Randomized algorithms for the approximation of matrices

Luis Rademacher The Ohio State University Computer Science and Engineering (joint work with Amit Deshpande, Santosh Vempala, Grant Wang)

Two topics

- Low-rank matrix approximation (PCA).
- Subset selection: Approximate a matrix using another matrix whose columns lie in the span of a few columns of the original matrix.

Motivating example: DNA microarray

- [Drineas, Mahoney] *Unsupervised* feature selection for classification
 - Data: table of gene expressions (features) v/s patients
 - Categories: cancer types
 - Feature selection criterion: Leverage scores (importance of a given feature in determining top principal components)
- Empirically: Leverage scores are correlated with "information gain", a supervised measure of influence. Somewhat unexpected.
- Leads to clear separation (clusters) from selected features.

In matrix form:

 A is m × n matrix, m patients, n genes (features), find

 $A \approx CX$,

where the columns of C are a few columns of A (so $X = C^+A$).

• They prove error bounds when columns of *C* are selected at random according to leverage scores (importance sampling).

(P1) Matrix approximation

- Given *m*-by-*n* matrix, find low rank approximation ...
- ... for some norm:

 $- \|A\|_F^2 = \sum_{ij} A_{ij}^2 \quad \text{(Frobenius)}$ $- \|A\|_2 = \sigma_{\max}(A) = \max_x \|Ax\| / \|x\| \quad \text{(spectral)}$

Geometric view

- Given points in Rⁿ, find subspace close to them.
- Error: Frobenius norm corresponds to sum of squared distances.

Classical solution

- Best rank-k approximation A_k in $\|\cdot\|_F^2$ and $\|\cdot\|_2$:
 - Top k terms of singular value decomposition (SVD): if $A = \sum_{i} \sigma_{i} u_{i} v_{i}^{T}$ then $A_{k} = \sum_{i=1}^{k} \sigma_{i} u_{i} v_{i}^{T}$
- Best k-dim. subspace: rowspan(A_k), i.e.
 Span of top k eigenvectors of A^TA.
- Leads to iterative algorithm. Essentially, in time mn^2 .

Want algorithm

- With better error/time trade-off.
- Efficient for very large data:
 - Nearly linear time
 - Pass efficient: if data does not fit in main memory, algorithm should not need random access, but only a few sequential passes.
- Subspace equal to or contained in the span of a few rows (actual rows are more informative than arbitrary linear combinations).

Idea [Frieze Kannan Vempala]

- Sampling rows.
 Uniform does not work (e.g. a single non-zero entry)
- By "importance": sample s rows, each independently with probability proportional to squared length.

[FKV]

Theorem 1. Let S be a sample of k/ϵ rows where

 $\mathbb{P}(row \ i \ is \ picked) \propto ||A_i||^2.$

Then the span of S contains the rows of a matrix \tilde{A} of rank k for which

$$\mathsf{E}(\|A - \tilde{A}\|_F^2) \cdot \|A - A_k\|_F^2 + \epsilon \|A\|_F^2.$$

This can be turned into an efficient algorithm: 2 passes, complexity $O(kmn/\epsilon)$.

(to compute \tilde{A} , SVD in span of S, which is fast because n becomes k/ϵ).

One drawback of [FKV]

 Additive error can be large (say if matrix is nearly low rank).
 Prefer relative error, something like

$$||A - \tilde{A}||_F^2 \cdot (1 + \epsilon) ||A - A_k||_F^2.$$

Several ways:

- [Har-Peled '06] (first linear time relative approximation)
- [Sarlos '06]: Random projection of rows onto a $O(k/\epsilon)$ -dim. subspace. Then SVD.
- [Deshpande R Vempala Wang '06] [Deshpande Vempala '06] Volume sampling (rough approximation) + adaptive sampling.

