

# Temporal Difference Updating without a Learning Rate



Marcus Hutter  
Canberra, ACT, 0200, Australia  
<http://www.hutter1.net/>



## Summary

- Reinforcement learning TD update:  $V_s^{t+1} = V_s^t + E_s^t \beta_t(s, s_{t+1})(r_t + \gamma V_{s_{t+1}}^t - V_s^t) \forall s$
- We derived learning rate  $\beta_t$  and eligibility trace  $E_s^t$  from statistical principles.
- In every setting that we have tested, superior performance & fewer parameters to tune.



Shane Legg  
6900, Lugano, Switzerland  
<http://www.vetta.org/shane/>



## Introduction & Background

### Abstract

In the field of reinforcement learning, temporal difference (TD) learning is perhaps the most popular way to estimate the future discounted reward of states. We derive an equation for TD learning from statistical principles. Specifically, we start with the variational principle and then bootstrap to produce an updating rule for discounted state value estimates. The resulting equation is similar to the standard equation for temporal difference learning with eligibility traces, so called TD( $\lambda$ ), however it lacks the parameter  $\alpha$  that specifies the learning rate. In the place of this free parameter there is now an equation for the learning rate that is specific to each state transition. We experimentally test this new learning rule against TD( $\lambda$ ) and find that it offers superior performance in various settings. Finally, we combine our update equation with both Watkins's Q( $\lambda$ ) and Sarsa( $\lambda$ ) and find that it again offers superior performance without a learning rate parameter.

## Temporal Difference Learning

- $s_k$ =state and  $r_k$ =reward in cycle  $k$ ,  $1 > \gamma$ =discount.
- Value of state  $s$  = expected future discounted reward = 
$$\bar{V}_s := \mathbf{E}\{r_k + \gamma r_{k+1} + \gamma^2 r_{k+2} + \dots | s_k = s\} = \mathbf{E}\{r_k + \gamma \bar{V}_{s_{k+1}} | s_k = s\} \quad (1)$$
- Given: History of states  $s_1, s_2, \dots, s_t$ , and observed rewards  $r_1, r_2, \dots, r_t$ .
- Goal: Compute estimate  $V_s^t$  of  $\bar{V}_s \forall s$ .
- Equation (1) suggests the following TD(0) learning algorithm 
$$V_{s_t}^{t+1} := V_{s_t}^t + \alpha (r_t + \gamma V_{s_{t+1}}^t - V_{s_t}^t),$$
- where  $\alpha$  is a parameter that controls the rate of learning.
- Shortcoming: At each time  $t$ , the value of only the last state  $s_t$  is updated.

## Eligibility Traces

- Improvement: TD with eligibility traces:
  - update value of all recently visited states proportional to their ...
- Eligibility:  $E_s^t := \gamma \lambda E_s^{t-1} + \delta_{s_t, s}$
- $\lambda$  controls the rate at which the eligibility trace is discounted.
- The TD( $\lambda$ ) update [Sutton 1988] is then, for all states  $s$ , 
$$V_s^{t+1} := V_s^t + \alpha E_s^t (r_t + \gamma V_{s_{t+1}}^t - V_s^t) \quad \forall s. \quad (2)$$
- TD makes intuitive sense and has various motivations [Sutton 1988], but is heuristic rather than a theoretically founded.

## The New Update Equation

### Main Novel Idea

- Idea: Derive a TD rule from statistical principles.
- Empirical future discounted reward of a state  $s_k$  at times  $k = 1, \dots, t$  is

$$v_k := \sum_{u=k}^{\infty} \gamma^{u-k} r_u,$$

- $v_k$  is unknowable, since it depends also on unknown future ( $u > t$ ) rewards.
- Goal: Estimate  $V_s^t$  close to  $\bar{V}_s$ .
- $\Rightarrow$  we would like  $V_s$  to be close to  $v_k$  for all  $k$  such that  $s = s_k$ .

- Also: Discount old evidence by  $\lambda \in (0, 1]$  in non-stationary environments:

$$\Rightarrow \text{Principle: minimize loss } L := \frac{1}{2} \sum_{k=1}^t \lambda^{t-k} (v_k - V_{s_k}^t)^2$$

- For stationary environments we may simply set  $\lambda = 1$  a priori.

