

# COMPARING DIRECT AND INDIRECT TEMPORAL-DIFFERENCE METHODS FOR ESTIMATING THE VARIANCE OF THE RETURN

---

Craig Sherstan, Dylan R. Ashley\*, Brendan Bennett\*, Kenny Young,  
Adam White, Martha White, Richard S. Sutton

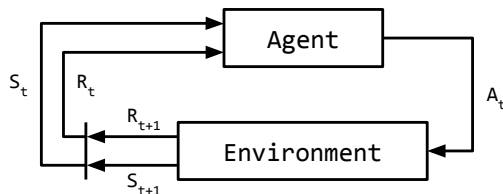
Reinforcement Learning and Artificial Intelligence Laboratory, University of Alberta

# BACKGROUND

---

# WHAT IS REINFORCEMENT LEARNING?

Reinforcement learning considers an agent interacting with an environment:



The function the agent uses to pick actions in states is known as the policy. Often the challenge is to find a "good" policy.

In reinforcement learning the return is defined as follows:

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

Often we want to maximize the **expected value** of the return. But "good" does depend on what we want.

## HOW CAN WE LEARN IT'S EXPECTED VALUE?

Temporal-difference (TD) methods have been fairly successful in tackling reinforcement learning problems so far. TD methods use predictions to update predictions.

One of the most straightforward TD methods is  $TD(\lambda)$ :

$$\delta_t = R_{t+1} + \gamma_{t+1} w_t^T x_{t+1} - w_t^T x_t$$

$$z_t = \gamma_t \lambda_t z_{t-1} + x_t$$

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

## SO WHAT'S THIS PRESENTATION ABOUT?

Recall what the return is:

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

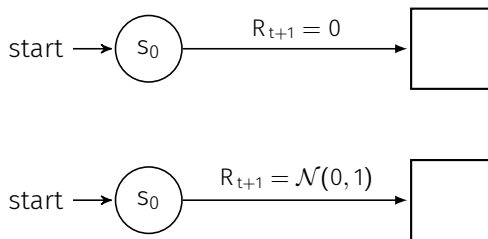
We're not limited to learning only its expected value. We could also learn more parts of its distribution such as its **variance**.

# MOTIVATION

---

## WHY MIGHT WE WANT TO LEARN ITS VARIANCE?

The variance might tell us things about the distribution that the expected value can't. Sometimes these things are **interesting**. For example it could differentiate these two domains:





## ANY BETTER REASONS?

The variance can give us **useful** information about the distribution. It can tell us how risky an action is to take in a state.

Humans take risk into decisions and don't necessarily act in a way that maximizes the expected value.

We can use an estimate of the variance to **learn** how to learn. For example here is an algorithm that uses an estimate of the variance to tune  $\lambda$  on the fly:

---

**Algorithm 2:**  $\lambda$ -greedy( $\mathbf{w}^{\text{err}}, \mathbf{w}^{\text{sq}}, \mathbf{w}_t, \mathbf{x}_t, \mathbf{x}_{t+1}, r_{t+1}, \rho_t$ )

