

Value Function Approximation:

- RL is usually used for large problems, eg.
- Backgammon - 10^{20} states
- Go - 10^{170} states
- Helicopter - continuous state space.
- We cannot (in almost all cases) create tables for states.
- So, we need some way to extend the model-free control methods for prediction and control.

Until now:

- Value function was represented as a lookup table
i.e. every state s has an entry $V(s)$ or
for every (s,a) pair, there was an entry $Q(s,a)$.

Problem with Large MDP's:

- 1) Too many states / action pairs to store in memory.
- 2) Too slow to learn the value of each state individually.

Solution for Large MDP's:

- Estimate value function with function approximation.

$$\hat{V}(s, w) \approx V_{\pi}(s) \text{ or } \hat{Q}(s, a, w) \approx Q_{\pi}(s, a).$$

- 1) Generalise from seen states to unseen states.
- 2) Update param w using MC or TD Learning.

Types of Value-Function Approximation:

→ We either approximate V or Q .

For V : $s \rightarrow [w] \rightarrow \hat{V}(s, w)$.

For Q :

① $s, a \rightarrow [w] \rightarrow \hat{Q}(s, a, w)$

② $s \rightarrow [w] \rightarrow \hat{Q}(s, a_1, w), \dots, \hat{Q}(s, a_n, w)$.

→ Which Function Approximator?

examples:

- 1) Linear combination of features
- 2) Neural Network.
- 3) DT's
- 4) KNN's. etc.

For RL, we need differentiable function approximators, so they are Linear combination and Neural Networks.

2 Methods for Function Approximation:

- 1) Incremental Methods
- 2) Batch Methods.

Gradient Descent

→ Let $J(w)$ be a differentiable function of parameter vector w .

$$\nabla_w J(w) = \begin{pmatrix} \frac{\partial J}{\partial w_1} \\ \vdots \\ \frac{\partial J}{\partial w_n} \end{pmatrix}_w.$$

→ Goal: to find a local minima of J .

So, we adjust w in the direction of -ve gradient.

$$\Delta w = -\frac{1}{2} \alpha \nabla_w J(w), \quad (\alpha \text{ is step-size})$$

Value Function Approx by Stochastic Gradient Descent:

Goal: Find param w which minimises MSE between $\hat{V}(s, w)$ and true value $f^* v_\pi(s)$.

$$J(w) = E_\pi \left[\left(v_\pi(s) - \hat{V}(s, w) \right)^2 \right].$$

Assume for now, that v_π is present (or given to you)

→ GD finds a local minimum:

$$\Delta w = -\frac{1}{2} \alpha \nabla_w J(w).$$

$$= \alpha E_\pi \left[(v_\pi(s) - \hat{V}(s, w)) \nabla_w \hat{V}(s, w) \right]$$

→ SGD samples the Gradient.

$$\Delta \omega = \alpha (V_{\pi}(s) - \hat{V}(s, \omega)) \nabla_{\omega} \hat{V}(s, \omega).$$

Feature Vector

Represent state by a feature vector.

$$x(s) = \begin{pmatrix} x_1(s) \\ \vdots \\ x_n(s) \end{pmatrix}$$

ex: 1) Distance of a robot from landmarks.

2) Piece configurations in Chess.

Linear Value Function Approximation:

→ Represent Value f^n by a linear combi of features.

$$\hat{V}(s, \omega) = x(s)^T \omega = \sum_{j=1}^n x_j(s) \omega_j$$

→ Objective function is quadratic in ω , i.e.

$$J(\omega) = E_{\pi} \left[(V_{\pi}(s) - x(s)^T \omega)^2 \right]$$

→ SGD will converge on global optimum.

$$\nabla_{\omega} \hat{V}(s, \omega) = x(s).$$

$$\Delta \omega = \alpha (V_{\pi}(s) - \hat{V}(s, \omega)) x(s).$$

So, update = ~~step~~ step-size \times prediction error \times feature value.

Table Lookup Features:

→ Table lookup is a special case of linear value function approximation

→ Using table lookup features:

$$x^{\text{table}}(s) = \begin{pmatrix} 1(s=s_1) \\ \vdots \\ 1(s=s_n) \end{pmatrix}$$

→ parameter \underline{w} gives value of each individual state

$$\hat{v}(s, w) = \begin{pmatrix} 1(s=s_1) \\ \vdots \\ 1(s=s_n) \end{pmatrix} \cdot \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix}.$$

Incremental Prediction Algos:

→ we assumed true value $f^* v_{\pi}(s)$ given by a supervisor.

→ But obviously this is cheating, we only have rewards.

→ In practice, we substitute a target for $v_{\pi}(s)$.

→ For MC, the target is G_t

$$\Delta w = \alpha (G_t - \hat{v}(s_t, w)) \nabla_w \hat{v}(s_t, w)$$

→ For TD(0), the target is the TD target $R_{t+1} + \gamma \hat{v}(s_{t+1}, w)$

$$\Delta w = \alpha (R_{t+1} + \gamma \hat{v}(s_{t+1}, w) - \hat{v}(s_t, w)) \nabla_w \hat{v}(s_t, w).$$

→ For TD(1), the target is the λ return G_t^λ

$$\Delta w = \alpha (G_t^\lambda - \hat{v}(s_t, w)) \nabla_w \hat{v}(s_t, w)$$

1) Monte Carlo:

→ G_t is an unbiased, noisy sample of true value $v_\pi(s_t)$.

→ We can therefore treat

$\langle s_1, G_1 \rangle, \langle s_2, G_2 \rangle, \dots, \langle s_T, G_T \rangle$ as training data in

Supervised Learning.

