



北京师范大学 珠海校区
BEIJING NORMAL UNIVERSITY AT ZHUHAI

前深度学习时代

马静



北京師範大學 珠海校区

BEIJING NORMAL UNIVERSITY AT ZHUHAI

CONTENT

- 01 传统推荐模型的演化关系
- 02 协同过滤
- 03 矩阵分解法
- 04 补充知识

We can read of things that happened
5,000 years ago in the Near East,
where people first learned to write.
But there are some parts of the world
where even now people cannot write.



北京师范大学 珠海校区
BEIJING NORMAL UNIVERSITY AT ZHUHAI



01 传统推荐模型的 演化关系



基于协同过滤算法
的推荐系统

1992

LR, FM, GBDT

2000

深度学习
推荐系统

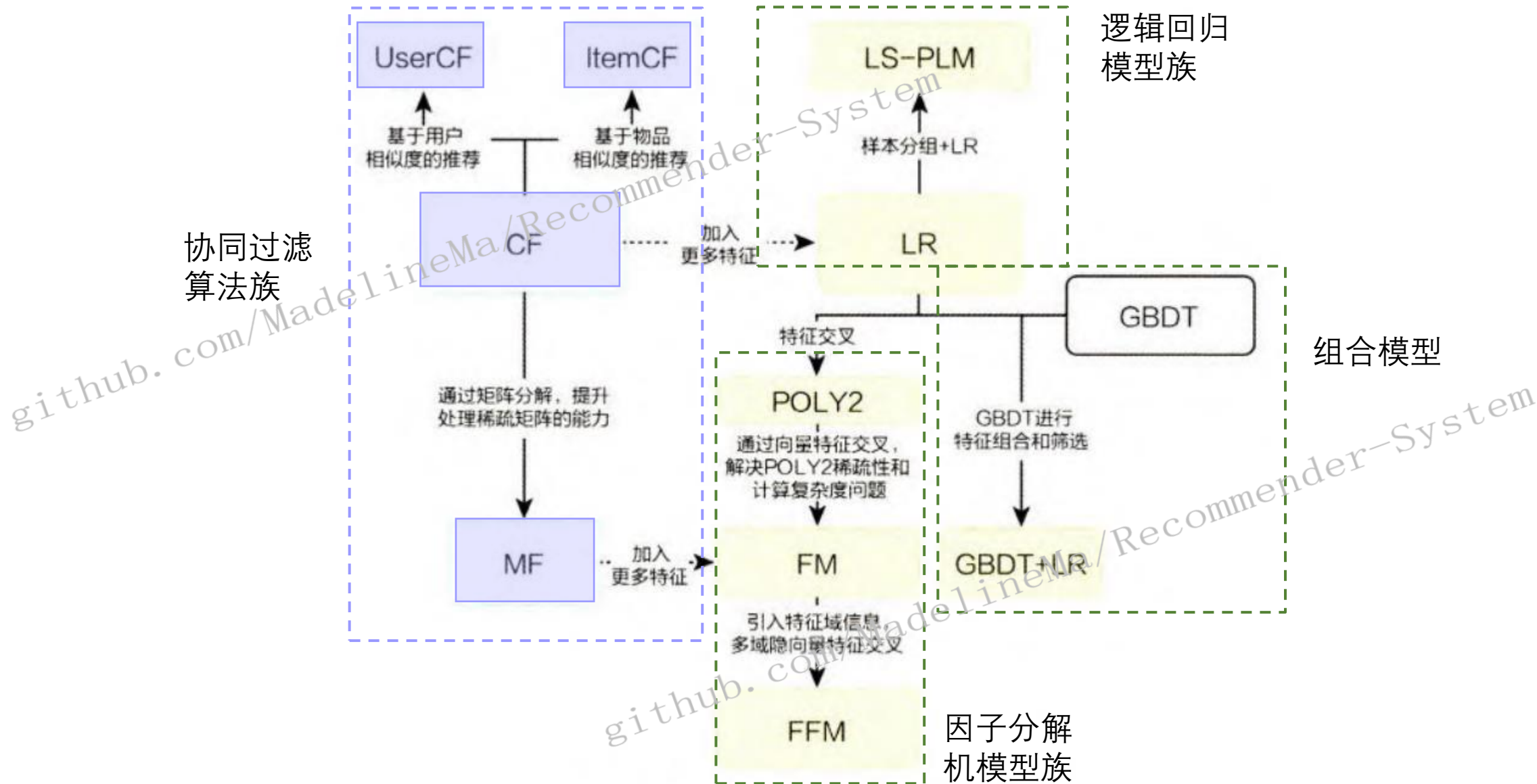
2015

可解释性强,
硬件环境要求低,
易于快速训练和部署

深度学习的基础
召回的前身



传统推荐模型的演化关系

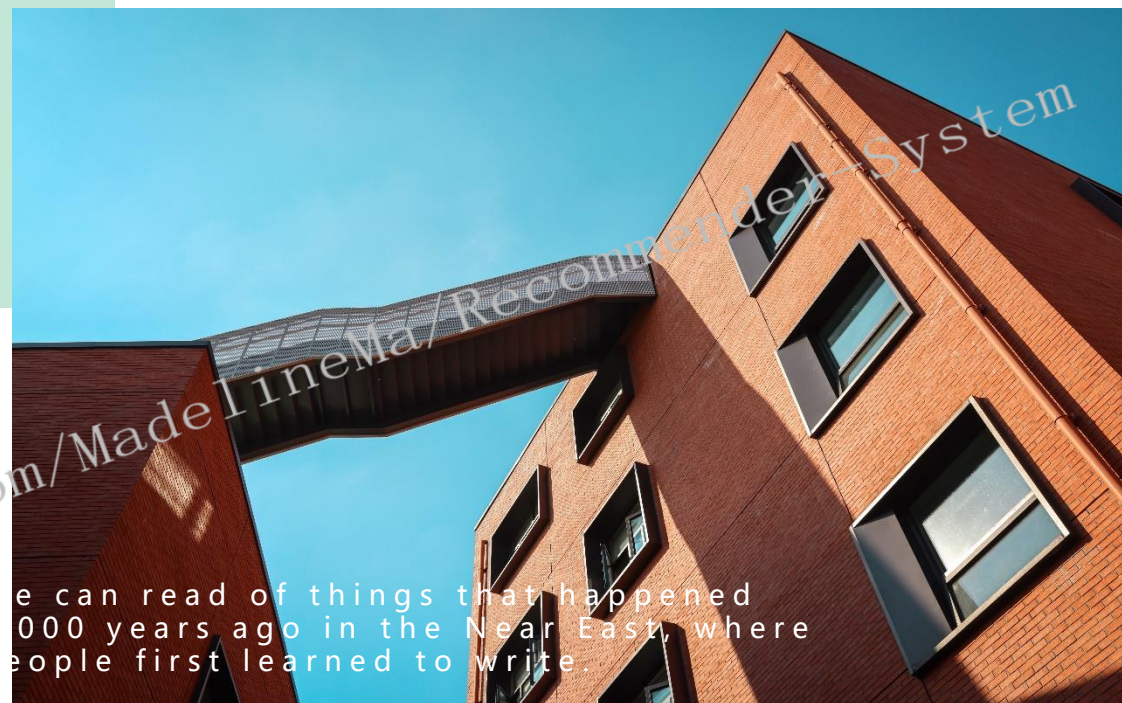




北京师范大学 珠海校区
BEIJING NORMAL UNIVERSITY AT ZHUHAI

02

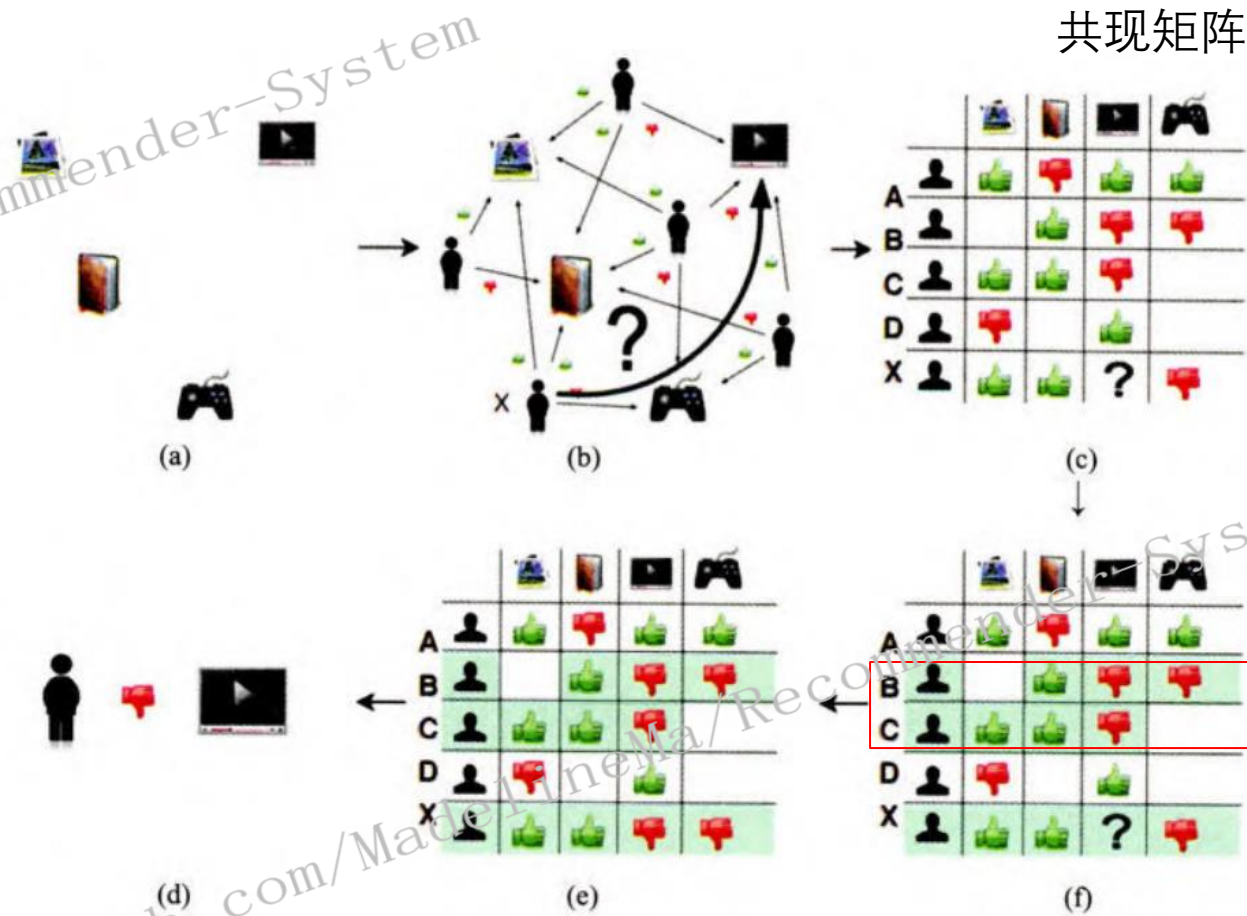
协同过滤-经典的推荐算法



We can read of things that happened
1000 years ago in the Near East, where
people first learned to write.



- 1992年，应用于邮件筛选系统，过滤不感兴趣的无用邮件；
- 2003年，Amazon应用于推荐系统；
- 定义：协同大家的反馈、评价和意见一起对海量的信息进行过滤，从中筛选出目标用户可能感兴趣的信息的推荐过程；
- 用户相似度计算与最终结果的排序不严谨。





用户相似度计算 -> 向量之间相似度计算

余弦相似度

$$\text{sim}(i, j) = \cos(i, j) = \frac{i \cdot j}{\|i\| \cdot \|j\|}$$

用户打分纠偏的皮尔逊相似度

$$\text{sim}(i, j) = \frac{\sum_{p \in P} (R_{i,p} - \bar{R}_i)(R_{j,p} - \bar{R}_j)}{\sqrt{\sum_{p \in P} (R_{i,p} - \bar{R}_i)^2} \sqrt{\sum_{p \in P} (R_{j,p} - \bar{R}_j)^2}}$$

