

When Recommender Systems Meet Network Representation Learning

王鸿伟

June 3, 2018



Outline

- Recommender systems (RS)
- Network representation learning (NRL)
 - GraphGAN [AAAI 18]
 - Knowledge graph embedding
- RS + NRL
 - One-by-one learning [WWW 18]
 - Alternate learning [NIPS 18 in submission]
 - Joint learning [CIKM 17][WSDM 18][CIKM 18 in submission]
 - Pros and cons

Recommender Systems

寻梦环游记 Coco (2017)

喜欢这部电影的人也喜欢



头脑特工队



帕丁顿熊2



你的名字。



龙猫



海边的曼彻斯特



摔跤吧！爸爸

Movie



机器学习 平装 - 2016年2月1日

周志华 (作者)

★★★★☆ 142 条商品评论 | 分享

显示所有 格式和版本

平装
¥73.30 ✓prime

促销信息: 满减 中文图书全场满99元赠书 共1个促销

配送至: 北京朝阳区: 现在有货

送达日期: 明天(12月24日), 请在6小时54分钟内下单并选择“快递送货上门”。
(精确送达时间请至结账页面查询)

销售配送: 由亚马逊直接销售和发货。

购买此商品的顾客也同时购买



统计学习方法
李航
★★★★☆ 158
平装
¥32.80 ✓prime



机器学习实战
哈林顿 (Pete ...
★★★★☆ 154
平装
¥47.30



深度学习(deep learning)
伊恩·古德费洛 (...
★★★★☆ 63
平装
¥113.50

Book

Recommender Systems

● **热门推荐** 华语 | 流行 | 摇滚 | 民谣



吃的饱饱

429万

胖友们，又到了屯脂肪的季节啦！



锦鲤歌单|考研静心|美梦成真

52万



10万

没机会现场看爱豆的演唱会？一张歌单过足瘾



新世相

11741

电台节目 新世相 | 我，斗战胜佛

Music

美食 **火锅** 西餐 日本菜 霸王餐 全部 >



蟹的网田屋 **订** **团**

722条点评 日本料理 百联世纪/96广场 ¥574/人



弦月窗.LUNET... **订** **团**

1440条点评 法国菜 淮海路 ¥490/人



黑门和牛 **团**

1701条点评 日式烧维/铁板烧 新天地 ¥464/人



尚9-滴水 **订** **团**

4354条点评 西餐 北外滩 ¥534/人



1886汽车主题德... **团**

9979条点评 西餐 外滩 ¥170/人



度小月 **团**

2898条点评 台湾菜 月星环球港 ¥79/人

Restaurant

Recommender Systems



QA



Video




News

Recommender Systems


Recommender systems (RS) intend to address the information explosion by finding a small set of items for users to meet their personalized interests and demands

RS Task


- Recommendation Task 1: Rating / Click-through rate (CTR) prediction




2	?	3	?
?	?	4	?
?	5	?	2
3	1	4	?



explicit feedback




1	?	0	?
?	?	1	?
?	0	?	1
0	1	0	?




implicit feedback

RS Task

- Recommendation Task 2: Top-K recommendation



2		3	
		4	
	5		2
3	1	4	



You may also like...

	Wolf Warriors
	Saving Private Ryan
	Hacksaw Ridge
	Forrest Gump

Collaborative Filtering

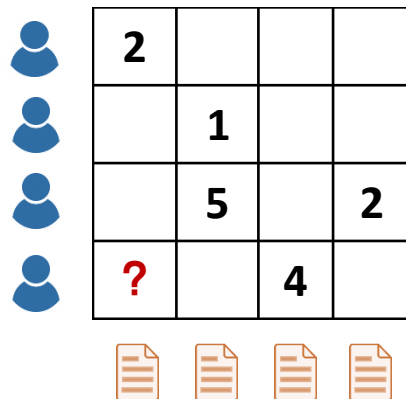
- **Collaborative filtering (CF)** considers users' historical interactions and makes recommendations based on their potential common preferences
 - Matrix factorization (MF)





$$R_{pq} = \mathbf{p}^\top \mathbf{q}$$


$$L = \| \mathbf{R} - \mathbf{P}^\top \mathbf{Q} \|_2^2 + \| \mathbf{P} \|_2^2 + \| \mathbf{Q} \|_2^2$$

CF fails to address...

