

1-step

2

3

4

5

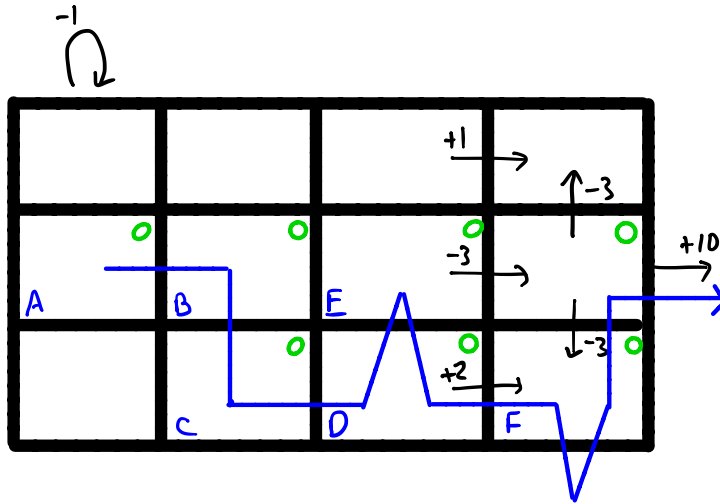
6

7

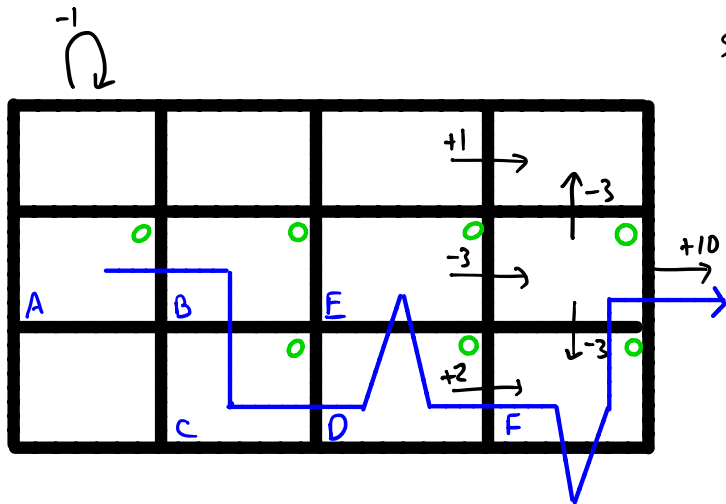
8

9

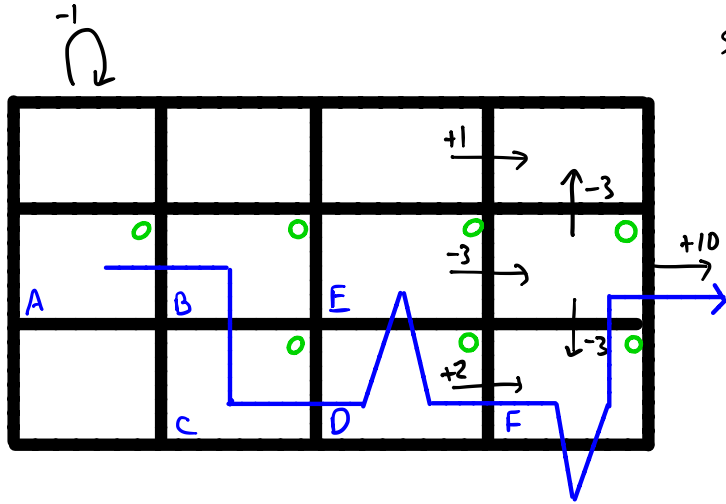
# Review: n-step returns



- A
- 1-step ?
  - 2-step ?
  - 3 " ?
  - 4 :
  - 5 :
  - 6
  - 7
  - 8
  - 9

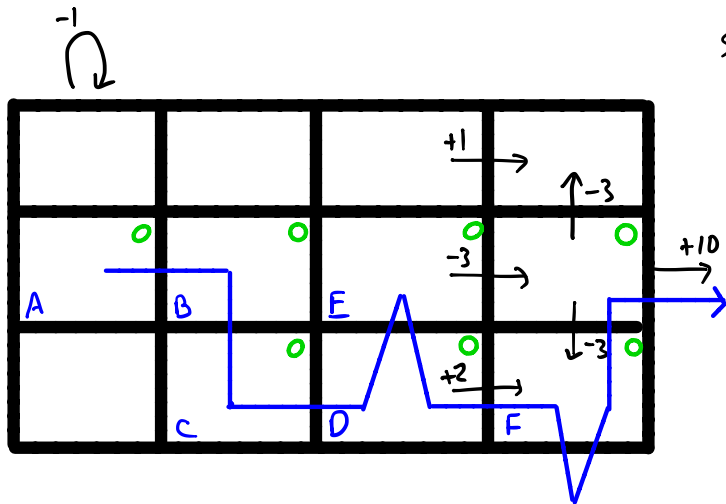


- |        |                 |
|--------|-----------------|
| 1-step | A               |
| 2      | 0               |
| 3      | 0               |
| 4      | 0               |
| 5      | 0               |
| 6      | 28 <sup>5</sup> |
| 7      | ?               |
| 8      | ?               |
| 9      | ?               |

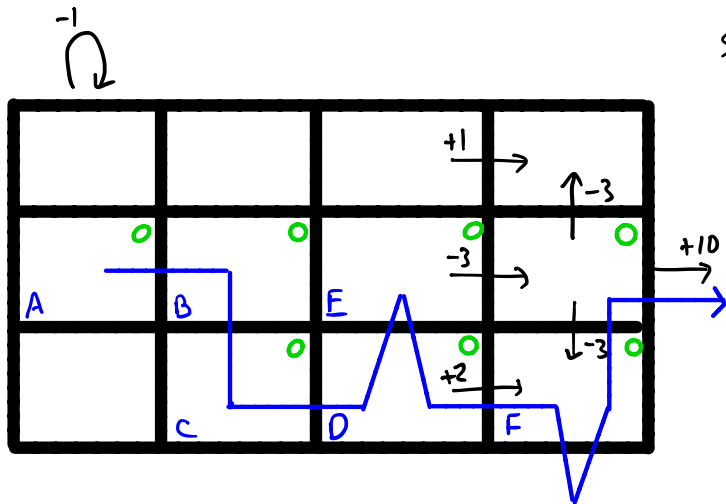


- 1-step A 0
- 2 0
- 3 0
- 4 0
- 5 0
- 6  $2x^5$
- 7  $2x^5 - x^6$
- 8  $2x^5 - x^6$
- 9  $2x^5 - x^6 + 10x^8$

B ?



