

On Sampling Top-K Recommendation Evaluation

Dong Li
dli12@kent.edu
Kent State University

Ruoming Jin
rjin1@kent.edu
Kent State University

Jing Gao
jgao@ilambda.com
iLambda

Zhi Liu
zliu@ilambda.com
iLambda

ABSTRACT

Recently, Rendle has warned that the use of sampling-based top- k metrics might not suffice. This throws a number of recent studies on deep learning-based recommendation algorithms, and classic non-deep-learning algorithms using such a metric, into jeopardy. In this work, we thoroughly investigate the relationship between the sampling and global top- K Hit-Ratio (HR, or Recall), originally proposed by Koren [2] and extensively used by others. By formulating the problem of aligning sampling top- k ($SHR@k$) and global top- K ($HR@K$) Hit-Ratios through a mapping function f , so that $SHR@k \approx HR@f(k)$, we demonstrate both theoretically and experimentally that the sampling top- k Hit-Ratio provides an accurate approximation of its global (exact) counterpart, and can consistently predict the correct winners (the same as indicate by their corresponding global Hit-Ratios).

CCS CONCEPTS

• **Information systems** → **Collaborative filtering; Recommender systems; Retrieval effectiveness.**

KEYWORDS

Recommender systems, top- k , hit ratio, recall, evaluation metric

ACM Reference Format:

Dong Li, Ruoming Jin, Jing Gao, and Zhi Liu. 2020. On Sampling Top-K Recommendation Evaluation. In *Proceedings of the 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '20), August 23–27, 2020, Virtual Event, CA, USA*. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3394486.3403262>

1 INTRODUCTION

Over the last few years, in both industry and academic research communities, many efforts have been taken to integrate deep learning into recommendation techniques [22]. Though the flourishing list of publications has demonstrated sizeable improvements over the classical non-deep (linear) approaches, several recent studies [3, 16] have sounded the alarm: The displayed success in recommendation may contribute to the weaker baseline. Some other factors, such as evaluation protocols and performance measures, together with choices of datasets, may also play roles in the potentially over-promising results [3].

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

KDD '20, August 23–27, 2020, Virtual Event, CA, USA

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-7998-4/20/08...\$15.00

<https://doi.org/10.1145/3394486.3403262>

Lately, Rendle [15] has noticed that, in recent deep learning-based recommendation studies, it is becoming popular [5, 7, 8, 11, 19–21] to adopt sampling-based criteria for top- k evaluation. Basically, instead of ranking all available items, which might be a very large list, these studies sample a smaller set of (irrelevant) items, and rank the relevant items against the sampled items. In his findings, he claimed that the typical process used top- k evaluation metrics, such as Recall/Precision (Hit-Ratio), Average Precision (AP) and nDCG, other than (AUC), are all “inconsistent” with respect to the exact metrics (even in expectation). He even suggests avoiding using the sampled metrics for top- k evaluation.

What does this mean to all the existing studies which use sampling top- k criteria? Have their results all become (sort-of) invalid? Does that mean we have to use all the items for any meaningful top- k evaluation? Clearly, a significant amount of efforts in the recommendation community is at risk here. To be able to firmly answer these questions, a better understanding of the sampling-based top- k metrics is much needed. In the meantime, a sampling approach, where acceptable, can be a useful tool for saving computational cost and speeding up evaluation time. While computational resources might not be a big problem for enormous mega-corporations, such as Google or Amazon, for many smaller, resource-constrained organizations, it may still be an issue. For instance, if valid, a sampling approach can be a quick way to help evaluate the promise of a given algorithm, screening for the eventual exact/global top- k evaluation.

A Bit of History on Sampling Top- k Evaluation: The sampling top- k method was initially suggested by Koren in the seminal work [10] as an approach to measure the success of top- k recommenders. Specifically, he uses 1000 additional random movies (which may include already-ranked ones) against the targeted movie i for a user. He ranks these 1001 movies by the predicted rating (relevance score), and he normalizes the ranking score between 0 and 1. Finally, he draws the cumulative distributions of all the users, with respect to the ranking score. Essentially, for any algorithm, at a given rank k , the value (the point in the cumulative distribution curve) is basically Hit-Ratio or Recall (at k) under sampling. Another highly cited work [2] has utilized this metric to evaluate the performance of a variety of recommendation algorithms on Top- N recommendation tasks.

This method was first adopted by deep learning-based recommendation papers in [6] and then in [7]. Here, the authors go beyond the top- k Hit-Ratio suggested by Koren [6, 7], extending to metrics such as Mean Reciprocal Rank (MRR) and nDCG. Since Koren, various other deep learning-based recommendation studies [5, 8, 11, 19–21] have adopted such sampling-based top- k evaluation metrics. In these studies, they typically sample only those “irrelevant” items (not scored by the users), unlike the work in Koren, which may sample relevant items, as well. The number of items sampled typically ranges from 100 to 1000.

Table 1: Notations

I	entire set of items, and its size $ I = N$
i_u	relevant item for user u in testing data
I_n	sampled item set (user-specific), composed of $n - 1$ sampled items and i_u
R	integer variable, referring to item rank position, in range $[1, N]$
R_u	rank of item i_u among I for user u
r^u	rank of item i_u among I_n for user u ; also denotes a random variable
r^R	$:= r^u$, for the group of users whose $R_u = R$
$HR@K$	global top- K hit-ratio (recall), Formula 1
$SHR@k$	sampling top- k hit-ratio (recall), Formula 3
f	mapping $SHR@k$ to $HR@K$, where $K = f(k)$
W_R	fraction of users where $R_u = R$, (Eq 2); also denotes user ranking distribution
p^u	$Pr(r^u \leq k)$, probability that $r^u \leq k$ among I_n for a user u
p^R	$:= p^u$, for the group of users whose $R_u = R$
p_u	$:= \frac{R_u - 1}{N - 1}$, probability for sampling an item that ranks higher than i_u for u

However, besides the latest study and warning by Rendle [15], there have been no studies on the statistical properties of sampling-based top- k evaluation metrics. Clearly, as indicated by Rendle [15], the sampling top- k metric is very different from global top- k metric. But do they relate to each other? Can sampling top- k reflect global (exact) top metrics? And how do we interpret the existing experimental results which use sampling-based metrics?

Our Contributions: To answer the above questions, we perform the first study to thoroughly investigate the relationship between the sampling and global top- k Hit-Ratio (HR, or Recall), originally proposed by Koren [10]. Top- k hit ratio is one of the most popular metrics used for evaluating almost all top- N recommenders [22]. Specifically, we made the following contribution:

- (Section 2) We formalize the problem of aligning sampling top- k ($SHR@k$) and global top- K ($HR@K$) Hit-Ratios through a mapping function, so that $SHR@k \approx HR@f(k)$, where f is the functions map of the k in the sampling to global top $f(k)$. We also prove the *Sampling Theorem*, which shows the sampling Hit-Ratio preserves the “dominating” property between global Hit-Ratios.
- (Sections 3 and 4) We develop novel methods to approximate function f , and we show that it is surprisingly approximately linear, even under non-linear computation (when k is large). Basically, the “sampling” location of the global top- K curve is almost equally intervalled ($f(i) - f(i - 1)$ is close to constant). In addition, we develop algorithm-specific mapping functions and discuss a list of key properties to help ensure the predictive power of sampling.
- (Section 5) We experimentally validate our mapping function f by comparing between the sampling and its global top- K counterparts, and we show that our function can provide a rather accurate estimate of its global top- $f(k)$. We also show that the sampling Hit-Ratio can accurately predict the same winners as the corresponding global Hit-Ratio.

2 PROBLEM AND SAMPLING THEOREM

2.1 Problem Formulation

Assume we split the entire dataset into training and testing. Let the testing dataset consist of M users and $N = |I|$ items, where I is the entire set of items. From the training, we will learn a recommendation algorithm A , which can rank a given set of items for a user. To simplify our discussion, we consider the *leave-one-out* strategy [1], where **each user has and only has one relevant**

item to be evaluated, though such treatment can be naturally generalized to the situation where a user may have more than one targeted item in the testing data [6, 13]. Table 1 highlights the key notations used in the rest of the paper.

2.1.1 Global Top- K Hit Ratio. Given a user u , and its relevant item i_u ($i_u \in I$) in the testing dataset, the recommendation algorithm A will calculate the relative rank of the relevant item i_u , denoted as R_u , among all available items I , $R_u = A(u, i_u, I)$.

