

Better and Simpler Error Analysis of the Sinkhorn-Knopp Algorithm for Matrix Scaling*

Deeparnab Chakrabarty¹ and Sanjeev Khanna²

1 Department of Computer Science, Dartmouth College, Hanover NH, USA
deeparnab@dartmouth.edu

2 Department of Computer and Information Science, University of Pennsylvania,
Philadelphia PA, USA
sanjeev@cis.upenn.edu

Abstract

Given a non-negative $n \times m$ real matrix A , the *matrix scaling* problem is to determine if it is possible to scale the rows and columns so that each row and each column sums to a specified target value for it. The matrix scaling problem arises in many algorithmic applications, perhaps most notably as a preconditioning step in solving linear system of equations. One of the most natural and by now classical approach to matrix scaling is the Sinkhorn-Knopp algorithm (also known as the RAS method) where one alternately scales either all rows or all columns to meet the target values. In addition to being extremely simple and natural, another appeal of this procedure is that it easily lends itself to parallelization. A central question is to understand the rate of convergence of the Sinkhorn-Knopp algorithm.

Specifically, given a suitable error metric to measure deviations from target values, and an error bound ε , how quickly does the Sinkhorn-Knopp algorithm converge to an error below ε ? While there are several non-trivial convergence results known about the Sinkhorn-Knopp algorithm, perhaps somewhat surprisingly, even for natural error metrics such as ℓ_1 -error or ℓ_2 -error, this is not entirely understood. In this paper, we present an elementary convergence analysis for the Sinkhorn-Knopp algorithm that improves upon the previous best bound. In a nutshell, our approach is to show (i) a simple bound on the number of iterations needed so that the KL-divergence between the current row-sums and the target row-sums drops below a specified threshold δ , and (ii) then show that for a suitable choice of δ , whenever KL-divergence is below δ , then the ℓ_1 -error or the ℓ_2 -error is below ε . The well-known Pinsker's inequality immediately allows us to translate a bound on the KL divergence to a bound on ℓ_1 -error. To bound the ℓ_2 -error in terms of the KL-divergence, we establish a new inequality, referred to as (KL vs ℓ_1/ℓ_2). This new inequality is a strengthening of the Pinsker's inequality that we believe is of independent interest. Our analysis of ℓ_2 -error significantly improves upon the best previous convergence bound for ℓ_2 -error.

The idea of studying Sinkhorn-Knopp convergence via KL-divergence is not new and has indeed been previously explored. Our contribution is an elementary, self-contained presentation of this approach and an interesting new inequality that yields a significantly stronger convergence guarantee for the extensively studied ℓ_2 -error.

1998 ACM Subject Classification F.2.1.3 Computations on matrices

Keywords and phrases Matrix Scaling, Entropy Minimization, KL Divergence Inequalities

Digital Object Identifier 10.4230/OASIS.SOSA.2018.4

* This work was supported in part by the National Science Foundation grants CCF-1552909 and CCF-1617851.



© Deeparnab Chakrabarty and Sanjeev Khanna;
licensed under Creative Commons License CC-BY
1st Symposium on Simplicity in Algorithms (SOSA 2018).

Editor: Raimund Seidel; Article No. 4; pp. 4:1–4:11

Open Access Series in Informatics



OASIS Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

1 Introduction

In the matrix scaling problem one is given an $n \times m$ non-negative matrix A , and positive integer vectors $\mathbf{r} \in \mathbb{Z}_{>0}^n$ and $\mathbf{c} \in \mathbb{Z}_{>0}^m$ with the same ℓ_1 norm $\sum_{i=1}^n \mathbf{r}_i = \sum_{j=1}^m \mathbf{c}_j = h$. The objective is to determine if there exist diagonal matrices $R \in \mathbb{R}^{n \times n}$ and $S \in \mathbb{R}^{m \times m}$ such that the i th row of the matrix RAS sums to \mathbf{r}_i for all $1 \leq i \leq n$ and the j th column of RAS sums to \mathbf{c}_j for all $1 \leq j \leq m$. Of special importance is the case when $n = m$ and $\mathbf{r} \equiv \mathbf{c} \equiv \mathbf{1}_n$, the n -dimensional all-ones vector – the $(\mathbf{1}, \mathbf{1})$ -matrix scaling problem wishes to scale the rows and columns of A to make it doubly stochastic. This problem arises in many different areas ranging from transportation planning [12, 26] to quantum mechanics [32, 1]; we refer the reader to a recent comprehensive survey by Idel [15] for more examples.

One of the most natural algorithms for the matrix scaling problem is the following Sinkhorn-Knopp algorithm [33, 34], which is known by many names including the RAS method [4] and the Iterative Proportional Fitting Procedure [30]. The algorithm starts off by multiplicatively scaling all the columns by the columns-sum times \mathbf{c}_j to get a matrix $A^{(0)}$ with column-sums \mathbf{c} . Subsequently, for $t \geq 0$, it obtains the $B^{(t)}$ by scaling each row of $A^{(t)}$ by the respective row-sum times \mathbf{r}_i , and obtain $A^{(t+1)}$ by scaling each column of $B^{(t)}$ by the respective column sums time \mathbf{c}_j . More precisely,

$$A_{ij}^{(0)} := \frac{A_{ij}}{\sum_{i=1}^n A_{ij}} \cdot \mathbf{c}_j \quad \forall t \geq 0, \quad B_{ij}^{(t)} := \frac{A_{ij}^{(t)}}{\sum_{j=1}^m A_{ij}^{(t)}} \cdot \mathbf{r}_i, \quad A_{ij}^{(t+1)} := \frac{B_{ij}^{(t)}}{\sum_{i=1}^n B_{ij}^{(t)}} \cdot \mathbf{c}_j$$

The above algorithm is simple and easy to implement and each iteration takes $O(\text{nnz}(A))$, the number of non-zero entries of A . Furthermore, it has been known for almost five decades [33, 34, 13, 35] that if A is (\mathbf{r}, \mathbf{c}) -scalable then the above algorithm asymptotically¹ converges to a right solution. More precisely, given $\varepsilon > 0$, there is some finite t by which one obtains a matrix which is “ ε -close to having row- and column-sums \mathbf{r} and \mathbf{c} ”.

