

FMDPs

$p(s', r | s, a)$ - opis stochastyczny
 $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$

Finale nutowyca stany $v_{\pi}(s) = \mathbb{E}[G_t | s_t = s]$
 Finale nutowyca akcje $q_{\pi}(s, a) = \mathbb{E}[G_t | s_t = s, a_t = a]$
 Równanie Bellmana i Bellman Optimality Principle

$$v_{\pi}(s) = \sum_a q_{\pi}(s, a) \pi(a | s)$$

$$q_{\pi}(s, a) = \sum_{s', r} p(s', r | s, a) [\tau + \gamma v_{\pi}(s')]$$

NOTE: Zastąpić τ to u. znowu uogólnienie dla $v_{\pi}(s)$ lub $q_{\pi}(s, a)$.

$$v^*(s) = \max_{\pi} v_{\pi}(s) = \max_a q^*(s, a)$$

$$q^*(s, a) = \max_{\pi} q_{\pi}(s, a) = \sum_{s', r} p(s', r | s, a) [\tau + \gamma v^*(s')]$$

Iterative Policy Evaluation

1. Initializacja $v_{\pi}^{(0)}(s)$ s.t. $v_{\pi}^{(0)}(terminal) = 0$
2. $v_{\pi}^{(k+1)}(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [\tau + \gamma v_{\pi}^{(k)}(s')]$

Policy Iteration

1. Initializacja $\pi_0(s)$ (deterministic policy)
2. Oblicz $v_{\pi_k}(s)$ konvergencja z IPE
3. Popraw $\pi_{k+1}(s) = \operatorname{argmax}_a q_{\pi_k}(s, a)$

Value Iteration

1. Initializacja $v_0(s)$ s.t. $v_0(terminal) = 0$
2. $v_{k+1}(s) = \max_a \sum_{s', r} p(s', r | s, a) [\tau + \gamma v_k(s')]$
3. Znowu politykę wybieramy.

Bandyty

Finale nutowyca akcje
 $N_t(a) = \sum_{i=1}^{t-1} [a_i = a]$
 $Q_t(a) = \frac{1}{N_t(a)} \sum_{i=1}^{t-1} r_i [a_i = a] \Rightarrow Q_{t+1}(a) = Q_t(a) + \frac{[a_t = a] - Q_t(a)}{N_t(a)}$

Exponential recency-weighted average

$$Q_{t+1}(a) = Q_t(a) + \alpha [a_t = a] (r_t - Q_t(a))$$

- ϵ -greedy $a_t = \operatorname{rand}(a)$ if $\operatorname{rand}() < \epsilon$ else $\operatorname{argmax}_a Q_t(a)$
- UCB $a_t = \operatorname{argmax}_a \{ Q_t(a) + c \sqrt{\frac{\ln t}{N_t(a)}} \}$
- Gradient approach $\pi(a) = \frac{\exp w_a}{\sum_{b \in A} \exp w_b}$
 $w_{t+1}(a) = w_t(a) + \alpha (r_t - \bar{r}_t) [a_t = a] - \pi_t(a)$

Thompson sampling

$$p_t(\theta_a | r_t) \propto p(r_t | \theta_a) p_t(\theta_a)$$

$$p_{t+1}(\theta_a) = p_t(\theta_a | r_t)$$

$$a_t = \operatorname{argmax}_a \mathbb{E}[r | \theta_a], \quad \theta_a \sim p_t(\theta_a)$$

Bernoulli bandits

$$p(r | \theta_a) \equiv \operatorname{Ber}(\theta_a)(r) = \theta_a^r (1 - \theta_a)^{1-r}$$

$$p_1(\theta_a) \equiv \operatorname{Beta}(\alpha_1, \beta_1) \propto \theta_a^{\alpha_1-1} \cdot (1 - \theta_a)^{\beta_1-1}$$

$$p_2(\theta_a) \equiv \theta_a^{\alpha_1} \cdot (1 - \theta_a)^{\beta_1-1} \cdot \theta_a^{\alpha_1-1} \cdot (1 - \theta_a)^{\beta_1-1} = \operatorname{Beta}(\alpha_1 + r_1, \beta_1 + 1 - r_1) = \operatorname{Beta}(\alpha_2, \beta_2)$$

MONTÉ CARLO

on-policy - oceniamy i ulepszamy tę samą politykę wg. której wykonujemy akcje

Importance sampling

Jeżeli wykonujemy off-policy to mamy 2 polityki:
 π - optymalna, b - behawioralna / eksploracyjna
 Ponieważ pseudosobowzowna myślenie dotyczy epizodu, a nie dla π i b więc musimy to naprawić

$$\alpha \leftarrow \left(\prod_{k=t}^{T-1} \frac{\pi(A_k | S_k)}{b(A_k | S_k)} \right) \cdot \alpha$$

Q-Learning niegdy optymalna polityka, ale w same treningu konflikt z polityką behawioralną (która może spadać z kitem). Jest więc lepszy wybór, gdyż mamy 2 tryby: eksploracyjny i oceniający i eksploracyjny nie są wzajemnie.

