

# Sublinear Time Orthogonal Tensor Decomposition

**Authored by:**

David Woodruff  
Zhao Song  
Huan Zhang

## Abstract

A recent work (Wang et. al., NIPS 2015) gives the fastest known algorithms for orthogonal tensor decomposition with provable guarantees. Their algorithm is based on computing sketches of the input tensor, which requires reading the entire input. We show in a number of cases one can achieve the same theoretical guarantees in sublinear time, i.e., even without reading most of the input tensor. Instead of using sketches to estimate inner products in tensor decomposition algorithms, we use importance sampling. To achieve sublinear time, we need to know the norms of tensor slices, and we show how to do this in a number of important cases. For symmetric tensors  $T = \sum_{i=1}^k \lambda_i u_i^{\otimes p}$  with  $\lambda_i \geq 0$  for all  $i$ , we estimate such norms in sublinear time whenever  $p$  is even. For the important case of  $p = 3$  and small values of  $k$ , we can also estimate such norms. For asymmetric tensors sublinear time is not possible in general, but we show if the tensor slice norms are just slightly below  $\|T\|_F$  then sublinear time is again possible. One of the main strengths of our work is empirical - in a number of cases our algorithm is orders of magnitude faster than existing methods with the same accuracy.

## 1 Paper Body

Tensors are a powerful tool for dealing with multi-modal and multi-relational data. In recommendation systems, often using more than two attributes can lead to better recommendations. This could occur, for example, in Groupon where one could look at users, activities, and time (season, time of day, week-day/weekend, etc.), as three attributes to base predictions on (see [13] for a discussion). Similar to low rank matrix approximation, we seek a tensor decomposition to succinctly store the tensor and to apply it quickly. A popular decomposition method is the canonical polyadic decomposition, i.e., the CAN-DECOMP/PARAFAC (CP) decomposition, where the tensor is decomposed

into a sum of rank-1 components [9]. We refer the reader to [23], where applications of CP including data mining, computational neuroscience, and statistical learning for latent variable models are mentioned. A natural question, given the emergence of large data sets, is whether such decompositions can be performed quickly. There are a number of works on this topic [17, 16, 7, 11, 10, 4, 20]. Most related to ours are several recent works of Wang et al. [23] and Tung et al. [18], in which it is shown how to significantly speed up this decomposition for orthogonal tensor decomposition using the randomized technique of linear sketching [15]. In this work we also focus on orthogonal tensor decomposition. The idea in [23] is to create a succinct sketch of the input tensor, from which one can then perform implicit tensor decomposition by approximating inner products in existing decomposition methods. Existing methods, like the power method, involve computing the inner product of a vector, which is now a rank-1 matrix, with another vector, which is now a slice of a tensor. Such inner products can

Full version appears on arXiv, 2017. Work done while visiting IBM Almaden. Supported by XDATA DARPA Air Force Research Laboratory contract FA8750-12-C-0323.

30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

be approximated much faster by instead computing the inner product of the sketched vectors, which have significantly lower dimension. One can also replace the sketching with sampling to approximate inner products; we discuss some sampling schemes [17, 4] below and compare them to our work.

### 1.1 Our Contributions

We show in a number of important cases, one can achieve the same theoretical guarantees in the work of Wang et al. [23] (which was applied later by Tung et al. [18]), in sublinear time, that is, without reading most of the input tensor. While previous work needs to walk through the input at least once to create a sketch, we show one can instead perform importance sampling of the tensor based on the current iterate, together with reading a few entries of the tensor which help us learn the norms of tensor slices. We use a version of  $\ell_2$ -sampling for our importance sampling. One source of speedup in our work and in Wang et al. [23] comes from approximating inner products in iterations in the robust tensor power method (see below). To estimate  $\langle \mathbf{u}, \mathbf{v} \rangle$  for  $n$ -dimensional vectors  $\mathbf{u}$  and  $\mathbf{v}$ , their work computes sketches  $S(\mathbf{u})$  and  $S(\mathbf{v})$  and approximates  $\langle \mathbf{u}, \mathbf{v} \rangle \approx \langle S(\mathbf{u}), S(\mathbf{v}) \rangle$ . Instead, if one has  $\mathbf{u}$ , one can sample coordinates  $i$  proportional to  $u_i^2$ , which is known as  $\ell_2$ -sampling [14, 8]. One estimates  $\langle \mathbf{u}, \mathbf{v} \rangle$  as  $\frac{1}{k} \sum_{i=1}^k \frac{v_i}{u_i^2}$ , which is unbiased and has variance  $O(\frac{\|\mathbf{v}\|_2^2}{k \|\mathbf{u}\|_2^2})$ . These guarantees are similar to those using sketching, though the constants are significantly smaller (see below), and unlike sketching, one does not need to read the entire tensor to perform such sampling.

