

Лекция 9

Обучение с подкреплением

Никита Юдин, iudin.ne@phystech.edu

Московский физико-технический институт
Физтех-школа прикладной математики и информатики

17 апреля 2024



Deep Deterministic Policy Gradient (DDPG)

Схема DQN имела принципиальный недостаток: мы не могли работать с непрерывными пространствами действий в силу необходимости постоянно считать операторы максимума и аргмаксимума

$$\max_a Q_\theta(s, a)$$

как для жадного выбора действия, так и для построения таргета в задаче регрессии.

Единственная возможная архитектура модели — приём действий на вход вместе с состояниями, тогда поиск аргмаксимума проводить, в целом, можно, но дорого: инициализируем a^0 случайно, и устраиваем градиентный подъём по входу в модель:

$$a^{k+1} = a^k + \alpha \nabla_a Q_\theta(s, a)|_{a=a^k}.$$

Итак, пусть $\pi_\omega(s)$ принимает на вход состояние s и выдаёт аргмаксимум текущей аппроксимации Q -функции, то есть будем добиваться $\pi_\omega(s) \approx \underset{a}{\operatorname{argmax}} Q_\theta(s, a)$. Понятно, как обучать такую сеть:

$$Q_\theta(s, \pi_\omega(s)) \rightarrow \max_\omega.$$

Весь алгоритм DQN оставляем неизменным с единственной модификацией, что на каждом батче также нужно сделать шаг оптимизации ω . При этом каждый раз, когда в схеме необходимо считать максимум или аргмаксимум Q_θ , используется $\pi_\omega(s)$. В стандартном алгоритме DQN нам было необходимо считать $\max_{a'} Q_{\theta^-}(s', a')$, и в дефолтной версии алгоритма использовалась таргет-сеть. Технически это означает, что для таргет-сети $Q_{\theta^-}(s', a')$ нам тоже нужно знать аргмаксимум, поэтому можно хранить старую версию вспомогательной функции $\pi_{\omega^-}(s)$.

Итого мы получили, что для жадного выбора действия используется $\pi_\omega(s)$ (отсюда такое обозначение этой «вспомогательной» функции — это фактически стратегия); а таргет для перехода $\mathbb{T} := (s, a, r, s')$ вычисляется по формуле

$$y(\mathbb{T}) := r + \gamma Q_{\theta^-}(s', \underset{a'}{\operatorname{argmax}} Q_{\theta^-}(s', a')) \approx r + \gamma Q_{\theta^-}(s', \pi_{\omega^-}(s)).$$

Такой алгоритм называется Deep Deterministic Policy Gradient (DDPG), и название может сбить с толку: а причём здесь policy gradient?

Теорема 1 — Deterministic Policy Gradient: В непрерывных пространствах действий в предположении дифференцируемости Q-функций по действиям:

$$\nabla_{\theta} J(\pi_{\theta}) = \frac{1}{1-\gamma} \mathbb{E}_{d \sim \pi_{\theta}(s)} (\nabla_{\theta} \pi_{\theta}(s))^{\top} \nabla_a Q^{\pi}(s, a)|_{a=\pi_{\theta}(s)}. \quad (1)$$

Доказательство.

$$\nabla_{\theta} V^{\pi_{\theta}}(s) = \{VQ \text{ уравнение}\} = \nabla_{\theta} \mathbb{E}_{a \sim \pi_{\theta}(s)} Q^{\pi_{\theta}}(s, a) = \nabla_{\theta} Q^{\pi_{\theta}}(s, \pi_{\theta}(s)) = (*).$$

Заметим, что в последнем выражении при малом изменении θ поменяется не только $\pi_{\theta}(s)$, но и сама оценочная функция $Q^{\pi_{\theta}}$. Считая, что якобиан функции $\mathbb{R}^n \rightarrow \mathbb{R}^m$ имеет размерность $n \times m$, и обозначая размерность действий как A , а размерность параметров θ буквой d , получаем следующие размерности матриц и векторов:

$$\nabla_{\theta} Q^{\pi_{\theta}}(s, a) \in \mathbb{R}^{d \times 1}, \quad \nabla_a Q^{\pi_{\theta}}(s, a) \in \mathbb{R}^{A \times 1}, \quad \nabla_{\theta} \pi_{\theta}(s) \in \mathbb{R}^{A \times d}.$$

Тогда продолжение вычисления выглядит так:

$$(*) = (\nabla_{\theta} \pi_{\theta}(s))^{\top} \nabla_a Q^{\pi}(s, a)|_{a=\pi_{\theta}(s)} + \nabla_{\theta} Q^{\pi_{\theta}}(s, a)|_{a=\pi_{\theta}(s)},$$

где последнее слагаемое — градиент $Q^{\pi_{\theta}}$ при фиксированном a по параметрам стратегии π_{θ} , которую он оценивает. Отдельно это слагаемое имеет вид:

$$\nabla_{\theta} Q^{\pi_{\theta}}(s, a) = \{QV \text{ уравнение} \} = \gamma \mathbb{E}_{s'} \nabla_{\theta} V^{\pi_{\theta}}(s').$$

Получаем рекурсивную формулу, аккуратно собирая которую, получим:

$$\nabla_{\theta} V^{\pi_{\theta}}(s) = \mathbb{E}_{\mathcal{T} \sim \pi | s_0=s} \sum_{t \geq 0} \gamma^t (\nabla_{\theta} \pi_{\theta}(s_t))^{\top} \nabla_a Q^{\pi}(s_t, a)|_{a=\pi_{\theta}(s_t)}.$$

