

Computationally Efficient Estimators for Dimension Reductions Using Stable Random Projections

Ping Li

Department of Statistical Science, Cornell University, Ithaca NY 14853, USA
 {pingli}@cornell.edu

Abstract. The¹ method of *stable random projections* is a tool for efficiently computing the l_α distances using low memory, where $0 < \alpha \leq 2$ is a tuning parameter. The method boils down to a statistical estimation task and various estimators have been proposed, based on the *geometric mean*, the *harmonic mean*, and the *fractional power* etc.

This study proposes the ***optimal quantile*** estimator, whose main operation is ***selecting***, which is considerably less expensive than taking fractional power, the main operation in previous estimators. Our experiments report that the *optimal quantile* estimator is nearly one order of magnitude more computationally efficient than previous estimators. For large-scale learning tasks in which storing and computing pairwise distances is a serious bottleneck, this estimator should be desirable.

In addition to its computational advantages, the *optimal quantile* estimator exhibits nice theoretical properties. It is more accurate than previous estimators when $\alpha > 1$. We derive its theoretical error bounds and establish the explicit (i.e., no hidden constants) sample complexity bound.

1 Introduction

The method of *stable random projections*[1,2,3], as an efficient tool for computing pairwise distances in massive high-dimensional data, provides a promising mechanism to tackle some of the challenges in modern machine learning. In this paper, we provide an easy-to-implement algorithm for *stable random projections* which is both statistically accurate and computationally efficient.

1.1 Massive High-dimensional Data in Modern Machine Learning

We denote a data matrix by $\mathbf{A} \in \mathbb{R}^{n \times D}$, i.e., n data points in D dimensions. Data sets in modern applications exhibit important characteristics which impose tremendous challenges in machine learning [4]:

- Modern data sets with $n = 10^5$ or even $n = 10^6$ points are not uncommon in supervised learning, e.g., in image/text classification, ranking algorithms for search engines, etc. In the unsupervised domain (e.g., Web clustering, ads clickthroughs, word/term associations), n can be even much larger.
- Modern data sets are often of ultra high-dimensions (D), sometimes in the order of millions (or even higher), e.g., image, text, genome (e.g., SNP), etc. For example, in image analysis, D may be $10^3 \times 10^3 = 10^6$ if using pixels as features, or $D = 256^3 \approx 16$ million if using color histograms as features.
- Modern data sets are sometimes collected in a dynamic streaming fashion.

¹ First draft Feb. 2008, slightly revised in June 2008. The results were announced in January 2008 at SODA'08 when the author presented the work of [2].

- Large-scale data are often heavy-tailed, e.g., image and text data.

Some large-scale data are dense, such as image and genome data. Even for data sets which are sparse, such as text, the absolute number of non-zeros may be still large. For example, if one queries “machine learning” (a not-too-common term) in Google.com, the total number of pagehits is about 3 million. In other words, if one builds a term-doc matrix at Web scale, although the matrix is sparse, most rows will contain large numbers (e.g., millions) of non-zero entries.

1.2 Pairwise Distances in Machine Learning

Many learning algorithms require a similarity matrix computed from pairwise distances of the data matrix $\mathbf{A} \in \mathbb{R}^{n \times D}$. Examples include clustering, nearest neighbors, multidimensional scaling, and kernel SVM (support vector machines). The similarity matrix requires $O(n^2)$ storage space and $O(n^2 D)$ computing time.

This study focuses on the l_α distance ($0 < \alpha \leq 2$). Consider two vectors $u_1, u_2 \in \mathbb{R}^D$ (e.g., the leading two rows in \mathbf{A}), the l_α distance between u_1 and u_2 is

$$d_{(\alpha)} = \sum_{i=1}^D |u_{1,i} - u_{2,i}|^\alpha. \quad (1)$$

Note that, strictly speaking, the l_α distance should be defined as $d_{(\alpha)}^{1/\alpha}$. Because the power operation $(\cdot)^{1/\alpha}$ is the same for all pairs, it often makes no difference whether we use $d_{(\alpha)}^{1/\alpha}$ or just $d_{(\alpha)}$; and hence we focus on $d_{(\alpha)}$.

The radial basis kernel (e.g., for SVM) is constructed from $d_{(\alpha)}$ [5,6]:

$$\mathbf{K}(u_1, u_2) = \exp \left(-\gamma \sum_{i=1}^D |u_{1,i} - u_{2,i}|^\alpha \right), \quad 0 < \alpha \leq 2. \quad (2)$$

When $\alpha = 2$, this is the Gaussian radial basis kernel. Here α can be viewed as a *tuning* parameter. For example, in their histogram-based image classification project using SVM, [5] reported that $\alpha = 0$ and $\alpha = 0.5$ achieved good performance. For heavy-tailed data, tuning α has the similar effect as term-weighting the original data, often a critical step in a lot of applications [7,8].

For popular kernel SVM solvers including the *Sequential Minimal Optimization* (SMO) algorithm[9], storing and computing kernels is the major bottleneck. Three computational challenges were summarized in [4, page 12]:

- **Computing kernels is expensive**
- **Computing full kernel matrix is wasteful** Efficient SVM solvers often do not need to evaluate all pairwise kernels.
- **Kernel matrix does not fit in memory** Storing the kernel matrix at the memory cost $O(n^2)$ is challenging when $n > 10^5$, and is not realistic for $n > 10^6$, because $O(10^{12})$ consumes at least 1000 GBs memory.

A popular strategy in large-scale learning is to evaluate distances **on the fly**[4]. That is, instead of loading the similarity matrix in memory at the cost of $O(n^2)$, one can load the original data matrix at the cost of $O(nD)$ and recompute pairwise distances on-demand. This strategy is apparently problematic when D

is not too small. For high-dimensional data, either loading the data matrix in memory is unrealistic or computing distances on-demand becomes too expensive.

Those challenges are not unique to kernel SVM; they are general issues in distanced-based learning algorithms. The method of *stable random projections* provides a promising scheme by reducing the dimension D to a small k (e.g., $k = 50$), to facilitate compact data storage and efficient distance computations.

1.3 Stable Random Projections

The basic procedure of *stable random projections* is to multiply $\mathbf{A} \in \mathbb{R}^{n \times D}$ by a random matrix $\mathbf{R} \in \mathbb{R}^{D \times k}$ ($k \ll D$), which is generated by sampling each entry r_{ij} i.i.d. from a symmetric stable distribution $S(\alpha, 1)$. The resultant matrix $\mathbf{B} = \mathbf{A} \times \mathbf{R} \in \mathbb{R}^{n \times k}$ is much smaller than \mathbf{A} and hence it may fit in memory.

Suppose a stable random variable $x \sim S(\alpha, d)$, where d is the scale parameter. Then its characteristic function (Fourier transform of the density function) is

$$\mathbb{E}(\exp(\sqrt{-1}xt)) = \exp(-d|t|^\alpha),$$

which does not have a closed-form inverse except for $\alpha = 2$ (normal) or $\alpha = 1$ (Cauchy). Note that when $\alpha = 2$, d corresponds to “ σ^2 ” (not “ σ ”) in a normal.

Corresponding to the leading two rows in \mathbf{A} , $u_1, u_2 \in \mathbb{R}^D$, the leading two rows in \mathbf{B} are $v_1 = \mathbf{R}^T u_1$, $v_2 = \mathbf{R}^T u_2$. The entries of the difference,

$$x_j = v_{1,j} - v_{2,j} = \sum_{i=1}^D r_{ij} (u_{1,i} - u_{2,i}) \sim S\left(\alpha, d_{(\alpha)} = \sum_{i=1}^D |u_{1,i} - u_{2,i}|^\alpha\right),$$

for $j = 1$ to k , are i.i.d. samples from a stable distribution with the scale parameter being the l_α distance $d_{(\alpha)}$, due to properties of Fourier transforms. For example, when $\alpha = 2$, a weighted sum of i.i.d. standard normals is also normal with the scale parameter (i.e., variance) being the sum of squares of all weights.

Once we obtain the stable samples, one can discard the original matrix \mathbf{A} and the remaining task is to estimate the scale parameter $d_{(\alpha)}$ for each pair.

Some applications of *stable random projections* are summarized as follows:

- **Computing all pairwise distances** The cost of computing all pairwise distances of $\mathbf{A} \in \mathbb{R}^{n \times D}$, $O(n^2 D)$, is significantly reduced to $O(nDk + n^2 k)$.
- **Estimating l_α distances online** For $n > 10^5$, it is challenging or unrealistic to materialize all pairwise distances in \mathbf{A} . Thus, in applications such as online learning, databases, search engines, and online recommendation systems, it is often more efficient if we store $\mathbf{B} \in \mathbb{R}^{n \times k}$ in the memory and estimate any distance *on the fly* if needed. Estimating distances online is the standard strategy in large-scale kernel learning[4]. With *stable random projections*, this simple strategy becomes effective in high-dimensional data.
- **Learning with dynamic streaming data** In reality, the data matrix may be updated overtime. In fact, with streaming data arriving at high-rate[1,10], the “data matrix” may be never stored and hence all operations (such as clustering and classification) must be conducted on the fly. The

method of *stable random projections* provides a scheme to compute and update distances on the fly in one-pass of the data; see relevant papers (e.g., [1]) for more details on this important and fast-developing subject.

- **Estimating entropy** The entropy distance $\sum_{i=1}^D |u_{1,i} - u_{2,i}| \log |u_{1,i} - u_{2,i}|$ is a useful statistic. A workshop in NIPS'03 (www.menem.com/~ilya/pages/NIPS03) focused on entropy estimation. A recent practical algorithm is simply using the difference between the l_{α_1} and l_{α_2} distances [11], where $\alpha_1 = 1.05$, $\alpha_2 = 0.95$, and the distances were estimated by *stable random projections*.

If one tunes the l_α distances for many different α (e.g., [5]), then *stable random projections* will be even more desirable as a cost-saving device.

2 The Statistical Estimation Problem

Recall that the method of *stable random projections* boils down to a statistical estimation problem. That is, estimating the scale parameter $d_{(\alpha)}$ from k i.i.d. samples $x_j \sim S(\alpha, d_{(\alpha)})$, $j = 1$ to k . We consider that a good estimator $\hat{d}_{(\alpha)}$ should have the following desirable properties:

- (Asymptotically) unbiased and small variance.
- Computationally efficient.
- Exponential decrease of error (tail) probabilities.

The *arithmetic mean* estimator $\frac{1}{k} \sum_{j=1}^k |x_j|^2$ is good for $\alpha = 2$. When $\alpha < 2$, the task is less straightforward because (1) no explicit density of x_j exists unless $\alpha = 1$ or $0+$; and (2) $E(|x_j|^t) < \infty$ only when $-1 < t < \alpha$.

2.1 Several Previous Estimators

Initially reported in arXiv in 2006, [2] proposed the *geometric mean* estimator

$$\hat{d}_{(\alpha),gm} = \frac{\prod_{j=1}^k |x_j|^{\alpha/k}}{\left[\frac{2}{\pi} \Gamma\left(\frac{\alpha}{k}\right) \Gamma\left(1 - \frac{1}{k}\right) \sin\left(\frac{\pi}{2} \frac{\alpha}{k}\right) \right]^k}.$$

where $\Gamma(\cdot)$ is the Gamma function, and the *harmonic mean* estimator

$$\hat{d}_{(\alpha),hm} = \frac{-\frac{2}{\pi} \Gamma(-\alpha) \sin\left(\frac{\pi}{2} \alpha\right)}{\sum_{j=1}^k |x_j|^{-\alpha}} \left(k - \left(\frac{-\pi \Gamma(-2\alpha) \sin(\pi \alpha)}{[\Gamma(-\alpha) \sin\left(\frac{\pi}{2} \alpha\right)]^2} - 1 \right) \right).$$

More recently, [3] proposed the *fractional power* estimator

$$\hat{d}_{(\alpha),fp} = \left(\frac{1}{k} \frac{\sum_{j=1}^k |x_j|^{\lambda^* \alpha}}{\frac{2}{\pi} \Gamma(1 - \lambda^*) \Gamma(\lambda^* \alpha) \sin\left(\frac{\pi}{2} \lambda^* \alpha\right)} \right)^{1/\lambda^*} \times \left(1 - \frac{1}{k} \frac{1}{2\lambda^*} \left(\frac{1}{\lambda^*} - 1 \right) \left(\frac{\frac{2}{\pi} \Gamma(1 - 2\lambda^*) \Gamma(2\lambda^* \alpha) \sin(\pi \lambda^* \alpha)}{[\frac{2}{\pi} \Gamma(1 - \lambda^*) \Gamma(\lambda^* \alpha) \sin\left(\frac{\pi}{2} \lambda^* \alpha\right)]^2} - 1 \right) \right),$$

where

$$\lambda^* = \underset{-\frac{1}{2\alpha} < \lambda < \frac{1}{2}}{\operatorname{argmin}} \frac{1}{\lambda^2} \left(\frac{\frac{2}{\pi} \Gamma(1 - 2\lambda) \Gamma(2\lambda \alpha) \sin(\pi \lambda \alpha)}{[\frac{2}{\pi} \Gamma(1 - \lambda) \Gamma(\lambda \alpha) \sin\left(\frac{\pi}{2} \lambda \alpha\right)]^2} - 1 \right).$$

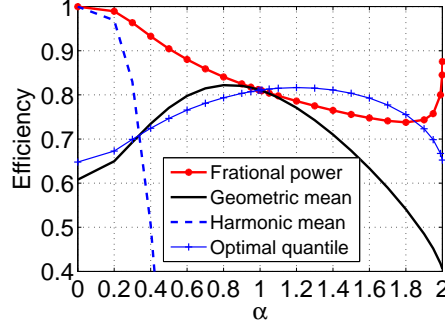


Fig. 1. The Cramér-Rao efficiencies (the higher the better, max = 100%) of various estimators, including the *optimal quantile* estimator proposed in this study.

All three estimators are unbiased or asymptotically (as $k \rightarrow \infty$) unbiased. Figure 1 compares their asymptotic variances in terms of the Cramér-Rao efficiency, which is the ratio of the smallest possible asymptotic variance over the asymptotic variance of the estimator, as $k \rightarrow \infty$.

The *geometric mean* estimator, $\hat{d}_{(\alpha),gm}$ exhibits tail bounds in exponential forms, i.e., the errors decrease exponentially fast:

$$\Pr\left(|\hat{d}_{(\alpha),gm} - d_{(\alpha)}| \geq \epsilon d_{(\alpha)}\right) \leq 2 \exp\left(-k \frac{\epsilon^2}{G_{gm}}\right).$$

The *harmonic mean* estimator, $\hat{d}_{(\alpha),hm}$, works well for small α , and has exponential tail bounds for $\alpha = 0+$.

The *fractional power* estimator, $\hat{d}_{(\alpha),fp}$, has smaller asymptotic variance than both the *geometric mean* and *harmonic mean* estimators. However, it does not have exponential tail bounds, due to the restriction $-1 < \lambda^* \alpha < \alpha$ in its definition. As shown in [3], it only has finite moments slightly higher than the 2nd order, when α approaches 2 (because $\lambda^* \rightarrow 0.5$), meaning that large errors may have a good chance to occur. We will demonstrate this by simulations.

