



СУББОТНИКИ  
Мини-конференции. Омск

# Реальная жизнь рекомендательных систем

Иван Гундырев  
Аналитик  
Codenetix

# Рекомендательные системы = RS

- Основная задача - спрогнозировать, какой контент (фильмы, музыка, книги, новости, веб-сайты) будет интересен пользователю.
- Для решения этой задачи используются - профили пользователей, история взаимодействия с контентом (покупки, просмотры, закладки, поисковые запросы), информация о контенте.

## RS для онлайн-кинотеатров:

- Пользователи – подписчики сервиса, посетители сайта.

## RS для онлайн-кинотеатров:

- Пользователи – подписчики сервиса, посетители сайта.
- Контент – фильмы, сериалы, TV-программы.

## RS для онлайн-кинотеатров:

- Пользователи – подписчики сервиса, посетители сайта.
- Контент – фильмы, сериалы, TV-программы.
- Задача – удержание пользователей (чтобы им было интересно смотреть и они продлевали подписку) + увеличение продаж контента.

Кто лидер и делает Rocket Science в этой области ?



Кто лидер и делает Rocket Science в этой области ?



[Artwork Personalization at Netflix](#)

[Data Science and the Art of Producing Entertainment at Netflix](#)

# Историческое замечание

- Широко известный конкурс [Netflix Prize](#) 2006-2009 поднял интерес к RS в целом и привлек многих исследователей к задачам, возникающим при разработке RS.





# Основные идеи RS

- 1. Фильтрация на основе содержания — выбор жанров, актеров и рейтингование контента, попадающего под эти условия.

# Основные идеи RS

- 1. Фильтрация на основе содержания — выбор жанров, актеров и рейтингование контента, попадающего под эти условия.
- 2. Коллаборативная фильтрация — используется история пользователей в прошлом (покупки, просмотры).

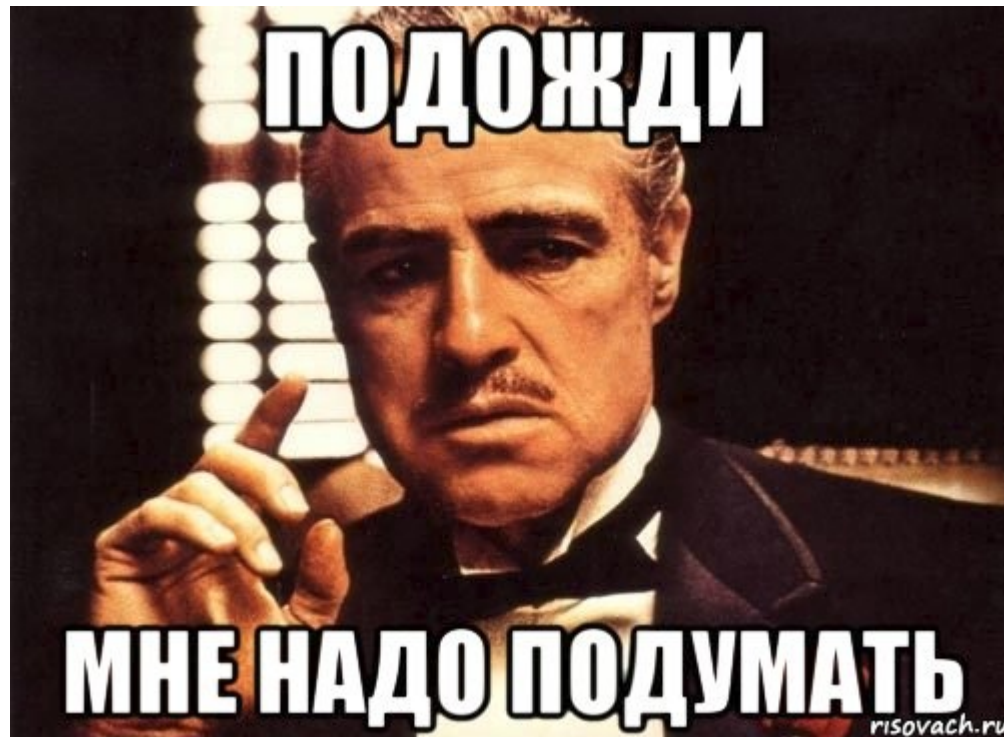
# Основные идеи RS

- 1. Фильтрация на основе содержания — выбор жанров, актеров и рейтингование контента, попадающего под эти условия.
- 2. Коллаборативная фильтрация — используется история пользователей в прошлом (покупки, просмотры).
- 3. Гибридная модель — сочетание моделей 1 и 2.

