

Многорукие бандиты

Селиханович Даниил, МФТИ ФУПМ, 572 группа

4 апреля 2018 г.

Постановка задачи

Многорукие бандиты. Имеется два «одноруких бандита» (так называют игровые автоматы с ручкой, дергая за которую получаем случайный выигрыш). Вероятность выиграть 1 франк на первом автомате $p_1 > 0$ (с вероятностью $1 - p_1$ выигрыш равен 0), а на втором — $p_2 > 0$. Обе вероятности неизвестны. Игрок может в любом порядке раз дергать за ручки «одноруких бандитов». Стратегией игрока является выбор ручки на каждом шаге, в зависимости от результатов всех предыдущих шагов, так чтобы суммарный выигрыш был бы максимальным. Приведите стратегию игрока. Предложите свое обобщение задачи.

Содержание

1	Динамическое программирование.	2
1.1	Испытания Бернулли.	2
1.2	Алгоритм поиска ожидаемой максимальной прибыли.	3
2	Заключение	14
3	Дополнение. Метод зеркального спуска в задачах стохастической выпуклой оптимизации.	15
3.1	Постановка задачи в терминах оптимизации.	15
3.2	Оценка регрета.	16
3.3	Сглаживание задачи.	17
3.4	Корректность оценок для сглаживаемой задачи.	18

1 Динамическое программирование.

1.1 Испытания Бернулли.

Постановка задачи: у нас имеется последовательность подбрасываний монетки, о «честности» которой нам ничего не известно. Требуется определить, какова вероятность получить орёл или решку в следующем подбрасывании.

Обозначим вероятность выпадения орла через q , случайные величины, соответствующие i -му подбрасыванию монетки, через X_i , а их значения через x_i : $x_i = 1$, если выпал орёл и $x_i = 0$ для решки. По формуле Байеса плотность вероятности $f(q) = P(q \mid X_1 = x_1, \dots, X_n = x_n)$ вычисляется так:

$$P(q \mid X_1 = x_1, \dots, X_n = x_n) = \frac{P(q) P(X_1 = x_1, \dots, X_n = x_n \mid q)}{P(X_1 = x_1, \dots, X_n = x_n)}. \quad (1)$$

Мы будем считать, что априорное распределение «честности» монетки равномерное, то есть пока мы не видели ни одного результата, мы вообще ничего не можем сказать о вероятности выпадения орла и считаем все такие вероятности одинаково возможными: $P(q) = 1$ при $0 \leq q \leq 1$ и $P(q) = 0$ иначе.

$$P(X_1 = x_1, \dots, X_n = x_n \mid q) = \prod_{i=1}^n q^{x_i} (1-q)^{1-x_i} = q^s (1-q)^{n-s}. \quad (2)$$

Здесь s - число выпадений орлов в эксперименте из n подбрасываний.

Дискретная формула полной вероятности в непрерывном случае, как в данной задаче, пример вид:

$$P(X_1 = x_1, \dots, X_n = x_n) = \int_0^1 q^s (1-q)^{n-s} dq = B(s+1, n-s+1) = \frac{\Gamma(s+1)\Gamma(n-s+1)}{\Gamma(n+2)} = \frac{s!(n-s)!}{(n+1)!}$$

Получили, что

$$P(q \mid X_1 = x_1, \dots, X_n = x_n) = \frac{(n+1)!}{s!(n-s)!} q^s (1-q)^{n-s} \quad (3)$$

Вычислим матожидание $E[q]$ выпадения орла в следующем подбрасывании монетки при условиях эксперимента $X_1 = x_1, \dots, X_n = x_n$:

$$\begin{aligned} E[q] &= \int_0^1 q f(q) dq = \int_0^1 q \frac{(n+1)!}{s!(n-s)!} q^s (1-q)^{n-s} dq = \frac{(n+1)!}{s!(n-s)!} \int_0^1 q^{s+1} (1-q)^{n-s} dq = \\ &= \frac{(n+1)!}{s!(n-s)!} B(s+2, n-s+1) = \frac{(n+1)!}{s!(n-s)!} \frac{\Gamma(s+2)\Gamma(n-s+1)}{\Gamma(n+3)} = \\ &= \frac{(n+1)!}{s!(n-s)!} \frac{(s+1)!(n-s)!}{(n+2)!} = \frac{s+1}{n+2}. \quad (4) \end{aligned}$$

1.2 Алгоритм поиска ожидаемой максимальной прибыли.

Пусть мы имеем k разных «одноруких бандитов». Будем обозначать состояние игрока после n_1 использований 1-ой ручки, из которых только w_1 были успешными, \dots, n_k использований k -ой ручки, из которых только w_k были успешными, через $S = \{(n_1, w_1), \dots, (n_k, w_k)\}$, а $V^*(S)$ через ожидаемый максимальный выигрыш для состояния с данной историей, если всего можно делать h выборов.

Тогда задача сводится к вычислению $V^*((0, 0), \dots, (0, 0))$ для данного числа выборов h . Базой рекурсии, которую мы хотим составить, будет служить очевидное равенство:

$$V^*(n_1, w_1), \dots, (n_k, w_k) = 0 \text{ для всех таких } n_1, \dots, n_k : \sum_{i=1}^k n_i = h \quad (5)$$

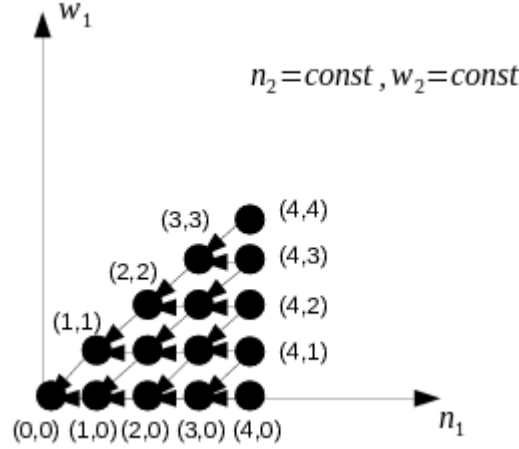
. Если обозначить через ρ_i ожидаемую вероятность того, что после истории испытаний на ручке i (n_i, w_i) игрок в состоянии S получит выигрыш после выбора ручки i , то можем выписать следующую рекурренту:

$$\begin{aligned} (6) \quad & V^*(n_1, w_1), \dots, (n_k, w_k) = \\ & = \max_{1 \leq i \leq k} \left(\rho_i \left\{ 1 + V^*(n_1, w_1), \dots, (n_i+1, w_i+1), \dots \right\} + (1 - \rho_i) V^*(n_1, w_1), \dots, (n_i+1, w_i), \dots \right) \end{aligned}$$

Согласно (4) можем считать, что $\rho_i = \frac{w_i+1}{n_i+2}$ в соответствии с условием задачи, так как вероятности выигрыша p_i могут быть произвольными ненулевыми.

