

High-dimensional optimization for the multi-spike tensor PCA problem.

Cédric Gerbelot (ENS Lyon).

Based on joint work with Gérard Ben Arous (Courant)
and Vanessa Piccolo (ENS Lyon \rightarrow EPFL).

Goal: recover r orthogonal spikes $v_1, \dots, v_r \in \mathbb{S}^{N-1}(1)$ for M noisy observations of the form

$$Y^p = W^p + \sum_{i=1}^r \sqrt{N^p} \lambda_i v_i^{\otimes p} \quad 1 \leq p \leq M, \quad r \text{ fixed, independent on } N.$$

where $W^p \in (\mathbb{R}^N)^{\otimes p}$ have i.i.d. subGaussian entries, $p \geq 2$ and $\lambda_1 \geq \dots \geq \lambda_r$ are the signal-to-noise ratios (SNRs). Assume p and r are known.

To solve this problem: gradient flow, (Langevin dynamics) and online SGD on the objective function defined by Gaussian maximum likelihood:

$$\mathcal{E}(X) = \frac{1}{M} \sum_{p=1}^M \left\| Y^p - \sum_{i=1}^r \lambda_i x_i^{\otimes p} \right\|_F^2, \quad \left(\text{if choose } \sum_{i=1}^r x_i^{\otimes p} \text{ doesn't change much} \right).$$

where $X = [x_1 | \dots | x_r] \in \mathbb{R}^{N \times r}$ is constrained to the Stiefel manifold

$$\text{St}_N(r, r) = \left\{ X \in \mathbb{R}^{N \times r} : X^T X = I_r \right\}.$$

Expanding $\mathcal{E}(X)$, we obtain

$$\mathcal{E}(X) = \underbrace{\frac{1}{M} \sum_{p=1}^M \sum_{i=1}^r \lambda_i \langle W^p, x_i^{\otimes p} \rangle}_{\text{noise part } H_0(X)} - \underbrace{\sum_{1 \leq i, j \leq r} \sqrt{N^p} \lambda_i \lambda_j m_{ij}^{p-1}(X)}_{\text{signal part } \phi(X)},$$

where

$m_{ij}^{p-1}(X) = \langle v_i, x_j \rangle$ is the correlation between v_i and x_j .

Upon appropriate control of the noise, can study autonomous, low-dimensional (r^2) dynamics on the $\{m_{ij}(X)\}_{1 \leq i, j \leq r} \Rightarrow$ effective dynamics.

Large body of recent works on single and multi-index models used as templates to understand high-dimensional, non-convex optimization. We had wonderful talks by Aukesh, summary statistics with Reza & Gerard, - - - etc - - -, talks by members of Florent's group, talk by Eshaan and collaborators, wonderful lecture by Bruno, etc - - -

In multi-index case, many works use various "oracle" modifications of dynamics and focus on achieving positive correlation with the target subspace.

Here, simple model but try to understand entire dynamics precisely: $\left(\begin{array}{l} \text{more modest than} \\ \text{two layer neural} \\ \text{network} \end{array} \right)$

- understand fixed points, no "oracle" modification:

- perfect recovery: $x_i = (1 - o(1)) v_i$
- recovery of a permutation $x_i = (1 - o(1)) v_{\sigma(i)}$, what permutation?
- recovery of the good subspace: $\text{dist}(XX^T, WW^T) = o(1)$.
- recovery of a rank deficient subspace.

- with what time and sample complexity, starting from uninformative initialization.

Background: the single spike tensor PCA problem.

Recover $v \in S^{d-1}(1)$ from noisy observations of the form

$$Y^p = W^p + \lambda \sqrt{n} v^{\otimes p} \quad 1 \leq p \leq P.$$

For $p=2$, matrix PCA [Johnstone 2002], for $p \geq 3$ introduced by [Montanari & Richard 2014]. For first order methods curvature of landscape is the main obstacle:
 $\hookrightarrow P \gg N^{P-1}$ is enough, conjectured N^{P-2} . (full batch).

- theoretical computer science & statistics : low degree polynomials, sum-of-squares ...
 what is hard to compute, with what algorithm, what method is optimal, in what sense
 [Hopkins, Shi & Steiner 2015], [Perry, Wein & Bandeira 2020]. Mention talks by Alex, Theo

- probability and mathematical physics : static and dynamical questions on high-dimensional random landscapes : complexity (# of critical points, ...), behaviour of Langevin dynamics on such landscapes, etc ... links with spin glass theory ...

