

# Taming Nonconvexity in Statistical and Reinforcement Learning



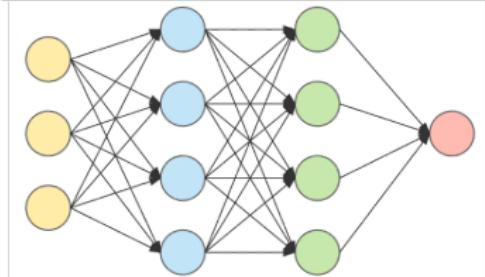
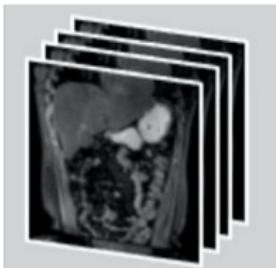
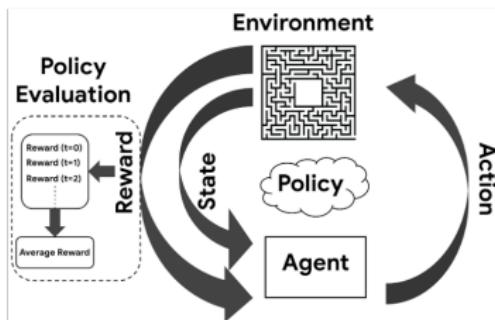
Yuxin Chen

Princeton University

# Nonconvex problems are everywhere

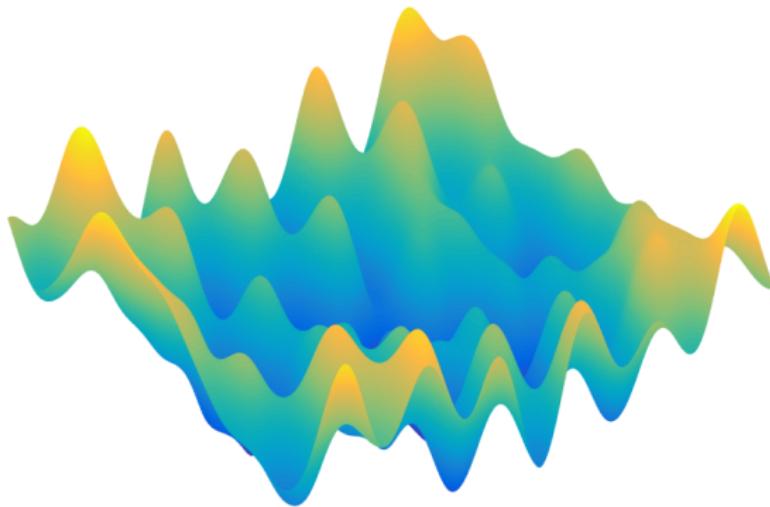
Empirical risk minimization is usually nonconvex

$$\text{minimize}_x \quad f(x; \text{data})$$



# Nonconvex optimization may be super scary

---

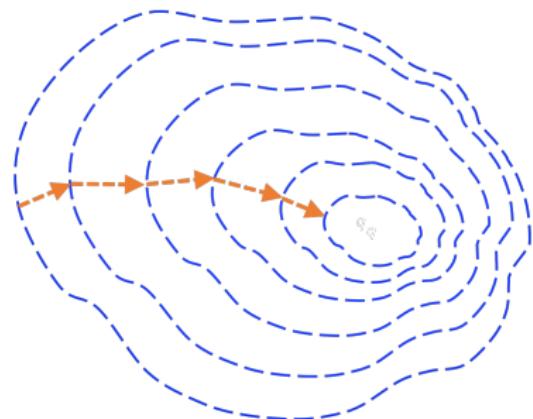
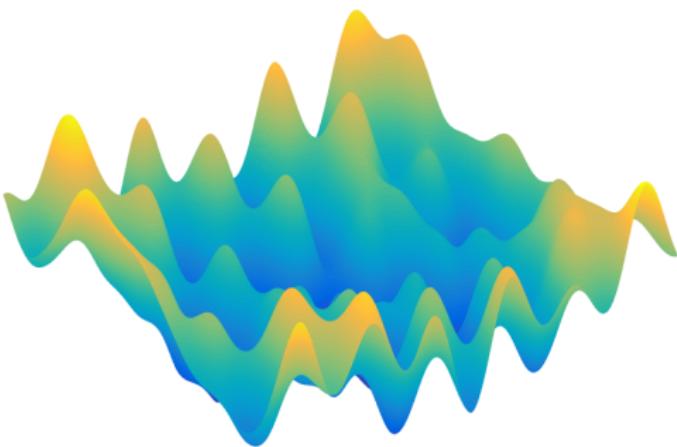


There may be bumps everywhere and exponentially many local optima

e.g. 1-layer neural net (Auer, Herbster, Warmuth '96; Vu '98)

# Nonconvex optimization may be super scary

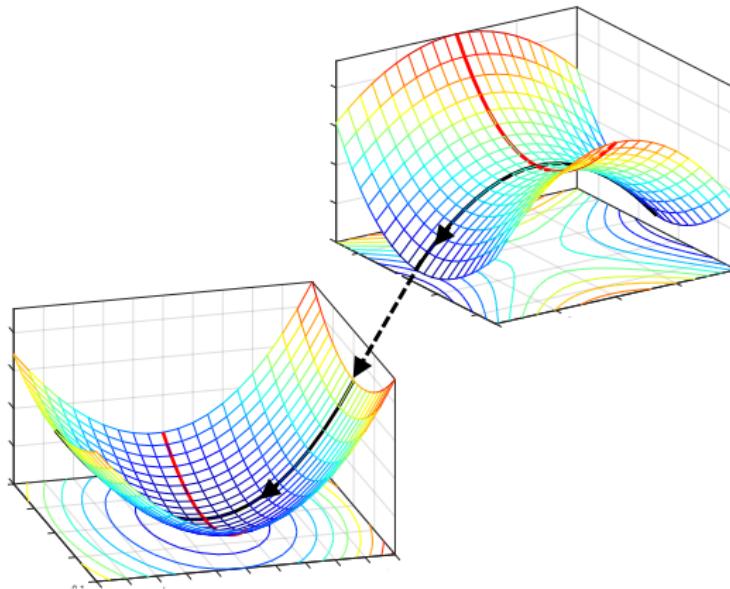
---



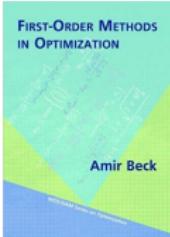
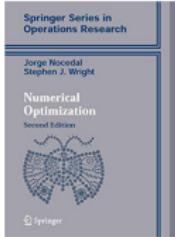
But they are solved on a daily basis via simple algorithms like  
*(stochastic) gradient descent*

# Towards demystifying nonconvex optimization

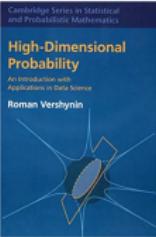
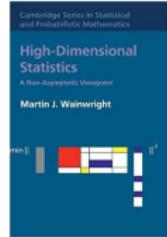
---



- Characterize optimization landscapes
- Exploit statistical tools to understand finite-sample behavior



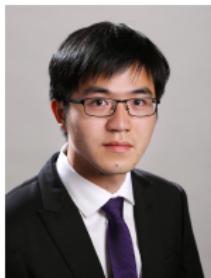
nonconvex optimization



(high-dimensional) statistics

1. Nonconvex statistical learning
  - *an efficient nonconvex algorithm for noisy tensor completion*
2. Nonconvex reinforcement learning
  - *an exponential lower bound for policy gradient methods & an efficient remedy*

# 1. Nonconvex optimization for tensor completion



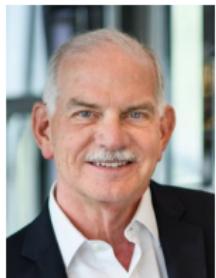
Changxiao Cai  
Princeton



Gen Li  
Tsinghua



Yuejie Chi  
CMU



H. Vincent Poor  
Princeton

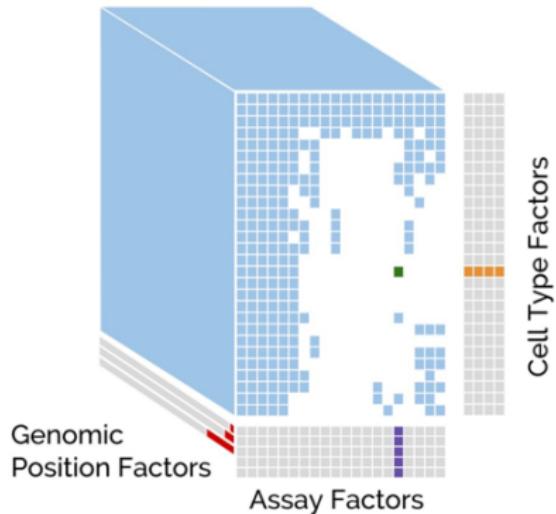
"Nonconvex low-rank tensor completion from noisy data," C. Cai, G. Li, H. Poor, Y. Chen, *Operations Research*, 2021+

"Subspace estimation from unbalanced and incomplete data matrices:  $\ell_{2,\infty}$  statistical guarantees," C. Cai, G. Li, Y. Chi, H. Poor, Y. Chen, *Annals of Statistics*, 2021+

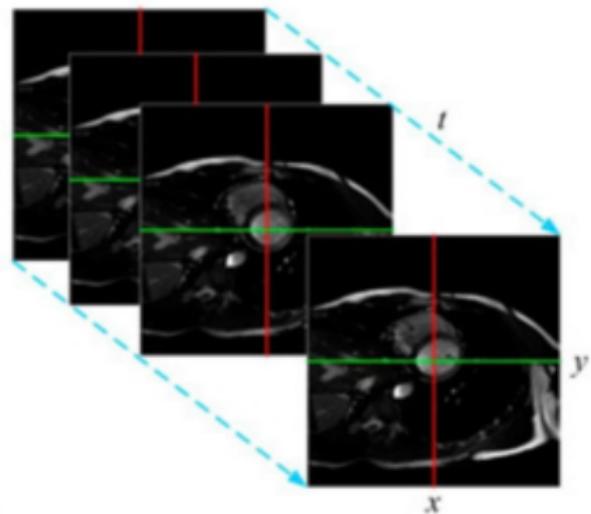
"Uncertainty quantification for nonconvex tensor completion," C. Cai, H. Poor, Y. Chen, *ICML*, 2020

# Ubiquity of high-dimensional tensor data

---



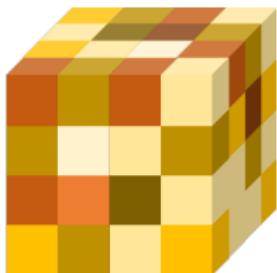
computational genomics  
— fig. credit: Schreiber et al. 19



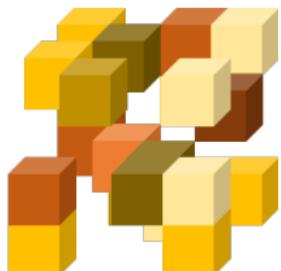
dynamic MRI  
— fig. credit: Liu et al. 17

# Imperfect data acquisition

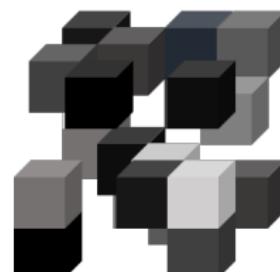
---



a tensor of interest

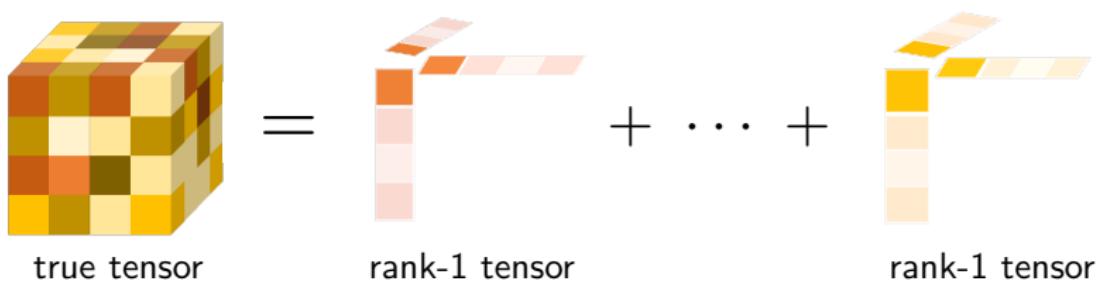


missing data



noise

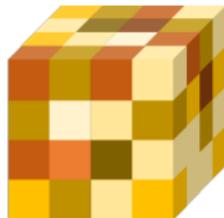
Key to enabling reliable reconstruction from incomplete data  
— exploiting **low CP-rank structure**



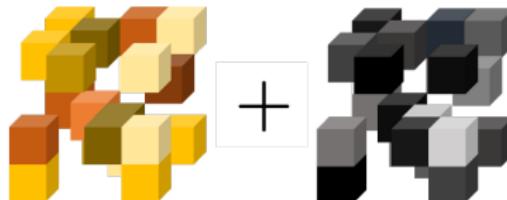
$$\mathbf{T}^* = \sum_{i=1}^r \mathbf{u}_i^* \otimes \mathbf{u}_i^* \otimes \mathbf{u}_i^*$$

