

Оптимизация в задачах обучения с подкреплением

Методы оптимизации в машинном обучении

Никита Юдин, iudin.ne@phystech.edu

Московский государственный университет имени М.В. Ломоносова
Факультет вычислительной математики и кибернетики
Кафедра математических методов прогнозирования

3 мая 2024

Марковский процесс принятия решений (MDP)

В обучении с подкреплением взаимодействия между агентом и окружающей средой часто описываются марковским процессом принятия решений (MDP). Различают:

- дисконтированный марковский процесс принятия решений (*γ -Discounted Markov Decision Process, DMDP*);
- MDP с усредненным вознаграждением (*infinite-horizon Average reward Markov Decision Process, AMDP*);
- эпизодический марковский процесс принятия решений (*H-episodic Markov Decision Process, HMDP*);
- другие, в том числе и *частично наблюдаемые*.

Марковский процесс принятия решений (MDP)

Марковский процесс принятия решений представляет собой систему, которая со временем ($t = 0, 1, 2, \dots$) претерпевает случайные изменения и обозначается кортежем $M = (\mathcal{S}, \mathcal{A}, p, r, \gamma)$ со следующими объектами:

- (i) \mathcal{S} – пространство состояний, $S := |\mathcal{S}|$ – количество уникальных состояний.
- (ii) \mathcal{A} – пространство действий, $A := |\mathcal{A}|$ – количество уникальных действий.
- (iii) $p(s, a; s')$ – вероятность перехода из состояния $s \in \mathcal{S}$ в момент времени t с определенным действием $a \in \mathcal{A}$ в состояние $s' \in \mathcal{S}$ в момент $(t + 1)$ (при этом $\sum_{s' \in \mathcal{S}} p(s, a; s') = 1$,
 $p(s, a; s') \equiv P(s'|s, a)$). Функция вероятности $p(\cdot)$ вместе с функцией вероятности $P(a|s)$ задают вероятности перехода для Марковского ядра, оно же ядро MDP.

Марковский процесс принятия решений (MDP)

(iv) Функция награды

$r_{\xi}(s, a) : \Omega \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, ($\mathbb{E}_{\xi} [r_{\xi}(s, a)] = r(s, a)$, где $\mathbb{E}[\cdot]$ – математическое ожидание). В зависимости от постановки задачи функция награды может зависеть от следующего за состоянием s состояния s' :

$$r_{\xi}(s, a; s') : \Omega \times \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1], \quad \mathbb{E}_{\xi} [r_{\xi}(s, a; s')] = r(s, a; s').$$

Стоит отметить, что мы предполагаем в общем случае стохастическую природу функции награды в зависимости от случайной величины $\xi \in \Omega$. При детерминированном вычислении награды относительно фиксированных (s, a) или (s, a, s') мы просто опускаем обозначение ξ в $r_{\xi}(\cdot)$ и математическое ожидание по нему.

Марковский процесс принятия решений (MDP)

- (iv) (продолжение) Здесь и далее используется предположение об ограниченности награды за каждое действие, поэтому без ограничений общности использованы приведённые выше определения функции $r(\cdot)$. В работе используются детерминированные относительно своих аргументов награды, если не оговорено иное.
- (v) $\gamma \in (0, 1]$ – коэффициент дисконтирования для DMDP, для AMDP $\gamma = 1$, но просто положив $\gamma = 1$ из DMDP не сделать AMDP, понадобится ещё усреднение суммарной награды агента за взаимодействие с MDP по времени.

Нередко рассматривается более общая форма $M = (\mathcal{S}, \mathcal{A}, p, r, \mu_0, \gamma)$, в которой μ_0 – вероятностное распределение начального состояния $s_0 \sim \mu_0$, при явном отсутствии μ_0 происходит обуславливание всех вычислений на $s_0 \in \mathcal{S}$.

MDP. Принятие решений

- Здесь и далее приводятся результаты для дискретных \mathcal{S} и \mathcal{A} с конечными мощностями, однако они могут быть обобщены на непрерывный случай заменой суммы по переменной в области её непрерывности на соответствующий интеграл по области.
- Стратегией принятия решений или политикой агента, принимающего решение в MDP, обозначим через символ π и присвоим ему отображения, задающие вероятностную меру на пространстве действий:

$\pi(a|s) \equiv P(a|s)$ в общем случае, $\hat{a} \sim \pi(\cdot|s)$ или $\pi(s) \sim \pi(\cdot|s)$;
 $\hat{a} := \pi(s)$ в случае вырожденного распределения: $P(\hat{a}|s) = 1$.

MDP. Ядро

- Введённое распределение позволяет явно записать Марковское ядро перехода между состояниями $s \mapsto s'$, оно же ядро MDP, его также корректно называть Марковским ядром, обусловленным политикой π :

$$P^\pi(s'|s) := \sum_{a \in \mathcal{A}} \pi(a|s) p(s, a; s').$$

- В процессах с конечным количеством состояний Марковское ядро можно задать с помощью матрицы:

$$P^\pi = \left(\sum_{a \in \mathcal{A}} \pi(a|s) p(s, a; s') \right)_{s \in \mathcal{S}, s' \in \mathcal{S}}, \quad s - \text{строка}, s' - \text{столбец}.$$

MDP. Ядро

В процессе взаимодействия с марковским процессом стратегия π собирает траекторию $\tau_t := (s_0, a_0, r_0, \dots, s_t, a_t, r_t)$. Её правдоподобие выражается следующим образом:

$$P(\tau_{H-1} | \pi) = \mu_0(s_0) \prod_{t=0}^{H-2} (\pi(a_t | s_t) p(s_{t+1}, a_t | s_t)) \pi(a_{H-1} | s_{H-1}).$$

DMDP. V -функция ценности

Для фиксированной политики и начального состояния $s_0 = s$ определяется V -функция значений (ценности) $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$ как дисконтированная сумма будущих вознаграждений:

$$V^\pi(s) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a(s_t)) \middle| \pi, s_0 = s \right],$$

где s_t – состояние системы в момент времени t , $a(s_t)$ – выбор действия в соответствии с политикой $\pi(\cdot)$. Это есть средняя награда по политике π , если агент начинает действовать в момент времени t из состояния s . Иногда индекс политики опускают: $V(s) := V^\pi(s)$.

Замечание

В определении V -функции в левой части выражения отсутствует обозначение t в силу однородности MDP.

DMDP. Q-функция ценности

Схожим образом задается Q-функция ценности $Q^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$:

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a(s_t)) \middle| \pi, s_0 = s, a_0 = a \right].$$

То же, что V функция, только теперь из состояния $s_t = s$ обязательно сначала совершается действие $a_t = a$. Иногда индекс политики опускают: $Q(s, a) := Q^\pi(s, a)$.

Замечание

В определении V - и Q - функции в левой части выражения отсутствует обозначение t в силу однородности MDP.

DMDP. Дисконтированная награда

Дисконтированная кумулятивная награда за эпизод длины $H - t$, $t = \overline{0, H - 1}$:

$$R_t^{H-1} := \sum_{j=t}^{H-1} \gamma^{j-t} r(s_j, a(s_j)) \text{ и } R_t := R_t^\infty := \sum_{j=t}^{\infty} \gamma^{j-t} r(s_j, a(s_j)).$$

При переходе к AMDP ($\gamma = 1$) наиболее естественным аналогом кумулятивной награды является среднее арифметическое наград по времени:

$$R_t^{H-1} := \frac{1}{H-t} \sum_{j=t}^{H-1} r(s_j, a(s_j)) \text{ и } R_t := R_t^\infty := \lim_{H \rightarrow \infty} \frac{1}{H-t} \sum_{j=t}^{H-1} r(s_j, a(s_j)).$$

DMDP. Дисконтированная награда

Мажоранта на $V^\pi(\cdot)$:

$$r(s, a) \in [0, 1] : \quad 0 \leq V^\pi(s) \leq \frac{1}{(1 - \gamma)} \quad \forall s \in \mathcal{S}, a \in \mathcal{A}.$$

Q^π -функция обладает той же мажорантой, что и $V^\pi(\cdot)$:

$$r(s, a) \in [0, 1] : \quad 0 \leq Q^\pi(s, a) \leq \frac{1}{(1 - \gamma)} \quad \forall s \in \mathcal{S}, a \in \mathcal{A}.$$

DMDP. Уравнения Беллмана

$$\begin{aligned} V^\pi(s) &= \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a(s_t)) \middle| \pi, s_0 = s \right] = \\ &= \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s' \in \mathcal{S}} p(s, a; s') [r(s, a) + \gamma V^\pi(s')] = \mathbb{E}_\pi [Q^\pi(s, a)] = \\ &= \mathbb{E}_\pi [r(s, a)] + \gamma \mathbb{E}_{p, \pi} [V^\pi(s')] = \mathbb{E}_\pi [r(s, a)] + \gamma \mathbb{E}_{p, \pi} [Q^\pi(s', a')] ; \end{aligned}$$

$$\begin{aligned} Q^\pi(s, a) &= \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a(s_t)) \middle| \pi, s_0 = s, a_0 = a \right] = \\ &= \sum_{s' \in \mathcal{S}} p(s, a; s') [r(s, a) + \gamma V^\pi(s')] = \\ &= r(s, a) + \gamma \mathbb{E}_p [V^\pi(s')] = r(s, a) + \gamma \mathbb{E}_{p, \pi} [Q^\pi(s', a')] , \\ &a' \sim \pi(\cdot | s'). \end{aligned}$$

DMDP. Задача RL

Цель задачи обучения с подкреплением (Reinforcement Learning, RL) – поиск политики, позволяющей получить максимальное кумулятивное вознаграждение в долгосрочной перспективе. В большинстве практических случаев задача RL формулируется как задача оптимизации следующего формата:

$$\pi^* \in \operatorname{Arg\,max}_{\pi \in \hat{\Pi}} \left\{ \mathbb{E}_{P(\tau_{H-1}|\pi)} \left[R_0^{H-1} \right] = \mathbb{E}_{s \sim \mu_0} [V^\pi(s)] \right\};$$
$$\hat{\Pi} := \left\{ \pi \left| \pi(a|s) \geq 0, \sum_{\hat{a} \in \mathcal{A}} \pi(\hat{a}|s) = 1, a \in \mathcal{A}, s \in \mathcal{S} \right. \right\}.$$

DMDP. Задача RL

Следующее утверждение [1] (детали: глава 3, утв. 21 и предшествующие) задаёт подкласс оптимальных политик, в рамках которого достаточно производить поиск интересующей π .

Пусть Π – набор всех нестационарных и рандомизированных политик. $V^\pi(s)$, $Q^\pi(s, a)$ зажаты между 0 и $\frac{1}{1-\gamma}$, следовательно, существуют конечные

$$V^*(s) := \sup_{\pi \in \Pi} \{V^\pi(s)\}, \quad Q^*(s, a) := \sup_{\pi \in \Pi} \{Q^\pi(s, a)\};$$

$\exists \pi$ – стационарная, детерминированная, такая, что $\forall s \in \mathcal{S}, a \in \mathcal{A}$:

$$V^\pi(s) = V^*(s), \quad Q^\pi(s, a) = Q^*(s, a),$$

а, значит, π – оптимальная политика.

DMDP. Задача RL

В данном утверждении мы можем легко заменить операцию \sup на операцию \max , как минимум, в случае наград с достижимыми верхними гранями:

$$V^*(s) = \max_{\pi \in \Pi} \{V^\pi(s)\},$$

где V^* – оптимальная функция ценности. Введём обозначение класса всех отображений, описывающих детерминированные политики в данном процессе:

$$\mathbb{A} = \{a(\cdot) \mid a : \mathcal{S} \mapsto \mathcal{A}\}.$$

DMDP. Уравнения Беллмана

Если воспользоваться принципом динамического программирования, то удаётся вывести уравнение оптимальности Беллмана на V -функцию ценности:

$$\begin{aligned} V^*(s) &= \max_{a(\cdot) \in \mathcal{A}} \left\{ \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a(s_t)) \right] \right\} = \\ &= \max_{a(\cdot) \in \mathcal{A}} \left\{ \mathbb{E} \left[r(s, a(s)) + \gamma \sum_{t=0}^{\infty} \gamma^t r(s_{t+1}, a(s_{t+1})) \right] \right\} = \\ &= \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \mathbb{E} [V^*(s')] \right\} = \\ &= \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') V^*(s') \right\}. \end{aligned}$$

DMDP. Уравнения Беллмана

Соответственно, мы можем провести аналогичные рассуждения для Q -функции ценности:

$$\begin{aligned} Q^*(s, a) &= \max_{\pi \in \Pi} \{ Q^\pi(s, a) \}; \\ Q^*(s, a) &= \mathbb{E} [r(s, a) + \gamma V^*(s') | s_0 = s, a_0 = a] = \\ &= r(s, a) + \gamma \mathbb{E} \left[\max_{a' \in \mathcal{A}} \{ Q^*(s', a') \} \right] = \\ &= r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') \max_{a' \in \mathcal{A}} \{ Q^*(s', a') \}. \end{aligned}$$

Критерий оптимальности относительно Q -функции ценности [1]
(детали: глава 3.1.10.).