### Symmetric Tensors:

As in [23], we focus on orthogonal tensor decomposition of symmetric tensors, though we explain the extension to the asymmetric case below. Symmetric tensors arise in engineering applications, for example, to represent the symmetric tensor field of stress, strain, and anisotropic conductivity. Another example is diffusion MRI in which one uses symmetric tensors to describe diffusion in the brain or

other parts of the body. In spectral methods symmetric tensors are exactly those that come up in Latent Dirichlet Allocation problems. Although one can symmetrize a tensor using simple matrix operations (see, e.g., [1]), we cannot do this in sublinear time. In orthogonal tensor decomposition of a symmetric matrix, there is an underlying  $n \times n \times \dots \times n$  tensor  $\mathcal{T} = \sum_{i=1}^p v_i v_i^T$ , and the input tensor is  $\mathcal{T} = \mathcal{T} + \mathcal{E}$ , where  $\|\mathcal{E}\|_F \leq \epsilon$ . We have  $\|v_i\|_2 = 1$ ,  $\|v_i\|_2 \leq \epsilon$  and that  $\{v_i\}_{i=1}^p$  is a set of orthonormal vectors. The goal is to reconstruct  $\mathcal{T}$  to the  $\epsilon$ . Our results naturally generalize approximations  $v_i$  to the vectors  $v_i$ , and approximations  $\mathcal{T}$  to tensors with different lengths in different dimensions. For simplicity, we first focus on order  $p = 3$ . In the robust tensor power method [1], one generates a random initial vector  $u$ , and performs  $T$  update steps  $u = T(I, u, u)/\|T(I, u, u)\|_2$ , where  $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$ .  $\mathcal{T}(I, u, u) = T_{1,j} u_j u_j, T_{2,j} u_j u_j, \dots, T_{n,j} u_j u_j$ .  $j=1 \dots p$

The matrices  $T_{1,j}, \dots, T_{n,j}$  are referred to as the slices. The vector  $u$  typically converges to the top eigenvector in a small number of iterations, and one often chooses a small number  $L$  of random initial vectors to boost confidence. Successive eigenvectors can be found by deflation. The algorithm and analysis immediately extend to higher order tensors. We use  $\epsilon$ -sampling to estimate  $T(I, u, u)$ . To achieve the same guarantees as in [23], for typical settings of parameters (constant  $k$  and several eigenvalue assumptions) naively one needs to take  $O(n^2)$   $\epsilon$ -samples from  $u$  for each slice in each iteration, resulting in  $O(n^3)$  time and destroying our sublinearity. We observe that if we additionally knew the squared norms  $\|T_{1,j}\|_F^2, \dots, \|T_{n,j}\|_F^2, \|T\|_F^2$

then we could take  $O(n^2)$   $\epsilon$ -samples in total, where we take  $k_i \epsilon$   $O(n^2)$   $\epsilon$ -samples from the  $T_{i,j}$ -th slice in expectation. Perhaps in some applications such norms are known or cheap to compute in a single pass, but without further assumptions, how can one obtain such norms in sublinear time?  $\|T\|_F^2$  If  $\mathcal{T}$  is a symmetric tensor, then  $T_{j,j,j} = \sum_{i=1}^p v_{i,j} v_{i,j} v_{i,j} + E_{j,j,j}$ . Note that if there were no noise,  $\|T_{j,j,j}\|_F^2 = \sum_{i=1}^p v_{i,j}^2$ , and so

$\|T_{j,j,j}\|_F^2$  is an approximation to  $\|T_{j,j,j}\|_F^2$  up to factors depending on  $k$  and the eigenvalues. However, there is indeed noise. To obtain non-trivial guarantees, the robust tensor power method assumes  $\|\mathcal{E}\|_2 = O(1/n)$ , where  $\|\mathcal{E}\|_2 = \sup_{u,v,w} \sum_{i,j,k} \mathcal{E}_{i,j,k} u_i v_j w_k$

which in particular implies  $\|E_{j,j,j}\| = O(1/n)$ . This assumption comes from the  $O(1/n)$  correlation of the random initial vector to  $v_1$ . This noise bound does not trivialize the problem; indeed,  $E_{j,j,j}$  can be chosen adversarially subject to  $\|E_{j,j,j}\| = O(1/n)$ , and if the  $v_i$  were random  $\|v_i\| = O(1/n^{3/2})$ , which is small enough unit vectors and the  $\|v_i\|$  and  $k$  were constant, then  $\|v_i\|$  to be completely masked by the noise  $E_{j,j,j}$ . Nevertheless, there is a lot of information about the 3 slice norms. Indeed, suppose  $k = 1$ ,  $\|v_1\| = O(1/n)$ , and  $k \cdot T \cdot k_F = 1$ . Then  $T_{j,j,j} = \|v_1\| + E_{j,j,j}$ ,  $\|T_{j,j,j}\| = O(1/n)$  and one can show  $\|T_{j,j,j}\| = O(1/n)$ . Again using that  $\|E_{j,j,j}\| = O(1/n)$ , this implies  $\|T_{j,j,j}\| = O(1/n)$  if and only if  $\|T_{j,j,j}\| = O(1/n)$ , and therefore one would notice this by reading  $T_{j,j,j}$ . There can only be  $O(n^{2/3})$  slices  $j$  for which  $\|T_{j,j,j}\| = O(1/n)$ , since  $k \cdot T \cdot k_F = 1$ . Therefore, for each of them we can afford to take  $O(n^{2/3})$  samples and still have an  $O(n^{2+2/3}) = O(n^{8/3})$  sublinear running time. The remaining slices all have  $\|T_{j,j,j}\| = O(n^{2/3})$ , and therefore if we also take  $O(n^{1/3})$  samples from every slice, we will also estimate the contribution to  $T(I, u, u)$  from these slices well. This is also a sublinear  $O(n^{2+1/3})$  number of samples. While the previous paragraph illustrates the idea for  $k = 1$ , for  $k = 2$  we need to read more than the  $T_{j,j,j}$  entries to decide how many samples to take from a slice. The analysis is more complicated because of sign cancellations. Even for  $k = 2$  we could have  $T_{j,j,j} = \|v_1\| + \|v_2\| + E_{j,j,j}$ ,  $\|T_{j,j,j}\| = O(1/n)$  and if  $\|v_1\| = \|v_2\|$  then we may not detect that  $\|T_{j,j,j}\| = O(1/n)$  is large. We fix this by also reading the entries  $T_{j,j,i}$ ,  $T_{j,i,j}$ , and  $T_{j,i,i}$  for every  $i$  and  $j$ . This is still only  $O(n^2)$  entries and so we are still sublinear time. Without additional assumptions, we only give a formal analysis of this for  $k \in \{1, 2\}$ . More importantly, if instead of third-order symmetric tensors we consider  $p$ -th order symmetric tensors for even  $p$ , we do not have sign cancellations. In this case we do not have any restrictions on  $k$  for estimating slice norms. One does need to show after deflation, the slice norms can still be estimated; this holds because the eigenvectors and eigenvalues are estimated sufficiently well. We also give several per-iteration optimizations of our algorithm, based on careful implementations of generating a sorted list of random numbers and random permutations. We find empirically (see below) that we are much faster per iteration than previous sketching algorithms, in addition to not having to read the entire input tensor in a preprocessing step.