# Setup

---



$T^*$



$T$

- unknown rank- $r$  tensor  $T^*$ :

$$T^* = \sum_{i=1}^r u_i^* \otimes u_i^* \otimes u_i^* \in \mathbb{R}^{d \times d \times d}$$

# Setup

---



- unknown rank- $r$  tensor  $\mathbf{T}^*$ :

$$\mathbf{T}^* = \sum_{i=1}^r \mathbf{u}_i^* \otimes \mathbf{u}_i^* \otimes \mathbf{u}_i^* \in \mathbb{R}^{d \times d \times d}$$

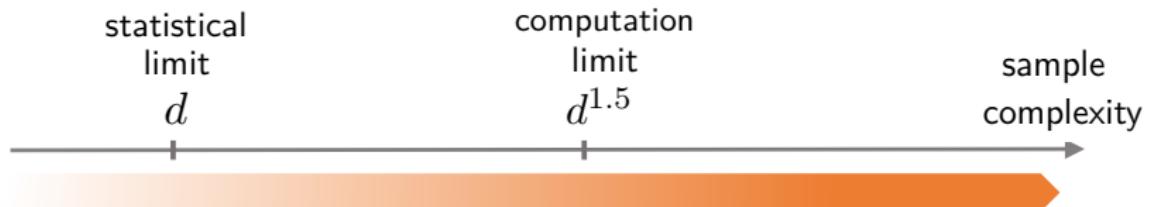
- incomplete & noisy observations over a random sampling set  $\Omega$ :

$$T_{i,j,k} = T_{i,j,k}^* + \text{noise}, \quad (i, j, k) \in \Omega$$

**Goal:** recover  $\{\mathbf{u}_i^*\}_{i=1}^r$  and  $\mathbf{T}^*$

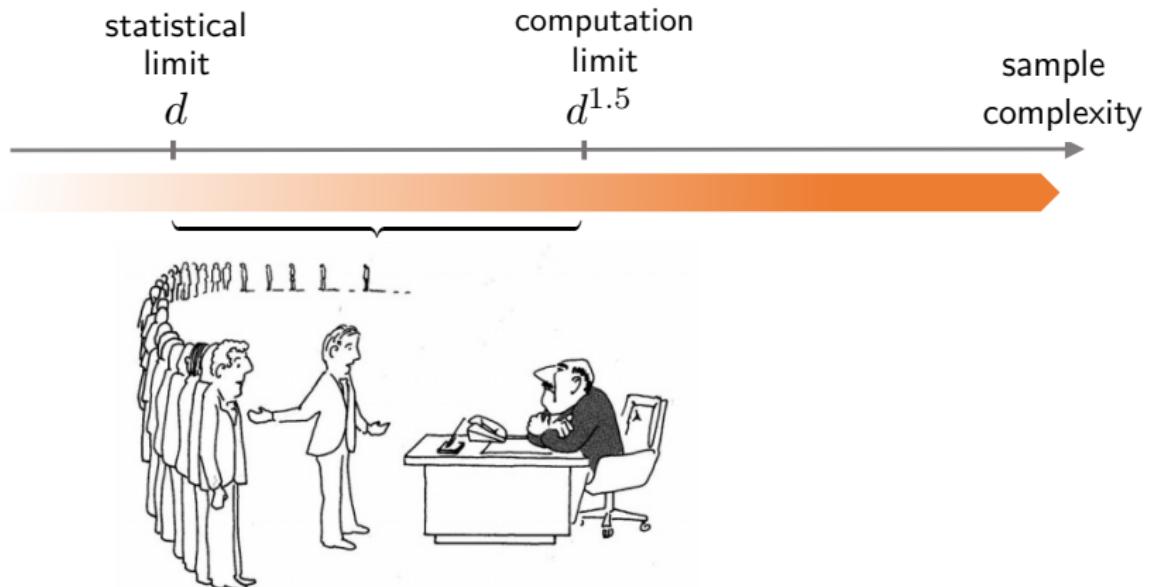
# Statistical-computational gap ( $r = O(1)$ )

---



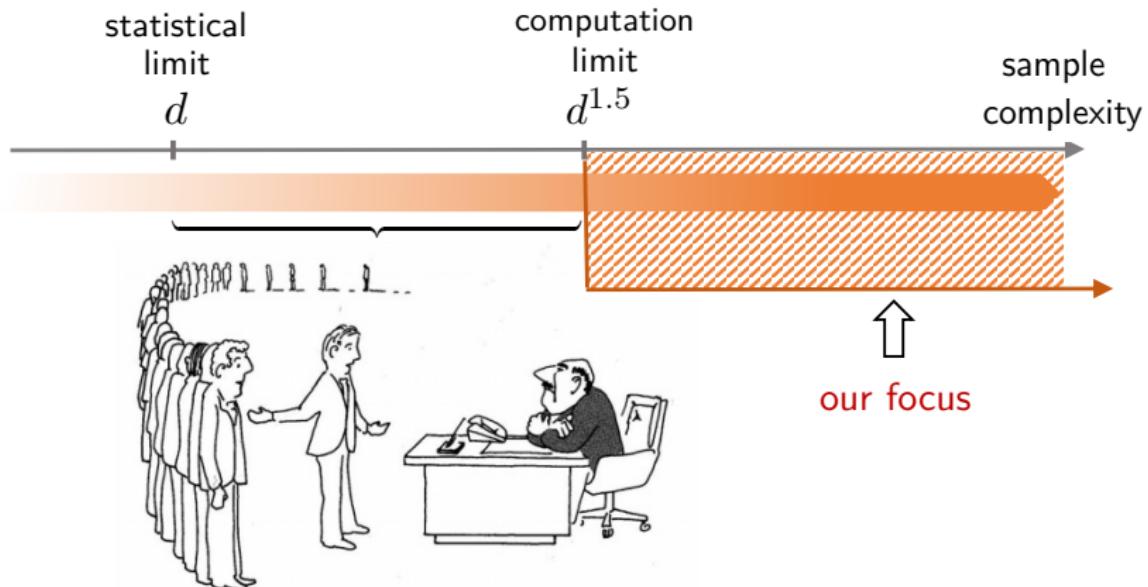
# Statistical-computational gap ( $r = O(1)$ )

---



*"I can't find an efficient algorithm, but neither can all these people."*

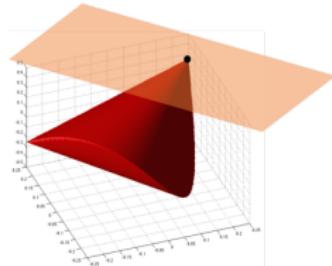
# Statistical-computational gap ( $r = O(1)$ )



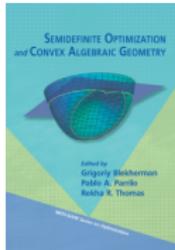
*"I can't find an efficient algorithm, but neither can all these people."*

# Prior art

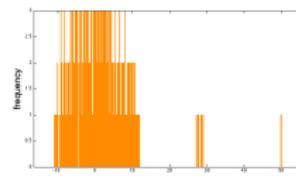
---



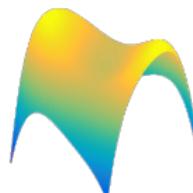
convex relaxation



sum-of-squares hierarchy



spectral methods



nonconvex optimization

- Gandy, Recht, Yamada '11
- Liu, Musalski, Wonka, Ye '12
- Kressner, Steinlechner, Vandereycken '13
- Xu, Hao, Yin, Su '13
- Romera-Paredes, Pontil '13
- Jain, Oh '14
- Huang, Mu, Goldfarb, Wright '15
- Barak, Moitra '16
- Zhang, Aeron '16
- Yuan, Zhang '16
- Montanari, Sun '16
- Kasai, Mishra '16
- Potechin, Steurer '17
- Dong, Yuan, Zhang '17
- Xia, Yuan '19
- Zhang '19
- ...

# Prior art ( $r = O(1)$ )

---



	algorithm	sample size	comput. cost	recovery type (noiseless)
Yuan, Zhang '16	tensor nuclear norm	$d$	<b>NP-hard</b>	exact
Xia, Yuan '17	spectral method + GD on manifold	$d^{3/2}$	<b>slow</b>	exact
Montanari, Sun '18	spectral method	$d^{3/2}$	$d^3$	<b>inexact</b>
Barak, Moitra '16	sum-of-squares	$d^{3/2}$	<b>slow</b> ( $d^{15}$ )	exact
Potechin et al. '17	sum-of-squares	$d^{3/2}$	<b>slow</b> ( $d^{10}$ )	exact

# Prior art ( $r = O(1)$ )

---



	algorithm	sample size	comput. cost	recovery type (noiseless)
Yuan, Zhang '16	tensor nuclear norm	$d$	<b>NP-hard</b>	exact
Xia, Yuan '17	spectral method + GD on manifold	$d^{3/2}$	<b>slow</b>	exact
Montanari, Sun '18	spectral method	$d^{3/2}$	$d^3$	<b>inexact</b>
Barak, Moitra '16	sum-of-squares	$d^{3/2}$	<b>slow</b> ( $d^{15}$ )	exact
Potechin et al. '17	sum-of-squares	$d^{3/2}$	<b>slow</b> ( $d^{10}$ )	exact

	algorithm	$\ell_2$ error (noisy)	$\ell_\infty$ error (noisy)
Xia, Yuan, Zhang '17	spectral method + tensor power method	<b>suboptimal</b>	<b>n/a</b>
Barak, Moitra '16	sum-of-squares	<b>suboptimal</b>	<b>n/a</b>

*Can we design an algorithm that is simultaneously  
sample-efficient, computationally fast, & minimax-optimal?*

## A nonconvex least squares formulation

---

$$\underset{\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_r] \in \mathbb{R}^{d \times r}}{\text{minimize}} \quad f(\mathbf{U}) := \underbrace{\sum_{(i,j,k) \in \Omega} \left\{ \left( \sum_{s=1}^r \mathbf{u}_s^{\otimes 3} \right)_{i,j,k} - T_{i,j,k} \right\}^2}_{\text{squared loss over observed entries}}$$

## A nonconvex least squares formulation

---

$$\underset{\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_r] \in \mathbb{R}^{d \times r}}{\text{minimize}} \quad f(\mathbf{U}) := \underbrace{\sum_{(i,j,k) \in \Omega} \left\{ \left( \sum_{s=1}^r \mathbf{u}_s^{\otimes 3} \right)_{i,j,k} - T_{i,j,k} \right\}^2}_{\text{squared loss over observed entries}}$$

- **pros:** statistically efficient *if we can find global solutions*

# A nonconvex least squares formulation

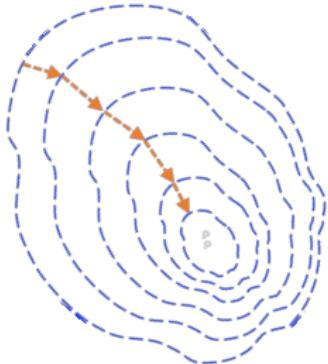
---

$$\underset{\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_r] \in \mathbb{R}^{d \times r}}{\text{minimize}} \quad f(\mathbf{U}) := \underbrace{\sum_{(i,j,k) \in \Omega} \left\{ \left( \sum_{s=1}^r \mathbf{u}_s^{\otimes 3} \right)_{i,j,k} - T_{i,j,k} \right\}^2}_{\text{squared loss over observed entries}}$$

- **pros:** statistically efficient *if we can find global solutions*
- **cons:** highly nonconvex  $\longrightarrow$  computationally challenging

# Gradient descent (GD) with random initialization?

$$\underset{\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_r] \in \mathbb{R}^{d \times r}}{\text{minimize}} \quad f(\mathbf{U}) := \sum_{(i,j,k) \in \Omega} \left\{ \left( \sum_{s=1}^r \mathbf{u}_s^{\otimes 3} \right)_{i,j,k} - T_{i,j,k} \right\}^2$$



- **initialize**  $\mathbf{U}^0$  randomly
- **gradient descent:** for  $t = 0, 1, \dots,$

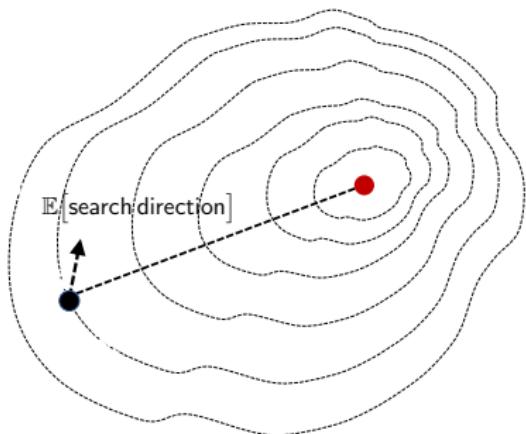
$$\mathbf{U}^{t+1} = \mathbf{U}^t - \eta_t \nabla f(\mathbf{U}^t)$$

— succeeds for phase retrieval (Chen et al. '18)

