

# TD<sub>γ</sub>

George Konidaris  
[gdk@csail.mit.edu](mailto:gdk@csail.mit.edu)

(joint work with Scott Niekum and Phil Thomas)



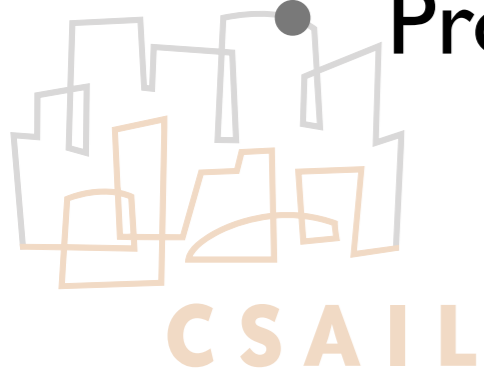
# Introduction

TD( $\lambda$ ): Dominant family of RL algorithms.

- Parameter  $\lambda$  used to mix unbiased, high-variance estimates with biased, low-variance estimates.
- $\lambda$  must be set manually.
- Up until now, we did not understand what it really does.

This talk:

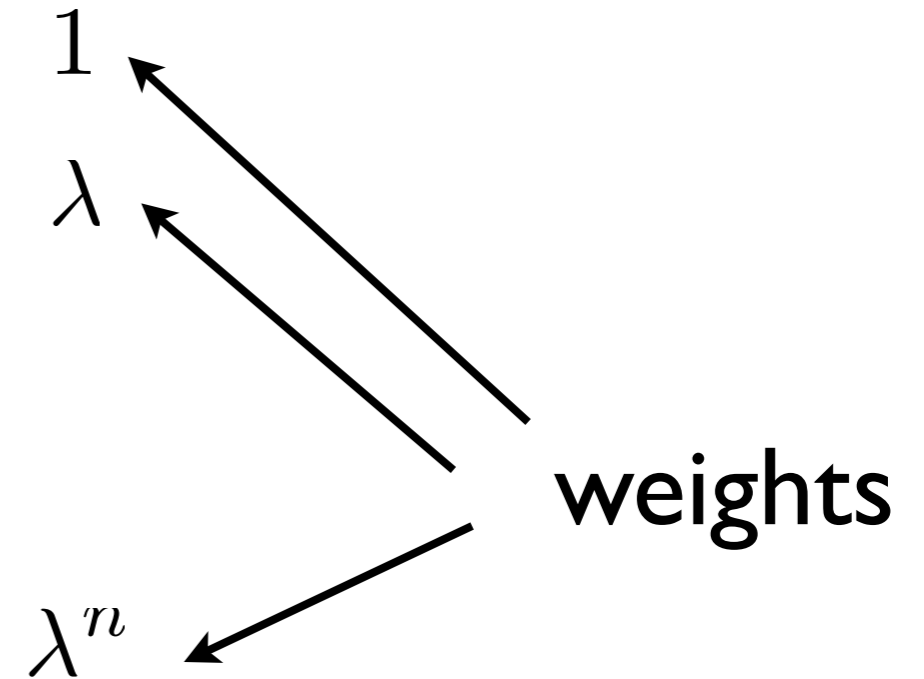
- Expose assumptions underlying TD( $\lambda$ ).
- Show that they are wrong!
- Propose another - parameter free - method, TD $_{\gamma}$ .



# TD( $\lambda$ )

Weighted sum:

$$\begin{aligned} R^{(1)} &= r_0 + \gamma V(s_1) \\ R^{(2)} &= r_0 + \gamma r_1 + \gamma^2 V(s_2) \\ &\cdot \\ &\cdot \\ &\cdot \\ R^{(n)} &= \sum_{i=0}^{n-1} \gamma^i r_i + \gamma^n V(s_n) \end{aligned}$$



Estimator:

$$R_{s_t}^\lambda = (1 - \lambda) \sum_{n=0}^{\infty} \lambda^n R_{s_t}^{(n+1)}$$



# TD( $\lambda$ )

This is called the  $\lambda$ -return.

