



One again  $S_0, A_0, R_1, S_1, A_1, R_2, \dots$

- Recall:

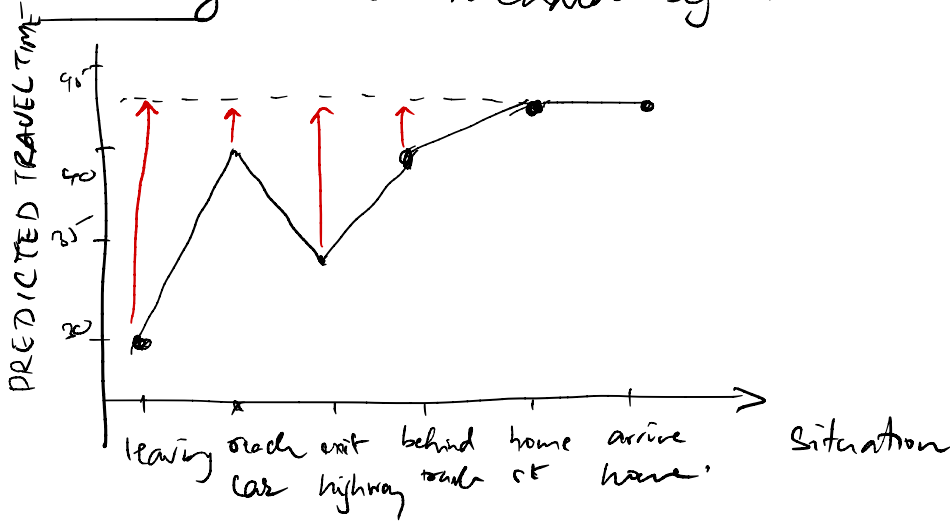
$$TD(0): \quad V(S_t) \leftarrow V(S_t) + \alpha \left( \overbrace{R_{t+1} + V(S_{t+1}) - V(S_t)}^{TD \text{ error}} \right)$$

Example:

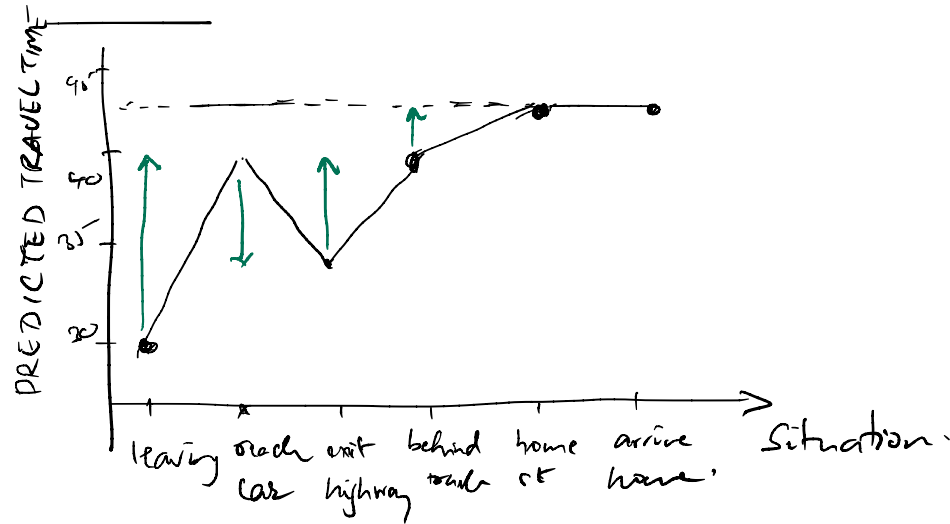
Driving Home example: (Barto & Sutton, David Silver)

State	Elapsed time	Predicted time to go	Predicted Total time
leaving office	0	30	30
reach car park	5	35	40
exit highway	20	15	35
behind truck	30	10	40
home street	40	3	43
arrive home	43	0	43

Changes recommended by MC.



Changes recommended by TD:



## Batch MC & TD:

MC & TD convergence;  $V(s) \rightarrow V_{\pi}(s)$  as experience  $\rightarrow \infty$ ;

- Batch selection for finite experience<sup>2</sup>

$$\left. \begin{array}{cccc} s_1^1 & a_1^1 & r_2^1 & \dots & s_T^1 \\ \vdots & & & & \\ s_1^k & a_1^k & r_2^k & \dots & s_T^k \end{array} \right\}$$

- Repeatedly sample episode  $k \in [1, K]$
- Apply MC or TD(0) to episode  $k$ .

