AMI TEA TIME TALKS 2020

- This talk will be recorded. If you want to remain anonymous, please join the meeting from an incognito window and keep your video off.
- Clarifications can be asked for right away. There is a Q-A session at the end of the talk for questions of more technical or philosophical kind.

LEARNING AND PLANNING IN AVERAGE-REWARD MDPS

Abhishek Naik abhishek.naik@ualberta.ca

w/ Yi Wan and Rich Sutton



OUTLINE

- Contributions
- Background
 - Problem setting
 - Related work
- Algorithms and Experiments
 - Control
 - Prediction
 - Centering
- Conclusions and Future Work

1. The first general proven-convergent off-policy modelfree *control* algorithm without reference states

- 1. The first general proven-convergent off-policy modelfree *control* algorithm without reference states
- 2. The first proven-convergent off-policy model-free *prediction* algorithm

- 1. The first general proven-convergent off-policy modelfree *control* algorithm without reference states
- 2. The first proven-convergent off-policy model-free *prediction* algorithm
- 3. A general technique to estimate the actual *centered* value function rather than the value function plus an offset

PROBLEM SETTING



Continuing problems

- Continuing problems
- Tabular representation

- Continuing problems
- Tabular representation
- Unichain MDPs

- Continuing problems
- Tabular representation
- Unichain MDPs



- Continuing problems
- Tabular representation
- Unichain MDPs

$$r(\pi) = \lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^{n} \mathbb{E}[R_t | S_0, A_{0:t-1} \sim \pi]$$

- Continuing problems
- Tabular representation
- Unichain MDPs

$$r(\pi) = \lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^{n} \mathbb{E}[R_t | S_0, A_{0:t-1} \sim \pi]$$

Differential value function

Reward rate

$$v_{\pi}(s) = \lim_{n \to \infty} \frac{1}{n} \sum_{k=1}^{n} \sum_{t=1}^{k} \mathbb{E}[R_t - r(\pi) | S_0 = s, A_{0:t-1} \sim \pi] \quad \forall s$$

- Continuing problems
- Tabular representation
- Unichain MDPs

$$r(\pi) = \lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^{n} \mathbb{E}[R_t | S_0, A_{0:t-1} \sim \pi]$$

Differential value function

$$v_{\pi}(s) = \lim_{n \to \infty} \frac{1}{n} \sum_{k=1}^{n} \sum_{t=1}^{k} \mathbb{E}[R_t - r(\pi) | S_0 = s, A_{0:t-1} \sim \pi] \quad \forall s$$

$$v(s) = \sum_{a} \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) \left[r - \overline{r} + v(s') \right] \quad \forall s$$



- Continuing problems
- Tabular representation
- Unichain MDPs

$$r(\pi) = \lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^{n} \mathbb{E}[R_t | S_0, A_{0:t-1} \sim \pi]$$

Differential value function

$$v_{\pi}(s) = \lim_{n \to \infty} \frac{1}{n} \sum_{k=1}^{n} \sum_{t=1}^{k} \mathbb{E}[R_t - r(\pi) | S_0 = s, A_{0:t-1} \sim \pi] \quad \forall s$$

$$v(s) = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a) \left[r - \bar{r} + v(s')\right] \quad \forall s$$
$$q(s,a) = \sum_{s',r} p(s',r \mid s,a) \left[r - \bar{r} + \max_{a'} q(s',a')\right] \quad \forall s,a$$

- Continuing problems
- Tabular representation
- Unichain MDPs

Reward rate

$$r(\pi) = \lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^{n} \mathbb{E}[R_t | S_0, A_{0:t-1} \sim \pi]$$

Differential value function

$$v_{\pi}(s) = \lim_{n \to \infty} \frac{1}{n} \sum_{k=1}^{n} \sum_{t=1}^{k} \mathbb{E}[R_t - r(\pi) | S_0 = s, A_{0:t-1} \sim \pi] \quad \forall s$$

$$v(s) = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a) [r - \overline{r} + v(s')] \quad \forall s$$
$$q(s,a) = \sum_{s',r} p(s',r \mid s,a) [r - \overline{r} + \max_{a'} q(s',a')] \quad \forall s,a$$

