

```

1 #####
2
3 # This Python 3 environment comes with many helpful analytics libraries installed
4 # It is defined by the kaggle/python docker image:
  https://github.com/kaggle/docker-python
5 # For example, here's several helpful packages to load in
6
7 import numpy as np # linear algebra
8 import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
9
10 # Input data files are available in the "../input/" directory.
11 # For example, running this (by clicking run or pressing Shift+Enter) will list the
  files in the input directory
12
13 from pandas import DataFrame
14 from sklearn.model_selection import train_test_split
15
16
17 import os
18 #print(os.listdir("../input"))
19
20 # Any results you write to the current directory are saved as output.
21
22 train_d = pd.read_csv("train.csv")
23 test = pd.read_csv("test.csv")
24 #train.describe(include="all")
25 count = 1
26
27 combine = [train_d, test]
28 #smoke = train_d['smoker_status']
29 for train in combine:
30     train['smoker_status_1'] = train['smoker_status'].replace(['__', '11', ',', '.',
31         '_', '[]', ', ', '>?', '?', '>', 'quit?', 'non>', 'N?A', '..',
32         '??', '?', 'N?a', '>'], 'ND')
33
34
35 sex_age = train["sex and age"].str.split(',').tolist()
36 #print(np.array(sex_age))
37 age = []
38 sex = []
39 #df2 = pd.DataFrame(np.array(sex_age), columns = ["sex", "age"])
40 #print(sex_age[0:7])
41 #sex_age = sex_age[~np.isnan(sex_age)]
42 #sex_age.fillna([0,0])
43 for ages in sex_age:
44     try:
45         if pd.isnull(ages):
46             #print('NAN')
47             age.append(0)
48             sex.append(0)
49             continue
50         #print(ages[1])
51         if float(ages[1]):
52             #print("HIIII")
53             age.append(ages[1])
54             sex.append(ages[0])
55     except ValueError:
56         #print(f'Error {ages}')
57         try:
58             if float(ages[1]):
59                 #print("HIIII")
60                 age.append(ages[1])
61                 sex.append(ages[0])
62         except ValueError:
63             try:
64                 if float(ages[0]):
65                     age.append(ages[0])
66                     sex.append(ages[1])
67             except:

```

```

68         age.append(0)
69         sex.append(0)
70     #print(age)
71     #print(sex)
72     #DataFrame.from_records(sex_age)
73     #df = pd.DataFrame({'age':age,'sex':sex})
74     #df2 = pd.DataFrame(np.array(age), columns = ["age"])
75
76     age_1 = [float(i) for i in age]
77     #print(age_1)
78     train['age'] = age_1
79     train['sex'] = sex
80
81     js_la = train["job_status and living_area"].str.split('?',n=1,expand = True)
82     train['temp_js'] = js_la[0].str.lower()
83     #print(train['temp_js'].head(5))
84     train['temp_la'] = js_la[1].str.lower()
85     #print(train['temp_la'].head(5))
86     #train['temp_js'].head(2572)
87     #train['temp_la'].unique()
88
89
90     train['temp_js'] = train['temp_js'].replace(['remote','remotee','r'], 'r')
91     train['temp_js'] = train['temp_js'].replace(['city','c'], 'c')
92     train['temp_js'] = train['temp_js'].replace(['government','govt.'], 'g')
93     train['temp_js'] = train['temp_js'].replace(['private_sector', 'privattte',
94     'private', 'private sector'], 'p')
95     train['temp_js'] = train['temp_js'].replace(['biz', 'business_owner', 'business
96     owner'], 'b')
97     train['temp_js'] = train['temp_js'].replace(['unemployed', '', 'n.a'], 'u')
