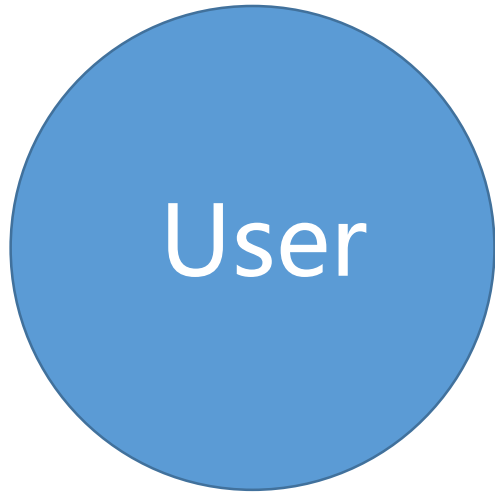


ARCF Introduction

--Xinpeng Guo

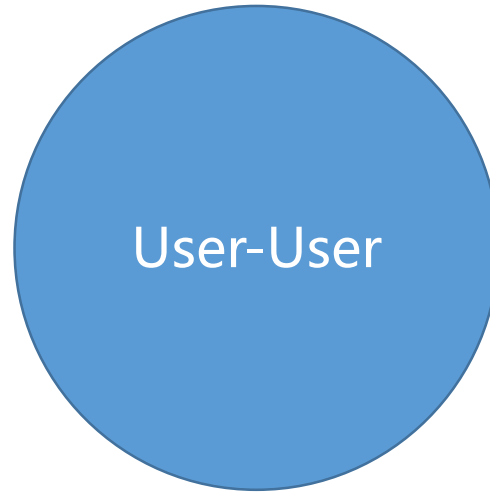
Basic Concepts

Two Basic Entity

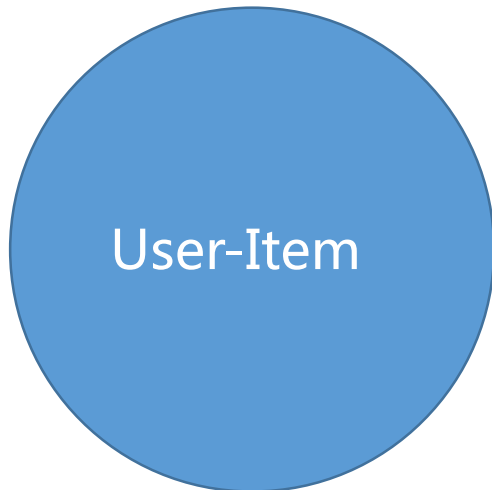


1. Song
2. Singer
3. Song List
4. MV
5. Album
- ...

Three Basic Relationships



1. Circle
2. Discover New People (Famous People)
3. People U may know (Social Relation Chain)

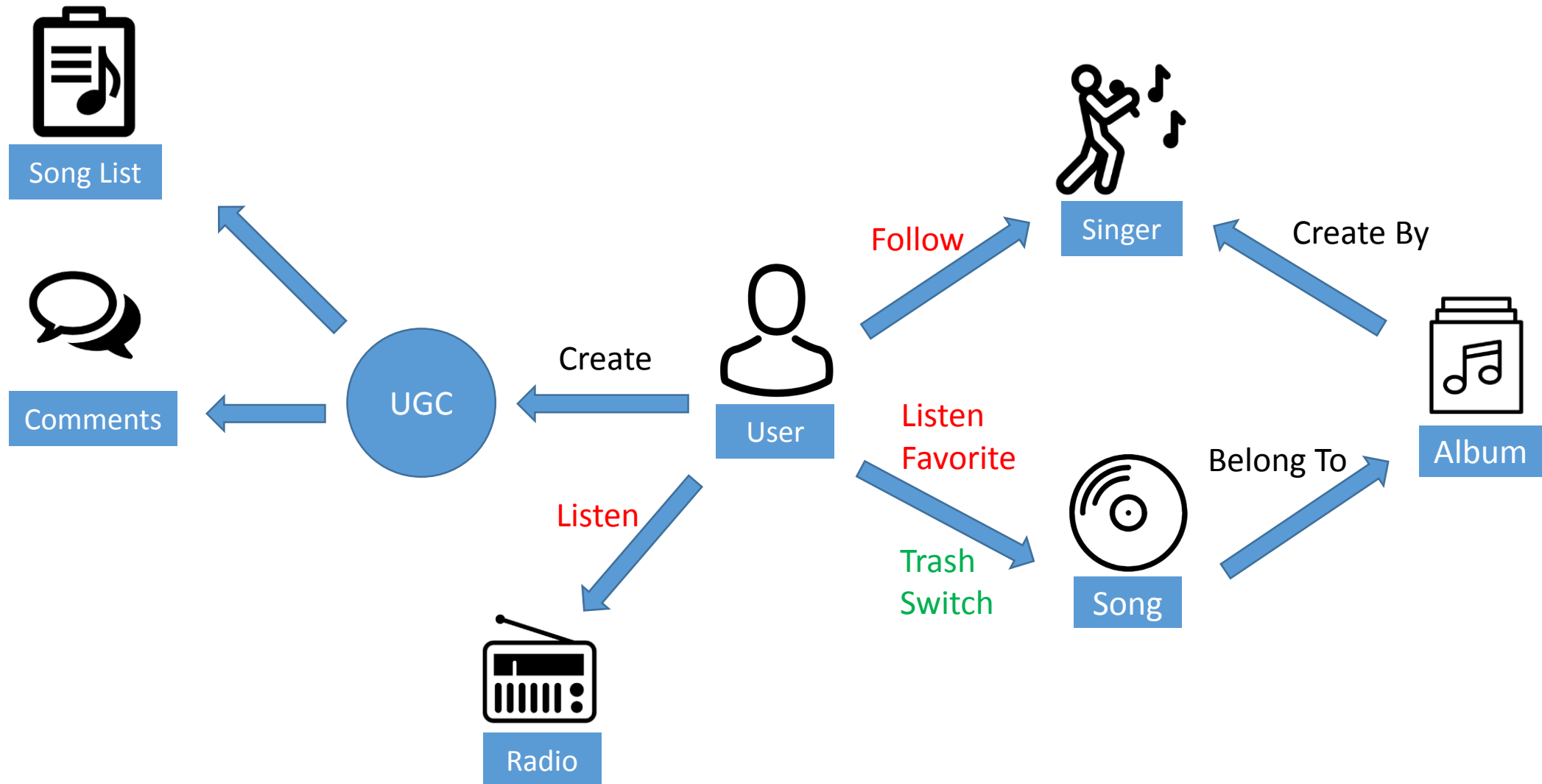


Guess U like



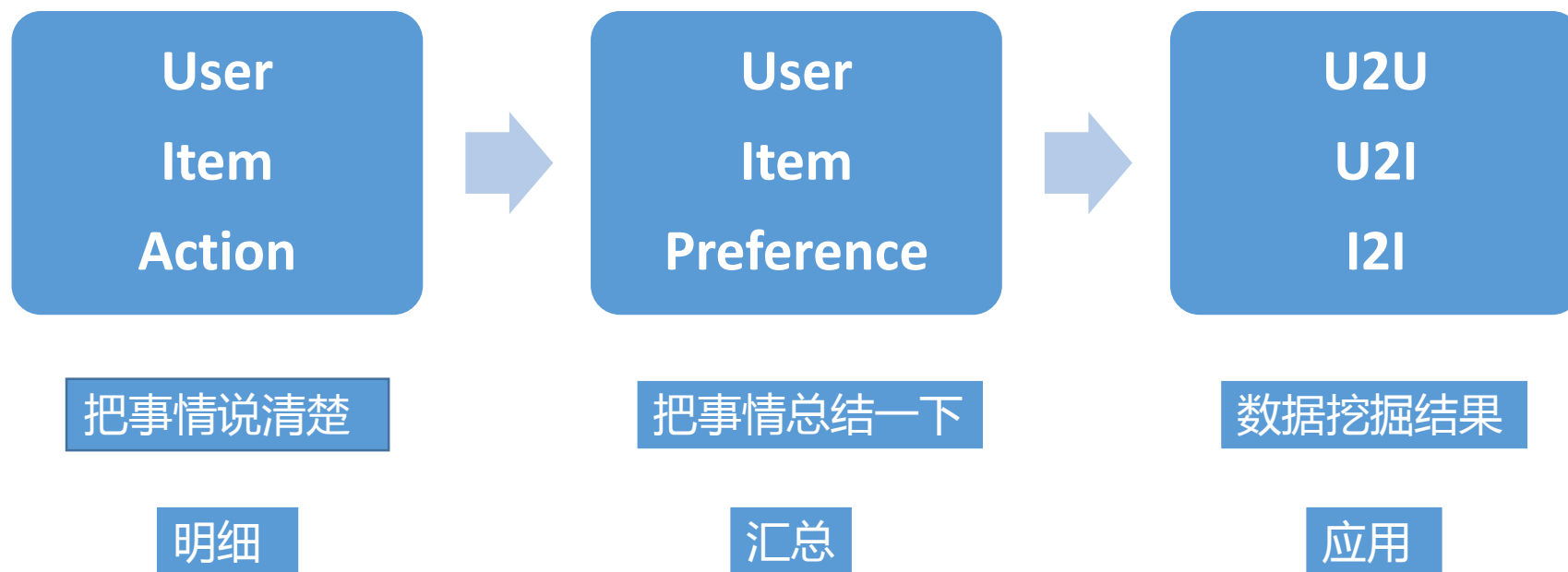
1. Similar Songs
2. Single Song Radio
3. Buy This Also Buy
4. View This Also View

Knowledge Graph



Data Flow

Three Steps



UIA=Who+Where+When+What

t_7d_imusic_user_item_action			
Field	Type	Instruction	Example
user_id	string	用户ID	12345
item_id	string	物品ID	54321
action	int	动作枚举值，1=点击，2=下载，3=收藏，4=删除，5=观看	2
location	int	位置，1000001=首页banner，100002=榜单页新歌榜	1000001
timestamp	bigint	linux时间戳	1504784445
