

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

Tensor Rank Properties

- ▶ Tensor rank does not satisfy many of the properties of matrix rank

Typical Rank and Generic Rank

- ▶ When there is only a single typical tensor rank, it is the *generic rank*

Uniqueness Sufficient Conditions

- ▶ Unlike the low-rank matrix case, the CP decomposition can be unique

Uniqueness Necessary Conditions

- ▶ Necessary conditions for uniqueness of the CP decomposition also exist

Degeneracy

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

Border Rank

- ▶ Degeneracy motivates an approximate notion of rank, namely *border rank*

Approximation by CP Decomposition

- ▶ Approximation via CP decomposition is a nonlinear optimization problem

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*

Properties of Alternating Least Squares for CP

Alternating Least Squares for Tucker Decomposition

- ▶ For Tucker decomposition, an analogous optimization procedure to ALS is referred to as *high-order orthogonal iteration (HOOI)*

Dimension Trees for ALS

- ▶ The cost of ALS can be reduced by amortizing computation common terms

Fast Residual Norm Calculation

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

Pairwise Perturbation Algorithm

- ▶ A route to further reducing the cost of ALS is to perform it approximately via *pairwise perturbation*

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

Gauss-Newton Algorithm

- ▶ ALS generally achieves linear convergence, while Newton-based methods can converge quadratically

Gauss-Newton for CP Decomposition

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

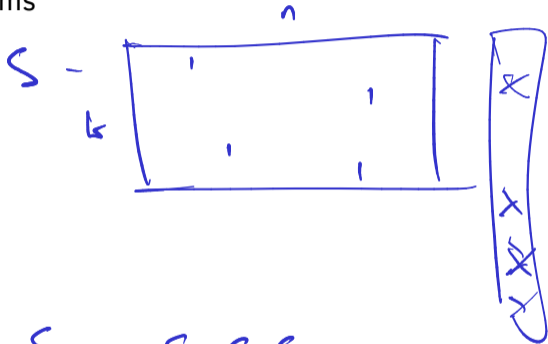
Gauss-Newton for CP Decomposition

- ▶ A step of Gauss-Newton requires solving a linear system with H

```
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])
```

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



$$S = S_1 \otimes S_2$$

$$S_u \quad \text{if} \quad u = u_1 \otimes u_2$$

$$S_u = S_{1,u_1} \otimes S_{2,u_2}$$

Random Projections

Accuracy of sketching techniques is theoretically characterized by statistical analysis

$$AB \approx AS^T SB$$

$$\langle u, v \rangle \approx \langle S_u, S_v \rangle$$

1st ensure $E[\langle S_u, S_v \rangle] = \langle u, v \rangle$, easy

challenge: bound variance of $\langle S_u, S_v \rangle$

↳ error bound
e.g., Chernoff bounds

$$\mathbb{P}\left[\left| \sum_{i=1}^n \pm x_i \right| \geq t\right]$$

Random Projections

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

$$\underbrace{u_1, \dots, u_n}_{S_A \leftarrow S_B} \in \mathbb{R}^d, \quad k < d, \quad S \in \mathbb{R}^{k \times d}, \quad S_{ij} \sim \text{Gaussian random}$$

$$\|Su_i - Su_j\|_2 = (1 \pm \epsilon) \|u_i - u_j\|_2 \quad \forall i, j \quad \text{with prob. } \frac{1}{2}$$

$$\text{So long as } k \geq \frac{16 \log n}{\epsilon^2 - \epsilon^3}$$

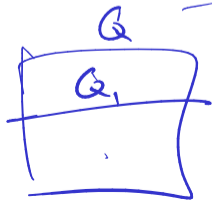
$$\|Su_i\|_2^2 \approx \|u_i\|_2^2$$

$$\|u_i - u_j\|_2^2 = \|u_i\|_2^2 + \|u_j\|_2^2 - 2 \langle u_i, u_j \rangle$$

Matrix Sketching

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

ideally S simply samples, so it's cheap to apply



random vector

$$\langle Su, Sv \rangle$$

\uparrow

$$|\bar{S}| = k$$

$$\sum_i u_i v_i \approx \frac{n}{k} \sum_{i \in \bar{S}} u_i v_i$$

$$Q \in \mathbb{R}^{n \times n}$$

$$\langle Qu, Qv \rangle = u^T Q^T Q v = \langle u, v \rangle$$

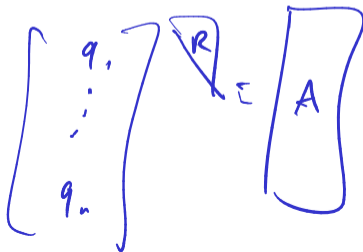
$$\rightarrow \sum_{i=1}^n u_i^{(q)} v_i^{(q)} \approx \frac{n}{k} \sum_{i=1}^k u_i^{(q)} v_i^{(q)}$$

