

$$P(\theta) \propto \prod_{i=1}^n P(s_i, a_i | s_{i-1}, a_{i-1})$$

这个算法, 是openai的默认算法. 所以特别叫. 参考<https://zhuanlan.zhihu.com/p/86306874>
自己写过程, 尽量把所有细节都写上.

$$\tau = \{s_1, a_1, s_2, a_2, \dots\}$$

$$\theta \xrightarrow{\text{取定后}} \pi \xrightarrow{\text{取定}} \tau \quad \text{给定一个 } s_1, \text{ 就会得到一个序列}$$

所以说序列跟 θ 有关. 给定 θ 后才能生成这个序列, 利用bg算法求解更新theta, 所以是on-policy的.

obj-func:

$$R(\tau) \triangleq \sum_{t=1}^T r_t$$

$$\bar{R}_\theta = \sum R(\tau) P_\theta(\tau)$$

$$P_\theta(\tau) \triangleq P(s_1) P_\theta(a_1 | s_1) P_\theta(s_2 | a_1, s_1) \dots$$

注意的一点是, 虽然每一个 θ 取定后 π 取定, 这个 π 对于给定的 s_t , 返回的 a_t 不是一个确定的值, 而是一个action的分布. 这样我们学习 θ , 让他改变这些个分布, 从而进行学习.

下面我们求

\bar{R}_θ

的梯度.

$$\begin{aligned}\nabla \bar{R}_\theta &= \sum_{\tau} R(\tau) \nabla p_\theta(\tau) \\ &= \sum_{\tau} R(\tau) p_\theta(\tau) \nabla (\log p_\theta(\tau)) \\ &= E_{\tau \sim p_\theta(\tau)} R(\tau) \nabla (\log p_\theta(\tau))\end{aligned}$$

蒙特采样 = $\frac{1}{N} \sum_{n=1}^N R(\tau^n) \nabla \log p_\theta(\tau^n)$

蒙特采样法求

改进为ppo 的off-learn算法

$$\nabla \bar{R}_\theta = E_{\tau \sim p_\theta(\tau)} [R(\tau) \nabla \log p_\theta(\tau)]$$

$$\begin{aligned}\nabla \bar{R}_\theta &= E_{\tau \sim p_{\theta'}(\tau)} \left[\frac{p_\theta(\tau)}{p_{\theta'}(\tau)} p_{\theta'}(\tau) \nabla \log p_{\theta'}(\tau) \right] \\ &= E_{\tau \sim p_{\theta'}(\tau)} \frac{\nabla p_\theta(\tau)}{p_{\theta'}(\tau)} A_{\theta'}(\tau)\end{aligned}$$

上面式子同事对 θ 等式两边做积分:

$$\bar{R}_\theta = E_{\tau \sim p_{\theta'}(\tau)} \frac{p_\theta(\tau)}{p_{\theta'}(\tau)} A_{\theta'}(\tau)$$

DQN:

为什么使用target net.

李宏毅老师说的是如果你的目标网络也是浮动的话, 收敛会波动, 收敛速度会变慢. 所以使用target net来 每隔多少个周期再更新权重会保证收敛速度.

off-policy: 里面的记忆库用的不是pi来生成的用来, 做pi的训练也没关系. 因为学的是reward. 还是有学习价值的.

李宏毅的double dqn也讲得非常细, 细节都到位了.

dqn的缺陷是没法计算continuous action.

$$a = \underset{a}{\operatorname{argmax}} Q(s, a)$$

这个方程对于连续情况没法解.

解决方法:

1. 抽样一堆action带入, 选择最优的, 也就是离散化.
2. 当作最优问题做梯度下降. 解a. --- 计算量挺大, 也不一定收敛.
3. 特别再设计一个网络来解这个方程.



最好的学习方法是, 讲给自己听, 或者讲给别人听. 然后看细节是否都能过去. 一个数学家提出的学习方法.