98     train['temp_js'] = train['temp_js'].replace(['leave', 'parental_leave', 'parental
99     leave'], 'l')
100
101     #print(train['temp_js'].head(5))
102
103     train['temp_la'] = train['temp_la'].replace(['remote','remotee','r'], 'r')
104     train['temp_la'] = train['temp_la'].replace(['city','c'], 'c')
105     train['temp_la'] = train['temp_la'].replace(['government','govt.'], 'g')
106     train['temp_la'] = train['temp_la'].replace(['private_sector', 'privattte',
107     'private', 'private sector'], 'p')
108     train['temp_la'] = train['temp_la'].replace(['biz', 'business_owner', 'business
109     owner'], 'b')
110     train['temp_la'] = train['temp_la'].replace(['unemployed', '', 'n.a'], 'u')
111     train['temp_la'] = train['temp_la'].replace(['leave', 'parental_leave', 'parental
112     leave'], 'l')
113
114
115     train['temp_la'] = train['temp_la'].replace(['null'], '')
116     train['temp_js'] = train['temp_js'].replace(['null'], '')
117     #print(train['temp_la'].head(5))
118
119     js = ['' for _ in range(len(train["job_status and living_area"]))]
120     la = ['' for _ in range(len(train["job_status and living_area"]))]
121
122     loc = ['r','c']
123     job = ['g','p','u','b','l']
124
125     #print(train['temp_la'][0])
126     for i in range(len(train["job_status and living_area"])):
127         if train['temp_la'][i] in loc:
128             la[i] = train['temp_la'][i]
129             #js[i] = train['temp_js'][i]
130         if train['temp_la'][i] in job:
131             #else:
132             #js[i] = train['temp_la'][i]
133             la[i] = train['temp_js'][i]
134         if train['temp_js'][i] in loc:
135             js[i] = train['temp_la'][i]
136             #la[i] = train['temp_js'][i]

```

```

131         if train['temp_js'][i] in job:
132             #la[i] = train['temp_la'][i]
133             js[i] = train['temp_js'][i]
134             '''if la[i] == 'p':
135                 print(train['temp_la'][i])
136                 print(train['temp_js'][i])
137                 pass'''
138         '''print(set(js))
139         print(set(la))
140         print(la.count('p'))
141         print(la.count('u'))'''
142         train['job'] = js
143         train['location'] = la
144
145
146         train['female'] = train['sex']
147         train['female'] = train['female'].replace(['F',' F', 'female', 'Female' , 'f'
148             , 'Other', ' Other', 'femalle', 'FEMALE' ],1)
149         train['female'] = train['female'].replace(['M',' M' , 'male' , 'MALE', 'm', 'mmale',
150             'MM', 'Male'],0)
151
152         train['male'] = train['sex']
153         train['male'] = train['male'].replace(['M',' M' , 'male' , 'MALE', 'm', 'mmale',
154             'MM', 'Male'],1)
155         train['male'] = train['male'].replace(['F',' F', 'female', 'Female' , 'f' , 'Other', '
156             Other', 'femalle', 'FEMALE' ],0)
157
158         train['remote'] = train['location']
159         train['remote'] = train['remote'].replace('r',1)
160         train['remote'] = train['remote'].replace(['c',' ', 'p'],0)
161
162         train['city'] = train['location']
163         train['city'] = train['city'].replace('c',1)
164         train['city'] = train['city'].replace(['r',' ', 'p'],0)
165
166         train['job'] = train['job'].replace(' ', 'u')
167
168         train['govt'] = train['job']
169         train['govt'] = train['govt'].replace('g',1)
170         train['govt'] = train['govt'].replace(['g', 'p', 'u', 'b', 'l'],0)
171
172         train['priv'] = train['job']
173         train['priv'] = train['priv'].replace('p',1)
174         train['priv'] = train['priv'].replace(['g', 'u', 'b', 'l'],0)
175
176         train['unemp'] = train['job']
177         train['unemp'] = train['unemp'].replace('u',1)
178         train['unemp'] = train['unemp'].replace(['g', 'p', 'b', 'l'],0)
179
180         train['business'] = train['job']
181         train['business'] = train['business'].replace('b',1)
182         train['business'] = train['business'].replace(['g', 'p', 'u', 'l'],0)
183
184         train['leave'] = train['job']
185         train['leave'] = train['leave'].replace('l',1)
186         train['leave'] = train['leave'].replace(['g', 'p', 'u', 'b'],0)
187
188         train['high_BP'] = train['high_BP'].fillna(0)
189         train['high_BP'] = train['high_BP'].replace('.',0)
190
191         train['average_blood_sugar'] = train['average_blood_sugar'].fillna(0)
192         train['BMI'] = train['BMI'].fillna(0)
193
194         train['heart_condition_detected_2017'] =
195         train['heart_condition_detected_2017'].fillna(0)
196         train['heart_condition_detected_2017'] =
197         train['heart_condition_detected_2017'].replace(['.', 'N?A', 'n.a', '.', ', ', 'n.a', 'N?A',
198             ''],0)

```

```

193
194 train['smoker_status_1'] = train['smoker_status_1'].fillna('ND')
195 train['smoker_status_1'] =
train['smoker_status_1'].replace(['.', '.', ', ', 'n.a', 'N?A', ''], 'ND')

196
197
198 #print(train["average_blood_sugar"].unique())
199 train['BMI'] = train['BMI'].replace(['?', '.', ', ', 'n.a', 'N?A', ''], 0)
200 train['BMI'] = train['BMI'].astype(float)
201 train['average_blood_sugar'] = train['average_blood_sugar'].astype(float)
202
203
204
205 #print(train['smoker_status_1'].unique())
206 # ['non-smoker' 'quit' 'active_smoker' nan 'ND']
207
208 train['active_smoker'] = train['smoker_status_1']
209 train['active_smoker'] = train['active_smoker'].replace('active_smoker', 1)
210 train['active_smoker'] =
train['active_smoker'].replace(['non-smoker', 'quit', 'active_smoker', 'ND'], 0)

211
212 train['non_smoker'] = train['smoker_status_1']
213 train['non_smoker'] = train['non_smoker'].replace('non-smoker', 1)
214 train['non_smoker'] =
train['non_smoker'].replace(['non-smoker', 'quit', 'active_smoker', 'ND'], 0)

215
216 train['quit'] = train['smoker_status_1']
217 train['quit'] = train['quit'].replace('quit', 1)
218 train['quit'] = train['quit'].replace(['non-smoker', 'quit', 'active_smoker', 'ND'], 0)
219
220 train['nd_smoker'] = train['smoker_status_1']
221 train['nd_smoker'] = train['nd_smoker'].replace('ND', 1)
222 train['nd_smoker'] =
train['nd_smoker'].replace(['non-smoker', 'quit', 'active_smoker', 'ND'], 0)

223
224
225 train['married'] = train['married'].replace('.', 0)
226 #print(train['age'].loc[pd.isnull(train['married'])])
227 #train['married'].loc[pd.isnull(train['married'])] = (lambda x: 1 if
train['age'].loc[pd.isnull(train['married'])] > 30 else 0)
228 train['married'].loc[pd.isnull(train['married'])] =
np.where(train['age'].loc[pd.isnull(train['married'])]>30, 1, 0)

229
230 from copy import deepcopy
231 import statistics
232
233 age_wo0 = deepcopy(age)
234 age_wo0[:] = (float(value) for value in age_wo0 if value != 0)
235
236 #print(set(sex))
237 #print(float(statistics.mean(age_wo0)))
238
239 train['age'] = train['age'].replace(0, float(statistics.mean(age_wo0)))
240
241 #print(float(statistics.mean(train['BMI'])))
242