Замечаем, что данная рекуррента с учётом начальных условий позволяет определить $V^*((0, 0), \dots, (0, 0))$ для данного числа выборов h , не проводя сам эксперимент. Тем самым возникает наивный алгоритм действий игрока, желающего получить ожидаемый максимальный выигрыш при игре с однорукими бандитами.

Рис. 1: Порядок заполнения таблицы



Алгоритм

1) предварительно вычислить таблицу $V^*((n_1, w_1), \dots, (n_k, w_k))$ для всех $n_1, \dots, n_k, w_1, \dots, w_k :$

$0 \leq w_i \leq n_i, \sum_{i=1}^k n_i \leq h$. Сложность заполнения таблицы - $O(k|V^*|)$, где $|V^*| =$

$O\left((h+1)^k \sum_{i=0}^h \binom{i+k-1}{i}\right) = O\left(h^k \binom{h+k}{k}\right)$ - размер заполняемой таблицы. В общем случае получаем экспоненциальное время заполнения от h , но при небольшом числе ручек k получаем полином от h - так, для двух ручек получим сложность заполнения таблицы $O(h^4)$.

2) на каждом шаге с данным фиксированным состоянием опыта $S = \{(n_1, w_1), \dots, (n_k, w_k)\}$ выбирать действие i , для которого достигается

$$\max_{1 \leq i \leq k} \left(\rho_i \left\{ 1 + V^*((n_1, w_1), \dots, (n_i+1, w_i+1), \dots) \right\} + (1-\rho_i) V^*((n_1, w_1), \dots, (n_i+1, w_i), \dots) \right).$$

3) переходить в новое состояние S_{new} .

Анализ алгоритма: при значениях k , сравнимых с h , время заполнения экспоненциально от времени игры, так что можем использовать только в тех случаях, когда нужно планировать игру не слишком далеко. По сути - полный перебор всех возможных событий. Зато таблица $V^*((n_1, w_1), \dots, (n_k, w_k))$ может использоваться многократно.

