# International Journal of Research Publication and Reviews

# A Case Study on Bank Loan by using MS Excel

*Pooja Kumawat [1], Mayuree Katara[2]*

**1,2Department of Computer Science & Engineering Vivekananda institute of technology Jaipur**
pooja.k@vitj.ac.in[1], **Katara.mayuree@vitj.ac.in[2]**

**ABSTRACT**

Data is being produced very quickly due to increase in information in everyday life. Vast volume of data generated from various organization. That is very difficult to analyze .Data created by different organizations such as healthcare database, banks, sales data, government databases etc. are example of large data processing. Analyzing, Processing and communicating such large amount of data are a challenge. It is very difficult task for analyzing and visualizing such large amount of data. Therefore some tools are required for analyzing such large amount of data. This paper focuses on data analysis and visualizing techniques. This paper shows a case study of bank loan.

*Index Terms— Analysis, EDA, visualization, MS Excel*

## 1. Introduction:

We are all always thinking about the future and what is expected to happen in the coming weeks, months, and even years, and in order to do so, a look into the past is required. Businesses must be able to see their development and the factors influencing their sales [1]. In this technological era of large-scale data, organizations must reconsider current techniques to better understand clients in order to achieve a competitive advantage in the market. Data is useless if it cannot be analyzed, comprehended, and applied in context [2].The primary goal here is to read and analyze the available datasets.

### I.A. WHAT IS DATA ANALYSIS?

Data analysis is a process that involves creating dataset, studying it, cleaning it and by removing null values, outliers and converting it to generate a useful result. By analyzing previous data helps you to make better decision in the future.

### I.B. Data visualization

Data visualization is a process which aims to communicate data effectively and clearly to the user through graphical representation. Effective and efficient data visualization is the key part of the discovery process. It is the intermediate between the human intuition and quantitative context of the data, thus an essential component of the scientific path from data into knowledge and understanding. It is a powerful new technology having a great potential to help researchers as well as companies for building revenue decision [3].

### I. C. Visualization Tools:

Data visualization tools organize and present information naturally. For simple visualization of data MS excel can be used many formats like bubble chart, line plots, histograms, tree map can be created. For plotting and graphing Matplotlib was used which has many packages for statistics, clustering and plotting. For creating graphics, processing was used widely. It allows uploading a data set and exploring it with pre build techniques. Spotfire is the best web client and analytical functionality tool. Qliktech may be the best visualization product having interactive drill-down capabilities. Tableau and Power bi has the excellent skill to interact with OLAP cubes. [18].

### I.D Data Visualization Techniques

Different data visualization techniques are used for data visualization.

- Box plots
- Histograms
- Heat maps

- Charts
- Tree maps
- Word Cloud/Network diagram

Box Plots

Box plot is a graph that shows how the values in your data are spread. It is a standard way presenting data based on this five attributes ("minimum", first quartile(Q1),median, third quartile(Q3) and "maximum).It can tell you about your outliers and their values. It can also show your data if it is symmetrical, how tightly your data is grouped and how your data is skewed.

**Summary of Box Plot**

| Minimum | Q1 -1.5*IQR |
|---|---|
| **First quartile (Q1/25th Percentile)"**: | The middle number between the smallest number (not the "minimum") and the median of the dataset |
| **Median (Q2/50th Percentile)"**: | The middle value of the dataset |
| **Third quartile (Q3/75th Percentile)"**: | The middle value between the median and the highest value (not the "maximum") of the dataset. |
| **Maximum"** | Q3 + 1.5*IQR |
| **interquartile range (IQR)** | 25th to the 75th percentile. |

**Histograms**

A histogram is a graphical representation of data that uses bars of various heights. Each bar in a histogram divides numbers into ranges. More data falls inside that range, as shown by taller bars. The form and distribution of continuous sample data are shown in a histogram.

It is a plot that enables you to identify and display the frequency distribution (shape) underneath a set of continuous data. This enables the examination of the data for outliers, skewness, and other factors such as outliers and underlying distribution (such as the normal distribution). It relates just one variable and accurately depicts the distribution of numerical data. Count the number of values that fall within each interval, which is achieved by dividing the whole range of values into bins or buckets.

