

Stochastic gradient descent in high dimensions for multi-spiked tensor PCA

Program “Random tensors and related topics”, IHP Paris

VANESSA PICCOLO (ENS de Lyon)

Based on joint work with GÉRARD BEN AROUS (Courant Institute, NYU) and CÉDRIC GERBELOT (Courant Institute, NYU)

October, 2 2024

Multi-spiked tensor PCA problem

Multi-spiked tensor PCA problem

Definition (Spiked tensor model [Johnstone 2001; Richard, Montanari 2014])

Observe a p -tensor $\mathbf{Y} \in (\mathbb{R}^N)^{\otimes p}$ of the form

$$\mathbf{Y} = \sum_{i=1}^r \sqrt{N} \lambda_i \mathbf{v}_i^{\otimes p} + \mathbf{W}, \quad (1)$$

where

- $p \geq 2$ and r is fixed,
- $\mathbf{W} \in (\mathbb{R}^N)^{\otimes p}$ is a p -tensor with i.i.d. Gaussian entries $W_{i_1, \dots, i_p} \sim \mathcal{N}(0, 1)$,
- $\lambda_1 \geq \dots \geq \lambda_r \geq 0$ are the signal-to-noise ratios (SNRs),
- $\mathbf{v}_1, \dots, \mathbf{v}_r \in \mathbb{S}^{N-1}$ are unknown, orthogonal signal vectors.

Multi-spiked tensor PCA problem

Definition (Spiked tensor model [Johnstone 2001; Richard, Montanari 2014])

Observe a p -tensor $\mathbf{Y} \in (\mathbb{R}^N)^{\otimes p}$ of the form

$$\mathbf{Y} = \sum_{i=1}^r \sqrt{N} \lambda_i \mathbf{v}_i^{\otimes p} + \mathbf{W}, \quad (1)$$

where

- $p \geq 2$ and r is fixed,
- $\mathbf{W} \in (\mathbb{R}^N)^{\otimes p}$ is a p -tensor with i.i.d. Gaussian entries $W_{i_1, \dots, i_p} \sim \mathcal{N}(0, 1)$,
- $\lambda_1 \geq \dots \geq \lambda_r \geq 0$ are the signal-to-noise ratios (SNRs),
- $\mathbf{v}_1, \dots, \mathbf{v}_r \in \mathbb{S}^{N-1}$ are unknown, orthogonal signal vectors.

Goal: Given M i.i.d. samples $(\mathbf{Y}^\ell)_{\ell \leq M}$ of the form (1), estimate $\mathbf{v}_1, \dots, \mathbf{v}_r$ (with high probability as $N \rightarrow \infty$).

Multi-spiked tensor PCA problem

Definition (Spiked tensor model [Johnstone 2001; Richard, Montanari 2014])

Observe a p -tensor $\mathbf{Y} \in (\mathbb{R}^N)^{\otimes p}$ of the form

$$\mathbf{Y} = \sum_{i=1}^r \sqrt{N} \lambda_i \mathbf{v}_i^{\otimes p} + \mathbf{W}, \quad (1)$$

where

- $p \geq 2$ and r is fixed,
- $\mathbf{W} \in (\mathbb{R}^N)^{\otimes p}$ is a p -tensor with i.i.d. Gaussian entries $W_{i_1, \dots, i_p} \sim \mathcal{N}(0, 1)$,
- $\lambda_1 \geq \dots \geq \lambda_r \geq 0$ are the signal-to-noise ratios (SNRs),
- $\mathbf{v}_1, \dots, \mathbf{v}_r \in \mathbb{S}^{N-1}$ are unknown, orthogonal signal vectors.

Goal: Given M i.i.d. samples $(\mathbf{Y}^\ell)_{\ell \leq M}$ of the form (1), estimate $\mathbf{v}_1, \dots, \mathbf{v}_r$ (with high probability as $N \rightarrow \infty$).

