

# CS 598 EVS: Tensor Computations

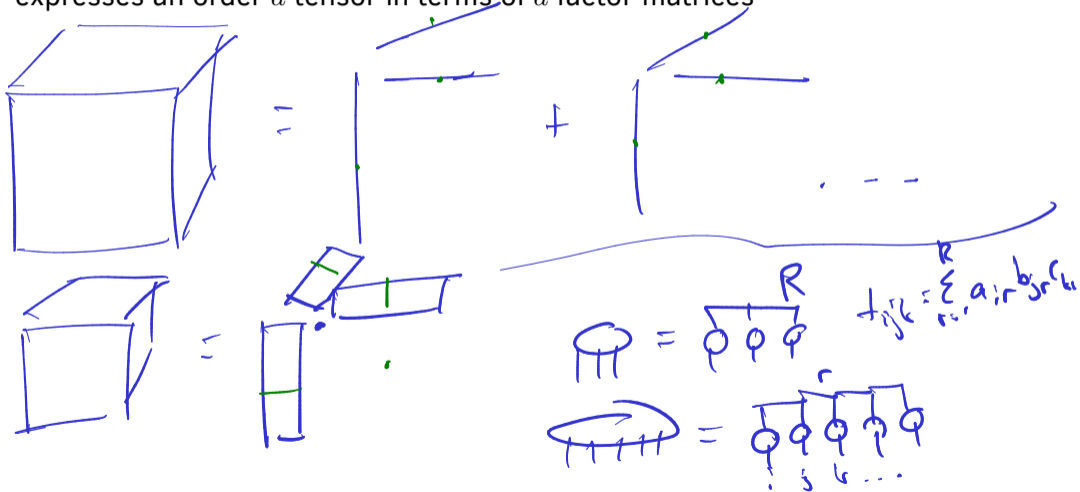
## Tensor Decomposition

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



# CP - CANDECOMP / PARAFAC

- canonical polyadic

1927 Hitchcock

fogou

$n \times n \times n$

## Applications

- bilinear algorithms (exact CP,  $R > n$ )
- data analysis (high-order clustering)
- modelling (CP completion)

→ performance of applications

for different parameters  
# param blocks size  
met. dir.

- low-order methods (compression)
- quantum chemistry

# Tensor Rank Properties

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

smallest  $R$



$$= \underbrace{\left[ \begin{matrix} \leftarrow & & & \leftarrow \\ \leftarrow & + & \dots & + & \leftarrow \\ \leftarrow & & & & \leftarrow \end{matrix} \right]}_R$$

$d=3$   
 $1R + 1R + 1R$   
 $3R \geq n^3$   
 $R \geq \frac{3n^2}{2}$

$$R \leq \left\lceil \binom{n}{d} / n \right\rceil$$

generic rank over  $\mathbb{C}$

$$T = \left[ \left[ \begin{matrix} 1 & & \\ & 1 & \\ & & -1 \end{matrix} \right], \left[ \begin{matrix} & & 1 \\ & & & 1 \\ & & & & 1 \end{matrix} \right] \right]$$

rank over  $\mathbb{R}$  is 3  
 rank over  $\mathbb{C}$  is 2

$$T = [A, B, C] = \left[ \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right]$$

$$K(T) = 1$$

$$T = u \otimes w$$

$$\| -u \| = \| u \| \| -w \|$$

multiple typical ranks over  $\mathbb{R}$

• e.g.  $n=2$ ,  $2 \times 2 \times 2$  tensors

• over  $\mathbb{R}$ , rank is 2 with 79% prob

3 with 21% prob

• for  $n=2$ , there are 2 typical ranks

• if there is only one typical rank, it is the generic rank

• over  $\mathbb{C}$ , for all  $n$ , for symmetric tensors, there is a generic rank

$$R = \left[ \binom{n}{d} \right]_{n!} \quad \binom{n+d-1}{d}$$

## Typical Rank and Generic Rank

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

• for  $n_1 \times n_2 \times n_3$  tensor

$$R \leq \min(n_1 n_2, n_2 n_3, n_1 n_3)$$

# Uniqueness Sufficient Conditions

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

$$A = UV^T = UMM^{-1}V^T \Rightarrow \text{not unique}$$

$$T = \begin{pmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{pmatrix} + \begin{pmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{pmatrix} + \dots$$

permutations

sufficient condition for uniqueness (scaling, permutations) is

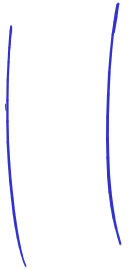
$$k_A + k_B + k_C \geq 2R + 2$$

$k_X$  = rank of  $X$ , largest  $k$  columns of  $X$  is linearly independent  
 $k_A, k_B, k_C \leq n$   
 $R \leq \frac{3n-2}{2}$

rank = 2



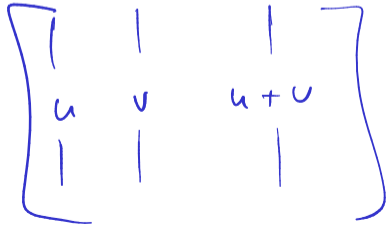
k-rank 0



k-rank 1



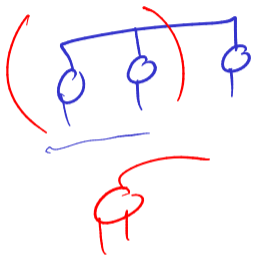
k-rank 2





# Uniqueness Necessary Conditions

- Necessary conditions for uniqueness of the CP decomposition also exist



$$\begin{aligned} \text{rank}(U \circ V) &\geq \text{rank}(U \otimes V) \\ &= \text{rank}(U) \text{rank}(V) \end{aligned}$$