Функция Q представляет собой оптимальную функцию ценности Q^* , если и только если она удовлетворяет уравнениям оптимальности Беллмана:

$$\begin{aligned} Q(s, a) &= \mathbb{E} \left[r(s, a(s)) + \gamma \max_{a' \in \mathcal{A}} \{ Q(s', a') \} \middle| s_0 = s, a_0 = a \right] = \\ &= \sum_{s' \in \mathcal{S}} p(s, a; s') \left[r(s, a) + \gamma \max_{a' \in \mathcal{A}} \{ Q(s', a') \} \right], \quad \forall s \in \mathcal{S}, a \in \mathcal{A}. \end{aligned}$$

Кроме того, детерминированная политика, определенная как

$$\pi(s) \in \operatorname{Arg} \max_{a \in \mathcal{A}} \{ Q^*(s, a) \},$$

есть оптимальная политика.

DMDP. Уравнения оптимальности Беллмана

Таким образом, для оптимальной политики π^* выполнены следующие соотношения [2]:

$$1) V^{\pi^*}(s) = \max_{a \in \mathcal{A}} \{Q^{\pi^*}(s, a)\}, \quad \forall s \in \mathcal{S};$$

$$2) Q^{\pi^*}(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') V^{\pi^*}(s'), \quad \forall s \in \mathcal{S}, a \in \mathcal{A};$$

$$3) V^{\pi^*}(s) = \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') V^{\pi^*}(s') \right\}, \quad \forall s \in \mathcal{S};$$

$$4) Q^{\pi^*}(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') \max_{a' \in \mathcal{A}} \{Q^{\pi^*}(s', a')\}, \quad s \in \mathcal{S}, a \in \mathcal{A}.$$

Соответствующая этим соотношениям детерминированная политика:

$$\pi^*(s) \in \operatorname{Arg max}_{a \in \mathcal{A}} \{Q^{\pi^*}(s, a)\} = \operatorname{Arg max}_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') V^{\pi^*}(s') \right\}.$$

DMDP. Уравнения оптимальности Беллмана

Основное свойство оптимальных V/Q -функций

Если V^*, Q^* – оптимальные функции, то функции $V^*(s) := V^*(s) + \alpha, Q^*(s, a) := Q^*(s, a) + \alpha, \forall s \in \mathcal{S}, a \in \mathcal{A}, \alpha \in \mathbb{R}$ так же соответствуют оптимальной политике.

$$\begin{aligned}\pi^*(s) \in \operatorname{Arg max}_{a \in \mathcal{A}} \{Q^*(s, a)\} &= \operatorname{Arg max}_{a \in \mathcal{A}} \{Q^*(s, a) + \alpha\} = \\ &= \operatorname{Arg max}_{a \in \mathcal{A}} \{Q^*(s, a)\} = \operatorname{Arg max}_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') V^*(s') \right\} = \\ &= \operatorname{Arg max}_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') \underbrace{(V^*(s') + \alpha)}_{=V^*(s')} \right\}, \forall s \in \mathcal{S}, \alpha \in \mathbb{R}\end{aligned}$$

DMDP. Уравнения оптимальности Беллмана

Для уравнения Беллмана верна связь с константным сдвигом функции награды $\hat{r}(s, a) := r(s, a) + \alpha(1 - \gamma)$, $\forall s \in \mathcal{S}, a \in \mathcal{A}, \alpha \in \mathbb{R}$:

$$\begin{aligned} Q^*(s, a) &= Q^*(s, a) + \alpha = (r(s, a) + \alpha(1 - \gamma)) + \\ &\quad + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') \max_{a' \in \mathcal{A}} \{Q^*(s', a') + \alpha\} = \\ &= \hat{r}(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') \max_{a' \in \mathcal{A}} \{Q^*(s', a')\} = \\ &= \hat{r}(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') V^*(s'); \end{aligned}$$

$$\begin{aligned} V^*(s) &= \max_{a \in \mathcal{A}} \left\{ (r(s, a) + \alpha(1 - \gamma)) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') (V^*(s') + \alpha) \right\} = \\ &= \max_{a \in \mathcal{A}} \left\{ \hat{r}(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') V^*(s') \right\} = \max_{a \in \mathcal{A}} \{Q^*(s, a)\}. \end{aligned}$$

DMDP. Уравнения оптимальности Беллмана

Аддитивное преобразование, не меняющее решение π^*

$\hat{r}(s, a) := \beta \left(r(s, a) + f(s) - \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') f(s') \right), \forall s \in \mathcal{S}, a \in \mathcal{A}, \beta > 0$, где $f : \mathcal{S} \mapsto \mathbb{R}$ – произвольное отображение.

$${}^a\pi^*(s) = \arg \max_{a \in \mathcal{A}} \{Q^*(s, a)\} = \arg \max_{a \in \mathcal{A}} \{\beta(Q^*(s, a) + f(s))\}, s \in \mathcal{S}.$$

$$\begin{aligned} Q^*(s, a) &:= \beta(Q^*(s, a) + f(s)) = \\ &= \beta(r(s, a) + f(s)) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') \max_{a' \in \mathcal{A}} \{\beta Q^*(s', a')\} = \\ &= \hat{r}(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') \max_{a' \in \mathcal{A}} \{\beta(Q^*(s', a') + f(s'))\} = \\ &= \hat{r}(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') \max_{a' \in \mathcal{A}} \{Q^*(s', a')\}. \end{aligned}$$

DMDP. Уравнения оптимальности Беллмана

Для функции награды, зависящей от $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$, аддитивное преобразование, инвариантное относительно оптимальной политики, выглядит проще: $\hat{r}(s, a, s') := \beta (r(s, a, s') + f(s) - \gamma f(s'))$.

$$\begin{aligned} Q^*(s, a) &:= \beta (Q^*(s, a) + f(s)) = \\ &= \sum_{s' \in \mathcal{S}} p(s, a; s') (\beta (r(s, a, s') + f(s) - \gamma f(s')) + \\ &\quad + \gamma \max_{a' \in \mathcal{A}} \{ \beta (Q^*(s', a') + f(s')) \}) = \\ &= \sum_{s' \in \mathcal{S}} p(s, a; s') \left(\hat{r}(s, a, s') + \gamma \max_{a' \in \mathcal{A}} \{ Q^*(s', a') \} \right). \end{aligned}$$

Предположения

Будем считать, что о нашей среде нам известны следующие функции:

- $p(s'|s, a)$ — вероятность попасть в состояние s' из состояния s с помощью действия a
- $r(s, a)$ — награда за выполнения действия a в состоянии s

Замечание

Эти допущения существенно сужают круг задач, которые мы можем решить, однако алгоритмы, предложенные в этой лекции будут оптимальными в данной постановке.

Построение алгоритма обучения π

Сам процесс можно разделить на два этапа: оценка качества текущей политики и поиск следующего приближения оптимальной политики.

Дополнительно предположим дискретность и конечность пространств \mathcal{S} и \mathcal{A} . Начнём с оценивания — вычислим $V^\pi(s)$ и $Q^\pi(s, a)$:

$$V^\pi(s) = \underbrace{\mathbb{E}_{\pi(a|s)} [r(s, a)]}_{:=u(s)} + \underbrace{\gamma \mathbb{E}_{\pi(a|s)} \mathbb{E}_{p(s'|s,a)} [V^\pi(s')]}_{:=F(V^\pi(s))}, \quad \forall s \in \mathcal{S}.$$

Представим выражение выше в виде матрично векторных операций.

V^π — вектор ценностей состояний. P — матрица вероятностей:

$P_{ss'} = \sum_{a \in \mathcal{A}} \pi(a|s) p(s'|s, a) = p(s'|s)$. u — вектор средних наград за шаг.

Построение алгоритма обучения π

$$\begin{aligned} V^\pi(s) &= \sum_{a \in \mathcal{A}} \pi(a|s) r(s, a) + \gamma \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s' \in \mathcal{S}} p(s'|s, a) V^\pi(s') = \\ &= u(s) + \gamma \sum_{s' \in \mathcal{S}} V^\pi(s') p(s'|s) \implies V^\pi = F(V^\pi) = u + \gamma P V^\pi. \end{aligned}$$

Полученное отображение F является сжимающим, для произвольных векторов V и W :

$$\begin{aligned} \|F(V) - F(W)\|_\infty &= \|u + \gamma P V - u - \gamma P W\|_\infty = \gamma \|P(V - W)\|_\infty \leq \\ &\leq \gamma \|P\|_\infty \|V - W\|_\infty = \gamma \|V - W\|_\infty, \end{aligned}$$

$$\|P\|_\infty = \max_{x \neq 0} \left\{ \frac{\|Px\|_\infty}{\|x\|_\infty} \right\} = \max_{\|x\|_\infty \leq 1} \max_{s \in \mathcal{S}} \left\{ \sum_{s' \in \mathcal{S}} p(s'|s) \underbrace{x_{s'}}_{=1} \right\} = 1.$$

Построение алгоритма обучения π

Получили алгоритм Iterative Policy Evaluation для оценки политики π с заданной точностью $\varepsilon > 0$:

1) Инициализировать $V(s), \forall s \in \mathcal{S}$;

2) Повторять в цикле:

2.1) $\Delta := 0$

2.2) для всех $s \in \mathcal{S}$:

$$\delta := V(s);$$

$$V(s) := \mathbb{E}_{\pi(a|s)} [r(s, a)] + \gamma \mathbb{E}_{\pi(a|s)} \mathbb{E}_{p(s'|s,a)} [V(s')] ;$$

$$\Delta := \max\{\Delta, |\delta - V(s)|\}.$$

2.3) Если $\Delta \leq \varepsilon$, то выход, иначе — переход на шаг 2.

Построение алгоритма обучения π

Теперь построим процедуру улучшения оценённой политики π :

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} [V^\pi(s')].$$

Определение

Политика $\hat{\pi} \succeq \pi$ (монотонно лучше политики π), если $V^{\hat{\pi}}(s) \geq V^\pi(s)$, $\forall s \in \mathcal{S}$.

Таким образом, изменяя действие для одного состояния $\hat{s} \in \mathcal{S}$ мы производим улучшение π :

$$\hat{\pi}(a|\hat{s}) = \delta \left\{ \arg \max_{\hat{a} \in \mathcal{A}} \{Q^\pi(\hat{s}, \hat{a})\} \right\} (a), \quad \hat{\pi}(\hat{s}) := \arg \max_{\hat{a} \in \mathcal{A}} \{Q^\pi(\hat{s}, \hat{a})\};$$

$$\hat{\pi}(a|s) = \pi(a|s), \quad \forall s \in \mathcal{S}, \quad s \neq \hat{s}.$$

То есть $Q^\pi(s, \hat{\pi}(s)) \geq V^\pi(s)$.

Построение алгоритма обучения π

Теорема об улучшении политики [1] (детали: глава 3.2.3, теорема 17)

Пусть π и π' – любая пара детерминированных политик, таких, что

$$\forall s \in \mathcal{S} \quad Q^{\pi}(s, \pi'(s)) \geq V^{\pi}(s). \quad (1)$$

Тогда π' должна быть не хуже, чем π , то есть ценность не хуже $\forall s \in \mathcal{S}$:

$$V^{\pi'}(s) \geq V^{\pi}(s). \quad (2)$$

Более того, если в каком-либо состоянии существует строгое (1), то и (2) должно быть строгим.

Построение алгоритма обучения π

Действительно,

$$\begin{aligned} V^\pi(s) &\leq Q^\pi(s, \underbrace{\hat{\pi}(s)}_{\text{детерминирован}}) = r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{p(s'|s,a)} [V^\pi(s')] \leq \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{p(s'|s,a)} [Q^\pi(s', \hat{\pi}(s'))] \leq \dots \leq \\ &\leq r(s, \hat{\pi}(s)) + \gamma \mathbb{E}_{p(s'|s,\hat{\pi}(s))} [r(s', \hat{\pi}(s')) + \gamma^2 \dots] = V^{\hat{\pi}}(s), \quad \forall s \in \mathcal{S}. \end{aligned}$$

Шаг для обновления всей политики

$$\pi_{\text{new}}(s) := \arg \max_{a \in \mathcal{A}} \{Q^{\pi_{\text{old}}}(s, a)\}, \quad \forall s \in \mathcal{S}.$$

Построение алгоритма обучения π

Если после очередного шага обновления политики получилось, что $Q^\pi(s, \hat{\pi}(s)) = V^\pi(s)$, $\forall s \in \mathcal{S}$, то это означает удовлетворение уравнению Беллмана:

$$\hat{\pi} = \pi;$$

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{p(s'|s,a)} [V^\pi(s')] , \quad \forall s \in \mathcal{S}, a \in \mathcal{A}.$$

Построение алгоритма обучения π

Мы получили алгоритм Policy Iteration:

- 1) Инициализируем $V(s)$ и $\pi(a|s)$ для всех $s \in \mathcal{S}$.
- 2) Оценить $V^\pi(s)$ для текущей π , используя Iterative Policy Evaluation.
- 3) Улучшаем политику:
 - 3.1) $\text{Flag} := \text{True}$;
 - 3.2) Для всех $s \in \mathcal{S}$:

$$a = \pi(s);$$

$$\pi(s) := \arg \max_{a \in \mathcal{A}} \{ r(s, a) + \gamma \mathbb{E}_{p(s'|s,a)} [V^\pi(s')] \};$$

если $a \neq \pi(s)$, то $\text{Flag} := \text{False}$.

