

Combinational Collaborative Filtering for Personalized Community Recommendation

Wen-Yen Chen
Computer Science
University of California
Santa Barbara, CA 93106
wychen@cs.ucsb.edu

Dong Zhang
Google Research, Beijing
No. 1 Zhongguancun E. Road
Beijing 100084, China
dongzhang@google.com

Edward Y. Chang
Google Research
Mountain View, CA 94043
edchang@google.com

ABSTRACT

Rapid growth in the amount of data available on social networking sites has made information retrieval increasingly challenging for users. In this paper, we propose a collaborative filtering method, *Combinational Collaborative Filtering* (CCF), to perform personalized community recommendations by considering multiple types of co-occurrences in social data at the same time. This filtering method fuses semantic and user information, then applies a hybrid training strategy that combines Gibbs sampling and Expectation-Maximization algorithm. To handle the large-scale dataset, parallel computing is used to speed up the model training. Through an empirical study on the Orkut dataset, we show CCF to be both effective and scalable.

Categories and Subject Descriptors

H.4.m [Information Systems Applications]: Miscellaneous

General Terms

Algorithms, Experimentation

Keywords

Collaborative filtering, probabilistic models, personalized recommendation

1. INTRODUCTION

Social networking products are flourishing. Sites such as MySpace, Facebook, and Orkut attract millions of visitors a day, approaching the traffic of Web search sites [1]. These social networking sites provide tools for individuals to establish communities, to upload and share user generated content, and to interact with other users. In recent articles, users complained that they would soon require a full-time employee to manage their sizable social networks. Indeed,

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. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

KDD'08, August 24–27, 2008, Las Vegas, Nevada, USA.
Copyright 2008 ACM 978-1-60558-193-4/08/08 ...\$5.00.

take Orkut as an example. Orkut enjoys 100+ million communities and users, with hundreds of communities created each day. A user cannot possibly view all communities to select relevant ones.

In this work, we tackle the problem of *community recommendation* for social networking sites. Such a problem fits in the framework of *collaborative filtering* (CF), which offers personal recommendations (of e.g., Web sites, books, or music) based on a user's profile and prior information-access patterns. What differentiates our work from prior work is that we propose a fusion method, which combines information from multiple sources. We name our method CCF for *Combinational Collaborative Filtering*. CCF views a community from two simultaneous perspectives: *a bag of users* and *a bag of words*. A community is viewed as a bag of participating users; and at the same time, it is viewed as a bag of words describing that community. Traditionally, these two views are independently processed. Fusing these two views provides two benefits. First, by combining *bags of users* with *bags of words*, CCF can perform *personalized* community recommendations, which the *bags of words* alone model cannot. Second, augmenting *bags of users* with *bags of words*, CCF achieves better personalized recommendations than the *bags of users* alone model, which may suffer from information sparsity.

A practical recommendation system must be able to handle large-scale datasets and hence demands scalability. We devise two strategies to speed up training of CCF. First, we employ a hybrid training strategy, which combines Gibbs sampling with the Expectation-Maximization (EM) algorithm. Our empirical study shows that Gibbs sampling provides better initialization for EM, and thus can help EM to converge to a better solution at a faster pace. Our second speedup strategy is to parallelize CCF to take advantage of the distributed computing infrastructure of modern data centers. Our scalability study on a real-world dataset of 312k active users and 109k popular communities demonstrates that our parallelization scheme on CCF is effective. CCF can achieve near-linear speedup on up to 200 distributed machines, attaining 116 times speedup over the training time of using one machine. CCF is attractive for supporting large-scale personalized community recommendations, thanks to its effectiveness and scalability.

1.1 Related Work

Several algorithms have been proposed to deal with either *bags of words* or *bags of users*. Specifically, Probabilistic Latent Semantic Analysis (PLSA) [7] and Latent Dirichlet

Allocation (LDA) [3] model document-word co-occurrence, which is similar to the *bags of words* community view. Probabilistic Hypertext Induced Topic Selection (PHITS) [4], a variant of PLSA, models document-citation co-occurrence, which is similar to the *bags of users* community view. However, a system that considers just bags of users cannot take advantage of content similarity between communities. A system that considers just bags of words cannot provide personalized recommendations: all users who joined the same community would receive the same set of recommendations. We propose CCF to model multiple types of data co-occurrence simultaneously. CCF’s main novelty is in fusing information from multiple sources to alleviate the information sparsity problem of a single source.

Several other algorithms have been proposed to model publication and email data¹. For instance, the *author-topic* (AT) model [11] employs two factors in characterizing a document: the document’s authors and topics. Modeling both factors as variables within a Bayesian network allows the AT model to group the words used in a document corpus into semantic topics, and to determine an author’s topic associations. For emails, the *author-recipient-topic* (ART) model [8] considers email recipient as an additional factor. This model can discover relevant topics from the sender-recipient structure in emails, and enjoys an improved ability to measure role-similarity between users. Although these models fit publication and email data well, they cannot be used to formulate personalized community recommendations, whereas CCF can.

1.2 Contribution Summary

In summary, this paper makes the following three contributions:

1. We propose CCF to effectively fuse multiple information sources. For community data, we illustrate that by fusing *bags of users* with *bags of words*, CCF can perform *personalized* community recommendations, which the *bags of words* alone model cannot. By adding *bags of words* to *bags of users*, CCF achieves better personalized recommendations than the *bags of user* alone model, which suffers from the information sparsity problem.
2. We devise a hybrid training method, which uses Gibbs sampling to seed EM. Empirical study shows that this hybrid training method can typically make EM converge to a better solution at a faster pace.
3. We parallelize CCF to achieve near-linear speedup on distributed machines. Empirical study shows that parallel CCF can deal with large-scale data. For instance, a training task that requires one day to complete on one machine takes only 14 minutes on 200 machines.

The remainder of the paper is organized as follows. In Section 2, we present CCF, including its model structure and semantics, hybrid training strategy, and parallelization scheme. In Section 3, we present our experimental results on both synthetic and Orkut datasets. We provide concluding remarks and discuss future work in Section 4.

¹We discuss only related model-based work since the model-based approach has been proven to be superior to the memory-based approach.