What do MC & TD converge to?

Ex: Two states; no discounting; episodes of experience;

A, 0, B, 0

B, 1

B, 1

B, 1

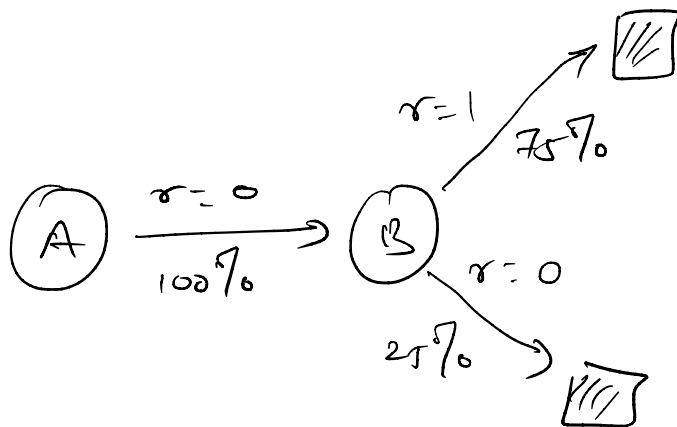
B, 1

B, 1

B, 1

B, 0

$$\begin{array}{l}
 \text{MC} \\
 \text{TD}
 \end{array}
 \left[
 \begin{array}{l}
 V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t)) \\
 V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))
 \end{array}
 \right.$$



MC converges to solution with min mean-sq error.

Best fit to observed returns;

$$\sum_{k=1}^K \sum_{t=1}^{T_k} (G_t^k - V(S_t^k))^2$$

for the example  $V(A) = 0;$

TD(0) converges to a solution of max likelihood Markov model;

$\langle S, A, \hat{P}, \hat{R}, \gamma \rangle$  which best fits the data.

$$P_{s,s'}^a = \frac{1}{N(s,a)} \sum_{k=1}^K \sum_{t=1}^{T_k} \mathbb{1} \left( \begin{matrix} s_t^k & a_t^k & s_{t+1}^k \\ \text{"} & \text{"} & \text{"} \\ s & a & s' \end{matrix} \right)$$

$$R_s^a = \frac{1}{N(s,a)} \sum_{k=1}^K \sum_{t=1}^{T_k} \mathbb{1} \left( \begin{matrix} s_t^k & a_t^k \\ \text{"} & \text{"} \\ s & a \end{matrix} \right) r_t^k.$$

for the example,  $V(A) = 0.75;$

ADV of TD - exploits Markov property.  
more effective in Markov env.

ADV of MC: Does not exploit Markov property.  
effective in non Markov env.

Last time:

TD(n);

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$

• n-step TD learning:

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t^n - V(S_t)).$$

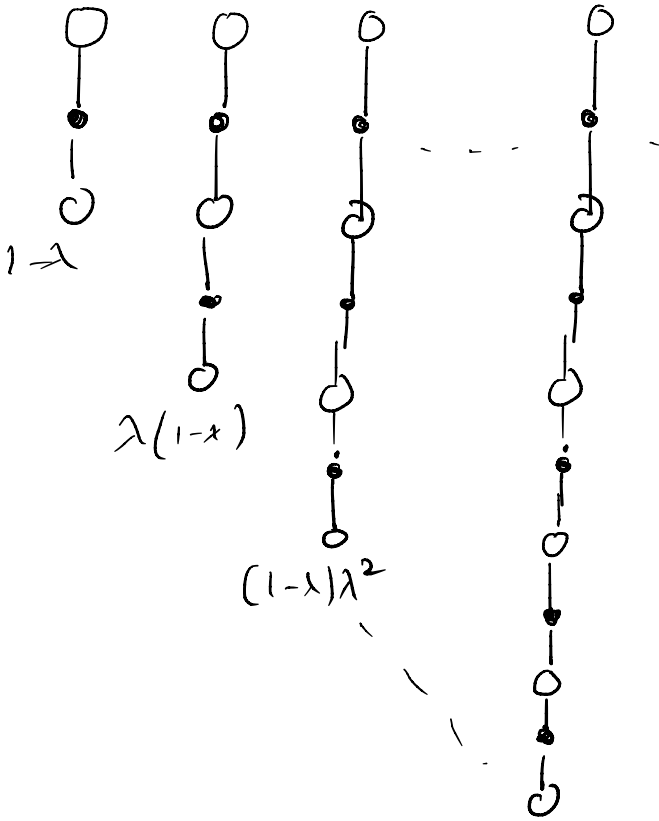
Now TD( $\lambda$ ):