- 4) Если $\text{Flag} = \text{True}$, то выход, иначе — шаг 2.

Построение алгоритма обучения π

Описанную ранее процедуру поиска оптимальной политики можно представить как последовательность монотонно улучшающихся политик и функций ценности:

$$\pi_0 \xrightarrow{E} V^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V^{\pi_1} \xrightarrow{I} \dots \xrightarrow{I} \pi^* \xrightarrow{E} V^*,$$

где \xrightarrow{E} обозначает оценку политики, \xrightarrow{I} – улучшение политики, то есть получен алгоритм итеративной оптимизации политики.

Построение алгоритма обучения π

Уравнение Беллмана относительно фиксированной политики по сути решается простым итеративным способом. Начальное приближение V_0 выбирается произвольно, а каждая последовательная итерация реализуется согласно уравнению Беллмана:

$$V_{t+1}(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \sum_{s' \in \mathcal{S}} p(s, a; s') (r(s, a) + \gamma V_t(s')) ,$$

где $V_t = V^\pi$ – фиксированная точка. Для получения каждого последующего приближения, V_{t+1} из V_t при итеративной оценке политики применяется та же операция к каждому состоянию s , и ее называют ожидаемым обновлением, а данный алгоритм – итерации функции ценности.

Построение алгоритма обучения π

Для решения уравнения оптимальности Беллмана можно проводить следующую процедуру:

$$V_{t+1}(s) = \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') V_t(s') \right\}, \quad \forall s \in \mathcal{S};$$
$$\pi_{t+1}(s) \in \text{Arg max}_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') V_{t+1}(s') \right\}, \quad \forall s \in \mathcal{S}.$$

Причём начальное значение V_0 может быть произвольным, а в качестве критерия останова может выступать $\|V_{t+1} - V_t\|_\infty \leq \varepsilon$.

Построение алгоритма обучения π

На предыдущем слайде предложен по сути частный случай алгоритма Policy Iteration — Value Iteration (вместо шагов 2 и 3 один шаг делаем):

- 1) Инициализируем $V(s)$ для всех $s \in \mathcal{S}$.
- 2) Повторять:
 - 2.1) $\Delta := 0$;
 - 2.2) Для всех $s \in \mathcal{S}$:

$$v = V(s)$$

$$V(s) := \max_{a \in \mathcal{A}} \{ r(s, a) + \gamma \mathbb{E}_{p(s'|s,a)} [V(s')] \}$$

$$\Delta := \max\{\Delta, |v - V(s)|\}$$

- 2.3) Если $\Delta \leq \varepsilon$, то выход, иначе — переход на шаг 2.

Мы по сути схлопнули Policy Iteration и Policy Improvement, не вычисляя π .

Построение алгоритма обучения π

В результате работы алгоритма Value Iteration оптимизированная политика вычисляется следующим образом:

$$\pi(s) = \arg \max_{a \in \mathcal{A}} \{ r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} [V(s')] \} .$$

Сходимость алгоритма Value Iteration

Для оптимальной V -функции выполнено соотношение

$$V^*(s) = V^{\pi^*}(s) = \max_{a \in \mathcal{A}} \{r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} [V^*(s')]\}, \quad \forall s \in \mathcal{S}.$$

Рассмотрим

$$V^{\pi_{t+1}}(s), \quad \pi_{t+1}(s) = \arg \max_{a \in \mathcal{A}} \{r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} [V^{\pi_{t+1}}(s')]\}.$$

Построим оценку на невязку по V -функции,

$$V^*(s) \geq V^{\pi_t}(s), \quad t \in \mathbb{Z}_+, \quad s \in \mathcal{S}:$$

$$\begin{aligned} \max_{s \in \mathcal{S}} \{|V^*(s) - V^{\pi_{t+1}}(s)|\} &= \max_{s \in \mathcal{S}} \left\{ \max_{a \in \mathcal{A}} \{r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} [V^*(s')]\} - \right. \\ &\quad \left. - \max_{a' \in \mathcal{A}} \{r(s, a') + \gamma \mathbb{E}_{p(s''|s, a')} [V^{\pi_t}(s'')]\} \right\} \leq \gamma \max_{\substack{s \in \mathcal{S} \\ a \in \mathcal{A}}} \{ \mathbb{E}_{p(s'|s, a)} [V^*(s')] - \\ &\quad - \mathbb{E}_{p(s''|s, a)} [V^{\pi_t}(s'')] \} = \gamma \max_{s \in \mathcal{S}, a \in \mathcal{A}} \{ \mathbb{E}_{p(s'|s, a)} [V^*(s') - V^{\pi_t}(s')] \} \leq \end{aligned}$$

Сходимость алгоритма Value Iteration

продолжим с последнего неравенства, раскрывая рекуррентную зависимость:

$$\begin{aligned} &\leq \gamma \max_{s' \in \mathcal{S}} \{V^*(s') - V^{\pi_t}(s')\} = \gamma \max_{s' \in \mathcal{S}} \{|V^*(s') - V^{\pi_t}(s')|\} \Rightarrow \\ &\Rightarrow \max_{s \in \mathcal{S}} \{|V^{\pi_t}(s) - V^*(s)|\} \leq \gamma^t \max_{s' \in \mathcal{S}} \{|V^{\pi_0}(s') - V^*(s')|\}. \end{aligned}$$

В качестве $V^{\pi_0}(s)$, $s \in \mathcal{S}$ можно взять произвольное начальное приближение, например, тождественную по всем состояниям константу. Имеем оценку относительно внешних итераций:

$$\max_{s \in \mathcal{S}} \{|V^{\pi_t}(s) - V^*(s)|\} = \mathcal{O}(\gamma^t), \quad t \in \mathbb{Z}_+, \quad \gamma \in (0, 1).$$

Value Iteration. Сложность решения

Введём следующий оператор $T : \mathbb{R}^S \mapsto \mathbb{R}^S$:

$$T(V)_s := \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') V(s') \right\}. \text{ Теперь обновление}$$

V -функции в Value Iteration выражается следующим образом:

$V_{t+1} := T(V_t)$, а V -функция кодируется вещественным вектором.

Теорема [3]

Существует DMDP, для которого последовательность оценок V -функции удовлетворяет:

$V_0 = 0_S$, $V_{n+1} \in \text{span} \{ V_0, \dots, V_n, T(V_0), \dots, T(V_n) \}$, $n \in \mathbb{Z}_+$, со следующим свойством $\forall n = \overline{0, N-1}$:

$$\|V_n - V^*\|_\infty \geq \frac{\gamma^n}{1 + \gamma}, \quad V^* = T(V^*).$$

Value Iteration. Сложность решения

Доказательство. Для произвольного $\gamma \in (0, 1)$ предложим DMDP с N состояниями и с одним действием. Награда для первого состояния $r_1 := 1$, для остальных состояний $r_i := 0, i = \overline{2, N}$. Действие из первого состояния оставляет в нём же, действие из $(i + 1)$ -го состояния переводит в i -ое. Оптимальное значение V -функции следующее:

$V^*(i) = \frac{\gamma^{i-1}}{1-\gamma}$. Рассмотрим последовательность векторов $(V_n)_{n \geq 0}$, $V_0 = 0_S$ и

$$V_{n+1} \in \text{span} \{ V_0, \dots, V_n, T(V_0), \dots, T(V_n) \}, n \geq 0.$$

Докажем через раскрытие рекурсии, что $\forall n \geq 0, i \in S$ имеем $V_n(i) = 0$, если $i \geq n + 1$. Это верно для $n = 0$, так как $V_0 = 0_S$. Предположим, что верно и для V_0, \dots, V_{n-1} . По определению оператора T и в силу того, что $r_i = 0, i \geq 2$, имеем $T(V_t)_i = 0$, если $i \geq t + 2, \forall t \leq n - 1$.

Value Iteration. Сложность решения

Следовательно, в силу $V_{n+1} \in \text{span} \{V_0, \dots, V_n, T(V_0), \dots, T(V_n)\}$ заметим, что $V_n(i) = 0$, если $i \geq n+1$, и мы доказали нашу рекурсию. $r_1 > 0$ – единственное, по сути мы доказали, что для любого метода первого порядка требуется $n - 1$ шаг для распространения награды первого состояния до состояния $1 \leq n \leq N$.

Value Iteration. Сложность решения

Теперь мы имеем для $1 \leq n \leq N - 1$:

$$\|V_n - T(V_n)\|_\infty = \|V_n - T(V_n) - (V^* - T(V^*))\|_\infty \quad (3)$$

$$\geq (1 - \gamma) \|V_n - V^*\|_\infty \quad (4)$$

$$= (1 - \gamma) \max_{1 \leq i \leq N} \{|V_n(i) - V^*(i)|\}$$
$$\geq (1 - \gamma) \max_{n+1 \leq i \leq N} \{|V_n(i) - V^*(i)|\}$$

$$\geq (1 - \gamma) \max_{n+1 \leq i \leq N} \{|V^*(i)|\} \quad (5)$$

$$\geq (1 - \gamma) \max_{n+1 \leq i \leq N} \left\{ \frac{\gamma^{i-1}}{1 - \gamma} \right\}$$
$$\geq \gamma^n,$$

где (3) следует из $V^* = T(V^*)$,

Value Iteration. Сложность решения

(4) следует из (6), и (5) следует из $V_n(i) = 0$ для $i \geq n + 1$. Можем заключить в силу (6):

$$\begin{aligned}\|V_n - V^*\|_\infty &\geq \frac{1}{1 + \gamma} \cdot \|V_n - T(V_n) - (V^* - T(V^*))\|_\infty = \\ &= \frac{1}{1 + \gamma} \cdot \|V_n - T(V_n)\|_\infty \geq \frac{\gamma^n}{1 + \gamma}.\end{aligned}$$

$\forall V, W \in \mathbb{R}^S$ (проверяется непосредственно) :

$$(1 - \gamma) \cdot \|V - W\|_\infty \leq \|(I - T)(V) - (I - T)(W)\|_\infty; \quad (6)$$

$$\|(I - T)(V) - (I - T)(W)\|_\infty \leq (1 + \gamma) \cdot \|V - W\|_\infty.$$

Утверждение доказано.

