

Tensor power method

Daniel Hsu

COMS 4772

1

Orthogonal tensor decompositions

2

Moments of hidden variable models

- ▶ Many hidden variable models have observable moments (perhaps after transformation) of the form

$$\sum_{i=1}^k w_i \cdot \boldsymbol{\mu}_i^{\otimes p}.$$

- ▶ Jennrich's algorithm: uses \mathbf{S} ($p = 2$) and \mathbf{T} ($p = 3$) to recover parameters (assuming $\mathbb{S} := \text{span}\{\boldsymbol{\mu}_i\}_{i=1}^k$ has dimension k).

3

Orthogonality

- ▶ Suppose all $w_i > 0$ and $\dim(\mathbb{S}) = k$.
- ▶ Then \mathbf{S} is psd and has rank k .
- ▶ \mathbf{S} defines inner product over \mathbb{S} in which $\{\sqrt{w_i}\boldsymbol{\mu}_i\}_{i=1}^k$ are *orthonormal*:

$$\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{S}^\dagger} := \mathbf{x}^\top \mathbf{S}^\dagger \mathbf{y}.$$

$$\langle \sqrt{w_i}\boldsymbol{\mu}_i, \sqrt{w_j}\boldsymbol{\mu}_j \rangle_{\mathbf{S}^\dagger} = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{if } i \neq j. \end{cases}$$

4

Whitening

- ▶ Can write $\mathbf{S}^\dagger := \mathbf{W}\mathbf{W}^\top$ with rank k matrix $\mathbf{W} \in \mathbb{R}^{d \times k}$ called “whitening transformation”:

$$\mathbf{S}(\mathbf{W}, \mathbf{W}) = \mathbf{W}^\top \mathbf{S} \mathbf{W} = \mathbf{I},$$

so $\{\mathbf{W}^\top(\sqrt{w_i}\boldsymbol{\mu}_i)\}_{i=1}^k$ is ONB in \mathbb{R}^k .

- ▶ Can also apply \mathbf{W} to higher-order tensors, e.g.,

$$\begin{aligned} \mathcal{T}(\mathbf{W}, \mathbf{W}, \mathbf{W}) &= \sum_{i=1}^k w_i \cdot (\mathbf{W}^\top \boldsymbol{\mu}_i)^{\otimes 3} \\ &= \sum_{i=1}^k \frac{1}{\sqrt{w_i}} \cdot (\mathbf{W}^\top(\sqrt{w_i}\boldsymbol{\mu}_i))^{\otimes 3}. \end{aligned}$$

5

Odeco tensors

- ▶ (Symmetric) *orthogonally decomposable (odeco) tensors*:

$$\sum_{i=1}^k \lambda_i \mathbf{v}_i^{\otimes p}$$

where $\lambda_i > 0$ and $\{\mathbf{v}_i\}_{i=1}^k$ is ONB.

- ▶ (Assume positivity of λ_i for simplicity.)
- ▶ Is the decomposition of an odeco tensor unique?
 - ▶ $p = 2$: **no**
 - ▶ $p \geq 3$: **yes**
- ▶ **Variational claim**: for $p \geq 3$, isolated local maximizers of degree- p homogeneous polynomial $f_{\mathcal{T}}(\mathbf{x}) := \mathcal{T}(\mathbf{x}, \mathbf{x}, \dots, \mathbf{x})$ over B^k are $\{\mathbf{v}_i\}_{i=1}^k$.

6

Variational characterization

- ▶ **Claim:** for $p \geq 3$, isolated local maximizers of $f_{\mathbf{T}}(\mathbf{x}) := \mathbf{T}(\mathbf{x}, \mathbf{x}, \dots, \mathbf{x})$ over B^k are $\{\mathbf{v}_i\}_{i=1}^k$.
- ▶ Observation: by orthogonality,

$$f_{\mathbf{T}}(\mathbf{v}_j) = \sum_{i=1}^k \lambda_i \langle \mathbf{v}_j, \mathbf{v}_i \rangle^p = \lambda_j.$$

- ▶ What about other vectors?
- ▶ May as well think of \mathbf{v}_i as i -th coordinate basis vector.

