

Метод определения подобия информационных единиц по неявным пользовательским предпочтениям в рекомендательных системах поддержки жизнеобеспечения

© С. А. Филиппов

© В. Н. Захаров

© С. А. Ступников

© Д. Ю. Ковалев

Институт проблем информатики ФИЦ ИУ РАН,
Москва

stanislav@philippov.ru

VZakharov@ipiran.ru
dm.kovalev@gmail.com

ssa@ipi.ac.ru

Аннотация

Целью данной работы является описание метода определения подобия информационных единиц посредством анализа данных о пользовательских предпочтениях. Метод является реализацией подхода Item-Item CF (коллаборативная фильтрация на основе подобия информационных единиц), который в свою очередь является одним из наиболее популярных подходов к построению современных рекомендательных систем. Исходными данными для коллаборативной фильтрации (другими словами для выявления пользовательских предпочтений) являются данные о пользовательской активности при просмотре страниц конкретных интернет-ресурсов (информационных единиц). Данные могут собираться как явным (оценки, опросы, рейтинги), так и неявным образом (протоколирование действий пользователей). Предложенный метод позволяет решить проблему холодного старта, т.е. выдачи рекомендаций в период отсутствия подробной информации о посетителе системы поддержки жизнеобеспечения (здесь и далее под такой системой подразумевается интернет-магазин), но при наличии неявных данных о маршрутах других посетителей системы. Метод опробован на реальных данных, полученных с действующего интернет-магазина Thaisoap, где подтвердил возможность своей применимости в рамках поставленной задачи. Работа выполнена при поддержке Министерства образования и науки РФ, уникальный идентификатор проекта RFMEFI60414X0139.

Введение

Одним из современных трендов в развитии Интернет является персонализация. Поисковые системы, социальные сети, форумы, новостные ресурсы и Интернет магазины стараются адаптировать внешний вид и содержимое (контент) своих страниц под нужды конкретных пользователей. По результатам исследования компании Evergage (www.evergage.com) в 2015 году персонализацию в реальном времени использовали 44% веб сайтов, 17% мобильных сайтов, 13% веб-приложений и 9% мобильных приложений [1]. При этом 78% тех, кто не использует персонализацию сейчас, утверждают, что планируют начать в течение следующих 12 месяцев. Увеличение вовлеченности посетителей, улучшение пользовательского опыта и повышение конверсии считаются самыми важными результатами ее применения.

Предоставление персонализированного контента пользователям позволяет существенно повысить эффективность сайтов, которая выражается в терминологии маркетинга таким показателем как конверсия (число посетителей, совершивших полезные действия к общему числу посетителей выраженное в процентах). Для качественной персонализации сайтов, ориентированных на работу с большой аудиторией пользователей, как правило, используется комплексный подход, сочетающий маркетинговые исследования и анализ поведения конкретных посетителей сайтов. Информацию о маркетинговых качествах посетителей можно получить, в том числе используя системы веб-аналитики, такие как Adobe Digital Marketing Suite или Google Analytics и Siteapps.com. Исходными данными для анализа поведения пользователей являются сведения об их активности, которые могут собираться явным или неявным образом. Явным образом получают результаты голосований и опросов, а также оценки, которые пользователи дают

Труды XVIII Международной конференции DAMDID/RCDL'2016 «Аналитика и управление данными в областях с интенсивным использованием данных», Ершово, 11-14 октября 2016

тем или иным объектам на сайтах. Основное количество информации о пользовательской активности собирается неявным образом посредством протоколирования его действий. Предметом отслеживания являются переходы пользователей по ссылкам на сайтах, время их пребывания на отдельных страницах, факты покупки товаров и услуг. Необходимо отметить, что, речь идет об огромных массивах данных, которые являются неоднородными и требующими отдельных подходов к интерпретации.

В сфере электронной коммерции основным инструментом персонализации контента являются рекомендательные системы, обеспечивающие автоматическую обработку данных о пользовательской активности и выработку рекомендаций на товары и услуги, которые могут быть интересны конкретным пользователям. При реализации рекомендательных систем широко используются методы интеллектуального анализа данных (Data Mining) [2].

Основной задачей, которую авторы данной работы ставили перед собой, является разработка комбинированного подхода к построению рекомендательных систем, обеспечивающего наиболее полное использование всех данных о посетителях интернет-магазинов с целью выработки рекомендаций, наиболее адекватно отражающих их ожидания (пертинентность предложения). Научно-практическая новизна работы заключается в идее комбинированного использования методов Item-Item CF и User-User CF, что позволяет минимизировать недостатки каждого из них и добиться более высокого качества работы рекомендательной системы в целом. Данная статья входит в серию статей по данной проблематике и посвящена описанию метода определения подобия информационных единиц по неявным пользовательским предпочтениям, который является вариантом реализации метода Item-Item CF. Данный метод позволяет вырабатывать приемлемые по качеству рекомендации в условиях, когда сведения о пользовательских предпочтениях отсутствуют, минимальны или слабо информативны. Для выявления групп подобных товаров используются методы кластеризации, что позволяет добиться хороших показателей качества и быстродействия в работе алгоритма.

1 Построение рекомендательных систем с использованием методов коллаборативной фильтрации

Основная задача рекомендательной интернет-системы – формирование контента, максимально соответствующего ожиданиям, в том числе неявным, конкретного пользователя. Для решения этой задачи в большинстве современных рекомендательных систем используется один из двух базовых подходов: коллаборативная фильтрация (collaborative filtering,

CF) и контентная фильтрация (content-based filtering, CbF) [3].

Метод контентной фильтрации фокусируется на выявлении объектов со схожими характеристиками по отношению к тем объектам, которые уже заинтересовали пользователя. При этом учитывается модель поведения пользователя и характеристики (контент) заинтересовавших его объектов. При выработке рекомендаций выявляются объекты со схожими характеристиками (контентом). Для эффективной работы метода контентной фильтрации, как правило, необходимо подробное описание характеристик объектов (так в проекте Music Genome Project музыкальный аналитик оценивает каждую композицию по сотням различных музыкальных характеристик), а также сведения о конкретном пользователе (например, ответы на конкретные вопросы в анкете).

В основе метода коллаборативной фильтрации лежит предположение о консервативности пользовательских предпочтений (т.е. пользователи, одинаково оценивающие определенные объекты, скорее всего аналогичным образом будут оценивать и новые объекты со сходными характеристиками) [4]. По существу, рекомендации базируются на автоматическом сотрудничестве множества пользователей и на выделении (методом фильтрации) тех пользователей, которые демонстрируют схожие предпочтения или шаблоны поведения. Таким образом, метод коллаборативной фильтрации вырабатывает рекомендации, основанные на модели предшествующего поведения пользователя и с учетом поведения пользователей со схожими характеристиками.

Наибольшее распространение в сфере электронной коммерции получили рекомендательные системы, использующие следующие варианты реализации метода коллаборативной фильтрации, а также их гибриды:

- коллаборативная фильтрация посредством анализа предпочтений групп пользователей со схожими интересами (User-User Collaborative Filtering, User-User CF);
- коллаборативная фильтрация посредством анализа взаимосвязей между объектами (Item-Item Collaborative Filtering, Item-Item CF);

Основными проблемами, связанными с реализацией и практическим использованием алгоритмов коллаборативной фильтрации, являются разреженность данных, проблема холодного старта и масштабируемость. Дополнительно к перечисленным проблемам можно отметить проблему ограничения разнообразия предложений. Рекомендательные системы, использующие коллаборативную фильтрацию, склонны предлагать товары уже пользующиеся популярностью, что создает проблемы для продвижения новых товаров и услуг [5].

В методе User-User CF определяется сходство между пользователями и в качестве рекомендаций пользователю выдается n самых часто покупаемых товаров k наиболее похожими на него покупателями. Для оценки степени схожести пользователей в плане их предпочтений могут использоваться различные функции сходства (метрики). Наиболее популярными среди них являются: евклидово расстояние, косинусная мера, расстояние Хэмминга, коэффициент корреляции Пирсона, коэффициент Танимото, Манхэттенское расстояние и некоторые другие [4, 6]. Определение рекомендаций методом User-User CF предполагает построение матрицы активности пользователей, каждая строка которой описывает действия конкретного пользователя применительно к конкретному объекту (категория, товар, услуга) на сайте. Действия пользователей могут обозначаться самыми различными способами. Например, это может быть бинарная информация о посещении или не посещении заданного ресурса данным пользователем, частота (или число) использований ресурса g пользователем u , стоимость или рейтинг, проставленный пользователем u для ресурса g и т.д. Таким образом, каждая строка матрицы активности представляет собой вектор оценок, соответствующих различным категориям товаров (тематический профиль пользователя). Профиль пользователя характеризует степень его интереса к каждой группе товаров. Для каждой пары «пользователь-объект (товар, услуга, действие)» в матрице активности вычисляется мера близости с использованием выбранной метрики [7].

