

Comparative Study of Neural Networks and Design of Experiments to the Classification of HIV status

WILBERT SIBANDA

Student Number: 21935009

- BSc (Life Sciences) University of Witwatersrand, Johannesburg
- BSc (Med) Hons Pharmacology, University of Cape Town
- MSc (Med) Pharmacy University of the Witwatersrand, Johannesburg

Thesis submitted for the degree of Doctor of Philosophy in Information Technology (IT) at the Vaal Triangle Campus of the North-West University

Promoter: **Prof. Philip Pretorius**
School of Information Technology
Vaal Triangle campus
North-West university
South Africa

Abstract

This research addresses the novel application of design of experiment, artificial neural networks and logistic regression to study the effect of demographic characteristics on the risk of acquiring HIV infection among the antenatal clinic attendees in South Africa.

The annual antenatal HIV survey is the only major national indicator for HIV prevalence in South Africa. This is a vital technique to understand the changes in the HIV epidemic over time. The annual antenatal clinic data contains the following demographic characteristics for each pregnant woman; age (herein called mother's age), partner's age (herein father's age), population group (race), level of education, gravidity (number of pregnancies), parity (number of children born), HIV and syphilis status.

This project applied a screening design of experiment technique to rank the effects of individual demographic characteristics on the risk of acquiring an HIV infection. There are various screening design techniques such as fractional or full factorial and Plackett-Burman designs. In this work, a two-level fractional factorial design was selected for the purposes of screening. In addition to screening designs, this project employed response surface methodologies (RSM) to estimate interaction and quadratic effects of demographic characteristics using a central composite face-centered and a Box-Behnken design.

Furthermore, this research presents the novel application of multi-layer perceptrons (MLP) neural networks to model the demographic characteristics of antenatal clinic attendees. A review report was produced to study the application of neural networks to modeling HIV/AIDS around the world. The latter report is important to enhance our understanding of the extent to which neural networks have been applied to study the HIV/AIDS pandemic.

Finally, a binary logistic regression technique was employed to benchmark the results obtained by the design of experiments and neural networks methodologies.

The two-level fractional factorial design demonstrated that HIV prevalence was highly sensitive to changes in the mother's age (15-55 years) and level of her education (Grades 0-13). The central composite face centered and Box-Behnken designs employed to study the individual and interaction effects of demographic characteristics on the spread of HIV in South Africa, demonstrated that HIV status of an antenatal clinic attendee was highly sensitive to changes in pregnant mother's age and her educational level. In addition, the interaction of the mother's age with other demographic characteristics was also found to be an important determinant of the risk of acquiring an HIV infection. Furthermore, the central

composite face centered and Box-Behnken designs illustrated that, individually the pregnant mother's parity and her partner's age had no marked effect on her HIV status. However, the pregnant woman's parity and her male partner's age did show marked effects on her HIV status in "two way interactions with other demographic characteristics".

The multilayer perceptron (MLP) sensitivity test also showed that the age of the pregnant woman had the greatest effect on the risk of acquiring an HIV infection, while her gravidity and syphilis status had the lowest effects. The outcome of the MLP modeling produced the same results obtained by the screening and response surface methodologies.

The binary logistic regression technique was compared with a Box-Behnken design to further elucidate the differential effects of demographic characteristics on the risk of acquiring HIV amongst pregnant women. The two methodologies indicated that the age of the pregnant woman and her level of education had the most profound effects on her risk of acquiring an HIV infection. To facilitate the comparison of the performance of the classifiers used in this study, a receiver operating characteristics (ROC) curve was applied. Theoretically, an ROC analysis provides tools to select optimal models and to discard suboptimal ones independent from the cost context or the classification distribution. SAS Enterprise Miner™ was employed to develop the required receiver-of-characteristics (ROC) curves.

To validate the results obtained by the above classification methodologies, a credit scoring add-on in SAS Enterprise Miner™ was used to build binary target scorecards comprised of HIV positive and negative datasets for probability determination. The process involved grouping variables using weights-of-evidence (WOE), prior to performing a logistic regression to produce predicted probabilities. The process of creating bins for the scorecard enables the study of the inherent relationship between demographic characteristics and an individual's HIV status. This technique increases the understanding of the risk ranking ability of the scorecard method, while offering an added advantage of being predictive.

