深度强化学习:使用深度神经网络作为强化 学习的函数近似器。



传统在线 O 学习问题: 1, 样本不满足 iid:

2, 无法证明收敛, 更新并非梯度下降。

经验回放+目标网络结合的Q学习步骤: 保存目标Q网络参数φ'←φ 2.利用某些策略收集样本 {(s_i, a_i, s_i', r_i)}, 加入样本池B 3 在经验池 B 采样 hatch 数据 3.在经验泡8 木叶Dotch19x20v K 4. 更新参数φ+φ-α∑_i ^{agg}(s_i,a_i)(Q_φ(s_i,a_i) (r(s_i,a_i)+ymax_a(Q_φ(s_i,a_i))

4. 更新参数 $\varphi \leftarrow \varphi - \alpha \sum_{i} \frac{dQ_{\varphi}}{d\varphi}(s_{i}a_{i})(Q_{\varphi}(s_{i},a_{i}) + r(s_{i},a_{i}) + r(s_{i},a_{i}) + r(s_{i},a_{i}))$

Batch 训练降低方差

利用ε-greedy策略执行动作α_i, 收集样本{(s_i, α_i, s'_i, r_i)}, 加入经验池Σ

K=1 N=1

- 2.在经验池B 中采样batch数据 $\{(s_j,a_j,s_j',\tau_j)\}$
- 3. 计算目标网络估计值 $y_j \leftarrow \begin{cases} y_j = r_j; \\ r(s_j,a_j) + \end{cases}$
- 4. 更新Q网络参数 $\varphi \leftarrow \varphi \alpha \sum_{i} \frac{dQ_{\varphi}}{d\varphi}(s_{j}a_{j})(\overline{Q_{\varphi}(s_{j},a_{j})} y_{j})$
- 5.更新目标网络参数 g/←g

多步回报得更低的方差。Q 值估计不准, TD 目标值中的 max 引入正向偏差,导致下一时 刻的目标值过估计,解:DDON 使用不同的网 络来分别计算目标 Q 网络值和选择动作; 当 前Q选动作;目标Q值估计目标值。

 $y_i \leftarrow r(s_i, a_i) + \gamma Q_{\omega'}(s'_i, \operatorname{argmax} Q_{\omega'}(s'_i, a'_i))$ 标准DQN

Double DQN $y_j \leftarrow r(s_j, a_j) + \gamma Q_{\varphi'}(s'_j, \operatorname{argmax} Q_{\varphi}(s'_i, a'_i))$ DON 中经验池均匀采样的效率不高,不能反 驶样本的价值·带优先级的经验同放,利用 ID 误差去衡量优先级,大优先级高。优先级

经验回放带来的问题: TD 误差对噪声敏感: ID 误差小的得不到更新: 过分关注 TD 误 差大的样本, 丧失样本多样性。引入重要性 采样权重来平衡"有偏"问题。前期注重优先

级高后期注重无偏性。

Dueling-DQN: 将Q函数分解成不依赖动作 的值函数V(s)依赖动作的优势函数A(s,a)。 分解的原因: 对于很多状态并不需要估计每 个动作的值,增加了 V 函数的学习机会; V

函数的污化性能好, 当有新动作加入时, 并 不需要重新学习;减少了Q 函数由于状态和

动作维度差导致的噪声和突变模仿学习

适用情况 1.提供多演示轨迹: 状态-动作序列 2.很容易演示, 轨迹容易收集。目标: 智能体 需要找到期望策略,使得该策略下状态-动作 轨迹分布尽可能匹配专家演示样本。输入: S,A, P (s'|s, a); 无R; 专家演示样本集 s0, a0, s1, a1, ...。利用监督学习获得期望策略 (状态到专家动作的映射)输入为演示样本集 中的状态,标签为演示样本集中的动作。行 为克隆: 专家演示样本集 $\{St, at | t \in T\}$,在专 家演示样本集训练得到期望策略π最小化误 设: 训练集和测试集的分布一致, 而行为克

隆存在的问题:复合误差,最终导致分布漂

移样本增广:累计误差的解决办法。数据聚

1. 训练 $\pi_{\theta}(a_t|s_t)$ 基于专家演示数据集 $\{(s_t, a_t|t \in T)$ 2. 执行 $\pi_{\theta}(a_t|s_t)$ 得到新轨迹 $\mathcal{T}_{\pi} = \{s_1, s_2, ... s_M\}$

