

Reinforcement Learning

Dynamic programming,
Bellman equations

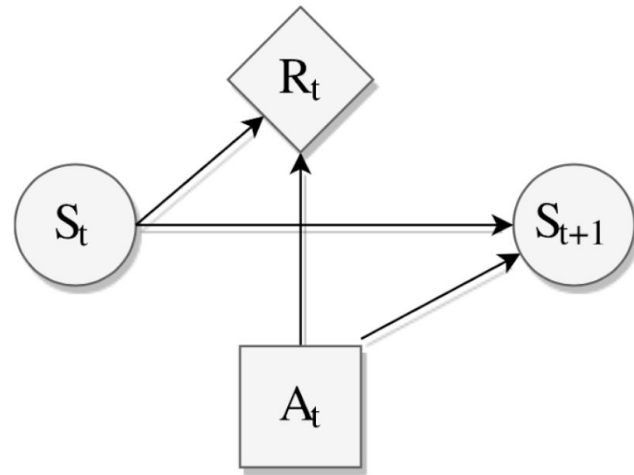
Александр Костин
telegramm: @Ko3tin
LinkedIn: [kostinalexander](#)

Recap. MDP

- s - состояние (наблюдение)
- a - действие
- r - награда за действие
- s' - следующее состояние

Свойство марковости:

$$p(r_{t+1}, s_{t+1} \mid s_0, a_0, r_0, \dots, s_t, a_t, r_t) = p(r, s_{t+1} \mid s_t, a_t)$$



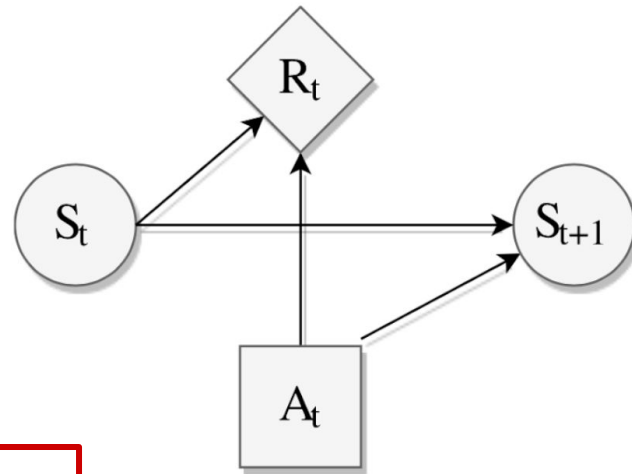
Recap. MDP

- s - состояние (наблюдение)
- a - действие
- r - награда за действие
- s' - следующее состояние

Свойство марковости:

$$p(r_{t+1}, s_{t+1} \mid s_0, a_0, r_0, \dots, s_t, a_t, r_t) = p(r, s_{t+1} \mid s_t, a_t)$$

динамика среды



Задача

Пусть известна динамика среды для какого-то MDP.

Как:

- Оценить агента?
- Улучшить агента?

Награда

$\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T, a_T, r_T)$ - траектория

$R = \sum_{t=0}^T r_t$ - награда за траекторию

Награда

$\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T, a_T, r_T)$ - траектория

$R = \sum_{t=0}^T r_t$ - награда за траекторию

$G_t = r_t + r_{t+1} + \dots + r_T$ - return на t шаге

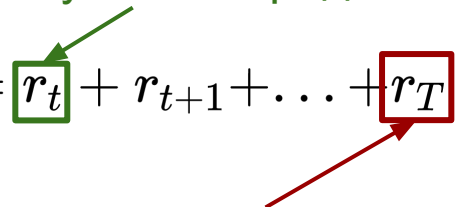
Награда

$\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T, a_T, r_T)$ - траектория

$R = \sum_{t=0}^T r_t$ - награда за траекторию

сиюминутная награда

$G_t = \boxed{r_t} + r_{t+1} + \dots + \boxed{r_T}$ - return на t шаге



отложенная награда

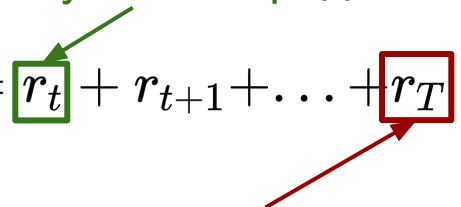
Награда

$\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T, a_T, r_T)$ - траектория

$R = \sum_{t=0}^T r_t$ - награда за траекторию

сиюминутная награда

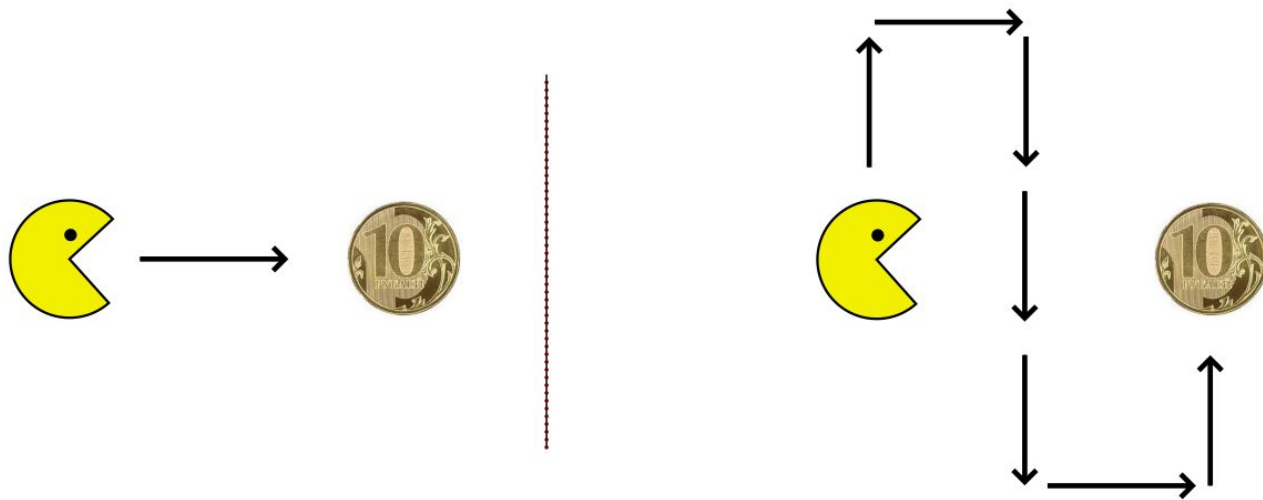
$G_t = \boxed{r_t} + r_{t+1} + \dots + \boxed{r_T}$ - return на t шаге



