

DESIGN AND DEVELOPMENT OF A DATA-MINING MODEL FOR CENSUS ANALYSIS WITH GEOGRAPHICAL INFORMATION SYSTEM (GIS)

Ogochukwu Clementina Okeke^{1*}, Esther Nneka Okonkwo²

^{1*} Department of Computer Science, Chukwuemeka Odumegwu Ojukwu University, Uli, Nigeria

² Department of Computer Science, Chukwuemeka Odumegwu Ojukwu University, Uli, Nigeria.

Abstract: A new model for National Population Commission and Survey in Nigeria was developed to obtain geo-spatial distribution of predicted demographic information. Model used in census analysis gives prediction without searching for existing patterns or trends and without geo-spatial distribution of predicted demographic information. The system is a desktop application based on states and behaviours of objects. Object-Oriented Analysis and Design Methodology was adopted in this study and Microsoft Visual Basic.NET was used to implement the system. This was supported with SQLite database, a free distributed database engine. The new model developed automated the process of searching for patterns in the census data. The system is a multi-line user input interface which allows users to input specific data input to make predictions. The output of this paper is shown on the Nigerian map comprising of thirty six states, table data and bar chart simultaneously.

Keywords: Data-mining, Geographical Information System (GIS), Nigeria, Population.

1.0 Introduction

Nigeria has had a long history of census taking. The first census was in 1866. The censuses of 1866, 1871, 1896, were restricted to Lagos Island and parts of Lagos Mainland. The censuses of 1901, 1911 and 1921 covered, in addition to Lagos a few more urban towns in the colony. Most of these censuses were actually population estimates. Although the census of 1952/53 was elaborate in organization, its non-simultaneity, which has implications for possible double counting is considered its weak points. The first post-independence census in Nigeria was carried out in 1962. This was out rightly cancelled and another conducted in 1963. The 1963 census results became the official figure and were used until the 1991 census was accepted. There was a census in 1973 but the results were declared unacceptable on account of massive inaccuracies. The 1991 census broke the myth of failed censuses in Nigeria. In 2006, Nigeria had its last census [1]. It provides the total characteristics in every town, village or /and locality. This information provides data for planning programmes in education. Information on housing is dispensable to planners and policy makers in evaluating housing conditions, estimating housing needs and formulating housing policies [2]. Most countries of the world conduct population and economic censuses at regular intervals. Population census information is of great value in planning public services for governments at all levels, such as cities, counties, provinces, and states. Both population and economic census data are also used by private companies or community organizations for various purposes, such as marketing studies, situating new factories or shopping malls, developing social service programs. Therefore, the application of data mining techniques to census data has great potential both in underpinning good public policy and in supporting business developments. However, mining census data is not straight forward and requires challenging methodological research [3]. Census analyses in Nigeria do not search for existing patterns or trends during prediction to bring out hidden and important attributes in census data. A lot of data are not discovered during analysis due to the model in existence. National Population Commission has conducted out a lot of predictions in past, having much data that requires mining for further analysis. In this paper a model that predicts and gives predicted information on the map of Nigeria using dot density was developed. The new model developed searches for existing patterns using selected number of data input to carry forecast. The data-mining algorithm is the mechanism that creates mining models. To create a model, an algorithm first analyzes a set of data, looking for specific patterns and trends [4]. The algorithm then uses the result of this analysis to define the parameters of the mining model to give geo-spatial distribution. The data-mining model this algorithm creates takes various forms including; Marital status, gender, unemployment, employment, number of registered births, number of registered deaths and total population in each States and give geo-spatial distribution of predicted demographic information. The data-mining extract these attributes out from pool of census database. Data-mining techniques are the result of a long process of research and product development [5]. This evolution began when business data was first stored on computers, continued with improvements in data access, and more recently, generated technologies that allow users to navigate through their data in real time. Data-mining takes evolutionary process beyond retrospective data access and navigation to prospective and proactive information delivery. Data-mining is ready for application in the business community because it is supported by three technologies that are now sufficiently mature, which are massive data collection, powerful multiprocessors computers and data-mining algorithms [5]. The concept of data-mining is becoming increasingly popular as a business information management tool where it is expected to reveal knowledge structures that can guide decisions in conditions of limited certainty. Recently, there has been increased interest in developing new analytic techniques specifically designed to address the issues relevant to business data-mining (e.g Classification Trees), but data-mining is still based on the conceptual principles of statistics including the traditional exploratory data analysis and modeling and it shares with them both some components of its general approaches and specific techniques.

Data mining is the computational process of discovering patterns in large data sets involving methods at the intersection of artificial intelligence, machine learning, statistics, and database systems [6]. The overall goal of the data mining process is to extract information from a data set and transform it into an understandable structure for further use. Data mining has two tasks as follows: Prediction method: Use some variables to predict unknown or future values of other variables [7]. Description method: Find human-interpretable patterns that describe the data

[7]. The first and simplest analytical step in data mining is to describe the data, summarize its statistical attributes (such as means and standard deviations), visually review it using charts and graphs, and look for potentially meaningful links among variables (such as values that often occur together). A good model should never be confused with reality but it can be a useful guide to understanding your business [8]. Data mining is increasingly popular because of the substantial contribution it can make. It can be used to control costs as well as contribute to revenue increases. Data-mining process consists of three stages: the initial exploration, model building or pattern identification with validation and verification and deployment (i.e. the application of the model to new data in order to generate predictions) [9]. A geographic information system or geographical information system (GIS) is a system designed to capture, store, manipulate, analyze, manage, and present all types of spatial or geographical data [10]. A GIS is an organized collection of computer hardware, software, geographic data, and personnel designed to efficiently capture, store, update, manipulate, analyze, and display all forms of geographically referenced information [11].

1.1 Statement of the Problem

The problems identified with census analyses in Nigeria are;

- a. Census map is visual representation of States, local governments and towns and there is no prediction of demographic information on it.
- b. Census demographic analysis is calculated without mining any pattern or trend of existing demographic information.
- c. Census figures are not stored in a data-mining warehouse; population distribution decisions, finding hidden patterns and relationship become a problem.
- d. Census maps do not contain any chart that shows how predictions are narrowed down to a particular State.