物品打分纠偏的皮尔逊相似度

$$\text{sim}(i, j) = \frac{\sum_{p \in P} (R_{i,p} - \bar{R}_p)(R_{j,p} - \bar{R}_p)}{\sqrt{\sum_{p \in P} (R_{i,p} - \bar{R}_p)^2} \sqrt{\sum_{p \in P} (R_{j,p} - \bar{R}_p)^2}}$$



最终结果排序

假设：目标用户与其相似用户的喜好是相似的

- 目标货品池打分：

$$R_{u,p} = \frac{\sum_{s \in S} (w_{u,s} \cdot R_{s,p})}{\sum_{s \in S} w_{u,s}}$$

- 对打分结果进行排序，选取高分货品进行推荐。

- ✗ UserCF 需要维护用户相似度矩阵表以便快速找出TopN相似用户，存储以及存储的增长对在线存储系统提出挑战；
- ✗ 历史数据稀疏，对冷启动用户以及反馈难获取的场景并不适用（酒店预订，大件商品低频购买应用）

阿里解决方式：离线计算Top200相似用户喜欢的相似货品，存为在线表格，以便实时快速查询

userId	SimiUserId	itemId	info
1110	2314	11348	...





货品相似度计算->共现矩阵的列向量相似度

1. 基于历史数据，构建以用户（假设用户总数为 m ）为行坐标，物品（物品总数为 n ）为列坐标的 $m \times n$ 维的共现矩阵。
2. 计算共现矩阵两两列向量间的相似性（相似度的计算方式与用户相似度的计算方式相同），构建维的物品相似度矩阵。
3. 获得用户历史行为数据中的正反馈物品列表。
4. 利用物品相似度矩阵，针对目标用户历史行为中的正反馈物品，找出相似的TopA个物品，组成相似物品集合。
5. 对相似物品集合中的物品，利用相似度分值进行排序，生成最终的推荐列表。



UserCF 与 ItemCF应用差异:

- 面向的对象: 快速得知与自己兴趣相似的人最近喜欢什么? UserCF: 新闻, 微博, 热点发现.
一段时间内更倾向于某类货品. ItemCF: 电商, 视频推荐.

CF的缺点:

- 泛化能力差, 热门商品以及反馈多的商品容易更热门, 新品和小众的商品快速沉底
(长尾问题以及处理稀疏向量能力弱)

$$\begin{matrix} A[0 & 0 & 0 & 1 & 1 & 0 & 1 & 0 & 1] \\ B[0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0] \\ C[0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0] \\ D[1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 1] \end{matrix} \Rightarrow \begin{matrix} A & B & C & D \\ \begin{matrix} A \\ B \\ C \\ D \end{matrix} \begin{bmatrix} - & 0.00 & 0.00 & 0.71 \\ 0.00 & - & 0.00 & 0.18 \\ 0.00 & 0.00 & - & 0.18 \\ 0.71 & 0.18 & 0.18 & - \end{bmatrix} \end{matrix}$$

- 仅利用用户和物品的交互信息, 无法有效引入用户基本信息, 商品基本信息, 上下文特征.



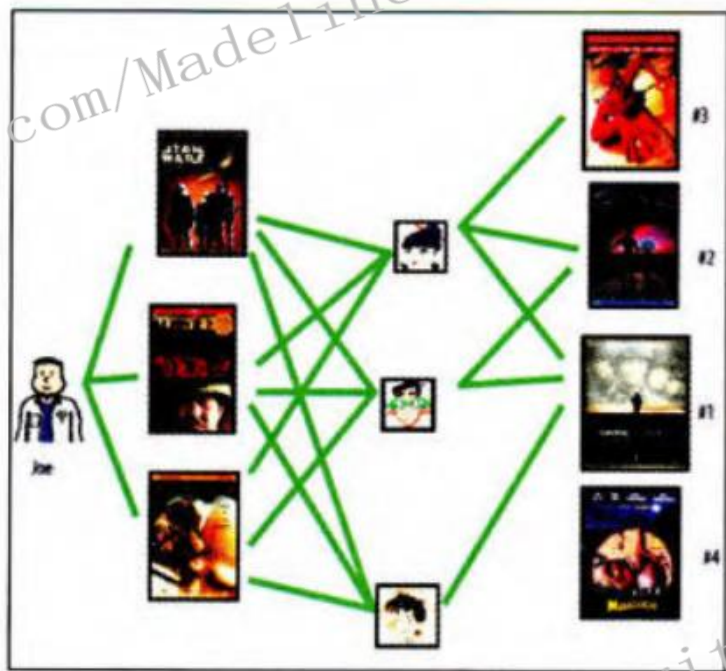
北京师范大学 珠海校区
BEIJING NORMAL UNIVERSITY AT ZHUHAI