---

```
// Use GTD to update  $\mathbf{w}^{\text{err}}$ 
 $\bar{g}_{t+1} \leftarrow \mathbf{x}_{t+1}^\top \mathbf{w}^{\text{err}}$ 
 $\delta_t \leftarrow r_{t+1} + \gamma_{t+1} \bar{g}_{t+1} - \mathbf{x}_t^\top \mathbf{w}^{\text{err}}$ 
 $\bar{\mathbf{e}}_t = \rho_t(\gamma_t \bar{\mathbf{e}}_{t-1} + \mathbf{x}_t)$ 
 $\mathbf{w}^{\text{err}} = \mathbf{w}^{\text{err}} + \alpha \delta_t \bar{\mathbf{e}}_t$ 
// Use VTD to update  $\mathbf{w}^{\text{sq}}$ 
 $\bar{r}_{t+1} \leftarrow \rho_t^2 r_{t+1}^2 + 2\rho_t^2 \gamma_{t+1} r_{t+1} \bar{g}_{t+1}$ 
 $\bar{\gamma}_{t+1} \leftarrow \rho_t^2 \gamma_{t+1}^2$ 
 $\bar{\delta}_t \leftarrow \bar{r}_{t+1} + \bar{\gamma}_{t+1} \mathbf{x}_{t+1}^\top \mathbf{w}^{\text{sq}} - \mathbf{x}_t^\top \mathbf{w}^{\text{sq}}$ 
 $\bar{\mathbf{z}}_t = \bar{\gamma}_t \bar{\mathbf{z}}_{t-1} + \mathbf{x}_t$ 
 $\mathbf{w}^{\text{sq}} = \mathbf{w}^{\text{sq}} + \alpha \bar{\delta}_t \bar{\mathbf{z}}_t$ 
// Compute  $\lambda$  estimate
 $\text{errsq} = (\bar{g}_{t+1} - \mathbf{x}_{t+1}^\top \mathbf{w}_t)^2$ 
 $\text{varg} = \max(0, \mathbf{x}_{t+1}^\top \mathbf{w}^{\text{sq}} - (\bar{g}_{t+1})^2)$ 
 $\lambda_{t+1} = \text{errsq} / (\text{varg} + \text{errsq})$ 
return  $\lambda_{t+1}$ 
```

---

# LEARNING THE VARIANCE

---

## HOW CAN WE LEARN ITS VARIANCE?

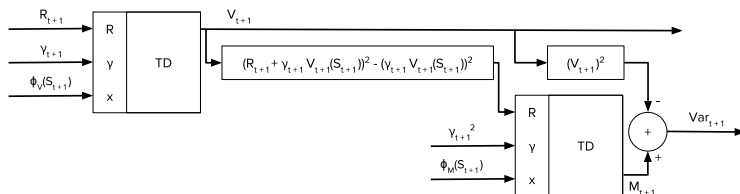
We can use this identity:

$$\text{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$$

If we are learning  $\mathbb{E}_\pi[G_t | S_t = s]$  then we can just learn  $\mathbb{E}_\pi[G_t^2 | S_t = s]$  on the side and use both our estimates to try to estimate  $\text{Var}_\pi(G_t | S_t = s)$ .

## WHAT WOULD THIS LOOK LIKE?

Using the identity  $\text{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$  one can estimate the variance using the following structure:



We can also use this identity:

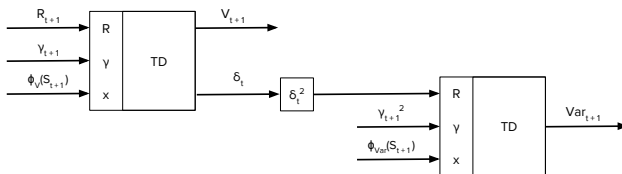
$$\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}[X])^2]$$

If we are learning  $\mathbb{E}_\pi[G_t | S_t = s]$  then we can approximate the variance using the following:

$$\text{Var}_\pi(G_t | S_t = s) \approx \mathbb{E}_\pi \left[ \delta_t^2 + \sum_{i=t+1}^{\infty} \left( \delta_i \prod_{j=t+1}^i \gamma_j \right)^2 \middle| S_t = s \right]$$

## SO HOW WOULD THIS LOOK?

Using the identity  $\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}[X])^2]$  one can estimate the variance using the following structure:



## WHAT WOULD BE THE UPDATE EQUATIONS FOR THIS?

Using TD( $\lambda$ ) with  $\bar{w}$  as the parameter vector for estimating the variance we obtain the following update equations:

$$\delta_t = R_{t+1} + \gamma_{t+1} w_t^T x_{t+1} - w_t^T x_t$$

$$z_t = \gamma_t \lambda_t z_{t-1} + x_t$$

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

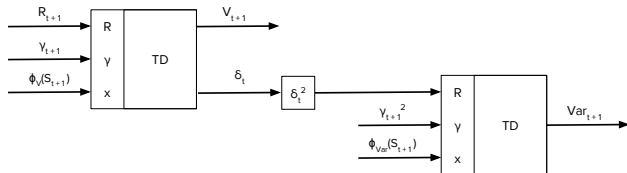
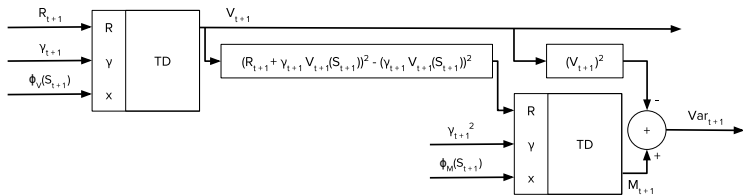
$$\bar{\delta}_t = \delta_t^2 + \gamma_{t+1}^2 \bar{w}_t^T x_{t+1} - \bar{w}_t^T x_t$$

$$\bar{z}_t = \gamma_t^2 \bar{\lambda}_t \bar{z}_{t-1} + x_t$$

$$\bar{w}_{t+1} = \bar{w}_t + \bar{\alpha}_{t+1} \bar{\delta}_t \bar{z}_t$$



# HOW DO THEY COMPARE VISUALLY?



# EMPIRICAL COMPARISON

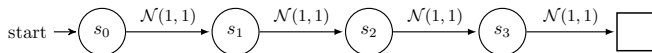
---

Ideally we want to know if the direct method

- is faster or slower to converge than the indirect method,
- is more robust or less robust to differences in the value and variance learner, and
- performs better or worse under linear function approximation.

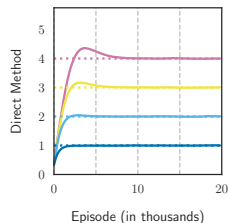
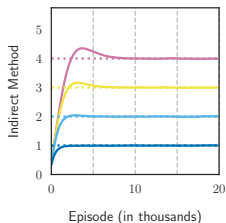
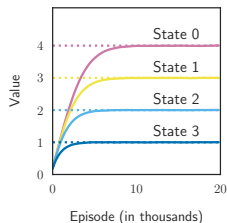
## WHAT IS THE SIMPLEST DOMAIN WE CAN COMPARE THEM ON?

We begin by comparing them on the following simple Markov chain with gaussian rewards:



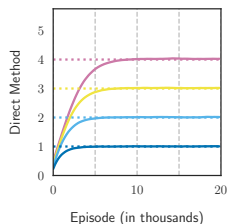
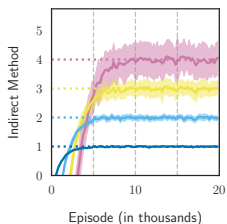
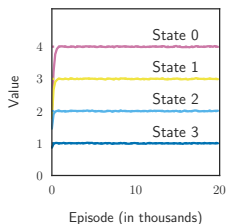
## HOW DO THEY COMPARE WHEN $\alpha = \bar{\alpha}$ ?

When  $\alpha = \bar{\alpha} = 0.001$  both perform roughly the same:



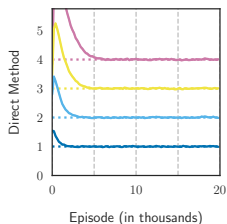
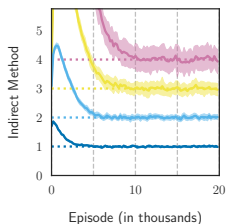
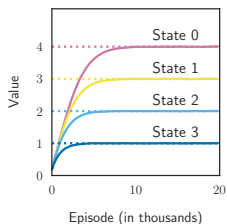
## WHAT ABOUT WHEN $\alpha > \bar{\alpha}$ ?

When  $\alpha = 0.01$  and  $\bar{\alpha} = 0.001$  the variance of the indirect method is higher:



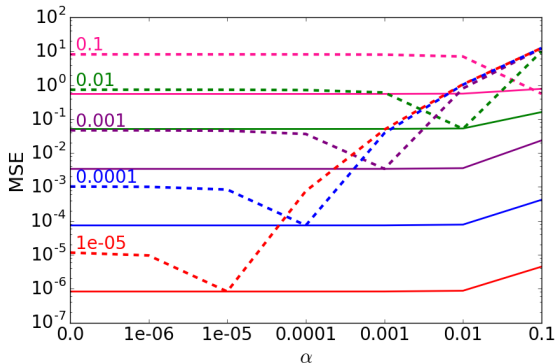
## SO WHAT ABOUT WHEN $\alpha < \bar{\alpha}$ ?

