RandomForestClassifier\_v004

February 1, 2020

© Thomas Robert Holy 2019 Version 1.0 Visit me on GitHub: https://github.com/trh0ly Kaggle Link: https://www.kaggle.com/c/dda-p2/leaderboard

1 RandomForestClassifier

from sklearn.pipeline import Pipeline

from sklearn.metrics import auc

import matplotlib.pyplot as plt

from sklearn.svm import SVC

# Werte von c betrachtet werden.

#-----

if c != None:

if c == None:

import datetime as dt

1.2 Hilfsfunktionen

#-----# Argumente:

from sklearn.linear\_model import LogisticRegression from sklearn.preprocessing import StandardScaler from sklearn.preprocessing import RobustScaler from sklearn.preprocessing import MinMaxScaler from sklearn.metrics import confusion\_matrix from sklearn.metrics import classification\_report from sklearn.metrics import roc\_curve, roc\_auc\_score

from sklearn.ensemble import RandomForestClassifier from sklearn.neighbors import KNeighborsClassifier from sklearn.tree import DecisionTreeClassifier from sklearn.preprocessing import PolynomialFeatures from sklearn.model\_selection import GridSearchCV from sklearn.model\_selection import StratifiedKFold from sklearn.model\_selection import cross\_val\_score from sklearn.model\_selection import train\_test\_split

from IPython.core.display import display, HTML from scipy.spatial.distance import euclidean

1.2.1 Funktion zur Betrachtung der Konfusinsmatrix

from sklearn.metrics.pairwise import manhattan\_distances

# -  $y_true$ : Zum DataFrame X gehörige Werte der Zielgröße

[2]: # Definition einer Funktion, welche eine Konfusionsmatrix und einen Klassifikationsreport # zurückgibt. Die Konfusionsmatrix kann, wenn ein Wert für c gegeben ist, für beliebige

# - X: DataFrame auf welchem die Prognose durchgefürt werden soll (ohne die Zielgröße)

# ---> Wenn None, dann wird die Konfusionsmatrix ohne die Einbeziehung von c bestimmt # ---> Wenn != None, dann wird die Konfusionsmatrix in Abhänqiqkeit von c bestimmt

# - class\_names: Bezeichnung für die Spalten des Dataframes (default=['0', '1'], mit 0 = negativ und 1 = positiv)

# - model: Modell auf Basis dessen die Konfusionsmatrix berechnet werden soll

def get\_confusion\_matrix(X, y\_true, model, class\_names=['0', '1'], c=None):

# Überführung in einen DataFrame für eine bessere Übersichtlichkeit

df\_conf\_mat = pd.DataFrame(conf\_mat, index=df\_index, columns=df\_cols)

return df\_conf\_mat, classification\_report(y\_true, y\_pred)

df\_index = pd.MultiIndex.from\_tuples([('Wahrer Wert', cn) for cn in class\_names])

[3]: | # Definition einer Funktion, welche auf Basis eines gegeben Modells und zweier zusammengehöriger

df\_cols = pd.MultiIndex.from\_tuples([('Prognose des Modells', cn) for cn in class\_names])

# Vorgelagerte Berechnung falls ein Wert für c gegeben ist # und die Konfusionsmatrix für ein gegebenes c anpasst

pred\_probability = model.predict\_proba(X) pred\_probability = pred\_probability >= c y\_pred = pred\_probability[:, 1].astype(int)

# Wenn kein Wert für c gegeben, dann führe Prognose

# lediglich auf Basis des Modells durch

conf\_mat = confusion\_matrix(y\_true, y\_pred)

y\_pred = model.predict(X)

#-----# Berechnet die Konfusionsmatrix

#-----

1.2.2 Funktion zur Betrachtung der ROC-Kurve

# In .csv speichern, wenn save=True

#-----

if manu\_name == False:

# Standardnamen wählen, wenn manu\_name == False

submission.to\_csv('./predicted\_values.csv', index=False)

# Standardnamen mit timestamp kombinieren, wenn manu\_name == True

if save == True:

0.0

1.1 Package Import

import pandas as pd

[1]: import numpy as np

if manu\_name == True: import datetime now = datetime.datetime.now() name = now.strftime(' $^{\prime\prime}_{Y}$ - $^{\prime\prime}_{m}$ - $^{\prime\prime}_{d}$ T $^{\prime\prime}_{H}$  $^{\prime\prime}_{M}$  $^{\prime\prime}_{S}$ ') + ('- $^{\prime\prime}_{0}$ 2d'  $^{\prime\prime}_{m}$  (now.microsecond / 10000)) submission.to\_csv('./predicted\_values\_' + str(name) + '.csv', index=False) return submission.head(), submission.loc[submission['Fehlerhaft'] == 1] 1.2.4 Funktion zum Filtern von Quantilen [5]: | # Definition einer Funktion, welche einen gegeben DataFrame # um untere und obere Quantile beschneiden kann #-----# Argumente: # - orignal\_df: DataFrame welcher bearbeitet werden soll # - quantile\_low: Unteres Quantil bis zu welchem orignal\_df beschnitten werden soll # - quantile\_high: Oberes Quantil welchem orignal\_df beschnitten werden soll # - colum\_to\_drop: Spalte des orignal\_df, welche während des Vorgangs gedroppt werden soll def filter\_my\_df(orignal\_df, quantile\_low, quantile\_high, colum\_to\_drop): # Spalte "colum\_to\_drop" aus dem Datensatz entfernen df\_filtered = orignal\_df.loc[:, orignal\_df.columns != colum\_to\_drop] # Quantil-DataFrame erzeugen quant\_df = df\_filtered.quantile([quantile\_low, quantile\_high]) # Quantil-DataFrame auf orignal\_df anweden df\_filtered = df\_filtered.apply(lambda x: x[(x>quant\_df.loc[quantile\_low,x.name]) & (x < quant\_df.loc[quantile\_high,x.name])],</pre> #----# Spalte "Fehlerhaft" dem gefiltertem DataFrame wieder anfügen df\_filtered = pd.concat([orignal\_df.loc[:,colum\_to\_drop], df\_filtered], axis=1) # Aus Beschneidung resultierende NaN-Werte bereinigen df\_filtered.dropna(inplace=True) return df\_filtered 1.3 Datensatz einlesen (bereinigigen) und betrachten 1.3.1 Datensatz einlesen Г61: | #----# Datensatz einlesen data = pd.read\_csv('train.csv', index\_col=0) 1.3.2 Optionale Datensatzbereinigung [7]: """ # Datensatz unterteilen df\_fehlerfrei = data.loc[data['Fehlerhaft'] == 0] df\_fehlerhaft = data.loc[data['Fehlerhaft'] == 1] #-----# Fehlerfreie Stückgüter colum\_to\_drop = 'Fehlerhaft' orignal\_df = df\_fehlerfrei low = .0 # Unteres Quantil high = .99 # Oberes Quantil  $df\_fehlerfrei\_filtered = filter\_my\_df(df\_fehlerfrei, \ low, \ high, \ colum\_to\_drop)$ # Fehlerhafte Stückgüter colum\_to\_drop = 'Fehlerhaft'  $orignal\_df = df\_fehlerhaft$ low = .018333 # Unteres Quantil high = 1. # Oberes Quantil  $df_fehlerhaft_filtered = filter_my_df(df_fehlerhaft, low, high, colum_to_drop)$ #-----# Teil-DataFrames zusammenführen  $data\_filtered = pd.concat([df\_fehlerhaft\_filtered, df\_fehlerfrei\_filtered], sort=False)$ [7]: "\n#-----\n# Fehlerfreie Stückgüter\ncolum\_to\_drop = 'Fehlerhaft'\norignal\_df = df\_fehlerfrei\nlow = .0 # Unteres Quantil \nhigh = .99 # Oberes Quantil\ndf\_fehlerfrei\_filtered = filter\_my\_df(df\_fehlerfrei, low, high, colum\_to\_drop)\n\n#----\n# Fehlerhafte .018333 # Unteres Quantil \nhigh = 1. # Oberes Quantil \ndf\_fehlerhaft\_filtered = filter\_my\_df(df\_fehlerhaft, low, high, colum\_to\_drop)\n\n#-----\n# Teil-DataFrames zusammenführen\ndata\_filtered = pd.concat([df\_fehlerhaft\_filtered, df\_fehlerfrei\_filtered], sort=False)\n" 1.3.3 Beschreibung der separierten Datensätze (Betrachtung Min-/ Maximum und Qunatile) [8]: """  $df_-fehlerfrei.describe()$ [8]: '\ndf\_fehlerfrei.describe()\n'  $df_{-}fehlerhaft.describe()$ [9]: '\ndf\_fehlerhaft.describe()\n' [10]: data\_new = data #\_filtered data\_new['Fehlerhaft'].value\_counts() 20208 [10]: 0 284 Name: Fehlerhaft, dtype: int64 1.3.4 Betrachtung Korrelationsmatrix [11]: data\_new = data #\_filtered # Für schnellere Laufzeit und mehr Übersicht in den Plots: Stichprobe der Daten abbilden data\_sample = data\_new.sample(2000, random\_state=28) # random\_state sorgt für reproduzierbare Stichprobe, sodass die Stichprobe für uns alle identisch ist \_ = pd.plotting.scatter\_matrix(data\_sample, c=data\_sample['Fehlerhaft'], cmap='seismic', figsize=(16, 20)) 40 Sensor\_1 0 -20 7.5 Sensor 2 2.5