However, the rate of convergence of this simple algorithm is still not fully understood. Since the rate depends on how we measure “ ε -closeness”, we look at two natural error definitions. For any t , let $\mathbf{r}^{(t)} := A^{(t)} \mathbf{1}_m$ denote the vector of row-sums of $A^{(t)}$. Similarly, we define $\mathbf{c}^{(t)} := B^{(t)\top} \mathbf{1}_n$ to be the vector of the column-sums of $B^{(t)}$. Note that $\sum_{i=1}^n \mathbf{r}_i^{(t)} = \sum_{j=1}^m \mathbf{c}_j^{(t)} = h$ for all t . The error of the matrix A_t (the error of matrix B_t similarly defined) is

$$\ell_1\text{-error} : \text{error}_1(A_t) := \|\mathbf{r}^{(t)} - \mathbf{r}\|_1 \quad \ell_2\text{-error} : \text{error}_2(A_t) := \|\mathbf{r}^{(t)} - \mathbf{r}\|_2$$

In this note, we give simple convergence analysis for both error norms. Our result is the following.

¹ Computationally, this asymptotic viewpoint is unavoidable in the sense that there are simple examples for which the unique matrix scaling matrices need to have irrational entries. For instance, consider the following example from Rothblum and Schneider [29]. The matrix is $\begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}$ with $\mathbf{r} \equiv \mathbf{c} \equiv [1, 1]^\top$.

The unique R and S matrices are $\begin{bmatrix} (\sqrt{2}+1)^{-1} & 0 \\ 0 & (\sqrt{2}+2)^{-1} \end{bmatrix}$ and $\begin{bmatrix} \sqrt{2} & 0 \\ 0 & 1 \end{bmatrix}$, respectively, giving

$$RAS = \begin{bmatrix} 2 - \sqrt{2} & \sqrt{2} - 1 \\ \sqrt{2} - 1 & 2 - \sqrt{2} \end{bmatrix}.$$

► **Theorem 1.** Given a matrix $A \in \mathbb{R}_{\geq 0}^{n \times m}$ which is (\mathbf{r}, \mathbf{c}) -scalable, and any $\varepsilon > 0$, the Sinkhorn-Knopp algorithm

1. in time $t = O\left(\frac{h^2 \ln(n\rho/\nu)}{\varepsilon^2}\right)$ returns a matrix A_t or B_t with ℓ_1 -error $\leq \varepsilon$.
2. in time $t = O\left(\rho h \ln(n\rho/\nu) \cdot \left(\frac{1}{\varepsilon} + \frac{1}{\varepsilon^2}\right)\right)$ returns a matrix A_t or B_t with ℓ_2 -error $\leq \varepsilon$.

Here $h = \sum_{i=1}^n \mathbf{r}_i = \sum_{j=1}^m \mathbf{c}_j$, $\rho = \max(\max_i \mathbf{r}_i, \max_j \mathbf{c}_j)$, and $\nu = \frac{\min_{i,j: A_{ij} > 0} A_{ij}}{\max_{i,j} A_{ij}}$.

For the special case of $n = m$ and $\mathbf{r} \equiv \mathbf{c} \equiv \mathbf{1}_n$, we get the following as a corollary.

► **Corollary 2.** Given a matrix $A \in \mathbb{Z}_{\geq 0}^{n \times n}$ which is $(\mathbf{1}, \mathbf{1})$ -scalable, and any $\varepsilon > 0$, the Sinkhorn-Knopp algorithm

1. in time $t = O\left(\frac{n^2 \ln n}{\varepsilon^2}\right)$ returns a matrix A_t or B_t with ℓ_1 -error $\leq \varepsilon$.
2. in time $t = O\left(n \ln n \cdot \left(\frac{1}{\varepsilon} + \frac{1}{\varepsilon^2}\right)\right)$ returns a matrix A_t or B_t with ℓ_2 -error $\leq \varepsilon$.

► **Remark.** To our knowledge, the ℓ_1 -error hasn't been explicitly studied in the literature, although for small $\varepsilon \in (0, 1)$ the same can be deduced from previous papers on matrix scaling [20, 14, 19, 16]. One of our main motivations to look at ℓ_1 -error arose from the connections to perfect matchings in bipartite graphs as observed by Linial, Samorodnitsky and Wigderson [20]. For the ℓ_2 error, which is the better studied notion in the matrix scaling literature, the best analysis is due to Kalantari et al [18, 19]. They give a $\tilde{O}(\rho h^2 / \varepsilon^2)$ upper bound on the number of iterations for the general problem, and for the special case when $m = n$ and the square matrix has positive permanent (see [18]), they give a $\tilde{O}(\rho(h^2 - nh + n) / \varepsilon^2)$ upper bound. Thus, for $(\mathbf{1}, \mathbf{1})$ -scaling, they get the same result as in Corollary 2. We get a quadratic improvement on h in the general case, and we think our proof is more explicit and simpler.

► **Remark.** Both parts of Theorem 1 and Corollary 2 are interesting in certain regimes of error. When the error ε is “small” (say, ≤ 1) so that $1/\varepsilon^2 \geq 1/\varepsilon$, then statement 2 of Corollary 2 implies statement 1 by Cauchy-Schwarz. However, this breaks down when ε is “large” (say $\varepsilon = \delta n$ for some constant $\delta > 0$). In that case, statement 1 implies that in $O(\ln n / \delta^2)$ iterations, the ℓ_1 -error is $\leq \delta n$, but Statement 2 only implies that in $O(\ln n / \delta^2)$ iterations, the ℓ_2 norm is $\leq \delta n$. This “large ℓ_1 -error regime” is of particular interest for an application to approximate matchings in bipartite graphs discussed below.

