

BITCOIN PRICE ANALYZE AND PREDICTION USING DATA SCIENCE PROCESS

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

by

KOLLIMALLI RAJA RAMESH (Reg. No. 38110253)
VELUGUBANTI NARENDRA (Reg. No. 38110356)



**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

MAY - 2022



SATHYABAMA

**INSTITUTE OF SCIENCE AND TECHNOLOGY
(DEEMED TO BE UNIVERSITY)**

Accredited "A" Grade by NAAC | 12B Status by UGC | Approved by AICTE

www.sathyabama.ac.in

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

BONAFIDE CERTIFICATE

This is to certify that this Project Report is the bonafide work of **KOLLIMALLI RAJA RAMESH (Reg. No: 38110253), VELUGUBANTI NARENDRA (Reg. No: 38110356)** who carried out the project entitled "**BITCOIN PRICE ANALYZE AND PREDICTION USING DATA SCIENCE PROCESS**" under my supervision from November 2021 to April 2022.

**Internal Guide
Dr. G.NAGARAJAN.,**

External Guide (if Applicable)

**Head of the Department
Dr. L. LAKSHMANAN M.E., Ph.D.,**

Submitted for Viva voce Examination held on _____

Internal Examiner

External Examiner

DECLARATION

I **KOLLIMALLI RAJA RAMESH, VELUGUBANTI NARENDRA** hereby declare that the Project Report entitled **BITCOIN PRICE ANALYZE AND PREDICTION USING DATA SCIENCE PROCESS** done by me under the guidance of **Dr. M. Maheswari M.E., Ph.D., (Internal)** is submitted in partial fulfillment of the requirements for the award of Bachelor of Engineering Degree in Computer Science and Engineering.

DATE: 30-03-2022

K.RAJA RAMESH, V.NARENDRA

PLACE: CHENNAI

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. Head of the Department, Dept. 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 Dr. G.NAGARAJAN 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

Bitcoin is a digital asset and a payment system that is used as a form of Internet currency. It allows for anonymous payment from one person to another and is therefore a preferred payment method for criminal actions on the Internet. Recently Bitcoin has received a lot of attention from the media and the public due to its recent price hike. The objective of this paper is to determine the predictable price direction of Bitcoin price. Machine learning models can likely give us the insight we need to learn about the future of Cryptocurrency. It will not tell us the future but it might tell us the general trend and direction to expect the prices to move. The proposed model is to build a machine learning model where the data is used to made to learn about the pattern in the dataset and the machine learning algorithm is used to predict the bitcoin price.

TABLE OF CONTENTS

CHAPTER No.	TITLE	PAGE No
	ABSTRACT	i
	LIST OF FIGURES	v
1	INTRODUCTION	1
	1.1 Objective of the project	1
	1.1.1 Necessity	1
	1.1.2 Software development method	1
	1.1.3 Layout of the document	1
	1.2 Overview of the designed project	2
2	LITERATURE SURVEY	3
	2.1 Literature Survey	3
3	AIM AND SCOPE OF THE PRESENT INVESTIGATION	11
	3.1 Project Proposal	11
	3.1.1 Mission	11
	3.1.2 Goal	11
	3.2 Scope of the Project	11
	3.3 Overview of the project	11
	3.4 Existing system	12
	3.4.1 Disadvantages	12
	3.5 Preparing the dataset	12
	3.6 Proposed system	13
	3.6.1 Exploratory Data Analysis of loan approval	13
	3.6.2 Data Wrangling	13
	3.6.3 Data collection	13
	3.6.4 Building the classification model	13
	3.6.5 Advantages	14
	3.8 Flow chart	15

4	EXPERIMENTAL OR MATERIALS AND METHODS; ALGORITHMS USED	16
4.1	System Study	16
4.1.1	System requirement specifications	16
4.2	System Specifications	16
4.2.1	Machine Learning Overview	16
4.2.2	Flask Overview	17
4.3	Steps to download & install Python	17
4.3.1	IDE Installation for python	17
4.3.2	Python File Creation	17
4.4	Python Libraries needed	17
4.4.1	Numpy library	18
4.4.2	Pandas library	18
4.4.3	Matplotlib library	18
4.4.4	Seaborn library	19
4.4.5	Scikit Learn library	19
4.4.6	Flask	19
4.5	Modules	20
4.6	UML diagrams	21
4.6.1	Use Case Diagram	21
4.6.2	Class Diagram	22
4.6.3	Activity Diagram	23
4.6.4	Sequence Diagram	24
4.6.5	Entity Relationship Diagram	25
4.7	Module Details	26
4.7.1	Data Pre-processing	26
4.7.2	Data Validation /Cleaning /Preparing Process	27
4.7.3	Exploration data analysis of visualization	28
4.7.4	Comparing Algorithm with prediction in the form of best accuracy result	29
4.7.5	Algorithm and Techniques	33
4.7.6	Deployment Using Flask	44

5	RESULTS AND DISCUSSION, PERFORMANCE ANALYSIS	46
	5.1 Performance Analysis	46
	5.2 Discussion	46
6	SUMMARY AND CONCLUSION	48
	6.1 Summary	48
	6.2 Conclusion	48
	6.3 Future Work	48
	REFERENCES	49
	APPENDIX	
	A. SOURCE CODE	50-60
	B. SCREENSHOTS	61-62

LIST OF FIGURES

FIGURE NO.	FIGURE NAME	PAGE NO.
3.1	Architecture of Proposed Model	11
3.2	Flow chart	12
4.1	System Architecture	21
4.2	Use Case Diagrams	21
4.3	Class Diagram	22
4.4	Activity Diagram	23
4.5	Sequence Diagram	24
4.6	Entity Relationship Diagram	25
4.7	LINEAR Regression	34
4.8	Random Forest Classifier	37
4.9	Decision Tree Classifier	38
4.10	Lasso Regression	40
4.11	Support Vector Classifier	42
4.12	Support Vector Classifier	42
4.13	Gradient Boost	43
6.1	Machine Learning Algorithms Accuracy	50
6.2	Home Page	60
6.3	Output Graph	61

CHAPTER-1

INTRODUCTION

1.1 OBJECTIVE OF THE PROJECT:

The goal is to develop a machine learning model for Bank Loan Approval Prediction, to potentially replace the updatable supervised machine learning classification models by predicting results in the form of best accuracy by comparing supervised algorithm.

1.1.1 *Necessity:*

This online bank loan approval system helps in overcoming the time management. This Application is very easy to use. It can work accurately and very smoothly in a different scenario. It reduces the effort workload and increases efficiency in work. In aspects of time value, it is worthy. In this website the user can check the loan status easily whether approved or not.

1.1.2 *Software development method:*

In many software applications program different methods and cases are followed such as, Waterfall model, Iterative model, Spiral model, V-model and Big Bang model. I used waterfall model in this application. I tried to use test case and case software approaches.

1.1.3 *Layout of the document:*

This documentation starts with formal introduction. After introduction analysis and design of the project are described. In analysis and design of the project have many parts such as project proposal, mission, goal, target audience, environment. Use cases and test cases are in chapter 2 and chapter 3 respectively. Finally, this documentation finished with result and Conclusion part.

1.2 OVERVIEW OF THE DESIGNED PROJECT:

At first, we take the dataset from our resource then we have to perform data-preprocessing, visualization methods for cleaning and visualizing the dataset respectively and we applied the Machine Learning algorithms on the dataset then we generate the pickle file for best algorithm and flask is used as user interface for displaying the result.

CHAPTER-2

LITERATURE SURVEY

2.1 LITERATURE SURVEY:

General

A literature review is a body of text that aims to review the critical points of current knowledge on and/or methodological approaches to a particular topic. It is secondary sources and discuss published information in a particular subject area and sometimes information in a particular subject area within a certain time period. Its ultimate goal is to bring the reader up to date with current literature on a topic and forms the basis for another goal, such as future research that may be needed in the area and precedes a research proposal and may be just a simple summary of sources. Usually, it has an organizational pattern and combines both summary and synthesis.

A summary is a recap of important information about the source, but a synthesis is a re-organization, reshuffling of information. It might give a new interpretation of old material or combine new with old interpretations or it might trace the intellectual progression of the field, including major debates.

Depending on the situation, the literature review may evaluate the sources and advise the reader on the most pertinent or relevant of them

Review of Literature Survey

Title : Enhancing Bitcoin Price Fluctuation Prediction Using Attentive LSTM and Embedding Network

Author: Yang Li, Zibin Zheng and Hong-Ning Dai

Year : 2019

Bitcoin has attracted extensive attention from investors, researchers,

regulators, and the media. A well-known and unusual feature is that Bitcoin's price often fluctuates significantly, which has however received less attention. In this paper, we investigate the Bitcoin price fluctuation prediction problem, which can be described as whether Bitcoin price keeps or reversals after a large fluctuation. In this paper, three kinds of features are presented for the price fluctuation prediction, including basic features, traditional technical trading indicators, and features generated by a Denoising autoencoder. We evaluate these features using an Attentive LSTM network and an Embedding Network (ALLEN). In particular, an attentive LSTM network can capture the time dependency representation of Bitcoin price and an embedding network can capture the hidden representations from related cryptocurrencies. Experimental results demonstrate that ALLEN achieves superior state-of-the-art performance among all baselines. Furthermore, we investigate the impact of parameters on the Bitcoin price fluctuation prediction problem, which can be further used in a real trading environment by investors

Title : A Research On Bitcoin Price Prediction Using Machine Learning Algorithms.

Author: Lekkala Sreekanth Reddy, Dr.P. Sriramya

Year : 2020

In this paper, we proposed to predict the Bitcoin price accurately taking into consideration various parameters that affect the Bitcoin value. By gathering information from different reference papers and applying in real time ,I found the advantages and disadvantages of bitcoin price prediction. Each and every paper has its own set of methodologies of bitcoin price prediction. Many papers has accurate price but some other don't, but the time complexity is higher in those predictions, so to reduce the time complexity here in this paper we use an algorithm linked to artificial intelligence named LASSO(least absolute shrinkage selection operator. The other papers used different algorithms like SVM(support vector machine),coinmarkupcap, Quandl, GLM, CNN(Convolutional Neural Networks)and RNN(Recurrent neural networks) etc.. which do not have a great time management, but in LASSO finding of the results from a larger database is quick and fast..so for this purpose we draw a comparison between other algorithms and the LASSO algorithm, this survey paper helps the upcoming researchers to make an impact in the their papers. The process happens in the paper is first moment of the research, we aim to understand and find daily trends in the Bitcoin market while gaining insight into optimal features surrounding Bitcoin price. Our data set consists of various features relating to the Bitcoin price and payment network over the course of every years, recorded daily. By preprocessing the dataset, we apply the some data mining techniques to reduce the noise of data. Then the second moment of our research, using the

available information, we will predict the sign of the daily price change with
highestpossibleaccuracy

Title : Bitcoin Price Prediction using Machine Learning.

Author: Mr. Shivam Pandey¹, Mr. Anil Chavan².

Year : 2021

In this paper, we attempt to predict the Bitcoin price accurately taking into consideration various parameters that affect the Bitcoin value. For the first phase of our survey, we aim to understand and identify daily trends in the Bitcoin market while gaining insight into optimal features surrounding Bitcoin price. For the second phase of our survey, using the available information, we will predict the sign of the daily price change with highest possible accuracy. Predicting the future will always be on the top of the list of uses for machine learning algorithms. Here in this project we have attempted to predict the prices of Bitcoins using two deep learning methodologies. This work focuses on the development of project based learning in the field of computer science engineering, by taking into account the problem definition, progression, student assessment and use of hands on activities based on use of learning algorithm to develop application.

Title : Forecasting cryptocurrency returns and volume using search engines

Author: Muhammad Ali Nasir¹, Toan Luu Duc Huynh.

Year : 2021

In the context of the debate on the role of cryptocurrencies in the economy as well as their dynamics and forecasting, this brief study analyzes the predictability of Bitcoin volume and returns using Google search values. We employed a rich set of established empirical approaches, including a VAR framework, a copulas approach, and non-parametric drawings, to capture a dependence structure. Using a weekly dataset from 2013 to 2017, our key results suggest that the frequency of Google searches leads to positive returns and a surge in Bitcoin trading volume. Shocks to search values have a positive effect, which persisted for at least a week. Our findings contribute to the debate on cryptocurrencies/Bitcoins and have profound implications in terms of understanding their dynamics, which are of special interest to investors and economic policymakers.

Title : Prediction for Loan Approval using Machine Learning Algorithm

Author: Ashwini S. Kadam, Shraddha R. Nikam, Ankita A. Aher, Gayatri V. Shelke, Amar S. Chandgude

Year : 2021

In our banking system, banks have many products to sell but main source of income of any banks is on its credit line. So they can earn from interest of those loans which they credits. A bank's profit or a loss depends to a large extent on loans i.e. whether the customers are paying back the loan or defaulting. By predicting the loan defaulters, the bank can reduce its Non-performing Assets. This makes the study of this phenomenon very important. Previous research in this era has shown that there are so many methods to study the problem of controlling loan default. But as the right predictions are very important for the maximization of profits, it is essential to study the nature of the different methods and their comparison. A very important approach in predictive analytics is used

to study the problem of predicting loan defaulters (i) Collection of Data, (ii) Data Cleaning and (iii) Performance Evaluation. Experimental tests found that the Naïve Bayes model has better performance than other models in terms of loan forecasting.

