# 12 Exercises

## 12.1 Question

Just as the return can be written recursively in terms of the first reward and itself one-step later (3.9), so can the $\lambda$-return. Derive the analogous recursive relationship from (12.2) and (12.1).

**Answer**

Equation 12.1:
$$G_{t_t+n} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \ldots + \gamma^{n-1} R_{t+n} + \gamma^n \hat{v}(S_{t+n}, w_{t+n-1})$$
Equation 12.2:
$$G_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_{t:t+n}$$

Let's start with expanding 12.2:
$$G_t^\lambda = (1-\lambda)[\lambda^0 G_{t:t+1} + \lambda^1 G_{t:t+2} + \lambda^2 G_{t:t+3} + \ldots]$$
$$G_t^\lambda = (1-\lambda)G_{t:t+1} + (1-\lambda)[\lambda^1 G_{t:t+2} + \lambda^2 G_{t:t+3} + \ldots]$$
$$G_t^\lambda = (1-\lambda)G_{t:t+1} + \lambda(1-\lambda)[\lambda^0 G_{t:t+2} + \lambda^1 G_{t:t+3} + \ldots]$$
$$G_t^\lambda = (1-\lambda)G_{t:t+1} + \lambda(1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_{t:t+1+n}$$

Let's not touch the left expression for now and focus on the second expression. We have a return of form $G_{t:t+1+n}$ but a return of form $G_{t+1:t+1+n}$ would be more useful. Let's try to express one in the form of the other.
$$G_{t+1:t+1+n} = R_{t+2} + \gamma R_{t+3} + \ldots + \gamma^{n-1} R_{t+1+n} + \gamma^n \hat{v}(s_{t+1+n}, w_{t+n}) \quad (12.1.1)$$
$$G_{t:t+1+n} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \ldots + \gamma^n R_{t+1+n} + \gamma^{n+1} \hat{v}(s_{t+1+n}, w_{t+n})$$
$$G_{t:t+1+n} = R_{t+1} + \gamma[R_{t+2} + \gamma R_{t+3} + \ldots + \gamma^{n-1} R_{t+1+n} + \gamma^n \hat{v}(s_{t+1+n}, w_{t+n})]$$
$$G_{t:t+1+n} = R_{t+1} + \gamma G_{t+1:t+1+n} \quad (12.1.2)$$

Going back and using equation (12.1.2):
$$G_t^\lambda = (1-\lambda)G_{t:t+1} + \lambda(1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1}(R_{t+1} + \gamma G_{t+1:t+1+n})$$
$$G_t^\lambda = (1-\lambda)G_{t:t+1} + \lambda(1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} R_{t+1} + \lambda(1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \gamma G_{t+1:t+1+n}$$

$G_t^\lambda = (1-\lambda)G_{t:t+1}+\lambda(1-\lambda)\sum_{n=1}^{\infty}\lambda^{n-1}R_{t+1}+\gamma\lambda(1-\lambda)\sum_{n=1}^{\infty}\lambda^{n-1}G_{t+1:t+1+n}$

We have a geometric series of form ($\sum_{n=1}^{\inf} r\lambda^{n-1} = \frac{r}{1-\lambda}$):
$G_t^\lambda = (1-\lambda)G_{t:t+1} + \lambda(1-\lambda)\frac{R_{t+1}}{(1-\lambda)} + \gamma\lambda(1-\lambda)\sum_{n=1}^{\infty}\lambda^{n-1}G_{t+1:t+1+n}$
$G_t^\lambda = (1-\lambda)G_{t:t+1} + \lambda R_{t+1} + \gamma\lambda(1-\lambda)\sum_{n=1}^{\infty}\lambda^{n-1}G_{t+1:t+1+n}$

The right most expression is a lambda return:
$G_t^\lambda = (1-\lambda)G_{t:t+1} + \lambda R_{t+1} + \gamma\lambda G_{t+1}^\lambda$

