# Generalization Bounds for Data-Driven Numerical Linear Algebra

**Peter Bartlett**

UC Berkeley

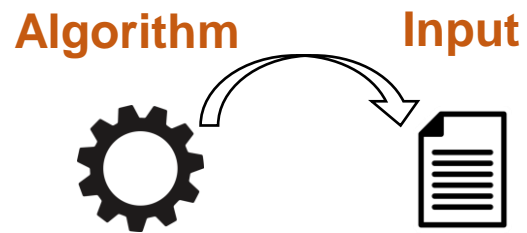**Piotr Indyk**

MIT

**Tal Wagner**
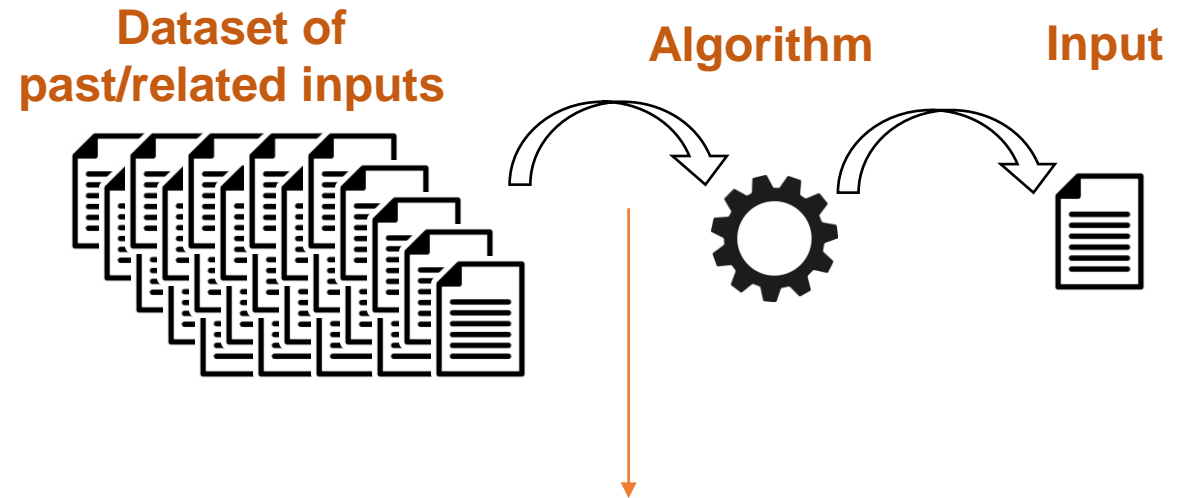
Microsoft Research

# Data-Driven Algorithms

Traditional algorithm design:

**Algorithm**         **Input**

Modern reality of algorithm design:

**Dataset of
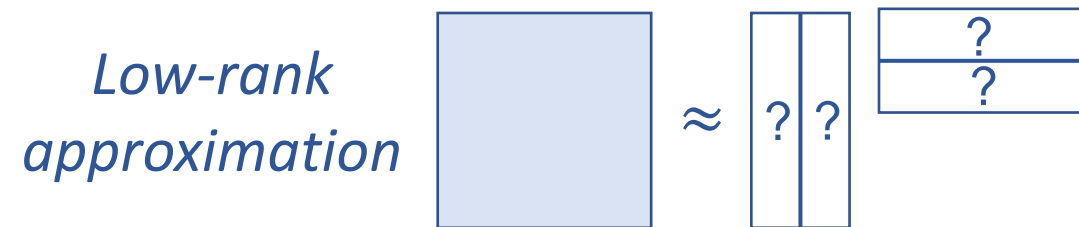past/related inputs**          **Algorithm**          **Input**

Goal:
- Use data to improve algorithm
- Automate using ML

# Numerical (or Efficient) Linear Algebra

- Problems in computational linear algebra:

*Regression*    [matrix] [?] $\approx$ [vector]     *Low-rank approximation*    [matrix] $\approx$ [?][?][? / ?]
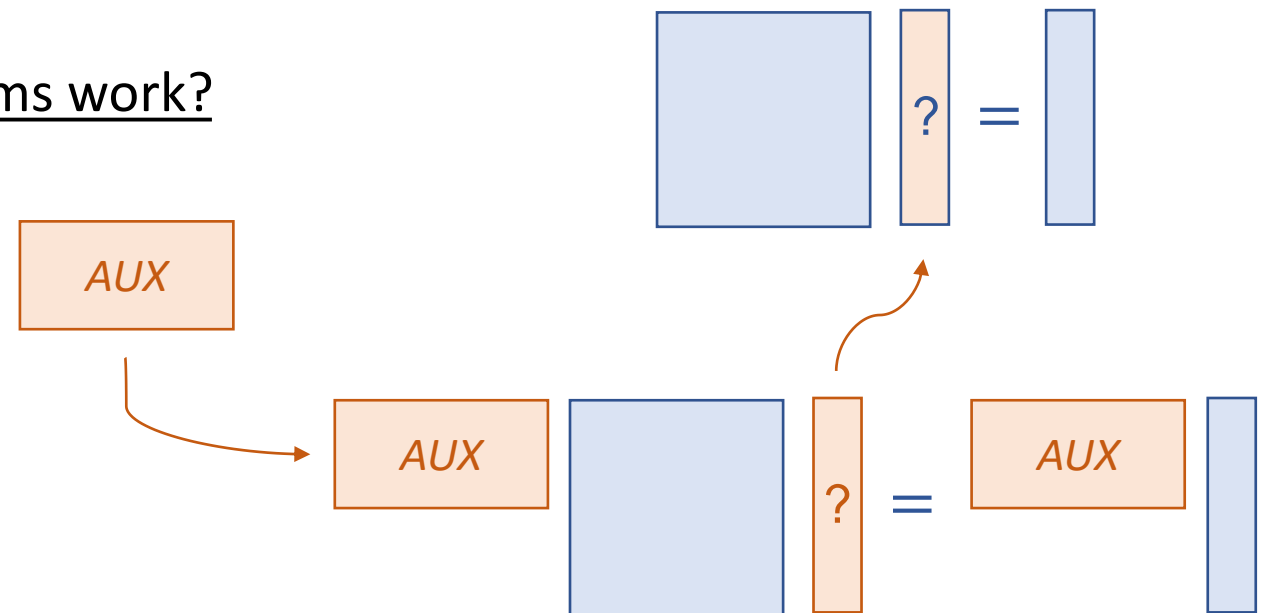
- Exact algorithms: SLOW 🥱
    - $\Omega(n^\omega)$ for an $n \times n$ matrix

- Approximate algorithms: Near-linear time! 😎
    - $\tilde{O}(n^2)$ for an $n \times n$ matrix, $\tilde{O}(\#nonzero\ entries)$ for a sparse matrix

# Data-Driven Numerical Linear Algebra

- How do numerical linear algebra algorithms work?
    - Choose auxiliary matrix
    - Use it to make problem smaller
    - Solve small problem
    - Use solution for large problem

- How do we choose the auxiliary matrix?
    - Traditionally: Either at random 🤪 or by elaborate heuristics 👷
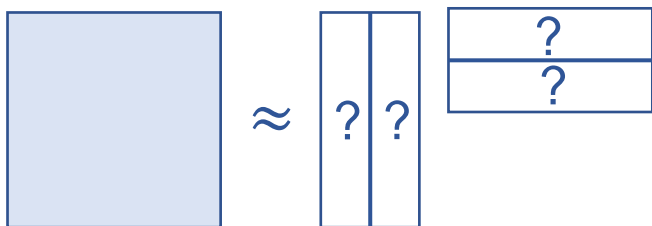    - Recently: **Learn it from data** (i.e., past inputs) 🤓
        - [Indyk-Vakilian-Yuan'19] [Ailon-Leibovich-Nair'20] [Luz-Galun-Maron-Basri-Yavneh'20] [Liu-Liu-Vakilian-Wan-Woodruff'20] [Indyk-Wagner-Woodruff'21]

