Q.3. Use the dataset Data1.csv to implement the list of classification models below in order to predict whether a tumor is malign or beginning. The dependent variable in this dataset is Class, it is equal to 2 if the tumor is benign and 4 if it is malign. The remaining columns are features that describe different characteristics of each tumor. Import the libraries In [1]: from sklearn import svm import pandas as pd from sklearn.model selection import train test split import numpy as np import matplotlib.pyplot as plt from sklearn import metrics import seaborn as sb Read the data from the csv. In [2]: dataset = pd.read csv('Data1.csv') x = dataset.iloc[:, :-1].valuesy = dataset.iloc[:, -1].valuesNow will split the dataset into the training set and test set In [3]: x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.2, random_state = 0) Feature scaling In [4]: from sklearn.preprocessing import StandardScaler sc = StandardScaler() x train = sc.fit transform(x train) x test = sc.transform(x test) 1. Logistic Regression Import the Logistic Regression library In [5]: from sklearn.linear model import LogisticRegression Create the Logistic Regression object. logisticRegr = LogisticRegression(random state=0) Training the model on the data, storing the information learned from the data Model is learning the relationship between digits (x_train) and labels (y_train). In [7]: logisticRegr.fit(x_train, y_train) LogisticRegression(random state=0) Out[7]: Predict labels for new data. Uses the information the model learned during the model training process. In [8]: logisticRegr.predict(x test[0].reshape(1,-1)) array([2], dtype=int64) Out[8]: Predict for Multiple Observations at Once In [9]: logisticRegr.predict(x test[0:10]) array([2, 2, 4, 4, 2, 2, 2, 4, 2, 2], dtype=int64) Out[9]: Make predictions on entire test data. In [10]: predictions = logisticRegr.predict(x test) Let's see the accuracy of the model using the score method: In [11]: score = logisticRegr.score(x test, y test) print("Accuracy: {:.2f} %".format(score*100)) Accuracy: 95.62 % In [12]: cm = metrics.confusion matrix(y test, predictions) print(cm) [[84 3] [3 47]] According to the results above the accuracy of the model is not guite good - 95.62%. There were 6 values out of 137 that were not predicted correctly. Using K-fold cross-validation method (cv = 10) will ensure that the accuracy score calculated only on one test set was not just a lucky occurence. To be more exactly in k-fold, with K=10, the dataset is devided into 10 subsets. At the first iteration the model is trained on the first 9 subsets and tested on the 10th subset. In the second iteration, the model is trained on the first 8 subsets and 10th, and is tested on the 9th subset. In this way there are 10 iterations and we obtain 10 accuracies of the model. Will print the mean accuracy score and the standard deviation of the score computed. In [13]: from sklearn.model selection import cross val score logisticRegr accuracies = cross val score(estimator = logisticRegr, X = x train, y = y train, cv = 10) print("Accuracy: {:.2f} %".format(logisticRegr accuracies.mean()*100)) print("Standard Deviation: {:.2f} %".format(logisticRegr accuracies.std()*100)) logisticRegr accuracies Accuracy: 96.70 % Standard Deviation: 1.97 % , 0.94545455, array([0.94545455, 0.96363636, 0.96363636, 1. Out[13]: , 0.96296296, 0.96296296, 0.98148148, 0.94444444]) The mean accurracy for 10 test sets is 96.7% that is not much bigger that the accuracy of the model calculated on only one train data set (95.62%). 