

# Effective Evidence Aggregation to Identify the Ranking Fraud for Mobile APPS Using KNN

K. Sujitha<sup>1</sup>, S. Radha Krishna<sup>2</sup>

M. Tech Scholar, Department of CSE, University College of Engineering Vizianagaram

Assistant Professor, Department of CSE, University College of Engineering Vizianagaram

**Abstract:** *Nowadays everyone having smart phone. We have to install so many applications in that smart phone. Presently, the advancements made in the mobile technology are the production of mobile apps. Due to the more number of mobile apps, there is a chance of fraud mobile applications is of greater existence. Ranking fraud is the key challenges persist in the mobile app market. It defines the report over the apps in order to place them in leader board. In this study, we consider the fraud activities in the mobile apps, as a significant one. We propose ranking scheme for recognizing the fraudulent apps. While the essentials of envisioning situating deception have been for the most part seen, there is obliged understanding and research here. To this end, in this paper, we dedicate an all-encompassing view of positioning extortion and propose a positioning misrepresentation location framework for versatile Apps, We first propose to just determine the positioning misrepresentation by mining the dynamic time frames, to be specific driving sessions, of versatile Apps. Such driving sessions can be used for identifying the neighbourhood irregularity rather than worldwide inconsistency of App rankings. All, we examine three sorts of copies, i.e., positioning based confirmations, evaluation based proofs and audit based proofs, by displaying Apps' positioning, rating and survey practices through measurable theory tests. What's more, we offer an enhancement based conglomeration strategy to incorporate every one of the confirmations for extortion recognition. Specifically, it is proposed to exactly discover the mining to posture blackmail the dynamic time frames, to be particular driving sessions, of compact Apps. The KNN algorithm is applied to enhance effectiveness and precision of the application, we approve the sufficiency of the proposed framework and demonstrate the versatility of the identification, calculation and some normality of positioning misrepresentation exercises.*

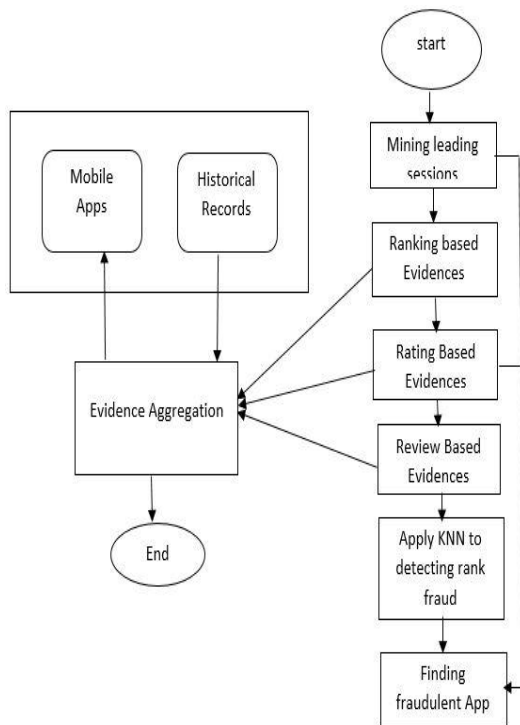
**Keywords:** Mobile Apps, Evidence Aggregation, Positioning misrepresentation, Classification, KNN

## 1. Introduction

Examining the data and obtaining out significant part out of it is really difficult and is the most important need. There is a huge amount of data available in the Information Industry. This data is of no use until it is changed into utile information. It is necessary to analyse this huge amount of data and extract useful information from it. Data mining is interdisciplinary subfields of computer science that can help us meet this need by providing tools to find out the important part or we can say knowledge from data. The reason for positioning misrepresentation in the mobile application market is for enrolling the app in the popularity list. While the significance of anticipating positioning misrepresentation has been generally perceived, there is restricted comprehension and exploration around there. In fact, the App leaderboard is a standout amongst the most essential routes for advancing portable Apps. A top most position on the leaderboard more popular is the app is the fact. The related works of this study classified into three categories. web positioning spam identification [4], [5], [6], online survey spam identification [7], [8], [9], and portable App proposal [10], [11], [12], [13], the issue of recognizing positioning misrepresentation for portable Apps is still under-investigated.