Applications to Parallel Algorithms for Bipartite Perfect Matching. As a corollary, we get the following application, first pointed by Linial et al [20], to the existence of perfect matchings in bipartite graphs. Let A be the adjacency matrix of a bipartite graph $G = (L \cup R, E)$ with $A_{ij} = 1$ iff $(i, j) \in E$. If G has a perfect matching, then clearly there is a doubly stochastic matrix X in the support of A . This suggests the algorithm of running the Sinkhorn-Knopp algorithm to A , and the following claim suggests when to stop. Note that each iteration can be run in $O(1)$ parallel time with m -processors.

► **Lemma 3.** If we find a column (or row) stochastic matrix Y in the support of A such that $\text{error}_1(Y) \leq \varepsilon$, then G has a matching of size $\geq n(1 - \varepsilon)$.

Proof. Suppose Y is column stochastic. Given $S \subseteq L$, consider $\sum_{i \in S, j \in \Gamma S} Y_{ij} = |S| + \sum_{i \in S} \left(\sum_{j=1}^n Y_{ij} - 1\right) \geq |S| - \sum_{i=1}^n \left|\sum_{j=1}^n Y_{ij} - 1\right| \geq |S| - n \cdot \text{error}(Y) \geq |S| - n\varepsilon$. On the other hand, $\sum_{i \in S, j \in \Gamma S} Y_{ij} \leq \sum_{j \in \Gamma S} \sum_{i=1}^n Y_{ij} = |\Gamma S|$. Therefore, for every $S \subseteq L$, $|\Gamma S| \geq |S| - n\varepsilon$. The claim follows by approximate Hall's theorem. ◀

► **Corollary 4** (Fast Parallel Approximate Matchings). Given a bipartite graph G of max-degree Δ and an $\varepsilon \in (0, 1)$, $O(\ln \Delta / \varepsilon^2)$ -iterations of Sinkhorn-Knopp algorithm suffice to distinguish

between the case when G has a perfect matching and the case when the largest matching in G has size at most $n(1 - \varepsilon)$.

Thus the approximate perfect matching problem in bipartite graphs is in NC for ε as small as polylogarithmic in n . This is not a new result and can indeed be obtained from the works on parallel algorithms for packing-covering LPs [21, 36, 3, 23], but the Sinkhorn-Knopp algorithm is arguably simpler.

1.1 Perspective

As mentioned above, the matrix scaling problem and in particular the Sinkhorn-Knopp algorithm has been extensively studied over the past 50 years. We refer the reader to Idel’s survey [15] and the references within for a broader perspective; in this subsection we mention the most relevant works.

We have already discussed the previously best known, in their dependence on h , analysis for the Sinkhorn-Knopp algorithm in Remark 1. For the special case of *strictly positive* matrices, better rates are known. Kalantari and Khachiyan [16] showed that for positive matrices and the $(\mathbf{1}, \mathbf{1})$ -scaling problem, the Sinkhorn-Knopp algorithm obtains ℓ_2 error $\leq \varepsilon$ in $O(\sqrt{n} \ln(1/\nu)/\varepsilon)$ -iterations; this result was extended to the general matrix scaling problem by Kalantari et al [19]. In a different track, Franklin and Lorenz [13] show that in fact the dependence on ε can be made logarithmic, and thus the algorithm has “linear convergence”, however their analysis² has a polynomial dependence of $(1/\nu)$. All these results use the positivity crucially and seem to break down even with one 0 entry.

The Sinkhorn-Knopp algorithm has polynomial dependence on the error parameter and therefore is a “pseudopolynomial” time approximation. We conclude by briefly describing bounds obtained by other algorithms for the matrix scaling problem whose dependence on ε is logarithmic rather than polynomial. Kalantari and Khachiyan [17] describe a method based on the ellipsoid algorithm which runs in time $O(n^4 \ln(n/\varepsilon) \ln(1/\nu))$. Nemirovskii and Rothblum [25] describe a method with running time $O(n^4 \ln(n/\varepsilon) \ln \ln(1/\nu))$. The first strongly polynomial time approximation scheme (with no dependence on ν) was due to Linal, Samoridnitsky, and Wigderson [20] who gave a $\tilde{O}(n^7 \ln(h/\varepsilon))$ time algorithm. Rote and Zachariasen [28] reduced the matrix scaling problem to flow problems to give a $O(n^4 \ln(h/\varepsilon))$ time algorithms for the matrix scaling problem. To compare, we should recall that Theorem 1 shows that our algorithm runs in time $O(\text{nnz}(A)h^2/\varepsilon^2)$ time.

Very recently, two independent works obtain vastly improved running times for matrix scaling. Cohen et al [9] give $\tilde{O}(\text{nnz}(A)^{3/2})$ time algorithm, while Allen-Zhu et al [2] give a $\tilde{O}(n^{7/3} + \text{nnz}(A) \cdot (n + n^{1/3}h^{1/2}))$ time algorithm; the tildes in both the above running times hide the logarithmic dependence on ε and ν . Both these algorithms look at the matrix scaling problem as a convex optimization problem and perform second order methods.

2 Entropy Minimization Viewpoint of the Sinkhorn-Knopp Algorithm

There have been many approaches (see Idel [15], Section 3 for a discussion) towards analyzing the Sinkhorn-Knopp algorithm including convex optimization and log-barrier methods [16, 19, 22, 5], non-linear Perron-Frobenius theory [24, 35, 13, 8, 16], topological methods [27, 6], connections to the permanent [20, 18], and the entropy minimization method [7, 10, 11, 14] which is what we use for our analysis.

² [13] never make the base of the logarithm explicit, but their proof shows it can be as large as $1 - 1/\nu^2$.