Matrix Sketching via Sampling

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

Leverage score sampling provides better accuracy guarantees

$$L_i(A) = \|q_i\|_2^2$$



$$Ax \approx b$$

Mixing Techniques

To circumvent leverage score sampling, we can mix rows randomly. Instead of

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

Approximate CP ALS using Random Sampling

- ▶ Another approach to approximating ALS is to sample the least-squares equations¹

$$A \approx T \cdot E$$

$$(B \circ C) \cdot A^T \approx T^T \cdot (1)$$

$$S = S_1 \otimes S_2$$

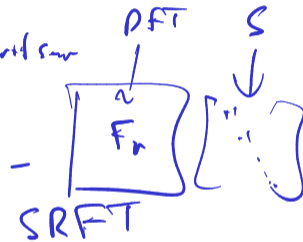
$$S(B \circ C) =$$

$$S = P(\bar{S}_1 \otimes S_2)$$

$$\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \text{with row } n$$

$$\bar{S}_1 =$$

$$\begin{bmatrix} S_1 & S_2 \end{bmatrix}$$



$$S^T \cdot (1) =$$

¹C. Battaglino, G. Ballard, T. G. Kolda, 2018

$$e_{ij}(B \otimes C) = e_i(B) e_j(C)$$

$$M = B \otimes C$$

$$m_{ij} = b_{ik} c_{kj}$$

$$e_{ij}(B \otimes C) = e_i(B) e_j(C)$$

$$Q^{(B)} R^{(B)} \otimes Q^{(C)} R^{(C)} =$$

$$Q^{(B)} \otimes Q^{(C)} (R^{(B)} \otimes R^{(C)})$$

Wort

CP Tensor Completion Gradient and Hessian

- ▶ The gradient of the tensor completion objective function is sparsified according to the set of observed entries

- ▶ 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

Coordinate Descent for CP Tensor Completion

- ▶ Coordinate descent avoids the need to solve linear systems of equations

Sparse Matrices

Solving $Ax = b$ or $A \approx W^T$

when A is sparse

- direct methods (LU, SVD, QR), because of fill (unless A is e.g. banded)
- iterative methods based on SpMV

y = Ax = $f_A(x)$
has cost $O(n + nnz(A))$

CSR - compressed sparse row

vals | array with nonzero values (nnz)
cols | array with column indices (nnz)
offs | array with row offsets (n)

alt.

COO

coord. val

vals (nnz)

cols (nnz)

rows (nnz)

SpMV(x)

for $i=1$ to n

[for $j=offs[i]$ to $offs[i+1]$

$y_i = vals[j] \cdot x[cols[j]]$]

Sparse Tensor Contractions

- ▶ Tensor completion and sparse tensor decomposition require operations on sparse tensors

CP-ALS touches T only for MTTKRP

MTTKRP \rightarrow TTM + low. order

- ▶ Sparse tensor contractions often correspond to products of *hypersparse* matrices, i.e., matrices with mostly zero rows

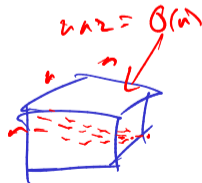
$T \in \mathbb{R}^{n_1 \times n_2 \times n_3}$

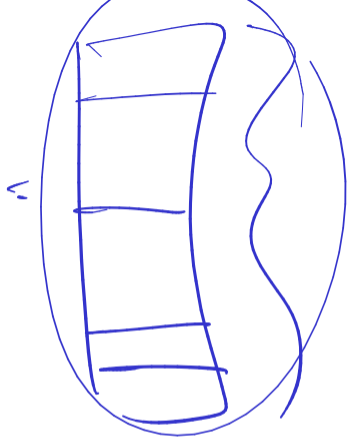
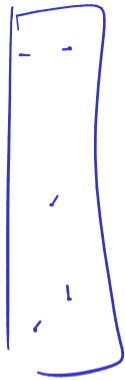


