## Markov Decision Processes(MDP)

is specified by:
• A set of states S　　　　• A set of actions A
• Initial state distribution p(s0)　• A reward function r(s,a)
• Markov kernel p(s'|s,a) → transition model

### Inverted Pendulum(倒立摆)
A set of states $S=[-1,1]\times[-1,1]\times[-8,8]$ 用正余弦表是为了解决角度的周期性
A set of actions $A=[-2,2]$
Initial state distribution: $p(s_0)=[\cos(U(-\pi,\pi)),\sin(U(-\pi,\pi)),U(-1,1)]$
Transition kernel $p(s'|s,a)=\delta(s'-[\cos(\theta_{n+1}),\sin(\theta_{n+1}),\dot\theta_{n+1}])$

$\theta_{n+1}=\min\left(\max\left(\theta_n+\left(\frac{3g}{2l}\sin(\theta)+\frac{3\tau}{ml^2}\right)\Delta t,-8\right),8\right),\ \ \theta_{n+1}=\theta_n+\dot\theta_n\Delta t$

A reward function: $r(s,a)=-(\theta^2+0.1\dot\theta^2+0.001\tau^2)$
惩罚角度偏离θ(竖直向上)；惩罚角速度过大；惩罚施加过大的力矩
A finite MDP:只有转换矩阵变了$P(S'|s,a)=\mathbb{P}(S_{t+1}=s'|S_t=s,A_t=a)$

### Tic-Tac-Toe
A set of states $S=\{0,1,2\}^9$代表棋盘上每个格子的内容
A set of actions $A=\{0,1,2,3,4,5,6,7,8\}$当前玩家选择下了的格子编号
Initial state distribution $p(S_0)=\delta(\theta_0)$ 所有元素都是 0 的 9 维向量
Transition probability $P(S'_{t+1}=s'|S_t=s,A_t=a)$ 取决于对手的策略
A reward function $r(s,a)$: e.g. 1 if win, 0 otherwise

### Value Functions in MDPs
Fix deterministic policy $\pi:S\to A$ 对于每一个状态，π 都会唯一确定一个动作 $a=\pi(s)$，而不是一个动作的概率分布。
MDP becomes a Markov chain with transition probabilities 动作 a 是确定的不用考虑动作 a，简化为 $P(S_{t+1}=s'|S_t=s)=P(s'|s,\pi(s))$
Value function $V^\pi(s)=\mathbb{E}[\sum_{t=0}^\infty\gamma^t r(S_t,\pi(S_t))|S_0=s]$累计折扣奖励
Recursive Definition 分解：即时奖励和下一状态折后价值期望
贝尔曼期望方程：$V^\pi(s)=r(s,\pi(s))+\gamma\sum_{s'}P(s'|s,\pi(s))V^\pi(s')$
矩阵形式： $V^\pi=r^\pi+\gamma\mathbf{T}^\pi V^\pi$

### Solving for Value Functions 求解 $V^\pi$ 的两种策略
1. Solving Linear System $V^\pi=(I-\gamma\mathbf{T}^\pi)^{-1}r^\pi$ 适用状态有限，数量少
2. Fixed Point Iteration
　Initialize $V_0^\pi$ ; For t=1:T do: $V_t^\pi=r^\pi+\gamma\mathbf{T}^\pi V_{t-1}^\pi$　:= $\mathcal{B}^\pi V_{t-1}^\pi$
　当这代收敛时，即 $V_t^\pi=V_{t-1}^\pi$，找到了 $V^\pi$ 的解，是定点。
　$\|V_t^\pi-V^\pi\|_\infty\le\gamma^t\|V_0^\pi-V^\pi\|_\infty$ 第 t 次迭代得到的 $V_t^\pi$ 与真实值的差收敛

### Picking Actions Given Value Functions
We knew V, could choose: $arg\max_a\left(r(s,a)+\gamma\sum_{s'}P(s'|s,a)V(s')\right)$贪婪策略

### Connecting Value Functions and Policies 价值函数与策略之间的联系
Every policy induces a value function:给定策略，得到$V^\pi$
Every value function induces a policy 给定 V，找到一个更好的策略

### Bellman's Identity 贝尔曼恒等式
$V^*(s)=\max_a\left(r(s,a)+\gamma\sum_{s'}P(s'|s,a)V^*(s')\right)$

### find value function satisfying:
Policy Iteration: 策略$\pi$开始，再用$V^\pi$改进得到$\pi'$,重复直到两者收敛
伪代码: Start with arbitrary policy π
　　　While not converged do
　　　　Compute value function $V^\pi$
　　　　Compute greedy policy $\pi_g$ w.r.t. $V^\pi$
　　　　Set $\pi\leftarrow\pi_g$
Convergence to optimal policy in $O(|S|^2|A|/1-\gamma)$ iterations

### Value Iteration: 只有价值更新，从$V_0$反复迭代，直到 $V_t$ 收敛对 $V^*$
$V_0(s)=\max_a r(s,a),\ V_{t+1}(s)=\max_a\left(r(s,a)+\gamma\sum_{s'}P(s'|s,a)V_t(s')\right)$
每次迭代只用，用上一步的价值函数$V_t$来计算更优的估计 $V_{t+1}$
伪代码: Initialize $V_0(s)=\max_a r(s,a)$
　　　Do
　　　　Define $Q_t(s,a)=r(s,a)+\gamma\sum_{s'}P(s'|s,a)V_{t-1}(s')$
　　　　Find$V_t(s)=\max_a Q_t(s,a)$
　　　While $\|V_t-V_{t-1}\|_\infty=\max_s|V_t(s)-V_{t-1}(s)|\le\epsilon$
　　　　Compute greedy policy $\pi_g$ w.r.t. $V_t$
Guaranteed to converge to $\epsilon$-optimal policy!

### Value vs. Policy Iteration

| | Policy | Value |
|---|---|---|
| Converges to | exact solution | $\epsilon$-optimal solution |
| complexity per iteration | $O(|S|^3)$ | $O(|S|^2|A|)$ |

## The Multi-Armed Bandit Problem
衡量指标: cumulative regret $\rho_T=\max_{t=1}^T\mathbb{E}[r_a]-\sum_{t=1}^T\mathbb{E}[r_a]$
在 T 个回合中，始终选择最佳臂获得的预期奖励总奖励，与实际获得的预期总奖励差值，目标是设计一个策略，使累积遗憾$\rho_T$最小化
Uniform Exploration: $\rho_T\ge T\left(\frac{K-1}{K}\right)\min_a\mathbb{E}[r^*-r_a]$

### Explore-First Algorithm:
Idea: approximate expectation empirically:对每个臂的真实 r 进行经验近似
1. Try each arm $N_a$ times: 估计$\mathbb{E}[r_a]\approx\frac{1}{N_a}\sum_{t=1}^{N_a}R_t/N_a=\bar\mu(a)$经验平均
2. Play arm with highest average for the remaining rounds
Hoeffding inequality 霍夫丁不等式
Assume $R_t\in[0,1]$. Then,
高概率保证，经验估计与真实值误差不超过$\epsilon$。$P\left(|\mu(a)-\bar\mu(a)|\le\sqrt{\frac{2\log(T)}{N_a}}\right)\ge1-\frac{2}{T^4}$
$N_a$太大，虽然估计准确，但会浪费大量的回合，从而增加累积遗憾$\rho_T$
$N_a$太小，对每个臂的估计不可靠
Theorem: Assume $R_t\in[0,1]$and pick $N_a=(T/K)^{2/3}(\log(T))^{1/3}$.
$\mathbb{E}[\rho_T]\le O(KT^{2/3}(\log(T))^{1/3})$说明这样选择$N_a$,累计遗憾有上界

### Epsilon-Greedy Algorithm　1-ε+ε/N
大多数时间($p=1-\epsilon$)选择当前最好的知识，但在小部分时间($p=\epsilon$)随机探索。
伪代码: For each round t = 1, ..., T do
　　　Toss coin with success probability $\epsilon_t$(以 $\epsilon_t$ 概率探索)
　　　If success then
　　　　Choose arm uniformly at random 均匀探索，收集信息
　　　Else(以 1-$\epsilon_t$ 的概率利用)
　　　　Choose arm with highest average 选当前$\mu(a)$最高的臂

### Incremental Updates: $\mu_{N_a+1}(a)=\mu_{N_a}(a)+(R_{N_a}-\mu_{N_a}(a))/N_a$
Theorem: Assume $R_t\in[0,1]$ and set $\epsilon_t=(K\log(t)/t)^{1/3}$.
$\mathbb{E}[\rho_T]\le O(KT^{2/3}(\log(T))^{1/3})$使用最优$\epsilon_t$时$\mathbb{E}[\rho_T]$的上界与探索优先算法中在最优 $N_a$ 下的上界具有相同的渐进增长率