Осталось только применить теорему об эквивалентной форме мат.ожидания по траекториям для

$$f(s, a) = (\nabla_{\theta} \pi_{\theta}(s))^{\top} \nabla_a Q^{\pi}(s, a)|_{a=\pi_{\theta}(s)}.$$



Сразу построим суррогатную функцию для такой формулы градиента:

$$\mathcal{L}_{\tilde{\pi}}(\theta) := \frac{1}{1 - \gamma} \mathbb{E}_{d_{\tilde{\pi}}(s)} Q^{\tilde{\pi}}(s, \pi_{\theta}(s)).$$

Действительно, если мы посчитаем градиент этой функции по θ , то мы просто получим формулу chain rule для оптимизации параметров стратегии через градиент Q-функции по действиям. Иными словами, градиент по параметрам детерминированной стратегии указывает просто проводить policy improvement: выбирать те действия, для которых Q-функция больше, используя её градиент по действиям. Если мы хотим построить Actor-Critic схему, воспользовавшись такой формулой, нам придётся аппроксимировать Q-функцию $Q_{\omega}(s, a) \approx Q^{\pi}(s, a)$ и явно использовать её градиент по действиям, надеясь на то, что $\nabla_a Q_{\omega}(s, a) \approx \nabla_a Q^{\pi}(s, a)$.

Итого, будем сэмплировать батч состояний из реплей буфера и делать шаг градиентного подъёма:

$$\theta \leftarrow \theta + \alpha \mathbb{E}_s (\nabla_{\theta} \pi_{\theta}(s))^{\top} \nabla_a Q^{\pi}(s, a)|_{a=\pi_{\theta}(s)},$$

где состояния s приходят из произвольного распределения (например, из реплей буфера). Это эквивалентно одному шагу градиентной оптимизации суррогатной функции:

$$\mathbb{E}_s Q^{\pi}(s, \pi_{\theta}(s)) \rightarrow \max_{\theta}. \quad (2)$$

Q -функцию $Q_{\omega}(s, a) \approx Q^{\pi}(s, a)$, необходимую для такой оптимизации, будем тоже учить в off-policy режиме с одношаговых оценок: ему для данной пары s, a требуется лишь сэмпл s' , поэтому такого критика можно обучать по переходам $\mathbb{T} := (s, a, r, s')$ из буфера на таргеты

$$y(\mathbb{T}) := r + \gamma Q_{\omega-}(s', \pi(s')).$$

Теорема

Предыдущие две схемы (вывод через DQN и через Policy Gradient) эквивалентны полностью.

Доказательство. Методом пристального взгляда. ■

В рассмотренной схеме из-за использования детерминированной стратегии, как и в DQN, возникает проблема exploration-exploitation-a. В непрерывных пространствах действий вместо ε -жадной стратегии возможно добавлять к выбранным стратегией действиям шум из нормального распределения:

$$a_t := \pi(s_t) + \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(0, \sigma^2 I).$$

Пример

Если действия робота — это направление движения (например, поворот руля управляемой машины), а один шаг в среде это доля секунды, странно проводить исследования, случайно «подёргиваясь» пару раз в секунду. Хочется целенаправленно смещать траекторию: если мы решили в целях исследования повернуть руль чуть правее, чем говорит наша детерминированная стратегия, следует сохранить это смещение руля вправо и в дальнейшем. Для моделирования этого шум должен быть скоррелированным: поэтому вместо независимого шума имеет смысл добавлять случайный процесс, колеблющийся вокруг нуля.

Определение

Шум Орнштейна — Уленбека (Ornstein–Uhlenbeck noise), в начале эпизода инициализированный нулём, задаётся рекурсивно как:

$\varepsilon_{t+1} := \alpha \varepsilon_t + \mathcal{N}(0, \sigma^2 I)$, где $\alpha \leq 1$ и σ — гиперпараметры.

Алгоритм 1: Deep Deterministic Policy Gradient (DDPG)

Гиперпараметры: B — размер мини-батчей, β — коэф. экспоненциального сглаживания для таргет-сеток, α, σ — параметры шума, $Q_\theta(s, a)$ — нейросетка с параметрами θ , $\pi_\omega(s)$ — детерминированная стратегия с параметрами ω , SGD-оптимизаторы, K — периодичность обновления весов таргета.

Инициализировать θ, ω произвольно

Положить $\theta^- := \theta$

Положить $\omega^- := \omega$

Инициализировать шум $\varepsilon := 0$

Пронаблюдать s_0

На очередном шаге t :

- 1 обновить шум $\varepsilon \leftarrow \alpha\varepsilon + \hat{\varepsilon}$, где $\hat{\varepsilon} \sim \mathcal{N}(0, \sigma^2 I)$
- 2 выбрать $a_t := \pi_\omega(s_t) + \varepsilon$
- 3 пронаблюдать $r_t, s_{t+1}, \text{done}_{t+1}$
- 4 добавить пятёрку $(s_t, a_t, r_t, s_{t+1}, \text{done}_{t+1})$ в реплей буфер
- 5 засэмплировать мини-батч размера B из буфера

6 сделать один шаг градиентного подъёма по ω :

$$\frac{1}{B} \sum_{s \in B} Q_{\theta}(s, \pi_{\omega}(s)) \rightarrow \max_{\omega}$$