Estimation task: Produce estimators $\mathbf{x}_1, \dots, \mathbf{x}_r$ attaining

- exact recovery of the spikes: $\mathbf{x}_i = (1 - o(1))\mathbf{v}_i$ for all $1 \leq i \leq r$,
- recovery of a permutation of the spikes: there exists a permutation $\sigma \in S_r$ such that $\mathbf{x}_i = (1 - o(1))\mathbf{v}_{\sigma(i)}$ for all $1 \leq i \leq r$,

The multi-spiked tensor PCA problem

Statistical procedure: Maximum Likelihood Estimator of $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_r]$ is given by a solution of

$$\begin{aligned} \text{minimize} \quad \mathcal{L}(\mathbf{X}; \mathbf{Y}) &= \sum_{i=1}^r \lambda_i \langle \mathbf{W}, \mathbf{x}_i^{\otimes p} \rangle - \sum_{1 \leq i, j \leq r} \sqrt{N} \lambda_i \lambda_j \langle \mathbf{v}_i, \mathbf{x}_j \rangle^{p-1} \\ \text{subject to} \quad \mathbf{X}^\top \mathbf{X} &= \mathbf{I}_r, \end{aligned} \tag{2}$$

where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_r] \in \mathbb{R}^{N \times r}$. The set $\text{St}(N, r) = \{\mathbf{X} \in \mathbb{R}^{N \times r} : \mathbf{X}^\top \mathbf{X} = \mathbf{I}_r\}$ is known as the **Stiefel manifold**.

The multi-spiked tensor PCA problem

Statistical procedure: Maximum Likelihood Estimator of $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_r]$ is given by a solution of

$$\begin{aligned} \text{minimize} \quad & \mathcal{L}(\mathbf{X}; \mathbf{Y}) = \sum_{i=1}^r \lambda_i \langle \mathbf{W}, \mathbf{x}_i^{\otimes p} \rangle - \sum_{1 \leq i, j \leq r} \sqrt{N} \lambda_i \lambda_j \langle \mathbf{v}_i, \mathbf{x}_j \rangle^{p-1} \\ \text{subject to} \quad & \mathbf{X}^\top \mathbf{X} = \mathbf{I}_r, \end{aligned} \tag{2}$$

where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_r] \in \mathbb{R}^{N \times r}$. The set $\text{St}(N, r) = \{\mathbf{X} \in \mathbb{R}^{N \times r} : \mathbf{X}^\top \mathbf{X} = \mathbf{I}_r\}$ is known as the **Stiefel manifold**.

Algorithmic approach: We need an algorithm for outputting this estimator $\hat{\mathbf{X}} = \operatorname{argmin}_{\mathbf{X} \in \text{St}(N, r)} \mathcal{L}(\mathbf{X}; \mathbf{Y})$.

The multi-spiked tensor PCA problem

Statistical procedure: Maximum Likelihood Estimator of $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_r]$ is given by a solution of

$$\begin{aligned} \text{minimize} \quad & \mathcal{L}(\mathbf{X}; \mathbf{Y}) = \sum_{i=1}^r \lambda_i \langle \mathbf{W}, \mathbf{x}_i^{\otimes p} \rangle - \sum_{1 \leq i, j \leq r} \sqrt{N} \lambda_i \lambda_j \langle \mathbf{v}_i, \mathbf{x}_j \rangle^{p-1} \\ \text{subject to} \quad & \mathbf{X}^\top \mathbf{X} = \mathbf{I}_r, \end{aligned} \tag{2}$$

where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_r] \in \mathbb{R}^{N \times r}$. The set $\text{St}(N, r) = \{\mathbf{X} \in \mathbb{R}^{N \times r} : \mathbf{X}^\top \mathbf{X} = \mathbf{I}_r\}$ is known as the **Stiefel manifold**.

Algorithmic approach: We need an algorithm for outputting this estimator $\hat{\mathbf{X}} = \text{argmin}_{\mathbf{X} \in \text{St}(N, r)} \mathcal{L}(\mathbf{X}; \mathbf{Y})$.