### Incremental Value Estimation 新估计=旧估计+步长(目标-旧估计)
Exact estimates if $\sum_{\alpha_{N_a}}^\infty=\infty$ and $\sum_{\alpha_{N_a}}^\infty\alpha_{N_a}^2<\infty$ 数据量大估计过更真实
In practice often: $\alpha_{N_a}=c\leftarrow$ exponential recency decay

### Optimism in the Face of Uncertainty 探索: 赋予不充分尝试的较高估计值
Idea: Use Hoeffding inequality to compute upper confidence bound (UCB)
1. For each round t = 1, … , T do
2. Choose arm that maximizes UCB $ucb(a)=\mu(a)+\sqrt{\frac{2\log(T)}{N_a}}$
3. Play arm and update UCB; Theorem:$E[\rho_T]\le O(\sqrt{KT\log(T)})$
Random Policies in Bandits
• Idea: Don't sample uniformly at random=>Boltzmann distribution
$$\pi_t(a)=\mathbb{P}(A_t=a)=e^{H_t(a)}/\sum_{k=1}^K e^{H_k(a)}$$
For each round t = 1, ..., T do
Sample arm $A_t\sim\pi_t$ 根据 Boltzmann 分布得到的概率$\pi_t$(a), 随机采样$A_t$

### Stochastic Gradient Ascent for Policies 随机梯度上升: 更新偏好值
Update $H_t$ using gradient ascent: 希望通过梯度上升来更新偏好值 $H_t(a)$
　　　　以最大化预期的回报 $\mathbb{E}[r_{A_t}]$
$H_{t+1}(a)=H_t(a)+\alpha\frac{\partial\mathbb{E}[r_{A_t}]}{\partial H_t(a)}=H_t(a)+\alpha\mathbb{E}\left[(r_{A_t}-R_t)(\mathbb{1}_{A_t=a}-\pi_t(a))\right]$
Approximate expectation using single sample:没法对 E 的精确值，用单次采样来近似$H_{t+1}(a)=H_t(a)+\alpha(r_{A_t}-R_t)(\mathbb{1}_{A_t=a}-\pi_t(a))$

### Comparison of Approaches UCB 是 最备棒和性能最好的算法

### Reinforcement Learning Setting
learn transition probabilities -- model-based RL
learn value function -- model-free RL

### Exploration-Exploitation in MDPs Similar problem as in bandits:
Always pick the best action(• Get good immediate reward • Can sacrifice long-term return and get stuck in sub-optimal actions 牺牲长期回报)
Always pick random action(• Allows to correctly estimate all probabilities and rewards可能准确估计• May give very poor returns 但是高遗憾)

### Epsilon-Greedy Action Selection in MDPs
The $R_{max}$ Algorithm(model-based, off-policy 可以用观测数据更新模型)
• Idea: optimism in the face of uncertainty
• If no estimate for r(s,a) is known → set it to $R_{max}$ 奖励不知道就设最大
• If no estimate for P(s'|s,a) is known → set P(s*|s,a)=1,r(s*,a)=$R_{max}$, where s* is a 'fairy tale' state 无转移概率, 假设转移到有 $R_{max}$的状态
• Iterate between solving estimated MDP and applying obtained policy long enough 用当前的乐观估计模型运行规划算法（如价值迭代或策略迭代）得到当前的最优策略 by the next $\frac{1}{\epsilon^2}$ steps．即用当前策略更新模型中(s,a)的估计，当被观察多次后，就会从乐观假设中移除
• Greedy optimistic selection leads to implicit exploration!追逐最大奖励的过程中，自然而然地发生对不确定区域的探索。

### How long is 'long enough'
1. How many samples per state-action pair 确保估计准确需要多少样本通过 Hoeffding inequality 来推断需要多少次观测样本 N，
2. How long it takes to 'traverse' the MDP 需要多长时间探索完所有状态
Lemma: Every T time steps, with high probability, either • A near optimal reward is obtained or At least one unknown state-action pair is visited
3. mixing time $T^\pi$ time until the Markov chain is close to steady state distribution 经过足多长时间步后，状态分布会变得接近平稳状态分布
Theory: With probability $1-\delta$, the $R_{max}$ Algorithm will reach an $\epsilon$-optimal policy in a number of time steps that is polynomial in $|S|,|A|,T,\frac{1}{\epsilon}$,and $\log(\frac{1}{\delta})$.

### Challenges with Model-Based RL
Memory requirement: For every s,a,s', we need to store the probability estimate $\hat P(s'|s,a)\in O(|S|^2|A|)$ for dense MDPs. 需要巨大内存
Computation time: need to repeatedly solve the estimated MDPs, e.g. using policy/value iteration.每次更新模型都要重新求解，时间会很长

### Monte-Carlo Estimates of Value Functions(model-free)
Idea: average over returns instead of estimating probabilities
(Return, G): $G_t=R_t+\gamma R_{t+1}+\gamma^2 R_{t+2}+...$; (Value, V)是对这个 G 的期望
Hoeffding inequality to bound estimation error for finite horizon tasks;
Truncation error 截断误差 in infinite horizon tasks decays exponentially;
Value of single state can be estimated independently 不需知道邻近状态 V
伪代码: For each episode n=1,...,N do
　　　roll-out policy π 在环境中用策略π运行一个完整的轨迹
　　　For each state $s\in S$ do:
　　　　G← return after reaching s first time: (first visit MC)
　　　　append G to returns(s):
　　　　V(s)←average(returns(s)):

### Monte-Carlo Estimates of State-Action Value Functions 拓展到 Q(s,a)
Q(s,a)是在 s 采取 a 之后，并永久遵循策略所获得的预期累积折扣回报
探索不足：Requires policy to play every action 需要所有 a 都被尝试到
→ not possible with deterministic policies 确定性策略π时无法得到其他 Q(s,a')
伪代码: 上述的把 s 换(s,a),G 换 q,V 换 Q

### Monte-Carlo Control
Idea: Use random initial state-action pair for exploration 从随机(s,a)开始

Do policy iteration with Monte-Carlo estimate instead of exact value function in terms 价值迭代循环中, 用 MC 估计来代替求解价值函数的精确解
MC Policy Evaluation(π→Q)通循得到每(s,a)的平均回报，得到估计$Q^\pi$
Greedy Policy Improvement(Q→π):根据$Q^\pi$, 在每个状态 s 下选择使 Q(s,a)最大化的动作 a，从而生成一个新的贪婪策略 $\pi$
On-policy method: data can only be used to estimate value function for policy generating the data!

### Monte-Carlo Control without Exploring Starts 不使用探索性开始
Exploring starts are unrealistic in many applications
→ use random policies: ε-greedy, Boltzmann
→ soft policies: every action retains minimal probability 保证任何状态 s 下的 a 都有非零概率被选中的策略: π(a|s)>0 for all s,a
Policy improvement theorem 策略改进定理 also holds for ε-greedy policies 过程:用软性策略产生数据，并用 MC 方法估计$Q^\pi$,再用$Q^\pi$找一个对$Q^\pi$贪婪的软性策略 $\pi'$, 重复直到收敛最优最优的软性策略 $\pi_\epsilon$

### Importance Sampling for Off-Policy RL
Reason: On-policy methods are not data efficient,改进后旧 data 不能 reuse
Idea: use weighted average of returns for policy $\mu$ to estimate value functions for different policy π用权重重来校正策略$\mu$的$G_t$, 像$\pi$产生的回报

$\rho_t=\frac{\prod_{k=1}^{t-1}\pi(A_k|S_k)P(S_{k+1}|S_k,A_k)}{\prod_{k=1}^{t-1}\mu(A_k|S_k)P(S_{k+1}|S_k,A_k)}=\prod_{k=t}^{t-1}\frac{\pi(A_k|S_k)}{\mu(A_k|S_k)}$ $V(s)=\frac{\sum_{n=1}^N\rho_t(n)G_{t(n)}(s)}{N(s)}$

$\rho_{t(n)}$: 第 n 个回合中，状态 s 首次访问时的重要性比率
$G_{t(n)}(s)$: 状态 s 首次访问后获得的回报。N(s): 访问状态 s 的总回合数。