3. 专家再次标记新轨迹 T_{π} 的标签 $\{a_1, a_2, ... a_M\}$

4. 数据聚合: T ← T ∪ T_n

5. 执行步骤1. 循环

合:

数据聚合的问题: 1.不安全/部分训练的策略,

2. 需要大量专家标记,与任务相关。模仿学习

经常面临:数据分配不匹配:专家非马尔可

夫行为。**专家多模态行为一些提升的办法:**

1 增广 2 从一个稳定的轨迹分布中采样 3

聚合.4 改讲模型提升精度 5. 多的数据。 逆强

化学习假设专家策略 π E是最优的,从专家演

示样本中恢复出奖赏函数R*

价值函数逼近

好处: 更少的参数: 泛化能力: 更少样本学精 确价值函数;多样性:多种特征表示和逼近器 结构。**价值函数逼近的类型:** 1. V(s);2. Q(s, a)3.O(s)输出所有的 O(s, a), ∀a

线性函数语近:

用特征向量 (Feature Vector)代表某一状态:

$$x(s)=[x_1(s), x_2(s), x_3(s)...x_n(s)]^T$$

$$\hat{V}(s, \mathbf{w}) = \mathbf{x}(s)^T \mathbf{w} = \sum_{j=1}^n \mathbf{x}_j(s) \mathbf{w}_j$$

均方误差(MSE)·逼近函数和直实价值函数之 间的误差(期望损失) d 智体在 π 的状态分布

$$J(w) = \int s \ d(s)(V\pi(s) - V^{\hat{}}(s,w))^2$$

对线性逼近器, 用样本近似期望损失:

$$\frac{1}{M} \sum_{m=1}^{M} (V_{\pi}(s_m) - \mathbf{x}(s_m)^T \mathbf{w})^2$$

Π 代表从 f(x)到 g(x|w)最佳匹配。满足:

$$\hat{V}(\mathbf{w}^*) = \Pi(V_{\pi})$$
s.t. $\mathbf{w}^* = \arg\min_{\mathbf{w}} \mathbb{E}_{\pi} \left[(V_{\pi}(s) - \mathbf{x}(s)^T \mathbf{w})^2 \right]$

$$\mathbf{w}^* = \left(\sum_{m=1}^T \mathbf{x}(s_m)\mathbf{x}(s_m)^T\right)^{-1} \sum_{m=1}^T \mathbf{x}(s_m) V_{\pi}(s_m)$$

最小二乘法优占, 计算过程简单 一次性求 解最佳权重.求逆,O(n3);当特征空间比较大 时, 计算量大, 矩阵可能不满秩, 易有误差

常见的特征表示方法

查表法: 离散化法: 粗糖编码: 圆形区域特征可 以重叠.径向基函数泛化特征值可以是 [0,1], 代表属于该特征的程度.

价值迭代 + 离散化方法

- 将状态空间离散化成状态子集 $\{S_k\}, k = 1, ..., K$ 定义有限动作集 $\{A_l\}, l = 1,...$
- 3: 初始化 $Q_{k,l}^{(0)}$ (e.g. $Q_{k,l}^{(0)}=0$), i=04: repeat {在第 i 次迭代} 5: 更新 Q 值
 - 1 P. P为确定型 $Q_{k,l}^{(i+1)} = \mathcal{R}_{k,l} + \gamma \max_{l} Q_{k',l'}^{(i)}, \quad \forall k, l$

6: $i \leftarrow i+1$ 7: $\mathbf{until} \parallel Q^{(i)} - Q^{(i-1)} \parallel < \varepsilon$

收敛性;离散化程度越高, 结果越精确,但是 存储空间和计算量增加; 维数灾(指数增)

Fitted Q Iteration (Q_learning)

Least-Squared FOI

```
1: 定义特征向量 \mathbf{x}, 和线性 \mathbf{Q} 逼近器 Q(s,a,\mathbf{w}_0)=\mathbf{x}^T\mathbf{w}_0, 初始化
    重 \mathbf{w}_0, i=0
采样一组数据 \{(s_k, a_k, r_{k+1}, s_{k+1})\}, k=1, \ldots, K
    repeat for all k = 1, \dots, K do
```

 $\mathfrak{L} \mathfrak{M} \mathbf{w}_{i+1} = \operatorname{arg min}_{\mathbf{w}} \frac{1}{K} \sum_{k=1}^{K} (q_k - \mathbf{x}(s_k, a_k)^T \mathbf{w})^2$

until w_i 收敛或迭代一定次数

傍略迭代 + 最小二乘 (闭解)

