

Analysis of Complains and Accidents using machine learning

A PROJECT REPORT

Submitted by

Akshat Khapra
Reg. No. 18MCA1112

in partial fulfillment for the award of the degree of

Master of Computer Applications



School of Computer Science and Engineering

Vellore Institute of Technology - Chennai Campus

Vandalur - Kelambakkam Road, Chennai - 600 127



School of Computer Science and Engineering

DECLARATION

I hereby declare that the project entitled **Analysis of Complains and Accidents using machine learning** submitted by me to the School of Computer Science and Engineering, Vellore Institute of Technology - Chennai Campus, 600 127 in partial fulfillment of the requirements of the award of the degree of **Master of Computer Applications** is a bona-fide record of the work carried out by me under the supervision of **Prof. M. Premalatha** . I further declare that the work reported in this project, has not been submitted and will not be submitted, either in part or in full, for the award of any other degree or diploma of this institute or of any other institute or University.

Place: Chennai
Date:

Signature of Candidate
(Akshat Khapra)



VIT[®]
Vellore Institute of Technology
(Deemed to be University under section 3 of UGC Act, 1956)
CHENNAI

School of Computer Science and Engineering

CERTIFICATE

This is to certify that the report entitled **Analysis of Complains and Accidents using machine learning** is prepared and submitted by **Akshat Khapra (Reg. No. 18MCA1112)** to Vellore Institute of Technology - Chennai Campus, in partial fulfillment of the requirement for the award of the degree of **Master of Computer Applications** is a bona-fide record carried out under my guidance. The project fulfills the requirements as per the regulations of this University and in my opinion meets the necessary standards for submission. The contents of this report have not been submitted and will not be submitted either in part or in full, for the award of any other degree or diploma and the same is certified.

Guide/Supervisor

Name: Prof. M. Premalatha
Date:

Head Of The Department

Name: Dr. M. Sivagami
Date:

Examiner

Name:
Date:

Examiner

Name:
Date:

(Seal of SCOPE)

Acknowledgement

I wish to express my sincere thanks to Dr. G.Viswanathan, Chancellor, for providing me an excellent academic environment and facilities for pursuing MCA program. I wish to express my sincere gratitude to Dr. shivgami, Program chair of MCA for providing me an opportunity to do my project work. I would like to express my gratitude to my guide Professor Premalatha m. who inspite of his busy schedule guided me in the correct path. Finally, I would like to thank my family and friends who motivated me during the course of my project work.

Akshat Khapra
Reg. No. 18MCA1112

Abstract

The street mishap has become a worldwide issue and set apart as the ninth unmistakable reason for death on the planet. Because of the huge number of street mishaps consistently, it has become a significant issue. It is completely unacceptable and disheartening to permit its resident to slaughter by street mishaps. Therefore, to deal with this overpowered circumstance, an exact examination is required. This paper has been done to predict in advance on time taken to solve the particular type of accident complain using machine learning approaches. We additionally make sense of those note worthy variables that clearly affect street mishaps and give some valuable proposals in regard to this issue. K Nearest Neighbors (KNN), Time Series and Random forest algorithms[RFA], These three administered learning procedures, to arrange the seriousness of mishaps and to forecast time taken to solve a particular type of complain which helps to arrange the manpower in advance on particular type of complain.

Contents

Declaration	i
Certificate	i
Acknowledgement	iii
Abstract	iv
1 Introduction	1
2 Overview / Literature Review	3
3 System Design	13
3.1 Feasibility Report	13
3.2 Machine Learning	13
3.3 Difference between Artifical Intelligence and Machine Learning .	14
3.4 Supervised Learning	15
3.5 Unspervised Learning	16
3.6 Semi Supervised Learning	16
3.7 Reinforcement Learning	17
3.8 How does Supervised Learning works	17
3.9 How does Supervised Training Works	18
3.10 How to evalute machine learning models	20
3.11 What are Neural Network and How they are Trained	20
3.12 Why is machine learning successful	21
3.13 What ways machine learning used for	22
3.14 Which services are available for machine learning	23
3.15 How Data are Process	25

3.16	Pre- Process of Data	27
3.17	Preprocessing Operation	29
3.18	Preprocessing Granularity	30
3.19	Window Aggregation During Traning and prediction	31
3.20	Machine Learning Pipeline on Google Cloud	32
3.21	Where To Do Processing	34
4	Implementation of System/ Methodology	40
4.1	Random Forest Algorithm	40
4.2	KNN algorithm	42
4.3	Time Series	44
4.4	What is time series and how it is used?	44
4.5	components	46
4.6	Time Series Models	46
4.7	Decomposition	47
4.8	Moving Average	48
4.9	CODE	49
5	Results and Discussions	59
6	Conclusion and Future Work	64
	REFERENCES	64

List of Figures

3.1	Machine Learning Epic	14
3.2	Supervised Learning Example	19
3.3	How Data are Trained	21
3.4	Data Flow Diagram	25
3.5	Pre-process	28
3.6	pre-process granularity	31
3.7	High Level Architecture	33
3.8	Data Flow Diagram	37
3.9	Transformation Training and Evaluation Data	38
3.10	Summary Table	39
4.1	Random Forest Algorithm	41
4.2	K-nearest neighbor	42
4.3	Time Series	45
4.4	Additive Model	47
4.5	Decomposition	48
4.6	Decomposition	49
5.1	KNN	59
5.2	Random Forest Algorithm	59
5.3	Time Series	60
5.4	Time Series	61
5.5	Time Series	61
5.6	Time Series	62
5.7	Time Series	62
5.8	Time Series	63
5.9	Time Series	63

Chapter 1

Introduction

- Background: The project is related to the accident complaints, In which analysis is done on the duration of the complaints. The data is about 52 * 364558[rows*columns]which has 52 different attributes and have different scenario through which analysis can be done. In the project all different types of complaints have been noted from which top ten complain are used for prediction, The prediction is done using kNN [k- nearest neighbor], Random Forest algorithm, and Naive Baye's algorithms.
- Statement: The project is related to the analysis of road accident complaints, In which the analysis is done on the duration of time taken to solve the cases,By which man power in the future will be managed in a systematic way. as well, agencies will have advance knowledge about the duration to solve the problems.
- Motivation: The project well in advance decides about the manpower used to solve the particular problem. By which human resources can be managed properly in the future.
- Post/Related work [Existing methods including pros and cons of the methods should be cited wherever possible].
- Challenges
 - find out the subject of data?
 - find out complain types?
 - find out the location type?

- Create a different subset of datasets and analyze the data?
 - Use label encoder and get dummies for the appropriate column?
 - find out x,y?
 - find out whether it is related to supervised or unsupervised algorithms?
- Essence of your approach [The work has done through various machine learning algorithms and the coding part is done through python programming in the jupyter notebook].
- Aim(s) and Objective(s)::The objective of the project is to perform an analysis of complaints and accidents using machine learning and to manage the manpower according to complaints. Also, the project describes the prediction of complaints in the future and the time taken to resolve if any case arises in future.

Chapter 2

Overview / Literature Review

- 1) This paper shows an overview of street mishap examination strategies in information mining. In the information mining there is no. of procedures accessible for grouping and arrangement, from those methods k-mean, affiliation rule, SVM, Weka device was utilized in already inquire about for street mishap examination. In our day by day life there are no. of mishap increments and it is large issue to us in light of the fact that no. of individuals passing and harmed for that improve the street transportation framework is required. In this exploration self association map (SOM) is utilized for discover a no. of example to examination the street mishap information which help to forecast the mishap reasons and improve the exactness of investigation contrast with k-implies grouping calculation. With the assistance of SOM, groups are made and investigate them. Self Organizing map technique depends on neural system, it is utilized as an unaided learning strategy. It will assist with improving examination exactness. In the paper there are two algorithms used to predict the accident and from that k-means gives better forecast. K-mean calculation is best one grouping calculation it is used in bunch examination in information digging for the most part utilized for content bunching it expect to cause bunch of same article or includes each group to have their diverse item , groups are distinctive with one another. It is important to characterize centroids to each bunch and the method of centroid structure is object of each group are like one another. SOM otherwise called kohonen highlight map, It is an Unsupervised neural system, It is utilized as a Clustering device of high-dimensional and complex information and it keeps up the topology of the dataset in this calculation preparing happens by means of rivalry between the neurons and

here difficult to appoint arrange hubs to explicit information classes ahead of time, It can be utilized for identifying comparability and degrees of closeness, It is expected that info design fall into adequately enormous particular groupings. In this paper Road Accident Data Analysis portray by utilizing the various strategies for information mining. In this paper the investigation of K-mean calculation is given. It is chiefly utilized for grouping of information, and the outcome got from this calculation shows the exactness however it has some impediment, it required more emphasis for process execute. Self association map give the best outcome as contrast with k-mean, it create design map which is straightforward and expectation.

- 2) Street mishaps are a significant issue these days that causes passing and incapacity everywhere throughout the world. As indicated by the factual records of World Health Organization, Oman positioned at the twelfth situation in Road car crashes internationally and first s position among the Gulf nations. The paper tends to an inside and out investigation that distinguishes the contributory elements, causes behind the street crashes and the evaluation of the components that influence the recurrence and seriousness of mishaps dependent on the accident information accessible. It additionally audit about the different components engaged with the foundation of various sorts of mishaps and the related procedures that are applied for an intensive examination in the field of interests. The assurance of different parameters that prompts the accident is researched by applying different measurable methods which causes us to acquire a solid end and reconstruct dependent on the necessities. The focused on examination additionally gives a superior comprehension to creating measures to improve their wellbeing execution. It additionally presents the creative methodology called AI based prescient investigation to anticipate the quantity of mishaps that may occur sooner or in near future. Streets are the primary methods for transport in Oman which thusly prompts considerable number mishaps. This paper broke down practically 50 of the decade information with the idealistic sign in different street mishap type and setback seriousness. The key discoveries rose up out of our investigation shows that the seriousness of accidents and the mishaps rates are decreased significantly over the period. The primary explanations for this decrease of number mishaps are because of the mindfulness made by the administration in the psyches of the individuals. Notwithstanding that, from April 2016, onwards ROP had made most vital strides towards petty criminal offense. The extreme laws and discipline are executed for street

petty criminal offense including fine range from 200 OMR to 3000 OMR and prison term. Oman needs a maintainable vehicle strategy to diminish the traffic volume and its serious results. From 2015 onwards, Public Bus Transport System was presented generally in Muscat Governorate. Waiting ideally in future 2017-2018, as information are reflected to show how much advantageous is the new law execution and whether it assumes an essential job in diminishing street mishap rate.

- 3) The emotional increment in street auto collisions on the planet is causing major issues in each part of human lives. The most significant and important nature of traffic qualities, causation investigation, and relationship between various causal components have been overlooked. Besides, the auto collision information is just used to direct a simple factual investigation and information mining endeavors which results just in examples and insights. The fundamental focuses of this street mishap information grouping are to distinguish the major and key factors that cause the street auto collision and structure arrangements and preventive activities that would diminish the mishap seriousness level. AI calculations are utilized to break down the information, extricate concealed examples, foresee the seriousness level of the mishaps and sum up the data in a helpful arrangement. In this work, we have applied diverse AI arrangement calculations and examined here the six calculations with high exactness and best grouping exhibitions, for example, Fuzzy-FARCHD, Random Forest, Hierarchal LVQ, RBF Network (Radial Basis Function Network), Multilayer Perceptron, and Naïve Bayes on street car crash informational index acquired from UK street car crash of the year 2016. The informational index contains data on all street mishap losses across Calder dale. The outcomes from our investigation show that Fuzzy-FARCHD calculation is compelling to arrange the dataset and accomplishes an exactness of 85.94. In this work, we have uncovered that Lighting Conditions, first Road Class and No., Number of vehicles are the key highlights in choosing the characteristics. Road traffic accidents are the key reason for serious injuries as well as the death of precious human lives. Discovering the factors that are related to the class values that are important to achieve an accurate result. The brutality of the accident problem is getting a catastrophic level becoming horrific shockingly and indicating that adequate measures have not been taken to prevent, control and/or lessen the appalling rate of the accident. Various scholars have tried to solve this issues but still, there are a lot of uncovered issues and gaps in

predicting accident severity and specifically discover the influential factors such as time and season in which the accident frequently happened. This makes the field of traffic accident analysis and prediction more challenging. The main target of this research work is to discover the potential data mining technology at road traffic accident data for making a classification model. The developed classification model could support decision makers, policy designer and traffic officers for making effective decisions in traffic control actions. Using the analyzed outcome, decision-makers can easily understand the accident mode, driver's behavior, time, road and weather conditions and other key factors that are causing traffic accidents resulting in the fatalities and serious injuries so as to articulate improved traffic safety control strategies. They also may utilize the predictive models to take a new initiative in road safety, accident prevention and to develop new policies in this regards. Hence, we can conclude that Fuzzy FARCHD algorithm can be applied in the machine learning tools to classify and predict the road accidents based on different attributes. For the future considering the frequently increasing size of the data sets, more features, and clusters, it's better to use deep learning techniques for improved classification and cluster of the data records.

- 4) As indicated by the car crash programmed notice framework, data on mishap event, area of mishap, level of injury of survivors of car crashes is transmitted to crisis salvage offices, which helps crisis clinical treatment at the mishap site and determination of proper clinical organizations. In the paper, we propose a model for expectation of the seriousness of car crash wounds pertinent to the programmed warning arrangement of car crashes. We utilize a choice tree for the prescient model. The utilization information depends on auto collision overview information and the NASS-CDS information, which gathers clinical records of mishap casualties. The prescient models created in this paper are relied upon to be helpful in the programmed car crash warning framework since they have a high precise relying upon the presentation investigation and supplementation of each model. In order to solve the social problems caused by traffic accidents, we constructed an injury grade prediction model of traffic accidents based on the machine learning model. In this paper, decision tree is used. The proposed method uses the NASS-CDS database for five years from 2011 to 2015 in the United States, and predicts the injury grade as 107 inputs and VAIS of the 11 of outputs. When the injury grade is predicted by the developed Decision Tree

what we made, it has the accuracy from the latter half of 40 percent to the latter half of 50 percent. Therefore, it is necessary to optimize the input parameters using the Genetic Algorithm (GA) and Principal Component Analysis (PCA) , and to improve the performance by optimizing the model parameters.

- 5) Street mishap information is the essential proportion of security without which the scale and nature of street wellbeing issues can't be set up with assurance. Accordingly, the presence of a solid mishap database is a key factor in the administration of street security. The nation's mishap information assortment framework is as yet conflicting and sporadic on the grounds that there is neither a uniform information assortment design nor a hearty arrangement of dependable and significant recovery of customary and precise information. Precise and far reaching mishap records are the reason for mishap examination. The powerful utilization of mishap records relies upon three variables, to be specific, the precision of the information, record maintenance and information investigation. The requirement for an exclusive requirement of occurrence detailing is a significant essential for the utilization of mishap records to create street security measures. On the off chance that the first occurrence report itself is poor, at that point the investigation and utilization of the outcomes will be poor. Incorrect and deficient mishap information make the outcomes fluffy, deluding, and not productive. Street mishaps are uncertain and sporadic occasions and their examination should know about the elements that influence them. Street mishaps are characterized by a lot of properties that are regularly extraordinary. The primary trouble of mishap information investigation is its heterogeneity. Hence, heterogeneity must be estimated through the examination of the information, or there will be consequences various connections among the information may remain covered up. Division is utilized to lessen information heterogeneity utilizing various estimates, for example, master information, however there is no assurance that this will bring about the best division of the gathering including street mishaps. Group investigation can assist with dividing street mishaps information. This paper presents a method for analyzing accident patterns of different types of accidents on roads, which uses k to represent clustering and association rules mining algorithms. The study used accidents in the Maharashtra road network in 2015 and 2016. K-means clustering finds five clusters based on attribute accident type, road type, light condition, and road characteristics. Association rule mining has been ap-

plied to every cluster as well as the whole data set to create rules. Strong rules are used for analysis with high lift values. The rules of every cluster disclose situations related to incidents within the cluster.

