

Онлайн ценообразование с помощью модели структурированных многоруких бандитов

Гунаев Руслан Гуламович

Московский физико-технический институт
Факультет управления и прикладной математики
Кафедра интеллектуальных систем

Научный руководитель Ю. Дорн

15 декабря 2020 г.

Актуальность исследования

В данный момент в теории многоруких бандитов не существует оптимального и точного алгоритма, способного восстанавливать точное распределение цен.

Цель

Предложить онлайн алгоритм многоруких бандитов для нахождения распределений цен на товары-заменители.

Определение

Товары-заменители – товары, находящиеся в одной категории, но разных цен, которые мы и хотим точно определять для увеличения прибыли магазина.

Стохастические многорукие бандиты

В стохастических многоруких бандитах учащийся должен выбрать ручку из конечного множества X в каждый из T раундов. Дергая ручку x в раунде t он получает положительный выигрыш $R_t(x) = r \in \mathcal{R}$ с вероятностью $P(x, r)$. Задача учащегося получить максимальный ожидаемый выигрыш.

Структурированные многорукие бандиты

Распределения выигрышей заранее неизвестны. Известно только то, что $P = (P(r, x))_{r \in \mathcal{R}, x \in X} \in \mathcal{P}$, которое является замкнутым выпуклым множеством с непустой внутренностью. Это знание позволяет ввести некоторую структуру на распределение, а именно $P(r = 0, x) \geq P(r = 0, x') \forall x \geq x'$.



Gershkov Alex, Moldovanu Benny, et al.

Revenue maximizing mechanisms with strategic customers and unknown, markovian demand.

In *European Semantic Web Conference*. Springer, 2016.



Weichao Mao, Zhenzhe Zheng, and Fan Wu.

Online pricing for revenue maximization with unknown time discounting valuations.

Twenty-Seventh International Joint Conference on Artificial Intelligence (IJCAI-18), pages 440–476.



Parys Van and Golrezaei Negin.

Optimal learning for structured bandits.

arXiv:2011.07738, 2018.

Определение

Имеется n товаров заменителей, p_1, \dots, p_n – распределения товаров, которые мы хотим узнать. Множество покупателей C . Каждый покупатель готов платить за каждый из товаров максимальную цену $\tilde{q}_1(c), \dots, \tilde{q}_n(c)$. Также введем аппроксимацию стоимостей товаров $q_1(t), \dots, q_n(t)$ в каждый момент времени t .

Задача

Найти функционал, увеличивая который, сможем получить максимальную прибыль

$$\text{Reg} = \Phi(\mathbf{q}, \tilde{\mathbf{q}}, T, C) \rightarrow \min_{\mathbf{q}}$$

$$\mathbb{E} \text{Reg} = T \sum_{c=1}^{|C|} v_1(c) - \mathbb{E} \sum_{t=1}^T \sum_{c=1}^{|C|} q_{b_c}(t), \quad b_c \in \{1, \dots, n\}.$$

$v_i(c)$ – пороговая оценка i -го товара для c покупателя.

Также пусть $\forall c \in C \rightarrow v_1(c) \geq \dots \geq v_n(c)$.

$b_c(t) = \operatorname{argmax}_k \{v_k(c) - q_k(t) : q_k(t) \leq v_k(c)\}$

Предложение 2

Определение

Политику π назовем равномерно хорошей, если для любого $\alpha > 0$ и для любого распределения \mathbf{q} выполнено

$$\limsup_{T \rightarrow \infty} \mathbb{E} \text{Reg}_{\pi}(T, \mathbf{q}) / T^{\alpha} = 0$$

$$\text{Rew}(\mathbf{q}) = \sum_{t=1}^T \sum_{c=1}^{|C|} q_{b_c}(t), \quad \text{Rew}^* = \arg \max_{\mathbf{q}} \text{Rew}(\mathbf{q}).$$

Нижняя граница регрета

Пусть π – произвольная равномерно хорошая политика.
Тогда

$$\liminf_{T \rightarrow \infty} \mathbb{E} \text{Reg}_{\pi}(T, \mathbf{q}) / \log(T) \geq C(\mathbf{q})$$

$$C(\mathbf{q}) = \inf_{\mathbf{q}} \Delta(\mathbf{q}), \quad \Delta(\mathbf{q}) = \text{Rew}^* - \text{Rew}(\mathbf{q})$$

Главная теорема

Для любого $0 < \varepsilon < 1/n$ существует оптимальная политика π такая, что

$$\limsup_{T \rightarrow \infty} \mathbb{E} [\text{Reg}_\pi(T, \mathbf{q})] / \log(T) \leq (1 + \varepsilon)C(\mathbf{q}) + O(\varepsilon).$$

DUSA

Предлагается найти двойственную задачу и воспользоваться алгоритмом двойной структуры, который необходимо подстроить под нашу задачу.