

Agenda

- TD (bootstrapping) with approximation
- Off-policy with approximation

Value function approximation II

TD

Konpat Preechakul
Chulalongkorn University
September 2019

Recap

Approximate value, approximate policy

We can show that:

$$v_\pi(s) \geq v^*(s) - \frac{2\epsilon}{1 - \gamma}$$

- $v^*(s)$ is the optimal policy performance
- $v_\pi(s)$ is our policy (using $q_\theta(s, a)$)
- ϵ the maximum error between $q_\theta(s, a)$ and $q^*(s, a)$
- Our policy has a lower bound depending on the error!

Intuitive interpretation

$$v_\pi(s) \geq v^*(s) - \frac{2\epsilon}{1-\gamma}$$

- Per-step error:

$$\textcircled{1} \quad v^*(s) - q^*(s, \pi(s)) \leq 2\epsilon$$

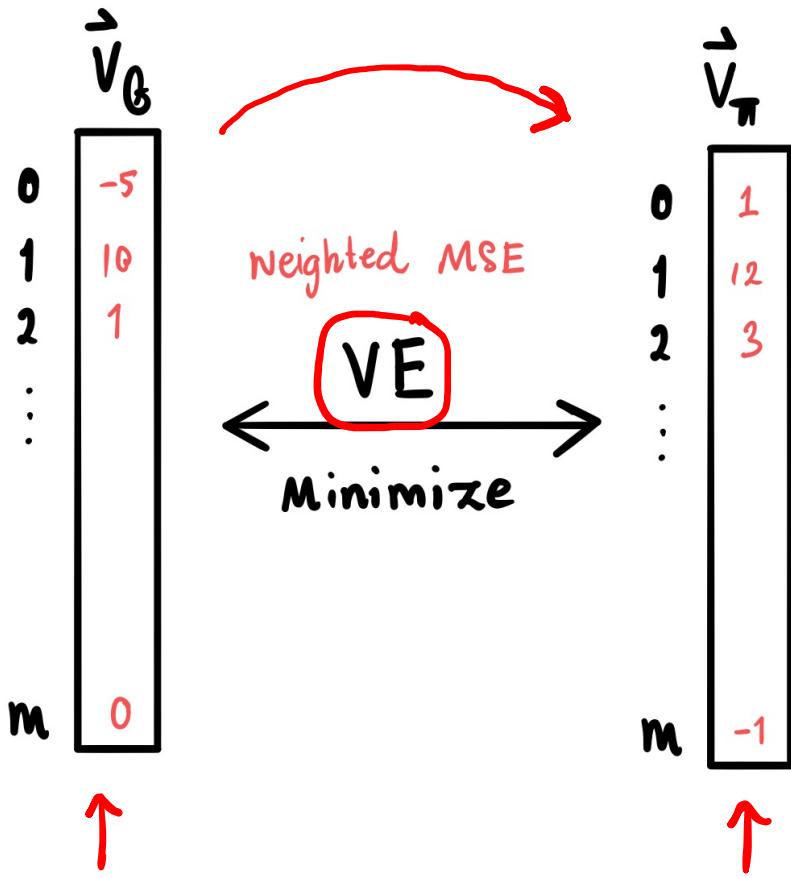
~~$q^*(s, \pi^*(s))$~~

- Trajectory error:

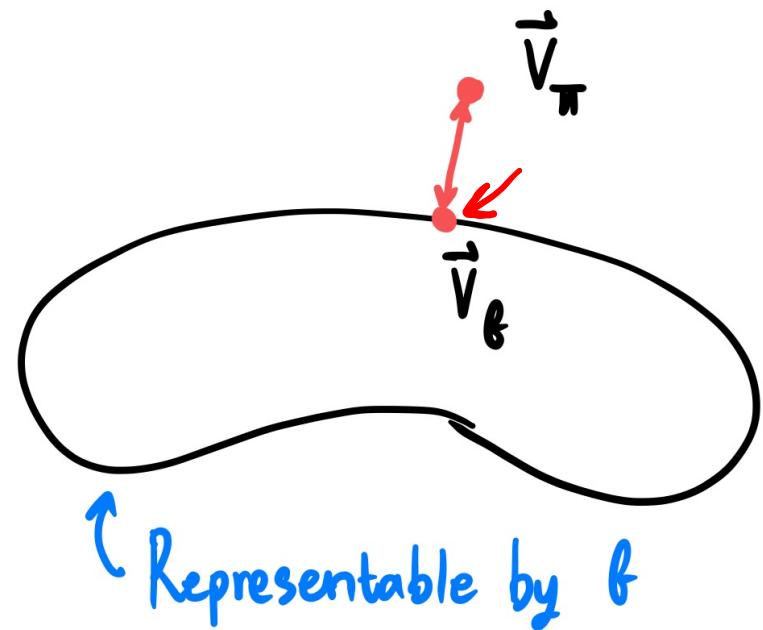
$$\textcircled{2} \quad v_\pi(s) \geq v^*(s) - \frac{2\epsilon}{1-\gamma}$$

$\gamma \uparrow \quad \gamma \downarrow$

Views of approximation



$$\operatorname{argmin}_{\theta} \|\vec{V}_\pi - \vec{V}_\theta\|_{P^\pi}^2$$



Value Error + SGD

$$\mathcal{L}(\theta) = \text{VE}(\theta) = \sum_s \underbrace{\text{P}^\pi(s)}_{\text{on-policy}} \left[\frac{1}{2} (\underbrace{G(s) - v_\theta(s)}_{\text{value error}})^2 \right]$$