Вопрос:

Какая RS лучше работает ?

Вопрос: Какая RS лучше работает ?



Вопрос: Какая RS лучше работает ?



**I don't know.**

- Ответ: Можно понять только после сравнительного тестирования.

Вопрос: Какая RS лучше работает ?



**I don't know.**

- Ответ: Можно понять только после сравнительного тестирования.
- Хорошо организованное А/В – тестирование очень важно в переговорах с заказчиком.

Важное замечание:

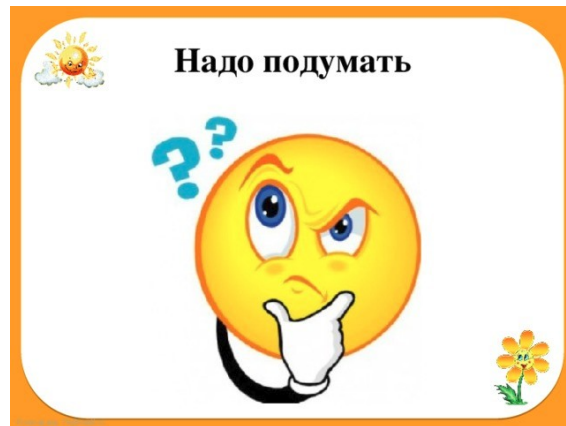


- Всегда реализуйте легкий способ смешивания (гибридизации) рекомендаций от разных моделей RS в различных пропорциях.



## Предварительный этап

- Что нужно для запуска алгоритмов ?



## Предварительный этап

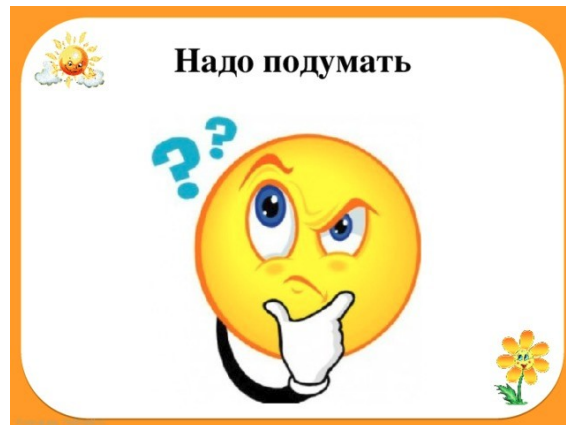
- Что нужно для запуска алгоритмов ?

Базовая таблица:

Строка – id пользователя,

Столбец – id контента,

Значение элемента таблицы – это  
оценка контента пользователем.



Чуть – чуть математики :-).

Обозначим через  $U$  – множество пользователей в системе,  
 $V$  – весь известный системе контент,  
 $M$  – возможные оценки контента (подмножество вещественных чисел),  
 $D$  – подмножество  $U \times V$  такое, что для любой пары  $(u, v) \in D$  известна оценка.  
 $NA$  – неизвестное значение.

Тогда базовая таблица RS:  $f: U \times V \rightarrow M \cup \{NA\}$  и  $f(D) \subset M$

Обычно для многих пар  $(u, v) \in U \times V$  значение  $f(u, v)$  неизвестно. Поскольку пользователь просматривал (покупал) лишь небольшую часть контента в системе.

**Задача RS** – по известному отображению  $f$  спрогнозировать оценки для таких пар (дозаполнить таблицу). И предложить контент с max оценками пользователю.

## Вся мощь линейной алгебры в помощь!

- Для решения задачи заполнения пропусков в базовой таблице один из самых широко используемых методов – классический SVD (Singular-value decomposition).

**SVD**

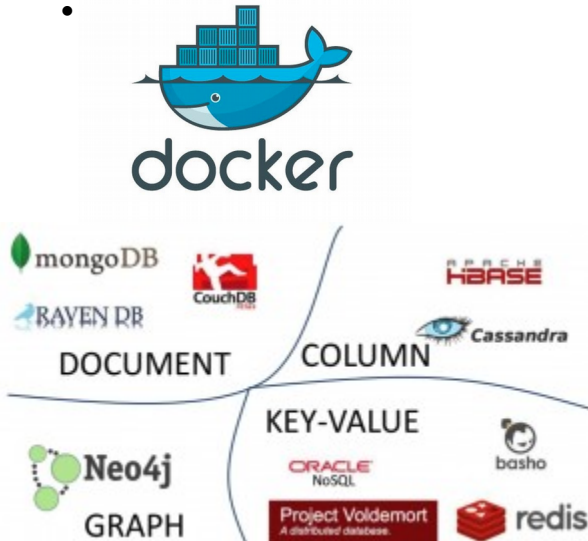
$$X_{[n \times m]} = U_{[n \times r]} S_{[r \times r]} (V_{[m \times r]})^T$$