**Asymmetric Tensors:** For asymmetric tensors, e.g., 3rd-order tensors of the form  $\sum_{i=1}^n u_i \otimes v_i \otimes w_i$ , it is impossible to achieve sublinear time in general, since it is hard to distinguish  $T = e_i \otimes e_j \otimes e_k$  for random  $i, j, k \in \{1, 2, \dots, n\}$  from  $T = 0$ . We make a necessary and sufficient assumption that all the entries of the  $u_i$  are less than  $n^{-\epsilon}$  for an arbitrarily small constant  $\epsilon > 0$ . In this case, all slice norms are  $O(n^{-\epsilon})$  and by taking  $O(n^{\epsilon})$  samples from each slice we achieve sublinear time. We can also apply such an assumption to symmetric tensors.

**Empirical Results:** One of the main strengths of our work is our empirical results. In each iteration we approximate  $T(I, u, u)$  a total of  $B$  times independently and take the median to increase our confidence. In the notation of [23],  $B$  corresponds to the number of independent sketches used. While the median works empirically, there are some theoretical issues with it

discussed in Remark 4. Also let  $b$  be the total number of ‘2’-samples we take per iteration, which corresponds to the sketch size in the notation of [23]. We found that empirically we can set  $B$  and  $b$  to be much smaller than that in [23] and achieve the same error guarantees. One explanation for this is that the variance bound we obtain via importance sampling is a factor of  $43 = 64$  smaller than in [23], and for  $p$ -th order tensors, a factor of  $4p$  smaller. To give an idea of how much smaller we can set  $b$  and  $B$ , to achieve roughly the same squared residual norm error on the synthetic data sets of dimension 1200 for finding a good rank-1 approximation, the algorithm of [23] would need to set parameters  $b = 216$  and  $B = 50$ , whereas we can set  $b = 10 \cdot 1200$  and  $B = 5$ . Our running time is 2.595 seconds and we have no preprocessing time, whereas the algorithm of [23] has a running time of 116.3 seconds and 55.34 seconds of preprocessing time. We refer the reader to Table 1 in Section 3. In total we are over 50 times faster. We also demonstrate our algorithm in a real-world application using real datasets, even when the datasets are sparse. Namely, we consider a spectral algorithm for Latent Dirichlet Allocation [1, 2] which uses tensor decomposition as its core computational step. We show a significant speedup can be achieved on tensors occurring in applications such as LDA, and we refer the reader to Table 2 in

3

Section 3. For example, on the wiki [23] dataset with a tensor dimension of 200, we run more than 5 times faster than the sketching-based method. Previous Sampling Algorithms: Previous sampling-based schemes of [17, 4] do not achieve our guarantees, because [17] uses uniform sampling, which does not work for tensors with spiky elements, while the non-uniform sampling in [4] requires touching all of the entries in the tensor and making two passes over it. Notation Let  $[n]$  denote  $\{1, 2, \dots, n\}$ . Let  $\otimes$  denote the outer product, and  $u \otimes u = u \otimes u$ . Let  $T \in \mathbb{R}^{n_1 \times \dots \times n_p}$ , where  $p$  is the order of tensor  $T$  and  $n_i$  is the dimension of ptensor  $T$ . Let  $\langle A, B \rangle$  denote inner product between two tensors  $A, B \in \mathbb{R}^{n_1 \times \dots \times n_p}$ , e.g.,  $\langle A, B \rangle = \sum_{i_1=1}^{n_1} \dots \sum_{i_p=1}^{n_p} A_{i_1 \dots i_p} B_{i_1 \dots i_p}$ . For a tensor  $A \in \mathbb{R}^{n_1 \times \dots \times n_p}$ ,  $A_{i_1 \dots i_p} = A_{i_1, i_2, \dots, i_p}$ . For random variable  $X$  let  $E[X]$  denote its expectation of  $X$  and  $V[X]$  its variance (if these quantities exist).

