

A Guided Tour of Chapter 11: Batch RL: Experience Replay, DQN, LSPI

Ashwin Rao

ICME, Stanford University

Incremental RL makes inefficient use of training data

- Incremental versus Batch RL in the context of fixed finite data
- Let's understand the difference for the simple case of MC Prediction
- Given fixed finite sequence of trace experiences yielding training data:

$$\mathcal{D} = [(S_i, G_i) | 1 \leq i \leq n]$$

- Incremental MC estimates $V(s; \mathbf{w})$ using $\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w})$ for each data pair:

$$\mathcal{L}_{(S_i, G_i)}(\mathbf{w}) = \frac{1}{2} \cdot (V(S_i; \mathbf{w}) - G_i)^2$$

$$\nabla_{\mathbf{w}} \mathcal{L}_{(S_i, G_i)}(\mathbf{w}) = (V(S_i; \mathbf{w}) - G_i) \cdot \nabla_{\mathbf{w}} V(S_i; \mathbf{w})$$

$$\Delta \mathbf{w} = \alpha \cdot (G_i - V(S_i; \mathbf{w})) \cdot \nabla_{\mathbf{w}} V(S_i; \mathbf{w})$$

- n updates are performed in sequence for $i = 1, 2, \dots, n$
- Uses update method of FunctionApprox for each data pair (S_i, G_i)
- Incremental RL makes inefficient use of available training data \mathcal{D}
- Essentially each data point is “discarded” after being used for update

Batch MC Prediction makes efficient use of training data

- Instead we'd like to estimate the Value Function $V(s; \mathbf{w}^*)$ such that

$$\begin{aligned} \mathbf{w}^* &= \arg \min_{\mathbf{w}} \frac{1}{n} \cdot \sum_{i=1}^n \frac{1}{2} \cdot (V(S_i; \mathbf{w}) - G_i)^2 \\ &= \arg \min_{\mathbf{w}} \mathbb{E}_{(S, G) \sim \mathcal{D}} \left[\frac{1}{2} \cdot (V(S; \mathbf{w}) - G)^2 \right] \end{aligned}$$

- This is the solve method of FunctionApprox on training data \mathcal{D}
- This approach to RL is known as *Batch RL*
- solve by doing updates with repeated use of available data pairs
- Each update using random data pair $(S, G) \sim \mathcal{D}$

$$\Delta \mathbf{w} = \alpha \cdot (G - V(S; \mathbf{w})) \cdot \nabla_{\mathbf{w}} V(S; \mathbf{w})$$

- This will ultimately converge to desired value function $V(s; \mathbf{w}^*)$
- Repeated use of available data known as *Experience Replay*
- This makes more efficient use of available training data \mathcal{D}

Batch TD Prediction makes efficient use of *Experience*

- In Batch TD Prediction, we have experience \mathcal{D} available as:

$$\mathcal{D} = [(S_i, R_i, S'_i) | 1 \leq i \leq n]$$

- Where (R_i, S'_i) is the pair of reward and next state from a state S_i
- So, Experience \mathcal{D} in the form of finite number of atomic experiences
- This is represented in code as an Iterable [TransitionStep [S]]
- Parameters updated with repeated use of these atomic experiences
- Each update using random data pair $(S, R, S') \sim \mathcal{D}$

$$\Delta \mathbf{w} = \alpha \cdot (R + \gamma \cdot V(S'; \mathbf{w}) - V(S; \mathbf{w})) \cdot \nabla_{\mathbf{w}} V(S; \mathbf{w})$$

- This is TD Prediction with Experience Replay on Finite Experience \mathcal{D}

Batch TD(λ) Prediction

- In Batch TD(λ) Prediction, given finite number of trace experiences

$$\mathcal{D} = [(S_{i,0}, R_{i,1}, S_{i,1}, R_{i,2}, S_{i,2}, \dots, R_{i,T_i}, S_{i,T_i}) | 1 \leq i \leq n]$$

- Parameters updated with repeated use of these trace experiences
- Randomly pick trace experience (say indexed i) $\sim \mathcal{D}$
- For trace experience i , parameters updated at each time step t :

$$\mathbf{E}_t = \gamma \lambda \cdot \mathbf{E}_{t-1} + \nabla_{\mathbf{w}} V(S_{i,t}; \mathbf{w})$$

$$\Delta \mathbf{w} = \alpha \cdot (R_{i,t+1} + \gamma \cdot V(S_{i,t+1}; \mathbf{w}) - V(S_{i,t}; \mathbf{w})) \cdot \mathbf{E}_t$$

Key Takeaways from this Chapter

- Batch RL makes efficient use of data
- DQN uses experience replay together with a frozen DNN avoiding pitfalls of time-correlation and semi-gradient
- LSTD is a fast linear-algebra calculation of Batch TD Prediction
- LSPI is an off-policy, experience-replay Control Algorithm using LSTD for Policy Evaluation