

CS 598 EVS: Tensor Computations

Tensor Decomposition

Edgar Solomonik

University of Illinois at Urbana-Champaign

CP Decomposition Rank

- ▶ The *canonical polyadic or CANDECOMP/PARAFAC (CP) decomposition* expresses an order d tensor in terms of d factor matrices
 - ▶ For $\mathcal{T} \in \mathbb{R}^{n_1 \times \dots \times n_d}$, a rank R CP decomposition is defined by matrices $U^{(i)} \in \mathbb{R}^{s_i \times R}$ so that

$$t_{i_1 \dots i_d} = \sum_{r=1}^R \prod_{j=1}^d u_{i_j r}^{(j)}$$

- ▶ The CP decomposition is also often denoted by

$$\mathcal{T} = \llbracket U^{(1)}, \dots, U^{(d)} \rrbracket$$

- ▶ First proposed by Hitchcock in 1927
- ▶ Given a tensor, the smallest R for which it has a CP decomposition is the *tensor rank*, also sometimes referred to as the CP rank or canonical rank
- ▶ Finding the CP rank and associated decomposition enables automatic derivation of bilinear algorithms among other applications reviewed later in this lecture, which mostly follows T. Kolda and B. Bader "Tensor Decompositions and Applications", SIAM Review 2009.

Tensor Rank Properties

- ▶ Tensor rank does not satisfy many of the properties of matrix rank
 - ▶ Rank of a real-valued tensor can be different over the complex field, e.g., for

$$\mathcal{T} = \left[\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \right]$$

which is a perfectly conditioned tensor, the rank over \mathbb{R} is 3 but over \mathbb{C} it is 2,

$$\mathcal{T} = \llbracket \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -i & i \end{bmatrix}, \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ i & -i \end{bmatrix}, \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ i & -i \end{bmatrix} \rrbracket$$

- ▶ The maximal possible rank of tensors of particular dimensions is often unequal to the **typical tensor rank**, which is any rank for which the set of tensors of a given size with that rank has positive volume
- ▶ $2 \times 2 \times 2$ tensors have typical ranks 2 (79%) and 3 (21%) over \mathbb{R} , and typical rank 2 over \mathbb{C}

Typical Rank and Generic Rank

- ▶ When there is only a single typical tensor rank, it is the *generic rank*
 - ▶ For decomposition over \mathbb{C} , tensors have a single generic rank
 - ▶ If we restrict to symmetric tensors of order d and dimension n , the generic rank over \mathbb{C} is

$$R = \left\lceil \binom{n+d-1}{d} / n \right\rceil$$

except when $(d, n) \in \{(3, 5), (4, 3), (4, 4), (4, 5)\}$ in which cases it should be increased by one

- ▶ This rank bound makes sense, as the total amount of information in a single factor matrix is $nR \approx \binom{n+d-1}{d}$, which matches the number of unique entries in the symmetric tensor
- ▶ For maximal rank of an $n_1 \times n_2 \times n_3$ tensor, the maximal rank is bounded (weakly) by $R \leq \min(n_1 n_2, n_1 n_3, n_2 n_3)$, which follows by the same intuition

Uniqueness Sufficient Conditions

- ▶ Unlike the low-rank matrix case, the CP decomposition can be unique
 - ▶ *In the matrix case, given $A = UV^T$, for any invertible M we can obtain a new factorization $A = UM(VM^{-1})^T$*
 - ▶ *In CP decomposition, the indeterminacy is generally limited to permutation of the R rank-1 factors and scaling of their components*
 - ▶ *Modulo permutation and scaling, strong conditions exist on uniqueness of the CP decomposition*
 - ▶ *Define the k -rank of a matrix as the maximum value of k such that any k columns of the matrix are linearly independent*
 - ▶ *For an order 3 tensor with CP decomposition $[[A, B, C]]$ where the factor matrices have k -ranks k_A, k_B , and k_C , a sufficient condition for uniqueness is*

$$k_A + k_B + k_C \geq 2R + 2$$

- ▶ *For order d tensors whose CP decomposition is composed of matrices with has k ranks k_1, \dots, k_d , the sufficient condition is*

$$\sum_{i=1}^d k_i \geq 2R + (d - 1)$$

Uniqueness Necessary Conditions

- ▶ Necessary conditions for uniqueness of the CP decomposition also exist
 - ▶ *A simple necessary condition for uniqueness is*

$$\min_{l \in \{1, \dots, d\}} \left(\prod_{i=1, i \neq l}^d \text{rank}(\mathbf{U}^{(i)}) \right) \geq R$$

- ▶ *This condition stems from the more general restriction that*

$$\min_{l \in \{1, \dots, d\}} \text{rank} \left(\bigodot_{i=1, i \neq l}^d \mathbf{U}^{(i)} \right) \geq R$$

- ▶ *When one of the d Khatri-Rao products is rank deficient, multiple (infinite) choices of the l th factor matrices must bring the residual to zero*

Degeneracy

- ▶ The best rank- k approximation may not exist, a problem known as *degeneracy* of a tensor
 - ▶ Consider a rank 3 tensor

$$\mathcal{T} = \mathbf{a}_1 \otimes \mathbf{b}_1 \otimes \mathbf{c}_2 + \mathbf{a}_1 \otimes \mathbf{b}_2 \otimes \mathbf{c}_1 + \mathbf{a}_2 \otimes \mathbf{b}_1 \otimes \mathbf{c}_1$$

where the factor matrices are defined to be orthogonal

- ▶ The tensor can be approximated arbitrarily closely by

$$\mathcal{W}^{(\alpha)} = \alpha(\mathbf{a}_1 + \frac{1}{\alpha}\mathbf{a}_2) \otimes (\mathbf{b}_1 + \frac{1}{\alpha}\mathbf{b}_2) \otimes (\mathbf{c}_1 + \frac{1}{\alpha}\mathbf{c}_2) - \alpha\mathbf{a}_1 \otimes \mathbf{b}_1 \otimes \mathbf{c}_1$$

in particular

$$\lim_{\alpha \rightarrow \infty} \|\mathcal{W}^{(\alpha)} - \mathcal{T}\| = 0$$

