

Д. М. Понизовкин

Влияние меры сходства на результативность РС

Аннотация. В данной работе рассматриваются две основные задачи традиционных РС, вводятся понятия степень релевантности и релевантность между контентом пользователей или объектов, вычисление которых производится с помощью меры сходства. От выбора меры сходства зависит выполнение транзитивности отношения релевантности. В работе показано влияние транзитивности на значение критерия качества, который характеризует качество работы РС. Показано, что не всегда свойство транзитивности выполняется в традиционных РС.

Ключевые слова и фразы: рекомендательная система, мера сходства, транзитивность, релевантность контента.

Введение

Традиционная рекомендательная система (далее РС) — это информационная система, целевые задачи которой состоят:

- (1) в поиске N объектов таких, что степень релевантности их контента и контента пользователя высока;
- (2) в выявлении степени релевантности контента некоторого объекта РС контенту пользователя.

Для упрощения изложения будем считать, что между пользователями и их контентом существует взаимно-однозначное соответствие. То же верно и для объектов. Далее по тексту контент пользователя и пользователь или контент объекта и объект — синонимичные и взаимозаменяемые понятия. К примеру, говоря, что объекты релевантны, имеется в виду, что релевантны их контенты. Здесь и далее контент — это некоторая структура u , которая хранит информацию о пользователе, или t , которая хранит информацию об объекте. Множество контентов пользователей обозначим U , $|U| = |U|$, объектов —

$T, |T| = |T|$. Между любыми двумя контентами *выполняется отношение релевантности \mathbf{R}* , если *степень их релевантности высока*. Говоря о степени релевантности или релевантности пользователя или объекта, будем предполагать степень релевантности или релевантность их контентов. Степень релевантности u^i и t^j задается самим пользователем РС оценкой v_j^i по некоторой шкале \mathcal{S} в момент работы с системой, однако не для каждой пары $(u^i, t^j) \exists v_j^i$. Степень релевантности может быть вычислена алгоритмически через задание *меры сходства* sim — функции, сопоставляющей паре контентов вещественное число, определяющее степень релевантности в области определения \mathcal{D}_{sim} меры сходства. Система (V, e, sim) ранжирует множество $V: \text{sim}(v_i, e) \geq \text{sim}(v_{i+1}, e)$.

Первая целевая задача РС называется **top-N** [1]. Она заключается в формировании подмножества $T_{\text{top}}^a = \{t \in T' \subset T | t \mathbf{R} u^a\}$, $|T_{\text{top}}^a| = N$, u^a — активный пользователь системы, для которого в данный момент решается задача РС.

Вторая целевая задача РС называется задачей *прогнозирования неизвестной оценки* v_p^a [2] и определяется как вычисление такой оценки vr_p^a , что $vr_p^a \approx v_p^a$.

Для оценки качества работы РС проводится тестирование, после которого производится расчет *критерия оценки*. Критерий оценки — функция, сопоставляющая результату вещественное число, по которому можно судить о качестве РС. Множество входных данных, как правило, разбивается на два подмножества: обучающее (будем индексировать его и его элементы символом 0), на котором производится тест, и тестовое (будем индексировать его и его элементы символом \times), с которым производится сравнение полученного результата.

1. Задача top-N в традиционных РС

1.1. Модель

Неформально модель задачи top-N традиционных РС определяется следующим *эвристическим утверждением*: если пользователю «нравится» объект t^1 и объект t^2 «похож» на объект t^1 , то пользователю «понравится» объект t^2 [1]. Понятие «похож» определяется традиционными РС на основании значений мер сходства $\mathcal{D}_{\text{sim}} = T \times T$, понятие «нравится» не определено, так как не введена $\mathcal{D}_{\text{sim}} = U \times T$. Следуя введенной в статье терминологии, утверждение примет вид:

$$(1) \quad u^a \mathbf{R}t^i \wedge t^i \mathbf{R}t^j \Rightarrow u^a \mathbf{R}t^j$$

Формула (1) описывает транзитивность отношения релевантности.

