

Logistic Regression (Module -9)

Instructions

Please share your answers filled inline in the word document. Submit Python code and R code files wherever applicable.

Please ensure you update all the details:

Name: Nukala Ayyappa Bharthwaj Batch Id: DSWDMCOS 21012022

Topic: Logistic Regression.

- 1. Business Problem
 - 1.1. Objective
 - 1.2. Constraints (if any)
- 2. Work on each feature of the dataset to create a data dictionary as displayed in the below image:

Name of Feature	Description	Туре	Relevance
ID	Customer ID	Quantitative, Nominal	Irrelevant, ID does not provide useful information
	8		

2.1 Make a table as shown above and provide information about the features such as its Data type and its relevance to the model building, if not relevant provide reasons and provide description of the feature.

Using R and Python codes perform:

- 3. Data Pre-processing
 - 3.1 Data Cleaning, Feature Engineering, etc.
 - 3.2 Outlier Imputation
- 4. Exploratory Data Analysis (EDA):
 - 4.1. Summary
 - 4.2. Univariate analysis
 - 4.3. Bivariate analysis
- 5. Model Building
 - 5.1 Build the model on the scaled data (try multiple options)
 - 5.2 Perform Logistic Regression model.
 - 5.3 Train and Test the data and compare accuracies by Confusion Matrix, plot ROC



AUC curve.

6. Briefly explain the model output in the documentation. Share the benefits/impact of the solution - how or in what way the business (client) gets benefit from the solution provided.

Note:

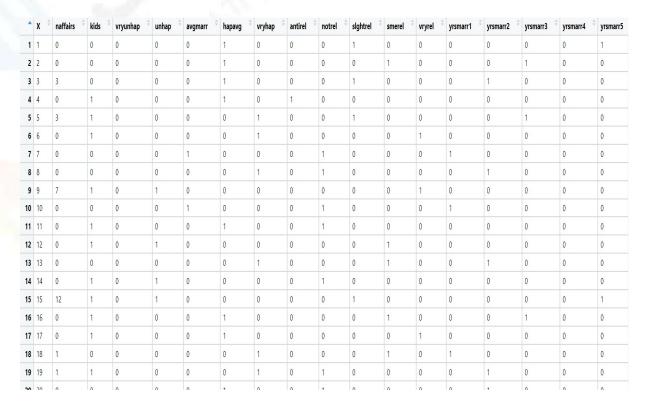
The assignment should be submitted in the following format:

- R code
- Python code
- Code Modularization should be maintained
- Documentation of the model building (elaborating on steps mentioned above)

Problem Statement: -

A psychological study has been conducted by a team of students at a University on married couples to determine the cause and effect on their married life and why they tend to have an extra marital affair, they have surveyed and collected a sample of data on which they would like to do further analysis to improve the relationship bond between couple, is it even possible to do so? Using your skills of Machine Learning apply Logistic Regression Model on the data and correctly classify whether a given person will have an affair or not given the set of attributes.

Convert naffairs column to Discreet Binary before proceeding with algorithm.





Sol:

Business Objective: To predict whether the person will have an external affair or not by using logistic regression model.

Constraints: Lack of analysis of affairs data of the people. **Data Types:** the data given and its types are as follows:

Name of feature	description	Data Type	Relevance
X	study number	Ordinal	Irrelevant since it is case number
Naffairs	Number of affairs a person have	Nominal	Relevant
Kids	Weather they have kids or not	Nominal	Relevant
Vryunhap	Weather they are very un happy or not	Nominal	Relevant
Unhap	Weather they are unhappy or not	Nominal	Relevant
Avgmarr	Weather avgmarr or not	Nominal	Relevant
Hapavg	Weather Hapavg or not	Nominal	Relevant
Vryhap	Weather they are very happy or not	Nominal	Relevant
antirel	Weather they are antirel or not	Nominal	Relevant
Notrel	Weather they are notrel or not	Nominal	Relevant
Slghtrel	Weather they are Sightrel or not	Nominal	Relevant
Smerel	Weather they are smerel or not	Nominal	Relevant
vryrel	Weather they are vryrel or not	Nominal	Relevant
yrsmarr1	Numbers of years married	Nominal	Relevant
Yrsmarr2	Numbers of years married	Nominal	Relevant
Yrsmarr3	Numbers of years married	Nominal	Relevant
Yrsmarr4	Numbers of years married	Nominal	Relevant
Yrsmarr5	Numbers of years married	Nominal	Relevant
Yrsmarr6	Numbers of years married	Nominal	Relevant

