

GLOBAL JOURNAL OF ENGINEERING SCIENCE AND RESEARCHES

DATA NOISE REDUCTION AND ANALYSIS FOR DOMAIN BASED LEAD GENERATION

Gaurav Kumar^{*1}, Sheba Selvam², Meghana R. Salagundi³ & Aastha Mishra⁴

^{*1,2,3&4}Computer Science of Engineering, Rajarajeswari College of Engineering, Bangalore, India

ABSTRACT

Traditional ways to spot sales leads like company surveys and marketing are manual, high priced and not scalable. Data collected from multiple sources are used as reference to gather the data could be used to recognize potential buyers automatically. In here, we use data science & machine learning to generate new leads for any new product, service or franchise launches, by a company in market. Businesses face challenges to cleanse the massive quantity of collected information. The information may be unstructured, unverified, and might contain heaps of noise within the information. We basically prepare the data, and identify the required features by feature engineering, analyze the search & engagement patterns of customers towards the products/services. We will be using the prepared data to train our machine learning algorithm (clustering or classification) to develop, test and fine tune recommendation models, which will be helpful to generate an effective lead for converting an existing customer to buy more products/services, or to identify new customers who might also be interested in a product/service. In this paper, we propose an intelligent recommendation engine, which can cleanse the silos of data and dynamically identify leads for new business products/services.

Keywords: data cleansing, machine learning, recommendation models, lead generation, clustering, classification.

I. INTRODUCTION

Data science, additionally referred to as data-driven science, is a knowledge base field of scientific strategies, processes, algorithms and systems to withdraw knowledge or perceptions from information in numerous forms, either structured or unstructured, almost like data processing. Data science could be a "concept to unify statistics, data analysis, machine learning and their connected methods" so as to "perceive and examine actual phenomena" with information.[3] It makes use of methods and theories extracted from several fields inside the broad areas of arithmetic, statistics, information science, and computer science, above all from the subdomains of machine learning, classification, cluster analysis, uncertainty quantification, computational science, data processing, databases, and visualization. Noise is "unrelated or insignificant or senseless data". For many existing information cleansing ways, the main focus is on the detection associated removal of noise (low-level information errors) that's the outcome of a defective data collection process. This demands a need to address this sort of noise is obvious because it is prejudicial to virtually any reasonably data analysis. However, normal data objects that are not relevant or hold a very weak percentage of relevance to a selected data analysis can even considerably hinder the data analysis, and so these objects ought to be additionally thought-about as noise, a minimum of within the context of a particular analysis. as an example, in document information sets that incorporates news stories, there are several stories that are solely associated with the other news stories. If the aim is to use clustering to search sturdy topics in an exceedingly set of documents, then the analysis can suffer unless moot and sapless relevant documents may be eliminated. Consequently, there's a desire for data cleansing techniques that take away each forms of noise. In some cases the quantity of noise in an exceedingly information set is comparatively tiny. Lead generation is that the action or method of recognizing and cultivating potential customers for a business's merchandise or services. Leads might come back from numerous sources or activities, as an example, digitally via the net, through personal referrals, through phone calls either by the corporate or telemarketers, through advertisements, and events.

II. PROBLEM DEFINITION

- Identifying new customers for any business is a tough task
- Most of the businesses follow on 2 approaches,

- Leveraging on existing customers by retargeting products/services to them (ex. you've bought an AC from a company, they'd promote you to also buy a TV from them)
- Bringing in new customers, which can be footfalls within the store, or website visits or first time customers.
- In here, we think that we could use multiple points of data being collected to identify prospective customers and target products that the customer might be interested in buying. Exposure to data will also help us in bringing new customers to the business. In summary, we generate leads, by
 - Leveraging on existing customers to buy more products,
 - Use data to identify new prospective customers.