2.2 The Issue of Computational Efficiency

In the definitions of $\hat{d}_{(\alpha),gm}$, $\hat{d}_{(\alpha),hm}$ and $\hat{d}_{(\alpha),fp}$, all three estimators require evaluating fractional powers, e.g., $|x_j|^{\alpha/k}$. This operation is relatively expensive, especially if we need to conduct this tens of billions of times (e.g., $n^2 = 10^{10}$).

For example, [5] reported that, although the radial basis kernel (2) with $\alpha = 0.5$ achieved good performance, it was not preferred because evaluating the square root was too expensive.

2.3 Our Proposed Estimator

We propose the *optimal quantile* estimator, using the q^* th smallest $|x_j|$:

$$\hat{d}_{(\alpha),oq} \propto (q^*\text{-quantile}\{|x_j|, j = 1, 2, \dots, k\})^\alpha, \quad (3)$$

where $q^* = q^*(\alpha)$ is chosen to minimize the asymptotic variance.

This estimator is computationally attractive because **selecting** should be much less expensive than evaluating fractional powers. If we are interested in $d_{(\alpha)}^{1/\alpha}$ instead, then we do not even need to evaluate any fractional powers.

As mentioned, in many cases using either $d_{(\alpha)}$ or $d_{(\alpha)}^{1/\alpha}$ makes no difference and $d_{(\alpha)}$ is often preferred because it avoids taking $(\cdot)^{1/\alpha}$ power. The radial basis kernel (2) requires $d_{(\alpha)}$. Thus this study focuses on $d_{(\alpha)}$. On the other hand, if we can estimate $d_{(\alpha)}^{1/\alpha}$ directly, for example, using (3) without the α th power, we might as well just use $d_{(\alpha)}^{1/\alpha}$ if permitted. In case we do not need to evaluate any fractional power, our estimator will be even more computationally efficient.

In addition to the computational advantages, this estimator also has good theoretical properties, in terms of both the variances and tail probabilities:

1. Figure 1 illustrates that, compared with the *geometric mean* estimator, its asymptotic variance is about the same when $\alpha < 1$, and is considerably smaller when $\alpha > 1$. Compared with the *fractional power* estimator, it has smaller asymptotic variance when $1 < \alpha \leq 1.8$. In fact, as will be shown by simulations, when the sample size k is not too large, its mean square errors are considerably smaller than the *fractional power* estimator when $\alpha > 1$.
2. The *optimal quantile* estimator exhibits tail bounds in exponential forms. This theoretical contribution is practically important, for selecting the sample size k . In learning theory, the generalization bounds are often loose. In our case, however, the bounds are tight because the distribution is specified.

The next section will be devoted to analyzing the *optimal quantile* estimator.

3 The Optimal Quantile Estimator

Recall the goal is to estimate $d_{(\alpha)}$ from $\{x_j\}_{j=1}^k$, where $x_j \sim S(\alpha, d_{(\alpha)})$, i.i.d. Since the distribution belongs to the scale family, one can estimate the scale parameter from quantiles. Due to symmetry, it is natural to consider the absolute values:

$$\hat{d}_{(\alpha),q} = \left(\frac{q\text{-Quantile}\{|x_j|, j = 1, 2, \dots, k\}}{q\text{-Quantile}\{|S(\alpha, 1)|\}} \right)^\alpha, \quad (4)$$

which is best understood by the fact that if $x \sim S(\alpha, 1)$, then $d^{1/\alpha}x \sim S(\alpha, d)$, or more obviously, if $x \sim N(0, 1)$, then $(\sigma^2)^{1/2}x \sim N(0, \sigma^2)$. By properties of order statistics [12], any q -quantile will provide an asymptotically unbiased estimator.

Lemma 1 provides the asymptotic variance of $\hat{d}_{(\alpha),q}$.

Lemma 1. Denote $f_X(x; \alpha, d_{(\alpha)})$ and $F_X(x; \alpha, d_{(\alpha)})$ the probability density function and the cumulative density function of $X \sim S(\alpha, d_{(\alpha)})$, respectively.

The asymptotic variance of $\hat{d}_{(\alpha),q}$ defined in (4) is

$$\text{Var}(\hat{d}_{(\alpha),q}) = \frac{1}{k} \frac{(q - q^2)\alpha^2/4}{f_X^2(W; \alpha, 1) W^2} d_{(\alpha)}^2 + O\left(\frac{1}{k^2}\right) \quad (5)$$

where $W = F_X^{-1}((q + 1)/2; \alpha, 1) = q\text{-Quantile}\{|S(\alpha, 1)|\}$.

Proof: See Appendix A. \square .

3.1 Optimal Quantile $q^*(\alpha)$

We choose $q = q^*(\alpha)$ so that the asymptotic variance (5) is minimized, i.e.,

$$q^*(\alpha) = \underset{q}{\operatorname{argmin}} g(q; \alpha), \quad g(q; \alpha) = \frac{q - q^2}{f_X^2(W; \alpha, 1) W^2}. \quad (6)$$

The convexity of $g(q; \alpha)$ is important. Graphically, $g(q; \alpha)$ is a convex function of q , i.e., a unique minimum exists. An algebraic proof, however, is difficult. Nevertheless, we can obtain analytical solutions when $\alpha = 1$ and $\alpha = 0+$.

Lemma 2. *When $\alpha = 1$ or $\alpha = 0+$, the function $g(q; \alpha)$ defined in (6) is a convex function of q . When $\alpha = 1$, the optimal $q^*(1) = 0.5$. When $\alpha = 0+$, $q^*(0+) = 0.203$ is the solution to $-\log q^* + 2q^* - 2 = 0$.*

Proof: See Appendix B. \square .

It is also easy to show that when $\alpha = 2$, $q^*(2) = 0.862$.

We denote the *optimal quantile estimator* by $\hat{d}_{(\alpha),oq}$, which is same as $\hat{d}_{(\alpha),q^*}$. For general α , we resort to numerical solutions, as presented in Figure 2.

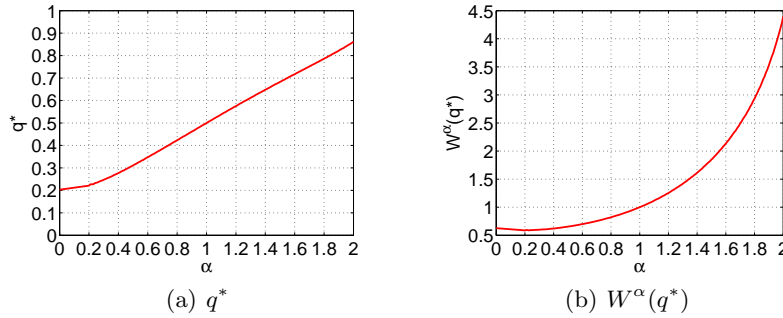


Fig. 2. (a) The optimal values for $q^*(\alpha)$, which minimizes asymptotic variance of $\hat{d}_{(\alpha),q}$, i.e., the solution to (6). (b) The constant $W^\alpha(q^*) = \{q^*\text{-quantile}\{|S(\alpha, 1)|\}\}^\alpha$.

3.2 Bias Correction

Although $\hat{d}_{(\alpha),oq}$ (i.e., $\hat{d}_{(\alpha),q^*}$) is asymptotically (as $k \rightarrow \infty$) unbiased, it is seriously biased for small k . Thus, it is practically important to remove the bias. The unbiased version of the *optimal quantile estimator* is

$$\hat{d}_{(\alpha),oq,c} = \hat{d}_{(\alpha),oq} / B_{\alpha,k}, \quad (7)$$

where $B_{\alpha,k}$ is the expectation of $\hat{d}_{(\alpha),oq}$ at $d_{(\alpha)} = 1$. For $\alpha = 1, 0+$, or 2 , we can evaluate the expectations (i.e., integrals) analytically or by numerical integrations. For general α , as the probability density is not available, the task is difficult and prone to numerical instability. On the other hand, since the Monte-Carlo simulation is a popular alternative for evaluating difficult integrals, a practical solution is to simulate the expectations, as presented in Figure 3.

Figure 3 illustrates that $B_{\alpha,k} > 1$, meaning that this correction also reduces variance while removing bias (because $\operatorname{Var}(x/c) = \operatorname{Var}(x)/c^2$). For example, when

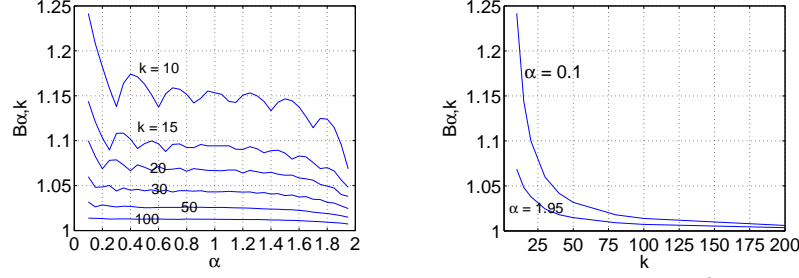


Fig. 3. The bias correction factor $B_{\alpha,k}$ in (7), obtained from 10^8 simulations for every combination of α (spaced at 0.05) and k . $B_{\alpha,k} = E(\hat{d}_{(\alpha),oq}; d_{(\alpha)} = 1)$.

$\alpha = 0.1$ and $k = 10$, $B_{\alpha,k} \approx 1.24$, which is significant, because $1.24^2 = 1.54$ implies a 54% difference in terms of variance, and even more considerable in terms of the mean square errors $MSE = \text{variance} + \text{bias}^2$.

$B_{\alpha,k}$ can be tabulated for small k , and absorbed into other coefficients, i.e., this does not increase the computational cost at run time. We fix $B_{\alpha,k}$ as reported in Figure 3. The simulations in Section 4 directly used those fixed $B_{\alpha,k}$ values.

3.3 Computational Efficiency

Figure 4 compares the computational costs of the *geometric mean*, the *fractional power*, and the *optimal quantile* estimators. The *harmonic mean* estimator was not included as it costs very similarly to the *fractional power* estimator.

We used the build-in function “pow” in gcc for evaluating the fractional powers. We implemented a “quick select” algorithm, which is similar to quick sort and requires on average linear time. For simplicity, our implementation used recursions and the middle element as pivot. Also, to ensure fairness, for all estimators, coefficients which are functions of α and/or k were pre-computed.

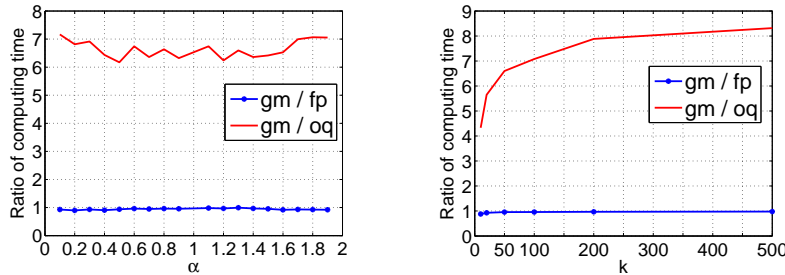


Fig. 4. Relative computational cost ($\hat{d}_{(\alpha),gm}$ over $\hat{d}_{(\alpha),oq,c}$ and $\hat{d}_{(\alpha),gm}$ over $\hat{d}_{(\alpha),fp}$), from 10^6 simulations at each combination of α and k . The left panel averages over all k and the right panel averages over all α . Note that the cost of $\hat{d}_{(\alpha),oq,c}$ includes evaluating the α th moment once.

Normalized by the computing time of $\hat{d}_{(\alpha),gm}$, we observe that relative computational efficiency does not strongly depend on α . We do observe that the ratio

of computing time of $\hat{d}_{(\alpha),gm}$ over that of $\hat{d}_{(\alpha),oq,c}$ increases consistently with increasing k . This is because in the definition of $\hat{d}_{(\alpha),oq}$ (and hence also $\hat{d}_{(\alpha),oq,c}$), it is required to evaluate the fractional power once, which contributes to the total computing time more significantly at smaller k .

Figure 4 illustrates that, (A) the *geometric mean* estimator and the *fractional power* estimator are similar in terms of computational efficiency; (B) the *optimal quantile* estimator is nearly one order of magnitude more computationally efficient than the *geometric mean* and *fractional power* estimators. Because we implemented a “naïve” version of “quick select” using recursions and simple pivoting, the actual improvement may be more significant. Also, if applications require only $d_{(\alpha)}^{1/\alpha}$, then no fractional power operations are needed for $\hat{d}_{(\alpha),oq,c}$ and the improvement will be even more considerable.

3.4 Error (Tail) Bounds

Error (tail) bounds are essential for determining k . The variance alone is not sufficient for that purpose. If an estimator of d , say \hat{d} , is normally distributed, $\hat{d} \sim N(d, \frac{1}{k}V)$, the variance suffices for choosing k because its error (tail) probability $\Pr(|\hat{d} - d| \geq \epsilon d) \leq 2 \exp\left(-k \frac{\epsilon^2}{2V}\right)$ is determined by V . In general, a reasonable estimator will be asymptotically normal, for small enough ϵ and large enough k . For a finite k and a fixed ϵ , however, the normal approximation may be (very) poor. This is especially true for the *fractional power* estimator, $\hat{d}_{(\alpha),fp}$.

Thus, for a good motivation, Lemma 3 provides the error (tail) probability bounds of $\hat{d}_{(\alpha),q}$ for any q , not just the optimal quantile q^* .

Lemma 3. Denote $X \sim S(\alpha, d_{(\alpha)})$ and its probability density function by $f_X(x; \alpha, d_{(\alpha)})$ and cumulative function by $F_X(x; \alpha, d_{(\alpha)})$. Given $x_j \sim S(\alpha, d_{(\alpha)})$, i.i.d., $j = 1$ to k . Using $\hat{d}_{(\alpha),q}$ in (4), then

$$\Pr\left(\hat{d}_{(\alpha),q} \geq (1 + \epsilon)d_{(\alpha)}\right) \leq \exp\left(-k \frac{\epsilon^2}{G_{R,q}}\right), \epsilon > 0, \quad (8)$$

$$\Pr\left(\hat{d}_{(\alpha),q} \leq (1 - \epsilon)d_{(\alpha)}\right) \leq \exp\left(-k \frac{\epsilon^2}{G_{L,q}}\right), 0 < \epsilon < 1, \quad (9)$$

$$\frac{\epsilon^2}{G_{R,q}} = -(1 - q) \log(2 - 2F_R) - q \log(2F_R - 1) + (1 - q) \log(1 - q) + q \log q, \quad (10)$$

$$\frac{\epsilon^2}{G_{L,q}} = -(1 - q) \log(2 - 2F_L) - q \log(2F_L - 1) + (1 - q) \log(1 - q) + q \log q, \quad (11)$$

$$W = F_X^{-1}((q + 1)/2; \alpha, 1) = q\text{-quantile}\{|S(\alpha, 1)|\},$$

$$F_R = F_X\left((1 + \epsilon)^{1/\alpha} W; \alpha, 1\right), \quad F_L = F_X\left((1 - \epsilon)^{1/\alpha} W; \alpha, 1\right).$$

As $\epsilon \rightarrow 0+$

$$\lim_{\epsilon \rightarrow 0+} G_{R,q} = \lim_{\epsilon \rightarrow 0+} G_{L,q} = \frac{q(1 - q)\alpha^2/2}{f_X^2(W; \alpha, 1)W^2}. \quad (12)$$

Proof: See Appendix C. \square

The limit in (12) as $\epsilon \rightarrow 0$ is precisely twice the asymptotic variance factor of $\hat{d}_{(\alpha),q}$ in (5), consistent with the normality approximation mentioned previously. This explains why we express the constants as ϵ^2/G . (12) also indicates that the tail bounds achieve the “optimal rate” for this estimator, in the language of large deviation theory.

