Support Vector Machines\_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 Support Vector Machines

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.

plt.ylim([0.0, 1.05])

plt.title('ROC-Kurve')

plt.show()

1.2.3 Funktion für das Submitten

plt.xlabel('False Positive Rate') plt.ylabel('True Positive Rate')

# - model: Modell auf Basis dessen die Prognose erfolgt

[4]: # Definition einer Funktion, welche das Submitten der Prognose auf dem Testdatensatz erleichert

# ---> Wenn None, dann wird die Prognose ohne die Berücksichtung von c vorgenommen # ---> Wenn != None, dann wird Prognose mit der Berücksichtung von c vorgenommen

plt.legend(loc="lower right")

import datetime as dt

1.2 Hilfsfunktionen

#-----# Argumente:

# - c:

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 Abhängigkeit 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

1.1 Package Import

import pandas as pd

[1]: import numpy as np

# ---> Wenn False, dann werden die prognostizierten Daten nicht gespeichert # ---> Wenn True, dann werden die prognostizierten Daten als .csv gespeichert # ---> Wenn None, dann wird ein nicht eindeutiger Standardname als Bezeichnung der .csv gewählt # ---> Wenn != None, dann wird die zu speichernde .csv mit einem timestamp versehen def submit(model, c=None, save=False, manu\_name=False): #-----# Testdatensatz einlesen X\_test = pd.read\_csv('test.csv', index\_col=0) #-----# Prognosewerte auf Sensordaten des Testdatensatzes und unter # Berücksichtung von c erzeugen if c != None: predicted\_test = (model.predict\_proba(X\_test) >= c)[:,1].astype(int) #-----# Prognosewerte auf Sensordaten des Testdatensatzes erzeugen # ohne c zu Berücksichtigen if c == None: predicted\_test = model.predict(X\_test)  ${\it\# Submission datensatz\ einlesen\ und\ prognostizierte\ Werte\ hineinschreiben}$ submission = pd.read\_csv('sample\_submission.csv') submission['Fehlerhaft'] = predicted\_test #-----# In .csv speichern, wenn save=True if save == True: # Standardnamen wählen, wenn manu\_name == False if manu\_name == False: submission.to\_csv('./predicted\_values.csv', index=False) #-----# Standardnamen mit timestamp kombinieren, wenn manu\_name == True if manu\_name == True: import datetime now = datetime.datetime.now() name = now.strftime('\%Y-\%m-\%dT\%H\%M\%S') + ('-\%02d' \% (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 [6]: | #-----# 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 Stückgüter\ncolum\_to\_drop = 'Fehlerhaft'\norignal\_df = df\_fehlerhaft\nlow = .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' [9]:  $df_-fehlerhaft.describe()$ [9]: '\ndf\_fehlerhaft.describe()\n' [10]: data\_new = data #\_filtered data\_new['Fehlerhaft'].value\_counts() [10]: 0 20208 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)) Sensor\_1 -20 7.5 Sensor 2 2.5 0.0 -2.5Sensor\_3 Sensor 4 5.0 2.5 0.0 -2.5Sensor\_5 -1510 Sensor\_6 -10Sensor\_7 -10-15 -20Sensor 8 -10-151.0 0.8 Fehlerhaft 0.6 0.4 0.2 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 KNeighborsClassifier 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 # - C: Parameter für die Regularisierung; je kleiner, desto stärker die Regularisierung # - kernel: Im algorithmus verwendeter Kernel # - gamma: Kernelkoeffizient # - jobs: Anzahl der Threads die für den Durchlauf zur Verfügung stehen # - gs\_scoring: Scoring Verfahren im Rahmen der GridSearch # - folts: Komplexität der Kreuzvalidierung def svm\_1(i, X, y\_true, my\_scaler=MinMaxScaler, jobs=-3, gs\_scoring='f1', folts=5, C=[0.1, 1, 10, 100, 1000], kernel=['linear', 'poly', 'rbf', 'sigmoid'], gamma=['scale']): #-----# Pipeline erzeugen prediction\_pipe = Pipeline([('scaler', my\_scaler()), ('add\_x\_square', PolynomialFeatures(degree=i)), ('classifier', SVC(cache\_size=7000, probability=True))]) # Parameter Grid param\_grid = [{'classifier': [SVC(cache\_size=7000, probability=True)], 'classifier\_\_C': C, #'classivier\_\_degree': degree, 'classifier\_\_gamma': gamma, 'classifier\_\_kernel': kernel} ] #-----# StratifiedKFold für 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 →gamma=['scale']) svm\_score Fitting 5 folds for each of 15 candidates, totalling 75 fits [Parallel(n\_jobs=10)]: Using backend LokyBackend with 10 concurrent workers. [Parallel(n\_jobs=10)]: Done 30 tasks | elapsed: 8.4s [Parallel(n\_jobs=10)]: Done 75 out of 75 | elapsed: 5.3min finished Modell 2 [15]: svm\_model2, svm\_score2 = svm\_1(1, X\_train, y\_train, my\_scaler=StandardScaler, jobs=10, gs\_scoring='f1', C=[0.1, 1, 10, 100, 1000], kernel=['linear', 'poly', 'rbf'], | →gamma=['scale']) svm\_score2 Fitting 5 folds for each of 15 candidates, totalling 75 fits  $[Parallel(n\_jobs=10)]: \ Using \ backend \ LokyBackend \ with \ 10 \ concurrent \ workers.$ [Parallel(n\_jobs=10)]: Done 30 tasks | elapsed: 29.8s [Parallel(n\_jobs=10)]: Done 75 out of 75 | elapsed: 20.3min finished [15]: 0.8807232413976462 [16]: svm\_model3, svm\_score3 = svm\_1(1, X\_train, y\_train, my\_scaler=RobustScaler, jobs=10, gs\_scoring='f1', C=[0.1, 1, 10, 100, 1000], kernel=['linear', 'poly', 'rbf'], |  $\rightarrow$ gamma=['scale']) svm\_score3 Fitting 5 folds for each of 15 candidates, totalling 75 fits [Parallel(n\_jobs=10)]: Using backend LokyBackend with 10 concurrent workers. [Parallel(n\_jobs=10)]: Done 30 tasks | elapsed: 48.9s [Parallel(n\_jobs=10)]: Done 75 out of 75 | elapsed: 39.6min finished Scoring Model 1 / Modell 2 print(model.best\_params\_) model = svm\_model2 print(model.best\_params\_) model = svm\_model3 print(model.best\_params\_) {'classifier': SVC(C=1000, cache\_size=7000, class\_weight=None, coef0=0.0, decision\_function\_shape='ovr', degree=3, gamma='scale', kernel='rbf', max\_iter=-1, probability=True, random\_state=None, shrinking=True, tol=0.001, verbose=False), 'classifier\_\_C': 1000, 'classifier\_\_gamma': 'scale', 'classifier\_\_kernel': 'rbf'} {'classifier': SVC(C=0.1, cache\_size=7000, class\_weight=None, coef0=0.0, decision\_function\_shape='ovr', degree=3, gamma='scale', kernel='linear', max\_iter=-1, probability=True, random\_state=None, shrinking=True, tol=0.001, verbose=False), 'classifier\_\_C': 0.1, 'classifier\_\_gamma': 'scale', 'classifier\_\_kernel': 'linear'} {'classifier': SVC(C=1, cache\_size=7000, class\_weight=None, coef0=0.0, decision\_function\_shape='ovr', degree=3, gamma='scale', kernel='rbf', max\_iter=-1, probability=True, random\_state=None, shrinking=True, tol=0.001, verbose=False), 'classifier\_\_C': 1, 'classifier\_\_gamma': 'scale', 'classifier\_\_kernel': 'rbf'} [33]: model = svm\_model2 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 Stückgut fehlerfrei Stückgut fehlerhaft Wahrer Wert Stückgut fehlerfrei 16164 Stückgut fehlerhaft 179 46 [34]: roc\_curve\_func(X\_train, y\_train, model) ROC-Kurve 1.0 0.8 0.6