-2.57.5 Sensor -2.5Sensor\_5 -15 10 Sensor\_6 -10Sensor 7 -20 Sensor 8 -151.0 Fehlerhaft 0.6 0.4 0.2 10 Sensor\_2 Sensor\_4 Sensor\_1 Sensor\_3 Fehlerhaft Sensor\_5 Sensor\_6 Sensor\_7 Sensor\_8 1.3.5 Dateinsatz in Traings- und Validierungsteil splitten [12]: X = data\_new.drop('Fehlerhaft', axis=1) y = data\_new['Fehlerhaft'] X\_train, X\_validierung, y\_train, y\_validierung = train\_test\_split(X, y, test\_size=0.2, random\_state=2121) 1.4 Modell aufstellen [13]: # Definition einer Funktion, welche eine Gittersuche mit einem RandomForestClassifier durchführt # und nach einer 5-fach Kreuzvalidierung das beste Modell zurückgibt # Argumente: # - i: Fügt X^i der Featurematrix hinzu # - X: DataFrame auf welchem die Prognose durchgefürt werden soll (ohne die Zielgröße) # - y\_true: Zum DataFrame X gehörige Werte der Zielgröße # - my\_scaler: Zu verwendender Scaler; per default MinMaxScaler; weitere Scaler: RobustScaler, Standardscaler # - max\_features: Anzahl der Features die einbezogen werden sollen (default=['auto']=sqrt(n\_features)) # - n\_estimators: Anzahl der "Bäume" im "Wald" # - jobs: Anzahl der Threads die für den Durchlauf zur Verfügung stehen # - gs\_scoring: Scoring Verfahren im Rahmen der GridSearch # - folts: Komplexität der Kreuzvalidierung #----def rndfrst\_1(i, X, y\_true, my\_scaler=MinMaxScaler, jobs=-3, gs\_scoring='f1', folts=5, n\_estimators=list(range(10, 200 + 1,10)), max\_features=None): #-----# Pipeline erzeugen prediction\_pipe = Pipeline([('scaler', my\_scaler()), ('add\_x\_square', PolynomialFeatures(degree=i)), ('classifier' , RandomForestClassifier(n\_jobs=jobs)) ]) #-----# Parameter Grid param\_grid = [{'classifier' : [RandomForestClassifier()], 'classifier\_\_n\_estimators' : n\_estimators, 'classifier\_\_max\_features' : [max\_features], 'classifier\_\_random\_state': [2111]} ] #---- ${\it\#\ Stratified KFold\ f\"{u}r\ unbalancierten\ Datensatz}$ scv = StratifiedKFold(n\_splits=folts) #-----# Gittersuche grid\_search = GridSearchCV( estimator=prediction\_pipe, param\_grid=param\_grid, scoring=gs\_scoring, cv=scv, verbose=True, n\_jobs=jobs, iid=False) #\_\_\_\_\_ model = grid\_search.fit(X,y\_true) return model, grid\_search.best\_score\_ 1.4.1 Modelaufruf und Scoring Modell 1 [14]: rndfrst\_model, rndfrst\_score = rndfrst\_1(1, X\_train, y\_train, my\_scaler=MinMaxScaler, jobs=-3, gs\_scoring='f1', folts=5, n\_estimators=list(range(10, 200 + 1,10)), →max\_features=None) rndfrst\_score Fitting 5 folds for each of 20 candidates, totalling 100 fits [Parallel(n\_jobs=-3)]: Using backend LokyBackend with 10 concurrent workers. [Parallel(n\_jobs=-3)]: Done 30 tasks | elapsed: 16.2s [Parallel(n\_jobs=-3)]: Done 100 out