

(s, a, r, s', a')

$Q_{t+1}^{\pi}(s, a) = Q_t^{\pi}(s, a) + \alpha_t (r(s, a) + \gamma Q_t^{\pi}(s', a') - Q_t^{\pi}(s, a))$

TD - learning SARSA → learn the fixed policy

$Q_{t+1}^{\pi^*}(s, a) = Q_t^{\pi^*}(s, a) + \alpha_t (r(s, a) + \gamma Q_t^{\pi^*}(s', a^*) - Q_t^{\pi^*}(s, a))$

$\hookrightarrow \underset{a}{\operatorname{argmax}} Q_t(s', a)$

Q-learning → learn optimal policy

$Q_{t+1}(s, a) = Q_t(s, a) + \alpha_t (r(s, a) + \gamma \sum_{a^*} \pi(a^* | s') Q(s', a^*) - Q_t(s, a))$

expectation Q-learning.

Q How to construct an alternate to exploring "at random"?

- "reward" exploring new states
- distribution should depend on current

$$Q(s, a) \rightarrow \frac{e^{n Q(s, a)}}{\sum_{a'} e^{n Q(s, a')}} =: \pi(a | s)$$

	a
1	2
s	1 - 0.1 0.1

$$\rightarrow \frac{e^{-0.1}}{e^{-0.1} + e^{0.1}}$$

$$\frac{e^{0.1}}{e^{-0.1} + e^{0.1}}$$

$n=1$

	a
2	~ 2 ~ 2
s	2 - 0.2 0.2

$$\rightarrow \frac{e^{0.3}}{e^{-0.2} + e^{0.2}}$$

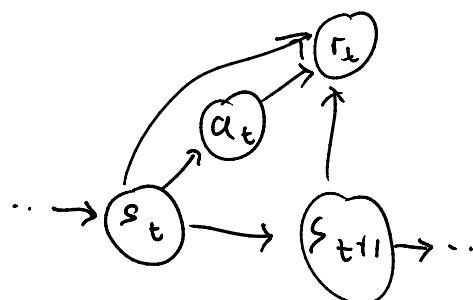
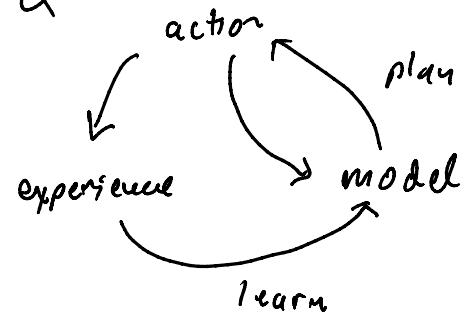
$$2 \quad 0.3 \quad 0.2 \quad \rightarrow \quad \frac{e^{0.3}}{e^{0.3} + e^{0.2}}$$

$$\frac{e^{0.2}}{e^{0.3} + e^{0.2}}$$

Experience Replay Buffer & Dual-Q

$n=1000$

$$\left[\begin{array}{cccccc} s_0, a_0, r_0, s_1, a_1 \\ s_1, a_1, r_1, s_2, a_2 \\ \vdots \\ s_n, a_n, r_n, s_{n+1}, a_{n+1} \end{array} \right]$$



$$\text{IPCL } s_{t+1} = s', r_t = r \mid s_t = s, a_t = a) \\ = B_{s,a}(s', r)$$

MLE

Q-learning has the danger of overestimating the value of certain actions from certain states.

→ double Q-learning

$$Q_{t+1}^{(1)}(s_t, a_t) = Q_t^{(1)}(s_t, a_t)$$

$$+ \alpha_t [r(s_t, a_t) + \gamma Q_t^{(2)}(s_{t+1}, a^*) - Q_t^{(1)}(s_t, a_t)]$$

↳ argmax_{a'} Q_t^{(1)}(s_{t+1}, a')

similar for Q^{(2)} (1) ↔ (2)

randomly choose to update $Q^{(1)}$ or $Q^{(2)}$

agnostic

model-free RL

- no model of the world
- learn value / policy from experience (s, a, r, s')

model-based RL

- learn model from experience
- plan / optimise using sims from model

Dyna RL

- learn model from experience
- learn & plan from simulations & experience
- treat experience & model sims as equiv.

