

Practical Activity 4.4.1: Logistic Regression

Logistic regression is a linear model for classification. Logistic regression is also known as logit regression, maximum-entropy classification (MaxEnt) or the log-linear classifier. More technical details can be found at https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression (https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression).

In this practical, we will use the Breast Cancer Wisconsin (Diagnostic) Dataset. The dataset is downloaded from the UCI Machine Learning repository [https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic)) ([https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic))).

The dataset consists of 569 samples. Each sample has 30 features. The features are computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe characteristics of the cell nuclei present in the image.

The task is to classify each sample into two classes Benign (1) or Malignant (0).

In [84]:

```
1 import pandas as pd
2 import numpy as np
3 from sklearn import preprocessing
4 import matplotlib.pyplot as plt
5 plt.rc("font", size=14)
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.model_selection import train_test_split
8 import seaborn as sns
9 from sklearn.feature_selection import RFE
10 from sklearn.linear_model import LogisticRegression
11 sns.set(style="white")
12 sns.set(style="whitegrid", color_codes=True)
```

In [85]:

```
1 # Load the dataset
2 data = pd.read_csv('wdbc.csv', header=0)
```

In [86]:

```
1 # inspect the data, drop rows with missing data
2 data = data.dropna()
3 print(data.shape)
4 data.head()
```

(569, 31)

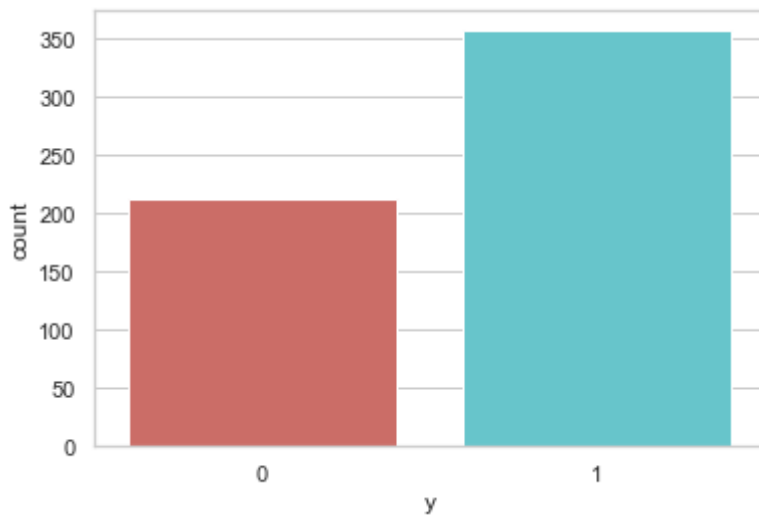
Out[86]:

	Feature1	Feature2	Feature3	Feature4	Feature5	Feature6	Feature7	Feature8	Feature9	Feature10
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	0.2419	0.1812
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	0.1812	0.2069
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	0.2597	0.1809
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	0.2597	0.1809
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	0.1809	0.2419

5 rows × 31 columns

In [87]:

```
1 #See label distribution
2 sns.countplot(x='y',data=data, palette='hls')
3 plt.show()
```



In [88]:

```
1 count_Benign = len(data[data['y']==1])
2 count_Malignant= len(data[data['y']==0])
3 Total = count_Benign + count_Malignant
4 pct_of_Ben = count_Benign/Total
5 print("Percentage of benign samples in the dataset:", pct_of_Ben*100)
6 pct_of_Malignant = count_Malignant/Total
7 print("Percentage of malignant samples in the dataset: ", pct_of_Malignant*100)
8
```

```
Percentage of benign samples in the dataset: 62.741652021089635
Percentage of malignant samples in the dataset: 37.258347978910365
```

We can see that the classes are slightly imbalanced. The number of samples available for benign class are significantly higher than the malignant class.

In [75]:

```
1 X = data.iloc[:, data.columns != 'y']
2 y = data.iloc[:, data.columns == 'y']
```

In [89]:

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
2                                                    random_state=0)
3 columns = X_train.columns
```

Feature scaling or Feature normalization

Feature normalization is a common pre-processing requirement for many machine learning models. In this example, we use the min-max normalization.

In [90]:

```
1 #Feature scaling or feature normalization
2 from sklearn.preprocessing import MinMaxScaler
3 scaler = MinMaxScaler(feature_range=(0, 1))
4 x_train_scaled = scaler.fit_transform(X_train)
```

In [91]:

```
1 X_train = pd.DataFrame(x_train_scaled, columns=X_train.columns)
2 x_test_scaled = scaler.fit_transform(X_test)
3 X_test = pd.DataFrame(x_test_scaled, columns=X_test.columns)
```

In [92]:

```
1 X_train = X_train.iloc[:,:]
2 X_test = X_test.iloc[:,:]
```

Feature selection

Different kinds of feature selection algorithms can be used select only the most important features. In this example we use the SlectKBest feature selection method.

In [93]:

```
1 from sklearn.feature_selection import SelectKBest
2 from sklearn.feature_selection import mutual_info_classif
3
4 skb = SelectKBest(score_func=mutual_info_classif, k=10)
5 sel_skb = skb.fit(X_train, y_train.values.ravel())
6 selected_features_ind = sel_skb.get_support()
7 #print('scores: ', sel_skb.scores_)
8 #print(sel_skb.index)
```

In [94]:

```
1 X_train_new = X_train.iloc[:, selected_features_ind]
2 X_test_new = X_test.iloc[:, selected_features_ind]
3
4 X_train_new.shape
```

Out[94]:

(398, 10)

Train and test a Logistic Regression model

In [51]:

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn import metrics
```

In [96]:

```
1 #Train a Logistic regression model
2 logreg = LogisticRegression()
3 logreg.fit(X_train_new, y_train.values.ravel())
```

Out[96]:

LogisticRegression()

In [97]:

```
1 #Test the model
2 y_pred = logreg.predict(X_test_new)
```

Results and Analysis

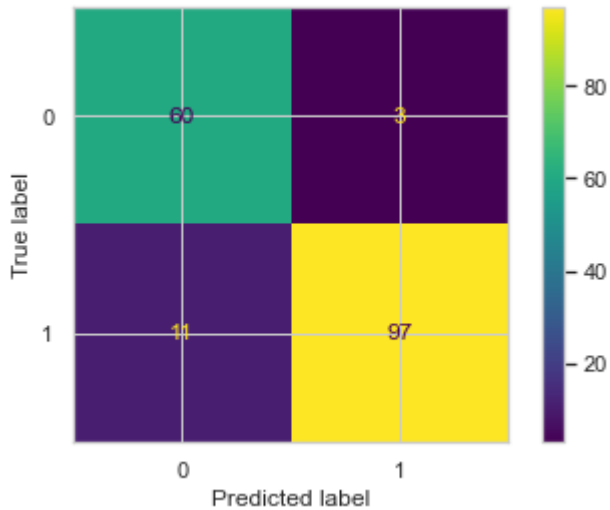
In [99]:

```
1 #Classification accuracy
2 print('Accuracy of the logistic regression classifier on test set: {:.2f}'
3       .format(logreg.score(X_test_new, y_test)))
```

Accuracy of the logistic regression classifier on test set: 0.92

In [100]:

```
1 # Confusion matrix
2 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
3 from sklearn.metrics import plot_confusion_matrix
4 confusion_matrix = confusion_matrix(y_test, y_pred)
5 #print(confusion_matrix)
6 disp = ConfusionMatrixDisplay(confusion_matrix)
7 disp.plot()
8 plt.show()
```



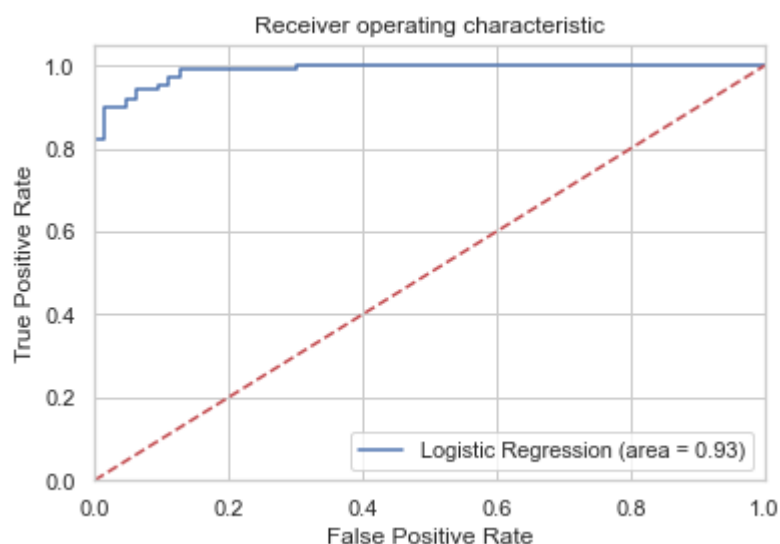
In [103]:

```
1 #Performance measures
2 from sklearn.metrics import classification_report
3 print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.85	0.95	0.90	63
1	0.97	0.90	0.93	108
accuracy			0.92	171
macro avg	0.91	0.93	0.91	171
weighted avg	0.92	0.92	0.92	171

In [101]:

```
1 #ROC Curve
2 from sklearn.metrics import roc_auc_score
3 from sklearn.metrics import roc_curve
4 logit_roc_auc = roc_auc_score(y_test, logreg.predict(X_test_new))
5 fpr, tpr, thresholds = roc_curve(y_test, logreg.predict_proba(X_test_new)[:,1])
6 plt.figure()
7 plt.plot(fpr, tpr, label='Logistic Regression (area = %0.2f)' % logit_roc_auc)
8 plt.plot([0, 1], [0, 1], 'r--')
9 plt.xlim([0.0, 1.0])
10 plt.ylim([0.0, 1.05])
11 plt.xlabel('False Positive Rate')
12 plt.ylabel('True Positive Rate')
13 plt.title('Receiver operating characteristic')
14 plt.legend(loc="lower right")
15 plt.savefig('Log_ROC')
16 plt.show()
```



Model tuning

You can perform the following experiments for experimentally tuning your model to achieve the best possible accuracy.

1. Use different types of feature scaling and analyze the classification accuracy of the model. <https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)
2. Vary parameter k (number of features) in the SelectKBest and see it's effect on the classification accuracy of the model. Can you experimentally find a good parameter k that gives the maximum accuracy?

3. Use a different kind of feature selection method such as Recursive Feature Elimination and see it's effect on the accuracy. https://scikit-learn.org/stable/modules/feature_selection.html (https://scikit-learn.org/stable/modules/feature_selection.html)
4. Use different parameters for the LogisticRegression() model and analyze the results. https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html (https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html)

In []:

1	
---	--