[Ben Arous, Lei, Montanari and Niza '19], [Ben Arous, Jaganathan & Gheissari 2020+]
 \hookrightarrow proved $N^{1/2}$ for (full-batch) gradient flow.

- statistical physics : similar to the above using different (heuristic) methods
 [Sarao, Urbani, Kizukawa & Zdeborova 2020+] . important to mention multispike extension, harder.

\longrightarrow Here, no search for an optimal algorithm, i.e. low-degree polynomials or spectral method, again, sandbox to understand high-dimensional, non-convex optimization.

Back to the multispike tensor PCA problem.

Gradient flow:

Maybe do this part on the end.

Recall

$$\inf_{X \in \mathcal{S}_N(1,1)} H_\bullet(X) + \phi(X)$$

Gradient flow :

$$\begin{cases} \dot{X}(t) = -\nabla_{\mathcal{S}} \mathcal{E}(X(t)) \\ X(0) = X_0 \sim \mathcal{U}(\mathcal{S}_N(1,1)) \text{ invariant measure on Stiefel.} \end{cases}$$

where

$$\nabla_{\mathcal{S}} \mathcal{E}(X(t)) = \nabla \mathcal{E}(X(t)) - \frac{1}{2} X(X^T \nabla \mathcal{E}(X(t)) + X^T \nabla \mathcal{E}(X)),$$

is the orthogonal projection of the Euclidean gradient $\nabla \ell(X(t))$ on the tangent space $T_X \text{St}(1, n)$. The main focus of this talk will be on online stochastic gradient descent, but let's give a few precisions on gradient flow:

$$\dot{X}(t) = -\nabla_{\text{St}} H_0(X) - \nabla_{\text{St}} \phi(X),$$

then
$$\dot{m}_{ij}(t) = -\langle v_i, \nabla_{\text{St}} H_0(X)_j \rangle - \langle v_i, \nabla_{\text{St}} \phi(X)_j \rangle.$$

Standard approach is to bound the noise by using uniform concentration on the gradient to bound $\sup_{X \in S^{N-1}(1)} \|\nabla_{\text{St}} H_0(X)\|_2$ gives the wrong exponent: N^{P-1} . To prove

N^{P-2} for full-batch gradient flow, need time dependent bound on $\langle v_i, \nabla_{\text{St}} H_0(X)_j \rangle$: adapt "bounding flow" method of [BAGT20, 20+] to multispike case.

→ first paper with Gerard & Vanessa → Langevin & gradient flow.

→ Interesting links with DRIFT, more of interest to probabilists, happy to talk about it more after.

→ Focus now on online SGD for rest of the talk.

Online SGD:

Single sample cost function $\ell(X(t))$: $\ell(X(t))$ with $n=1$.

$$\begin{cases} X(t+1) = R_{X(t)}(-\sum_n \nabla_{\text{St}} \ell(X(t))) \\ X(0) = X_0 \sim \mathcal{U}(\text{St}(1, n)) \end{cases}$$

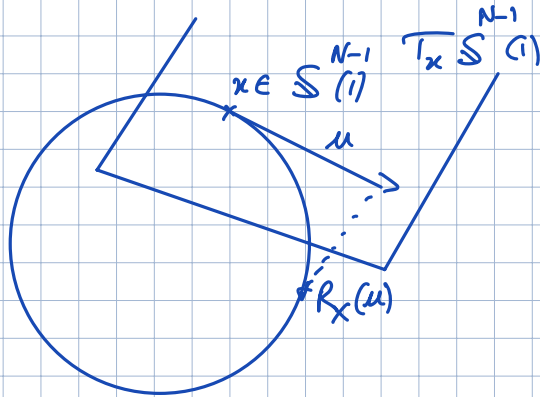
where

$$\nabla_{\text{St}} \ell(X(t)) = \nabla \ell(X(t)) - \frac{1}{2} X(X^T \nabla \ell(X(t)) + X^T \nabla \ell(X)),$$

is the orthogonal projection of the Euclidean gradient $\nabla \ell(X(t))$ on the tangent space $T_X \text{St}(1, n)$, and $R_X(U)$ is the polar retraction, i.e.

$$R_X(U) = (X+U)(I_N + UU^T)^{-1/2} \text{ for any } (X, U) \in (S_N^{+}(1, r), T_X S_N^{+}(1, r)).$$

Remark: similar to Cline's path, remember on the sphere



Project Euclidean gradient on tangent space, exit the manifold by taking a finite step, back to the manifold with the retraction. Book by N. Boumal "Optimization on smooth manifolds".