Рис. 2: Несимметричные вероятности

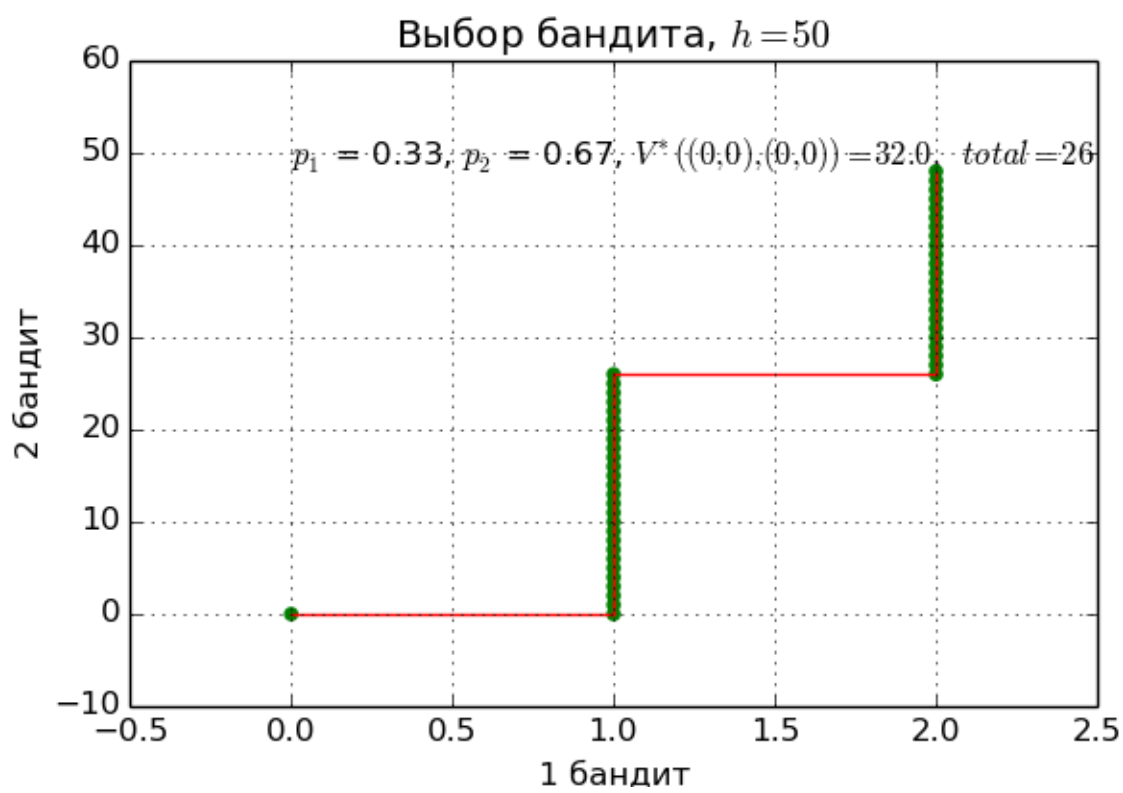


Рис. 3: Несимметричные вероятности

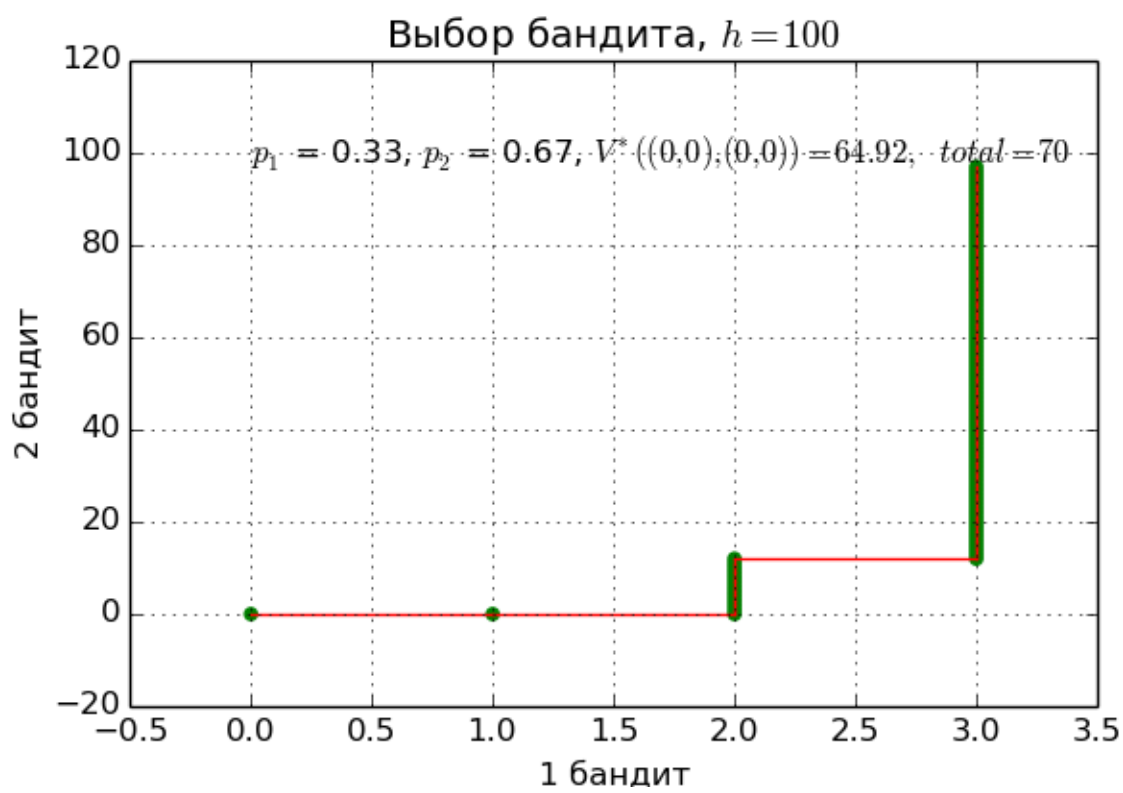


Рис. 4: Несимметричные вероятности

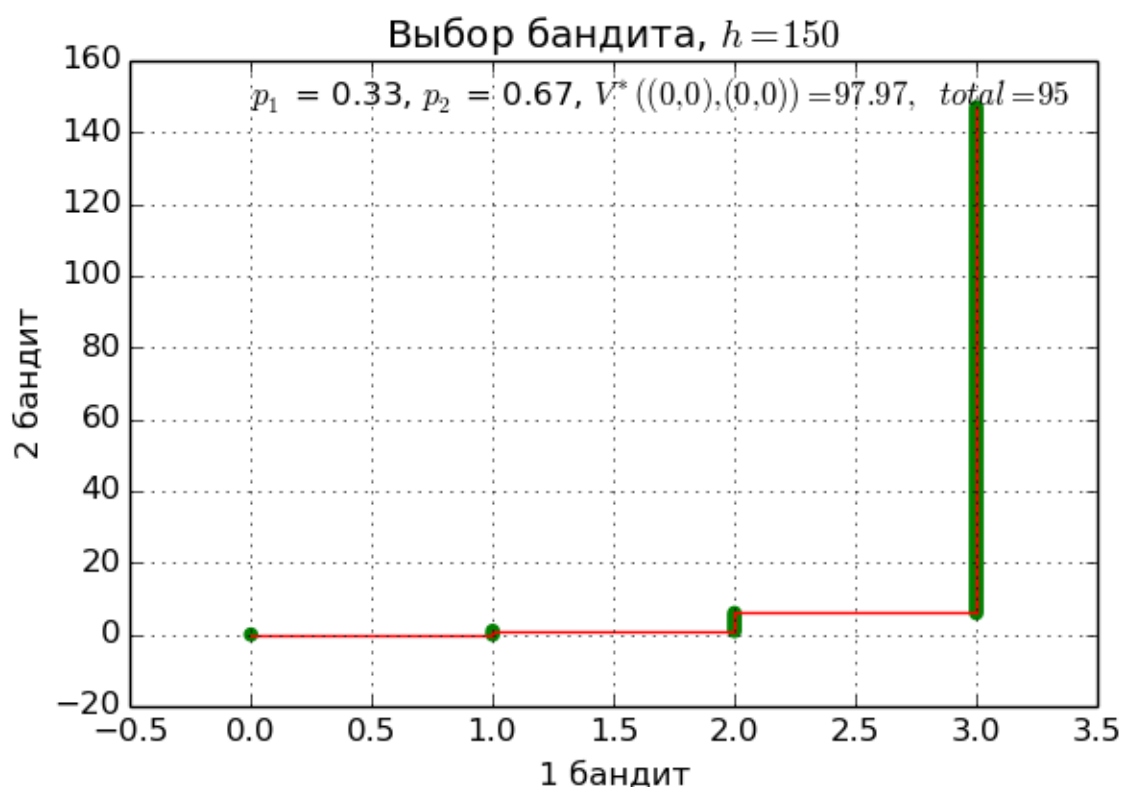


Рис. 5: Несимметричные вероятности