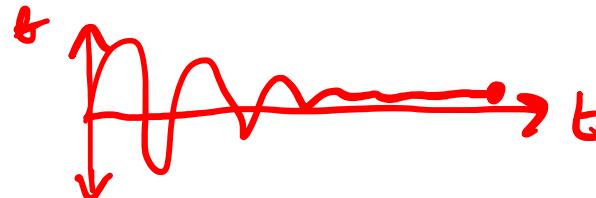
$$\nabla_\theta \mathcal{L}(\theta) = -(G - v_\theta(s)) \nabla_\theta v_\theta(s)$$

$$S_0 \sim \text{P}_{s_0}, S \sim \pi | S_0, G \sim \pi | S$$

$$\theta \leftarrow \theta + \alpha (G - v_\theta(S)) \nabla_\theta v_\theta(S)$$

Convergence and fixed point

Convergence



- A training process converges to a fixed point where there is no further progress

Fixed point

- The solution when the training process converges



TDwith approximation

Prediction and control

$V_{\pi}(s)$? π

optimal π
 $\pi \xrightarrow{\text{impr}} \pi' \xrightarrow{} \pi''$

Prediction

- Get V (estimate of the policy)
- A single policy in concern

$\uparrow V_E \rightarrow V_f(s)$

Control

- Prediction + improvement
- A series of policies

Linear vs non-linear

* Linear

$$\rightarrow v_\theta(s) = \theta^T \underline{\phi(s)}$$

feature vector

$$\nabla_\theta v_\theta(s) = \underline{\underline{\phi(s)}}$$

- Non-linear

$$v_\theta(s) = \underline{f_\theta(\phi(s))}$$

TD Prediction with approximation

$$\hat{V}_f(s)$$

Semi-gradient one-step TD

$$\text{VE}(\theta) = \mathbb{E}_{S \sim P^\pi(s)} \left[\frac{1}{2} (\hat{G}(s) - v_\theta(S))^2 \right]$$

$$\theta \leftarrow \theta + \alpha (v_\pi(s_t) - v_\theta(s_t)) \nabla_\theta v_\theta(s_t)$$

We don't have a return, we use "bootstrapping" target instead

$$② \quad \theta \leftarrow \theta + \alpha (r_{t+1} + \gamma v_\theta(s_{t+1}) - v_\theta(s_t)) \nabla_\theta v_\theta(s_t)$$

$$V_\pi(s) \approx r + \gamma \underbrace{V_\theta(s')}$$

Why we call “semi-gradient”?

$$\theta \leftarrow \theta + \alpha (R + \gamma v_\theta(S') - v_\theta(S)) \boxed{\nabla_\theta v_\theta(S)}$$

- The gradient is “incomplete”
- We assume that $v(s')$ is “independent” from theta
- This is a false assumption
- **Semi-gradient doesn’t share a usual SGD convergence guarantee**

Semi-gradient TD fixed point

$$\theta \leftarrow \theta + \alpha (R + \gamma v_\theta(S') - v_\theta(S)) \nabla_\theta v_\theta(S)$$

