

# **Policy Gradients: Optimality**

# Today

- Recap:
  - NPG convergence proof wrap up
  - remember compatible function approximation

# Today

- Recap:
  - NPG convergence proof wrap up
  - remember compatible function approximation
- Today:
  - What about function approximation?  
log linear policy classes and neural policy classes
  - PG methods have stronger guarantees when we have errors.

Recap

# Things to remember

For all  $\pi, \pi', s_0$ :

$$V^\pi(s_0) - V^{\pi'}(s_0) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{s_0}^\pi} \mathbb{E}_{a \sim \pi(\cdot | s)} [A^{\pi'}(s, a)]$$

$$\nabla_\theta J(\theta) := \frac{1}{1 - \gamma} \mathbb{E}_{s, a \sim d^{\pi_\theta}} [\nabla_\theta \ln \pi_\theta(a | s) Q^{\pi_\theta}(s, a)]$$

Today: we will use  $d_{s_0}^\pi$  for a state distribution measure.

(it should be clear from context how we use it).

$$d_{s_0}^\pi(s) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}(s_h = s | s_0, \pi)$$

$$d_{s_0}^\pi(s, a) = (1 - \gamma) \sum_{h=0}^{\infty} \gamma^h \mathbb{P}(s_h = s, a_h = a | s_0, \pi)$$

$$V^\pi(\mu) = E_{s \sim \mu}[V^\pi(s)]$$

$$d_\mu^\pi(s) = E_{s_0 \sim \mu}[d_{s_0}^\pi(s)]$$

# The Natural Policy Gradient

- Define  $\mathcal{F}_\rho^\theta$  as the (average) Fisher matrix on the family of distributions  $\{\pi_\theta(\cdot | s) | s \in \mathcal{S}\}$  as:  
$$\mathcal{F}_\rho^\theta := E_{s \sim d_\rho^{\pi_\theta}} E_{a \sim \pi_\theta(\cdot | s)} \left[ (\nabla \log \pi_\theta(a | s)) \nabla \log \pi_\theta(a | s)^\top \right] .$$
- The NPG algorithm performs gradient updates in this induced geometry:  
$$\theta^{(t+1)} = \theta^{(t)} + \eta F_\rho(\theta^{(t)})^\dagger \nabla_\theta V^{(t)}(\rho),$$
  
where  $M^\dagger$  denotes the Moore-Penrose pseudoinverse of  $M$ .
- Idea:
  - ‘stretch’ the corners of the simplex out to travel faster  
(as opposed to the log-barrier which keeps us away)

# NPG softmax case

(NPG as “soft” policy iteration)

- **Lemma:** (Softmax NPG as soft policy iteration) The NPG update is:

$$\theta^{(t+1)} = \theta^{(t)} + \frac{\eta}{1 - \gamma} A^{(t)}$$

- and so:

$$\pi^{(t+1)}(a | s) = \pi^{(t)}(a | s) \frac{\exp(\eta A^{(t)}(s, a)/(1 - \gamma))}{Z_t(s)},$$

where  $Z_t(s) = \sum_a \pi^{(t)}(a | s) \exp(\eta A^{(t)}(s, a)/(1 - \gamma))$ .

# Global convergence for NPG

- **Theorem:** Params:  $\theta^{(0)} = 0$  and  $\eta > 0$ . For all  $\rho$  and  $T > 0$ , we have:

$$V^{(T)}(\rho) \geq V^*(\rho) - \frac{\log A}{\eta T} - \frac{1}{(1 - \gamma)^2 T}.$$

- Setting  $\eta \geq (1 - \gamma)^2 \log A$ , NPG finds an  $\epsilon$ -opt policy when  $T \geq \frac{2}{(1 - \gamma)^2 \epsilon}$ .
- Iteration complexity has:
  - No dimension dependence (no dependence on  $S, A$ )
  - No dependence on start state measure  $\rho$  (and no “dist mismatch factor”)
  - No ‘flat gradient’ problem
- What about approx/estimation errors? (next lecture)



