

TITLE PAGE

**STATISTICAL MODELING OF STUDENTS  
ACHIEVEMENT**

A CASE STUDY OF GOVERNMENT DAY SECONDARY SCHOOL, MINNA

By

MUSA MUHAMMADU ISA

PGD/MCS/2003/2004/1123

**A PROJECT SUBMITTED TO THE DEPARTMENT OF  
MATHEMATICS/ COMPUTER SCIENCE.**

**IN PARTIAL FUL FILMENT OF THE REQUIEMENTS FOR THE AWARD OF A  
POSTGRADUATE DIPLOMA IN COMPUTE SCEINCE.**

FEDERAL UNIVERSITY OF TECHNOLOGY,

MINNA, NIGER STATE

**November 2004**

## **DECLARATION**

I hereby declare that this research work (project), a partial fulfillment of requirement for the award of post graduate diploma of computer science of computer science of, Federal University of Technology, Minna is a product of my research effort.

All sources of information collected for the project are clearly acknowledged by means of references.

MUSA MUHAMMADU ISA

SIGNATURE

DATE

## CERTIFICATION

I hereby certify that this is the original work of MR MUSA  
MUHAMMADU ISA PGD/MCS/2003/2004/1123. Submitted for the  
award of a postgraduate diploma in computer science of federal  
university of technology, Minna.

---

MALL. ISA AUDU  
PROJECCT SUPERVISOR

---

DATE

---

MR. L.N. EZEAKO  
HEAD OF DEPARTMENT

---

DATE

---

EXTERNAL EXAMINER

---

DATE

## **DEDICATION**

This project is dedicated to almighty Allah who has granted me the time and energy to go through my studies successful. It is dedicated to my beloved late mother Hajiya Fatima Kaka. May her humble soul rest in perfect peace, Ameen.

## ACKNOWLEDGEMENT

MY deep gratitude and appreciation to Almighty Allah who has been the source of my strength and health throughout the pursuit of my PGD. in computer science.

I am full of gratitude to my beloved parents whom their supports have helped me ever before I exist. in this line, I acknowledge the assistance of my brother Dr Mahmud and his family whom their support both morally and financially are immeasurable.

My sincere appreciation goes to my supervisor Mallam ISA AUDU who despite his tight schedule was always patient wit me, reading through the entire manuscript and made many valuable suggestions and corrections which p improved the style and quality of my work. Sincere gratitude is due to the H.O.D Maths/Computer science department MR L.N EZEAKO, DR AKINWANDE, DR YOMI AIYESIMI, PGD CORDINATOR DR U.Y ABUBAKAR, PRINCE BADMUS, MR JIYA MOHAMMED, MR HAKIMI and others lecturer of the department.

I am most grateful to my brothers; Mallam Abdullah Ndagaa, Aliyu Isa, Tanimu Mamman, my friends and class mates; Mohammed Abdulhakeem,

Alhaji Lawal Jibrin, Alhaji Sanusi Abubakar, Pharm.Umar Ndagi,  
Mohammed Abdulkadir, Zak Silas.

Constraints of space and time prevent mention of other friends and  
classmates whose contribution and assistance I cannot forget but cherish  
with silent gratitude

May Almighty Allah continue to bless us (Ameen).

## **ABSTRACT**

This project work which is aimed at modeling students achievement was carried out using regression analysis technique. Little introduction concerning case study was made.

Manual computation explanation and computation were carried out and subsequently; computerized method was carried out using statistical package for social sciences (SPSS). Results interpretations and graphs were later made and recommendation and conclusions made.

## TABLE OF CONTENT

TITLE PAGE	I
DECLARATION	II
CERTIFICATION	III
DEDICATION	IV
ACKNOWLEDGEMENT	V
TABLE OF CONTENT	VI-VII
ABSTRACT	VIII
<b><u>CHAPTER ONE</u></b>	
GENERAL INTRODUCTION	1-2
1.2 STATISTICS IN EDUCATION	2-3
1.3 CASE STUDY	3-6
1.4 OBJECTIVE OF THE STUDY	7
1.5 SCOPE AND LIMITATION	
1.6 DATA	7-8
1.7 LAYOUT OF THE STUDY	8-9
<b><u>CHAPTER TWO</u></b>	
LITERATURE REVIEW	10
2.1 REGRESSION ANALYSIS	10-11
2.1.1 PROBABILISTIC AND DETERMINISTIC MODELS	12-14
2.1.2 GENERAL FORM OF PROBABILISTIC MODEL	14
2.1.3 A FIRST ORDER (STRAIGHT-LINE) PROBABILISTIC MODEL	14-15
2.1.4 THE LEAST SQUARE APPROACH	15-18
2.2 MULTIPLE REGRESSION	18-19
2.3 ADVANTAGES OF USING REGRESSION ANALYSIS	19
2.4 AREAS OF APPLICATION OF REGRESSION ANAL	19
2.4.1 HEALTH	19
2.4.2 AGRICULTURAL SECTOR	20
2.4.3 BUSINESS SECTOR	20
2.5 RESIDUALS	21
2.6 MODEL TEST OF VARIABILITY	21-22
2.6.1 FITNESS OF MODEL / QUANTITATIVE DATA	22-23
2.6.2 TEST OF UTILITY	24-25



## CHAPTER ONE

### General Introduction

#### 1.1 Introduction

Just as problems are been expressed in words, figures, etc, so also solutions are expressed. But solutions are generally expressed logically; therefore logical expressions of solutions to problems are called models. Invariable, statistical techniques are unique models applied to problems solving with a view to making meaningful decision. Statistical techniques are applied in making predictions. These predictions though not usually exact but give close to exact solutions to the problems. These techniques can also be used in making estimations and graphical representation of data that eases the interpretation of voluminous set of data at a glance.

This project explains some areas of application of statistical tools employed in this project works that is multiple regression analysis in education. Among these areas is predicting students score based on his/her past performance.

The selected case study, which is Government day Secondary School Minna, has also been explained and hierarchical structure also presented in this chapter.

Objectives of the study have also been high lightened for easy understanding of what the project work is all about. Equally, the scope of the study has been

spelt out. The use of data that was used in carrying out the analysis and its type of method of collection was well explained in this chapter. And, lastly a layout of the study was also been presented.

## **1.2 STATISTICS IN EDUCATION**

A researcher is asked to decide between two different communities, which should be connected to a national grid. Among factors to be considered is which of the two uses electricity most at night. After several consultations and findings, the researcher comes up with a conclusion that the community with higher number of students should be connected. This judgment may be based on the fact that students make use of light at night for their assignments and revisions at night.

Through the Federal Office of Statistics, the Federal Government plans for infrastructure to be made available and allocated to each school based on the past population record of the students.

Universities usually have a minimum and maximum number of students it expects to gain admission into the school through JAMB based upon past and present statistics of students in the school. This among other reasons is to avoid over stretching the facilities of the school.

Graphical representation of average performance of various sets of students over a period of time easily conveys the message of the set that performed best and also worst.

Estimates on expected number of graduates for a nation at a particular time can be made based on past record of graduating students. A stable government through its ministry of education determines the work force expected and plans on how to absorb the labor force into its system.

Predictions and forecasting are among the major aspects of statistics study. With particular respect to educational sector, efficiency of different subject teacher can be predicted upon the students' performance. It is possible that the students all performed well in the two different subjects as a result of the two different teachers input. We may still to further to determine whose teaching technique is better. Infact, statistics can be applied in diverse aspects of education apart from the fact that it is a professional discipline of its own

### **3 CASE STUDY**

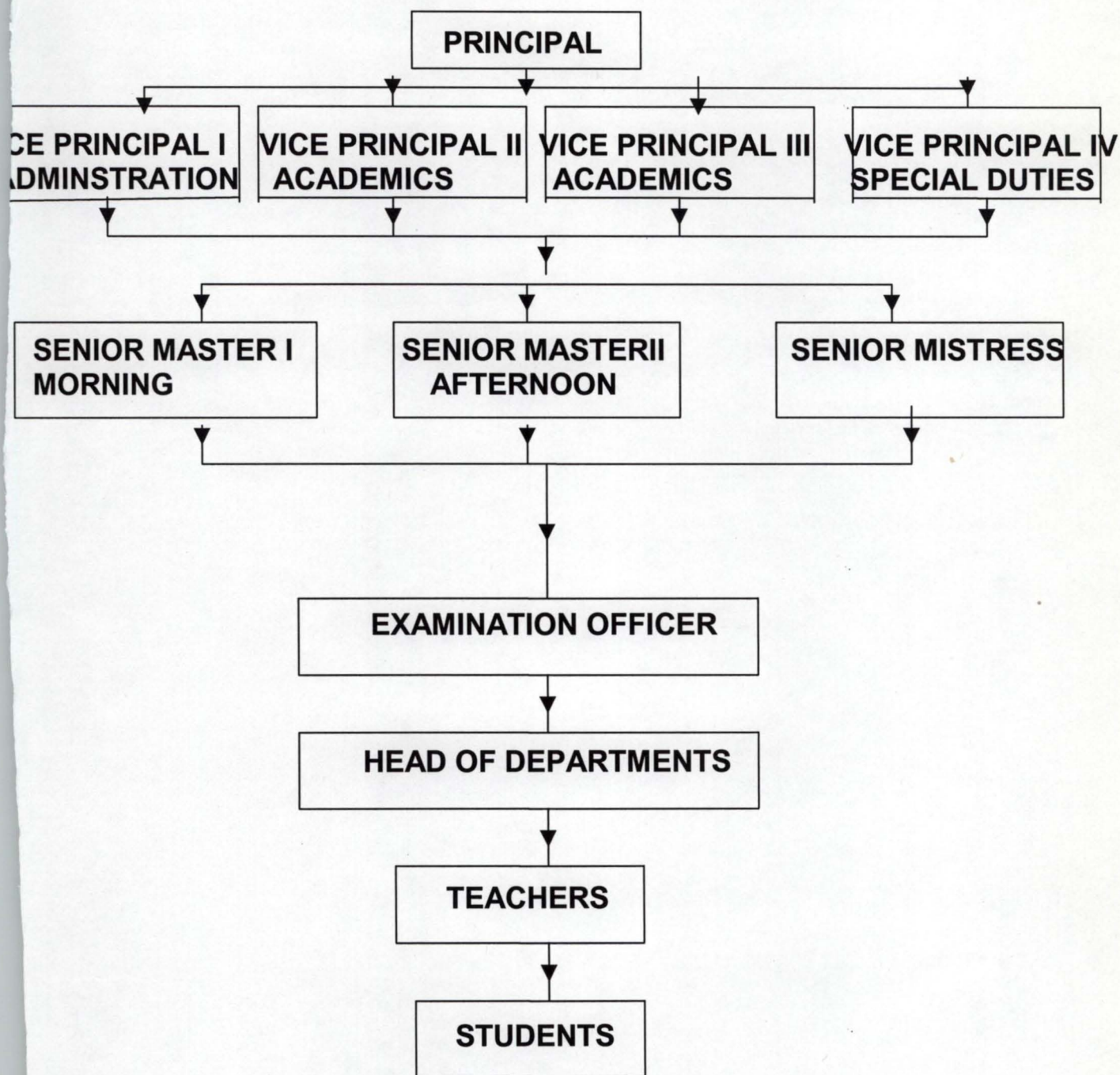
Government Secondary School Minna located along Bosso Road, Minna is the case study, the school is under the supervision of Niger state ministry of education and been regulated by national policy on education of the Federal Government of Nigeria.