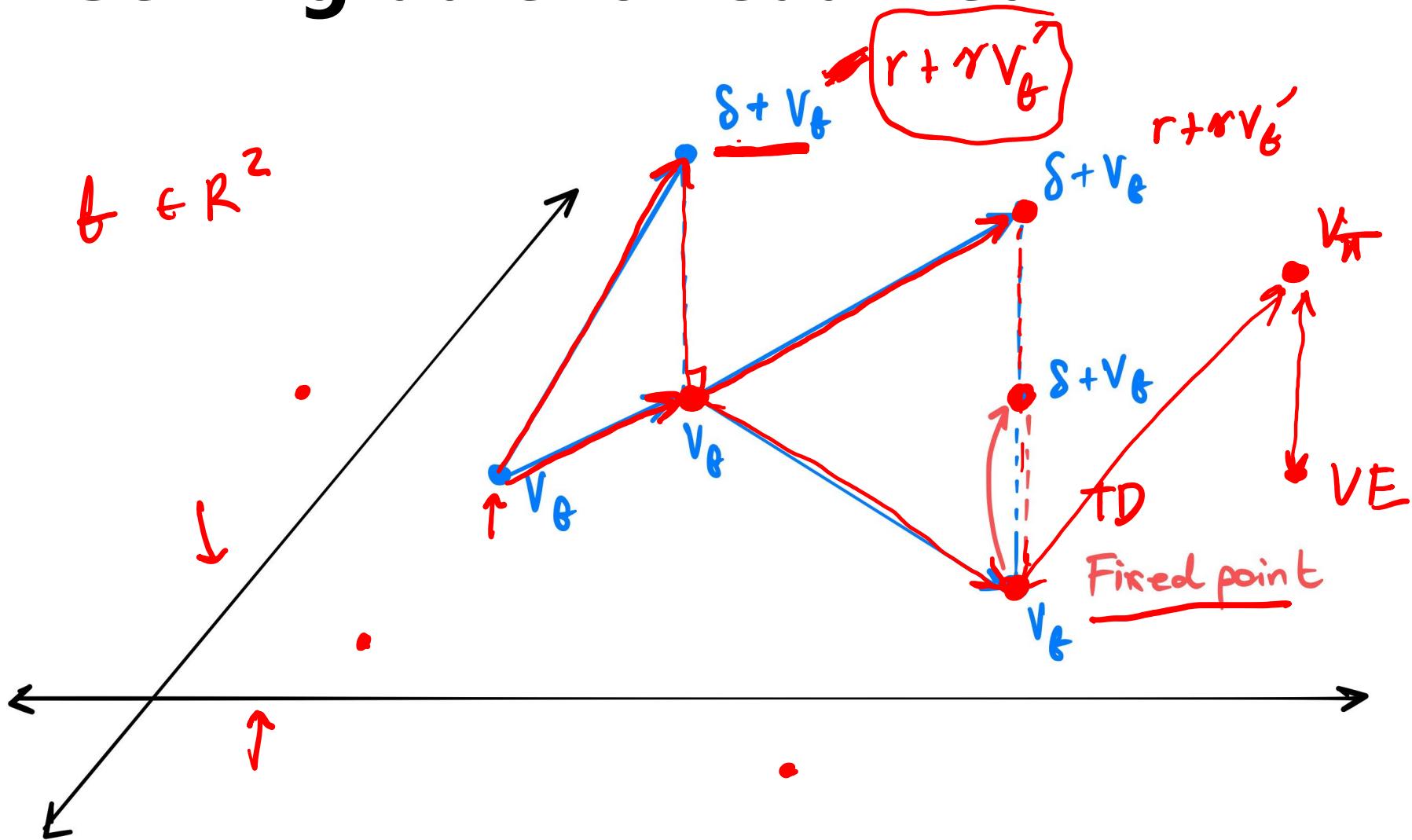
- No further update

$$\theta = \theta + \alpha \mathbb{E} [\delta \nabla_\theta v_\theta(S)]$$

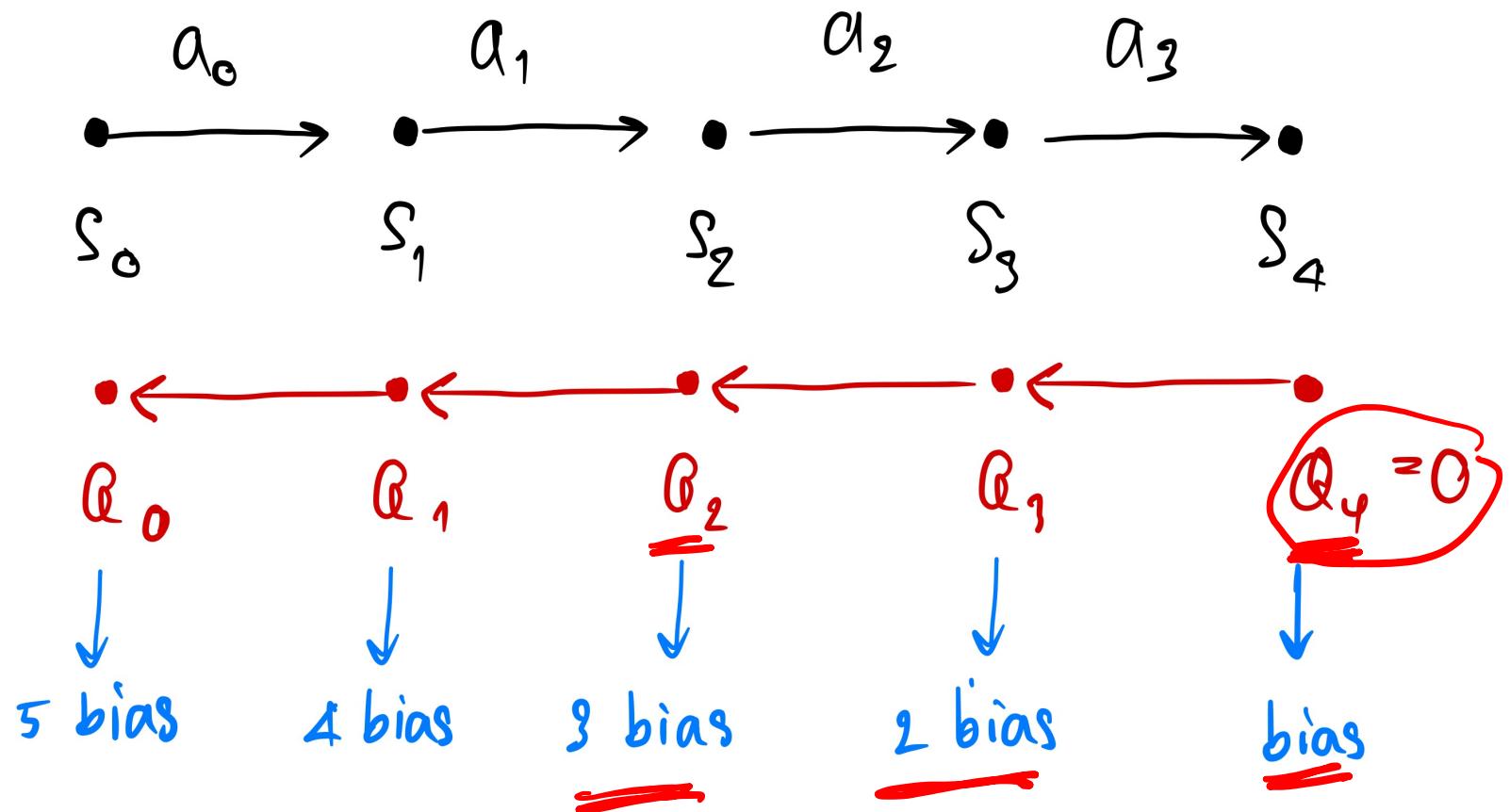
$$0 = \mathbb{E} [\delta \nabla_\theta v_\theta(S)]$$

- Not much could be said...