Linear models

$$V^\pi(s) \sim V_\theta = \theta' g$$

↑
learn these parameters

fixed

$$\begin{pmatrix} g_1(s) \\ g_2(s) \\ \vdots \\ g_d(s) \end{pmatrix}$$

$$\underline{V^\pi(s) \sim V_\theta(s)}$$

minimise Loss associated with L^2 error of the approximation

$$\ell(\theta) := \frac{1}{2} \mathbb{E}^{P_n} [(V^\pi(s) - V_\theta(s))^2]$$

$$\nabla_\theta \ell(\theta) = \mathbb{E}^{P_n} [(V^\pi(s) - V_\theta(s)) \nabla_\theta V_\theta(s)]$$

?

$$V^\pi(s) := \mathbb{E}^{P_n} [r_1 + \gamma r_2 + \gamma^2 r_3 + \dots | s=s]$$

$$= \mathbb{E}^{P_n} [r_1 + \gamma \underbrace{V^\pi(s')}_{?} | s=s]$$

$$\theta \leftarrow \theta - \gamma \nabla_{\theta} l(\theta) \rightarrow \text{gradient descent}$$

$\sim (r + \gamma V_{\theta}(s') - V_{\theta}(s)) P_{\theta} V_{\theta}(s)$

$$\boxed{\theta_{t+1} = \theta_t - \eta_t (r_t + \gamma V_{\theta}(s_{t+1}) - V_{\theta}(s_t)) \nabla_{\theta} V_{\theta}(s_t)}$$

$$\eta_1 > \eta_2 > \dots > 0 \quad \lim_{t \rightarrow \infty} \eta_t = 0 \quad \sum \eta_t = +\infty$$

$$\sum \eta_t^L < +\infty$$

$$\eta_t = \frac{c}{D+t}$$

$$l(\theta) = \frac{1}{T} \mathbb{E} [(Q(s, a) - Q_{\theta}(s, a))^2]$$

$$\nabla_{\theta} l(\theta) = \mathbb{E} [(\underline{Q}(s, a) - \underline{Q}_{\theta}(s, a)) \nabla_{\theta} \underline{Q}_{\theta}(s, a)]$$

$$\boxed{\theta_{t+1} = \theta_t - \eta_t [((r + \gamma Q_{\theta}(s', a')) - Q_{\theta}(s, a)) \nabla_{\theta} Q_{\theta}(s, a)]}$$

Function approximation associated with SARSA

$$\boxed{\theta_{t+1} = \theta_t - \eta_t [(r + \gamma Q_{\theta}(s', a^*)) - Q_{\theta}(s, a)) \nabla_{\theta} Q_{\theta}(s, a)]}$$

$\hookrightarrow \underset{a'}{\operatorname{argmax}} Q_{\theta}(s', a')$

FA for Q-learning

Double Q-learning with FA

$$\alpha = \alpha - \eta \Gamma (1 + \gamma \max_{a' \in A(s')} \rightarrow \underset{a' \in A(s')}{\operatorname{argmax}} Q_{\theta}(s', a'))$$

$$\Theta_{t+1} = \Theta_t - \eta_t \left[(r + \gamma Q_{\tilde{\Theta}}(s', a^*) - Q_{\Theta}(s, a)) \nabla_{\Theta} Q_{\Theta}(s, a) \right]$$

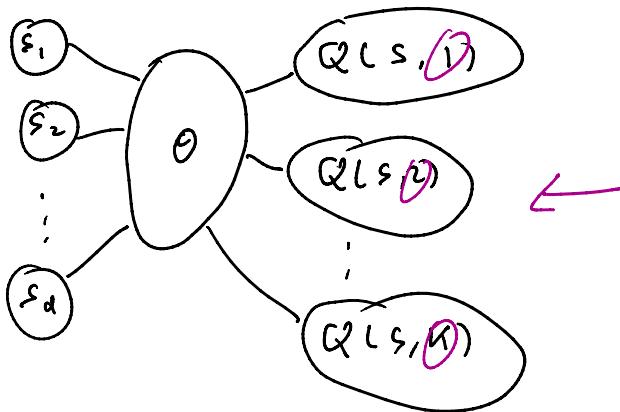
$\rightarrow \underset{a' \in A(s')}{\text{argmax}} Q_{\Theta}(s', a')$

$$\tilde{\Theta}_{t+1} = \tilde{\Theta}_t - \eta_t \left[(r + \gamma Q_{\tilde{\Theta}}(s', a^*)) - Q_{\tilde{\Theta}}(s, a) \right] \nabla_{\tilde{\Theta}} Q_{\tilde{\Theta}}(s, a)$$

$\hookrightarrow \underset{a' \in A(s)}{\text{argmax}} Q_{\tilde{\Theta}}(s', a')$

neural-net approximations

Deep Q-Learning



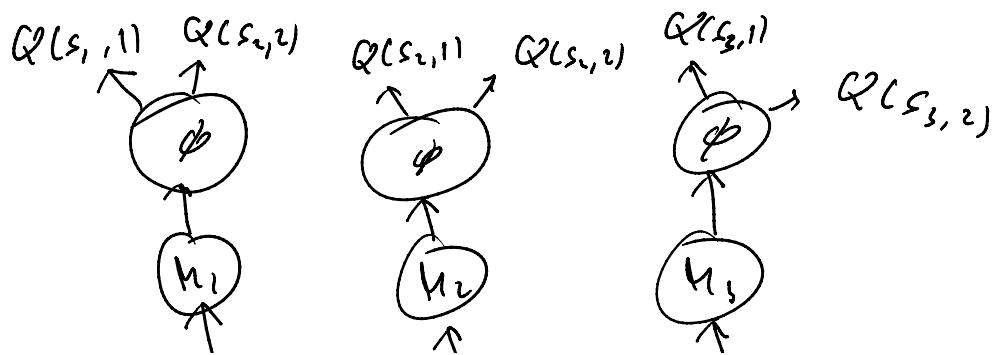
define a modified loss

$$L(\Theta) := \mathbb{E} \left[(r + \gamma Q_{\tilde{\Theta}}(s', a^*) - Q_{\Theta}(s, a))^2 \right]$$

