

Mining Challengers from Hefty Amorphous Datasets

SYEDA FATIHA BUTTUL¹, K. SHILPA²

¹PG Scholar, Dept of CSE, Shadan Women's College of Engineering and Technology, Hyderabad, TS, India,

²Associate Professor, Dept of CSE, Shadan Women's College of Engineering and Technology, Hyderabad, TS, India,

Abstract: A protracted line of studies has exhibited the vital hugeness of distinguishing and checking an association's rivals. Spurred with the aid of this trouble, the selling and administration group have targeting precise electronics for contender distinguishing proof and in reckoning up on procedures for breaking down recognized contenders. Surviving examination at the preceding has targeting mining similar articulations (e.g. "Thing A is surpassing to Item B") from the net or other printed sources. Even anyway such articulations can undoubtedly be markers of intensity, they are truant in divers spaces. For instance, reflect on consideration on the space of excursion bundles (e.g flight-inn auto blends). For this situation, things had allotted name by using which they'll be questioned or contrasted and every different. Further, the excessive indication of formal near confirmation can differ extraordinarily crosswise over areas. For eg, when contrasting logo names at the hale level (e.g. "Google versus Yahoo" or "Sony versus Panasonic"). It's miles to make sure that comparable precedents may be discovered by generally analyze the net. In any case, it is everything except hard to apprehend standard areas where such evidence is greatly limited, for eg: boot, adornments, lodgings, eateries, and furniture. Inspired by these weaknesses, other rationalization of the forcefulness among things, in angle of the sell parts that each can cover.

Keywords: Data Mining, Web Mining, Information Search And Retrieval, Electronic Commerce.

I. INTRODUCTION

Along line of research has shown the key significance of recognizing and observing a company's rivals [2]. Propelled by this issue, the promoting and administration network have concentrated on observational strategies for contender distinguishing proof [3], [4], [5], [6], [7], and on techniques for breaking down known contenders [8]. Surviving exploration on the previous has concentrated on mining similar articulations (e.g. "Thing An is superior to Item B") from the Web or other printed sources [9], [10], [12], [13],[14]. Despite the fact that such articulations can for sure be pointers of aggressiveness, they are missing in numerous spaces. For example, think about the space of excursion bundles (e.g flight-lodging auto blends). For this situation, things have no doled out name by which they can be questioned or contrasted and each other. Further, the recurrence of literary relative proof can fluctuate incredibly crosswise over spaces. For instance, when looking at mark names at the firm level (e.g. "Google versus Yahoo" or "Sony versus Panasonic"), it is without a doubt likely that relative examples can be found by basically questioning the web. Be that as it may, it is anything but difficult to recognize standard areas where such confirmation is to a great degree rare, for example, shoes, jewelry, inns, eateries, and furniture. Spurred by these deficiencies, we propose another formalization of the intensity between two things, in

light of the market portions that they can both cover. Formally:

Definition 1. [Competitiveness]: Let U be the number of inhabitants in every single conceivable client in a given market. We think about that as a thing I covers a client $u \in U$ in the event that it can cover the greater part of the client's prerequisites. At that point, the intensity between two things I, j is relative to the quantity of clients that they can both cover. Our aggressiveness worldview depends on the accompanying perception: the intensity between two things depends on whether they go after the consideration and business of similar gatherings of clients (i.e. a similar market fragments). For instance, two eateries that exist in various nations are clearly not aggressive, since there is no cover between their objective gatherings. Consider the illustration appeared in Fig.1.

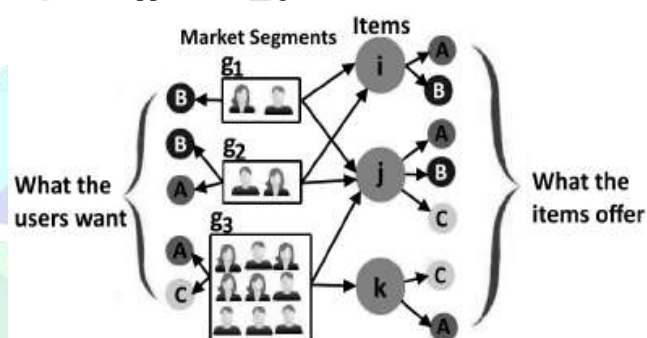


Fig.1. A (simplified) example of our competitiveness paradigm.

The figure illustrates the competitiveness between three items i, j and k . Each item is mapped to the set of features that it can offer to a customer. Three features are considered in this example: A, B and C. Even though this simple example considers only binary features (i.e. available/not available), our actual formalization accounts for a much richer space including binary, categorical and numerical features. The left side of the figure shows three groups of The figure delineates the aggressiveness between three things I, j and k . Everything is mapped to the arrangement of highlights that it can offer to a client. Three highlights are considered in this illustration: A, B and C. Despite the fact that this basic case thinks about just double highlights (i.e. accessible/not accessible), our real formalization represents a substantially more extravagant space including double, all out and numerical highlights. The left half of the figure indicates three gatherings of clients g_1, g_2 , and g_3 . Each gathering speaks to an alternate market portion. Clients are gathered in view of their inclinations concerning the highlights. For instance, the clients in g_2 are just inspired by highlights An and B. We watch that things I and k are not aggressive, since they just don't speak to similar gatherings of clients. Then again, j

contends with both I (for bunches g_1 and g_2) and k (for g_3). At long last, an intriguing perception is that j vies for 4 clients with I and for 9 clients with k. At the end of the day, k is a more grounded contender for j, since it guarantees a significantly bigger part of its piece of the overall industry than I.

This case represents the perfect situation, in which we approach the entire arrangement of clients in a given market, and in addition to particular market sections and their prerequisites. By and by, be that as it may, such data isn't accessible. With a specific end goal to conquer this, we portray a strategy for registering every one of the sections in a given market in light of mining extensive audit datasets. This technique enables us to operationalize our meaning of intensity and address the issue of finding the best k contenders of a thing in any given market. As we appear in our work, this issue presents noteworthy computational difficulties, particularly within the sight of expansive datasets with hundreds or thousands of things, for example, those that are frequently found in standard areas. We address these difficulties by means of an exceptionally adaptable system for top-k calculation, including a proficient assessment calculation and a proper file. Our work makes the accompanying commitments:

- A formal meaning of the aggressiveness between two things, in view of their interest to the different client sections in their market. Our approach beats the dependence of past work on rare near proof mined from content.
- A formal technique for the ID of the distinctive kinds of clients in a given market, and additionally for the estimation of the level of clients that have a place with each sort.
- An exceptionally adaptable structure for finding the best k contenders of a given thing in extensive datasets.

II. DEFINING COMPETITIVENESS

The commonplace client session on an audit stage, for example, Yelp, Amazon or Trip Advisor, comprises of the accompanying advances: 1) Specify every required component in a question. 2) Submit the inquiry to the site's web crawler and recover the coordinating things. 3) Process the surveys of the returned things and settle on a buy choice. In this setting, things that cover the client's necessities will be incorporated into the internet searcher's reaction and will seek her consideration. Then again, non-covering things won't be considered by the client and, in this way, won't have an opportunity to contend. Next, we display an illustration that stretches out this basic leadership procedure to a multi client setting. Consider a straightforward market with 3 lodgings I, j, k and 6 double highlights: bar, breakfast, exercise center, stopping, pool, wi-fi. Table 1 incorporates the estimation of every inn for each element. In this basic case, we accept that the market incorporates 6 fundamentally unrelated client fragments (types). Each fragment is spoken to by an inquiry that incorporates the highlights that are important to the clients incorporated into the portion. Data on each portion is given in Table 2. For example, the primary fragment incorporates 100 clients who are keen on stopping and wi-fi, while the second section incorporates 50 clients who are just inspired by stopping.