Semi-gradient visualized



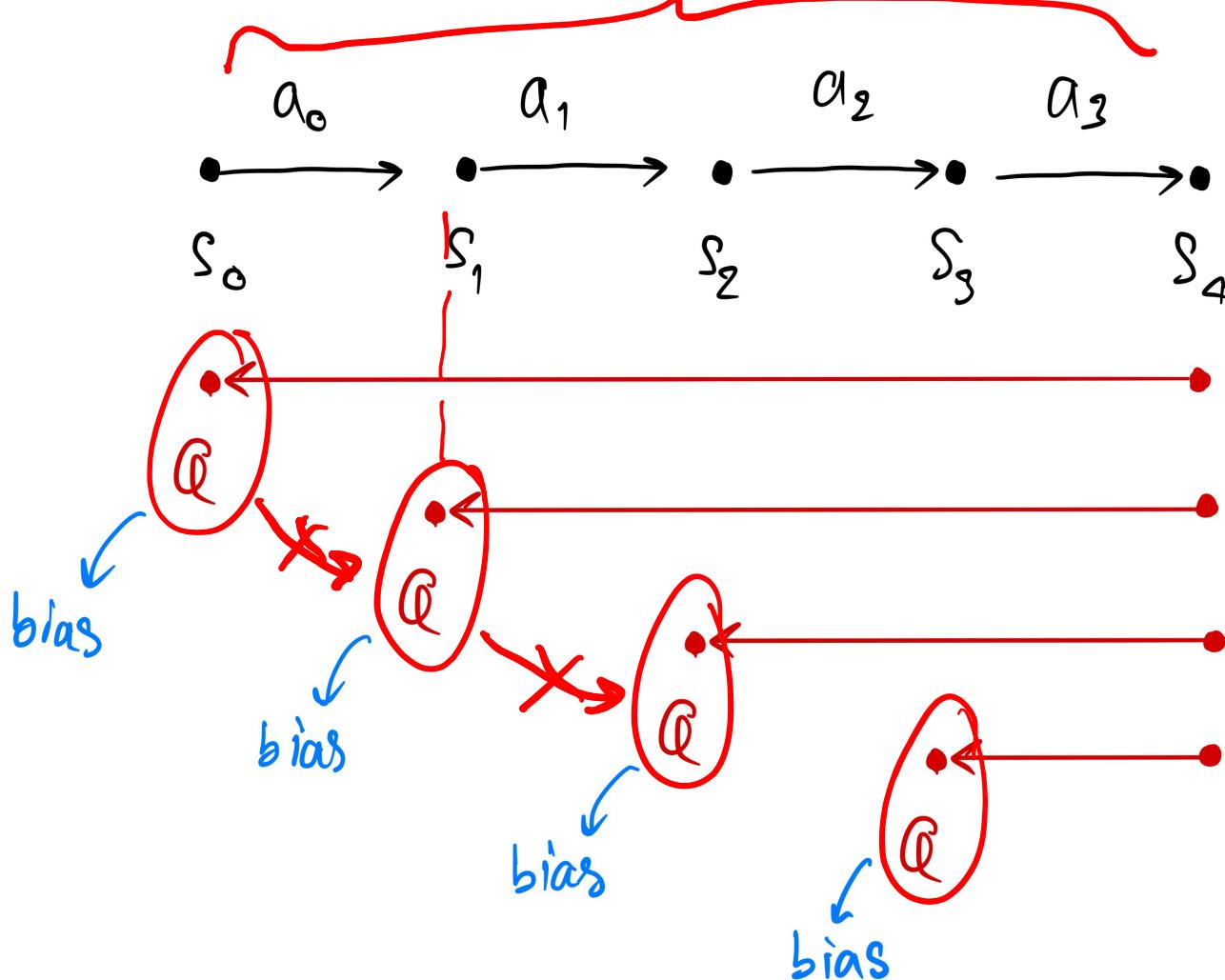
Solution of semi-gradient TD



Solution of VE

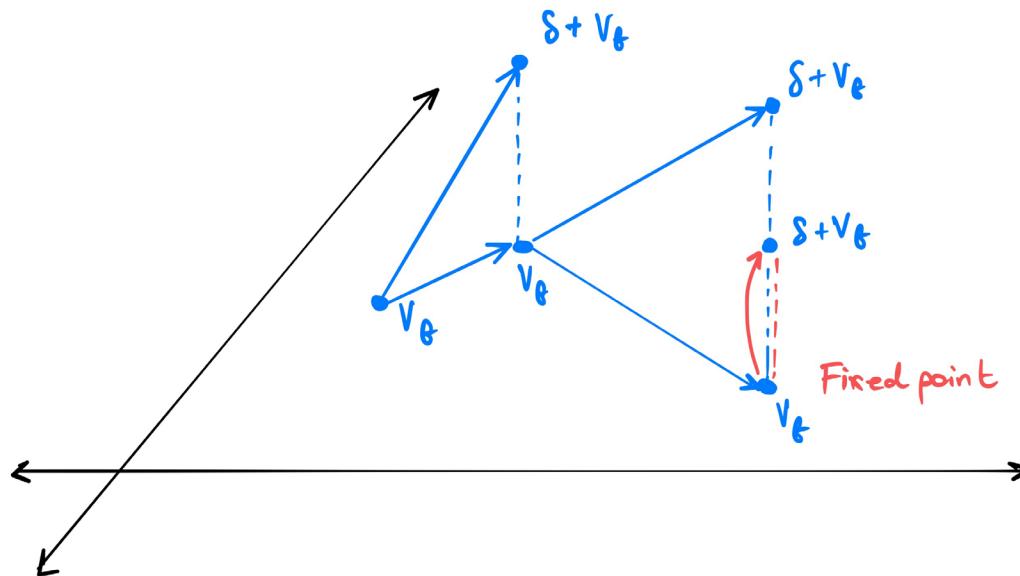
Best

g



Solution of semi-gradient TD

- Poorer solution than VE
- Propagation of errors
 - Each step incurs some error, many steps large error
 - Due to projection onto representable space