## A negative conjecture

---

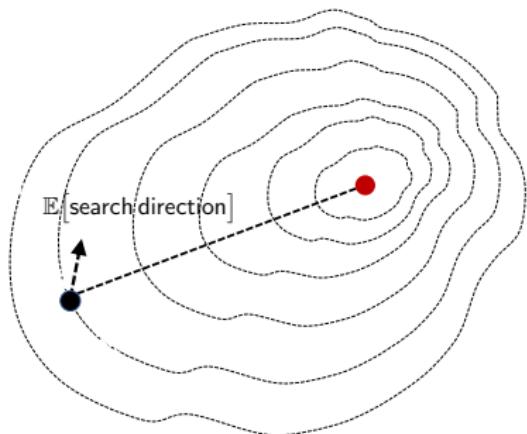
Randomly initialized GD does NOT work unless sample size  $> d^2$



## A negative conjecture

---

Randomly initialized GD does NOT work unless sample size  $> d^2$



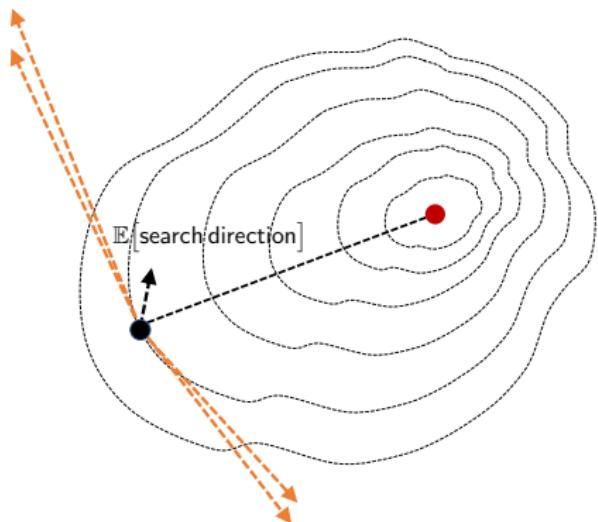
When sample size  $\asymp d^{1.5}$ :

- $\mathbb{E}[\text{search direction}]$  is desirable

# A negative conjecture

---

Randomly initialized GD does NOT work unless sample size  $> d^2$

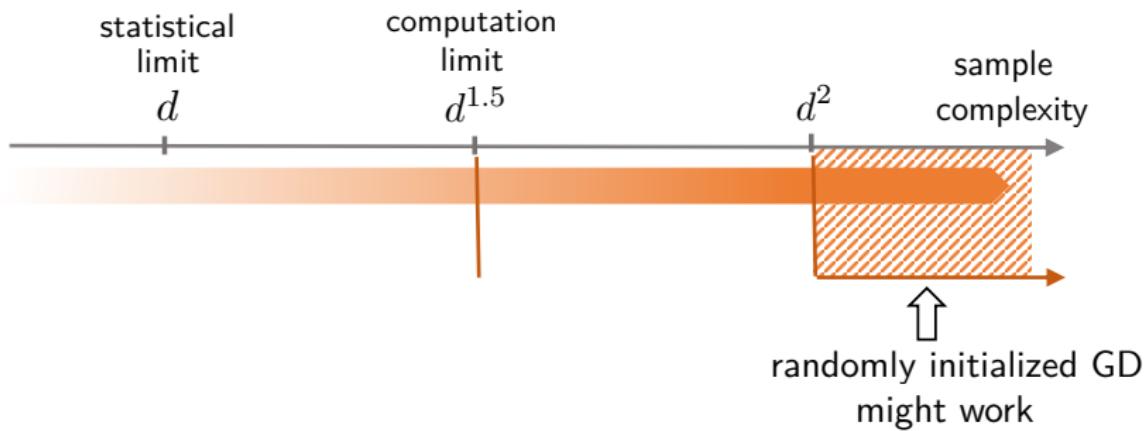


When sample size  $\asymp d^{1.5}$ :

- $E[\text{search direction}]$  is desirable
- **issue:** variance  $\gtrsim \sqrt{d} \text{ mean}^2$

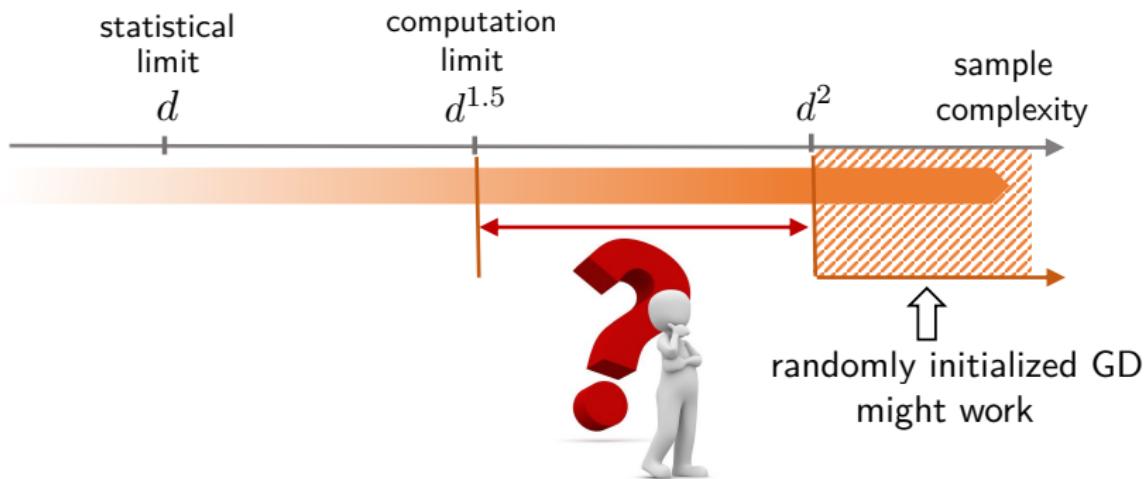
# A negative conjecture

Randomly initialized GD does NOT work unless sample size  $> d^2$

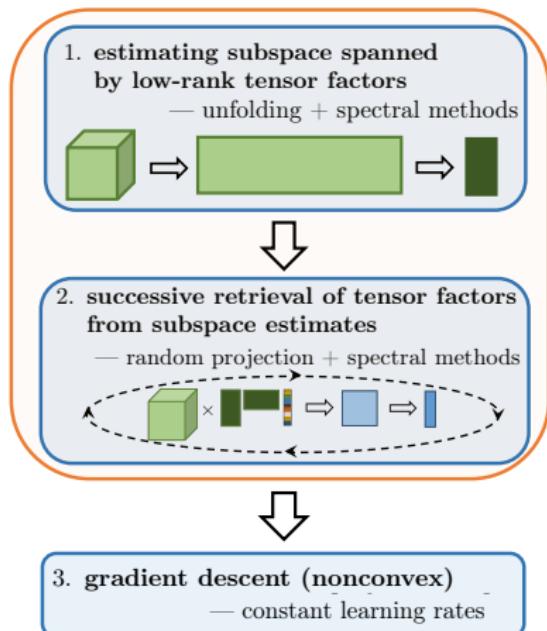


# A negative conjecture

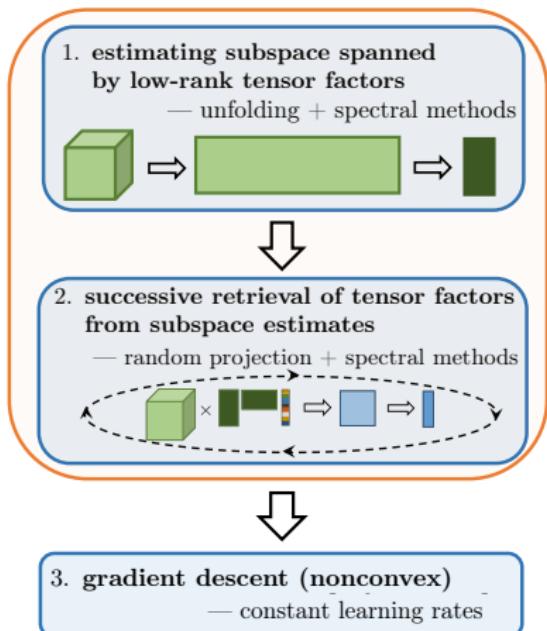
Randomly initialized GD does NOT work unless sample size  $> d^2$



# Our proposal: a two-stage nonconvex algorithm



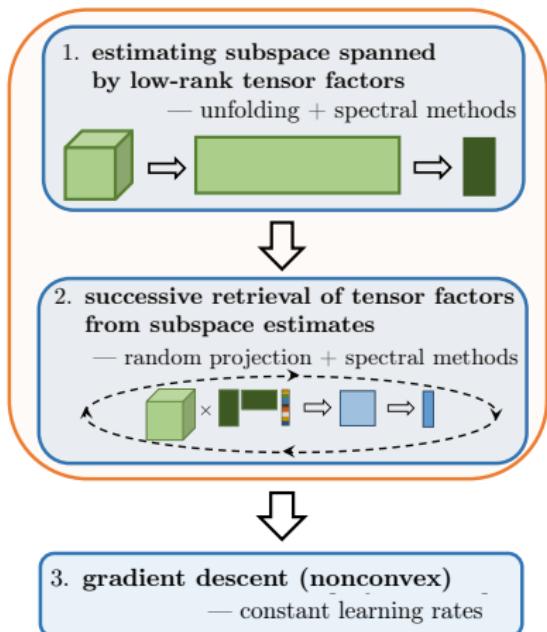
# Our proposal: a two-stage nonconvex algorithm



## 1. initialization: $U^0$

- estimate  $\text{span}\{\mathbf{u}_i^*\}$  via spectral method
- disentangle individual factors  $\{\mathbf{u}_i^*\}$  from subspace estimate

# Our proposal: a two-stage nonconvex algorithm



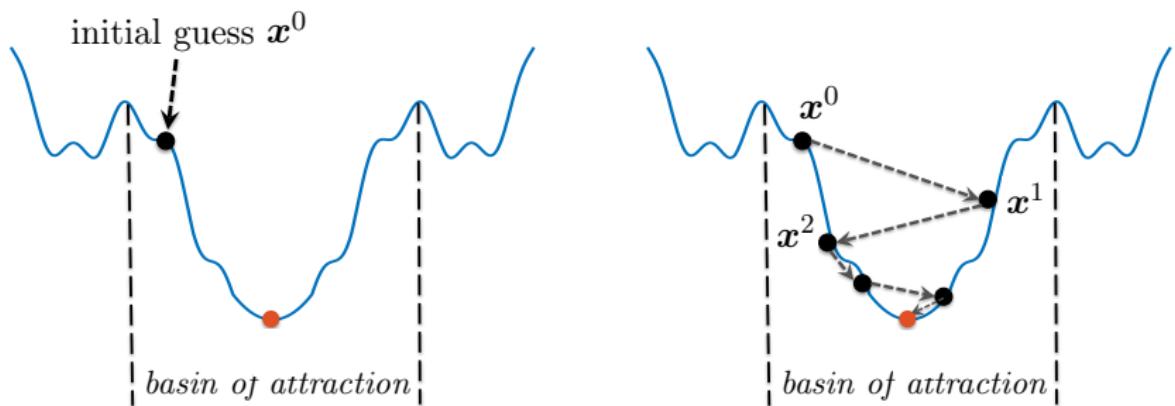
## 1. initialization: $\mathbf{U}^0$

- estimate  $\text{span}\{\mathbf{u}_i^*\}$  via spectral method
- disentangle individual factors  $\{\mathbf{u}_i^*\}$  from subspace estimate

## 2. gradient descent: for $t = 0, 1, \dots$

$$\mathbf{U}^{t+1} = \mathbf{U}^t - \eta \nabla f(\mathbf{U}^t)$$

# Rationale of two-stage approach

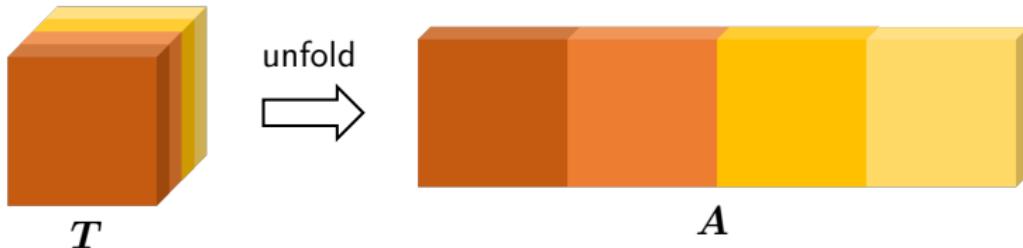


1. initialize within a local basin sufficiently close to global min  
  \underbrace{\hspace{10em}}\_{\text{(restricted) strongly convex}}
2. iterative refinement

# A bit more details about initialization

**Step 1.1:** estimate  $\text{span}\{\mathbf{u}_i^*\}_{1 \leq i \leq r} \longrightarrow \mathbf{U}_{\text{sub}}$

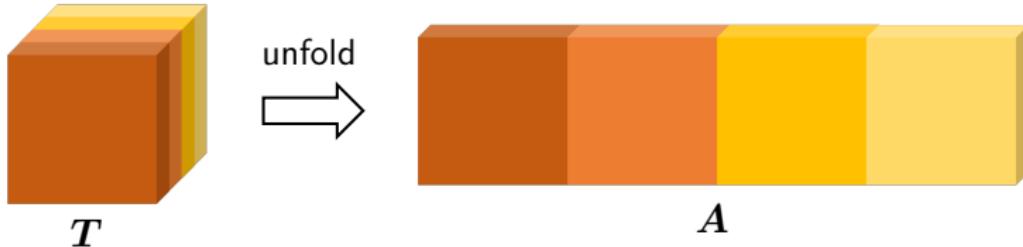
- matricization:  $\mathbf{A} = \text{unfold}(\mathbf{T})$
- estimate rank- $r$  subspace of  $\mathcal{P}_{\text{off-diag}}(\mathbf{A}\mathbf{A}^\top)$  (diagonal deletion)



