

# CS 598 EVS: Tensor Computations

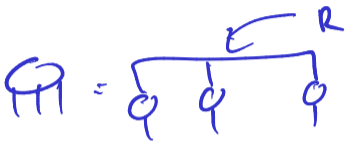
## Tensor Decomposition

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 (rank of a random tensor)

- 1 typical rank for  $\mathbb{C}$

- by d.f. arguments  $R \approx n^{d-1}$   
each factor is of size  $n \times R$

## Typical Rank and Generic Rank

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

# Uniqueness Sufficient Conditions

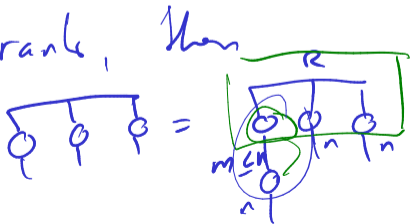
- ▶ Unlike the low-rank matrix case, the CP decomposition can be unique

- up to permutation of  $R$  rank-1 terms  
• and scaling of vectors in each

- uniqueness depends on rank of factors  
• also  $k$ -rank

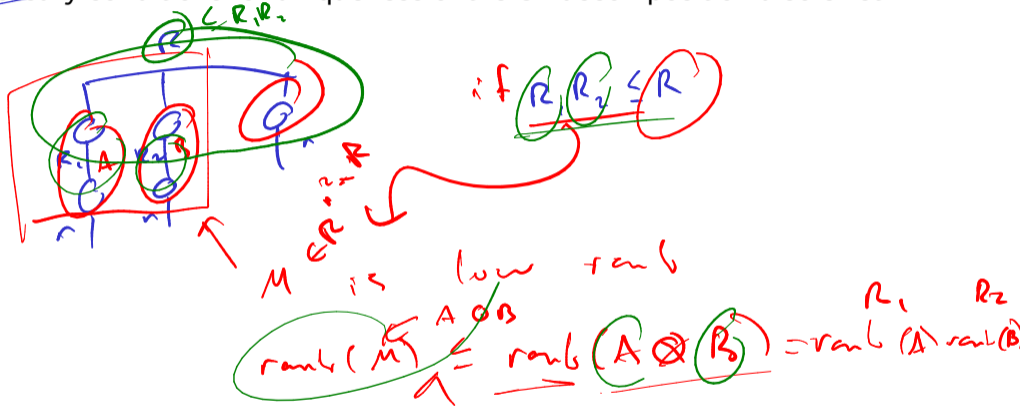
- if  $A, B, C$  are full rank, then

- $R \leq \underline{3n/2}$



# 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

↓  
in matrix case, given by <sup>truncated</sup> SVD

$$\|A - \bar{A}\|_F$$

$$\|T - \bar{T}\|$$

harder rank

$$\bar{T}, \min_{\text{rank}(T)=R}$$



$$\lim_{\lambda=1, \dots, \infty} \|T - T_{\lambda}\| = 0$$

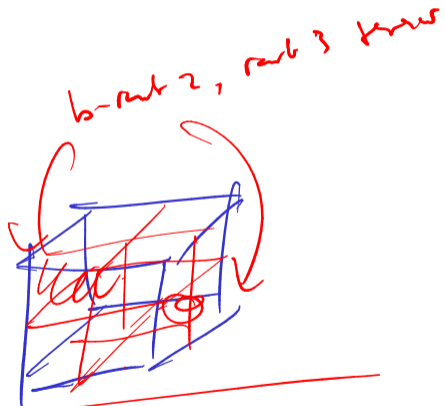
border rank of  $T$

as smallest  $R$ , s.t.  $\exists$  (conformal) sequence

$T_1, \dots$

s.t.


min  $\|T_i \oplus T\| = 0$   
i:1..

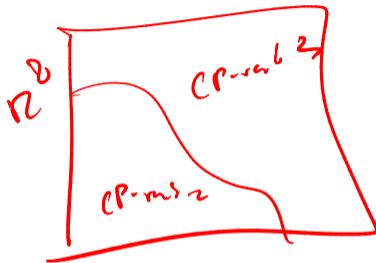




# Border Rank

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

$$R^{2 \times 2 \times 2} \sim R^8$$




# Approximation by CP Decomposition

- Approximation via CP decomposition is a nonlinear optimization problem

nonlinear least squares

$$\min_{u, v, w} \frac{1}{2} \| T - \underbrace{[u, v, w]}_F \|^2$$

$$f(u, v, w) = \delta$$

$$e(u, v, w) = \frac{1}{2} \sum_{ijk} (t_{ijk} - \underbrace{\sum_r u_{ir} v_{jr} w_{kr}}_e)^2 \quad \left( \sum_r u_{ir} v_{jr} w_{kr} \right) \left( \sum_s u_{is} v_{js} w_{ks} \right)$$

