

ОБОБЩЕНИЯ TD И ПЛАНИРОВАНИЕ

Сергей Николенко

Академия MADE — Mail.Ru

02 октября 2021 г.

Random facts:

- 2 октября 1187 г. Саладин после недолгой осады захватил Иерусалим, что послужило поводом к началу Третьего крестового похода
- 2 октября 1552 г. Иван Грозный вошёл в Казань и присоединил Казанское ханство к России
- 2 октября 1607 г. голландский очковый мастер Иоанн Липперсгей продемонстрировал в Гааге прототип оптического телескопа; Галилей направил его в небо только в 1609 году
- 2 октября 1836 г. Чарльз Дарвин вернулся из длившегося почти пять лет кругосветного путешествия на корабле «Бигль»
- 2 октября 1928 г. святой Хосемария Эскрива де Балагер создал Прелатуру Святого Креста Opus Dei, позже прославленную Дэном Брауном

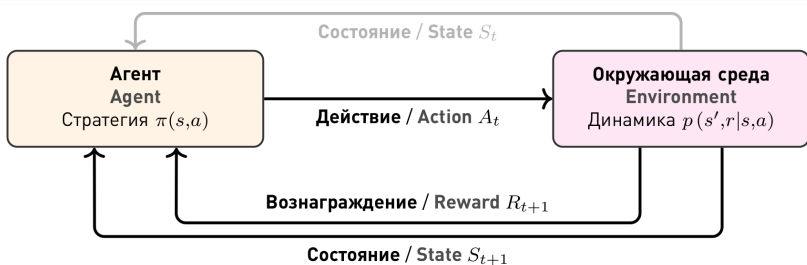
ОБУЧЕНИЕ С ПОДКРЕПЛЕНИЕМ

- В прошлый раз мы ввели основные понятия динамики марковских процессов принятия решений:
 - собственно динамику процесса:

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

- награды за каждый эпизод, начиная со времени t :

$$G_t = R_{t+1} + \gamma G_{t+1} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1};$$



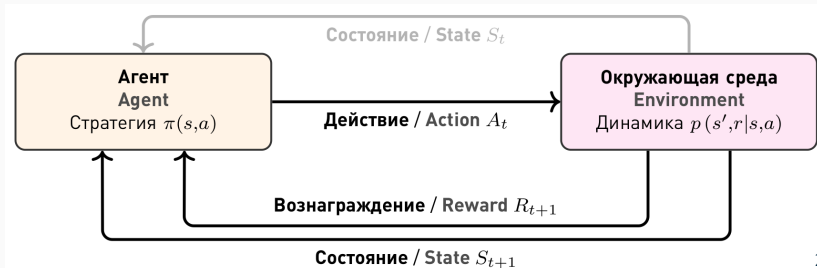
ОБУЧЕНИЕ С ПОДКРЕПЛЕНИЕМ

- Определили функции значений V и Q

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

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

- Выписали уравнения Беллмана и научились их решать.



ОСНОВНЫЕ ЗАДАЧИ

- Теоретически всё готово, но у нас много проблем:
 - уравнения знаем, но пока не знаем, как их решать, то есть как найти V^π для данного π ?
 - разных стратегий очень, очень много — как найти оптимальную стратегию поведения агента в данной модели и соответствующие V^* ?
 - но уравнений тоже не знаем — в реальности обычно P и R не даны, их тоже нужно обучить; как?
 - более того, их обычно даже записать не получится, слишком уж много состояний в любой реальной задаче... что делать?

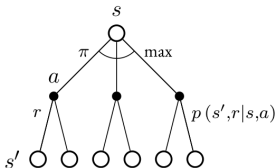


- Давайте есть слона по частям...