- ▶ Consequently, the best rank-2 approximation does not exist for this tensor, as in the limit $\mathcal{W}^{(\alpha)}$ converges to an order 3 tensor

Border Rank

- ▶ Degeneracy motivates an approximate notion of rank, namely *border rank*
 - ▶ *The border rank of a tensor \mathcal{T} is defined as the smallest R such that, for any $\epsilon > 0$, there exists a rank R tensor \mathcal{W} such that*

$$\|\mathcal{T} - \mathcal{W}\| < \epsilon$$

- ▶ *The border rank is always less than the rank of a tensor, but can also be smaller*
- ▶ *The concept of border rank has been intensively used to find fast bilinear algorithms for matrix multiplication*
- ▶ *The border rank and rank of the $4 \times 4 \times 4$ multiplication tensor are both 7, yielding Strassen's algorithm*
- ▶ *For the $9 \times 9 \times 9$ tensor defining multiplication of 3×3 matrices, determining rank and border rank is an open problem, the rank is between 19 and 23, while the border rank is between 14 and 21*

Approximation by CP Decomposition

- ▶ Approximation via CP decomposition is a nonlinear optimization problem
 - ▶ Given order d tensor \mathcal{T} with all dimensions equal to n , the rank- R CP approximation problem can be written as

$$\min_{\mathbf{U}^{(1)}, \dots, \mathbf{U}^{(d)} \in \mathbb{R}^{n \times R}} \frac{1}{2} \underbrace{\|\mathcal{T} - \llbracket \mathbf{U}^{(1)}, \dots, \mathbf{U}^{(d)} \rrbracket\|_F^2}_{\phi(\mathbf{U}^{(1)}, \dots, \mathbf{U}^{(d)})}$$

- ▶ The gradient of this objective function is

$$\nabla \phi = [d\phi/d\mathbf{U}^{(1)} \dots d\phi/d\mathbf{U}^{(d)}]$$

- ▶ Each component of the gradient has the form

$$\frac{d\phi}{d\mathbf{U}^{(i)}}(\mathbf{U}^{(1)}, \dots, \mathbf{U}^{(d)}) = \mathbf{U}^{(i)} *_{j=1, j \neq i}^d \mathbf{U}^{(j)T} \mathbf{U}^{(j)} - \underbrace{\mathcal{T}_{(i)} \overset{\circ}{\bigcirc}_{j=1, j \neq i} \mathbf{U}^{(j)}}_{\text{MTTKRP}}$$

- ▶ Unless R is very large, computing $\frac{d\phi}{d\mathbf{U}^{(i)}}(\mathbf{U}^{(1)}, \dots, \mathbf{U}^{(d)})$ is not much cheaper than minimizing ϕ w.r.t. $\mathbf{U}^{(i)}$ by solving for $\mathbf{U}^{(i)}$ in $\frac{d\phi}{d\mathbf{U}^{(i)}}(\mathbf{U}^{(1)}, \dots, \mathbf{U}^{(d)}) = \mathbf{0}$

Alternating Least Squares Algorithm

- ▶ The standard approach for finding an approximate or exact CP decomposition of a tensor is the *alternating least squares (ALS) algorithm*
 - ▶ Consider rank R decomposition of a tensor $\mathcal{T} \in \mathbb{R}^{n \times n \times n}$ over \mathbb{R}
 - ▶ A *sweep* takes as input $[[\mathbf{U}^{(k)}, \mathbf{V}^{(k)}, \mathbf{W}^{(k)}]]$ solves 3 quadratic optimization problems to obtain $[[\mathbf{U}^{(k+1)}, \mathbf{V}^{(k+1)}, \mathbf{W}^{(k+1)}]]$, updating each factor matrix in sequence, typically via the normal equations:

$$(\mathbf{V}^{(k)T} \mathbf{V}^{(k)} * \mathbf{W}^{(k)T} \mathbf{W}^{(k)}) \mathbf{U}^{(k+1)} = \mathbf{T}_{(1)}(\mathbf{V}^{(k)} \odot \mathbf{W}^{(k)})$$

$$(\mathbf{U}^{(k+1)T} \mathbf{U}^{(k+1)} * \mathbf{W}^{(k)T} \mathbf{W}^{(k)}) \mathbf{V}^{(k+1)} = \mathbf{T}_{(2)}(\mathbf{U}^{(k+1)} \odot \mathbf{W}^{(k)})$$

$$(\mathbf{U}^{(k+1)T} \mathbf{U}^{(k+1)} * \mathbf{V}^{(k+1)T} \mathbf{V}^{(k+1)}) \mathbf{W}^{(k+1)} = \mathbf{T}_{(3)}(\mathbf{U}^{(k+1)} \odot \mathbf{V}^{(k+1)})$$

- ▶ Residual decreases monotonically, since the subproblems in each subset of nR variables are quadratic
- ▶ Forming the linear equations has cost $O(dnR^2)$ while forming the right-hand-sides requires an MTTKRP with cost $O(n^d R)$

Properties of Alternating Least Squares for CP

- ▶ *CP-ALS achieves linear local convergence to local minima of our objective ϕ*
 - ▶ *this follows from the equivalence of the optimality conditions (vanishing gradient) and the ALS update rule*
 - ▶ *no global convergence guarantees are available, and in practice algorithm convergence can stagnate, typically due to the factor matrix iterates becoming ill-conditioned*
- ▶ *CP-ALS guarantees monotonic decrease in residual*
 - ▶ *the exact minimizer is found for each quadratic subproblem, which cannot be worse than the previous choice*
- ▶ *the equations for each subproblem are formed by a Khatri-Rao product, which makes subproblems amenable to fast approximate methods*