1.2 Objectives of the Paper

- a. To develop a model that accepts given probabilistic demographic information outputs a better prediction of future values.
- b. To develop a population density map of Nigeria using dot density.
- c. To develop a data-mining warehouse and obtain geo-spatial distribution attributes of a population.
- d. To make population distribution decisions, find hidden patterns and relationships from census database.
- e. To develop a model that makes use of chart to illustrate population predictions in each State in Nigeria.

1.3 Importance of the Paper

With the predictive capability of data-mining and geographical information system that gives geo-spatial distribution of population, Nigerian government will be more proactive in formulating economic and other developmental plans. The work will provide support for tactical and strategic business decisions. Geographical information system has the ability to separate information in layers, and then combine it with other layers of information which is used in research and for decision making tools.

1.4 Related Literatures

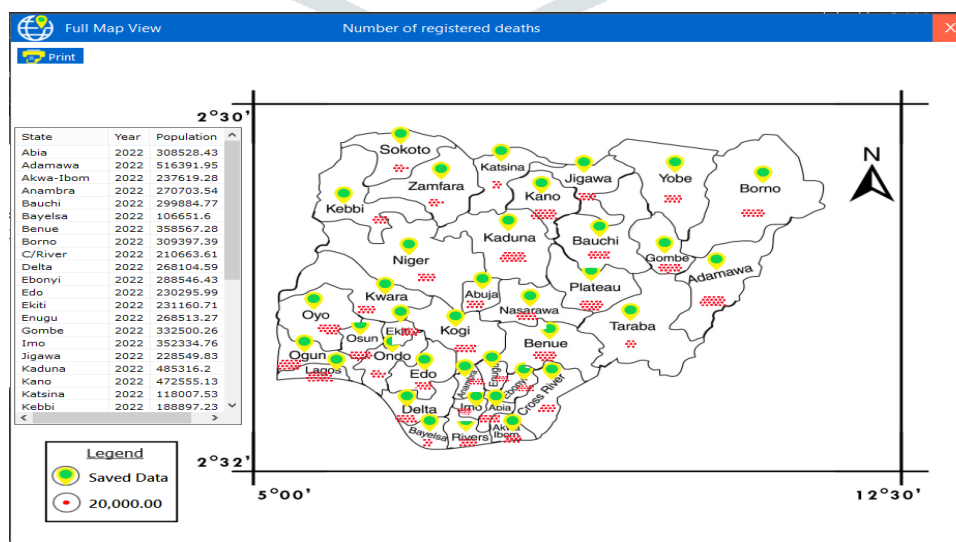
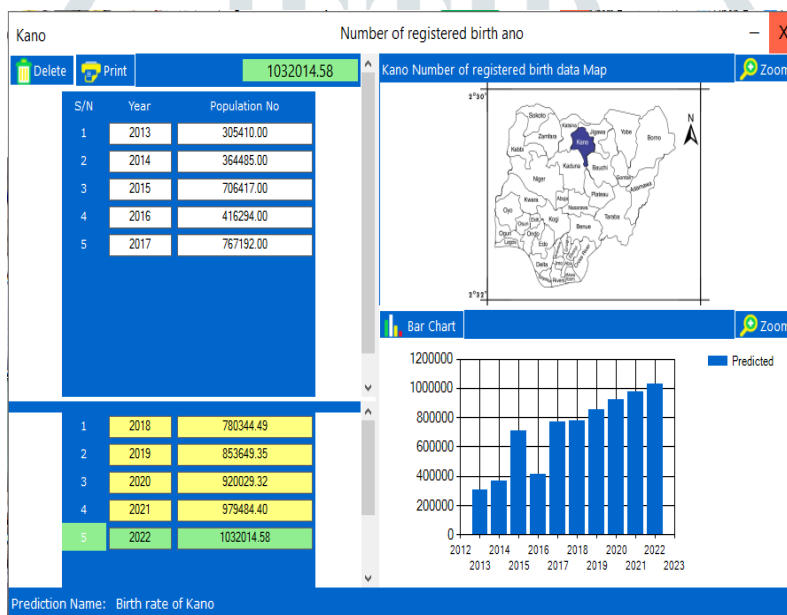
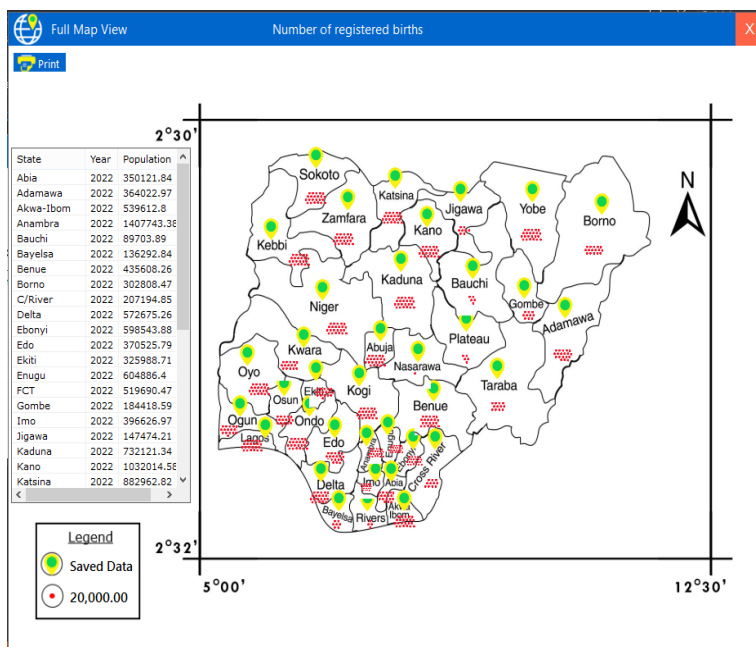
Ramezankhani et al. [12] Identified a low risk population for type 2 diabetes, Tehran lipid and glucose study using decision tree algorithm. The output was to find out number of people with Type 2 diabetes, Tehran and glucose level in future. Manoj and Madhu [13] applied a time series ARIMA model for predicting sugarcane production in India; their output was to forecast sugar cane production in future years. Md and Syed [14] used cohort-component method in forecasting future levels of fertility, mortality, sex composition, migration and other parameters and their output was to forecast factors that cause population change. David, Alan and Bob [15] applied Hamilton –Perry Method to forecast population using data of two recent censuses and the result showed how population increased in years to come. Kayode and Jimoh [16] used regression analysis with time series data to periodically forecast stock market prices stock market prices and have been proven complement to other numeric forecasting method and the result was to predict stock market prices in future and financial institution. Friedrich and Weixin [17] predicted literacy rate based on age disaggregated literacy data and demographic data collected from UN population using Global Age Specific Literacy Projection (GALP). Saigal and Mehrotra [18] predicted total birth using autoregressive integrated moving average (ARIMA) and output was predict number of births over given years. Rajasekhar and Karth [19] proposed hybrid support vector machine technique for weather prediction using Artificial Neural Network (ANN) and result yielded Euclidean distance on monthly mean of each year average temperature.

