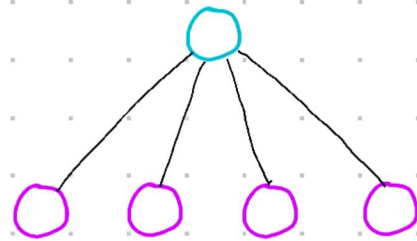


# Stochastic Alternating Least Squares For Tensor Decomposition.

Luke Oeding Joint Work with

Yanzhao Cao  
Somak Das

Hans Werner Van Wyk



Auburn University

ArXiv: 2004.12530

A  
B  
C  
D

A	G	T	-	-	-
A	C	T	-	-	-
A	G	T	-	-	-
A	G	C	-	-	-

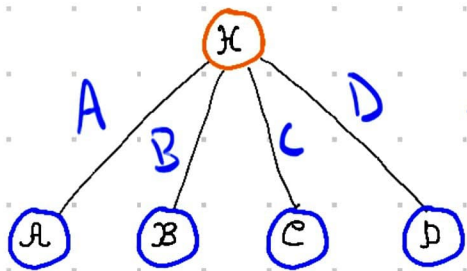
S

Data

Contingency table  
of occurrences of  
Quadruples of letters

Tabcd increased by 1  
whenever see a,b,c,d

# Motivating Example :



(so  $K=4$ )

Given Aligned Sequences of DNA from species that are assumed independent given a common ancestor

Determine the parameters of this model

$$\sum_s A_s \otimes B_s$$

The Model:  $P_{abcd} = P(A=a, B=b, C=c, D=d)$

$$= \sum_{h \in \mathcal{H}} P(A=a, B=b, C=c, D=d | H=h) \cdot P(H=h)$$

$$= \sum_{h \in \mathcal{H}} P(A=a | H=h) \cdot P(B=b | H=h) \cdot P(C=c | H=h) \cdot P(D=d | H=h) \cdot P(H=h)$$

$$P_{abcd} = \sum_{h \in \mathcal{H}} A_{ah} B_{bh} C_{ch} D_{dh} \cdot \alpha_h$$

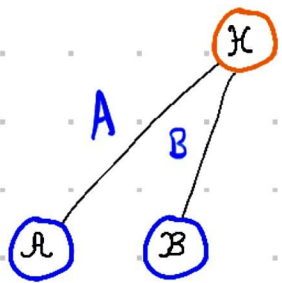
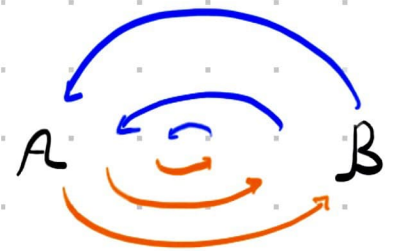
$$= \sum_{h \in \mathcal{H}} \alpha_h A_h \otimes B_h \otimes C_h \otimes D_h$$

The Data:

$T = (T_{abcd})$

A:	A	C	G	T	...
B:	A	C	G	T	
C:	G	T	A	T	
D:	A	C	C	T	

# 2-factor Example Alternating Least Squares



$$T = (T_{ab}) = \frac{1}{n} (C_{ab}) \quad n - \text{Samples}$$

$$P = A B^T + E$$

See [Li - Kunderman - Navasca 2013]

$$\text{Min}_{A, B} \| T - A B^T \|^2 + \frac{\lambda}{2} (\|A\|^2 + \|B\|^2)$$

Regularization

Not convex, in general very hard,

If have A approximate B: Regularized ordinary least squares

$$\hat{b}_u^t$$

$$(A^T A + \lambda I)^{-1} A^T T_u$$

If have B approximate A: ROLS

$$\hat{a}_v^t = (B^T B + \lambda I)^{-1} B^T T_v^t$$

Pseudo code

Guess  $A$

While  $\|T - AB^T\|^2 + \frac{\lambda}{2}(\|A\|^2 + \|B\|^2) > \epsilon$

[ Approximate  $\hat{B}$  holding  $A$  fixed  
 $B \leftarrow \hat{B}$   
Approximate  $\hat{A}$  holding  $B$  fixed  
 $A \leftarrow \hat{A}$