By the Bonferroni bound, it is easy to determine the sample size k

$$\Pr\left(|\hat{d}_{(\alpha),q} - d_{(\alpha)}| \geq \epsilon d_{(\alpha)}\right) \leq 2 \exp\left(-k \frac{\epsilon^2}{G}\right) \leq \delta/(n^2/2) \implies k \geq \frac{G}{\epsilon^2} (2 \log n - \log \delta).$$

Lemma 4. Using $\hat{d}_{(\alpha),q}$ with $k \geq \frac{G}{\epsilon^2} (2 \log n - \log \delta)$, any pairwise l_α distance among n points can be approximated within a $1 \pm \epsilon$ factor with probability $\geq 1 - \delta$. It suffices to let $G = \max\{G_{R,q}, G_{L,q}\}$, where $G_{R,q}, G_{L,q}$ are defined in Lemma 3.

The Bonferroni bound can be unnecessarily conservative. It is often reasonable to replace $\delta/(n^2/2)$ by δ/T , meaning that except for a $1/T$ fraction of pairs, any distance can be approximated within a $1 \pm \epsilon$ factor with probability $1 - \delta$.

Figure 5 plots the error bound constants for $\epsilon < 1$, for both the recommended *optimal quantile* estimator $\hat{d}_{(\alpha),oq}$ and the baseline *sample median* estimator $\hat{d}_{(\alpha),q=0.5}$. Although we choose $\hat{d}_{(\alpha),oq}$ based on the asymptotic variance, it turns out $\hat{d}_{(\alpha),oq}$ also exhibits (much) better tail behaviors (i.e., smaller constants) than $\hat{d}_{(\alpha),q=0.5}$, at least in the range of $\epsilon < 1$.

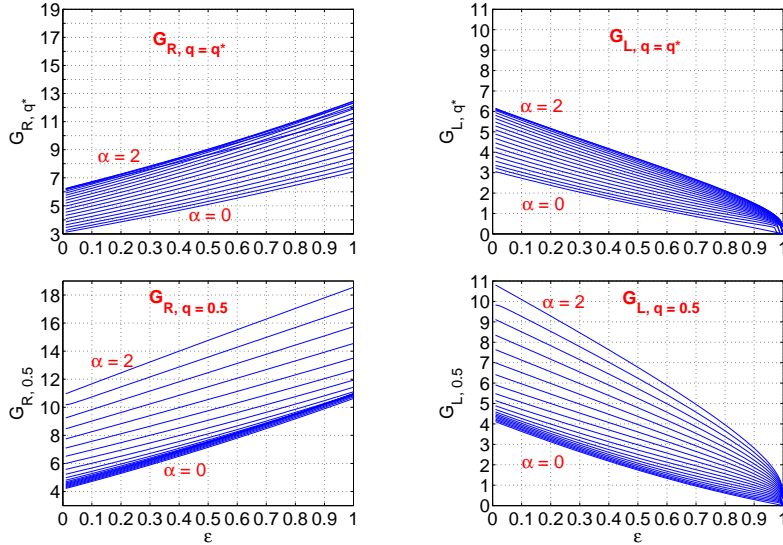


Fig. 5. Tail bound constants for quantile estimators; the lower the better. Upper panels: optimal quantile estimators $\hat{d}_{(\alpha),q^*}$. Lower panels: median estimators $\hat{d}_{(\alpha),q=0.5}$.

Consider $k = \frac{G}{\epsilon^2} (\log 2T - \log \delta)$ (recall we suggest replacing $n^2/2$ by T), with $\delta = 0.05$, $\epsilon = 0.5$, and $T = 10$. Because $G_{R,q^*} \approx 5 \sim 9$ around $\epsilon = 0.5$, we obtain

$k \approx 120 \sim 215$, which is still a relatively large number (although the original dimension D might be 10^6). If we choose $\epsilon = 1$, then approximately $k \approx 40 \sim 65$.

It is possible $k = 120 \sim 215$ might be still conservative, for three reasons: (A) the tail bounds, although “sharp,” are still upper bounds; (B) using $G = \max\{G_{R,q^*}, G_{L,q^*}\}$ is conservative because G_{L,q^*} is usually much smaller than G_{R,q^*} ; (C) this type of tail bounds is based on relative error, which may be stringent for small (≈ 0) distances.

In fact, some earlier studies on *normal random projections* (i.e., $\alpha = 2$) [13,14] empirically demonstrated that $k \geq 50$ appeared sufficient.

4 Simulations

We resort to simulations for comparing the finite sample variances of various estimators and assessing the more precise error (tail) probabilities.

One advantage of *stable random projections* is that we know the (manually generated) distributions and the only source of errors is from the random number generations. Thus, we can simply rely on simulations to evaluate the estimators without using real data. In fact, after projections, the projected data follow exactly the stable distribution, regardless of the original real data distribution.

Without loss of generality, we simulate samples from $S(\alpha, 1)$ and estimate the scale parameter (i.e., 1) from the samples. Repeating the procedure 10^7 times, we can reliably evaluate the mean square errors (MSE) and tail probabilities.

4.1 Mean Square Errors (MSE)

As illustrated in Figure 6, in terms of the MSE, the *optimal quantile* estimator $\hat{d}_{(\alpha),oq,c}$ outperforms both the *geometric mean* and *fractional power* estimators when $\alpha > 1$ and $k \geq 20$. The *fractional power* estimator does not appear to be very suitable for $\alpha > 1$, especially for α close to 2, even when the sample size k is not too small (e.g., $k = 50$). For $\alpha < 1$, however, the *fractional power* estimator has good performance in terms of MSE, even for small k .

4.2 Error(Tail) Probabilities

Figure 7 presents the simulated right tail probabilities, $\Pr(\hat{d}_{(\alpha)} \geq (1 + \epsilon)d_{(\alpha)})$, illustrating that when $\alpha > 1$, the *fractional power* estimator can exhibit very bad tail behaviors. For $\alpha < 1$, the *fractional power* estimator demonstrates good performance at least for the probability range in the simulations.

Thus, Figure 7 demonstrates that the *optimal quantile* estimator consistently outperforms the *fractional power* and the *geometric mean* estimators when $\alpha > 1$.

5 The Related Work

There have been many studies of *normal random projections* in machine learning, for dimension reduction in the l_2 norm, e.g., [14], highlighted by the Johnson-Lindenstrauss (JL) Lemma [15], which says $k = O(\log n / \epsilon^2)$ suffices when using normal (or normal-like, e.g., [16]) projection methods.

The method of *stable random projections* is applicable for computing the l_α distances ($0 < \alpha \leq 2$), not just for l_2 . [1, Lemma 1, Lemma 2, Theorem 3]

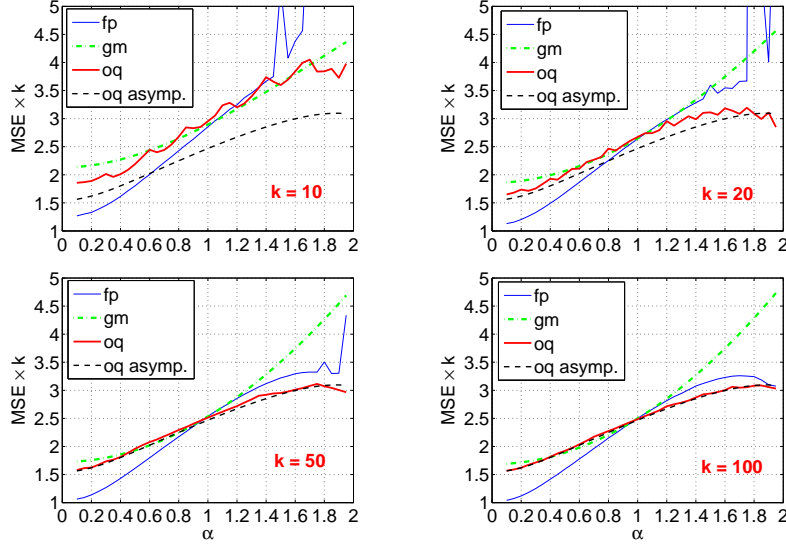


Fig. 6. Empirical mean square errors (MSE, the lower the better), from 10^7 simulations at every combination of α and k . The values are multiplied by k so that four plots can be at about the same scale. The MSE for the *geometric mean* (gm) estimator is computed exactly since closed-form expression exists. The lower dashed curves are the asymptotic variances of the *optimal quantile* (oq) estimator.

suggested the *median* (i.e., $q = 0.5$ quantile) estimator for $\alpha = 1$ and argued that the sample complexity bound should be $O(1/\epsilon^2)$ ($n = 1$ in their study). Their bound was not provided in an explicit form and required an “ ϵ is small enough” argument. For $\alpha \neq 1$, [1, Lemma 4] only provided a conceptual algorithm, which “is not uniform.” In this study, we prove the bounds for any q -quantile and any $0 < \alpha \leq 2$ (not just $\alpha = 1$), in explicit exponential forms, with no unknown constants and no restriction that “ ϵ is small enough.”

The quantile estimator for stable distributions was proposed in statistics quite some time ago, e.g., [17,18]. [17] mainly focused on $1 \leq \alpha \leq 2$ and recommended using $q = 0.44$ quantiles (mainly for the sake of smaller bias). [18] focused on $0.6 \leq \alpha \leq 2$ and recommended $q = 0.5$ quantiles.

This study considers all $0 < \alpha \leq 2$ and recommends q based on the minimum asymptotic variance. Because the bias can be easily removed (at least in the practical sense), it appears not necessary to use other quantiles only for the sake of smaller bias. Tail bounds, which are useful for choosing q and k based on confidence intervals, were not available in [17,18].

Finally, one might ask if there might be better estimators. For $\alpha = 1$, [19] proposed using a linear combination of quantiles (with carefully chosen coefficients) to obtain an asymptotically optimal estimator for the Cauchy scale parameter. While it is possible to extend their result to general $0 < \alpha < 2$ (requiring some non-trivial work), whether or not it will be practically better than the *optimal*

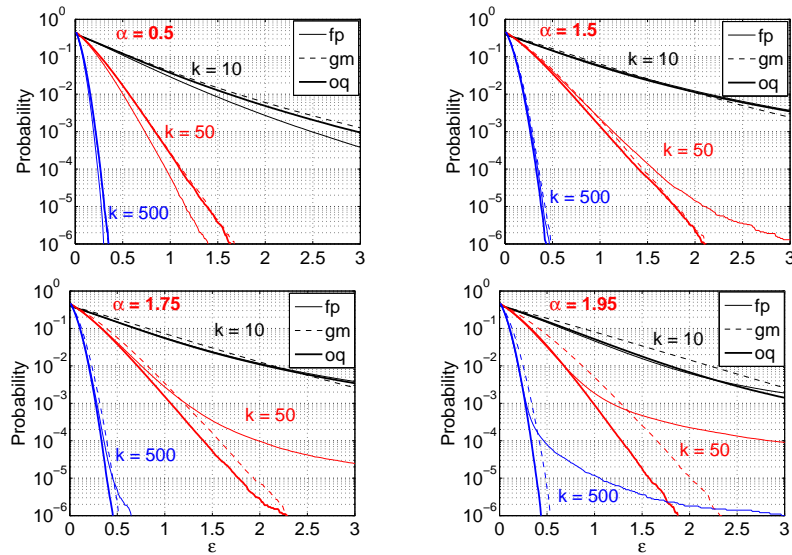


Fig. 7. The right tail probabilities (the lower the better), from 10^7 simulations at each combination of α and k .

quantile estimator is unclear because the extreme quantiles severely affect the tail probabilities and finite-sample variances and hence some kind of truncation (i.e., discarding some samples at extreme quantiles) is necessary. Also, exponential tail bounds of the linear combination of quantiles may not exist or may not be feasible to derive. In addition, the *optimal quantile* estimator is computationally more efficient.

6 Conclusion

Many machine learning algorithms operate on the training data only through pairwise distances. Computing, storing, updating and retrieving the “matrix” of pairwise distances is challenging in applications involving massive, high-dimensional, and possibly streaming, data. For example, the pairwise distance matrix can not fit in memory when the number of observations exceeds 10^6 (or even 10^5).

The method of *stable random projections* provides an efficient mechanism for computing pairwise distances using low memory, by transforming the original high-dimensional data into *sketches*, i.e., a small number of samples from α -stable distributions, which are much easier to store and retrieve.

This method provides a uniform scheme for computing the l_α pairwise distances for all $0 < \alpha \leq 2$. Choosing an appropriate α is often critical to the performance of learning algorithms. In principle, we can tune algorithms for many l_α distances; and *stable random projections* can provide an efficient tool.

To recover the original distances, we face an estimation task. Compared with previous estimators based on the *geometric mean*, the *harmonic mean*, or the *fractional power*, the proposed *optimal quantile* estimator exhibits two advantages. Firstly, the *optimal quantile* estimator is nearly one order of magnitude

more efficient than other estimators (e.g., reducing the training time from one week to one day). Secondly, the *optimal quantile* estimator is considerably more accurate when $\alpha > 1$, in terms of both the variances and error (tail) probabilities. Note that $\alpha \geq 1$ corresponds to a convex norm (satisfying the triangle inequality), which might be another motivation for using l_α distances with $\alpha \geq 1$.

One theoretical contribution is the explicit tail bounds for general quantile estimators and consequently the sample complexity bound $k = O(\log n/\epsilon^2)$. Those bounds may guide practitioners in choosing k , the number of projections. The (practically useful) bounds are expressed in terms of the probability functions and hence they might be not as convenient for further theoretical analysis. Also, we should mention that the bounds do not recover the optimal bound of the *arithmetic mean* estimator when $\alpha = 2$, because the *arithmetic mean* estimator is statistically optimal at $\alpha = 2$ but the *optimal quantile* estimator is not.

While we believe that applying *stable random projections* in machine learning has become straightforward, there are interesting theoretical issues for future research. For example, how theoretical properties of learning algorithms may be affected if the approximated (instead of exact) l_α distances are used?