We can further simplify by expanding the 1-step return:
$G_t^\lambda = (1-\lambda)(R_{t+1} + \gamma G_{t+1:t+1}) + \lambda R_{t+1} + \gamma\lambda G_{t+1}^\lambda$
$G_t^\lambda = (1-\lambda)(R_{t+1} + \gamma\hat{v}(s_{t+1}, w_t)) + \lambda R_{t+1} + \gamma\lambda G_{t+1}^\lambda$
$G_t^\lambda = (1-\lambda)\gamma\hat{v}(s_{t+1}, w_t) + (1-\lambda)R_{t+1} + \lambda R_{t+1} + \gamma\lambda G_{t+1}^\lambda$
$G_t^\lambda = (1-\lambda)\gamma\hat{v}(s_{t+1}, w_t) + R_{t+1} + \gamma\lambda G_{t+1}^\lambda$
$G_t^\lambda = R_{t+1} + (1-\lambda)\gamma\hat{v}(s_{t+1}, w_t) + \gamma\lambda G_{t+1}^\lambda$ where $t < T - 1$

if $t \geq T - 1$ then by equation 12.3:
$G_t^\lambda = G_t$

## 12.2  Question

The parameter $\lambda$ characterizes how fast the exponential weighting in Figure 12.2 falls off, and thus how far into the future the $\lambda$-return algorithm looks in determining its update. But a rate factor such as $\lambda$ is sometimes an awkward way of characterizing the speed of the decay. For some purposes it is better to specify a time constant, or half-life. What is the equation relating $\lambda$ and the half-life, $\tau_\lambda$, the time by which the weighting sequence will have fallen to half of its initial value? □

### Answer

Half of initial weighting sequence: $\frac{(1-\lambda)}{2}$
Weighting sequence at time $t + m$: $(1-\lambda)\lambda^m$
Then by definition: $(1-\lambda)\lambda^m = \frac{(1-\lambda)}{2}$
$\lambda^m = \frac{1}{2}$
$m = \log_\lambda \frac{1}{2} = -\log_\lambda 2 = \frac{\ln 2}{\ln \lambda}$
The actual time step is:
$\tau_\lambda = t + m = t + \frac{\ln 2}{\ln \lambda}$

## 12.3  Question

Some insight into how $TD(\lambda)$ can closely approximate the offline $\lambda$-return algorithm can be gained by seeing that the latter's error term (in brackets in (12.4)) can be written as the sum of TD errors (12.6) for a single fixed w.

Show this, following the pattern of (6.6), and using the recursive relationship for the $\lambda$-return you obtained in Exercise 12.1.

### Answer

Error term in (12.4) $(w_{t+1} = w_t + \alpha[G_t^\lambda - \hat{v}(S_t, w_t)]\Delta\hat{v}(S_t, w_t))$:
$G_t^\lambda - \hat{v}(S_t, w_t)$

TD error (12.6):
$\delta_t = R_{t+1} + \gamma\hat{v}(S_{t+1}, w_t) - \hat{v}(S_t, w_t)$

Recursive relationship for the $\lambda$-return obtained in Exercise 12.1:
$G_t^\lambda = R_{t+1} + (1 - \lambda)\gamma\hat{v}(S_{t+1}, w_t) + \gamma\lambda G_{t+1}^\lambda$

