

# Data Mining Using K-Mean Clustering and K-Nearest Neighbor Classification

**Prof. Shivangi D. Matieda, Prof. Mani Butwall, Prof. Diksha Durgapal**

Department of Computer Science and Engineering

ITM Universe, Vadodara, Gujrat, India

s.borasia@gmail.com, s.borasia@gmail.com, s.borasia@gmail.com

**Abstract:** *Understudies casual discussions via web-based networking media (e.g., Twitter, Facebook) shed regard into their instructive encounters—sentiments, emotions, and worries concerning the training procedure. data from such unenlightened things will offer profitable data to teach understudy learning. Examining such data, in any case, will be testing. The complexity of understudies' encounters mirrored from net based mostly life content needs human elucidation. In any case, the developing size of data knowledge requests programmed information investigation systems. During this paper, we have a tendency to design a piece method to include each subjective investigation and large-scale data mining ways. We have a tendency to focus on building understudies Twitter presents to get problems and problems in their instructive encounters. We have a tendency to originally direct a subjective examination on tests taken from around twenty-five thousand tweets known as building understudy's college life. We have a tendency to discover building understudy's expertise problems, as an example, substantial investigation load, absence of social commitment, and lack of sleep. In light of those outcomes, we have a tendency to actually do a multi-name arrangement calculation to order tweets mirroring understudies' problems.*

**Keywords:** Data Mining, Index Terms—K-Nearest Neighbor, K-Mean Algorithm, Facebook, Twitter

## I. INTRODUCTION

Web primarily based life locales, as an example, Twitter, Facebook, and You-Tube provide unbelievable settings to understudies to share their encounters, vent feeling and stress, and appearance for social help. On totally different web-based social networking destinations, understudies refer and share their normal experiences in a casual and easygoing approach. Understudies' advanced impressions provide immense live of verifiable learning and a totally different purpose of read for instructive analysts and professionals to grasp understudies' encounters outside the controlled homeroom condition. This comprehension will educate institutional basic leadership on mediations for in danger understudies, improvement of coaching quality, and in this way upgrade understudy tour of duty, maintenance, and achievement. The plenty of net primarily based life data gives possibilities to grasp understudies' encounters, yet in addition brings method challenges up in understanding internet primarily based life data for instructive functions. Simply envision the sheer data volumes, the various variety of web slang, and also the capriciousness of space and timing of understudies posting on the net, even as the multifaceted nature of understudies' encounters. Pure manual examination can't manage the systematically developing size of knowledge, whereas pure programmed calculations for the foremost half can't catch within and out significance within the data. Generally, instructive specialists are utilizing techniques, as an example, studies, interviews, center gatherings, and study hall exercises to collect data known with understudies' learning encounters. These ways area unit normally terribly tedious, thus can't be derived or rehashed with high return. the scale of such examinations is likewise for the most half restricted. Also, once angry concerning their encounters, understudies have to be compelled to speculate what they were considering and doing at some purpose before, which can have moved toward turning into clouded when a while. The rising fields of learning examination and instructive information mining (EDM) have focused on breaking down organized data got from course the board frameworks (CMS), study hall innovation use, or controlled web primarily based learning things to advise instructive basic leadership. In any case, as so much as we tend to may probably recognize, there is no exploration found to licitly mine and dissect student posted

content from uncontrolled areas on the social web with the affordable objective of understanding understudies' learning encounters. The examination objectives of this investigation area unit 1) to exhibit a piece method of on-line networking data sense-production for instructive purposes, incorporating each subjective examination and large scale data mining strategies as made public.

## **II. RELATED WORK**

Therefore cannot be derived or rehashed with high come back. the size of such examinations is likewise for the most 0.5 restricted. Also, once angry regarding their encounters, understudies have to be compelled to be compelled to take a position what they were considering and doing at some purpose before, which might have moved toward turning into clouded once a jiffy. The rising fields of learning examination and instructive information mining (EDM) have targeted on breaking down organized information got from course the board frameworks (CMS), study hall innovation use, or controlled web based totally learning things to advise instructive basic leadership. In any case, as most as we tend to tend to could most likely acknowledge, there is no exploration found to legitimately mine and dissect student posted content from uncontrolled areas on the social web with the cheap objective of understanding understudies' learning encounters. The examination objectives of this investigation square measure 1) to exhibit a bit technique of on-line networking information sense-production for instructive purposes, incorporating every subjective examination and largescale data processing ways as created public.