TABLE I: Hotels and Their Features

Name	Bar	Breakfast	Gym	Parking	Pool	Wi-Fi
Hilton	Yes	No	Yes	Yes	Yes	Yes
Marriot	Yes	Yes	No	Yes	Yes	Yes
Westin	No	Yes	Yes	Yes	No	Yes

Keeping in mind the end goal to quantify the opposition between any two lodgings, we have to recognize the quantity of clients that they can both fulfill. The outcomes are appeared in Table 3. The Hilton and the Marriot can cover portions q_1 , q_3 , and q_4 . Hence, they seek $(100 + 50 + 60)/660 \approx 32\%$ of the whole market. We watch this is the least aggressiveness accomplished for any combine, despite the fact that the two inns are likewise the most comparative. Truth be told, the most elevated aggressiveness is seen between the Marriot and the Westin, that seek 70% of the market. This is a basic perception that exhibits that comparability is certifiably not a decent intermediary for intensity. The clarification is natural. The accessibility of both a pool and a bar makes the Hilton and the Marriot more like each other and less like the Westin. In any case, neither of these highlights affects intensity. To begin with, the pool highlight isn't required by any of the clients in this market. Second, despite the fact that the accessibility of a bar is required by portion q_6 , none of the three lodgings can cover each of the three of this present fragment's prerequisites. Accordingly, none of the lodgings seek this specific section. Another instinctive perception is that the measure of the portion directly affects intensity. For instance, despite the fact that the Westin has a similar number of sections (4) with the other two lodgings, its aggressiveness with the Marriot is fundamentally higher. This is because of the extent of the q_5 portion, which is more than twofold the measure of q_4 .

TABLE II: Customer Segments

ID	Segment Size	Features of Interest
q_1	100	(parking, wi-fi)
q_2	50	(parking)
q_3	60	(wi-fi)
q_4	120	(gym, wi-fi)
q_5	250	(breakfast, parking)
q_6	80	(gym, bar, breakfast)

TABLE III: Common Segments For Restaurant Pairs

Restaurant Pairs	Common Segments	Common %
Hilton, Marriot	(q_1, q_2, q_3)	32%
Hilton, Westin	(q_1, q_2, q_3, q_4)	50%
Marriot, Westin	(q_1, q_2, q_3, q_5)	70%

The above illustration is restricted to double highlights. In this basic setting, it is insignificant to decide whether two things can both cover a component. Be that as it may, as we examine in detail in Section 2.1, the things in a market can have diverse kinds of highlights (e.g. numeric) that might be just somewhat secured by two things. Formally, let $p(q)$ be the level of clients spoke to by an inquiry q and let $V I, j q$ be the match astute scope offered by two things I and j to the space characterized by the highlights in q . At that point, we characterize the intensity amongst I and j in a market with a component subset F as takes after:

$$C_{\mathcal{F}}(i, j) = \sum_{q \in 2^{\mathcal{F}}} p(q) \times V_{i,j}^q, \quad (1)$$

This definition has an unmistakable probabilistic translation: given two things I, j , their aggressiveness $C_{\mathcal{F}}(i, j)$ speaks to the likelihood that the two things are incorporated into the thought set of an arbitrary client. This new definition has coordinate ramifications for customers, who regularly depend on proposal frameworks to enable them to pick one of a few applicant items. The capacity to quantify the aggressiveness between two things empowers the suggestion framework to deliberately choose the request in which things ought to be prescribed or the arrangements of things that ought to be incorporated together in a gathering proposal. For example, if an irregular client u indicates enthusiasm for a thing I , at that point she is additionally prone to be keen on the things with the most noteworthy $C_{\mathcal{F}}(i, \cdot)$ qualities. Such focused things are probably going to meet the criteria fulfilled by I and even cover extra parts of the component space. What's more, as the client u rates more things and the framework picks up a more exact perspective of her prerequisites, our intensity measure can be inconsequentially changed in accordance with consider just those highlights from \mathcal{F} (and just those esteem interims inside each component) that are important for u . This intensity based suggestion worldview is a takeoff from the standard approach that changes the weight (significance) of a thing j for a client u in light of the rating that u submits for things like j . As examined, this approach overlooks that (i) the similitude might be because of unimportant or paltry highlights and (ii) for a client who loves a thing I , a thing j that is far better than I with deference than the client's necessities (and in this way very unique) is a superior proposal hopeful than a thing j' that is profoundly like I . In the accompanying two areas we depict the calculation of the two essential segments of intensity: (1) the pairwise scope $V_{i,j}^q$ of an inquiry that incorporates parallel, unmitigated, ordinal or numeric highlights, and (2) the rate $p(q)$ of clients spoke to by each question q .

A. Pairwise Coverage

We start by characterizing the pairwise scope of a solitary component f . We at that point characterize the pairwise scope of a whole inquiry of highlights q .

Definition 2. [Pairwise Feature Coverage]: We characterize the pairwise scope $V_{i,j}^f$ of a component f by two things I, j as the level of every single conceivable estimation of f that can be secured by both I and j . Formally, given the arrangement of every conceivable incentive V^f for f , we characterize:

$$V_{i,j}^f = \frac{|\{v \in V^f : v \leq f[i] \wedge v \leq f[j]\}|}{|values(f)|} \quad (2)$$

where $v \leq f[i]$ speaks to that v is secured by the estimation of thing I for include f . Next, we depict the calculation of $V_{i,j}^f$ for various kinds of highlights.

Binary and Categorical Features: Categorical highlights take at least one qualities from a limited space. Cases of

single esteem highlights incorporate the brand of a computerized camera or the area of an eatery. Cases of multi-esteem highlights incorporate the luxuries offered by a lodging or the sorts of cooking offered by an eatery. Any clear cut component can be coded by means of an arrangement of paired highlights, with every double element demonstrating the (absence of) scope of one of the first element's conceivable qualities. In this basic setting, the element can be completely secured (if $f[i] = f[j] = 1$ or, identically, $f[i] \times f[j] = 1$), or not secured by any means. Formally, the pairwise scope of a parallel element f by two things i, j can be processed as takes after:

$$V_{i,j}^f = f[i] \times f[j] \quad (3)$$

Numeric Features: Numeric highlights take esteems from a pre-characterized run. From this time forward, without loss of sweeping statement, we consider numeric highlights that take esteems in $[0, 1]$, with higher qualities being best. The pairwise scope of a numeric element f by two things I and j can be effortlessly processed as the littlest (most exceedingly terrible) esteem accomplished for f by either thing. For example, think about two eateries I, j with values 0.8 and 0.5 for the component nourishment quality. Their pairwise scope in this setting is 0.5. Thoughtfully, the two things will seek any client who acknowledges a quality $y \leq 0.5$. Clients with higher measures would dispose of eatery j , which will never have an opportunity to seek their business. Formally, the pairwise scope of a numeric element f by two things I, j can be registered as takes after:

$$V_{i,j}^f = \min(f[i], f[j]) \quad (4)$$

Ordinal Features: Ordinal highlights take esteems from a limited arranged rundown. A trademark illustration is the well known five star scale used to assess the nature of an administration or item. For instance, think about that the estimations of two things I and j on the 5-star rating scale are $\star\star$ and $\star\star\star$, separately. Clients that request no less than 4 stars won't think about both of the two things, while clients that request no less than 3 stars will just think about thing j . The two things will therefore vie for all clients that will acknowledge 1 or 2 stars. In this way, as on account of numeric highlights, the pairwise scope for ordinal highlights is controlled by the most noticeably awful of the two qualities. In this illustration, given that the two things go after 2 of the 5 levels of the ordinal scale (1 and 2 stars), their aggressiveness is corresponding to $2/5 = 0.4$. Formally, the pairwise scope of an ordinal component f by two things I, j can be figured as takes after:

$$V_{i,j}^f = \frac{\min(f[i], f[j])}{|V^f|} \quad (5)$$

Pairwise Scope Of An Element Question: We presently talk about how scope can be stretched out to the inquiry level. Figure 2 pictures an inquiry q that incorporates two numeric highlights f_1 and f_2 . The figure likewise incorporates two focused things I and j , situated by their qualities for the two highlights: $f_1[i] = 0.3, f_2[i] = 0.3, f_1[j] = 0.2$, and $f_2[j] = 0.7$. We watch that the level of the 2-dimensional space that every thing covers is identical to the zone of the square shape

characterized by the start of the two tomahawks (0,0) and the thing's qualities for f1 and f2. For instance, the secured territory for thing I is $0.3 \times 0.3 = 0.09$, equivalent to 9% of the whole space. Essentially, the pairwise scope gave by the two things is equivalent to $0.2 \times 0.3 = 0.06$ (i.e. 6% of the market). Per our case, the pairwise scope of a given inquiry q by two things I, j can be estimated as the volume of the hyper-square shape characterized by the pairwise scope gave by the two things to each element $f \in q$. Formally:

$$V_{IJ}^q = \prod_{f \in q} V_{IJ}^f \quad (6)$$

Eq5 enables us to register the pairwise scope of any inquiry of highlights, as required by the meaning of intensity in Eq. 1.

