

Lecture 18 — 03-25

Instructor: Shashanka Ubaru

Scribe: Gabriel H. Brown

1 Leverage scores and coherence

Recall the definition of row leverage scores of a matrix $\mathbf{A} \in \mathbb{R}^{n \times r}$. If \mathbf{U} is an orthonormal basis for $\text{span}(\mathbf{A})$, then the i th leverage score is given by

$$\ell_i(\mathbf{A}) = \sup_{\mathbf{x}} \frac{(\mathbf{A}_i \mathbf{x})^2}{\|\mathbf{A} \mathbf{x}\|_2^2} = \|\mathbf{U}_{i:}\|_2^2, \quad i \in [n]. \quad (1)$$

One can then sample rows according to probabilities $p_i = \ell_i/r$. It is also possible to approximately compute leverage scores at a reduced complexity.

The coherence of \mathbf{A} , denoted by $\mu(\mathbf{A})$ is the *maximum leverage score*, that is,

$$\mu(\mathbf{A}) = \max_{i \in [n]} \ell_i(\mathbf{A}). \quad (2)$$

Coherence obeys the following inequalities: $\frac{r}{n} \leq \mu(\mathbf{A}) \leq 1$. The first follows from $1 = \sum_{i=1}^n p_i = \frac{1}{r} \sum_{i=1}^n \ell_i \Rightarrow r = \sum_{i=1}^n \ell_i \Rightarrow r \leq n\mu(\mathbf{A})$; the second follows from $\ell_i = \|\mathbf{U}_{i:}\|_2^2 = \|\mathbf{e}_i \mathbf{U}\|_2^2 \leq 1^2 \cdot 1^2$ by submultiplicativity. We say \mathbf{A} is **incoherent** if $\mu(\mathbf{A}) \approx \frac{r}{N}$.

One can use leverage scores to sample linear least squares problems, getting approximate solutions at a reduced cost.

Proposition 1. Given a matrix $\mathbf{A} \in \mathbb{R}^{n \times r}$ and a fixed vector $\mathbf{b} \in \mathbb{R}^n$, let $\mathbf{x}^* = \min_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{A} \mathbf{x} - \mathbf{b}\|_2$. Let $\mathbf{S} \in \mathbb{R}^{m \times n}$ be a sampling matrix with probabilities $p_i = \ell_i/r$, and $\mathbf{S}_{i*} = \mathbf{e}_j / \sqrt{mp_j}$ with $\mathbb{P}(j = i) = p_i$. If $m = O(r \log(r/\delta)/\varepsilon)$ and $\tilde{\mathbf{x}} = \min_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{S}(\mathbf{A} \mathbf{x} - \mathbf{b})\|_2$, then, with high probability,

$$\|\mathbf{A} \tilde{\mathbf{x}} - \mathbf{b}\|_2 \leq (1 + \varepsilon) \|\mathbf{A} \mathbf{x}^* - \mathbf{b}\|_2.$$

2 Leverage score sampling for CP-ALS

Our goal is to accelerate CP-ALS by using leverage score sampling on the least squares subproblems that arise for the approximate CP factor matrices of tensor \mathcal{X} (here order 3)

$$\min_{\mathbf{A}_1} \|(\mathbf{A}_3 \odot \mathbf{A}_2) \mathbf{A}_1^T - \mathbf{X}_{(1)}^T\|_F^2. \quad (3)$$

However, even approximately computing the leverage scores for $(\mathbf{A}_3 \odot \mathbf{A}_2)$ can be prohibitively expensive. But we can estimate/bound the leverage scores for this Khatri-Rao structured matrix in terms of the leverage scores of the matrices \mathbf{A}_2 and \mathbf{A}_3 .

Lemma 1 (Cheng, et al.: Theorem 3.2). $\mu(\mathbf{A} \odot \mathbf{B}) \leq \mu(\mathbf{A})\mu(\mathbf{B})$

This implies that if two matrices \mathbf{A}, \mathbf{B} are incoherent, then their Khatri-Rao product $\mathbf{A} \odot \mathbf{B}$ is also incoherent.

Motivated by Lemma 1, instead of sampling according to $p_k = \ell_k(\mathbf{A}_3 \odot \mathbf{A}_2)/r$ we will instead use $p_k = \ell_i(\mathbf{A}_3)\ell_j(\mathbf{A}_2)/r^2$, which requires only $\mathcal{O}((n_2 + n_3)r)$ work. Specifically, we'll use the following procedure:

- choose $i \sim p_i = \ell_i(\mathbf{A}_3)/r$
- choose $j \sim p_j = \ell_j(\mathbf{A}_2)/r$
- select row $k = i + (j - 1)n_3$.

