

# Reinforcement Learning

Multi Armed Bandits  
Contextual Bandits

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

# Ресар

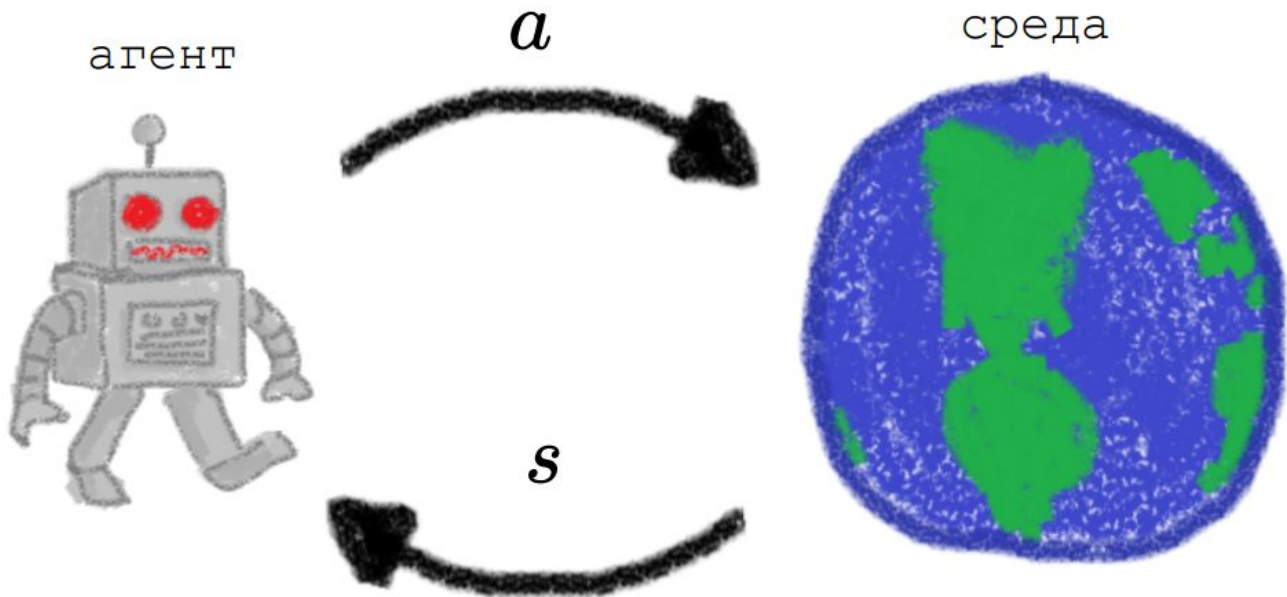
Политика агента:

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

$$a \sim \pi(a|s)$$

Цель агента:

$$\mathbb{E}_{p(\tau|\pi)}(\sum_{t=0}^T r_t) \rightarrow \max_{\pi}$$



# Ресар

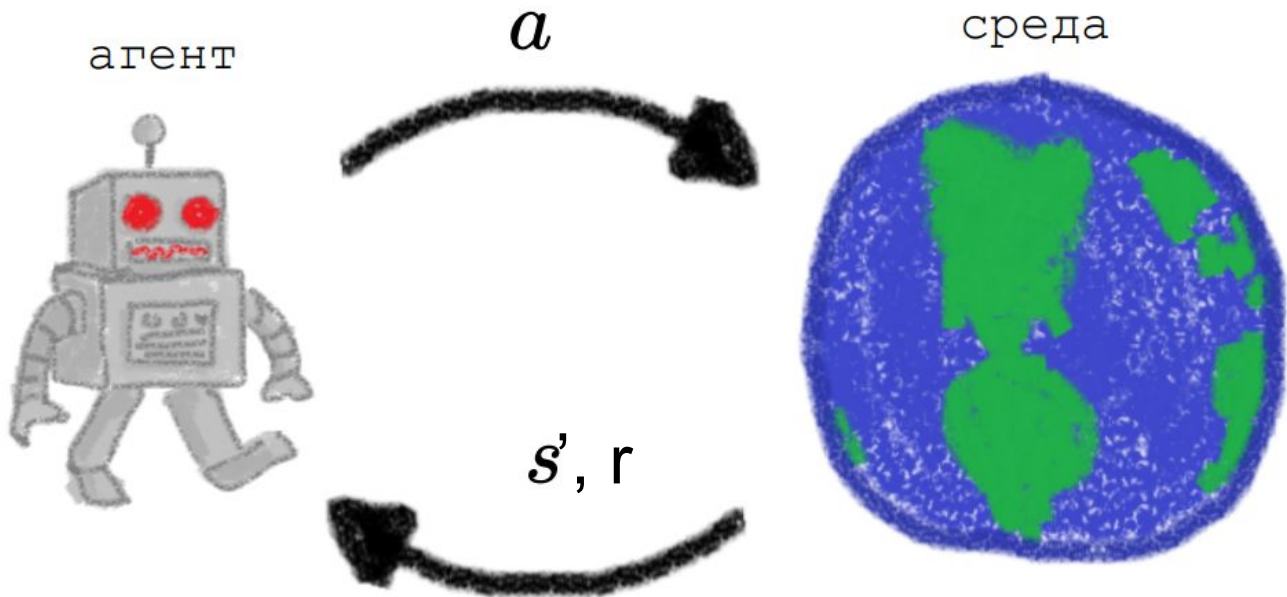
Политика агента:

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

$$a \sim \pi(a|s)$$

Цель агента:

$$\mathbb{E}_{p(\tau|\pi)}(\sum_{t=0}^T r_t) \rightarrow \max_{\pi}$$



# Рекомендация музыки

## Имеем:

- Есть много разных пользователей
- Есть много разной музыки

## Хотим:

- Рекомендовать музыку
- Пользователи продолжали пользоваться сервисом

Хорошая музыка      Плохая музыка



Пользователь

# What if

- $s'$  не зависит от  $s$  и  $a$
- все  $s$  одинаковы

Получаем многорукого бандита



# Многорукие бандиты



# Exploration vs Exploitation Dilemma

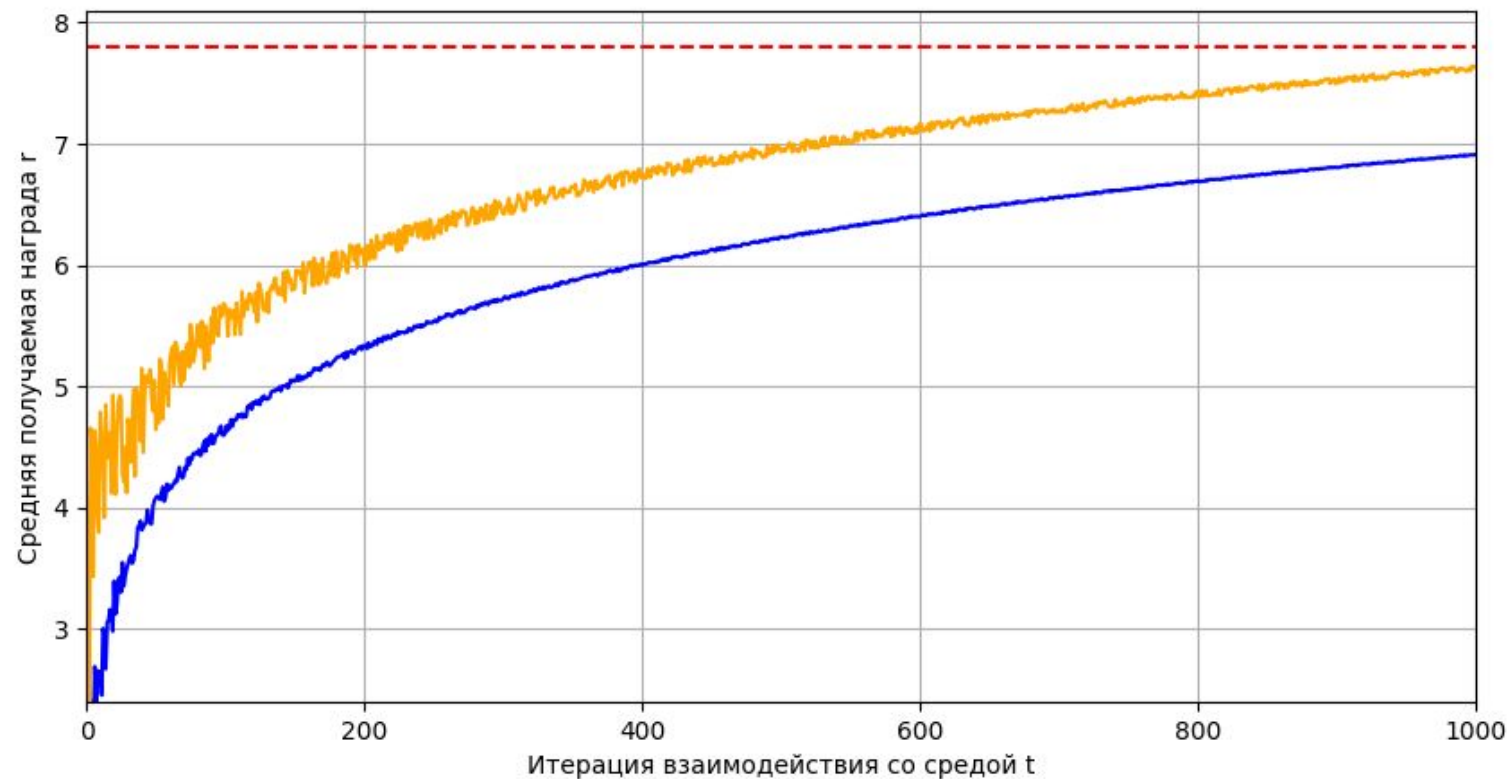
- Принятие решений в режиме онлайн предполагает фундаментальный выбор:
  - Exploitation. Использование имеющихся знаний для максимизации награды здесь и сейчас
  - Exploration. Исследовать среду. Собрать больше информации
- Лучшая долгосрочная стратегия может включать краткосрочные потери
- Чтобы добиться наилучшего результата в целом решения, необходимо собрать достаточно информации

