
A Fast Sampling Algorithm for Maximum Inner Product Search

Qin Ding

University of California, Davis

Hsiang-Fu Yu

Amazon

Cho-Jui Hsieh

University of California, Los Angeles

Abstract

Maximum Inner Product Search (MIPS) has been recognized as an important operation for the inference phase of many machine learning algorithms, including matrix factorization, multi-class/multi-label prediction and neural networks. In this paper, we propose Sampling-MIPS, which is the first sampling based algorithm that can be applied to the MIPS problem on a set of general vectors with both positive and negative values. Our Sampling-MIPS algorithm is efficient in terms of both time and sample complexity. In particular, by designing a two-step sampling with alias table, Sampling-MIPS only requires constant time to draw a candidate. In addition, we show that the probability of candidate generation in our algorithm is consistent with the true ranking induced by the value of the corresponding inner products, and derive the sample complexity of Sampling-MIPS to obtain the true candidate. Furthermore, the algorithm can be easily extended to large problems with sparse candidate vectors. Experimental results on real and synthetic datasets show that Sampling-MIPS is consistently better than other previous approaches such as LSH-MIPS, PCA-MIPS and Diamond sampling approach.

1 Introduction

Maximum Inner Product Search (MIPS) has been recognized as the main computational bottleneck of the inference phase of many machine learning algorithms, including matrix factorization for recommender systems (Koren et al., 2009), multi-class/multi-label classification (Babbar and Schölkopf, 2017; Yu et al., 2014b), and recurrent neural networks when the output vocabulary size is big (Sutskever

et al., 2014). Given a set of d -dimensional candidate vectors $\mathcal{H} = \{\mathbf{h}_1, \dots, \mathbf{h}_n\}$ and a query vector $\mathbf{w} \in \mathbb{R}^d$, the problem of MIPS is to find the vector in \mathcal{H} with the maximum inner product value with \mathbf{w} :

$$\operatorname{argmax}_{j=1, \dots, n} \mathbf{w}^T \mathbf{h}_j.$$

For instance, in matrix completion model for recommender systems, \mathcal{H} is the latent features for items, \mathbf{w} is the latent feature for a particular user, and the inner product $\mathbf{w}^T \mathbf{h}_j$ reflects the user’s interest in the j -th item. Therefore recommending top items is equivalent to solving the MIPS problem.

A naive linear search procedure that computes the inner product between \mathbf{w} and all $\{\mathbf{h}_j\}_{j=1}^n$ would require $O(nd)$ time, which is too slow for problems with large n . For example, in recommender systems, n (number of items) could easily go beyond a million and the number of features can also be very large. To solve this problem, many approximate algorithms have been proposed in the past to speed up MIPS, including LSH-based approaches and tree-based approaches (Bachrach et al., 2014; Neyshabur and Srebro, 2014; Ram and Gray, 2012; Shrivastava and Li, 2014a,b).

We study the sampling-based algorithm for speeding up the computation of MIPS. The main idea is to estimate the values of $\mathbf{w}^T \mathbf{h}_j$ using a fast sampling approach, and then rank the estimated scores to get the final candidates. Although sampling approach sounds reasonable for this task, it has not been successfully applied to MIPS due to the following reasons: 1) the hardness of handling negative values, 2) each sample has to be drawn efficiently without explicitly computing $\mathbf{w}^T \mathbf{h}_j$, and 3) lack of theoretical guarantee. By Cohen and Lewis (1999), a sampling-based approach was first discussed, but their algorithm only works for non-negative data matrices and so could not be used in most of the real applications. Recently Ballard et al. (2015) proposed a similar sampling approach for Maximum All-pairs Dot-Product (MAD) search problem, but also due to the hardness of handling negative values, they are only able to rank inner products based on absolute value of inner product $|\mathbf{w}^T \mathbf{h}_j|$. Moreover, their analysis works only for positive matrices.

In this paper, we propose a novel Sampling-MIPS algo-

rithm which can be used in not only non-negative cases, but also general cases with both positive and negative values. Our algorithm only requires $O(1)$ time for generating each sample, and we prove that the algorithm will output the correct ranking when sample size is large enough. Furthermore, our algorithm can be used for sparse data with still $O(1)$ time complexity per sample and space complexity proportional to number of nonzero elements in \mathcal{H} . Experimental results show that our proposed algorithm is superior compared to other state-of-the-art methods.

Notations: Throughout the paper, $\mathbf{w} \in \mathbb{R}^d$ is the query vector, $\mathcal{H} = \{\mathbf{h}_1, \dots, \mathbf{h}_n\}$ are the n vectors in the database. For simplicity, we use a n -by- d matrix $H = [\mathbf{h}_1 \dots \mathbf{h}_n]^T$ to denote the database, and $h_{jt} = (\mathbf{h}_j)_t$ is the t -th feature of j -th vector.

2 Related Work

2.1 Previous MIPS Algorithms

It has been shown by Bachrach et al. (2014) that the MIPS problem can be reduced to Nearest Neighbor Search (NNS) by adding one more dimension to the features. Based on this reduction, NNS algorithms such as Locality Sensitive Hashing (LSH) (Indyk and Motwani, 1998) and PCA-tree (a modification of KD-tree) (Bachrach et al., 2014; Ram and Gray, 2012) can be used for solving MIPS. Neyshabur and Srebro (2014); Shrivastava and Li (2014a,b) proposed similar ways to reduce MIPS to NNS and solve the resulting problem using LSH. In the experiments, we show our proposed algorithm is better than both LSH and tree based approaches. Furthermore, as pointed out by Yu et al. (2017), these algorithms have to pre-specify number of hashings and tree splits, thus cannot control the trade-off between speedup and precision in run time. In comparison, our algorithm can easily control this trade-off by changing number of samples used in run time.

