

Введение в машинное обучение. Лекция 8



Отмечаемся и оставляем отзывы.

Спасибо!

Содержание лекции



- 1. Обзор задач для рекомендательных систем
- 2. Данные
- 3. Виды рекомендательных систем
- 4. Метрики
- 5. Примеры

Постановка задачи



Требуется порекомендовать пользователю что-то (товар, услугу, продукт, контент), что сможет его заинтересовать.



Рекомендательные системы



Зачем Нужны?

Рекомендательные системы



- тах вероятность покупки
- тах матожидание прибыли
- товары из категории (long-tail)

Предмет рекомендации



Что рекомендовать?

Предмет рекомендации



Товары	Книги Фильмы Музыка Игры Приложения
Контент	Новости Сайты Статьи
Досуг	Рестораны Отели Театральные представления Выставки Туры
Социальные связи	друзья группы

Предмет рекомендации



Что рекомендовать?

- Товары (Amazon, Ozon, Aliexpress)
- Статьи (Arxhiv.org)
- Новости (Surfingbird, Яндекс.Дзен, Пульс)
- Изображения
- Видео (YouTube, Netflix)
- Люди (Linkedin, VK)
- Музыка (Yandex.Music, Spotify)

В целом рекомендовать можно что угодно

Netflix Prize



Задача:

 Необходимо улучшить алгоритм рекомендации фильмов на Netflix на 10%

Обучающие данные:

- 100 490 507 оценок (от 1 до 5)
- 480 189 пользователей
- 17 770 фильмов

Предсказать:

• 2 800 000 оценок (float)



Метрика:

RMSE

Netflix Prize



На открытой выборке:

Nº	Команда	СКО	Улучшение в %	Время отправки
1	The Ensemble	0.8553	10.10	2009-07-26 18:38:22
2	BellKor's Pragmatic Chaos	0.8554	10.09	2009-07-26 18:18:28

На **скрытой** выборке обе команды улучшили результаты на **10,06** %

Какие входные данные?



Users (Пользователи)

• Items (Товары, Музыка ..)

• Матрица взаимодействий (user-item)

Требуется



• Предсказать оценку user-a *u к товару і*

• Дать персональные рекомендации

• Найти похожие товары

Признаки пользователей





Признаки пользователей



- Пол
- Возраст
- Интересы (группы ok, vk)
- Посещенные страницы (title, description)
- Место нахождения
- Купленные/просмотренные товары
- Установленные приложения
- Любимые точки пользователей

И так далее..

Признаки item-ов





Признаки item-ов



- наименование
- стоимость
- текстовое описание
- категория
- изображение item-a

Зависит от предметной области

Матрица взаимодействий (user-item)



	фильм #1	фильм #2	фильм #3	фильм #4	фильм #5	фильм #6
Аня	4			5	2	3
Наташа		1	3		1	4
Светлана			4	2	2	
Юлиана	5		2			
Андрей			2		5	

Матрица взаимодействий (user-item)

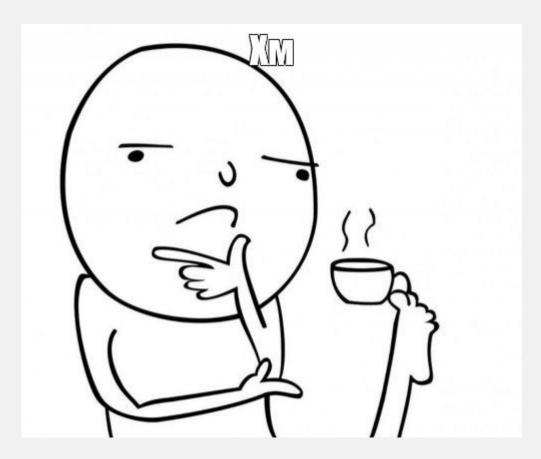


	фильм #1	фильм #2	фильм #3	фильм #4	фильм #5	фильм #6
Аня	4			5	2	3
Наташа		1	3		1	4
Светлана			4	2	2	
Юлиана	5		2			
Андрей			2	?	5	

Рекомендательные системы



Какие же они бывают?