### Value Estimation using Temporal Differences (model-free, on-policy)
Idea 1: approximate expectation using one sample: $V^\pi(s)\approx r(s,\pi(s))$
Idea 2: bootstrap using old estimate of $V^\pi(s)$: $V_{new}^\pi(s)=r(s,\pi(s))+\gamma V_{old}^\pi(s')$
TD(0) Learning $V_{new}^\pi(s)=r(s,\pi(s))+\gamma V_{old}^\pi(s')$ 用当前 V 估计真实$V^\pi(s)$
estimate can exhibit high variance 波动很大
mix old and new estimate:$V_{new}^\pi(s)=(1-\alpha_t)V_{old}^\pi(s)+\alpha_t(r(s,\pi(s))+\gamma V_{old}^\pi(s'))$
$V_{new}^\pi(s)=V_{old}^\pi(s)+\alpha_t(r(s,\pi(s))+\gamma V_{old}^\pi(s')-V_{old}^\pi(s))$
temporal difference error = new estimate – old estimate
伪代码: For each episode n=1,...,N do:
　　　initialize $S_0$:
　　　For each $t=0,...,T-1$ do:
　　　　apply action $A_t$ given by $\pi(S_t)$: 在当前状态下执行动作
　　　　observe $S_{t+1},R_t$: 观察返回奖励和下一状态
　　　　$V(S_t)\leftarrow V(S_t)+\alpha_t(R_t+\gamma V(S_{t+1})-V(S_t))$用观测更新$S_t$价值估计

### SARSA 　 On-Policy TD Control 　(s,a,r(s,a),s',a')
　　　$Q_{new}(s,a)=(1-\alpha_t)Q_{old}(s,a)+\alpha_t(r(s,a)+\gamma Q_{old}(s',a'))$
伪代码: For each episode n=1,...,N do:
　　　initialize $S_0$:
　　　determine policy π based on Q:
　　　For each $t=0,...,T-1$ do:
　　　　apply action $A_t$ for $S_t$ based on policy π:
　　　　observe $S_{t+1},R_t$: 观察返回奖励和下一状态
　　　　choose action $A_{t+1}$ for $S_{t+1}$ based on policy π:
　　　　$Q(S_t,A_t)\leftarrow Q(S_t,A_t)+\alpha_t(R_t+\gamma Q(S_{t+1},A_{t+1})-Q(S_t,A_t))$
收敛到最优 $Q^*$ 条件: $\sum_{t=1}^\infty\alpha_t=\infty$且$\sum_{t=1}^\infty\alpha_t^2<\infty$; all (s,a)被访问无穷多次

### Off-Policy State-Action Value Estimation
SARSA estimate:Q($S_{t+1},A_{t+1}$)当前策略实际选择并执行的动作
Off-policy estimate:$Q^*(S_{t+1},\pi(S_{t+1}))$采取目标策略认为最好的那个动作
Actions do not need to be picked according to π!生成数据的动作 $A_t$
不受π 策略 b 选择，不用与想要评估的目标策略所建议的动作一致

### Q-Learning 直接学习最优价值函数
Idea: use greedy policy in off-policy evaluation: a'是理论上最好的动作
　　　$Q^*(S_t,A_t)\leftarrow Q^*(S_t,A_t)+\alpha_t(R_t+\gamma\max_{a'}Q^*(S_{t+1},a')-Q^*(S_t,A_t))$
伪代码: For each episode n=1,...,N do:
　　　initialize S0:
　　　For each $t=0,...,T-1$ do:
　　　　apply action At for St based on $Q^*$:
　　　　observe $S_{t+1},R_t$:
　　　　$Q^*(S_t,A_t)\leftarrow Q^*(S_t,A_t)+\alpha_t(R_t+\gamma\max_{a'}Q^*(S_{t+1},a')-Q^*(S_t,A_t))$
多项式界限:log|S|(←number of states): 意味 Q̂-Learning 在状态空间较大的问题上具有较好的理论性能。log|A|(←number of actions):
1/ε(←ε←optimal): 与所需的精度成多项式关系。
log(1/δ)(←probability of guarantee): 与所需保证的概率成对数关系。
性质 Memory requirement: need to store $Q^*$ (s,a)对所有 s,a →O(|S||A|)
Computation time: O(|A|) per iteration a'∈$arg\max_a Q^*$(s,a')耗在找最优 a'

### Challenges of Tabular RL
MDP and RL polynomial in |S| and|A| in general 理论复杂度
key challenge: scaling up (Structured domains 结构) (Continuous domains)
→ Use function approximation for value functions 用参数化值估计

### TD-Learning as Stochastic Gradient Descent
Idea: TD-Learning looks like gradient descent update
$V_{new}^\pi(s)=V_{old}^\pi(s)+\alpha_t(r(s,\pi(s))+\gamma V_{old}^\pi(s')-V_{old}^\pi(s))$
'parameterize' value function: $V^\pi(s;\theta)$向$\theta$转移到参数$\theta$
$V^\pi(s;\theta)\leftarrow V^\pi(s;\theta)+\alpha_t(r(s,\pi(s))+\gamma V^\pi(s';\theta^{old})-V^\pi(s;\theta))$
→ stochastic semi-gradient descent 随机半梯度下降不求$\theta^{old}$这部分梯度
→ identical to TD learning if $\theta^{old}=\theta$ 使用固定目标网络来提高稳定性
缺点: Sample-inefficiency of model-free RL
目标函数不正确且不稳定: 大方差来自单样本本期近似

### Q-learning via Stochastic Gradient Descent
loss function for SGD:
$l(\theta;s,a)=\frac{1}{2}\left[(r(s,a)+\gamma E_{s'}[\max_{a'}Q^*(s',a';\theta^{old})]-Q^*(s,a;\theta)\right]$
update:$Q^*(s,a;\theta)\leftarrow Q^*(s,a;\theta)+\alpha_t(r(s,a)+\gamma\max_{a'}Q^*(s',a';\theta^{old})-Q^*(s,a;\theta))$

### Model-free RL with Function Approximation
Idea: SGD does allows using other parameterizations of value functions
Q-learning:$\theta\leftarrow\theta+\alpha_t\delta_B(s,a,s')\nabla_\theta Q^*(s,a)$收敛很慢

### Heuristics 启发式 for Deep Q-Learning
stabilize training by keeping $\theta^{old}$ over multiple SGD iterations 固定参数
maintain a data set of observed transitions
do SGD updates on mini batches of data → experience replay

### Drone:
State: Drone's position 位置; orientation 姿态; velocity; target command
Actions: rotor command 旋翼指令　Initial state distribution p(s0)
Markov transition model p(s'|s,a)
Reward function: Task reward (tracking command): $-\|v_t-v_{target}\|^2$,
Regularization reward (punish undesired jerky motion): $-|acc_t|^2$,
Survival reward (terminate if fail): $-c$ 在任务失败给的一个很大负奖励
Total reward: $r(s,a)=w_{task}\|v_t-v_{target}\|^2-w_{reg}|acc_t|^2$

### Quadrupedal robot:
State: Robot's joints state, base orientation, velocity, target command
Actions: joints command　Initial state distribution p(s0)
Markov transition model p(s'|s,a)
Reward function: Task reward (tracking command): $-\|v_t-v_{target}\|^2$
Regularization reward (punish undesired joints jerky motion): $-|\delta_t|^2$,
Survival reward (terminate if fail): $-c$
Total reward: $r(s,a)=w_{task}\|v_t-v_{target}\|^2-w_{reg}|\delta_t|^2$
Good controller = Maximize objective function

### MDP in GO
State: stones state　Actions: next stone position

### RL in: Go vs. Robot Control
• Observation state: Discrete vs. Continuous
• Action: Discrete vs. Continuous
• Transition model: Table/Matrix vs. Approximated function
• policy: Over finite state vs. Prob. density function Over continuous space

### Monte Carlo Tree Search (MCTS)
• Heuristic search algorithm for decision-making.启发式搜索算法
• Balances exploration and exploitation using simulations.平衡探索和利用
• No domain-specific knowledge required.无需领域知识，大量随机模拟
• Applications: Go, Chess, planning, etc.
1. Selection: Traverse tree using Upper Confidence Bound applied to Trees (UCT) to select promising node.在每个节点选择最有希望的子节点。
2. Expansion: Add a new child node.有未探索的节点，添加新的子节点
3. Simulation: Run rollout to get outcome.评估新子节点的价值
4. Backpropagation: Update values up the tree.沿路径更新访问次数，总 V

### Why Balances Exploration and Exploitation?
$$UCT(s,a)=\underbrace{Q(s,a)}_{\text{利用(Exploitation)}}+c\cdot\underbrace{\sqrt{\frac{\ln N(s)}{N(s,a)}}}_{\text{探索项(Exploration)}}$$
Q(s,a): Estimated value (average reward) 利用项(Exploitation)of taking action a from state s,选择价值高的
N(s): Total number of visits to state s, N(s,a):从状态 s 采取动作 a 的次数

### Advantages:
• Scalable and domain-agnostic. 不用存储各个状态空间而是构建搜索树
• No value function needed.无需先验评估, 用大量随机模拟的真实结果
• Strong performance with enough simulations.仿真次数多时性能强大

