# 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

# Random Projections

Accuracy of sketching techniques is theoretically characterized by statistical analysis

# Johnson-Lindenstrauss Lemma

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

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

# Matrix Sketching

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

# Matrix Sketching via Sampling

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

Leverage score sampling provides better accuracy guarantees

## 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
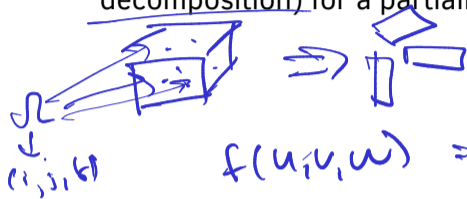
# Approximate CP ALS using Random Sampling

- ▸ Another approach to approximating ALS is to sample the least-squares equations[1]

---

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

# Tensor Completion

- The *tensor completion* problem seeks to build a model (e.g. CP or Tucker decomposition) for a partially-observed tensor
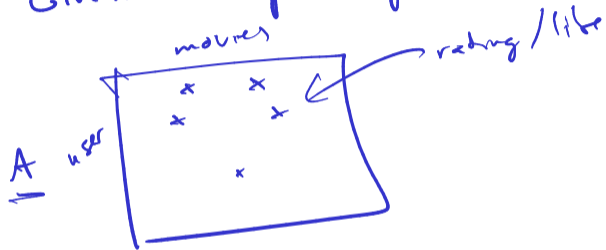
$$\sum_{r=1}^{R} u_{ir} v_{jr} w_{kr}$$

$$f(u,v,w) = \frac{1}{|\Omega|} \sum_{(i,j,k)\in\Omega} \left( t_{ijk} - \langle u_i, v_j, w_k \rangle \right)^2 + \lambda \|U\|_F^2 + \cdots$$

- The problem was partially popularized by the Netflix prize collaborative filtering problem

# Matrix Completion

Given a partially obs. matrix



$A$ user
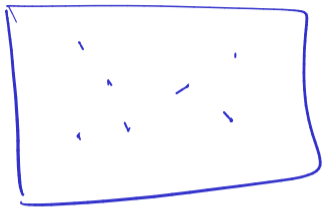
movies

rating / like

$\Omega$ — set of pairs of indices

$a_{ij} \quad \forall \, (i,j) \in \Omega$

---

## recommender system

Netflix prize

## Model



## Objective

$$f(U,V) = \frac{1}{|\Omega|} \sum_{(i,j) \in \Omega} \left( a_{ij} - \sum_{r=1}^{R} u_{ir} v_{jr} \right)^2 + \lambda \|U\|_F^2 + \lambda \|V\|_F^2$$

no closed form solution

but $f(U,V)$ is convex if we fix $U$ or $V$

# CP Tensor Completion Gradient and Hessian

$$\frac{\partial}{\partial u_i} \langle u_i, v_j, w_k \rangle = v_j * w_k$$

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

$$f(u,v,w) = \frac{1}{|\Omega|} \sum_{(i,j,k) \in \Omega} \left( t_{ijk} - \langle u_i, v_j, w_k \rangle \right)^2 + \lambda \|U\|_F^2 + \dots$$

$$\left( \frac{\partial f}{\partial u_i} \right) = \frac{1}{|\Omega|} \left( \sum_{(j,k) \in \Omega_i} \left( t_{ijk} - \langle u_i, v_j, w_k \rangle \right) \right) (v_j * w_k) + 2\lambda u_i$$

$$\Omega_i = \{ (j,k) : (i,j,k) \in \Omega \}$$

▸ 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

$$H_f^{(u,:)}(u,v,w) = \frac{1}{|\Omega|} \sum_{(j,k) \in \Omega_i} \underbrace{(v_j * w_k)(v_j * w_k)^T}_{v_j v_j^T \,*\, w_k w_k^T} + 2\lambda I$$

$$\Lambda_i = \text{all } (j, k) \in \{1, n\}^2$$

$$u = \frac{1}{5} \sum_j \sum_k \underline{v_k v_j^T} \times w_k w_j^T$$

$$u = \frac{1}{5} \left( \sum_j \underline{v_j v_j^T} \right) \left( \sum_v \underline{w_k w_k^T} \right)$$

$$= \frac{1}{5} \underline{\boxed{V^T V}} \times \underline{W^T W} = (V \odot W)^T (V \odot W)$$

fixed

independent

$$\sum_{(i,j,k) \in \Omega_i} t_{ijk} v_{jr} w_{kr}$$

define $\overline{T}$ so that $\overline{T}_{ijk} = t_{ijk} \quad \forall (i,j,k) \in \Omega$

$= 0 \quad$ otherwise

then

$$\sum_j \sum_k \overline{T}_{ijk} v_{jr} w_{kr} =$$

$\underbrace{\qquad\qquad\qquad}_{\text{sparse MTTKRP}}$

# Methods for CP Tensor Completion

$T \in \mathbb{R}^{n \times n \times q}$    $|\Omega| = nn \cdot z$

▸ ALS for tensor completion with CP decomposition incurs additional cost

1. To optimize $U$
   need to form $n$ Hessians. $H^{(ni)} \sim \sum_{(j,k) \in \Omega_i^1} v_j v_j^T \circ w_k w_k^T$
   and solve with                                                    $\pm 2\lambda I$