Рекомендательные системы



- Кластеризация пользователей
- Модель совстречаемости
- Summary-based (неперсональные)
- Content-based (основанные на описании товара)
- Коллаборативная фильтрация
- Матричная факторизация

Кластеризация пользователей



- Выберем меру схожести пользователей sim(u, v) по истории их оценок
- Объединим пользователей в кластеры так, чтобы похожие пользователи попали в один кластер
- Оценку пользователя объекту будем предсказывать как среднюю оценку всего кластера по этому объекту

Кластеризация пользователей



Проблемы:

Кластеризация пользователей



Проблемы:

• Нечего рекомендовать **новым/нетипичным пользователям**

• Если в кластере никто не оценивал объект, то предсказание сделать не получится

Модель совстречаемости (ассоциативные правила)



Похожими считаются товары, которые смотрят / кладут в корзину / покупают (нужное подчеркнуть) вместе.

С ЭТИМ ТОВАРОМ ЧАСТО ПОКУПАЮТ











Summary-based рекомендации



Неперсональные рекомендации – когда всем пользователям рекомендуется одно и тоже (например, топ товары).





He Summary-based рекомендации



Персональные рекомендации используют всю доступную информацию о клиенте.

Рекомендуем вам



11111 2 076,08 py6.



48 877,00 руб.



ППП 3 829,66 руб.



135,86 py6.



424,06 py6.



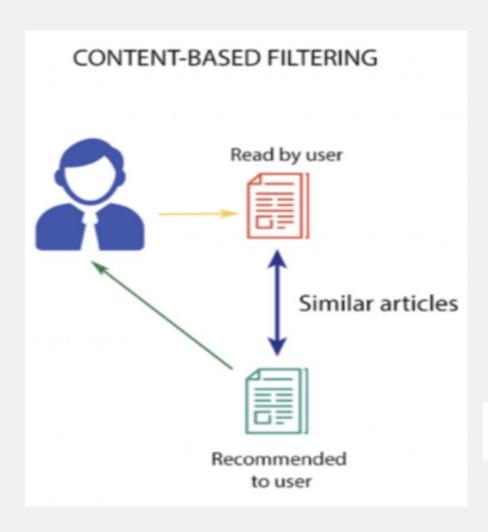
1 466,71 py6.



Описание товара сопоставляется с интересами пользователя, полученными из его предыдущих оценок.

Купил утюг - теперь его предлагают на всех сайтах с Desktop, в телефоне, в объявлении в подъезде и даже во сне!





Если есть хорошие признаковые описания пользователей и объектов (и только они), тогда

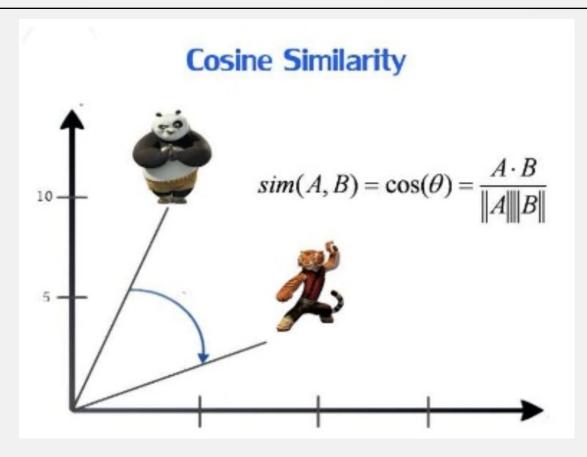
$$u \sim f_u$$
$$i \sim f_i$$

Можно решать как обычную задачу обучения с учителем

$$\{([f_u,f_i],r_{ui})\}$$

Цель:
$$u \to i_1, ..., i_k : \hat{r}_{ui_1} \ge \hat{r}_{ui_2} \ge ...$$





	утюг	увеличенное	отверстие	для	залива	воды	регулировка	температуры	режим	сухого	глажения	регулируемая	подача	пара
item1	1	0	0	0	0	0	1	1	0	0	0	0	0	0
item2	0	0	0	1	0	1	0	0	0	0	0	0	0	0
item3	1	0	0	0	0	0	0	0	0	0	0	0	0	0



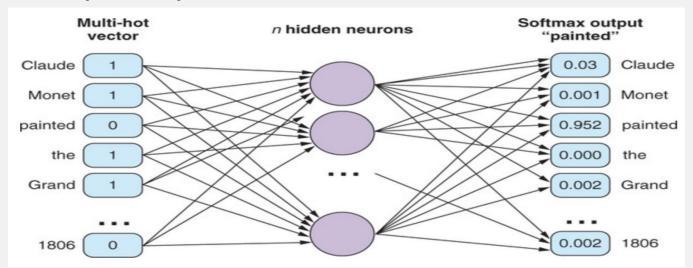
Что использовать?