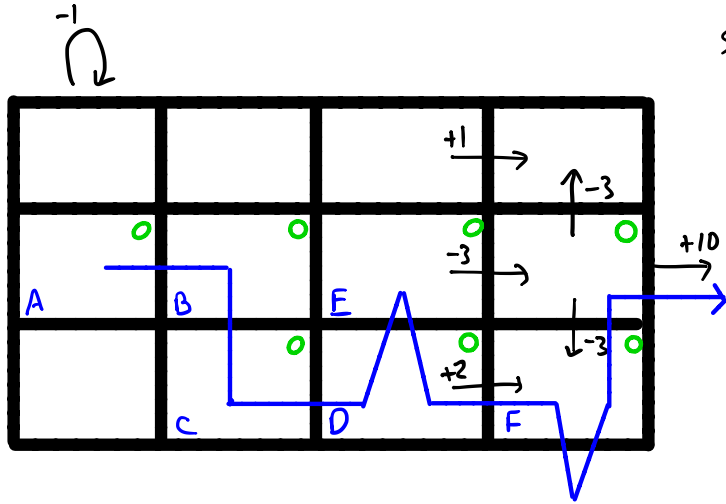
|        | A                    | B                    |
|--------|----------------------|----------------------|
| 1-step | 0                    | 0                    |
| 2      | 0                    | 0                    |
| 3      | 0                    | 0                    |
| 4      | 0                    | 0                    |
| 5      | 0                    | $2x^4$               |
| 6      | $2x^5$               | $2x^4 - x^5$         |
| 7      | $2x^5 - x^6$         | $2x^4 - x^5$         |
| 8      | $2x^5 - x^6$         | $2x^4 - x^5 + 10x^7$ |
| 9      | $2x^5 - x^6 + 10x^8$ |                      |



Stand  
 ↑  
 Clap  
 ↓  
 Wave

offline  $\lambda$ -return

|        | A   | B                                   |
|--------|---|-------------------------------------|
| 1-step | 0 $(1-\lambda)$                                 | 0                                   |
| 2      | 0 $(1-\lambda)\lambda$                          | 0                                   |
| 3      | 0 $(1-\lambda)\lambda^2$                        | 0                                   |
| 4      | 0 $(1-\lambda)\lambda^3$                        | 0                                   |
| 5      | 0 $(1-\lambda)\lambda^4$                        | $2\delta^4$                         |
| 6      | $2\delta^5$ $(1-\lambda)\lambda^5$              | $2\delta^4 - \delta^5$              |
| 7      | $2\delta^5 - \delta^6$ $(1-\lambda)\lambda^6$   | $2\delta^4 - \delta^5$              |
| 8      | $2\delta^5 - \delta^6$ $(1-\lambda)\lambda^7$   | $2\delta^4 - \delta^5 + 10\delta^7$ |
| 9      | $2\delta^5 - \delta^6 + 10\delta^8$ $\lambda^8$ |                                     |

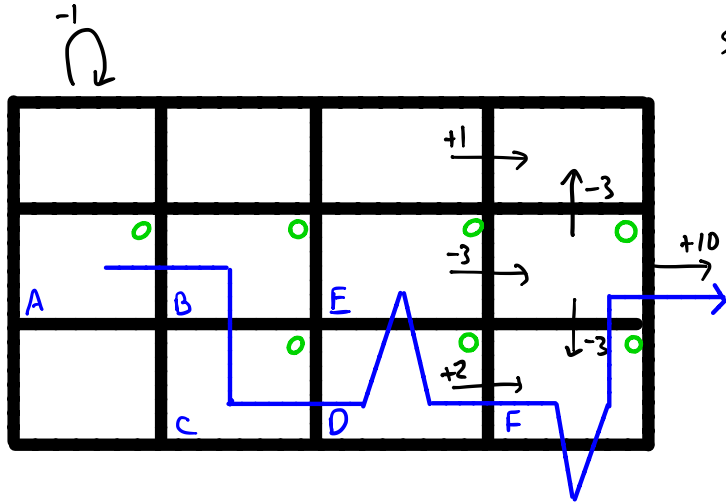


offline  $\lambda$ -return

|        | A  | B  |
|--------|--|--|
| 1-step | 0 (1- $\lambda$ )                                  | 0 (1- $\lambda$ )                                  |
| 2      | 0 (1- $\lambda$ ) $\lambda$                        | 0 (1- $\lambda$ ) $\lambda$                        |
| 3      | 0 (1- $\lambda$ ) $\lambda^2$                      | 0 (1- $\lambda$ ) $\lambda^2$                      |
| 4      | 0 (1- $\lambda$ ) $\lambda^3$                      | 0 (1- $\lambda$ ) $\lambda^3$                      |
| 5      | 0 (1- $\lambda$ ) $\lambda^4$                      | $2\delta^4$ (1- $\lambda$ ) $\lambda^4$            |
| 6      | $2\delta^5$ (1- $\lambda$ ) $\lambda^5$            | $2\delta^4 - \delta^5$ (1- $\lambda$ ) $\lambda^5$ |
| 7      | $2\delta^5 - \delta^6$ (1- $\lambda$ ) $\lambda^6$ | $2\delta^4 - \delta^5$ (1- $\lambda$ ) $\lambda^6$ |
| 8      | $2\delta^5 - \delta^6$ (1- $\lambda$ ) $\lambda^7$ | $2\delta^4 - \delta^5 + 10\delta^7$ $\lambda^7$    |
| 9      | $2\delta^5 - \delta^6 + 10\delta^8$ $\lambda^8$    | $\lambda^8$  |

What makes this "offline?"

Why use this bizarre weighting scheme?



offline  $\lambda$ -return

|        | A   | B   |
|--------|---|---|
| 1-step | 0 $(1-\lambda)$                               | 0 $(1-\lambda)$                               |
| 2      | 0 $(1-\lambda)\lambda$                        | 0 $(1-\lambda)\lambda$                        |
| 3      | 0 $(1-\lambda)\lambda^2$                      | 0 $(1-\lambda)\lambda^2$                      |
| 4      | 0 $(1-\lambda)\lambda^3$                      | 0 $(1-\lambda)\lambda^3$                      |
| 5      | 0 $(1-\lambda)\lambda^4$                      | 0 $(1-\lambda)\lambda^4$                      |
| 6      | $2\delta^5 (1-\lambda)\lambda^5$              | $2\delta^4 - \delta^5 (1-\lambda)\lambda^5$   |
| 7      | $2\delta^5 - \delta^6 (1-\lambda)\lambda^6$   | $2\delta^4 - \delta^5 (1-\lambda)\lambda^6$   |
| 8      | $2\delta^5 - \delta^6 (1-\lambda)\lambda^7$   | $2\delta^4 - \delta^5 + 10\delta^7 \lambda^7$ |
| 9      | $2\delta^5 - \delta^6 + 10\delta^8 \lambda^8$ | $\lambda^8$                                   |

What makes this "offline?"

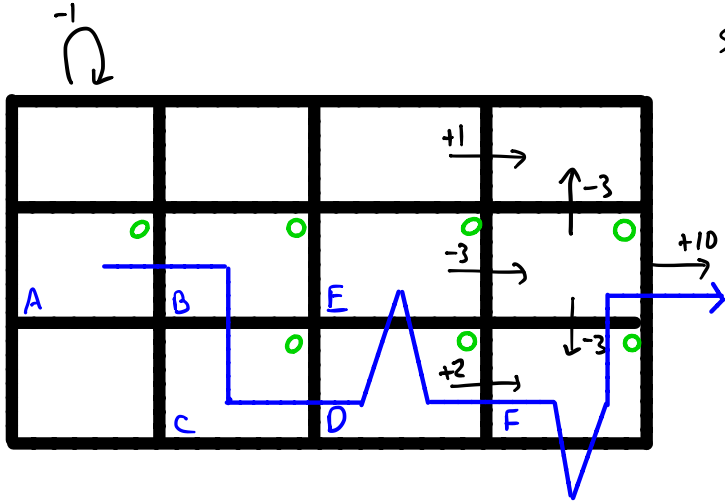
Why use this bizarre weighting scheme?