B. Estimating Query Probabilities

The meaning of intensity given in Eq. 1 considers the likelihood $p(q)$ that an arbitrary client will be spoken to by a particular inquiry of highlights q, for each conceivable question $q \in 2^F$. In this area, we portray how these probabilities can be assessed from genuine information. Highlight inquiries are an immediate portrayal of client inclinations. In a perfect world, we would approach the question logs of the stage's (e.g. Amazon's or Trip Advisor's) web crawler. By and by, be that as it may, the delicate and restrictive nature of such data makes it difficult for firms to share openly. Hence, we outline an estimation procedure that lone expects access to a bottomless asset: client surveys. Each survey incorporates a client's sentiments on a specific subset of highlights of the investigated thing. Surviving examination has more than once approved the utilization of audits to assess client inclinations as for various highlights in different spaces, for example, telephone applications [15], motion pictures [16], gadgets [17], and lodgings [18]. A unimportant approach is to assess the interest for each element independently, and after that total the individual appraisals at the subset level. Notwithstanding, this approach accepts highlight autonomy, a solid suspicion that would first must be approved crosswise over areas. To maintain a strategic distance from this supposition and catch conceivable element relationships, we consider every one of the highlights specified in each audit as a solitary question. We at that point register the recurrence of each question q in our survey corpus R, and separate it by the whole of the frequencies everything being equal. This gives us a gauge of the likelihood that an irregular client will be keen on precisely the arrangement of highlights incorporated into q. Formally:

$$p(q) = \frac{\text{freq}(q, R)}{\sum_{q \in 2^F} \text{freq}(q, R)} \quad (7)$$

In a perfect world, we would approach the arrangement of prerequisites of each conceivable client in presence. The greatest probability gauge of Eq. 6 would then register the correct event likelihood of any inquiry q. While this sort of worldwide access is doubtful, Eq. 6 can in any case convey precise appraisals if the quantity of surveys in R is sufficiently vast to precisely speak to the client populace. The handiness of the gauge or is in this manner controlled

by a basic inquiry: what number of audits do we have to accomplish exact evaluations? We address this inquiry in Section 5.7 of the tests, where we show our outcomes on datasets from various spaces.

C. Extending our Competitiveness Definition

Highlight Uniformity: Our aggressiveness definition accepts that client prerequisites are consistently circulated inside the esteem space of each component. This presumption enables us to fabricate a computational model for intensity, however by and by it may not generally be valid. For example, the quantity of clients requesting quality in $[0, 0.1]$ may be not quite the same as those requesting an incentive in $[0.4, 0.5]$. Besides, for absence of more exact data, it gives a traditionalist lower bound of our model's actual adequacy: approaching the dispersion of enthusiasm inside each element could just enhance the nature of our outcomes. On the off chance that such data was without a doubt accessible, at that point the gullible approach is to consider all conceivable intrigue interims mixes for every single conceivable question. From this time forward, we allude to these as expanded questions. Plainly, the quantity of conceivable broadened inquiries is exponential and renders the computational cost of any assessment calculation restrictive. This constraint can be tended to by arranging the dataset into a multi-dimensional lattice, where each component speaks to an alternate measurement. Every cell in the network speaks to an alternate broadened inquiry (i.e. an arrangement of highlights and an intrigue interim for each element). We would then be able to register the intensity between two things by essentially tallying the quantity of information focuses that fall in the cells that they can both cover.

We can likewise pre-process the totals of every cell disconnected with the prefix-entirety cluster procedure [19], and decrease the space unpredictability by means of approximations [20], [21] or multidimensional histograms [22], [23]. A parameter of the framework development process is the cell measure, with bigger cells giving up precision for proficiency. Practically speaking, this parameter will be controlled by the granularity of the information, and the professional's computational imperatives. **Highlight Importance:** A second suspicion of our intensity definition is that every one of the highlights in a question q are similarly imperative. Be that as it may, a client who presents the inquiry may think more about f1 than for f2. Similarly as with the instance of highlight consistency, the thought of such weights requires the accessibility of fitting information that is infrequently accessible by and by. In any case, we can address this constraint by broadening our meaning of pairwise scope. For example, consider that the component weights are in $[0, 1]$ and that the weights for f1 and f2 are $w_1 = 0.8$ and $w_2 = 0.4$, individually. We are then given two things I, j to such an extent that: $f_1[i] = 0.5, f_2[i] = 0.3, f_1[j] = 0.5, f_2[j] = 0.6$. According to our underlying definition, the pairwise scope of the 2-dimensional space by the two things is $\min(0.5, 0.5) \times \min(0.3, 0.6) = 0.5 \times 0.3 = 0.15$. On the off chance that we consider the component weights, the calculation progresses toward becoming: $(w_1 \times 0.5) \times (w_2 \times 0.3) = 0.048$. Formally, this expansion means the presentation of the

component weight as a multiplier for the right-hand side of Eq. 3. Note that, while this illustration incorporates just numeric highlights, a similar expansion for unmitigated and ordinal properties inconsequentially takes after.

III. FINDING THE TOP-K COMPETITORS

Given the meaning of the intensity in Eq. 1, we ponder the characteristic issue of finding the best k contenders of a given thing. Formally:

Problem1. [Top-k Competitors Problem]: We are presented with a market with a set of n items I and a set of features F . Then, given a single item $i \in I$, we want to identify the k items from I that maximize $C_F(i, \cdot)$. A credulous calculation would register the intensity amongst I and each conceivable hopeful. The intricacy of this savage power technique is plainly $\Theta(|I| \times |F| \times |I|)$, which can be effortlessly commanded by the power set factor and, as we exhibit in our investigations, is illogical for expansive datasets. One choice could be to play out the guileless calculation in an appropriated form. Indeed, even for this situation, be that as it may, we would require one string for each of the n^2 sets. This is a long way from inconsequential, in the event that one considers that n could gauge in the several thousands. Likewise, an innocent Map Reduce usage would confront the bottleneck of going everything through the reducer to represent the self-join incorporated into the calculation. By and by, the self join would need to be actualized by means of a modified method for lessen side joins, which is a non-trifling and very costly task [24]. These issues rouse us to present CMiner, a proficient correct calculation for Problem 1. With the exception of the production of our ordering instrument, each other part of CMiner can likewise be fused in a parallel arrangement. Initially, we characterize the idea of thing strength, which will help us in our investigation:

Definition 3. [Item Dominance]: Think about a market with an arrangement of things I and an arrangement of highlights F . At that point, we say that a thing $I \in I$ command san other thing $j \in I$, if $f[i] \geq f[j]$ for each component $f \in F$. Reasonably, a thing rules another on the off chance that it has better or equivalent qualities crosswise over highlights. We watch that, per Eq.1, any thing I that rules j additionally accomplishes the greatest conceivable aggressiveness with j , since it can cover the necessities of any client secured by j . This persuades us to use the horizon of the whole arrangement of things I . The horizon is an all around contemplated idea that speaks to the subset of focuses in a populace that are not ruled by some other point [25]. We allude to the horizon of an arrangement of things I as Sky (I).

