An Efficient, Sparsity-Preserving, Online Algorithm for Low-Rank Approximation

Dave Anderson and Ming Gu

Department of Mathematics, University of California, Berkeley

Background

Goal: low-rank approximations (LRAs) that are

- high-quality
- efficient
- sparsity-preserving

- can be updated

can be continued

spectrum-preserving highlights features

Approach: use the LU decomposition

$$\Pi_{1}\mathbf{A}\Pi_{2}^{T} = \begin{pmatrix} k & \mathbf{L}_{11} & \mathbf{L}_{21} & \mathbf{I}_{m-k} \end{pmatrix} \begin{pmatrix} \mathbf{U}_{11} & \mathbf{U}_{12} \\ \mathbf{S} \end{pmatrix} \approx \begin{pmatrix} \mathbf{L}_{11} \\ \mathbf{L}_{21} \end{pmatrix} (\mathbf{U}_{11} \mathbf{U}_{12}) \stackrel{\text{def}}{=} \widehat{\mathbf{L}}\widehat{\mathbf{U}}$$

Problem: the LU decomposition is unstable and both Π_1 and Π_2^T are needed. Computing both is expensive

Solution: use randomization to approximate complete pivoting by cheaply calculating Π_1 and Π_2^T

Previous Work

Deterministic Algorithms

• Truncated SVD, QR, ID, Constrained Nuclear Norm, others Randomized Algorithms for Faster LRA:

Black Box Algorithms [2, 6]

$$\mathbf{Q}\mathbf{R} = \mathbf{A}\Omega$$

$$\mathbf{A} \approx \mathbf{Q}^T \mathbf{Q} \mathbf{A}$$

Sampling Algorithms [4]

$$\mathbf{A} pprox \mathbf{C} \left(\mathbf{C}^\dagger \mathbf{A} \mathbf{R}^\dagger
ight) \mathbf{R}$$

Complexity:

$$O\left(\operatorname{nnz}\left(\mathbf{A}\right)\right) + n \cdot \operatorname{poly}\left(\frac{k}{\epsilon}\right)$$

• Error bound:

$$\|\mathbf{A} - \widehat{\mathbf{A}}\| \le (1 + \epsilon) \|\mathbf{A} - \mathbf{A}_k\|$$

...but $\epsilon \not\to 0$ because of the complexity

Truncated LU Theory

General Truncated LU Bounds

Theorem

Let $(\cdot)_s$ denote the rank-s truncated SVD for $s \leq k \ll m, n$. Then for any truncated LU factorization with Schur complement S:

$$\|\Pi_1 \mathbf{A} \Pi_2^T - \widehat{\mathbf{L}} \widehat{\mathbf{U}}\| = \|\mathbf{S}\|$$

for any norm, and

$$\|\Pi_1 \mathbf{A} \Pi_2^T - (\widehat{\mathbf{L}}\widehat{\mathbf{U}})_s\|_2 \le 2\|\mathbf{S}\|_2 + \sigma_{s+1}(\mathbf{A}).$$