Concurrent to this work there are some other recent methods proposed for MIPS. Zhang et al. (2018) also applied the same MIPS-to-NNS reduction and then solve the resulting NNS problem by a graph-based approach (Malkov and Yashunin, 2018). Chen et al. (2019) proposed another gumbel clustering approach for MIPS in NLP applications.

2.2 Greedy-MIPS

Most recently, Yu et al. (2017) developed a deterministic algorithm called Greedy-MIPS. Their algorithm is mainly based on the following assumption:

$$\mathbf{w}^T \mathbf{h}_j > \mathbf{w}^T \mathbf{h}_i \Leftrightarrow \max_{t=1, \dots, d} w_t h_{jt} > \max_{t=1, \dots, d} w_t h_{it}, \quad (1)$$

which means when $\mathbf{w}^T \mathbf{h}_j = \sum_{t=1}^d w_t h_{jt}$ is large, at least one element in the summation will also be large. Based

on this assumption they proposed a fast way to construct candidate set based on the ranking of $\max_{t=1, \dots, d} w_t h_{jt}$ for each j , and then use a linear search to get the final ranking among candidate set.

Clearly, Assumption (1) does not hold for general cases and Greedy-MIPS does not have performance guarantee in general. Nevertheless, Yu et al. (2017) showed that this Greedy-MIPS algorithm is quite powerful on real datasets. In Section 5.4, we conduct a comparison with Greedy-MIPS under different scenarios, showing that although Greedy-MIPS performs well on real datasets, when assumption (1) is violated it will produce poor results. This suggests that using Greedy-MIPS is unsafe. In contrary, our algorithm is guaranteed to output the correct top- K elements when the sample size is large enough.

2.3 Sampling-based approaches for Maximum Squared Inner Product Search

Our algorithm is the first sampling-based algorithm designed for MIPS with general input matrices. All the previous approaches are either restricted to non-negative input matrices or designed for another problem called Maximum Squared Inner Product Search (MSIPS).

Idea of sampling method for inner product search was firstly proposed by (Cohen and Lewis, 1999). The intuition behind this idea is to sample index j with probability $P(j|\mathbf{w}) \propto \mathbf{w}^T \mathbf{h}_j$. Unfortunately, this approach requires all the elements of \mathbf{w} and H to be non-negative in order to define such a probability distribution, thus the applicability of (Cohen and Lewis, 1999) is very limited. Recently, a similar approach called diamond sampling was extended by Ballard et al. (2015). However, also due to the need of defining a valid probability distribution, this method can only identify indices with the largest $|\mathbf{w}^T \mathbf{h}_j|$ or $(\mathbf{w}^T \mathbf{h}_j)^2$, which is called MSIPS problem. If both \mathbf{w} and H have non-negative entries or more generally, $\mathbf{w}^T \mathbf{h}_j \geq 0, \forall j$, MSIPS can be reduced to MIPS problem. But in reality, we frequently have matrices and queries with both positive and negative values. Thus to develop a proper sampling approach for MIPS with negative values with well-founded theory remains a challenge.

3 Proposed Algorithm

Similar to many previous methods, our algorithm has two phases. Before observing any query, we pre-process the database H and construct some data structure to facilitate the later sampling procedure. This is called the pre-processing phase and the time will not be counted into the per-query cost. In the query-dependent phase, the MIPS algorithm is conducted using both query \mathbf{w} and the pre-constructed data structure, and our main goal is to reduce the per-query cost in this phase.

Without loss of generality, in the paper, we assume that none of $\{\mathbf{h}_j\}$ and \mathbf{w} is the zero vector. In the following we first present Sampling-MIPS algorithm for non-negative values and then show how to handle problems with both positive and negative values.

3.1 Sampling-MIPS for non-negative values

First we assume the database H and query \mathbf{w} only have non-negative values. The main idea is to draw samples from the marginal distribution, where each index j is sampled with probability

$$P(j|\mathbf{w}) \propto \mathbf{w}^T \mathbf{h}_j.$$

If we can draw samples according to the above distribution, then with high probability the candidate with the maximum inner product value will be drawn most frequently. Moreover, we have total freedom to control the search quality in query time, since the quality can keep being improved by drawing more samples.

However, computing the distribution of $P(j|\mathbf{w})$ for all $j = 1, \dots, n$ will be extremely time-consuming. In order to sample from this marginal distribution without explicitly forming them, we define a joint distribution for the pair (j, t) over the region $\{1, 2, \dots, n\} \times \{1, 2, \dots, d\}$. Let $P(j, t|\mathbf{w}) \propto w_t h_{jt}$, then the marginal distribution of j can be expressed as

$$P(j|\mathbf{w}) = \sum_{t=1}^d P(j, t|\mathbf{w}) = \sum_{t=1}^d P(t|\mathbf{w})P(j|t, \mathbf{w}), \quad (2)$$

where

$$P(t|\mathbf{w}) = \sum_{j=1}^n P(j, t|\mathbf{w}) \propto \sum_{j=1}^n w_t h_{jt} = w_t s_t, \quad (3)$$

$$\begin{aligned} P(j|t, \mathbf{w}) &= \frac{P(j, t|\mathbf{w})}{P(t|\mathbf{w})} \\ &\propto \frac{w_t h_{jt}}{\sum_{j=1}^n w_t h_{jt}} = \frac{h_{jt}}{\sum_{j=1}^n h_{jt}} = \frac{h_{jt}}{s_t}. \end{aligned} \quad (4)$$

Note that we define $s_t = \sum_{j=1}^n h_{jt}$ to simplify the notation and we assume $s_t \neq 0, \forall t$. From equation (4), we can see that $P(j|t, \mathbf{w})$ is irrelevant to \mathbf{w} , so we can rewrite $P(j|t, \mathbf{w})$ as $P(j|t)$.

Based on (3) and (4), to draw a sample from $P(j|\mathbf{w})$, we can first draw t from a multinomial distribution, $P(t|\mathbf{w}) = \text{multinomial}([\frac{w_1 s_1}{\sum_{t=1}^d w_t s_t}, \dots, \frac{w_d s_d}{\sum_{t=1}^d w_t s_t}])$, and then draw a sample for j from multinomial distribution $P(j|t) = \text{multinomial}([\frac{h_{1t}}{s_t}, \dots, \frac{h_{nt}}{s_t}])$. The procedure is illustrated in Figure 1. In Figure 1, sampling $t \sim P(t|\mathbf{w})$ gives us $t = 4$, so we focus on the 4-th column of matrix H and sample $j \sim P(j|t = 4)$ which gives us $j = 2$, and this sampling procedure is equivalent to sampling j from the marginal distribution $P(j|\mathbf{w})$.

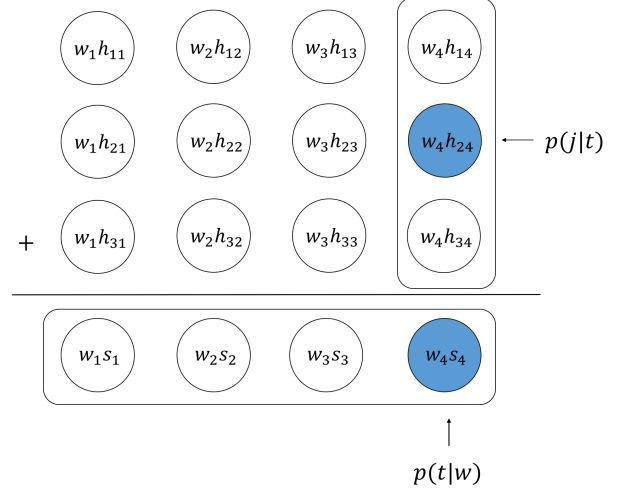


Figure 1: Illustration of the sampling procedure.

In order to quickly draw samples, we use alias method proposed by Walker (1977). For a given k -dimensional multinomial distribution, constructing an alias table requires $O(k)$ time and $O(k)$ space, and the alias table will allow us to draw each sample with only $O(1)$ time. From the above definitions, we can see that the distribution for sampling index j is irrelevant to \mathbf{w} , which means we can construct the alias table for $P(j|t)$ in the pre-processing phase. We present our algorithm for pre-processing phase in Algorithm 1.

Algorithm 1 Query-independent pre-processing

- 1: $s_t \leftarrow \sum_{j=1}^n h_{jt}, \forall t \in \{1, \dots, d\}$
 - 2: Use alias table to construct $P(j|t) \leftarrow \text{multinomial}([h_{1t}, \dots, h_{nt}]), \forall t$
-

In the query-dependent phase, we should first construct the alias table for $P(t|\mathbf{w})$, which requires $O(d)$ time. This will allow us to draw samples from $P(t|\mathbf{w})$ in constant time later. We show the procedure of query-dependent candidate screening in Algorithm 2. We use a vector of length n to count how many times each index is sampled. From line 3 to line 7 in Algorithm 2, we draw B samples and for each one we sample j using the two-step procedure mentioned above. Each time an index j is sampled, we increase the count of index j by 1.

After the sampling process, we sort all the n candidates according to the ranking of x_j 's and take top- C of them. We perform a naive-MIPS on these selected C candidates, i.e., we compute these C inner products and find which candidates have the maximum or top- K values. Details are shown in Algorithm 3.

Since this method is built on the fact that we sample index j with probability proportional to $\mathbf{w}^T \mathbf{h}_j$, we know that the algorithm will succeed (output the real top- K elements)

Algorithm 2 Query-dependent candidate screening

- 1: Use alias table to construct $\dots O(d)$
 $P(t|\mathbf{w}) \leftarrow \text{multinomial}([w_1s_1, \dots, w_d s_d])$
 - 2: $\mathbf{x} = [x_1, x_2, \dots, x_n] \leftarrow [0, \dots, 0]$
 - 3: **for** $b = 1, \dots, B$ **do** $\dots O(B)$
 - 4: Use alias method to sample $t \sim P(t|\mathbf{w})$
 - 5: Use alias method to sample $j \sim P(j|t)$
 - 6: $x_j \leftarrow x_j + 1$
 - 7: **end for**
 - 8: Pass \mathbf{x} to Algorithm 3
-

Algorithm 3 Final Prediction Phase

- 1: Find the top- C elements in \mathbf{x} , i.e., $|\mathcal{C}(\mathbf{w})| = C$ and
 $\mathcal{C}(\mathbf{w}) \leftarrow \{j | x_j \geq x_{j'}, \forall j' \notin \mathcal{C}(\mathbf{w})\}$ $\dots O(B)$
 - 2: **for** $j \in \mathcal{C}(\mathbf{w})$ **do** $\dots O(Cd)$
 - 3: Compute inner product $\mathbf{w}^T \mathbf{h}_j$
 - 4: **end for**
 - 5: Output: indexes of the top- K elements of
 $\{\mathbf{w}^T \mathbf{h}_{j_1}, \dots, \mathbf{w}^T \mathbf{h}_{j_C}\}$ $\dots O(C)$
-

with high probability if the sample size B is sufficiently large. We will provide the detailed theoretical analysis in Section 4.

3.2 Handling both positive and negative values

Next we discuss how to handle the case when both H and \mathbf{w} can be negative. To extend our approach to the general case, we need to make sure that we don't violate the principle of getting larger x_j for index j with larger inner product. In order to achieve this goal, we keep the probability distribution to be non-negative by taking absolute values of all the elements of H and \mathbf{w} , that is, $P(j, t|\mathbf{w}) \propto |w_t h_{jt}|$. We then use the following sampling procedure to estimate the expectation of $E_{P(j,t|\mathbf{w})}[\text{sgn}(w_t h_{jt})]$, which is proportional to $\mathbf{w}^T \mathbf{h}_j$ (see Lemma 1 in Section 4). In the pre-processing phase, we build alias tables for the probability distribution $P(j|t) = \frac{|h_{jt}|}{\sum_{j=1}^n |h_{jt}|}$ for each $t \in \{1, \dots, d\}$. Then in the query-dependent phase, we construct alias table with $P(t|\mathbf{w}) \propto \sum_{j=1}^n |w_t h_{jt}|$ to sample t given the current query. The marginal distribution then becomes

$$\begin{aligned}
 P(j|\mathbf{w}) &= \sum_{t=1}^d P(j, t|\mathbf{w}) \\
 &= \sum_{t=1}^d P(j|t)P(t|\mathbf{w}) \\
 &\propto \sum_{t=1}^d |w_t h_{jt}|. \tag{5}
 \end{aligned}$$

We sample t and then j for B times, but each time a pair of index (j, t) is sampled, we increase x_j by $\text{sgn}(w_t h_{jt})$

instead of 1 for non-negative cases. The detailed procedure and time complexity of each step are shown in Algorithm 4. The total time complexity is then $O(d + B + Cd)$, while the naive time complexity is $O(nd)$. Usually, the budget size C is much smaller than n and in practice $B \approx n$ gives good results, so the proposed algorithm is much faster than the naive search approach.

Algorithm 4 Sampling-MIPS for general cases

Pre-processing:

- 1: $s_t \leftarrow \sum_{j=1}^n |h_{jt}|, \forall t = 1, \dots, d$

- 2: Use alias table to construct

$$P(j|t) \leftarrow \text{multinomial}([|h_{1t}|, \dots, |h_{nt}|]), \forall t$$

Candidate Screening:

- 3: Use alias table to construct $\dots O(d)$

$$P(t|\mathbf{w}) \leftarrow \text{multinomial}([|w_1 s_1|, \dots, |w_d s_d|])$$

- 4: $\mathbf{x} = [x_1, x_2, \dots, x_n] \leftarrow [0, \dots, 0]$

- 5: Specify sample size B

- 6: **for** $b = 1, \dots, B$ **do** $\dots O(B)$

- 7: Use alias method to sample $t \sim P(t|\mathbf{w})$

- 8: Use alias method to sample $j \sim P(j|t)$

- 9: $x_j \leftarrow x_j + \text{sgn}(w_t h_{jt})$

- 10: **end for**

Final Prediction Phase:

$$\dots O(B + Cd)$$

- 11: See Algorithm 3.
-

The key idea of Sampling-MIPS is to construct the ‘‘scores’’ (x_j) to estimate the n inner products. Index with larger inner product will have larger ‘‘score’’ after the entire sampling procedure. Yet directly sampling j is hard, we define a probability distribution for (j, t) and sample indexes using the two-step procedure. This method is expected to succeed with high probability if sampled enough times. When B is large enough, the result won't change a lot. Intuitively, we know that the required best sample size B to achieve certain accuracy relies on the relative gap between the n inner products. We will give a detailed analysis of the error bound for Sampling-MIPS in Section 4.

4 Mathematical Analysis of Sampling-MIPS

In the following we analyze the Sampling-MIPS algorithm for general case (with both positive and negative values). Due to the space limits, the proofs are deferred to the supplementary material. Recall that for the general case, the marginal distribution of j is $P(j|\mathbf{w}) \propto \sum_{t=1}^d |w_t h_{jt}|$, and our goal is to prove that $E[x_j]$ will be proportional to $\mathbf{w}^T \mathbf{h}_j$. For convenience, we define $x_{jb} = \text{sgn}(w_t h_{jt})$ if index pair (j, t) is sampled in the b -th sampling step, and $x_{jb} = 0$ otherwise. Based on this definition, we have $x_j = \sum_{b=1}^B x_{jb}$. The mean and variance of x_{jb} and x_j can be established by the following lemmas.

Lemma 1. Define $S = \sum_{t=1}^d \sum_{j=1}^n |w_t h_{jt}|$ and $c_j = \mathbf{w}^T \mathbf{h}_j$, then $E[x_{jb}] = \frac{c_j}{S}$ and $E[x_j] = \frac{Bc_j}{S}$.

Lemma 2. Define $a_j = \sum_{t=1}^d |w_t h_{jt}|$, then $E[(x_{jb} - x_{mb})^2] = \frac{a_j + a_m}{S}, \forall j \neq m$. Therefore, from Lemma 1, we have

$$\text{Var}(x_{jb} - x_{mb}) = \frac{a_j + a_m}{S} - \frac{(c_j - c_m)^2}{S^2}$$

and $\text{Var}(x_j - x_m) = B \text{Var}(x_{jb} - x_{mb})$.

Theorem 1. If index m has the maximum inner product, we define for $j \in \{j : c_j < c_m\}$

$$\begin{aligned} \lambda_j &= \frac{c_m - c_j}{S}, \\ T_j &= \frac{(S + c_m - c_j)(c_m - c_j)}{S^2}, \\ M_j &= \frac{S^2}{(S + c_m - c_j)^2}, \\ \sigma_j^2 &= \text{Var}(x_{jb} - x_{mb}). \end{aligned}$$

Also let $\lambda = \min_{j \in \{c_j < c_m\}} \lambda_j$, $n' = \#\{j : c_j < c_m\}$ and assume $n' \neq 0$. If sample size B satisfies the following equation for all $j \in \{j : c_j < c_m\}$,

$$B \geq \frac{\log \frac{n'}{\delta}}{M_j \sigma_j^2 \left[\left(1 + \frac{T_j}{\sigma_j^2}\right) \log \left(1 + \frac{T_j}{\sigma_j^2}\right) - \frac{T_j}{\sigma_j^2} \right]}, \quad (6)$$

implying that when $B \sim O(\frac{d}{\lambda} \log n')$, then with error probability $\delta \in (0, 1)$ we have

$$P(x_m > x_j, \forall j \in \{c_j < c_m\}) \geq 1 - \delta.$$

This Theorem states that when number of samples $B = O(\frac{d}{\lambda} \log n')$, Sampling-MIPS will be able to distinguish candidate vector m from the rest candidates with smaller inner products. When candidate vector m has the maximum inner product, n' can be replaced by n , so Theorem 1 suggests that Sampling-MIPS can get precise results when $B = O(\frac{d}{\lambda} \log n)$. Also, when the probability for sampling different indexes are very close, we need larger sample size to distinguish them. This is reflected in Theorem 1 via the numerator of λ , $\min_{\{j: c_j < c_m\}} (c_m - c_j)$. Note that $\lambda < 1$ here.

Theorem 2. Assume $w_t h_{jt} \geq 0, \forall (j, t) \in \{1, \dots, n\} \times \{1, \dots, d\}$ and $c_j \stackrel{iid}{\sim} \text{Exp}(\beta)$. In Sampling-MIPS algorithm, if we take $C = B$, where C means we calculate inner products of indexes with top- C scores, then when $B \geq \frac{\rho}{\alpha} n$, where $\alpha \in (0, \frac{1}{2})$ and $\rho \in (0, \frac{\alpha d(n-1)}{(d+1)n})$,

$$P(\text{not identifying maximum inner product}) \leq O\left(\frac{1}{n^\rho}\right).$$

This result can also be generalized to inner products *i.i.d.* sampled from heavy-tailed distributions.

Definition 1. The distribution of a random variable X with c.d.f. $F(x)$ is said to have a heavy (right) tail if

$$\lim_{x \rightarrow +\infty} e^{\mu x} (1 - F(x)) = +\infty, \forall \mu > 0. \quad (7)$$

Note: Heavy-tailed distribution includes many common distributions like log-normal, log-cauchy, Lévy and Pareto distribution, etc.

Corollary 1. Assume $w_t h_{jt} \geq 0, \forall (j, t) \in \{1, \dots, n\} \times \{1, \dots, d\}$. If \mathcal{F} is a continuous, non-negative, heavy-tailed distribution, $c_j \stackrel{iid}{\sim} \mathcal{F}$ and $E[c_j^2] < \infty$, then when $B \geq \frac{\rho}{\alpha} n$, where $\alpha \in (0, 1)$ and $\rho \in (0, \frac{\alpha d(n-1)}{(d+1)n})$, take $C = B$ in Sampling-MIPS algorithm, we have

$$P(\text{not identifying maximum inner product}) \leq O\left(\frac{1}{n^\rho}\right).$$

Theorem 2 and Corollary 1 imply that we only need sample size B to be $\frac{\rho}{\alpha} n$ to make sure that the maximum inner product to be identified. The total time complexity of Sampling-MIPS is $O(d + B + Bd)$ and the time complexity of Naive-MIPS is $O(nd)$. When $B = \frac{\rho}{\alpha} n$,

$$d + B + Bd = d + \frac{\rho}{\alpha} (d+1)n < d + \frac{d(n-1)}{(d+1)n} n(d+1) = nd.$$

Therefore, the conditions of parameter α and ρ ensures that Sampling-MIPS is more efficient than Naive-MIPS.

Note that in Theorem 2, α can be any number in $(0, \frac{1}{2})$, but Corollary 1 only requires $\alpha \in (0, 1)$. ρ can be any number in $(0, \frac{\alpha d(n-1)}{(d+1)n})$. Sampling-MIPS algorithm is faster with smaller ρ , but can also yield larger error probability. So ρ should be chosen according to the trade-off between efficiency and tolerance of error.

5 Experimental Results

In this section, we compare the performance of Sampling-MIPS with other state-of-the-art algorithms. In order to have a fair comparison, all the algorithms are implemented in highly optimized C++ code using the platform provided by Yu et al. (2017). We also follow the same evaluation criteria used in (Yu et al., 2017) to measure the precision and speedup of MIPS algorithms. To compute the precision, we randomly select 2,000 queries and compute the average of $\text{prec}@K$ for $K = 1, 5, 10$, where $\text{prec}@K$ is defined as:

$$\text{prec}@K = \frac{\#\{\text{selected top-}K\} \cap \{\text{true top-20}\}}{K}.$$

The true top-20 candidates are identified by naive MIPS approach, and the selected top- K candidates are identified by approximate MIPS algorithms. For the speed of MIPS algorithms, we report the speedup of each algorithm over the naive linear search approach (computing all the inner

products). All the experiments are conducted on a server with Intel Xeon E5-2640 CPU and only one core is used for the experiments.

5.1 Datasets and parameter settings

We compare the results on several real recommender system datasets, including MovieLens-20M (Harper and Konstan, 2016) and Netflix. MovieLens-20M contains 27,278 movies and 138,493 users while Netflix contains 17,770 items and 480,189 users. For each dataset, we obtain the low-rank matrices W and H by matrix factorization model using the open source matrix factorization package LIBPMF (Yu et al., 2014a). We use $\lambda = 0.05$ for both MovieLens-20M and Netflix datasets to get latent features W and H with rank $d = 50, 100, 200$.

To test the performance of algorithms under different scenarios, we also conduct experiments on synthetic datasets. The candidate matrix H and query vector w are generated from $\mathcal{N}(0, 10)$, with $n = 200000$ and $d \in \{50, 100, 200\}$ (used in Figure 3).

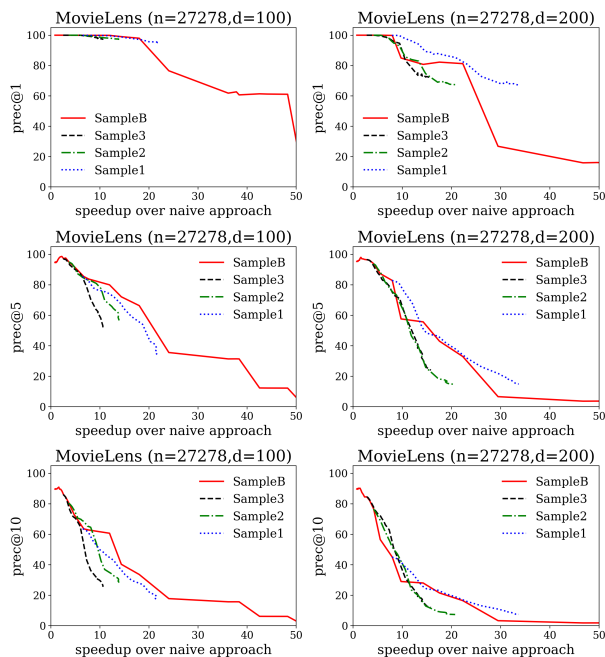


Figure 2: Comparison of four different parameter settings of Sampling-MIPS on MovieLens data.

5.2 Sampling-MIPS with different parameters

In Sampling-MIPS, there are two parameters that determine the trade-off between speedup and precision: sample size B and the size of budget $\mathcal{C}(w)$ denoted as C . We describe four methods of choosing B and C in this paper:

- SampleB: Set $C = B$ and B ranges from $0.95n$ to $0.001n$.

- Sample1: Fix sample size $B = 0.1n$ and set $C = r \times B$, where r ranges from 1 to 0.06.
- Sample2: Fix sample size $B = 0.2n$ and set $C = r \times B$, where r ranges from 1 to 0.06.
- Sample3: Fix sample size $B = 0.3n$ and set $C = r \times B$, where r ranges from 1 to 0.06.

The comparison of SampleB, Sample1, Sample2 and Sample3 on MovieLens-20M data are shown in Figure 2. In general we found that they are similar in most cases, but the approach of SampleB is the most stable one, so we choose SampleB and compare it with other MIPS algorithms in the following subsections.

5.3 Experimental results on real datasets

We report the comparisons on MovieLens-20M, Netflix and synthetic data with different dimensionality in Figure 3. It is easy to see that Sampling-MIPS outperforms PCA-MIPS, LSH-MIPS and Diamond-MSIPS in most of the cases. It is worthwhile to mention that on MovieLens-20M data with $d = 50, 100$, when speedup is 10, Sampling-MIPS yields $prec@1 \approx 100\%$ and $prec@5 \approx 80\%$, while none of the other methods can reach $prec@1 \geq 50\%$ and $prec@5 \geq 20\%$. See tables in the supplementary material for more details.

5.4 Comparison to Greedy-MIPS

Greedy-MIPS is the most competitive algorithm among all the previous approaches in practice. Although it does not have a theoretical guarantee, it usually yields better results than Sampling-MIPS and all the other approaches on real datasets (see second row of Figure 4). This means Assumption (1), although clearly not true in general, is valid to some extent on real datasets. However, as we will show in the following experiments, when Assumption (1) does not hold, Greedy-MIPS will perform poorly.

In the first row of Figure 4, we construct the synthetic data with irregular pattern which violates Assumption (1). The i -th row of candidate matrix H is generated from $\mathcal{N}(\frac{200000}{i}, \frac{i}{10})$ and query vector w is generated from $\mathcal{N}(1, 0.1)$, with $n = 200000$ and $d = 2000$. The results are shown in the first row of Figure 4. Clearly, in this setting Assumption (1) does not hold, so $prec@10$ of Greedy-MIPS drops to 40% even when speedup ≤ 5 . In comparison, our algorithm still gets almost perfect precision under the same speedup. In comparison, if we generate synthetic datasets with all the rows in H follow the same uniform distribution, the results in the third row of Figure 4 shows Greedy-MIPS performs well in this case since Assumption (1) is mostly satisfied.

5.5 Implementation for sparse matrices

Sampling-MIPS draws samples based on a probability distribution of each entry. As zero entries will never be sampled, we can improve the efficiency of sampling and reduce the memory cost by only considering the non-zeros. Therefore, it is nature that Sampling-MIPS will work especially good when dealing with sparse matrices. We utilize Eigen (version 3) (Guennebaud et al., 2010) to input data and implement Sampling-MIPS algorithm in a smarter way by only dealing with non-zero entries.

To test our implementation, we generate sparse candidate matrices $H \sim Uniform(0, 1)$ and sparse query vectors $w \sim \mathcal{N}(0, 10)$ with $n = 20000$ and $d \in \{200, 1000\}$. The experimental results in Figure 5 demonstrate that our algorithm performs well on sparse data, while many existing methods such as Greedy-MIPS cannot handle sparse data.