# A bit more details about initialization

**Step 1.1:** estimate  $\text{span}\{\mathbf{u}_i^*\}_{1 \leq i \leq r} \longrightarrow \mathbf{U}_{\text{sub}}$

- matricization:  $\mathbf{A} = \text{unfold}(\mathbf{T})$
- estimate rank- $r$  subspace of  $\mathcal{P}_{\text{off-diag}}(\mathbf{A}\mathbf{A}^\top)$  (diagonal deletion)



**Step 1.2:** retrieve tensor factors from subspace estimate

- generate a random vector  $\mathbf{g}$  from  $\mathbf{U}_{\text{sub}}$
- compute leading eigenvector of  $\mathbf{T} \otimes \mathbf{g} = \sum_i \langle \mathbf{u}_i^*, \mathbf{g} \rangle \mathbf{u}_i^* \mathbf{u}_i^{*\top} + \text{noise}$
- repeat ...

find the  $\mathbf{u}_i^*$  most aligned with  $\mathbf{g}$

# Assumptions

---

$$\mathbf{T}^* = \sum_{i=1}^r \mathbf{u}_i^* \otimes \mathbf{u}_i^* \otimes \mathbf{u}_i^* \in \mathbb{R}^{d \times d \times d}$$

- **random sampling:** each entry is observed independently with prob.  $p$

# Assumptions

---

$$\mathbf{T}^* = \sum_{i=1}^r \mathbf{u}_i^* \otimes \mathbf{u}_i^* \otimes \mathbf{u}_i^* \in \mathbb{R}^{d \times d \times d}$$

- **random sampling:** each entry is observed independently with prob.  $p$
- **random noise:** independent zero-mean sub-Gaussian noise with variance  $O(\sigma^2)$

# Assumptions

---

$$\mathbf{T}^* = \sum_{i=1}^r \mathbf{u}_i^* \otimes \mathbf{u}_i^* \otimes \mathbf{u}_i^* \in \mathbb{R}^{d \times d \times d}$$

- **random sampling:** each entry is observed independently with prob.  $p$
- **random noise:** independent zero-mean sub-Gaussian noise with variance  $O(\sigma^2)$
- **ground truth:** low-rank ( $r = O(1)$ ), well-conditioned, incoherent ( $\{\mathbf{u}_i^*\}$  are de-localized and not aligned)

## $\ell_2$ and $\ell_\infty$ theoretical guarantees

### Theorem 1 (Cai, Li, Poor, Chen '19)

*There exists some constant  $\rho < 1$  and some permutation matrix  $\Pi \in \mathbb{R}^{r \times r}$  s.t. with high prob., the  $t$ -th iterate satisfies*

$$\|\mathbf{U}^t \Pi - \mathbf{U}^*\|_{\text{F}} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{U}^*\|_{\text{F}}$$

$$\|\mathbf{T}^t - \mathbf{T}^*\|_{\text{F}} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{T}^*\|_{\text{F}}$$

## $\ell_2$ and $\ell_\infty$ theoretical guarantees

### Theorem 1 (Cai, Li, Poor, Chen '19)

*There exists some constant  $\rho < 1$  and some permutation matrix  $\Pi \in \mathbb{R}^{r \times r}$  s.t. with high prob., the  $t$ -th iterate satisfies*

$$\|\mathbf{U}^t \Pi - \mathbf{U}^*\|_{\text{F}} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{U}^*\|_{\text{F}}$$

$$\|\mathbf{T}^t - \mathbf{T}^*\|_{\text{F}} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{T}^*\|_{\text{F}}$$

$$\|\mathbf{U}^t \Pi - \mathbf{U}^*\|_{\infty} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{U}^*\|_{\infty}$$

$$\|\mathbf{T}^t - \mathbf{T}^*\|_{\infty} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{T}^*\|_{\infty}$$

*provided that sample size  $\gtrsim d^{1.5} \text{poly log}(d)$*

## $\ell_2$ and $\ell_\infty$ theoretical guarantees

### Theorem 1 (Cai, Li, Poor, Chen '19)

There exists some constant  $\rho < 1$  and some permutation matrix  $\Pi \in \mathbb{R}^{r \times r}$  s.t. with high prob., the  $t$ -th iterate satisfies

$$\|\mathbf{U}^t \Pi - \mathbf{U}^*\|_{\text{F}} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{U}^*\|_{\text{F}}$$

$$\|\mathbf{T}^t - \mathbf{T}^*\|_{\text{F}} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{T}^*\|_{\text{F}}$$

$$\|\mathbf{U}^t \Pi - \mathbf{U}^*\|_{\infty} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{U}^*\|_{\infty}$$

$$\|\mathbf{T}^t - \mathbf{T}^*\|_{\infty} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{T}^*\|_{\infty}$$

provided that sample size  $\gtrsim d^{1.5} \text{poly log}(d)$

- linear/geometric convergence  $\longrightarrow$  linear-time algorithm

## $\ell_2$ and $\ell_\infty$ theoretical guarantees

### Theorem 1 (Cai, Li, Poor, Chen '19)

There exists some constant  $\rho < 1$  and some permutation matrix  $\Pi \in \mathbb{R}^{r \times r}$  s.t. with high prob., the  $t$ -th iterate satisfies

$$\|\mathbf{U}^t \Pi - \mathbf{U}^*\|_{\text{F}} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{U}^*\|_{\text{F}}$$

$$\|\mathbf{T}^t - \mathbf{T}^*\|_{\text{F}} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{T}^*\|_{\text{F}}$$

$$\|\mathbf{U}^t \Pi - \mathbf{U}^*\|_{\infty} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{U}^*\|_{\infty}$$

$$\|\mathbf{T}^t - \mathbf{T}^*\|_{\infty} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{T}^*\|_{\infty}$$

provided that sample size  $\gtrsim d^{1.5} \text{poly log}(d)$

- near-optimal sample complexity (among poly-time algorithms)

## $\ell_2$ and $\ell_\infty$ theoretical guarantees

### Theorem 1 (Cai, Li, Poor, Chen '19)

There exists some constant  $\rho < 1$  and some permutation matrix  $\Pi \in \mathbb{R}^{r \times r}$  s.t. with high prob., the  $t$ -th iterate satisfies

$$\|\mathbf{U}^t \Pi - \mathbf{U}^*\|_{\text{F}} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{U}^*\|_{\text{F}}$$

$$\|\mathbf{T}^t - \mathbf{T}^*\|_{\text{F}} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{T}^*\|_{\text{F}}$$

$$\|\mathbf{U}^t \Pi - \mathbf{U}^*\|_{\infty} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{U}^*\|_{\infty}$$

$$\|\mathbf{T}^t - \mathbf{T}^*\|_{\infty} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{T}^*\|_{\infty}$$

provided that sample size  $\gtrsim d^{1.5} \text{poly log}(d)$

- near-optimal statistical accuracy (both  $\ell_2$  and  $\ell_\infty$ )

## $\ell_2$ and $\ell_\infty$ theoretical guarantees

### Theorem 1 (Cai, Li, Poor, Chen '19)

There exists some constant  $\rho < 1$  and some permutation matrix  $\Pi \in \mathbb{R}^{r \times r}$  s.t. with high prob., the  $t$ -th iterate satisfies

$$\|\mathbf{U}^t \Pi - \mathbf{U}^*\|_{\text{F}} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{U}^*\|_{\text{F}}$$

$$\|\mathbf{T}^t - \mathbf{T}^*\|_{\text{F}} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{T}^*\|_{\text{F}}$$

$$\|\mathbf{U}^t \Pi - \mathbf{U}^*\|_{\infty} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{U}^*\|_{\infty}$$

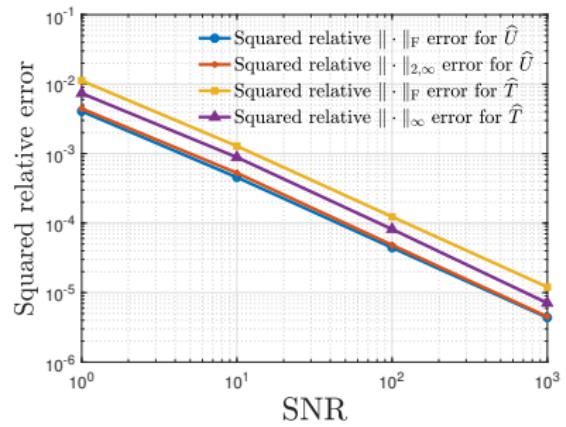
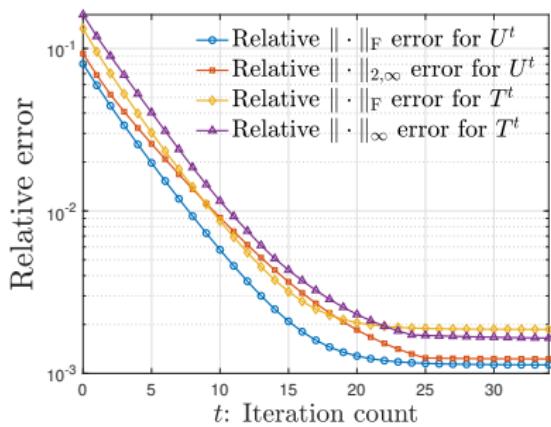
$$\|\mathbf{T}^t - \mathbf{T}^*\|_{\infty} \lesssim (\rho^t + \sigma \sqrt{d/p}) \|\mathbf{T}^*\|_{\infty}$$

provided that sample size  $\gtrsim d^{1.5} \text{poly log}(d)$

- no need of sample splitting  
→ can handle **complicated stat dependency** across iterations

# Numerical experiments

---



$$d = 100, r = 4, p = 0.1$$

## **Key proof idea: leave-one-out decoupling**

---

Leave out a small amount of randomness and re-run the algorithm

# Key proof idea: leave-one-out decoupling

Leave out a small amount of randomness and re-run the algorithm

- Stein '72
- El Karoui, Bean, Bickel, Lim, Yu '13
- El Karoui '15
- Javanmard, Montanari '15
- Zhong, Boumal '17
- Lei, Bickel, El Karoui '17
- Sur, Chen, Candès '17
- Abbe, Fan, Wang, Zhong '17
- Chen, Fan, Ma, Wang '17
- Ma, Wang, Chi, Chen '17
- Chen, Chi, Fan, Ma '18
- Ding, Chen '18
- Dong, Shi '18
- Chen, Liu, Li '19
- Chen, Fan, Ma, Yan '19
- Pananjady, Wainwright '19
- Ling '20
- Chen, Fan, Ma, Yan '20
- Agarwal, Kakade, Yang '20
- Abbe, Fan, Wang '20
- Li, Wei, Chi, Gu, Chen '20

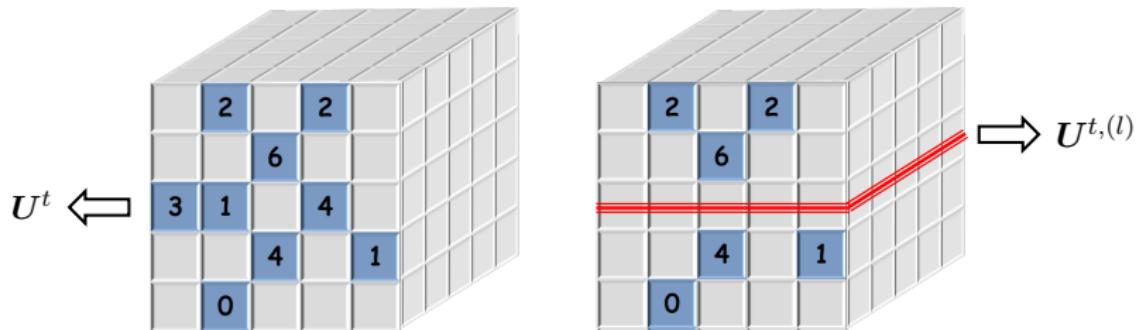
Foundations and Trends® in Machine Learning  
**Spectral Methods for Data Science: A Statistical Perspective**

Suggested Citation: Yuxin Chen, Yuejie Chi, Jianqing Fan and Cong Ma (2020), "Spectral Methods for Data Science: A Statistical Perspective", Foundations and Trends® in

4 Fine-grained analysis: $\ell_\infty$ and $\ell_{2,\infty}$ perturbation theory	126
4.1 Leave-one-out-analysis: An illustrative example . . . . .	127