243 train['BMI'] = train['BMI'].replace(0, float(statistics.mean(train['BMI']))+5.0)
244 train['average_blood_sugar'] =
train['average_blood_sugar'].replace(0, float(statistics.mean(train['average_blood_sugar'])))

245
246
247 train['TreatmentA'] = train['TreatmentA'].fillna(0)
248 train['TreatmentB'] = train['TreatmentB'].fillna(0)
249 train['TreatmentC'] = train['TreatmentC'].fillna(0)
250 train['TreatmentD'] = train['TreatmentD'].fillna(0)
251 train['TreatmentD'] = train['TreatmentD'].replace('0+E1860:E1868', 0)
252
253

```

```

254 train['high_BP'] = train['high_BP'].astype(float)
255 train['heart_condition_detected_2017'] =
256 train['heart_condition_detected_2017'].astype(float)
257 train['married'] = train['married'].astype(float)
258 train['average_blood_sugar'] = train['average_blood_sugar'].astype(float)
259 train['TreatmentA'] = train['TreatmentA'].astype(float)
260 train['TreatmentB'] = train['TreatmentB'].astype(float)
261 train['BMI'] = train['BMI'].astype(float)
262 train['TreatmentC'] = train['TreatmentC'].astype(float)
263 train['TreatmentD'] = train['TreatmentD'].astype(float)
264
265 train['age'] = train['age'].astype(float)
266 train['BMI'] = train['BMI'].astype(float)
267 train['female'] = train['female'].astype(float)
268 train['male'] = train['male'].astype(float)
269 train['remote'] = train['remote'].astype(float)
270 train['city'] = train['city'].astype(float)
271 train['govt'] = train['govt'].astype(float)
272 train['priv'] = train['priv'].astype(float)
273 train['unemp'] = train['unemp'].astype(float)
274 train['business'] = train['business'].astype(float)
275 train['leave'] = train['leave'].astype(float)
276 train['active_smoker'] = train['active_smoker'].astype(float)
277 train['non_smoker'] = train['non_smoker'].astype(float)
278 train['quit'] = train['quit'].astype(float)
279 train['nd_smoker'] = train['nd_smoker'].astype(float)
280
281 #for i in range(len(train['BMI'])):
282 #    BMI_SUM = (train['BMI'])
283 #print(train.columns)
284 #print(train['male'].unique())
285
286 if count == 1:
287     train.to_csv(r'D:\MS-UNSW\DataSoc_Datathon/pandas_train.csv')
288     count+=1
289
290 Drop = ['sex and age', 'job_status and
291 living_area', 'smoker_status', 'temp_js', 'temp_la', 'sex', 'job', 'location', 'smoker_statu
292 s_1']
293 train.drop(train[Drop], axis = 1, inplace = True)
294 #print(train.columns)
295
296 #train.to_csv(r'D:\MS-UNSW\DataSoc_Datathon/pandas.csv')
297 train.head()
298
299 train_d['stroke_in_2018'] = train_d['stroke_in_2018'].fillna(0)
300 train_d['stroke_in_2018'] =
301 train_d['stroke_in_2018'].replace(['.', ',', '?', 'n.a', 'nuLL', 'N?A', ''], 0)
302 train_d['stroke_in_2018'] = train_d['stroke_in_2018'].astype(float)
303
304 predictors = train_d.drop(['stroke_in_2018', 'id'], axis=1)
305 target = train_d["stroke_in_2018"]
306 x_train, x_val, y_train, y_val = train_test_split(predictors, target, test_size = 0.52,
307 random_state = 0)
308
309 # Perceptron
310 from sklearn.linear_model import Perceptron
311
312 perceptron = Perceptron()
313 perceptron.fit(x_train, y_train)
314 y_pred = perceptron.predict(x_val)
315 acc_perceptron = round(accuracy_score(y_pred, y_val) * 100, 2)
316 #print(acc_perceptron)
317
318 ids = test['id']

```

```
318 predictions = perceptron.predict(test.drop('id', axis=1))
319
320 #set the output as a dataframe and convert to csv file named submission.csv
321 output = pd.DataFrame({ 'id' : ids, 'stroke_in_2018': predictions })
322 output.to_csv('submission_per.csv', index=False)
323
324
```