

initialize, for $a = 1$ to k :

$$\begin{aligned} Q(a) &\leftarrow 0 \\ N(a) &\leftarrow 0 \end{aligned}$$

Loop forever:

$$\begin{aligned} A &\leftarrow \begin{cases} \arg\max_a Q(a) & \text{with probability } 1 - \varepsilon \\ \text{a random action} & \text{with probability } \varepsilon \end{cases} \\ R &\leftarrow \text{bandit}(A) \\ N(A) &\leftarrow N(A) + 1 \\ Q(A) &\leftarrow Q(A) + \frac{1}{N(A)} [R - Q(A)] \end{aligned}$$

$$\begin{aligned} Q_{n+1} &= (1 - \alpha)^n Q_1 + \sum_{i=1}^n \alpha (1 - \alpha)^{n-i} R_i \\ \sum_{n=1}^{\infty} \alpha_n(a) &= \infty \quad \text{and} \quad \sum_{n=1}^{\infty} \alpha_n^2(a) < \infty \\ A_t &\doteq \arg\max_a \left[Q_t(a) + c \sqrt{\frac{\ln t}{N_t(a)}} \right] \end{aligned}$$

Dynamics of the MDP

$$p(s', r | s, a) \doteq \Pr\{S_t = s', R_t = r \mid S_{t-1} = s, A_{t-1} = a\}$$

State-transition probabilities

$$p(s' | s, a) \doteq \Pr\{S_t = s' \mid S_{t-1} = s, A_{t-1} = a\} = \sum_{r \in \mathcal{R}} p(s', r | s, a).$$

Expected rewards state-action

$$r(s, a) \doteq \mathbb{E}[R_t \mid S_{t-1} = s, A_{t-1} = a] = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} p(s', r | s, a)$$

Expected rewards state-action-state

$$r(s, a, s') \doteq \mathbb{E}[R_t \mid S_{t-1} = s, A_{t-1} = a, S_t = s'] = \sum_{r \in \mathcal{R}} r \frac{p(s', r | s, a)}{p(s' | s, a)}.$$

State-value function for policy π

$$\begin{aligned} v_{\pi}(s) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right] \\ &= \mathbb{E}_{\pi}[R_{t+1} + \gamma G_{t+1} \mid S_t = s] \\ &= \sum_a \pi(a | s) \sum_{s'} \sum_r p(s', r | s, a) \left[r + \gamma \mathbb{E}_{\pi}[G_{t+1} | S_{t+1} = s'] \right] \\ &= \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) \left[r + \gamma v_{\pi}(s') \right], \quad \text{for all } s \in \mathcal{S}. \end{aligned}$$

Action-value function for policy π

$$q_{\pi}(s, a) \doteq \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a] = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

Bellman optimality equation

$$\begin{aligned} v_*(s) &= \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \max_a \sum_{s', r} p(s', r | s, a) \left[r + \gamma v_*(s') \right], \quad \text{or} \\ q_*(s, a) &= \mathbb{E} \left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right] \\ &= \sum_{s', r} p(s', r | s, a) \left[r + \gamma \max_{a'} q_*(s', a') \right], \end{aligned}$$

Some equations related to Bellman

$$\begin{aligned} v_{\pi}(s) &= \sum_a \pi(a | s) \sum_{s'} p(s' | s, a) [r(s, a) + v_{\pi}(s')] \\ q_{\pi}(s, a) &= \sum_{s'} p(s' | s, a) \left[r(s, a) + \gamma \sum_{a'} \pi(a' | s') q_{\pi}(s', a') \right] \end{aligned}$$

$$Q_t(a) \doteq \frac{\text{sum of rewards when } a \text{ taken prior to } t}{\text{number of times } a \text{ taken prior to } t} = \frac{\sum_{i=1}^{t-1} R_i \cdot \mathbb{1}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbb{1}_{A_i=a}}$$

Soft-max distribution

$$\Pr\{A_t = a\} \doteq \frac{e^{H_t(a)}}{\sum_{b=1}^k e^{H_t(b)}} \doteq \pi_t(a).$$

Action preferences

$$\begin{aligned} H_{t+1}(A_t) &\doteq H_t(A_t) + \alpha(R_t - \bar{R}_t)(1 - \pi_t(A_t)), \quad \text{and} \\ H_{t+1}(a) &\doteq H_t(a) - \alpha(R_t - \bar{R}_t)\pi_t(a), \quad \text{for all } a \neq A_t, \end{aligned}$$

Policy Evaluation: Sequence $\{v_k\}$ can be shown in general to converge to v_{π} as $k \rightarrow \infty$ under the same conditions that guarantee the existence of v_{π} . This algorithm is called **iterative policy evaluation**.