Необходимые данные для решения задачи и расчета критерия оценки:

- (1) $T_0^a = \{t_0^j | t_0^j \mathbf{R}u^a\}$
- (2) $T_x^a = \{t_x^j | t_x^j \mathbf{R}u^a\}$
- (3) $T_0^a \cap T_x^a = \emptyset$

1.2. Схема решения

Неформальное описание схемы решения определяется следующим образом: найти N объектов, наиболее «схожих» с теми, которые нравятся пользователю.

Схема решения:

$$(2) \quad T_{\text{top}}^a = \{t_1, \dots, t_N\} \subset (T \setminus T_0^a, T_0^a, \sum_{t_0 \in T_0^a} \text{sim})$$

Модель и решение реализованы, к примеру, компанией [Amazon](#) [3].

1.3. Пример традиционного решения задачи top-N

Контентом пользователя является вектор, отображающий информацию о том, какие объекты пользователь предпочел (поставил высокую оценку) за время работы с РС. $u^a = (v_1^a, \dots, v_{|T|}^a)$, где:

$$v_j^a = \begin{cases} 1, & \text{если } u^i \mathbf{R}t^j \\ 0, & \text{иначе.} \end{cases}$$

$t^j \in T_0$. К примеру, для системы Amazon, если $v_j^a = 1$, то пользователь приобрел товар j .

Контент объекта $t^j = (ct_1^j, \dots, ct_{nt}^j)$ — вектор некоторых характеристик. Множество характеристик объектов и их значения зависят от системы, и для описания решения их определение не является обязательным.

Типичной мерой сходства, используемой при решении задачи top-N является $\text{sim}(t^i, t^j) = \cos(\angle(t^i, t^j))$ — косинус угла между векторами [3].

Модель РС задается матрицей \mathcal{M} мер сходства размерности $|T| \times |T|$:

$$\mathcal{M}_{ij} = \begin{cases} \text{sim}(t^i, t^j), & i \neq j \\ 0, & \text{иначе} \end{cases}$$

Решение основано на (2) и принимает вид:

$$T_{\text{top}}^a = \left\{ t \mid \sum_{t_0} \text{sim}(t, t_0) \rightarrow \max \right\}$$

Шаги решения:

(1) $s = \mathcal{M} \times u^a = (s_1, \dots, s_{|T|})$, где

$$s_j = \begin{cases} s_j, & \text{если } s_j \in (s, 0, s_i) \\ 0, & \text{иначе} \end{cases}$$

в вектор s входит N максимальных элементов, остальные заменяются на ноль

(2) $T_{\text{top}}^a = \{t^j \mid s_j \neq 0\}$.

ALGORITHM 1. Построение матрицы мер сходства модели

```

1: for  $i = 1, i < |T|$  do ▷ Построение модели
2:   for  $j = i, j < |T|$  do
3:      $\mathcal{M}_{ij} = \text{sim}(t^i, t^j)$ 
4:      $\mathcal{M}_{ji} = \mathcal{M}_{ij}$ 
5:      $j = j + 1$ 
6:   end for
7:    $i = i + 1$ 
8: end for

```

ALGORITHM 2. Решение задачи *top-N*

```

1:  $s = (s_1, \dots, s_{|T|}) = \mathcal{M} \times u^a$ 
2:  $T_{\text{top}}^a = \{t^j \mid s_j \in (s, 0, s_j)\}$ 

```

1.4. Схема оценки качества решения задачи top-N. Примеры

Для того, чтобы определить качество решения задачи top-N, необходимо определить число объектов результирующей выборки, релевантных пользователю: $|\{t_{\text{top}}^j \in T_{\text{top}}^a | u^a \mathbf{R}t_{\text{top}}^j\}|$. Описанная модель и схема решения не использует $\mathcal{D}_{\text{sim}} = U \times T$, поэтому определение качества основано на (1): $t_{\times}^i \mathbf{R}t_{\text{top}}^j \wedge u^a \mathbf{R}t_{\times}^i \Rightarrow u^a \mathbf{R}t_{\text{top}}^j$. Существующие критерии E_{top} прямо-пропорциональны числу $|\{t_{\text{top}} | t_{\text{top}} \mathbf{R}t_{\times}\}|$.

