# Describing and Predicting Early Reviewers for Efficient Product Marketing on E -Commerce Sites

[1]M.Ravindar, [2]Dr. R. Jegadeesan [3]T.Anusha, [4]N.Manisha, [5]K.Saimanith, [6]M.Sowmya

[3,4,5,6]Final year student Computer Science & Engineering, [1,2]Associate professor-CSE

[1,2,3,4,5]Jyothishmathi Institute of Technology & Science, Karimnagar, India

***ABSTRACT***: Now a days, the decision of purchasing an online product is relying on the reviews it get .This paper emphasizes an marketing a product based on its reviews .These reviews are mainly grouped into three types they are premature reviews, matured reviews and late responding reviews. However, this division is done on the basis of time span of a product as it has starting stage, central stage and finishing stage in its sale process. Comparison of these reviews is done on real time E commerce websites for example amazon and yelp. It is known that the impact of premature reviews is high on the customers when compared with matured and late responding reviews. As the product is newly available in the market and based on the reviews the scope of the project in sustaining will be decided. The product is judged as a good product/bad product by depending on its attributes like rating ,number of positive responses, Number of product sold in a specific time and so on.....The implementation of this process is done by collecting  the sample datasets of reviews on e commerce websites.

*Index terms* : premature reviewer, early review, embedding model

## 1. INTRODUCTION

The e commerce websites are useful for users to share their views on online purchasing which contains more information while purchasing the product. We categorized the reviews based upon lifetime. Users who posts reviews early are considered as Premature reviewers. These are helpful in determining the success or failure of the product.[1][2] and hence the companies need to identify them.

The fundamental role of premature reviewers has attracted great attention for the sale of products [3]. Amazon advocated a premature review program that helps acquire pre-mature reviewers in case there are fewer revisions. Amazon Vine invites the most reliable people who tend to post useful comments. The screening of premature adopters in the diffusion of innovation has attracted much attention from the research community. The following are some of the fundamental dissemination processes: Attribute of innovation. Communication channel Social network structures [5]. Innovation studies have been carried out extensively in social networks [6] - [8]. For any given product, reviewers are viewed based on the timestamp as premature, mature, and late-response reviews. To analyze the characteristics of premature reviews, we consider two important metrics:

1) Premature reviews tend to assign a higher score.

2) Premature reviews tend to publish useful criticisms.

In addition, we explain the findings of herd behavior [14] [15] widely studied in economics and sociology [9] [11] which refers to the fact that individuals depend to a large extent on others when making decisions. We predicted premature reviews based on the multiplayer competition game [12] [13].

The task of Premature reviewer prediction has received very less attention, our contributions are summarized as follows:

--We first characterize Premature reviews on two real world datasets(Amazon ,Yelp).

--We quantitatively analyze the characteristics which is having a high impact on the product.

--Extensive experiments have stated that the effectiveness of our approach got  the prediction of Premature reviewers.

## EXISTING SYSTEM:

- Previous studies have extremely emphasized the development that people are powerfully influenced by the selections of others, which might be explained by herd behavior. The influence of early reviews on ulterior purchase is Understood as a special case of swarming result. Early reviews contain vital product evaluations from previous adopters, that are valuable reference resources for consequent purchase selections. As shown in existing papers, once customers use the merchandise evaluations of others to estimate product quality on the web, herd behavior happens within the on-line searching method.

- Chen et al. have projected to use dimensional representations to capture each intransitiveness and context info for modeling pair wise comparison relations.

**DISADVANTAGES OF EXISTING SYSTEM:**

- Early prediction model don't seem to be correct.

- Existing works depends on extracting opinions or distinguishing opinion targets (or holders) from review knowledge.

- Most of those studies square measure theoretical analysis at the macro level and there's an absence of quantitative investigations

**PROPOSED SYSTEM:**

• To model the behaviors of the first reviewers, we tend to develop a scrupulous one thanks to characterize the adoption method in two sets of revision data of real world giants. Additionally, given a product, the reviewers adjusted according to their time stamps for commercial companies and their revisions.

• In our work here, we tend to specialize mainly in 2 tasks, the main task is to investigate the general characteristics of the first reviewers compared to the general reviewers and laggards. We tend to characterize their rating behaviors and also the utility scores received from others and also the correlation of their reviews with the quality of the product. The second task is to find a prediction model that predicts the first reviewers of a product.

**ADVANTAGES OF PROPOSED SYSTEM:**

• A higher average score in the initial reviews probably points to a higher product quality.

• It is likely that a higher utility score of the initial reviews extends or decreases the quality of the product

• We provide a primary study to characterize the first reviewers of real-life giant data sets associated with associated e-commerce websites.

• We quantitatively analyze the characteristics of the first reviewers and their impact on product quality. Our empirical analysis supports a series of theoretical conclusions from social sciences and social sciences.

• We read the opinion publishing method as a multiplayer competition game associated with the development of a classification model based on the inlay for the prediction of the first reviewers. Our model will affect the inconvenience of cold start by incorporating faceted data of the merchandise.

• Extensive experiments on 2 giant real-world datasets, namely Amazon and Yelp have questioned the effectiveness of our approach to predicting the first reviewers.

**2. PRELIMINARIES**

The notations and tokens are used in the paper are mentioned below with their description in the following table

TABLE 1
Notations and Descriptions.

| Notations | Descriptions |
|---|---|
| $\mathcal{U}$ | a set of e-commerce users, $u \in \mathcal{U}$ |
| $\mathcal{P}$ | a set of e-commerce products, $p \in \mathcal{P}$ |
| $r, s$ | rating $r$ posted by a user with a timestamp $s$ on a product |
| $n_Y, n_N$ | the number of 'yes' votes and 'no' votes a review received |
| $d$ | a review $d$ is composed of $\langle u, p, r, s, n_Y, n_N \rangle$ |
| $c_p, t_p$ | the category label $c_p$ and title description $t_p$ of a product |
| $\mathcal{L}_p$ | a list of ordered reviews of a product $p$, $\mathcal{L}_p : d_1 \rightarrow d_2 \dots \rightarrow d_i \dots \rightarrow d_{N_p}$ |
| $\Delta_L^{(p)}$ | the leading gap for product $p$ |
| $\Delta_T^{(p)}$ | the trailing gap for product $p$ |
| $\Delta_M^{(p)}$ | the maximum interval for product $p$ |
| $\boldsymbol{v}_p, \boldsymbol{v}_u$ | low-dimensional representation vector of product $p$ and user $u$ |
| $\boldsymbol{v}_{t_p}, \boldsymbol{v}_{c_p}$ | title embedding and category embedding of product $p$ |
| $S(p, u)$ | the likelihood that user $u$ becomes an early reviewer of product $p$ |

A user ụΣŲ posts a review d for a product p Σ P with rating r and timestamps s. The votes may be $n_y$ $'yes'$ or $n_n$ $'no'$ from the users. A product has Category label $C_p$ and Title description $t_p$ .here $N_p$ is the number of reviews for a particular product. Hence we form a ordered review list as $L_p: d_1 \rightarrow d_2 \dots \rightarrow d_i \dots \rightarrow d_{N_p}$ on the basis of $s_i$ timestamp $d_i$