Bellman equations

- Continuing problems
- Tabular representation
- Unichain MDPs

Unique solution for \overline{r} , multiple solutions for v, q

$$r(\pi) = \lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^{n} \mathbb{E}[R_t | S_0, A_{0:t-1} \sim \pi]$$

Differential value function

Reward rate

$$v_{\pi}(s) = \lim_{n \to \infty} \frac{1}{n} \sum_{k=1}^{n} \sum_{t=1}^{k} \mathbb{E}[R_t - r(\pi) | S_0 = s, A_{0:t-1} \sim \pi] \quad \forall s$$

Bellman equations

$$v(s) = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a) \left[r - \overline{r} + v(s')\right] \quad \forall s$$
$$q(s,a) = \sum_{s',r} p(s',r \mid s,a) \left[r - \overline{r} + \max_{a'} q(s',a')\right] \quad \forall s,a$$

Average-reward <i>learning</i> algorithms	Prediction	Control
On-policy		
Off-policy		

Average-reward <i>learning</i> algorithms	Prediction	Control
On-policy	Average Cost TD (1999)	Actor-critic (2000, 2009)
Off-policy		R-learning (1993) Singh (1994) RVI Q-learning (2001) Gosavi (2004)

Average-reward <i>learning</i> algorithms	Prediction	Control
On-policy	Average Cost TD (1999)	Actor-critic (2000, 2009)
Off-policy	Differential TD-learning	R-learning (1993) Singh (1994) RVI Q-learning (2001) Gosavi (2004) Differential Q-learning

BACKGROUND

Average-reward <i>learning</i> algorithms	Prediction	Control
On-policy	Average Cost TD (1999)	Actor-critic (2000, 2009)
Off-policy	Differential TD-learning	R-learning (1993) Singh (1994) RVI Q-learning (2001) Gosavi (2004) Differential Q-learning

Many algorithms that minimize regret (UCRL2, POLITEX, Opt-QL, EE-QL)

BACKGROUND

Average-reward <i>learning</i> algorithms	Prediction	Control
On-policy	Average Cost TD (1999)	Actor-critic (2000, 2009)
Off-policy	Differential TD-learning	R-learning (1993) Singh (1994) RVI Q-learning (2001) Gosavi (2004) Differential Q-learning

Many algorithms that minimize regret, but: typically not off-policy as a behavioural policy needs to be specified, (UCRL2, POLITEX, Opt-QL, EE-QL) value estimates might be unbounded.



$$\mathbf{v}(s) = \sum_{a} \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) \left[r - \overline{r} + v(s') \right] \quad \forall s$$

$$\mathbf{v}(s) = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a) \left[r - \overline{r} + v(s')\right] \quad \forall s$$
$$\overline{r} = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a) \left[r - v(s) + v(s')\right] \quad \forall s$$

$$\mathbf{v}(s) = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a) \left[r - \overline{r} + v(s')\right] \quad \forall s$$
$$\overline{r} = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a) \left[r - v(s) + v(s')\right] \quad \forall s$$

$$\bar{R}_{t+1} = \bar{R}_t + \beta \left(R_{t+1} - V(S_t) + V(S_{t+1}) - \bar{R}_t \right)$$

$$\mathbf{v}(s) = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a) \left[r - \overline{r} + v(s')\right] \quad \forall s$$
$$\overline{r} = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a) \left[r - v(s) + v(s')\right] \quad \forall s$$

$$\bar{R}_{t+1} = \bar{R}_t + \beta \left(R_{t+1} - V(S_t) + V(S_{t+1}) - \bar{R}_t \right)$$

$$r(\pi) = \sum_{s} d_{\pi}(s) \sum_{a} \pi(a \,|\, s) \sum_{s',r} p(s',r \,|\, s,a) r$$

BACKGROUND

$$\mathbf{v}(s) = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a) \left[r - \overline{r} + v(s')\right] \quad \forall s$$
$$\overline{r} = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a) \left[r - v(s) + v(s')\right] \quad \forall s$$

$$\bar{R}_{t+1} = \bar{R}_t + \beta \left(R_{t+1} - V(S_t) + V(S_{t+1}) - \bar{R}_t \right)$$

$$r(\pi) = \sum_{s} d_{\pi}(s) \sum_{a} \pi(a \,|\, s) \sum_{s',r} p(s',r \,|\, s,a) r$$

