Prediction of bigmart sales using machine learning algorihms

Submitted in partial fulfillment of the requirements for the award of Bachelor of Engineering degree in Computer Science and Engineering

by

Mr.MallipeddiVineethGuptha(38110291) Mr.Gande Abhilash(38110006)



DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

SCHOOL OF COMPUTING

SATHYABAMA

INSTITUTE OF SCIENCE AND TECHNOLOGY (DEEMED TO BE UNIVERSITY)

Accredited with Grade "A" by NAAC

JEPPIAAR NAGAR, RAJIV GANDHI SALAI, CHENNAI - 600 119

MARCH - 2022



DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

BONAFIDE CERTIFICATE

This is to certify that this Project Report is the bonafide work of Mr.MallipeddiVineethGuptha(38110291),Mr.GANDEABHILASH(38110006) who carried out the project entitled "PREDICTIVE ANALYSIS OF BIGMART SALES USING MACHINE LEARNING ALGORITHMS"

under my supervision from October 2021 to May 2022

InternalGuide (Mrs.N.S.Usha,B.E,M.E,(Ph.D) CSE)

Head of the Department

Submitted for Viva voce Examination held on

External

DECLARATION

I Mr.Mallipeddi Vineeth Guptha, Mr.Gande Abhilash hereby declare that the Project Report entitled **PREDICTIVE ANALYSIS OF BIGMART SALES USING MACHINE LEARNING ALGORITHMS** done by me under the guidance of Mrs.N.S.Usha,B.E,M.E,(Ph.D) CSE (Internal) is submitted in partial fulfillment of the requirements for the award of Bachelor of Engineering / Technology degree In COMPUTER SCIENCE AND TECHNOLOGY.

DATE:

PLACE:

SIGNATURE OF THE CANDIDATE

ACKNOWLEDGEMENT

I am pleased to acknowledge my sincere thanks to **Board of Management** of **SATHYABAMA** for their kind encouragement in doing this project and for completing it successfully. I am grateful to them.

I convey my thanks to **Dr. T.Sasikala M.E., Ph.D**, **Dean**, School of Computing **Dr. L. Lakshmanan M.E., Ph.D.**, and **Dr.S.Vigneshwari M.E., Ph.D. Heads** of the Department of Computer Science and Engineering for providing me necessary support and details at the right time during the progressive reviews.

I would like to express my sincere and deep sense of gratitude to my Project Guide **Mrs.N.S.Usha,B.E,M.E,(Ph.D) CSE** for his valuable guidance, suggestions and constant encouragement paved way for the successful completion of my project work.

I wish to express my thanks to all Teaching and Non-teaching staff members of the **Department of Computer Science and Engineering** who were helpful in many ways for the completion of the project.

ABSTRACT

Nowadays shopping malls and Big Marts keep the track of their sales data of each and every individual item for predicting future demand of the customer and update the inventory management as well. These data stores basically contain a large number of customer data and individual item attributes in a data warehouse. Further, anomalies and frequent patterns are detected by mining the data store from the data warehouse. The resultant data can be used for predicting future sales volume with the help of different machine learning techniques for the retailers like Big Mart. In this paper, we propose a predictive model using XG boost Regressor technique for predicting the sales of a company like Big Mart and found that the model produces better performance as compared to existing models.

CONTENT

Chapter		List of content	Page no.
		Abstract	1
li		List of acronyms	
lii		List of figures	
1		Introduction	
2		Literature Survey	13
3		System Requirements	17
	3.1	Hardware Requirements	
	3.2	Software Requirements	
	3.3	Language Specification	
	3.4	History of python	
4		System Analysis	33
	4.1	Purpose	
	4.2	Scope	
	4.3	Existing System	
	4.4	Proposed System	
5		System Design	36
	5.1	Input Design	
	5.2	Output Design	
	5.3	Data Flow Diagram	
6		Module Implementations	42
	6.1	Modules	
	6.1.1	Data Collection	
	6.1.2	Data Set	
	6.1.3	Data Preparation	
	6.1.4	Model Selection	
	6.1.5	Analyze and Prediction	
	6.1.6	Accuracy on Test Set	
	6.1.7	Saving the Training Model	
7		System Implementation	47
8		System Testing	48
	8.1	Unit Testing	
	8.2	Integration Testing	

2

	8.3 8.4 8.5 8.6	Functional Testing System Test White Box Testing Black Box Testing	
9	8.7	Results And Discussion	53
10		Conclusion	56
11		References	57

LIST OF FIGURES:

FIGURE NO.	FIGURE NAME	PAGE NO.
7.1	System Architecture	47
9.2	Data collection from kaggle	54
9.4	prediction of sales	55
9.6	performance values of propose	ed 56
	model	

CHAPTER 1

INTRODUCTION

Big Mart is a big supermarket chain, with stores all around the country and its current board set out a challenge to all Data Scientist out there to help them create a model that can predict the sales, per product, for each store to give accurate results. Big Mart has collected sales data from Kaggle, for various products across different stores in different cities. With this information the corporation hopes we can identify the products and stores which play a key role in their sales and use that information to take the correct measures to ensure success of their business.

1.1 PROPOSED ALGORITHMS

Decision Trees are a type of Supervised Machine Learning (that is you explain what the input is and what the corresponding output is in the training data) where the data is continuously split according to a certain parameter. The tree can be explained by two entities, namely decision nodes and leaves. The leaves are the decisions or the final outcomes. And the decision nodes are where the data is split.



An example of a decision tree can be explained using above binary tree. Let's say you want to predict whether a person is fit given their information like age, eating habit, and physical activity, etc. The decision nodes here are questions like 'What's the age?', 'Does he exercise?', 'Does he eat a lot of pizzas'? And the leaves, which are outcomes like either 'fit', or 'unfit'. In this case this was a binary classification problem (a yes no type problem).