## 2.1 Products with complete lifetime :

A review time span of product is the time span between first and last review. According to notation $[S_1, S_{N_p}]$ .our observation window contains amazon and yelp datasets. Amazon dataset is of 18 years [may 1996 to july 2014] and yelp dataset is of 13 years [july2004 to january 2017].
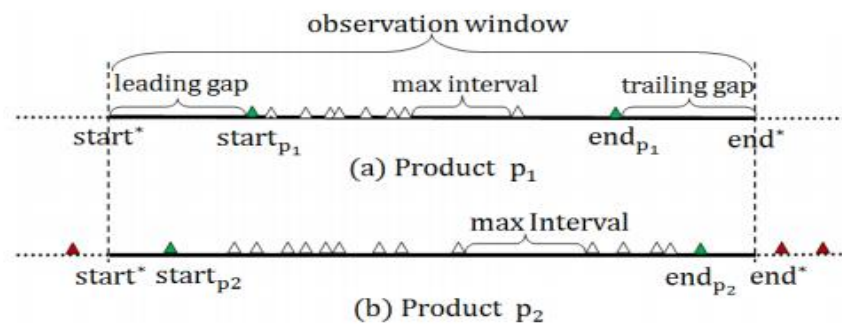


Fig 1

## 2.1.1 Determining complete review time span:

Our observation window [fig 1]contain leading and trailing gap. That is the timestamps of observation window is $[S_{start}\ and\ S_{end}]$. So the gap between $S_{start}$ and $S_1$[the first review] is called leading gap denoted by$\Delta_L{}^{(p)}$ and the same between $S_{end}$ and $S_{N_p}$[the last review] is called the trailing gap denoted by $\Delta_T{}^{(p)}$. The time span is obtained by

$$\Delta_M{}^{(p)} = \max_{i=1}{}^{N_p-1} (S_{i+1} - S_i).$$

Should satisfy, $\Delta_L{}^{(p)} > \Delta_M{}^{(p)}$ and $\Delta_T{}^{(p)} > \Delta_M{}^{(p)}$

## 2.1.2 Estimating the product lifetime:

According to the available information it is not possible to measure the accurate lifetime of a product but estimated time is measured based on the review time span.

## 2.2  Early reviewer identification:

Given a product lifetime studying the product lifetime of a product and dividing it into stages. At first we classify the reviewers as five types based on the timestamps they post their reviews. they are Innovators, early adopters, early majority, late majority and laggards[fig 2]. Following [5] we using classic Roger's bell curve to divide the product lifetime into five stages [fig 2]. Our dataset contain less number of Innovators and early adopters so later they were grouped into one stage premature reviewers, early majority and late majority grouped into matured reviewers and laggards as late-responding reviewers. The probabilities of these three stages are [0, 0.16], [0.16, 0.84] and [0.84, 1].
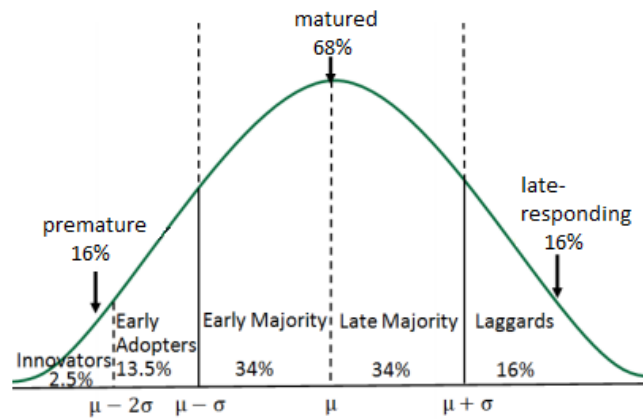
Fig 2

**Definition of early review and premature reviewer:**

If a user u posting review d for a product p with timestamps s falls in the probability range of [0, 0.16] then it is called early review and the reviewer is considered as premature reviewer. These particular reviews are product specific. The paper emphasizes on early reviews because these reviews show high impact on the marketing and also influences the adopters in adopting the products.

The following section 3 focuses on dataset, section 4 on analysis, section 5 on margin bases ranking model, experimental setup and results in section 6, related study and conclusion is contained by section 7 and section 8 respectively.
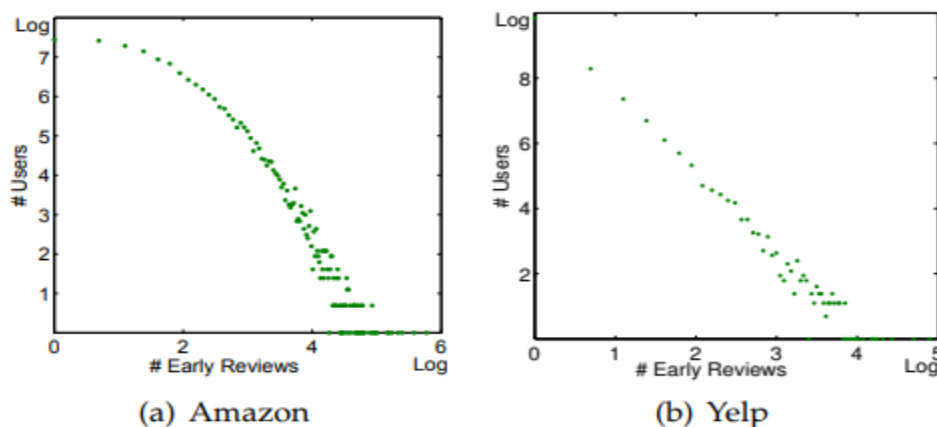


(a) Amazon          (b) Yelp

Fig 3

**3. DATA PREPARATION**

In this paper we tend to use Amazon[4],[17] and yelp datasets. Amazon originally contain 142.8 million products starting from May 1996 to July 2014 whereas Yelp contains 4.7 million products starting from July 2004 to January 2017. Each review are printed primarily based upon timestamp and rating. For any review, Amazon users will vote by mistreatment Yes(or) No buttons that indicates positive and negative angle and is recorded whereas Yelp cannot.

**3.1 Data improvement:**

Data improvement contains a pair of steps as follows:

**3.1.1 Preprocessing:**

In this method we tend to take away duplicate reviews, reviews from anonymous users, inactive users, unpopular products and we only keep the users who have denote between 5 to 10 reviews for every product.

### 3.1.2 Review sender Detection and Removal:

Firstly our focus are on early users. As shown in (23) the amount of spam reviews has been accrued and could be denote by review spammers that predicts false opinions on products and influence the shoppers. As spam reviews ends up in erronous conclusion we want to get rid of spammers as a part of our knowledge improvement method.

Here, we adopt the approach projected in (24) for removing spam reviews that considers 3 factors:

Early Deviation Spamming(ED)

Review Text Spamming (RT)

Time primarily based Spamming (TS).

