

[Speaker Note]: Will use TT in place of MPS sometimes.

HSDC : Constructing Tensor Trains / MPS

• Given: Tensor A in D dimensions, general fmt.

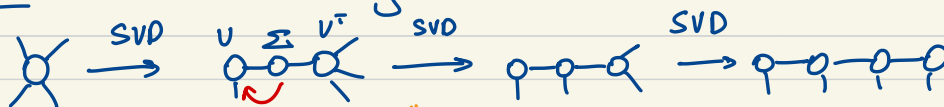
- Giant table (dense)
- Sparse (mostly zero)
- Black-box, on-demand access

Source Paper:

TT-Cross Approximation for Multidimensional Arrays, Oseledets & Tyrtysnikov.

• Produce: MPS approximation T w/ cores (T_1, T_2, \dots, T_D)

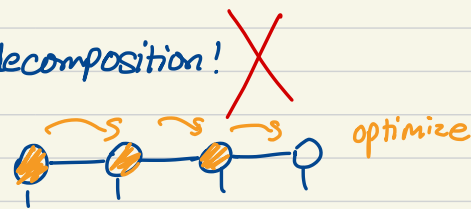
• Recall: Successive SVD algorithm:



Replace w/ R-SVD [if ranks low], but still "see" whole tensor at step 1.

• No sub-exp. runtime algs w/ guaranteed decomposition!

Heuristics: 1) Start w/ random MPS



2) for $i=1 \dots D$:

optimize core i

for $i=D \dots 1$:

optimize core i

optimize objective | core at a time.

3) Repeat step 2 until convergence.

Two perspectives

- Volume maximization
- L2 error minimization

TT-Cross Heuristic: - Simple & Fast (iterative)

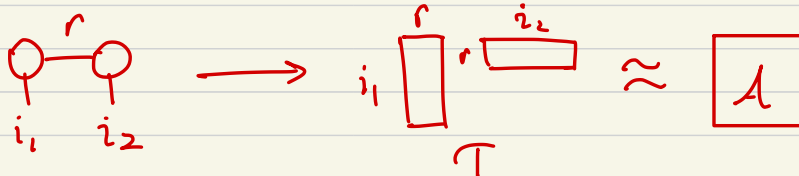
- Few guarantees

- Difficult to beat (when A has sufficient structure)

Backbone of functional tensor train algorithms.

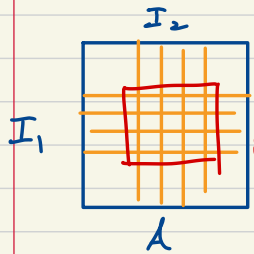
Begin w/ special case: 2D tensor train = low-rank matrix approximation problem.

Briefly pretend that $N \approx 100$.



A volume maximization perspective

• CUR decomposition (aka cross/skeleton decomposition):



• Select $I_1 \subseteq [N]$ rows of A , $I_2 \subseteq [N]$ columns.

Want "most" lin. indep. set of rows & cols.

Define $C = A[:, I_2]$
 $U = A[I_1, I_2]$, $A \approx CU^T R$
 $R = A[I_1, :]$

Moore-Penrose pseudo-inverse.

• If A exactly rank- R , then select any lin. indep. subset of rows I_1 , cols I_2 .

Equality holds: $A = CU^T R$.

Motivation: Suppose singular values beyond R are small.

• NOT EXACTLY low-rank: select I_1, I_2 to maximize $|\det U|$.