2. Decision Tree Classifier Let's create a Decision Tree Model using Scikit-learn. Import Decision Tree Classifier In [14]: from sklearn.tree import DecisionTreeClassifier Create Decision Tree classifer object. In [15]: dtc = DecisionTreeClassifier(random state=0) Train Decision Tree Classifer. In [16]: dtc = dtc.fit(x train, y train) Predict the response for test dataset. In [17]: y pred = dtc.predict(x test) Let's estimate, how accurately the classifier or model can predict whether a tumor is malign or bening. Accuracy can be computed by comparing actual test set values and predicted values. In [18]: print("Accuracy: {:.2f} %".format(metrics.accuracy_score(y_test, y_pred)*100)) Accuracy: 92.70 % In [19]: cm = metrics.confusion_matrix(y_test, y_pred) print(cm) [[80 7] [3 47]] We see that this model has a big accuracy that means it is good to use for predicting whether a tumor is malign or benign. The confusion matrix also shows that only 10 out of 137 objects are not predicted as they had to. Using K-fold cross-validation method (cv = 10) will ensure that the accuracy score calculated only on one test set was not just a lucky occurrence. To be more exactly in k-fold, with K=10, the dataset is devided into 10 subsets. At the first iteration the model is trained on the first 9 subsets and tested on the 10th subset. In the second iteration, the model is trained on the first 8 subsets and 10th, and is tested on the 9th subset. In this way there are 10 iterations and we obtain 10 accuracies of the model. Will print the mean accuracy score and the standard deviation of the score computed. In [20]: from sklearn.model_selection import cross_val_score dtc_accuracies = cross_val_score(estimator = dtc, X = x_train, y = y_train, cv = 10) print("Accuracy: {:.2f} %".format(dtc_accuracies.mean()*100)) print("Standard Deviation: {:.2f} %".format(dtc_accuracies.std()*100)) dtc_accuracies Accuracy: 93.59 % Standard Deviation: 2.21 % array([0.92727273, 0.90909091, 0.92727273, 0.96363636, 0.94545455, Out[20]: 0.92727273, 0.90740741, 0.92592593, 0.98148148, 0.94444444]) The mean accuracy calculated for 10 subsets of data is bigger with 0.89% that the accuracy of the model calculated on only one test set. 3. Random Forest Classifier (with nb_trees = 10) Import Random Forest Model. In [21]: from sklearn.ensemble import RandomForestClassifier Create a Gaussian Classifier. In [22]: rfc=RandomForestClassifier(n_estimators=10) Train the model using the training sets y_pred=clf.predict(X_test) In [23]: rfc.fit(x_train,y_train) RandomForestClassifier(n estimators=10) Out[23]: Predict the response for test dataset. In [24]: y_pred=rfc.predict(x test) After training, check the accuracy using actual and predicted values. In [25]: print("Accuracy: {:.2f} %".format(metrics.accuracy_score(y_test, y_pred)*100)) Accuracy: 95.62 % In [26]: cm = metrics.confusion matrix(y test, y pred) print(cm) [[84 3] [3 47]] We can observe that the accuracy of the model is pretty high - 96,35%. There are only 5 values that were not predicted to the correct Using K-fold cross-validation method (cv = 10) will ensure that the accuracy score calculated only on one test set was not just a lucky occurence. To be more exactly in k-fold, with K=10, the dataset is devided into 10 subsets. At the first iteration the model is trained on the first 9 subsets and tested on the 10th subset. In the second iteration, the model is trained on the first 8 subsets and 10th, and is tested on the 9th subset. In this way there are 10 iterations and we obtain 10 accuracies of the model. Will print the mean accuracy score and the standard deviation of the score computed. In [27]: from sklearn.model_selection import cross_val_score rfc_accuracies = cross_val_score(estimator = rfc, X = x_train, y = y_train, cv = 10) print("Accuracy: {:.2f} %".format(rfc accuracies.mean()*100)) print("Standard Deviation: {:.2f} %".format(rfc accuracies.std()*100)) rfc accuracies Accuracy: 97.07 % Standard Deviation: 2.03 % array([0.92727273, 0.96363636, 0.98181818, 0.98181818, 0.96363636, Out[27]: , 0.96296296, 0.96296296, 1. , 0.96296296]) Here we can observe that the mean value of the accuracy for the model calculated on 10 test datasets (96.15%) is bigger than the value of the accuracy on the model calculated on only one test dataset (95.62%). 