For clarity of exposition, focus on the unit quadrant where $m_{ij}^{(0)} \geq 0$ for all $1 \leq i, j \leq r$. Denote this normalized volume measure $\mathcal{U}_+(S_N^+(1, r))$.

Remark: can generate a sample from $\mathcal{U}(S_N^+(1, r))$ with

$$Y = X(X^T X)^{-1/2} \text{ and } X \text{ has i.i.d. entries.}$$

In high dimension, can roughly consider the overlap matrix

$$M_0 = V^T X_0 \text{ as i.i.d. } \mathcal{N}(0, \frac{1}{N}) \text{ random variables.}$$

To state our main result, we first need a definition:

Definition (Greedy maximum selection).

Let $A = (m_{ij}^{(p-2)}(X_0))_{1 \leq i, j \leq r} \in \mathbb{R}^{r \times r}$. We define the pairs $\{(i_k^*, j_k^*)\}_{k=1}^r$

recursively as

$$(i_k^*, j_k^*) = \operatorname{argmax}_{1 \leq i, j \leq n - (k-1)} [A^{(k-1)}]_{ij},$$

where $A^{(k-1)}$ is obtained by removing the rows i_1^*, \dots, i_{k-1}^* and the columns j_1^*, \dots, j_{k-1}^* from A . Greedy maximum selection of $A : (i_1^*, j_1^*) \dots (i_n^*, j_n^*)$.

For $p \geq 3$, we then have the following result: (will do $p=2$ if time permits).

Theorem: $X_0 \sim \mathcal{U}_+(\delta_N^{(1)})$. If $M \gg \log(N) N^{p-2}$, the online SGD with step size $\delta_N \ll \log(N)^{-1} N^{-\frac{p-1}{2}}$ produces an estimator X_T s.t., for all $k \in [n]$

$$|m_{i_k^*, j_k^*}^{(X_T)}| \xrightarrow[N \rightarrow \infty]{P} 1.$$

Remark:

- always recover a permutation
- theorem is asymptotic, but very robust to finite size effects (finite size in paper).
- perfect recovery if the SNRs are well separated.

Sketch of proof:

Output of online SGD at time t :

$$\begin{aligned} X_t &= R_{X_{t-1}} (-\delta_N \nabla_{\theta_t} \alpha(X_{t-1}; Y^t)) \\ &= (X_{t-1} - \delta_N \nabla_{\theta_t} \alpha(X_{t-1}; Y^t)) (\text{Id} + \delta_N^2 \nabla_{\theta_t}^2 \alpha(X_{t-1}; Y^t)^T \nabla_{\theta_t} \alpha(X_{t-1}; Y^t))^{-1/2} \\ &\approx X_{t-1} - \delta_N \nabla_{\theta_t} \alpha(X_{t-1}; Y^t) + \text{higher order terms in } \delta_N. \end{aligned}$$

↪ need to be careful to control these.

Corrections evolve according to

$$\begin{aligned}
 m_{ij}(t) &\simeq \langle v_i, (x_{t-1})_j^\top \rangle - \delta_N \langle v_i, \mathbb{P}_{\text{st}} \phi(x_{t-1}; \gamma^p) \rangle \\
 m_{ij}(t) &\simeq m_{ij}(0) - \delta_N \sum_{\ell=1}^t \langle v_i, (\mathbb{P}_{\text{st}} \phi(x_{\ell-1}; \gamma^p))_j \rangle \\
 &= \underbrace{m_{ij}(0)}_{\text{init.}} - \underbrace{\delta_N \sum_{\ell=1}^t \langle v_i, \mathbb{P}_{\text{st}} H^p(x_{\ell-1})_j^\top \rangle}_{\text{martingale error term}} - \underbrace{\delta_N \sum_{\ell=1}^t \langle v_i, (\mathbb{P}_{\text{st}} \phi(x_{\ell-1}))_j \rangle}_{\text{signal}} \\
 &\quad \text{of order } \delta_N \left(\sum_{\ell=1}^t \text{Var}[\langle v_i, \mathbb{P}_{\text{st}} H^p(x_{\ell-1})_j^\top \rangle] \right)^{1/2}
 \end{aligned}$$

→ governs the sample complexity by balancing with init + signal
(proof method pioneered in Tan & Vershynin '19, ISAGT '21.)

