import numpy as np # linear algebra import pandas as pd # data processing, CSV file I/O (e.g. pd.read csv) import seaborn as sns import matplotlib.pyplot as plt from sklearn.metrics import mean absolute error, accuracy score, recall score from sklearn.model selection import train test split from sklearn.preprocessing import StandardScaler from sklearn.preprocessing import MinMaxScaler from sklearn.preprocessing import StandardScaler from tensorflow import keras from tensorflow.keras import layers from sklearn.preprocessing import OrdinalEncoder from sklearn.ensemble import RandomForestClassifier #Creating a dataframe from the data heart data = pd.read csv(r"C:\Users\alexa\OneDrive\Documents\heart\heart.csv") #Looking at the different types of data stored in each column to determine if there are any categoricals that heart data.dtypes Out[13]: Age int64 object ChestPainType object RestingBP int64 Cholesterol int64 FastingBS int64 object RestingECG MaxHR int64 ExerciseAngina object Oldpeak float64 ST Slope object int64 HeartDisease dtype: object In [14]: #A quick view of the data heart data.head() Out[14]: ChestPainType RestingBP Cholesterol FastingBS RestingECG MaxHR ExerciseAngina Oldpeak ST_Slope HeartDisease Sex Age 0 Normal 0 0 40 M ATA 140 289 172 Ν 0.0 Up 1 49 NAP 160 180 0 Normal 156 1.0 Flat 0 0.0 0 2 37 M ATA 130 283 ST 98 Ν Up 3 48 ASY 138 214 0 Normal 108 1.5 Flat NAP 150 195 0 0.0 0 54 M Normal 122 Ν Up #Making a histogram that shows the distribution of the data heart data.hist(figsize = (10,15)) Out[15]: array([[<AxesSubplot:title={'center':'Age'}>, <AxesSubplot:title={'center':'RestingBP'}>, <AxesSubplot:title={'center':'Cholesterol'}>], [<AxesSubplot:title={'center':'FastingBS'}>, <AxesSubplot:title={'center':'MaxHR'}>, <AxesSubplot:title={'center':'Oldpeak'}>], [<AxesSubplot:title={'center':'HeartDisease'}>, <AxesSubplot:>, <AxesSubplot:>]], dtype=object) RestingBP Cholesterol Age 200 300 400 175 250 150 300 200 125 100 200 75 100 50 100 50 25 0 0 0 40 60 50 100 150 200 0 200 400 600 FastingBS MaxHR Oldpeak 200 700 350 175 600 300 150 500 250 125 400 200 100 300 150 75 200 100 50 100 50 25 0.00 0.25 0.50 0.75 1.00 100 150 200 -2 HeartDisease 500 400 300 200 100 0.25 0.50 0.75 1.00 0.00 In [16]: #Setting up features to be plotted in the next code block discrette feature = [i for i in heart_data.columns if heart_data[i].nunique() < 10]</pre> continuous_feature = [i for i in heart_data.columns if heart_data[i].nunique() > 10] print(f'discrette feature: {discrette feature}') print(f'continuous feature: {continuous feature}') discrette feature: ['Sex', 'ChestPainType', 'FastingBS', 'RestingECG', 'ExerciseAngina', 'ST Slope', 'HeartDise continuous feature: ['Age', 'RestingBP', 'Cholesterol', 'MaxHR', 'Oldpeak'] #Potting continuous features of the dataset as histograms showing distributions for both those who had heart fe plt.style.use('fivethirtyeight') i = 1plt.figure(figsize = (18,16))for feature in continuous feature: plt.subplot(3, 2, i) sns.histplot(x=heart data[feature],kde=True,bins = 50, hue = heart data.HeartDisease) plt.xlabel(feature, size = 12) plt.ylabel("Density", size = 12) i += 1 plt.show() HeartDisease HeartDisease 80 0 25 70 1 1 60 20 Density 40 Density 15 10 30 20 5 10 0 0 60 40 50 70 25 50 75 100 125 150 175 200 RestingBP Age HeartDisease HeartDisease 140 40 0 1 1 120 30 100 Density 00 Density 80 60 40 10 20 0 0 100 200 300 400 500 100 160 180 120 140 Cholesterol MaxHR 250 HeartDisease 0 ___1 200 Density 150 100 50 0 0 6 Oldpeak #Creating a visualization of the correlation between variables in the data set plt.figure(figsize=(12, 8)) sns.heatmap(heart data.corr(), annot=True) plt.xticks(rotation=45); 1.0 Age 1 0.25 -0.095 0.2 -0.380.26 0.28 0.8 RestingBP 0.25 0.1 0.07 -0.110.16 0.11 1 0.6 Cholesterol -0.095 0.1 1 -0.260.24 0.05 -0.230.4 0.053 FastingBS 0.2 0.07 -0.261 -0.130.27 0.2 -0.38 0.24 -0.4 MaxHR -0.11-0.13-0.16 1 0.0 Oldpeak 0.26 0.16 0.05 