4. K- Nearest Neighbors (K-NN) Import the KNeighborsClassifier library In [28]: from sklearn.neighbors import KNeighborsClassifier Create the object of KNeighborsClassifier type with number of neighbors =3. In [29]: neigh = KNeighborsClassifier(n neighbors=3) Train the KNeighbors Classifier. In [30]: neigh.fit(x train, y train) KNeighborsClassifier(n_neighbors=3) Out[30]: Predict the response for test dataset. In [31]: y pred = neigh.predict(x test) After training, check the accuracy using actual and predicted values. In [32]: print("Accuracy: {:.2f} %".format(metrics.accuracy score(y test, y pred)*100)) Accuracy: 97.08 % In [33]: cm = metrics.confusion_matrix(y_test, y_pred) print(cm) [[85 2] [2 48]] We can observe that this model is better that the previous ones. This has 97.08% of accuracy and only 4 values out of 137 are predicted wrongly. For sure is one of the best candidates for the best models. Using K-fold cross-validation method (cv = 10) will ensure that the accuracy score calculated only on one test set was not just a lucky occurence. To be more exactly in k-fold, with K=10, the dataset is devided into 10 subsets. At the first iteration the model is trained on the first 9 subsets and tested on the 10th subset. In the second iteration, the model is trained on the first 8 subsets and 10th, and is tested on the 9th subset. In this way there are 10 iterations and we obtain 10 accuracies of the model. Will print the mean accuracy score and the standard deviation of the score computed. In [34]: from sklearn.model_selection import cross_val_score neigh_accuracies = cross_val_score(estimator = neigh, X = x_train, print("Accuracy: {:.2f} %".format(neigh_accuracies.mean()*100)) print("Standard Deviation: {:.2f} %".format(neigh_accuracies.std()*100)) neigh accuracies Accuracy: 96.34 % Standard Deviation: 2.44 % array([0.90909091, 0.94545455, 0.98181818, 0.98181818, 0.94545455, Out[34]: , 0.96296296, 0.96296296, 0.98148148, 0.96296296]) We can observe that the accuracy for the model calculated on only one test dataset is bigger than the mean accuracy for the model calculated on 10 smaller test data sets with 0.74%. Here we can say that there was a little lucky occurences. 5. Naive Bayes Import the GaussianNB library In [35]: from sklearn.naive_bayes import GaussianNB Create the GaussianNB object and train it on the train data. In [36]: gnb = GaussianNB() gnb.fit(x_train, y_train) GaussianNB() Out[36]: Making predictions on the testing set In [37]: y_pred = gnb.predict(x_test) Comparing actual response values (y_test) with predicted response values (y_pred) In [38]: from sklearn import metrics print("Accuracy: {:.2f}%".format(metrics.accuracy_score(y_test, y_pred)*100)) Accuracy: 94.89% In [39]: cm = metrics.confusion matrix(y test, y pred) print(cm) [[80 7] [0 50]] The model has a pretty good accuracy- 94.89%. There are only 7 values that were predicted wrongly. Using K-fold cross-validation method (cv = 10) will ensure that the accuracy score calculated only on one test set was not just a lucky occurence. To be more exactly in k-fold, with K=10, the dataset is devided into 10 subsets. At the first iteration the model is trained on the first 9 subsets and tested on the 10th subset. In the second iteration, the model is trained on the first 8 subsets and 10th, and is tested on the 9th subset. In this way there are 10 iterations and we obtain 10 accuracies of the model. Will print the mean accuracy score and the standard deviation of the score computed. In [40]: from sklearn.model selection import cross val score gnb accuracies = cross val score(estimator = gnb, X = x_train, y = y_train, cv = 10) print("Accuracy: {:.2f} %".format(gnb accuracies.mean()*100)) print("Standard Deviation: {:.2f} %".format(gnb_accuracies.std()*100)) gnb_accuracies Accuracy: 96.16 % Standard Deviation: 1.91 % array([0.92727273, 0.96363636, 0.96363636, 0.96363636, 0.96363636, Out[40]: 0.96363636, 0.98148148, 0.94444444, 1. There can be observed that the accuracy score of the model calculated on only one test set is smaller that the accuracy mean calculated on 10 smaller data sets. (94.89% < 96.16%) 6. Support Vector Machine (SVM) Train the SVC classifier In [41]: clf = svm.SVC(random state=0) clf.fit(x train, y train) SVC(random_state=0) Out[41]: Make predictions and compare the results by displaying the predicted values of y next to the test values of y in a two-dimensional array. In [42]: y_pred = clf.predict(x_test) y_pred array([2, 2, 4, 4, 2, 2, 2, 4, 2, 2, 4, 2, 4, 2, 2, 4, 4, 4, 