

Quantifying Herding Effects in Crowd Wisdom

Ting Wang Dashun Wang Fei Wang
IBM T.J. Watson Research Center
Yorktown Heights, NY
{tingwang, dashun, fwang}@us.ibm.com

ABSTRACT

In many diverse settings, aggregated opinions of others play an increasingly dominant role in shaping individual decision making. One key prerequisite of harnessing the “crowd wisdom” is the independency of individuals’ opinions, yet in real settings collective opinions are rarely simple aggregations of independent minds. Recent experimental studies document that disclosing prior collective opinions distorts individuals’ decision making as well as their perceptions of quality and value, highlighting a fundamental disconnect from current modeling efforts: How to model social influence and its impact on systems that are constantly evolving? In this paper, we develop a mechanistic framework to model social influence of prior collective opinions (e.g., online product ratings) on subsequent individual decision making. We find our method successfully captures the dynamics of rating growth, helping us separate social influence bias from inherent values. Using large-scale longitudinal customer rating datasets, we demonstrate that our model not only effectively assesses social influence bias, but also accurately predicts long-term cumulative growth of ratings solely based on early rating trajectories. We believe our framework will play an increasingly important role as our understanding of social processes deepens. It promotes strategies to untangle manipulations and social biases and provides insights towards a more reliable and effective design of social platforms.

Categories and Subject Descriptors

H.2.8 [Database Management]: Database Applications—*Data mining*; J.4 [Social and Behavior Sciences]: Sociology

General Terms

Algorithms, Experimentation

Keywords

Crowd wisdom; social influence; herding effect

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
KDD’14, August 24–27, 2014, New York, NY, USA.
Copyright 2014 ACM 978-1-4503-2956-9/14/08 ...\$15.00.
<http://dx.doi.org/10.1145/2623330.2623720>.

1. INTRODUCTION

With the explosive growth of information, our decisions are increasingly relying on aggregated opinions contributed by others, with the belief that the aggregations over a large population can successfully harness the “wisdom of crowds” [22]. Indeed, rooting back to Galton [8], many studies have shown that collective opinions of a group are often closer to the truth than the answer of an individual to a question. While the crowd wisdom applies usefully to a spectrum of domains, ranging from product or service recommendation [10] and crowdsourcing [5, 20, 26] to stock markets and political elections [22], one key prerequisite of harnessing the crowd wisdom is the independency of individuals’ opinions [22]. Indeed, most if not all of the times, individuals are exposed to others’ opinions before forming and expressing their own. As concrete examples, we go to the theater after checking reviews of the movies online; we download songs from the hit list; we purchase products or go to restaurants after researching what others think about them. As a result, the market does not simply aggregate pre-existing individual preferences, but rather creates an environment rich in social influence.

Thanks to the availability of Web-based experiments, recent studies offered convincing evidence that social influence exerts important but counterintuitive effects on collective judgement [16, 19]. Indeed, through carefully designed control experiments in different settings, these studies demonstrate that disclosing prior collective opinions distorts individuals’ decision making as well as their perceptions of quality and value, creating herding effects that are irrational and pervasive, yet consequential to market outcome. Despite the significance of these results in experimental settings, there has been no quantitative framework to model social influence and its impact on systems that are constantly evolving. Indeed, models on collective intelligence, from majority voting to collaborating filtering to crowdsourcing [5], all assume independent crowds, representing a critical gap between modeling frameworks and empirical insights.

Here we develop a mechanistic framework to model social influence of prior collective opinions (e.g., product ratings) on subsequent individual decisions, namely, Herding Effect Aware Rating Dynamics Model (HEARD). Using 28 million ratings spanning over 18 years on over 1.7 million products from *Amazon* [15] as an exemplary case, we demonstrate that our method successfully captures the dynamics of rating growth across different product categories, allowing us to separate social biases introduced by prior ratings from the true values inherent to products. We further show that,

comparing with competing methods, our framework not only effectively detects the presence of social biases and gauges less biased values for any given product, but also accurately predicts the long-term cumulative growth of ratings through a scalable estimation model solely based on early rating trajectories. As a result, HEARD can also make testable predictions of collective response to artificial manipulations in rating systems, assisting in further testings through more systematic experiments.

To the best of our knowledge, this work represents one of the first few quantitative framework to model social influence biases introduced from prior opinions. We believe our method is of fundamental importance to studies of social processes, promotes new strategies in untangling manipulations and biases within social environments, and provides significant insights towards design of platforms that aggregate individual opinions, from electoral polling to market analysis to product recommendation.

The remainder of the paper will proceed as follows. Section 2 surveys relevant literature. Section 3 details the model design of HEARD and develops efficient inference algorithms to fit the model. Section 4 presents a scalable algorithm to predict the future rating growth based on HEARD. Section 5 empirically evaluates the proposed models and algorithms. The paper is concluded in Section 6.

2. RELATED WORK

In this section, we review three categories of related work, namely, social network induced influence, measuring social influence in experimental settings, and effect of semantics of prior opinions.

Social networks have attracted significant interest, partly due to the availability of large datasets in many domains. One active line of enquiry in social network studies is how behavior [1, 2], opinion [7], and information [13, 24, 3] spreads through social networks. It is conceivable that microscopic social interactions could induce influence that is visible on an aggregated level [25]. In a way, the process of generating collective opinions is similar to consensus formation [11]. For example, individuals may change their opinions after learning about what their friends think. This is supported by experimental results by Lorenz et al. [14], in which they demonstrated that even mild social interactions can significantly bias simple estimation tasks. Therefore, Das et al. [4] proposed a social sampling method that takes into account individuals’ influence from their social neighbors and arrives at a de-biased estimation of collective opinions. While this line of research shows that social interactions can exert influence on overall outcome, their focus on networks inevitably distinguishes themselves from our work. Indeed, often times, the population responsible for collective opinions are not interactive. You choose a restaurant, go watch a movie, or purchase a book, because of the opinions or reviews authored by people you do not know. Therefore, our work focuses on how to model social influence on a macroscopic level and hence predict the outcome of crowd wisdom.

On the other hand, there have been a number of experimental studies on measuring social influence within a population, thanks to the emergence of Web-based experiments. For example, Salganik et al. [19] implemented a music lab, where individuals download and rate songs with or without information about how good the songs are, and they demonstrated that increasing social influence could result in

differential outcomes for songs of similar quality. Muchnik et al. [16] ran a large-scale randomized experiment on a reddit like website, finding that prior ratings created significant bias in individual rating behavior, from turnout to binary choices. These studies confirmed experimentally that disclosing prior ratings can create strong herding effects that are irrational and pervasive, leading to significant bias that is consequential to collective outcome. At the same time, they also highlight a fundamental gap between experimental insights and modeling efforts. Our work directly addresses this gap: To the best of our knowledge, this work is among the first few attempt to quantitatively model the herding effects in crowd wisdom and develop effective mechanisms to factor out such bias in estimation.

Finally, there have been a number of interesting studies into the semantics of collective opinions, such as that analyze the text and social aspects of product reviews [10, 15, 21, 9]. While they are useful for review spam detection, customer sentiment analysis, product recommendation, and more, insights extracted from semantic features are, however, not mechanistic, hence not capable of projecting the full rating trajectories. Nevertheless, these studies are complementary to our work, in a sense that the useful semantic features learned can be integrated into our model in forms of prior belief of model parameters. Indeed, one shall see in next sections that incorporating such text and social information into the rating growth model would be a promising future direction.

Our work also draws connections to other modeling efforts. The design of our herding effects model is inspired by the multi-neuron coupled spiking phenomena [17]. The exponential additive generative mechanism has been applied in modeling latent topics for text [6]. Our work differs in proposing a more general form of generative model and developing scalable inference algorithms to fit the model.

3. MODEL AND ALGORITHM

In this section, we detail the design of the HEARD model and present efficient inference algorithms to fit the model. Concretely, we draw an analogy to the coupled spiking phenomena in a multi-neuron system to model the dynamics of rating growth and fit the model parameters using maximum likelihood estimation.

3.1 HEARD Model

Without loss of generality, we consider a discrete K -level rating system, which is extensively used by today’s online retailers; for example, *Amazon* adopts a one-to-five star rating system. Consider the sequence of ratings regarding a specific product, with $r_i \in \{1, 2, \dots, K\}$ being the i -th rating. We assume the first $(i - 1)$ ratings form the *history* for r_i : $\mathbf{x}_i = [x_{i,1}, x_{i,2}, \dots, x_{i,K}]^\top$, where $x_{i,k}$ represents the proportion of level- k ratings among the first $(i - 1)$ ratings. Clearly, $\sum_{k=1}^K x_{i,k} = 1$ for $i > 1$ and \mathbf{x}_1 is an all-zero vector. We intend to model how disclosing such rating history would influence individual rating behavior on r_i .

Intuitively, the generation of a new level- k rating is driven by multiple factors, including: the intrinsic product quality, the occurrence of preceding level- k ratings, and the history of other ratings. We can draw a close analogy to the spiking activities of a multi-neuron system [17]: the response (i.e., spike) generated by a neuron is jointly determined by the stimulus strength and the preceding spikes of this neuron

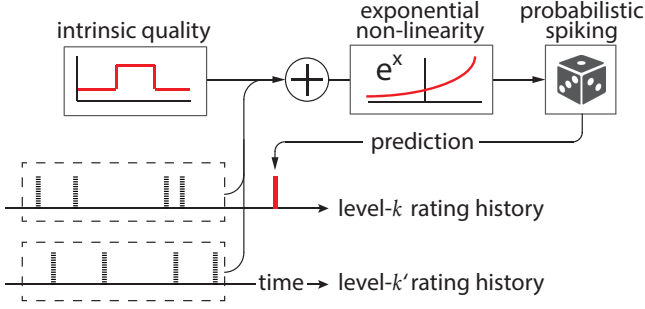


Figure 1: Illustration of HEARD model. The occurrence of a new level- k rating is jointly influenced by (i) the intrinsic quality of product, (ii) the preceding level- k ratings, and (iii) the history of other ratings.

and correlated neurons. We therefore introduce an *additive generative* model to describe the distribution of the i -th rating r_i over different levels:

$$Pr(r_i = k | \mathbf{x}_i) = \frac{\exp(\mu_k + f(i)\boldsymbol{\theta}_k^\top \mathbf{x}_i)}{\sum_{k'=1}^K \exp(\mu_{k'} + f(i)\boldsymbol{\theta}_{k'}^\top \mathbf{x}_i)} \quad (1)$$

This conditional distribution describes the likelihood of observing a level- k rating given rating history \mathbf{x}_i . In this general formulation, we have:

- $\boldsymbol{\mu} = [\mu_1, \mu_2, \dots, \mu_K]^\top \in \mathbb{R}^K$ represents the coefficients of an *intrinsic distribution*, which is assumed related to the true quality of the product.
- $f(\cdot)^1$ is the *magnitude function*, which describes the relationship between the strength of herding effects and the number of historical ratings; in particular, we have $f(1) = 0$.
- $\boldsymbol{\theta}_k \in \mathbb{R}^K$ weighs the different components of \mathbf{x}_i . Note that our model captures both positive and negative influence. Concretely, when the k' -th component $\theta_{k,k'} > 0$, the preceding level- k' ratings *excite* the occurrence of level- k ratings; while if $\theta_{k,k'} < 0$, the level- k' ratings *inhibit* the generation of new level- k ratings.

These factors are then integrated in an exponential function, as illustrated in Figure 1.

Note that here we ignore the time dimension in our model because various external factors may abruptly influence the temporal dynamics of rating growth, e.g., low price promotion, emergence of new products, advertisements, etc.

Let $\boldsymbol{\Theta} = [\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_K, \mathbf{u}]$ represent all the parameters. Both $\boldsymbol{\Theta}$ and magnitude function $f(\cdot)$ are estimated from data; in particular, $f(\cdot)$ is estimated from an infinite dimensional functional space. Next we elaborate their inference.

3.2 Model Inference

We assume regarding a specific product, a temporally ordered sequence of N ratings $\{r_i\}_{i=1}^N$ has been observed. Note that while we focus on the case of a single product for ease of presentation, the extension to multiple products is straightforward. For notational simplicity, we introduce a set of indicator variables $\mathbf{y}_i \in \{0, 1\}^K$ with $y_{i,k} = 1$ if $r^{(i)} = k$ and

¹In the following, we use f_i as a short notation of $f(i)$.

0 otherwise. Then the log-likelihood of parameters $\boldsymbol{\Theta}$ given this rating sequence is expressed as:

$$\begin{aligned} \mathcal{L}(\boldsymbol{\Theta}) &= \frac{1}{N} \log \prod_{i=1}^N Pr(r_i | \mathbf{x}_i, \boldsymbol{\Theta}) \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K y_{i,k} \log \frac{\exp(\mu_k + f_i \boldsymbol{\theta}_k^\top \mathbf{x}_i)}{\sum_{k'=1}^K \exp(\mu_{k'} + f_i \boldsymbol{\theta}_{k'}^\top \mathbf{x}_i)} \end{aligned}$$

We estimate the model parameters by minimizing the penalized log-likelihood function, which is defined as:

$$\mathcal{L}_\lambda(\boldsymbol{\Theta}) = -\mathcal{L}(\boldsymbol{\Theta}) + \frac{\lambda}{2} (\|\boldsymbol{\Theta}\|_F^2 + \mathcal{R}(f)) \quad (2)$$

where the first term represents the negative log-likelihood, the second term is a regularizer with λ being the balance parameter to prevent overfitting, and $\|\cdot\|_F$ denotes the matrix Frobenius norm. In particular, $\mathcal{R}(f)$ is a penalty term preferring smooth functions. Without prior knowledge, we use $\mathcal{R}(f) = \int_0^\infty (f'(t))^2 dt$, where $f'(\cdot)$ represents the derivative of $f(\cdot)$.

While $\mathcal{L}_\lambda(\boldsymbol{\Theta})$ appears similar to the softmax regression; it contains the integral of an unknown function and meanwhile all the parameters are coupled, which makes it difficult to directly apply off-the-shelf optimization methods (e.g., coordinate descent). Next we propose an iterative algorithm which optimizes $\mathcal{L}_\lambda(\boldsymbol{\Theta})$ by (i) constructing a surrogate function to decouple the parameters and (ii) applying an Euler-Lagrange equation to fit the unknown function.

More specifically, let $\boldsymbol{\Theta}^{(n)} = [\boldsymbol{\theta}_1^{(n)}, \boldsymbol{\theta}_2^{(n)}, \dots, \boldsymbol{\theta}_K^{(n)}, \boldsymbol{\mu}^{(n)}]$ denote the current parameter setting. We construct the following surrogate function $Q(\boldsymbol{\Theta}; \boldsymbol{\Theta}^{(n)})$, which is a tight upper bound of $\mathcal{L}_\lambda(\boldsymbol{\Theta})$:

$$\begin{aligned} Q(\boldsymbol{\Theta}; \boldsymbol{\Theta}^{(n)}) &= \frac{1}{N} \sum_i \sum_k \left(\phi_{i,k}^2 + \left(\beta_{i,k}^{(n)} - 2\phi_{i,k}^{(n)} - y_{i,k} \right) \phi_{i,k} \right) \\ &\quad - \frac{1}{NK} \sum_i \left(\sum_k \phi_{i,k} - 2 \sum_k \phi_{i,k}^{(n)} \right) \left(\sum_k \phi_{i,k} \right) \\ &\quad + \frac{\lambda}{2} (\|\boldsymbol{\Theta}\|_F^2 + \mathcal{R}(f)) + \frac{1}{N} \sum_i C_i^{(n)} \end{aligned} \quad (3)$$

where the terms $\phi_{i,k}$, $\phi_{i,k}^{(n)}$, $\beta_{i,k}^{(n)}$ and $C_i^{(n)}$ are defined below:

$$\begin{aligned} \phi_{i,k} &= \mu_k + f_i \boldsymbol{\theta}_k^\top \mathbf{x}_i \\ \phi_{i,k}^{(n)} &= \mu_k^{(n)} + f_i^{(n)} \boldsymbol{\theta}_k^{(n)\top} \mathbf{x}_i \\ \beta_{i,k}^{(n)} &= \frac{\exp(\phi_{i,k}^{(n)})}{\sum_{k'} \exp(\phi_{i,k'}^{(n)})} \\ C_i^{(n)} &= \sum_k \left(\phi_{i,k}^{(n)2} - \beta_{i,k}^{(n)} \phi_{i,k}^{(n)} \right) - \frac{1}{K} \left(\sum_k \phi_{i,k}^{(n)} \right)^2 \\ &\quad + \log \sum_k \exp(\phi_{i,k}^{(n)}) \end{aligned}$$

It is noted that $Q(\boldsymbol{\Theta}; \boldsymbol{\Theta}^{(n)})$ possesses the following desirable properties (details in Appendix):

$$\begin{cases} Q(\boldsymbol{\Theta}; \boldsymbol{\Theta}^{(n)}) \geq \mathcal{L}_\lambda(\boldsymbol{\Theta}) & \forall \boldsymbol{\Theta}, \boldsymbol{\Theta}^{(n)} \\ Q(\boldsymbol{\Theta}; \boldsymbol{\Theta}^{(n)}) = \mathcal{L}_\lambda(\boldsymbol{\Theta}^{(n)}) & \forall \boldsymbol{\Theta}^{(n)} \end{cases}$$

which imply that if $\boldsymbol{\Theta}^{(n+1)} = \arg \min_{\boldsymbol{\Theta}} Q(\boldsymbol{\Theta}; \boldsymbol{\Theta}^{(n)})$, then we must have $\mathcal{L}_\lambda(\boldsymbol{\Theta}^{(n)}) \geq \mathcal{L}_\lambda(\boldsymbol{\Theta}^{(n+1)})$. Therefore, minimizing

$Q(\Theta; \Theta^{(n)})$ with respect to Θ at each iteration will ensure that $\mathcal{L}_\lambda(\Theta)$ decreases monotonically.

Updating Parameters

The formulation above has the advantage that we can derive the closed form solution of Θ for $\arg \min_{\Theta} Q(\Theta; \Theta^{(n)})$. Specifically, by deriving the derivatives of $Q(\Theta; \Theta^{(n)})$ with respect to μ_k and $\theta_{k,k'}$ and set them to zero, we obtain their update rules as follows:

$$\mu_k^{(n+1)} = \frac{K \sum_i (y_{i,k} - \beta_{i,k}^{(n)}) + 2N(K-1)\mu_k^{(n)}}{2N(K-1) + NK\lambda} \quad (4)$$

$$\theta_{k,k'}^{(n+1)} = \frac{K \sum_i f_i x_{i,k'} (y_{i,k} - \beta_{i,k}^{(n)}) + 2(K-1) \sum_i f_i^2 x_{i,k'}^2 \theta_{k,k'}^{(n)}}{2(K-1) \sum_i f_i^2 x_{i,k'}^2 + NK\lambda} \quad (5)$$

Updating Magnitude Function

Next we derive the update rule for magnitude function $f(\cdot)$ by optimizing it in an infinite dimensional functional space. We extract the parts of $Q(\Theta; \Theta^{(n)})$ relevant to $f(\cdot)$ and then reformulate the problem of minimizing $Q(\Theta; \Theta^{(n)})$ with respect to $f(\cdot)$ as follows:

$$\min_{f \in L_1(\mathbb{R})} \sum_i A_i f_i^2 + \sum_i B_i f_i + \frac{\lambda}{2} \int_0^{+\infty} (f'(t))^2 dt \quad (6)$$

where terms A_i and B_i are defined below:

$$\begin{aligned} A_i &= \frac{1}{N} \sum_k (\theta_k^{(n)\top} \mathbf{x}_i)^2 - \frac{1}{NK} \left(\sum_k \theta_k^{(n)\top} \mathbf{x}_i \right)^2 \\ B_i &= \frac{1}{N} \sum_k (2\mu_k^{(n)} - 2\phi_{i,k}^{(n)} + \beta_{i,k}^{(n)} - y_{i,k}) \theta_k^{(n)\top} \mathbf{x}_i \\ &\quad + \frac{2}{NK} \left(\sum_k \theta_k^{(n)\top} \mathbf{x}_i \right) \left(\sum_k \phi_{i,k}^{(n)} - \sum_k \mu_k^{(n)} \right) \end{aligned}$$

Abusing the notation a little, we introduce two functions: $A(t) = A_i \mathbb{I}\{t \leq N \wedge t \in \mathbb{N}\}$ and $B(t) = B_i \mathbb{I}\{t \leq N \wedge t \in \mathbb{N}\}$, where $\mathbb{I}\{\cdot\}$ is the indicator function which returns 1 if the predicate is true and 0 otherwise.

Then the solution of the objective function as defined in Eqn.(6) must satisfy the Euler-Lagrange equation [27] (proof referred to Appendix):

$$2A(t)f(t) + B(t) - \lambda f''(t) = 0 \quad (7)$$

where $g''(\cdot)$ is the second order derivate of $g(\cdot)$.

Due to the discrete nature of the functions $A(t)$ and $B(t)$, we solve this differential equation numerically using a Seidal type iteration. Specifically, we discretize the differential equations over intervals of length 1:

$$\lambda(f_{i+1} - 2f_i + f_{i-1}) - 2A_i f_i - B_i = 0$$

Clearly from the equations above we can efficiently solve f_i for $i = 1, 2, \dots, N$. We may then perform curve fitting to extrapolate the values of f_i for $i > N$.

Complete Algorithm

To set a proper starting point for optimization, we consider the degenerated case where the prior ratings have no effect

on individual rating behaviors. Under this assumption, we have the following setting:

$$\begin{cases} \mu_k = \log \left(\frac{\sum_{i=1}^N y_{i,k}}{N} \right) & k = 1, 2, \dots, K \\ f_i = 0 & i = 1, 2, \dots, N \end{cases} \quad (8)$$

Meanwhile we initialize $\theta_1, \theta_2, \dots, \theta_K$ randomly.

Putting everything together, Algorithm 1 sketches the procedure of model inference. After initialization, it iterates between updating parameters Θ and solving magnitude function $f(\cdot)$ until the objective function converges.

Algorithm 1: Inference of HEARD Model

```

Input: rating history  $\{r_i\}_{i=1}^N$ 
Output: setting of parameters  $\Theta$  and function  $f$ 
// initialization
initialize  $\Theta$  and  $f$  according to Eqn.(8);
compute statistics  $\{x_i\}_{i=1}^N$ ;
// iterative optimization
while not converged yet do
  // update parameter
  for  $k = 1, 2, \dots, K$  do
    update  $\mu_k$  following Eqn.(4);
    for  $k' = 1, 2, \dots, K$  do
      update  $\theta_{k,k'}$  following Eqn.(5);
  // update magnitude function
  compute  $\{f_i\}_i$  by solving differential Eqn.(7);
return setting of  $\Theta$  and  $f$ ;

```

4. APPLICATION

In this section we show that equipped with the aforementioned HEARD model, we are able to answer a set of fundamental questions, including:

Debiasing: What is the intrinsic quality of a product after factoring out the herding effects from its collective ratings?

Prediction: Based on its rating history, can we predict the distribution of its next 100 ratings?

What-If Analysis: Given its current ratings, how would its future ratings be “herded” if we could “inject” in 50 five-star ratings?

Next we detail how to answer these questions.

4.1 Debiasing Collective Ratings

To the first question, recall that the HEARD model defined in Eqn.(1) comprises two additive components, namely, the intrinsic distribution and the herding effect distributions. The background intrinsic distribution as controlled by parameters $\{\mu_k\}$ is assumed related to the true quality of a product. Therefore, once we have estimated $\{\mu_k\}$ from the rating history of a product, we can then “debias” the collective ratings by factoring out the components attributed to the herding effects.

More concretely, abusing the notation a little, let $\boldsymbol{\mu} = [\mu_1, \mu_2, \dots, \mu_K]^\top$. Without the herding effects, each rating

is generated by the following unconditional categorical distribution:

$$\boldsymbol{\eta} = \frac{\exp(\boldsymbol{\mu})}{\sum_k \exp(\mu_k)} \quad (9)$$

which represents the intrinsic rating of the product.

The straightforward solution to estimating $\boldsymbol{\mu}$ of a given product is to directly fit the model parameters using its rating history as in Algorithm 1, which however may lead to overfitting. Instead, we introduce an “out-of-sample” extension. As will be detailed in Section 5, the herding effects often follow similar patterns for products of the same category (e.g., books). We therefore use the rating histories of a bulk of products in the same category to train *category-level* parameters $\{\boldsymbol{\theta}_k\}$ and magnitude function $f(\cdot)$. For the query product, we fix $\{\boldsymbol{\theta}_k\}$ and $f(\cdot)$ and focus on learning *product-level* parameter $\boldsymbol{\mu}$. As shown in Algorithm 2, this procedure is similar to Algorithm 1, except that at each iteration we only need to update $\boldsymbol{\mu}$.

Algorithm 2: Out of Sample Extension

Input: rating history $\{r_i\}_i$ of query product, setting of $\{\boldsymbol{\theta}_k\}_k$ and $f(\cdot)$
Output: setting of $\boldsymbol{\mu}$ for given product
// initialization
initialize $\boldsymbol{\mu}$ according to Eqn.(8);
compute statistics $\{\mathbf{x}_i\}_i$;
// iterative optimization
while not converged **yet do**
 // update parameter
 for $k = 1, 2, \dots, K$ **do**
 | update μ_k following Eqn.(4);
return $\boldsymbol{\mu}$;

4.2 Predicting Rating Growth

Another interesting question one may pose is: given the current rating history of a product, is it possible to predict the distribution of its future ratings? While it is of theoretical interest to discuss the statistical convergence properties of the rating distribution as the number of ratings approaches infinity; in real settings, most products during their lifetimes receive only limited number of ratings. We thus focus on a more concrete question as follows:

Given the first N ratings of a product, can we characterize the distribution of its next M ratings?

Let us first consider the herding effects-agnostic case, in which each rating is independently generated by the categorical distribution as defined in Eqn.(9). Under this assumption, the next M ratings follow a multinomial distribution; specifically, the expected number of level- k rating is given by $M\eta_k$ with variance $M\eta_k(1 - \eta_k)$.

Next we incorporate the herding effects. Recall that the distribution of the first $(i - 1)$ ratings is given by \mathbf{x}_i , which also corresponds to the history for the i -th rating. The transition probability from \mathbf{x}_i to \mathbf{x}_{i+1} can be described as below:

$$Pr\left(\mathbf{x}_{i+1} = \frac{i-1}{i}\mathbf{x}_i + \frac{\mathbf{e}_k}{i} \mid \mathbf{x}_i\right) = \frac{\exp(\phi_{i,k})}{\sum_{k'} \exp(\phi_{i,k'})} \quad (10)$$

where \mathbf{e}_k is a 1-of- K vector with the k -th element being 1.

Note that this transition rule essentially specifies a non-stationary Markov chain in which both the state space and the transition probability change from step to step. This setting is not amenable to exact inference; we thus resort to Monte Carlo methods [18].

Algorithm 3 sketches our prediction model. Starting with current rating distribution \mathbf{x}_{N+1} estimated from the given rating history, it iteratively samples the next rating distribution using the transition rule in Eqn.(10). Let $\{\mathbf{x}_{N+M+1}^{(i)}\}_{i=1}^L$ be the set of samples of target distribution \mathbf{x}_{N+M+1} and $\hat{\mathbf{x}}_{N+M+1} = \frac{1}{L} \sum_{i=1}^L \mathbf{x}_{N+M+1}^{(i)}$ be the expectation of target distribution. We can prove that for given thresholds ϵ and δ , if the sample size L satisfies the following condition:

$$L \geq \lfloor \frac{1}{2\epsilon^2} \log \frac{2}{\delta} \rfloor \quad (11)$$

then $|\hat{\mathbf{x}}_{N+M+1} - \mathbf{x}_{N+M+1}| \leq \epsilon \mathbf{1}$ with probability at least $1 - \delta$, where $\mathbf{1}$ denotes a K -dimensional all-ones vector (details given in Appendix).

It is also noted that Algorithm 3 features the complexity of $O(\frac{1}{\epsilon^2} \log(\frac{1}{\delta})MK)$, thereby scaling up to large M .

Algorithm 3: Prediction of Rating Growth

Input: rating history $\{r_i\}_{i=1}^N$, prediction range M , error threshold ϵ
Output: estimation of rating distribution \mathbf{x}_{N+M+1}
// initialization
estimate Θ and f by Algorithm 1;
compute required sample size L by Eqn.(11);
compute \mathbf{x}_{N+1} from $\{r_i\}_{i=1}^N$;
// random sampling
for $i = 1, 2, \dots, L$ **do**
 $\mathbf{x}_{N+1}^{(i)} \leftarrow \mathbf{x}_{N+1}$;
 for $j = 1, 2, \dots, M$ **do**
 | sample $\mathbf{x}_{N+1+j}^{(i)}$ according to Eqn.(10);
 store $\mathbf{x}_{N+M+1}^{(i)}$;
estimate $\mathbb{E}[\mathbf{x}_{N+M+1}]$ by $\frac{1}{L} \sum_{i=1}^L \mathbf{x}_{N+M+1}^{(i)}$;

4.3 What-If Analysis

The Markovian nature of the HEARD model also enables us to perform the “what-if” analysis. Concretely, given the current rating distribution \mathbf{x}_i , one may arbitrarily change \mathbf{x}_i to another distribution \mathbf{x}'_i to reflect any artificial conditions one wishes to “inject” in (e.g., a burst of 50 five-star ratings due to certain promotion campaigns). Staring from this new state \mathbf{x}'_i and applying the prediction method above, one may then gauge the consequences of the injected conditions by predicting the trends of future rating growth.

Such what-if analysis is especially valuable for a range of applications including market profitability estimation, budgeted advertising, and fraudulent manipulation detection.

5. EVALUATION

In this section we present an empirical evaluation on the efficacy of the proposed models and algorithms.

5.1 Experimental Setting

We start with introducing the datasets and the alternative techniques to be evaluated.

category	# products	# ratings	avg. # ratings	avg. rating	avg. entropy
Books	929,264	12,886,488	13.9	4.271	0.666
Music	556,814	6,396,350	11.5	4.410	0.555
Movies	212,836	7,850,072	36.9	3.944	0.955
Electronics	82,067	1,241,778	15.1	3.791	0.824
Total	1,780,981	28,374,688	15.9	4.253	0.673

Table 1: Summary of Amazon customer rating dataset.

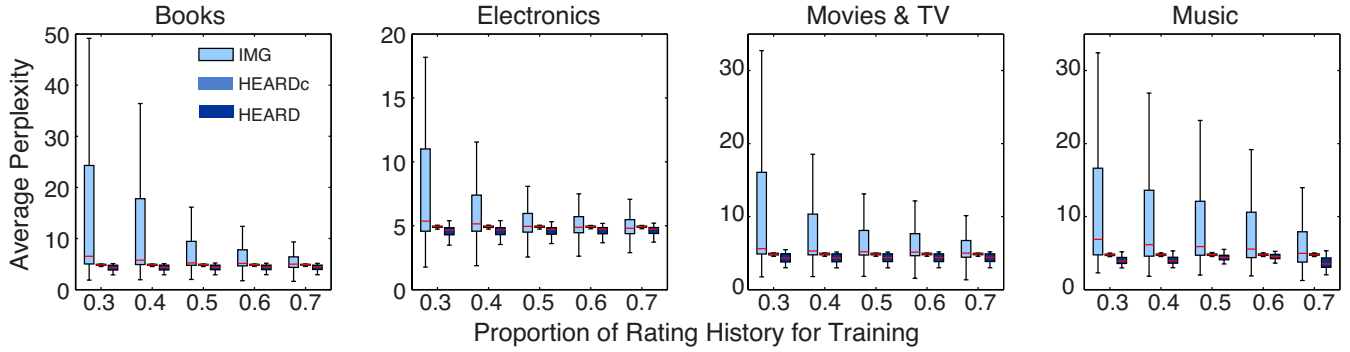


Figure 2: Accuracy of short-term prediction versus the length of rating history used for training.

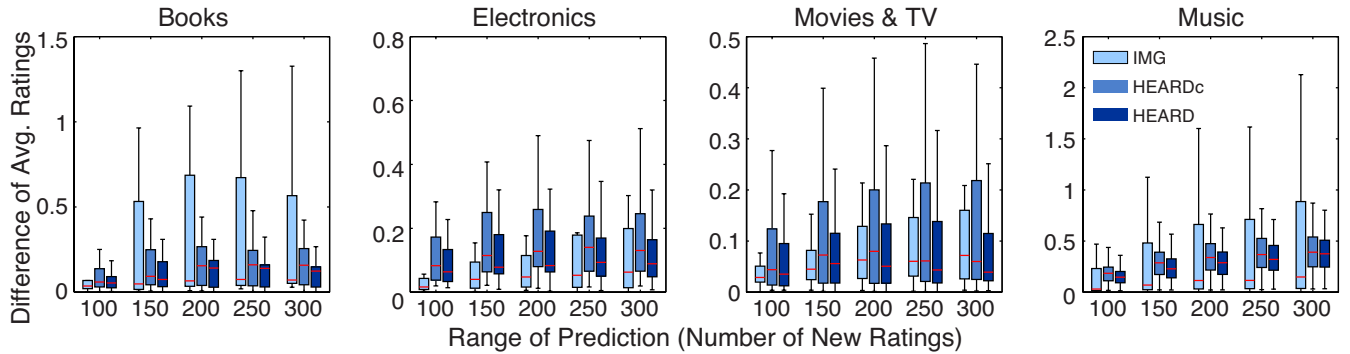


Figure 3: Accuracy of long-term prediction by rating growth models under varying range of prediction.

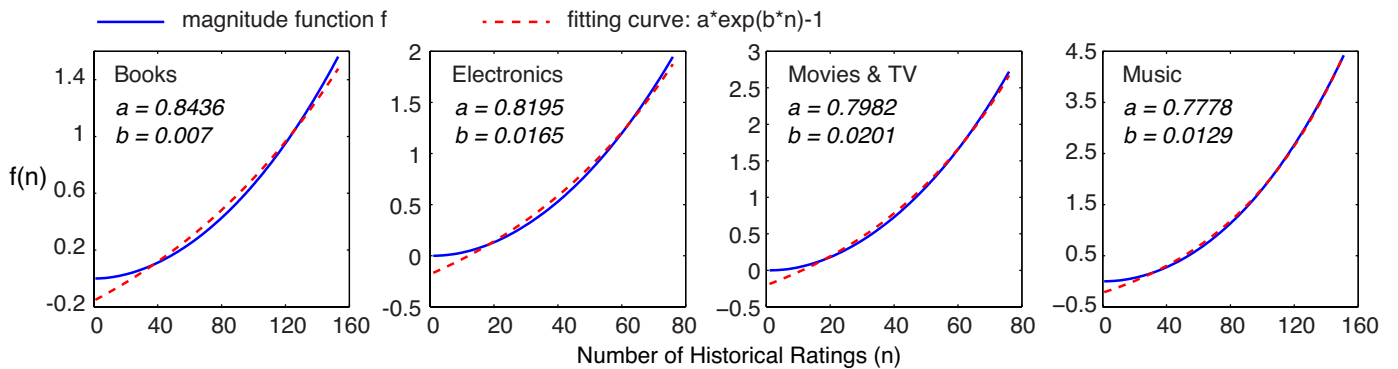


Figure 4: Estimated magnitude function $f(n)$ and fitting curve $a * \exp(b * n) - 1$.

Real Customer Rating Data

We evaluate different models using the real customer rating data collected from *Amazon*², which spans a period of approximately 18 years, including around 35 million ratings regarding about 2.4 million products [15]. In particular, we focus on the products in four major categories: *Books*, *Music*, *Movies & TV*, and *Electronics*, which cover over 72% of the total number of products in the collection. The statistics of this rating dataset is summarized in Table 1. It is noticed that these four categories demonstrate fairly diverse characteristics, for example, with average rating entropy ranging from 0.56 to 0.96.

Alternative Models

For comparison purposes, besides the HEARD model, we implemented two additional rating growth models:

- Independent Multinomial Generative model (IMG). It is the null hypothesis model, which assumes each new rating is generated according to a fixed multinomial distribution over different rating levels. This multinomial model is estimated from the rating history following the maximum likelihood principle.
- Constant HEARD model (HEARD_C). It is a simplified variant of HEARD, which follows the definition of Eqn.(1), except that the magnitude function is set as $f(x) = 1$ for $x > 1$; that is, it assumes the strength of herding effects stays constant regardless of the cumulative number of ratings.

We implemented all the models and associated algorithms in Matlab and conducted the experiments on a Linux box running 3.5GHz Intel i7 CPU and 16GB RAM. The default parameter setting is: $\lambda = 1$, $\delta = 0.05$, and $\epsilon = 0.01$.

5.2 Validating Rating Growth Models

In this first set of experiments, we intend to evaluate the validity of different rating growth models. For each product in the dataset, we partition its temporally ordered sequence of ratings into two subsequences as the training (i.e., rating history) and testing parts respectively. We use the rating history to train the rating growth models and let them predict the “future” ratings in the testing set. We compare their accuracy in both short-term and long-term prediction.

Short-Term Prediction

In short-term prediction, we vary the length of rating history (as the proportion of the entire rating sequence of a product) and measure the average perplexity of the prediction of the next 50 ratings by different models.

The results are shown in Figure 2. It is noticed that across all four product categories, HEARD and HEARD_C outperform IMG in terms of prediction accuracy. In particular, when only limited data (e.g., 30%) is available for training, the accuracy of IMG can be arbitrarily bad. This is attributed to the fact that the prediction of IMG relies on the overall statistics of the rating history of *each* product, which however has not emerged yet at this early stage. In contrast, HEARD leverages the rating histories of *all* the products in the same category to fit the model, thereby achieving high accuracy even when

²We would like to thank J. McAuley and J. Leskovec for sharing the dataset.

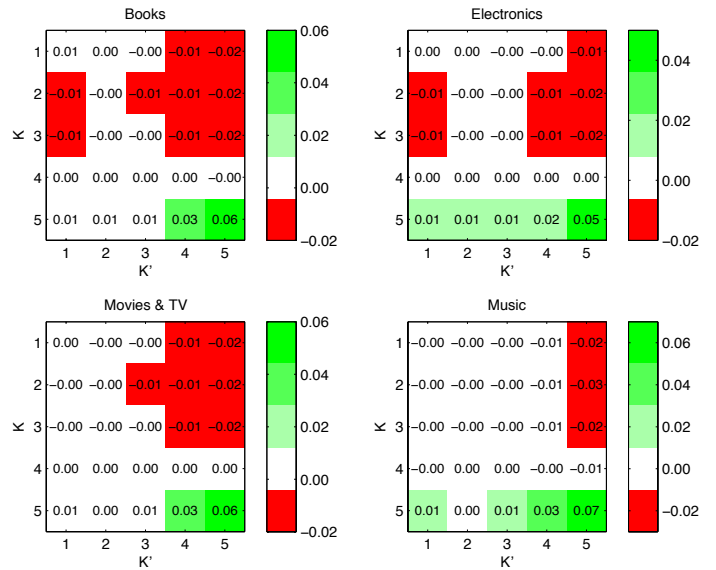


Figure 5: Heat maps of parameters $\{\theta_{k,k'}\}$ for each product category.

facing limited training data. This desirable property makes HEARD especially valuable for early-stage prediction, as we will discuss shortly.

It is also noticed that HEARD achieves higher accuracy than HEARD_C but with marginally larger variance. This is consistent with the fact that HEARD adopts a more complicated model than HEARD_C, which enables HEARD to model a wider range of herding effects but at the cost of slightly higher variance.

Long-Term Prediction

In long-term prediction, we select the products with at least 500 ratings and fix the length of rating history (for training) as 200. We then apply each model to predicting the rating distribution after the next M ratings (M is referred to as the *prediction range*). The accuracy is measured by the difference between predicted and actual average ratings.

The performance of different models is illustrated in Figure 3, wherein we vary prediction range M from 100 to 300. It is observed that compared with IMG and HEARD_C, the prediction accuracy of HEARD is much less sensitive to the setting of M . This can be explained as follows. First, the prediction of IMG depends on the simple statistics (i.e., fraction of ratings at different levels), which however may fluctuate significantly over a large time span. Second, as M increases, the change of the strength of herding effects can no longer be ignored as HEARD_C does.

We can thus conclude that HEARD achieves reliable accuracy in both short-term and long-term prediction tasks, implying that HEARD faithfully captures the growth dynamics of product ratings.

5.3 Characterizing Herding Effect

Next, equipped with the HEARD model as the analytical tool, we conduct a quantitative study on the herding effects observable in real customer rating data. More concretely, for each product category, we apply Algorithm 1 to fitting the model and examine the herding effects as characterized

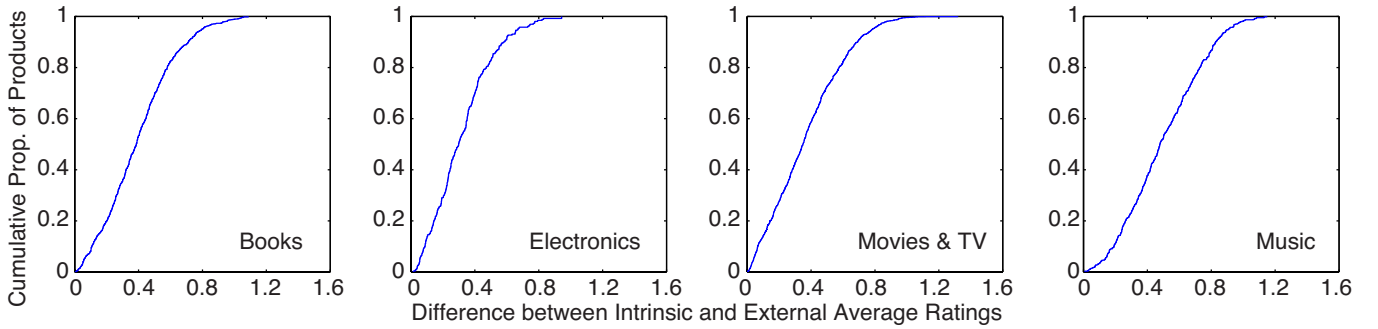


Figure 6: Cumulative proportion of products versus difference between intrinsic and external average ratings.

by the estimated magnitude function $f(\cdot)$ and category-level parameters $\{\theta_k\}_k$.

Strength of Herding Effect

Recall that $f(n)$ specifies the strength of herding effects as a function of the number of historical ratings n . Figure 4 illustrates the estimated $f(n)$ for each product category. We further apply curve fitting to $f(n)$ with an exponential model $a * \exp(b * n) - 1$ (a and b are parameters). Interestingly, the magnitude functions in all four categories tightly follow the exponential curves, despite their slightly different parameter settings of a and b .

This finding entails multi-fold implications: First, it confirms our intuition that the strength of herding effects evolves with the cumulative number of historical ratings. Second, it also echoes the results documented by existing experimental studies (e.g., [19]) on the nonlinear relationship between the predicability of individual behaviors and external influence. Third, most importantly, it provides a formula to explicitly quantify the strength of herding effects. For example, comparing the curves for the categories of *Books* and *Movies & TV*, it is observed that the herding effects is stronger in the category of *Movies & TV*, that is, customers are more easily to be influenced by prior ratings when purchasing *Movies & TV* products. Such information can be valuable for applications such as targeted advertising.

Mutual Influence

Now we proceed to examining parameters $\{\theta_k\}$. Recall that these parameters dictate the mutual influence between the ratings at different levels, concretely, with $\theta_{k,k'}$ specifying how preceding level- k' ratings may positively excite or negatively inhibit the generation of level- k ratings.

Figure 5 illustrates the heat maps of $\{\theta_k\}$ estimated for each product category. While each category has its unique traits, certain common patterns are observed. First, high ratings (e.g., five-star ratings) tend to stimulate new high ratings while inhibiting the generation of low ratings. Second, high ratings are more impactful than low ratings in influencing other ratings. These observations are consistent with the finding of the asymmetric herding effects of positive and negative prior opinions as reported in [16].

5.4 Case Studies

As discussed in Section 4, equipped with HEARD, we are able to perform various analytical tasks. In this set of ex-

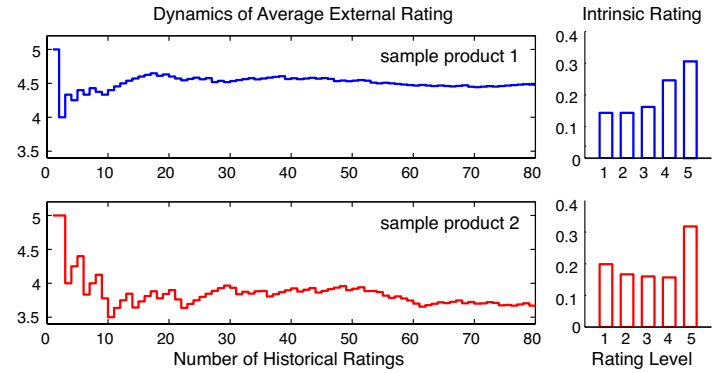


Figure 7: Two sample products with similar intrinsic ratings but with different rating growth histories, leading to significantly distinct external ratings.

periments, we showcase how HEARD helps answer two fundamental questions: (i) exposing the rating inherent to the quality of a product (i.e., “intrinsic rating”) by factoring out the herding effects from collective ratings, and (ii) performing predicative, what-if analysis by incorporating artificial conditions into the rating growth dynamics model.

Debiasing Collective Ratings

To understand the issue that the simple aggregated (or external) rating of a product deviates from its true quality, we apply HEARD to estimate the intrinsic ratings as in Section 4.1 and then measure for each product the difference between its intrinsic and external average ratings.

Figure 6 shows the cumulative proportion of products with respect to the difference between intrinsic and external ratings in each category. It is observed that in all the cases, over 50% products have their external ratings deviate at least 0.5 from their intrinsic ratings, which is significant considering that *Amazon* uses a five-level rating system.

Endowed with the capability of exposing the intrinsic rating of a product, we can then compare the true quality of two products without being misguided by their external ratings. Figure 7 showcases such an example, in which the dynamics of the average external ratings of two sample products is depicted. Despite that they differ significantly in their external ratings (about 0.9), their intrinsic ratings are indeed fairly similar as shown in the right plot. This is explained by

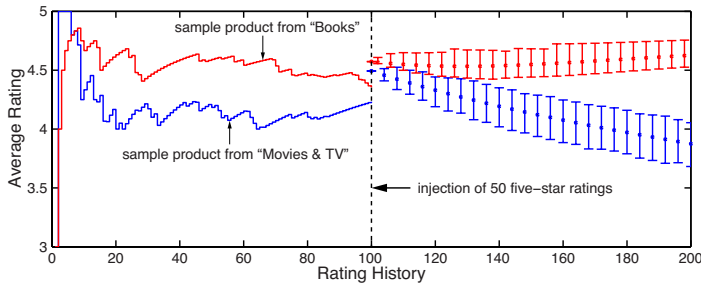


Figure 8: What-if analysis by incorporating artificial conditions into prediction model.

that sample product 2 experiences a sequence of low ratings at the early stage of its history, which considerably changes the dynamics of its rating growth. With the help of HEARD, however, we are able to maximally debias this type of influence caused by the herding effects.

What-If Analysis

As introduced in Section 4.3, the Markovian nature of the HEARD model enables us to perform predicative, “what-if” analysis by artificially incorporating desired conditions into the prediction model and analyze the consequences using simulation. For example, before deciding whether to invest in a promotion campaign for a product, market analysts may wish to estimate the long-term influence of the burst of high ratings due to the promotion.

Figure 8 shows one concrete example. We pick two sample products respectively from the categories of *Movies & TV* and *Books*, which have fairly close average ratings thus far. Now, assuming the promotion takes effect, we inject 50 five-star ratings into their rating histories. As shown in the right panel of Figure 8, the prediction by HEARD tells us: while both products experience similar short-term bursts in their popularity, in the long run the promotion has much longer-lasting influence on the sample product from the category of *Books*. It is clear that this provides valuable information for the decision making of market analysts.

5.5 Scalability

In the last set of experiments, we evaluate the scalability of HEARD. Specifically, for model inference, we measure the average execution time per product (i.e., rating sequence) by HEARD as the length of rating history (for training) varies; meanwhile, for future rating prediction, we measure its average execution time under varying setting of prediction range.

The results are depicted in Figure 9. It is clear that the execution time of HEARD grows approximately linearly with the rating sequence length and prediction range. This also confirms our theoretical analysis on the complexity of Algorithm 1 and Algorithm 3. We can thus conclude that HEARD scales up to large rating datasets.

6. CONCLUSION

This paper presented a quantitative framework to gauge the herding effects in collective opinions of individuals. We proposed HEARD, a mechanistic modeling framework for the growth dynamics of online product ratings, which explicitly accounts for the herding effects of prior customer opinions. Using massive customer rating datasets, we demonstrated

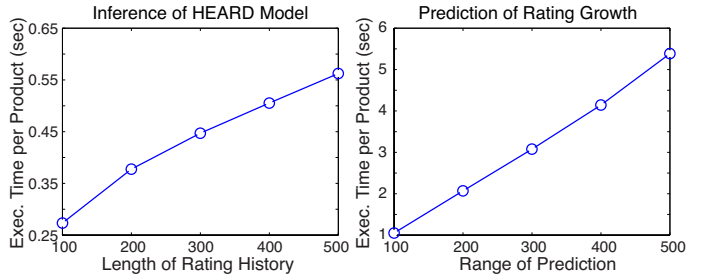


Figure 9: Average execution time per product by HEARD in model inference and rating prediction.

the efficacy of HEARD in capturing the dynamics of rating growth, quantifying social influence and debiasing collective ratings, and further performing what-if analysis against artificial manipulations. HEARD is not limited to product rating systems. Indeed, the mechanistic nature of HEARD makes it applicable for modeling the herding effects in other domains where social influence plays a role, from crowdsourcing and recommender systems to electoral polling.

This work also opens up several directions for future investigations. For example, recent work has shown that the temporal dynamics of collective response to a publication follows rather reproducible patterns, as citations can be captured by a mechanistic model [23]. Hence, incorporating the temporal dynamics in the rating growth model can be fruitful and could potentially shed new light on the nature of crowd wisdom. Furthermore, our framework is orthogonal to studies on the text and social aspects of product reviews and collective opinions, suggesting a rather promising direction by combining the two approaches. Lastly, the model makes falsifiable predictions for collective response against artificial manipulations, making it a viable candidate to assess and guide experimental studies, results of which could feed back to and improve the model with more accurate and realistic predictions.

7. REFERENCES

- [1] N. A. Christakis and J. H. Fowler. The spread of obesity in a large social network over 32 years. *New England journal of medicine*, 357(4):370–379, 2007.
- [2] N. A. Christakis and J. H. Fowler. The collective dynamics of smoking in a large social network. *New England journal of medicine*, 358(21):2249–2258, 2008.
- [3] P. Cui, S. Jin, L. Yu, F. Wang, W. Zhu, and S. Yang. Cascading outbreak prediction in networks: A data-driven approach. In *KDD*, 2013.
- [4] A. Das, S. Gollapudi, R. Panigrahy, and M. Salek. Debiasing social wisdom. In *KDD*, 2013.
- [5] A. P. Dawid and A. M. Skene. Maximum likelihood estimation of observer error-rates using the em algorithm. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 28(1):20–28, 1979.
- [6] J. Eisenstein, A. Ahmed, and E. P. Xing. Sparse additive generative models of text. In *ICML*, 2011.
- [7] J. H. Fowler, N. A. Christakis, Steptoe, and D. Roux. Dynamic spread of happiness in a large social network: longitudinal analysis of the framingham heart study social network. *BMJ: British medical journal*, pages 23–27, 2009.

- [8] F. Galton. Vox populi. *Nature*, 75(7):450, 1907.
- [9] G. Ganu, N. Elhadad, and A. Marian. Beyond the stars: Improving rating predictions using review text content. In *WebDB*, 2009.
- [10] M. Hu and B. Liu. Mining and summarizing customer reviews. In *KDD*, 2004.
- [11] S. Judd, M. Kearns, and Y. Vorobeychik. Behavioral dynamics and influence in networked coloring and consensus. *Proceedings of the National Academy of Sciences*, 107(34):14978–14982, 2010.
- [12] B. Krishnapuram, L. Carin, M. A. T. Figueiredo, and A. J. Hartemink. Sparse multinomial logistic regression: Fast algorithms and generalization bounds. *IEEE Trans. Pattern Anal. Mach. Intell.*, 27(6):957–968, June 2005.
- [13] D. Liben-Nowell and J. Kleinberg. Tracing information flow on a global scale using internet chain-letter data. *Proceedings of the National Academy of Sciences*, 105(12):4633–4638, 2008.
- [14] J. Lorenz, H. Rauhut, F. Schweitzer, and D. Helbing. How social influence can undermine the wisdom of crowd effect. *Proceedings of the National Academy of Sciences*, 108(22):9020–9025, 2011.
- [15] J. McAuley and J. Leskovec. Hidden factors and hidden topics: Understanding rating dimensions with review text. In *RecSys*, 2013.
- [16] L. Muchnik, S. Aral, and S. J. Taylor. Social influence bias: A randomized experiment. *Science*, 341(6146):647–651, 2013.
- [17] J. W. Pillow, J. Shlens, L. Paninski, A. Sher, A. M. Litke, E. J. Chichilnisky, and E. P. Simoncelli. Spatio-temporal correlations and visual signaling in a complete neuronal population. *Nature*, 454(7206):995–999, 2008.
- [18] C. P. Robert and G. Casella. *Monte Carlo Statistical Methods (Springer Texts in Statistics)*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2005.
- [19] M. J. Salganik, P. S. Dodds, and D. J. Watts. Experimental study of inequality and unpredictability in an artificial cultural market. *Science*, 311(5762):854–856, 2006.
- [20] V. S. Sheng, F. Provost, and P. G. Ipeirotis. Get another label? improving data quality and data mining using multiple, noisy labelers. In *KDD*, 2008.
- [21] H. Sun, A. Morales, and X. Yan. Synthetic review spamming and defense. In *KDD*, 2013.
- [22] J. Surowiecki. *The Wisdom of Crowds*. Anchor, 2005.
- [23] D. Wang, C. Song, and A.-L. Barabási. Quantifying long-term scientific impact. *Science*, 342(6154):127–132, 2013.
- [24] D. Wang, Z. Wen, H. Tong, C.-Y. Lin, C. Song, and A.-L. Barabási. Information spreading in context. In *WWW*, 2011.
- [25] T. Wang, M. Srivatsa, D. Agrawal, and L. Liu. Microscopic social influence. In *SDM*, 2012.
- [26] D. Zhou, J. C. Platt, S. Basu, and Y. Mao. Learning from the wisdom of crowds by minimax entropy. In *NIPS*, 2012.
- [27] K. Zhou, H. Zha, and L. Song. Learning triggering kernels for multi-dimensional hawkes processes. In *ICML*, 2013.

APPENDIX

Surrogate Function.

We first prove that the objective function $\mathcal{L}_\lambda(\Theta)$ as defined in Eqn.(2) and its surrogate function $Q(\Theta; \Theta^{(n)})$ as defined in Eqn.(3) satisfy the following relationships:

$$\begin{cases} Q(\Theta; \Theta^{(n)}) \geq \mathcal{L}_\lambda(\Theta) & \forall \Theta, \Theta^{(n)} \\ Q(\Theta^{(n)}; \Theta^{(n)}) = \mathcal{L}_\lambda(\Theta^{(n)}) & \forall \Theta^{(n)} \end{cases}$$

First, according to the definition of Eqn.(1), we have:

$$\begin{aligned} \mathcal{L}_\lambda(\Theta) &= \frac{1}{N} \sum_i \left(\log \sum_k \exp(\phi_{i,k}) - \sum_k y_{i,k} \phi_{i,k} \right) \\ &\quad + \frac{\lambda}{2} (\|\Theta\|_F^2 + \mathcal{R}(f)) \end{aligned}$$

We focus on the log-sum-exponential term $\log \sum_k \exp(\phi_{i,k})$ and apply the following quadratic upper bound [12]: for any vectors $\mathbf{u} \in \mathbb{R}^K$ and $\mathbf{v} \in \mathbb{R}^K$,

$$\begin{aligned} \log \sum_k \exp(u_k) &\leq \sum_k (u_k - v_k)^2 - \frac{1}{K} \left(\sum_k (u_k - v_k) \right)^2 \\ &\quad + \sum_k \frac{\exp(v_k)(u_k - v_k)}{\sum_{k'} \exp(v_{k'})} + \log \sum_k \exp(v_k) \end{aligned}$$

In our context, we replace the log-sum-exponential term of $\mathcal{L}_\lambda(\Theta)$ with its upper bound and substitute u_k with $\phi_{i,k}$ and v_k with $\phi_{i,k}^{(n)}$, which then leads to the result of $Q(\Theta; \Theta^{(n)}) \geq \mathcal{L}_\lambda(\Theta)$.

To prove $Q(\Theta^{(n)}; \Theta^{(n)}) = \mathcal{L}_\lambda(\Theta^{(n)})$, it is noted that the upper bound above is exact when $\mathbf{u} = \mathbf{v}$.

Euler-Lagrange Equation.

To derive Eqn.(7), it is first noticed that the optimization problem in Eqn.(6) can be rewritten as:

$$\min_{f \in L_1(\mathbb{R})} \int_0^\infty F(f, f') dt$$

where $F(f, f')$ is defined by:

$$F(f, f') = A_t \mathbb{I}\{t \in \mathbb{N}\} f(t)^2 + B_t \mathbb{I}\{t \in \mathbb{N}\} f(t) + \frac{\lambda}{2} (f'(t))^2$$

According to Euler-Lagrange equation, the solution of this problem satisfies the following differential equation:

$$\frac{\partial F}{\partial f} - \frac{d}{dt} \frac{\partial F}{\partial f'} = 0$$

By substituting F with the definition above, we get the differential equation in Eqn.(7).

Sample Size.

Here we derive the number of samples required for the given thresholds ϵ and δ as in Eqn.(11). Without loss of generality, consider the k -th element of \mathbf{x}_{N+M+1} , $x_{N+M+1,k}$. Following the Hoeffding's inequality, we have:

$$Pr(|\hat{x}_{N+M+1,k} - \mathbb{E}[\hat{x}_{N+M+1,k}]| \geq \epsilon) \leq \exp\left(-\frac{2L^2 \epsilon^2}{\sum_{i=1}^L (b_i - a_i)^2}\right)$$

where $x_{N+M+1,k}^{(i)}$ is bounded by $[a_i, b_i]$.

Notice that $\hat{\mathbf{x}}_{N+M+1}$ is an unbiased estimator of \mathbf{x}_{N+M+1} and $x_{N+M+1,k}^{(i)}$ is bounded by $[0, 1]$. By setting the right side of the inequality above to δ , we obtain Eqn.(11).