## Derivation of the Value Estimate

$$\frac{\partial L}{\partial V_s^t} = 0 \Leftrightarrow V_s^t N_s^t = \sum_{k=1}^t \lambda^{t-k} \delta_{s_k, s} v_k \quad (3)$$

- $N_s^t := \sum_{k=1}^t \lambda^{t-k} \delta_{s_k, s} =$  discounted state visit counter.
- $v_k$  has a self-consistency property:  $v_k = \sum_{u=k}^{t-1} \gamma^{u-k} r_u + \gamma^{t-k} v_t$ .
- Substituting this into (3) we get:  $V_s^t N_s^t = R_s^t + E_s^t v_t$  (\*)
- $E_s^t := \sum_{k=1}^t (\lambda \gamma)^{t-k} \delta_{s_k, s} =$  the eligibility trace of state  $s$ ,
- $R_s^t := \sum_{u=1}^{t-1} \lambda^{t-u} E_s^u r_u =$  discounted reward with eligibility.
- Bootstrap: Replace unknowable  $v_t$  in (\*) by known  $V_{s_t}^t$  and solve w.r.t.  $V_s^t$ :

$$\text{Value estimate: } V_s^t = \frac{1}{N_s^t} \left[ R_s^t + E_s^t \frac{R_{s_t}^t - E_{s_t}^t V_{s_t}^t}{1 - \gamma} \right] \quad (4)$$

## Incremental Update Rule

Rewrite (4) as incremental update rule HL( $\lambda$ ):

- $V_s^{t+1} = V_s^t + E_s^t \beta_t(s, s_{t+1}) (r_t + \gamma V_{s_{t+1}}^t - V_s^t) \quad \forall s \quad (**)$
- $\beta_t(s, s_{t+1}) := \frac{1}{N_s^t} \cdot \frac{1}{1 - \gamma E_{s_{t+1}}^t / N_{s_{t+1}}^t} =$  Learning rate.
- $N_s^{t+1} = \lambda N_s^t + \delta_{s_{t+1}, s}, \quad E_s^{t+1} = \lambda \gamma E_s^t + \delta_{s_{t+1}, s}$ .
- $V_s^0 = N_s^0 = E_s^0 = R_s^0 = 0$ .

Explanation:

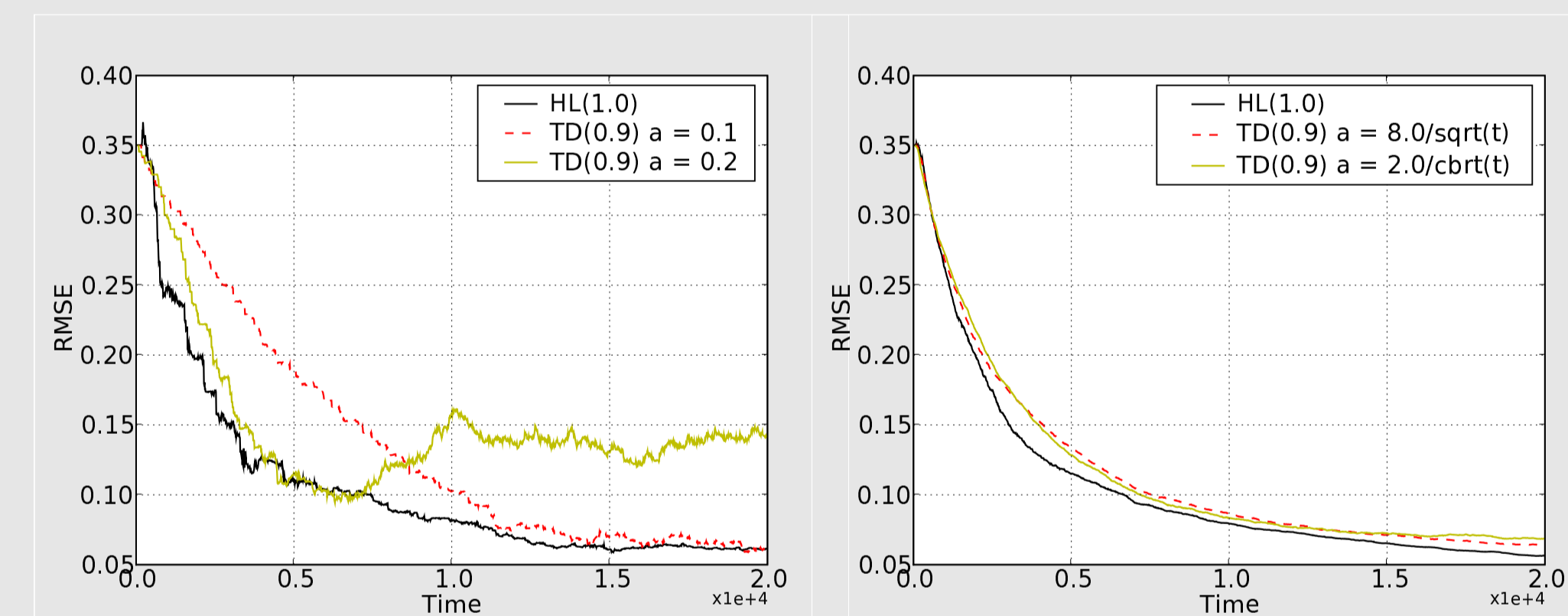
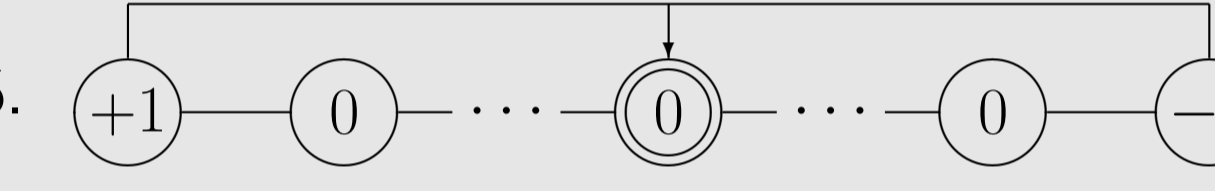
- (\*\*) is usual update equation for TD learning with eligibility trace (2),
- but with heuristic hand-tuned learning rate parameter  $\alpha$
- replaced by exact/automatic/dynamic  $\beta_t(s, s_{t+1})$ .
- First term in  $\beta_t$  is inversely proportional to the state visit counter  $N_s^t$ .
- The second term in  $\beta_t$  has bounded fluctuation between 1 and  $\frac{1}{1-\gamma}$ . It is "large" for recently visited states. It converges to 1 if and only if  $\lambda = 1$ .