We briefly describe the entropy minimization viewpoint. Given two non-negative matrices M and N let us define the *Kullback-Leibler* divergence³ between M and N as follows

$$\mathbf{D}(M, N) := \frac{1}{h} \sum_{1 \leq i \leq n} \sum_{1 \leq j \leq m} M_{ij} \ln \left(\frac{M_{ij}}{N_{ij}} \right) \quad (1)$$

with the convention that the summand is zero if both M_{ij} and N_{ij} are 0, and is ∞ if $M_{ij} > 0$ and $N_{ij} = 0$. Let Φ_r be the set of $n \times m$ matrices whose row-sums are \mathbf{r} and let Φ_c be the set of $n \times m$ matrices whose column sums are \mathbf{c} . Given matrix A suppose we wish to find the matrix $A^* = \arg \min_{B \in \Phi_r \cap \Phi_c} \mathbf{D}(B, A)$. One algorithm for this is to use the method of alternate projections with respect to the KL-divergence [7] (also known as I -projections [10]) which alternately finds the matrices in Φ_r and Φ_c closest in the KL-divergence sense to the current matrix at hand, and then sets the minimizer to be the current matrix. It is not too hard to see (see Idel [15], Observation 3.17 for a proof) that the above alternate projection algorithm is precisely the Sinkhorn-Knopp algorithm. Therefore, at least in this sense, the right metric to measure the distance to optimality is not the ℓ_1 or the ℓ_2 error as described in the previous section, but the rather the KL-divergence between the normalized vectors as described below.

Let $\pi_r^{(t)} := \mathbf{r}^{(t)}/h$ be the n -dimensional probability vector whose i th entry is $\mathbf{r}_i^{(t)}/h$; similarly define the m -dimensional vector $\pi_c^{(t)}$. Let π_r denote the n -dimensional probability vector with the i th entry being \mathbf{r}_i/h ; similarly define π_c . Recall that the KL-divergence between two probability distributions p, q is defined as $\mathbf{D}_{KL}(p||q) := \sum_{i=1}^n p_i \ln(q_i/p_i)$. The following theorem gives the convergence time for the KL-divergence.

► **Theorem 5.** *If the matrix $A \in \mathbb{R}_{\geq 0}^{n \times m}$ is (\mathbf{r}, \mathbf{c}) -scalable, then for any $\delta > 0$ there is a $t \leq T = \lceil \left(\frac{\ln(1+2n\rho/\nu)}{\delta} \right) \rceil$ with either $\mathbf{D}_{KL}(\pi_r || \pi_r^{(t)}) \leq \delta$ or $\mathbf{D}_{KL}(\pi_c || \pi_c^{(t)}) \leq \delta$.*

Proof. Let $Z := RAS$ be a matrix with row-sums \mathbf{r} and column-sums \mathbf{c} for diagonal matrices R, S . Recall A^0 is the matrix obtained by column-scaling A . Note that the minimum non-zero entry of A^0 is $\geq \nu/n$.

► **Lemma 6.** $\mathbf{D}(Z, A^0) \leq \ln(1 + 2n\rho/\nu)$ and $\mathbf{D}(Z, A^t) \geq 0$ for all t .

Proof. By definition,

$$\mathbf{D}(Z, A^{(t)}) = \frac{1}{h} \sum_{j=1}^m \sum_{i=1}^n Z_{ij} \ln \left(\frac{Z_{ij}}{A_{ij}^{(t)}} \right) = \frac{1}{h} \sum_{j=1}^m \mathbf{c}_j \sum_{i=1}^n \frac{Z_{ij}}{\mathbf{c}_j} \ln \left(\frac{Z_{ij}}{A_{ij}^{(t)}} \right)$$

For a fixed j , the vectors $\left(\frac{Z_{1j}}{\mathbf{c}_j}, \frac{Z_{2j}}{\mathbf{c}_j}, \dots, \frac{Z_{nj}}{\mathbf{c}_j} \right)$ and $\left(\frac{A_{1j}^{(t)}}{\mathbf{c}_j}, \frac{A_{2j}^{(t)}}{\mathbf{c}_j}, \dots, \frac{A_{nj}^{(t)}}{\mathbf{c}_j} \right)$ are probability vectors, and therefore the above is a sum of \mathbf{c}_j -weighted KL-divergences which is always non-negative. For the upper bound, one can use the fact (Inequality 27, [31]) that for any two distributions p and q , $D(p||q) \leq \ln(1 + \frac{\|p-q\|_2^2}{q_{\min}}) \leq \ln(1 + \frac{2}{q_{\min}})$ where q_{\min} is the smallest non-zero entry of q . For our purpose, we note that the minimum non-zero probability of the $A_j^{(0)}$ distribution being $\geq \nu/n\rho$. Therefore, the second summand is at most $\ln(1 + 2n\rho/\nu)$ giving us $D(Z, A^{(0)}) \leq \frac{1}{h} \sum_{j=1}^m \mathbf{c}_j \cdot \ln(1 + 2n\rho/\nu) = \ln(1 + 2n\rho/\nu)$. ◀

³ The KL-divergence is normally stated between two distributions and doesn't have the $1/h$ factor. Also the logarithms are usually base 2.

► **Lemma 7.**

$$\mathbf{D}(Z, A^{(t)}) - \mathbf{D}(Z, B^{(t)}) = \mathbf{D}_{KL}(\pi_{\mathbf{r}} || \pi_{\mathbf{r}}^{(t)}) \quad \text{and} \quad \mathbf{D}(Z, B^{(t)}) - \mathbf{D}(Z, A^{(t+1)}) = \mathbf{D}_{KL}(\pi_{\mathbf{c}} || \pi_{\mathbf{c}}^{(t)})$$

Proof. The LHS of the first equality is simply

$$\begin{aligned} \frac{1}{h} \sum_{j=1}^m \sum_{i=1}^n Z_{ij} \ln \left(\frac{B_{ij}^{(t)}}{A_{ij}^{(t)}} \right) &= \frac{1}{h} \sum_{j=1}^m \sum_{i=1}^n Z_{ij} \ln \left(\frac{\mathbf{r}_i}{\mathbf{r}_i^{(t)}} \right) \\ &= \frac{1}{h} \sum_{i=1}^n \ln \left(\frac{\mathbf{r}_i}{\mathbf{r}_i^{(t)}} \right) \sum_{j=1}^m Z_{ij} \\ &= \sum_{i=1}^n \left(\frac{\mathbf{r}_i}{h} \right) \cdot \ln \left(\frac{\mathbf{r}_i/h}{\mathbf{r}_i^{(t)}/h} \right) \end{aligned}$$

since $\sum_{j=1}^m Z_{ij} = \mathbf{r}_i$. The last summand is precisely $\mathbf{D}_{KL}(\pi_{\mathbf{r}} || \pi_{\mathbf{r}}^{(t)})$. The other equation follows analogously. ◀

