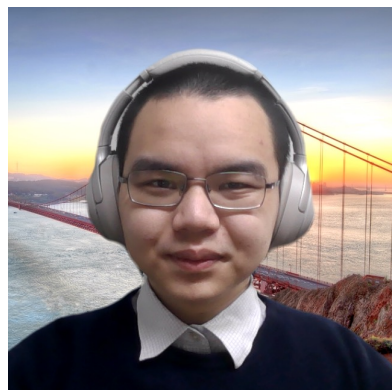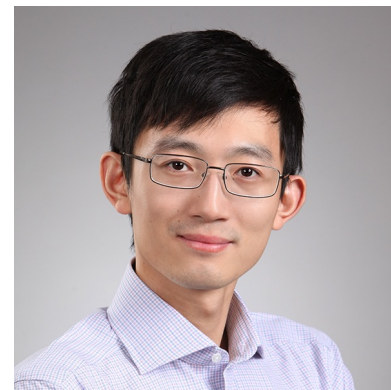# Provable Model-based Nonlinear Bandit and RL: Shelve Optimism, Embrace Virtual Curvature

Kefan Dong

Jiaqi Yang

Tsinghua University

Tengyu Ma

Stanford

# Toward a Theory for Deep RL

# Existing RL theory cannot apply to Neural Nets

- None of these give polynomial sample complexities for even one-layer NNs.

| | B-Rank | B-Complete | W-Rank | Bilinear Class (this work) |
|---|---|---|---|---|
| Tabular MDP | ✓ | ✓ | ✓ | ✓ |
| Reactive POMDP [Krishnamurthy et al., 2016] | ✓ | ✗ | ✓ | ✓ |
| Block MDP [Du et al., 2019a] | ✓ | ✗ | ✓ | ✓ |
| Flambe / Feature Selection [Agarwal et al., 2020b] | ✓ | ✗ | ✓ | ✓ |
| Reactive PSR [Littman and Sutton, 2002] | ✓ | ✗ | ✓ | ✓ |
| Linear Bellman Complete [Munos, 2005] | ✗ | ✓ | ✗ | ✓ |
| Linear MDPs [Yang and Wang, 2019, Jin et al., 2020] | ✓ ! | ✓ | ✓ ! | ✓ |
| Linear Mixture Model [Modi et al., 2020b] | ✗ | ✗ | ✗ | ✓ |
| Linear Quadratic Regulator | ✗ | ✓ | ✗ | ✓ |
| Kernelized Nonlinear Regulator [Kakade et al., 2020] | ✗ | ✗ | ✗ | ✓ |
| $Q^\star$ "irrelevant" State Aggregation [Li, 2009] | ✓ | ✗ | ✗ | ✓ |
| Linear $Q^\star/V^\star$ (this work) | ✗ | ✗ | ✗ | ✓ |
| RKHS Linear MDP (this work) | ✗ | ✗ | ✗ | ✓ |
| RKHS Linear Mixture MDP (this work) | ✗ | ✗ | ✗ | ✓ |
| Low Occupancy Complexity (this work) | ✗ | ✗ | ✗ | ✓ |
| $Q^\star$ State-action Aggregation [Dong et al., 2020] | ✗ | ✗ | ✗ | ✗ |
| Deterministic linear $Q^\star$ [Wen and Van Roy, 2013] | ✗ | ✗ | ✗ | ✗ |
| Linear $Q^\star$ [Weisz et al., 2020] | Sample efficiency is not possible | | | |

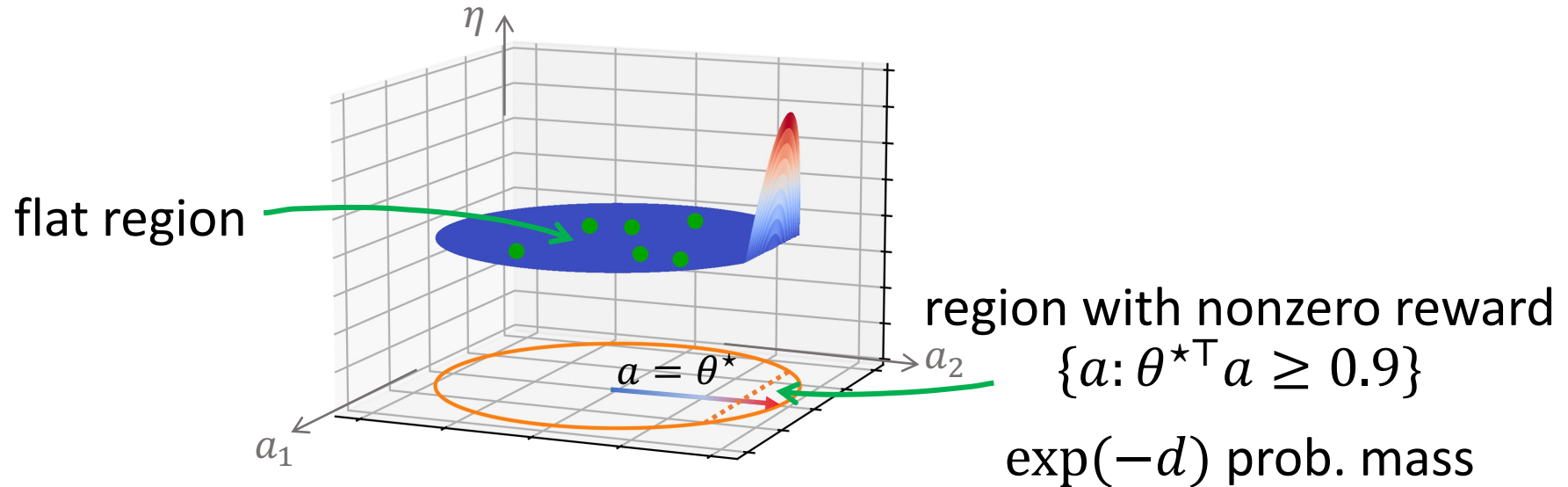Du, Simon S., et al. "Bilinear Classes: A Structural Framework for Provable Generalization in RL."

