

Distributed Rating Prediction in User Generated Content Streams

Sibren Isaacman, Stratis Ioannidis,
Augustin Chaintreau, and Margaret Martonosi

Technicolor Technical Report

Number: **PAC-2011-05-01**

First publication date: Sept. 2011

Abstract: Recommender systems predict user preferences based on a range of available information. For systems in which users generate streams of content (*e.g.*, blogs or periodically-updated newsfeeds), users may rate the produced content that they read, and be given accurate predictions about future content they are most likely to prefer. We design a distributed mechanism for predicting user ratings that avoids the disclosure of information to a centralized authority or an untrusted third party: users disclose the rating they give to certain content only to the user that produced this content.

We demonstrate how rating prediction in this context can be formulated as a matrix factorization problem. Using this intuition, we propose a distributed gradient descent algorithm for its solution that abides with the above restriction on how information is exchanged between users. We formally analyse the convergence properties of this algorithm, showing that it reduces a weighted root mean square error of the accuracy of predictions. Although our algorithm may be used many different ways, we evaluate it on the Netflix data set and prediction problem as a benchmark. In addition to the improved privacy properties that stem from its distributed nature, our algorithm is competitive with current centralized solutions. Finally, we demonstrate the algorithm's fast convergence in practice by conducting an online experiment with a prototype user-generated content exchange system implemented as a Facebook application.



1. INTRODUCTION

A considerable portion of web activity today can be attributed to direct interactions between online users. Blogs, social networks such as Facebook, micro-blogging applications such as Twitter, and video sharing sites such as YouTube have provided online platforms through which users author and share original content, establishing an online following that often overlaps with their real-life social circles.

The sheer volume of content generated daily in the blogosphere and on social networks makes identifying relevant and interesting content a challenging task. At present, providers of these services have deployed rating mechanisms through which a user can give feedback on content generated by other users. Facebook and Twitter have expanded their rating mechanisms to the blogosphere, competing with other traditional aggregators such as Digg, Reddit and StumbleUpon that offer similar rating interfaces. Access to such ratings allow such companies to improve their recommendations but also to profile users; such profiles are a resource that companies monetize, *e.g.*, through advertising.

On the other hand, the increased monetization of private data has been met by a sharp rise in privacy concerns within advocacy groups like the Electronic Frontier Foundation and regulatory bodies like the US Congress [2]. Privacy in general, and online privacy in particular, is recognized as a fundamental human right by laws such as the European Directive on Protection of Personal Data and the Electronic Communications Privacy Act in the U.S. Nevertheless, there are many reasons why online users readily disclose private information to the above companies. Behavioral economists have identified bounded rationality, immediate gratification, underestimating risk [1] and the paradox of control [5] as some of them. Nevertheless, the importance of respecting the privacy of online users has been recognized not only by the authorities and non-profit organizations but by companies as well. For example, Google and Microsoft have struggled to find compromises that respect the privacy of their users without significantly undermining their profits [29,31].

The challenge thus arising from this state of affairs is enabling users to view and access relevant, interesting, user-generated content without releasing their private information to untrusted third parties. In this paper, we propose a solution to this problem by designing a distributed rating prediction mechanism that allows users to restrict information sharing only among trusted parties.

More specifically, we consider a general system in which users generate and share streams of content items. For example, a content producer in such a system may maintain a blog, a Facebook wall or a Twitter feed. Updates generated by users (*i.e.*, new blog posts, wall entries or tweets) are shared with select subscribers, who may respond to received content by rating it. Note that, at present, such systems are centralized. However, this need not be the case in the system we consider here: users may generate and share both their content as well as their ratings with each other in a peer-to-peer fashion.

Our goal is to design a fully distributed mechanism that learns and predicts the ratings of content consumers—*i.e.*, the subscribers. More specifically, we would like to devise a collaborative filtering scheme operating under the constraint that information is shared only between (a) a content producer and (b) its subscribers. For example, a user rating certain content should share this rating *only* with the user that generated the content.

Again, the main reason motivating our focus on distributed

rating prediction is privacy: we wish to avoid the disclosure of ratings or, in fact, any user information (such as profiles), to a central authority or an untrusted third party. Another reason is scalability, as a distributed approach can scale as the number of content producers and consumers increases. Furthermore, since ratings are exchanged among content producer/consumer pairs, our algorithm naturally exploits the social relationships between them: if two friends subscribe to each other's feed, training our predictor based on content exchanges between them can exploit latent similarities between their interests.

Our contributions can be summarized as follows:

- We propose a mathematical model of a system for distributed sharing of user-generated content streams. Our model captures a variety of different applications, and incorporates correlations both on how producers deliver content and how consumers rate it.
- We illustrate that estimating the *probability distribution* of content ratings can be naturally expressed as a Matrix Factorization (MF) problem. This is in contrast to standard MF formulations that focus on estimating ratings directly, rather than their distribution. To the best of our knowledge, our work is the first to apply a MF technique in the context of rating prediction in user-generated content streams.
- Using the above intuition, we propose a decentralized rating prediction algorithm in which information is exchanged only across content producer/consumer pairs. Producers and consumers maintain their own individual profiles; a producer shares its profile *only* with consumers to which it delivers content, and consumers share a rating they give to an item, as well as their profile, *only* with the producer that generated it.
- In spite of the above restriction on how information is exchanged among users, our distributed prediction algorithm optimizes a global performance objective. In particular, we formally characterize the algorithm's convergence properties under our model, showing that it reduces a weighted mean square error of its rating distribution estimates.
- We validate our algorithm empirically. First, we use the Netflix data set as a benchmark to compare the performance of our distributed approach to offline centralized algorithms. Second, we developed a Facebook application that reproduces the main features of a peer-to-peer content exchange environment. Using a month-long experiment with 43 users we show that our algorithm predicts ratings accurately with limited user feedback.

The remainder of this paper is organized as follows. In Section 2 we briefly describe the relationship of our work to previous rating prediction mechanisms. In Section 3 we introduce our mathematical model, and in Section 4 we present our distributed prediction algorithm as well as our main results on its convergence. Our numerical evaluation over the Netflix data set and our case study using a Facebook application are in Sections 6 and 7, respectively. Finally, we conclude in Section 8.

2. RELATED WORK

A traditional approach to distributed rating prediction in peer-to-peer networks involves computing a similarity metric (*e.g.*, Pearson correlation) among neighboring peers; predicted ratings are obtained as a function of the ratings in a peer’s neighborhood, taking into account its similarity to its neighbors [9, 19, 22, 25, 30]. Neighbor selection is part of system design, as predictions improve when neighbors have high similarity scores. We depart considerably from these works by focusing on dynamic user generated content (UGC), rather than static content, as well as by providing formal performance guarantees on the prediction error.

Matrix factorization is popular among collaborative filtering methods due to its ability to scale over large datasets [4, 16, 27]. A simple method for obtaining a low-rank factorization of a rating matrix is gradient descent on the prediction RMSE, potentially with added regularization terms (see, *e.g.*, [27], and the references therein). Assuming that ratings follow the Gaussian distribution, minimizing the RMSE is equivalent to maximizing the likelihood of observed entries [26]. In this work, we do not rely on a Gaussian assumption; instead, we directly apply MF to the *probability distribution* of ratings in dynamic UGC streams. Moreover, our distributed algorithm behaves as gradient descent over a *weighted* RMSE metric, whose weights correspond to content delivery rates across different consumers.

Recent work on MF has demonstrated that if entries removed from an r -rank matrix are selected uniformly at random, the matrix can be reconstructed accurately with high probability. Candès and Recht [7] show that, under certain conditions, if at least $\Omega(n^{6/5}r \log n)$ entries of an $n \times n$ -sized matrix are known, the entire matrix can be retrieved w.h.p. by solving a semi-definite problem (a nuclear norm minimization); Candès and Tao [6] reduce this bound to $\Omega(nr \log n)$. Keshavan *et al.* [13] propose an alternative approach that reconstructs the matrix with the same bound on known entries as [6]: missing entries are replaced by zeros and a low-rank approximation through singular value decomposition (SVD) is attained; the error on the known entries is then minimized through gradient descent. In fact, simply getting a low-rank approximation through SVD on the matrix padded with zeros is known to reconstruct the original matrix w.h.p., albeit under looser lower bounds [3]. These methods have been extended to provide guarantees of matrix reconstruction w.h.p. in the presence of noise [14, 32].

These approaches yield stronger results than [27] as, due to the lack of convexity, gradient descent does not always converge to the original matrix, or even a minimizer of the prediction error. Nevertheless, gradient descent methods can reduce the prediction error even when entries are not removed uniformly at random—which is true for the system we study. Most importantly, the algorithms in [3, 6, 7, 13, 14, 32] are centralized and thus cannot be directly applied to the distributed setting considered in this work.

Although the above traditional MF methods do not focus on distributed implementations, there exist distributed algorithms for principal component analysis (PCA) of the adjacency matrix of a graph that restrict message exchanges only across adjacent nodes. These algorithms relate to our work as PCA can be used to construct an optimal low-rank approximation of the adjacency matrix—though, technically speaking, a MF need not be ortho-normalized, so PCA is harder than factorization. Kempe and McSherry [12] propose such an algorithm based on the power method, and Korada *et al.* [15] extend it to PCA over the expectation of the adjacency matrix of a random graph. Tomozei and Massoulié [28] consider a similar setting as [15] rely-

ing however on Oja’s method [23]. Unfortunately, such approaches do not apply to our setting as we *do not* factorize the random adjacency matrix restricting communication between nodes/users (matrix $A(k)$ in our model): we factorize the distribution of ratings, which are decoupled from the communication process. Moreover, in contrast to our algorithm, the above methods are not fully distributed, as ortho-normalizing principal components requires broadcasting or gossiping normalization factors across all nodes.

Privacy in collaborative filtering was previously addressed using homomorphic encryption [8], randomized perturbation [24] and concordance measure [18] techniques. Our approach, in contrast, relies on trust between participants: we restrict information exchanges only between trusted content producer-consumer pairs. Within this context, secure multi-party computation is not required to predict ratings without disclosing user ratings publicly.

3. SYSTEM MODEL

In this section, we present the mathematical model that we use in our analysis.

3.1 Content Sharing

We consider a set \mathcal{U} of users generating and sharing content in a peer-to-peer manner. A subset $\mathcal{N} \subseteq \mathcal{U}$ of all users, whom we call *producers*, generate a stream of items. For example, producers could maintain a blog, a news-feed or a twitter-feed. Time is divided into timeslots, and in each timeslot every producer generates a new content item (*i.e.*, a blog entry or a tweet) that is added to her locally-maintained feed; this is subsequently shared with other users in a set $\mathcal{M} \subseteq \mathcal{U}$, which we call *consumers*. Note that \mathcal{M} and \mathcal{N} may intersect, as a user may both produce and consume content.

Producers share items only with consumers that belong to their social circle and/or subscribe to their feed. Even so, items may fail to reach all such consumers, either because the producer shares them selectively or because the consumers fail to receive them (*e.g.*, because they do not observe the feed continuously).

We model this as follows. Let $a_{i,j}(k) \in \{0, 1\}$ be a binary random variable indicating whether the item generated by $i \in \mathcal{N}$ at timeslot k is delivered to $j \in \mathcal{M}$, and let $A(k) = [a_{i,j}(k)]_{i \in \mathcal{N}, j \in \mathcal{M}}$ be the corresponding $|\mathcal{N}| \times |\mathcal{M}|$ matrix. We make the following assumption:

ASSUMPTION 1. $\{A(k)\}_{k \in \mathbb{N}}$ is an *i.i.d.* sequence.