отложенная награда

Что важнее?

Дисконтирование награды



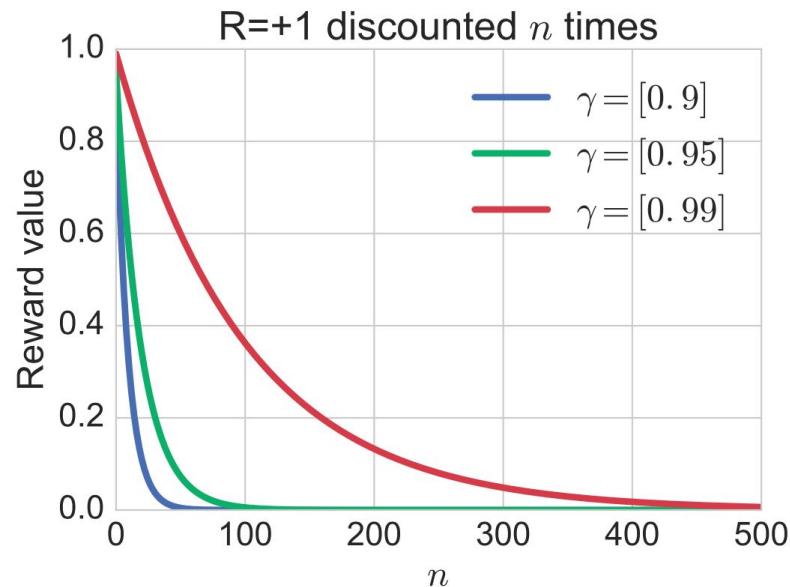
Будем уменьшать награду каждый шаг:

$$G_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^k r_{t+k} = \sum_{k=0}^{T-t} \gamma^k r_{t+k}$$

Что если эпизоды бесконечны?

$$G_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^k r_{t+k} = \sum_{k=0}^{T-t} \gamma^k r_{t+k}$$

- $R = +1$
- $T = +\infty$



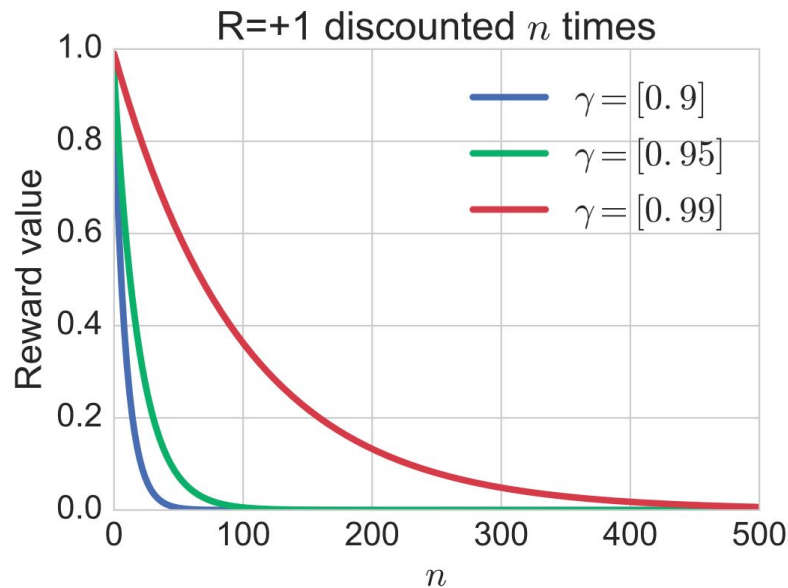
Что если эпизоды бесконечны?

$$G_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^k r_{t+k} = \sum_{k=0}^{T-t} \gamma^k r_{t+k}$$

- $R = +1$
- $T = +\infty$

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

γ	0.9	0.95	0.99
$\frac{1}{1-\gamma}$	10	20	100



Что если эпизоды бесконечны?