- 6) Street security has become a significant issue in the urban territories because of the high vehicle thickness. Street wellbeing can be improved by decreasing the mishaps. Street mishap causes traffic deterrent which has become terrible particularly in huge urban communities. Thusly, examining the street mishaps precisely can assist with taking care of the issue of car accidents. In our undertaking, we propose a half breed model that consolidates both KNearest Neighbor and Support Vector Machines calculation for street mishap investigation and expectation of mishap type, which depends on the progressive learning approach. The mishap types are delegated crash, smashed and drive, fire and slip. Our proposed model uses the blend of both KNN and SVM calculations with the verifiable datasets gathered from UCI Repository. This broke down information will be progressively valuable to propose better security measures to maintain a strategic distance from car accidents. We tentatively investigate the presentation of both KNN and SVM calculations utilizing R programming with huge mishap datasets. Results show that our cross breed model upgrades the precision of street mishap investigation. Breaking down the street mishaps with more precision can ease car accidents in urban regions. In this paper, we proposed a crossover model that applies both KNN and SVM calculations to examine street mishaps precisely and anticipate the kind of mishaps. These strategies and methods have been examined and incorporated for street mishap examination. Our strategy has taken the huge recorded information of street mishaps and anticipating the mishap type. We have additionally assessed the exhibition of KNN and SVM exclusively. In this way, we tentatively exhibited that our proposed model for street mishap examination has prevalent execution and accomplished the exactness of 92 percent.
- 7) Information mining is the way toward dissecting information from various area and summing up it into helpful data and it very well may be utilized for settling on keen business choice. It isn't explicit to any industry, applied in practically all regions to investigate the chance of concealed information. Significant information mining procedures are order, bunching, time arrangement investigation, affiliation rule mining and relapse. Grouping is a significant information mining strategy which investigates the information

and characterizes the information into a predefined set of classes. There are at least one AI calculations are accessible; they are straight relapse, strategic relapse, choice tree, SVM, Naïve Bayes, KNN, Random Forest and inclination boosting calculation. The principle point of this work is to break down the presentation of AI calculations utilizing R apparatus. The presentation measures are exactness, review, precision rate, genuine positive, bogus positive, blunder rate and level of effectively grouped cases. Trial results uncovers that, KNN outflanked different calculations with higher precision and a lower blunder rate. This examination paper shows that four causes: lethal mishap, significant injury mishap, minor injury mishap, all out mishaps and year savvy reports for the street mishaps in India. The rest of the part of the paper has seven areas. Segment II portrays the approach of this investigation work. Area III presents the characterization exactness. The end is given in segment IV. The fundamental point of this examination is to assess and explore order calculations. The street mishap dataset is utilized to test the presentation of the chose classifiers. The calculation which has the least mean outright mistake and higher exactness is picked as the best calculation. By considering various parameters of precision and the blunder rate, it is discovered that the KNN characterization calculation is the best calculation with a most extreme exactness of 93.7 than other grouping calculations. Trial results, street mishaps expanded in the time of 2016 in Tamilnadu, India.

- 8) Crash injury seriousness expectation is a promising exploration focus in rush hour gridlock wellbeing. Customarily, different factual techniques were utilized for demonstrating crash injury severities. As of late, AI based techniques are getting mainstream because of their great prescient exhibition. In any case, AI based models are normally condemned as they perform like a black-box. In this examination, we target looking at the prescient exhibition, including forecast exactness and estimation of variable significance, among different AI and measurable techniques with unmistakable demonstrating rationale for crash seriousness investigation. The accident seriousness, street geometry and traffic stream information were gathered at turnpike veer zones in Florida. We assessed two most usually utilized factual techniques which were requested probit (OP) model and multinomial logit (MNL) model, and four well known AI strategies including K-Nearest Neighbor (KNN), Decision Tree (DT), Random Forest (RF), and Support Vector Machine (SVM). The right expectation rate for each

crash seriousness level and the general right forecast rate were determined. The outcomes indicated that AI strategies had higher foreseeing exactness than measurable techniques, however they experienced the over-fitting issue. The RF strategy had the best expectation in generally speaking and serious accidents while OP was the most vulnerable one. We thought about factor significance on crash seriousness through an annoyance based affectability investigations. The outcomes demonstrated that the surmisings of variable significance from various techniques were not constantly reliable and should be given cautious consideration.

- 9) Car crashes are among the most basic issues confronting the world as they cause numerous passings, wounds, and fatalities just as financial misfortunes consistently. Exact models to anticipate the auto collision seriousness is a basic errand for transportation frameworks. This examination exertion sets up models to choose a lot of persuasive components and to develop a model for grouping the seriousness of wounds. These models are figured by different machine learning methods. Administered machine learning calculations, for example, AdaBoost, Logistic Regression (LR), Naive Bayes (NB), and Random Forests (RF) are actualized on auto collision information. Destroyed calculation is utilized to deal with information irregularity. The discoveries of this examination show that the RF model can be a promising apparatus for foreseeing the injury seriousness of car crashes. RF calculation has demonstrated preferable execution with 75 percent precision over LR with 74.5percent, NB with 73.1percent, and AdaBoost with 74 percent exactness. This examination explored the proficiency of the four AI calculations to construct classifiers that are exact and dependable. This incorporates the Random Forest (RF), Logistic Regression (LR), Naïve Bayesian Classifier (NB), and AdaBoost calculations. In view of the disarray lattice F1-Score, the test outcomes show that the Random Forest appeared to perform superior to different models. This exploration study shows that the calculations can anticipate mishaps with a 75.5 percent exactness. This investigation can help give valuable data to parkway architects and transportation creators to plan more secure streets. Further examinations ought to be done to gather related data and explore the effects of these components. It is suggested that Random timberland (the best model for foreseeing road collides with) be applied in checking Fatal and genuine wounds. The prescribed prescient model can be utilized to quickly and proficiently distinguish the Key factor causing car accidents. One con-

finement of the present examination is that a portion of the components (for example attributes of the driver, traveler, and person on foot, alongside traffic conditions) may effectly affect mishap seriousness and span, which are not considered in light of the absence of reasonable information.

- 10) Compelling information driven dynamic requires sufficient proof and ability. Sufficient proof is in some cases estimated through the accessibility of enough information and accessibility of enough sign in that information. That is, the information must contain significant credits that add to the event or non-occurrence of an occasion. Sufficient ability can be estimated by the capacity of chiefs, or of their subordinates, to have the option to break down and model utilizing this information. For governments the test will in general lie in the two circles, having sufficient proof and ability. In this paper, we use street traffic mishap (RTA) information to talk about a portion of these difficulties and explore displaying utilizing information from street traffic mishaps in South Africa. Right off the bat, we will take a gander at how RTA information is gathered and talk about a portion of the difficulties thusly of dealing with information. Second, we examine explanations behind gathering this information, talk about ebb and flow research to moderate RTAs in South Africa; and third, we present outcomes from the investigation and demonstrating of RTAs utilizing a dataset from Soshanguve, South Africa. In the wake of performing exploratory information examination to sum up the fundamental attributes inside the information, we saw a need to research the effect that uncompromising licenses (codes: C1, EC1 and EC) have on street wellbeing in South Africa. Prior to examining this effect, we first fragmented the information utilizing grouping as a starter venture to find relationship inside the factors in the information. The outcomes from the examination and displaying indicated that, despite the fact that a huge segment of drivers in Soshanguve, South Africa are in control of hard core licences, the uncompromising drivers drive motorcars/station wagons more than combis/mini busses and light conveyance trucks. We likewise found that these drivers supported the most genuine wounds in contrast with light obligation permit drivers. In this investigation, we dissected RTAs examples and patterns by utilizing grouping and affiliation rule mining. The examination utilized mishaps that happened in Soshanguve, South Africa in the years 2015-2016. We initiated our examination by first taking a gander at the information challenges in South Africa and difficulties we had in our information. We at that point utilized PAM grouping as a funda-

mental assignment before finding relationship inside the information. Affiliation rule mining was then applied on each bunch and the EUD. The outcomes uncovered that despite the fact that an enormous segment of drivers in Soshanguve, South Africa are in control of rock solid licenses, the hard core drivers drive engine vehicles/station wagons more than combis/mini busses, light conveyance trucks and different kinds of vehicles as appeared in Table II. We likewise found that these drivers supported the most genuine wounds in contrast with light obligation permit drivers (code: B and EB). In spite of the fact that the seriousness of a mishap is dependent on a scope of various components, this paper featured the significance to further investigate the impact that each heavy-duty licence code has on mishap seriousness. This will be altogether researched in future work. Additionally, in future work, we will apply ascription to address missing information issues experienced in the dataset utilized in the examination.

Chapter 3

System Design

3.1 Feasibility Report

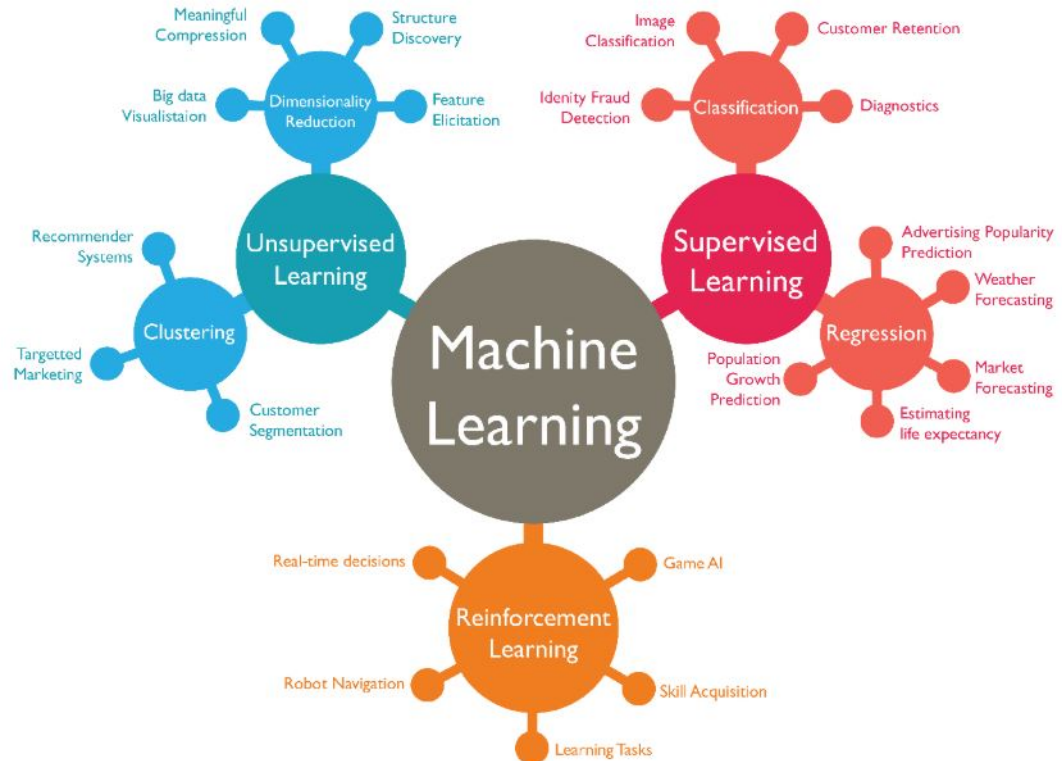
- The project marks on the analysis of road accident using machine learning approaches. The project has data of around 3 lakhs which is collected by one of the agency in new York city. Here our main moto is to do analysis on road accident complains and time taken to solve the particular type of complain, by that, we can strength the manpower on particular type of case. As well, management for cops can be done properly.

3.2 Machine Learning

At an exceptionally elevated level, AI is the way toward showing a PC framework how to make exact forecasts when taken care of information. Those expectations could be noting whether a bit of organic product in a photograph is a banana or an apple, spotting individuals going across the street before a self-driving vehicle, regardless of whether the utilization of the word book in a sentence identifies with a soft cover or a lodging reservation, whether an email is spam, or perceiving discourse precisely enough to create subtitles for a YouTube video. The key distinction from conventional PC programming is that a human engineer hasn't composed code that trains the framework how to differentiate between the banana and the apple. Rather an AI model has been instructed how to dependably separate between the organic products by being prepared on a lot of information, in

this example likely an enormous number of pictures marked as containing a banana or an apple.

Figure 3.1: Machine Learning Epic



3.3 Difference between Artificial Intelligence and Machine Learning

AI may have delighted in colossal achievement generally, however it is only one technique for accomplishing man-made consciousness. At the introduction of the field of AI during the 1950s, AI was characterized as any machine equipped for playing out an undertaking that would ordinarily require human insight. Computer based intelligence frameworks will for the

most part exhibit probably a portion of the accompanying attributes: arranging, getting the hang of, thinking, critical thinking, information portrayal, recognition, movement, and control and, to a lesser degree, social knowledge and imagination. Close by AI, there are different methodologies used to manufacture AI frameworks, including developmental calculation, where calculations experience arbitrary changes and blends between ages trying to "advance" ideal arrangements, and master frameworks, where PCs are customized with decides that permit them to copy the conduct of a human master in a particular area, for instance an autopilot framework flying a plane.

3.4 Supervised Learning

During preparing for administered learning, frameworks are presented to a lot of marked information, for instance pictures of written by hand figures explained to demonstrate which number they relate to. Given adequate models, an administered learning framework would figure out how to perceive the groups of pixels and shapes related with each number and in the long run have the option to perceive written by hand numbers, ready to dependably recognize the numbers 9 and 4 or 6 and 8.

Be that as it may, preparing these frameworks ordinarily requires immense measures of named information, with certain frameworks waiting be presented to a large number of guides to ace an errand.

Accordingly, the datasets used to prepare these frameworks can be tremendous, with Google's Open Images Dataset having around 9,000,000 pictures, its named video storehouse YouTube-8M connecting to 7,000,000 named recordings and ImageNet, one of the early databases of this sort, having in excess of 14 million arranged pictures. The size of preparing datasets keeps on developing, with Facebook as of late reporting it had assembled 3.5 billion pictures freely accessible on Instagram, utilizing hashtags connected to each picture as names. Utilizing one billion of these photographs to prepare a picture acknowledgment framework yielded record levels of precision - of 85.4 percent - on ImageNet's benchmark.

The arduous procedure of marking the datasets utilized in preparing is regularly done utilizing crowdworking administrations, for example, Amazon

Mechanical Turk, which gives access to an enormous pool of minimal effort work spread over the globe. For example, ImageNet was assembled more than two years by about 50,000 individuals, for the most part selected through Amazon Mechanical Turk. In any case, Facebook's methodology of utilizing openly accessible information to prepare frameworks could give an elective method of preparing frameworks utilizing billion-in number datasets without the overhead of manual marking.

3.5 Unsupervised Learning

In contrast, unsupervised learning tasks algorithms with identifying patterns in data, trying to spot similarities that split that data into categories. An example might be Airbnb clustering together houses available to rent by neighborhood, or Google News grouping together stories on similar topics each day. The algorithm isn't designed to single out specific types of data, it simply looks for data that can be grouped by its similarities, or for anomalies that stand out.

3.6 Semi Supervised Learning

The significance of enormous arrangements of named information for preparing AI frameworks may decrease after some time, because of the ascent of semi-regulated learning.

As the name proposes, the methodology blends directed and solo learning. The procedure depends after utilizing a limited quantity of named information and a lot of unlabelled information to prepare frameworks. The marked information is utilized to incompletely prepare an AI model, and afterward that somewhat prepared model is utilized to name the unlabelled information, a procedure called pseudo-naming. The model is then prepared on the subsequent blend of the marked and pseudo-named information.

The reasonability of semi-administered learning has been supported as of late by Generative Adversarial Networks (GANs), AI frameworks that can utilize marked information to produce totally new information, for instance making new pictures of Pokemon from existing pictures, which thusly can be utilized to help train an AI model.

Were semi-administered figuring out how to get as compelling as managed learning, at that point access to enormous measures of processing force may wind up being increasingly significant for effectively preparing AI frameworks than access to huge, marked datasets.

3.7 Reinforcement Learning

An approach to comprehend fortification learning is to consider how somebody may figure out how to play an outdated PC game just because, when they aren't acquainted with the standards or how to control the game. While they might be a finished fledgling, in the long run, by taking a gander at the connection between the catches they press, what occurs on screen and their in-game score, their exhibition will show signs of improvement and better.

A case of fortification learning is Google DeepMind's Deep Q-arrange, which has beaten people in a wide scope of vintage computer games. The framework is taken care of pixels from each game and decides different data about the condition of the game, for example, the separation between objects on screen. It at that point thinks about how the condition of the game and the activities it acts in game identify with the score it accomplishes.

Over the procedure of numerous patterns of playing the game, in the end the framework assembles a model of which activities will augment the score wherein situation, for example, on account of the computer game Breakout, where the oar ought to be moved to so as to block the ball.