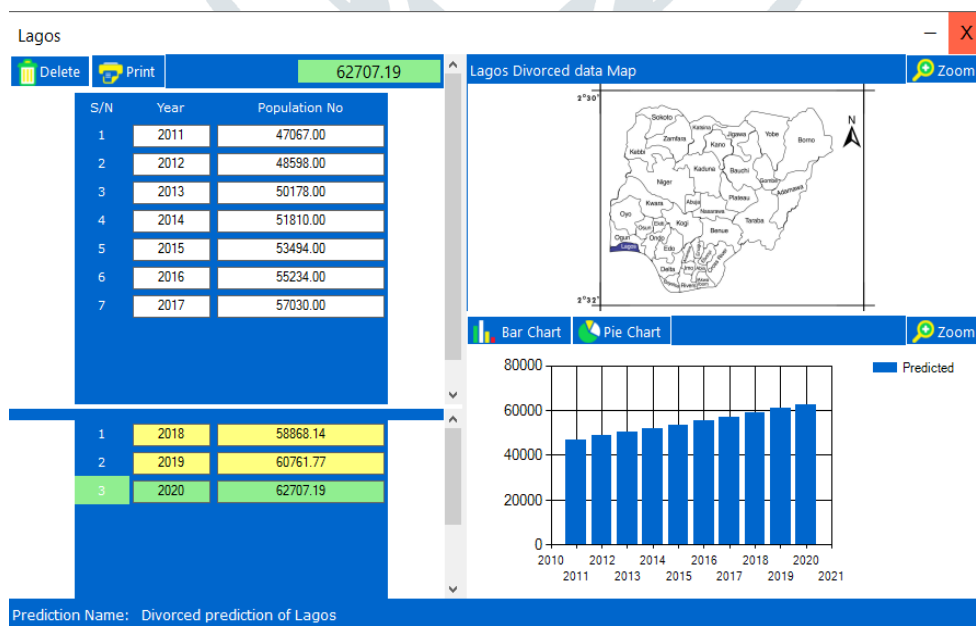
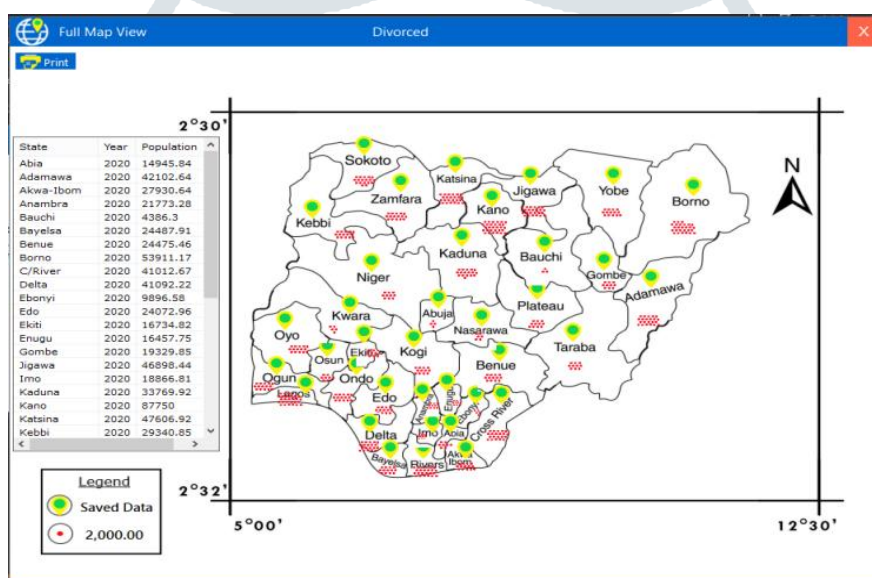
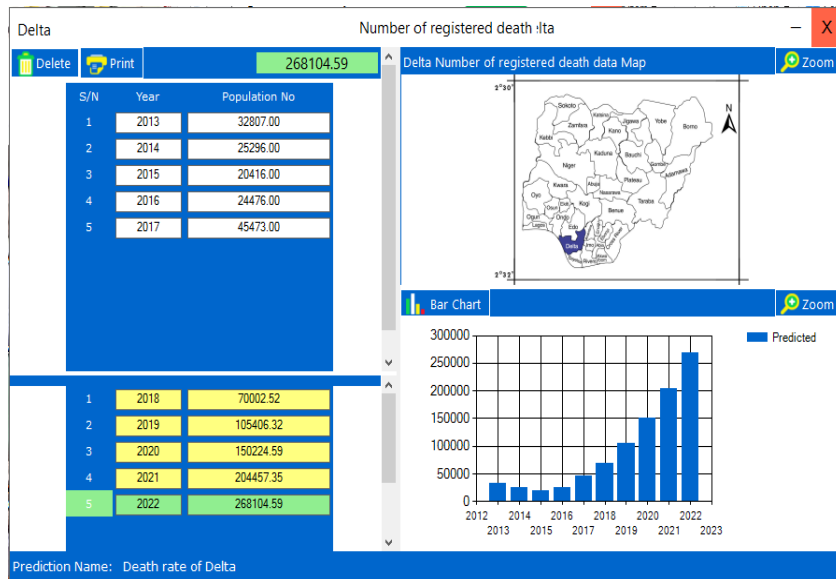
2.0 Methodology