The CMiner Algorithm: Next, we display CMiner, a correct calculation for finding the best k contenders of a given thing. Our calculation influences utilization of the horizon to pyramid keeping in mind the end goal to decrease the quantity of things that should be considered. Given that we just think about the best k contenders, we can incrementally process the score of every applicant and stop when it is ensured that the best k have risen. The pseudo code is given in Algorithm 1.

Discussion of CMiner: The information incorporates the arrangement of things I , the arrangement of highlights F , the thing of intrigue I , the number k of best contenders to recover, the set Q of inquiries and their probabilities, and the horizon pyramid DI . The calculation first recovers the things that overwhelm I , by means of $masters(i)$ (line 1). These things have the most extreme conceivable intensity with I . In the event that in any event k such things exist, we report those and finish up (lines 2-4). Else, we add them to $TopK$ and decrement our financial plan of k in like manner (line 5). The variable LB keeps up the most minimal lower bound from the current $topk$ set (line 6) and is utilized to prune hopefuls. In line 7, we introduce the arrangement of hopefuls X as the association of things in the first layer of the pyramid and the arrangement of things overwhelmed by those as of now in the $TopK$. This is accomplished through calling $GETSLAVES(TopK, DI)$. In each cycle of lines 8-17, CMiner encourages the arrangement of competitors X to the $UPDATETOPK()$ schedule, which prunes things in light of the LB edge. It at that point refreshes the $TopK$ set by means of the $MERGE()$ work, which identifies the things with the most astounding aggressiveness from $TopK \cup X$. This can be accomplished in straight time, since both X and $TopK$ are arranged. In line 13, the pruning limit LB is set to the most exceedingly terrible (least) score among the new $TopK$. At long last, $GETSLAVES()$ is utilized to extend the arrangement of hopefuls by including things that are ruled by those in X .

Discussion of UPDATETOPK(): This normal procedures the hopefuls in X and finds at most k competitors with the most noteworthy aggressiveness with I . The routine uses an information structure nearby $Top K$, actualized as an acquainted cluster: the score of every hopeful fills in as the key, while its id fills in as the esteem. The cluster is key-arranged, to encourage

Algorithm 1 CMiner

Input: Set of items I , Item of interest $i \in I$, feature space F , Collection $Q \in 2^F$ of queries with non-zero weights, skyline pyramid D_I , int k
Output: Set of top- k competitors for i

```

1:  $TopK \leftarrow masters(i)$ 
2: if ( $k \leq |TopK|$ ) then
3:   return  $TopK$ 
4: end if
5:  $k \leftarrow k - |TopK|$ 
6:  $LB \leftarrow -1$ 
7:  $X \leftarrow GETSLAVES(TopK, D_I) \cup D_I[0]$ 
8: while ( $|X| \neq 0$ ) do
9:    $X \leftarrow UPDATETOPK(k, LB, X)$ 
10:  if ( $|X| \neq 0$ ) then
11:     $TopK \leftarrow MERGE(TopK, X)$ 
12:    if ( $|TopK| = k$ ) then
13:       $LB \leftarrow WORSTIN(TopK)$ 
14:    end if
15:     $X \leftarrow GETSLAVES(X, D_I)$ 
16:  end if
17: end while
18: return  $TopK$ 

```

```

19: Routine UPDATETOPK( $k, LB, \mathcal{X}$ )
20:  $localTopK \leftarrow \emptyset$ 
21:  $low(j) \leftarrow 0, \forall j \in \mathcal{X}$ .
22:  $up(j) \leftarrow \sum_{q \in Q} p(q) \times V_{j,i}^q, \forall j \in \mathcal{X}$ .
23: for every  $q \in Q$  do
24:    $maxV \leftarrow p(q) \times V_{i,i}^q$ 
25:   for every item  $j \in \mathcal{X}$  do
26:      $up(j) \leftarrow up(j) - maxV + p(q) \times V_{i,j}^q$ 
27:     if ( $up(j) < LB$ ) then
28:        $\mathcal{X} \leftarrow \mathcal{X} \setminus \{j\}$ 
29:     else
30:        $low(j) \leftarrow low(j) + p(q) \times V_{i,j}^q$ 
31:        $localTopK.update(j, low(j))$ 
32:       if ( $|localTopK| \geq k$ ) then
33:          $LB \leftarrow WORSTIN(localTopK)$ 
34:       end if
35:     end if
36:   end for
37:   if ( $|\mathcal{X}| \leq k$ ) then
38:     break
39:   end if
40: end for
41: for every item  $j \in \mathcal{X}$  do
42:   for every remaining  $q \in Q$  do
43:      $low(j) \leftarrow low(j) + p(q) \times V_{i,j}^q$ 
44:   end for
45:    $localTopK.update(j, low(j))$ 
46: end for
47: return TOPK( $localTopK$ )

```

the computation of the k best items. The structure is automatically truncated so that it always contains at most k items. In lines 21-22 we initialize the lower and upper bounds. For every item $j \in \mathcal{X}$, $low(j)$ maintains the current competitiveness score of j as new queries are considered, and serves as a lower bound to the candidate's actual score. Each lower bound $low(j)$ starts from 0, and after the completion of UPDATETOPK(), it includes the true competitiveness score $C_{\mathcal{F}(i,j)}$ of candidate j with the focal item i . On the other hand, $up(j)$ is an optimistic upper bound on j 's competitiveness score. Initially, $up(j)$ is set to the maximum possible score (line 22). This is equal to $\sum_{q \in Q} p(q) \times V_{i,i}^q$, where $V_{i,i}^q$ is simply the coverage provided exclusively by i to q . It is then incrementally reduced toward the true $C_{\mathcal{F}(i,j)}$ value as follows. For every query $q \in Q$, $maxV$ holds the maximum possible competitiveness between item i and any other item for that query, which is in fact the coverage of i with respect to q . Then, for each candidate $j \in \mathcal{X}$, we subtract $maxV$ from $up(j)$ and then add to it the actual competitiveness between i and j for query q . If the upper bound $up(j)$ of a candidate j becomes lower than the pruning threshold LB , then j can be safely disqualified (lines 27-29).

Otherwise, $low(j)$ is updated and j remains in consideration (lines 30-31). After each update, the value of LB is set to the worst score in $localTopK$ (lines 32-33), to employ stricter pruning in future iterations. In the event that the quantity of competitors $|\mathcal{X}|$ turns out to be less or equivalent to k (line 37), the circle over the inquiries stops. This is an early-halting model: since we will likely recover the best k competitors in \mathcal{X} , having $|\mathcal{X}| \leq k$ implies that every outstanding hopeful ought to be returned. In lines 41-46 we finish the aggressiveness calculation of the rest of the

competitors and refresh local Topk in like manner. This happens after the finish of the first circle, keeping in mind the end goal to dodge pointless bound-checking and enhance execution. Unpredictability: If the thing of intrigue I is commanded by in any event k things, at that point these will be returned by masters(i). This progression should be possible in $O(k)$, by iteratively recovering k things that overwhelm I . Something else, the multifaceted nature of CMiner is controlled by UPDATETOPK(), which relies upon the quantity of things in the applicant set \mathcal{X} . In its least complex frame, in the k -th call of the strategy, the applicant set contains the whole k -th horizon layer, $\mathcal{D}_T[k]$.

According to Bentley et al. for n uniformly-distributed d -dimensional data points (items), the expected size of the skyline (1st layer) is $|\mathcal{D}_T[0]| = \Theta(\frac{\ln^{d-1} n}{(d-1)!})$. UPDATETOPK() will be called at most k times, each time fetching (at least) 1 new item, meaning that we will evaluate $O(k * \frac{\ln^{d-1} n}{(d-1)!})$ items. For each candidate, we need to iterate over the $|Q|$ queries and update the TopK structure with the new score, which takes $O(\log k)$ time using a Red-Black tree, for a total complexity of $O(|Q| * k * \log k * \frac{\ln^{d-1} n}{(d-1)!})$. In any case, as we talk about straightaway, this is a cynical investigation in light of the gullible supposition that every one of the k layers will be thought about totally. In practice, with the exception of the first layer, we only need to check a small fraction of the candidates in the skyline layers. For instance, in a uniform By and by, except for the main layer, we just need to check a little part of the applicants in the horizon layers. For example, in a uniform appropriation with successive layers of comparable size, the quantity of focuses to be considered will be in the request of k , since connections will be equitably disseminated among the horizon focuses.