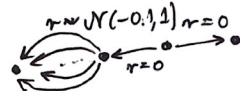
Expected SARSA

$$Q(s, A) \leftarrow Q(s, A) + \alpha [R + \gamma \sum_a \pi(a | s') Q(s', a) - Q(s, A)]$$

(Jest drugi człon obliczamy)

Maximisation Bias

Q-Learning (i trochę Expected SARSA) myślenie punktowe maksymalne nadmierne optymistyczne bo bierze max.



Dla tego punktowa Q-Learning będzie dłużej szukać w ten sposób maksymalnego rozwiązania

w maksymalnym momencie może być dwie (mamy nadzieję o dwóch wariantach). Algorytm będzie powoli dodawał wartości. Znowu max. bo jest to, że oceniamy nasz wybór na podstawie tego, co było do tego wyboru myśle (wybieramy ale na podstawie Q , a potem oceniamy ją Q)

Double Q-Learning

Mamy dwie funkcje Q_1, Q_2 , które aktualizujemy równo w każdym momencie

$$Q_1(s, A) \leftarrow Q_1(s, A) + \alpha [R + \gamma Q_2(s', \operatorname{argmax}_a Q_2(s', a)) - Q_1(s, A)]$$

TD

$$\delta_t = R_{t+1} + \gamma v(S_{t+1}) - v(S_t)$$

SARSA (on-policy TD(0) without)

For each episode:

Take action A using $Q(s, \cdot)$ (e.g. ϵ -greedy)
 Take action A' using $Q(s', \cdot)$ (e.g. ϵ -greedy)
 $Q(s, A) \leftarrow Q(s, A) + \alpha [R + \gamma Q(s', A') - Q(s, A)]$

Q-Learning (off-policy TD(0) without)

For each episode:

Take action A using $Q(s, \cdot)$ (e.g. ϵ -greedy)
 $Q(s, A) \leftarrow Q(s, A) + \alpha [R + \gamma \max_a Q(s', a) - Q(s, A)]$

SARSA znajduje optymalną politykę wśród polityk eksploracyjnych. Jest więc lepszym wyborem, gdyż mamy właśnie eksplorację (punktowa z kitem: SARSA znajduje siebie nieco gorzej, ale inaczej lepiej).

semi-quantitative SPRSA

60 2 drugo mesto treba imeti vsaj eno

2 polny masy polnych opytnych, a 2 osoby polnych, ktore same eksperyment.

II. (checking, bootstrap), to a purposeful analysis of the

Deadly Triad

Weeks 9-11: 11/11/19

$$\tilde{G}_t \leq K_{t+1} - K_t - K_{t+2} + K_{t+1} + \dots$$

→ psychische Probleme wie - psychische (z.B.)

$$\begin{aligned}\phi_{t+1} &= \phi_t + \alpha [G_t - \hat{v}(s_t; \phi_t)] \nabla_{\phi} \hat{v}(s_t; \phi_t) \\ \theta_{t+1} &= \theta_t + \alpha [G_t - \hat{v}(s_t; \phi_t)] \nabla_{\theta} \ln \pi(A_t | s_t; \theta_t)\end{aligned}$$
$$\begin{aligned}\phi_{t+1} &= \phi_t + \alpha [R_{t+1} + \gamma \hat{v}(s_{t+1}; \phi_t) - \hat{v}(s_t; \phi_t)] \nabla_{\phi} \hat{v}(s_t; \phi_t) \\ \theta_{t+1} &= \theta_t + \alpha [R_{t+1} + \gamma \hat{v}(s_{t+1}; \phi_t) - \hat{v}(s_t; \phi_t)] \nabla_{\theta} \ln \pi(\theta_t | s_t)\end{aligned}$$

2] into the set "molecules"

$$\theta^{t+1} = \theta^t + \alpha [G^t(\theta^t) - \nabla_{\theta} \ell(\theta^t)]$$

1. Wiederholung (Repetition)
 2. Verknüpfung (Association)
 3. Veranschaulichung (Visualization)
 4. Verknüpfung (Association)
 5. Wiederholung (Repetition)

maður upphaf, dæta nemiðna

explanatory. In order to understand the nature of the process, it is necessary to know the nature of the process.

$$-Q(s, A)$$

in purposeful n-norms spray (conspire)

Dyna-Q: $\boxed{\text{ENV}} \rightarrow \boxed{\text{real model}} \xrightarrow{\text{Re upd}}$ Policy / value

Model

Dyna - Q+ nagroda wykonana na
stanie 100% w 100% w 100%

ostatnego odniedzielną przeważną masę

~~_____~~ | POLICY GRADIENT

probable beipostedno poudlilni opytanij m

$$\theta_{t+1} = \theta_t + \alpha \frac{\partial J}{\partial \theta}, \quad J(\theta) = v_{\pi_\theta}(s_0)$$

$$\left. \begin{array}{l} \text{Th. 0} \\ \text{gracilencie} \\ \text{politychi} \end{array} \right| \quad \frac{\partial}{\partial \theta} \propto \mathbb{E} [G_{\theta} \nabla_{\theta} \ln \pi]$$

REINFORCE (Monte Carlo Policy Gradient)

$$2. \theta_{t+1} = \theta_t + \alpha G_t \nabla_{\theta} \ln \pi(A_t | S_t; \theta_t)$$

REINFORCE w/ baseline
12 dathone moreny ughwysten fth 2 baseline

$$C_{t+1} = C_t + (C_t - C_{t-1}) / \alpha$$