There are two main types of Decision Trees:

Classification trees (Yes/No types)

What we've seen above is an example of classification tree, where the outcome was a variable like 'fit' or 'unfit'. Here the decision variable is Categorical.

Regression trees (Continuous data types)

Here the decision or the outcome variable is Continuous, e.g. a number like 123.

Working

Now that we know what a Decision Tree is, we'll see how it works internally. There are many algorithms out there which construct Decision Trees, but one of the best is called as ID3 Algorithm. ID3 Stands for Iterative Dichotomiser 3.

Before discussing the ID3 algorithm, we'll go through few definitions.

- Decision Tree is a Supervised learning technique that can be used for both classification and Regression problems, but mostly it is preferred for solving Classification problems. It is a tree-structured classifier, where internal nodes represent the features of a dataset, branches represent the decision rules and each leaf node represents the outcome.
- In a Decision tree, there are two nodes, which are the Decision Node and Leaf Node. Decision nodes are used to make any decision and have multiple branches, whereas Leaf nodes are the output of those decisions and do not contain any further branches.

- The decisions or the test are performed on the basis of features of the given dataset.
- It is a graphical representation for getting all the possible solutions to a problem/decision based on given conditions.
- It is called a decision tree because, similar to a tree, it starts with the root node, which expands on further branches and constructs a tree-like structure.
- In order to build a tree, we use the CART algorithm, which stands for Classification and Regression Tree algorithm.
- A decision tree simply asks a question, and based on the answer (Yes/No), it further split the tree into subtrees.
- Below diagram explains the general structure of a decision tree:

Note: A decision tree can contain categorical data (YES/NO) as well as numeric data.



Why use Decision Trees?

There are various algorithms in Machine learning, so choosing the best algorithm for the given dataset and problem is the main point to remember while creating a machine learning model. Below are the two reasons for using the Decision tree:

- Decision Trees usually mimic human thinking ability while making a decision, so it is easy to understand.
- The logic behind the decision tree can be easily understood because it shows a tree-like structure.

Decision Tree Terminologies

□ **Root Node:** Root node is from where the decision tree starts. It represents the entire dataset, which further gets divided into two or more homogeneous sets.

□ **Leaf Node:** Leaf nodes are the final output node, and the tree cannot be segregated further after getting a leaf node.

□ **Splitting:** Splitting is the process of dividing the decision node/root node into subnodes according to the given conditions.

□ **Branch/Sub Tree:** A tree formed by splitting the tree.

□ **Pruning:** Pruning is the process of removing the unwanted branches from the tree.

□ **Parent/Child node:** The root node of the tree is called the parent node, and other nodes are called the child nodes.

How does the Decision Tree algorithm Work?

In a decision tree, for predicting the class of the given dataset, the algorithm starts from the root node of the tree. This algorithm compares the values of root attribute with the record (real dataset) attribute and, based on the comparison, follows the branch and jumps to the next node.

For the next node, the algorithm again compares the attribute value with the other subnodes and move further. It continues the process until it reaches the leaf node of the tree. The complete process can be better understood using the below algorithm: Competitive questions on Structures in Hindi

Keep Watching

- **Step-1:** Begin the tree with the root node, says S, which contains the complete dataset.
- Step-2: Find the best attribute in the dataset using Attribute Selection Measure (ASM).
- **Step-3:** Divide the S into subsets that contains possible values for the best attributes.
- **Step-4:** Generate the decision tree node, which contains the best attribute.
- Step-5: Recursively make new decision trees using the subsets of the dataset created in step -3. Continue this process until a stage is reached where you cannot further classify the nodes and called the final node as a leaf node.

Example: Suppose there is a candidate who has a job offer and wants to decide whether he should accept the offer or Not. So, to solve this problem, the decision tree starts with the root node (Salary attribute by ASM). The root node splits further into the next decision node (distance from the office) and one leaf node based on the corresponding labels. The next decision node further gets split into one decision node (Cab facility) and one leaf node. Finally, the decision node splits into two leaf nodes (Accepted offers and Declined offer). Consider the below diagram:



Attribute Selection Measures

While implementing a Decision tree, the main issue arises that how to select the best attribute for the root node and for sub-nodes. So, to solve such problems there is a technique which is called as **Attribute selection measure or ASM.** By this measurement, we can easily select the best attribute for the nodes of the tree. There are two popular techniques for ASM, which are:

- Information Gain
- Gini Index
- 1. Information Gain:
 - Information gain is the measurement of changes in entropy after the segmentation of a dataset based on an attribute.
 - It calculates how much information a feature provides us about a class.

- According to the value of information gain, we split the node and build the decision tree.
- A decision tree algorithm always tries to maximize the value of information gain, and a node/attribute having the highest information gain is split first. It can be calculated using the below formula:
- 1. Information Gain= Entropy(S)- [(Weighted Avg) *Entropy(each feature)

Entropy: Entropy is a metric to measure the impurity in a given attribute. It specifies randomness in data. Entropy can be calculated as:

 $Entropy(s) = -P(yes)\log 2 P(yes) - P(no) \log 2 P(no)$

Where,

- S= Total number of samples
- P(yes)= probability of yes
- P(no)= probability of no

2. Gini Index:

- Gini index is a measure of impurity or purity used while creating a decision tree in the CART(Classification and Regression Tree) algorithm.
- An attribute with the low Gini index should be preferred as compared to the high Gini index.
- It only creates binary splits, and the CART algorithm uses the Gini index to create binary splits.
- Gini index can be calculated using the below formula:

Gini Index= 1- $\sum_{j} P_{j}^{2}$

Pruning: Getting an Optimal Decision tree

Pruning is a process of deleting the unnecessary nodes from a tree in order to get the optimal decision tree.