2

## Main Results

We explain the details of our main results in this section. First, we state the importance sampling lemmas for our tensor application. Second, we explain how to quickly produce a list of random tuples according to a certain distribution needed by our algorithm. Third, we combine the first and the second parts to get a fast way of approximating tensor contractions, which are used as subroutines in each iteration of the robust tensor power method. We then provide our main theoretical results, and how to estimate the slice norms needed by our main algorithm. Importance sampling lemmas. Approximating an inner product is a simple application of importance sampling. Tensor contraction  $T(u, v, w)$  can be regarded as the inner product between two  $n^3$ -dimensional vectors, and thus importance sampling can be applied. Lemma 1 suggests that we can take a few

samples according to their importance, e.g., we can sample  $T_{i,j,k} u_i v_j w_k$  with probability  $u_i v_j w_k / (u_i^2 / k_{u,22} + v_j^2 / k_{v,22} + w_k^2 / k_{w,22})$ . As the number of samples is large enough, it will approximate the true tensor contraction  $\sum_{i,j,k} T_{i,j,k} u_i v_j w_k$  with small variance after a final rescaling. Lemma 1. Suppose random variable  $X = T_{i,j,k} u_i v_j w_k / (u_i^2 / k_{u,22} + v_j^2 / k_{v,22} + w_k^2 / k_{w,22})$  with probability  $p_i q_j r_k$  where  $p_i = u_i^2 / k_{u,22}$ ,  $q_j = v_j^2 / k_{v,22}$ , and  $r_k = w_k^2 / k_{w,22}$ , and we take  $L$  i.i.d. samples of  $X$ , PL denoted  $X_1, X_2, \dots, X_L$ . Let  $Y = \frac{1}{L} \sum_{i=1}^L X_i$ . Then (1)  $E[Y] = \sum_{i,j,k} T_{i,j,k} u_i v_j w_k$ , and (2)  $V[Y] \leq \frac{1}{L} \sum_{i,j,k} T_{i,j,k}^2 u_i^2 v_j^2 w_k^2$ . Similarly, we also have importance sampling for each slice  $T_{i,?,?}$ , i.e., each row of  $T$ . Lemma 2. For all  $i \in [n]$ , suppose random variable  $X_i = T_{i,j,k} v_j w_k / (q_j^2 / k_{v,22} + r_k^2 / k_{w,22})$  with probability  $q_j r_k$ , where  $q_j = v_j^2 / k_{v,22}$  and  $r_k = w_k^2 / k_{w,22}$ , and we take  $L_i$  i.i.d. samples of  $X_i$ , say PL  $X_{i1}, X_{i2}, \dots, X_{iL_i}$ . Let  $Y_i = \frac{1}{L_i} \sum_{i=1}^{L_i} X_{ii}$ . Then (1)  $E[Y_i] = \sum_{j,k} T_{i,j,k} v_j w_k$ , and (2)  $V[Y_i] \leq \frac{1}{L_i} \sum_{j,k} T_{i,j,k}^2 v_j^2 w_k^2$ . Generating importance samples in linear time. We need an efficient way to sample indices of a vector based on their importance. We view this problem as follows: imagine  $[0, 1]$  is divided into  $z$  bins with different lengths corresponding to the probability of selecting each bin, where  $z$  is the number of indices in a probability vector. We generate  $m$  random numbers uniformly from  $[0, 1]$  and see which bin each random number belongs to. If a random number is in bin  $i$ , we sample the  $i$ -th index of a vector. There are known algorithms [6, 19] to solve this problem in  $O(z + m)$  time. We give an alternative algorithm GENERATOR. Our algorithm combines Bentley and Saxe's algorithm [3] for efficiently generating  $m$  sorted random numbers in  $O(m)$  time, and Knuth's shuffling algorithm [12] for generating a random permutation of  $[m]$  in  $O(m)$  time. We use the notation  $CUMPROB(v, w)$  and  $CUMPROB(u, v, w)$  for the algorithm creating the distributions on  $2 \times 3 \times n$  and  $n \times n$  of Lemma 2 and Lemma 1, respectively. We note that naively applying previous algorithms would require  $z = O(n^2)$  and  $z = O(n^3)$  time to form these two distributions, but we can take  $O(m)$  samples from them implicitly in  $O(n + m)$  time. Fast approximate tensor contractions. We propose a fast way to approximately compute tensor contractions  $T(I, v, w)$  and  $T(u, v, w)$  with a sublinear number of samples of  $T$ , as shown in Algorithm 1 and Algorithm 2. Naively computing tensor contractions using all of the entries of  $T$  gives an exact answer but could take  $n^3$  time. Also, to keep our algorithm sublinear time, we never explicitly compute the deflated tensor; rather we represent it implicitly and sample from it. 4

Algorithm 1 Subroutine for approximate tensor contraction  $T(I, v, w)$  Algorithm 2 Subroutine for approximate tensor contraction  $T(u, v, w)$  1: function APPROXIMATE\_TENSOR\_CONTRACTION( $T, v, w, n, B, \{b_i\}$ ) 2:  $q, r \leftarrow CUMPROB(v, w)$  3: for  $d = 1$  to  $B$  do 4:  $L \leftarrow GENERATOR(\{b_i\}, q, r)$  5: for  $i = 1$  to  $n$  do (d) 6:  $s_i \leftarrow 0$  for  $\ell = 1$  to  $b_i$  do (j, k)  $\leftarrow L(i, \ell)$  7:  $s_i \leftarrow s_i + T_{i,j,k} v_j w_k$  8: 9:  $s_i \leftarrow s_i / b_i$

(d)

```

b v, w) i ? 10: T(I,
k
8: s(d) ? s(d) / bb b 9: T(u, v, w) ? median s(d)
1 qj rk Ti,j,k ?uj ? uk (d) median si /bbi , ?i ? [n] d?[B]
? si
1: function A PPROX T UVW(T, u, v, w, n, B, b b) 2: pe, qe, re ? C UM
P ROB(u, v, w) 3: for d = 1 ? B do 4: L ? G EN R AND T UPLES(bb, pe,
qe, re). 5: s(d) ? 0 6: for (i, j, k) ? L do 7: s(d) ? s(d) + pi qlj r Ti,j,k ?ui ?
uj ? uk
+
d?[B]
b 10: return T(u, v, w)
b v, w) 11: return T(I,

```

The following theorem gives the error bounds of A PPROX TI VW and A PPROX T UVW (in Algorithm 1 and 2). Let bbi be the number samples we take from slice  $i \in [n]$  in A PPROX TI VW, and let bb denote the total number of samples in our algorithm. Theorem 3. For  $T \in \mathbb{R}^{n \times n \times n}$  and  $u \in \mathbb{R}^n$  with  $\|u\|_2 = 1$ , define the number  $\|T(u)\|_2 = \|b \cdot u, u\|_2 = \|T(I, u, u)\|_2$ . For any  $b \geq 0$ , if  $\|T(u, u, u) - T(u, u, u)\|_2 \leq b$  and the vector  $\|T(u)\|_2 = \|T(I, bbi \otimes bk \cdot Ti, \cdot, \cdot)\|_2 / k \cdot T \cdot k^2$  then the following bounds hold  $1 : F F E[\|T(u)\|_2^2] = O(k \cdot T \cdot k^2 F / b)$ , and  $E[\|T(u)\|_2^2] = O(nk \cdot T \cdot k^2 F / b)$ . In addition, for any fixed  $\|T(u)\|_2 = 1$ ,  $E[\|T(u)\|_2^2] = O(k \cdot T \cdot k^2 F / b)$ .

(1)

Eq. (1) can be obtained by observing that each random variable  $\|T(u)\|_2^2$  is independent and so  $\sum_{i=1}^n \|T(u)\|_2^2 = \sum_{i=1}^n \|T(u)\|_2^2$ .  $(\sum_{i=1}^n \|T(u)\|_2^2) \cdot b \cdot F = b \cdot F \cdot V[\|T(u)\|_2^2] = \sum_{i=1}^n \|T(u)\|_2^2 \cdot b \cdot F$

Remark 4. In [23], the coordinate-wise median of B estimates to the  $T(I, v, w)$  is used to boost the success probability. There appears to be a gap [21] in their argument as it is unclear how to achieve (1) after taking a coordinate-wise median, which is (7) in Theorem 1 of [23]. To fix this, we instead pay a factor proportional to the number of iterations in Algorithm 3 in the sample complexity bb. Since we have expectation bounds on the quantities in Theorem 3, we can apply a Markov bound and a union bound across all iterations. This suffices for our main theorem concerning sublinear time below. One can obtain high probability bounds by running Algorithm 3 multiple times independently, and taking coordinate-wise medians of the output eigenvectors. Empirically, our algorithm works even if we take the median in each iteration, which is done in line 10 in Algorithm 1. Replacing Theorem 1 in [23] by our Theorem 3, the rest of the analysis in [23] is unchanged. Our Algorithm 3 is the same as the sketching-based robust tensor power method in [23], except for lines 10, 12, 15, and 17, where the sketching-based approximate tensor contraction is replaced by our importance sampling procedures A PPROX T UVW and A PPROX TI VW. Rather than use Theorem 2 of Wang et al. [23], the main theorem concerning the correctness of the robust tensor decomposition algorithm, we use a recent improvement of it by Wang and Anandkumar in Theorems 4.1





1000

1200

(b) Preprocessing time

Figure 1: Running time with growing dimension

for all  $i \in [k]$ ,  $t \in [T]$ , and  $\ell \in [L]$  and furthermore  $\sum_{\ell=1}^L C_1 \sum_{t=1}^T \min_{i \in [k]} \|T_{i,t,\ell}\|_F \leq \epsilon$ ,  $L \leq \max\{K_0, k\} \log(\max\{K_0, k\})$ , then with probability at least  $9/10$ , there exists a permutation  $\pi : [k] \rightarrow [k]$  such that  $b(\pi(i)) = C_2 \sum_{t=1}^T \sum_{\ell=1}^L \|T_{i,t,\ell}\|_F^2$ ,  $\forall i \in [k]$ . Combining the previous theorem with our importance sampling analysis, we obtain: Theorem 6 (Main). Assume the notation of Theorem 5. For each  $j \in [k]$ , suppose we take  $b(j) = \sum_{t=1}^T \sum_{\ell=1}^L \|T_{j,t,\ell}\|_F^2$  samples during the power iterations for recovering  $\sum_{i=1}^k b_i P P^T(j) j^T j^T \sum_{t=1}^T \sum_{\ell=1}^L \|T_{i,t,\ell}\|_F^2$  where  $b = \sum_{i=1}^k b_i$  and  $\sum_{i=1}^k b_i = 1$ . Then the output guarantees of Theorem 5 hold for Algorithm 3 with constant probability. Our total time is  $O(LT \sum_{j=1}^k b(j))$  and the space is  $O(nk)$ , where  $bb = \max_{j \in [k]} b(j)$ . In Theorem 3, if we require  $b_i = \sum_{t=1}^T \sum_{\ell=1}^L \|T_{i,t,\ell}\|_F^2$ , we need to scan the entire tensor to compute  $\sum_{t=1}^T \sum_{\ell=1}^L \|T_{i,t,\ell}\|_F^2$ , making our algorithm not sublinear. With the following mild assumption in Theorem 7, our algorithm is sublinear when sampling uniformly ( $b_i = b/n$ ) without computing  $\sum_{t=1}^T \sum_{\ell=1}^L \|T_{i,t,\ell}\|_F^2$ : Theorem 7 (Bounded slice norm). There is a constant  $\epsilon > 0$ , a constant  $\delta \in (0, 1]$  and a sufficiently small constant  $\epsilon > 0$ , such that, for any 3rd order tensor  $T = T^* + E \in \mathbb{R}^{n \times n \times n}$  with  $\text{rank}(T^*) \leq n^\delta$ ,  $\|E\|_F \leq \epsilon/n$ , if  $\sum_{t=1}^T \sum_{\ell=1}^L \|T_{i,t,\ell}\|_F^2 \leq n^\delta$  for all  $i \in [n]$ , and  $E$  satisfies (2), then Algorithm 3 runs in  $O(n^3 \delta)$  time. The condition  $\delta \in (0, 1]$  is a practical one. When  $\delta = 1$ , all tensor slices have equal Frobenius norm. The case  $\delta = 0$  only occurs when  $\sum_{t=1}^T \sum_{\ell=1}^L \|T_{i,t,\ell}\|_F^2 = 0$ ; i.e., all except one slice is zero. This theorem can also be applied to asymmetric tensors, since the analysis in [23] can be extended to them. For certain cases, we can remove the bounded slice norm assumption. The idea is to take a sublinear number of samples from the tensor to obtain upper bounds on all slice norms. In the full version, we extend the algorithm and analysis of the robust tensor power method to  $p \geq 3$  by replacing contractions  $T(u, v, w)$  and  $T(I, v, w)$  with  $T(u_1, u_2, \dots, u_p)$  and  $T(I, u_2, \dots, u_p)$ . As outlined in Section 1, when  $p$  is even, because we do not have sign cancellations we can show: Theorem 8 (Even order). There is a constant  $\epsilon > 0$  and a sufficiently small constant  $\delta > 0$ , such that, for any even order- $p$  tensor  $T = T^* + E \in \mathbb{R}^{n \times n \times \dots \times n}$  with  $\text{rank}(T^*) \leq n^\delta$ ,  $p \leq n^\delta$  and  $\|E\|_F \leq \epsilon/n$ . For any sufficiently large constant  $c_0$ , there exists a sufficiently small constant  $c > 0$ , for any  $\delta \in (0, c/k / (c_0 p^2 \ln(p^2/2))$  if  $E$  satisfies  $\sum_{t=1}^T \sum_{\ell=1}^L \|T_{i,t,\ell}\|_F^2 \leq c/(c_0 n)$ , Algorithm 3 runs in  $O(n p^2 \delta)$  time. 6