Theorem

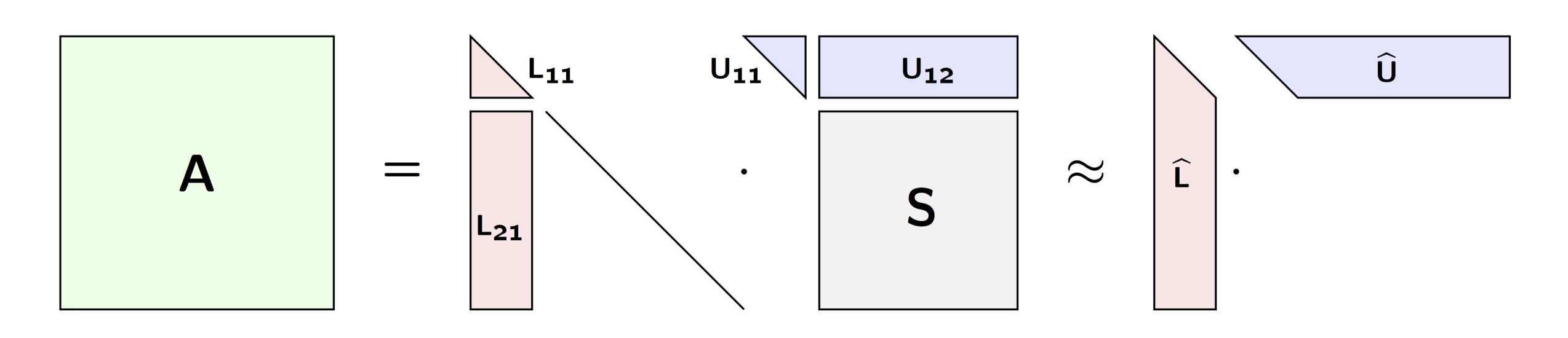
For a general rank-k truncated LU decomposition, we have for all $1 \leq 1$

$$\sigma_j(\mathbf{A}) \le \sigma_j(\widehat{\mathbf{L}}\widehat{\mathbf{U}}) \left(1 + \left(1 + \frac{\|\mathbf{S}\|_2}{\sigma_k(\widehat{\mathbf{L}}\widehat{\mathbf{U}})}\right) \frac{\|\mathbf{S}\|_2}{\sigma_j(\mathbf{A})}\right).$$

Theorem

$$\|\Pi_1 \mathbf{A} \Pi_2^T - \widehat{\mathbf{L}} \mathbf{M} \widehat{\mathbf{U}}\|_2 \le 2 \|\mathbf{S}\|_2$$
$$\|\Pi_1 \mathbf{A} \Pi_2^T - \widehat{\mathbf{L}} \mathbf{M} \widehat{\mathbf{U}}\|_F \le \|\mathbf{S}\|_F.$$

The Truncated LU Factorization



Spectrum-Revealing LU (SRLU)

SRLU = TRLUCP + SRP

Algorithm 1 Truncated Randomized LU with Complete Pivoting (TRLUCP)

Inputs: Data matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, target rank k, block size b, oversampling parameter $p \geq b$, random Gaussian matrix $\Omega \in \mathbb{R}^{p \times m}$, $\widehat{\mathbf{L}}$ and $\widehat{\mathbf{U}}$ are initially 0 matrices

Calculate random projection $\mathbf{R} = \Omega \mathbf{A}$

for $j = 0, b, 2b, \dots, k - b$ do

Perform column selection algorithm on ${f R}$ and swap columns of ${f A}$

Update block column of $\widehat{\mathbf{L}}$

Perform block LU with partial row pivoting and swap rows of A

Update block row of $\widehat{\mathbf{U}}$

Update **R**

end for

How did we do?

Let α be the largest element in the Schur complement:

$$\Pi_1 \mathbf{A} \Pi_2^T = \widehat{\mathbf{L}} \widehat{\mathbf{U}} + \begin{pmatrix} 0 & 0 \\ 0 & \mathbf{S} \end{pmatrix} = \begin{pmatrix} \mathbf{L}_{11} \\ \ell^T \\ \mathbf{L}_{31} \end{pmatrix} (\mathbf{U}_{11} \ u \ \mathbf{U}_{13}) + \begin{pmatrix} 0 & 0 \\ 0 & \begin{pmatrix} \alpha & s_{12}^T \\ s_{21} \ \mathbf{S}_{22} \end{pmatrix})$$

Compare what we have to what we missed:

Define:
$$\overline{\mathbf{A}}_{11} \stackrel{\text{def}}{=} \begin{pmatrix} \mathbf{L}_{11} \\ \ell^T \end{pmatrix} \begin{pmatrix} \mathbf{U}_{11} \ u \\ \alpha \end{pmatrix}$$
 and test: $\|(\overline{\mathbf{A}}_{11})^{-1}\|_{\text{MAX}} \stackrel{?}{\leq} \frac{f}{|\alpha|}$