Goal today: Understand thresholds (number of samples / steps needed) for SGD from random initializations to recover.

Online stochastic gradient descent (SGD)

Online SGD algorithm

Input: i.i.d. samples $(\mathbf{Y}^\ell)_{\ell \leq M}$, loss function $\mathcal{L}(\mathbf{X}; \mathbf{Y}^\ell)$, initial guess \mathbf{X}_0 , and step size $\delta_N > 0$.

Update:

$$\mathbf{X}_t = \mathcal{R}_{\mathbf{X}_{t-1}}(-\delta_N \nabla_{\text{St}} \mathcal{L}(\mathbf{X}_{t-1}; \mathbf{Y}^t)), \quad (3)$$

where $\mathcal{R}_{\mathbf{X}}: T_{\mathbf{X}}\text{St}(N, r) \rightarrow \text{St}(N, r)$ denotes a retraction map and

$$T_{\mathbf{X}}\text{St}(N, r) = \left\{ \mathbf{V} \in \mathbb{R}^{N \times r} : \mathbf{X}^\top \mathbf{V} + \mathbf{V}^\top \mathbf{X} = 0 \right\}$$

denotes the tangent space at $\mathbf{X} \in \text{St}(N, r)$. Here, we choose the polar retraction defined by

$$\mathcal{R}_{\mathbf{X}}(\mathbf{U}) = (\mathbf{X} + \mathbf{U}) \left(\mathbf{I}_r + \mathbf{U}\mathbf{U}^\top \right)^{-1/2}.$$

Moreover, for a function $f: \text{St}(N, r) \rightarrow \mathbb{R}$,

$$\nabla_{\text{St}} f(\mathbf{X}) = \nabla f(\mathbf{X}) - \frac{1}{2} \mathbf{X}(\mathbf{X}^\top \nabla f(\mathbf{X}) + \nabla f(\mathbf{X})^\top \mathbf{X}).$$

Output: \mathbf{X}_M

Main result for $p \geq 3$

Theorem (Recovery of a permutation of the spikes for $p \geq 3$)

Let \mathbf{X}_0 be uniformly distributed on $\text{St}(N, r)$. Assume that $M \gg \log(N)N^{p-2}$, and consider the online SGD started from \mathbf{X}_0 with step size $\delta_N \ll \log(N)^{-1}N^{-\frac{p-1}{2}}$. Then, after M steps, there exists a permutation $\sigma_* \in S_r$ such that for all $k \in [r]$,

$$|\langle \mathbf{v}_{\sigma_*(k)}, (\mathbf{X}_M)_k \rangle| \rightarrow 1 \quad \text{in probability.}$$

Theorem (Recovery of a permutation of the spikes for $p \geq 3$)

Let \mathbf{X}_0 be uniformly distributed on $\text{St}(N, r)$. Assume that $M \gg \log(N)N^{p-2}$, and consider the online SGD started from \mathbf{X}_0 with step size $\delta_N \ll \log(N)^{-1}N^{-\frac{p-1}{2}}$. Then, after M steps, there exists a permutation $\sigma_* \in S_r$ such that for all $k \in [r]$,

$$|\langle \mathbf{v}_{\sigma_*(k)}, (\mathbf{X}_M)_k \rangle| \rightarrow 1 \quad \text{in probability.}$$

- If $\mathbf{X} \in \mathbb{R}^{N \times r}$ is a matrix with i.i.d. entries $\mathcal{N}(0, 1)$, then $\mathbf{Y} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1/2}$ is uniformly distributed on $\text{St}(N, r)$ (Chikuse 1994).