## **III. LITERATURE SURVEY**

Online administrations provide a selection of probabilities for understanding human conduct through the massive mix informational indexes that their activity gathers. the information sets they gather don't elementary partner model or mirror the globe occasions. Amid this paper we tend to tend to utilize information from Foursquare, a well-liked space entry profit, to contend for the significance of breaking down net primarily based life as associate open rather than representational framework. Drawing on logs of all Foursquare registration quite eight we've got a propensity toes we tend to feature four selections of 4 square's utilization: the connection among visiting and registration, occasion registration, business impetuses to entry, and in finish absurd registration These focuses demonstrate anyway mammoth information investigation is laid low with the simplest shopper uses to it informal organizations are place. we recommend that the composition and usage of viable Social Learning Analytics (SLA) blessing essential difficulties and open doors for each investigation and endeavor, in three very important regards. the essential is that the learning scene is dreadfully angry at the current, in no small 0.5 as a result of innovative drivers. on-line social learning is ascending as an interesting advancement for a scope of reasons, that we are going to normally survey, to propel the develop of social learning. The second take a look at is to identify contrasting kinds of SLA and their connected advances and employments. we are going to normally refer five categories of investigative in respect to on-line social taking in; these examination zone unit either naturally social or are liberal. This sets the scene for a 3rd take a look at, that of execution examination that have instructive associated ethical uprightness in an extremely setting were power and administration over information region unit at once of essential significance. We will normally think about a number of the contemplations that learning investigation incite, and advocate that Social Learning Analytics may provide routes that forward. We will in general available returning to the drivers and patterns, and consider future inevitabilities that we are going to normally might even observe unroll as SLA instruments and administrations mature. Microblogging might be associate fashionable innovation nose to nose to person communication applications that offers shoppers an opportunity to distribute on-line short instant messages (e.g., yet two hundred characters) incessantly through information superhighway, SMS, moment electronic informing customers, and then forth. Microblogging can be an honest instrument within the pace and has typically augmented outstanding enthusiasm from the instruction network. This paper proposes associate exceptional utilization of content order for two sorts of microblogging questions asked amid a space, particularly necessary (i.e., inquiries that the trainer desires to manage within the class) and digressive queries. data-based outcomes and examination demonstrate that abuse personalization next to address content lands up in higher arrangement exactitude than misuse question message alone. it's moreover helpful to use the association amongst queries and on the market address materials equally in light-weight of the very fact that the connection be tween's inquiries asked amid associate address.

additionally, experimental outcomes moreover demonstrate that the disposal of stop-words lands up in higher relationship estimation among queries and lands up in higher arrangement accuracy. On the contrary hand, fusing understudies' votes on the inquiries doesn't enhance order precision, but an analogous element has been perceived to be powerful in network question respondent conditions for surveying question quality.[3] This paper builds up the necessity for in conjunction with on-line character administration procure in school man coaching, as a chunk of creating ready understudies for returning into the work showcase. It talks regarding the impact of on-line info on business, and presents unique meeting learning with regard to planning and innovation college man understudies' on-line temperament administration rehearses. The paper contends for the necessity to demonstrate understudy's on-line life securing and proposes a specific orchestrate on-line character administration that will be incorporated into college man curriculum.

#### IV. ALGORITHM

Clustering K-implies bunching may be a reasonably unsupervised realizing, that is used after you have unlabeled data (i.e., data while not characterized categories or gatherings). the target of this calculation is to find bunches within the data, with the number of gatherings spoken to by the variable K. The calculation works iteratively to allot each data {point information} point to at least one of K bunches keen about the highlights that area unit given. data focuses area unit sorted keen about highlight closeness. the results of the K-implies bunching calculation are: The centroids of the K bunches, which might be used to mark new data Marks for the preparation data (every data {point information} point is relegated to a solitary group) As opposition characterizing bunches before taking a goose at the knowledge, bunching permits you to find and investigate the gatherings that have framed naturally. The "Picking K" section below portrays however the number of gatherings may be resolved. each center of mass of a bunch may be a gathering of highlight esteems that characterize the following gatherings. Inspecting the center of mass embrace hundreds may be used to subjectively translate what form of gathering every bunch speaks to. The K-implies grouping calculation utilizes unvarying refinement to deliver a final outcome. The calculation inputs area unit the number of bunches K and therefore the informational index. The informational index is Associate in Nursing accumulation of highlights for each data {point information} point. The calculations begin with starting assessments for the K centroids, which might either be haphazardly created or randomly chosen from the informational index. The calculation at that time emphasizes between 2 stages:

##### 4.1 Data Assignment Step

Each centroid defines one of the clusters. In this step, each data point is assigned to its nearest centroid, based on the squared Euclidean distance. More formally, if  $c_i$  is the collection of centroids in set C, then each data point x is assigned to a cluster based on

$$\operatorname{argmin}_{c_i \in C} \operatorname{dist}(c_i, x)^2$$

where  $\operatorname{dist}(\cdot)$  is the standard (L2) Euclidean distance. Let the set of data point assignments for each its cluster centroid be  $S_i$ .