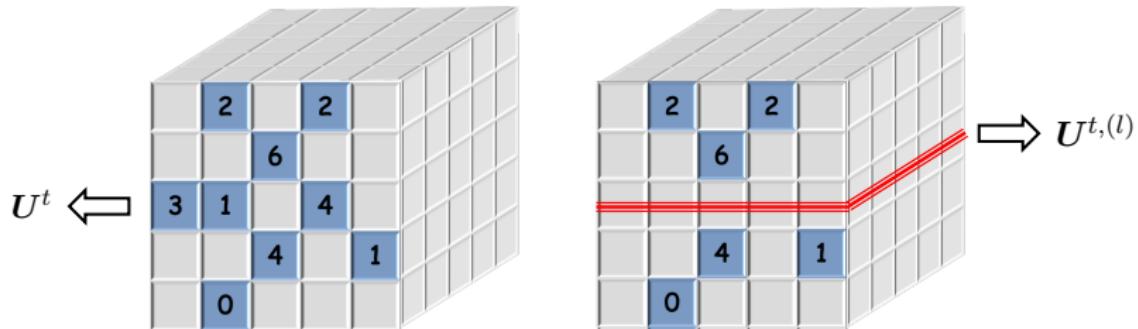
## Key proof idea: leave-one-out decoupling

For each  $1 \leq l \leq d$ , generate leave-one-out auxiliary iterates  $\{\mathbf{U}^{t,(l)}\}$  by replacing  $l^{\text{th}}$  slice with true values



## Key proof idea: leave-one-out decoupling

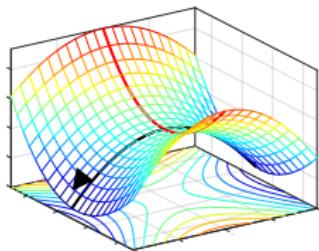
For each  $1 \leq l \leq d$ , generate leave-one-out auxiliary iterates  $\{\mathbf{U}^{t,(l)}\}$  by replacing  $l^{\text{th}}$  slice with true values



- exploit partial statistical independence
- exploit leave-one-out stability
- enable optimal  $\ell_\infty$  error control

# Summary of Part 1

---

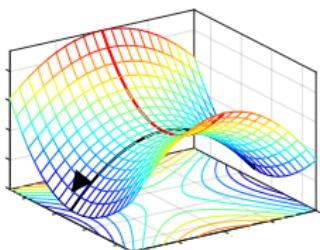


nonconvex  
optimization

- optimal estimation guarantees
- linear-time algorithm
- minimal sample size
- fine-grained uncertainty quantification

# Summary of Part 1

---



nonconvex  
optimization

- optimal estimation guarantees
- linear-time algorithm
- minimal sample size
- fine-grained uncertainty quantification

phase  
retrieval

matrix  
completion

ranking

blind  
deconvolution

joint image  
alignment

mixture  
models

reinforcement  
learning

## *2: Nonconvex optimization in reinforcement learning*



Gen Li  
Tsinghua



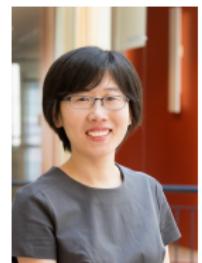
Shicong Cen  
CMU



Chen Cheng  
Stanford



Yuting Wei  
CMU

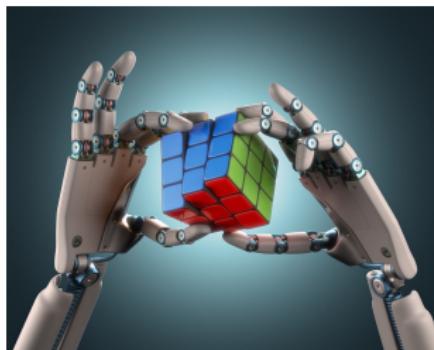


Yuejie Chi  
CMU

"Softmax policy gradient methods can take exponential time to converge," G. Li, Y. Wei, Y. Chi, Y. Gu, Y. Chen, arxiv:2102.11270, 2021

"Fast global convergence of natural policy gradient methods with entropy regularization," S. Cen, C. Cheng, Y. Chen, Y. Wei, Y. Chi, under revision, *Operations Research*, 2020

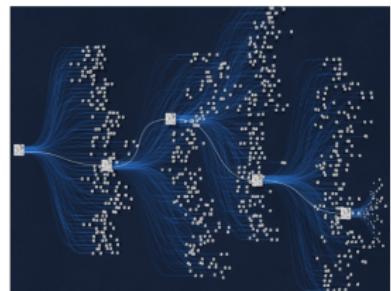
# Recent successes in reinforcement learning (RL)



*Policy optimization: a major contributor to recent RL advances*

# RL challenges

In RL, an agent learns by interacting with an unknown environment



- enormous state and action space
- delayed rewards
- credit assignments
- nonconvexity everywhere

# RL challenges

In RL, an agent learns by interacting with an unknown environment

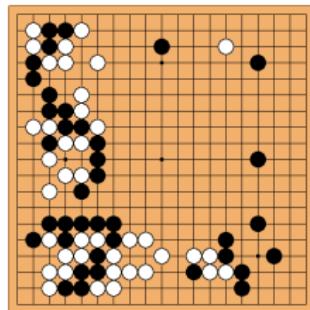
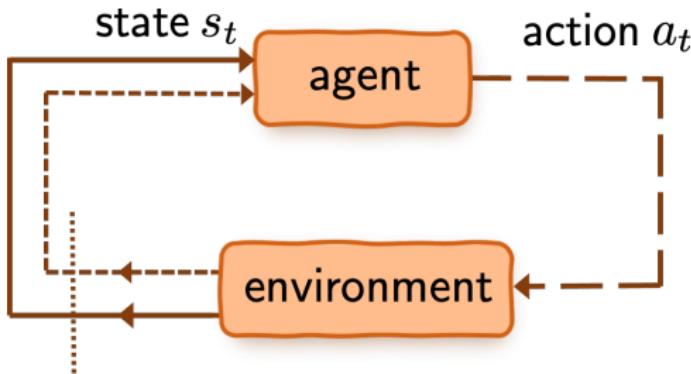


- enormous state and action space
- delayed rewards
- credit assignments
- nonconvexity everywhere

How to enable scalable and guaranteed RL despite nonconvexity?

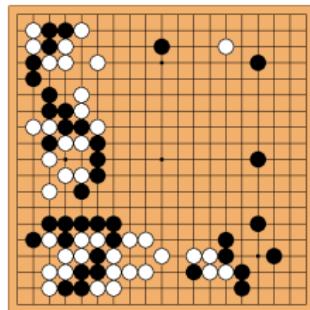
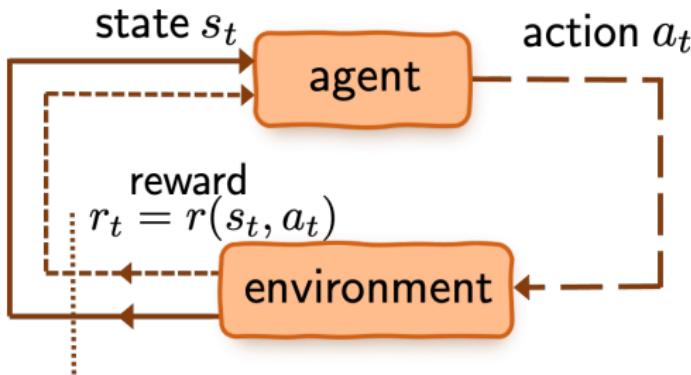
*Backgrounds: policy optimization for MDPs*

# Markov decision process (MDP)



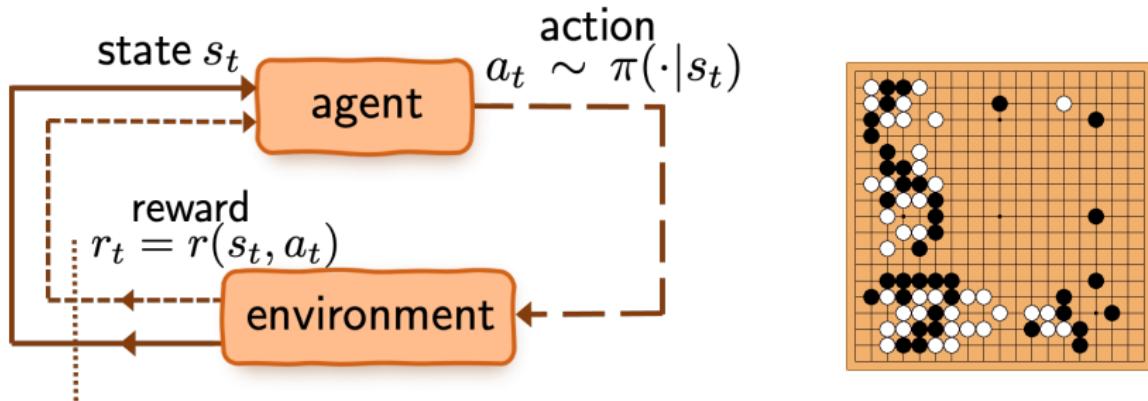
- $\mathcal{S}$ : state space
- $\mathcal{A}$ : action space

# Markov decision process (MDP)



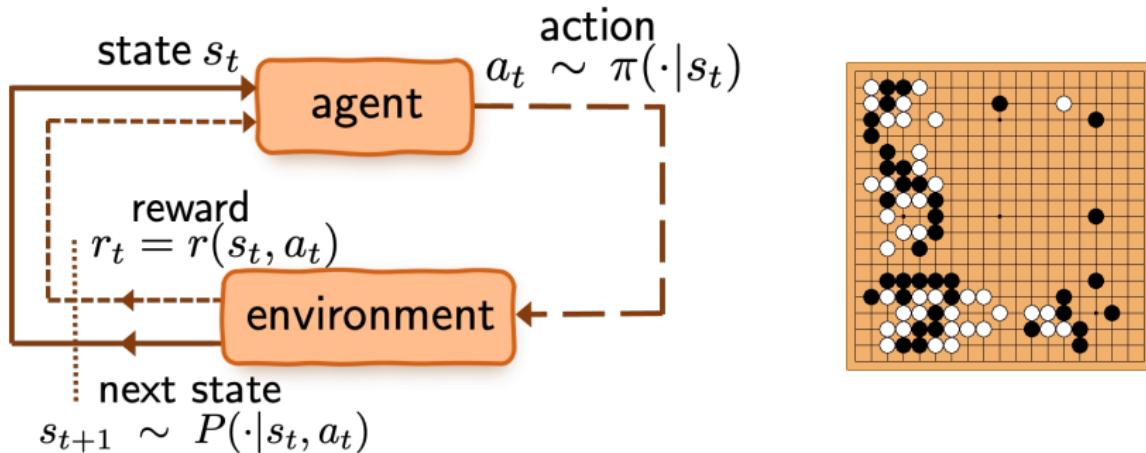
- $\mathcal{S}$ : state space
- $r(s, a) \in [0, 1]$ : immediate reward
- $\mathcal{A}$ : action space

# Markov decision process (MDP)



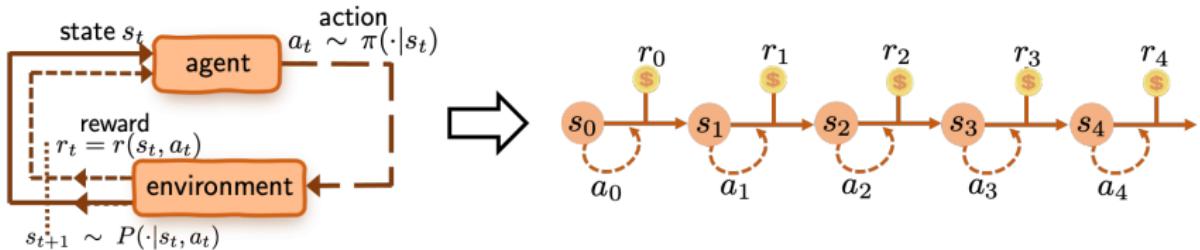
- $\mathcal{S}$ : state space
- $r(s, a) \in [0, 1]$ : immediate reward
- $\pi(\cdot | s)$ : policy (or action selection rule)
- $\mathcal{A}$ : action space

# Markov decision process (MDP)



- $\mathcal{S}$ : state space
- $r(s, a) \in [0, 1]$ : immediate reward
- $\pi(\cdot | s)$ : policy (or action selection rule)
- $P(\cdot | s, a)$ : transition probabilities
- $\mathcal{A}$ : action space

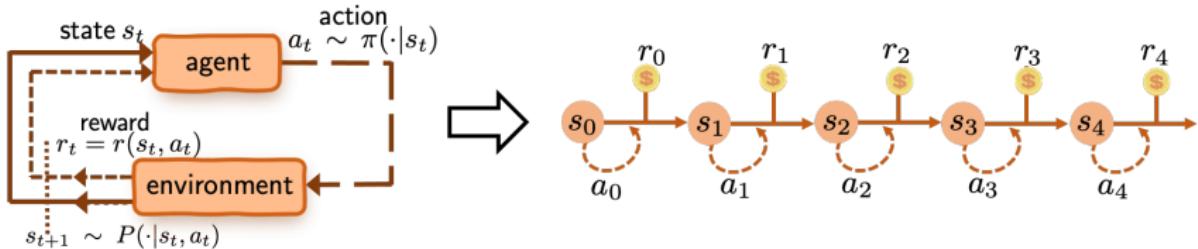
# Value function of policy $\pi$



cumulative discounted reward:  $V^\pi(s) := \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s \right], s \in \mathcal{S}$