As outlined in Section 1, for  $p = 3$  and small  $k$  we can take sign considerations into account: Theorem 9 (Low rank). There is a constant  $\epsilon > 0$  and a sufficiently small constant  $\delta > 0$  such that for any symmetric tensor  $T = T^* + E \in \mathbb{R}^{n \times n \times n}$  with  $E$  satisfying (2),  $\text{rank}(T^*) \leq 2$ , and  $\|E\|_F \leq \epsilon/n$ , then Algorithm 3 runs in  $O(n^3 \delta)$  time.

3

Experiments

3.1 Experiment Setup and Datasets Our implementation shares the same code base 1 as the sketching-based robust tensor power method proposed in [23]. We ran our experiments on an i7-5820k CPU with 64 GB of memory in singlethreaded mode. We ran two versions of our algorithm: the version with pre-scanning scans the full tensor to accurately measure per-slice Frobenius norms and make samples for each slice in proportion to its Frobenius norm in APPROX-TI-VW; the version without pre-scanning assumes that the Frobenius norm of each slice is bounded by  $n^{1/2} k T k_F$ ,  $\alpha \in (0, 1]$  and uses  $b/n$  samples per slice, where  $b$  is the total number of samples our algorithm makes, analogous to the sketch length  $b$  in [23]. Synthetic datasets. We first generated an orthonormal basis  $\{v_i\}_{i=1}^n$  and then computed the synthetic  $P_k$  tensor as  $T = \sum_{i=1}^n \lambda_i v_i^{\otimes 3}$ , with  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$ . Then we normalized  $T$  such that  $k T k_F = 1$ , and added a symmetric Gaussian noise tensor  $E$  where  $E_{ijl} \sim N(0, n^{-1.5})$  for  $i, j, l \in [n]$ . Then  $\alpha$  controls the noise-to-signal ratio and we kept it as 0.01 in all our synthetic tensors. For the eigenvalues  $\lambda_i$ , we generated three different decays: inverse decay  $\lambda_i = 1/i$ , inverse square decay  $\lambda_i = 1/i^2$ , and linear decay  $\lambda_i = 1 - i/n$ . We also set  $k = 100$  when generating tensors, since higher rank eigenvalues were almost indistinguishable from the added noise. To show the scalability of our algorithm, we generated tensors with different dimensions:  $n = 200, 400, 600, 800, 1000, 1200$ . Real-life datasets. Latent Dirichlet Allocation [5] (LDA) is a powerful generative statistical model for topic modeling. A spectral method has been proposed to solve LDA models [1, 2] and the most critical step in spectral LDA is to decompose a symmetric  $K \times K$  tensor with orthogonal eigenvectors, where  $K$  is the number of modeled topics. We followed the steps in [1, 18] and built a  $K \times K \times K$  tensor TLDA for each dataset, and then ran our algorithms directly on TLDA to see how it works on those tensors in real applications. In our experiments we keep  $K = 200$ . We used the two same datasets as the previous work [23]: Wiki and Enron, as well as four additional real-life datasets. We refer the reader to our GitHub repository 2 for our code and full results. 3.2 Results We considered running time and the squared residual norm to evaluate the performance of our  $P_k$  algorithms. Given a tensor  $T \in \mathbb{R}^{n \times n \times n}$ , let  $k T = \sum_{i=1}^k u_i \otimes v_i \otimes w_i k_F^2$  denote the squared residual norm where  $\{(\lambda_1, u_1, v_1, w_1), \dots, (\lambda_k, u_k, v_k, w_k)\}$  are the eigenvalue/eigenvectors obtained by the robust power method. To reduce the experiment time we looked only for the first eigenvalue and eigenvector, but our algorithm is capable of finding any number of eigenvalues/eigenvectors. We list the pre-scanning time as preprocessing time in tables. It only depends on the tensor dimension  $n$  and unlike the sketching based method, it does not depend on  $b$ . Pre-scanning time is very short, because it only requires one pass of sequential access to the tensor which is very efficient on hardware. Sublinear time verification. Our theoretical result suggests the total number of samples  $b_{no-prescan}$  for our algorithm without pre-scanning is  $n^{1/\alpha} (\alpha \in (0, 1])$  times larger than  $b_{prescan}$  for our algorithm with pre-scanning. But in experiments we observe that when  $b_{no-prescan} = b_{prescan}$  both algorithms achieve very similar accuracy, indicating that in practice  $\alpha \approx 1$ . Synthetic datasets. We ran our algorithm on a large number of synthetic tensors with different dimensions