Some more relevant work

- [Papadimitriou Raghavan Tamaki Vempala '98]: Introduced random projection for matrix approximation.
- [Achlioptas McSherry '01][Clarkson Woodruff '09] One-pass algorithm.
- [Woolfe Liberty Rokhlin Tygert '08] [Rokhlin Szlam Tygert '09] Random projection + power iteration to get *very fast practical* algorithms. Read survey [Halko Martinsson Tropp '09].
- D'Aspremont, Drineas, Ipsen, Mahoney, Muthukrishnan, ...

(P2) Algorithmic Problems: Volume sampling and subset selection

- Given *m*-by-*n* matrix, pick set of k rows at random with probability proportional to squared volume of k-simplex spanned by them and origin. [DRVW] (equivalently, squared
 - volume of parallelepiped determined by them)





Volume sampling

• Let S be k-subset of rows of A

- [k! vol(conv(0, A_s))]² = vol($\Box(A_s)$)² = det($A_s A_s^T$)(*)

- volume sampling for A is equivalent to: pick k by k principal minor "S × S" of A A^T with prob. proportional to $det(A_S A_S^T)$ $A A^T \stackrel{\leq}{=} \left\{ A_S A_S^T \right\}$
- For(*): complete A_s to a square matrix B by adding orthonormal rows, orthogonal to span(A_s).

$$\operatorname{vol}\Box(A_S)^2 = (\det B)^2 = \det(BB^T) = \det\begin{pmatrix}A_S A_S^T & 0\\ 0 & I\end{pmatrix} = \det(A_S A_S^T)$$

Original motivation:

- **Relative error** low rank matrix approximation [DRVW]:
 - S: k-subset of rows according to volume sampling
 - A_k: best rank-k approximation, given by principal components (SVD)
 - π_{s} : projection of rows onto rowspan(A_s)

$$\implies \mathsf{E}_{S}(\|A - \pi_{S}(A)\|_{F}^{2}) \cdot (k+1)\|A - A_{k}\|_{F}^{2}$$

- Factor "k+1" is best possible [DRVW]
- Interesting existential result (there exist k rows...). Alg.?
- Lead to linear time, pass-efficient algorithm for relative approximation of A_k [DV]. (1+ ϵ) in span of O^{*}(k/ ϵ) rows

Where does volume sampling come from?

• No self-respecting architect leaves the scaffolding in place after completing the building.

Gauss?

Where does volume sampling come from?

- Idea:
 - For picking k out of k + 1 points, k with maximum volume is optimal.
 - For picking 1 out of *m*, random according to squared length is better than max. length.
 - For k out of m, this suggest volume sampling.



Where does volume sampling come from?

- Why does the algebra work? Idea:
 - When picking 1 out of *m* random according to squared length, expected error is sum of squares of areas of triangles;

$$E(error) = \sum_{s} \frac{\|A_{s}\|^{2}}{\sum_{t} \|A_{t}\|^{2}} \sum_{i} d(A_{i}, span(A_{s}))^{2}$$

$$\int_{t} (A_{i}, max A_{s}) \left\{ \begin{array}{c} A_{i} \\ A_{i} \\ A_{s} \end{array} \right\}$$

– This sum corresponds to certain coefficient of the characteristic polynomial of AA^T

Later motivation [BDM,...]

- (row/column) Subset selection.
 A refinement of principal component analysis:
 Given a matrix A,
 - PCA: find k-dim subspace V that minimizes $\|A \pi_V(A)\|_F^2$
 - Subset selection: find V spanned by k rows of A.
 - Seemingly harder, combinatorial flavor.

(π projects rows)

Why subset selection?