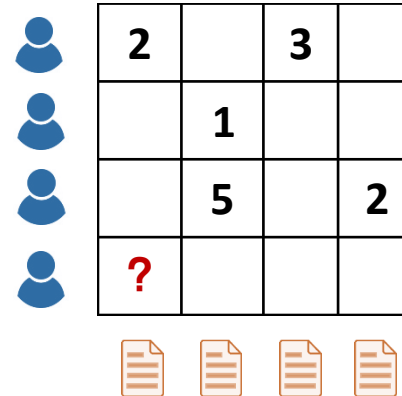
- Sparsity of user-item interactions
- Cold start problem








	2			
		1		
		5		2
	?		4	



sparsity

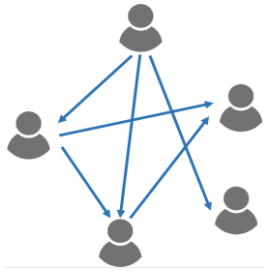


	2		3	
		1		
		5		2
	?			



cold start

CF + Side Information



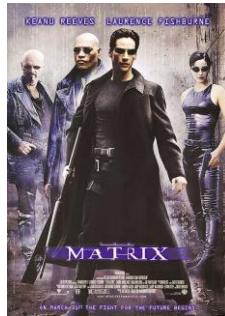
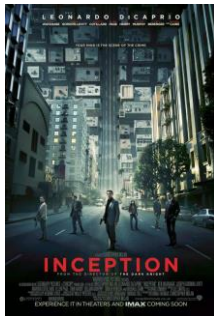
Social network

Alice
Female
California
...



iPhone X
2017
5.8 inch
\$999
...

User/item attributes



Multimedia



purchase



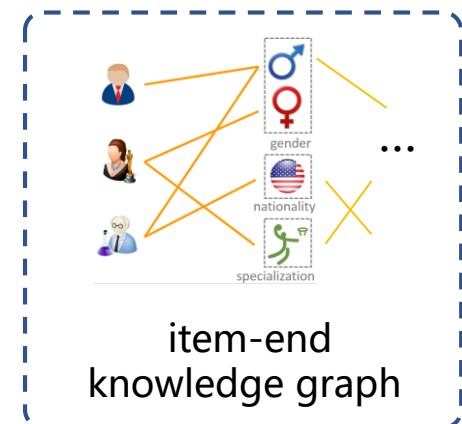
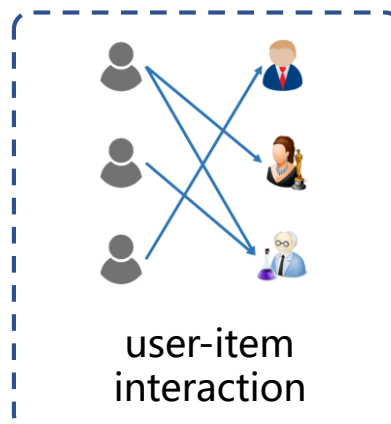
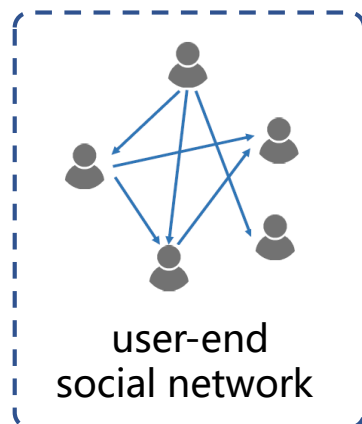
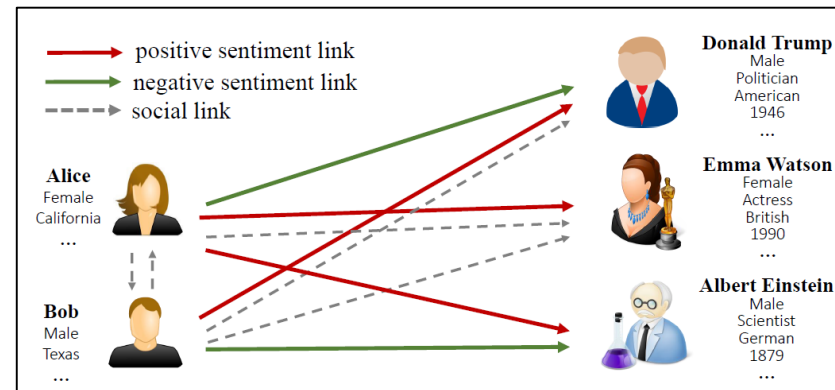
time: 20:10
location: Beijing
What else in carts: ...