• Average n-step returns over different n;

Ex:  $\frac{1}{2} G^{(2)} + \frac{1}{2} G^{(4)}$

• Can we combine such information from all time steps?

TD( $\lambda$ ): (From Silver)



$$G_T^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_T^n$$

$\lambda$  return used for update.

$$V(s_t) \leftarrow V(s_t) + \alpha (G_T^\lambda - V(s_t))$$

Needs entire episode.



## Eligibility traces:

- Which states should get credit for current return?

Frequency Heuristic

Assign credit to frequent states

Recency.

Assign credit to recent states;

Eligibility trace: Combines these heuristics;

$$E_0(s) = 0$$

$$E_t(s) = \gamma \lambda E_{t-1}(s) + \mathbb{1}(s_t = s)$$

- Idea: Keep an eligibility trace for all states;

Update  $V(s)$   $\forall s$ , in proportion to  $\delta_t$ , and  $E_t(s)$

Backward view TD( $\lambda$ ).

$$\delta_t = R_{t+1} + \gamma V(s_{t+1}) - V(s_t)$$

$$V(s) \leftarrow V(s) + \alpha \delta_t E_t(s).$$

• Last time - we saw function approximation.

⇓

unavoidable when we want to solve large problems;

Backgammon -  $10^{20}$  states;  
Go  $10^{170}$  states.

- idea:  $\hat{V}(s, \omega) \approx v_{\pi}(s)$ ; value function approximation  
 $\hat{Q}(s, a, \omega) \approx q_{\pi}(s, a)$  - Q value approximation

Example: Represent a state by a feature

vector  $\begin{pmatrix} x_1(s) \\ \vdots \\ x_n(s) \end{pmatrix}$

A linear approximation of value function is

$$\hat{V}(s, \omega) = \sum_{i=1}^n \omega_i x_i(s)$$

• Natural objective function to minimize

$$\mathcal{L}(\pi, \omega) \triangleq \mathbb{E}_{\pi} \left( (v_{\pi}(s) - x(s)^T \omega)^2 \right)$$

we don't know this unless in supervised learning!

Idea: Substitute a target for this;

$$\nabla_{\omega} \mathcal{L} = \mathbb{E}_{\pi} \left[ 2 \left( (v_{\pi}(s) - x(s)^T \omega) \cdot x(s) \right) \right]$$

$$\therefore \omega \leftarrow \omega - \alpha \left( v_{\pi}(s) - x(s)^T \omega \right) x(s).$$

Stochastic gradient using one sample;

In MC:  $v_{\pi}(s)$  replaced by  $G_t$ ;

In TD0: target replaced by  $R_{t+1} + \gamma \hat{V}(S_{t+1}, \omega)$

update:

$$\omega \leftarrow \omega - \alpha \left( R_{t+1} + \gamma \hat{V}(S_{t+1}, \omega) - \hat{V}(S_t, \omega) \right) x(s)$$

## Monte-Carlo + function approximation:

- $G_t$  - is an unbiased estimator of  $v_{\pi}(S_t)$ ;
- $\langle S_1, G_1 \rangle \langle S_2, G_2 \rangle \dots \langle S_T, G_T \rangle$  - thought of as supervised training data;

$$\Delta \omega \leftarrow \alpha (G_t - \hat{v}(S_t, \omega)) \nabla_{\omega} \hat{v}(S_t, \omega)$$

↑

works even if this is not a linear fn of  $\omega$ .