**Line Chart**

Line plot is used to plot the connection or dependence of one variable on another. To plot the relationship between the two variables, we can simply use the plot function.

**Bar Charts**

Bar charts are used for comparing different categories or groups. Values of categories can be showed with vertical or horizontal bars. Each bar represents the value length and height.

**Pie Chart**

This statistical chart is circular and uses slices to show numerical proportion. Each slide's arc length is inversely correlated with the quantity it represents in this diagram. Typically used to compare the components of a whole, they work best with few components and when text and percentages are utilized to explain the contents. However, because the human eye struggles to estimate areas and contrast visual angles, they can be challenging to comprehend.

**Scatter Charts**

A scatter plot, which is a two-dimensional graphic that represents the joint variation of two data items, is another popular visualization tool. Each marker (such as a dot, square, or plus sign) denotes a particular observation. The value for each observation is shown by the marker location. A scatter plot matrix, which is a set of scatter plots depicting every possible pairing of the measurements that are assigned to the visualization, is created when you assign more than two measures. Scatter plots are used to analyze the relationship between the variable X and Y.

**Bubble Charts**

In this visualization technique data points are replaced with bubbles and dimension of data is represented in the size of bubbles.

**Timeline Charts**

Timeline charts show events in chronological sequence, in whichever unit of time the data was collected, such as week, month, year,      or quarter, for example the development of a project, an advertising campaign, or the acquisition process. On a timeframe, it displays the order in which past or future events occurred.

**Tree Map**

For hierarchical data, a tree map is a sophisticated, area-based data visualization that can be challenging to exactly analyze. Bar charts  and other less complex visualizations are frequently used.

## II. LITERATURE REVIEW

The term visualization is an evolving study area, where many researchers have contributed from the last few decades. Various authors have proposed different techniques and technologies to support data visualization. This section elaborates about how the flow of research has been carried out by the authors and researchers from reputed journals and conferences.

In [4] the author has proposed a Sensor: Network based approach for storing, sharing, visualizing and analyzing data from multiple devices and to interact with each other and with the end user through an open REST- based API. The author has visualized the geographical location of the data stream which when clicked pops up a tabbed window containing different associated information.

In [5] the author has proposed a virtual reality platform for scientific data visualization, a tool for multi-dimensional data visualization using which scientist can interact with the data and their colleagues in the same space. The author has mapped data parameters in different data points, shapes, size, colors, XYZ axis and many more. The author has used iViz a visualization tool which can be run as a standalone application or in a web browser. The author has discussed about a framework of financial time series delivery and visualization which can be used in viewing the historical price movement of a stock [6]. Specialized binary tree (SB- tree) is used for representing the financial time series. Time series data server, SB-tree server and web service contains is the three major components which are distributed on different machines. The system can reduce data volume and can capture the critical points.

In [7] the author has proposed a dashboard for displaying data used for communicating and finding trends in laboratory operation. System is based on .NET scripts, SQL repository. The author depicts that data is collected from the multiple sources like admin, internet and user portal and are stored in database using XML layer, Adobe flash, Action Script, etc. Data is being visualized which is used for laboratory and staff management.

In [8] the author has used a concept of visual web mining for analyzing the web data. A tool named WET is been used for visualization which provides a set of visual metaphor that represents the structure of the websites. The websites exploration tool is used for exploring the websites and for giving the feedback to the website owner for the betterment of the website.

In [9] the author has used a concept for analyzing data for examining the trend and evaluating the eco-environment impact of three gorges project. VC.NET and ArcIMs is the development platform for information system. ArcSDE and oracle 10g are used for management and use of spatial data. The author introduces method and processing and storing the data generated from cross-region, cross-department. Visualization helps in enhancing the data analysis and data mining.

In [10] the author has discussed the problem in compliance management which becomes an obstacle for decision making for effective and efficient monitoring. The person should be provided with compliance software which will help in getting high level information about overall compliance status and low level problem regarding possible problems. The author has designed a dashboard for watching the compliance which avoids the obstacle and decision can be made effectively. In [11] the author has introduced a tool named SECONDA which is used for analyzing both individual and grouped evolution of projects and develops belonging to a software ecosystem, Visualization is implemented in java using JFREECHART libraries. The author has used GNOME ecosystem for studying, under SECONDA. It offers a dashboard for fast visual analysis of local and global matrixes that can be extracted from information stored in the repositories.