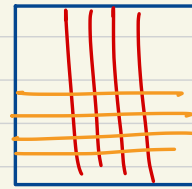
Intuition: Higher determinant \leftrightarrow Higher diversity of sampled rows & cols

NP-Hard. Settle for a greedy algorithm:

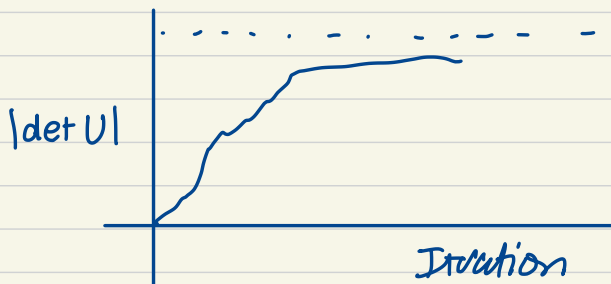
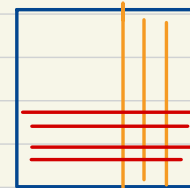
Init: Random I_1, I_2 .
Repeat: • Greedily exchange rows to increase $|\det U|$
 • Greedily exchange cols to increase $|\det U|$
Repeat.



Greedy swaps



Greedy swaps



• $|\det U|$ increases monotonically, may not converge to global maximum, but will converge.

[Other ways to do this:

- LU w/ column-pivoting
- If A : Ridge leverage scores
- p.s.d Determinantal Point Process Sampling.]

◦ In practice: Use QR decomposition. Numerical stability in case rank overestimated.

$$R = A(:, J)$$

$$C = A(I, :), \text{ compute } C = QT$$

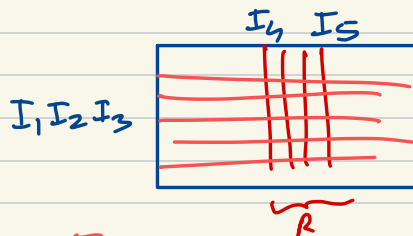
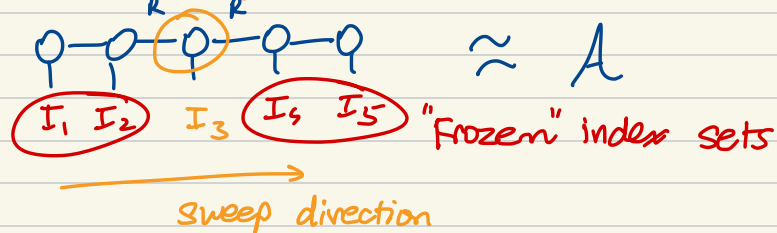
$$\text{then } A_{k+1} = A(:, J)$$



Adapt to general MPS: greedy volume maximization on the matricizations of the tensor.

→ Each set has cardinality R .

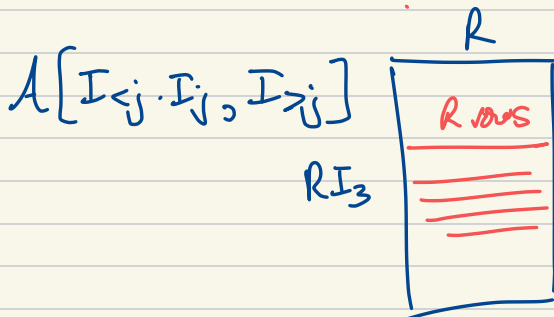
1) Start w/ random I_1, I_2, \dots, I_D . Sweep left to right, right to left:



Update Rows (Too many rows to select from!)
Restrict to a candidate set of left-nested indices:

I_1	I_2	I_3
0	3	?
5	4	?
2	7	?
1	2	?
⋮	⋮	⋮

max # of choices: $R I_3$.

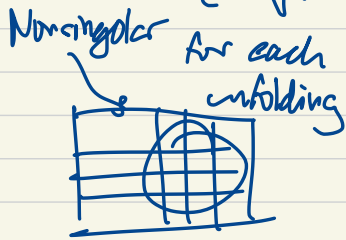


Run greedy maxvol row swap algorithm on this restricted matrix to update I_3

1) Caveat: In practice, run maxvol on the QR decomposition of the selected matrix; prevents the case where the ranks are over-estimated.

Problem: Given cores C_1, C_2, \dots, C_D from the index sets where

$$T[I_{e_j}, I_{s_j}] \text{ nonsingular for } \forall j_s$$



Construct

$$\hat{C}_k = C_k \times_3 \hat{A}_k$$

$$\hat{C}_1 = C_1 \hat{A}_1$$

Max evaluations of tensor: $R^2 I$ per core update, NIR^2 per sweep.

