

# Putting it all together

RL Superman: uses Monte-Carlo to estimate state-action values using  $\epsilon$ -greedy exploration with decaying  $\epsilon$



- Why super?
  - Learns  $Q$ -values from observed returns (doesn't require a model)
  - Estimates become better over time with more experience
  - Can choose the best action as  $\operatorname{argmax}_a Q^\pi(s, a)$
  - Starts out with exploration ( $\epsilon = 1$ ), but slowly becomes greedy ( $\epsilon \approx 0$ )
- But...
  - requires a lot of experience to get good estimates and policy
  - works for small finite MDPs only
  - applicable to episodic problems only
  - wastes time evaluating bad policies
- Solution: Temporal Difference Learning

# Temporal Differencing (SUBA6)

*“If one had to identify one idea as central and novel to reinforcement learning, it would undoubtedly be temporal-difference (TD) learning.”* (SUBA6)

- Monte Carlo
  - update values after an episode is done (based on returns  $R = \sum_t r_t$ )
- Temporal-difference learning
  - update values after each step (based on immediate reward  $r_t$ )

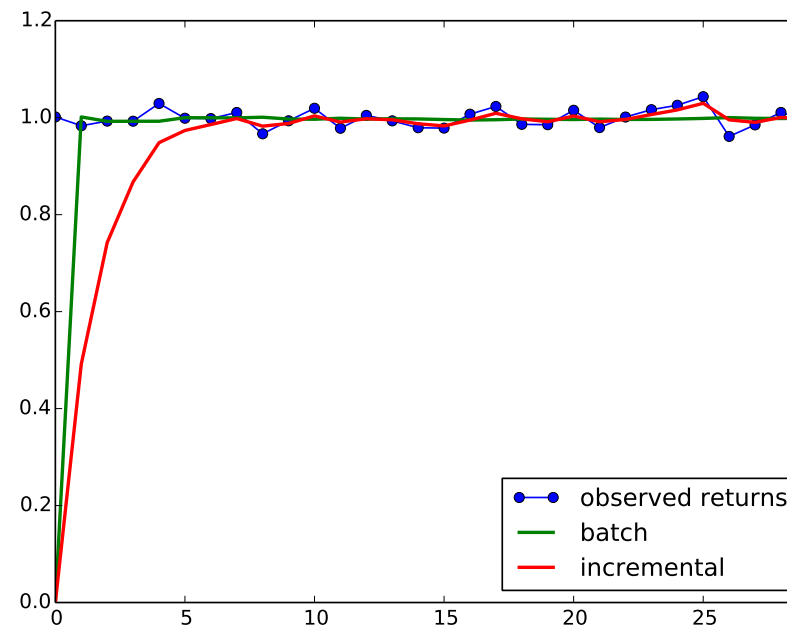
# Temporal Differencing (SUBA6)

*“If one had to identify one idea as central and novel to reinforcement learning, it would undoubtedly be temporal-difference (TD) learning.”* (SUBA6)

- Monte Carlo
  - update values after an episode is done (based on returns  $R = \sum_t r_t$ )
- Temporal-difference learning
  - update values after each step (based on immediate reward  $r_t$ )
- Explained in two steps
  - 1 Show incremental version of Monte Carlo (uses  $R$ )
  - 2 Show TD learning (uses  $r_t$ )