A Proof of Lemma 1

Denote $f_X(x; \alpha, d_{(\alpha)})$ and $F_X(x; \alpha, d_{(\alpha)})$ the probability density function and the cumulative density function of $X \sim S(\alpha, d_{(\alpha)})$, respectively. Similarly we use $f_Z(z; \alpha, d_{(\alpha)})$ and $F_Z(z; \alpha, d_{(\alpha)})$ for $Z = |X|$. Due to symmetry, the following relations hold

$$\begin{aligned} f_Z(z; \alpha, d_{(\alpha)}) &= 2f_X(z; \alpha, d_{(\alpha)}) = 2/d_{(\alpha)}^{1/\alpha} f_X(z/d_{(\alpha)}^{1/\alpha}; \alpha, 1), \\ F_Z(z; \alpha, d_{(\alpha)}) &= 2F_X(z; \alpha, d_{(\alpha)}) - 1 = 2F_X(z/d_{(\alpha)}^{1/\alpha}; \alpha, 1) - 1, \\ F_Z^{-1}(q; \alpha, d_{(\alpha)}) &= F_X^{-1}((q+1)/2; \alpha, d_{(\alpha)}) = d_{(\alpha)}^{1/\alpha} F_X^{-1}((q+1)/2; \alpha, 1). \end{aligned}$$

Let $W = q\text{-Quantile}\{|S(\alpha, 1)|\} = F_X^{-1}((q+1)/2; \alpha, 1)$ and $W_d = F_Z^{-1}(q; \alpha, d_{(\alpha)}) = d_{(\alpha)}^{1/\alpha} W$. Then, following known statistical results, e.g., [12, Theorem 9.2], the asymptotic variance of $\hat{d}_{\alpha, q}^{1/\alpha}$ should be

$$\begin{aligned} \text{Var}(\hat{d}_{\alpha, q}^{1/\alpha}) &= \frac{1}{k} \frac{q - q^2}{f_Z^2(W_d; \alpha, d_{(\alpha)}) W^2} + O\left(\frac{1}{k^2}\right) = \frac{1}{k} \frac{q - q^2}{d_{(\alpha)}^{-2/\alpha} f_Z^2(W; \alpha, 1) W^2} + O\left(\frac{1}{k^2}\right) \\ &= \frac{1}{k} \frac{q - q^2}{4d_{(\alpha)}^{-2/\alpha} f_X^2(W; \alpha, 1) W^2} + O\left(\frac{1}{k^2}\right). \end{aligned}$$

By ‘‘delta method,’’ i.e., $\text{Var}(h(x)) \approx \text{Var}(x) (h'(x))^2$,

$$\text{Var}(\hat{d}_{\alpha, q}) = \text{Var}(\hat{d}_{\alpha, q}^{1/\alpha}) \left(\alpha d_{(\alpha)}^{(\alpha-1)/\alpha} \right)^2 + O\left(\frac{1}{k^2}\right) = \frac{1}{k} \frac{(q - q^2)\alpha^2/4}{f_X^2(W; \alpha, 1) W^2} d_{(\alpha)}^2 + O\left(\frac{1}{k^2}\right).$$

B Proof of Lemma 2

First, consider $\alpha = 1$. In this case,

$$f_X(x; 1, 1) = \frac{1}{\pi} \frac{1}{x^2 + 1}, \quad W = F_X^{-1}((q+1)/2; 1, 1) = \tan\left(\frac{\pi}{2}q\right),$$

$$g(q; 1) = \frac{q - q^2}{\left(\frac{2}{\pi} \frac{1}{\tan^2(\frac{\pi}{2}q) + 1}\right)^2 \tan^2(\frac{\pi}{2}q)} = \frac{q - q^2}{\sin^2(\pi q)} \pi^2.$$

It suffices to study $L(q) = \log g(q; 1)$.

$$L'(q) = \frac{1}{q} - \frac{1}{1-q} - \frac{2\pi \cos(\pi q)}{\sin(\pi q)}, \quad L''(q) = -\frac{1}{q^2} - \frac{1}{(1-q)^2} + \frac{2\pi^2}{\sin^2(\pi q)}.$$

Because $\sin(x) \leq x$ for $x \geq 0$, it is easy to see that $\frac{\pi}{\sin(\pi q)} - \frac{1}{q} \geq 0$, and $\frac{\pi}{\sin(\pi q)} - \frac{1}{1-q} = \frac{\pi}{\sin(\pi(1-q))} - \frac{1}{1-q} \geq 0$. Thus, $L'' \geq 0$, i.e., $L(q)$ is convex and so is $g(q; 1) = e^{L(q)}$. Since $L'(1/2) = 0$, we know $q^*(1) = 0.5$.

Next we consider $\alpha = 0+$, using a fact [2] that as $\alpha \rightarrow 0+$, $|S(\alpha, 1)|^\alpha$ converges to $1/E_1$, where E_1 stands for an exponential distribution with mean 1.

Denote $h = d_{(0+)}$ and $z_j \sim h/E_1$. The sample quantile estimator becomes

$$\hat{d}_{(0+),q} = \frac{q\text{-Quantile}\{|z_j|, j = 1, 2, \dots, k\}}{q\text{-Quantile}\{1/E_1\}}.$$

In this case,

$$f_Z(z; h) = e^{-h/z} \frac{h}{z^2}, \quad F_Z^{-1}(q; h) = -\frac{h}{\log q},$$

$$\text{Var}\left(\hat{d}_{(0+),q}\right) = \frac{1}{k} \frac{1-q}{q \log^2 q} h^2 + O\left(\frac{1}{k^2}\right).$$

It is straightforward to show that $\frac{1-q}{q \log^2 q}$ is a convex function of q and the minimum is attained by solving $-\log q^* + 2q^* - 2 = 0$, i.e., $q^* = 0.203$.

C Proof of Lemma 3

Given k i.i.d. samples, $x_j \sim S(\alpha, d_{(\alpha)})$, $j = 1$ to k . Let $z_j = |x_j|$, $j = 1$ to k . Denote by $F_Z(t; \alpha, d_{(\alpha)})$ the cumulative density of z_j , and by $F_{Z,k}(t; \alpha, d_{(\alpha)})$ the empirical cumulative density of z_j , $j = 1$ to k .

It is the basic fact [12] about order statistics that $kF_{Z,k}(t; \alpha, d_{(\alpha)})$ follows a binomial, i.e., $kF_{Z,k}(t; \alpha, d_{(\alpha)}) \sim \text{Bin}(k, F_Z(t; \alpha, d_{(\alpha)}))$. For simplicity, we replace $F_Z(t; \alpha, d_{(\alpha)})$ by $F(t, d)$, $F_{Z,k}(t; \alpha, d_{(\alpha)})$ by $F_k(t, d)$, and $d_{(\alpha)}$ by d , in this proof.

Using the *original* binomial Chernoff bounds [20], we obtain, for $\epsilon' > 0$,

$$\begin{aligned} & \Pr(kF_k(t; d) \geq (1 + \epsilon')kF(t; d)) \\ & \leq \left(\frac{k - kF(t; d)}{k - (1 + \epsilon')kF(t; d)} \right)^{k - k(1 + \epsilon')F(t; d)} \left(\frac{kF(t; d)}{(1 + \epsilon')kF(t; d)} \right)^{(1 + \epsilon')kF(t; d)} \\ & = \left[\left(\frac{1 - F(t; d)}{1 - (1 + \epsilon')F(t; d)} \right)^{1 - (1 + \epsilon')F(t; d)} \left(\frac{1}{1 + \epsilon'} \right)^{(1 + \epsilon')F(t; d)} \right]^k, \end{aligned}$$

and for $0 < \epsilon' < 1$,

$$\begin{aligned} & \mathbf{Pr} (kF_k(t; d) \leq (1 - \epsilon')kF(t; d)) \\ & \leq \left[\left(\frac{1 - F(t; d)}{1 - (1 - \epsilon')F(t; d)} \right)^{1 - (1 - \epsilon')F(t; d)} \left(\frac{1}{1 - \epsilon'} \right)^{(1 - \epsilon')F(t; d)} \right]^k. \end{aligned}$$

Consider the general quantile estimator $\hat{d}_{(\alpha),q}$ defined in (4). For $\epsilon > 0$, (again, denote $W = q\text{-quantile}\{|S(\alpha, 1)|\}$),

$$\begin{aligned} & \mathbf{Pr} (\hat{d}_{(\alpha),q} \geq (1 + \epsilon)d) = \mathbf{Pr} (q\text{-quantile}\{|x_j|\}) \geq ((1 + \epsilon)d)^{1/\alpha} W \\ & = \mathbf{Pr} (kF_k((1 + \epsilon)^{1/\alpha} W; 1) \leq qk) = \mathbf{Pr} (kF_k(t; 1) \leq (1 - \epsilon')kF(t; 1)), \end{aligned}$$

where $t = (1 + \epsilon)^{1/\alpha} W$ and $q = (1 - \epsilon')F(t; 1)$. Thus

$$\begin{aligned} & \mathbf{Pr} (\hat{d}_{(\alpha),q} \geq (1 + \epsilon)d) \\ & \leq \left[\left(\frac{1 - F(((1 + \epsilon)^{1/\alpha} W; 1))}{1 - q} \right)^{1 - q} \left(\frac{F(((1 + \epsilon)^{1/\alpha} W; 1))}{q} \right)^q \right]^k = \exp \left(-k \frac{\epsilon^2}{G_{R,q}} \right). \end{aligned}$$

where

$$\begin{aligned} \frac{\epsilon^2}{G_{R,q}} &= -(1 - q) \log \left(1 - F \left((1 + \epsilon)^{1/\alpha} W; 1 \right) \right) \\ & - q \log \left(F \left((1 + \epsilon)^{1/\alpha} W; 1 \right) \right) + (1 - q) \log(1 - q) + q \log(q). \end{aligned}$$

For $0 < \epsilon < 1$,

$$\mathbf{Pr} (\hat{d}_{(\alpha),q} \leq (1 - \epsilon)d) = \mathbf{Pr} (kF_k((1 - \epsilon)^{1/\alpha} W; 1) \geq qk) = \mathbf{Pr} (kF_k(t; 1) \geq (1 + \epsilon')kF(t; 1)),$$

where $t = (1 - \epsilon)^{1/\alpha} W$ and $q = (1 + \epsilon')F(t; 1)$. Thus,

$$\begin{aligned} & \mathbf{Pr} (\hat{d}_{(\alpha),q} \leq (1 - \epsilon)d) \\ & \leq \left[\left(\frac{1 - F((1 - \epsilon)^{1/\alpha} W; 1)}{1 - q} \right)^{1 - q} \left(\frac{F((1 - \epsilon)^{1/\alpha} W; 1)}{q} \right)^q \right]^k = \exp \left(-k \frac{\epsilon^2}{G_{L,q}} \right), \end{aligned}$$

where

$$\begin{aligned} \frac{\epsilon^2}{G_{L,q}} &= -(1 - q) \log \left(1 - F \left((1 - \epsilon)^{1/\alpha} W; 1 \right) \right) \\ & - q \log \left(F \left((1 - \epsilon)^{1/\alpha} W; 1 \right) \right) + (1 - q) \log(1 - q) + q \log(q). \end{aligned}$$

Denote $f(t; d) = F'(t; d)$. Using L'Hospital's rule

$$\begin{aligned}
\lim_{\epsilon \rightarrow 0+} \frac{1}{G_{R,q}} &= \lim_{\epsilon \rightarrow 0+} \frac{-(1-q) \log \left(1 - F \left((1+\epsilon)^{1/\alpha} W; 1 \right) \right)}{\epsilon^2} \\
&\quad + \frac{-q \log \left(F \left((1+\epsilon)^{1/\alpha} W; 1 \right) \right) + (1-q) \log(1-q) + q \log(q)}{\epsilon^2} \\
&= \lim_{\epsilon \rightarrow 0+} \frac{f \left((1+\epsilon)^{1/\alpha} W; 1 \right) \frac{W}{\alpha} (1+\epsilon)^{1/\alpha-1}}{F \left((1+\epsilon)^{1/\alpha} W; 1 \right) \left(1 - F \left((1+\epsilon)^{1/\alpha} W; 1 \right) \right)} \times \frac{F \left((1+\epsilon)^{1/\alpha} W; 1 \right) - q}{2\epsilon} \\
&= \lim_{\epsilon \rightarrow 0+} \frac{\left(f \left((1+\epsilon)^{1/\alpha} W; 1 \right) \frac{W}{\alpha} (1+\epsilon)^{1/\alpha-1} \right)^2}{2F \left((1+\epsilon)^{1/\alpha} W; 1 \right) \left(1 - F \left((1+\epsilon)^{1/\alpha} W; 1 \right) \right)} \\
&= \frac{f^2(W; 1) W^2}{2q(1-q)\alpha^2}, \quad (q = F(W, 1)).
\end{aligned}$$

Similarly

$$\lim_{\epsilon \rightarrow 0+} G_{L,q} = \frac{2q(1-q)\alpha^2}{f^2(W; 1) W^2}.$$

To complete the proof, apply the relations on $Z = |X|$ in the proof of Lemma 1.

References

- Indyk, P.: Stable distributions, pseudorandom generators, embeddings, and data stream computation. *Journal of ACM* **53**(3) (2006) 307–323
- Li, P.: Estimators and tail bounds for dimension reduction in l_α ($0 < \alpha \leq 2$) using stable random projections. In: SODA. (2008) 10 – 19
- Li, P., Hastie, T.J.: A unified near-optimal estimator for dimension reduction in l_α ($0 < \alpha \leq 2$) using stable random projections. In: NIPS, Vancouver, BC, Canada (2008)
- Bottou, L., Chapelle, O., DeCoste, D., Weston, J., eds.: Large-Scale Kernel Machines. The MIT Press, Cambridge, MA (2007)
- Chapelle, O., Haffner, P., Vapnik, V.N.: Support vector machines for histogram-based image classification. *IEEE Trans. Neural Networks* **10**(5) (1999) 1055–1064
- Schölkopf, B., Smola, A.J.: Learning with Kernels. The MIT Press, Cambridge, MA (2002)
- Leopold, E., Kindermann, J.: Text categorization with support vector machines. how to represent texts in input space? *Machine Learning* **46**(1-3) (2002) 423–444
- Rennie, J.D., Shih, L., Teevan, J., Karger, D.R.: Tackling the poor assumptions of naive Bayes text classifiers. In: ICML, Washington, DC (2003) 616–623
- Platt, J.C.: Using analytic qp and sparseness to speed training of support vector machines. In: NIPS, Vancouver, BC, Canada (1998) 557–563
- Babcock, B., Babu, S., Datar, M., Motwani, R., Widom, J.: Models and issues in data stream systems. In: PODS, Madison, WI (2002) 1–16
- Zhao, H., Lall, A., Ogihara, M., Spatscheck, O., Wang, J., Xu, J.: A data streaming algorithm for estimating entropies of od flows. In: IMC, San Diego, CA (2007)
- David, H.A.: Order Statistics. Second edn. John Wiley & Sons, Inc., New York, NY (1981)
- Bingham, E., Mannila, H.: Random projection in dimensionality reduction: Applications to image and text data. In: KDD, San Francisco, CA (2001) 245–250
- Fradkin, D., Madigan, D.: Experiments with random projections for machine learning. In: KDD, Washington, DC (2003) 517–522
- Johnson, W.B., Lindenstrauss, J.: Extensions of Lipschitz mapping into Hilbert space. *Contemporary Mathematics* **26** (1984) 189–206
- Achlioptas, D.: Database-friendly random projections: Johnson-Lindenstrauss with binary coins. *Journal of Computer and System Sciences* **66**(4) (2003) 671–687
- Fama, E.F., Roll, R.: Parameter estimates for symmetric stable distributions. *Journal of the American Statistical Association* **66**(334) (1971) 331–338
- McCulloch, J.H.: Simple consistent estimators of stable distribution parameters. *Communications on Statistics-Simulation* **15**(4) (1986) 1109–1136
- Chernoff, H., Gastwirth, J.L., Johns, M.V.: Asymptotic distribution of linear combinations of functions of order statistics with applications to estimation. *The Annals of Mathematical Statistics* **38**(1) (1967) 52–72
- Chernoff, H.: A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *The Annals of Mathematical Statistics* **23**(4) (1952) 493–507

Computationally Efficient Estimators for Dimension Reductions Using Stable Random Projections

Ping Li

Department of Statistical Science, Cornell University, Ithaca NY 14853, USA
 {pingli}@cornell.edu

Abstract. The¹ method of *stable random projections* is a tool for efficiently computing the l_α distances using low memory, where $0 < \alpha \leq 2$ is a tuning parameter. The method boils down to a statistical estimation task and various estimators have been proposed, based on the *geometric mean*, the *harmonic mean*, and the *fractional power* etc.

