Randomization in Numerical Linear Algebra

Petros Drineas

Department of Computer Science Purdue University

For slides, etc. google drineas



<u>**Randomization and sampling</u>** allow us to design provably accurate algorithms for problems that are:</u>

> Massive

(matrices so large that can not be stored at all, or can only be stored in slow memory devices)

Computationally expensive or NP-hard

(combinatorial optimization problems, such as the Column Subset Selection Problem)

RandNLA in a slide

Randomized algorithms

• By (carefully) sampling rows/columns of a matrix, we can construct new, smaller matrices that are close to the original matrix (w.r.t. matrix norms) with high probability.

Here C consists of a few (rescaled) columns of A and R consists of the corresponding (rescaled) rows of B.

RandNLA in a slide

Randomized algorithms

• By preprocessing the matrix using almost any random matrix, we can sample rows/columns much less carefully (uniformly at random) and still get nice bounds with high probability.

This is equivalent to setting C = AX, where X is (say) a random Gaussian matrix and the number of columns in X is much smaller than the number of columns in A; similarly, $R = X^T B$.

RandNLA in a slide

Randomized algorithms

• By (carefully) sampling rows/columns of a matrix, we can construct new, smaller matrices that are close to the original matrix (w.r.t. matrix norms) with high probability.

• By preprocessing the matrix using "random projection" matrices, we can sample rows/columns much less carefully (uniformly at random) and still get nice bounds with high probability.

Matrix perturbation theory

• The resulting "sketches" behave similarly (e.g., in terms of singular values and singular vectors) to the original matrices thanks to the norm bounds.



Applications in **BIG DATA**

(Data Mining, Information Retrieval, Machine Learning, Bioinformatics, etc.)



Roadmap (Drineas & Mahoney, Lectures on RandNLA, Vol. 25, Amer. Math. Soc., 2018)

- > RandNLA approaches for regression problems
- RandNLA approaches for Principal Component Analysis (PCA)

(was also discussed in Andreas Stathopoulos' talk on Monday)

Roadmap

(Drineas & Mahoney, Lectures on RandNLA, Vol. 25, Amer. Math. Soc., 2018)

- RandNLA approaches for regression problems
- RandNLA approaches for Principal Component Analysis (PCA)

(was also discussed in Andreas Stathopoulos' talk on Monday)

Why regression and PCA?

Both problems are of paramount importance in Big (as well as in Tiny, Small, Medium, Massive, etc.) Data analysis.

Both problems are at the heart of multiple disciplines: Computer Science (Numerical Linear Algebra, Machine Learning), Applied Mathematics, and Statistics.

Both problems have a very rich history:

- Regression was introduced in the 1800s (Gauss, Legendre, etc.)
- > PCA was introduced in the 1900s (Pearson, Hotelling, etc.)

Problem definition and motivation

In data analysis applications one has <u>n observations</u> of the form:

$$y_i = y(t_i), i = 1, \dots, n$$

Model y(t) (unknown) as a linear combination of <u>d</u> basis functions:

$$y(t) \approx x_1 \phi_1(t) + \dots + x_d \phi_d(t)$$

A is an n-by-d "design matrix" (n >> d):

$$A_{ij} = \phi_j(t_i)$$

In matrix-vector notation,

$$y \approx Ax$$

Least-norm approximation problems

The linear measurement model:

$$y = Ax + \varepsilon \quad \begin{cases} y \text{ are the measurements} \\ x \text{ is the unknown} \\ \varepsilon \text{ is an error process} \end{cases}$$

In order to estimate x, solve:

$$\widehat{x} = \arg\min\|y - Ax\|$$

Application: data analysis in science

• First application: Astronomy

Predicting the orbit of the asteroid Ceres (in 1801!). Gauss (1809) -- see also Legendre (1805) and Adrain (1808). First application of "least squares optimization" and runs in O(nd²) time!

• Data analysis: Fit parameters of a biological, chemical, economical, physical, astronomical, social, internet, etc. model to experimental data.

Norms of common interest

Let y = b and define the residual: $r = Ax - b \in R^n$

Least-squares approximation:

minimize:
$$||Ax - b||_2^2 = r_1^2 + r_2^2 + \dots + r_n^2$$

Chebyshev or mini-max approximation:

minimize:
$$||Ax - b||_{\infty} = \max\{|r_1|, \dots, |r_n|\}$$

Sum of absolute residuals approximation:

minimize:
$$||Ax - b||_1 = |r_1| + |r_2| + \dots + |r_n|$$



We are interested in over-constrained least-squares problems, $n \gg d$.

We will briefly discuss under-constrained (n << d) and square (n \approx d) problems later. Typically, there is no x_{opt} such that $Ax_{opt} = b$. Want to find the "best" x_{opt} such that $Ax_{opt} \approx b$.

Exact solution to L₂ regression

Cholesky Decomposition:

If A is full rank and well-conditioned, decompose $A^TA = R^TR$, where R is upper triangular, and solve the normal equations: $R^TRx = A^Tb$.