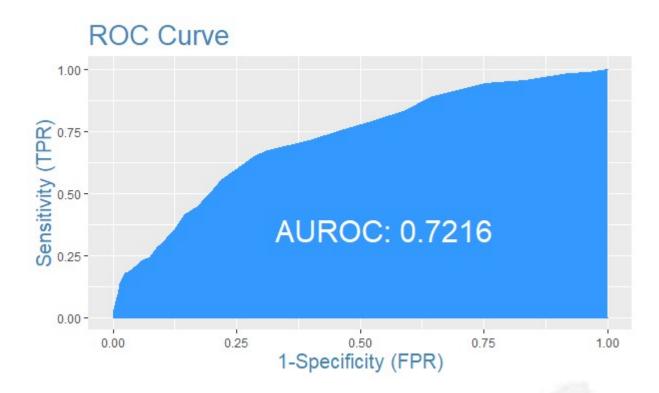
Data Pre-Processing: The number of affairs column in the given data is in the range of 0-16 and the same is converted into binary data. If number of affairs is >0 then they are treated as 1 and the remaining are treated as 0.

Since there are no missing values in the given data the same data can be used further to do analysis.

Logistic Regression: Initially linear regression model is applied to see the predicted values of the output and then the Generalized Regression model is applied to get the correct predicted values. Accuracy of the model is 78.06% and the cut of value is 56%



and the ROC curve the model is as follows:



ROC Curve for affairs data set

Problem Statement: -

In this time and age of widespread internet usage, effective and targeted marketing play a vital role, a marketing company would like to develop their strategy by analyzing their customer data and how effectively they can do targeted marketing, for this historical data has been collected of users clicking on ad given different factors such as age, location, time of activity and more. Perform Logistic Regression on the given data and classify the user who click's on ad's and who does not click on ad.

Dally Time Sport on Age		Assitions:	Daily Franse Diagra	Ac Topic Uni	Ob.	Niki	Coursy	Tinetano	Closed or Ac
885	x	81833.9	211 21	Cored Strawards	wyten		I Table	2116432112311	
9022	36	(8418)	192.79	Monitored national of	Neg Josi		1 Natu	20164-W 12936	
8947	25	5878534	235.1	Oparic tutore-live	Covidor		1 San Marino	201403-0120304	
74,15	25	540018	201	Train-Authoric recip	Neg Ter Lt		1 Mb	201641-1022119	
88.37	X	13989.99	285	Rhalisjáciá	Sout Farual		1 belot	2016/03/18 2:29:18	
\$935	25	58781.95	29.74	Statolecterion	Janiebeg		1 Novey	2019/05/09 14:33 17	
8891	20	5395285	23131	Etharous decicates	Bardysal		1 Nyone	en e	
68	4	266338	19170	Reactive local chall	Por Jefferja. y		1 Autolia	2016/2018 1:00:15	
74.52	х	\$980	22151	Configurable colum	WeeGsin		1 Gerats	201641-16153142	
99.08	x	5594232	10.0	Mandatay hamoger	Raniston		1 Ghana	2016/7/11/12/31	
658	4	4552251	122 82	Contractors routed o	Wellerdorte		1 Oran	2019/01/01 20:00 2	
8507	35	62-6131	28 8	Transactivities gibbs	East Transcook's		1 Baard	20164506 8:00 10	
19.57	4	5192830	10.0	Consider carbon	NotKalata.		1 Egot	2016/09/05 12:14:41	
79.52	28	55739.83	214 21	Sympisis feet this	ton fan		1 Boris av Heorga	2115 M &1 25 H Z	
428	35	30676	14356	Gasinotica way	WesWilliam		E Batistis	20164320 5:21:46	
8345	25	62162.25	141.6	Presidencementel	Son Travistiva		1 Sprin	201643-06-2-01:00	
55.16	37	23983 96	128.41	Situr oka iraz. 0	Wet Dyarting		1 Faledaia Terby	2015/01/2015/2014	
12/2	48	71581.38	107.53	latitive operations	Poittouti		1 Alpentan	2016454675058	
502	36	3198754	193	Commoderations	Josephia		1 Britis Indiar Coun	2016/2/17:53:56	
74.58	40	2916172	135.51	Acres at 207 pad	Niktor		1 Rusies Federalian	SHEEDER HOOF	
77.22	x	6400233	29.6	Chierdonic respo	Por Jacquel so		1 Cancern	2016/01/06 7:52:48	
84.56	X	6014537	2315	Steamined for 40	Laie Nook		1 Cancern	2019/05/05 10:22 3	
41.46	52	22635.7	16431	Mandatary dishlars	South John		t bard	2016/53/E-00 10	
17.25	×	6192872	23131	Figure onchol metr	Pangrood		1 5000	2012/2012	