TD( $\lambda$ ) is online — why is that preferable?

$TD(2)$  is online — why is that preferable?

1. Updates on every step: no memory
2. Updates equally distributed in time
3. Can be applied to continuing problems: even  $TD(1)=MC$

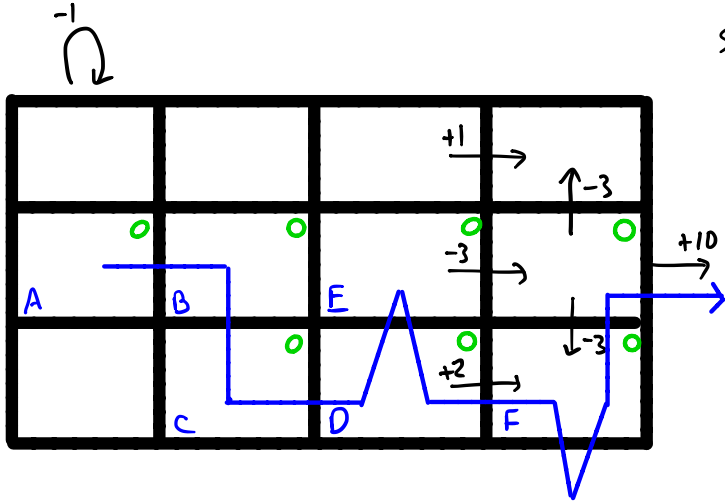


|   | z |   | δ |   |   |  |
|---|---|---|---|---|---|--|
| A | B | C | D | E | F |  |
|   | 0 | 0 | 0 | 0 | 0 |  |
| ? |   |   |   |   |   |  |
| ? |   |   |   |   |   |  |
| ⋮ |   |   |   |   |   |  |

|        | A  | B  |
|--------|--|--|
| 1-step | 0 (1-λ)  | 0 (1-λ)  |
| 2      | 0 (1-λ)λ   | 0 (1-λ)λ   |
| 3      | 0 (1-λ)λ <sup>2</sup>  | 0 (1-λ)λ <sup>2</sup>  |
| 4      | 0 (1-λ)λ <sup>3</sup>  | 0 (1-λ)λ <sup>3</sup>  |
| 5      | 0 (1-λ)λ <sup>4</sup>  | 2λ <sup>4</sup> (1-λ)λ <sup>4</sup>                              |
| 6      | 2λ <sup>5</sup> (1-λ)λ <sup>5</sup>                              | 2λ <sup>4</sup> -λ <sup>5</sup> (1-λ)λ <sup>5</sup>              |
| 7      | 2λ <sup>5</sup> -λ <sup>6</sup> (1-λ)λ <sup>6</sup>              | 2λ <sup>4</sup> -λ <sup>5</sup> (1-λ)λ <sup>6</sup>              |
| 8      | 2λ <sup>5</sup> -λ <sup>6</sup> (1-λ)λ <sup>7</sup>              | 2λ <sup>4</sup> -λ <sup>5</sup> +10λ <sup>7</sup> λ <sup>7</sup> |
| 9      | 2λ <sup>5</sup> -λ <sup>6</sup> +10λ <sup>8</sup> λ <sup>8</sup> |  |

$$\delta = R + \gamma v(s') - v(s)$$

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

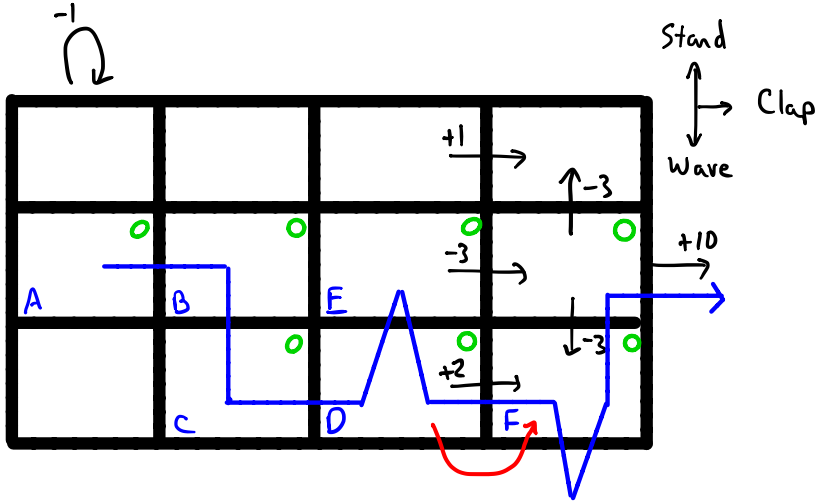


|                | z |   | s   |   |   |   |
|----------------|---|---|-----|---|---|---|
|                | A | B | C   | D | E | F |
| 1              | 0 | 0 | 0   | 0 | 0 | 0 |
| $\delta z$     |   |   |     |   |   |   |
| $(\delta z)^2$ |   |   |     |   |   |   |
| $\vdots$       |   |   |     |   |   |   |
| $(\delta z)^4$ | ? | ? | ... |   |   |   |

|        | A                                   |                        | B                                   |                        |
|--------|-------------------------------------|------------------------|-------------------------------------|------------------------|
| 1-step | 0                                   | $(1-\lambda)$          | 0                                   | $(1-\lambda)$          |
| 2      | 0                                   | $(1-\lambda)\lambda$   | 0                                   | $(1-\lambda)\lambda$   |
| 3      | 0                                   | $(1-\lambda)\lambda^2$ | 0                                   | $(1-\lambda)\lambda^2$ |
| 4      | 0                                   | $(1-\lambda)\lambda^3$ | 0                                   | $(1-\lambda)\lambda^3$ |
| 5      | 0                                   | $(1-\lambda)\lambda^4$ | 0                                   | $(1-\lambda)\lambda^4$ |
| 6      | $2\delta^5$                         | $(1-\lambda)\lambda^5$ | $2\delta^4$                         | $(1-\lambda)\lambda^5$ |
| 7      | $2\delta^5 - \delta^6$              | $(1-\lambda)\lambda^6$ | $2\delta^4 - \delta^5$              | $(1-\lambda)\lambda^6$ |
| 8      | $2\delta^5 - \delta^6$              | $(1-\lambda)\lambda^7$ | $2\delta^4 - \delta^5 + 10\delta^7$ | $\lambda^7$            |
| 9      | $2\delta^5 - \delta^6 + 10\delta^8$ | $\lambda^8$            |                                     |                        |