In [12] the author has proposed a system for monitoring the user exercising progress and presenting exercise parameters in relation to prescribed targets. This system can be used for monitoring the intensity of the levels recommended by the patients care provider. It uses a miniature wireless 3-axis acceleration tied on the wrist of the patient that transmits acceleration data. The dashboard allows graphical visualization of exercise progress in real time. The author introduces a system where the huge amount of data generated from the collaborative software development tool during the lifecycle of a project can be used to analyze the performance of the individual member, or a team or manager.

[13] They can analyze from different perspectives across different dimensions and visualized in different ways.

In [14] the author has proposed a dashboard which is an integration, validation and visualizing tool for natural language processing. The system helps the system integration team to integrate and validate the system; developers to profile each module and researchers to evaluate and compare the module with the earlier versions. It also supports execution of modules on heterogeneous platform with an easy to use graphical interface developed using eclipse RCP.

## III. Proposed Methodology

The primary goals of data visualizations are to analyze the data and communicate it to the user. Visualization's primary objective is to use graphics to connect information in a clear and effective way. We are putting forth a system that will examine and display sales data. The data will be graphed using various parameters to reflect various viewpoints. To find trends for future forecasts, data mining will be used. For analysis and visualizations, a data set from kaggle is used. We have taken the data set of bank loan application



Fig. III. a Data analysis process

III. A In this analysis following questions are solved.

- Present the overall approach of the analysis. Mention the problem statement and the analysis approach briefly

- Identify the missing data and use appropriate method to deal with it. (Remove columns/or replace it with an appropriate value)

- Identify if there are outliers in the dataset. Also, mention why you think it is an outlier. Again, remember that for this exercise, it is not necessary to remove any data points.

- Identify if there is data imbalance in the data. Find the ratio of data imbalance.

- Explain the results of univariate**,** segmented univariate**,** bivariate analysis, etc**.** in business terms.

- Find the top 10 correlation for the Client with payment difficulties and all other cases (Target variable). Note that you have to find the top correlation by segmenting the data frame w.r.t to the target variable and then find the top correlation for each of the segmented data and find if any insight is there. Say, there are 5+1(target) variables in a dataset: Var1, Var2, Var3, Var4, Var5, and Target. And if you have to find top 3 correlation, it can be: Var1 & Var2, Var2 & Var3, Var1 & Var3. Target variable will not feature in this correlation as it is a categorical variable and not a continuous variable which is increasing or decreasing.

- Include visualizations and summarize the most important results in the presentation.

## IV. Tech Stack Used

In order to execute this project Microsoft Excel is used. Once the data record is downloaded. After that various formulas and graphs are used to represent the data.