Alternating Least Squares for Tucker Decomposition

- ▶ For Tucker decomposition, an analogous optimization procedure to ALS is referred to as *high-order orthogonal iteration (HOOI)*
 - ▶ *Each component of the derivative of the Tucker approximation objective function with respect to the product of a factor matrix and the core tensor is a TTMc (as opposed to MTTKRP in the CP case)*

$$\psi(\mathbf{Z}, \mathbf{U}, \mathbf{V}, \mathbf{W}) = \frac{1}{2} \left(t_{ijk} - \sum_{pqr} z_{pqr} u_{ip} v_{jq} w_{kr} \right)^2$$

$$\frac{d\psi}{d(\mathbf{Z} \times_1 \mathbf{U})}(\mathbf{Z}, \mathbf{U}, \mathbf{V}, \mathbf{W}) = \sum_{j,k} t_{ijk} v_{jq} w_{kr} - \sum_{pq'r'} z_{pqr} u_{ip} \underbrace{\left(\sum_j v_{jq'} v_{jq} \right)}_{\delta(q,q')} \underbrace{\left(\sum_k w_{kr'} w_{kr} \right)}_{\delta(r,r')}$$

- ▶ *Consequently, we can find the minimizing $\mathbf{Z} \times_1 \mathbf{U}$ by SVD of the mode-1 unfolding of the TTMc $\mathcal{T} \times_2 \mathbf{V}^T \times_3 \mathbf{W}^T$, which is a $s \times R^2$ matrix*
- ▶ *Optimizing for a single factor matrix in this way costs $O(s^d R + s R^d)$*
- ▶ *A sweep of HOOI requires forming N such TTMcs and computing their SVDs*

Dimension Trees for ALS

- ▶ The cost of ALS can be reduced by amortizing computation common terms
 - ▶ *The cost of ALS is typically dominated by MTTKRPs, d of which are computed for each sweep, for $d = 3$,*

$$\mathbf{T}_{(1)}(\mathbf{V}^{(k)} \odot \mathbf{W}^{(k)}), \mathbf{T}_{(2)}(\mathbf{U}^{(k+1)} \odot \mathbf{W}^{(k)}), \mathbf{T}_{(3)}(\mathbf{U}^{(k+1)} \odot \mathbf{V}^{(k+1)})$$

- ▶ *Note that given $\mathcal{Z} = \mathcal{T} \times_3 \mathbf{W}^{(k)T}$, we can compute the first two MTTKRPs with $O(s^2 R)$ cost, since*

$$\sum_{j,l} t_{ijl} v_{jr}^{(k)} w_{lr}^{(k)} = \sum_j z_{ijr} v_{jr}^{(k)} \quad \text{and} \quad \sum_{j,l} t_{ijl} u_{ir}^{(k+1)} w_{lr}^{(k)} = \sum_i z_{ijr} u_{ir}^{(k+1)}$$

- ▶ *In general, we can reuse a single TTM to compute the next $d - 1$ sets of right-hand-sides (MTTKRPs) in ALS (in this sweep or the next sweep)*
- ▶ *The amortized cost of each ALS sweep (assuming Strassen-like matrix-multiplication algorithms are not used) is then given by $\frac{2d}{d-1} s^d R + O(ds^{d-1} R) + O(dR^3)$ where the final term comes from Cholesky factorization of the matrices $\mathbf{G}^{(i)} = \ast_{j=1, j \neq i}^d \mathbf{U}^{(j)T} \mathbf{U}^{(j)}$*

Fast Residual Norm Calculation

- ▶ Calculating the norm of the residual has cost $2ds^dR$, but can be done more cheaply within ALS

- ▶ *We can expand the square of the residual norm as follows*

$$\|\mathcal{T} - \llbracket U, V, W \rrbracket\|_F^2 = \|\mathcal{T}\|_F^2 - \sum_{ijk} t_{ijk} u_{ir} v_{jr} w_{kr} + \sum_{rs} \left(\sum_i u_{ir} u_{is} \right) \left(\sum_j v_{jr} v_{js} \right) \left(\sum_k w_{kr} w_{ks} \right)$$

- ▶ *The first term does not change across ALS sweeps, so it suffices to compute it once*
 - ▶ *The second term can be obtained cheaply from any MTTKRP with $d - 1$ of the factor matrices*
 - ▶ *The third term can be computed from Gram matrices that are needed for ALS anyway*

Pairwise Perturbation Algorithm

- ▶ A route to further reducing the cost of ALS is to perform it approximately via *pairwise perturbation*
 - ▶ *ALS convergence can be slow with factor matrices changing little at every iteration*
 - ▶ *In such cases, can leverage perturbative approximation to compute each MTTKRP approximately using old intermediates*
 - ▶ *Given order 3 tensor \mathcal{T} can compute intermediates*

$$\mathcal{Z}^{(1)} = \mathcal{T} \times_1 \mathbf{U}^{(k)T}, \mathcal{Z}^{(2)} = \mathcal{T} \times_2 \mathbf{V}^{(k)T}, \mathcal{Z}^{(3)} = \mathcal{T} \times_3 \mathbf{W}^{(k)T}$$

and approximate for subsequent iteration $k' > k$ so long as

$$\|\mathbf{U}^{(k')} - \mathbf{U}^{(k)}\|_F, \|\mathbf{V}^{(k')} - \mathbf{V}^{(k)}\|_F, \|\mathbf{W}^{(k')} - \mathbf{W}^{(k)}\|_F < \epsilon,$$

$$\sum_{j,l} t_{ijl} v_{jr}^{(k')} w_{lr}^{(k')} = \sum_j z_{ijr}^{(3)} v_{jr}^{(k')} + \sum_{j,l} t_{ijl} v_{jr}^{(k')} (w_{lr}^{(k')} - w_{lr}^{(k)}) \approx \sum_j z_{ijr}^{(3)} v_{jr}^{(k')}$$

- ▶ *For general order d tensors, each approximated sweep for $k' > k$ has cost $d(d-1)s^2R + O(dR^3)$*
- ▶ *For Tucker, dimension trees and pairwise perturbation work similarly*