3.8 How does Supervised Learning works

Everything starts with preparing an AI model, a numerical capacity able to do over and again altering how it works until it can make exact forecasts when given new information. Prior to preparing starts, you initially need to pick which information to accumulate and choose which highlights of the information are significant. A colossally streamlined case of what information highlights are is given in this explainer by Google, where an AI model is prepared to perceive the distinction among lager and wine, in light of two highlights, the beverages' shading and their alcoholic volume (ABV). Each

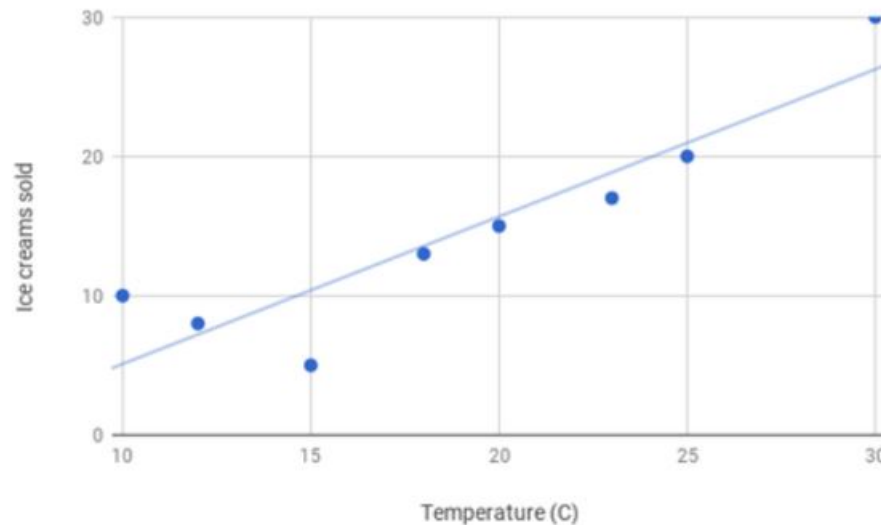
drink is marked as a lager or a wine, and afterward the applicable information is gathered, utilizing a spectrometer to quantify their shading and hydrometer to gauge their liquor content. A significant point to note is that the information must be adjusted, in this occasion to have a generally equivalent number of instances of lager and wine. The assembled information is then part, into a bigger extent for preparing, state around 70 percent, and a littler extent for assessment, state the staying 30 percent. This assessment information permits the prepared model to be tried to perceive how well it is probably going to perform on genuine information. Before preparing gets in progress there will be and large additionally be an information planning step, during which procedures, for example, deduplication, standardization and blunder remedy will be done. The subsequent stage will pick a proper AI model from the wide assortment accessible. Each have qualities and shortcomings relying upon the kind of information, for instance some are fit to taking care of pictures, some to content, and some to absolutely numerical information.

3.9 How does Supervised Training Works

Essentially, the preparation procedure includes the AI model consequently tweaking how it capacities until it can make exact expectations from information, in the Google model, effectively naming a beverage as lager or wine when the model is given a beverage's shading and ABV. A decent method to clarify the preparation procedure is to consider a model utilizing a straightforward AI model, known as direct relapse with slope drop. In the accompanying model, the model is utilized to evaluate what number of frozen yogurts will be sold dependent outwardly temperature. Envision taking past information indicating dessert deals and outside temperature, and plotting that information against one another on a dissipate chart - essentially making a dispersing of discrete focuses. To anticipate what number of desserts will be sold in future dependent on the outside temperature, you can draw a line that goes through the center of every one of these focuses, like the outline underneath.

When this is done, frozen yogurt deals can be anticipated at any temperature by finding where the line goes through a specific temperature and perusing off the comparing deals by then. Taking it back to preparing an AI model,

Figure 3.2: Supervised Learning Example



in this example preparing a direct relapse model would include altering the vertical position and slant of the line until it lies in the entirety of the focuses on the dissipate diagram. At each progression of the preparation procedure, the vertical separation of every one of these focuses from the line is estimated. On the off chance that an adjustment in incline or position of the line brings about the separation to these focuses expanding, at that point the slant or position of the line is altered in the contrary course, and another estimation is taken. Thusly, by means of numerous modest changes in accordance with the incline and the situation of the line, the line will continue moving until it in the end settles in a position which is a solid match for the circulation of every one of these focuses, as found in the video underneath. When this preparation procedure is finished, the line can be utilized to make exact forecasts for how temperature will influence frozen yogurt deals, and the AI model can be said to have been prepared. While preparing for increasingly complex AI models, for example, neural systems contrasts in a few regards, it is comparable in that it likewise utilizes an "inclination drop" approach, where the estimation of "loads" that change input information are over and again changed until the yield esteems delivered by the model are as close as conceivable to what is wanted.

3.10 How to evaluate machine learning models

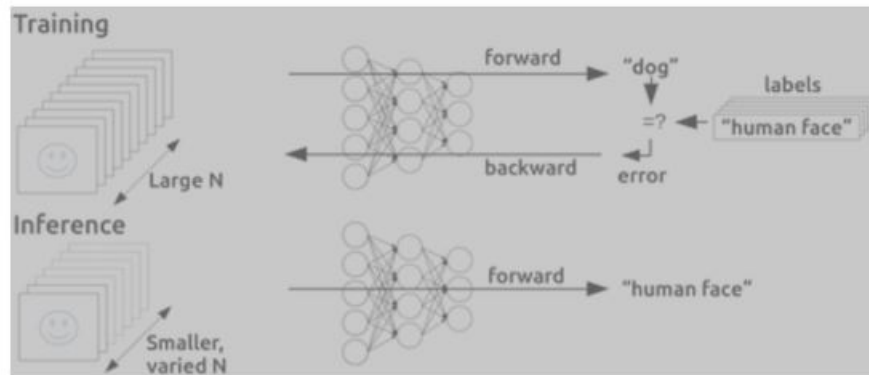
When preparing of the model is finished, the model is assessed utilizing the rest of the information that wasn't utilized during preparing, assisting with checking its certifiable execution. To additionally improve execution, preparing parameters can be tuned. A model may be adjusting the degree to which the "loads" are changed at each progression in the preparation procedure.

3.11 What are Neural Network and How they are Trained

A significant gathering of calculations for both managed and unaided AI are neural systems. These underlie quite a bit of AI, and keeping in mind that straightforward models like direct relapse utilized can be utilized to make expectations dependent on few information highlights, as in the Google model with brew and wine, neural systems are helpful when managing huge arrangements of information with numerous highlights. Neural systems, whose structure is inexactly propelled by that of the mind, are interconnected layers of calculations, called neurons, which feed information into one another, with the yield of the former layer being the contribution of the ensuing layer. Each layer can be thought of as perceiving various highlights of the general information. For example, consider the case of utilizing AI to perceive manually written numbers somewhere in the range of 0 and 9. The principal layer in the neural system may quantify the shade of the individual pixels in the picture, the subsequent layer could spot shapes, for example, lines and bends, the following layer may search for bigger parts of the composed number - for instance, the adjusted circle at the base of the number 6. This continues entirely through to the last layer, which will yield the likelihood that a given written by hand figure is a number somewhere in the range of 0 and 9. See progressively: Special report: How to execute AI and AI (free PDF) The system figures out how to perceive every part of the numbers during the preparation procedure, by step by step tweaking the significance of information as it streams between the layers of the system. This is conceivable because of each connection between layers having a joined weight, whose worth can be expanded or diminished to adjust that

connection's noteworthiness. Toward the finish of each preparation cycle the framework will look at whether the neural system's last yield is drawing nearer or further away based on what is wanted - for example is the system showing signs of improvement or more terrible at recognizing a written by hand number 6. To close the hole between the real yield and wanted yield, the framework will at that point work in reverse through the neural system, modifying the loads appended to these connections between layers, just as a related worth called inclination. This procedure is gotten back to proliferation. In the long run this procedure will choose values for these loads and predispositions that will permit the system to dependably play out a given undertaking, for example, perceiving manually written numbers, and the system can be said to have "realized" how to complete a particular errand

Figure 3.3: How Data are Trained



3.12 Why is machine learning successful

While AI is certainly not another procedure, enthusiasm for the field has detonated lately. This resurgence returns on the of a progression of advancements, with profound getting the hang of establishing new precedents for precision in territories, for example, discourse and language acknowledgment, and PC vision. What's made these victories conceivable are principally two components, one being the immense amounts of pictures, discourse, video and content that is available to specialists hoping to prepare

AI frameworks. Yet, considerably increasingly significant is the accessibility of immense measures of equal handling power, kindness of current designs preparing units (GPUs), which can be connected together into groups to shape AI powerhouses. Today anybody with a web association can utilize these bunches to prepare AI models, by means of cloud administrations gave by firms like Amazon, Google and Microsoft. As the utilization of AI has taken off, so organizations are currently making particular equipment custom-made to running and preparing AI models. A case of one of these custom chips is Google's Tensor Processing Unit (TPU), the most recent form of which quickens the rate at which AI models constructed utilizing Google's TensorFlow programming library can gather data from information, just as the rate at which they can be prepared. These chips are not simply used to prepare models for Google DeepMind and Google Brain, yet additionally the models that support Google Translate and the picture acknowledgment in Google Photo, just as administrations that permit the general population to assemble AI models utilizing Google's TensorFlow Research Cloud. The second era of these chips was uncovered at Google's I/O gathering in May a year ago, with a variety of these new TPUs ready to prepare a Google AI model utilized for interpretation in a fraction of the time it would take a variety of the top-end GPUs, and the as of late declared third-age TPUs ready to quicken preparing and deduction considerably further. As equipment turns out to be progressively particular and AI programming structures are refined, it's getting progressively normal for ML errands to be completed on purchaser grade telephones and PCs, as opposed to in cloud data centres. In the late spring of 2018, Google made a stride towards offering a similar nature of computerized interpretation on telephones that are disconnected as is accessible on the web, by turning out neighbourhood neural machine interpretation for 59 dialects to the Google Translate application for iOS and Android.

3.13 What ways machine learning used for

Machine learning systems are used all around us, and are a cornerstone of the modern internet. Machine-learning systems are used to recommend which product you might want to buy next on Amazon or video you want to may want to watch on Netflix. Every Google search uses multiple machine-

learning systems, to understand the language in your query through to personalizing your results, so fishing enthusiasts searching for "bass" aren't inundated with results about guitars. Similarly Gmail's spam and phishing-recognition systems use machine-learning trained models to keep your inbox clear of rogue messages. One of the most obvious demonstrations of the power of machine learning are virtual assistants, such as Apple's Siri, Amazon's Alexa, the Google Assistant, and Microsoft Cortana. Each relies heavily on machine learning to support their voice recognition and ability to understand natural language, as well as needing an immense corpus to draw upon to answer queries. But beyond these very visible manifestations of machine learning, systems are starting to find a use in just about every industry. These exploitations include: computer vision for driverless cars, drones and delivery robots; speech and language recognition and synthesis for chatbots and service robots; facial recognition for surveillance in countries like China; helping radiologists to pick out tumours in x-rays, aiding researchers in spotting genetic sequences related to diseases and identifying molecules that could lead to more effective drugs in healthcare; allowing for predictive maintenance on infrastructure by analysing IoT sensor data; underpinning the computer vision that makes the cashierless Amazon Go supermarket possible, offering reasonably accurate transcription and translation of speech for business meetings – the list goes on and on.

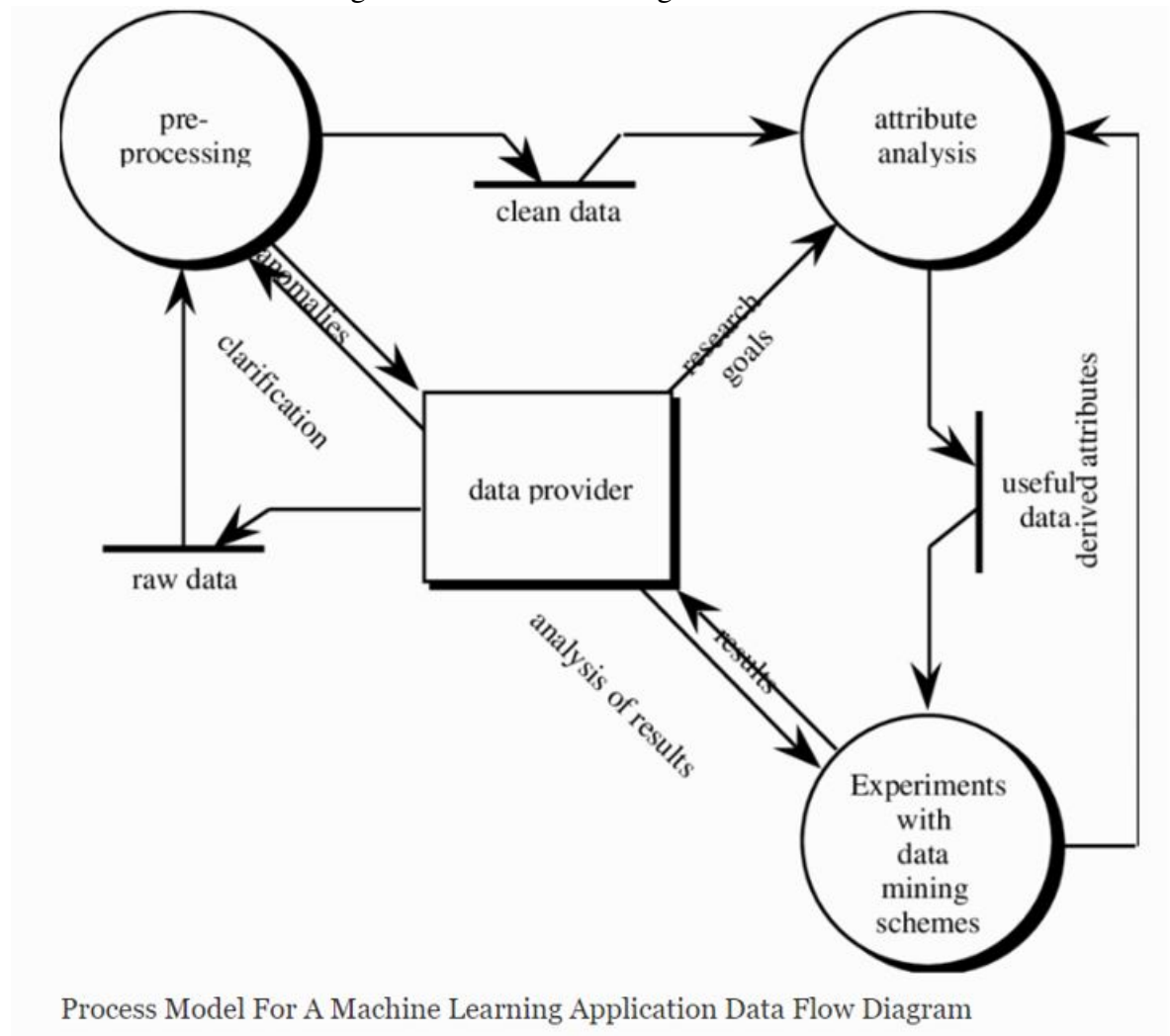
3.14 Which services are available for machine learning

The entirety of the significant cloud stages - Amazon Web Services, Microsoft Azure and Google Cloud Platform - give access to the equipment expected to prepare and run AI models, with Google letting Cloud Platform clients try out its Tensor Processing Units - custom chips whose structure is upgraded for preparing and running AI models. This cloud-based framework incorporates the information stores expected to hold the tremendous measures of preparing information, administrations to set up that information for examination, and perception instruments to show the outcomes plainly. More current administrations even smooth out the making of custom AI models, with Google as of late uncovering an assistance that robotizes the production of AI models, called Cloud AutoML. This simplified

assistance constructs custom picture acknowledgment models and requires the client to have no AI aptitude, like Microsoft's Azure Machine Learning Studio. In a comparable vein, Amazon as of late disclosed new AWS contributions intended to quicken the way toward preparing up AI models. For information researchers, Google's Cloud ML Engine is an overseen AI administration that permits clients to prepare, send and trade custom AI models dependent on Google's publicly released TensorFlow ML system or the open neural system structure Keras, and which currently can be utilized with the Python library sci-unit learn and XGBoost. Database administrators without a foundation in information science can utilize Google's BigQueryML, a beta help that permits administrators to call prepared AI models utilizing SQL orders, permitting forecasts to be made in database, which is less complex than trading information to a different AI and investigation condition. For firms that would prefer not to manufacture their own AI models, the cloud stages likewise offer AI-controlled, on-request benefits -, for example, voice, vision, and language acknowledgment. Microsoft Azure stands apart for the broadness of on-request benefits on offer, firmly followed by Google Cloud Platform and afterward AWS. In the interim IBM, close by its progressively broad on-request contributions, is additionally endeavoring to sell area explicit AI administrations planned for everything from medicinal services to retail, gathering these contributions under its IBM Watson umbrella. From the get-go in 2018, Google extended its AI driven administrations to the universe of publicizing, discharging a set-up of apparatuses for making increasingly compelling promotions, both computerized and physical. While Apple abhors a similar notoriety for front line discourse acknowledgment, normal language preparing and PC vision as Google and Amazon, it is putting resources into improving its AI administrations, as of late placing Google's previous boss accountable for AI and AI technique over the organization, including the advancement of its aide Siri and its on-request AI administration Core ML. In September 2018, NVIDIA propelled a joined equipment and programming stage intended to be introduced in datacenters that can quicken the rate at which prepared AI models can complete voice, video and picture acknowledgment, just as other ML-related administrations. The NVIDIA TensorRT Hyperscale Inference Platform utilizes NVIDIA Tesla T4 GPUs, which conveys up to 40x the exhibition of CPUs when utilizing AI models to make deductions from information, and the TensorRT programming stage, which is intended to streamline the presentation of prepared neural systems.