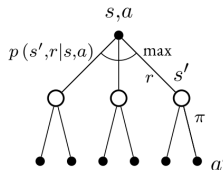
ДИАГРАММЫ ОБХОДА

- Вспомним, о чём мы говорили, и заодно введём новый вид картинок: *диаграммы обхода* (backup diagrams)

Уравнение Беллмана для V_*



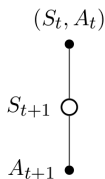
Уравнение Беллмана для Q_*



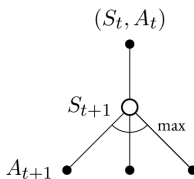
$TD(0)$



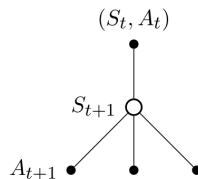
Sarsa





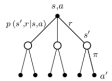


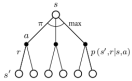
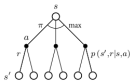
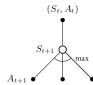
Q -обучение



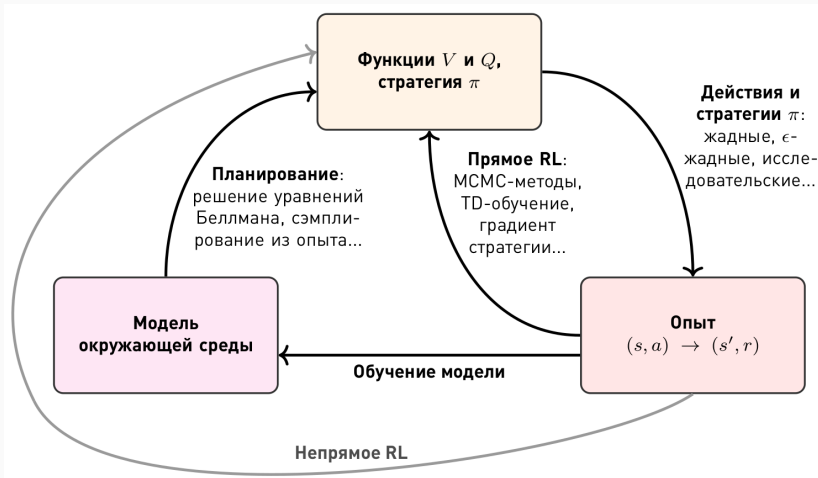
Sarsa с ожиданием



ТАБЛИЦА

Оценка стратегии π	Функция	Обновление в ожидании	Обновление по выборке
	$V_{\pi}(S_t) :=$	Уравнение Беллмана $\sum_{a \in \mathcal{A}} \pi(a S_t) \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r s, a) (r + \gamma V_{\pi}(s'))$ 	TD(0) $V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$ 
Управление, поиск π_*	$Q_{\pi}(S_t, A_t) :=$	Уравнение Беллмана $\sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r S_t, A_t) (r + \gamma \sum_{a' \in \mathcal{A}} \pi(a' s') Q_{\pi}(s', a'))$ 	Sarsa, Sarsa с ожиданием $Q(S_t, A_t) + \alpha (R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$ $Q(S_t, A_t) + \alpha (R_{t+1} + \gamma \sum_a \pi(a S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t))$  
Управление, поиск π_*	$V_*(S_t) :=$	Уравнение Беллмана $\max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r S_t, a) (r + \gamma V_*(s'))$ 	
	$Q_*(S_t, A_t) :=$	Уравнение Беллмана $\sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r S_t, A_t) \left(r + \gamma \max_{a'} Q_*(s', a') \right)$ 	Q-обучение $Q(S_t, A_t) + \alpha (R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t))$ 

ОБУЧЕНИЕ С ПОДКРЕПЛЕНИЕМ



ОБОБЩЕНИЯ И ДОПОЛНИТЕЛЬНЫЕ ЗАМЕЧАНИЯ

- Мы говорили о TD-алгоритмах, которые оценивают

$$G_t = R_{t+1} + \gamma V_t(S_{t+1}).$$

- Но можно же и дальше развернуть:

$$G_{t:t+2} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{t+1}(S_{t+1})$$

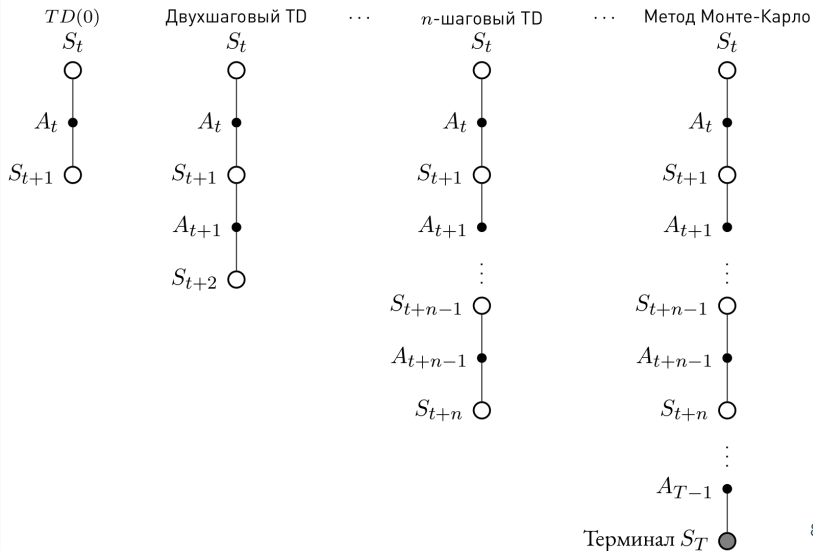
$$G_{t:t+3} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 V_{t+2}(S_{t+3})$$

$$\dots = \dots$$

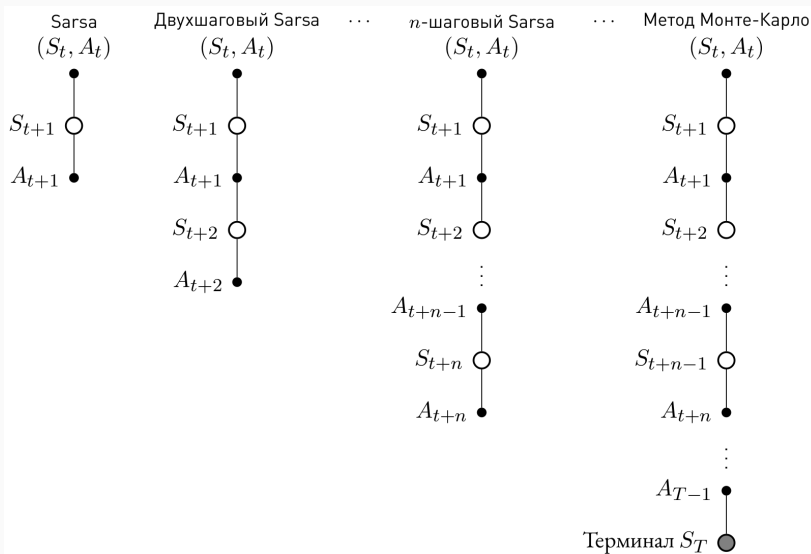
$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n-1}(S_{t+n})$$

- На бесконечности это сливается с методами Монте-Карло, потому что там мы обновляем на основе награды за весь эпизод, от конца к началу

- Обобщение TD-обучения: n -шаговые методы



- Обобщение TD-управления: n -шаговые методы



- И теперь появляются n -шаговые варианты всех алгоритмов.

Алгоритм n -step Sarsa (on-policy TD control):

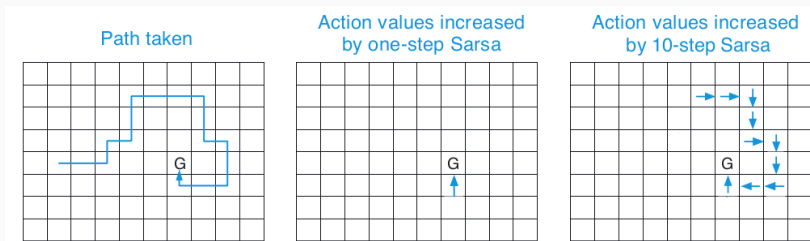
инициализировать случайно $Q(s, a)$; повторять до сходимости:

- инициализировать S_0 , выбрать A_0 по стратегии, полученной из Q (например, по ϵ -жадной стратегии); $T := \infty$;
- для каждого шага в эпизоде $t = 0, \dots, T$:
 - если $t < T$, то:
 - сделать действие A_t , получить R_{t+1}, S_{t+1} ;
 - если это терминальное состояние, то $T := t + 1$, а если нет, выбрать A_{t+1} в состоянии S_{t+1} по стратегии π ;
 - $\tau := t - n + 1$ (мы будем обновлять оценку на шаге τ)
 - если $\tau \geq 0$, то обновляем:
 - $G := \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$
 - если $\tau + n < T$, то $G := G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$;
 - $Q(S_\tau, A_\tau) := Q(S_\tau, A_\tau) + \alpha (G - Q(S_\tau, A_\tau))$
 - если обучаем π , то поменять π на ϵ -жадную по Q

- Та же Sarsa, но дальше заглядывает.

n -ШАГОВЫЕ АЛГОРИТМЫ

- Смысл здесь в том, что n -step алгоритмы обновляют сразу много значений по одной и той же награде, даже если всё остальное пока что нулевое:



- Всё то же самое можно переписать в терминах просто изменения TD-ошибки:

$$G_{t:t+n} = Q_{t-1}(S_t, A_t) + \sum_{k=t}^{\min(t+n, T)-1} \gamma^{k-t} (R_{k+1} + \gamma Q_k(S_{k+1}, A_{k+1}) - Q_{k-1}(S_k, A_k))$$

- Если хочется сделать off-policy, то можно сделать, конечно, с importance sampling, как мы обсуждали.
- А можно сделать поиск по дереву: идём обратно по дереву от (S_t, A_t) к концу эпизода; в каждый момент у нас есть одно действие, которое реально произошло, а для других берём bootstrap-оценки:

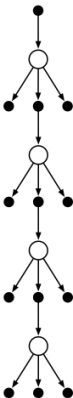
$$\begin{aligned} G_{t:t+1} &= R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) Q_t(S_{t+1}, a), \\ G_{t:t+2} &= R_{t+1} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1}) Q_{t+1}(S_{t+1}, a) \\ &\quad + \gamma \pi(A_{t+1}|S_{t+1}) (R_{t+2} + \gamma \sum_a \pi(a|S_{t+2}) Q_{t+1}(S_{t+2}, a)), \\ &\quad \dots \\ G_{t:t+n} &= R_{t+1} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1}) Q_{t+1}(S_{t+1}, a) \\ &\quad + \gamma \pi(A_{t+1}|S_{t+1}) G_{t+1:t+n}. \end{aligned}$$

- А можно попытаться нечто среднее сделать, выбирая то сэмплы, как в Sarsa, то апдейт по дереву, как выше:

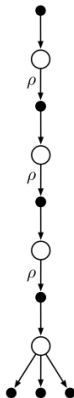
4-step
Sarsa



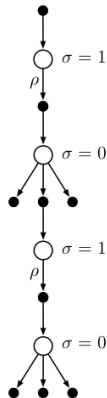
4-step
Tree backup



4-step
Expected Sarsa

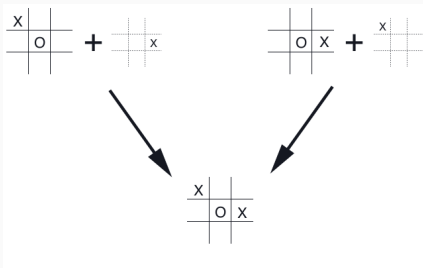


4-step
 $Q(\sigma)$



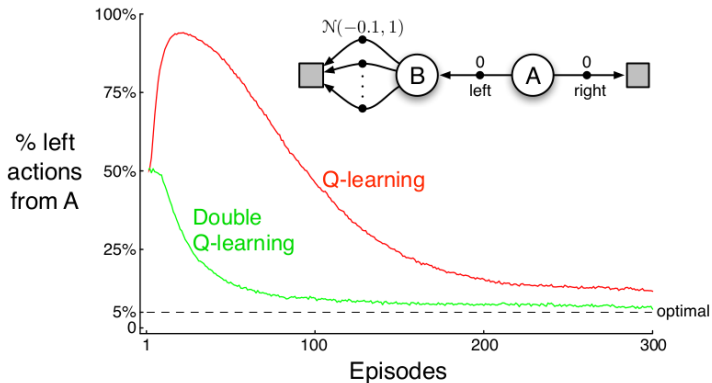
AFTERSTATES

- Часто обучают не по состояниям $V(s)$ и не по парам $Q(s, a)$, а по т.н. *afterstates* – состояниям после действия.
- Это хорошо, когда действия имеют немедленный эффект, а случайный процесс происходит уже потом.
- Крестики-нолики – многие пары приводят к одной и той же позиции.