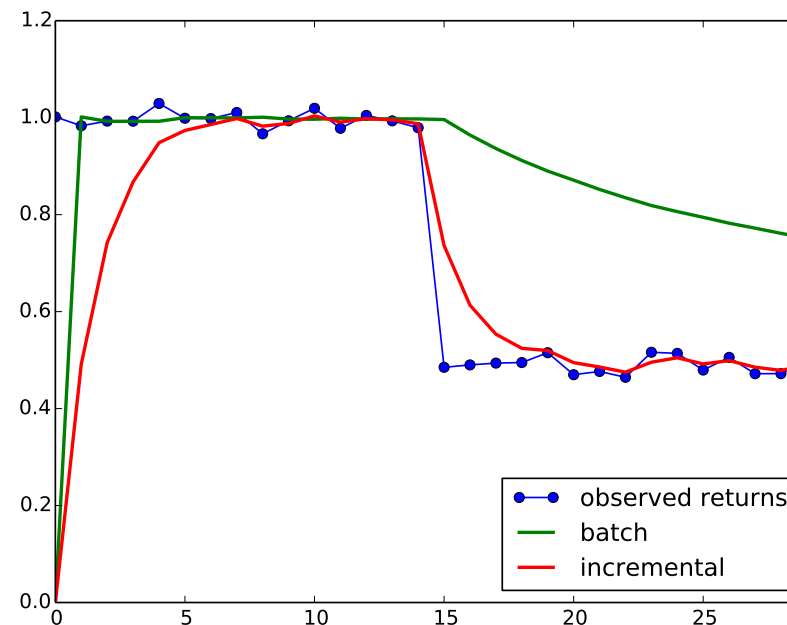
# Monte Carlo: Batch vs. Incremental

- Batch:  $V^\pi(s) = \text{average}(\text{Returns}(s)) = \frac{1}{N} \sum_{e=1}^N R(s)^e$ 
  - $N$  is the number of episodes
- Incremental updates:  $V^\pi(s) = V^\pi(s) + \alpha [R - V^\pi(s)]$ 
  - $0 < \alpha < 1$  is 'learning rate' ( $\alpha = 1$  would be batch,  $\alpha = 0$  would mean no learning)
  - $\alpha$  also known as 'step-size'



# Monte Carlo: Batch vs. Incremental

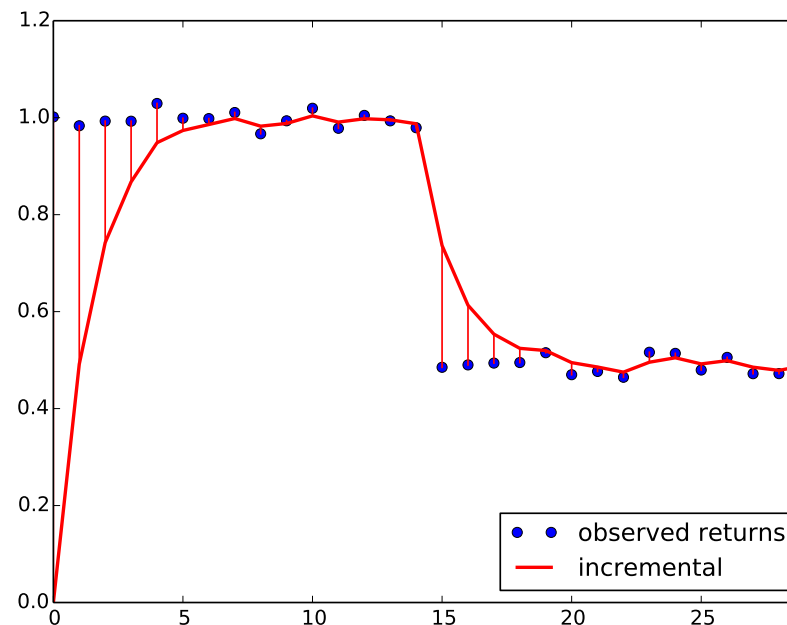
- Batch:  $V^\pi(s) = \text{average}(\text{Returns}(s)) = \frac{1}{N} \sum_{e=1}^N R(s)^e$ 
  - $N$  is the number of episodes
- Incremental updates:  $V^\pi(s) = V^\pi(s) + \alpha [R - V^\pi(s)]$ 
  - $0 < \alpha < 1$  is 'learning rate' ( $\alpha = 1$  would be batch,  $\alpha = 0$  would mean no learning)
  - $\alpha$  also known as 'step-size'



- Incremental better if environment is non-stationary (changes over time)

# Monte Carlo: Batch vs. Incremental

- Batch:  $V^\pi(s) = \text{average}(\text{Returns}(s)) = \frac{1}{N} \sum_{e=1}^N R(s)^e$ 
  - $N$  is the number of episodes
- Incremental updates:  $V^\pi(s) = V^\pi(s) + \alpha [R - V^\pi(s)]$ 
  - $0 < \alpha < 1$  is 'learning rate' ( $\alpha = 1$  would be batch,  $\alpha = 0$  would mean no learning)
  - $\alpha$  also known as 'step-size'



# From incremental MC to TD (SUBA6.1)

Definition

$$V^\pi(s) = E_\pi\{R_t | s_t = s\}$$

Incremental estimation from experience

$$V^\pi(s_t) = V^\pi(s_t) + \alpha[R_t - V(s_t)] \quad (\text{MC})$$

# From incremental MC to TD (SUBA6.1)

Definition

$$\begin{aligned} V^\pi(\mathbf{s}) &= E_\pi\{R_t | \mathbf{s}_t = \mathbf{s}\} \\ &= E_\pi\{\sum_{k=0}^T \gamma^k r_{t+k+1} | \mathbf{s}_t = \mathbf{s}\} \end{aligned}$$