3.15 How Data are Process

Figure 3.4: Data Flow Diagram



Data preprocessing is a vital advance in Machine Learning as the nature of information and the valuable data that can be gotten from it straightforwardly influences the capacity of our model to learn; in this way, it is critical that we preprocess our information before taking care of it into our model. This two-section investigates the subject of information building

and highlight designing for AI (ML). This initial segment examines best acts of preprocessing information in an AI pipeline on Google Cloud. The article centers around utilizing TensorFlow and the open source TensorFlow Transform (tf.Transform) library to get ready information, train the model, and serve the model for forecast. This part features the difficulties of preprocessing information for AI, and shows the alternatives and situations for performing information change on Google Cloud successfully. Machine learning (ML) helps in consequently discovering complex and conceivably valuable examples in information. These examples are consolidated in a ML model that would then be able to be utilized on new information focuses—a procedure called making forecasts or performing derivation. Building a ML model is a multistep procedure. Each progression presents its own specialized and calculated difficulties. In this two-section set, we center around the way toward choosing, changing, and expanding the source information to make incredible prescient signs to the objective (reaction) variable (in regulated learning undertakings). These tasks consolidate area information with information science methods. They are the quintessence of highlight building. The size of preparing datasets for genuine ML models can without much of a stretch reach or outperform the terabyte (TB) mark. Thus, you need huge scope information preparing systems so as to process these datasets proficiently and distributedly. Besides, when you utilize a ML model for making expectations, you need to apply similar changes that you utilized for the preparation information on the new information focuses, so that the live dataset is introduced to the ML model the way that the model anticipates. The primary article talks about these difficulties for various degrees of granularity of highlight building activities: case level, full-spend, and time-window conglomerations. Also, the article outlines the choices and situations for performing information change for ML on Google Cloud viably. The article likewise gives a review of TensorFlow Transform (tf.Transform), a library for TensorFlow that permits you to characterize both occasion level and full-go information change through information preprocessing pipelines. These pipelines are executed with Apache Beam, and as results they make antiques that let you apply indistinguishable changes during forecast from when the model is served.

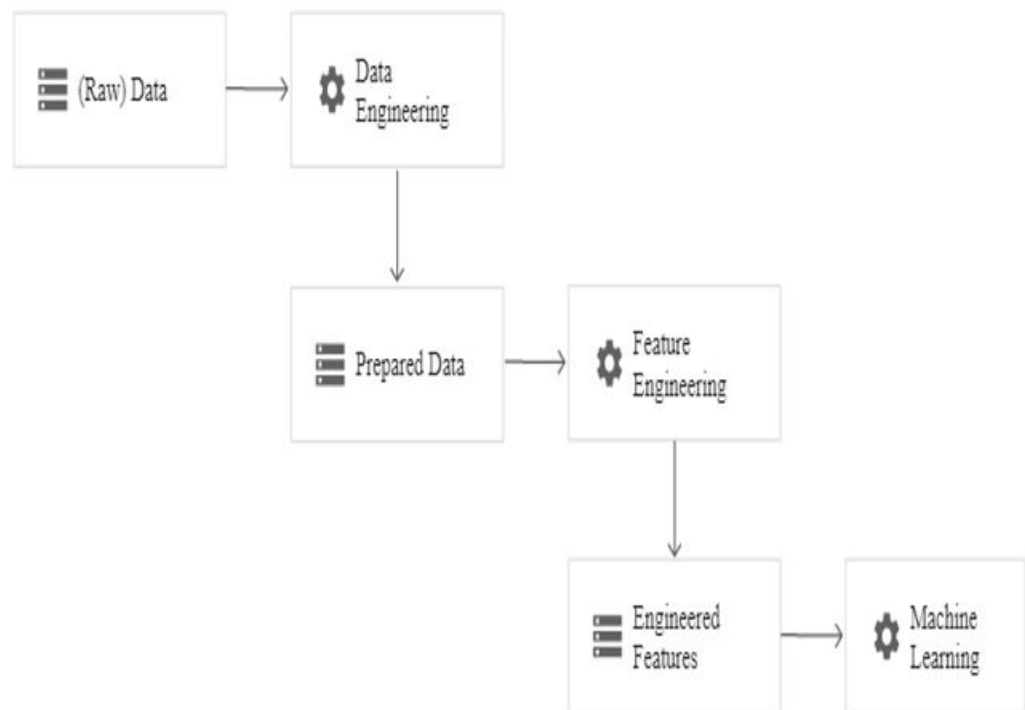
3.16 Pre- Process of Data

Preprocessing the information for ML includes the two information designing and highlight building. Information designing is the way toward changing over crude information into arranged information. Highlight designing at that point tunes the readied information to make the highlights expected by the ML model. These terms have explicit implications, as laid out in the accompanying rundown:

- **Raw Data:** This alludes to the information in its source structure, with no earlier groundwork for ML. Note that in this unique circumstance, the information may be in its crude structure (in an information lake) or in a changed structure (in an information distribution center). Changed information in an information distribution center may have been changed over from its unique crude structure to be utilized for examination, yet in this setting it implies that the information was not arranged explicitly for your ML task. Furthermore, information sent from spilling frameworks that in the end call ML models for forecasts is viewed as information in its crude structure.
- **Prepared Data:** This alludes to the dataset in the structure prepared for your ML task. Information sources have been parsed, joined, and put into a plain structure. Information has been collected and summed up to the correct granularity—for instance, each line in the dataset speaks to a remarkable client, and every section speaks to rundown data for the client, similar to the all out spent over the most recent a month and a half. On account of regulated learning errands, the objective element is available. Unessential sections have been dropped, and invalid records have been sifted through.
- **Engineered Feature:** This alludes to the dataset with the tuned highlights expected by the model—that is, playing out certain ML-explicit procedure on the sections in the readied dataset, and making new highlights for your model during preparing and forecast, as depicted later under Preprocessing activities. Models incorporate scaling numerical sections to an incentive somewhere in the range of 0 and 1, cutting qualities, and one-hot-encoding straight out highlights.

By and by, information from a similar source is regularly at various phases of availability. For instance, a field from a table in your in-

Figure 3.5: Pre-process



formation distribution center could be utilized legitimately as a built element. Simultaneously, another field in a similar table may need to experience changes before turning into a designed element. Also, information building and highlight designing tasks may be joined in similar information preprocessing step.

3.17 Preprocessing Operation

Information preprocessing incorporates different activities. Every activity intends to help AI manufacture better prescient models. For organized information, information preprocessing tasks incorporate the accompanying:

- * **Data Cleaning:** Evacuating or amending records with adulterated or invalid qualities from crude information, just as expelling records that are feeling the loss of an enormous number of sections.
- * **Instance selection and partitioning:** Choosing information focuses from the info dataset to make preparing, assessment (approval), and test sets. This procedure incorporates methods for repeatable irregular examining, minority classes oversampling, and separated apportioning.
- * **feature tuning:** Improving the nature of an element for ML, which incorporates scaling and normalizing numeric qualities, crediting missing qualities, cutting anomalies, and modifying values with slanted circulations.
- * **Representation Transformation:** Changing over a numeric component to an all out element (through bucketization), and changing over clear cut highlights to a numeric portrayal (through one-hot encoding, learning with tallies, scanty element embeddings, etc). A few models work just with numeric or straight out highlights, while others can deal with blended kind highlights. In any event, when models handle the two kinds, they can profit by various portrayal (numeric and unmitigated) of a similar component.
- * **Feature extraction:** Lessening the quantity of highlights by making lower-measurement, all the more impressive information portrayals utilizing strategies, for example, PCA, installing extraction,

and hashing.

- * Feature selection: Choosing a subset of the info highlights for preparing the model, and overlooking the immaterial or repetitive ones, utilizing channel or wrapper techniques. This can likewise include just dropping highlights if the highlights are feeling the loss of an enormous number of qualities.
- * Feature Construction: Making new highlights either by utilizing run of the mill methods, for example, polynomial extension (by utilizing univariate numerical capacities) or highlight crossing (to catch include associations). Highlights can likewise be developed by utilizing business rationale from the space of the ML use case.

3.18 Preprocessing Granularity

This section discusses the granularity of types of data transformations. It shows why this perspective is critical when preparing new data points for predictions using transformations that are applied on training data. Preprocessing and transformation operations can be categorized as follows, based on operation granularity:

Occurrence level changes during preparing and expectation. These are clear changes, where just qualities from a similar occasion (information point) are required for the change. For instance, this may incorporate cut-out the estimation of a component to some limit, polynomially growing another element, duplicating two highlights, or contrasting two highlights with make a Boolean banner. These changes must be applied indistinguishably during preparing and expectation, on the grounds that the model will be prepared on the changed highlights, not on the crude information esteems. On the off chance that the information isn't changed indistinguishably, the model carries on ineffectively in light of the fact that it is given information that has a circulation of qualities that it was not prepared with. For more data, see the conversation of preparing/serving slant in the preprocessing challenges segment later in this report. Full-pass changes during preparing, yet example level changes during expectation. In this situation, changes are stateful, in light of the fact that they need to utilize some pre-computed measurements to play out the change. During preparing,

you have to dissect the entire assemblage of preparing information to register amounts, for example, least, most extreme, mean, and difference for changing preparing information, assessment information, and new information at expectation time. For instance, to standardize a numeric component for preparing, you have to register its mean and its standard deviation over the entire of the preparation information. This is known as a full-pass activity. At that point when you serve the model for forecast, the estimation of another information point must be standardized to abstain from preparing/serving slant Accordingly, mean and SD values that are figured during preparing are utilized to modify the component esteem, which is a straightforward occasion level activity:

Figure 3.6: pre-process granulaity

$$value_{scaled} = (value_{raw} - \mu) \div \sigma$$

Transformations that fall into these categories are:

- MinMax scaling numerical features using *min* and *max* computed from the training dataset.
- Standard scaling (z-score normalization) numerical features using μ and σ computed on the training dataset.
- Bucketizing numerical features using quantiles.
- Imputing missing values using the median (numerical features) or the mode (categorical features).
- Converting strings (nominal values) to integers (indexes) by extracting all of the distinct values (vocabulary) of an input categorical feature.
- Counting the occurrence of a term (feature value) in all of the documents (instances) to calculate for TF-IDF.
- Computing the PCA of the input features to project the data into a lower dimensional space (with linearly dependent features).

3.19 Window Aggregation During Traning and prediction

This methodology includes making an element by summing up continuous qualities after some time. That is, the occurrences to total are characterized through fleeting window statements. For instance, envision that you need to prepare a model that assesses the taxi trip time

dependent on the traffic measurements for the course over the most recent 5 minutes, over the most recent 10 minutes, over the most recent 30 minutes, and at different interims. Another model is foreseeing the disappointment of a motor part dependent on the moving normal of temperature and vibration esteems processed in the course of the most recent 3 minutes. Despite the fact that these conglomerations can be readily disconnected for preparing, they must be registered progressively from an information stream during serving.

All the more decisively, when you are getting ready preparing information, if the collected worth isn't in the crude information, it is made during the information building stage. The crude information is normally put away in a database with configuration of . In the past models, substance would be the course portion identifier for the taxi courses and the motor part identifier for the motor disappointment. You can utilize windowing tasks to process and utilize the conglomeration includes as a contribution for your model preparing.

Be that as it may, when the model for continuous (on the web) expectation is being served, the model expects highlights got from the accumulated qualities as an info. Therefore, you can utilize a stream-handling innovation like Apache Beam to register the conglomerations on the fly from the continuous information focuses spilled into your framework. You can likewise play out extra component building (tuning) to these accumulations before preparing and expectation.

3.20 Machine Learning Pipeline on Google Cloud

This area talks about the structure squares of a common start to finish pipeline to prepare and serve TensorFlow ML models on Google Cloud utilizing oversight administrations. It additionally examines where you can execute various classifications of the information preprocessing activities, just as regular difficulties you may confront when you actualize such changes. Afterward, you will perceive how the TensorFlow Transform library helps address these difficulties. • In the wake of being imported, the crude information is put away in BigQuery (or in Cloud Storage, on account of pictures, docs, sound, video, etc).

Figure 3.7: High Level Architecture

HIGH LEVEL ARCHITECTURE:

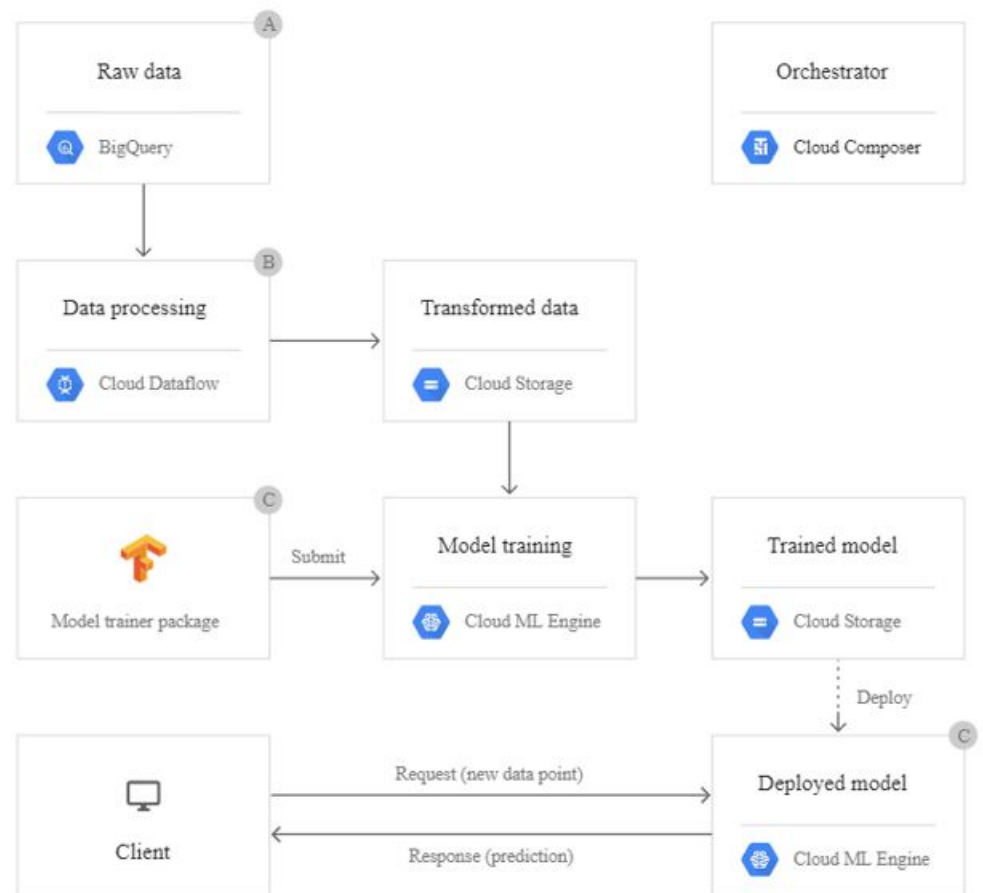


Figure 2

- Information designing (arrangement) and highlight building are executed at scale utilizing Dataflow. This produces ML-prepared preparing, assessment, and test sets that are put away in Cloud Storage. Preferably, these datasets are put away as records, which is the streamlined arrangement for TensorFlow calculations.
- A TensorFlow model mentor bundle is submitted to AI Platform, which utilizes the preprocessed information from the past strides to prepare the model. The yield of this progression is a prepared TensorFlow SavedModel that is sent out to Cloud Storage.
- The prepared TensorFlow model is conveyed to AI Platform as a microservice that has a REST API with the goal that it very well may be utilized for online expectations. A similar model can be utilized additionally for clump expectation employments.
- After the model is conveyed as a REST API, customer applications and inner frameworks can summon this API by sending demands with certain information focuses, and accepting reactions from the model with forecasts.
- For organizing and computerizing this pipeline, Cloud Composer can be utilized as a scheduler to conjure the information arrangement, model preparing, and model sending steps.