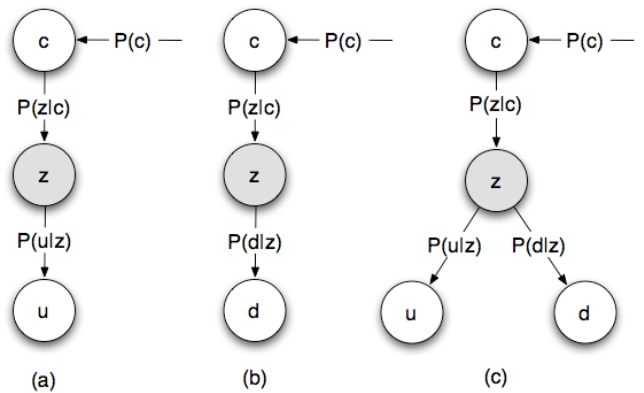


Figure 1: (a) Graphical representation of the Community-User (C-U) model. (b) Graphical representation of the Community-Description (C-D) model. (c) Graphical representation of Combinational Collaborative Filtering (CCF) that combines both bag of users and bag of words information.

2. CCF: COMBINATIONAL COLLABORATIVE FILTERING

We start by introducing the baseline models. We then show how our CCF model combines baseline models. Suppose we are given a collection of co-occurrence data consisting of communities $C = \{c_1, c_2, \dots, c_N\}$, community descriptions from vocabulary $D = \{d_1, d_2, \dots, d_V\}$, and users $U = \{u_1, u_2, \dots, u_M\}$. If community c is joined by user u , we set $n(c, u) = 1$; otherwise, $n(c, u) = 0$. Similarly, we set $n(c, d) = R$ if community c contains word d for R times; otherwise, $n(c, d) = 0$. The following models are latent aspect models, which associate a latent class variable $z \in Z = \{z_1, z_2, \dots, z_K\}$.

Before modeling CCF, we first model community-user co-occurrences (C-U), shown in Figure 1(a); and community-description co-occurrences (C-D), shown in Figure 1(b). Our CCF model, shown in Figure 1(c), builds on C-U and C-D models. The shaded and unshaded variables in Figure 1 indicate latent and observed variables, respectively. An arrow indicates a conditional dependency between variables.

2.1 C-U and C-D Baseline Models

The C-U model is derived from PLSA and for community-user co-occurrence analysis. The co-occurrence data consists of a set of community-user pairs (c, u) , which are assumed to be generated independently. The key idea is to introduce a latent class variable z to every community-user pair, so that community c and user u are rendered conditionally independent. The resulting model is a mixture model that can be written as follows:

$$P(c, u) = \sum_z P(c, u, z) = P(c) \sum_z P(u|z)P(z|c), \quad (1)$$

where z represents the topic for a community. For each community, a set of users is observed. To generate each user, a community is c chosen uniformly from the community set, then a topic z is selected from a distribution $P(z|c)$ that is specific to the community, and finally a user u is generated by sampling from a topic-specific distribution $P(u|z)$.

The second model is for community-description co-occurrence analysis. It has a similar structure to the C-U model with

the joint probability written as:

$$P(c, d) = \sum_z P(c, d, z) = P(c) \sum_z P(d|z)P(z|c), \quad (2)$$

where z represents the topic for a community. Each community's interests are modeled with a mixture of topics. To generate each description word, a community c is chosen uniformly from the community set, then a topic z is selected from a distribution $P(z|c)$ that is specific to the community, and finally a word d is generated by sampling from a topic-specific distribution $P(d|z)$.

Remark: One can model C-U and C-D using LDA. Since the focus of this work is on model fusion, we defer the comparison of PLSA vs. LDA to future work.

2.2 CCF Model

In the C-U model, we consider only *links*, i.e., the observed data can be thought of as a very sparse binary $M \times N$ matrix W , where $W_{i,j} = 1$ indicates that user i joins (or linked to) community j , and the entry is unknown elsewhere. Thus, the C-U model captures the linkage information between communities and users, but not the community content. The C-D model learns the topic distribution for a given community, as well as topic-specific word distributions. This model can be used to estimate how similar two communities are in terms of topic distributions. Next, we introduce our CCF model, which combines both the C-U and C-D.

For the CCF model (Figure 1(c)), the joint probability distribution over community, user, and description can be written as:

$$\begin{aligned} P(c, u, d) &= \sum_z P(c, u, d, z) \\ &= P(c) \sum_z P(u|z)P(d|z)P(z|c), \quad (3) \end{aligned}$$

The CCF model represents a series of probabilistic generative processes. Each community has a multinomial distribution over topics, and each topic has a multinomial distribution over users and descriptions, respectively.

2.2.1 Gibbs & EM Hybrid Training

Given the model structure, the next step is to learn model parameters. There are some standard learning algorithms, such as Gibbs sampling [6], Expectation-Maximization (EM) [5], and Gradient descent. For CCF, we propose a hybrid training strategy: We first run Gibbs sampling for a few iterations, then switch to EM. The model trained by Gibbs sampling provides the initialization values for EM. This hybrid strategy serves two purposes. First, EM suffers from a drawback in that it is very sensitive to initialization. A better initialization tends to allow EM to find a "better" optimum. Second, Gibbs sampling is too slow to be effective for large-scale datasets in high-dimensional problems [2]. A hybrid method can enjoy the advantages of Gibbs and EM.

Gibbs sampling

Gibbs sampling is a simple and widely applicable Markov chain Monte Carlo algorithm, which provides a simple method for obtaining parameter estimates and allows for combination of estimates from several local maxima of the posterior distribution. Instead of estimating the model parameters directly, we evaluate the posterior distribution on z and then use the results to infer $P(u|z)$, $P(d|z)$ and $P(z|c)$.

For each user-word pair, the topic assignment is sampled from:

$$P(z_{i,j} = k | u_i = m, d_j = n, \mathbf{z}_{-i,-j}, U_{-i}, D_{-j}) \propto \frac{C_{mk}^{UZ} + 1}{\sum_{m'} C_{m'k}^{UZ} + M} \frac{C_{nk}^{DZ} + 1}{\sum_{n'} C_{n'k}^{DZ} + V} \frac{C_{ck}^{CZ} + 1}{\sum_{k'} C_{ck'}^{CZ} + K}, \quad (4)$$

where $z_{i,j} = k$ represents the assignment of the i^{th} user and j^{th} description word in a community to topic k . $u_i = m$ represents the observation that the i^{th} user is the m^{th} user in the user corpus, and $d_j = n$ represents the observation that the j^{th} word is the n^{th} word in the word corpus. $\mathbf{z}_{-i,-j}$ represents all topic assignments not including the i^{th} user and the j^{th} word. Furthermore, C_{mk}^{UZ} is the number of times user m is assigned to topic k , not including the current instance; C_{nk}^{DZ} is the number of times word n is assigned to topic k , not including the current instance; C_{ck}^{CZ} is the number of times topic k has occurred in community c , not including the current instance.