Upon controlling the noise, can focus on population dynamics

Here, pop. dyn. reads:

$$\dot{m}_{ij}(t) = p \lambda_i \lambda_j m_{ij}^{p-1} - \overbrace{\frac{p}{2} \sum_{1 \leq k, \ell \leq r} \lambda_k m_{kj} m_{k\ell} m_{\ell j} (\lambda_j m_{kj}^{p-2} + \lambda_\ell m_{k\ell}^{p-2})}^{\text{orthogonal correction}}.$$

How do we analyze this? In similar fashion to $\mathcal{U}(S^{d-1}(1))$, one can show that for $x_0 \sim \mathcal{U}(\mathcal{B}(\tau, 1))$, we have for all $1 \leq i, j \leq r$:

$$m_{ij}(0) \simeq \frac{1}{\sqrt{N}} \quad \text{w.h.p.} \quad \left(\text{mention restriction to } \mathcal{U}_+(\mathcal{B}(1, 1)). \right)$$

Then, near initialization, we can write

$$\dot{m}_{ij}(t) \simeq p \lambda_i \lambda_j m_{ij}^{p-1}(t)$$

so that, for $p \geq 3$

$$m_{ij}(t) \simeq m_{ij}(0) \left(1 - (p-2) \lambda_i \lambda_j m_{ij}(0)^{p-2} t\right)^{-\frac{1}{p-2}}.$$

We now write

$$m_{ij}(0) = \frac{\sigma_{ij}}{\sqrt{N}} \quad \text{for some } \sigma_{ij} \text{ of order one.}$$

Then, the first hitting time of $\{m_{ij} \geq \varepsilon\}$ is given by

$$T_{\varepsilon}^{(ij)} \simeq \frac{1 - \left(\frac{\sigma_{ij}}{\sqrt{N}}\right)^{p-2}}{\lambda_i \lambda_j (p-2) \sigma_{ij}^{p-2}} N^{\frac{p-2}{2}}.$$

The key observation is to see that, for any $(i, j), (i', j')$, if $\lambda_i \lambda_j m_{ij}(0)^{p-2} > (1+\delta) \lambda_{i'} \lambda_{j'} m_{i'j'}(0)^{p-2}$ for some δ of order one, then m_{ij} reaches ε before $m_{i'j'}$ can escape its original scale, i.e.

$$m_{i'j'}(T_{\varepsilon}^{(ij)}) \leq \frac{C}{\sqrt{N}}.$$

Thus m_{ij} will trigger the correction term on $m_{i'j'}$, $m_{i'j'}$ before it can move too much and send it decreasing near zero. Since the σ_{ij} are roughly i.i.d. standard normal, can always find an ordering of the $\lambda_i \lambda_j m_{ij}(0)^{p-2}$ verifying a sufficient separation. Thus, largest $\lambda_i \lambda_j m_{ij}(0)^{p-2}$ will rise, eliminate all those sharing a line and column index, and so on and so forth. Drawing on blackboard.
+ small fluctuations.

Remark: partitioning $\frac{1}{\sqrt{N}}$, ε , $1-\varepsilon$ insufficient, here, need sequence of hitting times

$$T_{ij}^{(m)} = \inf t \geq 0 : m_{ij}(X) \geq \frac{C^m}{\sqrt{N}} \quad \text{for suitable sequence } m(N).$$

→ Sharpen control on all error terms, show stability of ordering defined by GTS.

→ adds a logarithmic factor from strong Markov + union bound.

→ show simulation, explain that there are smaller fluctuations that need to be controlled, remind presence of the noise. Simulations for two dimensions, then for more.

for $p=2$:

Theorem:

Assume $X_0 \in \mathcal{U}_+(\mathcal{S}_N(1,1))$ and $\lambda_i = \lambda_{i+1}(1+k_i)$ for $k_i > 0$ of order 1. If $N \gg \log(N)^2 N^{\frac{1}{2}(1-\frac{\lambda_1}{\lambda_i})}$, $S_N \ll \log(N)^{-1} N^{-1+\frac{\lambda_1}{2\lambda_i}}$, then for all $i \in [r]$:

$$|m_{i\bar{i}}(X_{\bar{i}})| \xrightarrow[N \rightarrow \infty]{P} 1.$$

- reach the global minimizer (mention Ruckert cost function and stable manifold).
- sequential elimination harder to show: all m_{ij} escape the scale of $\frac{1}{\sqrt{N}}$.
- more sensitive to finite size effect.
- if all λ 's are equal, monotone dynamics on eigenvalues of $G = M^T \bar{M}$, subspace recovery.

If time, show simulations for $p=2$ and give intuition on the bounding flow.