# Neural Net Bandit: A Simplification

- Reward function $\eta(\theta, a)$
  - $\theta \in \Theta$: model parameter
  - $a \in \mathcal{A}$: continuous action
- Linear bandit: $\eta(\theta, a) = \theta^\top a$
- Neural net bandit: $\eta(\theta, a) = \mathrm{NN}_\theta(a)$
- Realizable and deterministic reward setting:
  - Agent observes ground-truth reward $\eta(\theta^\star, a)$ after playing action $a$
- Goal: finding the best action
$$a^\star = \mathrm{argmax}_{a \in \mathcal{A}} \, \eta(\theta^\star, a)$$

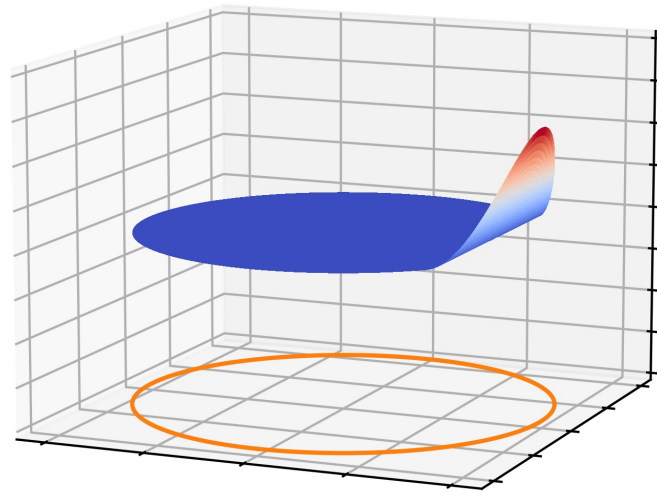# Neural Net Bandit is Statistically Hard!

- $\Theta, \mathcal{A}$: unit $\ell_2$-ball in $\mathbb{R}^d$

- $\eta(\theta, a) = \text{relu}(\theta^\top a - 0.9), \quad a^\star = \underset{||a||_2 \leq 1}{\text{argmax}} \, \text{relu}(\theta^{\star\top} a - 0.9) = \theta^\star$



flat region

region with nonzero reward
$\{a : \theta^{\star\top} a \geq 0.9\}$

$\exp(-d)$ prob. mass

needle in a haystack!

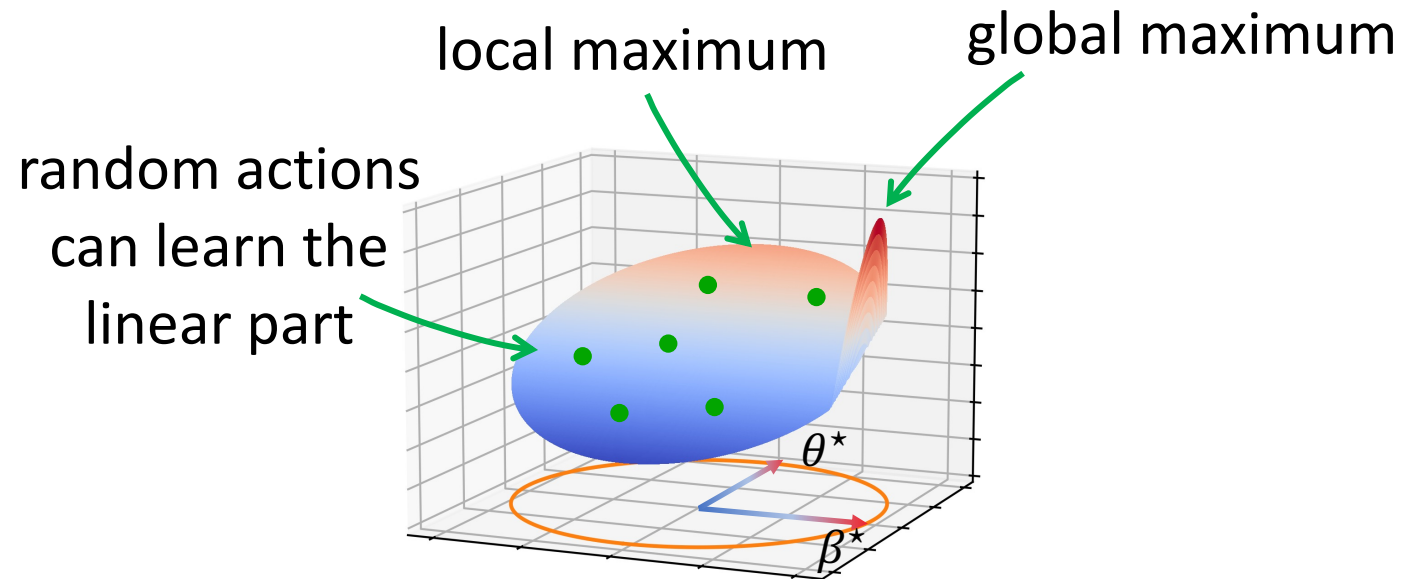# Neural Net Bandit is <span style="color:red">Statistically</span> Hard!