We analyze the computational complexity of Gibbs sampling in CCF. In Gibbs sampling, one needs to compute the posterior probability

$$P(z_{i,j} = k | u_i = m, d_j = n, \mathbf{z}_{-i,-j}, U_{-i}, D_{-j})$$

for user-word pairs ($M \times L$) within N communities, where L is the number of words in community description (Note $L \geq V$). Each $P(z_{i,j} = k | u_i = m, d_j = n, \mathbf{z}_{-i,-j}, U_{-i}, D_{-j})$ consists of K topics, and requires a constant number of arithmetic operations, resulting in $O(K \cdot N \cdot M \cdot L)$ for a single Gibbs sampling. During parameter estimation, the algorithm needs to keep track of a topic-user ($K \times M$) count matrix, a topic-word ($K \times V$) count matrix, and a community-topic ($N \times K$) count matrix. From these count matrices, we can estimate the topic-user distributions $P(u_m|z_k)$, topic-word distributions $P(d_n|z_k)$ and community-topic $P(z_k|c_c)$ by:

$$\begin{aligned} P(u_m|z_k) &= \frac{C_{mk}^{UZ} + 1}{\sum_{m'} C_{m'k}^{UZ} + M}, \\ P(d_n|z_k) &= \frac{C_{nk}^{DZ} + 1}{\sum_{n'} C_{n'k}^{DZ} + V}, \\ P(z_k|c_c) &= \frac{C_{ck}^{CZ} + 1}{\sum_{k'} C_{ck'}^{CZ} + K}, \quad (5) \end{aligned}$$

where $P(u_m|z_k)$ is the probability of containing user m in topic k , $P(d_n|z_k)$ is the probability of using word n in topic k , and $P(z_k|c_c)$ is the probability of topic k in community c . The estimation of parameters by Gibbs sampling replaces the random seeding in EM's initialization step.

Expectation-Maximization algorithm

The CCF model is parameterized by $P(z|c)$, $P(u|z)$, and $P(d|z)$, which are estimated using the EM algorithm to fit the training corpus with community, user, and description by maximizing the log-likelihood function:

$$L = \sum_{c,u,d} n(c, u, d) \log P(c, u, d), \quad (6)$$

$$n(c, u, d) = n(c, u)n(c, d) = \begin{cases} R & \text{if community } c \text{ has user } u \\ & \text{and contains word } d \text{ for } R \text{ times;} \\ 0 & \text{otherwise.} \end{cases} \quad (7)$$

Starting with the initial parameter values from Gibbs sampling, the EM procedure iterates between Expectation (E) step and Maximization (M) step:

- **E-step:** where the probability that a community c has user u and contains word d explained by the latent variable z is estimated as:

$$P(z|c, u, d) = \frac{P(u|z)P(d|z)P(z|c)}{\sum_{z'} P(u|z')P(d|z')P(z'|c)}, \quad (8)$$

- **M-step:** where the parameters $P(u|z)$, $P(d|z)$, and $P(z|c)$ are re-estimated to maximize L in Equation (6):

$$P(u|z) = \frac{\sum_{c,d} n(c, u, d)P(z|c, u, d)}{\sum_{c,u',d} n(c, u', d)P(z|c, u', d)}, \quad (9)$$

$$P(d|z) = \frac{\sum_{c,u} n(c, u, d)P(z|c, u, d)}{\sum_{c,u,d'} n(c, u, d')P(z|c, u, d')}, \quad (10)$$

$$P(z|c) = \frac{\sum_{u,d} n(c, u, d)P(z|c, u, d)}{\sum_{u,d,z'} n(c, u, d)P(z'|c, u, d)}. \quad (11)$$

We analyze the computational complexity of the E-step and the M-step. In the E-step, one needs to compute the posterior probability $P(z|c, u, d)$ for M users, N communities, and V words. Each $P(z|c, u, d)$ consists of K values, and requires a constant number of arithmetic operations to be computed, resulting in $O(K \cdot N \cdot M \cdot V)$ operations for a single E-step. In the M-step, the posterior probabilities are accumulated to form the new estimates for $P(u|z)$, $P(d|z)$ and $P(z|c)$. Thus, the M-step also requires $O(K \cdot N \cdot M \cdot V)$ operations. Typical values of K in our experiments range from 28 to 256. The community-user (c, u) and community-description (c, d) co-occurrences are highly sparse, where $n(c, u, d) = n(c, u) \times n(c, d) = 0$ for a large percentage of the triples (c, u, d) . Because the $P(z|c, u, d)$ term is never separated from the $n(c, u, d)$ term in the M-step, we do not need to compute $P(z|c, u, d)$ for $n(c, u, d) = 0$ in the E-step. We compute only $P(z|c, u, d)$ for $n(c, u, d) \neq 0$. This greatly reduces computational complexity.

2.2.2 Parallelization

The parameter estimation using Gibbs sampling and the EM algorithm described in the previous sections can be divided into parallel subtasks [9].

Parallel Gibbs sampling

We distribute the computation among machines based on community IDs. Thus, each machine i only deals with a specified subset of communities c_i , and aware of all users u and descriptions d . We then perform Gibbs sampling simultaneously on each machine independently. For each machine, given the current state of all but one variable $z_{i,j}$, the posterior probability is sampled using Equation (4). We then calculate a number of local count variables, including $C_{m_i k}^{UZ}$, and $C_{n_i k}^{DZ}$. After a single pass through the data on each machine, we perform a global update to merge back to a single set of count variables: $C_{m k}^{UZ}$ and $C_{n k}^{DZ}$. We summarize the process in Algorithm 1.

