

---

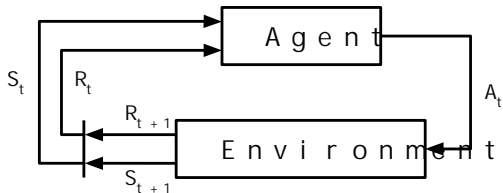
\*

\*

Reinforcement Learning and Artificial Intelligence La



Reinforcement learning considers an environment:



The function the agent uses to pick action policy. Often the challenge is to find a

In reinforcement learning the return is

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

Often we want to maximize the return.  
"good" does depend on what we want.

Temporal - difference (TD) methods have been tackling reinforcement learning problems by using predictions to update predictions.

One of the most straightforward TD methods

$$\delta_t = R_{t+1} + \gamma V_{t+1} - V_t$$

$$Z_t = \gamma \lambda Z_{t-1} + \mathbf{x}_t$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_{t+1} \delta_t Z_t$$

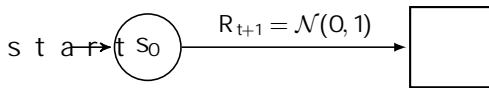
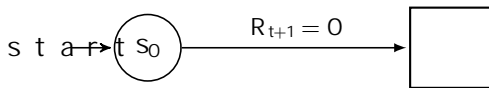
Recall what the return is:

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

We're not limited to learning only  $V$  it's learn more parts of its distribution as  $V$  is a function of  $s$  and  $a$ .



The variance might tell us things about  
expected value can't. Sometimes for these  
example it could differentiate these two





The variance  $\sigma^2$  can give information about the  
can tell us how risky an action is to take.

Humans take risk into decisions and do  
that that maximizes the expected value.

We can use an estimate of  $\lambda$  to improve variance in the example here is an algorithm that uses TD to do this:

---

**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  $\lambda_{t+1}$ 

```

---

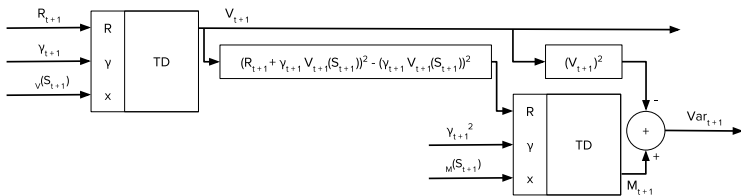


We can use this identity:

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

If we are  $\mathbb{E}_\pi[G_t | S_t = s]$  then we can just learn  $\mathbb{E}_\pi[G_t^2 | S_t = s]$  on the side and use both our estimates  $\hat{G}_t | S_t = s$ .

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



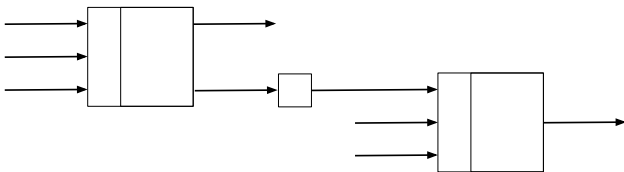
We can also use this identity:

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

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

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

Using the identity  $\sigma^2(X) = E[(X - E[X])^2]$  one can estimate variance using the following structure



Using TD(0) with the parameter vector for the variance we obtain the following update

$$\delta_t = R_{t+1} + \gamma_{t+1} \mathbf{w}_t^T \mathbf{x}_{t+1} - \mathbf{w}_t^T \mathbf{x}_t$$

$$\mathbf{z}_t = \gamma_t \lambda_t \mathbf{z}_{t-1} + \mathbf{x}_t$$

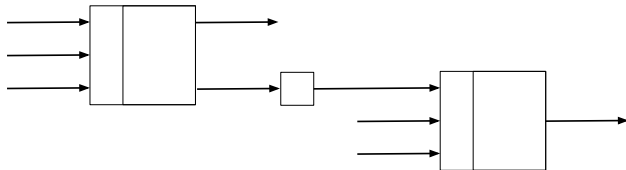
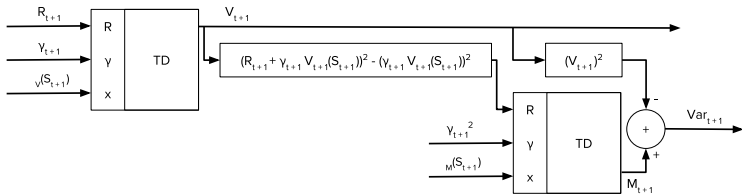
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_{t+1} \delta_t \mathbf{z}_t$$

$$\bar{\delta}_t = \delta_t^2 + \gamma_{t+1}^2 \bar{\mathbf{w}}_t^T \mathbf{x}_{t+1} - \bar{\mathbf{w}}_t^T \mathbf{x}_t$$

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

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



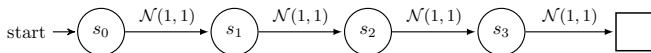




I d e a l l y w e w a n t t o k n o w i f t h e d i r e c t m e

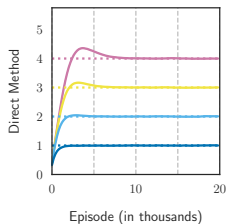
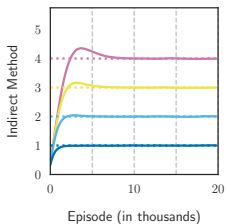
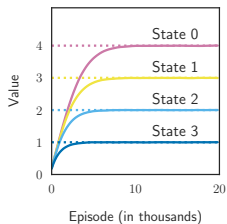
- i s f a s t e r o r s l o w e r t o c o n v e r g e t h a n
- i s m o r e r o b u s t o r l e s s r o b u s t t o d i f f
- v a r i a n c e l e a r n e r , a n d
- p e r f o r m s b e t t e r o r w o r s e u n d e r l i n e a

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



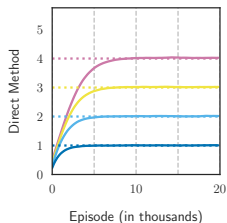
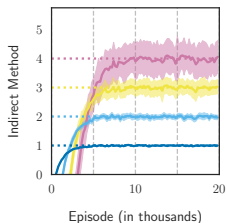
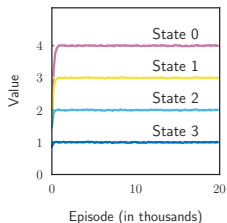
$$\alpha = \bar{\alpha}$$

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



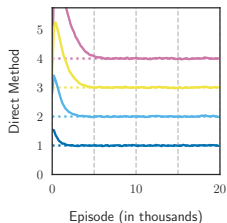
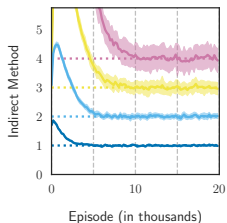
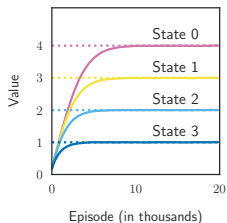
$$\alpha > \bar{\alpha}$$

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

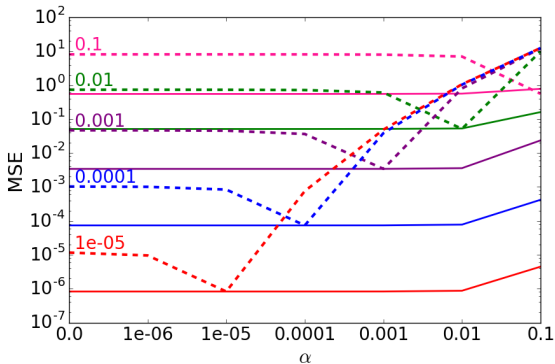


$$\alpha < \bar{\alpha}$$

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



In this domain we only see the two perfect sizes are equal (note that the dotted line method and the solid line represents t



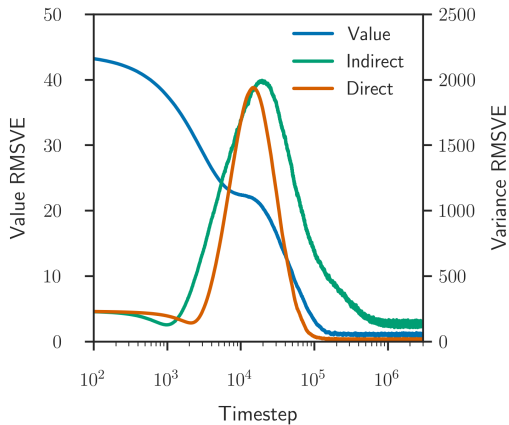


We use the following domain previously indirect method:

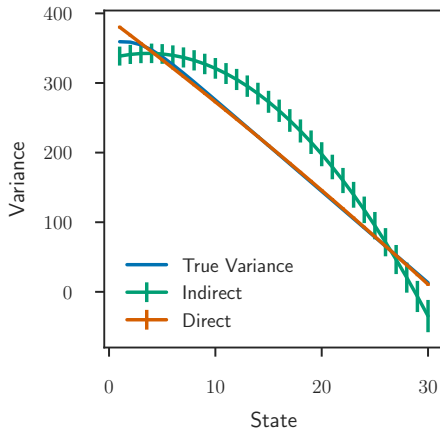


For each  $i$ , we use  $\mu(s_i) = [1, i/30]^T$  for our value estimate and  $\sigma_2(s_i) = [1, i/30, (i/30)^2]^T$  for our variance estimate.

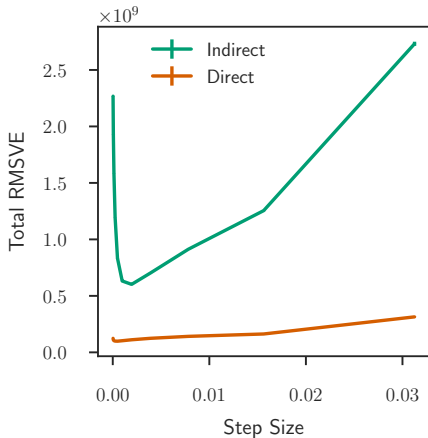
Here the direct method vastly outperforms



The direct method reaches a much better  
much less variance in its variance est



In this domain, the direct method is much more accurate than the indirect method for a wide range of step sizes.



We have described a method of directly  
the return using temporal - difference  
learning the variance of the return ca

- tell **l** user **i**s **n** **f** **i** **n** **g** information about our doma
- tell **l** user **i**s **n** **f** **i** **n** **g** information about the distri  
and
- can be used **e** **a** **t** **h** **o** **w** to learn.

We have furthermore shown evidence that

- learns just as fast as the standard method,
- is more robust to inconsistencies in the data learner, and
- exhibits superior performance under linear approximation.