## Experimental Evaluation

- Root mean square error (RMSE) between estimate  $V_s^t$  and exact value  $\bar{V}_s$ , averaged over  $s$ , are plotted.
- Best value of  $\lambda$  was used.

### A simple Markov Process

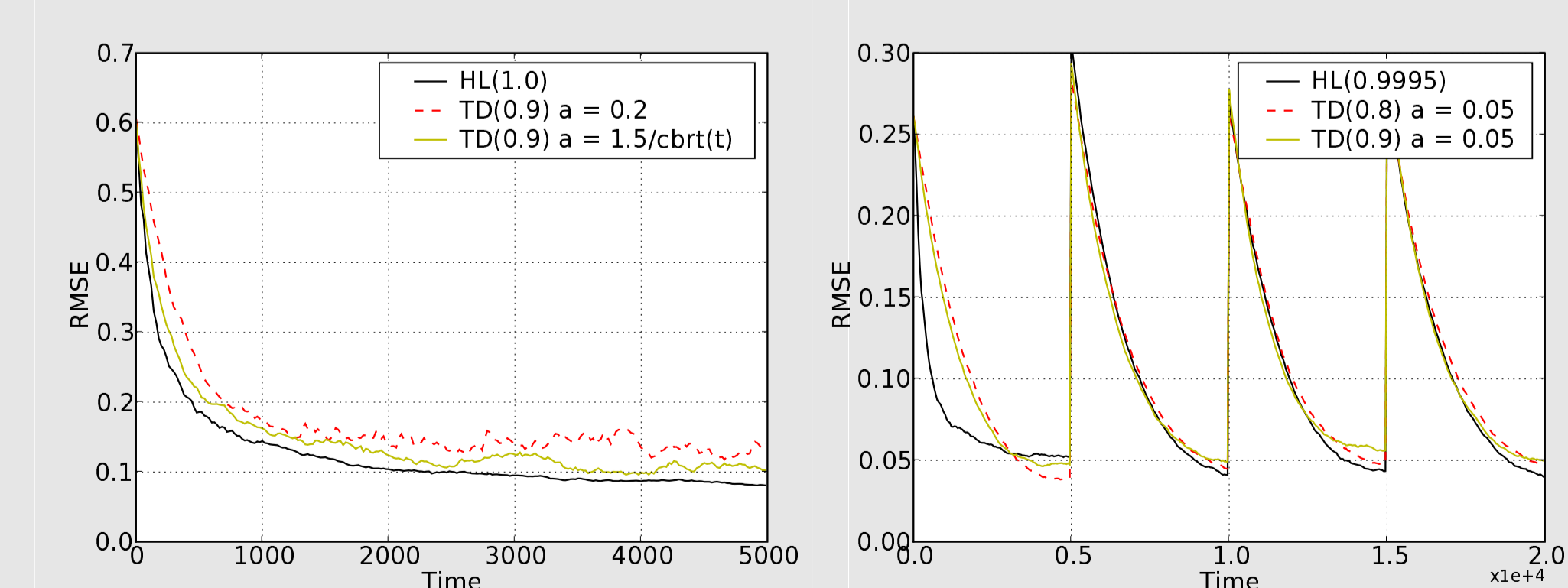
- Linear Markov chain with
- 51 states. Start at  $s = 25$ .
- No reward, except ...
- Reward +1 at  $s = 0$ , reward -1 at  $s = 50$ . In both cases reset to  $s = 25$ .