III. ARCHITECTURE DIAGRAM

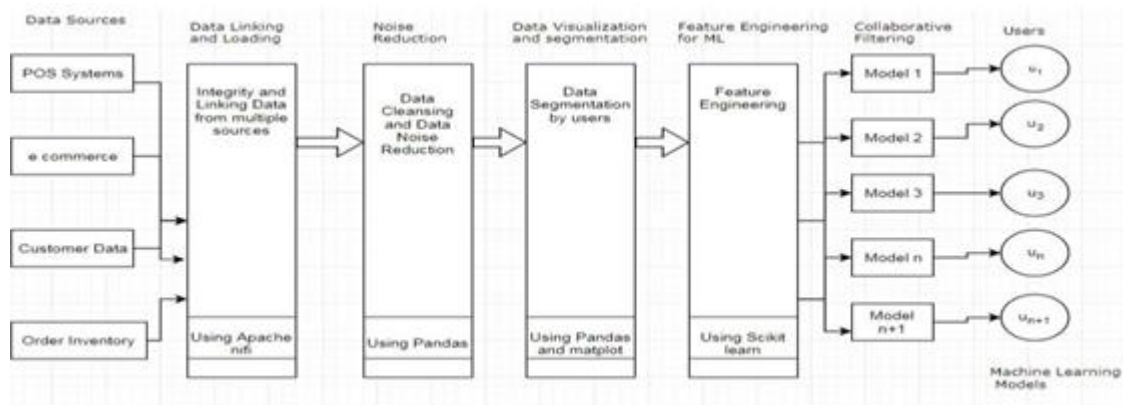


Fig. 1. Architectural diagram

The proposed architectural diagram consists of the following modules

- ☞ Data Sources.
- ☞ Data Linking and Loading.
- ☞ Noise Reduction.
- ☞ Data Visualization and Segmentation.
- ☞ Feature Engineering for Machine Learning.
- ☞ Recommendation models.

A. Data sources

In PC programming, source information or information source is the essential area from where information comes. The information source can be a database, a dataset, a spreadsheet or even hard-coded information. At the point when information is shown in a page or application, in a section push design or different configurations, the information is recovered from its information source and introduced in the arrangement characterized in the code.

Here the data sources are from a variety of systems that contain the customer information. In order to gather the data we use various sources like POS Systems, E-Commerce, Customer data, Data Inventory etc.

B. Data linking and loading

The connecting and stacking are expert by a program which is differently called the linker, or the loader, or the connecting loader. In spite of the fact that connecting and stacking are theoretically discrete, it is exceedingly normal that a solitary program consolidate those capacities. So "connecting" and "stacking" mean a similar thing, when talking casually, despite the fact that they are separate activities. In this proposed system, we first load all the data obtained from various form of data sources and then combine them by linking them with one another. This process of linking is the most important aspect of our system as it helps us analyze the customer behavior and their

pattern of searching for their item of interest, which in turn helps us come up with an effective recommendation model for the users.

C. Noise reduction

Evacuating objects that are noise is an essential objective of information cleaning as commotion prevents most kinds of information examination. Most existing information cleaning techniques center around expelling noise that is the aftereffect of low-level information mistakes that outcome from a blemished information gathering process, yet information protests that are superfluous or just feebly important can likewise essentially prevent information investigation. Consequently, if the objective is to improve the information examination however much as could reasonably be expected.

these articles ought to likewise be considered as noise, in any event regarding the fundamental investigation. Subsequently, there is a requirement for information cleaning methods that evacuate the two sorts of commotion. Since informational collections can contain a lot of noise, these systems likewise should have the capacity to dispose of a possibly vast portion of the information. Once the data is loaded and linked the next important step is to clean this data that contain many irrelevant data that create noise. In order to clean this noise we make use of many python libraries which are found in the Anaconda Python Framework. This frame work contains Jupiter editor that helps us work on the cleaning of the data by making use of certain python functions present in its libraries. It helps us work on the linked data directly and displays the results immediately.

D. Data visualization and segmentation

Data visualization presents how a given performance is spread by item (Region in this example). Interesting here to be able to compare at the same time the performance spreads across two analytical axes. This visualization is an efficient way to see the best customers for instance and their relative weight in your overall business. Once the unwanted noise is removed from the data the next step towards building a recommendation model is to segment these set of information according to the various category of customer available. This data is segmented based on the customer's choice, preference and priorities. This segmented data helps us identify potential customers for any business. This data is obtained by analyzing the customer's pattern of purchasing goods, or searching for interested items, or even the price range the customers are comfortable with. This set of segmented data that has different information for different customer helps us build a unique recommendation model for every customer.