Linear Regression model combines these factors to form judgement and might be calculated as

$S(u) = \alpha SED(u) + \beta SRT(u) + \gamma ST\, S(u).$

where , SED(u), SRT (u) and ST S(u) are the scores of spamming behaviour .

$\alpha, \beta, \gamma$ are the tuning loads for consolidating these three variables and we have $\alpha + \beta + \gamma = 1$. In our tests, we experimentally set $\alpha = \beta = \gamma = 1/3$.We finally known 4.65%, 4.53% of spam users in Amazon and Yelp datasets severally
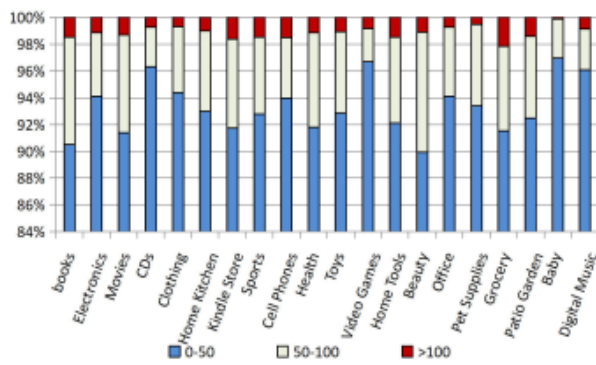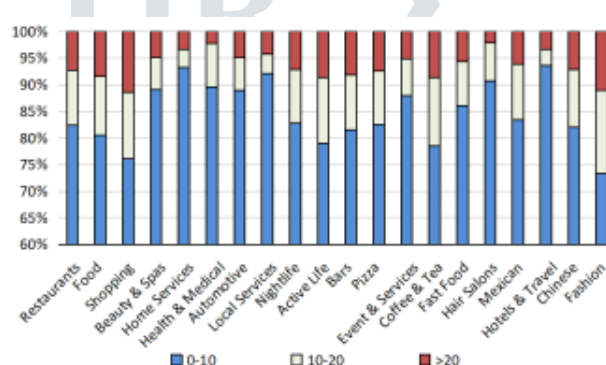


Fig:4                    Fig:5

### 3.2 Basic applied math Analysis:

Using the categorization thresholds in fig a pair of we tend to divide reviews as Premature, Mature and Late responding reviews. For any product, review with premature label is taken into account as premature review and therefore the user is considered as premature reviewer.

### Premature reviews gift power-law likelihood distribution:

In fig3. the graph represents variety of users v/s number of premature reviews. Each figure indicates overwhelming majority of users who acted as premature reviewers. About 70-85% users of Amazon and Yelp are thought-about as premature reviewers.

### Product class influences users enthusiasm of adopting new merchandise:

We more examine statistics by considering every product. Our dataset contain merchandise from 20 main classes. For each class we tend to cipher early reviews. We tend to descritize users into 3 bins supported total variety of early reviews:$(0,50)(50,100)(100,+\infty)$in Amazon and $(0,10)(10,20)(20,+\infty)$in yelp as shown in fig 4,5. Different product classes tends to urge completely different variety of premature reviews from users.

## 4.QUANTITATIVELY ANALYZING THE CHARACTERISTICS OF EARLY REVIEWERS

As early adopters are important to the diffusion of innovations [5]. Premature reviewers play a key role in future product adoptions. There has been a lack of quantitative analysis of the correlations between the premature reviewers and product adoptions on large datasets, i.e., Amazon and Yelp. In this section, we study how early reviewers are different from others and how they impact product popularity.

### 4.1 Characteristics of Early Reviewers:

To understand the difference between premature reviewers from others we compare the premature reviews with overall reviews and helpfulness scores voted by others. Using the categorization methodology discussed in Section 2, we tend to assign every review into one among the three classes outlined in Fig 2 .The rating score of every review is in a five-star scale. For helpfulness, in Amazon dataset, we tend to count the amount of Yes and No votes respectively and so normalize them to the range of [0; 1].While in Yelp dataset, helpfulness can be gained by clicking the helpful button. We tend to count the amount of Useful as the review's helpfulness score. Given the 3categories of reviews, we compute the average ratings and helpfulness scores in every review class.

**Premature reviewers tend to assign a higher average rating score:**

In Fig6: we observe that premature reviews are additional doubtless to come with the next rating score than those from the opposite two categories. Note that we have removed spam reviews since their ratings tend to be extreme, either too high or too low.

**Premature reviewers tend to post more helpful reviews:**

Based on the average helpful scores of reviews by the 3 categories in fig 7. Amazon dataset having both Yes and No votes of reviews but we only consider yes as the helpfulness scores while in Yelp dataset ,we use the number of useful votes as the helpfulness score indicates that Premature reviews ate more important than others  .This might be effect by the growth of time reviews :premature reviews themselves tend to receive more attention .To reduce the effect of time span ,in Amazon dataset ,we report both Yes and No counts and normalize these votes by categories in Table 2.The higher normalized Yes votes are premature reviews are be likely to post more useful reviews .Fig 8 shows box plot for the distribution of reviews length in 3 categories. It is observed that on average premature reviews are longer than other 2 categories .By inspecting reviews we found that premature reviews tends to score more helpfulness compared to others.
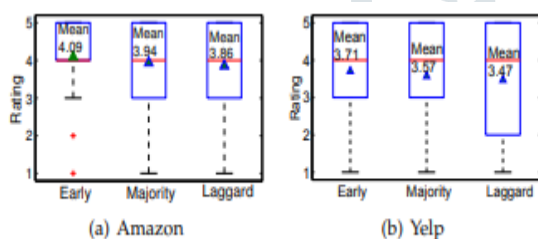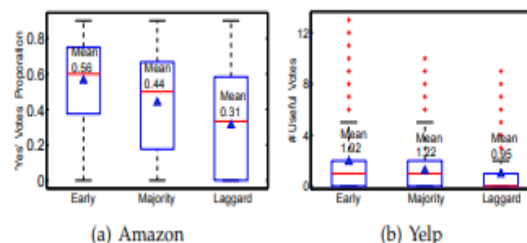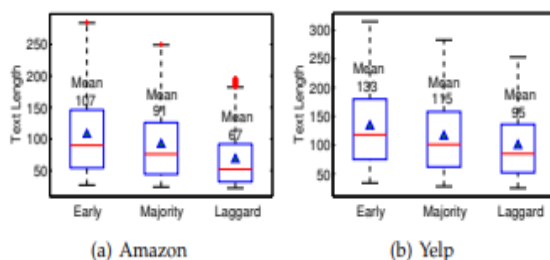


(a) Amazon  (b) Yelp

Fig:6



(a) Amazon  (b) Yelp

Fig:7



(a) Amazon  (b) Yelp

Fig:8

**TABLE 2**
Comparisons of the helpfulness scores by the three categories of reviews in the Amazon dataset.

| Categories | No. of 'Yes' | No. of 'No' | Normalized 'Yes' | Normalized 'No' |
|---|---|---|---|---|
| Early | $15.5 \pm 0.21$ | $3.28 \pm 0.05$ | $0.72 \pm 0.002$ | $0.28 \pm 0.002$ |
| Majority | $4.98 \pm 0.05$ | $2.28 \pm 0.02$ | $0.68 \pm 0.001$ | $0.32 \pm 0.001$ |
| Laggards | $2.23 \pm 0.04$ | $1.44 \pm 0.03$ | $0.67 \pm 0.003$ | $0.33 \pm 0.003$ |

**Connection with personality variables theory:**

Our previous findings may find relevance in well-known principles in the theory of personality variables, which mainly studies how innovation extends over time between members [8]. The theory emphasizes two important features of early adopters. Principle on personality variables: the premature adopter shares a more favorable attitude towards changes than later

adopters. Principle on the behavior of communication: premature adopters have a greater degree of opinion leadership than later adopters. We can relate our findings to the theory of personality variables in the following way:

1) the highest average rating scores can be considered as a favorable attitude toward the products;

2) the most useful votes of premature reviews given by others can be seen as a proxy measure of opinion leadership.

