

Temporal Difference Learning

Sargur N. Srihari

srihari@cedar.buffalo.edu

Topics in Temporal Difference Learning

1. Temporal difference learning

Temporal Difference

- Temporal difference is an important concept at the heart of the Q-learning algorithm
 - It is how everything comes together in Q-learning

- The temporal difference is defined as:

$$TD(a,s) = \underbrace{r(s,a) + \gamma \max_{a'} Q_n(s',a')} - \underbrace{Q_{n-1}(s,a)}$$

- The first element is what we get after taking an action and the second element is the previous Q-value
- Thus $Q_n(s,a) = Q_{n-1}(s,a) + \alpha TD_n(a,s)$
 - Where α is the learning rate
 - Also
 - $Q_n(s,a) = Q_{n-1}(s,a) + \alpha (r(s,a) + \gamma \max_{a'} Q(s',a') - Q_{n-1}(s,a))$

Temporal difference (TD) learning

- The Q learning algorithm iteratively reduces discrepancy between adjacent states
 - Thus it is a special case of temporal difference algorithms
- Whereas the training rule
$$\hat{Q}_n(s, a) \leftarrow (1 - \alpha_n) \hat{Q}_{n-1}(s, a) + \alpha_n \left[r + \gamma \max_{a'} \hat{Q}_{n-1}(s', a') \right]$$
 - Reduces the difference between estimated \hat{Q} values of a state s and its immediate successor s'
- We could just as well design an algorithm that reduces discrepancies between the state and its more distant descendants or ancestors

Temporal difference over n -steps

- Let $Q^{(1)}(s_t, a_t)$ denote the training value calculated by the one-step look-ahead

$$Q^{(1)}(s_t, a_t) \leftarrow r_t + \gamma \max_a \hat{Q}(s_{t+1}, a)$$

- One alternative to compute a training value for $Q(s_t, a_t)$ is to base it on the observed reward for two steps

$$Q^{(2)}(s_t, a_t) \leftarrow r_t + \gamma r_{t+1} + \gamma^2 \max_a \hat{Q}(s_{t+2}, a)$$

- Or in general

$$Q^{(n)}(s_t, a_t) \leftarrow r_t + \gamma r_{t+1} + \dots + \gamma^{(n-1)} r_{t+n-1} + \gamma^n \max_a \hat{Q}(s_{t+n}, a)$$

Blending the training estimates

- TD(λ) is a method to blend the alternative training estimates obtained from various lookahead distances

$$Q^\lambda(s_t, a_t) \leftarrow (1 - \lambda) \left[Q^{(1)}(s_t, a_t) + \lambda Q^{(2)}(s_t, a_t) + \lambda^2 Q^{(3)}(s_t, a_t) + \dots \right]$$

- An equivalent recursive definition is

$$Q^\lambda(s_t, a_t) \leftarrow r_t + \gamma \left[(1 - \lambda) \max_a \hat{Q}(s_t, a_t) + \lambda Q^\lambda(s_{t+1}, a_{t+1}) \right]$$

- If we choose $\lambda=0$ we have our original training estimate $Q^{(1)}$, which considers only one-step discrepancies in the \hat{Q} estimates
 - As λ is increased the algorithm places increasing emphasis on discrepancies based on more distant lookaheads

Applications of TD Learning

- In some settings training will be more efficient if more distant lookaheads are considered
- $TD(\lambda)$ approach was used in the TD-Gammon program for playing backgammon
- Received attention in neuroscience
 - firing rate of dopamine neurons mimic error function in the algorithm
 - study schizophrenia
 - consequences of pharmacological manipulations of dopamine on learning