

**A Project Report
on
“Predict diseases (heart attacks) with AI”
by**

Swarali Dinkar Borde

1. Introduction:

An estimated 17 million people die of CVDs (Cardiovascular disease), particularly heart attacks and strokes, in the world every year. Cardiac ailments killed more Indians in 2016 (28%) than any other non-communicable disease, said a new study published in the September 2018 issue of health journal, The Lancet. These are double the numbers reported in 1990 when heart disease caused 15% of deaths in India. Today we will try to build a heart attack predictor. Based on some diagnostically measured parameters we will predict who among the subjects under consideration, are on high risk of heart attack. This can revolutionize the healthcare system and help save many many lives.

2. Objective:

To build a heart attack predictor based on some diagnostically measured parameters.

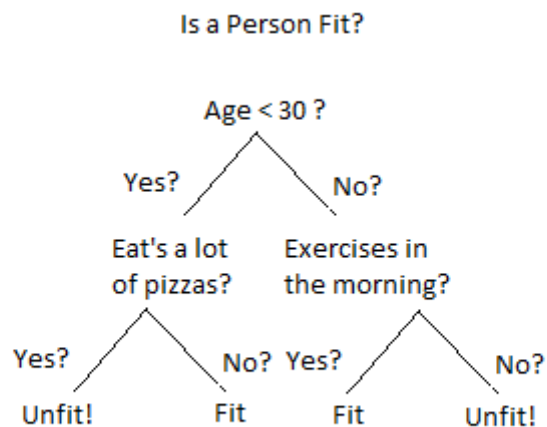
3. Problem Definition:

Predicting heart attacks using predictive analysis models of machine Learning like decision tree,logostic regression.

4.Related Theory:

4.1.Decision Tree:

Decision Trees are a type of Supervised Machine Learning (that is you explain what the input is and what the corresponding output is in the training data) where the data is continuously split according to a certain parameter. The tree can be explained by two entities, namely decision nodes and leaves. The leaves are the decisions or the final outcomes. And the decision nodes are where the data is split.



An example of a decision tree can be

explained using above binary tree. Let's say you want to predict whether a person is fit given their information like age, eating habit, and physical activity, etc. The decision nodes here are questions like 'What's the age?', 'Does he exercise?', 'Does he eat a lot of pizzas'? And the leaves, which are outcomes like either 'fit', or 'unfit'. In this case this was a binary classification problem (a yes no type problem).

4.1.1.Entropy:

Entropy, also called as Shannon Entropy is denoted by $H(S)$ for a finite set S , is the measure of the amount of uncertainty or randomness in data.

$$H(S) = \sum_{x \in X} p(x) \log_2 \frac{1}{p(x)}$$

Intuitively, it tells us about the predictability of a certain event. Example, consider a coin toss whose probability of heads is 0.5 and probability of tails is 0.5. Here the entropy is the highest possible, since there's no way of determining what the outcome might be. Alternatively, consider a coin which has heads on both the sides, the entropy of such an event can be predicted perfectly since we know beforehand that it'll always be heads. In other words, this event has **no randomness** hence its entropy is zero.

In particular, lower values imply less uncertainty while higher values imply high uncertainty.

4.1.2.Information Gain:

Information gain is also called as Kullback-Leibler divergence denoted by $IG(S,A)$ for a set S is the effective change in entropy after deciding on a particular attribute A . It measures the relative change in entropy with respect to the independent variables.