- MC - converges to a local opt
- MC " even if the function approximator is non linear;
- Can do TD learning with function value approximation  
TD target  $R_{t+1} + \gamma \hat{v}(S_{t+1}, \omega)$  current  $\omega$   
Biased sample of  $v_{\pi}(S_t)$ .
- Can still do supervised learning as before

$\langle S_1, R_1 + \gamma \hat{V}(S_2, \omega) \rangle, \langle S_2, R_2 + \gamma \hat{V}(S_3, \omega) \rangle, \dots$

Apply supervised learning to above;

$$TD(0): \Delta \omega = \alpha (R + \gamma \hat{V}(S', \omega) - \hat{V}(S, \omega)) \nabla_{\omega} \hat{V}(S, \omega).$$

$\omega \leftarrow \omega - \alpha \Delta \omega = \omega - \alpha \Phi^T x(S)$ , in the linear case;

Thm: Linear TD(0) converges to global optimum;

Can do the same with  $TD^{\lambda}$

$$\Delta \omega = \alpha (\overset{\lambda}{G}_t - \hat{V}(S_t, \omega)) \nabla_{\omega} \hat{V}(S_t, \omega)$$

• In backward linear TD( $\lambda$ ):

$$S_t = R_{t+1} + \gamma \hat{V}(S_{t+1}, \omega) - \hat{V}(S_t, \omega)$$

$$E_t = \underset{\substack{\uparrow \\ \text{vector}}}{\gamma \lambda E_{t-1}} + \underset{\substack{\uparrow \\ \text{a vector}}}{x(S_t)}$$

$$\Delta w = \alpha \delta_t E_t;$$

Last time: Compatible function approximator for state-value function.

• In general:  $\hat{v}(s, A, w) \approx v_{\pi}(s, A)$ .

• Minimize least square error;

• Examples:  $x(s, A) = \begin{pmatrix} x_1(s, A) \\ \vdots \\ x_n(s, A) \end{pmatrix}$

output value a linear function of these features,

• Use stochastic gradient as before;

• Targets for gradient updates: Like before;

$$G_t, R_{t+1} + \gamma \hat{v}(s_{t+1}, A_{t+1}, w), \delta_t \dots$$

- The previous algorithms - a single sample used as an unbiased estimator.

## BATCH REINFORCEMENT LEARNING:

- Given value function approximation  $\hat{v}(s, \omega) \approx v_{\pi}(s)$  and experience  $\mathcal{D} = \{ \langle s_1, v_1^{\pi} \rangle, \langle s_2, v_2^{\pi} \rangle, \dots, \langle s_T, v_T^{\pi} \rangle \}$

GOAL: Find parameters  $\omega$  which give the best fitting value function  $\hat{v}(s, \omega)$ .

## LEAST SQUARE:

$$LS(\omega) = \sum_{t=1}^T (v_t^{\pi} - \hat{v}(s_t, \omega))^2$$

### Algorithm:

- Sample state, value from experience  $\langle s, v^{\pi} \rangle \sim \mathcal{D}$
- Apply SGD update 
$$\Delta \omega = \alpha (v^{\pi} - \hat{v}(s, \omega)) \nabla_{\omega} \hat{v}(s, \omega)$$



- Converges to  $\omega^* = \underset{\omega}{\operatorname{argmin}} LS(\omega)$

## Experience replay in DQN:

- Take action  $a_t$  according to  $\epsilon$ -greedy.
- Store  $(s_t, a_t, r_{t+1}, s_{t+1})$  in  $\mathcal{D}$ -replay memory
- Sample random mini-batch of transitions  $(s, a, r, s')$  from  $\mathcal{D}$
- Compute Q-learning targets w.r.t old  $\omega$
- OPTIMIZE MSE b/w Q-network & Q-learning targets.

$$L_i(\omega_i) = \mathbb{E}_{s, a, r, s' \in \mathcal{D}_i} \left[ \left( r + \gamma \max_{a'} Q(s', a', \omega) - Q(s, a, \omega_i) \right)^2 \right]$$