Title : Social signals and algorithmic trading of Bitcoin

Author: David Garcia and Frank Schweitzer.

Year : 2015

The availability of data on digital traces is growing to unprecedented sizes, but inferring actionable knowledge from large-scale data is far from being trivial. This is especially important for computational finance, where digital traces of human behaviour offer a great potential to drive trading strategies. We contribute to this by providing a consistent approach that integrates various datasources in the design of algorithmic traders. This allows us to derive insights into the principles behind the profitability of our trading strategies. We illustrate our approach through the analysis of Bitcoin, a cryptocurrency known for its large price fluctuations. In our analysis, we include economic signals of volume and price of exchange for USD, adoption of the Bitcoin technology and transaction volume of Bitcoin. We add social signals related to information search, word of mouth volume, emotional valence and opinion polarization as expressed in tweets related to Bitcoin for more than 3 years. Our analysis reveals that increases in opinion polarization and exchange volume precede rising Bitcoin prices, and that emotional valence precedes opinion polarization and rising exchange volumes. We apply these insights to design algorithmic trading strategies for Bitcoin, reaching very high profits in less than a year. We verify this high profitability with robust statistical methods that take into account risk and trading costs, confirming the longstanding hypothesis that trading-based social media sentiment has the potential to yield positive returns on investment.

CHAPTER-3

AIM AND SCOPE OF THE PRESENT INVESTIGATION

3.1 PROJECT PROPOSAL:

The project proposal is the term of documents. A project can describe the project proposal. It is the set of all plans of a project. Like, how the software works, what are the steps to complete the entire projects, and what are the software requirements and analysis for this project. In my project, I am doing all the steps and also risk and reward and other project dependencies in the project proposal.

3.1.1 Mission:

The goal is to develop a machine learning model for Bitcoin Prediction, to potentially replace the updatable supervised machine learning classification models by predicting results in the form of best accuracy by comparingsupervised algorithm.

3.1.2 Goal:

The goal is to develop a machine learning model for Bitcoin Price Prediction.

3.2 SCOPE OF THE PROJECT:

Bitcoin is an internet-based medium of exchange in the form of digital assets which uses cryptographic functions to conduct financial transactions.Bitcoin leverage block chain technology to gain decentralization, transparency,and immutability. The main scope of the project is to finding the accuracy,Minimize the error rate and getting result from the flask framework deployment.

3.3 OVERVIEW OF THE PROJECT:

All historic open, high, low, close, trading volume and market cap info for all Bitcoin.I've had to go over the code with a fine tooth comb to get it compatible with CRAN so there have been significant enhancements to howsome of thefield conversions have been undertaken and the data being cleaned.This should eliminate a few issues around number formatting orunexpected handling of scientific notations.

3.4 EXISTING SYSTEM:

Bitcoin is a digital asset and a payment system that is used as a form of Internet currency. It allows for anonymous payment from one person to another and is therefore a preferred payment method for criminal actions on the Internet. Recently Bitcoin has received a lot of attention from the media and the public due to its recent price hike. The objective of this paper is to determine the predictable price direction of Bitcoin price. Machine learning models can likely give us the insight we need to learn about the future of Cryptocurrency. It will not tell us the future but it might tell us the general trend and direction to expect the prices to move. The proposed model is to build a machine learning model where the data is used to made to learn about the pattern in the dataset and the machine learning algorithm is used to predict the bitcoin price.

3.4.1 Disadvantages:

- 1.They had made only data analysis and they did not build a predicting model.
- 2.The classification model was not discussed and performance metrics like accuracy are not calculated.

3.5 PREPARING THE DATASET:

This dataset was created in order to build models for bitcoin price prediction.it contains

- The price of bitcoin [USD]
- The total number of bitcoin confirmed transactions per day
- Average transaction fees in USD per bitcoin transaction [USD]
- Google bitcoin trends search
- Gold ounce price [USD]
- Oil WTI price [USD]
- M2 money supply in the USA
- SP500 close index
- The time period is between 12.2014 - 04.2020

3.6 PROPOSED SYSTEM:

3.6.1 Exploratory Data Analysis of Bitcoin price

Multiple datasets from different sources would be combined to form a generalized dataset, and then different machine learning algorithms would be applied to extract patterns and to obtain results with maximum accuracy.

3.6.2 Data Wrangling

In this section of the report will load in the data, check for cleanliness, and then trim and clean given dataset for analysis. Make sure that the document steps carefully and justify for cleaning decisions.

3.6.3 Data collection

The data set collected for predicting given data is split into Training set and Test set. Generally, 7:3 ratios are applied to split the Training set and Test set. The Data Model which was created using machine learning algorithms are applied on the Training set and based on the test result accuracy, Test set prediction is done.

3.6.4 Building the classification model

The prediction of Price of bitcoin, ML algorithm prediction model is effective because of the following reasons: It provides better results in classification problem.

- It is strong in preprocessing outliers, irrelevant variables, and a mix of continuous, categorical and discrete variables.
- It produces out of bag estimate error which has proven to be unbiased in many tests and it is relatively easy to tune with.

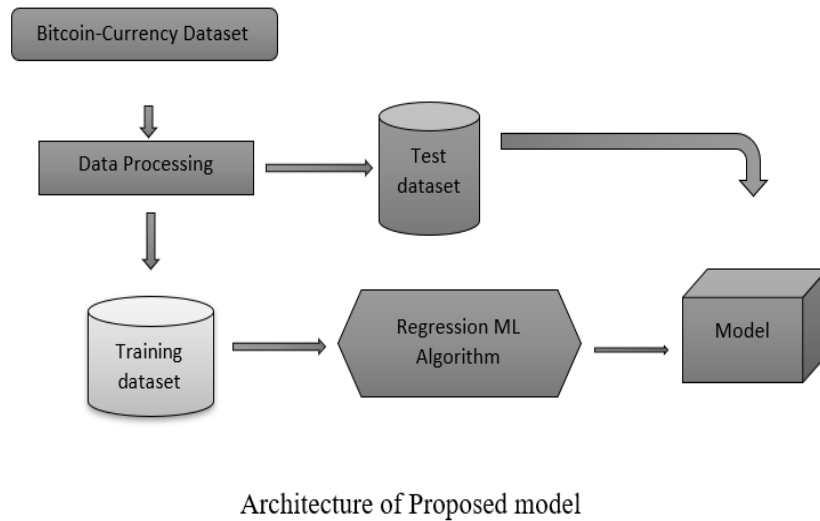


Fig:3.1:Architecture of Proposed model

3.6.5 Advantages:

- The model can be used to predict the bitcoin future. Performance metrics like accuracy, recall and precision can be calculated.
- Bitcoin future may be predicted and the investments can be made wisely.

3.8 FLOW CHART:

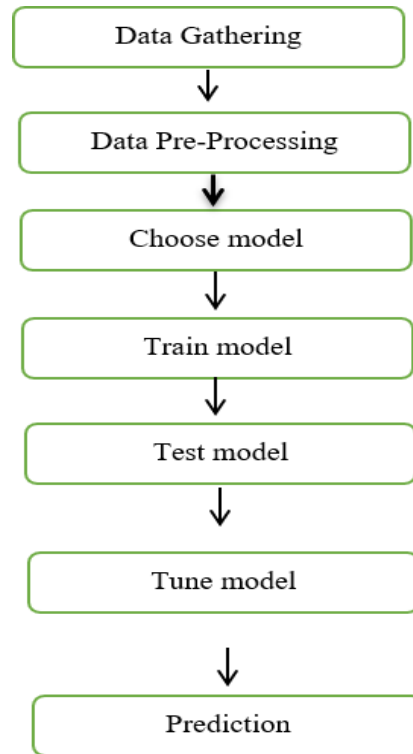


Fig:3.2: FLOW CHART

CHAPTER-4

EXPERIMENTAL OR MATERIALS AND METHODS

ALGORITHMS USED

4.1 SYSTEM STUDY:

To develop this model we use new modern technologies which are Machine Learning using Python for predicting and Flask is for user interface.

4.1.1 System requirement specifications:

a) Hardware requirements:

- Processor : Intel
- RAM : 2GB
- Hard Disk : 80GB

b) Software requirements:

- OS : Windows
- Framework : Flask
- Technology : Machine Learning using Python
- Web Browser : Chrome, Microsoft Edge
- Code editor : Visual Studio Code, Google Colab, Anaconda or Jupyter notebook.

4.2 SYSTEM SPECIFICATIONS:

4.2.1 Machine Learning Overview:

Machine learning is a field of study that looks at using computational algorithms to turn empirical data into usable models. The machine learning field grew out of traditional statistics and artificial intelligences communities. Through their business processes immense amounts of data have been and will be collected. This has provided an opportunity to re-invigorate the statistical and computational approaches to autogenerate useful models from data. Machine learning algorithms can be used to (a) gather understanding of the cyber phenomenon that produced the data under study, (b) abstract the

understanding of underlying phenomena in the form of a model, (c) predict future values of a phenomena using the above-generated model, and (d) detect anomalous behavior exhibited by a phenomenon under observation.

4.2.2 Flask Overview:

Flask is an API of Python that allows us to build up web-applications. It was developed by Armin Ronacher. Flask's framework is more explicit than Django's framework and is also easier to learn because it has less base code to implement a simple web-Application.

4.3 STEPS TO DOWNLOAD & INSTALL PYTHON:

Download the Latest version of the **Python** executable installer (<https://www.python.org/downloads/>). Watch the PIP list where pip is the package installer for python. Now upgrade the pip and setuptools using the command

Pip install --upgrade pip and Pip install --upgrade setuptools

4.3.1 IDE INSTALLATION FOR PYTHON

IDE stands for Integrated Development Environment. It is a GUI (Graphical User Interface) where programmers write their code and produce the final products. Best IDE is Pycharm. So download the pycharm new version and install the software (<https://www.jetbrains.com/pycharm/download/>)

4.3.2 PYTHON FILE CREATION

GO To FILE MENU > CREATE > NEW > PYTHON FILE
>(Name Your Python File as "BITCOIN PRICE PRIDITION" > SAVE

4.4 PYTHON LIBRARIES NEEDED

There are many libraries in python. In those we only use few main libraries needed.

4.4.1 **NUMPY LIBRARY**

NumPy is an open-source numerical Python library. NumPy contains a multi- dimensional array and matrix data structures. It can be utilized to perform a number of mathematical operations on arrays such as trigonometric, statistical, and algebraic routines like mean, mode, standard deviation etc...,

Installation- (<https://numpy.org/install/>)

```
pip install NUMPY
```

Here we mainly use array, to find mean and standard deviation.

4.4.2 **PANDAS LIBRARY**

Pandas is a high-level data manipulation tool developed by Wes McKinney. It is built on the Numpy package and its key data structure is called the DataFrame. DataFrames allow you to store and manipulate tabular data in rows of observations and columns of variables. There are several ways to create a DataFrame.

Installation- (https://pandas.pydata.org/getting_started.html)

```
pip install PANDAS
```

Here we use pandas for reading the csv files, for grouping the data, for cleaning the data using some operations.

4.4.3 **MATPLOTLIB LIBRARY**

Matplotlib is a comprehensive library for creating static, animated, and interactive visualizations in Python. Matplotlib makes easy things easy and hard things possible. Use interactive figures that can zoom, pan, update, visualize etc.,

Installation- (<https://matplotlib.org/users/installing.html>)

```
pip install Matplotlib
```

Here we use pyplot mainly for plotting graphs.

matplotlib.pyplot is a collection of functions that make matplotlib work like MATLAB. Each pyplot function makes some change to a figure: e.g., creates a figure, creates a plotting area in a figure, plots some lines in a plotting area, decorates the plot with labels, etc.

4.4.4 SEABRON LIBRARY

Seaborn package was developed based on the Matplotlib library. It is used to create more attractive and informative statistical graphics. While seaborn is a different package, it can also be used to develop the attractiveness of matplotlib graphics.

Installation-(<https://seaborn.pydata.org/installing.html>)

```
pip install Seaborn
```

4.4.5 SCIKIT-LEARN LIBRARY

Scikit-learn is a free machine learning library for the Python. It features various algorithms like support vector machine, random forests, regression and k-neighbors, and it also supports Python numerical and scientific libraries like NumPy and SciPy.

```
Pip install Scikit-Learn
```

Installation-(<https://scikit-learn.org/stable/install.html>)

Here use scikit-learn's regression methods for prediction purpose.

4.4.6 FLASK

Flask is an API of Python that allows us to build up web-applications. It was developed by Armin Ronacher. Flask's framework is more

explicit than Django's framework and is also easier to learn because it has less base code to implement a simple web-Application.

```
pip install flask
```

Here we use flask for the user-interface.

4.5 MODULES:

A modular design reduces complexity, facilitates change (a critical aspect of software maintainability), and results in easier implementation by encouraging parallel development of different part of system. Software with effective modularity is easier to develop because function may be compartmentalized and interfaces are simplified. Software architecture embodies modularity that is software is divided into separately named and addressable components called modules that are integrated to satisfy problem requirements.

Modularity is the single attribute of software that allows a program to be intellectually manageable. The five important criteria that enable us to evaluate a design method with respect to its ability to define an effective modular design are: Modular decomposability, Modular Comps ability, Modular Understand ability, Modular continuity, Modular Protection.