[14]: svm\_model, svm\_score = svm\_1(1, X\_train, y\_train, my\_scaler=MinMaxScaler, jobs=10, gs\_scoring='f1', C=[0.1, 1, 10, 100, 1000], kernel=['linear', 'poly', 'rbf'], | [14]: 0.8876821136601146 [16]: 0.8877796901893287 [17]: model = svm\_model [33]: True Positive Rate 0.4 0.2 ROC-Kurve (AUC = 0.95827) 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 0.99 0.86 0.92 225 accuracy 1.00 16393 0.99 0.93 16393 macro avg 0.96 weighted avg 1.00 1.00 1.00 16393 Scoring auf Validerungsdatensatz [31]: model = svm\_model2 c = 0.35class\_names = ['Stückgut fehlerfrei', 'Stückgut fehlerhaft'] confusion\_matrix2, report2 = get\_confusion\_matrix(X\_validierung, y\_validierung, model, class\_names, c) confusion\_matrix2 Prognose des Modells [31]: Stückgut fehlerfrei Stückgut fehlerhaft Wahrer Wert Stückgut fehlerfrei 4037 Stückgut fehlerhaft 17 42 [32]: roc\_curve\_func(X\_validierung, y\_validierung, model) ROC-Kurve 1.0 0.8 True Positive Rate 0.6 0.4 ROC-Kurve (AUC = 0.95000) 0.0 0.2 0.4 0.6 0.8 0.0 1.0 False Positive Rate [23]: print(report2) precision recall f1-score support

1.00

0.73

0.86

1.00

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

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

1.00

0.83

1.00

0.91

1.00

Prognose des Modells

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

1.5.3 Ausgabe DataFrame mit als defekt klassifizierten Stückgütern im Testdatensat

4040

4099

4099

4099

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

Stückgut fehlerfrei Stückgut fehlerhaft

42

1

4037

17

59

1.00

0.96

0.98

1.00

accuracy

macro avg

weighted avg

1.5 Submit

c = 0.35

model = svm\_model2

confusion\_matrix3

Wahrer Wert Stückgut fehlerfrei

1.5.2 Submit der Prognose

ID Fehlerhaft

 $submission\_fehlerhaft$ 

35272

35026

35341 35102

35186

[175 rows x 2 columns]

14779 35316 14796 35054 14809 35435 14871 21379 14886 35019

0

0

0

0

ID Fehlerhaft

1

1

 $submission\_head$ 

2801

23197

9209

1 13844

4 31773

Stückgut fehlerhaft

11 11 11

HHHH

[35]:

[35]:

[36]:

[36]:

0

[37]: """

16 55

154

166

241

[37]: