



Московский государственный университет имени М. В. Ломоносова

Факультет вычислительной математики и кибернетики

Кафедра математических методов прогнозирования

Коваленко Павел Антонович

# Использование случайных блужданий по графу для повышения качества рекомендаций

МАГИСТЕРСКАЯ ДИССЕРТАЦИЯ

**Научный руководитель:**

д.ф-м.н., профессор

Дьяконов Александр Геннадьевич

Москва, 2019

# Содержание

<b>1</b>	<b>Введение</b>	<b>2</b>
<b>2</b>	<b>Обзор предметной области</b>	<b>2</b>
2.1	Математическая постановка задачи рекомендаций . . . . .	3
2.2	Метрики качества рекомендаций . . . . .	3
2.3	SVD . . . . .	3
2.4	SVD для ранжирования . . . . .	4
2.5	Random walks . . . . .	4
2.6	Random walks + SVD . . . . .	5
<b>3</b>	<b>Описание датасетов</b>	<b>6</b>
3.1	MovieLens . . . . .	6
3.2	Million Playlist Dataset . . . . .	6
3.3	Сравнение датасетов . . . . .	6

# 1 Введение

Рекомендательные системы — это семейство алгоритмов, предназначенных для выделения из широкого множества объектов подмножества, релевантного поданному на вход запросу.

Вода: что такое рекомендательные системы, примеры задач.

Частым случаем использования рекомендательных систем является построения персонализированных рекомендаций для пользователя некоторого сервиса. Например, рекомендации фильмов, музыки, товаров в интернет-магазине и т.д. В этом случае запросом является пользователь, а множеством объектов для рекомендации — все доступные книги, музыка и товары соответственно. Далее в работе будет рассматриваться задача рекомендаций объектов (или документов) пользователю, однако все идеи легко переносятся на более общий случай.

## 2 Обзор предметной области

Рекомендательные системы разделяются в зависимости от используемых ими данных и предположений на контентные (content-based filtering) и коллаборативные (collaborative filtering). [ссылка]

Контентные модели используют признаковое описание пользователя и объекта для предсказания релевантности объекта к пользователю. Таким образом, задача рекомендаций сводится к задаче регрессии меры релевантности или классификации на релевантные и не релевантные объекты.

Коллаборативные модели действуют в предположении, что похожим пользователям нравятся похожие объекты. Они часто используются в случаях, когда у пользователей и объектов нет признакового описания, только история взаимодействия. В частности, к коллаборативным моделям относятся SVD и случайные блуждания, которые будут рассмотрены далее.

Данная классификация не покрывает все возможные случаи. Существуют модели, которые используют похожесть между пользователями совместно с признаками пользователей и объектов, например, факторизационные машины [ссылка или опи-

сание]. Также часто коллаборативные модели используются как признаки для обучения контентных моделей [ссылки на победителей рексисов].

## 2.1 Математическая постановка задачи рекомендаций

Имеется множество пользователей (запросов)  $U = \{u\}$  и множество объектов (документов)  $I = \{i\}$ ,  $|U| = N$ ,  $|I| = M$ . Для некоторых пар  $(u, i)$  известна оценка  $r_{ui}$ , которую пользователь  $u$  поставил объекту  $i$ . Требуется заполнить пропуски, то есть для всех пар  $(u, i)$  предсказать оценку  $\hat{r}_{ui}$ . Другая возможная постановка — для каждого пользователя  $u$  найти  $K$  максимально релевантных объектов  $i_1, \dots, i_K$ , то есть объектов с наибольшей предсказанной оценкой.

Часто рассматривают задачи с неявным откликом — все  $r_{ui} = 1$ . Например, если человек слушает музыку, то мы знаем, какие треки он слушал, и считаем, что они ему нравятся, но ничего не знаем про все остальные треки. В таком случае обычно считают, что все пары  $(u, i)$ , которых нет в выборке (то есть все треки, которые пользователь не слушал), являются отрицательными примерами, для них считаем  $r_{ui} = 0$ . Далее в работе будет рассматриваться только задача рекомендаций с неявным откликом.

Для удобства часто вводят матрицу  $R = \{r_{ui}\} \in \mathbb{R}^{N \times M}$ , считая, что не все ее ячейки могут быть заполнены, или заполняя пропуски нулями.

## 2.2 Метрики качества рекомендаций

// P@k, R@k, MAP@k, AUC?, nDCG?

## 2.3 SVD

Данная модель предполагает, что есть некоторое  $d$ -мерное вещественное пространство, и всем пользователям  $u$  и объектам  $i$  соответствуют вектора в этом пространстве  $p_u$  и  $q_i$ , называемые латентными векторами. Близость (релевантность) между пользователем и объектом определяется как скалярное произведение между их векторами:  $\hat{r}_{ui} = \langle p_u, q_i \rangle$ . Задача состоит в поиске латентных векторов пользователей и объектов, для которых среднеквадратичное отклонение на объектах обучающей

выборки является минимальным:

$$L(p, q) = \sum_{u, i} (r_{ui} - \langle p_u, q_i \rangle)^2 \rightarrow \min_{p, q}$$

Данная модель часто называется Singular Vector Decomposition (SVD), поскольку в случае полностью заполненной матрицы  $R$  оптимальными значениями для  $p$  и  $q$  являются соответствующие столбцы из сингулярного разложения матрицы  $R$ :  $R = P\Sigma Q^T$  [ссылка]. Также для поиска латентных векторов можно использовать алгоритм Alternating Least Squares, основанный на итеративной оптимизации сначала пользовательских, потом объектных векторов. Поэтому описанная выше модель также встречается под названием ALS. [ссылка на статью Миши]