- At  $\lambda=0$  we get TD, at  $\lambda=1$  we get MC.
- Intermediate values of  $\lambda$  usually best.

Results in a *family* of algorithms.

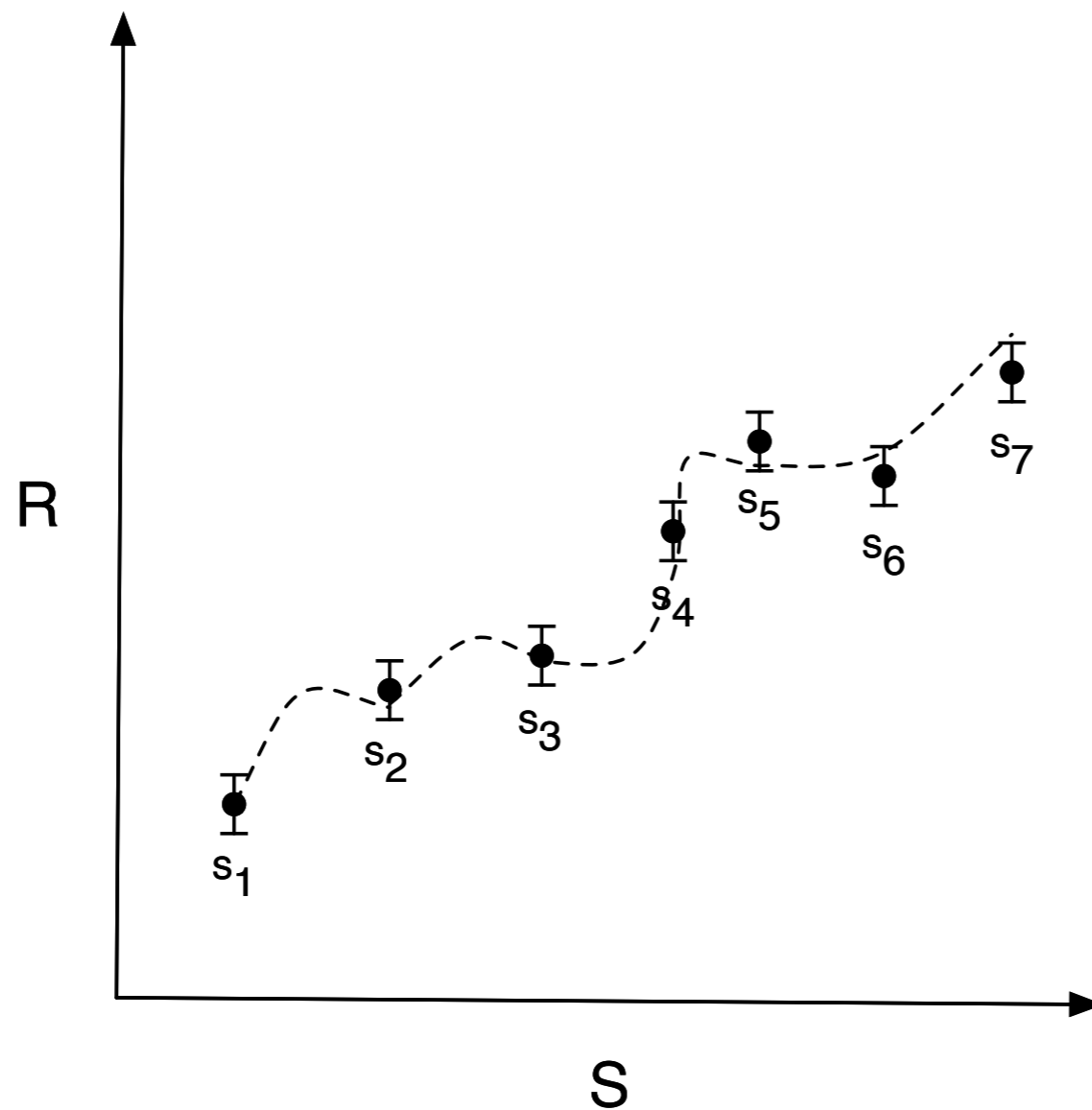
- Update rules via error metric using  $\lambda$ -return.
- Used almost exclusively, unchanged, since 1988.
- Original paper has ~3000 citations.