$$\frac{de}{du} (u^*, v^*, w^*) = 0 \quad \left| \begin{aligned} \left( \frac{\partial e}{\partial u} \right)_{ir} &= - \sum_{jk} (t_{ijk} v_{jr} w_{kr}) \\ &+ v_{jr} w_{kr} \left( \sum_s u_{is} v_{js} w_{ks} \right) \end{aligned} \right.$$

$$\sum_{j,k} t_{ijk} v_{jr} w_{kr} = \sum_{j,k,s} v_{jr} w_{kr} v_{js} w_{ks} u_{is}$$


---

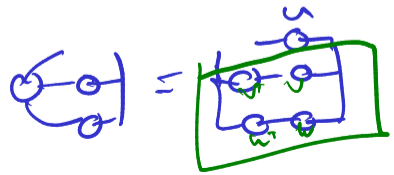
MTHRP

$$= \left( \sum_j v_{jr} v_{js} \right) \left( \sum_k w_{kr} w_{ks} \right) u_{is}$$

$$T_{(1)}(W \circ V) = \left( \left( \underbrace{V^T V} + \underbrace{W^T W} \right) u^+ \right)^T$$

$$u^+ \left( \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \right) \left( \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \right)^T$$

$$\frac{\partial L}{\partial u} = \begin{matrix} \text{---} \text{---} \text{---} \\ | \quad | \quad | \\ \text{---} \text{---} \text{---} \end{matrix} + \begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \end{bmatrix} = 0$$



# Alternating Least Squares Algorithm

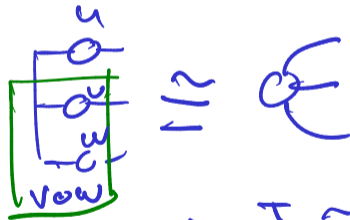
- ▶ The standard approach for finding an approximate or exact CP decomposition of a tensor is the *alternating least squares (ALS) algorithm*

$$U' = \min_U \ell(U, V, W)$$

$$V' = \min_V \ell(U', V, W)$$

$$W' = \min_W \ell(U', V', W)$$

optimality conditions  
= normal eq's of



A diagram of a 3D tensor represented as a rectangular prism. The top edge is labeled 'u', the front-left vertical edge is labeled 'v', and the front-right vertical edge is labeled 'w'. A green rectangular box is drawn around the front face, which is the v-w plane. A blue circle is drawn around the top edge, which is the u dimension.

$$(vow) u^T \approx T_{(u)}$$

## Properties of Alternating Least Squares for CP

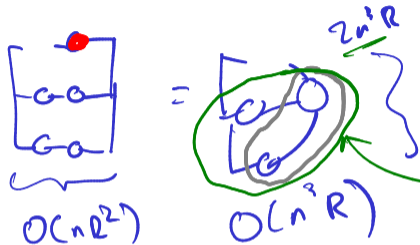
- eqs are defined by Khatri-Rao products
  - amenable to optimizations
- convergence
  - locally (near local minima of  $\epsilon$ )  
linear convergence
  - no global guarantee
  - convergence can stagnate due to iterates becoming ill-conditioned  $(u_i, v_i, w_i)$
- stationary point of ALS, is a local minima or saddle point

## Alternating Least Squares for Tucker Decomposition

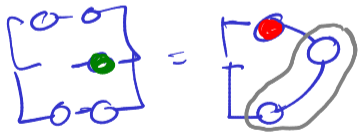
- ▶ For Tucker decomposition, an analogous optimization procedure to ALS is referred to as *high-order orthogonal iteration (HOOI)*