Policy Improvement

New greedy policy, π_0

$$\begin{aligned} \pi'(s) &\doteq \arg\max_a q_{\pi}(s, a) \\ &= \arg\max_a \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \arg\max_a \sum_{s', r} p(s', r | s, a) \left[r + \gamma v_{\pi}(s') \right], \end{aligned}$$

$$\begin{aligned} v_{\pi'}(s) &= \max_a \mathbb{E}[R_{t+1} + \gamma v_{\pi'}(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \max_a \sum_{s', r} p(s', r | s, a) \left[r + \gamma v_{\pi'}(s') \right]. \end{aligned}$$

$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} \dots \xrightarrow{I} \pi_* \xrightarrow{E} v_*.$$

1. Initialization

$V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation

Loop:

$\Delta \leftarrow 0$

Loop for each $s \in \mathcal{S}$:

$v \leftarrow V(s)$

$V(s) \leftarrow \sum_{s', r} p(s', r | s, \pi(s)) [r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

3. Policy Improvement

policy-stable \leftarrow true

For each $s \in \mathcal{S}$:

old-action $\leftarrow \pi(s)$

$\pi(s) \leftarrow \arg\max_{s'} \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$

If *old-action* $\neq \pi(s)$, then *policy-stable* \leftarrow false

If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

Value Iteration

$$\begin{aligned} v_{k+1}(s) &= \max_a \mathbb{E}[R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \max_a \sum_{s', r} p(s', r | s, a) \left[r + \gamma v_k(s') \right], \end{aligned}$$

Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

$\Delta \leftarrow 0$

Loop for each $s \in \mathcal{S}$:

$v \leftarrow V(s)$

$V(s) \leftarrow \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$

Output a deterministic policy, $\pi \approx \pi_*$, such that

$\pi(s) = \arg\max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$

Monte Carlo ES (Exploring Starts), for estimating $\pi \approx \pi_*$

Initialize:

$\pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all $s \in \mathcal{S}$

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

$\text{Returns}(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

Loop forever (for each episode):

Choose $S_0 \in \mathcal{S}$, $A_0 \in \mathcal{A}(S_0)$ randomly such that all pairs have probability > 0

Generate an episode from S_0, A_0 , following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $\text{Returns}(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$

$\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$

On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

$\text{Returns}(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $\text{Returns}(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$

$A^* \leftarrow \arg\max_a Q(S_t, a)$

For all $a \in \mathcal{A}(S_t)$:

$$\pi(a | S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

$$\text{Probability of nongreedy actions} \quad \frac{\varepsilon}{|\mathcal{A}(s)|}$$

If the step-size parameters, α , are not constant, then the estimate Q_n is a weighted average of previously received rewards with a weighting different from that given by Equation 2.6. What is the weighting on each prior reward for the general case, analogous to Equation 2.6, in terms of the sequence of step-size parameters?

$$Q_{n+1} = Q_n + \alpha_n (R_n - Q_n) \\ = \alpha_n R_n + (1 - \alpha_n) [Q_{n-1} + \alpha_{n-1} (R_{n-1} - Q_{n-1})] \\ = \dots \\ = \sum_{i=1}^n \alpha_i R_i \prod_{j=i+1}^n (1 - \alpha_j) + Q_1 \prod_{i=1}^n (1 - \alpha_i)$$

- If $Q_1 = 0$, is Q_n (for $n > 1$) biased? Justify your answer with brief words and equations!

(b) Equation 2.6

$$Q_{n+1} = (1-\alpha)Q_n + \alpha \frac{1}{(1-\alpha)^{n-1}} R_1 \\ \text{when } Q_1 = 0 \\ Q_{n+1} = \frac{\alpha}{(1-\alpha)} (1-\alpha)^{n-1} R_1 \\ \Rightarrow E(Q_{n+1}) = E\left(\sum_{i=1}^n \alpha_i (1-\alpha)^{n-i} R_i\right) \\ = \sum_{i=1}^n \alpha_i (1-\alpha)^{n-i} E(R_i) \\ \Rightarrow \text{When } \sum_{i=1}^n \alpha_i (1-\alpha)^{n-i} = 1 \\ E(Q_{n+1}) = E(R_1)$$

(c) Equation is unbiased

- Derive condition(s) for Q_1 for when Q_n will be unbiased.

(c) Q_n will be unbiased if $\sum_{i=1}^n \alpha_i (1-\alpha)^{n-i} = 1$

- Show that Q_n is an unbiased estimator as n (which is often referred to as asymptotically unbiased).

(d) $Q_{n+1} = (1-\alpha)Q_n + \sum_{i=1}^n \alpha_i (1-\alpha)^{n-i} R_i$

$$\lim_{n \rightarrow \infty} Q_{n+1} = \lim_{n \rightarrow \infty} (1-\alpha)Q_n + \lim_{n \rightarrow \infty} \sum_{i=1}^n \alpha_i (1-\alpha)^{n-i} R_i \\ = \lim_{n \rightarrow \infty} \alpha \cdot \frac{1-(1-\alpha)^{n+1}}{1-(1-\alpha)} \cdot R_1 \\ = \lim_{n \rightarrow \infty} \alpha \cdot \frac{1-(1-\alpha)^{n+1}}{\alpha} \cdot R_1 \\ = \lim_{n \rightarrow \infty} (1-\alpha)^{n+1} \cdot R_1 \\ = R_1 \\ \Rightarrow \lim_{n \rightarrow \infty} E(Q_{n+1}) = E(R_1)$$

- Why should we expect that the exponential recency-weighted average will be biased in practice?

For a , we are able to control the learning rate by modify α , the learning process could be speeded up by increasing α and vice versa. For Q_1 , the agent could start explore more in the beginning rather than always choosing the best choice by increase Q_1 .

- Show that in the case of two actions, the soft-max distribution is the same as that given by the logistic, or sigmoid, function often used in statistics and artificial neural networks

$$4. \text{ Softmax: } \sigma_a = \frac{\exp(H_a(a))}{\sum_{b=1}^B \exp(H_b(b))}$$

case of two action:

$$\sigma_a = \frac{\exp(H_a(a))}{\exp(H_a(a)) + \exp(H_a(b))} = \frac{1}{1 + \exp(H_a(b) - H_a(a))}$$

$$\text{if } Q = - (H_a(b) - H_a(a))$$

$$\text{then } \sigma_a = \frac{1}{1 + \exp(-Q)} = \sigma(Q) \rightarrow \text{sigmoid}$$

- UCB also produce spikes in the very beginning in both the two reproduced figures. Explain in your own words why the spikes appear (both the sharp increase and sharp decrease)

- In the optimistic initialization, the first few steps are not really random. Instead, we loop through all the actions multiple times, picking our large optimistic value at random. In the first round, all actions have equal optimization percentages (10%). The next action will have a spike, because on average the best action will have the largest value. This will be repeated with decreasing values until the effect of the optimal initial value fades away.

* Expand the Bellman equation for the 2 states in the recycling robot, for an arbitrary policy $\pi(a|s)$, discount factor γ , and domain parameters $\alpha, \beta, r_{\text{search}}, r_{\text{wait}}$ as described in the example.

