Problem 3: Be Discrete author: David Nogales Perez Code for obtaining the dataset #install library to import dataset #!pip3 install pmlb #from pmlb import fetch\_data #appendicitis\_data = fetch\_data('appendicitis') #appendicitis\_data.to\_csv('appendicitis.csv', index=False) Importing necessary libraries and defining important functions import pandas as pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns from sklearn.model\_selection import train\_test\_split,RandomizedSearchCV,GridSearchCV,cross\_val\_score from sklearn.preprocessing import StandardScaler,KBinsDiscretizer from sklearn.naive\_bayes import GaussianNB from sklearn.linear\_model import LogisticRegression from sklearn.neighbors import KNeighborsClassifier from sklearn.metrics import classification\_report,confusion\_matrix,ConfusionMatrixDisplay,confusion\_matrix sns.set(rc={'figure.figsize':(10, 7.5)}) sns.set\_color\_codes("deep") sns.set\_style("whitegrid") def preprocess\_data(X\_train, X\_test,bins=0,new\_strategy=''): X\_train\_preprocessed = X\_train.copy(deep=True) X\_test\_preprocessed = X\_test.copy(deep=True) #Discretize variables according to parameters if (bins!=0): kBinDiscretizer = KBinsDiscretizer(n\_bins=bins, encode='onehot-dense', strategy=new\_strategy) kBinDiscretizer.fit(X\_train) X\_train\_preprocessed = kBinDiscretizer.transform(X\_train\_preprocessed) kBinDiscretizer.fit(X\_test) X\_test\_preprocessed = kBinDiscretizer.transform(X\_test\_preprocessed) #Scale variables scale = StandardScaler() X\_train\_preprocessed = scale.fit\_transform(X\_train\_preprocessed) X\_test\_preprocessed = scale.transform(X\_test\_preprocessed) return X\_train\_preprocessed,X\_test\_preprocessed def generate\_result(X\_p\_train, X\_p\_test, y\_train, y\_test,n\_neighbors=5): models = { "naive\_bayes":GaussianNB(), "logistic\_regression":LogisticRegression(), "KNN":KNeighborsClassifier(n\_neighbors=n\_neighbors) results\_numeric = {} results\_object = {} for name, model in models.items(): current\_model = model.fit(X\_p\_train,y\_train) y\_pred = current\_model.predict(X\_p\_train) model\_cvs = cross\_val\_score(model,X\_p\_train,y\_train,cv=5) model\_test = current\_model.score(X\_p\_test, y\_test) model\_cls\_report = classification\_report(y\_train, y\_pred,output\_dict=True) model\_conf\_matrix = confusion\_matrix(y\_train, y\_pred, labels=current\_model.classes\_) result\_numeric = { "cross\_validation\_error\_avg":model\_cvs.mean(), "test\_error":model\_test, result\_object = { "confusion\_matrix":model\_conf\_matrix, "classification\_report":model\_cls\_report results\_numeric[name] = result\_numeric results\_object[name] = result\_object return results\_numeric,results\_object def tune\_knn\_hyperparameters(X\_p\_train,y\_train): param\_grid = {'n\_neighbors': np.arange(1,20)} knn\_best = GridSearchCV(KNeighborsClassifier(), param\_grid, cv=5) knn\_best.fit(X\_p\_train, y\_train) print("GridSearch: n\_neightbors ",knn\_best.best\_params\_['n\_neighbors']) return knn\_best.best\_params\_['n\_neighbors'] def results\_object\_to\_dataframe(results\_object): clfs\_dfs = [] conf\_mat\_dfs = [] for \_,result in results\_object.items(): conf\_mat\_df = pd.DataFrame(result["confusion\_matrix"]) clf\_rep\_df = pd.DataFrame(result["classification\_report"]) clfs\_dfs.append(clf\_rep\_df) conf\_mat\_dfs.append(conf\_mat\_df) clfs\_df = pd.concat(clfs\_dfs,keys=results\_object.keys()).round(2) conf\_mat\_df = pd.concat(conf\_mat\_dfs,keys=results\_object.keys()) return clfs\_df,conf\_mat\_df def generate\_heatmaps\_from\_dataframe(df,index\_labels,ax): indexes = range( len(index\_labels)) for index in indexes: sns.heatmap(df.loc[index\_labels[index]],annot=True,ax=ax[index]) ax[index].set\_title(index\_labels[index]) #appendicitis\_data = fetch\_data('appendicitis') appendicitis\_data = pd.read\_csv('appendicitis.csv') appendicitis\_data Out[ ]: At1 At2 At3 At4 At5 At6 At7 target **0** 0.213 0.554 0.207 0.000 0.000 0.749 0.220 **1** 0.458 0.714 0.468 0.111 0.102 0.741 0.436 **2** 0.102 0.518 0.111 0.056 0.022 0.506 0.086 **3** 0.187 0.196 0.105 0.056 0.029 0.133 0.085 **4** 0.236 0.804 0.289 0.111 0.066 0.756 0.241 **101** 0.449 0.875 0.523 0.083 0.076 0.920 0.487 **102** 0.102 0.000 0.022 0.000 0.000 0.000 0.017 **103** 0.409 0.875 0.482 0.306 0.259 0.914 0.443 **104** 0.427 0.804 0.474 0.056 0.048 0.836 0.437 **105** 0.462 0.911 0.551 0.167 0.154 0.931 0.500 106 rows × 8 columns a) Split data into train and test (60%/40%) In [ ]: X = appendicitis\_data.drop("target",axis=1) y = appendicitis\_data.target # Split data 60%-40% into training set and test set X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.40, random\_state=42) sns.histplot(y) Out[]: <AxesSubplot:xlabel='target', ylabel='Count'> 70 60 50 30 20 10 8.0 0.0 0.2 0.4 0.6 1.0 target The histogram of the target variable reveals the imbalance in our dataset. b) Preprocess adequately for näive bayes, logistic regression and KNN, train and obtain: • the crossvalidation error the test error confusion matrix • classification report For this step, we first preprocess the data, then explore the KNN hyperparameters and, lastly, we generate the baseline results. In [ ]: X\_p\_train,X\_p\_test = preprocess\_data(X\_train, X\_test) n\_neighbors = tune\_knn\_hyperparameters(X\_p\_train,y\_train) results\_numeric,results\_object = generate\_result(X\_p\_train,X\_p\_test, y\_train, y\_test,n\_neighbors) baseline\_results = pd.DataFrame(results\_numeric) baseline\_results.T.sort\_values('test\_error').round(3) GridSearch: n\_neightbors 7 Out[ ]: cross\_validation\_error\_avg test\_error naive\_bayes 0.888 0.814 **KNN** 0.888 0.860 logistic\_regression 0.890 0.884 As we can see in our preliminary results, all models have desirable scores. Although Näive Bayes presents some overfitting. In [ ]: clfs\_df,conf\_mat\_df =results\_object\_to\_dataframe(results\_object) clfs\_df Out[ ]: accuracy macro avg weighted avg naive\_bayes precision 0.94 0.80 0.89 0.67 0.89 **recall** 0.92 0.73 0.83 0.89 0.89 **f1-score** 0.93 0.70 0.89 0.81 0.89 **support** 52.00 11.00 0.89 63.00 63.00 logistic\_regression precision 0.93 0.88 0.92 0.92 0.90 **recall** 0.98 0.64 0.92 0.92 0.81 **f1-score** 0.95 0.74 0.92 0.85 0.92 **support** 52.00 11.00 0.92 63.00 63.00 0.92 0.81 KNN precision 0.70 0.89 0.89 **recall** 0.94 0.64 0.89 0.79 0.89 **f1-score** 0.93 0.80 0.67 0.89 0.89 **support** 52.00 11.00 0.89 63.00 63.00 The classification reports of each models are very similar, but looking closely, Logistic Regression offers a slighly better precision for both types of patients. Although that might not be that relevant since the number of patients with appendicitis (label 1) is small compared to the patients without it (label 0) In [ ]: fig, ax = plt.subplots(nrows=1, ncols=3,figsize=(20, 5)) index\_labels = list(results\_object.keys()) generate\_heatmaps\_from\_dataframe(conf\_mat\_df,index\_labels,ax) fig.suptitle("Baseline Confusion Matrices", fontsize=16) Text(0.5, 0.98, 'Baseline Confusion Matrices') **Baseline Confusion Matrices** logistic\_regression KNN naive\_bayes 51  $\circ$ - 30 - 25 - 20 0 0 0 The confusion matrices of each of the models show similar results as well as their respective classification reports. In this case we can observe that Näive Bayes, which has the highest recall in the report, predicts 3 false negatives which only differs with the other models by 1. This means that each 10% of the recall is due to a single patient with appendicitis classified without it. c) Use KbinsDiscretizer to discretize continuous data. Apply the following discretizations: quantile kmeans Obtaining a discretization into 2,3,4 values.(Use onehot-dense as encoding to set each value as attribute) Fit the models and compare results. • Has this process affected to the quality of the models? Which one would you chose and why? For each strategy and number of bins we preprocess the data with the given strategy and number of bins, obtain the KNN hyperparameters and generate the results for later inspection. In [ ]: #Discretize variables and generate results for various parameters strategies = ['quantile','kmeans']  $n_{bins_{list}} = [2,3,4]$ discretization\_results = {} for strategy in strategies: aux\_results = {} for n\_bins in n\_bins\_list: X\_p\_train,X\_p\_test = preprocess\_data(X\_train, X\_test,bins=n\_bins,new\_strategy=strategy) n\_neighbors = tune\_knn\_hyperparameters(X\_p\_train,y\_train) results,\_ = generate\_result(X\_p\_train,X\_p\_test, y\_train, y\_test,n\_neighbors) aux\_results[n\_bins] = results discretization\_results[strategy] = aux\_results GridSearch: n\_neightbors 3 GridSearch: n\_neightbors 9 GridSearch: n\_neightbors 2 GridSearch: n\_neightbors 11 GridSearch: n\_neightbors 7 GridSearch: n\_neightbors 11 Transform the results to a Dataframe for easier manipulation. In [ ]: baseline\_results\_n = baseline\_results.T.copy(deep=True) baseline\_results\_n.index = pd.MultiIndex.from\_product([[0], baseline\_results\_n.index ], names=['bins', 