## V.  Results

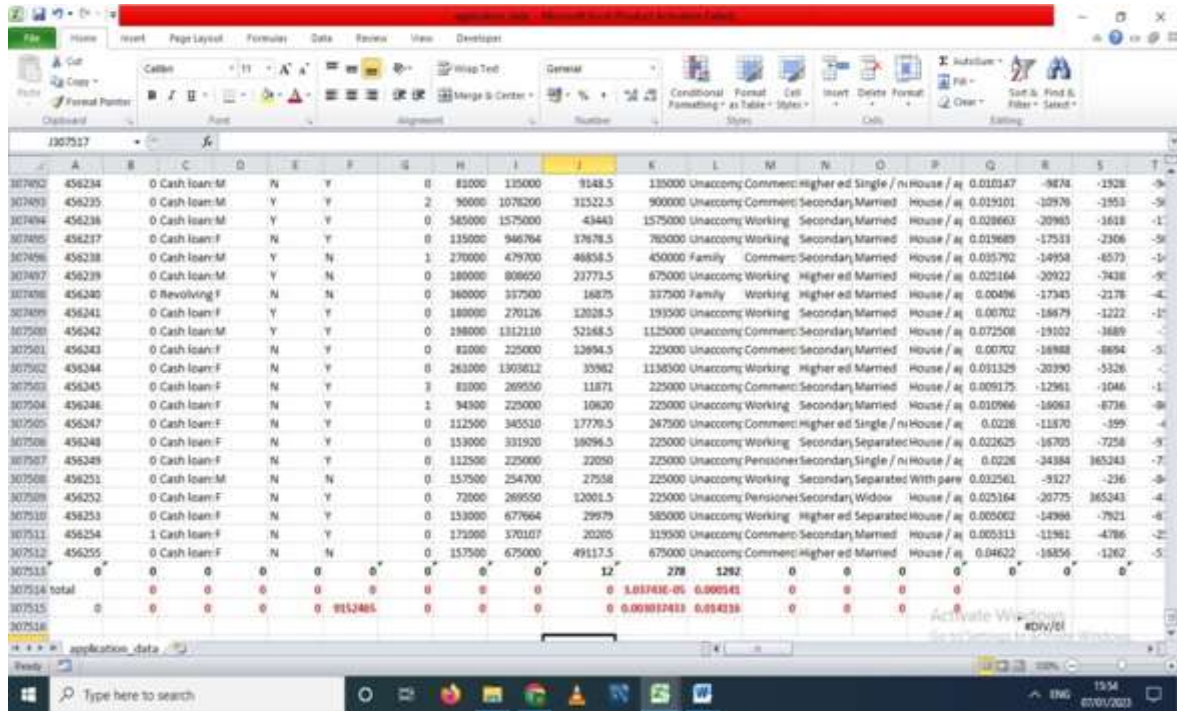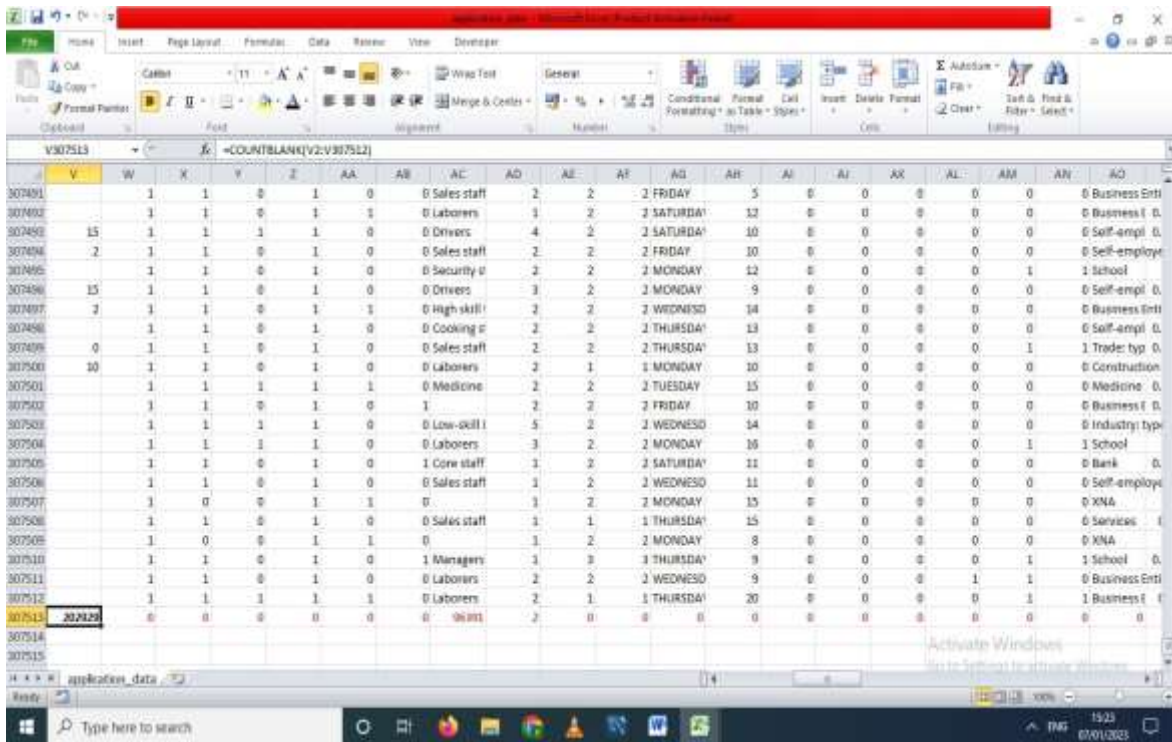