Guarantees

- Converges to best fit TT? ~~X~~
- Converges? ~~X~~

Notice: no monotonically increasing objective.

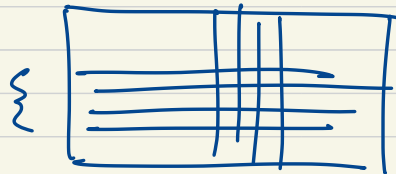
But in practice, usually does converge.

part of single core

- "Fixed point" property?



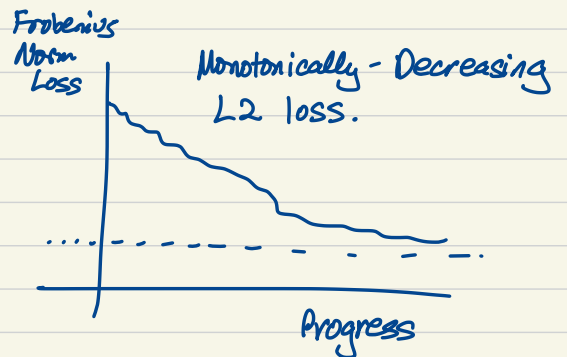
Algorithm finds a lin. indep. set of columns



Alternative view: Alternating Least-Squares

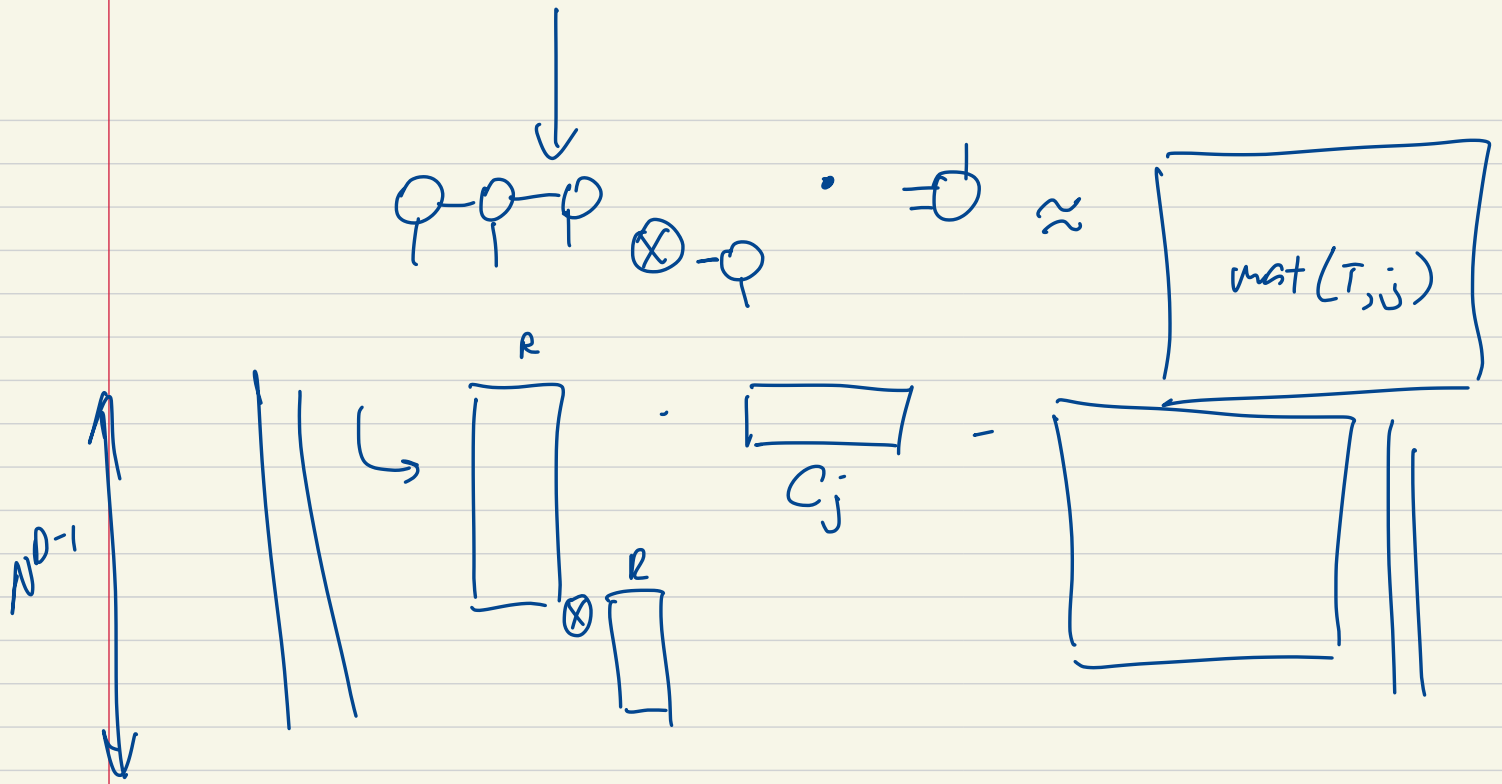
for $i = 1 \dots D$
optimize core i to minimize $\|T - A\|_F$.

