

→ 一般稀疏

1. 效用矩阵 *utility matrix* : user 和 item, 每行对应元素值代表当前用户对当前项喜欢程度.

		item		
		H_1	H_2	H_3
user	A	1	0	1
	B	.	.	.
	C	.	.	.
	D	.	.	.

基于内容过程:

- ① 建立项模型, 用一系列特征向量代表项. → 使用 TF-IDF 打分, 挑选分高的.
- ② 项模型表示: 构建(特征-值)成对的项模型, 用效用矩阵中每一行构建反映用户偏好的用户模型.
- ③ 构建用户模型.

1. 推荐系统的三种方法:

- 1. content-based 基于内容
- 2. collaborative 协作式
- 3. latent factor based 基于潜在因素 → 无需其他用户数据

→ 关注项的属性, 项之间的相似度通过计算它们属性间的相似度确定.

① content-based: 向用户X推荐与用户X评价较高的商品类似的商品.

eg: 推荐同演过的电影.

模型

→ 项模型

eg: $\begin{cases} \text{movie: 演员, 标题, 导演} \\ \text{text: 一系列关键词} \end{cases}$

item profile: 项模型, 用于代表该项的重要特征的一条或多条记录.

① 文档特征: 将具有最高 TF-IDF 得分的 n 个词作为一篇文档的特征. 或词语

② 度量两文档相似度: Δ 文档词之间 Jaccard 距离. Δ 向量间余弦距离.

① TF-IDF:

feature

item

t_{ij} : term i 在 doc j 中的频率

$$TF_{ij} = \frac{t_{ij}}{\max_k t_{kj}}$$

Doc profile: 一系列 TF-IDF 分数

$$IDF_i = \log \frac{N}{n_i}$$

n_i : 含 i 的 doc 数目. 最高的词.

N : doc 总数

$$\text{TF-IDF score: } w_{ij} = TF_{ij} \times IDF_i$$

④ user profile: 关联 item profile 的加权平均

user profile $\rightarrow X$ item profile $\rightarrow i$

$$u(X, i) = \cos(X, i) = \frac{X \cdot i}{\|X\| \cdot \|i\|}$$

→ 余弦距离
$$\frac{\sum_{i=1}^n x_i \cdot y_i}{\sqrt{\sum_{i=1}^n x_i^2} \cdot \sqrt{\sum_{i=1}^n y_i^2}}$$

② collaborative Filtering. 找到与用户A 相似度高的k个人, 将k个人喜欢的推荐.

计算相似度: → 适用于数据只有0,1

Jacard 距离: $\left\{ \begin{array}{l} J \text{ 相似度: } \frac{\text{交集大小}}{\text{并集大小}} \\ J \text{ 距离: } 1 - \frac{\text{交集大小}}{\text{并集大小}} \end{array} \right.$ 关注集合而不关注具体数值

余弦距离: $\frac{r_i \cdot r_j}{\|r_i\| \cdot \|r_j\|} = \cos(r_i, r_j)$ cos 值越大代表角度越小, 越近.

归一化: 每一个值要减平均值后计算. → 比较适合复杂如评分

皮尔逊系数:
$$r = \frac{\sum (x - \bar{x}) \cdot (y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2} \cdot \sqrt{\sum (y - \bar{y})^2}}$$

预测: r_x : User x's ratings 的 vector.
(对项目).

N : 与 x 最相似的 k 个用户的集合, 这些用户对项目 i 进行评分.

User-user

$$r_{xi} = \frac{1}{k} \sum_{y \in N} r_{yi} \quad \rightarrow \text{平均}$$

$$s_{xy} = \text{sim}(x, y)$$

$$r_{xi} = \frac{\sum_{y \in N} s_{xy} \cdot r_{yi}}{\sum_{y \in N} s_{xy}} \quad \rightarrow \text{加权相似度.}$$

两种

item-item

估计 i 的 rating by 与 i 相似的 item.

$$r_{xi} = \frac{\sum_{j \in N(i;x)} s_{ij} \cdot r_{xj}}{\sum_{j \in N(i;x)} s_{ij}}$$

$N(i;x)$: x 用户中相似于 i 的集合.
就是 nearest neighbors.

常归一化(归一系数).

$$r_{xi} = b_{xi} + \frac{\sum_{j \in N(i;x)} s_{ij} \cdot (r_{xj} - b_{xj})}{\sum_{j \in N(i;x)} s_{ij}}$$

baseline

$$b_{xi} = \mu + b_x + b_i$$

μ : overall mean
 b_x : 用户 x 的评分偏差
(用户 x 评分平均数 - μ)
 b_i : movie i 的评分偏差

评估: RMSE

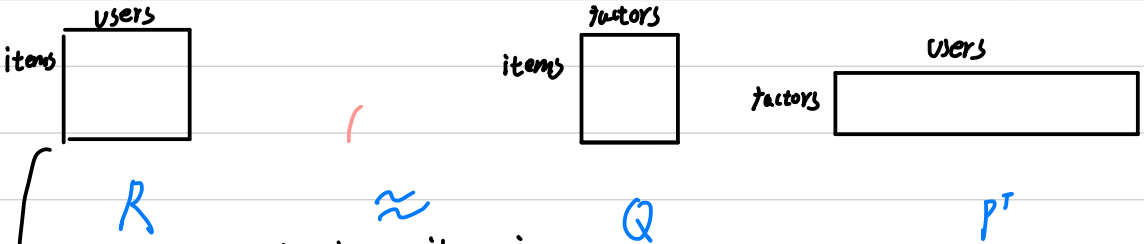
$$\sqrt{\sum_{x,i} (r_{x,i} - \hat{r}_{x,i}^*)^2} : r_{x,i} \text{ 为预测, } \hat{r}_{x,i}^* \text{ 为真实.}$$