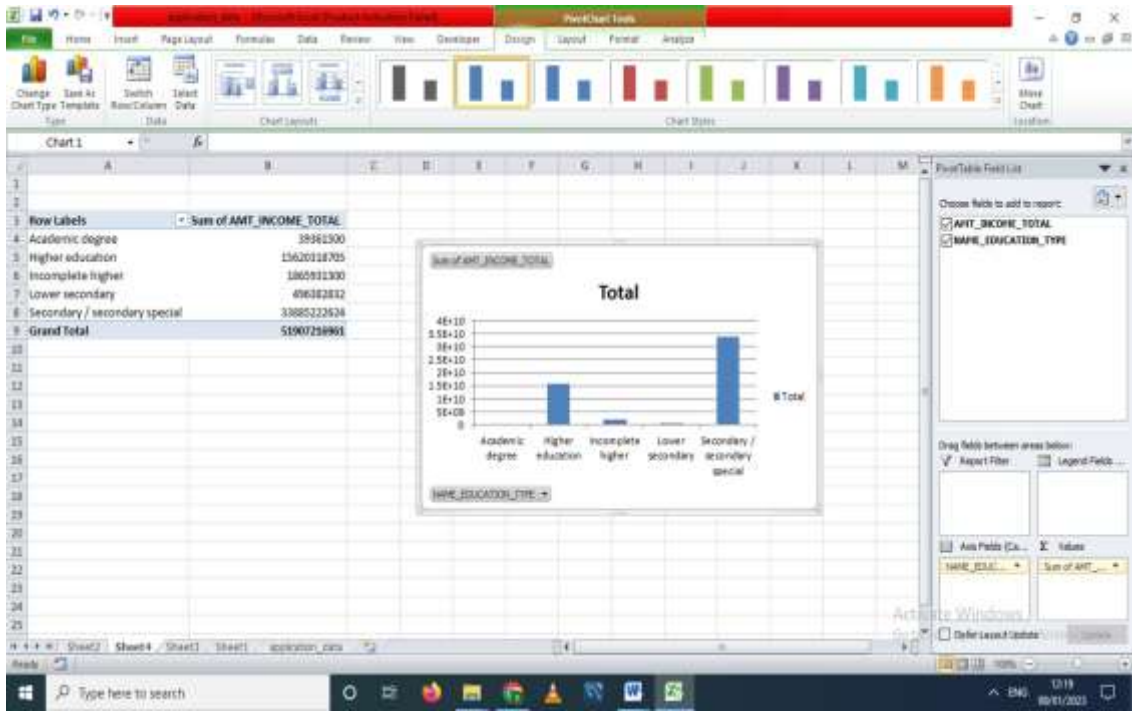We can count the no of blank cell by using