# TD( $\lambda$ )

- What are the implicit assumptions that lead to this estimator?



# Linear Least-Squares

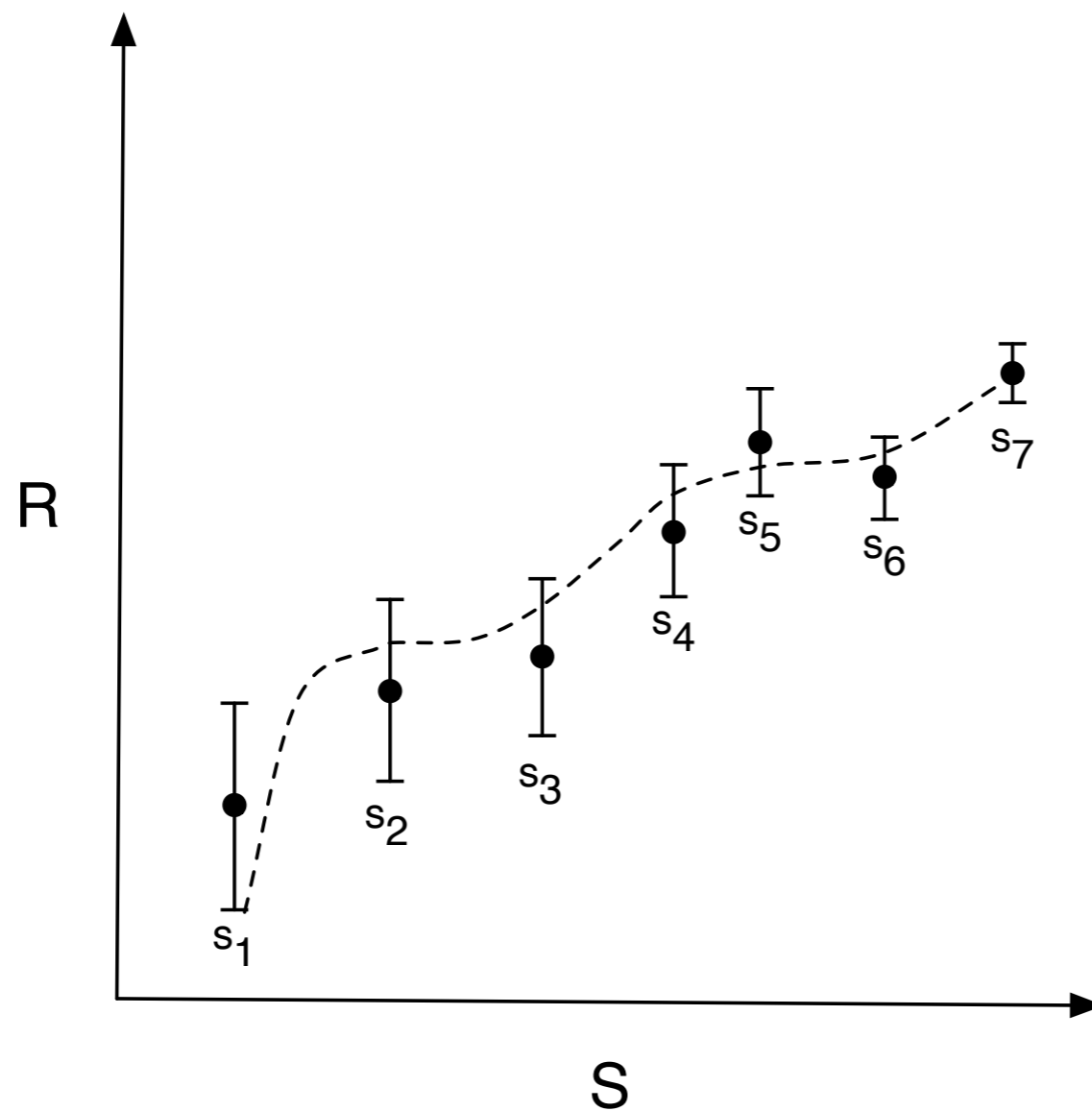


mean  
 $\mathbf{w} \cdot \Phi_i(s)$

variance  
 $\beta^{-1}$

$$e_i = \sum_{j=1}^m [\mathbf{w} \cdot \Phi_i(\mathbf{s}_j) - R_j]^2$$

# Weighted Linear Least-Squares



mean  
 $\mathbf{w} \cdot \Phi_i(s)$

variance  
 $(\rho^{m-i} \beta)^{-1}$

$$e_i = \sum_{j=1}^m \rho^{(m-j)} [\mathbf{w} \cdot \Phi_i(\mathbf{s}_j) - R_j]^2$$

# TD( $\lambda$ )

Consider the following assumptions:

- Each  $n$ -step rollout is independent.
- Each  $n$ -step rollout is normally distributed with mean of the true return.
- Variance of  $n$ -step rollout is  $k(n)$ .



# TD( $\lambda$ )

Likelihood:

$$\mathcal{L}(\hat{R}_{s_t} | R_{s_t}^{(1)}, \dots, R_{s_t}^{(n)}; k) = \prod_{n=1}^L \mathcal{N}(R_{s_t}^{(n)} | \hat{R}_{s_t}, k(n))$$

Maximizing the log likelihood:

$$\hat{R}_{s_t} = \frac{\sum_{n=1}^L k(n)^{-1} R_{s_t}^{(n)}}{\sum_{n=1}^L k(n)^{-1}}$$

this is the  $\lambda$ -return, where:

$$k(n) = \lambda^{-n}$$
$$L \rightarrow \infty$$



# TD( $\lambda$ )

Therefore,  $\lambda$ -return is the estimator you get given three assumptions:

- Normal distribution of return estimates.
- Independence of rollouts.
- Variance of rollouts increases geometrically with common ratio  $1/\lambda$ .

**All three of these assumptions are false.**

# On Rollout Variance

Let's let the first two slide, and consider the variance of an  $n$ -step sample return:

$$\begin{aligned} \text{Var} \left[ R_{s_t}^{(n)} \right] &= \text{Var} \left[ R_{s_t}^{(n-1)} - \gamma^{n-1} V(s_{t+n-1}) + \gamma^{n-1} r_{t+n-1} + \gamma^n V(s_{t+n}) \right] \\ &= \text{Var} \left[ R_{s_t}^{(n-1)} \right] + \gamma^{2(n-1)} \text{Var} \left[ V(s_{t+n-1}) - (r_{t+n-1} + \gamma V(s_{t+n})) \right] \\ &\quad + 2 \text{Cov} \left[ R_{s_t}^{(n-1)}, -\gamma^{n-1} V(s_{t+n-1}) + \gamma^{n-1} r_{t+n-1} + \gamma^n V(s_{t+n}) \right]. \end{aligned}$$

First thing to notice:

- Variance increases from  $n-1$  to  $n$  **additively**.
- We assume the covariance away.

# New Variance Model

We obtain:

$$\text{Var} \left[ R_{s_t}^{(n)} \right] \approx \text{Var} \left[ R_{s_t}^{(n-1)} \right] + \gamma^{2(n-1)} \text{Var} \left[ \overset{\text{TD error at step } n}{V(s_{t+n-1}) - (r_{t+n-1} + \gamma V(s_{t+n}))} \right]$$

We assume the TD error variance is the same everywhere, set its value to  $K$ .

Simple model of the variance of an  $n$ -step sample of return:

$$k(n) = \sum_{i=1}^n \gamma^{2(i-1)} K.$$



# TD<sub>γ</sub>

Resulting estimator:

$$R_{st}^\gamma = \frac{\cancel{\kappa}^{-1} \sum_{n=1}^L (\sum_{i=1}^n \gamma^{2(i-1)})^{-1} R_{st}^{(n)}}{\cancel{\kappa}^{-1} \sum_{n=1}^L (\sum_{i=1}^n \gamma^{2(i-1)})^{-1}} = \sum_{n=1}^L w(n, L) R_{st}^{(n)}$$

where

$$w(n, L) = \frac{(\sum_{i=1}^n \gamma^{2(i-1)})^{-1}}{\sum_{n=1}^L (\sum_{i=1}^n \gamma^{2(i-1)})^{-1}}$$

Parameter-free!