3.21 Where To Do Processing

- * BigQuery: BigQuery SQL contents can be utilized as a source question for the Dataflow preprocessing pipeline. This is the information preparing step in Figure 2. For instance, envision that a framework will be utilized in Canada, yet the information distribution center has exchanges from around the globe. Sifting to get Canadian-just preparing information is best done in BigQuery. It is conceivable to execute case level changes, stateful full-pass changes, and window conglomerations include changes in a BigQuery SQL content to set up the preparation information. Notwithstanding, this isn't suggested. On the off chance that you are conveying the model for online expectations, you need to repeat the SQL preprocessing procedure on the crude information focuses that you produced from differ-

ent frameworks and that will be sent model's API. At the end of the day, you have to execute the rationale twice. The first run through is in SQL to preprocess preparing information in BigQuery. The other time is in the rationale of the application that devours the model to preprocess online information focuses for expectation. For instance, if your customer application is written in Java, you have to reimplement the rationale in Java. This can acquaint blunders due with usage inconsistencies. (See the conversation of preparing/serving slant in the preprocessing challenges segment later in this record.) What's more, keeping up two unique executions is additional overhead. At whatever point you change the rationale in SQL to preprocess the preparation information, you have to change the Java usage in like manner to preprocess information at spending time in jail.

On the off chance that you are utilizing your model just for clump forecast (scoring) utilizing AI Platform group expectation, and if your information for scoring is sourced from BigQuery, it is possible to execute these preprocessing tasks as a component of the BigQuery SQL content. All things considered, you can utilize the equivalent preprocessing SQL content to get ready both preparing and scoring information. This suggests you are utilizing assistant tables to store amounts required by stateful changes, for example, means and differences to scale numerical highlights. It additionally implies expanded multifaceted nature in the SQL contents, and complex reliance among preparing and the scoring SQL contents. The accompanying code bits show theoretical instances of BigQuery SQL for information groundwork for preparing and forecast. In the main content (the preparation content), the mean and the standard deviation for fields f1 and f2 are put away in the details table. At that point you get the training data table together with the details table to change your preparation information.

- * Data Flow: you can execute computationally costly preprocessing activities in Apache Beam, and run them at scale utilizing Dataflow, which is a completely overseen autoscaling administration for cluster and stream information preparing. Dataflow can perform example level changes, stateful full pass changes, and window total element changes. Specifically, if your ML mod-

els expect an information include like total number of clicks last 90sec, Apache Beam windowing capacities can figure these highlights dependent on conglomerating the estimations of time windows of constant (gushing) occasions information (for instance, clicks). In the previous conversation of granularity of changes, this was alluded to as "Window collections during preparing and serving." Figure 3 shows the job of Dataflow in preparing stream information for close to constant forecasts. Basically, occasions (information focuses) are ingested into Pub/Sub. Dataflow expends these information focuses, registers highlights dependent on totals after some time, and calls the conveyed ML model API for expectations. Forecasts are then sent to an outbound Pub/Sub line. From that point, the expectations can be devoured by downstream (checking or control) frameworks or pushed back (for instance, as warnings) to the first mentioning customer. Another choice is to store the expectations in a low-idleness information store like Cloud Bigtable for continuous getting. Cloud Bigtable can likewise be utilized to amass and store these constant totals so they can be looked into when required for forecast. Figure 3 shows this situation. A similar Apache Beam execution can be utilized to bunch process preparing information that originates from a disconnected datastore like BigQuery and stream-process constant information for serving on the web expectations.

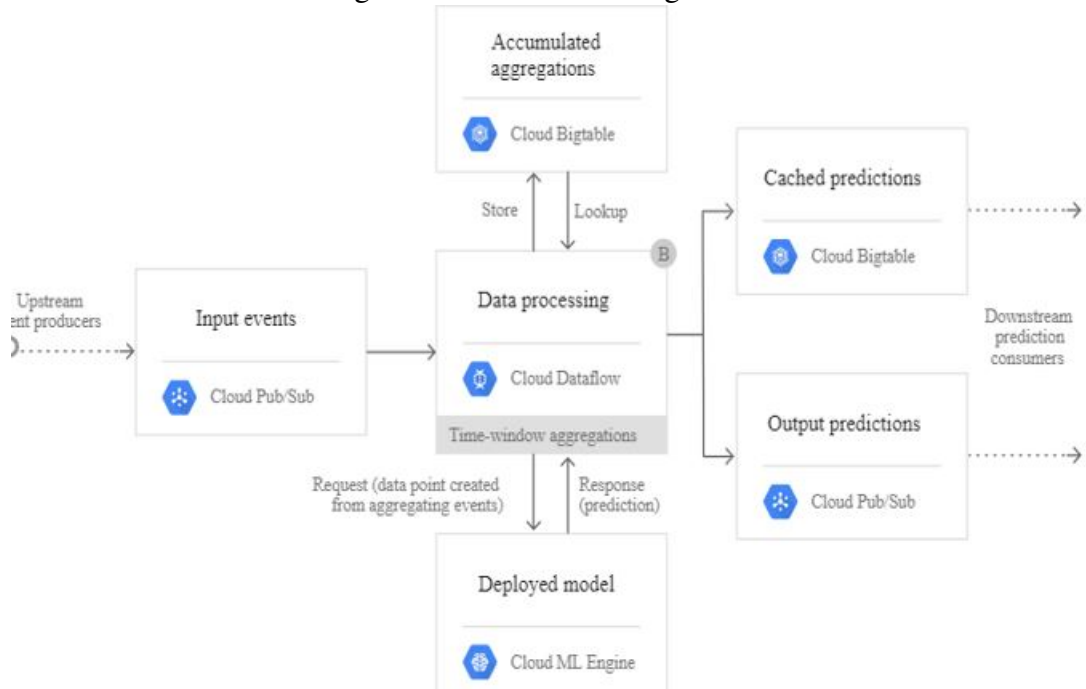
- * TensorFlow: As should be obvious in Figure 2, you can execute information preprocessing and change activities in the TensorFlow model itself. As appeared in the figure, the preprocessing you execute for preparing the TensorFlow model turns into a necessary piece of the model when the model is sent out and conveyed for forecasts. TensorFlow changes can be practiced in one of the accompanying ways:

Broadening your base feature columns (utilizing crossed column, embedding column, bucketized column, etc).

Executing the entirety of the occurrence level change rationale in a capacity that you bring in each of the three information capacities: train input fn, eval input fn, and serving input fn.

On the off chance that you are making custom estimators, placing the code in the model fn work.

Figure 3.8: Data Flow Diagram



On the off chance that you utilize a similar change rationale code in the serving input fn work, which characterizes the serving interface of your SavedModel for online forecast, it guarantees that similar changes utilized for getting ready preparing information will be applied on new expectation information focuses during serving.

Figure 3.9: Transformation Training and Evaluation Data

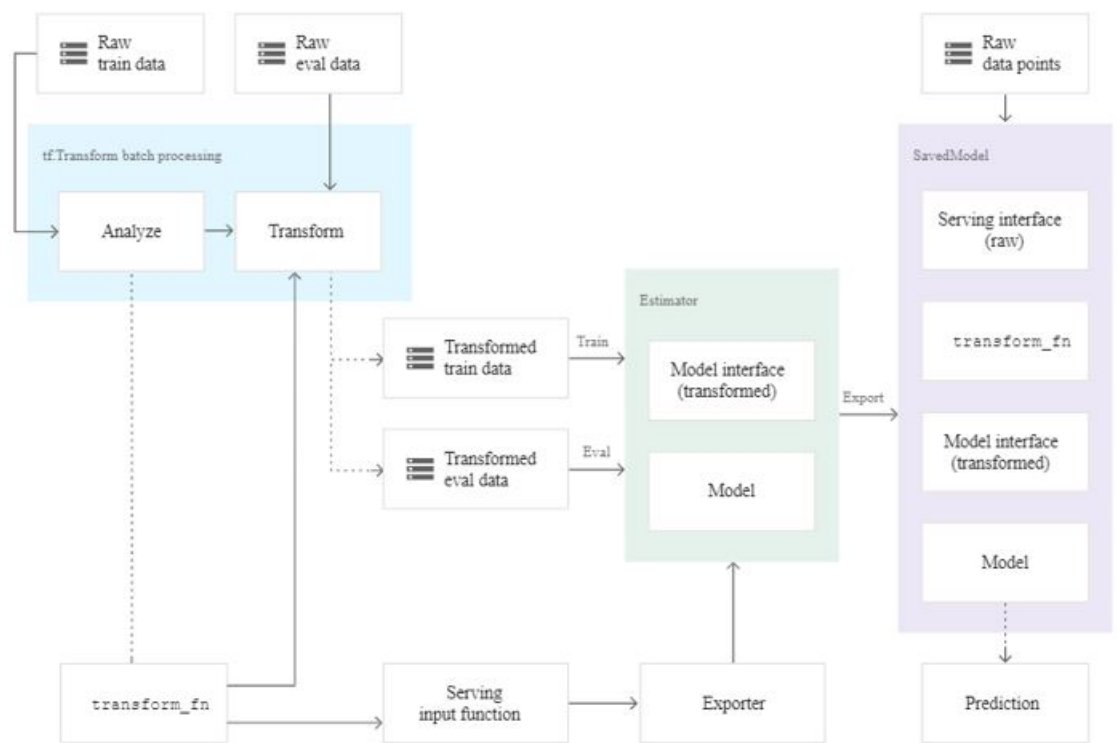


Figure 4

TRANSFORMATION TRAINING AND EVALUATION DATA:

Figure 3.10: Summary Table

	Instance-level (stateless transformations)		Full-pass during training instance-level during serving (stateful transformations)		Real-time (window) aggregations during training and serving (streaming transformations)	
	Batch scoring	Online prediction	Batch scoring	Online prediction	Batch scoring	Online prediction
BigQuery (SQL)	OK The same transformation implementation is applied on data during training and batch scoring.	Possible to process training data but not recommended Results in training/serving skew, because you process serving data using different tools.	Possible to use statistics computed using BigQuery for instance-level batch/online transformations. Not easy; you must maintain a stats store to be populated during training and used during prediction.	N/A Aggregates like these computed based on real-time events.	Possible to process training data but not recommended. Results in training/serving skew, because you process serving data using different tools.	
Dataflow (Apache Beam)		OK if data at serving time comes from Pub/Sub to be consumed by Dataflow. Otherwise, results in training/serving skew.	Possible to use statistics computed using Dataflow for instance-level batch/online transformations. Not easy; you must maintain a stats store to be populated during training and used during prediction.		OK The same Apache Beam transformation is applied on data during training (batch) and serving (stream).	
Dataflow (Apache Beam + TFT)	OK The same transformation implementation is applied to data during training and batch scoring.	Recommended Avoids training/serving skew and prepares training data up front.	Recommended Transformation logic + computed statistics during training are stored as a tf.Graph that's attached to the exported model for serving.			
TensorFlow* (input_fn & serving_input_fn)	Possible but not recommended For training & prediction efficiency, it's better to prepare the training data up front.	Possible but not recommended For training efficiency, it's better to prepare the training data up front.	Not Possible		Not Possible	

Chapter 4

Implementation of System/ Methodology

4.1 Random Forest Algorithm

- What is Random Forest Algorithm ?
- Random Forest, similar to its name infers, comprises of an enormous number of individual choice trees that work as an outfit. Every individual tree in the arbitrary timberland lets out a class expectation and the class with the most votes turns into our model's forecast (see figure beneath).

Perception of a Random Forest Model Making a Prediction

The central idea driving arbitrary woodland is a basic however amazing one — the shrewdness of groups. In information science talk, the explanation that the irregular timberland model works so well is:

An enormous number of generally uncorrelated models (trees) working as a council will outflank any of the individual constituent models.

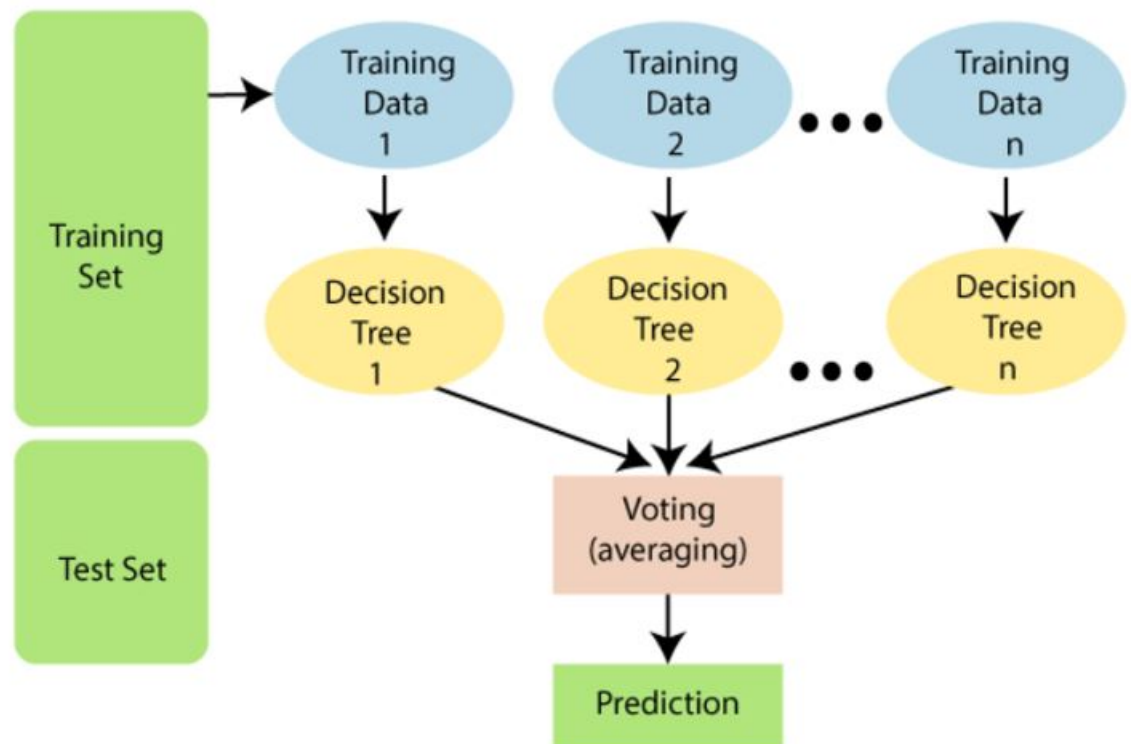
The low relationship between's models is the key. Much the same as how speculations with low connections (like stocks and securities) meet up to shape a portfolio that is more noteworthy than the aggregate of its parts, uncorrelated models can deliver outfit forecasts that are more exact than any of the individual expectations. The explanation behind this superb impact is that the trees shield each other from their individual mistakes (as long as they don't continually all blunder a similar way). While a few trees might

not be right, numerous different trees will be correct, so as a gathering the trees can move in the right course. So the essentials for arbitrary woodland to perform well are:

- There needs to be some actual signal in our features so that models built using those features do better than random guessing.
- The predictions (and therefore the errors) made by the individual trees need to have low correlations with each other. The distinction between Random Forest calculation and the choice tree calculation is that in Random Forest, the procedure es of finding the root hub and parting the component hubs will run arbitrarily.