- $\Theta, \mathcal{A}$: unit $\ell_2$-ball in $\mathbb{R}^d$

- $\eta(\theta, a) = \mathrm{relu}(\theta^\top a - 0.9)^2$



smoothed version

# Neural Net Bandit is Statistically Hard!

- Convergence to a global maximum is generally statistically intractable

- Existing RL theory cannot apply to NNs because they aim for global maximum

local maximum

global maximum

random actions can learn the linear part

$\theta^\star$

$\beta^\star$

$$\eta((\theta, \beta), a) = \theta^\top a + 20 \cdot \text{relu}(\beta^\top a - 0.9)$$

needle in a haystack!

# A New Paradigm for Bandit/RL

1. Convergences to local maxima for general instances     ☞  This talk

2. Analysis of the landscape of the true reward $\eta(\theta^\star, \cdot)$

# Main Results

- Theorem (informal): Under Lipschitz assumptions on $\eta$, there exists an algorithm that converges to a $\epsilon$-approximate local maxima in $\tilde{O}(R(\Theta)\epsilon^{-8})$.

  measures hardness of online learning w.r.t. model class

- Similar results for nonlinear RL (with many more assumptions and stochastic policies.)

# Baseline: Zero-order Optimization for Bandit

- True reward $f(a) = \eta(\theta^\star, a)$

- Zero-order optimization:
  - estimate gradient $\nabla f(a)$ by finite difference
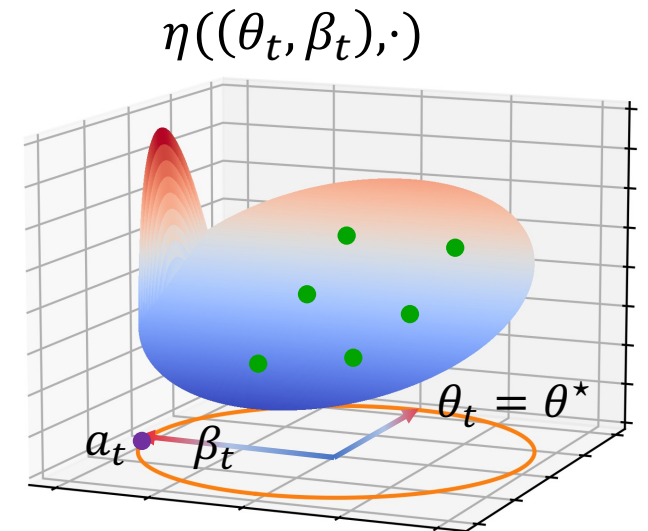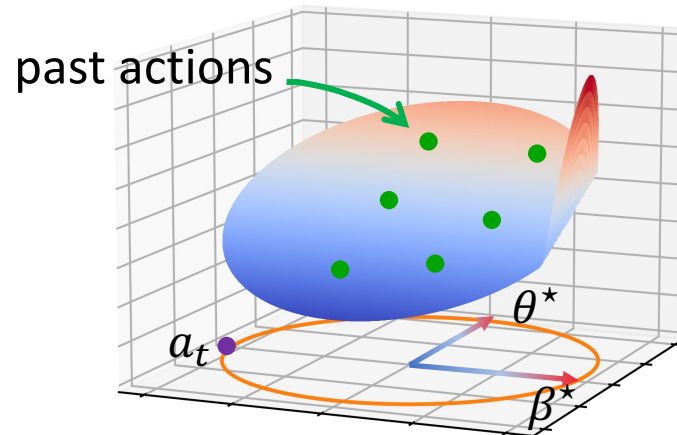
    For $\xi \sim \mathcal{N}(0, I)$ and $\epsilon > 0$,
    $$\frac{1}{\epsilon} \mathbb{E}\big[\xi\big(f(a + \epsilon\xi) - f(a)\big)\big] \approx \mathbb{E}[\xi\xi^\top \nabla f(a)] = \nabla f(a)$$