Following the pattern of (6.6):
$G_t^\lambda - \hat{v}(S_t, w_t) = R_{t+1} + (1 - \lambda)\gamma\hat{v}(S_{t+1}, w_t) + \gamma\lambda G_{t+1}^\lambda - \hat{v}(S_t, w_t)$
$G_t^\lambda - \hat{v}(S_t, w_t) = R_{t+1} + (1-\lambda)\gamma\hat{v}(S_{t+1}, w_t) + \gamma\lambda G_{t+1}^\lambda - \hat{v}(S_t, w_t) - \gamma\hat{v}(S_{t+1}, w_t) + \gamma\hat{v}(S_{t+1}, w_t)$
$G_t^\lambda - \hat{v}(S_t, w_t) = \delta_t + (1 - \lambda)\gamma\hat{v}(S_{t+1}, w_t) + \gamma\lambda G_{t+1}^\lambda - \gamma\hat{v}(S_{t+1}, w_t)$
$G_t^\lambda - \hat{v}(S_t, w_t) = \delta_t - \lambda\gamma\hat{v}(S_{t+1}, w_t) + \gamma\lambda G_{t+1}^\lambda$
$G_t^\lambda - \hat{v}(S_t, w_t) = \delta_t - \lambda\gamma\hat{v}(S_{t+1}, w_t) + \gamma\lambda[R_{t+2} + (1 - \lambda)\gamma\hat{v}(S_{t+2}, w_{t+1}) + \gamma\lambda G_{t+2}^\lambda + \gamma\hat{v}(S_{t+2}, w_{t+1}) - \gamma\hat{v}(S_{t+2}, w_{t+1})]$
$G_t^\lambda - \hat{v}(S_t, w_t) = \delta_t + \gamma\lambda[\delta_{t+1} + (1-\lambda)\gamma\hat{v}(S_{t+2}, w_{t+1}) + \gamma\lambda G_{t+2}^\lambda - \gamma\hat{v}(S_{t+2}, w_{t+1})]$
$G_t^\lambda - \hat{v}(S_t, w_t) = \delta_t + \gamma\lambda[\delta_{t+1} - \lambda\gamma\hat{v}(S_{t+2}, w_{t+1}) + \gamma\lambda G_{t+2}^\lambda]$

for t = T-1, the inner-most expression will be:
$\delta_{T-2} - \gamma\lambda\hat{v}(S_{T-1}, w_{T-2}) + \gamma\lambda G_{T-1}^\lambda = \delta_{T-2} - \gamma\lambda\hat{v}(S_{T-1}, w_{T-2}) + \gamma\lambda G_t$
$= \delta_{T-2} - \gamma\lambda\hat{v}(S_{T-1}, w_{T-2}) + \gamma\lambda(R_T + \gamma\hat{v}(S_T, w_{T-1}))$
$= \delta_{T-2} + \gamma\lambda(R_T + \gamma\hat{v}(S_T, w_{T-1}) - \hat{v}(S_{T-1}, w_{T-2}))$
$= \delta_{T-2} + \gamma\lambda\delta_{T-1}$

Now we can write whole expression as a proper sum:
$G_t^\lambda - \hat{v}(S_t, w_t) = \sum_{k=t}^{T-1}(\lambda\gamma)^{(k-t)}\delta_k$

## 12.4   Question

Use your result from the preceding exercise to show that, if the weight updates over an episode were computed on each step but not actually used to change the weights (w remained fixed), then the sum of TD($\lambda$)'s weight

updates would be the same as the sum of the offline $\lambda$-return algorithm's updates.

## Answer

Equation 12.5:
$z_{-1} = 0$
$z_t = \gamma\lambda z_{t-1} + \Delta\hat{v}(S_t, w_t)$ where $0 \leq t \leq T$

Let's try to write as a sum:
$z_t = \gamma\lambda z_{t-1} + \Delta\hat{v}(S_t, w)$
$z_t = \gamma\lambda(\gamma\lambda z_{t-2} + \Delta\hat{v}(S_{t-1}, w)) + \Delta\hat{v}(S_t, w)$
$z_t = \sum_{k=0}^{t}(\gamma\lambda)^{k-t}\Delta\hat{v}(S_k, w)$

Sum of TD($\lambda$)'s weight updates:
$\alpha\sum_{t=0}^{\inf}\delta_t z_t = \alpha\sum_{t=0}^{\inf}\delta_t\sum_{k=0}^{t}(\gamma\lambda)^{t-k}\Delta\hat{v}(S_k, w)$

Expanding the sum:
t=0 $\alpha\delta_0[(\gamma\lambda)^0\Delta\hat{v}(S_0, w)]$
t=1 $\alpha\delta_1[(\gamma\lambda)^1\Delta\hat{v}(S_0, w) + (\gamma\lambda)^0\Delta\hat{v}(S_1, w)]$
t=2 $\alpha\delta_2[(\gamma\lambda)^2\Delta\hat{v}(S_0, w) + (\gamma\lambda)^1\Delta\hat{v}(S_1, w) + (\gamma\lambda)^0\Delta\hat{v}(S_2, w)]$
t=3 $\alpha\delta_3[(\gamma\lambda)^3\Delta\hat{v}(S_0, w) + (\gamma\lambda)^2\Delta\hat{v}(S_1, w) + (\gamma\lambda)^1\Delta\hat{v}(S_2, w) + (\gamma\lambda)^0\Delta\hat{v}(S_3, w)]$