extend	string	扩展位置（如观看时长长度等），分号分隔	60
Partition	Type	Instruction	Sample
ds	string	日期，格式：YYYY-MM-DD	"2017-9-7"
item_type	int	物品类型，1=歌曲，2=歌手，3=歌单，4=mv	1

三级分隔符	分号
	逗号
	冒号

UIP = Who + How + What

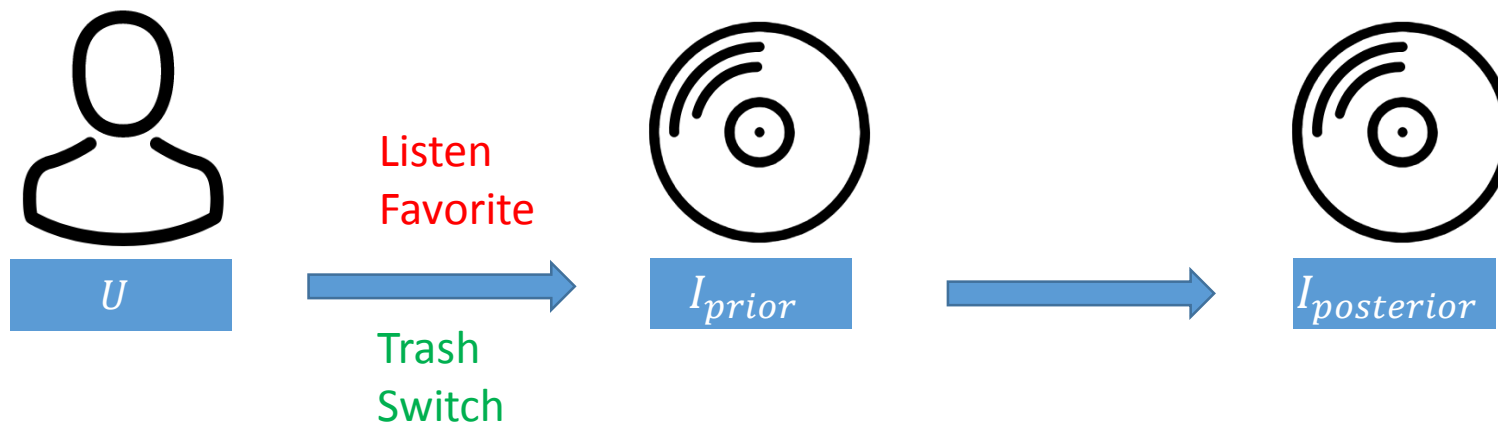
Real Matrix

t_7d_imusic_user_item_preference			
Field	Type	Instruction	Example
user_id	string	用户ID	12345
item_id	string	物品ID	54321
preference	double	分值 (0-5分)	1.432
Partition	Type	Instruction	Sample
ds	string	日期，格式：YYYY-MM-DD	"2017-9-7"
item_type	int	物品类型，1=歌曲，2=歌手，3=歌单，4=mv	1

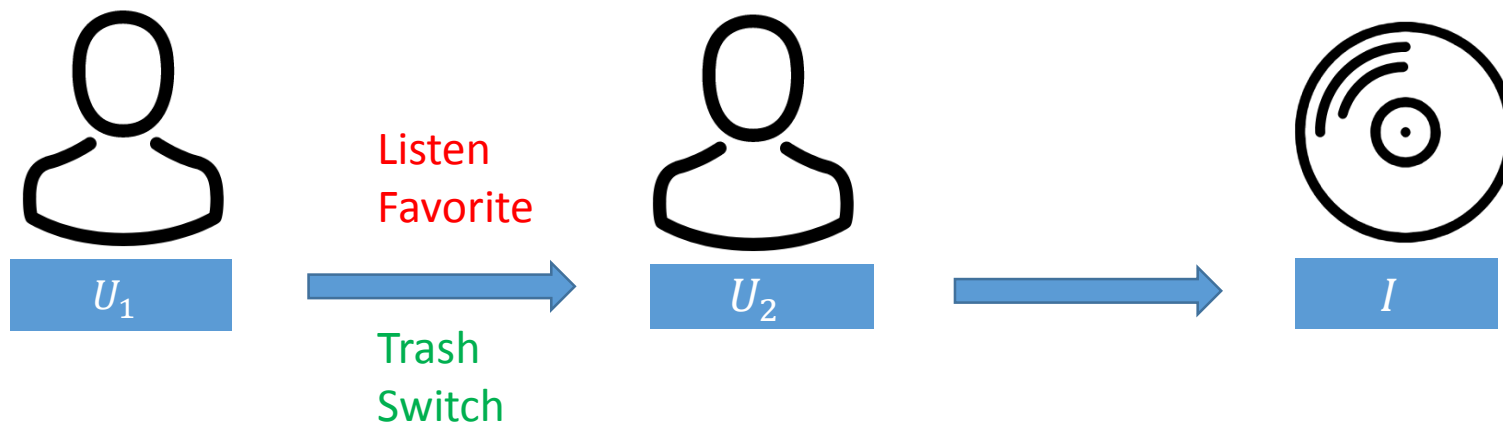
Bool Matrix

t_7d_imusic_user_item_preference_bool_matrix			
Field	Type	Instruction	Example
user_id	string	用户ID	12345
item_id	string	物品ID	54321
preference	int	1喜欢，0不喜欢/不知道	1.432
Partition	Type	Instruction	Sample
ds	string	日期，格式：YYYY-MM-DD	"2017-9-7"
item_type	int	物品类型，1=歌曲，2=歌手，3=歌单，4=mv	1