$r_{2,4}$

Latent Factor Model (潜在因素):

eg SVD.

$$r_{2,4} = Q_{2,1} \cdot P_{1,4}$$

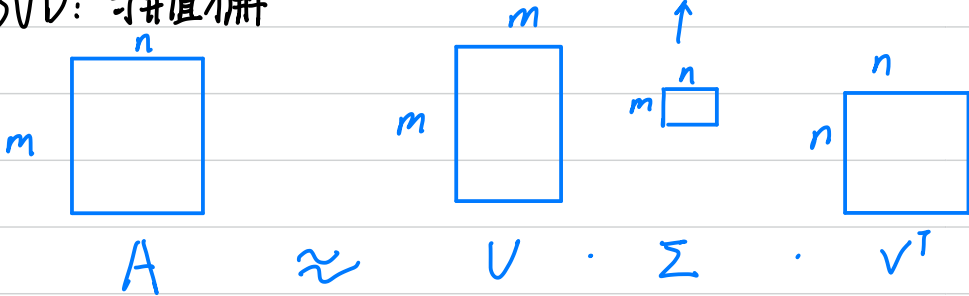


user x for item i .

$$\text{估计缺失值: } \hat{r}_{x,i} = q_i \cdot p_x = \sum_j q_{i,j} \cdot p_{x,j}$$

q_i : Q 的行 i
 p_x : P^T 的列 x

SVD: 奇异值分解



在使用中 $A=R$

$Q=U$

$P^T = \Sigma V^T$

$$(A \approx Q P^T) \\ (r_{x,i} \approx q_i \cdot p_x)$$

找 P, Q : $\min_{P,Q} \sum_{(i,x) \in R} (r_{x,i} - q_i \cdot p_x)^2 \rightarrow \text{让其最小}$

解决过拟合问题

↓

$$\min_{P, Q} \underbrace{\sum_{\text{train}} (r_{xi} - q_i \cdot p_x)^2}_{\text{error}} + \underbrace{\left[\lambda_1 \sum_x \|p_x\|^2 + \lambda_2 \sum_i \|q_i\|^2 \right]}_{\text{length}}$$

↓
用户设定正则化参数.

考虑用户 bias 和电影 bias, 有 $r_{xi} = \mu + \underbrace{b_x}_{\text{user bias}} + \underbrace{b_i}_{\text{movie bias}} + q_i \cdot p_x$

补充: 效用矩阵拆分的原理: 考虑到用户一般可能只对很小一部分特征感兴趣, 故将低维度的矩阵以分析特征. 它们在非零元素上会很相近.

$M = UV$ \rightarrow 且乘积的值为0可以预测M中空值

RMSE: 衡量矩阵拆分后, UV 和 M 的相近程度.

计算过程: ① M 与 UV , 计算每一个元素差的平方, 求和.

② 开根号

③ $\frac{1}{|R|}$

UV分解过程:

思路: 选择任意 U, V . 反复调整 U 与 V 使 RMSE 越来越小.

过程:

① 矩阵 M 的预处理

② U 和 V 的初始化

③ U 和 V 元素优化的排序

预处理: $\left\{ \begin{array}{l} \text{每个 } m_{ij} \text{ 减去用户 } i \text{ 评分的平均值, 再减去项 } j \text{ 的平均值.} \\ \text{每个 } m_{ij} \text{ 减去用户 } i \text{ 和项 } j \text{ 平均评分的均值.} \end{array} \right.$

梯度的计算: X 为输入, θ 为模型参数, y 为输出.

梯度是一个向量, 每个分量都表示损失函数对相应参数的偏导.

$$g_{\theta} = \text{np.dot}(X.T, \text{np.dot}(X, \theta) - y) / \text{len}(y)$$

SGD vs GD: $\left\{ \begin{array}{l} \text{GD 用总样本平均值来更新参数} \\ \text{SGD 随机选一个样本更新参数.} \end{array} \right.$

PPT里 r_{xi} 为 $\begin{array}{c} \text{user} \\ \boxed{} \\ \text{movie} \end{array}$ 矩阵中用户 x 对 movie i 的评分.
 \rightarrow 防止过拟合.

η 为学习率用控制步长. λ 为正则化参数 ∇ 为梯度.

$\nabla Q = [\Delta_{q-i}] \rightarrow$ 偏导数是一个向量, 存放每个元素偏导.

$[q-i]$ 代表每一个 i 的偏导集合.