I.e., deliveries are independent and identically distributed across time. Let $\lambda_{i,j} = \mathbb{E}[a_{i,j}]$ be the probability that i delivers an item to j . If j does not subscribe to i ’s feed, then $\lambda_{i,j} = 0$. Different consumers may receive items from a feed with different probabilities; our model thus allows heterogeneity in how content is targeted by producers and how often consumers fail to observe it. Note that Assumption 1 *does not imply* that deliveries between different user pairs are independent. *E.g.*, $A(k)$ may be such that an item delivered to Alice is always also delivered to Bob.

3.2 Content Ratings

Whenever a producer i delivers content to a consumer j , the consumer provides some feedback to i in the form of a *rating*. We denote by \mathcal{O} the set of possible ratings provided by consumers. In general, ratings are application-dependent. For example, consumers may indicate through an appropriate interface whether they liked, disliked or were neutral towards the content, so that $\mathcal{O} = \{+, -, \emptyset\}$. Con-

sumers may also indicate their interest on a scale from 1 (lowest interest) to 5 (highest interest), *i.e.*, $\mathcal{O} = \{1, 2, 3, 4, 5\}$.

Let $\mathcal{I}_{i,j} = \{k \in \mathbb{N} : a_{i,j}(k) = 1\}$ be the set of timeslots at which i delivers content to j . W.l.o.g., for all $k \in \mathcal{I}_{i,j}$, j provides a rating to i within the duration of the timeslot k . Depending on the application, lack of feedback can be modelled either as a failed delivery ($a_{i,j}(k) = 0$) or an additional rating (element in \mathcal{O}). For $k \in \mathcal{I}_{i,j}$, let $r_{i,j}(k) \in \mathcal{O}$ be the rating given by j to i 's content. We assume that:

ASSUMPTION 2. $\{r_{i,j}(k)\}_{k \in \mathcal{I}_{i,j}}$ is an *i.i.d.* sequence.

Note that ratings *need not* be independent across users. For example, Alice and Charlie may always give the same rating to an item from Bob.

3.3 Rating Distributions

Denote by $\tilde{\pi}_{i,j}^o$, $o \in \mathcal{O}$, the probability that $r_{i,j}(k) = o$ for $k \in \mathcal{I}_{i,j}$. For every rating $o \in \mathcal{O}$, let

$$\tilde{\Pi}^o = [\tilde{\pi}_{i,j}^o]_{i \in \mathcal{N}, j \in \mathcal{M}}. \quad (1)$$

This is a $|\mathcal{N}| \times |\mathcal{M}|$ matrix; each element corresponds to a producer/consumer pair i, j , and contains the probability that j gives rating o to an item from i . Our goal is to correctly estimate the probability matrices $\tilde{\Pi}^o$, for all $o \in \mathcal{O}$. *I.e.*, for any producer/consumer pair and any rating, we wish to find the probability that the consumer will react to content generated by the producer by providing this rating. Most importantly, we wish to do so in a distributed fashion, by restricting information exchanges only directly between producers and consumers.

Note that every time a consumer rates a content item the rating is in effect a sample from the distribution defined by the matrices $\tilde{\Pi}^o$, $o \in \mathcal{O}$. However, due to the heterogeneity of the process $A(k)$, samples are not obtained at the same rate. In fact, if $\lambda_{i,j} = 0$, the distribution of ratings of the (i, j) producer/consumer pair is *never* sampled; this relates the estimation of $\tilde{\Pi}^o$ to matrix completion, as certain entries of $\tilde{\Pi}^o$ are missing and cannot be directly observed.

Contrary to traditional matrix completion, missing entries of $\tilde{\Pi}^o$ are not selected uniformly at random—their absence is determined by, *e.g.*, the feeds to which a consumer subscribes. However, just as in traditional matrix completion, it is very natural to make the following assumption:

ASSUMPTION 3. *The probability matrices $\tilde{\Pi}^o$ are low-rank.*

This assumption implies that MF techniques can be applied to estimate these matrices; indeed our algorithm, presented in Section 4, exploits this relationship. Note that the low-rank property holds for the *rating distribution*, rather than the ratings themselves—as each producer generates an infinite number of items, this distinction is necessary.

Assumption 3 can be interpreted as a consequence of a generative latent factor model and the total probability theorem. As such, it is indeed very natural in the context of our system. We illustrate this below.

3.4 A Low-Rank Latent Factor Model

In this section we give an example of how the low-rank property of the matrices $\tilde{\Pi}^o$ may manifest by making additional assumptions on how users generate and rate content. This is only for the sake of illustration and to motivate our approach: our main results (Theorems 1 to 3) rely only on the assumptions we have made so far (in particular, Assumptions 1 and 2).

Suppose that content items generated by producers are grouped by similarity with respect to some features, thus forming a partition of the “content universe” into categories. For example, categories may pertain to topics (*e.g.*, news, music, sports, *etc.*). Formally, assume the existence of a set $\tilde{\mathcal{F}}$, whose elements we refer to as categories, such that every content item generated by a producer belongs to a single category. When $i \in \mathcal{N}$ generates new item, this item belongs to category $f \in \tilde{\mathcal{F}}$ with probability

$$\tilde{p}_{i,f} \geq 0, \quad \text{where } \sum_{f \in \tilde{\mathcal{F}}} \tilde{p}_{i,f} = 1, \quad (2)$$

independently of any categories of items the producer has generated in the past. Moreover, when $j \in \mathcal{M}$ views an item in category $f \in \tilde{\mathcal{F}}$, j provides rating $o \in \mathcal{O}$ with probability

$$\tilde{q}_{j,f}^o \geq 0, \quad \text{where } \sum_{o \in \mathcal{O}} \tilde{q}_{j,f}^o = 1, \quad (3)$$

independently of any ratings the user has given in the past.

Then, from the total probability theorem, the probability that j gives rating o when viewing content from i is:

$$\tilde{\pi}_{i,j}^o = \sum_{f \in \tilde{\mathcal{F}}} \tilde{p}_{i,f} \tilde{q}_{j,f}^o = \langle \tilde{p}_i, \tilde{q}_j^o \rangle, \quad o \in \mathcal{O}, \quad (4)$$

or, in matrix form,

$$\tilde{\Pi}^o = \tilde{P} \cdot (\tilde{Q}^o)^T, \quad o \in \mathcal{O}, \quad (5)$$

where $\tilde{P} = [\tilde{p}_{i,f}]_{i \in \mathcal{N}, f \in \tilde{\mathcal{F}}}$ and $\tilde{Q}^o = [\tilde{q}_{j,f}^o]_{j \in \mathcal{M}, f \in \tilde{\mathcal{F}}}$. *I.e.*, the matrices $\tilde{\Pi}^o$, $o \in \mathcal{O}$, admit a $|\tilde{\mathcal{F}}|$ -rank decomposition and, as such, their ranks are at most $|\tilde{\mathcal{F}}|$. Thus, *if the number of categories is small, the matrices $\tilde{\Pi}^o$ are low-rank*. Moreover, the l.h.s. matrix \tilde{P} is the same in all $|\mathcal{O}|$ decompositions. The low-rank property is thus a consequence of the existence of content categories and the total probability theorem.

4. DISTRIBUTED RATING PREDICTION

The rating prediction problem in the context of this work is to correctly estimate the probability matrices $\tilde{\Pi}^o$, $o \in \mathcal{O}$, in a distributed fashion. Below, we first formulate this problem as a MF problem and then present our distributed gradient-descent mechanism for its solution.

4.1 Problem Formulation as Matrix Factorization

To estimate $\tilde{\Pi}^o$ through MF, we construct low-dimensional *profiles* of producers and consumers. The inner product of two such profiles yields an estimate of the rating probabilities $\tilde{\pi}_{i,j}^o$. More specifically, let $d \in \mathbb{N}$ be a small integer such that $d \ll \min(|\mathcal{N}|, |\mathcal{M}|)$ and denote by \mathcal{F} the set $\{1, 2, \dots, d\}$. For now, we make no assumptions on how d relates to the rank of $\tilde{\Pi}^o$ (*i.e.*, the number of “categories” $|\tilde{\mathcal{F}}|$ in the example given by (5)).

Each producer $i \in \mathcal{N}$ maintains a vector $p_i \in [0, 1]^d$ that satisfies (2), which we call the *production profile* of i . Similarly, each consumer $j \in \mathcal{M}$ maintains $|\mathcal{O}|$ vectors $q_j^o \in \mathbb{R}_+^d$, one for every rating $o \in \mathcal{O}$, that satisfy (3). These $|\mathcal{O}|$ vectors constitute the *consumption profile* of j :

$$q_j = (q_j^{o_1}, q_j^{o_2}, \dots, q_j^{o_{|\mathcal{O}|}}).$$

Given the above profiles, our estimate of the probability $\tilde{\pi}_{i,j}^o$ (the probability that when j views content generated by i it gives rating o)—is computed as follows:

$$\tilde{\pi}_{i,j}^o = \sum_{f \in \mathcal{F}} p_{i,f} q_{j,f}^o = \langle p_i, q_j^o \rangle, \quad o \in \mathcal{O}. \quad (6)$$

Note the similarity between Equations (6) and (4). The constraints (2) and (3) immediately imply that the inner

products in (6) are non-negative and $\sum_{o \in \mathcal{O}} \pi_{i,j}^o = 1$, *i.e.*, the latter indeed constitute a probability distribution.

Our goal is then to devise an algorithm for finding profiles p_i, q_j such that the prediction $\pi_{i,j}^o$ is as close to $\tilde{\pi}_{i,j}^o$ as possible. More formally, we wish to solve the following optimization problem:

RATING PREDICTION

$$\text{Minimize } E = \sum_{i \in \mathcal{N}, j \in \mathcal{M}} \lambda_{i,j} \sum_{o \in \mathcal{O}} |\tilde{\pi}_{i,j}^o - \pi_{i,j}^o|^2, \quad (7a)$$

$$\text{subject to: } p_i \in D_1, \quad i \in \mathcal{N}, \text{ and} \quad (7b)$$

$$q_j \in D_2, \quad j \in \mathcal{M}, \quad (7c)$$

where D_1 and D_2 are the sets of profiles that satisfy (2) and (3), respectively. We call the objective function E in (7a) the *error* of our estimate. It corresponds to the error of a rating selected uniformly at random among ratings within a timeslot. Its minimization is equivalent to minimizing a weighted root mean square error (RMSE), with weights equal to the delivery rates $\lambda_{i,j}$. In particular, when $\lambda_{i,j} = 0$ (*i.e.*, when i never delivers content to j), then E does not account for the distance between $\pi_{i,j}^o$ and $\tilde{\pi}_{i,j}^o$. This is not a bug but a useful feature: we never have to predict how j would react to content from i unless it receives such content.

Assumption 3 implies that there exists a small dimension, namely $|\tilde{\mathcal{F}}|$, such that if $d \geq |\tilde{\mathcal{F}}|$ the minimum error will be zero. When $d < |\tilde{\mathcal{F}}|$, there may not exist profiles yielding $E = 0$. In other words, *underestimating* the number of categories may preclude achieving a zero error; this is possible as the number of ‘‘categories’’ is often unknown.

4.2 A Distributed Learning Algorithm

Our distributed algorithm for solving RATING PREDICTION is specified in Figure 1. It is fully distributed, and ensures that a consumer discloses the rating of a content item only to the producer that generated it. Moreover, producers share their profiles only with consumers that subscribe to their feeds, and vice-versa.