Theorem (Recovery of a permutation of the spikes for $p \geq 3$)

Let \mathbf{X}_0 be uniformly distributed on $\text{St}(N, r)$. Assume that $M \gg \log(N)N^{p-2}$, and consider the online SGD started from \mathbf{X}_0 with step size $\delta_N \ll \log(N)^{-1}N^{-\frac{p-1}{2}}$. Then, after M steps, there exists a permutation $\sigma_* \in S_r$ such that for all $k \in [r]$,

$$|\langle \mathbf{v}_{\sigma_*(k)}, (\mathbf{X}_M)_k \rangle| \rightarrow 1 \quad \text{in probability.}$$

- If $\mathbf{X} \in \mathbb{R}^{N \times r}$ is a matrix with i.i.d. entries $\mathcal{N}(0, 1)$, then $\mathbf{Y} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1/2}$ is uniformly distributed on $\text{St}(N, r)$ (Chikuse 1994).
- The permutation is determined by $\lambda_i \lambda_j \langle \mathbf{v}_i, (\mathbf{X}_0)_j \rangle^{p-2}$.

Theorem (Recovery of a permutation of the spikes for $p \geq 3$)

Let \mathbf{X}_0 be uniformly distributed on $\text{St}(N, r)$. Assume that $M \gg \log(N)N^{p-2}$, and consider the online SGD started from \mathbf{X}_0 with step size $\delta_N \ll \log(N)^{-1}N^{-\frac{p-1}{2}}$. Then, after M steps, there exists a permutation $\sigma_* \in S_r$ such that for all $k \in [r]$,

$$|\langle \mathbf{v}_{\sigma_*(k)}, (\mathbf{X}_M)_k \rangle| \rightarrow 1 \quad \text{in probability.}$$

- If $\mathbf{X} \in \mathbb{R}^{N \times r}$ is a matrix with i.i.d. entries $\mathcal{N}(0, 1)$, then $\mathbf{Y} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1/2}$ is uniformly distributed on $\text{St}(N, r)$ (Chikuse 1994).
- The permutation is determined by $\lambda_i \lambda_j \langle \mathbf{v}_i, (\mathbf{X}_0)_j \rangle^{p-2}$.
- If the SNRs $\lambda_1, \dots, \lambda_r$ are sufficiently separated, then we have **exact recovery** of the spikes.

Theorem (Recovery of a permutation of the spikes for $p \geq 3$)

Let \mathbf{X}_0 be uniformly distributed on $\text{St}(N, r)$. Assume that $M \gg \log(N)N^{p-2}$, and consider the online SGD started from \mathbf{X}_0 with step size $\delta_N \ll \log(N)^{-1}N^{-\frac{p-1}{2}}$. Then, after M steps, there exists a permutation $\sigma_* \in S_r$ such that for all $k \in [r]$,

$$|\langle \mathbf{v}_{\sigma_*(k)}, (\mathbf{X}_M)_k \rangle| \rightarrow 1 \quad \text{in probability.}$$

- If $\mathbf{X} \in \mathbb{R}^{N \times r}$ is a matrix with i.i.d. entries $\mathcal{N}(0, 1)$, then $\mathbf{Y} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1/2}$ is uniformly distributed on $\text{St}(N, r)$ (Chikuse 1994).
- The permutation is determined by $\lambda_i \lambda_j \langle \mathbf{v}_i, (\mathbf{X}_0)_j \rangle^{p-2}$.
- If the SNRs $\lambda_1, \dots, \lambda_r$ are sufficiently separated, then we have **exact recovery** of the spikes.
- Regardless of the values of the SNRs, **recovery of a permutation of the spikes is always possible**, provided a sample complexity of order $\log(N)N^{p-2}$.