  - $\Omega(d)$ sample complexity

- Our key idea: leverage model extrapolation

# Model-based UCB Does Even Not Converge To Local Max

$$a_t, \theta_t = \underset{a \in \mathcal{A}, \theta \in \Theta_t}{\mathrm{argmax}} \ \eta(\theta, a)$$

$$\underbrace{\qquad\qquad}_{\text{confidence region}}$$

- $\Theta_t$ pins down $\theta^\star$ but has no clue about $\beta^\star$
- UCB keeps guessing $\beta_t$
- and choses $a_t = \beta_t$



$\eta((\theta_t, \beta_t), \cdot)$

past actions

$\eta((\theta, \beta), a) = \theta^\top a + 20 \cdot \mathrm{relu}(\beta^\top a - 0.9)$

- UCB over-explores and doesn't converge in polynomial steps
- In partice, deep RL methods with optimism also over-explore

# Reviewing the Analysis of UCB

1. Optimization (high virtual reward):

$$\text{by optimism, } \eta(\theta_t, a_t) \geq \eta(\theta^\star, a^\star)$$

2. Extrapolation (in average):

$$\sum_{t=1}^{T} \big(\eta(\theta_t, a_t) - \eta(\theta^\star, a_t)\big)^2 \leq \sqrt{\underbrace{\dim_E(\Theta)} \cdot T}$$

Eluder dimension

- $1 + 2 \Rightarrow \quad \eta(\theta^\star, a_t) \rightarrow \eta(\theta^\star, a^\star)$

- Step 2 fails for neural net models because $\dim_E(\Theta) \approx \exp(d)$

# Re-Prioritizing the Two Steps

1. Extrapolation by online learning (OL) oracles:

$$\mathbb{E}\left[\sum_{t=1}^{T}\left(\eta(\theta_t, a_t) - \eta(\theta^\star, a_t)\right)^2\right] \leq \sqrt{R(\Theta)\mathrm{T}\,\mathrm{polylog}(\mathrm{T})}$$

OL oracle outputs a distribution of $\theta_t$

Sequential Rademacher Complexity

[Rakhlin-Sridharan-Tewari'15]

- For finite hypothesis $\Theta$, $R(\Theta) = \log|\Theta|$
- For neural nets:
    $R(\Theta) = \mathrm{poly}(d)$    vs.    Eluder dim $= \exp(d)$
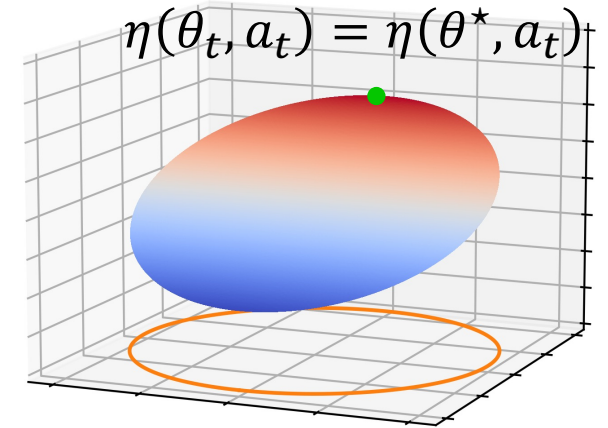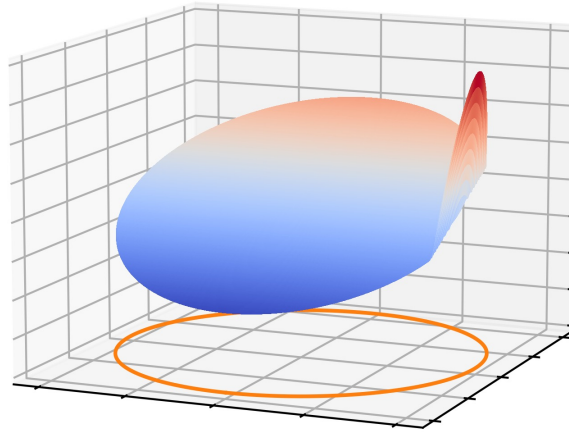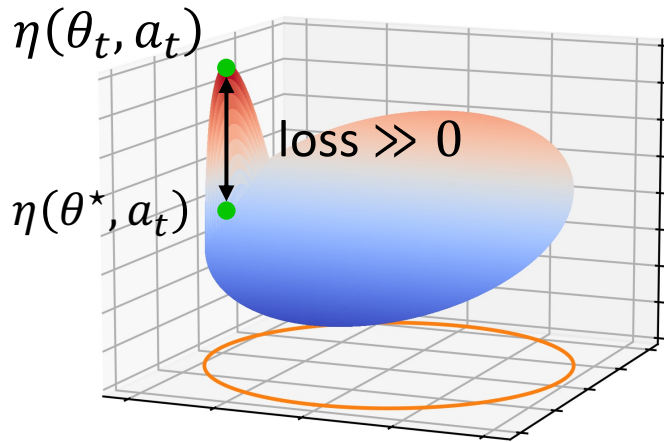
    or the weight norm

Rakhlin, Alexander, Karthik Sridharan, and Ambuj Tewari. "Online learning via sequential complexities." 2015.