and different eigengaps. Table 1 shows results for a tensor with 1200 dimensions with 100 non-zero eigenvalues decaying as  $\frac{1}{i^2}$ . To reach roughly the same residual norm, the running time of our algorithm is over 50 times faster than that of the sketching-based robust tensor power method, thanks to the fact that we usually need a relatively small  $B$  and  $b$  to get a good residual, and the hidden constant factor in the running time of sampling is much smaller than that of sketching. Our algorithm scales well on large tensors due to its sub-linear nature. In Figure 1(a), for the sketching-based method we kept  $b = 216$ ,  $B = 30$  for  $n \leq 800$  and  $B = 50$  for  $n > 800$  (larger  $n$  requires more sketches to observe a reasonable recovery). For our algorithm, we chose  $b$  and  $B$  such

<http://yining-wang.com/fftlda-code.zip> [https://github.com/huanzhang12/sampling\\_tensor\\_decomp/](https://github.com/huanzhang12/sampling_tensor_decomp/)

that for each  $n$ , our residual norm is on-par or better than the sketching-based method. Our algorithm needs much less time than the sketching-based one over all dimensions. Another advantage of our algorithm is that there are zero or very minimal preprocessing steps. In Figure 1(b), we can see how the preprocessing time grows to prepare sketches when the dimension increases. For applications where only the first few eigenvectors are needed, the preprocessing time could be a large overhead. Real-life datasets Due to the small tensor dimension (200), our algorithm shows less speedup than the sketching-based method. But it is still 2 - 6 times faster in each of the six real-life datasets, achieving the same squared residual norm. Table 2 reports results for one of the datasets in many different settings of  $(b, B)$ . Like in synthetic datasets, we also empirically observe that the constant  $b$  in importance sampling is much smaller than the  $b$  used in sketching to get the same error guarantee.

$b$   
 $b$   
 $b$

Sketching based robust power method:  $n = 1200$ ,  $\frac{1}{i^2}$  Squared residual norm Running time (s) Preprocessing time (s)  $B$  10 30 50 10 30 50 10 30 50  
210 1.010 1.014 0.5437 0.6114 2.423 4.374 5.361 15.85 26.08 212 1.020 0.2271  
0.1549 1.344 4.563 8.022 5.978 17.23 28.31 214 0.1513 0.1097 0.1003 4.928 15.51  
27.87 8.788 24.72 40.4 16 2 0.1065 0.09242 0.08936 22.28 69.7 116.3 13.76 34.74  
55.34 Importance sampling based robust power method (without prescanning):  
 $n = 1200$ ,  $\frac{1}{i^2}$  Squared residual norm Running time (s) Preprocessing time  
(s)  $B$  10 30 50 10 30 50 10 30 50 5n 0.08684 0.08637 0.08639 2.595 8.3 15.46  
0.0 0.0 0.0 10n 0.08784 0.08671 0.08627 4.42 13.68 25.84 0.0 0.0 0.0 20n 0.08704  
0.08700 0.08618 8.02 24.51 46.37 0.0 0.0 0.0 30n 0.08697 0.08645 0.08625 11.63  
35.35 66.71 0.0 0.0 0.0 40n 0.08653 0.08664 0.08611 15.19 46.12 87.24 0.0 0.0  
0.0 1 Importance sampling based robust power method (with prescanning):  $n$   
 $= 1200$ ,  $\frac{1}{i^2}$  Squared residual norm Running time (s) Preprocessing time (s)  
 $B$  10 30 50 10 30 50 10 30 50 5n 0.08657 0.08684 0.08636 3.1 10.47 18 2.234  
2.236 2.234 10n 0.08741 0.08677 0.08668 5.427 17.43 30.26 2.232 2.233 2.233 20n  
0.08648 0.08624 0.08634 9.843 31.42 54.49 2.226 2.226 2.226 30n 0.08635 0.08634  
0.08615 14.33 45.4 63.85 2.226 2.224 2.227 40n 0.08622 0.08652 0.08619 18.68  
59.32 82.83 2.225 2.225 2.225

