SVM - classification SHUMBUL ARIFA \ 181CO152 Task Performing Kernel SVM on a classification dataset. 1. Linear Kernel 2. Polynomial Kernel 3. Radial Basis Function (RBF) kernel import numpy as np In [1]: import pandas as pd import seaborn as sns import matplotlib as plt df = pd.read_csv("Movie_classification.csv") In [2]: df.head() Out[2]: Marketing Production Multiplex Lead Budget Movie_length Lead_Actress_rating Director_rating Producer_rating Critic_ı expense expense coverage Actor_Rating 20.1264 59.62 0.462 36524.125 138.7 7.825 8.095 7.910 7.995 20.5462 69.14 0.531 35668.655 152.4 7.505 7.650 7.440 7.470 2 20.5458 69.14 0.531 39912.675 134.6 7.485 7.570 7.495 7.515 0.542 38873.890 3 20.6474 59.36 119.3 6.895 7.035 6.920 7.020 21.3810 59.36 0.542 39701.585 127.7 6.920 7.070 6.815 7.070 In [3]: df.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 506 entries, 0 to 505 Data columns (total 19 columns): # Column Non-Null Count Dtype - - -0 Marketing expense 506 non-null float64 Production expense float64 1 506 non-null 2 Multiplex coverage 506 non-null float64 3 506 non-null float64 Budget 4 Movie_length 506 non-null float64 506 non-null float64 Lead_ Actor_Rating 6 506 non-null float64 Lead_Actress_rating 7 Director_rating 506 non-null float64 8 Producer_rating 506 non-null float64 9 Critic_rating 506 non-null float64 10 Trailer_views 506 non-null int64 3D_available 506 non-null object 11 12 Time_taken 494 non-null float64 506 non-null float64 13 Twitter_hastags 506 non-null object 14 Genre 506 non-null 15 Avg_age_actors int64 int64 16 Num_multiplex 506 non-null 17 Collection 506 non-null int64 506 non-null Start_Tech_Oscar int64 dtypes: float64(12), int64(5), object(2) memory usage: 75.2+ KB **Data Cleaning and preprocessing** time_taken has some missing values 2. 3D_available and Genre -> object type (string) In [4]: mean = df['Time_taken'].mean() Out[4]: 157.39149797570855 In [5]: df['Time_taken'].fillna(value = mean, inplace = True) df.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 506 entries, 0 to 505 Data columns (total 19 columns): Column Non-Null Count Dtype 506 non-null 0 Marketing expense float64 Production expense 506 non-null float64 Multiplex coverage 506 non-null float64 3 506 non-null float64 Budget 4 Movie_length 506 non-null float64 Lead_ Actor_Rating 506 non-null float64 Lead_Actress_rating 506 non-null float64 7 506 non-null float64 Director_rating 8 Producer_rating 506 non-null float64 9 Critic_rating 506 non-null float64 int64 10 Trailer_views 506 non-null 11 3D_available 506 non-null object 12 Time_taken 506 non-null float64 506 non-null float64 Twitter_hastags 14 506 non-null object Genre Avg_age_actors 506 non-null 15 int64 16 Num_multiplex 506 non-null int64 Collection 506 non-null int64 18 Start_Tech_Oscar 506 non-null int64 dtypes: float64(12), int64(5), object(2) memory usage: 75.2+ KB In [6]: ## 3D-available and Genre # ### Using dummy variable creation # df = pd.get_dummies(df, columns = ["3D_available", "Genre"]) obj_df = df.select_dtypes(include=['object']).copy() obj_df.head() Out[6]: Genre 3D_available YES Thriller 1 NO Drama NO Comedy 3 YES Drama NO Drama ## if any null value is present in those rows In [7]: obj_df[obj_df.isnull().any(axis=1)] Out[7]: 3D_available Genre # ## if it was present in column "c" # obj_df["c"].value_counts() # obj_df = obj_df.fillna({"c": "NEW_NAME"}) In [9]: ## replace {"YES": 1, "NO": 0}, cleanup_nums = {"3D_available": "Genre": {"Thriller": 0, "Drama": 1, "Comedy": 2, "Action": 3}} In [10]: | ## replace only once! df = df.replace(cleanup_nums) df.info() ## done <class 'pandas.core.frame.DataFrame'> RangeIndex: 506 entries, 0 to 505 Data columns (total 19 columns): Column Non-Null Count Dtype 0 Marketing expense 506 non-null float64 Production expense 506 non-null float64 Multiplex coverage 506 non-null float64 3 506 non-null float64 Budget 506 non-null float64 Movie_length Lead_ Actor_Rating 506 non-null float64 506 non-null Lead_Actress_rating float64 7 float64 Director_rating 506 non-null 8 Producer_rating 506 non-null float64 Critic_rating 506 non-null float64 Trailer_views 506 non-null int64 506 non-null int64 3D_available Time_taken 506 non-null float64 Twitter_hastags 506 non-null float64 506 non-null int64 14 Genre 15 Avg_age_actors 506 non-null int64 506 non-null Num_multiplex int64 Collection 506 non-null int64 18 Start_Tech_Oscar 506 non-null int64 dtypes: float64(12), int64(7) memory usage: 75.2 KB X_y split In [11]: | X = df.loc[:,df.columns!="Start_Tech_Oscar"] # All cols except collection type(X) Out[11]: pandas.core.frame.DataFrame In [12]: X.head() Out[12]: Marketing Production Multiplex Lead Budget Movie_length Lead_Actress_rating Director_rating Producer_rating Critic_i Actor_Rating expense expense coverage 20.1264 59.62 0.462 36524.125 138.7 7.825 8.095 7.910 7.995 1 20.5462 69.14 0.531 35668.655 152.4 7.505 7.650 7.440 7.470 20.5458 69.14 0.531 39912.675 134.6 7.485 7.570 7.495 7.515 3 20.6474 59.36 0.542 38873.890 119.3 6.895 7.035 6.920 7.020 21.3810 59.36 0.542 39701.585 127.7 6.920 7.070 6.815 7.070 In [13]: | X.shape Out[13]: (506, 18) In [14]: y = df["Start_Tech_Oscar"] type(y) Out[14]: pandas.core.series.Series In [15]: y.head() Out[15]: 0 1 0 2 1 3 1 Name: Start_Tech_Oscar, dtype: int64 In [16]: y.shape Out[16]: (506,) **Test-Train Split** from sklearn.model_selection import train_test_split In [17]: In [18]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 0) X_train.head() In [19]: Out[19]: Marketing Production Multiplex Lead Budget Movie_length Lead_Actress_rating Director_rating Producer_rating Critic Actor_Rating expense expense coverage 220 27.1618 0.493 38612.805 67.40 162.0 8.485 8.640 8.485 8.670 0.587 33113.355 71 23.1752 76.62 91.0 7.280 7.400 7.290 7.455 22.2658 0.572 38312.835 240 64.86 127.8 6.755 6.935 6.800 6.840 0.476 33396.660 6 21.7658 70.74 140.1 7.065 7.265 7.150 7.400 417 538.8120 91.20 0.321 29463.720 162.6 9.135 9.305 9.095 9.165 In [20]: y_test.head() Out[20]: 329 0 371 1 219 0 403 0 78 Name: Start_Tech_Oscar, dtype: int64 In [21]: X_train.shape Out[21]: (404, 18) **Standardizing Data** Coverting mean and variance close to 0 and 1, for each variable. SVM only gives correct result when we standardize our data! Ways: StandardScaler, MinMax scaler In [22]: from sklearn.preprocessing import StandardScaler In [23]: sc = StandardScaler().fit(X_train) X_train_std = sc.transform(X_train) X_test_std = sc.transform(X_test) In [24]: X_test_std # here, we only need to std our X data, not y Out[24]: array([[-4.08358690e-01, -1.12872913e+00, 1.12308956e+00, -8.87385815e-01, 1.15409837e-03], [7.19251107e-01, 9.98884403e-01, -6.52839787e-01, ..., -1.15123717e+00, 6.08961586e-01, 2.97217905e+00], [-4.02574884e-01, 3.96108293e-01, 5.11537670e-02, ..., -1.47614099e+00, 