- TF-IDF
- Word2Vec
- Doc2Vec

$$w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$$

 $tf_{i,j}$ = number of occurrences of i in j df_i = number of documents containing iN = total number of documents

Сжатие размерности: SVD/ PCA





Преимущества	Недостатки
Независимость от данных других пользователей	Когда появляется новый пользователь с недостатком данных о его транзакциях, мы не можем качественно делать рекомендации.
Нет проблемы "холодного старта" для новых предметов, т.к. используя признаки предметов, мы можем легко находить похожие предметы	Формирование четких групп похожих продуктов может ограничить рекомендации других продуктов. Мы можем снова и снова рекомендовать лишь малое подмножество из всех продуктов
Результаты рекомендаций интерпретируемы	Если информация о продуктах ограниченна, трудно различать предметы и группировать их. В результате качество рекомендаций будет низким



Если известна лишь статистика:

$$\{(u,i,r_{ui})\}$$

нет содержательных признаков!

Решение на статистике поведения лучше, чем на описаниях!

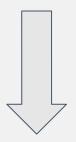


Item-based



Предложи к твоим товарам другие связанные с ними товары (чехол к телефону, чипсы к пиву и т.д.)

User-based



Предложи товары или услуги, которые есть у твоих друзей (похожих пользователей)



По пользователям (User-based)

$$\hat{r}_{ui} = \bar{r}_u + \frac{\sum_{v} sim(u, v)(r_{vi} - \bar{r}_v)}{\sum_{v} sim(u, v)}$$

По товарам (Item-based)

$$\hat{r}_{ui} = \bar{r}_i + \frac{\sum_{j} sim(i,j)(r_{uj} - \bar{r}_j)}{\sum_{j} sim(i,j)}$$



Что лучше? item-based или user-based?



Что лучше? item-based или user-based?

• Когда пользователей много (почти всегда), задача поиска ближайшего соседа становится плохо вычислимой



Что лучше? item-based или user-based?

- Когда пользователей много (почти всегда), задача поиска ближайшего соседа становится плохо вычислимой
- Оценка близости товаров гораздо более точная, чем оценка близости пользователей



Что лучше? item-based или user-based?

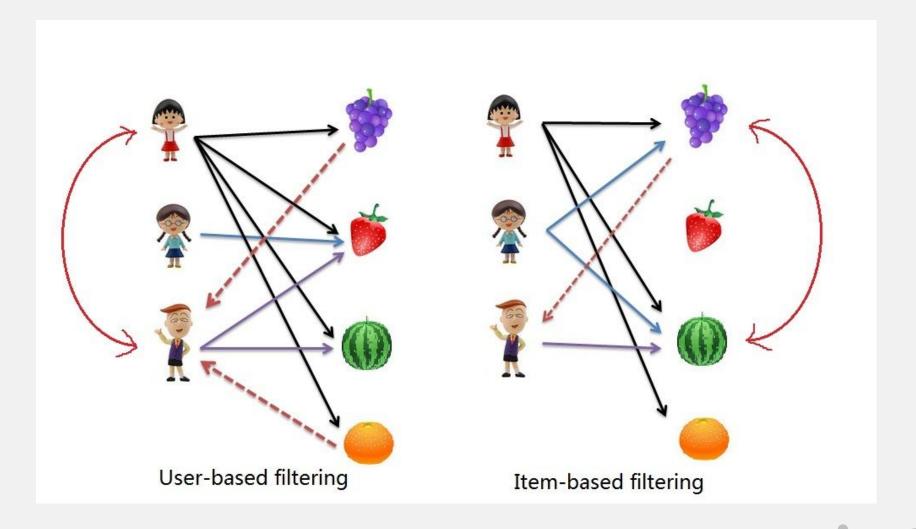
- Когда пользователей много (почти всегда), задача поиска ближайшего соседа становится плохо вычислимой;
- Оценка близости товаров гораздо более точная, чем оценка близости пользователей;
- В user-based варианте описания пользователей, как правило, сильно разрежены (товаров много, оценок мало).
 Проблема: сколько соседей не бери, список товаров, которые в итоге можно порекомендовать, получается очень небольшим;



Что лучше? item-based или user-based?

- Когда пользователей много (почти всегда), задача поиска ближайшего соседа становится плохо вычислимой;
- Оценка близости товаров гораздо более точная, чем оценка близости пользователей;
- В user-based варианте описания пользователей, как правило, сильно разрежены (товаров много, оценок мало).
 Проблема: сколько соседей не бери, список товаров, которые в итоге можно порекомендовать, получается очень небольшим;
- Предпочтения пользователя могут меняться со временем, но описание товаров штука гораздо более устойчивая.