Для поиска рекомендаций конкретному пользователю на основании его поведенческого профиля используются три основных подхода: основанный на соседстве (memory based), основанный на модели (model based) и гибридный подход (hybrid). В современных коммерческих системах наибольшее распространение получили гибридный подход и подход, основанный на использовании моделей (алгоритмы кластеризации, байесовские сети доверия, латентные семантические модели) [3, 9]. Для выявления групп пользователей со схожими характеристиками часто используются различные алгоритмы кластеризации.

Метод Item-Item CF исторически появился как альтернатива методу User-User CF, призванная повысить производительность рекомендательных систем для тех магазинов, где число покупателей существенно превышает количество наименований товаров в каталоге [8]. Первоначально данный метод был предложен компанией Amazon для решения следующих основных проблем подхода User-User CF: проблема холодного старта и проблема частого обновления данных о пользовательской активности. Проблема холодного старта существенно снижает качество работы рекомендательной системы вследствие отсутствия данных о предпочтениях новых (или мало активных) пользователей. Проблема частого обновления данных о

пользовательской активности (в случае компании Amazon речь идет о миллионах покупателей) резко снижает производительность рекомендательной системы в целом.

Основная идея метода Item-Item CF заключается в группировке информационных единиц (товары, услуги, действия) имеющих сходные оценки пользователей (рейтинги). Рекомендации вырабатываются по следующему принципу: пользователю оценившему объект X высоко будет предложен объект Y , который высоко оценили другие пользователи, также высоко оценившие и объект X . Использование метода Item-Item CF позволяет повысить качество рекомендаций для новых пользователей (нет критической зависимости от данных о пользовательских предпочтениях), а также значительно повышает производительность рекомендательной системы в случае, когда количество пользователей существенно превышает количество объектов (характеристики объектов меняются реже). При этом качество рекомендаций в среднем выше, чем в случае использования подхода, основанного на анализе пользовательских профилей. Для вычисления попарной близости информационных единиц могут использоваться те же метрики, что и в случае с парами «пользователь-объект» (часто используется косинусная или модифицированная косинусная меры). Для поиска рекомендаций на основании матрицы объектов часто используются весовые функции и методы регрессионного анализа. Одним из перспективных методов решения задачи Item-Item CF является метод Item2Vec [10]. Тем не менее для большинства интернет-магазинов подход, связанный с рекомендациями по рейтингам, слабо применим в силу отсутствия возможности мотивировать пользователей определять рейтинг информационных единиц (покупатели приходят из поисковых систем и товарных каталогов, делают нужную им покупку и уходят, чтобы больше никогда не вернуться). И встает задача, как в таких условиях сделать рекомендацию (информационное предложение), на которую откликнется пользователь.

2 Определение подобия (кластеров) информационных единиц по неявным пользовательским предпочтениям

В целях решения задачи формирования рекомендации с уместной информацией в условиях недостаточности знаний о пристрастиях пользователей авторами предлагается использовать метод, в основе которого лежит расчёт близости пар и последующая группировка (кластеризация) информационных единиц на основе данных пользователей, последовательно просматривающих несколько товаров. При отсутствии данных предлагается использовать обычные классификаторы с учётом цены и параметров объектов, список «Новинки», а также матрицу «С этим товаром покупают» (аксессуары,

дополняющие основную покупку). При данном подходе явное участие пользователей интернет-магазина в формировании рейтинга товаров не требуется.

Первым шагом алгоритма является построение матрицы подобию информационных единиц, где и по вертикали, и по горизонтали присутствуют все информационные единицы интернет-магазина. Заполнение матрицы происходит по следующему правилу: если пользователь последовательно просмотрел два товара, то вес подобию в матрице для этих двух товаров увеличивается на 1.

Для обработки матрицы в целях выявления групп информационных единиц, которые являются близкими по своим оценкам подобию, из всех известных алгоритмов кластеризации в результате проведенного моделирования был выбран современный производительный алгоритм Affinity Propagation. Одним из преимуществ данного алгоритма является отсутствие необходимости предварительной оценки оптимального количества кластеров [11].

