zeroc<>de learning

Learning Data Analytics Made Easy

USER GUIIDE

BINARY LOGIT ANALYSIS



1.	MODEL- BINARY LOGIT ANALYSIS
2.	ALL ABOUT LEFT PANEL
3.	DATA INPUT AND OVERVIEW TAB
4.	DATA SUMMARY TAB
5.	DATA EXPLORATION TAB
6.	SUMMARY LOGIT TAB
7.	PREDICTED PROB VS ACTUAL TAB
8.	ACCURACY & ROC TAB

BINARY LOGIT REGRESSION ANALYSIS

Binary logit Regression, Binary logistic regression (LR) is a regression model where the target variable is binary, that is, it can take only two values, 0 or 1.Logistic regression is easier to implement, interpret, and very efficient to train. If the number of observations is lesser than the number of features, Logistic Regression should not be used, otherwise, it may lead to overfitting. It makes no assumptions about distributions of classes in feature space. It is named 'Logistic Regression' because its underlying technique is quite the same as Linear Regression. The term "Logistic" is taken from the Logit function that is used in this method of classification.



LEFT PANEL

(INPUT AREA)

Browse	No file selected
Data Selec	ction
Apply Chan	ges
Advance (Options

	Overview	Data Summary	Data Exploration	Summary Logit	
	Predicted Pro	obabilty vs. Actual	Accuracy, ROC &	AUC Prediction New Data	
1	In statistics, the logistic regression (or binary logit model) is used to model the probability of a certain class or event existing such as pass/fail, win/lose, alive/dead or healthy/sick.				
	This application panel on the let csv file (comma in csv format ar	n requires one data ft) and upload the c a delimited file), so i nd then proceed. M	input from the user. sv data input file. Not f you don't have csv i ake sure you have top	To do so, click on 'Browse' (te that this application can read input data file, first convert you o row as variable names.	n the I only data
J	Once csv file is uploaded successfully, variables in the data file will reflect in the 'Data Selection' panel on the left. Now you can select dependent variable (Y Variable) from drop- down menu. By default all other remaining variables will be selected as explanatory variables (X variables). If you want to drop any variable from explanatory variables, just uncheck that variable and it will be dropped from the model analysis. If any of the explanatory variables is a factor variable, you can define that variable as factor variable by selecting that variable in the last list of variables.				
	Binary logit cla between two cl	assification model asses (0 and 1 outc	trains better when omes).	observations are equal distri	outed

OPERATIONAL

ANALYSIS TAB (MAIN PANEL)

LEFT PANEL (INP)



DATA INPUT (UPLOADING DATASET)

- Click on browse
- Select the datafile that is in the form of csv format.(Ex program.csv)
- Browse the file and select the data to train your model for prediction.
- Top rows of the dataset should be of 'variable names'.

Data Exploration and Descriptive Statistics

ø			
Data Input			
Note: input d	ata should be in csv format		
Upload inpu	Upload input data (csv file with header)		
Browse	Dataset.csv		
	Dataset.csv		

OVERVIEW TAB

This tab provides you with relevant study resources, tutorials, sample datasets and a short overview to start with, which helps you understand and comprehend your data correctly. This tab also provides you the basic idea about regression analysis, gives sample data and provides the description about regression.



DATA SUMMARY TAB

It is very important to understand our data completely to infer meaningful insights and to get an overview of all the data points as a whole, but it is quite impossible to analyze thousand data points manually.

The 'Data Summary' option enables you to get a comprehensive evaluation through statistical measures that help us form the basis of our analysis. It will display all the 'descriptive analytics' measures including mean, median, standard deviation, variance etc. for all the data variables present in the dataset. we can review the uploaded data and the contents of it, A brief summary of the data can be seen it includes range of data values, minimum and maximum value missing and null values etc.