- expectation is over randomness of MDP & policy  $\pi$

# Value function of policy $\pi$

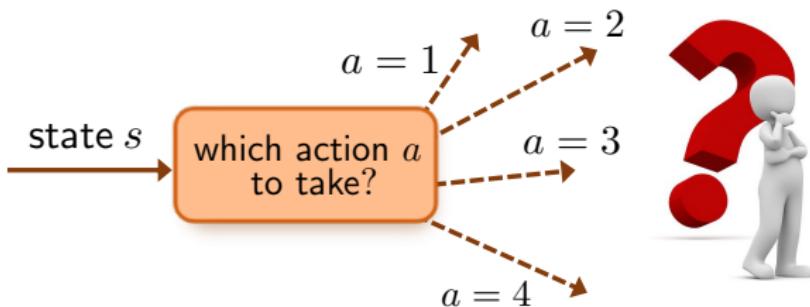


cumulative discounted reward:  $V^\pi(s) := \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s \right], s \in \mathcal{S}$

- expectation is over randomness of MDP & policy  $\pi$
- $\gamma \in [0, 1]$ : discount factor
  - take  $\gamma \rightarrow 1$  to approximate long-horizon MDPs
  - **effective horizon:**  $\frac{1}{1-\gamma}$

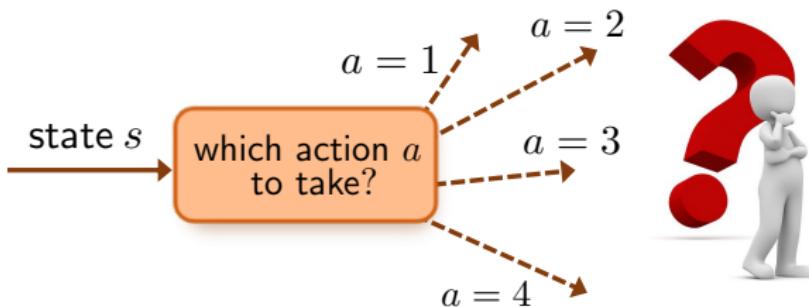
# Optimal policy and optimal value

---



- **goal:** find optimal policy  $\pi^*$  that maximizes value functions
- optimal value function:  $V^*(s) := \max_{\pi} V^{\pi}(s)$  for all  $s \in \mathcal{S}$

# Optimal policy and optimal value



- **goal:** find optimal policy  $\pi^*$  that maximizes value functions
- optimal value function:  $V^*(s) := \max_{\pi} V^{\pi}(s)$  for all  $s \in \mathcal{S}$

How to accomplish it via nonconvex optimization algorithms?

# Policy optimization

---

Given state distribution  $s \sim \rho$   
(e.g. uniform)

$$\max_{\pi} V^{\pi}(\rho) := \mathbb{E}_{s \sim \rho} [V^{\pi}(s)]$$

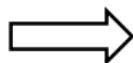
# Policy optimization

---

Given state distribution  $s \sim \rho$   
(e.g. uniform)

$$\max_{\pi} V^{\pi}(\rho) := \mathbb{E}_{s \sim \rho} [V^{\pi}(s)]$$

parameterize



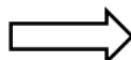
$$\max_{\theta} V^{\pi_{\theta}}(\rho) := \mathbb{E}_{s \sim \rho} [V^{\pi_{\theta}}(s)]$$

# Policy optimization

Given state distribution  $s \sim \rho$   
(e.g. uniform)

$$\max_{\pi} V^{\pi}(\rho) := \mathbb{E}_{s \sim \rho} [V^{\pi}(s)]$$

parameterize



$$\max_{\theta} V^{\pi_{\theta}}(\rho) := \mathbb{E}_{s \sim \rho} [V^{\pi_{\theta}}(s)]$$

$$\pi_{\theta}(a|s) = \frac{\exp(\theta(s, a))}{\sum_a \exp(\theta(s, a))}$$

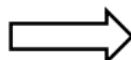
softmax parameterization

# Policy optimization

Given state distribution  $s \sim \rho$   
(e.g. uniform)

$$\max_{\pi} V^{\pi}(\rho) := \mathbb{E}_{s \sim \rho} [V^{\pi}(s)]$$

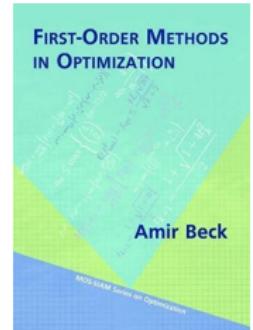
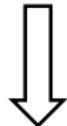
parameterize



$$\max_{\theta} V^{\pi_{\theta}}(\rho) := \mathbb{E}_{s \sim \rho} [V^{\pi_{\theta}}(s)]$$

$$\pi_{\theta}(a|s) = \frac{\exp(\theta(s, a))}{\sum_a \exp(\theta(s, a))}$$

softmax parameterization

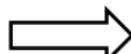


# Policy optimization

Given state distribution  $s \sim \rho$   
(e.g. uniform)

$$\max_{\pi} V^{\pi}(\rho) := \mathbb{E}_{s \sim \rho} [V^{\pi}(s)]$$

parameterize



$$\max_{\theta} V^{\pi_{\theta}}(\rho) := \mathbb{E}_{s \sim \rho} [V^{\pi_{\theta}}(s)]$$

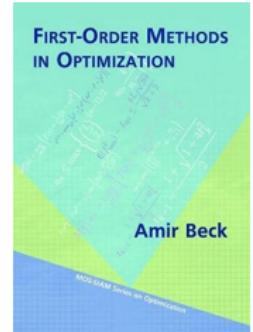
$$\pi_{\theta}(a|s) = \frac{\exp(\theta(s, a))}{\sum_a \exp(\theta(s, a))}$$

softmax parameterization

Policy gradient method (Sutton et al. '00)

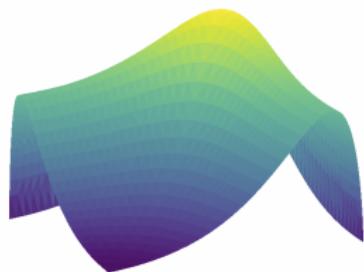
$$\theta^{(t+1)} = \theta^{(t)} + \eta \nabla_{\theta} V^{(t)}(\rho), \quad t = 0, 1, \dots$$

- $\eta$ : learning rate



# Does policy gradient (PG) method converge?

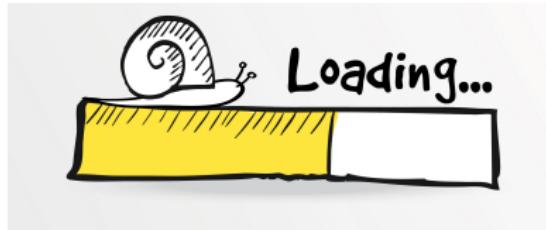
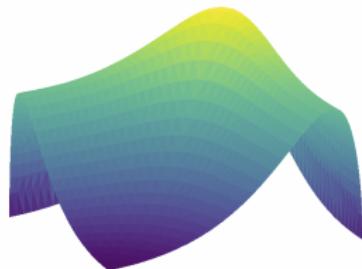
---



- (Agarwal et al. '19) Softmax PG converges to global opt as  $t \rightarrow \infty$

# Does policy gradient (PG) method converge?

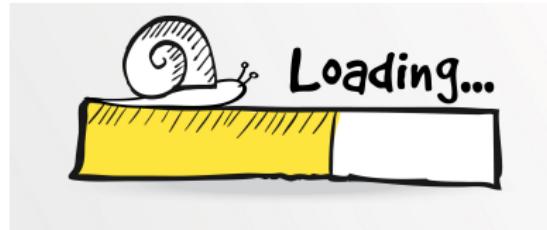
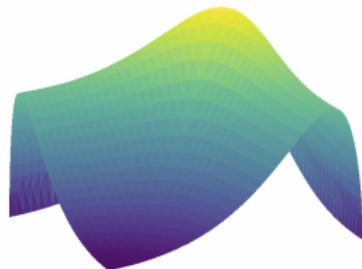
---



- (Agarwal et al. '19) Softmax PG converges to global opt as  $t \rightarrow \infty$

However, “asymptotic convergence” might mean “taking forever”

# Does policy gradient (PG) method converge?

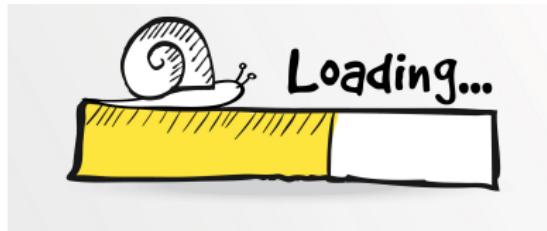
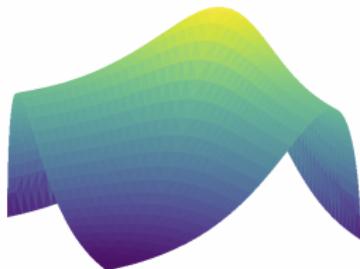


- (Agarwal et al. '19) Softmax PG converges to global opt as  $t \rightarrow \infty$
- (Mei et al. '20) Softmax PG converges to global opt in

$$O\left(\frac{1}{\varepsilon}\right) \text{ iterations}$$

However, “asymptotic convergence” might mean “taking forever”

# Does policy gradient (PG) method converge?



- (Agarwal et al. '19) Softmax PG converges to global opt as  $t \rightarrow \infty$
- (Mei et al. '20) Softmax PG converges to global opt in

$$c(|\mathcal{S}|, |\mathcal{A}|, \frac{1}{1-\gamma}, \dots) O(\frac{1}{\varepsilon}) \text{ iterations}$$

However, “asymptotic convergence” might mean “taking forever”

# A negative message

---

## Theorem 2 (Li, Wei, Chi, Chen '21)

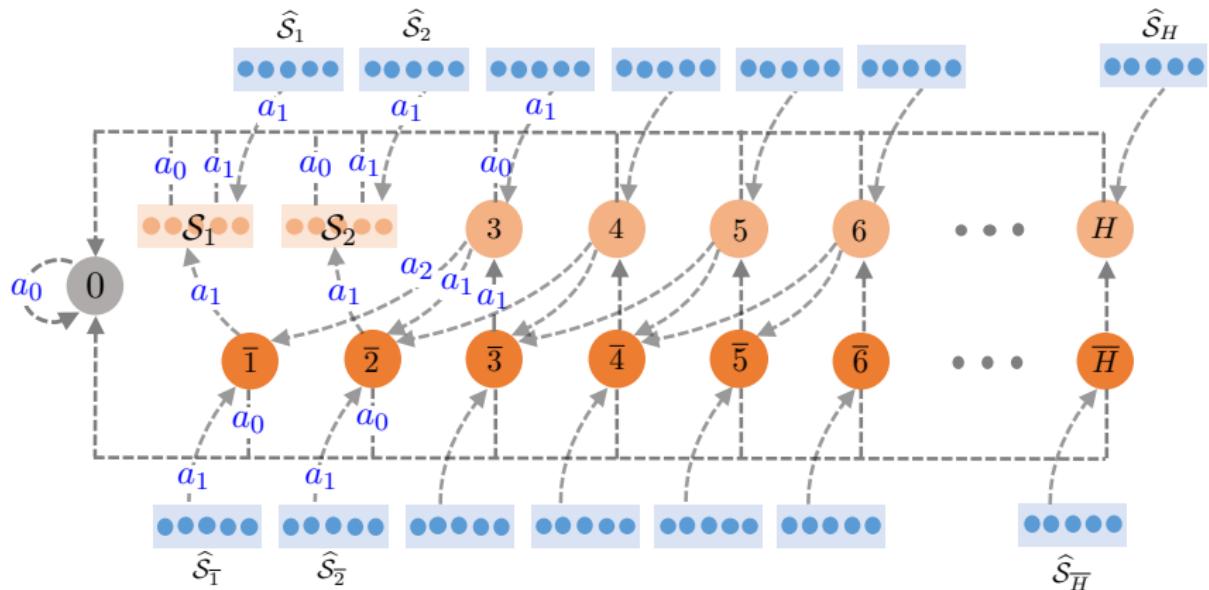
There exists an MDP s.t. it takes softmax PG at least

$$\frac{1}{\eta} |\mathcal{S}|^{1.5^{\Theta(\frac{1}{1-\gamma})}} \text{ iterations}$$

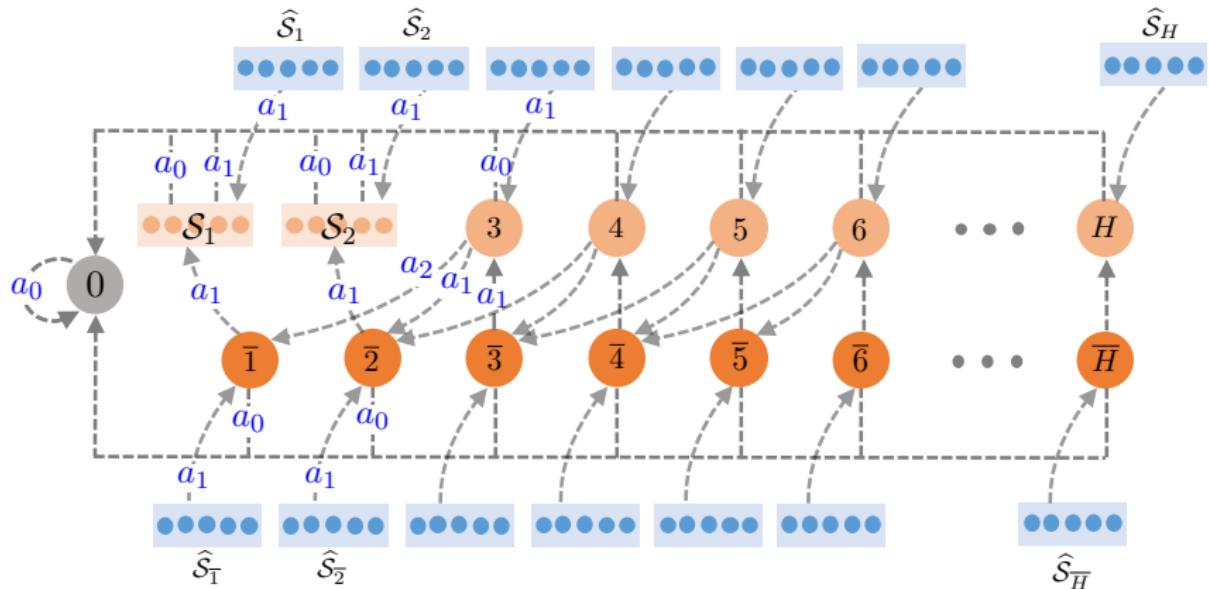
to achieve  $\|V^{(t)} - V^*\|_\infty \leq 1/2$  (even with infinite samples)

- Softmax PG method can take **exponential time** to converge (in problems w/ large state space & long effective horizon)!

