

Decompositions of moment tensors

Daniel Hsu

COMS 4772

1

Tensor decompositions

2

High-dimensional support recovery from moments

- ▶ Random vector \mathbf{X} , supported on k distinct points $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_k \in \mathbb{R}^d$,

$$\mathbb{P}(\mathbf{X} = \mathbf{z}_t) = w_t > 0.$$

- ▶ Can we learn the parameters $\{(w_t, \mathbf{z}_t)\}_{t=1}^k$ from moments?
- ▶ **Moments:**

$$\begin{aligned}\mathbb{E}(\mathbf{X}) &= \sum_{t=1}^k w_t \cdot \mathbf{z}_t, \\ \mathbb{E}(\mathbf{X} \otimes \mathbf{X}) &= \sum_{t=1}^k w_t \cdot \mathbf{z}_t \otimes \mathbf{z}_t, \\ \mathbb{E}(\mathbf{X} \otimes \mathbf{X} \otimes \mathbf{X}) &= \sum_{t=1}^k w_t \cdot \mathbf{z}_t \otimes \mathbf{z}_t \otimes \mathbf{z}_t.\end{aligned}$$

3

Insufficiency of low-order moments

- ▶ The following looks like eigenvalue decomposition:

$$\mathbb{E}(\mathbf{X} \otimes \mathbf{X}) = \sum_{t=1}^k w_t \cdot \mathbf{z}_t \otimes \mathbf{z}_t,$$

but $\{\mathbf{z}_t\}_{t=1}^k$ need not be orthogonal.

- ▶ Possible to have

$$\sum_{t=1}^k w_t \cdot \mathbf{z}_t \otimes \mathbf{z}_t = \sum_{t=1}^k \tilde{w}_t \cdot \tilde{\mathbf{z}}_t \otimes \tilde{\mathbf{z}}_t$$

for different parameters $\{(w_t, \mathbf{z}_t)\}_{t=1}^k$ and $\{(\tilde{w}_t, \tilde{\mathbf{z}}_t)\}_{t=1}^k$.

- ▶ **Note:** additional constraints *could* make decomposition unique (e.g., separability in NMF).

4

Jennrich's algorithm (1970)

- ▶ Define

$$\mathbf{S} := \mathbb{E}(\mathbf{X} \otimes \mathbf{X}) = \sum_{t=1}^k w_t \cdot \mathbf{z}_t \otimes \mathbf{z}_t,$$

$$\mathbf{T} := \mathbb{E}(\mathbf{X} \otimes \mathbf{X} \otimes \mathbf{X}) = \sum_{t=1}^k w_t \cdot \mathbf{z}_t \otimes \mathbf{z}_t \otimes \mathbf{z}_t.$$

- ▶ **Key assumption:** $\mathbf{Z} := [\mathbf{z}_1 | \mathbf{z}_2 | \dots | \mathbf{z}_k]$ has rank k .
- ▶ **Main idea:** combine \mathbf{S} and \mathbf{T} to form a diagonalizable matrix whose eigenvectors are the parameters.
- ▶ Can write

$$\mathbf{S} = \mathbf{Z}\mathbf{W}\mathbf{Z}^\top$$

where $\mathbf{W} = \text{diag}(w_1, w_2, \dots, w_k)$.

- ▶ What about \mathbf{T} ?

5

“Flattening” a tensor to a matrix

- ▶ Recall multilinear function $\mathbf{T}(\mathbf{u}, \mathbf{v}, \mathbf{w}) = \sum_{i,j,k} T_{i,j,k} u_i v_j w_k$.
- ▶ Can also think of $\mathbf{T}: \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$: (j, k) -th entry of $\mathbf{T}(\mathbf{u})$ is

$$\mathbf{T}(\mathbf{u})_{j,k} = \mathbf{T}(\mathbf{u}, \mathbf{e}_j, \mathbf{e}_k).$$

(Like “currying” in functional programming.)

- ▶ Since $\mathbf{T} = \sum_{t=1}^k w_t \cdot \mathbf{z}_t \otimes \mathbf{z}_t \otimes \mathbf{z}_t$,

$$\begin{aligned} \mathbf{T}(\mathbf{u}) &= \sum_{t=1}^k w_t \langle \mathbf{z}_t, \mathbf{u} \rangle \cdot \mathbf{z}_t \otimes \mathbf{z}_t \\ &= \mathbf{Z}\mathbf{D}_\mathbf{u}\mathbf{W}\mathbf{Z}^\top \end{aligned}$$

where $\mathbf{D}_\mathbf{u} = \text{diag}(\langle \mathbf{z}_1, \mathbf{u} \rangle, \langle \mathbf{z}_2, \mathbf{u} \rangle, \dots, \langle \mathbf{z}_k, \mathbf{u} \rangle)$.