The above two lemmas easily imply the theorem. If for all $0 \leq t \leq T$, both $\mathbf{D}_{KL}(\pi_{\mathbf{r}} || \pi_{\mathbf{r}}^{(t)}) > \delta$ and $\mathbf{D}_{KL}(\pi_{\mathbf{c}} || \pi_{\mathbf{c}}^{(t)}) > \delta$, then substituting in Lemma 7 and summing we get $\mathbf{D}(Z, A^{(0)}) - \mathbf{D}(Z, A^{(T+1)}) > T\delta > \ln(1 + 2n\rho/\nu)$ contradicting Lemma 6. ◀

Theorem 1 follows from Theorem 5 using connections between the KL-divergence and the ℓ_1 and ℓ_2 norms. One is the following famous Pinsker's inequality which allows us to easily prove part 1 of Theorem 1. Given any two probability distributions p, q ,

$$\mathbf{D}_{KL}(p || q) \geq \frac{1}{2} \cdot \|p - q\|_1^2 \quad (\text{Pinsker})$$

Proof of Theorem 1, Part 1. Apply (Pinsker) on the vectors $\pi_{\mathbf{r}}$ and $\pi_{\mathbf{r}}^{(t)}$ to get

$$\mathbf{D}_{KL}(\pi_{\mathbf{r}} || \pi_{\mathbf{r}}^{(t)}) \geq \frac{1}{2h^2} \|\mathbf{r}^{(t)} - \mathbf{r}\|_1^2$$

Set $\delta := \frac{\varepsilon^2}{2h^2}$ and apply Theorem 5. In $O\left(\frac{h^2 \ln(n\rho/\nu)}{\varepsilon^2}\right)$ time we would get a matrix with $\delta > \mathbf{D}_{KL}(\pi_{\mathbf{r}} || \pi_{\mathbf{r}}^{(t)})$ which from the above inequality would imply $\|\mathbf{r}^{(t)} - \mathbf{r}\|_1 \leq \varepsilon$. ◀

To prove Part 2, we need a way to relate the ℓ_2 norm and the KL-divergence. In order to do so, we prove a different lower bound which implies Pinsker's inequality (with a worse constant), but is significantly stronger in certain regimes. To the best of our knowledge this is a new bound which may be of independent interest in other domains. Below we state the version which we need for the proof of Theorem 1, part 2. This is an instantiation of the general inequality Lemma 9 which we prove in Section 3.

► **Lemma 8.** *Given any pair of probability distributions p, q over a finite domain, define $\mathcal{A} := \{i : q_i > 2p_i\}$ and $\mathcal{B} := \{i : q_i \leq 2p_i\}$. Then,*

$$\mathbf{D}_{KL}(p || q) \geq (1 - \ln 2) \cdot \left(\sum_{i \in \mathcal{A}} |q_i - p_i| + \sum_{i \in \mathcal{B}} \frac{(q_i - p_i)^2}{p_i} \right) \quad (\text{KL vs } \ell_1/\ell_2)$$

Proof of Theorem 1, Part 2. We apply Lemma 8 on the vectors $\pi_{\mathbf{r}}$ and $\pi_{\mathbf{r}}^{(t)}$. Lemma 8 gives us

$$\begin{aligned} \mathbf{D}_{KL}(\pi_{\mathbf{r}}|\pi_{\mathbf{r}}^{(t)}) &\geq C \cdot \left(\frac{1}{h} \sum_{i \in A} |\mathbf{r}_i^{(t)} - \mathbf{r}_i| + \frac{1}{h} \sum_{i \in B} \frac{(\mathbf{r}_i^{(t)} - \mathbf{r}_i)^2}{\mathbf{r}_i} \right) \\ &\geq \frac{C}{h} \left(\sum_{i \in A} |\mathbf{r}_i^{(t)} - \mathbf{r}_i| + \frac{1}{\rho} \sum_{i \in B} (\mathbf{r}_i^{(t)} - \mathbf{r}_i)^2 \right) \end{aligned}$$

where $C = 1 - \ln 2$. If the second summand in the parenthesis of the RHS is $\geq \frac{1}{2} \|\mathbf{r}^{(t)} - \mathbf{r}\|_2^2$, then we get $\mathbf{D}_{KL}(\pi_{\mathbf{r}}|\pi_{\mathbf{r}}^{(t)}) \geq \frac{C}{2\rho h} \|\mathbf{r}^{(t)} - \mathbf{r}\|_2^2$. Otherwise, we have $\mathbf{D}_{KL}(\pi_{\mathbf{r}}|\pi_{\mathbf{r}}^{(t)}) \geq \frac{C}{\sqrt{2}h} \|\mathbf{r}^{(t)} - \mathbf{r}\|_2$, where we used the weak fact that the sum of some positive numbers is at least the square-root of the sum of their squares. In any case, we get the following

$$\mathbf{D}_{KL}(\pi_{\mathbf{r}}|\pi_{\mathbf{r}}^{(t)}) \geq \min \left(\frac{C}{2\rho h} \|\mathbf{r}^{(t)} - \mathbf{r}\|_2^2, \frac{C}{\sqrt{2}h} \|\mathbf{r}^{(t)} - \mathbf{r}\|_2 \right) \quad (2)$$

To complete the proof of part 2 of Theorem 1, set $\delta := \frac{C}{2\rho h(\frac{1}{\varepsilon} + \frac{1}{\varepsilon^2})}$ and apply Theorem 5. In $O(\rho h \ln(n\rho/\nu) \cdot (\frac{1}{\varepsilon} + \frac{1}{\varepsilon^2}))$ time we would get a matrix with $\delta \geq \mathbf{D}_{KL}(\pi_{\mathbf{r}}|\pi_{\mathbf{r}}^{(t)})$. If the minimum of the RHS of (2) is the first term, then we get $\|\mathbf{r}^{(t)} - \mathbf{r}\|_2^2 \leq \varepsilon^2$ implying the ℓ_2 -error is $\leq \varepsilon$. If the minimum is the second term, then we get $\|\mathbf{r}^{(t)} - \mathbf{r}\|_2 \leq \frac{\varepsilon}{\sqrt{2}\rho} < \varepsilon$ since $\rho \geq 1$. \blacktriangleleft

3 New Lower Bound on the KL-Divergence

We now establish a new lower bound on KL-divergence which yields (KL vs ℓ_1/ℓ_2) as a corollary.