Therefore, our analysis results are consistent with the theory of personality variables and provide empirical evidence to the latter

## 4.2 The Impact on Product Popularity :

In this subsection, we investigate how premature reviews impact product popularity. As we don't have the actual product purchase transactions in our datasets. However ,the customers usually write the reviews after purchasing the product we can predict the popularity based upon the reviews posted in majority stages for every product. As reviews in later group introduce noises, we remove them for popularity value calculation. For ease of analysis, we first discretize both kinds of continuous scores into disjoint value intervals. For ratings, we use four bins: [1; 2], (2; 3], (3; 4], (4; 5] in Amazon and Yelp datasets. For helpfulness scores, we discretize the helpfulness scores of [0,1] into two consecutive bins in Amazon dataset, i.e., A : [0; 0:5], B : (0:5; 1]; In Yelp dataset, the helpful score is to show  as the number of Useful votes. We first compute the median and then use the median to discretize a helpfulness score into two bins ,namely A : [0; median] and B : (median;+1).
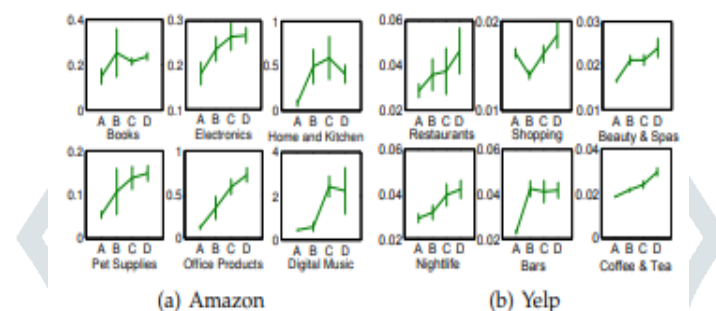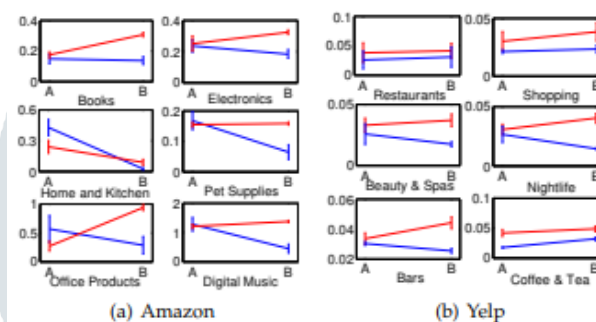


(a) Amazon      (b) Yelp

Fig:9



(a) Amazon      (b) Yelp

Fig:10

**A higher average rating score of premature reviews is likely to indicate a higher product popularity**.

Given a product, we first calculate its average rating score of its premature reviews, and then assign it to the corresponding rating bin defined above. We present the average daily popularity over the  products in different product categories by rating bins in Fig 9. The figure shows an upward trend with increase of the average rating scores from premature reviews

**A higher helpfulness score of premature reviews is likely to increase or decrease product popularity**.

Different from rating scores, a high helpfulness score does not necessarily indicate positive opinions towards a product. If a negative review gives very good reasons behind its negative feedback, the review is likely to be perceived as helpful by other customers and hence would receive a fairly high number of Yes . Hence, to examine the impact of review helpful scores, we need to discriminate between two types of reviews which are indicated with stars. We report the daily popularity with different helpful bins in Figure 10. It is interesting to observe that a higher useful score on positive reviews as well as negative reviews.

.

**Connection with the herd behaviour theory :**

In economics and humanism, herd behavior is an important concept [9]–[11], which describes that individuals are strongly influenced by the decisions of others in some situations. It  emphasizes the influence from existing adoption behaviours of others, especially the early adopters. It has been verified that online herd behavior occurs when people use the product assessment of others to indicate product quality on the Web [10]. Our findings can be explained with herd behaviour theory in a specific settings ,where existing on web selection data is reflected via rated reviews .We describe the influence of early adopters using two metrics ,namely average rating and helpful scores .Our work provides clear quantitative evidence  to the herd behavior through the analysis on two real-world e-commerce datasets. We also observe that a negative impact of Premature adopters on product sales when they assign low rating score to the products.

## 5. PREDICTING PREMATURE REVIEWERS

As Premature reviews are important to product popularity. These predicting Premature reviews has 2 benefits: First, identifying premature reviewers is helpful to monitor and manage early promotion. Second, premature reviewers are very likely to be the actual adopters of a product, leading to direct purchase.

### 5.1 Problem Formulation:

Given a product p and a candidate user set $U_p$: $\{u_1, u_2,..., u\ N_p\}$, produce a top-K list of users from $U_p$, who would post reviews on p at the early stage of product p in market can be determined as a ranking problem. We propose to use a ranking function S(p,u) measures the likelihood that user u becomes an premature reviewer of product p. Each training instance consists of a product pi with a complete lifetime $L_i$: $(u_1^{(i)} + s_1^{(i)}) \rightarrow (u_2^{(i)} + s_2^{(i)}) \rightarrow$…………...

$(U_{NP}^{(i)} + S_{NP}^{(i)})$ is an ordered list of reviewers $\{u_j^{(i)}\}$ on $p_i$ by the timestamps $\{s_j^{(i)}\}$.A major challenge is that our task is a cold-start ranking problem. Inspired by previous cold-start recommendation algorithms [20], a product p is with a category label cp and a title description tp are used to learn product representations or embeddings as will be discussed in Section 5.2.

Given a product p and two candidate users u and $u^1$, we seek to model the partial order between them. We consider the review posting process as multiplayer competition [26].We use u $>p^{ui}$ denote that user u has an premature review timestamp than u0 for product p. It has been explored for community question answering [27] and player ranking [21].

### 5.2  A Margin-based Embedding Model for Predicting Early Reviewers:

The essence of this task is to model the partial order between two candidate users u and u0 given a product p. We can cast the total order ranking problem into a pairwise comparison problem in distributed representation learning [12], [13].