2. for RHS need MTTKRP with $T$ (observed entries of $T$)
   (1) has cost $O(\sum |\Omega_i| R^2) = O(nn z R^2)$ (2) $O(nnz \cdot R)$
   $+ \text{solve} \; O(nR^3)$

▸ Alternative methods for tensor completion include coordinate descent and stochastic gradient descent

• Coord. desc. minimize for some $u_{ir}$ instead of $u_i$ at a time (less program / subproblem), but updates cheaper

• SGD: consider subgradients

consider $\partial u_i^{(new)} = u_i - \frac{1}{|\Omega|} \cdot \eta \cdot (t_{ijk} - \langle u_i, v_j, w_k \rangle)(v_j * w_k)$
(random $(i,j,k)$)

# Coordinate Descent for CP Tensor Completion

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

$$f(u_i, v_i, w) = \frac{1}{|\Omega|} \sum_{(i,j,k) \in \Omega} \left( t_{ijk} - \langle u_i, v_j, w_k \rangle \right)^2 + \lambda \|U\|_F^2 + \cdots \qquad \sum_{r=1}^{R} u_{ir} v_{jr} v_{br}$$

$$\frac{\partial f}{\partial u_{ir}} = \frac{1}{|\Omega|} \sum_{(i,j,k) \in \Omega_i} \left( t_{ijk} - \langle u_i, v_j, w_k \rangle \right) \left( \qquad v_{jr} w_{br} \right)$$

$$+ 2\lambda u_{ir}$$

$$p_{ijk}^{(r)} = t_{ijk} - \langle u_i, v_j, w_k \rangle + u_{ir} v_{jr} w_{br}$$

$$= t_{ijk} - \sum_{r=1}^{R} u_{ir} v_{jr} w_{br}$$

$$= \frac{1}{|\Omega|} \sum_{(i,j,k) \in \Omega_i} \left( p_{ijk}^{(r)} + u_{ir} v_{jr} v_{rr} \right) \left( v_{jr} w_{br} \right)^{r \neq r} + 2\lambda u_{ir}$$

$$u_{ir} = \frac{1}{|\Omega|2\lambda} \sum_{(j,k)\in\Omega_i} p^{(n)}_{ijk} v_{jr} v_{kr} + \frac{1}{u_{ir} \, |\Omega|2\lambda} \sum_{(j,l)\in\Omega_i} v_{jr} w_{lr}$$

$$= \dots \quad \frac{\sum_{(j,k)\in\Omega_i} \overbrace{p^{(n)}_{ijk}}^{\text{residual-like term}} v_{jr} w_{kr}}{1 + \frac{1}{|\Omega|2\lambda} \sum_{(j,l)\in\Omega_i} v_{jr} w_{lr}}$$

CCD

CCD++



diff. ordering choices

1st,3rd  2nd

nth

cyclic coordinate descent

or

"ALS" ordering

coordinate

Suppose we've updated $u_{ir}$ $\forall i$ (one column)

with

residual $\left| p_{ijk}^{(r)} = (t_{ijk} - \langle u_{ir}, v_{jr}, w_{kr}\rangle) + u_{ir}v_{jr}w_{kr} \right)$

next we may consider $p_{ijk}^{(r+1)} = p_{ijk}^{(r)} - u_{ir}v_{jr}w_{kr}$
$$+ u_{irn}v_{jrn}w_{krn}$$

CCD cost per entry update is $O(|\Omega|)$, so for all factors/cols.

$$O(|\Omega| R)$$



$$\boxed{\cdot} = \boxed{\cdot} - \cancel{\phantom{}} + \cancel{\phantom{}}$$

$O(|\Omega|)$ cost

# 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 $f^{(\mathcal{T})}(x,y) = T_{(1)}(x \odot y)$ 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 $u_{ir} = \sum_{j,k} t_{ijk} v_{jr} w_{kr}$ 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

$$\min_{u,v,w \in S} || T - [\![ u,v,w ]\!] ||_F$$

$$\downarrow$$

$$u,v,w \geq 0 \quad \text{nonnegativity}$$

$$u^T u \approx I \quad \text{(orthogonality or near orthogonality)}$$

$$|| u_i - u_{i,1} || \leq C \quad \text{(continuity)}$$

# Nonnegative Tensor Factorization

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

proximal optimization method

$$\min_{u \in S} \left( f(u) + \tfrac{1}{2}\|u - x\|_2 \right)$$

$u \in S$ — $\mathbb{R}_+^r$

Optimal (outside $S$)

# 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

$$\min_{u,v \in \mathbb{R}_+^{r \times \varrho}}$$

$$f(u,v) = \| T - uv^T \|_2^2$$

$$\frac{\partial f}{\partial u_i} = (T - uv^T) v_i = 0 = \rho v_i - u_i v_i^T v_i \Rightarrow$$

$$\rho = T - uv^t + u_i v_i^T$$

$$u_i = \frac{\rho v_i}{v_i^T v_i}$$

$$u_i = \left| \frac{\rho v_j}{v^T v_i} \right|_+ \qquad u_i = \frac{\rho v_i}{v_i^T v_i}$$

$$\left( |x|_+ \right)_{(i)} = \begin{cases} x_i & \text{if } x_i \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

# 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