##### 4.2 Centroid Update Step

In this step, the centroids are recomputed. This is done by taking the mean of all data points assigned to that centroid's cluster.

$$c_i = \frac{1}{|S_i|} \sum_{x_j \in S_i} x_j$$

The calculation emphasizes between stages 1 and 2 till a ceasing criterion is met (i.e., no data focuses modification teams, the whole lot of the separations is proscribed, or some most extreme variety of emphases is come back to). This calculation is ensured to satisfy to associate outcome. the result could be a section ideal (for example not very the foremost ideal result), implying that measuring quite one keep running of the calculation with randomized starting centroids could provide a superior result.

### 4.3 K Nearest Neighbors

A case is ordered by a larger part vote of its neighbors, with the case being doled out to the class most basic among its K closest neighbors estimated by a separation work. In the event that  $K = 1$ , at that point the case is basically relegated to the class of its closest neighbor.

#### Distance functions

Euclidean	$\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$
Manhattan	$\sum_{i=1}^k  x_i - y_i $
Minkowski	$\left( \sum_{i=1}^k ( x_i - y_i ^q) \right)^{1/q}$

It has to be compelled to likewise be detected that each one amongst the 3 separation measures area unit legitimate for nonstop factors. within the case of downright factors, the playacting separation should be used. It in addition raises the problem of institutionalization of the numerical factors somewhere within the vary of zero and one once there's a mix of numerical and every one out factor within the dataset.

#### Hamming Distance

$$D_H = \sum_{i=1}^k |x_i - y_i|$$

$$x = y \Rightarrow D = 0$$

$$x \neq y \Rightarrow D = 1$$

X	Y	Distance
Male	Male	0
Male	Female	1

Picking the ideal incentive for K is best done by first investigating the information. All in all, a huge K esteem is progressively exact as it decreases the general commotion yet there is no assurance. Cross-approval is another approach to reflectively decide a decent K esteem by utilizing a free dataset to approve the K esteem. Generally, the ideal K for most datasets has been between 3-10. That produces much preferable outcomes over 1NN.

### V. DATASET TO DETECT STUDENT PROBLEM

There were thirty-five,598 novel posts within the Purdue tweet gathering. we have a tendency to took associate discretionary example of one,000 posts, and located near to five plc. of those posts were talking regarding building problems. Our motivation here was to tell apart the modest range of posts that replicate planning understudies' problems. The contrasts between #engineering-Problems informational assortment and Purdue informational index is that the last contains loads smaller range of positive examples to be known, and its "others" category has increasingly differing substance. later on, to create the preparation set higher suits the Purdue informational index, we have a tendency to took

associate discretionary example of five,000 posts from the Purdue informational index, enclosed them into the two,785 engineering issues posts, and marked them as "others". below five plc. positive examples during this classification do not impact the adequacy of the ready model. we have a tendency to thus used the seven,785 posts as contribution to organize the multi-mark classifier. Since no extra human effort is needed, and classifier is extraordinarily effective as so much as calculation time, the model getting ready here caused no extra expense. Table one demonstrates the highest most plausible words in each classification positioned utilizing the contingent chance.

Category	Top 25 words
<i>Heavy Study Load</i>	hour, homework, exam, day, class, work, negtoken, problem, study, week, toomuch, all, lab, still, out, time, page, library, spend, today, long, school, due, engineer, already
<i>Lack of Social Engagement</i>	negtoken, Friday, homework, out, study, work, weekend, life, class, engineer, exam, drink, break, Saturday, people, social, lab, spend, tonight, watch, game, miss, party, sunny, beautiful
<i>Negative Emotion</i>	hate, f***, shit, exam, negtoken, week, class, hell, engineer, suck, study, hour, homework, time, equate, FML, lab, sad, bad, day, feel, tired, damn, death, hard
<i>Sleep Problems</i>	sleep, hour, night, negtoken, bed, allnight, exam, homework, nap, coffee, time, study, more, work, class, dream, ladyengineer, late, week, day, long, morning, wake, awake, no-sleep
<i>Diversity Issues</i>	girl, class, only, negtoken, guy, engineer, Asia, professor, speak, English, female, hot, kid, more, toomuch, walk, people, teach, understand, chick, China, foreign, out, white, black

Table 1: Sample dataset

## VI. RESULT

In total, we have tested 200 posts which were collected randomly from Facebook. Among those 200 posts, 118 are political and 82 are non-political. Table VI shows the detection results. Among 118 political posts our algorithm correctly detects 99 posts and among 82 non-political posts it correctly detects 67 posts. So, the overall accuracy is 83%. the table shows the analysis and measurement of posts.