$$U2I_{posterior} = UI_{prior} + I_{prior}2I_{posterior}$$



$$U_1 \rightarrow I = U_1 \rightarrow U_2 + U_2 \rightarrow I$$



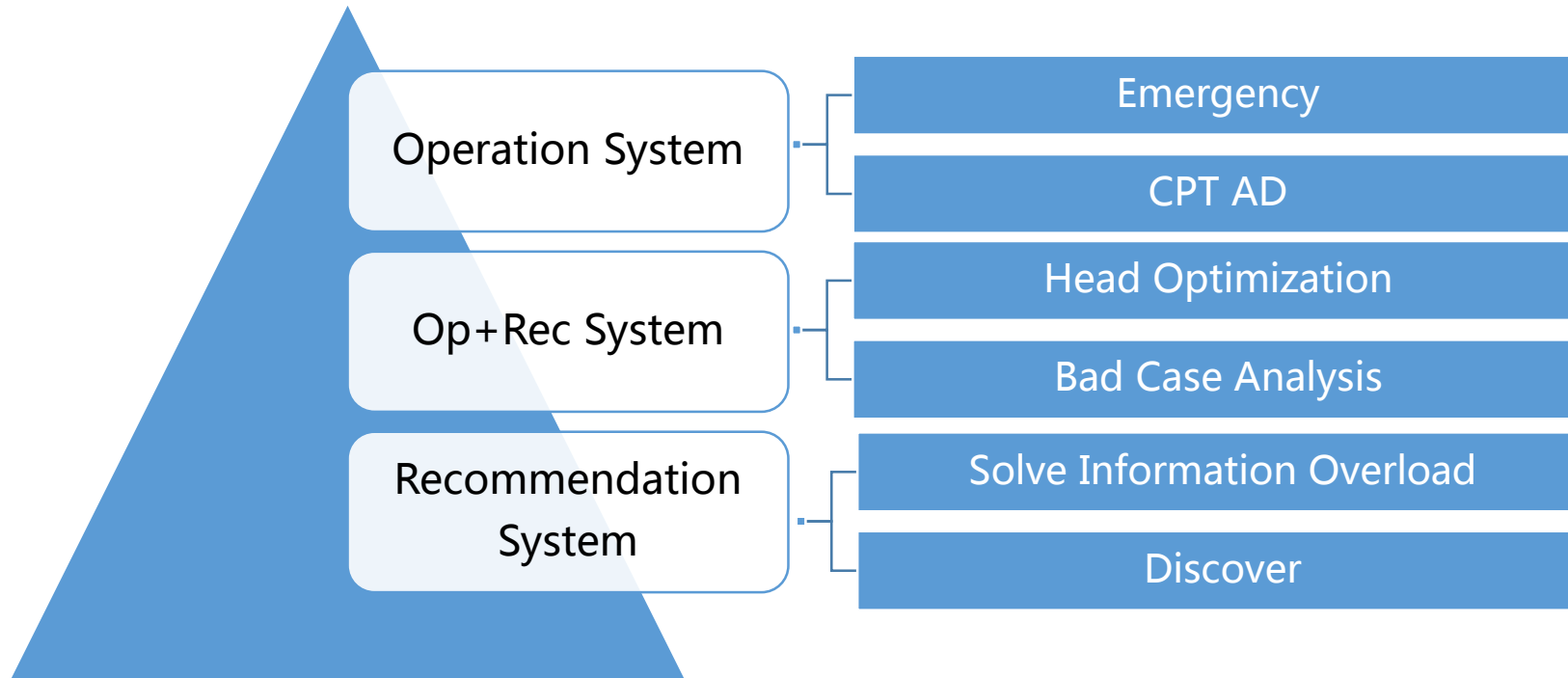
App Use
Relation Chain

Two Thoughts

- Action Based : CF
- Content Based : Profile + pCTR

Architecture

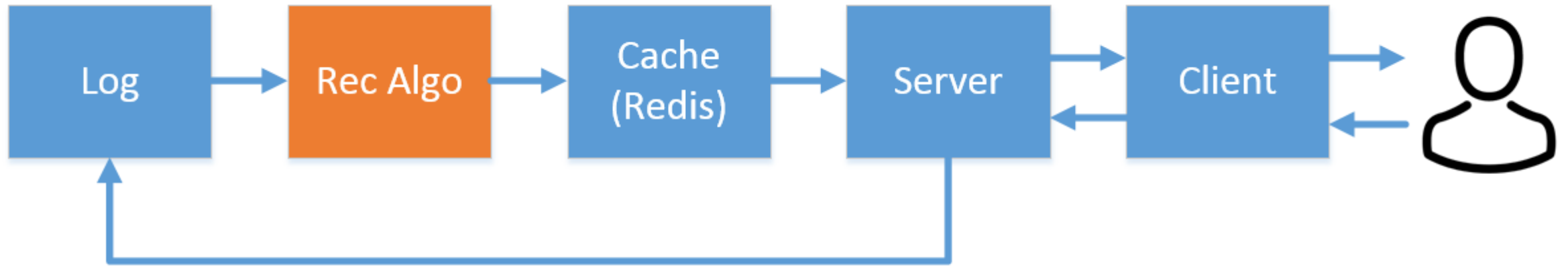
Three Layers



Framework

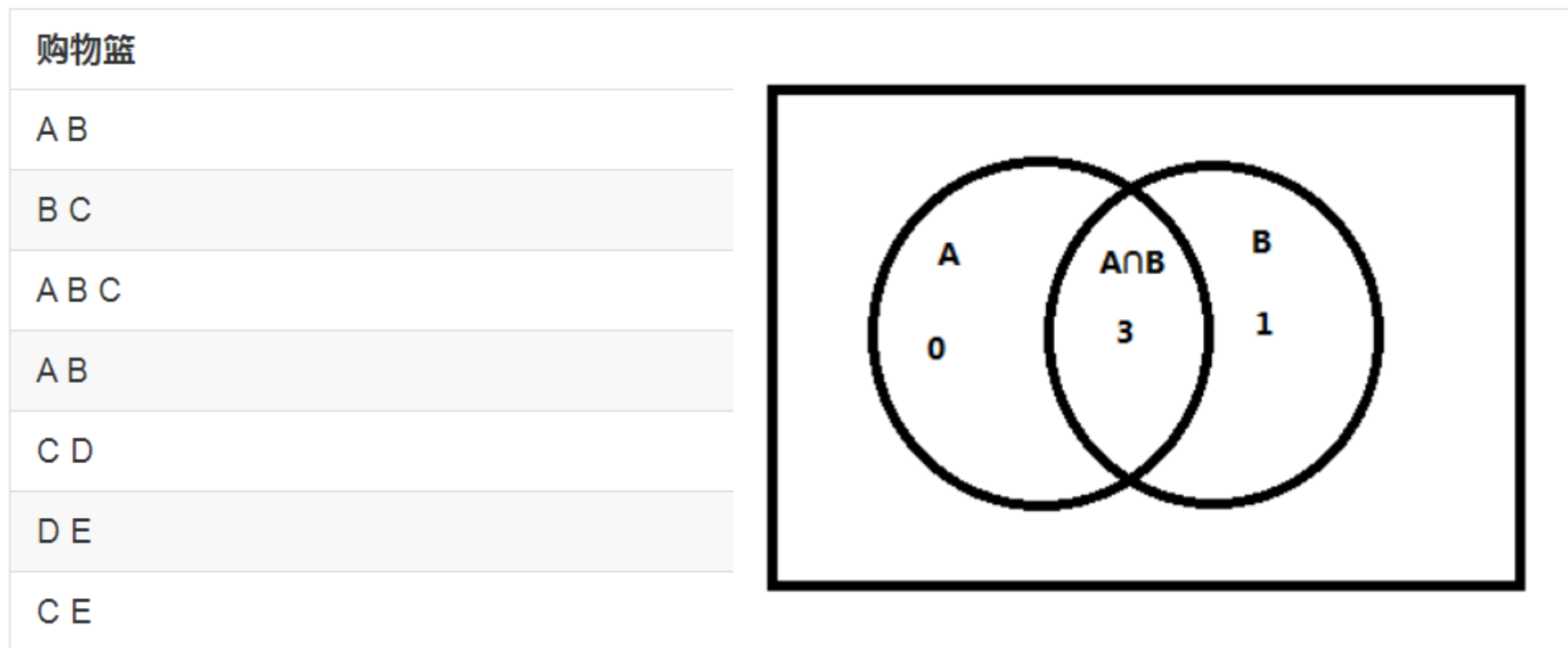


Data Closed Loop

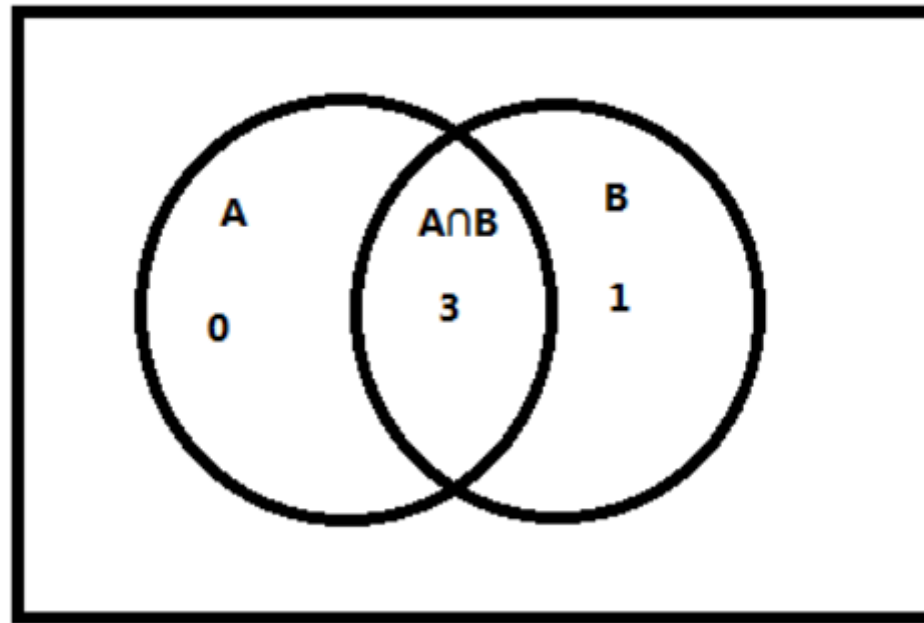