Примеры критериев E_{top} :

- прежде, чем задать примеры, введем функцию

$$r(i) = \begin{cases} 1, & \exists t_{\times} t_{\text{top}}^i \mathbf{R}t_{\times} \\ 0, & \text{иначе.} \end{cases}$$

- Точность $P = \frac{1}{N} \cdot \sum_{i=1}^N r(i)$;
- $P@K = \frac{1}{K} \cdot \sum_{i=1}^K r(i)$;
- $AveP = \frac{1}{\sum_l r(i)} \sum_{K=1}^N P@K \cdot r(K)$;

Результат решения качественный, когда $E_{\text{top}} \rightarrow \max$. При нарушении свойства транзитивности релевантности $E_{\text{top}} \not\rightarrow \max$.

1.5. Надежность решения и транзитивность релевантности

Схема решения *не гарантирует* $E_{\text{top}} \rightarrow \max$, так как не гарантирует выполнения (1) в общем случае:

- пользователь предпочитает различные объекты;
- с течением времени вкусы пользователя могут смениться.

Эти причины могут привести к тому, что $T_0 \neq \{t_0 | t_0 \mathbf{R}t_{\times}\}$, за счет чего $E_{\text{top}} \not\rightarrow \max$.

2. Задача прогнозирования традиционных РС

2.1. Модель

Неформально модель задачи прогнозирования в традиционных РС определяется следующим *эвристическим утверждением*: пользователи с «похожими» предпочтениями в прошлом, будут иметь

«похожие» предпочтения в будущем. Предпочтения для задачи прогнозирования выражаются оценкой, а понятие «похож» определяется традиционными РС на основании значений мер сходства $\mathcal{D}_{\text{sim}} = U \times U$. Следуя введенной в статье терминологии, утверждение примет вид:

$$(3) \quad u^1 \mathbf{R} u^2 \Leftrightarrow \forall t^i : v_i^1 \approx v_i^2$$

Модель заключается в выявлении множества *socedей* активного пользователя \mathcal{N}_p^a :

$$\mathcal{N}_p^a = \{u^i | u^i \mathbf{R} u^a\} \wedge \exists v_p^i, i = 1..N$$

На основании оценок этих пользователей составляется прогнозная оценка.

Это направление разрабатывалось, к примеру, разрабатывались исследователями, принимавшими участие в соревновании **Netflix Prize**. Необходимые данные для решения задачи и расчета критерия оценки:

- (1) $U_0^a = \{v_j^a\}$
- (2) $U_{\times}^a = \{v_j^a\}$
- (3) $U_0^a \cap U_{\times}^a = \emptyset$

2.2. Схема решения

$$\mathcal{N}_p^a = \{u^1, \dots, u^N\} \subset (U, u^a, \text{sim}), \exists v_p^i, i = 1..N$$

Прогнозная оценка:

$$vp_p^a \approx pr\{v_p^i\}, u^i \in \mathcal{N}_p^a$$

pr — некоторая функция, сопоставляющая множеству оценок прогнозную оценку. Примеры функции pr :

- Средняя оценка:

$$\frac{1}{\mathcal{N}_p^a} \cdot \sum_{u^i \in \mathcal{N}_p^a} v_p^i$$

- Средняя взвешенная оценка:

$$\frac{\sum_{i \in \mathcal{N}_p^a} \text{sim}(u^a, u^i) \cdot v_p^i}{\sum_{u^i \in \mathcal{N}_p^a} |\text{sim}(u^a, u^i)|}$$