This study proposes the ***optimal quantile*** estimator, whose main operation is ***selecting***, which is considerably less expensive than taking fractional power, the main operation in previous estimators. Our experiments report that the *optimal quantile* estimator is nearly one order of magnitude more computationally efficient than previous estimators. For large-scale learning tasks in which storing and computing pairwise distances is a serious bottleneck, this estimator should be desirable.

In addition to its computational advantages, the *optimal quantile* estimator exhibits nice theoretical properties. It is more accurate than previous estimators when $\alpha > 1$. We derive its theoretical error bounds and establish the explicit (i.e., no hidden constants) sample complexity bound.

1 Introduction

The method of *stable random projections*[1,2,3], as an efficient tool for computing pairwise distances in massive high-dimensional data, provides a promising mechanism to tackle some of the challenges in modern machine learning. In this paper, we provide an easy-to-implement algorithm for *stable random projections* which is both statistically accurate and computationally efficient.

1.1 Massive High-dimensional Data in Modern Machine Learning

We denote a data matrix by $\mathbf{A} \in \mathbb{R}^{n \times D}$, i.e., n data points in D dimensions. Data sets in modern applications exhibit important characteristics which impose tremendous challenges in machine learning [4]:

- Modern data sets with $n = 10^5$ or even $n = 10^6$ points are not uncommon in supervised learning, e.g., in image/text classification, ranking algorithms for search engines, etc. In the unsupervised domain (e.g., Web clustering, ads clickthroughs, word/term associations), n can be even much larger.
- Modern data sets are often of ultra high-dimensions (D), sometimes in the order of millions (or even higher), e.g., image, text, genome (e.g., SNP), etc. For example, in image analysis, D may be $10^3 \times 10^3 = 10^6$ if using pixels as features, or $D = 256^3 \approx 16$ million if using color histograms as features.
- Modern data sets are sometimes collected in a dynamic streaming fashion.

¹ First draft Feb. 2008, slightly revised in June 2008. The results were announced in January 2008 at SODA'08 when the author presented the work of [2].

- Large-scale data are often heavy-tailed, e.g., image and text data.

Some large-scale data are dense, such as image and genome data. Even for data sets which are sparse, such as text, the absolute number of non-zeros may be still large. For example, if one queries “machine learning” (a not-too-common term) in Google.com, the total number of pagehits is about 3 million. In other words, if one builds a term-doc matrix at Web scale, although the matrix is sparse, most rows will contain large numbers (e.g., millions) of non-zero entries.

1.2 Pairwise Distances in Machine Learning

Many learning algorithms require a similarity matrix computed from pairwise distances of the data matrix $\mathbf{A} \in \mathbb{R}^{n \times D}$. Examples include clustering, nearest neighbors, multidimensional scaling, and kernel SVM (support vector machines). The similarity matrix requires $O(n^2)$ storage space and $O(n^2 D)$ computing time.

This study focuses on the l_α distance ($0 < \alpha \leq 2$). Consider two vectors $u_1, u_2 \in \mathbb{R}^D$ (e.g., the leading two rows in \mathbf{A}), the l_α distance between u_1 and u_2 is

$$d_{(\alpha)} = \sum_{i=1}^D |u_{1,i} - u_{2,i}|^\alpha. \quad (1)$$

Note that, strictly speaking, the l_α distance should be defined as $d_{(\alpha)}^{1/\alpha}$. Because the power operation $(\cdot)^{1/\alpha}$ is the same for all pairs, it often makes no difference whether we use $d_{(\alpha)}^{1/\alpha}$ or just $d_{(\alpha)}$; and hence we focus on $d_{(\alpha)}$.

The radial basis kernel (e.g., for SVM) is constructed from $d_{(\alpha)}$ [5,6]:

$$\mathbf{K}(u_1, u_2) = \exp \left(-\gamma \sum_{i=1}^D |u_{1,i} - u_{2,i}|^\alpha \right), \quad 0 < \alpha \leq 2. \quad (2)$$

When $\alpha = 2$, this is the Gaussian radial basis kernel. Here α can be viewed as a *tuning* parameter. For example, in their histogram-based image classification project using SVM, [5] reported that $\alpha = 0$ and $\alpha = 0.5$ achieved good performance. For heavy-tailed data, tuning α has the similar effect as term-weighting the original data, often a critical step in a lot of applications [7,8].

For popular kernel SVM solvers including the *Sequential Minimal Optimization* (SMO) algorithm[9], storing and computing kernels is the major bottleneck. Three computational challenges were summarized in [4, page 12]:

- **Computing kernels is expensive**
- **Computing full kernel matrix is wasteful** Efficient SVM solvers often do not need to evaluate all pairwise kernels.
- **Kernel matrix does not fit in memory** Storing the kernel matrix at the memory cost $O(n^2)$ is challenging when $n > 10^5$, and is not realistic for $n > 10^6$, because $O(10^{12})$ consumes at least 1000 GBs memory.

A popular strategy in large-scale learning is to evaluate distances **on the fly**[4]. That is, instead of loading the similarity matrix in memory at the cost of $O(n^2)$, one can load the original data matrix at the cost of $O(nD)$ and recompute pairwise distances on-demand. This strategy is apparently problematic when D

is not too small. For high-dimensional data, either loading the data matrix in memory is unrealistic or computing distances on-demand becomes too expensive.

Those challenges are not unique to kernel SVM; they are general issues in distanced-based learning algorithms. The method of *stable random projections* provides a promising scheme by reducing the dimension D to a small k (e.g., $k = 50$), to facilitate compact data storage and efficient distance computations.

1.3 Stable Random Projections

The basic procedure of *stable random projections* is to multiply $\mathbf{A} \in \mathbb{R}^{n \times D}$ by a random matrix $\mathbf{R} \in \mathbb{R}^{D \times k}$ ($k \ll D$), which is generated by sampling each entry r_{ij} i.i.d. from a symmetric stable distribution $S(\alpha, 1)$. The resultant matrix $\mathbf{B} = \mathbf{A} \times \mathbf{R} \in \mathbb{R}^{n \times k}$ is much smaller than \mathbf{A} and hence it may fit in memory.

Suppose a stable random variable $x \sim S(\alpha, d)$, where d is the scale parameter. Then its characteristic function (Fourier transform of the density function) is

$$\mathbb{E}(\exp(\sqrt{-1}xt)) = \exp(-d|t|^\alpha),$$

which does not have a closed-form inverse except for $\alpha = 2$ (normal) or $\alpha = 1$ (Cauchy). Note that when $\alpha = 2$, d corresponds to “ σ^2 ” (not “ σ ”) in a normal.

Corresponding to the leading two rows in \mathbf{A} , $u_1, u_2 \in \mathbb{R}^D$, the leading two rows in \mathbf{B} are $v_1 = \mathbf{R}^T u_1$, $v_2 = \mathbf{R}^T u_2$. The entries of the difference,

$$x_j = v_{1,j} - v_{2,j} = \sum_{i=1}^D r_{ij} (u_{1,i} - u_{2,i}) \sim S\left(\alpha, d_{(\alpha)} = \sum_{i=1}^D |u_{1,i} - u_{2,i}|^\alpha\right),$$

for $j = 1$ to k , are i.i.d. samples from a stable distribution with the scale parameter being the l_α distance $d_{(\alpha)}$, due to properties of Fourier transforms. For example, when $\alpha = 2$, a weighted sum of i.i.d. standard normals is also normal with the scale parameter (i.e., variance) being the sum of squares of all weights.

Once we obtain the stable samples, one can discard the original matrix \mathbf{A} and the remaining task is to estimate the scale parameter $d_{(\alpha)}$ for each pair.

Some applications of *stable random projections* are summarized as follows:

- **Computing all pairwise distances** The cost of computing all pairwise distances of $\mathbf{A} \in \mathbb{R}^{n \times D}$, $O(n^2 D)$, is significantly reduced to $O(nDk + n^2 k)$.
- **Estimating l_α distances online** For $n > 10^5$, it is challenging or unrealistic to materialize all pairwise distances in \mathbf{A} . Thus, in applications such as online learning, databases, search engines, and online recommendation systems, it is often more efficient if we store $\mathbf{B} \in \mathbb{R}^{n \times k}$ in the memory and estimate any distance *on the fly* if needed. Estimating distances online is the standard strategy in large-scale kernel learning[4]. With *stable random projections*, this simple strategy becomes effective in high-dimensional data.
- **Learning with dynamic streaming data** In reality, the data matrix may be updated overtime. In fact, with streaming data arriving at high-rate[1,10], the “data matrix” may be never stored and hence all operations (such as clustering and classification) must be conducted on the fly. The

method of *stable random projections* provides a scheme to compute and update distances on the fly in one-pass of the data; see relevant papers (e.g., [1]) for more details on this important and fast-developing subject.

- **Estimating entropy** The entropy distance $\sum_{i=1}^D |u_{1,i} - u_{2,i}| \log |u_{1,i} - u_{2,i}|$ is a useful statistic. A workshop in NIPS'03 (www.menem.com/~ilya/pages/NIPS03) focused on entropy estimation. A recent practical algorithm is simply using the difference between the l_{α_1} and l_{α_2} distances [11], where $\alpha_1 = 1.05$, $\alpha_2 = 0.95$, and the distances were estimated by *stable random projections*.

If one tunes the l_α distances for many different α (e.g., [5]), then *stable random projections* will be even more desirable as a cost-saving device.

2 The Statistical Estimation Problem

Recall that the method of *stable random projections* boils down to a statistical estimation problem. That is, estimating the scale parameter $d_{(\alpha)}$ from k i.i.d. samples $x_j \sim S(\alpha, d_{(\alpha)})$, $j = 1$ to k . We consider that a good estimator $\hat{d}_{(\alpha)}$ should have the following desirable properties:

- (Asymptotically) unbiased and small variance.
- Computationally efficient.
- Exponential decrease of error (tail) probabilities.

The *arithmetic mean* estimator $\frac{1}{k} \sum_{j=1}^k |x_j|^2$ is good for $\alpha = 2$. When $\alpha < 2$, the task is less straightforward because (1) no explicit density of x_j exists unless $\alpha = 1$ or $0+$; and (2) $E(|x_j|^t) < \infty$ only when $-1 < t < \alpha$.

2.1 Several Previous Estimators

Initially reported in arXiv in 2006, [2] proposed the *geometric mean* estimator

$$\hat{d}_{(\alpha),gm} = \frac{\prod_{j=1}^k |x_j|^{\alpha/k}}{\left[\frac{2}{\pi} \Gamma\left(\frac{\alpha}{k}\right) \Gamma\left(1 - \frac{1}{k}\right) \sin\left(\frac{\pi}{2} \frac{\alpha}{k}\right) \right]^k}.$$

where $\Gamma(\cdot)$ is the Gamma function, and the *harmonic mean* estimator

$$\hat{d}_{(\alpha),hm} = \frac{-\frac{2}{\pi} \Gamma(-\alpha) \sin\left(\frac{\pi}{2} \alpha\right)}{\sum_{j=1}^k |x_j|^{-\alpha}} \left(k - \left(\frac{-\pi \Gamma(-2\alpha) \sin(\pi \alpha)}{[\Gamma(-\alpha) \sin\left(\frac{\pi}{2} \alpha\right)]^2} - 1 \right) \right).$$

More recently, [3] proposed the *fractional power* estimator

$$\hat{d}_{(\alpha),fp} = \left(\frac{1}{k} \frac{\sum_{j=1}^k |x_j|^{\lambda^* \alpha}}{\frac{2}{\pi} \Gamma(1 - \lambda^*) \Gamma(\lambda^* \alpha) \sin\left(\frac{\pi}{2} \lambda^* \alpha\right)} \right)^{1/\lambda^*} \times \left(1 - \frac{1}{k} \frac{1}{2\lambda^*} \left(\frac{1}{\lambda^*} - 1 \right) \left(\frac{\frac{2}{\pi} \Gamma(1 - 2\lambda^*) \Gamma(2\lambda^* \alpha) \sin(\pi \lambda^* \alpha)}{[\frac{2}{\pi} \Gamma(1 - \lambda^*) \Gamma(\lambda^* \alpha) \sin\left(\frac{\pi}{2} \lambda^* \alpha\right)]^2} - 1 \right) \right),$$

where

$$\lambda^* = \underset{-\frac{1}{2\alpha} < \lambda < \frac{1}{2}}{\operatorname{argmin}} \frac{1}{\lambda^2} \left(\frac{\frac{2}{\pi} \Gamma(1 - 2\lambda) \Gamma(2\lambda \alpha) \sin(\pi \lambda \alpha)}{[\frac{2}{\pi} \Gamma(1 - \lambda) \Gamma(\lambda \alpha) \sin\left(\frac{\pi}{2} \lambda \alpha\right)]^2} - 1 \right).$$

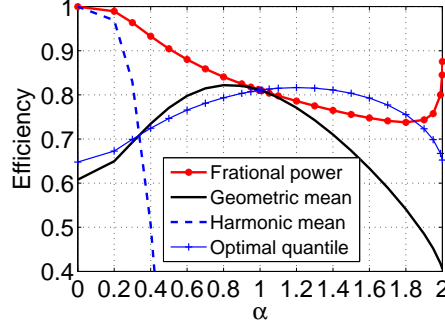


Fig. 1. The Cramér-Rao efficiencies (the higher the better, max = 100%) of various estimators, including the *optimal quantile* estimator proposed in this study.

All three estimators are unbiased or asymptotically (as $k \rightarrow \infty$) unbiased. Figure 1 compares their asymptotic variances in terms of the Cramér-Rao efficiency, which is the ratio of the smallest possible asymptotic variance over the asymptotic variance of the estimator, as $k \rightarrow \infty$.

The *geometric mean* estimator, $\hat{d}_{(\alpha),gm}$ exhibits tail bounds in exponential forms, i.e., the errors decrease exponentially fast:

$$\Pr\left(|\hat{d}_{(\alpha),gm} - d_{(\alpha)}| \geq \epsilon d_{(\alpha)}\right) \leq 2 \exp\left(-k \frac{\epsilon^2}{G_{gm}}\right).$$

The *harmonic mean* estimator, $\hat{d}_{(\alpha),hm}$, works well for small α , and has exponential tail bounds for $\alpha = 0+$.