# OL Oracle Extrapolates Optimally

Extrapolation error: $\sum_{t=1}^{T}\big(\eta(\theta_t, a_t) - \eta(\theta^\star, a_t)\big)^2$

UCB: loss $\gg 0$

Ground truth

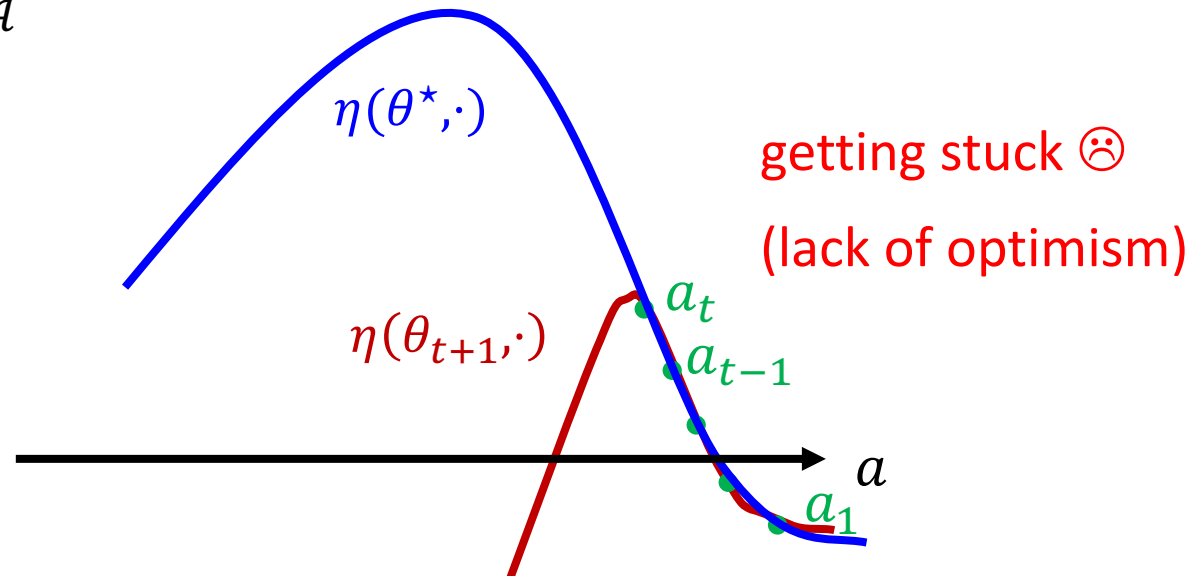OL oracle: loss = 0
(for most of the times)

# Re-Prioritizing the Two Steps

1. Extrapolation by online learning (OL) oracles

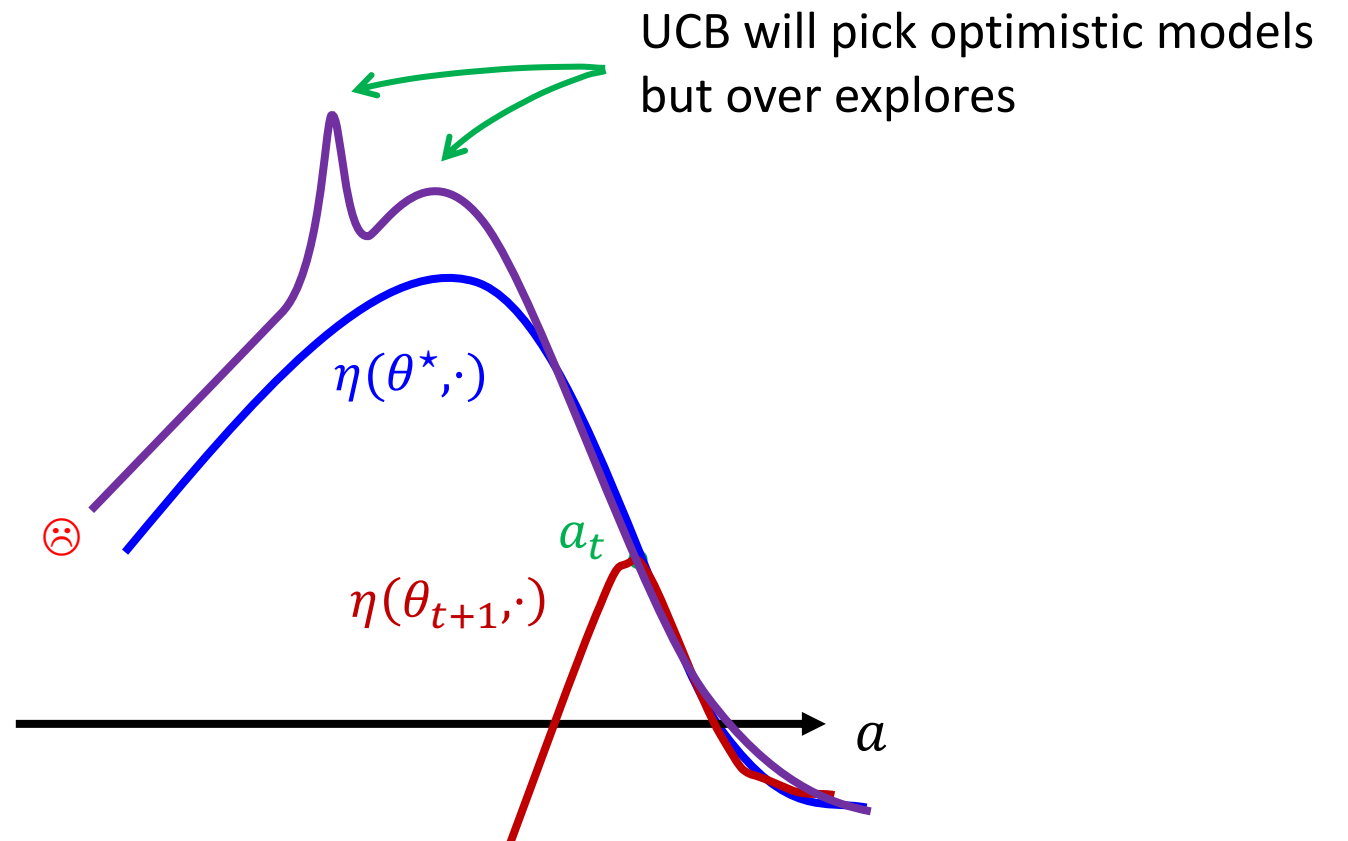$$\mathbb{E}\left[\sum_{t=1}^{T}\left(\eta(\theta_t, a_t) - \eta(\theta^\star, a_t)\right)^2\right] \leq \sqrt{R(\Theta) T \text{ polylog(T)}}$$