**Count blank (A3:A307512)**

**Fig. V.1 Counting no. of blank cells.**



**Total blank cell: sum(=K307513/F307515)*100**

Fig. V.II Counting the no. of blank cells.

**V. B Find the quartile using quartile function.**

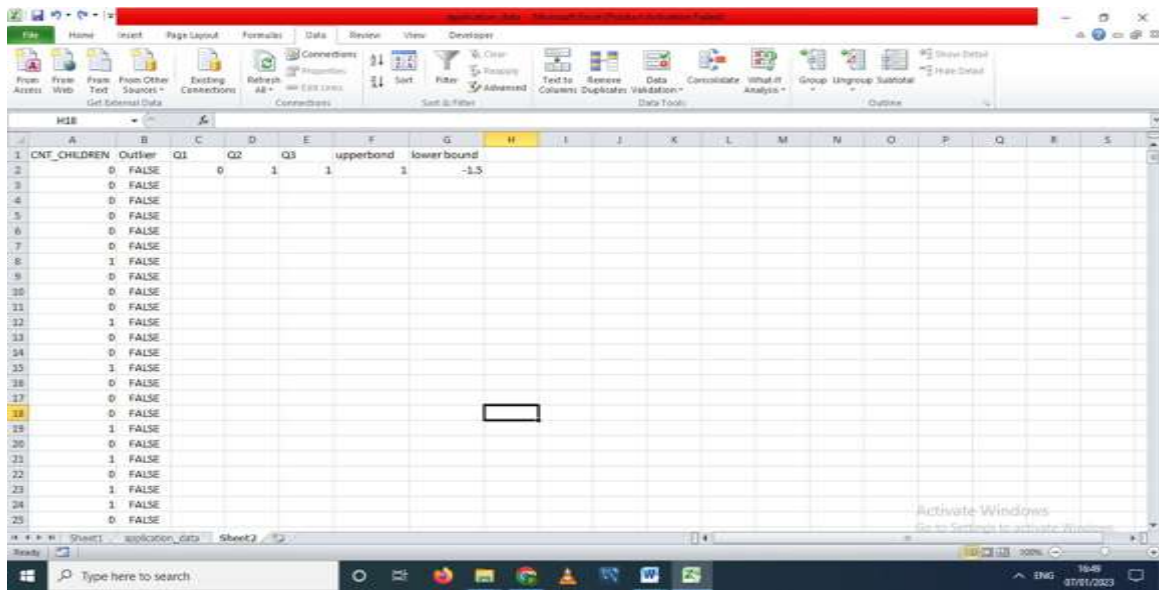In this fig we have calculated the lower bound and upper bound by using quartile function.



Fig V.II Finding Outlier.