→ So, if our loss was MSE, the update would be

$$\Delta w = \alpha (G_t - \hat{v}(s_t, w)) \nabla_w \hat{v}(s_t, w)$$

~~see~~

2) TD learning:

→ The TD target $R_{t+1} + \gamma \hat{v}(s_{t+1}, w)$ is a biased sample of true value $v_\pi(s_t)$.

→ we can still apply supervised learning to "training data".

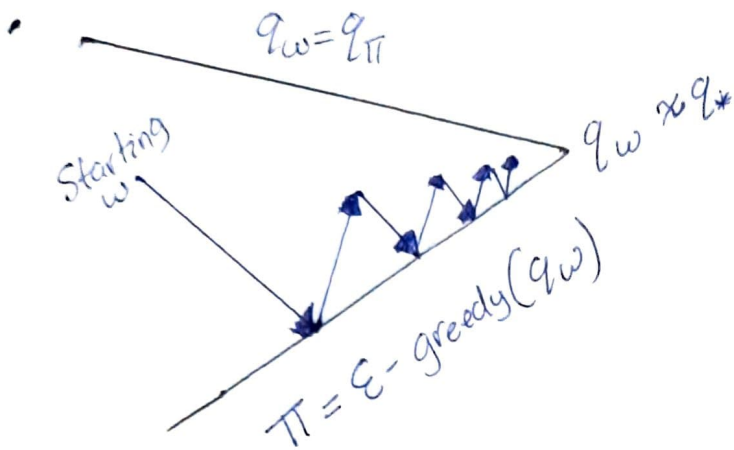
$$(s_1, R_2 + \gamma \hat{v}(s_2, w)), \dots, (s_{t-1}, R_t)$$

ex: For MSE, $\Delta w = \alpha (R + \gamma \hat{v}(s', w) - \hat{v}(s, w)) \nabla_w \hat{v}(s, w)$

Control with Value Function Approximation:

→ Policy Evaluation: Approximate policy evaluation, $\hat{q}(\cdot, \cdot, \omega) \approx q_\pi$

→ Policy Improvement: ϵ -greedy policy improvement.



Action Value Function Approximation:

→ Approximate the action-value function

$$\hat{q}(s, A, \omega) \approx q_\pi(s, A)$$

→ Minimise MSE between approximate action value $f^\pi \hat{q}(s, A, \omega)$ and true action-value $f^\pi q_\pi(s, A)$

$$J(\omega) = E_\pi \left[\left(q_\pi(s, a) - \hat{q}(s, a, \omega) \right)^2 \right]$$

→ Use SGD to find a local minimum.

→ Represent state and action by a feature vector.

$$X(s, A) = \begin{pmatrix} x_1(s, a) \\ \vdots \\ x_n(s, a) \end{pmatrix}$$

Simplest Case: Linear Approximator

$$\hat{q}(s, a, \omega) = x(s, a)^T \omega = \sum_{j=1}^n x_j(s, a) \omega_j$$

→ Do the same as with Value f^n approximation to approximate $q(s, a)$.

Policy Based RL

- Until now, we were approximating the value function or the Q-value function.
- The policy was generated directly from the value function.

- Now, we will directly parametrise / approximate the policy.

$$\pi_{\theta}(s, a) = P[a | s, \theta]$$

Value Based

- Learnt Value function

- Implicit Policy

Policy Based

- No value function

- Learnt policy

Actor Critic

- learn both value f^* and policy

Advantages of Policy-Based

- Better convergence properties
- Effective in high-dimensional or continuous action-spaces
↓
most important reason
- Can learn stochastic policies.

Disadvantages.

- Typically converge to a local rather than global optimum
- Evaluating a policy is typically inefficient and high variance.

Q) Why would we need stochastic policies?

ex: Rock-paper-Scissors.

→ A deterministic policy is easily exploited.

Policy Objective Functions:

Goal: Given policy $\pi_\theta(s, a)$ with parameters θ , find θ .

→ But how do we measure the quality of a policy π_θ ?

1) In episodic environments we can use the start value.

$$J_1(\theta) = V^{\pi_\theta}(s_1) = E_{\pi_\theta}[V]$$

a) In continuing environments we can use the average value

$$J_{avV}(\theta) = \sum_s d^{\pi_\theta}(s) V^{\pi_\theta}(s)$$

or the average reward per time step

$$J_{avR}(\theta) = \sum_s d^{\pi_\theta}(s) \sum_a \pi_\theta(s, a) R_s^a$$

→ where $d^{\pi_\theta}(s)$ is stationary distribution of Markov chain for π_θ .

Policy Optimisation:

- Policy Based RL is an optimisation problem.
- Find θ that maximises $J(\theta)$.

Policy Gradient:

- Let $J(\theta)$ be any policy objective function.
- Policy gradient algos search for a local maxima in $J(\theta)$ by ascending the gradient of the policy w.r.t params θ .

$$\Delta \theta = \alpha \nabla_{\theta} J(\theta);$$

where $\nabla_{\theta} J(\theta) = \begin{pmatrix} \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_n} \end{pmatrix}(\theta)$ α is step-size.

Computing Gradients by Finite Differences:

- To evaluate policy gradient of $\pi_{\theta}(s, a)$

For each dim $k \in [1, n]$,

- Estimate the k^{th} p.d of J w.r.t θ as follows,

$$\frac{\partial J}{\partial \theta_k} \approx \frac{J(\theta + \epsilon U_k) - J(\theta)}{\epsilon}, \text{ where } U_k \text{ is } (0, 0, \dots, \underset{\substack{\uparrow \\ k^{\text{th}} \text{ position}}}{1}, \dots, 0)$$