

## Eligibility Traces

- basics of RL
- TD( $\lambda$ ),  $\lambda$  is an eligibility trace
- + easily combinable with any TD (Sutton, Barto, etc.)
- $\lambda$  unifies & generalizes TD & MC methods
- $\lambda = 0$  MC,  $\lambda = 1$  TD
- + in between, method often better than either end
- how diff from n-step methods?
- an eligibility trace  $z_t \in \mathbb{R}^d$
- x when a weight in  $w$  participates, its component in  $z$  bumps up, then fades away according to  $\lambda$ , called trace decay param.
- \* computationally more efficient
- \* learning is uniform  $\lambda$  in a single step rather than time delayed
- easier to implement
- + most also require forward view, what not always available
- + can do almost, if not, exactly, the same, using backward views
- look at state values, prediction, extend action values, then control

## The $\lambda$ -return

- $G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^n R_{t+n} + \gamma^n \hat{v}(S_{t+n-1}, w_{t+n-1})$
- valid update target is not just any n-step, but also any any n-step return for diff n's
- + e.g.  $\frac{1}{2} G_{t:t+2} + \frac{1}{2} G_{t:t+1}$
- + produce compact errors & guaranteed convergence
- x e.g. using 1-step &  $\infty$ -step return
- any update of simple = compound update
- + updates can only occur when terminal component is done
- TD( $\lambda$ ) only n-step updates
- $G_t^\lambda = (1-\lambda) \sum_{n=0}^{\infty} \lambda^n G_{t:t+n}$
- +  $(1-\lambda)$  normalizes to ensure sum = 1
- x each return fades with  $\lambda$
- x can reach post-termination state
- $G_t^\lambda = (1-\lambda) \sum_{n=0}^{T-t-1} \lambda^n G_{t:t+n} + \lambda^{T-t} G_t$
- off-line  $\lambda$ -return algorithm
- $w_{t+1} = w_t + \alpha [G_t^\lambda - \hat{v}(S_t, w_t)] \nabla \hat{v}(S_t, w_t)$
- performance of  $\lambda$  similar to choice of n
- forward-view model

## TD( $\lambda$ )

- computationally inspired backwards view
- 3 improvements
- 1) update at every step, not end. Better, same estimate
- 2) computation distributed across time, not bulk at end
- 3) continuing & episodic application
- $z_t$  initialized with 0 values

$$z_{t+1} = z_t + \gamma \lambda z_{t+1} + \nabla \hat{v}(S_t, w_t)$$

- + in linear approx,  $\nabla \hat{v}$  is just  $\hat{x}$ , so  $z_t$  is fading sum of on feature vector
- trace indicates eligibility of components in weight vector for underlying learning from a reinforcing event.

$$\delta_t = R_{t+1} + \gamma \hat{v}(S_{t+1}, w_t) - \hat{v}(S_t, w_t)$$

$$w_{t+1} = w_t + \alpha \delta_t z_t$$

Semi-gradient TD( $\lambda$ ) for  $\hat{v} \approx v_w$

inputs:

- policy  $\pi$  for env
- differentiable func  $\hat{v}: S \times \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $\hat{v}(\text{term}) = 0$
- also params:

- $\alpha > 0$
- $\lambda \in [0, 1]$
- $w \in \mathbb{R}^d$  arbitrarily

for each ep

init  $S$

$z = 0$

for each step in ep:

$$A = \pi(S)$$

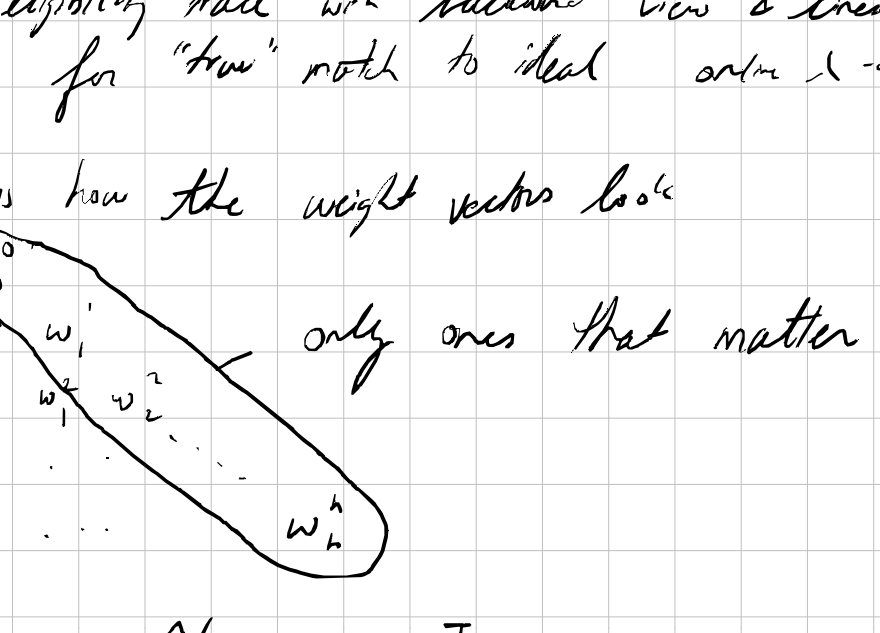
$$S', R = \text{env.step}(A)$$

$$z = \gamma \lambda z + \nabla \hat{v}(S, w)$$

$$\delta = R + \gamma \hat{v}(S', w) - \hat{v}(S, w)$$

$$w = w + \alpha \delta z$$

$$S = S'$$



- backwords, or mechanistic, view
- smaller, or more distant states are given less credit for contributing to current error

- for undiscounted episodic MC action,  $\lambda=1$   $\gamma=1$
- + if  $\lambda=1$ , collect TD(1)
- + TD(1) is more general, easier to apply
- x can do discounting, can do online, can do incremental
- x TD(1) can stay away from bad actions mid learn, MC cannot

- linear TD( $\lambda$ ) converges in on-policy, at
- $$\overline{V_E}(w_{t+1}) \leq \frac{1-\gamma\lambda}{1-\gamma} \min_w \overline{V_E}(w)$$
- for continuous, discounted.
- + approaches minimum error as  $\lambda \rightarrow 1$ , but impractical

## n-step Truncated $\lambda$ -return

### Methods

- problem:  $\lambda$ -return still unknown till end of ep
- + in continuing case, never known
- + however, discount of future rewards falls off by  $\gamma\lambda$ , so some approximation would be good
- truncated  $\lambda$ -return

$$G_{t:t+h}^\lambda = (1-\lambda) \sum_{n=0}^{h-t-1} \lambda^n G_{t:t+n} + \lambda^{h-t} G_{t:t+h}$$

- + simple return up to horizon h
- family of n-step  $\lambda$ -returns when all  $1 \leq h \leq n$  returns are used instead of just nth return
- + for state values, called Truncated TD( $\lambda$ ) or TTD( $\lambda$ )

- x  $w_{t+n} = w_{t+n-1} + \alpha [G_{t:t+n}^\lambda - \hat{v}(S_t, w_{t+n-1})] \nabla \hat{v}(S_t, w_{t+n-1})$

- efficient impl can be done using k-returns:
- $$G_{t:t+k}^\lambda = \hat{v}(S_t, w_{t+1}) + \sum_{i=t+1}^{t+k-1} (\gamma\lambda)^{i-t} \delta_i$$
- $$\delta_i = R_{i+1} + \gamma \hat{v}(S_{i+1}, w_t) - \hat{v}(S_i, w_{t-1})$$

## Relating Updates: Online $\lambda$ -return

### Algorithm

- how to find the right n?
- + updates are done every step, with better info
- + t+1, we know  $G_{0:1}$
- + t+2, we know  $G_{0:2}$ ,  $G_{1:2}$
- x weight vector changes at each horizon,  $w_t^h$

$$h=1, w_1^1 = w_0^1 + \alpha [G_{0:1}^\lambda - \hat{v}(S_0, w_0^1)] \nabla \hat{v}(S_0, w_0^1)$$

$$h=2, w_1^2 = w_0^2 + \alpha [G_{0:2}^\lambda - \hat{v}(S_0, w_0^2)] \nabla \hat{v}(S_0, w_0^2)$$