# MDP construction for our lower bound



# MDP construction for our lower bound

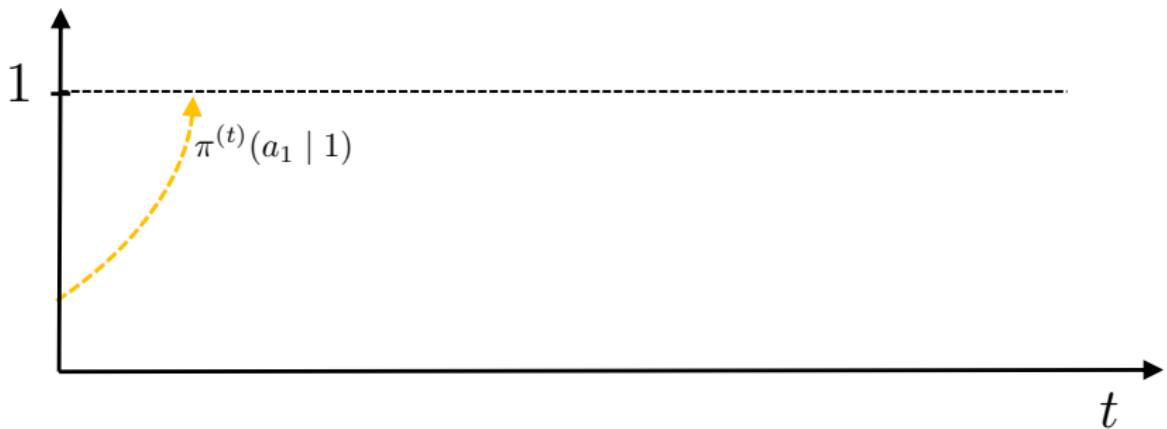


Key design ingredients: for  $3 \leq s \leq H \asymp \frac{1}{1-\gamma}$ ,

- $a_1$  is optimal action
- delayed rewards
- $\pi^{(t)}(a_1|s)$  keeps decreasing until  $\pi^{(t)}(a_1|s-2) \approx 1$

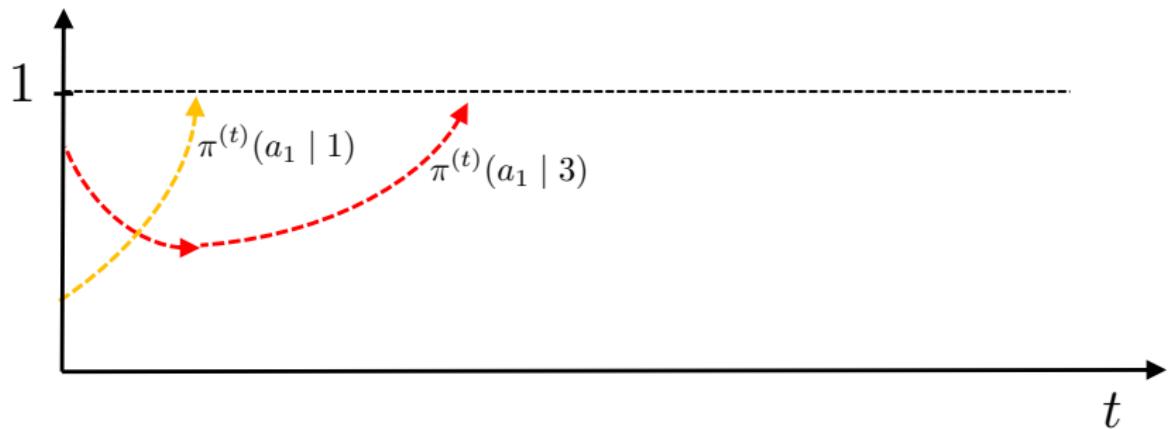
# What is happening in our constructed MDP?

---



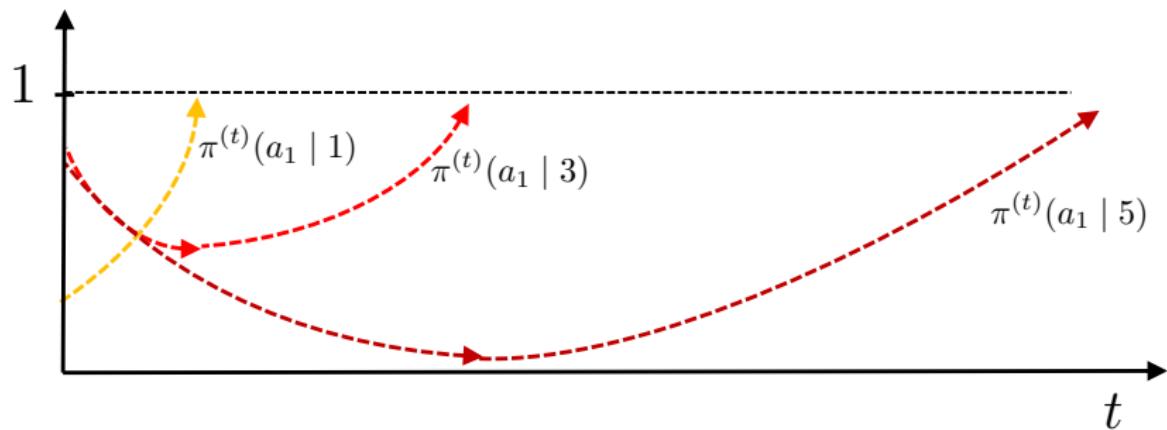
# What is happening in our constructed MDP?

---



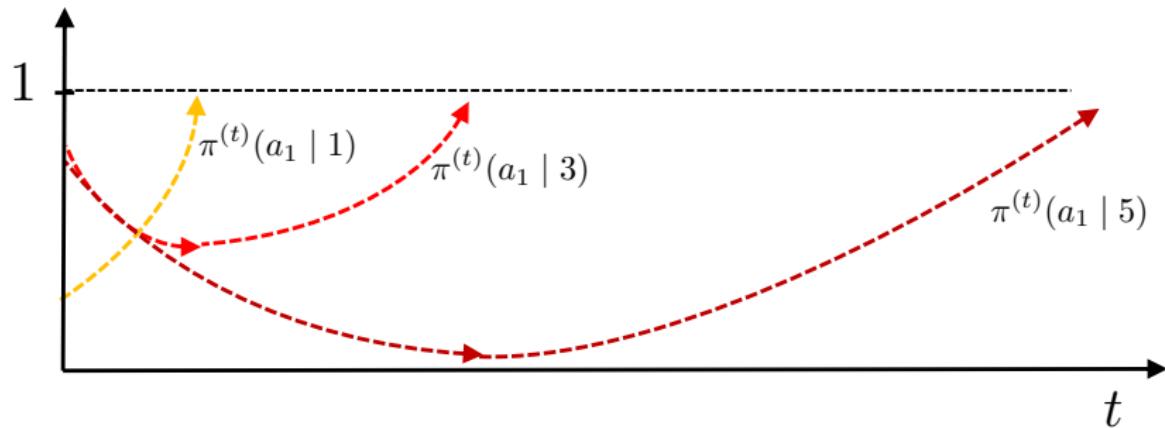
# What is happening in our constructed MDP?

---



**observation:** convergence time for state  $s$  grows geometrically as  $s \uparrow$

# What is happening in our constructed MDP?



**observation:** convergence time for state  $s$  grows geometrically as  $s \uparrow$

$$\text{convergence-time}(s) \gtrsim (\text{convergence-time}(s - 2))^{1.5}$$

## Booster 1: entropy regularization

---

*accelerate convergence by regularizing objective function*

$$V_\tau^\pi(s_0) := \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t (r_t - \tau \log \pi(a_t | s_t)) \mid s_0 \right]$$

## Booster 1: entropy regularization

---

accelerate convergence by regularizing objective function

$$\begin{aligned} V_\tau^\pi(s_0) &:= \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t (r_t - \color{blue}{\tau \log \pi(a_t|s_t)}) \mid s_0 \right] \\ &= V^\pi(s) + \frac{\color{blue}{\tau}}{1-\gamma} \mathbb{E}_{s \sim d_s^\pi} \underbrace{\left[ - \sum_a \pi(a|s) \log \pi(a|s) \mid s_0 \right]}_{\text{Shannon entropy}} \end{aligned}$$

- $\tau$ : regularization parameter
- $d_s^\pi$ : certain marginal distribution

# Booster 1: entropy regularization

accelerate convergence by regularizing objective function

$$\begin{aligned} V_\tau^\pi(s_0) &:= \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t (r_t - \color{blue}{\tau \log \pi(a_t|s_t)}) \mid s_0 \right] \\ &= V^\pi(s) + \frac{\color{blue}{\tau}}{1-\gamma} \mathbb{E}_{s \sim d_s^\pi} \underbrace{\left[ - \sum_a \pi(a|s) \log \pi(a|s) \mid s_0 \right]}_{\text{Shannon entropy}} \end{aligned}$$

- $\tau$ : regularization parameter
- $d_s^\pi$ : certain marginal distribution

entropy-regularized value maximization

$$\underset{\theta}{\text{maximize}} \quad V_{\color{red}{\tau}}^{\pi_\theta}(\rho) := \mathbb{E}_{s \sim \rho} [V_{\color{red}{\tau}}^{\pi_\theta}(s)]$$

## Our negative message remains . . .

### Theorem 3 (Li, Wei, Chi, Chen '21)

*There is an MDP s.t. it takes entropy-regularized softmax PG at least*

$$\min \left\{ \exp \left( \Theta \left( \frac{1}{\varepsilon} \right) \right), \frac{1}{\eta} |\mathcal{S}|^{2^{\Theta(\frac{1}{1-\gamma})}} \right\} \text{ iterations}$$

*to achieve  $\|V^{(t)} - V^*\|_\infty \leq \varepsilon$  (even with infinite samples)*

- Softmax PG method with entropy regularization can still take **exponential time** to converge!