Let us consider the global top- K Hit-Ratio (or recall) metric:

$$HR@K = \frac{1}{M} \sum_{u=1}^M \mathbf{1}_{R_u \leq K} = \sum_{R=1}^N W_R \cdot \mathbf{1}_{R \leq K} \quad (1)$$

Here, $\mathbf{1}_X$ is the indicator function of event X ($\mathbf{1}_X = 1$ iff X is true and 0 others), and W_R is the frequency of users with item i_u rank in position R :

$$W_R = \frac{1}{M} \sum_{u=1}^M \mathbf{1}_{R_u=R} \quad (2)$$

2.1.2 Sampling Top- k Hit Ratio. Next, let us revisit the top- k Hit-Ratio under sampling. For a given user u and the relevant item i_u , we first sample $n - 1$ items from the entire set of items I , forming the subset I_n (including i_u). Let the relative rank of i_u among I_n be denoted as $r^u = A(u, i_u, I_n)$. Note that r^u is a random variable depending on the sampling set I_n .

Given this, the sampling top- k Hit-Ratio can be defined as $SHR@k$:

$$SHR@k = \frac{1}{M} \sum_{u=1}^M Z^u, \quad Z^u \sim \text{Bernoulli}(p^u = Pr(r^u \leq k)) \quad (3)$$

where, Z^u is a random variable for each user u , and follows a Bernoulli distribution with probability $p^u = Pr(r^u \leq k)$.

Now, recall that we are trying to study the relation between $SHR@k$ and $HR@K$.

2.2 Sampling

We note that the population sum $\sum_{u=1}^M Z^u$ is a Poisson binomial distributed variable (a sum of M independent Bernoulli distributed variables). Its mean and variance will simply be sums of the mean and variance of the n Bernoulli distributions:

$$\mu = \sum_{u=1}^M p^u, \quad \sigma^2 = \sum_{u=1}^M p^u (1 - p^u)$$

Given this, the expectation and variance of $SHR@k$:

$$E[SHR@k] = \frac{1}{M} \sum_{u=1}^M p^u = \sum_{R=1}^N W_R \cdot p^R \quad (4)$$

$$Var[SHR@k] = \frac{1}{M^2} \sum_{u=1}^M p^u (1 - p^u) = \frac{1}{M} \sum_{R=1}^N W_R \cdot p^R (1 - p^R) \quad (5)$$

The probability for users who are in the same group ($R_u = R$), share the same p^u , will be denoted by p^R and W_R is defined in equation 2.

To define $p^u = Pr(r^u \leq k)$ precisely, let us consider the two commonly used types of sampling (with and without replacements). **Sampling with replacement (Binomial Distribution):**

For a given user u , let X^u denote the number of sampled items that are ranked in front of relevant item i_u :

$$X^u = \sum_{i=1}^{n-1} X_i^u, \quad X_i^u \sim \text{Bernoulli}(p_u = \frac{R_u - 1}{N - 1})$$

where X_i^u is a Bernoulli random variable for each sampled item i : $X_i^u = 1$ if item i has rank range in $[1, R_u - 1]$ (p_u is the corresponding probability) and $X_i^u = 0$ if i is located in $[R_u + 1, N]$. Thus, X^u follows binomial distribution:

$$X^u \sim \text{Binomial}(n - 1, p_u = \frac{R_u - 1}{N - 1}) \quad (6)$$

And the random variable $r^u = X^u + 1$, and we have

$$p^u = \text{CDF}(k; n - 1, p_u) = \Pr(r^u \leq k) = \begin{cases} \sum_{l=0}^{k-1} \binom{n-1}{l} p_u^l (1-p_u)^{n-1-l}, & R_u \geq k \\ 1, & R_u < k \end{cases}$$

Sampling without replacement (Hypergeometric Distribution): If we sample $n - 1$ items from the total $N - 1$ items without replacement, and the total number of successful cases is $R_u - 1$, then let X^u be the random variable for the number of items appearing in front of relevant item i_u ($r^u = X^u + 1$):

$$X^u \sim \text{Hypergeometric}(N - 1, R_u - 1, n - 1) \\ p^u = \text{CDF}(k; N - 1, R_u - 1, n - 1) = \Pr(r^u \leq k) = \begin{cases} \sum_{l=0}^{k-1} \frac{\binom{R_u-1}{l} \binom{N-R_u}{n-1-l}}{\binom{N-1}{n-1}}, & R_u \geq k \\ 1, & R_u < k \end{cases}$$

It is well-known that, under certain conditions, the hypergeometric distribution can be approximated by binomial distribution. We will focus on using binomial distribution for analysis, and we will validate the results on hypergeometric distribution experimentally.

2.3 A Functional View of $HR@K$ and $SHR@k$

To better understand the relationship between $HR@K$ (global top- K Hit-Ratio) and $SHR@k$ (the sampling version), it is beneficial to take a functional view of them. Let \mathcal{R} be the random variable for the user's item rank, with probability mass function $\Pr(\mathcal{R} = R)$; then, $HR@K$ is simply the empirical cumulative distribution of \mathcal{R} (\widehat{Pr}):

$$HR@K = \widehat{Pr}(\mathcal{R} \leq K), \quad W_R = \widehat{Pr}(\mathcal{R} = R) \quad (7)$$

For $SHR@k$, its direct meaning is more involved and will be examined below. For now, we note that $SHR@k$ is a function of k varying from 1 to n , where $n - 1$ is the number of sampled items.

Figure 1a displays the curves of functional fitting of empirical accumulative distribution $HR@K$ (aka the global top- K Hit-Ratio, varying K from 1 to $N = 3706$), for 6 representative recommendation algorithms (3 classical and 3 deep learning methods), on the MovieLens 1M dataset. To observe the performance of these methods more closely when the K is small, we first highlight K from 1 to 350, and then again from 60 to 120.

Figure 1b displays the curves of functional fitting of function $SHR@K$ (the sampling top- k Hit-Ratio, varying k from 1 to $n = 100$) with $n - 1$ samples, under sampling with replacement, for the same 6 representative recommendation algorithms on the same dataset. Similarly, we highlight k from 1 to 10, and then again from 2 to 5.

How can the sampling Hit-Ratio curves help to reflect what happened in the global curves? Before we consider the more detailed relationship between them, we introduce the following results:

THEOREM 2.1 (SAMPLING THEOREM). *Let us assume we have two global Hit-Ratio curves (empirical cumulative distribution), $HR_1@K$ and $HR_2@K$, and assume one curve dominates the other one, i.e., $HR_1@K \geq HR_2@K$ for any $1 \leq K \leq N$; then, for their corresponding sampling curve at any k for any size of sampling, we have*

$$E(SHR_1@k) \geq E(SHR_2@k)$$

PROOF. Recall Equation 4: $E(SHR@k) = \sum_{R=1}^N W_R \cdot p^R = \sum_{u=1}^M \Pr(r^u \leq k)$. Let us assign each user u the weight $\Pr(r^u \leq k)$ for both curves, HR_1 and HR_2 . Now, let us build a bipartite graph by connecting any u in the HR_1 with user v in HR_2 , if $R_u \leq R_v$. We can then apply Hall's marriage theorem to claim there is a one-to-one matching between users in HR_1 to users in HR_2 , such that $R_u \leq R_v$, and $\Pr(r^u \leq k) \geq \Pr(r^v \leq k)$. (To see that, use the fact that $\sum_{R=1}^K W_R^{(1)} \geq \sum_{R=1}^K W_R^{(2)}$, where $W_R^{(1)}$ and $W_R^{(2)}$ are the empirical probability mass distributions of user-ranks, or equivalently, $\sum_{R=K}^N W_R^{(1)} \leq \sum_{R=K}^N W_R^{(2)}$. Thus, any subset in HR_1 is always smaller than its neighbor set $N(HR_1)$ in (HR_2) . Given this, we can observe that the theorem holds. \square

The above theorem shows that, under the strict order of global Hit-Ratio curves (though it may be quite applicable for searching/evaluating better recommendation algorithms, such as in Figure 1), sampling hit ratio curves can maintain such order.

However, this theorem does not explain the stunning similarity, shapes and trends shared by the global and their corresponding sampling curves. Basically, the detailed performance differences among different recommendation algorithms seem to be well-preserved through sampling. However, unless $n \approx N$, $SHR@k$ does not correspond to $HR@K$ (as in what is being studied by Rendel [15]).