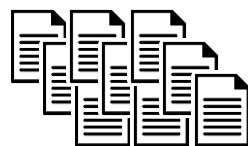# Data-Driven Numerical Linear Algebra: In Action

### Problem:

Low-rank approximation (LRA)
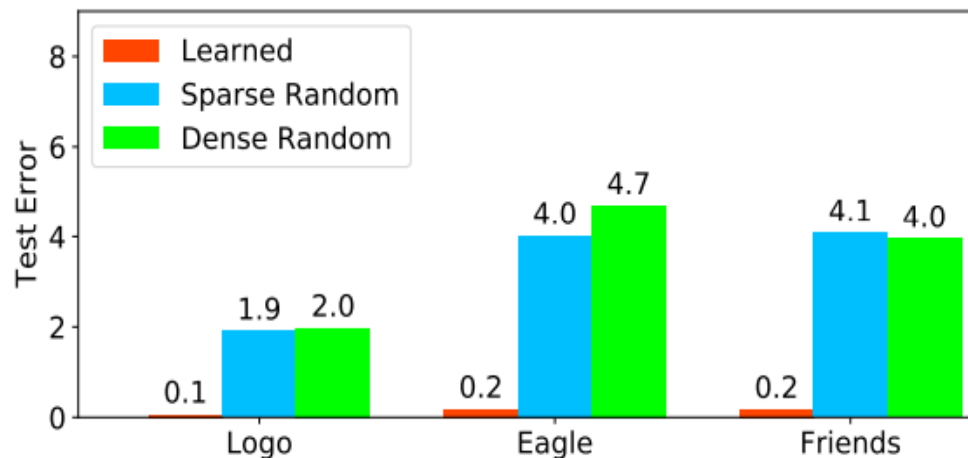
### Setting:

Learning the auxiliary matrix
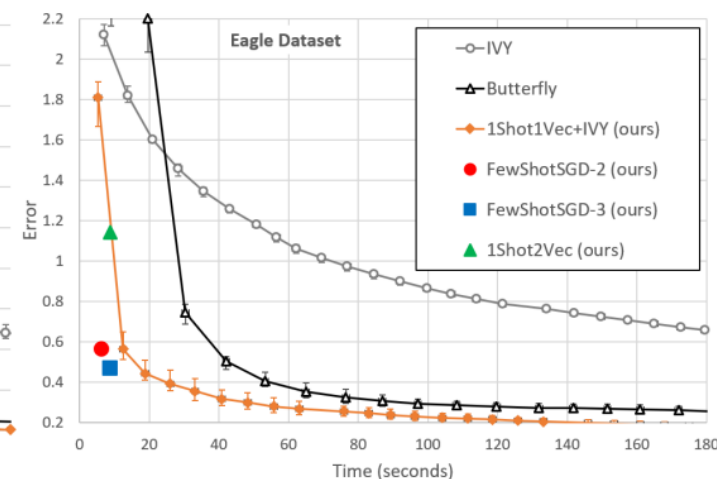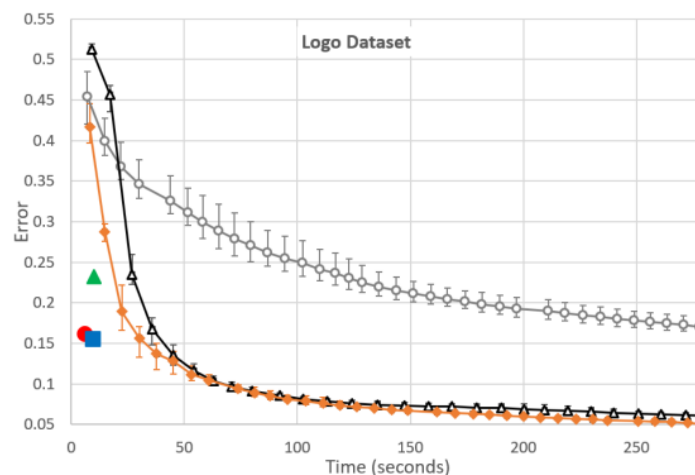


Given a training set of input matrices

Use it to learn an auxiliary matrix $S$

Evaluate $S$ by using it for fast (near-linear time) LRA on a test set of input matrices



[Indyk-Vakilian-Yuan NeurIPS'19]

[Indyk-Wagner-Woodruff NeurIPS'21]

# Data-Driven Algorithms: Theory?

- Can we **provably** learn good algorithms from past inputs?