51 states. Discount  $\gamma = 0.99$ . Constant learning rate  $a = \alpha$ . HL( $\lambda$ ) signif. outperforms TD( $\lambda$ ). ( $\alpha_t \propto 1/t$  very poor).

- Similar results for larger and smaller Markov chains, and different  $\gamma$ .

## Random & Non-Stationary Markov Process



Random 50 state Markov process. Discount  $\gamma = 0.9$ . 21 state non-stationary Markov process, where reward function changed 5000 steps.

## New Reinforcement Learning Algorithm

### Algorithm for Markov Decision Processes

- Include actions: Markov process  $\sim$  Markov Decision Process (MDP).
- TD( $\lambda$ )  $\sim$  Sarsa( $\lambda$ ) [Rummery 1994].
- Combining Sarsa( $\lambda$ ) with HL( $\lambda$ )  $\Rightarrow$  our HLS( $\lambda$ ) algorithm.

### Algorithm HLS( $\lambda$ )

Initialise  $Q(s, a) = 0$ ,  $N(s, a) = 1$  and  $E(s, a) = 0$  for all  $s, a$   
Initialise  $s$  and  $a$

repeat

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

Choose  $a'$  by using  $\epsilon$ -greedy selection on  $Q(s', \cdot)$

$\Delta \leftarrow r + \gamma Q(s', a') - Q(s, a)$

$E(s, a) \leftarrow E(s, a) + 1$

$N(s, a) \leftarrow N(s, a) + 1$

for all  $s, a$  do

$\beta((s, a), (s', a')) \leftarrow \frac{1}{N(s, a)} \cdot \frac{1}{1 - \gamma E(s', a') / N(s', a')}$