AR—Bool Matrix

Basket Analysis



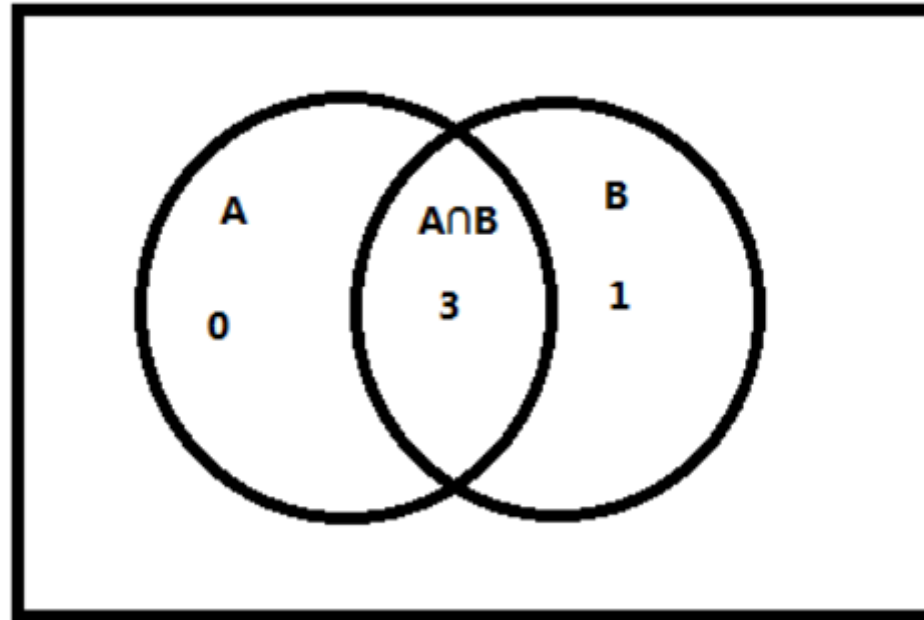
AR Metrics—Pair Frequency



$$\text{Count}(A, B) = 3$$

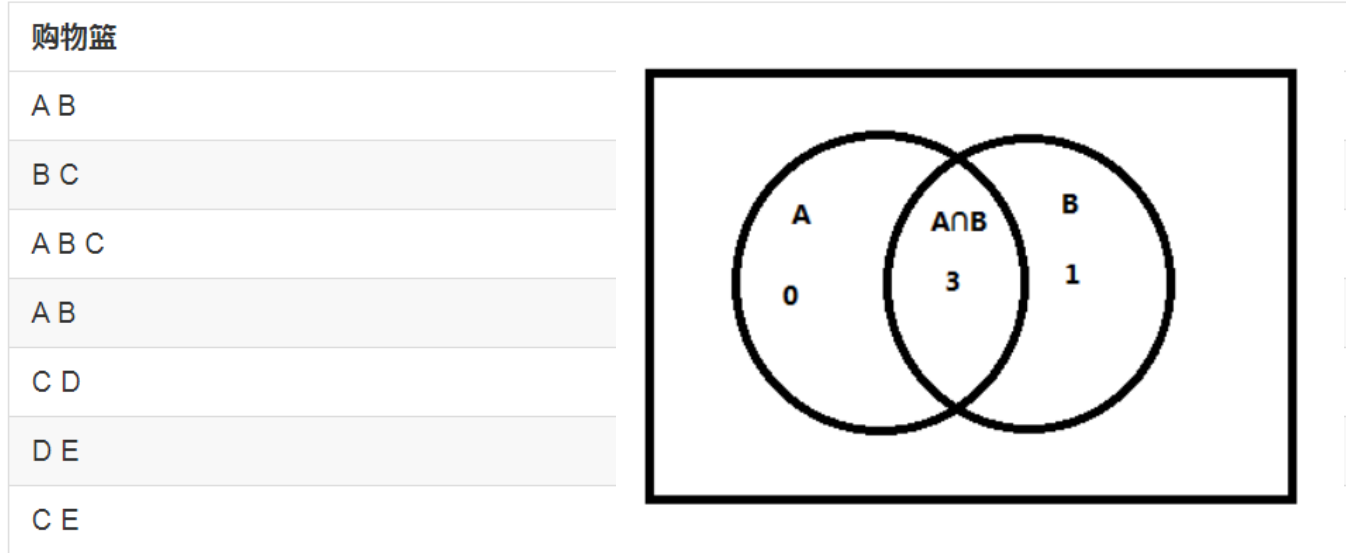


AR Metrics—Jaccard



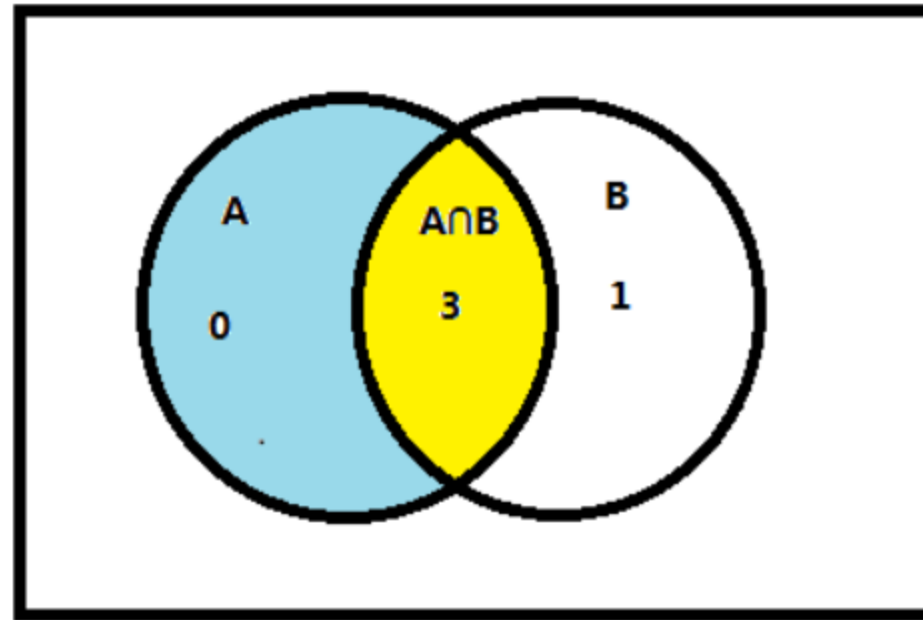
$$\text{Jaccard}(A, B) = \text{count}(A \cap B) / \text{count}(A \cup B) = 3/4$$

AR Metrics—Support



$$\text{Support (A, B)} = \text{count}(A \cap B) / \text{count}(\text{ALL}) = 3 / 7$$

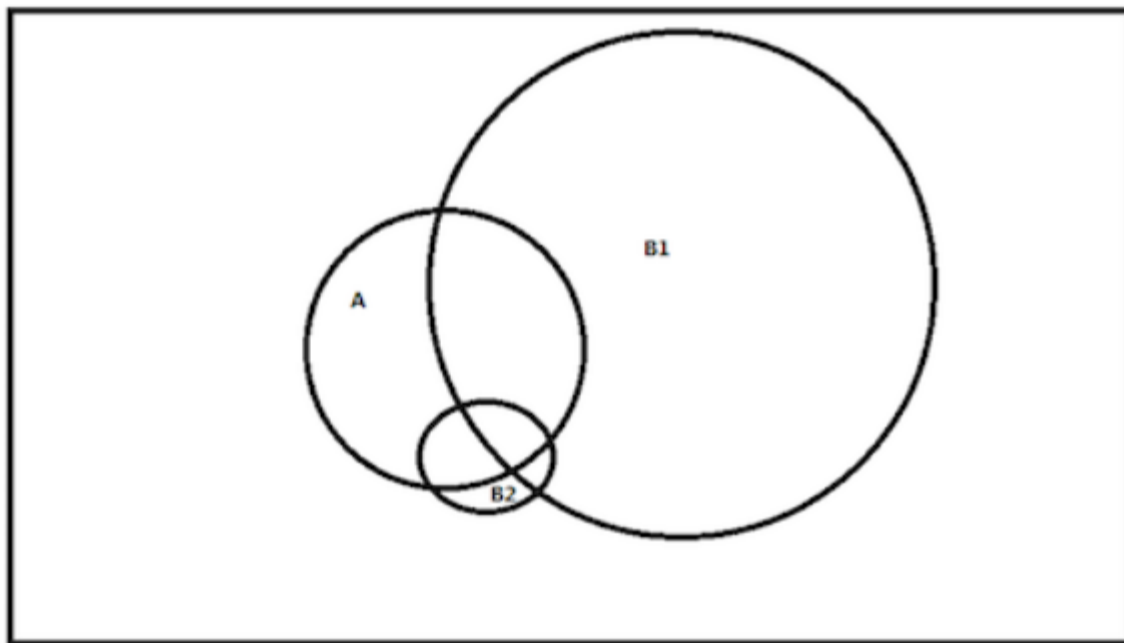