$$\max_{\mathbf{x} \in \mathbb{R}^k} \sum_{i=1}^k \lambda_i x_i^p \quad \text{s.t.} \quad \sum_{i=1}^k x_i^2 \leq 1.$$

- ▶ If both x_1 and x_2 are non-zero, then

$$\lambda_1 x_1^p + \lambda_2 x_2^p < \lambda_1 x_1^2 + \lambda_2 x_2^2 \leq \max\{\lambda_1, \lambda_2\}.$$

- ▶ Hence, better to only have a *single* non-zero entry.
- ▶ I.e., better to have $\mathbf{x} = \mathbf{v}_i$ for some i . □

7

Tensor power method

8

Optimality condition

$$\max_{\mathbf{x} \in \mathbb{R}^k} \sum_{i=1}^k \lambda_i \langle \mathbf{x}, \mathbf{v}_i \rangle^p \quad \text{s.t.} \quad \sum_{i=1}^k x_i^2 \leq 1.$$

- ▶ Lagrangian:

$$\mathcal{L}(\mathbf{x}, \lambda) := \sum_{i=1}^k \lambda_i \langle \mathbf{x}, \mathbf{v}_i \rangle^p - \frac{p}{2} \lambda (\|\mathbf{x}\|_2^2 - 1).$$

- ▶ First-order optimality condition:

$$p \sum_{i=1}^k \lambda_i \langle \mathbf{x}, \mathbf{v}_i \rangle^{p-1} \mathbf{v}_i - p \lambda \mathbf{x} = \mathbf{0}.$$

- ▶ I.e.,

$$\underbrace{\mathbf{T}(\mathbf{x}, \dots, \mathbf{x}, I)}_{p-1 \text{ times}} = \sum_{i=1}^k \lambda_i \langle \mathbf{x}, \mathbf{v}_i \rangle^{p-1} \mathbf{v}_i = \lambda \mathbf{x}.$$

- ▶ Maximizer must be an “eigenvector” of degree-($p-1$) map.

9

Fixed-point iteration algorithm

- ▶ Consider map from first-order condition:

$$\phi_{\mathbf{T}}(\mathbf{x}) := \mathbf{T}(\underbrace{\mathbf{x}, \dots, \mathbf{x}}_{p-1 \text{ times}}, I).$$

- ▶ Goal: find $\mathbf{x} \in S^{k-1}$ that is fixed under

$$\mathbf{x} \mapsto \frac{\phi_{\mathbf{T}}(\mathbf{x})}{\|\phi_{\mathbf{T}}(\mathbf{x})\|_2}.$$

- ▶ “Tensor power method” (De Lathauwer et al, 2000):
 - ▶ Repeatedly apply $\phi_{\mathbf{T}}$ to initial $\mathbf{x}^{(0)} \in S^{k-1}$ (and re-normalize).
- ▶ Question: Does it find the \mathbf{v}_i ?

10

Example

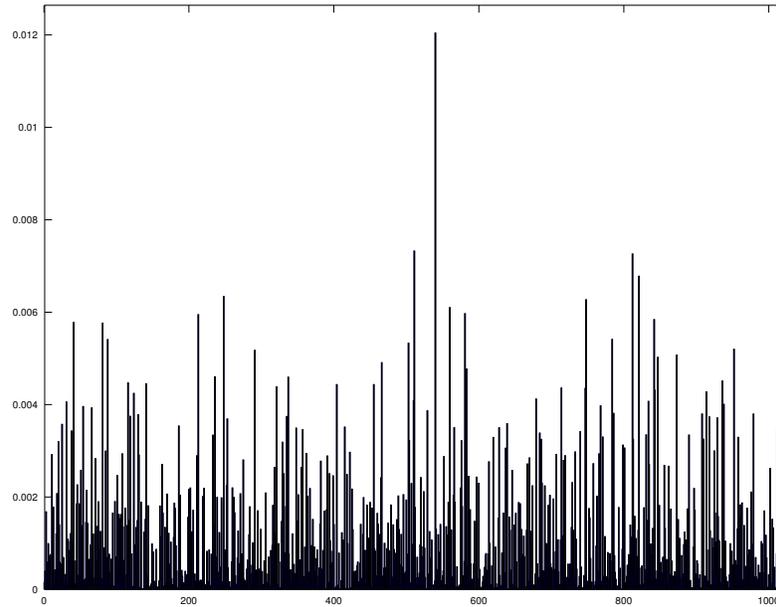


Figure 1: $\langle \mathbf{x}^{(0)}, \mathbf{v}_i \rangle$ for $i = 1, 2, \dots, 1024$