In more detail, whenever producer $i \in \mathcal{N}$ delivers content to consumer $j \in \mathcal{M}$, the following interactions take place. First, in addition to the content item, i sends to j its profile p_i . Second, the consumer views and rates the content item with a rating $r_{i,j} \in \mathcal{O}$. Third, the consumer reports to the producer (a) the rating $r_{i,j}$ for this content, as well as (b) its consumption profile q_j . We assume that consumer j reports its rating and consumption profile to i instantaneously, upon receipt of the item. Nevertheless, our results can be directly extended to the case where these are reported with an arbitrary delay within the current timeslot (see Appendix A).

Upon exchanging the above information, i and j update their profiles as follows:

$$p_i \leftarrow p_i + \gamma \sum_{o \in \mathcal{O}} (\mathbb{1}_{r_{i,j}=o} - \langle p_i, q_j^o \rangle) q_j^o \quad (8a)$$

$$q_j^o \leftarrow q_j^o + \gamma (\mathbb{1}_{r_{i,j}=o} - \langle p_i, q_j^o \rangle) p_i, \quad o \in \mathcal{O} \quad (8b)$$

where $\gamma = \gamma(k)$ is the learning rate. We assume that γ satisfies the following properties:

$$\gamma(k) \geq 0, \quad \sum_{k=1}^{\infty} \gamma(k) = \infty, \quad \text{and} \quad \sum_{k=1}^{\infty} (\gamma(k))^2 < \infty, \quad (9)$$

which hold, *e.g.*, for $\gamma(k) = 1/k$.

Finally, at the end of the timeslot, after (8a) and (8b) have been applied at every encounter the producers and consumers further update their profiles by forcing them to

Producer i at timeslot k :

i generates new content item

$\gamma \leftarrow \gamma(k)$

for every pair i, j s.t. $a_{i,j}(k) = 1$:

i sends its item and p_i to j .

i receives $r_{i,j}$ and q_j from j .

$p_i \leftarrow p_i + \gamma(k) \sum_{o \in \mathcal{O}} (\mathbb{1}_{r_{i,j}=o} - \langle p_i, q_j^o \rangle) q_j^o$

$p_i \leftarrow \Pi_{D_1}(p_i)$

Consumer j at timeslot k :

$\gamma \leftarrow \gamma(k)$

for every pair i, j s.t. $a_{i,j}(k) = 1$:

j receives item and p_i from i .

j rates item with $r_{i,j} \in \mathcal{O}$.

j sends $r_{i,j}$ and q_j to i .

for every $o \in \mathcal{O}$:

$q_j^o \leftarrow q_j^o + \gamma(\mathbb{1}_{r_{i,j}=o} - \langle p_i, q_j^o \rangle) p_i$.

$q_j \leftarrow \Pi_{D_2}(q_j)$.

Figure 1: Distributed learning algorithm.

satisfy (2) and (3):

$$p_i \leftarrow \Pi_{D_1}(p_i), \quad i \in \mathcal{N} \quad (10a)$$

$$q_j \leftarrow \Pi_{D_2}(q_j), \quad j \in \mathcal{M} \quad (10b)$$

where $\Pi_{D_1} : \mathbb{R}^{|\mathcal{F}|} \rightarrow D_1$, $\Pi_{D_2} : \mathbb{R}^{|\mathcal{O}| \times |\mathcal{F}|} \rightarrow D_2$ are the orthogonal projections to D_1 and D_2 , respectively. Due to the simple geometry of D_1, D_2 , there are known algorithms that compute these projections efficiently [20, 21]. In particular, Π_{D_1} can be computed in $O(|\mathcal{F}|)$ steps, while Π_{D_2} can be computed in $O(|\mathcal{O}||\mathcal{F}|)$ steps.

4.3 Convergence Properties

In this section, we state our main result: the dynamics of our algorithm *always* lead to a decrease in E . This important property ensures that updating the production and consumption profiles improves our predictions, even if we underestimate the number of dimensions (*i.e.*, when $d < |\tilde{\mathcal{F}}|$). We also characterize the limit points of this system and prove that it is locally stable around an optimal solution. In the rest of this section, we state formally the above results. All proofs are included in the Appendix.

Our argument is as follows. First, we show that the trajectories of the stochastic system we study asymptotically converge to the solution of a system of ordinary differential equations. Second, we prove that the error E always decreases; this, in turn, implies local stability properties. We note that our results rely on Assumptions 1 and 2 but not on Assumption 3: in particular, the weighted error E always decreases *even if the matrices $\tilde{\Pi}^o$, $o \in \mathcal{O}$, are not low-rank*. Nevertheless, as discussed in Section 4.1, when Assumption 3 fails to hold, profiles under which E is zero may not exist.

Convergence to the solution of a system of ODEs.

Consider the partial gradients

$$\nabla_{p_i} E = [\partial E / \partial p_{i,f}]_{f \in \mathcal{F}}, \quad \nabla_{q_j} E = [\partial E / \partial q_{j,f}^o]_{f \in \mathcal{F}, o \in \mathcal{O}},$$

and the following system of ODEs:

$$\frac{dp_i}{dt} = -\nabla_{p_i} E + z_{D_1}, \quad i \in \mathcal{N}, \quad (11a)$$

$$\frac{dq_j}{dt} = -\nabla_{q_j} E + z_{D_2}, \quad j \in \mathcal{M}, \quad (11b)$$

where z_{D_1}, z_{D_2} , defined formally in Appendix A, are the minimum forces required to keep the evolution of the system within the constraint set $D_1 \times D_2$. In other words, the ODE

(11) can be viewed as a *projected gradient descent* for the RATING PREDICTION problem.

Let $p_i(k), q_j(k)$, be the production and consumption profiles of $i \in \mathcal{N}, j \in \mathcal{M}$, respectively, at timeslot $k \in \mathbb{N}$. We extend these to continuous-time functions $p_i : \mathbb{R}_+ \rightarrow D_1, q_j : \mathbb{R}_+ \rightarrow D_2$ as follows. Let $t_k = \sum_{i=1}^k \gamma(k)$ and define for all i and j the continuous-time interpolated processes:

$$p_i(s) = p_i(k), \quad q_j(s) = q_j(k), \text{ for } s \in [t_k; t_{k+1}].$$

Theorem 1 below establishes that after sufficiently long time, the continuous-time interpolated trajectories of our system can be arbitrarily well approximated by solutions to (11). Moreover, any limit points of our system must also be limit points of the ODE (11). Its proof can be found in Appendix A.

THEOREM 1. *Consider the interpolated processes*

$$[p_i(t), q_j(t)]_{i \in \mathcal{N}, j \in \mathcal{M}}, \quad t \in \mathbb{R}_+.$$

Then, for any $T > 0$ there exist solutions of (11)

$$[p_i^*(t), q_j^*(t)]_{i \in \mathcal{N}, j \in \mathcal{M}}, \quad t \in \mathbb{R}_+,$$

such that, as $k \rightarrow \infty$, the quantity

$$\sup_{t_k \leq \tau \leq t_{k+T}} \left(\sum_i \|p_i(\tau) - p_i^*(\tau)\|_\infty + \sum_j \|q_j(\tau) - q_j^*(\tau)\|_\infty \right)$$

converges to 0, in probability. In addition, $[p_i, q_j]_{i \in \mathcal{N}, j \in \mathcal{M}}$ converge, in probability, to the limit set of ODE (11).

Error reduction, local stability.

Armed with the above description of system dynamics, the following theorem establishes that under any solution of (11) the error E decreases with time.

THEOREM 2. *If, $p_i(t), q_j(t), t \in \mathbb{R}_+$, evolve according to (11), then $\frac{dE}{dt} \leq 0$.*

Hence, the evolution of (11) pushes the system in the right direction, reducing the error function E . The above result is true *irrespective* of whether the user profiles have the same dimension as the rank of $\tilde{\Pi}^\circ$. Even if the ranks of the probability matrices are underestimated, the dynamics of (11) still push towards a decrease. The proof can be found in Appendix B.

Moreover, Theorem 2 has an implication regarding the stability of (11) around the minimizers of the error E . The following corollary follows immediately from Theorem 2 and Lyapunov's stability theorem.

COROLLARY 1. *Let $x^* = [p_i^*, q_j^*]_{i \in \mathcal{N}, j \in \mathcal{M}}$ be a solution of the RATING PREDICTION problem (7). Then, (11) is locally stable in the sense of Lyapunov around x^* .*

In other words, if the system starts close to a global minimum x^* , it is guaranteed to never drift away from it. Whenever E is locally convex, we can make a stronger statement:

THEOREM 3. *Let $x^* = [p_i^*, q_j^*]_{i \in \mathcal{N}, j \in \mathcal{M}}$ local minimum of E and assume that E is convex in a neighborhood around x^* . Then, there exists a δ such that if (11) is within a δ -neighborhood of x^* , it will converge to a point x s.t. $E(x) = E(x^*)$.*

Hence, if the ODE starts from profiles close enough to x^* , it is guaranteed the converge to profiles at an error at least as good as the error of x^* .

5. SYSTEM EXTENSIONS

5.1 Improving Predictions using Item Profiles

Recall that content items are sent to consumers along with the production profile of their producer. When an item is propagated among several consumers, rather than modifying these profiles, it may be preferable to adapt them *as the item is propagated from one consumer to the next*.

In particular, the system can be extended so that content items generated by a producer a are associated with a profile vector $t \in [0, 1]^d$, that is initialized to $t = p_a$ when the item is generated. This profile is delivered along with the item to every consumer that it passes through. However, instead of remaining static, the item profile is adapted through (8a) and (10): after the k -th viewing of the item, say at consumer b , the profile can be adapted as follows

$$t(k+1) = \Pi_{D_1} \left(t + \gamma \sum_{o \in \mathcal{O}} (\mathbb{1}_{r_{a,b}^o} - \langle t, q_b \rangle) q_b^o \right) \quad (12)$$

To gain some intuition as to why this would improve predictions, recall that each generated item belongs to a certain category. At the time of its generation, our prediction of which category an item belongs to is given by \tilde{p}_a , namely, the inherent probability distribution its producer. However, as an item travels through the system, the ratings of users can be used to enforce our belief on the category the item belongs to. The update rule (12) aims at exploiting the additional knowledge gained by consumer reactions.

A formal analysis of joint dynamics of a system in which publisher, consumer *and* item profiles are adapted is quite intricate and beyond the scope of the present paper. Nevertheless, we verified the performance of such a combined evolution through our experimental study in Section 7, where item profiles are introduced into the system's design.

5.2 Using Predictions for Recommendations

In this section we show that when the predicted rating distribution is correctly estimated we can generate optimal recommendations.

More specifically, we design two recommendation algorithms, each aiming to attain a different objective. The first maximizes the fraction of content items that receive a given rating, subject to a constraint on the number of items shown to the consumer. The second maximizes the number of content items shown to the user, given a lower bound on the fraction of content items that receive a given rating. Below, we discuss these two objectives in more detail, and then present the two algorithms that indeed meet these objectives.

5.2.1 Recommendation Objectives

Assume that a sequence of content items, labeled as $k = 1, 2, 3, \dots$ arrive at a given consumer. When a new item arrives, a recommendation algorithm decides whether to display (*i.e.*, recommend) this item to the consumer.

Let $N(k)$ be the number of items that the algorithm has recommended to the consumer up to and including the k -th item. For $o \in \mathcal{O}$, $N^o(k)$ is the number of items that were rated o . Obviously, $N^o(k) \leq N(k)$, as the user reacts only to items recommended by our algorithm. Let

$$r = \lim_{k \rightarrow \infty} \frac{N(k)}{k} \quad r^o = \lim_{k \rightarrow \infty} \frac{N^o(k)}{N(k)}$$

be the fraction of items recommended to the user and the fraction of these items that receive rating o , respectively.