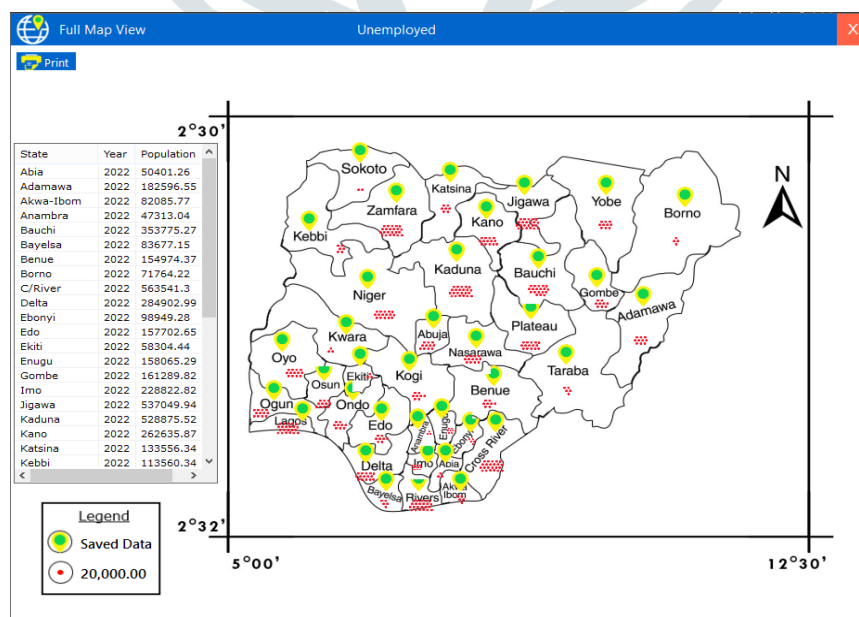
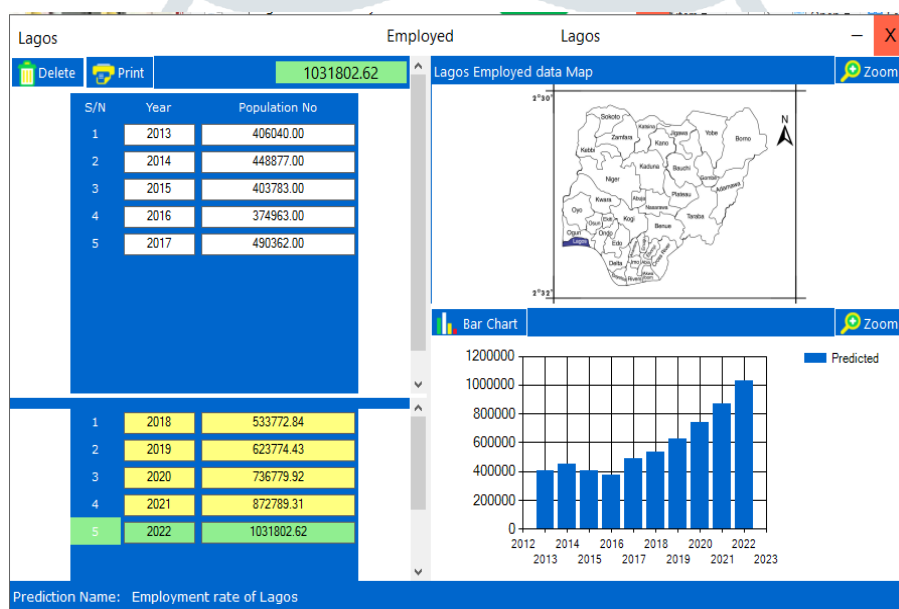
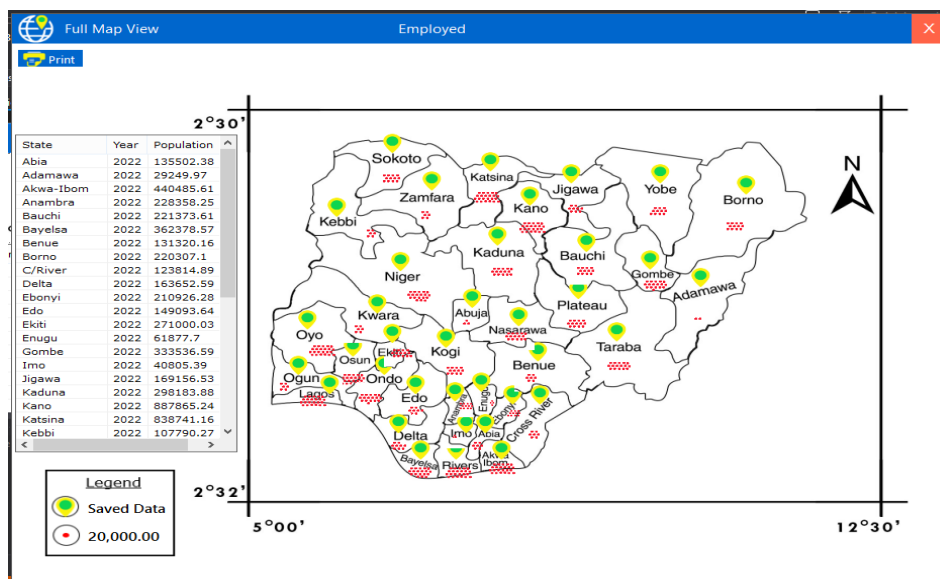
Object-Oriented Analysis and Design Methodology (OOADM) was used. The model used in analysis was $P_t = (X'X + \ln KI)^{-1} X'Y$. Where P_t is number of people at a future time, $X'X$ is a matrix with rows and columns; a three by three matrix, $\ln KI$ is an identity matrix with three by three matrix with a value 0.000099995 and Y is a centered n-vector that is three by one matrix. The calculations were done using Microsoft Excel to obtain a constant, and time coefficients which is $P_t = \beta_0 + \beta_1 t + \beta_2 t^2$

