

Title

Methods Bootstrap and Cross-Validation, in the Construction of the Model of Vector Machine of Support for the Prediction of the Defaulting of the Banwest Bank's Clients

Participant

REYES RAMIREZ GIANCARLOS

Teacher

RICHARD FERNANDO FERNÁNDEZ VÁSQUEZ

Index

ABSTRACT	¡Error! Marcador no definido.
CHAPTER I	¡Error! Marcador no definido.
ANTECEDENTS	¡Error! Marcador no definido.
1.1 Researchings.....	¡Error! Marcador no definido.
CHAPTER II	¡Error! Marcador no definido.
PROBLEM STATEMENT	¡Error! Marcador no definido.
2.1 Description of the problem.....	¡Error! Marcador no definido.
2.2 Formulation of the problem	¡Error! Marcador no definido.
2.2.1 General Problem.....	¡Error! Marcador no definido.
2.2.2 Specific Problems.....	¡Error! Marcador no definido.
2.3 RESEARCHING AIMS	¡Error! Marcador no definido.
2.3.1 General Aim.....	¡Error! Marcador no definido.
2.3.2 Specific Aims	¡Error! Marcador no definido.
2.4 HYPOTHESIS STATEMENT	¡Error! Marcador no definido.
2.4.1 General Hypothesis	¡Error! Marcador no definido.
2.4.2 Specific Hypothesis.....	¡Error! Marcador no definido.
2.5 JUSTIFICATION	¡Error! Marcador no definido.
CHAPTER III	¡Error! Marcador no definido.
FRAMEWORK	¡Error! Marcador no definido.
3.1 SUPERVISED MODELS	¡Error! Marcador no definido.

3.2 Binary Logistic Regression	¡Error! Marcador no definido.
3.2.1. Method of Estimation	¡Error! Marcador no definido.
3.2.2. Significance Constrast	¡Error! Marcador no definido.
3.2.3. Contrast of kindness adjustment of Hosmer and Lemeshow	¡Error! Marcador no definido.
3.2.4. Clasification Table	¡Error! Marcador no definido.
3.4. DECISIÓN TREE	¡Error! Marcador no definido.
3.5. ARTIFICIALS NEURONAL NETWORKS (ANNS)	¡Error! Marcador no definido.
3.5.1 The ANNs of Radial Base:	¡Error! Marcador no definido.
3.5.2. Architecture	¡Error! Marcador no definido.
3.6. VECTOR MACHINE OF SUPPORT	¡Error! Marcador no definido.
3.7. TECHNIQUES OF RESAMPLING	¡Error! Marcador no definido.
3.7.1CROSS-VALIDATION.....	¡Error! Marcador no definido.
3.7.2 BOOTSTRAP	¡Error! Marcador no definido.
3.8 ROC CURVE/ GINI INDEX	¡Error! Marcador no definido.
3.9 BASIC TERMINOLOGY	¡Error! Marcador no definido.
CHAPTER IV.....	¡Error! Marcador no definido.
METHODOLOGY	¡Error! Marcador no definido.
4.1 POPULATION IN STUDY	¡Error! Marcador no definido.
4.2 SOURCES OF INFORMATION	¡Error! Marcador no definido.
4.3 DEFINITION OF VARIABLES	¡Error! Marcador no definido.
1. BASIC DATA	¡Error! Marcador no definido.
2. CREDIT DATA.....	¡Error! Marcador no definido.
4.4 DESIGNE OF SAMPLING AND PREPARATION DATA	¡Error! Marcador no definido.
4.5 STATISTICAL PROCESSING	¡Error! Marcador no definido.
GANTT CHART	¡Error! Marcador no definido.
COST BUDGET.....	¡Error! Marcador no definido.
RESULTS.....	¡Error! Marcador no definido.
BIBLIOGRAPHIC REFERENCES	¡Error! Marcador no definido.

Abstract

This research has as its main aim to determine the best method of resampling, in the construction model vector machine of support for the prediction of the defaulting of the Banwest Bank's clients, in order to analyze different scenarios that allow us to analyze

the efficiency of the Model of Vector Machine of support, through historical information of Banwest Bank corresponding to the period of July and September of 2016.