7 для каждого перехода $\mathbb{T} := (s, a, r, s', \text{done})$ посчитать таргет:

$$y(\mathbb{T}) := r + \gamma(1 - \text{done})Q_{\theta^{-}}(s', \pi_{\omega^{-}}(s'))$$

8 сделать один шаг градиентного спуска по θ :

$$\frac{1}{B} \sum_{\mathbb{T}} (Q_{\theta}(s, a) - y(\mathbb{T}))^2 \rightarrow \min_{\theta}$$

9 обновить таргет-сети, если $t \bmod K = 0$:

$$\theta^{-} \leftarrow (1 - \beta)\theta^{-} + \beta\theta$$

$$\omega^{-} \leftarrow (1 - \beta)\omega^{-} + \beta\omega$$

Алгоритм 2: Twin Delayed DDPG (TD3)

Гиперпараметры: B — размер мини-батчей, N — периодичность обновления весов стратегии, α, σ — параметры шума, $\hat{\sigma}, c$ — параметры шума для добавки к действиям для таргета, β — коэф. экспоненциального сглаживания для таргет-сеток, $Q_{\theta_1}(s, a), Q_{\theta_2}(s, a)$ — нейросетки с параметрами θ_1, θ_2 , $\pi_{\omega}(s)$ — детерминированная стратегия с параметрами ω , SGD-оптимизаторы.

Инициализировать $\theta_1, \theta_2, \omega$ произвольно

Инициализировать таргет-сетки $\theta_1^- := \theta_1, \theta_2^- := \theta_2, \omega^- := \omega$

Инициализировать шум $\varepsilon_0 := 0$

Пронаблюдать s_0

На очередном шаге t :

- 1 посчитать шум $\varepsilon_t := \alpha \varepsilon_{t-1} + \varepsilon$, где $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$
- 2 выбрать $a_t := \pi_{\omega}(s_t) + \varepsilon_t$
- 3 пронаблюдать $r_t, s_{t+1}, \text{done}_{t+1}$
- 4 добавить пятёрку $(s_t, a_t, r_t, s_{t+1}, \text{done}_{t+1})$ в реплей буфер
- 5 засэмплировать мини-батч размера B из буфера

6 для каждого перехода $\mathbb{T} := (s, a, r, s', \text{done})$ посчитать таргет:

$$\varepsilon' \sim \text{clip}(\mathcal{N}(0, \hat{\sigma}I), -c, c)$$

$$y(\mathbb{T}) := r + \gamma(1 - \text{done}) \min_{i \in \{1, 2\}} Q_{\theta_i^-}(s', \pi_{\omega^-}(s') + \varepsilon')$$

7 сделать один шаг градиентного спуска по θ_1 и θ_2 :

$$\frac{1}{B} \sum_{\mathbb{T}} (Q_{\theta_1}(s, a) - y(\mathbb{T}))^2 \rightarrow \min_{\theta_1}$$

$$\frac{1}{B} \sum_{\mathbb{T}} (Q_{\theta_2}(s, a) - y(\mathbb{T}))^2 \rightarrow \min_{\theta_2}$$

8 если $t \bmod N = 0$:

- сделать один шаг градиентного подъёма по ω :

$$\frac{1}{B} \sum_{s \in B} Q_{\theta_1}(s, \pi_{\omega}(s)) \rightarrow \max_{\omega}$$

- обновить таргет-сети:

$$\theta_1^- \leftarrow (1 - \beta)\theta_1^- + \beta\theta_1$$

$$\theta_2^- \leftarrow (1 - \beta)\theta_2^- + \beta\theta_2$$

$$\omega^- \leftarrow (1 - \beta)\omega^- + \beta\omega$$

Обучение стохастических политик

Определение

Скажем, что для параметризации $\pi_\theta(a \mid s)$ применим **репараметризационный трюк** (reparameterization trick), если сэмплирование $a \sim \pi_\theta(a \mid s)$ эквивалентно сэмплированию шума из некоторого не зависящего от параметров распределения $\varepsilon \sim p(\varepsilon)$ и его дальнейшего детерминированного преобразования $a = f_\theta(s, \varepsilon)$.

Обучение стохастических политик

Пример 1: Пусть наша стратегия параметризована нормальным распределением:

$$\pi_{\theta}(a \mid s) := \mathcal{N}(\mu_{\theta}(s), \sigma_{\theta}(s)^2 I).$$

Тогда для неё применим репараметризационный трюк: сэмплирование действий эквивалентно $a := \mu_{\theta}(s) + \varepsilon \odot \sigma_{\theta}(s)$, где $\varepsilon \sim \mathcal{N}(0, I)$, \odot — поэлементное перемножение.

Пример

Семейство детерминированных стратегий $\pi_{\theta}(s)$ тоже можно считать таким «вырожденным» примером параметризаций, для которой можно проворачивать репараметризационный трюк: просто шум ε считаем «пустым».