### Limitations:
• High computational cost.计算密集型
• Struggles with sparse rewards and rollout policy design. 奖励稀疏，波动

### MCTS in Practice
• Used in: AlphaGo, robotics planning.
• Variants: PUCT (AlphaZero), 嵌套 Nested MCTS, MCTS + RL.
• Rollout policy and tree depth are critical for success.模拟策略和树深度

### It cannot be Applied to Robot Control Problem:
Continuous space is too costly to build such tree structure.

### Model-based RL
循环:generate samples→fit a model to estimate return→improve the policy
Model: estimate p(s'|s,a) Supervised Learning $\min_p\sum_i\|f_\theta(s_i,a_i)-s_i'\|^2$
目标: min 预测的下一状态与实际下一状态之间的误差平方, 优化$\pi_\theta(a|s)$
• Sample efficient (Leverage models)优点: 需要更少的真实环境数据
• More complex to implement(Learn models, then plan based on model)
• Sensitive to model errors(Less robust)劣势: 复杂, 对模型错误敏感

### Model-free RL
Replay Buffer<=Interaction(s,a,r)=>Model-free Agent=>Replay Buffer
Actor 根据当前的策略π向环境发出动作, 环境返回 s'和 r, 这些经验被记录并送入回放缓冲区, 智能体从回放缓冲区中随机抽取小批量数据, 价值函数和执行器的更新完全基于这些数据进行学习更新, 用于指导下一步的行为
• Simple to implement • Robust to model errors
• Lots of samples are required(Without models)

### Model-based RL vs. Model-free RL
Model-based RL: Learns Mode; Uses Planning(MPC); more sample efficient
Model-free RL: Simpler setup; Less sensitive(robust)

### Robot Reward Design
Run to the target position $r_{pos\_track}=|p-p_{target}|^2$ $r_{vel\_track}=|v-v_{target}|^2$
Not work in real robots: • Unsmooth motion, hard to track by hardware
• Too large contact force, easy to damage • Not user-friendly motion
Punish Undesired Motion: (Shake; Base tilts; Too fast joints motion; Large contact force)
设计奖励策略, 要让机器人稳定运动• Friendly for approx.

### Some Common Reward Function
• Quadratic 二次 function $R(s,a)=-(x-x^*)^2$
• Exponential function $R(s,a)=\exp(-\alpha(x-x^*)^2)$

### Discrete Reward 特定的、不连续的事件, 通常是一个固定的正值或负值
Termination Reward-正向奖励(Successfully finish the task; Usually a discrete positive reward)
Truncated Reward-负向奖励(Reset when failure happens; Fail to finish within time limits; Largely deviate from ideal states; discrete negative)

### Why Truncated Reward? 内容很多，但感觉没用
• Avoid learning from outlier states 避免从异常状态中学习• Data collected may not represent meaningful behavior• Strong rewards (positive or negative) in such outliers, it could bias the policy toward unrealistic behaviors
• Prevent instability in value function learning( • Very large or extreme rewards can cause exploding gradients or unstable updates.• Truncating rewards smooths learning by keeping values within a bounded range.)
• Encourage robustness and Generalization 鲁棒泛化(• By ignoring extreme states, the agent learns policies that are effective in the most relevant parts of the state space.• This avoids overfitting to rare, unlikely transitions.)

### Why Sim-to-Real mismatch?(• Actuators performance mismatch to ideal models• Sensors noises• Physical engines inaccuracy in simulation • Robot kinematic model inaccuracy• Unexpected deploying scenarios(e.g. perturbation force))

### Improve Sim-to-Real Transfer Success 提高成功率• Reduce sim-to-real gap(Improve modelling in simulation: Better calibration; Data-driven method)• Improve robustness of control policy during training(Survive against sim-to-real mismatches)

### Domain Randomization 域随机化(提高成功率的一种方法)
Technique that randomly inject noises or perturbation to the system during learning(效果:). 随机 robot can meet as much as possible potential scenarios; Emerge generalized skills adapted to environment changes)

### Enhance Controller with DR(Sensor noises; Sensor delay / latency; Wind perturbation; Dynamic model mismatch; Actuator performance mismatch)

### DR is not the larger the better(why?Too conservative behavior if too large; Fail to learn the behavior)DR settings need robot experience & knowledge.

What's the Trend in RL for Robots?(• More modality: from blind to perspective • More advanced network architecture: from simple MLP network to memory to attention• More training complexity: from simple environments to more diverse environments• More skills: from single task to multiple tasks performance

Test: Reward Shaping; Constrained MDPs (strict)
Training: Constrained Predictive Control; Control Barrier functions (strict)

### Safety via Reward Shaping
Idea: define reward function that assigns large negative values to unsafe behavior: $r(s,a)=r_{task}(s,a)+r_{safety}(s,a)$
Strength: Standard RL algorithms applicable; Risk-sensitive RL variants(变体).Weakness: No strong guarantees; Complicated reward design.

### Risk-sensitive RL
Safety often considers low-probability events, average expectations can weaken the impact of unsafe returns. Risk-sensitivity desired.
$V^\pi(s)=\mathbb{E}_\beta[\sum_{k=0}^\infty\gamma^k r(s_k,\pi(s_k))]$
Entropic Risk $R_\beta[X]=\log(E[\exp(-\beta X)])/\beta$

### Constrained MDPs(CMDP)
Idea: explicitly represent safety conditions as constraints in RL
A set of states $\mathbb{S}$; A set of actions $\mathbb{A}$; Initial state distribution $p(s_0)$;
Markov kernel $p(s'|s,a)$; A reward function $r_{task}(s,a)$ for the task
$r_{task}(s,a)$,越大越好;A reward function $r_{safety}(s,a)$ for constraints eg binary

### Policy optimization problem(on-policy):
$\pi(s)=arg\max_\pi(r_{task}(s,a)+yE_{s'}[V_{task}^\pi(s')])$ 全局; $arg\max_\pi E[V_{task}^\pi(s)]$
such that $r_{safety}(s,a)+yE_{s'}[V_{safety}^\pi(s')]\ge c$ 全局; $E[V_{safety}^\pi(s)]\ge c$
constraint is defined via a value function:
$V_{safety}^\pi(s)=r_{safety}(s,a)+yE_{s'}[V_{safety}^\pi(s')]$
Temporal difference learning is applicable:
$V_{safety}(s)\leftarrow V_{safety}(s)+\alpha(r_{safety}+\gamma V_{safety}(s')-V_{safety}(s))$

### 方法 1.Constrained Policy Optimization (CPO) (policy iteration):
$\pi_{k+1}=arg\max E[A_{task}^\pi(s,a)]$
$s.t.E[V_{safety}^\pi(s_0)]-E[A_{safety}^\pi(s,a)]/(1-\gamma)\ge c$
$\bar D_{kl}(\pi\|\pi_k)\le\delta$ (新策略和旧策略的 KL 距离要小,不能变化太快)
Where,
$s_0$: initial state; Advantage functions: $A^\pi(s,a)=Q^\pi(s,a)-V^\pi(s)$ (动作危险与否); $1/(1-\gamma)$ is the tightening of-policy;
This guarantees:(1) Policy update is stable;(2) Will not suddenly perform dangerous actions;(3) The on-policy nature remains good.[trust region]

### 方法 2.Primal-dual method/Lagrangian relaxation:
$$\min_{\lambda\ge0} d(\lambda)$$
$d(\lambda)=\max_\pi L(\pi,\lambda)=\max_\pi E[V_{task}(s)]+\lambda(E[V_{safety}(s)]-c)$
For each round $t=1,...,T$ do
simulate a trajectory using current policy
estimate primal gradient $\nabla_\pi L(\pi,\lambda)$
estimate dual gradient $\nabla d(\lambda)=E[V_{safety}]-c$
gradient update: Gradient descent for $\lambda$, but gradient ascent for $\pi$
update value function estimates
CPO and Lagrangian methods rely on accurate constraint value functions ($E[V_{safety}(s)]-c$). If true constraint value is below the threshold but its estimate is above, Lagrange multiplier update is in wrong direction. Larger estimation errors in off-policy RL due to distribution shift become problematic.
不能用在 safe off-policy RL 上

## Column 1

### Control Barrier Functions(CBF)

h(x): >0, inside constraint set; 0 at constraint set boundary; <0 outside constraint set. $\pi_{safe}(s) = \min_a \|a - \pi(s)\|^2$

Condition for positive value (safety filter):

$s.t.\ E[h(s')] \geq \alpha h(s), \alpha \in [0,1], s' \sim p(\cdot|s,a)$

Allows handling new test-time constraints: + CBF implicitly represents a safe back-up strategy; - Can cause distribution shift. of on-policy rollouts Integration into training by lumping safety filter into environment; + RL agent can adapt to it; - Model knowledge required with little data May not feasible. =>learn from expert data, RL policy roll-outs

### Constrained predictive control(CPC)

Idea: optimize over sequence of actions.

$$\min_a E\left[\sum_{k=0}^{K} \gamma^k r(s_k, a_k) + \gamma^{k+1} V(s_{k+1})\right]$$

Only apply first action, optimize again; Model predictive control(MPC) Challenge: Computational complexity of online optimization.