$$\delta = R + \gamma v(s') - v(s)$$

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



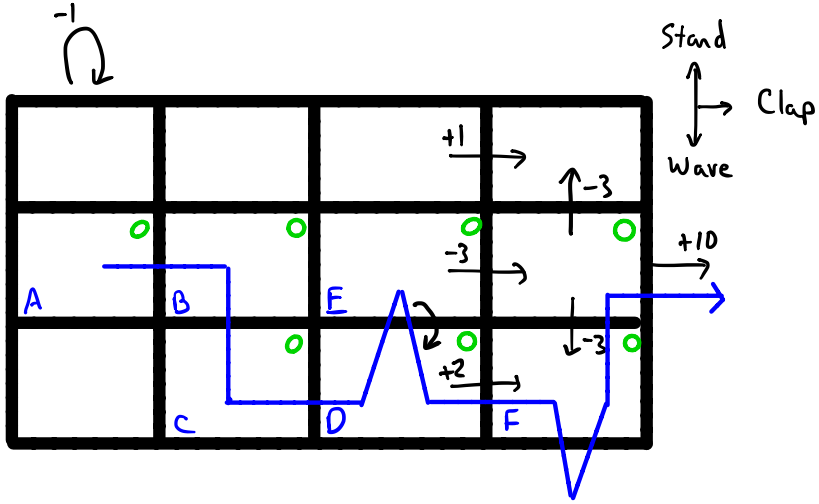
$w = v(A)$        $TD(s)$        $v(F)$

|                | $z$ |   |   |   |   |   |          |
|----------------|-----|---|---|---|---|---|----------|
|                | A   | B | C | D | E | F | $\delta$ |
| 1              | 0   | 0 | 0 | 0 | 0 | 0 | ?        |
| $\delta z$     |     |   |   |   |   |   | ...      |
| $(\delta z)^2$ |     |   |   |   |   |   | ?        |
| $\vdots$       |     |   |   |   |   |   |          |
| $(\delta z)^4$ |     |   |   |   |   |   |          |
| ?              |     |   |   |   |   |   |          |
| .              |     |   |   |   |   |   |          |

|        | A   | B   |
|--------|---|---|
| 1-step | 0 $(1-\lambda)$                                 | 0 $(1-\lambda)$                                 |
| 2      | 0 $(1-\lambda)\lambda$                          | 0 $(1-\lambda)\lambda$                          |
| 3      | 0 $(1-\lambda)\lambda^2$                        | 0 $(1-\lambda)\lambda^2$                        |
| 4      | 0 $(1-\lambda)\lambda^3$                        | 0 $(1-\lambda)\lambda^3$                        |
| 5      | 0 $(1-\lambda)\lambda^4$                        | 0 $(1-\lambda)\lambda^4$                        |
| 6      | $2\delta^5$ $(1-\lambda)\lambda^5$              | $2\delta^4$ $(1-\lambda)\lambda^4$              |
| 7      | $2\delta^5 - \delta^6$ $(1-\lambda)\lambda^6$   | $2\delta^4 - \delta^5$ $(1-\lambda)\lambda^5$   |
| 8      | $2\delta^5 - \delta^6$ $(1-\lambda)\lambda^7$   | $2\delta^4 - \delta^5 + 10\delta^7$ $\lambda^7$ |
| 9      | $2\delta^5 - \delta^6 + 10\delta^8$ $\lambda^8$ |   |

$$\delta = R + \gamma v(s') - v(s)$$

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



$w = v(A)$   
?

$T0(z)$

$v(F)$

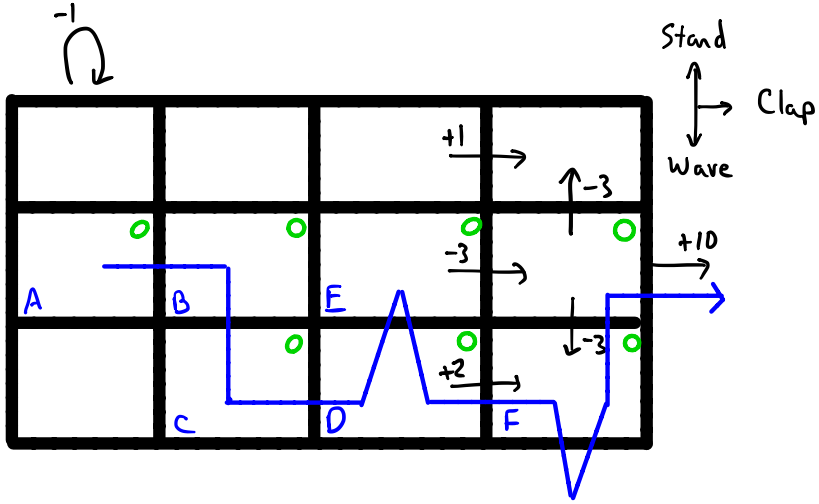
|                | $z$ |                |                |                    |            |   |          |
|----------------|-----|----------------|----------------|--------------------|------------|---|----------|
|                | A   | B              | C              | D                  | E          | F | $\delta$ |
| 1              | 0   | 0              | 0              | 0                  | 0          | 0 | 0        |
| $\delta z$     |     |                |                |                    |            |   | 0        |
| $(\delta z)^2$ |     |                |                |                    |            |   | 0        |
| $\vdots$       |     |                |                |                    |            |   | $\vdots$ |
| $(\delta z)^4$ |     | $(\delta z)^3$ | $(\delta z)^2$ | $(\delta z)$       | 1          | 0 | 0        |
| $(\delta z)^5$ |     | $(\delta z)^4$ | $(\delta z)^3$ | $(\delta z)^2 + 1$ | $\delta z$ | 0 | 2        |

↑ accumulating trace

|        | A   | B   |
|--------|---|---|
| 1-step | 0 $(1-\lambda)$                                 | 0 $(1-\lambda)$                                 |
| 2      | 0 $(1-\lambda)\lambda$                          | 0 $(1-\lambda)\lambda$                          |
| 3      | 0 $(1-\lambda)\lambda^2$                        | 0 $(1-\lambda)\lambda^2$                        |
| 4      | 0 $(1-\lambda)\lambda^3$                        | 0 $(1-\lambda)\lambda^3$                        |
| 5      | 0 $(1-\lambda)\lambda^4$                        | 0 $(1-\lambda)\lambda^4$                        |
| 6      | $2\delta^5$ $(1-\lambda)\lambda^5$              | $2\delta^4$ $(1-\lambda)\lambda^4$              |
| 7      | $2\delta^5 - \delta^6$ $(1-\lambda)\lambda^6$   | $2\delta^4 - \delta^5$ $(1-\lambda)\lambda^5$   |
| 8      | $2\delta^5 - \delta^6$ $(1-\lambda)\lambda^7$   | $2\delta^4 - \delta^5 + 10\delta^7$ $\lambda^7$ |
| 9      | $2\delta^5 - \delta^6 + 10\delta^8$ $\lambda^8$ |   |

