
Column Subset Selection



Joel A. Tropp

Applied & Computational Mathematics
California Institute of Technology
`jtropp@acm.caltech.edu`

Thanks to B. Recht (Caltech, IST)

Column Subset Selection

$$A = \left[\begin{array}{c} \text{orange} \\ \text{green} \\ \text{cyan} \\ \text{yellow} \\ \text{tan} \\ \text{blue} \\ \text{brown} \\ \text{purple} \\ \text{dark blue} \\ \text{grey} \\ \text{red} \\ \text{light blue} \end{array} \right]$$

$$\mathcal{T} = \{ \text{orange} \quad \text{green} \quad \text{blue} \quad \text{brown} \quad \text{dark blue} \quad \text{red} \}$$

$$A_{\mathcal{T}} = \left[\begin{array}{c} \text{orange} \\ \text{green} \\ \text{blue} \\ \text{brown} \\ \text{dark blue} \\ \text{red} \end{array} \right]$$

Spectral Norm Reduction

Theorem 1. [Kashin–Tzafriri] *Suppose the n columns of \mathbf{A} have unit ℓ_2 norm. There is a set τ of column indices for which*

$$|\tau| \geq \frac{n}{\|\mathbf{A}\|^2} \quad \text{and} \quad \|\mathbf{A}_\tau\| \leq C.$$

Examples:

☛ \mathbf{A} has identical columns. Then $|\tau| \geq 1$.

☛ \mathbf{A} has orthonormal columns. Then $|\tau| \geq n$.

Spectral Norm Reduction

Theorem 1. [Kashin–Tzafriri] *Suppose the n columns of \mathbf{A} have unit ℓ_2 norm. There is a set τ of column indices for which*

$$|\tau| \geq \frac{n}{\|\mathbf{A}\|^2} \quad \text{and} \quad \|\mathbf{A}_\tau\| \leq C.$$

Theorem 2. [T 2007] *There is a randomized, polynomial-time algorithm that produces the set τ .*

Overview:

- ☞ Randomly select columns
- ☞ Remove redundant columns

Random Column Selection: Intuitions

- ☞ Random column selection reduces norms
- ☞ A random submatrix gets “its share” of the total norm
- ☞ Submatrices with small norm are ubiquitous
- ☞ Random selection is a form of regularization
- ☞ Added benefit: Dimension reduction

Example: What Can Go Wrong

$$A = \left[\begin{array}{c|c} \begin{array}{c} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{array} & \begin{array}{c} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{array} \end{array} \right]$$

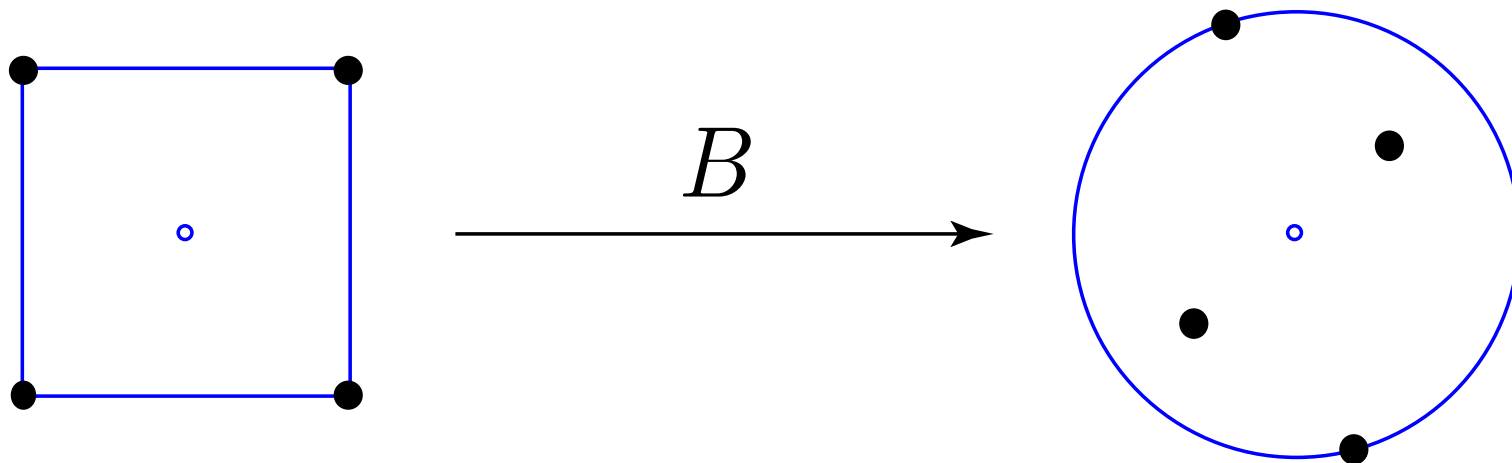
$$A_{\tau} = \left[\begin{array}{c|c} \begin{array}{c} 1 \\ 1 \\ 1 \end{array} & \begin{array}{c} 1 \\ 1 \\ 1 \end{array} \end{array} \right]$$

$$\|A\| = \|A_{\tau}\| = \sqrt{2} \quad \Longrightarrow \quad \textit{No reduction!}$$

The $(\infty, 2)$ Operator Norm

Definition 3. The $(\infty, 2)$ operator norm of a matrix B is

$$\|B\|_{\infty, 2} = \max \{ \|Bx\|_2 : \|x\|_{\infty} = 1 \}.$$



Proposition 4. If B has s columns, then the best general bound is

$$\|B\|_{\infty, 2} \leq \sqrt{s} \|B\|.$$

Random Reduction of $(\infty, 2)$ Norm

Lemma 5. *Suppose the n columns of \mathbf{A} have unit ℓ_2 norm. Draw a uniformly random subset σ of columns whose cardinality*

$$|\sigma| = \frac{2n}{\|\mathbf{A}\|^2}.$$

Then

$$\mathbb{E} \|\mathbf{A}_\sigma\|_{\infty,2} \leq C\sqrt{|\sigma|}.$$

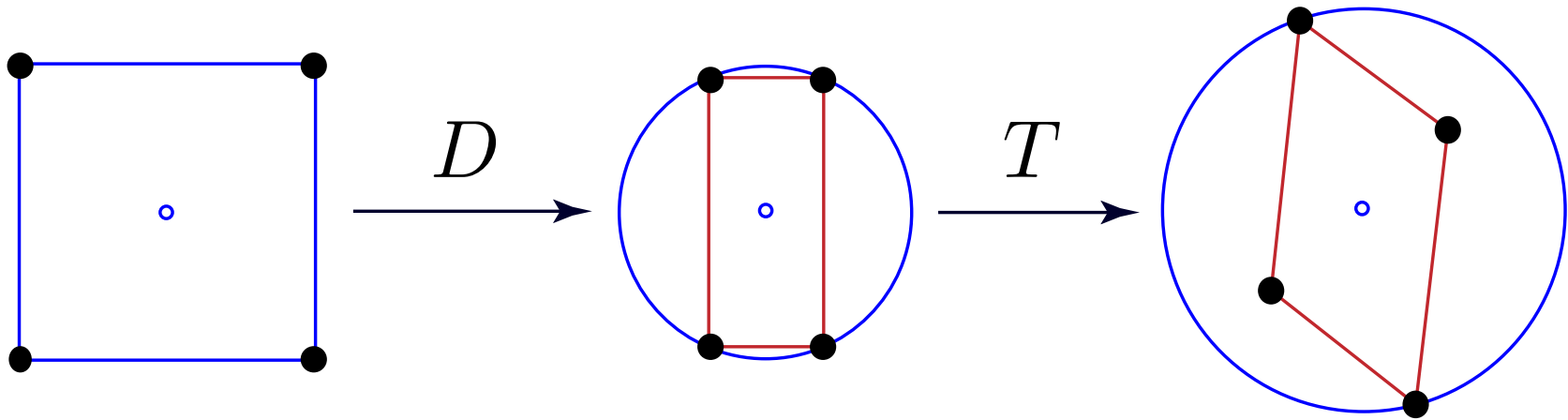
🐛 **Problem:** How can we use this information?

Pietsch Factorization

Theorem 6. [Pietsch, Grothendieck] *Every matrix B can be factorized as $B = TD$ where*

• D is diagonal and nonnegative with $\text{trace}(D^2) = 1$, and

• $\|B\|_{\infty,2} \leq \|T\| \leq \sqrt{\pi/2} \|B\|_{\infty,2}$



Pietsch and Norm Reduction

Lemma 7. *Suppose B has s columns. There is a set τ of column indices for which*

$$|\tau| \geq \frac{s}{2} \quad \text{and} \quad \|B_\tau\| \leq \sqrt{\pi} \cdot \frac{1}{\sqrt{s}} \|B\|_{\infty,2}.$$

Proof. Consider a Pietsch factorization $B = TD$. Select

$$\tau = \{j : d_{jj}^2 \leq 2/s\}.$$

Since $\sum d_{jj}^2 = 1$, Markov's inequality implies $|\tau| \geq s/2$. Calculate

$$\|B_\tau\| = \|TD_\tau\| \leq \|T\| \cdot \|D_\tau\| \leq \sqrt{\pi/2} \|B\|_{\infty,2} \cdot \sqrt{2/s}.$$

Proof of Kashin–Tzafriri

- Suppose the n columns of \mathbf{A} have unit ℓ_2 norm
- Lemma 5 provides (random) σ for which

$$|\sigma| = \frac{2n}{\|\mathbf{A}\|^2} \quad \text{and} \quad \|\mathbf{A}_\sigma\|_{\infty,2} \leq C\sqrt{|\sigma|}$$

- Lemma 7 applied to $\mathbf{B} = \mathbf{A}_\sigma$ yields a subset $\tau \subset \sigma$ for which

$$|\tau| \geq \frac{|\sigma|}{2} \quad \text{and} \quad \|\mathbf{B}_\tau\| \leq \sqrt{\pi} \cdot \frac{1}{\sqrt{|\sigma|}} \cdot \|\mathbf{B}\|_{\infty,2}$$

- Simplify

$$|\tau| \geq \frac{n}{\|\mathbf{A}\|^2} \quad \text{and} \quad \|\mathbf{A}_\tau\| \leq C\sqrt{\pi}$$

- *Note:* This is almost an algorithm

Pietch and Eigenvalues

☛ Consider a matrix B with Pietch factorization $B = TD$

☛ Suppose $\|T\| \leq \alpha$

☛ Calculate

$$\begin{aligned} B = TD &\implies \|Bx\|_2^2 = \|TDx\|_2^2 && \forall x \\ &\implies \|Bx\|_2^2 \leq \alpha^2 \|Dx\|_2^2 && \forall x \\ &\implies x^*(B^*B)x \leq \alpha^2 \cdot x^*D^2x && \forall x \\ &\implies x^* [B^*B - \alpha^2 D^2] x \leq 0 && \forall x \\ &\implies \lambda_{\max}(B^*B - \alpha^2 D^2) \leq 0 \end{aligned}$$

Pietsch is Convex

- 🐼 *Key new idea:* Can find Pietsch factorizations by convex programming

$$\min \lambda_{\max}(\mathbf{B}^* \mathbf{B} - \alpha^2 \mathbf{F})$$

$$\text{subject to } \mathbf{F} \text{ diagonal, } \mathbf{F} \geq \mathbf{0}, \text{ trace}(\mathbf{F}) = 1$$

- 🐼 If value at \mathbf{F}_\star is nonpositive, then we have a factorization

$$\mathbf{B} = (\mathbf{B}\mathbf{F}_\star^{-1/2}) \cdot \mathbf{F}_\star^{1/2} \quad \text{with} \quad \|\mathbf{B}\mathbf{F}_\star^{-1/2}\| \leq \alpha$$

- 🐼 Proof of Kashin–Tzafriri offers target value for α
- 🐼 Can also perform binary search to approximate minimal value of α

An Optimization over the Simplex

☛ Express $F = \text{diag}(\mathbf{f})$

☛ Constraints delineate the probability simplex:

$$\Delta = \{\mathbf{f} : \text{trace}(\mathbf{f}) = 1 \quad \text{and} \quad \mathbf{f} \geq \mathbf{0}\}$$

☛ Objective function and its subdifferential:

$$J(\mathbf{f}) = \lambda_{\max}(\mathbf{B}^* \mathbf{B} - \alpha^2 \text{diag}(\mathbf{f}))$$

$$\partial J(\mathbf{f}) = \text{conv} \left\{ -\alpha^2 |\mathbf{u}|^2 : \mathbf{u} \text{ top evec. } \mathbf{B}^* \mathbf{B} - \alpha^2 \text{diag}(\mathbf{f}), \|\mathbf{u}\|_2 = 1 \right\}$$

☛ Obtain

$$\min J(\mathbf{f}) \quad \text{subject to} \quad \mathbf{f} \in \Delta$$

Entropic Mirror Descent

1. Initialize $\mathbf{f}^{(1)} \leftarrow s^{-1}\mathbf{e}$ and $k \leftarrow 1$
2. Compute a subgradient: $\boldsymbol{\theta} \in \partial J(\mathbf{f}^{(k)})$
3. Determine step size:

$$\beta_k \leftarrow \sqrt{\frac{2 \log s}{k \|\boldsymbol{\theta}\|_\infty^2}}$$

4. Update variable:

$$\mathbf{f}^{(k+1)} \leftarrow \frac{\mathbf{f}^{(k)} \circ \exp\{-\beta_k \boldsymbol{\theta}\}}{\text{trace}(\mathbf{f}^{(k)} \circ \exp\{-\beta_k \boldsymbol{\theta}\})}$$

5. Increment $k \leftarrow k + 1$, and return to 2.

References: [Eggermont 1991, Beck–Teboulle 2003]

Other Formulations

• Modified primal to simultaneously identify α

$$\begin{aligned} \min \quad & \lambda_{\max}(\mathbf{B}^* \mathbf{B} - \alpha^2 \mathbf{F}) + \alpha^2 \\ \text{subject to} \quad & \mathbf{F} \text{ diagonal, } \mathbf{F} \geq \mathbf{0}, \quad \text{trace}(\mathbf{F}) = 1, \quad \alpha \geq 0 \end{aligned}$$

• Dual problem is the famous MAXCUT SDP:

$$\max \langle \mathbf{B}^* \mathbf{B}, \mathbf{Z} \rangle \quad \text{subject to} \quad \text{diag}(\mathbf{Z}) = \mathbf{e}, \quad \mathbf{Z} \succcurlyeq \mathbf{0}$$

Related Results

Theorem 8. [Bourgain–Tzafriri 1991] *Suppose the n columns of A have unit ℓ_2 norm. There is a set τ of column indices for which*

$$|\tau| \geq \frac{cn}{\|A\|^2} \quad \text{and} \quad \kappa(A_\tau) \leq \sqrt{3}.$$

Examples:

☛ A has identical columns. Then $|\tau| \geq 1$.

☛ A has orthonormal columns. Then $|\tau| \geq cn$.

Related Results

Theorem 8. [Bourgain–Tzafriri 1991] *Suppose the n columns of \mathbf{A} have unit ℓ_2 norm. There is a set τ of column indices for which*

$$|\tau| \geq \frac{cn}{\|\mathbf{A}\|^2} \quad \text{and} \quad \kappa(\mathbf{A}_\tau) \leq \sqrt{3}.$$

Theorem 9. [T 2007] *There is a randomized, polynomial-time algorithm that produces the set τ .*

To learn more...

E-mail:

✉ jtropp@acm.caltech.edu

Web: <http://www.acm.caltech.edu/~jtropp>

Papers in Preparation:

- ✉ T, “Column subset selection, matrix factorization, and eigenvalue optimization”
- ✉ T, “Paved with good intentions: Computational applications of matrix column partitions”
- ✉ ...