### Online Optimization for MPC:

• Constrained Cross-Entropy Method (CEM): Randomly sample multiple action sequences (GPU parallelization!); Simulate their future trajectories (rollout); Select the top n% samples with the best performance (elite set); Update the distribution and repeat sampling.

Sequential Quadratic Programming (SQP 顺序二次规划): Linearize transitions and constraints; Approximately equivalent to a Quadratic Program (QP); Repeatedly iterate to approach the optimal solution.

Non-convex collocation / shooting methods: Taking system dynamics as the constraint of optimization variables. Suitable for deterministic systems(确定性系统)

### LEC8
### Partially Observable Markov Decision Process (POMDP)

A set of states $S$ (状态集); A set of actions $A$ (动作集)
A set of observations $O$ (观察集); Initial state distribution $p(s_0)$
Markov kernel $p(s'|s,a)$; State observation probabilities $p(o|s,a)$ (观察概率)
A reward function $r_{task}(s,a)$ (奖励函数)

### The challenge of solving POMDPs

Problem: A single observation does not provide full information of the state of a POMDP 仅依赖当前观察$a_k=\pi(o_k)$。当前无法区分 由状态不同但观察相同的情况→Agents need memory!

### Planning with Believe States

Idea: Iteratively estimate states from observations
→ Bayesian updates based on $p(s'|s,a)$ and $p(o|s,a)$
→ exact updates for special cases, e.g. finite state-action spaces
Believe state estimation => Kalman filtering
$b_k=p(s_k|o_0,a_0,\ldots,a_{k-1},o_k)$: $b_k$ 是概率分布，表示在给定迄今为止所有观察和动作历史的条件下，智能体在k时刻处于真实状态$s_k$的概率，即使真实状态 $s$ 是离散的，信念状态$b_k$也是一个连续变量的向量(概率分布)
→ Markovian, but continuous believe state
$b_k=P(s_k|h_k)=\{o_0,a_0,o_1,a_1,\ldots,a_{k-1},o_k\}$, 基于旧b_k和a_k预测先验$b_{k+1}$
$b_{k+1}(s)=P(s'|h_k,a_k)=\sum_s P(s'|s,a_k) \cdot P(s|h_k)=\sum_s P(s'|s,a_k) \cdot b_k(s)$
→ 接收$o_{k+1}$后更新$s=s$ 得到后验信念状态$b_{k+1}$
$b_{k+1}=P(s'|h_k,a_k,o_{k+1})=\frac{P(o_{k+1}|s',a_k)b_{k+1}(s')}{P(o_{k+1}|h_k,a_k)}$

Learning value functions:
策略的目标是最大化在信念状态$b_k$下的期望累积折扣奖励
$E\left[\sum_{k=0}^{\infty}\sum_{s\in S} \gamma^k b_k(s)r(s,a_k)\right]$ $\sum_{s\in S} b_k(s)r(s,a_k)$: 给定$b_k$时，执行$a_k$的即时期望 r
$a_k=\pi(b_k)$ $b_k=g(b_{k-1},a_{k-1},o_k)$ →requires model of transition and observation probabilities→tractable only for small problems in practice

### Approximate Planning with Believe States 精确太难

MLE approximation:最大似然，用信念状态中最有可能的真实状态代替
$E\left[\sum_{k=0}^{\infty}\gamma^k r(argmax_s b_k(s),a_k)\right]$, $a_k=\pi(argmax_s b_k(s))$, $b_k=g(b_{k-1},a_{k-1},o_k)$
从连续信念状态空间近似回离散的真实状态空间，就可以用 MDP 求解
→ standard planning problem with continuous states:理论上仍然是连续的
→ Much more: variants of value iteration, point-based methods, policy search, Monte-Carlo methods 针对 POMDP 的许多方法
POMDPs are powerful, but solving them is very complicated!

### RL for POMDPs

Model-based approaches(• Learn models for transition and output probabilities→ difficult) 学习 POMDP 的转移概率和观察，非常难
Model-free approaches(• Finite window of past observations and actions 依赖过去 L 个时间步→) counter examples can be found 更早的历史无效

### Recurrent Neural Networks

$h_t$ 依赖于前一时刻的隐藏状态$h_{t-1}$和当前时刻的输入$O_t$
Neural network becomes cyclic!→RNNs have memory!
technical difficulties:(back-propagation; Exploding/vanishing gradients)
Common realizations of RNNs for RL: LSTM, GRU

### Model-based: Deep Variational Reinforcement Learning, DVRL

Idea: Use RNNs in a model-based approach to determine believe

### Model-free: RNNs as parameterization

Idea: Use RNNs as encoder in
value function and
policy parameterization
用$h_t$充当信念状态的近似
Design considerations: 没有一刀切方案
• RNN variant (LSTM, GRU, …)
→ no significant difference
• Sequence length during training



## Column 2

(short: ~5, long ~100+)
• Architecture → separate encoders seem to work better
• Model-free RL algorithm → indicators for benefits of off-policy
• Encoder inputs → reward signal can be beneficial 奖励也作为输入
• End-to-end? → additional information is accessible for training
Performance Comparison(在没有速度信息情况下，从历史判断速度)
No method consistently outperforms all others 不同算法不同任务表现不一
Model-based method is comp. Expensive 基于模型的方法计算成本高
Difference between methods insignificant"(方法间的差异微乎其微)
### Important Special Cases of POMDPs (Specialized algorithms exist)
Meta RL problems→learning from different tasks without knowing about them 从观察到的奖励和环境反馈中推断出它所处的任务
Robust RL problems→unobservable state affecting transitions 抵抗扰动
Temporal credit assignment→delayed rewards

### Visuomotor Policy Training 视觉运动策略训练（输入是图片）
convolutional neural network layer for low dim. Representation 用 CNN
Challenge: Training difficult/data-hungry
Design considerations:
• RNN necessary? → finite/no history can work 有时可以省略 RNN
• End-to-end? → additional information is accessible for training 额外信息
• CNN pre-training? → unclear 不一定要与训练 CNN
• Many additional tricks for training CNNs
• Additional supervision can be used in training(auxiliary local policies 子任务的策略); expert demonstrations 用人类或高表现策略的演示数据