6

Combining the second and third moments

- ▶ Notation: \mathbf{A}^\dagger is Moore-Penrose pseudoinverse of \mathbf{A} .
- ▶ Pick vector \mathbf{u} somehow, and form $\mathbf{T}(\mathbf{u})\mathbf{S}^\dagger$.
 - ▶ $\mathbf{T}(\mathbf{u})\mathbf{S}^\dagger = (\mathbf{Z}\mathbf{D}_u\mathbf{W}\mathbf{Z}^\top)(\mathbf{Z}\mathbf{W}\mathbf{Z}^\top)^\dagger = \mathbf{Z}\mathbf{D}_u\mathbf{Z}^\dagger$
- ▶ Extract eigenvectors of $\mathbf{T}(\mathbf{u})\mathbf{S}^\dagger$ with non-zero eigenvalues.
 - ▶ If $\langle \mathbf{z}_t, \mathbf{u} \rangle \neq 0$, then \mathbf{z}_t is eigenvector with non-zero eigenvalue.
 - ▶ **Caveat:** can only get \mathbf{z}_t up to scaling $\sigma_t \mathbf{z}_t$ for $\sigma_t \neq 0$; but corresponding eigenvalue is $\langle \mathbf{z}_t, \mathbf{u} \rangle$, so

$$\sigma_t = \frac{\langle \sigma_t \mathbf{z}_t, \mathbf{u} \rangle}{\langle \mathbf{z}_t, \mathbf{u} \rangle}.$$

- ▶ Eigendecomposition is unique (up to scaling factors) as long as eigenvalues are distinct.
 - ▶ This holds if, e.g., \mathbf{u} chosen uniformly at random from S^{d-1} .

7

Examples

8

Simple topic model

- ▶ Bag-of-words model for documents.
 - ▶ Document: seq. of tokens $X^{(1)}, X^{(2)}, \dots$ from $\{1, 2, \dots, d\}$.
 - ▶ Model parameters: $\{(w_t, \boldsymbol{\mu}_t)\}_{t=1}^k$, where $(w_1, w_2, \dots, w_k) \in \Delta^{k-1}$ and $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_k \in \Delta^{d-1}$.
- ▶ Generative process:
 - ▶ Pick topic Y from $\{1, 2, \dots, k\}$ with

$$\mathbb{P}(Y = t) = w_t.$$

- ▶ Given Y , draw tokens iid from distribution $\boldsymbol{\mu}_Y$.
- ▶ One-hot encode $X^{(i)} = j$ as vector $\mathbf{X}^{(i)} = \mathbf{e}_j \in \mathbb{R}^d$.
- ▶ **Claim:**

$$\mathbb{E}(\mathbf{X}^{(1)} \otimes \mathbf{X}^{(2)}) = \sum_{t=1}^k w_t \cdot \boldsymbol{\mu}_t^{\otimes 2},$$

$$\mathbb{E}(\mathbf{X}^{(1)} \otimes \mathbf{X}^{(2)} \otimes \mathbf{X}^{(3)}) = \sum_{t=1}^k w_t \cdot \boldsymbol{\mu}_t^{\otimes 3}.$$

9

Simple topic model: third-order moments

$$\mathbf{T} = \mathbb{E}(\mathbf{X}^{(1)} \otimes \mathbf{X}^{(2)} \otimes \mathbf{X}^{(3)})$$

$$\begin{aligned} T_{i,j,k} &= \mathbb{P}(X^{(1)} = i \wedge X^{(2)} = j \wedge X^{(3)} = k) \\ &= \sum_t w_t \cdot \mathbb{P}(X^{(1)} = i \wedge X^{(2)} = j \wedge X^{(3)} = k \mid Y = t) \\ &= \sum_t w_t \cdot \mu_{t,i} \cdot \mu_{t,j} \cdot \mu_{t,k} \end{aligned}$$

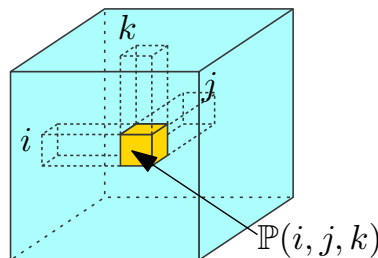


Figure 1: Third-order moments

10

Mixtures of spherical Gaussians

- ▶ Model parameters: $\{(w_t, \boldsymbol{\mu}_t, \sigma_t^2)\}_{t=1}^k$
 - ▶ $(w_1, w_2, \dots, w_k) \in \Delta^{k-1}$;
 - ▶ $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_k \in \mathbb{R}^d$;
 - ▶ $\sigma_1^2, \sigma_2^2, \dots, \sigma_k^2 > 0$.
- ▶ For simplicity, assume $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2 = \sigma^2$.
- ▶ $\mathbf{X} = \mathbf{Y} + \sigma \mathbf{Z}$, where \mathbf{Y} and \mathbf{Z} are independent, and
 - ▶ $\mathbb{P}(\mathbf{Y} = \boldsymbol{\mu}_t) = w_t$ for each $t \in \{1, 2, \dots, k\}$,
 - ▶ $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

