Trabalho A2 - ML Carlos César de O. Fonseca import pandas as pd import matplotlib.pyplot as plt import numpy as np from funcoes import evalu Contextualizando O meu trabalho é sobre Dados que foram coletados apartir de uma imagem digitalizada de um aspirado por agulha fina (PAAF) de uma massa mamária. Eles descrevem características dos núcleos celulares presentes na imagem. Apartir dai, temos vários dados quantitativos que representam algumas coisas e temos uma coluna de dados qualitativa, os dados do Diagnostico. Esses dados a principio era uma coluna de B's e M's, onde B significava Benigno, enquanto M significava Maligno. Porém para rodar os modelos e desenvolver uma forma de avaliar se o tumor é benigno ou maligno dado apenas os dados quantitativos da interpretação da imagem. O objetivo ao fazer a analise dos dados era prever prever os malignos como malignos e os benignos como benignos ou seja a métrica de interesse era a acurácia. Antes de falar de acurácia, vamos definir algumas coisas: ullet VP (Verdadeiro positivo) será quando o modelo prever M e target for M• FP (Falso Positivo) será quando o modelo prever M porém o target é Bullet FN (Falso Negativo) será quando o modelo prever B enquanto o real é M• VN (Verdadeiro Negativo) será quando o modelo prever B e realmente for BPortanto a acurácia seria:  $Acuracia = \frac{VP + VN}{VP + FP + FN + VN}$ Esse valor é legal de ser analisado e também diz bastante sobre a quantidade de acertos que o modelo tem. Porém, acho que seria bom analisar outras métricas como Precisão e Recall.  $Precis\~{a}oM = rac{VP}{VP + FP}$  $Recall M = rac{VP}{VP + FN}$ Acho que dados os meus dados é tão importante observar o Recall de Malignos, quanto analisar a Acurácia. Os meus dados tem 357 benignos enquanto tem 212 malignos como target, acho que não é um problema para predizer, porém é necessário ter cuidado na hora de analisar os resultados, pois os dados cerca de  $\frac{3}{5}$  são Benignos enquanto apenas cerca  $\frac{2}{5}$  são malignos, tem um desbalanciamento, porém não é tão grande e acredito que apenas tomando cuidado com as medidas de avaliação seja possível desenvolver o melhor modelo. Além disso, acho importante analisar o Recall dos Malignos, pois a previsão de Falsos Negativos é mais negativo que os Falsos Positivos, FP acarretam alguns danos como uma ansiedade, preocupação das pessoas que foram diagnosticadas e da família dela, porém o FN pode ser pior, a pessoa em questão acredita que o tumor é benigno e para o tratamento, levando em um aumento da doença e numa menor chance de cura, possivelmente levando a óbito. Resumo dos Passos... • Ler os dados Trocar M por 1 e B por 0 na variável diagnostico Separar os dados entre treino e teste Testar alguns modelos com os dados Escolher o melhor modelo e chegar a uma conclusão Leitura dos Dados Aqui abrimos os dados que já utilizamos no trabalho 1 e convertemos B para 0 e M para 1: df2 = pd.read csv('cancer.csv') df2['diagnosis'] = df2['diagnosis'].replace('M', 1) df2['diagnosis'] = df2['diagnosis'].replace('B', 0) df2.diagnosis.astype('int64') df2.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 569 entries, 0 to 568 Data columns (total 12 columns): Non-Null Count Dtype Column int64 569 non-null 0 id int64 1 diagnosis 569 non-null radius mean 569 non-null float64 3 texture mean 569 non-null float64 569 non-null 4 perimeter mean float64 float64 5 area mean 569 non-null 6 smoothness mean 569 non-null float64 7 compactness mean 569 non-null float64 8 concavity\_mean 569 non-null float64 concave points\_mean 569 non-null float64 10 symmetry\_mean 569 non-null float64 11 fractal dimension mean 569 non-null float64 dtypes: float64(10), int64(2) memory usage: 53.5 KB Observando os valores que usaremos para prever, ou seja, o X: df2.iloc[:,2:] radius\_mean texture\_mean perimeter\_mean area\_mean smoothness\_mean