N-step semi-gradient

$$\textcircled{1} \quad \theta \leftarrow \theta + \alpha \underbrace{(r_{t+1} + \gamma v_\theta(s_{t+1}) - v_\theta(s_t))}_{\text{Tree backup}} \nabla_\theta v_\theta(s_t)$$

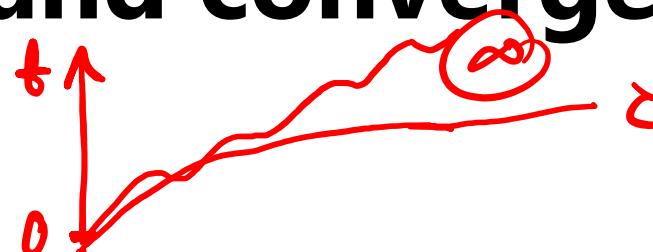
$$\textcircled{2} \quad \theta \leftarrow \theta + \alpha (g_{t:t+n} - v_\theta(s_t)) \nabla_\theta v_\theta(s_t)$$

Tree backup , ↗

- Replace the target with an n-step return
- You could use any kind of target here
- All returns are semi-gradients (except full-return)

Stability and convergence

Stability



* Weights don't explode to **infinity**

- Proof is easier

Convergence

- Weights converge to a fixed point where objective function is minimized
- * Convergence != good fixed point

$$VE > \text{gen gradient TD}$$

Stability of semi-gradient linear TD(0)



- We describe the update in terms of “matrix” multiplication

$$\theta_{t+1} = M\theta_t + c$$

- We show that the matrix is a “**contraction**” mapping
 - Output is smaller than the input vector

$$\|M\theta\| < \|\theta\|$$

Stability of semi-gradient linear TD(0)

$$\theta \leftarrow \theta + \alpha (r_{t+1} + \gamma v_\theta(s_{t+1}) - v_\theta(s_t)) \nabla_\theta v_\theta(s_t)$$
$$\theta_{t+1} = \theta_t + \alpha (r_{t+1} + \gamma \theta_t^T x_{t+1} - \theta_t^T x_t) x_t$$

Goal:

$$\theta_{t+1} = M\theta_t + c$$

$$\|M\theta\| < \|\theta\|$$

Stability of semi-gradient linear TD(0)

$$\theta_{t+1} = M\theta_t + c$$

① $\theta_{t+1} = \theta_t + \alpha (r_{t+1} + \gamma \theta_t^T x_{t+1} - \theta_t^T x_t) x_t$

$\alpha (r_{t+1} x_t + \gamma x_t \theta_t^T x_{t+1} - x_t \theta_t^T x_t)$

$\alpha (x_t + \theta_t^T (\gamma x_{t+1} - x_t))$

$\lambda (r_{t+1} x_t - (x_t \theta_t^T (x_t - \gamma x_{t+1})))$

$$\theta_{t+1} = b_t + \left[d(r_{t+1} x_t) - x_t (x_t - \gamma x_{t+1})^T b_t \right]$$

Stability of semi-gradient linear TD(0)

$$\theta_{t+1} = M\theta_t + c$$

$$\begin{aligned}\theta_{t+1} &= \theta_t + \alpha (r_{t+1}x_t - x_t(x_t - \gamma x_{t+1})^T \theta_t) \\&= \theta_t + \underbrace{\alpha r_{t+1}x_t}_b - \underbrace{\alpha x_t(x_t - \gamma x_{t+1})^T}_{A^+} \underbrace{\theta_t}_b \\&= (\underbrace{I - \alpha A^+}_M) \theta_t + \underbrace{\alpha b}_{\text{bias}}\end{aligned}$$

Stability of semi-gradient linear TD(0) — on-policy

$$\mathbb{E}\theta_{t+1} = (I - \alpha A)\theta_t + \alpha b \quad b = \mathbb{E}_{\pi}^{\pi}[r_{t+1}x_t] \\ A = \mathbb{E}_{\pi}^{\pi}[x_t(x_t - \gamma x_{t+1})^T]$$

- $I - \alpha A$ is a contraction
- By showing that $I - \alpha A$ has eigenvalues between 0 and 1
- By showing that A is positive-definite matrix
- **Under on-policy assumption**

Semi-gradient linear TD fixed-point

$$\mathbb{E}\theta_{t+1} = \theta_t + \alpha b - \alpha A\theta_t$$

$$\mathbb{E} = x_t, x_{t+1} \sim \mathbb{P}^\pi, R \sim \pi$$

$$b = \mathbb{E} [r_{t+1}x_t]$$

Signs of on-policy

$$A = \mathbb{E} [x_t(x_t - \gamma x_{t+1})^T]$$

No further update

$$\theta \leftarrow \theta + \alpha b - \alpha A\theta$$

$$0 = \underline{\alpha b - \alpha A\theta}$$

$$0 = b - A\theta$$

$$b = A\theta$$

$$A^{-1}b = \theta$$

LSTD

Semi-gradient TD property

$$\mathbb{E}\theta_{t+1} = \theta_t + \alpha\mathbb{E} [\delta_t \nabla_\theta v_\theta(s_t)]$$