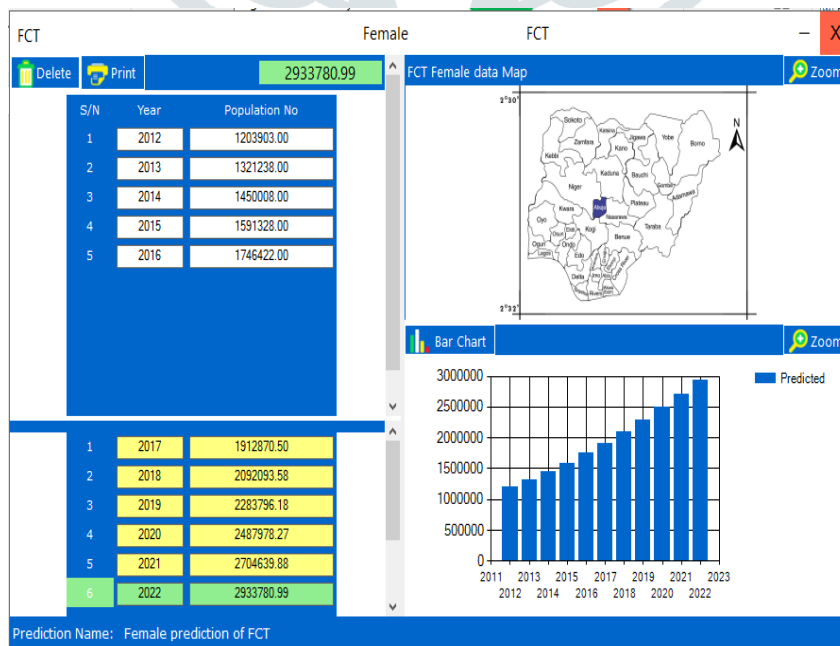
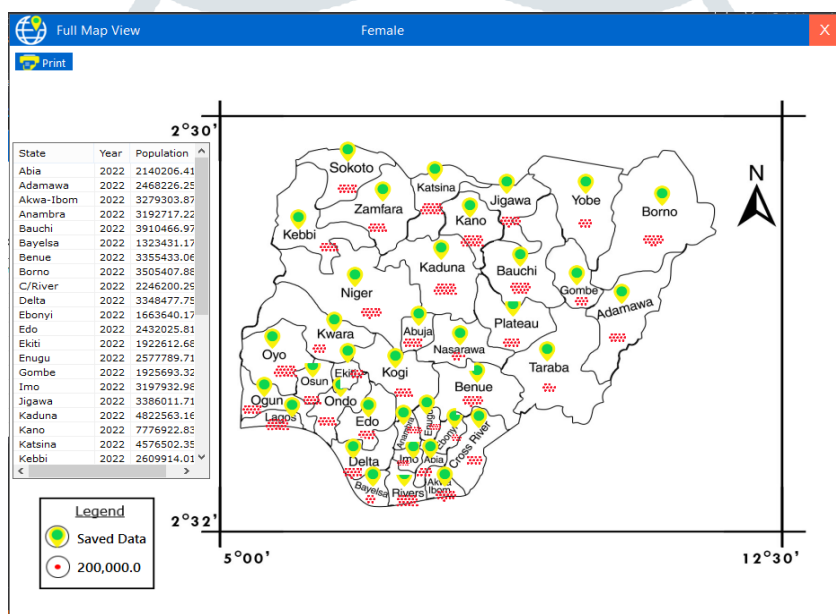
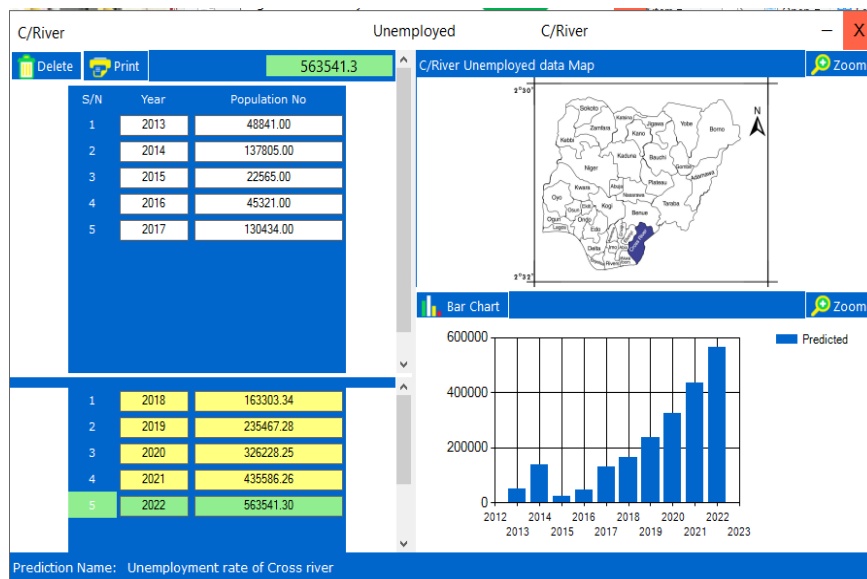
3.0 Results