We will refer to r and r° as the *recommendation rate* and the *o-rate*, respectively.

Constrained Bitrate.

The objective of our first recommendation algorithm is to maximize r° subject to the constraint that $r = \rho$, for some $\rho \in [0, 1]$. We call this the *constrained-bitrate* objective.

$$\text{Maximize: } r^\circ, \quad (13a)$$

$$\text{subject to: } r = \rho, \quad (13b)$$

for some $\rho \in (0, 1)$. Intuitively, the constrained-bitrate objective is most suitable when the application requires a steady, controlled rate of items of being presented to the consumer. This is natural, *e.g.*, in a network where bandwidth or the users attention span is limited: given that a device (or the consumer herself) cannot process more than a fraction r of all incoming items, (13) aims to maximize the fraction of displayed items receiving rating o .

Maximum Bitrate.

The constrained-bitrate objective is suitable when the user desires to see a specific, steady volume of content items. However, the above objective does not give guarantees on how good r° can be: it will be the maximum possible, but that will largely depend on the quality of the content arriving at the user. In a system where bandwidth or attention span is not an issue, it makes sense giving guarantees in terms of r° , rather than r :

$$\text{Maximize: } r, \quad (14a)$$

$$\text{subject to: } r^\circ \geq \rho, \quad (14b)$$

where, again, $\rho \in (0, 1)$. This objective guarantees that the *o-rate* of the content shown to the consumer will be at least ρ . It is more suitable than (13) when the user prefers guarantees on the quality—rather than the volume—of the content recommended to her; subject to these guarantees, the consumer wishes to see as much content as possible.

5.2.2 Optimal Recommendation Algorithms

In this section we present two algorithms for recommending items in that achieve objectives (13) and (14), respectively. These recommendation algorithms have access to a prediction algorithm like the one presented in Section 4: for every item k , the recommendation algorithm can obtain $\pi^\circ(k)$, $o \in \mathcal{O}$, the probability that the consumer gives rating o to this item.

The recommendation algorithm that solves the optimization problem (13) is the following. At each point in time, the algorithm maintains a threshold τ . When a new item arrives, the recommendation algorithm obtains from the prediction algorithm of Section 4 the quantity π° , *i.e.*, the probability that the item receives a rating $o \in \mathcal{O}$. The item is then recommended to the user if and only if $\pi^\circ \geq \tau$.

The threshold τ is updated when the k -th item arrives in the following manner: it decreases whenever the quantity $\pi^\circ(k)$ is below τ and increases whenever $\pi^\circ(k)$ is above τ . In particular, after the viewing of the k -th item,

$$\tau(k+1) = \tau(k) + \gamma(k) (\mathbb{1}_{\pi^\circ(k) \geq \tau(k)} - \rho) \quad (15)$$

where $\gamma(k)$ is a decreasing gain factor satisfying (9).

The recommendation algorithm that solves (14) is similar. A threshold τ is again used to determine whether an item with a given π° will be recommended to the user. The difference between the two algorithms lies in how this threshold is updated: here, the threshold increases when the item is

shown and receives rating o , and decreases under other ratings. More specifically, given that an item is shown to the consumer, denote by $r \in \mathcal{O}$ the rating the consumer provided. Then the threshold τ is updated as follows

$$\tau(k+1) = \tau(k) + \gamma(k) (\mathbb{1}_{\pi^\circ(k) \geq \tau(k)} \wedge r(k) = + - \rho \mathbb{1}_{\pi^\circ(k) \geq \tau(k)}) \quad (16)$$

where $\gamma(k)$ again a decreasing gain factor satisfying (9). Intuitively, the threshold increases when the item is shown and receives rating o , and decreases under other ratings.

5.2.3 Optimality of Recommendations

We establish in this section the optimality of the recommendation algorithms introduced in Section 5.2.2. We assume here that the prediction algorithm has already converged, and gives perfect predictions of the rating distribution, *i.e.*, $\pi^\circ = \tilde{\pi}^\circ$; an evaluation of the two recommendation algorithms under imperfect predictions can be found in Section 6.

The algorithm will be optimal among the class of recommendation algorithms that decide whether to display the k -th item based on $\pi^\circ(k)$. Formally, these algorithms can be defined through a function $x(\pi^\circ)$, which is 1 if the algorithm shows the item given that the approval probability is π° , and 0 otherwise. We will say that a recommendation algorithm is a *threshold algorithm* if $x(\pi^\circ) = \mathbb{1}_{\pi^\circ \geq \tau}$ for some threshold $\tau \in [0, 1]$.

The following theorem states that the algorithm meeting the constrained bitrate objective is a threshold algorithm; moreover, the recommendation algorithm (15) eventually converges to this threshold.

THEOREM 4. *Assume that the sequence $\pi^\circ(k)$, $k = 1, \dots$ is i.i.d. Further assume that the c.d.f. $\mathbf{P}(\pi^\circ < z)$ is Lipschitz continuous and strictly monotone. Then, the optimal solution of the constrained-bitrate optimization problem is a threshold algorithm $\mathbb{1}_{\pi^\circ \geq \tau^*}$ and (15) converges to τ^* w.p.1.*

The proof of this Theorem can be found in Appendix D.

A similar result holds for the iterative algorithm (16). In this case, the problem (14) may not always be feasible: for example, if the items submitted never receive the rating o , the *o-ratio* cannot be bounded from below. Nevertheless, provided that the problem is feasible, there exists an optimal solution that is of a threshold algorithm.

THEOREM 5. *Assume that the sequence $\pi^\circ(i)$, $i = 1, \dots$ is i.i.d.. Further assume that the c.d.f. $\mathbf{P}(\pi^\circ < z)$ is Lipschitz continuous and strictly monotone. Then, if (14) is feasible, it admits an optimal solution that is a threshold algorithm $\mathbb{1}_{\pi^\circ \geq \tau^*}$. Moreover, if τ^* is this policy's threshold and $\gamma(k) = \frac{1}{k}$ then and (16) converges to τ^* w.p.1.*

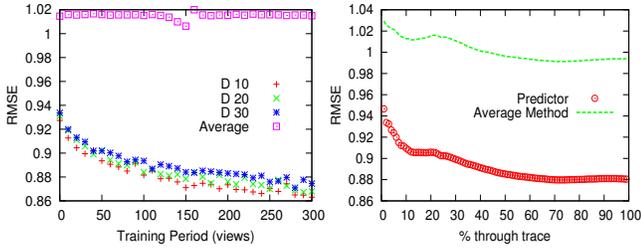
This theorem is proved in Appendix E

6. COMPARISON ON NETFLIX DATA SET

In this section, we test our rating prediction algorithm on the Netflix dataset. The wide use of the dataset as a benchmark allows us to implicitly compare the performance of our algorithm to the one achieved by state-of-the-art algorithms.

6.1 Netflix Data Set

In 2006, Netflix announced a competition for recommendation systems, and released a dataset on which competitors could train their algorithms. The Netflix dataset consists of pairs of anonymized movies and anonymized users. Each trace entry includes a timestamp, the user ID, and the user's rating (on an integer scale of 1 to 5).



(a) RMSE v. training time (b) RMSE Evolution

Figure 2: (a) RMSE as a function of the training period. (b) RMSE evolution through the course of the Netflix trace. Most of the drop in RMSE occurs in the first 10% of the trace.

The dataset includes both publicly-available training data, for which ratings were provided, and a testing dataset, for which ratings were not disclosed. If for every movie in the test set, we simply always predict its average rated value from the training set, this approach would yield a root mean square error of 1.0540 on the test set. The winning team of the Netflix Prize challenge generated predictions with RMSE of 0.8572 [4], a 10% improvement of the RMSE of Cinematch (0.9525), the algorithm designed by Netflix engineers.

We apply our algorithm as follows. Each movie is given a production profile $p_m \in [0, 1]^d$. Similarly, each user is given a consumption profile $q_u \in [0, 1]^{5 \times d}$, corresponding to ratings with one, two, three, four, or five stars respectively (*i.e.*, $\mathcal{O} = \{1, 2, 3, 4, 5\}$). The Netflix dataset is arranged chronologically and as users rate movies, p for the movie and q for the user are updated according to (8a) and (8b) with each rating and are subsequently projected to D_1 and D_2 according to (10).

6.2 Evaluating our Predictions on Netflix

We evaluate the performance of our predictions of user behavior in two ways, as discussed below.

RMSE of rating prediction.

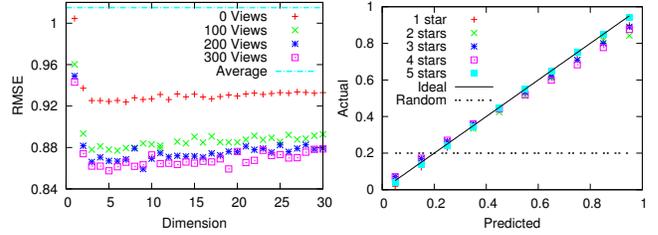
Prior to each rating event, we predict the rating that a user, u , will give to a movie, m . We make this prediction by reporting the expected rating based on our estimation of the rating probabilities, *i.e.*,

$$\text{PredictedRating}(u, m) = \sum_{o \in \mathcal{O}} o \cdot \langle q_u^o, p_m \rangle.$$

We can then compare our prediction to the user’s rating for this movie as recorded in the dataset, to calculate an RMSE.

Our algorithm starts with a randomly selected p for each movie and a random q for each user and adapts them as users rate movies. To account for a training period, we discard some of the early predictions before computing the RMSE. More precisely, we compute an RMSE for predictions with training period more than k by including predictions for which either the movie or the user profile has been adapted at least k times.

Although we do not know how the Cinematch algorithm (Netflix’s baseline) performs on the *training* dataset, we know that on the test set it performs roughly 10% better than the naïve algorithm that guesses the average rating for each movie. Therefore, to evaluate our effectiveness, we compare our RMSE against this “average” algorithm applied to the training dataset (100,480,507 ratings that 480,189 users gave to 17,770 movies). We know the RMSE of this



(a) RMSE v. dimension (b) Rating Distribution

Figure 3: (a) Effect of dimensionality on the RMSE, which rapidly decreases until the time is insufficient to train all dimensions. (b) Comparison between predicted and empirical rating distributions.

“average” algorithm for both training (1.015) and test (1.05) datasets, and so we can compare to these numbers.

Figure 2(a) shows the improvement in RMSE as we vary the training period for different numbers of dimensions d of the vectors. Regardless of the length of the training period, the RMSE of the “average” method remains at 1.015. Again, though we do not know how the original Netflix algorithm (Cinematch) functioned and how it performs on the training dataset, we expect this centralized solution to achieve roughly 10% improvement (9.18 RMSE over the training dataset). Our algorithm outperforms this value with a training period of only 10 iterations; a 15% improvement can be reached with training period of 300 views. Even when the training period is set to zero (*i.e.*, we include all predictions in the RMSE), it improves on the “average” method by 8% for all values of d , only slightly worse than Cinematch. Thus, our technique performs at least as well as some of the leading systems, even though our algorithm functions in a completely distributed manner.