T may have only $\Theta(n)$ nonzeros

$T^{(i)}$ is hypersparse (most rows are 0)

• CSR is not an efficient format (size $\underline{n^2 + n \cdot z(i)}$)





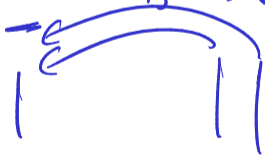
$$z = f^{(n)}(x, y)$$

$$z_i = \sum_{j,k} t_{ijk} x_j y_k$$

for all numbers (i, j, k, t_{ijk})

for $r \in \mathbb{R}$

$$z_{ir} = t_{ijk} x_j y_k$$



all-at-once contraction

since

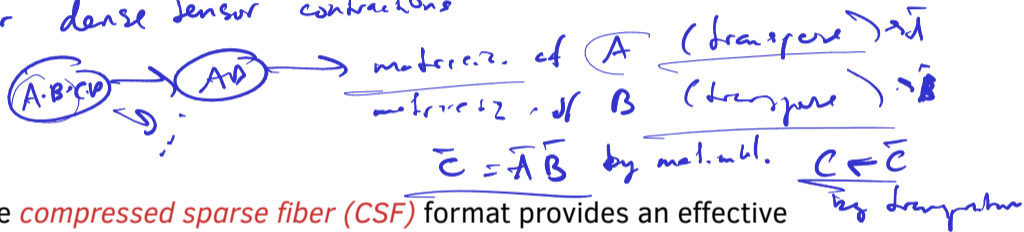
$$\begin{matrix} T x_1 \otimes x_2 y \\ \text{ad} \otimes \text{ad} \\ \hline (T x_1 \otimes x_2) y \end{matrix} \text{ all-at-once}$$

} sufficient to check
our numbers on
any form

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

for dense tensor contractions



- ▶ The compressed sparse fiber (CSF) format provides an effective representation for sparse tensors

assumes index order (e.g. CSR via CSC)
 and constructs tree (of height equal to tensor order)
 where each nonzero is a leaf

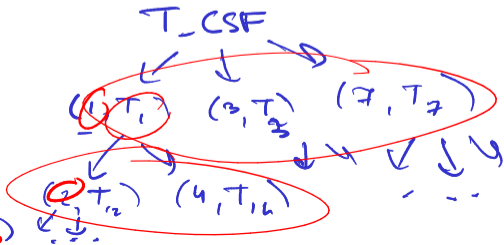
Operations in Compressed Format

- ▶ CSF permits efficient execution of important sparse tensor kernels
 - ▶ Analogous to CSR format, which enables efficient implementation of the sparse matrix vector product
 - ▶ where $\text{row}[i]$ 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 $f^{(\mathcal{T})}(x, y) = \mathbf{T}_{(1)}(x \odot y)$ 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 $u_{ir} = \sum_{j,k} t_{ijk} v_j \cdot w_{kr}$ 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]
```

$$\sum_{j,k} t_{ijk} w_{kr} v_j = \sum_j \left(\sum_k t_{ijk} w_{kr} \right) v_j$$

- ▶ 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



$$n^3 R + n^2 R^2$$

$$\underline{n^3 R} \quad \uparrow \text{refactoring} \quad \text{sit space}$$

loop 3

All-at-once Contraction

- ▶ When working with sparse tensors, it is often more efficient to contract multiple operands in an all-at-once fashion

Constrained Tensor Decomposition

- ▶ Many applications of tensor decomposition in data science, feature additional structure, which can be enforced by constraints

Nonnegative Tensor Factorization

- ▶ *Nonnegative tensor factorization (NTF)*, such as CP decomposition with $\mathcal{T} \geq 0$ and $\mathbf{U}, \mathbf{V}, \mathbf{W} \geq 0$ are widespread and a few classes of algorithms have been developed

Nonnegative Matrix Factorization

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

Coordinate Descent for NMF and NTF

- ▶ Coordinate descent gives optimal closed-form updates for variables in NMF and NTF

Generalized Tensor Decomposition

- ▶ Aside from addition of constraints, the objective function may be modified by using different elementwise loss functions

- ▶ Some loss function admit ALS-like algorithms, while others may require gradient-based optimization