Class Level	Previous System	Proposed System
Heavy Study Load	67	78
Lack of Social Engagement	45	54
Negative Emotion	43	50
Sleep Problem	45	56
Diversity Issues	20	30

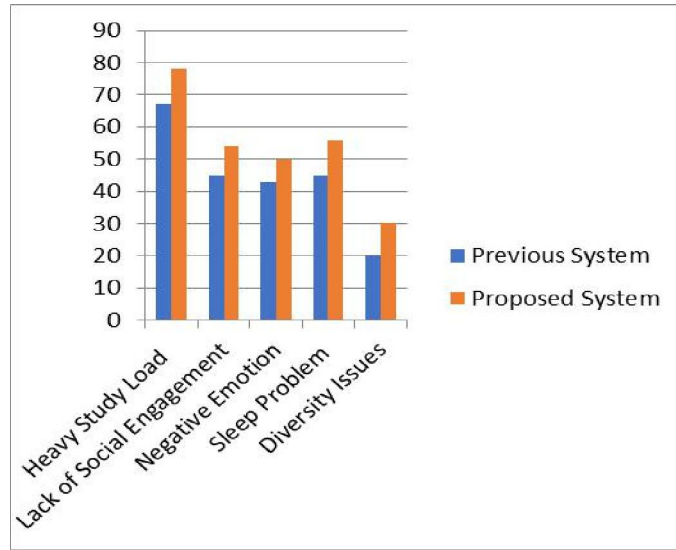


Figure: Throughput Result Graph

**Table:** Analysis and Measurement Posts.

Class Level	Previous System	Proposed System
Heavy Study Load	257	344
Lack of Social Engagement	135	167
Negative Emotion	52	65
Sleep Problem	154	142
Diversity Issues	47	37

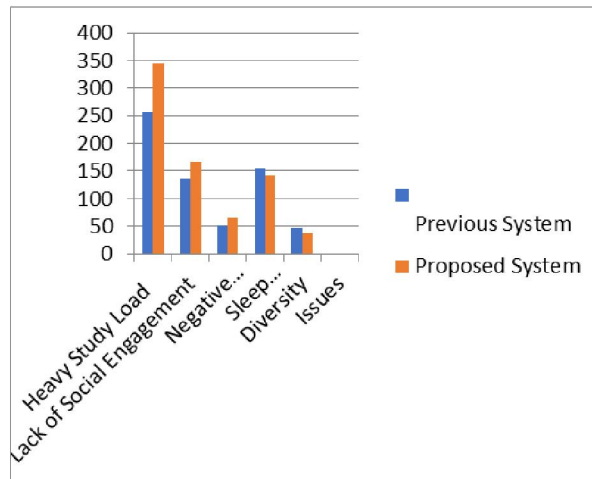


Figure: Efficiency Result Graph

**Table:** Analysis and Measurement Posts.

Class Level	Previous System	Proposed System
Heavy Study Load	87	97
Lack of Social Engagement	76	89
Negative Emotion	86	77
Sleep Problem	40	56
Diversity Issues	50	60

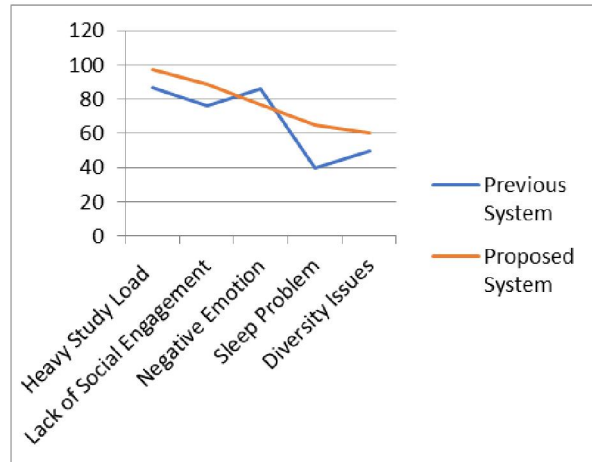


Figure: Accuracy Result Graph

**Table:** Analysis and Measurement Posts

Class Level	Previous System	Proposed System
Heavy Study Load	56	67
Lack of Social Engagement	60	76
Negative Emotion	65	70
Sleep Problem	50	60
Diversity Issues	34	47