Specifically, the Web ranking spam mentions to any careful actions which bring to selected Web pages an unjustifiable favorable relevance or importance. In this, the problem of unsupervised web spam detection is studied. The second category is focused on detecting online review spam. Here we identified several representative behaviors of review spammers and model these behaviors to detect the

spammers. Here we aims to detect users generating spam reviews or review spammers. They identify several characteristic behaviors of review spammers and model these behaviors so as to detect the spammers. Finally, the third category includes the studies on mobile App recommendation, which is based on user's App usage records to build a preference matrix instead of using explicit user ratings. But in this paper, we aim to construct up a positioning misrepresentation identification system for portable Apps. Such test can be regarded as identifying the nearby irregularity rather than worldwide peculiarity of portable Applications. Second, because of the enormous number of versatile Apps, it is hard to physically mark positioning extortion for each App, so it is imperative to sustain a scalable way to consequently recognize positioning extortion without utilizing any benchmark information. At final, because of the dynamic way of outline rankings, it is difficult to recognise and affirm the proofs connected to positioning extortion, which goads us to find a few verifiable misrepresentation examples of versatile Apps as proofs.



**Figure 1:** The framework in ranking fraud detection system for mobile Apps

In particular, we first propose a straightforward yet compelling calculation to identify the main sessions of each App in light of its chronicled positioning records. At that stage, with the investigation of Apps' positioning practices, we see that the fake Apps frequently have diverse positioning examples in every driving session thought about with ordinary Apps. Thus, we describe some misrepresentation confirmations from Apps' authentic positioning records, and create three capacities to concentrate such positioning based misrepresentation confirmations. No matter, the positioning based confirmations can be influenced by App designers' notoriety and some true blue advertising effort, for instance, "constrained time markdown". As an event, it is not enough to just utilize positioning based proofs.

Along these lines, we encourage propose two sorts of extortion confirmations in view of Apps' evaluating and audit history, which mirror some inconsistency designs from Apps' authentic rating what's more, audit records. What's more, we build up an unsupervised prove collection strategy to coordinate these three sorts of proofs for assessing the validity of driving sessions from portable Apps. Figure 1 demonstrates the system of our positioning misrepresentation location framework for portable Apps. It is significant that every one of the confirmations are separated by demonstrating Apps' positioning, rating and audit practices through measurable theories tests.

## 2. Describing the leading sessions for mobile apps

From [1] First, ranking fraud does not always occur in the whole life cycle of an App, so we need to find the time when fraud happens. Such challenge can be considered as finding the local anomaly instead of a global anomaly of mobile

Apps. Second, due to the vast number of mobile Apps, it is difficult to manually label ranking fraud for each App, so it is significant to sustain a scalable path to automatically detect ranking fraud without using any benchmark data. Lastly, due to the dynamic nature of chart rankings, it is not easy to identify and confirm the evidence linked to ranking fraud, which prompts us to discover some implicit fraud patterns of mobile Apps as evidence. Our deliberate observation uncovers that mobile Apps are not constantly ranked high in the leader board, but merely in some leading events, which form different leading sessions. In another language, ranking fraud usually comes about in these leading sessions. Thus, detecting ranking fraud of mobile Apps is actually to detect ranking fraud within leading sessions on mobile Apps. Specifically, Observe that we use a ranking threshold  $K^*$  which is usually smaller than  $K$  here because  $K$  maybe very heavy (e.g., more than 1000), and the ranking records beyond  $K^*$  (e.g., 300) are not really useful for detecting the ranking manipulations. Furthermore, we likewise discover that some Apps have several adjacent leading events which are near to each other and build a leading session

### Mining leading sessions

There are two main steps for mining primary sessions. Formost, we need to discover leading events from the App's Historical ranking records. Second we need to merge contiguous leading events for building leading sessions.

#### Pseudo code of mining leading sessions for a given App

```

1:    $E_s = \phi; e = \phi; t_{start}^e = 0;$ 
2:   for each  $i \in [1, |R_a|]$  do
3:     if  $r_i^a \leq K^*$  and  $t_{start}^e = 0$  then
4:        $t_{start}^e = t_i;$ 
5:     else if  $r_i^a > K^*$  and  $t_{start}^e \neq 0$  then
6:       //found one event;
7:        $t_{end}^e \neq t_{i-1}; e = \langle t_{start}^e, t_{end}^e \rangle;$ 
8:       if  $E_s = \phi$  then
9:          $E_s \cup = e; t_{start}^s = t_{start}^e; t_{end}^s = t_{end}^e;$ 
10:      else if  $(t_{start}^e - t_{end}^s) < \phi$  then
11:         $E_s \cup = e; t_{end}^s = t_{end}^e;$ 
12:      else then
13:        // found one session;
14:         $s = \langle t_{start}^s, t_{end}^s, E_s \rangle;$ 
15:         $S_a \cup = s; s = \phi$  is a new session;
16:         $E_s = \{e\}; t_{start}^e = t_{start}^e; t_{end}^e = t_{end}^e;$ 
17:         $t_{start}^e = 0; e = \phi$  is a new leading event;
18:      return  $S_a$ 
  