Given that it is important for financial institutions count with predictive models for defaults, we seek to raise an optional way to predict the already stated variable. Therefore, through this application, it will allow us to show that the Cross-Validation method can be more powerful than the Bootstrap method, and this is due to the fact that it has higher GINI and sensitivity indicators.

Key Words

Support Vector Machine, Cross Validation, Bootstrap, Defaulting

CONCLUSIONS

The best resampling method, for the construction of the vector machine of support, for the prediction of the defaulting of the Banwest Bank's clients turn out to be the Bootstrap's model, it shows us an indicator of sensibility of the 86.78%, for which reason is best in 8.9% in comparison to the Cross-Validation (77.80%), which means, through the Vector Machine of Support, with the Bootstrap's model we was achieved to classify the clients of the delinquent portfolio.

With the Boruta's algorithm, it was determined that the factors most influence in the defaulting of the delinquent portfolio's clients are the loan state, Category of the Superintendence of Bank and Insurances, Internal Segmentation, Total Debt, Type of Credit, Disbursement Amount, Interest Rate, Instalment Rates, Balance Amount, Days past due, Overdue Balance.

The scenario, which under the model of the vector machine of support, for the prediction of the defaulting of the Banwest Bank's clients, through the Bootstrap's model, showed better efficiency, considering 40 Bootstrap's samples, for which it was presented a Gini indicator of 77.78%. Whereas the capacity to predict the defaulting event was 86.78% with an estimation error of 891.123.

The scenario, which under the model of the vector machine of support, for the prediction of the defaulting of the Banwest Bank's clients, through the Cross-Validation method, presented a better efficiency, considering 10 partitions, for which it was presented a Gini indicator of 77.26%. Whereas the capacity to predict the defaulting event was 77.80% with an estimation error of 1547.45.

Title

Predicting of Success of the Telemarketing Campaign through Deposits to Term to the Clients of the Company "Bankuni", Through the Comparison of the Random Forest's Technique and the Naive Bayes' Technique

Participant

CHERO CAJUSOL CARLOS ENRIQUE

Teacher

AMELIDA PINEDO

Index

ABSTRACT (SPANISH VERSION) _____	¡Error! Marcador no definido.
ABSTRACT (ENGLISH VERSION) _____	¡Error! Marcador no definido.
Introduction _____	¡Error! Marcador no definido.
CHAPTER I _____	¡Error! Marcador no definido.
1. ANTECEDENTS _____	¡Error! Marcador no definido.
CHAPTER II _____	¡Error! Marcador no definido.
2. THE PROBLEM _____	¡Error! Marcador no definido.
2.1. SUBJECT _____	¡Error! Marcador no definido.
2.2. DESCRIPTION OF THE PROBLEM _____	¡Error! Marcador no definido.
2.3. FORMULATION OF THE PROBLEM _____	¡Error! Marcador no definido.
2.3.1. General Problem _____	¡Error! Marcador no definido.
2.3.2. Problemas Específicos: _____	¡Error! Marcador no definido.
2.4. AIMS _____	¡Error! Marcador no definido.
2.4.1. General Aim _____	¡Error! Marcador no definido.
2.4.2. Specific Aims _____	¡Error! Marcador no definido.
2.5. HYPOTHESIS _____	¡Error! Marcador no definido.
2.5.1. General Hypothesis _____	¡Error! Marcador no definido.
2.5.2. Specific Hypothesis _____	¡Error! Marcador no definido.
2.6. JUSTIFICATION _____	¡Error! Marcador no definido.
2.7. OPERACIONATION OF THE VARIABLES _____	¡Error! Marcador no definido.
2.8. MATRIX OF CONSISTENCE _____	¡Error! Marcador no definido.

