

1.

a)

we know that softmax can be denoted as:

$$k = \frac{e^{Qt(a)/t}}{\sum_{i=1}^n e^{Qt(i)/t}}$$

we can get the

$$\ln K = \frac{Qt(a)}{t} - \ln \left(e^{\frac{Qt(a)}{t}} \left(1 + \sum_{i=1}^n \frac{e^{Qt(i)/t}}{e^{Qt(a)/t}} \right) \right)$$

when $t \Rightarrow 0$, $e^{Qt(i)/t} \ll e^{\frac{Qt(a)}{t}}$. So $\frac{e^{Qt(i)/t}}{e^{Qt(a)/t}} \Rightarrow 0$, $\sum \frac{e^{Qt(i)/t}}{e^{Qt(a)/t}} \Rightarrow 0$

$$\text{So } \ln K \Rightarrow \frac{Qt(a)}{t} - \ln \left(e^{\frac{Qt(a)}{t}} \right) = \frac{Qt(a)}{t} - \frac{Qt(a)}{t} = 0;$$

So $k \Rightarrow 1$. Therefore, when $t \rightarrow 0$, it becomes the same as greedy action selection, always choose the biggest one.

b)

we know that softmax can be denoted as:

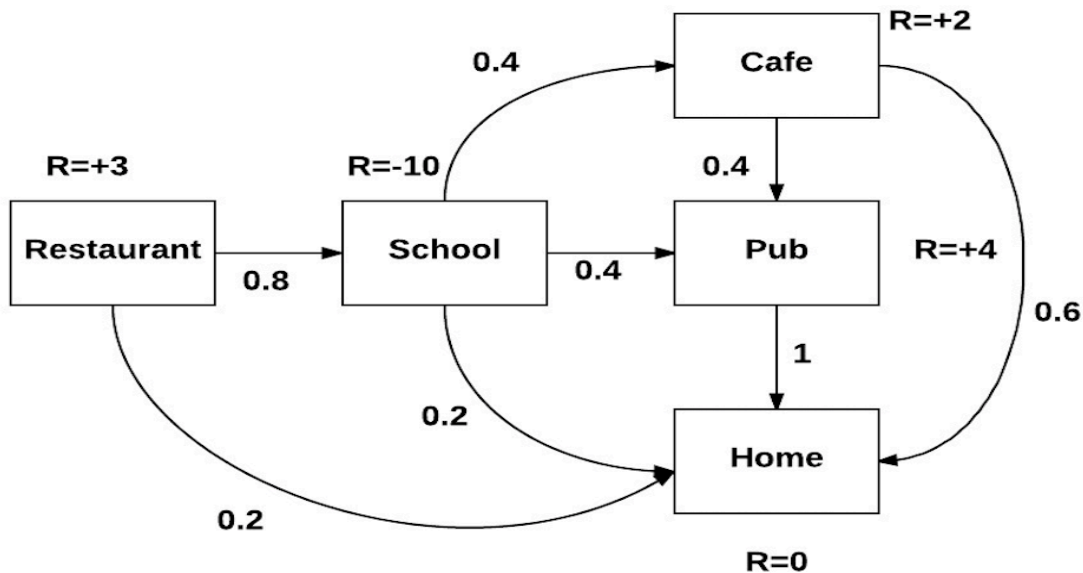
$$k = \frac{e^{Qt(a)/t}}{\sum_{i=1}^n e^{Qt(i)/t}}$$

So for two actions, the probability of either one can be denoted as:

$$P_i = \frac{e^{\frac{Q_i}{t}}}{e^{\frac{Q_1}{t}} + e^{\frac{Q_2}{t}}} = \frac{1}{1 + e^{-(Q_1 - Q_2)/t}}$$

So it becomes a Logistic function.

. 2. MRP



let's take Restaurant :0, School :1, Café: 2, Pub: 3, Home:4 .
 Transition Probability matrix is P, the Reward function is R.

$$P = \begin{bmatrix} 0.8 & 0.2 \\ 0.4 & 0.4 & 0.2 \\ 0.4 & 0.6 \\ 1 \\ 0 \end{bmatrix}$$

$$R = \begin{bmatrix} 3 \\ -10 \\ 2 \\ 4 \\ 0 \end{bmatrix}$$

According to Bellman Equation:

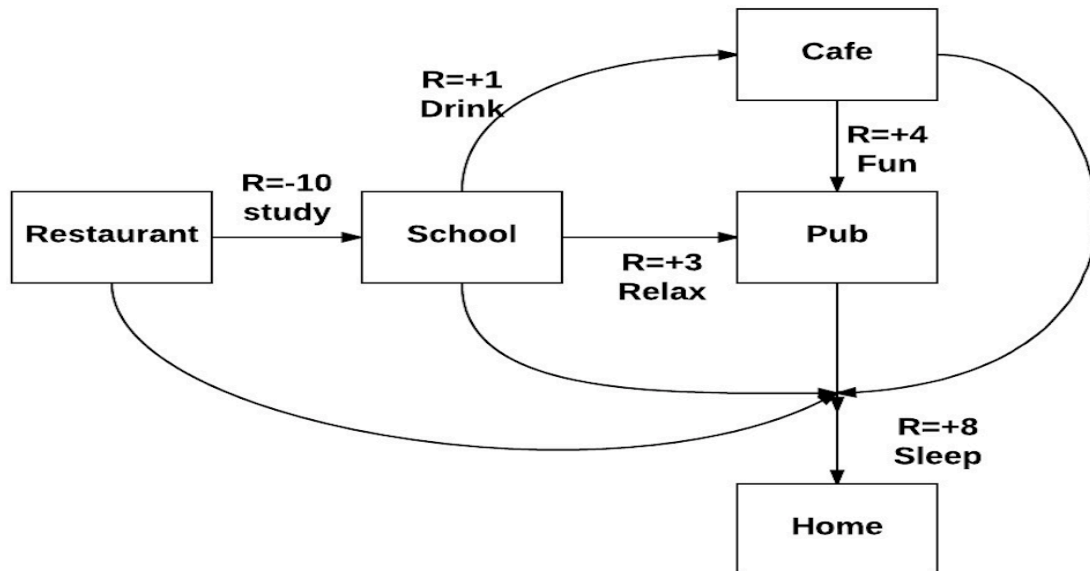
$$V = R + rPV$$

$$V = (I - rP)^{-1}R$$

$$(I - rP)^{-1} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0.4 & 0.56 & 1 \\ 0 & 0 & 1 & 0.4 & 1 \\ 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$V = \begin{bmatrix} 3 \\ -10 \\ 2 \\ 4 \\ 0 \end{bmatrix}$$

3.MDP



For our MDP, the $r = 1$, and has a stochastic policy. $\pi(a|s)$ is the same. Transition Probability matrix is P , the Reward function is R . let's take Restaurant :0, School :1, Café: 2, Pub: 3, Home:4 . Follow the rule:

$$\mathcal{R}_s^\pi = \sum_{a \in \mathcal{A}} \pi(a|s) \mathcal{R}_s^a$$

$$P^\pi = \begin{bmatrix} 0.5 & 0.5 \\ 0.33 & 0.33 & 0.33 \\ 0.5 & 0.5 \\ 1 \\ 0 \end{bmatrix} \quad R^\pi = \begin{bmatrix} -1 \\ 4 \\ 6 \\ 8 \\ 0 \end{bmatrix}$$

According to Bellman Equation:

$$V_\pi = (I - \gamma P^\pi)^{-1} \mathcal{R}^\pi$$

$$(I - \gamma P^\pi)^{-1} = \begin{bmatrix} 10.50 & 1.66 & 5.02 & 4.97 & 5.09 \\ 0 & 1 & 0.33 & 0.49 & 0.99 \\ 0 & 0 & 1 & 0.5 & 1 \\ 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \quad V^\pi = \begin{bmatrix} 3.997 \\ 9.994 \\ 10 \\ 8 \\ 0 \end{bmatrix}$$

4. Derive the weighted-average update rule (5.5) from (5.4). Follow the pattern of the derivation of the unweighted rule (2.3).

Let $C_{n+1} = C_n + W_{n+1}$, where $C_0 = 0$;

$$\begin{aligned} V_{n+1} &= \frac{\sum_{k=1}^n W_k * G_k}{\sum_{k=1}^n W_k} = \frac{W_n * G_n + V_n * (C_n - W_n)}{C_n} = V_n + \frac{W_n * G_n - V_n * W_n}{C_n} \\ &= V_n + \frac{W_n}{C_n} (G_n - V_n) \end{aligned}$$