2. High virtual reward:

best attempt: $a_t = \underset{a \in \mathcal{A}}{\text{argmax}} \, \mathbb{E}[\eta(\theta_t, a)]$

# Re-Prioritizing the Two Steps



UCB will pick optimistic models but over explores

$\eta(\theta^\star, \cdot)$

$\eta(\theta_{t+1}, \cdot)$

$a_t$

$a$
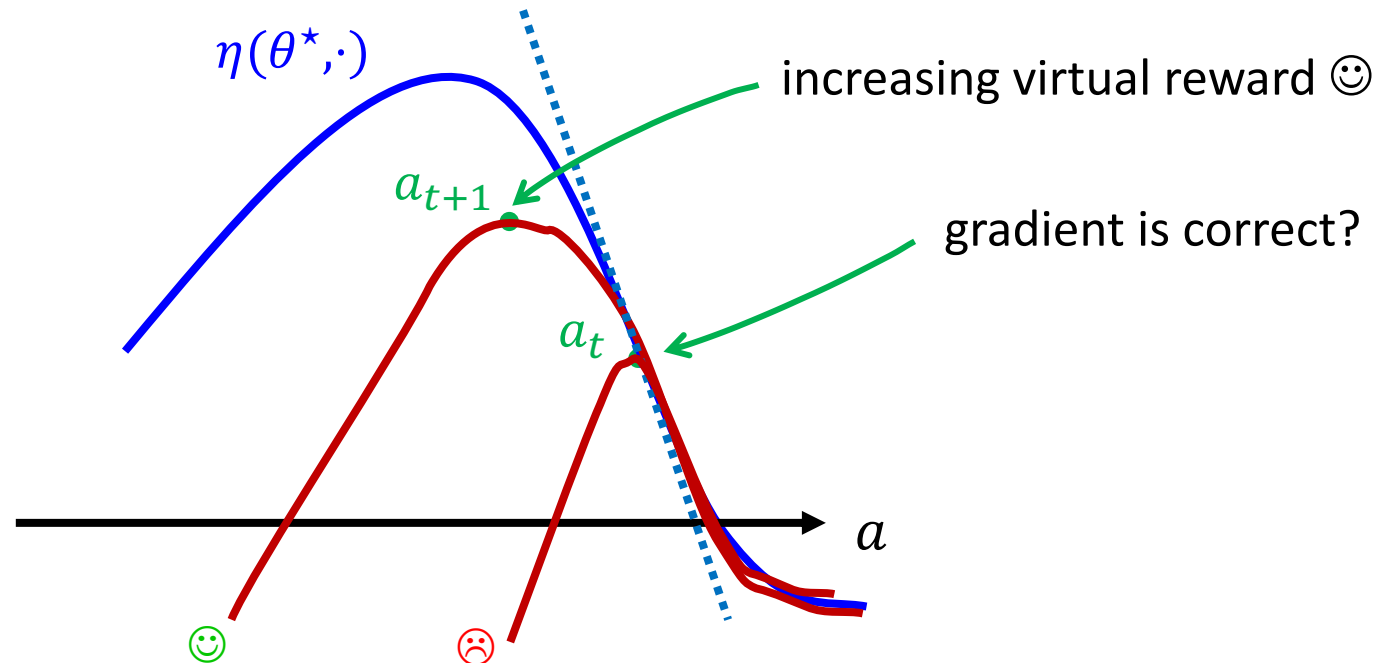
# Embrace Virtual Curvature

- Need the online learner to work harder to guarantee an increasing virtual reward