Examples

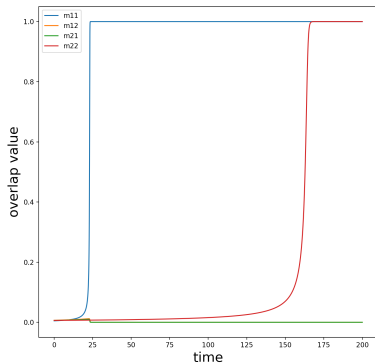
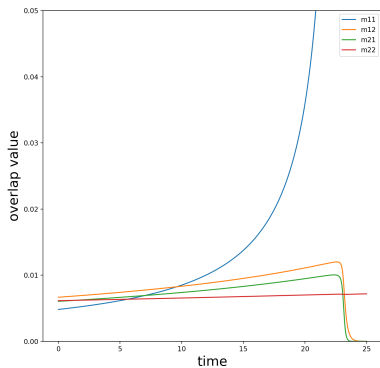


Figure 1: Evolution of the correlations $\{m_{ij} = \langle \mathbf{v}_i, \mathbf{x}_j \rangle, 1 \leq i, j \leq 2\}$ under the population dynamics for $\rho = 3$, $\lambda_1 = 3$ and $\lambda_2 = 1$.

Examples

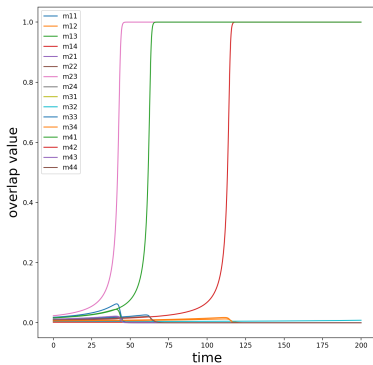
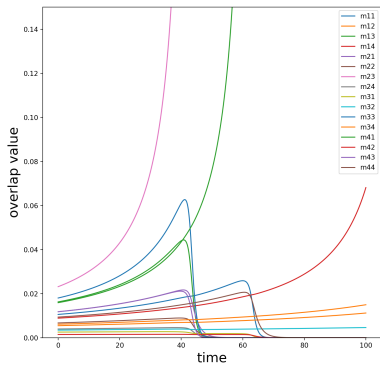


Figure 2: Evolution of the correlations $\{m_{ij} = \langle \mathbf{v}_i, \mathbf{x}_j \rangle, 1 \leq i, j \leq 4\}$ under the population dynamics for $\rho = 3, \lambda_1 = \dots = \lambda_4 = 1$.

Examples

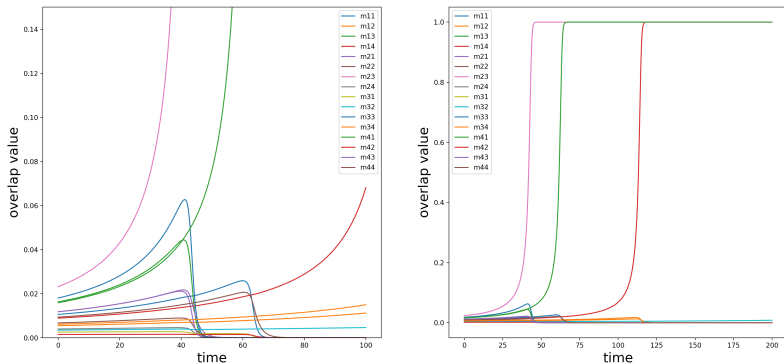


Figure 2: Evolution of the correlations $\{m_{ij} = \langle \mathbf{v}_i, \mathbf{x}_j \rangle, 1 \leq i, j \leq 4\}$ under the population dynamics for $\rho = 3, \lambda_1 = \dots = \lambda_4 = 1$.

Sequential elimination phenomenon: The correlations $\{\langle \mathbf{v}_{\sigma_*(k)}, \mathbf{x}_k \rangle\}_{k=1}^r$ increase one by one, sequentially eliminating those that share a row or column index.

Summary

- The number of samples required for online SGD from random initializations to recover scales as $\log(N)N^{p-2}$;
- For $p \geq 3$, **recovery of a permutation** of the spikes is always achievable, even when the SNRs are equal;
- The hidden vectors are recovered sequentially in a process we term **sequential elimination**: once a correlation exceeds a critical threshold, all correlations sharing a row or column index become sufficiently small, allowing the next correlation to grow and become macroscopic.