Rewrite as an algorithm:

Algorithm 2 Spectrum-Revealing Pivoting (SRP)

Input: Truncated LU factorization $\mathbf{A} \approx \widehat{\mathbf{L}}\widehat{\mathbf{U}}$, tolerance f > 1while $\|\overline{\mathbf{A}}_{11}^{-1}\|_{\max} > \frac{f}{|\alpha|}$ do

Set α to be the largest element in **S** (or find an approximate α using **R**) Swap row and column containing α with row and column of largest element in

Update truncated LU factorization

end while

Improving Quality

Use a CUR technique:

$$\Pi_1 \mathbf{A} \Pi_2^T \approx \widehat{\mathbf{L}} \underbrace{(\widehat{\mathbf{L}}^{\dagger} \Pi_1 \mathbf{A} \Pi_2^T \widehat{\mathbf{U}}^{\dagger})}_{\mathbf{M}} \widehat{\mathbf{U}}$$

M is dense, but small

Online Updating

Begin with an SRLU factorization:

$$\mathbf{A}\mathbf{\Pi}_2^T = egin{pmatrix} \mathbf{L}_{11} \ \mathbf{L}_{21} \ \mathbf{I} \end{pmatrix} egin{pmatrix} \mathbf{U}_{11} \ \mathbf{U}_{12} \ \mathbf{S} \end{pmatrix}$$

Receive new data $\mathbf{B}\Pi_2^T = (\mathbf{B}_1 \ \mathbf{B}_2)$ and combine:

$$egin{pmatrix} egin{pmatrix} \mathbf{I}_{1} \mathbf{A} \mathbf{\Pi}_{2}^{T} \\ \mathbf{B} \mathbf{\Pi}_{2}^{T} \end{pmatrix} = egin{pmatrix} \mathbf{L}_{11} & & & \\ \mathbf{L}_{21} & \mathbf{I} & & \\ \mathbf{B}_{1} \mathbf{U}_{11}^{-1} & & \mathbf{I} \end{pmatrix} egin{pmatrix} \mathbf{U}_{11} & \mathbf{U}_{12} & & \\ & \mathbf{S} & & \\ \mathbf{B}_{2} - \mathbf{B}_{1} \mathbf{U}_{11}^{-1} \mathbf{U}_{12} \end{pmatrix}$$

Perform Spectrum-Revealing Pivoting

Theoretical Results

Theorem

SRLU Bounds

For $j \leq k$ and $\gamma = O(fk\sqrt{mn})$, SRP produces a rank-k SRLU factor*ization with*

$$\|\Pi_1 \mathbf{A} \Pi_2^T - \widehat{\mathbf{L}} \widehat{\mathbf{U}}\|_2 \le \gamma \sigma_{k+1} (\mathbf{A}),$$

$$\|\Pi_1 \mathbf{A} \Pi_2^T - (\widehat{\mathbf{L}} \widehat{\mathbf{U}})_i\|_2 \le \sigma_{j+1} (\mathbf{A}) \left(1 + 2\gamma \frac{\sigma_{k+1}(\mathbf{A})}{\sigma_{j+1}(\mathbf{A})}\right)$$

Theorem

For $1 \le j \le k$, SRP produces a rank-k SRLU factorization with

$$\frac{\sigma_{j}(\mathbf{A})}{1 + \tau \frac{\sigma_{k+1}(\mathbf{A})}{\sigma_{j}(\mathbf{A})}} \leq \sigma_{j}(\widehat{\mathbf{L}}\widehat{\mathbf{U}}) \leq \sigma_{j}(\mathbf{A}) \left(1 + \tau \frac{\sigma_{k+1}(\mathbf{A})}{\sigma_{j}(\mathbf{A})}\right)$$

for $\tau \leq O\left(mnk^2f^3\right)$.