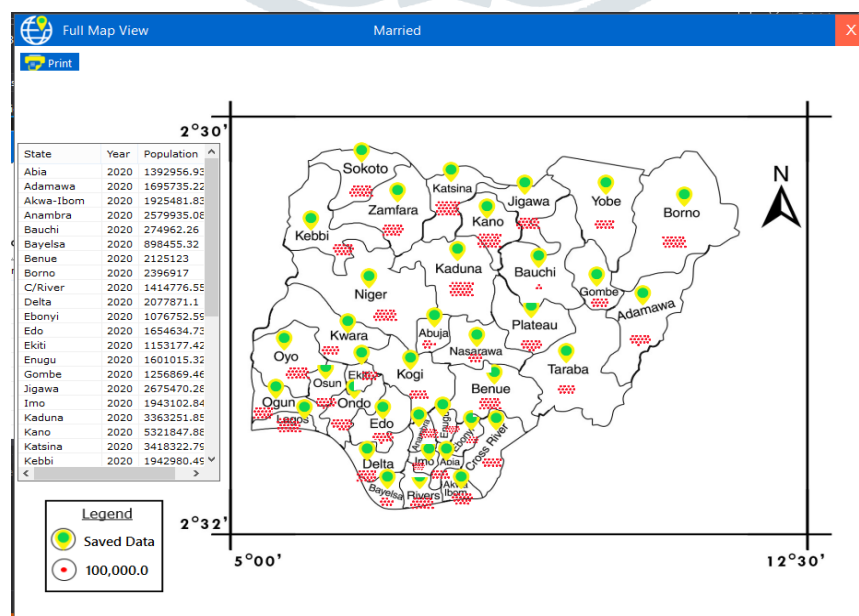
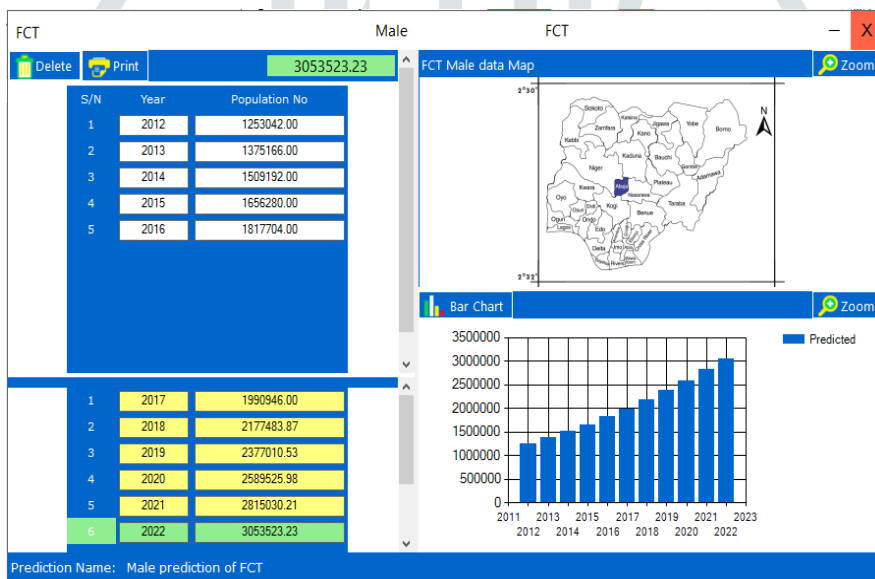
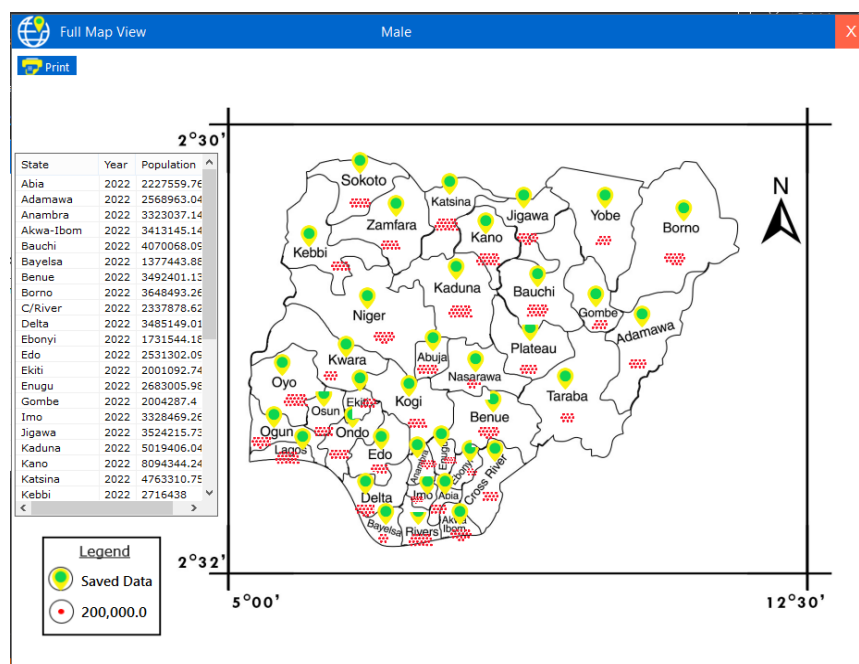
The output of this paper are in three forms; dot density map, table data showing number of data inputs and predicted data with year and bar chart illustrating prediction based on the predicted figures.