Those observations hold on other datasets and recommendation algorithms as well, not only on this dataset. Thus, intuitively and through the above experiments, we may conjecture that it is the overall curve $HR@K$ that is being approximated by $SHR@k$. Since these functions are defined on different domain sizes N vs n , we need to define such approximation carefully and rigorously.

2.4 Mapping Function f

To explain the similarity between the global and sampling top- k Hit-Ratio curves, we hypothesize that there exists a function $f(k)$ such that the relation $SHR@k \approx HR@f(k)$ holds for different ranking algorithms on the same dataset. In a way, the sampling metric $SHR@k$ is like "signal sampling" [14], where the global metrics between top 1 to N are sampled (and approximated) at only $f(1) < f(2) < \dots < f(n)$ locations, which corresponds to $SHR@k$ ($k = 1, 2, \dots, n$). In general, $f(k) \neq k$ (when $n < N$) ([15]).

In order to identify such a mapping function, let us take a look at the error between $SHR@k$ and $HR@f(k)$: $|SHR@k - HR@f(k)|$

$$\leq |SHR@k - E[SHR@k]| + |E[SHR@k] - HR@f(k)| \quad (8)$$

Thanks to the Hoeffding's bound, we observe,

$$\Pr(|SHR@k - E[SHR@k]| \geq t) \leq 2 \exp(-2Mt^2)$$

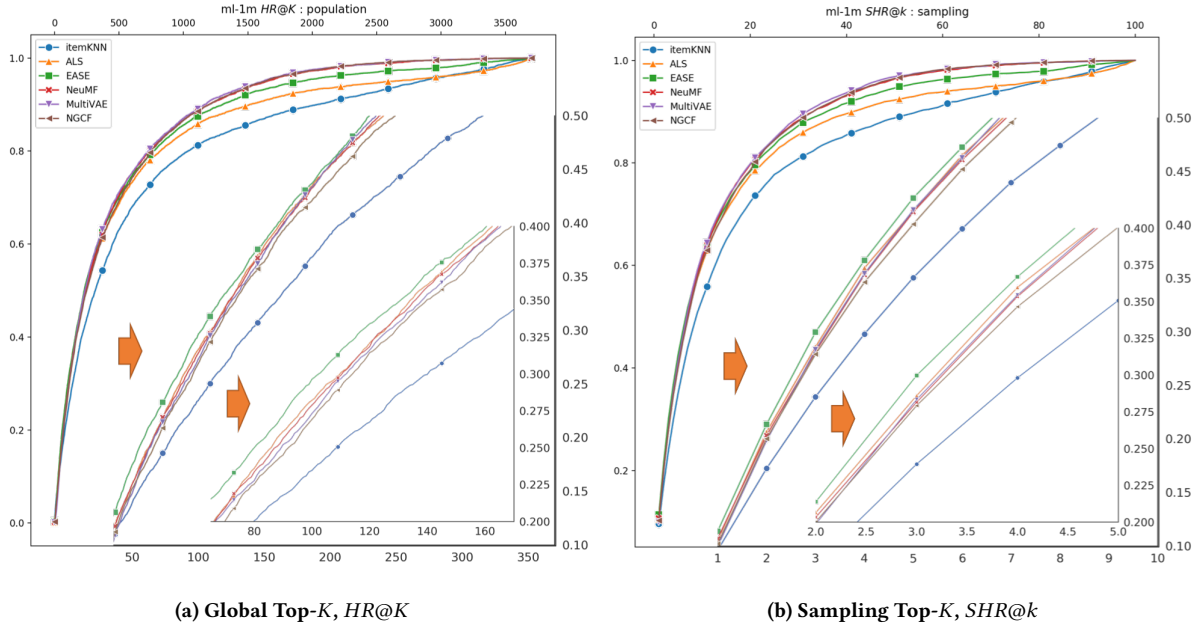


Figure 1: Global vs Sampling Top k Hit-Ratio on MovieLens 1M dataset (ml-1m). To display the details clearly, we zoom in at different range scales. Compare two figures, we can easily conclude that sampling evaluation maintains the same trend as global evaluation for different algorithms even at small error range.

This can be a rather tight bound, due to the large number of users in the population. For example, if $M = 30K$, $t = 0.01$:

$$\Pr(|SHR@k - E[SHR@k]| \geq 0.01) \leq 0.005$$

If we want to look at more closely, we may use the law of large numbers and utilize the variance in equation 5 for deducing the difference between $SHR@k$ and its expectation. Overall, for a large user population, the sampling top- k Hit-Ratio will be tightly centered around its mean. Furthermore, if the user number is, indeed, small, an average of multiple sampling results can reduce the variance and error. In the publicly available datasets, we found that one set of samples is typically very close to the average of multiple runs.

Given this, our problem is **how to find the mapping function f , such as $|E[SHR@k] - HR@f(k)|$ can be minimized (ideally close to or equal to 0**. Note that f should work for all the k (from 1 to n), and it should be independent of algorithms on the same dataset.

3 APPROXIMATING MAPPING FUNCTION f

Baseline: To start, we may consider the following naive mapping function. We notice that for any n ,

$$E(X^u) = (n-1)p_u = (n-1)\frac{R_u - 1}{N-1} = E(r_u) - 1$$

When the sample n is large, we simply use the indicator function $1_{E(r_u) \leq k}$ to approximate and replace $\Pr(r_u \leq k)$. Thus,

$$R_u \leq \frac{k-1}{n-1} * (N-1) + 1 = f(k) \quad (9)$$

Since the indicator function ($1_{E(r_u) \leq k}$) is a rather crude estimation of the CDF of r_u at k , this only serves as a baseline for our approximation of the mapping function f .

Approximation Requirements: Before we introduce more carefully designed approximations of the mapping function f , let us take a close look of the expectation of the sampling top- k Hit-Ratio $E(SHR@k)$ and $HR@f(k)$. Figure 2 shows how the user probability mass function W_R works with the step function (indicator function) $1_{R_u \leq f(k)}$, and p_u (assuming a hypergeometric distribution), to generate the global top- K and sampling Hit-Ratios.

We make the following observations (as well as requirements): *Existence of mapping function f for each individual HR curve:* Given any k , assuming $HR@f(k)$ is a continuous cumulative distribution function (i.e., assuming that there is no jump/discontinuity on the CDF, and that $f(k)$ is a real value), then, there is $f(k)$ such that $HR@f(k) = E(SHR@k)$.

In our problem setting, where $f(k)$ is integer-valued and ranges between 0 and N , the best $f(k)$, theoretically, is

$$f(k) = \arg \min |HR@f(k) - E(SHR@k)|$$

Mapping function f for different HR curves: Since our main purpose is for $SHR@k$ to be comparable across different recommendation algorithms, we prefer $f(k)$ to be the same for different HR curves (on the same dataset). Thus, by comparing different $SHR@k$, we can infer their corresponding Hit-Ratio HR at the same $f(k)$ location. Recall Figures 1 shows that the sampling Hit-Ratio curves are comparable with respect to their respective counterparts, and suggests that such a mapping function, indeed, may exist.

But how does this requirement coexist with the first requirement of the minimal error of individual curves? We note that, for most

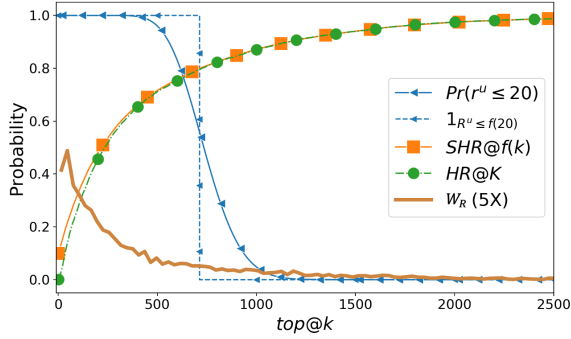


Figure 2: Curve Relationship. $HR@K$ is the top- K global hit-ratio; $SHR@f(k)$ is the sampling top- k hit-ratio shown in global scale; $W_R(5X)$ is the empirical user ranking distribution, where we multiply by 5 for displaying purpose.

of the recommendation algorithms, their overall Hit-Ratio curve $HR@K$, and the empirical probability mass function W_R (Formula 2), are actually fairly similar. From another viewpoint, if we allow individual curves to have different optimal $f(k)$, the difference (or shift) between them is rather small and does not affect the performance comparison between them, using the sampling curves $SHR@k$. We will make this case more rigorously in Section 4, and we refer to this problem as the *sampling correspondence alignment problem*.

