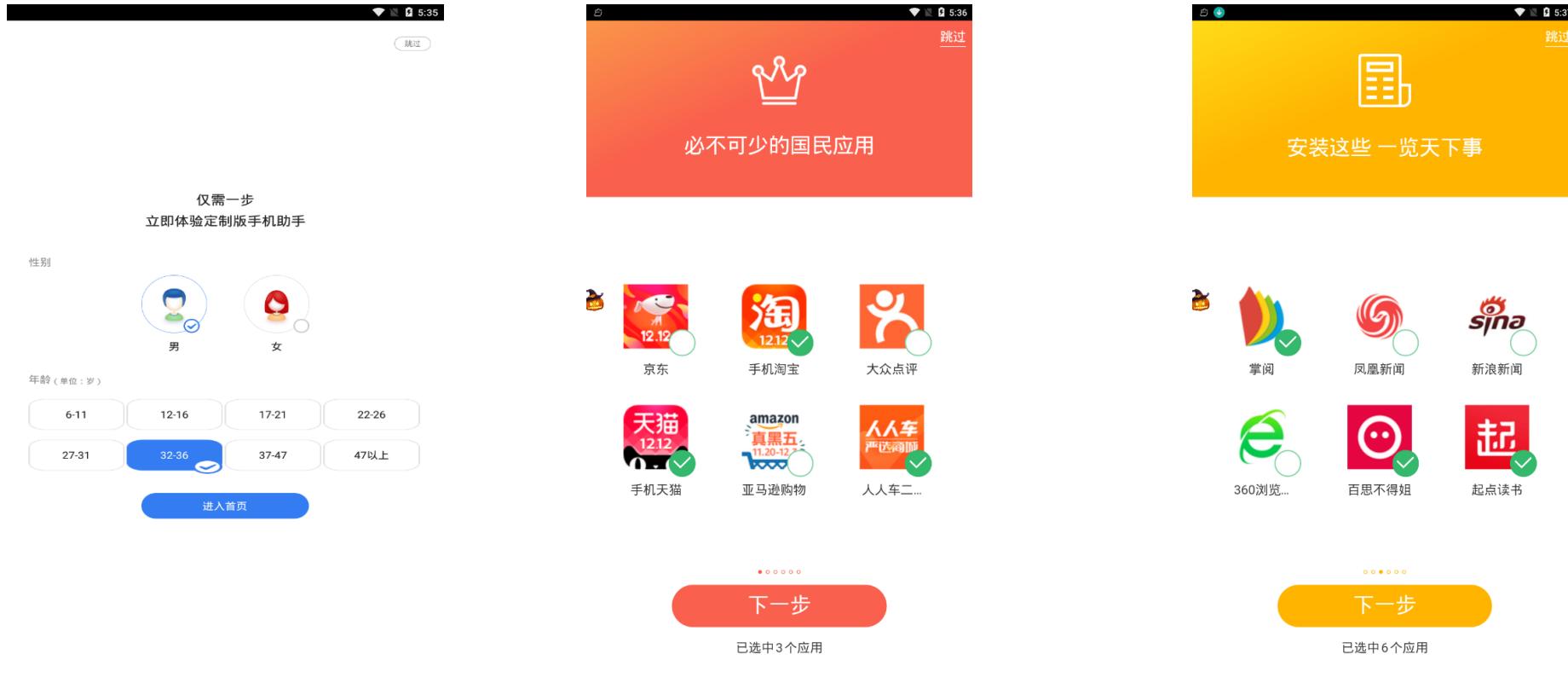


Recommendation System

陈浩

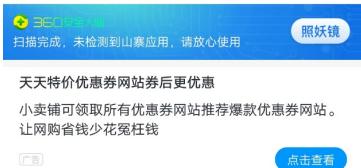
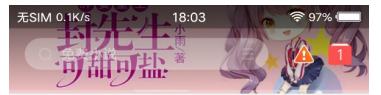
2020-12-11

Example



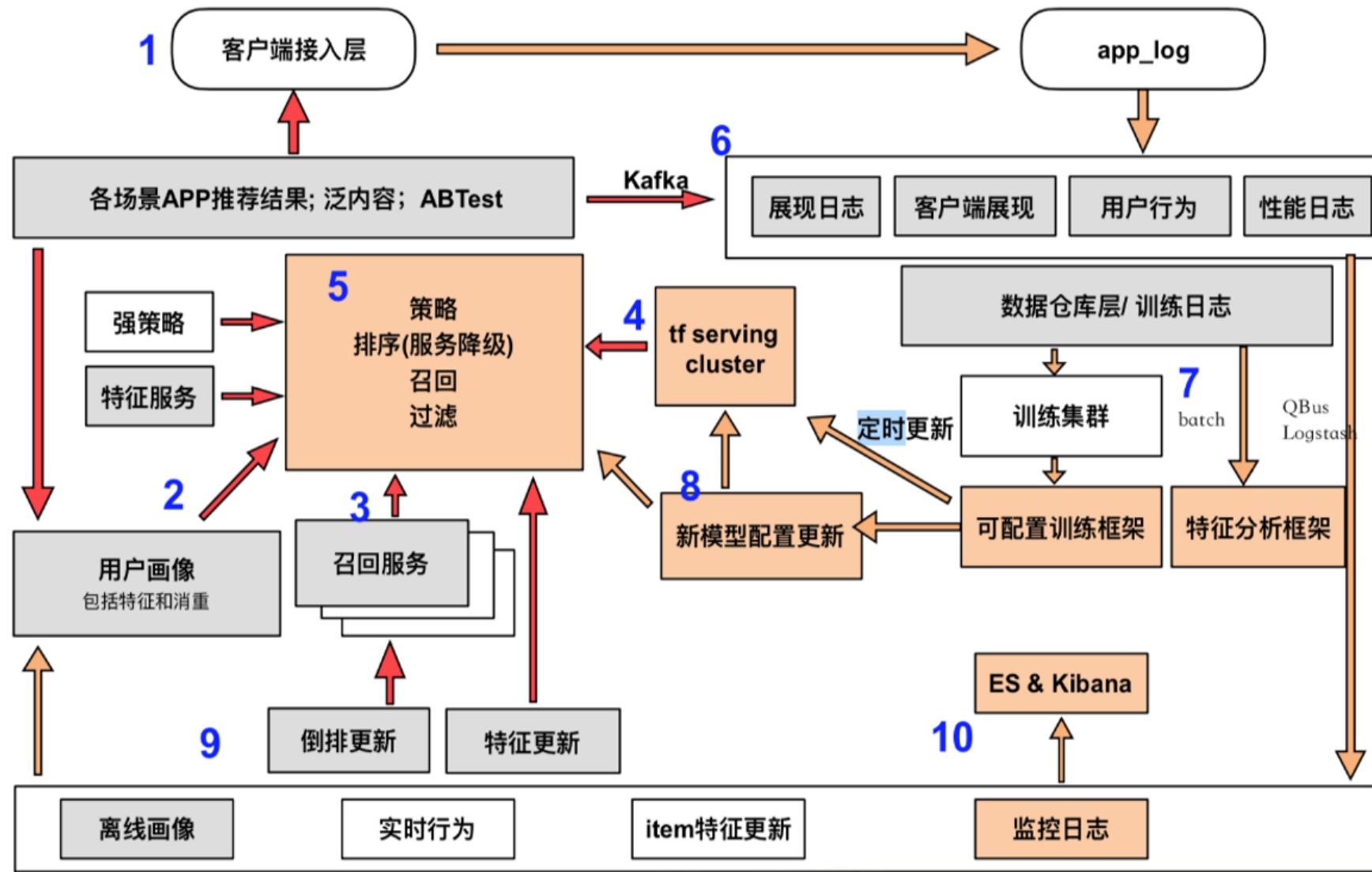
基于用户选择年龄、性别和地域信息推荐购物、旅游、新闻、工具、游戏四类二级类别应用
User 特征：用户的年龄、性别、地域、网络环境、渠道、手机型号品牌，
Item 特征：应用的名称、标签、简介、评论以及历史点击、下载、搜索、安装等信息