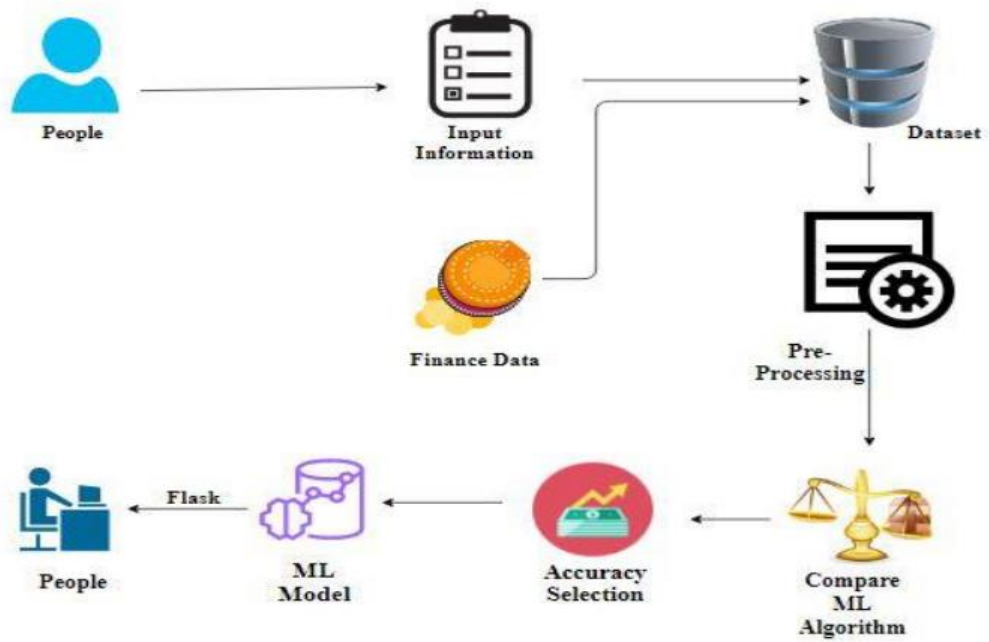


Fig:4.1: SYSTEM ARCHITECTURE

4.6 UML DIAGRAMS

4.6.1 Use Case Diagram

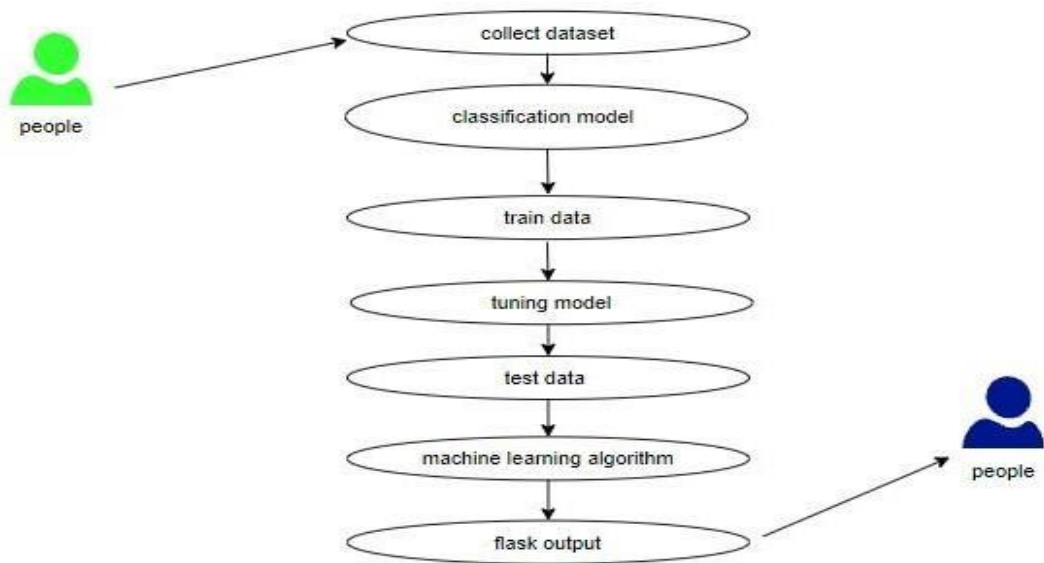


Fig:4.2: USE CASE DIAGRAM

Use case diagrams are considered for high level requirement analysis of a system. So when the requirements of a system are analyzed the functionalities are captured in use cases. So, it can say that uses cases are nothing but the system functionalities written in an organized manner.

4.6.2 Class Diagram

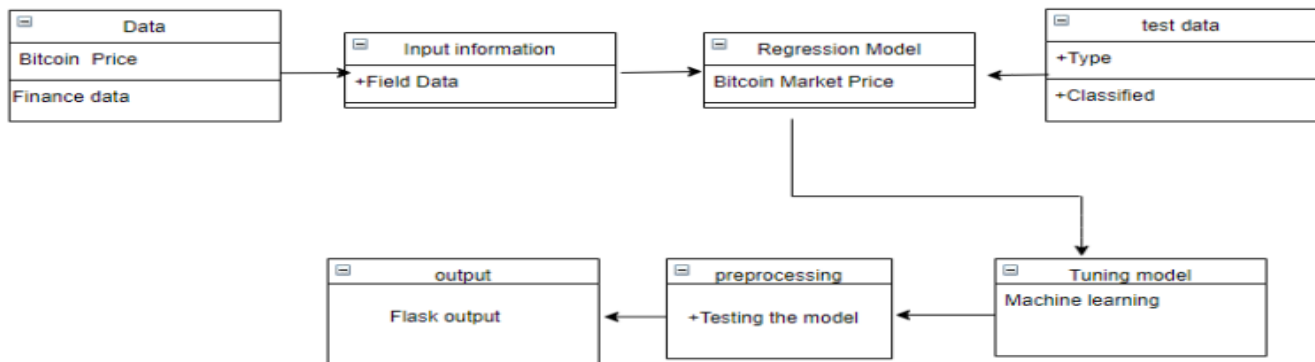


Fig:4.3: CLASS DIAGRAM

Class diagram is basically a graphical representation of the static view of the system and represents different aspects of the application. So a collection of class diagrams represent the whole system. The name of the class diagram should be meaningful to describe the aspect of the system. Each element and their relationships should be identified in advance Responsibility (attributes and methods) of each class should be clearly identified for each class minimum number of properties should be specified and because, unnecessary properties will make the diagram complicated. Use notes whenever required to describe some aspect of the diagram and at the end of the drawing it should be understandable to the developer/coder. Finally, before making the final version, the diagram should be drawn on plain paper and rework as many times as possible to make it correct.

4.6.3 Activity Diagram

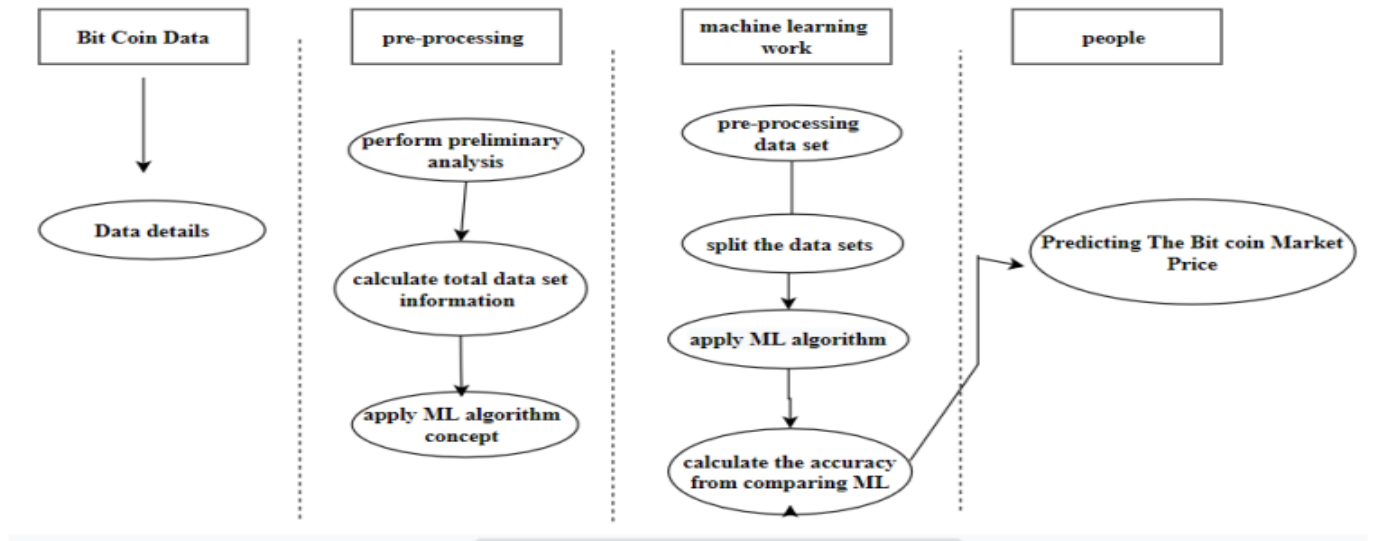


Fig:4.4: ACTIVITY DIAGRAM

Activity is a particular operation of the system. Activity diagrams are not only used for visualizing dynamic nature of a system but they are also used to construct the executable system by using forward and reverse engineering techniques. The only missing thing in activity diagram is the message part. It does not show any message flow from one activity to another. Activity diagram is some time considered as the flow chart. Although the diagrams looks like a flow chart but it is not. It shows different flow like parallel, branched, concurrent and single.

4.6.4 Sequence Diagram

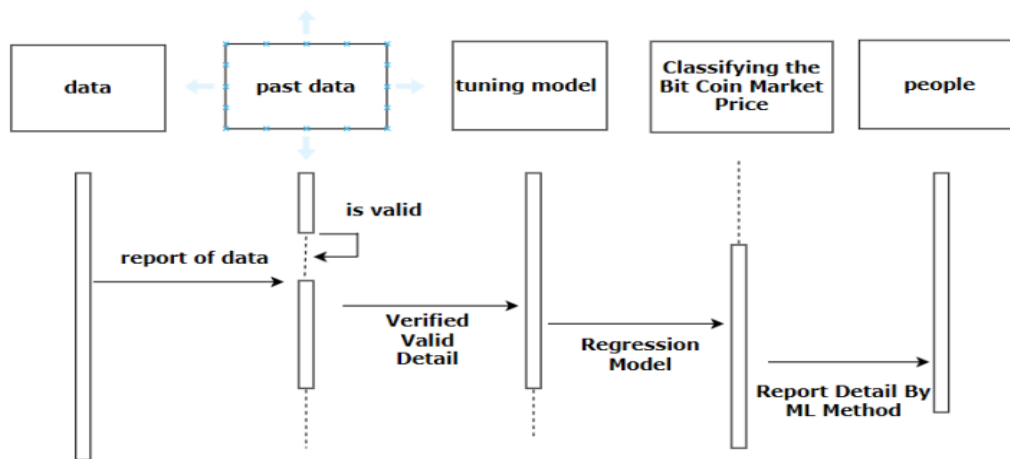


Fig:4.5: SEQUENCE DIAGRAM

Sequence diagrams model the flow of logic within your system in a visual manner, enabling you both to document and validate your logic, and are commonly used for both analysis and design purposes. Sequence diagrams are the most popular UML artifact for dynamic modeling, which focuses on identifying the behavior within your system. Other dynamic modeling techniques include activity diagramming, communication diagramming, timing diagramming, and interaction overview diagramming. Sequence diagrams, along with class diagrams and physical data models are in my opinion the most important design-level models for modern business application development.

4.6.5 Entity Relationship Diagram (ERD)

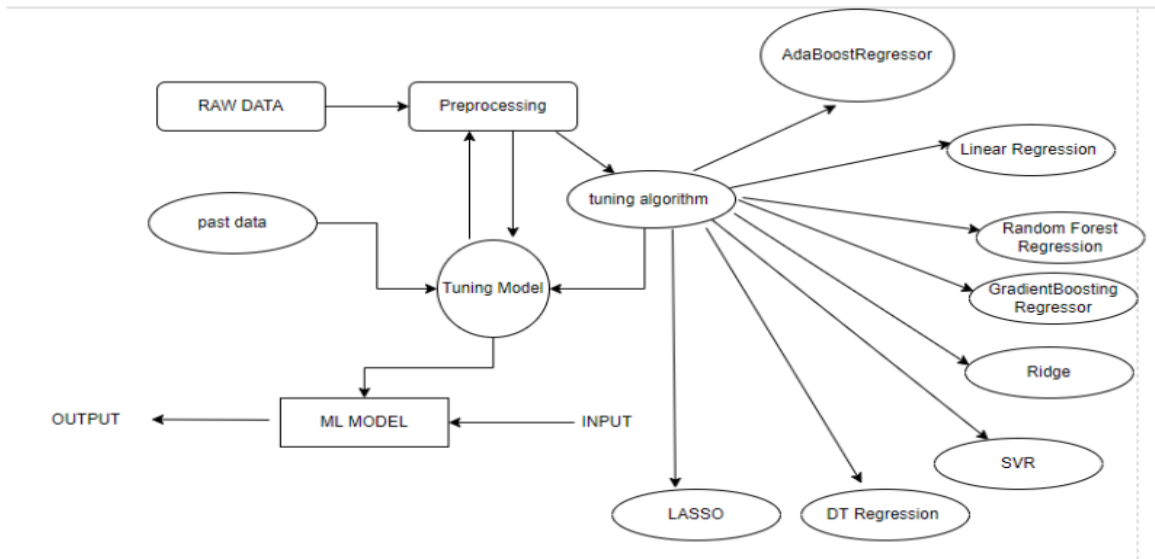


Fig:4.6: ENTITY RELATIONSHIP DIAGRAM

An entity relationship diagram (ERD), also known as an entity relationship model, is a graphical representation of an information system that depicts the relationships among people, objects, places, concepts or events within that system. An ERD is a data modelling technique that can help define business processes and be used as the foundation for a relational database. Entity relationship diagrams provide a visual starting point for database design that can also be used to help determine information system requirements throughout an organization. After a relational database is rolled out, an ERD can still serve as a referral point, should any debugging or business process re-engineering be needed later.