### LEC9 Increase Data Efficiency
### Inverse Reinforcement Learning(IRL)
Assumption: demonstrations(示范) are from optimal policy:知策略推奖励
$E\left[\sum_{t=0}^{\infty}\gamma^t r(s_t)|\pi^*\right] \geq E\left[\sum_{t=0}^{\infty}\gamma^t r(s_t)|\pi\right],\forall \pi$
在推断出的奖励函数下，专家策略获得的期望累积回报，必须大于等于任何其他策略获得的期望累积回报
Issues: Reward ambiguity: trivial solutions(琐碎解): $r(s_t)=0$; Relies on expert optimality; Computational complexity (calculate all $\pi$).
### Featurization of Reward Functions 限制奖励函数的搜索空间
Idea: Linearly parameterize the rewards $r(s)=w^T\phi(s)$.权重未知
$E\left[\sum_{t=0}^{\infty}\gamma^t r(s_t)|\pi\right]=E\left[\sum_{t=0}^{\infty}\gamma^t w^T\phi(s_t)|\pi\right]=w^T\mu(\pi)$
Find $w^*$ such that $w^T\mu(\pi^*)\geq w^T\mu(\pi)$ for all $\pi$
### Max-Margin Formulation
$$\min \|w\|^2 + c\varepsilon$$
$s.t.\ w^T\mu(\pi^*)\geq w^T\mu(\pi)+m(\pi,\pi^*)-\varepsilon,\forall\pi$
$m(\pi,\pi^*)$: the gap between the strategy $\pi$ and the expert strategy $\pi^*$;
$\varepsilon$: slack variable.
伪代码: Initialize weights $w$
For $t=1,\ldots,T$
obtain $\pi(s)$ by solving forward RL(often approximate solutions in practice)
problem with $r(s)=w^T\phi(s)$
do sub-gradient update step for $w$(non-differentiable), $\varepsilon$
### Max-Entropy Feature Expectation Matching
Insight: Similarity of policies can be measured independently from the weights $w$.(只要策略的特征期望$\mu(\pi)$和专家策略$\mu(\pi^*)$足够接近，那么在任何线性奖励函数 $w$ 下，$\pi$的价值和专家的价值也会很接近。)
Idea: Parameterize problem in terms of path probabilities $P(\zeta)$→ Maximize entropy of distributions over paths. Trajectory $\zeta=(s_0,a_0,s_1,a_1,\ldots)$
$\max_P -\sum_\zeta P(\zeta)\log P(\zeta)$ resolves ambiguity
$s.t.\ \sum_\zeta P(\zeta)\mu(\zeta)=\mu(\pi^*)$ sum of features along path
Solution via log-likelihood maximization of maximum entropy distribution
$P(\zeta)=exp(w^T\mu(\zeta))/z(\theta)=exp(w^T\mu(\zeta))/\sum_\zeta exp(w^T\mu(\zeta))$
### Imitation Learning (just want to learn a policy)
Use fixed data set $\mathcal{D}=\{(s_i,a_i)_{i=1\ldots N}\}$ to do empirical loss minimization
$\hat\pi^*=\arg\min \sum_{i=1}^N l(\pi,s_i,a_i)$, i.e. $l(\pi,s_i,a_i)=-\log \pi(a_i|s)$ or $\|\pi(s)-a\|^2$
### Behavior Cloning
Data is generally not uniformly distributed across the state space!
Expert policy $\pi^*$ induces a state distribution $p(\pi^*)$
$\arg\min \sum_i l(\pi,s_i,a_i)\approx \arg\min E_{s\sim p(\pi^*)}[l(\pi,s,\pi^*(s))]$ 近似期望损失
Distribution Shift→进入专家未见过的新状态，策略表现差，误差累积
Training distribution $s\sim p_{\pi^*}$; Testing distribution $s\sim p_{\hat\pi}$
Distribution shift occurs between training and testing!
compounding errors drop returns by at most $\varepsilon T^2$,
$\varepsilon$: training loss; T: number of time steps.
### On-Policy Imitation Learning
Idea: Generate data iteratively to converge to test distribution.
=>Employ mixing techniques.伪代码如下:
Initialize $\hat\pi(s),\mathcal{D}$
For $i=1,\ldots,N$
$\pi_i=\beta\pi^*+(1-\beta_i)\hat\pi$ (Data Aggregation (DAgger) method)
rollout policy $\pi i$ to generate trajectory $\tau=\{s_0,s_1,\ldots\}$
query expert to generate data $\mathcal{D}_i=\{(s_0,\pi^*(s_0)),(s_1,\pi^*(s_1)),\ldots\}$
aggregate data sets $\mathcal{D}\leftarrow \mathcal{D}\cup\mathcal{D}_i$
retrain $\hat\pi$ using $\mathcal{D}$
For sufficiently large $N$, returns for the policy $\hat\pi$ obtained from DAgger drop by at most $\varepsilon T$.
Downsides:(→policy needs to be retrained →needs access to expert policy)
### Multi-Modality
parameterize the policy $e.g., \hat\pi=\mathcal{N}(\mu(s),\Sigma(s))$由神经网络参数化
### Diffusion Models
Idea: Interpret imitation learning as fitting a probability distribution.
$\pi_\theta$ is Neural network.



## Column 3

• Learn/predict trajectories→increase temporal consistency; • Temporal convolution as encoder→capture temporal locality; • Include rewards in trajectories→optimization via conditioning; • Goal/initial state conditioning→fix some states (inpainting)

reverse/denoising process



### LEC10
### State value function $V^\pi(s)=E_\pi\left[\sum_t \gamma^t r_t \mid s_0=s\right]$ 给定策略$\pi$下，从状态 $s$ 开始，智能体所能获得的期望累积折扣奖励
### Action value function $Q^\pi(s,a)=E_\pi\left[\sum_{t=0}^\infty \gamma^t r_t|s_0=s,a_0=a\right]$给定$\pi$下，从$s$开始，首先采取$a$，从下一时刻开始遵循$\pi$，所获的期望累积折扣奖励
### Connection between $V^\pi(s)$and $Q^\pi(s,a)$: $V^\pi(s)=\sum_a \pi(a\mid s)Q^\pi(s,a)$
$V^\pi(s)$等于在所有可能动作下 a 的$Q^\pi(s,a)$的加权平均。
### Bellman Optimality 找到对每个状态s的最优值$V^*(s)$，前提是：takes the best possible action a and then continues acting optimally in the future.选择使$r(s,a)+\gamma E[V^*(s')]$最大化的动作 a，并假设下一状态的值$V^*(s')$也最优
### Optimal Q-value Function $Q^*(s,a)=\max_a E\left[\sum_{t=0}^\infty \gamma^t r_t \mid s_0=s,a_0=a\right]$
Bellman equation:
$$Q^*(s,a)=\sum_{s'} P(s'\mid s,a)\left[r(s,a)+\gamma \max_{a'} Q^*(s',a')\right]$$
### Value Iteration
$$Q_{i+1}(s,a)=E_{s'\sim P(\cdot|s,a)}\left[r(s,a)+\gamma \max_{a'} Q_i(s',a')|s,a\right]$$
Problem:Not scalable to cover / enumerate all the (s,a) pairs for Q(s,a)
### Solution in Q Learning
Use a function approximator: $Q(s,a;\theta)\approx Q^*(s,a)$不存储每一个值
If the function is neural network • Deep Q Learning
### Curse of Dimensionality 维度灾难
• Traditional RL methods fail in large state/action spaces • Robotics tasks often have continuous states and actions • Tabular methods become infeasible • Need for function approximation
### Motivation for Deep Q-Learning (DQN)
• Q-learning works well in small, discrete state spaces • Fails with high-dimensional inputs (e.g., images, robotics) • Deep neural networks approximate Q-values→ Enabled breakthroughs in Atari and robotics
### Deep Q Learning, DQN
Optimal Q function:
$$Q^*(s,a)=\mathbb{E}_{s'}\left[r(s,a)+\gamma \max_{a'} Q^*(s',a')|s,a\right]$$
Loss function for SGD:
$$l(\theta;s,a)=\frac{1}{2}\left(\underbrace{r(s,a)+\gamma \mathbb{E}_{s'}\left[\max_{a'} Q(s',a';\theta^{old})\right]}_{\text{目标 Q 值(Target Q-Value)}}-\underbrace{Q(s,a;\theta)}_{\text{当前 Q 值(Current Q-Value)}}\right)^2$$
Gradient update:
$\nabla_\theta l(\theta;s,a)$
### DQN Examples (Atari)
• Objective: Complete the game with the highest score
• State: Raw pixel inputs of the game state
• Action: Game controls e.g. Left, Right
• Reward: Score increase/decrease at each time step
1. 输入(Current state $s_t$) 84 × 84 × 4 stack of last 4 frames 最后四帧堆叠预处理: after RGB → grayscale conversion, downsampling, and cropping
2. CNN（16 8×8 conv, stride 4）（32 4×4 conv, stride 2）
3. Fully Connected Layers, FC 展开输入 FC 来进行抽象和 Q 值计算
4. Output Layer FC-4 (Q-values): $Q(s_t,a_1),Q(s_t,a_2),Q(s_t,a_3),Q(s_t,a_4)$
CNN 处理视觉输入（解决维度灾难），用 RNN 的思想（通过堆叠多帧）来引入时序信息，FC 层输出所有可能动作的 Q 值，选出最高的动作
### DQN Examples on Robots
•Observation: [target position in x axis (in pixels); target in y axis (in pixels); current joint 1 position (in radians); current joint 2 position (in radians)]
•Action Spaces (Discrete): [Hold current joint angle value; Increment joint 1; Decrement joint 1; Increment joint 2; Increment joint 1 and joint 2; Decrement joint 1 and joint 2]
### Training the Q-network: Experience Replay 经验回放
• Learning from batches of consecutive samples is problematic [Samples are correlated => inefficient learning 样本相关; Current Q-network parameters determines next training samples (e.g. if maximizing action is to move left, training samples will be dominated by samples from left-hand side) => can lead to bad feedback loops 陷入局部最优，形成不健康的负反馈循环]
• Address these problems using experience replay [Continually update a replay memory table of transitions ($s_t,a_t,r_t,s_{t+1}$) as game (experience) episodes are played 持续更新回放内存表; Train Q-network on random minibatches of transitions from the replay memory, instead of consecutive samples 从回放内存中采样 minibatch 训练,打破相关性,高效利用数据]
• Each transition can also contribute to multiple weight updates => data efficiency 样本可多次利用
Problem: The Q-function can be very complicated!But the policy can be much simpler. 学习高维状态空间下的 Q 函数的精确动作困难，但策略简单