$$IG(S, A) = H(S) - H(S, A)$$

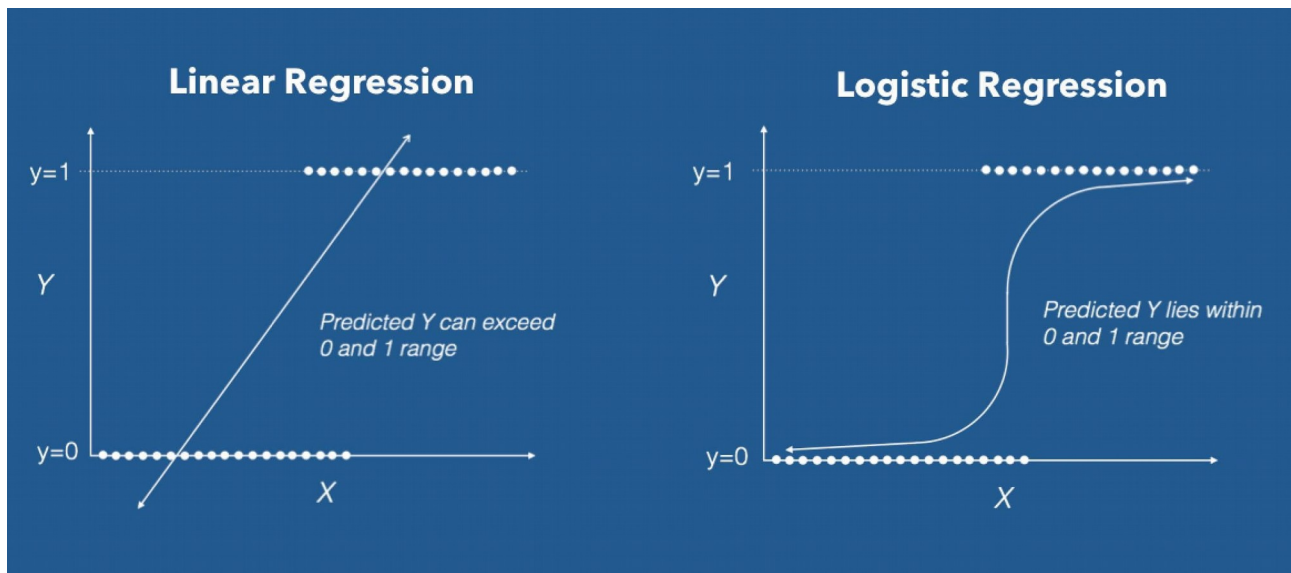
Alternatively,

$$IG(S, A) = H(S) - \sum_{i=0}^n P(x) * H(x)$$

where $IG(S, A)$ is the information gain by applying feature A . $H(S)$ is the Entropy of the entire set, while the second term calculates the Entropy after applying the feature A , where $P(x)$ is the probability of event x .

4.1.Logistic Regression:

Logistic Regression is a Machine Learning algorithm which is used for the classification problems, it is a predictive analysis algorithm and based on the concept of probability.



Logistic regression is a classification algorithm used to assign observations to a discrete set of classes. Some of the examples of classification problems are Email spam or not spam, Online transactions Fraud or not Fraud, Tumor Malignant or Benign. Logistic regression transforms its output using the logistic sigmoid function to return a probability value.

Activities

LibreOffice Calc

Mon Nov 4, 6:41 PM

100%

heart_disease_data.csv - LibreOffice Calc

FileEditViewInsertFormatStylesSheetDataToolsWindowHelp

LibreOffice Calc icons

LibreOffice Calc icons

Search: Liberation Sans, 10, Bold, Italic, Underline, Paragraph, Styles, Lists, Tables, Grid, Print, Zoom, etc.

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V
age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target								
63	1	3		145	233	1	0	150	0	2.3	0	0	1	1							
37	1	2		130	250	0	1	187	0	3.5	0	0	2	1							
41	0	1		130	204	0	0	172	0	1.4	2	0	2	1							
56	1	1		120	236	0	1	178	0	0.8	2	0	2	1							
57	0	0		120	354	0	1	163	1	0.6	2	0	2	1							
57	1	0		140	192	0	1	148	0	0.4	1	0	1	1							
56	0	1		140	294	0	0	153	0	1.3	1	0	2	1							
44	1	1		120	263	0	1	173	0	0	2	0	3	1							
52	1	2		172	199	1	1	162	0	0.5	2	0	3	1							
57	1	2		150	168	0	1	174	0	1.6	2	0	2	1							
54	1	0		140	239	0	1	160	0	1.2	2	0	2	1							
48	0	2		130	275	0	1	139	0	0.2	2	0	2	1							
49	1	1		130	266	0	1	171	0	0.6	2	0	2	1							
64	1	3		110	211	0	0	144	1	1.8	1	0	2	1							
58	0	3		150	283	1	0	162	0	1	2	0	2	1							
50	0	2		120	219	0	1	158	0	1.6	1	0	2	1							
58	0	2		120	340	0	1	172	0	0	2	0	2	1							
66	0	3		150	226	0	1	114	0	2.6	0	0	2	1							
43	1	0		150	247	0	1	171	0	1.5	2	0	2	1							
69	0	3		140	239	0	1	151	0	1.8	2	2	2	1							
59	1	0		135	234	0	1	161	0	0.5	1	0	3	1							
44	1	2		130	233	0	1	179	1	0.4	2	0	2	1							
42	1	0		140	226	0	1	178	0	0	2	0	2	1							
61	1	2		150	243	1	1	137	1	1	1	0	2	1							
40	1	3		140	199	0	1	178	1	1.4	2	0	3	1							
71	0	1		160	302	0	1	162	0	0.4	2	2	2	1							
59	1	2		150	212	1	1	157	0	1.6	2	0	2	1							