2.9. DEFINITION OF THE VARIABLES _____	¡Error! Marcador no definido.
CHAPTER III _____	¡Error! Marcador no definido.
3. FRAME WORK _____	¡Error! Marcador no definido.
3.1. TECHNIQUES TO USE: _____	¡Error! Marcador no definido.
3.1.1. RANDOM FOREST _____	¡Error! Marcador no definido.
3.1.2. NAIVE BAYES _____	¡Error! Marcador no definido.
CHAPTER IV _____	¡Error! Marcador no definido.
4. METODOLOGY _____	¡Error! Marcador no definido.
4.1. STANDPOINT, TYPE and RESEARCH LEVEL _____	¡Error! Marcador no definido.
4.1.1. STANDPOINT _____	¡Error! Marcador no definido.
4.1.2. TYPE OF RESEARCH _____	¡Error! Marcador no definido.
4.1.3. RESEARCH LEVEL _____	¡Error! Marcador no definido.
4.1.4. DESIGN OF THE RESEARCH _____	¡Error! Marcador no definido.
4.2. SAMPLE DESIGNE _____	¡Error! Marcador no definido.
4.2.1. POPULATION AND SAMPLE _____	¡Error! Marcador no definido.
4.3. SOURCE OF INFORMATION _____	¡Error! Marcador no definido.
4.4. PLAN OF INFORMATION PROCESS _____	¡Error! Marcador no definido.
CHAPTER V _____	¡Error! Marcador no definido.
5. ANALYSIS AND INTERPRETATION OF OUTCOMES _____	¡Error! Marcador no definido.
5.1. DATA PREPARATION _____	¡Error! Marcador no definido.
5.2. RANDOM FOREST _____	¡Error! Marcador no definido.
5.3. NAIVE BAYES _____	¡Error! Marcador no definido.
5.4. MODELS EVALUATION _____	¡Error! Marcador no definido.
CHAPTER V _____	¡Error! Marcador no definido.
6. CONCLUSIONS AND RECOMENDATIONS _____	¡Error! Marcador no definido.
6.1. CONCLUSIONS _____	¡Error! Marcador no definido.
6.2. RECOMENDATIONS _____	¡Error! Marcador no definido.
BIBLIOGRAPHY _____	¡Error! Marcador no definido.
ATTACHED _____	¡Error! Marcador no definido.
ATTACHED A: COST AND BUDGET _____	¡Error! Marcador no definido.
ATTACHED B: GANTT GRAPHIC _____	¡Error! Marcador no definido.

Abstract

The BankUNI Company is a financial institution whose main activities are to provide credit services to customers, receiving savings deposits, etc. In a context of customer loyalty it is necessary to implement relational marketing tools that help improve the company's position in the financial sector in the city of Lima.

It is for this reason that this research has focused on an analysis of the success of the telemarketing campaign offering term deposits to its customers in order to increase customer loyalty.

Thus the study of the problem: What data mining technique, Naive Bayes or Random Forests has better predictive ability to determine the success of the telemarketing campaign term deposits?, This predictability will be measured by indicators such as Specificity, sensitivity and the ROC curve.

KEYWORDS: Random Forest, Naive Bayes, specificity, sensitivity, Gini Index.

Conclusions

- The Random Forest Model turned out to be appropriate that showed positive indicators to classify if a customer will access to the subscription of a term deposit after being managed through the telemarketing's campaign.
- The most significant variables to the Random Forest model was the duration of the last call in seconds and the contacts of the last month to have success the telemarketing's campaign.
- In the Naive Bayes's model, the most significant variables was the Numbers of days that has passed after the customer have been contacted by a previous company and the duration of the last call in seconds.

Title

Identification of the Main Factors than Influence in the Mortal Car Accidents In Metropolitan Lima, by te Outcomes of the Importance of Variables of the Models Of Data Mining, Random Forest, Boosting and Cart's Decisions Tree - National Census of Police Stations 2014

Participant

Tarazona Tocto Sherly Sindia

Teacher

RICHARD FERNANDO FERNÁNDEZ VÁSQUEZ

Index

ABSTRACT

INTRODUCTION

CHAPTER I.....	12
ANTECEDENTS.....	12
1.1 Researches.....	12
CHAPTER II.....	14
PROBLEM STATEMENT.....	15
2.1 Description of the problem.....	15
2.2 Formulation of the problem.....	16
2.3 Objetives of the research.....	16
2.4 Hypothesis of the research.....	16
2.5 Justification.....	18
2.6 Matrix of consistence.....	18
CHAPTER III.....	22
FRAME WORK.....	22

