Going from the forward view to
  the backward view for TD($\lambda$):

Idea:- Store an additional memory variable
      for each state, called the eligibility trace,
      $e_t(s)$, also called the e-trace.

  → Quantifies how much s should be updated
    if there is a TD error in the current step

$$e_t(s) = \begin{cases} \gamma\lambda e_{t-1}(s) & \text{if } s \neq S_t \\ \gamma\lambda e_{t-1}(s)+1 & \text{otherwise} \end{cases}$$

called accumulating traces

    (Alternative is "Replacing traces".)

  $e_{-1}(s) = 0$ for all s

TD($\lambda$) Algorithm (Backward View)

$e_1(s) = 0$ for all $s$, for each episode

$\lambda = 0 \implies$ TD

$\lambda = 1 \implies$ MC $\longleftarrow$ allows MC updates without waiting to the end of the episode.

$0 < \lambda < 1 \implies$ TD($\lambda$)

$\delta = R_t + v(S_{t-1}) - v(S_t)$

$\forall s : e(s) \leftarrow \gamma \lambda e(s)$

$e(S_t) \leftarrow e(S_t) + 1$

$\forall s : v(s) \leftarrow v(s) + \alpha \delta e(s)$

$\longleftarrow$ Easy to apply to more general case (weight updates)

Proof of equiv. of forward + backward views:

[No hard math, just long!]

Want to show: $\swarrow$ Assumes same actions, rewards, + transitions.

$\forall s \in \mathcal{S} : \sum_{t=0}^{\ell-1} \Delta V_t^B(s) = \sum_{t=0}^{\ell-1} \Delta V_t^F(s) I_{ss_t}$

where $\ell$ is the length of the episode.

Notation: $\delta_t \triangleq R_t + \gamma V(S_{t+1}) - V(S_t)$

$I_{ss_t} \triangleq 1$ if $s = s_t$ and $0$ otherwise

$e_t(s) = \sum_{k=0}^{t} (\gamma \lambda)^{t-k} I_{ss_k}$

$\Delta V_t^F(s) \triangleq$ update at time $t$ of $V_t(s)$ $= \alpha (G_t^\lambda - V_t(s))$ (only $S_t$ is updated)

according to the Forward view

$\Delta V_t^B(s) \triangleq$ likewise for the Backward view $= \alpha \delta_t e_t(s)$

$$\sum_{t=0}^{l-1} \Delta V_t^B(s)$$

$$= \sum_{t=0}^{l-1} \alpha \delta_t e_t(s)$$

$$= \sum_{t=0}^{l-1} \alpha \delta_t \sum_{k=0}^{t} (\gamma\lambda)^{t-k} I_{SS_k}$$

$$= \sum_{k=0}^{l-1} \alpha \delta_k \sum_{t=0}^{k} (\gamma\lambda)^{k-t} I_{SS_t}$$

$$= \sum_{k=0}^{l-1} \alpha \sum_{t=0}^{k} (\gamma\lambda)^{k-t} I_{SS_t} \delta_k$$

$$= \sum_{t=0}^{l-1} \alpha \sum_{k=t}^{l-1} (\gamma\lambda)^{k-t} I_{SS_t} \delta_k$$

$$= \sum_{t=0}^{l-1} \alpha I_{SS_t} \sum_{k=t}^{l-1} (\gamma\lambda)^{k-t} \delta_k$$

Now, RHS for a single update:

$$\Delta V_t^F(s) = \alpha\left(G_t^\lambda - V_t(S_t)\right)$$

$$\frac{1}{\alpha}\Delta V_t^F(s) = -V_t(S_t) + G_t^\lambda$$

$$= -V_t(S_t) + (1-\lambda)\lambda^0\left(R_t + \gamma V_t(S_{t+1})\right)$$
$$+ (1-\lambda)\lambda^1\left(R_t + \gamma R_{t+1} + \gamma^2 V_t(S_{t+2})\right)$$
$$+ \dots$$

$$= -V_t(S_t) + (1-\lambda)\sum_{i=0}^{\infty}\lambda^i R_t$$
$$+ (1-\lambda)\gamma\sum_{i=0}^{\infty}\lambda R_{t+1}$$
$$+ \quad \vdots$$
$$+ \left[V_t \text{ terms}\right]$$

$$= -V_t(S_t) + \sum_{k=0}^{\infty}(\gamma\lambda)^k R_{t+k} + \sum_{k=0}^{\infty}(1-\lambda)\lambda^k \gamma^{k+1} V_t(S_{t+k+1})$$

$$= -V_t(S_t) + \sum_{k=0}^{\infty}(\gamma\lambda)^k\left(R_{t+k} + \gamma V_t(S_{t+k+1}) - \gamma\lambda V_t(S_{t+k+1})\right)$$

$$= \sum_{k=0}^{\infty}(\gamma\lambda)^k\left(R_{t+k} + \gamma V_t(S_{t+k+1}) - V_t(S_{t+k})\right)$$

$$= \sum_{k=0}^{\infty}(\gamma\lambda)^k \delta_k$$

The proof works if $\delta_t$ always uses $V_t$. But real updates don't do that in the Backward view. Using the actual $V$, the two views are equal up to a difference of $O(\alpha^2)$.

$$\sum_{t=0}^{l-1} \Delta V_t^F(S_t) I_{SS_t} = \sum_{t=0}^{l-1} \alpha I_{SS_t} \sum_{k=t}^{l-1} (\gamma\lambda)^{k-t} \delta_k$$

Coming up: In some cases can adjust Backward View a little and get the Forward + Backward equivalent.

Will also apply $TD(\lambda)$ idea to Sarsa (Sarsa $(\lambda)$) and q-learning $(q(\lambda))$.