- It converges to **TD fixed point** in linear case with on-policy

$$b = \mathbb{E} [r_{t+1}x_t]$$

$$\theta = A^{-1}b$$

$$A = \mathbb{E} [x_t(x_t - \gamma x_{t+1})^T]$$

- * **TD fixed point is considered to be a “good” fixed point**
- * It **doesn’t** converge in non-linear case even with on-policy

Summary

- VE update

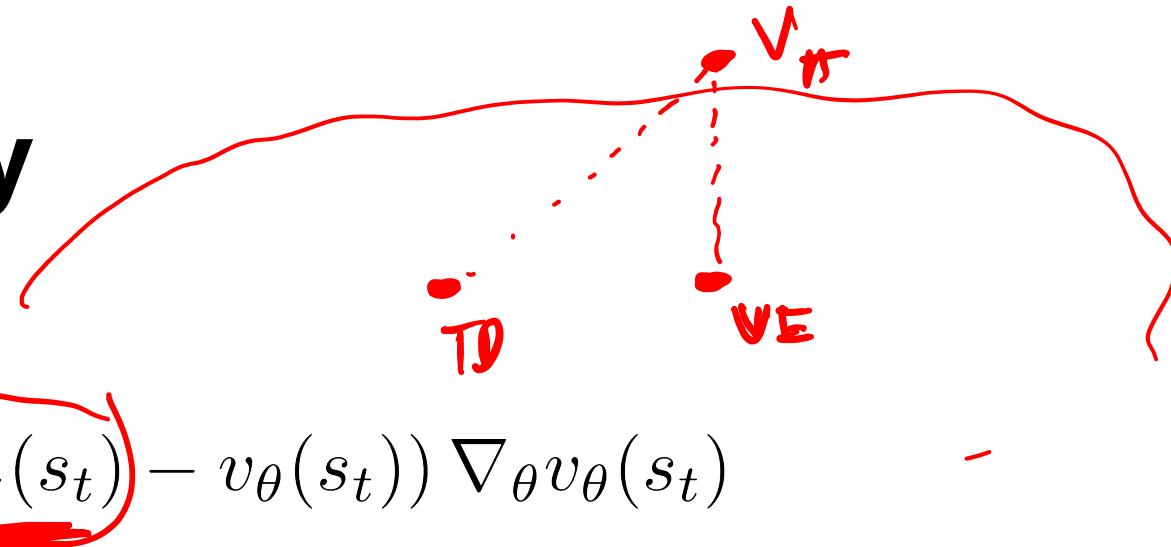
$$\theta \leftarrow \theta + \alpha (v_\pi(s_t) - v_\theta(s_t)) \nabla_\theta v_\theta(s_t)$$

- Semi-gradient TD update

$$\theta \leftarrow \theta + \alpha (r_{t+1} + \gamma v_\theta(s_{t+1}) - v_\theta(s_t)) \nabla_\theta v_\theta(s_t)$$

semi-gradient

- Converges to TD fixed point with on-policy
(only linear case)



TD control with approximation

SARSA with approximation

$$① \theta \leftarrow \theta + \alpha (r_{t+1} + \gamma v_\theta(s_{t+1}) - v_\theta(s_t)) \nabla_\theta v_\theta(s_t)$$

10,000 ~~v_θ(s_{t+1})~~ *q_θ(s_{t+1})*

for until *q_θ* is stable do

1) take action according to $q_\theta(s, a)$ *epsilon greedy*

2) collect $(\underline{s}, \underline{a}, \underline{r}, \underline{s'}, \underline{a'})$

3) $\delta = \underline{r} + \gamma \underline{q_\theta(s', a')} - \underline{q_\theta(s, a)}$

$\theta \leftarrow \theta + \alpha \delta \nabla_\theta q_\theta(s, a)$

end for

- We need to approximate $q_\theta(s, a)$

Implementation considerations

Implementation considerations

- ① Architecture design
- ② Co-adaptation nature of on-policy learning
- ③ Sample correlation
- ④ Forgetting problems