Figure 2(b) shows how the RMSE evolves throughout the course of the netflix trace. Our technique makes 50% of its improvements in RMSE during the first 10% of the trace, indicating that it performs well even in cold start situations. The algorithm consistently tracks the performance of the “average” method while demonstrating a consistent improvement of 8%-12%. Figure 3(a) shows the effect of modifying the dimensionality of the prediction vectors on the RMSE for different training periods. As the dimensionality increases, we see a large drop; the RMSE quickly reaches a plateau between $d = 3$ and $d = 10$, and then increases slowly. We see the same trend for all training periods, with longer training periods tending to be flatter.

Thus, choosing the dimensionality is a tradeoff. On the one hand, a minimum number is needed to capture the diversity of the system, while on the other hand, too many dimensions makes the algorithm slow to learn and requires a longer training period. In practice, the best number of dimensions can be small (typically, 5 or 10), which also means that very little meta-data (5 to 10 floating-point numbers) must be transferred with each piece of content.

Predicted vs. empirical rating distributions.

We have seen that our algorithm performs well in terms of rating prediction RMSE. However, it can also be used to estimate the rating probability distribution. We now evaluate whether our estimated distribution is observed in practice across all users and movies: *e.g.*, for all events where we predict that rating o is obtained with probability 0.3, do we observe that this rating is seen 30% of the time? This is, in

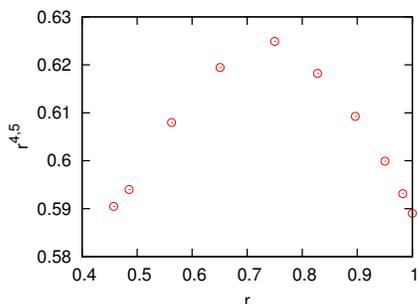


Figure 4: Fraction of movies rated 4 or 5 ($r^{4,5}$) against total fraction of movies viewed (r).

effect, a measure of the precision of the rating predictor.

We measure the goodness of fit of our distribution to the empirical distribution in Figure 3(b). To compute it, we first bin our predicted probability into 10 bins (0 to 0.1, 0.1 to 0.2, etc.). We then compare the bin value to the actual rate of occurrence in the bin, for different ratings. Note that for all ratings, the square correlation coefficient R^2 is above 0.98 indicating a very good match. The slope of the best fit line is thus nearly 1. Another metric is the distance of our result to the line $y = x$ which would represent an ideal predictor. Random guessing, for instance, has a mean distance of 0.415. Our rating predictor brings a large improvement, with distance errors as small as 6×10^{-4} .

This accuracy implies that the system can do more than provide a single estimated rating for a specific movie and user, it can actually characterize its behavior on a finer grain which may help provide additional confidence. As an example, it is interesting to note that users are more predictable when they express strong opinions (5 stars, 1 stars), slightly less predictable when they are neutral, and even less predictable when they indicate 4 or 2 stars.

6.3 Evaluating Recommendations

Another application of our prediction technique is to apply a threshold as in Section 5.2 to limit the bytes transferred to the user while maximizing the fraction of content they enjoy (here defined as all movies they rated 4 or 5).

Figure 4 plots the ratio of movies viewed and the fraction of these that are enjoyed for different threshold values. A low threshold value corresponds to almost all movies being shown (on the right on the Figure). Using our prediction probability and a larger threshold allows to both reduce the volume of movies shown, until about 75% while achieving a maximum fraction of 63% content enjoyed. However, applying larger threshold does not improve this fraction since learning begins to slow.

This indicates that using the rating prediction is effective to select the most pertinent content to display to the user, providing both an improvement in bandwidth as well as an improved user experience. On the other hand, it is important to make sure that users get exposed to sufficient content for the algorithm to train.

7. CASE STUDY: WEBDOSE

The evaluation of our rating prediction algorithm on the Netflix data set demonstrates that it can, indeed, be used to make accurate predictions. However, it does not illustrate its behavior in a distributed, heterogeneous environment.

Therefore, we implemented our system as a Facebook application (“WEBDOSE”), which allows Facebook users to view,

rate, and share content in the form of web pages. We used Facebook as a distribution platform to ensure that the application would be immediately available to a wide range of users without requiring extensive development overhead.

7.1 Application Description

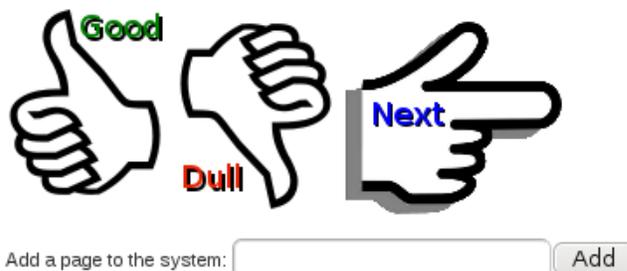


Figure 5: Function buttons in the WebDose application. These buttons are displayed along with the web page to be rated by the user.

User interface.

When users log into the system, WEBDOSE presents them with a screenshot of a web page and three rating icons: “thumbs up,” “thumbs down,” or “next” (see Fig. 5). We interpret each of these ratings as positive, negative or neutral, respectively (*i.e.*, $\mathcal{O} = \{+, -, \emptyset\}$). By design, users are not shown a new page until they have clicked a button and thus rated the page they are currently viewing. Additionally, users have the option to “Add” (*i.e.*, produce) content to the system by typing a url into a text box at the bottom of their screen. WEBDOSE, thus, follows the same producer/consumer model used throughout this paper.

Emulating Content Sharing.

We emulate content sharing through discrete “contact events” between users. When a user logs into the application (or refreshes the page), she experiences a “contact” with one user that has also logged in to WEBDOSE within the previous three hours. If no such user exists, no contact occurs. A user may experience many contacts during a browsing session as other users log in and trigger their own “contact events.”

To account for social relationships, when multiple people have logged in the system in the past 3 hours, we bias contact events between users that are Facebook friends. If multiple possible swaps exist, we choose to swap with a random friend half the time and the other half of the time we choose randomly among all users (including friends).

Content is “produced” when a user adds a web page via the provided interface. Each user is limited to viewing web-pages in a “cache” whose contents change dynamically. When two users “meet”, each item in each user’s cache is randomly reassigned to the other user with probability 0.5. Web pages thus perform a “random walk” through WEBDOSE users, and there is no centralized control dictating how they move.

Behind-the-scenes prediction.

The first time a user u logs into the system, she is assigned a random production profile p_u and a random consumption profile q_u . The dimension of these profiles is set to $d = 10$; this will be examined in Section 7.3.

We also incorporate item profiles in WEBDOSE. When an

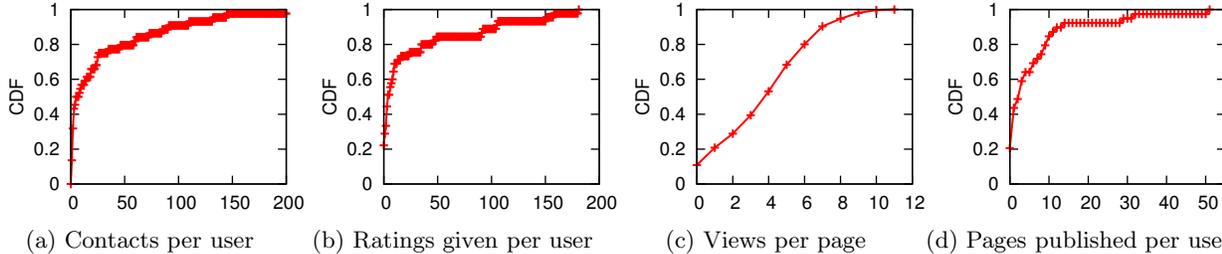


Figure 6: CDFs of various facets of user behavior throughout the experiment. In (a) the most active 50% experience a wide range of contacts. In (b) we see that 20% of users rate more than 50 pages each, indicating high use. In (c) we observe that no page is viewed more than 11 times but half the pages are viewed more than 4 times. Finally, (d) indicates that the majority of content is published by a small number of users.

item is propagated among several consumers, rather than modifying these profiles, it is preferable to adapt them *as the item is propagated from one consumer to the next*. The system is extended so that content items generated by a producer a are associated with a profile vector $t \in [0, 1]^d$, that is initialized to $t = p_a$ when the item is generated. This profile is delivered along with the item to every consumer that it passes through. However, instead of remaining static, the item profile is adapted through (8a) and (10). A formal analysis of joint dynamics of a system in which publisher, consumer *and* item profiles are adapted is quite intricate and beyond the scope of the present paper. Nevertheless, we verify the performance of such a combined evolution through the WEBDOSE experiment.

As web pages are rated by other users of the system, t_w and q_u are updated according to Eq. (8a) and (8b) respectively. Both the predicted as well as the actual rating are logged. Web pages in the system carry an identifier of the web page’s producer, *i.e.*, the user that added it to WEBDOSE. When a user rates a page, the system stores the producer and rating locally. Subsequently, when the two users meet, the web page consumer informs the web page producer of the rating given to its content. The producer then uses this information to update its profile p_u according to Algorithm (8a). All updates are followed by an immediate projection to D_1, D_2 , as in (10).

Recommendation.

WEBDOSE uses our prediction algorithm to select pages to show to a user. Each time a page is shown, WEBDOSE first calculates π^+ for each as-yet-unrated web page in the user’s cache. The page with the highest π^+ is then shown to the user. Once that page has been rated, π^+ is recalculated for each unrated page and the new highest page is shown. We chose this form of recommendation to ensure that a user is always able to view all of the content in their cache, and that WEBDOSE would not be hindered from training user’s profiles through pages not viewed.

7.2 WEBDOSE Experiment

During a 33 day period, 43 users spanning 12 time zones registered for WEBDOSE and viewed or added 326 web pages. Although these small numbers make learning user preferences difficult, we will show in Section 7.3 that the system was able to adapt well. Users were free to log in as frequently as they wished and could rate as many pages as they wanted. Once a page was rated, it was never shown to that user again. As an incentive to add quality content to the system, we provided each user with statistics of the ratings their content received, as well as their ranking in

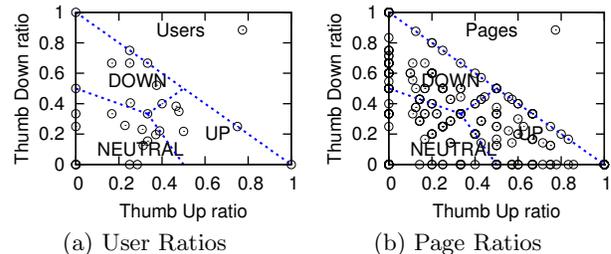


Figure 7: Regions are marked to indicate dominant user behavior. In (a) users are marked according to their ratio of thumb ups to thumb downs. Users generally fall into the thumbs down or neutral areas. However, (b) shows that this is not a property of the pages themselves as pages fill the space uniformly.

“thumbs-up” and “thumbs-down” received.

User activity.

Here we examine some of the trends in the way that users interact with this system. Because users are given few instructions beyond “thumb it up if you like it, down if you don’t, and remember to add pages,” we feel that the broad trends from this deployment would hold in any application built with a similar model.

Figure 6(a) provides the CDF of the number of “contact-events” that users experienced. Half of the users had more than 10 contacts and 20% had over 50. Figure 6(b) shows the CDF of the ratings given by each user. The most active user rated over 180 pages; 25% of users rated more than 20 pages. This is consistent with a real mobile environment: most users look for web pages occasionally, while an active few regularly check for pages.

For the system to function as intended the pages must be viewed frequently enough that the item profile t is trained. Figure 6(c) shows the frequency with which pages are viewed. The high number of pages in the system (over 300) compared to the number of users means that each page is viewed by only a limited number of users. Over the course of our experiment, no page was viewed more than 11 times. However, more than half the pages were viewed more than 4 times, which proves to be sufficient.