E. Feature engineering for machine learning

Feature engineering is the way toward utilizing area information of the information to make includes that make machine learning calculations work. Feature engineering is basic to the use of machine learning, and is both troublesome and costly. The requirement for manual element building can be forestalled via robotized highlight learning. An element is a trait or property shared by the majority of the free units on which investigation or expectation is to be finished. Any trait could be an element, as long as it is helpful to the model. The reason for an element, other than being a property, would be significantly simpler to comprehend with regards to an issue. A component is a trademark that may help when tackling the issue. The features in our fragmented information are imperative to the suggestion show that we make utilization of and will impact the outcomes that will be accomplished. The quality and amount of the features will have incredible effect on whether the model is great or not. The better the features are, the better the outcome is picking the correct features are still imperative. Better features can deliver less complex and more adaptable models, and they regularly yield better outcomes. The calculations we utilized are extremely standard for Kagglers.

F. Recommendation models

A recommender framework or a suggestion framework is a subclass of data sifting framework that tries to anticipate the "rating" or "inclinations" a client would provide for a thing. Suggestion models have turned out to be progressively prevalent as of late, and are used in an assortment of zones including motion pictures, music, news, books, examine articles, look questions, social labels, and items when all is said in done. There are additionally suggestion models for specialists, teammates, jokes, eateries, pieces of clothing, money related administrations, extra security, and Twitter pages. Recommender frameworks regularly create a rundown of proposals in one of two courses – through community oriented sifting or through substance based

separating (otherwise called the identity based approach) Collaborative separating approaches manufacture a model from a client's past conduct (things beforehand obtained or chose as well as numerical evaluations given to those things) and also comparative choices made by different clients. This model is then used to anticipate things (or appraisals for things) that the client may have an enthusiasm for. Content-based sifting approaches use a progression of discrete attributes of a thing keeping in mind the end goal to suggest extra things with comparative properties. Hence this recommendation model in our proposed system helps us come up with a unique solution provided to every customer to enhance their way of searching for items they might be interested in and also purchase them. This recommendation model not only helps the customers but also helps the companies by marketing their products efficiently.

IV. RELATED WORK

The planning of a novel book proposal framework is discussed in [1]. Pursuers will be diverted to the proposal pages when they can't locate the required book through the library bibliographic recovery framework. The suggestion pages contain all the fundamental and growing book data for pursuers to allude to. Pursuers can suggest a book on these pages, and the proposal information will be dissected by the suggestion framework to settle on logical obtaining choice. It can meet the book requests of pursuers to the best degree, in the meantime it can spare the acquisition financing for the library. The exact impact still ought to be confirmed through long haul trial. We ought to explore the book use in the course database and change the book esteem model and duplicate number model as indicated by the measurable information.

An unstructured idea of data influences the extraction to assignment of trigger occasions troublesome is discussed in

It represent the issue of trigger occasions extraction as an order issue and create strategies for learning trigger occasion classifiers utilizing existing arrangement techniques. We exhibit strategies to naturally produce the preparation information required to take in the classifiers. We additionally propose a technique for highlight deliberation that utilizations named element acknowledgment to take care of the issue of information shortage. We score and rank the trigger occasions extricated from ETAP for simple perusing. Our investigations demonstrate the adequacy of the technique and accordingly set up the possibility of programmed prospective customer age utilizing the Web data. It utilizes highlight deliberation to address the issue of information shortage and to accomplish speculation. It gives a novel method those aides in distinguishing the correct level of abstraction. We need to know every one of the varieties to the reference of the organization. This data isn't generally accessible and computerized techniques to decide varieties of an organization name should be created.