'model']) quantile\_df = pd.concat({ k: pd.DataFrame.from\_dict(v, 'index') for k, v in discretization\_results['quantile'].items() axis=0) kmeans\_df = pd.concat({ k: pd.DataFrame.from\_dict(v, 'index') for k, v in discretization\_results['kmeans'].items() }, axis=0) results\_df = pd.concat([ quantile\_df, kmeans\_df,baseline\_results\_n ], keys=['quantile','kmeans','baseline'],axis=0) results\_df.index.names = ['strategy','bins','model'] sorted\_results = results\_df.sort\_values('test\_error',ascending=False).round(3) sorted\_results cross\_validation\_error\_avg test\_error Out[ ]: strategy bins model 0 logistic\_regression baseline 0.890 0.884 KNN 4 0.921 0.884 kmeans 3 logistic\_regression 0.764 0.884 KNN 0.888 0.860 baseline KNN 0.873 0.860 kmeans KNN 0.905 0.860 quantile 4 3 logistic\_regression 0.735 0.837 2 logistic\_regression 0.873 0.814 0.888 0.814 baseline naive\_bayes KNN 2 0.873 0.814 kmeans 0.814 logistic\_regression 0.872 KNN 0.888 0.814 quantile KNN 0.873 0.814 4 logistic\_regression 0.791 0.844 kmeans naive\_bayes 0.779 0.791 quantile 4 logistic\_regression 0.744 0.810 naive\_bayes 0.744 0.812 naive\_bayes 0.721 0.619 kmeans 0.698 naive\_bayes 0.697 quantile naive\_bayes 0.442 kmeans 3 0.415 naive\_bayes 0.419 2 0.447 Looking at the Dataframe with the final results, one can see that the models with better scores are mainly Logistic Regression and KNN with different strategies In [ ]: ax = sorted\_results.plot(rot=90,xticks=np.arange(0,21)) plt.suptitle("Cross Validation and Test errors of Trained Models") Out[]: Text(0.5, 0.98, 'Cross Validation and Test errors of Trained Models') Cross Validation and Test errors of Trained Models cross\_validation\_error\_avg test\_error 0.9 0.8 0.7 0.6 0.5 0.4 (kmeans, 2, KNN) (kmeans, 4, logistic\_regression) logistic\_regression) 2, logistic\_regression) logistic\_regression) (quantile, 2, KNN) strategy,bins,model Looking at the above graph, one can see that several models present some degree of overfitting. Interestingly the baseline model Logistic Regression is the best model overall, with the best score and little overfitting. Following the Occam's razor principle one could pick the baseline logistic regression model since it's the simplest model overall, but given the sensitive nature of our dataset it would be advisable to explore more alternatives to improve the model. fig, ax = plt.subplots(nrows=1, ncols=2,figsize=(15, 7),sharey=True) fig.suptitle("Scores and Bins Using Different Strategies", fontsize=16) df = sorted\_results.loc["quantile"].sort\_values('model') sns.lineplot(x='bins', y="test\_error", hue="model", data=df,ax=ax[0]).set\_title("Quantile Strategy") df = sorted\_results.loc["kmeans"].sort\_values('model') sns.lineplot(x='bins', y="test\_error", hue="model", data=df, ax=ax[1]).set\_title("Kmeans Strategy") Text(0.5, 1.0, 'Kmeans Strategy') Scores and Bins Using Different Strategies Quantile Strategy Kmeans Strategy 0.9 0.8 0.7 model model logistic\_regression naive\_bayes 3.00 3.25 2.50 3.50 3.75 2.50 2.75 3.50 3.75 2.75 bins bins Here we can see how the binarization affects the different models. For example Näive Bayes using the quantile strategy decreases its accuracy when we increase the number of bins whereas using Kmeans as the discretization strategy follows an inverse trend which with sufficient bins maybe could improve the model but for the bins tested with either strategy we don't achieve better results than the baseline. For Logistic regression we can see that with 3 bins we find a local maximum in both strategies. Also when using the Kmeans strategy we obtain the same result as the KNN model with less bins. Lastly, we can see that by discretizing the variables and using more bins we seem to always improve the KNN model. d) Considering the interpretability of the models: • do you think that has advantages to work with discretized data? • why? As we have seen previously, working with discretized data improves the performance of some models in certain conditions (Logistic regression with 3 bins, KNN with more bins) but for others (Bayes) it worsens the performance. Furthermore, given that the best model in our tests doesn't even require to discretize data shows that the advantage is minimal to none (at least with the number of bins tested). Maybe exploring a wider range of bins would allow us to be able to obtain some meaningful results. But with the current results it doesn't seem to have a significant impact.

**Assignment 3**