QR Decomposition:

Slower but numerically stable, esp. if A is rank-deficient. Write A = QR, and solve $Rx = Q^{T}b$.

Singular Value Decomposition:

Most expensive, but best if A is very ill-conditioned. Write $A = U\Sigma V^{T}$, in which case: $\mathbf{x}_{opt} = A^{+}b = V\Sigma^{-1}U^{T}b$.

Complexity is $O(nd^2)$, but constant factors differ.

Projection of b on the subspace spanned by the columns of A

$$Z_{2}^{2} = \|b\|_{2}^{2} - \|AA^{+}b\|_{2}^{2}$$
$$x_{opt} = A^{+}b$$

Pseudoinverse of A

Algorithm: Sampling for L₂ regression

(Drineas, Mahoney, Muthukrishnan SODA 2006, Sarlos FOCS 2007, Drineas, Mahoney, Muthukrishnan, & Sarlos NumMath2011)

$$\mathcal{Z}_{2}^{2} = \min_{x \in \mathbb{R}^{d}} \|b - Ax\|_{2}^{2} = \|b - Ax_{opt}\|_{2}^{2}$$



Algorithm

- Compute a probability distribution over the rows of A (p_i, i=1...n, summing up to one).
- 2. In r i.i.d. trials pick r rows of A and the corresponding elements of b with respect to the p_i .

(Rescale sampled rows of A and sampled elements of b by $(1/(rp_i)^{1/2})$.)

3. Solve the induced problem.

Algorithm: Sampling for L₂ regression

(Drineas, Mahoney, Muthukrishnan SODA 2006, Sarlos FOCS 2007, Drineas, Mahoney, Muthukrishnan, & Sarlos NumMath2011)

$$\tilde{\mathcal{Z}}_2^2 = \min_{x \in \mathbb{R}^d} \left\| \tilde{b} - \tilde{A}x \right\|_2^2 = \left\| \tilde{b} - \tilde{A}\tilde{x}_{opt} \right\|_2^2$$





- Compute a probability distribution over the rows of A (p_i, i=1...n, summing up to one).
- 2. In r i.i.d. trials pick r rows of A and the corresponding elements of b with respect to the p_i .

(Rescale sampled rows of A and sampled elements of b by $(1/(rp_i)^{1/2})$.)

3. Solve the induced problem.

Algorithm: Sampling for L₂ regression

(Drineas, Mahoney, Muthukrishnan SODA 2006, Sarlos FOCS 2007, Drineas, Mahoney, Muthukrishnan, & Sarlos NumMath2011)

$$\tilde{\mathcal{Z}}_2^2 = \min_{x \in \mathbb{R}^d} \left\| \tilde{b} - \tilde{A}x \right\|_2^2 = \left\| \tilde{b} - \tilde{A}\tilde{x}_{opt} \right\|_2^2$$

Algorithm



- 1. Compute a probability distribution over the rows of A (p_i , i=1...n, summing up to one).
- 2. In r i.i.d. trials pick r rows of A and the corresponding elements of b with respect to the p_i .

(Rescale sampled rows of A and sampled elements of b by $(1/(rp_i)^{1/2})$.)

3. Solve the induced problem.

We will now discuss the p_i 's: our work introduced the notion of the leverage scores.

Leverage scores: tall & thin matrices

Let A be a (full rank) n-by-d matrix with n>>d whose SVD is:

$$\begin{pmatrix} A \\ A \end{pmatrix} = \begin{pmatrix} U \\ U \\ d \times d \end{pmatrix} \begin{pmatrix} \Sigma \\ d \times d \end{pmatrix} \begin{pmatrix} V^T \\ d \times d \end{pmatrix}$$
$$n \times d \quad n \times d$$

- \succ The matrix U contains the left singular vectors of A.
- \succ The columns of U are pairwise orthogonal and normal.
- This is NOT the case for rows of U: all we know is that the Euclidean norms of its rows are between zero and one.

Leverage scores: tall & thin matrices

Let A be a (full rank) n-by-d matrix with n>>d whose SVD is:



The (row) leverage scores can now be used to sample rows from A to create a sketch.

Computing leverage scores

Drineas, Magdon-Ismail, Mahoney, and Woodruff ICML 2012, JMLR 2012

> <u>Trivial</u>: via the Singular Value Decomposition

 $O(nd^2)$ time for n-by-d matrices with n>d.

> <u>Non-trivial</u>: relative error $(1+\epsilon)$ approximations for all leverage scores.