The following are the modules of the project, which is planned in aid to complete the project with respect to the proposed system, while overcoming existing system and also providing the support for the future enhancement.

4.7 MODULE DETAILS:

4.7.1 *Data Pre-processing*

Validation techniques in machine learning are used to get the error rate of the Machine Learning (ML) model, which can be considered as close to the true error rate of the dataset. If the data volume is large enough to be representative of the population, you may not need the validation techniques. However, in real-world scenarios, to work with samples of data that may not be a true representative of the population of given dataset. To finding the missing value, duplicate value and description of data type whether it is float variable or integer. The sample of data used to provide an unbiased evaluation of a model fit on the training dataset while tuning model hyper parameters.

The evaluation becomes more biased as skill on the validation dataset is incorporated into the model configuration. The validation set is used to evaluate a given model, but this is for frequent evaluation. It as machine learning engineers use this data to fine-tune the model hyper parameters. Data collection, data analysis, and the process of addressing data content, quality, and structure can add up to a time-consuming to-do list. During the process of data identification, it helps to understand your data and its properties; this knowledge will help you choose which algorithm to use to build your model.

A number of different data cleaning tasks using Python Pandas library and specifically, it focuses on probably the biggest data cleaning task, missing values and it able to more quickly clean data. It wants to spend less time cleaning data, and more time exploring and modeling.

Some of these sources are just simple random mistakes. Other times, there can be a deeper reason why data is missing. It's important to understand these different types of missing data from a statistics point of view. The type of missing data will influence how to deal with filling in the missing values and to detect missing values, and do some basic imputation and detailed statistical approach for dealing with missing data. Before, joint into code, it's important to understand the sources of missing data. Here are some typical reasons why data is missing:

- User forgot to fill in a field.
- Data was lost while transferring manually from a legacy database.
- There was a programming error.
- Users chose not to fill out a field tied to their beliefs about how the results would be used or interpreted.

Variable identification with Uni-variate, Bi-variate and Multi-variate analysis:

- import libraries for access and functional purpose & read the given dataset
- General Properties of Analyzing the given dataset
- Display the given dataset in the form of data frame
- show columns
- shape of the data frame
- To describe the data frame
- Checking data type and information about dataset
- Checking for duplicate data
- Checking Missing values of data frame
- Checking unique values of data frame
- Checking count values of data frame
- Rename and drop the given data frame
- To specify the type of values
- To create extra columns

4.7.2 Data Validation/ Cleaning/Preparing Process

Importing the library packages with loading given dataset. To analyzing the variable identification by data shape, data type and evaluating the missing values, duplicate values. A validation dataset is a sample of data held back from training your model that is used to give an estimate of model skill while tuning model's and procedures that you can use to make the best use of validation and test datasets when evaluating your models. Data cleaning / preparing by rename the given dataset and drop the column etc. to analyze the uni-variate, bi-variate and multi-variate process.

The steps and techniques for data cleaning will vary from dataset to dataset. The primary goal of data cleaning is to detect and remove errors and anomalies to increase the value of data in analytics and decision making.

MODULE DIAGRAM



GIVEN INPUT EXPECT OUTPUT

input: data

output: removing noisy data

4.7.3 Exploration data analysis of visualization

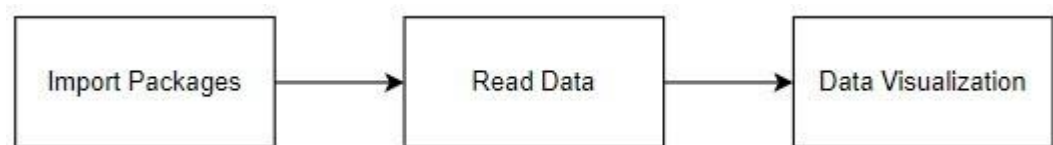
Data visualization is an important skill in applied statistics and machine learning. Statistics does indeed focus on quantitative descriptions and estimations of data. Data visualization provides an important suite of tools for gaining a qualitative understanding. This can be helpful when exploring and getting to know a dataset and can help with identifying patterns, corrupt data, outliers, and much more. With a little domain knowledge, data visualizations can be used to express and demonstrate key relationships in plots and charts that are more visceral and stakeholders than measures of association or significance. Data visualization and exploratory data analysis are whole fields themselves and it will recommend a deeper dive into some the books mentioned at the end.

Sometimes data does not make sense until it can look at in a visual form, such as with charts and plots. Being able to quickly visualize of data samples and others is an important skill both in applied statistics and in applied machine learning. It will discover the many types of plots that you will need to know when visualizing data in

Python and how to use them to better understand your own data.

- How to chart time series data with line plots and categorical quantities with bar charts.
- How to summarize data distributions with histograms and box plots.

MODULE DIAGRAM



GIVEN INPUT EXPECT OUTPUT

input: data

output: visualized data

4.7.4 Comparing Algorithm with prediction in the form of best accuracy result

It is important to compare the performance of multiple different machine learning algorithms consistently and it will discover to create a test harness to compare multiple different machine learning algorithms in Python with scikit-learn. It can use this test harness as a template on your own machine learning problems and add more and different algorithms to compare. Each model will have different performance characteristics. Using resampling methods like cross validation, you can get an estimate for how accurate each model may be on unseen data. It needs to be able to use these estimates to choose one or two best models from the suite of models that you have created. When have a new dataset, it is a good idea to visualize the data using different techniques in order to look at the data from different perspectives. The same idea applies to model selection. You should use a number of different ways of looking at the estimated accuracy of your machine learning algorithms in order to choose the one or two to finalize. A way to do this is

to use different visualization methods to show the average accuracy, variance and other properties of the distribution of model accuracies.

In the next section you will discover exactly how you can do that in Python with scikit-learn. The key to a fair comparison of machine learning algorithms is ensuring that each algorithm is evaluated in the same way on the same data and it can achieve this by forcing each algorithm to be evaluated on a consistent test harness.

Pre-processing refers to the transformations applied to our data before feeding it to the algorithm. Data Preprocessing is a technique that is used to convert the raw data into a clean data set. In other words, whenever the data is gathered from different sources it is collected in raw format which is not feasible for the analysis. To achieving better results from the applied model in Machine Learning method of the data has to be in a proper manner. Some specified Machine Learning model needs information in a specified format, for example, Random Forest algorithm does not support null values. Therefore, to execute random forest algorithm null values have to be managed from the original raw data set. And another aspect is that data set should be formatted in such a way that more than one Machine Learning and Deep Learning algorithms are executed in given dataset.

In the example below these 7 different algorithms are compared:

- Random Forest
- Decision Tree Classifier
- Linear Regression
- Support Vector Classifier
- Lasso Regression
- Gradient Boosting

The K-fold cross validation procedure is used to evaluate each algorithm, importantly configured with the same random seed to ensure that the same splits to the training data are performed and that each algorithm is evaluated in precisely

the same way. Before that comparing algorithm, Building a Machine Learning Model using install Scikit-Learn libraries. In this library package have to done preprocessing, linear model with logistic regression method, cross validating by KFold method, ensemble with random forest method and tree with decision tree classifier. Additionally, splitting the train set and test set. To predicting the result by comparing accuracy.

False Positives (FP): A person who will pay predicted as defaulter. When actual class is no and predicted class is yes. E.g. if actual class says this passenger did not survive but predicted class tells you that this passenger will survive.

False Negatives (FN): A person who default predicted as payer. When actual class is yes but predicted class in no. E.g. if actual class value indicates that this passenger survived and predicted class tells you that passenger will die.

True Positives (TP): A person who will not pay predicted as defaulter. These are the correctly predicted positive values which means that the value of actual class is yes and the value of predicted class is also yes. E.g. if actual class value indicates that this passenger survived and predicted class tells you the same thing.

True Negatives (TN): A person who default predicted as payer. These are the correctly predicted negative values which means that the value of actual class is no and value of predicted class is also no. E.g. if actual class says this passenger did not survive and predicted class tells you the same thing.

Prediction result by accuracy:

Logistic regression algorithm also uses a linear equation with independent predictors to predict a value. The predicted value can be anywhere between negative infinity to positive infinity. It needs the output of the algorithm to be classified variable data. Higher accuracy predicting result is logistic regression model by comparing the best accuracy.

True Positive Rate (TPR) = $TP / (TP + FN)$

False Positive Rate (FPR) = $FP / (FP + TN)$

Accuracy: The Proportion of the total number of predictions that is correct otherwise overall how often the model predicts correctly defaulters and non-defaulters.

Accuracy calculation:

Accuracy = $(TP + TN) / (TP + TN + FP + FN)$

Accuracy is the most intuitive performance measure and it is simply a ratio of correctly predicted observation to the total observations. One may think that, if we have high accuracy then our model is best. Yes, accuracy is a great measure but only when you have symmetric datasets where values of false positive and false negatives are almost same.

Precision: The proportion of positive predictions that are actually correct.

Precision = $TP / (TP + FP)$

Precision is the ratio of correctly predicted positive observations to the total predicted positive observations. The question that this metric answer is of all passengers that labeled as survived, how many actually survived? High precision relates to the low false positive rate. We have got 0.788 precision which is pretty good.

Recall: The proportion of positive observed values correctly predicted. (The proportion of actual defaulters that the model will correctly predict)

Recall = $TP / (TP + FN)$

Recall (Sensitivity) - Recall is the ratio of correctly predicted positive observations to the all observations in actual class - yes.

F1 Score is the weighted average of Precision and Recall. Therefore, this score takes both false positives and false negatives into account. Intuitively it is not as

easy to understand as accuracy, but F1 is usually more useful than accuracy, especially if you have an uneven class distribution. Accuracy works best if false positives and false negatives have similar cost. If the cost of false positives and false negatives are very different, it's better to look at both Precision and Recall.

General Formula:

$$F\text{- Measure} = 2TP / (2TP + FP + FN)$$

F1-Score Formula:

$$F1 \text{ Score} = 2 * (\text{Recall} * \text{Precision}) / (\text{Recall} + \text{Precision})$$

4.7.5 ALGORITHM AND TECHNIQUES

Algorithm Explanation

In machine learning and statistics, classification is a supervised learning approach in which the computer program learns from the data input given to it and then uses this learning to classify new observation. This data set may simply be bi-class (like identifying whether the person is male or female or that the mail is spam or non-spam) or it may be multi-class too. Some examples of classification problems are: speech recognition, handwriting recognition, bio metric identification, document classification etc. In Supervised Learning, algorithms learn from labeled data. After understanding the data, the algorithm determines which label should be given to new data based on pattern and associating the patterns to the unlabeled new data.

Used Python Packages:

sklearn:

- In python, sklearn is a machine learning package which include a lot of ML algorithms.
- Here, we are using some of its modules like `train_test_split`, `DecisionTreeClassifier` or `Logistic Regression` and `accuracy_score`.

NumPy:

- It is a numeric python module which provides fast maths functions for calculations.
- It is used to read data in numpy arrays and for manipulation purpose.

Pandas:

- Used to read and write different files.
- Data manipulation can be done easily with data frames.

Matplotlib:

- Data visualization is a useful way to help with identify the patterns from given dataset.
- Data manipulation can be done easily with data frames.