AR Metrics—Confidence



$$\text{Confidence}(A \rightarrow B) = P(A \cap B | A) = P(\text{yellow} | \text{yellow} + \text{blue}) = 1.0$$

AR Metrics—Confidence Pain Point

【置信度的痛点】 但是推荐过程中，使用会出现偏热的现象，是因为后验B太热门导致的。

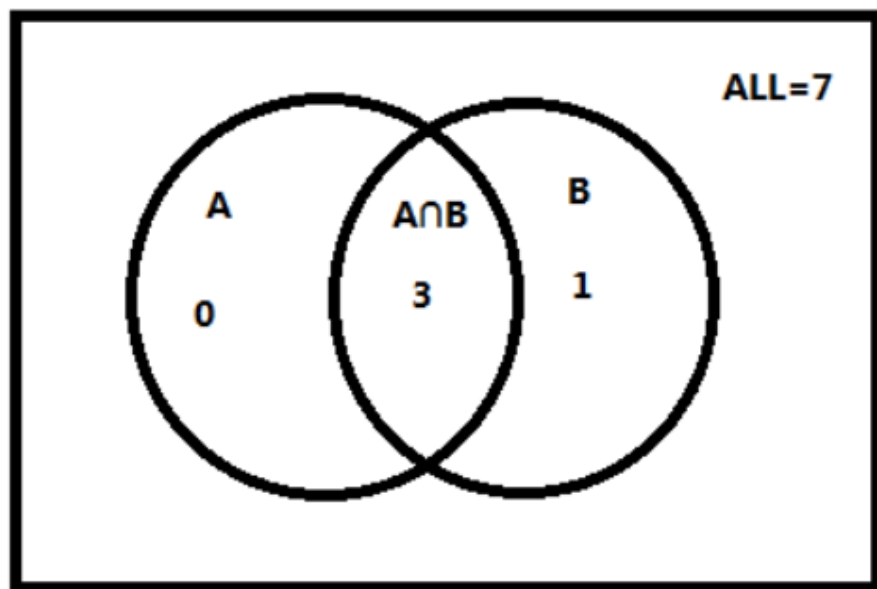


A是用户的先验（比如是手管的app）； B1是一款热门软件（比如微信）； B2是一款相对冷门的软件（比如同步助手）；

可以发现由于B1的体量太大，导致 $\text{conf}(B1|A)$ 很大； 由于B2的体量较小，导致 $\text{conf}(B2|A)$ 偏小； 但是明显应该推荐B2，因为B2都快和A重合进去了，应该对B1的体量进行惩罚。