The school was established in November 1979, with Mallam Mammam Batati as its current principal .the principal has the overall say of the school. He ensures that the school is working towards its stated goals. He acts as intermediary between the ministry of education and other staff of the school. Next to the principal are four vice principals, they are vice principal administration I, vice principal academics II, vice principal academics III, and vice principal special duties IV. The vice principal I by name Mallam Mohammad Gbaka receive directives from the principal directly. He is responsible for all administrative activities of the school in conduction with the principal. The two vice principal academics II and III are Mallam Mohammed Gana and Mrs. Eunice Gana respectively. They are in charge of assigning teachers to subjects and classed they are to teach. They make sure that each subject teacher forms his/her lesson note. Lastly, vice principal special duties IV by name Hajiya Aishatu M.Z. Kolo whose duties whose duties are to settle disputes between teachers and students and sometimes teachers and parents of students.

Immediately after the vice-principals are the senior masters and senior mistress. The school has senior master I morning, senior master II, afternoon

and senior mistress of the school. The two senior masters I and II are Mallam Sani Mohammed and Mallan Mohammed Adamu respectively. The senior mistress is by name Mrs. Elizabeth Yahya. They act as intermediary between the principal and teachers of the school in dissemination of information. They together prepare timetable of teachers to be on duty for each week.

# ORGANOGRAM



## **OBJECTIVES OF THE STUDY**

- i) To formulate statistical model of students achievement
- ii) Estimate the parameters in models
- ii) Fit the model and use the model for prediction
- iii) Screen variables to determine which have significant effect on the scores
- iv) Arrive at the most effective prediction equation.

## **1.5 SCOPES AND LIMITATION OF THE STUDY**

- i) The data used are for only two classes which are science class and social science classes
- ii) Only simple and multiple regressions are explained and used.

## **1.6 DATA**

The data used for this project work are students' examination scores for Science and social science class in Government Day Secondary School Minna.

Scores of students are entered into the record sheets by each subject teacher teaching that particular subject and these records are collected and stored in the examination office for future uses.

Some students after been issued their results indulge in alterations. Complaints of such cases by the parents are been attended to by the examination officer who makes use of the record sheets to confirm any discrepancies.

Some teachers mistakenly skip entering some student's scores. Such complaints are been dealt with by the clerks in the exams office consulting the record sheet of that particular class.

The school uses the record sheets to evaluate the average performance of a particular class or set. This is used as tool in measuring the performance of the students and also to ascertain whether the aim of the school is achieved. Record sheets are made available to students of higher learning who make use of the data to carry out their research.

## **.7 LAYOUT OF THE PROJECT**

Chapter one gives an insight on the need for statistics in educational system . a brief introduction was made explaining what the project work is all about. The case study was briefly explained along with organization structure diagram made. Objectives and scope of the study were also clearly spelt out and also the type of data used was explained.



Chapter two is a literature review, which explained some of the challenges faced during the project work. An introduction into simple and multiple regression is also in the chapter. This chapter explained some areas of application of regression analysis, and some advantages of using it. Test to confirm the variability of the tool used, which is regression analysis, is also explained.

Chapter three is concerned with the methodology of the existing system. This implies the manual computation of regression analysis. Some other techniques applied such as matrix model to solution were also explained.

Chapter four explains vividly the syntax used by the application program towards archiving the manually computed results in chapter three were made. Explanations on computed results by the application package were also made clearly.

Lastly, chapter five explains further some of the results in chapter four. In addition, summary of chapter three and four were made and a conclusion reached based on their comparism.

## **CHAPTER TWO**

### **LITERATURE REVIEW**

Isyaku (1997) in his work on effect of rainfall and temperature on rice production employed regression analysis technique to carry out his analysis. He wrote a program in QBasic to compute the regression for effect of annual rainfall on annual yield per bag. He also modeled effect of annual temperature on annual yield. Only simple regression was used in his analysis.

In this research work, program will not be used but a statistical package known as SPSS package, version 10 will be used to compute the regression analysis. Multiple regression will also be explained and employed in modeling linear equation for students performance.

#### **.1 Regression Analysis**

Based on the case study (G.D.S.S, MINNA), supervisors from secondary schools board are always interested in whether students who do well in term papers also do well in the final examination and whether students who did poorly in one also did poorly in the other. They are always interested too in knowing whether population of students per class affects the students performance in term papers and examinations or not. It is equally expected

that a student who is good in physics should also be good in mathematics. Efficiency of teachers can be determined based on some parameters. Based on past records of students performance, the instructor can estimate the final examination and conclude whether they will do well or not. Proper use of regression technique can give answer to these questions and many more of such related questions.

Regression is the process of fitting a straight line to two or more continuous variables, say  $x$ 's and ' $y$ '. Regression can be said to as a set of techniques, which utilizes the presence of an association between two variables to predict the value of one (the dependent variable or primary variable) from those of another (the independent or secondary variables).

The process of estimating and using a straight-line relationship is referred to as simple linear regression analysis. Basically, there are two methods of fitting as a straight line to the two variables. There are however other approaches to fitting s straight line but we will be restricted to the two below:

## I.1 PROBABILISTIC AND DETERMINISTIC MODEL

If we were to construct a model that suggests an exact relationship between variables, it would be called a deterministic model. For example, if we believe that  $y$ , the students' performance in final examination will be exactly 2 times  $x$ , his/her performance, then

$$Y=2x$$

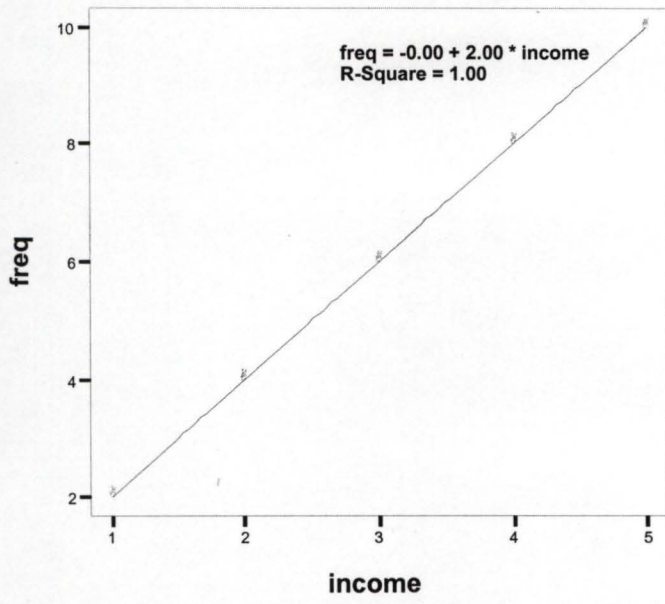
This represents a deterministic relationship between the variable  $y$  and  $x$ . and this implies that  $y$  can always be determined exactly when the value of  $x$  is known. This method does not give room for error. In other words, its values ( $y$ ) are exact. If on the other hand, we believe that there will be unexplained variation in the final

examination performance maybe caused by important but unincluded variables or by random phenomenon, we then cannot use the deterministic model but another model. This model accounts for this random error.

Therefore, if we assume that the final examination performance is related to the percentage of term paper then

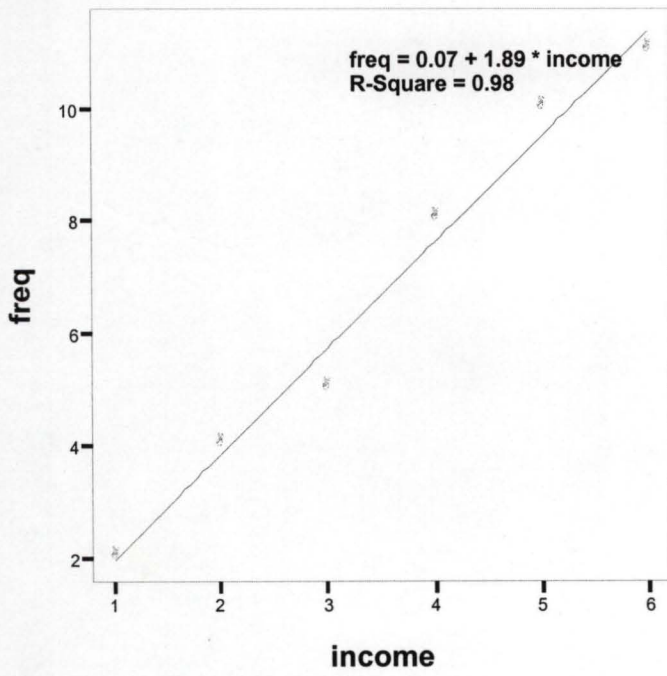
$$Y=2x + \text{Random Error}$$

Fig1



Linear Regression

fig2



Linear Regression

In fig1, all the values most fall exactly on the line because a deterministic model does not give room for error. On the other hand, fig2 shows that some points fall off the straight line implying a probabilistic model, giving room for error.