Tall & thin matrices:



<u>Running time:</u> $O(nd\epsilon^{-2} polylog(n/\epsilon)).$

Theorem

If the p_i are the row leverage scores of A, then, with probability at least 0.8,

$$\|b - Ax_{opt}\|_{2} \le \|b - A\tilde{x}_{opt}\|_{2} \le (1 + \epsilon) \|b - Ax_{opt}\|_{2}$$

The sampling complexity (the value of r) is

$$r = O\left(\frac{d}{\epsilon} + d\ln d\right)$$

Proof: a structural result

Consider the over-constrained least-squares problem:

$$\mathcal{Z}_2^2 = \min_{x \in \mathbb{R}^d} \|b - Ax\|_2^2 = \|b - Ax_{opt}\|_2^2$$

and the "sketched" (or "preconditioned") problem

$$\tilde{\mathcal{Z}}_{2}^{2} = \min_{x \in \mathbb{R}^{d}} \|X(b - Ax)\|_{2}^{2} = \|Xb - XA\tilde{x}_{opt}\|_{2}^{2}$$

<u>Recall</u>: A is n-by-d with n »> d; X is r-by-n with r << n.

- > Think of XA as a "sketch" of A.
- Our approach (using the leverage scores) focused on sketches of A that consist of (rescaled) rows of A.
- > More general matrices X are possible and have been heavily studied.

 $\mathcal{Z}_{2}^{2} = \min_{x \in \mathbb{R}^{d}} \|b - Ax\|_{2}^{2} = \|b - Ax_{opt}\|_{2}^{2} \qquad \tilde{\mathcal{Z}}_{2}^{2} = \min_{x \in \mathbb{R}^{d}} \|X(b - Ax)\|_{2}^{2} = \|Xb - XA\tilde{x}_{opt}\|_{2}^{2}$

Let U_A be the n-by-d matrix of the left singular vectors of A. If X satisfies (constants are somewhat arbitrary):

$$b^{\perp} = b - U_A U_A^T b \qquad \sigma_{min}^2 \left(X U_A \right) \ge 1/\sqrt{2}$$
$$\left\| U_A^T X^T X b^{\perp} \right\|_2^2 \le \epsilon \mathcal{Z}_2^2/2,$$

then,

$$\|A\tilde{x}_{opt} - b\|_{2} \leq (1+\epsilon)\mathcal{Z}_{2}$$
$$\|x_{opt} - \tilde{x}_{opt}\|_{2} \leq \frac{1}{\sigma_{min}(A)}\sqrt{\epsilon}\mathcal{Z}_{2}$$

Constructions for X

If X is a sampling-and-rescaling matrix formed using the row leverage scores of the matrix A, then both conditions are satisfied.

(I.e., an r-by-n matrix whose t-th row has a single non-zero entry indicating, and rescaling, the row of A that was sampled at the t-th trial.)

- Interestingly, many other matrices X satisfy both conditions: e.g., X can be a matrix whose entries are:
 - Random Gaussians (up to normalization).
 - Random signs (up to normalization).
 - > The randomized Hadamard transform and its variants.
 - > The input sparsity transform of Clarkson & Woodruff.

The "heart" of the proof

At the heart of all proofs in this line of research lies the following observation:



Then, we can prove that with probability at least $1-\delta$:

$$\left\| U_A^T U_A - U_A^T X^T X U_A \right\|_2 = \left\| I - U_A^T X^T X U_A \right\|_2 \le \varepsilon$$

It follows that, for all i: $\sqrt{1-\varepsilon} \leq \sigma_i \left(X U_A \right) \leq \sqrt{1+\varepsilon}$

The "heart" of the proof (cont'd)

<u>Recall</u>: with probability at least $1-\delta$:

$$\left\| U_A^T U_A - U_A^T X^T X U_A \right\|_2 = \left\| I - U_A^T X^T X U_A \right\|_2 \le \varepsilon$$

It follows that, for all i: $\sqrt{1-\varepsilon} \leq \sigma_i \left(X U_A \right) \leq \sqrt{1+\varepsilon}$

- The sampling complexity is r=O(d ln d).
- Proving the above inequality is (now) routinely done via matrix concentration inequalities (at least in most cases).
- > Early proofs were very complicated and not user-friendly.



A lot of follow-up work, including:

- Avron, Maymounkov, and Toledo SISC 2010: Blendenpik, a solver that uses the "sketch" XA as a preconditioner, combined with an iterative least-squares solver. Beats LAPACK by a factor of four in essentially all over-constrained leastsquares problems.
 - Iyer, Avron, Kollias, Inechein, Carothers, and Drineas JCS 2016: an evaluation of Blendenpik on terascale matrices in Rensselaer's BG/Q; again factor four-to-six speedups compared to Elemental's QR-based solver.
- Drineas, Mahoney, Woodruff, and collaborators (SODA 2008, SIMAX 2009, SODA 2013, SIMAX 2016): general p-norm regression, beyond Euclidean norm.
- Clarkson and Woodruff STOC 2013: relative error algorithms for overconstrained least-squares regression problems in input sparsity time using a novel construction for the sketching matrix X.