The guarantee for this procedure is as follows.

Theorem 1 (Larsen and Kolda: Theorem 6). *Let $\mathbf{A}_i \in \mathbb{R}^{n_i \times r}$, $\mathbf{X}_{(1)} \in \mathbb{R}^{n_2 n_3 \times n_1}$ and consider the linear least squares problem*

$$\arg \min_{\mathbf{A}_1} \|(\mathbf{A}_3 \odot \mathbf{A}_2)\mathbf{A}_1^T - \mathbf{X}_{(1)}^T\|_F^2.$$

with optimal solution \mathbf{A}_1^* . Now let $\tilde{\mathbf{A}}_1$ be the optimal solution to the problem

$$\arg \min_{\mathbf{A}_1} \|(\mathbf{S}(\mathbf{A}_3 \odot \mathbf{A}_2)\mathbf{A}_1^T - \mathbf{S}\mathbf{X}_{(1)}^T)\|_F^2.$$

where $\mathbf{S} \in \mathbb{R}^{s \times n_2 n_3}$ is the leverage score sampling matrix which samples according to the procedure described above.

If $s = r^4 \max\{1700 \log(r/\delta), 1/(\delta\epsilon)\}$, then

$$\Pr \left[\|(\mathbf{A}_3 \odot \mathbf{A}_2)\tilde{\mathbf{A}}_1^T - \mathbf{X}_{(1)}^T\|_F^2 \leq (1 + \epsilon) \|(\mathbf{A}_3 \odot \mathbf{A}_2)(\mathbf{A}_1^*)^T - \mathbf{X}_{(1)}^T\|_F^2 \right] \geq 1 - \delta.$$

Larsen and Kolda also suggest additional practical tips for efficient implementation:

- **hybrid approach:** deterministically include all rows whose leverage scores/probabilities are above some threshold and randomly sample from the remaining rows; using the hybrid strategy, they demonstrate equally good or better decompositions with the same number of total samples
- **unfoldings:** never form $\mathbf{X}_{(i)}^T$ explicitly if \mathcal{X} is sparse, instead precompute linear indices for every nonzero for each mode to directly form a sparse unfolding/right hand side; this requires $3\text{nnz}(\mathcal{X})$ extra memory
- **estimate residual:** calculating the residual is necessary to determine when the approximation is sufficiently converged, but computing the residual can take many times longer than updating all three factor matrices; therefore, as a *practical hack* with no theoretical guarantees the authors suggest estimating the residual based on a random sample of tensor elements (using a stratified sampling to correct for problems introduced by sparsity)

We conclude this discussion by comparing the complexities of the two main kernels in CP-ALS: computation of the residual and the formation and solution of (one) linear least squares problem. Here s_{fit} is the user specified number of elements used to sample and estimate the residual and j corresponds to the mode being updated by the least squares subproblem.

operation	sparse/dense	complexity (big \mathcal{O})	sampled complexity
residual	dense	$rn_1n_2n_3$	rs_{fit}
	sparse	$r\text{nnz}(\mathcal{X})$	rs_{fit}
least squares (mode j)	dense	$rn_1n_2n_3$	$3sr + sr^2 + srn_j$
	sparse	$rn_1n_2n_3$	$3sr + sr^2 + r\text{nnz}(\mathbf{X}_{(j)})$

3 Troubles with CP decomposition

Finally, we summarize some mathematical troubles that plague the CP decomposition, stressing that these are independent of algorithm.

- **ill-posedness:** from de Silva and Lim there is no guarantee that the best rank k approximation exists; for example no rank 3 $2 \times 2 \times 2$ tensor has a best rank 2 approximation, and a random $m \times n \times p$ tensor has no best rank 2 approximation with probability; in general this rules out the possibility of a theorem like Echart-Young-Mirsky for tensors
- **complexity of determining rank:** given a tensor \mathcal{T} , determining its real rank is NP hard (and so are many other tensor problems, see Lim)
- **many local minima:** the best rank k approximation problem is non-convex and non-linear, and so may have many local minima; no known results suggest that the standard algorithms frequently find “good” local minima