As we just extend the best k things in each progression, roughly k new things will be assessed straightaway, making the cost of UPDATETOPK() in consequent calls $O(|Q| * k * \log k)$. Given that this cost is paid for each of the (at most) $k-1$ emphases after the first, the aggregate cost moves toward becoming $O(|Q| * (k^2 + \frac{\ln^{d-1} n}{(d-1)!}) * \log k)$. As we appear in our analyses, the real dispersions found in genuine datasets consider substantially speedier calculations. In the accompanying area, we depict a few speed-ups that can accomplish critical reserve funds by and by. Regarding space, the UPDATETOPK() strategy acknowledges $|\mathcal{X}|$ things as information and works on that set alone, bringing about $O(|\mathcal{X}|)$ space. For every thing in \mathcal{X} , we keep up its lower and upper bound, which is still $O(|\mathcal{X}|)$. As we emphasize over the questions, we refresh those qualities and dispose of things, diminishing the required space, conveying it closer to $O(k)$. Since the TopK structure dependably contains k sections, the space of CMiner is controlled by \mathcal{X} , which is at its most extreme when we recover the first horizon layer (line 7). Our supposition that the essential horizon fits in memory is sensible and shared by earlier takes a shot at horizon calculations [25].

IV. BOOSTING THE CMINER ALGORITHM

Next, we describe several improvements that we have applied to CMiner in order to achieve computational savings while maintaining the exact nature of the algorithm.

A. Query Ordering

Our unpredictability examination depends on the preface that CMiner assesses all questions Q for every competitor thing j . Notwithstanding, this supposition gullibly overlooks the calculation's pruning capacity, which depends on utilizing lower and upper limits on intensity scores to dispose of hopefuls early. Next, we demonstrate to significantly enhance the calculation's pruning viability by deliberately choosing the handling request of questions (line 23 of CMiner). CMiner utilizes the accompanying refresh rules for the lower and upper limits for an applicant j :

$$low(j) \leftarrow low(j) + p(q) \times V_{i,j}^q \quad (8)$$

$$up(j) \leftarrow up(j) - p(q) \times V_{i,i}^q + p(q) \times V_{i,j}^q \quad (9)$$

By expanding the sequences and using the initial values $low(j) = 0$ and $up(j) = C_{\mathcal{F}}(i, i)$, we can re-write the

$$bounds: low^m(j) = \sum_{l=1}^m p(q_l) \times V_{i,j}^{q_l} \\ up^m(j) = C_{\mathcal{F}}(i, i) - \sum_{l=1}^m p(q_l) \times V_{i,i}^{q_l} + \sum_{l=1}^m p(q_l) \times V_{i,j}^{q_l}, \text{ where } low^m(j)$$

and $up^m(j)$ are the values of the bounds after considering the m th query q_m . We can then define a recursive function $T(j) = up(j) - low(j)$ as follows:

$$T(j) \leftarrow T(j) - p(q) \times V_{i,i}^q \quad (9)$$

$T(j)$ catches the room for mistakes for the intensity between the thing of intrigue I and a competitor j . As more inquiries are assessed and the two limits are refreshed, the edge diminishes. At last, it ends up equivalent to zero when we have the last $C_{\mathcal{F}}(i, j)$ score. We estimate that the capacity to limit this edge speedier can expand the quantity of pruned competitors because of the presence of stricter limits in early emphases. Given Eq. 7 and 8, the estimation of $T(j)$ subsequent to considering m inquiries can be re-composed as takes after:

$$T^m(j) = C_{\mathcal{F}}(i, i) - \sum_{l=1}^m p(q_l) \times V_{i,i}^{q_l}, \quad (10)$$

where q_l is the l th question handled by the calculation. Given Eq. 10, unmistakably we can ideally limit the edge between the lower and upper limits on the intensity of a competitor by preparing questions in diminishing request of their $p(q) \times V_{i,i}^q$ qualities. We allude to this requesting plan as COV. We assess the computational investment funds accomplished by COV in Section 5.4 of our trials, where we likewise contrast it and elective methodologies.

B. Improving UPDATETOPK() and GETSLAVES()

In this segment we portray a few upgrades to the CMiner's two fundamental schedules. We execute these changes into an improved calculation, which we allude to as CMiner++. We incorporate this form in our trial assessment, where we contrast its productivity and that of CMiner, and in addition to that of different baselines. Despite the fact that CMiner can successfully prune low quality applicants, a noteworthy bottleneck inside the UPDATETOPK() work is the calculation of the last aggressiveness score between every hopeful and the thing of intrigue I (lines 41-46). Accelerating this calculation can tremendously affect the productivity of our calculation. Next, We illustrate this with

an example. Assume that items are defined in a 4-dimensional space with features f_1, f_2, f_3, f_4 . Without loss of generality, we assume that all features are numeric. We also consider 3 queries $q_1 = (f_1, f_2, f_3)$, $q_2 = (f_2, f_3, f_4)$ and $q_3 = (f_2, f_4)$ with probabilities $w(q_1)$, $w(q_2)$, and $w(q_3)$, respectively. In order to compute the competitiveness between two items i and j , we need to consider all queries and, according to Eq. 5, compute $V_{i,j}^{q_1} = V_{i,i}^{f_1} \times V_{i,j}^{f_2} \times V_{i,j}^{f_3}$, $V_{i,j}^{q_2} = V_{i,j}^{f_2} \times V_{i,j}^{f_3} \times V_{i,j}^{f_4}$, and $V_{i,j}^{q_3} = V_{i,j}^{f_2} \times V_{i,j}^{f_4}$. Given that the three items include common sequences of factors, we would like to avoid repeating their computation, when possible. First, we sort all features according to their frequency in the given set of queries. In our example, the order is: f_4, f_2, f_3, f_1 . In this order, (f_2, f_3) becomes a common prefix for q_1 and q_2 , whereas f_2 is a common prefix for all 3 queries.

We then build a prefix-tree to ensure that the computation of such common prefixes is only completed once. For instance, the computation of $V_{i,j}^{f_2} \times V_{i,j}^{f_3}$ is done only once and used for both q_1 and q_2 . The tree is utilized as a part of lines 41-46 of CMiner to assist the calculation of the intensity between the thing of intrigue and the rest of the competitors in X . This change is propelled by Huffman encoding, whereby visit images (includes for our situation) are nearer to the root, so they are encoded with less bits. Note that Huffman encoding is ideal if the images free of each other, similar to the case in our own setting. The GETSLAVES() strategy is utilized to broaden the arrangement of applicants by including the things that are overwhelmed by those in a gave set (lines 7 and 15). Thus forward, we allude to this as the dominator set. A gullible execution would incorporate all things that are ruled by no less than one thing in the dominator set. In any case, as expressed in Lemma 1, if a thing j is ruled by a thing j' , at that point the aggressiveness of j with anything of intrigue can't be higher than that of j' . This infers things that are commanded by the k -th best thing of the given set will have an aggressiveness score lower than the present k -th score and will in this manner not be incorporated into the last outcome. In this way, we just need to extend the best $k - 1$ things and just those that have not been extended as of now amid a past emphasis. What's more, the GETSLAVES() technique can be additionally enhanced by utilizing the lower bound LB (the score of the k -th best competitor) as takes after: instead of returning all the items that are dominated by those in the dominator set, we only have to consider a dominated item j if $C_{\mathcal{F}}(j, j) > LB$. This is due to the fact that the competitiveness between i and j is upper-bounded by the minimum coverage achieved by either of the two items (over all queries), i.e., $C_{\mathcal{F}}(i, j) \leq \min(C_{\mathcal{F}}(i, i), C_{\mathcal{F}}(j, j))$. Therefore, an item with a coverage $\leq LB$ cannot replace any of the items in the current TopK.