# Improvement Lower Bound

- **Lemma:** For the iterates  $\pi^{(t)}$  generated by the NPG, we have for all distributions  $\mu$ :  
$$V^{(t+1)}(\mu) - V^{(t)}(\mu) \geq \frac{(1 - \gamma)}{\eta} E_{s \sim \mu} \log Z_t(s) \geq 0.$$

# NPG Conv. Proof, Part 1

- $d^\star$  as shorthand for  $d_\rho^\star$ ;  $\pi_s$  as shorthand for the vector of  $\pi(\cdot | s)$
- By the performance difference lemma,

$$V^{\pi^\star}(\rho) - V^{(t)}(\rho) = \frac{1}{1-\gamma} E_{s \sim d^\star} \sum_a \pi^\star(a | s) A^{(t)}(s, a)$$

$$= \frac{1}{\eta} E_{s \sim d^\star} \sum_a \pi^\star(a | s) \log \frac{\pi^{(t+1)}(a | s) Z_t(s)}{\pi^{(t)}(a | s)}$$

$$= \frac{1}{\eta} E_{s \sim d^\star} \left( \text{KL}(\pi_s^\star || \pi_s^{(t)}) - \text{KL}(\pi_s^\star || \pi_s^{(t+1)}) + \sum_a \pi^\star(a | s) \log Z_t(s) \right)$$

$$= \frac{1}{\eta} E_{s \sim d^\star} \left( \text{KL}(\pi_s^\star || \pi_s^{(t)}) - \text{KL}(\pi_s^\star || \pi_s^{(t+1)}) + \log Z_t(s) \right),$$

# NPG Conv. Proof, Part 2

- By the improvement lemma  $V^{(t+1)}(\rho) \geq V^{(t)}(\rho)$ . Hence,

$$V^{\pi^*}(\rho) - V^{(T-1)}(\rho) \leq \frac{1}{T} \sum_{t=0}^{T-1} (V^{\pi^*}(\rho) - V^{(t)}(\rho))$$

$$= \frac{1}{\eta T} \sum_{t=0}^{T-1} E_{s \sim d^*} (\text{KL}(\pi_s^* || \pi_s^{(t)}) - \text{KL}(\pi_s^* || \pi_s^{(t+1)})) + \frac{1}{\eta T} \sum_{t=0}^{T-1} E_{s \sim d^*} \log Z_t(s)$$

$$\leq \frac{E_{s \sim d^*} \text{KL}(\pi_s^* || \pi^{(0)})}{\eta T} + \frac{1}{\eta T} \sum_{t=0}^{T-1} E_{s \sim d^*} \log Z_t(s).$$

# NPG Conv. Proof, Part 2

- By the improvement lemma  $V^{(t+1)}(\rho) \geq V^{(t)}(\rho)$ . Hence,

$$V^{\pi^*}(\rho) - V^{(T-1)}(\rho) \leq \frac{1}{T} \sum_{t=0}^{T-1} (V^{\pi^*}(\rho) - V^{(t)}(\rho))$$

$$= \frac{1}{\eta T} \sum_{t=0}^{T-1} E_{s \sim d^*} (\text{KL}(\pi_s^* || \pi_s^{(t)}) - \text{KL}(\pi_s^* || \pi_s^{(t+1)})) + \frac{1}{\eta T} \sum_{t=0}^{T-1} E_{s \sim d^*} \log Z_t(s)$$

$$\leq \frac{E_{s \sim d^*} \text{KL}(\pi_s^* || \pi^{(0)})}{\eta T} + \frac{1}{\eta T} \sum_{t=0}^{T-1} E_{s \sim d^*} \log Z_t(s).$$

- By the improvement lemma (applied with  $d^*$  as the distribution), we have:

$$\frac{1}{\eta} E_{s \sim d^*} \log Z_t(s) \leq \frac{1}{1 - \gamma} \left( V^{(t+1)}(d^*) - V^{(t)}(d^*) \right)$$

which gives us a bound on  $E_{s \sim d^*} \log Z_t(s)$ .