► **Lemma 9.** *Let p and q be two distributions over a finite n -element universe. For any fixed $\theta > 0$, define the sets $\mathcal{A}_\theta := \{i \in [n] : q_i \geq (1 + \theta)p_i\}$ and $\mathcal{B}_\theta = [n] \setminus \mathcal{A}_\theta = \{i \in [n] : q_i \leq (1 + \theta)p_i\}$. Then we have the following inequality*

$$\mathbf{D}_{KL}(p||q) \geq \left(1 - \frac{\ln(1 + \theta)}{\theta}\right) \cdot \left(\sum_{i \in \mathcal{A}_\theta} |q_i - p_i| + \frac{1}{\theta} \sum_{i \in \mathcal{B}_\theta} p_i \left(\frac{q_i - p_i}{p_i}\right)^2\right) \quad (3)$$

When $\theta = 1$, we get (KL vs ℓ_1/ℓ_2).

A Comparison of (Pinsker) and (KL vs ℓ_1/ℓ_2): To see why (KL vs ℓ_1/ℓ_2) generalizes (Pinsker) with a weaker constant, note that

$$\|p - q\|_1^2 = \left(\sum_{i \in \mathcal{A}} |q_i - p_i| + \sum_{i \in \mathcal{B}} |q_i - p_i|\right)^2 \leq 2 \left(\sum_{i \in \mathcal{A}} |q_i - p_i|\right)^2 + 2 \left(\sum_{i \in \mathcal{B}} p_i \frac{|q_i - p_i|}{p_i}\right)^2$$

The first parenthetical term above, since it is ≤ 1 , is at most the first summation in the parenthesis of (KL vs ℓ_1/ℓ_2). The second parenthetical term above, by Cauchy-Schwarz, is at most the second summation in the parenthesis of (KL vs ℓ_1/ℓ_2). Thus (KL vs ℓ_1/ℓ_2) implies $\mathbf{D}_{KL}(p||q) \geq \frac{(1 - \ln 2)}{2} \|p - q\|_1^2$. On the other hand, the RHS of (KL vs ℓ_1/ℓ_2) can be much larger than that of (Pinsker). For instance, suppose $p_i = 1/n$ for all i , $q_1 = 1/n + 1/\sqrt{n}$, and for $i \neq 1$, $q_i = 1/n - \frac{1}{(n-1)\sqrt{n}}$. The RHS of (Pinsker) is $\Theta(1/n)$ while that of (KL vs ℓ_1/ℓ_2) is $\Theta(1/\sqrt{n})$ which is the correct order of magnitude for $\mathbf{D}_{KL}(p||q)$.

Proof of Lemma 9: We need the following fact which follows from calculus; we provide a proof later for completeness.

► **Lemma 10.** *Given any $\theta > 0$, define $a_\theta := \frac{\ln(1+\theta)}{\theta}$ and $b_\theta := \frac{1}{\theta} \left(1 - \frac{\ln(1+\theta)}{\theta}\right)$. Then,*

- *For $t \geq \theta$, $(1+t) \leq e^{a_\theta t}$*
- *For $t \leq \theta$, $(1+t) \leq e^{t-b_\theta t^2}$*

Define $\eta_i := \frac{q_i - p_i}{p_i}$. Note that $\mathcal{A}_\theta = \{i : \eta_i > \theta\}$ and \mathcal{B}_θ is the rest. We can write the KL-divergence as follows

$$\mathbf{D}_{KL}(p||q) := \sum_{i=1}^n p_i \ln(p_i/q_i) = - \sum_{i=1}^n p_i \ln(1 + \eta_i)$$

For $i \in \mathcal{A}_\theta$, since $\eta_i > \theta$, we upper bound $(1 + \eta_i) \leq e^{a_\theta \eta_i}$ using Fact 10. For $i \in \mathcal{B}_\theta$, that is $\eta_i \leq \theta$, we upper bound $(1 + \eta_i) \leq e^{\eta_i - b_\theta \eta_i^2}$ using Fact 10. Lastly, we note $\sum_i p_i \eta_i = 0$ since p, q both sum to 1, implying $\sum_{i \in \mathcal{B}_\theta} p_i \eta_i = - \sum_{i \in \mathcal{A}_\theta} p_i \eta_i$. Putting all this in the definition above we get

$$\mathbf{D}_{KL}(p||q) \geq -a_\theta \cdot \sum_{i \in \mathcal{A}_\theta} p_i \eta_i - \sum_{i \in \mathcal{B}_\theta} p_i \eta_i + b_\theta \sum_{i \in \mathcal{B}_\theta} p_i \eta_i^2 = (1 - a_\theta) \sum_{i \in \mathcal{A}_\theta} p_i \eta_i + b_\theta \sum_{i \in \mathcal{B}_\theta} p_i \eta_i^2$$

The proof of inequality (3) follows by noting that $b_\theta = \frac{1-a_\theta}{\theta}$. ◀

Proof of Lemma 10. The proof of both facts follow by proving non-negativity of the relevant function in the relevant interval. Recall $a_\theta = \ln(1+\theta)/\theta$ and $b_\theta = \frac{1}{\theta}(1 - a_\theta)$. We start with the following three inequalities about the log-function.

$$\text{For all } z > 0, \quad z + z^2/2 > (1+z) \ln(1+z) > z \quad \text{and} \quad \ln(1+z) > z - z^2/2 \quad (4)$$

The third inequality in (4) implies $a_\theta > 1 - \theta/2$ and thus, $b_\theta < 1/2$. The first inequality in (4) implies $a_\theta < \frac{1+\frac{\theta}{2}}{1+\theta}$ which in turn implies $b_\theta > 1/2(1+\theta)$. For brevity, henceforth let us lose the subscript on a_θ and b_θ .