Note ROLS steps don't increase error

## Options For ROLS

- use a closed form solution
- use SVD
- use an iterative method (Gauss-Seidel)
- use a stochastic method
- Reformulate to take advantage of sparsity

# Scale Ambiguity and Regularization

$$T \approx AB^T$$

$$UU^T = I$$

$$= (AU)(U^TB^T)$$

$$T = (A \text{ Diag}(d_1 \dots d_n)) (\text{Diag}(d_1^{-1} \dots d_n^{-1}) B)$$

Option 1: Restrict cols of  $A$ ,  $B$   
to be all on a unit sphere -

Option 2: Regularize

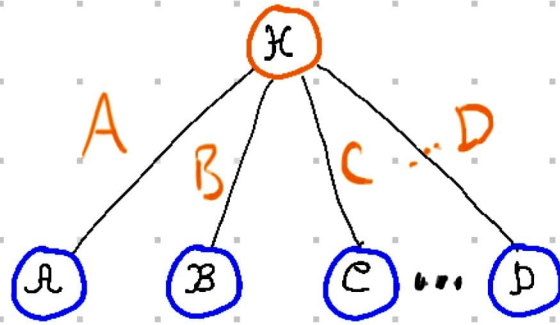
Introduce penalty term

$$+ \frac{\lambda}{2} (\|A\|^2 + \|B\|^2)$$

$$\|AD\|^2 \leq \|A\|^2 \|D\|^2$$

$$\|D^{-1}B\|^2 \leq \|D\|^2 \|B\|^2$$

d-factors



[Carroll-Chang 1970]

$$(T_{abcd}) \approx (P_{abcd}) = A \circ B \circ C \circ \dots \circ D$$

"highly" non-linear  $= \sum_{h \in H} \tau_h A_h \otimes B_h \dots \otimes D_h$

Minimize  $\|T - A \circ B \circ C \circ \dots \circ D\|^2 + \frac{1}{2}(\|A\|^2 + \|B\|^2 + \|C\|^2 + \dots + \|D\|^2)$   
A, B, C, ... D

Guess A, B, C

While objective > tol

D ← solve ROLS for D with A, B, C fixed

A ← solve ROLS for A with B, C, D fixed

B ← solve ROLS for B with A, C, D fixed

C ← solve ROLS for C with A, B, D fixed

# Stochastic Methods

## Warmup Real Time Averaging

Given  $\{x_1, x_2, \dots, x_n, \dots\} = x$

estimate the average (Expected Value)

$$E(x) \approx \frac{1}{n} \sum x_i = y_n$$

Sample and update  $y_{n+1} = \frac{n}{n+1} y_n + \frac{1}{n+1} x_{n+1}$

## Deterioration of Sparsity

$\vec{x}_i$  only with few non-zero terms

Large sample  $\rightarrow$  Dense average.

# Basic Newton Method for Minimization

Assume  $f(x)$  is  $C^2$  (use 2<sup>nd</sup> order Taylor approximates)

guess  $x_1$

$$\text{Iterate } x_{k+1} = x_k - \frac{f'(x_k)}{f''(x_k)} \epsilon$$

until  $\|f(x_k)\|^2$  stops improving

Learning Rate  $\gamma$

$$x_{k+1} = x_k + \gamma \epsilon$$

## Higher Dimensions:

$$f'(x) = \nabla f(x) = g(x)$$

$$f''(x) = \nabla^2 f(x) = H(x)$$

guess  $x_1$

$$\boxed{\text{update}} \quad x_{k+1} = x_k - \gamma [H(x)]^{-1} g$$



# Stochastic Newton

$$H_i = \sigma^2 f_i$$

$$g_i = \sigma f_i$$

$$\text{If } f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x)$$

$n$  large

$$\text{update } x_{k+1} = x_k - \frac{1}{n} \left( \sum_{i=1}^n H_i(x_k) \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n g_i(x_k) \right)$$

hard to compute

Instead only update a sample of the directions

Keep a sequence

$$x_k^1 \dots x_k^n$$

← subseq.