## 1.2 GENERAL FORM OF PROBABILISTIC MODEL

$Y = \text{Deterministic component} + \text{Random Error}$

Where  $y$  is the variable of interest. We always assume that the mean value of the random error equals '0'. This implies that the mean value of  $y$ ,  $E(y)$  equals the deterministic component of the model. That is

$E(y) = \text{Deterministic component.}$

## 1.3 A FIRST ORDER (STRAIGHT LINE) PROBABILISTIC MODEL

$$Y = \beta_0 + \beta_1 x + \epsilon$$

Where  $y$  = Dependent or response variable (variable to be modeled)

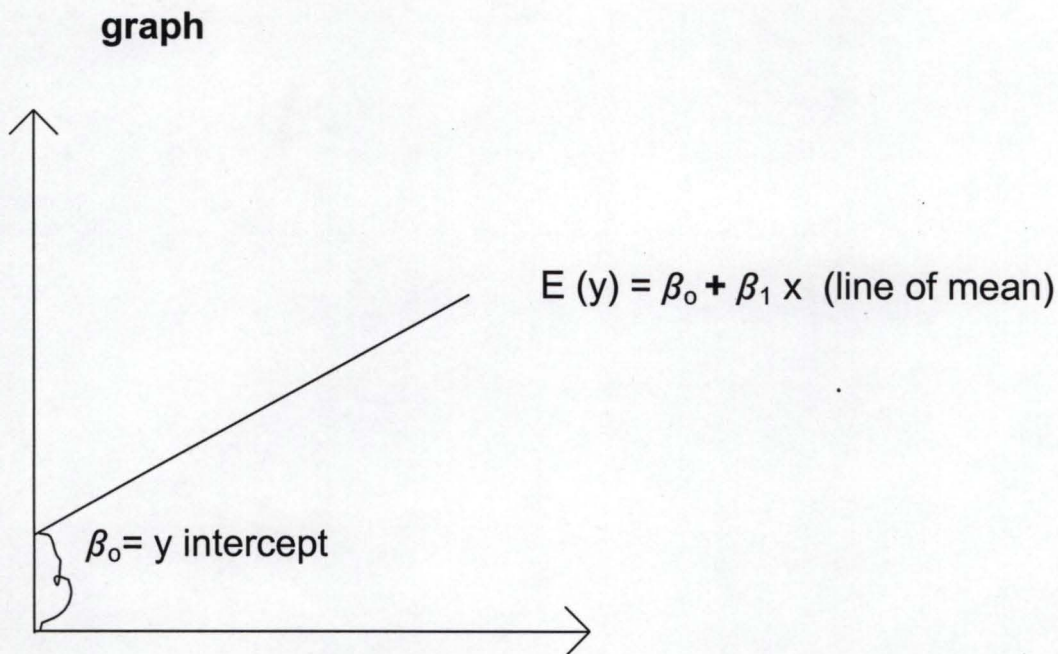
$X$  = Independent or predictor variable (variable used as predictor of  $y$ )

$E(y) = \beta_0 + \beta_1 x = \text{deterministic component}$

$\epsilon$  = Random error component or residual

$\beta_0$  =  $y$  intercept of the line, the point at which the line cuts through the  $Y$ -axis

$\beta_1$  = Slope of the intercept, that is the amount of increase (or decrease in the deterministic component of  $y$  for every 1 unit in  $x$ )



Note that in the probabilistic model, the deterministic component is referred to as the line of mean, because the mean of  $y$ ,  $E(y)$  is equal to the straight-line component of the model.

That is

$$E(y) = \beta_0 + \beta_1 x$$

#### 1.4 THE LEAST SQUARE APPROACH

This method is used to determine the regression line using mathematical method called the least square principle. It determines a regression equation

by minimizing the sum of squares of the vertical distance between the actual Y values and the predicted values  $\hat{y}$ .

$$\hat{Y} = a + bx \text{-----(1)}$$

$\hat{Y}$  is the fitted estimate of Y based on a, b, and x

'a' is the Y-intercept. It is the estimated value of Y when  $x = 0$ . It can also be said to as the estimated values of Y where the regression line crosses the Y axis when x is equal to zero

Mathematically,

$$\hat{a} = \bar{y} - b\bar{X} \quad \text{or}$$

b is the slope (that is, the average change in y per unit change in x).

Another name for b is the regression coefficient

Mathematically,

Where

$x_i$  = value of the independent variables

$y_i$  = values of the dependent variables

n = Number of items in the sample

The above equation (1) is the same for equation of a straight-line graph in general and  $\hat{a}$  is the intercept (the average value of y when  $x=0$ ). The said equation is also referred to as linear equation because when the values are plotted and joined together, it gives us a straight-line graph. The equation is a simple, two variable regressions; this is because only two variables can be analyzed with this method.



The accuracy of an estimate solemnly depends on the extent to which the regression equation and its corresponding graph actually fit the data. The regression line is always the line of best fit. One way of obtaining the line of best fit is to plot the values and ion as many points as possible with a single straight line. We then locate where this line cuts the Y-axis to give us the constant 'a'

Line of best fit is a situation of drawing the graph going by some criteria. Among, we could adopt would be to minimize the magnitude of the error involved, that is the difference between the values of the points on the diagram (Y) and the values of the points on the regression line we have drawn (y) and the values of point on the regression line we have drawn. From that, we can see that line of best fit is that line which minimizes the sum of square s of the errors.

The values of  $\hat{a}$  and b will be unknown in almost all practical applications of regression analysis. The process of developing a model, estimating the unknown parameters and using the model can be viewed as the five step procedures stated below

- 1) Postulate the deterministic component of the model that related the mean  $E(y)$ , to the independent variable x.
- 2) Use the sample data to estimate unknown parameters in the model

- 3) Specify the probabilistic distribution of the random error term and estimate the standard deviation of this distribution.
- 4) Statistically, evaluate the usefulness of the model.
- 5) When satisfied that the model is useful, use it for prediction, estimation and other purpose.

The above-mentioned steps will be applied in a later chapter of this project.

## 2.2 MULTIPLE REGRESSION

A Multiple regression analysis involves fitting the model to data set, testing the utility of the model and using it for estimation and prediction. It is a technique developed for modeling the mean of  $y$  as a function of two or more independent variables. As such, probabilistic models that include more than one independent variable are called multiple regression models.

The general form of this model is

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

Which minimizes the sum of square error (SSE)

$$SSE = \sum (y - \hat{y})^2$$

The primary difference between fitting the simple and multiple regression models is computational difficulty. The  $(k+1)$  simultaneous linear equations that must be solved to find the estimated coefficients of

$\beta_0, \beta_1, \beta_2, \dots, \beta_k$  are difficult and sometimes nearly impossible to solve with manual computation. However, with computer and computer software such as SAS, MINITAB and SPSS, the computation becomes easier.

### 2.3 Advantages Of Using Regression Analysis

- 1) In the first place, we square the errors, which takes care of the negative signs and as such avoids errors canceling out the positive ones.
- 2) Squaring tends to emphasize the larger errors and minimizing criterion means that we are trying to avoid in the larger errors.
- 3) All points (variables) are taken into consideration, hence equal representation.

## 4 AREAS OF APPLICATION OF REGRESSION ANALYSIS

4.1 **HEALTH:** - Health researcher may want to find out the relationship between hemoglobin and creatinine to see if hemoglobin and creatinine depends on renal function (as measured by the serum creatinine Concentration). Can patients with renal failure be easily exposed to anemia? Does the amount of drugs sold in a pharmaceutical locality depend on? the number of patients who visit the local hospital? The researcher is trying to find out the rat of self-medication.

The regression technique can be used in diverse areas of health.

2 **AGRICULTURAL SECTOR:** - in a highly mechanized farm, projection on yield can be made based on the quality of seeds planted, number of bags of fertilizer applied, mean annual rainfall, number of hectares covered, e.t.c. Different brands of fertilizer applied on a piece of farmland at different time can be analyzed to know the better of the fertilizer brands that yield better and quality produce.

2.4.3 **BUSINESS SECTOR:** -The manager of an industry may want to know which of the factors affects the gross profit of the firm among the following variables; the number of employees, consecutive dividends, beginning inventory, etc.

Circulation manager of a newspaper may want to confirm what variables relate to the number of subscriptions to the paper. Variables such as metropolitan population, the advertising budget of the paper and median family income in the metropolitan area. With proper regression analysis, solution will be obtained.

Regression analysis techniques can be employed for analysis in diverse area of human endeavor and areas of study (discipline).

## 5 RESIDUALS

When regression equation is used to determine the value of a variable 'y' from those one or more independent variable(s) 'x', the estimates  $\hat{y}$  will usually fall short of complete accuracy, geometrically, speaking, the data points will not all fall precisely upon the straight line, plane or hyper plane specified by the regression. The discrepancies  $(y - \hat{y})$  on the predicted variable are known as residual. When using regression methods, the study of the residuals is of great importance, because they form the basis for measures of the accuracy of the estimates and of the extent to which the regression model gives a good account of the data in question.

Mathematically, Residual

$$S_y = \sqrt{\frac{\sum (y - \hat{y})^2}{n-1}}$$

### Model Test Of Variability

Confirmatory tests are carried out to ascertain the overall utility of a model. Despite the fact that a t-test can be conducted in checking the overall utility of a model, it has some lapsed. Fitting a model to two or more quantitative variability and using t-test to check its variability may include a large number of insignificant variables and exclude some useful ones. Hence, if we want to

test the utility of a multiple regression model, we will need a more comprehensive test that combines all the  $B_i$ 's. Also we may like to find some statistical quantities that measure how well the model fits the data.