Consider the function $f(t) = e^{at} - (1+t)$. Note that $f'(t) = ae^{at} - 1$ which is increasing in t since $a > 0$. So, for any $t \geq \theta$, we have $f'(t) \geq ae^{a\theta} - 1 = \frac{(1+\theta) \ln(1+\theta)}{\theta} - 1 \geq 0$, by the second inequality in (4). Therefore, f is increasing when $t \geq \theta$. The first part of Fact 10 follows since $f(\theta) = 0$ by definition of a .

Consider the function $g(t) = e^{t(1-bt)} - (1+t)$. Note that $g(0) = g(\theta) = 0$. We break the argument in two parts: we argue that $g(t)$ is strictly positive for all $t \leq 0$, and that $g(t)$ is strictly positive for $t \in (0, \theta)$. This will prove the second part of Fact 10.

The first derivative is $g'(t) = (1 - 2bt)e^{t(1-bt)} - 1$ and the second derivative is $g''(t) = e^{t(1-bt)} ((1 - 2bt)^2 - 2b)$. Since $b < 1/2$, we have $2b < 1$, and thus for $t \leq 0$, $g''(t) > 0$. Therefore, g' is strictly increasing for $t \leq 0$. However, $g'(0) = 0$, and so $g'(t) < 0$ for all $t < 0$. This implies g is strictly decreasing in the interval $t < 0$. Noting $g(0) = 0$, we get $g(t) > 0$ for all $t < 0$. This completes the first part of the argument.

For the second part, we first note that $g'(\theta) < 0$ since $b > \frac{1}{2(1+\theta)}$. That is, g is strictly decreasing at θ . On the other hand g is increasing at θ . To see this, looking at g' is not enough since $g'(0) = 0$. However, $g''(0) > 0$ since $b < 1/2$. This means that 0 is a strict (local) minimum for g implying g is increasing at 0. In sum, g vanishes at 0 and θ , and is increasing at 0 and decreasing at θ . This means that if g does vanish at some $r \in (0, \theta)$, then it must vanish once again in $[r, \theta)$ for it to be decreasing at θ . In particular, g' must vanish three times in $(0, \theta)$ and thus four times in $[0, \theta)$ since $g'(0) = 0$. This in turn

implies g'' vanishes three times in $[0, \theta)$ which is a contradiction since g'' is a quadratic in t multiplied by a positive term.

We end by proving (4). This also follows the same general methodology. Define $p(z) := (1+z)\ln(1+z) - z$ and $q(z) := p(z) - z^2/2$. Differentiating, we get $p'(z) = \ln(1+z) > 0$ for all $z > 0$, and $q'(z) = \ln(1+z) - z < 0$ for all $z > 0$. Thus, p is increasing, and q is decreasing, in $(0, \infty)$. The first two inequalities of (4) follow since $p(0) = q(0) = 0$. To see the third inequality, define $r(z) = \ln(1+z) - z + z^2/2$ and observe $r'(z) = \frac{1}{1+z} - 1 + z = \frac{z^2}{1+z}$ which is > 0 if $z > 0$. Thus r is strictly increasing, and the third inequality of (4) follows since $r(0) = 0$. ◀

References

- 1 Scott Aaronson. Quantum computing and hidden variables. *Phys. Rev. A*, 71:032325, Mar 2005. doi:10.1103/PhysRevA.71.032325.
- 2 Zeyuan Allen Zhu, Yuanzhi Li, Rafael Oliveira, and Avi Wigderson. Much faster algorithms for matrix scaling. *58th IEEE Annual Symposium on Foundations of Computer Science, FOCS*, 2017. URL: <http://arxiv.org/abs/1704.02315>.
- 3 Zeyuan Allen Zhu and Lorenzo Orecchia. Using optimization to break the epsilon barrier: A faster and simpler width-independent algorithm for solving positive linear programs in parallel. In *Proceedings of the Twenty-Sixth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA, San Diego, CA, USA*, pages 1439–1456, 2015.
- 4 Michael Bacharach. Estimating nonnegative matrices from marginal data. *International Economic Review*, 6(3):294–310, 1965. URL: <http://www.jstor.org/stable/2525582>.
- 5 H. Balakrishnan, Inseok Hwang, and C. J. Tomlin. Polynomial approximation algorithms for belief matrix maintenance in identity management. In *2004 43rd IEEE Conference on Decision and Control (CDC) (IEEE Cat. No.04CH37601)*, volume 5, pages 4874–4879 Vol.5, Dec 2004. doi:10.1109/CDC.2004.1429569.
- 6 R.B. Bapat and T.E.S. Raghavan. An extension of a theorem of Darroch and Ratcliff in loglinear models and its application to scaling multidimensional matrices. *Linear Algebra and its Applications*, 114:705–715, 1989. Special Issue Dedicated to Alan J. Hoffman. doi:10.1016/0024-3795(89)90489-8.
- 7 L.M. Bregman. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Mathematical Physics*, 7(3):200–217, 1967. doi:10.1016/0041-5553(67)90040-7.
- 8 Richard A Brualdi, Seymour V Parter, and Hans Schneider. The diagonal equivalence of a nonnegative matrix to a stochastic matrix. *Journal of Mathematical Analysis and Applications*, 16(1):31–50, 1966. doi:10.1016/0022-247X(66)90184-3.
- 9 Michael B. Cohen, Aleksander Madry, Dimitris Tsipras, and Adrian Vladu. Matrix scaling and balancing via box constrained newton’s method and interior point methods. *58th IEEE Annual Symposium on Foundations of Computer Science, FOCS*, 2017. URL: <http://arxiv.org/abs/1704.02310>.
- 10 Imre Csiszar. I-divergence geometry of probability distributions and minimization problems. *Ann. Probab.*, 3(1):146–158, 02 1975. doi:10.1214/aop/1176996454.
- 11 Imre Csiszar. A geometric interpretation of darroch and ratcliff’s generalized iterative scaling. *Ann. Statist.*, 17(3):1409–1413, 09 1989. doi:10.1214/aos/1176347279.
- 12 W. Edwards Deming and Frederick F. Stephan. On a least squares adjustment of a sampled frequency table when the expected marginal totals are known. *Ann. Math. Statist.*, 11(4):427–444, 12 1940. doi:10.1214/aoms/1177731829.