Finally, we examine the process by which web pages are added to the system. Figure 6(d) examines the number of pages added to the system by each user. The top 20% of users each produce more than 10 pages with the highest producer contributing around 50 web pages (16% of all web

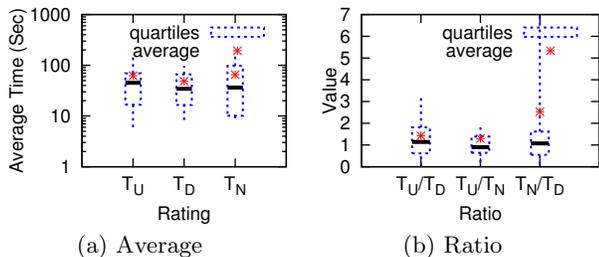


Figure 8: In these boxplots, the ends of the whiskers denote the 5th and 95th percentiles of users while the box denotes the 25th and 75th percentiles. The dark bar is the median and the star is the mean. In (a) we look at the raw times taken by users to rate a page. (b) examines the ratio of times taken by a single user. Thumbing down a page is faster than thumbs up and neutral, both of which take roughly the same amount of time.

pages in the system). It is not necessarily true that the highest producers are also the highest consumers. The relatively low number of high producers indicate that it is crucial that item profiles are introduced, otherwise all pages from the same producer will have the same profile, which would slow learning. Item profiles allow the system to learn and distinguish different pages from the same producer, accelerating convergence by training individual web pages while the production profiles train.

User rating behavior.

Here, we focus on how users rated the pages in WEBDOSE. Define the “Thumb up (down)” ratio to be the fraction of all pages that a user rates positively (negatively). Figure 7 examines these ratios by users and by pages.

Figure 7(a) shows that the majority of users favor neutral ratings, followed closely by those that mostly rate pages thumbs down. With most people thus being generally ambivalent or negative on most content, it is important that the system makes accurate predictions.

Figure 7(b) shows the rating ratios given to individual pages. Despite the generally negative user behavior, pages themselves fill the interest space uniformly. The wide range of topics serves to train the system as well as demonstrate the need for accurate recommendation.

We next examine the speed with which users provide ratings to the system. Intuitively, a user will spend more time looking at a web page that she finds enjoyable and quickly reject pages that do not interest her. Figure 7.2 shows comparisons of the speed with which users make decision about various ratings.

We define T_U, T_D , and T_N as the time taken to rate a page “up,” “down,” or “neutral,” respectively. In Figure 8(a) we see the distribution of times taken to rate a page. “Thumb up” indeed takes the longest to select, followed by neutral. By examining the ratio of these quantities, as in Figure 8(b), we see that 70% of users thumb a page up slower than they thumb a page down (*i.e.* $\frac{T_U}{T_D} > 1$). Thus, one may be able to judge the relative like or dislike of a page simply by observing how long the user takes to make a decision. Interestingly, it may in principle be possible to remove ratings completely, relying only on time to judge a users interest, or to nuance the ratings with timing information.

7.3 Predictive Power

	<i>a priori</i>		convergence	
	R^2	slope	R^2	slope
Thumb Up	0.43	0.29	0.98	0.95
Thumb Down	0.80	0.63	0.98	1.20
Neutral	0.96	1.07	0.95	1.29

Table 1: R^2 values and slopes for thumb rate predictions both *a priori* and after convergence.

WEBDOSE Prediction accuracy.

To assess the correctness of predictions in WEBDOSE, we repeated the evaluation of our predictors with respect to the two metrics we considered over the Netflix data set, the RMSE and the goodness of fit of the predicted distribution.

Figure 9(a) displays the RMSE of our predictions as a function of the dimension d of our profiles. The dimension 10 was used in the actual experiment and other values were obtained by running our algorithm on the collected trace.

To compute the RMSE, we associate the values $-1, 0, +1$ to $-, \emptyset, +$ respectively. Before the rating of a web page w by a user u , we again generate our estimate of the rating as

$$\text{PredictedRating}(u, v) = \sum_{o \in \{-1, 0, +1\}} o \cdot \langle t_w, q_u^o \rangle,$$

where t_w, q_u the item and consumption profiles, respectively. As in Netflix, we again observe the error quickly decreases up to a dimension of 10, after which it remains constant. We note that the RMSE for $d = 10$ is 0.819, 42% less than the RMSE obtained by predictions obtained by the running average ratio of each rating.

The goodness of fit of our distribution to the empirical distribution is illustrated in Figure 9(b) (using the same metrics from Sec. 6.2). The values for R^2 and the slope of the best fit lines are provided in Table 1.

The outcomes “Thumb down” and “neutral” are both predicted accurately, with slopes above 0.63 and exhibiting good fits (R^2 above 0.8). The “thumbs-up” predictions are considerably worse, with a slope of 0.29. This is not surprising, given that most ratings in the system were for “thumbs-down” and “neutral”: the experiment does not provide enough data to train quickly enough for “thumbs-up” ratings.

Since WEBDOSE has far more complicated dynamics than the simple model we analysed, it is interesting to investigate the convergence of our prediction scheme. For this reason, we repeated the goodness of fit test using instead the content and consumption profiles *at the end* of the experiment. The result is shown in Figure 9(c), with R^2 vales and slopes again found in Table 1. It is quite clear that the content and consumption profiles have adapted to fit the empirical distribution of ratings very well.

8. CONCLUSIONS

User generated content is a primary factor to the success of the social web as it connects users through content sharing. It presents unique challenges for collaborative filtering due to a large volume and the sensitivity of sharing information beyond one’s immediate circle of acquaintances. This paper proves that information exchange can be deployed only among trusted pairs while providing strong guarantees on the rating prediction error, thanks to a fully distributed and asynchronous learning algorithm.

Our results expand the scope of recommender systems to operate on top of any social content sharing platform, which unveils important research questions. The case where trust is not reciprocal (such as the follower relationship within Twitter) remains an interesting open case. Our algorithm

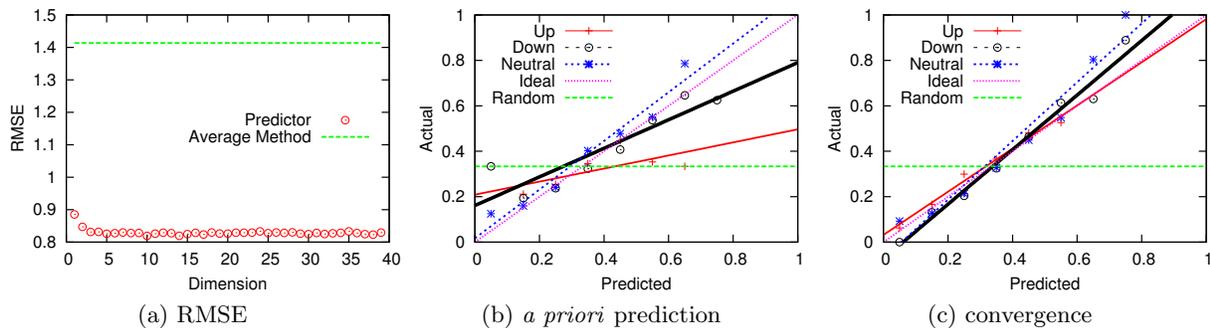


Figure 9: Predictions of the thumb rate of pages. (a) shows the effect of the dimension chosen at the outset of the experiment. As the dimension increases, the error decreases. Our selection of 10 occurs below the knee of the curve. In (b) and (c) pages are binned according to estimated percentage of “thumbing” them in a given direction. (b) shows *a priori* prediction and moderate accuracy, with slopes from 0.29 to 1.07. (c) shows convergence of π at the end of the experiment and has slopes from 0.95 to 1.3. In all figures, the green line gives a comparison to other predictors.

could leverage secure multiparty computation, to provide at a smaller cost the privacy guarantee offered by centralized schemes. Finally, our model can be used to analyze how social proximity, captured through “rate of delivery” for producer-consumer pairs, impacts the efficiency of learning. Both questions are interesting open problems.

9. ACKNOWLEDGEMENTS

This work was partially funded by the European Commission under the FIRE SCAMPI project (FP7-IST-258414) and by the French National Research Agency (ANR) under the PROSE project (VERSO program).

10. REFERENCES

- [1] ACQUISTI, A., AND GROSSKLAGS, J. What can behavioral economics teach us about privacy? In *Digital Privacy: Theory, Technologies and Practices* (2007), pp. 363–377.
- [2] ANGWIN, J. US seeks web privacy ‘bill of rights’. *Wall Street Journal* (Dec. 17th 2010).
- [3] AZAR, Y., FIAT, A., KARLIN, A., MCSHERRY, F., AND SAIA, J. Spectral analysis of data. In *STOC* (2001).
- [4] BELL, R. M., AND KOREN, Y. Scalable collaborative filtering with jointly derived neighborhood interpolation weights. In *ICDM* (2007).
- [5] BRANDIMARTE, L., ACQUISTI, A., AND LOEWENSTEIN, G. Privacy concerns and information disclosure: An illusion of control hypothesis. In *CIST* (2010).
- [6] CANDÈS, E., AND TAO, T. The power of convex relaxation: Near-optimal matrix completion. *IEEE Trans. Inform. Theory* 56, 5 (2009), 2053–2080.
- [7] CANDÈS, E. J., AND RECHT, B. Exact matrix completion via convex optimization. *Found. of Comput. Math.* 9 (2008), 717–772.
- [8] CANNY, J. Collaborative filtering with privacy via factor analysis. In *SIGIR* (2002).
- [9] CASTAGNOS, S., AND BOYER, A. Personalized Communities in a Distributed Recommender System. In *ECIR* (2007), pp. 343–355.
- [10] HARAUX, A. How to differentiate the projection on a convex set in Hilbert space. Some applications to variational inequalities. *J. Math. Soc. Japan* 29 (1977), 615–631.
- [11] IOANNIDIS, S., AND MASSOULIÉ, L. Surfing the blogosphere: Optimal personalized strategies for searching the web. In *Proc. IEEE Infocom* (2010).
- [12] KEMPE, D., AND MCSHERRY, F. A decentralized algorithm for spectral analysis. *Journal of Computer and System Sciences* 74, 1 (feb 2008).
- [13] KESHAVAN, R., MONTANARI, A., AND OH, S. Matrix completion from a few entries. *Trans. Inform. Theory* (2010).
- [14] KESHAVAN, R., MONTANARI, A., AND OH, S. Matrix completion from noisy entries. *JMLR* (2010).
- [15] KORADA, S. B., MONTANARI, A., AND OH, S. Gossip PCA. *SIGMETRICS* (2011).
- [16] KOREN, Y. Collaborative filtering with temporal dynamics. In *KDD* (2009).
- [17] KUSHNER, H. J., AND YIN, G. *Stochastic Approximation and Recursive Algorithms and Applications*, 2nd ed. Springer, 2003.
- [18] LATHIA, N., HAILES, S., AND CAPRA, L. Private distributed collaborative filtering using estimated concordance measures. In *RecSys* (Oct. 2007).
- [19] LIU, K., BHADURI, K., DAS, K., NGUYEN, P., AND KARGUPTA, H. Client-side web mining for community formation in peer-to-peer environments. *WEBKDD* (2006).
- [20] MACULAN, N., SANTIAGO, C., MACAMBIRA, E., AND JARDIM, M. An $O(n)$ algorithm for projecting a vector on the intersection of a hyperplane and a box in \mathbb{R}^n . *J. of Opt. Th. and App.* 117, 3 (2003), 553–574.
- [21] MICHELOT, C. A finite algorithm for finding the projection of a point onto the canonical simplex of \mathbb{R}^n . *J. of Opt. Th. and App.* 50, 1 (1986), 195–200.
- [22] MILLER, B., KONSTAN, J., AND RIEDL, J. PocketLens. *ACM Transactions on Information Systems* 22 (2004), 437–476.
- [23] OJA, E. On stochastic approximation of the eigenvectors and eigenvalues of the expectation of a random matrix. *J. of Math. Anal. and App.*, 106 (1985), 69–84.
- [24] POLAT, H., AND DU, W. Privacy-preserving collaborative filtering using randomized perturbation techniques. In *IEEE ICDM* (2003), pp. 625–628.
- [25] RUFFO, G., AND SCHIFANELLA, R. A peer-to-peer recommender system based on spontaneous affinities. *ACM Transactions on Internet Technology* 9 (2009).