# NPG Conv. Proof, Part 3

$$\begin{aligned} V^{\pi^*}(\rho) - V^{(T-1)}(\rho) &\leq \frac{E_{s \sim d^*} \text{KL}(\pi_s^* || \pi^{(0)})}{\eta T} + \frac{1}{\eta T} \sum_{t=0}^{T-1} E_{s \sim d^*} \log Z_t(s) \\ &\leq \frac{E_{s \sim d^*} \text{KL}(\pi_s^* || \pi^{(0)})}{\eta T} + \frac{1}{(1-\gamma)T} \sum_{t=0}^{T-1} \left( V^{(t+1)}(d^*) - V^{(t)}(d^*) \right) \\ &= \frac{E_{s \sim d^*} \text{KL}(\pi_s^* || \pi^{(0)})}{\eta T} + \frac{V^{(T)}(d^*) - V^{(0)}(d^*)}{(1-\gamma)T} \\ &\leq \frac{\log A}{\eta T} + \frac{1}{(1-\gamma)^2 T}. \end{aligned}$$

# Policy Gradient: Examples of Policy Parameterization (discrete actions)

## 1. Softmax Policy for Tabular MDPs:

$$\theta_{s,a} \in \mathbb{R}, \forall s, a \in S \times A$$

$$\pi_{\theta}(a | s) = \frac{\exp(\theta_{s,a})}{\sum_{a'} \exp(\theta_{s,a'})}$$

## 2. Log Linear Policy (e.g., for linear MDPs):

Feature vector  $\phi(s, a) \in \mathbb{R}^d$ , and parameter  $\theta \in \mathbb{R}^d$

$$\pi_{\theta}(a | s) = \frac{\exp(\theta^{\top} \phi(s, a))}{\sum_{a'} \exp(\theta^{\top} \phi(s, a'))}$$

## 3. Neural Policy:

Neural network  
 $f_{\theta} : S \times A \mapsto \mathbb{R}$

$$\pi_{\theta}(a | s) = \frac{\exp(f_{\theta}(s, a))}{\sum_{a'} \exp(f_{\theta}(s, a'))}$$

# NPG & Compatible Function Approximation

- Let  $w^\star$  denote the following minimizer of the “compatible function approximation” error:

$$w^\star \in \operatorname{argmin}_w E_{s \sim d_\mu^{\pi_\theta}} E_{a \sim \pi_\theta(\cdot|s)} \left[ \left( A^{\pi_\theta}(s, a) - w \cdot \nabla_\theta \log \pi_\theta(a|s) \right)^2 \right]$$

- Lemma:** Let  $\widehat{A}^{\pi_\theta}(s, a)$  be the best linear predictor of  $A^{\pi_\theta}(s, a)$  using  $\nabla_\theta \log \pi_\theta(a|s)$ , i.e.  $\widehat{A}^{\pi_\theta}(s, a) := w^\star \cdot \nabla_\theta \log \pi_\theta(a|s)$ . We have:

$$\nabla_\theta V^{\pi_\theta}(\mu) = \frac{1}{1-\gamma} E_{s \sim d_\mu^{\pi_\theta}} E_{a \sim \pi_\theta(\cdot|s)} \left[ \nabla_\theta \log \pi_\theta(a|s) \widehat{A}^{\pi_\theta}(s, a) \right]$$

We can use  $\widehat{A}^{\pi_\theta}(s, a)$  instead of  $A^{\pi_\theta}(s, a)$ .