11

Example

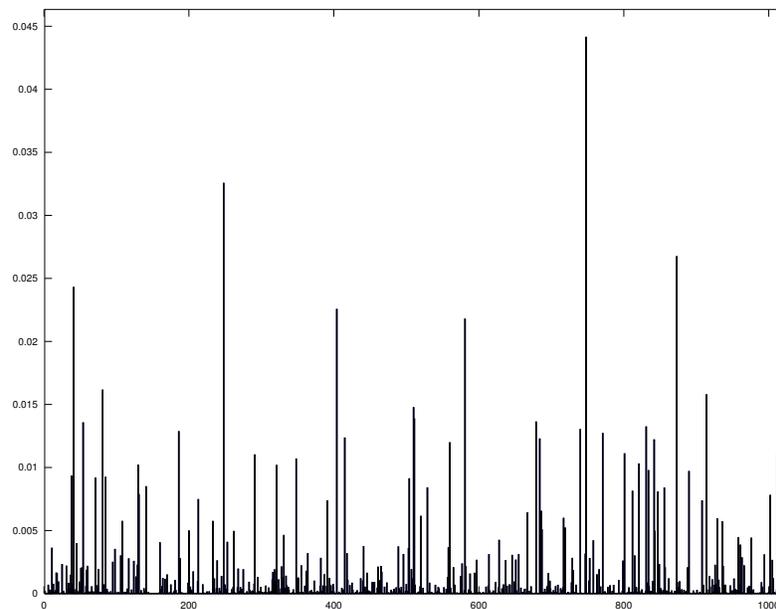


Figure 2: $\langle \mathbf{x}^{(1)}, \mathbf{v}_i \rangle$ for $i = 1, 2, \dots, 1024$

12

Example

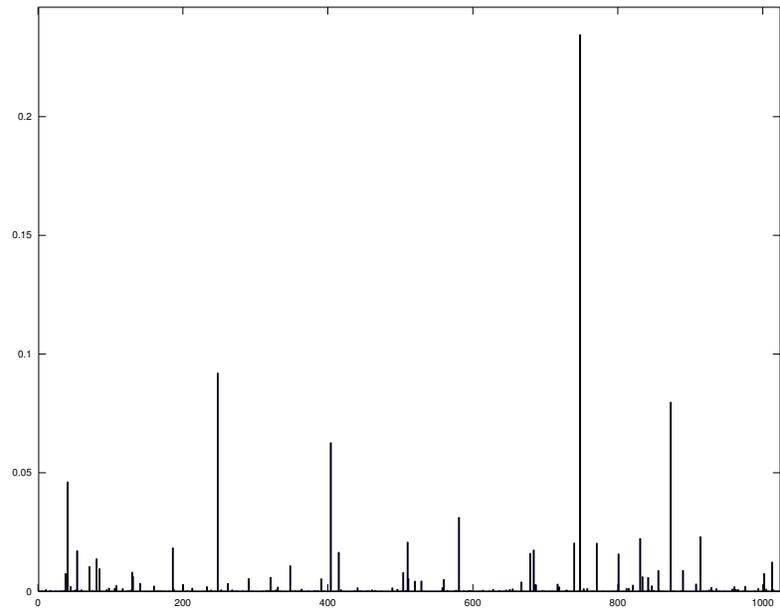


Figure 3: $\langle \mathbf{x}^{(2)}, \mathbf{v}_i \rangle$ for $i = 1, 2, \dots, 1024$