 $\bar{R}_{t+1} = \bar{R}_t + \beta(R_{t+1} - \bar{R}_t)$

BACKGROUND

$$\mathbf{v}(s) = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a) \left[r - \overline{r} + v(s')\right] \quad \forall s$$
$$\bar{r} = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a) \left[r - v(s) + v(s')\right] \quad \forall s$$

$$\bar{R}_{t+1} = \bar{R}_t + \beta \left(R_{t+1} - V(S_t) + V(S_{t+1}) - \bar{R}_t \right)$$

$$r(\pi) = \sum_{s} d_{\pi}(s) \sum_{a} \pi(a \,|\, s) \sum_{s',r} p(s',r \,|\, s,a) r$$

 $\overline{\bar{R}}_{t+1} = \overline{\bar{R}}_t + \beta(R_{t+1} - \overline{\bar{R}}_t)$

CONTROL

CONTROL

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \max_a Q_t(S_{t+1}, a) - Q_t(S_t, A_t)$$

CONTROL

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \max_a Q_t(S_{t+1}, a) - Q_t(S_t, A_t)$$
$$Q_{t+1}(S_t, A_t) \doteq Q_t(S_t, A_t) + \alpha_t \delta_t$$
$$\bar{R}_{t+1} \doteq \bar{R}_t + \eta \alpha_t \delta_t$$

CONTROL

Differential Q-learning

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \max_a Q_t(S_{t+1}, a) - Q_t(S_t, A_t)$$
$$Q_{t+1}(S_t, A_t) \doteq Q_t(S_t, A_t) + \alpha_t \delta_t$$
$$\bar{R}_{t+1} \doteq \bar{R}_t + \eta \alpha_t \delta_t$$

CONTROL

Differential Q-learning

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \max_a Q_t(S_{t+1}, a) - Q_t(S_t, A_t)$$
$$Q_{t+1}(S_t, A_t) \doteq Q_t(S_t, A_t) + \alpha_t \delta_t$$
$$\bar{R}_{t+1} \doteq \bar{R}_t + \eta \alpha_t \delta_t$$

Theorem 1 (sketch)

If 1) the MDP is unichain,
2) the stepsizes are decreased appropriately,
3) all the state-action pairs are updated infinite number of times,
4) the maximum ratio of the update frequencies is finite,

then the Differential Q-learning algorithm converges a.s.: \bar{R}_t to $r(\pi^*)$, Q_t to a solution of the Bellman optimality equation.

CONTROL

Algorithm 1: Differential Q-learning (one-step off-policy control)

Input: The policy b to be used (e.g., ϵ -greedy) **Algorithm parameters:** step size α, η

- 1 Initialize $Q(s, a) \forall s, a; \overline{R}$ arbitrarily (e.g., to zero)
- 2 Obtain initial S
- 3 while still time to train do
- 4 $A \leftarrow action given by b for S$
- 5 Take action A, observe R, S'

$$\delta = R - \bar{R} + \max_a Q(S', a) - Q(S, A)$$

$$O(S, A) = O(S, A) + \alpha \delta$$

7
$$Q(S, A) = Q(S, A) + \alpha \delta$$

8
$$R = R + \eta \alpha \delta$$

9
$$S =$$

- 10 **end**
- 11 return Q
CONTROL

Algorithm 1: Differential Q-learning (one-step off-policy control)

Input: The policy b to be used (e.g., ϵ -greedy) Algorithm parameters: step size α, η

- 1 Initialize $Q(s, a) \forall s, a; \overline{R}$ arbitrarily (e.g., to zero)
- 2 Obtain initial S
- 3 while still time to train do
- 4 $A \leftarrow action given by b for S$
- 5 Take action A, observe R, S'

$$\delta = R - \bar{R} + \max_a Q(S', a) - Q(S, A)$$

7
$$Q(S,A) = Q(S,A) + \alpha q$$

8
$$R = R + \eta \alpha \delta$$

9 $S = S'$

10 end

11 return Q

RVI Q-learning

$$Q_{t+1}(S_t, A_t) = Q_t(S_t, A_t) + \alpha_t \Big(R_{t+1} - f(Q_t) \\ + \max_a Q_t(S_{t+1}, a) - Q_t(S_t, A_t) \Big)$$