$$\begin{array}{c} X \\ \left( \begin{array}{cccc} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & \\ \vdots & \vdots & \ddots & \\ x_{m1} & & & x_{mn} \end{array} \right) \\ m \times n \end{array} = \begin{array}{c} U \\ \left( \begin{array}{ccc} u_{11} & \dots & u_{1r} \\ \vdots & \ddots & \\ u_{m1} & & u_{mr} \end{array} \right) \\ m \times r \end{array} \begin{array}{c} S \\ \left( \begin{array}{ccc} s_{11} & 0 & \dots \\ 0 & \ddots & \\ \vdots & & s_{rr} \end{array} \right) \\ r \times r \end{array} \begin{array}{c} V^T \\ \left( \begin{array}{ccc} v_{11} & \dots & v_{1n} \\ \vdots & \ddots & \\ v_{r1} & & v_{rn} \end{array} \right) \\ r \times n \end{array}$$

- $X$ :  $m \times n$  matrix (e.g.,  $m$  users,  $n$  videos)
- $U$ :  $m \times r$  matrix ( $m$  users,  $r$  concepts)
- $S$ :  $r \times r$  diagonal matrix (strength of each 'concept') ( $r$ : rank of the matrix)
- $V$ :  $r \times n$  matrix ( $n$  videos,  $r$  concepts)

# Немного о технологиях

- Почти наверняка будет использоваться: облачная инфраструктура, docker контейнеры, NoSQL базы данных, системы мониторинга – Grafana.



# Что же делать?

- Как жить, дядя Мить?



# Что же делать?

- Как жить, дядя Мить?
- 
- 
- Кадры, овладевшие техникой, решают всё!





## Кино для взрослых ≠ Adult Movie

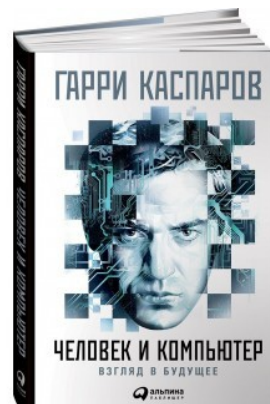


- 1) Какой фильм скрывается под знаком вопроса ?
- 2) Какой фильм “лишний” в этой подборке ?



# Познавательные (не технические) книги

- Эрик Сигель “Просчитать будущее. Кто кликнет, купит, соврет или умрет”.
- Сет Стивенс-Давидовиц “Все лгут. Поисковики, Big Data и Интернет знают о вас все”.
- Гарри Каспаров “Человек и компьютер. Взгляд в будущее”.

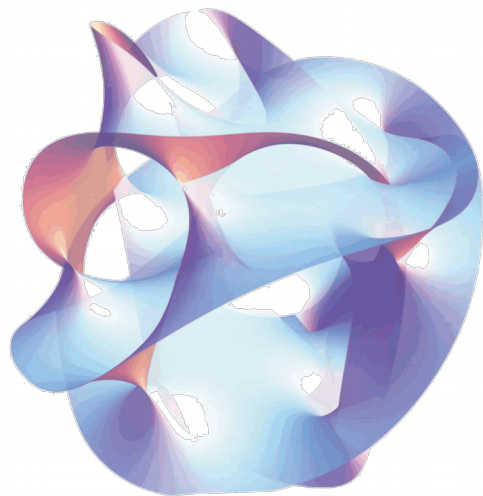


## Познавательные (технические + много математики) видео

- [Разработка рекомендательной системы для интерактивного телевидения](#)
- [Тестирование рекомендательных систем, или как проверить то - не знаю что](#)
- [Машинное обучение в Дзене](#)
- [Построение рекомендательной системы на Python](#)
- [Какие задачи решает команда рекомендаций в Avito](#)

## Пара слов о работе с унаследованным кодом

Картинки на которые можно смотреть пока докладчик, что-то рассказывает



# Вопросы ?

Гундырев Иван, Аналитик,  
Codenetix



You have to make the  
good out of the bad  
because that is all you  
have got to make it out of

Robert Penn Warren

PICTUREQUOTES.COM



PICTUREQUOTES