for $i = D-1 \dots 1$
" " " "



Just linear least-squares problem looks like this:





Thm: There exist sketching algorithms to solve the LSTSQ problem to residual accuracy ϵ w.h.p. $(1-\delta)$ in poly-time.

Proof: Apply a Johnson-Lindenstrauss transform matrix to both sides.

[Note: MPS random projections are useful for a bunch of things].

HSDC Seminar Round 2

Agenda

- Recap: MPS-Cross & Data Plots
- TT-ALS & Modifications to TT-ALS
- Random projections for MPS & new results.

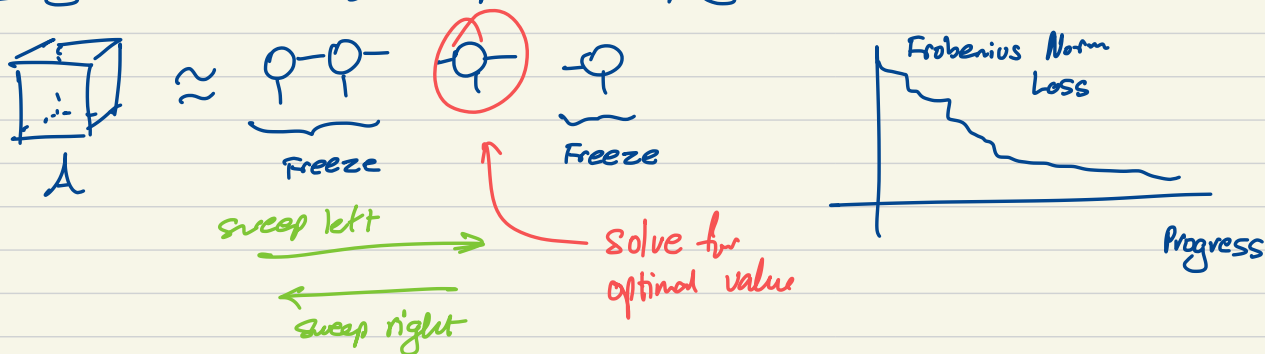
◦ Last time: Volume maximization to construct TT-tensor.

- Manual on each matricization of tensor T , use indices to construct tensor. Start w/ random index sets from tensor and then refine.

◦ Problem: Local minima, only works for the low-error case.

[Show graphs of TT-cross heuristic]

◦ Today: Linear least-squares approach / Adapting the tensor-train rank.



Advantages: • Monotonically decreasing objective
• Parallel, simple

• Can be adapted to deal w/ missing entries.