Linear Regression:

Linear Regression is a machine learning algorithm based on supervise learning. Linear regression performs the task to predict a dependent variable value (y) based on a given independent variable (x). So, this regression technique finds out a linear relationship between x (input) and y(output).failure. It performs a regression task. Regression models a target prediction value based on independent variables. It is mostly used for finding out the relationship between variables and forecasting. Different regression models differ based on the kind of relationship between dependent and independent variables, they are considering and the number of independent variables being usedLinear regression performs the task to predict a dependent variable value (y) based on a given independent variable (x). So, this regression technique finds out a linear relationship between x (input) and y(output). Hence, the name is Linear Regression. In the figure above, X (input) is the work experience and Y (output) is the work experience and Y (output) is the salary of a person. The regression line is the best fit line for our model. Linear regression is one of the easiest and most popular Machine Learning algorithms. It is a statistical method that is used for predictive analysis. Linear regression makes predictions for continuous/real

or numeric variables such as sales, salary, age, product price, etc. Linear regression algorithm shows a linear relationship between a dependent (y) and one or more independent (x) variables, hence called as linear regression. Since linear regression shows the linear relationship, which means it finds how the value of the dependent variable is changing according to the value of the independent variable. The linear regression model provides a sloped straight line representing the relationship between the variables.

```
• MEAN ABSOLUTE ERROR VALUE IS : 1067.2353549011561  
  
MEAN SQUARED ERROR VALUE IS : 1902495.778078313  
  
MEDIAN ABSOLUTE ERROR VALUE IS : 918.413006370537  
  
ACCURACY RESULT OF LINEAR REGRESSION IS : 87.15064247291966  
  
R2_SCORE VALUE IS : 0.8707448052684265
```

MODULE DIAGRAM

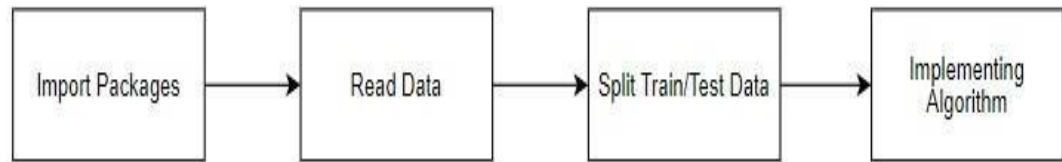


Fig:4.7: LOGISTIC REGRESSION

GIVEN INPUT EXPECT OUTPUT

input: data

output: getting accuracy

Random Forest Classifier:

Random forests or random decision forests are an ensemble learning method for classification, regression and other tasks, that operate by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees. Random decision forests correct for decision trees' habit of over fitting to their training set. Random forest is a type of supervised machine learning algorithm based on ensemble learning. Ensemble learning is a type of learning where you join different types of algorithms or same algorithm multiple times to form a more powerful prediction model. The random forest algorithm combines multiple algorithm of the same type i.e. multiple decision *trees*, resulting in a *forest of trees*, hence the name "Random Forest". The random forest algorithm can be used for both regression and classification tasks.

The following are the basic steps involved in performing the random forest algorithm:

- Pick N random records from the dataset.
- Build a decision tree based on these N records.
- Choose the number of trees you want in your algorithm and repeat steps 1 and 2.

In case of a regression problem, for a new record, each tree in the forest predicts a value for Y (output). The final value can be calculated by taking the average of all the values predicted by all the trees in forest. Or, in case of a classification problem, each tree in the forest predicts the category to which the new record belongs. Finally, the new record is assigned to the category that wins the majority vote.

```
... MEAN ABSOLUTE ERROR VALUE IS : 286.4136750830562

MEAN SQUARED ERROR VALUE IS : 220700.95360825915

MEDIAN ABSOLUTE ERROR VALUE IS : 138.39609999999948

ACCURACY RESULT OF RANDOM FOREST REGRESSOR IS : 98.40131103640768

R2_SCORE VALUE IS : 0.9840131103507547
```

Fig:4.8: RANDOM FOREST CLASSIFIER

GIVEN INPUT EXPECT OUTPUT

input: data

output: getting accuracy

Decision Tree Classifier:

It is one of the most powerful and popular algorithm. Decision-tree algorithm falls under the category of supervised learning algorithms. It works for both continuous as well as categorical output variables. Assumptions of Decision tree:

- At the beginning, we consider the whole training set as the root.
- Attributes are assumed to be categorical for information gain, attributes are assumed to be continuous.
- On the basis of attribute values records are distributed recursively.
- We use statistical methods for ordering attributes as root or internal node.

Decision tree builds classification or regression models in the form of a tree structure. It breaks down a data set into smaller and smaller subsets while at the

same time an associated decision tree is incrementally developed. A decision node has two or more branches and a leaf node represents a classification or decision. The topmost decision node in a tree which corresponds to the best predictor called root node. Decision trees can handle both categorical and numerical data. Decision tree builds classification or regression models in the form of a tree structure. It utilizes an if-then rule set which is mutually exclusive and exhaustive for classification. The rules are learned sequentially using the training data one at a time. Each time a rule is learned, the tuples covered by the rules are removed. This process is continued on the training set until meeting a termination condition. It is constructed in a top-down recursive divide-and-conquer manner. All the attributes should be categorical. Otherwise, they should be discretized in advance. Attributes in the top of the tree have more impact towards in the classification and they are identified using the information gain concept. A decision tree can be easily over-fitted generating too many branches and may reflect anomalies due to noise or outliers.

MEAN ABSOLUTE ERROR VALUE IS : 382.8996677740864

MEAN SQUARED ERROR VALUE IS : 505679.82296345517

MEDIAN ABSOLUTE ERROR VALUE IS : 127.29999999999927

ACCURACY RESULT OF DECISION TREE REGRESSOR IS : 96.33746358788717

R2_SCORE VALUE IS : 0.9633701286949756

Lasso Regression:

Lasso regression is a regularization technique. It is used over regression methods for a more accurate prediction. This model uses shrinkage. Shrinkage is where data values are shrunk towards a central point as the mean. The lasso procedure encourages simple, sparse models (i.e. models with fewer parameters). The word “LASSO” stands for Least Absolute Shrinkage and Selection Operator. It is a statistical formula for the regularisation of data models and feature selection. Lasso regression is a regularization technique. It is used over regression methods for a more accurate prediction. This model uses shrinkage. Shrinkage is where data values are shrunk towards a central point as the mean. The lasso procedure encourages simple, sparse models (i.e. models with fewer parameters). This particular type of regression is well-suited for models showing high levels of multicollinearity or when you want to automate certain parts of model selection, like variable selection/parameter elimination. Lasso Regression uses L1 regularization technique (will be discussed later in this article). It is used when we have more number of features because it automatically performs feature selection. The L1 regularization performed by Lasso, causes the regression coefficient of the less contributing variable to shrink to zero or near zero. Lasso regression uses shrinkage, where the data values are shrunk towards a central point such as the mean value. The Lasso penalty shrinks or reduces the coefficient value towards zero. The less contributing variable is therefore allowed to have a zero or near-zero coefficient.

MEAN ABSOLUTE ERROR VALUE IS : 1075.5567469858108

MEAN SQUARED ERROR VALUE IS : 1953393.127966158

MEDIAN ABSOLUTE ERROR VALUE IS : 895.3240026988433

ACCURACY RESULT OF LASSO REGRESSION IS : 87.52728082813397

R2_SCORE VALUE IS : 0.8739646288984018

Support Vector Classifier

Support Vector Classifier or SVC is one of the most popular Supervised Learning algorithms, which is used for Classification as well as Regression problems. Primarily, it is used for Classification problems in Machine Learning. The goal of the SVC algorithm is to create the best line or decision boundary that can segregate n-dimensional space into classes. So that we can easily put the new data point in the correct category in the future. This best decision boundary is called a hyperplane. SVC chooses the extreme points/vectors that help in creating the hyperplane. These extreme cases are called as support vectors, and hence algorithm is termed as Support Vector Machine.

MEAN ABSOLUTE ERROR VALUE IS : 3298.218492879275

MEAN SQUARED ERROR VALUE IS : 15338696.058602957

MEDIAN ABSOLUTE ERROR VALUE IS : 2988.528175301004

ACCURACY RESULT OF SVR IS : 1.183743451702468

R2_SCORE VALUE IS : 0.010328119678867376

Fig:4.12: SUPPORT VECTOR CLASSIFIER

Gradient Boost

Gradient boosting is a machine learning technique used in regression and classification tasks, among others. It gives a prediction model in the form of an ensemble of weak prediction models, which are typically decision trees. When a decision tree is the weak learner, the resulting algorithm is called gradient-boosted trees; it usually outperforms random forest.[1][2][3] A gradient-boosted trees model is built in a stage-wise fashion as in other boosting methods, but it generalizes the other methods by allowing optimization of an arbitrary differentiable loss function

Generic gradient boosting at the m -th step would fit a decision tree to pseudo-residuals. Let L be the number of its leaves. The tree partitions the input space into disjoint regions and predicts a constant value in each region. Using the indicator notation, the output of for input x can be written as the sum.

```
MEAN ABSOLUTE ERROR VALUE IS : 373.78721511434054  
  
MEAN SQUARED ERROR VALUE IS : 341418.86746489047  
  
MEDIAN ABSOLUTE ERROR VALUE IS : 203.53474726439526  
  
ACCURACY RESULT OF GradientBoostingRegressor IS : 97.69316837085633  
  
R2_SCORE VALUE IS : 0.9768040682624889
```

Fig:4.13: GRADIENT BOOST

4.7.6 Deployment Using Flask (Web Framework):

Flask is a micro web framework written in Python. It is classified as a micro-framework because it does not require particular tools or libraries. It has no database abstraction layer, form validation, or any other components where pre-existing third-party libraries provide common functions. However, Flask supports extensions that can add application features as if they were implemented in Flask itself.

Extensions exist for object-relational mappers, form validation, upload handling, various open authentication technologies and several common framework related tools.

Flask was created by Armin Ronacher of Pocoo, an international group of Python

enthusiasts formed in 2004. According to Ronacher, the idea was originally an April Fool's joke that was popular enough to make into a serious application. The name is a play on the earlier Bottle framework.

When Ronacher and Georg Brand created a bulletin board system written in Python, the Pocoo projects Werkzeug and Jinja were developed.

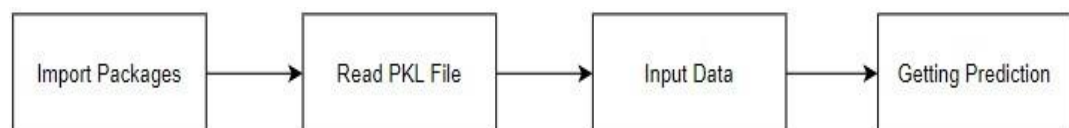
In April 2016, the Pocoo team was disbanded and development of Flask and related libraries passed to the newly formed Pallets project.

Flask has become popular among Python enthusiasts. As of October 2020, it has second most stars on GitHub among Python web-development frameworks, only slightly behind Django, and was voted the most popular web framework in the Python Developers Survey 2018.

The micro-framework Flask is part of the Pallets Projects, and based on several others of them.

Flask is based on Werkzeug, Jinja2 and inspired by Sinatra Ruby framework, available under BSD licence. It was developed at pocoo by Armin Ronacher. Although Flask is rather young compared to most Python frameworks, it holds a great promise and has already gained popularity among Python web developers. Let's take a closer look into Flask, so-called "micro" framework for Python.

MODULE DIAGRAM



GIVEN INPUT EXPECTED OUTPUT

input : data values

output : predicting output

CHAPTER-5

RESULTS AND DISCUSSION, PERFORMANCE ANALYSIS

5.1 PERFORMANCE ANALYSIS:

Website performance optimization, the focal point of technologically superior website designs is the primary factor dictating success for Loan approval process. After all, unimpressive website performance kills admission process when the torture of waiting for slow Web pages to load frustrates visitors into seeking alternatives – impatience is a digital virtue! And also the ml algorithms used in our project will give the best accurate result to the user for Loan approval.

We created the following six chapter in-depth speed optimization guide to show you how important it is to have a fast loading, snappy website! Countless research papers and benchmarks prove that optimizing your sites' speed is one of the most affordable and highest ROI providing investments!

Lightning-fast page load speed amplifies visitor engagement, retention, and boosts sales. Instantaneous website response leads to higher conversion rates, and every 1 second delay in page load decreases customer satisfaction by 16 percent, page views by 11 percent and conversion rates by 7 percent according to a recent Aberdeen Group research.

5.2 DISCUSSION:

The utilisation of BitConduite's entity-based analysis, as well as the necessity of exploratory analytical approaches, were highlighted during the preliminary evaluation. Despite this, we discovered two key flaws in our technique and the entity clustering mechanism we use. Approach. Because our exploratory method demands variable activity measurement calculation, scalability is limited. While we made every effort to reduce data processing times, the most significant bottleneck is on-the-fly clustering, which can take minutes in the worst-case situation (if many of entities are clustered at a time). In this case, a method of progressive clustering [19] would be useful. The process of combining entities is referred to as aggregation of entities.

Although entity-based analysis is important, there is some uncertainty about entities because there are no methods for evaluating the quality of entity aggregation. Mixing services (also called as tumblers) make address aggregation more difficult by hiding the connections between addresses for privacy reasons. Furthermore, entity aggregation is computationally expensive and, in our case, memory intensive.

CHAPTER-6

SUMMARY AND CONCLUSION

6.1 SUMMARY:

This project objective is to predict the Loan Approval of the user. So this online banking loan approval system will reduce the paper work and reduce the wastage of bank asserts and efforts and also saves the valuable time of the customer.

6.2 CONCLUSION:

Data cleansing and processing were the first steps in the study, followed by missing value detection, exploratory analysis, and model construction and evaluation. A better accuracy score on a public test set will be discovered for the highest accuracy. The BITCOIN Market Price may be found with the aid of this software. Data cleansing and processing were the first steps in the study, followed by missing value detection, exploratory analysis, and model construction and evaluation. A better accuracy score on a public test set will be discovered for the highest accuracy. The BITCOIN Market Price may be found with the aid of this software.

6.3 FUTURE WORK:

- Bitcoin Market Price prediction to connect with AI model.
- To make this procedure more efficient, you may use a web application or a desktop programme to display the prediction result.
- To reduce the amount of time and effort required to implement in an AI system.environment.