- Gupta & Roughgarden (2016):

  - View as statistical learning problem

  - Prove upper bounds on (real-valued analogs of) VC dimension

    ⇒ PAC-learning generalization bounds on number of training samples

- **This work**: Bounds for all existing data-driven numerical linear algebra algorithms

| Reference | Algorithm | Problem | Algorithmic "family" |
|---|---|---|---|
| [Indyk-Vakilian-Yuan'19] | **IVY** | LRA | Sketching |
| [Ailon-Leibovich-Nair'20] | **Butterfly LRA** | LRA | Sketching |
| [Liu-Liu-Vakilian-Wan-Woodruff'20] | **Multisketch LRA** | LRA | Sketching |
| [Indyk-Wagner-Woodruff'21] | **Few-shot LRA** | LRA | Sketching |
| [Luz-Galun-Maron-Basri-Yavneh'20] | **Learned AMG** | Regression | Multigrid |

*\* All bounds are near-proportional to the number of learned parameters*

# Prior and Related Work

- Gupta & Roughgarden (ITCS 2016, SICOMP 2017):
  - Initiated framework
  - Upper bound technique for **greedy heuristics** and **local search** algorithms

- Balcan, DeBlasio, Dick, Kingsford, Sandholm, Vitercik (STOC 2021):
  - General upper bound technique
  - Applications for **pattern matching** and **mechanism design** algorithms
  - Does not work for the linear algebra algorithms we consider

# Review:
# Statistical Learning

(for data-driven algorithms, but also in general)

# Data-Driven Algorithms: Setting

**A loss minimization problem:**

- Inputs: $x \in X$

- Algorithms: $\mathcal{L} = \{L_\rho : \rho \in \mathbb{R}^n\}$, parameterized by $\rho \in \mathbb{R}^n$

- Losses: Identify $L_\rho$ with a map $L_\rho : X \to [0,1]$ that maps inputs to losses
  - $L_\rho(x)$ is the loss of solving for $x$ with parameters $\rho$

**Our case:** The low-rank approximation (LRA) problem

- Inputs: $X$ is the set of matrices $A \in \mathbb{R}^{n \times n}$ with $\|A\|_F = 1$
- Algorithms: $\mathcal{L} = \{L_S : S \in \mathbb{R}^{m \times n}\}$, parameterized by auxiliary matrices $S \in \mathbb{R}^{m \times n}$
- Loss: $L_S(A) = \left\|A - \widetilde{A}_S\right\|_F^2$, where $\widetilde{A}_S$ is the LRA of $A$ computed with aux. matrix $S$

# Statistical Learning and ERM

**Statistical learning:** Suppose we have a distribution $D$ over $X$

- Goal: Estimate the best parameters for $D$

$$\rho^* = \mathrm{argmin}_{\rho \in \mathbb{R}^n} \mathbb{E}_{x \in D}\left[L_\rho(x)\right]$$

- Method: Draw $s$ samples $x_1, \dots, x_s \sim D$ and use Empirical Risk Minimization (ERM)

$$\hat{\rho} = \mathrm{argmin}_{\rho \in \mathbb{R}^n} \frac{1}{s} \sum_{i=1}^{s} L_\rho(x_i)$$

- We say $\mathcal{L} = \{L_\rho : \rho \in \mathbb{R}^n\}$ is $(\epsilon, \delta)$-*learnable* with $s$ samples (by ERM) if

$$\Pr_{x_1,\dots,x_s \sim D}\left[\mathbb{E}_{x \in D}\left[L_{\hat{\rho}}(x)\right] \leq \mathbb{E}_{x \in D}\left[L_{\rho^*}(x)\right] + \epsilon\right] \geq 1 - \delta$$