0.053 -0.161 0.4 -0.2 HeartDisease -0.4From this visualization we can see that the variables most strongly correlated to heart disease are old peak and max heart rate. #Setting the target for prediction and features that will be used to train the network y = heart data['HeartDisease'] X = heart data.drop(['HeartDisease'], axis=1) train_X, test_X, train_y, test_y = train_test_split(X, y, train_size=0.8, test_size=0.2) In [19]: #Using ordinal encoder to convert categoricals to numerical values for network s = (heart data.dtypes == 'object') object_cols = list(s[s].index) label_train_X = train_X.copy() label_test_X = test_X.copy() ordinal encoder = OrdinalEncoder() label_train_X[object_cols] = ordinal_encoder.fit_transform(train_X[object_cols]) label_test_X[object_cols] = ordinal_encoder.transform(label_test_X[object_cols]) #Checking that encoding worked label train X.info() <class 'pandas.core.frame.DataFrame'> Int64Index: 734 entries, 498 to 790 Data columns (total 11 columns): Non-Null Count Dtype # Column --- ---------0 Age 734 non-null int64 1 Sex 734 non-null float64 2 ChestPainType 734 non-null float64
3 RestingBP 734 non-null int64
4 Cholesterol 734 non-null int64
5 FastingBS 734 non-null int64
6 RestingECG 734 non-null float64
7 MaxHR 734 non-null int64 8 ExerciseAngina 734 non-null float64 9 Oldpeak 734 non-null float64 10 ST_Slope 734 non-null float64 10 ST_Slope dtypes: float64(6), int64(5) memory usage: 68.8 KB #Scaling the data for use in the network scaler = StandardScaler() label_train_X = scaler.fit_transform(label_train_X) label_test_X = scaler.transform(label_test_X) #Building the structure of the model model = keras.Sequential([layers.Dense(32, activation = 'swish', input shape = [11]), layers.Dropout(.3), layers.Dense(32, activation = 'swish'), layers.Dropout(.3), layers.Dense(32, activation = 'swish'), layers.Dropout(.3), layers.Dense(1, activation = 'sigmoid'),]) #Adding methods for optimizing and evaluating the model model.compile(optimizer = 'adam', loss = 'binary_crossentropy', metrics = ['binary_accuracy'] In [24]: #Adding a callback that will stop the model at best values. early stopping = keras.callbacks.EarlyStopping(patience=5, min delta=0.001, restore best weights=True, #Training the model history = model.fit(label train X, train y, validation data=(label test X, test y), batch size=576, epochs=200, callbacks=[early stopping], #Plotting the learning curves history df = pd.DataFrame(history.history) history df.loc[:, ['loss', 'val loss']].plot(title="Cross-entropy") history df.loc[:, ['binary accuracy', 'val binary accuracy']].plot(title="Accuracy") Epoch 1/200 2/2 [===========] - 0s 68ms/step - loss: 0.7262 - binary accuracy: 0.4237 - val loss: 0.6953 - val_binary_accuracy: 0.5326 Epoch 2/200 2/2 [==============] - 0s 6ms/step - loss: 0.7075 - binary_accuracy: 0.4837 - val_loss: 0.6829 - val_binary_accuracy: 0.6250 Epoch 3/200 2/2 [=============] - 0s 6ms/step - loss: 0.6853 - binary accuracy: 0.5463 - val loss: 0.6710 - val binary accuracy: 0.7228 Epoch 4/200 - val_binary_accuracy: 0.7663 Epoch 5/200 - val_binary_accuracy: 0.7935 Epoch 6/200 - val_binary_accuracy: 0.7989 Epoch 7/200 - val binary_accuracy: 0.8098 Epoch 8/200 - val_binary_accuracy: 0.8152 Epoch 9/200 - val_binary accuracy: 0.8098 Epoch 10/200 - val_binary_accuracy: 0.8098 Epoch 11/200 - val_binary_accuracy: 0.8152 Epoch 12/200 - val binary accuracy: 0.8152 Epoch 13/200 2/2 [========== 0.5745 - binary_accuracy: 0.7779 - val_loss: 0.5521 - val_binary_accuracy: 0.8152 Epoch 14/200 2/2 [============================] - 0s 6ms/step - loss: 0.5584 - binary_accuracy: 0.7766 - val_loss: 0.5393 - val binary accuracy: 0.8152 Epoch 15/200 - val_binary_accuracy: 0.8207 Epoch 16/200 - val binary accuracy: 0.8207 Epoch 17/200 - val_binary_accuracy: 0.8207 Epoch 18/200 - val_binary_accuracy: 0.8207 Epoch 19/200 2/2 [====== ======] - 0s 6ms/step - loss: 0.5061 - binary_accuracy: 0.8202 - val_loss: 