Приведенный метод кластеризации был опробован на тестовом массиве данных, предоставленных интернет-магазином Thaisoap. Магазин ориентирован на продажу натуральной тайской косметики и кокосового масла. Каталог товаров магазина содержит более 1 500 наименований товаров, которые разбиты на 180 классов (44 корневых классов, 136 подклассов). Ежедневно магазин посещают в среднем около 1 500 посетителей и проводят на нем (в среднем) порядка 11 минут каждый (на каждого посетителя приходится в среднем 28 переходов по ссылкам). Исходные данные охватывают период в один квартал (IV квартал 2015 года), в котором каталог товаров был неизменен.

p#3	p#4	p#5	p#6	p#7	p#8	p#9	p#10	p#11	p#12	p#13	p#14
p#1	1.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000
p#4	0.00000000	1.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000
p#5	0.00000000	0.00000000	1.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000
p#6	0.00000000	0.00000000	0.00000000	1.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000
p#7	0.00000000	0.00000000	0.00000000	0.00000000	1.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000
p#8	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	1.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000
p#9	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	1.00000000	0.00000000	0.00000000	0.00000000	0.00000000
p#10	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	1.00000000	0.00000000	0.00000000	0.00000000
p#11	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	1.00000000	0.00000000	0.00000000
p#12	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	1.00000000	0.00000000
p#13	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	1.00000000
p#14	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000	0.00000000

Рисунок 1 Матрица подобию товаров

На основе указанных данных была построена матрица подобию по всему временному периоду. На рисунке 1 представлен фрагмент получившейся матрицы подобию товаров для всех товаров из каталога (значения нормированы). Всего в каталоге на данный момент присутствует 1522 товара. Как видно из рисунка матрица сильно разрежена, так как для многих пар товаров оценка подобию отсутствует (т.е. в течение анализируемого периода времени

пользователи не интересовались некоторыми товарами из каталога).

В результате обработки матрицы подобию по алгоритму Affinity Propagation (с использованием статистического пакета R) была построена гистограмма расстояний. Результаты работы алгоритма представлены на рисунке 2 в виде кластерной тепловой карты (размерность карты 1522 на 1522). Преобладание одного цвета на карте обусловлено тем фактом, что в тестовой выборке данных для большинства пар товаров не определена оценка подобию (т.е. пользователи не интересовались данными товарами в течение рассматриваемого в тестовой выборке периода времени).

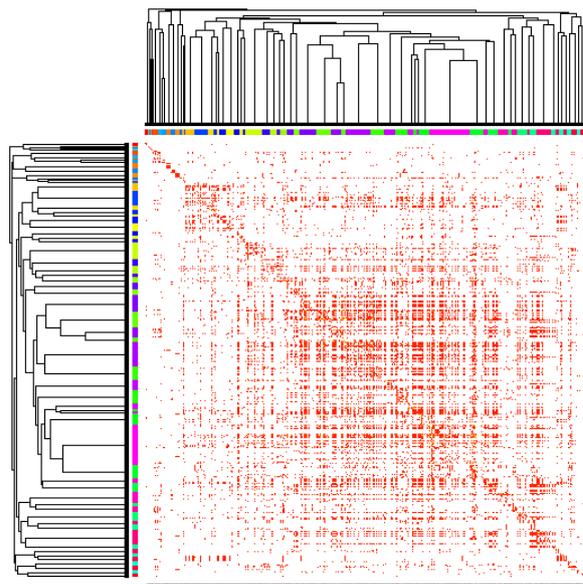


Рисунок 2 Кластерная тепловая карта

Всего алгоритм выделил 64 кластера, наиболее крупными из которых являются кластера с номерами 5 (75 объектов), 8 (44 объекта), 10 (30 объектов), 19 (27 объектов) и 55 (31 объект).

Качество работы алгоритма можно оценить на примере кластера номер 5, описание которого представлено в таблице 1. В частности, видно, что для референсной информационной единицы (массажное кокосовое масло) в кластер подобию попали товары на основе кокосового масла или косвенно ассоциирующиеся с кремами и маслами для ухода за телом.

Аналогичные результаты показывает и исследование других полученных кластеров. Таким образом, метод определения подобию информационных единиц выполняет возложенные на него задачи: формируется рекомендация из информационных единиц (товаров), уместных по отношению к товару, который заинтересовал неизвестного посетителя в данный конкретный момент времени.

Таблица 1 Детализация кластера номер 5

Кластер	Референсная информационная единица	Примеры товаров из кластера
ID: 5 Size: 75	ID: 76. Нерафинированное 100% массажное кокосовое масло "Citronella" Tropicana, 100 мл.	ID: 43. Кокосовое масло Tropicana 1 литр, нерафинированное ID: 51. Кокосовое масло нерафинированное Tropicana в аптекарском флаконе, 90 мл. ID: 466. Восстанавливающий кокосовый ЛОСЬОН для тела Tropicana "Sweet Coconut" (без парабенов), 200 мл. ID: 624. Маска-эксфолиант для лица "Морской коллаген" Artiscent, 100 мл. ID: 1234. Мини-набор Шампунь и Кондиционер для волос "Золотой шелк с экстрактом шелковицы"

Заключение

Персонализация контента интернет-ресурсов на сегодня является одним из активно развивающихся направлений ИТ-индустрии. Важнейшими результатами ее применения являются увеличение вовлеченности посетителей, улучшение пользовательского опыта и повышение конверсии. Персонализация контента в сфере электронной коммерции выражается в адресном предложении товаров, а также услуг конкретным пользователям и реализуется посредством рекомендательных систем. Современные рекомендательные системы обеспечивают обработку огромных массивов данных о пользовательской активности с целью формирования предсказаний для конкретных пользователей в момент запроса.