A too-large tree increases the risk of overfitting, and a small tree may not capture all the important features of the dataset. Therefore, a technique that decreases the size of the learning tree without reducing accuracy is known as Pruning. There are mainly two types of tree **pruning** technology used:

• Cost Complexity Pruning

• Reduced Error Pruning.

Advantages of the Decision Tree

- It is simple to understand as it follows the same process which a human follow while making any decision in real-life.
- It can be very useful for solving decision-related problems.
- It helps to think about all the possible outcomes for a problem.
- There is less requirement of data cleaning compared to other algorithms.

Disadvantages of the Decision Tree

- The decision tree contains lots of layers, which makes it complex.
- It may have an overfitting issue, which can be resolved using the Random Forest algorithm.
- For more class labels, the computational complexity of the decision tree may increase.

0

CHAPTER 2 LITERATURE SURVEY

1) A comparative study of linear and nonlinear models for aggregate retails sales forecasting

AUTHORS: Ching Wu Chu and Guoqiang Peter Zhang

The purpose of this paper is to compare the accuracy of various linear and nonlinear models for forecasting aggregate retail sales. Because of the strong seasonal fluctuations observed in the retail sales, several traditional seasonal forecasting methods such as the time series approach and the regression approach with seasonal dummy variables and trigonometric functions are employed. The nonlinear versions of these methods are implemented via neural networks that are generalized nonlinear functional approximators. Issues of seasonal time series modeling such as deseasonalization are also investigated. Using multiple cross-validation samples, we find that the nonlinear models are able to outperform their linear counterparts in out-of-sample forecasting, and prior seasonal adjustment of the data can significantly improve forecasting performance of the neural network model. The overall best model is the neural network built on deseasonalized time series data. While seasonal dummy variables can be useful in developing effective regression models for predicting retail sales, the performance of dummy regression models may not be robust. Furthermore, trigonometric models are not useful in aggregate retail sales forecasting.

2) Sustainable development and management in consumer electronics using soft computation

AUTHORS: Wang, Haoxiang

Combination of Green supply chain management, Green product deletion decision and green cradle-to-cradle performance evaluation with Adaptive-Neuro-Fuzzy Inference System (ANFIS) to create a green system. Several factors like design process, client

specification, computational intelligence and soft computing are analysed and emphasis is given on solving problems of real domain. In this paper, the consumer electronics and smart systems that produce nonlinear outputs are considered. ANFIS is used for handling these nonlinear outputs and offer sustainable development and management. This system offers decision making considering multiple objectives and optimizing multiple outputs. The system also provides efficient control performance and faster data transfer.

3) Data Mining based Prediction of Demand in Indian Market for Refurbished Electronics

AUTHORS: Suma, V., and Shavige Malleshwara Hills

There has been an increasing demand in the e-commerce market for refurbished products across India during the last decade. Despite these demands, there has been very little research done in this domain. The real-world business environment, market factors and varying customer behavior of the online market are often ignored in the conventional statistical models evaluated by existing research work. In this paper, we do an extensive analysis of the Indian e-commerce market using data-mining approach for prediction of demand of refurbished electronics. The impact of the real-world factors on the demand and the variables are also analyzed. Real-world datasets from three random e-commerce websites are considered for analysis. Data accumulation, processing and validation is carried out by means of efficient algorithms. Based on the results of this analysis, it is evident that highly accurate prediction can be made with the proposed approach despite the impacts of varying customer behavior and market factors. The results of analysis are represented graphically and can be used for further analysis of the market and launch of new products.

4) Forecasting Monthly Sales Retail Time Series: A Case Study

AUTHORS: Giuseppe Nunnari, Valeria Nunnari

This paper presents a case study concerning the forecasting of monthly retail time series recorded by the US Census Bureau from 1992 to 2016. The modeling problem is tackled in two steps. First, original time series are de-trended by using a moving windows averaging approach. Subsequently, the residual time series are modeled by Non-linear Auto-Regressive (NAR) models, by using both Neuro-Fuzzy and Feed-Forward Neural Networks approaches. The goodness of the forecasting models, is objectively assessed by calculating the bias, the mae and the rmse errors. Finally, the model skill index is calculated considering the traditional persistent model as reference. Results show that there is a convenience in using the proposed approaches, compared to the reference one.

5) Multiple Linear Regression Analysis of the Overlay Accuracy Model Zone

AUTHORS: Zone-Ching Lin, Wen-Jang Wu

The multiple linear regression method was used to analyze the overlay accuracy model and study the feasibility of using linear methods to solve parameters of nonlinear overlay equations. The methods of analysis include changing the number of sample points to derive the least sample number required for solving the accurate estimated parameter values. Besides, different high-order lens distortion parameters were ignored, and only the various modes of low-order parameters were regressed to compare their effects on the overlay analysis results. The findings indicate that given a sufficient number of sample points, the usage of multiple linear regression analysis to solve the high-order nonlinear overlay accuracy model containing seventh-order lens distortion parameters is feasible. When the estimated values of low-order overlay distortion parameters are far greater than those of high-order lens distortion parameters, excellent overlay improvement can still be obtained even if the high-order lens distortion parameters are ignored. When the overlay at the four corners of image field obviously exceeds that near the center of image field, it is found, through simulation, that the seventh-order parameters overlay model established in this paper has to be corrected by high-order lens distortion parameters to significantly improve the overlay accuracy.