Value Iteration. Сложность решения

Из доказанного утверждения вместе с утверждением о сходимости Value Iteration следует оптимальность алгоритма Value Iteration:

$$\frac{\gamma^t}{1-\gamma} = \gamma^t \|V_0 - V^*\|_\infty \geq \|V_t - V^*\|_\infty \geq \frac{\gamma^t}{1+\gamma}, \quad t \in \mathbb{Z}_+, \quad \gamma \in (0, 1);$$
$$\|V_t - V^*\|_\infty = \Theta(\gamma^t), \quad V_0 = 0_N.$$

LP-релаксация MDP. AMDP [4]

Для AMDP LP-задача вводится следующим образом:

$$V^*(s) = \max_{a(\cdot) \in \mathcal{A}} \left\{ \lim_{H \rightarrow \infty} \frac{1}{H} \mathbb{E} \left[\sum_{t=0}^{H-1} r(s_t, a_t(s_t)) \middle| s_0 = s \right] \right\},$$

где H – эпизодическое ограничение, то есть максимальная длина эпизода. В случае эпизодов конечной длины предел опускается и используется максимальное значение H . Для политики $\pi(a|s)$ можно определить стационарное распределение:

$$v_\pi(s') = \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} p(s, a; s') \pi(a|s) v_\pi(s), \quad s' \in \mathcal{S},$$

которое соответствует своему вектору из вероятностей

$$v_\pi = (v_\pi(s))_{s \in \mathcal{S}}.$$

LP-релаксация MDP. AMDP [4]

И если MDP равномерно эргодично, то:

$$V^\pi := V(\pi) := \lim_{H \rightarrow \infty} \frac{1}{H} \mathbb{E} \left[\sum_{t=0}^{H-1} r(s_t, a_t(s_t)) \right] = \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} r(s, a) \pi(a|s) v_\pi(s).$$

Напоминаем, что в данном случае равномерная эргодичность соответствует:

$$\max_{i=1,S} \{ \| (P^\pi)^n e_i - v_\pi \|_\infty \} \rightarrow 0, \quad n \rightarrow \infty, \quad e_i^\top = (0, \dots, 0, \underbrace{1}_i, 0, \dots, 0).$$

LP-релаксация MDP. AMDP [4]

Вводится распределение действий по состояниям –

$\mu(s, a) = v_\pi(s)\pi(a|s)$, следовательно, можно переписать задачу поиска оптимальной политики в AMDP как задачу LP со смыслом оценки ценности политики по распределению μ :

$$\max_{\mu \in \Delta^{\mathcal{S} \times \mathcal{A}}} \left[V(\mu) = \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} r(s,a) \mu(s,a) = \langle r, \mu \rangle : \sum_{b \in \mathcal{A}} \mu(s', b) = \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} p(s, a; s') \mu(s,a), s' \in \mathcal{S} \right];$$

$$\Delta^{\mathcal{S} \times \mathcal{A}} = \left\{ \mu : \mu(s,a) \geq 0, \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \mu(s,a) = 1 \right\}, \pi_\mu(a|s) = \frac{\mu(s,a)}{\sum_{b \in \mathcal{A}} \mu(s,b)}.$$

LP-релаксация MDP. AMDP [4]

Данную задачу можно напрямую переписать в матричной форме:

$$\begin{aligned} & \max_{\mu \in \Delta^{\mathcal{S} \times \mathcal{A}}} \langle r, \mu \rangle ; \\ & s.t. (\hat{I} - P)\mu = 0. \end{aligned}$$

Единичная матрица \hat{I} имеет нестандартный формат: это прямоугольная матрица размера $S \times (SA)$, на каждой строке $s \in \mathcal{S}$ только элементы, соответствующие паре (s, a) , $a \in \mathcal{A}$, равняются единице, остальные элементы данной строки равняются нулю, то есть на каждой строке \hat{I} ровно A единиц. У матрицы P размера $S \times (SA)$ в каждом столбце $(s, a) \in \mathcal{S} \times \mathcal{A}$ записано распределение $P(\cdot | s, a)$.

LP-релаксация MDP. AMDP [4]

Для этой задачи LP напрямую строится двойственная задача, с условием, что $\mu \geq 0$, которая имеет смысл оценки ценности оптимальной политики через V -функцию:

$$\begin{aligned} \min_{\bar{V} \in \mathbb{R}, V \in \mathbb{R}^{|S|}} \quad & \bar{V}; \\ \text{s.t.} \quad & r - \bar{V} \cdot 1_{SA} - (\hat{I} - P)^{\top} V \leq 0. \end{aligned}$$

Таким образом, имеет место уравнение оптимальности Беллмана со средним вознаграждением:

$$V(s) = \max_{a \in \mathcal{A}} \left\{ r(s, a) - V^* + \sum_{s' \in \mathcal{S}} p(s, a; s') V(s') \right\}, \quad V^* = \langle r, \mu^* \rangle,$$

полученное из ограничений вида неравенства:

$$\hat{I}^{\top} V \geq r - \bar{V} \cdot 1_{SA} + P^{\top} V.$$

LP-релаксация MDP. AMDP [4]

Таким образом, для AMDP мы получили аналог уравнения Беллмана — уравнение Пуассона:

$$V^\pi + V^\pi(s) = \mathbb{E}_\pi [r(s, a)] + \mathbb{E}_{p, \pi} [V^\pi(s')] , \quad \forall s \in \mathcal{S}.$$

Множество его решений:

$$\{(V^\pi(\cdot) + ce(\cdot), V^\pi) : c \in \mathbb{R}\} , \quad e(s) = 1, \quad \forall s \in \mathcal{S}.$$

Критерий оптимальной политики:

$$V^* + V^*(s) = \max_{a \in \mathcal{A}} \{r(s, a) + \mathbb{E}_p [V^*(s')]\} , \quad \forall s \in \mathcal{S};$$

$$\pi^* \in \operatorname{Argmax}_{\pi \in \Pi} \{V^\pi\} , \quad V^* \equiv V^{\pi^*}.$$

LP-релаксация MDP. AMDP [4]

Уравнение Пуассона также существует для Q -функции ценности:

$$V^\pi + Q^\pi(s, a) = r(s, a) + \mathbb{E}_{p, \pi} [Q^\pi(s', a')] , \quad \forall s \in \mathcal{S}, a \in \mathcal{A}.$$

Множество его решений:

$$\{(Q^\pi(\cdot) + c\bar{e}(\cdot), V^\pi) : c \in \mathbb{R}\} , \quad \bar{e}(s, a) = 1, \quad \forall s \in \mathcal{S}, a \in \mathcal{A}.$$

Критерий оптимальной политики:

$$V^* + Q^*(s, a) = r(s, a) + \mathbb{E}_p \left[\max_{a' \in \mathcal{A}} \{Q^*(s', a')\} \right] , \quad \forall s \in \mathcal{S}, a \in \mathcal{A};$$

$$\pi^* \in \operatorname{Argmax}_{\pi \in \Pi} \{V^\pi\} .$$

LP-релаксация MDP. DMDP [4]

Для DMDP задача LP записывается в следующем виде (q – распределение начального состояния μ_0 в виде вектора):

$$\begin{aligned} \min_{V \in \mathbb{R}^{|\mathcal{S}|}} \quad & \langle q, V \rangle; \\ \text{s.t.} \quad & r - (\hat{I} - \gamma P)^\top V \leq 0. \end{aligned}$$

И ей соответствует такая двойственная задача:

$$\begin{aligned} \max_{\mu \in \Delta^{\mathcal{S} \times \mathcal{A}}} \quad & \langle r, \mu \rangle; \\ \text{s.t.} \quad & (\hat{I} - \gamma P)\mu = q. \end{aligned}$$

LP-релаксация MDP. DMDP [5]

Существует также постановка задачи LP для ограниченного DMDP – Constrained Markov Decision Process, CMDP:

$$\begin{aligned} & \max_{\mu \in \Delta^{\mathcal{S} \times \mathcal{A}}} \langle r, \mu \rangle ; \\ & s.t. (\hat{I} - \gamma P)\mu = q, \quad D\mu \geq c. \end{aligned}$$

По сравнению с предыдущими задачами линейного программирования вводится дополнительно аффинное ограничение вида неравенства:
 $D\mu \geq c$.

LP-релаксация MDP. DMDP [5]

Заметим, что вместо обозначенного ранее распределения $\mu(s, a) = \nu_\pi(s)\pi(a|s)$ может быть полезно рассмотреть:

$$\mu(s, a) := \mu^\pi(s, a) = \mathbb{E}_{s_0 \sim \mu_0} \left[\sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s, a_t = a | s_0) \right].$$

Или даже масштабированную сумму сверху в виде корректно определённой вероятностной меры:

$$\begin{aligned} \mu(s, a) &:= \tilde{\mu}^\pi(s, a) = (1 - \gamma)\mu^\pi(s, a) = \\ &= \mathbb{E}_{s_0 \sim \mu_0} \left[(1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s, a_t = a | s_0) \right]. \end{aligned}$$

LP-релаксация MDP. DMDP [5]

В обоих случаях получается одна и та же политика:

$$\pi(a|s) = \frac{\mu^\pi(s, a)}{\sum_{b \in \mathcal{A}} \mu^\pi(s, b)} = \frac{\tilde{\mu}^\pi(s, a)}{\sum_{b \in \mathcal{A}} \tilde{\mu}^\pi(s, b)}.$$

Связь решения DMDP с выпуклой оптимизацией

Обозначим за $v \in \mathbb{R}^S$ вектор, кодирующий V -функцию ценности. Тогда решение уравнения $v = T(v)$ соответствует поиску стационарной точки некоторой функции $f : \mathbb{R}^S \mapsto \mathbb{R}$ со следующим градиентом и (суб)гессианом:

$$\nabla f(v) := v - T(v) = (I - T)(v) =: F(v);$$

$$\nabla^2 f(v) := I - \gamma \cdot P^{\pi^v} =: \partial F(v) \succeq 0;$$

$$T(v) = r^{\pi^v} + \gamma \cdot P^{\pi^v} v, \quad \left(r^{\pi^v}\right)_s := r(s, a_s^v), \quad s \in S,$$

$$\pi^v(s) := a_s^v \in \operatorname{Arg\,max}_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \sum_{s' \in S} p(s, a; s') V(s') \right\}.$$

Таким образом, можно ввести формально следующую функцию:

$$f(v) := \frac{1}{2} \left\langle \left(I - \gamma \cdot P^{\pi^v} \right) v, v \right\rangle - \langle r^{\pi^v}, v \rangle + \text{const}.$$

Связь решения DMDP с выпуклой оптимизацией [6]

Свойства функции $f(\cdot)$, $\forall v, w \in \mathbb{R}^S$:

$$(1 - \gamma) \cdot \|v - w\|_\infty \leq \|\nabla f(v) - \nabla f(w)\|_\infty \leq (1 + \gamma) \cdot \|v - w\|_\infty;$$

$$\frac{1 - \gamma}{\sqrt{S}} \|v - w\|_2 \leq \|\nabla f(v) - \nabla f(w)\|_\infty \leq (1 + \gamma) \cdot \|v - w\|_2;$$

$$\frac{1 - \gamma}{\sqrt{S}} \|v - w\|_2 \leq \|\nabla f(v) - \nabla f(w)\|_2 \leq \sqrt{S}(1 + \gamma) \cdot \|v - w\|_2.$$

Связь решения DMDP с выпуклой оптимизацией [6]

Мы с помощью введённых обозначений можем переписать алгоритм Policy Iteration как шаг метода Ньютона для решения нелинейного уравнения $F(v) = 0$, $v \in \mathbb{R}^S$:

$$v_{t+1} := v_t - (\partial F(v_t))^{-1} F(v_t), \quad t \in \mathbb{Z}_+.$$

Данная процедура не обладает глобальной сходимостью.

Связь решения DMDP с выпуклой оптимизацией [6]

Сглаженный оператор оптимальности Беллмана, $\beta > 0$:

$$T_{\beta}(v)_s := \frac{1}{\beta} \ln \left(\sum_{a \in \mathcal{A}} \exp \left(\beta \left(r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s, a; s') V(s') \right) \right) \right), \forall s \in \mathcal{S}.$$

Если $v_{\beta}^* = T_{\beta}(v_{\beta}^*)$ и $v^* = T(v^*)$, то:

$$\|v_{\beta}^* - v^*\|_{\infty} \leq \frac{\gamma \ln(A)}{\beta(1 - \gamma)}.$$

Метод Ньютона для решения $F_{\beta}(v) := v - T_{\beta}(v) = 0$ сходится квадратично для любого начального приближения.

Q-обучение

Несложно заметить, что если

$$V^*(s) = \max_{a \in \mathcal{A}} Q(s, a),$$

то Q^* -функция должна удовлетворять Q -уравнению:

$$Q(s, a) = \sum_{s' \in \mathcal{S}} p(s, a; s') \left(r(s, a; s') + \gamma \max_{a' \in \mathcal{A}} Q(s', a') \right),$$

в текущем случае рассматривается уравнение Беллмана для более общего процесса MDP, в котором награда зависит уже от (s, a, s') , то есть добавилось ещё и следующее за s состояние s' .

Q-обучение

С такой зависимостью наград удобнее рассматривать траекторию τ_{H-1} политики π как набор четвёрок (s, a, r, s') со следующим правдоподобием:

$$\tau_{H-1} = (s_0, a_0, r_0, s_1, a_1, r_1, s_2, \dots, a_{H-1}, r_{H-1}, s_H);$$

$$P(\tau_{H-1}|\pi) = \mu_0(s_0) \prod_{t=0}^{H-1} (\pi(a_t|s_t)p(s_t, a_t; s_{t+1})).$$

Q-обучение

Данное уравнение Беллмана может быть решено методом простых итераций; если смотреть на $Q = \{Q(s, a)\}_{s \in \mathcal{S}, a \in \mathcal{A}}$ как на вектор, то можно записать в операторном виде $Q = F(Q)$ (метод простых итераций будет иметь вид $Q_{t+1} = F(Q_t)$), где по определению оператор в правой части F является сжимающим с коэффициентом γ в норме Чебышева:

$$\max_{s \in \mathcal{S}, a \in \mathcal{A}} |F(\tilde{Q}(s, a)) - F(Q(s, a))| \leq \gamma \max_{s \in \mathcal{S}, a \in \mathcal{A}} |\tilde{Q}(s, a) - Q(s, a)|,$$

\tilde{Q} и Q – произвольные.