Sol:

Business Objective: To predict whether who clicks on add or not by using logistic

regression model.

Constraints: Lack of analysis of advertisement data of the users.

Data Types: the data given and its types are as follows:

Name of feature	description	Data Type	Relevance
Age	Age of the user	Ratio	Relevant
Area_Income	Area income of the particular add	Ratio	Relevant
Daily Internet Usage	Daily internet usage of the user	Ratio	Relevant
Ad_Topic_Line	Topic line of the add	Nominal	Relevant
City	City of the user	Nominal	Relevant
Male	Whether the user is male or not	Nominal	Relevant
Country	Country of the user	Nominal	Relevant
Timestamp	Time stamp of the add	Ratio	Relevant
Clicked_on_Ad	Whether the custome seed the add	Nominal	Relevant
	or not		

Data Pre-Processing: Catagorical data is converted into numeric data so that it will be used for doing the analysis.

Logistic Regression: Initially linear regression model is applied to see the predicted values of the output and then the Generalized Regression model is applied to get the correct predicted values. Accuracy of the model is 86%. All the values are predicted corrected based on the given variables if the variables of the data changes then there will be change in the accuracy of the data.

Problem Statement: -

Prediction of election results has become trivial in these days, the outcome variable is (0/1) and the other factors that affect a candidate win or loss is amount of money spent, popularity and more. Perform Logistic Regression on the dataset and classify the candidates.

Sol:

Business Objective: To predict whether the candidate will win or lose the in election by using logistic regression model.

Constraints: Lack of analysis of election data of the members.

Data Types: the data given and its types are as follows:

Name of feature	description	Data Type	Relevance
Election-id	Id of the member	Ordinal	Irrelevant since it is I'd
		100	number
Result	Result for the member whether	Nominal	Relevant
	won or not		
Year	Age of the member	Nominal	Relevant
Amount Spent	Amount spent by the member	Nominal	Relevant



Popularity Rank	Popularity rank of the member	Nominal	Relevant
. opanancy mann	. openancy rank or the member		1.0.0.0.0

Data Pre-Processing: Since election id is not useful for the analysis that column is removed and the data is used for further analysis.

Logistic Regression: Initially linear regression model is applied to see the predicted values of the output and then the Generalized Regression model is applied to get the correct predicted values. Accuracy of the model is 100%. All the values are predicted corrected based on the given variables if the variables of the data changes then there will be change in the accuracy of the data.