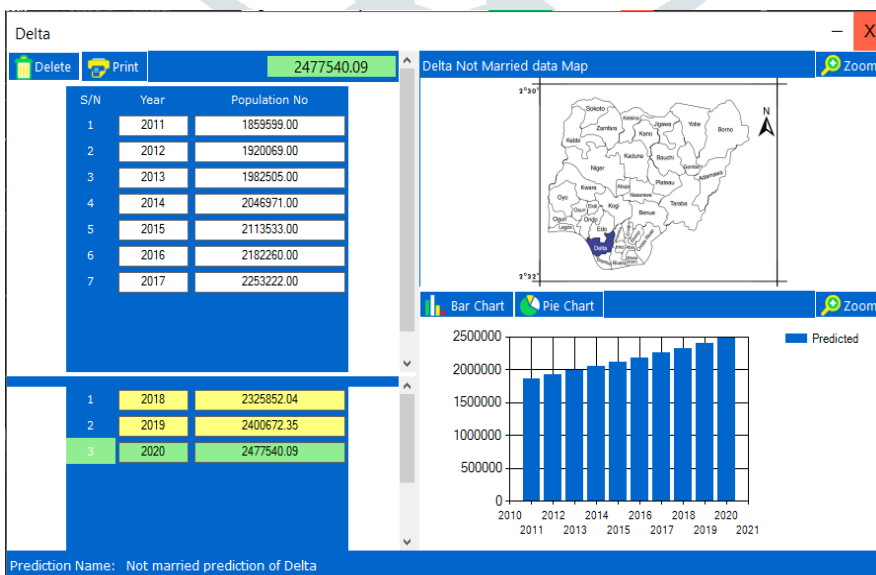
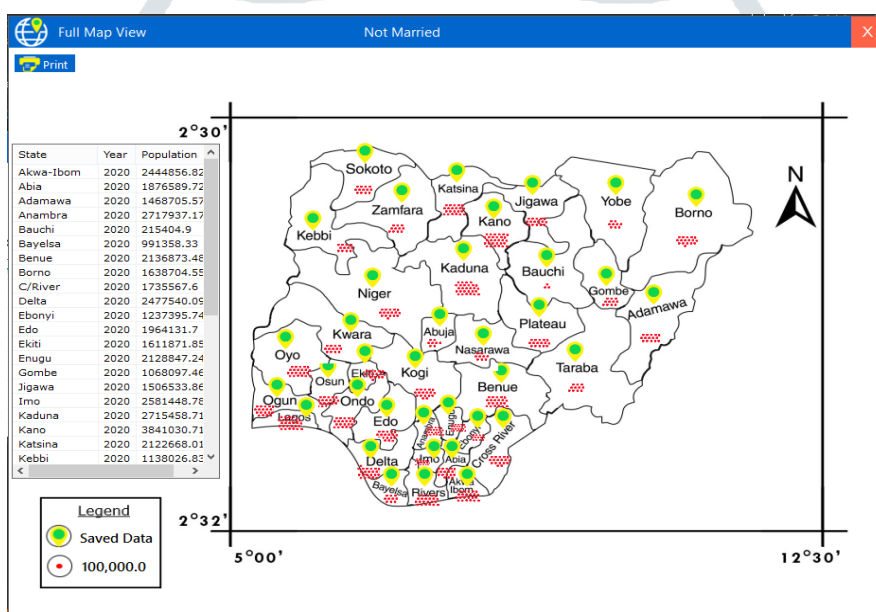
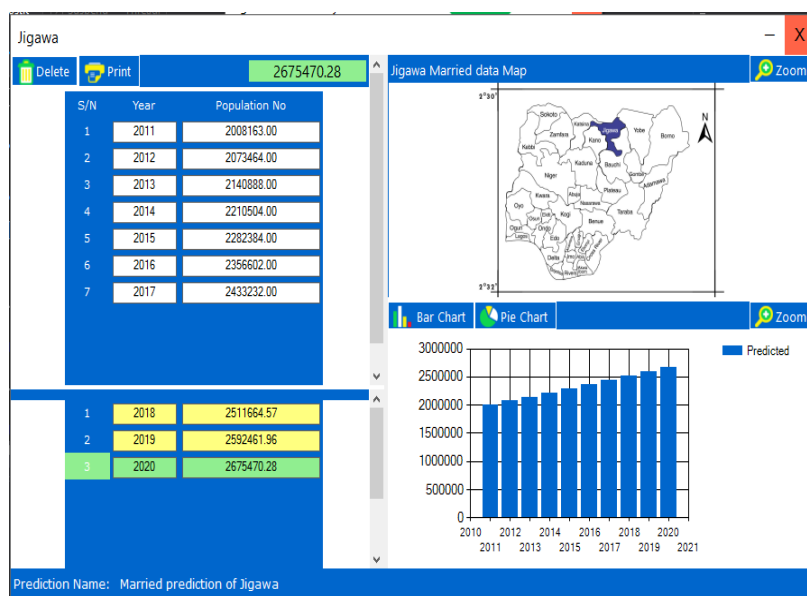