Pairwise Perturbation Second Order Correction

- ▶ When approximating a tensor using CP, the partially converged CP factors can sometimes be used in place of the tensor to accelerate cost
 - ▶ *An example of this idea is the second order correction to pairwise perturbation*
 - ▶ *The truncated term in the pairwise perturbation algorithm can be approximated by approximating the tensor \mathcal{T} with its current approximate CP decomposition*

$$\sum_{j,l} t_{ijl} v_{jr}^{(k')} (w_{lr}^{(k')} - w_{lr}^{(k)}) \approx \sum_{j,l,q} u_{iq}^{(k')} v_{jq}^{(k')} w_{lq}^{(k')} v_{jr}^{(k')} (w_{lr}^{(k')} - w_{lr}^{(k)})$$

- ▶ *These contractions can be evaluated with cost $O(sR^2)$*

$$\sum_{j,l,q} u_{iq}^{(k')} v_{jq}^{(k')} w_{lq}^{(k')} v_{jr}^{(k')} (w_{lr}^{(k')} - w_{lr}^{(k)}) = \sum_q u_{iq}^{(k')} \left(\sum_j v_{jq}^{(k')} v_{jr}^{(k')} \right) \left(\sum_l w_{lq}^{(k')} (w_{lr}^{(k')} - w_{lr}^{(k)}) \right)$$

Gauss-Newton Algorithm

- ▶ ALS generally achieves linear convergence, while Newton-based methods can converge quadratically
 - ▶ *Derive these by casting CP as a nonlinear least squares problem,*

$$\phi(\mathbf{x}) = \frac{1}{2} \underbrace{\|\mathbf{y} - \mathbf{f}(\mathbf{x})\|^2}_{\mathbf{r}(\mathbf{x})}$$

- ▶ *Newton's method computes $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{H}_\phi(\mathbf{x})^{-1} \nabla \phi(\mathbf{x})$*
- ▶ *For nonlinear least squares problems, the gradient and Hessian are*

$$\nabla \phi(\mathbf{x}) = \mathbf{J}_r^T(\mathbf{x}) \mathbf{r}(\mathbf{x}),$$

$$\mathbf{H}_\phi(\mathbf{x}) = \mathbf{J}_r^T(\mathbf{x}) \mathbf{J}_r(\mathbf{x}) + \sum_i r_i(\mathbf{x}) \mathbf{H}_{r_i}(\mathbf{x})$$

- ▶ *The Gauss-Newton method approximates $\mathbf{H}_\phi(\mathbf{x}) \approx \mathbf{J}_r^T(\mathbf{x}) \mathbf{J}_r(\mathbf{x})$, so*

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{s}^{(k)}, \quad \mathbf{s}^{(k)} = (\mathbf{J}_r^T(\mathbf{x}^{(k)}) \mathbf{J}_r(\mathbf{x}^{(k)}))^{-1} \mathbf{J}_r^T(\mathbf{x}^{(k)}) \mathbf{r}(\mathbf{x}^{(k)}),$$

$$\mathbf{J}_r(\mathbf{x}^{(k)}) \mathbf{s}^{(k)} \cong \mathbf{r}(\mathbf{x}^{(k)})$$

Gauss-Newton for CP Decomposition

- ▶ CP decomposition for order $d = 3$ tensors ($d > 3$ is similar) minimizes

$$\phi(\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)}) = \frac{1}{2} \sum_{ijk} \left(t_{ijk} - \sum_{r=1}^R u_{ir}^{(1)} u_{jr}^{(2)} u_{kr}^{(3)} \right)^2$$

- ▶ *The Gauss-Newton approximate Hessian is $dnR \times dnR$,*

$$\mathbf{H} = \begin{bmatrix} \mathbf{H}^{(1,1)} & \dots & \mathbf{H}^{(1,d)} \\ \vdots & \ddots & \vdots \\ \mathbf{H}^{(d,1)} & \dots & \mathbf{H}^{(d,d)} \end{bmatrix}, \text{ where } \mathbf{H}^{(q,q)} = \mathbf{G}^{(n,n)} \otimes \mathbf{I}$$

$$\text{while for } q \neq p, \quad h_{krlz}^{(q,p)} = u_{kz}^{(q)} u_{lr}^{(p)} g_{rz}^{(q,p)},$$

$$\text{where in both cases } g_{rz}^{(n,p)} = \prod_{m=1, m \neq q, p}^d \left(\sum_i u_{ir}^{(m)} u_{iz}^{(m)} \right)$$

Gauss-Newton for CP Decomposition

- ▶ A step of Gauss-Newton requires solving a linear system with \mathbf{H}
 - ▶ Cholesky of \mathbf{H} requires $O(d^2 n^2 R^2)$ memory and cost $O(d^3 n^3 R^3)$
 - ▶ Matrix-vector product with \mathbf{H} can be computed with cost $O(d^2 n R^2)$
 - ▶ Can use CG method with implicit matrix-vector product¹
 - ▶ Can formulate product $\mathbf{u} = \mathbf{H}\mathbf{v}$ using tensor contractions each with cost $O(nR^2)$

```
u = []
for q in range(d):
    u.append(zeros((n,R)))
    for p in range(d):
        if q == p:
            u[q] += einsum("rz,kz->kr",G[q,p],v[p])
        else:
            u[q] += einsum("kz,lr,rz,lz->kr", \
                           U[q],U[p],G[q,p],v[p])
```

¹P. Tichavsky, A. H. Phan, and A. Cichocki, 2013

Matrix Sketching

Randomized methods provide accurate approximate solutions to linear least squares problems, which can be applied to accelerate ALS, as well as more basic problems