Example

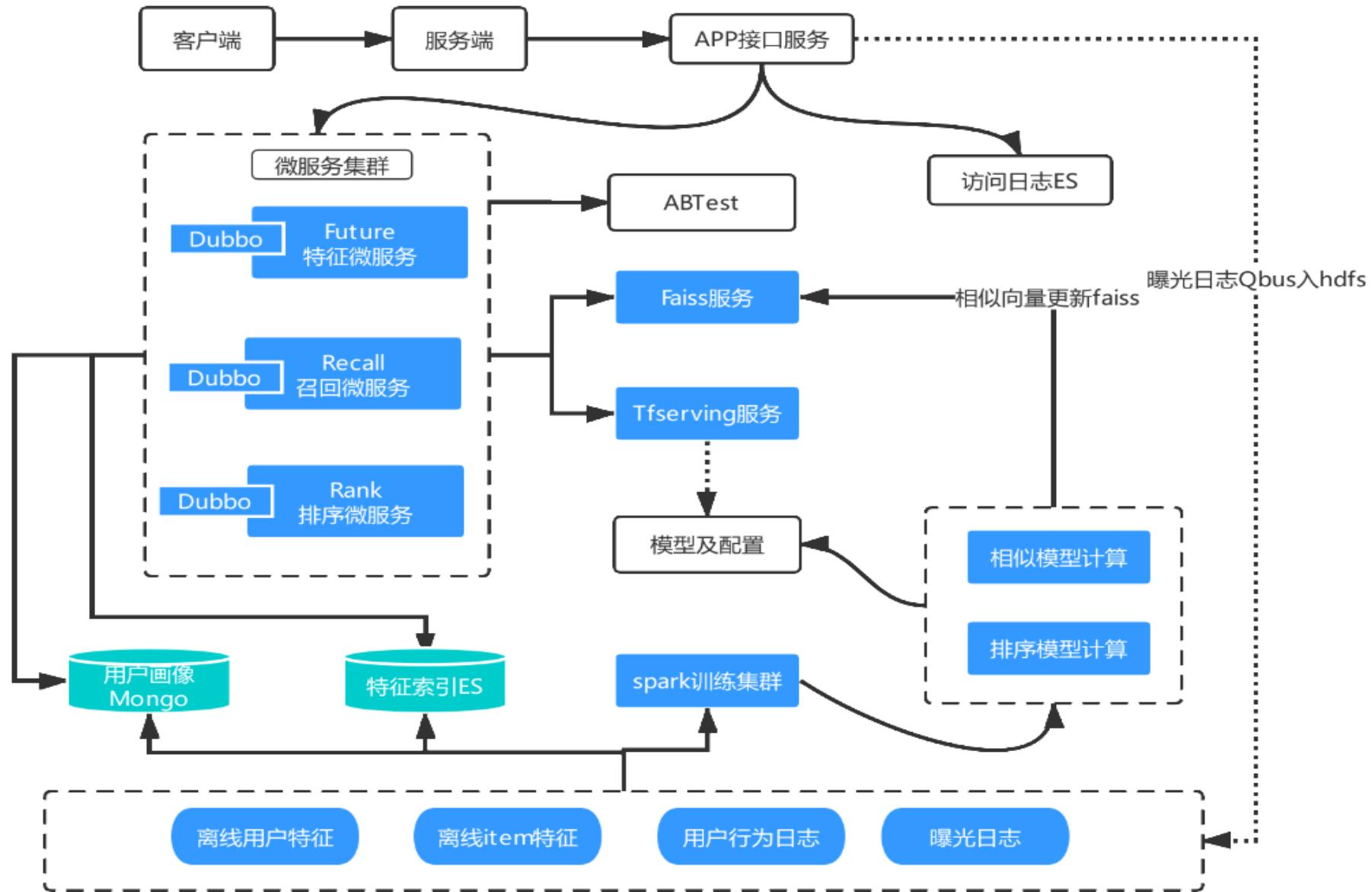


今日热点
弹窗
详情页
安装完成
搜索词
Topic List
开屏
搜索结果
猜你喜欢
下了还会下
.....

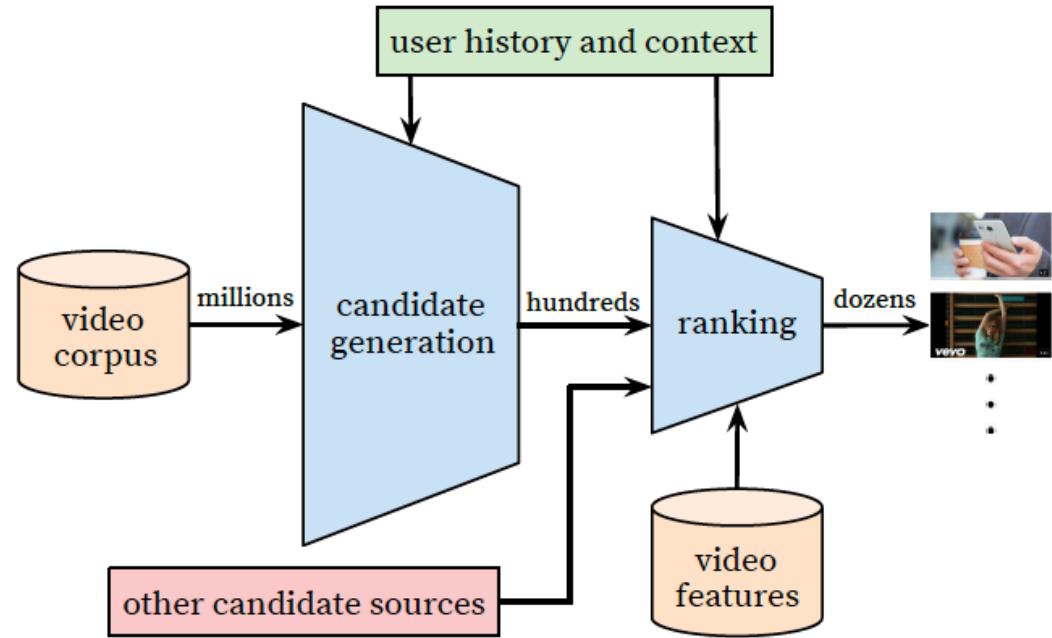
Engine Framework



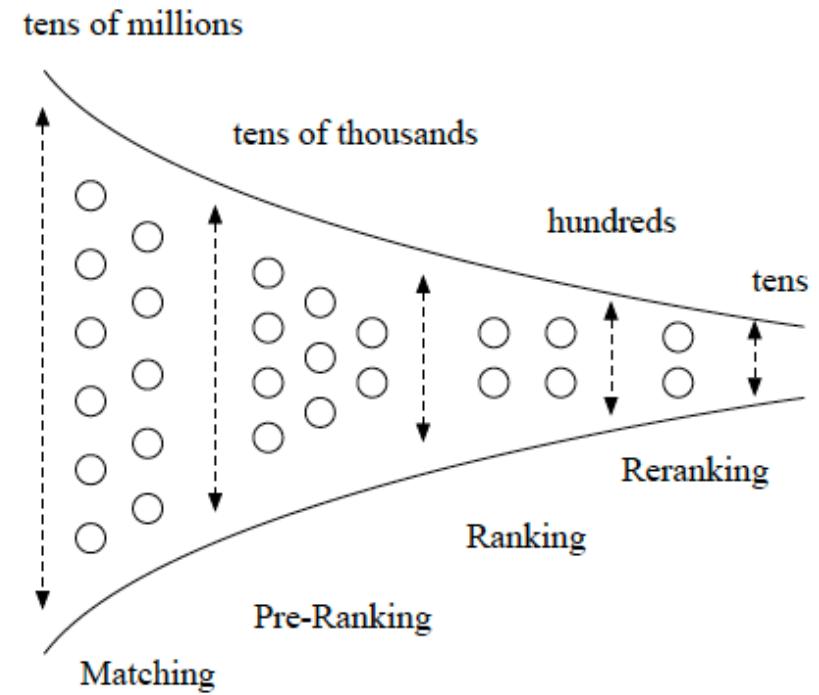
Algorithm Framework



Recommendation System Stage



Google 2016



Ali 2019



Match - TGI



Target Group Index (目标群体指数) (冷启动)

$$TGI = \frac{\text{目标群体中具有某一特征的群体所占比例}}{\text{总体中具有相同特征的群体所占比例}} * \text{标准数} 100$$

人数比例	男性	女性
总人数	60%	40%
王者荣耀	90%	10%
美图秀秀	20%	80%

仅需一步
立即体验定制版手机助手

性别

男

女

年龄 (单位:岁)

6-11 12-16 17-21 22-26

27-31 32-36 37-47 47以上

[进入首页](#)