# TD $_{\gamma}$

Weighted sum:

$$R^{(1)} = r_0 + \gamma V(s_1)$$

$$R^{(2)} = r_0 + \gamma r_1 + \gamma^2 V(s_2)$$

$$R^{(3)} = r_0 + \gamma r_1 + \gamma^2 r_2 + \gamma^3 V(s_3)$$

•

•

•

$$R^{(n)} = \sum_{i=0}^{n-1} \gamma^i r_i + \gamma^n V(s_n)$$

weights

$$\frac{1}{1}$$

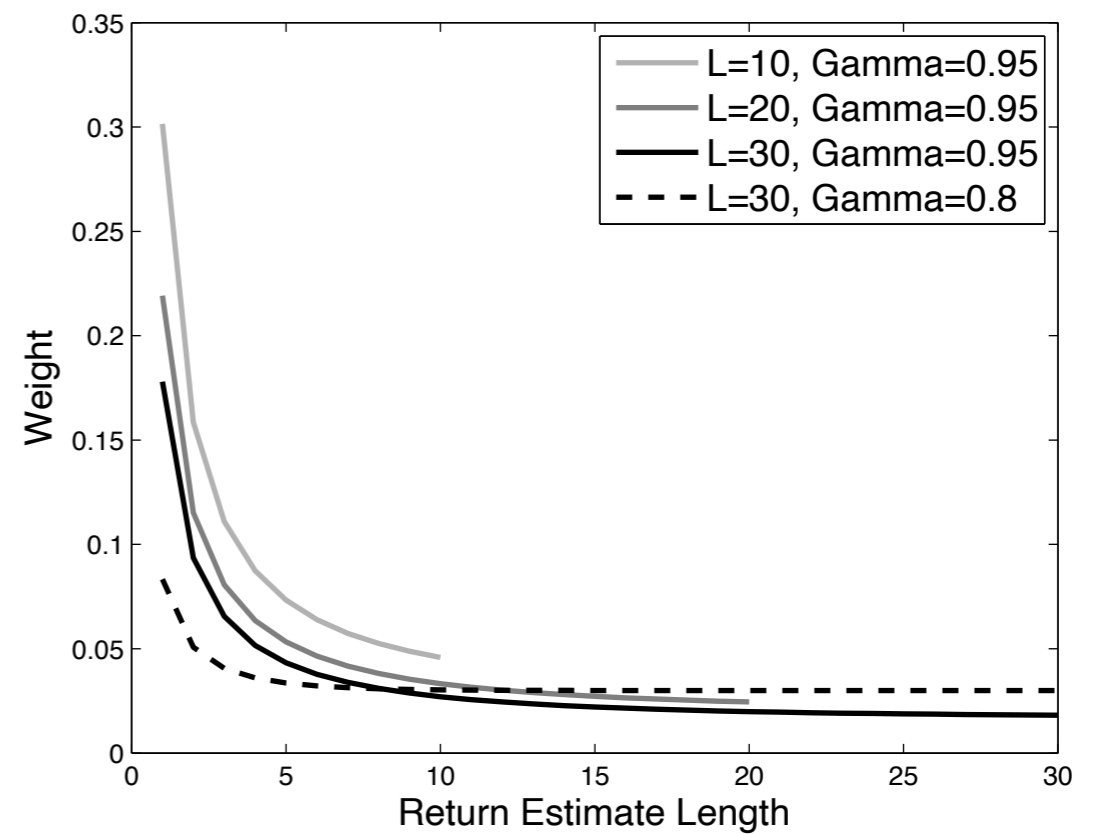
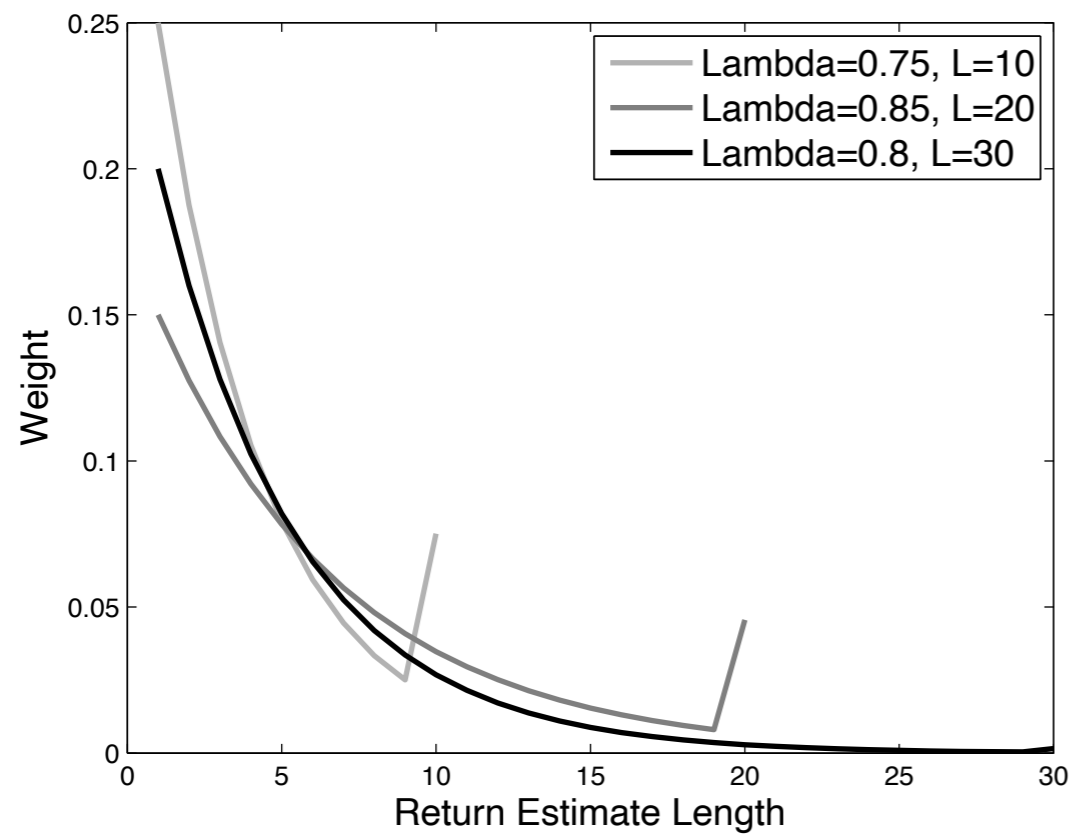
$$\frac{1}{1 + \gamma^2}$$