$$\delta = R + \gamma v(s') - v(s)$$

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



$w = v(A)$   
 $2\alpha(\delta z)^5$

$v(B)$   
 $2\alpha(\delta z)^4$

$TD(z)$

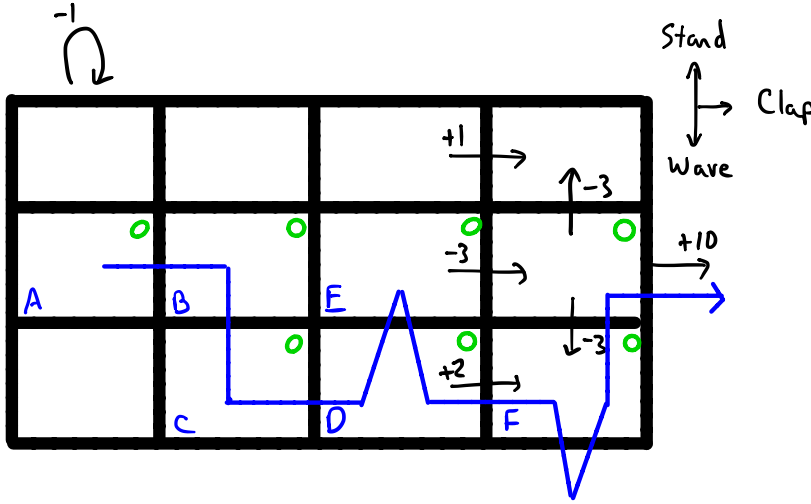
$v(F)$   
 $0$

|                |                |                |                    |            |   |          |
|----------------|----------------|----------------|--------------------|------------|---|----------|
|                | $z$            |                |                    |            |   |          |
| A              | B              | C              | D                  | E          | F | $\delta$ |
| 1              | 0              | 0              | 0                  | 0          | 0 | 0        |
| $\delta z$     |                |                |                    |            |   | $\vdots$ |
| $(\delta z)^2$ |                |                |                    |            |   | $\vdots$ |
| $(\delta z)^5$ | $(\delta z)^4$ | $(\delta z)^3$ | $(\delta z)^2 + 1$ | $\delta z$ | 0 | 2        |
| ?              | ...            |                |                    |            |   |          |

|        |   |   |
|--------|---|---|
|        | A   | B   |
| 1-step | 0 $(1-\lambda)$                                 | 0 $(1-\lambda)$                                 |
| 2      | 0 $(1-\lambda)\lambda$                          | 0 $(1-\lambda)\lambda$                          |
| 3      | 0 $(1-\lambda)\lambda^2$                        | 0 $(1-\lambda)\lambda^2$                        |
| 4      | 0 $(1-\lambda)\lambda^3$                        | 0 $(1-\lambda)\lambda^3$                        |
| 5      | 0 $(1-\lambda)\lambda^4$                        | 0 $(1-\lambda)\lambda^4$                        |
| 6      | $2\delta^5$ $(1-\lambda)\lambda^5$              | $2\delta^4$ $(1-\lambda)\lambda^4$              |
| 7      | $2\delta^5 - \delta^6$ $(1-\lambda)\lambda^6$   | $2\delta^4 - \delta^5$ $(1-\lambda)\lambda^5$   |
| 8      | $2\delta^5 - \delta^6$ $(1-\lambda)\lambda^7$   | $2\delta^4 - \delta^5 + 10\delta^7$ $\lambda^7$ |
| 9      | $2\delta^5 - \delta^6 + 10\delta^8$ $\lambda^8$ |   |

$$\delta = R + \gamma v(s') - v(s)$$

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



$w = v(A)$   
 $2\alpha(\delta\lambda)^5$   
 ?

$v(B)$   
 $2\alpha(\delta\lambda)^4$

$TD(\lambda)$

$v(F)$   
 0  
 ?

|                     | $z$ |   |   |   |   |   |          |
|---------------------|-----|---|---|---|---|---|----------|
|                     | A   | B | C | D | E | F | $\delta$ |
| 1                   | 0   | 0 | 0 | 0 | 0 | 0 | 0        |
| $\delta\lambda$     |     |   |   |   |   |   | $\vdots$ |
| $(\delta\lambda)^2$ |     |   |   |   |   |   | $\vdots$ |
| $\vdots$            |     |   |   |   |   |   | $\vdots$ |
| $(\delta\lambda)^5$ |     |   |   |   |   |   | 2        |
| $(\delta\lambda)^6$ |     |   |   |   |   | 1 | -1       |

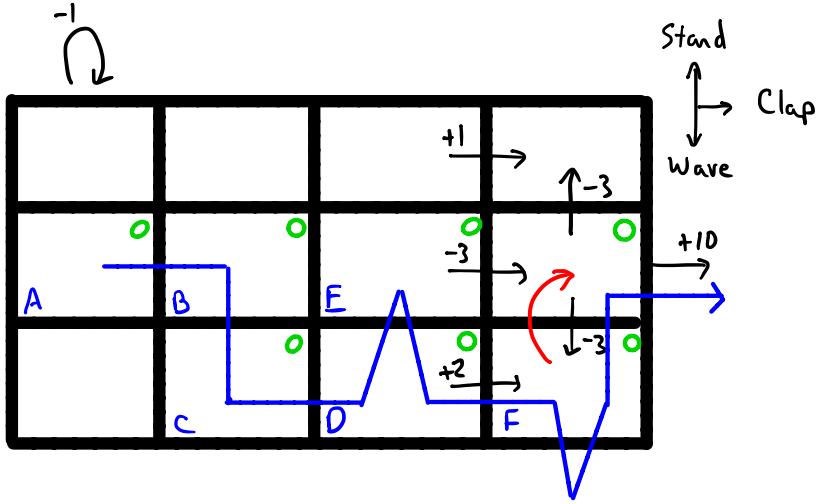
<snip>

|        | A   | B   |
|--------|---|---|
| 1-step | 0 $(1-\lambda)$                                 | 0 $(1-\lambda)$                                 |
| 2      | 0 $(1-\lambda)\lambda$                          | 0 $(1-\lambda)\lambda$                          |
| 3      | 0 $(1-\lambda)\lambda^2$                        | 0 $(1-\lambda)\lambda^2$                        |
| 4      | 0 $(1-\lambda)\lambda^3$                        | 0 $(1-\lambda)\lambda^3$                        |
| 5      | 0 $(1-\lambda)\lambda^4$                        | 0 $(1-\lambda)\lambda^4$                        |
| 6      | $2\delta^5$ $(1-\lambda)\lambda^5$              | $2\delta^4$ $(1-\lambda)\lambda^4$              |
| 7      | $2\delta^5 - \delta^6$ $(1-\lambda)\lambda^6$   | $2\delta^4 - \delta^5$ $(1-\lambda)\lambda^5$   |
| 8      | $2\delta^5 - \delta^6$ $(1-\lambda)\lambda^7$   | $2\delta^4 - \delta^5 + 10\delta^7$ $\lambda^7$ |
| 9      | $2\delta^5 - \delta^6 + 10\delta^8$ $\lambda^8$ |   |

$$\delta = R + \gamma v(s') - v(s)$$

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





$TD(z)$

$w = v(A)$   
 $2\alpha(\delta z)^5$   
 $2\alpha(\delta z)^5 - \alpha(\delta z)^6$   
 ?