## 6.1 Fitness Of A Model On Quantitative Data

The method for finding the fitness of a model on quantitative data can be measured using a statistical technique called multiple coefficient of determination,  $R^2$  and it is expressed as

$$R^2 = 1 - \frac{SSE}{SS_{yy}}$$

Where  $SSE = \sum (y - \hat{y})^2$

$$SS_{yy} = \sum y^2 - \frac{(\sum y)^2}{n} = \frac{\text{Explained variability}}{\text{Total variability}}$$

$R^2$  is interpreted as the fraction of the sample variation of the  $y$  values (measured by  $SS_{yy}$ ) that is explained by the least square prediction equation.

Thus,  $R^2 = 0$  implies a complete lack of fit of the model to the data and  $R^2 = 1$  implies a complete perfect fit with the model passing through every data point.

However, we should note that  $0 < R^2 < 1$ .  $R^2$  is a sample statistics that tells us how well the model fits that the data and therefore represents a measure of the usefulness of the entire model.

A large value of  $R^2$  computed from the sample data does not necessarily tell us that the model provides a good fit to all of the data points in the population.

Cases of models that comprise of more than two parameters will always provide a perfect fit. Also,  $R^2$  is always equal to 1 for model with more than two parameters. As such, when dealing with more than two parameters  $\beta_i$ , a model adequacy; the adjusted multiple coefficient of determination denoted by  $R_a^2$  is preferred.

Mathematically,

$$\begin{aligned} R_a^2 &= 1 - \frac{(n-1)(SSE)}{n-(k+1) SS_{yy}} \\ &= 1 - \frac{n-1}{n-(k+1)} * (1 - R^2) \end{aligned}$$

Note that  $R_a^2$  is always less than or equal to  $R^2$ , though they have similar interpretation.  $R^2$  and  $R_a^2$  are only sample statistics and as such are not reliable in judging usefulness of the model. Therefore, a test of utility to the multiple regression models is required.

## 2 Test Of Utility

Test of utility involves a test of hypothesis encompassing all the  $\beta$  parameters (except  $\beta_0$ ) in a model.

Hypotheses can be said to as statement about a population developed for the purpose of testing. Hypothesis testing; a procedure based on a sample evidence and probability theory to determine whether the hypothesis is reasonable statement.

Null hypothesis, denoted as  $H_0$ ; is a statement about the values of a population parameter, the alternative hypothesis ( $H_1$ ) describes what you will conclude if you reject the null hypothesis.

Alternative hypothesis is a statement that is accepted if the sample data provide enough evidence that the null is false.

In making hypothesis formulation for multiple regressions, we would test:

$$H_0: \beta_1 = \beta_2 = \beta_3 = 0$$

$H_1$ : at least one of the coefficients is non-zero

To confirm the utility of the null hypothesis, an F-statistics test can be computed and used

Mathematically

$$F\text{-statistics} = \frac{SS_{yy} - SSE/k}{SSE/[n-(k+1)]}$$



$$= \frac{R^2 / k}{(1 - R^2) / [n - (k + 1)]}$$

We then compare the calculated F-statistic to a tabulated F-value with k degree of freedom (df) in the numerator and  $[n - (k + 1)]$  df in the denominator. If the calculated value exceeds the tabulated value, we reject the null hypothesis and conclude that at least one of the model coefficients  $\beta_1, \beta_2, \beta_3, \dots, \beta_k$  is non zero.

Therefore, this global F-test indicates that the first order model

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

is usually for predicting the dependent variable for example say final examination score.

# CHAPTER THREE

## SYSTEM ANALYSIS AND DESIGN

### 3.1 Introduction

As earlier explained, this project work will be based on regression analysis technique. The set of data to be used are exams scores of students in the senior classes of science and social science classes. The data variables used are mathematics, physics and chemistry scores from the science class. While data variables used for social science class are Government, Economics and commerce subject.

### 3.2 Techniques-:

In carrying out the regression analysis, certain variables (subjects) are made dependent while some independent variables.

For the science class, the criteria for selecting the dependent variable is that a student who is good in physics must be very good in mathematics which was used for simple regression. This implies also that for a student to be good in physics, he/she must be well knowledgeable in mathematics subject. Equally, for the multiple regressions, physics and chemistry were selected as independent variables since a student who is good in those subjects is

also expected to be very good in mathematics. The same reason goes for the social science class where Economics was selected as dependent variable with Government and Commerce subjects as independent variables.

### 3 Least square computation

#### 3.1 Least-Square Computations for sheet1

Computation using the least square approach as earlier explained involves a dependent variable with only one independent variable. From score Sheet1, assuming we use Maths score as dependent variable and Physics scores as independent variable. Maths performance score can to some extent be estimate using the least square approach as thus-:

Let Maths score =  $\hat{y}$

Physics score =  $x_i$

Then  $\hat{y} = \hat{a} + b x_i$

$\hat{a} = - b x_i$

$$b = \frac{n \sum xy - (\sum x)(\sum y)}{n (\sum x^2) - (\sum x)^2}$$

$$b = \frac{44 \cdot 80189 - 1781 \cdot 1923}{44 \cdot 78463 - (1781)^2}$$

$$b = 0.3689 \approx 0.37$$

$$\hat{a} = \frac{1923}{44} - 0.3689 \frac{1781}{44}$$

$$\hat{a} = 28.77$$

$$\hat{y} = 28.77 + 0.37 x_i$$

Accordingly, a linear equation can be deduced for mathematics score using chemistry scores

$$b = \frac{44 * 72820 - 1595 * 1923}{44 * 63635 - (1595)^2}$$

$$= \frac{13689}{255915}$$

$$= 0.535 \approx 0.54$$

$$\hat{a} = \frac{1923}{44} - 0.535 \frac{1595}{44}$$

$$\hat{y} = 24.31 + 0.54 * x_i$$

### 3.3.2 Residual Computation For Least Square

The residual values are used to estimate the random error and to check the regression assumptions. The residual values are plotted on the vertical axis against the variable (dependent)  $y_i$ , on the horizontal axis

$$S_y = \sqrt{\frac{\sum (y - \hat{y})^2}{n-1}}$$

$$= \sqrt{\frac{1580.624}{44 - 1}}$$

$$= \sqrt{36.76}$$

$$= 6.06$$

The above computed residual is for mathematics with Physics as independent variable.

We then computed residual for mathematics with chemistry as independent.

$$\begin{aligned}
 S_y &= \sqrt{\frac{\sum (y - \hat{y})^2}{n-1}} \\
 &= \sqrt{\frac{1930.94}{4}} = \sqrt{44.9} = 6.70
 \end{aligned}$$

### 3.3 Least Square Computations For Sheet2

We know consider sheet2 by changing the variables that is making use of economics marks as dependent variable  $y$  and commerce independent

$$\hat{y} = \hat{a} + b x_i$$

$$\begin{aligned}
 b &= \frac{n \sum xy - (\sum x)(\sum y)}{n(\sum x^2) - (\sum x)^2} \\
 &= \frac{46 * 93177 - 2205 * 1909}{46 * 120127 - 2205^2} \\
 &= \frac{76797}{663817} = 0.116
 \end{aligned}$$

$$\begin{aligned}
 \hat{a} &= y - b x_i \\
 &= \frac{1909}{46} - 0.166 * \frac{2205}{46} \\
 &= 41.50 - 5.56 \quad \hat{a} = 39.9
 \end{aligned}$$

We can equally derive a linear equation for eons score using government score as independent marks

$$\hat{y} = \hat{a} + b x_i$$

$$b = \frac{n \sum xy - (\sum x)(\sum y)}{n (\sum x^2) - (\sum x)^2}$$

$$b = \frac{46 * 88050 - (1932)(1909)}{46 * 99436 - 1932^2}$$

$$= \frac{362112}{841432}$$

$$= 0.43$$

$$\hat{a} = y - b x_i$$

$$= \frac{1909}{46} - 0.43 * \frac{1932}{46} = 41.5 - 18.06$$

$$\hat{a} = 23.44$$

$$\hat{y} = 23.44 + 0.43 * \text{Gov't score}$$

### 3.4 Residual Computation For Least Square Of Sheet2

$$\begin{aligned} S_y &= \sqrt{\frac{\sum (y - \hat{y})^2}{n-1}} \\ &= \sqrt{\frac{12868}{45}} \\ &= \sqrt{285.96} = 14.10 \end{aligned}$$

### 3.4 Multiple Regression Computation

From record sheet1, we will have more than one independent variable since we are computing a multiple regression. Mathematics score will be our dependent variable with physics and chemistry scores as independent variables denoted by  $x_1$  and  $x_2$  respectively. Note that more than two independent variables can be used but the computation becomes more complex and in some cases impossible.

To compute and obtain the linear equation for three variables, we generate the following three equations using

$$\hat{y} = \hat{a} + b_1x_1 + b_2x_2 \quad \text{giving rise to the three equations as follows}$$

$$\Sigma y = an + b_1 \Sigma x_1 + b_2 \Sigma x_2$$

$$\Sigma x_1y = a \Sigma x_1 + b_1 \Sigma x_1^2 + b_2 \Sigma x_1x_2$$

$$\Sigma x_2y = a \Sigma x_2 + b_1 \Sigma x_1x_2 + b_2 \Sigma x_2^2$$

The unknown variables are computed using any convenient method to solving matrix problem. We then substitute the various values from the table given.