```

We denote each leading event  $e$  and session  $s$  as tuples  $\langle t_{start}^e, t_{end}^e \rangle$  and  $\langle t_{start}^s, t_{end}^s, E_s \rangle$  respectively, where  $E_s$  is the set of leading events in session  $s$ . specifically, we first excerpt individual leading event  $e$  for the given app  $a$ . from step 2 to step 7 indicates extraction of individual leading event from the beginning time. For each evoked individual leading event  $e$ , we determine the time span between  $e$  and the current leading session  $s$  to determine whether they belong to the same leading session or not. Particularly, if  $(t_{start}^e - t_{end}^s) < \phi$   $e$  will be deliberated as a new leading session from Step 8 to step 16 indicates how to construct

leading session. Thus this algorithm can identify leading events and sessions by scanning a's historical records only once.

### 3. Describing the evidences for ranking fraud detection

From [1] Identifying different evidences for ranking fraud detection is apply on output of mining leading session algorithm. Step by step apply three evidences applied are ranking based, rating based and review based. We build up an unsupervised proof conglomeration technique to coordinate these three sorts of confirmations for assessing the believability of driving sessions from portable Apps.

#### 3.1 Ranking based evidences

The positioning based confirmations are valuable for positioning misrepresentation discovery By examining the Apps' chronicled positioning records, we watch that Apps' positioning practices in a main occasion dependably fulfil a particular positioning example, which comprises of three distinctive positioning stages, inparticular, rising stage, keeping up stage and subsidence stage. In particular, in every driving occasion, an App's positioning first increments to a pinnacle position in the pioneer board (i.e., rising stage), then keeps such pinnacle position for a period (i.e., looking after stage), lastly diminishes till the end of the occasion (i.e., subsidence stage).

We propose some positioning based marks of driving sessions to build extortion confirmations for positioning misrepresentation recognition.

##### Evidence: 1:

We use two shape parameters  $\theta_1$  and  $\theta_2$  to quantify the ranking figures of therising phase and the recession phase of App a's leading event e, which can be computed by

$$\theta_1^e = \arctan\left(\frac{r^* - r_{\bar{e}}^e}{t_{\bar{e}}^e - t_{\bar{e}}^e}\right), \theta_2^e = \arctan\left(\frac{r^* - r_{\bar{e}}^e}{t_{\bar{e}}^e - t_{\bar{e}}^e}\right)$$

Naturally, a broad  $\theta_1$  may exhibit that the App has been thump to a high rank inside a brief traverse, and a significant  $\theta_2$  may demonstrate that the App has dropped from a high rank to the base inside a brief time span. Along these lines, a primary session, which has furthermore driving events with broad  $\theta_1$  and  $\theta_2$  values, has higher probability of having situating deception. Here, we define a fraud signature us for a leading session as follows:

$$\bar{\theta} = \frac{1}{|E_s|} \sum_{e \in E_s} (\theta_1^e + \theta_2^e)$$

where  $|E_s|$  is the number of leading events in session s. we can calculate the p-value by

$$(\mathcal{N}(\mu_{\bar{\theta}}, \sigma_{\bar{\theta}}) \geq \bar{\theta}_s) = 1 - \frac{1}{2} (1 + \operatorname{erf}\left(\frac{\bar{\theta}_s - \mu_{\bar{\theta}}}{\sigma_{\bar{\theta}} \sqrt{2}}\right)),$$

Where erf(x) is the Gaussian Error Function as follows

$$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$$

Intuitively, a guiding session with a smaller p-value P has more prospect to reject Hypothesis 0 and accept Hypothesis