Incremental estimation from experience

$$V^\pi(\mathbf{s}_t) = V^\pi(\mathbf{s}_t) + \alpha[R_t - V(\mathbf{s}_t)] \quad (\text{MC})$$



# From incremental MC to TD (SUBA6.1)

Definition

$$\begin{aligned}
 V^\pi(\mathbf{s}) &= E_\pi\{R_t | \mathbf{s}_t = \mathbf{s}\} \\
 &= E_\pi\{\sum_{k=0}^T \gamma^k r_{t+k+1} | \mathbf{s}_t = \mathbf{s}\} \\
 &= E_\pi\{r_{t+1} + \gamma \sum_{k=1}^T \gamma^{k-1} r_{t+k+1} | \mathbf{s}_t = \mathbf{s}\}
 \end{aligned}$$

Incremental estimation from experience

$$V^\pi(\mathbf{s}_t) = V^\pi(\mathbf{s}_t) + \alpha[R_t - V(\mathbf{s}_t)] \quad (\text{MC})$$

# From incremental MC to TD (SUBA6.1)

Definition

$$\begin{aligned}
 V^\pi(\mathbf{s}) &= E_\pi\{R_t | \mathbf{s}_t = \mathbf{s}\} \\
 &= E_\pi\{\sum_{k=0}^T \gamma^k r_{t+k+1} | \mathbf{s}_t = \mathbf{s}\} \\
 &= E_\pi\{r_{t+1} + \gamma \sum_{k=1}^T \gamma^{k-1} r_{t+k+1} | \mathbf{s}_t = \mathbf{s}\} \\
 &= E_\pi\{r_{t+1} + \gamma V^\pi(\mathbf{s}_{t+1}) | \mathbf{s}_t = \mathbf{s}\}
 \end{aligned}$$

Bellman equation

Incremental estimation from experience

$$V^\pi(\mathbf{s}_t) = V^\pi(\mathbf{s}_t) + \alpha[R_t - V(\mathbf{s}_t)] \quad (\text{MC})$$

# From incremental MC to TD (SUBA6.1)

Definition

$$\begin{aligned}
 V^\pi(\mathbf{s}) &= E_\pi\{R_t | \mathbf{s}_t = \mathbf{s}\} \\
 &= E_\pi\{\sum_{k=0}^T \gamma^k r_{t+k+1} | \mathbf{s}_t = \mathbf{s}\} \\
 &= E_\pi\{r_{t+1} + \gamma \sum_{k=1}^T \gamma^{k-1} r_{t+k+1} | \mathbf{s}_t = \mathbf{s}\} \\
 &= E_\pi\{r_{t+1} + \gamma V^\pi(\mathbf{s}_{t+1}) | \mathbf{s}_t = \mathbf{s}\}
 \end{aligned}$$

Bellman equation

Incremental estimation from experience

$$V^\pi(\mathbf{s}_t) = V^\pi(\mathbf{s}_t) + \alpha[R_t - V(\mathbf{s}_t)] \quad (\text{MC})$$

# From incremental MC to TD (SUBA6.1)

Definition

$$\begin{aligned}
 V^\pi(\mathbf{s}) &= E_\pi\{R_t | \mathbf{s}_t = \mathbf{s}\} \\
 &= E_\pi\{\sum_{k=0}^T \gamma^k r_{t+k+1} | \mathbf{s}_t = \mathbf{s}\} \\
 &= E_\pi\{r_{t+1} + \gamma \sum_{k=1}^T \gamma^{k-1} r_{t+k+1} | \mathbf{s}_t = \mathbf{s}\} \\
 &= E_\pi\{r_{t+1} + \gamma V^\pi(\mathbf{s}_{t+1}) | \mathbf{s}_t = \mathbf{s}\} \\
 &\text{Bellman equation}
 \end{aligned}$$

Incremental estimation from experience

$$V^\pi(\mathbf{s}_t) = V^\pi(\mathbf{s}_t) + \alpha[R_t - V(\mathbf{s}_t)] \quad (\text{MC})$$

$$V^\pi(\mathbf{s}_t) = V^\pi(\mathbf{s}_t) + \alpha[r_{t+1} + \gamma V(\mathbf{s}_{t+1}) - V(\mathbf{s}_t)] \quad (\text{TD})$$