- Estimating the curvature: learn $\theta_t$ such that

  1. $\eta(\theta_t, a_t) \approx \eta(\theta^\star, a_t)$
  2. $\nabla_a \eta(\theta_t, a_t) \approx \nabla_a \eta(\theta^\star, a_t)$
  3. $\nabla_a^2 \eta(\theta_t, a_t) \approx \nabla_a^2 \eta(\theta^\star, a_t)$
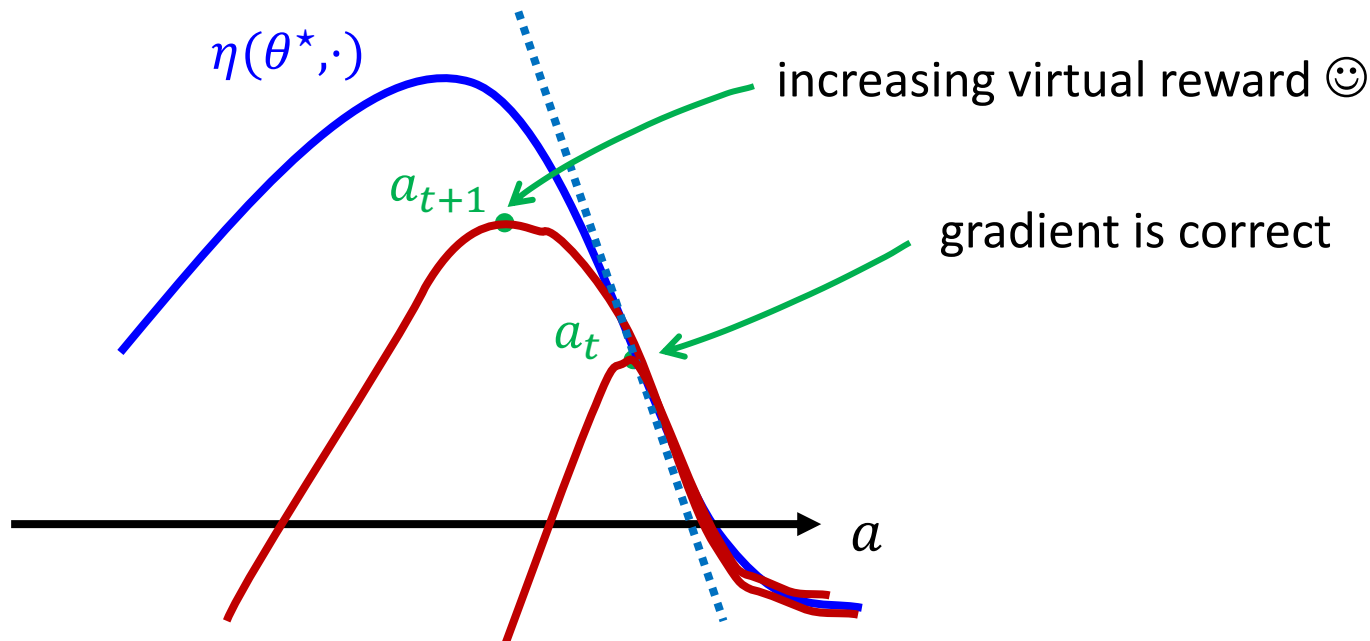
# Virtual Improvements With Curvature Estimates

- Assume $\eta(\theta, \cdot)$ is smooth (with bounded 2$^{\text{nd}}$ derivative)

- $\eta(\theta^\star, a_{t+1}) \approx \mathbb{E}_{\theta_{t+1}}[\eta(\theta_{t+1}, a_{t+1})]$

$$\geq \mathbb{E}_{\theta_{t+1}}[\eta(\theta_{t+1}, a_t)] + \Omega(||\mathbb{E}_{\theta_{t+1}}[\nabla\eta(\theta_{t+1}, a_t)]||_2^2)$$

$$\approx \eta(\theta^\star, a_t) + \Omega(||\nabla\eta(\theta^\star, a_t)||_2^2)$$

OL guarantee

$\eta(\theta^\star, \cdot)$

increasing virtual reward ☺

$a_{t+1}$

$a_t$

gradient is correct

"descent lemma":

$$\max_a f(a) \geq f(a_0) + \Omega(||\nabla f(a_0)||_2^2)$$

$a$

# Virtual Improvements With Curvature Estimates

- Consider online learning problem with the idealized $\overline{\ell}_t$

$$\overline{\ell}_t(\theta) = \big(\eta(\theta, a_t) - \eta(\theta^\star, a_t)\big)^2 + \big(\eta(\theta, a_{t-1}) - \eta(\theta^\star, a_{t-1})\big)^2$$
$$+ ||\nabla\eta(\theta, a_{t-1}) - \nabla\eta(\theta^\star, a_{t-1})||_2^2$$

- ViOL (**Vi**rtual Ascent with **O**nline Model **L**earner)
1. Use OL to minimize losses $\overline{\ell}_t$ (assuming $\overline{\ell}_t$ is accessible) and get a distribution of $\theta_t$
2. Take $a_t = \text{argmax}_a \ \mathbb{E}_{\theta_t}[\eta(\theta_t, a)]$