В данной работе изложен метод определения подобия информационных единиц по неявным пользовательским предпочтениям в рекомендательных системах поддержки жизнеобеспечения на основе упрощенной метрики близости пар информационных единиц по алгоритму

кластеризации Affinity Propagation. Метод проверен на данных интернет-магазина Thaisoap и показал по результатам высокий уровень уместности информации в формируемой рекомендации.

Таким образом описанный метод класса Item-Item CF вполне применим для новых (или малоактивных) пользователей. При этом по мере накопления данных о предпочтениях пользователей рекомендуются отдавать большее предпочтение методам класса User-User CF, которые дают тем более точные предсказания чем более подробны данные о пользовательской активности.

Литература

- [1] Почему персонализация контента это еще не веб-персонализация // Статья в сети Интернет, URL: <http://lpgenerator.ru/blog/2016/03/19/pochemu-personalizaciya-kontenta-eto-eshe-ne-veb-personalizaciya/>
- [2] С.А.Филиппов, В.Н.Захаров, С.А.Ступников, Д.Ю.Ковалев Подходы к повышению пертинентности информационного предложения в медиасервисах на основе обработки больших объемов данных // XVII международная конференция «Аналитика и управление данными в областях с интенсивным использованием данных» DAMDID/RCDL'2015, Октябрь 13-16, Обнинск, 2015, с. 224-228..
- [3] М. Тим Джонс Рекомендательные системы: Часть 1. Введение в подходы и алгоритмы // Статья в сети Интернет, URL: <http://www.ibm.com/developerworks/ru/library/os-recommender1/>, 2013.
- [4] Xiaoyuan Su, Taghi M. Khoshgoftaar A survey of collaborative filtering techniques // Advances in Artificial Intelligence, Volume 2009 (2009), Article ID 421425, 19p.
- [5] Fleder D., Hosanagar K. Blockbuster Culture's Next Rise or Fall: The Impact of Recommender Systems on Sales Diversity // Management Science, Vol. 55, No. 5, May 2009, pp. 697-712.
- [6] В.А. Лексин Технология персонализации на основе выявления тематических профилей пользователей и ресурсов Интернет // ВКР Магистра, Вычислительный Центр им. А.А. Дородницына РАН, 2007.
- [7] Брейкин Е. А. Рекомендательная система на основе коллаборативной фильтрации // Молодой ученый. — 2015. — №13. — С. 31-33.
- [8] Greg Linden, Brent Smith and Jeremy York Amazon.com recommendations: Item-to-Item Collaborative Filtering // Industry Report, IEEE INTERNET COMPUTING, 2003.
- [9] Барсегян А.А., Куприянов М.С., Степаненко В.В., Холод И.И. Методы и модели анализа данных: OLAP и Data Mining // СПб.: БХВ-Петербург, 2004. — 336 с.

- [10] Barkan O., Koenigstein N. Item2Vec: Neural Item Embedding for Collaborative Filtering // arXiv preprint arXiv:1603.04259, Mar 2016.
- [11] Brendan J. Frey, Delbert Dueck Clustering by passing messages between data points // Science 16 Feb 2007 Vol. 315, Issue 5814, pp. 972-976, DOI: 10.1126/science.1136800

Determination of similarity of information items based on implicit user preferences in life-support recommender systems

Stanislav A. Philippov, Victor N. Zakharov,
Sergey A. Stupnikov, Dmitriy Yu. Kovalev

The purpose of this paper is to describe the method for determining the similarity of the information items

through the analysis of user preference data. The method is an implementation approach known as Item-Item CF (collaborative filtering based on the similarity of the information items), which in turn is one of the most popular approaches to the construction of modern recommender systems. Initial data for collaborative filtering are the data about users' activity when they are browsing web resources. Data can be collected as explicit (evaluations, surveys, ratings) and implicit (logging of users' actions). The proposed method solves the problem of cold start using implicit data about the routes of other users. The method was tested on real data from existing online store Thaisoap, which confirmed the possibility of its applicability in the framework of the task. A unique identifier of the project supported by the Ministry of education and science of the RF is RFMEFI60414X0139.