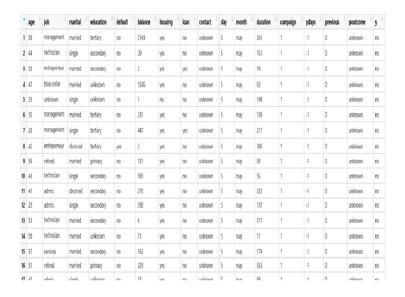
•	Election.id [‡]	Result [‡]	Year [‡]	Amount.Spent *	Popularity.Rank
1	NA	NA	NA	NA	NA
2	122	0	32	3.81	3
3	315	1	48	6.32	2
4	201	1	51	3.67	1
5	965	0	40	2.93	4
6	410	1	52	3.60	1
7	150	0	35	4.20	4
8	743	1	39	5.66	2
9	612	1	42	4.32	3
10	206	1	44	3.26	3
11	792	0	50	4.52	4

Problem Statement:

In Financial Institutions getting their customers to do a fixed deposit in the banks is a vital and at most important for the bank as they bank uses it and pays an interest amount to those deposited customers. To ask every customer for a term deposit is not viable as well as time consuming process, can you come up with a Logistic Regression model to predict customers who will do a term deposit or not.

The output variable in the dataset is Y which is discreet and binary. Snapshot of the dataset is given below.





Sol:

Business Objective: To predict the customer who will do term deposit or not by using

logistic regression model.

Constraints: Lack of analysis of Customer data.

Data Types: the data given and its types are as follows:

Name of feature	description	Data Type	Relevance
Age	Age of the customer	Ratio	Relevant
Default	Whether customer is default or not	Nominal	Relevant
Balance	Balance of the customer account	Ratio	Relevant
Housing	Whether the customer is housing or not	Nominal	Relevant
Loan	Whether the customer is having loan or not	Nominal	Relevant
Duration	Duration of the loan	Ratio	Relevant
Campaign	Whether the customer is campaign or not	Nominal	Relevant
Pdays	Whether the customer is paid in days or not	Nominal	Relevant
Previous	Whether the customer is is previous or not	Nominal	Relevant
poutfailure	Whether the customer is poutfailure or not	Nominal	Relevant
poutother	Whether the customer is poutother or not	Nominal	Relevant
poutsuccess	Whether the customer is poutsuccess or not	Nominal	Relevant
poutunknown	Whether the customer is poutunknown or not	Nominal	Relevant
con_cellular	Whether the customer is having coc_cellur or not	Nominal	Relevant
con_telephone	Whether the customer is having con_telph or not	Nominal	Relevant
con_unknown	Whether the customer is having unk_cont or not	Nominal	Relevant
divorced	Whether the customer is divorced or not	Nominal	Relevant
married	Whether the customer is married or not	Nominal	Relevant
single	Whether the customer is single or not	Nominal	Relevant
joadmin.	Whether the customer job is admin or not	Nominal	Relevant
joblue.collar	Whether the customer job is blue.col or not	Nominal	Relevant
joentrepreneur	Whether the customer job is entrepreneur or not	Nominal	Relevant
johousemaid	Whether the customer job is house maid or not	Nominal	Relevant
jomanagement	Whether the customer job is mang. Or not	Nominal	Relevant
joretired	Whether the customer job is retired or not	Nominal	Relevant
joself.employed	Whether the customer job is employed or not	Nominal	Relevant
joservices	Whether the customer job is service or not	Nominal	Relevant
jostudent	Whether the customer student or not	Nominal	Relevant



jotechnician	Whether the customer job is technician or not	Nominal	Relevant
jounemployed	Whether the customer job is unemployed or not	Nominal	Relevant
jounknown	Whether the customer job is unknown or not	Nominal	Relevant
У	Output variable whether the customer paid or	Nominal	Relevant
	not		

Data Pre-Processing: All the given data is numeric some of them are continues and some are binary data, it is used for doing the further analysis.

Logistic Regression: Initially linear regression model is applied to see the predicted values of the output and then the Generalized Regression model is applied to get the correct predicted values. Accuracy of the model is 90.01% and the cut of value is 46.9% and the ROC curve the model is as follows:

