

# Fast Generating A Large Number of Gumbel-Max Variables

Yiyan Qi<sup>1</sup>, Pinghui Wang<sup>2,1,\*</sup>, Yuanming Zhang<sup>1</sup>, Junzhou Zhao<sup>1,\*</sup>, Guangjian Tian<sup>3</sup>,  
Xiaohong Guan<sup>2,1,4</sup>

<sup>1</sup>MOE Key Laboratory for Intelligent Networks and Network Security, Xi'an Jiaotong University, Xi'an, China

<sup>2</sup>Shenzhen Research Institute of Xi'an Jiaotong University, Shenzhen, China

<sup>3</sup>Huawei Noah's Ark Lab, Hong Kong

<sup>4</sup>Department of Automation and NLIST Lab, Tsinghua University, Beijing, China

{qiyiyan,zhangyuanming}@stu.xjtu.edu.cn, {phwang,xhguan}@mail.xjtu.edu.cn,  
junzhou.zhao@xjtu.edu.cn, Tian.Guangjian@huawei.com

## ABSTRACT

The well-known Gumbel-Max Trick for sampling elements from a categorical distribution (or more generally a nonnegative vector) and its variants have been widely used in areas such as machine learning and information retrieval. To sample a random element  $i$  (or a Gumbel-Max variable  $i$ ) in proportion to its positive weight  $v_i$ , the Gumbel-Max Trick first computes a Gumbel random variable  $g_i$  for each positive weight element  $i$ , and then samples the element  $i$  with the largest value of  $g_i + \ln v_i$ . Recently, applications including similarity estimation and graph embedding require to generate  $k$  independent Gumbel-Max variables from high dimensional vectors. However, it is computationally expensive for a large  $k$  (e.g., hundreds or even thousands) when using the traditional Gumbel-Max Trick. To solve this problem, we propose a novel algorithm, *FastGM*, that reduces the time complexity from  $O(kn^+)$  to  $O(k \ln k + n^+)$ , where  $n^+$  is the number of positive elements in the vector of interest. Instead of computing  $k$  independent Gumbel random variables directly, we find that there exists a technique to generate these variables in descending order. Using this technique, our method *FastGM* computes variables  $g_i + \ln v_i$  for all positive elements  $i$  in descending order. As a result, *FastGM* significantly reduces the computation time because we can stop the procedure of Gumbel random variables computing for many elements especially for those with small weights. Experiments on a variety of real-world datasets show that *FastGM* is orders of magnitude faster than state-of-the-art methods without sacrificing accuracy and incurring additional expenses.

## CCS CONCEPTS

• **Mathematics of computing** → **Probabilistic algorithms**; • **Information systems** → **Similarity measures**; • **Theory of computation** → **Sketching and sampling**.

## KEYWORDS

Gumbel-Max Trick, Sketching, Graph embedding

\*Corresponding Author.

This paper is published under the Creative Commons Attribution 4.0 International (CC-BY 4.0) license. Authors reserve their rights to disseminate the work on their personal and corporate Web sites with the appropriate attribution.

WWW '20, April 20–24, 2020, Taipei, Taiwan

© 2020 IW3C2 (International World Wide Web Conference Committee), published under Creative Commons CC-BY 4.0 License.

ACM ISBN 978-1-4503-7023-3/20/04.

<https://doi.org/10.1145/3366423.3380160>

## ACM Reference Format:

Yiyan Qi, Pinghui Wang, Yuanming Zhang, Junzhou Zhao, Guangjian Tian, and Xiaohong Guan. 2020. Fast Algorithm for Generating A Large Number of Gumbel-Max Variables. In *Proceedings of The Web Conference 2020 (WWW '20)*, April 20–24, 2020, Taipei, Taiwan. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3366423.3380160>

## 1 INTRODUCTION

The Gumbel-Max Trick [1] is a popular technique for sampling elements from a categorical distribution (or more generally a nonnegative vector). Given a nonnegative vector  $\vec{v} = (v_1, \dots, v_n)$  where each element  $v_i \in \mathbb{R}_{\geq 0}$ , the Gumbel-Max Trick computes a random variable  $s(\vec{v})$  as

$$s(\vec{v}) = \arg \max_{i \in N_{\vec{v}}^+} g_i + \ln v_i,$$

where  $N_{\vec{v}}^+ \triangleq \{i: v_i > 0, i = 1, \dots, n\}$  is the set of indices of positive elements in  $\vec{v}$ ,  $g_i \triangleq -\ln(-\ln a_i)$  and  $a_i$  is a random variable drawn from the uniform distribution  $\text{UNI}(0, 1)$ . We call  $s(\vec{v})$  a Gumbel-Max variable of vector  $\vec{v}$  and the probability of selecting  $i$  as the Gumbel-Max variable is  $P(s(\vec{v}) = i) = \frac{v_i}{\sum_{j=1}^n v_j}$ . The Gumbel-Max Trick and its variants have been used widely in many areas.

**Similarity estimation.** Similarity estimation lies at the core of many data mining and machine learning applications, such as web duplicate detection [2, 3], collaborate filtering [4], and association rule learning [5]. To efficiently estimate the similarity between two vectors, several algorithms [6–9] compute  $k$  random variables  $-\frac{\ln a_{i,1}}{v_i}, \dots, -\frac{\ln a_{i,k}}{v_i}$  for each positive element  $v_i$  in  $\vec{v}$ , where  $a_{i,1}, \dots, a_{i,k}$  are independent random variables drawn from the uniform distribution  $\text{UNI}(0, 1)$ . Then, these algorithms build a sketch (or called *Gumbel-Max sketch* in this paper) of vector  $\vec{v}$  consisting of  $k$  registers, and each register records  $s_j(\vec{v})$  where

$$s_j(\vec{v}) = \arg \min_{i \in N_{\vec{v}}^+} -\frac{\ln a_{i,j}}{v_i}, \quad 1 \leq j \leq k.$$

We find that  $s_j(\vec{v})$  is exactly a Gumbel-Max variable of vector  $\vec{v}$  as  $\arg \min_{i \in N_{\vec{v}}^+} -\frac{\ln a_{i,j}}{v_i} = \arg \max_{i \in N_{\vec{v}}^+} \ln v_i - \ln(-\ln a_{i,j})$ . Let  $\mathbb{1}(x)$  be an indicator function. Yang et al. [6–8] use  $\frac{1}{k} \sum_{j=1}^k \mathbb{1}(s_j(\vec{u})) = s_j(\vec{v})$  to estimate the *weighted Jaccard similarity* of two nonnegative vectors  $\vec{u}$  and  $\vec{v}$  which is defined by

$$\mathcal{J}_W(\vec{u}, \vec{v}) \triangleq \frac{\sum_{i=1}^n \min\{u_i, v_i\}}{\sum_{i=1}^n \max\{u_i, v_i\}}.$$

Recently, Moulton et al. [9] prove that the expectation of estimate  $\frac{1}{k} \sum_1^k \mathbb{1}(s_j(\vec{u}) = s_j(\vec{v}))$  actually equals the *probability Jaccard similarity*, which is defined by

$$\mathcal{J}_{\mathcal{P}}(\vec{u}, \vec{v}) \triangleq \sum_{i \in N_{\vec{v}, \vec{u}}^+} \frac{1}{\sum_{l=1}^n \max\left(\frac{u_l}{u_i}, \frac{v_l}{v_i}\right)}.$$

Here,  $N_{\vec{v}, \vec{u}}^+ \triangleq \{i: v_i > 0 \wedge u_i > 0, i = 1, \dots, n\}$  is the set of indices of positive elements in both  $\vec{v}$  and  $\vec{u}$ . Compared with the weighted Jaccard similarity  $\mathcal{J}_{\mathcal{W}}$ , Moulton et al. demonstrate that the probability Jaccard similarity  $\mathcal{J}_{\mathcal{P}}$  is scale-invariant and more sensitive to changes in vectors. Moreover, each function  $s_j(\vec{v})$  maps similar vectors into the same value with high probability. Therefore, similar to the regular locality-sensitive hashing (LSH) schemes [10–12], one can use these Gumbel-Max sketches to build LSH index for fast similarity search in a large dataset, which is capable to search similar vectors for any query vector in sub-linear time.

**Graph embedding.** Recently, graph embedding attracts a lot of attention. A variety of graph embedding methods have been developed to transform a graph into a low-dimensional space in which each node is represented by a low-dimensional vector and meanwhile the proximities of nodes are preserved. With node embeddings, many off-the-shelf data mining and machine learning algorithms can be applied on graphs such as node classification and link prediction. Lately, Yang et al. [13] reveal that existing methods are computationally intensive especially for large graphs. To address this challenge, Yang et al. build a Gumbel-Max sketch described above for each row of the graph’s self-loop-augmented adjacency matrix and use these sketches as the first-order node embeddings of nodes. To capture the high-order node proximity, they propose a fast method *NodeSketch* to recursively generate  $r$ -order node embeddings (i.e., Gumbel-Max sketches) based on the graph’s self-loop-augmented adjacency matrix and  $(r - 1)$ -order node embeddings.

**Machine learning.** Lorberbom et al. [14] use the Gumbel-Max Trick to reparameterize discrete variant auto-encoding (VAE). Buchnik et al. [15] apply the Gumbel-Max Trick to select training examples, and sub-epochs of sampled examples preserve rich structural properties of the full training data. These studies show that the trick not only significantly accelerates the convergence of the loss function, but also improves the accuracy. Other examples include reinforcement learning [16], and integer linear programs [17]. Lately, Eric et al. [18] extend the Gumbel-Max Trick to embed discrete Gumbel-Max variables in a continuous space, which enables us to compute the gradients of these random variables easily. This technique has also been used for improving the performance of neural network models such as Generative Adversarial Networks (GAN) [19] and attention models [20].

Despite the wide use of the Gumbel-Max Trick in various domains, it is expensive for the above algorithms to handle large dimensional vectors. Specifically, the time complexity of generating a Gumbel-Max sketch consisting of  $k$  independent Gumbel-Max variables is  $O(n^+k)$ , where  $n^+ = |N^+|$  is the number of positive elements in the vector of interest. In practice,  $k$  is usually set to be hundreds or even thousands when selecting training samples [15], estimating probability Jaccard similarity [9], and learning *NodeSketch* graph embedding [13]. To solve this problem, we propose a novel method,

*FastGM*, to fast compute a Gumbel-Max sketch. The basic idea behind *FastGM* can be summarized as follows. For each element  $v_i > 0$  in  $\vec{v}$ , we find that  $k$  random variables  $-\frac{\ln a_{i,1}}{v_i}, \dots, -\frac{\ln a_{i,k}}{v_i}$  can be generated in ascending order. That is, we can generate a sequence of  $k$  tuples  $\left(-\frac{\ln a_{i,i_1}}{v_i}, i_1\right), \dots, \left(-\frac{\ln a_{i,i_k}}{v_i}, i_k\right)$ , where  $-\frac{\ln a_{i,i_1}}{v_i} < \dots < -\frac{\ln a_{i,i_k}}{v_i}$  and  $i_1, \dots, i_k$  is a random permutation of integers  $1, \dots, k$ . We propose to sort  $kn^+$  random variables of all  $n^+$  positive elements and compute these variables sequentially in ascending order. Then, we model the procedure of computing the Gumbel-Max sketch  $(s_1(\vec{v}), \dots, s_k(\vec{v}))$  of vector  $\vec{v}$  as a Balls-and-Bins model. Specifically, randomly throw balls one by one into  $k$  empty bins, where each ball is assigned a random variable in ascending order. When no bins are empty, we early stop the procedure, and then each  $s_j(\vec{v})$ ,  $j = 1, \dots, k$ , records the random variable of the first ball thrown into bin  $j$ . We summarize our main contributions as:

- We introduce a simple Balls-and-Bins model to interpret the procedure of computing the Gumbel-Max sketch of vector  $\vec{v}$ , i.e.,  $(s_1(\vec{v}), \dots, s_k(\vec{v}))$ . Using this stochastic process model, we propose a novel algorithm, called *FastGM*, to reduce the time complexity of computing the Gumbel-Max sketch  $(s_1(\vec{v}), \dots, s_k(\vec{v}))$  from  $O(n^+k)$  to  $O(k \ln k + n^+)$ , which is achieved by avoiding calculating all  $k$  variables  $-\frac{\ln a_{i,1}}{v_i}, \dots, -\frac{\ln a_{i,k}}{v_i}$  for each  $i \in N_{\vec{v}}^+$ .
- We conduct experiments on a variety of real-world datasets for estimating probability Jaccard similarity and learning *NodeSketch* graph embeddings. The experimental results demonstrate that our method *FastGM* is orders of magnitude faster than the state-of-the-art methods without incurring any additional cost.

The rest of this paper is organized as follows. Section 2 summarizes related work. The problem formulation is presented in Section 3. Section 4 presents our method *FastGM*. The performance evaluation and testing results are presented in Section 5. Concluding remarks then follow.

## 2 RELATED WORK

### 2.1 Jaccard Similarity Estimation

Broder et al. [11] proposed the first sketch method *MinHash* to compute the Jaccard similarity of two sets (or binary vectors). *MinHash* builds a sketch consisting of  $k$  registers for each set. Each register uses a hash function to keep track of the set’s element with the minimal hash value. To further improve the performance of *MinHash*, [5, 21, 22] developed several memory-efficient methods. Li et al. [23] proposed *One Permutation Hash* (OPH) to reduce the time complexity of processing each element from  $O(k)$  to  $O(1)$  but this method may exhibit large estimation errors because of the empty buckets. To solve this problem, several densification methods [24–27] were developed to set the registers of empty buckets according to the values of non-empty buckets’ registers.

Besides binary vectors, a variety of methods have also been developed to estimate generalized Jaccard similarity on weighted vectors. For vectors consisting of only nonnegative integer weights, Haveliwala et al. [28] proposed to add a corresponding number of replications of each element in order to apply the conventional *MinHash*. To handle more general real weights, Haeupler et al. [29] proposed to generate another additional replication with probability

that equals the floating part of an element’s weight. These two algorithms are computationally intensive when computing hash values of massive replications for elements with large weights. To solve this problem, [30, 31] proposed to compute hash values only for few necessary replications (i.e., “active indices”). ICWS [32] and its variations such as 0-bit CWS [33], CCWS [34], PCWS [35], I<sup>2</sup>CWS [36] were proposed to improve the performance of CWS [31]. The CWS algorithm and its variants all have the time complexity of  $O(n^+k)$ , where  $n^+$  is the number of elements with positive weights. Recently, Otmár [37] proposed another efficient algorithm *BagMinHash* for handling high dimensional vectors. *BagMinHash* is faster than ICWS when the vector has a large number of positive elements, e.g.,  $n^+ > 1,000$ , which may not hold for many real-world datasets. The above methods all estimate the weighted Jaccard similarity. Ryan et al. [9] proposed a Gumbel-Max Trick based sketching method,  $\mathcal{P}$ -MinHash, to estimate another novel Jaccard similarity metric, *probability Jaccard similarity*. They also demonstrated that the probability Jaccard similarity is scale-invariant and more sensitive to changes in vectors. However, the time complexity of  $\mathcal{P}$ -MinHash processing a weighted vector is  $O(n^+k)$ , which is not feasible for high-dimensional vectors.

## 2.2 Graph Embedding

DeepWalk [38] employed truncated random-walks to transform a network into sequences of nodes and learns node embedding using skip-gram model [39]. The basic idea behind DeepWalk is that two nodes should have similar embedding when they tend to co-occur in short random-walks. Node2vec [40] extended DeepWalk by introducing two kinds of search strategies, i.e., breadth- and depth-first search, into random-walk. LINE [41] explicitly defined the first-order and second-order proximities between nodes, and learned node embedding by minimizing the two proximities. GraRep [42] extended LINE to high-order proximities. Qiu et al. [43] unified the above methods into a general matrix factorization framework. Many recent works also learned embeddings on attributed graphs [44–47], partially labeled graphs [48–51], and dynamic graphs [52–55]. The above methods are computationally intensive and so many are prohibitive for large graphs. To solve this problem, recently, Yang et al. [13] developed a novel embedding method, NodeSketch, which builds a Gumbel-Max sketch for each row  $i$  of the graph’s self-loop-augmented adjacency matrix and uses it as the low-order node embedding of node  $i$ . To capture the high-order node proximity, NodeSketch recursively generates  $r$ -order node embeddings (i.e., Gumbel-Max sketches) based on the graph’s self-loop-augmented adjacency matrix and  $(r-1)$ -order node embeddings. NodeSketch requires time complexity  $O(n^+k)$  to compute the Gumbel-Max sketch of a node, which is expensive when existing a large number of nodes in the graph.

## 3 PROBLEM FORMULATION

We first introduce some notations and then formulate our problem. For a nonnegative vector  $\vec{v} = (v_1, \dots, v_n)$  and each element  $v_i \geq 0$ , let  $\vec{v}^* = (v_1^*, \dots, v_n^*)$  be the normalized vector of  $\vec{v}$ , where

$$v_i^* \triangleq \frac{v_i}{\sum_{j=1}^n v_j}, \quad i = 1, \dots, n.$$

Let  $N_{\vec{v}}^+ \triangleq \{i: v_i > 0, i = 1, \dots, n\}$  be the set of indices of positive elements in  $\vec{v}$ , and  $n_{\vec{v}}^+ \triangleq |N_{\vec{v}}^+|$  be its cardinality.

For each  $i = 1, \dots, n$ , we independently draw  $k$  random samples  $a_{i,1}, \dots, a_{i,k}$  from the uniform distribution  $\text{UNI}(0, 1)$ . Note that  $a_{i,1}, \dots, a_{i,k}$  are the same for different vectors. Given a nonnegative vector  $\vec{v}$ , we aim to fast compute its Gumbel-Max sketch  $\vec{s}(\vec{v}) = (s_1(\vec{v}), \dots, s_k(\vec{v}))$ , where

$$\begin{aligned} s_j(\vec{v}) &\triangleq \arg \max_{i \in N_{\vec{v}}^+} \ln v_i - \ln(-\ln a_{i,j}) \\ &\triangleq \arg \min_{i \in N_{\vec{v}}^+} -\frac{\ln a_{i,j}}{v_i}. \end{aligned}$$

To compute the Gumbel-Max sketches of a large collection of vectors (e.g., bag-of-words representations of documents), the straightforward method (also used in NodeSketch [13]) first instantiates variables  $a_{i,1}, \dots, a_{i,k}$  from  $\text{UNI}(0, 1)$  for each index  $i = 1, \dots, n$ . Then, for each nonnegative vector  $\vec{v}$ , it enumerates each  $i \in N_{\vec{v}}^+$  and compute  $-\frac{\ln a_{i,1}}{v_i}, \dots, -\frac{\ln a_{i,k}}{v_i}$ . The above method requires memory space  $O(nk)$  to store all  $[a_{i,j}]_{1 \leq i \leq n, 1 \leq j \leq k}$ , and time complexity  $O(kn_{\vec{v}}^+)$  to obtain the Gumbel-Max sketch  $\vec{s}(\vec{v})$  of each vector  $\vec{v}$ . We note that  $k$  is usually set to be hundreds or even thousands, therefore, the straightforward method costs a very large amount of memory space and time when the vector of interest has a large dimension, e.g.,  $n = 10^9$ . To reduce the memory cost, one can easily use hash techniques or random number generators with specific seeds (e.g., consistent random number generation methods in [24–26]) to generate each of  $a_{i,1}, \dots, a_{i,k}$  on the fly, which do not require to calculate and store variables  $[a_{i,j}]_{1 \leq i \leq n, 1 \leq j \leq k}$  in memory.

To address the computational challenge, in this paper, we propose a method that reduces the time complexity of computing sketch  $\vec{s}(\vec{v})$  from  $O(kn_{\vec{v}}^+)$  to  $O(k \ln k + n_{\vec{v}}^+)$ . In the follows, when no confusion raises, we simply write  $s_j(\vec{v})$  and  $n_{\vec{v}}^+$  as  $s_j$  and  $n^+$  respectively.

## 4 OUR METHOD

In this section, we first introduce the basic idea behind our method FastGM. Specially, we find that, for each element  $v_i > 0$  in vector  $\vec{v}$ , the variables  $-\frac{\ln a_{i,1}}{v_i}, \dots, -\frac{\ln a_{i,k}}{v_i}$  can be computed in ascending order and this procedure can be viewed as a Balls-and-Bins model. Then, we derive our model BBM-Mix to randomly put balls one by one into  $k$  empty bins, where each ball is assigned with a random variable in ascending order. Based on BBM-Mix, we next introduce our method FastGM to compute the Gumbel-Max sketch  $(s_1(\vec{v}), \dots, s_k(\vec{v}))$  of vector  $\vec{v}$ . Specifically, we model this procedure as randomly throwing balls arrived at different rates into  $k$  empty bins and each bin records the timestamp of the first arrived ball. When no bins are empty, we early stop the procedure, and then each  $s_j(\vec{v}), j = 1, \dots, k$ , records the random variable of the first ball thrown into bin  $j$ . At last, we discuss the time and space complexity of our method.

### 4.1 Basic Idea

In Figure 1, we provide an example of generating a Gumbel-Max sketch of a vector  $\vec{v} = (0.3, 0.1, 0.05, 0.05, 0.2, 0.07, 0.1, 0.03)$  to illustrate our basic idea, where we have  $n = 8$  and  $k = 10$ . Note that

$-\frac{\ln a_{i,j}}{v_i}$	$j=1$									$j=10$
$i=1$	7.3331	0.5148	5.4761	<b>0.1682</b>	<b>0.3682</b>	6.4619	7.8593	<b>2.1604</b>	2.0626	<b>0.0393</b>
	8.9781	31.3752	<b>2.8180</b>	6.9800	5.4039	5.0094	11.2953	11.2301	<b>0.0833</b>	0.7313
	2.4642	7.3742	21.2113	5.6169	10.8000	27.3706	5.2378	3.6822	5.6149	17.8556
$\dots$	12.0248	<b>0.4248</b>	17.4156	5.9572	3.1348	22.5307	29.0386	4.2071	0.3948	159.6455
	4.9848	6.3068	9.2985	0.9248	16.7813	<b>4.5591</b>	<b>1.5031</b>	2.9260	7.2455	3.0728
	22.4079	28.7365	2.8524	26.4955	1.7385	12.8640	5.2427	19.0820	9.1086	22.4502
	8.1884	3.7793	4.7013	7.8240	8.9715	9.5142	1.9349	3.8483	29.6742	15.1739
$i=8$	<b>1.4931</b>	3.1637	10.1017	16.0368	110.4560	16.4874	6.2923	48.4639	9.2852	37.3817

**Figure 1: An example of computing  $k$  independent Gumbel-Max variables  $s_1, \dots, s_k$  of a vector  $\vec{v} = (0.3, 0.1, 0.05, 0.05, 0.2, 0.07, 0.1, 0.03)$ , where  $k = 10$ . The Gumbel-Max variable  $s_j$  equals the index of the smallest element (i.e., the red and bold one) in the  $j$ -th column of matrix  $\left[-\frac{\ln a_{i,j}}{v_i}\right]_{1 \leq i \leq 8, 1 \leq j \leq 10}$ .**

we aim to fast compute each  $s_j = \arg \min_{1 \leq i \leq 8} -\frac{\ln a_{i,j}}{v_i}$ ,  $1 \leq j \leq 10$ . i.e., the index of the minimum element in each column  $j$  of matrix  $\left[-\frac{\ln a_{i,j}}{v_i}\right]_{1 \leq i \leq 8, 1 \leq j \leq 10}$ . We generate matrix  $\left[-\frac{\ln a_{i,j}}{v_i}\right]_{1 \leq i \leq 8, 1 \leq j \leq 10}$  based on the traditional Gumbel-Max Trick and mark the minimum element (i.e., the red and bold one indicating the Gumbel-Max variable) in each column  $j$ . We find that Gumbel-Max variables tend to equal index  $i$  with large weight  $v_i$ . For example, among the values of all Gumbel-Max variables  $s_1, \dots, s_{10}$ , index 1 with  $v_1 = 0.3$  appears 4 times, while index 3 with  $v_3 = 0.05$  never occurs. Furthermore, let  $R = 25$ , and  $R_i = \lceil R v_i^* \rceil$ , where  $\vec{v}^* = (v_1^*, \dots, v_8^*)$  is the normalized vector of  $\vec{v}$ . We have  $R_1 = 8, R_2 = 3, R_3 = 2, R_4 = 2, R_5 = 5, R_6 = 2, R_7 = 3$ , and  $R_8 = 1$ . We also find that each Gumbel-Max variable occurs as one of a row  $i$ 's Top- $R_i$  minimal elements. For example, the four Gumbel-Max variables occurring in the 1-st row are all among the Top- $R_1$  (i.e., Top-8) minimal elements. Based on the above insights, we derive our method FastGM, of which elements in each row can be generated in ascending order. As a result, we can early stop the computation when all the Gumbel-Max variables are acquired. Take Figure 1 as an example, compared with the straightforward method that requires to compute all  $nk = 80$  random variables, we compute  $s_1, \dots, s_k$  by only obtaining Top- $R_i$  minimal elements of each row  $i$ , which significantly reduces the computation to around  $\sum_{i=1}^8 R_i = 26$ .

## 4.2 A Building Block of Our Method FastGM

In this section, we introduce our model BBM-Mix to generate  $k$  random variables in ascending order for each positive element  $v_i$  of vector  $\vec{v}$ . We define variables

$$b_{i,j} = -\frac{\ln a_{i,j}}{v_i}, \quad j = 1, \dots, k. \quad (1)$$

We easily observe that  $b_{i,1}, \dots, b_{i,k}$  are equivalent to  $k$  independent random variables generated according to the exponential distribution  $\text{EXP}(v_i)$ . It is also well known that the first arrival time of a Poisson process with rate  $v_i$  is a random variable following the exponential distribution  $\text{EXP}(v_i)$ . Therefore, each  $b_{i,j}$  can also be viewed as the first arrival time of a Poisson process  $P_{i,j}$  with rate  $v_i$

and all Poisson processes  $\{P_{i,j} : i = 1, \dots, n, j = 1, \dots, k\}$  are independent with each other. Next, we show that  $b_{i,j}$  can be generated via the following *Balls-and-Bins Model* (in short, *BBM*).

**4.2.1 Basic BBM.** In the basic BBM, balls arrive independently according to a Poisson process  $\mathcal{P}_i$  with rate  $kv_i$ . When a ball arrives at time  $x$ , we select a bin  $j$  from  $k$  bins at random, and then put the ball into bin  $j$  as well as set  $b_{i,j} = \min(b_{i,j}, x)$ . The register  $b_{i,j}$  is used to record the timestamp of the first ball arriving in bin  $j$ . When no urn is empty, we stop the process because all  $b_{i,1}, \dots, b_{i,k}$  will not change anymore. Then, we easily find that the arrival of balls at any bin  $j$  is a Poisson process  $P_{i,j}$  with rate  $v_i$  because all balls are randomly split into  $k$  bins and the Poisson process is splittable [56]. The sequence of inter-arrival times of Poisson process  $\mathcal{P}_i$  are independent and they are identically distributed exponential random variables with mean  $\frac{1}{kv_i}$ . Therefore, the above BBM can be simulated as the following model.

**4.2.2 BBM-Hash.** We use a variable  $x_i$  to record the time when the latest ball arrives. Initialize  $x_i = 0$  and  $b_{i,j} = +\infty, j = 1, 2, \dots, k$  and then repeat the following Steps 1 to 3 until no bin is empty:

- Step 1:** Generate a random variable  $u$  according to the uniform distribution  $\text{UNI}(0, 1)$ ;
- Step 2:** Compute  $x_i = x_i - \frac{\ln u}{kv_i}$ ;
- Step 3:** Select a number  $j$  from  $\{1, \dots, k\}$  at random, and then put a ball into bin  $j$  as well as set  $b_{i,j} = \min\{b_{i,j}, x_i\}$ .

Clearly, it may require more than  $k$  iterations to fill all these  $k$  bins. To compute the number of required iterations is exactly the coupon collector's problem (see [57], Chapter 5.4.2) and we find that  $O(k \ln k)$  iterations are required in expectation. To reduce the number of iterations, we propose another model *BBM-Permutation*.

**4.2.3 BBM-Permutation.** For the above *BBM-Hash*, at Step 3, a nonempty bin  $j$  may be selected and the value of  $b_{i,j}$  will not change. Therefore, *BBM-Hash* may need more than one iterations to encounter an empty bin especially when few bins are empty. We use a variable  $m_i$  to keep track of the number of empty bins and  $k - m_i$  is the number of filled bins. Let  $z$  denote the number of iterations (balls) to encounter an empty bin. We easily find that  $z$  is a geometric random variable with success probability  $\frac{m_i}{k}$  when

---

**Algorithm 1:** Pseudo code of BBM-Mix, where function `GetNextBalls(i)` used in line 4 is defined in Algorithm 2.

---

```

/*  $v_i$  is the  $i^{\text{th}}$  element of vector  $\vec{v}$  */
Input :  $v_i$ 
Output:  $b_{i,1}, \dots, b_{i,k}$ 
1  $x_i \leftarrow 0; z_i \leftarrow 0; m_i \leftarrow k;$ 
2  $(\pi_{i,1}, \dots, \pi_{i,k}) \leftarrow (1, \dots, k);$ 
3 while  $m_i > 0$  do
4    $x_i, c \leftarrow \text{GetNextBalls}(i);$ 
5   if  $b_{i,c}$  is empty then
6      $b_{i,c} \leftarrow \frac{x_i}{k v_i};$ 

```

---



---

**Algorithm 2:** Pseudo code of `GetNextBalls(i)`.

---

```

/*  $x_i$  is the timestamp of the last arrived ball
   and  $c$  is the index of assigned bin. */
Input :  $i$ 
Output:  $x_i, c$ 
/*  $z_i$  is the accumulated number of arrived balls
   and we use it to guarantee the consistency. */
1  $seed \leftarrow i || z_i;$ 
2  $u \leftarrow \text{RandUNI}();$ 
3 if  $m_i > \phi_k$  then
4   /* RandInt(k) returns a number from  $\{1, 2, \dots, k\}$ 
     at random. */
5    $j \leftarrow \text{RandInt}(k);$ 
6    $x \leftarrow -\ln u;$ 
7    $z \leftarrow 1;$ 
8 else
9    $z \leftarrow \left\lfloor \frac{\ln u}{\ln(1-m_i/k)} \right\rfloor + 1;$ 
10  /* RandGamma(z, 1) returns a random variable that
     is gamma-distributed with shape  $z$  and scale
     1. */
11   $x \leftarrow \text{RandGamma}(z, 1);$ 
12   $j \leftarrow \text{RandInt}(m_i);$ 
13 if  $j < m_i$  then
14    $\text{Swap}(\pi_{i,j}, \pi_{i,m_i});$ 
15    $m_i \leftarrow m_i - 1;$ 
16  $x_i \leftarrow x_i + x;$ 
17  $z_i \leftarrow z_i + z;$ 
18  $c \leftarrow \pi_{i,m_i};$ 

```

---

$m_i < k$ , that is,

$$P(z = l) = \left(1 - \frac{m_i}{k}\right)^{l-1} \frac{m_i}{k}, \quad l = 1, 2, \dots$$

The results of these  $z$  iterations is equivalent to the following procedure: Generate  $z$  independent random variables  $u_1, \dots, u_z$  according to the uniform distribution  $\text{UNI}(0, 1)$ , set  $x_i = x_i - \sum_{l=1}^z \ln u_l$  at

Step 2, and randomly select an empty bin and set  $b_{i,j} = \min\{b_{i,j}, x_i\}$  at Step 3.

Furthermore, because  $-\ln u_1, \dots, -\ln u_z$  are  $z$  independent and identically distributed exponential random variables with mean 1, the sum of these  $z$  random variables, i.e.,  $-\sum_{l=1}^z \ln u_l$ , is exactly a random variable that is distributed according to the Gamma distribution  $\text{Gamma}(z, 1)$  with shape  $z$  and scale 1. Therefore, we can directly generate a random variable  $x$  according to  $\text{Gamma}(z, 1)$  and  $x$  has the same probability distribution as  $-\sum_{l=1}^z \ln u_l$ . This can significantly reduce the computational cost of generating the random variable  $-\sum_{l=1}^z \ln u_l$  when  $z$  is large.

Then, we derive the model BBM-Permutation, which outputs  $b_{i,1}, \dots, b_{i,k}$  with the same statistical distribution as BBM-Hash. We use a vector  $\vec{\pi}_i = (\pi_{i,1}, \dots, \pi_{i,k})$  to record the index of empty and nonempty bins, which is initialized to  $\vec{\pi}_i = (1, \dots, k)$ . Specially, the first  $m_i$  elements  $\pi_{i,1}, \dots, \pi_{i,m_i}$  are used to keep track of all remaining empty bins at the current time and the rest  $k - m_i$  elements  $\pi_{i,m_i+1}, \dots, \pi_{i,k}$  are all current nonempty bins. In addition, we initialize  $m_i = k$  and  $x_i = 0$ . Then, we repeat the following Steps 1 to 5 until no bin is empty.

- Step 1:** Generate a random variable  $u$  according to the uniform distribution  $\text{UNI}(0, 1)$ ;
- Step 2:** Set  $z = 1$  when  $m_i = k$ . Otherwise, we compute a variable  $z = \left\lfloor \frac{\ln u}{\ln(1-m_i/k)} \right\rfloor + 1$ , which is a geometric random variable with success probability  $\frac{m_i}{k}$ . Here we generate  $z$  using the knowledge that  $\left\lfloor \frac{\ln u}{\ln(1-p)} \right\rfloor + 1$  is a geometric random variable with success probability  $p$  [56];
- Step 3:** Generate a random variable  $x$  according to Gamma distribution  $\text{Gamma}(z, 1)$ ;
- Step 4:** Compute  $x_i = x_i + \frac{x}{k v_i}$ ;
- Step 5:** Select a number  $j$  from  $\{1, \dots, m_i\}$  at random, and then put  $z$  balls into bin  $\pi_{i,j}$  as well as set  $m_i = m_i - 1$  and  $b_{\pi_{i,j}} = x_i$ . In addition, we swap the values of  $\pi_{i,j}$  and  $\pi_{i,m_i}$ .

Unlike BBM-Hash, more than one balls may occur at the same time and all these balls are put into the same bin selected from the current empty bins at random. BBM-Permutation requires exactly  $k$  iterations to fill all bins. In the end,  $\vec{\pi}_i$  is exactly a random permutation of integers  $1, \dots, k$ . We design Step 5 inspired by the Fisher-Yates shuffle [58], which is a popular algorithm used for generating a random permutation of a finite sequence. Compared with BBM-Hash using  $64k$  bits to store  $b_{i,1}, \dots, b_{i,k}$ , BBM-Permutation also requires  $k \log k$  bits to store  $\pi_{i,1}, \dots, \pi_{i,k}$ . Next, we elaborate our model BBM-Mix, which can be used to further accelerate the speed of BBM-Permutation.

**4.2.4 BBM-Mix.** Let  $t_H$  denote the average computational cost required for one iteration of BBM-Hash. For BBM-Hash, when there exist  $m_i$  empty bins, we easily find that on average  $\frac{k}{m_i}$  iterations are required to encounter the next empty bin. Thus, the average computational cost of looking for and filling the next empty bin is  $\frac{k t_H}{m_i}$  and the cost increases as the number of current empty bins  $m_i$  decreases. We also let  $t_P$  denote the average computational cost required for one iteration of BBM-Permutation. For each iteration, BBM-Hash

requires to generate two uniform random variables, while BBM-Permutation requires to generate two uniform random variables and one Gamma random variable. Therefore, BBM-Permutation requires more computations than BBM-Hash to complete an iteration, i.e.,  $t_p$  is larger than  $t_H$ . However, BBM-Permutation requires only one iteration to fill an empty bin selected at random. As a result, we use a variable

$$\phi_k = \frac{kt_H}{t_p},$$

to determine whether BBM-Hash is faster than BBM-Permutation to fill an empty bin for a specific  $m_i$ . Specifically, we use BBM-Hash when  $m_i > \phi_k$  and otherwise BBM-Permutation. In our experiments, we find that  $\frac{t_H}{t_p} \approx 10$ , so we set  $\phi_k = \frac{k}{10}$ .

We can further accelerate the model by simplifying some operations. Specially, we first reduce the division operations at BBM-Hash's Step 2 and BBM-Permutation's Step 4 by computing  $x_i = x_i - \ln u$  and  $x_i = x_i + x$  instead of  $x_i = x_i - \frac{\ln u}{kv_i}$  and  $x_i = x_i + \frac{x}{kv_i}$  respectively, and at the end we enlarge  $b_{i,1}, \dots, b_{i,k}$  by  $\frac{1}{kv_i}$  times. Based on the above improvements, we derive our model BBM-Mix and the Pseudo code of BBM-Mix is given in Algorithm 1. Initialize  $\vec{\pi}_i = (\pi_{i,1}, \dots, \pi_{i,k}) = (1, \dots, k)$ ,  $m_i = k$ ,  $x_i = 0$ , and  $\phi_k = \frac{t_H k}{t_p}$ . We run the following procedure until no bin is empty.

**Step 1:** Generate a random variable  $u$  according to the uniform distribution  $\text{UNI}(0, 1)$ .

**Step 2:** If  $m_i > \phi_k$ , go to Step 3. Otherwise, go to Step 5;

**Step 3:** Compute  $x_i = x_i - \ln u$ ;

**Step 4:** Select a number  $j$  from set  $\{1, \dots, k\}$  at random. If  $j < m_i$  (i.e., now bin  $\pi_{i,j}$  is empty), we put a ball into bin  $\pi_{i,j}$ , and set  $b_{i,\pi_{i,j}} = x_i$ . In addition, we also swap the values of  $\pi_{i,j}$  and  $\pi_{i,m_i}$ , and then set  $m_i = m_i - 1$ . After this step, go to Step 1;

**Step 5:** Compute  $z = \left\lfloor \frac{\ln u}{\ln(1-m_i/k)} \right\rfloor + 1$ , which is a geometric random variable with success probability  $\frac{m_i}{k}$ ;

**Step 6:** Generate a random variable  $x$  according to the Gamma distribution  $\text{Gamma}(z, 1)$ ;

**Step 7:** Compute  $x_i = x_i + x$ ;

**Step 8:** Select a number  $j$  from  $\{1, \dots, m_i\}$  at random, and then put  $z$  balls into bin  $\pi_{i,j}$  as well as set  $m_i = m_i - 1$  and  $b_{\pi_{i,j}} = x_i$ . In addition, swap the values of  $\pi_{i,j}$  and  $\pi_{i,m_i}$ . When finishing this step, go to Step 1.

We easily find that BBM-Mix is faster than both BBM-Hash and BBM-Permutation while the generated variables  $b_{i,1}, \dots, b_{i,k}$  also have the same statistical distribution as those of BBM-Hash and BBM-Permutation. Last, we would like to point out that BBM-Mix has the same space complexity as BBM-Permutation.

### 4.3 Our Method FastGM

Based on the model BBM-Mix, generating a Gumbel-Max sketch of a vector can be equivalently viewed as the procedure of throwing  $kn_{\vec{v}}^+$  balls generated by  $n_{\vec{v}}^+$  Poisson processes into  $k$  bins, where each ball is assigned with a random variable in ascending order. Before introducing our method FastGM in detail, we naturally have the following two fundamental questions for the design of FastGM:

**Question 1.** How to fast search balls with the smallest timestamps to fill bins from these  $n_{\vec{v}}^+$  Poisson processes?

---

**Algorithm 3:** Pseudo code of our FastGM, where function  $\text{GetNextBalls}(i)$  used in lines 10 and 24 is defined in Algorithm 2.

---

**Input :**  $\vec{v} = (v_1, \dots, v_n)$ ,  $k$ ,  $\phi_k$   
**Output :**  $\vec{s} = (s_1, \dots, s_k)$

```

1  $R \leftarrow 0$ ;  $k^* \leftarrow k$ ;  $(y_1, \dots, y_k) \leftarrow (-1, \dots, -1)$ ;
2 foreach  $i \in N_{\vec{v}}^+$  do
3    $x_i \leftarrow 0$ ;  $z_i \leftarrow 0$ ;  $m_i \leftarrow k$ ;
4    $(\pi_{i,1}, \dots, \pi_{i,k}) \leftarrow (1, \dots, k)$ ;
   /* The following part is LinearFill */
5 while  $k^* \neq 0$  do
6    $R \leftarrow R + \Delta$ ;
7   foreach  $i \in N_{\vec{v}}^+$  do
8      $R_i \leftarrow \lceil Rv_i^* \rceil$ ;
9     while  $m_i > 0$  and  $z_i < R_i$  do
10       $x_i, c \leftarrow \text{GetNextBalls}(i)$ ;
11       $b_i \leftarrow \frac{x_i}{kv_i}$ ;
12      if  $y_c < 0$  then
13         $y_c \leftarrow b_i$ ;  $s_c \leftarrow i$ ;
14         $k^* \leftarrow k^* - 1$ ;
15      else if  $b_i < y_c$  then
16         $y_c \leftarrow b_i$ ;  $s_c \leftarrow i$ ;
   /* The following part is FastPrune */
17  $j^* \leftarrow \arg \max_{j=1, \dots, k} y_j$ ;
18  $N \leftarrow N_{\vec{v}}^+$ ;
19 while  $N$  is not empty do
20    $R \leftarrow R + \Delta$ ;
21   foreach  $i \in N$  do
22      $R_i \leftarrow \lceil Rv_i^* \rceil$ ;
23     while  $m_i > 0$  and  $z_i < R_i$  do
24       $x_i, c \leftarrow \text{GetNextBalls}(i)$ ;
25       $b_i \leftarrow \frac{x_i}{kv_i}$ ;
26      if  $b_i > y_{j^*}$  then
27         $N \leftarrow N - \{i\}$ ;
28        break;
29      if  $b_i < y_c$  then
30         $y_c \leftarrow b_i$ ;  $s_c \leftarrow i$ ;
31        if  $c == j^*$  then
32           $j^* \leftarrow \arg \max_{j=1, \dots, k} y_j$ ;

```

---

**Question 2.** How to early stop a Poisson process  $\mathcal{P}_i$ ,  $i \in N_{\vec{v}}^+$ ?

We first discuss Question 1. We note that balls of different Poisson processes  $\mathcal{P}_i$  arrive at different rates  $kv_i$ . Recall the example in Figure 1, the basic idea behind the following technique is that process  $\mathcal{P}_i$  with high rate  $kv_i$  is more likely to produce balls with the smallest timestamps (i.e. Gumbel-Max variables). Specially, when  $z$  balls have been generated, let  $x_{i,z}$  denote the time of the latest ball occurred in our method BBM-Mix. We easily find that  $x_{i,z}$  can be represented as the sum of  $z$  identically distributed exponential

random variables with mean  $\frac{1}{kv_i}$ . Therefore, the expectation and variance of variable  $x_{i,z}$  are computed as

$$\mathbb{E}(x_{i,z}) = \frac{z}{kv_i}, \quad \text{Var}(x_{i,z}) = \frac{z}{k^2v_i^2}. \quad (2)$$

We find that  $\mathbb{E}(x_{i,z})$  is  $l$  times smaller than  $\mathbb{E}(x_{j,z})$  when  $v_i$  is  $l$  times larger than  $v_j$ . To obtain the first  $R$  balls of the joint of all Poisson processes  $\mathcal{P}_i$ ,  $i \in N_{\vec{v}}^+$ , we let each process  $\mathcal{P}_i$  generate  $R_i = \lceil Rv_i^* \rceil$  balls. Then, we have  $\sum_{i=1}^n R_i \approx R$ . For all  $i \in N_{\vec{v}}^+$ , their  $x_{i,r_i}$  approximately have the same expectation.

$$\mathbb{E}(x_{i,r_i} | R) \approx \frac{R}{k \sum_{j=1}^n v_j}, \quad i \in N_{\vec{v}}^+. \quad (3)$$

Therefore, the  $R$  balls with the smallest timestamps are expected to be obtained.

Then, we discuss Question 2, which is inspired by the ascending-order random variables for each arrived ball. We let each bin  $j$  use two registers  $y_j$  and  $s_j$  to keep track of information of the ball with the smallest timestamp among its currently received balls, where  $y_j$  records the ball's timestamp and  $s_j$  records the ball's origin, i.e., the ball is coming from Poisson process  $\mathcal{P}_{s_j}$ . When all bins  $1, \dots, k$  are filled with at least one ball, we let  $y^*$  keep track of the maximum value of  $y_1, \dots, y_k$ , i.e.,

$$y^* = \max_{j=1, \dots, k} y_j.$$

Then, we can stop Poisson process  $\mathcal{P}_i$  when a ball getting from  $\mathcal{P}_i$  has a timestamp the larger than  $y^*$  because the timestamps of the subsequent balls from  $\mathcal{P}_i$  are also larger than  $y^*$ , which will not change any  $y_1, \dots, y_k$  and  $s_1, \dots, s_k$ .

Based on the above two discussions, we develop our final method FastGM to fast generate  $k$  Gumbel-Max variables  $\vec{s}_{\vec{v}} = (s_1, \dots, s_k)$  of any nonnegative vector  $\vec{v}$ . As shown in Algorithm 3, FastGM consists of two modules: LinearFill and FastPrune. LinearFill is designed to quickly search balls with smallest timestamps arriving from all processes  $\mathcal{P}_1, \dots, \mathcal{P}_n$  and fill all bins  $1, \dots, k$ . When no bin is empty, we start the FastPrune module to early stop each Poisson process  $\mathcal{P}_i$ ,  $i \in N_{\vec{v}}^+$ . We perform the procedure of FastPrune because Poisson Process  $\mathcal{P}_i$  may also get balls with timestamps smaller than  $y^*$  and the balls may change the values of  $y_j$  and  $s_j$  for some bins  $j$  after the procedure of LinearFill. Next, we introduce these two modules in detail.

• **LinearFill Module:** This module fast search balls with smallest timestamps, and consists of the following steps:

**Step 1:** Iterate on each  $i \in N_{\vec{v}}^+$  and repeat function GetNextBall( $i$ ) in Algorithm 2 until it has received not less than  $\lceil Rv_i^* \rceil$  balls since the beginning of the algorithm. Meanwhile, each bin  $j$  uses registers  $y_j$  and  $s_j$  to keep track of information of its received ball having the smallest timestamp, where  $y_j$  records the ball's timestamp and  $s_j$  records the index of the Poisson process where the ball comes from;

**Step 2:** If there exist any empty bins, we increase  $R$  by  $\Delta$  and then repeat Step 1. Otherwise, we stop the LinearFill procedure.

For simplicity, we set the parameter  $\Delta = k$ . In our experiments, we find that the value of  $\Delta$  has a small effect on the performance of FastGM.

• **FastPrune Module:** When all bins  $1, \dots, k$  have been filled by at least one ball. We start the FastPrune module, which mainly consists of the following two steps:

**Step 1.** Compute  $y^* = \max_{j=1, \dots, k} y_j$ .

**Step 2.** For each  $i \in N_{\vec{v}}^+$ , we repeat function GetNextBall( $i$ ) in Algorithm 2 to generate balls. When the ball has a timestamp larger than  $y^*$ , we terminate Poisson process  $\mathcal{P}_i$ . Note that  $y_1, \dots, y_k$  and  $s_1, \dots, s_k$  are also updated by received balls at this step. Therefore,  $y^*$  may also decrease with the number of arriving balls, which accelerates the speed of terminating all Poisson processes  $\mathcal{P}_i$ ,  $i \in N_{\vec{v}}^+$ .

## 4.4 Complexity

**Space Complexity.** For a nonnegative vector  $\vec{v}$  with  $n_{\vec{v}}^+$  positive elements, our method FastGM requires  $k \log k$  bits to store  $\vec{\pi}_i$  for each  $i \in N_{\vec{v}}^+$ , and in summary,  $n_{\vec{v}}^+ k \log k$  bits are desired. In addition,  $64k$  bits are desired for storing  $y_1, \dots, y_k$  (we use 64-bit floating-point registers to record  $y_1, \dots, y_k$ ), and  $k \log n$  bits are required for storing  $s_1, \dots, s_k$ , where  $n$  is the size of the vector. However, the additional memory is released immediately after computing the sketch, and is far smaller than the memory for storing the generated sketches of massive vectors (e.g. documents). Therefore, FastGM requires  $n_{\vec{v}}^+ k \log k + 64k + k \log n$  bits when generating  $k$  Gumbel-Max variables  $\vec{s}(\vec{v}) = (s_1, \dots, s_k)$  of  $\vec{v}$ .

**Time Complexity.** We easily find that a nonnegative vector and its normalized vector have the same Gumbel-Max sketch. For simplicity, therefore we analyze the time complexity of our method for only normalized vectors. Let  $\vec{v}^* = (v_1^*, \dots, v_n^*)$  be a normalized and nonnegative vector. Define

$$\tilde{y}_j = \min_{i \in N_{\vec{v}^*}^+} -\frac{\ln a_{i,j}}{v_i^*}, \quad j = 1, \dots, k,$$

$$\tilde{y}^* = \max_{j=1, \dots, k} \tilde{y}_j.$$

At the end of our FastPrune procedure, each register  $y_j$  used in the procedure equals  $\tilde{y}_j$  and register  $y^*$  equals  $\tilde{y}^*$ . Because  $\frac{\ln a_{i,j}}{v_i^*} \sim \text{EXP}(v_i^*)$ , we easily find that each  $y_j$  follows the exponential distribution  $\text{EXP}(\sum_{i=1}^n v_i^*) = \text{EXP}(1)$ . From [59], we then easily have

$$\mathbb{E}(\tilde{y}^*) = \sum_{m=1}^k \frac{1}{m} = \ln k + \gamma, \quad \text{Var}(\tilde{y}^*) = \sum_{m=1}^k \frac{1}{m^2} < \sum_{m=1}^{\infty} \frac{1}{m^2} = \frac{\pi^2}{6},$$

where  $\gamma \approx 0.577$  is the Euler-Mascheroni constant. From Chebyshev's inequality, we have

$$P(|\tilde{y}^* - \mathbb{E}(\tilde{y}^*)| \geq \alpha \text{Var}(\tilde{y}^*)) \leq \frac{1}{\alpha^2}.$$

Therefore,  $\tilde{y}^* \leq \mathbb{E}(\tilde{y}^*) + \alpha \text{Var}(\tilde{y}^*)$  happens with a large probability when  $\alpha$  is large. In other words, the random variable  $\tilde{y}^*$  can be upper bounded by  $\mathbb{E}(\tilde{y}^*) + \alpha \text{Var}(\tilde{y}^*)$  with a large probability. Next, we derive the expectation of  $x_{i,R}$  after the first  $R$  balls generated by our method FastGM have been put into  $k$  bins. For each Poisson process  $\mathcal{P}_i$ ,  $i \in N_{\vec{v}}^+$ , from equations (2) and (3), we find that its last produced ball among these first  $R$  balls has a timestamp  $x_{i,R_i}$  with the expectation  $\mathbb{E}(x_{i,R_i} | R) \approx \frac{R}{k}$ . When  $R = k(\mathbb{E}(\tilde{y}^*) + \alpha \text{Var}(\tilde{y}^*)) < k(\ln k + \gamma + \frac{\alpha \pi^2}{6})$ , the probability of  $\mathbb{E}(x_{i,R_i}) > \tilde{y}^*$  is almost 1

for large  $\alpha$ , e.g.,  $\alpha > 10$ . Therefore, we find that after putting the first  $O(k \ln k)$  balls into the  $k$  bins, each Poisson process  $\mathcal{P}_i$  is expected to be early terminated and so we are likely to acquire all the Gumbel-Max variables. We also note that each positive element has to be enumerated in the FastPrune model, therefore, the total time complexity of our method FastGM is  $O(k \ln k + n \frac{1}{\alpha})$ .

## 5 EVALUATION

We evaluate our method FastGM with the state-of-the-arts on two tasks: **task 1**) probability Jaccard similarity estimation, and **task 2**) network embedding. All algorithms run on a computer with a Quad-Core Intel(R) Xeon(R) CPU E3-1226 v3 CPU 3.30GHz processor. To demonstrate the reproducibility of the experimental results, we make our source code publicly available<sup>1</sup>.

### 5.1 Dataset

We conduct experiments on both real-world and synthetic datasets. For task 1, we run experiments on six real-world datasets: Real-sim [60], Rcv1 [61], Webspam [62], Libimseti [63], Last.fm [64], and MovieLens [65]. In detail, Real-sim [60], Rcv1 [61], and Webspam [62] are datasets of web documents from different resources, where each vector represents a document and each entry in the vector refers to the TF-IDF score of a specific word for the document. Libimseti [63] is a dataset of ratings between users on the Czech dating site, where each vector refers to a user and each entry records the user's rating to another one. Last.fm [64] is a dataset of listening history, where each vector represents a song and each entry in the vector is the number of times the song has been listened to by a specific user. MovieLens [65] is a dataset of movie ratings, where each vector is a user and each entry in the vector is that user's rating to a specific movie.

For task 2, we perform experiments on four real-world graphs: YouTube [66], Email-EU [67], Twitter [68], and WikiTalk [69]. In detail, YouTube [66] is a dataset of friendships between YouTube users. Email-EU [67] is an email communication network, where nodes represent individual persons and edges are the communications between persons. Twitter [68] is the following network collected from twitter.com and WikiTalk [69] is the communication network of the Chinese Wikipedia. The statistics of all the above datasets are summarized in Table 1.

### 5.2 Baseline

For task 1, we compare our method with  $\mathcal{P}$ -MinHash [9] on probability Jaccard similarity estimation to evaluate the performance of FastGM. To highlight the efficiency of FastGM, we further compare FastGM with the state-of-the-art weighted Jaccard similarity estimation method, BagMinHash [37]. Notice that BagMinHash estimates a different metric and thus we only show results on efficiency. For task 2, we compare our method with the state-of-the-art embedding algorithm NodeSketch [13]. Specially, we use FastGM to replace the module of generating Gumbel-Max sketches (i.e., node embeddings) in NodeSketch. In our experiments, we vary the size of node embedding  $k$  and set decay weight  $a = 0.005$  and order of proximity  $r = 5$  as suggested in the original paper of NodeSketch.

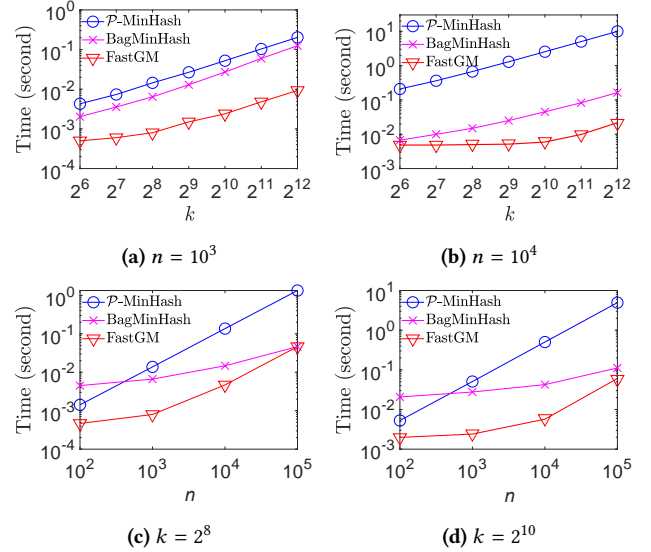
<sup>1</sup><https://github.com/qyy0180/FastGM>

**Table 1: Statistics of all used datasets.**

(a) Datasets used in task 1		
Dataset	#Vectors	#Features
Real-sim [60]	72,309	20,958
Rcv1 [61]	20,242	47,236
Webspam [62]	350,000	16,609,143
Libimseti [63]	220,970	220,970
Last.fm [64]	992	1,085,612
MovieLens [65]	69,878	80,555

(b) Datasets used in task 2		
Dataset	#Nodes	#Edges
YouTube [66]	1,138,499	2,990,443
Email-EU [67]	265,214	420,045
Twitter [68]	465,017	834,797
WikiTalk [69]	1,219,241	2,284,546



**Figure 2: (Task 1) Efficiency of FastGM compared with  $\mathcal{P}$ -MinHash and BagMinHash on synthetic vectors, where each element in the vector is randomly selected from UNI(0,1).**

### 5.3 Metric

For task 1, we use the sketching time and root mean square error (RMSE) to measure the efficiency and the effectiveness respectively. For task 2, besides the efficiency comparison with the original NodeSketch, we also evaluate the effectiveness of our method on two popular applications of graph embedding: node classification and link prediction. Similar to [13], we use the Macro and Micro F1 scores to measure the performance of node classification, and Precision@K and Recall@K to evaluate the performance of link prediction, where Precision@K (resp. Recall@K) is the precision (resp. recall) on top K high-similarity testing node pairs. All experimental results are empirically computed from 100 independent runs by default.



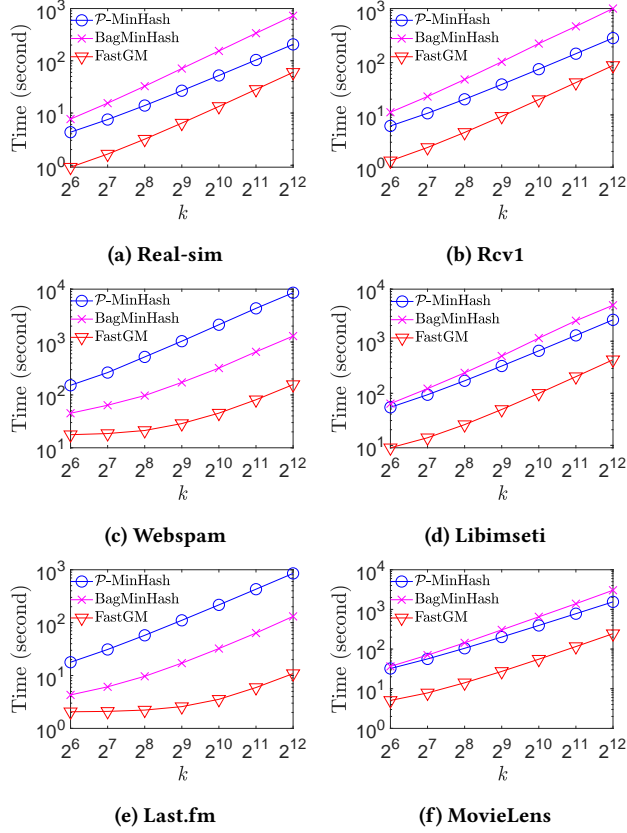


Figure 3: (Task 1) Efficiency of FastGM compared with  $\mathcal{P}$ -MinHash for different  $k$ .

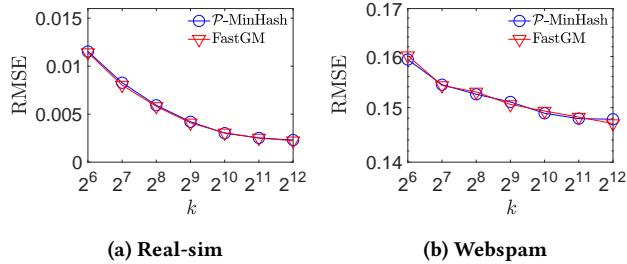


Figure 4: (Task 1) Accuracy of FastGM compared with  $\mathcal{P}$ -MinHash for different  $k$ .

#### 5.4 Probability Jaccard Similarity Estimation

We conduct experiments on both synthetic and real-world datasets for task 1. Specially, we first use synthetic weighted vectors to evaluate the performance of FastGM for vectors with different dimensions. Then, we show results on a variety of real-world datasets.

**Results on synthetic vectors.** In this experiment, we also compare our method with another state-of-the-art algorithm BagMinHash [37], which is used for estimating weighted Jaccard similarity. We note that BagMinHash estimates an alternative similarity metric. However, a lot of experiments and theoretical analysis [9] have

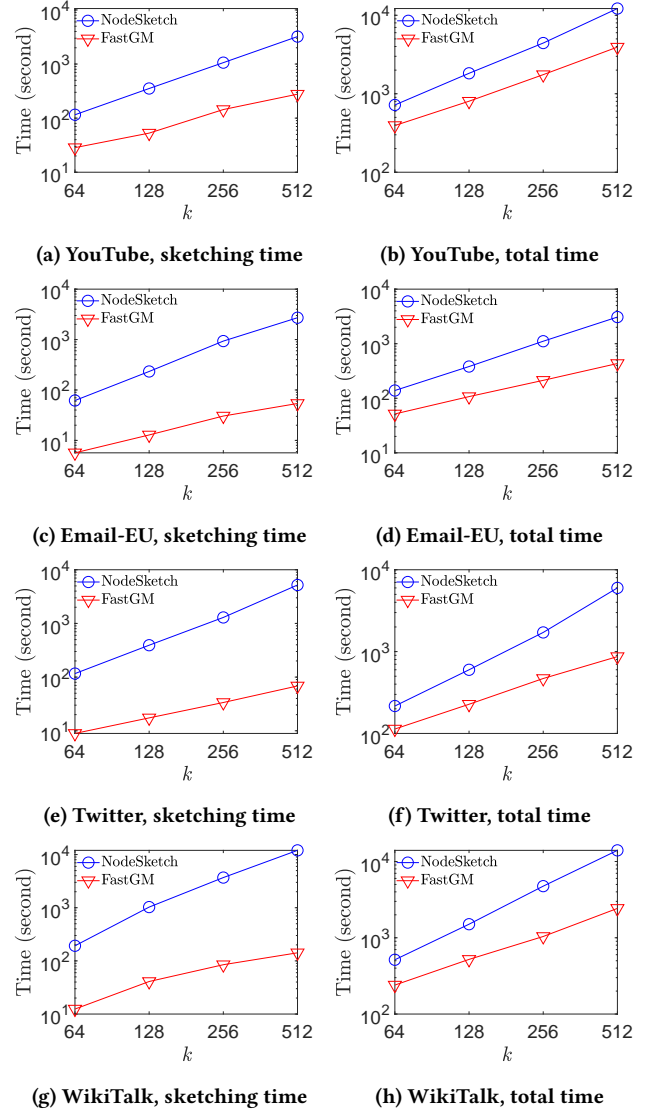
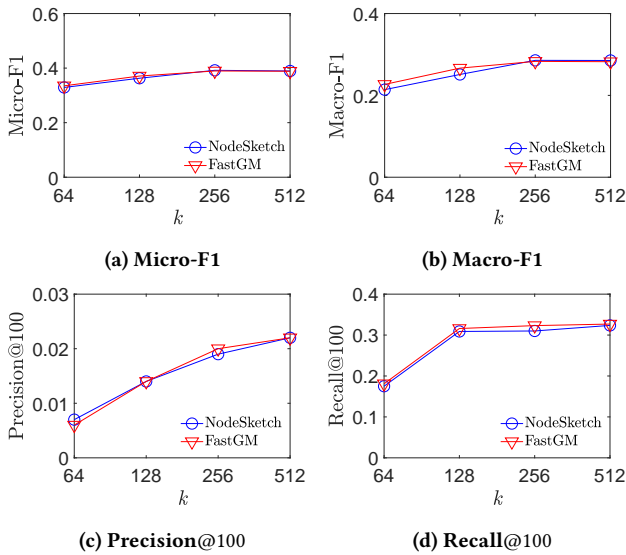


Figure 5: (Task 2) Efficiency of our method FastGM in comparison with NodeSketch for different  $k$ .

shown that weighted Jaccard similarity and probability Jaccard similarity usually have similar performance on many applications such as fast searching similar set. We conduct experiments on weighted vectors with uniform-distribution weights. Without loss of generality, we let  $n_{\vec{v}}^+ = n$  for each vector, i.e., all elements of each vector are positive. As shown in Figures 2 (a) and (b), when  $n = 10^3$ , FastGM is 13 and 22 times faster than BagMinHash and  $\mathcal{P}$ -MinHash respectively. As  $n$  increases to  $10^4$ , the improvement becomes 8 and 125 times respectively. Especially, the sketching time of our method is around 0.02 seconds when  $n = 10^4$  and  $k = 2^{12}$ , while BagMinHash and  $\mathcal{P}$ -MinHash take over 0.15 and 2.5 seconds for sketching respectively. Figures 2 (c) and (d) show the running time of all competitors for different  $n$ . Our method FastGM is 3 to 100 times faster than  $\mathcal{P}$ -MinHash for different  $n$ . Compared with



**Figure 6: (Task 2, YouTube) Accuracy of FastGM and NodeSketch. (a)(b): node classification; (c)(d): link prediction.**

BagMinHash, FastGM is about 10 times faster when  $n = 1,000$ , and is comparable as  $n$  increases to 100,000. It indicates that our method FastGM significantly outperforms BagMinHash for vectors having less than 100,000 positive elements, which are prevalent in real-world datasets. In addition, we also conduct experiments on weighted vectors with exponential-distribution weights and omit similar results here.

**Results on real-world datasets.** Next, we show results on the real-world datasets in Table 1. Figure 3 exhibits the sketching time of all algorithms. We see that our method outperforms  $\mathcal{P}$ -MinHash and BagMinHash on all the datasets and the improvement increases as  $k$  increases. On sparse datasets such as Real-sim, Rcv1, and MovieLens, FastGM is about 8 and 12 faster than  $\mathcal{P}$ -MinHash and BagMinHash respectively. BagMinHash is even slower than  $\mathcal{P}$ -MinHash on these datasets. On datasets Webspam and Last.fm, we note that FastGM is 55 and 80 times faster than  $\mathcal{P}$ -MinHash respectively. Figure 4 shows the estimation error of FastGM and  $\mathcal{P}$ -MinHash on datasets Real-sim and Webspam. Due to the large number of vector pairs, we here randomly select 100,000 pairs of vectors from each dataset and report the average RMSE. We note that both algorithms give similar accuracy, which is coincident with our analysis. We omit similar results on other datasets.

## 5.5 Graph Embedding

We compare FastGM with the regular NodeSketch to demonstrate the efficiency of our method on the graph embedding task. Specially, we show the sketching time and the total time for different  $k$  (i.e., the size of node embeddings), where the sketching time refers to the accumulated time of computing Gumbel-Max sketches from different orders of the self-loop-augmented adjacent matrix. As shown in Figure 5, we note that our method FastGM gives significant improvement on all datasets. In detail, on dataset WikiTalk, FastGM gives an improvement of 16 times at  $k = 2^6$  and the improvement

increases to 84 times at  $k = 2^9$  on sketching time, which results in a gain of up to 5 times improvements for the total time of learning node embeddings.

We also conduct experiments on two popular applications of graph embedding, i.e., node classification and link prediction. All experimental settings are the same as [13]. For node classification, we randomly select 10% nodes as the training set and others as the testing set. Then, we build a one-vs-rest SVM classifier based on the training set. For link prediction, we randomly drop out 20% edges from the original graph as the testing set and learn the embeddings based on the remaining graph. We predict the potential edges by generating a ranked list of node pairs. For each node pair, we use the Hamming similarity of their embeddings to generate the ranked list. Due to the massive size of node pairs, we randomly sample  $10^5$  pairs of nodes for evaluation. We report Precision@100 and Recall@100. Figure 6 shows the results on node classification as well as link prediction on dataset YouTube. We notice that our method FastGM gives similar accuracy compared with NodeSketch. We omit the results of the other three datasets.

## 6 CONCLUSIONS AND FUTURE WORK

In this paper, we develop a novel algorithm FastGM to fast compute a nonnegative vector’s Gumbel-Max sketch, which consists of  $k$  independent Gumbel-Max variables. We prove that FastGM generates Gumbel-Max sketch with the same quality as the traditional Gumbel-Max Trick but reduces time complexity from  $O(n^+k)$  to  $O(k \ln k + n^+)$ , where  $n^+$  is the number of the vector’s positive elements. We conduct a variety of experiments on two tasks: Probability Jaccard similarity estimation and graph embedding, and the experimental results demonstrate that our method FastGM is orders of magnitude faster than the state-of-the-arts without sacrificing accuracy. In the future, we plan to extend FastGM to vectors consisting of elements arriving in a streaming fashion.

## ACKNOWLEDGMENT

The research presented in this paper is supported in part by National Key R&D Program of China (2018YFC0830500), Shenzhen Basic Research Grant (JCYJ20170816100819428), National Natural Science Foundation of China (61922067, U1736205, 61902305), MoE-CMCC “Artificial Intelligence” Project (MCM20190701), Natural Science Basic Research Plan in Shaanxi Province of China (2019JM-159), Natural Science Basic Research Plan in Zhejiang Province of China (LGG18F020016).

## REFERENCES

- [1] R Duncan Luce. *Individual choice behavior: A theoretical analysis*. Courier Corporation, 1959.
- [2] Monika Henzinger. Finding near-duplicate web pages: a large-scale evaluation of algorithms. In *SIGIR*, pages 284–291. ACM, 2006.
- [3] Gurmeet Singh Manku, Arvind Jain, and Anish Das Sarma. Detecting near-duplicates for web crawling. In *WWW*, pages 141–150. ACM, 2007.
- [4] Yoram Bachrach, Ely Porat, and Jeffrey S Rosenschein. Sketching techniques for collaborative filtering. In *IJCAI*, 2009.
- [5] Michael Mitzenmacher, Rasmus Pagh, and Ninh Pham. Efficient estimation for high similarities using odd sketches. In *WWW*, pages 109–118, 2014.
- [6] Dingqi Yang, Bin Li, and Philippe Cudré-Mauroux. Poisketch: Semantic place labeling over user activity streams. Technical report, Université de Fribourg, 2016.
- [7] Dingqi Yang, Bin Li, Laura Rettig, and Philippe Cudré-Mauroux. Histosketch: Fast similarity-preserving sketching of streaming histograms with concept drift.