- Lemma:** We have that  $F_\mu(\theta)^\dagger \nabla_\theta V^\theta(\mu) = \frac{1}{1-\gamma} w^\star$ ,

The NPG direction is the weights  $w^\star$

# Today:

Natural Policy Gradient and Approximation



# Examples:

NPG and variants for log-linear policy classes

# NPG & Log Linear Policy Classes

- Feature vector  $\phi(s, a) \in \mathbb{R}^d$ ,  $\pi_\theta(a | s) = \frac{\exp(\theta^\top \phi(s, a))}{\sum_{a'} \exp(\theta^\top \phi(s, a'))}$

# NPG & Log Linear Policy Classes

- Feature vector  $\phi(s, a) \in \mathbb{R}^d$ ,  $\pi_\theta(a | s) = \frac{\exp(\theta^\top \phi(s, a))}{\sum_{a'} \exp(\theta^\top \phi(s, a'))}$
- We have:  
 $\nabla_\theta \log \pi_\theta(a | s) = \bar{\phi}_{s,a}^\theta$ , where  $\bar{\phi}_{s,a}^\theta = \phi_{s,a} - E_{a' \sim \pi_\theta(\cdot | s)}[\phi_{s,a'}]$ .

# NPG & Log Linear Policy Classes

- Feature vector  $\phi(s, a) \in \mathbb{R}^d$ ,  $\pi_\theta(a | s) = \frac{\exp(\theta^\top \phi(s, a))}{\sum_{a'} \exp(\theta^\top \phi(s, a'))}$
- We have:  
 $\nabla_\theta \log \pi_\theta(a | s) = \bar{\phi}_{s,a}^\theta$ , where  $\bar{\phi}_{s,a}^\theta = \phi_{s,a} - E_{a' \sim \pi_\theta(\cdot | s)}[\phi_{s,a'}]$ .
- The NPG update:  
 $\theta \leftarrow \theta + \eta w_\star$ ,  $w_\star \in \operatorname{argmin}_w E_{s \sim d_\rho^{\pi_\theta}, a \sim \pi_\theta(\cdot | s)} \left[ \left( A^{\pi_\theta}(s, a) - w \cdot \bar{\phi}_{s,a}^\theta \right)^2 \right]$ .

# NPG & Log Linear Policy Classes

- Feature vector  $\phi(s, a) \in \mathbb{R}^d$ ,  $\pi_\theta(a | s) = \frac{\exp(\theta^\top \phi(s, a))}{\sum_{a'} \exp(\theta^\top \phi(s, a'))}$
- We have:  
 $\nabla_\theta \log \pi_\theta(a | s) = \bar{\phi}_{s,a}^\theta$ , where  $\bar{\phi}_{s,a}^\theta = \phi_{s,a} - E_{a' \sim \pi_\theta(\cdot | s)}[\phi_{s,a'}]$ .
- The NPG update:  
 $\theta \leftarrow \theta + \eta w_\star$ ,  $w_\star \in \operatorname{argmin}_w E_{s \sim d_\rho^{\pi_\theta}, a \sim \pi_\theta(\cdot | s)} \left[ \left( A^{\pi_\theta}(s, a) - w \cdot \bar{\phi}_{s,a}^\theta \right)^2 \right]$ .
- Equivalently, for the same  $w_\star$ ,  
 $\pi(a | s) \leftarrow \frac{\pi(a | s) \exp(w_\star \cdot \phi_{s,a})}{Z_s}$   
( $Z_s$  is the normalizing constant.) Using  $\bar{\phi}$  or  $\phi$  result in the same update for  $\pi$ .

# Q-NPG: use Q rather A

(a little nice to interpret for analysis)

- Still log linear class.

# Q-NPG: use Q rather A

(a little nice to interpret for analysis)

- Still log linear class.
- The Q-NPG update:

$$\theta \leftarrow \theta + \eta w_{\star}, \quad w_{\star} \in \operatorname{argmin}_w E_{s \sim d_{\rho}^{\pi_{\theta}}, a \sim \pi_{\theta}(\cdot|s)} \left[ \left( Q^{\pi_{\theta}}(s, a) - w \cdot \phi_{s,a}^{\theta} \right)^2 \right].$$

# Q-NPG: use Q rather A

(a little nice to interpret for analysis)

- Still log linear class.