- Two-state MDP
 - state 0 is transient



- Two-state MDP
 - state 0 is transient



- Reference state-action pair: (0, a)
 - i.e., f(Q) = Q(0, a)
- $\alpha = 0.01, \eta = 1$

Two-state MDP



state 0 is transient



- Reference state-action pair: (0, a)
 - i.e., f(Q) = Q(0, a)

• $\alpha = 0.01, \eta = 1$

Two-state MDP

 $\begin{array}{c} 0.9 \\ +1 \\ 0 \\ b \\ -10 \end{array}$

state 0 is transient



- Reference state-action pair: (0, a)
 - i.e., f(Q) = Q(0, a)
- $\alpha = 0.01, \, \eta = 1$

pair: (0, a)

i.e., f(Q) = Q(0, a)

 $\alpha = 0.01, \eta = 1$

Two-state MDP

0.1 0.9 а +2 +1 0 b -10

state 0 is transient



RVI Q-learning diverges if reference state is transient.

environment CONTROL

Access Control Queueing Task



environment CONTROL

Access Control Queueing Task



environment CONTROL

Access Control Queueing Task



EXPERIMENT

- α ∈ {0.0015625, 0.00625, 0.025, 0.1, 0.4}
- ▶ $\eta \in \{0.125, 0.25, 0.5, 1, 2\}$
- $\epsilon = 0.1$
- 80,000 steps
- 30 runs

- Reference states:
 - 0, 2, 4, 6, 8, 10
 free servers
 - priority 8
 - accept

EXPERIMENT

- ▶ $\alpha \in \{0.0015625, 0.00625, 0.025, 0.1, 0.4\}$
- η ∈ {0.125, 0.25, 0.5, 1, 2}
- $\bullet \ \epsilon = 0.1$
- 80,000 steps
- 30 runs

- Reference states:
 - 0, 2, 4, 6, 8, 10
 free servers
 - priority 8
 - accept



EXPERIMENT CONTROL

Sensitivity analysis



EXPERIMENT CONTROL

Sensitivity analysis



- RVI Q-learning's performance depends significantly on the choice of the reference state.
- Differential Q-learning's performance varies only slightly over a wide range of parameter values.

$$\delta_t \doteq R_{t+1} - \bar{R}_t + V_t(S_{t+1}) - V_t(S_t)$$

$$\delta_t \doteq R_{t+1} - \bar{R}_t + V_t(S_{t+1}) - V_t(S_t)$$

$$V_{t+1}(S_t) \doteq V_t(S_t) + \alpha_t \rho_t \delta_t$$
$$\bar{R}_{t+1} \doteq \bar{R}_t + \eta \alpha_t \rho_t \delta_t$$

PREDICTION

Differential TD-learning

$$\delta_t \doteq R_{t+1} - \bar{R}_t + V_t(S_{t+1}) - V_t(S_t)$$

$$V_{t+1}(S_t) \doteq V_t(S_t) + \alpha_t \rho_t \delta_t$$
$$\bar{R}_{t+1} \doteq \bar{R}_t + \eta \alpha_t \rho_t \delta_t$$

PREDICTION

Differential TD-learning

$$\delta_t \doteq R_{t+1} - \bar{R}_t + V_t(S_{t+1}) - V_t(S_t)$$

$$V_{t+1}(S_t) \doteq V_t(S_t) + \alpha_t \rho_t \delta_t$$
$$\bar{R}_{t+1} \doteq \bar{R}_t + \eta \alpha_t \rho_t \delta_t$$

Theorem 2 (sketch)

- If 1) the MDP is unichain,
 - 2) the stepsizes are decreased appropriately,
 - 3) all the states are updated infinite number of times,
 - 4) the maximum ratio of the update frequencies is finite,
 - 5) b covers all the actions that π may choose in all states,

then the Differential TD-learning algorithm converges a.s.: \bar{R}_t to $r(\pi)$, V_t to a solution of the Bellman equation.