1. This means it has more chance of committing ranking fraud. Thus, we define the evidence as

$$\Psi_1(s) = 1 - \mathbb{P}(\mathcal{N}(\mu_{\bar{\theta}}, \sigma_{\bar{\theta}}) \geq \bar{\theta}_s)$$

##### Evidence2

The Apps with ranking fraud often experience a short maintaining phase with high ranking posts in each running event. Thus, if we denote the maintaining phase of a leading event e as  $\Delta t_m^e = (t_c^e - t_b^e + 1)$ , and the average rank in this maintaining phase as  $\bar{r}_m^e$  we can define a fraud signature  $\chi_s$  for each leading session as follows:

$$\chi_s = \frac{1}{|E_s|} \sum_{e \in E_s} \frac{r^* - \bar{r}_m^e}{\Delta t_m^e}$$

We assume  $\chi_s$  follows the Gaussian distribution,  $\chi_s \sim \mathcal{N}(\mu_{\chi}, \sigma_{\chi})$  where  $\mu_{\chi}$  and  $\sigma_{\chi}$  can be discovered by the MLE method from the observations of  $\chi_s$  in all Apps' classical leading sessions. Then, we can calculate the evidence by

$$\Psi_2(s) = 1 - \mathbb{P}(\mathcal{N}(\mu_{\chi}, \sigma_{\chi}) \geq \chi_s)$$

##### Evidence 3:

The act of leading events in a leading session, i.e.,  $|E_s|$  is also a strong signature of ranking fraud. Since  $|E_s|$  always has discrete values, we suggest to leverage the Poisson approximation to calculate the p-value with the above hypotheses. Specifically, we assume  $|E_s|$  follows the Poisson distribution,  $|E_s| \sim \mathcal{P}(\lambda_s)$  where the criterion  $\lambda_s$  can be learnt by the MLE method from the observations of  $|E_s|$  in all Apps' historical leading sessions. Then, we can calculate the p-value as follows

$$\mathbb{P}(\mathcal{P}(\lambda_s) \geq |E_s|) = 1 - e^{-\lambda_s} \sum_{i=0}^{|E_s|} \frac{(\lambda_s)^i}{i!}$$

Therefore we can compute the evidence by

$$\Psi_3(s) = 1 - \mathbb{P}(\mathcal{P}(\lambda_s) \geq |E_s|)$$

#### 3.2 Rating based evidences

The ranking based evidences are useful for ranking fraud detection. Yet, sometimes, it is not sufficient to only use ranking based evidences. Specifically, after an App has been printed, it can be ordered by any user who downloaded it. Indeed, user rating is one of the most significant features of App advertisement. An App which has higher rating may downloaded by many users and can also be ranked higher in the Leaderboard. Therefore, rating manipulation is also an important perspective on ranking fraud.

##### Evidence 4

This evidence is based on rating levels given by the user., each rating can be classified into  $|L|$  discrete rating levels, e.g., 1 to 5, here rating level denotes the user preferences of the app, Based on that Then, we use the Cosine similarity between  $P(l_i | \mathcal{R}_{s,a})$  and  $P(l_i | \mathcal{R}_a)$  to estimate the change as follows:

$$D(s) = \frac{\sum_{i=1}^{|L|} P(l_i | \mathcal{R}_{s,a}) \times P(l_i | \mathcal{R}_a)}{\sqrt{\sum_{i=1}^{|L|} P(l_i | \mathcal{R}_{s,a})^2} \times \sqrt{\sum_{i=1}^{|L|} P(l_i | \mathcal{R}_a)^2}}$$

Therefore, if  $D(s)$  value of leading session has significantly lower compared with other leading sessions of Apps in the leaderboard, it has high probability of having ranking fraud. We use the Gaussian approximation to compute the p-value with the above hypotheses. Specifically, we assume  $D(s)$  follows the Gaussian distribution,  $D(s) \sim \mathcal{N}(\mu_D, \sigma_D)$  where  $\mu_D$  and  $\sigma_D$  can be learnt by the MLE method from the observations of  $D(s)$  in all Apps' historical leading sessions. Hence, we can calculate the evidence by

$$\Psi_4(s) = 1 - \mathbb{P}(\mathcal{N}(\mu_D, \sigma_D) \leq D(s)).$$

**Evidence 5:**

The average rating in a specific leading session should be coherent with the fair value of all historical military ranks. We define a fraud signature  $\Delta R_s$  for each leading session as follows:

$$\Delta R_s = \frac{\bar{R}_s - \bar{R}_a}{\bar{R}_a}, (s \in a)$$

Where  $\bar{R}_s$  is the average rating in leading session  $s$ , and  $\bar{R}_a$  is the average historical rating of App  $a$ . we can compute the evidence by

$$\Psi_5(s) = 1 - \mathbb{P}(\mathcal{N}(\mu_R, \sigma_R) \geq \Delta R_s),$$

**3.3 Review based evidences**

Easily review manipulation is one of the most significant perspective of App positioning misrepresentation. Most of the App stores also allow users to write some textual note as App reviews. Such reviews can reflect the personal perceptions and usage experiences of existing users for particular mobile Apps. Imposters often post mock reviews in the leading sessions of a specific App in order to inflate the App downloads, and so propel the App's ranking position in the leaderboard. Here we propose two fraud evidences based on Apps' review behaviours in leading a session for detecting ranking fraud.

**Evidence 6:**

We propose to leverage topic modelling to excerpt the latent topics of reviews. Specifically, here we take over the widely used Latent Dirichlet Allocation (LDA) model for learning latent semantic topics. If we denote the reviews in leading sessions of a as  $C_{s,a}$  we can use the KL-divergence to estimate the change of topic distributions between  $C_a$  and  $C_{s,a}$

$$D_{KL}(s||a) = \sum_k P(z_k | C_{s,a}) \ln \frac{P(z_k | C_{s,a})}{P(z_k | C_a)}$$

Where  $P(z_k | C_a)$  and  $P(z_k | C_{s,a}) \propto P(z_k) \prod_{w \in C_{s,a}} P(w | z_k)$  obtained through the LDA training process. Then, we can compute the evidence by

$$\Psi_6(s) = 1 - \mathbb{P}(\mathcal{N}(\mu_{KL}, \sigma_{KL}) \geq D_{KL}(s||a))$$

**Evidence 7**

For most, for each review  $c$  in leading session  $s$ , we remove all stop words (e.g., "of", "the") and normalize verbs and adjectives (e.g., "plays  $\rightarrow$  play", "better  $\rightarrow$  good"). Second, we build a normalized words vector  $\vec{w}_c = \text{dim}[n]$  for each review  $c$ , where  $n$  indicates the number of all unique normalized words in all reviews of  $s$ . To be specific, here we

have  $\text{dim}[i] = \frac{\text{freq}_{i,c}}{\sum_i \text{freq}_{i,c}} (1 \leq i \leq n)$  where  $\text{freq}_{i,c}$  is the frequency of the  $i$ th word in  $c$ . Finally, we can calculate the similarity between two reviews  $c_i$  and  $c_j$  by the Cosine similarity  $\text{Cos}(\vec{w}_{c_i}, \vec{w}_{c_j})$ . Thus, the fraud signature  $\text{Sim}(s)$  can be computed by

$$\text{Sim}(s) = \frac{2 \times \sum_{1 \leq i < j \leq N_s} \text{Cos}(\vec{w}_{c_i}, \vec{w}_{c_j})}{N_s \times (N_s - 1)},$$

Where  $N_s$  is the number of reviews during leading sessions. Possibly, the higher value of  $\text{sim}(s)$  indicates more same/near-duplicate reviews in  $s$ . Here, we use the Gaussian approximation to compute the p-value with the above hypotheses. Hence, we can calculate the evidence by

$$\Psi_7(s) = 1 - \mathbb{P}(\mathcal{N}(\mu_{sim}, \sigma_{sim}) \geq \text{Sim}(s))$$

**Evidence Aggregation:**

Here we calculate the final evidence  $\Psi^*(s)$  as a linear combination of all the evidences we calculated earlier

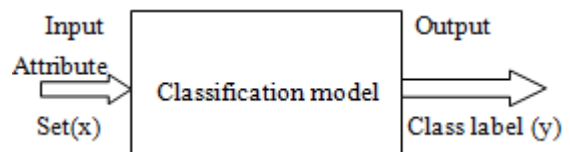
$$\Psi^*(s) = \sum_{i=1}^N w_i \times \Psi_i(s)$$