# Какая стратегия исследования лучше

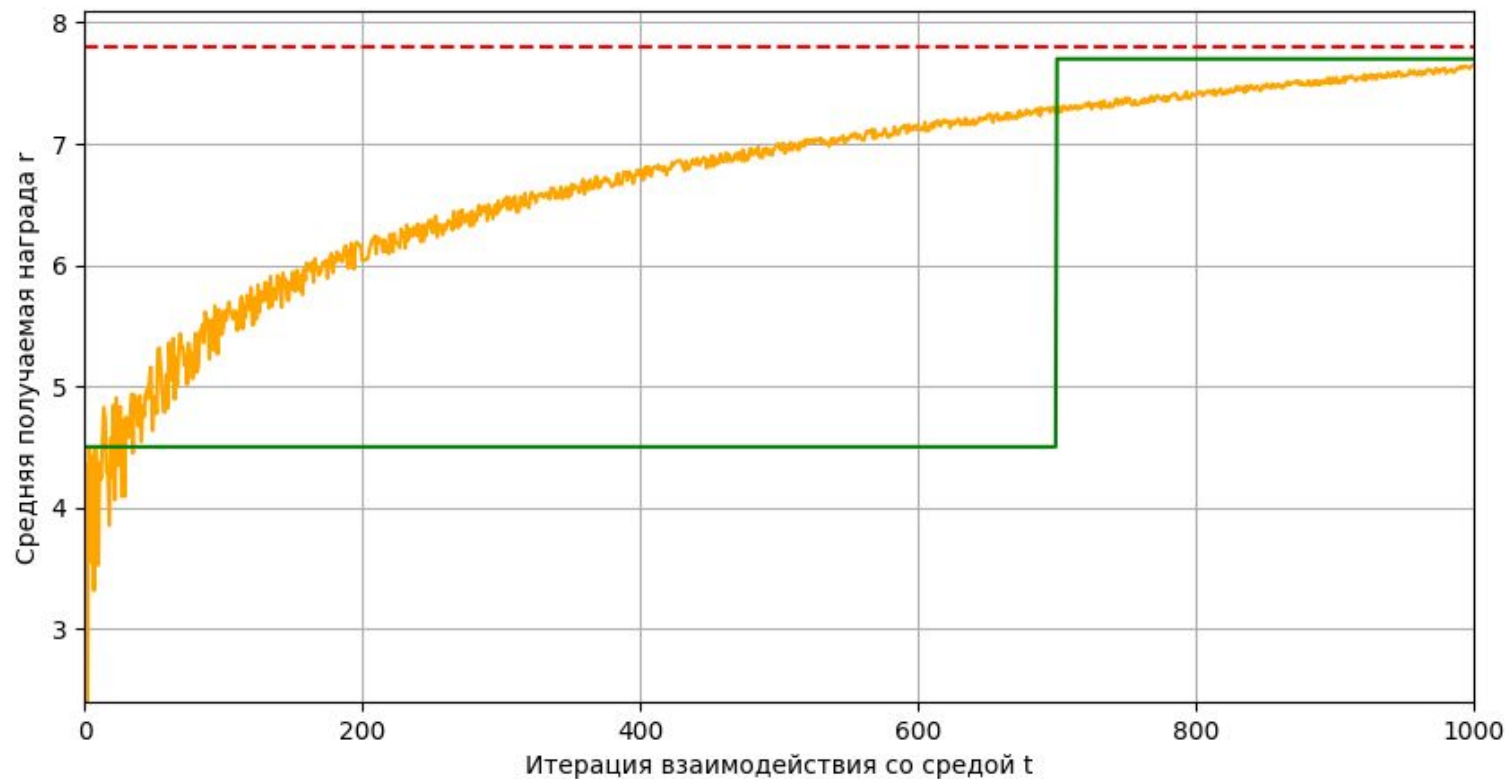
Идеи?