Q-обучение

Основная идея Q-обучения заключается в замене невычислимой правой части в уравнении $Q_{t+1} = F(Q_t)$ на ее вычислимую несмещенную оценку:

$$Q_{t+1}(s, a) = Q_t(s, a) + \alpha_t(s, a) \left(r(s, a; s'(s, a)) + \gamma \max_{a' \in \mathcal{A}} Q_t(s'(s, a), a') - Q_t(s, a) \right), \quad (7)$$

где $s'(s, a)$ – положение процесса на шаге $t + 1$, если на шаге t процесс был в состоянии s и было выбрано действие a , то параметр $0 < \alpha_t(s, a) \leq 1$, иначе $\alpha_t(s, a) = 0$. Правая часть (7) следует из перехода в $s'(s, a)$, $\{Q_t(s, a)\}_{s \in \mathcal{S}, a \in \mathcal{A}}$ известно с прошлой итерации (можно посчитать $\max_{a' \in \mathcal{A}} Q_t(s'(s, a), a')$).

Q-обучение

Вознаграждение $r(s, a; s'(s, a))$ получается при переходе из состояния s при действии a в $s'(s, a)$, в ненаблюдаемых случаях $\alpha_t(s, a) = 0$, то есть значение r не интересно. Подход Q-обучения в отличие от ранее рассмотренных полагается на сэмплирование непосредственно траекторий из MDP, что освобождает от знания всего пространства $\mathcal{S} \times \mathcal{A}$ в каждый конечный момент времени, при этом описанный процесс вычисления Q_{t+1} всё также реализован через сжимающее отображение, гарантирующее асимптотическую сходимость к оптимальной политике.

Теорема [7]

Если при стратегии $a(s)$ с вероятностью 1 каждая пара (s, a) будет неограниченное число раз встречаться на бесконечном горизонте наблюдения, то при

$$\sum_{t=0}^{\infty} \alpha_t(s, a) = \infty, \quad \sum_{t=0}^{\infty} \alpha_t^2(s, a) < \infty,$$

следует сжимаемость (7):

$$\lim_{t \rightarrow \infty} Q_t(s, a) = Q(s, a), \quad V^*(s) = \max_{a \in \mathcal{A}} Q(s, a).$$

Таким образом, после достаточно большого числа шагов, даже в отсутствие какой-либо информации об управляемом марковском процессе, можно определить оптимальную стратегию

$$a(s) \in \operatorname{Arg} \max_{a \in \mathcal{A}} Q(s, a).$$

Q-обучение

Расширение Q-обучения для θ -параметрической аппроксимации функций ($Q_\theta; \theta \in \mathbb{R}^d$) выражается через следующее соотношение:

$$\theta_{t+1} = \theta_t + \alpha_t(s, a) \left\{ r(s, a; s'(s, a)) + \gamma \max_{a' \in \mathcal{A}} Q_{\theta_t}(s'(s, a), a') - Q_{\theta_t}(s, a) \right\} \nabla_\theta Q_{\theta_t}(s_t, a_t).$$

Для линейной параметризации $Q_\theta = \theta^\top \varphi$, $\varphi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ обновление параметров выглядит следующим образом [8]:

$$\begin{aligned} 1 : \delta &\leftarrow r(s_t, a_t; s_{t+1}) + \gamma \cdot \max_{a' \in \mathcal{A}} \theta^\top \varphi(s_{t+1}, a') - \theta^\top \varphi(s_t, a_t); \\ 2 : \theta &\leftarrow \theta + \alpha_t(s_t, a_t) \cdot \delta \cdot \varphi(s_t, a_t). \end{aligned}$$

Основные подходы в Deep Reinforcement Learning

- *On-policy* алгоритмы — алгоритмы, которые оценивают и улучшают ту же самую политику, которую используют для выбора действий (Target Policy = Behavior Policy).
- *Off-policy* алгоритмы — алгоритмы, которые оценивают и улучшают одну политику, а для выбора действий используют другую политику (Target Policy \neq Behavior Policy).

Основные подходы в Deep Reinforcement Learning

DQN

- off-policy
- одношаговое оценивание политики (смещённая)
- ϵ -жадная политика
- учим оптимальную Q -функцию Q^*

Policy Gradient

- on-policy
- оценка до конца эпизода (большая дисперсия)
- обучение явной политики как распределения $\pi(a|s)$
- учим V^π

Первый способ вывода policy gradient

Постановка задачи

Рассмотрим θ -параметризованное семейство политик $\pi(a|s, \theta)$ или $\pi_\theta(a|s)$. Тогда будем максимизировать следующую величину:

$$V^\pi := J(\theta) := \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t \geq 0} \gamma^t r_t \right] \rightarrow \max_{\theta}.$$

Первый способ вывода policy gradient

Подсчет градиента

$$V^\pi := \mathbb{E}_{a \sim \pi_\theta} Q^{\pi_\theta}(s, a) = \int_A \pi_\theta(a|s) Q^{\pi_\theta}(s, a) da,$$

тогда

$$\begin{aligned} \nabla_\theta J(\theta) &= \nabla_\theta \int_A \pi_\theta(a|s) Q^{\pi_\theta}(s, a) da = \int_A \nabla_\theta [\pi_\theta(a|s) Q^{\pi_\theta}(s, a)] da = \\ &= \int_A \nabla_\theta \pi_\theta(a|s) Q^\pi(s, a) da + \int_A \pi(a|s) \nabla_\theta Q^{\pi_\theta}(s, a) da. \end{aligned}$$

Первый способ вывода policy gradient

Трюк производной логарифма

$$\nabla_{\theta} p_{\theta}(x) = p_{\theta}(x) \frac{\nabla_{\theta} p_{\theta}(x)}{p_{\theta}(x)} = p_{\theta}(x) \nabla_{\theta} \log p_{\theta}(x)$$

Первое слагаемое

$$\begin{aligned} \int_A \nabla_{\theta} \pi_{\theta}(a|s) Q^{\pi}(s, a) da &= \int_A \pi_{\theta}(a|s) \nabla_{\theta} \log \pi_{\theta}(a|s) Q^{\pi}(s, a) da = \\ &= \mathbb{E}_a \nabla_{\theta} \log \pi_{\theta}(a|s) Q^{\pi}(s, a) \end{aligned}$$

Первый способ вывода policy gradient

Второе слагаемое

$$\int_A \pi(a|s) \nabla_{\theta} Q^{\pi_{\theta}}(s, a) da = \mathbb{E}_a \nabla_{\theta} Q^{\pi_{\theta}}(s, a)$$

Итого

$$\nabla_{\theta} J(\theta) = \mathbb{E}_a [\nabla_{\theta} \log \pi_{\theta}(a|s) Q^{\pi}(s, a) + \nabla_{\theta} Q^{\pi_{\theta}}(s, a)]$$

Замечание

Мы смогли выразить градиент V -функции через градиент Q -функции, попробуем сделать наоборот.

Первый способ вывода policy gradient

Q через V

$$\begin{aligned}\nabla_{\theta} Q^{\pi_{\theta}}(s, a) &= \nabla_{\theta} r(s, a) + \nabla_{\theta} \gamma \mathbb{E}_{s'} V^{\pi_{\theta}}(s') = \\ &= \nabla_{\theta} \gamma \int_{\mathcal{S}} V^{\pi_{\theta}}(s') p(s' | s, a) ds' = \gamma \mathbb{E}_{s'} \nabla_{\theta} V^{\pi_{\theta}}(s')\end{aligned}$$

Подставляя одно в другое

$$\nabla_{\theta} J(\theta) = \mathbb{E}_a \mathbb{E}_{s'} [\nabla_{\theta} \log \pi_{\theta}(a | s) Q^{\pi}(s, a) + \gamma \nabla_{\theta} V^{\pi_{\theta}}(s')]$$

Первый способ вывода policy gradient

Замечание

Получили что-то в духе уравнения Беллмана. В правой части стоит матожидание по действию a , совершаемому из состояния s , и по следующему состоянию s' . Раскрутим рекурсию до бесконечности и получим:

Окончательный результат *Policy Gradient Theorem*

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi} \sum_{t \geq 0} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) Q^{\pi}(s_t, a_t),$$

где τ — траектория согласно политике.

Второй способ вывода policy gradient

Постановка задачи

Рассмотрим θ -параметризованное семейство политик $\pi(a|s, \theta)$ или $\pi_\theta(a|s)$, которое порождает траектории τ с вероятностями $p_\theta(\tau)$. Тогда будем максимизировать следующую величину:

$$V^\pi := J(\theta) := \mathbb{E}_{\tau \sim p_\theta(\tau)} \left[\sum_{t \geq 0} \gamma^t r_t \right] \rightarrow \max_{\theta}.$$

Обозначим награду на траектории τ за $r(\tau) := \sum_{t \geq 0} \gamma^t r_t$.

Второй способ вывода policy gradient

Вычисляем градиент

По определению матожидания:

$$J(\theta) := \mathbb{E}_{\tau \sim p_{\theta}(\tau)}[r(\tau)] = \int p_{\theta}(\tau) r(\tau) d\tau,$$

тогда

$$\nabla_{\theta} J(\theta) = \int \nabla_{\theta} p_{\theta}(\tau) r(\tau) d\tau.$$

Трюк производной логарифма

$$\nabla_{\theta} p_{\theta}(\tau) = p_{\theta}(\tau) \frac{\nabla_{\theta} p_{\theta}(\tau)}{p_{\theta}(\tau)} = p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau)$$

Второй способ вывода policy gradient

Вычисляем градиент

Используем трюк

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \int \nabla_{\theta} p_{\theta}(\tau) r(\tau) d\tau = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) \cdot r(\tau) d\tau = \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log p_{\theta}(\tau) r(\tau)].\end{aligned}$$

Вероятность траектории

$$p_{\theta}(\tau) = p_{\theta}(s_0, a_0, \dots, s_T, a_T, s_{T+1}) = p(s_0) \cdot \prod_{t=0}^T (\pi_{\theta}(a_t | s_t) p(s_{t+1} | a_t, s_t))$$

Второй способ вывода policy gradient

Логарифм вероятности траектории

$$\log p_{\theta}(\tau) = \log p(s_0) + \sum_{t=0}^T [\log \pi_{\theta}(a_t | s_t) + \log p(s_{t+1} | a_t, s_t)]$$

Градиент логарифма вероятности траектории

$$\nabla_{\theta} \log p_{\theta}(\tau) = \nabla_{\theta} \sum_{t=0}^T [\log \pi_{\theta}(a_t | s_t)] + \underbrace{\nabla_{\theta} \sum_{t=0}^T [\log p(s_{t+1} | a_t, s_t)]}_{\text{в среднем, нулевой вектор по трюку производной логарифма}}$$

Второй способ вывода policy gradient

Конечная формула второго способа

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left(\sum_{t \geq 0} \nabla_{\theta} [\log \pi_{\theta}(a_t | s_t) r(\tau)] \right)$$

Вывод

Видим, что мы более простым способом получили очень похожую формулу, но с суммарной наградой за игры вместо Q -функции из первого способа. Математически эти формы будут эквивалентны, то есть равны, как интегралы, но их Монте-Карло оценки могут начать вести себя совершенно по-разному.

Эквивалентность

Потихоньку идем к эквивалентности

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\left(\sum_{t \geq 0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left(\sum_{\hat{t} \geq 0} \gamma^{\hat{t}} r_{\hat{t}} \right) \right] = \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \sum_{t \geq 0} \sum_{\hat{t} \geq 0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \gamma^{\hat{t}} r_{\hat{t}},\end{aligned}$$

выпишем одно из слагаемых этой двойной суммы:

$$j_t := \sum_{\hat{t} \geq 0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \gamma^{\hat{t}} r_{\hat{t}}.$$

Эквивалентность

Замечание

Видим, что на слагаемое $j_t = \sum_{\hat{t} \geq 0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \gamma^{\hat{t}} r_{\hat{t}}$, отвечающее за градиент решения выбрать a_t в момент времени t , влияют не только награды после принятия этого решения ($\hat{t} \geq t$), но и награды из прошлого ($\hat{t} < t$), то есть некая величина, на которую наше только что принятое решение никак не могло повлиять.

Пример почему это плохо

До момента времени $t_{\text{near end}} \approx T$ агент мог выполнять максимально хорошие действия и $\sum_{\hat{t} < t_{\text{near end}}} \gamma^{\hat{t}} r_{\hat{t}}$ очень большая величина, а вот решение на $t_{\text{near end}}$ может быть просто ужасным и после него награда всегда минимально возможная. Но на градиенте это плохое решение мы не увидим, потому что в сумме награда всё ещё большая.