Keywords:

Factorial, Central composite face-centered, Box-Behnken, multilayer perceptron, binary logistic regression, HIV, demographic characteristics.

Declaration

I, Wilbert Sibanda, hereby declare that the thesis entitled '**Comparative Study of the novel application of design of experiments (DOE) and Neural Networks (NN) to the classification of HIV status of antenatal clinic attendees in South Africa**' is my work. No plagiarism has taken place and due acknowledgements and references were given.

Wilbert Sibanda

15/11/2013

Acknowledgements

I would like to extend special thanks to the following people and/or institutions for their contributions towards this project.

- My beloved wife Cathrine (*Katie*) and my children, Lorraine (*Masi*), Janice (*Mute-Mute*), Njabulo (*Draad*) and Vuyo (*Pumpkin*).
- Prof. Philip Pretorius (my PhD supervisor)
- Medical Research Council of South Africa for the doctoral funding
- South African Centre for Epidemiological Modeling and Analysis (SACEMA), for the doctoral funding and constant workshops that sharpened my skills in epidemiological modeling.
- North-West university (Vaal Triangle campus) for graduate funding

Dedicated to the Memory of my Late Parents

Mr. and Mrs. Ores Leonard Sibanda

Table of Contents

| | |
|---|------------|
| Abstract | ii |
| Acknowledgements | vi |
| Chapter 1: Introduction | |
| 1.1. Background..... | 1 |
| 1.2. Problem statement..... | 5 |
| 1.3. Aims and objectives..... | 6 |
| Chapter 2: Research Methodology | |
| 2.1. Data Exploration..... | 11 |
| 2.2. Design of experiments..... | 13 |
| 2.3. Artificial neural networks..... | 15 |
| 2.4. Logistic regression..... | 23 |
| 2.5. Comparison of models using ROC curves..... | 28 |
| Chapter 3: Study Plan | |
| 3.1. Introduction..... | 37 |
| 3.2. Research Outline | 37 |
| 3.3. Software tools..... | 38 |
| 3.4. Methods..... | 38 |
| Chapter 4: Results | |
| 4.1. Data Exploration..... | 65 |
| 4.2. Screening design..... | 71 |
| 4.3. Response surface methodology..... | 76 |
| 4.4. Comparison of two response surface methodologies (RSM).. | 79 |
| 4.5. Comparison of RSM and binary logistic regression..... | 89 |
| 4.6. Application of multilayer perceptron (MLP)..... | 101 |
| 4.7. Using ROC curves to compare models..... | 105 |
| 4.8. Data sources for the models..... | 120 |

| | |
|---|------------|
| Chapter 5: Conclusion | |
| 5.1. Screening design..... | 126 |
| 5.2. Response surface methodology..... | 126 |
| 5.3. Comparison of two response surface methodologies (RSM)..... | 128 |
| 5.4. Comparison of RSM and binary logistic regression..... | 130 |
| 5.5. Application of multilayer perceptron (MLP)..... | 130 |
| 5.6. Model comparison using ROC curves..... | 131 |
| 5.7. Building scorecard using weights-of-evidence (woe)..... | 133 |
| Chapter 6: Implications of Research Findings..... | 134 |
| 6.1. What do models mean?..... | 134 |
| Bibliography..... | 136 |
| Annexure A: Screening design publication..... | 140 |
| Annexure B: Response surface methodology (RSM) publication..... | 146 |
| Annexure C: Comparison of two RSM methodologies publication..... | 169 |
| Annexure D: Comparison of RSM and BLR publication..... | 181 |
| Annexure E: MLP for epidemiological modeling publication..... | 206 |
| Annexure F: A review of neural networks..... | 212 |
| Annexure G: Development and validation of an HIV risk scorecard model..... | 220 |
| Annexure H: ROC curves to compare models..... | 227 |