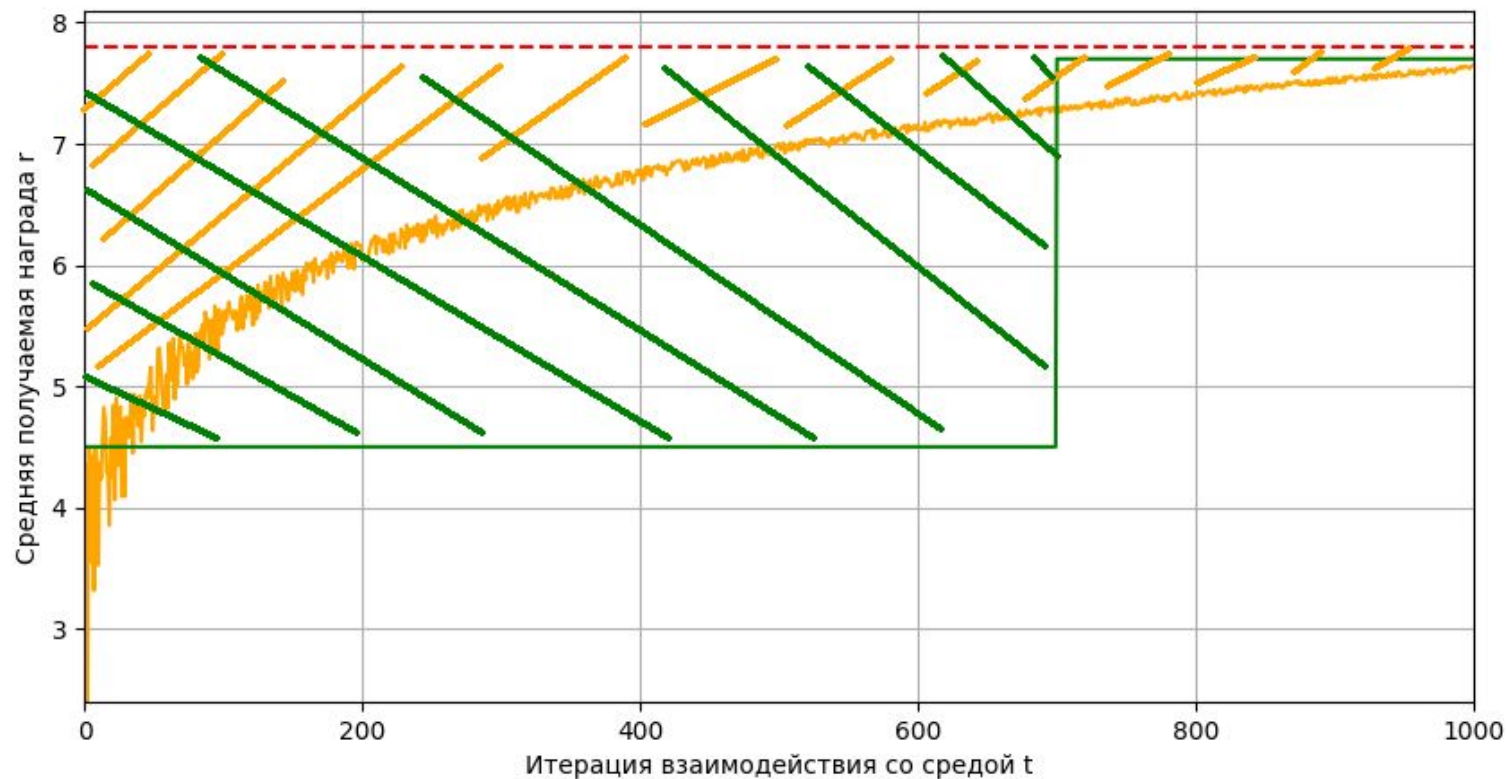
# Какая стратегия исследования лучше



# Какая стратегия исследования лучше



# Какая стратегия исследования лучше



# Постановка задачи многорукого бандита

Предположим, что эпизод бесконечен, но в среде имеется только одно состояние. Поэтому связи между последовательными действиями нет. Агент постоянно сталкивается с выбором между различными действиями.

# Постановка задачи многорукого бандита

Предположим, что эпизод бесконечен, но в среде имеется только одно состояние. Поэтому связи между последовательными действиями нет. Агент постоянно сталкивается с выбором между различными действиями.

Многорукий бандит – это пара  $\langle \mathbf{R}, \mathbf{A} \rangle$ :

- $\{\mathbf{r}_a \in \mathbf{R} \mid \mathbf{a} \in \mathbf{A}\}$  - набор распределений вознаграждений
- На каждом шаге  $t$  агент выбирает  $\mathbf{a}_t$  и получает вознаграждение  $r_a$

Задача агента максимизировать награду:

$$R = E_{\pi} \left[ \sum_{t=1}^T r(a_t) \right]$$

# Постановка задачи многорукого бандита

Исследование: найти действие, дающее наибольшую награду

Ценность действия:  $Q(a) = E[r_t \mid a_t = a]$

Оптимальная награда:  $V^* = \max_a Q(a)$

Regret:  $E_\pi[V^* - Q(a)] \geq 0$

# Постановка задачи многорукого бандита

Исследование: найти действие, дающее наибольшую награду

Ценность действия:  $Q(a) = E[r_t \mid a_t = a]$

Оптимальная награда:  $V^* = \max_a Q(a)$

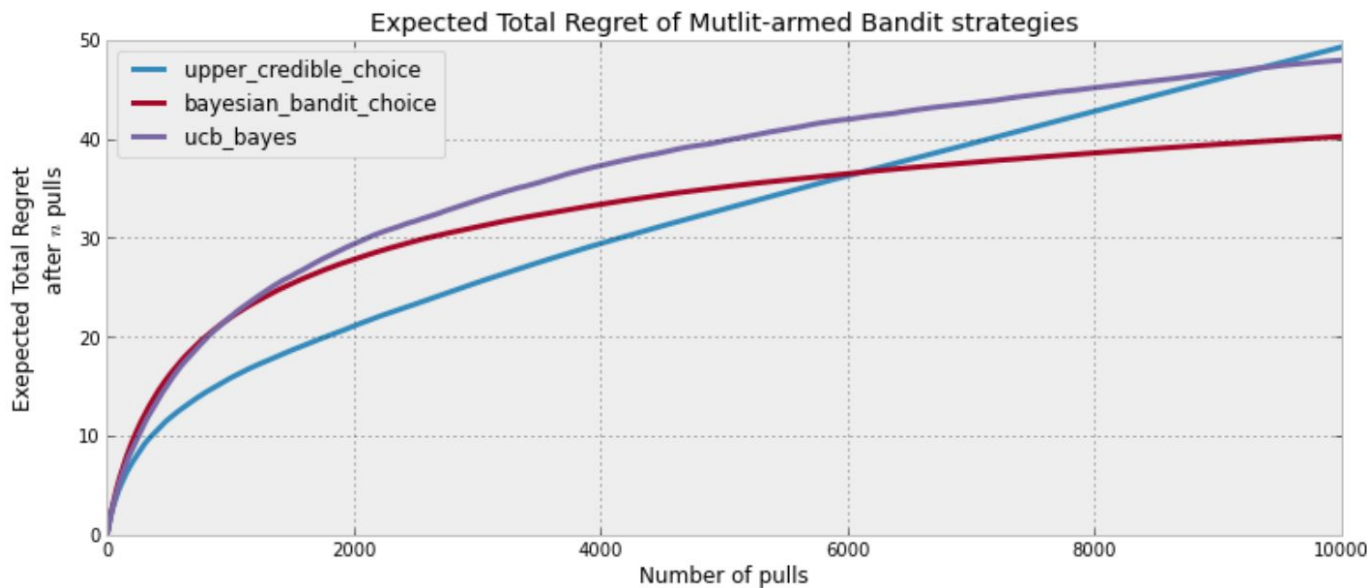
Regret:  $E_\pi[V^* - Q(a)] \geq 0$

Total Regret:  $E_\pi \sum_{t=1}^T [V^* - Q(a)] \rightarrow \min_\pi \iff E_\pi \sum_{t=1}^T [r_t] \rightarrow \max_\pi$

# Regret

Regret - недополученная награда

$$\eta = E_{\pi}[V^{\star} - Q(a)]$$





# $\epsilon$ -greedy Policy

Идея: давайте брать жадное действие с вероятностью  $(1 - \epsilon)$ , и произвольное с вероятностью  $\epsilon$

$$Q(a) = \frac{\sum r_a}{N_a}$$

$$\pi(a) = \begin{cases} \arg \max(Q(a)) & \text{if } p < 1 - \epsilon \\ \frac{\epsilon}{|A|} & \text{otherwise} \end{cases}$$

- Такой метод всегда продолжает изучение среды
- Линейный regret при  $T \rightarrow \infty$

# Оптимальность исследования

А линейный regret - это плохо или хорошо?

# Оптимальность исследования

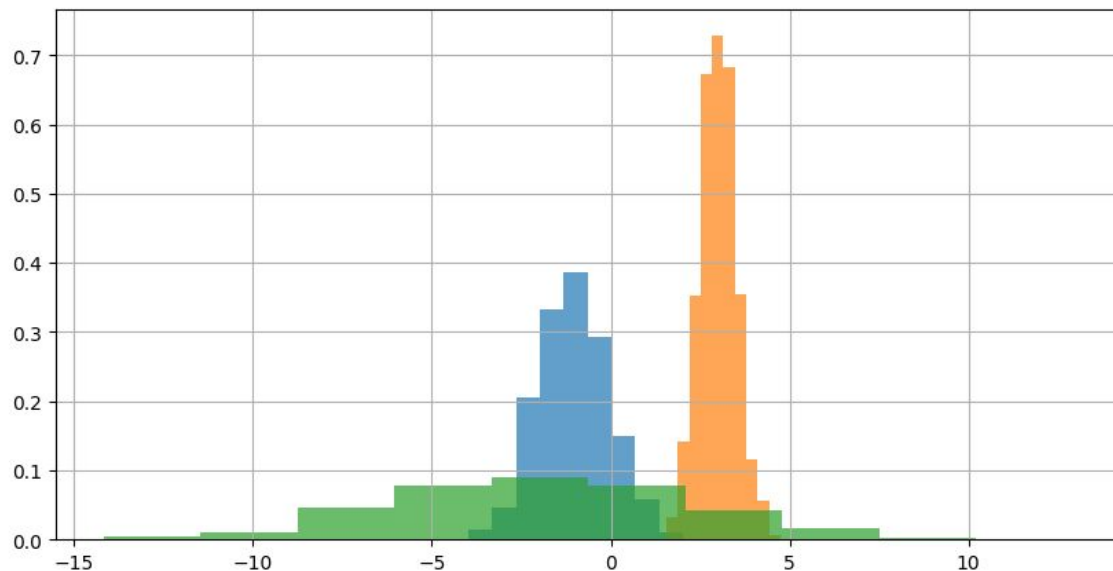
А линейный regret - это плохо или хорошо?

Теорема (Lai and Robbins):

$$\lim_{t \rightarrow \infty} \text{Regret} \geq \log t \sum_{a \mid \eta_a > 0} \frac{\eta_a}{KL(r_a \parallel r_{a^*})}$$

# $\epsilon$ -greedy Policy

Очевидно, что нет особого смысла дергать синюю ручку на шаге исследования.



# Thompson Sampling

1. Хотим пробовать действия пропорционально вероятностям, с которыми они дадут наибольшую награду
2. Зададим априорное распределение на ручках бандита
3. Дергаем ручку бандита с вероятностями:

$$\pi(a \mid h_t) = P[Q(a) > Q(a'), a \neq a' \mid h]$$

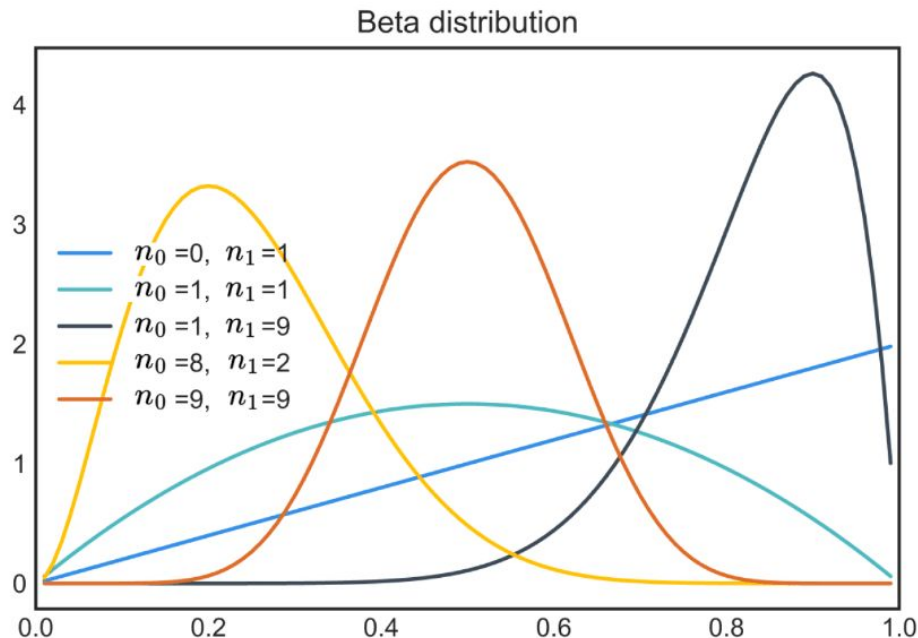
4. Считаем апостериорное распределение по формуле Байеса
5. Возвращаемся к шагу 2

У этого метода логарифмический regret

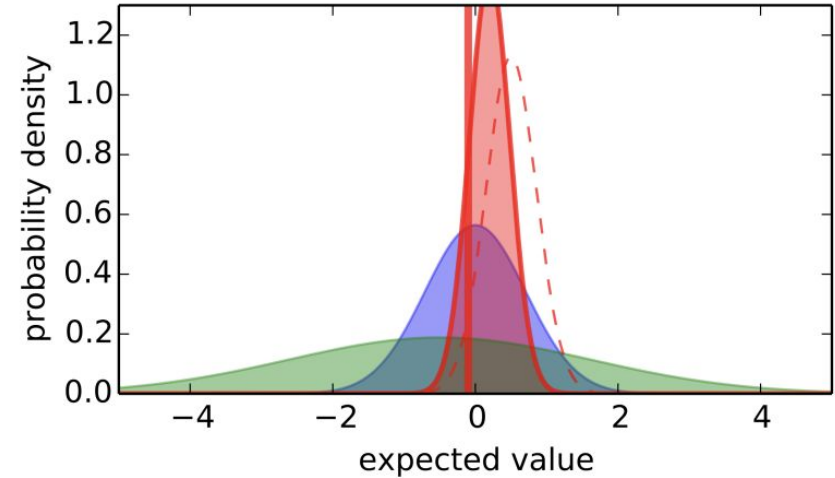
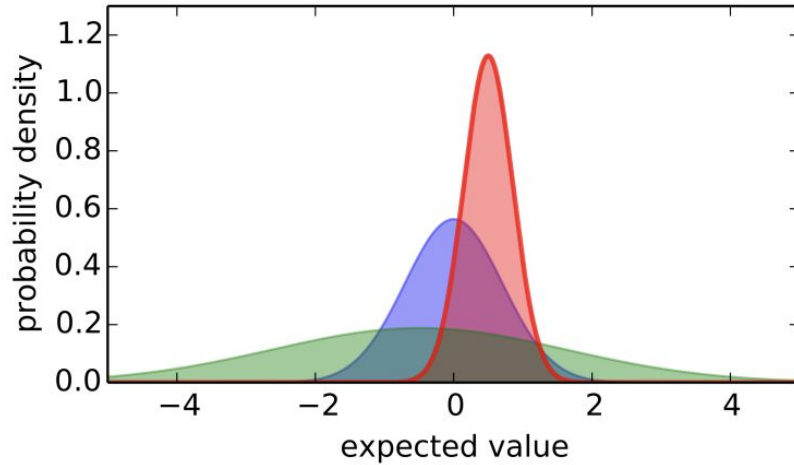
# Thompson Sampling

Случай бернуллиевского бандита:

- $r_a = \begin{cases} 1 & \text{if } p > \theta_a \\ 0 & \text{otherwise} \end{cases}$
- $E[Q(a)] = \theta_a$
- Априори:  $\theta \sim I[0, 1]$
- Апостериори:  
 $P[Q(a)] = \text{Beta}(n_1; n_0)$