**Algorithm 1** Deep Q-learning with Experience Replay
Initialize replay memory $\mathcal{D}$ to capacity $N$
Initialize action-value function Q with random weights
for episode = 1, M do
Initialise sequence $s_1=\{x_1\}$ and preprocessed sequenced $\phi_1=\phi(s_1)$
for $t=1,T$ do
With probability $\epsilon$ select a random action $a_t$
otherwise select $a_t=\max_a Q^*(\phi(s_t),a;\theta)$
Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$
Set $s_{t+1}=s_t,a_t,x_{t+1}$ and preprocess $\phi_{t+1}=\phi(s_{t+1})$
Store transition $(\phi_t,a_t,r_t,\phi_{t+1})$ in $\mathcal{D}$
Sample random minibatch of transitions $(\phi_j,a_j,r_j,\phi_{j+1})$ from $\mathcal{D}$
Set $y_j=\begin{cases} r_j & \text{for terminal } \phi_{j+1}\\ r_j+\gamma \max_{a'} Q(\phi_{j+1},a';\theta) & \text{for non-terminal } \phi_{j+1}\end{cases}$
Perform a gradient descent step on $(y_j-Q(\phi_j,a_j;\theta))^2$ according to equation 3
end for
end for

## Column 4

Value-based: Difficult with continuous action spaces; Indirect policy improvement via max operation; Need for explicit $\arg\max$ and exploration-exploitation balancing
### Policy-based methods directly learn:$\pi_\theta(a|s)$给定状态 s 时采取动作 a 概率
目标函数 $J(\theta)=\mathbb{E}_{\pi_\theta}\left[\sum_t \gamma^t r_t\right]$ 最大化策略 直接通过调整参数$\theta$来增加预期奖励，而不是通过学习一个 值函数来间接推导出策略
### Policy Gradient Theorem
Gradient of expected return: $\nabla_\theta J(\theta)=\mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(a|s) Q^{\pi_\theta}(s,a)]$
$\nabla_\theta \log \pi_\theta(a|s)$: 选择动作a的对数概率的梯度; $Q^{\pi_\theta}(s,a)$该动作有多好
$\nabla_\theta \log \pi_\theta(a|s)=\frac{\nabla_\theta \pi_\theta(a|s)}{\pi_\theta(a|s)}$
$J(\theta)=\sum_s d^{\pi_\theta}(s)\sum_a \pi_\theta(a|s)Q^{\pi_\theta}(s,a)$ Such that $\nabla_\theta\pi_\theta(a|s)=\pi_\theta(a|s)\nabla_\theta \log \pi_\theta(a|s)$
$d^{\pi_\theta}(s)$:$\pi_\theta$下各个状态的平稳分布,即 agent 在长期运行中处于状态 s 的概率
对求微分:$\nabla_\theta J(\theta)=\sum_s d^{\pi_\theta}(s)\sum_a \nabla_\theta \pi_\theta(a|s)Q^{\pi_\theta}(s,a)$ 只对$\pi_\theta(a|s)$微分
为什么不对 Q 和 d 求微分:(Simplified to avoid over-complexity)
$\nabla_\theta J(\theta)=\sum_s d^{\pi_\theta}(s)\sum_a \pi_\theta(a|s)\nabla_\theta \log \pi_\theta(a|s)Q^{\pi_\theta}(s,a)$
转化为期望=>$\nabla_\theta J(\theta)=\mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(a|s)Q^{\pi_\theta}(s,a)]$
### Intuition Behind REINFORCE
• If an action leads to higher reward → increase probability
• If an action leads to lower reward → decrease probability
• Similar to trial-and-error learning • Policy learns from sampled trajectories
### REINFORCE Algorithm
REINFORCE uses the empirical return from the episode 使用经验回报实际观察到的累积奖励 taking $a_t$ in $s_t$: $G_t=\sum_{k=t}^T \gamma^{k-t} r_k$ 替换 Q
Monte Carlo policy gradient method
$$\nabla_\theta J(\theta)\approx \frac{1}{N}\sum_{i=1}^N \sum_{t=1}^{T_i} \nabla_\theta \log \pi_\theta(a_t^i|s_t^i)G_t^i$$
对 N 条计算取平均以近似期望
### Algorithm Steps: • Initialize policy parameters $\theta$ • Run episodes using $\pi_\theta(a|s)$ • Compute returns $Gt$ (Monte Carlo cumulative reward) • Update policy: $\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi_\theta(a_t|s_t)G_t$ • Repeat until convergence
### Main Drawbacks
• High Variance in Gradient Estimates (问题:The gradient depends on full-episode returns $G_t=\sum_{k=t}^T \gamma^{k-t} r_k$。后果: noisy and unstable learning, especially for long or stochastic environments。问题:Variance grows with episode length → slower convergence 高方差导致学习收敛速度慢)
• Inefficiency in Sample Usage ( REINFORCE is a pure on-policy, Monte Carlo algorithm)就被抛弃。off-policy 可重用历史数据
• Delayed Credit Assignment 更新必须等到整个轨迹完成后，不能实时学 $G_t$

| Concept | REINFORCE | Off-policy (e.g. DQN) |
|---|---|---|
| Sampling | On-policy (from current π) | Off-policy (reuse old data) |
| Replay buffer | ❌ Not compatible | ✅ Essential |
| Why not | ❌ Distribution mismatch → biased gradient | ✅ Learning to just approximate the Q-function |
| Data efficiency | Low | High |
| Update target | $G_t$ (full return) | $r+\gamma Q(s,a)$ |
| **Combination:** | **stability + flexibility** | |

| Method | Strength | Weakness |
|---|---|---|
| Value-based (e.g., DQN) | Low variance, stable learning | Not suitable for continuous actions; indirect policy improvement |
| Policy-based (e.g., REINFORCE) | Works for continuous action spaces; directly optimizes policy | High variance, slow convergence |

### Actor–Critic
### Policy-based part: Actor 演员是策略网络，负责学习和执行策略
Learns a parameterized policy $\pi_\theta(a|s)$, improving action selection via gradient ascent: $\nabla_\theta J(\theta)=\mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(a|s)A]$
通过梯度改进动作选择，目标是最大化期望回报
### Value-based part : Critic 评论家是一个价值网络，负责评估演员的策略
Learns a value function to estimate future returns, reducing variance of the gradient.学一个值函数估计来估计未来回报，将这个估计值用于降低策略梯度的方差
### Motivation: REINFORCE 用 Gt 估计引起的高方差可以用 learned value estimate 替代Gt 解决→ Keeps the gradient unbiased but reduces variance
### Advantage Function: How much better is this action than average?
• If $A(s,a)>0$ :increase probability of this a • If $A(s,a)<0$ :decrease
• If $A(s,a)=0$ :no change (action is average)
Make A(s,a) the most informative and stable signal for updating the policy
### Why• Using Q(s,a) alone → high variance->Includes shared baseline(to A)
• Subtracting a baseline doesn't change the expectation 引入基线，期望不变
• The best baseline V(s)-> No action discrimination-> Ignores a differences 减去 V(s) 可以消除奖励中与动作无关的、共同的、高噪声的部分=>Variance reduction effect 方差降低 Reduce random fluctuations around the mean reward; Training becomes more stable and converge faster)
### Practical Considerations (Learning rate tuning is critical; Reward normalization helps stability; Exploration vs exploitation trade-off; Use entropy regularization to encourage exploration)
### Limitations(Poor exploration; On-policy inefficiency; Instability)
### LEC11
### Entropy-Regularized RL(Introduce entropy term to encourage exploration)
$$J(\theta)=E_{\pi_\theta}\left[\sum_t \gamma^t (r_t+\alpha\mathcal{H}(\pi(a|s_t)))\right]$$
Entropy: $\mathcal{H}(\pi(a|s))=-E_{\pi}[\log \pi(a|s)]$
$\alpha$(temperature) controls exploration - exploitation trade-off 大则探索性强
### Soft Bellman Equations $V^\pi(s)=E_{a\sim\pi}[Q^\pi(s,a)-\alpha\log \pi(a|s)]$ Soft value adds entropy bonus → encourages policies that remain uncertain and exploratory
### Comparison to Previous Methods
Q-Learning=>Q-function=>Simple, stable=>Discrete only
REINFORCE=>Policy=>continuous spaces=>High variance
Actor–Critic=>Policy,Value=>Lower variance=>Lacks exploration