AR Metrics—Lift

【提升度】提升度是置信度B1体量惩罚的其中一种方法，含义为：



$$\text{lift}(B|A) = \text{conf}(B|A) / P(B) = p(A \cap B) / (P(A) * P(B)) = \text{count}(A, B) * ALL / (\text{count}(A) * \text{count}(B)) = (3 * 7) / (3 * 4) = 21/12$$

其中，lift等于1，表示先验知识A的知道与否对B的概率没有影响。lift大于1，表示促进作用；lift小于1，表示抑制作用；

发现， $\text{lift}(A|B) = \text{lift}(B|A)$ ，即lift是对称的。

AR Metrics—KULR/IR

Given two itemsets, A and B , the **Kulczynski** measure of A and B (abbreviated as **Kulc**) is defined as

$$Kulc(A, B) = \frac{1}{2}(P(A|B) + P(B|A)). \quad (6.11)$$

“Among the all_confidence, max_confidence, Kulczynski, and cosine measures, which is best at indicating interesting pattern relationships?”

To answer this question, we introduce the **imbalance ratio (IR)**, which assesses the imbalance of two itemsets, A and B , in rule implications. It is defined as

$$IR(A, B) = \frac{|sup(A) - sup(B)|}{sup(A) + sup(B) - sup(A \cup B)}, \quad (6.13)$$

AR Metrics—LLR

其他：【卡方值】和lift特性非常相似的还有卡方值，其和lift指标二选一即可，此处不再介绍。【LLR】全称，log-likelihood ratio，是根据熵的变化计算相似度的一种方式。对于item项A和B来说，有：

	Event A	Everything but A
Event B	A and B together (k_{11})	B, but not A (k_{12})
Everything but B	A without B (k_{21})	Neither A nor B (k_{22})

LLR的计算公式为：

$$LLR = 2 \sum(k) (H(k) - H(\text{rowSums}(k)) - H(\text{colSums}(k)))$$

其中，熵的计算方法为：

$$H = \text{function}(k) \{N = \text{sum}(k); \text{return} (\text{sum}(k/N * \log(k/N + (k==0))))\}$$

由公式可知：

1. 满足交换律， $LLR(a, b) = LLR(b, a)$
2. 表达式恒大于0，0表示不相关，>0表示相关，可能正相关，也可能负相关。

Sim—Real Matrix

Sim Metrics—Cos

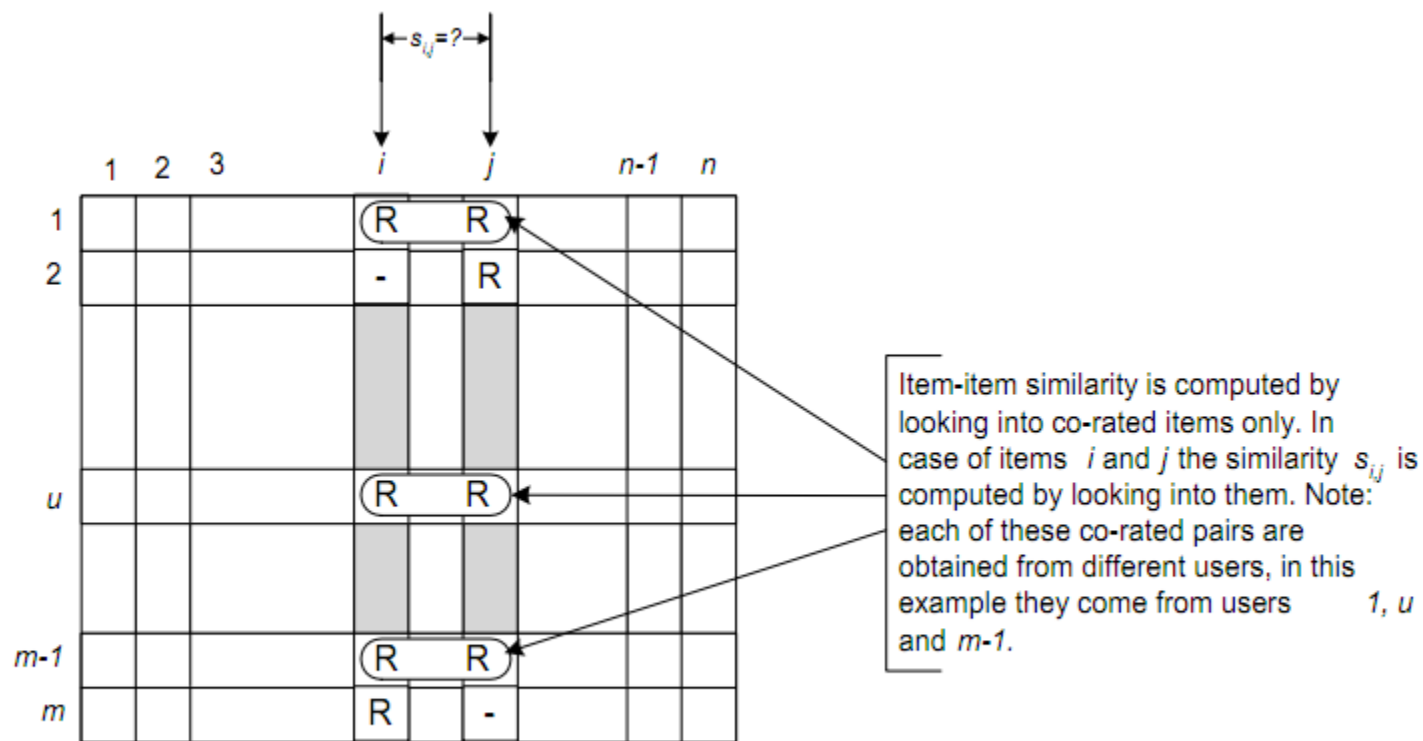


Figure 2: Isolation of the co-rated items and similarity computation

【cos】 上文已经提到了cos，此处的cos和上文的不同在于，user-item矩阵可以是实数型的，而关联矩阵必须是0-1矩阵。假设A和B是两个item，有 $\cos(A, B) = A * B / (|A| * |B|)$

Sim Metrics—Adjust Cos

【用adjust-cos去掉user-bias】 如果考虑到每user的打分标准都不同，user1喜欢打高分，user2喜欢打低分，应该把用户的打分标准减掉，有 $\text{adjust_cos}(A, B) =$

$$\frac{\sum_{u \in U} (R_{u,i} - \bar{R}_u)(R_{u,j} - \bar{R}_u)}{\sqrt{\sum_{u \in U} (R_{u,i} - \bar{R}_u)^2} \sqrt{\sum_{u \in U} (R_{u,j} - \bar{R}_u)^2}}.$$

