# Temporal Difference Lambda

Given a sample path : $\{S_0, S_1, S_2, \ldots S_T\}$, with state space $S_t \in \{s_1, s_2, s_3, \ldots s_M\}$, where $t \in [0,T]$, the forward view and backward view of *temporal difference lambda* respectively are :

$$\Delta_{fwd} = \sum_{t=0}^{\infty} \alpha(G_{\lambda t}(s) - V(S_t))1_{S_t=s} \qquad \text{\textit{any limited path can be converted to } $\infty$ \textit{ path by repeating last state}}$$

$$\Delta_{back} = \sum_{t=0}^{\infty} \alpha \delta_t E_t(s)$$

$$\text{where} \quad G_t^{\lambda}(s) = (1-\lambda)\sum_{n=0}^{\infty} \lambda^n G_t^{(n)}(s) \qquad \text{\textit{where is (1-$\lambda$) normalization factor as } } \sum_{n=0}^{\infty} \lambda^n = (1-\lambda)^{-1}$$

$$G_t^{(n)}(s) = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \ldots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$

$$\text{and} \quad \delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \qquad \text{\textit{which is one step look ahead error}}$$

$$E_t(s) = (\lambda\gamma)E_{t-1}(s) + 1_{S_t=s} \qquad \text{\textit{called eligibility trace}}$$

Both views are equivalent. The forward view gives the intuition, while the backward view gives the implementation. Prior to proof, We have to unwire the recursive-form eligibility trace into iterative-form :

$$E_t(s) = (\lambda\gamma)E_{t-1}(s) + 1_{S_t=s} \qquad \text{where } s \in \{s_1, s_2, s_3, \ldots s_M\}$$

$$= (\lambda\gamma)(\lambda\gamma E_{t-2}(s) + 1_{S_{t-1}=s}) + 1_{S_t=s}$$

$$= (\lambda\gamma)^2 E_{t-2}(s) + (\lambda\gamma)1_{S_{t-1}=s} + 1_{S_t=s}$$

$$= (\lambda\gamma)^3 E_{t-3}(s) + (\lambda\gamma)^2 1_{S_{t-2}=s} + (\lambda\gamma)1_{S_{t-1}=s} + 1_{S_t=s}$$

$$= \ldots$$

$$= (\lambda\gamma)^t E_{t-t}(s) + \sum_{n=1}^{t} (\lambda\gamma)^{t-n} 1_{S_n=s}$$

$$= \sum_{n=0}^{t} (\lambda\gamma)^{t-n} 1_{S_n=s} \qquad \text{since } E_0(s) = 1_{S_0=s}$$

i.e.

$$E_0(s) = 1_{S_0=s}$$

$$E_1(s) = 1_{S_1=s} + (\lambda\gamma)1_{S_0=s}$$

$$E_2(s) = 1_{S_2=s} + (\lambda\gamma)1_{S_1=s} + (\lambda\gamma)^2 1_{S_0=s}$$

$$E_3(s) = 1_{S_3=s} + (\lambda\gamma)1_{S_2=s} + (\lambda\gamma)^2 1_{S_1=s} + (\lambda\gamma)^3 1_{S_0=s}$$

Lets start proving from forward view, there are 3 for loops, starting from outermost :
- backup starting from different states along one sample path
- sum of different TDs
- sum of different returns in one sample path

We have to group these 3 loops in forward view into 2 loops in backward view.

$$\Delta_{fwd} = \alpha\left[(1-\lambda)\times\left(\begin{array}{l}+\lambda^0(R_1+\gamma V(S_1))\\+\lambda^1(R_1+\gamma R_2+\gamma^2 V(S_2))\\+\lambda^2(R_1+\gamma R_2+\gamma^2 R_3+\gamma^3 V(S_3))\\+\lambda^3(R_1+\gamma R_2+\gamma^2 R_3+\gamma^3 R_4+\gamma^4 V(S_4))\\+\ldots\end{array}\right)-V(S_0)\right]1_{S_0=s} +$$

$$\alpha\left[(1-\lambda)\times\left(\begin{array}{l}+\lambda^0(R_2+\gamma V(S_2))\\+\lambda^1(R_2+\gamma R_3+\gamma^2 V(S_3))\\+\lambda^2(R_2+\gamma R_3+\gamma^2 R_4+\gamma^3 V(S_4))\\+\lambda^3(R_2+\gamma R_3+\gamma^2 R_4+\gamma^3 R_5+\gamma^4 V(S_5))\\+\ldots\end{array}\right)-V(S_1)\right]1_{S_1=s} +\ldots$$

The target is to group all $R_1$ terms together, all $R_2$ terms together, and so on, making use of $\sum_{n=0} \lambda^n = (1-\lambda)^{-1}$. Thus removes 1 loop.

