# Reinforcement Learning: LSTD, LSTDQ, LSPI, Fitted Q

B. Ravindran

# LSTD & LSTDQ

$$v^\pi \in \mathbb{R}^N$$
[N States]

$\hat{v}^{opt}$

$\hat{v}*$

$$\Phi \in \mathbb{R}^k$$
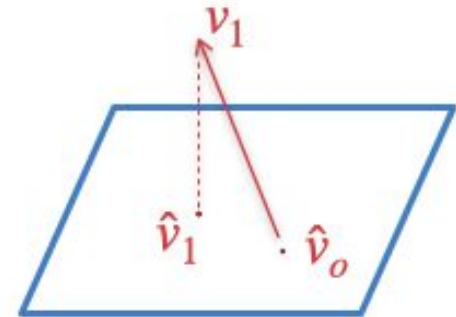[k << N]

An overview of linear function approximations:

- We look for projection of real value function $v^\pi$ to the space spanned by feature vectors $\Phi$

- Given the true $v^\pi$ values, $\hat{v}^{opt}$ is the projection. Since, we do not know true $v^\pi$, we converge at $\hat{v}*$ following the procedure.

"Why the gap?"

Procedure:

- Start with $\hat{v}_o$

- Apply Update $v_1 = \hat{T}_\pi \hat{v}_0$          Final Update (each iteration)

- Project back $\hat{v}_1 = \mathbb{P}v_1$

- Repeat till convergence

$v_1$

$\hat{v}_1$   $\hat{v}_o$

# LSTD & LSTDQ

Setting up as a traditional 'Regression Problem':

Traditional Data Representation

Equivalent for Value Function Approximation

Fixed point solution (Desired   ):   $(\mathbb{P}T_\pi)\hat{v}^\pi = \hat{v}^\pi$

$$\hat{v}^\pi = \underline{\Phi(\Phi^T\Phi)^{-1}\Phi}\underline{(R^\pi + \gamma P^\pi\hat{v}^\pi)}$$

Projection     Bellman

Similarly

$$\hat{Q}^\pi = \Phi(\Phi^T\Phi)^{-1}\Phi(R^\pi + \gamma P^\pi\hat{Q}^\pi)$$

$$\Phi\theta^\pi = \Phi(\Phi^T\Phi)^{-1}\Phi(R^\pi + \gamma P^\pi\Phi\theta^\pi)$$

Therefore, The Least Square solution:

$$\theta^\pi = (\Phi^T(\Phi - \gamma P^\pi\Phi))^{-1}\Phi^T R^\pi$$

NOTE:

- We do NOT have the model

- This solution treats every sample as an equal contributor to the error.

- We wish to prioritise some states over others.

Weighted Version:

$$\underline{\Phi^T w(\Phi - \gamma P^\pi\Phi)}\theta^\pi = \underline{\Phi^T w R^\pi}$$

A     b

W - Weight function over all states
(eg: steady-state probability distribution)

# LSTD & LSTDQ

Weighted Version:

$$\Phi^T w (\Phi - \gamma P^\pi \Phi) \theta^\pi = \Phi^T w R^\pi$$

$$\underbrace{\Phi^T w (\Phi - \gamma P^\pi \Phi) \theta^\pi}_{A} = \underbrace{\Phi^T w R^\pi}_{b}$$

We wish to estimate A & B matrices directly from samples rather than, from the model based P & R matrices.

$$A = \Phi^T w (\Phi - \gamma P^\pi \Phi)$$

$$A = \Sigma_s \phi(s) w(s) \left( \Phi(s) - \gamma \Sigma_{s'} P(s, \pi(s), s') \right)^T$$

$$A = \Sigma_s w(s) \Sigma_{s'} P(s, \pi(s), s') \left[ \phi(s) \left( \phi(s) - \gamma \phi(s') \right)^T \right]$$

$$b = \Phi^T w R^\pi$$

$$b = \Sigma_s \Phi(s) w(s) \Sigma_{s'} P(s, \pi(s), s') R(s, \pi(s), s')$$

# LSTD & LSTDQ

Given Samples of form:

$$D = \{(s_i, a_i, r_i, s_i') \mid i = 1, 2, 3, \ldots L\}$$

[i -> sample index; L - no. of samples]

Estimate from samples:

$$\tilde{A} = \frac{1}{L} \Sigma_{i=1}^{L} [\phi(s_i) \, (\phi(s_i) - \gamma\phi(s_i'))^T]$$

$$\tilde{B} = \frac{1}{L} \Sigma_{i=1}^{L} \phi(s_i) r_i$$

Action version:

$$\tilde{A} = \frac{1}{L} \Sigma_{i=1}^{L} [\phi(s_i) \, (\phi(s_i) - \gamma\phi(s_i', \pi(s_i')))^T]$$

$$\tilde{B} = \frac{1}{L} \Sigma_{i=1}^{L} \phi(s_i, a_i) r_i$$

Now, we can solve for
and thereby the value function.

$$\tilde{A} \cdot \theta^\pi = \tilde{b}$$

$$\pi_{t+1}(s) = argmax_a \, \hat{Q}^{\pi_t}(s, a)$$

$$\pi_{t+1} = argmax \, \Phi \hat{\theta}^{\pi_t}$$

- Solve for (some) 'j' no of states
- Generate training data for classifier (finite actions) or regressor (continuous actions)

$a_1$      $\phi_1(s_1), \phi_2(s_1), \ldots \ldots, \phi_k(s_1)$

$a_2$      $\phi_1(s_2), \phi_2(s_2), \ldots \ldots, \phi_k(s_2)$

⋮         ⋮

$a_j$      $\phi_1(s_j), \phi_2(s_j), \ldots \ldots, \phi_k(s_j)$

Training data

Issues:
Actions and States might be continuous and complex
Maximise each time to pick action?
Solve for all the states?

Sample for 'j' states

$$(s_1, s_2, \ldots \ldots s_j) \quad -> States$$
$$(a_1, a_2, \ldots \ldots a_j) \quad -> \pi_{t+1}(s)$$

# Fitted-Q Iteration

- Start with few samples. Generate Targets (Q-Learning targets).
  Train(Fit) function approximator (Not necessarily linear) to predict the targets.
- Use function approximator to form policy to generate new samples.
- Repeat.

Begin with $\hat{Q}_o$ $\longrightarrow$ Generate Training Data:

$$< s_i, a_i, s_i', r_i >$$
$$< \Phi(s_i, a_i), target_i >$$
$$target_i = r_i + \gamma max_a \hat{Q}_o(s_i, a)$$

Iterate

Estimate $\hat{Q}_1$ $\longleftarrow$ Feed to Function Approximator