V. EXPERIMENTAL EVALUATION

In this segment we portray the trials that we led to assess our approach. All trials were finished on a work area with a Quad-Core 3.5 GHz Processor and 2GB RAM.

A. Datasets and Baselines

Our tests incorporate four datasets, which were gathered for the reasons for this undertaking. The datasets were purposefully chosen from various spaces to depict the cross-area appropriateness of our approach.

Notwithstanding the full data on every thing in our datasets, we likewise gathered the full arrangement of surveys that were accessible on the source site. These audits were utilized to (1) evaluate questions probabilities, as portrayed in Section 2.2 and (2) extricate the assessments of analysts on particular highlights. The very referred to strategy by Ding et al. is utilized to change over each audit to a vector of sentiments, where every assessment is characterized as an element extremity mix (e.g. service+, nourishment). The level of audits on a thing that express a positive sentiment on a particular element is utilized as the element's numeric incentive for that thing. We allude to these as conclusion highlights. Incorporates graphic measurements for each dataset, while a point by point depiction is given underneath.

CAMERAS: This dataset incorporates 579 computerized cameras from Amazon.com. We gathered the full arrangement of surveys for every camera, for an aggregate of 147192 audits. The arrangement of highlights incorporates the determination (in MP), shade speed (in a flash), zoom (e.g. 4x), and cost. It likewise incorporates assessment includes on manual, photographs, video, plan, streak, center, menu alternatives, lcd screen, measure, highlights, focal point, guarantee, hues, adjustment, battery life, determination, and cost.

HOTELS: This informational collection incorporates 80799 survey child 1283 lodgings from Booking.com. The arrangement of highlights incorporates the offices, exercises, and administrations offered by the lodging. Every one of the three of these multi-straight out highlights are accessible on the site. The dataset additionally incorporates feeling highlights on area, administrations, neatness, staff, and solace.

RESTAURANTS: This informational collection incorporates 30821 audits on 4622 New York City eateries from TripAdvisor.com. The arrangement of highlights for this dataset incorporates the food writes and dinner composes (e.g. lunch, supper) offered by the eatery, and in addition the action composes (e.g. drinks, parties) that it is useful for. Each of the three of these multi-straight out highlights are accessible on the site. The dataset additionally incorporates feeling highlights on nourishment, benefit, esteem for-cash, climate, and cost.

RECIPES: This dataset incorporates 100000 formulas from Sparkrecipes.com. It likewise incorporates the full arrangement of audits on every formula, for a sum of 21685 surveys. The arrangement of highlights for every formula incorporates the quantity of calories, and the accompanying wholesome data, estimated in grams: fat, cholesterol, sodium, potassium, carb, fiber, protein, vitamin A, vitamin B12, vitamin C, vitamin E, calcium, copper, folate, magnesium, niasin, phosphorus, riboflavin, selenium, thiamin, zinc. All data is straightforwardly accessible on the site.

B. Evaluating Comparative Methods

Past work on contender mining has been founded on similar proof between two things, found in various kinds of content information. In any case, these methodologies

depend on the supposition that such similar confirmation can be found in plenitude in the accessible information. In this trial, we assess this supposition on our four datasets. For each match of things in each dataset, we report (1) the quantity of audits that specify the two things and (2) the quantity of surveys that incorporate an immediate correlation between the two things. We concentrate such near confirmation in view of the association of "aggressive proof" dictionaries utilized by past work. Given two things I and j, the dictionary incorporates the accompanying near examples: I as opposed to j, I dissimilar to j, I contrasted and j, contrast I with j, I beat(s) j, I surpasses j, I outperform(s) j, lean toward I to j, I than j, I same as j, I like j, I better than j, I superior to j, I more terrible than, I more than j, I not as much as j, I versus j. We display the outcomes in Table 4, in which we report the normal number of discoveries for each match of things in each dataset.

TABLE IV: Evidence on Comparative Methods

	Co-occurrence	Comparative
Cameras	1.7	1.2
Hotels	0.06	0.02
Restaurants	0.09	0.04
Recipes	0	0

The outcomes check that strategies in view of near confirmation are totally ineffectual in numerous spaces. Truth be told, notwithstanding for CAMERAS, the dataset with the biggest tally, prove was constrained to few sets. In particular, the normal number of times that any two particular cameras seem together in a similar survey is 1.7. What's more, just 1.2 of these co-event were really relative, a number that is dreadfully low to take into consideration a sure assessment of intensity. This exhibits the meager condition of near confirmation in genuine information, which incredibly restricts the appropriateness of any approach that depends on such proof. These discoveries additionally persuade our work, which has no requirement for this kind of data.

C. Computational Time

In this investigation we contrast the speed of CMiner and that of the two baselines (Naive and GMiner), and additionally with that of the improved CMiner++ calculation. Specifically, we utilize every calculation to figure the arrangement of best k contenders for every thing in our datasets. Initially, the Naive calculation reliably reports the same computational time paying little mind to k, since it innocently registers the intensity of each and every thing in the corpus regarding the objective thing. Consequently, any insignificant varieties in the required time are because of the way toward keeping up the best k set. All in all, Naive is beaten by the two different calculations, and is aggressive for substantial estimations of k for the HOTELS dataset. The last case can be credited to the modest number of inquiries and things incorporated into this dataset, which confine the capacity of more advanced calculations to significantly prune the space when the quantity of required contenders is vast. For the CAMERAS dataset, CMiner and GMiner, display relatively indistinguishable running circumstances. This is because of (1) the simple huge number of unmistakable inquiries for

this dataset (14779), which fills in as a computational bottleneck for CMiner and (2) fills the profoundly grouped structure of the thing populace, which incorporates 579 things. A more profound investigation uncovers that GMiner distinguishes and normal of 443.63 thing gatherings (i.e. gatherings of indistinguishable things) per inquiry. This implies the calculation spares (on desire) a sum of $(579 - 443) \times 14779 = 2009944$ scope calculations for every question, enabling it to be focused to the generally unrivaled CMiner. Indeed, for alternate datasets, CMiner shows a reasonable preferred standpoint. This preferred standpoint is amplified for the RECIPES dataset, which is the most crowded of the four as far as included things. The test on this dataset additionally outlines the versatility of the approach as for k. For the HOTELS and RESTAURANTS datasets, despite the fact that the computational time of CMiner seems to ascend as k increments for the other three datasets, it never goes over 0.035 seconds. For the CAMERAS dataset, the extensive number of considered inquiries has an antagonistic of the adaptability of CMiner, since it brings about bigger number of required calculations for bigger estimations of k. This finding rouses us to consider pruning the arrangement of inquiries by dispensing with those that have a low likelihood.

We investigate this bearing in the examination introduced in Section 5.6. At last, we watch that the upgraded CMiner++ calculation reliably beat the various methodologies, crosswise over datasets and estimations of k . The upside of CMiner++ is expanded for bigger estimations of k , which enable the calculation to profit by its enhanced pruning. This confirms the utility of the upgrades depicted in Section 4.2 and exhibits that viable pruning can prompt an execution that far surpasses the most pessimistic scenario unpredictability examination of CMiner.

VI. RESULTS

Results of this paper is as shown in bellow Figs.2 to 10.

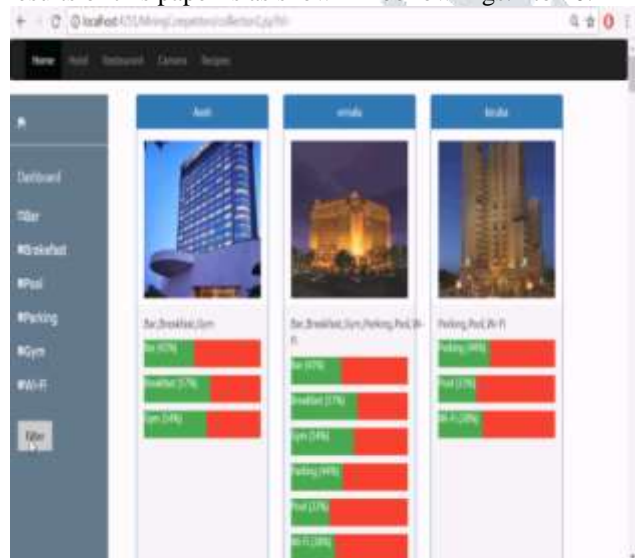


Fig.2. Hotels.