2.3. Пример традиционного решения задачи прогнозирования

ALGORITHM 3. Построение множества соседей

```

1:  $\mathcal{N}_p^a = \emptyset$ 
2: for  $i = 1, |U|$  do
3:   if  $u^a \mathbf{R}u^i \wedge \exists v_p^i$  then
4:      $\mathcal{N}_p^a = \mathcal{N}_p^a \cup \{u^i\}$ 
5:      $i = i + 1$ 
6:   end if
7:   if  $i \geq N$  then
8:     выход из цикла
9:   end if
10: end for

```

ALGORITHM 4. Решение задачи прогнозирования при f , являющей простой взвешенной суммой

```

1:  $vp_p^a = 0$ 
2:  $divider = 0$ 
3: for  $i = 1, N$  do
4:    $s = |\text{sim}(u^a, u^i)|$ 
5:    $vp_p^a = vp_p^a + s \cdot v_p^u$ 
6:    $divider = divider + s$ 
7:    $i = i + 1$ 
8: end for
9:  $vp_p^a = \frac{vp_p^a}{divider}$ 

```

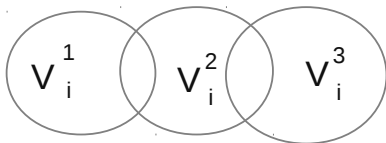


Рис. 1. Транзитивность «близких» оценок

2.4. Схема оценки качества, критерий E_p

Расчет оценки качества основан на (3): так как прогнозная оценка строится на оценках релевантных пользователей, то она должна быть приблизительно равна реальной. Для оценки качества определяется разность прогнозной и реальной. Примеры критерия E_p :

$$(1) MAE = \frac{1}{M} \sum_{p=1}^M |vp_p^a - v_p^a|$$

$$(2) RMSE = \sqrt{\frac{1}{M} \sum_{p=1}^M (vp_p^a - v_p^a)^2}$$

где M — число прогнозных оценок.

Результат решения является качественным, если $E_p \rightarrow \min$. Традиционные РС основывают решение на предположении, что, если $v_p^a \approx v_p^i \approx v_p^j \approx pr(\{v_p^a, v_p^j\}) \Rightarrow |vp_p^a - v_p^a| \rightarrow \min$. Предполагается, что модель и схема решения гарантирует выполнение $E_p \rightarrow \min$, так как

$$(4) u^a \mathbf{R}u^i \wedge u^a \mathbf{R}u^j \Rightarrow u^i \mathbf{R}u^j$$

2.5. Надежность решения и транзитивность релевантности

Схема решения, в общем случае, не гарантирует выполнение $E_p \rightarrow \min$:

- Отношения \approx не транзитивно: $\forall t^i : v_i^1 \approx v_i^2 \wedge \forall t^i : v_i^2 \approx v_i^3 \not\Rightarrow \forall t^i : v_i^1 \approx v_i^3$ Эта ситуация проиллюстрирована на Рис 1. Таким образом, если $u^1, u^2 \in \mathcal{N}^a \not\Rightarrow \forall t^i : v_i^1 \approx v_i^2 \Rightarrow E_p \not\rightarrow \min$
- Не всякая мера сходства гарантирует выполнение транзитивности. К примеру, в качестве функции sim часто используются коэффициенты корреляции, для которых транзитивность может не выполняться [4, 5]:

$$u^1 \mathbf{R}u^2 \wedge u^2 \mathbf{R}u^3 \not\Rightarrow u^1 \mathbf{R}u^3, u^i \in \mathcal{N}$$

ТАБЛИЦА 1. Сравнение качества традиционных РС

Тип РС	MAE	NMAE	RMSE
Традиционная РС, Пирсон	0.85	0.21	1.09
Реформированная, Хэмминг	0.65	0.16	0.96

- Если предпочтения пользователей схожи, то это не означает, что их оценки схожи [7]. Данное обстоятельство является следствием того, что характеры пользователей могут отличаться. Лояльные пользователи ставят оценки не ниже определенной, критики — наоборот, поэтому разница оценок может быть существенной, несмотря на релевантность пользователей.