0.4760 - val binary accuracy: 0.8207 Epoch 20/200 2/2 [=============] - 0s 6ms/step - loss: 0.4844 - binary accuracy: 0.8202 - val loss: 0.4647 - val binary accuracy: 0.8207 Epoch 21/200 - val_binary_accuracy: 0.8261 Epoch 22/200 - val_binary_accuracy: 0.8315 Epoch 23/200 - val binary_accuracy: 0.8261 Epoch 24/200 - val binary accuracy: 0.8261 Epoch 25/200 - val_binary accuracy: 0.8261 Epoch 26/200 - val_binary_accuracy: 0.8261 Epoch 27/200 - val_binary_accuracy: 0.8261 Epoch 28/200 - val binary accuracy: 0.8261 Epoch 29/200 - val_binary_accuracy: 0.8315 Epoch 30/200 - val_binary_accuracy: 0.8370 Epoch 31/200 2/2 [=============] - 0s 6ms/step - loss: 0.4125 - binary accuracy: 0.8379 - val loss: 0.3974 - val binary_accuracy: 0.8370 Epoch 32/200 - val binary accuracy: 0.8478 Epoch 33/200 - val_binary_accuracy: 0.8478 Epoch 34/200 - val_binary_accuracy: 0.8533 Epoch 35/200 - val_binary_accuracy: 0.8533 Epoch 36/200 - val_binary accuracy: 0.8533 Epoch 37/200 - val_binary_accuracy: 0.8533 Epoch 38/200 - val_binary_accuracy: 0.8533 Epoch 39/200 - val binary_accuracy: 0.8533 Epoch 40/200 - val binary accuracy: 0.8533 Epoch 41/200 - val_binary_accuracy: 0.8533 Epoch 42/200 - val_binary_accuracy: 0.8478 Epoch 43/200 2/2 [=============] - 0s 7ms/step - loss: 0.3905 - binary accuracy: 0.8488 - val loss: 0.3858 - val_binary_accuracy: 0.8424 Epoch 44/200 2/2 [=============] - 0s 10ms/step - loss: 0.3724 - binary_accuracy: 0.8420 - val_loss: 0.3853 - val binary accuracy: 0.8424 Epoch 45/200 ======] - Os 8ms/step - loss: 0.3968 - binary_accuracy: 0.8420 - val_loss: 0.3847 2/2 [== - val_binary_accuracy: 0.8424 Epoch 46/200 2/2 [============================] - 0s 7ms/step - loss: 0.3848 - binary_accuracy: 0.8433 - val_loss: 0.3837 - val_binary_accuracy: 0.8424 Epoch 47/200 =======] - 0s 8ms/step - loss: 0.3853 - binary accuracy: 0.8338 - val loss: 0.3826 2/2 [====== - val binary_accuracy: 0.8424 Epoch 48/200 - val binary accuracy: 0.8478 Epoch 49/200 - val_binary_accuracy: 0.8478 Epoch 50/200 - val_binary_accuracy: 0.8478 Epoch 51/200 2/2 [====== - val_binary_accuracy: 0.8533 Epoch 52/200 - val binary accuracy: 0.8533 Epoch 53/200 - val_binary_accuracy: 0.8478 Epoch 54/200 - val_binary_accuracy: 0.8478 Epoch 55/200 =======] - 0s 6ms/step - loss: 0.3859 - binary accuracy: 0.8447 - val loss: 0.3741 2/2 [====== - val binary_accuracy: 0.8533 Epoch 56/200 - val binary accuracy: 0.8533 Epoch 57/200 =========] - 0s 5ms/step - loss: 0.3516 - binary_accuracy: 0.8488 - val_loss: 0.3724 2/2 [======= - val_binary_accuracy: 0.8533 Epoch 58/200 - val_binary_accuracy: 0.8533 Epoch 59/200 - val_binary_accuracy: 0.8533 Epoch 60/200 - val binary accuracy: 0.8533 Epoch 61/200 - val_binary_accuracy: 0.8533 Epoch 62/200 - val_binary_accuracy: 0.8533 Epoch 63/200 =======] - 0s 6ms/step - loss: 0.3571 - binary_accuracy: 0.8488 - val_loss: 0.3688 2/2 [========== - val_binary_accuracy: 0.8533 Epoch 64/200 - val binary accuracy: 0.8533 Epoch 65/200 - val_binary_accuracy: 0.8533 Epoch 66/200 - val_binary_accuracy: 0.8533 Epoch 67/200 2/2 [====== - val_binary_accuracy: 0.8533 Out[24]: <AxesSubplot:title={'center':'Accuracy'}> Cross-entropy loss 0.7 val_loss 0.6 0.5 0.4 30 40 50 0 10 20 60 Accuracy 0.8 0.7 0.6 0.5 binary accuracy val_binary_accuracy 10 40 60 20 30 50 The network predicts for heart failure with an accuracy of around 85% #Creating a random forest classifier to compare to the neural network RFC = RandomForestClassifier(random state = 1) RFC.fit(label train X, train y) RFC predictions = RFC.predict(label test X) RFC_acc = accuracy score(test y, RFC predictions) RFC_recall = recall_score(test_y, RFC_predictions) print(RFC acc) print(RFC recall) 0.875 0.9207920792079208 The random forest classifier appears to have a higher accuracy than the network at 87.5% accuracy. Although the difference is small the

random forrest was easier to setup and required less computation time as well which are both important factors.