**4. Classification**

Classification techniques in information mining are capable of swearing of heavy quality of information. It can be applied to predict categorical class labels and classifies data based on the training set and class labels and it can be applied for classifying newly available information. The term could cover any context in which some decision or forecast is prepared on the basis of currently available data. Classification procedure is recognized method for repeatedly making such decisions in new places. Here if we assume that problem is a concern with the construction of a process that will be given to a continuing sequence of cases in which each new instance must be allotted to one of a set of predefined classes on the footing of observed characteristics of information. Founding of a classification procedure of a set of information for which the exact classes are known in advance is termed as pattern recognition or supervised learning. Some of the most critical problems arising in science, industry, and commercialism can be called as classification or decision problems. All groups have some aims in common

**4.1 Classification Algorithm**

Categorization is one of the Data Mining techniques that is mainly used to study a given information set and takes each instance of it and assigns this instance to a particular course of study such that classification error will be misplaced. It is applied to extract models that accurately define important data classes within the given data set.



**Figure 2:** Classification as the task of mapping an input attributes set  $x$  into its class label  $y$ .

**Definition of Classification**

Compartmentalization is the task of leaning a target function  $f$  that maps each attribute set  $x$  to one of the predefined classes  $y$ . The objective function is usually known as classification model. Classification model is useful for the following function.

**Descriptive modelling** is a classification model can function as an explanatory tool to distinguish between objects of different divisions.

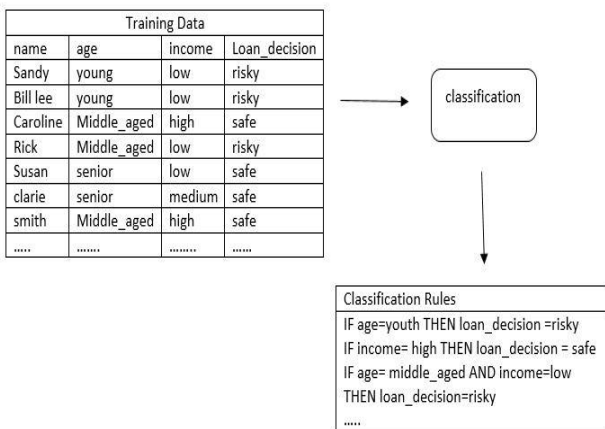
**Predictive modelling** is a classification model can likewise be applied to predict the category label of unknown records. Classification techniques are almost suited for predicting or describing datasets with binary or nominal categories.

**The Data Classification process includes two steps:**

- Constructing the Classifier or Model
- Utilizing Classifier for Classification

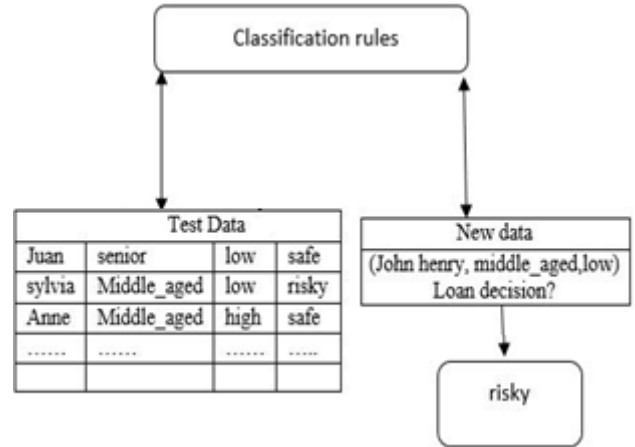
**Constructing the Classifier or Model:**

- This measure is the learning step or the learning phase.
- In this step the classification algorithms build the classifier.
- The classifier is construction from the training set made up of database tables and their associated class labels.
- Each tuples that constitutes the training band is referred to as a category or division. These boards can also be named to as sample, object or data points.



**Utilizing classifier for classification**

- In this measure, the classifier is applied for classification. Here the test data is applied to reckon the accuracy of classification rules. The classification rules can be enforced to the new data tuples if the accuracy is regarded from the applications satisfactory.



- To get with the mining driving sessions is used to obtain driving occasions from the application's chronicled positioning records and after that, it combines nearby driving occasions for building driving sessions. At that point, the positioning based proof dissect the underlying properties of driving occasions for separating misrepresentation confirmations. Audit based confirmation is used to check the surveys of the diligence. The K-nearest neighbors (KNN) calculation is used to raise strength and precision of the diligence. These all proofs are consolidated for recognizing the extortion applications.