### 5.2.1 Modeling the Pairwise Comparison:

We can define the objective function S(p,u) as an inner product between user and product embeddings, i.e.,

$$S(p,u) = v_p^T.v_u$$

In the embedding space, it is expected that $v_p^T.v_u > v_p^T.v\ u^1$ when u $>p^{ui}$. Given the original training set A = $\{(p_i, l_i)\}$, we first transform them into a set of partial order pairs T = $\{$ u $>p^{ui}$ |u,u$^i \in L_p\}$, where $L_p$ is the reviewer list of product p. To learn such embeddings, we minimize a margin-based ranking criterion [17] over the training set

$$L(T) = \sum_{u\ >pui\ \Sigma\ T} [m+S(p,u^1)-S(p,u)]_+,$$

$$= \sum_{u\ >pui\ \Sigma\ T} [m+ v_{u1}^T.v_p - v_u^T.v_p]_+,$$

where $[x]_+ = \max(0,x)$ and m is the margin coefficient set to 0.1 in our experiments. The objective function in Eq. 2 is very intuitive. When u $>p^{ui}$ and S(p,u) < S(p,u$^1$), there would incur a cost. We would like to optimize the objective function by trying to fit all the partial order pairs u $>p^{ui}$.

### 5.2.2 Learning the Product Embeddings:

A major problem with the above objective function is that the learning of product embeddings relies on the past review data. We can utilize the learned word to derive the product embeddings in current cold start setting. To achieve this, a simple method will be to aggregate the embeddings in the title description of a product from doc2vec model[23].The doc2vec can be used to summarize the information of entire text in a document. In doc2vec,a doc ID is incorporated into the context of word i.e. Pr(wi|wi−c : wi+c,tp), where tp represents the title of a product p. To incorporate both the title and category label, we model the generative probability Pr(wi|wi−c : wi+c,tp,cp).

**Algorithm 1** The learning algorithm for user embeddings.

**Input** training instances $\mathcal{T} = \{u \succ_p u' | u, u' \in \mathcal{U}\}$,
　　　　products embeddings set $\{v_p\}$,
　　　　learning rate $\lambda$,
　　　　margin coefficient $m$,
　　　　embedding dimensions $L$.

**Output** user embeddings $\{v_u | \forall u \in \mathcal{U}\}$

**Procedure:**
1: initialize user embeddings:
2:　　$v_u \leftarrow uniform(-\frac{6}{\sqrt{L}}, \frac{6}{\sqrt{L}}), \forall u \in \mathcal{U}$
3:　　$v_u \leftarrow v_u / \|v_u\|, \forall u \in \mathcal{U}$
4: **loop**
5:　　sample a training instance $\langle u \succ_p u' \rangle \in \mathcal{T}$ **do**
6:　　update user embeddings:
7:　　　　$v_u := v_u - \frac{\partial \ell(\mathcal{T})}{\partial v_u}$,
8:　　　　$v_{u'} := v_{u'} - \frac{\partial \ell(\mathcal{T})}{\partial v_{u'}}$.
9: **until** convergence

### 5.2.3 Learning the User Embeddings:

To learn the embedding parameters in Eq. 2, we can simply apply the Gradient Stochastic Descent (SGD) to update the user's incrustations {vu} and the product incorporation {vp}. To handle the cold start problem, we incorporate the title and category information to learn the embeddings vp product in advance. The detailed optimization procedure is described in Algorithm 1. All inlays for users are initially initialized randomly according to a uniform distribution, the strategy proposed in [24]. In each main iteration of the algorithm, a training triplet hp, u, u0i, where we have the partial order up to u0, from the training set of the classification function based on the margin l (T). The total number of product revisions (n) and total number of comparison pairs generated O (n2). Therefore, we maintain all comparison pairs that include a premature reviewer, while others are selected with random sampling.

## 6 EXPERIMENTS ON EARLY REVIEWER PREDICTION

In this section, we have a tendency to conduct experiments to judge our proposed margin-based embedding model for early reviewer prediction.

**TABLE 3**
Statistics of the evaluation sets in early reviewer prediction. ANRU and ANRP are the abbreviations of Average Number of Reviews posted by each User and Average Number of Reviews received by each Product.

| Dataset | #Product | #User | #Pairs | ANRU | ANRP |
|---------|----------|--------|-----------|------|------|
| Amazon | 12,814 | 16,355 | 3,122,797 | 18 | 23 |
| Yelp | 2,545 | 3,912 | 282,718 | 14 | 22 |

### 6.1 Datasets:

Since it's unreliable to incorporate users or product with only a few reviews for analysis, we have a tendency to take away the product that are related to but fifty reviews in Amazon dataset and ten reviews in Yelp dataset, and users UN agency announce but fifty reviews in Amazon dataset and ten reviews in Yelp dataset. The statistics of the information sets employed in our experiment are shown in Table three. Note that "#Pairs" indicates the whole variety of comparison pairs that may be generated in our analysis set following the strategy mentioned in Section five.2. Given a product, though its associated reviews in our analysis set are solely a set of all reviews found regarding this product within the original dataset, the temporal arrangement of those reviews (and the corresponding reviewers) remains the identical. we have a tendency to assign the class labels to reviewers based mostly on the initial dataset and use them as our ground truth.

### 6.2 Evaluation metrics:

Given a product, every candidate methodology can turn out Associate in Nursing ordered list of users. Hence, we have a tendency to adopt 3 ranking-based metrics for analysis of predicting results.

**Overlapping Ratio relation at rank k** (OR@k). Given the expected ordered list of users for a product, OR@k is outlined as:

OR@k = |L(k) ∩ G(k) | k ,

where  L(k) and G(k) denote the sets of users came by a candidate methodology and obtained by sorting consistent with actual timestamps for the primary k reviewers severally. Note that once k is larger than the particular variety of early reviewers given a product, G(k) would contain users  who are not early reviewers.

**Hit Ratio  at rank k (Hit@k).**

Given the expected ordered list of users for a product, Hit@k is defined as:

$$\text{Hit@k} = \frac{\sum_{i=1}^{k} I(p, ui)}{Np(E)}$$

where I(p, ui) returns one if $u_i$ was an premature reviewer  for product p in original dataset, and zero otherwise; and $N_p^{(E)}$ is the actual variety of early reviewers for product p.

**Ratio of Correct Comparison Pairs (RCCP).**

Since our model is trained from comparison pairs, we have a tendency to use RCCP jointly to measure the standard of pair classification, which is described as:

$$\text{RCCP} = \frac{\#correctly\ predicted\ pairs}{\#test\ pairs}$$

Note that we have a tendency to don't adopt ranking-based correlation coefficient as analysis metrics (e.g., Spearman or Edward Kendall Tau). For our task, the standard of high predictions for premature reviewers   are increasingly critical to consider. Henceforth, we essentially utilize the previously mentioned measurements for best k positioning.

**6.3 Methods to Compare for Early Reviewer Prediction:**

Our  task  is to anticipate who will turn out to be early analysts of an item. We think about three sorts of techniques for comparisons: measurements based strategies, rivalry based models what's more, our edge based inserting positioning model.