- 13 Joel Franklin and Jens Lorenz. On the scaling of multidimensional matrices. *Linear Algebra and its Applications*, 114:717–735, 1989. Special Issue Dedicated to Alan J. Hoffman. doi:10.1016/0024-3795(89)90490-4.
- 14 Leonid Gurvits and Peter N. Yianilos. The deflation-inflation method for certain semidefinite programming and maximum determinant completion problems. Technical report, NEC Research Institute, 4 Independence Way, Princeton, NJ 08540, 1998.
- 15 Martin Idel. A review of matrix scaling and Sinkhorn’s normal form for matrices and positive maps. *ArXiv e-prints*, 2016. arXiv:1609.06349.
- 16 Bahman Kalantari and Leonid Khachiyan. On the rate of convergence of deterministic and randomized RAS matrix scaling algorithms. *Oper. Res. Lett.*, 14(5):237–244, 1993. doi:10.1016/0167-6377(93)90087-W.
- 17 Bahman Kalantari and Leonid Khachiyan. On the complexity of nonnegative-matrix scaling. *Linear Algebra and its Applications*, 240:87–103, 1996. doi:10.1016/0024-3795(94)00188-X.
- 18 Bahman Kalantari, Isabella Lari, Federica Ricca, and Bruno Simeone. On the complexity of general matrix scaling and entropy minimization via the RAS algorithm. *Technical Report, n.24, Department of Statistics and Applied probability, La Sapienza University, Rome*, 2002.
- 19 Bahman Kalantari, Isabella Lari, Federica Ricca, and Bruno Simeone. On the complexity of general matrix scaling and entropy minimization via the RAS algorithm. *Math. Program.*, 112(2):371–401, 2008. doi:10.1007/s10107-006-0021-4.
- 20 Nathan Linial, Alex Samorodnitsky, and Avi Wigderson. A deterministic strongly polynomial algorithm for matrix scaling and approximate permanents. *Combinatorica*, 20(4):545–568, 2000. doi:10.1007/s004930070007.
- 21 Michael Luby and Noam Nisan. A parallel approximation algorithm for positive linear programming. In S. Rao Kosaraju, David S. Johnson, and Alok Aggarwal, editors, *Proceedings of the Twenty-Fifth Annual ACM Symposium on Theory of Computing, May 16-18, 1993, San Diego, CA, USA*, pages 448–457. ACM, 1993. doi:10.1145/167088.167211.
- 22 Sally M Macgill. Theoretical properties of biproportional matrix adjustments. *Environment and Planning A*, 9(6):687–701, 1977. doi:10.1068/a090687.
- 23 Michael W. Mahoney, Satish Rao, Di Wang, and Peng Zhang. Approximating the Solution to Mixed Packing and Covering LPs in Parallel $O(\varepsilon^{-3})$ Time. In *43rd International Colloquium on Automata, Languages, and Programming, ICALP 2016, July 11-15, 2016, Rome, Italy*, pages 52:1–52:14, 2016.
- 24 MV Menon. Reduction of a matrix with positive elements to a doubly stochastic matrix. *Proceedings of the American Mathematical Society*, 18(2):244–247, 1967.
- 25 Arkadi Nemirovskii and Uriel Rothblum. On complexity of matrix scaling. *Linear Algebra and its Applications.*, 302:435–460, 1999.
- 26 Juan de Dios Ortúzar and Luis G. Willumsen. *Modelling Transport*. John Wiley & Sons, Ltd, 2011. doi:10.1002/9781119993308.fmatter.
- 27 T.E.S. Raghavan. On pairs of multidimensional matrices. *Linear Algebra and its Applications*, 62:263–268, 1984. doi:10.1016/0024-3795(84)90101-0.
- 28 Günter Rote and Martin Zachariasen. Matrix scaling by network flow. In *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA ’07*, pages 848–854, Philadelphia, PA, USA, 2007. Society for Industrial and Applied Mathematics. URL: <http://dl.acm.org/citation.cfm?id=1283383.1283474>.
- 29 Uriel Rothblum and Hans Schneider. Scalings of matrices which have prespecified row sums and column sums via optimization. *Linear Algebra and its Applications*, 114:737–764, 1989. Special Issue Dedicated to Alan J. Hoffman.
- 30 Ludger Ruschendorf. Convergence of the iterative proportional fitting procedure. *Ann. Statist.*, 23(4):1160–1174, 1995. doi:10.1214/aos/1176324703.

- 31 Igal Sason and Sergio Verdú. Upper bounds on the relative entropy and rényi divergence as a function of total variation distance for finite alphabets. In *2015 IEEE Information Theory Workshop - Fall (ITW), Jeju Island, South Korea, October 11-15, 2015*, pages 214–218. IEEE, 2015. doi:10.1109/ITWF.2015.7360766.
- 32 Erwin Schrödinger. Über die umkehrung der naturgesetze. *Preuss. Akad. Wiss., Phys.-Math. Kl.*, 1931.
- 33 Richard Sinkhorn. Diagonal equivalence to matrices with prescribed row and column sums. *The American Mathematical Monthly*, 74(4):402–405, 1967. URL: <http://www.jstor.org/stable/2314570>.
- 34 Richard Sinkhorn and Paul Knopp. Concerning nonnegative matrices and doubly stochastic matrices. *Pacific J. Math.*, 21(2):343–348, 1967. URL: <https://projecteuclid.org:443/euclid.pjm/1102992505>.
- 35 George W. Soules. The rate of convergence of sinkhorn balancing. *Linear Algebra and its Applications*, 150:3–40, 1991. doi:10.1016/0024-3795(91)90157-R.
- 36 Neal E. Young. Sequential and parallel algorithms for mixed packing and covering. In *42nd Annual Symposium on Foundations of Computer Science, FOCS, Las Vegas, Nevada, USA*, pages 538–546, 2001.