← iterate

Sample  $S \subseteq \{1, \dots, n\}$

Only need to recompute  $H_i(x_k)$  and  $g_i(x_k)$

for  $i \in S$

# CP Decomposition from Samples

$k$  fixed

$\# \mathcal{H} = k$

Suppose  $\mathcal{X}$  is a random tensor, we observe samples  $\mathcal{X}^1, \dots, \mathcal{X}^m$

And want to find a CP decomp

factor matrices

$$\hat{\mathcal{X}} = \sum_{l=1}^k a_l^i \otimes \dots \otimes a_l^d = \underline{[A_1 \dots A_d]} = \underline{[x]}$$

Solving

$$\arg \min_{\hat{\mathcal{X}} \in \mathbb{R}^{d \times k}}$$

$$\mathbb{E} [\|\mathcal{X} - \hat{\mathcal{X}}\|^2]$$

ave decomp

equivalently,

$$\arg \min_{A_1, \dots, A_d \in \mathbb{R}^{k \times n}} \left\| \mathbb{E} [\mathcal{X}] - \underline{[A_1 \dots A_d]} \right\|^2$$

decomp ave tensor

We don't have access to  $\mathcal{X}$ , only samples of  $\mathcal{X}$

Regularized Version:  $x = (A_1 \dots A_d)$

$$\arg \min_{x \in \mathbb{R}^{k \times d}} \mathbb{E} \left( \|\mathcal{X} - \underline{[x]}\|^2 + \lambda \|x\|^2 \right) \quad \lambda > 0$$

# SALS

## Block Stochastic Newton minimization

[Maehara, Hayashi, Kawarabayashi 2016]

(Other instances of Randomized tensor decomp)

[Kolda - Hong, Barttaglino - Ballard - Kolda 2018]

(Dense w/  
sparse samples)

(Randomized methods for OLS)

Fix  
Rank say  $r$

Given Samples  $\underline{X}^1, \dots, \underline{X}^m$  of  $\underline{X}$

Find a minimizer for (regularized) CP decomp of  $E[\underline{X}]$

Variables  $x = [A_1, \dots, A_d]$   $\leftarrow A_i$  is  $n \times r$

① Guess  $x^1$

For  $k=1, 2, \dots$  do

Generate Random Sample  $\underline{X}^k = [\underline{X}^{k,1}, \dots, \underline{X}^{k,m_k}]$

For  $i=1 \dots p$  do

Compute sample gradient  $\tilde{g}^{k,i}$  and Hessian  $H^{k,i}$

Compute step size  $\alpha^{k,i}$

Update block  $i$  :

$$x_i^{k+1} = x_i^k - \alpha^{k,i} (H^{k,i})^{-1} \tilde{g}^{k,i}$$

$$\text{So } x^{k,i} = \left( \underline{x}_1^{k+1}, \dots, \underline{x}_i^{k+1}, \underline{x}_{i+1}^k, \dots, \underline{x}_p^k \right)$$

End for

$$x^{k+1} = x^{k,p}$$

End for

## Theorem (Cao-Das-Oeding-Van Wyk)

If the observed data is bounded,

Then the SALS algorithm for the regularized CP decomposition of the average of a random tensor converges <sup>A</sup> to a minimizer

Bounded data Assumption

$$\|X\| \leq M \quad \text{a.s. on } \Omega$$

Boundedness and multilinear

→ Lipschitz Continuity of gradients and Hessians

→ Bounded invertibility of Hessians

→ Bounds on expectations of iterates

$$\mathbb{E} \|x_i^{k+1}\|$$

Regularization ensures:

- ① Minimizers exist
- ② Hessians stay positive definite (and invertible)
- ③ Iterates are bounded in terms of bounds on  $\lambda$ .

Numerical Experiments study effects of

Noise

Learning Rate

Batch Size

Complexity