The *fractional power* estimator, $\hat{d}_{(\alpha),fp}$, has smaller asymptotic variance than both the *geometric mean* and *harmonic mean* estimators. However, it does not have exponential tail bounds, due to the restriction $-1 < \lambda^* \alpha < \alpha$ in its definition. As shown in [3], it only has finite moments slightly higher than the 2nd order, when α approaches 2 (because $\lambda^* \rightarrow 0.5$), meaning that large errors may have a good chance to occur. We will demonstrate this by simulations.

2.2 The Issue of Computational Efficiency

In the definitions of $\hat{d}_{(\alpha),gm}$, $\hat{d}_{(\alpha),hm}$ and $\hat{d}_{(\alpha),fp}$, all three estimators require evaluating fractional powers, e.g., $|x_j|^{\alpha/k}$. This operation is relatively expensive, especially if we need to conduct this tens of billions of times (e.g., $n^2 = 10^{10}$).

For example, [5] reported that, although the radial basis kernel (2) with $\alpha = 0.5$ achieved good performance, it was not preferred because evaluating the square root was too expensive.

2.3 Our Proposed Estimator

We propose the *optimal quantile* estimator, using the q^* th smallest $|x_j|$:

$$\hat{d}_{(\alpha),oq} \propto (q^*\text{-quantile}\{|x_j|, j = 1, 2, \dots, k\})^\alpha, \quad (3)$$

where $q^* = q^*(\alpha)$ is chosen to minimize the asymptotic variance.

This estimator is computationally attractive because **selecting** should be much less expensive than evaluating fractional powers. If we are interested in $d_{(\alpha)}^{1/\alpha}$ instead, then we do not even need to evaluate any fractional powers.

As mentioned, in many cases using either $d_{(\alpha)}$ or $d_{(\alpha)}^{1/\alpha}$ makes no difference and $d_{(\alpha)}$ is often preferred because it avoids taking $(\cdot)^{1/\alpha}$ power. The radial basis kernel (2) requires $d_{(\alpha)}$. Thus this study focuses on $d_{(\alpha)}$. On the other hand, if we can estimate $d_{(\alpha)}^{1/\alpha}$ directly, for example, using (3) without the α th power, we might as well just use $d_{(\alpha)}^{1/\alpha}$ if permitted. In case we do not need to evaluate any fractional power, our estimator will be even more computationally efficient.

In addition to the computational advantages, this estimator also has good theoretical properties, in terms of both the variances and tail probabilities:

1. Figure 1 illustrates that, compared with the *geometric mean* estimator, its asymptotic variance is about the same when $\alpha < 1$, and is considerably smaller when $\alpha > 1$. Compared with the *fractional power* estimator, it has smaller asymptotic variance when $1 < \alpha \leq 1.8$. In fact, as will be shown by simulations, when the sample size k is not too large, its mean square errors are considerably smaller than the *fractional power* estimator when $\alpha > 1$.
2. The *optimal quantile* estimator exhibits tail bounds in exponential forms. This theoretical contribution is practically important, for selecting the sample size k . In learning theory, the generalization bounds are often loose. In our case, however, the bounds are tight because the distribution is specified.

The next section will be devoted to analyzing the *optimal quantile* estimator.

3 The Optimal Quantile Estimator

Recall the goal is to estimate $d_{(\alpha)}$ from $\{x_j\}_{j=1}^k$, where $x_j \sim S(\alpha, d_{(\alpha)})$, i.i.d. Since the distribution belongs to the scale family, one can estimate the scale parameter from quantiles. Due to symmetry, it is natural to consider the absolute values:

$$\hat{d}_{(\alpha),q} = \left(\frac{q\text{-Quantile}\{|x_j|, j = 1, 2, \dots, k\}}{q\text{-Quantile}\{|S(\alpha, 1)|\}} \right)^\alpha, \quad (4)$$

which is best understood by the fact that if $x \sim S(\alpha, 1)$, then $d^{1/\alpha}x \sim S(\alpha, d)$, or more obviously, if $x \sim N(0, 1)$, then $(\sigma^2)^{1/2}x \sim N(0, \sigma^2)$. By properties of order statistics [12], any q -quantile will provide an asymptotically unbiased estimator.

Lemma 1 provides the asymptotic variance of $\hat{d}_{(\alpha),q}$.

Lemma 1. Denote $f_X(x; \alpha, d_{(\alpha)})$ and $F_X(x; \alpha, d_{(\alpha)})$ the probability density function and the cumulative density function of $X \sim S(\alpha, d_{(\alpha)})$, respectively.

The asymptotic variance of $\hat{d}_{(\alpha),q}$ defined in (4) is

$$\text{Var}(\hat{d}_{(\alpha),q}) = \frac{1}{k} \frac{(q - q^2)\alpha^2/4}{f_X^2(W; \alpha, 1) W^2} d_{(\alpha)}^2 + O\left(\frac{1}{k^2}\right) \quad (5)$$

where $W = F_X^{-1}((q + 1)/2; \alpha, 1) = q\text{-Quantile}\{|S(\alpha, 1)|\}$.

Proof: See Appendix A. \square .

3.1 Optimal Quantile $q^*(\alpha)$

We choose $q = q^*(\alpha)$ so that the asymptotic variance (5) is minimized, i.e.,

$$q^*(\alpha) = \underset{q}{\operatorname{argmin}} g(q; \alpha), \quad g(q; \alpha) = \frac{q - q^2}{f_X^2(W; \alpha, 1) W^2}. \quad (6)$$

The convexity of $g(q; \alpha)$ is important. Graphically, $g(q; \alpha)$ is a convex function of q , i.e., a unique minimum exists. An algebraic proof, however, is difficult. Nevertheless, we can obtain analytical solutions when $\alpha = 1$ and $\alpha = 0+$.

Lemma 2. *When $\alpha = 1$ or $\alpha = 0+$, the function $g(q; \alpha)$ defined in (6) is a convex function of q . When $\alpha = 1$, the optimal $q^*(1) = 0.5$. When $\alpha = 0+$, $q^*(0+) = 0.203$ is the solution to $-\log q^* + 2q^* - 2 = 0$.*

Proof: See Appendix B. \square .

It is also easy to show that when $\alpha = 2$, $q^*(2) = 0.862$.

We denote the *optimal quantile estimator* by $\hat{d}_{(\alpha),oq}$, which is same as $\hat{d}_{(\alpha),q^*}$. For general α , we resort to numerical solutions, as presented in Figure 2.

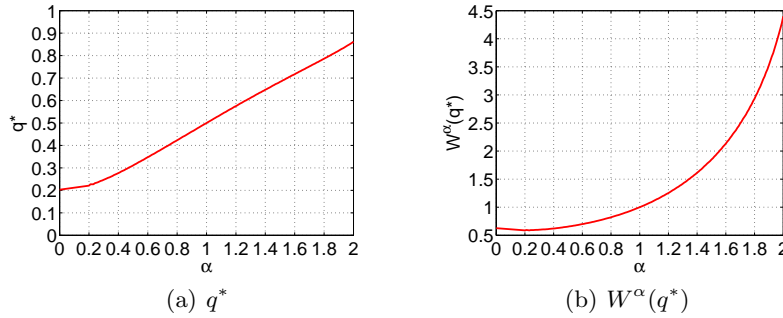


Fig. 2. (a) The optimal values for $q^*(\alpha)$, which minimizes asymptotic variance of $\hat{d}_{(\alpha),q}$, i.e., the solution to (6). (b) The constant $W^\alpha(q^*) = \{q^*\text{-quantile}\{|S(\alpha, 1)|\}\}^\alpha$.

3.2 Bias Correction

Although $\hat{d}_{(\alpha),oq}$ (i.e., $\hat{d}_{(\alpha),q^*}$) is asymptotically (as $k \rightarrow \infty$) unbiased, it is seriously biased for small k . Thus, it is practically important to remove the bias. The unbiased version of the *optimal quantile estimator* is

$$\hat{d}_{(\alpha),oq,c} = \hat{d}_{(\alpha),oq} / B_{\alpha,k}, \quad (7)$$

where $B_{\alpha,k}$ is the expectation of $\hat{d}_{(\alpha),oq}$ at $d_{(\alpha)} = 1$. For $\alpha = 1, 0+$, or 2 , we can evaluate the expectations (i.e., integrals) analytically or by numerical integrations. For general α , as the probability density is not available, the task is difficult and prone to numerical instability. On the other hand, since the Monte-Carlo simulation is a popular alternative for evaluating difficult integrals, a practical solution is to simulate the expectations, as presented in Figure 3.

Figure 3 illustrates that $B_{\alpha,k} > 1$, meaning that this correction also reduces variance while removing bias (because $\operatorname{Var}(x/c) = \operatorname{Var}(x)/c^2$). For example, when

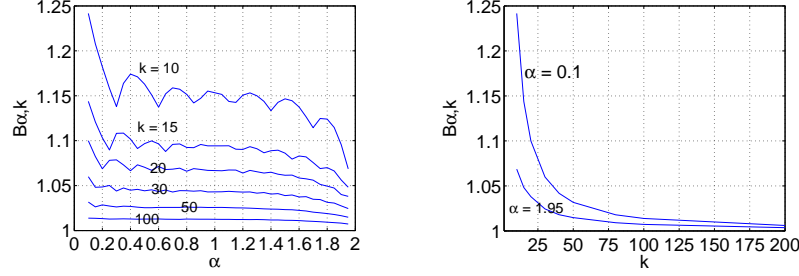


Fig. 3. The bias correction factor $B_{\alpha,k}$ in (7), obtained from 10^8 simulations for every combination of α (spaced at 0.05) and k . $B_{\alpha,k} = E(\hat{d}_{(\alpha),oq}; d_{(\alpha)} = 1)$.

$\alpha = 0.1$ and $k = 10$, $B_{\alpha,k} \approx 1.24$, which is significant, because $1.24^2 = 1.54$ implies a 54% difference in terms of variance, and even more considerable in terms of the mean square errors $MSE = \text{variance} + \text{bias}^2$.

$B_{\alpha,k}$ can be tabulated for small k , and absorbed into other coefficients, i.e., this does not increase the computational cost at run time. We fix $B_{\alpha,k}$ as reported in Figure 3. The simulations in Section 4 directly used those fixed $B_{\alpha,k}$ values.

3.3 Computational Efficiency

Figure 4 compares the computational costs of the *geometric mean*, the *fractional power*, and the *optimal quantile* estimators. The *harmonic mean* estimator was not included as it costs very similarly to the *fractional power* estimator.

We used the build-in function “pow” in gcc for evaluating the fractional powers. We implemented a “quick select” algorithm, which is similar to quick sort and requires on average linear time. For simplicity, our implementation used recursions and the middle element as pivot. Also, to ensure fairness, for all estimators, coefficients which are functions of α and/or k were pre-computed.

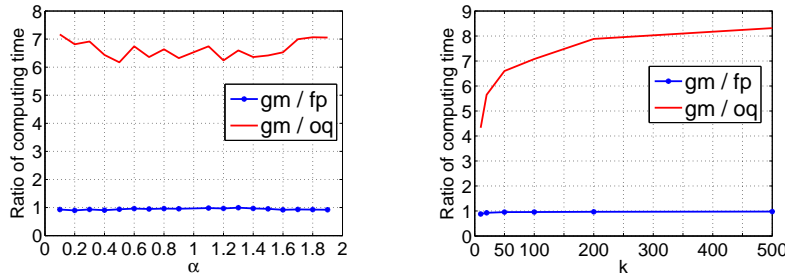


Fig. 4. Relative computational cost ($\hat{d}_{(\alpha),gm}$ over $\hat{d}_{(\alpha),oq,c}$ and $\hat{d}_{(\alpha),gm}$ over $\hat{d}_{(\alpha),fp}$), from 10^6 simulations at each combination of α and k . The left panel averages over all k and the right panel averages over all α . Note that the cost of $\hat{d}_{(\alpha),oq,c}$ includes evaluating the α th moment once.

Normalized by the computing time of $\hat{d}_{(\alpha),gm}$, we observe that relative computational efficiency does not strongly depend on α . We do observe that the ratio

of computing time of $\hat{d}_{(\alpha),gm}$ over that of $\hat{d}_{(\alpha),oq,c}$ increases consistently with increasing k . This is because in the definition of $\hat{d}_{(\alpha),oq}$ (and hence also $\hat{d}_{(\alpha),oq,c}$), it is required to evaluate the fractional power once, which contributes to the total computing time more significantly at smaller k .

Figure 4 illustrates that, (A) the *geometric mean* estimator and the *fractional power* estimator are similar in terms of computational efficiency; (B) the *optimal quantile* estimator is nearly one order of magnitude more computationally efficient than the *geometric mean* and *fractional power* estimators. Because we implemented a “naïve” version of “quick select” using recursions and simple pivoting, the actual improvement may be more significant. Also, if applications require only $d_{(\alpha)}^{1/\alpha}$, then no fractional power operations are needed for $\hat{d}_{(\alpha),oq,c}$ and the improvement will be even more considerable.

3.4 Error (Tail) Bounds

Error (tail) bounds are essential for determining k . The variance alone is not sufficient for that purpose. If an estimator of d , say \hat{d} , is normally distributed, $\hat{d} \sim N(d, \frac{1}{k}V)$, the variance suffices for choosing k because its error (tail) probability $\Pr(|\hat{d} - d| \geq \epsilon d) \leq 2 \exp\left(-k \frac{\epsilon^2}{2V}\right)$ is determined by V . In general, a reasonable estimator will be asymptotically normal, for small enough ϵ and large enough k . For a finite k and a fixed ϵ , however, the normal approximation may be (very) poor. This is especially true for the *fractional power* estimator, $\hat{d}_{(\alpha),fp}$.

Thus, for a good motivation, Lemma 3 provides the error (tail) probability bounds of $\hat{d}_{(\alpha),q}$ for any q , not just the optimal quantile q^* .

Lemma 3. Denote $X \sim S(\alpha, d_{(\alpha)})$ and its probability density function by $f_X(x; \alpha, d_{(\alpha)})$ and cumulative function by $F_X(x; \alpha, d_{(\alpha)})$. Given $x_j \sim S(\alpha, d_{(\alpha)})$, i.i.d., $j = 1$ to k . Using $\hat{d}_{(\alpha),q}$ in (4), then

$$\Pr\left(\hat{d}_{(\alpha),q} \geq (1 + \epsilon)d_{(\alpha)}\right) \leq \exp\left(-k \frac{\epsilon^2}{G_{R,q}}\right), \epsilon > 0, \quad (8)$$

$$\Pr\left(\hat{d}_{(\alpha),q} \leq (1 - \epsilon)d_{(\alpha)}\right) \leq \exp\left(-k \frac{\epsilon^2}{G_{L,q}}\right), 0 < \epsilon < 1, \quad (9)$$

$$\frac{\epsilon^2}{G_{R,q}} = -(1 - q) \log(2 - 2F_R) - q \log(2F_R - 1) + (1 - q) \log(1 - q) + q \log q, \quad (10)$$

$$\frac{\epsilon^2}{G_{L,q}} = -(1 - q) \log(2 - 2F_L) - q \log(2F_L - 1) + (1 - q) \log(1 - q) + q \log q, \quad (11)$$

$$W = F_X^{-1}((q + 1)/2; \alpha, 1) = q\text{-quantile}\{|S(\alpha, 1)|\},$$

$$F_R = F_X\left((1 + \epsilon)^{1/\alpha} W; \alpha, 1\right), \quad F_L = F_X\left((1 - \epsilon)^{1/\alpha} W; \alpha, 1\right).$$

As $\epsilon \rightarrow 0+$

$$\lim_{\epsilon \rightarrow 0+} G_{R,q} = \lim_{\epsilon \rightarrow 0+} G_{L,q} = \frac{q(1 - q)\alpha^2/2}{f_X^2(W; \alpha, 1)W^2}. \quad (12)$$

Proof: See Appendix C. \square

The limit in (12) as $\epsilon \rightarrow 0$ is precisely twice the asymptotic variance factor of $\hat{d}_{(\alpha),q}$ in (5), consistent with the normality approximation mentioned previously. This explains why we express the constants as ϵ^2/G . (12) also indicates that the tail bounds achieve the “optimal rate” for this estimator, in the language of large deviation theory.