Parallel EM algorithm

The parallel EM algorithm can be applied in a similar fashion. We describe the procedure below and provide the pseudocode in Algorithm 2.

Algorithm 1: Parallel Gibbs Sampling of CCF

Input: $N \times M$ community-user matrix
 $N \times V$ community-description matrix;
 I : number of iterations
 P : number of machines
Output: $P(u|z)$, $P(d|z)$, $P(z|c)$
Variables:
 x_{ic} : the i^{th} row of comm-user matrix with comm id c
 y_{ic} : the i^{th} row of comm-word matrix with comm id c

for $i = 0$ **to** $N - 1$ **do**
 Load x_{ic} into machine $c\%P$
 Load y_{ic} into machine $c\%P$
end
Gibbs sampling initialization
for $iter = 0$ **to** $I - 1$ **do**
 foreach $\langle user, word \rangle$ *pair* **do**
 Each machine i executes Gibbs sampling as in Equation (4)
 Each machine i calculates individual counts $C_{m_i k}^{UZ}$, $C_{n_i k}^{DZ}$ and $C_{c_i k}^{CZ}$
 end
 Master machine computes the followings and broadcasts to each machine:
 $C_{m k}^{UZ} = C_{m k}^{UZ} + \sum_i (C_{m_i k}^{UZ} - C_{m k}^{UZ})$
 $C_{n k}^{DZ} = C_{n k}^{DZ} + \sum_i (C_{n_i k}^{DZ} - C_{n k}^{DZ})$
end

- **E-step:** Each machine i computes the $P(z|c_i, u, d)$ values, the posterior probability of the latent variables z given communities c_i , users u and descriptions d , using the current values of the parameters $P(z|c_i)$, $P(u|z)$ and $P(d|z)$. As this posterior computation can be performed locally, we avoid the need for communications between machines in the E-step.
- **M-step:** Each machine i recomputes the parameters $P(z|c_i)$, $P(u|z)$ and $P(d|z)$ using the previously calculated values $P(z|c_i, u, d)$. This requires communication between machines, as $P(u|z)$ and $P(d|z)$ must be broadcasted and distributed among the machines. Since some intermediate parameters are pre-computed locally, and are subsequently used to compute the final values, the communication overhead can be dramatically reduced.

We analyze the computational and communication complexities for both algorithms using distributed machines. Assuming that there are P machines, the computational complexity of each training algorithm reduces to $O((K \cdot N \cdot M \cdot L)/P)$ (for Gibbs) and $O((K \cdot N \cdot M \cdot V)/P)$ (for EM) since P machines share the computations simultaneously. For communication complexity, two variables need to be distributed among P machines for computations in next iteration: $C_{m k}^{UZ}$, $C_{n k}^{DZ}$ in Gibbs sampling, and $P(u|z)$, $P(d|z)$ in EM. Thus broadcasting will take up to $O(P \cdot K \cdot (M + V))$.

2.2.3 Inference

Once we have learned the model parameters, we can infer three relationships using Bayesian rules, namely user-community relationship, community similarity, and user similarity. We derive these three relationships as follows:

- **User-community relationship:** Communities can

Algorithm 2: Parallel EM algorithm of CCF

Input: $N \times M$ community-user matrix; $N \times V$ community-description matrix; K : number of topics; I : number of iterations; P : number of machines; $P(u|z)$, $P(d|z)$, $P(z|c)$ of Gibbs sampling

Output: $P(u|z)$, $P(d|z)$, $P(z|c)$

Variables:
 x_{ic} : the i^{th} row of comm-user matrix with comm id c
 y_{ic} : the i^{th} row of comm-word matrix with comm id c

Load $P(u|z)$, $P(d|z)$, $P(z|c)$ of Gibbs sampling

for $i = 0$ **to** $N - 1$ **do**
 Load x_{ic} into machine $c\%P$
 Load y_{ic} into machine $c\%P$
end

for $iter = 0$ **to** $I - 1$ **do**
 for $k = 0$ **to** $K - 1$ **do**
 E-step:
 Each machine i computes $P(z_k|u, c_i, d)$
 M-step:
 Each machine i computes $P(z_k|c_i)$, $P(u_i|z_k)$, $P(d_i|z_k)$
 Master machine computes the followings and broadcasts to each machine:
 $P(u|z_k) = \sum_i P(u_i|z_k)$,
 $P(d|z_k) = \sum_i P(d_i|z_k)$
 end
end

be ranked for a given user according to $P(c_j|u_i)$, *i.e.* which communities should be recommended for a given user? Communities with top ranks and communities that the user has not yet joined are good candidates for recommendations. $P(c_j|u_i)$ can be calculated using Equation (12):

$$\begin{aligned} P(c_j|u_i) &= \frac{\sum_z P(c_j, u_i, z)}{P(u_i)} \\ &= \frac{P(c_j) \sum_z P(u_i|z)P(z|c_j)}{P(u_i)} \\ &\propto \sum_z P(u_i|z)P(z|c_j), \end{aligned} \quad (12)$$

where we assume that $P(c_j)$ is a uniform prior for simplicity.

- **Community similarity:** Communities can also be ranked for a given community according to $P(c_j|c_i)$, *i.e.* which communities should be recommended for a given community? We calculate $P(c_j|c_i)$ using Equation (13):

$$\begin{aligned} P(c_j|c_i) &= \frac{\sum_z P(c_j, c_i, z)}{P(c_i)} \\ &= \frac{\sum_z P(c_j|z)P(c_i|z)P(z)}{P(c_i)} \\ &= P(c_j) \sum_z \frac{P(z|c_j)P(z|c_i)}{P(z)} \\ &\propto \sum_z \frac{P(z|c_j)P(z|c_i)}{P(z)}, \end{aligned} \quad (13)$$

where we assume that $P(c_j)$ is a uniform prior for simplicity.

- **User similarity:** Users can be ranked for a given user according to $P(u_j|u_i)$, *i.e.* which users should be recommended for a given user? Similarly, we can calculate $P(u_j|u_i)$ using Equation (14):

$$\begin{aligned} P(u_j|u_i) &= \frac{\sum_z P(u_j, u_i, z)}{P(u_i)} \\ &= \frac{\sum_z P(u_j|z)P(u_i|z)P(z)}{P(u_i)} \\ &= P(u_j) \sum_z \frac{P(z|u_j)P(z|u_i)}{P(z)} \\ &\propto \sum_z \frac{P(z|u_j)P(z|u_i)}{P(z)}, \end{aligned} \quad (14)$$

where we assume that $P(u_j)$ is a uniform prior for simplicity.