Теорема

Если для $\pi_\theta(a \mid s)$ применим репараметризационный трюк, то:

$$\nabla_\theta J(\pi_\theta) = \frac{1}{1 - \gamma} \mathbb{E}_{d_{\pi_\theta}(s)} \mathbb{E}_{\varepsilon \sim p(\varepsilon)} \nabla_\theta f_\theta(s, \varepsilon) \nabla_a Q^\pi(s, a) \big|_{a=f_\theta(s, \varepsilon)}. \quad (3)$$

Доказательство. Полностью полностью повторяет вывод теоремы 1. ■

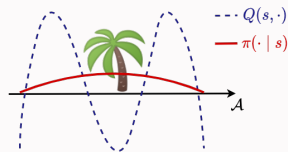
Таким образом, все идеи DDPG расширяются на этот случай. Убирая из формулы градиента частоты посещения состояний и переходя к policy iteration схеме, получаем следующий функционал:

$$\mathbb{E}_s \mathbb{E}_{a \sim \pi_\theta(a|s)} Q_\omega(s, a) \rightarrow \max_{\theta}, \quad (4)$$

где состояния берутся из буфера. В силу репараметризационного трюка мы легко справимся со взятием градиента на практике:

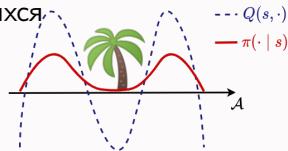
$$\nabla_{\theta} \mathbb{E}_{a \sim \pi_\theta(a|s)} Q_\omega(s, a) = \mathbb{E}_{\varepsilon \sim p(\varepsilon)} \nabla_{\theta} Q_\omega(s, f_\theta(s, \varepsilon)).$$

Пример 2: Представьте, что вы хотите объехать дерево. Вы можете объехать его справа, можете слева. Критик сообщает вам высокие значения и слева, и справа, и оптимизируя (4) в классе гауссиан, можно получить стратегию, которая с наибольшей вероятностью выбирает действие «врезаться в дерево».



Пример 3: Например, можно использовать смесь гауссиан. Тогда при использовании K компонент смеси актёр для данного состояния s выдаёт следующие величины: K суммирующихся в единицу чисел $w_i(s, \theta)$, а также K векторов $\mu_i(s, \theta), \sigma_i(s, \theta)$, где $i \in 1, 2, \dots, K$. Итоговое распределение полагается

$$\pi_{\theta}(a \mid s) := \sum_{i=1}^K w_i(s, \theta) \mathcal{N}(\mu_i(s, \theta), \sigma_i(s, \theta)^2 I).$$



Утверждение

Градиент (4) по параметрам актёра θ равен

$$\nabla_{\theta} = \mathbb{E}_{a \sim \pi_{\theta}(a|s)} \nabla_{\theta} \log \pi_{\theta}(a | s) (Q_{\omega}(s, a)(s, a) - b(s)) ,$$

где $b(s)$ — бэйзлайн, произвольная функция от состояний.

Теорема (Формула замены переменной в плотности)

Пусть $\pi_\theta(a \mid s) := g(u)$, где $g: \mathbb{R}^A \rightarrow \mathbb{R}^A$, и $u \sim \mu_\theta(u \mid s)$. Тогда:

$$\log \pi_\theta(a \mid s) = \log \mu_\theta(a \mid s) - \log |\det \nabla_u g|. \quad (5)$$

Без доказательства. ■

Определение

Энтропией распределения $\pi(a)$ называется

$$\mathcal{H}(\pi(a)) := -\mathbb{E}_{\pi(a)} \log \pi(a). \quad (6)$$

Пример

Например, для функции $g(u) := \tanh(u)$ якобиан $\nabla_u g$ есть диагональная матрица (поскольку преобразование поэлементное), и его определитель равен покомпонентному произведению:

$$\begin{aligned}\det \nabla_u g(u) &= \prod_{i=1}^A \nabla_{u_i} \tanh(u_i) = \\ &= \{ \text{производная гиперболического тангенса} \} = \prod_{i=1}^A (1 - \tanh^2(u_i))\end{aligned}$$

Заметим, что все компоненты положительные ($\tanh(u_i) \leq 1$), поэтому модуль из формулы (5) брать не нужно. Подставляя в (5), получаем окончательно:

$$\log \pi_{\theta}(a \mid s) = \log \mu_{\theta}(a \mid s) - \sum_{i=1}^A \log(1 - \tanh^2(u_i)).$$

Soft Actor-Critic

Определение

Задачей **Maximum Entropy RL** является максимизация функционала

$$J_{\text{soft}}(\pi) := \mathbb{E}_{\mathcal{T} \sim \pi} \sum_{t \geq 0} \gamma^t [r_t + \alpha \mathcal{H}(\pi(\cdot \mid s))] \rightarrow \max_{\pi}, \quad (7)$$

где α — гиперпараметр, называемый **температурой** (temperature).

Soft Actor-Critic

Утверждение

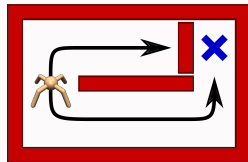
Задача (7) эквивалентна

$$J_{\text{soft}}(\pi) := \mathbb{E}_{\mathcal{T} \sim \pi} \sum_{t \geq 0} \gamma^t [r_t - \alpha \log \pi(a_t | s_t)] \rightarrow \max_{\pi},$$

Доказательство. Следует из определения энтропии (6); ведь мат.ожидание по действиям присутствует в мат.ожидании по траекториям.