Data Summary of Selected Y and X Varaibles	
Note: maximum 2,000 observations randomly selected, see advance options in the panel on the left.	minimum value
\$Dimensions [1] 30 2 \$Summary \$Summary\$Numeric.data YearExperience Salary min 1.1000 37731.00 max 10.5000 122391.00 range 9.4000 84660.00 median 4.7000 65237.00 mean 5.3133 76003.00 var 8.0536 751550960.41 std.dev 2.8379 27414.43 \$Summary\$factor.data NULL	maximum value , range between data values ,mean ,median ,mode with standard deviation that is the terms of statistics
Note on Scientific Notation (-3.21e+02 = -3.21x10^2 = -321) - Wikipedia 'data.frame': 30 obs. of 2 variables: \$ YearsExperience: num 1.1 1.3 1.5 2 2.2 2.9 3 3.2 3.2 3.7	
\$ Salary : num 39343 46205 37731 43525 39891	
Missing Data Rows (Sample) Note: to impute or drop missing values (if any) check advance options in the panel on the left.	Info about missing values
\$MissingDataRows [1] YearsExperience Salary	

It also segregates dataset variables into respective data types, such as integer, whole numbers, character etc.



Use the left panel to transform selected variables as per the requirement of analysis , correspondingly the data summary will also change.

DATA EXPLORATION TAB

Exploration Data is the representation of data through use of common graphics, such charts, plots, infographics, as animations. These even and visual displays of information communicate complex data relationships and data-driven insights in a way that is easy to understand.

1)HISTOGRAM: A histogram is a graph that shows the frequency of numerical data using rectangles.

The height of a rectangle (the vertical

axis) represents the distribution frequency of a variable (the amount, or how often that variable appears). Histograms give a rough sense of the density of of the underlying distribution the data. and often for density estimation, estimating the probability density function of the underlying variable. The total area of a histogram used for probability density is always normalized to 1. If the length of the intervals on the x-axis are all 1, then a histogram is identical to a relative frequency plot.

PAIR PLOTS : A pairplot plot a pairwise relationships in a dataset. The pairplot function creates a grid of Axes such that each variable in data will by shared in the y-axis across a single row and in the x-axis across a single column. That creates plots as shown above Pair plot is used to understand the best set of features to explain a relationship between two variables or to form the most separated clusters. It also helps to form some simple classification models by drawing some simple lines or make linear separation in our dataset.



05

2)CORRELATION TABLE : A two-way tabulation of the relations between correlates; row headings are the scores on one variable and column headings are the scores on the second variables and a cell shows how many times the score on that row was associated with the score in that column.

A correlation matrix is a table showing correlation coefficients between variables. Each cell in the table shows the correlation between two variables. A correlation matrix is used to summarize data, as an input into a more advanced analysis, and as a diagnostic for advanced analyses.

Example: A positive correlation is a relationship between two variables in which both variables move in the same direction. Therefore, when one variable increases as the other variable increases, or one variable decreases while the other decreases. An example of positive correlation would be height and weight.

Use the left panel to modify/deal with the outliers identified here.

SUMMARY LOGIT TAB



The summary tab provides the overall analysis result. The residual deviance tells us how well the response variable can be predicted by a model with p predictor variables. The lower the value, the better the model is able to predict the value of the response variable. Pearson residuals are defined as the standardized distances between the observed and expected responses, and deviance residuals are defined as the signed square root of the individual contributions to the model deviance. If we use a generalized linear model (GLM) to model the relationship, deviance is a measure of goodness of fit: the smaller the deviance, the better the fit. An important assumption of logistic regression is that the errors (residuals) of the model are approximately normally distributed. The observed values on the response variable cannot be normally distributed themselves, because Y is binary.

CORRELATION PLOT : The correlation coefficient is the specific measure that quantifies the strength of the linear relationship between two variables in a correlation analysis. The coefficient is what we symbolize with the r in a correlation report. A correlation analysis provides information on the strength and direction of the linear relationship between two variables, while a simple linear regression analysis estimates parameters in a linear equation that can be used to predict values of one variable based on the other.

