


 mohitKhanna1411 / COMP9321_20T1_UNSW[Code](#) [Issues](#) [Pull requests](#) [Actions](#) [Projects](#) [Wiki](#) [Security](#) [Insights](#) master ▾

...

[COMP9321_20T1_UNSW](#) / [Assignment_3](#) / [z5266543.py](#) / [Code](#) Jump to ▾

mohitKhanna1411 ready to push

 History 1 contributor

Raw

Blame



226 lines (188 sloc) | 9.85 KB

```
1 # Author: Mohit Khanna
2 # Student ID: z5266543
3 # Platform: Mac
4
5 import os
6 import ast
7 import sys
8 import pandas as pd
9 from sklearn import metrics
10 from scipy.stats import pearsonr
11 from sklearn.ensemble import RandomForestRegressor, GradientBoostingClassifier
12
13
14 zid = os.path.splitext(os.path.basename(__file__))[0]
15
16 #####
17 #                               Downloaded List
18 holiday_month = list(['04', '05', '06', '11', '12'])
19 elite_list_production_companies = list(['Warner Bros', 'Sony Pictures', 'Walt Disney', 'Universal Pictures',
20                                         'Pixar Animation', 'New Line Cinema', 'Twentieth Century Fox Film Cor
21                                         'Paramount Pictures', 'Marvel Studios', 'Columbia Pictures'])
22 elite_list_cast = list(['Samuel L. Jackson', 'Harrison Ford', 'Robert Downey Jr', 'Tom Hanks', 'Morgan Freeman',
23                         'Scarlett Johansson', 'Andy Serkis', 'Anthony Daniels', 'Tom Cruise', 'Eddie Murphy',
24                         'Stanley Tucci', 'Idris Elba', 'Johnny Depp', 'Ian McKellan', 'Bradley Cooper', 'Don
25                         'Vin Diesel', 'Michael Caine', 'Gary Oldman', 'Dwayne Johnson', 'Zoe Saldana', 'Stell
26                         'Robin Williams', 'Woody Harrelson', 'Cate Blanchett', 'Robert DeNiro', 'Bruce Willis
27                         'Emma Watson', 'Will Smith', 'Matt Damon', 'Chris Pratt', 'Chris Evans', 'Chris Hemsw
28                         'Cameron Diaz', 'Liam Neeson', 'Orlando Bloom', 'Steve Carell', 'Simon Pegg', 'Mark
29                         'Helena Bonham Carter', 'Mark Wahlberg', 'Owen Wilson', 'Julia Roberts', 'Ralph Fienn
30                         'Carrie Fisher', 'Elizabeth Banks', 'Forest Whitaker', 'Ben Stiller', 'Adam Sandler',
31 elite_list_crew = list(['Gore Verbinski', 'David Yates', 'Ridley Scott', 'J.J. Abrams', 'George Lucas', 'Chri
32                         'Clint Eastwood', 'Tim Burton', 'Ron Howard', 'Christopher Nolan', 'James Cameron', '
33                         'Peter Jackson', 'Michael Bay', 'Steven Spielberg', 'Kevin Feige', 'Kathleen Kennedy'
34                         'Jerry Bruckheimer', 'Neal H. Moritz', 'Frank Marshall', 'Charles Roven', 'Lorenzo di
35                         'Ian Bryce', 'Lauren Shuler Donner', 'Avi Arad', 'Christopher Meledandri', 'Janet Hea
```

```
36         'Steven Spielberg', 'Christopher Markus', 'Stephen McFeely', 'Andrew Stanton', 'Steve
37         'David Koepp', 'Lawrence Kasdan', 'Fran Walsh', 'Terry Rossio', 'Ted Elliott'])
38 #####
39
40
41 #####
42 #                                     Data Pre-processing
43 def normalize_columns(data):
44     if isinstance(data, list):
45         res = list()
46         for i in range(len(data)):
47             datum = data[i]
48             res.append(datum['name'])
49         return ','.join(res)
50
51
52 def normalize_production_countries(data):
53     if "United States of America" in data:
54         if len(data.split(",")) == 1:
55             return "USA"
56         else:
57             return "USA,Other_countries"
58     else:
59         return "Other_countries"
60
61
62 def normalize_spoken_languages(data):
63     if "English" in data:
64         if len(data.split(",")) == 1:
65             return "English"
66         else:
67             return "English,Other_lang"
68     else:
69         return "Other_lang"
70
71
72 def normalize_keywords(data):
73     length = len(data.split(','))
74     if length <= 5:
75         return 1
76     elif length >= 6 and length <= 15:
77         return 2
78     elif length >= 16:
79         return 3
80
81
82 def normalize_runtime(data):
83     if data <= 60:
84         return 1
85     elif data >= 61 and data <= 100:
86         return 2
87     elif data >= 101 and data <= 140:
88         return 3
89     elif data >= 141 and data <= 180:
90         return 4
91     elif data >= 181:
```

```
92         return 5
93
94
95 def normalize_production_companies(data):
96     if any(e in data for e in elite_list_production_companies):
97         if len(data.split(',')) == 1:
98             return "Elite Group"
99         else:
100             return "Elite Group,Other_companies"
101     else:
102         return "Other_companies"
103
104
105 def normalize_gender(data):
106     male = 0
107     female = 0
108     neutral = 0
109     if isinstance(data, list):
110         for i in range(len(data)):
111             datum = data[i]
112             if datum['gender'] == 1:
113                 female += 1
114             elif datum['gender'] == 2:
115                 male += 1
116             else:
117                 neutral += 1
118
119     return pd.Series([male, female, neutral], index=['Male', 'Female', 'Neutral'])
120
121
122 def normalize_cast_crew(data, elite_list):
123     ctr = 0
124     for e in elite_list:
125         if e in data:
126             ctr += 1
127     return ctr
128
129
130 def preprocess(dataframe):
131     df = dataframe.copy()
132     columns_to_normalize = ['genres', 'production_countries', 'cast', 'crew',
133                             'keywords', 'production_companies', 'spoken_languages']
134     for c in columns_to_normalize:
135         dataframe[c] = dataframe[c].apply(lambda data: ast.literal_eval(data))
136     for c in columns_to_normalize:
137         df[c] = dataframe[c].apply(normalize_columns)
138     df['runtime'] = df['runtime'].apply(normalize_runtime)
139     df['production_countries'] = df['production_countries'].apply(
140         normalize_production_countries)
141     df['production_companies'] = df['production_companies'].apply(
142         normalize_production_companies)
143     df['n_genres'] = df['genres'].apply(
144         lambda data: len(data.split(',')))
145     df['spoken_languages'] = df['spoken_languages'].apply(
146         normalize_spoken_languages)
147     df['keywords'] = df['keywords'].apply(normalize_keywords)
```

```

148 df['n_cast'] = df['cast'].apply(lambda data: len(data.split(',')))
149 df['n_crew'] = df['crew'].apply(lambda data: len(data.split(',')))
150 df['cast'] = df['cast'].apply(normalize_cast_crew, args=(elite_list_cast,))
151 df['crew'] = df['crew'].apply(normalize_cast_crew, args=(elite_list_crew,))
152 df[['Male', 'Female', 'Others']]
153     ] = dataframe['cast'].apply(normalize_gender)
154 df['release_date'] = df['release_date'].apply(
155     lambda data: 1 if data[5:7] in holiday_month else 0)
156 df['homepage'] = df['homepage'].apply(
157     lambda data: 0 if (data != data) else 1)
158 df['tagline'] = df['tagline'].apply(
159     lambda data: 0 if (data != data) else 1)
160 dummies_for_columns = ['genres', 'production_countries',
161                        'production_companies', 'spoken_languages']
162 for c in dummies_for_columns:
163     df = pd.concat([df.drop(c, axis=1), df[c].str.get_dummies(
164         sep=",", axis=1)], axis=1)
165
166     return df, dataframe
167 #####
168
169
170 if __name__ == "__main__":
171
172     df_train = pd.read_csv(sys.argv[1])
173     df_train, o_df_train = preprocess(df_train)
174     df_test = pd.read_csv(sys.argv[2])
175     df_test, o_df_test = preprocess(df_test)
176
177     columns_to_drop = ['revenue', 'movie_id', 'rating',
178                       'original_title', 'original_language', 'overview', 'status']
179
180 #####
181 #                                     PART - 1 REGRESSION
182
183     x_train = df_train.drop(columns_to_drop, axis=1).values
184     y_train_revenue = df_train['revenue'].values
185     x_test = df_test.drop(columns_to_drop, axis=1).values
186     y_test_revenue = df_test['revenue'].values
187     movie_ids = df_test['movie_id'].values
188
189     rfr_model = RandomForestRegressor(random_state=0)
190     rfr_model.fit(x_train, y_train_revenue)
191     y_predicted_revenue = rfr_model.predict(x_test)
192
193     msr = metrics.mean_squared_error(y_test_revenue, y_predicted_revenue)
194     pcc, _ = pearsonr(y_predicted_revenue, y_test_revenue)
195
196     pd.DataFrame({'movie_id': movie_ids, 'predicted_revenue': y_predicted_revenue}, columns=[
197         'movie_id', 'predicted_revenue']).to_csv(zid + '.PART1.output.csv', index=False)
198     pd.DataFrame([[zid, round(msr, 2), round(pcc, 2)]], columns=['zid', 'MSR', 'correlation']).to_csv(
199         zid + '.PART1.summary.csv', index=False)
200 #####
201
202 #####
203 #                                     PART - 2 CLASSIFICATION