03 矩阵分解法

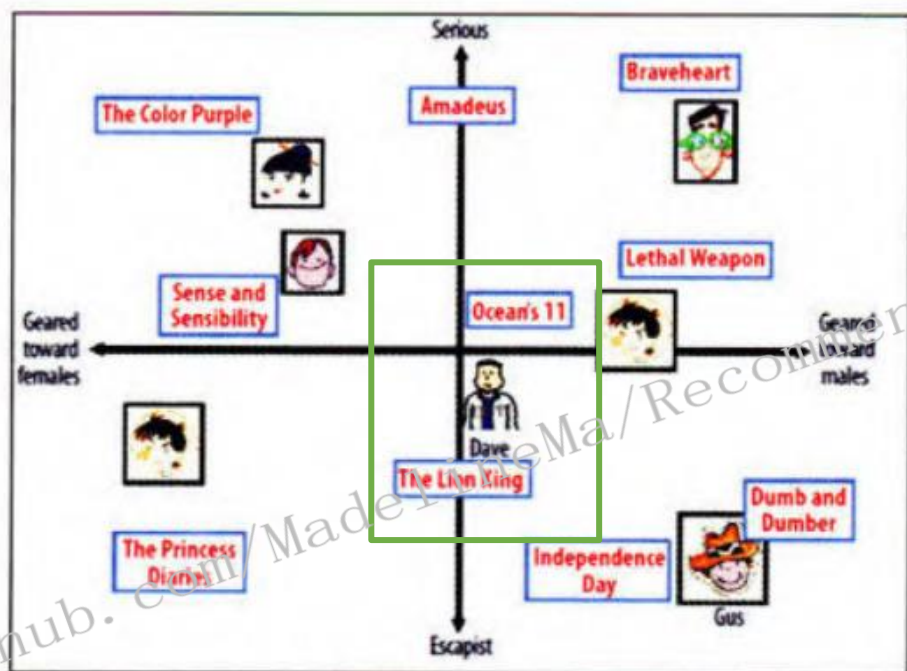
We can read of things that happened 5,000 years ago in the Near East, where people first learned to write.



- CF的缺点：泛化能力差，长尾问题以及处理稀疏向量能力弱；
- 2006年，Netflix Prize Challenge中，MF算法大放异彩；
- 隐向量的引入；