Find heart_disease_data

☐ Find All☐ Formatted Display☐ Match Case

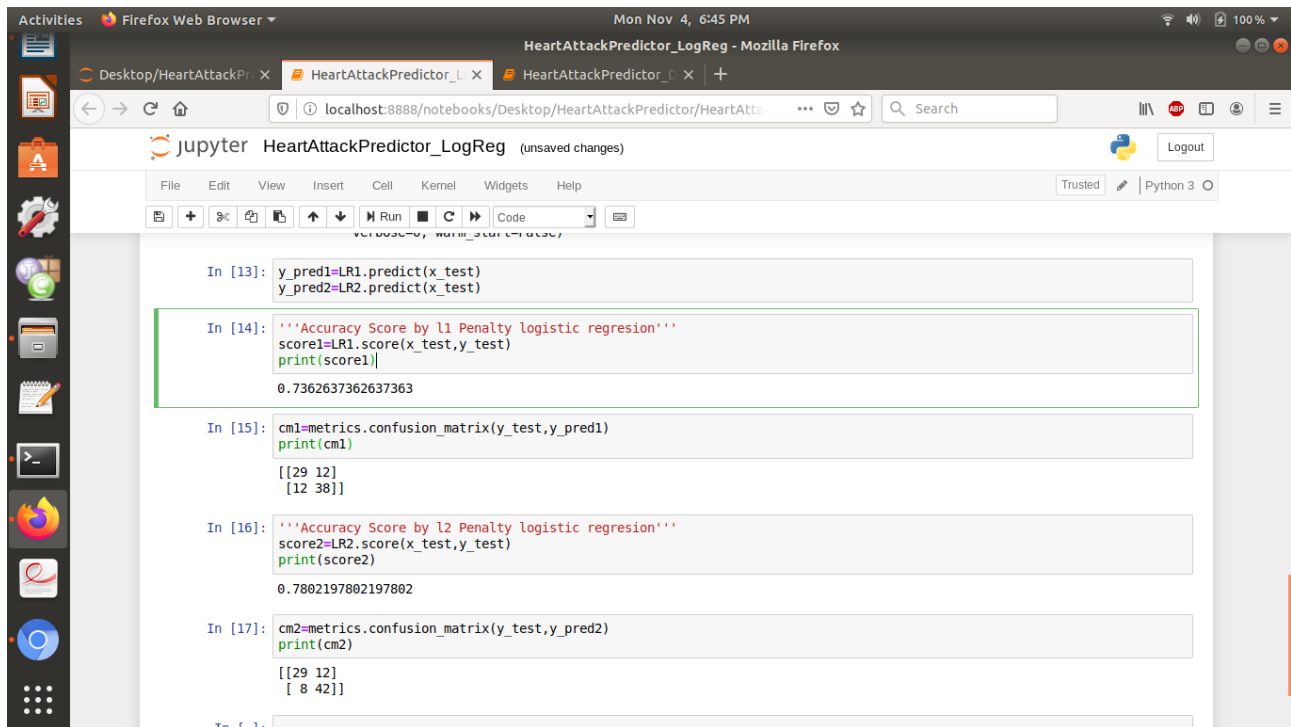
Sheet 1 of 1

DefaultEnglish (India)