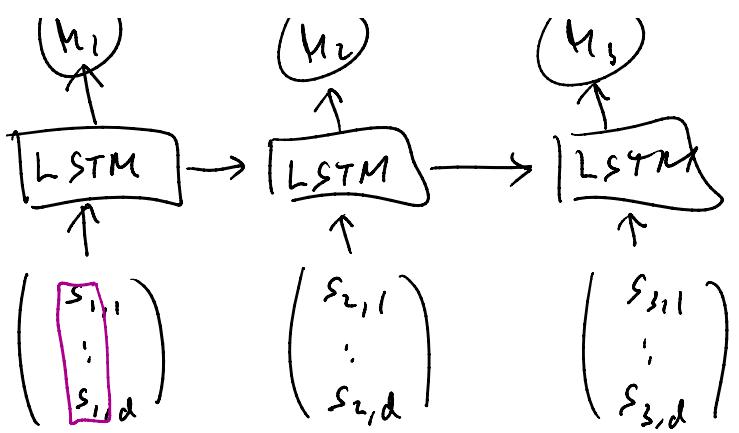
$\rightarrow \underset{a' \in A(s')}{\text{argmax}} Q_{\Theta}(s', a')$

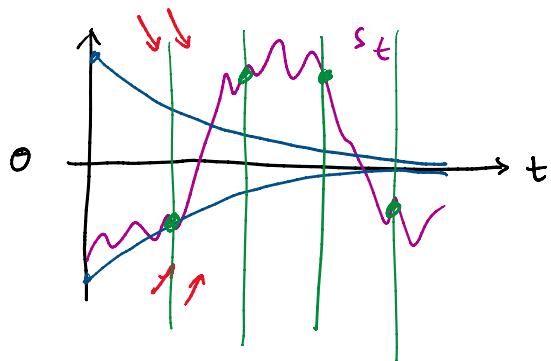
$$\approx \frac{1}{N} \sum_{n=1}^N (r^{(n)} + \gamma Q_{\tilde{\Theta}}(s'^{(n)}, a'^{(n)}) - Q_{\Theta}(s^{(n)}, a^{(n)}))^2$$

$\nabla_{\Theta} L(\Theta) \rightarrow$ use backprop to compute



$$\Theta = \{\phi, \mathcal{L}(\Theta)\}$$





Ornstein-Uhlenbeck (OU)

$$S = (S_t)_{t \geq 0}$$

$$dS_t = \kappa(\theta - S_t)dt + \sigma dW_t$$

$W = (W_t)_{t \geq 0}$ a Brownian motion

$$\dot{S}_t = \kappa(\theta - S_t)$$

$$S_t = \theta + (S_0 - \theta)e^{-\kappa t}$$

$$S_t = \theta + (S_0 - \theta)e^{-\kappa t} + \sigma \int_0^t e^{-\kappa(t-u)} dW_u$$

$$S_{t+\Delta t} \Big|_{\mathcal{F}_t} = \underbrace{\left(\theta + (S_t - \theta) e^{-\kappa \Delta t} \right)}_{\text{mean}} + \underbrace{\sigma \int_t^{t+\Delta t} e^{-\kappa(t+\Delta t-u)} dW_u}_{\text{error}}$$

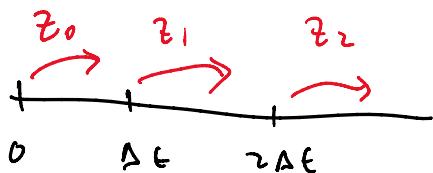
~N(0; \int_t^{t+\Delta t} e^{-2\kappa(u-t)} du)

$$= \alpha + \beta S_t + \tilde{\sigma} Z_t$$

$\tilde{\sigma} = \sqrt{\frac{1-e^{-2\kappa\Delta t}}{2\kappa}}$

$$Z_t \sim N(0, 1) \text{ iid}$$

AR(1)



state = (S_t, I_t) if finite horizon,
include t as a dimension of state variable

$$I_t \in \{+1, 0, -1\} := A$$

$$a_t \in A$$



$$r_t = a_t(S_{t+\Delta t} - S_t) - \lambda |I_t - a_t| S_t$$

S_t

I_t

t

$t + \Delta t$

$$r_t = a_t(S_{t+\Delta t} - S_t) - \lambda |I_t - a_t| S_t$$
$$S_{t+\Delta t} = \alpha + \beta S_t + \tilde{\sigma} Z_t$$
$$I_{t+\Delta t} = a_t$$