Evacuating objects that are clamor is an essential objective of information cleaning as commotion obstructs most kinds of information examination are explained in [3]. Most existing information cleaning strategies center around expelling commotion that is the consequence of low-level information mistakes that outcome from a flawed information gathering process, however information questions that are unimportant or just pitifully pertinent can likewise essentially obstruct information examination. Thus, there is a requirement for information cleaning strategies that expel the two kinds of commotion. This paper investigates four methods proposed for commotion expulsion to upgrade information examination within the sight of high clamor levels. Three of these strategies depend on conventional exception recognition methods: separate based, bunching based, and an approach in light of the Local Outlier Factor (LOF) of a question. The other system, which is another technique that we are proposing, is hyper inner circle based information cleaner (HCleaner). Our exploratory outcomes demonstrate that these strategies can give better bunching execution and higher quality affiliation designs as the measure of commotion being expelled increments, in spite of the fact that HCleaner for the most part prompts better grouping execution and higher quality relationship than the other three techniques for parallel information. This might be commotion because of defects in the information accumulation process or clamor that comprises of superfluous or pitifully significant information protest. Given that HCleaner, CCleaner, and the LOF based strategy were each the best in various circumstances, it could be helpful to consider a voting plan that consolidates these three procedures, however it isn't actualized to that broaden.

A tourism RS that is based on its recommendations on data dynamically aggregated and extrapolated from the Facebook check-in data is explained in [4]. In addition, the so-called “cold-start” problem has been resolved by using users’ Friends’ check- in data to analyze ongoing Facebook activity and update user profiles in the system. An RS is an information filtering system that analyses user preferences and adapts its functions to individual users and finding user interests or preferences is therefore a key process of an RS. The information here is directly useful for user attraction preference analysis and significantly benefits tourism industries. In this the approach differs from existing approaches presented in the literature by overcoming the problems by collecting information from individual users and Friends available in Facebook. Here the approach relies solely on user check-in activities which only represent visits to specific places without indicating like or dislike preferences.

There are a few approaches to take care of the issue, for example, classification list, web index and suggestion motor. Among them web index and proposal motor are more well known is expressed in [5]. In genuine circumstance a proposal motor has been turned out to be much more powerful than web index in online business and customized data recovery. Incidentally, as a result of the utilization of portable Internet, online short video has turned into an exceptionally prevalent type of video these days. A decent proposal framework can enable individuals to achieve the substance they to need effectively. This article presents the plan of a strong proposal motor particular for online short video. It is a decent proposal framework that can enable individuals to achieve the substance they to need effectively. Need to enhance the present calculations. Enhance the stage with new strategies suggestion assessment.

A center system is mapping of client look Goals in regulated reference classes with unsupervised data recovered from web which makes semi-administered approach, a best approach is expressed in [6]. Calculation advancement centers are around improved execution in sense to discover pertinent data instead of time multifaceted nature. The result is Recommendation Engine, which suggests client in view of client navigate logs. In this examination a precise Literature review of 25 articles help to discover inquire about degree and in this way procedure to understand this issues and difficulties. Investigation of inquiry logs helps in finding look objectives for equivocal questions, for which proposed framework groups criticism sessions. What truly customer watches over differing inquiries revelation of reasonable pre-unmistakable pursuit references classes is testing and irrational.

The issue of information irregularities while incorporating informational collections from different self-ruling social databases is addressed by [7]. We begin by belligerence that the semantics of incorporating potentially conflicting information is normally caught by the maximal predictable subsets of the arrangement of all data contained in the gathered information. In view of this thought, we propose a basic and natural semantically structure, called the coordinated social analytics which is an expansion of the established social math, for controlling and questioning conceivably conflicting information. We additionally contend that JEexible social model gives a fairly powerless inquiry dialect. We at that point demonstrate that for the databases with just a single key the adaptable model gives a right incorporation of conflicting information. We propose a basic and natural semantic structure; called the incorporated social math which is an expansion of the traditional social analytics, for controlling and questioning perhaps conflicting information. There are various issues which have been left open in this paper. The first is to locate a powerful calculation for inquiry assessment. For those questions which are proportionate to articulations in adaptable social polynomial math, strategies created in adaptable social model can be connected. In any case, since adaptable social model is fairly feeble, we likely need to search somewhere else for such calculation. Another issue is to expand this coordinated social model for different sorts of invalid esteems.