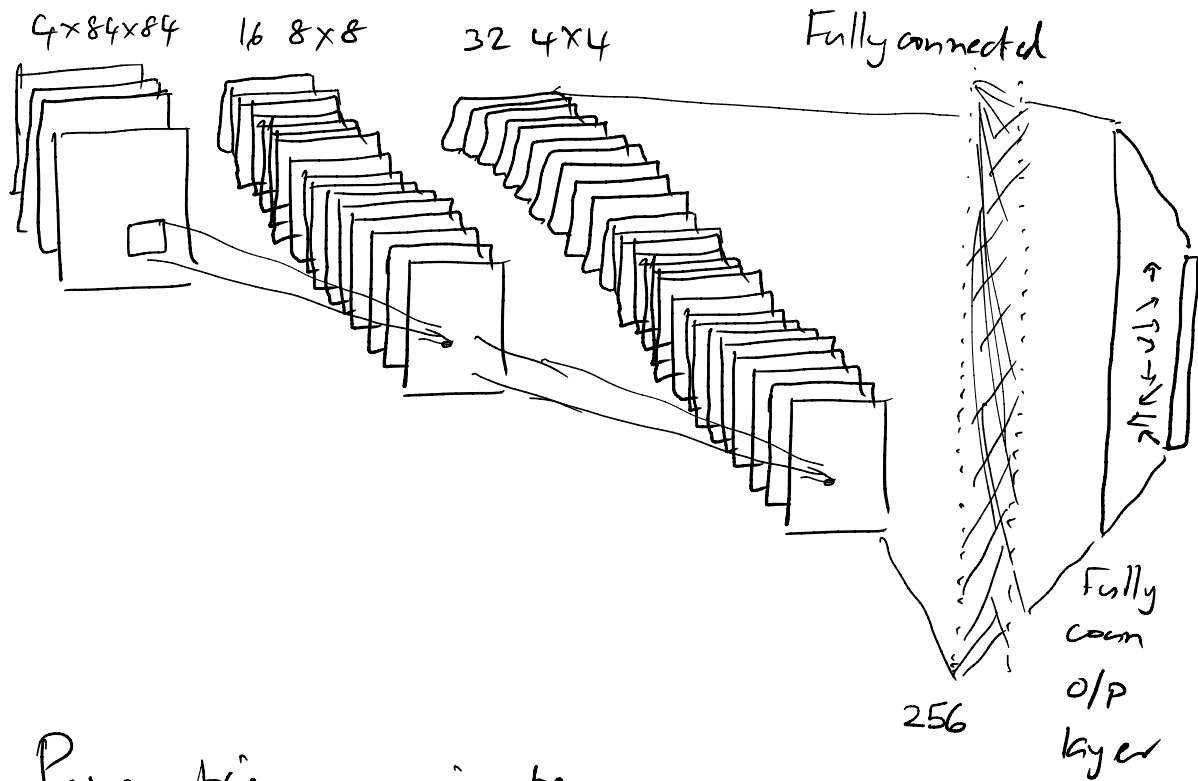
Use stochastic gradient descent;

DQN: Nature - 2015; Human level control through deep RL

Sets up a DQN to approximate the optimal action-value function

$$Q^*(s, a) = \max_{\pi} \mathbb{E} \left[ r_t + \gamma V_{t+1} + \dots \mid s_t = s, a_t = a, \pi \right]$$

- \* Experience replay - biologically inspired. Since it randomizes over the data it removes correlations in the observation sequences.
- \* Uses an iterative update, adjusting action-value towards target values which are only periodically updated and so reducing correlations with the target.
- End to end learning of values  $Q(s, a)$  from pixels



Parametrize approximate  
action value function  $Q(s, a, \theta;)$

$\theta_i$  - are weights of the DQN.

- Perform experience replay by storing agents experiences  $e_t = (s_t, a_t, r_t, s_{t+1})$  at each time step  $t$  in  $\mathcal{D}_t = \{e_1, \dots, e_t\}$ .

Q-learning update at iteration  $i$  uses loss fn

$$L_i(\theta_i) = \mathbb{E}_{(s, a, r, s') \sim \mathcal{U}(\mathcal{D})} \left[ \left( r + \gamma \max_{a'} Q(s', a'; \bar{\theta}_i) - Q(s, a; \theta_i) \right)^2 \right]$$

Network parameters used to compute the target at iteration  $i$ .

Parameters of the network in iteration  $i$

- Target network parameters  $\bar{\theta}_i$  are updated with the Q-network parameters  $\theta_i$  every  $C$  steps and are held fixed b/w individual updates.

## Algorithm:

- At each time step, agent selects  $a_t$ , a legal game action and passes it to the emulator.
- Emulator modifies its internal state & game score.
- The agent observes an image from the emulator, the vector of pixel values representing current screen. Agent receives reward  $r_t$ , the change in game score.