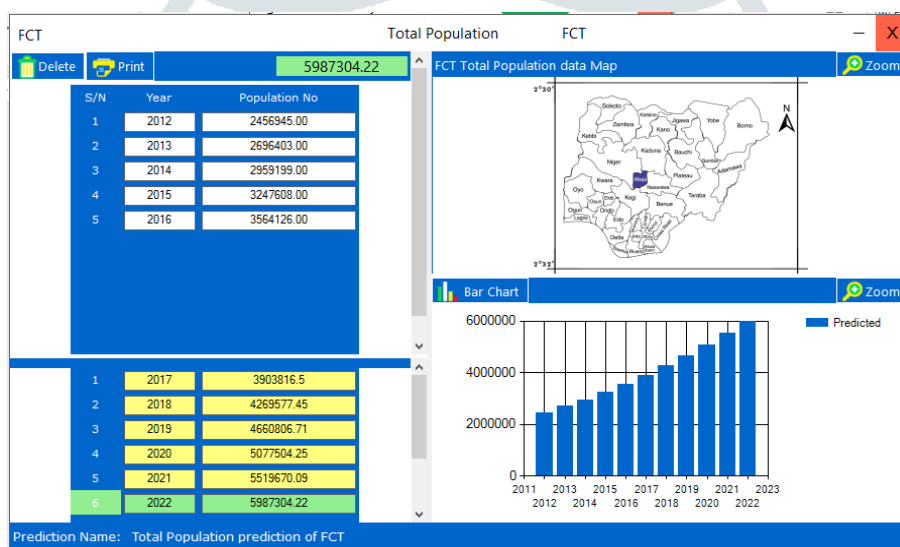
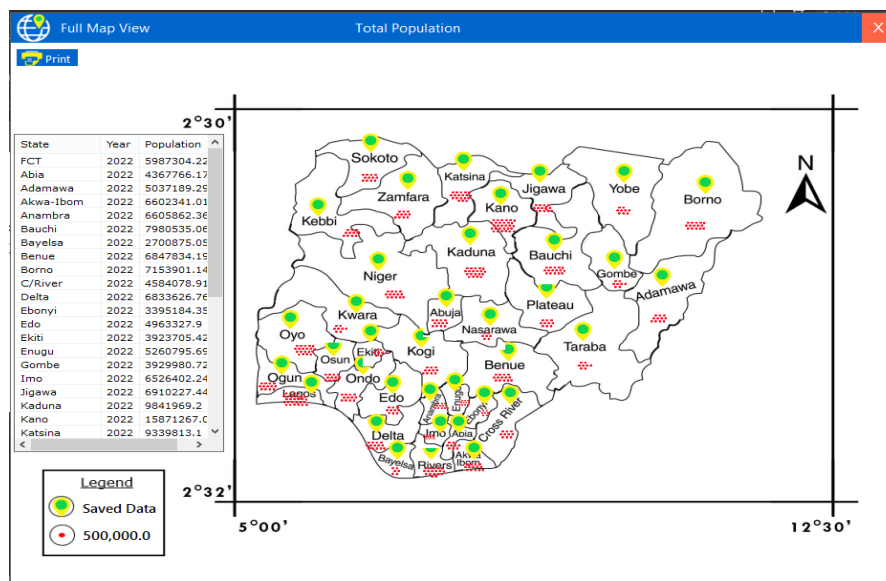












4.0 Contributions to Knowledge

National Population Commission has conducted a lot of censuses and surveys in the past. Therefore have a lot of repository of data that have not been mined for further analysis. This new model will assist in mining these repositories of data and bring out certain socio-economic trends, relationships and characteristics.

5.0 Conclusion

The new model is more effective in census analysis. It gives a better prediction than the existing ones not considering outliers and unknown parameters. The potentials of data mining for census results are enormous and it was harnessed in this research. This paper has made a huge difference to what has been in existence and thereby have added to the body of knowledge. Hence a call for the adoption of the new model is a right step in the right direction. Over estimation and under estimation which is prone to existing systems has been taken care of by adopting a narrow range smoothening constant. An appreciable level of computer knowledge to enable users successfully make prediction with the new system is required. Set of inputs which are known as observations (data input) are required by the system before it can make any prediction. Additionally, the number of observations cannot be less than five (5). These inputs can be entered via the keyboard by the user. Another significance of this study is that predicted demographic information is shown on the map.

REFERENCES

- [1] Census Enumerator's Manual (2006). National Population Commission Abuja: Nigeria.
- [2] Gaughan, A. E., Stevens, F. R., Linard, C., Jia, P., & Tatem, A. J. (2013). High resolution population distribution maps for Southeast Asia in 2010 and 2015. PloS one, 8(2), e55882.
- [3] Mohamed, M.G (2009).Scientific Data mining and Knowledge Discovery: Principles and Foundations. Springer: New York.
- [4] Micheal, B and Gordon, L (2012).Data mining Techniques.3rd Edition. Wiley: United States.
- [5] Kurt, T. (2012). Introduction to Data Mining, Wiley: Asia.
- [6] Acem, S. (2006).Machine Learning Algorithm and Data-mining. Cambridge University Press. London.
- [7] Ian. H.W , Eibe, F and Mark, A, H (2011). Data mining Practical Machine Learning Tools and Techniques.3rd Edition. Morgan Kaufmann: United States.

- [8] Micheal, R. B and David .J. H (2007).Intelligent Data Analysis.2nd Edition. Springer: New York.
- [9] Witten, Ian, H and Frank, E. (2011). Data mining: Practical Machine Learning Tools and Techniques.3rd Edition. CC Press: Elsevier.
- [10] Kenneth E. F. and Margaret, L.(2015).The Geographer's Craft Project, Department of Geography. The University of Colorado: Boulder.
- [11] GoodChild, M. F (2010).Twenty years of progress: GIScience: Journal of spatial Information science.Vol 2, pp 10-12.
- [12] Ramezankhani, A, Pournik, O, Shahrabi, A, Khalili, D, Azizi, F and Hadaegh, F. (2014).Applying decision tree for identification of a low risk population for type 2 Diabetes. Tehran lipid and Glucose study. International Journal on Diabetes Research and Clinical Practice. Elsevier Ireland Ltd. Vol 105, pg 391-398.
- [13]Manoj, K and Madhu, A. (2012).International Journal of Studies in Business and Economics.Pg 81- 94.
- [14]Md. M and Syed, S. H. (2012). Population Forecasts for Bangladesh, Using a Bayesian Methodology, Journal of Health, Population and Nutrition (JHPN) , 30(4): 456– 463.
- [15]David, A. S, Alan, S and Bob, S. (2010). Forecasting the Population of Census tracts by age,and sex. Springer Population Research and Policy Review, Volume 29, Issue 1, pp 47-63.
- [16]Kayode, J and Jimoh, M (2011).Used regression analysis with times series to periodically forecast stock. Advances in Information Mining Volume 4,Issue 1, pp57-66
- [17]Friedrich, H and Weixin L. (2013).ADULT AND YOUTH LITERACY- National, Regional and Global trends, 1985-2015.UIS Information Paper.
- [18]Saigal, S. and Mehrotra D. (2012). Performance Comparison of time series data using predictive data mining techniques. Advances in Information Mining Volume 4, Issue 1,pp.-57-66.
- [19]Rajasekhar,N and Karth T. V. R. (2014).Hybrid SVM Data Mining Techniques for Weather Data Analysis of Krishna District of Andhra Region. International Journal of Research in Computer and Communication Technology, Vol 3, Issue 7.