REFERENCES:

- [1] .Andrade de Oliveira, L. Enrique ZÃ¡rate and M. de Azevedo Reis; C. NeriNobre, —The use of artificial neural networks in the analysis and prediction of stock prices,II in IEEE International Conference on Systems, Man, and Cybernetics, 2011, pp. 2151-2155.
- [2] Application of Random Forest Algorithm on Feature Subset Selection and Classification and RegressionII ;Jitendra Kumar Jaiswal, Rita Samikannu IEEE 2017.
- [3] Gulani Senthuran, Malka N. Halgamuge, " Prediction of Crptocurrency Market Price Using Deep Learning and Blockchain Information: Bitcoin and Ethereum".M. A. Sheikh, A. K. Goel and T. Kumar, "An Approach for Prediction of Loan Approval using Machine Learning Algorithm," 2020 International Conference onElectronics and Sustainable Communication Systems (ICESC), 2020, pp. 490- 494, doi: 10.1109/ICESC48915.2020.9155614.
- [4] .Huisu Jang and Jaewook Lee, —An Empirical Study on Modelling and Prediction of Bitcoin Prices with Bayesian Neural Networks based on Blockchain Information,II in IEEE Early Access Articles, 2017, vol. 99, pp. 1-1. [3] F.
- [5] .Lekkala Sreekanth Reddy,Dr.P.Sriramya, "A Research On Bitcoin Price Prediction Using Machine Learning Algorithms".
- [6] Mr. Shivam Pandey1, Mr.Anil Chavan2, miss. Dhanashree Paraskar3, Prof. Sareen Deore"Bitcoin Price Prediction using Machine Learning"
- [7] Variable selection using the Lasso-Cox model with Bayesian regularization, Wenxin Lu ; Zhuliang Yu ; ZhenghuiGu ; Jinhong Huang ; Wei Gao ; Haiyu Zhou
- [8] Yang Li 1,2 , Zibin Zheng 1,2,* and Hong-Ning Dai, " Enhancing Bitcoin Price Fluctuation Prediction Using Attentive LSTM and Embedding Network".Vincenzo Moscato, Antonio Picariello, Giancarlo Sperl , A benchmark of machine learning approaches for credit score prediction, Expert Systems with Applications, Volume 165, 2021, 113986, ISSN 0957-4174.

```
Q1 = df.quantile(0.25)
Q3 = df.quantile(0.75)
IQR = Q3 - Q1
```

APPENDIX:

A. SOURCE CODE:

jupyter Notebook

```
import warnings
warnings.filterwarnings('ignore')
#import library packages
import pandas as p
import numpy as n
#Load given dataset
data = p.read_csv("Bitcoin.csv")
data.head()
data.shape
df=data.dropna()
df.head()
df.shape
#show columns
df.columns
```

```
df.rename(columns={"BTC price [USD]":"BTC_price_USD"},inplace=True)
df.rename(columns={"n-transactions":"n_transactions"},inplace=True)
df.rename(columns={"fee [USD]":"fee_USD"},inplace=True)
df.rename(columns={"btc search trends":"btc_search_trends"},inplace=True)
df.rename(columns={"Gold price[USD]":"Gold_price_USD"},inplace=True)
df.rename(columns={"SP500 close index":"SP500_close_index"},inplace=True)
df.rename(columns={"Oil WTI price[USD]":"Oil_WTI_price_USD"},inplace=True)
df.rename(columns={"M2(Not seasonally adjusted)[1e+09 USD]":"M2_money_supply_USA"},inplace=True)
df.columns
```

```
df['n_transactions'].unique()
df['fee_USD'].unique()
df["btc_search_trends"].unique()
df["Gold_price_USD"].unique()
df["SP500_close_index"].unique()
df["Oil_WTI_price_USD"].unique()
df["M2_money_supply_USA"].unique()
#To describe the dataframe
df.describe()
```

```
#Checking datatype and information about dataset
df.info()
#Checking for duplicate data
df.duplicated()
#find sum of duplicate data
sum(df.duplicated())
```

```
#Checking sum of missing values
df.isnull().sum()
df.columns
df['BTC_price_USD'].value_counts()print("Minimum BTC_price_USD value
is:",
```

```
df.BTC_price_USD.min())
print("Maximum BTC_price_USD value is:", df.BTC_price_USD.max())
p.Categorical(df['btc_search_trends']).describe()
p.Categorical(df['BTC_price_USD']).describe()
df.corr()
df.head()
df.columns
df.head(20)
```

```
import warnings
warnings.filterwarnings('ignore')
#import library packages
import pandas as p
import matplotlib.pyplot as plt
import seaborn as s
import seaborn as s
import numpy as n
data = p.read_csv("Bitcoin
```

```
df.rename(columns={"BTC price [USD]":"BTC_price_USD"},inplace=True)
df.rename(columns={"n-transactions":"n_transactions"},inplace=True)
df.rename(columns={"fee [USD]":"fee_USD"},inplace=True)
df.rename(columns={"btc search trends":"btc_search_trends"},inplace=True)
df.rename(columns={"Gold price[USD]":"Gold_price_USD"},inplace=True)
df.rename(columns={"SP500 close index":"SP500_close_index"},inplace=True)
df.rename(columns={"Oil WTI price[USD]":"Oil_WTI_price_USD"},inplace=True)
df.rename(columns={"M2(Not seasonally adjusted)[1e+09
USD]":"M2_money_supply_USA"},inplace=True)
```

```
df.columns
#Histogram Plot
df['n_transactions'].hist(figsize=(7,6), color='yellow', alpha=0.7)
plt.xlabel('n_transactions')
plt.ylabel('No of n_transactions')
plt.title('Number of n_transactions')
#Histogram Plot
df['fee_USD'].hist(figsize=(7,6), color='pink', alpha=0.7)
plt.xlabel('fee_USD')
plt.ylabel('No of fee_USD')
```

```

#Histogram Plot
df['btc_search_trends'].hist(figsize=(7,6), color='red', alpha=0.7)
plt.xlabel('btc_search_trends')
plt.ylabel('No of btc_search_trends')
plt.title('Number of btc_search_trends')
#Histogram Plot
df['Gold_price_USD'].hist(figsize=(7,6), color='blue', alpha=0.7)
plt.xlabel('Gold_price_USD')
plt.ylabel('No of Gold_price_USD')
plt.title('Number of Gold_price_USD')
#Histogram Plot
df['SP500_close_index'].hist(figsize=(7,6), color='black', alpha=0.7)
plt.xlabel('SP500_close_index')
plt.ylabel('No of SP500_close_index')
plt.title('Number of SP500_close_index')
#Propagation by variable
def PropByVar(df, variable):
    dataframe_pie = df[variable].value_counts()
    ax = dataframe_pie.plot.pie(figsize=(10,10), autopct='%1.2f%%', fontsize = 12)
    ax.set_title(variable + ' \n', fontsize = 15)
    return n.round(dataframe_pie/df.shape[0]*100,2)
PropByVar(df, 'btc_search_trends')
plt.boxplot(df['n_transactions'])
plt.show()
plt.boxplot(df['fee_USD'])
plt.show()plt.boxplot(df['SP500_close_index'])
plt.show()
# Heatmap plot diagram
fig, ax = plt.subplots(figsize=(15,10))
s.heatmap(df.corr(), annot=True)
fig, ax = plt.subplots(figsize=(16,8))
ax.scatter(df['btc_search_trends'],df['BTC_price_USD'])
ax.set_xlabel('btc_search_trends')
ax.set_ylabel('BTC_price_USD')
plt.show()
df.columns
from sklearn.preprocessing import LabelEncoder
var_mod = ['Date', 'BTC_price_USD', 'n_transactions', 'fee_USD',
           'btc_search_trends', 'Gold_price_USD', 'SP500_close_index',
           'Oil_WTI_price_USD', 'M2_money_supply_USA']
le = LabelEncoder()
for i in var_mod:
    df[i] = le.fit_transform(df[i]).astype(int)
df.head()
#preprocessing, split test and dataset, split response variable
X = df.drop(labels='BTC_price_USD', axis=1)
#Response variable
y = df.loc[:, 'BTC_price_USD']
#We'll use a test size of 30%. We also stratify the split on the response
variable, which is very important to do because there are so few fraudulent
transactions.
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3)

```

```

#import library packages
import pandas as p
import matplotlib.pyplot as plt
import seaborn as s
import numpy as n
#Load given dataset
data = p.read_csv("Bitcoin.csv")
import warnings
warnings.filterwarnings('ignore')
df=data.dropna()
df.head()
df.columns
df.rename(columns={"BTC price [USD]":"BTC_price_USD"},inplace=True)
df.rename(columns={"n-transactions":"n_transactions"},inplace=True)
df.rename(columns={"fee [USD]":"fee_USD"},inplace=True)
df.rename(columns={"btc search trends":"btc_search_trends"},inplace=True)
df.rename(columns={"Gold price[USD]":"Gold_price_USD"},inplace=True)
df.rename(columns={"SP500 close index":"SP500_close_index"},inplace=True)
df.rename(columns={"Oil WTI price[USD]":"Oil_WTI_price_USD"},inplace=True)
df.rename(columns={"M2(Not seasonally adjusted)[1e+09
USD]":"M2_money_supply_USA"},inplace=True)
df.info()
del df['Date']
df.columns
df.head()
from sklearn.metrics import
mean_absolute_error,mean_squared_error,r2_score,explained_variance_score,median_abso
lute_error
X = df.drop(labels='BTC_price_USD', axis=1)
#Response variable
y = df.loc[:, 'BTC_price_USD']
#We'll use a test size of 30%. We also stratify the split on the response variable,
which is very important to do because there are so few fraudulent transactions.
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3)

```

Linear regression

```

from sklearn.linear_model import LinearRegression
linR= LinearRegression()
linR.fit(X_train,y_train)
predictL = linR.predict(X_test)
MAE= (mean_absolute_error(y_test,predictL))
print('MEAN ABSOLUTE ERROR VALUE IS :',MAE)
print(" ")
MSE=(mean_squared_error(y_test,predictL))
print('MEAN SQUARED ERROR VALUE IS :',MSE)
print(" ")
MedianAE=(median_absolute_error(y_test,predictL))
print('MEDIAN ABSOLUTE ERROR VALUE IS :',MedianAE)
print(" ")
EVS=(explained_variance_score(y_test,predictL)*100)
print('ACCURACY RESULT OF LINEAR REGRESSION IS :',EVS)
print(" ")
R2=(r2_score(y_test,predictL))
print('R2 SCORE VALUE IS :',R2)

```

Gradient boosting regressor

```
from sklearn.ensemble import GradientBoostingRegressor
gb= GradientBoostingRegressor()
gb.fit(X_train,y_train)

predictgb = gb.predict(X_test)

MAE= (mean_absolute_error(y_test,predictgb))
print('MEAN ABSOLUTE ERROR VALUE IS :',MAE)
print(" ")

MSE=(mean_squared_error(y_test,predictgb))
print('MEAN SQUARED ERROR VALUE IS :',MSE)
print(" ")

MedianAE=(median_absolute_error(y_test,predictgb))
print('MEDIAN ABSOLUTE ERROR VALUE IS :',MedianAE)
print(" ")

EVS=(explained_variance_score(y_test,predictgb)*100)
print('ACCURACY RESULT OF GradientBoostingRegressor IS :',EVS)
print(" ")

R2=(r2_score(y_test,predictgb))
print('R2_SCORE VALUE IS :',R2)
print(" ")
from sklearn.linear_model import LinearRegression
reg= LinearRegression()
reg.fit(X_train, y_train)
print('Coefficients: \n', reg.coef_)
print('Variance score: {}'.format(reg.score(X_test, y_test)))
plt.style.use('fivethirtyeight')
plt.scatter(reg.predict(X_train), reg.predict(X_train) - y_train,
            color = "green", s = 10, label = 'Train data')
plt.scatter(reg.predict(X_test), reg.predict(X_test) - y_test,
            color = "blue", s = 10, label = 'Test data')
plt.hlines(y = 0, xmin = 0, xmax = 50, linewidth = 2)
plt.legend(loc = 'upper right')
plt.title("Residual errors")
plt.show()
```

```
#import library packages
import pandas as p
import matplotlib.pyplot as plt
import seaborn as s
import numpy as n
#Load given dataset
data = p.read_csv("Bitcoin.csv")
import warnings
warnings.filterwarnings('ignore')
df=data.dropna()
df.head()
```

```

df.rename(columns={"BTC price [USD]":"BTC_price_USD"},inplace=True)
df.rename(columns={"n-transactions":"n_transactions"},inplace=True)
df.rename(columns={"fee [USD]":"fee_USD"},inplace=True)
df.rename(columns={"btc search trends":"btc_search_trends"},inplace=True)
df.rename(columns={"Gold price[USD]":"Gold_price_USD"},inplace=True)
df.rename(columns={"SP500 close index":"SP500_close_index"},inplace=True)
df.rename(columns={"Oil WTI price[USD]":"Oil_WTI_price_USD"},inplace=True)
df.rename(columns={"M2(Not seasonally adjusted)[1e+09
USD]":"M2_money_supply_USA"},inplace=True)
df.info()
del df['Date']
df.head()
df.tail()
from sklearn.metrics import
mean_absolute_error,mean_squared_error,r2_score,explained_variance_score,median
_absolute_error
X = df.drop(labels='BTC_price_USD', axis=1)
#Response variable
y = df.loc[:, 'BTC_price_USD']
#We'll use a test size of 30%. We also stratify the split on the response
variable, which is very important to do because there are so few fraudulent
transactions.
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3)

