

Ch2.

ex. 与魔咒机

1. 贪婪算法

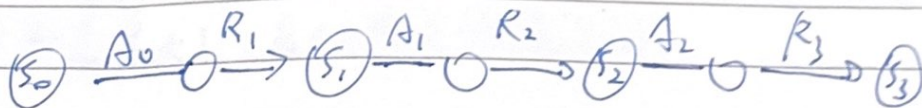
2. 平衡了 exploit - explore

$\epsilon$ -greedy  $\begin{cases} \epsilon & \text{random choice} \\ 1-\epsilon & \text{choose the best} \end{cases}$

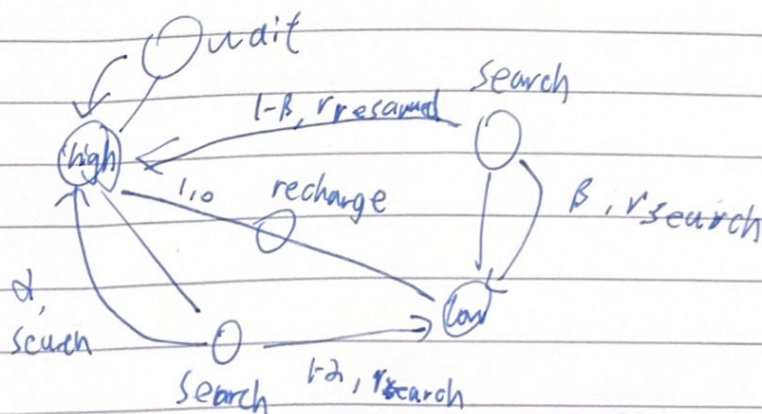
MDP (Markov decision Process)

① should be Markov chain.

ex

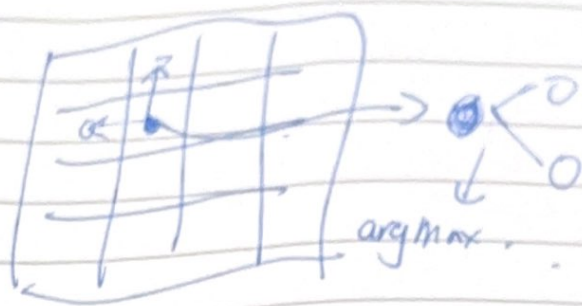


ex



D 12.

best policy = max / all sub state' policy + transfer



缺点: 状态空间指数级增长

优点: 状态空间小时相当快(比直接暴力更快  
很多)

ex. 0-1 背包: 给出  $n$  个有价值的物品与  $k$  元钱, 如何取舍  
使得收益最大! (有价值的物品价值不同)

动态规划方法

动态规划 { DP: 最优子结构以推导全局解  
MC: 大样本以拟合目标函数

on-policy:  $\epsilon$ -greedy,

重要性采样 (off-policy)

How to use  $E_b(G)$  to estimate  $E_{\pi}(G)$ ?

$$P_{\pi}(\{(S_i, A_i)\}_{i=1}^T) = \prod_{i=1}^T \pi(A_i | S_i) p(S_{i+1} | S_i, A_i)$$

$$\therefore p_{t:T-1} = \frac{\prod_{i=t}^T \pi(A_i | S_i)}{\prod_{i=t}^T b(A_i | S_i)}$$

一般重要性采样

$$V_{\pi}(s) \approx \frac{\sum_{k=1}^N p_{\pi}^k(s) : T^k(s) r^k(s)}{N}$$

$$\text{其中, } T^k(s) = \min \{i : S_i^k = s\}$$



加权重复抽样

$$V_n(s) = \frac{\sum_k^n p_k(s) : T^{-1} G_{p^k}(s)}{\sum_k^n p_k(s) : T^{-1}}$$

(有偏估计, 但  $b \rightarrow 0$ )

\* 增量实现:

$$V_n = \frac{\sum_{k=0}^{n-1} W_k C_k}{\sum_{k=0}^{n-1} W_k}$$

$$\rightarrow V_{n+1} = V_n + \frac{W_n}{C_n} (G_n - V_n)$$

$$\text{其中 } C_{n+1} = C_n + W_{n+1}$$

Ch3.

TD 方法 (时间差分法)

27601

$$\left\{ \begin{array}{l} \text{MC: } V(s_t) \leftarrow V(s_t) + \alpha (G_t - V(s_t)) \\ \text{TD: } V(s_t) \leftarrow V(s_t) + \alpha [R_{t+1} + \gamma V(s_{t+1}) - V(s_t)] \end{array} \right.$$

相比于 DP, TD 采用估计值而非真实值去更新模型,  
相比于 MC, TD 采用 one-step 去估计模型.

上述为 TD(0).

Sarsa: on-policy TD

TD(0) for  $V(\text{action})$ :

$$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \alpha [R_{t+1} + \gamma Q(s_{t+1}, A_{t+1}) - Q(s_t, A_t)]$$

Sarsa:  $s_t, A_t, R_{t+1}, s_{t+1}, A_{t+1}$ .