3. Сравнение качества традиционных РС при использовании различных мер сходства

Для проведения тестов РС была взята база данных [MovieLens 1M](#). Описание данных:

- число пользователей — 6040;
- число объектов — 3952;
- $v_j^i \in \{1, 2, 3, 4, 5\}$;

На этих данных была решена задача прогнозирования в традиционной модели РС с использованием коэффициента корреляции Пирсона в качестве меры сходства. Полученные результаты сравнивались с традиционной РС, работающей с контентными, которые представляют собой нечеткие множества. Описание подобного представления данных приведено в [6]. Для этого была проведена небольшое реформирование контентов пользователей: $v_j^i = \frac{(v_j^i - 1)}{4}$ и использованием расстояния Хэмминга в качестве меры сходства. *Кроме данной реформации данных никаких других настроек системы с целью улучшения результата не производилось.*

4. Заключение

Для получения хороших показателей критериев качества оценки РС, необходимо, чтобы выполнялось свойство транзитивности релевантности контентов. Однако данное свойство выполняется не всегда в традиционных РС. Качество традиционных РС можно повысить

с помощью небольшой реформации данных и использовании мер сходства, обладающих метрическими свойствами.

Список литературы

- [1] M. Deshpande, G. Karypis, “Item-based top-N recommendation algorithms”, *ACM Transactions on Information System*, **22**:1, pp. 143–177 ↑ 56.
- [2] X. Su, T.M. Khoshgoftaar, “A Survey of Collaborative Filtering Techniques”, *Advances in Artificial Intelligence*, **2009**, pp. 19 ↑ 56.
- [3] G. Linden, B. Smith, J. York, “Amazon.com recommendations: item-to-item collaborative filtering”, *IEEE Internet Computing*, **7**:1 (2003), pp. 76–80 ↑ 57.
- [4] A.E.C. Sotos, S. Vanhoof, W.V. Noortgate, P. Oughena, “The Non-Transitivity of Pearson’s Correlation Coefficient: An Educational Perspective”, Proceedings of the 56th Session of the ISI, *Bulletin of the ISI*, **62**, pp. 4609–4613 ↑ 62.
- [5] P. Burger, *Non-Transitivity Property of Correlation*, <http://ec2-54-245-224-94.us-west-2.compute.amazonaws.com/wordpress/2013/05/24/non-transitivity-property-of-correlation/> ↑ 62.
- [6] С. А. Амелкин, Д. П. Позновкин, «Математическая модель задачи top-N для контентных рекомендательных систем», *Известия МГТУ МАМИ*, **2**, с. 26–31 ↑ 63.
- [7] R. Jin, L. Si, C. Zhai, J. Callan, “Collaborative filtering with decoupled models for preferences and ratings”, Proc. of the 12th international conference on Information and knowledge management, pp. 309–316 ↑ 63.

Рекомендовал к публикации

д.т.н. В. И. Гурман

Об авторе:



Денис Михайлович Позновкин

Автор статьи — аспирант, интересом исследований которого является область рекомендательных систем.

e-mail:

denis.ponizovkin@gmail.com

Образец ссылки на эту публикацию:

Д. М. Позновкин. *Влияние меры сходства на результативность РС* // Программные системы: теория и приложения: электрон. научн. журн. 2014. Т. 5, № 5(23), с. 55–65.

URL http://psta.psiras.ru/read/psta2014_5_55-65.pdf

Denis Ponyzovkin. *Quality of recommender systems and transitivity of content's relevance.*

ABSTRACT. In this paper, we propose relevance and degree of relevance of user's or item's contents. Heuristic motivations of traditional recommender systems assume that there is transitivity property of relevance. We introduce importance of transitivity property, its influence on values of evaluation metrics and show that transitivity not performed in general case for traditional recommender systems. (*in Russian*).

Key Words and Phrases: traditional recommender system, similarity measure, transitivity, content's relevance, evaluation metric.