TGI	男性	女性
王者荣耀	150	25
美图秀秀	33.33	200

$$[\text{王者荣耀}, \text{男性}] TGI = \frac{90\%}{60\%} * 100 = 150$$

Match - Collaborative Filtering

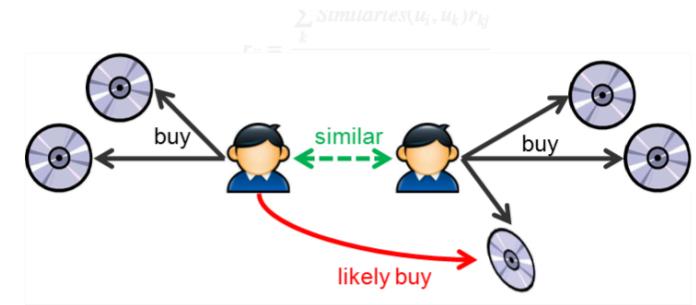
	风景	书	电影	游戏
用户1	喜欢	不喜欢	喜欢	喜欢
用户2		喜欢	不喜欢	不喜欢
用户3	喜欢	喜欢	不喜欢	
用户4	不喜欢		喜欢	
用户5	喜欢	喜欢	?	不喜欢

$$\text{Sim}(u_i, u_k) = \frac{\sum_j (r_{ij} - r_i)(r_{kj} - r_k)}{\sqrt{\sum_j (r_{ij} - r_i)^2} \sqrt{\sum_j (r_{kj} - r_k)^2}}$$

$$r_{ij} = \frac{\sum_k \text{Sim}(u_i, u_k)r_{kj}}{\text{number of ratings}}$$

优点：新颖性、相似推荐

缺点：冷启动、稀疏性、扩展性

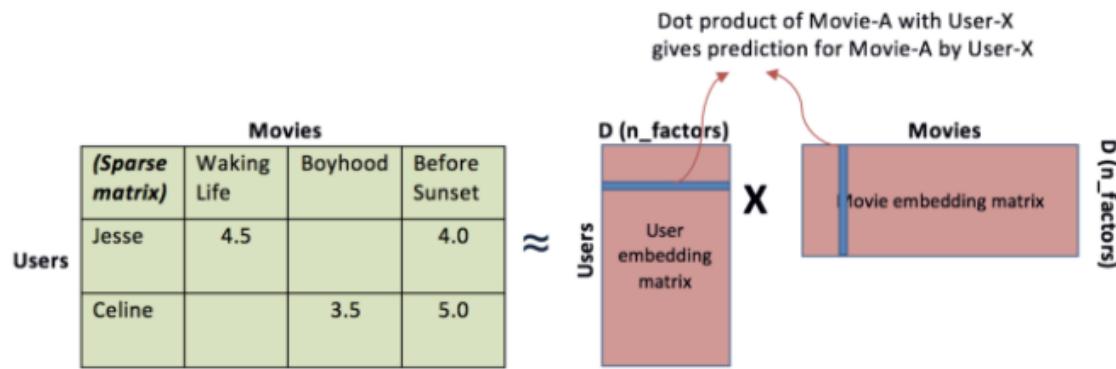


User-based CF

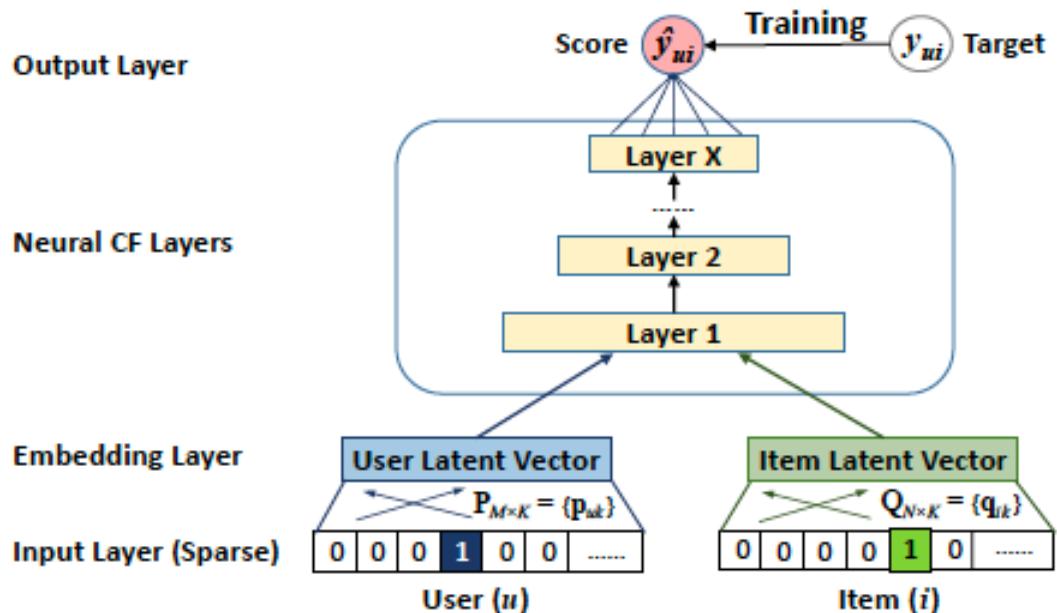


Item-based CF

Match - MF vs NCF

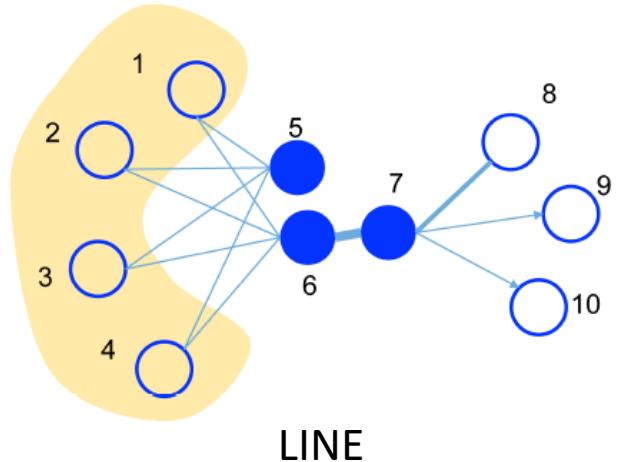


$$\min_{p,q,b_u,b_i} \sum (r_{ui} - (p_u^T q_i + \mu + b_u + b_i))^2$$

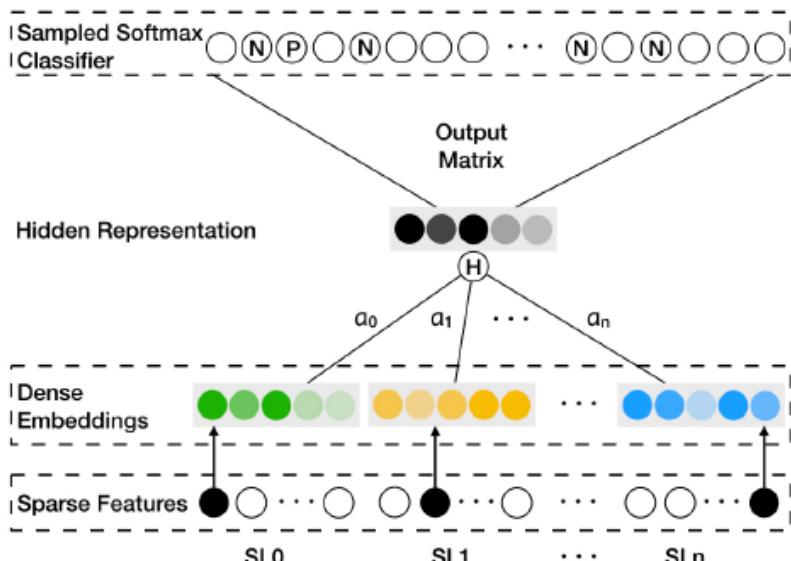


$$\begin{aligned}
 \mathbf{z}_1 &= \phi_1(\mathbf{p}_u, \mathbf{q}_i) = \begin{bmatrix} \mathbf{p}_u \\ \mathbf{q}_i \end{bmatrix} \\
 \phi_2(\mathbf{z}_1) &= a_2 (\mathbf{W}_2^T \mathbf{z}_1 + \mathbf{b}_2) \\
 &\dots \\
 \phi_L(\mathbf{z}_{L-1}) &= a_L (\mathbf{W}_L^T \mathbf{z}_{L-1} + \mathbf{b}_L) \\
 \hat{y}_{ui} &= \sigma(\mathbf{h}^T \phi_L(\mathbf{z}_{L-1}))
 \end{aligned}$$