Average: Sum: 0

6. Test Results and Analysis:

6.1 Results using Logistic Regrssion:



The screenshot shows a Jupyter Notebook interface in a Firefox browser. The notebook is titled "HeartAttackPredictor_LogReg" and is running on a local server at localhost:8888. The notebook contains several code cells that perform logistic regression analysis. The first cell (In [13]) shows the prediction of y values using LR1 and LR2 models. The second cell (In [14]) shows the accuracy score for LR1, which is 0.7362637362637363. The third cell (In [15]) shows the confusion matrix for LR1, which is [[29 12], [12 38]]. The fourth cell (In [16]) shows the accuracy score for LR2, which is 0.7802197802197802. The fifth cell (In [17]) shows the confusion matrix for LR2, which is [[29 12], [8 42]].

```
In [13]: y_pred1=LR1.predict(x_test)
y_pred2=LR2.predict(x_test)

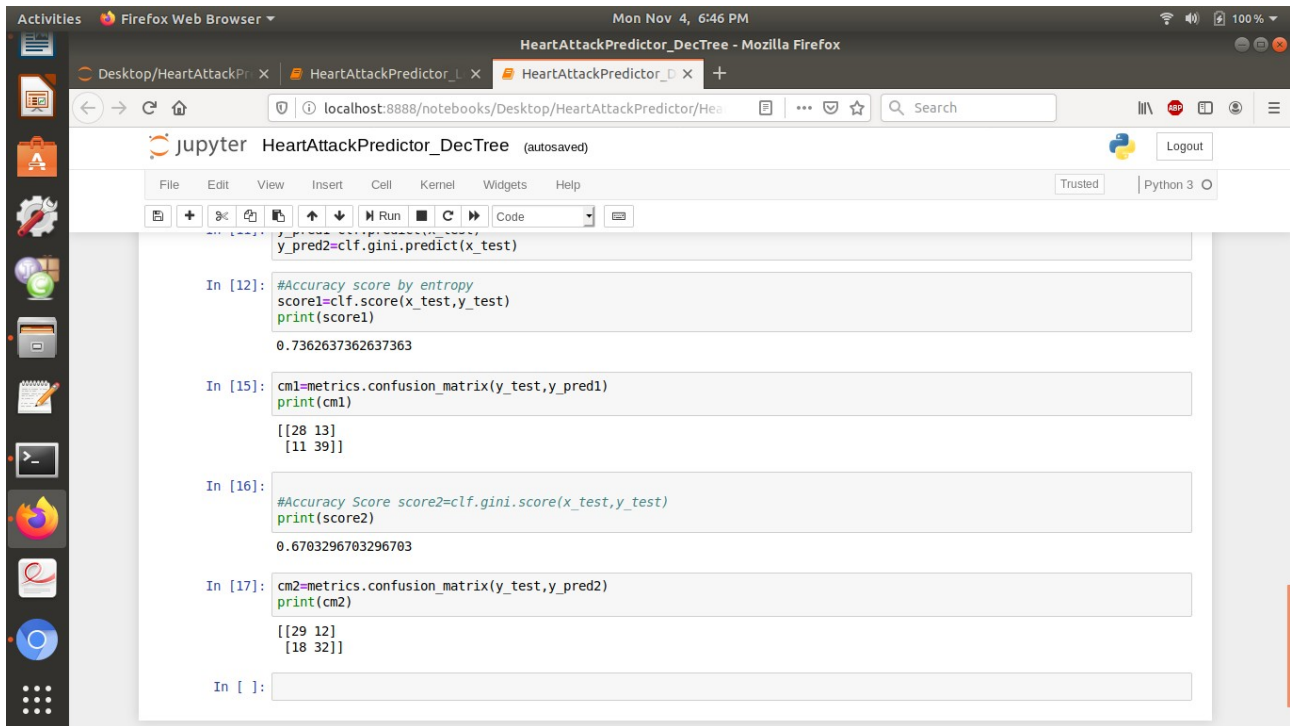
In [14]: '''Accuracy Score by l1 Penalty logistic regresion'''
score1=LR1.score(x_test,y_test)
print(score1)
0.7362637362637363

In [15]: cm1=metrics.confusion_matrix(y_test,y_pred1)
print(cm1)
[[29 12]
 [12 38]]

In [16]: '''Accuracy Score by l2 Penalty logistic regresion'''
score2=LR2.score(x_test,y_test)
print(score2)
0.7802197802197802

In [17]: cm2=metrics.confusion_matrix(y_test,y_pred2)
print(cm2)
[[29 12]
 [ 8 42]]
```

6.2: Results using Decision Tree Classifier:



```
HeartAttackPredictor_DecTree (autosaved)
File Edit View Insert Cell Kernel Widgets Help Trusted Python 3
y_pred2=clf.gini.predict(x_test)

In [12]: #Accuracy score by entropy
score1=clf.score(x_test,y_test)
print(score1)
0.7362637362637363

In [15]: cm1=metrics.confusion_matrix(y_test,y_pred1)
print(cm1)
[[28 13]
 [11 39]]

In [16]: #Accuracy Score score2=clf.gini.score(x_test,y_test)
print(score2)
0.6703296703296703

In [17]: cm2=metrics.confusion_matrix(y_test,y_pred2)
print(cm2)
[[29 12]
 [18 32]]

In [ ]:
```

7.Conclusion:

Thus I have successfully predicted the risk of heart attack using decision tree and logistic regression. Logistic regression with L2 penalty gave highest accuracy of 78%.