

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

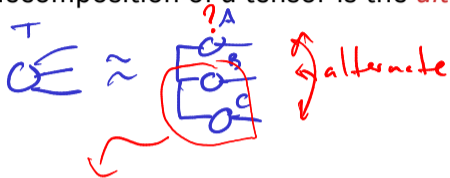
- ▶ Degeneracy motivates an approximate notion of rank, namely *border rank*

Approximation by CP Decomposition

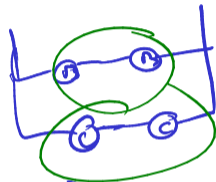
- ▶ Approximation via CP decomposition is a nonlinear optimization problem

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*



$$(B \odot C) A^T \approx T_{(i)}$$



$$\min_A \left\| (B \odot C) A^T - T_{(i)}^T \right\|_F$$

$O(n_1 r^2 + r^3)$
 $O(n_1 r^2)$
 $O(n^3 r)$
 $T \in \mathbb{R}^{n_1 \times n_2 \times n_3}$
 $CP \text{ rank } R$

$$\underbrace{(B \odot C)^T (B \odot C)}_{B^T B \times C^T C} A^+ = \underbrace{(B \odot C)^T T_{(i)}^T}_{M^T + K R P}$$

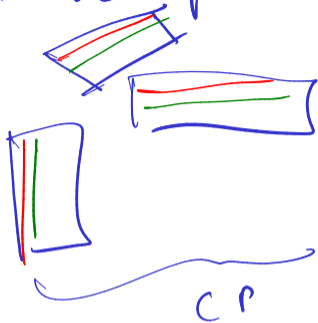
Properties of Alternating Least Squares for CP

+ monotonic decrease of objective

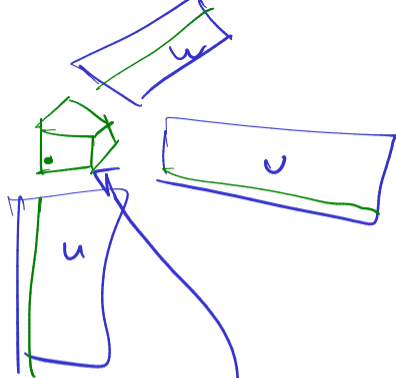
+ stationary point is local minima

- linear (local) convergence rate

Tucker decomposition



$$t_{ijk} = \sum_{r=1}^R a_{ir} b_{jr} c_{kr}$$



$$t_{ijk} = \sum_{r=1}^R \sum_{s=1}^R \sum_{q=1}^R z_{rpq} u_{ir} v_{js} w_{kq}$$

