

Function Approximate:

$$G_0 \rightarrow |S| = 10^{170}$$

Helicopter \rightarrow state is cts.
Drone

Instead of represent $V(s)$, $Q(s, a)$ in a look up table /
Tabular case

we use a function approximation.

$\hat{V}(s; w)$ the value fn is determined by parameter w .

$\hat{q}(s, a; w)$

- Even if the data doesn't see the state, using fn approx will generalise to unseen state.

i). $s \rightarrow \boxed{w} \rightarrow V(s; w)$
input
is determined by w

$$s \rightarrow \boxed{w} \rightarrow Q(s, a; w)$$

$$s \rightarrow \boxed{w} \rightarrow \begin{array}{l} Q(s, a_1; w) \\ Q(s, a_2; w) \\ \vdots \\ Q(s, a_{|A|}; w) \end{array}$$

- Type of approx:
- linear combination of feature.
 - NN (nonlinear fun approx)
 - Fourier / wavelet bases.
- linear approx*

Q: How to determine the parameter ω ?

1. Recall SGD:

$$\min_{\omega} J(\omega)$$

$$w_{k+1} = w_k - \alpha \nabla_w J(w_k)$$

$$\min_{\omega} \mathbb{E} J(\omega, x)$$

$$x \sim p(x)$$

$$w_{k+1} = w_k - \alpha \nabla_w J(w_k, x_k) \in \text{stochastic GD}.$$

2. Imagine in policy prediction case, we know $V^\pi(s)$

$$\min_{\omega} J(\omega) = \mathbb{E}_{s} \left[\frac{1}{2} (V^\pi(s) - \hat{V}(s; \omega))^2 \right]$$

$$GD: \quad w_{k+1} = w_k + \alpha_k \mathbb{E} \left[(V^\pi(s) - \hat{V}(s; w_k)) \nabla_w \hat{V}(s; w_k) \right]$$

$$SGD: \quad w_{k+1} = w_k + \alpha_k (V^\pi(s_t) - \hat{V}(s_t; w_k)) \nabla_w \hat{V}(s_t; w_k)$$

e.g. $\hat{V}(s; \omega)$ is represented in $\{x_j(s)\}_{j=1}^n$

$$\hat{V}(s; \omega) = \sum_{j=1}^n w_j x_j(s) = (x_1(s), \dots, x_n(s)) \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix} = (\pi(s))^T w$$

$\downarrow \mathbb{R}^n$

$$GD: w_{k+1} = w_k + \gamma_k \mathbb{E}[(\hat{V}^\pi(s) - \pi(s)^T w_k) \pi(s)]$$

$$SGD: w_{k+1} = w_k + \gamma_k (\hat{V}^\pi(s_t) - \pi(s_t)^T w_k) \pi(s_t)$$

e.g. Tabular is a special case of linear approx:

with $x_i(s) = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$:- component

$$\hat{V}(s; \omega) = \sum_{i=1}^{|S|} x_i(s) w_i$$

3. However, in RL, we don't really know the true $V^\pi(s)$.

i). We can substitute $V^\pi(s)$ by an unbiased estimate $G_t^{(\infty)}$

$$G_t^{(\infty)} = r_t + \gamma r_{t+1} + \dots + \gamma^T r_{t+\tau}$$

MC learning: $w_{k+1} = w_k + \gamma_k \left(G_t^{(\infty)} - \hat{V}(s_t, w_k) \right) \nabla_w \hat{V}(s_t, w_k)$

\downarrow

$\boxed{\mathbb{E}[G_t^{(\infty)}] = V^\pi(s_t)}$

↳ an unbiased estimate for the true $\nabla J(w)$.

2). TD learning: substitute $G_t^{(1)} = r_t + \gamma \hat{V}(s_{t+1}; w_k)$.

$$w_{k+1} = w_k + \alpha_k \left(G_t^{(1)} - \hat{V}(s_t; w_k) \right) \nabla_w \hat{V}(s_t, w_k)$$

* Note: here $G_t^{(1)}$ is not necessarily an unbiased estimate for $V(s)$.