1.51479578e-01, 4.45267254e-02], [-3.98260097e-01, -8.58124181e-01, 8.94207776e-01, ..., -7.45107395e-01, -1.01128719e+00, -4.21729015e-01], [-3.99342792e-01, -7.63765430e-02, 5.81321752e-01, ..., -2.93820817e+00, -9.92225442e-01, 5.97527720e-01], [-4.00880712e-01, -3.67026306e-01, 3.11892120e-01, ..., 7.16959787e-01, -4.10842057e-01, -3.02454291e-01]]) All decimals, scales of values changed -> uniform scale Now, we can perform SVM **Performing SVM Classification Linear Kernel** In [25]: from sklearn.svm import SVC ## C - classification svc = SVC(kernel='linear', C=0.01) In [26]: | svc.fit(X_train_std, y_train) Out[26]: SVC(C=0.01, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=3, gamma='scale', kernel='linear', max_iter=-1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False) **Predict values using trained model** In [27]: y_test_pred = svc.predict(X_test_std) y_train_pred = svc.predict(X_train_std) In [28]: y_test_pred 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]) **Model Performance** In [29]: **from sklearn.metrics import** classification_report, accuracy_score, confusion_matrix In [30]: | svc.score(X_test, y_test) Out[30]: 0.43137254901960786 In [31]: confusion_matrix(y_test, y_test_pred) Out[31]: array([[7, 37], [1, 57]]) In [32]: | accuracy_score(y_test, y_test_pred) Out[32]: 0.6274509803921569 In [33]: | svc.n_support_ ## What's this?? ### number of support vectors ### for 1st class (0) and 2nd class (1) RESPECTIVELY Out[33]: array([186, 191], dtype=int32) **Grid Search** To choose the ost accurate parameter C In [34]: from sklearn.model_selection import GridSearchCV In [35]: params = $\{ C': (0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50, 100, 150, 500, 1000) \}$ In [36]: svc_op = SVC(kernel='linear') In [37]: | svm_grid_lin = GridSearchCV(svc_op, params, n_jobs=-1, cv=10, verbose=1, scoring='accuracy') In [38]: | svm_grid_lin.fit(X_train_std, y_train) Fitting 10 folds for each of 14 candidates, totalling 140 fits [Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers. [Parallel(n_jobs=-1)]: Done 88 tasks | elapsed: [Parallel(n_jobs=-1)]: Done 140 out of 140 | elapsed: 1.3min finished Out[38]: GridSearchCV(cv=10, error_score=nan, estimator=SVC(C=1.0, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=3, gamma='scale', kernel='linear', max_iter=-1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False), iid='deprecated', n_jobs=-1, param_grid={'C': (0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50, 100, 150, 500, 1000)}, pre_dispatch='2*n_jobs', refit=True, return_train_score=False, scoring='accuracy', verbose=1) In [39]: svm_grid_lin.best_params_ ## What's this?? ### get the best parameter!!! Out[39]: {'C': 0.5} In [40]: linsvm_clf = svm_grid_lin.best_estimator_ Out[40]: SVC(C=0.5, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=3, gamma='scale', kernel='linear', max_iter=-1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False) In [41]: | accuracy_score(y_test, linsvm_clf.predict(X_test_std)) Out[41]: 0.6078431372549019 Polynomial Kernel In [51]: poly_svm = SVC(kernel = 'poly', degree = 0.1, C=0.1) poly_svm.fit(X_train_std, y_train) Out[51]: SVC(C=0.1, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=0.1, gamma='scale', kernel='poly', max_iter=-1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False) In [52]: | y_train_pred = poly_svm.predict(X_train_std) y_test_pred = poly_svm.predict(X_test_std) In [53]: | accuracy_score(y_test, y_test_pred) Out[53]: 0.5686274509803921 In [54]: | poly_svm.n_support_ Out[54]: array([186, 186], dtype=int32) In [55]: params = {'C':(0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50, 100, 150, 500, 1000)} svm_grid_poly = GridSearchCV(poly_svm, params, n_jobs=-1, cv=10, verbose=1, scoring='accuracy') In [56]: | svm_grid_poly.fit(X_train_std, y_train) Fitting 10 folds for each of 14 candidates, totalling 140 fits [Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers. [Parallel(n_jobs=-1)]: Done 68 tasks | elapsed: 2.4s [Parallel(n_jobs=-1)]: Done 140 out of 140 | elapsed: 2.7s finished Out[56]: GridSearchCV(cv=10, error_score=nan, estimator=SVC(C=0.1, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=0.1, gamma='scale', kernel='poly', max_iter=-1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False), iid='deprecated', n_jobs=-1, param_grid={'C': (0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50, 100, 150, 500, 1000)}, pre_dispatch='2*n_jobs', refit=True, return_train_score=False, scoring='accuracy', verbose=1) In [57]: | svm_grid_poly.best_params_ Out[57]: {'C': 0.001} In [58]: polysvm_clf = svm_grid_poly.best_estimator_ polysvm_clf Out[58]: SVC(C=0.001, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=0.1, gamma='scale', kernel='poly', max_iter=-1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False) In [59]: | accuracy_score(y_test, polysvm_clf.predict(X_test_std)) Out[59]: 0.5686274509803921 Radial Kernel In [62]: rad_svm = SVC(kernel = 'rbf', gamma = 0.1, C=10) rad_svm.fit(X_train_std, y_train) Out[62]: SVC(C=10, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=3, gamma=0.1, kernel='rbf', max_iter=-1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False) In [63]: y_train_pred = rad_svm.predict(X_train_std) y_test_pred = rad_svm.predict(X_test_std) In [64]: | accuracy_score(y_test, y_test_pred) Out[64]: 0.5294117647058824 In [65]: rad_svm.n_support_ Out[65]: array([150, 171], dtype=int32) In [67]: params = {'C':(0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50), 'gamma': (0.001, 0.01, 0.1, 0.5, 1)} svm_grid_rad = GridSearchCV(rad_svm, params, n_jobs=-1, cv=10, verbose=1, scoring='accuracy') In [68]: svm_grid_rad.fit(X_train_std, y_train) Fitting 10 folds for each of 50 candidates, totalling 500 fits [Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers. | elapsed: [Parallel(n_jobs=-1)]: Done 68 tasks [Parallel(n_jobs=-1)]: Done 500 out of 500 | elapsed: 4.5s finished Out[68]: GridSearchCV(cv=10, error_score=nan, estimator=SVC(C=10, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=3, gamma=0.1, kernel='rbf', max_iter=-1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False), iid='deprecated', n_jobs=-1, param_grid={'C': (0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1, 5, 10, 50), 'gamma': (0.001, 0.01, 0.1, 0.5, 1)}, pre_dispatch='2*n_jobs', refit=True, return_train_score=False, scoring='accuracy', verbose=1) In [69]: svm_grid_rad.best_params_ Out[69]: {'C': 50, 'gamma': 0.001} In [70]: radsvm_clf = svm_grid_rad.best_estimator_ radsvm_clf Out[70]: SVC(C=50, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=3, gamma=0.001, kernel='rbf', max_iter=-1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False) In [72]: | accuracy_score(y_test, radsvm_clf.predict(X_test_std)) Out[72]: 0.5980392156862745