- ***Question***: *What is the smallest number of samples $s$ that suffices?*

# VC-Dimension and Fat Shattering Dimension

**Definition**: Let $\mathcal{L}$ be a family of functions $X \to \{0,1\}$.

- A set $x_1, \ldots, x_s \in X$ is **shattered** by $\mathcal{L}$ if for every $I \subset \{1, \ldots, s\}$, there is $L \in \mathcal{L}$ s.t.:

$$L(x_i) = 1 \Leftrightarrow i \in I.$$

- The **VC-dimension** $\text{VCdim}(\mathcal{L})$ of $\mathcal{L}$ is the size of the largest shattered set.

**Definition**: Let $\mathcal{L}$ be a family of functions $X \to [0,1]$. Let $\gamma \geq 0$.

- A set $x_1, \ldots, x_s \in X$ is $\gamma$-**fat shattered** by $\mathcal{L}$ if there are thresholds $r_1, \ldots, r_s \in \mathbb{R}$, such that for every $I \subset \{1, \ldots, s\}$, there is $L \in \mathcal{L}$ s.t.:

$$i \in I \Rightarrow L(x_i) > r_i + \gamma \quad \text{and} \quad i \notin I \Rightarrow L(x_i) < r_i - \gamma$$

- The $\gamma$-**fat shattering dimension** $\text{fat}_\gamma(\mathcal{L})$ of $\mathcal{L}$ is the size of the largest $\gamma$-fat shattered set.

**Classical learning theory**: The sample complexity of $(\epsilon, \delta)$-learning $\mathcal{L}$ (by ERM) is proportional to the $\gamma$-fat shattering dimension with $\gamma = \Theta(\epsilon)$.
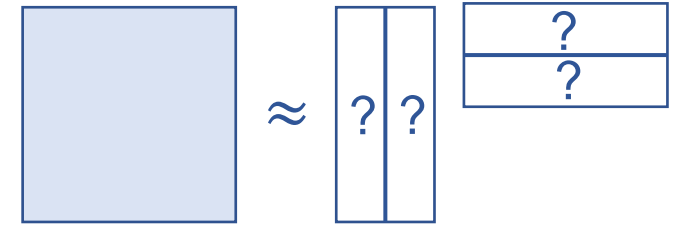
# Review:
# Fast Low-Rank Approximation

(with data-driven algorithms, but also in general)

# Low-Rank Approximation (LRA)

- **Problem:**
  - Input: $A \in \mathbb{R}^{n \times n}$ with $\|A\|_F = 1$, target rank $k \ll n$
  - Goal: $\widetilde{A}$ of rank $k$ that approximately minimizes $\left\|A - \widetilde{A}\right\|_F^2$

- **Exact solution: SVD**
  - Returns: $[A]_k$ such that $\|A - [A]_k\|_F^2 = \min_{\widetilde{A} \text{ of rank } k} \left\|A - \widetilde{A}\right\|_F^2$
  - Runtime: $O(n^\omega)$ 🥱

# Efficient Low-Rank Approximation

- The **SCW** algorithm [Sarlos'06, Clarkson-Woodruff'09,13]:
    - Pick an auxiliary matrix $S \in \mathbb{R}^{m \times n}$, where $k \le m \ll n$
    - Project $A$ onto $rows(SA)$
    - Return: Best rank-$k$ approximation of projected $A$
$$\widetilde{A}_S = \left[ A(SA)^\dagger (SA) \right]_k$$



Reminder: $M^\dagger M$ is the orthogonal projection matrix on the row space $M$

- <u>Lemma</u>: $\widetilde{A}_S$ can be computed in time $mult(S, A) + O(m^2 n)$ and space $O(mn)$.
    - By the formula: $\widetilde{A}_S = [AV]_k V^T$ where $SA = U \Sigma V^T$

- <u>Theorem</u>: If $S$ has $m = \tilde{O}(k/\epsilon)^2$ rows s.t.:
    - Each column one uniformly random non-zero
    - Each non-zero is uniform in $\{1, -1\}$

Then, whp, $\left\| A - \widetilde{A}_S \right\|_F^2 \le (1 + \epsilon) \cdot \| A - [A]_k \|_F^2$.