end for

for all  $s, a$  do

$Q(s, a) \leftarrow Q(s, a) + \beta((s, a), (s', a')) E(s, a) \Delta$

$E(s, a) \leftarrow \gamma \lambda E(s, a)$

$N(s, a) \leftarrow \lambda N(s, a)$

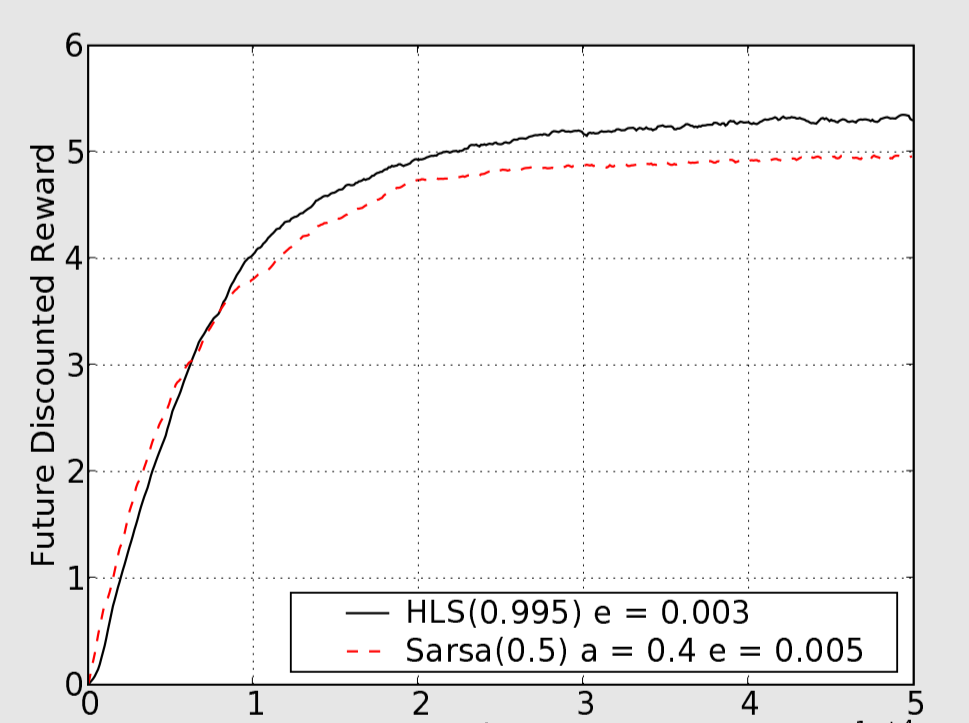
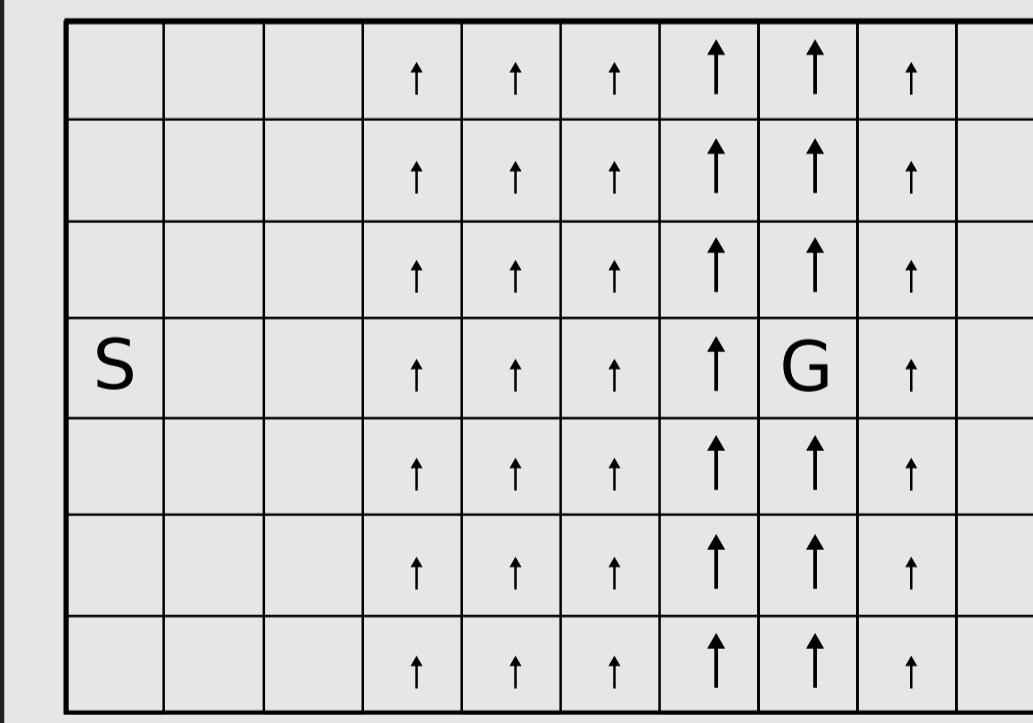
end for

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

until end of run

### Windy Gridworld

- Move from start state  $S$  to goal state  $G$  on a windy grid [Sutton 1998].
- Once the goal state  $G$  is reached, the agent jumps back  $S$  and receives a reward of 1 ( $\gamma = 0.99$ ).



- Despite considerable effort in tuning the parameters of Sarsa( $\lambda$ ), HLS( $\lambda$ ) consistently beats Sarsa( $\lambda$ ).
- Similar results for Q( $\lambda$ ) [Watkins 1989] but much smaller gap.

## Conclusions

- We have mathematically derived the equation for setting the learning rate in temporal difference learning with eligibility traces.
- The equation replaces the learning rate  $\alpha$ , which is a free parameter that normally has to be experimentally tuned by hand.
- In every setting that we have tested [(non)stationary Markov (Decision) processes], our new update equation has produced superior results.  $\Rightarrow$  superior performance and fewer parameters to tune