Follow-up (cont'd)

- Pilanci and Wainwright IEEE TIF 2015, JMLR 2016, SIOPT 2017: A novel iterative sketching-based method (Hessian sketch) to solve over-constrained least-squares regression problems over convex bodies.
- Paul, Magdon-Ismail, and Drineas NIPS 2015, Derezinski and Warmuth NIPS 2017, AISTATS 2018, COLT 2018, JMLR 2018: Adaptive and volume sampling approaches to construct the sketching matrix X.
- Alaoui and Mahoney NIPS 2015, Cohen, Musco, Musco, and collaborators STOC 2015, SODA 2017, FOCS 2017: ridge leverage scores, a smooth and regularized generalization of the leverage scores.
- Chowdhuri, Yang, and Drineas ICML 2018: structural conditions for underconstrained problems (n « d case); a preconditioned Richardson-like solver for such problems; check our paper for a detailed discussion on prior work for such under-constrained problems.

Related work: the "square" case

The "square" case: solving systems of linear equations

- Almost optimal relative-error approximation algorithms for Laplacian and, more generally, Symmetric Diagonally Dominant (SDD) matrices
 - Pioneered by Spielman and Teng, major contributions later by Miller, Koutis, Peng, and many others.
 - Roughly speaking, the proposed methods are iterative preconditioned solvers where the preconditioner is a sparse version of the original graph.
 - This sparse graph is constructed by sampling edges of the original graph with probability proportional to their *leverage scores*, which in the context of graphs are called *effective resistances*.
- <u>Still open:</u> progress beyond Laplacians.
 - Results by Peng Zhang and Rasmus Kyng (FOCS 2017) indicate that such progress might be challenging.
- Check Koutis, Miler, and Peng CACM 2012 for a quick intro.



- > RandNLA approaches for regression problems
- > RandNLA approaches for Principal Component Analysis (PCA)

PCA: An example in human genetics

Single Nucleotide Polymorphisms: the most common type of genetic variation in the genome across different individuals.

They are known locations at the human genome where two alternate nucleotide bases (alleles) are observed (out of A, C, G, T).

SNPs

individuals

... AG CT GT GG CT CC CC CC AG AG AG AG AG AG AA CT AA GG GG CC GG AG CG AC CC AA CC AA GG TT AG CT CG CG CG AT CT CT AG CT AG GG GT GA AG GG TT TT GG TT CC CC CC CC GG AA AG AG AG AA CT AA GG GG CC GG AA GG AA CC AA CC AA GG TT AG TT GG GG GG TT TT CC GG TT GG GG TT GG AA GG TT TT GG TT CC CC CC CC GG AA AG AG AG AA CT AA GG GG CC AG AG CG AC CC AA CC AA GG TT AG CT CG CG CG AT CT CT AG CT AG GG GT GA AG GG TT TT GG TT CC CC CC CC GG AA AG AG AA AG CT AA GG GG CC AG AG CG AC CC AA CC AA GG TT AG CT CG CG CG AT CT CT AG CT AG GT GT GA AG GG TT TT GG TT CC CC CC CC GG AA AG AG AG AA CC GG AA CC CA AG GG CC AC CC AA CG AA GG TT AG CT CG CG CG AT CT CT AG CT AG GT GT GA AG GG TT TT GG TT CC CC CC CC GG AA GG GG GG AA CT AA GG GG CT GG AG CC CC CG AA CC AA GG TT AG CT CG CG CG AT CT CT AG CT AG GT TT GG AA GG TT TT GG TT CC CC CC CG CC AG AG AG AG AG AA CT AA GG GG CT GG AG CC CC CG AA CC AA GT TT AG CT CG CG CG AT CT CT AG CT AG GT TGG AA GG TT TT GG TT CC CC CC CG GC AA AG AG AG AG AA CT AA GG GG CT GG AG CC CC CG AA CC AA GT TT AG CT CG CG CG AT CT CT AG CT AG GT TGG AA GG TT TT GG TT CC CC CC CC GG AA AG AG AG AA TT AA GG GG CC AG AG CC CC CG AA CC AA GT TT AG CT CG CG CG AT CT CT AG CT AG GT TT GG AA GG TT TT GG TT CC CC CC CC GG AA AG AG AG AA TT AA GG GG CC AG AG CC AA CC AA CG AA GT TT AG TT GG GG TT TT CC GG TT GG GT TT GG AA ...

Typical sizes: tens of thousands of individuals and hundreds of thousands of SNPs.



45 Yakut

HGDP data

- 1,033 samples
- 7 geographic regions
- 52 populations

HapMap Phase 3 data

- 1,207 samples
- 11 populations

Matrix dimensions:

2,240 subjects (rows) 447,143 SNPs (columns)

We will apply PCA (i.e., SVD on a suitably rescaled covariance matrix) to visualize and/or analyze the data. SVD: formal definition

$$\begin{array}{c} A \\ m \times n \end{array} \right) = \left(\begin{array}{c} U \\ m \times \rho \end{array} \right) \cdot \left(\begin{array}{c} \Sigma \\ \rho \times \rho \end{array} \right) \cdot \left(\begin{array}{c} V \\ \rho \times n \end{array} \right)^{T}$$

 ρ : rank of A

U (V): orthogonal matrix containing the left (right) singular vectors of A.