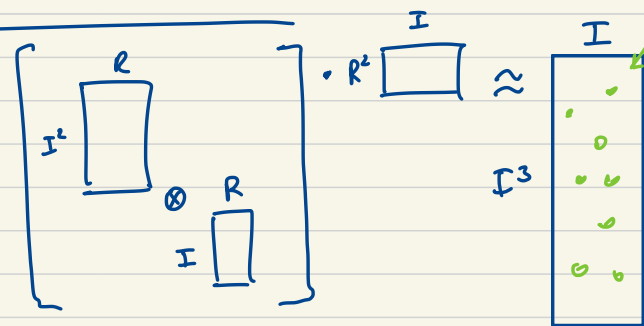
Disadvantages: • Too expensive w/ out modifications
• Needs whole tensor*

Randomized algs address this

sparse set of known entries

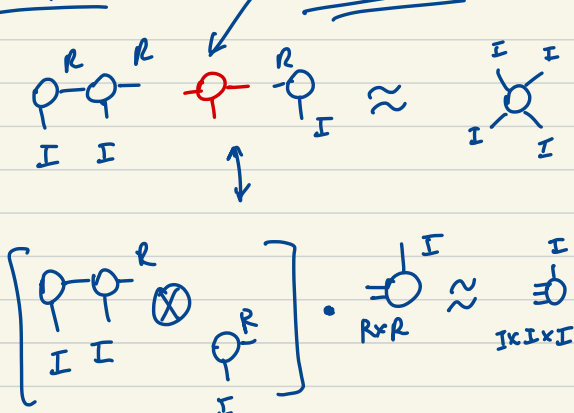
+

"Less chance" of local minimum.



◦ TT-ALS:

Solve for this



Cost to solve this: $O(I^4 R^2 + I R^4)$
[by QR decomposition of LHS, cost to multiply by Q , backsolve w/ Kronecker product R matrix].

Might as well perform TT-SVD, so this is not that useful.

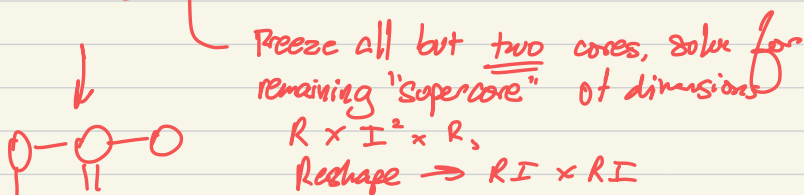
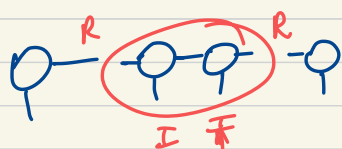
[Note: this alternating scheme also used for eigenvalue problems where both matrix & vector have MPO/MPS form: DMRG algorithm]

Modifications:

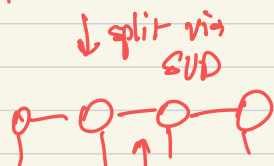
- Missing entries: Solve for $[T_1, \dots, T_0]$ representation that minimizes loss w.r.t. sparse subset of known entries, representation generalizes to the unknown entries.

Solve for each core in time $O(\text{rank}(A)R^2)$, system is smaller.

- Adaptive Rank (also applies to TT-cross):



Freeze all but two cores, solve for remaining "supercore" of dimensions $R \times I^2 \times R$, Reshape $\rightarrow RI \times RI$

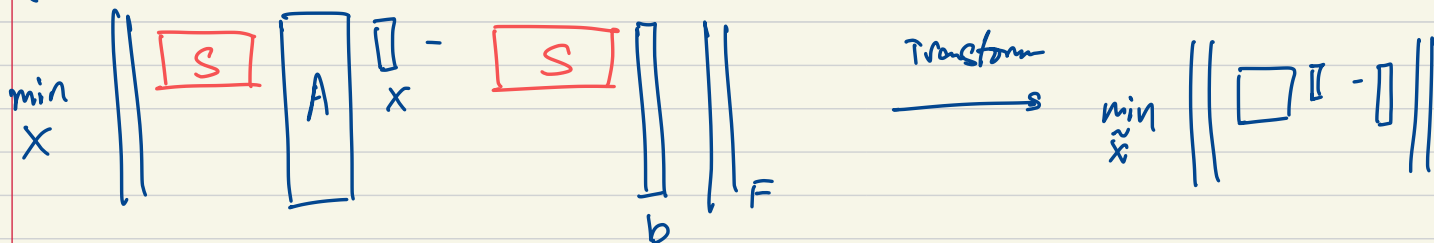


split via SVD

Compute SVD, truncate at appropriate threshold.