13

Example

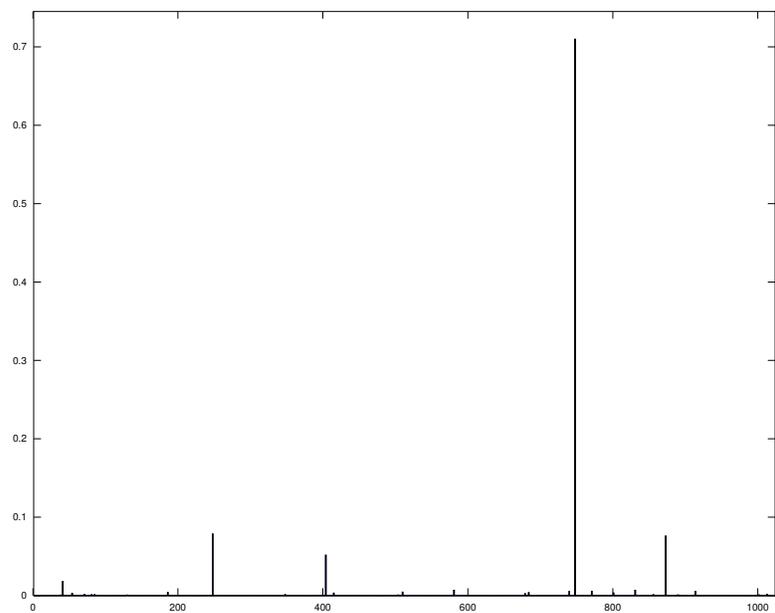


Figure 4: $\langle \mathbf{x}^{(3)}, \mathbf{v}_i \rangle$ for $i = 1, 2, \dots, 1024$

14

Example

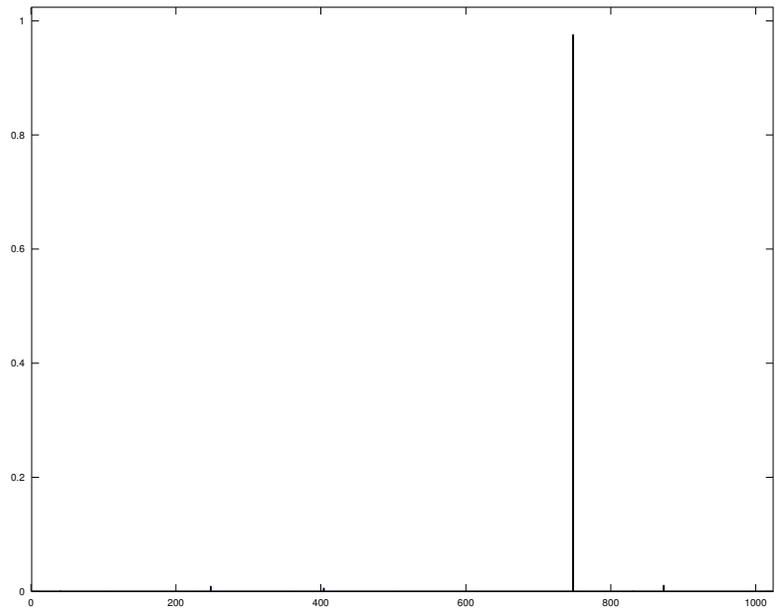


Figure 5: $\langle \mathbf{x}^{(4)}, \mathbf{v}_i \rangle$ for $i = 1, 2, \dots, 1024$

15

Example

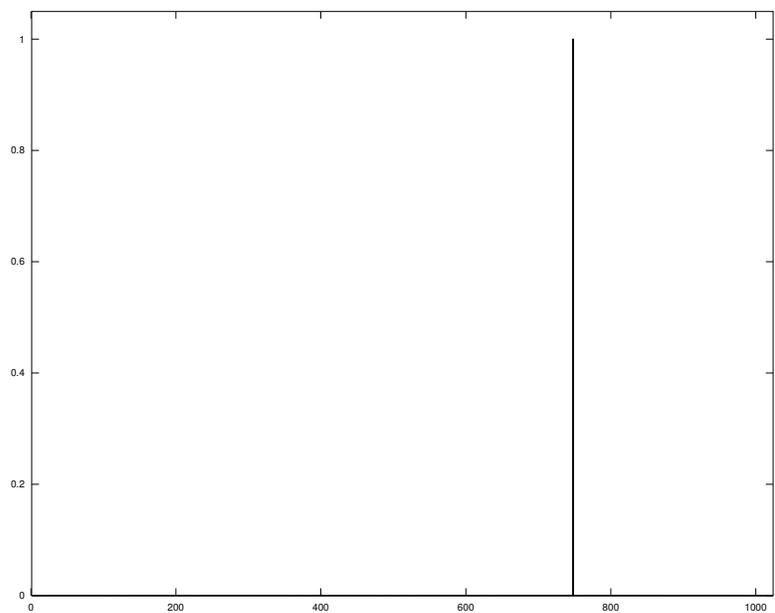


Figure 6: $\langle \mathbf{x}^{(5)}, \mathbf{v}_i \rangle$ for $i = 1, 2, \dots, 1024$