4, 4, 2, 2, 2, Out[42]: 4, 2, 4, 4, 2, 2, 2, 4, 2, 4, 4, 2, 2, 2, 4, 4, 2, 4, 2, 2, 2, 2, 2, 2, 2, 4, 2, 2, 4, 2, 4, 2, 2, 2, 4, 4, 2, 4, 2, 2, 2, 2, 2, 2, 2, 2, 2, 4, 4, 2, 2, 2, 2, 2, 2, 4, 2, 2, 4, 2, 4, 2, 4, 2, 4, 2, 4, 2, 4, 4, 2, 4, 2, 4, 4, 2, 4, 4, 4, 4, 2, 2, 2, 4, 4, 2, 2, 4, 4, 2, 2, 4, 2, 2, 4, 2, 2, 2, 4, 2, 2, 2, 4, 2, 2, 4, 4, 2, 4, 2, 4, 4, 2, 4, 2, 2, 4, 2], dtype=int64) Create and print the Confusion Matrix and the accuracy score of the model and interpret the two. In [43]: from sklearn.metrics import confusion matrix, accuracy score cm = confusion_matrix(y_test, y_pred) print (cm) print("SVC model accuracy: {:.2f}%".format(accuracy_score(y_test, y_pred)*100)) [[82 5] [1 49]] SVC model accuracy: 95.62% We can observe that this is a good model with an accuracy of 95.62%. There are 6 values that were not predicted correctly. Using K-fold cross-validation method (cv = 10) will ensure that the accuracy score calculated only on one test set was not just a lucky occurence. To be more exactly in k-fold, with K=10, the dataset is devided into 10 subsets. At the first iteration the model is trained on the first 9 subsets and tested on the 10th subset. In the second iteration, the model is trained on the first 8 subsets and 10th, and is tested on the 9th subset. In this way there are 10 iterations and we obtain 10 accuracies of the model. Will print the mean accuracy score and the standard deviation of the score computed. In [44]: from sklearn.model selection import cross val score svm_accuracies = cross_val_score(estimator = clf, X = x_train, y = y_train, cv = 10) print("Accuracy: {:.2f} %".format(svm_accuracies.mean()*100)) print("Standard Deviation: {:.2f} %".format(svm_accuracies.std()*100)) svm accuracies Accuracy: 96.53 % Standard Deviation: 1.91 % array([0.92727273, 0.96363636, 0.96363636, 0.98181818, 0.96363636, Out[44]: 0.96363636, 0.98148148, 0.96296296, 1. , 0.9444444)) We can observe that the accuracy score calculated for a only one test set is smaller with 0.92% than the accuracy mean calculated for 10 test sets. Q.4. Choose the best performing model based on the results from performing the k-fold cross-validation. Discuss your choice. In [55]: lg_accuracy = logisticRegr_accuracies.mean()*100 lg_std = logisticRegr_accuracies.std()*100 dtc_accuracy = dtc_accuracies.mean()*100 dtc_std = dtc_accuracies.std()*100 rfc_accuracy = rfc_accuracies.mean()*100 rfc_std = rfc_accuracies.std()*100 neigh_accuracy = neigh_accuracies.mean()*100 neigh_std = neigh_accuracies.std()*100 gnb_accuracy = gnb_accuracies.mean()*100 gnb_std = gnb_accuracies.std()*100 svm_accuracy = svm_accuracies.mean()*100 svm_std = svm_accuracies.std()*100 print("Mean accuracy for Logistic Regression model: {:.2f} %, Std: {:.2f}.".format(lg_accuracy, lg_std)) print("Mean accuracy for Decision Tree Classifier model: {:.2f}%, Std: {:.2f}.".format(dtc_accuracy, dtc_std)) print("Mean accuracy for Random Forest Classifier with 10 nb_trees: {:.2f} %, Std: {:.2f}.".format(rfc_accuracy print("Mean accuracy for K-Nearest Neighbors: {:.2f} %, Std: {:.2f}.".format(neigh_accuracy, neigh_std)) print("Mean accuracy for Naive Bayes: {:.2f} %, Std: {:.2f}.".format(gnb_accuracy, gnb_std)) print("Mean accuracy for Support Vector Machine: {:.2f} %, Std: {:.2f}.".format(svm_accuracy, svm_std)) lg1 = lg_accuracy-lg_std dtc1 = dtc_accuracy - dtc_std rfc1 = rfc_accuracy - rfc_std neigh1 = neigh_accuracy - neigh_std gnb1 = gnb_accuracy - gnb_std svm1 = svm_accuracy - svm_std map = {"Logistic Regression":lg1, "Decision Tree Classifier" : dtc1, "Random Forest Classifier" : rfc1, "K-Near "Naive Bayes" : gnb1, "Support Vector Machine" : svm1} print ("\n Below we can see the results for all 6 models of the difference calculation between mean accuracy ar print(map) max = max(map.values()) print("The best result is: {:.2f}".format(max)) Mean accuracy for Logistic Regression model: 96.70 %, Std: 1.97. Mean accuracy for Decision Tree Classifier model: 93.59%, Std: 2.21. Mean accuracy for Random Forest Classifier with 10 nb_trees: 97.07 %, Std: 2.03. Mean accuracy for K-Nearest Neighbors: 96.34 %, Std: 2.44. Mean accuracy