New adaptively determined rank.

MPS Random Projections



$$\min_x \|Ax - b\|_2^2$$

$$\min_{\tilde{x}} \|SA\tilde{x} - Sb\|_2^2$$

Expensive. too many rows!

Soln: Apply sketching / sampling matrix to system; [S has fewer rows than columns].

Choose S so that $\|A\tilde{x} - b\|_2^2 \leq (1 + \epsilon) \min_x \|Ax - b\|_2^2$

How to do this? [Subject of active research w/ well-known results]

• Turns out S can be (works w/ high probability)

- An i.i.d Gaussian matrix [but MM ruins speedup]
- A sparse ± 1 random matrix w/ 1 nonzero per column [or several nonzero]
- An MPS where every core has [normalized] i.i.d Gaussian random entries

Statistical \star
leverage scores

- A sampling matrix where each row selected w/ probability α
 $l_i = A_i; (A^T A)^+ A_i^T$.



(subspace embedding, low distortion)

Key property: Let $A = U \Sigma V^T$ be SVD of A . Want

$K(SU)$ as small as possible.

Intuition: vectors that are orthogonal in original space almost orthogonal in transformed space.



"Embeds" subspace of A into a smaller subspace while preserving distances between vectors.

[These sketches have other uses, e.g. MPS-rounding].

Key Challenge: How to sketch when A is in tensor train format, B is in general A^m mat (matricized tensor?)

Output Of Sketch: Row count $O(R^2 \log(\dots) / \epsilon \delta)$ if column count is R^2 .

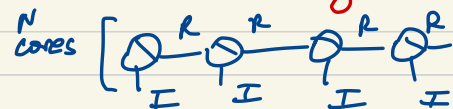
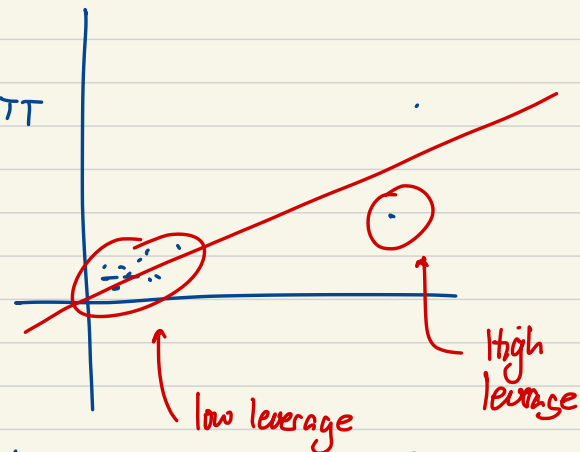
Results on Leverage Score Sampling

- Facts: Only $O(\frac{R^2 \log(R)}{\epsilon} / \frac{R}{\delta})$ rows needed for TT problem.

Leverage scores do not depend on obs. matrix B (!)

- Expensive to compute leverage scores. For A w/ I^N rows & R^2 columns,

$O(I^N R^4)$ to compute scores.



New Result: Can sample 1 row according to leverage score distribution of a TT in canonical form in time $O(NR^2 \log I)$, which is time required to form that row.

Time Complexity of ALS: $O(NIR^n)$

↓ Sketching

$$\tilde{O}(N(NR^n + R^6))$$

Hides log factors

ok when R is small,
space usage of dense
TT grows quadratically in
rank R .

- Compromise: "See" whole tensor eventually, but make progress faster than TT-SVD
(Full-batch vs. Stochastic gradient descent)