```

```
204 df_train['runtime'] = o_df_train['runtime']
205 df_test['runtime'] = o_df_test['runtime']
206
207 x_train = df_train[['runtime', 'budget', 'cast', 'crew']].values
208 x_test = df_test[['runtime', 'budget', 'cast', 'crew']].values
209 y_train_rating = df_train['rating'].values
210 y_test_rating = df_test['rating'].values
211
212 gbc_classifier = GradientBoostingClassifier()
213 gbc_classifier.fit(x_train, y_train_rating)
214 y_predicted_rating = gbc_classifier.predict(x_test)
215
216 reports = metrics.classification_report(
217     y_test_rating, y_predicted_rating, output_dict=True)
218 average_precision = round(reports['macro avg']['precision'], 2)
219 average_recall = round(reports['macro avg']['recall'], 2)
220 accuracy = round(reports['accuracy'], 2)
221
222 pd.DataFrame({'movie_id': movie_ids, 'predicted_rating': y_predicted_rating}, columns=[
223     'movie_id', 'predicted_rating']).to_csv(zid + '.PART2.output.csv', index=False)
224 pd.DataFrame([zid, average_precision, average_recall, accuracy], columns=[
225     'zid', 'average_precision', 'average_recall', 'accuracy']).to_csv(zid + '.PART2.summary.csv')
226 #####
```