By the Bonferroni bound, it is easy to determine the sample size k

$$\Pr\left(|\hat{d}_{(\alpha),q} - d_{(\alpha)}| \geq \epsilon d_{(\alpha)}\right) \leq 2 \exp\left(-k \frac{\epsilon^2}{G}\right) \leq \delta/(n^2/2) \implies k \geq \frac{G}{\epsilon^2} (2 \log n - \log \delta).$$

Lemma 4. Using $\hat{d}_{(\alpha),q}$ with $k \geq \frac{G}{\epsilon^2} (2 \log n - \log \delta)$, any pairwise l_α distance among n points can be approximated within a $1 \pm \epsilon$ factor with probability $\geq 1 - \delta$. It suffices to let $G = \max\{G_{R,q}, G_{L,q}\}$, where $G_{R,q}, G_{L,q}$ are defined in Lemma 3.

The Bonferroni bound can be unnecessarily conservative. It is often reasonable to replace $\delta/(n^2/2)$ by δ/T , meaning that except for a $1/T$ fraction of pairs, any distance can be approximated within a $1 \pm \epsilon$ factor with probability $1 - \delta$.

Figure 5 plots the error bound constants for $\epsilon < 1$, for both the recommended *optimal quantile* estimator $\hat{d}_{(\alpha),oq}$ and the baseline *sample median* estimator $\hat{d}_{(\alpha),q=0.5}$. Although we choose $\hat{d}_{(\alpha),oq}$ based on the asymptotic variance, it turns out $\hat{d}_{(\alpha),oq}$ also exhibits (much) better tail behaviors (i.e., smaller constants) than $\hat{d}_{(\alpha),q=0.5}$, at least in the range of $\epsilon < 1$.

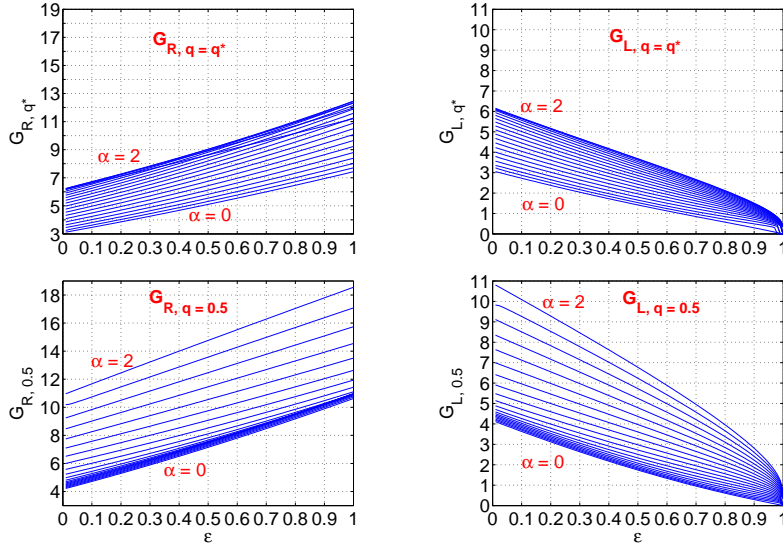


Fig. 5. Tail bound constants for quantile estimators; the lower the better. Upper panels: optimal quantile estimators $\hat{d}_{(\alpha),q^*}$. Lower panels: median estimators $\hat{d}_{(\alpha),q=0.5}$.

Consider $k = \frac{G}{\epsilon^2} (\log 2T - \log \delta)$ (recall we suggest replacing $n^2/2$ by T), with $\delta = 0.05$, $\epsilon = 0.5$, and $T = 10$. Because $G_{R,q^*} \approx 5 \sim 9$ around $\epsilon = 0.5$, we obtain

$k \approx 120 \sim 215$, which is still a relatively large number (although the original dimension D might be 10^6). If we choose $\epsilon = 1$, then approximately $k \approx 40 \sim 65$.

It is possible $k = 120 \sim 215$ might be still conservative, for three reasons: (A) the tail bounds, although “sharp,” are still upper bounds; (B) using $G = \max\{G_{R,q^*}, G_{L,q^*}\}$ is conservative because G_{L,q^*} is usually much smaller than G_{R,q^*} ; (C) this type of tail bounds is based on relative error, which may be stringent for small (≈ 0) distances.

In fact, some earlier studies on *normal random projections* (i.e., $\alpha = 2$) [13,14] empirically demonstrated that $k \geq 50$ appeared sufficient.

4 Simulations

We resort to simulations for comparing the finite sample variances of various estimators and assessing the more precise error (tail) probabilities.

One advantage of *stable random projections* is that we know the (manually generated) distributions and the only source of errors is from the random number generations. Thus, we can simply rely on simulations to evaluate the estimators without using real data. In fact, after projections, the projected data follow exactly the stable distribution, regardless of the original real data distribution.

Without loss of generality, we simulate samples from $S(\alpha, 1)$ and estimate the scale parameter (i.e., 1) from the samples. Repeating the procedure 10^7 times, we can reliably evaluate the mean square errors (MSE) and tail probabilities.

4.1 Mean Square Errors (MSE)

As illustrated in Figure 6, in terms of the MSE, the *optimal quantile* estimator $\hat{d}_{(\alpha),oq,c}$ outperforms both the *geometric mean* and *fractional power* estimators when $\alpha > 1$ and $k \geq 20$. The *fractional power* estimator does not appear to be very suitable for $\alpha > 1$, especially for α close to 2, even when the sample size k is not too small (e.g., $k = 50$). For $\alpha < 1$, however, the *fractional power* estimator has good performance in terms of MSE, even for small k .

4.2 Error(Tail) Probabilities

Figure 7 presents the simulated right tail probabilities, $\Pr(\hat{d}_{(\alpha)} \geq (1 + \epsilon)d_{(\alpha)})$, illustrating that when $\alpha > 1$, the *fractional power* estimator can exhibit very bad tail behaviors. For $\alpha < 1$, the *fractional power* estimator demonstrates good performance at least for the probability range in the simulations.

Thus, Figure 7 demonstrates that the *optimal quantile* estimator consistently outperforms the *fractional power* and the *geometric mean* estimators when $\alpha > 1$.

5 The Related Work

There have been many studies of *normal random projections* in machine learning, for dimension reduction in the l_2 norm, e.g., [14], highlighted by the Johnson-Lindenstrauss (JL) Lemma [15], which says $k = O(\log n / \epsilon^2)$ suffices when using normal (or normal-like, e.g., [16]) projection methods.

The method of *stable random projections* is applicable for computing the l_α distances ($0 < \alpha \leq 2$), not just for l_2 . [1, Lemma 1, Lemma 2, Theorem 3]

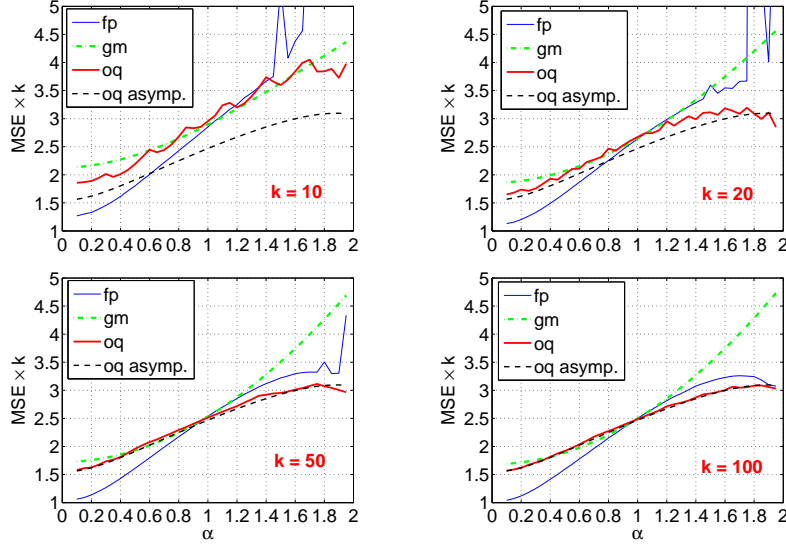


Fig. 6. Empirical mean square errors (MSE, the lower the better), from 10^7 simulations at every combination of α and k . The values are multiplied by k so that four plots can be at about the same scale. The MSE for the *geometric mean* (gm) estimator is computed exactly since closed-form expression exists. The lower dashed curves are the asymptotic variances of the *optimal quantile* (oq) estimator.

suggested the *median* (i.e., $q = 0.5$ quantile) estimator for $\alpha = 1$ and argued that the sample complexity bound should be $O(1/\epsilon^2)$ ($n = 1$ in their study). Their bound was not provided in an explicit form and required an “ ϵ is small enough” argument. For $\alpha \neq 1$, [1, Lemma 4] only provided a conceptual algorithm, which “is not uniform.” In this study, we prove the bounds for any q -quantile and any $0 < \alpha \leq 2$ (not just $\alpha = 1$), in explicit exponential forms, with no unknown constants and no restriction that “ ϵ is small enough.”

The quantile estimator for stable distributions was proposed in statistics quite some time ago, e.g., [17,18]. [17] mainly focused on $1 \leq \alpha \leq 2$ and recommended using $q = 0.44$ quantiles (mainly for the sake of smaller bias). [18] focused on $0.6 \leq \alpha \leq 2$ and recommended $q = 0.5$ quantiles.

This study considers all $0 < \alpha \leq 2$ and recommends q based on the minimum asymptotic variance. Because the bias can be easily removed (at least in the practical sense), it appears not necessary to use other quantiles only for the sake of smaller bias. Tail bounds, which are useful for choosing q and k based on confidence intervals, were not available in [17,18].

Finally, one might ask if there might be better estimators. For $\alpha = 1$, [19] proposed using a linear combination of quantiles (with carefully chosen coefficients) to obtain an asymptotically optimal estimator for the Cauchy scale parameter. While it is possible to extend their result to general $0 < \alpha < 2$ (requiring some non-trivial work), whether or not it will be practically better than the *optimal*

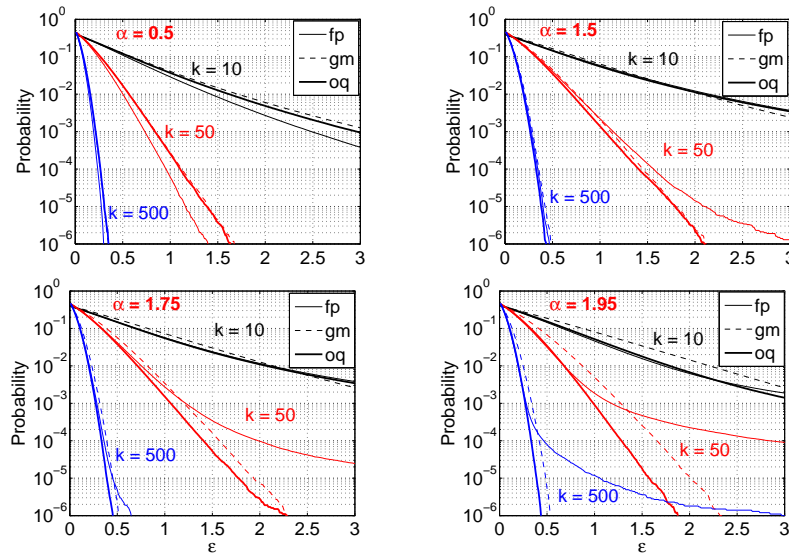


Fig. 7. The right tail probabilities (the lower the better), from 10^7 simulations at each combination of α and k .

quantile estimator is unclear because the extreme quantiles severely affect the tail probabilities and finite-sample variances and hence some kind of truncation (i.e., discarding some samples at extreme quantiles) is necessary. Also, exponential tail bounds of the linear combination of quantiles may not exist or may not be feasible to derive. In addition, the *optimal quantile* estimator is computationally more efficient.

6 Conclusion

Many machine learning algorithms operate on the training data only through pairwise distances. Computing, storing, updating and retrieving the “matrix” of pairwise distances is challenging in applications involving massive, high-dimensional, and possibly streaming, data. For example, the pairwise distance matrix can not fit in memory when the number of observations exceeds 10^6 (or even 10^5).

The method of *stable random projections* provides an efficient mechanism for computing pairwise distances using low memory, by transforming the original high-dimensional data into *sketches*, i.e., a small number of samples from α -stable distributions, which are much easier to store and retrieve.

This method provides a uniform scheme for computing the l_α pairwise distances for all $0 < \alpha \leq 2$. Choosing an appropriate α is often critical to the performance of learning algorithms. In principle, we can tune algorithms for many l_α distances; and *stable random projections* can provide an efficient tool.

To recover the original distances, we face an estimation task. Compared with previous estimators based on the *geometric mean*, the *harmonic mean*, or the *fractional power*, the proposed *optimal quantile* estimator exhibits two advantages. Firstly, the *optimal quantile* estimator is nearly one order of magnitude

more efficient than other estimators (e.g., reducing the training time from one week to one day). Secondly, the *optimal quantile* estimator is considerably more accurate when $\alpha > 1$, in terms of both the variances and error (tail) probabilities. Note that $\alpha \geq 1$ corresponds to a convex norm (satisfying the triangle inequality), which might be another motivation for using l_α distances with $\alpha \geq 1$.

One theoretical contribution is the explicit tail bounds for general quantile estimators and consequently the sample complexity bound $k = O(\log n/\epsilon^2)$. Those bounds may guide practitioners in choosing k , the number of projections. We should mention that

While we believe that applying *stable random projections* in machine learning has become straightforward, there are interesting theoretical issues for future research. For example, how theoretical properties of learning algorithms may be affected if the approximated (instead of exact) l_α distances are used?