Коллаборативная фильтрация Метрики схожести



Κοςνηγς Μερα
$$similarity = cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum\limits_{i=1}^n A_i \times B_i}{\sqrt{\sum\limits_{i=1}^n (A_i)^2} \times \sqrt{\sum\limits_{i=1}^n (B_i)^2}}$$

Коэффициент корреляции Пирсона
$$\sin(u,v) = \frac{\sum_{i} (r_{ui} - \bar{r}_{u})(r_{vi} - \bar{r}_{v})}{\sqrt{\sum_{i} (r_{ui} - \bar{r}_{u})^{2}} \sqrt{\sum_{i} (r_{vi} - \bar{r}_{v})^{2}}}$$

Евклидово расстояние

$$d(a,b) = \sqrt{\sum_{k=1}^{n} (a_k - b_k)^2}$$

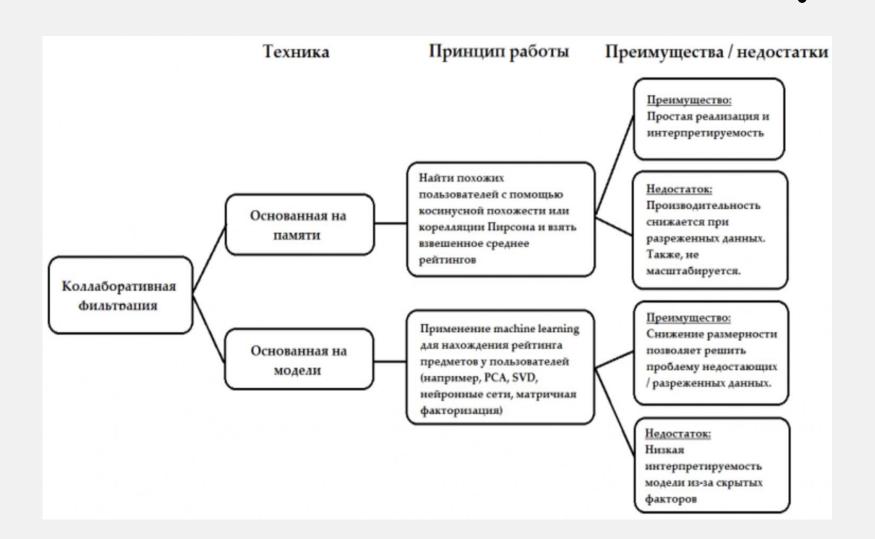
Коэффициент Жаккара

$$d(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

Манхэттенское расстояние и т.д.

$$d(a,b) = \sum_{k=1}^{n} |a_k - b_k|$$





Подход на основе памяти



Запоминаем матрицу interaction-ов, рекомендации составляются путем запроса данного пользователя к остальной части матрицы полезности.





В теореме о сингулярном разложении утверждается, что у любой матрицы \mathbf{A} размера $n \times m$ существует разложение в произведение трех матриц: \mathbf{U} , $\mathbf{\Sigma}$ и \mathbf{V} :

$$\boldsymbol{U}\boldsymbol{U}^T = \boldsymbol{I}_n, \quad \boldsymbol{V}\boldsymbol{V}^T = \boldsymbol{I}_m,$$

$$\Sigma = \mathsf{diag} \big(\lambda_1, \dots, \lambda_{\min(n,m)} \big), \quad \lambda_1 \geqslant \dots \geqslant \lambda_{\min(n,m)} \geqslant 0$$



$$\lambda_{d+1}, \dots, \lambda_{\min(n,m)} := 0$$

$$A' = U' \times \Sigma' \times V'^T$$



$$\lambda_{d+1}, \dots, \lambda_{\min(n,m)} := 0$$

$$\mathbf{A}'_{n \times m} = \mathbf{U}'_{n \times d} \times \mathbf{\Sigma}'_{d \times d} \times \mathbf{V'}^{T}_{d \times m}$$



$$\lambda_{d+1}, \dots, \lambda_{\min(n,m)} := 0$$

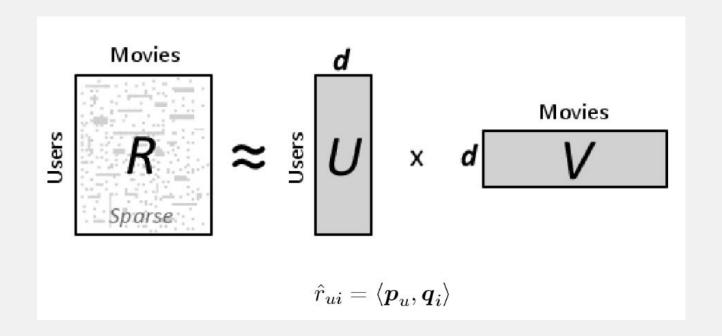
$$\mathbf{A}'_{n \times m} = \mathbf{U}'_{n \times d} \times \mathbf{\Sigma}'_{d \times d} \times \mathbf{V'}^{T}_{d \times m}$$

Полученная матрица **A'** хорошо приближает исходную матрицу **A** и, более того, является **наилучшим** низкоранговым приближением с точки зрения средне-квадратичного отклонения.

Подход на основе памяти



Одной из наиболее распространенных реализаций подхода на основе модели является **матричная факторизация**.



Скрытые представления пользователей и товаров!



$$\begin{split} \hat{r}_{ui}(\Theta) &= \boldsymbol{p}_u^T \boldsymbol{q}_i, \\ \Theta &= \{\boldsymbol{p}_u, \boldsymbol{q}_i \mid u \in U, i \in I\} \end{split}$$

<u>Нужно:</u> подобрать параметры так, чтобы на тех оценках, которые у нас уже есть, ошибка была как можно меньше

$$\mathbf{E}_{(u,i)}(\hat{r}_{ui}(\Theta) - r_{ui})^2 \to \min_{\Theta}$$

Чего-то не хватает...



$$\underbrace{\sum_{(u,i)\in\mathcal{D}} \left(\hat{r}_{ui}(\Theta) - r_{ui}\right)^2}_{\text{качество на обучающей выборке}} + \underbrace{\lambda \sum_{\theta\in\Theta} \theta^2}_{\text{регуляризация}} \to \min_{\Theta}$$

$$J(\Theta) = \sum_{(u,i) \in \mathcal{D}} (\boldsymbol{p}_u^T \boldsymbol{q}_i - r_{ui})^2 + \lambda (\sum_u \|\boldsymbol{p}_u\|^2 + \sum_i \|\boldsymbol{q}_i\|^2)$$



Как найти минимум функции, зависящей от большого количества переменных?

Правильно, нам потребуется градиент

$$\nabla J(\Theta) = \left(\frac{\partial J}{\partial \theta_1}, \frac{\partial J}{\partial \theta_2}, \dots, \frac{\partial J}{\partial \theta_n}\right)^T$$

Самый известный метод оптимизации функций — **градиентный спуск**



Для чего нужен градиент?

Градиент в какой-нибудь конкретной точке — это такой вектор, направленный в ту сторону, куда больше всего растет наша функция.

Соответственно, чтобы минимизировать наш функционал нужно..



Для чего нужен градиент?

Градиент в какой-нибудь конкретной точке — это такой вектор, направленный в ту сторону, куда больше всего растет наша функция.

Соответственно, чтобы минимизировать наш функционал нужно.. двигаться в сторону антиградиента

$$\Theta_{t+1} = \Theta_t - \eta \nabla J(\Theta)$$

Примеры!



Пример в jupyter notebook



Проблемы рекомендательных систем



При большом числе user-ов или item-ов поиск ближайших сущностей нецелесообразно осуществлять полным перебором.

Что делать?

Использовать приближенные методы поиска ближайших соседей:

HNSW, ANNOY, LSH



Желаемые свойства рекомендация



Разнообразие (diversity) ~ непохожие на другие товары из списка

Плохо: к ноутбуку только ноутбуки того же производителя

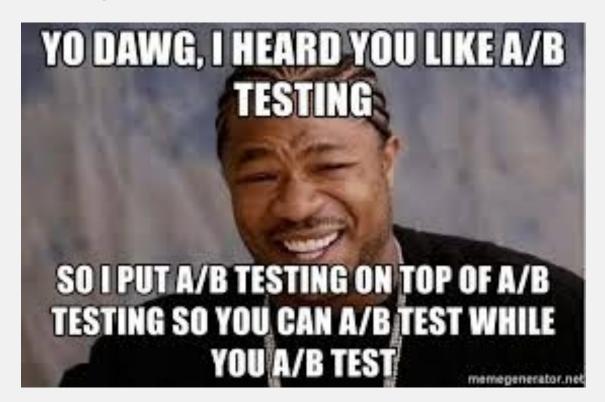
Новизна (novelty) ~ для пользователя Плохо: каждый день одно и то же

Доверие ~ обосновать рекомендацию «с товаром покупают», «скидка за комплект», ...