Note on Logisti	c Regression	- Wikipedia			
Nodel predicts	log-odds of 2	X = 1			
Call:					
glm(formula	= formula,	family = bin	nomial,	data =	mydata()
Deviance Re:	siduals:				
Min	10	Median		30	Ma
-2.333e-04	-2.100e-08	-2.100e-08	-2.100)e-08	2.518e-04
Coefficients	s: (2 not de	fined becau	se of si	ingulari	ties)
	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	2.625e+01	3.570e+05	0.000	1.00	0
PassengerId	-2.044e-01	3.312e+01	-0.006	0.99	5
Survived	3.816e+00	1.865e+04	0.000	1.00	0
Sex	-5.213e+01	1.916e+04	-0.003	0.99	8
Age	-1.505e+02	2.650e+04	-0.006	0.99	5
Fare	-3.357e+02	6.954e+04	-0.005	0.99	6
Pclass_1	7.996e+01	1.700e+04	0.005	0.99	6
Pclass_2	-2.784e+01	1.949e+04	-0.001	0.99	9
Pclass_3	NA	NA	NA	N	A
Family_size	1.448e+02	6.960e+04	0.002	0.99	8
Title_1	4.369e+01	4.170e+04	0.001	0.99	9
Title_2	6.373e+01	1.499e+05	0.000	1.00	0
Title_3	-2.749e+01	6.160e+04	0.000	1.00	0
Title_4	NA	NA	NA	N	A
Emb_1	-3.231e+01	3.593e+05	0.000	1.00	0
Emb_2	-1.410e+01	3.587e+05	0.000	1.00	0
Emb_3	-8.577e+01	3.596e+05	0.000	1.00	0

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 1.5348e+01 on 791 degrees of freedom Residual deviance: 1.8442e-07 on 777 degrees of freedom AIC: 30

Number of Fisher Scoring iterations: 25



Estimated confidence interval of beta (X) coefficient under normal distribution

This tab evaluates the correlation coefficients between variables and represents them through a correlation map as shown, where each cell depicts acorrelation between many variables. The size and colorof the circles in each cell depict the degree of correlation, the larger the size and darker the color shade; the higher is the correlation.



PREDICTED PROB VS ACTUAL TAB

Logistic Regression is the statistical fitting of an s-curve logistic or logit function to a dataset in order to calculate the probability of the occurrence of a specific categorical event based on the values of a set of independent variables. Unlike linear regression which outputs continuous number values, logistic regression transforms its output using the logistic sigmoid function to return a probability value which can then be mapped to two or more discrete classes.

Theoretical probability uses math to predict the outcomes. Just divide the favorable outcomes by the possible outcomes. Experimental probability is based on observing a trial or experiment, counting the favorable outcomes, and dividing it by the total number of times the trial was performed.

Predicted Probability : Logistic regression analysis predicts the odds of an outcome of a categorical variable based on one or more predictor variables. A categorical variable is one that can take on a limited number of values, levels, or categories, such as "valid" or "invalid".

Actual Probability : In statistics, the actual value is the value that is obtained by observation or by measuring the available data. It is also called the observed value. The predicted value is the value of the variable predicted based on the regression analysis.



Use the left panel to impute or drop the missing values identified here

ACCURACY, ROC & AUC TAB



Accuracy : accuracy = correct_predictions / total_predictions. Accuracy is the proportion of correct predictions over total predictions. This is how we can find the accuracy with logistic regression: score = LogisticRegression.score(X_test, y_test). The most basic diagnostic of a logistic regression is predictive accuracy. To understand this we need to look at the prediction-accuracy table (also known as the classification table, hit-miss table, and confusion matrix).

Not only does binary logistic regression allow you to assess how well your set of variables predicts your categorical dependent variable and determine the "goodness-of-fit" of your model as does regular linear regression, but also it provides a summary of the accuracy of the classification of cases,

Confusion Matrix : A confusion matrix of binary classification is a two by two table formed by counting of the number of the four outcomes of a binary classifier. We usually denote them as TP, FP, TN, and FN instead of "the number of true positives", and so on. Predicted, Theoretical probability uses math to predict the outcomes. Just divide the favorable outcomes by the possible outcomes. It gives information about errors made by the classifier and the types of errors that are being made. It reflects how a classification model is disorganized and confused while making predictions.

AUC : The Area Under the ROC curve (AUC) is an aggregated metric That evaluates how well a logistic regression model classifies positive and Negative outcomes at all possible cutoffs. It can range from 0.5 to 1, and the larger it is the better.

ROC : ROC curves in logistic regression are used for determining the best cutoff value for predicting whether a new observation is a "failure" (0) or a "success" (1)The ROC curve shows the trade-off between sensitivity (or TPR) and specificity (1 - FPR). Classifiers that give curves closer to the top-left corner indicate a better performance. As a baseline, a random classifier is expected to give points lying along the diagonal (FPR = TPR).

