

# unxx8hr4k

August 1, 2023

```
[7]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LogisticRegression
from sklearn.preprocessing import StandardScaler
```

```
[8]: from google.colab import drive
drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force\_remount=True).

```
[9]: df=pd.read_csv("/content/drive/MyDrive/mydatasets/C4_framingham.csv")
df
```

```
[9]:
```

	male	age	education	currentSmoker	cigsPerDay	BPMeds	\
0	1	39	4.0	0	0.0	0.0	
1	0	46	2.0	0	0.0	0.0	
2	1	48	1.0	1	20.0	0.0	
3	0	61	3.0	1	30.0	0.0	
4	0	46	3.0	1	23.0	0.0	
...	...	...	...	...	...	...	
4233	1	50	1.0	1	1.0	0.0	
4234	1	51	3.0	1	43.0	0.0	
4235	0	48	2.0	1	20.0	NaN	
4236	0	44	1.0	1	15.0	0.0	
4237	0	52	2.0	0	0.0	0.0	

  

	prevalentStroke	prevalentHyp	diabetes	totChol	sysBP	diaBP	BMI	\
0	0	0	0	195.0	106.0	70.0	26.97	
1	0	0	0	250.0	121.0	81.0	28.73	
2	0	0	0	245.0	127.5	80.0	25.34	
3	0	1	0	225.0	150.0	95.0	28.58	
4	0	0	0	285.0	130.0	84.0	23.10	
...	...	...	...	...	...	...	...	
4233	0	1	0	313.0	179.0	92.0	25.97	
4234	0	0	0	207.0	126.5	80.0	19.71	

4235	0	0	0	248.0	131.0	72.0	22.00
4236	0	0	0	210.0	126.5	87.0	19.16
4237	0	0	0	269.0	133.5	83.0	21.47

	heartRate	glucose	TenYearCHD
0	80.0	77.0	0
1	95.0	76.0	0
2	75.0	70.0	0
3	65.0	103.0	1
4	85.0	85.0	0
...	...	...	...
4233	66.0	86.0	1
4234	65.0	68.0	0
4235	84.0	86.0	0
4236	86.0	NaN	0
4237	80.0	107.0	0

[4238 rows x 16 columns]

```
[10]: df.head()
```

```
[10]:
```

	male	age	education	currentSmoker	cigsPerDay	BPMeds	prevalentStroke	\
0	1	39	4.0	0	0.0	0.0	0	
1	0	46	2.0	0	0.0	0.0	0	
2	1	48	1.0	1	20.0	0.0	0	
3	0	61	3.0	1	30.0	0.0	0	
4	0	46	3.0	1	23.0	0.0	0	

	prevalentHyp	diabetes	totChol	sysBP	diaBP	BMI	heartRate	glucose	\
0		0	0	195.0	106.0	70.0	26.97	80.0	77.0
1		0	0	250.0	121.0	81.0	28.73	95.0	76.0
2		0	0	245.0	127.5	80.0	25.34	75.0	70.0
3		1	0	225.0	150.0	95.0	28.58	65.0	103.0
4		0	0	285.0	130.0	84.0	23.10	85.0	85.0

	TenYearCHD
0	0
1	0
2	0
3	1
4	0

# 1 Data Cleaning and Data Preprocessing

```
[11]: df.dropna(inplace=True)
      df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 3656 entries, 0 to 4237
Data columns (total 16 columns):
#   Column                Non-Null Count  Dtype
---  -
0   male                  3656 non-null  int64
1   age                   3656 non-null  int64
2   education              3656 non-null  float64
3   currentSmoker          3656 non-null  int64
4   cigsPerDay              3656 non-null  float64
5   BPMeds                 3656 non-null  float64
6   prevalentStroke        3656 non-null  int64
7   prevalentHyp           3656 non-null  int64
8   diabetes               3656 non-null  int64
9   totChol                3656 non-null  float64
10  sysBP                  3656 non-null  float64
11  diaBP                  3656 non-null  float64
12  BMI                    3656 non-null  float64
13  heartRate              3656 non-null  float64
14  glucose                3656 non-null  float64
15  TenYearCHD             3656 non-null  int64
dtypes: float64(9), int64(7)
memory usage: 485.6 KB
```

```
[12]: df.describe()
```

```
[12]:
```

	male	age	education	currentSmoker	cigsPerDay	\
count	3656.000000	3656.000000	3656.000000	3656.000000	3656.000000	
mean	0.443654	49.557440	1.979759	0.489059	9.022155	
std	0.496883	8.561133	1.022657	0.499949	11.918869	
min	0.000000	32.000000	1.000000	0.000000	0.000000	
25%	0.000000	42.000000	1.000000	0.000000	0.000000	
50%	0.000000	49.000000	2.000000	0.000000	0.000000	
75%	1.000000	56.000000	3.000000	1.000000	20.000000	
max	1.000000	70.000000	4.000000	1.000000	70.000000	

  

	BPMeds	prevalentStroke	prevalentHyp	diabetes	totChol	\
count	3656.000000	3656.000000	3656.000000	3656.000000	3656.000000	
mean	0.030361	0.005744	0.311543	0.027079	236.873085	
std	0.171602	0.075581	0.463187	0.162335	44.096223	
min	0.000000	0.000000	0.000000	0.000000	113.000000	
25%	0.000000	0.000000	0.000000	0.000000	206.000000	

50%	0.000000	0.000000	0.000000	0.000000	234.000000
75%	0.000000	0.000000	1.000000	0.000000	263.250000
max	1.000000	1.000000	1.000000	1.000000	600.000000

	sysBP	diaBP	BMI	heartRate	glucose \
count	3656.000000	3656.000000	3656.000000	3656.000000	3656.000000
mean	132.368025	82.912062	25.784185	75.730580	81.856127
std	22.092444	11.974825	4.065913	11.982952	23.910128
min	83.500000	48.000000	15.540000	44.000000	40.000000
25%	117.000000	75.000000	23.080000	68.000000	71.000000
50%	128.000000	82.000000	25.380000	75.000000	78.000000
75%	144.000000	90.000000	28.040000	82.000000	87.000000
max	295.000000	142.500000	56.800000	143.000000	394.000000

	TenYearCHD
count	3656.000000
mean	0.152352
std	0.359411
min	0.000000
25%	0.000000
50%	0.000000
75%	0.000000
max	1.000000

```
[13]: df.columns
```

```
[13]: Index(['male', 'age', 'education', 'currentSmoker', 'cigsPerDay', 'BPMeds',
        'prevalentStroke', 'prevalentHyp', 'diabetes', 'totChol', 'sysBP',
        'diaBP', 'BMI', 'heartRate', 'glucose', 'TenYearCHD'],
        dtype='object')
```

```
[14]: feature_matrix = df.iloc[:,0:15]
      target_vector = df.iloc[:,-1]
```

```
[15]: fs = StandardScaler().fit_transform(feature_matrix)
      logr = LogisticRegression()
      logr.fit(fs,target_vector)
```

```
[15]: LogisticRegression()
```

```
[16]: observation=[[1,2,3,4,5,6,7,8,9,10,11,12,13,14,15]]
      prediction = logr.predict(observation)
      print(prediction)
```

```
[1]
```

```
[17]: logr.classes_
```

```
[17]: array([0, 1])
```

```
[18]: logit.predict_proba(observation)
```

```
[18]: array([[2.21478351e-04, 9.99778522e-01]])
```