3.1 Previous Techniques.....	22
3.1.1 Chi square correlation.....	23
3.1.2 Test of homogeneity of variances.....	24
3.1.2.1 Bartlett’s Test.....	25
3.1.2.2 Levene’s Test.....	25
3.2 Techniques to use.....	27
3.2.1 Random forest’s Model.....	27
3.2.1.1 Error estimate with Random Forest.....	27
3.2.2 Boosting’s Model.....	28
3.2.2.1 AdaBoost.M1 Algorithm.....	29
3.2.3 Decisions Tree.....	29
3.2.3.1 Segments with tree models.....	29
3.2.3.2 Segmentation Algorithm.....	30
3.2.3.2.1 CHAID Algorithm.....	30
3.2.3.2.2 CART Algorithm.....	31
3.2.3.2.3 CHAID vs CART.....	33
3.2.4 Methods for dealing with unbalanced data sets	
3.2.4.1 Undersampling.....	34
3.2.4.2 Oversampling.....	34
3.2.4.3 Smote.....	34
3.2.5 Indicators used to the comparison of models.....	35
3.2.5.1 Sensibility.....	35
3.2.5.2 Specificity.....	35
3.2.5.3 Overall rating.....	35
3.2.5.4 Youden Index.....	35
3.2.5.5 Ratio of true-positives.....	35
3.2.5.6 Ratio of true-negatives.....	35
3.2.5.7 Euclidean Distance.....	35
3.2.5.8 ROC Curve.....	35
3.2.5.9 Gini Index.....	35
3.2.6 Adjustments by balancing the distribution of the dependent variable.....	36
3.2.7 K-Fold Cross-Validation.....	37
3.3 Basic Terminology.....	37

CHAPTER IV	38
METODOLOGY.....	38
4.1 Population in study.....	38
4.2 Sources in study.....	38
4.3 Type of research.....	38
4.4 Definition of the variables.....	38
4.4.1 Dependent Variable.....	39
4.4.2 Independents Variables.....	39
4.5 Samples Design and Preparation of the data.....	41
4.6 Statistical Procedure.....	42
CHAPTER IV	44
OUTCOMES	43
5.1 Analysis of variables.....	43
5.1.1 Univariate analysis of independent variables.....	43
5.2 Modelling.....	52
5.2.1 Boosting Model.....	53
5.2.2 Random Forest Model.....	56
5.2.3 Decisions Tree’s CART Model.....	58
5.3 Comparison of the indicators of Models Performance.....	61
CONCLUSIONS.....	64
RECOMENDATIONS.....	65
BIBLIOGRAPHIC REFERENCES.....	66

Abstract

Every year, more than 1.2 million people die in the world in traffic accidents according to the World Health Report 2015; Likewise, a report from the Andean Community of Nations (CAN) published in 2013 placed Peru as the country with the highest traffic accident rate per 100,000 vehicles, followed by Bolivia, Colombia and Ecuador.

In 2014, one of the last National Census of Commissaries was held, which contains a section on Traffic Accidents. As a result of the census it was possible

to know that 45% of the total traffic accidents in the country occur in Lima Metropolitan.

Addressing the problem of fatal traffic accidents is a complex study because there are many factors involved, but at the same time interesting. The present investigation is an application in the field of the study on road accidents. The study started from the problem What are the main factors that influence the fatal traffic accidents in Metropolitan Lima, through the importance of the variables of the Data Mining models: Random Forest, Boosting, and CART Decision Tree?

To answer the research problem we used the use of Data Mining Models such as Boosting, Random Forest and CART Decision Tree, whose results showed a table of importance of variables from which it was possible to identify that the type of vehicle used (Road or avenue), Invasion of the lane and disrespect to the sign of transit by the driver, are the main factors that influence in the fatal result of the accidents of Transit in Metropolitan Lima.

Keywords

Fatal Traffic Accident, Sensitivity, Gini, Cross Validation, Boosting, Random Forest, Decision Tree CART, Smooth

Conclusions

- ✓ The main factors that influence in the fatality of the car accidents in Metropolitan Lima, which they were obtained from the outcomes of data mining techniques: Random Forest, Boosting and Decisions Tree, were: That the car accident occurred in the type of highway or avenue, it was because of the contempt for the transit signal by the driver, the type of the vehicle was a rural van (combi), that it was caused for invade a track, that the type of vehicle was a mototaxi and/or three-wheeled motorcycle
- ✓ The values of the performance indicators of the models used are: Gini higher to 50% for all models and higher to 70% to the Boosting Model, Sensibility higher to 55% for all the models, and Specificity higher to 60% for all the models.