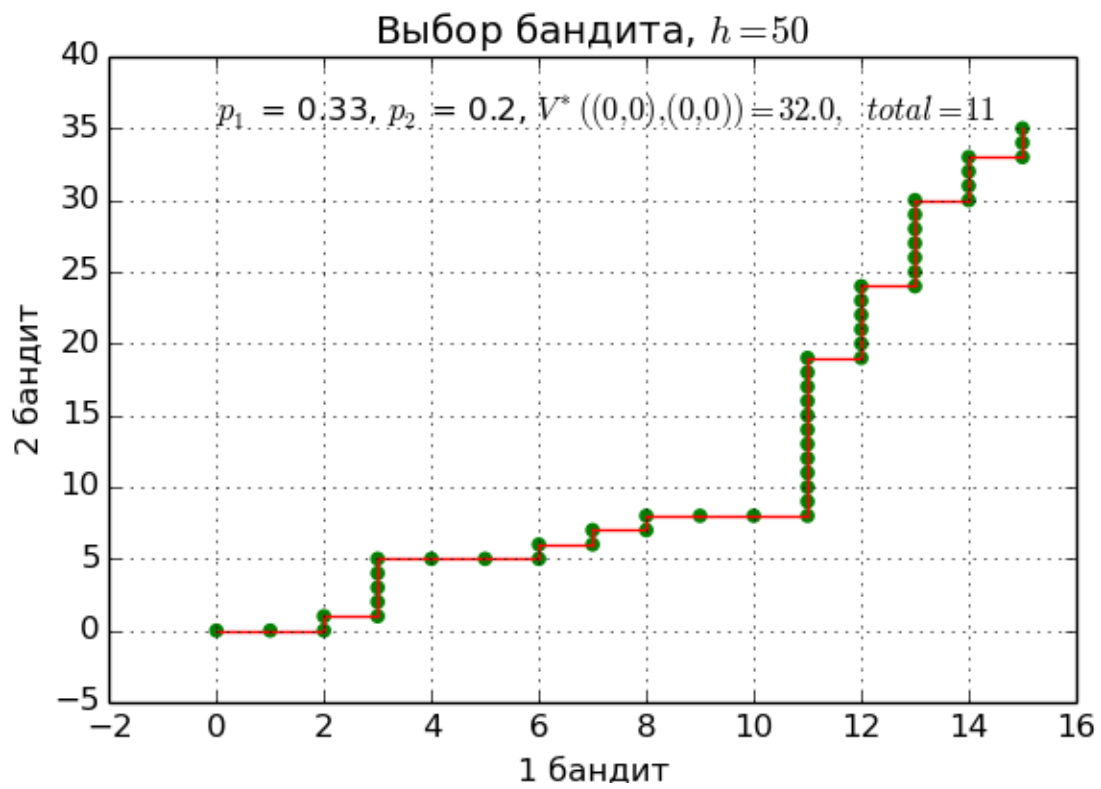


Рис. 6: Несимметричные вероятности

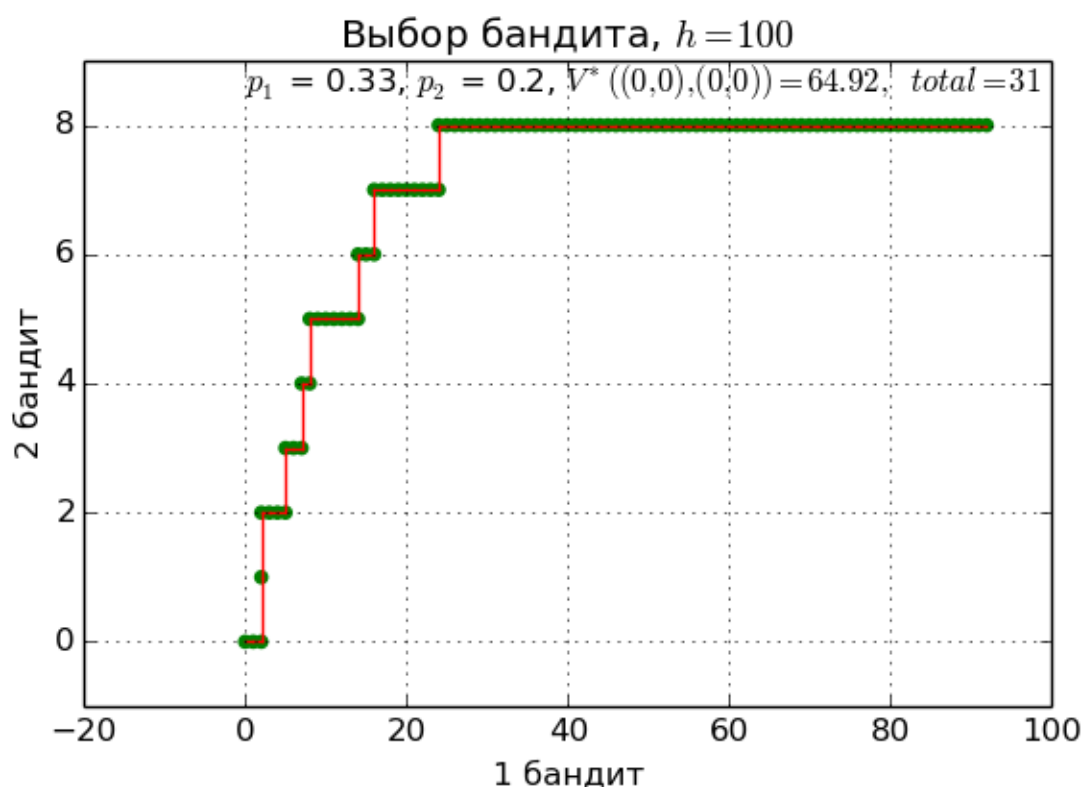


Рис. 7: Несимметричные вероятности

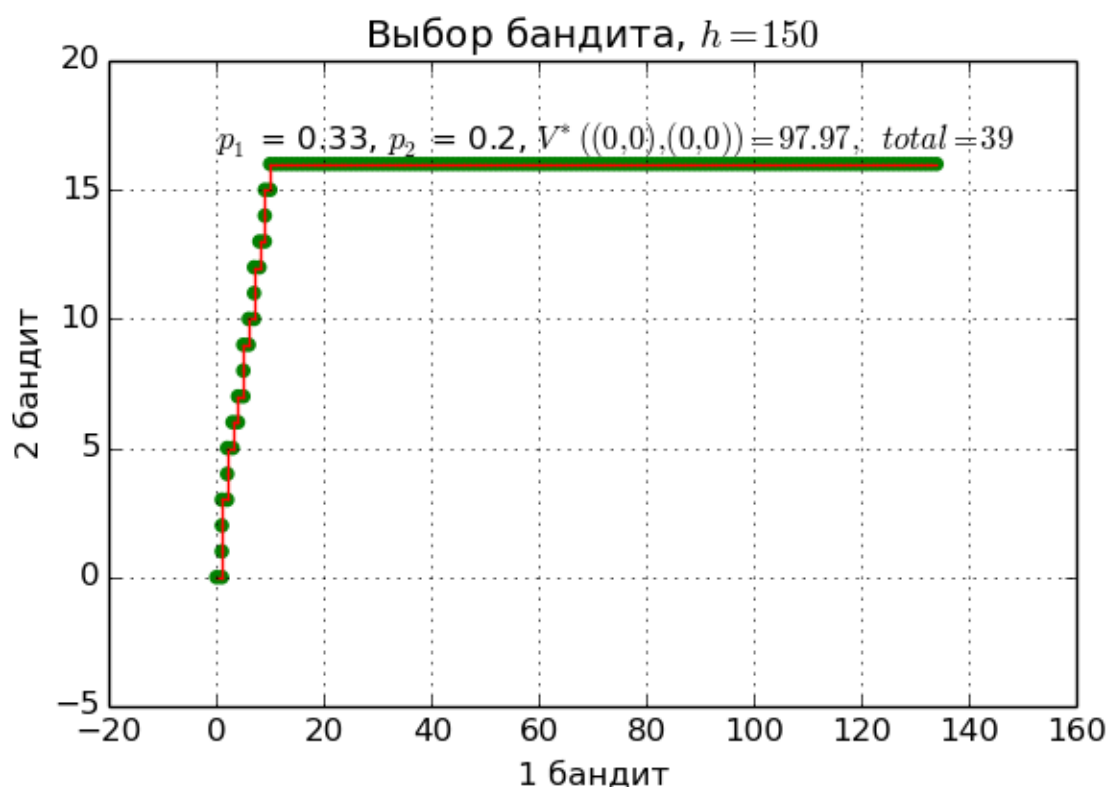


Рис. 8: Симметричные вероятности

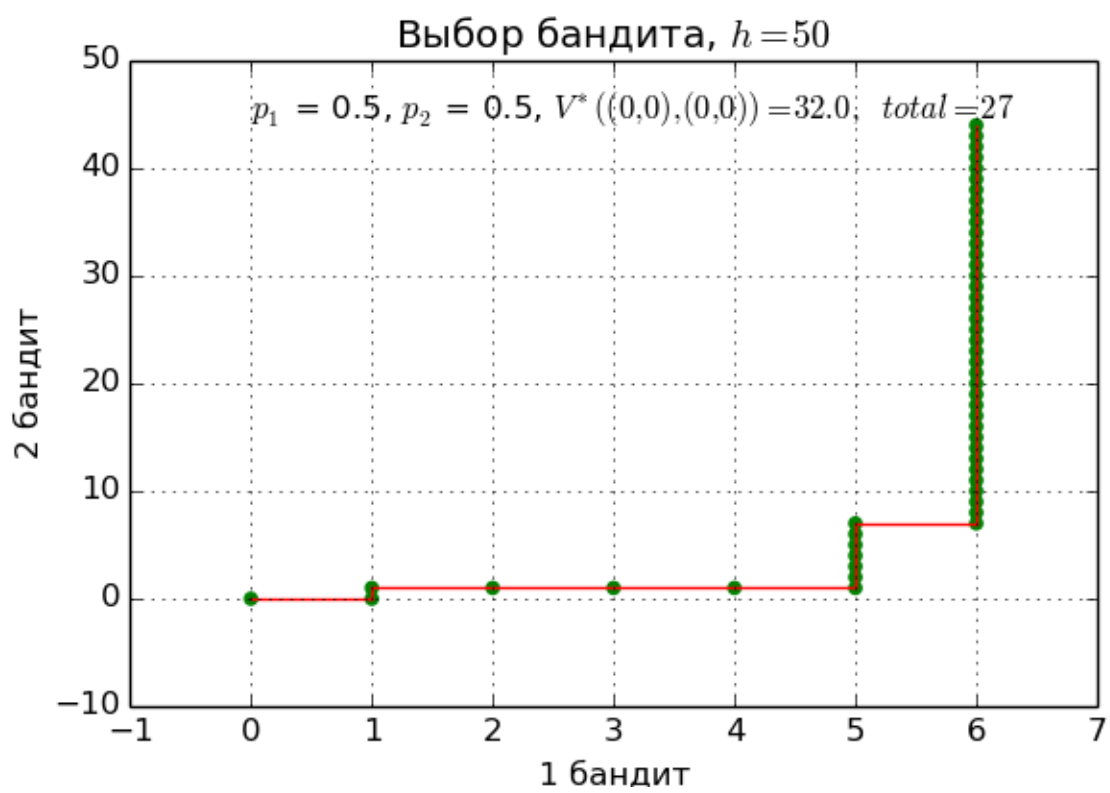