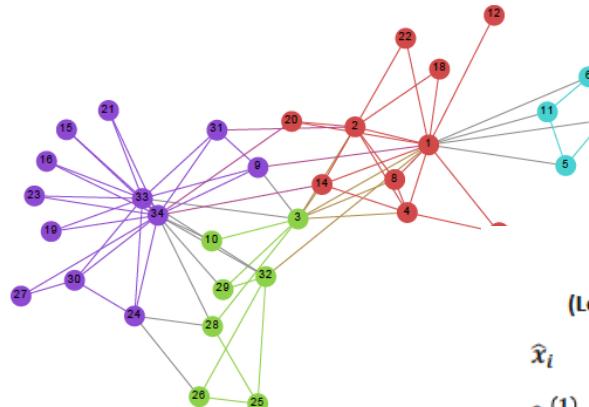
Match - Graph Embedding



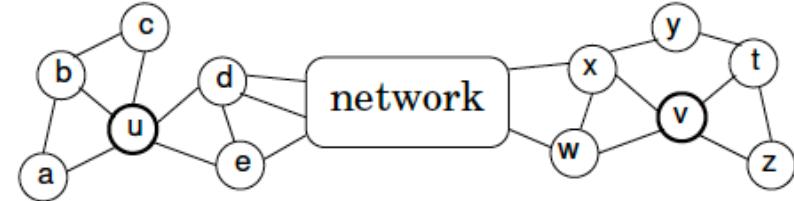
LINE



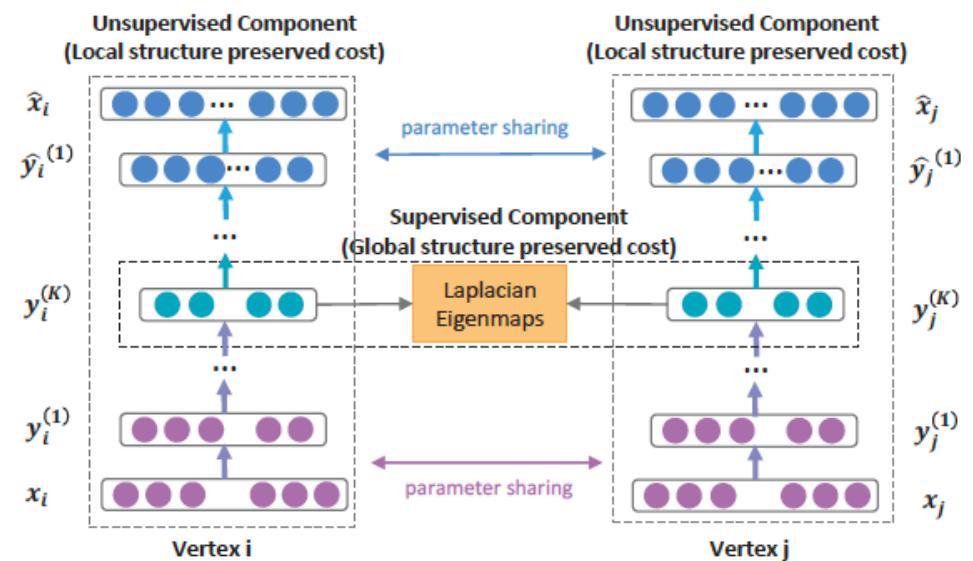
EGES



Deepwalk



Struct2vec



SDNE

Match - Node2vec

兼顾了深度优先和广度优先的策略

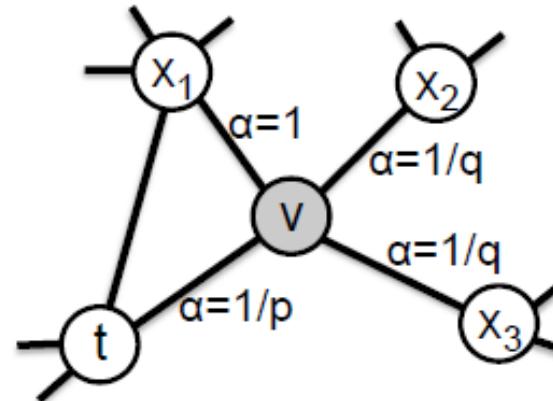
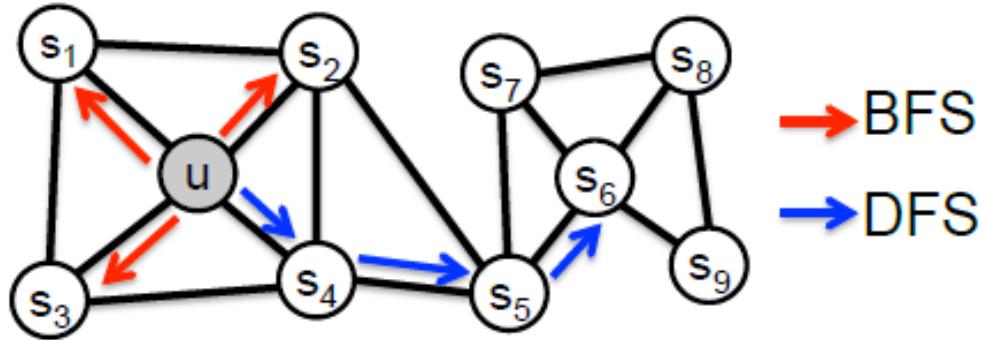
Data : 有偏的随机游走策略，将空间数据转化为序列数据

给定节点 v ，访问下一个节点 x 的概率为：

$$p(c_t = x | c_{t-1} = v) = \begin{cases} \frac{\pi_{vx}}{Z} & (v, x) \in E \\ 0 & otherwise \end{cases} \quad Z \text{ 为归一化常数}$$

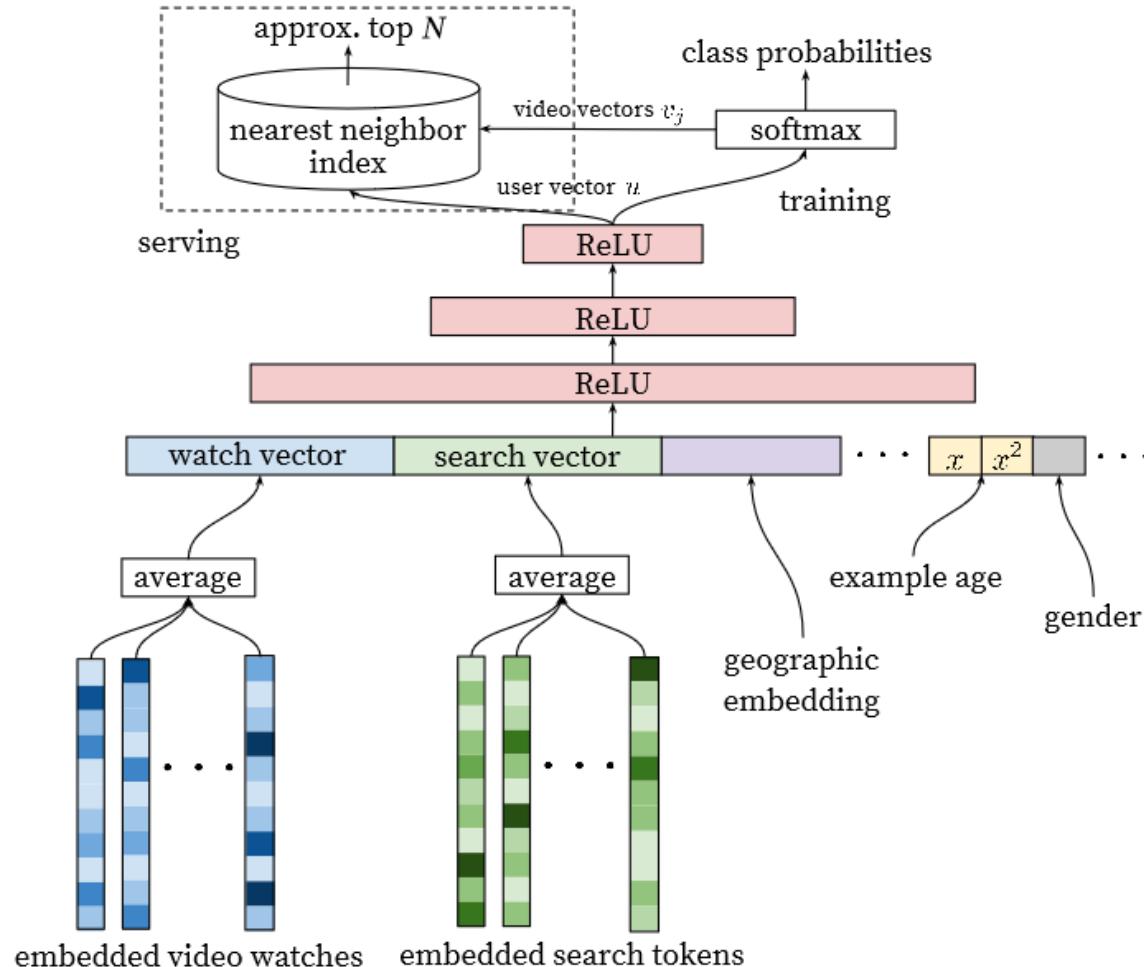
node2vec利用两个超参数 p 和 q 来控制随机游走的策略，假设随机游走经过边 (t, v) 到达节点 x 则 $\pi_{vx} = \alpha_{pq}(t, x) \cdot w_{vx}$

$$\alpha_{pq} = \begin{cases} \frac{1}{p} & if \quad d_{tx} = 0 \\ 1 & if \quad d_{tx} = 1 \quad d_{tx} \text{ 为节点 } t \text{ 到节点 } x \text{ 之间的最短路径距离} \\ \frac{1}{q} & if \quad d_{tx} = 2 \end{cases}$$