$$\frac{1}{1 + \gamma^2 + \gamma^4}$$

$$\frac{1}{1 + \gamma^2 + \gamma^4 + \dots + \gamma^{2n}}$$



# $TD_\gamma$ vs. $TD(\lambda)$



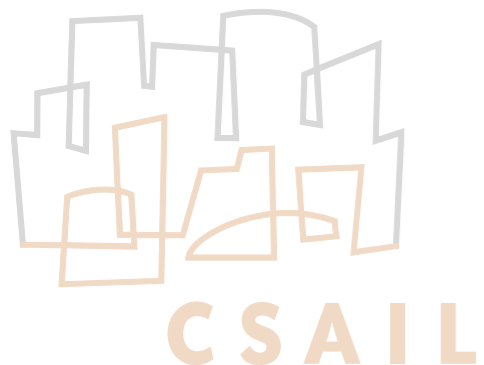
# The Catch

$$w(n, L) = \frac{(\sum_{i=1}^n \gamma^{2(i-1)})^{-1}}{\sum_{n=1}^L (\sum_{i=1}^n \gamma^{2(i-1)})^{-1}}$$

Normalizing constant is a function of episode length.  
Differs for each state.

$\lambda$ -return avoids this because it assumes episode is infinitely long, and sum of weights tends to a constant.

$$R_{s_t}^\lambda = (1 - \lambda) \sum_{n=0}^{\infty} \lambda^n R_{s_t}^{(n+1)}$$





# So:

TD $\gamma$  kills  $\lambda$  and replaces TD( $\lambda$ ) if:

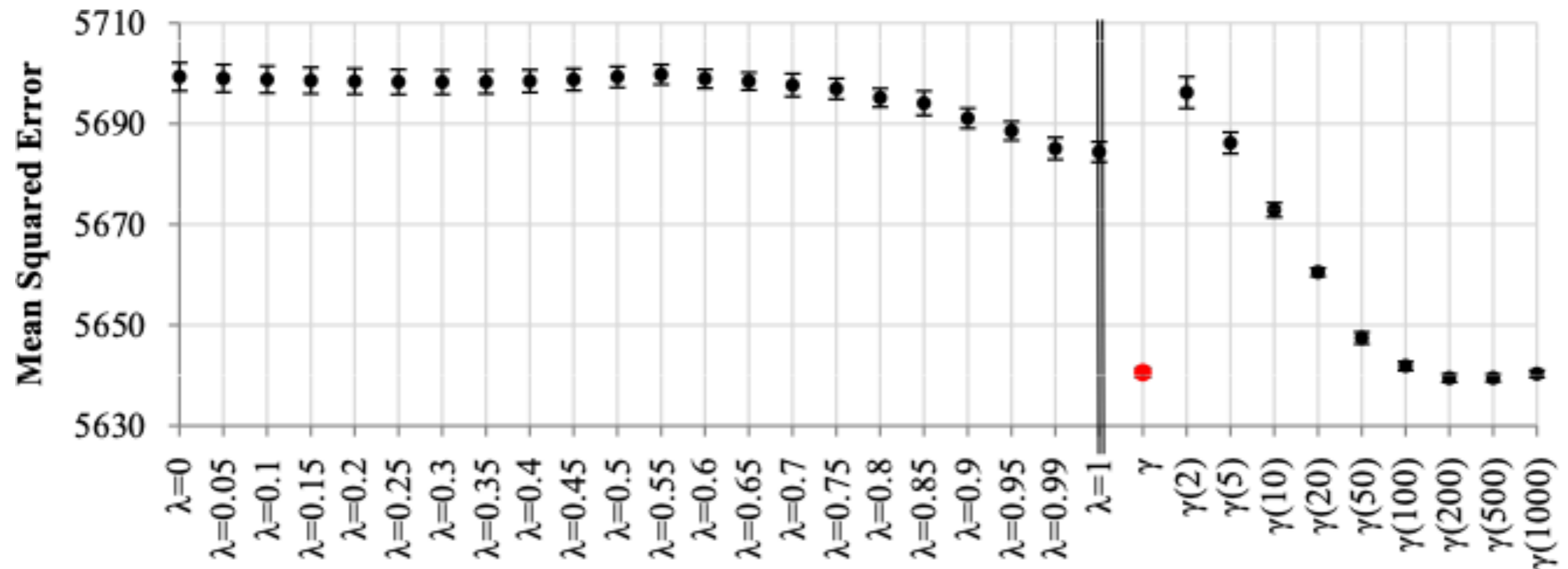
- We can be incremental episode-wise.
- We can process in a batch.

But not if we must be incremental transition-wise.

- We impose capacity  $C$ .
- Use the first  $C$  rollouts.
- Normalizer not known, except for last  $C-1$  steps.

# Results

## Acrobot



Similar for another 4 domains.

- TD $\gamma$  beats TD( $\lambda$ ) for any value of  $\lambda$  (4/5)
- Intermediate values of  $C$  do very well.

# Future Work

Better model of the variance.

Model that affords a completely incremental implementation.

Other members of the TD $\gamma$  family:

- LSTD $\gamma$
- Sarsa $\gamma$
- GQ $\gamma$

Account for covariance: TD(Omega)