of 100 | elapsed: 2.4min finished [14]: 0.8946162130675438 Modell 2 [15]: rndfrst\_model2, rndfrst\_score2 = rndfrst\_1(1, X\_train, y\_train, my\_scaler=StandardScaler, jobs=-3, gs\_scoring='f1', folts=5, n\_estimators=list(range(10, 200 + 1,10)), rndfrst\_score2 Fitting 5 folds for each of 20 candidates, totalling 100 fits [Parallel(n\_jobs=-3)]: Using backend LokyBackend with 10 concurrent workers. [Parallel(n\_jobs=-3)]: Done 30 tasks | elapsed: 18.6s [Parallel(n\_jobs=-3)]: Done 100 out of 100 | elapsed: 2.4min finished →max\_features=None) rndfrst\_score3 Fitting 5 folds for each of 20 candidates, totalling 100 fits [Parallel(n\_jobs=-3)]: Using backend LokyBackend with 10 concurrent workers. [Parallel(n\_jobs=-3)]: Done 30 tasks | elapsed: 10.3s [Parallel( $n_jobs=-3$ )]: Done 100 out of 100 | elapsed: 1.4min finished Scoring Model 1 / Modell 2 print(model.best\_params\_) model = rndfrst\_model2 print(model.best\_params\_) model = rndfrst\_model3 print(model.best\_params\_) {'classifier': RandomForestClassifier(bootstrap=True, class\_weight=None, criterion='gini', max\_depth=None, max\_features=None, max\_leaf\_nodes=None, min\_impurity\_decrease=0.0, min\_impurity\_split=None, min\_samples\_leaf=1, min\_samples\_split=2, min\_weight\_fraction\_leaf=0.0, n\_estimators=20, n\_jobs=None, oob\_score=False, random\_state=2111, verbose=0, warm\_start=False), 'classifier\_\_max\_features': None, 'classifier\_\_n\_estimators': 20, 'classifier\_\_random\_state': 2111} {'classifier': RandomForestClassifier(bootstrap=True, class\_weight=None, criterion='gini', max\_depth=None, max\_features=None, max\_leaf\_nodes=None, min\_impurity\_decrease=0.0, min\_impurity\_split=None, min\_samples\_leaf=1, min\_samples\_split=2, min\_weight\_fraction\_leaf=0.0, n\_estimators=20, n\_jobs=None, oob\_score=False, random\_state=2111, verbose=0, warm\_start=False), 'classifier\_\_max\_features': None, 'classifier\_\_n\_estimators': 20, 'classifier\_\_random\_state': 2111} {'classifier': RandomForestClassifier(bootstrap=True, class\_weight=None, criterion='gini', max\_depth=None, max\_features=None, max\_leaf\_nodes=None, min\_impurity\_decrease=0.0, min\_impurity\_split=None, min\_samples\_leaf=1, min\_samples\_split=2, min\_weight\_fraction\_leaf=0.0, n\_estimators=20, n\_jobs=None, oob\_score=False, random\_state=2111, verbose=0, warm\_start=False), 'classifier\_\_max\_features': None, 'classifier\_\_n\_estimators': 20, 'classifier\_\_random\_state': 2111} c = 0.35class\_names = ['Stückgut fehlerfrei', 'Stückgut fehlerhaft'] confusion\_matrix1, report1 = get\_confusion\_matrix(X\_train, y\_train, model, class\_names, c) confusion\_matrix1 Prognose des Modells