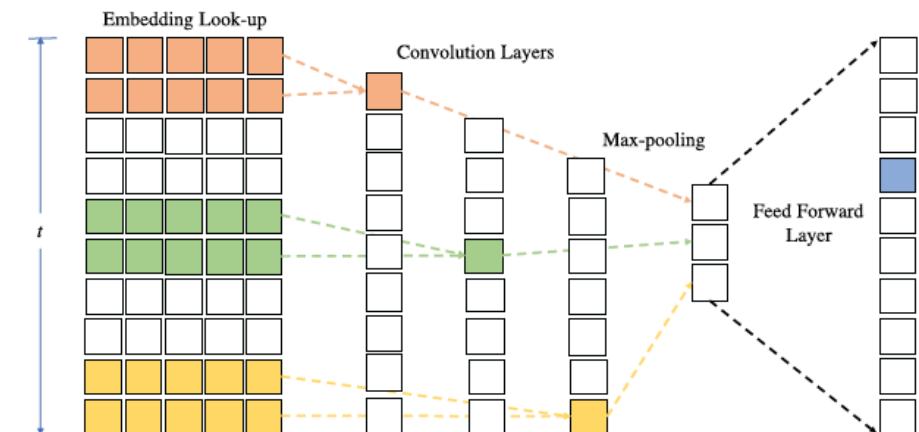


p 控制重复访问刚刚访问过的节点的概率
 q 控制偏向BFS或者偏向DFS的概率

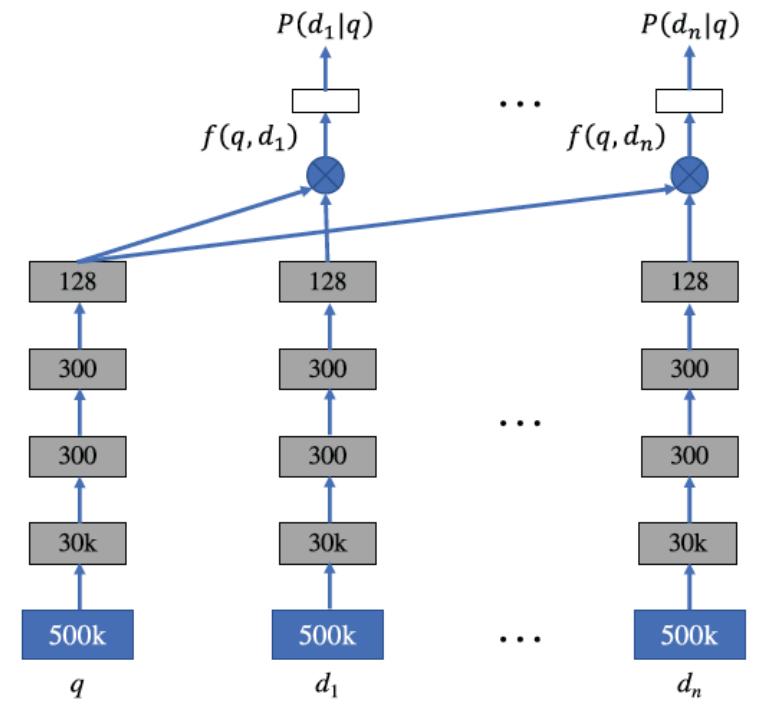
Match - DNN



Google 2016



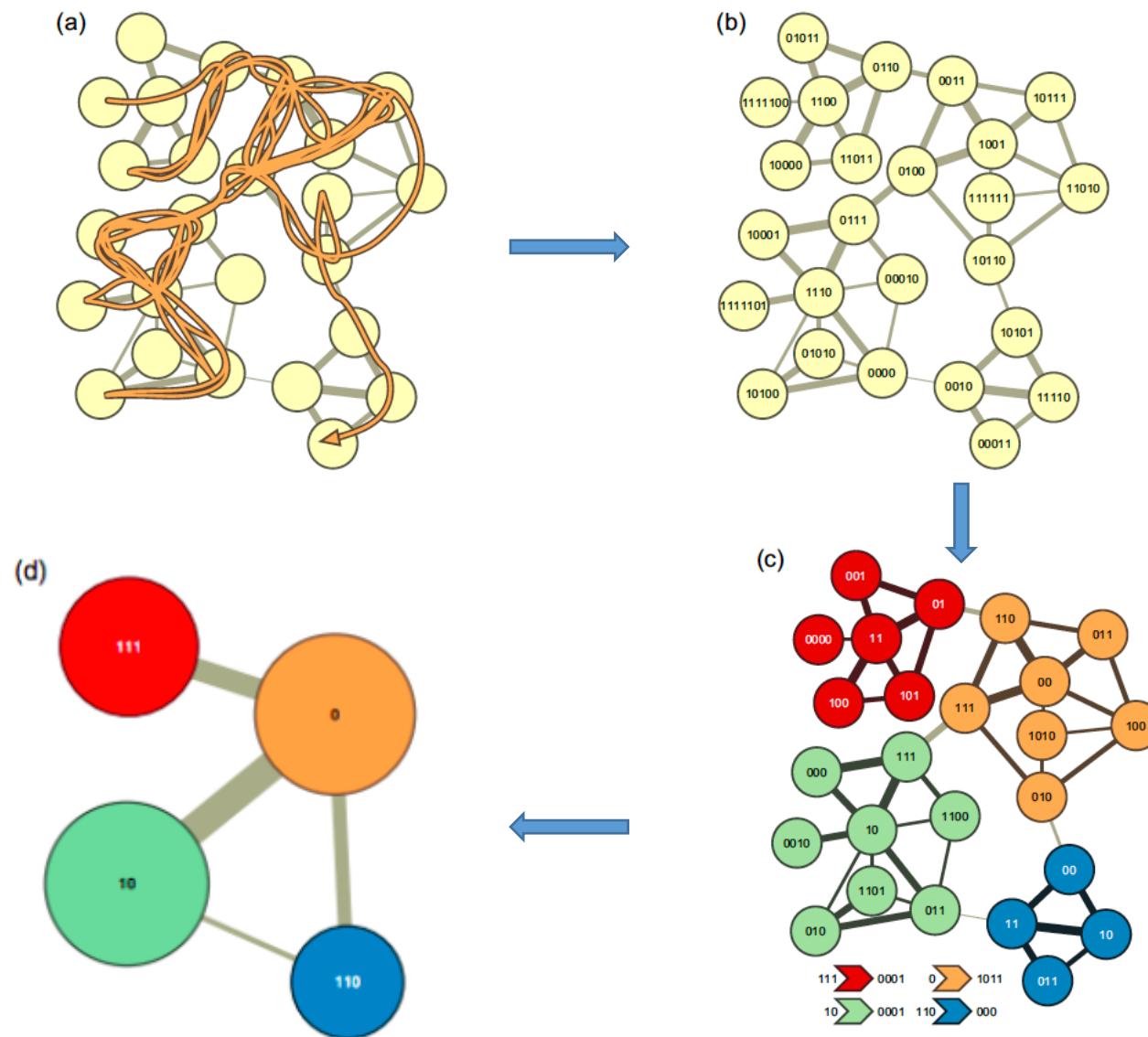
Caser



DSSM

Match - Others

- 热门
 - 广告
 - 聚类 (InfoMap)
 - 关联规则挖掘
 - 序列召回
 - 最近邻
 - 树模型



InfoMap

Rank - FM

$$x = [x_{user}; x_{item}] = [x_1, x_2, \dots, x_n]$$

$$y \in [0, 1]$$

$$Linear: \hat{y}(x) = w_0 + \sum_{i=1}^n w_i x_i$$

$$LR: \hat{y}(x) = \frac{1}{1 + w_0 \exp(-w^T x)}$$

$$FM: \hat{y}(x) = w_0 + \sum_{i=1}^n w_i x_i + \sum_{i=1}^n \sum_{j=i+1}^n \langle v_i, v_j \rangle x_i x_j$$