# Dimension Trees for ALS

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



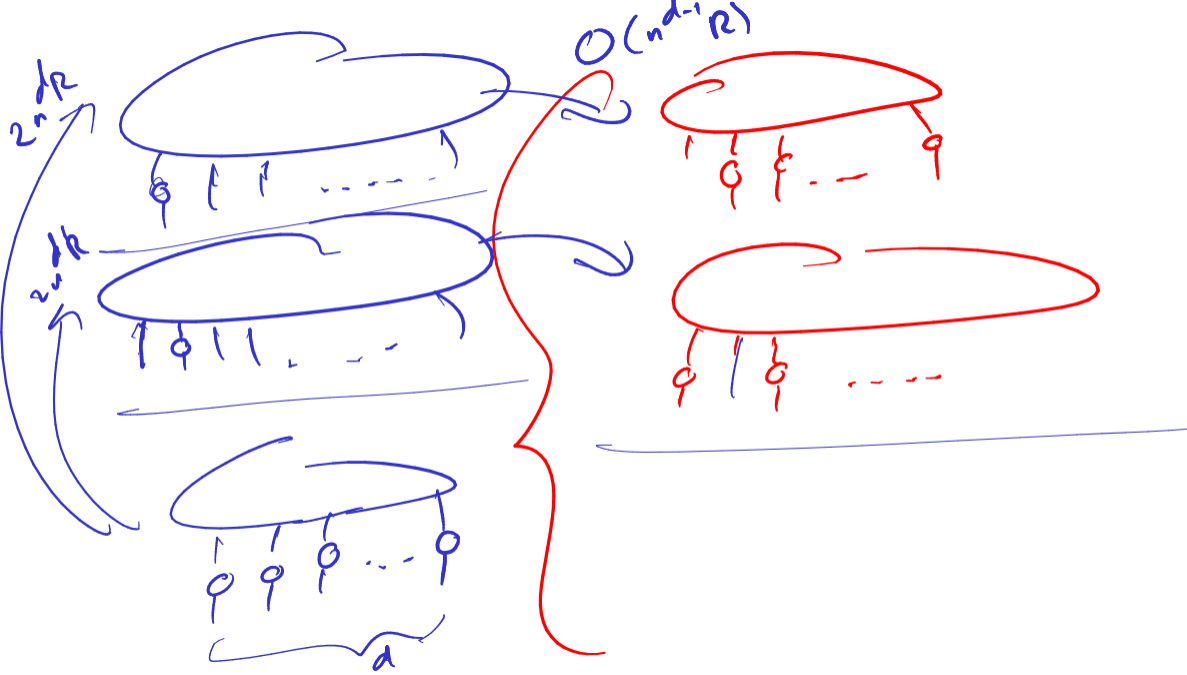
MTTRP =  $\tau + M$  + cheap



order 3: naive  $6n^3 R$   
 d.v. tree  $4n^3 R$



order d: naive  $d 2n^d R$   
 d.v. tree  $4n^d R$





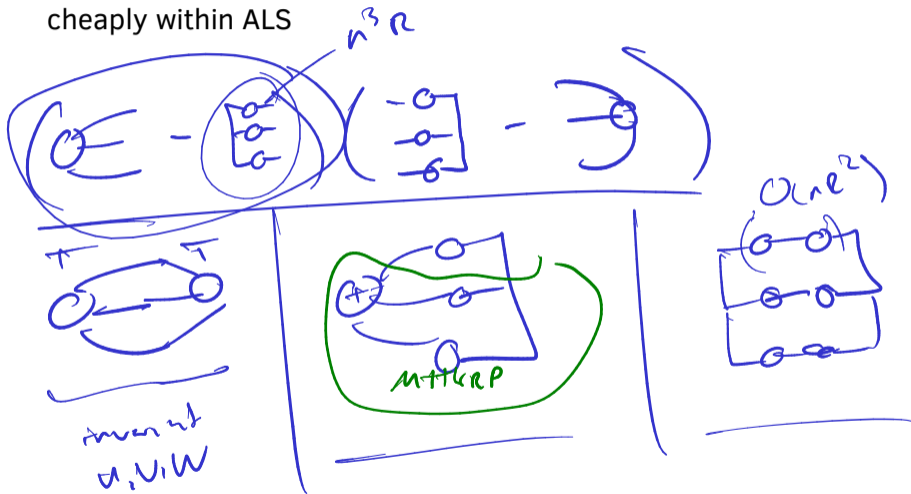


$d-1$  sweeps  
with  $d$  TTMs



# Fast Residual Norm Calculation

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





$$\underline{f_T(u, v)} = \underline{T_{\lambda_1} \times T_{\lambda_2}}$$

## 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

## Approximate CP ALS using Random Sampling

- ▶ Another approach to approximating ALS is to sample the least-squares equations<sup>1</sup>

---

<sup>1</sup>C. Battaglino, G. Ballard, T. G. Kolda, 2018

## 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])
```



## 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 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