11

Mixtures of spherical Gaussians: moments

- ▶ $\mathbb{E}(\mathbf{X}) = \mathbb{E}(\mathbf{Y}) = \sum_{t=1}^k w_t \cdot \boldsymbol{\mu}_t$
- ▶ $\mathbb{E}(\mathbf{X}^{\otimes 2}) = \sum_{t=1}^k w_t \cdot \boldsymbol{\mu}_t^{\otimes 2} + \sigma^2 \mathbf{I}$
 - ▶ How to estimate σ^2 ? Can look at smallest eigenvalue of centered second-order moment (i.e., covariance matrix).
- ▶ $\mathbb{E}(\mathbf{X}^{\otimes 3}) = \sum_{t=1}^k w_t \cdot \boldsymbol{\mu}_t^{\otimes 3} + f(\sigma^2, \mathbb{E}(\mathbf{X}))$
- ▶ **Upshot:** can estimate $\sum_{t=1}^k w_t \cdot \boldsymbol{\mu}_t^{\otimes 2}$ and $\sum_{t=1}^k w_t \cdot \boldsymbol{\mu}_t^{\otimes 3}$ from observable quantities.

12

Poisson topic model (following Canny, 2004)

- ▶ Bag-of-paragraphs model for documents.
 - ▶ Document: sequence of paragraphs (each a bag-of-words) $\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots$
- ▶ Generative process:
 - ▶ Pick topic affinities $\mathbf{H} = (H_1, H_2, \dots, H_k)$ from some non-Gaussian product distribution on $\mathbb{R}_{\geq 0}^k$.
 - ▶ Given \mathbf{H} , draw iid paragraph word count vectors $\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots$ where

$$\mathbf{X}^{(i)} \mid \mathbf{H} \sim \bigotimes_{j=1}^d \text{Poi}(\lambda_j)$$

and

$$\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \dots, \lambda_d) = \sum_{t=1}^k H_t \boldsymbol{\mu}_t.$$

- ▶ (Similar to Independent Component Analysis and Latent Dirichlet Allocation.)

13

Poisson topic model: moments

- ▶ First moment: $\boldsymbol{\mu} := \mathbb{E}(\mathbf{X}^{(1)}) = \sum_{t=1}^k \mathbb{E}(H_t) \cdot \boldsymbol{\mu}_t$.
- ▶ Second moment:

$$\begin{aligned} & \mathbb{E}[(\mathbf{X}^{(1)} - \boldsymbol{\mu}) \otimes (\mathbf{X}^{(2)} - \boldsymbol{\mu})] \\ &= \mathbb{E} \left[\mathbb{E}[(\mathbf{X}^{(1)} - \boldsymbol{\mu}) \otimes (\mathbf{X}^{(2)} - \boldsymbol{\mu}) \mid \mathbf{H}] \right] \\ &= \mathbb{E} \left[\left(\sum_{t=1}^k (H_t - \mathbb{E}(H_t)) \cdot \boldsymbol{\mu}_t \right)^{\otimes 2} \right] \\ &= \sum_{t=1}^k \text{var}(H_t) \cdot \boldsymbol{\mu}_t^{\otimes 2}. \end{aligned}$$

- ▶ Third moment:

$$\mathbb{E}[(\mathbf{X}^{(1)} - \boldsymbol{\mu}) \otimes (\mathbf{X}^{(2)} - \boldsymbol{\mu}) \otimes (\mathbf{X}^{(3)} - \boldsymbol{\mu})] = \sum_{t=1}^k \text{skew}(H_t) \cdot \boldsymbol{\mu}_t^{\otimes 3}.$$

14

Multiview models

- ▶ Both topic models are examples of *multiview models*.
- ▶ Observables include multiple “views” $\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \mathbf{X}^{(3)}, \dots$ that are conditionally independent given a hidden state H .
- ▶ **Example: hidden Markov model**
 - ▶ Hidden state sequence forms Markov chain (on finite state space): $H^{(1)} \rightarrow H^{(2)} \rightarrow H^{(3)} \rightarrow \dots$
 - ▶ Observation sequence $X^{(1)}, X^{(2)}, X^{(3)}, \dots$ where $X^{(t)}$ is independent of all other variables conditional on $H^{(t)}$.
 - ▶ **Multiview structure:** observables $X^{(1)}, X^{(2)}, X^{(3)}$ are conditionally independent given $H^{(t)}$.

15

Multiview moments

- ▶ View- i conditional means $\{\mu_t^{(i)}\}_{t=1}^k$,

$$\mu_t^{(i)} := \mathbb{E}[\mathbf{X}^{(i)} \mid H = t].$$

- ▶ Moments:

$$\mathbb{E}[\mathbf{X}^{(1)} \otimes \mathbf{X}^{(2)} \otimes \mathbf{X}^{(3)}] = \sum_{t=1}^k w_t \cdot \mu_t^{(1)} \otimes \mu_t^{(2)} \otimes \mu_t^{(3)}$$

where $w_t = \mathbb{P}(H = t)$.

- ▶ Not symmetric
 - ▶ Can be made symmetric using second-order moments.
 - ▶ Or, use asymmetric version of Jennrich’s algorithm.

16