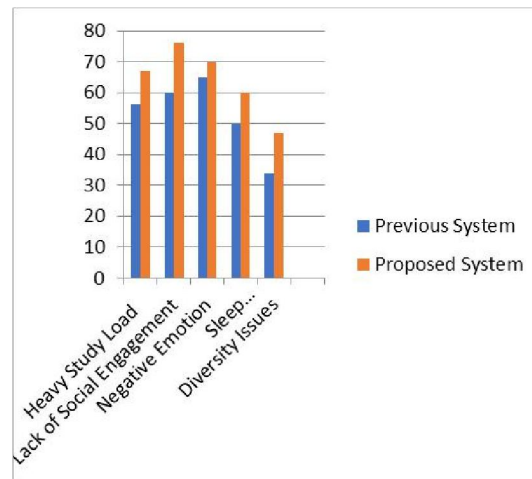


Figure: Speed Result Graph



#### **VII. CONCLUSION**

This examination investigates the already instrumented area on Facebook thus on comprehend building understudies' encounters, incorporating each subjective techniques and massive scale data mining techniques. In our investigation, through a subjective substance examination, we have a tendency to found that building understudies square measure to an excellent extent battling with the substantial investigation load, and don't seem to be able to manage it effectively. Overwhelming investigation load prompts varied results together with absence of social commitment, rest problems, and alternative mental and physical medical problems. varied understudies feel building is exhausting and laborious, that prompts absence of inspiration to concentrate and negative feelings. good selection problems to boot uncover culture clashes and culture generalizations existing among planning understudies. increasing over the subjective bits of data, we have a tendency to dead and assessed a multi-mark classifier to acknowledge planning understudy problems from Purdue University. This symbol is connected as a checking system to tell apart at risk understudies at a specific school over the long-term while not rehashing the manual work each currently and once more.

#### **REFERENCES**

- [1]. P. D. Turney, Thumbs up or thumbs down?: semantic orientation applied to unsupervised classification of reviews, in Proceedings of the 40th annual meeting on association for computational linguistics, pp. 417-424, Association for Computational Linguistics, 2002.
- [2]. A. Pak and P. Paroubek. Twitter as a Corpus for Sentiment Analysis and Opinion Mining". In Proceedings of the Seventh Conference on International Language Resources and Evaluation, 2010, pp.1320-1326.
- [3]. Po-Wei Liang, Bi-Ru Dai, Opinion Mining on Social MediaData", IEEE 14th International Conference on Mobile Data Management, Milan, Italy, June 3- 6, 2013, pp 91-96, ISBN: 978-1-494673-6068-<http://doi.ieeecomputersociety.org/10.1109/MDM.2013>.
- [4]. T. T. Thet, I.-C. Na, C. S. Khoo, and S. Shakhikumar, "Sentiment analysis of movie reviews on discussion boards using a linguistic approach," in Proceedings of the 1st international CIKM workshop on Topic-sentiment analysis for mass opinion. ACM, 2009, pp. 81-84.
- [5]. Hussein, D.-M.E.D.M. A survey on sentiment analysis challenges. Journal of King Saud University Engineering Sciences (2016), <http://dx.doi.org/10.1016/j.jksues.2016.04.002>
- [6]. A. Kowcika and Aditi Guptha Sentiment Analysis for Social Media, International Journal of Advanced Research in Computer Science and Software Engineering, 216-221, Volume 3, Issue 7, July 2013.
- [7]. G. Vinodini and RM.Chandrashekar, Sentiment Analysis and Opinion Mining: A Survey, International Journal of Advanced Research in Computer Science and Software Engineering, 283-294, Volume 2, Issue 6, June 2012.
- [8]. Apoorv Agarwal Boyi Xie Ilia Vovsha Owen Rambow Rebecca Passonneau, Sentiment Analysis of Twitter Data, Columbia University, New York.
- [9]. Pablo Gamallo and Marcos Garcia A Naive-Bayes Strategy for Sentiment Analysis on English Posts Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014), pages 171-175, Dublin, Ireland, Aug 23-24 2014.
- [10]. Harry Zhang "The Optimality of Naive Bayes", FLAIRS2004 conference. (available online: PDF (<http://www.cs.unb.ca/profs/hzhang/publications/FLAIRS04ZhangH.pdf>)).