Let's sum vertically:
$s = S_0$ $\alpha\Delta\hat{v}(S_0, w)\sum_{k=0}^{\inf}(\gamma\lambda)^k\delta_k$
$s = S_1$ $\alpha\Delta\hat{v}(S_1, w)\sum_{k=0}^{\inf}(\gamma\lambda)^k\delta_{k+1}$
$s = S_2$ $\alpha\Delta\hat{v}(S_2, w)\sum_{k=0}^{\inf}(\gamma\lambda)^k\delta_{k+2}$
$s = S_3$ $\alpha\Delta\hat{v}(S_3, w)\sum_{k=0}^{\inf}(\gamma\lambda)^k\delta_{k+3}$

$\alpha\sum_{t=0}^{\inf}\delta_t z_t = \sum_{t=0}^{\inf}\Delta\hat{v}(S_t, w)\sum_{k=0}^{\inf}(\lambda\gamma)^k\delta_{k+t}$

Start inner index k from t so that it is similar to result of question 12.3:
$\alpha\sum_{t=0}^{\inf}\delta_t z_t = \sum_{t=0}^{\inf}\Delta\hat{v}(S_t, w)\sum_{k=t}^{\inf}(\lambda\gamma)^{k-t}\delta_t$

Now we can use result of question 12.3
$\alpha\sum_{t=0}^{\inf}\delta_t z_t = \sum_{t=0}^{\inf}\Delta\hat{v}(S_t, w)(G_t^\lambda - \hat{v}(S_t, w_t))$

## 12.5 Question

Several times in this book (often in exercises) we have established that returns can be written as sums of TD errors if the value function is held constant. Why is (12.10) another instance of this? Prove (12.10).

### Answer

Equation 12.9 (h replaced with t+n):
$$G_{t:t+n}^{\lambda} = (1 - \lambda) \sum_{k=1}^{n-1} \lambda^{k-1} G_{t:t+k} + \lambda^{n-1} G_{t:t+n}$$

Get rid of $(1 - \lambda)$:
$$G_{t:t+n}^{\lambda} = \sum_{k=1}^{n-1} \lambda^{k-1} G_{t:t+k} - \sum_{k=1}^{n-1} \lambda^{k} G_{t:t+k} + \lambda^{n-1} G_{t:t+n}$$

Make the sums similar :
$$G_{t:t+n}^{\lambda} = \sum_{k=1}^{n} \lambda^{k-1} G_{t:t+k} - \sum_{k=1}^{n-1} \lambda^{k} G_{t:t+k}$$
$$G_{t:t+n}^{\lambda} = \sum_{k=0}^{n-1} \lambda^{k} G_{t:t+k+1} - \sum_{k=1}^{n-1} \lambda^{k} G_{t:t+k}$$
$$G_{t:t+n}^{\lambda} = \sum_{k=0}^{n-1} \lambda^{k} G_{t:t+k+1} - \sum_{k=0}^{n-1} \lambda^{k} G_{t:t+k} - (-G_{t:t+0})$$
$$G_{t:t+n}^{\lambda} = \sum_{k=0}^{n-1} \lambda^{k} G_{t:t+k+1} - \sum_{k=0}^{n-1} \lambda^{k} G_{t:t+k} + \hat{v}(S_t, w)$$
$$G_{t:t+n}^{\lambda} = \sum_{k=0}^{n-1} \lambda^{k} [G_{t:t+k+1} - G_{t:t+k}] + \hat{v}(S_t, w)$$