In this section, we will focus on studying dataset-independent mapping functions, and we will discuss the algorithm-specific and dataset-specific mapping function in the next section.

3.1 Boundary Condition Approximation

Consider that sampling with replacement, for any individual user, X^u from equation 6, obeys binomial distribution. Apply the general case of bounded variables Hoeffding's inequality:

$$Pr(|X^u - E[X^u]| \geq t) \leq 2e^{-\frac{2t^2}{n-1}}$$

since $r^u = X^u + 1$, and $E[r^u] = E[X^u] + 1 = (n-1)p_u + 1$:

$$\begin{cases} Pr(r^u \geq (n-1)p_u + 1 + t) \leq 2e^{-\frac{2t^2}{n-1}} \\ Pr(r^u \leq (n-1)p_u + 1 - t) \leq 2e^{-\frac{2t^2}{n-1}} \end{cases} \quad (10)$$

The above inequalities indicate that r^u is restricted around its expectation within the range defined by t .

The second term of error in equation 8 can be written as:

$$\begin{aligned} E[SHR@k] - HR@f(k) \\ = - \sum_{R=1}^{f(k)} W_R \cdot Pr(r^R \geq k+1) + \sum_{R=f(k)+1}^N W_R \cdot Pr(r^R \leq k) \end{aligned} \quad (11)$$

where, $r^R = r^u$ for $R_u = R$. For some relatively large t (compared to $\sqrt{n-1}$), the probability in equation 10 can come extremely close to 0. Based on this fact, if we would like to limit the first term $Pr(r^R \leq k+1)$ to approach 0, $k+1$ must be greater than $(n-1)p_u + 1 + t$. And similar to the second term, we have:

$$\begin{cases} r^u \geq k+1 \geq (n-1)\frac{R-1}{N-1} + 1 + t, & R = 1, \dots, f^l(k) \\ r^u \leq k \leq (n-1)\frac{R-1}{N-1} + 1 - t, & R = f^u(k) + 1, \dots, N \end{cases}$$

where $f^l(k)$ and $f^u(k)$ are the lower bound and upper bound for $f(k)$, respectively. Explicitly,

$$f^l(k) \leq (k-t) \cdot \frac{N-1}{n-1} + 1, \quad f^u(k) \geq (k+t-1) \cdot \frac{N-1}{n-1} \quad (12)$$

Given this, let the average of these two for f :

$$f(k) = \lfloor \frac{f^l(k) + f^u(k)}{2} \rfloor = \lfloor (k - \frac{1}{2}) \frac{N-1}{n-1} + \frac{1}{2} \rfloor \quad (13)$$

Note that, although this formula appears similar to our baseline Formula 9, the difference between them is actually pretty big ($\approx \frac{1}{2} \frac{N-1}{n-1}$). As we will show in the experimental results, this formula is remarkably effective in reducing the error $|HR@f(k) - SHR@k|$.

3.2 Beta Distribution Approximation

In this approach, we try to directly minimize $HR@f(k) - E[SHR@k]$ to 0, and this is equivalent to:

$$\sum_{R=1}^N W_R \cdot \mathbf{1}_{R \leq f(k)} = \sum_{R=1}^N W_R \cdot Pr(r^R \leq k) \quad (14)$$

In order to get a closed-form solution of $f(k)$ from the above equation, we leverage the Beta distribution $Beta(a, 1)$ to represent the user ranking distribution W_R , inspired by [12]: $W_R = \frac{1}{\mathcal{B}(a, 1)} (\frac{R-1}{N-1})^{a-1} \frac{1}{N-1}$, where a is a constant parameter and $\frac{1}{N-1}$ is the constant for discretized Beta distribution. Note that $\frac{R-1}{N-1}$ normalizes the user rank R_u from $[1, N]$, to $[0, 1]$. Especially, when $a < 1$, this distribution can represent exponential distribution, which can help provide fit for the HR distribution. Figure 3 illustrates the Beta distribution fitting of W_R .

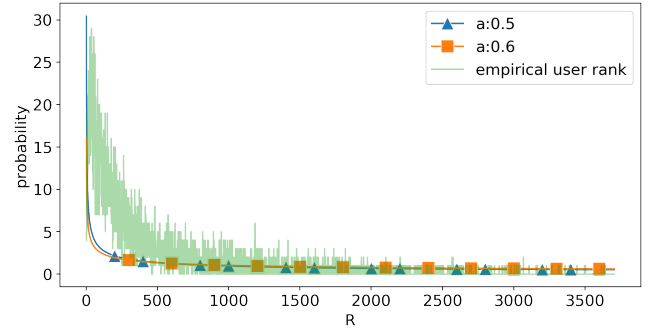


Figure 3: Beta distributions and empirical user rank distribution W_R

The left term of the equation 14 is denoted as \mathcal{L}_k

$$\begin{aligned} &= \sum_{R=1}^N W_R \cdot \mathbf{1}_{R \leq f(k)} = \sum_{R=1}^{f(k)} W_R = \frac{1}{\mathcal{B}(a, 1)} \sum_{R=1}^{f(k)} (\frac{R-1}{N-1})^{a-1} \cdot \frac{1}{N-1} \\ &= \frac{1}{\mathcal{B}(a, 1)} \sum_{x=0}^{\frac{f(k)-1}{N-1}} x^{a-1} \cdot \Delta x \quad \text{where, } x = \frac{R-1}{N-1}, \text{ and } \Delta x = \frac{1}{N-1}, \\ &\approx \frac{1}{\mathcal{B}(a, 1)} \int_0^{\frac{f(k)-1}{N-1}} x^{a-1} dx = \frac{1}{a\mathcal{B}(a, 1)} \left[\frac{f(k)-1}{N-1} \right]^a \end{aligned}$$

Considering sampling with replacement, then the right term is denoted as:

$$\mathcal{R}_k = \sum_{i=0}^{k-1} \binom{n-1}{i} \sum_{R=1}^N W_R \left(\frac{R-1}{N-1} \right)^i \left(1 - \frac{R-1}{N-1} \right)^{n-i-1}$$

Calculate the difference:

$$\begin{aligned} \mathcal{R}_{k+1} - \mathcal{R}_k &= \binom{n-1}{k} \sum_{R=1}^N W_R \left(\frac{R-1}{N-1} \right)^k \left(1 - \frac{R-1}{N-1} \right)^{n-1-k} \\ &\approx \binom{n-1}{k} \frac{1}{\mathcal{B}(a, 1)} \int_{x=0}^1 x^{a+k-1} (1-x)^{n-1-k} dx \\ &= \binom{n-1}{k} \frac{1}{\mathcal{B}(a, 1)} \mathcal{B}(a+k, n-k) = \frac{1}{\mathcal{B}(a, 1)} \frac{\Gamma(n)}{\Gamma(n+a)} \frac{\Gamma(k+a)}{\Gamma(k+1)} \end{aligned}$$

Based on above equations: $\mathcal{L}_{k+1} - \mathcal{L}_k \approx \mathcal{R}_{k+1} - \mathcal{R}_k$, we have (we denote the mapping function as $f(k; a)$ for parameter a).

$$\begin{aligned} [f(k+1; a) - 1]^a - [f(k; a) - 1]^a \\ = a[N-1]^a \binom{n-1}{k} \mathcal{B}(a+k, n-k) \end{aligned} \quad (15)$$

Then we have the following recurrent formula:

$$\boxed{f(k; a) = \left[a[N-1]^a \binom{n-1}{k} \mathcal{B}(a+k, n-k) + [f(k-1; a) - 1]^a \right]^{1/a} + 1} \quad (16)$$

where we have $f(1)$ by considering $\mathcal{L}_1 = \mathcal{R}_1$:

$$f(1; a) = (N-1) [a \mathcal{B}(a, n)]^{1/a} + 1 \quad (17)$$

3.3 Properties of Recurrent function f

In the following, we enumerate a list of interesting properties of this recurrent formula of f based on Beta distribution.