Q-learning : off-policy TD

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

收敛到  $q^*$ . 类似 Sarsa.

Expected Sarsa:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma E[Q(S_{t+1}, A_{t+1}) | S_{t+1}]$$

$$- Q(S_t, A_t)] \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \sum_a \pi(a | S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t)]$$



Ch4.

DQN (Deep Q-Networks)

Q-learning 的深度学习版本

更新规则同 Q-learning

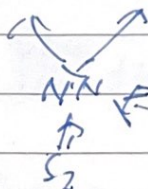
$$\text{loss: } L(\theta) = E[| \text{Target } Q - Q(s, a; \theta) |^2]$$

$$\text{Target } Q = r + \gamma \max_{a'} Q(s', a'; \theta)$$

$2(Q_{\text{real}} - Q_{\text{estimated}})$

$$Q(s') \text{ 观察 } R + \gamma^* \max[Q(s', a_1), Q(s', a_2)]$$

$$Q(s_2, a_1) \quad Q(s_2, a_2) \quad Q(s_2) \text{ 估计}$$



非常手写的 - 一块知识..

Chs. Ch6.

Policy Gradient 方法

直接更新策略网络的 DQN.

$$L(\theta) = E[r_1 + \gamma r_2 + \gamma^2 r_3 + \dots | \pi(\cdot, \theta)].$$

思路: 构造一个好的动作评判指标.

loss:  $L(\theta) = E \log \pi(a|s, \theta) f(s, a)$

$$\nabla_{\theta} E_{\pi}[f(x)] = \nabla_{\theta} \sum_{\pi} p(x) f(x) \quad \text{def.}$$

$$= \sum_{\pi} \nabla_{\theta} p(x) f(x) \quad \text{Swap sum and gradient}$$

$$= \sum_{\pi} p(x) \frac{\nabla_{\theta} p(x)}{p(x)} f(x) \quad \text{both multiply and divide by } p(x)$$

$$= \sum_{\pi} p(x) \nabla_{\theta} \log p(x) f(x) \quad \text{use the fact that } \nabla_{\theta} \log(z) = \frac{1}{z} \nabla_{\theta} z$$

$$= E_{\pi}[f(x) \nabla_{\theta} \log p(x)] \quad \text{def of exp.}$$

off-policy version

$\left\{ \begin{array}{l} \text{Greedy-GQ} \quad \text{过于前向, 没看状态. 略} \\ \text{GTD} \\ \text{DP-PG.} \end{array} \right.$



Ch7.

Entropy: 熵高方法.

$$H(x) \triangleq E[I(x)] = - \sum_{k \in x} p(k) \log_2 p(k)$$

QDRL:

$$\pi^* = \arg \max_{\pi} E_{(s_t, a_t) \sim p_{\pi}} \left[ \sum_t \gamma^t R(s_t, a_t) \right]$$

Max-Entropy:

$$\pi^* = \arg \max_{\pi} E_{(s_t, a_t) \sim p_{\pi}} \left[ \underbrace{\sum_t \gamma^t R(s_t, a_t)}_{\text{reward}} + \underbrace{\alpha H(\pi(\cdot | s_t))}_{\text{Entropy}} \right]$$

核心思想: 不选到任意一个有用的 action, 有用的 trajectory.

SAC  $\left\{ \begin{array}{l} \text{一个策略网络} \\ \text{两个Q网络} \end{array} \right.$

Ch 8.

ex: Coffee-making robot.

challenges:

1. Diversity of problems
2. high-dimensional and multimodal goals
3. Scalability
4. Instability
5. Optimality

Inverse RL (逆向强化学习) aka IRL.

exp / Soft Q Imitation

SQIL

PaJD

GAIL

other paradigms:

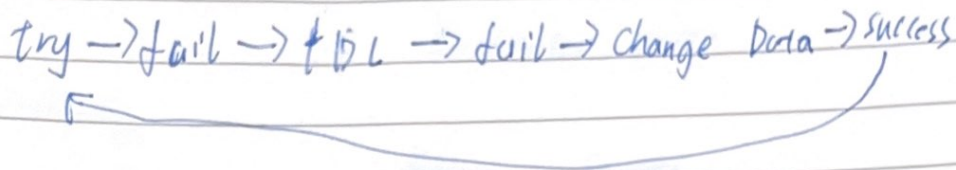
~~Bottom~~ exp / meta-learning

Transfer Learning

⋮

## Ch9. Practical Reinforcement Learning

### ① RL Project Life Cycle.



### ②. RL project definition.

1. Sequential
2. Strategic

本章是一些职业和发展的形而上内容, 不再赘述.



Ch 10. Operational RL.

frame work

Deployment

Abstractions

Log file.

Simplicity

Parameter server

Observability

Policy servers

Scalability

Replay buffers

minimal dependencies

Actors and critics

Documentation and examples

Learning

Included algorithms

Data augmentation.

Ancillary algorithms

Support

Commercial offerings.

Ch 11. Summary.

前言: 区别于 DL 的一种工具, 仅此而已.