

# Эффективное обучение кликовых моделей

Шугаепов Ильнур

научный руководитель: PhD И.Е. Марков

СПб АУ НОЦНТ РАН

13 июня 2017 г.

Кликовые модели описывают поведение пользователя на SERP<sup>1</sup> с точки зрения последовательности наблюдаемых и скрытых случайных величин.

## Определение

*Кликовая модель  $M$ , обученная на множестве сессий  $S$  алгоритмом  $A$ , есть пара  $\langle Q, I \rangle$ , где  $Q$  - набор запросо-зависимых параметров,  $I$  - набор запросо-независимых параметров.*

---

<sup>1</sup>Search Engine Result Page

## Position-Based Model (PBM):

- Случайные величины
  - $E_r = \mathcal{I}$ (пользователь изучил сниппет ранга  $r$ )
  - $A_{uq} = \mathcal{I}$ (пользователь привлечен  $u$ , при запросе  $q$ )
- Параметры

$$P(E_r = 1) = \gamma_r$$

$$P(A_{uq} = 1) = \alpha_{uq}$$

*Обучение* есть процесс вычисления параметров модели по *click log*.

Методы обучения кликовых моделей:

- Метод максимального правдоподобия (MLE) (все случайные величины наблюдаемы)
- EM-алгоритм (есть скрытые случайные переменные)

Кликовые модели, в описании которых участвуют скрытые случайные переменные, позволяют описать более сложные сценарии поведения пользователя на SERP.

## Недостатки:

- Весь набор параметров модели и *click log* зачастую не помещается целиком в оперативной памяти современного компьютера;
- Итеративное вычисление параметров является довольно медленным.

**Цель:** Сделать существующий подход к EM-алгоритму обучения кликовых моделей более эффективным

**Задачи:**

- 1 Распределить параметры кликовой модели и данные о поисковых сессиях между узлами кластера таким образом, чтобы обеспечить корректность работы EM-алгоритма
- 2 Сделать обучение в рамках узла кластера более эффективным
- 3 Уменьшить насколько возможно коммуникационную сложность распределенного алгоритма обучения

## Определение

Пусть  $\mathcal{Q} = \{q_1, q_2, \dots, q_l\}$  - есть множество уникальных поисковых запросов из множества  $\mathcal{S}$ . Тогда распределенная кликовая модель есть вектор  $\mathcal{M} = (M_{q_1}, M_{q_2}, \dots, M_{q_l})$ , где  $M_q$  - кликовая модель обученная на  $S_q$  алгоритмом  $\mathcal{A}$ , и  $S_q$  - множество поисковых сессий инициированных запросом  $q$ . И  $\forall i, j (i \neq j \Rightarrow Q_{q_i} \cap Q_{q_j} = \emptyset)$ .



Подготовка данных:

Представление  $\mathcal{S}$  в виде  $(\mathcal{S}_{q_1}, \mathcal{S}_{q_2}, \dots, \mathcal{S}_{q_l})$ .

Общая схема  $(t + 1)$ -ой итерации алгоритма  $\mathcal{A}_{MR}$ :

Fit (Map)

$$(M_q^{(t)}, \mathcal{S}_q) \xrightarrow{\mathcal{A}} M_q'^{(t+1)}$$

Extract (Map)

$$M_q'^{(t+1)} \rightarrow I_q'^{(t+1)}$$

Merge (Reduce)

$$(I_{q_1}'^{(t+1)}, I_{q_2}'^{(t+1)}, \dots, I_{q_l}'^{(t+1)}) \rightarrow I^{(t+1)}$$

Update (Map)

$$(M_q'^{(t+1)}, I^{(t+1)}) \rightarrow \langle Q_q^{(t+1)}, I^{(t+1)} \rangle =: M_q^{(t+1)}$$

Обучение в рамках узла кластера можно сделать более эффективным если представить EM-алгоритм в виде операций над матрицами.

## Предположение

$$\forall q \in \mathcal{Q} \forall r = \overline{1, n} \forall s, s' \in \mathcal{S}_q \left( u_r^{(s)} = u_r^{(s')} \right)$$

Предположения не верно для примерно 15% уникальных запросов в Yandex Relevance Prediction Challenge <sup>2</sup>.

---

<sup>2</sup>[https://academy.yandex.ru/events/data\\_analysis/relpred2011/](https://academy.yandex.ru/events/data_analysis/relpred2011/)

- Получено представление EM-алгоритм для РВМ в виде операций над матрицами как в частном так и в общем случаях.
- EM-алгоритм для частного случая обладает меньшей пространственной и временной сложностью.

Агрегация по рангу:

Пусть  $\mathcal{S} := \{s_1, s_2, \dots, s_T\}$ , тогда построим множества

$$\mathcal{R}_r := \left\{ \left( q_1, u_r^{(s_1)}, c_r^{(s_1)} \right), \left( q_2, u_r^{(s_2)}, c_r^{(s_2)} \right), \dots, \left( q_T, u_r^{(s_T)}, c_r^{(s_T)} \right) \right\},$$

$r = 1, 2, \dots, n,$

где  $n$  - максимальный ранг.

Для PBM, UBM<sup>3</sup> верно, что:

## Теорема

- 1 Модель  $\mathcal{M} = (M_1, M_2, \dots, M_n)$  - есть допустимая распределенная модель;
- 2 Построение параметров модели  $M_r^{(t+1)}$  алгоритмом  $\mathcal{A}$  зависит только от  $(M_r^{(t)}, \mathcal{R}_r)$ .

Процесс обучения может быть представлен в виде MapReduce фаз.

---

<sup>3</sup>A. Chuklin, I.Markov and M. de Rijke, Click Models for Web Search, 2015

- Реализована Hadoop MapReduce версия EM-алгоритма обучения для PBM,UBM,CCM,DBN<sup>4</sup>  
[https://bitbucket.org/Sh\\_Nur/pyclick2/branch/hadoop](https://bitbucket.org/Sh_Nur/pyclick2/branch/hadoop)
- Получено векторно-матричное представление формул для PBM
- Получен альтернативный подход к распределенному обучению PBM, UBM, не требующий обмена данными между узлами кластера
- Получено векторно-матричное представление EM-алгоритма для PBM,UBM в рамках альтернативного подхода

---

<sup>4</sup>A. Chuklin, I.Markov and M. de Rijke, Click Models for Web Search, 2015

Спасибо за внимание!