```

Random forest regression

```

from sklearn.ensemble import RandomForestRegressor
RF= RandomForestRegressor()
RF.fit(X_train,y_train)

predictR = RF.predict(X_test)

MAE= (mean_absolute_error(y_test,predictR))
print('MEAN ABSOLUTE ERROR VALUE IS :',MAE)
print(" ")

MSE=(mean_squared_error(y_test,predictR))
print('MEAN SQUARED ERROR VALUE IS :',MSE)
print(" ")

MedianAE=(median_absolute_error(y_test,predictR))
print('MEDIAN ABSOLUTE ERROR VALUE IS :',MedianAE)
print(" ")

EVS=(explained_variance_score(y_test,predictR)*100)
print('ACCURACY RESULT OF RANDOM FOREST REGRESSOR IS :',EVS)
print(" ")

R2=(r2_score(y_test,predictR))
print('R2_SCORE VALUE IS :',R2)
print(" ")

```


Decission tree regressor

```
from sklearn.tree import DecisionTreeRegressor
DT= DecisionTreeRegressor()
DT.fit(X_train,y_train)

predictD = DT.predict(X_test)

MAE= (mean_absolute_error(y_test,predictD))
print('MEAN ABSOLUTE ERROR VALUE IS :',MAE)
print(" ")

MSE=(mean_squared_error(y_test,predictD))
print('MEAN SQUARED ERROR VALUE IS :',MSE)
print(" ")

MedianAE=(median_absolute_error(y_test,predictD))
print('MEDIAN ABSOLUTE ERROR VALUE IS :',MedianAE)
print(" ")

EVS=(explained_variance_score(y_test,predictD)*100)
print('ACCURACY RESULT OF DECISION TREE REGRESSOR IS :',EVS)
print(" ")

R2=(r2_score(y_test,predictD))
print('R2_SCORE VALUE IS :',R2)
print(" ")

import joblib
joblib.dump(RF,'rf.pkl')
import joblib
joblib.dump(RF,'dt.pkl')
```

support vector regressor

```
from sklearn.svm import SVR
svmA= SVR()
svmA.fit(X_train,y_train)

predictS = svmA.predict(X_test)

MAE= (mean_absolute_error(y_test,predictS))
print('MEAN ABSOLUTE ERROR VALUE IS :',MAE)
print(" ")

MSE=(mean_squared_error(y_test,predictS))
```

```

#import library packages
import pandas as p
import matplotlib.pyplot as plt
import seaborn as s
import numpy as n
#Load given dataset
data = p.read_csv("Bitcoin.csv")
import warnings
warnings.filterwarnings('ignore')
df=data.dropna()
df.head()
df.columns
df.rename(columns={"BTC price [USD]":"BTC_price_USD"},inplace=True)
df.rename(columns={"n-transactions":"n_transactions"},inplace=True)
df.rename(columns={"fee [USD]":"fee_USD"},inplace=True)
df.rename(columns={"btc search trends":"btc_search_trends"},inplace=True)
df.rename(columns={"Gold price[USD]":"Gold_price_USD"},inplace=True)
df.rename(columns={"SP500 close index":"SP500_close_index"},inplace=True)
df.rename(columns={"Oil WTI price[USD]":"Oil_WTI_price_USD"},inplace=True)
df.rename(columns={"M2(Not seasonally adjusted)[1e+09 USD]":"M2_money_supply_USA"},inplace=True)
df.info()
del df['Date']
df.columns
df.head()
df.tail()
from sklearn.metrics import
mean_absolute_error,mean_squared_error,r2_score,explained_variance_score,median
_absolute_error
X = df.drop(labels='BTC_price_USD', axis=1)
#Response variable
y = df.loc[:, 'BTC_price_USD']
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3)

```

```

#import library packages
import pandas as p
import matplotlib.pyplot as plt
import seaborn as s
import numpy as n
import warnings
warnings.filterwarnings('ignore')
df=data.dropna()
df.head()
df.columns
df.rename(columns={"BTC price [USD]":"BTC_price_USD"},inplace=True)
df.rename(columns={"n-transactions":"n_transactions"},inplace=True)
df.rename(columns={"fee [USD]":"fee_USD"},inplace=True)
df.rename(columns={"btc search trends":"btc_search_trends"},inplace=True)
df.rename(columns={"Gold price[USD]":"Gold_price_USD"},inplace=True)
df.rename(columns={"SP500 close index":"SP500_close_index"},inplace=True)
df.rename(columns={"Oil WTI price[USD]":"Oil_WTI_price_USD"},inplace=True)
df.rename(columns={"M2(Not seasonally adjusted)[1e+09
USD]":"M2_money_supply_USA"},inplace=True)
df.info()
del df['Date']
df.columns
df.head()
df.tail()
from sklearn.metrics import
mean_absolute_error,mean_squared_error,r2_score,explained_variance_score,median_absolu
lute_error
X = df.drop(labels='BTC_price_USD', axis=1)
#Response variable
y = df.loc[:, 'BTC_price_USD']
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3)

```

Ridge

```

from sklearn.linear_model import Ridge
rid= Ridge()
rid.fit(X_train,y_train)

predictri = rid.predict(X_test)

MAE= (mean_absolute_error(y_test,predictri))
print('MEAN ABSOLUTE ERROR VALUE IS :',MAE)
print(" ")

MSE=(mean_squared_error(y_test,predictri))

```

Ada boost regressor

```
from sklearn.ensemble import AdaBoostRegressor
ad= AdaBoostRegressor()
ad.fit(X_train,y_train)

predictad= ad.predict(X_test)

MAE= (mean_absolute_error(y_test,predictad))
print('MEAN ABSOLUTE ERROR VALUE IS :',MAE)
print(" ")

MSE=(mean_squared_error(y_test,predictad))
print('MEAN SQUARED ERROR VALUE IS :',MSE)
print(" ")

MedianAE=(median_absolute_error(y_test,predictad))
print('MEDIAN ABSOLUTE ERROR VALUE IS :',MedianAE)
print(" ")

EVS=(explained_variance_score(y_test,predictad)*100)
print('ACCURACY RESULT OF AdaBoostRegressor IS :',EVS)
print(" ")

R2=(r2_score(y_test,predictad))
print('r2_score',R2)
print(" ")
```

Home page interface:

BITCOIN PRICE PREDICTIONS

-0.07

fee_USD

btc_search_trends

Gold_price_USD

SP500_close_index

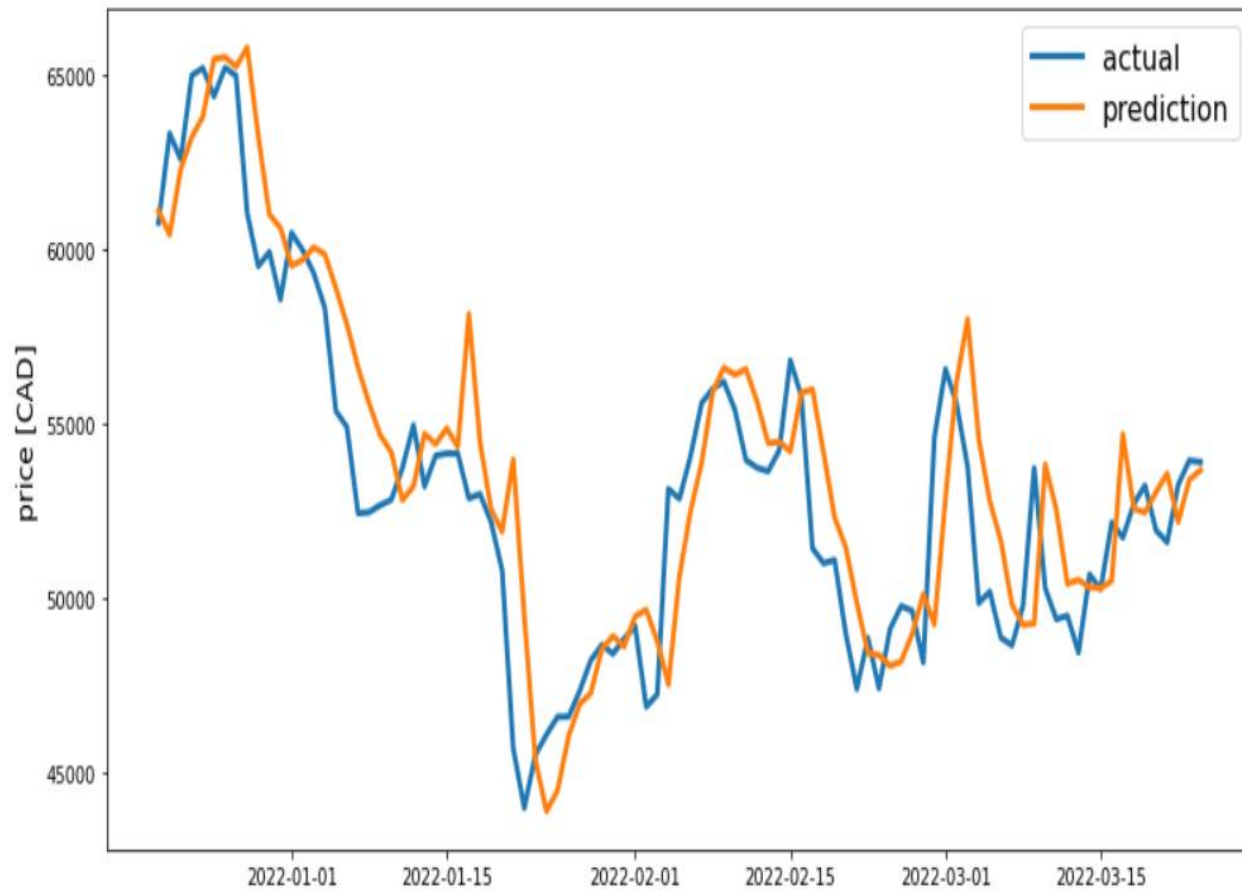
Oil_WTI_price_USD

M2_money_supply_USA

Predict

BITCOIN PRICE ANALYZE AND PREDICTION USING DATA SCIENCE PROCESS

OUTPUT GRAPH:



RE-2022-7203-plag-report

by Research Experts - Turnitin Report

Submission date: 03-Mar-2022 02:27AM (UTC-0600)

Submission ID: 1775424123

File name: RE-2022-7203.docx (147.75K)

Word count: 2576

Character count: 14296

BITCOIN PRICE ANALYZE AND PREDICTION USING DATA SCIENCE PROCESS

K. Raja Ramesh¹, V. Narendra², DR.G.Nagarajan³,

^{1,2} Student, ³Assistant professor, CSE, Sathyabama Institute of Science and Technology, Chennai, India

ABSTRACT:

Bitcoin is an advanced resource and an installment framework that is utilized as a type of Internet cash. It takes into consideration unknown installment starting with one individual then onto the next and is along these lines a favored installment strategy for criminal activities on the Internet. As of late Bitcoin has gotten a ton of consideration from the media and the general population because of its new value climb. The target of this paper is to decide the anticipated value heading of Bitcoin cost. AI models can almost certainly give us the knowledge we really want to find out with regards to the eventual fate of Cryptocurrency. It won't let us know the future yet it may let us know the overall pattern and course to anticipate that the costs should move. The proposed model is to construct an AI model where the information is utilized to made to find out with regards to the example in the dataset and the AI calculation is utilized to foresee the bitcoin cost.

Keywords: Keywords: Machine Learning, Bitcoin , Data Science

INTRODUCTION

MACHINE LEARNING

The goal of artificial intelligence (AI) is to predict the future based on the analysis of the past. Computerized reasoning (AI) ML is a kind of AI that allows PCs to learn without being updated. The core of artificial intelligence (AI) is the development of computer programmes that can adapt to new knowledge, as well as the nuts and bolts of machine learning. Preparation and forecasting need the use of specific computations. Uses preparation data to make projections on other test data based on the calculation's use of this preparation data.

RELATED WORK

In response to the increased interest in Bitcoin as a financial and social phenomena, new methods and methodologies for analysing Bitcoin data have emerged. This section presents the most similarly relevant past approaches for visual analysis of Bitcoin data. Many websites display Bitcoin blockchain data in easy-to-understand graphics. For example, blockchain.info [5] shows the current Bitcoin market value and the number of transactions each block. The majority of these websites provide information in the form of rudimentary charts that look like stock charts, presumably for the benefit of investors. Only a few technologies enable more in-

depth visual inspections of various it coin features. Both Bitcoin techniques are limited to certain subsets of transactions and focus on specific categories of businesses (mining pools and trading platforms), whereas BitConduite allows for an exploratory investigation of transactions from all types of stakeholders. Finally, visual methods for analysing user behaviour on the Bitcoin blockchain are rare and only cover the ability to answer prepared questions [11]. The goal of BitConduite is to provide a comprehensive, long-term, entity-centred perspective for exploratory activity analysis that no other method has ever offered.