$$\Delta_{fwd} = \begin{aligned} &\alpha\left[(1-\lambda)\times\begin{pmatrix}(1-\lambda)^{-1}R_1 + (1-\lambda)^{-1}(\lambda\gamma)R_2 + (1-\lambda)^{-1}(\lambda\gamma)^2 R_3 + (1-\lambda)^{-1}(\lambda\gamma)^3 R_4 + ...\\ + \gamma[(\lambda\gamma)^0 V(S_1) + (\lambda\gamma)^1 V(S_2) + (\lambda\gamma)^2 V(S_3) + (\lambda\gamma)^3 V(S_4) + ...]\end{pmatrix} - V(S_0)\right]1_{S_0=s} + \\ &\alpha\left[(1-\lambda)\times\begin{pmatrix}(1-\lambda)^{-1}R_2 + (1-\lambda)^{-1}(\lambda\gamma)R_3 + (1-\lambda)^{-1}(\lambda\gamma)^2 R_4 + (1-\lambda)^{-1}(\lambda\gamma)^3 R_5 + ...\\ + \gamma[(\lambda\gamma)^0 V(S_2) + (\lambda\gamma)^1 V(S_3) + (\lambda\gamma)^2 V(S_4) + (\lambda\gamma)^3 V(S_5) + ...]\end{pmatrix} - V(S_1)\right]1_{S_1=s} + ... \end{aligned}$$

$$= \begin{aligned} &\alpha\left[\begin{pmatrix}R_1 + (\lambda\gamma)R_2 + (\lambda\gamma)^2 R_3 + (\lambda\gamma)^3 R_4 + ...\\ + \gamma[(\lambda\gamma)^0 V(S_1) + (\lambda\gamma)^1 V(S_2) + (\lambda\gamma)^2 V(S_3) + (\lambda\gamma)^3 V(S_4) + ...]\\ -[(\lambda\gamma)^1 V(S_1) + (\lambda\gamma)^2 V(S_2) + (\lambda\gamma)^3 V(S_3) + (\lambda\gamma)^4 V(S_4) + ...]\end{pmatrix} - V(S_0)\right]1_{S_0=s} + \\ &\alpha\left[\begin{pmatrix}R_2 + (\lambda\gamma)R_3 + (\lambda\gamma)^2 R_4 + (\lambda\gamma)^3 R_5 + ...\\ + \gamma[(\lambda\gamma)^0 V(S_2) + (\lambda\gamma)^1 V(S_3) + (\lambda\gamma)^2 V(S_4) + (\lambda\gamma)^3 V(S_5) + ...]\\ -[(\lambda\gamma)^1 V(S_2) + (\lambda\gamma)^2 V(S_3) + (\lambda\gamma)^3 V(S_4) + (\lambda\gamma)^4 V(S_5) + ...]\end{pmatrix} - V(S_1)\right]1_{S_1=s} + ... \end{aligned}$$

Then group all terms according to $(\lambda\gamma)$, $(\lambda\gamma)^2$, $(\lambda\gamma)^3$ and so on.

$$\Delta_{fwd} = \alpha\begin{bmatrix}+(\lambda\gamma)^0(R_1 + \gamma V(S_1) - V(S_0))\\ +(\lambda\gamma)^1(R_2 + \gamma V(S_2) - V(S_1))\\ +(\lambda\gamma)^2(R_3 + \gamma V(S_3) - V(S_2))\\ +(\lambda\gamma)^3(R_4 + \gamma V(S_4) - V(S_3))\\ +...\end{bmatrix}1_{S_0=s} + \alpha\begin{bmatrix}+(\lambda\gamma)^0(R_2 + \gamma V(S_2) - V(S_1))\\ +(\lambda\gamma)^1(R_3 + \gamma V(S_3) - V(S_2))\\ +(\lambda\gamma)^2(R_4 + \gamma V(S_4) - V(S_3))\\ +(\lambda\gamma)^3(R_5 + \gamma V(S_5) - V(S_4))\\ +...\end{bmatrix}1_{S_1=s} + ...$$

$$= \begin{bmatrix}+\alpha(R_1 + \gamma V(S_1) - V(S_0))\times(1_{S_0=s})\\ +\alpha(R_2 + \gamma V(S_2) - V(S_1))\times(1_{S_1=s} + (\lambda\gamma)1_{S_0=s})\\ +\alpha(R_3 + \gamma V(S_3) - V(S_2))\times(1_{S_2=s} + (\lambda\gamma)1_{S_1=s} + (\lambda\gamma)^2 1_{S_0=s})\\ +\alpha(R_4 + \gamma V(S_4) - V(S_3))\times(1_{S_3=s} + (\lambda\gamma)1_{S_2=s} + (\lambda\gamma)^2 1_{S_1=s} + (\lambda\gamma)^3 1_{S_0=s})\\ +...\end{bmatrix}$$

$$= \sum_{t=0}^{\infty}\alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))E_t(s)$$

$$= \sum_{t=0}^{\infty}\alpha\delta_t E_t(s) \qquad\qquad where\ \delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$

$$= \Delta_{back}$$

Therefore, it is equivalent to backward view temporal difference lambda.