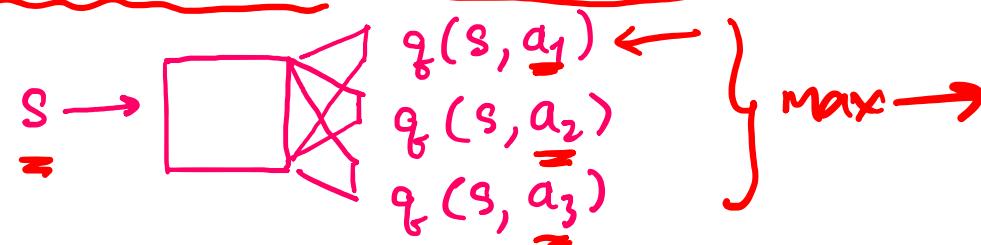
Architecture design

① Layer type

- Conv layers for image inputs

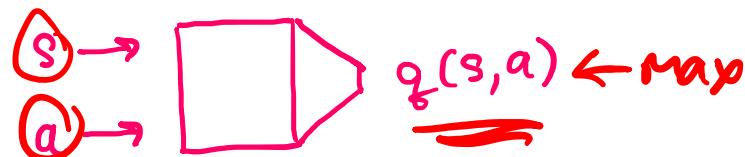
② Prediction head for Q function

* Discrete actions = multi-head

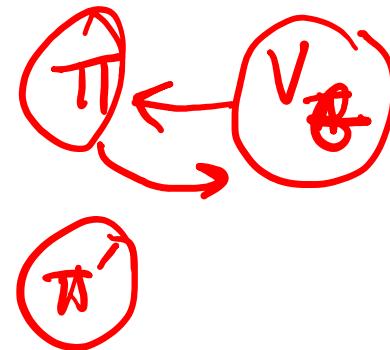


$$q(s, a) \rightarrow$$

- Continuous action = single-head



Co-adaptation nature



Supervised learning formulation

minimize $\hat{\theta}$

$$E_{(x,y) \sim D} [\lambda(x, y; \theta)] \quad D \sim d$$

$V_\theta \rightarrow \theta \rightarrow \theta'$

$$\pi \rightarrow D_\theta$$

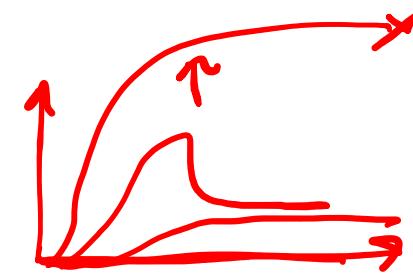
Reinforcement learning formulation

minimize $\hat{\theta}$

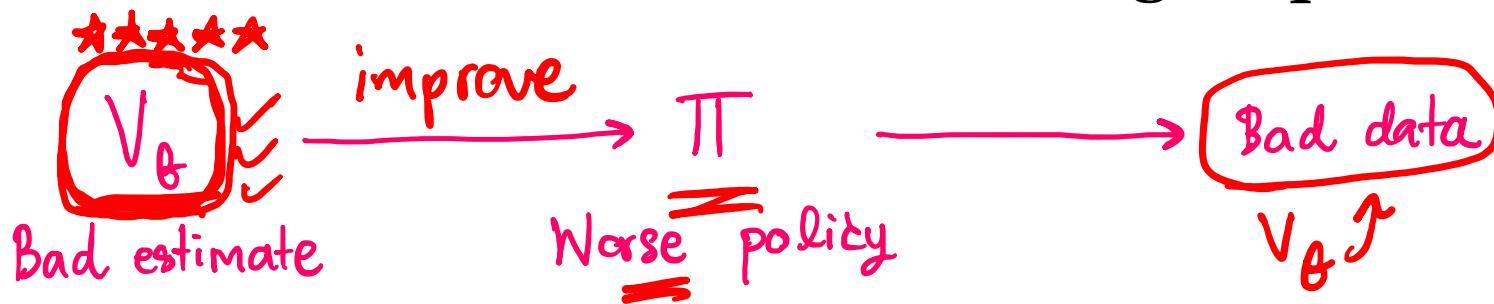
$$E_{(s,a,r,s') \sim D_\theta} [\lambda(\dots; \theta)]$$

$\theta \rightarrow \pi \rightarrow \hat{\theta} \rightarrow \pi \rightarrow \theta \rightarrow \pi$

Co-adaptation nature



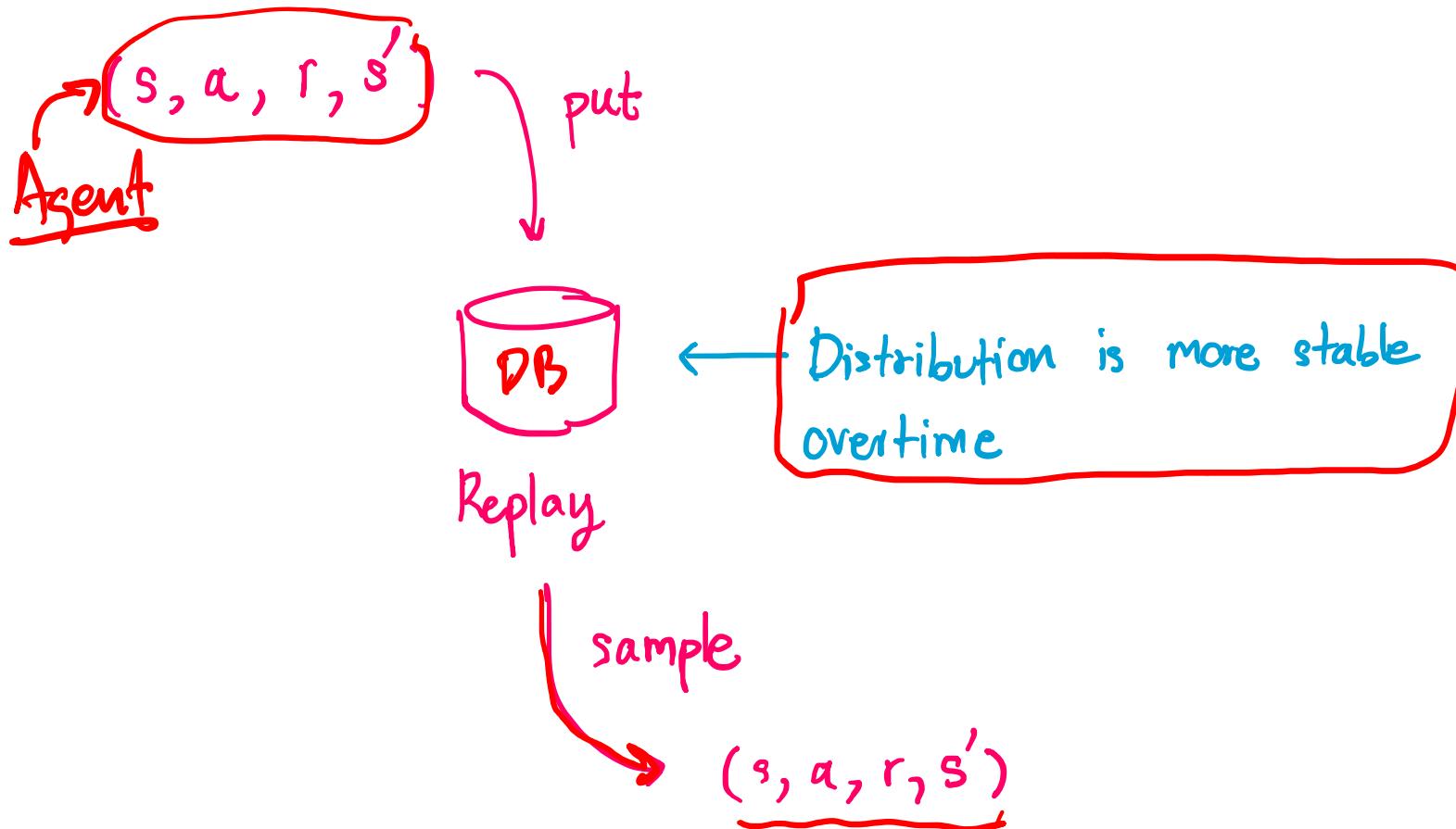
- The data distribution is constantly changing
 - Because the on-policy is constantly changing
- This could lead to unstable learning loop



