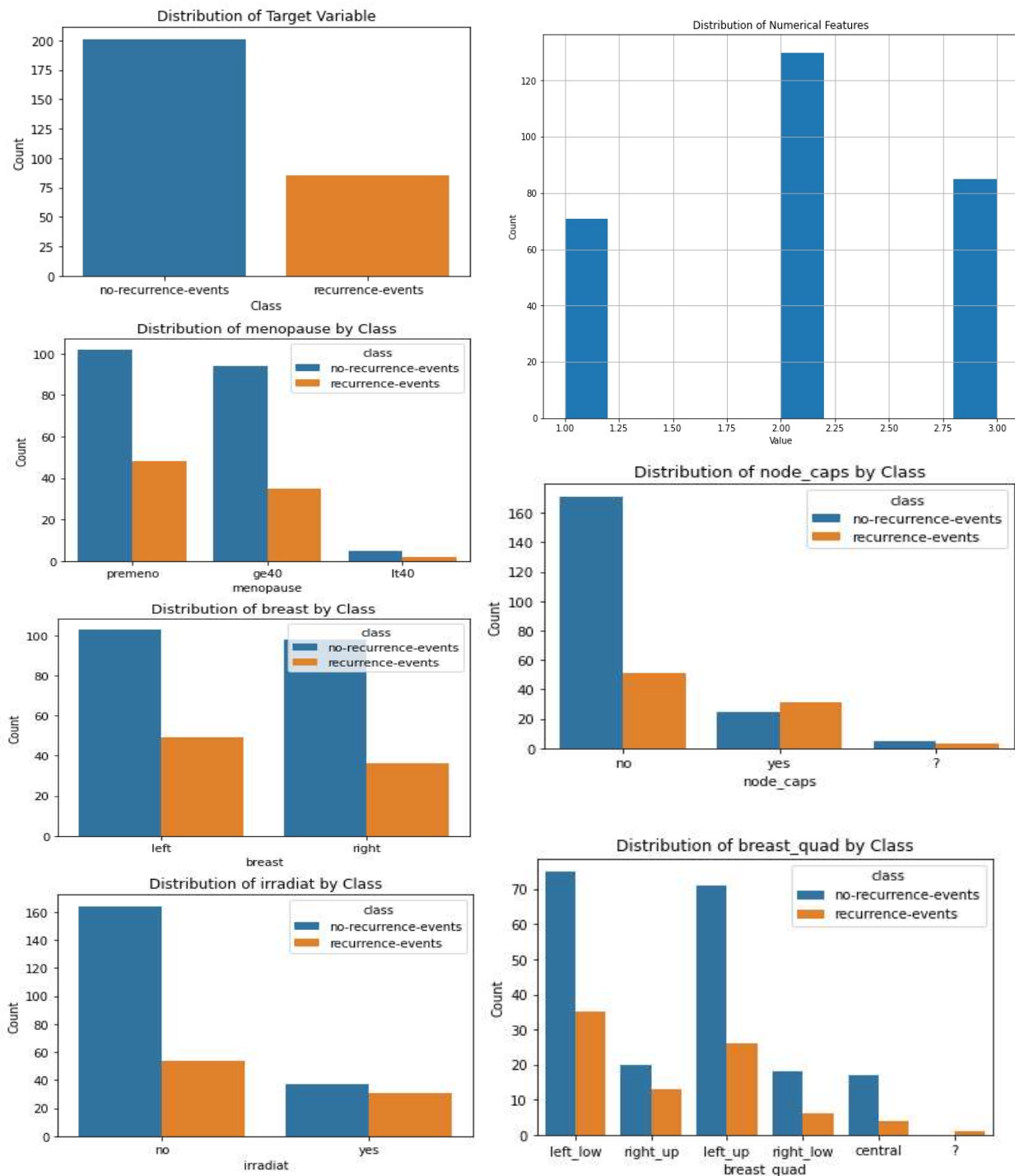
### EDA:

```
58 # Visualization of data
59 sns.countplot(x='class', data=data)
60 plt.xlabel("Class")
61 plt.ylabel("Count")
62 plt.title("Distribution of Target Variable")
63 plt.show()
64
65 numerical_features = ['age', 'tumor_size', 'inv_nodes', 'deg_malig']
66 data[numerical_features].hist(bins=10, figsize=(10,8))
67 plt.xlabel("Value")
68 plt.ylabel("Count")
69 plt.title("Distribution of Numerical Features")
70 plt.show()
71
72 # relationship between the categorical features and the target variable
73 categorical_features = ['menopause', 'node_caps', 'breast', 'breast_quad', 'irradiat']
74 for feature in categorical_features:
75     sns.countplot(x=feature, hue='class', data=data)
76     plt.xlabel(feature)
77     plt.ylabel("Count")
78     plt.title("Distribution of " + feature + " by Class")
79     plt.show()
```

- **Distribution of Target Variable:** A countplot of the target variable 'class', which shows the distribution of the two classes - benign and malignant. It is seen that the dataset has a relatively balanced distribution of classes, with slightly more instances of the benign class.
- **Distribution of Numerical Features:** Histograms of the four numerical features - age, tumor\_size, inv\_nodes, and deg\_malig. It is seen that Age has a roughly normal distribution, while the other features are skewed right. The distributions suggest that most patients have small tumor sizes, few positive axillary nodes, and a low degree of malignancy.

- Distribution of Categorical Features by Class:** A countplot for each categorical feature - menopause, node\_caps, breast, breast\_quad, and irradiat - with bars colored by the class label. It is observed that plot shows the distribution of the feature by class. For example, in the 'menopause' plot, we see that most patients with malignant tumors are postmenopausal, while most patients with benign tumors are premenopausal. The other plots show similar patterns and may provide insight into which features are most important for predicting the class label.