From record sheet1

$$1923 = 44a + 1781b_1 + 1595 b_2$$

$$80189 = 1781a + 78463 b_1 + 67747 b_2$$

$$72820 = 1595a + 67747 b_1 + 63635 b_2$$

Extracting out the coefficients to form a corresponding matrix as A, to use the inverse method

$$A = \begin{bmatrix} 44 & 1781 & 1595 \\ 1781 & 78463 & 67747 \\ 1595 & 67747 & 63635 \end{bmatrix}$$

$$AX = b$$

$$\text{But } b = A^{-1}X$$

To find  $A^{-1}$ ,

$$|A| = \begin{vmatrix} 44 & 1781 & 1595 \\ 1781 & 78463 & 67747 \\ 1595 & 67747 & 63635 \end{vmatrix}$$

$$|A| = 44(403336996) - 1781(5277470) + (-4491078) \\ = 1184384344$$

We then find and form the matrix of cofactors 'N'

$$\begin{bmatrix} 403336996 & -5277470 & -4491078 \\ -5277470 & 255915 & 191093 \\ -4491078 & -140173 & 280411 \end{bmatrix} = N$$



There after, we transpose the matrix of cofactors

$$N^T = \begin{bmatrix} 403336996 & -5277470 & -4491078 \\ 5277470 & 255915 & -140173 \\ -4491078 & 191093 & 280411 \end{bmatrix}$$

$$A^{-1} = \frac{1}{|A|} * N^T$$

$$= \frac{1}{1184384344} * \begin{bmatrix} 403336996 & -5277470 & -4491078 \\ 5277470 & 255915 & -140173 \\ -4491078 & 191093 & 280411 \end{bmatrix} = \text{Adj.A}$$

$$\frac{1}{A} * \begin{bmatrix} 403336996 & -5277470 & -4491078 \\ 5277470 & 255915 & -140173 \\ -4491078 & 191093 & 280411 \end{bmatrix} \begin{bmatrix} 1923 \\ 80189 \\ 72820 \end{bmatrix} = \begin{bmatrix} a \\ b_1 \\ b_2 \end{bmatrix}$$

$$\hat{a} = 654.8693819 + (-357.3122559) + (-276.1268347)$$

$$\hat{a} = 21.43$$

$$b_1 = -8.568649916 + 17.32678082 + (-8.618315424)$$

$$b_1 = 0.140$$

$$b_2 = 0.46$$

Therefore,

$$\hat{y} = 21.43 + 0.14b_1 + 0.46b_2$$

### 3.4.1 Residual For Multiple Regression

$$\begin{aligned} S_y &= \sqrt{\frac{\sum (y-\hat{y})^2}{n-1}} \\ &= \sqrt{\frac{1580.624}{44-1}} \\ &= \sqrt{36.754} = 6.06 \end{aligned}$$

### 3.4.2 Multiple Coefficient Of Determination

$$R^2 = 1 - \frac{SSE}{SS_{yy}}$$

$$\text{Where } SSE = \sum (y-\hat{y})^2$$

$$= 3335.159$$

$$SS_{yy} = \sum y^2 - \frac{(\sum y)^2}{n}$$

$$= 87379 - \frac{1923^2}{44}$$

$$= 3335.159$$

$$= 1 - \frac{1580.624}{3335.159}$$

$$= 1 - 0.4739$$

Therefore,

$$R^2 = 0.53$$

### 3.4 Computation For Adjusted Coefficient Of Determination

$$R_a^2 = 1 - \frac{[n-1](SSE)}{n-(k+1) SS_{yy}}$$
$$= 1 - \frac{n-1}{n-(k+1)} * (1-R^2)$$

Where k is the number of independent variables

$$R_a^2 = 1 - \left[ \frac{44-1}{44-(2+1)} \right] \quad (1-0.53)$$

$$= 1 - \left[ \frac{43}{41} \right] \quad (0.47)$$

$$= 1-0.493 \quad = 0.507$$

### 3.3.5 Test Of Hypothesis

$$H_0: b_1 = b_2$$

$$H_1: b_1 \neq b_2$$

$$F\text{-statistics} = \frac{SS_{yy} - SSE/k}{SSE/[n-(k+1)]}$$
$$= \frac{R^2 / k}{(1-R^2) / [n-(k+1)]}$$
$$= \frac{0.53 / 2}{(1-0.53)/[44-(2+1)]}$$

$$= \frac{0.265}{0.01146} = 23$$

### 3.5 Analysis Of Variance

An experiment is a planned procedure that gathers comparative data under controlled conditions. When experiments are performed, conclusions are reached after analyzing the results, based on the data with two or more independent variables; ANOVA is used to determine the difference between the averages.

The research questions are

- 1) How strongly are variables associated?
- 2) Can score on a target variable (or category membership, if the variables are qualitative) be predicted from data on other variables?
- 3) Is a difference (between averages) significant?

The analysis of variance (ANOVA) is actually a whole set of technique, each based upon a model of how the data were generated and culminated in tests that are appropriate for that are appropriate for that particular model only. As such, it is important to identify ANOVA correctly, in order to choose the right test.

In ANOVA, a factor is a set of related conditions or categories.

The conditions or categories making up a factor are known as levels.

Some factors are said to be between, this implies that participants are

treated under the same conditions. Other factor can be said to as within subjects, that is, participant is treated under all the conditions (levels) making up the factor.

In the course of this project work, I am trying to analyze the relationship between  $\mu_1, \mu_2, \mu_3, \dots, \mu_k$ . And since we are making use of sample, the data will be in a tabular form as shown below

				<b>Total</b>
$X_{11}$	$X_{12}$	$\dots$	$X_{1j}$	$\dots$
			$X_{1c}$	$X_{1.}$
$X_{21}$	$X_{22}$	$\dots$	$X_{2j}$	$\dots$
			$X_{2c}$	$X_{2.}$
$\dots$	$\dots$	$\dots$	$\dots$	$\dots$
			$\dots$	$\dots$
$X_{r1}$	$X_{r2}$	$\dots$	$X_{rj}$	$\dots$
			$X_{rc}$	$X_{r.}$
Total	$\bar{X}_{.1}$	$\bar{X}_{.2}$	$\dots$	$\bar{X}_{..}$

One way classification model assumes a hypothesis  $\mu_1 = \mu_2 = \dots = \mu_k$ .

The model is then applied to actual data to find out whether the hypothesis is true. An example is that given below

$X_{i1}$	$X_{i2}$	$X_{i3}$
$X_{11}$	$X_{12}$	$X_{13}$
$\dots$	$\dots$	$\dots$
$X_{41}$	$X_{42}$	$X_{43}$
$\Sigma X_{i1}$	$\Sigma X_{i2}$	$\Sigma X_{i3}$

Example,

$$\mu_{.2} = \frac{\sum X_{i2}}{N_2} = \frac{\sum X_{i2}}{4}$$

$$\mu_{..} = \frac{\sum \sum X_{ij}}{N_1+N_2+N_3} = \frac{\sum X_{ij}}{12}$$

The interpretation is the same as before. If the hypothesis is rejected, we assume that the  $\mu$ 's are independent. We can only say that they are not significantly different keeping in mind that, generally small samples cannot discover small differences. Lets now see how a researcher can reach conclusion about  $\mu$ 's from the comparism of variances of we have k sub-population whose sizes are  $N_1, N_2, N_3, \dots, N_k$  and whose means are  $\mu_{.1}, \mu_{.2}, \dots, \mu_{.k}$ , the pooled variance is found by

$$\delta^2_p = \frac{\sum (X_{i1} - \mu_{.1})^2 + \sum (X_{i2} - \mu_{.2})^2 + \dots + \sum (X_{ik} - \mu_{.k})^2}{N_1 + N_2 + N_3 + \dots + N_k} \quad \text{-----1}$$

Where  $\mu_{.j}$  are column means, the total variance then will be computed by the formula

$$\delta^2_t = \frac{\sum \sum (X_{ij} - \mu_{..})^2}{N_1 + N_2 + \dots + N_k} \quad \text{-----2}$$

Between  $\mu_{ij}$  values

$$\delta^2_{U_j} = \frac{\sum (U_{.j} - U_{..})^2}{k} \quad \text{-----3}$$

Example1

$X_{i1}$	$X_{i2}$	$X_{i3}$
1	1	1
5	5	5
3	3	3
7	7	7

Hypothesis

$$H_0: \mu_{.j} = \mu \quad \text{and} \quad \sigma^2_t = \sigma^2_p$$

$$H_1: \mu_{.1} \neq \mu_{.2} \neq \mu_{.3} \neq \mu \quad \text{and} \quad \sigma^2_t > \sigma^2_p$$

Recall that the sample mean is

$$\sigma^2_{\bar{X}} = \frac{\sigma^2}{n} \text{-----4}$$

From which the population variance is obtained by writing

$$\sigma^2 = n \sigma^2_{\bar{X}} \text{-----5}$$

And the value of the pooled variance is according to equation (1)

$$\sigma^2_p = \frac{20+20+20}{4+4+4} = 5$$

We can also have another example as below

Example (2)

$X_{i1}$	$X_{i2}$	$X_{i3}$
1	11	21
5	15	25
3	13	23
7	17	27

The total variance using equation (2) given  $\mu_{..} = 4$  is

$$\sigma_t^2 = \frac{\sum \sum (x_i - 4)^2}{12} = 5$$

Thus, it is obvious that these two variances are equal because in the problem at hand, the variance for each column are the same and the mean  $\mu_{.j}$  for each column are also the same. In this case,  $\sigma_{\mu.j}^2$  given in equation (3) is equal to zero. Therefore, comparison of these two variance (that is  $\sigma_t^2$  and  $\sigma_p^2$ ) will certainly tell us that the column means  $\mu_{.j}$  are the same.

The relation between the same variance is different in example (2). The variance for individual columns remain the same so that  $\sigma_1^2 = \sigma_2^2 = \sigma_3^2 = 5$ . Therefore, the value of pooled variance  $\sigma_p^2$  computed by equation (1) remains the same. However, the grand mean now is

$$\mu_{..} = \frac{\sum \sum X_{ij}}{N_1 + N_2 + N_3} = 14$$

Therefore, the total variance is

$$\begin{aligned} \sigma_t^2 &= 860 / 12 \\ &= 71.67 \end{aligned}$$

which is much greater than the pooled variance  $\sigma_p^2 = 5$  in example 1.



Again in example 2, the pooled variance is equal to the column variance.

However, the total variance is greater than the pooled variance and this is an indication that column means must be different. Thus, the knowledge and the comparison of these two variances are enough to tell us that the column means are different. We come to this conclusion without looking at the values of column means.