compactness\_mean concavity 0 17.99 10.38 122.80 1001.0 0.11840 0.27760 ( 20.57 17.77 132.90 1326.0 0.08474 0.07864 ( 2 19.69 21.25 130.00 1203.0 0.10960 0.15990 ( 3 20.38 77.58 386.1 0.14250 0.28390 11.42 ( 0.10030 0.13280 4 20.29 14.34 135.10 1297.0 ( 564 21.56 22.39 142.00 1479.0 0.11100 0.11590 ( 0.09780 565 20.13 28.25 131.20 1261.0 0.10340 ( 28.08 566 16.60 108.30 858.1 0.08455 0.10230 ( 20.60 0.11780 567 29.33 140.10 1265.0 0.27700 ( 24.54 0.05263 0.04362 568 7.76 47.92 181.0 569 rows × 10 columns Os dados qualitativos do diagnostico, ou seja, o y: In [4]: df2.iloc[:,1] Out[4]: 1 3 1 564 565 1 566 1 567 1 568 Name: diagnosis, Length: 569, dtype: int64 Treino e Teste Pegando o X e o y e separando 70 dos dados para treinar os modelos, enquanto os outros 30 estão destinados para testar o modelo e nós dizer uma métrica. from sklearn import model selection X train, X test, y train, y test = model selection.train test split(df2.iloc[:,2:], df2.iloc[:,1], test size=0.3, random state=0) resultados = pd.DataFrame(columns=(['modelo', 'acuracia', 'recall'])) **Testar Modelos** Para efeito de organização e tempo de execução para organizar este arquivo vou deixar o Grid Search apenas para o modelo KNN, porém foi usado para encontrar os melhores parametros para todos os modelos. **Naive Bayes** from sklearn.naive bayes import GaussianNB gnb = GaussianNB().fit(X train, y train) resultados = evalu(gnb, X test, y test, 'Naive Bayes', resultados) Resultados: precision recall f1-score support Benigno 0.92 0.93 0.92 108 0.87 0.86 0.86 Maligno 63 accuracy 0.90 171 macro avg 0.89 0.89 0.89 171 weighted avg 0.90 0.90 0.90 171 ROC: 0.8915343915343915 Acurácia: 0.9005847953216374 Precisão : 0.7991753577492118 Recall : 0.8571428571428571 f1 : 0.864 Matriz de Confusão 100 80 100 8 Benigno 60 40 9 Maligno 20 Maligno Benigno Predição KNN from sklearn.neighbors import KNeighborsClassifier E = KNeighborsClassifier() P = {'n neighbors':[5,10,15,20,25,35,45,50,55,65,75],'weights':['uniform','distance'] k = model selection.GridSearchCV(estimator=E, param grid=P, scoring='recall', n jobs=-1).fit(X train, y train) k.best estimator Out[9]: KNeighborsClassifier(weights='distance') knn = KNeighborsClassifier(weights='distance').fit(X train, y train) resultados = evalu(knn, X test, y test, 'KNN', resultados) Resultados: precision recall f1-score support Benigno 0.90 0.94 0.92 108 Maligno 0.90 0.83 0.86 63 0.90 accuracy 171 macro avg 0.90 0.88 0.89 171 weighted avg 0.90 0.90 0.90 171 ROC: 0.8849206349206349 Acurácia: 0.9005847953216374 Precisão: 0.8043384322876156 Recall: 0.8253968253968254 f1: 0.8595041322314049 Matriz de Confusão 100 80 102 6 Benigno Verdadeiro 60 40 11 Maligno 20 Maligno Benigno Predição **Logistic Regression** from sklearn.linear\_model import LogisticRegression LR = LogisticRegression(C=7.70299999999999, multi class='ovr', penalty='11', solver='liblinear', tol=0.0033333334000000005).fit(X train, y train resultados= evalu(LR, X test, y test, 'Reg. Logistica', resultados) Resultados: precision recall f1-score 0.97 Benigno 0.95 0.96 Maligno 0.95 0.90 0.93 63 accuracy 0.95 171 0.95 0.94 0.94 171 macro avg weighted avg 0.95 0.95 0.95 171 ROC: 0.9384920634920634 Acurácia: 0.9473684210526315 Precisão: 0.8946115288220551 Recall : 0.9047619047619048 f1 : 0.9268292682926829 Matriz de Confusão 100 105 3 Benigno 80 Verdadeiro 60 40 6 Maligno 20 Maligno Benigno Predição Arvore de Decisão from sklearn.tree import DecisionTreeClassifier dtc = DecisionTreeClassifier(random state=42, splitter='random' ).fit(X train, y train) resultados = evalu(dtc, X test, y test, 'Arvore de Decisão', resultados) Resultados: precision recall f1-score support 0.91 0.96 0.94 Maligno 0.93 0.84 0.88 63 0.92 171 accuracy 0.92 0.90 0.91 macro avg 0.92 171 weighted avg 0.92 0.92 ROC: 0.9021164021164022 Acurácia : 0.9181286549707602 Precisão : 