PREDICTION

Algorithm 3: Differential TD-learning (one-step off-policy prediction)

- **Input:** The policy π to be evaluated, and b to be used Algorithm parameters: step sizes α , η
- Initialize $V(s) \forall s, \overline{R}$ arbitrarily (e.g., to zero)
- 2 while still time to train do
- $\begin{array}{c|c}3 & A \leftarrow \text{action given by } b \text{ for } S \\ 4 & \text{Take action } A, \text{ observe } R, S' \end{array}$

$$\delta = R - \bar{R} + V(S') - V(S)$$

$$\rho = \frac{\pi(A|S)}{b(A|S)}$$

7 $V(S) = V(S) + \alpha \rho \delta$

$$\begin{array}{c|c} \mathbf{s} & R = R + \eta \alpha \rho \delta \\ \mathbf{s} & S = S' \end{array}$$

- 10 end
- 11 return V

PREDICTION

Algorithm 3: Differential TD-learning (one-step off-policy prediction) **Input:** The policy π to be evaluated, and b to be used Algorithm parameters: step sizes α, η Initialize $V(s) \forall s$, \overline{R} arbitrarily (e.g., to zero) while still time to train do $\mathbf{2}$ $A \leftarrow action given by b for S$ 3 Take action A, observe R, S'4 $\delta = R - \bar{R} + V(S') - V(S)$ 5 $\rho = \frac{\pi(A|S)}{b(A|S)}$ Average Cost TD-learning 6 $V(S) = V(S) + \alpha \rho \delta$ 7 $\bar{R}_{t+1} \doteq \bar{R}_t + \eta \alpha_t (R_{t+1} - \bar{R}_t)$ $\bar{R} = \bar{R} + \eta \alpha \rho \delta$ 8 S = S'9 end 1011 return V

PREDICTION

Algorithm 3: Differential TD-learning (one-step off-policy prediction) **Input:** The policy π to be evaluated, and b to be used Algorithm parameters: step sizes α, η Initialize $V(s) \forall s$, \overline{R} arbitrarily (e.g., to zero) while still time to train do $\mathbf{2}$ $A \leftarrow action given by b for S$ 3 Take action A, observe R, S'4 $\delta = R - \bar{R} + V(S') - V(S)$ 5 $\rho = \frac{\pi(A|S)}{b(A|S)}$ Average Cost TD-learning 6 $V(S) = V(S) + \alpha \rho \delta$ 7 $\bar{R}_{t+1} \doteq \bar{R}_t + \eta \alpha_t (R_{t+1} - \bar{R}_t)$ $\bar{R} = \bar{R} + \eta \alpha \rho \delta$ 8 S = S'9 (restricted to on-policy) end 1011 return V

environment PREDICTION

Two Loop Task



EXPERIMENT

- $\pi_0 = [0.5, 0.5], b_0 = [0.9, 0.1]$
- ▶ $\eta \in \{0.125, 0.25, 0.5, 1, 2\}$
- $\epsilon = 0.1$
- 10,000 steps
- 30 runs

environment PREDICTION

Two Loop Task



EXPERIMENT

- $\pi_0 = [0.5, 0.5], b_0 = [0.9, 0.1]$
- *α* ∈ {0.025, 0.05, 0.1, 0.2, 0.4}
- ▶ $\eta \in \{0.125, 0.25, 0.5, 1, 2\}$
- $\epsilon = 0.1$
- 10,000 steps
- 30 runs
- Evaluation metric:
 - RMSVE

$$\|v - v_{\pi}\|_{d_{\pi}}$$

environment PREDICTION

Two Loop Task



EXPERIMENT

- $\pi_0 = [0.5, 0.5], b_0 = [0.9, 0.1]$
- *α* ∈ {0.025, 0.05, 0.1, 0.2, 0.4}
- ▶ $\eta \in \{0.125, 0.25, 0.5, 1, 2\}$
- $\bullet \ \epsilon = 0.1$
- 10,000 steps
- 30 runs
- Evaluation metric:
 - RMSVE

 $\inf_{c} \|v - (v_{\pi} + ce)\|_{d_{\pi}}$

(Tsitsiklis and Van Roy, 1999)

RESULTS PREDICTION



Learning curves

results **PREDICTION**

Sensitivity analysis



RESULTS PREDICTION

Sensitivity analysis



Differential TD-learning converges faster for a wide range of parameters.

results **PREDICTION**

Sensitivity analysis



Differential TD-learning converges faster for a wide range of parameters.

results **PREDICTION**

Sensitivity analysis



- Differential TD-learning converges faster for a wide range of parameters.
- Differential TD-learning works in the off-policy setting as well.