[15]: 0.8966538306224028 [16]: rndfrst\_model3, rndfrst\_score3 = rndfrst\_1(1, X\_train, y\_train, my\_scaler=RobustScaler, jobs=-3, gs\_scoring='f1', folts=5, n\_estimators=list(range(10, 200 + 1,10)), [16]: 0.8966538306224028 [17]: model = rndfrst\_model [18]: model = rndfrst\_model [18]: Stückgut fehlerfrei Stückgut fehlerhaft Wahrer Wert Stückgut fehlerfrei 16163 Stückgut fehlerhaft 225 0 [19]: roc\_curve\_func(X\_train, y\_train, model) ROC-Kurve 1.0 0.8 True Positive Rate 0.6 0.4 0.2 ROC-Kurve (AUC = 0.99999) 0.0 0.0 0.2 0.4 0.6 0.8 1.0 False Positive Rate [20]: print(report1) precision recall f1-score support 1.00 1.00 16168 1.00 0.98 1.00 0.99 225 accuracy 1.00 16393 macro avg 0.99 1.00 0.99 16393 weighted avg 1.00 1.00 1.00 16393

Scoring auf Validerungsdatensatz [21]: model = rndfrst\_model c = 0.35class\_names = ['Stückgut fehlerfrei', 'Stückgut fehlerhaft'] confusion\_matrix2, report2 = get\_confusion\_matrix(X\_validierung, y\_validierung, model, class\_names, c) confusion\_matrix2 [21]: Prognose des Modells Stückgut fehlerfrei Stückgut fehlerhaft Wahrer Wert Stückgut fehlerfrei 4032 Stückgut fehlerhaft 14 45 [22]: roc\_curve\_func(X\_validierung, y\_validierung, model) ROC-Kurve 1.0 0.8 True Positive Rate 0.6 0.4 0.2 ROC-Kurve (AUC = 0.91304)

0.0 0.2 0.4 0.6 0.8 False Positive Rate [23]: print(report2) recall f1-score support precision 0 1.00 1.00 1.00 4040 0.85 0.76 0.80 59

1

0.99

0.90

0.99

accuracy

macro avg

weighted avg

1.5 Submit

c = 0.35

HHHH

[25]: """

[26]:

model = rndfrst\_model

confusion\_matrix3

c)\nconfusion\_matrix3\n"

1.5.2 Submit der Prognose

 $submission\_fehlerhaft$ 

[26]: '\nsubmission\_fehlerhaft\n

manu\_name=True)\nsubmission\_head\n'

 $submission\_head$ 

[24]: """

0.92

0.99

0.88

0.99

class\_names = ['Stückgut fehlerfrei', 'Stückgut fehlerhaft']

'Stückgut fehlerhaft']\nconfusion\_matrix3, report3 =

1.5.1 Kontrolle Modellwahl (Modell 1 oder 2) anhand der Konfusionsmatrix

[24]: "\nmodel = rndfrst\_model\nc = 0.35\n\nclass\_names = ['Stückgut fehlerfrei',

get\_confusion\_matrix(X\_validierung, y\_validierung, model, class\_names,

[25]: '\nsubmission\_head, submission\_fehlerhaft = submit(model, c, save=True,

1.5.3 Ausgabe DataFrame mit als defekt klassifizierten Stückgütern im Testdatensatz

 $submission\_head, \ submission\_fehlerhaft = submit(model, \ c, \ save=True, \ manu\_name=True)$ 

4099

4099

4099

 $confusion\_matrix3$ ,  $report3 = get\_confusion\_matrix(X\_validierung, y\_validierung, model, class\_names, c)$