## From incremental MC to TD (SUBA6.1)

$$\begin{aligned}
 V^\pi(s) &= E_\pi\{R_t | s_t = s\} \\
 &= E_\pi\{\sum_{k=0}^T \gamma^k r_{t+k+1} | s_t = s\} \\
 &= E_\pi\{r_{t+1} + \gamma \sum_{k=1}^T \gamma^{k-1} r_{t+k+1} | s_t = s\} \\
 &= E_\pi\{r_{t+1} + \gamma V^\pi(s_{t+1}) | s_t = s\} \\
 &\text{Bellman equation}
 \end{aligned}$$

Incremental estimation from experience

$$V^\pi(s_t) = V^\pi(s_t) + \alpha[R_t - V(s_t)] \quad (\text{MC})$$

$$V^\pi(s_t) = V^\pi(s_t) + \alpha[r_{t+1} + \gamma V(s_{t+1}) - V(s_t)] \quad (\text{TD})$$

Initialize  $V(s)$  arbitrarily,  $\pi$  to the policy to be evaluated

Repeat (for each episode):

  Initialize  $s$

  Repeat (for each step of episode):

$a \leftarrow$  action given by  $\pi$  for  $s$

    Take action  $a$ ; observe reward,  $r$ , and next state,  $s'$

$V(s) \leftarrow V(s) + \alpha[r + \gamma V(s') - V(s)]$

$s \leftarrow s'$

  until  $s$  is terminal

Figure : Policy evaluation with  $TD(0)$ .

# From incremental MC to TD (SUBA6.1)

Definition

$$V^\pi(s) = E_\pi\{R_t | s_t = s\}$$

$$= E_\pi\{\sum_{k=0}^T \gamma^k r_{t+k+1} | s_t = s\}$$

$$= E_\pi\{r_{t+1} + \gamma \sum_{k=1}^T \gamma^{k-1} r_{t+k+1} | s_t = s\}$$

$$= E_\pi\{r_{t+1} + \gamma V^\pi(s_{t+1}) | s_t = s\}$$

Bellman equation

Incremental estimation from experience

$$V^\pi(s_t) = V^\pi(s_t) + \alpha [R_t - V(s_t)] \quad (\text{MC})$$

$$V^\pi(s_t) = V^\pi(s_t) + \alpha [r_{t+1} + \gamma V(s_{t+1}) - V(s_t)] \quad (\text{TD})$$

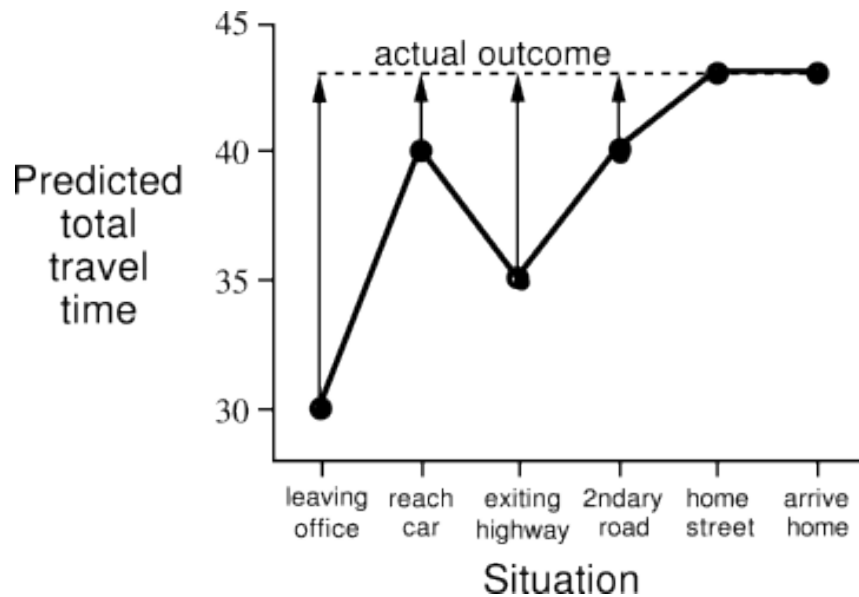


Figure : Driving home (MC)