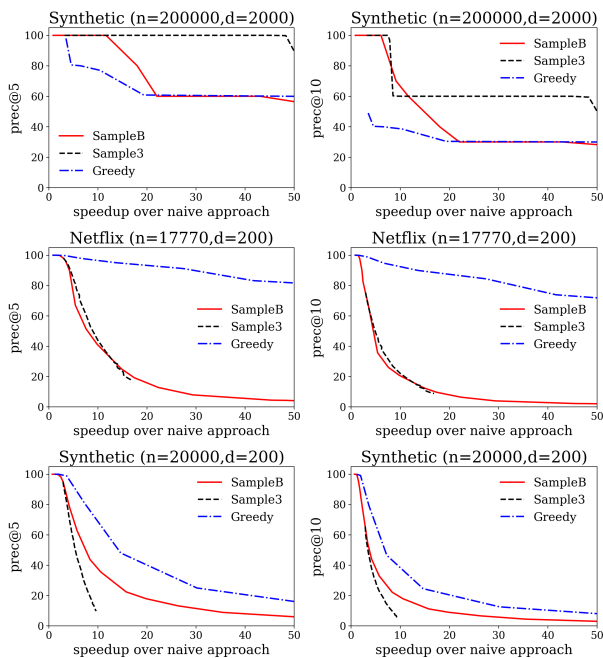


Figure 4: Comparison of Sampling-MIPS and Greedy-MIPS on Netflix and synthetic datasets.

6 Conclusion

We propose a Sampling-MIPS algorithm for approximate maximum inner product search. The core idea is to develop a novel way to sample from the distribution according to the magnitude of inner product, but without explicitly constructing such distribution. Each sample only costs constant time, and we show that the algorithm can recover the true ranking of inner products when sample size is large enough. Experimental results on real and synthetic datasets show that our algorithm significantly outperforms most of previous approaches including PCA-MIPS, LSH-

MIPS and Diamond-MSIPS. Compared to Greedy-MIPS, our algorithm does not rely on any vulnerable assumptions and is safer to be applied to all kinds of datasets. Moreover, our algorithm can be perfectly adapted to data of sparse format.

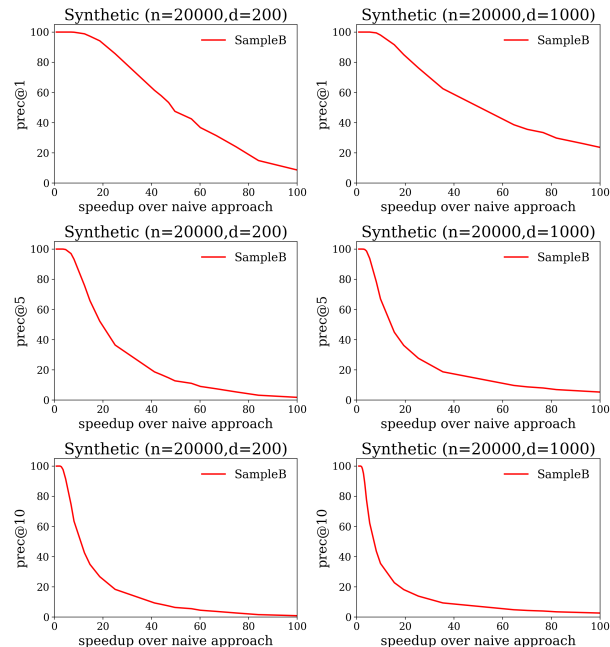


Figure 5: Results of Sampling-MIPS for data of sparse format with $n = 20000$ and $d \in \{200, 1000\}$.

7 Acknowledgement

We are grateful to Xiaodong Li and James Sharpnack for helpful comments. This work was partially supported by NSF IIS-1719097, Intel, Google Cloud and AITRICS.