3. EXPERIMENTAL RESULTS

We divided our experiments into two parts. The first part was conducted on a relatively small synthetic dataset with ground truth to evaluate the Gibbs & EM hybrid training strategy. The second part was conducted on a large, real-world dataset to test out CCF's performance and scalability. Our experiments were run on up to 200 machines at our distributed data centers. While not all machines are identically configured, each machine is configured with a CPU faster than 2GHz and memory larger than 4GBytes.

3.1 Gibbs + EM vs. EM

To precisely account for the benefit of Gibbs & EM over the EM-only training strategy, we used a synthetic dataset where we know the ground truth. The synthetic dataset consists of 5,000 documents with 10 topics, a vocabulary size 10,000, and a total of 50,000,000 word tokens. The true topic distribution over each document was pre-defined manually as the ground truth. We conducted the comparisons using the following two training strategies: (1) EM-only training (without Gibbs sampling as initialization) where the number of EM iterations is 10 through 70 respectively, (2) Gibbs+EM training where the number of Gibbs sampling iterations is 5, 10, 15 and 20, and the number of EM iterations is 10 through 70, respectively. We used Kullback-Leibler divergence (K-L divergence) to evaluate model performance since the K-L divergence is a good measure for the difference between the true topic distribution (P) and the estimated topic distribution (Q) defined as follows:

$$D_{KL}(P||Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}. \quad (15)$$

The smaller the K-L divergence is, the better the estimated topic distribution approximates the true topic distribution.

Figure 2 compares the average K-L divergences over 10 runs. It shows that more rounds of Gibbs sampling can help EM reach a solution that enjoys a smaller K-L divergence. Since each iteration of Gibbs sampling takes longer than EM, we must also consider *time*. Figure 3 shows the values of K-L divergence as a function of the training time. We can make two observations. First, given a large amount of time,

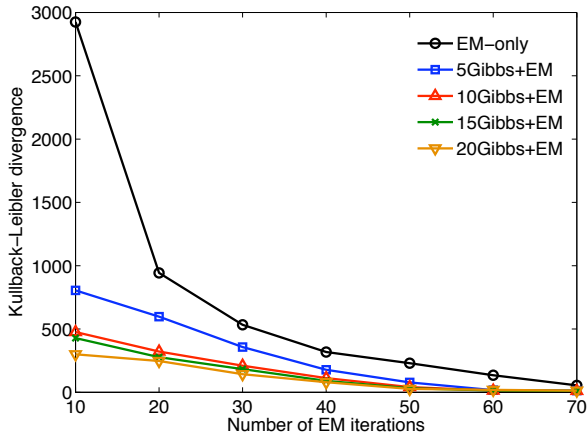


Figure 2: The Kullback-Leibler divergence as a function of the number of iterations.

both EM and the hybrid scheme can reach very low K-L divergence. On this dataset, when the training time exceeded 350 seconds, the value of K-L divergence approached zero for all strategies. Nevertheless, on a large dataset, we cannot afford a long training time, and the Gibbs & EM hybrid strategy provides a earlier point to stop training, and hence reduces the overall training time.

The second observation is on the number of Gibbs iterations. As shown in both figures, running more iterations of Gibbs before handing over to EM takes longer to yield a better initial point for EM. In other words, spending more time in the Gibbs stage can save time in the EM stage. Figure 3 shows that the best performance was produced by 10 iterations of Gibbs sampling before switching to EM. Finding the “optimal” switching point is virtually impossible in theory. However, the figure shows that different Gibbs iterations can all outperform the EM-only strategy to obtain a better solution early, and a reasonable number of Gibbs iterations can be obtained through an empirical process like our experiment. Moreover, the figure shows that a range of number of iterations can achieve similar K-L divergence (e.g., at time 250). This indicates that though an empirical process may not be able to pin down the “optimal” number of iterations (because of e.g., new training data arrival), the hybrid scheme can work well on a range of Gibbs-sampling iterations.

3.2 The Orkut Dataset

Orkut is an extremely active community site with more than two billion page views a day world-wide. The dataset we used was collected on July 26, 2007, which contains two types of data for each community: community membership information and community description information. We restrict our analysis to English communities only. We collected 312,385 users and 109,987 communities². The number of entries in the community-user matrix, effectively, the number of community-user pairs, is 35,932,001. As the density is around 0.001045, this matrix is extremely sparse. Figure 4(a) shows a distribution of the number of users per community. About 52% of all communities have less than 100

²All user data were anonymized, and user privacy is safeguarded, as performed in [10].

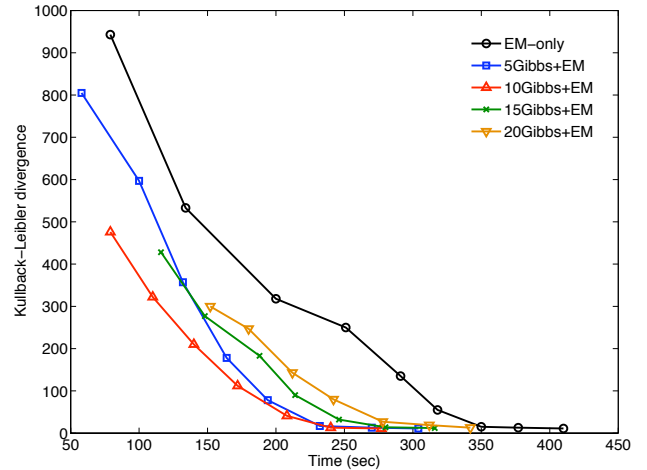


Figure 3: The Kullback-Leibler divergence as a function of the training time.

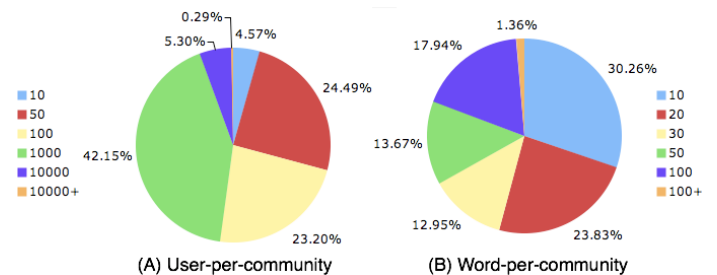


Figure 4: (a) Distribution of the number of users per community, and (b) distribution of the number of description words per community.

users, whereas 42% of all communities have more than 100 but less than 1,000 users.