16

Review: matrix power method

- ▶ Tensor power method for $p = 2$ is “matrix power method”:

$$\mathbf{x}^{(t)} := \mathbf{M}(\mathbf{I}, \mathbf{x}^{(t-1)}) = \mathbf{M}\mathbf{x}^{(t-1)}.$$

- ▶ If $\lambda_1 > \lambda_2$, and $\mathbf{x}^{(0)}$ not orthogonal to \mathbf{v}_1 , then angle $\theta^{(t)}$ between $\mathbf{x}^{(t)}$ and \mathbf{v}_1 decreases to zero at linear rate.

- ▶ Write $c_i := \langle \mathbf{x}^{(0)}, \mathbf{v}_i \rangle$ for $i = 1, 2, \dots, k$, so

$$\mathbf{x}^{(0)} = c_1 \mathbf{v}_1 + c_2 \mathbf{v}_2 + \dots + c_k \mathbf{v}_k,$$

$$\mathbf{x}^{(t)} = \lambda_1^t c_1 \mathbf{v}_1 + \lambda_2^t c_2 \mathbf{v}_2 + \dots + \lambda_k^t c_k \mathbf{v}_k.$$

- ▶ $\cos^2(\theta^{(t)})$:

$$\frac{\langle \mathbf{x}^{(t)}, \mathbf{v}_1 \rangle^2}{\|\mathbf{x}^{(t)}\|_2^2} = \frac{c_1^2 \lambda_1^{2t}}{\sum_{i=1}^k (c_i \lambda_i^t)^2} \geq \frac{1}{1 + \left(\frac{\lambda_2}{\lambda_1}\right)^{2t} \frac{1-c_1^2}{c_1^2}}.$$

- ▶ $p = 2$ behavior very different from $p \geq 3$.

17

Tensor power method ($p = 3$)

- ▶ Re-number components so that

$$\lambda_1 |c_1| \geq \lambda_2 |c_2| \geq \dots.$$

- ▶ Then

$$\mathbf{x}^{(1)} = \sum_{i=1}^k \lambda_i \langle \mathbf{x}^{(0)}, \mathbf{v}_i \rangle^2 \mathbf{v}_i = \sum_{i=1}^k \lambda_i c_i^2 \mathbf{v}_i.$$

- ▶ Coefficient c_i in $\mathbf{x}^{(0)}$ is *squared* in $\mathbf{x}^{(1)}$.
- ▶ If $\lambda_1 |c_1| > \lambda_2 |c_2|$, then angle between $\mathbf{x}^{(t)}$ and \mathbf{v}_1 decreases to zero at *quadratic rate*:

$$\frac{\langle \mathbf{x}^{(t)}, \mathbf{v}_1 \rangle^2}{\|\mathbf{x}^{(t)}\|_2^2} \geq \frac{1}{1 + \left(\frac{\lambda_2 |c_2|}{\lambda_1 |c_1|}\right)^{2^{t+1}} \sum_{i=2}^k \left(\frac{\lambda_i}{\lambda_1}\right)^2}.$$

- ▶ Note: which vector we called \mathbf{v}_1 depends on $\mathbf{x}^{(0)}$!