The returns can be written in the form of the other. One can expand both returns to show that:
$$G_{t:t+k+1} = G_{t:t+k} + \gamma^{k} R_{k+1} + \gamma^{k+1} \hat{v}(S_{k+1}, w) - \gamma^{k} \hat{v}(S_k, w)$$
$$G_{t:t+k+1} = G_{t:t+k} + \gamma^{k}(R_{k+1} + \gamma \hat{v}(S_{k+1}, w) - \hat{v}(S_k, w))$$
$$G_{t:t+k+1} = G_{t:t+k} + \gamma^{k} \delta_{t+k}$$

Apply the last finding to the equation:
$$G_{t:t+n}^{\lambda} = \sum_{k=0}^{n-1} \lambda^{k} [G_{t:t+k} + \gamma^{k} \delta_{t+k} - G_{t:t+k}] + \hat{v}(S_t, w)$$
$$G_{t:t+n}^{\lambda} = \sum_{k=0}^{n-1} \lambda^{k} \gamma^{k} \delta_{t+k} + \hat{v}(S_t, w)$$

Make $\delta$ indexed with i in the same way as in 12.10:
$$G_{t:t+n}^{\lambda} = \sum_{k=t}^{t+n-1} \lambda^{k-t} \gamma^{k-t} \delta_k + \hat{v}(S_t, w)$$
$$G_{t:t+n}^{\lambda} = \sum_{i=t}^{t+n-1} \lambda^{i-t} \gamma^{i-t} \delta_i + \hat{v}(S_t, w)$$

## 12.6 Question

Modify the pseudocode for Sarsa($\lambda$) to use dutch traces (12.11) without the other distinctive features of a true online algorithm. Assume linear function approximation and binary features.

# Answer

**Sarsa($\lambda$) with binary features and linear function approximation**
**for estimating $\mathbf{w}^\top \mathbf{x} \approx q_\pi$ or $q_*$**

Input: a function $\mathcal{F}(s, a)$ returning the set of (indices of) active features for $s, a$
Input: a policy $\pi$ (if estimating $q_\pi$)     <span style="color:red">x(S,A) gives related feature vector</span>
Algorithm parameters: step size $\alpha > 0$, trace decay rate $\lambda \in [0, 1]$
Initialize: $\mathbf{w} = (w_1, \ldots, w_d)^\top \in \mathbb{R}^d$ (e.g., $\mathbf{w} = \mathbf{0}$), $\mathbf{z} = (z_1, \ldots, z_d)^\top \in \mathbb{R}^d$

Loop for each episode:                            <span style="color:red">zp previous value of $z$, $\in \mathbb{R}^d$</span>
   Initialize $S$
   Choose $A \sim \pi(\cdot|S)$ or $\varepsilon$-greedy according to $\hat{q}(S, \cdot, \mathbf{w})$
   $\mathbf{z} \leftarrow \mathbf{0}$  <span style="color:red">zp = 0</span>
   Loop for each step of episode:
      Take action $A$, observe $R, S'$
      $\delta \leftarrow R$
      Loop for $i$ in $\mathcal{F}(S, A)$:
         $\delta \leftarrow \delta - w_i$
         ~~$z_i \leftarrow z_i + 1$~~     <span style="color:blue">dutch traces</span>     ~~(accumulating traces)~~
         ~~or $z_i \leftarrow 1$~~   <span style="color:red">z = z + (1-αγλ zp x(S,A))x(S,A)</span>   ~~(replacing traces)~~
      If $S'$ is terminal then:
         $\mathbf{w} \leftarrow \mathbf{w} + \alpha\delta\mathbf{z}$
         Go to next episode
      Choose $A' \sim \pi(\cdot|S')$ or near greedily $\sim \hat{q}(S', \cdot, \mathbf{w})$
      Loop for $i$ in $\mathcal{F}(S', A')$: $\delta \leftarrow \delta + \gamma w_i$
      $\mathbf{w} \leftarrow \mathbf{w} + \alpha\delta\mathbf{z}$  <span style="color:red">zp = z</span>
      $\mathbf{z} \leftarrow \gamma\lambda\mathbf{z}$
      $S \leftarrow S'; A \leftarrow A'$

6