Рис. 9: Симметричные вероятности

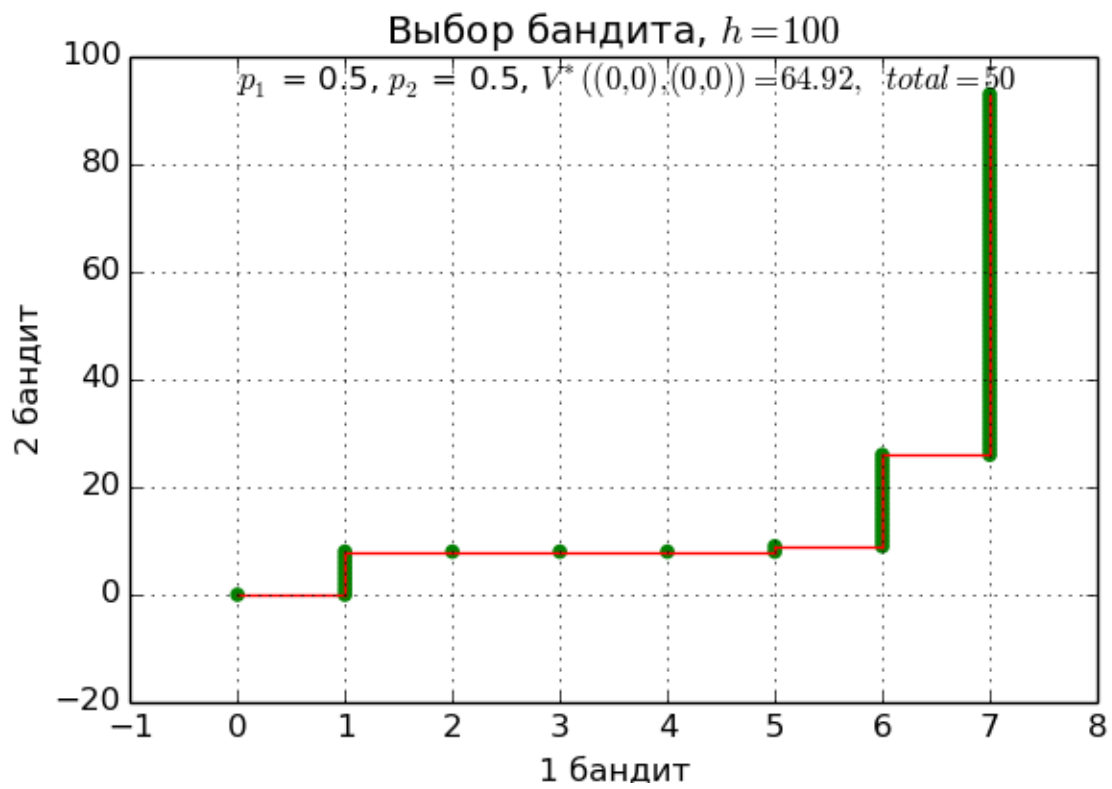
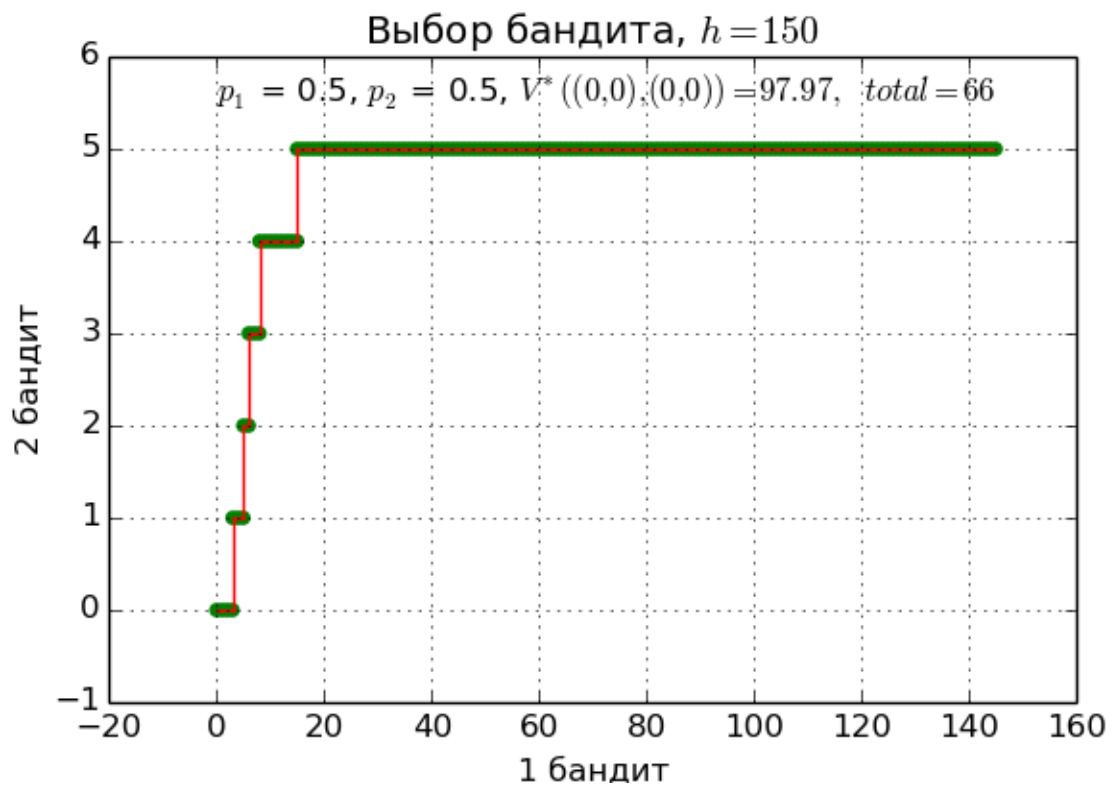


Рис. 10: Симметричные вероятности



2 Заключение

Выражаю благодарность Гасникову Александру Владимировичу и Сувориковой Александре за помощь в исследовании задачи.