CENTERING

Recall:
$$v(s) = \sum_{a} \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [R_{t+1} - \bar{r} + v(s')] \quad \forall s$$

CENTERING

Recall:
$$v(s) = \sum_{a} \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [R_{t+1} - \bar{r} + v(s')] \quad \forall s$$

Solutions: $v = v_{\pi} + ce$

CENTERING

Recall:
$$v(s) = \sum_{a} \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [R_{t+1} - \overline{r} + v(s')] \quad \forall s$$

Solutions:
$$v = v_{\pi} + ce$$

Lemma

$$d_{\pi}^T v_{\pi} = 0,$$

i.e., the average of the differential value function is zero.

CENTERING

Recall:
$$v(s) = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s', r \mid s, a) [R_{t+1} - \bar{r} + v(s')] \quad \forall s$$

Solutions: $v = v_{\pi} + ce$
Lemma
 $d_{\pi}^{T} v_{\pi} = 0,$
i.e., the average of the differential value function is zero.

 \implies there is only one *centered* differential value function
MOTIVATION

CENTERING

Recall:
$$v(s) = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s', r \mid s, a) [R_{t+1} - \bar{r} + v(s')] \quad \forall s$$

Solutions: $v = v_{\pi} + ce$
Lemma
 $d_{\pi}^{T} v_{\pi} = 0,$
i.e., the average of the differential value function is zero.

 \implies there is only one *centered* differential value function

 $v = v_{\pi} + ce$

MOTIVATION

CENTERING

Recall:
$$v(s) = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s', r \mid s, a) [R_{t+1} - \bar{r} + v(s')] \quad \forall s$$

Solutions: $v = v_{\pi} + ce$
Lemma
 $d_{\pi}^{T} v_{\pi} = 0,$
i.e., the average of the differential value function is zero.

 \implies there is only one *centered* differential value function

$$v = v_{\pi} + ce$$

$$\Rightarrow c = d_{\pi}^{T} v$$

MOTIVATION

CENTERING

Recall:
$$v(s) = \sum_{a} \pi(a \mid s) \sum_{s',r} p(s', r \mid s, a) [R_{t+1} - \bar{r} + v(s')] \quad \forall s$$

Solutions: $v = v_{\pi} + ce$
Lemma
 $d_{\pi}^{T} v_{\pi} = 0,$
i.e., the average of the differential value function is zero.

 \implies there is only one *centered* differential value function

$$v = v_{\pi} + ce$$
$$\implies c = d_{\pi}^{T}v$$
$$r(\pi) = d_{\pi}^{T}r_{\pi}$$

ALGORITHM

CENTERING

$$\begin{split} \delta_t &\doteq R_{t+1} - \bar{R}_t + V_t(S_{t+1}) - V_t(S_t) \\ V_{t+1}(S_t) &\doteq V_t(S_t) + \alpha_t \rho_t \delta_t \\ \bar{R}_{t+1} &\doteq \bar{R}_t + \eta \alpha_t \rho_t \delta_t \end{split} \begin{array}{l} \text{System 1} \\ \end{split}$$

ALGORITHM

CENTERING

$$\delta_{t} \doteq R_{t+1} - \bar{R}_{t} + V_{t}(S_{t+1}) - V_{t}(S_{t})$$

$$V_{t+1}(S_{t}) \doteq V_{t}(S_{t}) + \alpha_{t}\rho_{t}\delta_{t}$$

$$\bar{R}_{t+1} \doteq \bar{R}_{t} + \eta\alpha_{t}\rho_{t}\delta_{t}$$
System 1

$$\Delta_{t} \doteq V_{t}(S_{t}) - \bar{V}_{t} + F_{t}(S_{t+1}) - F_{t}(S_{t})$$

$$F_{t+1}(S_{t}) \doteq F_{t}(S_{t}) + \beta_{t}\rho_{t}\Delta_{t}$$

$$\bar{V}_{t+1} \doteq \bar{V}_{t} + \kappa\beta_{t}\rho_{t}\Delta_{t}$$
System 2