Эквивалентность

Теорема

Для произвольного распределения $\pi_\theta(a)$ верно:

$$\mathbb{E}_{a \sim \pi_\theta(a)} \nabla_\theta \log \pi_\theta(a) = 0.$$

Доказательство

$$\begin{aligned} \mathbb{E}_{a \sim \pi_\theta(a)} \nabla_\theta \log \pi_\theta(a) &= \mathbb{E}_{a \sim \pi_\theta(a)} \frac{\nabla_\theta \pi_\theta(a)}{\pi_\theta(a)} = \\ &= \int_A \nabla_\theta \pi_\theta(a) da = \nabla_\theta \int_A \pi_\theta(a) da = \nabla_\theta 1 = 0 \end{aligned}$$

Эквивалентность

Теорема — Принцип причинности

При $\hat{t} < t$:

$$\mathbb{E}_{\tau \sim \pi_\theta} \nabla_\theta \log \pi_\theta(a_t | s_t) \gamma^{\hat{t}} r_{\hat{t}} = 0.$$

Доказательство

По теореме выше:

$$\begin{aligned} \mathbb{E}_{\tau \sim \pi_\theta} \nabla_\theta [\log \pi_\theta(a_t | s_t) \gamma^{\hat{t}} r_{\hat{t}}] &= \\ \mathbb{E}_{a_1, s_1, \dots, s_{\hat{t}}, a_{\hat{t}}} \mathbb{E}_{s_{\hat{t}+1}, a_{\hat{t}+1}, \dots, s_t, a_t, \dots} [\nabla_\theta \log \pi_\theta(a_t | s_t) \gamma^{\hat{t}} r_{\hat{t}}] &= \\ = \mathbb{E}_{a_1, s_1, \dots, s_{\hat{t}}, a_{\hat{t}}} [\gamma^{\hat{t}} r_{\hat{t}} \cdot \mathbb{E}_{s_{\hat{t}+1}, a_{\hat{t}+1}, \dots, s_t, a_t, \dots} \nabla_\theta \log \pi_\theta(a_t | s_t)] &= 0. \end{aligned}$$

Эквивалентность

Вывод

В формуле полученной вторым способом из суммы можно убрать все слагаемые с $\hat{t} < t$, поскольку они после взятия математического ожидания обратятся в нуль. Плюс ко всему, вычеркивание этих слагаемых уменьшит дисперсию при оценке градиента, полученного вторым способом, по методу Монте-Карло.

Замечание

Дисконтирование в среде идет с самого начала, поэтому дисконтирующий фактор при слагаемых $\hat{t} \geq t$ своей степени не поменяет, а значит его можно переписать как:

$$\sum_{\hat{t} \geq t} \gamma^{\hat{t}} r_{\hat{t}} = \gamma^t \sum_{\hat{t} \geq t} \gamma^{\hat{t}-t} r_{\hat{t}} =: \gamma^t r_t(\tau).$$

Эквивалентность

Приближаясь к эквивалентной форме

На данный момент имеем:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \sum_{t \geq 0} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) r_t(\tau).$$

Неформальное доказательство эквивалентности

$r_t(\tau)$ — очень похож на Q -функцию, поскольку является её несмещенной Монте-Карло оценкой, а в формуле выше всё равно берется матожидание, поэтому формулы из первого способа и из второго — одно и то же.

Эквивалентность

Теорема об эквивалентности

Следующие формулы эквивалентны:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi} \sum_{t \geq 0} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) Q^{\pi}(s_t, a_t),$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \sum_{t \geq 0} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) r_t(\tau).$$

Доказательство

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \sum_{t \geq 0} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) r_t(\tau) = \\&= \sum_{t \geq 0} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) r_t(\tau) = \\&= \sum_{t \geq 0} \mathbb{E}_{a_0, s_1, \dots, s_t, a_t} \mathbb{E}_{s_{t+1}, a_{t+1}, \dots} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) r_t(\tau) = \\&= \sum_{t \geq 0} \mathbb{E}_{a_0, s_1, \dots, s_t, a_t} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \mathbb{E}_{s_{t+1}, a_{t+1}, \dots} r_t(\tau) = \\&= \sum_{t \geq 0} \mathbb{E}_{a_0, s_1, \dots, s_t, a_t} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) Q^{\pi}(s_t, a_t) = \\&= \mathbb{E}_{\tau \sim \pi} \sum_{t \geq 0} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) Q^{\pi}(s_t, a_t)\end{aligned}$$

Физический смысл

Заметим

Оказывается, градиент нашего функционала имеет вид градиента взвешенных логарифмов правдоподобий. Чтобы ещё лучше увидеть это, рассмотрим **суррогатную функцию** — другой функционал, который будет иметь в точке текущих значений параметров стратегии π такой же градиент, как и $J(\theta)$:

Определение суррогатной функции

$$\mathcal{L}_{\tilde{\pi}}(\theta) := \mathbb{E}_{\tau \sim \tilde{\pi}} \sum_{t \geq 0} \gamma^t \log \pi_{\theta}(a|s) Q^{\tilde{\pi}}(s, a)$$

Физический смысл

Что дальше?

Получили суррогатную функцию от двух стратегий: стратегии π_θ , которую мы оптимизируем, и ещё одной стратегии $\tilde{\pi}$. Давайте рассмотрим эту суррогатную функцию в точке θ такой, что эти две стратегии совпадают: $\pi_\theta = \tilde{\pi}$, и посмотрим на градиент при изменении θ , только одной из них. Буквально мы «заморозим» оценочную Q-функцию, и «заморозим» распределение, из которого приходят пары (s, a) .

Утверждение

$$\nabla_\theta \mathcal{L}_{\tilde{\pi}}(\theta) |_{\tilde{\pi}=\pi_\theta} = \nabla_\theta J(\theta)$$

Физический смысл

Доказательство

Поскольку мат.ожидание по траекториям не зависит в этой суррогатной функции от θ , то градиент просто можно пронести внутрь:

$$\nabla_{\theta} \mathcal{L}_{\tilde{\pi}}(\theta)|_{\tilde{\pi}=\pi_{\theta}} = \mathbb{E}_{\tau \sim \tilde{\pi}} \sum_{t \geq 0} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a|s)|_{\tilde{\pi}=\pi_{\theta}} Q^{\tilde{\pi}}(s, a).$$

В точке θ такой, что $\pi_{\theta} = \tilde{\pi}$ верно, что $p(\tau|\tilde{\pi}) \equiv p(\tau|\pi_{\theta})$ и $Q^{\tilde{\pi}}(s, a) = Q^{\pi}(s, a)$; следовательно, значение градиента в этой точке совпадает со значением формулы для $\nabla_{\theta} J(\theta)$.

Вывод

Значит, направление максимизации $J(\theta)$ в текущей точке θ просто совпадает с направлением максимизации этой суррогатной функции! Таким образом, можно считать, что в текущей точке мы на самом деле «как бы» максимизируем, а это уже в чистом виде логарифм правдоподобия каких-то пар (s, a) , для каждой из которых дополнительно выдан «вес» в виде значения $Q^\pi(s, a)$.

Физический смысл

Например

Если в машинном обучении в задачах регрессии и классификации мы для данной выборки (x, y) максимизировали правдоподобие:

$$\sum_{(x,y)} \log p(y|x, \theta) \rightarrow \max_{\theta},$$

то теперь в RL, когда выборки нет, мы действуем по-другому: мы сэмплируем сами себе входные данные s и примеры выходных данных a , выдаём каждой паре какой-то «кредит доверия», некую скалярную оценку хорошести, выраженную в виде $Q^{\pi}(s, a)$, и идём в направлении максимизации:

$$\sum_{(x,y)} \log p(y|x, \theta) Q^{\pi}(s, a) \rightarrow \max_{\theta}.$$

Способ получения

Monte-Carlo

Воспользуемся первой формулой для подсчета *Policy Gradient*:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi} \sum_{t \geq 0} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) Q^{\pi}(s_t, a_t),$$

в которой заменим всё неизвестное на оценку по Монте-Карло.

Способ получения

Что будем заменять?

- $\mathbb{E}_{\tau \sim \pi} \rightarrow$ сыграем несколько полных игр при помощи текущей стратегии π — алгоритм будет *on-policy*.
- $Q^\pi(s_t, a_t) \rightarrow r(\tau)$ — можно сказать, что мы воспользовались вторым способом подсчета *Policy Gradient* с Монте-Карло оценкой $\mathbb{E}_{\tau \sim \pi}$, что не удивляет, ведь подходы, как мы уже показали, эквиваленты.

Итоговый алгоритм

Reinforce

Гиперпараметры: N — количество игр, $\pi(a|s, \theta)$ — стратегия с параметрами θ , SGD -оптимизатор.

0. Произвольно инициализируем θ .

На очередном шаге t

1. Играем N игр $\tau_1, \tau_2, \dots, \tau_N \sim \pi$.
2. Для каждого t в каждом τ_i считаем $r_t(\tau) := \sum_{\hat{t} \geq t} \gamma^{\hat{t}-t} r_{\hat{t}}$.
3. Считаем оценку градиента:

$$\nabla_{\theta} J(\pi) := \frac{1}{N} \sum_{\tau} \sum_{t \geq 0} \gamma^t \nabla_{\theta} \log \pi(a_t | s_t, \theta) r_t(\tau).$$

4. Делаем шаг градиентного подъёма по θ , используя $\nabla_{\theta} J(\pi)$.

Итоговый алгоритм

Недостатки

1. Для одного шага градиентного подъёма нам необходимо играть несколько игр **до конца** при помощи текущей стратегии.
2. Колоссальная дисперсия нашей оценки градиента — на практике дожидаться каких-то результатов от такого алгоритма в сколько-то сложных задачах не получится.

Два типа стохастики

Мотивация

До этого мы часто работали с функционалами вида $\mathbb{E}_{\tau \sim \pi} \sum_{t \geq 0} \gamma^t f(s_t, a_t)$, где f — какая-то функция от пар состояние-действие.

Стохастика

В MDP есть два вида стохастики:

- 1) **внешняя** связанная со случайностью в самой среде и неподконтрольная агенту; она заложена в функции переходов $p(s'|s, a)$;
- 2) **внутренняя**, связанная со случайностью в стратегии самого агента; она заложена в $\pi(a|s)$. Это стохастика нам подконтрольна при обучении.

Два типа стохастики

Мотивация

Матожидание $\mathbb{E}_{\tau \sim \pi}$ плохо тем, что мат.ожидания по внешней и внутренней стохастике чередуются. При этом во время обучения из внешней стохастики мы можем только получать сэмплы, поэтому было бы здорово переписать наш функционал как-то так, чтобы он имел вид мат.ожидания по всей внешней стохастике.

Два типа стохастики

Утверждение

Состояния, которые встречается агент со стратегией π , приходят из некоторой стационарной марковской цепи.

Доказательство

Выпишем вероятность оказаться на очередном шаге в состоянии s' , если мы используем стратегию π :

$$p(s'|s) = \int_A \pi(a|s)p(s'|s, a)da.$$

Эта вероятность не зависит от времени и от истории, следовательно, цепочка состояний образует марковскую цепь.

Два типа стохастики

Допустим, начальное состояние s_0 фиксировано. Обозначим вероятность оказаться в состоянии s в момент времени t при использовании стратегии π как $p(s_t = s | \pi)$.

Определение

Для данного MDP и политики π **state visitation frequency** называется:

$$\mu_\pi(s) := \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t \geq 0}^T p(s_t = s | \pi).$$

Введём ещё один, «дисконтированный счётчик посещения состояний» для стратегии взаимодействия π . При дисконтировании отпадают проблемы с нормировкой.

Два типа стохастики

Определение

Для данного MDP и политики π **discounted state visitation distribution** называется

$$d_{\pi}(s) := (1 - \gamma) \sum_{t \geq 0} \gamma^t p(s_t = s | \pi).$$

Утверждение

State visitation distribution есть распределение на множестве состояний, то есть:

$$\int_S d_{\pi}(s) ds = 1.$$

Два типа стохастики

Доказательство

$$\begin{aligned}\int_S d\pi(s) ds &= \int_S (1 - \gamma) \sum_{t \geq 0} \gamma^t p(s_t = s | \pi) ds = \\ &= (1 - \gamma) \sum_{t \geq 0} \gamma^t \int_S p(s_t = s | \pi) ds = 1\end{aligned}$$

Два типа стохастики

Теорема

Для произвольной функции $f(s, a)$:

$$\mathbb{E}_{\tau \sim \pi} \sum_{t \geq 0} \gamma^t f(s_t, a_t) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_\pi(s)} \mathbb{E}_{a \sim \pi(a|s)} f(s, a).$$

Начало доказательства

$$\begin{aligned} \mathbb{E}_{\tau \sim \pi} \sum_{t \geq 0} \gamma^t f(s_t, a_t) &= \sum_{t \geq 0} \gamma^t \mathbb{E}_{\tau \sim \pi} f(s_t, a_t) = \\ &= \sum_{t \geq 0} \gamma^t \int \int_S \int_A p(s_t = s, a_t = a | \pi) f(s, a) da ds = \end{aligned}$$

Два типа стохастики

Продолжение доказательства

$$\begin{aligned} &= \sum_{t \geq 0} \gamma^t \int_S \int_A p(s_t = s | \pi) \pi(a | s) f(s, a) da ds = \\ &= \sum_{t \geq 0} \gamma^t \int_S p(s_t = s | \pi) \mathbb{E}_{\pi(a|s)} f(s, a) ds = \\ &= \int_S \sum_{t \geq 0} \gamma^t p(s_t = s | \pi) \mathbb{E}_{\pi(a|s)} f(s, a) ds = \\ &= \int_S \frac{d\pi(s)}{1 - \gamma} \mathbb{E}_{\pi(a|s)} f(s, a) ds = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_\pi(s)} \mathbb{E}_{\pi(a|s)} f(s, a) \end{aligned}$$

Два типа стохастики

Пример

$$J(\pi) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_\pi(s)} \mathbb{E}_{\pi(a|s)} r(s, a)$$

Пример

$$\nabla_\theta J(\pi) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_\pi(s)} \mathbb{E}_{\pi(a|s)} \nabla_\theta \log \pi_\theta(a|s) Q^\pi(s, a)$$

Уменьшаем дисперсию Reinforce. Baseline.