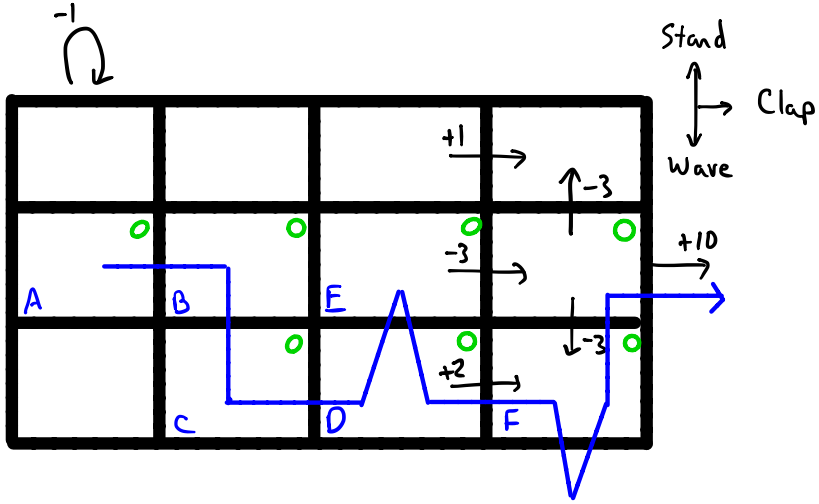
$v(F)$   
 0  
 $-\alpha$

|                |                |   |   |                |   |  |          |
|----------------|----------------|---|---|----------------|---|--|----------|
|                | $z$            |   |   |                |   |  |          |
| A              | B              | C | D | E              | F |  | $\delta$ |
| 1              | 0              | 0 | 0 | 0              | 0 |  | 0        |
| $\delta z$     |                |   |   |                |   |  | $\vdots$ |
| $(\delta z)^2$ |                |   |   |                |   |  | $\vdots$ |
| $\vdots$       |                |   |   |                |   |  | 2        |
| $(\delta z)^5$ | $(\delta z)^4$ |   |   | $\delta z$     | 0 |  | -1       |
| $(\delta z)^6$ | $(\delta z)^5$ |   |   | $(\delta z)^2$ | 1 |  | ?        |

$\langle \text{snip} \rangle$

|        |   |   |
|--------|---|---|
|        | A   | B   |
| 1-step | 0 $(1-\lambda)$                                 | 0 $(1-\lambda)$                                 |
| 2      | 0 $(1-\lambda)\lambda$                          | 0 $(1-\lambda)\lambda$                          |
| 3      | 0 $(1-\lambda)\lambda^2$                        | 0 $(1-\lambda)\lambda^2$                        |
| 4      | 0 $(1-\lambda)\lambda^3$                        | 0 $(1-\lambda)\lambda^3$                        |
| 5      | 0 $(1-\lambda)\lambda^4$                        | 0 $(1-\lambda)\lambda^4$                        |
| 6      | $2\delta^5$ $(1-\lambda)\lambda^5$              | $2\delta^4$ $(1-\lambda)\lambda^4$              |
| 7      | $2\delta^5 - \delta^6$ $(1-\lambda)\lambda^6$   | $2\delta^4 - \delta^5$ $(1-\lambda)\lambda^5$   |
| 8      | $2\delta^5 - \delta^6$ $(1-\lambda)\lambda^7$   | $2\delta^4 - \delta^5 + 10\delta^7$ $\lambda^7$ |
| 9      | $2\delta^5 - \delta^6 + 10\delta^8$ $\lambda^8$ |   |

$\delta = R + \gamma v(s') - v(s)$   
 $w = w + \alpha \delta z$



TD(z)

$w = v(A)$

$2\alpha(\delta z)^5$

$2\alpha(\delta z)^5 - \alpha(\delta z)^6$

$2\alpha(\delta z)^5 - \alpha(\delta z)^6 + \alpha^2(\delta z)^7$

$v(F)$

0

- $\alpha$

|                | z              |                |                           |                |   | \delta   |
|----------------|----------------|----------------|---------------------------|----------------|---|----------|
| A              | B              | C              | D                         | E              | F |          |
| 1              | 0              | 0              | 0                         | 0              | 0 | 0        |
| $\delta z$     |                |                |                           |                |   | ...      |
| $(\delta z)^2$ |                |                |                           |                |   | 2        |
| ...            |                |                |                           |                |   | -1       |
| $(\delta z)^5$ | $(\delta z)^4$ | $(\delta z)^3$ | $(\delta z)^2 + 1$        | $\delta z$     | 0 | $\alpha$ |
| $(\delta z)^6$ | $(\delta z)^5$ | $(\delta z)^4$ | $(\delta z)^3 + \delta z$ | $(\delta z)^2$ | 1 |          |
| $(\delta z)^7$ |                |                |                           |                |   |          |

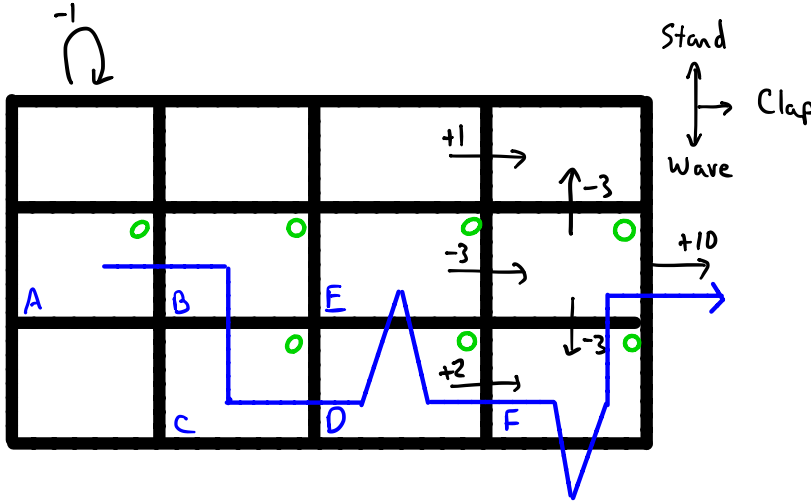
<snip>

different from off line case

|        | A                                   |                        | B                                   |                        |
|--------|-------------------------------------|------------------------|-------------------------------------|------------------------|
| 1-step | 0                                   | $(1-\lambda)$          | 0                                   | $(1-\lambda)$          |
| 2      | 0                                   | $(1-\lambda)\lambda$   | 0                                   | $(1-\lambda)\lambda$   |
| 3      | 0                                   | $(1-\lambda)\lambda^2$ | 0                                   | $(1-\lambda)\lambda^2$ |
| 4      | 0                                   | $(1-\lambda)\lambda^3$ | 0                                   | $(1-\lambda)\lambda^3$ |
| 5      | 0                                   | $(1-\lambda)\lambda^4$ | 0                                   | $(1-\lambda)\lambda^4$ |
| 6      | 0                                   | $(1-\lambda)\lambda^5$ | $2\delta^4$                         | $(1-\lambda)\lambda^5$ |
| 7      | $2\delta^5 - \delta^6$              | $(1-\lambda)\lambda^6$ | $2\delta^4 - \delta^5$              | $(1-\lambda)\lambda^6$ |
| 8      | $2\delta^5 - \delta^6$              | $(1-\lambda)\lambda^7$ | $2\delta^4 - \delta^5 + 10\delta^7$ | $\lambda^7$            |
| 9      | $2\delta^5 - \delta^6 + 10\delta^8$ | $\lambda^8$            |                                     |                        |