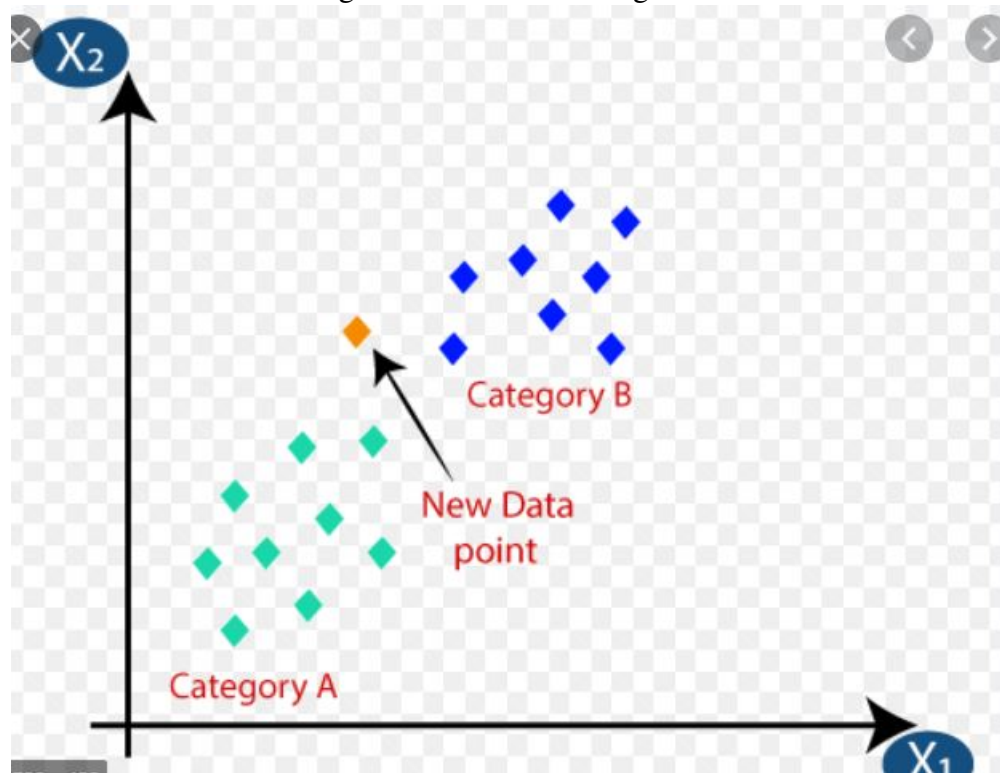
Figure 4.1: Random Forest Algorithm

The below diagram explains the working of the Random Forest algorithm:



4.2 KNN algorithm

Figure 4.2: K-nearest neighbor



order is one of the most key and basic grouping strategies and ought to be one of the main decisions for a characterization study when there is almost no earlier information about the dissemination of the information. K-closest neighbor order was created from the need to perform discriminant examination when solid parametric appraisals of likelihood densities are obscure or hard to decide. effortlessness, keeps on performing genuinely well for enormous preparing sets. It basically depends just on the most essential presumption hidden all forecast: that perceptions with comparative attributes will in general have comparable results. Closest Neighbor strategies allot an anticipated an incentive to another perception dependent on the majority or mean (here and there weighted) of its k "Closest Neighbors" in the preparation set. Given a vast measure of information, any perception will have many "neighbors" that are discretionarily close concerning every

single estimated trademark, and the inconstancy of their results will give as exact an expectation as is hypothetically conceivable excepting a consummately and totally determined model. Be that as it may, given that we never have a boundless measure of information, the genuine utility of this asymptotic property is flawed, particularly for humble datasets. Tragically, on the grounds that expectations depend entirely on an assortment of put away perceptions, it is computationally and memory-serious and touchy to the scourge of dimensionality. See the paper by Friedman for an increasingly complete conversation (Freidman, 1994). - Nearest Neighbors (KNN) is a standard AI strategy that has been stretched out to enormous scope information mining endeavors. The thought is that one uses a lot of preparing information, where every datum point is described by a lot of factors. Theoretically, each point is plotted in a high-dimensional space, where every hub in the space compares to an individual variable. At the point when we have another (test) information point, we need to discover the K closest neighbors that are nearest (ie, generally "comparable" to it). The number K is regularly picked as the square foundation of N, the absolute number of focuses in the preparation informational collection. (Along these lines, if N is 400, $K = 20$).

KNN is thoughtfully basic and has the benefit of being nonparametric. That is, the strategy can be utilized in any event, when the factors are straight out—however on the off chance that you are utilizing numeric factors in the blend, it is ideal to normalize them (see Section 4.4 of part: Core Technologies: Machine Learning and Natural Language Processing) to take out contrasts in scale. The test is that when the quantity of information focuses is enormous (eg, an online book retailer has a huge number of books), uncommon techniques must be utilized to quickly look through the space and locate the "most comparative" things.

Normally, some type of precomputation is utilized for instance, ordering. What's more, as opposed to utilizing all the information focuses, chose information focuses that are illustrative of individual groups ("models") might be utilized to encourage the inquiry against another thing, and afterward the precomputed neighbors of the most comparative model are additionally shown. Correspondingly, endeavoring to decrease the quantity of measurements with a technique like SVD/LSI and afterward plotting the information focuses in the diminished variable space may bring about huge gains in execution.

4.3 Time Series

A time series is a sequence of numerical data points in successive order. In investing, a time series tracks the movement of the chosen data points, such as a security's price, over a specified period of time with data points recorded at regular intervals. There is no minimum or maximum amount of time that must be included, allowing the data to be gathered in a way that provides the information being sought by the investor or analyst examining the activity. Time series examination can be helpful to perceive how a given resource, security, or financial variable changes after some time. It can likewise be utilized to analyze how the progressions related with the picked information direct look at toward shifts in different factors over a similar timespan.

For instance, assume you needed to examine a period arrangement of day by day shutting stock costs for a given stock over a time of one year. You would acquire a rundown of all the end costs for the stock from every day for as far back as year and show them in sequential request. This would be a one-year every day shutting value time arrangement for the stock.

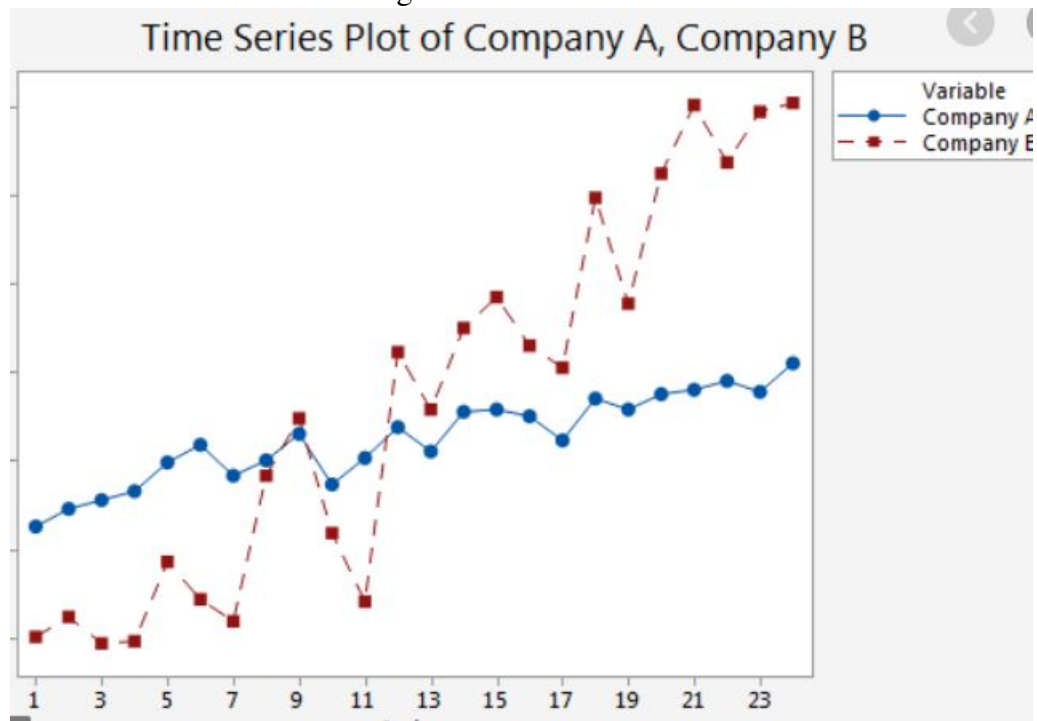
Digging somewhat more profound, you may break down time arrangement information with specialized investigation apparatuses to know whether the stock's time arrangement shows any regularity. This will assist with deciding whether the stock experiences pinnacles and troughs at standard occasions every year. Examination around there would require taking the watched costs and relating them to a picked season. This can incorporate customary schedule seasons, for example, summer and winter, or retail seasons, for example, special seasons.

On the other hand, you can record a stock's offer value changes as it identifies with a monetary variable, for example, the joblessness rate. By corresponding the information focuses with data identifying with the chose monetary variable, you can watch designs in circumstances showing reliance between the information focuses and the picked variable.

4.4 What is time series and how it is used?

Arrangement is a grouping of information focuses in ordered succession, frequently assembled in standard interims. Time arrangement examination can be applied to any factor that changes after some time and as a rule, normally information focuses that are nearer together are more comparable than those further separated.

Figure 4.3: Time Series



4.5 components

- **LEVEL:** When you read about the "level" or the "level record" of time arrangement information, it's alluding to the mean of the arrangement.
- **NOISE:** All time arrangement information will have clamor or irregularity in the information focuses that aren't corresponded with any clarified patterns. Clamor is unsystematic and is present moment.
- **SEASONALITY:** If there are standard and unsurprising changes in the arrangement that are connected with the schedule – could be quarterly, week by week, or even days of the week, at that point the arrangement incorporates a regularity part. Note that regularity is space explicit, for instance land deals are typically higher in the late spring months versus the winter months while standard retail as a rule tops during the year's end. Likewise, not unsurpassed arrangement have an occasional part, as referenced for sound or video information.
- **TREND:** When alluding to the "pattern" in time arrangement information, it implies that the information has a drawn out direction which can either be drifting the positive or negative way. A case of a pattern would be a drawn out increment in an organization's business information or system use.
- **Cycle:** Repeating periods that are not identified with the schedule. This incorporates business cycles, for example, monetary downturns or developments or salmon run cycles, or even sound documents which have cycles, yet aren't identified with the schedule in the week by week, month to month, or yearly sense.

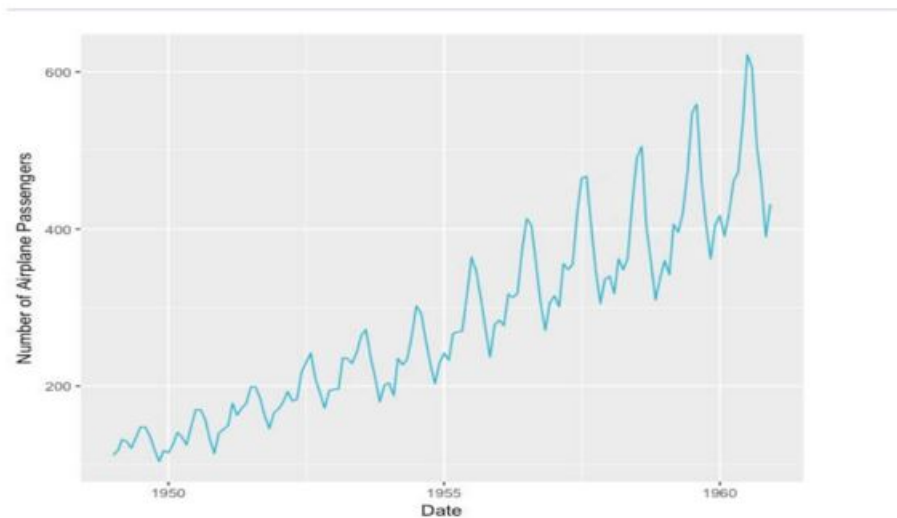
4.6 Time Series Models

At the point when an arrangement contains a pattern, regularity, and clamor, at that point you can characterize that arrangement by the manner in which those parts connect with one another. These communications can be diminished to what is called either a multiplicative or added substance time arrangement.

A multiplicative time arrangement is the point at which the changes in the time arrangement increment after some time and is subject to the degree of

the arrangement: Multiplicative Model: Time series = t (trend) * s (seasonality) * n (noise) Consequently, the regularity of the model would increment with the level after some time. In the diagram underneath, you can see that the regularity of plane travelers increments as the level increments: Addi-

Figure 4.4: Additive Model



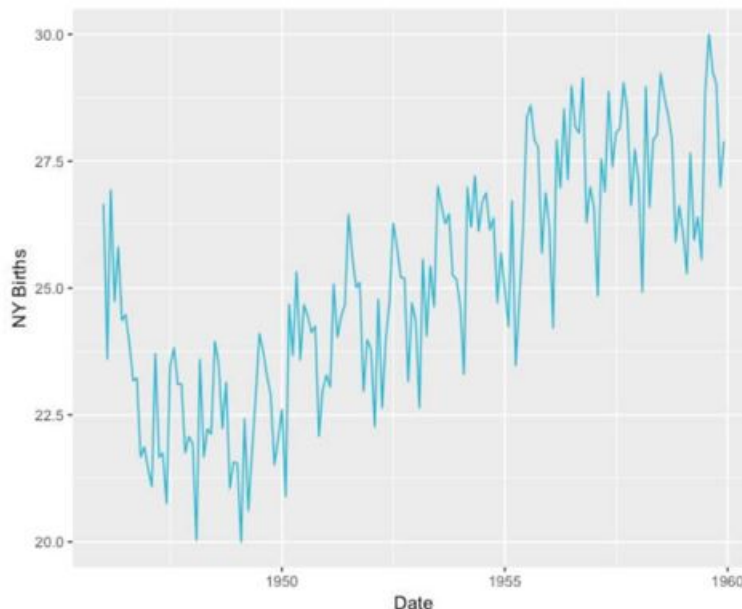
An additive model is when the fluctuations in the time series stay constant over time:

tive Model: Time series = t (trend) + s (seasonality) + n (noise) So an added substance model's regularity ought to be steady from year to year and not identified with the expansion or decline in the level after some time. The following is a diagram of births in New York where you can see that in spite of the fact that the level builds, the regularity remains the equivalent:

4.7 Decomposition

Decay is the deconstruction of the arrangement information into its different segments: pattern, cycle, commotion, and regularity when those exist. Two distinct sorts of great decay incorporate multiplicative and added substance deterioration.

Figure 4.5: Decomposition

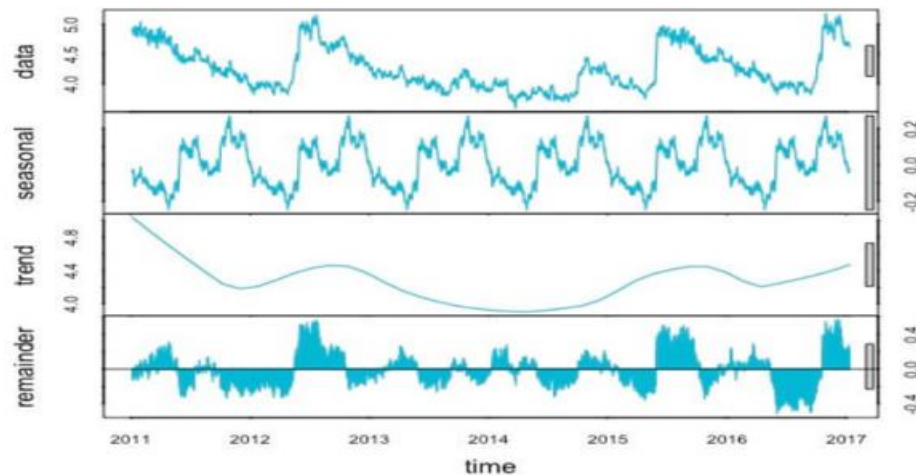


The reason for decay is to detach the different segments so you can see them each independently and perform investigation or gauging without the impact of clamour or regularity. For instance, on the off chance that you needed to just view the pattern of a land arrangement, you would need to expel the regularity found in the information, the commotion because of haphazardness, and any cycles, for example, monetary development. The following is a figure indicating the different parts of the home loan time arrangement examination including the first information, the regularity, pattern, and clamour or "leftover portion":

4.8 Moving Average

We previously discussed how there is arbitrariness, or clamor in our information, and how to isolate it with deterioration, however once in a while we just need to diminish the commotion in our information and smooth out the vacillations from the commotion so as to more readily gauge future information focuses. A moving normal model use the normal of the information focuses that exist in a particu-

Figure 4.6: Decomposition



lar covering subsection of the arrangement. A normal is taken from the principal subset of the information, and afterward it is pushed ahead to the following information point while dropping out the underlying information point. A moving normal can give you data about the present patterns, and decrease the measure of clamor in your information. Regularly, it is a preprocessing step for guaging.

Moving midpoints are utilized by financial specialists and brokers for investigating momentary patterns in securities exchange information, while SMA's are being utilized in social insurance to all the more likely comprehend current patterns in medical procedures, and even break down quality control in medicinal services suppliers.

While there are a wide range of moving normal models, we'll spread the basic moving normal (SMA), autoregressive incorporated moving normal (ARIMA), and exponential smoothing.