**V. C Data Imbalance**



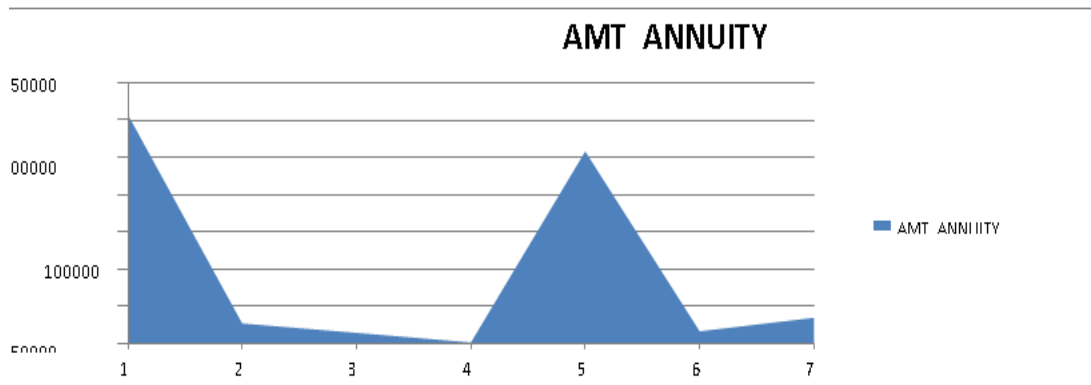Fig. V.1 Data Imbalance

V.D  Univarte Analysis
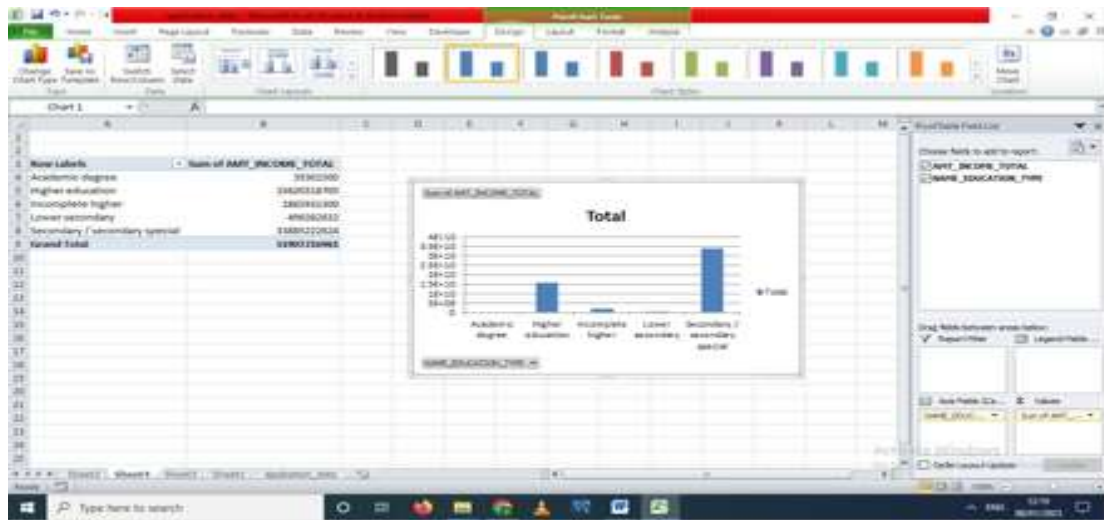


Fig V.IV Univarte analysis

V.D Bivarte Analysis:



Fig. V.1 Bivarte Analysis

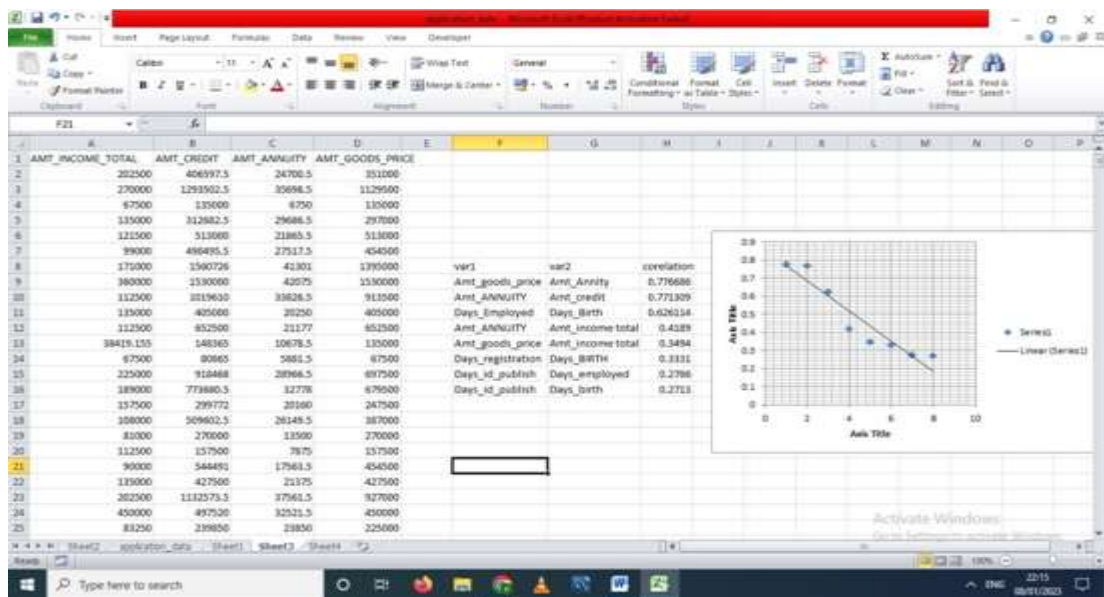V.E  Corelation Between two Variable. = **=CORREL(A2:A47532, B2:B725)**



Fig. V.1 Correlation between two variable

## VI. Conclusions

We have examined various strategies, techniques, and tools in this study, all of which have drawbacks of their own. We have discussed numerous papers that gave us a general understanding of the kind of system needed in the modern era for the analysis and visualizations of data so that business owners and investors may make wise decisions and produce profit. We have suggested a technique where data is processed before being imported and saved in a database. With varied factors and dimensions, this data is visualized. with the help of which the we can easily filter the data . We have suggested a method for importing and storing data.

## VII. Future Scope

In our upcoming work, we'll employ a variety of cutting-edge   visualization approaches and compile all the graphs and charts onto a single dashboard to aid users in quickly making decisions and producing income.

### REFERENCES

[1] Jiawei Han, Micheline Kamber and Jian Pei, "Data Mining Concepts and Techniques",Third edition, MK Publications, 2009.

[2] M. Tennekes and E. de Jonge, "Top-down Data Analysis with Treemaps," in Proceedings of the International Conference on Information Visualization Theory and Applications (IVAPP'11), pp. 236–241, March 2011.