Table 1: Synthetic tensor decomposition using the robust tensor power method. We use an order-3 normalized dense tensor with dimension  $n = 1200$  with  $\epsilon = 0.01$  noise added. We run sketching-based and sampling-based methods to find the first eigenvalue and eigenvector by setting  $L = 50$ ,  $T = 30$  and varying  $B$  and  $b$ . Sketching based robust power method: dataset wiki,  $kTk2F = 2.135e+07$  Squared residual norm Running time (s) Preprocessing time (s)  $B$  10 30 10 30 10 30  $b$  210 2.091e+07 1.951e+07 0.2346 0.8749 0.1727 0.2535 211 1.971e+07 1.938e+07 0.4354 1.439 0.2408 0.3167 212 1.947e+07 1.930e+07 1.035 2.912 0.4226 0.4275 213 1.931e+07 1.927e+07 2.04 5.94 0.5783 0.6493 14 2 1.928e+07 1.926e+07 4.577 13.93 1.045 1.121 Importance sampling based robust power method (without prescanning): dataset wiki,  $kTk2F = 2.135e+07$  Squared residual norm Running time (s) Preprocessing time (s)  $B$  10 30 10 30 10 30  $b$  5n 1.931e+07 1.928e+07 0.3698 1.146 0.0 0.0 10n 1.931e+07 1.929e+07 0.5623 1.623 0.0 0.0 20n 1.935e+07 1.926e+07 0.9767 2.729 0.0 0.0 30n 1.929e+07 1.926e+07 1.286 3.699 0.0 0.0 40n 1.928e+07 1.925e+07 1.692 4.552 0.0 0.0 Importance sampling based robust power method (with prescanning): dataset wiki,  $kTk2F = 2.135e+07$  Squared residual norm Running time (s) Preprocessing time (s)  $B$  10 30 10 30 10 30  $b$  5n 1.931e+07 1.930e+07 0.4376 1.168 0.01038 0.01103 10n 1.928e+07 1.930e+07 0.6357 1.8 0.0104 0.01044 20n 1.931e+07 1.927e+07 1.083 2.962 0.01102 0.01042 30n 1.929e+07 1.925e+07 1.457 4.049 0.01102 0.01043 40n 1.929e+07 1.925e+07 1.905 5.246 0.01105 0.01105

Table 2: Tensor decomposition in LDA on the wiki dataset. The tensor is generated by spectral LDA with dimension  $200 \times 200 \times 200$ . It is symmetric but not normalized. We fix  $L = 50$ ,  $T = 30$  and vary  $B$  and  $b$ .

8

## 2 References

- [1] A. Anandkumar, R. Ge, D. Hsu, S. M. Kakade, and M. Telgarsky. Tensor decompositions for learning latent variable models. *JMLR*, 15(1):2773?2832, 2014.
- [2] A. Anandkumar, Y.-k. Liu, D. J. Hsu, D. P. Foster, and S. M. Kakade. A spectral algorithm for latent dirichlet allocation. In *NIPS*, pages 917?925, 2012.
- [3] J. L. Bentley and J. B. Saxe. Generating sorted lists of random numbers. *ACM Transactions on Mathematical Software (TOMS)*, 6(3):359?364, 1980.
- [4] S. Bhojanapalli and S. Sanghavi. A new sampling technique for tensors. *CoRR*, abs/1502.05023, 2015.
- [5] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *JMLR*, 3:993?1022, 2003.
- [6] K. Bringmann and K. Panagiotou. Efficient sampling methods for discrete distributions. In *International Colloquium on Automata, Languages, and Programming*, pages 133?144. Springer, 2012.
- [7] J. H. Choi and S. Vishwanathan. Dfacto: Distributed factorization of tensors. In *NIPS*, pages 1296?1304, 2014.
- [8] K. L. Clarkson, E. Hazan, and D. P. Woodruff. Sublinear optimization for machine learning. *J. ACM*, 59(5):23, 2012.
- [9] R. A. Harshman. Foundations of the parafac procedure: Models and conditions for an explanatory multi-modal factor analysis. 16(1):84, 1970.
- [10] F. Huang, N. U. N, M. U. Hakeem, P. Verma, and A. Anandkumar. Fast de-

tection of overlapping communities via online tensor methods on gpus. CoRR, abs/1309.0787, 2013. [11] U. Kang, E. E. Papalexakis, A. Harpale, and C. Faloutsos. Gigatensor: scaling tensor analysis up by 100 times - algorithms and discoveries. In KDD, pages 316?324, 2012. [12] D. E. Knuth. The art of computer programming. vol. 2: Seminumerical algorithms. addisonwesley. Reading, MA, pages 229?279, 1969. [13] A. Moitra. Tensor decompositions and their applications, 2014. [14] M. Monemizadeh and D. P. Woodruff. 1-pass relative-error lp -sampling with applications. In SODA, pages 1143?1160, 2010. [15] N. Pham and R. Pagh. Fast and scalable polynomial kernels via explicit feature maps. In KDD, pages 239?247, 2013. [16] A. H. Phan, P. Tichavsk?, and A. Cichocki. Fast alternating LS algorithms for high order CANDECOMP/PARAFAC tensor factorizations. IEEE Transactions on Signal Processing, 61(19):4834?4846, 2013. [17] C. E. Tsourakakis. MACH: fast randomized tensor decompositions. In SDM, pages 689?700, 2010. [18] H.-Y. F. Tung, C.-Y. Wu, M. Zaheer, and A. J. Smola. Spectral methods for the hierarchical dirichlet process. 2015. [19] A. J. Walker. An efficient method for generating discrete random variables with general distributions. ACM Transactions on Mathematical Software (TOMS), 3(3):253?256, 1977. [20] C. Wang, X. Liu, Y. Song, and J. Han. Scalable moment-based inference for latent dirichlet allocation. In ECML-PKDD, pages 290?305, 2014. [21] Y. Wang. Personal communication. 2016. [22] Y. Wang and A. Anandkumar. Online and differentially-private tensor decomposition. CoRR, abs/1606.06237, 2016. [23] Y. Wang, H.-Y. Tung, A. J. Smola, and A. Anandkumar. Fast and guaranteed tensor decomposition via sketching. In NIPS, pages 991?999, 2015.