A substantial Twitter data-set, we break down the attributes of client and pattern profiles and assess the nature of the profiles with regards to a customized news suggestion framework is spoken in [8]. In this paper, we explore how singular Twitter exercises can be misused to construe individual interests and produce semantic client profiles that can be re-utilized likewise by different applications than Twitter. To acquire semantically significant ideas for speaking to the clients' advantages, our Twitter-based User Modeling Framework takes into consideration the age of various kinds of semantic client profiles. Individual client intrigue profiles are more critical for the news article

proposal process than open patterns. By entwining pattern and client profiles we prevail in additionally enhancing the suggestion quality. This idea is limited to twitter application as it were.

A lead age library configuration is an iterative procedure of choosing exacerbates that would make perfect medication hopefuls which is portrayed by [9]. We propose a visual examination show for successful assessment of different compound choice techniques. This model joins a grouping technique and an arrangement of assessment measurements for understanding the present condition of information to empower basic leadership and process customization between cycles. Visual investigation display was intended to help comprehend the present condition of information amid a solitary emphasis to empower basic leadership and process customization before the following cycle is executed. More extensive assortment of issues identified with sedate revelation ought to be approved.

Recommender frameworks have turned into a vital research field since 1990's and its applications incorporates a few areas which is spoken in [10]. The utilization of proposals in human services is a wide zone which suggests the patients about their wellbeing. Upon immense measure of information winning in human services area, this information can be handled utilizing enormous information devices to convey a significant forecast to the patients. The forecasts and suggestions will be more precise since we are managing tremendous measure of information. In addition it cautions the client from the event of sickness and takes fundamental activities previously it happens. In this way recommender framework alongside enormous information will guarantee arrangement that is winning in human services area. The whole of information identified with the patient and their prosperity constitutes the "Huge Data" issue in the social insurance industry. At the point when the information measure is huge we can foresee the consequences of specific event effortlessly. It likewise especially focuses on expectation and analysis thyroid issue in ladies. it gives a novel proposal motor for forecast of thyroid issue in ladies are not created.

sale-arranged online-shop administration bolster technique is proposed is portrayed in [11]. In this strategy, utilizing a deal administration site page, online-shop proprietors can effectively and dependably develop appealing deal pages reliable with their DB, utilizing various leveled deal learning made by a product specialist and adjusted by them. In particular, the programmed development and intuitive changes of offer pages and additionally the programmed/intelligent refresh of DB for every deal item gathering should be possible reliably and progressively, synchronized with the deal. It gives programmed age/show/intuitive remedy of a business screen. Programmed or intuitive advancement/affirmation of the items DB should be possible in one stroke or progressively as per a business period. It doesn't choose sentences through delaying the showed review screen subsequent to setting deals headings on the business administration page. Produce a business show consequently after rectification of sentences, their shading/textual style/measure/traits, (for example, adornment of words) or changing/including outlines through intuitive capacity keys (catches).

V. IMPLEMENTATION

The proposed system introduces the concept of building unique recommendation models for a variety of customers that it would serve. In the first step we gather data from various sources with the help of customer interaction in various ways like, customer search pattern, login details, storing of customer data for later purchases, orders placed by the customers, their likes and dislikes and so on. All these data are collected from various sources and loaded in one place, so that the manipulation on these sets of data i.e. a dataset becomes more convenient to handle. Here we use the dataset of a grocery store to identify our potential customers and hence enhance the leads. The datasets are now linked together according to the various parameters set for differentiating the data under different groups and are linked to one another. The information that are linked here are based on the various categories of customers i.e. according to their items of interest, to the price range a customer is comfortable in placing orders and so on. Once these datasets are categorised and the set of potential customers that a company might be interested in have been identified we now make use of our Anaconda Python Framework which comes with a number of advantages to perform the noise reduction. Noise here means the unwanted data that may cause some errors to the inputs of the recommendation models. Jupyter is one such editor present in Anaconda framework that helps us remove the unwanted noise from these categorised datasets. Jupyter notebook uses a number of inbuilt functions and libraries to