Список литературы

- [1] *Баяндина А. С., Гасников А.В., Гулиев Ф.Ш., Лагуновская А.А.* Безградиентные двухточечные методы решения задач стохастической негладкой выпуклой оптимизации при наличии малых шумов не случайной природы: <https://arxiv.org/ftp/arxiv/papers/1701/1701.03821.pdf>.
- [2] *Гасников А. В., Нестеров Ю.Е., Спокойный В.Г.* Об эффективности одного метода рандомизации зеркального спуска в задачах онлайн оптимизации: <https://arxiv.org/pdf/1410.3118.pdf>.
- [3] *Гасников А. В., Крымова Е.А., Лагуновская А.А., Усманова И.Н., Федоренко Ф.А.* Стохастическая онлайн оптимизация. Одноточечные и двухточечные нелинейные многорукие бандиты. Выпуклый и сильно выпуклый случаи: <https://arxiv.org/ftp/arxiv/papers/1509/1509.01679.pdf>.
- [4] *Гасников А. В., Крымова Е.А., Лагуновская А.А., Усманова И.Н., Федоренко Ф.А.* Безградиентные прокс-методы с неточным оракулом для негладких задач выпуклой стохастической оптимизации на симплексе: <https://arxiv.org/ftp/arxiv/papers/1412/1412.3890.pdf>.
- [5] *Аникин А. С., Гасников А.В., Дзуреченский П.Е., Тюрин А.И., Чернов А.В.* Двойственные подходы к задачам минимизации сильно выпуклых функционалов простой структуры при аффинных ограничениях: <https://arxiv.org/ftp/arxiv/papers/1602/1602.01686.pdf>.
- [6] Лекция Зорича В.А про концентрацию меры:
[http : //www.mathnet.ru/php/seminars.phtml?option_lang = rus&presentid = 14604](http://www.mathnet.ru/php/seminars.phtml?option_lang=rus&presentid=14604).
- [7] *Николенко С. И., Тулупьев А.Л.* Самообучающиеся системы. Российская академия наук. Санкт-Петербургский институт информатики и автоматизации РАН. Москва, Издательство МЦНМО, 2009.
- [8] *Bubeck S., Cesa-Bianchi N.* Regret analysis of stochastic and nonstochastic multi-armed bandit problems, *Foundation and Trends in Machine Learning*, 5:1 (2012) 1–122; *Lugosi G., Cesa-Bianchi N.* Prediction, learning and games. New York: Cambridge University Press, 2006.
- [9] *Lev Bogolubsky, Pavel Dvurechensky, Alexander Gasnikov, Gleb Gusev, Yurii Nesterov, Andrey Raigorodskii, Aleksey Tikhonov, Maksim Zhukovskii* Learning Supervised PageRank with Gradient-Free Optimization Methods: <https://arxiv.org/pdf/1411.4282.pdf>.

Спасибо за внимание!

3 Дополнение. Метод зеркального спуска в задачах стохастической выпуклой оптимизации.

3.1 Постановка задачи в терминах оптимизации.

Рассмотрим смежную задачу о нелинейных бандитах с шумом, в которой выбор ручки приносит нам случайные потери на шаге k ручки $i(k)$ ($1 \leq i \leq n$), равные $f_{i(k)}^k$, зависящие от номера шага, номера ручки и от того, какой стратегии мы придерживались до шага k включительно. Наша стратегия на шаге k описывается вектором распределения вероятностей $x^k \in S_n(1) = \{x \geq 0 \mid \sum_{i=1}^n x_i = 1\}$, согласно которому мы независимо ни от чего выбираем ручку, которую будем дёргать. Все, чем мы располагаем на шаге k , это вектором

$$\left((x^1, i(1), f_{i(1)}^1); \dots; (x^{k-1}, i(k-1), f_{i(k-1)}^{k-1}) \right).$$

Мы считаем, что потери на k -ом шаге f^k зависят от x^k (но не от результата разыгрывания из распределения x^k), зависят от (x^1, \dots, x^{k-1}) и от результатов соответствующих разыгрываний, а также зависят от (f^1, \dots, f^{k-1}) . Таким образом, целью является организовать процедуру выбора ручек, чтобы ожидаемые суммарные потери были бы минимальны.

Предполагаем, что функция $f(x^k)$ выпуклая на $Q = S_n(1)$ и равна потерям после выбора на k -ом шаге вектора вероятностей x^k . Однако игрок не может наблюдать значение самой функции $f(x)$, а видит лишь зашумленную реализацию значений этой функции:

$$\tilde{f}(x^k, \xi^k) = f(x^k, \xi^k) + \delta(x^k, \xi^k), \quad (7)$$

$$\mathbb{E}_{\xi^k}(f(x^k, \xi^k)) = f(x^k), |\delta(x^k, \xi^k)| < \delta. \quad (8)$$

Предполагается, что $\{\xi^k\}_k$ - независимые одинаково распределённые случайные величины (их реализации одинаковы для всех игроков), функция $f(x, \xi)$ как функция x в μ_0 -окрестности выпуклого множества Q является:

- **(9)** выпуклой;
- удовлетворяющей условию **(10)** $|f(x, \xi) - f(y, \xi)| \leq M \|y - x\|_2 \forall x, y, \xi$.