3). TD(λ): $w_{k+1} = w_k + \alpha_k \left(G_t^\lambda - \hat{V}(s_t; w_k) \right) \nabla_w \hat{V}(s_t, w_k)$

$$G_t^\lambda = (1-\lambda) \sum_{i=0}^{\infty} \lambda^i G_t^{(i)}$$

e.g. linear approx:

$$\underline{w_{k+1} = w_k + \alpha_k \left(r_t + \gamma x(s_{t+1})^T w_k - x(s_t)^T w_k \right) x(s_t)}$$

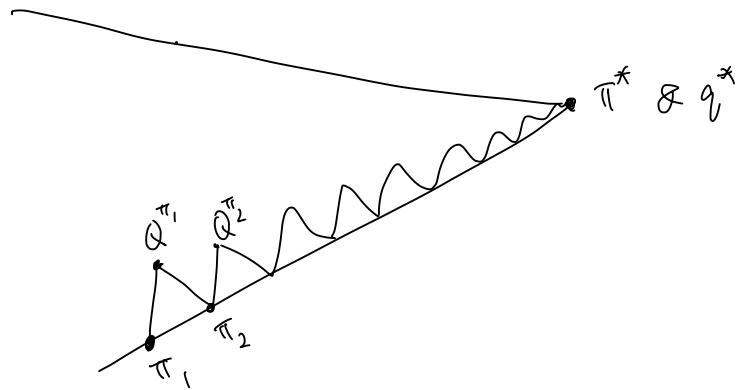
$$\begin{aligned} w_{k+1} &= w_k + \alpha_k \underbrace{d_t}_{\downarrow} x(s_t) \\ &\text{TD: } r_t + \gamma \hat{V}(s_{t+1}; w_k) - \hat{V}(s_t; w_k) \end{aligned}$$

"Convergence Thm" MC \rightarrow converges to a local min for both linear & nonlinear approx

TD \rightarrow converge (close to) a local min
 $\curvearrowleft \gamma$
 for linear approx.

TD(λ) \rightarrow similar as TD.

function approximation for control.



Algo : update $q(s, a; \omega)$ with one step update + ϵ -greedy policy improvement

(SARSA in fcn Approx) .

Initialize s, a, ω

Observe s', r

$$\omega \leftarrow \omega + \alpha (r + \gamma \hat{q}(s', a'; \omega) - \hat{q}(s, a; \omega)) \Delta \hat{q}(s, a; \omega)$$

④ Gradient TD.

- Prediction :

$$V(s) = r(s) + \gamma \mathbb{E}[V(s_{t+1}) | s_t = s]$$

$$\underset{\Delta}{\text{Bellman}} \quad V(s) = r(s) + \gamma \mathbb{E}[V(s_{t+1}) | s_t = s]$$

Bellman

$$\underset{\Delta}{\text{TD}} V = r + \gamma P^\pi V$$

\hat{L} has contraction property .

$V^\pi(s) = \underset{\Delta}{\text{TD}} V^\pi(s) \leftarrow \text{View } V^\pi(s) \text{ as the } \underline{\text{fixed point}} \text{ of the operator}.$

- Instead of view the solution $V^*(s)$ as the fixed pt of T . we can view it as the minimizer of $\int (V(s) - \bar{V}V(s))^2 ds$

- When $V_\theta(s)$ is parametrized by $\theta \in \mathbb{R}^d$

$$\min_{\theta} \mathbb{E}_s \left(r(s) - \gamma \mathbb{E}[V_\theta(s_{t+1}) | s_t = s] - V_\theta(s) \right)^2$$

- One can solve it by GD.

$$\theta_{k+1} = \theta_k - \beta_k \underbrace{\mathbb{E}_s \left[r(s) - \gamma \mathbb{E}[V_\theta(s_{t+1}) | s_t = s] - V_\theta(s) \right]}_{(*)} \underbrace{\left[-\gamma \mathbb{E}[\nabla_\theta V_\theta(s_{t+1}) | s_t = s] - \nabla_\theta V_\theta(s) \right]}_{(*)}$$