LEMMA 3.1 (LOCATION OF LAST POINT). *For any a , all $f(n)$ converge to N : $f(n) = N$.*

PROOF. We note that

$$\begin{aligned} \sum_{k=0}^{n-1} \binom{n-1}{k} \mathcal{B}(a+k, n-k) &= \int_0^1 \sum_{k=0}^{n-1} \binom{n-1}{k} t^{a+k-1} (1-t)^{n-k-1} dt \\ &= \int_0^1 t^{a-1} \left[\sum_{k=0}^{n-1} \binom{n-1}{k} t^k (1-t)^{n-k-1} \right] dt = \int_0^1 t^{a-1} dt = \frac{1}{a} \end{aligned}$$

Add up equation 15 from $k=1$ to $n-1$, we have:

$$\begin{aligned} \frac{[f(n) - 1]^a - [f(1) - 1]^a}{a[N-1]^a} &= \sum_{k=1}^{n-1} \binom{n-1}{k} \mathcal{B}(a+k, n-k) \\ &= \frac{1}{a} - \binom{n-1}{0} \mathcal{B}(a, n) \quad \text{then, we have } f(n) = N \end{aligned}$$

□

Uniform Distribution and Linear Map: When the parameter $a=1$, the Beta distribution degenerates to the uniform distribution. From equations 17 and 16, we have another simple linear map:

$$f(k; a=1) = k \frac{N-1}{n} + 1 \quad (18)$$

Even though the user rank distribution is quite different from the uniform distribution, we found that this formula provides a reasonable approximation for the mapping function, and generally, better than the Naive formula 9. More interestingly, we found that when a ranges from 0 to 1 (as they express an exponential-like distribution), they actually are quite close to this linear formula.

Approximately Linear: When we take a close look at the $f(k; a)$ sequences $f(1; a), f(2; a), \dots, f(k; a)$ for different parameters a from 0 to 1, we find that when k is large, $f(k; a)$ all gets very close to $f(k; 1)$ (the linear map function for the uniform distribution). Figures 4 show the relative difference of all $f(k; a)$ sequences for $a = 0.2, 0.6, 0.8$ with respect to $a = 1$, i.e., $[f(k; a) - f(k; a=1)]/f(k; a=1)$. Basically, they all converge quickly to $f(k; a=1)$ as k increases.

To observe this, let us take a look at their $f(k)$ locations when k is getting large. To simplify our discussion, let $g(k) = f(k) - 1$, and then we have

$$[g(k+1)]^a - [g(k)]^a = a(N-1)^a \frac{\Gamma(n)}{\Gamma(n+a)} \frac{\Gamma(k+a)}{\Gamma(k+1)}$$

$$\left[\frac{g(k+1)}{g(k)} \right]^a = 1 + \left[\frac{(N-1)k}{g(k)n} \right]^a \frac{a}{k}$$

$$\text{when } n \text{ and } k \text{ are large, } \lim_{n \rightarrow \infty} \frac{\Gamma(n)}{\Gamma(n+a)} = \frac{1}{n^a}$$

$$\left[\frac{g(k+1)}{g(k)} \right] = \left(1 + \left[\frac{(N-1)k}{g(k)n} \right]^a \frac{a}{k} \right)^{1/a} \approx 1 + \left[\frac{(N-1)k}{g(k)n} \right]^a \frac{1}{k}$$

When $g(k) = (N-1) \frac{k}{n}$, the above equation holds $\frac{g(k+1)}{g(k)} = 1 + \frac{1}{k}$, and this suggests they are all quite similar to the linear map $f(k; a=1)$ for the uniform distribution.

By looking at the difference $f(k; a) - f(k-1; a)$, we notice we will get very close to the constant $\frac{N-1}{n} = f(k; 1) - f(k-1; 1)$ even when k is small. To verify this, let y_{k+1}

$$= [f(k+1; a) - 1]^a - [f(k; a) - 1]^a = a[N-1]^a \frac{\Gamma(n)}{\Gamma(n+a)} \frac{\Gamma(k+a)}{\Gamma(k+1)}$$

Then we immediately observe:

$$\frac{y_{k+1}}{y_k} = \frac{a[N-1]^a \frac{\Gamma(n)}{\Gamma(n+a)} \frac{\Gamma(k+a)}{\Gamma(k+1)}}{a[N-1]^a \frac{\Gamma(n)}{\Gamma(n+a)} \frac{\Gamma(k-1+a)}{\Gamma(k)}} = 1 + \frac{a-1}{k}$$

Thus, after only a few iterations for $f(k; a)$, we have found that their (powered) difference will get close to being a constant.

4 DATASET-SPECIFIC MAPPING FUNCTION

In the last section, we introduced some generic mapping functions, which can be used across different datasets. As we will show in Section 5, choosing some generic a , e.g., $a=0.5$, can provide quite accurate approximation; the differences between $HR@f(k)$ and $SHR@k$ are quite small. However, why can such a generic a be so effective? Furthermore, can we design a better mapping function which can leverage the inherent characteristics of recommendation algorithm performance on different datasets?

By answer these questions, we found some further interesting properties of the sampling Hit-Ratio metrics: They can be a rather robust (or safe) measure to select the best algorithm! In layman's terms, if the sampling metric SHR shows a good improvement of an

¹This also holds when $a > 1$, but since the Hit-Ratio, aka the user rank distribution, is typically very different from these settings, we do not discuss them here.

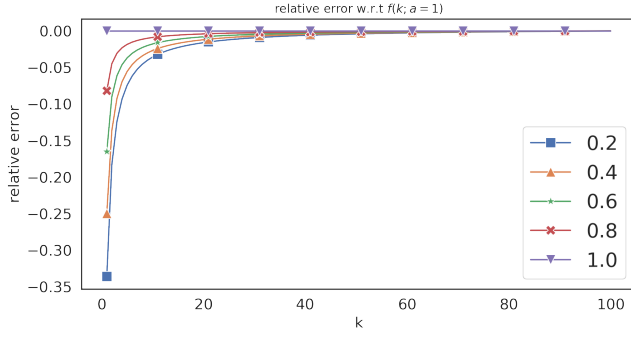


Figure 4: Relative error w.r.t. $f(k; a = 1)$.ml-1m dataset, $n = 100$.

algorithm over others, then there is a good chance that it may perform even better in the global top- K Hit-Ratio metrics. Conversely, if an algorithm under-performed in the sampling measure, it may be even worse in the global measure.

4.1 Optimizing a (Algorithm-Specific Mapping)

Recall that we use $Beta(a, 1)$ to fit W_R . Since W_R , the accumulative distribution of which is HR , is clearly unknown to us, we will try to use SHR to represent it instead. But how does it relate to HR ? Based on our earlier discussion, it can be approximated as $HR@f(k) \approx SHR@k$. Especially, if we leverage the $Beta$ distribution, when a is given, we can use the aforementioned $f(k; a)$ for our purpose. Once we have such mapping, we can then use $SHR@k$ (more precisely $\widehat{HR}@f(k) = SHR@k$), to help fit the beta distribution and consequently find a new parameter a .

Given this, we can utilize the following iterative procedure to identify the optimal a , which uses a maximal likelihood approach to help fit the beta distribution:

$$a^{(i+1)} = \frac{-M}{\sum_{u=1}^M \ln\left(\frac{f(r^u; a^{(i)}) - 1}{N-1}\right)} \quad (19)$$

where $\mathcal{L} := \prod_{u=1}^M a \left(\frac{f(r^u; a)-1}{N-1}\right)^{a-1} \frac{1}{N-1}$ is the likelihood function based on Beta distribution, and if we take the derivative of the log-likelihood, we have the above formula as the optimal parameter a . We can start with any reasonable a , such as $a = 1$ or $a = 0.5$. This procedure can be considered as a simplified EM algorithm. Our experiments show it converges very quickly (in two or three iterations) to some fixed point.

4.2 Sampling Correspondence Alignment Problem and its Remedy

Now, for different recommendation algorithms on the same data with the same sampling size n , each produces their corresponding sampling Hit-Ratio $SHR@k$, and each will produce different parameters a using the above method (Formula 19), which leads to a different mapping function $f(k; a)$. This leads to the *sampling correspondence alignment problem*: for any fixed k under sampling,

different algorithms' $SHR@k$ (with different a) measures their corresponding HR at different location $f(k; a)$. Given this, can sampling $SHR@k$ still be meaningful for performance evaluation?

Remedy #1: The difference between a is very small: Through extensive experimental evaluation, we found that on the same dataset, the optimal a of different recommendation algorithms are actually very close to one another (See Table 2, sampling size $n = 100$). In most of these datasets, their a 's difference is within 0.01.