Оценка качества рекомендательных систем



- offline-тестирование с помощью ретротестов
- А/В тестирование





Случай непрерывных оценок

Название	Формула	Описание
MAE (Mean Absolute Error)	E(P-R)	Среднее абсолютное отклонение
MSE (Mean Squared Error)	$E(\left P-R\right ^2)$	Среднеквадратичная ошибка
RMSE (Root Mean Squared Error)	$\sqrt{E(\left P-R\right ^2)}$	Корень из среднеквадратичной ошибки



Случай дискретных оценок (бинарные)

Название	Формула	Описание
Precision	$\frac{TP}{TP+FP}$	Доля рекомендаций, понравившихся пользователю
Recall	$rac{TP}{TP+FN}$	Доля интересных пользователю товаров, которая показана
F1-Measure	$\frac{2PR}{P+R}$	Среднее гармоническое метрик Precision и Recall. Полезно, когда заранее невозможно сказать, какая из метрик важнее
ROC AUC		Насколько высока концентрация интересных товаров в начале списка рекомендаций
Precision@N		Метрика Precision, посчитанная на Тор-N записях
Recall@N		Метрика Recall, посчитанная на Тор-N записях
AverageP		Среднее значение Precision на всем списке рекомендаций



recommender system precision: $P = \frac{\text{# of our recommendations that are relevant}}{\text{# of items we recommended}}$

recommender system recall: $r = \frac{\text{# of our recommendations that are relevant}}{\text{# of all the possible relevant items}}$

P@k и R@k — это просто precision и recall, рассчитанные с учетом только подмножества рекомендаций от ранга 1 до n.

AP@N =
$$\frac{1}{m} \sum_{k=1}^{N} (P(k) \text{ if } k^{th} \text{ item was relevant}) = \frac{1}{m} \sum_{k=1}^{N} P(k) \cdot rel(k)$$

rel~(k) — это просто индикатор (0/1), который сообщает нам, был ли релевантен этот k-й элемент, а P~(k) — точность @k



В пользовательской сессии были показаны рекомендации товаров **A**, **Б** и **B**. Куплены при этом были товары **Б** и **Д**. Знатоки, внимание вопрос:





В пользовательской сессии были показаны рекомендации товаров **A**, **Б** и **B**. Куплены при этом были товары **Б** и **Д**. Знатоки, внимание вопрос:

• Чему равен **R@3** для рекомендаций в этом примере?





В пользовательской сессии были показаны рекомендации товаров **A**, **Б** и **B**. Куплены при этом были товары **Б** и **Д**. **Знатоки**, **внимание вопрос**:

- Чему равен **R@3** для рекомендаций в этом примере?
- A P@3?



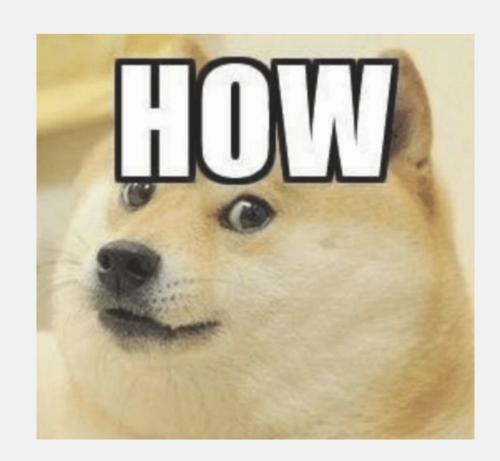


Можно придумывать свои метрики. Как?

Пример:

Есть item2item товарные рекомендации. Для каждого товара имеется 10 похожих.

Как оценить их качество оффлайн?



Условный precision



Имеем 10 рекомендаций к каждому товару. Фиксируем 1 товар.

Получаем следующую метрику:

$$pr(i) = \frac{1}{10} \frac{1}{|I_i|} \sum_{u \in I_i} |I_u \cap rec_i|,$$

где rec_i — рекомендации для товара $i,\,I_u$ — множество товаров посещенных пользователем $u,\,I_i$ — множество пользователей, имеющих взаимодействие с товаром i



Журавлёв Вадим

v.zhuravlyov@corp.mail.ru