- PCA unsatisfactory:
 - top components are linear combinations of rows (all rows, generically). Many applications prefer *individual*, most relevant rows, e.g.:
 - feature selection in machine learning
 - linear regression using only most relevant independent variables
 - out of thousands of genes, find a few that explain a disease

Known results

- [Deshpande-Vempala] Polytime k!-approximation to volume sampling, by adaptive sampling:
 - pick a row with probability proportional to squared length
 - project all rows orthogonal to it
 - repeat
- Implies for random k-subset S with that distribution:

$$\mathsf{E}_{S}(\|A - \pi_{S}(A)\|_{F}^{2}) \cdot (k+1)! \|A - A_{k}\|_{F}^{2}$$

Known results

- [Boutsidis, Drineas, Mahoney] Polytime randomized algorithm to find k-subset S: $\|A - \pi_S(A)\|_F^2 \cdot O(k^2 \log k) \|A - A_k\|_F^2$
- [Gu-Eisenstat] Deterministic algorithm,

$$||A - \pi_S(A)||_2^2 \cdot (1 + f^2 k(n-k)) ||A - A_k||_2^2$$

in time O((m + n log_f n)n²) Spectral norm:

 $||A||_2 = \sup_{x \in \mathbb{R}^n} ||Ax|| / ||x||$

Known results

- Remember, volume sampling equivalent to sampling k by k minor "S × S" of AA^T with probability proportional to (*) $det(A_S A_S^T)$
- [Goreinov, Tyrtishnikov, Zamarashkin] Maximizing
 (*) over S is good for subset selection.
- [Çivril , Magdon-Ismail] [see also Koutis '06] But maximizing is NP-hard, even approximately to within an exponential factor.

Results

- Volume sampling: Polytime exact alg. O(mn^ω log n) (arithmetic ops.) (some ideas earlier in [Houges Krishnapur Peres Virág])
- Implies alg. with optimal approximation for subset selection under Frobenius norm. Can be derandomized by conditional expectation. O(mn^ω log n)
- 1+ε approximations to the previous 2 algorithms in nearly linear time, using volume-preserving random projection [M Z].

Results

• Observation: Bound in Frobenius norm easily implies bound in spectral norm:

$$||A - \pi_S(A)||_2^2 \cdot ||A - \pi_S(A)||_F^2$$

$$\cdot (k+1)||A - A_k||_F^2$$

$$\cdot (k+1)(n-k)||A - A_k||_2^2$$

using

$$||A||_F^2 = \sum_i \sigma_i^2 \qquad ||A||_2^2 = \sigma_{\max}^2$$

 $\sigma_{\max} = \sigma_1 \ge \sigma_2 \ge \cdots \ge \sigma_n \ge 0$ are the singular values of A

Comparison for subset selection

Find S s.t. $||A - \pi_S(A)||_?^2 \cdot ?||A - A_k||_?^2$

	Frobenius norm sq	Spectral norm sq	Time (assuming m>n) ω: exponent of matrix mult.	
[D R V W]	k+1		Existential	
[Despande Vempala]	(k+1)!		kmn	R
[Gu Eisenstat]		1+k(n-k)	Existential	
[Gu Eisenstat]		1+f ² k(n-k)	((m + n log _f n)n ²	D
[Boutsidis Drineas Mahoney]	k² log k	k ² (n-k) log k (F implies spectral)	mn²	R
[Desphande R]	k+1 (optimal)	(k+1)(n-k)	kmn ^ω log n	D
[Desphande R]	(1+ε)(k+1)	(1+ε) (k+1)(n-k)	$O^*(mnk^2/\epsilon^2 + m k^2 \omega + 1/\epsilon^2 \omega)$	R

 Want (w.l.o.g.) k-tuple S of rows of m by n matrix A with probability

$$\frac{\det(A_S A_S^T)}{\sum_{S' \in [m]^k} \det(A_{S'} A_{S'}^T)}$$

 Idea: pick S=(S₁, S₂, ..., S_k) in sequence. Need marginal distribution of S₁ to begin:

$$\mathbb{P}(S_1 = i) = \frac{\text{tuples with } S_1 = i}{\text{all tuples}} = \frac{\sum_{S' \in [m]^k, S'_1 = i} \det(A_{S'} A_{S'}^T)}{\sum_{S' \in [m]^k} \det(A_{S'} A_{S'}^T)}$$