- ▶ Consider a linear least squares problem with $\mathbf{A} \in \mathbb{R}^{m \times n}$

$$\mathbf{A}\mathbf{x} \cong \mathbf{b}$$

- ▶ We seek to find a sketch matrix $\mathbf{S} \in \mathbb{R}^{s \times m}$, $s < m$ so as to obtain an approximate solution $\hat{\mathbf{x}} \approx \mathbf{x}$ via

$$\mathbf{S}\mathbf{A}\hat{\mathbf{x}} \cong \mathbf{S}\mathbf{b}$$

- ▶ We choose \mathbf{S} from a random distribution so that $\mathbf{S}\mathbf{A}$ is efficient to compute, and with high-probability, the residual is not much larger than optimal,

$$\frac{\|\mathbf{A}\hat{\mathbf{x}} - \mathbf{b}\|_2 - \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2}{\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2} \leq \epsilon$$

Random Projections

Accuracy of sketching techniques is theoretically characterized by statistical analysis

- ▶ *To ascertain accuracy of sketching for matrix/vector products and linear least squares, we want to bound error of $\langle Su, Sv \rangle$ relative to $\langle u, v \rangle$ for some set of pairs u, v*
- ▶ *Generally, its easy to enforce that the result computed in expectation agrees with the exact solution*
- ▶ *To bound error, it then suffices to bound the variance of the random variable, e.g., via Chernoff bounds*

Johnson-Lindenstrauss Lemma

The Johnson-Lindenstrauss lemma is a powerful tool for obtaining error bounds in a projected vector space

$$SA\hat{x} \cong Sb$$

- ▶ For $\mathbf{u}_1, \dots, \mathbf{u}_n$ if $S \in \mathbb{R}^{k \times d}$ is elementwise Gaussian-random, then if $k \geq 9 \log n / (\epsilon^2 - \epsilon^3)$, then with probability $1/2$, all projected distances $\|\mathbf{u}_i - \mathbf{u}_j\|_2$ are preserved to within a relative error of ϵ
- ▶ Since the squared distance $\|\mathbf{u}_i - \mathbf{u}_j\|_2^2 = \|\mathbf{u}_i\|_2^2 + \|\mathbf{u}_j\|_2^2 - 2\langle \mathbf{u}_i, \mathbf{u}_j \rangle$ and the projected vector norms are also preserved accurately w.h.p., it follows that we can also sketch $\langle \mathbf{u}_i, \mathbf{u}_j \rangle$ accurately
- ▶ For a more comprehensive review of these results and derivations, see the following notes by Michael Mahoney <https://arxiv.org/pdf/1608.04481.pdf>

Matrix Sketching

The best choice of sketch matrix depends on the desired accuracy and the structure of A

- ▶ *the matrix A may be sparse or structured, e.g., for ALS A is a Khatri-Rao product, while for HOOI, A is a Kronecker product of orthogonal matrices*
- ▶ *for dense matrices, random Gaussian sketches or sketches that can be applied with FFT or FFT-like algorithms are fastest*
- ▶ *for sparse matrices, ideally we wish to sample the rows of A or at least minimally mix them, for which we need a sparse S*
- ▶ *the sketched solve can also be applied as a preconditioner in order to improve accuracy*

Matrix Sketching via Sampling

Uniform sampling of rows is insufficient to obtain good accuracy in general

- ▶ Consider an A with a single row that is orthogonal to the rest, if it is not sampled, all of x spanned by this row will be missed in \hat{x}
- ▶ The component of the right-hand-side in the span of the columns of $A = QR$ is $QQ^T b$
- ▶ After sketching, we reduce the row space, and can capture less of the right-hand side $SQ(SQ)^\dagger Sb$
- ▶ The largest components of $QQ^T b$ will depend on the largest row-norms of Q

Leverage score sampling provides better accuracy guarantees

- ▶ The leverage scores of A are given by $l_i = \|q_i\|_2^2$ where q_i is the i th row of Q
- ▶ We can define S to sample the rows of A with probability l_i/n , since $\sum_i l_i = n$
- ▶ With this choice, a sample size of $s = O(n/\epsilon^2)$ suffices to get within $1 + \epsilon$ of the optimal residual
- ▶ While reducing row dimension from m to $O(n/\epsilon^2)$ is beneficial for small n , computing leverage scores is expensive (requires QR factorization)

Mixing Techniques

To circumvent leverage score sampling, we can mix rows randomly

- ▶ *Choosing elements of S from a Gaussian distribution, results in intermixing of rows of Q*
- ▶ *In effect, this corresponds to uniform sampling with close-to-uniform leverage scores*
- ▶ *The resulting sample complexity is as good as leverage-score-based sampling*

Instead of choosing elements of S randomly, pseudo-random distributions allow S to be applied more rapidly

- ▶ *Choose $S = PFD$ where D is a random diagonal sign matrix, F is the DFT or Hadamard transform, and P performs uniform-random sampling*
- ▶ *Application of orthogonal matrix FD results in a random intermixing of rows, does not modify solution x*
- ▶ *P performs sampling of a matrix with close-to-uniform leverage scores*
- ▶ *Same sample size complexity as Gaussian mixing and leverage score sampling*

Approximate CP ALS using Random Sampling

- ▶ Another approach to approximating ALS is to sample the least-squares equations²
 - ▶ Recall, that in the least squares view, in the first step of ALS we seek to optimize U in $\mathcal{T} \approx \llbracket U, V, W \rrbracket$, yielding

$$(V \odot W)U^T \cong \mathbf{T}_{(1)}^T$$

- ▶ Instead, we can sample S of the equations and solve

$$S(V \odot W)U^T \cong S\mathbf{T}_{(1)}^T$$

where the rows of $S \in \mathbb{R}^{S \times n^2}$ are sampled uniformly (with replacement) from the identity matrix

- ▶ This approach yields cost $O(S \times R^2 + n^2 SR)$ per iteration, as opposed to the usual $O(n^3 R + R^3)$
- ▶ An variation of this method is to mix the equations, e.g., by picking $S = S'(FD_1 \otimes FD_2)D$ where D_1 and D_2 are random diagonal sign matrices, F is the FFT matrix, and S' is a sampling matrix as before
- ▶ The FFTs can then be applied to the tensor as a preprocessing step