$$U_{\pi}(\text{high}) = \pi(\text{search}|\text{high}) \cdot \gamma [r_{\text{search}} + \gamma U_{\pi}(\text{high})] \\ + \pi(\text{search}|\text{high}) \cdot \beta [\alpha \gamma U_{\pi}(\text{low}) + (1-\alpha) \gamma U_{\pi}(\text{high})] \\ + \pi(\text{wait}|\text{high}) \cdot \gamma [r_{\text{wait}} + \gamma U_{\pi}(\text{high})] \\ = \pi(\text{search}|\text{high}) \cdot \alpha \gamma [r_{\text{search}} + \gamma U_{\pi}(\text{high})] \\ + \pi(\text{search}|\text{high}) \cdot (1-\alpha) \gamma [r_{\text{search}} + \gamma U_{\pi}(\text{low})] \\ + \pi(\text{wait}|\text{high}) \cdot [r_{\text{wait}} + \gamma U_{\pi}(\text{high})] \\ U_{\pi}(\text{low}) = \pi(\text{wait}|\text{low}) \cdot \gamma [r_{\text{wait}} + \gamma U_{\pi}(\text{low})] \\ + \pi(\text{search}|\text{low}) \cdot \gamma [\alpha \gamma U_{\pi}(\text{low}) + (1-\alpha) \gamma U_{\pi}(\text{high})] \\ + \pi(\text{recharge}|\text{low}) \cdot \gamma [r_{\text{recharge}} + \gamma U_{\pi}(\text{high})] \\ + \pi(\text{search}|\text{low}) \cdot \beta [\alpha \gamma U_{\pi}(\text{low}) + (1-\alpha) \gamma U_{\pi}(\text{high})] \\ = \pi(\text{wait}|\text{low}) \cdot [r_{\text{wait}} + \gamma U_{\pi}(\text{low})] \\ + \pi(\text{search}|\text{low}) \cdot \beta [\alpha \gamma U_{\pi}(\text{low}) + (1-\alpha) \gamma U_{\pi}(\text{high})] \\ + \pi(\text{recharge}|\text{low}) \cdot \gamma U_{\pi}(\text{high}) \\ + \pi(\text{search}|\text{low}) \cdot (1-\beta) \cdot [-3 + \gamma U_{\pi}(\text{low})]$$

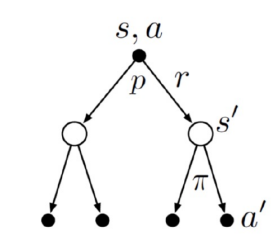
* You should now have two linear equations involving two unknowns, $v(\text{high})$ and $v(\text{low})$, as well as involving the policy $\pi(a|s)$, γ , and the domain parameters. Let $\alpha = 0.8, \beta = 0.6, \gamma = 0.9, r_{\text{search}} = 10, r_{\text{wait}} = 3$. Consider the policy $\pi(\text{search}|\text{high}) = 1, \pi(\text{wait}|\text{low}) = 0.5$, and $\pi(\text{recharge}|\text{low}) = 0.5$.

$$U_{\pi}(\text{high}) = \pi(\text{search}|\text{high}) \cdot \alpha \gamma [r_{\text{search}} + \gamma U_{\pi}(\text{high})] \\ + \pi(\text{search}|\text{high}) \cdot (1-\alpha) \gamma [r_{\text{search}} + \gamma U_{\pi}(\text{low})] \\ + \pi(\text{wait}|\text{high}) \cdot [r_{\text{wait}} + \gamma U_{\pi}(\text{high})] \\ = 1 \times 0.8 \times [10 + 0.9 \gamma U_{\pi}(\text{high})] \\ + 1 \times 0.2 \times [10 + 0.9 \gamma U_{\pi}(\text{low})] \\ U_{\pi}(\text{low}) = 8 + 0.72 U_{\pi}(\text{high}) + 3 + 0.18 U_{\pi}(\text{low}) \\ 0.18 U_{\pi}(\text{high}) = 10 + 0.18 U_{\pi}(\text{low}) \quad \text{①}$$

$$U_{\pi}(\text{low}) = \pi(\text{wait}|\text{low}) \cdot [r_{\text{wait}} + \gamma U_{\pi}(\text{low})] \\ + \pi(\text{search}|\text{low}) \cdot \beta [\alpha \gamma U_{\pi}(\text{low}) + (1-\alpha) \gamma U_{\pi}(\text{high})] \\ + \pi(\text{recharge}|\text{low}) \cdot \gamma U_{\pi}(\text{high}) \\ + \pi(\text{search}|\text{low}) \cdot (1-\beta) \cdot [-3 + \gamma U_{\pi}(\text{low})] \\ = 0.5 \times [3 + 0.9 \gamma U_{\pi}(\text{low})] \\ + 0.5 \times 0.9 \gamma U_{\pi}(\text{high}) \\ U_{\pi}(\text{low}) = 1.5 + 0.45 U_{\pi}(\text{low}) + 0.45 U_{\pi}(\text{high}) \\ 0.55 U_{\pi}(\text{low}) = 1.5 + 0.45 U_{\pi}(\text{high}) \\ 11 U_{\pi}(\text{low}) = 30 + 9 U_{\pi}(\text{high}) \quad \text{②}$$