$$u^T u = I$$

No SVD

+ direct / one-shot

+ exact if Tucker is exact

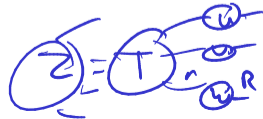
$$T_{(1)}^T = U M_{(1)}^T$$



sequentially connected No SVD

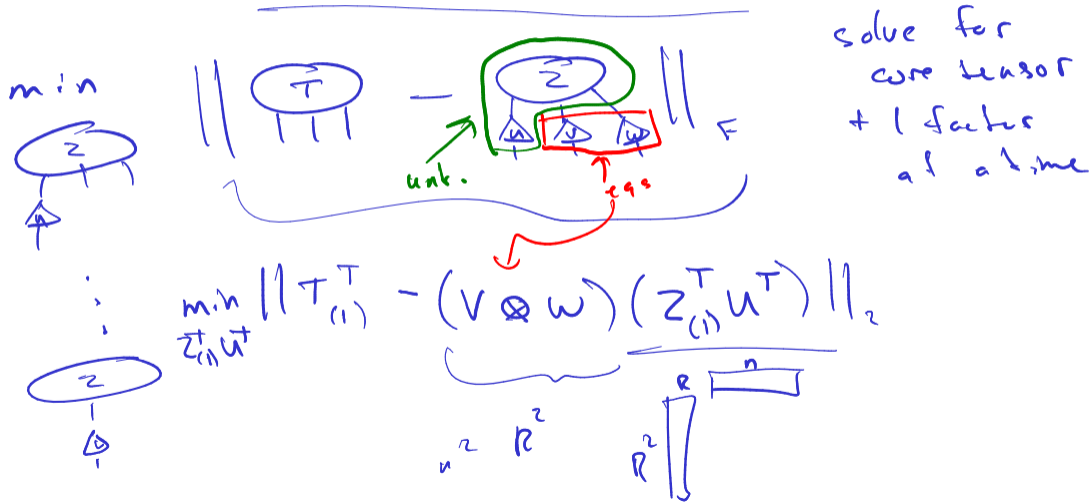
$$T_{(1)}^T = U M_{(1)}^T$$

$$T_{(2)}^T = V N_{(2)}^T$$



Alternating Least Squares for Tucker Decomposition

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



$$\min_{X, \text{rank}(X)=R} \| T_{(1)}^T - \boxed{(U \otimes W) X} \|_F$$

orthogonal

$$(U \otimes W)^T (U \otimes W) = I$$

$$U^T U \otimes W^T W$$

$$(U \otimes W)^T T_{(1)}^T$$

$$\min_X \| G^T B - X \|_F$$

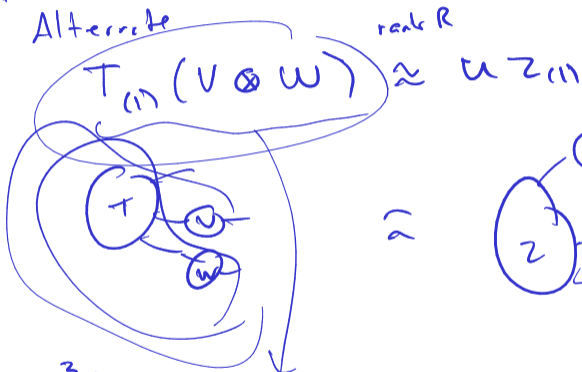
$$\| B - Q X \|_F =$$

HOOI Algorithm

• monotonic program

Iterate

Alternate



$$\frac{S^3 R}{S^2 R}$$

TTMc

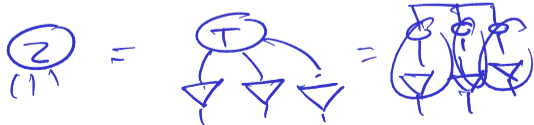
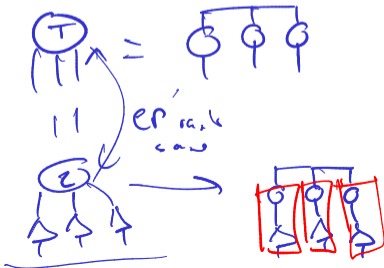
tensor-tensor matrix chain

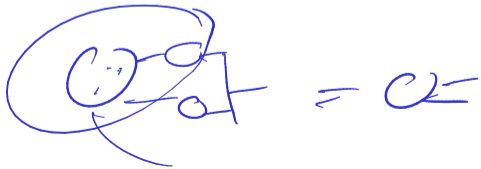
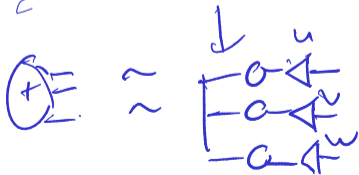
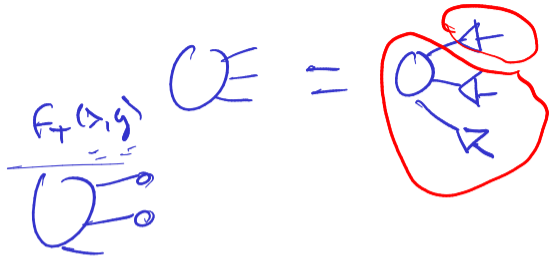
Tucker \leftrightarrow CP