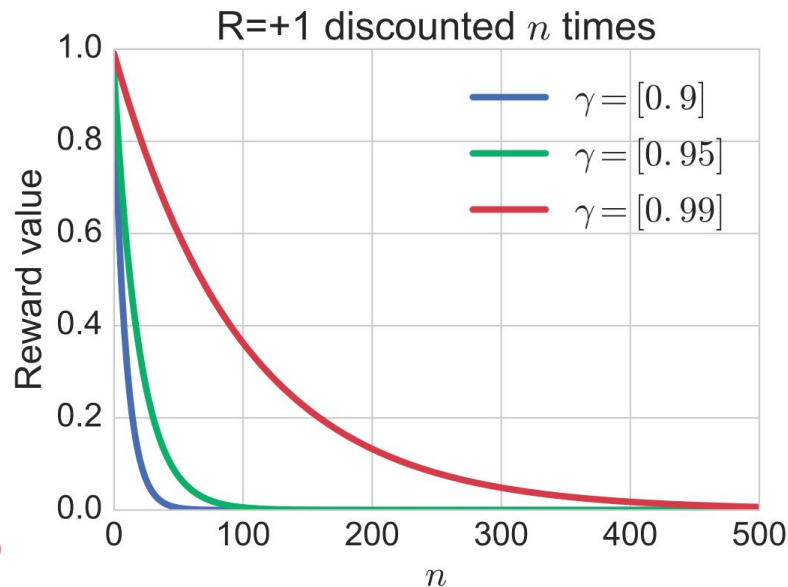
$$G_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^k r_{t+k} = \sum_{k=0}^{T-t} \gamma^k r_{t+k}$$

- $R = +1$
- $T = +\infty$

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

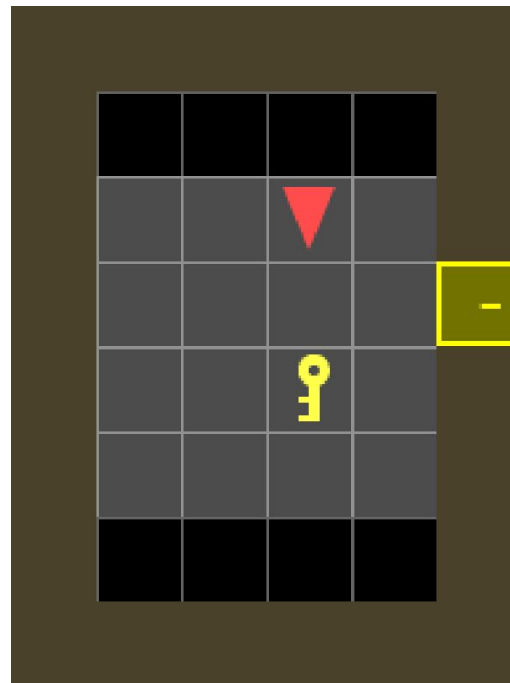
γ	0.9	0.95	0.99
$\frac{1}{1-\gamma}$	10	20	100

Дисконтирование
меняет оптимальную
политику!



Дисконтирование как модель эффекта

- Агент получает награду только при выходе из двери
- Ключ - крайне важная часть выполнения задания
- За поднятие ключа агент не получает награды



Дисконтирование как модель эффекта

Пусть γ - вероятность сохранения эффекта одного действия на другое

$$\begin{aligned} G_t = & (1 - \gamma)r_t \\ & + (1 - \gamma)\gamma(r_t + r_{t+1}) \\ & + (1 - \gamma)\gamma^2(r_t + r_{t+1} + r_{t+2}) \\ & \dots \\ & + \gamma^k \sum_{k=0}^{T-t} r_{t+k} \end{aligned}$$

Дисконтирование как модель эффекта

Пусть γ - вероятность сохранения эффекта одного действия на другое

The diagram illustrates the decomposition of the discount factor $(1 - \gamma)$ in the context of the Bellman optimality equation. It shows how the immediate reward is discounted by the probability of the effect terminating, and the discounted future rewards are discounted by the probability of the effect persisting.

Вероятность прекращения эффекта (Probability of effect termination) is represented by the red arrows pointing to the terms $(1 - \gamma)$ in the equation.

Вероятность сохранения эффекта (Probability of effect preservation) is represented by the green arrows pointing to the terms γ in the equation.

$$\begin{aligned} G_t = & (1 - \gamma)r_t \\ & + (1 - \gamma)\gamma(r_t + r_{t+1}) \\ & + (1 - \gamma)\gamma^2(r_t + r_{t+1} + r_{t+2}) \\ & \dots \\ & + \gamma^k \sum_{k=0}^{T-t} r_{t+k} \end{aligned}$$

Дисконтирование как модель эффекта

Пусть γ - вероятность сохранения эффекта одного действия на другое

Вероятность прекращения эффекта

$$\begin{aligned} G_t &= (1 - \gamma)r_t \\ &\quad + (1 - \gamma)\gamma(r_t + r_{t+1}) \\ &\quad + (1 - \gamma)\gamma^2(r_t + r_{t+1} + r_{t+2}) \\ &\quad \dots \\ &\quad + \gamma^k \sum_{k=0}^{T-t} r_{t+k} \\ &= r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^{T-t} r_T \end{aligned}$$

Вероятность сохранения эффекта

Оптимальность политики

Агент взаимодействует со средой в соответствии с политикой $\pi(\mathbf{a}|\mathbf{s})$