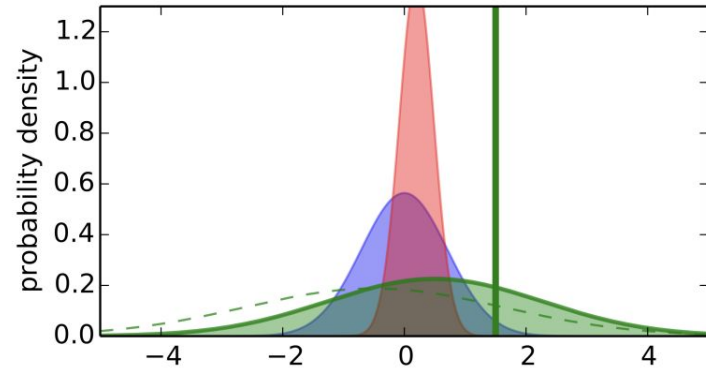
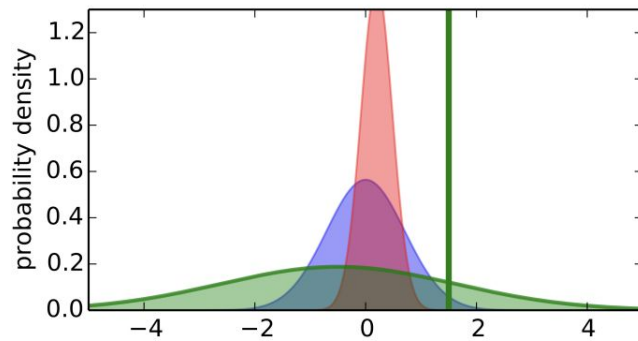
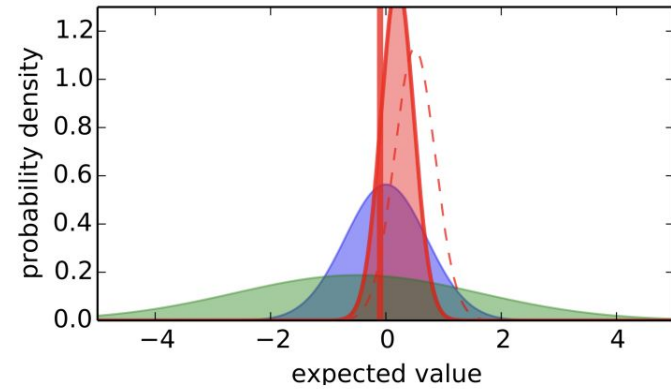
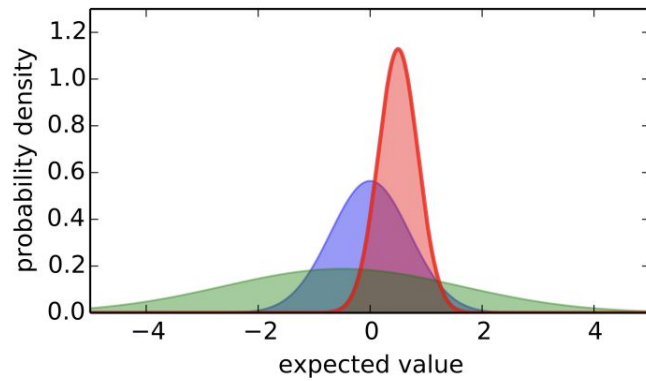


# Optimism in face of uncertainty



- Какое действие выбрать?
- Чем более неуверенны мы в отношении ценности действия, тем важнее изучить это действие.
- Оно может оказаться лучшим действием

# Optimism in face of uncertainty





# Upper Confidence Bound

- Для каждого действия **a** считаем верхнюю границу доверительного интервала **U(a)**
- Ширина интервала зависит от **N(a)** количества использований действия **a**
  - Действие почти не использовалось -> верхняя граница **U(a)** будет огромной. (Неуверенная оценка среднего)
  - Действие использовалось часто -> верхняя граница **U(a)** будет маленькой. (Уверенная оценка среднего)
- В качестве действия выбирается:

$$a_t = \arg \max_{a \in A} [Q_t(a) + U_t(a)]$$

# Как оценить верхнюю границу?

Неравенство Хёфдинга для произвольного распределения награды  $\mathbf{r}$ , но  $\mathbf{r} \in [0,1]$ :

- $P[Q_t(a) + U_t(a) < Q(a)] \leq e^{-2N_t(a)U_t(a)^2}$
- Пусть  $U_t(a) = \sqrt{\frac{-\log p}{2N_t(a)}}$  тогда  $e^{-2N_t(a)U_t(a)^2} = p$
- Возьмем  $p = \frac{1}{t}$  тогда  $U_t(a) = \sqrt{\frac{\log t}{2N_t(a)}}$

# UCB

1. Инициализируем распределения оценки среднего (например  $Q_t(a) \sim N(0,1)$ )
2. Выбираем действие  $a_t = \arg \max_{a \in A} \left[ Q_t(a) + c \sqrt{\frac{\log t}{N_t(a)}} \right]$
3. Обновляем оценку среднего  $Q_t(a)$  и счетчики  $t$  и  $N_t(a)$
4. Повторяем с шага 2