0.8407128933444723 Recall : 0.8412698412698413 f1 : 0.88333333333333334 Matriz de Confusão 100 104 80 Benigno Verdadeiro 60 40 10 Maligno 20 Maligno Benigno Predição Random Forest from sklearn.ensemble import RandomForestClassifier rf = RandomForestClassifier(criterion='entropy', max\_depth=6, n\_estimators=200, random state=42).fit(X train, y train) resultados = evalu(rf, X test, y test, 'Random Forest', resultados) Resultados: precision recall f1-score support 0.96 Benigno 0.94 0.95 108 Maligno 0.89 0.94 0.91 63 0.94 accuracy 171 0.93 0.93 0.94 171 macro avg weighted avg 0.94 0.94 0.94 171 ROC: 0.9358465608465609 Acurácia : 0.935672514619883 Precisão : 0.8605731500468342 Recall: 0.9365079365079365 f1 : 0.9147286821705426 Matriz de Confusão 100 80 101 Benigno 60 Verdadeiro 40 Maligno 20 Benigno Maligno Predição **SGDC Classifier** In [14]: from sklearn.linear model import SGDClassifier sgdc = SGDClassifier(alpha=1, loss='squared hinge', max iter=20, penalty='11').fit(X to the squared hinge') resultados = evalu(sgdc, X test, y test, 'SGDC', resultados) C:\Users\carlo\anaconda3\lib\site-packages\sklearn\linear\_model\\_stochastic\_gradient.p y:574: ConvergenceWarning: Maximum number of iteration reached before convergence. Con sider increasing max iter to improve the fit. warnings.warn("Maximum number of iteration reached before " Resultados: precision recall f1-score support 0.90 0.84 0.87 108 Benigno 0.84 Maligno 0.76 0.80 63 accuracy 0.84 171 0.83 0.84 0.83 171 macro avq 0.84 0.84 171 weighted avg ROC: 0.8419312169312169 Acurácia : 0.8421052631578947 Precisão : 0.6954409834109082 Recall : 0.8412698412698413 f1 : 0.7969924812030075 Matriz de Confusão 80 91 Benigno 70 60 Verdadeiro 50 10 Maligno 30 20 10 Benigno Maligno Predição Testando algo diferente... Aqui, vamos testar algumas ideias e utilizar os modelos de Naive Bayes (porque é simples) e Random Forest (porque foi a que teve o melhor desempenho associando acurácia e recall) para ver se elas são mais eficientes que os modelos que temos até agora. Ideia 1 Raio, Area e Perimetro são altamente correlacionados, como eu disse no primeiro trabalho, acho que Area e Perimetro são derivados do Raio, então vou testar alguns modelos sem essas variavéis. df3 = df2.drop(['perimeter\_mean', 'area\_mean'], axis=1) X\_train1, X\_test1, y\_train1, y\_test1 = model\_selection.train\_test\_split(df3.iloc[:,2: df3.iloc[:,1] test\_size=0.3, random\_state=0 gnb1 = GaussianNB().fit(X\_train1, y\_train1) resultados = evalu(gnb1, X\_test1, y\_test1, 'NB - 1', resultados) Resultados: precision recall f1-score support 0.90 0.92 0.91 108 Beniano Maligno 0.83 0.87 0.85 63 0.89 171 accuracy 0.88 0.89 macro avg 0.88 171 weighted avg 0.89 0.89 0.89 171 ROC: 0.8855820105820106 Precisão: 0.7742968532442217 Recall : 0.873015873015873 f1 : 0.8527131782945736 Matriz de Confusão 80 97 11 Benigno 70 60 Verdadeiro 50 40 Maligno 30 20 Benigno Maligno Predição rf1 = RandomForestClassifier(max\_depth=5, n\_estimators=500, random\_state=42).fit(X\_transfer) resultados = evalu(rf1, X\_test1, y\_test1, 'RF - 1', resultados) Resultados: recall f1-score precision support 0.94 0.94 0.94 108 Benigno 0.90 0.89 0.90 Maligno 63 0.92 171 accuracy 0.92 0.92 0.92 macro avg 171 weighted avg 0.92 0.92 0.92 171 Acurácia: 0.9239766081871345 Precisão: 0.8438030560271647 f1 : 0.8959999999999999 Matriz de Confusão 100 102 6 Benigno · Verdadeiro 60 40 Maligno 20 Benigno Maligno Predição Ideia 2 E se em vez de tirar os correlacionados entre si, retirar os que tem menor correlação com a resposta. Nessa hipotese, eu retirei as variaveis que tinha uma correlação menor que 0.4 com a variavel a ser predita. df4 = df2.drop(['fractal\_dimension\_mean','symmetry\_mean','smoothness\_mean'],axis=1) X\_train2, X\_test2, y\_train2, y\_test2 = model\_selection.train\_test\_split(df4.iloc[:,2: test size=0.3 gnb2 = GaussianNB().fit(X\_train2, y\_train2) resultados = evalu(gnb2, X\_test2, y\_test2, 'NB - 2', resultados) Resultados: recall f1-score precision support 0.92 0.93 0.92 Benigno 108 Maligno 0.86 0.87 0.86 63 accuracy 0.90 171 0.89 0.89 0.89 171 macro avg weighted avg 0.90 0.90 171 ROC: 0.8915343915343915 Acurácia: 0.9005847953216374 Precisão: 0.7991753577492118 Recall : 0.8571428571428571 f1 : 0.864 Matriz de Confusão 100 100 Benigno Verdadeiro 60 40 Maligno 20 Maligno Benigno Predição rf2 = RandomForestClassifier(max depth=4, n estimators=500, random state=42).fit(X tra resultados = evalu(rf2, X\_test2, y\_test2, 'RF - 2', resultados) Resultados: recall f1-score precision support 0.96 0.94 0.95 108 0.89 0.94 Maligno 0.91 63 0.94 171 accuracy macro avg 0.93 0.94 0.93 171 0.94 171 weighted avg 0.94 0.94 ROC: 0.9358465608465609 Acurácia: 0.935672514619883 Precisão: 0.8605731500468342 Recall: 0.9365079365079365 f1: 0.9147286821705426 Matriz de Confusão 100 80 101 7 Benigno Verdadeiro 60 40 Maligno 20 Benigno Maligno Predição Ideia 3 E se juntar a Ideia 1 com a Ideia 2 df5 = df2.drop(['perimeter\_mean','area\_mean','fractal\_dimension\_mean','symmetry\_mean' X\_train3, X\_test3, y\_train3, y\_test3 = model\_selection.train\_test\_split(df5.iloc[:,2: test size=0.3 gnb3 = GaussianNB().fit(X\_train3, y\_train3) resultados = evalu(gnb3, X\_test3, y\_test3, 'NB - 3', resultados) Resultados: recall precision f1-score support 0.93 0.93 Benigno 0.93 108 Maligno 0.87 0.87 0.87 63 accuracy 0.91 171 0.90 0.90 0.90 171 macro avg weighted avg 0.91 0.91 0.91 171 ROC: 0.8994708994708994 Acurácia: 0.9064327485380117 Precisão: 0.8089403402686611 Recall: 0.873015873015873 f1 : 0.8730158730158731 Matriz de Confusão 100 100 Benigno Verdadeiro 60 40 Maligno 20 Maligno Benigno Predição rf3 = RandomForestClassifier(max\_depth=4, n\_estimators=500, random\_state=42).fit(X\_transfer) resultados = evalu(rf3, X\_test3, y\_test3, 'RF - 3', resultados) Resultados: precision recall f1-score support 0.94 0.94 0.94 108 0.89 0.90 Maligno 0.90 63 0.92 171 accuracy 0.92 0.92 0.92 171 macro avg 171 weighted avg 0.92 0.92 0.92 ROC: 0.919973544973545 Acurácia: 0.9239766081871345 Precisão: 0.840891290726817 Recall : 0.9047619047619048 f1: 0.8976377952755906 Matriz de Confusão 100 80 7 101 Benigno Verdadeiro 60 40 Maligno 20 Maligno Benigno Predição In [24]: resultados Out[24]: modelo recall acuracia 0.900585 0.857143 0 Naive Bayes KNN 0.900585 0.825397 1 2 Reg. Logistica 0.947368 0.904762 Arvore de Decisão 0.918129 0.841270 4 Random Forest 0.935673 0.936508 5 SGDC 0.842105 0.841270 6 NB - 1 0.888889 0.873016 RF - 1 0.923977 0.888889 7 NB - 2 0.900585 0.857143 8 RF - 2 0.935673 0.936508 9 10 NB - 3 0.906433 0.873016 11 RF - 3 0.923977 0.904762 Conclusão Acredito que os modelos foram satisfatórios, e eu diria que a Random Forest, foi o modelo que melhor previu os dados. Isso porque mesmo não tendo a melhor acurácia, teve o melhor recall e teve uma acurácia muito próxima do modelo que teve a maior. Além disso, eu queria resaltar um importante ponto, o modelo deu certo tem uma taxa de acerto bem grande ao predizer se o tumor é maligno ou benigno, porém ele é inútil para pessoas "normais", para ele funcionar relativamente bem, as pessoas que vão sofrer a avaliação do modelo tem estar minimamente padronizadas com os dados que foram dados para modelar a questão. Como os dados são de exames de cancer de mama com pessoas que já tem o cancer diagnosticado, não é conveniente chegar e em pessoas que não estão sentido nada e fazer um teste com elas, o modelo com certeza perderia a eficiencia, pois os pre reqisitos na hora de colher os dados não estaram sendo compridos. **Outros:** Estão no GITHUB: O arquivo funcoes contem as funções que foram usadas para avaliar o modelo. O notebook com todos os passos, (aqui eu resumi um pouco de código para não poluir tanto) #FIM