List of Tables

| | |
|--|----|
| Table 2.1: Significance of moments in properties of distributions..... | 14 |
| Table 2.2: Determining the number of hidden layers..... | 25 |
| Table 2.3: Global classification..... | 40 |
| Table 2.4: Example of Scorecard..... | 46 |
| Table 3.1: The Fractional Factorial design matrix..... | 53 |
| Table 3.2: Factor Levels..... | 54 |
| Table 3.3: The central composite face-centered matrix..... | 56 |
| Table 3.4: Degrees of freedom of different errors..... | 57 |
| Table 3.5: Factor Levels..... | 59 |
| Table 3.6: Central composite face-centered design..... | 61 |
| Table 3.7: Box-Behnken..... | 62 |
| Table 3.8: Degrees of freedom of different errors..... | 63 |
| Table 3.9: Factor levels..... | 66 |
| Table 3.10: The Box-Behnken matrix..... | 68 |
| Table 3.11: Degrees of freedom for Box-Behnken design matrix..... | 69 |
| Table 3.12: Factor Levels..... | 71 |
| Table 3.13: Specifications of variables for the MLP technique..... | 74 |
| Table 3.14: Data Tagging..... | 75 |
| Table 4.1: Descriptive statistics of the 2007 antenatal HIV seroprevalence data..... | 81 |
| Table 4.2: Basic Statistics of the demographic variables..... | 81 |

| | |
|---|-----|
| Table 4.3: Basic Statistics of the demographic variables..... | 83 |
| Table 4.4: Basic Statistics of the demographic variables..... | 84 |
| Table 4.5: Basic Statistics of the demographic variables..... | 85 |
| Table 4.6: Basic Statistics of the demographic variables..... | 86 |
| Table 4.7: Basic Statistics of the demographic variables..... | 87 |
| Table 4.8: Pearson’s correlation coefficients of the demographic variables..... | 88 |
| Table 4.9: Predictive model generated by the screening design..... | 89 |
| Table 4.10: Constrained optimization results..... | 94 |
| Table 4.11: Model summary statistics: Small composite Hartley method..... | 94 |
| Table 4.12: Fit statistics for the CCF..... | 95 |
| Table 4.13: Fit statistics for the orthogonal central composite face-centered design..... | 97 |
| Table 4.14: Predictive models..... | 100 |
| Table 4.15: ANOVA Results..... | 101 |
| Table 4.16: Final equations of the central composite and Box-Benhken designs..... | 106 |
| Table 4.17: Sequential model sum of squares for the Box-Behnken design..... | 110 |
| Table 4.18: Model summary statistics for the Box-Behnken design..... | 110 |
| Table 4.19: Pearson’s Chi-Square test..... | 111 |
| Table 4.20: Deviance values..... | 111 |
| Table 4.21: AKAIKE Information Criterion (AIC)..... | 111 |
| Table 4.22: Schwarz criterion (SC)..... | 111 |
| Table 4.23: $-2\log L$ | 112 |
| Table 4.24: ANOVA Results for the Box-Behnken design..... | 112 |
| Table 4.25: Likelihood Ratio (LR), Wald and Score Tests..... | 113 |
| Table 4.26: Final Equation from Box-Behnken design..... | 118 |

| | |
|--|------------|
| Table 4.27: Maximum likelihood estimates..... | 121 |
| Table 4.28: Cross-validation results..... | 125 |
| Table 4.29: Classification..... | 129 |
| Table 4.30: Interpretation of Information values (IV)..... | 135 |
| Table 4.31: Binning of the variables..... | 136 |
| Table 4.32: Regression coefficients from the scorecard node..... | 137 |
| Table 4.33: Selected variables from the final scorecard..... | 137 |
| Table 4.34: Confusion matrix at 25% Threshold..... | 138 |

List of Figures

| | |
|--|----|
| Fig.2.1: A taxonomy of neural network architectures..... | 16 |
| Fig. 2.2: A schematic representation of an MLP with three layers..... | 17 |
| Fig. 2.3: Logistic transfer function..... | 18 |
| Fig. 2.4: A three-dimensional error plot..... | 19 |
| Fig. 2.5: The back-propagation algorithm..... | 20 |
| Fig. 2.6: The Sigmoid function..... | 22 |
| Fig. 2.7: The hyperbolic tangent activation function..... | 22 |
| Fig.2.8: The linear activation function..... | 23 |
| Fig. 2.9: The logit and probit transformations..... | 25 |
| Fig. 2.10: Confusion matrix..... | 29 |
| Figure 2.11: Confusion matrices | 29 |
| Fig. 2.12: ROC curves..... | 30 |
| Fig. 2.13: K-S test..... | 31 |
| Fig. 2.14: Theta (θ) = Area under the Curve..... | 32 |
| Fig. 3.1: Research Study Plan..... | 37 |
| Fig. 3.2: Demographic characteristics studied by the screening design..... | 40 |
| Fig. 3.3: Demographic characteristics studied by the Central Composite Face design... | 43 |
| Fig. 3.4: Standard errors..... | 46 |
| Fig. 3.5: Fraction of design space (FDS) plot of the standard error the design space.... | 47 |
| Fig. 3.6: Demographic characteristics studied by the CCF and BBD designs..... | 49 |
| Fig. 3.7: Standard error plot of the CCF and BBD designs respectively..... | 52 |
| Fig. 3.8: FDS plots of standard errors of the CCF and BBD designs..... | 53 |