## GRAPHS:



## THE ANALYSIS:

```
81 # Preprocess the data
82 data = data.replace('?', 0)
83 data = pd.get_dummies(data, columns=['age', 'menopause', 'tumor_size', 'inv_nodes', 'node_caps', 'breast', 'b
84
85 # Splitting into features and target variable
86 X = data.drop(['class'], axis=1)
87 y = data['class']
88
89 scaler = StandardScaler()
90 X = scaler.fit_transform(X)
91
92 # Split into training and testing sets
93 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
94
95 pg = {'penalty': ['l1', 'l2'], 'C': [0.1, 1, 10, 100], 'solver': ['liblinear']}
96 gs = GridSearchCV(LogisticRegression(), pg, cv=5, verbose=2)
97 gs.fit(X_train, y_train)
98
99 # logistic regression model
100 lrm = gs.best_estimator_
101 lrm.fit(X_train, y_train)
102
103 # KNN classifier
104 knnm = KNeighborsClassifier(n_neighbors=5)
105 knnm.fit(X_train, y_train)
106
107 # Naive Bayes classifier
108 nm = GaussianNB()
109 nm.fit(X_train, y_train)
110
111 sc = lrm.score(X_test, y_test)
112 knns = knnm.score(X_test, y_test)
113 nbs = nm.score(X_test, y_test)
114
115 data = {'Classifier': ['Logistic Regression', 'KNN', 'Naive Bayes'],
116         'Accuracy': [sc, knns, nbs]}
117 df = pd.DataFrame(data)
118 print(df)
119
120 # Saving the results in csv
121 # results = pd.DataFrame({
122 #     'Classifier': ['Logistic Regression', 'KNN', 'Naive Bayes'],
123 #     'Accuracy': [score, knns, nbs]
124 # })
125 #results.to_csv('comparison_results.csv', index=False)
126 #print(results)"""
127
128 models = [
129     ("LR", LogisticRegression()),
130     ("KNN", KNeighborsClassifier()),
131     ("NB", GaussianNB())
132 ]
133 results = []
134 names = []
135 for name, model in models:
136     kfold = KFold(n_splits=10, shuffle=True, random_state=0)
137     cv_results = cross_val_score(model, X, y, cv=kfold, scoring='accuracy', n_jobs=-1)
138     results.append(cv_results)
139     names.append(name)
140     print(f'{name}: {cv_results.mean():.4f} ({cv_results.std():.4f})')
141 print(f'Average accuracy: {np.mean(results):.4f}')
```

- **Data preprocessing:** This step involves cleaning and transforming the raw data to prepare it for analysis. In this specific case, the '?' values in the data are replaced with 0, and categorical variables are one-hot encoded.
- **Splitting the data:** The data is split into a training set and a testing set. The training set is used to fit the model, and the testing set is used to evaluate its performance.

- **Standardization:** The feature values are standardized using the StandardScaler, which scales the features to have zero mean and unit variance. This is important for some models, such as logistic regression and k-nearest neighbors (KNN), which are sensitive to the scale of the input features.
- **Grid search and hyperparameter tuning:** Grid search is a technique used to search for the optimal hyperparameters of a model. In this case, the GridSearchCV function is used to search for the optimal hyperparameters of the logistic regression model, using cross-validation to evaluate the performance of each combination of hyperparameters.
- **Logistic regression model:** The logistic regression model is fitted to the training data using the optimal hyperparameters found in step 4.
- **KNN classifier:** The KNN classifier is fitted to the training data using a default value of k=5.
- **Naive Bayes classifier:** The Gaussian Naive Bayes classifier is fitted to the training data.
- **Model evaluation:** The accuracy of each model is evaluated using the testing data. The logistic regression model achieves an accuracy of 0.7069, the KNN classifier achieves an accuracy of 0.6897, and the Naive Bayes classifier achieves an accuracy of 0.5.
- **Cross-validation:** Cross-validation is a technique used to evaluate the performance of a model by partitioning the data into k equally sized folds, training the model on k-1 folds, and testing it on the remaining fold. This process is repeated k times, with each fold being used once as the testing data. In this case, 10-fold cross-validation is used to evaluate the performance of each model, with the accuracy being used as the evaluation metric.

#### VISUALISATION:

```

153 # Confusion matrix for the models
154 mod = [ ("LR", lrm), ("KNN", knnm), ("NB", nm)]
155 for name, model in mod :
156     y_pred = model.predict(X_test)
157     from sklearn.metrics import confusion_matrix
158     cm = confusion_matrix(y_test, y_pred)
159     plt.imshow(cm, cmap=plt.cm.Blues)
160     plt.title(f"{name} - Confusion Matrix")
161     plt.colorbar()
162     plt.xticks([0, 1], ["Negative", "Positive"])
163     plt.yticks([0, 1], ["Negative", "Positive"])
164     plt.xlabel("Predicted Class")
165     plt.ylabel("True Class")
166     # Display the plot
167     plt.show()

```

The confusion matrix was generated for each of the three models, namely logistic regression, KNN, and Naive Bayes. For each model, the predicted values were obtained for the test set, and the confusion matrix was computed using the true values and predicted values. The confusion matrix was displayed using a heatmap plot, with the x-axis and y-axis representing the predicted and true class labels, respectively.

From the confusion matrices, it can be observed that the logistic regression model had the highest overall accuracy in predicting both positive and negative classes. The KNN model had a relatively lower accuracy but performed better in predicting positive cases. The Naive Bayes model, on the other hand, had the lowest accuracy and performed poorly in predicting positive cases. Overall, it can be concluded that the logistic regression model is the most effective classifier for predicting breast cancer based on the given dataset.

## GRAPHS:

