

计算第3个批量，即当 $i = (k-1)M \sim kM$ ，在所有随机状态中从第 $(k-1)M$ 个到第 kM 个属于某状态聚集 (x, a) 的状态的个数

$$M_k(x, a, y) = \sum_{i=(k-1)M}^{kM} \chi_i(x, a, y)$$

计算第k个批量，即当 $i = (k-1)M \sim kM$ ，在所有随机状态中从第 $(k-1)M$ 个到第 kM 个属于某状态聚集 (x, a, y) 的状态的个数

计算第1个批量，即当 $i = 0 \sim M$ ，在所有随机状态中从第0个到第M个属于某状态聚集 (x, a, y) 的状态的个数

计算第2个批量，即当 $i = 2 \sim 2M$ ，在所有随机状态中从第M个到第2M个属于某状态聚集 (x, a, y) 的状态的个数

计算第3个批量，即当 $i = 2M \sim 3M$ ，在所有随机状态中从第2M个到第3M个属于某状态聚集 (x, a, y) 的状态的个数

……

计算第k个批量，即当 $i = (k-1)M \sim kM$ ，在所有随机状态中从第 $(k-1)M$ 个到第 kM 个属于某状态聚集 (x, a, y) 的状态的个数

Then the Q-value of (x, a) after the k^{th} batch is given by:

$$\begin{aligned} Q_{k+1}(x, a) &= (1 - M_k(x, a)\alpha_k(x, a))Q_k(x, a) + M_k(x, a)\alpha_k(x, a) \left[\frac{R_k^a(x)}{M_k(x, a)} + \gamma \sum_{y \in Y} \frac{M_k(x, a, y)}{M_k(x, a)} \max_a' Q_k(y, a') \right] \\ &= Q_k(x, a) + M_k(x, a)\alpha_k(x, a) \left[\frac{R_k^a(x)}{M_k(x, a)} + \gamma \sum_{y \in Y} \frac{M_k(x, a, y)}{M_k(x, a)} \max_a' Q_k(y, a') - Q_k(x, a) \right] \end{aligned}$$

在第k个批量中，learning rate等于 $M_k(x, a)\alpha_k(x, a)$ ，相当于一次同时处理了所有随机状态里面从 $(k-1)M$ 到 kM 属于某聚集状态的多个随机状态，效率提高了 $M_k(x, a)$ 倍，所以学习率乘以 $M_k(x, a)$ 。而作为第k个批量打包来说，这个“包”里的每个聚集的回报需要除以 $M_k(x, a)$ ，而a动作对应的下个状态的聚集 (y, a') 有 $M_k(x, a, y)$ 个，所以需要乘以 $M_k(x, a, y)$ 再除以 $M_k(x, a)$ ，即乘以某个聚集中a动作对应的下个状态聚集Y发生的概率。根据下面经典的Q-learning update rule和前面有关的说明，可以推导出上面有关批量Q值计算的结果。

Q-learning: Off-policy TD Control

One of the early breakthroughs in reinforcement learning was the development of an off-policy TD control algorithm known as Q-learning (Watkins, 1989) defined by

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

关于Q-learning，核心是采取某种策略使得Q值越大越好。在当前状态下更新当前状态的Q值，预测接下来发生的某个动作导致的下一个状态（这个动作可能与当前动作不同，所以下一个状态也有可能不同）所带回来的回报，从估计角度看， $Q(S_t, A_t)$ 是当前状态的Q值；从现实角度看，回报在当前状态并没有得到，但是如果达到下一个状态了，回报就是实际存在的，既然是“預判”，所以下个状态的Q最大值需要乘以一个衰减系数再加上回报 $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$ ，估计与现实的差距与学习效率有关，学习效率快，则更新当前状态的Q值速度越快，所以差距需要乘以学习效率。

In this case, the learned action-value function, Q, directly approximates q_* , the optimal action-value function, independent of the policy being followed. This dramatically simplifies the analysis of the algorithm and enabled early convergence proofs. The policy still has an effect in that it determines which state-action pairs are visited and updated. However, all that is required for correct convergence is that all pairs continue to be updated. This is a minimal requirement in the sense that any method guaranteed to find optimal behavior in the general case must require it. Under this assumption and a variant of the usual stochastic approximation conditions on the sequence of step-size parameters, Q has been shown to converge with probability 1 to q_* .

Q-learning (off-policy TD control) for estimating $\pi=\pi_*$

```
Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$   

Initialize  $Q(s, a)$ , for all  $s \in S^*, a \in A(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$   

Loop for each episode:  

    Initialize  $S$   

    Loop for each step of episode:  

        Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)  

        Take action  $A$ , observe  $R, s'$   

         $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$   

         $S \leftarrow S'$   

    until  $S$  is terminal
```

<Reinforcement Learning, An Introduction> Richard S. Sutton and Andrew G. Barto

令 \bar{Q} 是公式(2)的解：

https://blogs.cuit.columbia.edu/zp2130/reinforcement_learning_with_soft_state_aggregation/

$$\bar{Q}(x, a) = \bar{R}^a(x) + \gamma \sum_{y \in Y} \bar{P}^a(x, y) \max_{a' \in A} \bar{Q}(y, a')$$

定义

$$F_k(x, a) = \frac{R_k^a(x)}{M_k(x, a)} + \gamma \sum_{y \in Y} \frac{M_k(x, a, y)}{M_k(x, a)} \max_{a'} Q_k(y, a') - \bar{Q}(x, a) \quad (\text{Define } F_k)$$

把上面推导第k个批量的Q值的公式复制到这里方便对比看定义的 $F_k(x, a)$ 有什么特点：和下面公式后面方括号里的内容几乎一样，当然，从强化学习概念的角度看，严格来说，因为是第k个批量，所以下面公式里Q值并不是所有随机状态序列里收敛的那个Q值，但可以基于两者高度相似性作为推导公式的入手点。

事实上， F_k 即是TD error。

$$\delta_t \doteq R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$

"Finally, note that the quantity in brackets in the TD(0) update is a sort of error, measuring the difference between the estimated value of S_t and the better estimate $R_{t+1} + \gamma V(S_{t+1})$. This quantity, called the **TD error**, arises in various forms throughout reinforcement learning.

$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

TD(0), or one-step TD, because it is a special case of the TD(λ) and n-step TD methods developed in Chapter 12 and Chapter 7.

Notice that the TD error at each time is the error in the estimate *made at that time*. Because the TD error depends on the next state and next reward, it is not actually available until one time step later. That is, δ_t is the error in $V(S_t)$, available at time $t+1$.

$$\begin{aligned} Q_{k+1}(x, a) &= (1 - M_k(x, a)\alpha_k(x, a))Q_k(x, a) + M_k(x, a)\alpha_k(x, a) \left[\frac{R_k^a(x)}{M_k(x, a)} + \gamma \sum_{y \in Y} \frac{M_k(x, a, y)}{M_k(x, a)} \max_{a'} Q_k(y, a') \right] \\ &= Q_k(x, a) + M_k(x, a)\alpha_k(x, a) \left[\frac{R_k^a(x)}{M_k(x, a)} + \gamma \sum_{y \in Y} \frac{M_k(x, a, y)}{M_k(x, a)} \max_{a'} Q_k(y, a') - Q_k(x, a) \right] \\ &= Q_k(x, a) + M_k(x, a)\alpha_k(x, a) \left[\frac{R_k^a(x)}{M_k(x, a)} + \gamma \sum_{y \in Y} \frac{M_k(x, a, y)}{M_k(x, a)} \max_{a'} Q_k(y, a') - Q_k(x, a) \right] \end{aligned}$$

如果能够证明方括号里面的值是极小的，即
 $\forall \epsilon > 0, \exists M_\epsilon < \infty$ such that $\epsilon_k^M < \epsilon$ with probability $1 - \epsilon$. $Q_k(x, a) \rightarrow Q_\infty(x, a)$, where $|Q_\infty(x, a) - \bar{Q}(x, a)| \leq C\epsilon$

则能证明Q值收敛于某个值。

为了说明 $F_k(x, a)$ 有个极小值，论文推导出了如下形式，我把这形式称为 $(\text{Quantity}_F)_k$

$$F_k(x, a) = \gamma \sum_y \frac{M_k(x, a, y)}{M_k(x, a)} [V_k(y) - \bar{V}(y)] + \left[\frac{R_k^a(x)}{M_k(x, a)} - R_{0, \infty}(x) \right] + \gamma \sum_y \left[\left(\frac{M_k(x, a, y)}{M_k(x, a)} - P_{0, \infty}^a(x, y) \right) \bar{V}(y) \right] \quad (\text{Quantity}_F_k)$$

上面公式中可能有排版错误，分母中X应该是x。

$$\begin{aligned} F_k(x, a) &= \gamma \sum_y \frac{M_k(x, a, y)}{M_k(x, a)} [V_k(y) - \bar{V}(y)] + \left[\frac{R_k^a(x)}{M_k(x, a)} - R_{0, \infty}(x) \right] + \gamma \sum_y \left[\left(\frac{M_k(x, a, y)}{M_k(x, a)} - P_{0, \infty}^a(x, y) \right) \bar{V}(y) \right] \\ &= \gamma \sum_y \frac{M_k(x, a, y)}{M_k(x, a)} [V_k(y)] + \left[\frac{R_k^a(x)}{M_k(x, a)} - R_{0, \infty}(x) \right] + \gamma \sum_y \left[\left(\frac{M_k(x, a, y)}{M_k(x, a)} - P_{0, \infty}^a(x, y) \right) \bar{V}(y) \right] \\ &= \frac{R_k^a(x)}{M_k(x, a)} + \gamma \sum_y \frac{M_k(x, a, y)}{M_k(x, a)} [V_k(y)] - R_{0, \infty}(x) + \gamma \sum_y \left[\left(\frac{M_k(x, a, y)}{M_k(x, a)} - P_{0, \infty}^a(x, y) \right) \bar{V}(y) \right] \\ &= A - B \\ A &= \frac{R_k^a(x)}{M_k(x, a)} + \gamma \sum_y \frac{M_k(x, a, y)}{M_k(x, a)} [V_k(y)] \end{aligned}$$

$$B = R_{0, \infty}(x) + \gamma \sum_y \left[\left(\frac{M_k(x, a, y)}{M_k(x, a)} - P_{0, \infty}^a(x, y) \right) \bar{V}(y) \right]$$

因为在随机过程产生的所有随机状态组成的序列中存在某聚集(x,y)的概率和回报都收敛：

$$P_{0,\infty}^a(x, y) = \bar{P}^a(x, y)$$

and

$$P_{0,\infty}^a(x) = \bar{R}^a(x)$$

根据前面 \bar{Q} 是公式 (2) 的解：

$$\bar{Q}(x, a) = \bar{R}^a(x) + \gamma \sum_{y \in Y} \bar{P}^a(x, y) \max_{a' \in A} \bar{Q}(y, a')$$

如果

$$V_k(x) = \max_a Q_k(x, a) : \text{所有随机状态中的某聚集中第k批量的Q值最大值}$$

则可知： $V(y) = \max_a Q(y, a)$ 则可知： $\bar{V}(y) = \max_a \bar{Q}(y, a)$

推导得：

$$F_k(x, a) = \frac{\bar{R}_k^a(x)}{M_k(x, a)} + \gamma \sum_y \frac{M_k(x, a, y)}{M_k(x, a)} \max_a Q(y, a') - \bar{Q}(x, a) \quad (\text{Define_} F_k)$$

当然，在论文中从公式(Define F_k)到(Quantity F_k)的顺序正好与这里的推导相反。这里是(Quantity F_k)推导到(Define F_k)。所以我们可以接着分析 F_k 的大小

$$F_k(x, a) = \gamma \sum_y \frac{M_k(x, a, y)}{M_k(x, a)} [V_k(y) - \bar{V}(y)] + \left[\frac{\bar{R}_k^a(x)}{M_k(x, a)} - R_{0,\infty}(x) \right] + \gamma \sum_y \left[\left(\frac{M_k(x, a, y)}{M_k(x, a)} - P_{0,\infty}^a(x, y) \right) \bar{V}(y) \right] \quad (\text{Quantity_} F_k)$$

Norms and Metrics

$$\begin{aligned} \text{http://www.u.arizona.edu/~mwalker/econ519/Econ519LectureNotes/} \\ \text{Norms\&Metrics.pdf} \\ \text{http://www.u.arizona.edu/~mwalker/econ519/Econ519LectureNotes/} \\ \text{Norms\&Metrics.pdf} \end{aligned}$$

Euclidean Norm

欧几里得范数

$$\|x\| = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$$

Euclidean metric (l² metric)

$$\begin{aligned} d_{l^2}((x_1, \dots, x_n), (y_1, \dots, y_n)) &:= \sqrt{(x_1 - y_1)^2 + \dots + (x_n - y_n)^2} \\ &= \left(\sum_{i=1}^n (x_i - y_i)^2 \right)^{1/2} \end{aligned}$$

If n=2, then $d_{l^2}((1,6),(4,2)) = \sqrt{(3^2 + 4^2)} = 5$. This metric corresponds to the geometric distance between the two points $(x_1, x_2, \dots, x_n), (y_1, y_2, \dots, y_n)$, as given by Pythagoras' theorem.毕达哥拉斯定理

Euclidean distance between two points x, x' $\in \mathbf{R}^n$

$$d(x, x') := \|x, x'\|$$

Taxicab metric (l¹ metric)

$$\begin{aligned} d_{l^1}((x_1, x_2, \dots, x_n), (y_1, y_2, \dots, y_n)) &:= |x_1 - y_1| + \dots + |x_n - y_n| \\ &= \sum_{i=1}^n |x_i - y_i| \end{aligned}$$

Because it models the distance a taxi-cab would have to traverse to get from one point to another if the cab was only allowed to move in cardinal directions (north, south, east, west) and not diagonally.出租车只能沿着东西走向或者南北走向的街道开，不能从房子上或者从房子里开过去（虽然两点间直线最短）。

$$d_{l^2}(x, y) \leq d_{l^1}(x, y) \leq \sqrt{n} d_{l^1}(x, y)$$

最简单形式的理解可以是直角三角形直角边之和大于斜边。

$$\begin{aligned} \|F_k(x, a)\| &\leq \gamma \|V_k - \bar{V}\| + \left\| \frac{\bar{R}_k^a(x)}{M_k(x, a)} - R_{0,\infty}^a(x) \right\| + \gamma \left\| \sum_y \left[\frac{M_k(x, a, y)}{M_k(x, a)} - P_{0,\infty}^a(x, y) \right] \right\| \\ &\leq \gamma \|V_k - \bar{V}\| + C \epsilon_k^M \end{aligned}$$

where

$$\epsilon_k^M = \max \left\{ \left| \frac{\bar{R}_k^a(x)}{M_k(x, a)} - R_{0,\infty}^a(x) \right|, \gamma \left| \sum_y \left(\frac{M_k(x, a, y)}{M_k(x, a)} - P_{0,\infty}^a(x, y) \right) \right| \right\}$$

 ϵ_k^M 是 l¹ metric 大于等于 norm (l¹ metric)

也就是说

$$\|F_k(x, a)\| < \text{一个极小值}$$

根据前面有关定义 F_k 的原因，为什么要证明 F_k 小于一个极小值的原因的说明，也就证明了 Q 值收敛于某个值。

简单介绍有关基于贝叶斯理论推导的依靠采样策略的聚类化的Q值计算方法及状态聚合(或状态聚类)与状态聚类分别在概率论与数理统计和强化学习的概念和区别

现实世界里，以机械手臂（甚至具有传感功能的“八爪鱼”）为例，可能有无数状态(State)和动作(Action)（我个人理解“八爪鱼”触手可能具有的传感功能包括：温度感应，压力感应，等等），实际工作中要求机械手臂的灵度越高越好，例如1000+自由度，从强化学习(Reinforcement Learning, RL)中数学角度看，状态空间(State Space)非常大，所以普遍认为，用更简洁的表达方式而不是查表(lookup table)来表示(Scaling)强化学习到现实问题是非常关键的。在这片论文发表之前，几乎所有的强化学习理论用查找表研究该领域的问题，而这篇论文针对结合近似函数(function approximation)和强化学习时产生的问题，提出解决方案，比如从依靠采样策略(sampling policy)产生聚类(clustering)端来计算Q值，而不是从状态端出发的查表方法。因为作者提出有别于聚类的聚类理论(软聚类，即某状态通过聚类概率而非完全从属于某聚类)，用贝叶斯理论将有关“状态”条件下“聚类”的概率转换为有关“聚类条件下的状态”概率计算Q值，也许可以这样理解：这样解决了只能从状态端出发的查找表来计算Q值的小算量过大的问题，也就是说可以从聚类端出发的策略计算Q值，降低了计算量，从而提高了计算效率，节省计算时间和成本，为此作者通过缜密的逻辑思路和精准的用词，层层推进，教科书式地提出了以下几个研究成果：

1. 一个近似函数(function approximator, FA) (在这个post中权且翻译成“近似者函数”，与通常的近似函数区分)，这个函数是基于状态聚合(state aggregation)²的一个简单的扩展，也就是说状态以聚类概率(clustering probability)³属于某个聚类，故论文名：“软”状态聚合(soft state aggregation)⁴；
2. 一个对于任意的但是固定的软状态聚合的强化学习的收敛理论。既然提出一个函数，而在“聚类”的层面上强化学习的问题可能属于非马尔可夫过程，所以让明这个近似者函数的收敛性是必要的，本文依靠率收敛数说明提出的软状态聚类的强化学习从数学概率论与数理统计的角度看也是正确的；
3. 关于这种软状态聚合对在线强化学习的一种新颖的直观理解，这里利用了贝叶斯公式简单形式和全概率公式结合形式，将有关聚类概率(状态发生的条件下聚类的频率)的计算动作值函数Q值转换为有关状态发生的概率率的计算动作值函数Q值，也可以将这种思路理解为是某种逆向概念：正是因为通过采样策略(sampling policy)来计算Q值，而不是查表(lookup table case)，(在本文关键字“贝叶斯公式”，有多个地方解释了该术语及真在论文中的应用)
4. 一种新的探索式的适应性(adaptive)软状态聚合算法，这个算法通过探索软状态聚合的非离散本质，找到改良的简洁的软状态聚合的表达式。

总之，这4步内容互相关联，教科书式地展开，层层递进，用词精准，逻辑缜密。首先提出一个函数，接着证明收敛性，最后说明此函数的应用。

前面提到海量状态空间，为了研究方便且有效率，有必要简化表达它，**状态聚合**²(State Aggregation)的概念就被提出来了。这篇论文提出了状态以聚类概率从属于多个聚类的形式，有必要化近似函数的形式，也就是把状态分组打包，每个状态聚类与通常常用的聚类方法，论文中用cluster，而非aggregation来表示聚类。而我在这篇post中将state aggregation翻译成状态聚合，也是为了避免在数学上混淆两者。

1. **近似函数(function approximation)**: 将状态空间(和所有的模型参数)以一种更简洁的方式近似化，比如线性(向量化)，神经网络，决定树，等等。
2. **状态聚合(state aggregation)**: 一个普遍使用的简洁表达状态的方式，它是一种简单的规范化近似函数的形式，也就是把状态分组打包，每个组有一个估计值(权重向量里的一个分量)，一个状态的值就是组里的分量。此状态更新，只要单独更新那个分量。
3. **聚类概率 clustering probability**: 每一个状态以某一个概率率从属于某一个聚类，而这个概率叫聚类概率(clustering probability)。我觉得可能因为作者不想在数学层面混淆两类概念，为了区分aggregation，而另取了名字cluster，其二者区别在于cluster强调状态以某个聚类概率(clustering probability)而非100%完全从属于某一个聚类(cluster)，换句话叫aggregation，换句说聚类概率不是这个概率率的“概率”，而是指允许每个状态从属于多个聚类，而cluster有。我把aggregation翻译成聚类，cluster翻译成聚类，也是为了避免在数学的“概率”层面区别这两者，虽然本质上含义可能差不多。
4. **“软”状态聚合(soft state aggregation)**: 状态以聚类概率而非完全从属于某一个聚类cluster，也就是说允许每个状态只从属于一个聚类cluster。一个聚类的值可以通过与聚类概率的正比关系推广或规范generalize到所有的状态。也被称作“状态聚类”(State Clustering)。

State aggregation is a simple form of generalizing function approximation in which states are grouped together, with one estimated value (one component of the weight vector w) for each group. The value of a state is estimated as its group's component, and when the state is updated, that component alone is updated.

Function Approximation: approximate the state space (and all model parameters) with a more compact one!

- Reduction in # of states (computation and space)
- More importantly: generalization to unseen states!

Types of (value / action-value) function approximation:

- Linear
- Neural network
- Decision tree

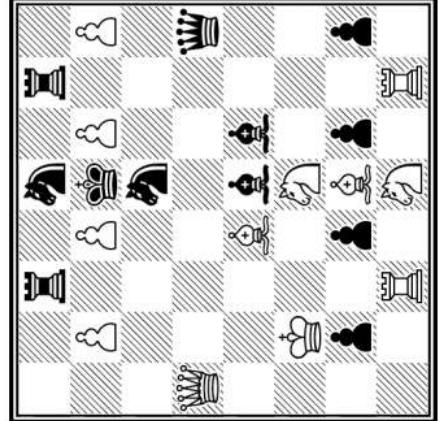
...

Function approximation
Finding optimal $\theta \rightarrow$ knowledge of value for ALL states!

$$v_\theta(s) = \theta_1 x_1(s) + \theta_2 x_2(s) + \dots + \theta_n x_n(s) = \theta^T x(s)$$

10⁴⁰ states are mapped to linear function over n "important" features, i.e.

1. Number of white pieces – black pieces
2. Distance between kings
3. Etc.



Function approximation idea – generalization and efficiency

- Gradient descent approximation to estimate value/Q functions
- gradient descent to optimize the optimal Q-function directly

402-lec20 (<http://blogs.cuit.columbia.edu/zp2130/files/2019/02/402-lec20.pdf>)**Markov Decision Process (MDP)**

Markov Reward Process, definition:

- Tuple (S, P, R, A, γ) where
 - S = states, including start state
 - A = set of possible actions
 - P = transition matrix $P_{SS'}^a = P_r[S_t+1 = s' | S_t = s, A_t = a]$
 - R = reward function, $R_S^a = E[R_{t+1} | S_t = s, A_t = a]$
 - $\gamma \in [0, 1]$ = discount factor

• Return

$$G_t = \sum_{i=1 \text{ to } \infty} R_{t+i} \gamma^{i-1}$$

• Goal: take actions to maximize expected return

Policies

The Markovian structure \rightarrow best action depends **only** on current state.