- In *IEEE ICDM*, pages 545–554. IEEE, 2017.
- [8] Dingqi Yang, Bin Li, Laura Rettig, and Philippe Cudré-Mauroux. D2 histosketch: discriminative and dynamic similarity-preserving sketching of streaming histograms. *IEEE TKDE*, pages 1–1, 2018.
  - [9] Ryan Moulton and Yunjiang Jiang. Maximally consistent sampling and the jaccard index of probability distributions. *arXiv preprint arXiv:1809.04052*, 2018.
  - [10] Aristides Gionis, Piotr Indyk, and Rajeev Motwani. Similarity search in high dimensions via hashing. In *PVLDB*, pages 518–529, 1999.
  - [11] Andrei Z Broder, Moses Charikar, Alan M Frieze, and Michael Mitzenmacher. Min-wise independent permutations. *J. Comput. Syst. Sci.*, 60(3):630–659, June 2000.
  - [12] Moses S Charikar. Similarity estimation techniques from rounding algorithms. In *STOC*, pages 380–388, 2002.
  - [13] Dingqi Yang, Paolo Rosso, Bin Li, and Philippe Cudre-Mauroux. Nodsketch: Highly-efficient graph embeddings via recursive sketching. In *SIGKDD*, 2019.
  - [14] Guy Liorbom, Andreea Gane, Tommi Jaakkola, and Tamir Hazan. Direct optimization through argmax for discrete variational auto-encoder. *arXiv preprint arXiv:1806.02867*, 2018.
  - [15] Eliav Buchnik, Edith Cohen, Avinatan Hasidim, and Yossi Matias. Self-similar epochs: Value in arrangement. In *ICML*, pages 841–850, 2019.
  - [16] Michael Oberst and David Sonntag. Counterfactual off-policy evaluation with gumbel-max structural causal models. *arXiv preprint arXiv:1905.05824*, 2019.
  - [17] Carolyn Kim, Ashish Sabharwal, and Stefano Ermon. Exact sampling with integer linear programs and random perturbations. In *AAAI*, 2016.
  - [18] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
  - [19] Matt J Kusner and José Miguel Hernández-Lobato. Gans for sequences of discrete elements with the gumbel-softmax distribution. *arXiv preprint arXiv:1611.04051*, 2016.
  - [20] Yi Tay, Anh Tuan Luu, and Siu Cheung Hui. Multi-pointer co-attention networks for recommendation. In *SIGKDD*, pages 2309–2318. ACM, 2018.
  - [21] Ping Li and Arnd Christian König. b-bit minwise hashing. In *WWW*, pages 671–680, 2010.
  - [22] Pinghui Wang, Yiyan Qi, Yuanming Zhang, Qiaozhu Zhai, Chenxu Wang, John C. S. Lui, and Xiaohong Guan. A memory-efficient sketch method for estimating high similarities in streaming sets. In *SIGKDD*, pages 25–33, 2019.
  - [23] Ping Li, Art B. Owen, and Cun-Hui Zhang. One permutation hashing. In *NIPS*, pages 3122–3130, 2012.
  - [24] Anshumali Shrivastava and Ping Li. Improved densification of one permutation hashing. In *UAI*, pages 732–741, 2014.
  - [25] Anshumali Shrivastava and Ping Li. Densifying one permutation hashing via rotation for fast near neighbor search. In *ICML*, pages 557–565, 2014.
  - [26] Anshumali Shrivastava. Optimal densification for fast and accurate minwise hashing. In *ICML*, pages 3154–3163, 2017.
  - [27] Søren Dahlgaard, Mathias Bæk Tejs Knudsen, and Mikkel Thorup. Fast similarity sketching. In *FOCS*, pages 663–671. IEEE, 2017.
  - [28] Taher Haveliwala, Aristides Gionis, and Piotr Indyk. Scalable techniques for clustering the web. 2000.
  - [29] Bernhard Haeupler, Mark Manasse, and Kunal Talwar. Consistent weighted sampling made fast, small, and easy. *arXiv preprint arXiv:1410.4266*, 2014.
  - [30] Sreenivas Gollapudi and Rina Panigrahy. Exploiting asymmetry in hierarchical topic extraction. In *CIKM*, pages 475–482. ACM, 2006.
  - [31] Mark Manasse, Frank McSherry, and Kunal Talwar. Consistent weighted sampling. Technical report, June 2010.
  - [32] Sergey Ioffe. Improved consistent sampling, weighted minhash and L1 sketching. In *ICDM*, pages 246–255, 2010.
  - [33] Ping Li. 0-bit consistent weighted sampling. In *SIGKDD*, pages 665–674, 2015.
  - [34] Wei Wu, Bin Li, Ling Chen, and Chengqi Zhang. Canonical consistent weighted sampling for real-value weighted min-hash. In *ICDM*, pages 1287–1292, 2016.
  - [35] Wei Wu, Bin Li, Ling Chen, and Chengqi Zhang. Consistent weighted sampling made more practical. In *WWW*, pages 1035–1043, 2017.
  - [36] Wei Wu, Bin Li, Ling Chen, Chengqi Zhang, and Philip Yu. Improved consistent weighted sampling revisited. *IEEE TKDE*, 2018.
  - [37] Otmár Ertl. Bagminhash-minwise hashing algorithm for weighted sets. In *SIGKDD*, pages 1368–1377. ACM, 2018.
  - [38] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social representations. In *SIGKDD*, pages 701–710, 2014.
  - [39] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *NIPS*, pages 3111–3119, 2013.
  - [40] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *SIGKDD*, pages 855–864, 2016.
  - [41] Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. Line: Large-scale information network embedding. In *WWW*, pages 1067–1077, 2015.
  - [42] Shaosheng Cao, Wei Lu, and Qiongkai Xu. Grarep: Learning graph representations with global structural information. In *CIKM*, pages 891–900, 2015.
  - [43] Jiezhong Qiu, Yuxiao Dong, Hao Ma, Jian Li, Kuansan Wang, and Jie Tang. Network embedding as matrix factorization: Unifying deepwalk, line, pte, and node2vec. In *WSDM*, pages 459–467, 2018.
  - [44] Jundong Li, Harsh Dani, Xia Hu, Jiliang Tang, Yi Chang, and Huan Liu. Attributed network embedding for learning in a dynamic environment. In *CIKM*, pages 387–396, 2017.
  - [45] Vachik S Dave, Baichuan Zhang, Pin-Yu Chen, and Mohammad Al Hasan. Neural-brane: Neural bayesian personalized ranking for attributed network embedding. *DSE*, 4(2):119–131, 2019.
  - [46] Stephen Bonner, Ibad Kureshi, John Brennan, Georgios Theodoropoulos, Andrew Stephen McGough, and Boguslaw Obara. Exploring the semantic content of unsupervised graph embeddings: an empirical study. *DSE*, 4(3):269–289, 2019.
  - [47] Lin Lan, Pinghui Wang, Junzhou Zhao, Jing Tao, John CS Lui, and Xiaohong Guan. Improving network embedding with partially available vertex and edge content. *Information Sciences*, 512:935–951, 2020.
  - [48] Cunhao Tu, Weicheng Zhang, Zhiyuan Liu, and Maosong Sun. Max-margin deepwalk: discriminative learning of network representation. In *IJCAI*, pages 3889–3895, 2016.
  - [49] Juzheng Li, Jun Zhu, and Bo Zhang. Discriminative deep random walk for network classification. In *ACL*, pages 1004–1013, 2016.
  - [50] Zhilin Yang, William W Cohen, and Ruslan Salakhutdinov. Revisiting semi-supervised learning with graph embeddings. In *ICML*, pages 40–48, 2016.
  - [51] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv.org*, 2016.
  - [52] Giang Hoang Nguyen, John Boaz Lee, Ryan A Rossi, Nesreen K Ahmed, Eunye Koh, and Sunghul Kim. Continuous-time dynamic network embeddings. In *Companion of the the Web Conference*, pages 969–976, 2018.
  - [53] Wenchao Yu, Wei Cheng, Charu C. Aggarwal, Kai Zhang, Haifeng Chen, and Wei Wang. Netwalk: A flexible deep embedding approach for anomaly detection in dynamic networks. In *SIGKDD*, pages 2672–2681, 2018.
  - [54] Lun Du, Yun Wang, Guojie Song, Zhicong Lu, and Junshan Wang. Dynamic network embedding: An extended approach for skip-gram based network embedding. In *IJCAI*, pages 2086–2092, 2018.
  - [55] Yiyan Qi, Jiefeng Cheng, Xiaojun Chen, Reynold Cheng, Albert Bifet, and Pinghui Wang. Discriminative streaming network embedding. *KBS*, 2019.
  - [56] Michael Mitzenmacher and Eli Upfal. *Probability and Computing: Randomized Algorithms and Probabilistic Analysis*, volume 160. Cambridge University Press, New York, NY, USA, 2005.
  - [57] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to algorithms*, volume 6. MIT Press, Cambridge, MA, 2nd edition, 2001.
  - [58] Frank Fisher, Ronald A. and Yates. *Statistical tables for biological, agricultural and medical research (3rd ed.)*. Oliver & Boyd, London, UK, 1948.
  - [59] Variance of the maximum of n independent exponentials. <https://math.stackexchange.com/questions/3175307/variance-of-the-maximum-of-n-independent-exponentials>.
  - [60] Wu Wei, Bin Li, Chen Ling, and Chengqi Zhang. Consistent weighted sampling made more practical. In *WWW*, pages 1035–1043, 2017.
  - [61] David D. Lewis, Yiming Yang, Tony G. Rose, and Li Fan. Rcv1: A new benchmark collection for text categorization research. *JMLR*, 5(2):361–397, 2004.
  - [62] De Wang, Danesh Irani, and Calton Pu. Evolutionary study of web spam: Webb spam corpus 2011 versus webb spam corpus 2006. In *COLLABORATECOM*, 2012.
  - [63] Libimseti.cz network dataset – KONECT, April 2017.
  - [64] Last.fm song network dataset – KONECT, April 2017.
  - [65] Movielens 10m network dataset – KONECT, April 2017.
  - [66] Lei Tang and Huan Liu. Scalable learning of collective behavior based on sparse social dimensions. pages 1107–1116, 2009.
  - [67] Jure Leskovec, Jon Kleinberg, and Christos Faloutsos. Graph evolution: Densification and shrinking diameters. *TKDD*, 1(1):2, 2007.
  - [68] Munmun De Choudhury, Yu-Ru Lin, Hari Sundaram, K. Selcuk Candan, Lexing Xie, and Aisling Kelliher. How does the data sampling strategy impact the discovery of information diffusion in social media? In *ICWSM*, pages 34–41, 2010.
  - [69] Jun Sun, Jerome Kunegis, and Steffen Staab. Predicting user roles in social networks using transfer learning with feature transformation. In *ICDMW*, pages 128–135. IEEE, 2016.