perform the noise reduction and displays the results immediately. The cleansed data is now ready to be segmented for the users to build the recommendation models. Next step is collaborative filtering (CF), it is a procedure utilized by recommender systems. Collaborative filtering has two detects, a limited one and a more broad one. In the more up to date, smaller sense, collaborative filtering is a strategy for making programmed forecasts (filtering) about the interests of a client by gathering inclinations or taste data from numerous clients (teaming up). The hidden suspicion of the collaborative filtering approach is that if a man A has an indistinguishable supposition from a man B on an issue, A will probably have B's assessment on an unexpected issue in comparison to that of a haphazardly picked individual. For instance, a collaborative filtering suggestion framework for TV tastes could influence forecasts about which TV to demonstrate a client should like given a fractional rundown of that client's tastes (likes or dislikes). Note that these expectations are particular to the client; however utilize data gathered from numerous clients. This varies from the less difficult approach of giving a normal (non-particular) score for everything of enthusiasm, for instance in view of its number of votes. In the more broad sense, collaborative filtering is the way toward filtering for data or examples utilizing methods including coordinated effort among various operators, perspectives, information sources, etc. Applications of collaborative filtering ordinarily include extensive informational collections. Collaborative filtering strategies have been connected to a wide range of sorts of information including: detecting and observing information, for example, in mineral investigation, ecological detecting over huge zones or various sensors; money related information, for example, monetary administration foundations that incorporate numerous budgetary sources; or in electronic trade and web applications where the attention is on client information, and so forth. The rest of this discourse centers on collaborative filtering for client information, albeit a portion of the techniques and methodologies may apply to the next significant applications also. The next step includes the use of Feature engineering which helps us to extract the essential features from the segmented data. These features may include checking the features, deciding what features are to be taken into consideration, checking whether the feature will work with our model and so on according to the datasets fed as input to it. We now use Machine learning algorithms to build and complete our recommendation models for the customers. These machine learning algorithms are the applications of Artificial intelligence as they self-learning and learn through experience. The unique recommendation models serve the customers by providing them the updates about the products they might be interested in and also the companies who want to find their potential customers to improve the leads of the company and also market their products.

VI. RESULTS AND DISCUSSIONS

- A. Collected data: In Fig. 2 and Fig. 3 data is collected from various sources that include surveys and public search patterns etc. The data of imdb and grocery store collected is in “.csv” format i.e. rows and column. The collected data is generally unstructured and unformatted. There are few examples of datasets which has been collected from online sources like IMDB (international movie data base) and grocery store dataset. It consists of a number of columns and attribute and only few of them are needed.



title	year	genre	director	cast
The Godfather	1972	Drama	Francis Ford Coppola	Al Pacino, Marlon Brando
The Godfather Part II	1974	Drama	Francis Ford Coppola	Al Pacino, Al Pacino, Marlon Brando
The Godfather Part III	1978	Drama	Francis Ford Coppola	Al Pacino, Al Pacino, Marlon Brando

Fig 2 Imdb collected data in .csv form



Fig. 3. Grocery datasets collected data

B. Noise reduction in data: fig. 4 shows an output of noise reduction functions done in jupyter notebook. The result of data noise reduction is a formatted structured data with optimized information. It removes the unwanted data that is not needed for the system purpose and gives only the important data. As in the analyzed data there are so many columns in original data but only five columns are needed so rest unwanted columns are deleted and new dataset is created on which all the manipulations are done. In the given example unwanted attributes are removed and only needed attributes like budget, revenue, movie id and the releasing dates etc.



Fig. 4. Result of imdb cleansed data

C. Analyzed data: This is the feature of cleaned data that tells the customer behavior and interest. It identifies all the features of data. It provides a key perspective on which the recommendations can be assigned. Fig. 5 shows the kernel density graph for the data. Fig. 6 shows histogram analysis on the dataset.