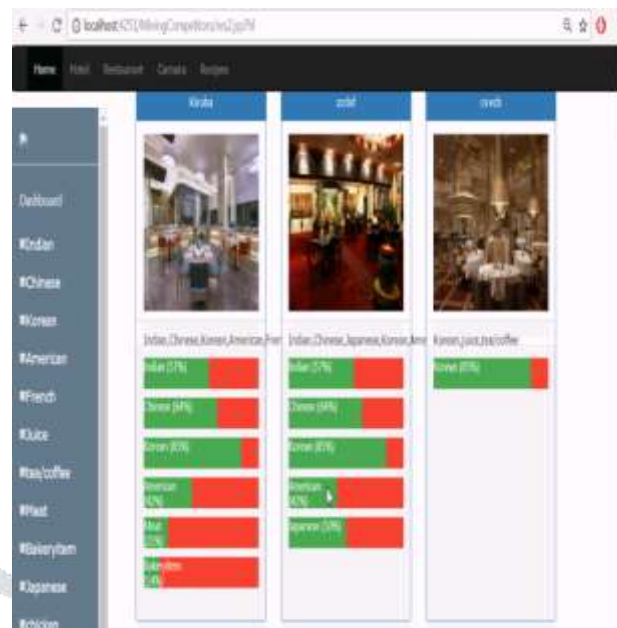


Fig.3. Restaurants.

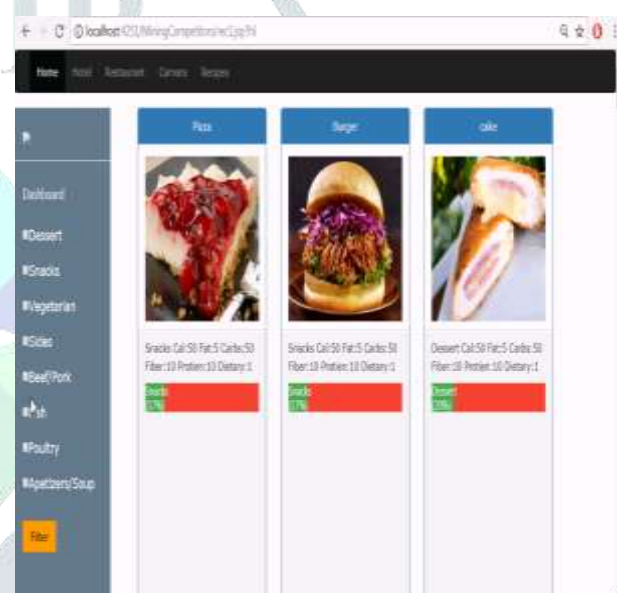


Fig.4. Recipes.

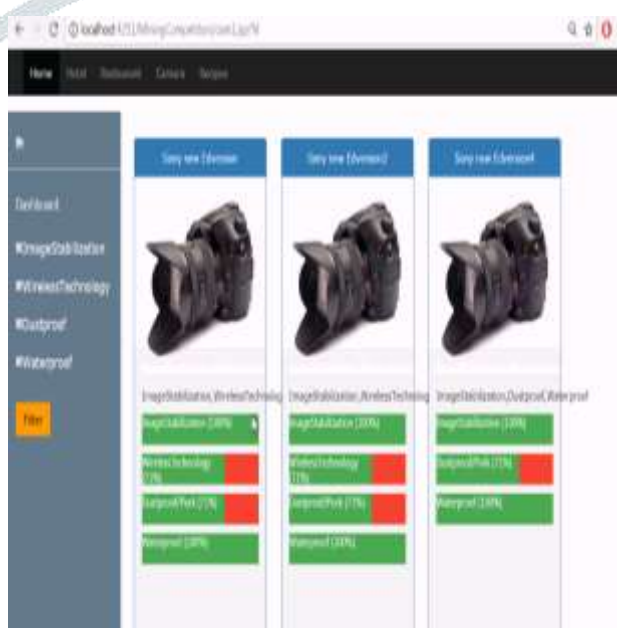


Fig.6. Cameras.

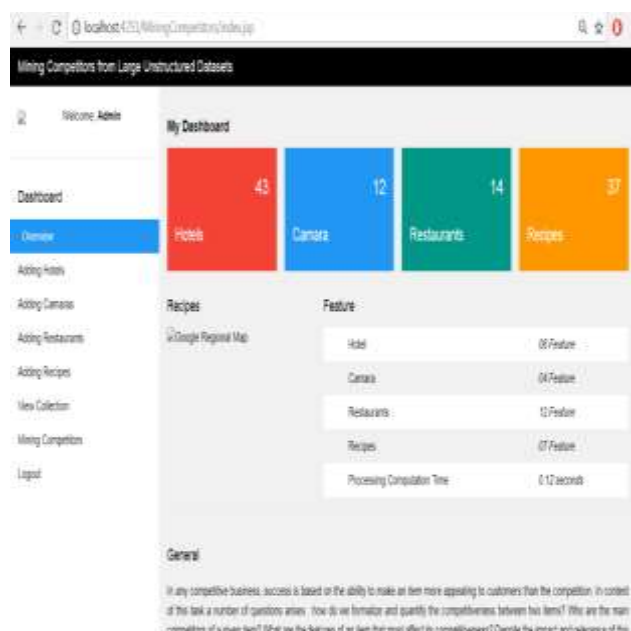


Fig.7. Overview of hotels, restaurants, recipes and cameras.

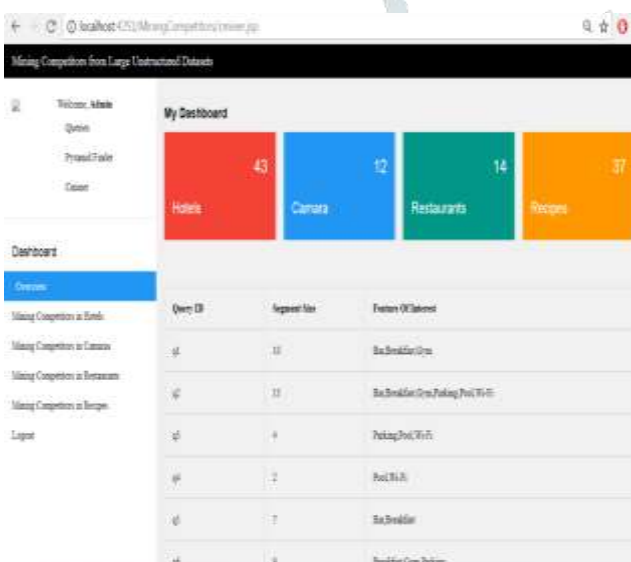


Fig.8. Queries.

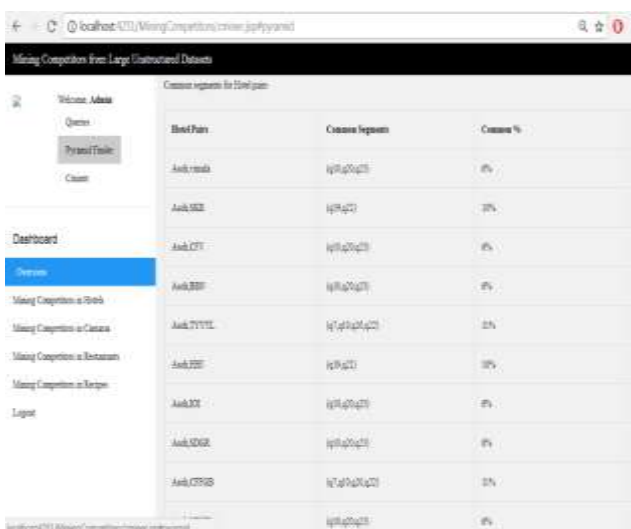


Fig.9. Pyramid finder-common segments for hotel pairs.