- SGD:

unbiased estimate for $(*)$

$$\begin{aligned} & \left(r(s_t) - \gamma \mathbb{E}[V_\theta(s_{t+1}) | s_t] - V_\theta(s_t) \right) \left(-\gamma \mathbb{E}[\nabla_\theta V(s_{t+1}) | s_t] - \nabla_\theta V(s_t) \right) \\ &= \left[r(s_t) - V_\theta(s_t) \right] \left(-\gamma \mathbb{E}[\nabla_\theta V(s_{t+1}) | s_t] \right) - \gamma \mathbb{E}[V_\theta(s_{t+1}) | s_t] \left(-\nabla_\theta V(s_t) \right) \\ &\quad + \gamma^2 \mathbb{E}[V_\theta(s_{t+1}) | s_t] \mathbb{E}[\nabla_\theta V(s_{t+1}) | s_t] \end{aligned}$$

$$\begin{aligned} & \text{unbiased } \left(r(s_t) - V_\theta(s_t) \right) \left(-\gamma \nabla_\theta V(s_{t+1}) \right) - \gamma V_\theta(s_{t+1}) \left(-\nabla_\theta V(s_t) \right) \\ &+ \gamma^2 V_\theta(s_{t+1}) \nabla_\theta V(s_{t+1}') \end{aligned}$$

where s_{t+1} & s'_{t+1} are two independent sample starting from s_t .

① when the underlying dynamics is deterministic.

$$S_{t+1} = S_t'$$

② When the underlying dynamics is stochastic, when we only given one trajectory data,

$\{S_t\}_{t=0}^T$, then it is hard to find two samples from S_t

"Double 'sampling' problem"

One way to avoid "DS"

$$\min_{\theta} \mathbb{E}_s \left(r(s) - \gamma \mathbb{E}[V_\theta(S_{t+1}) | S_t = s] - V_\theta(s) \right)^2$$

$$\min_{\theta} f^2(\theta)$$

$$\approx \min_{\theta} \max_y f(\theta) y - \frac{1}{2} y^2$$

$$\min_{\theta} \max_y \mathbb{E}_s \left(r(s) - \gamma \mathbb{E}[V_\theta(S_{t+1}) | S_t = s] - V_\theta(s) \right) y - \frac{1}{2} y^2$$

given fixed y :

$$\theta_{k+1} = \theta_k - \beta_k \left[(-\gamma \mathbb{E}[\nabla V_\theta(S_{t+1}) | S_t = s] - \nabla_\theta V_\theta(s)) y - \frac{1}{2} y^2 \right]$$

$$\left\{ \begin{array}{l} \theta_{k+1} = \theta_k - \beta_k \left[(-\gamma \nabla V_\theta(S_{t+1}) - \nabla_\theta V_\theta(s_t))y - \frac{1}{2} y^2 \right] \\ y_{k+1} = y_k + \beta_k (r(s) - \gamma V_\theta(S_{t+1}) - V_\theta(s_t))y_k - \frac{1}{2} y_k^2 \end{array} \right.$$

Another way: Borrowing from the future.

$$S_t, S_{t+1}, S_{t+2}, S_{t+3}, \dots$$

$$\hat{S}'_{t+1} = S_t + (S_{t+2} - S_{t+1})$$

\downarrow
This could be a good approximation when
the transition dynamics is smooth.

$$dS_t = \mu(S_t, a_t)dt + \sigma(S_t, a_t)dB_t$$

$$\|\nabla \mu\|, \|\nabla \sigma\| \leq L$$

$$\theta_{k+1} = \theta_k - \beta_k \left[r(s_t) - \gamma V_\theta(S_{t+1}) - V_\theta(s_t) \right] \left(-\gamma \nabla_\theta V(S_t + S_{t+2} - S_{t+1}) - \nabla_\theta V(s_t) \right)$$

- Control:

$$\min_{\theta} \mathbb{E}_{s,a} \left(r(s,a) - \gamma \max_a \mathbb{E} \left[Q_\theta(S_{t+1}, a) \mid S_t = s, a_t = a \right] - Q_\theta(s, a) \right)^2$$

Convergence for few APPROX.

(Prediction)	Tabular	(near)	non-(near)
MC	✓	✓	✓
TD(0)	✓	✓	✗
TD(λ)	✓	✓	✗
DTD	✓	✓	✓
(control)			
MC	✓	✓	✗
SARSA	✓	✓	✗
Q-learning	✓	✗	✗
DTD	✓	✓	✗

- Disadvantage of Value-based method.

$$\max_a Q_\theta(s, a)$$

could be hard to solve
when $|A| \gg |S|$ or its action
space

- Policy ∇ :

Value-based Algo

- Learn Value fn
- Implicit policy based on value fn

Policy-based

- No value fn
- Learn policy