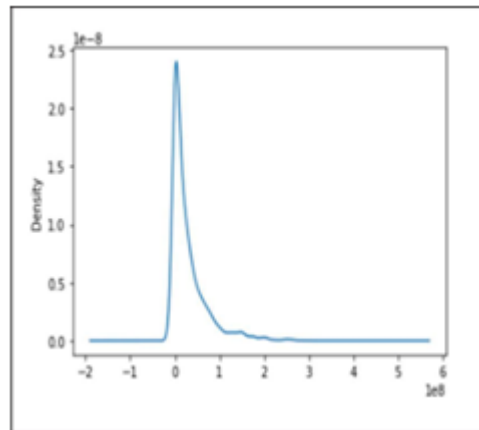


Fig. 5. Kernel density graph of imdb

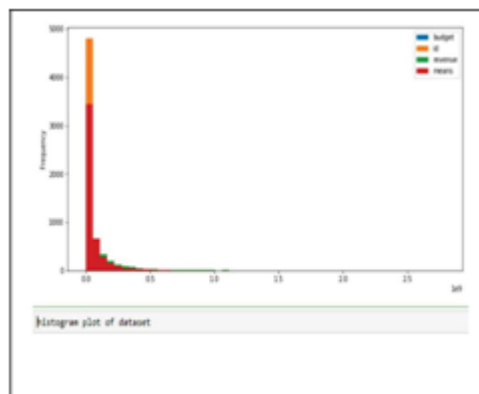


Fig. 6. Histogram plot of imdb data

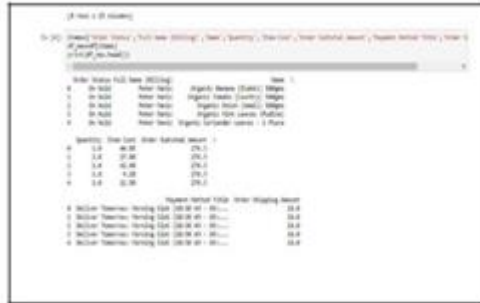


Fig. 7. Analysis of collected grocery

D. Customer segregation: Fig.7 shows the various groups of customers with their different values for the company. Some customers are high valued customer while some are low valued based on their investment for the company and benefit of company from them. Fig. 8 analyzes the data features and comes with answers of lead generation regarding questions.



Fig. 8. Psychographic segmentation of grocery dataset

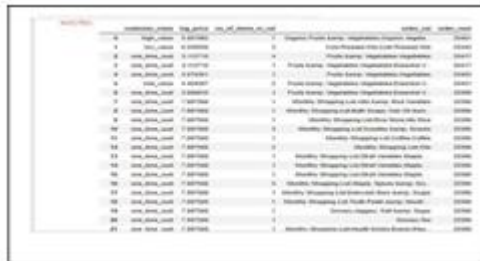


fig. 9. Psychographic segmentation of customers

E. Recommendation models: These are policy sales, offers and benefits that are decided by company for their customers. It is made to interact with various customers and engage them into company product purchasing according to their value for the company and their area of interest. In this segregated customers in different class of value will be provided different recommendations accordingly. It is the model trained by machine learning algorithms such as classification algorithms like NB classifiers shown in fig. 10 and tested with other models shown in fig. 11 also for the accuracy of predictions. The models are trained such that for a new entry there will of no need of segregation again and again, it should be done automatically by the machines and companies should treat them according to the class of segmentation. Finally all the classifier models are compared graphically shown in fig. 12 that gives the average accuracy approximately 80%.

```

dividing data for training & test
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state = 0)

Training the NB Classifier
from sklearn.naive_bayes import MultinomialNB
clf = MultinomialNB()
clf.fit(X_train, y_train)

Testing the model
# Model testing
for i, (x_test, y_test) in enumerate(zip(X_test, y_test)):
    prediction = clf.predict(x_test)
    print('Actual: %s Predicted: %s' % (y_test[i], prediction[i]))
    if not prediction[i] == y_test[i]:
        print(i)
    
```

fig. 10. Training and testing of NB classifier model

```

Checking confusion matrix
from sklearn.metrics import confusion_matrix, accuracy_score
cm = confusion_matrix(y_test, y_pred)
print(cm)

from sklearn.metrics import precision_score, recall_score
from sklearn.metrics import mean_absolute_error
from sklearn.metrics import mean_squared_error
from sklearn.metrics import r2_score
from sklearn.metrics import mean_absolute_percentage_error
print('Precision: %s' % precision_score(y_test, y_pred))
print('Recall: %s' % recall_score(y_test, y_pred))
print('MAE: %s' % mean_absolute_error(y_test, y_pred))
print('MSE: %s' % mean_squared_error(y_test, y_pred))
print('R2: %s' % r2_score(y_test, y_pred))
print('MAPE: %s' % mean_absolute_percentage_error(y_test, y_pred))

Testing with few other classifiers as well to check accuracy
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.neighbors import KNeighborsClassifier