- * Off-policy with more stable data distribution helps

The diagram shows four boxes labeled D_B , $D_{B'}$, $D_{B''}$, and $D_{B'''}$ stacked vertically, with a wavy line underneath them, representing a more stable data distribution compared to the on-policy case.

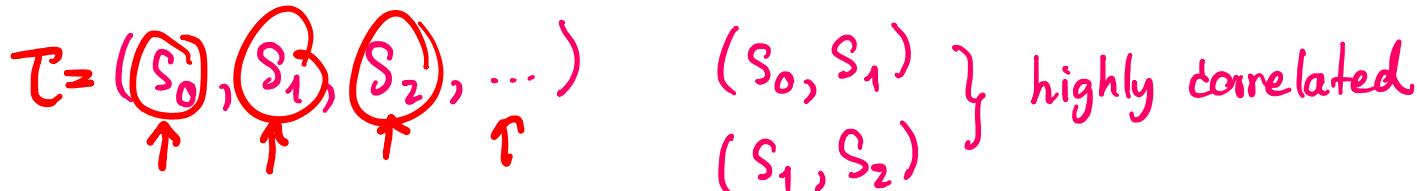
Off policy with replay



Sample correlation

- On-policy sample is highly correlated

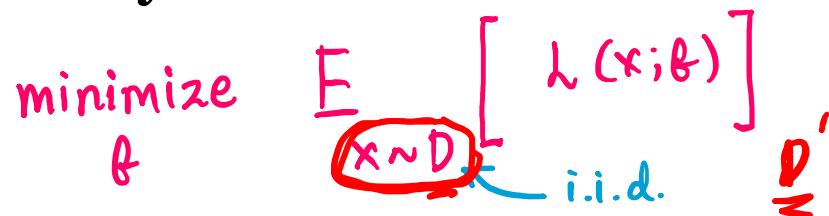
$T = (S_0, S_1, S_2, \dots)$ $\begin{matrix} (S_0, S_1) \\ (S_1, S_2) \end{matrix} \quad \left. \begin{matrix} \text{highly correlated} \end{matrix} \right\}$



- SGD with independent assumption doesn't work very well

$$\underset{\beta}{\text{minimize}} \quad E_{x \sim D} \left[L(x; \beta) \right]$$

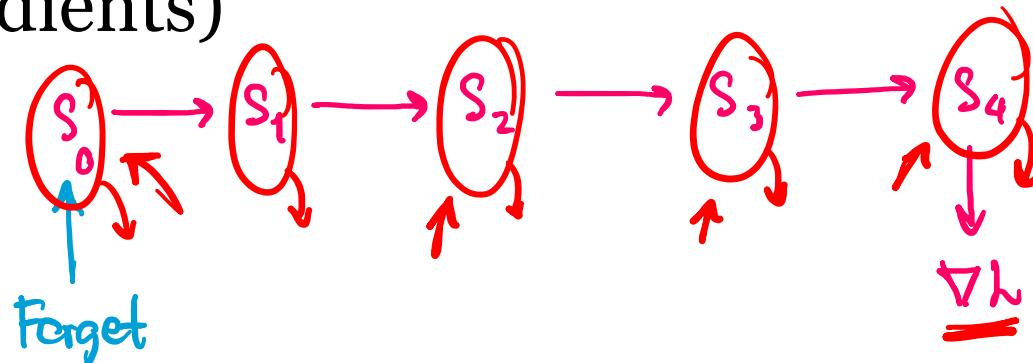
i.i.d.



- It might converge to sub-optimal minima
- Very low learning rate is needed otherwise

Forgetting

- If the gradient is not representative (correlated gradients)



- To reduce we might need very small learning rate