根据 pi 采样?如果 π 是确定策略, 需要探索

的动作,提泛化能力,如果π是随机策略,并

目 $\pi(a|s) > 0$, 可以 at $\sim \pi(st)$

最小二乘策略迭代 LSPI:

给定策略 π_0 , 定义特征向量 $\mathbf{x}(s,a)$, 初始化 i=0

: $\infty < x < w < a_0$ 、 人人では同量 $X(s, u_t, r_{t+1}, s_{t+1})$ } : (東略评估:) 計对策略 π_i 使用 LSTD-Q 计算权重 \mathbf{w}_i , 得到 策略的线性 Q 函数通道 $\mathbf{Q}_i(s, \mathbf{a}, \mathbf{w}_i) = \mathbf{x}^T(s, \mathbf{a}) \mathbf{w}_i$: (策略提升:) 提取食心策略 $\pi_{t+1}(s) = \arg\max_a Q_i(s, a, \mathbf{w}_i)$

until π. 收敛或迭代一定次数

approximate VI/PI 优缺点

优点:收敛;缺点:每次迭代的计算量大,线性最

小二乘法求样本特征矩阵的逆;通常 PI 的 迭代次数要比 VI 迭代次数少

预测学习 + 随机梯度下降法

梯度下降 MC 预测算法

```
给定策略 \pi, 定义价值逼近器 \hat{V}(s, \mathbf{w}), 初始化 \mathbf{w}
repeat {对每个 episode}
```

根据 π 采样轨迹 $\{s_0, a_0, r_1, s_1, \ldots, s_T\}$ repeat $\{$ 对 episode 中每个首次访问的 $s_t\}$

计算回报 $G_t = r_{t+1} + \gamma r_{t+2} + \dots$ 计算更新量 $\Delta \mathbf{w} = \alpha \left(G_t - \hat{V}(s_t, \mathbf{w}) \right) \nabla_{\mathbf{w}} \hat{V}(s_t, \mathbf{w})$

更新权重 $\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}$

until until

梯度下降 TD 预测学习,

1: 给定策略 π , 定义价值逼近器 $\tilde{V}(s, \mathbf{w})$, 初始化 \mathbf{w} , $s_t = s_0$

loop 选择执行动作 $a_t \sim \pi(s_t)$, 观测 r_{t+1}, s_{t+1} 计算更新量 $\hat{v}' \sim \hat{v}' \sim \hat{v} \cdot \hat{\mathbf{w}} - \hat{V}(s_t, \mathbf{w})$ $\Delta \mathbf{w} = \alpha (r_{t+1} + \gamma \hat{V}(s_{t+1}, \mathbf{w}) - \hat{V}(s_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{V}(s_t, \mathbf{w})$ 更新权重 $\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}$

7: end loop

对于前向 TD(λ): 上式为

$$\Delta \mathbf{w} = \alpha (\mathbf{G}_t^{\lambda} - \hat{V}(s_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{V}(s_t, \mathbf{w})$$

梯度下降 TD(λ)(后向)算法:

計算 TD 误差 $\delta_t = r_{t+1} + \gamma \hat{V}(s_{t+1}, \mathbf{w}) - \hat{V}(s_t, \mathbf{w})$ 更新資格達 $e_t = \gamma \lambda e_{t-1} + \nabla_{\mathbf{w}} \hat{V}(s_t, \mathbf{w})$ 更新収重 $\mathbf{w} \leftarrow \mathbf{w} + \alpha \delta_t e_t$

梯度下降 MC 控制算法

定义 Q 函数逼近器 $\hat{Q}(s,a,\mathbf{w})$, 初始化 w

及又 国 政連元章 Q(s,a,w), 初知化 教 検取出 ϵ 貪心策略 $\pi = \epsilon \operatorname{greedy}(Q(w))$ repeat { 村寿 γ episode 根据 π 采样軌迹 $\{s_0, a_0, r_1, s_1, \ldots, s_T\}$ repeat { 対 episode 中寿个首次访问的 (s_t, a_t) }

・ 日本 いっぱい $\{s_t, a_t\}\}$ 计算母根 $G_t = r_{t+1} + \gamma r_{t+2} + \dots$ 计算更新 $\Delta \mathbf{w} = \alpha(G_t - \hat{Q}(s_t, a_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s_t, a_t, \mathbf{w})$ 更新权重 $\mathbf{w} \leftarrow \mathbf{w} + \Delta \mathbf{w}$

更新策略 $\pi = \epsilon$ -greedy($\hat{Q}(\mathbf{w})$)