for Naive Bayes: 96.16 %, Std: 1.91. Mean accuracy for Support Vector Machine: 96.53 %, Std: 1.91. Below we can see the results for all 6 models of the difference calculation between mean accuracy and standard deviation. {'Logistic Regression': 94.73053901089192, 'Decision Tree Classifier': 91.38701005825551, 'Random Forest Classi fier': 95.03858734142017, 'K-Nearest Neighbors': 93.90304295607386, 'Naive Bayes': 94.2457764262916, 'Support V ector Machine': 94.61833923684192} The best result is: 95.04 We can observe that the best mean accuracy is for Random Forest Classifier model. Q. 5. For the chosen best performing model select two hyper parameters you would like to tune using the grid search method. Learn what the two chosen hyper parameters do to your model. Look into the description of hyper parameters for your respective model by using the scikit-org API https://scikitlearn.org/stable/modules/classes.html. Discuss here the two hyper parameters you chose and how they impact the model. Choose between 2 and 4 different values for each of these hyper parameters to test with grid search. Explain your choice and, if you have any, make some expectations (you could speculate for example which values of the chosen hyper parameters you believe will make the model perform best and why). Use grid search to find the best model. Print the accuracy of this model and the hyper parameter values of the best model. Do these values confirm your intuition? Firstly let's see what parameters has the classifier of Random Forest Classifier model. In [46]: rfc.get_params().keys() dict keys(['bootstrap', 'ccp alpha', 'class weight', 'criterion', 'max depth', 'max features', 'max leaf node Out[46]: s', 'max samples', 'min impurity decrease', 'min samples leaf', 'min samples split', 'min weight fraction lea f', 'n estimators', 'n jobs', 'oob score', 'random state', 'verbose', 'warm start']) In [47]: params = {"n_estimators": [1, 2, 3, 4, 5, 6, 7, 8, 9, 10], "criterion": ["gini", "entropy"], "max_depth": [1, 2, 3, 4, 5, 6, 7, 8, 9, None], "max_features": ["auto", "sqrt", "log2"]} Description of some hyper parameters that I chose: 1. n_estimators - Number of trees in the forest. 2. criterion - The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain. Gini Impurity is a measurement of the likelihood of an incorrect classification of a new instance of a random variable, if that new instance were randomly classified according to the distribution of class labels from the data set. Entropy is a measure of information that indicates the disorder of the features with the target. 3. max_depth - The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples. 4. max features - The number of features to consider when looking for the best split: - If int, then consider max_features features at each split. - If float, then max features is a fraction and round(max features * n features) features are considered at each split. If "auto", then max_features=sqrt(n_features). - If "sqrt", then max features=sqrt(n features) (same as "auto"). - If "log2", then max features=log2(n features). - If None, then max_features=n_features. In [48]: from sklearn.model selection import GridSearchCV grid search = GridSearchCV(estimator = rfc, param grid= params, cv = 10,scoring = 'accuracy', n jobs = -1)grid_search.fit(x_train, y_train) best_accuracy = grid_search.best_score_ best parameters = grid search.best params print("Best Accuracy: {:.2f} %".format(best accuracy*100)) print("Best Parameters:", best_parameters) Best Accuracy: 97.98 % Best Parameters: {'criterion': 'gini', 'max_depth': 3, 'max_features': 'log2', 'n_estimators': 9} We can observe that for the best accuracy we must have the values for parameters as: "criterion" = "gini", "max_depth" = 3, "max_features" = sqrt, "n estimators": 6. This means that the number of trees in the forest must be 3. The maximum depth of the tree is 3. The criterion for measure the quality of the split is "gini". The nuumber of features to consider when looking for the best split is "sqrt" of the total number of features. Using these parameters will obtain an accuracy of 97.98%.