[3] P. Hoek, "Parallel Arc Diagrams: Visualizing Temporal Interactions," Journal of Social Structure, vol. 12, 2011.

[4] Vipul Gupta, Arshan Porsohi and Poornaprajna Udupi, "Sensor Network: An Open Data Exchange for the Web of Things," in Proceedings of 8th IEEE Conference on Pervasive Computing and Communication Workshop (PERCOM), 2010.

[5] Ciro Donalek, S.G. Djorgovski, Alex cioc, Anwell Wang, Jerry Zhang, Elizabeth Lawler, Stacy Yeh, Ashish Mahbal, Matthew Graham, Andrew, Drake Scott Davidoff and Jeffrey S. Norris, "Immersive and Collaboration Data visualization using Virtual Reality Platforms," in Proceedings of IEEE International Conference on Big Data, 2014.

[6] Tak-chung Fu, Fu-lai Chung and Chun-fai Lam, "Adaptive Data Delivery Framework for Financial Time Series Visualization," in Proceedings of IEEE International Conference on Mobile Business (ICMB), pp. 267-273, 2005.

[7] Eric Martin and Vincenzo Di Bernardo, "Enterprise Dashboard Tools for Management of Share-use University Laboratory," in Proceedings of University Conveinment Industry Micro (UCIM), 2008.

[8] Victor Pascual-cid, "An Information Visualization System for the Understanding of Web Data," in Proceedings of IEEE Symposium on Information Visualization (INFOVIS), 2008.

[9] Liang Zhu, Bing-Fang, Yue-min Zhou, Xin-hui ma and Leidong yang, "Researches on Eco-environment Data Visualization for Three Gorgeous Project," in Proceedings of IEEE International Symposium on Information Engineering Electronic Commerce (IEEC), 2009.

[10] Thorben sander, Malthias kehlenbeeck and Michael H. Breithner, "Visualization of Automated Compliance Monitoring and Repairing," in Proceedings of IEEE Database and Expert System Applications (DEXA), 2012.

[11] Javier Perez, Romuald Deshayes, Methieu Goeminne and Tom Mens, "SECONDA: Software Ecosystem Analysis Dashboard," in International Conference on Software Maintenance and Reengineering (CSMR), 2010.

[12] Cheol Jeong and Joseph Finkelstein, "Computer-Assisted Upper Extremity Training Using Interactive Biking Exercise (iBikE) Platform," in Proceedings of IEEE Conference on Engineering in Medicine and Biological Society (EMBC), 2012.

[13] Fleni Stroulia, Isaac Matichuk, Fabio Rocha, Ken Baver,  "Interactive Exploration of collaboration software     development    data," in IEEE International Conference on Software Maintenance (ICSM), 2013.

[14] Pavan Kumar, Rashid Ahmad, B.D. Chaudary and Mukul K. Sinha, "Enriched dashboard:-An Integration and Visualization Tool for Distributed NLP System on Heterogeneous Platform," in International Conference on Computer Science and its Applications (ICCSA), 2013.

[15] Ma Xin-hui, Wu Bing-fang, Zhu Liang, etc. "Study on Methods of Hydro-environment Data Visualization for Three Gorges Project" [J], Remote sensing informatics, 2007, pp. 28-31.

[16] M. Lungu, M. Lanza, T. GˆÕrba, and R. Robbes, "The small project observatory: Visualizing software ecosystems," Science of Computer Programming, vol. 75, no. 4, pp. 264 – 275, 2010.

[17]R. Chang, M. Ghoniem, R. Kosara, W. Ribarsky, J. Yang, E. Suma, C. Ziemkiewicz, D. Kern, and A. Sudjianto, Wire-vis: Visualization of categorical, time-varying data from financial transactions. In IEEE Symposium on Visual Analytics Science and Technology, pages 155–162, 2007.

[18] "The 37 Best Tools for Visualization," http://www.creativebloq.com/design-tools/data-visualization712402.Harjeet Kaur, Varsha Sahni, and Dr. Manju Bala, "A survey of reactive, proactive, and hybrid routing protocols in MANET: a review", International Journal of Computer Science and Information Technologies, vol. 4, no. 3, pp. 498-500, 2013