



**Statistical Model To Identify Factors Leading To
Autism In Children With Application On Pediatric
Hospital Ain Shams University**

A thesis submitted in partial fulfillment of the requirements for the
Master Degree in Applied Statistics

By

Kholoud Ahmed Maher Hamed Ahmed

Demonstrator - Statistics, Mathematics & Insurance Department
Faculty of Commerce – Ain Shams University

Under supervision of

Prof. Dr. Medhat Mohamed Ahmed Abdel Aal

Professor – Statistics, Mathematics & Insurance Department
Faculty of Commerce – Ain Shams University

Prof. Dr. Menan Abd El Maqsoud Rabie

Professor – Psychiatry Department
Faculty of Medicine – Ain Shams University

Dr. Essam Fawzy Aziz

Associate Professor – Statistics, Mathematics &
Insurance Department
Faculty of Commerce – Ain Shams University

2017



Approval sheet

Title of thesis: Statistical Model To Identify Factors Leading To Autism In Children With Application On Pediatric Hospital Ain Shams University

Academic Degree: M.Sc. in applied Statistics

Name of Student: Kholoud Ahmed Maher Hamed Ahmed

The thesis submitted in partial fulfillment of the requirements for the Master Degree in Applied Statistics has been approved by:

Examination Committee

- 1- Professor Dr. Ibrahim Hassan Ibrahim**
Professor – Statistics, Mathematics & Insurance Department
Faculty of Commerce – Helwan University
- 2- Professor Dr. Amr Ibrahim Abd El-Rahman El-Atraby**
Dean of Faculty of Commerce – Ain Shams University
Professor – Statistics, Mathematics & Insurance Department
Faculty of Commerce – Ain Shams University
- 3- Professor Dr. Medhat Mohamed Ahmed Abdel Aal**
Professor – Statistics, Mathematics & Insurance Department
Faculty of Commerce – Ain Shams University
- 4- Professor Dr. Menan Abd El Maqsood Rabie**
Professor – Psychiatry Department
Faculty of Medicine – Ain Shams University
- 5- Dr. Essam Fawzy Aziz**
Associate Professor – Statistics, Mathematics & Insurance Dept.
Faculty of Commerce – Ain Shams University

Date of thesis Defense 11 / 1 / 2017

Approval date / / 2017

To the Soul of my father.....

Kholoud Ahmed Maher

Acknowledgement

First and foremost, I would like to thank ALLAH for giving me the opportunity and strength to accomplish this research work.

I strongly acknowledge the contribution of **Prof. DR. Ibrahim Hassan Ibrahim** for his valuable time, patience, great support and professionalism which added a great value to the study.

I would like to express my sincere thanks and extreme gratitude to **Prof. DR. Amr Ibrahim Abdelrahman Elatraby** for all what I learnt from him throughout my academic years and I would like to thank him for agreeing to be a member of this committee and for his valuable time and advices which added a great value to the study.

No words can be sufficient to express my deep gratitude and appreciation to **Prof. DR. Medhat Mohamed Ahmed Abdel Aal** for granting me the privilege of working under his supervision and for his great support, valuable advices and continuous encouragement during all my academic years and throughout this work. Without his guidance and assistance I could never have been able to complete my thesis successfully.

I am so grateful and greatly thankful to my respected supervisor **Prof DR. Menan Abdel Maqsoud Rabie** for her generous support in the medical aspects of this work also for her great effort, careful supervision, continuous advice and guidance.

I sincerely would like to express my deepest gratitude and thanks to my respected supervisor **Prof. DR. Essam Fawzy Aziz** for

granting me the privilege of working under his supervision also for his careful supervision, continuous encouragement, and great support. His sincere effort and help will never be forgotten.

Many special thanks go to my colleagues for their help and goodwill.

Last but not least, I would like to express my warm gratitude to my mother and sister for their trust, kindness and unfailing support also I dedicate this thesis to the soul of my father may ALLAH grant him highest paradise (Ameen). Very special thanks to my beloved husband Ahmed and my children Mohammed and Omar who have always loved me, supported me and who have made untold number of sacrifices to help me accomplish this work.

Kholoud Ahmed Maher

Abstract

Kholoud Ahmed Maher Hamed Ahmed

**Statistical Model To Identify Factors Leading To Autism In
Children With Application On Pediatric Hospital Ain Shams
University**

Master Degree in Applied Statistics

Ain shams University, Faculty of Commerce

Statistics, Mathematics & Insurance Department

2017.