Мотивация

При стохастической оптимизации ключевым фактором является дисперсия оценки градиента. Когда мы заменяем мат.ожидания на Монте-Карло оценки, дисперсия увеличивается. Понятно, что замена Q-функции — выинтегрированных будущих наград — на её Монте-Карло оценку в REINFORCE повышало дисперсию. Однако, в текущем виде основной источник дисперсии заключается в другом.

Уменьшаем дисперсию Reinforce. Baseline.

Причина большой дисперсии

Градиент логарифма правдоподобия в среднем равен нулю. Это значит, что если для данного s мы выдаём некоторое распределение $\pi(a|s)$, для увеличения вероятностей в одной области A нужно данный вес θ_i параметризации увеличивать, а в другой области — уменьшать. В среднем «магнитуда изменения» равна нулю. Но у нас в Монте-Карло оценке только $a \sim \pi(a|s)$, и для него направление изменения домножится на кредит: на нашу оценку $Q^\pi(s, a)$. Если эта оценка в одной области 100, а в другой 1000 — дисперсия получаемых значений $\nabla_\theta \log \pi_\theta(a|s) Q^\pi(s, a)$ становится колоссальной.

Вывод

Кредит надо центрировать «умным» нулем!

Уменьшаем дисперсию Reinforce. Baseline.

Утверждение

Для произвольной функции $b(s) : S \rightarrow \mathbb{R}$, называемой бэйзлайном, верно:

$$\nabla_{\theta} J(\pi) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{\pi}(s)} \mathbb{E}_{\pi(a|s)} \nabla_{\theta} \log \pi_{\theta}(a|s) (Q^{\pi}(s, a) - b(s)).$$

Доказательство

Добавленное слагаемое есть ноль в силу формулы теоремы о среднем градиента логарифма.

Уменьшаем дисперсию Reinforce. Baseline.

Замечание

Это верно для произвольной функции от состояний и становится неверно, если вдруг бэйзлайн $b(s)$ начинает зависеть от a . Мы вольны выбрать бэйзлайн произвольно; он не меняет среднего значения оценок градиента, но изменяет дисперсию.

Уменьшаем дисперсию Reinforce. Baseline.

Теорема

Бэйзлайном, максимально снижающим дисперсию Монте-Карло оценок формулы градиентов, является

$$b^*(s) := \frac{\mathbb{E}_a \|\nabla_\theta \log \pi(a, s)\|_2^2 Q^\pi(s, a)}{\mathbb{E}_a \|\nabla_\theta \log \pi(a, s)\|_2^2}$$

Проблема

Практическая ценность результата невысока. Знать норму градиента для всех действий a вычислительно будет труднозатратно даже в дискретных пространствах действий.

Уменьшаем дисперсию Reinforce. Baseline.

На практике

$$\begin{aligned} b^*(s) &:= \frac{\mathbb{E}_a \|\nabla_\theta \log \pi(a, s)\|_2^2 Q^\pi(s, a)}{\mathbb{E}_a \|\nabla_\theta \log \pi(a, s)\|_2^2} = \\ &= [\|\nabla_\theta \log \pi(a, s)\|_2^2 \approx \text{const}(a)] \approx \mathbb{E}_a Q^\pi(s, a) = V^\pi(s) \end{aligned}$$

Итого

$$\nabla_\theta J(\pi) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_\pi(s)} \mathbb{E}_{\pi(a|s)} \nabla_\theta \log \pi_\theta(a|s) A^\pi(s, a)$$

Уменьшаем дисперсию Reinforce. Baseline.

Итого

$$\nabla_{\theta} J(\pi) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\pi}(s)} \mathbb{E}_{\pi(a|s)} \nabla_{\theta} \log \pi_{\theta}(a|s) A^{\pi}(s, a)$$

Определение

Для данного MDP **Advantage-функцией** политики π называется

$$A^{\pi}(s, a) := Q^{\pi}(s, a) - V^{\pi}(s).$$

Схемы «актор-критик»

- Хотим оптимизировать параметры стратегии при помощи формулы градиента, не доигрывая эпизоды до конца.
- Введём вторую сетку, которая будет «оценивать» наши собственные решения — критика (critic). Нейросеть, моделирующую стратегию, соответственно будем называть актёром или актором (actor), и такие алгоритмы, в которых обучается как модель критика, так и модель актора, называются Actor-Critic.

Схемы «актор-критик»

- В качестве критика обычно учат именно V -функцию.
- Возможность не обучать сложную Q^* является одним из преимуществ подхода прямой оптимизации $J(\theta)$.
- Из соображений эффективности:

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s'} V^\pi(s') \approx r(s, a) + \gamma V_\phi(s'), \quad s' \sim p(s' | s, a).$$

Схемы «актор-критик». Bias-variance trade-off

Собираемся вместо честного advantage подставить некоторую его оценку (advantage estimator) и провести таким образом credit assingment:

$$\nabla_{\theta} J(\pi) \approx \frac{1}{1 - \gamma} \mathbb{E}_{d_{\pi}(s)} \mathbb{E}_{a \sim \pi(a|s)} \nabla_{\theta} \log \pi_{\theta}(a | s) \underbrace{\Psi(s, a)}_{\approx A^{\pi}(s, a)}.$$

В контексте policy gradient, речь напрямую идёт о дисперсии и смещении оценок градиента.

$\Psi(s, a)$	Дисперсия	Смещение
$R_t - V_{\phi}(s)$	высокая	нет
$r(s, a) + \gamma V_{\phi}(s') - V_{\phi}(s)$	низкая	большое

Схемы «актор-критик». Построение оценки Q-функции

$$Q^\pi(s, a) \approx \sum_{t=0}^{N-1} \gamma^t r^{(t)} + \gamma^N V_\phi(s^{(N)}).$$

Для credit assingment-а N -шаговой оценкой Advantage, или N -шаговой временной разностью:

$$\Psi_{(N)}(s, a) := \sum_{t=0}^{N-1} \gamma^t r^{(t)} + \gamma^N V_\phi(s^{(N)}) - V_\phi(s).$$

С ростом N дисперсия такой оценки увеличивается: всё больший фрагмент траектории мы оцениваем по Монте-Карло, нам становятся нужны сэмплы $a_{t+1} \sim \pi(a_{t+1} \mid s_{t+1})$, $s_{t+2} \sim \pi(s_{t+2} \mid s_{t+1}, a_{t+1})$, \dots , $s_{t+N} \sim \pi(s_{t+N} \mid s_{t+N-1}, a_{t+N-1})$.

Схемы «актор-критик». Построение оценки Q-функции

Определение

Для пар s_t, a_t из роллаута $s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_N$ длины N будем называть оценкой максимальной длины (max trace estimation) оценку с максимальным заглядыванием в будущее: для Q-функции

$$y^{\text{MaxTrace}}(s_t, a_t) := \sum_{\hat{t}=t}^{N-1} \gamma^{\hat{t}-t} r_{\hat{t}} + \gamma^{N-t} V^{\pi}(s_N), \quad (8)$$

для Advantage функции, соответственно:

$$\Psi^{\text{MaxTrace}}(s_t, a_t) := y^{\text{MaxTrace}}(s_t, a_t) - V^{\pi}(s_t). \quad (9)$$

Generalized Advantage Estimation (GAE)

Решение дилеммы bias-variance trade-off подсказывает теория $TD(\lambda)$ оценки. Нужно применить формулу $TD(\lambda)$ и просто заансамблировать N -шаговые оценки разной длины:

Определение

GAE-оценкой Advantage-функции называется ансамбль многошаговых оценок, где оценка длины N берётся с весом λ^{N-1} , где $\lambda \in (0, 1)$ — гиперпараметр:

$$\Psi_{\text{GAE}}(s, a) := (1 - \lambda) \sum_{N \geq 0} \lambda^{N-1} \Psi_{(N)}(s, a).$$

Как мы помним, при $\lambda \rightarrow 0$ такая GAE-оценка соответствует одношаговой оценке; при $\lambda = 1$ GAE-оценка соответствует Монте-Карло оценке Q-функции.

Generalized Advantage Estimation (GAE)

В текущем виде в формуле суммируются все N -шаговые оценки вплоть до конца эпизода. В реальности собранные роллауты могут прерваться в середине эпизода: допустим, для данной пары s, a через M шагов роллаут «обрывается». Тогда на практике используется чуть-чуть другим определением GAE-оценки: если мы знаем $s^{(M)}$, но после этого эпизод ещё не доигран до конца, мы пользуемся формулой $TD(\lambda)$ и оставляем от суммы только «доступные» слагаемые:

$$\Psi_{\text{GAE}}(s, a) := \sum_{t \geq 0}^{M-1} \gamma^t \lambda^t \Psi_{(1)}(s^{(t)}, a^{(t)}). \quad (10)$$

Компромисс между смещением и разбросом

Дана траектория $s, r, s', r', s'', r'' \dots s^{(M)}$ по политике π и приближение $V^\pi(s)$

выполним оценку advantage (credit assignment) для пары s, a (хорошее ли решение принято было?)

Для актора:

$$\nabla := \rho(\theta) \nabla_{\theta} \log \pi_{\theta}(a | s) \underbrace{\Psi(s, a)}_{\substack{\text{оценка} \\ \text{advantage}}}$$

Для критика:

$$\underbrace{y_Q}_{\substack{\text{целевое значение} \\ \text{для регрессии}}} := \Psi(s, a) + V(s)$$

	$\Psi(s, a)$	Смещение	Разброс
Монте-Карло	$\Psi_{(\infty)}(s, a) := r + \gamma r' + \gamma^2 r'' + \dots - V(s)$	0	высокий
N -шагов	$\Psi_{(N)}(s, a) := r + \gamma r' + \dots + \gamma^N V(s^{(N)}) - V(s)$	промежуточное	промежуточный
1-шаг	$\Psi_{(1)}(s, a) := r + \gamma V(s') - V(s)$	высокое	низкий

Проблема: выбор N .

Generalized Advantage Estimation (GAE)

Утверждение

Формула (10) эквивалентна следующему ансамблю N -шаговых оценок:

$$\Psi_{\text{GAE}}(s, a) = (1 - \lambda) \sum_{N \geq 0} \lambda^{N-1} \Psi_{(N)}(s, a) + \lambda^{M-1} \Psi_{(M)}(s, a).$$

В такой «обрезанной» оценке $\lambda = 1$ соответствует оценке максимальной длины (9), а $\lambda = 0$ всё ещё даст одношаговую оценку.