Now, assuming that the data in example 2 are not population but sample data. If so, then lieu to comparing  $\sigma_t^2$  with  $\sigma_p^2$ , we will compare (with the same result)

$$S_p^2 = \frac{\sum (X_{i1} - \bar{X}_{.1})^2 + \sum (X_{i2} - \bar{X}_{.2})^2 + \dots + \sum (X_{ik} - \bar{X}_{.k})^2}{n_1 + n_2 + n_3 + \dots + n_k - k} \quad \text{-----6}$$

with

$$S_b^2 = \frac{n \sum (\bar{X}_j - \bar{X}_{..})^2}{k-1} \quad \text{-----7}$$

where  $\bar{X}_{..}$  is the grand mean of sample data, computed in the same manner as  $\mu_{..}$  above. Equation 7 assumes that sample size  $n$  are equal. If they are not, the formula becomes

$$S_b^2 = \frac{\sum n_j (\bar{X}_j - \bar{X}_{..})^2}{k-1} \quad \text{-----8}$$

Without  $n$  and  $n_j$  in equation 7 and 8 respectively, the equation resembles the formula in equation 3. It actually becomes

$$\sigma_x^2 = \frac{\sum (\bar{X}_j - \bar{X}_{..})^2}{k-1} \quad \text{-----9}$$

The formula in equation 9 measures the dispersion of sample means  $\bar{X}_j$

and it is an estimate of  $\delta^2_x$  given in equation 4 above. If both  $\delta^2_x$  and the sample size  $n$  are known,  $\delta^2_x$  can be used as an estimate of population variance

$\delta^2 = n \delta^2_{\bar{x}}$ . If such an estimate has to be made from sample data, the variance of the population is estimated by equation 7, which is obtained by multiplying equation 9 by  $n$ . This variance, that is  $S_b^2$  in equation 7 is generally an estimate of  $\delta^2 + n\delta^2_{\mu_j}$ . However, it has already been noted that if the column means  $\mu_j$  are equal to zero, then  $\delta^2_{\mu_j}$  is equal to zero. In such a case,  $S_b^2$  in equation 7 (or equation 8) becomes an estimator of  $\delta^2$  in equation 6. If on the other hand; the column means are not the same, then  $S_b^2$  becomes an estimator of  $\delta^2 + n\delta^2_{\mu_j}$ . It will therefore be significantly greater than the pooled variance  $S_b^2$  as a result, we have the following fundamental rules in the analysis of variance: if  $S_b^2$  is significantly greater (and not significantly different) than  $S_p^2$ , column  $\mu$ 's are significantly different. A significantly greater  $S_p^2$  indicates that the variance  $\delta^2_{\mu_j}$  is significantly different from zero, which means that the  $\mu_j$  are significantly different. If on the other hand,  $S_b^2$  is not significantly greater than  $S_p^2$ , the column means are not significantly different.

It is useful to note that there is an additive relationship between the numerators of equation 6, 7 and 8. It can be shown as

$$\sum_i \sum_j (X_{ij} - X_{i.})^2 = \sum_i \sum_j (X_{ij} - X_{.j})^2 + n \sum_j (X_{.j} - X_{..})^2$$

In example 2, we have

$$860 = 4(200)$$

And in example 1

$$60 = 60 + 0$$

These numerators are called sum of squares. The first is "Total Sum of Squares" or  $SS_t$ , the second, "Sum of Square Within" or  $SS_w$ , and the third, the "Sum of Square Between" or  $SS_b$  and the Mean Square Between  $MS_b$ ;  $MS_b$  can be computed as

$$MS_b = SS_b / v_1$$

Where  $v_1$  is the number of degree of freedom. it is taken from the denominator in equation 7 that is  $v_1 = k - 1$ . The variance  $S_p^2$  is called mean square within or  $MS_w$  and can be computed as

$$MS_w = SS_w / v_2$$

The value of  $v_2$  comes from the denominator in equation 6. It is

$$v_2 = \sum n_j - k$$

Using the students score sheets however, computing these various sums of square becomes very tedious, cumbersome and time consuming. For a voluminous data, it becomes impossible manually. This is why I have decided to only explain it computational procedure with example.

The analysis of variance in this classification consists of comparing  $MS_b$

With  $SS_w$ . If  $MS_b$  is significantly greater than  $SS_w$ , the investigator concludes that mean  $\mu_j$  are significantly different. The comparison is done with the help of F-test. The computed F-value is obtained from

$$F^c = SS_b / MS_w \text{ and it is compared with } F_{(\alpha, v_1, v_2)} \text{ which can be found.}$$

These basic relationships are usually put together in tabular form

ANOVA TABLE

	SS	$v_i$	Ms	F
Between	$SS_b$	$V_1$	$MS_b$	$SS_b / MS_w$
Within	$SS_w$	$v_2$	$MS_w$	
Total	$SS_t$	$V_3$		

$$SS_b = n \sum \sum (X_j - X_{..})^2$$

$$SS_w = \sum \sum (X_{ij} - X_{.j})^2$$

$$SS_t = \sum \sum (X_j - X_{..})^2$$

$$V_3 = V_1 + V_2$$

## CHAPTER FOUR

### DESIGN AND IMPLEMENTATION

Since a program was not developed, but rather a standard application package known as Statistical Package for Social Sciences (SPSS) was used in modeling the students' achievement in their examination.

Syntax towards achieving this objective is as below:

#### 4.1 HOW TO INPUT DATA INTO THE DATA EDITOR

An operating system such as Windows creates a computing environment, or interface, within which the user accesses or runs applications such as SPSS and word processor. Windows creates a graphical interface, in which applications and other user-accessible software appear on the screen as small graphical symbols known as icons. When icon is activated and opened, commands can be given by using a device such as a mouse from lists of choices known as menus. When windows are first accessed, the screen shows an array in which several icons appear against a homogenous colored background known as the desktop. SPSS is opened by double clicking on the spss icon, a dialog box with the question:

“What would you like to do?” appears.

For more experienced users, this dialog box can be stopped from appearing automatically on opening in a non-network computer by clicking on the don't show this dialog box in the future. If the user is entering new data, click on the Type in data radio button and then on **OK** to bring the Variable View version of the SPSS Data Editor to the screen as shown below. If the **Data View** appears first, simply click on the tab labeled **Variable View** at the bottom left hand side of the window to open the Variable View

To name the variables English, Maths, Physics, type into the first cell of the first row in the column labeled Name of the Variable View and then press the down arrow key to move the highlighting down to the cell below. The type Maths, continue as explained earlier to type in Physics. The Variable View will appear as in below.

Once the user has completed the details of the variable in Variable View, then clicking on the Data View tab at the bottom of the window will bring on to the screen the Data View version.

In an SPSS data set, each row represents a single case, or subject and each column represents a variable. The cell to the left is highlighted and typing goes into the cell. Use the directional keys to navigate through the work sheet to enter values for each variable.

### 1.3 SELECTING A PROCEDURE FROM THE MENU

After entering the data into the data view, clicking **Analyze** initiates a pull down menu. Pointing the mouse on **REGRESSION** brings a sub-menu, then click on **Linear**.

#### SELECTING A VARIABLE FOR THE ANALYSIS

To select variables to regress after clicking on Linear, dialog box appears. The three variable names English, Maths and Physics will appear in the left-hand box. It is important to be clear about which variable is the dependent variable and which is the independent variable. Transfer these variable names into the appropriate boxes in the dialog box by clicking on the variable name and then on ►. This is a case of one independent variable which is an example of least square method. It is recommended to request additional descriptive Statistics and a residual analysis. To obtain descriptive statistics, click on the **Statistics** button on the dialog box to open the **Linear Regression: statistics** dialog box. Click on the **Descriptive** checkbox. A residual is the difference between the actual value of the dependent variable and its predicted value using the regression equation. Analysis of the residuals gives a measure of how good the prediction is and whether there are any cases, which are so discrepant that, they might be considered as outliers and so dropped

form the from the analysis. Click on the Casewise diagnostics check box to obtain a listing of any exceptionally large residuals. Click on **Continue** to return to the Linear Regression dialog box

Information about residuals is obtained by clicking on the Plot button in the dialog box to open the **Linear Regression: Plots** dialog box

Since systematic patterns between the predicted values and the residual can indicate possible violations of the assumption of linearity, it is recommended that a plot of the standardized residual (**\*ZRESID**) against the standardized predicted values **\*ZPRED** is also requested by transferring **\*ZRESID** into the Y: box and **\*ZPRED** into the box X: box.

Click on **Continue** to return to the Linear Regression dialog box

Click on **OK** to return the regression command.

## **OUTPUT FOR SIMPLE REGRESSION**

The various tables and charts in the output are listed in the left-hand pane of SPSS Viewer. Since the information about outliers in the Casewise Diagnostics table may affect whether the regression analysis is aborted and rerun after such outliers have been removed from the data set; it is recommended that this table is scrutinized first. It can be selected directly by moving the cursor to Casewise Diagnostics in the left-hand pane and clicking the left-hand mouse button.



**Variables Entered/Removed**

Model	Variables Entered	Variables Removed	Method
1	PHYSICS	.	Enter

a All requested variables entered.