Table 2: Beta Parameter

Model	ml-1m	yelp	pinterest-20	citeulike
NeuMF	0.3685	0.3059	0.2820	0.2601
MultivAE	0.3681	0.2977	0.2764	0.2449
EASE	0.3684	0.2972	0.2806	0.2532
itemKNN	0.4079	0.2885	0.2782	0.2519

Remedy #2: The difference between $f(k; a)$ for slightly different a is very small: When we put slightly different a into the mapping function, and observe their difference, $f(k; a) - f(k; a')$ is also very small. Figure 5 shows the difference of mapping functions for a ranges from 0.24 to 0.34 with the mapping function at $a = 0.3$, i.e., $f(k; a) - f(k; a = 0.3)$. We see that their absolute location difference is less than 1; i.e., with the original sampling location on the scale from 1 to N , for the same k , their correspondence location difference is less than 1. Thus, we generally can use any of the a obtained from one of the recommendation algorithms on a dataset as the choosing parameters for all the recommendation algorithms. In fact, this also suggests the parameter a is an inherent parameter for each dataset when using the existing (competitive) recommendation algorithms. Thus, we can have dataset-specific a (without worrying the algorithm-specific a).

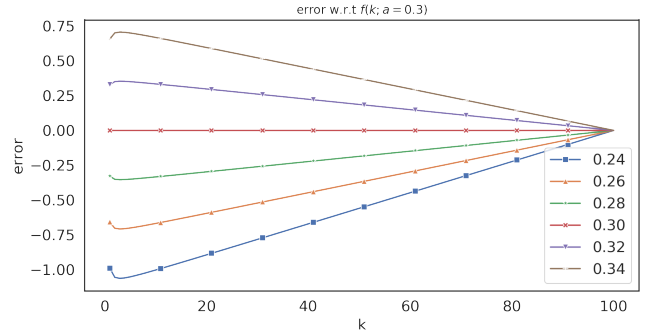


Figure 5: Error of $f(k; a)$ w.r.t $f(k; 0.3)$, a ranges from 0.24 to 0.34. ml-1m dataset, $n = 100$.

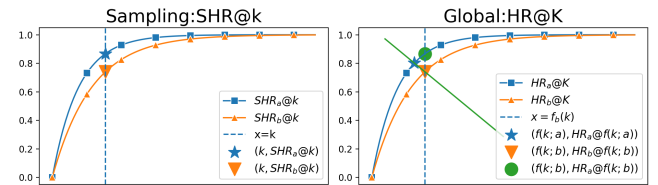


Figure 6: Sampling Effect

Remedy #3: When a is smaller, so is $f(k; a)$: An interesting discovery can be made by observing Figures 5 and 4: When $a > b$, $f(k; a) \leq f(k; b)$ for any k . This holds even if a and b are not close. When a and b are closer, then their difference has become smaller; but as long as a is larger than b , the mapping function $f(k; a)$ always corresponds to an earlier location than $f(k; b)$. Intuitively, when a is smaller, this corresponds to more users with higher rank, i.e., their Hit-Ratio $HR@K$ (accumulative distribution) under beta distribution is consistently better than the larger b . Mathematically, this is equivalent to saying that $f(k; a)$ is a monotone increasing function with respect to a (for $a > 0$). Even though we can numerically observe this, the rigorous proof of its monotonicity remains an open problem.

The implication of such a monotone property is quite interesting and likely useful: If the sampling metrics show a good improvement of an algorithm over another algorithm, then there is a good chance that it may perform even better in the global top- K Hit-Ratio metric, as it may actually correspond to an earlier (or smaller) $f(k)$: assuming $SHR_a@k > SHR_b@k$, and $f(k; a) < f(k; b)$,

$$HR_a@f(k; b) \geq HR_a@f(k; a) \approx SHR_a@k \geq SHR_b@k \approx HR_b@f(k; b)$$

Figure 6 illustrates this effect. This helps explain why the sampling Hit-Ratio is very effective in choosing the winner (or loser) of different recommendation algorithms. They ensure the correct prediction when two recommendation algorithms perform very differently (not covered by Remedy #2).

5 EXPERIMENTAL RESULTS

In this section, we experimentally study the sampling hit ratio $SHR@k$ and its corresponding global hit ratio $HR@K$ through different mapping functions f . Specifically, we aim to answer:

- (Question 1) How does the dataset-independent mapping function f help align $SHR@k$ with respect to $HR@f(k)$?
- (Question 2) How do different factors, such as the top- k location, varying the effect of the sampling scheme, and the sampling size affect the results?
- (Question 3) How does the algorithm-specific mapping function f compare with other mapping functions?
- (Question 4) How can the random sampling hit ratios be used to identify the winners of recommendation algorithms with respect to the corresponding global hit ratio?

In this section, we only focus on Question 1 and Question 4. The discussion of Question 2 and 3, together with their experimental results and the experimental setup will be put into the Appendix A. Due to the space limitation, we only report representative results here, and additional experimental results are openly available at ².

Aligning Sampling and Global Hit Ratio SHR and HR : In this experiment, we provide two dual views of the alignment between sampling and global hit ratio (curves). Here, we report only the MultiVAE result on dataset ml-1m in Figures 7. We use four different dataset-independent mapping functions, the linear, bound, $\beta@1$ and $\beta@0.5$, for the curve alignment. Figure 7a maps the sampling curve $SHR@k$ to the global top- K view by mapping $SHR@k$ to location $f(k)$ in the population/global top K view, and compares them with the global $HR@K$ (population curve). Figure 7b maps

the global curve $HR@K$ to the sampling top- k view by $HR@K$ to the location k (where $K = f(k)$) in the sampling top k view, and then compares them with the sample $SHR@k$ (sampling curve). We observe that both bound and $\beta@0.5$ achieve the best results from both views. The same observation holds on other recommendation algorithms and datasets.

Predicting Winners (and Relative Performance): In Table 3, we demonstrate the effectiveness of using sampling hit ratio $SHR@k$ to predict the recommendation algorithm performance when using the global hit ratio $HR@f(k)$. We compare the performance of the three most competitive recommendation methods on the four commonly used datasets. We also vary the k from 1 to 50 for sampling size $n - 1 = 99$. Throughout all these cases, $SHR@k$ and $HR@f(k)$ all consistently predict the same and correct winners. Specifically, if an algorithm has the highest $SHR@k$, then it has the highest $HR@f(k)$ as well. In fact, their relative orders are also mostly consistent besides the winners.

6 CONCLUSION AND DISCUSSION

In this work, we provide a thorough investigation of the sampling top- k Hit-Ratio and show how it “samples” the global Hit-Ratio in a way similar to “signal sampling”. Theoretically and empirically, we demonstrate the predictive power of sampling metrics, in terms of both approximating corresponding global Hit-Ratio, and predicting the relative performance between different algorithms.

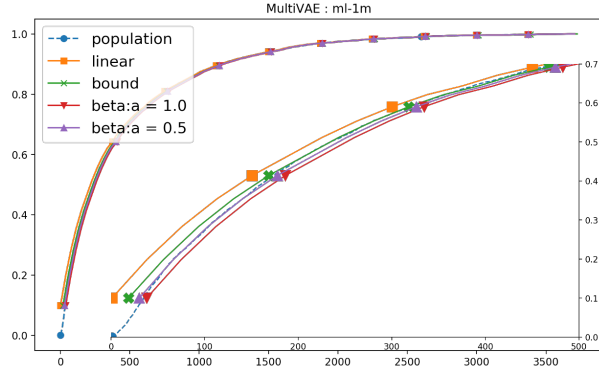
We would like to point out that the mapping function serves as a scaffold for us to understand how sampling works, with respect to the global Hit-Ratio curve. It provides us with a basic tool to help verify the accuracy of the metric being observed from the sampling, with respect to the global curve. Following our theoretical investigation and experimental testing/evaluating of different mapping functions (Section 5), we can safely use the sampling Hit-Ratio metrics without worrying about the mapping function.

However, there are several interesting open questions which need further investigation. (1) Can other sampling-based metrics, such as average precision and nDCG, have such properties as the Hit-Ratio (recall/precision)? (2) Can other distributions fit Hit-Ratio curves better than the beta distribution? (3) Can the monotone property of the mapping function (with respect to the parameter a in beta distribution) hold beyond Beta distribution? We plan to (and also welcome the community to) investigate these questions.