Успешность стратегии игрока, сыгравшего $N \gg 1$ раз, характеризуется величиной:

$$\text{Regret}_f(\{x^k\}_{k=0}^{N-1}) = \mathbb{E} \left\{ \frac{1}{N} \sum_{k=0}^{N-1} f(x^k) \right\} - \min_{x \in Q} f(x) \quad (11)$$

Пусть $x^k = x_*$ - оптимальная стратегия игрока и игрок начинает в состоянии x^0 . Введём прокс-функцию $d(x) \geq 0$ ($d(x^0) = 0$), которая предполагается сильно выпуклой относительно евклидовой нормы, и определим $R^2 = V(x_*, x^0)$, где прокс-расстояние (расстояние Брэгмана) определяется формулой:

$$V(x, z) = d(x) - d(z) - \langle \nabla d(z), x - z \rangle. \quad (12)$$

На каждой итерации можно один раз обратиться к оракулу за зашумленным значением стохастического градиента $\nabla_x \tilde{f}(x^k, \eta^k)$. Пусть для любых $\tilde{N} \leq N$ выполнено (Θ^{k-1} - сигма алгебра, порождённая случайными величинами $\eta^1, \dots, \eta^{k-1}$):

$$\sup_{\{x^k = x^k(\xi^1, \dots, \xi^{k-1})\}_{k=1}^{\tilde{N}} \in Q} \mathbb{E} \left\{ \frac{1}{\tilde{N}} \sum_{k=1}^{\tilde{N}} \left\langle \mathbb{E}_{\eta^k} \left\{ \nabla_x f(x^k, \eta^k) - \nabla_x \tilde{f}(x^k, \eta^k) \mid \Theta^{k-1} \right\}, x^k - x_* \right\rangle \right\} \leq \sigma; \quad (13)$$

$$\mathbb{E}_{\eta^k} \left\{ \nabla_x f(x^k, \eta^k) \right\} = \nabla f(x^k), \mathbb{E}_{\eta^k} \left\{ \left\| \nabla_x \tilde{f}(x^k, \eta^k) \right\|_2^2 \right\} \leq \tilde{M}^2, \quad (14)$$

где $\{\eta^k\}_{k=0}^{N-1}$ - независимые одинаково распределённые случайные величины.

3.2 Оценка регрета.

Метод зеркального спуска состоит в выборе на очередном шаге:

$$(15) \quad x^{k+1} = \text{Mirr}_{x^k} \left(h \nabla_x \tilde{f}(x^k, \eta^k) \right), \text{Mirr}_{x^k}(v) = \arg \min_{x \in Q} \left\{ \langle v, x - x^k \rangle + V(x, x^k) \right\},$$

где шаг спуска h будет впоследствии выбран.

Основное свойство МЗС:

$$(16) \quad 2V(x, x^{k+1}) \leq 2V(x, x^k) + 2h \langle \nabla_x \tilde{f}(x^k, \eta^k), x - x^k \rangle + h^2 \left\| \nabla_x \tilde{f}(x^k, \eta^k) \right\|_2^2. \text{ Тогда можем получить:}$$

$$(17) \quad f(x^k) - f(x) \leq_{(9)} \langle \nabla f(x^k), x^k - x \rangle =_{(14)} \langle \mathbb{E}_{\eta^k} \left\{ \nabla_x f(x^k, \eta^k) \mid \Theta^{k-1} \right\}, x^k - x \rangle \leq_{(8)}$$

$$\leq \langle \mathbb{E}_{\eta^k} \left\{ \nabla_x f(x^k, \eta^k) \mid \Theta^{k-1} \right\} - \nabla_x \tilde{f}(x^k, \eta^k), x^k - x \rangle + \frac{1}{h} \left(V(x, x^k) - V(x, x^{k+1}) \right) + \frac{h}{2} \left\| \nabla_x \tilde{f}(x^k, \eta^k) \right\|_2^2.$$

Берём условное математическое ожидание $E_{\eta^k} \left\{ \cdot \mid \Theta^{k-1} \right\}$ от обеих частей неравенства (17) и воспользуемся неравенством (14):

$$(18) \quad f(x^k) - f(x) \leq \langle \mathbb{E}_{\eta^k} \left\{ \nabla_x f(x^k, \eta^k) - \nabla_x \tilde{f}(x^k, \eta^k) \mid \Theta^{k-1} \right\}, x^k - x \rangle + \frac{1}{h} \left(V(x, x^k) - E_{\eta^k} \left\{ V(x, x^{k+1}) \mid \Theta^{k-1} \right\} \right) + \frac{h}{2} \tilde{M}^2.$$

Суммируем последнее неравенство по $k = 0, \dots, N-1$, делим на N обе части, а затем берём полное математическое ожидание от обеих частей и полагаем $x = x_* = \min_{x \in Q} f(x)$ и если такое x^* не единственно, тогда выбираем то, которое доставляет минимум $R^2 = V^*(x_*, x^0)$:

$$(19) \quad \mathbb{E} \left\{ \frac{1}{N} \sum_{k=0}^{N-1} f(x^k) \right\} - f(x_*) \leq \mathbb{E} \left\{ \frac{1}{N} \sum_{k=0}^{N-1} \left\langle \mathbb{E}_{\eta^k} \left\{ \nabla_x f(x^k, \eta^k) - \nabla_x \tilde{f}(x^k, \eta^k) \mid \Theta^{k-1} \right\}, x^k - x_* \right\rangle \right\} + \frac{1}{hN} \left(V(x_*, x^0) - V(x_*, x^N) \right) + \frac{h}{2} \tilde{M}^2.$$

В силу того, что расстояние Брэгмана всегда неотрицательно, то с учётом (13) получаем оценку на регрет в зависимости от шага h в МЗС:

$$\text{Regret}_f\left(\{x^k\}_{k=0}^{N-1}\right) \leq \sigma + \frac{R^2}{hN} + \frac{h}{2}\widetilde{M}^2. \quad (20)$$

Здесь $R^2 = V(x_*, x^0)$. Желаем минимизировать оценку по h , значит выбираем

$$h = \sqrt{\frac{2R^2}{N\widetilde{M}^2}} \quad (21)$$

При такой оценке получаем

$$\text{Regret}_f\left(\{x^k\}_{k=0}^{N-1}\right) \leq \sigma + \sqrt{\frac{2R^2\widetilde{M}^2}{N}} \quad (22)$$

Значит, после

$$N = \frac{2R^2\widetilde{M}^2}{\varepsilon^2} \quad (23)$$

испытаний мы получим оценку

$$\text{Regret}_f\left(\{x^k\}_{k=0}^{N-1}\right) \leq \sigma + \varepsilon. \quad (24)$$

3.3 Сглаживание задачи.

Пусть $\tilde{e} \in RB_2^n(1)$ - случайный вектор, равномерно распределённый на шаре единичного радиуса в норме $\|\cdot\|_2$ в R^n . Сгладим исходную функцию $f(x, \tilde{\eta})$ с помощью локального усреднения по шару радиуса $\mu > 0$ ($\mu < \mu_0$), который будет выбран позже, а \tilde{e} и $\tilde{\eta}$ предполагаются независимыми случайными величинами:

$$f^\mu(x, \tilde{\eta}) = \mathbb{E}_{\tilde{e}}\left\{f(x + \mu\tilde{e}, \tilde{\eta})\right\} = \frac{1}{V_{RB_2^n(1)}} \int_{RB_2^n(1)} f(x + \mu t, \tilde{\eta}) dt, \quad (25)$$

$$f^\mu(x) = \mathbb{E}_{\tilde{\eta}}\left\{f^\mu(x, \tilde{\eta})\right\}. \quad (26)$$

Из условия (10) получаем:

$$f^\mu(x, \tilde{\eta}) - f(x, \tilde{\eta}) = \frac{1}{V_{RB_2^n(1)}} \int_{RB_2^n(1)} \left(f(x + \mu t, \tilde{\eta}) - f(x, \tilde{\eta})\right) dt \leq \frac{1}{V_{RB_2^n(1)}} \int_{RB_2^n(1)} M\mu\|t\|_2 dt \leq M\mu$$