$\delta = R + \gamma v(s') - v(s)$

$w = w + \alpha \delta z$



$W = v(A)$   
 $2\alpha(\delta\lambda)^5$   
 $2\alpha(\delta\lambda)^5 - \alpha(\delta\lambda)^6$   
 $2\alpha(\delta\lambda)^5 - \alpha(\delta\lambda)^6 + \alpha^2(\delta\lambda)^7$

$v(F)$   
 $0$   
 $-\alpha$

$TD(0) \leftrightarrow TD(1)$   
 $TD \dots \lambda \dots MC$

Builds up incrementally to almost same result as off line forward view

|                     | $z$                 |                     |                     |                                     |                     |   |          |
|---------------------|---------------------|---------------------|---------------------|-------------------------------------|---------------------|---|----------|
|                     | A                   | B                   | C                   | D                                   | E                   | F | $\delta$ |
| A                   | 1                   | 0                   | 0                   | 0                                   | 0                   | 0 | 0        |
| $\delta\lambda$     | $(\delta\lambda)^2$ |                     |                     |                                     |                     |   | $\vdots$ |
| $(\delta\lambda)^3$ |                     | $(\delta\lambda)^4$ |                     |                                     |                     |   | 2        |
| $(\delta\lambda)^4$ |                     | $(\delta\lambda)^5$ |                     |                                     |                     |   | -1       |
| $(\delta\lambda)^5$ |                     |                     | $(\delta\lambda)^3$ | $(\delta\lambda)^2 + 1$             | $\delta\lambda$     | 0 | $\alpha$ |
| $(\delta\lambda)^6$ |                     |                     | $(\delta\lambda)^4$ | $(\delta\lambda)^3 + \delta\lambda$ | $(\delta\lambda)^2$ | 1 |          |
| $(\delta\lambda)^7$ |                     |                     |                     |                                     |                     |   |          |

<snip>

|        | A                                   |                        | B                                   |                        |
|--------|-------------------------------------|------------------------|-------------------------------------|------------------------|
| 1-step | 0                                   | $(1-\lambda)$          | 0                                   | $(1-\lambda)$          |
| 2      | 0                                   | $(1-\lambda)\lambda$   | 0                                   | $(1-\lambda)\lambda$   |
| 3      | 0                                   | $(1-\lambda)\lambda^2$ | 0                                   | $(1-\lambda)\lambda^2$ |
| 4      | 0                                   | $(1-\lambda)\lambda^3$ | 0                                   | $(1-\lambda)\lambda^3$ |
| 5      | 0                                   | $(1-\lambda)\lambda^4$ | 0                                   | $(1-\lambda)\lambda^4$ |
| 6      | 0                                   | $(1-\lambda)\lambda^5$ | $2\delta^4$                         | $(1-\lambda)\lambda^5$ |
| 7      | $2\delta^5 - \delta^6$              | $(1-\lambda)\lambda^6$ | $2\delta^4 - \delta^5$              | $(1-\lambda)\lambda^6$ |
| 8      | $2\delta^5 - \delta^6$              | $(1-\lambda)\lambda^7$ | $2\delta^4 - \delta^5 + 10\delta^7$ | $\lambda^7$            |
| 9      | $2\delta^5 - \delta^6 + 10\delta^8$ | $\lambda^8$            |                                     |                        |

different from off line case

$$\delta = R + \gamma v(s') - v(s)$$

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

offline  $\lambda$ -return

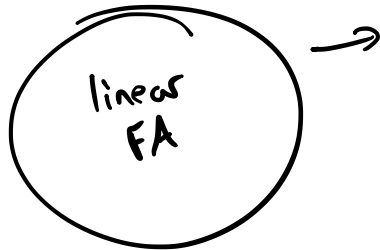
$TD(\lambda)$

offline  $\lambda$ -return

TD( $\lambda$ )

(Truncated  $\lambda$ -return  
Online  $\lambda$ -return)

True Online TD( $\lambda$ )



offline  $\lambda$ -return

$TD(\lambda)$

Truncated  $\lambda$ -return

Online  $\lambda$ -return

True Online  $TD(\lambda)$

## Online $\lambda$ -return

$$h = 1 : \quad \mathbf{w}_1^1 \doteq \mathbf{w}_0^1 + \alpha [G_{0:1}^\lambda - \hat{v}(S_0, \mathbf{w}_0^1)] \nabla \hat{v}(S_0, \mathbf{w}_0^1), \quad \leftarrow ?$$

$$h = 2 : \quad \mathbf{w}_1^2 \doteq \mathbf{w}_0^2 + \alpha [G_{0:2}^\lambda - \hat{v}(S_0, \mathbf{w}_0^2)] \nabla \hat{v}(S_0, \mathbf{w}_0^2), \\ \mathbf{w}_2^2 \doteq \mathbf{w}_1^2 + \alpha [G_{1:2}^\lambda - \hat{v}(S_1, \mathbf{w}_1^2)] \nabla \hat{v}(S_1, \mathbf{w}_1^2),$$

$$h = 3 : \quad \mathbf{w}_1^3 \doteq \mathbf{w}_0^3 + \alpha [G_{0:3}^\lambda - \hat{v}(S_0, \mathbf{w}_0^3)] \nabla \hat{v}(S_0, \mathbf{w}_0^3), \\ \mathbf{w}_2^3 \doteq \mathbf{w}_1^3 + \alpha [G_{1:3}^\lambda - \hat{v}(S_1, \mathbf{w}_1^3)] \nabla \hat{v}(S_1, \mathbf{w}_1^3), \\ \mathbf{w}_3^3 \doteq \mathbf{w}_2^3 + \alpha [G_{2:3}^\lambda - \hat{v}(S_2, \mathbf{w}_2^3)] \nabla \hat{v}(S_2, \mathbf{w}_2^3).$$

How differs from TD(0)?

## Online $\lambda$ -return

$$h = 1 : \quad \mathbf{w}_1^1 \doteq \mathbf{w}_0^1 + \alpha [G_{0:1}^\lambda - \hat{v}(S_0, \mathbf{w}_0^1)] \nabla \hat{v}(S_0, \mathbf{w}_0^1), \quad \leftarrow \text{No change}$$

$$h = 2 : \quad \mathbf{w}_1^2 \doteq \mathbf{w}_0^2 + \alpha [G_{0:2}^\lambda - \hat{v}(S_0, \mathbf{w}_0^2)] \nabla \hat{v}(S_0, \mathbf{w}_0^2), \\ \mathbf{w}_2^2 \doteq \mathbf{w}_1^2 + \alpha [G_{1:2}^\lambda - \hat{v}(S_1, \mathbf{w}_1^2)] \nabla \hat{v}(S_1, \mathbf{w}_1^2),$$

$$h = 3 : \quad \mathbf{w}_1^3 \doteq \mathbf{w}_0^3 + \alpha [G_{0:3}^\lambda - \hat{v}(S_0, \mathbf{w}_0^3)] \nabla \hat{v}(S_0, \mathbf{w}_0^3), \\ \mathbf{w}_2^3 \doteq \mathbf{w}_1^3 + \alpha [G_{1:3}^\lambda - \hat{v}(S_1, \mathbf{w}_1^3)] \nabla \hat{v}(S_1, \mathbf{w}_1^3), \\ \mathbf{w}_3^3 \doteq \mathbf{w}_2^3 + \alpha [G_{2:3}^\lambda - \hat{v}(S_2, \mathbf{w}_2^3)] \nabla \hat{v}(S_2, \mathbf{w}_2^3).$$

How differs from TD(0)?