Пример 4: Представим, что у агента есть два пути, и по мере углубления награда на каждом пути одинаково растёт. Пусть первый путь заканчивается тупиком и суммарно позволяет набрать не более +100, а на втором тупик стоит чуть дальше и даёт +110. Во время обучения агент может уловить награду вдоль первого пути и учиться углубляться в него, игнорируя исследование второго пути, даже если агент умеет набирать там награду как на первом; за счёт бонуса за наиболее стохастичную стратегию агент мотивирован в течение обучения в начале эпизодов случайно выбирать между двумя путями. То есть, энтропийный бонус помогает избегать подобных «локальных максимумов» в среде.



Можно считать, что в данном фреймворке мы на самом деле лишь чуть-чуть модифицировали награду в среде:

$$r_{\text{soft}}(s, a) := r(s, a) - \alpha \log \pi(a \mid s). \quad (8)$$

Утверждение

Оптимальной детерминированной стратегии может не существовать.

Доказательство. В MDP, где награда всегда 0, оптимальна стратегия с максимальной энтропией. ■

Теорема (Мягкие уравнения связи)

$$Q_{\text{soft}}^{\pi}(s, a) = r(s, a) + \gamma \mathbb{E}_{s'} V_{\text{soft}}^{\pi}(s') \quad (9)$$

$$V_{\text{soft}}^{\pi}(s) = \mathbb{E}_{\pi(a|s)} [Q_{\text{soft}}^{\pi}(s, a) - \log \pi(a | s)] \quad (10)$$

Доказательство. По определению с учётом договорённости. ■

Теорема (Мягкие уравнения Беллмана (soft Bellman equations))

$$Q_{\text{soft}}^{\pi}(s, a) = r(s, a) + \gamma \mathbb{E}_{s'} \mathbb{E}_{a'} [Q_{\text{soft}}^{\pi}(s', a') - \log \pi(a' | s')] \quad (11)$$

$$V_{\text{soft}}^{\pi}(s) = \mathbb{E}_a [r(s, a) - \log \pi(a | s) + \gamma \mathbb{E}_{s'} V_{\text{soft}}^{\pi}(s')] \quad (12)$$

Теорема

Операторы, стоящие в правой части мягких уравнений Беллмана, являются сжимающими с коэффициентом γ по метрике

$$\rho(V_1, V_2) := \max_s |V_1(s) - V_2(s)|.$$

Доказательство. Покажем для мягкой V -функции. Пусть $\mathfrak{B}_{\text{soft}}$ — оператор, стоящий в правой части (12), и пусть даны две V -функции V_1, V_2 . Тогда:

$$|[\mathfrak{B}_{\text{soft}} V_1](s) - [\mathfrak{B}_{\text{soft}} V_2](s)| = \gamma |\mathbb{E}_a \mathbb{E}_{s'} [V_1(s') - V_2(s')]|,$$

поскольку энтропия стратегии π вместе с наградой за шаг одинакова для V_1 и V_2 и потому сокращается. Дальше, как и для обычных V -функций, можно просто оценить это выражение сверху $\gamma \rho(V_1, V_2)$, заканчивая доказательство. ■

Утверждение

Метод простой итерации сходится к единственному решению мягких уравнений Беллмана из любого начального приближения.

Итак, мы уже можем сразу построить процедуру обучения критика. Рассмотрим обучение $Q_\omega(s, a) \approx Q_{\text{soft}}^\pi(s, a)$ с одношаговых оценок в off-policy режиме, то есть будем просто решать мягкое уравнение Беллмана (11). Тогда для заданного перехода $\mathbb{T} = (s, a, r, s')$ целевая переменная строится как

$$y(\mathbb{T}) := r + \gamma \mathbb{E}_{a' \sim \pi(a' | s')} [Q_\omega(s', a') - \log \pi(a' | s')].$$

В таком варианте таргеты для критика (Q-функция с параметрами ω) и для вспомогательной V-функции выглядят следующим образом: для заданного перехода $\mathbb{T} = (s, a, r, s')$, взятого из буфера, генерируем a_π из текущей версии стратегии π и запоминаем вероятность $\pi(a_\pi | s)$, после чего вычисляем несмещённые оценки правых частей уравнений связи (10) и (9):

$$y_Q(\mathbb{T}) := r + \gamma V_\psi(s'),$$

$$y_V(\mathbb{T}) := Q_\omega(s, a_\pi) - \log \pi(a_\pi | s).$$

Теорема 2 — Soft Policy Improvement: Пусть стратегии π_1 и π_2 таковы, что для всех состояний s выполняется:

$$\mathbb{E}_{\pi_2(a|s)} Q_{\text{soft}}^{\pi_1}(s, a) + \mathcal{H}(\pi_2(\cdot | s)) \geq V_{\text{soft}}^{\pi_1}(s),$$

тогда $\pi_2 \succeq \pi_1$ с учётом энтропийного бонуса; если хотя бы для одного s неравенство выполнено строго, то $\pi_2 \succ \pi_1$.