с одной стороны. С другой стороны из выпуклости $f(x, \tilde{\eta})$ как функции аргумента x неравенство Йенсена даёт:

$$f^\mu(x, \tilde{\eta}) = \frac{1}{V_{RB_2^n(1)}} \int_{RB_2^n(1)} f(x + \mu t, \tilde{\eta}) dt \geq f\left(x + \int_{RB_2^n(1)} \mu t dt, \tilde{\eta}\right) = f(x, \tilde{\eta})$$

Таким образом, имеют место оценки:

$$0 \leq f^\mu(x, \tilde{\eta}) - f(x, \tilde{\eta}) \leq M\mu. \quad (27)$$

Возьмём матожидание $\mathbb{E}_{\tilde{\eta}}$ от всех частей данного неравенства:

$$0 \leq f^\mu(x) - f(x) \leq M\mu. \quad (28)$$

Тогда если предположить, что:

$$M\mu \leq \frac{\varepsilon}{2} \quad (29)$$

и

$$\text{Regret}_{f^\mu}(\{x^k\}_{k=0}^{N-1}) \leq \frac{\varepsilon}{2} \quad (30)$$

то можем получить оценку изначального регрета:

$$\begin{aligned} (31) \quad \text{Regret}_f(\{x^k\}_{k=0}^{N-1}) &= \mathbb{E}\left\{\frac{1}{N} \sum_{k=0}^{N-1} f(x^k)\right\} - \min_{x \in Q} f(x) \leq \\ &\leq \underbrace{\mathbb{E}\left\{\frac{1}{N} \sum_{k=0}^{N-1} f^\mu(x^k)\right\} - \min_{x \in Q} f^\mu(x)}_{\text{Regret}_{f^\mu}(\{x^k\}_{k=0}^{N-1}) \leq \frac{\varepsilon}{2}} + \frac{\varepsilon}{2} \leq \varepsilon. \end{aligned}$$

Итак, при выполнении (29) оценка регрета для функции f^μ с точностью $\frac{\varepsilon}{2}$ позволяет оценить регрет для функции f с точностью ε .

3.4 Корректность оценок для сглаживаемой задачи.

Введём аналог зашумлённого стохастического градиента $\nabla_x \tilde{f}(x^k, \eta^k)$ из предыдущего пункта:

$$\begin{aligned} (31) \quad \nabla_x \tilde{f}(x, \eta, e) &= \frac{\text{Vol}(RS_2^n(\mu))}{\text{Vol}(RB_2^n(\mu))} \left(\tilde{f}(x + \mu e, \eta) - \tilde{f}(x, \eta) \right) e = \\ &= \frac{n \frac{\pi^{\frac{n}{2}}}{\Gamma(\frac{n}{2}+1)} \mu^{n-1}}{\frac{\pi^{\frac{n}{2}}}{\Gamma(\frac{n}{2}+1)} \mu^n} \left(\tilde{f}(x + \mu e, \eta) - \tilde{f}(x, \eta) \right) e = \frac{n}{\mu} \left(\tilde{f}(x + \mu e, \eta) - \tilde{f}(x, \eta) \right) e, \end{aligned}$$

где $e \in RS_2^n(1)$, т. е. случайный вектор e равномерно распределён на сфере радиуса 1 в евклидовой норме. Считаем, что разыгрывание e происходит независимо ни от чего. Аналогично можно определить незашумленную оценку стохастического градиента $\nabla_x f(x, \eta)$:

$$\nabla_x f(x, \eta, e) = \frac{n}{\mu} \left(f(x + \mu e, \eta) - f(x, \eta) \right) e. \quad (32)$$

Основное свойство такого градиента состоит в справедливости первой оценки в (14) и оно следует из векторного варианта формулы Стокса:

$$\mathbb{E}_e \left\{ \nabla_x f(x, \tilde{\eta}, e) \right\} = \nabla_x f^\mu(x, \tilde{\eta}) \quad (33)$$

Тогда оцениваем

$$\begin{aligned}
& \mathbb{E} \left\{ \frac{1}{N} \sum_{k=1}^N \left\langle \nabla_x f(x^k, \eta^k) - \nabla_x \tilde{f}(x^k, \eta^k), x^k - x_* \right\rangle \right\} = \\
& = \mathbb{E} \left\{ \frac{1}{N} \sum_{k=1}^N \left\langle \frac{n}{\mu} \left(f(x^k + \mu e_k, \eta^k) - f(x^k, \eta^k) \right) e_k - \frac{n}{\mu} \left(\tilde{f}(x^k + \mu e_k, \eta^k) - \tilde{f}(x^k, \eta^k) \right) e_k, x^k - x_* \right\rangle \right\} = \\
& = \mathbb{E} \left\{ \frac{1}{N} \sum_{k=1}^N \left\langle \frac{n}{\mu} \left(f(x^k + \mu e_k, \eta^k) - \tilde{f}(x^k + \mu e_k, \eta^k) \right) e_k + \frac{n}{\mu} \left(\tilde{f}(x^k, \eta^k) - f(x^k, \eta^k) \right) e_k, x^k - x_* \right\rangle \right\}.
\end{aligned}$$

Можем воспользоваться оценкой типа:

$$\langle x + y, z \rangle \leq \left| \langle x + y, z \rangle \right| = \left| \langle x, z \rangle + \langle y, z \rangle \right| \leq \left| \langle x, z \rangle \right| + \left| \langle y, z \rangle \right|$$

и в силу (7) получить:

$$\mathbb{E} \left\{ \frac{1}{N} \sum_{k=1}^N \left\langle \nabla_x f(x^k, \eta^k) - \nabla_x \tilde{f}(x^k, \eta^k), x^k - x_* \right\rangle \right\} \leq \frac{2\delta n}{\mu} \mathbb{E} \left\{ \frac{1}{N} \sum_{k=1}^N \left| \langle e_k, x^k - x_* \rangle \right| \right\}. \quad (34)$$

Если записать неравенство **(19)** с учётом неотрицательности регрета для произвольного k , то можно получить оценку при выбранном шаге h из (21):

$$\mathbb{E} \left(\|x^k - x_*\|_2^2 \right) \leq V(x_*, x^0) + \frac{Nh^2 \widetilde{M}^2}{2} = R^2 + R^2 = 2R^2 \quad (35)$$

Оказывается, явление концентрации равномерной меры на многомерной сфере позволяет из этой оценки и независимости $e_k \in RS_2^n(1)$ от $x^k - x_*$ сделать вывод:

$$\mathbb{E} \left\{ \frac{1}{N} \sum_{k=1}^N \left| \langle e_k, x^k - x_* \rangle \right| \right\} \leq \frac{2R}{\sqrt{n}},$$

что позволяет выбрать

$$\sigma = \frac{4\delta R \sqrt{n}}{\mu} \quad (36)$$

с учётом (34).