## 2.4 SVD для ранжирования

В случаях, когда не требуется предсказать оценку, а только для каждого пользователя упорядочить объекты по релевантности, используют в качестве метрики качества не среднеквадратичное отклонение оценки, а триплетный лосс, который можно ввести следующим образом [ссылка]:

$$L(u, i_p, i_n) = \log(\sigma(\langle p_u, q_{i_p} \rangle - \langle p_u, q_{i_n} \rangle))$$

где  $u$  — пользователь,  $i_p, i_n$  — два объекта, таких что  $r_{u, i_p} > r_{u, i_n}$ ,  $p_u$  — латентный вектор пользователя,  $q_{i_p}, q_{i_n}$  — латентные вектора объектов. Итоговый функционал можно записать в следующем виде:

$$L(p, q) = \sum_{u, i, i'} \mathbb{I}[r_{u, i} > r_{u, i'}] \log(\sigma(\langle p_u, q_i \rangle - \langle p_u, q_{i'} \rangle)) \rightarrow \min_{p, q}$$

## 2.5 Random walks

Пусть задан ориентированный граф. Неформально определить случайное блуждание можно следующим образом. Есть некий агент, который перемещается по графу, стартуя в некоторой вершине  $v_0$ . Далее в каждый момент времени  $t$  он перемещается из вершины  $v_t$  в случайную вершину  $v_{t+1}$ , соседнюю с вершиной  $v_t$ . Вероятности переходов между вершинами являются частью алгоритма, обычно вероятности

равные или зависят от степени вершины. Таким образом получаем бесконечную последовательность вершин  $\{v_k\}$ , в которой соседние вершины связаны ребром.

Чтобы использовать случайные блуждания для задачи рекомендаций, построим двудольный граф. Вершинами первой доли будут пользователи, второй — объекты. Между пользователем и объектом есть ребро, если пользователь положительно оценил объект (например, прослушал трек).

Чтобы построить рекомендации для пользователя  $u$ , запустим из соответствующей пользователю вершины  $K$  случайных блужданий, каждое из которых остановим после нечетного числа шагов  $L$ . В этом случае в силу двудольности графа все блуждания завершатся в вершинах, соответствующих некоторым объектам. После этого можно ввести меру близости между пользователем  $u$  и объектом  $i$  как число случайных блужданий, завершившихся в вершине  $i$ , и выбрать объекты, в которых заканчивается наибольшее число путей.

В случае  $L = 1$  получатся абсолютно точные и бесполезные рекомендации — случайное блуждание всегда будет останавливаться в объектах, уже оцененных пользователем.

## 2.6 Random walks + SVD

Одной из проблем при обучении SVD является низкая плотность матрицы  $R$ , из-за чего алгоритму требуется много шагов оптимизации для нахождения оптимальных латентных векторов [хорошо бы ссылку на пруф]. В статье [ссылка!] предложен гибридный подход, совмещающий преимущества случайных блужданий и SVD.

Используется логистический функционал, описанный в 2.4. Латентные вектора обучаются стохастическим градиентным спуском на батчах, сформированных следующим образом. Для каждого примера из батча случайно выбирается пользователь, для него выбрать случайный объект в качестве отрицательного примера, а в качестве положительного — конечную вершину случайного блуждания, начавшегося в вершине пользователя. Таким образом, обучающая выборка помимо объектов, которые пользователь явно оценил, пополняется похожими на них объектами. За счет этого понижается разреженность матрицы  $R$ .

Эксперименты, проведенные в [??], показывают, что за счет пополнения выборки удается повысить качество рекомендаций, измеренное как  $\text{precision@10}$ ,  $\text{recall@10}$  и  $\text{MAP@10}$ .

## 3 Описание датасетов

### 3.1 MovieLens

Данный датасет является выгрузкой оценок с сайта `movielens.org`, где пользователи могут ставить оценки фильмам, которые они просмотрели, и получать на основе этих оценок рекомендации для просмотра [ссылка]. Датасет предоставляется в четырех вариантах: `MovieLens 100K`, `MovieLens 1M`, `MovieLens 10M` и `MovieLens 20M`, которые содержат  $10^5$ ,  $10^6$ ,  $10^7$  и  $2 \times 10^7$  оценок соответственно. Оценки — целые или полуцелые числа от 1.0 до 5.0. Распределение оценок представлено на рис. ???. Во всех четырех датасетах есть только пользователи, оценившие не менее 20 фильмов.

// Неплохо бы нарисовать распределение числа оценок на пользователя и на фильм.

### 3.2 Million Playlist Dataset

Данный датасет был собран сервисом онлайн-стриминга музыки `spotify.com` для соревнования RecSys Challenge 2018 [ссылка на рексис, ссылка на датасет]. Объектами в данном датасете являются аудио-треки, а запросами — плейлисты, составленные пользователями сервиса. В датасете представлен 1 миллион плейлистов, каждый из которых содержит не менее 5 и не более 250 треков.

Этот датасет является примером датасета с неявным откликом — для каждого плейлиста известно только, какие треки ему точно релевантны (то есть содержатся в нем), а про все остальные треки не известно ничего. Поэтому все треки, не включенные в плейлист, считаются не релевантными этому плейлисту.

### 3.3 Сравнение датасетов

\* Количество запросов

- \* Количество объектов
- \* Количество оценок
- \* Плотность матрицы оценок