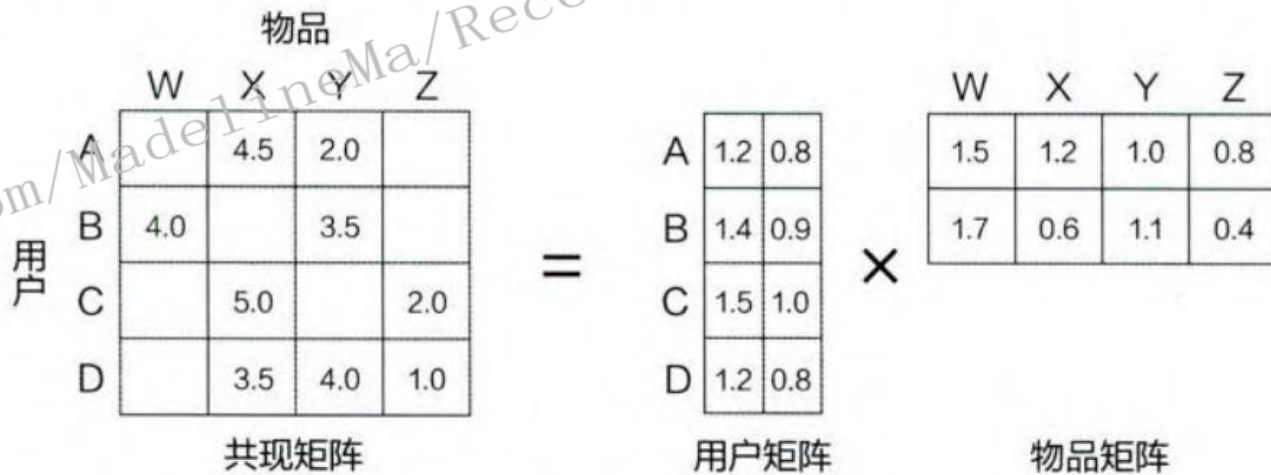
(a) 协同过滤算法原理图



(b) 矩阵分解算法原理图



- 目标：用隐向量表达用户和物品，还要保证相似的用户及用户可能喜欢的物品的距离相近。
- 用户和物品的隐向量是通过分解协同过滤生成的共现矩阵得到的。



- 将 $m \times n$ 维的共现矩阵及分解为 $m \times k$ 维的用户矩阵 U , $k \times n$ 维的物品矩阵 V 相乘的形式。
- k 的大小决定了隐向量表达能力的强弱和泛化能力强弱？
- 最终打分 $\hat{r}_{ui} = \mathbf{q}_i^T \mathbf{p}_u$

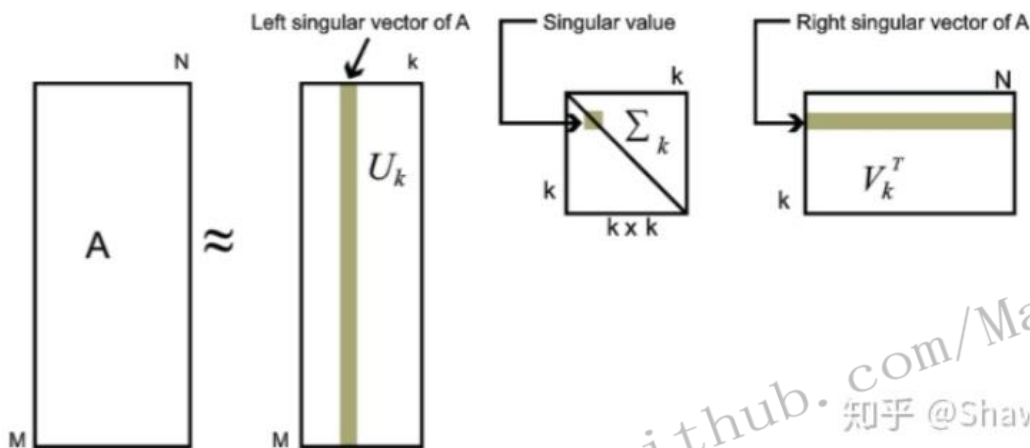


× 特征值分解：只能用于方阵，不适用于互联网场景，why?

× 奇异值分解：要求原始的共现矩阵是稠密的，复杂度为 $O(mn^2)$ ，不适用，Why?

假设矩阵 M 是一个 $m \times n$ 的矩阵，则一定存在一个分解 $M = U \Sigma V^T$ ，其中 U 是 $m \times m$ 的正交矩阵， V 是 $n \times n$ 的正交矩阵， Σ 是 $m \times n$ 的对角阵。

取对角阵 Σ 中较大的 k 个元素作为隐含特征，删除其他维度及 U 和 V 中对应的维度，矩阵 M 被分解为 $M \approx U_{m \times k} \Sigma_{k \times k} V_{k \times n}^T$



U_k : 用户矩阵的第 k 列被对角阵的第 k 行 k 列的元素进行整行放缩;

U_m : 取用户矩阵第 m 行隐向量，每个元素 j 逐个被对角阵的 j 行 j 列元素放缩，再与货品矩阵的列相乘。



✓ 梯度下降法：让原始评分与用户向量和物品向量之积的差尽量小

$$\min_{\mathbf{q}^*, \mathbf{p}^*} \sum_{(u,i) \in K} (r_{ui} - \mathbf{q}_i^T \mathbf{p}_u)^2$$

$$\min_{\mathbf{q}^*, \mathbf{p}^*} \sum_{(u,i) \in K} (r_{ui} - \mathbf{q}_i^T \mathbf{p}_u)^2 + \lambda(\|\mathbf{q}_i\| + \|\mathbf{p}_u\|)^2$$

正则化的overfitting, Why?

