# 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

$$r(x) \leftarrow \text{residual}$$
$$\uparrow$$
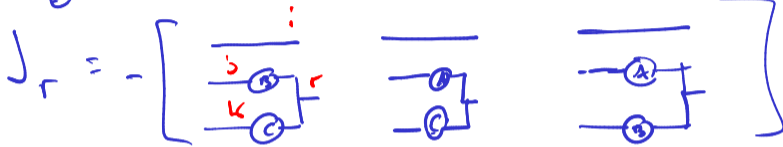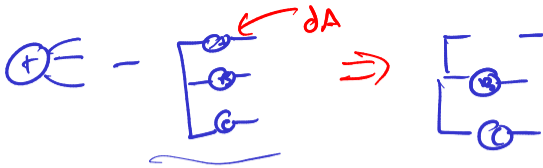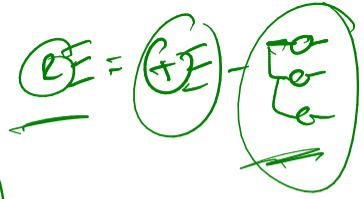$$\text{parameters}$$

$$J_r(x)$$

$$J_r \quad s \stackrel{\sim}{=} r(x)$$

$$\left( J_r^T J_r \right) s = J_r^+ \, r(x)$$

$$r(A, B, C) = \boxed{r} \left\{ \begin{array}{l} \end{array} \right. - \left[ \begin{array}{c} A \\ B \\ C \end{array} \right] \quad \xrightarrow{dA} \quad \Rightarrow \quad \left[ \begin{array}{c} B \\ C \end{array} \right]$$

$$\in \mathbb{R}^{n_A 3 n R}$$

$$J_r = - \left[ \begin{array}{ccc} \left[\begin{array}{c} B \\ C \end{array}\right] & \left[\begin{array}{c} A \\ C \end{array}\right] & \left[\begin{array}{c} A \\ B \end{array}\right] \end{array} \right]$$

$$J_r^T J_r = \left[ \begin{array}{ccc} \end{array} \right.$$

$$J_r^T r = \begin{bmatrix} \text{(R)} \to \text{(B)} \to \text{(C)} \\ \text{(R)} \to \text{(A)} \to \text{(C)} \\ \text{(R)} \to \text{(A)} \to \text{(B)} \end{bmatrix}$$

$R^T = (J^T J) - \begin{bmatrix} e \\ e \\ e \end{bmatrix}$

CG:

$J^T J x$

$n^T R$
$+ n R^2$

# 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

$$\min_x \|Ax - b\|_2$$

$A \in \mathbb{R}^{m \times n}$
$n \ll m$

$$\frac{\|Ax - b\|_2 - \|A\hat{x} - b\|_2}{\|Ax - b\|_2} \leq \epsilon$$

$S$ – sketch matrix
$S \in \mathbb{R}^{k \times m}$   $k \ll m$

$$\min_{\hat{x}} \|SA\hat{x} - Sb\|_2$$

$$\sigma(SA) \approx \sigma(A)$$

$$(size(A) + poly(n)) \log(1/\epsilon)$$

**Accuracy?**

direct ≈ exact A
iterative (CG)
$O(nnz(A) \log(1/\epsilon))$
$\min_{(S nnz(A))}$
with randomized
$O(size(A) + poly(n) poly(\frac{1}{\epsilon}))$
preconditioner

**structure of A?**

A can be sparse
↳ want S to be sparse
A can be structured
CP-ALS: $A = U \odot V$
HOOI: $A = U \otimes V$
$V^T V = U^T U = I$

# 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

$$A \in \begin{bmatrix} & & \\ & x \times x & \\ & & \end{bmatrix} \leftarrow \text{large in norm or linearly independent from other rows}$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & x_1 & x_1 \\ \vdots & x_2 & 0 \\ \vdots & & \ddots \\ 0 & & \ddots \end{bmatrix}$$

Leverage score sampling provides better accuracy guarantees

$$A = QR \qquad P = QQ^T \qquad Ax - b \perp \underline{QQ^+b} \qquad \text{define } S \text{ to sample ith row}$$

$$\ell_i = \| q_i \|_2^2 \qquad \ell_i = P_{ii} \qquad SQRx \cong Sb \qquad \text{of } A \text{ with prob.}$$

$$\underset{\text{ith row of } Q}{\uparrow} \qquad \qquad SQ(SQ)^+ Sb \qquad \frac{\ell_i}{n}$$

$$k = O(n/\epsilon^2) \rightarrow (1 \pm \epsilon) \text{ accuracy in result}$$

# Mixing Techniques

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

$$S \rightarrow \text{Gaussian random square}$$

$$SA = PP \begin{bmatrix} \frac{S}{s} \end{bmatrix} A$$

uniform samples

nearly-uniform leverage scores

$$O(n^2/\epsilon^2)$$

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

uniform samples      random sign matrix

$$S = \begin{bmatrix} \int \int \begin{bmatrix} F_n \end{bmatrix} \end{bmatrix} \begin{bmatrix} -1 & & \\ & -1 & \\ & & +1 \end{bmatrix} \qquad s = O(n/\epsilon^2)$$

subsampled

$$SA = P \; S_n D \; A$$

# 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 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

# Methods for CP Tensor Completion

- ALS for tensor completion with CP decomposition incurs additional cost

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

# 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 $f^{(\mathcal{T})}(x, y) = T_{(1)}(x \odot y)$ 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

# 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

# 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