18

Initialization of tensor power method

- ▶ Convergence of tensor power method requires gap between largest and second-largest $\lambda_i |\langle \mathbf{x}^{(0)}, \mathbf{v}_i \rangle|^{p-2}$.
- ▶ Bad initialization:
 - ▶ Suppose $\mathbf{T} = \sum_{i=1}^k \mathbf{v}_i^{\otimes p}$ and $\mathbf{x}^{(0)} = \mathbf{v}_1 + \mathbf{v}_2$:

$$\begin{aligned}\phi_{\mathbf{T}}(\mathbf{x}^{(0)}) &= \langle \mathbf{x}^{(0)}, \mathbf{v}_1 \rangle^p \mathbf{v}_1 + \langle \mathbf{x}^{(0)}, \mathbf{v}_2 \rangle^p \mathbf{v}_2 \\ &= \mathbf{v}_1 + \mathbf{v}_2.\end{aligned}$$

- ▶ But bad initialization points comprise measure-zero set.

19

Recovering all components

- ▶ Power method with $\mathbf{T} = \sum_{i=1}^k \lambda_i \mathbf{v}_i^{\otimes p}$ returns some \mathbf{v}_i .
 - ▶ Can also get λ_i via $\lambda_i = \mathbf{T}(\mathbf{v}_i, \mathbf{v}_i, \dots, \mathbf{v}_i)$.
- ▶ What about other components $j \neq i$?
 - ▶ “Deflation”: replace \mathbf{T} with $\mathbf{T}' := \mathbf{T} - \lambda_i \mathbf{v}_i^{\otimes p}$ so that

$$\mathbf{T}' = \sum_{j \neq i} \lambda_j \mathbf{v}_j^{\otimes p}.$$

- ▶ Can do this “inside” power method:

$$\mathbf{T}'(\mathbf{x}, \dots, \mathbf{x}, \mathbf{I}) = \mathbf{T}(\mathbf{x}, \dots, \mathbf{x}, \mathbf{I}) - \lambda_i \langle \mathbf{x}, \mathbf{v}_i \rangle^{p-1} \mathbf{v}_i.$$

- ▶ Implicitly tries to make power method (with \mathbf{T}') converge to something orthogonal to \mathbf{v}_i .
- ▶ Caveat: don't have \mathbf{v}_i and λ_i *exactly*, but only up to some small error, e.g.,

$$\|\hat{\mathbf{v}}_i - \mathbf{v}_i\|_2 \leq \varepsilon, \quad |\hat{\lambda}_i - \lambda_i| \leq \varepsilon'.$$

20

Error analysis

21

Nearly odecos tensors

- ▶ Suppose we have $\hat{\mathbf{T}} = \mathbf{T} + \mathbf{E}$ for some odecos $\mathbf{T} = \sum_{i=1}^k \lambda_i \mathbf{v}_i^{\otimes p}$ and (symmetric) “error tensor” \mathbf{E} with $\|\mathbf{E}\|_2 \leq \epsilon$, i.e.,

$$\max_{\mathbf{u} \in S^{k-1}} |\mathbf{E}(\mathbf{u}, \mathbf{u}, \dots, \mathbf{u})| \leq \epsilon.$$

- ▶ Matrix case ($p = 2$): ($\lambda_1 \geq \lambda_2 \geq \dots$)
 - ▶ Top eigenvalue/eigenvector ($\hat{\lambda}, \hat{\mathbf{v}}$) of $\hat{\mathbf{T}}$.
 - ▶ $\hat{\lambda}$ approximates λ_1 :
$$|\hat{\lambda} - \lambda_1| \leq \epsilon.$$
 - ▶ But need $\epsilon < \lambda_1 - \lambda_2$ for $\hat{\mathbf{v}}$ to approximate \mathbf{v}_1 (Davis-Kahan).

22

Nearly odecos tensors ($p \geq 3$)

- ▶ Higher-order case ($p \geq 3$):

- ▶ Maximum of $f_{\hat{\mathbf{T}}}$ approximates some λ_i , i.e.,

$$\left| \max_{\mathbf{u} \in S^{k-1}} \hat{\mathbf{T}}(\mathbf{u}, \mathbf{u}, \dots, \mathbf{u}) - \lambda_i \right| \leq \epsilon.$$

- ▶ Maximizers $\hat{\mathbf{v}}$ of $f_{\hat{\mathbf{T}}}$ also approximate some λ_i , i.e.,

$$\|\hat{\mathbf{v}} - \mathbf{v}_i\|_2 \leq O\left(\frac{\epsilon}{\lambda_i} + \left(\frac{\epsilon}{\lambda_i}\right)^2\right).$$