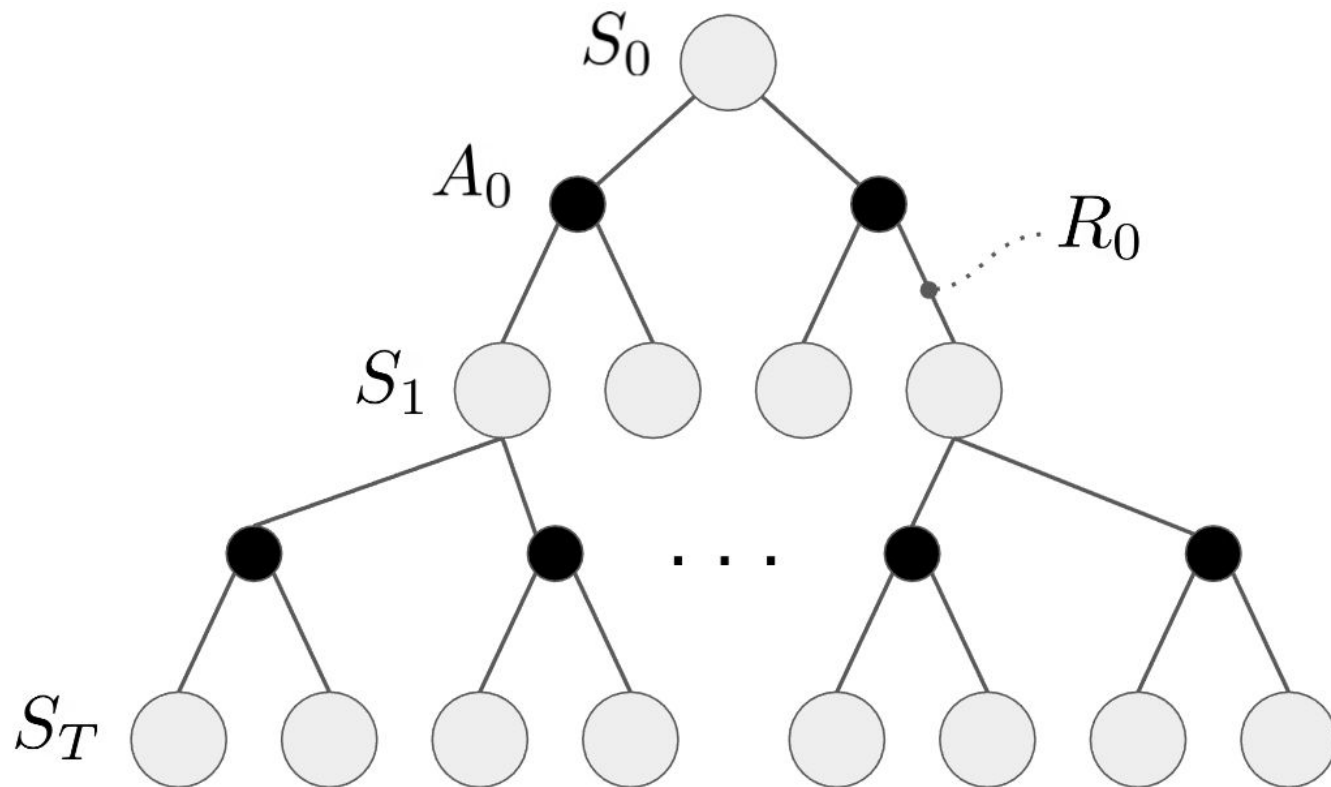
Как итог получилась траектория τ :

$$\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T, a_T, r_T), a_t \sim \pi(\cdot | s_t), (s_t, r_t) \sim p(\cdot | s_{t-1}, a_{t-1})$$

Оптимальная политика максимизирует ожидание \mathbf{G} за траекторию:

$$\begin{aligned} J(\pi) &= E_{\pi}[G_0] = E_{p(\tau | \pi)}[G_0] = \\ &= E_{s_0 \sim p(s_0)} \left[E_{a \sim \pi(\cdot | s_0)} \left[r_0 + E_{s_1 \sim p(\cdot | s_0, a_0)} \left[E_{a_1 \sim \pi(\cdot | s_1)} [\gamma r_1 + \dots] \right] \right] \right] \end{aligned}$$

Как найти оптимальную политику?



Ценность состояния

Насколько хорошо агенту пребывать в некотором состоянии \mathbf{s} и следовать политике π ?

$$V_{\pi}(s) = E[G_t \mid s_t = s] = E_{\pi} \left[\sum_{k=0}^{T-t} \gamma^k r_{t+k} \mid s_t = s \right]$$

Если состояние \mathbf{s} конечное, то $\mathbf{V}_{\pi}(\mathbf{s})=0$

Ценность состояния

$$V_{\pi}(s) = E_{\pi}[G_t \mid s_t = s] = E_{\pi}[r_t + \gamma G_{t+1} \mid s_t = s] =$$

Ценность состояния

$$\begin{aligned} V_{\pi}(s) &= E_{\pi}[G_t \mid s_t = s] = E_{\pi}[r_t + \gamma G_{t+1} \mid s_t = s] = \\ &= \sum_a \pi(a \mid s) \sum_{r, s'} p(r_t = r, s_{t+1} = s' \mid s_t = s, a_t = a, s_{t-1} = s'', a_{t-1} = a'', r_{t-1} = r'', \dots) [r + \gamma E_{\pi}[G_{t+1} \mid s_{t+1} = s']] = \end{aligned}$$

Ценность состояния

$$\begin{aligned} V_{\pi}(s) &= E_{\pi}[G_t \mid s_t = s] = E_{\pi}[r_t + \gamma G_{t+1} \mid s_t = s] = \\ &= \sum_a \pi(a \mid s) \sum_{r, s'} p(r_t = r, s_{t+1} = s' \mid s_t = s, a_t = a, s_{t-1} = s'', a_{t-1} = a'', r_{t-1} = r'', \dots) [r + \gamma E_{\pi}[G_{t+1} \mid s_{t+1} = s']] = \\ &= \sum_a \pi(a \mid s) \sum_{r, s'} p(r_t = r, s_{t+1} = s' \mid s_t = s, a_t = a) [\dots] = \end{aligned}$$