## Column 5

SAC=>Policy + Soft Q=>Stable, efficient, exploratory=>complex(offpolicy)
### Deep Deterministic Policy Gradient (DDPG)
Make the actor policy deterministic, but still learn via gradients from a critic.
• Pros(Works well in high-dimensional continuous spaces; Off-policy → high sample reuse; Faster convergence than stochastic actor methods
• Cons[Deterministic → poor exploration (needs noise injection); Prone to overestimation bias (one critic); Sensitive to hyperparameters → unstable]
### Unstable Policy Updates=>trust region methods.
### Trust Region Optimization Constrain the policy update so that the new policy is not too different from the old one $\max_\theta J(\theta)$ s.t. $D_{KL}(\pi_{old}(\cdot|s)\|\pi_\theta(\cdot|s))\leq\delta$目标仍然是 $J(\theta)$，约束条件是，旧策略与新策略的 KL 散度不能超过阈值
KL 散度是分布之间的信息差 $D_{KL}(P\|Q)=\int P(x)\log \frac{P(x)}{Q(x)}dx$ 非负,不对称
$D_{KL}(P\|Q)$: 衡量了使用 $Q$ 来近似 $P$ 时的信息损失。
$D_{KL}(Q\|P)$: 衡量了使用 $Q$ 来近似 $Q$ 时的信息损失。所以不等
### TRPO Objective Function 重要性采样校正后的策略梯度目标
$L_{TRPO}(\theta)=E_{s,a\sim\pi_{old}}\left[\frac{\pi_\theta(a|s)}{\pi_{old}(a|s)}A^{\pi_{old}}(s,a)\right]$
调整$\pi_\theta(a|s)$的概率 Subject to $E_{s\sim\pi_{old}}[D_{KL}(\pi_{old}(\cdot|s)\|\pi_\theta(\cdot|s))]\leq\delta$
如果$A^{\pi_{old}}(s,a)>0$(动作好)，增大比值，增加新策略 $\pi_\theta$ 选择动作 $a$ 的概率
如果$A^{\pi_{old}}(s,a)<0$(动作坏)，减小比值，降低新策略 $\pi_\theta$ 选择动作 $a$ 的概率
### Importance Sampling 没有 P(x)的样本，权重改写为 $Q(x)$P(x)/Q(x)
$E_{x\sim P}[f(x)]=\int P(x)f(x)=\int Q(x)\frac{P(x)}{Q(x)}f(x)=E_{x\sim Q}\left[\frac{P(x)}{Q(x)}f(x)\right]$
### Why Needs It: 使用重要性采样来修正新策略和旧策略间的偏差，转换为旧策略$\pi_{old}$下期望。可以在不采集新策略的情况下，计算新策略的近似值
### TRPO Derivation Intuition = policy gradient with a geometry-aware step
Linearize the surrogate objective (first-order)即 $A$ 给一个向上改进的方向
Quadratically approximate the KL constraint (second-order)KL 二阶展开
Solve via natural gradient $\theta_{t+1}=\theta_t+\sqrt{2\delta/(g^TH^{-1}g)}H^{-1}g$
Detail: $L_{TRPO}(\theta)\approx L(\theta_k)+g^T(\theta-\theta_k)$, $D_{KL}(\pi_{old}//\pi_\theta)\approx\frac{1}{2}(\theta-\theta_k)^TH(\theta-\theta_k)$
$g=\nabla_\theta L_{TRPO}(\theta)$ Gradient of the surrogate objective
$H$: KL 在新目标展开时，二阶海森矩阵,it defines the trust region shape.
We use the Lagrangian to solve it:$L(\theta,\lambda)=g^T(\theta-\theta_k)-\lambda(\frac{1}{2}(\theta-\theta_k)^TH(\theta-\theta_k)-\delta)$
Setting the gradient w.r.t θ to zero: $\nabla_\theta L(\theta,\lambda)=g-\lambda H(\theta-\theta_k)=0$, $\theta_k=\theta+H^{-1}g/\lambda$
Plug back into the constraint: $\lambda=\sqrt{2\delta/(g^TH^{-1}g)}$
So, $\theta_{k+1}=\theta_k+\sqrt{2\delta/(g^TH^{-1}g)}H^{-1}g$
### Implementation of TRPO(Collect trajectories under $\pi_{old}$; Estimate $A(s,a)$;
Compute gradient $g$; Compute $H$ and solve $H^{-1}g$ via conjugate gradient 计算量最大; Update parameters with line search to satisfy KL constraint)
Line Search(防止线性近似和二阶近似不准确，所以线性搜索一个步长
full step size:$\alpha=\sqrt{2\delta/(g^TH^{-1}g)}$ 沿自然梯度方向$\Delta\theta=H^{-1}g$ 减小，选择满足约束且带来最大回报改进的步长$\theta_{k+1}=\theta_k+\alpha\Delta\theta$
### Advantages(Stable updates; Theoretically justified monotonic improvement; Natural gradient direction; Strong baseline method)
### Limitations(Complex to implement; High computational cost;Hard to tune $\delta$; Hard to scale to large networks)
### Proximal Policy Optimization, PPO TRPO's lightweight, firstorder cousin (•Avoids second-order optimization--no Hessian) • Achieves similar stability and performance • Easy to implement with standard gradient descent)
### Clipped Objective: $L^{CLIP}(\theta)=E_t[\min(r_t(\theta)A_t, clip(r_t(\theta),1-\epsilon,1+\epsilon)A_t)]$
$r_t(\theta)=\pi_\theta(a|s)/\pi_{old}(a|s)$ Don't let the prob. ratio deviate too far from 1
### Full Loss Function: $L_{Total}(\theta)=L^{CLIP}(\theta)+c_1 L^{VF}-c_2\mathcal{H}(\pi)$
Value Function Loss:$L^{VF}$损失 MSE Entropy Loss:$-\mathcal{H}(\pi)$ 没符号的是最大化
### Steps(Collect trajectories under $\pi_{old}$; Estimate advantages (GAE); Compute ratio $r_t(\theta)$; Optimize clipped objective with SGD; Update $\pi_{old}\leftarrow\pi_\theta$ after several epochs) 优点 Simple to code, robust to hyperparameters
### Generalized Advantage Estimation, GAE 更准确地估计 A,以减小方差
$\hat{A}_t=\sum_{l=0}^\infty (\gamma\lambda)^l\delta_{t+l}$, $\delta_t=r_t+\gamma V(S_{t+1})-V(S_t)$ 单步 TD 残差加加权求和, $\delta_t=r_t+\gamma V(S_t)$
$\lambda$ controls bias–variance tradeoff ($\lambda=0$ low bias, $1=$low variance) PPT 说$\gamma$
$\lambda=0$: Normal A; $A_t=r_t+\gamma V(S_{t+1})-V(S_t)$; Simple, low variance; - High bias
$\lambda=1$: Monte Carlo A; $A_t=\sum_l \gamma^l \delta_{t+l}$; - Unbiased; - High variance
$\lambda=0-1$: GAE; $\hat{A}_t=\sum_{l=0}^\infty (\gamma\lambda)^l\delta_{t+l}$; Bias-variance tradeoff; Controlled by $\lambda$
Explanation: (0. Relies on imperfect critic's short prediction; 1. Uses full, noisy returns; 0-1: Smooth compromise between both)
### Compare:Constraint;Optimization;Computation;Performance;Implemented
TRPO:KL-divergence (hard constraint);Second-order (natural gradient); Expensive; Very stable; Theory-focused works
PPO:Clipping (soft constraint); First-order (SGD); Lightweight; Almost as stable, faster; Practically all modern RL (OpenAI Baselines)
### Why PPO? (Stability:Clipping prevents large updates; Simplicity:No need for KL constraint or FIM; Efficiency:Uses multiple epochs of on-policy data; Exploration:Maintains entropy term; Universality:discrete & continuous)
### Limitations of PPO (On-policy Inefficiency; Constraints guarantee; Sensitivity to hyperparameters; Value function instability; Limited exploration; Black-box system)
### 基础公式以及方法分类
$V^*(s)=\max_a Q^*(s,a)$
Bellman Optimality Equation:
星号代表所有可能策略下最大期望回报
$Q^*(s,a)=r(s,a)+\gamma\sum_{s'} P(s'|s,a)V^*(s')$
Bellman Equation:
$V^\pi(s)=\sum_a \pi(a|s)Q^\pi(s,a)$
对固定策略，用于评估策略价值
$Q^\pi(s,a)=\sum_{s',r} p(s',r|s,a)\left[r+\gamma\sum_{a'}\pi(a'|s')Q^\pi(s',a')\right]$
$V^\pi(s)=\sum_a \pi(a|s)\sum_{s',r} p(s',r|s,a)[r+\gamma V^\pi(s')]$
$Q^*(s,a)=\sum_{s',r} p(s',r|s,a)[r+\gamma \max_{a'} Q^*(s',a')]$
$V^*(s)=\max_a \sum_{s',r} p(s',r|s,a)[r+\gamma V^*(s')]$