When  $\alpha = 0.001$  and  $\bar{\alpha} = 0.01$  the variance of the indirect method is higher and the direct method is more stable:



## DOES THIS RESULT GENERALIZE?

In this domain we only see the two perform similarly when both step sizes are equal (note that the dotted line represents the indirect method and the solid line represents the direct method):





## WHAT ABOUT UNDER FUNCTION APPROXIMATION?

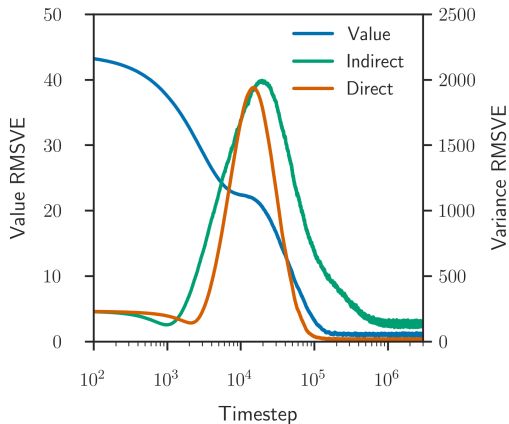
We use the following domain previously used to evaluate the indirect method:



For each state  $s_i$  we use  $\phi(s_i) = [1, i/30]^T$  for our value estimator and  $\phi_2(s_i) = [1, i/30, (i/30)^2]^T$  for our variance estimator.

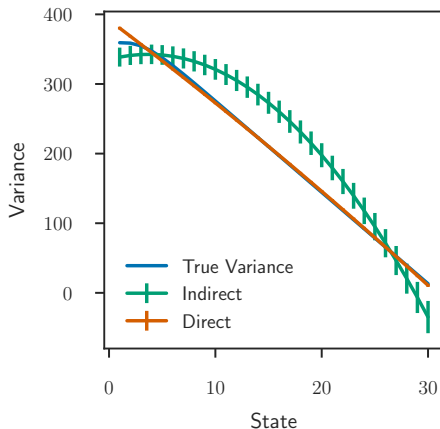
## HOW DO THEY PERFORM ON THIS DOMAIN?

Here the direct method vastly outperforms the indirect method:



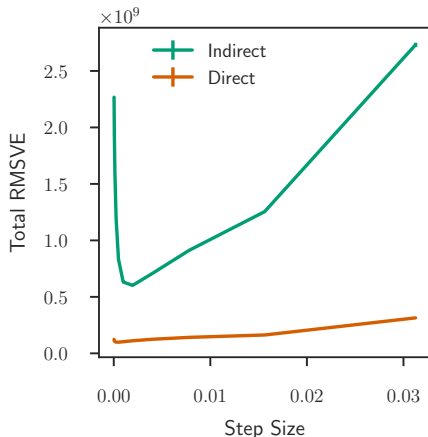
## WHAT IS THE QUALITY OF THE SOLUTIONS REACHED?

The direct method reaches a much better fixed point and exhibits much less variance in its variance estimates:



## WHAT IS THE PARAMETER SENSIVITY?

In this domain, the direct method is much less sensitive to the choice of step sizes than the indirect method:



We have described a method of directly estimating the variance of the return using temporal-difference methods. We have argued that learning the variance of the return can

- tell us **interesting** information about our domain,
- tell us **useful** information about the distribution of our return, and
- can be used to **learn** how to learn.

We have furthermore shown evidence that the direct method:

- learns just as fast and occasionally **faster** than the indirect method,
- is **more robust** to inconsistencies in the value and variance learner, and
- exhibits substantially **better** performance under linear function approximation.

QUESTIONS?