

AI-Powered Socioeconomic Prediction of Lifespan

Rustom Ichhaporia [rustomi2@illinois.edu]*

2020-12-12

Abstract

How long will you live? This age-old question has extensive implications in the billions of risk estimations made by individuals planning for the future every day. Although never certain, a stronger approximation of an individual's lifespan can enable more reliable future planning and greater sense of stability

unfortunately when running the `fmin` optimization function, the program crashes after the first of 150 loops with the error:

5.1 Limitations

¹ One significant drawback of the approach taken to estimating

1 Background

2 Dataset Selection

3 Preprocessing

The appendix of this report contains the code for training the model and saving the results in a file. It does not include the code for statistical plots.

4 Modeling

5 Results

Last week, I registered for HAL access so that I could run the hyperparameter optimization script remotely because my computer overheated and could not run it for the appropriate number of trials. I was able to upload my files and run the first half of the script, but

*This research project was completed during my time as a research intern at the Illinois Risk Lab (<https://irisklabuiuc.wixsite.com/>) during the Fall of 2020. My research was a part of the AI-Powered Lifecycle Financial Planning project, which is still under development. I appreciate the help of Dr. Runhuan Feng, Dr. Frank Quan, Dr. Yong Xie, Dr. Linfeng Zhang, and my fellow interns throughout the process. Thank you!

¹<https://github.com/microsoft/LightGBM/issues/2696>

6 Appendix

6.0.0.1 This section will have three columns

script.py

```
1  '''Imports'''
2
3  import pandas as pd
4  import numpy as np
5  import seaborn as sns
6  import matplotlib.pyplot as plt
7
8  from sklearn.linear_model import LogisticRegression
9  from sklearn.linear_model import LogisticRegressionCV
10 from sklearn.model_selection import train_test_split
11 from sklearn.metrics import classification_report
12 from sklearn.metrics import roc_auc_score
13 from sklearn.metrics import precision_recall_curve
14 from sklearn.metrics import plot_precision_recall_curve
15 from sklearn.metrics import plot_roc_curve
16
17 from imblearn.over_sampling import SMOTE
18
19 sns.set()
20
21 '''Preprocessing'''
22
23 df_raw = pd.read_csv('data/11.csv')
24 print('Data successfully loaded.')
25
26 df_raw = df_raw.drop(columns=['smok100', 'agesmk', 'smokstat', 'smokhome', 'curruse', 'everuse'])
27
28 df_raw['indmort'] = df_raw['inddea'][(df_raw['inddea'] == 1) & (df_raw['indalg'] == 1)]
29 df_raw['indmort'] = df_raw['indmort'].fillna(0)
30
31 used_numerical = ['age', 'hhnum']
32 used_ordinal = ['povpct', 'adjinc']
33 used_categorical = ['stater', 'pob', 'sex', 'race', 'urban', 'smsast']
34 used_special = ['wt', 'indmort']
35
36 used_features = used_numerical + used_ordinal + used_categorical + used_special
37
38 df_raw = df_raw[used_features]
39
40 df_raw[used_categorical] = df_raw[used_categorical].astype('category')
```

```
41
42 df_raw = df_raw.dropna(axis=0)
43
44 df = pd.get_dummies(df_raw)
45
46 X = df.drop(columns=['indmort'])
47 y = df['indmort']
48
49 '''Sampling'''
50
51 X_train, X_test, y_train, y_test = train_test_split(X, y)
52
53 print('Proportion of data from minority class before SMOTE: ', y_train.sum() / y_train.count())
54 X_train, y_train = SMOTE().fit_resample(X_train, y_train)
55 print('Proportion of data from minority class after SMOTE: ', y_train.sum() / y_train.count())
56
57 '''Modeling'''
58
59 model = LogisticRegressionCV(scoring='roc_auc', random_state=0)
60 model.fit(X_train, y_train)
61 print(classification_report(model.predict(X_test.drop(columns=['indmort'])), y_test))
62
63 pred_probs = model.predict_proba(X_test.drop(columns=['wt', 'indmort']))
64
65 print(classification_report(np.round(pred_probs + 0.25), y_test))
```