LogisticRegression_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 LogisticRegression

1.1 Package Import

import pandas as pd

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.

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 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):

X_test = pd.read_csv('test.csv', index_col=0)

#-----

Testdatensatz einlesen

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

[1]: import numpy as np

#-----# 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 einer LogisticRegression 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_range: Regulaisierungsstärke; kleinere Werte führen zu einer stärkeren Regularisierung # - solver: Algorithmus der zur Problemlösung verwendet wird # - penalty: Bei der Regularisierung verwendete Norm # - jobs: Anzahl der Threads die für den Durchlauf zur Verfügung stehen # - gs_scoring: Scoring Verfahren im Rahmen der GridSearch # - folts: Komplexität der Kreuzvalidierung \rightarrow logspace(-4, 4, 20)): #-----# Pipeline erzeugen prediction_pipe = Pipeline([('scaler', my_scaler()), ('add_x_square', PolynomialFeatures(degree=i)), ('classifier', LogisticRegression(n_jobs=jobs, max_iter=10000)) # Parameter Grid param_grid = [{'classifier' : [LogisticRegression(max_iter=10000)], 'classifier__penalty': penalty, 'classifier__C': c_range, 'classifier__solver': solver}, #-----# StratifiedKFold für unbalancierten Datensatz scv = StratifiedKFold(n_splits=folts) #-----

def logReg_1(i, X, y_true, my_scaler=MinMaxScaler, jobs=-3, gs_scoring='f1', folts=5, solver=['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga'], penalty=['12'], c_range=np. # Gittersuche grid_search = GridSearchCV(estimator=prediction_pipe, param_grid=param_grid, scoring=gs_scoring, cv=scv, verbose=True, n_jobs=jobs, iid=False) # Fit model = grid_search.fit(X,y_true) return model, grid_search.best_score_ 1.4.1 Modelaufruf und Scoring →'saga'], penalty=['12'], c_range=np.logspace(-4, 4, 20)) logReg_score Fitting 5 folds for each of 100 candidates, totalling 500 fits [Parallel(n_jobs=-3)]: Using backend LokyBackend with 10 concurrent workers. [Parallel(n_jobs=-3)]: Done 30 tasks | elapsed: 1.9s [Parallel(n_jobs=-3)]: Done 500 out of 500 | elapsed: 49.3s finished Modell 2 \rightarrow 'saga'], penalty=['12'], c_range=np.logspace(-4, 4, 20)) logReg_score2 Fitting 5 folds for each of 100 candidates, totalling 500 fits [Parallel(n_jobs=-3)]: Using backend LokyBackend with 10 concurrent workers. [Parallel(n_jobs=-3)]: Done 60 tasks | elapsed: 1.3s [Parallel(n_jobs=-3)]: Done 326 tasks | elapsed: 1.9min [Parallel(n_jobs=-3)]: Done 500 out of 500 | elapsed: 2.3min finished \rightarrow 'saga'], penalty=['12'], c_range=np.logspace(-4, 4, 20)) logReg_score3 Fitting 5 folds for each of 100 candidates, totalling 500 fits [Parallel(n_jobs=-3)]: Using backend LokyBackend with 10 concurrent workers. [Parallel(n_jobs=-3)]: Done 60 tasks | elapsed: 1.6s [Parallel(n_jobs=-3)]: Done 309 tasks | elapsed: 1.7min [Parallel(n_jobs=-3)]: Done 500 out of 500 | elapsed: 2.0min finished Scoring Model 1 / Modell 2 print(model.best_params_) model = logReg_model2 print(model.best_params_) model = logReg_model3 print(model.best_params_) {'classifier': LogisticRegression(C=206.913808111479, class_weight=None, dual=False, $\label{limit} \verb|fit_intercept=True|, intercept_scaling=1, l1_ratio=None|,$ max_iter=10000, multi_class='warn', n_jobs=None, penalty='12', random_state=None, solver='newton-cg', tol=0.0001, verbose=0, warm_start=False), 'classifier__C': 206.913808111479, 'classifier_penalty': 'l2', 'classifier_solver': 'newton-{'classifier': LogisticRegression(C=0.615848211066026, class_weight=None, dual=False, fit_intercept=True, intercept_scaling=1, l1_ratio=None, max_iter=10000, multi_class='warn', n_jobs=None, penalty='12', random_state=None, solver='newton-cg', tol=0.0001, verbose=0, warm_start=False), 'classifier__C': O.615848211066026, 'classifier__penalty': 'l2', 'classifier__solver': 'newton-{'classifier': LogisticRegression(C=0.23357214690901212, class_weight=None, dual=False, fit_intercept=True, intercept_scaling=1, l1_ratio=None, max_iter=10000, multi_class='warn', n_jobs=None, penalty='12', random_state=None, solver='newton-cg', tol=0.0001, verbose=0, warm_start=False), 'classifier__C': 0.23357214690901212, 'classifier__penalty': '12', 'classifier__solver': 'newtoncg'} c = 0.35