• Remember characteristic polynomial:

$$p_{AA^T}(x) = \det(xI - AA^T) = \sum_i c_i (AA^T) x^i$$
$$|c_{m-k}(AA^T)| = \sum_{S \subseteq [m], |S|=k} \det(A_S A_S^T)$$

• So, for
$$C_i = A - \pi_{A_i}(A)$$

$$\mathbb{P}(S_1 = i) = \frac{\sum_{S' \in [m]^k, S'_1 = i} \det(A_{S'}A_{S'}^T)}{\sum_{S' \in [m]^k} \det(A_{S'}A_{S'}^T)}$$

$$= \frac{(k-1)! ||A_i||^2 \sum_{S' \subseteq [m], |S'| = k} \det((C_i)_{S'}(C_i)_{S'}^T)}{k! \sum_{S' \subseteq [m], |S'| = k} \det(A_{S'}A_{S'}^T)}$$

$$= \frac{||A_i||^2 |c_{m-k+1}(C_iC_i^T)|}{k|c_{m-k}(AA^T)|}$$
intuition for numerator:

$$|\Box(A_1, A_2, A_3)| = ||A_1|| |\Box(\pi_{A_1^\perp}(A_2, A_3))|$$

- So: $\mathbb{P}(S_1 = i) = \frac{\|A_i\|^2 |c_{m-k+1}(C_i C_i^T)|}{k |c_{m-k}(AA^T)|}$ $C_i = A - \pi_{A_i}(A)$
- Can be computed in polytime.
- After S₁, project rows orthogonal to picked row, repeat marginal computation for S₂,... (use intuition for numerator)
- "flops": k * m * (m²n + m^ωlog m)

Faster: (we assume m > n)

-use $p(AA^T) = x^{m-n}p(A^TA)$

as A^TA is n by n (smaller than m by m)

– use rank-1 updates:

$$C_i = A - \frac{1}{\|A_i\|^2} A A_i A_i^T,$$

 $C_i^T C_i = A^T A - \frac{A^T A A_i A_i^T}{\|A_i\|^2} - \frac{A_i A_i^T A^T A}{\|A_i\|^2} + \frac{A_i A_i^T A^T A A_i A_i^T}{\|A_i\|^4}.$ - total flops: mn² + km(n² + n^{\omega} log n) = k m n^{\omega} log n

Even faster

 Volume sampling only cares about volumes of k-subsets,

⇒ can get 1+ε approximation using a volume preserving random projection [Magen, Zouzias] (generalizing Johnson Lindenstrauss, not same as Feige).

Even faster

- [Magen Zouzias]: For any A ∈ R^{m×n}, 1 ≤ k ≤ n, ε < 1/2 there is a d =O(k² ε⁻² log m) s.t. for all S, k-subset of [m]: det A_SA^T_S · det Ã_SÃ^T_S · (1 + ε) det A_SA^T_S,
 à = AG, G ∈ R^{n×d} random Gaussian matrix, scaled
 k=1 is JL
 - This as preprocessing implies 1+ε volume sampling in time
 O^{*}(mnk²/ε² + m k^{2ω + 1}/ε^{2ω})

Recent news

- [Boutsidis, Drineas, Magdon-Ismail (FOCS '11)]: improved subset selection using [BSS].
- [Guruswami, Sinop (SODA '12)]:
 - Volume sampling in time $O(kmn^2)$
 - Relative $(1 + \epsilon)$ matrix approximation with one round of $r = k - 1 + k/\epsilon$ rows of volume sampling. More precisely, for S a sample of size $r \ge k$ according to volume sampling:

$$\mathsf{E}_{S}(\|A - \pi_{S}(A)\|_{F}^{2}) \cdot \frac{r+1}{r+1-k} \|A - A_{k}\|_{F}^{2}$$

Open question

• For a subset of rows *S* according to volume sampling, lower bound (in expectation?):

$\sigma_{\min}(A_S)$

– Want A_S to be well conditioned.