- [26] SALAKHUTDINOV, R., AND MNIH, A. Probabilistic matrix factorization. *Advances in Neural Information Processing Systems 20* (2008), 1257–1264.
- [27] TAKÁCS, G., PILÁSZY, I., NÉMETH, B., AND TIKK, D. Scalable collaborative filtering approaches for large recommender systems. *JMLR 10* (2009), 623–656.
- [28] TOMOZEI, D.-C., AND MASSOULIÉ, L. Distributed user profiling via spectral methods. In *SIGMETRICS* (2010).
- [29] VASCELLARO, J. E. Google agonizes on privacy as ad world vaults ahead. *Wall Street Journal* (Aug. 10th 2010).
- [30] WANG, J., POWELSE, J., LAGENDIJK, R., AND REINDERS, M. Distributed collaborative filtering for peer-to-peer file sharing systems. In *SAC* (2006).
- [31] WINGFIELD, N. Microsoft quashed effort to boost online privacy. *Wall Street Journal* (Aug. 2nd 2010).
- [32] ZHOU, Z., WRIGHT, J., LI, X., CANDÈS, E. J., AND MA, Y. Stable principal component pursuit. In *ISIT* (2010).

APPENDIX

A. PROOF OF THEOREM 1

Our proof follows the standard method outlined in Kushner and Yin [17]. Assume that the duration of each timeslot is T . We denote the time at which the transmission from i to j by

$$t_{i \rightsquigarrow j}(k) \in [0, T) \cup \{\infty\},$$

By convention, $t_{i \rightsquigarrow j}(k) = \infty$ denotes that no such transmission takes place within the k -th slot. Put differently, $\{t_{i \rightsquigarrow j}(k) < \infty\} = \{a_{i,j}(k) = 1\}$.

We first precisely define the “minimum forces” z_{D_1} and z_{D_2} . Let $D \subseteq \mathbb{R}^d$ be a closed convex subset of \mathbb{R}^d , and consider two vectors $x \in \mathbb{R}^d$, $p \in D$. Following Kushner and Yin [17], we define $z_D(p, x) \in \mathbb{R}^d$ as follows.

$$z_D(p, x) = \lim_{\delta \rightarrow 0^+} \frac{\Pi_D(p + \delta x) - (p + \delta x)}{\delta},$$

where $\Pi_D : \mathbb{R}^d \rightarrow D$ is the orthogonal projection to D . For D closed and convex, the above limit is guaranteed to exist (see, e.g., [10]). The quantity $z_D(p, x)$ is called the minimum additional force to maintain p within D because, if “force” x is applied to p , i.e., $\frac{dp}{dt} = x$ then, for small $\delta > 0$

$$p + \delta x + \delta z_D(p, x) \simeq \Pi_D(p + \delta x).$$

Using the above notation, ODE (11) can be written as

$$\begin{aligned} \frac{dp_i}{dt} &= -\nabla_{p_i} E + z_{D_1}(p_i, -\nabla_{p_i} E), & a \in \mathcal{N}, \\ \frac{dq_j}{dt} &= -\nabla_{q_j} E + z_{D_2}(q_j, -\nabla_{q_j} E). & b \in \mathcal{M}, \end{aligned}$$

Note that, within a timeslot, as users encounter each other, their profiles change through (8a) and (8b). As a result, when a producer i encounters a consumer j , their profiles may differ from the ones they had in the beginning of the timeslot. The following lemma however states that, if we assume that during encounters profile had their initial values, we introduce only a very small error compared to how the actual process behaves:

LEMMA 1. *Let $p_i(\tau)$, $q_j(\tau)$, $\tau \in [0, T)$, the values of the profiles of users $i \in \mathcal{N}$ and $j \in \mathcal{M}$ at some time τ since the beginning of the timeslot. Denote by $E(\tau) = \{(a_i \rightsquigarrow b_i) :$*

$t_{i \rightsquigarrow b_i} \leq \tau\}$ the set of all encounter events that have taken place until time τ . Then, for any $\tau \in [0, T)$,

$$p_i(\tau) = p_i(0) + \gamma \sum_{j: i \rightsquigarrow j \in E(\tau, o)} (\mathbb{1}_{r_{i,j}=o} - \langle p_i(0), q_j^o(0) \rangle) q_j^o(0) + \beta_i, \quad (17a)$$

$$q_j^o(\tau) = q_j^o(0) + \gamma \sum_{i: i \rightsquigarrow j \in E(\tau)} (\mathbb{1}_{r_{i,j}=o} - \langle p_i(0), q_j^o(0) \rangle) p_i(0) + \beta_j, \quad (17b)$$

where β_i, β_j are random variables s.t. $\|\beta_i\|_2, \|\beta_j\|_2 \leq K\gamma^2$, where K depends only on the system dimensions $|\mathcal{N}|, |\mathcal{M}|, |\mathcal{O}|$ and $|\mathcal{F}|$.

PROOF. We will prove the lemma by induction on the encounters between users. W.l.o.g., we assume that encounters occur at distinct times: encounters that occur simultaneously can be ordered arbitrarily if they do not involved the same user, or by the order within which (8a) and (8b) take place, otherwise.

For the induction basis, observe that (17) holds at $\tau = 0$, before any encounters have taken place. Suppose that it also holds at τ^- , right before an encounter $i' \rightsquigarrow j'$ takes place. Then, it still holds time τ^+ , right after the encounter took place, for all $i \neq i'$ and all $j \neq j'$. Observe that, because $p_i \in D_1$ and $q_j \in D_2$, for all $i \in \mathcal{N}$ and all $j \in \mathcal{M}$

$$|\mathbb{1}_{r_{i,j}=o} - \langle p_i, q_j^o \rangle| \leq 1, \quad \|q_j^o\|_2 \leq |\mathcal{F}|, \quad \|p_i\|_2 \leq 1$$

As a result, using the induction hypothesis we can show that

$$|\langle p_{i'}(\tau^-), q_{j'}^o(\tau^-) \rangle - \langle p_{i'}(0), q_{j'}^o(0) \rangle| = O(\gamma)$$

where the constants in $O(\gamma)$ depend only on the system dimensions. By (8a) we have that

$$\begin{aligned} p_{i'}(\tau^+) &= p_{i'}(\tau^-) + \sum_{o'} (\mathbb{1}_{r_{i',b'}=o'} - \langle p_{i'}(\tau^-), q_{j'}^o(\tau^-) \rangle) q_{j'}^o(\tau^-) \\ &= p_i(0) + \gamma \sum_{j: i \rightsquigarrow j \in E(\tau, o)} (\mathbb{1}_{r_{i,j}=o} - \langle p_i(0), q_j^o(0) \rangle) q_j^o(0) + \beta_{i'} + \\ &\quad + \gamma \sum_{o'} (\mathbb{1}_{r_{i',b'}=o'} - \langle p_{i'}(0), q_{j'}^o(0) \rangle + O(\gamma)) ((q_{j'}^o(0) + O(\gamma))) \end{aligned}$$

where the constants in the $O(\gamma)$ terms only depend on the system dimensions. As, by the induction hypothesis, $\beta_{i'} = O(\gamma^2)$, the induction step follows; a similar argument can be made for the consumption profile q_j^o . \square

Note that the above proof by induction can be easily extended to the case where consumers report their ratings and consumption profiles at an arbitrary time after $t_{i \rightsquigarrow j}$.

Lemma 17 implies that the evolution of profiles from one timeslot to the next can be written as:

$$\begin{aligned} p_i(k+1) &= \Pi_{D_1}(p_i(k) + \gamma(k) Y_i(k) + \beta_i(k)), \\ q_j(k+1) &= \Pi_{D_2}(q_j(k) + \gamma(k) Y_j(k) + \beta_j(k)) \end{aligned}$$

where

$$Y_i = \sum_j \mathbb{1}_{a_{i,j}(k)=1} \sum_{o \in \mathcal{O}} (\mathbb{1}_{r_{i,j}=o}(k) - \langle p_i(k), q_j^o(k) \rangle) \cdot q_j^o(k)$$

and $Y_j = (Y_j^{o_1}, Y_j^{o_2}, \dots, Y_j^{o_{|\mathcal{O}|}})$ such that

$$Y_j^o = \sum_i \mathbb{1}_{a_{i,j}(k)=1} (\mathbb{1}_{r_{i,j}=o}(k) - \langle p_i(k), q_j^o(k) \rangle) \cdot p_j(k).$$

while $\beta_i(k), \beta_j(k)$ are random variables s.t. $\|\beta_i(k)\|_2 = O(\gamma^2(k)), \|\beta_j(k)\|_2 = O(\gamma_k^2)$. By the independence of propagation from the content categories, we have that

$$\mathbb{E}[\mathbb{1}_{i \rightsquigarrow j}(k) \cdot \mathbb{1}_{r_{i,j}=o}(k)] = \lambda_{i,j} \cdot \tilde{\pi}_{i,j}^o(k).$$

As a result $\mathbb{E}[Y_i] = \nabla_{p_i} E$ and $\mathbb{E}[Y_j] = \nabla_{q_j} E$. Moreover, by the fact that $p_i \in D_1$ and $q_j \in D_2$, $\mathbb{E}[\|Y_i(k)\|_2^2] < \infty$ and $\mathbb{E}[\|Y_j(k)\|_2^2] < \infty$ uniformly on k . Thus, the assumptions (A2.1)-(A2.5) of Theorem 2.3 in Chapter 5 of Kushner and Yin [17] are satisfied, and the theorem follows. \square

B. PROOF OF THEOREM 2

To prove Theorem 2, we will make use of the following lemma (Lemma 5 in [11]) concerning the orthogonal projection on the positive simplex.

LEMMA 2. Let $D = \{x \in \mathbb{R}_+^d : \sum_{i=1}^d x_i = 1\}$ and, for $x \in D$, denote by $I_0(x) = \{i : x_i = 0\}$ the set of zero-valued coordinates of x . Then, there exists a set $B \subset I_0(x)$ such that $z_D(x, y) = z$ where

$$z_i = \begin{cases} -y_i, & i \in B \\ -\frac{1}{|B|} \sum_{j \in \bar{B}} y_j, & i \in \bar{B}. \end{cases}$$

Moreover, for every $i \in B$, $y_i - \frac{1}{|B|} \sum_{j \in \bar{B}} y_j < 0$.