1: until

梯度下降 Sarsa 算法

1: 定义 Q 函数逼近器 Q(s, a, w), 初始化 w.

$$\begin{split} \pi &= \epsilon\text{-greedy}(\hat{Q}(\mathbf{w})), \ s_t = s_0, \ t = 0 \\ \text{采样动作} \ a_t &\sim \pi(s_t) \ \text{并执行}, \ \text{观测} \ (r_{t+1}, s_{t+1}) \end{split}$$

3: **loop**

 $\delta_t = r_{t+1} + \gamma \hat{Q}(s_{t+1}, a_{t+1}, \mathbf{w}) - \hat{Q}(s_t, a_t, \mathbf{w})$

更新权重 $\mathbf{w} \leftarrow \mathbf{w} + \alpha \delta_t \nabla_{\mathbf{w}} \hat{Q}(s_t, a_t, \mathbf{w})$ 更新策略 $\pi = \epsilon\text{-greedy}(\hat{Q}(\mathbf{w}))$

梯度下降 Sarsa(λ) 算法

更新資格述 $e_t = \gamma \lambda e_{t-1} + \nabla_{\mathbf{w}} \hat{Q}(s_t, a_t, \mathbf{w})$ 更新权重 $\mathbf{w} \leftarrow \mathbf{w} + \alpha \delta_t e_t$

梯度下降的 0 学习算法

 $\delta_t = r_{t+1} + \gamma \max \hat{Q}(s_{t+1}, a', \mathbf{w}) - \hat{Q}(s_t, a_t, \mathbf{w})$

策略梯度

策略 RL 好处:更好的收敛性:有效解决大规模 动作集或连续动作空间问题:能够学习随机 策略。缺点:收敛到局部最优;策略评估费力且 方差大.**策略逼近器不同形式:** $p(a|s)=pi(a|s,\theta)$. $a=pi(s, \theta), a\sim \mathcal{N}(\mu(s, \theta), \sum(s, \theta))$.不同优化目 标。优化问题:无梯度优化 (爬山法):梯度优化 有限差分法(FD):对策略参数 θ 每个维度增加 微小扰动计算.

解析法求策略梯度:采样 N 条轨迹,用样本近 似策略梯度: $\nabla_{\theta}J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \Big(\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_{i}|s_{t}) \Big) \Big(\sum_{t=0}^{T-1} r_{t+1} \Big)$

REINFORCE 算法 (蒙特卡洛策略梯度算法):

利用时间连贯性改变解析法的 $\nabla I(\theta)$: 时刻之 前的奖励与 t 时刻的策略无关, 所以不影响 t 时刻的梯度更新;同时将上述扩展到连续运

行场景。实际算法不可能产生无穷长的轨迹

1: 定义策略逼近器 π_{θ} , 初始化参数 θ , 初始状态分布 drepeat 根据 d 和 π_{θ} 采样 N 条軌途 $\tau_i = \{s_0^i, a_0^i, r_1^i, s_1^i, \dots\}$

for all $i = 1, \dots, N$ do

截断计算。

for all $t=0,1,\dots$ do 计算回报 $G_t^i=r_{t+1}^i+\gamma r_{t+2}^i+\dots$ 计算更新量 $\Delta\theta=\Delta\theta+\alpha\nabla_\theta\log\pi_\theta(s_t^i,a_t^i)G_t^i$ end for

end for 更新策略参数 $\theta \leftarrow \theta + \Delta \theta$ until 达到一定迭代次数或策略无明显的提升

Actor-Critic: Q 函数逼近器 $Q_w(s,a)$ 为 Critic,

策略逼近器 $\pi_{\theta}(s,a)$ 为 Actor.Critic 更新价值 函数的权重;Actor利用O(降低方差)更新策略

梯度。Actor 的策略梯度基于 Critic 定. 故让

β大一些.Critic 学的快一些.