- The Q-NPG update:

$$\theta \leftarrow \theta + \eta w_{\star}, \quad w_{\star} \in \operatorname{argmin}_w E_{s \sim d_{\rho}^{\pi_{\theta}}, a \sim \pi_{\theta}(\cdot | s)} \left[ \left( Q^{\pi_{\theta}}(s, a) - w \cdot \phi_{s,a}^{\theta} \right)^2 \right].$$

- Equivalently, for the same  $w_{\star}$ ,

$$\pi(a | s) \leftarrow \frac{\pi(a | s) \exp(w_{\star} \cdot \phi_{s,a})}{Z_s}$$

( $Z_s$  is the normalizing constant.)



# Approximate Q-NPG

(e.g. we use samples to estimate Q)

- For a state-action distribution  $\nu$ , define:

$$L(w; \theta, \nu) := E_{s,a \sim \nu} \left[ (Q^{\pi_\theta}(s, a) - w \cdot \phi_{s,a})^2 \right].$$

# Approximate Q-NPG

(e.g. we use samples to estimate Q)

- For a state-action distribution  $\nu$ , define:

$$L(w; \theta, \nu) := E_{s,a \sim \nu} \left[ (Q^{\pi_\theta}(s, a) - w \cdot \phi_{s,a})^2 \right].$$

- With an on-policy state action measure starting with  $s_0, a_0 \sim \nu$ . Shorthand:

$$d^{(t)}(s, a) := d_\nu^{\pi^{(t)}}(s, a)$$

# Approximate Q-NPG

(e.g. we use samples to estimate Q)

- For a state-action distribution  $\nu$ , define:

$$L(w; \theta, \nu) := E_{s,a \sim \nu} \left[ (Q^{\pi_\theta}(s, a) - w \cdot \phi_{s,a})^2 \right].$$

- With an on-policy state action measure starting with  $s_0, a_0 \sim \nu$ . Shorthand:

$$d^{(t)}(s, a) := d_\nu^{\pi^{(t)}}(s, a)$$

- The approximate version:

$$\theta^{(t+1)} = \theta^{(t)} + \eta w^{(t)}, \text{ where } w^{(t)} \approx \operatorname{argmin}_{\|w\|_2 \leq W} L(w; \theta^{(t)}, d^{(t)}),$$

# Approximate Q-NPG

(e.g. we use samples to estimate Q)

- For a state-action distribution  $\nu$ , define:

$$L(w; \theta, \nu) := E_{s,a \sim \nu} \left[ (Q^{\pi_\theta}(s, a) - w \cdot \phi_{s,a})^2 \right].$$

- With an on-policy state action measure starting with  $s_0, a_0 \sim \nu$ . Shorthand:

$$d^{(t)}(s, a) := d_\nu^{\pi^{(t)}}(s, a)$$

- The approximate version:

$$\theta^{(t+1)} = \theta^{(t)} + \eta w^{(t)}, \text{ where } w^{(t)} \approx \operatorname{argmin}_{\|w\|_2 \leq W} L(w; \theta^{(t)}, d^{(t)}),$$

- Equivalently,

$$\pi^{(t+1)}(a | s) \leftarrow \frac{\pi^{(t)}(a | s) \exp(w^{(t)} \cdot \phi_{s,a})}{Z_s}$$

# Error Analysis of NPG

(and variants)

# NPG regret lemma

- Consider the update rule:  $\theta^{(t+1)} = \theta^{(t)} + \eta w^{(t)}$   
(starting with  $\pi^{(0)}$  being the uniform policy).

# NPG regret lemma

- Consider the update rule:  $\theta^{(t+1)} = \theta^{(t)} + \eta w^{(t)}$   
(starting with  $\pi^{(0)}$  being the uniform policy).
- **Lemma:** (NPG Regret Lemma)  
Fix any comparison policy  $\tilde{\pi}$  and a state distribution  $\rho$ .  
Assume  $\log \pi_\theta(a | s)$  (for all  $s, a$ ) is a  $\beta$ -smooth function of  $\theta$ .  
Consider **an arbitrary sequence of weights**  $w^{(0)}, \dots, w^{(T)}$ , s.t.  $\|w^{(t)}\|_2 \leq W$ . Define:  
 $\text{err}_t = E_{s \sim \tilde{d}} E_{a \sim \tilde{\pi}(\cdot | s)} [A^{(t)}(s, a) - w^{(t)} \cdot \nabla_\theta \log \pi^{(t)}(a | s)]$ .