 Σ : diagonal matrix containing the singular values of A.

Let σ_1 , σ_2 , ..., σ_ρ be the entries of Σ .

Computing the SVD takes $O(\min\{mn^2, m^2n\})$ time.

The top k left/right singular vectors/values can be computed faster using iterative methods.



HGDP data

- 1,033 samples
- 7 geographic regions
- 52 populations

HapMap Phase 3 data

- 1,207 samples
- 11 populations

Matrix dimensions:

2,240 subjects (rows) 447,143 SNPs (columns)

PCA on the above data returns:

Paschou, Lewis, Javed, & Drineas (2010) J Med Genet



• Top two Principal Components (PCs or eigenSNPs)

(Lin and Altman (2005) Am J Hum Genet)

- Very good correlation between geography and the top two eigenSNPs.
- Mexican population seems out of place: we move to the top three PCs.



Not altogether satisfactory: the principal components are linear combinations of all SNPs, and - of course - can not be assayed!

Can we find **actual SNPs** that capture the information in the singular vectors? Formally: spanning the same subspace.



- PCA plots of genetic data from multiple populations around the Mediterranean Sea indicate that the Mediterranean acted as a "**barrier**" during the colonization of Europe from our species.
- Using PCA (and many other analyses) we proposed what is a called a **maritime route** for the colonization of Europe.
- Interpreting the singular vectors is, again, tricky; we identified SNPs (and genes) that capture the information in the singular vectors.

ARTICLE

Genetics of the peloponnesean populations and the theory of extinction of the medieval peloponnesean Greeks

George Stamatoyannopoulos^{*,1}, Aritra Bose², Athanasios Teodosiadis³, Fotis Tsetsos², Anna Plantinga⁴, Nikoletta Psatha⁵, Nikos Zogas⁶, Evangelia Yannaki⁶, Pierre Zalloua⁷, Kenneth K Kidd⁸, Brian L Browning^{4,9}, John Stamatoyannopoulos^{3,10}, Peristera Paschou¹¹ and Petros Drineas²



ACHAEA

PCA identifies and extracts genetic micro-structure at very local levels and small geographical distances.

Consider, for example, Peloponnesos.

Genetics of the peloponnesean populations and the theory of extinction of the medieval peloponnesean Greeks



<u>George Stamatoyannopoulos*,1</u>, Aritra Bose², Athanasios Teodosiadis³, Fotis Tsetsos², Anna Plantinga⁴, Nikoletta Psatha⁵, Nikos Zogas⁶, Evangelia Yannaki⁶, Pierre Zalloua⁷, Kenneth K Kidd⁸, Brian L Browning^{4,9}, John Stamatoyannopoulos^{3,10}, Peristera Paschou¹¹ and Petros Drineas²

SVD: computational time

Computing large SVDs: computational time

• In commodity hardware (e.g., a 4GB RAM, dual-core laptop), using MatLab 7.0 (R14), the computation of the SVD of the dense 2,240-by-447,143 matrix A <u>takes about 12 minutes</u>.

• Computing this SVD is not a one-liner, since we can not load the whole matrix in RAM (runs out-of-memory in MatLab); we compute the eigendecomposition of AA^{T} .

• In 2010, we had to compute **1,200 SVDs** on matrices of dimensions (approx.) 1,200-by-450,000 for a full leave-one-out cross-validation experiment.

(Drineas, Lewis, & Paschou (2010) PLoS ONE)

• To compare mtDNA derived from 37 ancient Minoan bones to 120 extant and ancient populations we ran (multiple) SVDs on (approx.) 14,000-by-14,000 matrices.

(Hughey, Paschou, Drineas, et al. (2013) Nat Comm)

• Current population genetics datasets generate 1,000,000-by-1,000,000 matrices (Bose et al. (2018) TeraPCA package.)

SVD: computational time

Computing large SVDs: computational time

• In commodity hardware (e.g., a 4GB RAM, dual-core laptop), using MatLab 7.0 (R14), the computation of the SVD of the dense 2,240-by-447,143 matrix A <u>takes about 12 minutes</u>.

• Computing this SVD is not a one-liner, since we can not load the whole matrix in RAM (runs out-of-memory in MatLab); we compute the eigendecomposition of AA^{T} .

• In 2010, we had to compute **1,200 SVDs** on matrices of dimensions (approx.) 1,200-by-450,000 for a full leave-one-out cross-validation experiment.

(Drineas, Lewis, & Paschou (2010) PLoS ONE)

• To compare mtDNA derived from 37 ancient Minoan bones to 120 extant and ancient populations we ran (multiple) SVDs on (approx.) 14,000-by-14,000 matrices.

(Hughey, Paschou, Drineas, et al. (2013) Nat Comm)

• Current population genetics datasets generate 1,000,000-by-1,000,000 matrices. (Bose et al. (2018) TeraPCA package.)

• Running time is <u>always</u> a concern, <u>but</u> machine-precision is <u>not</u> necessary!

• Data are noisy and approximate singular vectors work well in many settings.





- > It is easy to see that X = $\Sigma_k V_k^T = U_k^T A$.
- SVD has strong optimality properties.
- > The columns of U_k are linear combinations of up to all columns of A.

The CX decomposition Drineas, Mahoney, & Muthukrishnan (2008) SIAM J Mat Anal Appl Mahoney & Drineas (2009) PNAS



Why?

If A is an subject-SNP matrix, then selecting representative columns is equivalent to selecting representative SNPs to capture the same structure as the top eigenSNPs.

We want c as small as possible!



Easy to prove that optimal $X = C^{+}A$. (C^{+} is the Moore-Penrose pseudoinverse of C.) Thus, the challenging part is to find good columns (SNPs) of A to include in C.

From a mathematical perspective, this is a hard combinatorial problem, closely related to the so-called Column Subset Selection Problem (CSSP).

The CSSP has been heavily studied in Numerical Linear Algebra.

Relative-error Frobenius norm bounds Drineas, Mahoney, & Muthukrishnan (2008) SIAM J Mat Anal Appl

Given an m-by-n matrix A, there exists an $O(mn^2)$ algorithm that picks

at most O((k/ ϵ^2) log (k/ ϵ)) columns of A

such that with probability at least .9

$$\|A - CX\|_F = \left\|A - CC^{\dagger}A\right\|_F \le (1 + \varepsilon) \|A - A_k\|_F$$

<u>Notation</u>: $||X||_F^2 = \sum_{i,j} X_{ij}^2$

The algorithm

- <u>Input:</u> m-by-n matrix A,
 - $0 < \epsilon < .5$, the desired accuracy
- <u>Output:</u> C, the matrix consisting of the selected columns

Sampling algorithm

- Compute probabilities p_j summing to 1.
- Let c = O($(k/\epsilon^2) \log (k/\epsilon)$).

• In c i.i.d. trials pick columns of A, where in each trial the j-th column of A is picked with probability p_j .

• Let C be the matrix consisting of the chosen columns.

Subspace sampling (Frobenius norm)

$$\begin{pmatrix} A_k \\ m \times n \end{pmatrix} = \begin{pmatrix} U_k \\ m \times k \end{pmatrix} \cdot \begin{pmatrix} \Sigma_k \end{pmatrix} \cdot \begin{pmatrix} V_k^T \\ K \end{pmatrix}$$

 $V_{\rm k}\!\!:$ orthogonal matrix containing the top k right singular vectors of A.

 Σ_k : diagonal matrix containing the top k singular values of A.

Remark: The rows of V_k^{T} are orthonormal vectors, but its columns $(V_k^{T})^{(i)}$ are not.

Leverage score sampling:

$$p_{j} = \frac{\left\| \left(V_{k}^{T} \right)^{(j)} \right\|_{2}^{2}}{k}$$
Normalization s.t. the p_{j} sum up to 1

Subspace sampling (Frobenius norm)

$$\begin{pmatrix} A_k \\ m \times n \end{pmatrix} = \begin{pmatrix} U_k \\ m \times k \end{pmatrix} \cdot \begin{pmatrix} \Sigma_k \end{pmatrix} \cdot \begin{pmatrix} V_k^T \\ K \end{pmatrix} \cdot \begin{pmatrix} V_k^T \end{pmatrix} \cdot \begin{pmatrix} W_k^T \end{pmatrix} \cdot \begin{pmatrix}$$

 V_k : orthogonal matrix containing the top k right singular vectors of A.

 Σ_{k} : diagonal matrix containing the top k singular values of A.

Remark: The rows of V_k^{T} are orthonormal vectors, but its columns $(V_k^{T})^{(i)}$ are not.





SNPs by chromosomal order

Paschou et al (2007; 2008) PLoS Genetics; Paschou et al (2010) J Med Genet; Drineas et al (2010) PLoS One Hughey, Paschou, Drineas, et al. (2013) Nat Comm; Paschou, Drineas, et al. PNAS 2014;

Selecting PCA SNPs for individual assignment to four continents (Africa, Europe, Asia, America)



SNPs by chromosomal order

Paschou et al (2007; 2008) PLoS Genetics; Paschou et al (2010) J Med Genet; Drineas et al (2010) PLoS One Hughey, Paschou, Drineas, et al. (2013) Nat Comm; Paschou, Drineas, et al. PNAS 2014;

Approximating leverage scores

<u>Can we approximate the leverage scores fast?</u>

<u>Theorem</u>: Given any m-by-n matrix A with m > n, we can approximate its leverage scores (where k is the target rank) with relative error accuracy in

O(mnk log m) time,

as opposed to the - trivial - $O(mn^2)$ time.