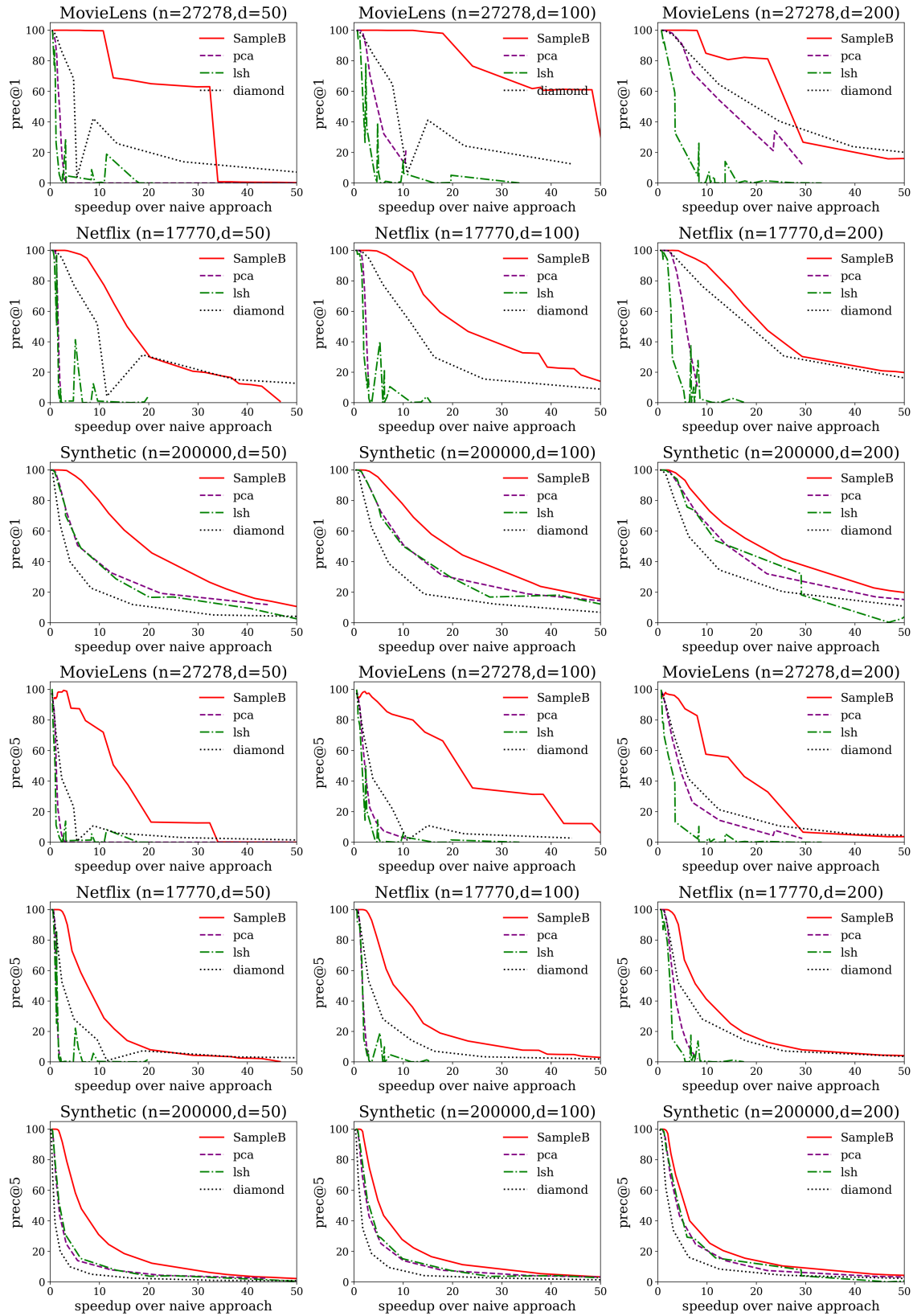


Figure 3: Comparison of Sampling-MIPS to PCA-MIPS, LSH-MIPS, Diamond-MSIPS on MovieLens, Netflix and synthetic datasets. The synthetic datasets have their H and w generated from $\mathcal{N}(0, 10)$.

References

- Rohit Babbar and Bernhard Schölkopf. Dismec: distributed sparse machines for extreme multi-label classification. In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*, pages 721–729. ACM, 2017.
- Yoram Bachrach, Yehuda Finkelstein, Ran Gilad-Bachrach, Liran Katzir, Noam Koenigstein, Nir Nice, and Ulrich Paquet. Speeding up the xbox recommender system using a euclidean transformation for inner-product spaces. In *Proceedings of the 8th ACM Conference on Recommender systems*, pages 257–264. ACM, 2014.
- Grey Ballard, Tamara G Kolda, Ali Pinar, and C Seshadhri. Diamond sampling for approximate maximum all-pairs dot-product (mad) search. In *Data Mining (ICDM), 2015 IEEE International Conference on*, pages 11–20. IEEE, 2015.
- George Bennett. Probability inequalities for the sum of independent random variables. *Journal of the American Statistical Association*, 57(297):33–45, 1962.
- Patrick H Chen, Si Si, Sanjiv Kumar, Yang Li, and Cho-Jui Hsieh. Learning to screen for fast softmax inference on large vocabulary neural networks. In *ICLR*, 2019.
- Edith Cohen and David D Lewis. Approximating matrix multiplication for pattern recognition tasks. *Journal of Algorithms*, 30(2):211–252, 1999.
- Gaël Guennebaud, Benoît Jacob, et al. Eigen v3. <http://eigen.tuxfamily.org>, 2010.
- F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 5(4):19, 2016.
- Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In *Proceedings of the Thirtieth Annual ACM Symposium on Theory of Computing*, pages 604–613, 1998.
- Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8), 2009.
- Yury A Malkov and Dmitry A Yashunin. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE transactions on pattern analysis and machine intelligence*, 2018.
- Behnam Neyshabur and Nathan Srebro. On symmetric and asymmetric lshs for inner product search. *arXiv preprint arXiv:1410.5518*, 2014.
- Parikshit Ram and Alexander G Gray. Maximum inner-product search using cone trees. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 931–939. ACM, 2012.
- Anshumali Shrivastava and Ping Li. Asymmetric lsh (alsh) for sublinear time maximum inner product search (mips). In *Advances in Neural Information Processing Systems*, pages 2321–2329, 2014a.
- Anshumali Shrivastava and Ping Li. Improved asymmetric locality sensitive hashing (alsh) for maximum inner product search (mips). *arXiv preprint arXiv:1410.5410*, 2014b.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112, 2014.
- Alastair J Walker. An efficient method for generating discrete random variables with general distributions. *ACM Transactions on Mathematical Software (TOMS)*, 3(3):253–256, 1977.
- Hsiang-Fu Yu, Cho-Jui Hsieh, Si Si, and Inderjit S Dhillon. Parallel matrix factorization for recommender systems. *Knowledge and Information Systems*, 41(3):793–819, 2014a.
- Hsiang-Fu Yu, Prateek Jain, Purushottam Kar, and Inderjit Dhillon. Large-scale multi-label learning with missing labels. In *International conference on machine learning*, pages 593–601, 2014b.
- Hsiang-Fu Yu, Cho-Jui Hsieh, Qi Lei, and Inderjit S Dhillon. A greedy approach for budgeted maximum inner product search. In *Advances in Neural Information Processing Systems*, pages 5459–5468, 2017.
- Minjia Zhang, Wenhan Wang, Xiaodong Liu, Jianfeng Gao, and Yuxiong He. Navigating with graph representations for fast and scalable decoding of neural language models. In *Advances in Neural Information Processing Systems*, pages 6311–6322, 2018.