

• CME 241 Notes:

- Incremental RL Makes inefficient use of training data
 - ↳ Incremental ML estimates $V(s; \mathbf{w})$ using $\nabla_{\mathbf{w}} L(\mathbf{w})$ ← Gradient descent
 - ↳ You would call update method in function approx.
 - ↳ Inefficient use of each data point. Point is discarded after each update

• Batch Monte Carlo Prediction

- Given batch data (s_i, g_i)
 - ↳ Use solve method of FunctionApprox
 - ↳ Given data, find best fitting weights
 - ↳ You can store batch of data in database and randomly sample points with replacement
 - ↳ Use the same datapoints again and again.
 - ↳ Converges to weight \mathbf{w}^*
 - ↳ Reuse of atomic or trace experience is experience replay
 - ↳ Makes efficient use of training data D
 - $D = [(s_i, r_i, s'_i) | 1 \leq i \leq n]$
 - ↳ In code, these are "Transition Step"
 - ↳ Experience Replay on finite experiences D

• Batch TD(λ) Prediction

- $D = [(s_{i,1}, r_{i,1}, s_{i,2}, r_{i,2}, \dots) | 1 \leq i \leq n]$
 - ↳ Randomly pick trace experiences from D

The Deep Q-Networks (DQN) Control Algorithm

- Uses Experience Replay and fixed Q-learning targets ← Two innovations
- One of the first works to use experience replay
- At each t from each episode:
 - ① Given state s_t , take action a_t according to ϵ -greedy policy extracted from Q-network values $Q(s_t, a; w)$
 - ② Obtain reward r_{t+1} and next state s_{t+1}
 - ③ Store atomic experience $(s_t, a_t, r_{t+1}, s_{t+1})$ in replay memory D
 - ④ Sample random minibatch of atomic experiences
 - ⑤ Update Q-network parameters w using Q-learning targets build on "frozen" parameters w^* of target network

↑ This is updated every π steps.
- Clever techniques / variants not explicitly proven = "hacky"

Least Squares RL Prediction

- Batch RL Prediction for general function approx. is iterative
- Uses experience replay and gradient descent.
- We can use a direct solve (without gradients) for linear function approx

Least Squares Monte Carlo

$$\mathcal{L}(w) = \frac{1}{2n} \sum_{i=1}^n (\phi(s_i)^T \cdot w - b_i)^2$$

→ Set gradient of \mathcal{L} wrt $w = 0$ and solve for w^*

w^* is solved by $A^T \cdot b$

→ See slide 8 for development of A, b

→ Least Squares Temporal Difference: Slide 9

A is formed from outer product of features (see A update rule slide)

→ Updated incrementally for efficient Sherman - Morrison Inverse

→ Least Squares Algorithm limited to linear function approx.

→ Recommended if you can get few features and linear function approx.

→ In reality, this is impractical

→ Q-learning is off policy (experience policy coming from database of past experience).

→ Approach also uses linear function approx.

→ Semi-gradient as in TD avoids bootstrapping, which would complete the deadly triangle.

→ LSTD(γ) is a batch form of TD(γ) using eligibility traces

Accumulates $\nabla_w V(s; w) = t(s)$

E_i = Eligibility trace at atomic experience i :

$$A \leftarrow A + E_i \cdot (t(s) - \gamma \cdot t(s'))^T$$

↑
Tracks how many times feature occurs in the path π timesteps weighted by recency.

Implement LSTD(γ) algorithm, building A, b incrementally

• Least Squares RL Control

• Policy Evaluation as Least-Squares Q-value Prediction

• Greedy Policy Improvement

• Least-Squares Policy Iteration (slide 12-13)

① Get $Q(s, a; w)$ using function approximation with initial feature, weight

④ Get Greedy policy to maximize Q

⑤ Sample minibatch of experiences

⑥ Goal of iteration to solve weights w^* to minimize $J(w)$
↳ Using Semigradient
↳ See slide 13

• We learned about LSPI batch for Q-value function (LSPIQ)

↳ This sits inside LSPI for control problem.

↳ LSPI is a linear function approximation

↳ Fluctuates about optimum in linear function approximation

• LSPI for Optimal Exercise of American Options

• LSPI can be prediction of Complex American Option

• Really, it is never optimal to exercise early for an American Call Option

• LSPI converges much faster than Longstaff-Schwartz algorithm

• Only use linear function approx when you are holding the option.

↳ When exercising, you know the true value

• Each data point in atomic experience database is:

(s, c, o, s')

Only store actions (left) action option
You will not get a reward when you hold.

• GitHub project on exotic options