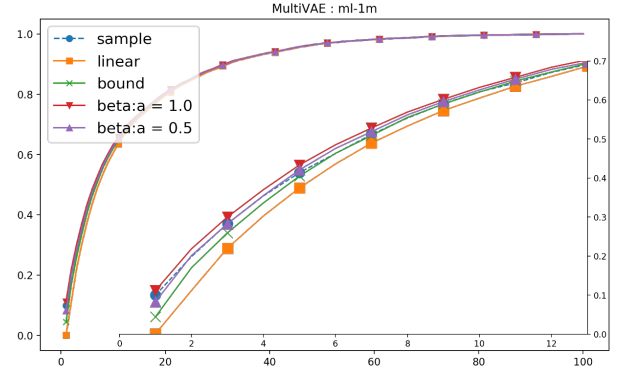
REFERENCES

- [1] Immanuel Bayer, X. He, B. Kanagal, and S. Rendle. 2017. A Generic Coordinate Descent Framework for Learning from Implicit Feedback. In *WWW'17*.
- [2] Paolo Cremonesi, Yehuda Koren, and Roberto Turrin. 2010. Performance of Recommender Algorithms on Top-n Recommendation Tasks. In *RecSys'10*.
- [3] Maurizio Ferrari Dacrema, P. Cremonesi, and D. Jannach. 2019. Are We Really Making Much Progress? A Worrying Analysis of Recent Neural Recommendation Approaches. In *RecSys'19*.
- [4] Mukund Deshpande and George Karypis. 2004. Item-Based Top-N Recommendation Algorithms. *ACM Trans. Inf. Syst.* (2004).
- [5] Travis Ebesu, Bin Shen, and Yi Fang. 2018. Collaborative Memory Network for Recommendation Systems. In *SIGIR'18*.
- [6] Ali Mamdouh Elkahky, Y. Song, and X. He. 2015. A Multi-View Deep Learning Approach for Cross Domain User Modeling in Recommendation Systems. In *WWW'15*.
- [7] Xiangnan He, L. Liao, H. Zhang, L. Nie, X. Hu, and T. Chua. 2017. Neural Collaborative Filtering. In *WWW'17*.
- [8] Binbin Hu, C. Shi, W. X. Zhao, and P. S. Yu. 2018. Leveraging Meta-Path Based Context for Top- N Recommendation with A Neural Co-Attention Model. In

²<https://github.com/dli12/KDD20-On-Sampling-Top-K-Recommendation-Evaluation>



(a) Hit Ratios from the Global Top- K . Except the "population" curve, all others map the sample hit-ratio to global scale by different mapping functions.



(b) Hit Ratios from the Sampling Top- k . Except the "sample" curve, all others map the global/population hit-ratio to sample scale by different mapping functions.

Figure 7: The right part of each figure is the zoom in version which display the relation of curves more clear. dataset setting: ml-1m, $n = 100$

Table 3: Predicting Winners (and Relative Performance)

ml-1m								
k	NeuMF			MultiVAE			EASE	
	SHR	bound	$\mathcal{B}@.5$	SHR	bound	$\mathcal{B}@.5$	SHR	bound
1	0.208	0.150	0.205	0.211	0.160	0.216	0.223	0.173
2	0.326	0.311	0.347	0.343	0.319	0.355	0.349	0.334
5	0.548	0.555	0.566	0.555	0.560	0.576	0.564	0.573
10	0.715	0.726	0.731	0.717	0.729	0.733	0.720	0.733
20	0.850	0.863	0.864	0.854	0.866	0.867	0.847	0.859
50	0.972	0.979	0.979	0.972	0.980	0.979	0.955	0.961

pinterest-20								
k	NeuMF			MultiVAE			EASE	
	SHR	bound	$\mathcal{B}@.5$	SHR	bound	$\mathcal{B}@.5$	SHR	bound
1	0.273	0.201	0.278	0.316	0.241	0.321	0.289	0.222
2	0.436	0.409	0.448	0.479	0.456	0.495	0.452	0.425
5	0.701	0.708	0.722	0.729	0.735	0.747	0.705	0.712
10	0.874	0.883	0.886	0.887	0.894	0.897	0.868	0.877
20	0.965	0.968	0.968	0.970	0.973	0.973	0.959	0.963
50	0.993	0.993	0.993	0.994	0.994	0.994	0.989	0.989

yelp								
k	NeuMF			MultiVAE			EASE	
	SHR	bound	$\mathcal{B}@.5$	SHR	bound	$\mathcal{B}@.5$	SHR	bound
1	0.234	0.184	0.248	0.262	0.215	0.283	0.275	0.228
2	0.369	0.360	0.392	0.404	0.398	0.429	0.418	0.410
5	0.593	0.600	0.611	0.626	0.635	0.645	0.632	0.642
10	0.753	0.760	0.764	0.775	0.781	0.783	0.777	0.784
20	0.881	0.885	0.886	0.892	0.894	0.895	0.886	0.888
50	0.975	0.976	0.976	0.974	0.975	0.975	0.957	0.957

criteo								
k	NeuMF			MultiVAE			EASE	
	SHR	bound	$\mathcal{B}@.5$	SHR	bound	$\mathcal{B}@.5$	SHR	bound
1	0.399	0.427	0.505	0.488	0.554	0.615	0.464	0.511
2	0.583	0.606	0.632	0.675	0.707	0.727	0.641	0.674
5	0.767	0.774	0.783	0.828	0.839	0.845	0.802	0.811
10	0.871	0.878	0.880	0.910	0.914	0.916	0.880	0.886
20	0.940	0.944	0.944	0.961	0.962	0.963	0.934	0.936
50	0.987	0.987	0.987	0.992	0.991	0.991	0.969	0.968

KDD'18.

- [9] Y. Hu, Y. Koren, and C. Volinsky. 2008. Collaborative filtering for implicit feedback datasets. In *ICDM'08*.
- [10] Yehuda Koren. 2008. Factorization Meets the Neighborhood: A Multifaceted Collaborative Filtering Model. In *KDD'08*.
- [11] Walid Krichene, N. Mayoraz, S. Rendle, L. Zhang, X. Yi, L. Hong, Ed H. Chi, and J. R. Anderson. 2019. Efficient Training on Very Large Corpora via Gramian Estimation. In *ICLR'2019*.
- [12] Wentian Li, Pedro Miramontes, and Cocho Germinal. 2010. Fitting Ranked Linguistic Data with Two-Parameter Functions. (2010).
- [13] Dawen Liang, R. G. Krishnan, M. D. Hoffman, and T. Jebara. 2018. Variational Autoencoders for Collaborative Filtering. In *WWW'18*.
- [14] K. Deegha Rao. 2018. *Signals and Systems*. Springer International Publishing.
- [15] Steffen Rendle. 2019. Evaluation Metrics for Item Recommendation under Sampling. arXiv:1912.02263
- [16] Steffen Rendle, Li Zhang, and Yehuda Koren. 2019. On the Difficulty of Evaluating Baselines: A Study on Recommender Systems. (2019).
- [17] Harald Steck. 2019. Embarrassingly Shallow Autoencoders for Sparse Data. *WWW'19* (2019).

- [18] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. 2019. Neural Graph Collaborative Filtering. *SIGIR'19* (2019).
- [19] Xiang Wang, D. Wang, C. Xu, X. He, Y. Cao, and T. Chua. 2019. Explainable Reasoning over Knowledge Graphs for Recommendation. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI2019*.
- [20] Longqi Yang, Eugene Bagdasaryan, Joshua Gruenstein, Cheng-Kang Hsieh, and Deborah Estrin. 2018. OpenRec: A Modular Framework for Extensible and Adaptable Recommendation Algorithms. In *WSDM'18*.
- [21] Longqi Yang, Y. Cui, Y. Xuan, C. Wang, S. J. Belongie, and D. Estrin. 2018. Unbiased Offline Recommender Evaluation for Missing-Not-at-Random Implicit Feedback. In *RecSys'18*.
- [22] Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. 2019. Deep Learning Based Recommender System: A Survey and New Perspectives. *ACM Comput. Surv.* (2019).

7 ACKNOWLEDGMENTS

The research was partially supported by a sponsorship research agreement between Kent State University and iLambda, Inc.