- ✓ The error rates of the build models, obtained through the cross-valid model k-fold, with $k=10$, are fewer to 9% for the Boosting Model, fewer to the 15% for the Random Forest model and fewer to 25% for the CART Decisions Tree Model.

- ✓ The variables (top 3) that listed in the importance table of variables of the Boosting Model are: Type of path of occurrence of the car accident (highway), type major vehicle involved (truck-rural bus) and disrespect to the transit signal by the driver; for the Random Forest Model are: Type of major vehicle involved (van - rural van), Type of transport (public) and the car accident by shove; finally to the Decisions Tree are: Type of route of accident occurrence (highway), type of minor vehicle involved (three-wheels motorcycle taxi) and for invasion in the opposite track.

Title

Comparison Between Neuronal Network Techniques and the Decision Tree for the Relapse Prediction of Young Offenders using Information of the Census in 2016

Participant

Adama Espirilla Jhonny Ivan

Teacher

RICHARD FERNANDO FERNÁNDEZ VÁSQUEZ

Index

ABSTRACT (SPANISH).....	¡Error! Marcador no definido.
ABSTRAC (ENGLISH).....	¡Error! Marcador no definido.
DEDICATORY	¡Error! Marcador no definido.
APPRECIATION	¡Error! Marcador no definido.
CHAPTER I	¡Error! Marcador no definido.
ANTECEDENTS.....	¡Error! Marcador no definido.
CHAPTER II	¡Error! Marcador no definido.
PROBLEM STATEMENT.....	¡Error! Marcador no definido.
2.1. Description of the problem	¡Error! Marcador no definido.
2.2. FORMULATION OF THE PROBLEM	¡Error! Marcador no definido.
2.2.1. General Problem.....	¡Error! Marcador no definido.
2.2.2. Specifics Problems	¡Error! Marcador no definido.
2.3. INVESTIGATION AIMS.....	¡Error! Marcador no definido.
2.3.1. General Aim	¡Error! Marcador no definido.
2.3.2. Specifics Aims	¡Error! Marcador no definido.
2.4. HYPOTHESIS STATEMENT.....	¡Error! Marcador no definido.
2.4.1. General Hypothesis	¡Error! Marcador no definido.
2.4.2. Specific Hypothesis:.....	¡Error! Marcador no definido.
2.5. JUSTIFICATION.....	¡Error! Marcador no definido.
2.6. CONSISTENCE MATRIX:	¡Error! Marcador no definido.
CHAPTER III	¡Error! Marcador no definido.

FRAMEWORK	¡Error! Marcador no definido.
3.1. GENERAL DEFINITIONS:.....	¡Error! Marcador no definido.
3.2. PRELIMINARY TESTS	¡Error! Marcador no definido.
3.3. DECISIONS TREE	¡Error! Marcador no definido.
3.4. NEURONAL NETWORKS:.....	¡Error! Marcador no definido.
3.5. SENSITIVITY AND ESPECIFICITY	¡Error! Marcador no definido.
3.6. THE ROC CURVE.....	¡Error! Marcador no definido.
3.7. MEAN SQUARE ERROR	¡Error! Marcador no definido.
3.8. OVER-SAMPLING	¡Error! Marcador no definido.
CHAPTER IV	¡Error! Marcador no definido.
METODOLOGY.....	¡Error! Marcador no definido.
4.1. Type of investigation:.....	¡Error! Marcador no definido.
4.2. Level of investigation:	¡Error! Marcador no definido.
4.3. Designe of the investigation:	¡Error! Marcador no definido.
4.4. Population in study:	¡Error! Marcador no definido.
4.5. Analisys unit	¡Error! Marcador no definido.
4.6. Information sources:	¡Error! Marcador no definido.
4.7. Sampling design and data preparation:.....	¡Error! Marcador no definido.
CHAPTER V	¡Error! Marcador no definido.
OUTCOMES	¡Error! Marcador no definido.
5.1. Interpretation of influential variables:.....	¡Error! Marcador no definido.
5.2. Data cleaning	¡Error! Marcador no definido.
5.3. Descriptive Analysis:.....	¡Error! Marcador no definido.
5.4. Over-sampling balance:.....	¡Error! Marcador no definido.
5.5. Modelling:.....	¡Error! Marcador no definido.
5.6. Cross-validation:.....	¡Error! Marcador no definido.
5.7. Comparison techniques.....	¡Error! Marcador no definido.
5.9. Profile of the young repeat offenders.....	¡Error! Marcador no definido.
CONCLUSIONS:.....	¡Error! Marcador no definido.
RECOMENDATIONS.....	¡Error! Marcador no definido.
BIBLIOGRAPHIC REFERENCES	¡Error! Marcador no definido.
ATTACHED	¡Error! Marcador no definido.