$S$

| ±1 | ±1 | 0 | 0 | 0 | ±1 | 0 |
|----|----|----|----|----|----|----|
| 0 | 0 | 0 | ±1 | 0 | 0 | ±1 |
| 0 | 0 | ±1 | 0 | ±1 | 0 | 0 |

$A$

# Data-Driven Low-Rank Approximation

- The **SCW** algorithm: $\text{SCW}_k(S, A) = \left[ A(SA)^\dagger(SA) \right]_k$
  - Loss: $L_S(A) = \| A - \text{SCW}_k(S, A) \|_F^2$

- **Oblivious** auxiliary matrix $S \in \mathbb{R}^{m \times n}$:
  - Each column one uniformly random non-zero
  - Each non-zero is uniform in $\{1, -1\}$

$$S = \begin{bmatrix} 1 & -1 & 0 & 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & -1 & 0 & 0 \end{bmatrix}$$

- **Data-driven** auxiliary matrix $S \in \mathbb{R}^{m \times n}$ [Indyk-Vakilian-Yuan'19]:
  - Each column one uniformly random non-zero
  - Its value is a trainable parameter (learned via SGD)

$$S = \begin{bmatrix} \rho_1 & \rho_2 & 0 & 0 & 0 & \rho_6 & 0 \\ 0 & 0 & 0 & \rho_4 & 0 & 0 & \rho_7 \\ 0 & 0 & \rho_3 & 0 & \rho_5 & 0 & 0 \end{bmatrix}$$

- **How many samples do we need to ERM-learn $\{L_S(A)\}$?**

# Our Results

**Theorem – Fat shattering dimension of SCW:**

- <u>Upper bound</u>: The $\epsilon$-fat shattering dimension of learned SCW is
$$O\big(n \cdot (m + k \log(n/k) + \log(1/\epsilon))\big),$$
with $A \in \mathbb{R}^{n \times n}, S \in \mathbb{R}^{m \times n}$, and low rank $k$.

- <u>Lower bound</u>: The $\epsilon$-fat shattering dimension of learned SCW is $\Omega(n)$, if $\epsilon < 1/(2\sqrt{k})$.

- Techniques apply to all other existing data-driven linear algebra algorithms.

**Corollary – Sample complexity of learning SCW:**

- Learning SCW with ERM requires at most $\tilde{O}(\epsilon^{-2} n \cdot m)$ samples.

- Learning SCW with ERM requires at least $\Omega(\epsilon^{-2} n/k)$ samples.

- Learning SCW with any method requires at least $\Omega(n + \epsilon^{-1})$ samples.

# Remarks

- <u>What if training set has nothing to do with test set?</u>
  - **Safeguarding** [Indyk-Vakilian-Yuan'19]:
    - Vertically augment the learned $S$ a random $S'$
    - Number of rows is only doubled
    - <u>Guarantees</u>: $S$ is at least as good as $S'$ on <u>any</u> input matrix $A$

- <u>What about the running time of computing the best $S$ for the sample?</u>
  - In practice: Use stochastic gradient descent (SGD) on training set
  - Not known to provably converge to empirical risk minimizer

# Proof Overview

# The Goldberg-Jerrum'95 Framework

**Definition:** A **GJ-algorithm** is a deterministic algorithm on real-valued inputs, with two types of operations:
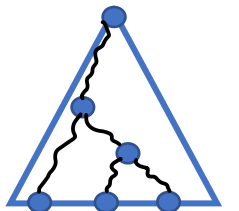
- <u>Arithmetic</u>: $v'' = v \odot v'$ where $\odot \in \{+, -, \times, \div\}$
- <u>Conditional</u>: "if $v \geq 0$ then ... else ..."