Доказательство. Полностью аналогично доказательству в обычном случае. Покажем, что $V_{\text{soft}}^{\pi_2}(s) \geq V_{\text{soft}}^{\pi_1}(s)$ для любого s :

$$\begin{aligned}
 V_{\text{soft}}^{\pi_1}(s) &\leq \{\text{по построению } \pi_2\} \leq \mathbb{E}_{\pi_2(a|s)} Q_{\text{soft}}^{\pi_1}(s, a) + \mathcal{H}(\pi_2(\cdot | s)) = \\
 &= \{\text{связь QV (9)}\} = \mathbb{E}_{\pi_2(a|s)} [r + \mathcal{H}(\pi_2(\cdot | s)) + \gamma \mathbb{E}_{s'} V_{\text{soft}}^{\pi_1}(s')] \leq \\
 &\leq \{\text{применяем это же неравенство рекурсивно}\} = \mathbb{E}_{\pi_2(a|s)} [r + \mathcal{H}(\pi_2(\cdot | s)) + \\
 &\quad + \mathbb{E}_{s'} \mathbb{E}_{\pi_2(a'|s')} [\gamma r' + \gamma \mathcal{H}(\pi_2(\cdot | s')) + \gamma^2 \mathbb{E}_{s''} V_{\text{soft}}^{\pi_1}(s'')]] \leq \\
 &\leq \{\text{раскручиваем цепочку далее}\} \leq \dots \leq \mathbb{E}_{\mathcal{T} \sim \pi_2 | s_0=s} \sum_{t \geq 0} \gamma^t r_t + \gamma^t \mathcal{H}(\pi_2(\cdot | s_t)) = \\
 &= \{\text{по определению мягкой V-функции}\} = V_{\text{soft}}^{\pi_2}(s)
 \end{aligned}$$

Если для какого-то s неравенство из условия теоремы было выполнено строго, то для него первое неравенство в этой цепочке рассуждений выполняется строго, и, значит, $V_{\text{soft}}^{\pi_2}(s) > V_{\text{soft}}^{\pi_1}(s)$. ■

Таким образом, аналог общей схемы Generalized Policy Iteration в задаче Maximum Entropy RL выглядит так:

- **Soft Policy Evaluation** заключается в обучении аппроксимации мягкой оценочной функции $Q_\omega \approx Q_{\text{soft}}^\pi$ для текущей стратегии π ;
- **Soft Policy Improvement** заключается в оптимизации следующего функционала (для разных состояний s):

$$\mathbb{E}_{\pi(a|s)} Q_\omega(s, a) + \mathcal{H}(\pi(\cdot | s)) \rightarrow \max_{\pi}. \quad (13)$$

Теорема 3: Задача (13) эквивалентна задаче

$$\text{KL}(\pi_\theta(\cdot | s) \parallel \exp Q_\omega(s, \cdot)) \rightarrow \min_\theta,$$

где $\exp Q_\omega(s, \cdot)$ — ненормированное распределение над действиями.

Доказательство. Обозначим нормировочную константу распределения $\exp Q_\omega(s, \cdot)$ как $Z_\omega(s) := \int_{\mathcal{A}} \exp Q_\omega(s, a) da$. Тогда:

$$\text{KL}(\pi_\theta(\cdot | s) \parallel \exp Q_\omega(s, \cdot)) = \mathbb{E}_{a \sim \pi_\theta(a|s)} \log \pi_\theta(a | s) - Q_\omega(s, a) + \log Z_\omega(s).$$

Осталось заметить, что при домножении на минус получим (13): первое слагаемое есть энтропия, а третье слагаемое не зависит от оптимизируемых параметров θ :

$$\mathbb{E}_{a \sim \pi_\theta(a|s)} \log Z_\omega(s) = \log Z_\omega(s) = \text{const}(\theta). \quad \blacksquare$$

Теорема (Вид жадной стратегии)

Максимальное значение (13) достигается на стратегии

$$\pi(a \mid s) \propto \exp Q_{\omega}(s, a). \quad (14)$$

Доказательство. Следует из теоремы 3, поскольку минимум KL-дивергенции достигается в нуле на совпадающих распределениях. ■

Пример 5: Пусть наша стратегия параметризована гауссианой (см. пример 1). Для неё можно проводить репараметризационный трюк и также можно аналитически посчитать энтропию:

$$\mathcal{H}(\mathcal{N}(\mu, \sigma^2 I)) = \sum_{i=1}^A \log \sigma_i.$$

Итого, формула (13) в таком случае получается следующей:

$$\mathbb{E}_{\varepsilon \sim \mathcal{N}(0, I)} Q_{\omega}(s, \mu_{\theta}(s) + \sigma_{\theta}(s) \odot \varepsilon) + \sum_{i=1}^A \log \sigma_i(s, \theta) \rightarrow \max_{\theta}.$$

Пример 6: Пусть наша стратегия параметризована смесью гауссиан (см. пример 3). Тогда для неё не применим репараметризационный трюк, и сложно аналитически посчитать энтропию. Тогда придётся применять REINFORCE, и формула градиента (13) получается следующей:

$$\nabla_{\theta} = \mathbb{E}_{a \sim \pi_{\theta}(a|s)} \nabla_{\theta} \log \pi_{\theta}(a | s) [Q_{\omega}(s, a) - \log \pi_{\theta}(a | s) - b(s)],$$

где $b(s)$ — бэйзлайн, произвольная функция от состояний. Имеет смысл делать её близкой к среднему значению $Q_{\omega}(s, a) - \log \pi_{\theta}(a | s)$, то есть хорошо выбирать $b(s) := V_{\omega}(s) = \mathbb{E}_a [Q_{\omega}(s, a) - \log \pi_{\theta}(a | s)]$.