6.3.1 Simple Statistics-based Methods :

A direct way to deal with this undertaking is to compute the number of times (or the proportion) that a client has gone about as an early commentator in history information. Instinctively, if a client has posted numerous early audits before, she is additionally likely to post early audits on another item. So we utilize the following measurements to gauge clients' positioning score of being early commentators.

• NR: Rank the clients basically dependent on the Number of Surveys (NR) that they have recently posted.

• NER: Rank the clients dependent on the Number of times that a client has recently gone about as an Early Reviewer (NER).

• PER: Rank the clients dependent on the Proportion that a client has gone about as an Early Reviewer (PER). PER is characterized as:

PER(u) = NER(u) /NR(u)

• SPER: Rank the clients dependent on Smoothed PER. The PER may be one-sided when NR is little. We propose to utilize the Smoothed Proportion that a client goes about as an Premature Reviewer (SPER), which is characterized as:

$$\text{SPER(u)} = \frac{NR(u)}{NR(u)+P} PER(u)$$

$$+\frac{p}{NR(u)+P}\text{PER}_{avg}$$

where ρ = 1 |U | · P u∈U NR(u), and P ERavg = 1 |U | · P u∈U P ER(u)

The above measurements based techniques are just ready to create a solitary rank list of clients for every one of the items, which can't use pair wise examinations and the item data. We further propose rivalry based models what's more, edge based implanting positioning model.

### 6.3.2 Competition-based Models:

The challenge based techniques take rivalry connection into thought, which we use in our errand of foreseeing premature analysts. We think about four strategies for examination.

• TS: [26]: TrueSkill is a Bayesian ability rating framework which is intended to figure the relative aptitude levels of players in multiplayer amusements. It accept that the down to earth aptitude dimension of every contender u pursues a ordinary circulation $N(\mu u, \sigma 2\ u\ )$, where $\mu$ is the average aptitude level and $\sigma$ is the estimation vulnerability. In our trial, we set the underlying estimations of the ability level $\mu$ and the standard deviation $\sigma$ of every player to the default esteems utilized in [21].

• SVMComp: [27]: The SVMComp display learns the weight of every client dependent on pair wise correlations utilizing the exemplary Support Vector Machine (SVM). Given a two-player rivalry k with a champ u and a failure u 0 , there are two preparing occasions created: $y_a = 1$, $x_a[u] = 1$, $x_a^1[u\ 0\ ] = -1$ and $y_a^1 = 0$, $x_a^1[u] = -1$, $x_a^1[u^1] = 1$. We utilize the toolbox SVM Lib Linear with direct part.

• B-T [31]: Bradley-Terry (B-T) show is a likelihood show that can foresee the result of a correlation. It learns a scalar parameter for every one of the player from memorable pair wise examination information. These parameters as a rule speak to the positions or qualities of people, with higher positions favored for the success over lower positions in future examinations. Following the technique in [25], we utilize the most extreme probability estimation to acquire the quality $\gamma u$ of client u.

• B-C [25]: The above techniques utilize a solitary number to speak to a player, which is somewhat oversimplified. In differentiate, Blade-Chest (B-C) demonstrate learns a multidimensional portrayal for every player from pair wise correlations. We receive the open-source code to actualize this model5.

The above four models consider pair wise examinations between clients, however despite everything it cannot use the data from the item side. At the end of the day, the halfway request of two clients continues as before for every one of the items in these models. Thus we propose our edge basedimplanting  positioning model which includes both challenge comparisons and the data of the items and learns the portrayal of clients consequently.

### 6.3.3 Margin-based Embedding Model :

This is our proposed Margin-based Embedding Ranking Display (MERM) proposed in Section 5.2. As far as anyone is concerned, no past examinations connected implanting models for predicting premature commentators. Our model can portray both client correlation relations and the data from the item side. Henceforth, it is relied upon to give preferable execution over the above standard strategies. As of now, we for the most part use the title and class data. It will be direct to  use different sorts of item data as the setting of a challenge between two clients.

For every one of the techniques, we report the execution utilizing five-overlap cross-approval. Note that we split information into five folds dependent on items, i.e., the whole audits of an item are either in the preparation set or test set. The parameters of a strategy are streamlined utilizing cross-approval. In B-C, we set the quantity of measurements of edge and chest vectors to 200 and 300 in Amazon and Yelp datasets separately. In our model MERM, we likewise set the quantity of installing measurements $2L = 200$ and $2L = 300$ in Amazon and Yelp datasets individually. For every item, we consider all the clients who have posted an audit of it as applicant clients. To make the assessment progressively reasonable, we additionally test five times of "negative" clients who did not audit the objective item however survey different items in a similar classification.

TABLE 4
Performance comparison on the results of early reviewer prediction.

| Datasets | Amazon | | | | | Yelp | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Models | OR@5 | OR@10 | Hit@5 | Hit@10 | RCCP | OR@5 | OR@10 | Hit@5 | Hit@10 | RCCP |
| NR | 0.0910 | 0.1416 | 0.1105 | 0.2088 | 50.15% | 0.0704 | 0.1187 | 0.0605 | 0.1110 | 55.26 % |
| NER | 0.1018 | 0.1516 | 0.1260 | 0.2131 | 61.17% | 0.0810 | 0.0982 | 0.1134 | 0.2052 | 60.53% |
| PER | 0.1114 | 0.1577 | 0.1334 | 0.2218 | 64.96% | 0.0738 | 0.0896 | 0.0971 | 0.1794 | 56.21% |
| SPER | 0.1125 | 0.1614 | 0.1353 | 0.2261 | 65.31% | 0.0763 | 0.1025 | 0.1060 | 0.2149 | 57.27% |
| B-T | 0.0931 | 0.1437 | 0.1120 | 0.2050 | 64.31% | 0.0864 | 0.0939 | 0.1044 | 0.1859 | 59.89% |
| B-C | 0.1132 | 0.1635 | 0.1361 | 0.2390 | 62.23% | 0.0931 | 0.1016 | 0.1120 | 0.1952 | 59.36% |
| TS | 0.1265 | 0.1720 | 0.1450 | 0.2465 | 67.54% | 0.0904 | 0.1013 | 0.1350 | 0.2300 | 59.82% |
| SVMComp | 0.1283 | 0.1747 | 0.1483 | 0.2503 | 67.97% | 0.0955 | 0.1045 | 0.1341 | 0.2201 | 60.13% |
| MERM | 0.1524* | 0.2273* | 0.1665* | 0.2823* | 69.25%* | 0.1212* | 0.1333* | 0.1462* | 0.2360* | 68.57%* |

Note: " * " indicates the statistically significant improvements (i.e., two-side t-test with $p < 0.01$ ) over the best baseline.