We have that:

$$\min_{t < T} \left\{ V^{\tilde{\pi}}(\rho) - V^{(t)}(\rho) \right\} \leq \frac{1}{1 - \gamma} \left( W \sqrt{\frac{2\beta \log A}{T}} + \frac{1}{T} \sum_{t=0}^{T-1} \text{err}_t \right).$$

(using  $\eta = \sqrt{2 \log A / (\beta W^2 T)}$ )

# NPG regret lemma

- Consider the update rule:  $\theta^{(t+1)} = \theta^{(t)} + \eta w^{(t)}$   
(starting with  $\pi^{(0)}$  being the uniform policy).
- **Lemma:** (NPG Regret Lemma)  
Fix any comparison policy  $\tilde{\pi}$  and a state distribution  $\rho$ .  
Assume  $\log \pi_\theta(a | s)$  (for all  $s, a$ ) is a  $\beta$ -smooth function of  $\theta$ .  
Consider **an arbitrary sequence of weights**  $w^{(0)}, \dots, w^{(T)}$ , s.t.  $\|w^{(t)}\|_2 \leq W$ . Define:  
 $\text{err}_t = E_{s \sim \tilde{d}} E_{a \sim \tilde{\pi}(\cdot | s)} [A^{(t)}(s, a) - w^{(t)} \cdot \nabla_\theta \log \pi^{(t)}(a | s)]$ .

We have that:

$$\min_{t < T} \left\{ V^{\tilde{\pi}}(\rho) - V^{(t)}(\rho) \right\} \leq \frac{1}{1 - \gamma} \left( W \sqrt{\frac{2\beta \log A}{T}} + \frac{1}{T} \sum_{t=0}^{T-1} \text{err}_t \right).$$

(using  $\eta = \sqrt{2 \log A / (\beta W^2 T)}$ )

- Proof: Mirror descent style of analysis + Perf. Difference Lemma



# Approximate Q-NPG

(e.g. we use samples to estimate Q)

- The approximate version:

$$\theta^{(t+1)} = \theta^{(t)} + \eta w^{(t)}, \text{ where } w^{(t)} \approx \operatorname{argmin}_{\|w\|_2 \leq W} L(w; \theta^{(t)}, d^{(t)}),$$

# Approximate Q-NPG

(e.g. we use samples to estimate Q)

- The approximate version:

$$\theta^{(t+1)} = \theta^{(t)} + \eta w^{(t)}, \text{ where } w^{(t)} \approx \operatorname{argmin}_{\|w\|_2 \leq W} L(w; \theta^{(t)}, d^{(t)}),$$

- Error Decomposition:

$$L(w^{(t)}; \theta^{(t)}, d^{(t)}) = \underbrace{L(w^{(t)}; \theta^{(t)}, d^{(t)}) - L(w_{\star}^{(t)}; \theta^{(t)}, d^{(t)})}_{\text{Excess risk}} + \underbrace{L(w_{\star}^{(t)}; \theta^{(t)}, d^{(t)})}_{\text{Approximation error}}$$

where  $w_{\star}^{(t)} \in \operatorname{argmin}_{\|w\|_2 \leq W} L(w; \theta^{(t)}, d^{(t)})$

# Q-NPG Conv Rate w/ Estimation Error (no approx error)

- Suppose **no approx error**:  $L(w_{\star}^{(t)}; \theta^{(t)}, d^{(t)}) = 0$

Suppose the **excess risk**:

$$L(w^{(t)}; \theta^{(t)}, d^{(t)}) - L(w_{\star}^{(t)}; \theta^{(t)}, d^{(t)}) \leq \epsilon_{\text{stat}},$$