- ▶ Output of power method: depends on initialization

$$\mathbf{x}^{(0)} = c_1 \mathbf{v}_1 + c_2 \mathbf{v}_2 + \dots + c_k \mathbf{v}_k.$$

- ▶ E.g., if all $\lambda_i \in [\Omega(1), O(1)]$, then need $\max_i c_i^2 \gg \epsilon$ to get

$$|\hat{\lambda} - \lambda_i| \leq O(\epsilon), \quad \|\hat{\mathbf{v}} - \mathbf{v}_i\|_2 \leq O(\epsilon)$$

for some component i , after $O(\log(k) + \log \log(1/\epsilon))$ iterations.

23

Error from deflation

- ▶ Since $(\hat{\lambda}, \hat{\mathbf{v}})$ obtained from $\hat{\mathbf{T}} = \mathbf{T} + \mathbf{E}$ is not exactly $(\lambda_i, \mathbf{v}_i)$ for any component i of \mathbf{T} , “deflation” introduces some error:

$$\begin{aligned} \hat{\mathbf{T}}' &:= \hat{\mathbf{T}} - \hat{\lambda} \hat{\mathbf{v}}^{\otimes p} \\ &= \sum_{j=1}^k \lambda_j \mathbf{v}_j^{\otimes p} + \mathbf{E} - \hat{\lambda} \hat{\mathbf{v}}^{\otimes p} \\ &= \sum_{j \neq i} \lambda_j \mathbf{v}_j^{\otimes p} + \mathbf{E} + (\lambda_i \mathbf{v}_i^{\otimes p} - \hat{\lambda} \hat{\mathbf{v}}^{\otimes p}) \\ &=: \mathbf{T}' + \mathbf{E} + \mathbf{E}_i. \end{aligned}$$

- ▶ **Danger:** $\|\mathbf{E}_i\|_2$ can be as large as $\|\mathbf{E}\|_2$.
 - ▶ So “error” has doubled?

24

Analysis of deflation error

- ▶ (For simplicity, assume all $\lambda_i = 1$.)
- ▶ Deflation error: $\mathbf{E}_i = \mathbf{v}_i^{\otimes p} - \hat{\mathbf{v}}^{\otimes p}$
 - ▶ All we know is that $\|\hat{\mathbf{v}} - \mathbf{v}_i\|_2 \leq O(\epsilon)$.
- ▶ Consider a unit vector \mathbf{u} orthogonal to \mathbf{v}_i :

$$\begin{aligned}\|\phi_{\mathbf{E}_i}(\mathbf{u})\|_2 &= \|\langle \mathbf{u}, \mathbf{v}_i \rangle^{p-1} \mathbf{v}_i - \langle \mathbf{u}, \hat{\mathbf{v}} \rangle^{p-1} \hat{\mathbf{v}}\|_2 \\ &= \|\langle \mathbf{u}, \hat{\mathbf{v}} \rangle^{p-1} \hat{\mathbf{v}}\|_2 \\ &= |\langle \mathbf{u}, \hat{\mathbf{v}} - \mathbf{v}_i \rangle|^{p-1} \\ &\leq O(\epsilon^{p-1}).\end{aligned}$$

- ▶ Therefore, for such \mathbf{u} ,

$$\|\phi_{\mathbf{E} + \mathbf{E}_i}(\mathbf{u})\|_2 \leq (1 + O(\epsilon^{p-2}))\epsilon.$$

- ▶ When $p \geq 3$, errors due to deflation have **lower-order effect** on ability to approximate remaining components.
 - ▶ Not true for $p = 2$.

25

Recap

- ▶ High-order moments
 - ▶ Get parameter identifiability.
 - ▶ But for very high-order moments, estimation may be difficult.
- ▶ Higher-than-order-2 moments in high-dimensions
 - ▶ Can still get parameter identifiability in many cases.
 - ▶ Arrangement in higher-order tensor facilitates reasoning/computation.
- ▶ Higher-than-order-2 tensors
 - ▶ Most computational problems (that were easy for matrices) become hard.
 - ▶ But when there is a lot of structure, some computational issues are better than in matrix case!

26