b Dependent Variable: MATHS

**Model Summary**

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.510	.260	.242	7.67

a Predictors: (Constant), PHYSICS

b Dependent Variable: MATHS

**ANOVA**

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	867.438	1	867.438	14.764	.000
	Residual	2467.721	42	58.755		
	Total	3335.159	43			

a Predictors: (Constant), PHYSICS

b Dependent Variable: MATHS

### Coefficients

Model	Unstandardized Coefficients		Standardized Coefficients		t	Sig.
	B	Std. Error	Beta			
1 (Constant)	28.771	4.055			7.096	.000
PHYSICS	.369	.096	.510		3.842	.000

a Dependent Variable: MATHS

### Casewise Diagnostics

Case Number	Std. Residual	MATHS
4	3.446	74

a Dependent Variable: MATHS

### Residuals Statistics

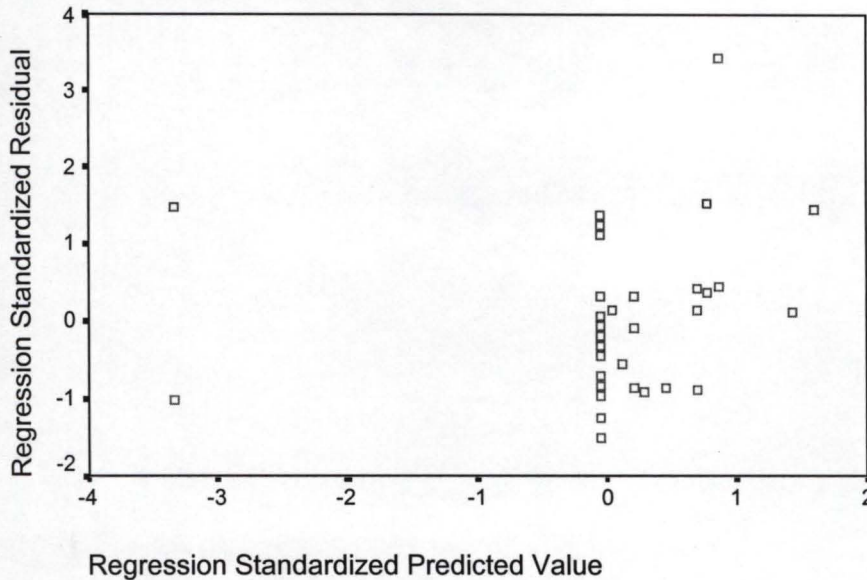
	Minimum	Maximum	Mean	Std. Deviation	N
Predicted Value	28.77	50.91	43.70	4.49	44
Residual	-11.53	26.41	3.23E-16	7.58	44
Std. Predicted Value	-3.325	1.604	.000	1.000	44
Std. Residual	-1.504	3.446	.000	.988	44

a Dependent Variable: MATHS

## Charts

### Scatterplot

Dependent Variable: MATHS



### Result Interpretation Of Simple Regression

The table of casewise diagnostic in the output shows that case 4 with a score of 74 for physics score is the only outlier with an absolute value greater than 3.

This outlier has to be eliminated and the regression analyses run again.

A more reliable regression analysis can be obtained by eliminating any outliers using the select case procedure described below

Choose Data,

Select cases...

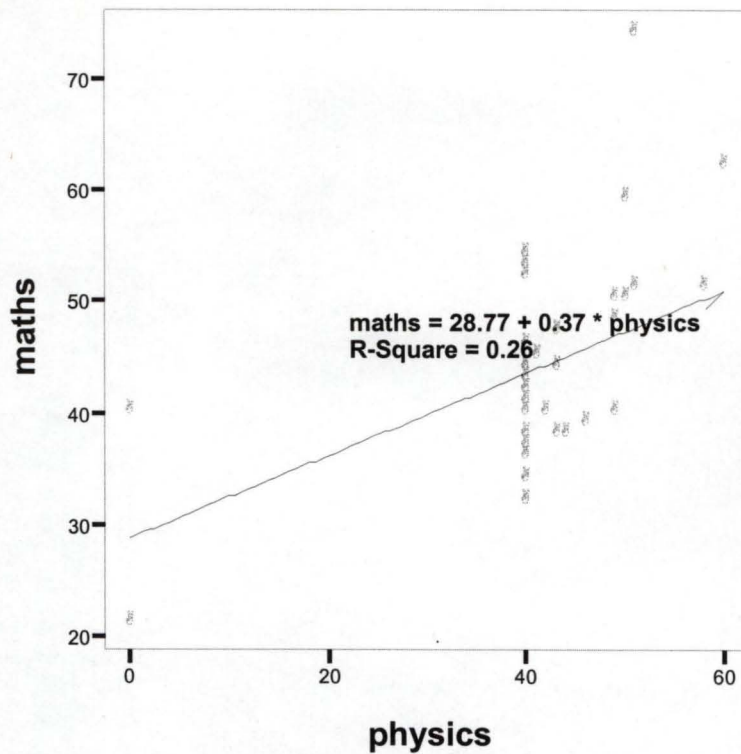
To open the select dialog box. Click on the **if** condition is satisfied radio button and define the condition as Case  $\neq$  4 (the symbol  $\neq$  means 'not equal to')

Clicks on continue and then **OK** to deselect this case to remove this case.

When the regression analysis is run again after deleting the original output to avoid confusion, there is no table of Casewise Diagnoses since no case have outlying residuals.

### **FITTING A SIMPLE REGRESSION LINE**

To fit a simple regression line, click on **Graph** to obtain a pull down menu, point to **Interactive** and a sub menu appears. Click on **Scatterplot** to obtain a dialog box as. Click on the **3D-coordinate** at the top right side of the dialog box, a pull sub menu appears, then click on **2D-corddinate**. Click and drag on the intended dependent and independent variables into Y-axis and X-axis respectively. Click **OK** to run the graph.



Linear Regression

## PROCEDURE FOR MULTIPLE REGRESSION

In the linear Regression dialog box, transfer the variable name Maths into the **Dependent Variable:** and Physics and chemistry into the **Independent Variables:** box by highlighting them and clicking on the appropriate button. To obtain descriptive statistics, click on the Statistics button as shown before to open the **Linear Regression: Statistics** dialog box. Click on the **Descriptive** checkbox. Click on **OK** to run the regression.]

## Output For Multiple Regressions

### Variables Entered/Removed

Model	Variables Entered	Variables Removed	Method
1	CHEMISTRY, PHYSICS	.	Enter

a All requested variables entered.

b Dependent Variable: MATHS

### Model Summary

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.725	.526	.503	6.21

a Predictors: (Constant), CHEMISTRY, PHYSICS

b Dependent Variable: MATHS

### ANOVA

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	1754.752	2	877.376	22.761	.000
	Residual	1580.407	41	38.547		
	Total	3335.159	43			

a Predictors: (Constant), CHEMISTRY, PHYSICS

b Dependent Variable: MATHS

## Coefficients

Model	Unstandardized Coefficients		Standardized Coefficients		t	Sig.
	B	Std. Error	Beta			
1 (Constant)	21.430	3.623			5.915	.000
PHYSICS	.140	.091	.193		1.532	.133
CHEMISTRY	.458	.096	.605		4.798	.000

a Dependent Variable: MATHS

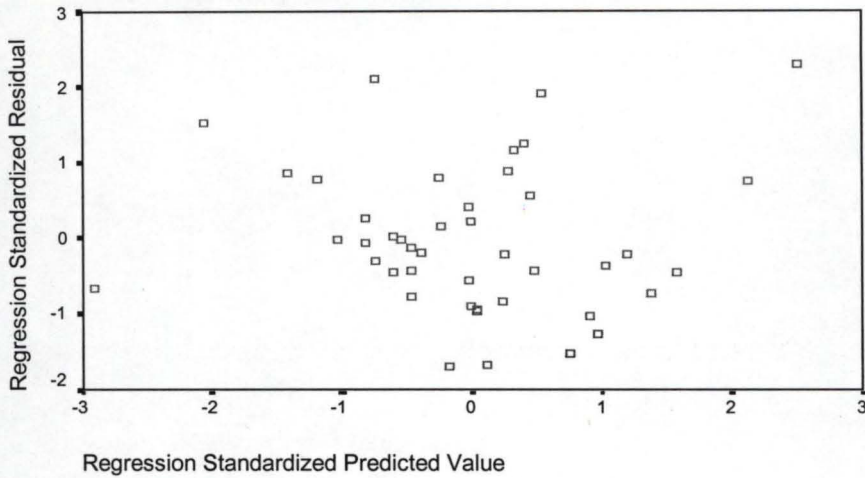
## Residuals Statistics

	Minimum	Maximum	Mean	Std. Deviation	N
Predicted Value	25.10	59.73	43.70	6.39	44
Residual	-10.61	14.27	4.52E-15	6.06	44
Std. Predicted Value	-2.913	2.508	.000	1.000	44
Std. Residual	-1.708	2.299	.000	.976	44

a Dependent Variable: MATHS

## Scatterplot

Dependent Variable: MATHS



## FITTING A MULTIPLE REGRESSION LINE MODEL

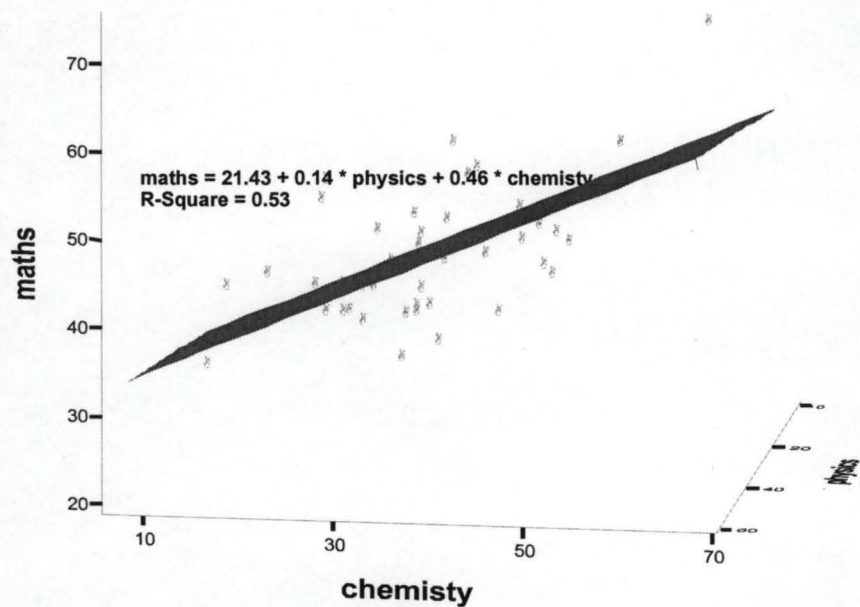
The syntax is as thus:

Click on **Graph** to obtain a pull down menu, point to **Interactive** and a sub menu appear. Click on **Scatterplot** to obtain a dialog box. Click and drag on the intended dependent and independent variables into the coordinates. Click **OK** to run the graph.



# MULTIPLE REGRESSION GRAPH

Linear Regression



## CHAPTER FIVE

### SUMMARY AND CONCLUSION

#### 5.1 Limitation Of The Study

Derivations of the three linear equation that resulted into the matrix formulation was not explained.

This project work is geared towards attempting modeling students scores for only two independent variables.

#### 5.2 Discussion

For the simple regression result, the first tables are the requested and correlation coefficients for the remaining 43 cases in output 2.