## Online $\lambda$ -return

$$h = 1: \mathbf{w}_1^1 \doteq \mathbf{w}_0^1 + \alpha [G_{0:1}^\lambda - \hat{v}(S_0, \mathbf{w}_0^1)] \nabla \hat{v}(S_0, \mathbf{w}_0^1),$$

← No change

$$h = 2: \cancel{\mathbf{w}_1^2 \doteq \mathbf{w}_0^2 + \alpha [G_{0:2}^\lambda - \hat{v}(S_0, \mathbf{w}_0^2)] \nabla \hat{v}(S_0, \mathbf{w}_0^2)},$$
$$\mathbf{w}_2^2 \doteq \mathbf{w}_1^1 + \alpha [G_{0:2}^\lambda - \hat{v}(S_1, \mathbf{w}_1^1)] \nabla \hat{v}(S_1, \mathbf{w}_1^1),$$

$G_{0:2}^\lambda$  with weights  $(1-\lambda), (1-\lambda)\lambda$   
(online truncates:  $(1-\lambda), \lambda$ )

$$h = 3: \mathbf{w}_1^3 \doteq \mathbf{w}_0^3 + \alpha [G_{0:3}^\lambda - \hat{v}(S_0, \mathbf{w}_0^3)] \nabla \hat{v}(S_0, \mathbf{w}_0^3),$$
$$\mathbf{w}_2^3 \doteq \mathbf{w}_1^3 + \alpha [G_{1:3}^\lambda - \hat{v}(S_1, \mathbf{w}_1^3)] \nabla \hat{v}(S_1, \mathbf{w}_1^3),$$
$$\mathbf{w}_3^3 \doteq \mathbf{w}_2^3 + \alpha [G_{2:3}^\lambda - \hat{v}(S_2, \mathbf{w}_2^3)] \nabla \hat{v}(S_2, \mathbf{w}_2^3).$$

How differs from TD(2)?

## Online $\lambda$ -return

$$h = 1: \quad \mathbf{w}_1^1 \doteq \mathbf{w}_0^1 + \alpha [G_{0:1}^\lambda - \hat{v}(S_0, \mathbf{w}_0^1)] \nabla \hat{v}(S_0, \mathbf{w}_0^1),$$

← No change

$$h = 2: \quad \cancel{\mathbf{w}_1^2 \doteq \mathbf{w}_0^2 + \alpha [G_{0:2}^\lambda - \hat{v}(S_0, \mathbf{w}_0^2)] \nabla \hat{v}(S_0, \mathbf{w}_0^2)},$$

$$\mathbf{w}_2^2 \doteq \cancel{\mathbf{w}_1^1} + \alpha [G_{\cancel{0:2}}^\lambda - \hat{v}(S_1, \cancel{\mathbf{w}_1^1})] \nabla \hat{v}(S_1, \cancel{\mathbf{w}_1^1}),$$

$G_{0:2}^\lambda$  with weights  $(1-\lambda), (1-\lambda)\lambda$   
 (online truncates:  $(1-\lambda), \lambda$ )

$$h = 3: \quad \cancel{\mathbf{w}_1^3 \doteq \mathbf{w}_0^3 + \alpha [G_{0:3}^\lambda - \hat{v}(S_0, \mathbf{w}_0^3)] \nabla \hat{v}(S_0, \mathbf{w}_0^3)},$$

$$\cancel{\mathbf{w}_2^3 \doteq \mathbf{w}_1^3 + \alpha [G_{1:3}^\lambda - \hat{v}(S_1, \mathbf{w}_1^3)] \nabla \hat{v}(S_1, \mathbf{w}_1^3)},$$

$$\mathbf{w}_3^3 \doteq \cancel{\mathbf{w}_2^2} + \alpha [G_{\cancel{2:3}}^\lambda - \hat{v}(S_2, \cancel{\mathbf{w}_2^2})] \nabla \hat{v}(S_2, \cancel{\mathbf{w}_2^2}).$$

$G_{0:3}^\lambda$  with weights  $(1-\lambda), (1-\lambda)\lambda, (1-\lambda)\lambda^2$

How differs from TD(2)?

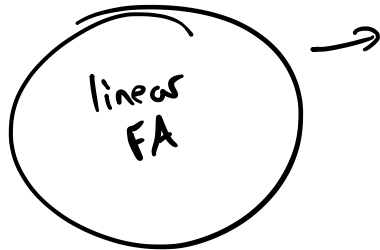
## Online $\lambda$ -return

$$h = 1: \quad \mathbf{w}_1^1 \doteq \mathbf{w}_0^1 + \alpha [G_{0:1}^\lambda - \hat{v}(S_0, \mathbf{w}_0^1)] \nabla \hat{v}(S_0, \mathbf{w}_0^1),$$

$$h = 2: \quad \mathbf{w}_1^2 \doteq \mathbf{w}_0^2 + \alpha [G_{0:2}^\lambda - \hat{v}(S_0, \mathbf{w}_0^2)] \nabla \hat{v}(S_0, \mathbf{w}_0^2),$$
$$\mathbf{w}_2^2 \doteq \mathbf{w}_1^2 + \alpha [G_{1:2}^\lambda - \hat{v}(S_1, \mathbf{w}_1^2)] \nabla \hat{v}(S_1, \mathbf{w}_1^2),$$

$$h = 3: \quad \mathbf{w}_1^3 \doteq \mathbf{w}_0^3 + \alpha [G_{0:3}^\lambda - \hat{v}(S_0, \mathbf{w}_0^3)] \nabla \hat{v}(S_0, \mathbf{w}_0^3),$$
$$\mathbf{w}_2^3 \doteq \mathbf{w}_1^3 + \alpha [G_{1:3}^\lambda - \hat{v}(S_1, \mathbf{w}_1^3)] \nabla \hat{v}(S_1, \mathbf{w}_1^3),$$
$$\mathbf{w}_3^3 \doteq \mathbf{w}_2^3 + \alpha [G_{2:3}^\lambda - \hat{v}(S_2, \mathbf{w}_2^3)] \nabla \hat{v}(S_2, \mathbf{w}_2^3).$$

True Online TD( $\lambda$ ) computes this  
fully incrementally using the "Dutch" trace  
(see Sutton's slides for illustration)



offline  $\lambda$ -return

$TD(\lambda)$

Truncated  $\lambda$ -return

Online  $\lambda$ -return

True Online  $TD(\lambda)$