4.9 CODE

```
In [1]: import numpy as np
import pandas as pd
from sklearn.preprocessing import LabelEncoder
from sklearn import model_selection
pd.__version__
```

```
In [2]: file='311_Service_Requests_from_2010_to_Present.csv'
        f_read=pd.read_csv(file,sep=',', error_bad_lines=False, index_col=False, dtype='unicode')

In [3]: f_read.shape
Out[3]: (364558, 53)
```

```
In [4]: f_read.head()
Out[4]:
```

	Unique Key	Created Date	Closed Date	Agency	Agency Name	Complaint Type	Descriptor	Location Type	Incident Zip	Incident Address	...	Bridge Highway Name	Bridge Highway Direction	Road Ramp	Bric Highway Segment
0	32310363	12/31/2015 23:59	1/1/2016 0:55	NYPD	New York City Police Department	Noise - Street/Sidewalk	Loud Music/Party	Street/Sidewalk	10034	VERMILYEA AVENUE	...	NaN	NaN	NaN	N
1	32309934	12/31/2015 23:59	1/1/2016 1:26	NYPD	New York City Police Department	Blocked Driveway	No Access	Street/Sidewalk	11105	27-07 23 AVENUE	...	NaN	NaN	NaN	N
2	32309159	12/31/2015 23:59	1/1/2016 4:51	NYPD	New York City Police Department	Blocked Driveway	No Access	Street/Sidewalk	10458	2897 VALENTINE AVENUE	...	NaN	NaN	NaN	N
3	32305098	12/31/2015 23:57	1/1/2016 7:43	NYPD	New York City Police Department	Illegal Parking	Commercial Overnight Parking	Street/Sidewalk	10461	2940 BAISLEY AVENUE	...	NaN	NaN	NaN	N
4	32306529	12/31/2015 23:56	1/1/2016 3:24	NYPD	New York City Police Department	Illegal Parking	Blocked Sidewalk	Street/Sidewalk	11373	87-14 57 ROAD	...	NaN	NaN	NaN	N

5 rows x 53 columns

This data is categorical Data

```
In [5]: column_list = f_read.columns

In [6]: new_col_list= []
        for i in column_list:
            new_col_list.append(i.replace(" ", "_"))

In [7]: ## This will change column name having space
        f_read.columns = new_col_list

In [8]: f_read.columns
Out[8]: Index(['Unique_Key', 'Created_Date', 'Closed_Date', 'Agency', 'Agency_Name',
               'Complaint_Type', 'Descriptor', 'Location_Type', 'Incident_Zip',
               'Incident_Address', 'Street_Name', 'Cross_Street_1', 'Cross_Street_2',
               'Intersection_Street_1', 'Intersection_Street_2', 'Address_Type',
               'City', 'Landmark', 'Facility_Type', 'Status', 'Due_Date',
               'Resolution_Description', 'Resolution_Action_Updated_Date',
               'Community_Board', 'Borough', 'X_Coordinate_(State_Plane)',
               'Y_Coordinate_(State_Plane)', 'Park_Facility_Name', 'Park_Borough',
               'School_Name', 'School_Number', 'School_Region', 'School_Code',
               'School_Phone_Number', 'School_Address', 'School_City', 'School_State',
               'School_Zip', 'School_Not_Found', 'School_or_Citywide_Complaint',
               'Vehicle_Type', 'Taxi_Company_Borough', 'Taxi_Pick_Up_Location',
               'Bridge_Highway_Name', 'Bridge_Highway_Direction', 'Road_Ramp',
               'Bridge_Highway_Segment', 'Garage_Lot_Name', 'Ferry_Direction',
               'Ferry_Terminal_Name', 'Latitude', 'Longitude', 'Location'],
              dtype='object')
```

```

In [9]: f_read["Address_Type"].unique()
Out[9]: array(['ADDRESS', nan, 'INTERSECTION', 'LATLONG', 'BLOCKFACE',
              'PLACENAME'], dtype=object)

In [10]: complaint_list = f_read['Complaint_Type'].unique()
In [11]: complaint_list
Out[11]: array(['Noise - Street/Sidewalk', 'Blocked Driveway', 'Illegal Parking',
              'Derelict Vehicle', 'Noise - Commercial',
              'Noise - House of Worship', 'Posting Advertisement',
              'Noise - Vehicle', 'Animal Abuse', 'Vending', 'Traffic',
              'Drinking', 'Bike/Roller/Skate Chronic', 'Panhandling',
              'Noise - Park', 'Homeless Encampment', 'Urinating in Public',
              'Graffiti', 'Disorderly Youth', 'Illegal Fireworks',
              'Ferry Complaint', 'Agency Issues', 'Squeegee', 'Animal in a Park'],
              dtype=object)

In [12]: f_read.head()
Out[12]:

```

	Unique_Key	Created_Date	Closed_Date	Agency	Agency_Name	Complaint_Type	Descriptor	Location_Type	Incident_Zip	Incident_Address	...	Bridge_Hig
0	32310363	12/31/2015 23:59	1/1/2016 0:55	NYPD	New York City Police Department	Noise - Street/Sidewalk	Loud Music/Party	Street/Sidewalk	10034	71 VERMILYEA AVENUE	...	
1	32309934	12/31/2015 23:59	1/1/2016 1:26	NYPD	New York City Police Department	Blocked Driveway	No Access	Street/Sidewalk	11105	27-07 23 AVENUE	...	
2	32309159	12/31/2015 23:59	1/1/2016 4:51	NYPD	New York City Police Department	Blocked Driveway	No Access	Street/Sidewalk	10458	2897 VALENTINE AVENUE	...	
3	32305098	12/31/2015 23:57	1/1/2016 7:43	NYPD	New York City Police Department	Illegal Parking	Commercial Overnight Parking	Street/Sidewalk	10461	2940 BAISLEY AVENUE	...	
4	32306529	12/31/2015 23:56	1/1/2016 3:24	NYPD	New York City Police Department	Illegal Parking	Blocked Sidewalk	Street/Sidewalk	11373	87-14 57 ROAD	...	

5 rows × 53 columns

```

In [13]: f_read.isnull().sum()
Out[13]:
Unique_Key          0
Created_Date         0
Closed_Date        2381
Agency              0
Agency_Name         0
Complaint_Type       0
Descriptor          6501
Location_Type        133
Incident_Zip        2998
Incident_Address    51699
Street_Name         51699
Cross_Street_1      57188
Cross_Street_2      57805
Intersection_Street_1  313438
Intersection_Street_2  314046
Address_Type         3252
City                2997
Landmark            364183
Facility_Type       2389
Status              0
Due_Date            3
Resolution_Description 0
Resolution_Action_Updated_Date 2402
Community_Board      0
Borough             0
X_Coordinate_(State_Plane) 4030
Y_Coordinate_(State_Plane) 4030
Park_Facility_Name   0
Park_Borough         0
School_Name          0
School_Number        0
School_Region        1
School_Code          1
School_Phone_Number  0
School_Address       0

```

```
In [14]: city_complaint_type=f_read.groupby(['City', 'Complaint_Type'])
```

Display the city type and complain Together

```
In [15]: city_complaint_type.size()
```

```
Out[15]: City      Complaint_Type      freq
ARVERNE  Animal Abuse      46
          Blocked Driveway  50
          Derelict Vehicle  32
          Disorderly Youth  2
          Drinking          1
          Graffiti         1
          Homeless Encampment 4
          Illegal Parking   62
          Noise - Commercial 2
          Noise - House of Worship 14
          Noise - Park      2
          Noise - Street/Sidewalk 29
          Noise - Vehicle   10
          Panhandling       1
          Traffic           1
          Urinating in Public 1
          Vending           1
ASTORIA  Animal Abuse      170
          Bike/Roller/Skate Chronic 16
          Blocked Driveway  3436
          Derelict Vehicle  426
          Disorderly Youth  5
          Drinking          43
          Graffiti         4
```

```
In [16]: complaint_type = f_read.groupby('Complaint_Type')
```

```
In [17]: complaint_analysis=complaint_type.size()
df=complaint_analysis.to_frame().reset_index()
```

```
In [18]: df.columns=['Complaint_Type', 'FREQ']
df.columns
```

```
Out[18]: Index(['Complaint_Type', 'FREQ'], dtype='object')
```

```
In [19]: sort_1=df.sort_values('FREQ',ascending=False)
sort_1
```

```
Out[19]:
```

	Complaint_Type	FREQ
4	Blocked Driveway	100881
12	Illegal Parking	92679
16	Noise - Street/Sidewalk	51692
13	Noise - Commercial	44109
5	Derelict Vehicle	21661
17	Noise - Vehicle	19352
1	Animal Abuse	10541
21	Traffic	5198
10	Homeless Encampment	4879
23	Vending	4192
15	Noise - Park	4109
7	Drinking	1409
14	Noise - House of Worship	1070
19	Posting Advertisement	681
22	Urinating in Public	641
3	Bike/Roller/Skate Chronic	478
18	Panhandling	327
6	Disorderly Youth	315

```
In [20]: top10=sort_1.head(10)
top10
```

```
Out[20]:
```

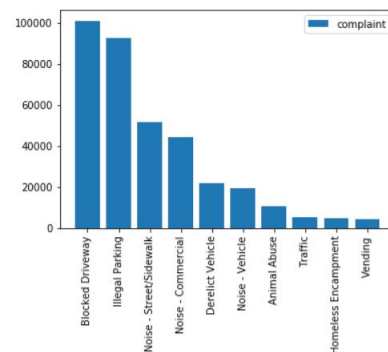
	Complaint_Type	FREQ
4	Blocked Driveway	100881
12	Illegal Parking	92679
16	Noise - Street/Sidewalk	51692
13	Noise - Commercial	44109
5	Derelect Vehicle	21661
17	Noise - Vehicle	19352
1	Animal Abuse	10541
21	Traffic	5198
10	Homeless Encampment	4879
23	Vending	4192

Visualize the complaint types

Display the major complaint types and their count

```
In [21]: import matplotlib.pyplot as plt
%matplotlib inline
```

```
In [22]: x=range(10)
plt.bar(top10.Complaint_Type, top10.FREQ)
plt.xticks(x, top10.Complaint_Type, rotation='vertical')
plt.legend(['complaint'])
plt.show()
```



```
In [23]: f_read.dtypes
```

Out[23]:	Unique_Key	object
	Created_Date	object
	Closed_Date	object
	Agency	object
	Agency_Name	object
	Complaint_Type	object
	Descriptor	object
	Location_Type	object
	Incident_Zip	object
	Incident_Address	object
	Street_Name	object
	Cross_Street_1	object
	Cross_Street_2	object
	Intersection_Street_1	object
	Intersection_Street_2	object
	Address_Type	object
	City	object
	Landmark	object
	Facility_Type	object
	Status	object
	Due_Date	object
	Resolution_Description	object
	Resolution_Action_Updated_Date	object
	Community_Board	object
	Borough	object
	X_Coordinate_(State_Plane)	object
	Y_Coordinate_(State_Plane)	object
	Park_Facility_Name	object
	Park_Borough	object
	School_Name	object
	School_Number	object
	School_Region	object
	School_Code	object
	School Phone Number	object

converting into date format

```
1 [24]: f_read['Created_Date'] = pd.to_datetime(f_read.Created_Date )
f_read['Closed_Date'] = pd.to_datetime(f_read.Closed_Date )
f_read['Due_Date'] = pd.to_datetime(f_read.Due_Date)
f_read['Resolution_Action_Updated_Date'] = pd.to_datetime(f_read.Resolution_Action_Updated_Date)

1 [25]: f_read['Call_closing_time'] = (f_read['Closed_Date'] - f_read['Created_Date']).dt.seconds/3600
f_read['Resolution_time'] = (f_read['Resolution_Action_Updated_Date'] - f_read['Created_Date']).dt.seconds/3600
```

```
In [26]: f_read.isnull().sum()
```

Out[26]:	Unique_Key	0
	Created_Date	0
	Closed_Date	2381
	Agency	0
	Agency_Name	0
	Complaint_Type	0
	Descriptor	6501
	Location_Type	133
	Incident_Zip	2998
	Incident_Address	51699
	Street_Name	51699
	Cross_Street_1	57188
	Cross_Street_2	57805
	Intersection_Street_1	313438
	Intersection_Street_2	314046
	Address_Type	3252
	City	2997
	Landmark	364183
	Facility_Type	2389
	Status	0

```
In [27]: f_read["Address_Type"].unique()
f_read["Address_Type"] = f_read["Address_Type"].mode()[0]
```

```
In [28]: f_read['City']=f_read['City'].fillna('Unknown')
a=f_read['Facility_Type'].mode()[0]
f_read['Facility_Type']=f_read['Facility_Type'].fillna(a)
f_read['Location_Type']=f_read['Location_Type'].fillna(f_read['Facility_Type'].mode()[0])
```



```
In [29]: f_read.head()
```

```
Out[29]:
```

	Unique_Key	Created_Date	Closed_Date	Agency	Agency_Name	Complaint_Type	Descriptor	Location_Type	Incident_Zip	Incident_Address	...
0	32310363	2015-12-31 23:59:00	2016-01-01 00:55:00	NYPD	New York City Police Department	Noise - Street/Sidewalk	Loud Music/Party	Street/Sidewalk	10034	71 VERMILYEA AVENUE	...
1	32309934	2015-12-31 23:59:00	2016-01-01 01:26:00	NYPD	New York City Police Department	Blocked Driveway	No Access	Street/Sidewalk	11105	27-07 23 AVENUE	...
2	32309159	2015-12-31 23:59:00	2016-01-01 04:51:00	NYPD	New York City Police Department	Blocked Driveway	No Access	Street/Sidewalk	10458	2897 VALENTINE AVENUE	...
3	32305098	2015-12-31 23:57:00	2016-01-01 07:43:00	NYPD	New York City Police Department	Illegal Parking	Commercial Overnight Parking	Street/Sidewalk	10461	2940 BAISLEY AVENUE	...
4	32306529	2015-12-31 23:56:00	2016-01-01 03:24:00	NYPD	New York City Police Department	Illegal Parking	Blocked Sidewalk	Street/Sidewalk	11373	87-14 57 ROAD	...

5 rows × 56 columns

```
In [30]: f_read.columns
```

```
Out[30]:
```

```
Index(['Unique_Key', 'Created_Date', 'Closed_Date', 'Agency', 'Agency_Name',
      'Complaint_Type', 'Descriptor', 'Location_Type', 'Incident_Zip',
      'Incident_Address', 'Street_Name', 'Cross_Street_1', 'Cross_Street_2',
      'Intersection_Street_1', 'Intersection_Street_2', 'Address_Type',
      'City', 'Landmark', 'Facility_Type', 'Status', 'Due_Date',
      'Resolution_Description', 'Resolution_Action_Updated_Date',
      'Community_Board', 'Borough', 'X_Coordinate (State Plane)',
      'Y_Coordinate (State Plane)', 'Park_Facility_Name', 'Park_Borough',
      'School_Name', 'School_Number', 'School_Region', 'School_Code',
      'School_Phone_Number', 'School_Address', 'School_City', 'School_State',
      'School_Zip', 'School_Not_Found', 'School_or_Citywide_Complaint',
      'Vehicle_Type', 'Taxi_Company_Borough', 'Taxi_Pick_Up_Location',
      'Bridge_Highway_Name', 'Bridge_Highway_Direction', 'Road_Ramp',
      'Bridge_Highway_Segment', 'Garage_Lot_Name', 'Ferry_Direction',
      'Ferry_Terminal_Name', 'Latitude', 'Longitude', 'Location',
      'Resolution_Action_Updated_Date', 'Call_closing_time', 'Resolution_time'],
      dtype='object')
```

```
In [31]: data_set = f_read.loc[:,['Complaint_Type', 'Location_Type', 'Address_Type', 'Facility_Type', 'City', 'Status', 'Call_closing_time', 'Resolution_time']]
```

```
In [32]: data_set.head()
```

```
Out[32]:
```

	Complaint_Type	Location_Type	Address_Type	Facility_Type	City	Status	Call_closing_time	Resolution_time
0	Noise - Street/Sidewalk	Street/Sidewalk	ADDRESS	Precinct	NEW YORK	Closed	0.933333	0.933333
1	Blocked Driveway	Street/Sidewalk	ADDRESS	Precinct	ASTORIA	Closed	1.450000	1.450000
2	Blocked Driveway	Street/Sidewalk	ADDRESS	Precinct	BRONX	Closed	4.866667	4.866667
3	Illegal Parking	Street/Sidewalk	ADDRESS	Precinct	BRONX	Closed	7.766667	7.766667
4	Illegal Parking	Street/Sidewalk	ADDRESS	Precinct	ELMHURST	Closed	3.466667	3.466667

```
In [33]: lbe=LabelEncoder()
```

```
In [34]: data_set['Complaint_Type']=lbe.fit_transform(data_set['Complaint_Type'])
data_set['Location_Type']=lbe.fit_transform(data_set['Location_Type'])
data_set['Address_Type']=lbe.fit_transform(data_set['Address_Type'])
data_set['City']=lbe.fit_transform(data_set['City'])
data_set['Facility_Type']=lbe.fit_transform(data_set['Facility_Type'])
data_set['Status']=lbe.fit_transform(data_set['Status'])
```

```
In [35]: data_set.head()
```

```
Out[35]:
```

	Complaint_Type	Location_Type	Address_Type	Facility_Type	City	Status	Call_closing_time	Resolution_time
0	16	15	0	0	33	1	0.933333	0.933333
1	4	15	0	0	1	1	1.450000	1.450000
2	4	15	0	0	6	1	4.866667	4.866667
3	12	15	0	0	6	1	7.766667	7.766667
4	12	15	0	0	13	1	3.466667	3.466667

```
In [36]: data_set.describe()
```

```
Out[36]:
```