(Drineas, Mahoney, Magdon-Ismail, & Woodruff ICML '12 JMLR '12)

Improvement: leverage scores can be computed in O(nnz(A) k) time!

Clarkson and Woodruff (STOC '13): introduced a sparse random projection; Mahoney and Meng (STOC '13): provided a better analysis for the above result; Nelson and Huy (FOCS '13): provided the best known analysis for the above result; Boutsidis and Woodruff (STOC '14): applications to many RandNLA problems. Sobczyk and Gallopoulos '17: block iterative methods for fast estimation

Selecting fewer columns

Problem

How many columns do we need to include in the matrix C in order to get relative-error approximations ?

<u>Recall</u>: with $O((k/\epsilon^2) \log (k/\epsilon))$ columns, we get (subject to a failure probability)

$$\left\| A - CC^{\dagger}A \right\|_{F} \le (1+\epsilon) \left\| A - A_{k} \right\|_{F}$$

Deshpande & Rademacher (FOCS '10): with exactly k columns, we get

$$\left|A - CC^{\dagger}A\right\|_{F} \le \sqrt{k} \left\|A - A_{k}\right\|_{F}$$

What about the range between k and O(k log(k))?

Selecting fewer columns (cont'd)

(Boutsidis, Drineas, & Magdon-Ismail, FOCS 2011 and SICOMP 2014)

Question:

What about the range between k and O(k log(k))?

Answer:

A relative-error bound is possible by selecting $s=2k/\epsilon$ columns!

<u>Technical breakthrough;</u>

A combination of sampling strategies with a novel approach on column selection, inspired by the work of Batson, Spielman, & Srivastava (STOC '09) on graph sparsifiers.

- The running time is $O((mnk+nk^3)\epsilon^{-1})$.
- Simplicity is gone...

Lower bounds and alternative approaches

Deshpande & Vempala, RANDOM 2006

A relative-error approximation necessitates at least k/ϵ columns.

<u>Guruswami & Sinop, SODA 2012</u>

Alternative approaches, based on volume sampling, guarantee

(r+1)/(r+1-k) relative error bounds.

This bound is asymptotically optimal (up to lower order terms).

The proposed deterministic algorithm runs in $O(rnm^3 \log m)$ time, while the randomized algorithm runs in $O(rnm^2)$ time and achieves the bound in expectation.

Guruswami & Sinop, FOCS 2011

Applications of column-based reconstruction in Quadratic Integer Programming.

Musco, Musco, Cohen, Woodruff, and collaborators

Multiple articles in STOC, FOCS, SODA, NIPS, ICML in 2016 and 2017 on ridge leverage scores and other approaches.



To get highly accurate approximations to singular vectors, use iterative methods.

1. Block subspace iteration

Given an m-by-n matrix A and a positive integer q, compute

$$K = \left(AA^T\right)^q AX$$

where X is an n-by-p (with $p \approx k$) random matrix, e.g., a random Gaussian matrix.

Compute the best rank-k approximation to A within the subspace spanned by the columns of K (much easier to do than it sounds...): denote it by \tilde{A}_k .

Iterative methods for PCA (cont'd)

1. Block subspace iteration

Given an m-by-n matrix A and a positive integer q, compute

 $K = \left(AA^T\right)^q AX$

where X is an n-by-p (with $p \approx k$) random matrix, e.g., a random Gaussian matrix.

Compute the best rank-k approximation to A within the subspace spanned by the columns of K (much easier to do than it sounds...): denote it by \tilde{A}_k .

- Strong bounds can be proven for the Frobenius and spectral norms of the matrix $A \tilde{A}_k$.
- We implemented block subspace iteration to approximate the top singular vectors of terascale matrices arising in population genetics in:

A. Bose, V. Kalantzis, E. Kontopoulou, M. Elkadi, P. Paschou, and P. Drineas, "TeraPCA: a fast and scalable method to study genetic variation in tera-scale genotypes", under review, Genome Biology, 2018.

Iterative methods for PCA

2. Block Krylov methods

Given an m-by-n matrix A (of rank ρ) and a positive integer q, compute

$$K = \left[AX, \left(AA^{T}\right)AX, \left(AA^{T}\right)^{2}AX, \ldots, \left(AA^{T}\right)^{q}AX\right]$$

where X is an n-by-p (with $p \approx k$) random matrix, e.g., a random Gaussian matrix.

Compute the best rank-k approximation to A within the subspace spanned by the columns of K (much easier to do than it sounds): denote it by \tilde{A}_k .

Assume a gap g(>0) between the k and (k+1)-st singular values (can be relaxed):

$$\sigma_k \ge (1+g)\,\sigma_{k+1} > 0$$

Iterative methods for PCA

2. Block Krylov methods

Given an m-by-n matrix A (of rank ρ) and a positive integer q, compute

$$K = \left[AX, \left(AA^{T}\right)AX, \left(AA^{T}\right)^{2}AX, \ldots, \left(AA^{T}\right)^{q}AX\right]$$

where X is an n-by-p (with $p \approx k$) random matrix, e.g., a random Gaussian matrix.