Ценность состояния

$$\begin{aligned} V_{\pi}(s) &= E_{\pi}[G_t \mid s_t = s] = E_{\pi}[r_t + \gamma G_{t+1} \mid s_t = s] = \\ &= \sum_a \pi(a \mid s) \sum_{r, s'} p(r_t = r, s_{t+1} = s' \mid s_t = s, a_t = a, s_{t-1} = s'', a_{t-1} = a'', r_{t-1} = r'', \dots) [r + \gamma E_{\pi}[G_{t+1} \mid s_{t+1} = s']] = \\ &= \sum_a \pi(a \mid s) \sum_{r, s'} p(r_t = r, s_{t+1} = s' \mid s_t = s, a_t = a) [\dots] = \\ &= \sum_a \pi(a \mid s) \sum_{r, s'} p(r, s' \mid s, a) [r + \gamma E_{\pi}[G_{t+1} \mid s_{t+1} = s']] = \\ &= \sum_a \pi(a \mid s) \sum_{r, s'} p(r, s' \mid s, a) [r + \gamma V_{\pi}(s')] \end{aligned}$$

Ценность состояния

$$\begin{aligned}
 V_{\pi}(s) &= E_{\pi}[G_t \mid s_t = s] = E_{\pi}[r_t + \gamma G_{t+1} \mid s_t = s] = \\
 &= \sum_a \pi(a \mid s) \sum_{r,s'} \boxed{p(r_t = r, s_{t+1} = s' \mid s_t = s, a_t = a, s_{t-1} = s'', a_{t-1} = a'', r_{t-1} = r'', \dots)} [r + \gamma E_{\pi}[G_{t+1} \mid s_{t+1} = s']] = \\
 &= \sum_a \pi(a \mid s) \sum_{r,s'} \boxed{p(r_t = r, s_{t+1} = s' \mid s_t = s, a_t = a)} [\dots] = \\
 &= \sum_a \pi(a \mid s) \sum_{r,s'} p(r, s' \mid s, a) [r + \gamma \boxed{E_{\pi}[G_{t+1} \mid s_{t+1} = s']}] = \\
 V_{\pi}(s) &= \sum_a \boxed{\pi(a \mid s)} \sum_{r,s'} \boxed{p(r, s' \mid s, a)} [r + \gamma \boxed{V_{\pi}(s')}]
 \end{aligned}$$

→ Стохастичность политики
→ Стохастичность среды

Ценность действия

Насколько хорошо агенту совершить действие **a** в некотором состоянии **s** и следовать политике π ?

$$Q_{\pi}(s, a) = E[G_t \mid s_t = s, a_t = a] = E_{\pi} \left[\sum_{k=0}^{T-t} \gamma^k r_{t+k} \mid s_t = s, a_t = a \right]$$

Если состояние **s** конечное, то $Q_{\pi}(s, a) = 0$, для любого **a**

Ценность действия

$$\begin{aligned} Q_{\pi}(s, a) &= E_{\pi}[G_t \mid s_t = s, a_t = a] = E_{\pi}[r_t + \gamma G_{t+1} \mid s_t = s, a_t = a] = \\ &= \sum_{r, s'} p(r, s' \mid s, a) [r + \gamma E_{\pi}[G_{t+1} \mid s_{t+1} = s']] = \\ Q_{\pi}(s, a) &= \sum_{r, s'} p(r, s' \mid s, a) [r + \gamma V_{\pi}(s')] \end{aligned}$$


Связь Q и V функций

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

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

Связь Q и V функций

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

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


Связь Q и V функций

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

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

$$V_{\pi}(s) = \sum_a \pi(a | s) Q_{\pi}(s, a) = E_{a \sim \pi(\cdot | s)} [Q_{\pi}(s, a)]$$

$$Q_{\pi}(s, a) = \sum_{r, s'} p(r, s' | s, a) \left[r + \gamma \sum_a \pi(a | s') Q_{\pi}(s', a) \right]$$

Уравнения Беллмана для матожидания

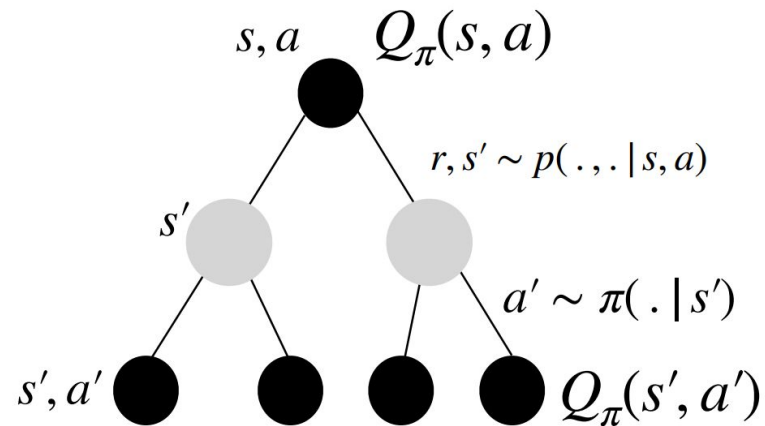
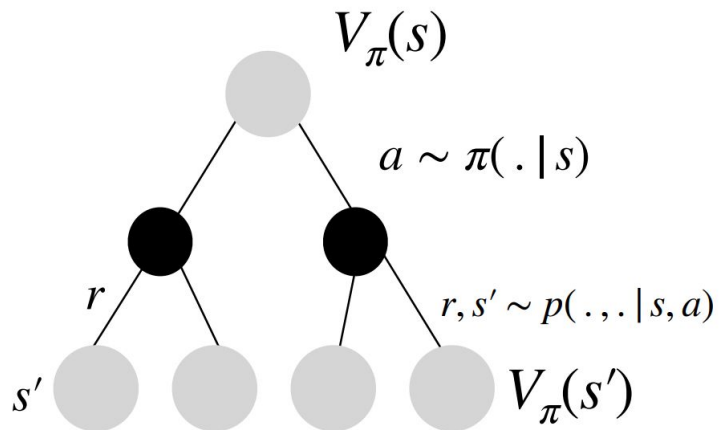
Для $V(s)$:

$$\begin{aligned} V_{\pi}(s) &= \sum_a \pi(a | s) \sum_{r, s'} p(r, s' | s, a) [r + \gamma V_{\pi}(s')] \\ &= E_{\pi}[r_t + \gamma V_{\pi}(s_{t+1}) | s_t = s] \end{aligned}$$

Для $Q(s,a)$:

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

Уравнения Беллмана для матожидания



Уравнения Беллмана для оптимальности

Вспомним, что наша цель - найти политику максимизирующую суммарную награду.

Определим отношение на множестве политик:

$$\pi \geq \pi' \iff V_{\pi}(s) \geq V_{\pi'}(s) \forall s \in \mathcal{S}$$

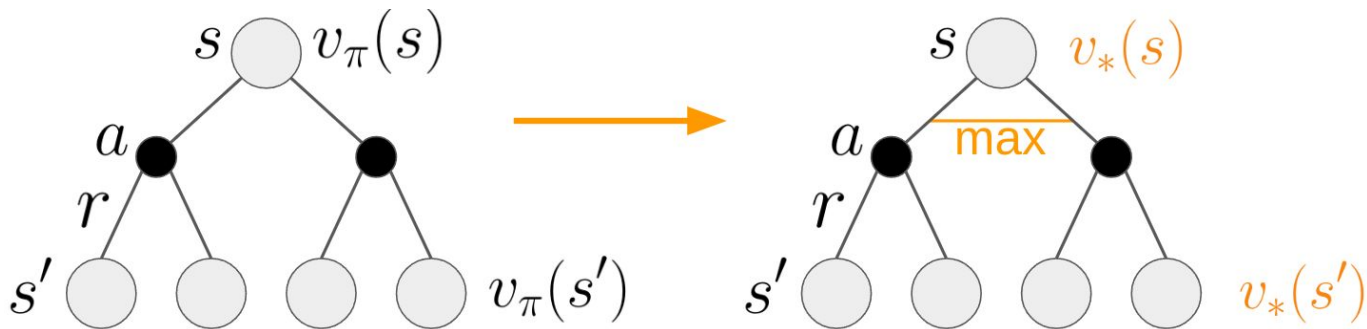
Лучшая политика π^* не хуже любой другой и имеет Q и V функции:

$$V^*(s) = \max_{\pi} V_{\pi}(s)$$

$$\pi^*(a | s) = [a = \operatorname{argmax}_{a'} Q^*(s, a')]$$

$$Q^*(s, a) = \max_{\pi} Q_{\pi}(s, a)$$

Уравнения Беллмана для оптимальности

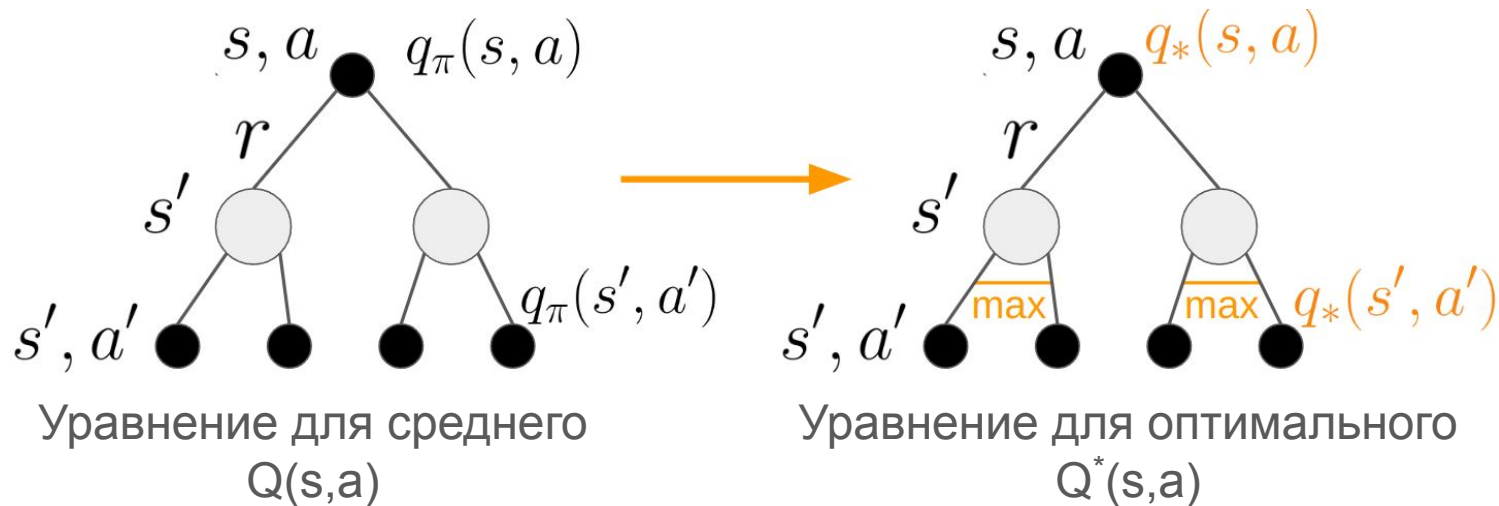


Уравнение для среднего
 $V(s)$

Уравнение для оптимального
 $V^*(s)$

$$\begin{aligned} V^* &= \max_a \sum_{r, s'} p(r, s' | s, a) [r + \gamma V^*(s')] \\ &= \max_a E[r_t + \gamma V^*(s_{t+1}) | s_t = s, a_t = a] \end{aligned}$$

Уравнения Беллмана для оптимальности



$$\begin{aligned} Q^{*}(s, a) &= E \left[r + \gamma \max_{a'} Q^{*}(s_{t+1}, a') \mid s_t = t, a_t = a \right] \\ &= \sum_{s', r} p(s', r \mid s, a) \left[r + \gamma \max_{a'} Q^{*}(s', a') \right] \end{aligned}$$

Сложность вычисления Q и V функций

- Можно заметить, что для всех 4-х выражений есть N неизвестных и N уравнений -> можно решить СЛАУ

$$V^{\pi}(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} r(s, a) + \gamma \mathbb{E}_{s' \sim p(s'|s)} V^{\pi}(s')$$

Сложность вычисления Q и V функций

- Можно заметить, что для всех 4-х выражений есть N неизвестных и N уравнений -> можно решить СЛАУ

$$V^{\pi}(s) = u(s) + \gamma \mathbb{E}_{s' \sim p(s'|s)} V^{\pi}(s')$$

Сложность вычисления Q и V функций

- Можно заметить, что для всех 4-х выражений есть N неизвестных и N уравнений -> можно решить СЛАУ

$$V^\pi(s) = u(s) + \gamma \mathbb{E}_{s' \sim p(s'|s)} V^\pi(s')$$

Относительно V все линейно

$$V = U + \gamma PV$$

Сложность вычисления Q и V функций

- Можно заметить, что для всех 4-х выражений есть N неизвестных и N уравнений -> можно решить СЛАУ

$$V^\pi(s) = u(s) + \gamma \mathbb{E}_{s' \sim p(s'|s)} V^\pi(s')$$

Относительно V все линейно

$$V = U + \gamma P V$$

$$(I - \gamma P)V = U$$

Сложность вычисления Q и V функций

- Можно заметить, что для всех 4-х выражений есть N неизвестных и N уравнений -> можно решить СЛАУ

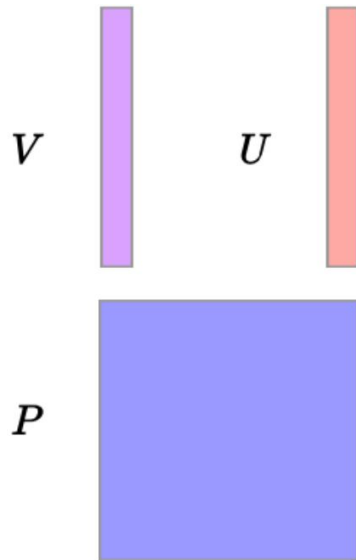
$$V^\pi(s) = u(s) + \gamma \mathbb{E}_{s' \sim p(s'|s)} V^\pi(s')$$

Относительно V все линейно

$$V = U + \gamma P V$$

$$(I - \gamma P)V = U$$

$$V = (I - \gamma P)^{-1} U$$



Сложность вычисления Q и V функций

- Можно заметить, что для всех 4-х выражений есть N неизвестных и N уравнений -> можно решить СЛАУ

$$V^\pi(s) = u(s) + \gamma \mathbb{E}_{s' \sim p(s'|s)} V^\pi(s')$$

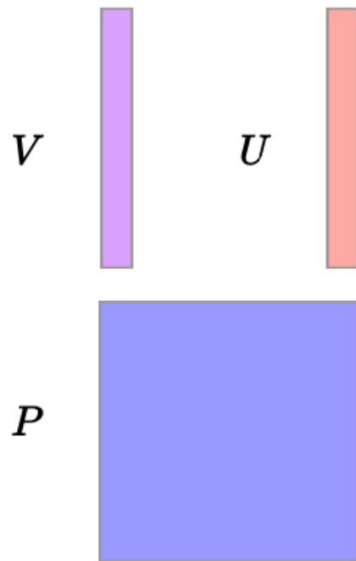
Относительно V все линейно

$$V = U + \gamma P V$$

$$(I - \gamma P)V = U$$

$$V = (I - \gamma P)^{-1} U$$

Сложность порядка $|s|^3$!



Операторы Беллмана

Оператор Беллмана для ожидания V:

$$[T^\pi V](s) = E_{r,s' \mid s,a=\pi(s)} [r + \gamma V(s')]$$

Оператор Беллмана для ожидания Q:

$$[T^\pi Q](s, a) = E_{r,s' \mid s,a} [r + \gamma E_{a' \sim \pi(s')} [Q(s', a')]]$$

Оператор Беллмана для оптимальности V:

$$[TV](s) = \max_a E_{r,s' \mid s,a} [r + \gamma V(s')]$$

Оператор Беллмана для оптимальности Q:

$$[T^\pi Q](s, a) = E_{r,s' \mid s,a} \left[r + \gamma \max_{a'} [Q(s', a')] \right]$$