Lemma: If online learning for $\overline{\ell}_t$ has regret

$$\mathbb{E}\left[\sum_{t=1}^T \overline{\ell}_t(\theta_t) - \min_\theta \sum_{t=1}^T \overline{\ell}_t(\theta)\right] = \mathbb{E}\left[\sum_{t=1}^T \overline{\ell}_t(\theta_t)\right] = o(T)$$

Then $a_t$ converges to a critical point of the reward $\eta(\theta^\star, \cdot)$

# Learning Gradients With Model Extrapolation

$$\bar{\ell}_t(\theta) = \big(\eta(\theta, a_t) - \eta(\theta^\star, a_t)\big)^2 + \big(\eta(\theta, a_{t-1}) - \eta(\theta^\star, a_{t-1})\big)^2$$
$$+ ||\nabla\eta(\theta, a_{t-1}) - \underbrace{\nabla\eta(\theta^\star, a_{t-1})}||_2^2$$

<span style="color:red">not observed</span>

- $||\nabla\eta(\theta, a) - \nabla\eta(\theta^\star, a)||_2^2 = \mathbb{E}_u[\langle\nabla\eta(\theta, a) - \nabla\eta(\theta^\star, a), u\rangle^2]$

where $u \sim \mathcal{N}(0, I)$

- Directional gradient $\langle\nabla\eta(\theta^\star, a), u\rangle$ can be computed by two actions

$$\langle\nabla\eta(\theta^\star, a), u\rangle \approx \frac{\eta(\theta^\star, a + \alpha u) - \eta(\theta^\star, a)}{\alpha} \quad (\alpha \to 0)$$

- Similarly to Johnson–Lindenstrauss, it requires complexity($\Theta$) samples of $u$.
- Zero order optimization requires $\Omega(d)$ samples.

# Algorithm and Theorem

$$\ell_t(\theta) = \big(\eta(\theta, a_t) - \eta(\theta^\star, a_t)\big)^2 + \big(\eta(\theta, a_{t-1}) - \eta(\theta^\star, a_{t-1})\big)^2$$
$$+ \langle \nabla\eta(\theta, a_{t-1}) - \nabla\eta(\theta^\star, a_{t-1}), u_t \rangle^2$$

- ViOL (**Vi**rtual Ascent with **O**nline Model **L**earner)
1. Sample $u_t \sim \mathcal{N}(0, I)$
2. Use OL to minimize losses $\ell_t$ and get a distribution of $\theta_t$
3. Take $a_t = \text{argmax}_a \, \mathbb{E}_{\theta_t}[\eta(\theta_t, a)]$

- Theorem (informal): Under Lipschitz assumptions on $\eta$, ViOL converges to a $\epsilon$-approximate local maxima in $\tilde{O}(R(\Theta)\epsilon^{-8})$.

# Instantiations

- Linear bandit with structured model family: $\eta(\theta, a) = \theta^\top a$
  - $\Theta$ is finite: $\text{poly}(\log |\Theta|)$ sample complexity
  - $\Theta$ contains $s$-sparse vectors: $\text{poly}(s, \log d)$ sample complexity
  - local maximum are global because $\eta(\theta^\star, \cdot)$ is concave.
  - only hold for <span style="color:red">deterministic reward</span>

- Neural net bandit: $\eta(W, a) = w_2^\top \sigma(W_1 a)$
  - assume $O(1)$ norms bounds on $||w_2||_1, \|W_1\|_{\infty \to \infty}$
  - $R(W) \leq \tilde{O}(1)$
  - sample complexity for local max = $\tilde{O}(1)$
  - Local maximum are global for input-concave neural nets

# First-cut Extension to RL

| RL | Bandit with Continous Actions |
| --- | --- |
| Dynamics $T_\theta$ | Model parameter $\theta$ |
| Policy $\pi_\psi$ | Action $a$ |
| Total return $\eta(T_\theta, \pi_\psi)$ | Reward function $\eta(\theta, a)$ |

- Caveat: $\{\eta(T_\theta, \cdot): \theta \in \Theta\}$ has high complexity
- A result for stochastic policies (with many Lipschitz conditions)

$$|\eta(\theta, \psi) - \eta(\theta^\star, \psi)|^2 \lesssim \mathbb{E}_{s,a \sim T_{\theta^\star}, \pi_\psi}[\|T_\theta(s,a) - T_{\theta^\star}(s,a)\|^2]$$

$$\|\nabla\eta(\theta, \psi) - \nabla\eta(\theta^\star, \psi)\|^2 \lesssim \mathbb{E}_{s,a \sim T_{\theta^\star}, \pi_\psi}[\|T_\theta(s,a) - T_{\theta^\star}(s,a)\|^2]$$

- Summary: online learning $T_\theta$ implies predicting the curvature of $\eta$

# Summary

- Global convergence for nonlinear models is <span style="color:red">statistically</span> intractable

- ViOL: convergence to a local maximum with sample complexity that only depends on the model class complexity

- Check out our paper for more detail: https://arxiv.org/abs/2102.04168

# Thank you for your attention ☺