ALGORITHM

CENTERING

$$\delta_{t} \doteq R_{t+1} - \bar{R}_{t} + V_{t}(S_{t+1}) - V_{t}(S_{t})$$

$$V_{t+1}(S_{t}) \doteq V_{t}(S_{t}) + \alpha_{t}\rho_{t}\delta_{t}$$

$$\bar{R}_{t+1} \doteq \bar{R}_{t} + \eta\alpha_{t}\rho_{t}\delta_{t}$$
System 1

$$\begin{split} \Delta_t &\doteq V_t(S_t) - \bar{V}_t + F_t(S_{t+1}) - F_t(S_t) \\ F_{t+1}(S_t) &\doteq F_t(S_t) + \beta_t \rho_t \Delta_t \\ \bar{V}_{t+1} &\doteq \bar{V}_t + \kappa \beta_t \rho_t \Delta_t \end{split} \qquad \begin{array}{l} \textbf{System 2} \\ \textbf{System 2} \\ \textbf{System 3} \\ \textbf{System 4} \\ \textbf{System 3} \\ \textbf{System 4} \\ \textbf{System 5} \\ \textbf{System 6} \\ \textbf{System 6}$$

Theorem 3 (sketch)

If the previous assumptions hold, then the Centered Differential TD-learning algorithm converges a.s.: \bar{R}_t to $r(\pi)$, $V_t - \bar{V}_t e$ to the centered differential value function

ENVIRONMENT

CENTERING

Two Loop Task

EXPERIMENT

- ▶ $\beta \in \{0.025, 0.05, 0.1, 0.2, 0.4\}$
- ▶ $\kappa \in \{0.125, 0.25, 0.5, 1, 2\}$
- $\epsilon = 0.1$
- 10,000 steps
- 30 runs



ENVIRONMENT

CENTERING

Two Loop Task



EXPERIMENT

- ▶ $\beta \in \{0.025, 0.05, 0.1, 0.2, 0.4\}$
- ▶ $\kappa \in \{0.125, 0.25, 0.5, 1, 2\}$
- $\epsilon = 0.1$
- 10,000 steps
- 30 runs

- Evaluation metric:
 - RMSVE

$$\|v-v_{\pi}\|_{d_{\pi}}$$

(the usual one)

results **CENTERING**



Learning curves

results **CENTERING**



Sensitivity analysis

The first general proven-convergent off-policy model-free control algorithm without reference states

- The first general proven-convergent off-policy model-free control algorithm without reference states
- The first proven-convergent off-policy model-free prediction algorithm

- The first general proven-convergent off-policy model-free control algorithm without reference states
- The first proven-convergent off-policy model-free prediction algorithm
- A general technique to estimate the actual centered value function rather than the value function plus an offset

- The first general proven-convergent off-policy model-free control algorithm without reference states
- The first proven-convergent off-policy model-free prediction algorithm
- A general technique to estimate the actual centered value function rather than the value function plus an offset
 - All of our learning algorithms are fully online, and all of our planning algorithms are fully incremental

- The first general proven-convergent off-policy model-free control algorithm without reference states
- The first proven-convergent off-policy model-free prediction algorithm
- A general technique to estimate the actual centered value function rather than the value function plus an offset
 - All of our learning algorithms are fully online, and all of our planning algorithms are fully incremental
 - Empirically, the use of the temporal-difference error generally results in faster learning in the domains tested, and reliance on a reference state generally results in slower learning and risks divergence.

Extension of these tabular algorithms to function approximation

- Extension of these tabular algorithms to function approximation
 - currently working on linear FA, both learning and planning

Extension of these tabular algorithms to function approximation

- currently working on linear FA, both learning and planning
- deadly triad :/

Extension of these tabular algorithms to function approximation

- currently working on linear FA, both learning and planning
- deadly triad :/

Extension to SMDPs so they can be used with temporal abstractions like options

Extension of these tabular algorithms to function approximation

- currently working on linear FA, both learning and planning
- deadly triad :/

Extension to SMDPs so they can be used with temporal abstractions like options

Does centering the features stabilize (off-policy) learning?



Paper: <u>https://arxiv.org/abs/2006.16318</u> Code: <u>https://github.com/abhisheknaik96/average-reward-methods</u>

INTUITION: WHY USING TD ERROR MIGHT BE BETTER

INTUITION: WHY USING TD ERROR MIGHT BE BETTER

Example 10.8 from the book