Autism spectrum disorders (ASDs) represent a group of neurodevelopmental disorders characterized by impairments in verbal and non-verbal communication, social withdrawal and stereotypical behaviors, which may or may not be associated with cognitive deficits, self-injurious behaviors and other neurological comorbidities. The increase of ASDs prevalence cannot be fully related to advances in diagnostics or sudden genetic shifts whereas there is a growing agreement among clinicians and scientists that ASDs resulted from an interaction between genetic and environmental factors.

One environmental factor that has received great attention is the body burden of heavy toxic metals as mercury, lead and aluminum. Heavy metals exposure is an increasing global problem and many previous studies demonstrated that heavy metals induce deleterious effects in humans.

This study drew a comparison between the results obtained on a given set of data gathered on a sample of Egyptian autistic children

against age and sex matched healthy controls using different statistical and data mining techniques being represented in parametric and nonparametric methods which are Logistic regression, Discriminant analysis, Classification and regression tree (CART), Artificial neural network and Random forests to determine the possible risk factors that may lead to autism and to reach the ideal model to be used as a tool to predict autism occurrence.

The variables studied in this thesis are House age, Age at conception, Consanguinity, Aluminum pans, Hair Lead level, Hair Mercury level, Dental amalgam, Passive smoking, Fish consumption and Gender. A split sample cross validation method is used to assess the validity of the classification models and it splits the sample into training and testing samples representing 75.7% and 24.3% of the total sample.

The traditional statistical techniques applied in this study represented in Forward Stepwise logistic regression and Stepwise discriminant analysis revealed that they have the same final model variables which are: House age, Age at conception, Aluminum pans, Consanguinity, and Hair lead level. Also, The data mining non parametric statistical techniques represented in Artificial neural network, Classification and regression tree (CART) and Random forests revealed better classification accuracy where the results sobtained according to the training data set showed that the artificial neural network has the best performance in the establishment of the prediction model and its classification accuracy is 82.08%, ***it is the recommended classifier for the diagnosis of autism.*** Moreover, the Classification and Regression trees (CART) as well as Random Forest models have the same classification accuracy which is equal to

78.30%. While the traditional models represented in Forward Stepwise logistic regression and Stepwise discriminant analysis models have accuracy rate of 77.40% and 76.40% respectively showing less classification accuracy than those of the nonparametric data mining techniques.

Furthermore, the study showed that the classification accuracy of the cross validated group cases (Testing data set) for Logistic regression and Discriminant analysis is 76.50% and 79.40% while for Artificial Neural Networks, Classification and Regression Trees (CART) and Random forest is 79.41%, 76.47%, 79.41% respectively.

Key words: Autism, House age, Age at conception, Consanguinity, Hair Lead level, Hair Mercury level, Dental amalgam, Logistic regression, Discriminant analysis, Artificial Neural Networks, Classification and Regression Trees (CART) and Random forest.

Table of Contents

List of Tables	IV
List of Figures	VI
List of Abbreviations	VIII
Chapter 1: Thesis Outline	
1.1 Introduction.....	1
1.2 Nature of the Problem.....	4
1.3 Importance of the Study.....	5
1.4 Objectives of the Study.....	6
1.5 Sources of Data.....	7
1.6 Variables of the Study.....	7
1.7 Hypothesis of the Study.....	8
1.8 Review of Literature.....	8
Chapter 2: Overview on Autism and Background on Statistical Techniques	
2.1 Overview on Autism	19
2.2 Logistic Regression	29
2.2.1 Odds.....	29
2.2.2 Odds Ratio.....	30
2.2.3 Assumptions of Logistic Regression.....	31
2.2.4 Overall Model Evaluation.....	31
2.2.4.1 Pearson Chi-Square Goodness of Fit.....	31
2.2.4.2 Homser-Lemshow Goodness of Fit.....	32
2.2.4.3 Deviance Goodness of Fit.....	33
2.2.5 Statistical Significance of Individual Regression Coefficients	
2.2.5.1 Likelihood Ratio Test.....	34
2.2.5.2 Wald Test.....	35