| | |
|---|----|
| Fig. 3.9: VDG graphs of CCF and BBG designs..... | 54 |
| Fig. 3.10: Demographic characteristics studied by BBD and binary logistic regression..... | 56 |
| Fig. 3.11: 3D Plot of standard errors of the Box-Behnken design..... | 58 |
| Fig. 3.12: FDS plot of the standard error over the BBD design space..... | 59 |
| Fig. 3.13: Demographic characteristics studied by MLP..... | 62 |
| Fig. 4.1: HIV frequency..... | 65 |
| Fig. 4.2: Syphilis frequency..... | 66 |
| Fig. 4.3: Frequency of HIV infection by syphilis state..... | 66 |
| Fig. 4.4: Frequency of HIV infection by pregnant woman's age..... | 68 |
| Fig. 4.5: Frequency of HIV infection by gravidity..... | 68 |
| Fig. 4.6: Frequency of HIV infection by parity..... | 69 |
| Fig. 4.7: Frequency of HIV infection by education..... | 69 |
| Fig. 4.8: Frequency of HIV infection by male partner's age | 70 |
| Fig. 4.9: Fractional factorial coefficient | 71 |
| Fig. 4.10: The Lenth plot of the effect of demographic characteristics on HIV risk..... | 72 |
| Fig. 4.11: Normal plot..... | 73 |
| Fig. 4.12: Normal probability plot of errors..... | 74 |
| Fig. 4.13: Plot of residuals against predicted values..... | 75 |
| Fig. 4.14: Plot of residuals against experimental cases..... | 75 |
| Fig. 4.15: Central composite face-centered coefficient plot..... | 77 |
| Fig. 4.16: Surface plot of father's age and mother's age on HIV..... | 78 |
| Fig. 4.17: Normal plot of residuals..... | 78 |

| | |
|--|------------|
| Fig. 4.18: Normal plots of residuals for the CCF and Box-Behnken design..... | 81 |
| Fig. 4.19: Plot of residuals vs fitted response for the CCF and BBD..... | 82 |
| Fig. 4.20: Plots of residuals vs observation order for CCF and BBD..... | 83 |
| Fig. 4.21: Plot of leverage of points for the CCF and BBD..... | 84 |
| Fig. 4.22: Coefficient Plot of demographic characteristics..... | 85 |
| Fig. 4.23: Main effects plot..... | 86 |
| Fig. 4.24: Interactions Plot..... | 87 |
| Fig. 4.25: 3D Response surface plots of CCF and BBD designs..... | 88 |
| Fig. 4.26: Normal plot of residuals for the BBD..... | 93 |
| Fig. 4.27: Plot of residuals vs fitted response for the BBD..... | 94 |
| Fig. 4.28: Plot of residuals vs observation order for the BBD..... | 94 |
| Fig. 4.29: Deviance residuals from the logistic regression..... | 95 |
| Fig. 4.30: Pearson residuals from the logistic regression..... | 95 |
| Fig. 4.31: Plot of Leverage of points for the BBD..... | 96 |
| Fig. 4.32: Plot of leverage points for the binary logistic regression model..... | 96 |
| Fig. 4.33: Main Effects Plot..... | 97 |
| Fig. 4.34: Interactions plot..... | 97 |
| Fig. 4.35: Coefficient plot of main and interaction effects..... | 99 |
| Fig.4.36: Coefficient plot of the main effects of logistic regression..... | 100 |