ДВОЙНОЕ ОБУЧЕНИЕ

- В TD-обучении есть проблема: наша целевая стратегия – это жадная стратегия по текущей Q , т.е. мы используем максимум по оценкам, чтобы оценить максимальное значение, а это вносит смещение в оценки.



- Чтобы это поправить, можно использовать двойное обучение:
 - поддерживаем две стратегии Q_1 и Q_2 ;
 - каждый раз бросаем монетку, выбираем, какую брать стратегию для обновления, а какую для цели:

$$Q_1(S, A) := Q_1(S, A) + \alpha \left(R + \gamma \max_a Q_2(S', a) - Q_1(S, A) \right)$$

или

$$Q_2(S, A) := Q_2(S, A) + \alpha \left(R + \gamma \max_a Q_1(S', a) - Q_2(S, A) \right).$$

- (van Hasselt, 2010): у Q-обучения основное уравнение апдейта

$$Q(S_t, A_t) := Q(S_t, A_t) + \alpha_t \left(R_{t+1} + \gamma \max_a Q_t(S_{t+1}, a) - Q(S_t, A_t) \right)$$

- И по умолчанию Q-обучение использует максимум из имеющихся оценок $Q_t(S_{t+1}, a)$.
- Т.е. получается, что мы оцениваем максимум из нескольких ожиданий, $\max_i \mathbb{E}[X_i]$, через максимум из оценок каждого из них.
- А на самом деле это будет несмещённая оценка для ожидания максимума $\mathbb{E}[\max_i X_i]$!

DOUBLE Q-LEARNING

- Поэтому ван Хассельт предложил использовать Double Q-learning: давайте возьмём два разных независимых набора сэмплов для X_i , и будем использовать один из них в тот момент, когда делаем апдейт Q-обучения для другого.
- Тогда (при независимых оценках) всё будет сходиться.

Algorithm 1 Double Q-learning

```
1: Initialize  $Q^A, Q^B, s$ 
2: repeat
3:   Choose  $a$ , based on  $Q^A(s, \cdot)$  and  $Q^B(s, \cdot)$ , observe  $r, s'$ 
4:   Choose (e.g. random) either UPDATE(A) or UPDATE(B)
5:   if UPDATE(A) then
6:     Define  $a^* = \arg \max_a Q^A(s', a)$ 
7:      $Q^A(s, a) \leftarrow Q^A(s, a) + \alpha(s, a) (r + \gamma Q^B(s', a^*) - Q^A(s, a))$ 
8:   else if UPDATE(B) then
9:     Define  $b^* = \arg \max_a Q^B(s', a)$ 
10:     $Q^B(s, a) \leftarrow Q^B(s, a) + \alpha(s, a) (r + \gamma Q^A(s', b^*) - Q^B(s, a))$ 
11:   end if
12:    $s \leftarrow s'$ 
13: until end
```

ПЛАНИРОВАНИЕ

- В прошлый раз мы ввели основные понятия динамики марковских процессов принятия решений:
 - собственно динамику процесса:

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

- награды за каждый эпизод, начиная со времени t :

$$G_t = R_{t+1} + \gamma G_{t+1} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1};$$

- функцию значения для состояний и пар состояние-действие:

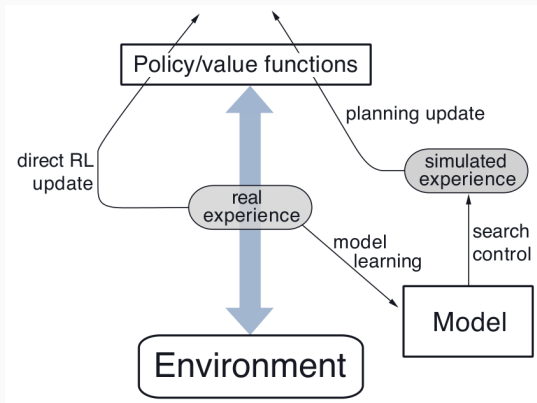
$$V_{\pi}(s) = \mathbb{E}_{\pi}[G_t \mid S_t = s] = \mathbb{E}_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s\right],$$

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

и их оптимальные варианты V_* , Q_* ;

- Мы выписывали уравнения Беллмана на V и Q и научились их решать.
- Кроме того, методами Монте-Карло мы умеем обучать модель окружающей среды.
- А TD-обучением мы можем обучать одновременно и модель, и оптимальную стратегию.
- *Планирование* (planning) — это то, как использовать модель, чтобы лучше обучать V и Q .
- Как нам это сделать?..

- Базовая идея: давайте сэмплировать новый опыт из модели окружающей среды, использовать симуляции.
- Пример – архитектура Dyna:



- В простейшем случае можно делать ещё несколько обновлений, скажем, Q-обучения, исходя из модели.

Алгоритм Dyna-Q:

- инициализировать случайно π , $Q(s, a)$, $M(s, a)$ (модель);
- повторять (цикл внутри эпизода, эпизоды тоже повторять):
 - для состояния S выбрать A по ϵ -жадной стратегии из Q ;
 - сделать действие A , пронаблюдать R и S' ;
 - $Q(S, A) := Q(S, A) + \alpha (R + \gamma \max_a Q(S', a) - Q(S, A))$;
 - добавить R, S' в $M(S, A)$;
 - повторить k раз:
 - (S, A) — случайная пара, которую уже наблюдали раньше (т.е. для неё есть $M(S, A)$);
 - породить R, S' по $M(S, A)$;
 - $Q(S, A) := Q(S, A) + \alpha (R + \gamma \max_a Q(S', a) - Q(S, A))$.
- Dyna-Q случайно обновляет k значений $Q(s, a)$ по имеющейся модели; k можно выбирать из соображений имеющихся ресурсов.

- Конечно, ещё лучше будет обновлять не случайно.

Алгоритм обхода по приоритетам (prioritized sweeping):

- инициализировать π , $Q(s, a)$, $M(s, a)$ (модель), PQ (очередь);
- повторять (цикл внутри эпизода, эпизоды тоже повторять):
 - для состояния S выбрать A по ϵ -жадной стратегии из Q ;
 - сделать A , пронаблюдать R и S' , добавить R, S' в $M(S, A)$;
 - $P := |R + \gamma \max_a Q(S', a) - Q(S, A)|$; если $P > \theta$ (порог), добавить (S, A) в PQ с приоритетом P ;
 - повторить k раз, пока $PQ \neq \emptyset$:
 - взять $(S, A) := \text{Head}(PQ)$, породить R, S' по $M(S, A)$;
 - $Q(S, A) := Q(S, A) + \alpha (R + \gamma \max_a Q(S', a) - Q(S, A))$.
 - для каждой пары (\tilde{S}, \tilde{A}) , из которой модель попадает в S :
 - * \tilde{R} – предсказанная моделью награда для $(\tilde{S}, \tilde{A}, S)$;
 - * $P := |\tilde{R} + \gamma \max_a Q(S, a) - Q(\tilde{S}, \tilde{A})|$ (приоритет);
 - * если $P > \theta$, то добавить (\tilde{S}, \tilde{A}) в PQ с приоритетом P .
- Это для детерминированного окружения, для случайного можно взвесить оценки вероятностей попасть в S .