# Q-NPG Conv Rate w/ Estimation Error (no approx error)

- Suppose **no approx error**:  $L(w_{\star}^{(t)}; \theta^{(t)}, d^{(t)}) = 0$   
Suppose the **excess risk**:  
 $L(w^{(t)}; \theta^{(t)}, d^{(t)}) - L(w_{\star}^{(t)}; \theta^{(t)}, d^{(t)}) \leq \epsilon_{\text{stat}}$
- **Conditioning**: suppose  $\|\phi_{s,a}\|_2 \leq 1$  and, **for the initial measure  $\nu$** ,  
 $\sigma_{\min}\left(E_{s,a \sim \nu}[\phi_{s,a} \phi_{s,a}^{\top}]\right) = \lambda_{\min}, \quad \kappa = 1/\lambda.$

# Q-NPG Conv Rate w/ Estimation Error (no approx error)

- Suppose **no approx error**:  $L(w_{\star}^{(t)}; \theta^{(t)}, d^{(t)}) = 0$

Suppose the **excess risk**:

$$L(w^{(t)}; \theta^{(t)}, d^{(t)}) - L(w_{\star}^{(t)}; \theta^{(t)}, d^{(t)}) \leq \epsilon_{\text{stat}},$$

- **Conditioning**: suppose  $\|\phi_{s,a}\|_2 \leq 1$  and, **for the initial measure  $\nu$** ,

$$\sigma_{\min}\left(E_{s,a \sim \nu}[\phi_{s,a} \phi_{s,a}^{\top}]\right) = \lambda_{\min}, \quad \kappa = 1/\lambda.$$

- **Theorem**: Fix any state distribution  $\rho$ ; **any comparator policy  $\pi^{\star}$**  (not necessarily optimal).

With  $\eta$  set appropriately and under the above assumptions, we have that:

$$E \left[ \min_{t < T} \left\{ V^{\pi^{\star}}(\rho) - V^{(t)}(\rho) \right\} \right] \leq \frac{BW}{1-\gamma} \sqrt{\frac{2 \log A}{T}} + \sqrt{\frac{4A}{(1-\gamma)^3} \left( \kappa \cdot \epsilon_{\text{stat}} \right)}$$

# Q-NPG Conv Rate with Approx+Est. Errors

- Suppose the **excess risk** and **approx error** are bounded as:

$$L(w^{(t)}; \theta^{(t)}, d^{(t)}) - L(w_{\star}^{(t)}; \theta^{(t)}, d^{(t)}) \leq \epsilon_{\text{stat}},$$

$$L(w_{\star}^{(t)}; \theta^{(t)}, d^{(t)}) \leq \epsilon_{\text{approx}},$$

# Q-NPG Conv Rate with Approx+Est. Errors

- Suppose the **excess risk** and **approx error** are bounded as:

$$L(w^{(t)}; \theta^{(t)}, d^{(t)}) - L(w_{\star}^{(t)}; \theta^{(t)}, d^{(t)}) \leq \epsilon_{\text{stat}},$$

$$L(w_{\star}^{(t)}; \theta^{(t)}, d^{(t)}) \leq \epsilon_{\text{approx}},$$

- Conditioning: suppose  $\|\phi_{s,a}\|_2 \leq 1$  and, **for the initial measure  $\nu$ ,**

$$\sigma_{\min}\left(E_{s,a \sim \nu}[\phi_{s,a} \phi_{s,a}^{\top}]\right) = \lambda_{\min}, \quad \kappa = 1/\lambda.$$

# Q-NPG Conv Rate with Approx+Est. Errors

- Suppose the **excess risk** and **approx error** are bounded as:

$$L(w^{(t)}; \theta^{(t)}, d^{(t)}) - L(w_{\star}^{(t)}; \theta^{(t)}, d^{(t)}) \leq \epsilon_{\text{stat}},$$

$$L(w_{\star}^{(t)}; \theta^{(t)}, d^{(t)}) \leq \epsilon_{\text{approx}},$$

- Conditioning: suppose  $\|\phi_{s,a}\|_2 \leq 1$  and, **for the initial measure  $\nu$** ,

$$\sigma_{\min}\left(E_{s,a \sim \nu}[\phi_{s,a} \phi_{s,a}^{\top}]\right) = \lambda_{\min}, \quad \kappa = 1/\lambda.$$

- **Theorem:** Fix any state distribution  $\rho$ ; **any comparator policy  $\pi^{\star}$**  (not necessarily optimal).