Abstract

Peru currently presents problems with juvenile delinquency, for this study is aimed at determining the factors that are part of the young offenders using decision tree techniques and neural networks, the theoretical information and the database received from the first census National to youth centers held on August 14 this year, which had 5 modules. The treatment of the variables was performed using data imputation methods, the sampling methodology was used for the treatment of asymmetry, and to verify that the parameters were not overestimated the cross-validation was performed, the results of the tree of the Which has been proved the best stability of the errors worked in the validation of the cross, in addition to its approximation of the indicators that had with the neural networks, finally the indicators for the techniques and the selection of the most adequate technique were determined It was the tree-making technique for its indicators.

Conclusions

- ✓ For the obtaining of the statistic technique was obtained that the most appropriate is the decision tree for his parameters of stability in the cross validation and his indicators.
- ✓ The indicators found in the neuronal network were:

INDICATORS	NEURONALNETWORK
sensitivity	66.2%
specificity	68.1%
accuracy	66.7%
VP	52.3%
VN	47.7%
GINI	38.2%

✓ The indicators found in the decision tree were:

INDICATORS	DECISION TREE
sensitivity	68.1%
specificity	68.0%
accuracy	67.2%
VP	53.1%
VN	46.9%
GINI	54.5%

Title

Effects of truncation and censoring of the dependent variable in 2-stage multilevel regression models

Participant

Huanca Huamaní Milagros Luz

Teacher

RICHARD FERNANDO FERNÁNDEZ VÁSQUEZ

Index

CHAPTER I.....	1
ANTECEDENTS.....	1
1.1. Investigations.....	1
1.2. General aspects.....	4
CHAPTER II.....	7
PROBLEM STATEMENT.....	7
2.1. Problem's description.....	7
2.2. Problem's formulation.....	8
2.2.1. General problem.....	8
2.2.2. Specific problems.....	8
2.3. Objectives.....	8
2.3.1. General objective.....	9
2.3.2. Specific objectives.....	9
2.4. Hypotesis.....	9
2.4.1. General hypotesis.....	9
2.4.2. Specific hypotesis.....	10
2.5. Justification.....	10
2.6. Consistence matrix.....	12
CHAPTER III.....	13

THEORETICAL FRAMEWORK.....	13
3.1. Using techniques.....	13
3.1.1. Model of multilevel regression.....	13
3.2.2. Truncated and censored data.....	30
3.3. Basic terminology.....	32
CHAPTER IV.....	34
METHODOLOGY.....	34
4.1. Population survey.....	34
4.2. Information sources.....	35
4.3. Definition of variables.....	35
4.4. Statistics procedures.....	36
CHAPTER V.....	37
RESULTS.....	37
5.1. Descriptive analysis of the variables.....	37
5.1.1. Descriptive analysis of the stage 1.....	37
5.1.2. Descriptive analysis of the stage 2.....	39
5.2. Analysis of the model indicators.....	43
5.2.1. Analysis indicators of the stage 1.....	43
5.2.2. Analysis indicators of the stage 2.....	46
CONCLUSIONS.....	51
RECOMMENDATIONS.....	53
BUDGET AND FINANCE.....	54
GANTT DIAGRAM.....	55
BIBLIOGRAPHIC REFERENCES.....	56
ANNEXES.....	57

Abstract

It has demonstrated that the information censored and truncated affects the estimation of the parameters by the MCO in the models of linear regression (Zuehlke & Kassekert, 2008). Moreover, we know that the multilevel models constitute the methodology of analysis more appropriate for treat information "Hierarchical" or "Nested" (Murillo, 2008). The multilevel models are considered extensions of the models of classics linear regression.