Combine ① and ②

$$\Rightarrow 14 U_{\pi}(\text{high}) = 500 + 9 U_{\pi}(\text{low}) \Rightarrow 77 U_{\pi}(\text{low}) = 210 + 63 U_{\pi}(\text{high}) \\ \Rightarrow U_{\pi}(\text{low}) = 67 \quad 36.5 U_{\pi}(\text{low}) = 2460 \\ \Rightarrow U_{\pi}(\text{high}) = 79 \\ 14 U_{\pi}(\text{high}) = 500 + 603$$



q_{π} backup diagram

* Give an equation for v_{π} in terms of q_{π} and π .

$$U_{\pi} = E_{\pi}[G_t | S_t = s] \\ = \sum_{a'} E_{\pi}[G_t | S_t = s, A_t = a] P[A_t = a | S_t = s] \\ = \sum_{a'} q_{\pi}(s, a) \pi(a|s)$$

* Give an equation for q_{π} in terms of v_{π} and the four argument \mathbf{d} .

(b) $q(s, a) = E_{\pi}[G_t | S_t = s, A_t = a]$

$$= \sum_{s'} \gamma [G_t + \gamma V_{\pi}(s')] P[S_{t+1} = s' | S_t = s, A_t = a] \\ = \sum_{s'} \gamma [G_t + \gamma V_{\pi}(s')] P[S_{t+1} = s' | S_t = s, A_t = a] \\ = \sum_{s'} \gamma [r_t + \gamma V_{\pi}(s')] P[S_{t+1} = s' | S_t = s, A_t = a] \\ = \sum_{s'} \gamma [r_t + \gamma V_{\pi}(s')] P(\omega_t | s, a)$$

* What is the Bellman equation for action values, that is, for q_{π} ? It must give the action value $q_{\pi}(s, a)$ in terms of the action values, $q_{\pi}(s', a')$, of possible successors to the state-action pair (s, a) .

$$q_{\pi}(s, a) = \sum_{s'} P(s' | s, a) [r + \gamma U_{\pi}(s')] \\ = \sum_{s'} P(s' | s, a) [r + \sum_{a'} \pi(a' | s') q_{\pi}(s', a')]$$

s	a	s'	$p(s' s, a)$	$r(s, a, s')$
high	search	high	α	r_{search}
high	search	low	$1 - \alpha$	r_{search}
low	search	high	$1 - \beta$	-3
low	search	low	β	r_{search}
high	wait	high	1	r_{wait}
high	wait	low	0	-
low	wait	high	0	-
low	wait	low	1	r_{wait}
low	recharge	high	1	0
low	recharge	low	0	-

$$G_t = \sum_{k=0}^{\infty} \gamma^k = \frac{1}{1-\gamma}$$

$$\sum_{k=0}^{\infty} y^k = 1 + y + y^2 + y^3 + \dots =$$

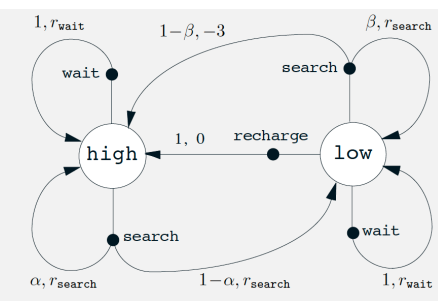
$$= 1 + y(1 + y + y^2 + y^3 + \dots) = 1 + y \left(\sum_{k=0}^{\infty} y^k \right)$$

$$\sum_{k=0}^{\infty} y^k = -\frac{1}{y-1} = \frac{1}{1-y}$$

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+1+k}$$

$$G'_t = \sum_{k=0}^{\infty} \gamma^k (R_{t+1+k} + c) = \sum_{k=0}^{\infty} \gamma^k R_{t+1+k} + \sum_{k=0}^{\infty} c \cdot \gamma^k =$$

$$= G_t + \frac{c}{1-\gamma}$$



s	a	s'	r	$p(s', r s, a)$
high	search	high	r_{search}	α
high	search	low	r_{search}	$1 - \alpha$
low	search	high	-3	$1 - \beta$
low	search	low	r_{search}	β
high	wait	high	r_{wait}	1
low	wait	low	r_{wait}	1
low	recharge	high	0	1