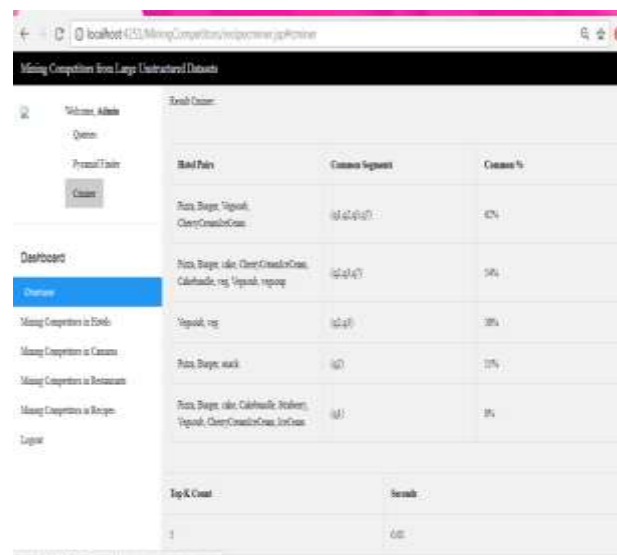


Fig.10. CMiner-top k-count.

VII. CONCLUSION

A formal criticalness of forcefulness between both things, affirmed both quantitatively and emotionally. Formalization is appropriate transversely over spaces, crushing the insufficiencies of past approaches. Consider distinctive fragments that have been by and large slighted previously, for eg, the circumstance of the matters in the multi-dimensional component space furthermore, the tendencies and finishes customers. Work familiarizes an end-with end approach for mining's such knowledge from huge datasets of clients reviews. In perspective of definition, Watched out for the computationally troublesome troubles of finding the best k contenders of a given factor. The proposed structure is capable and appropriated to spaces with broad masses of things. The capabilities of methodology was checked by strategies for a tests appraisal on authentic datasets from different spaces. Examinations additionally uncovered that solitary no of reviews is satisfactory to positively get to the specific sorts of clients in a given market, too the measure of clients that had put with each kind. In future, Formalization is pertinent transversely finished spaces, beating the inadequacies of past methodologies. Consider various variables that had been to a great extent neglected previously, for eg, the circumstance of things in multi-dimensional component space and the tendencies and supposition of the customers.

VIII. REFERENCES

- [1]GEORGE VALKANAS ; THEODOROS LAPPAS ; DIMITRIOS GUNOPILOS (SEPT. 1 2017).MINING COMPETITORS FROM LARGE UNSTRUCTURED DATASETS. VOLUME: 29,IEEE COMPUTER SOCIETY
- [2]M.E.Porter,CompetitiveStrategy:TechniquesforAnalyzing Industries and Competitors. Free Press, 1980.
- [3] R. Deshpande and H. Gatingon, "Competitive analysis," Marketing Letters, 1994.
- [4] B. H. Clark and D. B. Montgomery, "Managerial Identification of Competitors," Journal of Marketing, 1999.
- [5] W. T. Few, "Managerial competitor identification: Integrating the categorization, economic and organizational identity perspectives," Doctoral Dissertation, 2007.
- [6] M. Bergen and M. A. Peteraf, "Competitor identification and competitor analysis: a broad-based

managerial approach,” Managerial and Decision Economics, 2002.

[7] J. F. Porac and H. Thomas, “Taxonomic mental models in competitor definition,” The Academy of Management Review, 2008.

[8] M.-J. Chen, “Competitor analysis and inter firm rivalry: Toward a theoretical integration,” Academy of Management Review, 1996.

[9] R. Li, S. Bao, J. Wang, Y. Yu, and Y. Cao, “Cominer: An effective algorithm for mining competitors from the web,” in ICDM, 2006.

[10] Z. Ma, G. Pant, and O. R. L. Sheng, “Mining competitor relationships from online news: A network-based approach,” Electronic Commerce Research and Applications, 2011.

[11] R. Li, S. Bao, J. Wang, Y. Liu, and Y. Yu, “Web scale competitor discovery using mutual information,” in ADMA, 2006.

[12] S. Bao, R. Li, Y. Yu, and Y. Cao, “Competitor mining with the web,” IEEE Trans. Knowl. Data Eng., 2008.

[13] G. Pant and O. R. L. Sheng, “Avoiding the blind spots: Competitor identification using web text and linkage structure,” in ICIS, 2009.

[14] D. Zelenko and O. Semin, “Automatic competitor identification from public information sources,” International Journal of Computational Intelligence and Applications, 2002.

[15] R. Decker and M. Trusov, “Estimating aggregate consumer preferences from online product reviews,” International Journal of Research in Marketing, vol. 27, no. 4, pp. 293–307, 2010.

[16] C. W.-K. Leung, S. C.-F. Chan, F.-L. Chung, and G. Ngai, “A probabilistic rating inference framework for mining user preferences from reviews,” World Wide Web, vol. 14, no. 2, pp. 187–215, 2011.

[17] K. Lerman, S. Blair-Goldensohn, and R. McDonald, “Sentiment summarization: evaluating and learning user preferences,” in ACL, 2009, pp. 514–522.

[18] E. Marrese-Taylor, J. D. Velásquez, F. Bravo-Marquez, and Y. Matsuo, “Identifying customer preferences about tourism products using an aspect-based opinion mining approach,” Procedia Computer Science, vol. 22, pp. 182–191, 2013.

[19] C.-T. Ho, R. Agrawal, N. Megiddo, and R. Srikant, “Range queries in olap data cubes,” in SIGMOD, 1997, pp. 73–88.

[20] Y. L. Wu, D. Agrawal, and A. El Abbadi, “Using wavelet decomposition to support progressive and approximate range-sum queries over data cubes,” in CIKM, ser. CIKM ’00, 2000, pp. 414–421.

[21] D. Gunopulos, G. Kollios, V. J. Tsotras, and C. Domeniconi, “Approximating multi-dimensional aggregate range queries over real attributes,” in SIGMOD, 2000, pp. 463–474.

[22] M. Muralikrishna and D. J. DeWitt, “Equi-depth histograms for estimating selectivity factors for multi-dimensional queries,” in SIGMOD, 1988, pp. 28–36.

[23] N. Thaper, S. Guha, P. Indyk, and N. Koudas, “Dynamic multidimensional histograms,” in SIGMOD, 2002, pp. 428–439.

[24] K.-H. Lee, Y.-J. Lee, H. Choi, Y. D. Chung, and B. Moon, “Parallel data processing with mapreduce: a survey,” AcM SIGMOD Record, vol. 40, no. 4, pp. 11–20, 2012.

[25] S. Borsányi, D. Kossmann, and K. Stocker, “The skyline operator,” in ICDE, 2001.

[26] D. Papadias, Y. Tao, G. Fu, and B. Seeger, “An optimal and progressive algorithm for skyline queries,” ser. SIGMOD ’03.

Author’s Details:

Ms. SYEDA FATIHA BUTTUL has completed her project from Shadan Women’s College Of Engineering And Technology, Khairthabad, JNTUH University Hyderabad. Presently, she is pursuing her Masters in Computer Science and Engineering (C.S.E) from Shadan Women’s College Of Engineering And Technology, Khairthabad, Hyderabad, TS. India.

Ms. K. SHILPA has completed B.E (I.T) from M.V.S.R Engineering college, Osmania University, Hyderabad. M.tech (C.S.E) from Aurora Technology and Research Institute, JNTU University, Hyderabad. Currently working as an Assistant Professor of IT department in Shadan Women’s College Of Engineering And Technology, Hyderabad, TS. India.