Алгоритм 3: Soft Actor-Critic (SAC)

Гиперпараметры: B — размер мини-батчей, β — параметр экспоненциального сглаживания таргет-сети, α — температура, $\pi_\theta(a \mid s) := \mathcal{N}(\mu_\theta(s), \sigma_\theta(s)^2 I)$ — гауссова стратегия с параметрами θ , $Q_{\omega_1}(s, a), Q_{\omega_2}(s, a)$ — две нейросети с параметрами ω_1 и ω_2 , $V_\psi(s)$ — нейросетка с параметрами ψ , SGD-оптимизаторы.

Инициализировать $\theta, \omega_1, \omega_2, \psi$ произвольно

Инициализировать таргет-сеть $\psi^- := \psi$

Пронаблюдать s_0

На очередном шаге t :

- 1 выбрать $a_t \sim \pi_\theta(a_t \mid s_t)$
- 2 пронаблюдать $r_t, s_{t+1}, \text{done}_{t+1}$
- 3 добавить пятёрку $(s_t, a_t, r_t, s_{t+1}, \text{done}_{t+1})$ в реплей буфер
- 4 засэмплировать мини-батч размера B из буфера
- 5 для каждого s из батча засэмплировать шума $\varepsilon(s) \sim \mathcal{N}(0, I)$ и посчитать $\mu(s, \theta), \sigma(s, \theta)$ стратегии π_θ

6 посчитать оценку градиента по параметрам стратегии:

$$\nabla_{\theta} := \frac{1}{B} \sum_{s \in B} \nabla_{\theta} \left[\alpha \sum_{i=1}^A \log \sigma_i(s, \theta) + \min_{i=1,2} Q_{\omega_i}(s, \mu_{\theta}(s) + \sigma_{\theta}(s) \odot \varepsilon(s)) \right]$$

7 делаем шаг градиентного подъёма по θ , используя ∇_{θ}

8 для каждого перехода $\mathbb{T} := (s, a, r, s', \text{done})$ засэмплировать $a_{\pi} \sim \pi_{\theta}(a_{\pi} \mid s)$ и сохранить вероятности $\pi_{\theta}(a_{\pi} \mid s)$

9 посчитать таргеты:

$$y_V(\mathbb{T}) := \min_{i=1,2} Q_{\omega_i}(s, a_{\pi}) - \alpha \log \pi_{\theta}(a_{\pi} \mid s)$$

$$y_Q(\mathbb{T}) := r + \gamma V_{\psi^{-}}(s')$$

10 посчитать лоссы:

$$\text{Loss}_V(\psi) := \frac{1}{B} \sum_{\mathbb{T}} (V_\psi(s') - y_V(\mathbb{T}))^2$$

$$\text{Loss}_{Q1}(\omega_1) := \frac{1}{B} \sum_{\mathbb{T}} (Q_{\omega_1}(s, a) - y_Q(\mathbb{T}))^2$$

$$\text{Loss}_{Q2}(\omega_2) := \frac{1}{B} \sum_{\mathbb{T}} (Q_{\omega_2}(s, a) - y_Q(\mathbb{T}))^2$$

11 делаем шаг градиентного спуска по ψ , ω_1 и ω_2 , используя $\nabla_\psi \text{Loss}_V(\psi)$, $\nabla_\omega \text{Loss}_{Q1}(\omega_1)$ и $\nabla_{\omega_2} \text{Loss}_{Q2}(\omega_2)$ соответственно

12 обновляем таргет-сеть: $\psi^- \leftarrow (1 - \beta)\psi^- + \beta\psi$

Сформулируем критерий оптимальности Беллмана в задаче Maximum Entropy RL. По аналогии с обычным случаем, введём оптимальные оценочные функции:

$$V_{\text{soft}}^*(s) := \max_{\pi} V_{\text{soft}}^{\pi}(s)$$

$$Q_{\text{soft}}^*(s, a) := \max_{\pi} Q_{\text{soft}}^{\pi}(s, a)$$

Теорема (Вид оптимальной стратегии для Maximum Entropy RL)

Оптимальной является единственная стратегия

$$\pi(a \mid s) \propto \exp Q_{\text{soft}}^*(s, a)$$

Доказательство. Откажемся от стационарности и будем рассматривать задачу поиска оптимальной стратегии $\pi_t(a \mid s_0)$ в предположении, что в будущем мы сможем набрать максимально возможную награду $Q_{\text{soft}}^*(s, a, t)$:

$$\mathbb{E}_{\pi_t(a|s)} [Q_{\text{soft}}^*(s, a, t) - \log \pi_t(a \mid s)] \rightarrow \max_{\pi_t(a|s)}$$

Аналогично теореме 3 про soft policy improvement, можно заметить, что с точностью до константы, не зависящей от π_t , оптимизируемое выражение есть:

$$- \text{KL} \left(\pi_t(a \mid s) \parallel \frac{\exp Q_{\text{soft}}^*(s, a, t)}{Z(s, t)} \right) \rightarrow \max_{\pi_t(a|s)},$$

где $Z(s, t)$ — нормировочная константа $\exp Q_{\text{soft}}^*(s, a, t)$. Понятно, что оптимум достигается в нуле на $\pi_t(a \mid s)$, совпадающей с этим распределением.

Дальнейшее рассуждение строится как раньше: Q_{soft}^* от времени не зависит по определению, поэтому оптимальная стратегия получается стационарной, следовательно

$$\pi(a \mid s) \propto \exp Q_{\text{soft}}^*(s, a)$$

Заметим, что в силу однозначного определения $Q_{\text{soft}}^*(s, a)$, такая стратегия в принципе единственна в отличие от обычной задачи RL.