FM模型

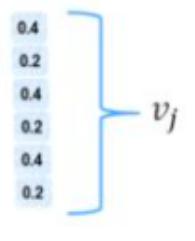
$$\langle v_i, v_j \rangle := \sum_{f=1}^k v_{i,f} \cdot v_{j,f}$$

v_1	0.3	0.2	0.6	0.8	0.1	0.2
v_2	0.1	0.8	0.6	0.8	0.4	0.6
v_3	0.4	0.2	0.7	0.2	0.1	0.2
v_4	0.1	0.2	0.6	0.8	0.5	0.2
v_{n-1}	0.3	0.2	0.6	0.8	0.1	0.2
v_n	0.6	0.8	0.9	0.8	0.4	0.6



$w_{i,j} = \langle v_i, v_j \rangle \neq 0$
even if 训练数据中 $x_i x_j = 0$
only if 在训练数据中存在 k 使得 $x_i x_k \neq 0$

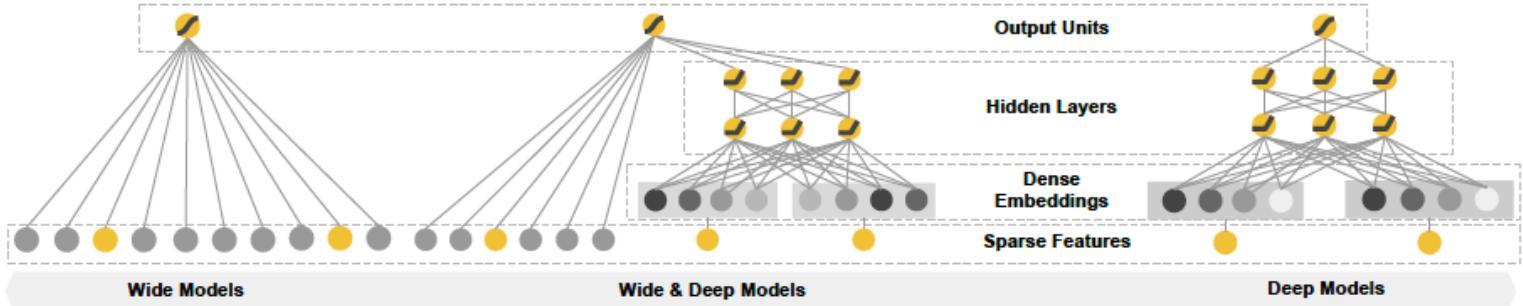
FM模型泛化能力强



Rank - Wide&Deep

$$w_{wide} = \frac{1}{1 + w_0 \exp(-w^T x)}$$

$$\hat{y}(x) = \sigma \left(w_{wide}^T [x, \phi(x)] + w_{deep}^T a^{(l_f)} + b \right)$$



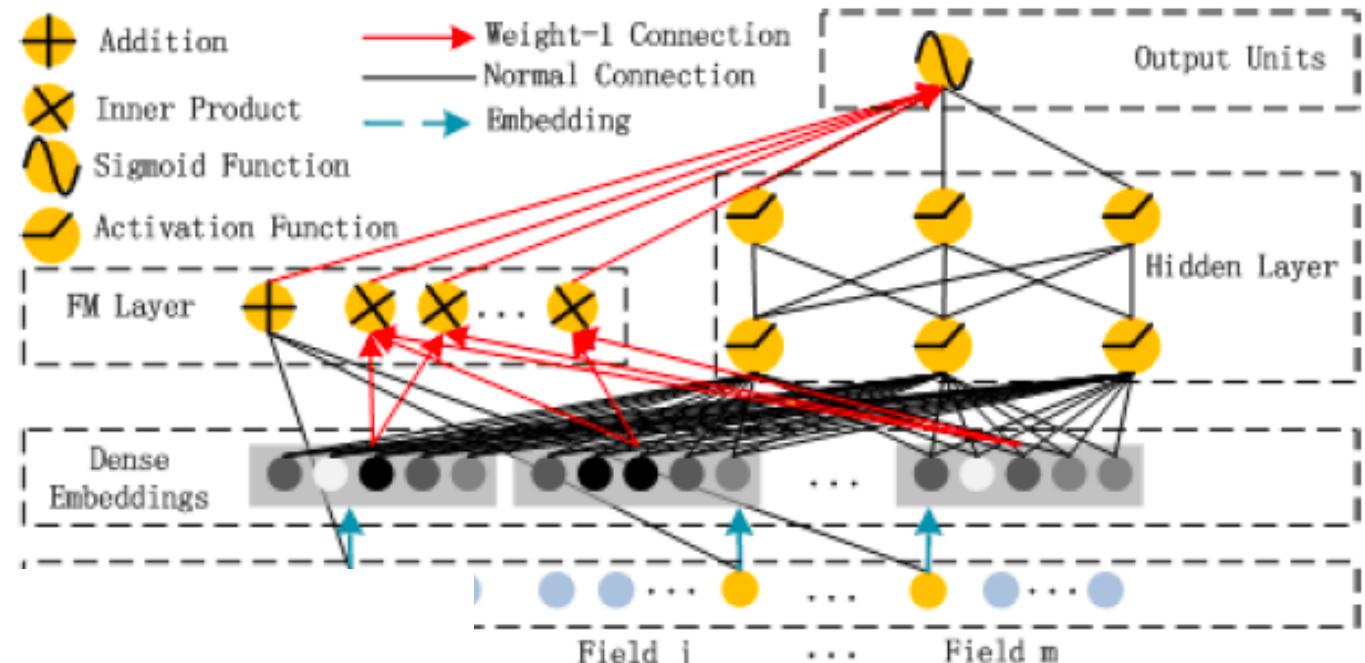
- Wide: Memorization, 基于人工先验知识, 特征交叉, 喂入Wide侧, 让Wide侧能够记住这些规则
- Deep: Generalization, 普通DNN, embedding categorical feature 让DNN学习到这些特征之间的深层交叉, 以增强扩展能力
- 人工特征交叉
样本 x 中的特征 $gender=female$ 和特征 $language=en$ 同时为1
新的组合($gender=female, language=en$)才为1
- Optimizer: Follow- the-regularized-leader (FTRL) Algorithm

Rank - DeepFM

- 自动的二阶特征交叉
- FM和Deep部分共享embedding
- 参数量虽多，效果优于wide&deep

$$\hat{y}(x) = \sigma(y_{FM} + y_{DNN})$$

$$y_{FM} = \langle w, x \rangle + \sum_{j_1=1}^d \sum_{j_2=j_1+1}^d \langle V_i, V_j \rangle x_{j_1} \cdot x_{j_2}$$