### 6.4 Results and Analysis:

We present the outcomes on premature commentator forecast in Table 4. It tends to be seen that the most straightforward gauge of positioning clients dependent on the quantity of audits posted previously (NR) plays out the most exceedingly awful. It

shows that clients posted a vast number of audits are not really dynamic in ahead of schedule selection of items. NER improves over NR, which appears that a client who has gone about as an early analyst for other items before is bound to receive new items in the future. PER, outflanks NER in Amazon dataset, while fails to meet expectations NER in Yelp dataset. The smoothed PER, i.e., SPER, performs superior to PER. The two comparison based baselines B-T and B-C beat the insights based strategies just now and again, and don't yield noteworthy improvement. These outcomes are reliable with the finding recently detailed in [22] that a straightforward proportion based technique functions admirably when the preparation information is adequately vast. Overall, B-C performs superior to B-T. Rather than utilizing a solitary esteem, B-C embraces a vectorized portrayal for modeling the player quality. Besides, the two competition based strategies TS and SVM Comp enhance all the above baselines. In spite of the fact that SVM Comp is marginally superior to TS, there is no noteworthy distinction between them. TS is an exemplary challenge display for portraying the player quality, while SVM Comp has been appeared to be compelling in QA master discovering undertaking [27]. These two techniques perform best among our baselines.

Our proposed model MERM accomplishes noteworthy improvement in contrast with every one of the baselines. Thought about with different baselines which just measure the earliness level of a client with a solitary esteem, MERM learns the multi dimensional portrayal of clients from relative sets.Despite the fact that B-C likewise embraces a multi-dimensional representation for demonstrating player quality, it doesn't perform very well in our undertaking. A conceivable reason is that B-C needs to learn more parameters (i.e., both sharp edge vectors and chest vectors); while, in our datasets, the correlation sets for preparing are scanty. The key contrast of MERM is that it learns item embeddings additionally dependent on the side data involving both the title and classification data of items. It adequately extends both item and client embeddings into the equivalent ceaseless space for direct correlation and positions clients by streamlining an edge based positioning goal work in an item reliant way. In our second arrangements of tests, we further look at the effect of the measure of preparing information on the outcomes of early commentator expectation. We present the consequences of Amazon dataset, the aftereffects of Yelp dataset are comparable and are excluded here. By fixing the test information at 20%, we fluctuate the staying 80% preparing information at five distinct parts :   {20%, 40%, 60%, 80%, 100%}. The outcomes are exhibited in Figure 11(a). Generally, we see that every one of the strategies endure from execution drop with the decrement of preparing information. Our technique MERM performs commonly superior to other strategies with any measure of preparing information. We likewise shift the quantity of measurements (i.e., 2L) for client and item portrayal in B-C and MERM, and report the outcomes in Figure 11(b). It tends to be seen that the dimensionality of 200 yields the best execution.
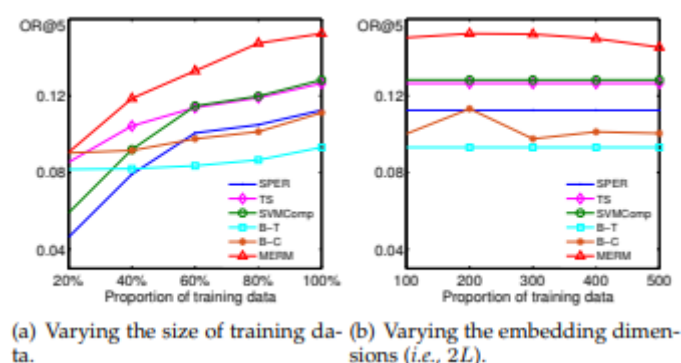


(a) Varying the size of training data.   (b) Varying the embedding dimensions (*i.e.*, 2L).

Fig:11

# 7. RELATED WORK

Our current study is mainly related to the three lines of research they are "early adopter detection , modelling comparison based preference, distributed representation learning".

## 7.1Early adopter detection

The term "early adopter" indicates from the classic theory for diffusions of innovations[5]. The further process of early adopters has been studied in sociology and economics. The early adopters are basically given priority in trend prediction, viral marketing, product promotion. Moreover, early adopters is clearly related to herding effect [3], [9], [11], [14-15], which describes are strongly influenced by decisions of other such as in stock market, decision making, social marketing and product success.

For eg: in product marketing, consumers select popular brands as they believe that popularity indicates better quality [9]. some further investigations also reveal that product evaluations from previous reviewers, such as star ratings, sales volume and customer's choices [9]. Three fundamental elements of this process are attributes of an innovation, communication and social

networking [5]. The theoretical analysis of macro level[2], [26] with rapid growth of online social platforms and including resource- constrained network [6] (or)retweet network [7], user click graphs [8] and text based innovation network.

## 7.2 Modelling comparision- based preference

Comparision – based preference has studied for several years [27], [28] and classic approacher and methods[29]. By this preference, we can perform and ranking task. For eg: In information retrieval (IR), the ranking for a list of candidate items with selected features [30], and paint wise, pair wise and list wise methods [31] are the categories for rank approaches. We use competition based ranking methods in games and matches, where the skill level of each player [32]-[33] and skill rating of an individual player. For eg : two player model [34], true skill system [21], Bradley-terry model [35], [36]. A part from this, many studies try to model the expert skill in a particular field  of a user using composition based ranking approach, for eg: community questioning and answering platforms [22], [38] and crowdsourcing systems [37], [39].

## 7.3 Distributed representation learning :

Distributed representations learning  has been used in various application areas like natural language processing (NLP), speech recognition and computer vision. For eg, in NLP, several semantic embedding models have been proposed including word embedding [12], phrase embedding [40], sentence embedding [41], word embedding models such as word2vec [12]. Specially, word2vec has two major model architectures as skipgram (SG) and continuous bag of words (CBOW). SG predicts current word as contexts. Based on word2vec, doc2vec[41] incorporates the document into word2vec model. In addition to this model, the various applications in other fields as network analysis[42] and recommendation [43], [44], [45].

## 8. CONCLUSION

In this paper, we've studied the novel task of Premature reviewer describing and prediction on 2 real-world on-line review datasets. Our empirical analysis strengthens a series of theoretical conclusions from social science and political economy. We found that associate degree premature reviewer tends to assign a better average rating score; associate degreed an premature reviewer tends to post additional useful reviews. Our experiments conjointly indicate that premature reviewers' ratings and their received helpfulness scores are probably to influence product quality at a later stage. we've adopted a competition-based viewpoint to model the review posting method, and developed a margin based embedding ranking model (MERM) for predicting Premature reviewers in an exceedingly cold-start setting.

In our current work, the review content isn't considered. within the future, we are going to explore effective ways in which in incorporating review content into our premature reviewer prediction model. Also, we've not studied the communication channel and social network structure in diffusion of innovations partially thanks to the issue in getting the relevant data from our review knowledge. we are going to look into alternative sources of information like Flixster within which social networks is extracted and do additional perceptive analysis. Currently, we tend to specialize in the analysis and prediction of early reviewers, whereas there remains a very important issue to address, i.e., a way to improve product promoting with the identified early reviewers. we are going to investigate this task with real e-commerce cases together with e-commerce companies within the future.

## REFERENCES:

[1] R. Peres, E. Muller, and V. Mahajan, "Innovation diffusion and new product growth models: A critical review and research directions," International Journal of Research in Marketing, vol. 27, no. 2, pp. 91 – 106, 2010.