DATA MODULES

Data Collection Module Data Pre-Processing module Evaluation Module Prediction Module

Data Collection Module:

In this module, the raw data collected by a big mart will be pre-processed for missing data, anomalies and outlier. An algorithm will then be trained to construct a model on that data. It is a system in which three functions are combined. It is used to extract and transform the data from one database into an appropriate format.

Data Pre-Processing Module:

The dataset used is Big Mart sales result and there are total 9 attributes such as Item_Type, item_MRP, Item_fat_content, Outlet_size, Outlet_type, Item_weight, Item_visiblity, Outlet_Establishment_Year, Outlet_Location_Type. Item Outlet Sales is the target variable and the other remaining attributes are independent variable. The pre-processing of data is a method for preparing and adapting raw data to a model of learning. This is the first and significant step to construct a machine learning model. Real-world data generally contain noise, missing values and may not be used in an unusable format especially for machine learning models.

Evaluation Module:

Evaluation of the model is the vital part of creating an efficient machine learning model. Therefore it is important to create a model and get suggestions from it in terms of metrics. It will take and continue until we achieve good accuracy according to the value obtained from metric improvements. Evaluation metrics describe one model's results. The ability to distinguish between model outcomes is an important feature of the evaluation metrics. Here, we used Root Mean Squared Error(RMSE)metric for evaluation process,

Prediction module:

We propose a predictive model using XG boost Regressor technique for predicting the sales of a company like Big Mart and found that the model produces better performance as compared to existing models

CHAPTER 3

SYSTEM REQUIREMENTS

3.1 HARDWARE REQUIREMENTS:

\triangleright	System	:	Pentium i3 Processor.
	Hard Disk	:	500 GB.
	Monitor	:	15" LED
	Input Devices	:	Keyboard, Mouse
\triangleright	Ram	:	4 GB

3.2 SOFTWARE REQUIREMENTS:

\triangleright	Operating system	:	Windows 10.
	Coding Language	:	Python
	Web Framework	:	Flask

3.3 SOFTWARE ENVIRONMENT

Python:

Python is a high-level, interpreted, interactive and object-oriented scripting language. Python is designed to be highly readable. It uses English keywords frequently where as other languages use punctuation, and it has fewer syntactical constructions than other languages.

- Python is Interpreted Python is processed at runtime by the interpreter. You do not need to compile your program before executing it. This is similar to PERL and PHP.
- **Python is Interactive** You can actually sit at a Python prompt and interact with the interpreter directly to write your programs.
- **Python is Object-Oriented** Python supports Object-Oriented style or technique of programming that encapsulates code within objects.

 Python is a Beginner's Language – Python is a great language for the beginner-level programmers and supports the development of a wide range of applications from simple text processing to WWW browsers to games.

3.4 History of Python

Python was developed by Guido van Rossum in the late eighties and early nineties at the National Research Institute for Mathematics and Computer Science in the Netherlands.

Python is derived from many other languages, including ABC, Modula-3, C, C++, Algol-68, SmallTalk, and Unix shell and other scripting languages.

Python is copyrighted. Like Perl, Python source code is now available under the GNU General Public License (GPL).

Python is now maintained by a core development team at the institute, although Guido van Rossum still holds a vital role in directing its progress.

Python Features

Python's features include -

- Easy-to-learn Python has few keywords, simple structure, and a clearly defined syntax. This allows the student to pick up the language quickly.
- Easy-to-read Python code is more clearly defined and visible to the eyes.
- Easy-to-maintain Python's source code is fairly easy-to-maintain.

- A broad standard library Python's bulk of the library is very portable and cross-platform compatible on UNIX, Windows, and Macintosh.
- Interactive Mode Python has support for an interactive mode which allows interactive testing and debugging of snippets of code.
- Portable Python can run on a wide variety of hardware platforms and has the same interface on all platforms.
- Extendable You can add low-level modules to the Python interpreter. These modules enable programmers to add to or customize their tools to be more efficient.
- Databases Python provides interfaces to all major commercial databases.
- GUI Programming Python supports GUI applications that can be created and ported to many system calls, libraries and windows systems, such as Windows MFC, Macintosh, and the X Window system of Unix.
- Scalable Python provides a better structure and support for large programs than shell scripting.

Apart from the above-mentioned features, Python has a big list of good features, few are listed below –

- It supports functional and structured programming methods as well as OOP.
- It can be used as a scripting language or can be compiled to bytecode for building large applications.

- It provides very high-level dynamic data types and supports dynamic type checking.
- It supports automatic garbage collection.
- It can be easily integrated with C, C++, COM, ActiveX, CORBA, and Java.

<u>CHAPTER 4</u> <u>SYSTEM ANALYSIS</u>

4.1 PURPOSE

The purpose of this document is predicting big mart sales using machine learning algorithms. In detail, this document will provide a general description of our project, including user requirements, product perspective, and overview of requirements, general constraints. In addition, it will also provide the specific requirements and functionality needed for this project - such as interface, functional requirements and performance requirements.

4.2 SCOPE

The scope of this SRSdocument persists for the entire life cycle of the project. This document defines the final state of the software requirements agreed upon by the customers and designers. Finally at the end of the project execution all the functionalities may be traceable from the SRSto the product. The document describes the functionality, performance, constraints, interface and reliability for the entire life cycle of the project.