²C. Battaglino, G. Ballard, T. G. Kolda, 2018

Tensor Completion

- ▶ The *tensor completion* problem seeks to build a model (e.g., CP decomposition) for a partially-observed tensor
 - ▶ *Completion differs from decomposition of a sparse tensor with zeros for entries that are unobserved, as the CP decomposition would be fitting the zeros, which we do not want*
 - ▶ *For an order three tensor $\mathcal{T} \in \mathbb{R}^{n \times n \times n}$, given a set of observed entries t_{ijk} for $(i, j, k) \in \Omega$, we seek to minimize*

$$f(\mathbf{U}, \mathbf{V}, \mathbf{W}) = \sum_{(i,j,k) \in \Omega} (t_{ijk} - \sum_r u_{ir} v_{jr} w_{kr})^2 + \lambda^2 (\|\mathbf{U}\|_2^2 + \|\mathbf{V}\|_2^2 + \|\mathbf{W}\|_2^2)$$

- ▶ The problem was partially popularized by the Netflix prize collaborative filtering problem
 - ▶ *This problem involved building a model for predicting user ratings of movies, given the set of movies they have already rated, with each rating corresponding to a tuple (user, movie, date), which can be enumerated in a tensor*
 - ▶ *Aside from this problem, tensor completion is a natural generalization of the matrix completion problem, which is intensively used and studied in machine learning*

CP Tensor Completion Gradient and Hessian

- ▶ The gradient of the tensor completion objective function is sparsified according to the set of observed entries
 - ▶ Lets restrict attention to optimizing for the i th row of the first factor matrix, define Ω_i so that $(j, k) \in \Omega_i$ if $s(i, j, k) \in \Omega$, then

$$\phi(\mathbf{u}_i) = \sum_{(j,k) \in \Omega_i} (t_{ijk} - \langle \mathbf{u}_i, \mathbf{v}_j, \mathbf{w}_k \rangle)^2 + \lambda \|\mathbf{u}_i\|_2^2 \quad \text{where} \quad \langle \mathbf{x}, \mathbf{y}, \mathbf{z} \rangle = \sum_r x_r y_r z_r$$

- ▶ Consider the derivative with respect to the i th row of the first factor matrix

$$\nabla \phi(\mathbf{u}_i) = \frac{d\phi}{d\mathbf{u}_i}(\mathbf{u}_i) = 2 \sum_{(j,k) \in \Omega_i} (\mathbf{v}_j * \mathbf{w}_k) (\langle \mathbf{u}_i, \mathbf{v}_j, \mathbf{w}_k \rangle - t_{ijk}) + 2\lambda \mathbf{u}_i$$

- ▶ ALS for tensor decomposition solves quadratic optimization problem for each row of each factor matrix, in the completion case, Newton's method on these subproblems yields different Hessians

- ▶ The Hessian $\mathbf{H}_i^{(\phi)}$ depends on the set of entries $\Omega_i = \{(j, k) : \exists(i, j, k) \in \Omega\}$,

$$\mathbf{H}_i^{(\phi)} = \frac{d^2\phi}{d\mathbf{u}_i d\mathbf{u}_i}(\mathbf{u}_i) = \sum_{(j,k) \in \Omega_i} (\mathbf{v}_j * \mathbf{w}_k) (\mathbf{v}_j * \mathbf{w}_k)^T + 2\lambda \mathbf{I}$$

Methods for CP Tensor Completion

- ▶ ALS for tensor completion with CP decomposition incurs additional cost
 - ▶ For each $(i, j, k) \in \Omega$, need to accumulate $(\mathbf{v}_j * \mathbf{w}_k)(\mathbf{v}_j * \mathbf{w}_k)^T$ to $\mathbf{H}_i^{(\phi)}$
 - ▶ While the n outer products can be amortized with cost $O(nR^2)$, no easy way to do so for their partial sums, leading to cost $O(|\Omega|R^2)$
- ▶ Alternative methods for tensor completion include coordinate descent and stochastic gradient descent
 - ▶ Stochastic gradient descent (SGD) would compute subgradients for each (i, j, k) which are summands in the sum over Ω_i in $\nabla\phi(\mathbf{u}_i)$
 - ▶ SGD can be implemented efficiently, by computing a sum over a random set of subgradients at a time, via subsampling of Ω
 - ▶ Coordinate descent optimizes an entry of each factor matrix at a time
 - ▶ Variants of coordinate descent select different orderings of entries to optimize, e.g., alternating among columns of factor matrices then factor matrices or vice versa

Coordinate Descent for CP Tensor Completion

- ▶ Coordinate descent avoids the need to solve linear systems of equations
 - ▶ *The coordinate-wise objective function is*

$$\psi(u_{ir}) = \sum_{(j,k) \in \Omega_i} (\rho_{ijk}^{(r)} - u_{ir} v_{jr} w_{kr})^2 + \lambda u_{ir}^2 \text{ where } \rho_{ijk}^{(r)} = t_{ijk} - \langle \mathbf{u}_i, \mathbf{v}_j, \mathbf{w}_k \rangle + u_{ir} v_{jr} w_{kr}$$

above $\rho_{ijk}^{(r)}$ is equal to an entry of the residual tensor with the r th rank-one component of the CP decomposition excluded

- ▶ *Taking its derivative, we obtain*

$$\psi'(u_{ir}) = -2 \sum_{(j,k) \in \Omega_i} v_{jr} w_{kr} (\rho_{ijk}^{(r)} - u_{ir} v_{jr} w_{kr}) + 2\lambda u_{ir}$$

- ▶ *Setting this derivative to zero, we can solve for u_{ir}*

$$u_{ir}^{(new)} = \frac{\sum_{(j,k) \in \Omega_i} v_{jr} w_{kr} \rho_{ijk}^{(r)}}{\lambda + \sum_{(j,k) \in \Omega_i} v_{jr}^2 w_{kr}^2}$$

- ▶ *This can be implemented efficiently by keeping track of a residual tensor and obtaining $\rho_{ijk}^{(r)}$ as a modification thereof when working on the r th column*