EXPECTED VS. SAMPLE UPDATES

- Вот ещё один взгляд на то, чем мы занимались – expected updates vs. sample updates:

Value
estimated

Expected updates
(DP)

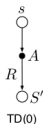
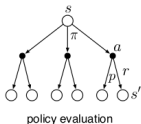
Sample updates
(one-step TD)

Value
estimated

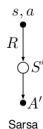
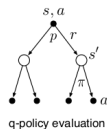
Expected updates
(DP)

Sample updates
(one-step TD)

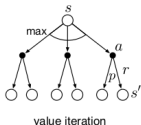
$v_{\pi}(s)$



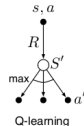
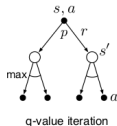
$q_{\pi}(s, a)$



$v_{*}(s)$



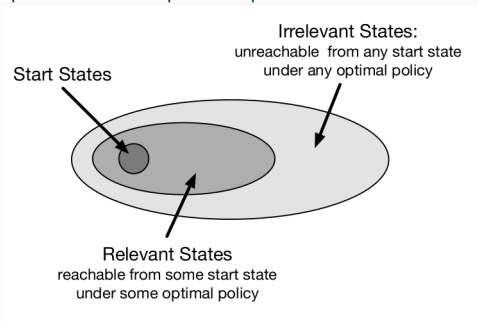
$q_{*}(s, a)$



- Expected update проходит по всем следующим состояниям.
- Sample update выбирает одно (случайно или из опыта).
- Expected, конечно, лучше, но и дороже, в целом sample даже выгоднее может получиться.

EXPECTED VS. SAMPLE UPDATES

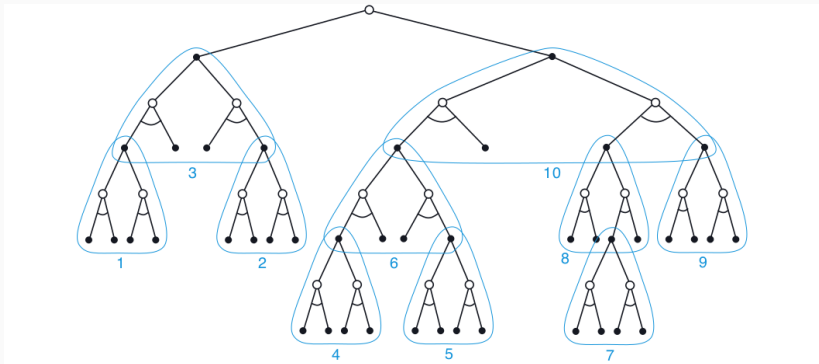
- Как лучше распределить апдейты (т.е. фактически имеющийся вычислительный ресурс)?
- Trajectory sampling: лучше сэмплировать сразу траекториями по реальной стратегии. Так мы не будем тратить время на недостижимые или очень-редко-достижимые состояния.
- Например, RTDP (real-time dynamic programming): обновляем по уравнениям Беллмана, но только состояния, которые реально встречались в траекториях.



- А можно использовать планирование прямо при выборе действия, т.е. фактически как часть стратегии (decision-time planning).
- Простейшее такое планирование мы знаем: когда мы обучали $V(s)$, а не $Q(s, a)$, потом для получения стратегии π по функции V нужно было перебрать возможные следующие состояния и выбрать лучшее.
- Как это обобщить и улучшить?..

ПЛАНИРОВАНИЕ В ТЕКУЩЕМ МОМЕНТЕ

- Правильно, это поиск на несколько «ходов» в глубину! Точнее говоря, поиск получается в ширину. :)



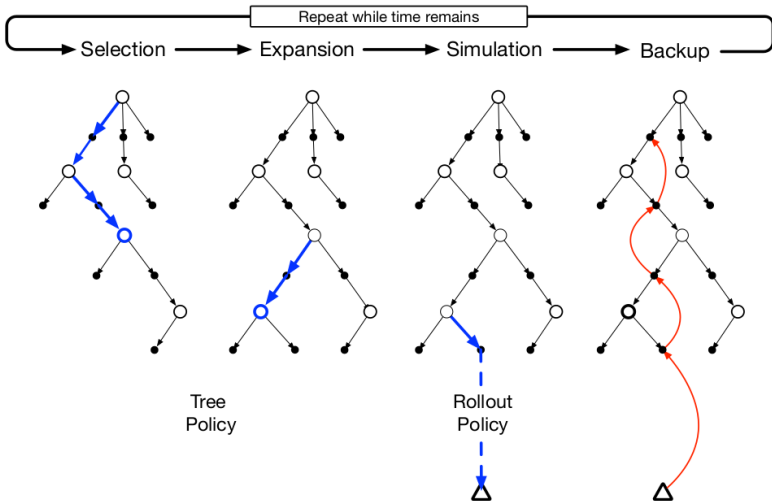
- Heuristic search: делаем поиск в ширину на несколько ходов вперёд, выбирая наилучшее действие (кстати, почему это называется «минимакс», если тут только max?)