Теорема

$$V_{\text{soft}}^*(s) = \log \int_{\mathcal{A}} \exp Q_{\text{soft}}^*(s, a) da \quad (15)$$

Доказательство. Мы знаем, что оптимальная стратегия имеет вид $\pi^*(a | s) = \frac{\exp Q_{\text{soft}}^*(s, a)}{Z(s)}$, где

$$Z(s) := \int_{\mathcal{A}} \exp Q_{\text{soft}}^*(s, a) da$$

является нормировочной константой. Посчитаем энтропию такого распределения:

$$\begin{aligned} -\mathbb{E}_{\pi^*(a|s)} \log \pi^*(a | s) &= \int_{\mathcal{A}} \frac{\exp Q_{\text{soft}}^*(s, a)}{Z(s)} (\log Z(s) - Q_{\text{soft}}^*(s, a)) da = \\ &= \log Z(s) - \mathbb{E}_{\pi^*(a|s)} Q_{\text{soft}}^*(s, a). \end{aligned}$$

Подставим оптимальную стратегию в мягкое VQ уравнение (10), которое справедливо в том числе и для оптимальной стратегии:

$$\begin{aligned} V_{\text{soft}}^*(s) &= \mathbb{E}_{\pi^*(a|s)} [Q_{\text{soft}}^*(s, a) - \log \pi^*(a | s)] = \\ &= \mathbb{E}_{\pi^*(a|s)} Q_{\text{soft}}^*(s, a) - \mathbb{E}_{\pi^*(a|s)} Q_{\text{soft}}^*(s, a) + \log Z(s) = \\ &= \log Z(s) \end{aligned}$$

Вспоминая определение $Z(s)$, получаем доказываемое.



Утверждение

$$Q_{\text{soft}}^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s'} V_{\text{soft}}^*(s')$$

Утверждение (Мягкое уравнение оптимальности Беллмана (soft Bellman optimality equations))

$$Q_{\text{soft}}^*(s, a) = r(s, a) + \gamma \mathbb{E}_{s'} \log \int_{\mathcal{A}} \exp Q_{\text{soft}}^*(s', a') da' \quad (16)$$

Теорема

Оператор, стоящий в правой части уравнения (16), является сжимающим с коэффициентом γ , и, следовательно, метод простой итерации решения этой системы уравнений сходится из любого начального приближения к единственной неподвижной точке.

Доказательство. Пусть даны две Q-функции, Q_1, Q_2 , и пусть

$$\rho(Q_1, Q_2) := \max_{s,a} |Q_1(s, a) - Q_2(s, a)| < \varepsilon.$$

Тогда:

$$\log \int_{\mathcal{A}} \exp Q_1(s, a) da \leq \log \int_{\mathcal{A}} \exp(Q_2(s, a) + \varepsilon) da = \varepsilon + \log \int_{\mathcal{A}} \exp Q_2(s, a) da.$$

Аналогично можно показать, что

$$\log \int_{\mathcal{A}} \exp Q_1(s, a) da \geq -\varepsilon + \log \int_{\mathcal{A}} \exp Q_2(s, a) da$$

Пусть $\mathfrak{B}_{\text{soft}}$ — оператор, стоящий в правой части (16). Тогда:

$$\begin{aligned} & |[\mathfrak{B}_{\text{soft}} Q_1](s, a) - [\mathfrak{B}_{\text{soft}} Q_2](s, a)| = \\ & = \gamma |\mathbb{E}_{s'} \left(\log \int_{\mathcal{A}} \exp Q_1(s', a') da' - \log \int_{\mathcal{A}} \exp Q_2(s', a') da' \right)| \leq \gamma \varepsilon. \end{aligned}$$

Таким образом, $\rho(\mathcal{B}_{\text{soft}} Q_1, \mathcal{B}_{\text{soft}} Q_2)$ уменьшилось по крайней мере в γ раз по сравнению с $\rho(Q_1, Q_2)$.



Утверждение

Если $\pi(a | s) \propto \exp Q_{\text{soft}}^{\pi}(s, a)$, то она оптимальна.

Доказательство. Q-функция такой стратегии удовлетворяет мягкому уравнению оптимальности Беллмана (16) и в силу единственности его решения совпадает с $Q_{\text{soft}}^(s, a)$.*



Теперь можно построить аналог DQN для задачи Maximum Entropy RL, называемым **Soft Q-learning**. В нём нет отдельной модели актёра, и текущая стратегия просто полагается жадной по отношению к текущему критику. Это также можно интерпретировать как моделирование **Soft Value Iteration** — решение мягкого уравнения оптимальности (16). Тогда таргет для перехода $\mathbb{T} := (s, a, r, s')$ строится как

$$y(\mathbb{T}) := r + \gamma \log \int_{\mathcal{A}} \exp Q_{\theta^-}(s', a') da',$$

где θ^- — параметры таргет-сети.

Теорема (Soft Policy Gradient)

$$\nabla_{\theta} J_{\text{soft}}(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\pi_{\theta}}(s)} \mathbb{E}_{a \sim \pi_{\theta}(a|s)} \nabla_{\theta} \log \pi_{\theta}(a | s) Q_{\text{soft}}^{\pi}(s, a) + \alpha \nabla_{\theta} \mathcal{H}(\pi_{\theta}(\cdot | s)).$$

Доказательство. Аналогично доказательству в обычном случае. ■