4.3 EXISTING SYSTEM:

- Auto-Regressive Integrated Moving Average, (ARMA) Auto-Regressive Moving Average, have been utilized to develop a few deals forecast standards. Be that as it may, deals anticipating is a refined issue and is influenced by both outer and inside factors, and there are two significant detriments to the measurable technique as set out in A. S. Weigend et A mixture occasional quantum relapse approach and (ARIMA) Auto-Regressive Integrated Moving Average way to deal with every day food deals anticipating were recommend by N. S. Arunraj and furthermore found that the exhibition of the individual model was moderately lower than that of the crossover model.
- E. Hadavandi utilized the incorporation of "Genetic Fuzzy Systems (GFS)" and information gathering to conjecture the deals of the printed circuit board. In their paper, K-means bunching delivered K groups of all information records. At that point, all bunches were taken care of into autonomous with a data set tuning and rule-based extraction ability.
- Perceived work in the field of deals gauging was done by P.A. Castillo, Sales estimating of new distributed books was done in a publication market the executives setting utilizing computational techniques. "Artificial neural organizations" are additionally utilized nearby income estimating. Fluffy Neural Networks have been created with the objective of improving prescient effectiveness, and the Radial "Base Function Neural Network (RBFN)" is required to have an incredible potential for anticipating deals.

DISADVANTAGES OF EXISTING SYSTEM:

- Complex models like neural networks are overkill for simple problems like regression.
- Existing system models prediction analysis which gives less accuracy.
- Forecasting methods and applications contains Lack of Data and short life cycles. So some of the data like historical data, consumer-oriented markets face uncertain demands, can be prediction for accurate result.

4.4 PROPOSED SYSTEM:

- The objective of this proposed system is to predict the future sales from given data of the previous year's using Decision Tree Regression
- Another objective is to conclude the best model which is more efficient and gives fast and accurate result by using Decision Tree Regression.
- To find out key factors that can increase their sales and what changes could be made to the product or store's characteristics.
- Experts also shown that a smart sales forecasting program is required to manage vast volumes of data for business organizations.
- We are predicting the accuracy for Decision Tree Regression. Our predictions help big marts to refine their methodologies and strategies which in turn helps them to increase their profit. The results predicted will be very useful for the executives of the company to know about their sales and profits. This will also give them the idea for their new locations or Centre's of Bigmart

ADVANTAGES OF PROPOSED SYSTEM:

- Business assessments are based on the speed and precision of the methods used to analyze the results. The Machine Learning Methods presented in this research paper should provide an effective method for data shaping and decision-making.
- New approaches that can better identify consumer needs and formulate marketing plans will be implemented.
- The outcome of machine learning algorithm will help to select the most suitable demand prediction algorithm and with the aid of which BigMart will prepare its marketing campaigns.

CHAPTER 5

SYSTEM DESIGN

5.1 INPUT DESIGN

The input design is the link between the information system and the user. It comprises the developing specification and procedures for data preparation and those steps are necessary to put transaction data in to a usable form for processing can be achieved by inspecting the computer to read data from a written or printed document or it can occur by having people keying the data directly into the system. The design of input focuses on controlling the amount of input required, controlling the errors, avoiding delay, avoiding extra steps and keeping the process simple. The input is designed in such a way so that it provides security and ease of use with retaining the privacy. Input Design considered the following things:

- > What data should be given as input?
- > How the data should be arranged or coded?
- > The dialog to guide the operating personnel in providing input.
- > Methods for preparing input validations and steps to follow when error occur.

5.2 OUTPUT DESIGN

A quality output is one, which meets the requirements of the end user and presents the information clearly. In any system results of processing are communicated to the users and to other system through outputs. In output design it is determined how the information is to be displaced for immediate need and also the hard copy output. It is the most important and direct source information to the user. Efficient and intelligent output design improves the system's relationship to help user decision-making.

The output form of an information system should accomplish one or more of the following objectives.

- Convey information about past activities, current status or projections of the
- Future.
- Signal important events, opportunities, problems, or warnings.
- Trigger an action.
- Confirm an action

5.3 DATA FLOW DIAGRAM

23

- The DFD is also called as bubble chart. It is a simple graphical formalism that can be used to represent a system in terms of input data to the system, various processing carried out on this data, and the output data is generated by this system.
- 2. The data flow diagram (DFD) is one of the most important modeling tools. It is used to model the system components. These components are the system process, the data used by the process, an external entity that interacts with the system and the information flows in the system.
- DFD shows how the information moves through the system and how it is modified by a series of transformations. It is a graphical technique that depicts information flow and the transformations that are applied as data moves from input to output.
- 4. DFD is also known as bubble chart. A DFD may be used to represent a system at any level of abstraction. DFD may be partitioned into levels that represent increasing information flow and functional detail.



UML DIAGRAMS

UML stands for Unified Modeling Language. UML is a standardized generalpurpose modeling language in the field of object-oriented software engineering. The standard is managed, and was created by, the Object Management Group.

The goal is for UML to become a common language for creating models of object oriented computer software. In its current form UML is comprised of two major components: a Meta-model and a notation. In the future, some form of method or process may also be added to; or associated with, UML.

The Unified Modeling Language is a standard language for specifying, Visualization, Constructing and documenting the artifacts of software system, as well as for business modeling and other non-software systems.

The UML represents a collection of best engineering practices that have proven successful in the modeling of large and complex systems.

The UML is a very important part of developing objects oriented software and the software development process. The UML uses mostly graphical notations to express the design of software projects.

GOALS:

The Primary goals in the design of the UML are as follows:

- 1. Provide users a ready-to-use, expressive visual modeling Language so that they can develop and exchange meaningful models.
- 2. Provide extendibility and specialization mechanisms to extend the core concepts.
- 3. Be independent of particular programming languages and development process.
- 4. Provide a formal basis for understanding the modeling language.
- 5. Encourage the growth of OO tools market.
- 6. Support higher level development concepts such as collaborations, frameworks, patterns and components.
- 7. Integrate best practices.

USE CASE DIAGRAM

A use case diagram in the Unified Modeling Language (UML) is a type of behavioral diagram defined by and created from a Use-case analysis. Its purpose is to present a

graphical overview of the functionality provided by a system in terms of actors, their goals (represented as use cases), and any dependencies between those use cases. The main purpose of a use case diagram is to show what system functions are performed for which actor. Roles of the actors in the system can be depicted.