[2] L. A. Fourt and J. W. Woodlock, "Early prediction of market success for new grocery products." Journal of Marketing, vol. 25, no. 2, pp. 31 – 38, 1960.

[3] Jegadeesan,R., Sankar Ram , and J.Abirmi "Implementing  Online Driving License Renewal by Integration of Web Orchestration and Web Choreogrphy" International journal of Advanced Research trends in Engineering and Technology (IJARTET) ISSN:2394-3785 (Volume-5, Issue-1, January  2018

[4] B. W. O, "Reference group influence on product and brand purchase decisions," Journal of Consumer Research, vol. 9, pp. 183–194, 1982.

[5] J. J. McAuley,C. Targett, Q. Shi, and A. van den Hengel, "Imagebasedrecommendationsonstylesandsubstitutes,"inSIGIR,2015, pp. 43–52.

[6] Jegadeesan,R., Sankar Ram,N. "Energy-Efficient Wireless Network   Communication with Priority Packet Based QoS Scheduling", Asian Journal of Information Technology(AJIT) 15(8): 1396-1404,2016 ISSN: 1682-3915,Medwell Journal,2016

[7] Jegadeesan,R.,Sankar Ram,N. "Energy Consumption Power Aware Data Delivery in Wireless Network", Circuits and Systems, Scientific Research Publisher,2016

[8] E.M.Rogers,Diffusion of Innovations. New York: The Rise of High Technology Culture, 1983.

[9] K. Sarkar and H. Sundaram, "How do we find early adopters who will guide a resource constrained network towards a desired distribution of behaviors?" in CoRR, 2013, p. 1303.

[10] D. Imamori and K. Tajima, "Predicting popularity of twitter accounts through the discovery of link-propagating early adopters," in CoRR, 2015, p. 1512.

[11] Jegadeesan,R., Sankar Ram "Defending Wireless Sensor Network using Randomized Routing "International Journal of Advanced Research in Computer Science and Software Engineering Volume 5, Issue 9, September 2015 ISSN: 2277 128X  Page | 934-938

[12] I. Mele, F. Bonchi, and A. Gionis, "The early-adopter graph and its application to web-page recommendation," in CIKM, 2012, pp. 1682–1686.

[13] Y.-F. Chen, "Herd behavior in purchasing books online," Computers in Human Behavior, vol. 24(5), pp. 1977–1992, 2008.

[14] Banerjee, "A simple model of herd behaviour," Quarterly Journal of Economics, vol. 107, pp. 797–817, 1992.

[13] Jegadeesan,R.,T.Karpagam, Dr.N.Sankar Ram , "Defending Wireless Network using Randomized Routing Process" International journal of Emerging Research in management and Technology ISSN: 2278-9359 (Volume-3, Issue-3) March  2014

[14] A. S. E, "Studies of independence and conformity: I. a minority of one against a unanimous majority," Psychological monographs: General and applied, vol. 70(9), p. 1, 1956.

[15] T. Mikolov, K. Chen, G. S. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," in ICLR, 2013.

[16] A. Bordes, N. Usunier, A. Garc´ıa-Dur´an, J. Weston, and O. Yakhnenko, "Translating embeddings for modeling multirelational data," in NIPS, 2013, pp. 2787–2795.

[17] A. S. E, "Studies of independence and conformity: I. a minority of one against a unanimous majority," Psychological monographs: General and applied, vol. 70(9), p. 1, 1956.

[18] Vijayalakshmi, Balika J Chelliah and Jegadeesan,R.,  February-2014 "SUODY-Preserving Privacy in Sharing Data with Multi-Vendor for Dynamic Groups" Global journal of Engineering,Design & Technology. G.J. E.D.T.,Vol.3(1):43-47  (January-February, 2014)  ISSN: 2319 –7293

[19] V. G. D. W. Shih-Lun Tseng, Shuya Lu, "The effect of herding behavior on online review voting participation," in AMCIS, 2017.

[20] S. M. Mudambi and D. Schuff, "What makes a helpful online review? a study of customer reviews on amazon.com," in MIS Quarterly, 2010, pp. 185–200.

[21] J. J. McAuley, R. Pandey, and J. Leskovec, "Inferring networks of substitutable and complementary products." in KDD, 2015, pp. 785–794.

[22] E. Gilbert and K. Karahalios, "Understanding deja reviewers." in CSCW, 2010, pp. 225–228.

[23] E.-P. Lim, V.-A. Nguyen, N. Jindal, B. Liu, and H. W. Lauw, "Detecting product review spammers using rating behaviors," in CIKM, 2010, pp. 939–948.

[24] C. Wang and D. M. Blei, "Collaborative topic modeling for recommending scientific articles," in SIGKDD, 2011, pp. 448–456.

[25] Jegadeesan,R., Sankar Ram, M.S.Tharani  (September-October, 2013) "Enhancing File Security by Integrating Steganography Technique in Linux Kernel"  Global journal of Engineering,Design & Technology  G.J. E.D.T., Vol. 2(5): Page No:9-14  ISSN: 2319 – 7293

[26] R. Herbrich, T. Minka, and T. Graepel, "Trueskill: A bayesian skill rating system," in NIPS, 2006, pp. 569–576.

[27] J. Liu, Y.-I. Song, and C.-Y. Lin, "Competition-based user expertise score estimation," in SIGIR, 2011, pp. 425–434.

[28] Q. V. Le and T. Mikolov, "Distributed representations of sentences and documents," in ICML, 2014, pp. 1188–1196.

[29] Jegadeesan,R., Sankar Ram   October -2013 "ENROUTING TECHNICS USING DYNAMIC WIRELESS NETWORKS" International Journal of Asia Pacific Journal of Research Ph.D Research Scholar 1, Supervisor2,  VOL -3  Page No: Print-ISSN-2320-5504   impact factor 0.433

[30] Y. B. Xavier Glorot, "Understanding the difficulty of training deep feed forward neural networks," in AISTATS, 2010, pp. 249–256.

[31] S. Chen and T. Joachims, "Modeling intransitivity in matchup and comparison data," in WSDM, 2016, pp. 227–236.

[32] N. Meade and T. Islam, "Modelling and forecasting the diffusion of innovation a25-yearreview,"InternationalJournalofForecasting, vol. 22, no. 3, pp. 519 – 545, 2006.

[33] R.D.Luce, "Individual choice behavior a theoretical analysis," in john Wiley and Sons, 1959.

[34] L.L.Thurstone, "A law of comparative judgment," Psychological review, vol. 34, no. 4, p. 273, 1927.

[35] M. Cattelan, "Models for paired comparison data: A review with emphasis on dependent data," Statistical Science, vol. 27, no. 3, pp. 412–433, 2012.

[36] .Jegadeesan,R.,Sankar Ram M.Naveen Kumar  JAN 2013  "Less Cost Any Routing With Energy Cost  Optimization" International Journal of Advanced Research in Computer Networking,Wireless and Mobile Communications.Volume-No.1:  Page no: Issue-No.1  Impact Factor = 1.5