Sim Metrics—Pearson

【用pearson去掉item-bias】 如果把A和B的向量归一化去掉bias，再做cos，即pearson的相似度：

$$\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Sim Metrics—Hybrid

1. 得到user-item matrix打分矩阵
2. 去掉matrix bias
3. 去掉user bias
4. 去掉item bias
5. 做cos

Similarity

Similarity Metrics--Bool or Real		
Bool Matrix Only	Both	Real Matrix Only
Support	Cos	Adjust cos
Jaccard		Pearson
Confidence		
Lift		
KULC		
LLR		
Conviction		

Similarity Metrics--Symmetry or Not	
Symmetry	Non Symmetry
Support	Confidence
Jaccard	
Lift	
KULC	
IR	
Cos	
Adjust Cos	
Pearson	

XX CFs

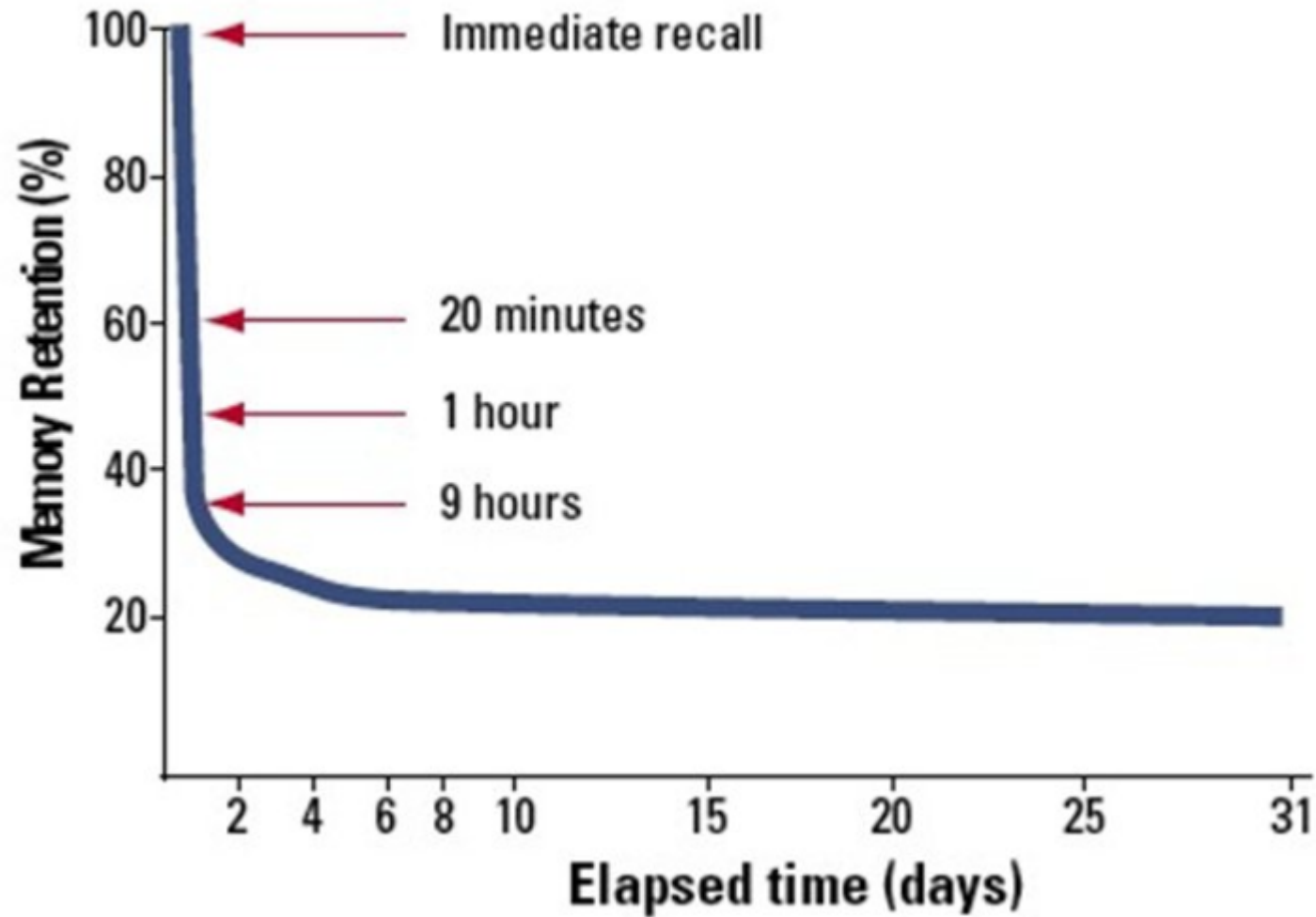
- Item2Vector CF
- Content CF
- LDA CF
- Simhash CF
- Node2Vector CF
- DNN Embedding CF

Refinement Prior

Action Rating

Action	Score
Play	1分
Favor	2分
Download	2分
Switch	-1'分
Trash	-2'分
Comment	3分
Search	2分
Share	2分