**4.2 KNN Algorithm**

The  $k$ -nearest neighbor classifier cuts out hyper spheres in the space of instances by assigning the majority class of the  $k$ -nearest instances according to a defined metric (e.g., Euclidean distance). It is very sensitive to the curse of the dimensionality. The classification performance strongly depends upon the used metric. Moreover, a small value of  $k$  results in chaotic boundaries and makes the method very sensitive to outliers.

( $k$ -NN) is a method for classifying objects based on closest training examples in the feature space.  $k$ -NN is a type of lazy learning, or instance-based learning where the function is only approximated locally and this instance based learning is amongst the simplest of all machine learning algorithms: Whenever we have a new point to classify, we find its  $K$  nearest neighbors from the training data and then we find the distance from new point to  $k$  nearest neighbors. A new point is classified by a majority vote of its neighbors. If  $k = 1$ , then the object is simply assigned to the class of its nearest neighbour

1. Store the output values of the  $M$  nearest neighbours to query scenario  $q$  in vector  $R = \{r^1, \dots, r^M\}$  By recurring the following loop  $M$  times:
  - a. Go to the next  $s^i$  in the data set, where  $i$  is the current iteration within the domain  $\{1, \dots, P\}$
  - b. if  $q$  is not set or  $q < d(q, s^i): q \leftarrow d(q, s^i), t \leftarrow 0^i$
  - c. Loop until we reach the end of the dataset (i.e.  $i = P$ )
  - d. Store  $q$  into vector  $c$  and  $t$  into vector  $r$
2. Calculate the arithmetic mean output across  $r$  as follows:

$$\bar{r} = \frac{1}{M} \sum_{i=1}^M r_i$$

3. Return  $\bar{r}$  as the output value for the query scenario  $q$



**References**

- [1] Hengshu Zhu, Hui Xiong, Enhong Chen and Yong Ge, "Discovery of Ranking Fraud for Mobile Apps" ON KNOWLEDGE AND DATA ENGINEERING, VOL. 27, NO. 1, JANUARY 2015
- [2] (2012). [Online]. Available: <https://developer.apple.com/news/index.php?id=02062012a>
- [3] (2012). [Online]. Available: <http://venturebeat.com/2012/07/03/apples-crackdown-on-app-ranking-manipulation/>
- [4] N. Spirin and J. Han, "Survey on web spam detection: Principles and algorithms," SIGKDD Explor. Newslett., vol. 13, no. 2, pp. 50–64, May 2012.
- [5] A. Ntoulas, M. Najork, M. Manasse, and D. Fetterly, "Detecting spam web pages through content analysis," in Proc. 15th Int. Conf. World Wide Web, 2006, pp. 83–92.
- [6] E.-P. Lim, V.-A. Nguyen, N. Jindal, B. Liu, and H. W. Lauw, "Detecting product review spammers using rating behaviors," in Proc. 19th ACM Int. Conf. Inform. Knowl. Manage., 2010, pp. 939–948.
- [7] B. Zhou, J. Pei, and Z. Tang, "A spamicity approach to web spam detection," in Proc. SIAM Int. Conf. Data Mining, 2008, pp. 277–288.
- [8] S. Xie, G. Wang, S. Lin, and P. S. Yu, "Review spam detection via temporal pattern discovery," in Proc. 18th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 2012, pp. 823–831.
- [9] Z. Wu, J. Wu, J. Cao, and D. Tao, "HySAD: A semi-supervised hybrid shilling attack detector for trustworthy product recommendation," in Proc. 18th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 2012, pp. 985–993.
- [10] B. Yan and G. Chen, "AppJoy: Personalized mobile application discovery," in Proc. 9th Int. Conf. Mobile Syst., Appl., Serv., 2011, pp. 113–126.
- [11] K. Shi and K. Ali, "Getjar mobile application recommendations with very sparse datasets," in Proc. 18th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 2012, pp. 204–212.
- [12] H. Zhu, E. Chen, K. Yu, H. Cao, H. Xiong, and J. Tian, "Mining personal context-aware preferences for mobile users," in Proc. IEEE 12th Int. Conf. Data Mining, 2012, pp. 1212–1217.
- [13] H. Zhu, H. Cao, E. Chen, H. Xiong, and J. Tian, "Exploiting enriched contextual information for mobile app classification," in Proc. 21st ACM Int. Conf. Inform. Knowl. Manage., 2012, pp. 1617–1621.