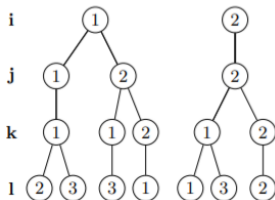
Sparse Tensor Contractions

- ▶ Tensor completion and sparse tensor decomposition require operations on sparse tensors
 - ▶ *In many publicly available sparse tensor datasets, the density is extremely low, e.g., 10^{-7} , i.e., there can be $O(n)$ nonzeros in interesting $n \times n \times n$ tensors*
 - ▶ *For both decomposition and completion, tensor sparsity does not generally imply sparsity of CP or Tucker factors, and these are typically assumed to be dense*
- ▶ Sparse tensor contractions often correspond to products of *hypersparse* matrices, i.e., matrices with mostly zero rows
 - ▶ *Consider TTM with a $n \times n \times n$ tensor \mathcal{T} containing $O(n)$ nonzeros, $\mathbf{T}_{(1)}^T \mathbf{M}$, the matrix $\mathbf{T}_{(1)}^T$ has $O(n)$ nonzeros, but n^2 rows, while $\mathbf{T}_{(1)}^T \mathbf{M}$ has $O(n)$ dense rows and all other $O(n^2)$ rows are zero*
 - ▶ *To reduce sparse tensor contractions to sparse matrix multiplication kernels, need support for hypersparse matrix formats (e.g., compressed sparse-row (CSR) format would require $\Theta(n^2)$ storage for $\mathbf{T}_{(1)}$) and ideally specialized formats for matrices such as $\mathbf{T}_{(1)}^T \mathbf{M}$ (e.g., dense matrix consisting of nonzero rows and vector of row indices)*

Sparse Tensor Formats

- ▶ The overhead of transposition, and non-standard nature of the arising sparse matrix products, motivates sparse data structures for tensors that are suitable for tensor contractions of interest
 - ▶ *Particularly important, especially for tensor decomposition, are MTTKRP (suffices to CP ALS) and TTMc (suffices for HOOI)*
 - ▶ *TTM is also prevalent, but is a less attractive primitive in the sparse case than MTTKRP and TTMc, as these yield dense, low-order outputs, while the output of TTM can be sparse and larger than the starting tensor*
- ▶ The **compressed sparse fiber (CSF)** format provides an effective representation for sparse tensors
 - ▶ *CSF can be visualized as a tree (diagram taken from original CSF paper, by Shaden Smith and George Karpis, IA³, 2015)*

i	j	k	l
1	1	1	2
1	1	1	3
1	2	1	3
1	2	2	1
2	2	1	1
2	2	1	3
2	2	2	2



Operations in Compressed Format

- ▶ CSF permits efficient execution of important sparse tensor kernels
 - ▶ Analogous to CSR format, which enables efficient implementation of the sparse matrix vector product
 - ▶ where `row[i]` stores a list of column indices and nonzeros in the i th row of A

```
for i in range(n):  
    for (a_ij,j) in row[i]:  
        y[i] += a_ij * x[j]
```

- ▶ In CSF format, a multilinear function evaluation $f^{(\mathcal{T})}(\mathbf{x}, \mathbf{y}) = \mathbf{T}_{(1)}(\mathbf{x} \odot \mathbf{y})$ can be implemented as

```
for (i,T_i) in T_CSF:  
    for (j,T_ij) in T_i:  
        for (k,t_ijk) in T_ij:  
            z[i] += t_ijk * x[j] * y[k]
```

MTTKRP in Compressed Format

- ▶ MTTKRP and CSF pose additional implementation opportunities and challenges
 - ▶ MTTKRP $u_{ir} = \sum_{j,k} t_{ijk} v_{jr} w_{kr}$ can be implemented by adding a loop over r to our code for $f^{(\mathcal{T})}$, but would then require $3mr$ operations if m is the number of nonzeros in \mathcal{T} , can reduce to $2mr$ by amortization

```
for (i,T_i) in T_CSF:
    for (j,T_ij) in T_i:
        for r in range(R):
            f_ij = 0
            for (k,t_ijk) in T_ij:
                f_ij += t_ijk * w[k,r]
            u[i,r] = f_ij * v[j,r]
```

- ▶ However, this amortization is harder (requires storage or iteration overheads) if the index i is a leaf node in the CSF tree
- ▶ Similar challenges in achieving good reuse and obtaining good arithmetic intensity arise in implementation of other kernels, such as TTMc

All-at-once Contraction

- ▶ When working with sparse tensors, it is often more efficient to contract multiple operands in an all-at-once fashion
 - ▶ *Given chain of matrix products $ABC \dots$, dimension of overall iteration space scales with number of matrices, but by contracting pairwise, obtain cubic cost in matrix dimension with linear dependence on number of matrices*
 - ▶ *A case when such pairwise contraction is not a good idea, is the sampled dense-dense matrix-multiplication (SDDMM),*

$$c_{ij} = \sum_r^R a_{ij} u_{ir} v_{jr} \Leftrightarrow \mathbf{C} = \mathbf{A} * (\mathbf{UV}^T)$$

where \mathbf{A} is sparse with m nonzeros, while \mathbf{U} and \mathbf{V} are dense

- ▶ *Since the sparsity pattern of \mathbf{C} is the same as of \mathbf{A} , suffices to iterate over nonzeros of \mathbf{A} and multiply each by inner product of a row of \mathbf{U} and a row of \mathbf{V} , with cost $O(mR)$*
- ▶ *Pairwise contraction is inefficient, contracting first \mathbf{A} with \mathbf{U} would yield a third-order intermediate, while contracting \mathbf{U} with \mathbf{V}^T would have cost $O(n^2R)$*
- ▶ *Generalizing SDDMM to higher order gives the tensor-times tensor-product (TTTP), an application of which is computing the residual in tensor completion*