При  $c=\sqrt{2}$  у этого метода логарифмический regret ([Auer et al., 2002](#))

# Контекстные бандиты

- $s'$  не зависит от  $s$  и  $a$
- состояния  $s$  бывают разными
- награда ручки  $a$  зависит от состояния  $s$

Получаем контекстного  
многорукого бандита



# Рекомендация музыки

## Имеем:

- Есть много разных пользователей
- Есть много разной музыки

## Хотим:

- Рекомендовать музыку
- Пользователи продолжали пользоваться сервисом

Хорошая музыка

Плохая музыка



Пользователь

# Постановка задачи контекстного бандита

Исследование: найти действие, дающее наибольшую награду

Ценность действия:  $Q(s, a) = E[r_t \mid s_t = s; a_t = a]$

Оптимальная награда:  $V^*(s) = \max_a Q(s, a)$

Regret:  $E_\pi[V^*(s) - Q(s, a)] \geq 0$

Total Regret:  $E_\pi \sum_{t=1}^T [V^*(s) - Q(s, a)] \rightarrow \min_\pi \iff E_\pi \sum_{t=1}^T [r_t] \rightarrow \max_\pi$

# LinUCB

- Давайте объединим состояние **s** и действие **a** в виде вектора контекста **c**
- И предположим, что ценность действия **a** линейно зависит от **c**:

$$Q(s_t, a_t) = c_t^T \theta^* + \epsilon_t; \quad \text{где } \epsilon_t \sim N(0, 1)$$

- Минимизация regret'a принимает вид:

$$\begin{aligned} \max_{\pi} \left( E \sum_{t=1}^T r_t \right) &\Leftrightarrow \min_{\pi} E \left[ \sum_{t=1}^T \max_{a \in A} \langle c^T; \theta^* \rangle - \sum_{t=1}^T Q_t \right] \\ &\Leftrightarrow \min_{\pi} E \left[ \sum_{t=1}^T \max_{a \in A} \langle c^T - c_t^T; \theta^* \rangle \right] \end{aligned}$$

# Линейная регрессия

Предположим на шаге  $t$  имеем несколько пар  $\{\mathbf{c}; \mathbf{Q}\}$  - наш датасет для обучения линейной регрессии:

$$\hat{\theta}_t = \arg \min_{\theta \in R^d} \sum_{k=1}^{t-1} (Q_k - c_k^T \theta)^2 + \frac{\lambda}{2} \|\theta\|_2^2$$

Решение:

$$\hat{\theta}_t = V_{t-1}^{-1} \sum_{k=1}^{t-1} c_k Q_k$$

$$\text{где } V_{t-1} = \sum_{k=1}^{t-1} c_k c_k^T + \lambda I_d$$



# Optimism in face of uncertainty

- В обычном UCB мы искали верхнюю грань доверительного интервала
- В линейной модели мы ищем вектор параметров  $\theta$  в эллипсоиде  $M \in \mathbb{R}^d$ :

$$a_t = \arg \max_{a \in A} c^T \hat{\theta}_t \Rightarrow a_t = \arg \max_{a \in A} \max_{\theta \in M_t} c^T \theta$$

# Confidence Ellipsoid

Пусть 
$$\beta_t(\delta) = \lambda + \sqrt{2\log\left(\frac{1}{\delta}\right) + d\log\left(1 + \frac{t}{\lambda d}\right)}$$

тогда 
$$M_t(\delta) = \left\{ \theta \in R^d : \left\| \theta^* - \hat{\theta}_t \right\|_{V_{t-1}} \leq \beta_{t-1}(\delta) \right\}$$

задает доверительный эллипсоид для  $\theta^*$  с уровнем уверенности  $1-\delta$

# LinUCB

- На входе: вероятность  $\delta$ , размерность  $d$ , регуляризация  $\lambda$
- Инициализация:  $b = 0_{R^d}$ ;  $V = \lambda I$ ;  $\hat{\theta} = 0_{R^d}$
- for  $t > 0$ :
  - Получаем состояние  $\mathbf{s}$
  - Вычисляем:  $\beta_{t-1}(\delta) = \lambda + \sqrt{2\log\left(\frac{1}{\delta}\right) + d\log\left(1 + \frac{t-1}{\lambda d}\right)}$
  - for  $a$  in  $A$ :
    - Собираем вектор  $\mathbf{c}$
    - Вычисляем:  $UCB(a) = \mathbf{c}^T \hat{\theta} + \beta_{t-1} \sqrt{\mathbf{c}^T V^{-1} \mathbf{c}}$
  - $a_t = \arg \max UCB(a)$
  - Обновляем:  $V = V + \mathbf{c}_t \mathbf{c}_t^T$

$$b = b + r_t \mathbf{c}_t$$

$$\hat{\theta} = V^{-1} b$$

Regret порядка  $d\sqrt{2}$

# Полезные ссылки

- 1) [D-LinUCB](#) - бандит для динамической среды
- 2) [Neural Contextual Bandits](#) - линейные бандиты с нейронками для исправления нелинейности
- 3) [Лекция по LinUCB](#)