

## **ИСПОЛЬЗОВАНИЕ КОЛЛАБОРАТИВНОЙ ФИЛЬТРАЦИИ В РЕКОМЕНДАТЕЛЬНЫХ СИСТЕМАХ ЭЛЕКТРОННОЙ КОММЕРЦИИ**

**Д.Н. АЛЕКСЕЕВ<sup>1</sup>, А.И. БАЛЕНКО<sup>2\*</sup>**

<sup>1</sup> *магістрант кафедри вычислительной техники и программирования, НТУ «ХПИ», Харьков, УКРАИНА*

<sup>2</sup> *доцент кафедри вычислительной техники и программирования, канд. техн. наук, НТУ «ХПИ», Харьков, УКРАИНА*

\* *email: alexibalenko@gmail.com*

На данный момент развитие систем электронной коммерции, оказания услуг и распространения контента достигло той точки, когда просмотр всего объема доступной информации становится практически невыполнимой задачей для рядового пользователя. В связи с этим возникает необходимость предоставить каждому отдельному пользователю рекомендации, относительно тех товаров, услуг или контента, в которых он может быть наиболее заинтересован.

Целью данной работы является построение рекомендательной системы-модуля интернет-сервиса, позволяющего подобрать наиболее подходящие, с точки зрения применяемых методик, для пользователя предложения.

В качестве инструмента подбора выбран метод коллаборативной фильтрации основанный на соседстве. Считается, что система знает каждого пользователя и имеет возможность связать его текущий сеанс со всеми его предыдущими сеансами, что позволяет получать данные о его интересах и оценках.

Выбранный алгоритм предполагает сбор данных о продуктах, которые были просмотрены, приобретены или оценены пользователем для дальнейшего построения вектора предпочтений и подбора продуктов на основании сравнения этого вектора с векторами предпочтений других пользователей.

Для выполнения поиска продуктов выполняется построение матрицы оценок товаров пользователями, состоящей из векторов предпочтений каждого пользователя.

После построения матрицы выполняем поиск наиболее похожих на текущего пользователя пользователей, для этого, на текущем этапе, была выбрана косинусная мера (*sim*):

$$sim(x, y) = \frac{\sum_{i \in P_{u_x, u_y}} r_{u_x, i} r_{u_y, i}}{\sqrt{\sum_{i \in P_{u_x, u_y}} r_{u_x, i}^2} * \sqrt{\sum_{i \in P_{u_x, u_y}} r_{u_y, i}^2}},$$

где  $P_{u_x, u_y}$  – підмножество продуктів  $i \in I$ , котрі оцінили обидва користувача,  $r_{u_x, i}$  і  $r_{u_y, i}$  – оцінки користувачів  $x$  і  $y$  для продукту  $i$ . Також до векторів користувачів  $u_x$  і  $u_y$  застосовується Евклідова нормалізація, проєктувальна їх на одиничний круг. Подібність користувачів обчислюється при допомозі скалярного добутку – косинуса кута між точками позначеними векторами. Так як оцінки користувачів додативні – результат обмежений  $[0, 1]$ .

Після розрахунку подібності кожного  $u_x, u_y \in U$  виконується передбачення оцінок для кожного продукту, котрий не був оцінений користувачем  $u \in U$ . Для розрахунку передбачень пропонується використовувати підхід зважених сумм, котрий розраховується по формулі:

$$r_{u_x, i} = \bar{r}_{u_x} + \frac{\sum_{u_y \in R_{u_x, i}} (R_{u_y, i} - \bar{r}_{u_x}) sim(u_x, u_y)}{\sqrt{\sum_{u_y \in R_{u_x, i}} sim(u_x, u_y)}},$$

где  $R_{u_x, i}$  – підмножество користувачів  $u_y \in U$ , різних від  $u_x$ , оцінивших продукт  $i$ , а  $\bar{r}_{u_x}$  – середня оцінка продукту користувачем  $u_x$ . Підхід зважених сумм приймає середні оцінки сусідів активного користувача і зважує кожну з них відповідно до подібності сусіда і активного користувача.

В результаті вибирається  $n$  елементів  $i \in I$  з найбільшою передбаченою оцінкою  $R_{u_x, i}$ . Так як передбачені оцінки показують відповідність релевантності продукту для активного користувача, ми вибираємо перші  $n$  найбільш високо оцінених продуктів з результату розрахунку зважених сумм.

Для реалізації вищеприведеного алгоритму був вибраний мовний програмування Python. Для виконання розподілених обчислень була вибрана бібліотека Apache Spark. Для створення кластера і виконання розподілених обчислень на ньому вибраний сервіс Amazon EMR.