- Policy = mapping from state to distribution over actions

$$\pi : S \times A \mapsto \Delta(A), \pi(a | s) = P_r[A_t = a | S_t = s]$$

• Given a policy, the MDP reduces to a Markov Reward Process

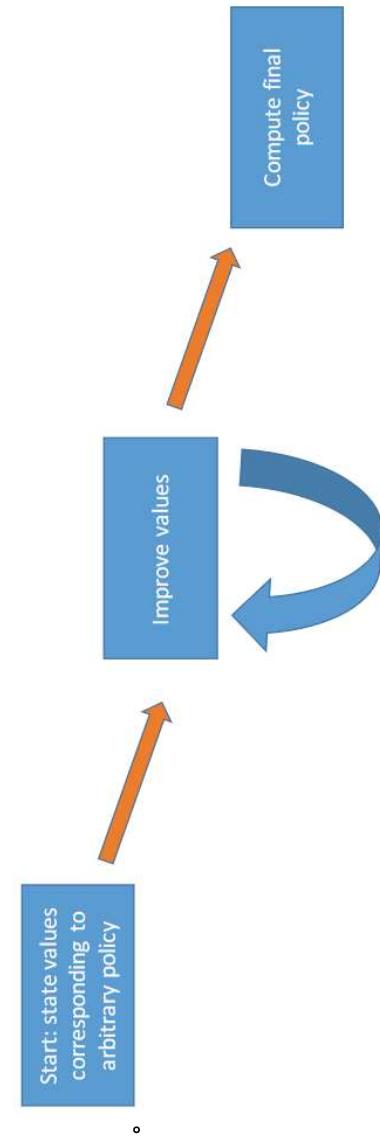
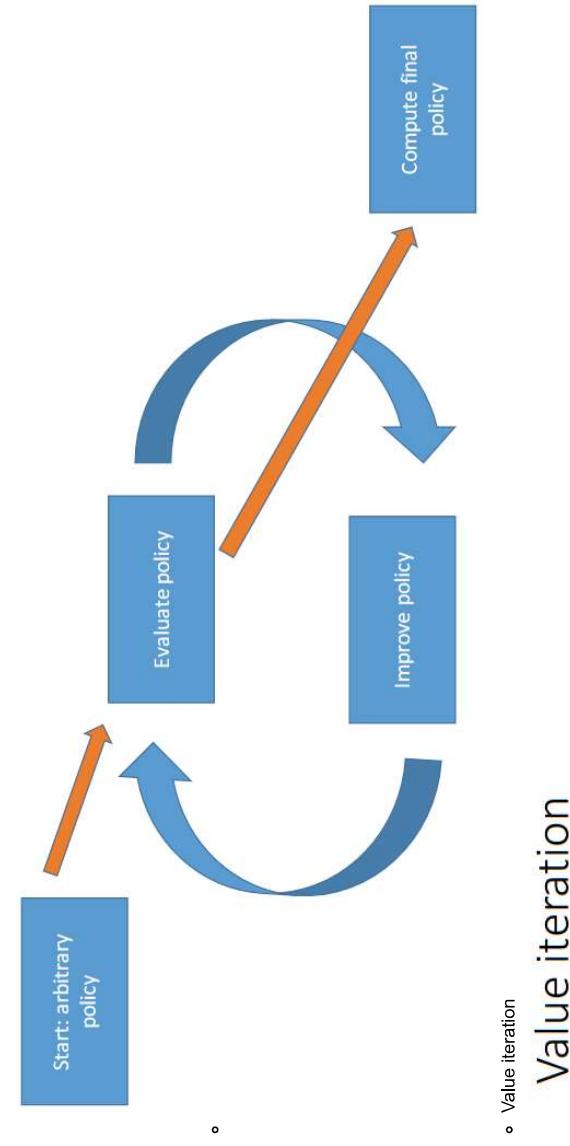
Bellman optimality equations

- Bellman equation:
 - $v_*(s) = \max_a v_a(s)$
 - implies Bellman optimality equations:

$$q_*(s, a) = R_S^a + \gamma \sum_{S'} P_{SS'}^a \max_{a'} \{q_*(s', a')\}$$

$$v_*(s) = \max_a \left\{ R_S^a + \gamma \sum_{S'} P_{SS'}^a v_x(S') \right\}$$

- Iterative methods based on the Bellman equations: dynamic programming
 - Policy iteration

Policy iteration**Temporal-Difference Learning**

A simple every-visit Monte Carlo method suitable for nonstationary environments is

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]$$

$$NewEstimate \leftarrow OldEstimate + StepSize [Target - OldEstimate]$$

The expression [Target – OldEstimate] is an error in the estimate.

where G_t is the actual return following time t, and α is a constant step-size parameterThe target for the Monte Carlo update is G_t .**TD(0), or one-step TD, because it is a special case of the TD(λ) and n-step TD methods**

The simplest TD method makes the update

$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

immediately on transition to S_{t+1} and receiving R_{t+1} .The target for the TD update is $R_{t+1} + \gamma V(S_{t+1})$

Tabular TD(0) for estimating v_{π} Input: the policy π to be evaluatedAlgorithm parameter: step size $\alpha \in (0, 1]$ Initialize $V(s)$, for all $s \in S^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop for each episode:

Initialize S

Loop for each step of episode:

 $A \leftarrow$ action given by π for s Take action A , observe R , s' ~~**Invalid Equation**~~ $S \leftarrow S'$ until s is terminal

<Reinforcement Learning, An Introduction> Richard S. Sutton and Andrew G. Barto

2. RL and Soft State Aggregation

这部分应用了定理1 的结论来提供两个例子的收敛性。

1. 用Q-learning和近似值函数(FA)解方程

2. 用TD(0)和近似值函数(FA)决定一个固定决策的价值函数。

在线强化学习，对于learning agent而言，transition probability和payoff function都是未知的。

2.1 Case

Q-learning and Fix Soft State Aggregation

这部分作者开发出了一个在聚集中分别对Q-learning的Bellman equation”。作者假设agent按照稳定的随机策略，也就是给每一个状态指定一个非零概率，即在每一个状态下执行每个动作都是非零概率。

1. 用Q-learning和近似值函数(FA)解方程

2. 用TD(0)和近似值函数(FA)决定一个固定决策的价值函数。

在线强化学习，对于learning agent而言，transition probability和payoff function都是未知的。

Corollary 1

利用了贝叶斯公式的扩展形式，即贝叶斯公式简单形式和全概率公式的结合形式，将有~~关~~聚集概率（状态发生的条件下聚集的概率）的计算动作值函数Q值转换为有关在聚集中聚集的条件下的概率的计算动作值函数Q值，也可以将这种思路理解为是某种逆向概念。

为了对比，将前文提到的计算动作值函数的Q值的公式和由贝叶斯公式转换而来的新的计算动作值函数的Q值的公式放在一起。

$$Q(x, a) = \bar{R}^a(x) + \gamma \sum_{y \in Y} \bar{P}^a(x, y) \max_{a' \in A} Q(y, a') \quad (2)$$

where

$$\begin{aligned} \bar{P}^a(x, y) &= P_{0, \infty}^a(x, y) \\ \bar{R}^a(x) &= R_{0, \infty}^a(x) \end{aligned}$$

where for all s , $P^{\pi}(s)$: steady-state probability of being in state s .

这是一个贝叶斯理论Bayes theorem，全概率公式law of total probability的结合形式，或者理解成Bayes' theorem Extended form。

分子分母都可以理解为有关条件概率的计算。首先，执行某个动作或者说成为某个状态 s 在一个概率，然后，状态 s 属于某聚集中存在一个聚集概率，所以总体看是先满足状态稳定的条件下再考虑此状态属于某个聚集，也就是说这是有关条件概率的计算。分子是在所有可能的下一步的状态条件下的聚集的条件概率分别与下一步状态概率乘积的总和，也就是聚集的概率，而分子当前状态的概率乘以在当前状态的条件下聚集的概率（条件概率），它们的结果是当前状态与某聚集的联合概率，分子分母的比值就是在聚集中聚集的条件下状态的概率，也就是说某聚集中存在某状态的概率，这与某状态从属于某聚集的概率（聚集概率，聚集概率也称条件概率，它是在某状态的条件下，属于聚集的概率）的概念互为逆向。接下来简单介绍贝叶斯公式及其扩展形式方便理解作者提出的“Bellman equation”。

条件概率(Conditional probability)

$$P(A | B) = \frac{P(A \cap B)}{P(B)}, \text{ if } P(B) \neq 0$$

$$P(B | A) = \frac{P(A \cap B)}{P(A)}, \text{ if } P(A) \neq 0$$

where $P(A \cap B)$ is 联合概率 joint probability of both A and B being true.

全概率公式(law of total probability):

$$\bar{R}^a(x) = \sum_s P^{\pi}(s | x) R^a(s)$$

$$\bar{P}^a(x, y) = \sum_s P^{\pi}(s | x) P^a(s, y)$$

for some partition $\{B_i\}$ of the sample space, the event space is given or conceptualized in terms of $P(B_i)$ 若事件 $B_1, B_2, B_3, \dots, B_n$ 是样本空间O的一个划分，则：

$$P(A) = \sum_{i=1}^n P(A \cap B_i)$$

and $P(A | B_j)$. It is then useful to compute $P(A)$ using the law of total probability. 又因为条件概率公式，可进一步得：

$$P(A) = \sum_{i=1}^n P(A | B_i) P(B_i)$$

$$V(x) = \sum_s P^{\pi}(s | x) \left[R^{\pi}(s) + \gamma \sum_y P^{\pi}(s, y) V(y) \right]$$

Proof:

TD(0)是马尔可夫决策过程Q-learning的特殊情况，也就是说每个状态只有一个简单（也可能是随机的）动作，也就是说不用考虑的多种可能性，即不存在比较各种动作下的Q值的最大值的情况，用V(y)代替公式(3)Max取值部分即可。

Case 2: TD(0) and Fixed Soft State Aggregation

贝叶斯理论 (Bayes' theorem)

1/15/25 4:55 PM

3. Adaptive State Aggregation

在聚集成数固定的前提下，找到好的聚类概率(clustering probabilities)来是高简洁代表 \bar{x} (compact representation)，从而获得更好的价值函数的近似值。在推理论和2里表明，RL在聚类空间(cluster space)可以找到零Bellman误差(zero Bellman error)的解。相应地，在状态空间，Bellman误差一般来说不说不为0。

总之，好的聚类化是以降低在马尔可夫决策过程状态中的Bellman误差为表现形式的。

聚类概率

$$P(x|s; \theta) = \frac{e^{\theta(x, s)}}{\sum_{x'} e^{\theta(x', s)}}$$

where

θ = weight between state s and cluster x .

Bellman error at state s given parameter θ (a matrix) is

$$\begin{aligned} J(s|\theta) &= V(s|\theta) - \left[R^\pi(s) + \gamma \sum_{s'} P^\pi(s, s') V(s'|\theta) \right] \\ &= \left[\sum_x P(x|s; \theta) V(x|\theta) \right] - \left[R^\pi(s, s) + \gamma \sum_{s'} P^\pi(s, s') \sum_x P(x|s'; \theta) V(x|\theta) \right] \end{aligned}$$

通过聚类概率将状态端 s 转换为聚类端 x

由原先的从单状态计算Q值，到现在从聚类计算Q值。

$$\frac{\partial J^2(s|\theta)}{\partial \theta(y, s)} = 2J(s|\theta) [P(y|s; \theta)(1 - \gamma P^\pi(s, s'))(V(y|\theta) - V(s|\theta))]$$

$$\frac{\partial J(s|\theta)}{\partial \theta(y, s)} = P(y|s; \theta)(1 - \gamma P^\pi(s, s'))(V(y|\theta) - V(s|\theta)) \quad (\text{How do I get this?})$$

Refer to Policy Gradient Methods for Reinforcement Learning with Function Approximation (https://blogs.cuit.columbia.edu/zp2130/policy_gradient_methods_for_reinforcement_learning_with_function_approximation/)

$$\frac{\partial P(x|s; \theta)}{\partial \theta(y, s)} = \frac{e^{\theta(x, s)} \sum_{x'} e^{(x', s)} - e^{\theta(x, s)} \sum_{x'} e^{(x', s)}}{(\sum_{x'} e^{(x', s)})^2} = 0 \quad (\text{Is it correct?})$$



useful information online:

<https://morvanzhou.github.io/tutorials/machine-learning/ml-intro/4-03-q-learning/> (<https://morvanzhou.github.io/tutorials/machine-learning/ml-intro/4-03-q-learning/>)

在看论文的过程中，看懂一点论文的状态动作值函数Q值越来越大。:)

edit (https://blogs.cuit.columbia.edu/zp2130/reinforcement_learning_with_soft_state_aggregation/edit)
Author: Z Pei (https://blogs.cuit.columbia.edu/zp2130/reinforcement_learning_with_soft_state_aggregation/)

Other posts

- Symbolic Netlist to Innovus-friendly Netlist (https://blogs.cuit.columbia.edu/zp2130/symbolic_netlist_to_innovus-friendly_netlist/)
- Finite-Sample Convergence Rates for Q-Learning and Indirect Algorithms (https://blogs.cuit.columbia.edu/zp2130/finite-sample_convergence_rates_for_q-learning_and_indirect_algorithms/)
- Solving H-Horizon, Stationary Markov Decision Problems In Time Proportional To Log(H) (https://blogs.cuit.columbia.edu/zp2130/paul_tseng_1990/)
- Randomized Linear Programming Solves the Discounted Markov Decision Problem In Nearly-Linear (Sometimes Sublinear) Run Time (https://blogs.cuit.columbia.edu/zp2130/randomized_linear_programming_solves_the_discounted_markov_decision_problem_in_nearly-linear_sometimes_sublinear_run_time/)
- KL Divergence (https://blogs.cuit.columbia.edu/zp2130/kl_divergence/)
- The Asymptotic Convergence-Rate of Q-Learning (https://blogs.cuit.columbia.edu/zp2130/the_asymptotic_convergence_rate_of_q-learning/)
- Hierarchical Apprenticeship Learning, with Application to Quadruped Locomotion (https://blogs.cuit.columbia.edu/zp2130/hierarchical_apprenticeship_learning_with_application_to_quadruped_locomotion/)
- Policy Gradient Methods (https://blogs.cuit.columbia.edu/zp2130/policy_gradient_methods/)
- Actor-Critic Algorithms for Hierarchical Markov Decision Processes (https://blogs.cuit.columbia.edu/zp2130/actor-critic_algorithms_for_hierarchical_markov_decision_processes/)
- Hierarchical Deep Reinforcement Learning: Integrating Temporal Abstraction and Intrinsic Motivation (https://blogs.cuit.columbia.edu/zp2130/hierarchical_deep_reinforcement_learning_integrating_temporal_abstraction_and_intrinsic_motivation/)