  Where $v, v'$ are either inputs or values previously computed by the algorithm.

**Theorem** [Goldberg-Jerrum'95]:

- Suppose there is a GJ-algorithm that takes $x \in X$, $\rho \in \mathbb{R}^n$ and $r \in \mathbb{R}$, and returns TRUE iff $L_\rho(x) \geq r$. Suppose it has running time $T$.

- Then, $\forall \gamma$, the $\gamma$-fat shattering dimension of $\mathcal{L} = \{L_\rho : \rho \in \mathbb{R}^n\}$ is $O(nT)$.

  <u>Proof sketch</u>: The GJ algorithm partitions $\mathbb{R}^n$ into constant sign regions with polynomial boundaries. Classical theorems on polynomials [Milnor'64, Warren'68] bound the number of sign regions.

# Goldberg-Jerrum: First Attempt

- Goal: GJ-algorithm for the SCW loss, $L_S(\boldsymbol{A}) = \left\| \boldsymbol{A} - \left[ \boldsymbol{A}(\boldsymbol{SA})^\dagger (\boldsymbol{SA}) \right]_k \right\|_F^2$.

- Need two steps:
  1. Projection (compute $\boldsymbol{M}^\dagger \boldsymbol{M}$ for a matrix $\boldsymbol{M}$)
  2. Best rank-$k$ approximation (computing $[\boldsymbol{M}]_k$ for a matrix $\boldsymbol{M}$)

- Problem: How to compute $[\boldsymbol{M}]_k$ with a GJ-algorithm (only arithmetic operations)?
- Solution: Approximate by the Power Method
  - $[\boldsymbol{M}]_k \approx \boldsymbol{Z}\boldsymbol{Z}^\dagger \boldsymbol{M}$ with $\boldsymbol{Z} = (\boldsymbol{M}\boldsymbol{M}^T)^q \boldsymbol{M}\boldsymbol{\Pi}$ and gaussian $\boldsymbol{\Pi} \in \mathbb{R}^{n \times k}$  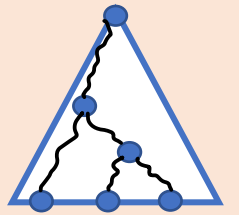[Rokhlin-Szlam-Tygert'10, Halko-Martinsson-Tropp'11, Boutsidis-Drineas-MagdonIsmail'14, Woodruff'14, Witten-Candes'15, Musco-Musco'15]

- However:
  - Power Method is *slow* 🥱 $q$ iterations take $qkn^{O(1)}$ time (here $q = O(\epsilon^{-1}\log(n/\epsilon))$)
  - Power Method is *randomized* 🤪 derandomizing $\boldsymbol{\Pi}$ takes $(n/k)^k$ time

# Refined Goldberg-Jerrum: Definitions

**Definition:** A **GJ-algorithm** is a deterministic algorithm on real-valued inputs, with two types of operations:

- <u>Arithmetic</u>: $v'' = v \odot v'$ where $\odot \in \{+, -, \times, \div\}$
- <u>Conditional</u>: "if $v \geq 0$ then … else …"

Where $v, v'$ are either inputs or values previously computed by the algorithm.

<u>Observe</u>: Every value computed by a GJ-algorithm is a rational function of the inputs.

**Definition:**

- The **degree** of a GJ-algorithm is the maximum degree of any rational function it computes.
- The **predicate complexity** of a GJ-algorithm is the number of **_distinct_** rational functions in its conditional statements.

# Refined Goldberg-Jerrum: Theorem
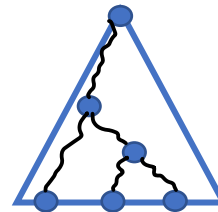
**Theorem**: Suppose we have,

- A GJ-algorithm that takes $x \in X$, $\rho \in \mathbb{R}^n$ and $r \in \mathbb{R}$, and returns TRUE iff $L_\rho(x) \geq r$.
- Suppose it has **degree** $\Delta$ and **predicate complexity** $P$.
- Then, $\forall \epsilon$, the $\epsilon$-fat shattering dimension of $\mathcal{L} = \{L_\rho : \rho \in \mathbb{R}^n\}$ is $O(n \log(\Delta P))$.