[14]: logReg_model, logReg_score = logReg_1(1, X_train, y_train, my_scaler=MinMaxScaler, jobs=-3, gs_scoring='f1', folts=5, solver=['newton-cg', 'lbfgs', 'liblinear', 'sag', | [15]: logReg_model2, logReg_score2 = logReg_1(1, X_train, my_scaler=StandardScaler, jobs=-3, gs_scoring='f1', folts=5, solver=['newton-cg', 'lbfgs', 'liblinear', 'sag', | [15]: 0.8811243461271822 [16]: logReg_model3, logReg_score3 = logReg_1(1, X_train, y_train, my_scaler=RobustScaler, jobs=-3, gs_scoring='f1', folts=5, solver=['newton-cg', 'lbfgs', 'liblinear', 'sag', | [16]: 0.8811243461271822 [17]: model = logReg_model [18]: model = logReg_model class_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 16160 Stückgut fehlerhaft ROC-Kurve 1.0 0.8 True Positive Rate 0.6 0.4 0.2 ROC-Kurve (AUC = 0.97360) 0.0 0.2 0.4 0.6 0.0 0.8 1.0 False Positive Rate [20]: print(report1) precision recall f1-score support 1.00 1.00 1.00 16168 0.96 0.82 0.88 225 16393 accuracy 1.00 macro avg 0.98 0.91 0.94 16393 weighted avg 1.00 1.00 1.00 16393 Scoring auf Validerungsdatensatz [21]: model = logReg_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 Prognose des Modells Stückgut fehlerfrei Stückgut fehlerhaft Wahrer Wert Stückgut fehlerfrei 4038 Stückgut fehlerhaft 13 46 [22]: roc_curve_func(X_validierung, y_validierung, model) ROC-Kurve 1.0

[18]: [19]: roc_curve_func(X_train, y_train, model) [21]: 0.8 True Positive Rate 0.6 0.4 0.2 ROC-Kurve (AUC = 0.97029) 0.0 0.4 0.6 0.2 0.8 False Positive Rate [23]: print(report2) recall f1-score precision support 0 1.00 1.00 1.00 4040 0.96 0.78 0.86 59 1.00 4099 accuracy 0.98 0.89 0.93 4099 macro avg weighted avg 1.00 1.00 1.00 4099 1.5 Submit 1.5.1 Kontrolle Modellwahl (Modell 1 oder 2) anhand der Konfusionsmatrix [24]: """ $model = logReg_model$ c = 0.35class_names = ['Stückgut fehlerfrei', 'Stückgut fehlerhaft'] $confusion_matrix3$, $report3 = get_confusion_matrix(X_validierung, y_validierung, model, class_names, c)$ $confusion_matrix3$ [24]: "\nmodel = logReg_model\nc = 0.35\n\nclass_names = ['Stückgut fehlerfrei', 'Stückgut fehlerhaft']\nconfusion_matrix3, report3 = get_confusion_matrix(X_validierung, y_validierung, model, class_names, c)\nconfusion_matrix3\n"

1

1.5.2 Submit der Prognose

submission_fehlerhaft

[26]: '\nsubmission_fehlerhaft\n'

manu_name=True) \nsubmission_head \n'

 $submission_head$

submission_head, submission_fehlerhaft = submit(model, c, save=True, manu_name=True)

[25]: '\nsubmission_head, submission_fehlerhaft = submit(model, c, save=True,

1.5.3 Ausgabe DataFrame mit als defekt klassifizierten Stückgütern im Testdatensatz

[25]: """

[26]: """