We have that

$$\begin{aligned} \frac{dE}{dt} &= \sum_i \langle \nabla_{p_i} E, \frac{dp_i}{dt} \rangle + \sum_j \langle \nabla_{q_j} E, \frac{dq_j}{dt} \rangle \\ &\stackrel{(11)}{=} \sum_i [-\langle \nabla_{p_i} E, \nabla_{p_i} E \rangle + \langle \nabla_{p_i} E, z_{D_1}(p_i, -\nabla_{p_i} E) \rangle] \\ &\quad \sum_j [-\langle \nabla_{q_j} E, \nabla_{q_j} E \rangle + \langle \nabla_{q_j} E, z_{D_2}(q_j, -\nabla_{q_j} E) \rangle] \end{aligned}$$

From Lemma 2, for every $i \in \mathcal{N}$ there exists a $B_i \subseteq \mathcal{F}$ such that $\langle \nabla_{p_i} E, z_{D_1}(p_i, \nabla_{p_i} E) \rangle$ equals

$$\sum_{f \in B_i} (\partial E / \partial p_{i,f})^2 + \frac{1}{|B_i|} \left(\sum_{f \in \bar{B}_i} \partial E / \partial p_{i,f} \right)^2.$$

Similarly, for every $j \in \mathcal{M}$ and every $f \in \mathcal{F}$ there exists a $B_{j,f} \subset \mathcal{O}$ such that $\langle \nabla_{q_j} E, z_{D_2}(q_j, -\nabla_{q_j} E) \rangle$ equals

$$\sum_f \left[\sum_{o \in B_{j,f}} (\partial E / \partial q_{i,f}^o)^2 + \frac{1}{|B_{j,f}|} \left(\sum_{o \in \bar{B}_{j,f}} \partial E / \partial q_{i,f}^o \right)^2 \right].$$

Using the above, dE/dt becomes:

$$\begin{aligned} \frac{dE}{dt} &= - \sum_i \sum_{f \in \bar{B}_i} \left(\frac{\partial E}{\partial p_{i,f}} - \frac{1}{|B_i|} \sum_{f' \in \bar{B}_i} \frac{\partial E}{\partial p_{i,f'}} \right)^2 \\ &\quad - \sum_{j,f} \sum_{o \in \bar{B}_{j,f}} \left(\frac{\partial E}{\partial q_{i,f}^o} - \frac{1}{|B_{j,f}|} \sum_{o' \in \bar{B}_{j,f}} \frac{\partial E}{\partial q_{i,f}^{o'}} \right)^2 \leq 0. \quad \square \end{aligned} \quad (18)$$

C. PROOF OF THEOREM 3

We begin by stating the Karush-Kuhn-Tucker (KKT) conditions for OUTPUT PREDICTION. The Lagrangian of (7) is

$$\begin{aligned} L &= E + \sum_i \mu_i \left(\sum_f p_{i,f} - 1 \right) + \sum_{j,f} \nu_{j,f} \left(\sum_o q_{j,f}^o - 1 \right) \\ &\quad - \sum_{i,f} \xi_{i,f} p_{i,f} - \sum_{j,f,o} \zeta_{j,f}^o q_{j,f}^o \end{aligned}$$

where $\mu_i, \nu_{j,f}, \xi_{i,f}, \zeta_{j,f}^o$ are the Lagrangian multipliers of the constraints (7b) and (7c). Thus, the KKT conditions can be written as

$$\frac{\partial E}{\partial p_{i,f}} + \mu_i - \xi_{i,f} = 0, \quad \xi_{i,f} \geq 0, \quad \xi_{i,f} p_{i,f} = 0, \quad \forall a, f \quad (19a)$$

$$\frac{\partial E}{\partial q_{i,f}^o} + \nu_{j,f} - \zeta_{j,f}^o = 0, \quad \zeta_{j,f}^o \geq 0, \quad \zeta_{j,f}^o q_{j,f}^o = 0, \quad \forall b, f, o \quad (19b)$$

$$p_i \in D_1, \quad q_j \in D_2, \quad \forall a, b \quad (19c)$$

The following lemma holds

LEMMA 3. Let $\{p_i^*, q_j^*\}_{i \in \mathcal{N}, j \in \mathcal{M}}$ be profiles for which $dE/dt = 0$. Then, there exist $\mu_i, \nu_{j,f}, \xi_{i,f}, \zeta_{j,f}^o$ for which the KKT conditions are satisfied.

PROOF. From (18) and Lemma 2, if $dE/dt = 0$ at then for every $i \in \mathcal{N}$ there exists a $B_i \subset \mathcal{F}$ such that, for $\mu_i = -\frac{1}{|B_i|} \sum_{f' \in \bar{B}_i} \frac{\partial E(x^*)}{\partial p_{i,f'}}$

$$\begin{aligned} p_{i,f}^* &= 0, \quad E(x^*) / \partial p_{i,f} + \mu_i > 0, & \forall f \in B_i, \\ \partial E(x^*) / \partial p_{i,f} + \mu_i, & & \forall f \in \bar{B}_i. \end{aligned}$$

Similarly, for every $j \in \mathcal{M}$ and $f \in \mathcal{F}$ there exists a $B_{j,f} \subset \mathcal{O}$ such that for $\nu_{j,f} = -\frac{1}{|B_{j,f}|} \sum_{o' \in \bar{B}_{j,f}} \frac{\partial E(x^*)}{\partial q_{j,f}^{o'}}$

$$\begin{aligned} (q_{j,f}^o)^* &= 0, \quad E(x^*) / \partial q_{j,f}^o + \nu_{j,f} > 0, & \forall o \in B_{j,f}, \\ \partial E(x^*) / \partial q_{j,f}^o + \nu_{j,f} &= 0 & \forall o \in \bar{B}_{j,f}. \end{aligned}$$

It is easy to verify that p_i^*, q_j^* , along with the above values $\mu_i, \nu_{j,f}$ and the values

$$\xi_{i,f} = \max(0, \frac{\partial E(x^*)}{\partial p_{i,f}} + \mu_i), \quad \zeta_{j,f}^o = \max(0, \frac{\partial E(x^*)}{\partial q_{j,f}^o} + \nu_{j,f})$$

satisfy the KKT conditions (19). \square

Suppose now that x^* is a local minimum and that E is locally convex at a δ -neighborhood around this minimum. Then, by Theorem 2, there exists a $\delta' \leq \delta$ such that if the ODE (11) starts from a δ' -neighborhood of x^* it will remain in this neighborhood. Consider now the problem (7) with the additional constraint $\|x - x^*\|_2 < \delta'$. From Lemma (3), that any limit point of (11) will satisfy the KKT constraints of this problem with a 0 lagrange multiplier for the added condition. Hence, any limit point of (11) must attain $E(x^*)$, as E is convex in this δ neighborhood. \square

D. PROOF OF THEOREM 4

We first show that the optimal recommendation policy is a threshold policy. Let ν be the probability measure of π_+ , that is, for every Borel set A , $\mathbf{P}(\pi_+ \in A) = \int_A d\nu(s)$. Since the sequence $\pi_+(i)$, $i = 1, \dots$ is i.i.d., the constrained-bitrate problem is equivalent to maximizing $\frac{1}{\rho} \int_0^1 s x(s) d\nu(s)$ subject to

$$\int_0^1 x(s) d\nu(s) = \rho \quad (20)$$

Consider two probability measures ν_0 and ν_1 that are absolutely continuous w.r.t. ν . Then, by the Radon-Nikodym theorem, there exist functions ℓ_0, ℓ_1 such that $\nu_0(A) = \int_A \ell_0(s) d\nu(s)$ and $\nu_1(A) = \int_A \ell_1(s) d\nu(s)$. The Neyman-Pearson lemma states the function x that maximizes $\mathbb{E}_1[1 - x(\pi_+)]$, subj. to $\mathbb{E}_0[x(\pi_+)] = \alpha$, has the following form $x(s) = \mathbb{1}_{\ell_1(s) > k \ell_0(s)}$, for some $k > 0$. Take $\ell_0(s) = 1$ and $\ell_1(s) = s/c$, where $c = \int_0^1 s d\nu(s)$. Then the Neyman-Pearson lemma

gives us that the function x that maximizes $\frac{1}{c} \int_0^1 sx(s)d\nu(s)$, subj. to $\int_0^1 x(s)d\nu(s) = \rho$, is given by $x = \mathbb{1}_{s \geq kc}$, for some k . Since ρ, c are positive constants, the above optimization problem has the same optimal solution as the constrained-bitrate objective, so a solution to the latter is indeed a threshold policy. The threshold $\tau^* = kc$ computed by the constraint (20), which implies that τ^* is the solution of the equation $\int_{\tau^*}^1 d\nu(s) = \rho$. Note that, by our hypothesis $\mathbf{P}(\pi_+ \geq \tau)$ is continuous and strictly monotone in $[0, 1]$, and the above equation has a unique solution for $\alpha \in [0, 1]$.

We now prove that (15) converges to the aforementioned τ^* . Note that $\mathbb{E}[\mathbb{1}_{\pi_+ \geq \tau}^2] \leq 1$ and let $F(\tau) = \mathbf{P}(\pi_+ \geq \tau)$, which is Lipschitz continuous. Hence, by Theorem 2.3 in Chapter 5 of [17], if $\gamma(k) = \frac{1}{k}$ the limit points (15) are a subset of the limit points of the ODE $\dot{\tau} = F(\tau) - \rho$. Consider the candidate Lyapunov function $L(\tau) = (\tau - \tau^*)^2$. Then $\frac{dL}{dt} = 2(\tau - \tau^*)(F(\tau) - \rho) \leq 0$ by the definition of τ^* and the monotonicity of F ; moreover, by the strict monotonicity of F , the above is an equality iff $\tau = \tau^*$, hence $\tau \rightarrow \tau^*$. \square

E. PROOF OF THEOREM 5

We first show that there exists a threshold strategy that is optimal. Let again ν be the probability measure of π_+ . Then (14) is equivalent to maximizing $\int_0^1 x(s)d\nu(s)$ subject to

$$\int_0^1 sx(s)d\nu(s) - \rho \int_0^1 x(s)d\nu(s) \geq 0 \quad (21)$$

Suppose that the problem is feasible, and let $x^*(s)$ be an optimal recommendation strategy. Let $opt = \int_0^1 x^*(s)d\nu(s)$. By Theorem 5 there exists a threshold strategy x s.t.

$$\int_0^1 x(s)d\nu(s) = opt$$

and

$$\int_0^1 sx(s)d\nu(s) \geq \int_0^1 sx^*(s)d\nu(s) \geq \rho \cdot opt = \rho \int_0^1 x(s)d\nu(s).$$

Hence, x is also an optimal strategy. The threshold can again be computed through (21) as the solution of the following equation $G(\tau) = \int_{\tau}^1 (s - \rho)d\nu(s) = 0$. Observe (*e.g.*, by differentiating), that under our assumptions on the c.d.f. of π_+ , $G(\tau)$ is Lipschitz, strictly increasing for $\tau \leq \rho$ and decreasing for $\tau > \rho$. Given that it is zero for $\tau = 1$, we deduce that $\tau_* \in [0, \rho]$ and that $G(\tau) < 0$ for $\tau < \tau^*$ and $G(\tau) > 0$ for $\tau \geq \tau^*$. Moreover, we also deduce that the necessary and sufficient condition for the existence of τ^* and, hence, the feasibility of (14), is that $G(0) \leq 0$. The same arguments as in the proof of Theorem 4 can be used to show that the limit points of (16) are, with probability one, a subset of the limit points of the ODE $\dot{\tau} = G(\tau)$. Using our above observations on G , it is easy to show that $(\tau - \tau^*)^2$ is a Lyapunov function of the above ODE, and that τ thus $\tau \rightarrow \tau^*$. \square