<u>Observe</u>: Runtime $T$ implies $\Delta, P \leq 2^T$.

- Thus, refines the previous theorem, *"runtime $T \Rightarrow$ fat-dim $O(nT)$"*.

<u>Why does it help?</u>

- $q$ Power Method iterations: time $qkn^{O(1)}$, **but** degree $O(q)$
  - $nT = qkn^{O(1)}$ **but** $n \log \Delta = O(n \log q)$
- Minimum of $t$ numbers: time $O(t)$, **but** predicate complexity $\binom{t}{2}$
  - $nT = O(nt)$ **but** $n \log P = O(n \log t)$
  - Derandomizing the power method: $t = \binom{n}{k}$, thus $nT = O(n(n/k)^k)$ **but** $n \log P = O(nk \log(n/k))$

# Refined GJ-Algorithm for SCW

- New goal:

  - GJ-algorithm for the SCW loss, $L_S(A) = \left\| A - \left[ A(SA)^\dagger (SA) \right]_k \right\|_F^2$

  - With efficient **degree** and **predicate complexity**.

- Need two steps:

  1. Projection (compute $M^\dagger M$ for a matrix $M$)
  2. Best rank-$k$ approximation (computing $[M]_k$ for a matrix $M$)

# Step 1: Computing Projection Matrices

**Lemma 1:** Given $M \in \mathbb{R}^{m \times n}$, there is a GJ-algorithm for computing $M^\dagger M$ with degree $O(m)$ and predicate complexity $2^m$.

**Proof**:

- <u>Fact 1</u>: If the rows of $N$ form a basis for the rows of $M$, then $M^\dagger M = N^T (NN^T)^{-1} N$.

- <u>Fact 2</u> (e.g., [Csanky'76]): There are algorithms that use only arithmetic operations for
    **(i)** checking if a matrix has full rank,
    **(ii)** inverting an invertible matrix.

   Their degree (as GJ-algorithms) for $m \times m$ matrices is $O(m)$.

- **GJ-Algorithm**: Try all $2^m$ subsets of rows of $M$ (predicate complexity $2^m$) to find a basis $N$ (use Fact 2**(i)** to check it is a basis). Invert $NN^T$ (with Fact 2**(ii)**). Return $N^T (NN^T)^{-1} N$.

# Step 2: Derandomized Power Method

**Lemma 2**: Given $M \in \mathbb{R}^{n \times n}$, there is a GJ-algorithm for $\epsilon$-approximating $[M]_k$, with degree $O(k\epsilon^{-1} \log(n/\epsilon))$ and predicate complexity $2^k \binom{n}{k}^2$.

*Remark: Approximation is why we needed the gap $\gamma$ in the definition of fat shattering.*

**Proof** (sketch):

- Fact 3: Starting with a gaussian $\Pi \in \mathbb{R}^{n \times k}$, $q = O(\epsilon^{-1} \log(n/\epsilon))$ Power Methods iterations suffice to $\epsilon$-approximate $[M]_k$ [Musco-Musco'15], by the formula:

$$[M]_k \approx_\epsilon ZZ^\dagger M \text{ where } Z = (MM^T)^q M\Pi$$

- Fact 4: $\Pi$ can be derandomized with $k$-subsets of the standard basis in $\mathbb{R}^n$.

- **GJ algorithm:** Try all $\binom{n}{k}$ subset of $\mathbb{R}^n$ as the initial matrix $\Pi$. For each one, compute $ZZ^\dagger M$ (using the previous lemma for $ZZ^\dagger$, which blows up the degree by $k$ and the predicate complexity by $2^k$). pick the one that minimizes the LRA error $\left\| M - ZZ^\dagger M \right\|_F^2$.