## Our negative message remains . . .

### Theorem 3 (Li, Wei, Chi, Chen '21)

*There is an MDP s.t. it takes entropy-regularized softmax PG at least*

$$\min \left\{ \exp \left( \Theta \left( \frac{1}{\varepsilon} \right) \right), \frac{1}{\eta} |\mathcal{S}|^{2^{\Theta(\frac{1}{1-\gamma})}} \right\} \text{ iterations}$$

*to achieve  $\|V^{(t)} - V^*\|_\infty \leq \varepsilon$  (even with infinite samples)*

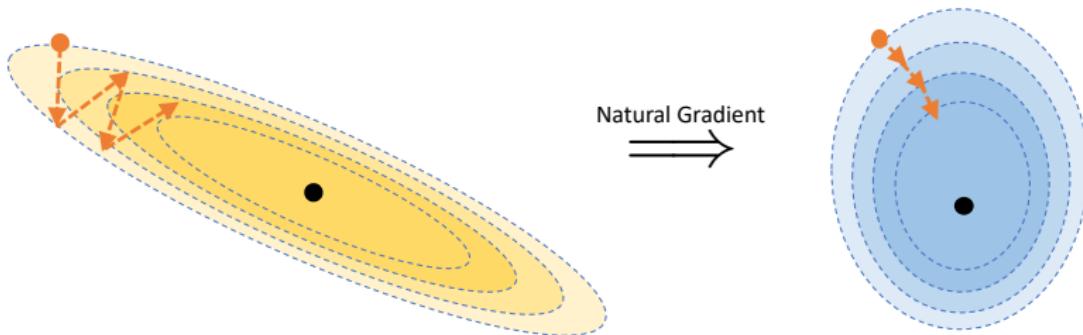
- Softmax PG method with entropy regularization can still take **exponential time** to converge!
- (Mei et al. '20) entropy-regularized softmax PG converges in

$$c(|\mathcal{S}|, |\mathcal{A}|, \frac{1}{1-\gamma}, \dots) O(\frac{1}{\varepsilon}) \text{ iterations}$$

## Booster 2: natural policy gradient (NPG)

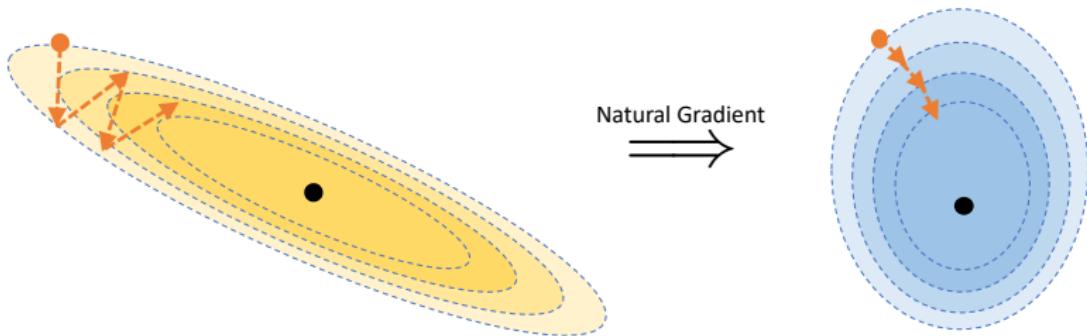
---

*precondition gradients to improve search directions ...*



## Booster 2: natural policy gradient (NPG)

*precondition gradients to improve search directions ...*



NPG method (Kakade '02)

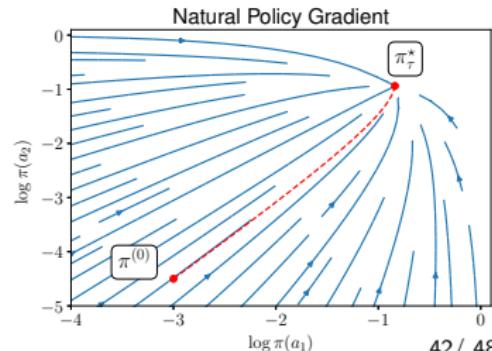
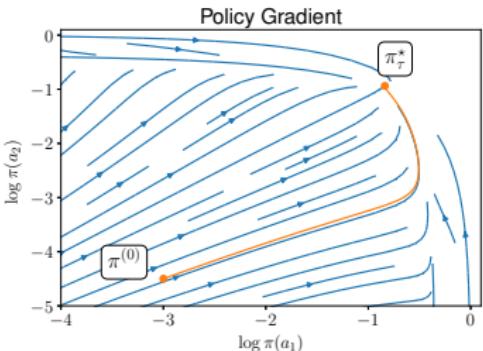
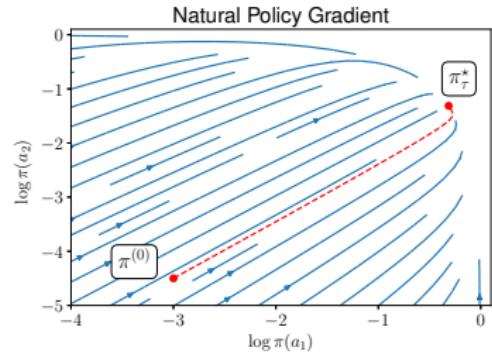
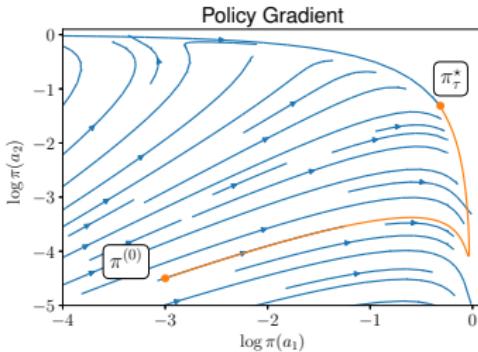
$$\theta^{(t+1)} = \theta^{(t)} + \eta (\mathcal{F}_\rho^\theta)^\dagger \nabla_\theta V_\tau^{(t)}(\rho), \quad t = 0, 1, \dots$$

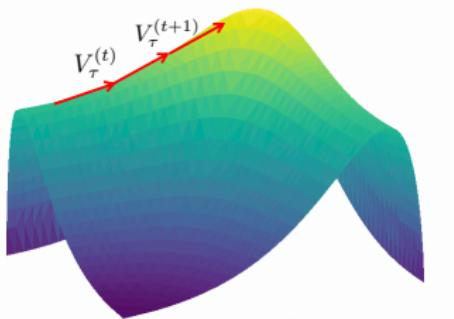
- $\mathcal{F}_\rho^\theta := \mathbb{E} \left[ (\nabla_\theta \log \pi_\theta(a|s)) (\nabla_\theta \log \pi_\theta(a|s))^\top \right]$ : Fisher info

# Entropy-regularized natural gradient helps!

A toy bandit example: 3 arms with rewards 1, 0.9 and 0.1

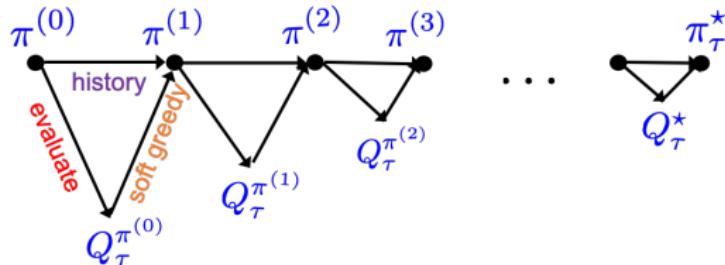
increase regularization





*How to characterize the efficiency of  
entropy-regularized NPG in tabular settings?*

# Linear convergence with exact gradients



exact access to gradients:

## Theorem 4 (Cen, Cheng, Chen, Wei, Chi '20)

If  $\eta \leq (1 - \gamma)/\tau$  and  $\tau \leq \frac{(1-\gamma)\varepsilon}{4 \log |\mathcal{A}|}$ , entropy-regularized NPG achieves

$$\|V^* - V^{(t+1)}\|_\infty \leq C_1 \gamma (1 - \eta \tau)^t + \varepsilon/2, \quad t = 0, 1, \dots$$

## Implications: iteration complexity

---

number of iterations needed to reach  $\|V^* - V^{(t)}\|_\infty \leq \varepsilon$ :

- **general learning rates** ( $0 < \eta < \frac{1-\gamma}{\tau}$ ):

$$\frac{1}{\eta\tau} \log \left( \frac{2C_1\gamma}{\varepsilon} \right)$$

- **optimal choice** ( $\eta = \frac{1-\gamma}{\tau}$ ):

$$\frac{1}{1-\gamma} \log \left( \frac{2C_1\gamma}{\varepsilon} \right)$$

## Implications: iteration complexity

---

number of iterations needed to reach  $\|V^* - V^{(t)}\|_\infty \leq \varepsilon$ :

- **general learning rates** ( $0 < \eta < \frac{1-\gamma}{\tau}$ ):

$$\frac{1}{\eta\tau} \log \left( \frac{2C_1\gamma}{\varepsilon} \right)$$

- **optimal choice** ( $\eta = \frac{1-\gamma}{\tau}$ ):

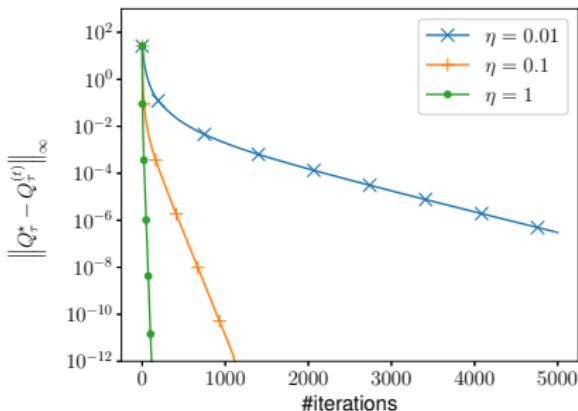
$$\frac{1}{1-\gamma} \log \left( \frac{2C_1\gamma}{\varepsilon} \right)$$

Nearly dimension-free global linear convergence!

# Regularized NPG vs. unregularized NPG

regularized NPG

$\tau = 0.001$

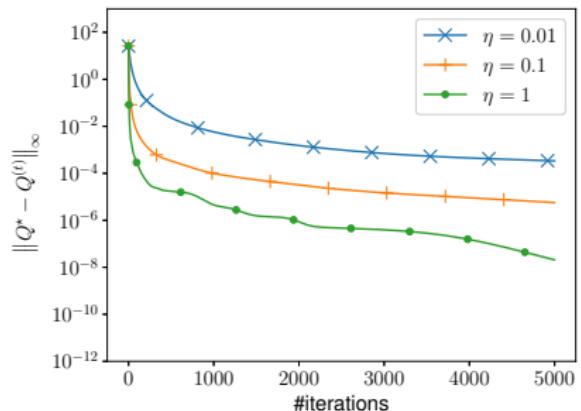


$$\text{linear rate: } \frac{1}{\eta\tau} \log\left(\frac{1}{\varepsilon}\right)$$

**Ours**

unregularized NPG

$\tau = 0$



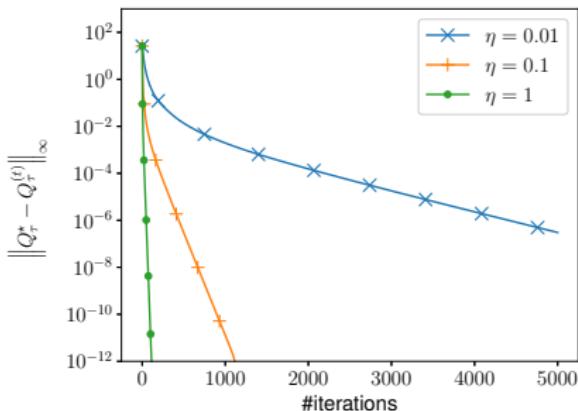
$$\text{sublinear rate: } \frac{1}{\min\{\eta, (1-\gamma)^2\}\varepsilon}$$

(Agarwal et al. '19)

# Regularized NPG vs. unregularized NPG

regularized NPG

$\tau = 0.001$

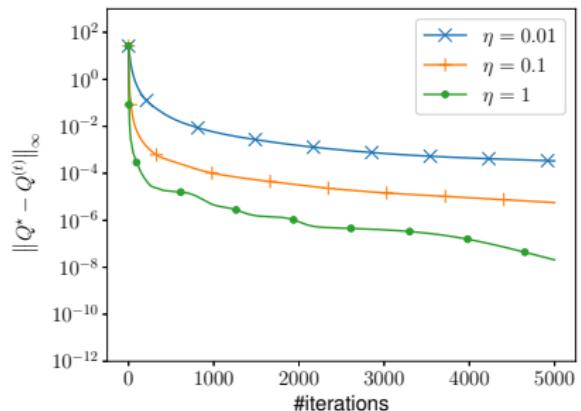


$$\text{linear rate: } \frac{1}{\eta\tau} \log\left(\frac{1}{\varepsilon}\right)$$

**Ours**

unregularized NPG

$\tau = 0$



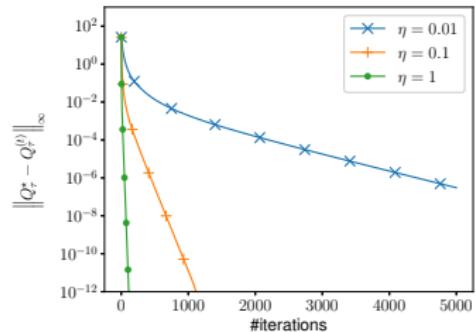
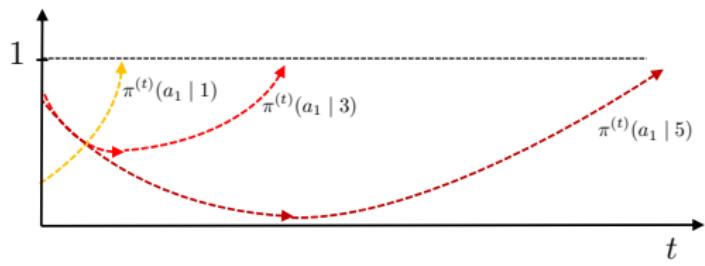
$$\text{sublinear rate: } \frac{1}{\min\{\eta, (1-\gamma)^2\}\varepsilon}$$

(Agarwal et al. '19)

Entropy regularization enables faster convergence!

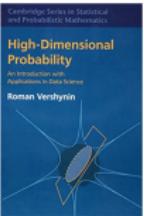
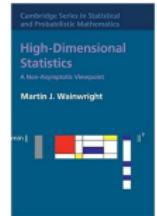
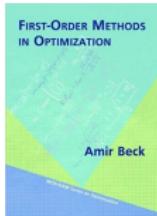
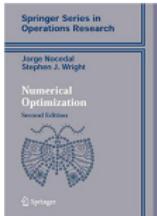
# Summary of Part 2

---



- Softmax policy gradient can take exponential time to converge
- Entropy regularization & natural gradients help!

# Concluding remarks



nonconvex optimization

(high-dimensional) statistics