SEQUENCE DIAGRAM

A sequence diagram in Unified Modeling Language (UML) is a kind of interaction diagram that shows how processes operate with one another and in what order. It is a construct of a Message Sequence Chart. Sequence diagrams are sometimes called event diagrams, event scenarios, and timing diagrams.



CHAPTER 6

6.1 MODULES:

- Data Collection
- Dataset
- Data Preparation
- Model Selection
- Analyze and Prediction
- Accuracy on test set
- Saving the Trained Model

MODULES DESCSRIPTION:

6.1.1 Data Collection:

This is the first real step towards the real development of a machine learning model, collecting data. This is a critical step that will cascade in how good the model will be, the more and better data that we get, the better our model will perform.

There are several techniques to collect the data, like web scraping, manual interventions and etc.

Predictive Analysis for Big Mart Sales Using Machine Learning Algorithms

Data set Link: https://www.kaggle.com/shivan118/big-mart-sales-prediction-datasets

6.1.2 Dataset:

The dataset consists of 8523 individual data. There are 12 columns in the dataset, which are described below.

1. **Item***Identifier* ----Unique product ID 2.ItemWeight ----Weight of product 3.ItemFatContent ----Whether the product is low fat or not 4.Item Visibility ---- The % of the total display area of all products in a store allocated to the particular product 5.ItemType ----The category which to the product belongs 6.ltem*MRP* -----Price product Maximum Retail (list price) of the 7.OutletIdentifier ----Unique ID store 8.Outlet EstablishmentYear ---- The year in which the store was established 9.Outlet Size ---- The size of the store in terms of ground area covered **10.OutletLocation** Type ---- The type of city in which the store is located 11.*OutletType ---- Whether the outlet is just a grocery store or some sort ofsupermarket

12.Item*Outlet***Sales** ---- sales of the product in t particular store. This is the outcome variable to be predicted.

2.

6.1.3 Data Preparation:

Wrangle data and prepare it for training. Clean that which may require it (remove duplicates, correct errors, deal with missing values, normalization, data type conversions, etc.)

28

Randomize data, which erases the effects of the particular order in which we collected and/or otherwise prepared our data

Visualize data to help detect relevant relationships between variables or class imbalances (bias alert!), or perform other exploratory analysis

Split into training and evaluation sets

6.1.4 Model Selection:

We used decision tree regression machine learning algorithm , We got a accuracy of 95.7% on test set so we implemented this algorithm.

Decision tree regression

Decision Tree is a decision-making tool that uses a flowchart-like tree structure or is a model of decisions and all of their possible results, including outcomes, input costs, and utility. Decision-tree algorithm falls under the category of supervised learning algorithms. It works for both continuous as well as categorical output variables. The branches/edges represent the result of the node and the nodes have either:

Conditions [Decision Nodes]

Result [End Nodes]

The branches/edges represent the truth/falsity of the statement and take makes a decision based on that in the example below which shows a decision tree that evaluates the smallest of three numbers:

Decision Tree Regression: Decision tree regression observes features of an object and trains a model in the structure of a tree to predict data in the future to produce meaningful continuous output. Continuous output means that the output/result is not discrete, i.e., it is not represented just by a discrete, known set of numbers or values.

6.1.5 Analyze and Prediction:

In the actual dataset, we chose only 9 features:

1.ItemWeight ----Weight of product 2.ItemFatContent ----Whether the product is low fat or not 3.Item Visibility ---- The % of the total display area of all products in a store allocated the particular to product The 4.ItemType ----category to which the product belongs **5.Item***MRP* -----Price Maximum Retail (list of price) the product 6.OutletEstablishmentYear ---- The year in which the store was established 7.Outlet Size ---- The size of the store in terms of ground area covered 8.OutletLocation Type ---- The type of city in which the store is located 9.*OutletType ---- Whether the outlet is just a grocery store or some sort of supermarket

6.1.6 Accuracy on test set:

We got an accuracy of 95.80% on test set.

6.1. 7 Saving the Trained Model:

Once you're confident enough to take your trained and tested model into the productionready environment, the first step is to save it into a .h5 (or) .pkl file using a library like pickle .

Make sure you have pickle installed in your environment.

Next, let's import the module and dump the model into .pkl file

CHAPTER 7

SYSTEM IMPLEMENTATION:

7.1 SYSTEM ARCHITECTURE:

The system architectural design is the design process for identifying the subsystems making up the system and framework for subsystem control and communication. The goal of the architectural design is to establish the overall structure of software system.



CHAPTER-8 SYSTEM TESTING

The purpose of testing is to discover errors. Testing is the process of trying to discover every conceivable fault or weakness in a work product. It provides a way to check the functionality of components, sub assemblies, assemblies and/or a finished product It is the process of exercising software with the intent of ensuring that the

Software system meets its requirements and user expectations and does not fail in an unacceptable manner. There are various types of test. Each test type addresses a specific testing requirement.

TYPES OF TESTS

8.1 Unit testing

Unit testing involves the design of test cases that validate that the internal program logic is functioning properly, and that program inputs produce valid outputs. All decision branches and internal code flow should be validated. It is the testing of individual software units of the application .it is done after the completion of an individual unit before integration. This is a structural testing, that relies on knowledge of its construction and is invasive. Unit tests perform basic tests at component level and test a specific business process, application, and/or system configuration. Unit tests ensure that each unique path of a business process performs accurately to the documented specifications and contains clearly defined inputs and expected results.