LITERATURE SURVEY

Using LSTM and Embedding Networks to Improve Bitcoin Price Fluctuation Prediction Yang Li, Zibin Zheng, and Hong-Ning Dai in the year 2019

Bitcoin has drawn in broad consideration from financial backers, specialists, controllers, and the media. A notable and uncommon component is that Bitcoin's cost regularly changes fundamentally, which has anyway gotten less consideration. In this paper, we explore the Bitcoin value variance expectation issue, which can be depicted as whether Bitcoin value keeps or inversions after a huge change. Essential components, typical specialised swapping pointers, and elements generated by a Denoising autoencoder are all highlighted in this article. By making use of an Attentive LSTM organisation and an Embedding Network, we evaluate these aspects (ALEN). Specifically, a mindful LSTM organization can catch the time reliance portrayal of Bitcoin cost and an inserting organization can catch the concealed portrayals from related digital currencies. Test results exhibit that ALEN accomplishes predominant cutting edge

execution among all baselines. Besides, we examine the effect of boundaries on the Bitcoin value change expectation issue, which can be additionally utilized in a genuine exchanging climate by financial backers.

Using Machine Learning Algorithms to Predict Bitcoin Prices, Dr.P. Sriramy, Lekkala Sreekanth Reddy and others in the year 2020

We proposed a strategy for reliably estimating the Bitcoin cost in this study, which took into account a number of factors that influence the Bitcoin value. I followed down the benefits and obstructions of bitcoin value expectation by using social event data from multiple reference papers and applied it gradually. Various publications utilise a number of methods to estimate bitcoin's future worth. Because many papers have exact prices but others don't, this article employs the LASSO (least outright shrinkage choice administrator) algorithm, which is based on man-made conscious, to decrease the time complexity in this document's expectations. Different papers used different calculations such as SVM (support vector machine), coin mark up cap, Quandl, GLM, CNN (Convolutional Neural Networks), and RNN (Recurrent neural organisations) and so on, which don't live it up administration, but finding the outcomes from a larger data set is quick and easy in LASSO. As a result, we make a comparison of several calculations with the LASSO calculation. We expect to comprehend and monitor day by day drifts in the Bitcoin market while gaining knowledge into optimal highlights encompassing Bitcoin cost in the interaction that occurs in the paper, which is the first snapshot of the investigation. Every day, we update our informational index with new highlights related to the Bitcoin price and payment network from

previous years. A few data mining techniques, such as preprocessing, allow us to reduce the amount of complexity in the dataset before we run the analysis. Our second snapshot, based on the data available to us, will provide the most accurate indicator of the daily change in value with the greatest degree of certainty.

Bitcoin Price Prediction using
Machine Learning, Mr. Shivam Pandey1,
Mr. Anil Chavan.in the year 2021

In this paper, we endeavor to anticipate the Bitcoin cost precisely thinking about different boundaries that influence the Bitcoin esteem. For the primary period of our overview, we plan to comprehend and distinguish every day drifts in the Bitcoin market while acquiring knowledge into ideal highlights encompassing Bitcoin cost. For the second period of our overview, utilizing the accessible data, we will foresee the indication of the every day value change with most noteworthy conceivable exactness. Anticipating the future will forever be on the top of the rundown of employments for AI algorithms. It is our goal in this project to forecast Bitcoin's expenses by combining two powerful learning algorithms. Using problem description, movement, student evaluation, and the use of hands-on exercises depending on the usage of learning calculation to stimulate application, this work aims to enhance venture-based learning in the area of software engineering design.

EXISTING SYSTEM:

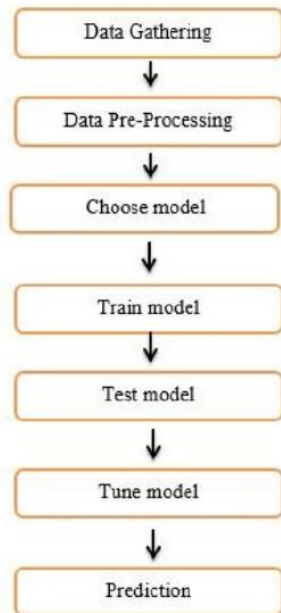
Bitcoin is a computerized resource and an installment framework that is utilized as a type of Internet cash. It takes into

account mysterious installment starting with one individual then onto the next and is along these lines a favored installment technique for criminal activities on the Internet. As of late Bitcoin has gotten a great deal of consideration from the media and the general population because of its new value climb. The target of this paper is to decide the anticipated value heading of Bitcoin cost. AI models can almost certainly give us the knowledge we really want to find out with regards to the fate of Cryptocurrency. It won't let us know the future however it may let us know the overall pattern and course to anticipate that the costs should move. The proposed model is to fabricate an AI model where the information is utilized to made to find out with regards to the example in the dataset and the AI calculation is utilized to foresee the bitcoin cost.

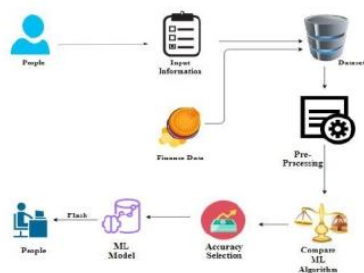
PROPOSED SYSTEM:

The media and the general public have been paying close attention to Bitcoin's recent price climb. Here, we're looking for patterns in bitcoin's price movement that might be used to forecast its future. The full dataset will be analysed using the supervised machine learning technique (SMLT) to identify variables, perform univariate, bivariate, and multivariate analysis, and examine missing value treatments, data validation, and data cleaning/preparation and visualisation. Based on our findings, model parameter sensitivity may be effectively assessed in terms of prediction accuracy. This paper proposes a machine learning-based approach and compares several machine learning methods against the given dataset. The model can be used to predict the bitcoin future. Performance metrics like accuracy, recall and precision can be calculated. Bitcoin future may be predicted and the investments can be made wisely.

WORKFLOW DIAGRAM



SYSTEM ARCHITECTURE



Module description:

Data Pre-processing

There are a number of approaches for AI to acquire the error margin of the machine learning model, which may be deemed to be close to the actual error rate of the dataset. If the amount of information is large enough to be representative of the general population, the approval techniques may not be required. Even in real-world settings, data that may or may not be a reliable indicator of the population in a dataset must be dealt with. Finding and resolving the value of an information type, whether a float variable or a whole number. The example of information was used to provide a fair evaluation of model fit on the preparation dataset when tuning model hyperparameters.

Data analysis of visualization

Accurate data interpretation is critical in the field of applied measurements and artificial intelligence (AI). When it comes to insights, quantitative representations and evaluations of data are by far the most important consideration. The collection of tools provided by information representation is critical in achieving a consensus on a subjective level. If you're trying to figure out how to analyse a dataset, this may help you identify trends, degenerate information, and more. Information representations may be used to convey and show essential connections in plots and outlines that are more instinctive and partners than proportions of affiliation or relevance with a little amount of information. Information representations. An in-depth look at the books referenced at the conclusion of this article would be recommended for anyone who are interested in learning more about

information representation and exploratory information research as separate topics.

ALGORITHM EXPLANATION

Characterization is a method used in artificial intelligence and insights in which a computer programme learns from the data that is sent to it and then uses that knowledge to create new perspectives. When it comes to determining if a person is male or female, or whether they've received spam, this data gathering may be limited to one kind of classification (such as determining whether the mail is spam or non-spam). Discourse recognition, penmanship acknowledgment, biometric distinguishing evidence, archive grouping, and so on are examples of characterisation issues. Marked data is used to make estimates in supervised learning. After understanding the information, the calculation figures out which name ought to be given to new information dependent on example and partners the examples to the unlabeled new information.

RESULT

This section contains the findings of the **interaction logs, user experience questionnaires, and participant-posted research questions.** Interaction Logbooks The **interaction logs** of each participant provide a rough estimate of the system that they are utilising. Specific characteristics were used or not used. It was possible to extract data. the **number of times participants performed a logged action** versus the number of times it occurred (through a mouse click). counted the total amount of clicks a proportions summary. Given the wide range of study subjects that participants were interested in, it's not unexpected that the proportions of interactions for each participant vary greatly in Fig. 8.

The filter view was used in a wide range of ways (10%–50%), and the entity browser was also used in a wide range of ways (0%–55%). **The transaction**

view received the least amount of interaction.

There are various reasons for its decreased use:

Some **research questions may not require** a detailed examination of specific businesses and their interactions. We didn't document all timeline interactions because **the view was at the bottom of the page** (e. g. varying the time range). The **interaction patterns** show that all participants frequently shifted between the filter and tree views (**filter tree filter tree filter tree filter tree filter**), implying that having them displayed next to each other visually reinforced their association. All participants, with the exception of P4, followed our proposed BitConduite approach (filter tree cluster). This participant's workflow was interrupted due to technical difficulties. Filter cluster entity browser was another common sequence utilised by participants P3, P4, and P6, which meant they didn't change the tree view and instead relied on the automatically picked group. We can infer from this data that participants used BitConduite as planned and used similar data-gathering techniques. The participants were asked to rate BitConduite's usefulness. There was no one available. Unnecessarily complicated system One individual took part in the discussion. They would require, they agreed, and one was neutral. additional encouragement to employ the method In the same way. Participants agreed that there were a lot of things they needed to learn prior to its application. With the exception of one, everyone was I was comfortable with the system and found it to be simple to operate.

```
... MEAN ABSOLUTE ERROR VALUE IS : 286.4136750630562
MEAN SQUARED ERROR VALUE IS : 220700.95360825915
MEDIAN ABSOLUTE ERROR VALUE IS : 138.39609999999948
ACCURACY RESULT OF RANDOM FOREST REGRESSOR IS : 98.40131103640768
R2_SCORE VALUE IS : 0.9840131103507547
```

DISCUSSIONS

The utilisation of BitConduite's entity-based analysis, as well as the necessity of exploratory analytical approaches, were highlighted during the preliminary evaluation. Despite this, we discovered two key flaws in our technique and the entity clustering mechanism we use. Approach. Because our exploratory method demands variable activity measurement calculation, scalability is limited. While we made every effort to reduce data processing times, the most significant bottleneck is on-the-fly clustering, which can take minutes in the worst-case situation (if many of entities are clustered at a time). In this case, a method of progressive clustering [19] would be useful. The process of combining entities is referred to as aggregation of entities. Although entity-based analysis is important, there is some uncertainty about entities because there are no methods for evaluating the quality of entity aggregation. Mixing services (also called as tumblers) make address aggregation more difficult by hiding the connections between addresses for privacy reasons. Furthermore, entity aggregation is computationally expensive and, in our case, memory intensive.

CONCLUSION

Data cleansing and processing were the first steps in the study, followed by missing value detection, exploratory analysis, and model construction and evaluation. A better accuracy score on a public test set will be discovered for the highest accuracy. The BITCOIN Market Price may be found with the aid of this software.

FUTURE WORK

- Bitcoin Market Price prediction to connect with AI model.
- To make this procedure more efficient, you may use a web application or a desktop programme to display the prediction result.
- To reduce the amount of time and effort required to implement in an AI system.

REFERENCES:

- [1]. Gulani Senthuran, Malka N. Halgamuge, " Prediction of Crptocurrency Market Price Using Deep Learning and Blockchain Information: Bitcoin and Ethereum".
- [2]. Yang Li 1,2 , Zibin Zheng 1,2,* and Hong-Ning Dai, " Enhancing Bitcoin Price Fluctuation Prediction Using Attentive LSTM and Embedding Network".
- [3]. Lekkala Sreekanth Reddy, Dr.P. Sriramy, " A Research On Bitcoin Price Prediction Using Machine Learning Algorithms".
- [4]. Mr. Shivam Pandey¹, Mr.Anil Chavan², Miss. Dhanashree Paraskar³, Prof. Sareen Deore, " Bitcoin Price Prediction using Machine Learning".
- [5]. Muhammad Ali Nasir^{1*} , Toan Luu Duc Huynh² , Sang Phu Nguyen³ and Duy Duong³, " Forecasting cryptocurrency returns and volume using search engines"

RE-2022-7203-plag-report

ORIGINALITY REPORT

11%

SIMILARITY INDEX

6%

INTERNET SOURCES

10%

PUBLICATIONS

0%

STUDENT PAPERS

PRIMARY SOURCES

- | | | |
|---|--|----|
| 1 | Christoph Kinkeldey, Jean-Daniel Fekete, Tanja Blascheck, Petra Isenberg. "BitConduite: Exploratory Visual Analysis of Entity Activity on the Bitcoin Network", IEEE Computer Graphics and Applications, 2021
Publication | 5% |
| 2 | arxiv.org
Internet Source | 2% |
| 3 | "Advances in Power Systems and Energy Management", Springer Science and Business Media LLC, 2021
Publication | 2% |
| 4 | www.ijstr.org
Internet Source | 1% |
| 5 | Diddi Priyanka, Diddi Anusha, T. Anandhi, P. Indria, E. Brumancia, R. M. Gomathi. "Chapter 25 Prediction of Chronic Kidney Disease by Best Accuracy Using Supervised Classification Machine Learning Approach", Springer Science and Business Media LLC, 2022
Publication | 1% |

6

scholars.ln.edu.hk

Internet Source

<1 %

7

Christoph Kinkeldey, Jean-Daniel Fekete, Tanja Blascheck, Petra Isenberg. "BitConduite: Exploratory Visual Analysis of Entity Activity on the Bitcoin Network", IEEE Computer Graphics and Applications, 2022

Publication

<1 %

8

camp-geval.com

Internet Source

<1 %

Exclude quotes On

Exclude matches Off

Exclude bibliography On