For the community description data, after applying down-casing, stopword filtering, and word stemming, we obtained a vocabulary of 191,034 unique English words. The distribution of the number of description words per community is displayed in Figure 4(b). On average, there are 27.64 words in each community description after processing. In order to establish statistical significance of the findings, we repeated all experiments 10 times with different random seeds and parameters, such as the number of latent aspects (ranging from 28 to 256), the number of Gibbs sampling iterations (ranging from 10 to 30) and the number of EM iterations (ranging from 100 to 500). The reported results are the average performance over all runs.

Results

Community Recommendation: $P(c_j|u_i)$

We use two standard measures from information retrieval to measure the recommendation effectiveness: *precision* and *recall*, defined as follows:

$$Precision = \frac{|\{recommendation\ list\} \cap \{joined\ list\}|}{|\{recommendation\ list\}|},$$

$$Recall = \frac{|\{recommendation\ list\} \cap \{joined\ list\}|}{|\{joined\ list\}|}. \quad (16)$$

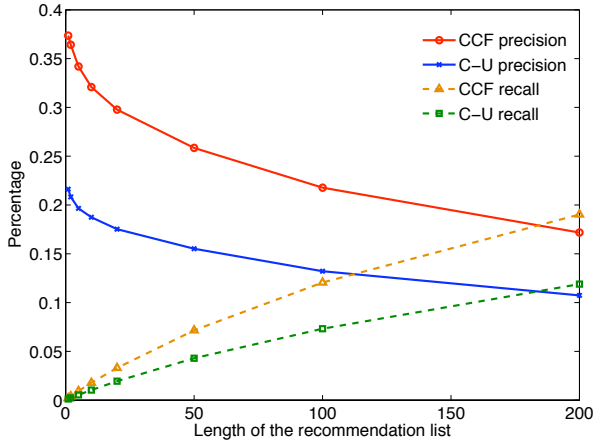


Figure 5: The precision and recall as functions of the length (up to 200) of the recommendation list.

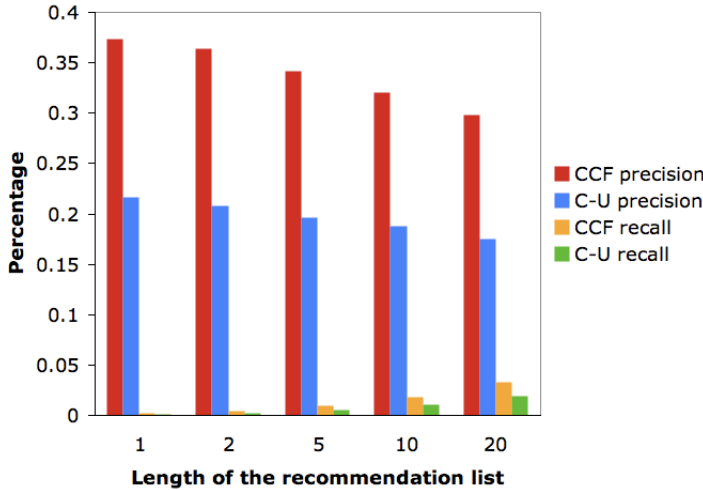


Figure 6: The precision and recall as functions of the length (up to 20) of the recommendation list.

Precision takes all recommended communities into account. It can also be evaluated at a given cut-off rank, considering only the topmost results recommended by the system. As it is possible to achieve higher recall by recommending more communities (note that a recall of 100% is trivially achieved by recommending all communities, albeit at the expense of having low precision), we limit the size of our community recommendation list to at most 200.

To evaluate the results, we randomly deleted one joined community for each user in the community-user matrix from the training data. We evaluated whether the deleted community could be recommended. This evaluation is similar to *leave-one-out*. Figure 5 shows the precision and recall as functions of the length (up to 200) of the recommendation list for both C-U and CCF. We can see that CCF always outperforms C-U for all lengths. Figure 6 presents precision and recall for the top 20 recommended communities. As both precision and recall of CCF are nearly twice higher than those of C-U, we can conclude that CCF enjoys better prediction accuracy than C-U. This is because C-U only considers community-user co-occurrence, whereas CCF con-

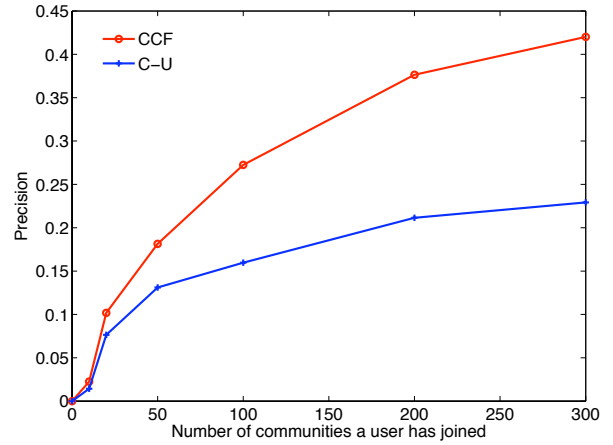


Figure 7: The precision as a function of the number of communities a user has joined. Here, the length of the recommendation list is fixed at 20.

Table 1: The comparison results of the three models using Normalized Mutual Information (NMI).

| Model | C-U | C-D | CCF |
|-------|--------|--------|--------|
| NMI | 0.4508 | 0.3127 | 0.4526 |

siders users, communities, and descriptions. By taking into other views into consideration, the information is denser for CCF to achieve higher prediction accuracy.