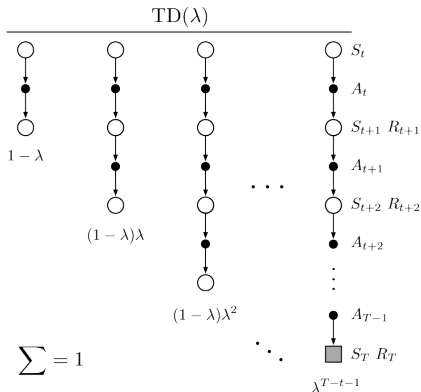
Generalized Advantage Estimation (GAE)

Шаг	Обновление	$\Psi_{(1)}(s, a)$	$\Psi_{(2)}(s, a)$	$\Psi_{(3)}(s, a)$...	$\Psi_{(N)}(s, a)$
0	$\Psi_{(1)}(s, a)$	1	0	0		0
1	$\Psi_{(1)}(s, a) + \gamma\lambda\Psi_{(1)}(s', a')$	$1 - \lambda$	λ	0		0
2	$\Psi_{(1)}(s, a) + \gamma\lambda\Psi_{(1)}(s', a') + (\gamma\lambda)^2\Psi_{(1)}(s'', a'')$	$1 - \lambda$	$(1 - \lambda)\lambda$	λ^2		0
\vdots						
N	$\sum_{t \geq 0}^N (\gamma\lambda)^t \Psi_{(1)}(s^{(t)}, a^{(t)})$	$1 - \lambda$	$(1 - \lambda)\lambda$	$(1 - \lambda)\lambda^2$		λ^N

Эквивалентные формы обновлений TD(λ)

$$\sum_{t=0}^{\infty} (\gamma\lambda)^t \Psi_{(1)}(s^{(t)}, a^{(t)}) = (1 - \lambda) \sum_{N=1}^{\infty} \lambda^{N-1} \Psi_{(N)}(s, a)$$

Generalized Advantage Estimation (GAE)



Что если для некоторой пары s, a нам известно будущее только на T шагов вперёд?

$$\Psi^{\text{GAE}}(s, a) := \sum_{t=0}^T (\gamma \lambda)^t \Psi_{(1)}(s^{(t)}, a^{(t)})$$

Используемое на практике уравнение:

$$\Psi^{\text{GAE}}(s_t, a_t) = \Psi_{(1)}(s_t, a_t) + \lambda \gamma (1 - \text{done}_{t+1}) \Psi^{\text{GAE}}(s_{t+1}, a_{t+1})$$

Generalized Advantage Estimation (GAE)

В коде формула (10) очень удобна для рекурсивного подсчёта оценки; также для практического алгоритма осталось учесть флаги done_t . Формулы подсчёта GAE-оценки для всех пар (s, a) из роллаута $s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_N$ приобретают такой вид:

$$\begin{aligned}\Psi_{\text{GAE}}(s_{N-1}, a_{N-1}) &:= \Psi_{(1)}(s_{N-1}, a_{N-1}) \\ \Psi_{\text{GAE}}(s_{N-2}, a_{N-2}) &:= \Psi_{(1)}(s_{N-2}, a_{N-2}) + \\ &\quad + \gamma \lambda (1 - \text{done}_{N-2}) \Psi_{\text{GAE}}(s_{N-1}, a_{N-1}) \\ &\vdots \\ \Psi_{\text{GAE}}(s_0, a_0) &:= \Psi_{(1)}(s_0, a_0) + \gamma \lambda (1 - \text{done}_0) \Psi_{\text{GAE}}(s_1, a_1)\end{aligned}$$

Заметим, что эти формулы очень похожи на расчёт кумулятивной награды за эпизод, где «наградой за шаг» выступает $\Psi_{(1)}(s, a)$.

Обучение критика

Воспользуемся идеей перехода к регрессии, которую мы обсуждали раньше в контексте DQN. Нам нужно просто решать методом простой итерации уравнение Беллмана:

$$V_{\phi_{k+1}}(s) \leftarrow \mathbb{E}_a [r + \gamma \mathbb{E}_{s'} V_{\phi_k}(s')] .$$

Воспользуемся преимуществами on-policy режима и поймём, что мы можем поступить точно также, как с оценкой Q-функции в формуле градиента: решать многошаговое уравнение Беллмана вместо одношагового. Например, можно выбрать любое N -шаговое уравнение и строить целевую переменную как

$$y := r + \gamma r' + \gamma^2 r'' + \dots + \gamma^N V_{\phi_k}(s^{(N)}) . \quad (11)$$

Обучение критика

Тогда если мы оцениваем Advantage как

$$\Psi(s, a) = y - V_{\phi}(s),$$

где y — некоторая оценка Q-функции, то y же является и таргетом для V-функции, и наоборот. Используя функцию потерь MSE с таким таргетом, мы как раз и учим среднее значение наших оценок Q-функции, то есть бэйзлайн.

Конечно же, мы можем использовать и GAE-оценку (10) Advantage, достаточно «убрать бэйзлайн»:

$$Q^{\pi}(s, a) = A^{\pi}(s, a) + V^{\pi}(s) \approx \Psi_{\text{GAE}}(s, a) + V_{\phi}(s).$$

Утверждение

Таргет $\Psi_{\text{GAE}}(s, a) + V_\phi(s)$ является несмещённой оценкой правой части «ансамбля» уравнений Беллмана:

$$V_\phi(s) = (1 - \lambda) \sum_{N>0} \lambda^{N-1} [\mathfrak{B}^N V_\phi](s),$$

где \mathfrak{B} — оператор Беллмана для V -функции.

Доказательство.

По определению, поскольку $\Psi_{(N)}(s, a) + V(s)$ является несмещённой оценкой правой части N -шагового уравнения Беллмана (т. е. несмещённой оценкой $[\mathfrak{B}^N V^\pi](s)$), а

$$(1 - \lambda) \sum_{N>0} \lambda^{N-1} (\Psi_{(N)}(s, a) + V(s)) = \Psi_{\text{GAE}}(s, a) + V(s).$$



Обучение критика

Делаем несколько шагов взаимодействия со средой, собирая таким образом роллаут некоторой длины N ; считаем для каждой пары s, a некоторую оценку Q-функции $y(s, a)$, например, оценку максимальной длины (8); оцениваем Advantage каждой пары как $\Psi(s, a) := y(s, a) - V_\phi(s)$; далее по Монте-Карло оцениваем градиент по параметрам стратегии

$$\nabla_\theta J(\pi) \approx \frac{1}{N} \sum_{s,a} \nabla_\theta \log \pi_\theta(a | s) \Psi(s, a)$$

и градиент для оптимизации критика (допустим, критик — Q-функция):

$$\text{Loss}^{\text{critic}}(\phi) = \frac{1}{N} \sum_{s,a} (y(s, a) - V_\phi(s))^2$$

Advantage Actor-Critic (A2C)

Гиперпараметры: M — количество параллельных сред, N — длина роллаутов, $V_\phi(s)$ — нейросеть с параметрами ϕ , $\pi_\theta(a | s)$ — нейросеть для стратегии с параметрами θ , SGD или другой оптимизатор первого порядка.

Инициализировать θ, ϕ

На каждом шаге:

- 1 в каждой параллельной среде собрать роллаут длины N , используя стратегию π_θ :

$$s_0, a_0, r_0, s_1, \dots, s_N$$

Advantage Actor-Critic (A2C)

- 2 для каждой пары s_t, a_t из каждого роллаута посчитать оценку Q-функции максимальной длины, игнорируя зависимость оценки от ϕ :

$$Q(s_t, a_t) := \sum_{\hat{t}=t}^{N-1} \gamma^{\hat{t}-t} r_{\hat{t}} + \gamma^{N-t} V_{\phi}(s_N)$$

- 3 вычислить лосс критика:

$$\text{Loss}^{\text{critic}}(\phi) := \frac{1}{MN} \sum_{s_t, a_t} (Q(s_t, a_t) - V_{\phi}(s_t))^2$$

- 4 делаем шаг градиентного спуска по ϕ , используя $\nabla_{\phi} \text{Loss}^{\text{critic}}(\phi)$

Advantage Actor-Critic (A2C)

5 вычислить градиент для актора:

$$\nabla_{\theta}^{\text{actor}} := \frac{1}{MN} \sum_{s_t, a_t} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) (Q(s_t, a_t) - V_{\phi}(s_t))$$

6 сделать шаг градиентного подъёма по θ , используя $\nabla_{\theta}^{\text{actor}}$

Off-policy оценка advantage

Данной траектории $s_0, r_0, s_1, r_1, s_2, r_2 \dots s_M$ по политике μ и приближению $V^\pi(s)$ требуется произвести оценку advantage (credit assignment) для пары состояние-действие s_0, a_0 в off-policy режиме: $\mu \neq \pi$.

Было бы замечательно воспользоваться GAE:

$$\sum_{t \geq 0} (\gamma \lambda)^t \Psi_{(1)}(s_t, a_t),$$

но $\Psi_{(1)}(s_t, a_t)$ зависит от случайных величин: $a_0, s_0, a_1, s_2, \dots s_{t+1}$.

Внимание!

Если $\pi(a_0 | s_0) = 0$, то ничего не получится.



Воспользуемся коррекцией с помощью выборки по значимости!

$$\begin{aligned}\Psi &= \sum_{t \geq 0} (\gamma \lambda)^t \left(\prod_{\hat{t}=1}^{\hat{t}=t} \frac{\pi(a_{\hat{t}} | s_{\hat{t}}) \cancel{p(s_{\hat{t}+1} | s_{\hat{t}}, a_{\hat{t}})}}{\mu(a_{\hat{t}} | s_{\hat{t}}) \cancel{p(s_{\hat{t}+1} | s_{\hat{t}}, a_{\hat{t}})}} \right) \Psi_{(1)}(s_t, a_t) = \\ &= \sum_{t \geq 0} (\gamma \lambda)^t \left(\prod_{\hat{t}=1}^{\hat{t}=t} \frac{\pi(a_{\hat{t}} | s_{\hat{t}})}{\mu(a_{\hat{t}} | s_{\hat{t}})} \right) \Psi_{(1)}(s_t, a_t), \quad \prod_{\hat{t}=1}^0 \equiv 1.\end{aligned}$$

Непрактично: очень высокая дисперсия!

- «Затухающий» след: $\mu(a|s) \gg \pi(a|s)$:
 - типичная ситуация μ делает примитивные случайные действия, которые π редко совершает. Не лечится.
- «Взрывающийся» след: $\mu(a|s) \ll \pi(a|s)$:
 - μ выбранное действие с малой $\mu(a|s)$, но *вероятное* для π .

Причина большой дисперсии.

Присвоение ценности: общий вид

Давайте перепишем ценность следующим образом:

$$\Psi = \sum_{t \geq 0} \gamma^t \left(\prod_{i=1}^{i=t} c_i \right) \Psi_{(1)}(s_t, a_t), \quad \prod_{i=1}^0 \equiv 1,$$

где c_i коэффициенты «отжига следа»:

Название оценки	Коэффициенты c_i	Возникающая проблема
GAE	λ	только on-policy
Одношаговая	0	большое смещение
Выборка по значимости	$\lambda \frac{\pi(a_i s_i)}{\mu(a_i s_i)}$	легко «взрывается»

Retrace: основная теорема

$$\Psi = \sum_{t \geq 0} \gamma^t \left(\prod_{i=1}^{i=t} c_i \right) \Psi_{(1)}(s_t, a_t), \quad \prod_{i=1}^0 \equiv 1.$$

Теорема о Retrace

В режиме on-policy возможен выбор произвольного коэффициента $c_i \in [0, 1]$, в off-policy режиме можно выбрать произвольный коэффициент в следующем интервале

$$c_i \in \left[0, \frac{\pi(a_i | s_i)}{\mu(a_i | s_i)} \right].$$

- «затухающий» след: ничего не поделаешь;
- «взрывающийся» след: если вес из выборки по значимости больше 1, то применяем клиппинг!

Retrace: финальный результат

$$\Psi = \sum_{t \geq 0} \gamma^t \left(\prod_{i=1}^{i=t} c_i \right) \Psi_{(1)}(s_t, a_t), \quad \prod_{i=1}^0 \equiv 1,$$

где

$$c_i := \lambda \min \left(1, \frac{\pi(a_i | s_i)}{\mu(a_i | s_i)} \right).$$

Используется в:

- off-policy RL алгоритмах для теоретически корректных многошаговых целевых значений;
 - ($\lambda = 1$, потому что оно быстро затухает).
- дистрибутивных on-policy RL системах, где данные о градиенте от некоторых серверов могут задерживаться на несколько итераций обновления.

Источники I

- [1] S. Ivanov, “Reinforcement learning textbook,” arXiv preprint arXiv:2201.09746, 2022.
- [2] D. Bertsekas, Reinforcement learning and optimal control. Athena Scientific, 2019.
- [3] V. Goyal and J. Grand-Clement, “A first-order approach to accelerated value iteration,” Operations Research, vol. 71, no. 2, pp. 517–535, 2023.
- [4] M. L. Puterman, Markov decision processes: discrete stochastic dynamic programming. John Wiley & Sons, 2014.

Источники II

- [5] T. Liu, R. Zhou, D. Kalathil, P. Kumar, and C. Tian, “Policy optimization for constrained mdps with provable fast global convergence,” [arXiv preprint arXiv:2111.00552](#), 2021.
- [6] J. Grand-Clément, “From convex optimization to mdps: A review of first-order, second-order and quasi-newton methods for mdps,” [arXiv preprint arXiv:2104.10677](#), 2021.
- [7] J. N. Tsitsiklis, “Asynchronous stochastic approximation and q-learning,” [Machine learning](#), vol. 16, pp. 185–202, 1994.
- [8] H. R. Maei, C. Szepesvári, S. Bhatnagar, and R. S. Sutton, “Toward off-policy learning control with function approximation.,” in [ICML](#), vol. 10, pp. 719–726, 2010.