Theorem

$$\|\Pi_1 \mathbf{A} \Pi_2^T - \widehat{\mathbf{L}} \mathbf{M} \widehat{\mathbf{U}}\|_2 \le 2\gamma \sigma_{k+1} (\mathbf{A})$$

$$\|\Pi_1 \mathbf{A} \Pi_2^T - \widehat{\mathbf{L}} \mathbf{M} \widehat{\mathbf{U}}\|_F \le \omega \sigma_{k+1} (\mathbf{A}),$$

where $\gamma = O(fk\sqrt{mn})$ is the same as in Theorem 4, and $\omega =$ O(fkmn).

Theorem

If $\sigma_i^2(\mathbf{A}) > 2\|\mathbf{S}\|_2^2$ then

$$\sigma_{j}(\mathbf{A}) \geq \sigma_{j}(\widehat{\mathbf{L}}\mathbf{M}\widehat{\mathbf{U}}) \geq \sigma_{j}(\mathbf{A})\sqrt{1-2\gamma\left(\frac{\sigma_{k+1}(\mathbf{A})}{\sigma_{j}(\mathbf{A})}\right)^{2}}$$

for $\gamma = O\left(mnk^2f^2\right)$ and f is an input parameter controlling a tradeoff of quality vs. speed as before.

Experiments

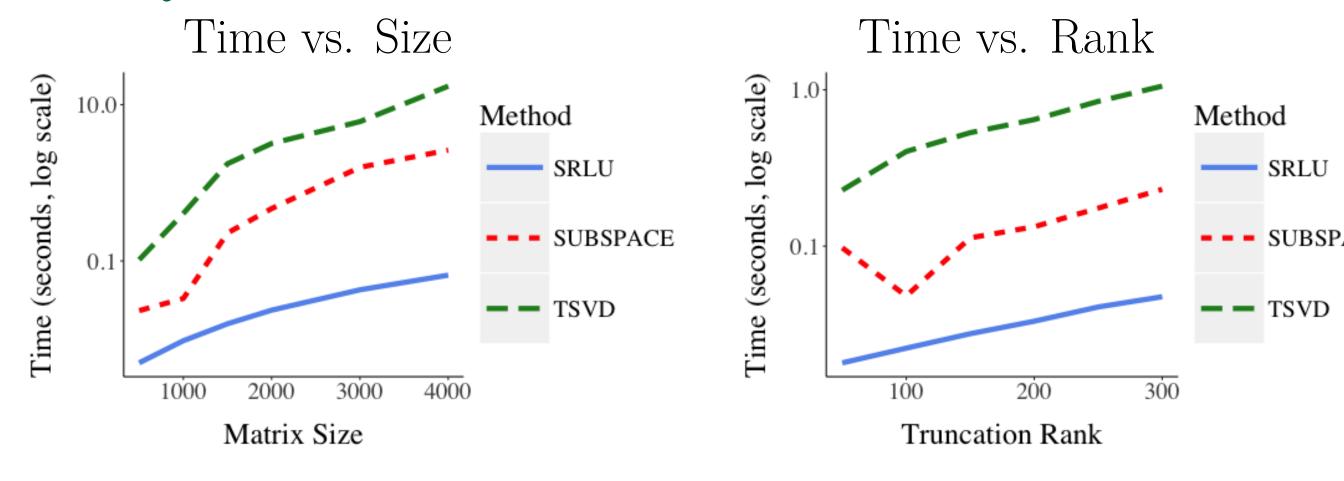
Accuracy: 0.8 Decay Rate 0.95 Decay Rate --- SRLU --- CUR --- CUR **—** TSVD **—** TSVD

Roughly a constant amount of oversampling is needed to match the accuracy of the Truncated SVD