Figure 7 depicts the relationship between the precision of the recommendation for a user and the number of communities that the user has joined. The more communities a user has joined, the better both C-U and CCF can predict the user’s preferences. For users who joined around 100 communities, the precision is about 15% for C-U and 27% for CCF. However, for users who joined just 20 communities, the precision is about 7% for C-U, and 10% for CCF. This is not surprising since it is very difficult for latent-class statistical models to generalize from sparse data. For large-scale recommendation systems, we are unlikely to ever have enough direct data with sufficient coverage to avoid sparsity. However, at the very least, we can try to incorporate indirect data to boost our performance, just as CCF does by using bags of words information to augment bags of users information. *Remark:* Because of the nature of leave-one-out, our experimental result can only show whether a joined community could be recovered. The low precision/recall reflects this necessary, restrictive experimental setting. (This setting is necessary for objectivity purpose as we cannot obtain ground-truth of all users’ future preferences.) The key observation from this study is not the absolute precision/recall values, but is the relative performance between CCF and C-U.

Community Similarity: $P(c_j|c_i)$

We next report the results of community similarities calculated by the three models. We used *community category* (available at Orkut websites) as the ground-truth for clustering communities. We also assigned each community an estimated label for the latent aspect with the highest probability value. We treated communities with the same es-

Table 2: The top recommended users using the C-U and CCF models for the query user “79”. The number of communities that “79” joined is 339. Note that the “Communities” field contains three numbers: the first number n is the total number of communities a user joined; the second number k is the number of overlapping communities between the recommended user and the query user, and the last number is percentage of $\frac{k}{n}$.

| | Rank 1 st | | Rank 2 nd | | Rank 3 rd | |
|-------|----------------------|------------------|----------------------|------------------|----------------------|-----------------|
| Model | User ID | Communities | User ID | Communities | User ID | Communities |
| C-U | 2390 | 551 (102, 18.5%) | 8207 | 456 (100, 21.9%) | 6734 | 494 (95, 19.2%) |
| CCF | 7931 | 518 (106, 20.5%) | 10968 | 680 (102, 15.0%) | 6776 | 680 (91, 13.4%) |

timated label as members of the same *community cluster*. We then compared the difference between community clusters and categories using the Normalized Mutual Information (NMI).

NMI between two random variables CAT (category label) and CLS (cluster label) is defined as $NMI(CAT; CLS) = \frac{I(CAT; CLS)}{\sqrt{H(CAT)H(CLS)}}$, where $I(CAT; CLS)$ is the mutual information between CAT and CLS . The entropies $H(CAT)$ and $H(CLS)$ are used for normalizing the mutual information to be in the range $[0, 1]$. In practice, we made use of the following formulation to estimate the NMI score [12, 13]:

$$NMI = \frac{\sum_{s=1}^K \sum_{t=1}^K n_{s,t} \log \left(\frac{n \cdot n_{s,t}}{n_s \cdot n_t} \right)}{\sqrt{\left(\sum_s n_s \log \frac{n_s}{n} \right) \left(\sum_t n_t \log \frac{n_t}{n} \right)}}, \quad (17)$$

where n is the number of communities, n_s and n_t denote the numbers of community in category s and cluster t , $n_{s,t}$ denotes the number of community in category s as well as in cluster t . The NMI score is 1 if the clustering results perfectly match the category labels and 0 for a random partition. Thus, the larger this score, the better the clustering results.

Table 1 shows that CCF slightly outperforms both C-U and C-D models, which indicates the benefit of incorporating two types of information.

User Similarity: $P(u_j|u_i)$

An interesting application is friend suggestion: finding users similar to a given user. Using Equation (14), we can compute user similarity for all pairs of users. From these values, we derive a ranking of the most similar users for a given query user. Due to privacy concerns, we were not able to obtain the friend graph of each user to evaluate accuracy. Table 2 shows an example of this ranking for a given user.

“Similar” users typically share a significant percentage of commonly-joined communities. For instance, the query user also joined 18.5% of the communities joined by the top user ranked by C-U, compared to 20.5% for CCF. It is encouraging to see that CCF’s top ranked user has more overlap with the query user than C-U’s top ranked user does. We believe that, again, incorporating the additional word co-occurrences has improved information density and hence yields higher prediction accuracy.

3.3 Runtime Speedup

In analyzing runtime speedup for parallel training, we trained CCF with 20 latent aspects, 10 Gibbs sampling, and 20 EM iterations. As the size of a dataset is large, a single machine cannot store all the data— $P(u|z)$, $P(d|z)$, $P(z|c)$, and $P(z|c, u, d)$ —in its local memory, we cannot obtain the

Table 3: Runtime comparisons for different number of machines.

| Machines | Time (sec.) | Speedup |
|----------|-------------|---------|
| 10 | 9,233 | 10 |
| 20 | 4,326 | 21.3 |
| 50 | 2,280 | 40.5 |
| 100 | 1,014 | 91.1 |
| 200 | 796 | 116 |

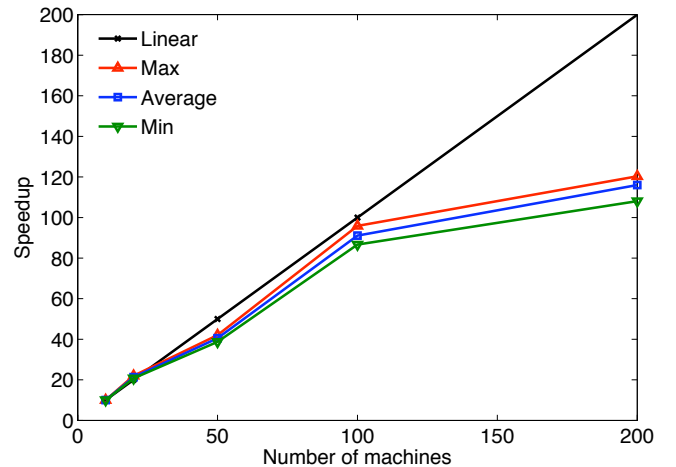


Figure 8: Speedup analysis for different number of machines.

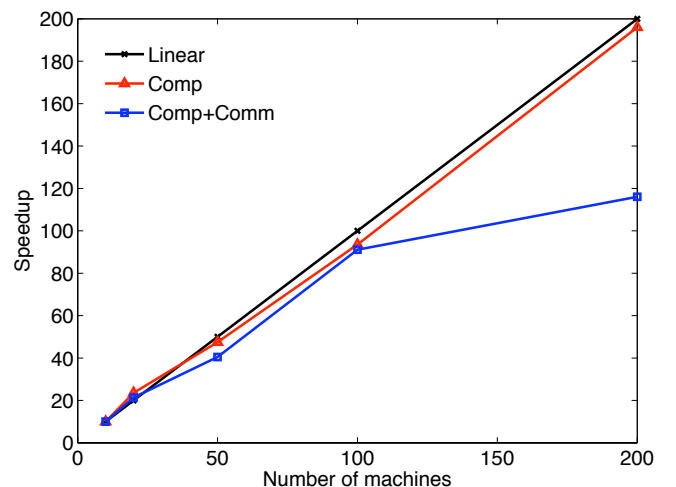


Figure 9: Speedup and overhead analysis.