```

fig. 11. Testing with other classifier model



fig. 12. Accuracy check graph of different models

VII. CONCLUSION

Business uses multiple sources of information to gather the data to identify potential customers which is a challenging task. This leads to collection of unwanted data that may hinder the building of good recommendation models for customers. Amongst the recommendation models available at present, the proposed system provides better accuracy in terms of recommendations to the customers and companies to enhance their leads by cleansing and segmenting the customers based on their interest. This is achieved by using collaborative filtering and the knowledge of machine learning by the proposed system.

REFERENCES

1. Binge Cui, Xin Chen. An online book recommendation system based on web service. In 2009 Sixth International Conference on Fuzzy Systems and Knowledge Discovery.
2. Ganesh Ramakrishnan, Sachindra Joshi, SumitNegi, Raghu Krishnapuram, SreeramBalakrishnan. Automatic Sales Lead Generation from Web Data. In Proceedings of the 22nd International Conference on Data Engineering (ICDE'06)8-7695-2570-9/06 \$20.00 © 2006 IEEE.

3. HuiXiong, Member, IEEE, GauravPandey, Michael Steinbach, Member, IEEE, and Vipin Kumar, Fellow, IEEE. Enhancing data analysis with noise removal. *IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING*, VOL. X, NO. X, XXX 200X.
4. K. KESORN , W. JURAPHANTHONG, AND A. SALAIWARAKUL. Personalized attraction recommendation system for tourists through check-In data. Received October 26, 2017, accepted November 24, 2017, date of publication November 29, 2017. Date of current version December 22, 2017.
5. Nan Yang, Sanxing Cao, Yu Liang, Zhengzheng Liu. Recommendation Engine for Online Short Video. In 2016 4th Intl Conf on Applied Computing and Information Technology/3rd Intl Conf on Computational Science/Intelligence and Applied Informatics/1st Intl Conf on Big Data, Cloud Computing, Data Science & Engineering.
6. OmkarTodkar, Prof.S.Z.Gawali, AniketD.Kadam. Recommendation Engine Feedback Session Strategy for Mapping User Search Goals. In 2016 IEEE.
7. Phan Minh Dung. Integrating data from possibly inconsistent databases. 0-8186-7505-5/06 2006 IEEE.
8. Qi Gao, Fabian Abel, Geert-Jan Houben, Ke Tao. Interweaving Trend and User Modeling for Personalized News Recommendation. In 2011 IEEE/WIC/ACM international conferences on web intelligence and intelligent agent technology
9. Shawn Konecni, Jianping Zhou, Georges Grinstein University of Massachusetts Lowell {skonecni@cs.uml.edu, jzhou@cs.uml.edu, grinstein@cs.uml.edu}. A Visual Analytics Model Applied to Lead Generation Library Design in Drug Discovery. In © 2009 IEEE DOI 10.1109/IV.2009.75.
10. Shobana.V1, N.Kumar2. A personalized recommendation engine for prediction of disorders using data analytics. In IEEE International Conference on Innovations in Green Energy and Healthcare Technologies (ICIGEHT'17).
11. Yoshitaka Sakurai¹, Takashi Kawabe¹, Takahiko Sakai¹, Kouhei Takada¹, Setsuo Tsuruta¹, Mizuno Yoshiyuki, School of Information Environment, Tokyo Denki University, Chiba, Japan, {ysakurai, tsuruta}@sie.dendai.ac.jp Kyoto Women's University, Kyoto, Japan. A sale-oriented online-shop management support method for e-commerce. In 2010 IEEE DOI 10.1109/SITIS.2010.54.