8.2 Integration testing

Integration tests are designed to test integrated software components to determine if they actually run as one program. Testing is event driven and is more concerned with the basic outcome of screens or fields. Integration tests demonstrate that although the components were individually satisfaction, as shown by successfully unit testing, the combination of components is correct and consistent. Integration testing is specifically aimed at exposing the problems that arise from the combination of components.

8.3 Functional test

Functional tests provide systematic demonstrations that functions tested are available as specified by the business and technical requirements, system documentation, and user manuals.

Functional testing is centered on the following items:

Valid Input : identified classes of valid input must be accepted.

Invalid Input : identified classes of invalid input must be rejected.

Functions : identified functions must be exercised.

Output : identified classes of application outputs must be exercised.

Systems/Procedures: interfacing systems or procedures must be invoked.

Organization and preparation of functional tests is focused on requirements, key functions, or special test cases. In addition, systematic coverage pertaining to identify Business process flows; data fields, predefined processes, and successive processes must be considered for testing. Before functional testing is complete, additional tests are identified and the effective value of current tests is determined.

8.4 System Test

System testing ensures that the entire integrated software system meets requirements. It tests a configuration to ensure known and predictable results. An example of system testing is the configuration oriented system integration test. System testing is based on process descriptions and flows, emphasizing pre-driven process links and integration points.

8.5 White Box Testing

White Box Testing is a testing in which in which the software tester has knowledge of the inner workings, structure and language of the software, or at least its purpose. It is purpose. It is used to test areas that cannot be reached from a black box level.

8.6 Black Box Testing

Black Box Testing is testing the software without any knowledge of the inner workings, structure or language of the module being tested. Black box tests, as most other kinds of tests, must be written from a definitive source document, such as specification or requirements document, such as specification or requirements document. It is a testing in which the software under test is treated, as a black box .you cannot "see" into it. The test provides inputs and responds to outputs without considering how the software works.

Unit Testing:

34

Unit testing is usually conducted as part of a combined code and unit test phase of the software lifecycle, although it is not uncommon for coding and unit testing to be conducted as two distinct phases.

Test strategy and approach

Field testing will be performed manually and functional tests will be written in detail.

Test objectives

- All field entries must work properly.
- Pages must be activated from the identified link.
- The entry screen, messages and responses must not be delayed.

Features to be tested

- Verify that the entries are of the correct format
- No duplicate entries should be allowed
- All links should take the user to the correct page.

Integration Testing

Software integration testing is the incremental integration testing of two or more integrated software components on a single platform to produce failures caused by interface defects.

The task of the integration test is to check that components or software applications, e.g. components in a software system or – one step up – software applications at the company level – interact without error.

Test Results: All the test cases mentioned above passed successfully. No defects encountered.

8.7 Acceptance Testing

User Acceptance Testing is a critical phase of any project and requires significant participation by the end user. It also ensures that the sytem meets the functional requirements.

Test Results: All the test cases mentioned above passed successfully. No defects encountered.

CHAPTER 9

RESULTS AND DISCUSSION

SalesPrediction

Home Login Upload



FIG 9.1 IOG IN PAGE

SalesPrediction

Home Login Upload Preview



FIG 9.2 DATA COLLECTION FROM KAGGLE

4597	NCJ05	18.700	Low Fat	0.046050	Health and Hygiene	151.9682	OUT013	1987
4598	FDS10	19.200	Low Fat	0.035385	Snack Foods	180.3318	OUT017	2007
4599	FDK34		Low Fat	0.038340	Snock	240.1564	OUT027	1985

<								>
			FIG 9.3 PREPOC	ESSED DATASET	OF BIGMART S	ALES		
	+							
±5000	predict							
llabe	M Small	(A The Complete C pr.,	🤨 El Main Chamility 🧟	ta https://conten.orly.	Oliviti Linear Regr.	C (TTTE) EASY WAY	S Skopw	* 📒 DB
Pr	edic	tion		Home L	ogin Upload	Prediction Cho	irt	



FIG 9.4 PREDICTION OF SALES OF A PRODUCT

sPr	edic	tion	10	lome	Login	Upload	Prediction	Chart	Performance analysis		
Date 1	N Crue I	No The Complete Con-	O IEE Main Charrier	01	insultra	elin	Citral Low Rep.		pada want. 🚦 Repe	2.1	1 OP4
1.5005/	hirt.										
*	+										





FIG 9.5 COMPARISION GRAPH ON SALES OF ITEM TYPE

🖬 🖉 🖬 🖬 💭 🖬 🖉 🖷

~ C 11

sPr	edic	tion	Hor	me	Login	Upload	Prediction	Chart	Performance analysis		
fute #	H Croài	 The Complete C pr 	EEMan Ownshy	0	60gu//cs	attau.ioty	(Attitut) Circler Regi-	0 (ItthE	EADA MAR 🕤 SAyee	*	Ciń
±50003	arlomer.	-									
*	+										

Performance analysis

Mean Absolute Error: 0.31

Mean Squared Error: 0.53

R^2 Score: 0.9596

Accuracy score: 0.95636

🖬 🔎 🖬 💭 🐂 🕘 🔮 📕 📮 🔍 📲 🔍 📲

FIG 9.6 PERFORMANCE VALUES OF PROPOSED MODEL

CHAPTER 10

CONCLUSION

In this work, the effectiveness of Decision Tree Regression on the data on revenue and review of, best performance-algorithm, here propose software to using regression approach for predicting the sales centered on sales data from the past the accuracy of linear regression prediction can be enhanced with this method, and Decision Tree Regression can be determined. So, we can conclude Decision Tree Regression gives the better prediction with respect to Accuracy.