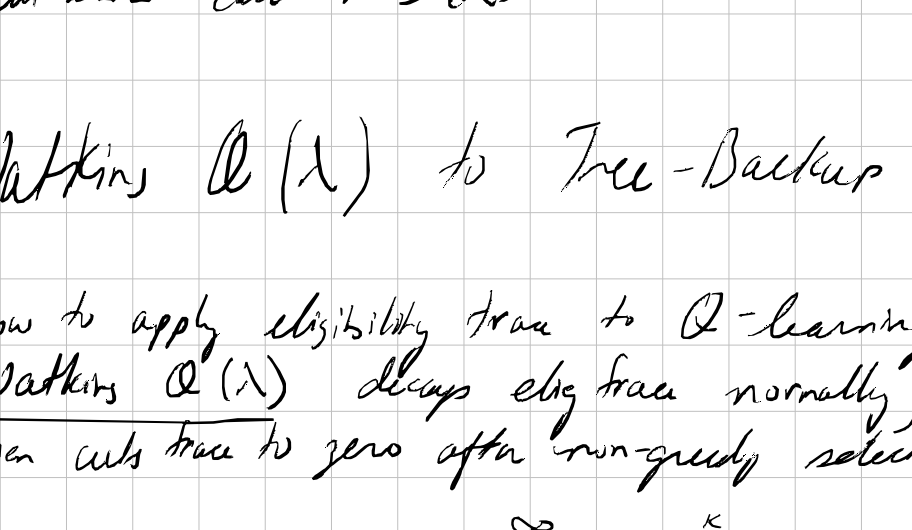
$$w_2^2 = w_1^2 + \alpha [G_{1:2}^\lambda - \hat{v}(S_1, w_1^2)] \nabla \hat{v}(S_1, w_1^2)$$

$$w_{t+1}^h = w_t^h + \alpha [G_{t:t+h}^\lambda - \hat{v}(S_t, w_t^h)] \nabla \hat{v}(S_t, w_t^h)$$

- thus,  $w_t \approx w_t^h$ , online  $\lambda$ -return algo
- strictly more computationally complex because it maintains the weight vector between every timestep
- \* adv: better judgement at all time, even end of episode

## True Online TD( $\lambda$ )

- use eligibility trace with backward view & linear approx for "true" math to deal online  $\lambda$ -alg
- this is how the weight vectors look



- in linear,  $\hat{v}(S, w) = w^T \times (S)$ , then:
- $$w_{t+1} = w_t + \alpha \delta_t z_t + \alpha (w_t^T x_t - w_{t-1}^T x_t) (z_t - x_t)$$
- $$x_t = \lambda(S_t)$$
- $$z_t = \gamma \lambda z_{t-1} + (1 - \alpha \gamma \lambda z_{t-1}^T x_t) x_t$$

- technically called a dutch trace
- + previous trace called accumulating trace
- skipping Dutch Traces in MC (history learn)

## Sarsa( $\lambda$ )

- estimate  $\hat{q}(S_t, w)$ ?
- pretty straightforward, just use  $G_{t:t+n}$  for  $\hat{q}$ , then update for offline  $\lambda$ -return
- $$w_{t+1} = w_t + \alpha [G_{t:t+n}^\lambda - \hat{q}(S_t, A_t, w)] \nabla \hat{q}(S_t, A_t, w)$$
- $$G_t^\lambda = G_{t:\infty}^\lambda$$

- for Sarsa, same update,
- $$w_{t+1} = w_t + \alpha \delta_t z_t$$
- with  $z_{t+1} = 0$
- $$z_t = \gamma \lambda z_{t-1} + \nabla \hat{q}(S_t, A_t, w_t)$$

- true online TD( $\lambda$ ) exists, therefore, so does true online Sarsa( $\lambda$ )
- also truncated version, called forward Sarsa( $\lambda$ ) for model free control in Markov games

## Variable $\lambda$ & $\gamma$

- TD also, fixed form!
- $\lambda, \gamma$  aren't scalars, they are functions of state and action!
- +  $\lambda_t, \gamma_t$

$$\lambda: S \times A \rightarrow [0, 1] \quad \lambda_t = \lambda(S_t, A_t)$$

$$\gamma: S \rightarrow [0, 1] \quad \gamma_t = \gamma(S_t)$$

- x  $\gamma_t$  is now the termination function
- x now return is

$$G_t = R_{t+1} + \gamma_{t+1} G_{t+1}$$

$$= R_{t+1} + \gamma_{t+1} R_{t+2} + \gamma_{t+1} \gamma_{t+2} R_{t+3} + \dots$$

$$= \sum_{k=t}^{\infty} \left( \prod_{i=t+1}^k \gamma_i \right) R_{k+1}$$

$$\prod_{k=t}^{\infty} \gamma_k = 0 \quad \text{to assure finite sums}$$

- convenience of definition: episodic setting precludes in a single stream of experience
- + no post-termination states, start distributions, of termination times

- only affects solution strategy
- $$G_t^{\lambda, \gamma} = R_{t+1} + \gamma_{t+1} ((1-\lambda_{t+1}) \hat{v}(S_{t+1}, w_t) + \lambda_{t+1} G_{t+1}^{\lambda, \gamma})$$
- $$\lambda_{t+1} G_{t+1}^{\lambda, \gamma}$$