- Другой способ реализовать такое планирование: rollouts.
- Начиная с текущего состояния, симулируем несколько траекторий по текущей стратегии, а потом оцениваем действия, усредняя результаты этих симуляций.
- Когда оценки становятся достаточно точными, можно сделать ход, а потом опять строить rollouts из следующего состояния.
- Это пробовал ещё Тезауро для нарда, и получалось очень хорошо, прямо неожиданно хорошо.
- Но ещё лучше совместить одно и другое и получить...

- ...Monte Carlo tree search (MCTS)! Это один из ключевых компонентов, например, в прогрессе алгоритмов для го между 2005 (слабый любитель) до 2015 (6 дан), а потом к AlphaGo, а потом и к AlphaZero — там везде есть MCTS.
- Мы строим дерево, в котором оцениваем листья через rollouts и выбираем, когда раскрыть лист дальше, т.е. итеративно:
 - выбираем лист дерева, действие в нём и следующее за ним состояние;
 - раскрываем всевозможные действия в этом состоянии;
 - делаем rollouts исходя из этих действий, используя их результаты для обновления оценок действий на пути к корню дерева.

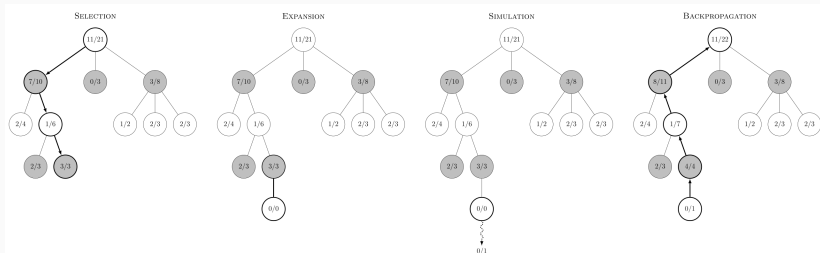
ПЛАНИРОВАНИЕ В ТЕКУЩЕМ МОМЕНТЕ

- Вот такая картинка получается:



ПЛАНИРОВАНИЕ В ТЕКУЩЕМ МОМЕНТЕ

- Выбирать листья и действия нужно по статистике побед/поражений:

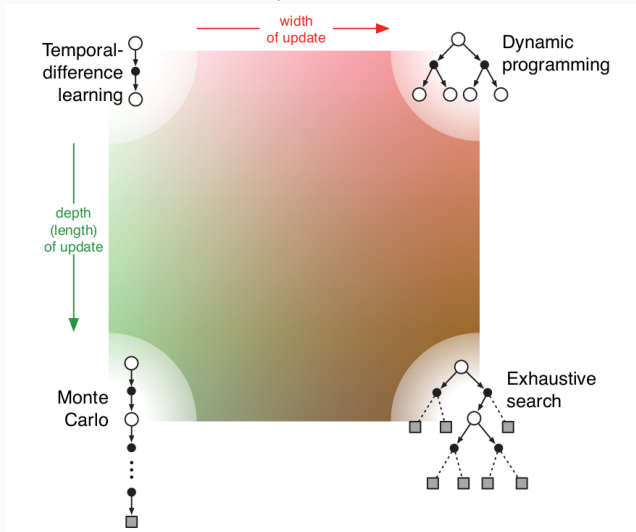


- Например, по методу UCT (upper confidence bounds applied to trees): выбираем для rollout действие, максимизирующее

$$\frac{w_i}{n_i} + c \sqrt{\frac{\ln N_i}{n_i}}.$$

НАШИ МЕТОДЫ

- И ещё один взгляд-иллюстрация:



Спасибо за внимание!