A Proof of Lemma 1

Denote $f_X(x; \alpha, d_{(\alpha)})$ and $F_X(x; \alpha, d_{(\alpha)})$ the probability density function and the cumulative density function of $X \sim S(\alpha, d_{(\alpha)})$, respectively. Similarly we use $f_Z(z; \alpha, d_{(\alpha)})$ and $F_Z(z; \alpha, d_{(\alpha)})$ for $Z = |X|$. Due to symmetry, the following relations hold

$$f_Z(z; \alpha, d_{(\alpha)}) = 2f_X(z; \alpha, d_{(\alpha)}) = 2/d_{(\alpha)}^{1/\alpha} f_X(z/d_{(\alpha)}^{1/\alpha}; \alpha, 1),$$

$$F_Z(z; \alpha, d_{(\alpha)}) = 2F_X(z; \alpha, d_{(\alpha)}) - 1 = 2F_X(z/d_{(\alpha)}^{1/\alpha}; \alpha, 1) - 1,$$

$$F_Z^{-1}(q; \alpha, d_{(\alpha)}) = F_X^{-1}((q+1)/2; \alpha, d_{(\alpha)}) = d_{(\alpha)}^{1/\alpha} F_X^{-1}((q+1)/2; \alpha, 1).$$

Let $W = q\text{-Quantile}\{S(\alpha, 1)\} = F_X^{-1}((q+1)/2; \alpha, 1)$ and $W_d = F_Z^{-1}(q; \alpha, d_{(\alpha)}) = d_{(\alpha)}^{1/\alpha} W$. Then, following known statistical results, e.g., [12, Theorem 9.2], the asymptotic variance of $\hat{d}_{\alpha, q}^{1/\alpha}$ should be

$$\begin{aligned} \text{Var}\left(\hat{d}_{\alpha, q}^{1/\alpha}\right) &= \frac{1}{k} \frac{q - q^2}{f_Z^2(W_d; \alpha, d_{(\alpha)}) W^2} + O\left(\frac{1}{k^2}\right) = \frac{1}{k} \frac{q - q^2}{d_{(\alpha)}^{-2/\alpha} f_Z^2(W; \alpha, 1) W^2} + O\left(\frac{1}{k^2}\right) \\ &= \frac{1}{k} \frac{q - q^2}{4d_{(\alpha)}^{-2/\alpha} f_X^2(W; \alpha, 1) W^2} + O\left(\frac{1}{k^2}\right). \end{aligned}$$

By ‘‘delta method,’’ i.e., $\text{Var}(h(x)) \approx \text{Var}(x) (h'(x))^2$,

$$\text{Var}\left(\hat{d}_{\alpha, q}\right) = \text{Var}\left(\hat{d}_{\alpha, q}^{1/\alpha}\right) \left(\alpha d_{(\alpha)}^{(\alpha-1)/\alpha}\right)^2 + O\left(\frac{1}{k^2}\right) = \frac{1}{k} \frac{(q - q^2)\alpha^2/4}{f_X^2(W; \alpha, 1) W^2} d_{(\alpha)}^2 + O\left(\frac{1}{k^2}\right).$$

B Proof of Lemma 2

First, consider $\alpha = 1$. In this case,

$$\begin{aligned} f_X(x; 1, 1) &= \frac{1}{\pi} \frac{1}{x^2 + 1}, \quad W = F_X^{-1}((q+1)/2; 1, 1) = \tan\left(\frac{\pi}{2}q\right), \\ g(q; 1) &= \frac{q - q^2}{\left(\frac{2}{\pi} \frac{1}{\tan^2(\frac{\pi}{2}q) + 1}\right)^2 \tan^2(\frac{\pi}{2}q)} = \frac{q - q^2}{\sin^2(\pi q)} \pi^2. \end{aligned}$$

It suffices to study $L(q) = \log g(q; 1)$.

$$L'(q) = \frac{1}{q} - \frac{1}{1-q} - \frac{2\pi \cos(\pi q)}{\sin(\pi q)}, \quad L''(q) = -\frac{1}{q^2} - \frac{1}{(1-q)^2} + \frac{2\pi^2}{\sin^2(\pi q)}.$$

Because $\sin(x) \leq x$ for $x \geq 0$, it is easy to see that $\frac{\pi}{\sin(\pi q)} - \frac{1}{q} \geq 0$, and $\frac{\pi}{\sin(\pi q)} - \frac{1}{1-q} = \frac{\pi}{\sin(\pi(1-q))} - \frac{1}{1-q} \geq 0$. Thus, $L'' \geq 0$, i.e., $L(q)$ is convex and so is $g(q; 1) = e^{L(q)}$. Since $L'(1/2) = 0$, we know $q^*(1) = 0.5$.

Next we consider $\alpha = 0+$, using a fact [2] that as $\alpha \rightarrow 0+$, $|S(\alpha, 1)|^\alpha$ converges to $1/E_1$, where E_1 stands for an exponential distribution with mean 1.

Denote $h = d_{(0+)}$ and $z_j \sim h/E_1$. The sample quantile estimator becomes

$$\hat{d}_{(0+),q} = \frac{q\text{-Quantile}\{|z_j|, j = 1, 2, \dots, k\}}{q\text{-Quantile}\{1/E_1\}}.$$

In this case,

$$f_Z(z; h) = e^{-h/z} \frac{h}{z^2}, \quad F_Z^{-1}(q; h) = -\frac{h}{\log q},$$

$$\text{Var}(\hat{d}_{(0+),q}) = \frac{1}{k} \frac{1-q}{q \log^2 q} h^2 + O\left(\frac{1}{k^2}\right).$$

It is straightforward to show that $\frac{1-q}{q \log^2 q}$ is a convex function of q and the minimum is attained by solving $-\log q^* + 2q^* - 2 = 0$, i.e., $q^* = 0.203$.

C Proof of Lemma 3

Given k i.i.d. samples, $x_j \sim S(\alpha, d_{(\alpha)})$, $j = 1$ to k . Let $z_j = |x_j|$, $j = 1$ to k . Denote by $F_Z(t; \alpha, d_{(\alpha)})$ the cumulative density of z_j , and by $F_{Z,k}(t; \alpha, d_{(\alpha)})$ the empirical cumulative density of z_j , $j = 1$ to k .

It is the basic fact[12] about order statistics that $kF_{Z,k}(t; \alpha, d_{(\alpha)})$ follows a binomial, i.e., $kF_{Z,k}(t; \alpha, d_{(\alpha)}) \sim \text{Bin}(k, F_Z(t; \alpha, d_{(\alpha)}))$. For simplicity, we replace $F_Z(t; \alpha, d_{(\alpha)})$ by $F(t, d)$, $F_{Z,k}(t; \alpha, d_{(\alpha)})$ by $F_k(t, d)$, and $d_{(\alpha)}$ by d , in this proof.

Using the *original* binomial Chernoff bounds [20], we obtain, for $\epsilon' > 0$,

$$\begin{aligned} & \Pr(kF_k(t; d) \geq (1 + \epsilon')kF(t; d)) \\ & \leq \left(\frac{k - kF(t; d)}{k - (1 + \epsilon')kF(t; d)} \right)^{k - k(1 + \epsilon')F(t; d)} \left(\frac{kF(t; d)}{(1 + \epsilon')kF(t; d)} \right)^{(1 + \epsilon')kF(t; d)} \\ & = \left[\left(\frac{1 - F(t; d)}{1 - (1 + \epsilon')F(t; d)} \right)^{1 - (1 + \epsilon')F(t; d)} \left(\frac{1}{1 + \epsilon'} \right)^{(1 + \epsilon')F(t; d)} \right]^k, \end{aligned}$$

and for $0 < \epsilon' < 1$,

$$\begin{aligned} & \Pr(kF_k(t; d) \leq (1 - \epsilon')kF(t; d)) \\ & \leq \left[\left(\frac{1 - F(t; d)}{1 - (1 - \epsilon')F(t; d)} \right)^{1 - (1 - \epsilon')F(t; d)} \left(\frac{1}{1 - \epsilon'} \right)^{(1 - \epsilon')F(t; d)} \right]^k. \end{aligned}$$

Consider the general quantile estimator $\hat{d}_{(\alpha),q}$ defined in (4). For $\epsilon > 0$, (again, denote $W = q\text{-quantile}\{|S(\alpha, 1)|\}$),

$$\begin{aligned} \Pr(\hat{d}_{(\alpha),q} \geq (1+\epsilon)d) &= \Pr(q\text{-quantile}\{|x_j|\} \geq ((1+\epsilon)d)^{1/\alpha} W) \\ &= \Pr(kF_k((1+\epsilon)^{1/\alpha} W; 1) \leq qk) = \Pr(kF_k(t; 1) \leq (1-\epsilon')kF(t; 1)), \end{aligned}$$

where $t = (1+\epsilon)^{1/\alpha} W$ and $q = (1-\epsilon')F(t; 1)$. Thus

$$\begin{aligned} &\Pr(\hat{d}_{(\alpha),q} \geq (1+\epsilon)d) \\ &\leq \left[\left(\frac{1 - F((1+\epsilon)^{1/\alpha} W; 1)}{1-q} \right)^{1-q} \left(\frac{F((1+\epsilon)^{1/\alpha} W; 1)}{q} \right)^q \right]^k = \exp\left(-k \frac{\epsilon^2}{G_{R,q}}\right). \end{aligned}$$

where

$$\begin{aligned} \frac{\epsilon^2}{G_{R,q}} &= -(1-q) \log\left(1 - F((1+\epsilon)^{1/\alpha} W; 1)\right) \\ &\quad - q \log\left(F((1+\epsilon)^{1/\alpha} W; 1)\right) + (1-q) \log(1-q) + q \log(q). \end{aligned}$$

For $0 < \epsilon < 1$,

$$\Pr(\hat{d}_{(\alpha),q} \leq (1-\epsilon)d) = \Pr(kF_k((1-\epsilon)^{1/\alpha} W; 1) \geq qk) = \Pr(kF_k(t; 1) \geq (1+\epsilon')kF(t; 1)),$$

where $t = (1-\epsilon)^{1/\alpha} W$ and $q = (1+\epsilon')F(t; 1)$. Thus,

$$\begin{aligned} &\Pr(\hat{d}_{(\alpha),q} \leq (1-\epsilon)d) \\ &\leq \left[\left(\frac{1 - F((1-\epsilon)^{1/\alpha} W; 1)}{1-q} \right)^{1-q} \left(\frac{F((1-\epsilon)^{1/\alpha} W; 1)}{q} \right)^q \right]^k = \exp\left(-k \frac{\epsilon^2}{G_{L,q}}\right), \end{aligned}$$

where

$$\begin{aligned} \frac{\epsilon^2}{G_{L,q}} &= -(1-q) \log\left(1 - F((1-\epsilon)^{1/\alpha} W; 1)\right) \\ &\quad - q \log\left(F((1-\epsilon)^{1/\alpha} W; 1)\right) + (1-q) \log(1-q) + q \log(q). \end{aligned}$$

Denote $f(t; d) = F'(t; d)$. Using L'Hospital's rule

$$\begin{aligned} \lim_{\epsilon \rightarrow 0+} \frac{1}{G_{R,q}} &= \lim_{\epsilon \rightarrow 0+} \frac{-(1-q) \log\left(1 - F((1+\epsilon)^{1/\alpha} W; 1)\right)}{\epsilon^2} \\ &\quad + \frac{-q \log\left(F((1+\epsilon)^{1/\alpha} W; 1)\right) + (1-q) \log(1-q) + q \log(q)}{\epsilon^2} \\ &= \lim_{\epsilon \rightarrow 0+} \frac{f((1+\epsilon)^{1/\alpha} W; 1) \frac{W}{\alpha} (1+\epsilon)^{1/\alpha-1}}{F((1+\epsilon)^{1/\alpha} W; 1) (1 - F((1+\epsilon)^{1/\alpha} W; 1))} \times \frac{F((1+\epsilon)^{1/\alpha} W; 1) - q}{2\epsilon} \\ &= \lim_{\epsilon \rightarrow 0+} \frac{\left(f((1+\epsilon)^{1/\alpha} W; 1) \frac{W}{\alpha} (1+\epsilon)^{1/\alpha-1}\right)^2}{2F((1+\epsilon)^{1/\alpha} W; 1) (1 - F((1+\epsilon)^{1/\alpha} W; 1))} \\ &= \frac{f^2(W; 1) W^2}{2q(1-q)\alpha^2}, \quad (q = F(W, 1)). \end{aligned}$$

Similarly

$$\lim_{\epsilon \rightarrow 0+} G_{L,q} = \frac{2q(1-q)\alpha^2}{f^2(W; 1) W^2}.$$

To complete the proof, apply the relations on $Z = |X|$ in the proof of Lemma 1.

References

1. Indyk, P.: Stable distributions, pseudorandom generators, embeddings, and data stream computation. *Journal of ACM* **53**(3) (2006) 307–323
2. Li, P.: Estimators and tail bounds for dimension reduction in l_α ($0 < \alpha \leq 2$) using stable random projections. In: SODA. (2008) 10 – 19
3. Li, P., Hastie, T.J.: A unified near-optimal estimator for dimension reduction in l_α ($0 < \alpha \leq 2$) using stable random projections. In: NIPS, Vancouver, BC, Canada (2008)
4. Bottou, L., Chapelle, O., DeCoste, D., Weston, J., eds.: *Large-Scale Kernel Machines*. The MIT Press, Cambridge, MA (2007)
5. Chapelle, O., Haffner, P., Vapnik, V.N.: Support vector machines for histogram-based image classification. *IEEE Trans. Neural Networks* **10**(5) (1999) 1055–1064
6. Schölkopf, B., Smola, A.J.: *Learning with Kernels*. The MIT Press, Cambridge, MA (2002)
7. Leopold, E., Kindermann, J.: Text categorization with support vector machines. how to represent texts in input space? *Machine Learning* **46**(1-3) (2002) 423–444
8. Rennie, J.D., Shih, L., Teevan, J., Karger, D.R.: Tackling the poor assumptions of naive Bayes text classifiers. In: ICML, Washington, DC (2003) 616–623
9. Platt, J.C.: Using analytic qp and sparseness to speed training of support vector machines. In: NIPS, Vancouver, BC, Canada (1998) 557–563
10. Babcock, B., Babu, S., Datar, M., Motwani, R., Widom, J.: Models and issues in data stream systems. In: PODS, Madison, WI (2002) 1–16
11. Zhao, H., Lall, A., Ogihara, M., Spatscheck, O., Wang, J., Xu, J.: A data streaming algorithm for estimating entropies of od flows. In: IMC, San Diego, CA (2007)
12. David, H.A.: *Order Statistics*. Second edn. John Wiley & Sons, Inc., New York, NY (1981)
13. Bingham, E., Mannila, H.: Random projection in dimensionality reduction: Applications to image and text data. In: KDD, San Francisco, CA (2001) 245–250
14. Fradkin, D., Madigan, D.: Experiments with random projections for machine learning. In: KDD, Washington, DC (2003) 517–522
15. Johnson, W.B., Lindenstrauss, J.: Extensions of Lipschitz mapping into Hilbert space. *Contemporary Mathematics* **26** (1984) 189–206
16. Achlioptas, D.: Database-friendly random projections: Johnson-Lindenstrauss with binary coins. *Journal of Computer and System Sciences* **66**(4) (2003) 671–687
17. Fama, E.F., Roll, R.: Parameter estimates for symmetric stable distributions. *Journal of the American Statistical Association* **66**(334) (1971) 331–338
18. McCulloch, J.H.: Simple consistent estimators of stable distribution parameters. *Communications on Statistics-Simulation* **15**(4) (1986) 1109–1136
19. Chernoff, H., Gastwirth, J.L., Johns, M.V.: Asymptotic distribution of linear combinations of functions of order statistics with applications to estimation. *The Annals of Mathematical Statistics* **38**(1) (1967) 52–72
20. Chernoff, H.: A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *The Annals of Mathematical Statistics* **23**(4) (1952) 493–507