- + bootstrapping from state values (S's)
- + add as for action values
- + in English

"the return at time t is the immediate reward, undiscounted, unaffected by bootstrapping, plus a possible second term to the extent we are not discounting at the next state (i.e.  $\gamma_{t+1}$  is not 0, e.g. for term). If we aren't, we use some ratio of bootstrapping or  $\lambda$ -return of t+1"

- Expected Sarsa:
- $$G_t^{\lambda, \gamma} = R_{t+1} + \gamma_{t+1} ((1-\lambda_{t+1}) \bar{V}(S_{t+1}) + \lambda_{t+1} G_{t+1}^{\lambda, \gamma})$$
- $$\bar{V}_t(S) = \sum_a \pi(a|S) \hat{q}(S, a, w)$$

## Off Policy Traces With Control

### Variates

- incorporate importance sampling
- + only applicable for bootstrapped

$$G_t^{\lambda, \gamma} = p_t (R_{t+1} + \gamma_{t+1} ((1-\lambda_{t+1}) \hat{v}(S_{t+1}, w_t) + \lambda_{t+1} G_{t+1}^{\lambda, \gamma}))$$

$$+ (1-p_t) \hat{v}(S_t, w_t)$$

- can still be approximated by
- $$\delta_t^S = R_{t+1} + \gamma_{t+1} \hat{v}(S_{t+1}, w_t) - \hat{v}(S_t, w_t)$$
- $$G_t^{\lambda, \gamma} \approx \hat{v}(S_t, w_t) + \sum_{k=t}^{\infty} \delta_k^S \prod_{i=t+1}^k \gamma_i \lambda_i$$

- can easily be used to a forward view update
- $$w_{t+1} = w_t + \alpha (\delta_t^S - \hat{v}(S_t, w_t)) \nabla \hat{v}(S_t, w_t)$$
- $$\approx w_t + \alpha p_t \left( \sum_{k=t}^{\infty} \delta_k^S \prod_{i=t+1}^k \gamma_i \lambda_i \right) \nabla \hat{v}(S_t, w_t)$$

- eligibility trace with semi-gradient update
- creates a general on/off policy TD( $\lambda$ )
- + off-policy is just not guaranteed stability

- for action value, eligibility trace with Sarsa( $\lambda$ ) is state of art
- $\lambda=1$  still weaker than MC to learn updates
- + control out in expected value
- x there does exist a method for exact equivalence call PTD( $\lambda$ )

## Watkins Q( $\lambda$ ) to Tree-Backup( $\lambda$ )

- how to apply eligibility trace to Q-learning
- Watkins Q( $\lambda$ ) delays elig trace normally, then sets trace to zero after non-greedy selection
- $$G_t^{\lambda, \gamma} \approx \hat{q}(S_t, A_t, w_t) + \sum_{k=t}^{\infty} \delta_k^Q \prod_{i=t+1}^k \gamma_i \lambda_i \pi(A_i|S_i)$$

- final form TD( $\lambda$ )
- $$z_t = \gamma \lambda \tau(A_t|S_t) z_{t-1} + \nabla \hat{q}(S_t, A_t, w_t)$$

## Stable Off-policy Methods With

### Traces

- all use Gradient TD or Explicit TD with linear func. approx.
- GTD( $\lambda$ ) like TD(0)
- + learns  $w_t \rightarrow \hat{v}(S, w) = w_t^T \times (S) \approx v_{\pi}(S)$
- + even if data is from off policy b
- + updates

$$w_{t+1} = w_t + \alpha \delta_t^S z_t - \alpha \gamma_{t+1} (1 - \lambda_{t+1}) (z_t^T x_t) x_{t+1}$$

$$v_{t+1} = v_t + \beta \delta_t^S z_t - \beta (v_t^T x_t) x_t$$

- GTD( $\lambda$ ) is GTD( $\lambda$ ) for action values
- HTD( $\lambda$ ) is hybrid of GTD( $\lambda$ ) & TD( $\lambda$ )
- Explicit TD( $\lambda$ ) contains one step Explicit TD

## Implementation Issues

- most of  $z$  is zero, few nonzero
- generally not a problem

## Conclusion

- eligibility trace + TD error = efficient, incremental way of going from one step TD to full MC
- online is the biggest win

- can go from how t relates to t+1 (forward view) to how t+1 relates to t (backward view)

- computationally expensive but good for forward learning and delayed rewards on sparse data