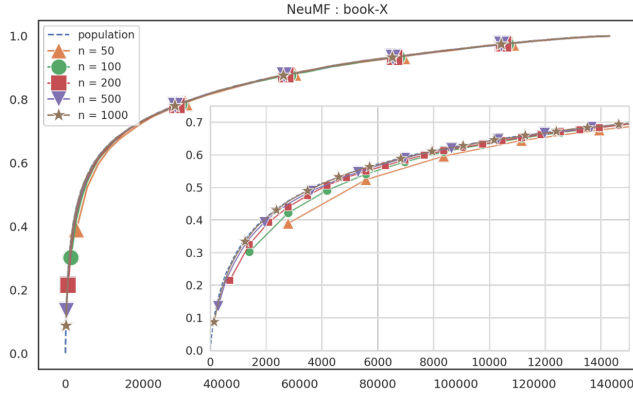


Figure 8: Global/Population HR curves compared with sample ones with different sample size. Experiment set-up: dataset:book-X, model:NeuMF

Table 4: Dataset Statistics

Dataset	Interactions	Users	Items	Sparsity
ml-1m	1,000,209	6,040	3,706	95.53%
pinterest-20	1,463,581	55,187	9,916	99.73%
citeulike	204,986	5,551	16,980	99.78%
yelp	696,865	25,677	25,815	99.89%
book-X	786,690	11,325	139,331	99.95%

A EXPERIMENTS AND REPRODUCIBILITY

Experimental Setup: We use four of the most commonly used datasets for recommendation studies, together with the book-X dataset, with a relatively large number of total items, for evaluating the effect of sampling size. The characteristics of these datasets are in Table 4.

For the different recommendation algorithms, we use some of the most well-known and the state-of-the-art algorithms, including three non-deep-learning options: itemKNN [4]; ALS [9]; and EASE [17]; and three deep learning options: NCF [7]; MultiVAE [13]; and NGCF [18]. For each recommendation algorithm on a particular dataset, we report the result of only one sample run, as we found they are very close to the average of multiple run results. The default sampling method is sampling with replacement, unless explicitly stated. Also, we collect both sampling hit ratio $SHR@k$ (k from 1 to n , the sampling size) and global hit ratio $HR@K$ (K from

1 to N , the number of total items). In the figures, we also use *population* for $HR@K$ (global) and *sample* for $SHR@k$ curves. For the different mapping functions, we consider the linear (Formula 9), the bound (Formula 13), $Beta@1$ ($f(k; a = 1)$, Formula 18), $Beta@0.5$ ($f(k; a = 0.5)$, Formula 16), and $Beta@P$ for the algorithm-specific mapping (Formula 19). Below, we will report our experimental findings for the aforementioned questions 2 and 3 in Section 5.

Key Factor Analysis (Top k , Sampling Factors and Sampling Size): To take a close look at different mapping functions, we listed the sampling hit ratio SHR at $k = 1, 2, 5, 10, 20, 50$ for sampling size $n - 1 = 99$ and their corresponding global hit ratio HR at $f(k)$ locations based on mapping functions, bound, $beta@1$ and $beta@0.5$ in Table 5. We show similar results for sampling size $n - 1 = 999$. In addition, we also consider three different sampling schemes: sampling with replacement (binom); sampling without replacement (hyper); and sampling without replacement using only irrelevant items (actual). We made the following observations. (1) The differences between the sampling and the global $|SHR@k - HR@f(k)|$ are fairly small, with the bound and $Beta@0.5$ more accurate; (2) When k becomes larger, the results are more accurate. (3) The results on the sampling with replacement and without replacement are very close to each other. Sampling with only irrelevant items, and the Global, which ranks only the irrelevant items, both lead to higher hit ratios. But the mapping function works equally well for this situation. (4) When sampling size increases (from $n = 100$ to $n = 1000$), the error also reduces. We further confirm this using Figure 8, which varies n from 50 to 1000 on book-X dataset with more 139K items. When n increases, the sampling hit ratio curve converges to the global hit ratio (population) rather quickly.

Dataset-independent vs algorithm-specific mapping: Table 7 compares different dataset-independent mapping functions with the algorithm-specific one, $Beta@P$. We compare both overall average absolute error ($\sum_{k=1}^n |HR@f(k) - SHR@k|/n$), the overall average relative error ($\sum_{k=1}^n (|HR@f(k) - SHR@k|/SHR@k)/n$, the errors at the top 1, i.e., $|HR@f(1) - SHR@1|$ and $|HR@f(1) - SHR@1|/SHR@1$, and the errors between the top 2 and 10. We observe that the algorithm-specific mapping function $Beta@P$ achieves the most minimal errors, over all. However, the dataset-independent measures, such as $Beta@0.5$, obtain comparable results and perform better when k is small. We believe one of the underlying reasons is that $Beta@P$ aims to fit all the users (including those on the lower rank), and the fitting of Beta distribution has limitations. Thus, we consider that the dataset-independent mapping functions, such as $Beta@0.5$, can be a relatively cost-effective way to align sampling and global top- k hit ratio curves.

Table 5: Varying Sampling Top k on MultiVAE:ml-1m:n=100

binom(with replacement)					hyper(without replacement)					actual(training data rank bottom)				
k	SHR	bound	beta@1	beta@0.5	k	SHR	bound	beta@1	beta@0.5	k	SHR	bound	beta@1	beta@0.5
1	0.1031	0.0442	0.1116	0.0829	1	0.1012	0.0442	0.1116	0.0829	1	0.2101	0.1608	0.2540	0.2161
2	0.2041	0.1700	0.2185	0.2012	2	0.1964	0.1700	0.2185	0.2012	2	0.3349	0.3199	0.3810	0.3555
5	0.4157	0.4050	0.4334	0.4209	5	0.4119	0.4050	0.4334	0.4209	5	0.5533	0.5603	0.5889	0.5768
10	0.6200	0.6204	0.6323	0.6258	10	0.6157	0.6204	0.6323	0.6258	10	0.7164	0.7295	0.7391	0.7334
20	0.7992	0.8012	0.8050	0.8028	20	0.7975	0.8012	0.8050	0.8028	20	0.8531	0.8661	0.8690	0.8674
50	0.9651	0.9682	0.9684	0.9682	50	0.9661	0.9682	0.9684	0.9682	50	0.9757	0.9800	0.9803	0.9798

Table 6: Varying Sampling Top k on MultiVAE:ml-1m:n=1000

binom(with replacement)					hyper(without replacement)					actual(training data rank bottom)				
k	SHR	bound	beta@1	beta@0.5	k	SHR	bound	beta@1	beta@0.5	k	SHR	bound	beta@1	beta@0.5
10	0.1091	0.0998	0.1116	0.1089	10	0.1096	0.0998	0.1116	0.1089	10	0.2417	0.2407	0.2540	0.2495
20	0.2126	0.2124	0.2185	0.2174	20	0.2126	0.2124	0.2185	0.2174	20	0.3682	0.3704	0.3810	0.3770
50	0.4281	0.4306	0.4334	0.4318	50	0.4290	0.4306	0.4334	0.4318	50	0.5776	0.5861	0.5889	0.5874
100	0.6293	0.6303	0.6316	0.6316	100	0.6315	0.6303	0.6316	0.6316	100	0.7286	0.7381	0.7389	0.7389
200	0.8036	0.8043	0.8050	0.8046	200	0.8060	0.8043	0.8050	0.8046	200	0.8570	0.8684	0.8690	0.8689
500	0.9672	0.9682	0.9684	0.9684	500	0.9677	0.9682	0.9684	0.9684	500	0.9747	0.9800	0.9803	0.9803

Table 7: Mapping Functions Comparison

MultiVAE:ml-1m							NeuMF:yelp						
Dataset	abs	rel	abs@1	rel@1	abs@2-10	rel@2-10	Dataset	abs	rel	abs@1	rel@1	abs@2-10	rel@2-10
Linear	0.0058	0.0226	0.0992	0.9934	0.0391	0.1238	Linear	0.0057	0.0164	0.2329	0.9886	0.0300	0.0632
Bound	0.0024	0.0103	0.0596	0.5970	0.0100	0.0363	Bound	0.0019	0.0046	0.0606	0.2571	0.0061	0.0125
Beta@1	0.0034	0.0074	0.0118	0.1177	0.0165	0.0440	Beta@1	0.0038	0.0070	0.0414	0.1756	0.0221	0.0413
Beta@0.5	0.0019	0.0043	0.0169	0.1692	0.0052	0.0112	Beta@0.5	0.0020	0.0030	0.0032	0.0134	0.0112	0.0197
Beta@0.2	0.0018	0.0065	0.0356	0.3566	0.0060	0.0207	Beta@0.2	0.0015	0.0027	0.0253	0.1075	0.0057	0.0097
Beta@P	0.0018	0.0052	0.0250	0.2504	0.0052	0.0146	Beta@P	0.0015	0.0024	0.0145	0.0617	0.0067	0.0106