Output model summary shows the value for MULTIPLE **R** which in the case of just one independent variable has the same absolute value as the correlation coefficient  $r$  listed in output of descriptive statistics. The other statistics listed are **R square** (a positively biased estimate of the proportion of the variance of the dependent variable accounted for by regression). **Adjusted R Square** (corrects this bias and therefore has a lower value), and **Standard Error** (the standard deviation of the residuals).

Output ANOVA contains the regression **ANOVA**, which is intended to test whether there really is a linear relationship between the variables by forming an F-ratio of the mean square for regression to the residual mean square.

In this example, the value of F in the ANOVA is highly significant. It should be noted, however, that only an examination of the scatterplot of the variables could confirm that the relationship between two variables is genuinely linear.

Output coefficient is the kernel of the regression analysis, because it contains the regression equation. The value of the **regression coefficient** and constant are given in column **B** of the table. The equation is therefore,

$$\text{Predicted mathematics score} = 28.771 + 0.369 * (\text{physics score})$$

Thus person with Physics score of 40 would be predicted to score

$$28.771 + 0.369 * 40 = 43.531$$

Notice from the data that the person who did score 40 in Physics actually scored 42 in mathematics .the residual is, therefore,

$$43.531 - 40 = +3.531$$

Other statistics listed are Std. Error, which is the standard error of the regression coefficient B, Beta which is the beta weight showing the change in the dependent variable (expressed in standard deviating units) that would be produced by a positive increment of one standard deviation in the independent variable, t which is a t-test for testing the regression coefficient for significance and sig. which is the P-

value of  $t$  (here .00 means  $<0.005$ , that is,  $t$  is significant well beyond the 0.0 level for the variable Physics).

Output residual is a table of statistics relating to the residuals. The variables predicted value comprises the Unstandardized predicted values, the variable Residual is the set of Unstandardized residuals, the variable std. predicted value (identified as \*ZPRED in the plots dialog box in table model summary) contains the standardized predicted values (that is predicted values transformed to a scale with mean 0 and SD 1). The scatterplot shows no obvious pattern, thereby confirming that the assumptions of linearity and homogeneity of variance have been met. If the cloud of points were crescent-shaped or funnel-shaped, further screening of the data (or abandoned of the analysis) would be necessary.

For the multiple regression result, recall that when one independent variable physics was used to predict maths score, the value of  $R$  was 0.510. With  $R=0.725$ , we see that the answer to the question of whether adding more independent variables improves the predictive power of the regression equation is certainly 'yes', in fact the improvement is remarkable. The ANOVA output shows that the regression is still highly significant ( $p<0.01$ ). But what about the second question? Do both new variables contribute substantially to the predictive power of the regression equation?

From column B in the section headed unstandardized coefficients, we see that the multiple regression equation of maths score upon physics and chemistry score is

$$Y=21.430 + 0.14*\text{physics} +0.458*\text{chemistry score}$$

The beta weight tell us about the relative importance of the independent variables, because each gives the number of standard deviations change on the dependent variable that will be produced by a change of one standard deviation on the independent variable concerned.

Chemistry makes by far the highest contribution, because a change of one standard deviation on that variable produces a change of 0.608 standard deviation on maths score, whereas a change of one standard deviation in physics produces an increase of only 0.193 of a standard deviation in maths score. This equally signifies that the variable (predictor) with the largest beta weight also has the largest correlation with the dependent variable.

### 5.3 Summary

Computed results in chapter three can be used to form the models. These models formulated in turn are used to estimate students' scores. Using the linear equation in page 24

$$\text{Maths score} = 28.77 + 0.37 * \text{Physics score}$$

Substituting the various values of physics score into the linear equation gives an estimated corresponding Mathematics score for that particular student.

Example, for the student who scored 58 marks in Physics will have an estimated Mathematics score as

$$\begin{aligned}\text{Maths score} &= 28.77 + 0.37 * 58 \\ &= 28.77 + 21.46 \\ &= 50.23 \\ &\approx 50\end{aligned}$$

Substituting higher values of physics score results to a better mark or pass score equivalent to maths. Likewise, substituting a lower marks of physics results to a failure of mathematics.

Also from page 25

$$\text{Maths score} = 24.31 + 0.54 * x_i$$

Example, a student who scored marks in chemistry will have an estimated or predicted maths score s thus,

$$\begin{aligned}\text{Maths score} &= 24.31 + 0.54 * 53 \\ &= 24.31 + 28.62 \\ &= 52.93 \\ &\approx 53\end{aligned}$$

## 5.4 Conclusion

The various results computed manually to obtain the linear regressions are almost as those of the computed using SPSS. Only slight variations occur. The

estimate obtained using SPSS results give a better estimate to the corresponding dependent variable.

In the multiple regression models, better estimates are made of the dependent variable compared to those of simple regression.

In general, we are able to identify that Chemistry makes by far the highest contribution.

## **5.5 Recommendation**

Manual computation in anyway is not as easy as the computer-applied method. Once data has been imputed into the Data Editor of the SPSS, results are computed for analysis in matter of fractions after issuing the necessary command without mistakes.

Multiple regression is preferred to simple regression, this is because estimates made using the multiple regression are closer to the dependent variables compared to their corresponding linear regression.

## Reference

- Adamu S.O. (1985) : STATISTICS FOR BEGINNERS  
Second Edition  
Evans Brother (Nigeria Publisher)  
LTD
- Attely I.G Brokes C.J, And Loxton .S.M (1979): FUNDAMENTAL OF MATHEMATICS  
AND STATISTICS  
Published By Bro  
(Norwich) Ltd.
- Benjamin Fruchter (1987): FUNDAMENTAL STATISTICS IN  
Guilford J.P PSYCHOLOGY AND EDUCATION  
Third Edition  
Printed By B& Jo Enterprises Pte  
Ltd Spore
- Frank J.Kohout (1974) STATISTICS FOR SOCIAL  
SCIENTISTS  
Published By John Wiley & Sons, Inc
- Harry Franck & (1995): STATISTICS  
Steven C.Althoen Concept And Application  
Cambridge University Press
- Ronald J, Guilford (1986): INTRODUCTORY STATISTIC  
Thomas H.Wonnacott Third Edition
- Terry, Sinich (1990): STATISTICS BY EXAMPLE  
4<sup>th</sup> Edition  
Dellan Publishing Company,  
Sanfranscisco.
- Michel S.Kramer (1984): CLINICAL EPIDEMIOLOGY AND  
BIOSTATISTICS  
2<sup>nd</sup> Edition  
Springer-Verlag Publication



senior1

	econs	agric	engsquar	matsquar	physquar	engmaths	engphys
1	71	50	5184	2601	3364	3672	4176
2	69	67	4624	3844	3600	4216	4080
3	52	54	3364	2304	2401	2784	2842
4	62	75	5625	5476	2601	5550	3825
5	64	62	5041	1764	1600	2982	2840
6	69	46	3721	2209	1849	2867	2623
7	57	48	4489	1849	1600	2881	2680
8	51	46	3844	2916	1600	3348	2480
9	55	52	5041	1600	2401	2840	3479
10	35	47	3249	1521	2116	2223	2622
11	55	46	4096	2601	2601	3264	3264
12	68	60	4356	2500	2500	3300	3300
13	43	41	3136	2304	2401	2688	2744
14	59	45	2809	3481	2500	3127	2650
15	45	42	4761	2809	1600	3657	2760
16	54	37	2601	2500	2401	2550	2499
17	68	41	4624	1936	1600	2992	2720
18	51	49	4624	1936	1849	2992	2924
19	66	35	3249	2809	1600	3021	2280
20	52	38	3721	2116	1600	2806	2440
21	0	39	2704	1156	1600	1768	2080
22	63	48	3600	1444	1600	2280	2400
23	39	44	2601	2025	1681	2295	2091
24	50	53	4356	1444	1849	2508	2838
25	25	45	2304	1600	1600	1920	1920
26	45	44	3969	2704	1600	3276	2520
27	0	47	4624	1600	0	2720	0
28	84	42	3136	1444	1849	2128	2408
29	39	38	2401	1369	1600	1813	1960
30	35	31	4624	1444	1936	2584	2992
31	45	32	2401	1600	1600	1960	1960
32	44	46	2304	1296	1600	1728	1920
33	40	38	2500	1600	1600	2000	2000
34	48	35	2116	1600	1600	1840	1840
35	47	29	3136	1849	1600	2408	2240
36	0	34	2704	1024	1600	1664	2080
37	0	46	3249	1369	1600	2109	2280
38	28	36	2401	1600	1600	1960	1960
39	40	0	2025	1681	1600	1845	1800

## seniorb

	case	english	maths	agric	econs	biology	chemist
1	1	74	56	84	90	63	85
2	2	71	64	76	88	66	64
3	3	72	57	61	82	54	59
4	4	65	44	48	54	49	48
5	5	62	57	52	63	34	69
6	6	59	37	69	89	43	56
7	7	39	41	43	57	46	70
8	8	73	34	55	63	42	43
9	9	61	43	63	49	29	41
10	10	73	30	45	60	32	43
11	11	58	28	55	56	29	52
12	12	57	48	41	70	38	72
13	13	63	34	50	53	27	54
14	14	48	37	50	70	29	45
15	15	71	29	72	53	21	48
16	16	74	48	65	61	21	39
17	17	51	27	26	59	25	46
18	18	44	28	40	48	5	46
19	19	52	38	30	52	32	61
20	20	59	28	36	51	34	50
21	21	61	40	68	62	12	39
22	22	40	23	0	62	42	30
23	23	68	35	41	44	30	44
24	24	52	39	54	44	16	84
25	25	63	30	61	45	8	36
26	26	40	24	45	40	34	34
27	27	42	33	41	58	28	48
28	28	45	25	44	35	31	44
29	29	49	24	40	42	17	30
30	30	54	21	0	46	16	30
31	31	44	29	42	64	31	38
32	32	30	13	41	65	13	16
33	33	44	30	60	42	18	25
34	34	50	26	35	55	31	32
35	35	40	36	35	63	28	33
36	36	33	32	41	44	17	32
37	37	44	27	40	47	21	15
38	38	41	24	36	42	18	28
39	39	59	26	0	27	0	15