What was done in this investigation is study the marginal effect of the dependent variable censored and truncated in the models of multilevel regression of 2 periods in the next criterions: Akaike information criterion (AIC), Akaike information criterion corrected (AICC) and conditional Akaike information criterion (CAIC) in the models of multilevel regression of 2 periods. This under the diagram of 2 stages: in the first stage worked with a sample of

200 searches and different proportions of truncation and censoring as 10, 20 and 30%. Keeping in mind that is consider a high censoring since 20%. In the second stage was fixed a proportion of truncation and censoring of 10% for different sample sides as 50, 200 and 600 registers.

It was concluded that the fact of the dependent variable is censored and truncated, it has effects (in different percentage) about Akaike information criterion (AIC), Akaike information criterion corrected (AICC) and conditional Akaike information criterion (CAIC), these values tend to decrease up to 12% in the case of truncation and increase up to 190% in the case of censoring.

Conclusions

- ✓ The fact that the dependent variable is censored or truncated has effects (In different percentages) about Akaike information criterion (AIC), Akaike information criterion corrected (AICC) and conditional Akaike information criterion (CAIC), these values tend to decrease up to 12% in the case of truncation and increase up to 190% in the case of censoring.
- ✓ As much for the stage 1 where the sample was fixed with 200 registers and was made different proportions of truncation as 10, 20 and 30%, as for the stage 2 where was fixed the proportion of truncation in 10% for different samples sides

as 50, 200 and 600 registers , were obtained values of Akaike information criterion (AIC), Akaike information criterion corrected (AICC) and conditional Akaike information criterion

(CAIC) below the obtained values in the block control, where it doesn't exist any type of modification of information. The difference in percentage values become being up to 12% (When the sample is 50).

- ✓ The effect of the censoring of the dependent variable in the values of Akaike information criterion (AIC), Akaike information criterion corrected (AICC) and conditional Akaike information criterion (CAIC) as much for the stage 1 where the sample was fixed with 200 registers and was made in different proportions of censoring of 10, 20 and 30%, as the stage 2 where was fixed the proportions of truncation in 10% for different samples sides as 50, 200 and 600 registers, is too negative so these values of information criterions with respect to the values in the control block where wasn't made censoring have a percentage difference from up to 190% (For a sample of 600).

Title

Classification of scam with the model of Neuronal Networks using a response surface to select the independent variables that optimize the malign connections in simulated datums of the United States Armed Forces in 2016

Participant

Huamani Gonzales L. Dayana

Teacher

Richard Fernando Fernández Vásquez

Abstract

Based on the fundamentals and techniques of data mining, models can be designed and elaborated that allow finding clandestine behaviors that are easily detected at first sight. In particular the usefulness of data mining in this area lies in a series of techniques, algorithms and methods that imitate the human characteristic of learning to be able to extract new knowledge from experiences. These characteristics can be of vital importance to be applied in information security through the detection of intruders. This paper aims to show the contribution to the information security using the Response Surface as a technique that with the experience of the investigator can prevent the new modalities of information theft.

Keywords

Denial of services, ANN, intrusions, data mining, model, prediction.

Index

ABSTRACT	¡Error! Marcador no definido.
CONCLUSIONES.....	¡Error! Marcador no definido.
ABSTRAC.....	¡Error! Marcador no definido.

Conclusions

The comparison of the Neuronal Networks, using or not using the surface response previously it wasn't clear and conclusive; in determining that the model of Neuronal Networks without the use of the Surface response got a Gini=95% whereas that the model that used the protocol variables, flag, logged, count range as part to optimize the malign connection got a Gini=98%.

It was determined the model of second order is the most appropriate, since it was compared the R-adjust in both models y got a 62.03% of prediction of the variable malign connection, that is precisely the variable that we want to maximize.

It was obtained that the protocolo variables, flag, logged and count range are the variables that optimize the malign connections in the model of first order, subsequently with the ascend method more pronounced we make the calculus of the model of second order in which we got as significant variables a protocol, flag, logged, count range, protocol*flag, protocolo*loged, loged*rango_count y flag_loged