if \forall rank $R < n$

CP rank of Z
is R

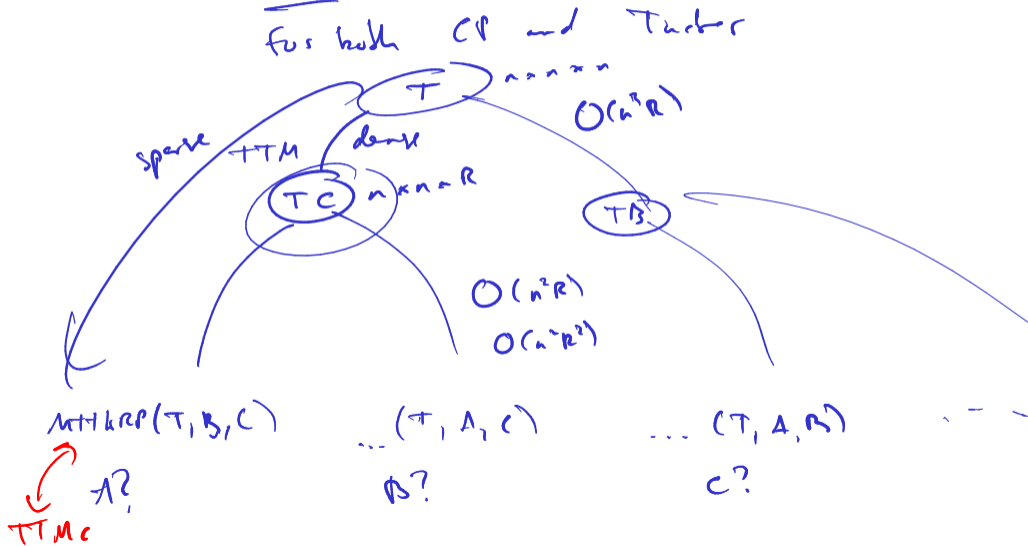
core tensor for rank R
Tucker





Dimension Trees for ALS

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



Fast Residual Norm Calculation

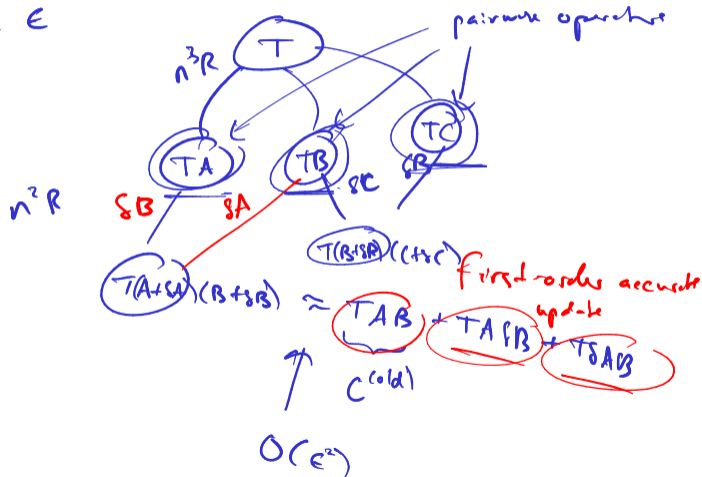
- ▶ Calculating the norm of the residual has cost $2ds^dR$, but can be done more cheaply within ALS

Pairwise Perturbation Algorithm

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

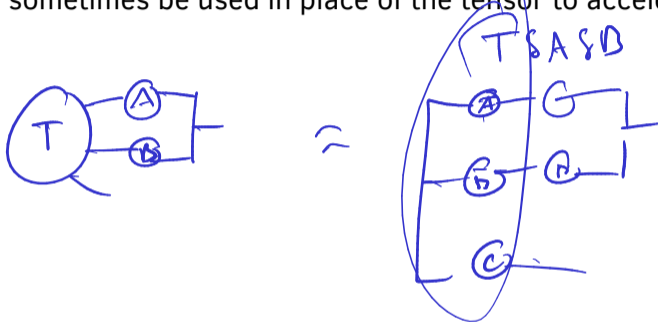
$$\|A^{(i+1)} - A^{(i)}\|_F \leq \epsilon$$

Function ind-ty



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



Approximate CP ALS using Random Sampling

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

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

Gauss-Newton Algorithm

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

$$f(A, B, C) = \frac{1}{2} \left(\begin{array}{c} \oplus \\ \ominus \\ \ominus \\ \ominus \end{array} \right) \left(\begin{array}{c} \oplus \\ \ominus \\ \ominus \\ \ominus \end{array} \right)^T$$

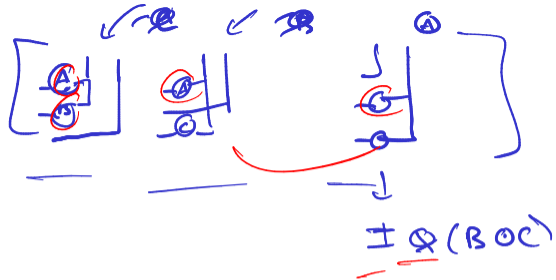
$$f = \frac{1}{2} r^T r$$

$$\nabla_r f = r$$

$$\nabla_r = \begin{bmatrix} \frac{\partial r}{\partial A} \\ \frac{\partial r}{\partial B} \\ \frac{\partial r}{\partial C} \\ \dots \end{bmatrix}$$

$r(A, B, C)$
 \swarrow SR
 \nwarrow n^3

$$\underbrace{J_r^T J_r}_{\text{Hessian}} x = J_r^T r$$

 $J_r =$


$$x = \begin{bmatrix} \text{vec}(A) \\ \text{vec}(b) \\ \text{vec}(C) \end{bmatrix} \in \mathbb{R}^{3nR}$$

$$J_r^T J_r \in \mathbb{R}^{\frac{3nR \times 3nR}{n^2 R^2}}$$

direct method has cost $O(n^3 R^3)$

$$J_r^T J_r$$

$$= \frac{df}{dA}$$

$$I \otimes (B^T B + C^T C)$$

$$J_r \in \mathbb{R}^{3 \times 3 \times n \mathbb{R}}$$

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])
```


Coordinate Descent for CP Tensor Completion

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

Sparse Tensor Contractions

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

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

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

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

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

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