Time Decay



Session Split

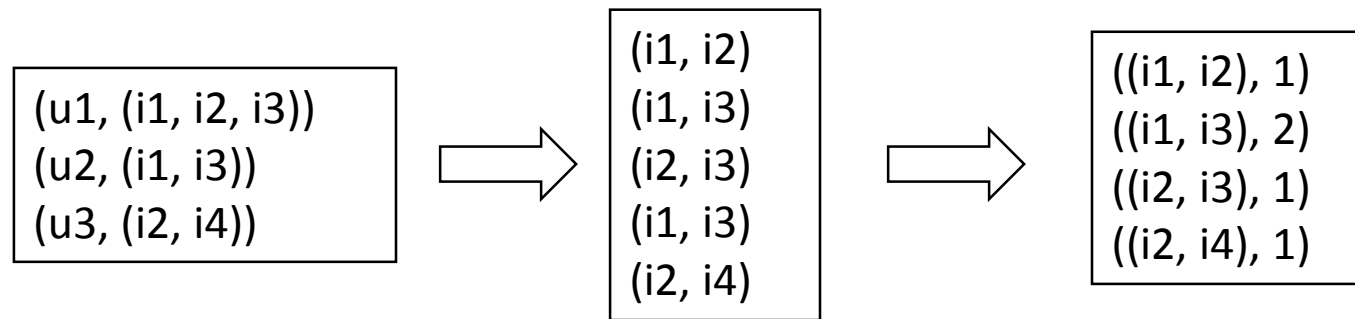
私人电台



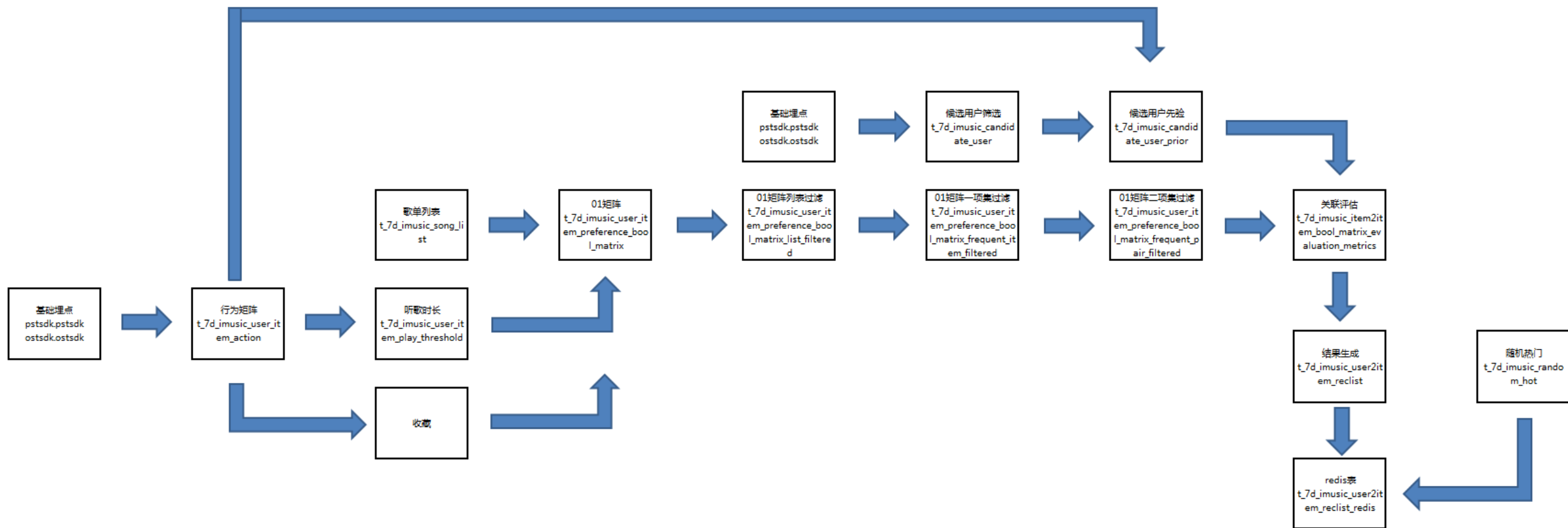
最近30分钟无行为

Best Practice

Data Sparsity

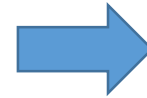
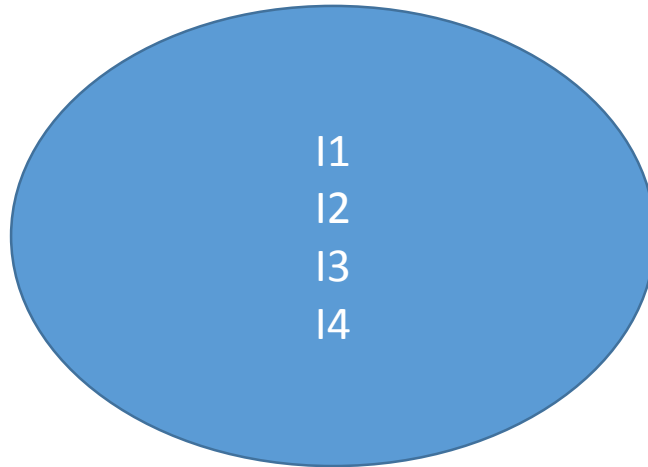


Data Flow

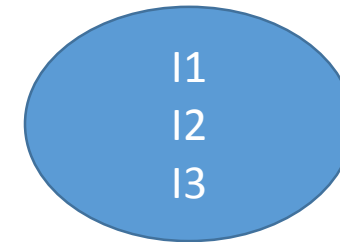


Prior Posterior Diff Set --1

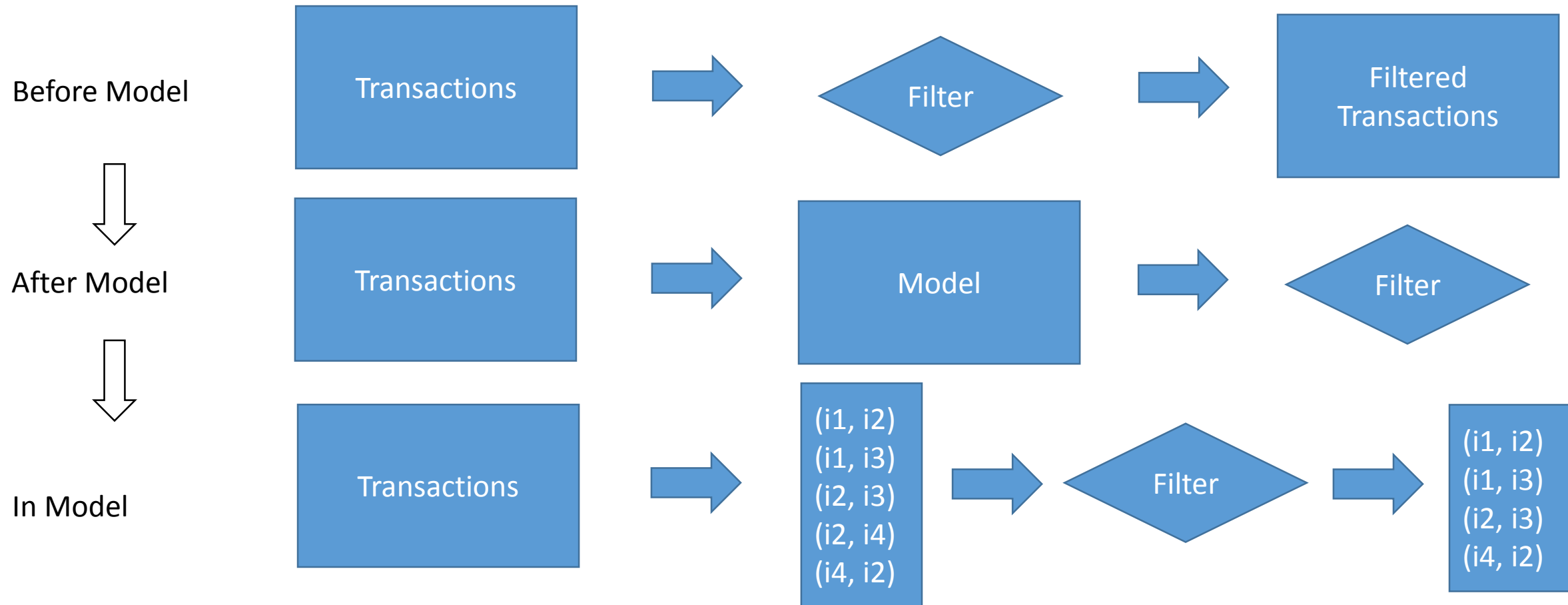
Prior Set



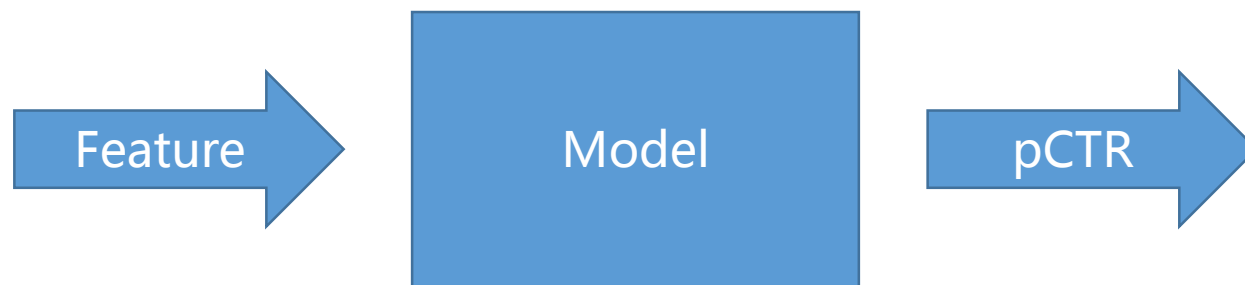
Posterior Set



Prior Posterior Diff Set --2



Online



pCTR 实时 = Feature 实时 or Model 实时

Q&A