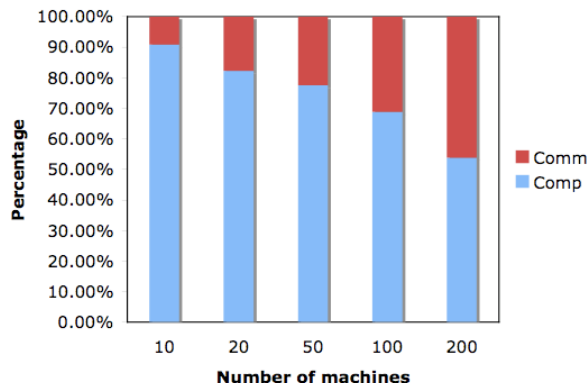


Figure 10: Runtime (Computation and Communication) composition analysis.

running time of CCF on one machine. Therefore, we use the runtime of 10 machines as the baseline and assume that 10 machines can achieve 10 times speedup. This assumption is reasonable as we will see shortly that our parallelization scheme can achieve linear speedup on up to 100 machines. Table 3 and Figure 8 report the runtime speedup of CCF using up to 200 machines. The Orkut dataset enjoys a linear speedup when the number of machines is up to 100. After that, adding more machines receives diminishing returns. This result led to our examination of overheads for CCF, presented next.

No parallel algorithm can infinitely achieve linear speedup because of the Amdahl’s law. When the number of machines continues to increase, the communication cost starts to dominate the total running time. The running time consists of two main parts: computation time (Comp) and communication time (Comm). Figure 9 shows how Comm overhead influences the speedup curves. We draw on the top the computation only line (Comp), which approaches the linear speedup line. The speedup deteriorates when communication time is accounted for (Comp + Comm). Figure 10 shows the percentage of Comp and Comm in the total running time. As the number of machines increases, the communication cost also increases. When the number of machines exceeds 200, the communication time becomes even larger than the computation time.

Though the Amdahl’s law eventually kicks in to forbid a parallel algorithm to achieve infinite speedup, our empirical study draws two positive observations.

1. When the dataset size increases, the “saturation” point of the Amdahl’s law is deferred, and hence we can add more machines to deal with larger sets of data.
2. The speedup that can be achieved by parallel CCF is very significant to enable near-real-time recommendations. As shown in the table, the parallel scheme reduces the training time from one day to less than 14 minutes. The parallel CCF can be run every 14 minutes to produce a new model to adapt to new access patterns and new users.

4. CONCLUDING REMARKS

We have introduced a generative graphical model, Combinational Collaborative Filtering (CCF), for collaborative filtering based on both bags of words and bags of users information. CCF uses a hybrid training strategy that combines

Gibbs sampling with the EM algorithm. The model trained by Gibbs sampling provides better initialization values for EM than random seeding. We also presented the parallel computing required to handle large-scale data sets. Experiments on a large Orkut data set demonstrate that our approaches successfully produce better quality recommendations, and accurately cluster relevant communities/users with similar semantics.

There are a couple of directions for future research. First, we would consider expanding CCF to incorporate more types of co-occurrence data. More types of co-occurrence data would help to overcome sparsity problem and make better recommendation. Second, in our analysis, the community-user pair value equals one, *i.e.* $n(u_i, c_j) = 1$ (if user u_i joins community c_j). An interesting extension would be to give this count a different value, *i.e.* $n(u_i, c_j) = f$, where f is the frequency of the user u_i visiting the community c_j . We are currently parallelizing LDA and will compare LDA and PLSA as the choice of our baseline algorithm in the future.

5. ACKNOWLEDGEMENT

The authors would like to thank Ellen Spertus for preparing the Orkut dataset and Jon Chu for helpful discussions. The first author is supported by NSF under grant II-0535085.

6. REFERENCES

- [1] Alexa internet. <http://www.alexacom/>.
- [2] D. M. Blei and M. I. Jordan. Variational methods for the Dirichlet process. In *Proc. of the 21st ICML Conference*, pages 373–380, 2004.
- [3] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.
- [4] D. Cohn and H. Chang. Learning to probabilistically identify authoritative documents. In *Proc. of the 17th ICML Conference*, pages 167–174, 2000.
- [5] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1):1–38, 1977.
- [6] S. Geman and D. Geman. Stochastic relaxation, gibbs distributions, and the bayesian restoration of images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6:721–741, 1984.
- [7] T. Hofmann. Probabilistic latent semantic analysis. In *Proc. of the 15th UAI Conference*, pages 289–296, 1999.
- [8] A. McCallum, A. Corrada-Emmanuel, and X. Wang. The author-recipient-topic model for topic and role discovery in social networks: Experiments with enron and academic email. Technical report, Computer Science, University of Massachusetts Amherst, 2004.
- [9] D. Newman, A. Asuncion, P. Smyth, and M. Welling. Distributed inference for latent Dirichlet allocation. In *NIPS*, 2007.
- [10] E. Spertus, M. Sahami, and O. Buyukkocuten. Evaluating similarity measures: a large-scale study in the orkut social network. In *Proc. of the 11th ACM SIGKDD Conference*, pages 678–684, 2005.
- [11] M. Steyvers, P. Smyth, M. Rosen-Zvi, and T. Griffiths. Probabilistic author-topic models for information discovery. In *Proc. of the 10th ACM SIGKDD Conference*, pages 306–315, 2004.
- [12] A. Strehl and J. Ghosh. Cluster ensembles – a knowledge reuse framework for combining multiple partitions. *Journal on Machine Learning Research*, 3:583–617, 2002.
- [13] S. Zhong and J. Ghosh. Generative model-based clustering of documents: a comparative study. *Knowledge and Information Systems*, 8:374–384, 2005.