FUTURE WORK:

In future, the forecasting sales and building a sales plan can help to avoid unforeseen cash flow and manage production, staff and financing needs more effectively. In future work we can also consider with the ARIMA model which shows the time series graph.

CHAPTER-11

REFERENCES

[1] Ching Wu Chu and Guoqiang Peter Zhang, "A comparative study of linear and nonlinear models for aggregate retails sales forecasting", Int. Journal Production Economics, vol. 86, pp. 217- 231, 2003.

[2] Wang, Haoxiang. "Sustainable development and management in consumer electronics using soft computation." Journal of Soft Computing Paradigm (JSCP) 1, no. 01 (2019): 56.- 2.

[3] Suma, V., and Shavige Malleshwara Hills. "Data Mining based Prediction of Demand in Indian Market for Refurbished Electronics." Journal of Soft Computing Paradigm (JSCP) 2, no. 02 (2020): 101-110 [4] Giuseppe Nunnari, Valeria Nunnari, "Forecasting Monthly Sales Retail Time Series: A Case Study", Proc. of IEEE Conf. on Business Informatics (CBI), July 2017.

[5]https://halobi.com/blog/sales-forecasting-five-uses/.

[6] Zone-Ching Lin, Wen-Jang Wu, "Multiple Linear Regression Analysis of the Overlay Accuracy Model Zone", IEEE Trans. On Semiconductor Manufacturing, vol. 12, no. 2, pp. 229 – 237, May 1999.

[7] O. Ajao Isaac, A. Abdullahi Adedeji, I. Raji Ismail, "Polynomial Regression Model of Making Cost Prediction In Mixed Cost Analysis", Int. Journal on Mathematical Theory and Modeling, vol. 2, no. 2, pp. 14 – 23, 2012.

[8] C. Saunders, A. Gammerman and V. Vovk, "Ridge Regression Learning Algorithm in Dual Variables", Proc. of Int. Conf. on Machine Learning, pp. 515 – 521, July 1998.IEEE TRANSACTIONS ON INFORMATION THEORY, VOL. 56, NO. 7, JULY 2010 3561.

[9] "Robust Regression and Lasso". Huan Xu, Constantine Caramanis, Member, IEEE, and Shie Mannor, Senior Member, IEEE. 2015 International Conference on Industrial Informatics-Computing Technology, Intelligent Technology, Industrial Information Integration."An improved Adaboost algorithm based on uncertain functions". Shu Xinqing School of Automation Wuhan University of Technology. Wuhan, China Wang Pan School of the Automation Wuhan University of Technology Wuhan, China. [10] Xinqing Shu, Pan Wang, "An Improved Adaboost Algorithm based on Uncertain Functions", Proc. of Int. Conf. on Industrial Informatics – Computing Technology, Intelligent Technology, Industrial Information Integration, Dec. 2015.

[11] A. S. Weigend and N. A. Gershenfeld, "Time series prediction: Forecasting the future and understanding the past", Addison-Wesley, 1994.

[12] N. S. Arunraj, D. Ahrens, A hybrid seasonal autoregressive integrated moving average and quantile regression for daily food sales forecasting, Int. J. Production Economics 170 (2015) 321-335P

[13] D. Fantazzini, Z. Toktamysova, Forecasting German car sales using Google data and multivariate models, Int. J. Production Economics 170 (2015) 97-135.

[14] X. Yua, Z. Qi, Y. Zhao, Support Vector Regression for Newspaper/Magazine Sales Forecasting, Procedia Computer Science 17 (2013) 1055–1062.

[15] E. Hadavandi, H. Shavandi, A. Ghanbari, An improved sales forecasting approach by the integration of genetic fuzzy systems and data clustering: a Case study of the printed circuit board, Expert Systems with Applications 38 (2011) 9392–9399.

42

[16] P. A. Castillo, A. Mora, H. Faris, J.J. Merelo, P. GarciaSanchez, A.J. Fernandez-Ares, P. De las Cuevas, M.I. Garcia-Arenas, Applying computational intelligence methods for predicting the sales of newly published books in a real editorial business management environment, Knowledge-Based Systems 115 (2017) 133-151.

[17] R. Majhi, G. Panda and G. Sahoo, "Development and performance evaluation of FLANN based model for forecasting of stock markets". Expert Systems with Applications, vol. 36, issue 3, part 2, pp. 6800-6808, April 2009.

[18] Pei Chann Chang and Yen-Wen Wang, "Fuzzy Delphi and back propagation model for sales forecasting in PCB industry", Expert systems with applications, vol. 30,pp. 715-726, 2006.

[19] R. J. Kuo, Tung Lai HU and Zhen Yao Chen "application of radial basis function neural networks for sales forecasting", Proc. Of Int. Asian Conference on Informatics in control, automation, and robotics, pp. 325-328, 2009.

[20] R. Majhi, G. Panda, G. Sahoo, and A. Panda, "On the development of Improved Adaptive Models for Efficient Prediction of Stock Indices using Clonal-PSO (CPSO) and PSO Techniques", International Journal of Business Forecasting and Market Intelligence, vol. 1, no. 1, pp.50-67, 2008.

43

[21]Suresh K and Praveen O, "Extracting of Patterns Using Mining Methods Over Damped Window," 2020 Second International Conference on Inventive Research in Computing Applications (ICIRCA), Coimbatore, India, 2020, pp. 235-241, DOI:

10.1109/ICIRCA48905.2020.9182893.

[22] Shobha Rani, N., Kavyashree, S., & Harshitha, R. (2020). ObjectDetection in Natural Scene Images Using Thresholding Techniques.Proceedings of the International Conference on Intelligent Computing andControl Systems, ICICCS 2020, Iciccs, 509–515.

[23] https://www.kaggle.com/brijbhushannanda1979/bigmartsalesdata.