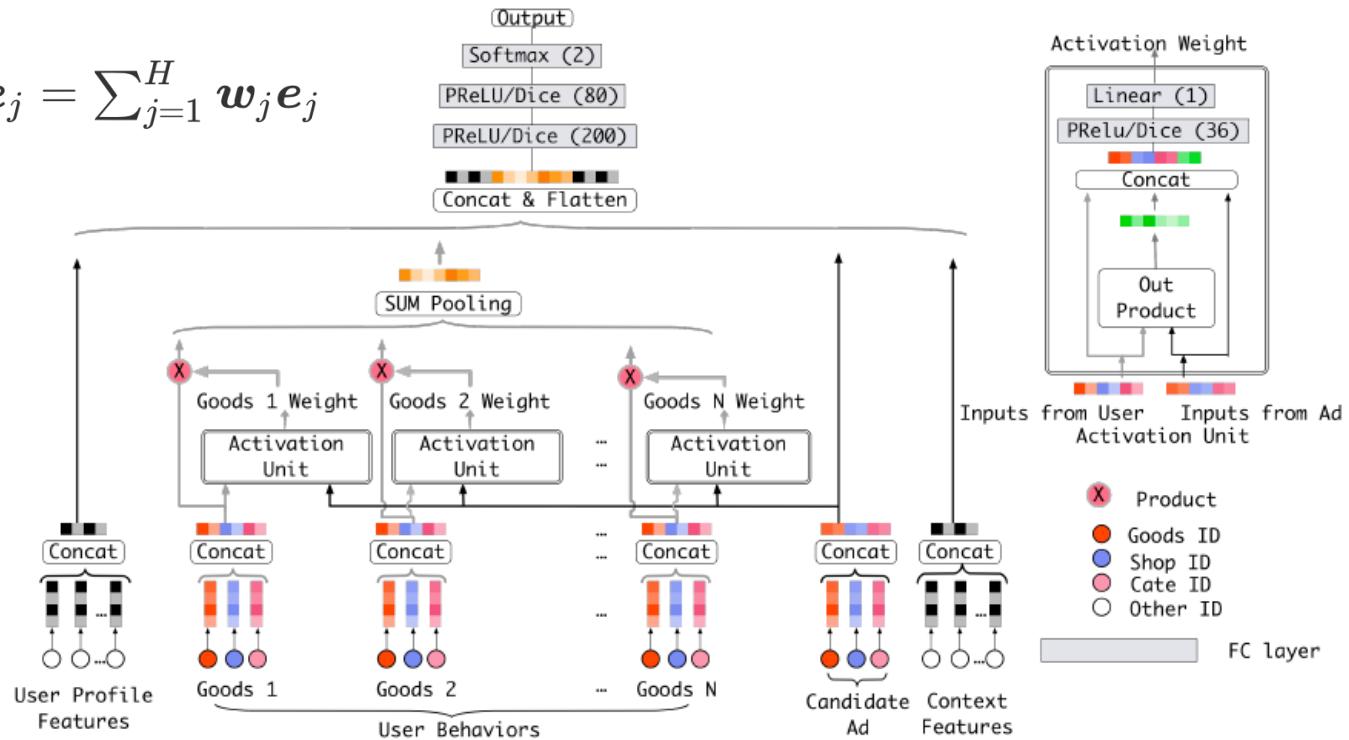


Rank - DIN

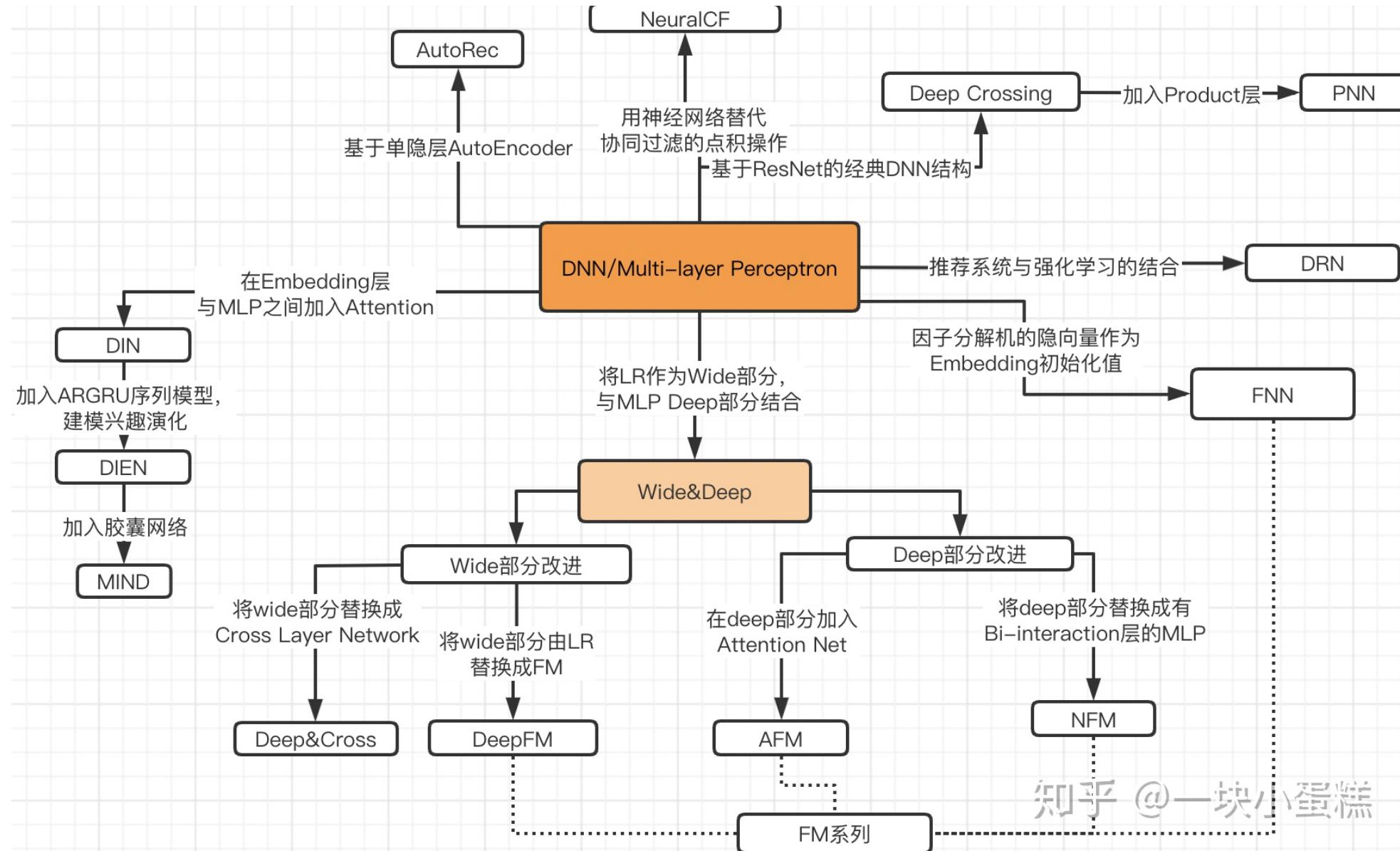
The use of fixed-length vector will be a bottleneck, which brings difficulty for Embedding&MLP methods to capture users diverse interests effectively from rich historical behaviors.

$$\mathbf{v}_U(A) = f(\mathbf{v}_A, \mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_H) = \sum_{j=1}^H a(\mathbf{e}_j, \mathbf{v}_A) \mathbf{e}_j = \sum_{j=1}^H \mathbf{w}_j \mathbf{e}_j$$

- 适合电商场景
- 考虑到用户的历史行为特征
- 考虑到用户的兴趣的变化



Rank - Others



Multi Task Learning- ESMM

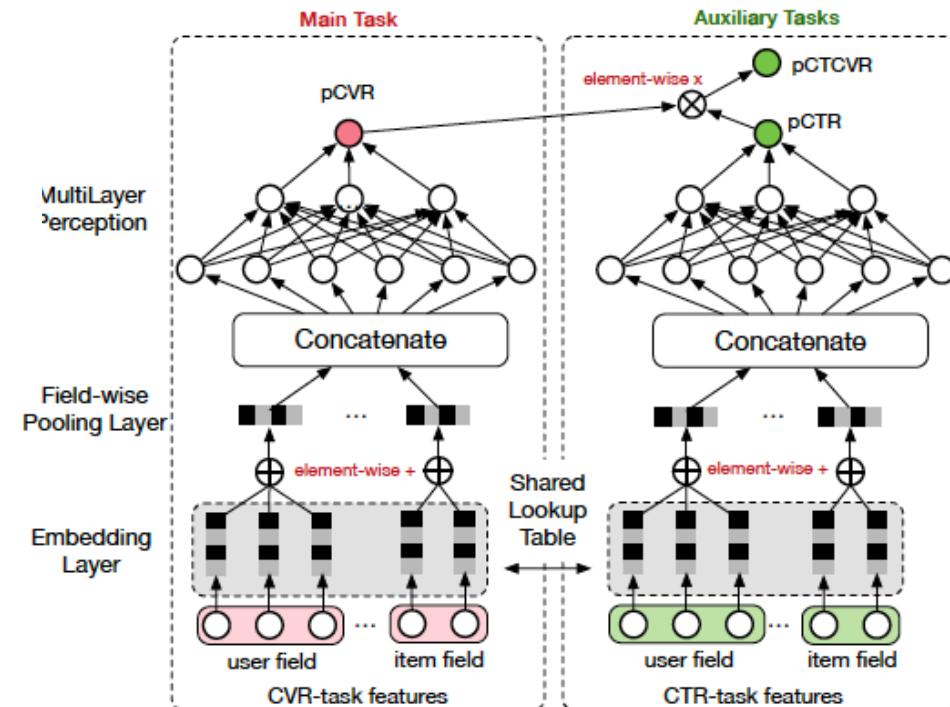
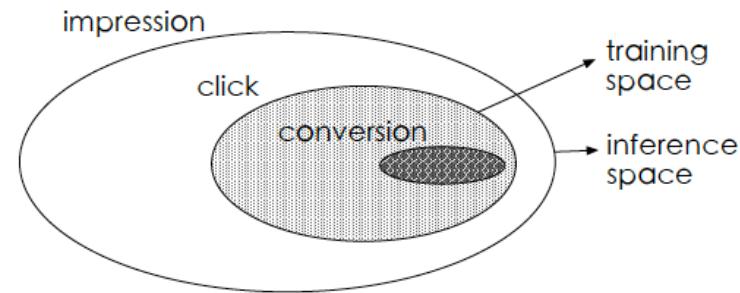
CTR : 从曝光到点击