Compute the best rank-k approximation to A within the subspace spanned by the columns of K (much easier to do than it sounds): denote it by \tilde{A}_k .

■ Assume a gap g(>0) between the k and (k+1)-st singular values (can be relaxed):

$$\sigma_k \ge (1+g)\,\sigma_{k+1} > 0$$

Bottom p-k singular vectors of A

$$\left\| V_{k,\perp}^T X \right\|_F^2 \le \gamma_2^2 (\rho - k)$$

Also assume (γ_1 and γ_2 are constants):

 $\sigma_{\min}^2\left(V_k^T X\right) \ge \gamma_1^2$ and

Iterative methods for PCA

2. Block Krylov methods

Given an m-by-n matrix A (of rank ρ) and a positive integer q, compute $M = \begin{bmatrix} AX, (AA^{T}) AX, (AA^{T})^{2} AX, \dots, (AA^{T})^{q} AX \end{bmatrix}$ $q = O\left(\frac{\log(\rho/\epsilon)}{\sqrt{g}}\right)$

where X is an n-by-p (with $p \approx k$) random matrix, e.g., a random Gaussian matrix.

Compute the best rank-k approximation to A within the subspace spanned by the columns of K (much easier to do than it sounds): denote it by \tilde{A}_k . Then,

$$\begin{aligned} \left\| A - \tilde{A}_k \right\|_F &\leq \| A - A_k \|_F + \epsilon \sigma_{k+1} \\ \left\| A - \tilde{A}_k \right\|_2 &\leq \| A - A_k \|_2 + \epsilon \sigma_{k+1} \end{aligned}$$

RandNLA and optimization

• Primal dual interior point methods necessitate solving least-squares problems (projecting the gradient on the null space of the constraint matrix in order to remain feasible).

(Dating back to the mid/late 1980's and work by Karmarkar, Ye, Freund)

- Can we solve these least squares problems approximately using random sampling/random projections?
- <u>Modern approaches</u>: primal/dual interior point methods iterate along an approximation to the Newton direction and tolerate (mild) infeasibilities. A system of linear equations must be solved.

(inexact interior point methods: work by Bellavia, Steihaug, etc.)

- <u>Well-known by practitioners</u>: the number of iterations in interior point methods is <u>not</u> the bottleneck, but the computational cost of solving a linear system is.
- <u>Goal:</u> Use sampling/random projection approaches to design efficient precoditioners to solve systems of linear equations that arise in primal-dual interior point methods faster.

Progress by Roosta & Mahoney (ArXiv 2016, 2017 multiple papers on subsampled second-order methods).



"Randomization is arguably the most exciting and innovative idea to have hit linear algebra in a long time." (Avron et al. (2010) SISC)

> RandNLA workshop, Simons Institute for the Theory of Computing, UC Berkeley, Foundations of Data Science, Sep 2018 <u>https://simons.berkeley.edu/data-science-2018-1</u>



RandNLA course, PCMI Summer School on Mathematics of Data, Jul 2016 <u>Drineas & Mahoney, Lectures on RandNLA, Vol. 25, Amer. Math. Soc., 2018</u>

> Highlighted at the Workshops on Algorithms for Modern Massive Datasets (MMDS) 2006, 2008, 2010, 2012, 2014, and 2016.

http://mmds-data.org/

Gene Golub SIAM Summer School (G2S3), Δελφοί, Greece, June 2015
<u>http://scgroup19.ceid.upatras.gr/g2s32015/</u>

> Invited tutorial at SIAM ALA 2015

> RandNLA workshop in FOCS 2012

http://ieee-focs.org/focs2012/workshops/RandomNLA/







<u>Stoupa</u>

Stoupa is where the (real) story of Zorba the Greek took place.





<u>Stoupa</u>

Stoupa is where the (real) story of Zorba the Greek took place.

Nomitsis (this is where I come from...)

Church of Agioi Anargyroi and Church of the Metamorphosis: a thousand years old.



h) tower-houses

he Greek took place.

e Metamorphosis: a



oogle



<u>Stoupa</u>

Stoupa is where the (real) story of Zorba the Greek took place.

Nomitsis (this is where I come from...)

Church of Agioi Anargyroi and Church of the Metamorphosis: a thousand years old.

<u>Oitylo</u>

First mentioned in Homer's Iliad.



l/modern) tower-houses

f Zorba the Greek took place.

urch of the Metamorphosis: a

oogle



<u>Stoupa</u>

Stoupa is where the (real) story of Zorba the Greek took place.

Nomitsis (this is where I come from...)

Church of Agioi Anargyroi and Church of the Metamorphosis: a thousand years old.

<u>Oitylo</u>

First mentioned in Homer's Iliad.

<u>Diros</u>

> Diros caves: partially submerged an underground river.