| | |
|--|-----|
| Fig. 4.37: 3D Response surface plot of the BBD design..... | 100 |
| Fig. 4.38: Mean performance as a function of the hidden unit..... | 101 |
| Fig. 4.39: MSE as a function of the training iteration number..... | 102 |
| Fig. 4.40: Classification performance vs iteration number..... | 102 |
| Fig. 4.41: Sensitivity test results..... | 104 |
| Fig. 4.42: Comparison of modeling techniques using SAS Enterprise Miner..... | 105 |
| Fig. 4.43: Response threshold chart for non-coded data..... | 107 |
| Fig. 4.44: Diagnostic Charts at 25% threshold level..... | 108 |
| Fig. 4.45: Plot of correctly classified individuals across different threshold levels..... | 109 |
| Fig. 4.46: Threshold based accuracy plot..... | 110 |
| Fig. 4.47: Cumulative lift charts..... | 110 |
| Fig. 4.48: Non-Cumulative lift charts..... | 111 |
| Fig. 4.49: ROC curves..... | 112 |
| Fig. 4.50: Cumulative percentage response chart..... | 112 |
| Fig. 4.51: The noncumulative captured HIV positive..... | 117 |
| Fig. 4.52: ROC chart for the scorecard..... | 118 |
| Fig. 4.53: Empirical odds plot..... | 119 |
| Fig. 4.54: Kolmogorov-Smirnov plot..... | 119 |
| Fig. 4.55: Data source for fractional factorial design..... | 120 |
| Fig. 4.56: Data source for CCF design..... | 121 |
| Fig. 4.57: Data sources for CCF and BBD designs..... | 122 |
| Fig. 4.58: Data sources for BBD and logistic regression models..... | 123 |
| Fig. 4.59: Data sources for the multilayer perceptron design..... | 124 |
| Fig. 4.60: Data sources for the HIV risk scorecard | 125 |

| | |
|---|------------|
| Fig. 5.1: Lenth's plot..... | 126 |
| Fig. 5.2: CCF coefficient plot..... | 127 |
| Fig. 5.3: Perturbation plot..... | 128 |
| Fig. 5.4: Predicted versus observed values..... | 128 |
| Fig. 5.5: Coefficient plot of CCF and BBD designs..... | 129 |
| Fig. 5.6: Actual vs predicted HIV risk for BBD design..... | 129 |
| Fig. 5.7: Actual vs predicted HIV risk for CCF design..... | 130 |
| Fig. 5.8: Cumulative percentage response chart..... | 131 |
| Fig. 5.9: Noncumulative percentage response plot..... | 131 |
| Fig. 5.10: Receiver Operating Characteristic curve (ROC)..... | 132 |
| Fig. 5.11: Response threshold (FullFactorial, DOE, Tree, NN and Reg)..... | 132 |
| Fig. 5.12: Schematic representation of a blackbox..... | 133 |

List of Scientific Publications in Peer-reviewed journals

Response surface modeling and optimization to elucidate the differential effects of demographic characteristics on HIV prevalence in South Africa.

Journal: 2012 Computers and Industrial Engineering 42
Publisher: CIE & SAIIE
Date of Publication: 2012

Artificial neural networks- A review of applications of neural networks in the modeling of HIV epidemic.

Journal: International Journal of Computer Applications (0975-8887)
Date of Publication: 2012.

Novel application of multi-layer perceptrons (MLP) neural networks to model HIV.

Journal: International Journal of Computer Applications (0975-8887)
Date of Publication: 2011

Application of two-level fractional factorial design to determine and optimize the effect of demographic characteristics on HIV prevalence.

Journal: International Journal of Computer Applications (0975-8887)
Date of Publication: 2011

Comparative Study of the Application of Box-Behnken Designs (BBD) and Binary Logistic Regression (BLR) to study the effect of demographic characteristics on HIV risk in South Africa.

Journal: International Journal of Applied Medical Sciences

Application of Central Composite Face-Centered (CCF) and Box-Behnken Designs (BBD) to study the effect of demographic characteristics on HIV risk in South Africa.

Journal: Network Modeling and Analysis in Health Informatics and Bioinformatics

Application of ROC curves to compare neural networks, logistic regression and decision trees for modelling the causal relationship between demographic characteristics and the risk of acquiring infection using antenatal seroprevalence data.

Abstract accepted: 42nd Annual Conference of the Operations Research Society of South Africa

Development and validation of an HIV risk scorecard model.

Submitted: International Symposium on Network Enabled Health Informatics, Biomedicine and Bioinformatics (HI-BI-BI 2013), Niagara Falls, Canada