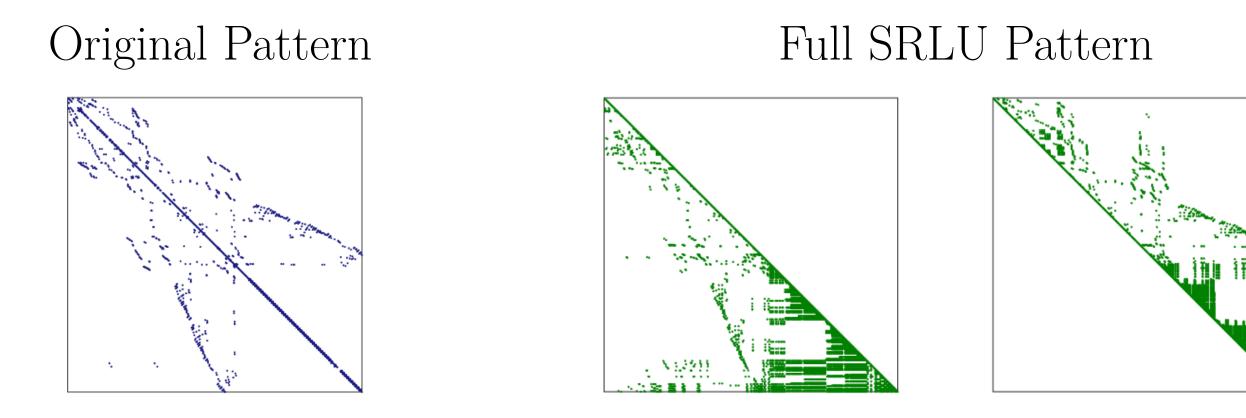
Truncation Rank

Experiments (Continued)

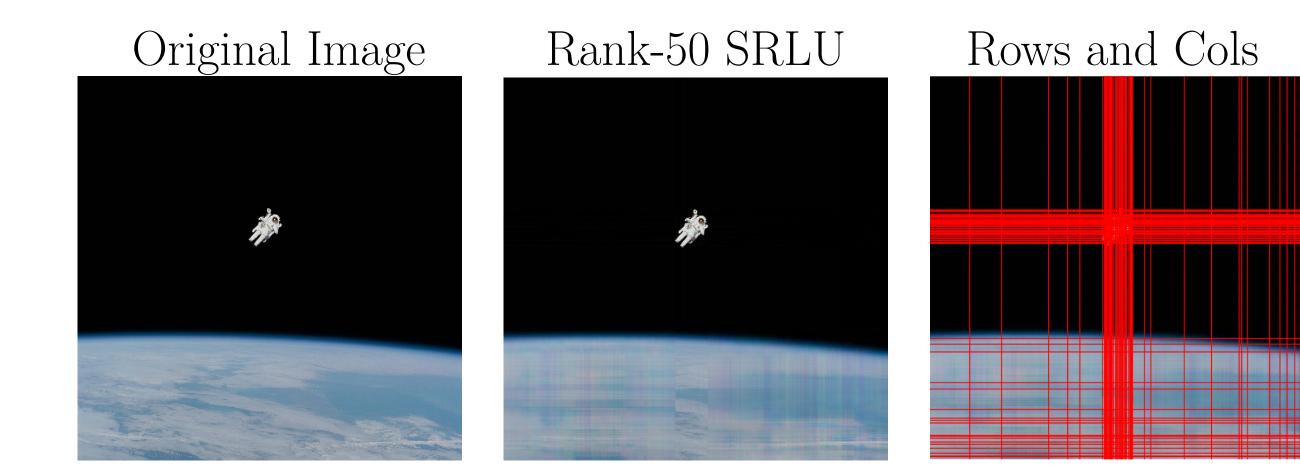
Efficiency:



Sparsity Preservation:



Feature Selection:



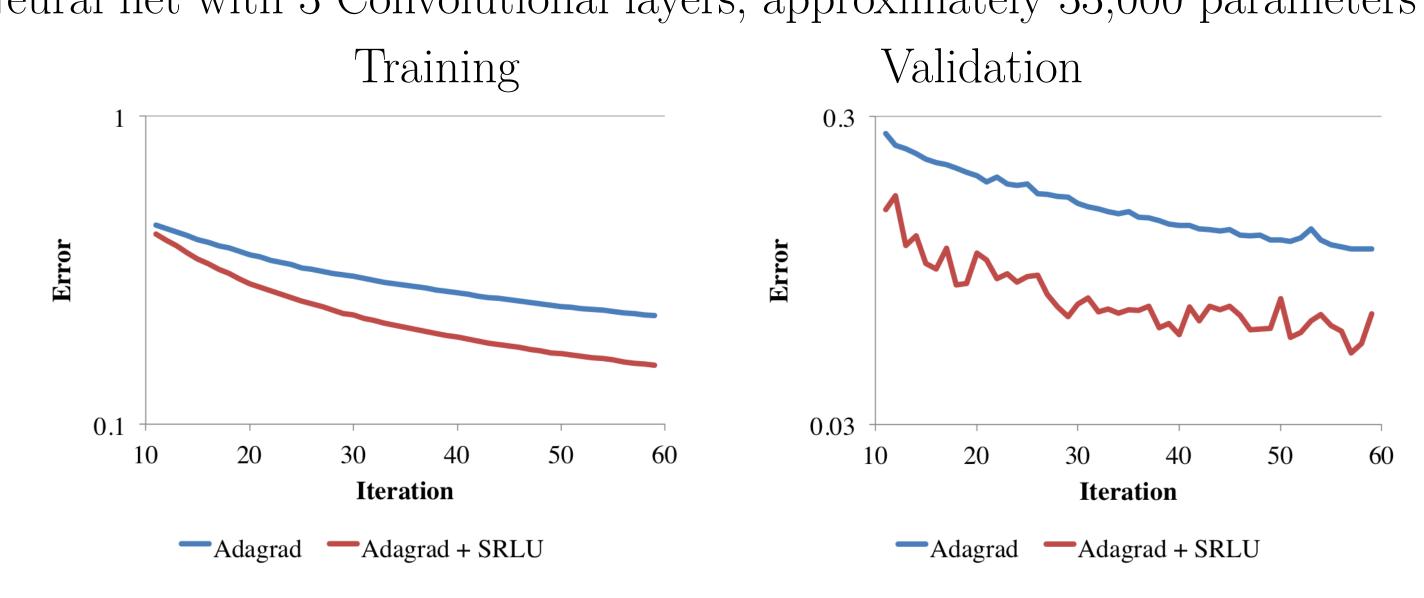
Online Updating:

SRLU finds top words in the Enron email corpus when processed in blocks

Going Forward

Efficient Optimizers: Approximation of Full Adagrad

- Based on recent work by Krummenacher et. al. [3]
- Use SRLU now (column pivots applied to rows too)
- Symmetric implementation
- Neural net with 3 Convolutional layers, approximately 33,000 parameters



References

- 1] Clarkson K. L. and Woodruff, D. P., Low Rank Approximation and Regression in Input Sparsity Time, CoRR, abs/1207.6365, 2012.
- [2] Gu, M. Subspace Iteration Randomization and Singular Value Problems, SIAM J. Scientific Computing, 37(3), 2015.
- [3] Krummenacher, G. and McWilliams, B., Kilcher, Y., Buhmann, J. M., and Meinshausen, N, Scalable Adaptive Stochastic Optimization Using Random Projections, NIPS, pp. 1750-1758, 2016.
- [4] Mahoney, M. W. and Drineas, P., CUR Matrix Decompositions for Improved Data Analysis, PNAS, 106(3), pp. 697-702, 2009.
- [5] Melgaard, C. and Gu, M., Gaussian Elimination with Randomized Complete Pivoting, CoRR, abs/1511.08528, 2015.
- [6] Sarlós, T., Improved Approximation Algorithms for Large Matrices Via Random Projections, FOCS, pp. 143-152, 2006.