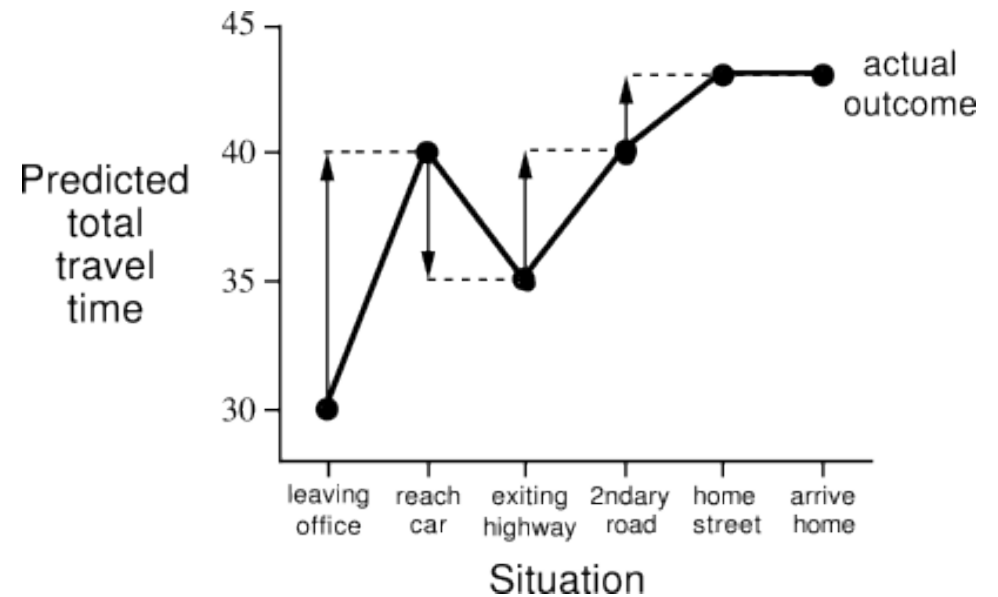


Figure : Driving home (TD)

# TD learning for Q-Values (SARSA) (SUBA6.4)

$$V^\pi(\mathbf{s}_t) = V^\pi(\mathbf{s}_t) + \alpha[r_{t+1} + \gamma V(\mathbf{s}_{t+1}) - V(\mathbf{s}_t)] \quad (1)$$

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = Q^\pi(\mathbf{s}_t, \mathbf{a}_t) + \alpha[r_{t+1} + \gamma Q(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}) - Q(\mathbf{s}_t, \mathbf{a}_t)] \quad (2)$$

# TD learning for Q-Values (SARSA) (SUBA6.4)

$$V^\pi(\mathbf{s}_t) = V^\pi(\mathbf{s}_t) + \alpha[r_{t+1} + \gamma V(\mathbf{s}_{t+1}) - V(\mathbf{s}_t)] \quad (1)$$

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = Q^\pi(\mathbf{s}_t, \mathbf{a}_t) + \alpha[r_{t+1} + \gamma Q(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}) - Q(\mathbf{s}_t, \mathbf{a}_t)] \quad (2)$$

- Incremental update requires  $(\mathbf{s}_t, \mathbf{a}_t, r_{t+1}, \mathbf{s}_{t+1}, \mathbf{a}_{t+1})$ 
  - Known as “SARSA”
  - If  $\mathbf{s}_{t+1}$  is terminal, then  $Q(\mathbf{s}_{t+1}, \mathbf{a}_{t+1})$  is zero



# TD learning for Q-Values (SARSA) (SUBA6.4)

$$V^\pi(\mathbf{s}_t) = V^\pi(\mathbf{s}_t) + \alpha[r_{t+1} + \gamma V(\mathbf{s}_{t+1}) - V(\mathbf{s}_t)] \quad (1)$$

$$Q^\pi(\mathbf{s}_t, \mathbf{a}_t) = Q^\pi(\mathbf{s}_t, \mathbf{a}_t) + \alpha[r_{t+1} + \gamma Q(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}) - Q(\mathbf{s}_t, \mathbf{a}_t)] \quad (2)$$

- Incremental update requires  $(\mathbf{s}_t, \mathbf{a}_t, r_{t+1}, \mathbf{s}_{t+1}, \mathbf{a}_{t+1})$ 
  - Known as “SARSA”
  - If  $\mathbf{s}_{t+1}$  is terminal, then  $Q(\mathbf{s}_{t+1}, \mathbf{a}_{t+1})$  is zero

Initialize  $Q(s, a)$  arbitrarily

Repeat (for each episode):

  Initialize  $s$

  Choose  $a$  from  $s$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

  Repeat (for each step of episode):

    Take action  $a$ , observe  $r, s'$

    Choose  $a'$  from  $s'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma Q(s', a') - Q(s, a)]$

$s \leftarrow s'; a \leftarrow a';$

  until  $s$  is terminal

Figure : SARSA: an on-policy TD control algorithm