2.2.6 Predictive Accuracy and Discrimination of the Logistic Regression	36
2.2.6.2 Classification Table	36
2.2.6.2 Discrimination with Roc Curves.....	36
2.2.7 Validation of the Logistic Regression.....	37
2.3 Discriminant Analysis	
2.3.1 Stages of Conducting Discriminant Analysis.....	38
2.4 Artificial Neural Networks.....	44
2.4.1 Architecture of ANN.....	44
2.4.2 Types of Neural Networks Architecture.....	47
2.4.2.1 Single-Layer Feed Forward Networks.....	47
2.4.2.2 Multi-Layer Feed Forward Networks.....	47
2.4.2.3 Recurrent Networks.....	48
2.4.3. Neural Network Learning (Training) Algorithm	
2.4.3.1 Supervised Learning.....	50
2.4.3.2 Unsupervised Learning.....	50
2.4.3.3 Reinforcement Learning.....	50
2.4.4 The Back propagation.....	51
2.4.5 Cross validation.....	53
2.4.5.1 K- Fold Cross Validation.....	53
2.4.5.2 Leave-One-Out Cross Validation.....	53
2.4.5.3 Jackknife Resampling.....	54
2.4.5.4 Bootstrap Resampling.....	54
2.5 Classification and Regression Trees (CART).....	55
2.5.1 Steps of Conducting CART.....	56
2.6 Random Forests.....	60

Chapter 3: Application of Statistical Techniques

3.1 Variables of the Study.....	63
3.2 Logistic Regression.....	65
3.3 Discriminant Analysis.....	72
3.4 Artificial Neural Network.....	78
3.5 Classification and Regression Trees.....	85
3.6 Random Forests.....	94

Chapter 4: Results, Conclusions and Future Work.....103

References.....	110
------------------------	------------

List of Tables

Table (2.1) Activation Functions.....	46
Table (3.1) Case processing summary.....	65
Table (3.2) Variables in the Equation.....	67
Table (3.3) Model Summary.....	68
Table (3.4) Hosmer and Lemeshow Test.....	69
Table (3.5) Classification Table of Forward Stepwise Logistic Regression.....	70
Table (3.6) Area under Curve of Logistic Regression.....	71
Table (3.7) Variables entered the equation.....	73
Table (3.8) Variables entered / Removed.....	73
Table (3.9) Canonical Discriminant Function Coefficients.....	74
Table (3.10) Significance of the Discriminant Function.....	75
Table (3.11) Eigen Value and Canonical Correlation.....	75
Table (3.12) Classification Table of Stepwise Discriminant Analysis.....	76
Table (3.13) Area under Curve of Discriminant Analysis.....	77
Table (3.14) Summary of Active Networks.....	79
Table (3.15) Classification Summary Samples: Train.....	79
Table (3.16) Classification Summary Samples: Validation.....	80
Table (3.17) Sensitivity analysis Samples: Train, Validation.....	83
Table (3.18) Area under Curve of ANN Samples: Train, Validation.....	84
Table (3.19) Tree Structure.....	88
Table (3.20) Classification Matrix Analysis Sample.....	90
Table (3.21) Classification Matrix Test Sample.....	92
Table (3.22) Risk Estimates.....	95

Table (3.23) Classification matrix Training set sample.....	98
Table (3.24) Classification matrix Test set sample.....	100
Table (4.1) Comparison Table for the data analysis approaches.....	103.
Table (4.2) Predictor Importance of CART.....	105
Table (4.3) Predictor Importance of Random Forests.....	105

List of Figures

Figure (2.1) Logistic Regression Function.....	31
Figure (2.2) Univariate Representation of Discriminant Z score.....	39
Figure (2.3) Biological and Artificial Neuron Design.....	45
Figure (2.4) Feed Forward Network with a Single Layer of Neurons.....	47
Figure (2.5) Feed Forward Network with one hidden and one output Layer.....	48
Figure (2.6) Recurrent Neural Network.....	49
Figure (2.7) Back propagation Neural Network with one hidden layer.....	51
Figure (2.8) K- Fold Cross Validation.....	53
Figure (2.9) Leave One-Out Validation.....	54
Figure (3.1) ROC Curve of Logistic Regression.....	71
Figure (3.2) ROC Curve of Discriminant Analysis.....	77
Figure (3.3) Histogram of Model Accuracy for the Training Subset.....	81
Figure (3.4) Histogram of Model Accuracy for the Validation Subset.....	82
Figure (3.5) ROC Curve of Artificial Neural Network.....	84
Figure (3.6) Tree Diagram.....	87
Figure (3.7) Importance plot of predictor variables in CART Analysis.....	89
Figure (3.8) Bivariate Histogram of the classification Accuracy of Training Sample in CART analysis.....	91
Figure (3.9) Bivariate Histogram of the classification Accuracy of Test Sample in CART analysis.....	93

Figure (3.10) Summary of Random Forest.....96
Figure (3.11) Importance plot of Random Forest.....97
Figure (3.12) Bivariate histogram of the classification accuracy of
training sample in Random Forest analysis.....99
Figure (3.13) Bivariate histogram of the classification accuracy of
test sample in Random Forest analysis.....101