$$\begin{aligned} T &= [U, V, W] \\ \text{rank}(U \circ V) &= R \\ \text{rank}(U \circ W) & \\ \text{rank}(V \circ W) & \end{aligned}$$

$$\text{rank}(U) \text{rank}(V) \geq R$$

# Degeneracy

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

$$T = \underline{a_1 \circ b_1 \circ c_2} + \underline{a_1 \circ b_2 \circ c_1} + \underline{a_2 \circ b_1 \circ c_1}$$

CP rank of  $T = 3$

$$T \approx \alpha(a_1 + \frac{a_2}{\alpha}) \circ \alpha(b_1 + \frac{b_2}{\alpha}) \circ \alpha(c_1 + \frac{c_2}{\alpha}) - \alpha a_1 \circ b_1 \circ c_1$$
$$= T + O(\frac{1}{\alpha}) = \tilde{T}_\alpha$$

$$\lim_{\alpha \rightarrow \infty} \|T - \tilde{T}_\alpha\| = 0$$

# Border Rank

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

- border rank  $R$  implies we can find a sequence  $T_\alpha$

$$\text{CP-rank}(T_\alpha) = R$$

$$\lim_{\alpha \rightarrow \infty} \|T - T_\alpha\| = 0$$

$$\alpha \rightarrow \infty$$

- for bilinear problem tensor  $T$  representing  $3 \times 3$  times  $3 \times 3$  mat.-mult.

- border rank is in  $[14, 21]$

- rank is in  $[19, 23]$

# Approximation by CP Decomposition

- Approximation via CP decomposition is a nonlinear optimization problem

nonlinear least squares

$$\min_{A, B, C} \frac{1}{2} \| \text{tens}(A, B, C) - \sum_{r=1}^R a_r b_r c_r \|_2^2$$

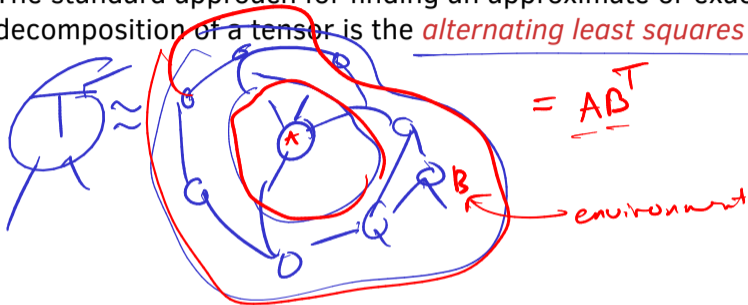
$\frac{\partial f}{\partial A} = \text{MTRP} + \dots$

$\frac{\partial f}{\partial A} = T_{(1)}(C \circ B) + \dots$

$\min_{A, B, C} f(A, B, C) \Rightarrow \frac{\partial f}{\partial A} = 0$

# 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



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

$$A^T = \underline{(B^T B)^{-1}} \underline{B^T T_{(C)}}$$

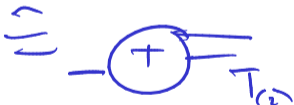
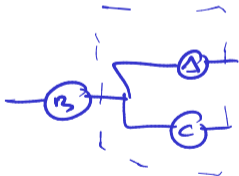
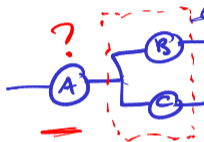
MITKPP

# ALS

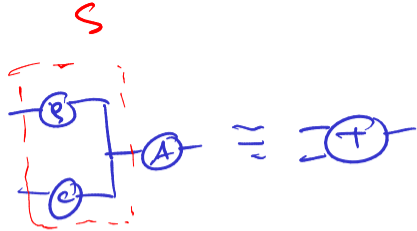
linear ✓

alternating least squares

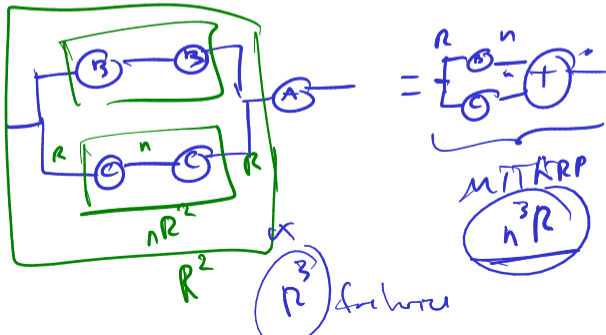
block-diagonal produced before squaring



T\_{(2)}



$$S^T S A = S^T T_{\text{in}}$$



$$\underline{\underline{R_n^2}}$$

## 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

## 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])
```



## Tensor Completion

- ▶ The *tensor completion* problem seeks to build a model (e.g., CP decomposition) for a partially-observed tensor
  
  
  
  
  
  
  
  
  
  
- ▶ The problem was partially popularized by the Netflix prize collaborative filtering problem

# 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 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  $U, V, W \geq 0$  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