CVR : 从点击到购买的转化

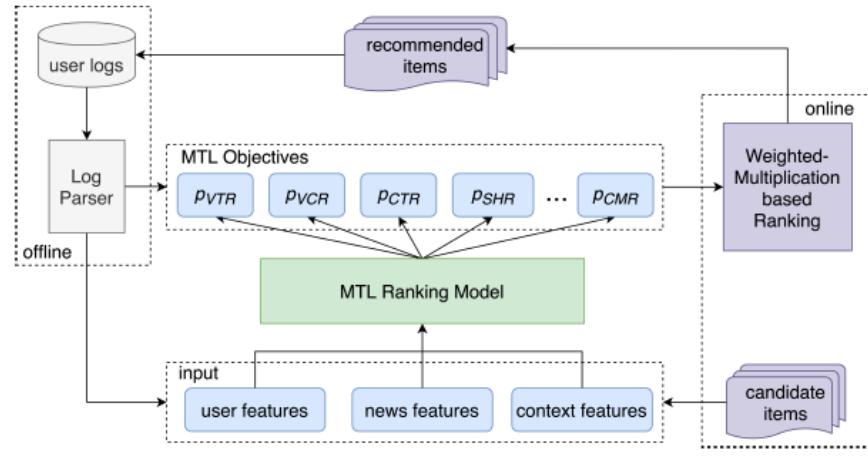
$$\underbrace{p(y = 1, z = 1 \mid x)}_{pCTCVR} = \underbrace{p(y = 1 \mid x)}_{pCTR} \times \underbrace{p(z = 1 \mid y = 1, x)}_{pCVR}$$

$$p(z = 1 \mid y = 1, x) = \frac{p(y=1,z=1|x)}{p(y=1|x)}$$

$$\begin{aligned} L(\theta_{cvr}, \theta_{ctr}) &= \sum_{i=1}^N l(y_i, f(x_i; \theta_{ctr})) \\ &+ \sum_{i=1}^N l(y_i \& z_i, f(x_i; \theta_{ctr}) \times f(x_i; \theta_{cvr})) \end{aligned}$$



Multi Task Learning- PLE

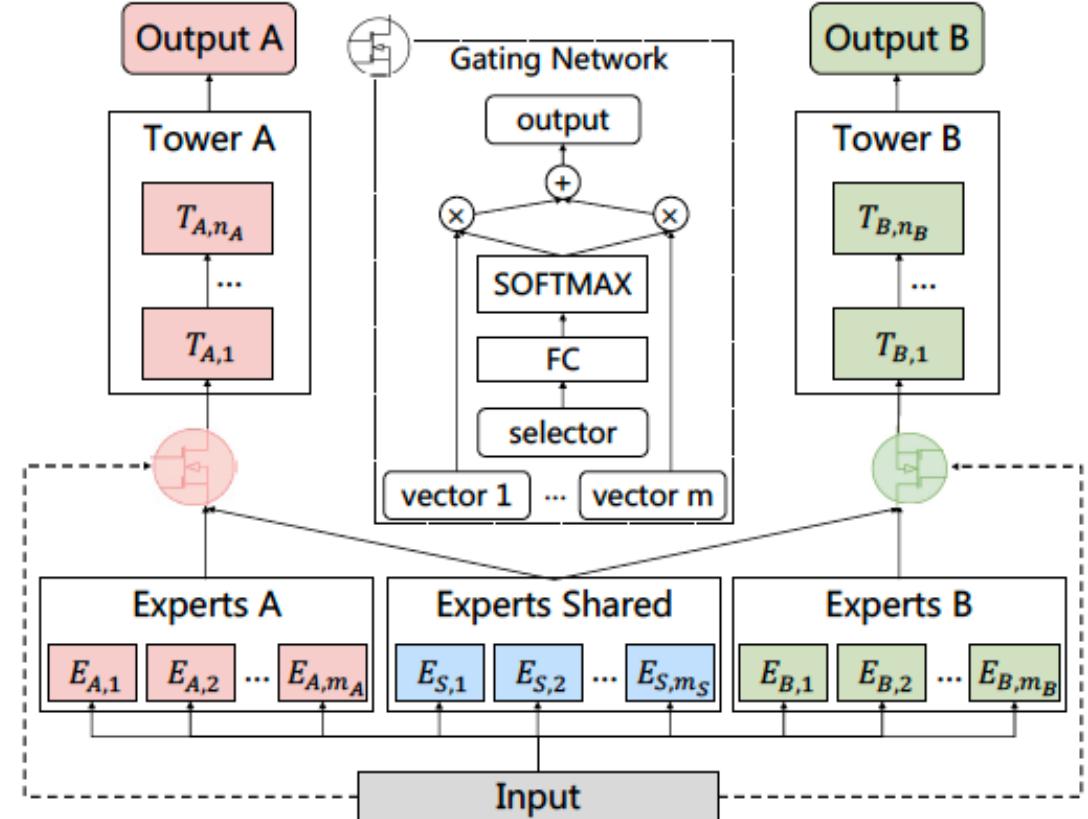


we have observed an interesting seesaw phenomenon that performance of one task is often improved by hurting the performance of some other tasks.

$$S^k(x) = \left[E_{(k,1)}^T, E_{(k,2)}^T, \dots, E_{(k,m_k)}^T, E_{(s,1)}^T, E_{(s,2)}^T, \dots, E_{(s,m_s)}^T \right]^T$$

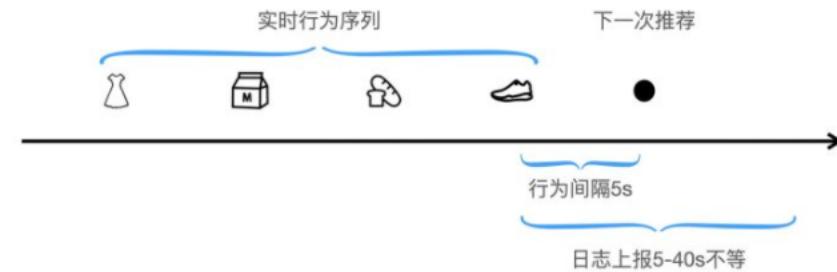
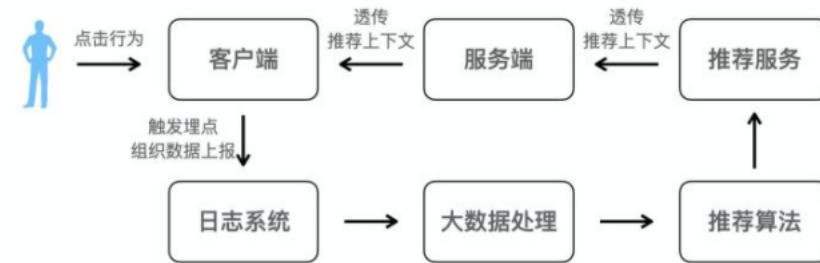
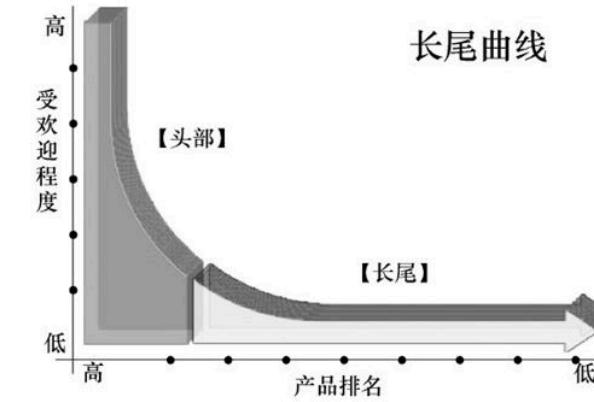
$$w^k(x) = \text{Softmax}(W_g^k x)$$

$$g^k(x) = w^k(x)S^k(x)$$



Problem

- 冷启动问题
- 评价指标不一
- Item 长尾问题
- E&E, exploration & exploitation 兴趣探索
- Data：
 - 数据量不足和数据稀疏
 - 数据清洗不干净或者过于干净
 - 数据分布的不一致
- Model：
 - 线下AUC涨，线上CTR跌
 - 特征/数据出现穿越
 - 实时性问题
- 短期被人高估，长期被人低估



谢 谢 !