With  $\eta$  set appropriately and under the above assumptions, we have that:

$$E \left[ \min_{t < T} \left\{ V^{\pi^{\star}}(\rho) - V^{(t)}(\rho) \right\} \right] \\ \leq \frac{BW}{1-\gamma} \sqrt{\frac{2 \log A}{T}} + \sqrt{\frac{4A}{(1-\gamma)^3} \left( \kappa \cdot \epsilon_{\text{stat}} + \left\| \frac{d^{\star}}{\nu} \right\|_{\infty} \cdot \epsilon_{\text{approx}} \right)}$$



# NPG & Neural Policy Classes

- Neural net  $f_\theta : S \times A \mapsto \mathbb{R}$ , Policy:

$$\pi_\theta(a | s) = \frac{\exp(f_\theta(s, a))}{\sum_{a'} \exp(f_\theta(s, a'))}$$

# NPG & Neural Policy Classes

- Neural net  $f_\theta : S \times A \mapsto \mathbb{R}$ , Policy:

$$\pi_\theta(a | s) = \frac{\exp(f_\theta(s, a))}{\sum_{a'} \exp(f_\theta(s, a'))}$$

- We have:

$$\nabla_\theta \log \pi_\theta(a | s) = g_\theta(s, a), \quad \text{where} \quad g_\theta(s, a) = \nabla_\theta f_\theta(s, a) - E_{a' \sim \pi_\theta(\cdot | s)}[\nabla_\theta f_\theta(s, a')].$$

# NPG & Neural Policy Classes

- Neural net  $f_\theta : S \times A \mapsto \mathbb{R}$ , Policy:

$$\pi_\theta(a | s) = \frac{\exp(f_\theta(s, a))}{\sum_{a'} \exp(f_\theta(s, a'))}$$

- We have:

$$\nabla_\theta \log \pi_\theta(a | s) = g_\theta(s, a), \quad \text{where } g_\theta(s, a) = \nabla_\theta f_\theta(s, a) - E_{a' \sim \pi_\theta(\cdot | s)}[\nabla_\theta f_\theta(s, a')].$$

- The NPG update rule is:

$$\theta \leftarrow \theta + \eta w_\star, \quad w_\star \in \operatorname{argmin}_w E_{s \sim d_\rho^{\pi_\theta}, a \sim \pi_\theta(\cdot | s)} \left[ \left( A^{\pi_\theta}(s, a) - w \cdot g_\theta(s, a) \right)^2 \right]$$

# NPG & Neural Policy Classes

- Neural net  $f_\theta : S \times A \mapsto \mathbb{R}$ , Policy:

$$\pi_\theta(a | s) = \frac{\exp(f_\theta(s, a))}{\sum_{a'} \exp(f_\theta(s, a'))}$$

- We have:

$$\nabla_\theta \log \pi_\theta(a | s) = g_\theta(s, a), \quad \text{where } g_\theta(s, a) = \nabla_\theta f_\theta(s, a) - E_{a' \sim \pi_\theta(\cdot | s)}[\nabla_\theta f_\theta(s, a')].$$

- The NPG update rule is:

$$\theta \leftarrow \theta + \eta w_\star, \quad w_\star \in \operatorname{argmin}_w E_{s \sim d_\rho^{\pi_\theta}, a \sim \pi_\theta(\cdot | s)} \left[ \left( A^{\pi_\theta}(s, a) - w \cdot g_\theta(s, a) \right)^2 \right]$$