Сжимающее свойство

Сжатие: $d(T(v), T(u)) \leq \gamma d(v, u)$, $\gamma \in [0, 1)$

Оператор: $[TV](s) = \max_a E_{r,s' | s,a} [r + \gamma V(s')]$

Обозначим $a^* = \arg \max_a E_{r,s' | s,a} [r + \gamma V(s')]$, тогда для любого \mathbf{s} :

$$\begin{aligned} [Tv](s) - [Tu](s) &\leq r(s, a^*) + \gamma E_{s' | s,a} [v(s')] - r(s, a^*) - \gamma E_{s' | s,a} [u(s')] \\ &= \gamma E_{s' | s,a} [v(s') - u(s')] \\ &\leq \gamma E_{s' | s,a} [|v(s') - u(s')|] \\ &\leq \gamma \max_{s'} |v(s') - u(s')| \\ &= \gamma \|v - u\|_\infty \end{aligned}$$

А взяв максимум по левой части получим:

$$\|Tv - Tu\|_\infty \leq \gamma \|v - u\|_\infty$$

Generalized Policy Iteration

- 1) Policy Evaluation
- 2) Policy Improvement

Policy Evaluation

$$V_{\pi}(s) = \sum_a \pi(a | s) \sum_{r, s'} p(r, s' | s, a) [r + \gamma V_{\pi}(s')] \quad \text{Решение СЛАУ}$$



$$V_{k+1}^{\pi}(s) = [T^{\pi} V_k](s) = \sum_a \pi(a | s) \sum_{r, s'} p(r, s' | s, a) [r + \gamma V_k^{\pi}(s')]$$

Итерационный метод

Policy Evaluation

Input π , the policy to be evaluated

Initialize an array $V(s) = 0$, for all $s \in \mathcal{S}^+$

Repeat

$\Delta \leftarrow 0$

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

$v \leftarrow V(s)$

Bellman expectation operator for $V^\pi(s)$

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

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

until $\Delta < \theta$ (a small positive number)

Output $V \approx v_\pi$

Policy Improvement

Давайте совершать действия жадно относительно $Q^\pi(s,a)$:

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

Policy Improvement

Давайте совершать действия жадно относительно $Q^\pi(s,a)$:

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

Этот процесс гарантированно произведет лучшую политику:

Если $Q^\pi(s, \pi'(s)) \geq V^\pi(s)$ для любого s

тогда $V^{\pi'}(s) \geq V^\pi(s)$

следовательно $\pi' \geq \pi$

Policy Improvement

Давайте совершать действия жадно относительно $Q^\pi(s,a)$:

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

Если $\pi' = \pi \rightarrow V^{\pi'} = V^\pi$ и V^π удовлетворяет уравнению оптимальности Беллмана:

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

Generalized Policy Iteration

Как скомбинировать Policy Iteration и Policy Improvement?

Generalized Policy Iteration

- 1) Evaluate given policy
- 2) Improve policy greedy w.r.t. to its value function

Policy Iteration:

- 1) Evaluate policy until convergence
- 2) Improve policy

Value Iteration:

- 1) Evaluate policy only with single iteration
- 2) Improve policy

Policy Iteration

1. Initialization

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

2. Policy Evaluation

Repeat

$\Delta \leftarrow 0$

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

Bellman expectation operator for $V^\pi(s)$

$v \leftarrow V(s)$

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

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

until $\Delta < \theta$ (a small positive number)

3. Policy Improvement

policy-stable \leftarrow true

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

Almost Bellman optimality operator for $V^\pi(s)$

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

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

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

v_k for the
random policy

greedy policy
w.r.t. v_k

$k = 0$

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

	↔	↔	↔
↔	↔	↔	↔
↔	↔	↔	↔
↔	↔	↔	

random
policy

$k = 1$

0.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0

	←	↔	↔
↑	↔	↔	↔
↔	↔	↔	↓
↔	↔	→	

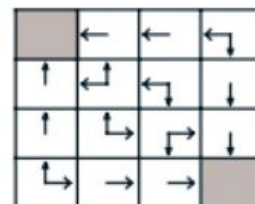
$k = 2$

0.0	-1.7	-2.0	-2.0
-1.7	-2.0	-2.0	-2.0
-2.0	-2.0	-2.0	-1.7
-2.0	-2.0	-1.7	0.0

	←	←	↔
↑	↖	↔	↓
↑	↔	↘	↓
↔	→	→	

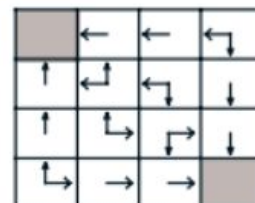
$k = 3$

0.0	-2.4	-2.9	-3.0
-2.4	-2.9	-3.0	-2.9
-2.9	-3.0	-2.9	-2.4
-3.0	-2.9	-2.4	0.0



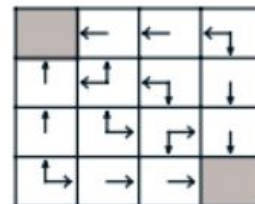
$k = 10$

0.0	-6.1	-8.4	-9.0
-6.1	-7.7	-8.4	-8.4
-8.4	-8.4	-7.7	-6.1
-9.0	-8.4	-6.1	0.0



$k = \infty$

0.0	-14.	-20.	-22.
-14.	-18.	-20.	-20.
-20.	-20.	-18.	-14.
-22.	-20.	-14.	0.0



optimal
policy

Value Iteration

Initialize array V arbitrarily (e.g., $V(s) = 0$ for all $s \in \mathcal{S}^+$)

Repeat

$\Delta \leftarrow 0$

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 $\Delta < \theta$ (a small positive number)

Bellman **optimality**
equation for $v(s)$

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

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