	Complaint_Type	Location_Type	Address_Type	Facility_Type	City	Status	Call_closing_time	Resolution_time
count	364558.000000	364558.000000	364558.0	364558.0	364558.000000	364558.000000	362177.000000	362156.000000
mean	10.271307	13.88279	0.0	0.0	19.757262	1.006803	3.860239	3.859210
std	5.301311	3.42997	0.0	0.0	15.153827	0.141963	3.785565	3.785543
min	0.000000	0.00000	0.0	0.0	0.000000	0.000000	0.000000	0.000000
25%	4.000000	15.00000	0.0	0.0	7.000000	1.000000	1.250000	1.250000
50%	12.000000	15.00000	0.0	0.0	11.000000	1.000000	2.650000	2.650000
75%	13.000000	15.00000	0.0	0.0	33.000000	1.000000	5.166667	5.166667
max	23.000000	18.00000	0.0	0.0	53.000000	3.000000	23.983333	23.983333

```
In [37]: data_set.dtypes
```

```
Out[37]: Complaint_Type      int32
Location_Type      int32
Address_Type      int32
Facility_Type      int32
City              int32
Status            int32
Call_closing_time  float64
Resolution_time    float64
dtype: object
```

```
In [38]: data_set["Time_bucket"] = np.where(data_set["Resolution_time"]<1,0,
                                             np.where(data_set["Resolution_time"]<3,1,
                                             np.where(data_set["Resolution_time"]<6,2,
                                             np.where(data_set["Resolution_time"]<10,3,
                                             np.where(data_set["Resolution_time"]<12,4,
                                             np.where(data_set["Resolution_time"]<24,5,
                                             np.where(data_set["Resolution_time"]<36,6,7))
                                             ))))
```

```
In [39]: data_set.isnull().sum()
```

```
Out[39]: Complaint_Type      0
Location_Type      0
Address_Type      0
Facility_Type      0
City              0
Status            0
Call_closing_time  2381
Resolution_time    2402
Time_bucket        0
dtype: int64
```

```
In [40]: data_set['Call_closing_time']=data_set['Call_closing_time'].fillna(data_set['Call_closing_time'].median())
```

```
In [41]: data_set['Resolution_time']=data_set['Resolution_time'].fillna(data_set['Resolution_time'].median())
```



```
In [42]: x=data_set.iloc[:,0:8]
         y=data_set.iloc[:,8]
```

```
In [43]: x
```

```
Out[43]:
```

	Complaint_Type	Location_Type	Address_Type	Facility_Type	City	Status	Call_closing_time	Resolution_time
0	16	15	0	0	33	1	0.933333	0.933333
1	4	15	0	0	1	1	1.450000	1.450000
2	4	15	0	0	6	1	4.866667	4.866667
3	12	15	0	0	6	1	7.766667	7.766667
4	12	15	0	0	13	1	3.466667	3.466667
5	12	15	0	0	7	1	1.900000	1.900000
6	12	15	0	0	33	1	1.966667	1.966667
7	4	15	0	0	6	1	1.800000	1.800000
8	12	15	0	0	26	1	8.566667	8.566667
9	4	15	0	0	7	1	1.400000	1.400000
10	4	15	0	0	24	1	7.816667	7.816667
11	4	15	0	0	6	1	11.133333	11.133333
12	16	15	0	0	6	1	2.483333	2.500000

```
In [44]: x.isnull().sum()
```

```
Out[44]: Complaint_Type    0
         Location_Type    0
         Address_Type    0
         Facility_Type    0
         City            0
         Status          0
         Call_closing_time 0
         Resolution_time  0
         dtype: int64
```

```
In [45]: xtrain,xtest,ytrain,ytest = model_selection.train_test_split(x,y,test_size=0.3,random_state=42)
         print(xtrain.shape)
         print(ytrain.shape)
         print(xtest.shape)
         print(ytest.shape)

(255190, 8)
(255190,)
(109368, 8)
(109368,)
```

```
In [46]: from sklearn.neighbors import KNeighborsClassifier
         from sklearn.preprocessing import StandardScaler
         ss= StandardScaler()
         x_train=ss.fit_transform(xtrain)
         x_test=ss.fit_transform(xtest)
```

```
C:\Users\akshat\Anaconda3\lib\site-packages\sklearn\preprocessing\data.py:625: DataConversionWarning: Data with input dtype int
32, float64 were all converted to float64 by StandardScaler.
    return self.partial_fit(X, y)
C:\Users\akshat\Anaconda3\lib\site-packages\sklearn\base.py:462: DataConversionWarning: Data with input dtype int32, float64 we
re all converted to float64 by StandardScaler.
    return self.fit(X, **fit_params).transform(X)
C:\Users\akshat\Anaconda3\lib\site-packages\sklearn\preprocessing\data.py:625: DataConversionWarning: Data with input dtype int
32, float64 were all converted to float64 by StandardScaler.
    return self.partial_fit(X, y)
C:\Users\akshat\Anaconda3\lib\site-packages\sklearn\base.py:462: DataConversionWarning: Data with input dtype int32, float64 we
re all converted to float64 by StandardScaler.
    return self.fit(X, **fit_params).transform(X)
```

```
In [47]: import math
print(len(y))
math.sqrt(len(ytest))

364558

Out[47]: 330.7083307084961

In [48]: category = KNeighborsClassifier(n_neighbors=20)
category.fit(x_train,ytrain)

Out[48]: KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                             metric_params=None, n_jobs=None, n_neighbors=20, p=2,
                             weights='uniform')

In [49]: ypredict=category.predict(x_test)
print(ypredict)

[2 2 2 ... 0 1 3]

In [50]: from sklearn.metrics import accuracy_score
accuracy_score(ytest,ypredict)

Out[50]: 0.9894484675590667
```

Chapter 5

Results and Discussions

Figure 5.1: KNN
[< < < ... 0 1 3]

```
In [52]: ► from sklearn.metrics import accuracy_score
          accuracy_score(ytest,ypredict)

Out[52]: 0.9894484675590667
```

Figure 5.2: Random Forest Algorithm

```
In [69]: ► print("model accuracy:", accuracy_score(ytest, y_pred)* 100)

model accuracy: 99.35904469314607
```

Figure 5.3: Time Series

```
In [79]: ▶ ## plotting of the crime count against the date  
time_se.plot()
```

```
Out[79]: <matplotlib.axes._subplots.AxesSubplot at 0x1adc4ae2d30>
```

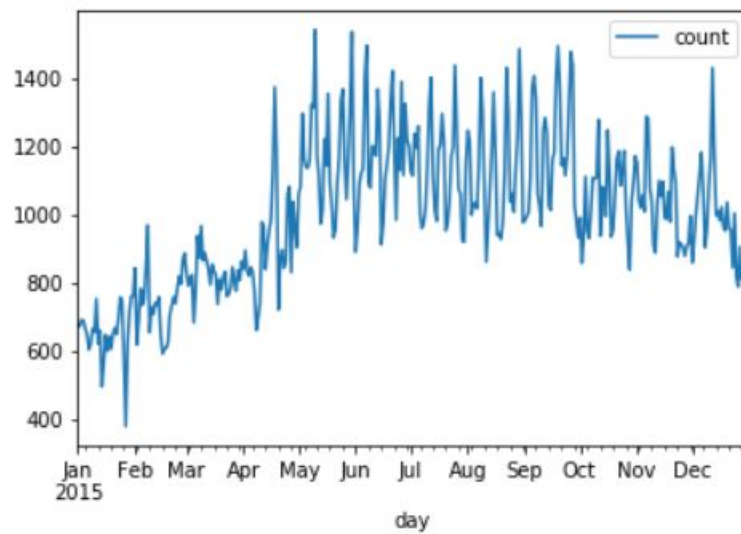


Figure 5.4: Time Series

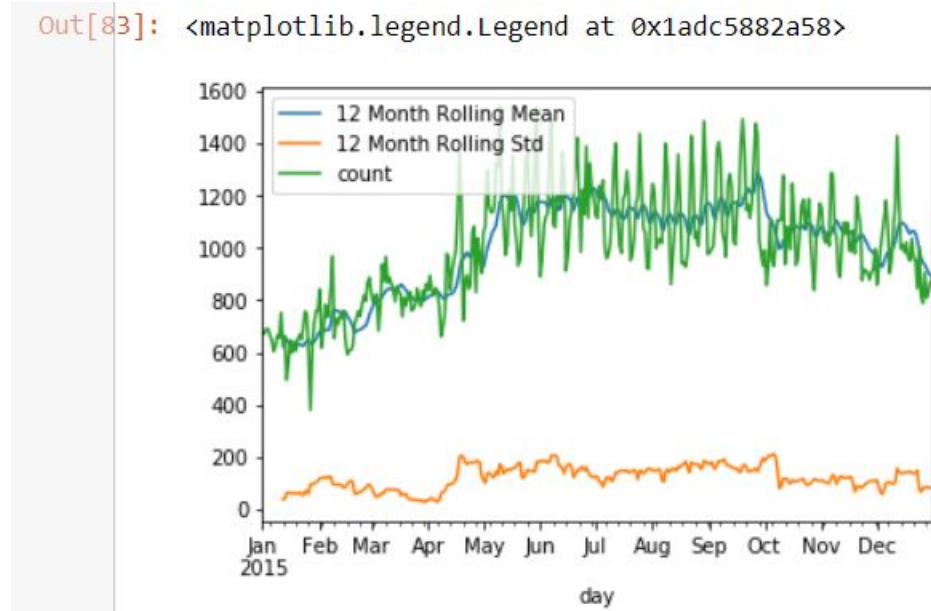


Figure 5.5: Time Series

<Figure size 432x288 with 0 Axes>

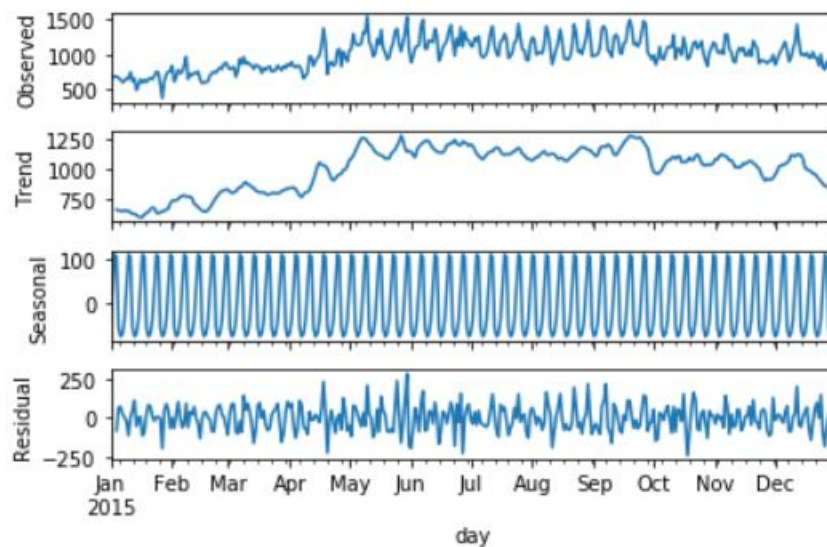


Figure 5.6: Time Series

```
Out[85]: (-2.079247645891997,
          0.25294067957638056,
          13,
          351,
          {'1%': -3.44911857009962,
           '5%': -2.8698097654570507,
           '10%': -2.5711757061225153},
          4221.601886007955)
```

Figure 5.7: Time Series

```
Augmented Dickey-Fuller Test:
ADF Test Statistic : -2.079247645891997
p-value : 0.25294067957638056
#Lags Used : 13
Number of Observations Used : 351
weak evidence against null hypothesis, time series is non-stationary
```

```
] : ▶ time_se['count'].shift(1)
```

```
t[88]: day
2015-01-01    NaN
2015-01-02    675.0
2015-01-03    672.0
2015-01-04    691.0
2015-01-05    689.0
...
2015-12-27    788.0
2015-12-28    905.0
2015-12-29    809.0
2015-12-30    830.0
2015-12-31    865.0
Name: count, Length: 365, dtype: float64
```

Figure 5.8: Time Series

Out[92]:

	count	count_first_diff
day		
2015-01-01	675	NaN
2015-01-02	672	-3.0
2015-01-03	691	19.0
2015-01-04	689	-2.0
2015-01-05	665	-24.0

```
[93]: time_se.count_first_diff.plot()
```

Out[93]: <matplotlib.axes._subplots.AxesSubplot at 0x1adc8753668>

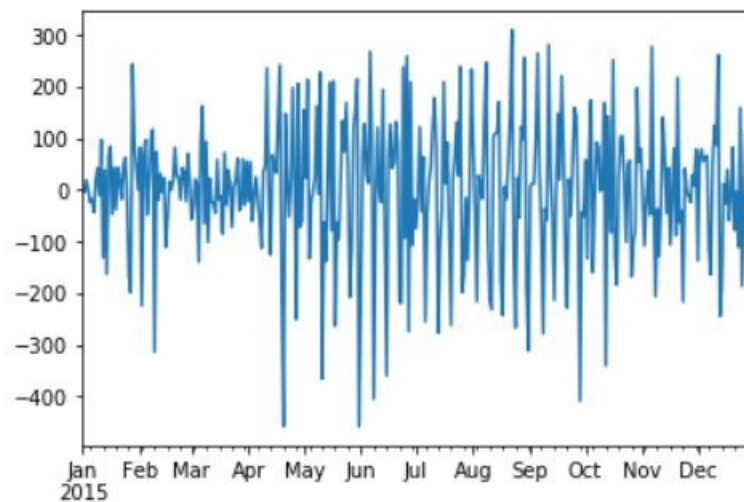


Figure 5.9: Time Series

Augmented Dickey-Fuller Test:

ADF Test Statistic : -9.209768555517952

p-value : 1.9016230790746454e-15

#Lags Used : 12

Number of Observations Used : 351

strong evidence against the null hypothesis, reject the null hypothesis. Data is

Chapter 6

Conclusion and Future Work

The project has done with machine learning approaches. Their are mainly three algorithms have been used named Random forest Algorithm(RFA),Time Series and K-nearest neighbor(KNN). Among three algorithms random forest has shown excel forecast upto 99 percent accuracy. Also, KNN has shown good accuracy of 98 percent. The project has predicted,the time taken or the duration to solve particular type of complaints. Moreover, Random Forest Algorithm has the best analysis and can be used for further prediction.

Bibliography

- [1] Prajakta S. kasbe Research Scholar, M. Tech Computer Science Engineering, G. H. Rasoni College of Engineering, Nagpur, India. kasbe prajakta.ghrcemtechcse@raisoni.net
- [2] ben.george@hct.edu.om Muscat, Oman Higher College of Technology Department of Information Technology mBen George Ephre girija@hct.edu.om Muscat, Oman Higher College of Technology Department of Information Technology Girija Narasimhan
- [3] Bulbula Kumeda, Fengli Zhang, Fan Zhou School of Information and Software Engineering, University of Electronic Science and Technology of China, Chengdu, China, 610054 e-mail: bekumeda@gmail.com, fzhang, fan.zhou@uestc.edu.cn
- [4] Bulbula Kumeda, Fengli Zhang, Fan Zhou School of Information and Software Engineering, University of Electronic Science and Technology of China, Chengdu, China, 610054 e-mail: bekumeda@gmail.com, fzhang, fan.zhou@uestc.edu.cn
- [5] Priyanka A. Nandurge Nagaraj V. Dharwadkar Department of Computer Science and Engineering Department of Computer Science and Engineering Rajarambapu Institute of Technology, Islampur India Rajarambapu Institute of Technology, Islampur India priyankanandurge@gmail.com nagaraj.dharwadkar@ritindia.edu
- [6] J.Saravanakumar A.Anigo Merjora Department of CSE, Department of CSE, Jeppiaar Maamallan Engineering College, Jeppiaar Maamallan Engineering College, Anna university ,Chennai India.. Anna university ,Chennai India. saravanakumar2005@gmail.com anigomerjora@gmail.com

-
- [7] ulbula Kumeda, Fengli Zhang, Fan Zhou School of Information and Software Engineering, University of Electronic Science and Technology of China, Chengdu, China, 610054 e-mail: bekumeda@gmail.com, fzhang, fan.zhou@uestc.edu.cn
 - [8] Ms. E. Suganya Ph.D Research Scholar Department of Computer Science Bharathiar University, Coimbatore-641046 elasugan1992@gmail.com
 - [9] Jian Zhang¹, Zhibin Li¹, Ziyuan Pu², Chengcheng Xu¹ ¹School of Transportation, Southeast University, Nanjing, 211189 China ²Department of Civil and Environmental Engineering, University of Washington, Seattle, 98195, United States Corresponding
 - [10] pho Mokoatle North West University Council for Scientific and Industrial Research (CSIR) Defence, Peace, Safety and Security Pretoria, South Africa Email: mmokoatle@csir.co.za