Constrained Tensor Decomposition

- ▶ Many applications of tensor decomposition in data science, feature additional structure, which can be enforced by constraints
 - ▶ *A basic and common constraint is nonnegativity of factor matrices, which often makes sense when working with a tensor that is nonnegative (e.g., count data)*
 - ▶ *Most of the methods we've discussed can be generalized to handle nonnegativity, e.g., one could perform ALS by solving each subproblem subject to nonnegativity constraints*
 - ▶ *Another common constraint is factor matrix orthogonality, which can be incorporated similarly into subproblems*
 - ▶ *For symmetric tensors, repeating factors are often desired, which can be formulated via constraints or by using an appropriate method (two good alternatives are ALS with subiterations to converge updates to repeated factors, or Gauss-Newton, which automatically preserves repeating factors when working with a symmetric tensor)*

Nonnegative Tensor Factorization

- ▶ *Nonnegative tensor factorization (NTF)*, such as CP decomposition with $\mathcal{T} \geq 0$ and $U, V, W \geq 0$ are widespread and a few classes of algorithms have been developed
 - ▶ *Optimization for one of U, V , or W (while the other two are fixed) is a convex optimization problem*
 - ▶ *Many methods based on alternating optimization/updates in the style of ALS*
 - ▶ *A basic approach is to 'clip' result of ALS step so that each factor matrix is nonnegative after update*
 - ▶ *Block coordinate descent (BCD) methods* update on or more columns of U, V , or W based on a coordinate-descent-like update rule
 - ▶ *Proximal gradient methods* are multicolumn BCD methods, which approximately solve each subproblem by minimizing a constrained objective derived based on a proximally projected gradient
 - ▶ *All-at-once methods that update all factor matrices, such as Gauss-Newton with an augmented Lagrangian objective function (sequential quadratic programming)*

Nonnegative Matrix Factorization

- ▶ NTF algorithms with alternating updates have a close correspondence with alternating update algorithms for *Nonnegative matrix factorization (NMF)*

- ▶ *The rank- r NMF problem is to find, given matrix $\mathbf{A} \in \mathbb{R}_+^{n \times n}$, the minimizer $\mathbf{U}, \mathbf{V} \in \mathbb{R}_+^{n \times r}$ to*

$$f^{(\mathbf{A})}(\mathbf{U}, \mathbf{V}) = \|\mathbf{A} - \mathbf{U}\mathbf{V}^T\|_F$$

- ▶ *Solutions to NMF are not easy to compute directly as for the unconstrained low-rank matrix factorization case*
- ▶ *Methods often minimize $\phi_{\mathbf{V}}^{(\mathbf{A})}(\mathbf{U}) = f^{(\mathbf{A})}(\mathbf{U}, \mathbf{V})$ and $\psi_{\mathbf{U}}^{(\mathbf{A})}(\mathbf{V}) = f^{(\mathbf{A})}(\mathbf{U}, \mathbf{V})$ in an alternating fashion via block coordinate descent*
- ▶ *Alternating optimization problems for NTF are essentially the same as in NMF, for*

$$g^{(\mathcal{T})}(\mathbf{U}, \mathbf{V}, \mathbf{W}) = \|\mathcal{T} - \llbracket \mathbf{U}, \mathbf{V}, \mathbf{W} \rrbracket\|_F$$

minimizing $\phi_{\mathbf{V}, \mathbf{W}}^{(\mathcal{T})}(\mathbf{U}) = g^{(\mathcal{T})}(\mathbf{U}, \mathbf{V}, \mathbf{W})$ is the same as the NMF subproblem

$$\phi_{\mathbf{V} \odot \mathbf{W}}^{(\mathbf{T}_{(1)})}(\mathbf{U}) = f^{(\mathbf{T}_{(1)})}(\mathbf{U}, \mathbf{V} \odot \mathbf{W})$$

Coordinate Descent for NMF and NTF

- ▶ Coordinate descent gives optimal closed-form updates for variables in NMF and NTF
 - ▶ We can write an optimization subproblem for a single column \mathbf{u}_i as minimizing

$$\phi_i^{(\mathbf{A})}(\mathbf{u}_i) = \|\mathbf{A} - \sum_{r=1}^R \mathbf{u}_r \mathbf{v}_r^T\|_2 \quad \text{s.t.} \quad \mathbf{u}_i \geq 0$$

$$\mathbf{u}_i^{\text{new}} = \left| \mathbf{u}_i + \frac{\mathbf{A} \mathbf{v}_i - \mathbf{U} \mathbf{V}^T \mathbf{v}_i}{\mathbf{v}_i^T \mathbf{v}_i} \right|_+$$

where $\mathbf{y} = |\mathbf{x}|_+$ gives $y_i = x_i$ if $x_i > 0$ and $y_i = 0$ otherwise

- ▶ Given $\boldsymbol{\rho}^{(i)} = \mathbf{A} - \mathbf{U} \mathbf{V}^T + \mathbf{u}_i \mathbf{v}_i^T = \mathbf{A} - \sum_{j \neq i} \mathbf{u}_j \mathbf{v}_j^T$, if columns of \mathbf{V} are normalized, we can alternatively write the update as

$$\mathbf{u}_i^{\text{new}} = \left| \frac{\boldsymbol{\rho}^{(i)} \mathbf{v}_i}{\mathbf{v}_i^T \mathbf{v}_i} \right|_+$$

Generalized Tensor Decomposition

- ▶ Aside from addition of constraints, the objective function may be modified by using different elementwise loss functions
 - ▶ *The standard loss function is $(x - m)^2$ where x is an element of the tensor and m is its approximation via CP*
 - ▶ *For count data, the Poisson loss function $m - x \log(m)$ may be more appropriate, and typically comes along with nonnegativity constraints*
 - ▶ *Other distributions and loss functions of interest include Gamma, Rayleigh, Bernoulli, and NegBinom (see D. Hong, T. Kolda, J. Duersch SIAM Review 2020)*
- ▶ Some loss function admit ALS-like algorithms, while others may require gradient-based optimization
 - ▶ *Can compute (sub-)gradients given any loss function, by differentiating as necessary*
 - ▶ *For Poisson, like for the standard loss function, ALS subproblems may be solved explicitly, allowing more robust convergence*