✓ 面试常见题：什么是overfitting？哪些降低overfitting的方法？L1和L2优缺点是什么？两种方法在收敛的角度看差异在哪里？L2好在哪里？

1) 对用户求偏导，得到的结果为 $2(\mathbf{r}_{ui} - \mathbf{q}_i^T \mathbf{p}_u) \mathbf{p}_u - 2\lambda \mathbf{q}_i$

2) 对货品求偏导，得到的结果为 $2(\mathbf{r}_{ui} - \mathbf{q}_i^T \mathbf{p}_u) \mathbf{q}_i - 2\lambda \mathbf{p}_u$

3) 更新梯度： $\mathbf{p}_u \leftarrow \mathbf{p}_u - \gamma ((\mathbf{r}_{ui} - \mathbf{q}_i^T \mathbf{p}_u) \mathbf{p}_u - \lambda \mathbf{q}_i)$ $\mathbf{q}_i \leftarrow \mathbf{q}_i - \gamma ((\mathbf{r}_{ui} - \mathbf{q}_i^T \mathbf{p}_u) \mathbf{p}_u - \lambda \mathbf{q}_i)$

4) 终止条件: n次或损失的变化几乎不变.



- ✓ 梯度下降法：让原始评分与用户向量和物品向量之积的差尽量小
- ✓ 将正则化项加入损失函数来保持模型稳定的做法也可以做如下理解。对于加入了正则化项的损失函数来说，模型权重的值越大，损失函数越大。梯度下降是朝着损失（Loss）小的方向发展的，因此正则化项其实是希望在尽量不影响原模型与数据集之间损失的前提下，使模型的权重变小，权重的减小自然会让模型的输出波动更小，从而达到让模型更稳定的目的。



MF解决了CF的缺点?

- CF的缺点：泛化能力差，长尾问题以及处理稀疏向量能力弱；
- 2006年，Netflix Prize Challenge中，MF算法大放异彩；
- 隐向量的引入，使任意的用户和物品之间都可以得到预测分值；
- 而隐向量的生成过程其实是对共现矩阵进行全局拟合的过程，因此隐向量其实是利用全局信息生成的，有更强的泛化能力；
- 而对协同过滤来说，如果两个用户没有相同的历史行为，两个物品没有相同的人购买，那么这两个用户和两个物品的相似度都将为0. Why?

✓ 泛化能力强.

✓ 空间复杂度低，存储低.

✓ Embedding思想的前身

x 不方便加入用户、商品、上下文的基本信息进行综合考虑.



1. 正则化的资料阅读与深入思考.
2. CF和MF的资料阅读和深入思考.
3. Embedding课题小组可趁热打铁组织报告或论文.

Reference:

[L1正则化与L2正则化 - 知乎 \(zhihu.com\)](#)

[l1 相比于 l2 为什么容易获得稀疏解? - 知乎 \(zhihu.com\)](#)

[最小推荐系统: Neural Collaborative Filtering \(NeuralCF\) - 知乎 \(zhihu.com\)](#)

[万物皆Embedding, 从经典的word2vec到深度学习基本操作item2vec - 知乎 \(zhihu.com\)](#)

[一个值得讨论的问题:word2vec与SVD/LSA等的关系 - 知乎 \(zhihu.com\)](#)

[\[NLP\] 秒懂词向量Word2vec的本质 - 知乎 \(zhihu.com\)](#)



北京師範大學 珠海校区

BEIJING NORMAL UNIVERSITY AT ZHUHAI

THANKS

DESIGNED BY 2xh