1: 定义 Q 函数逼近器 $Q_{\mathbf{w}}(s,a)$, 策略逼近器 $\pi_{\theta}(s,a)$, 初始化 \mathbf{w} , 2: repeat 3: 采样动作 $a_t \sim \pi_{\theta}(s_t, a_t)$ 并执行, 观测 r_{t+1}, s_{t+1}

$$\begin{split} & \star \star + \eta \eta \tau_{\mathbf{k}} \ a_{t} \sim \pi_{\theta}(s_{t}, a_{t}) \ \mathcal{H} \mathcal{H}(1, \mathcal{H}(\mathbf{k}), \mathbf{k}, \mathbf{k}) \ \mathbf{f} \\ & \delta = r_{t+1} + \gamma Q_{\mathbf{w}}(s_{t+1}, a_{t+1}) - Q_{\mathbf{w}}(s_{t}, a_{t}) \\ & \mathbf{w} \leftarrow \mathbf{w} + \beta \delta \nabla_{\mathbf{w}} Q_{\mathbf{w}}(s_{t}, a_{t}) \\ & \theta \leftarrow \theta + \alpha \nabla_{\theta} \log \pi_{\theta}(s_{t}, a_{t}) Q_{\mathbf{w}}(s_{t}, a_{t}) \end{split}$$

 $t \leftarrow t + 1$ 8: until

AC算法Q和pi满足条件则收:Q和pi兼容;

Q 的权重 w 等于真实 $Q_{\pi_{\theta}}$ 的最小均方误差解

策略梯度引入基准:在回报上减去一个

baseline 不改变梯度. b 值可以方差降到最低。

对于 episodic:b= V_0 .连续 b (s_t) = $V_{\pi_{\theta}}(s_t)$ 。**优势** 函数: $A_{\pi_{\theta}}(s,a) = Q_{\pi_{\theta}}(s,a) - V_{\pi_{\theta}}(s)$.

1: 定义 V 函数逼近器 $V_{\mathbf{w}}(s)$, 策略逼近器 $\pi_{\theta}(s,a)$, 初始化 \mathbf{w} ,

heta, $s_t = s_0$, t = 0repeat 采样动作 $a_t \sim$

epeat
$$\begin{split} & \kappa_{k} \text{样 动作} \ a_{t} \sim \pi_{\theta}(s_{t}, a_{t}) \ \text{并 执行, 观测} \ r_{t+1}, s_{t+1} \\ & \delta = r_{t+1} + \gamma V_{\mathbf{w}}(s_{t+1}) - V_{\mathbf{w}}(s_{t}) \\ & \mathbf{w} \leftarrow \mathbf{w} + \beta \delta \nabla_{\mathbf{w}} V_{\mathbf{w}}(s_{t}) \end{split}$$

 $\theta \leftarrow \theta + \alpha \delta \nabla_{\theta} \log \pi_{\theta}(s_t, a_t)$

自然梯度:最速上升方向不再根据欧式距离

决定。而根据 J 的策略概率分布π决定 KL 散 度在概率分布空间上描述两个分布的距离。

确定型 AC:以上算法为有限动作集 MDPs 即 $a \sim \pi_{\theta}(s)$. 本算法的策略则是确定性的

 $a=\pi_{\theta}(s)$.动作的好坏又 Q(s, $\pi_{\theta}(s)$)反应.

1: 定义 Q 函数逼近器 $Q_{\mathbf{w}}(s,a)$, 策略逼近器 $\pi_{\theta}(s,a)$, 初始化 \mathbf{w}

 $\begin{array}{ll} \theta, s_{t} = \infty, s - \omega \\ 2 \cdot \text{repeat} \\ 3 \cdot \mathcal{R}H \sin \theta \cdot a_{t} = \pi_{\theta}(s_{t}) + N_{t} \not + \mathcal{W}h_{t}, \mathcal{R} \mathcal{B}(\tau_{t+1}, s_{t+1}) \\ 4 \cdot \mathcal{H}\# \text{ TD } \not \in \mathbb{R}^{L}, \delta = \tau_{t+1} + \gamma Q_{\mathbf{w}}(s_{t+1}, \pi_{\theta}(s_{t+1})) - Q_{\mathbf{w}}(s_{t}, a_{t}) \\ 5 \cdot \mathcal{L}\# \text{ Critic: } \mathbf{w} \leftarrow \mathbf{w} + \beta \delta \nabla_{\mathbf{w}} Q_{\mathbf{w}}(s_{t}, a_{t}) \\ 6 \cdot \mathcal{L}\# \text{ Actor: } \theta \leftarrow \theta + \alpha \nabla_{\alpha} Q_{\mathbf{w}}(s_{t}, a_{t}) \Big|_{\alpha = \pi_{\theta}(s_{t})} \nabla_{\theta} \pi_{\theta}(s_{t}) \\ \dots & \ddots & \ddots & \ddots & \ddots \\ \end{array}$

5: 更新权重 $\mathbf{w} \leftarrow \mathbf{w} + \alpha \delta_t \nabla_{\mathbf{w}} \hat{Q}(s_t, a_t, \mathbf{w})$