Contexts

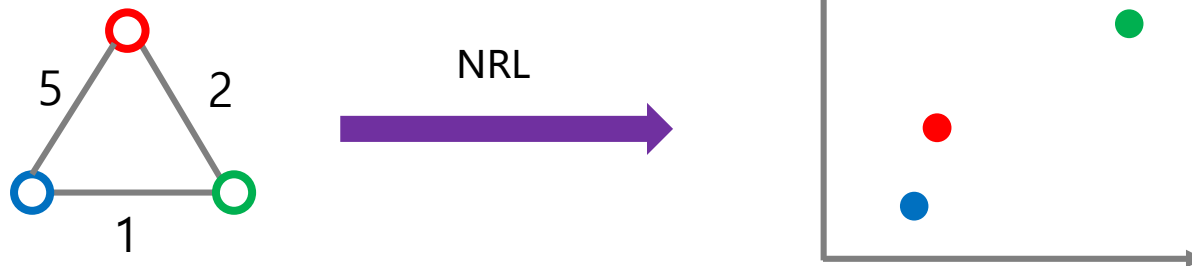
Side Information with Network Structure

Celebrity recommendation in Weibo



Network Representation Learning

- **Network representation learning (NRL)** tries to embed each node of a network into a low-dimensional vector space, which preserves the structural similarities or distances among the nodes in the network
- $G(V, E) \rightarrow \mathbf{E} \in \mathbb{R}^{|V| \times d}$
- NRL can be viewed as a **dimension reduction** technology
- a.k.a network embedding / graph representation learning / graph embedding
- **Applications:** node classification, link prediction, clustering, anomaly detection, social network analysis, etc.

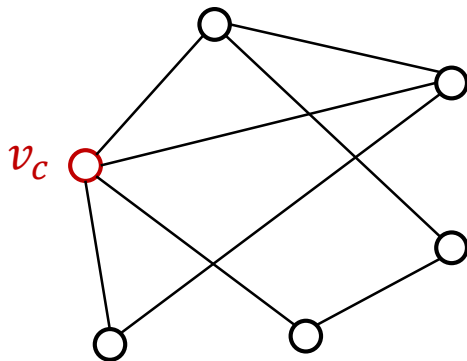


Network Representation Learning

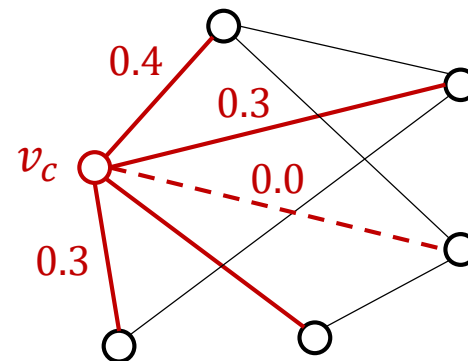
- **Traditional dimension reduction methods**
 - **PCA** (principle component analysis)
 - **LDA** (linear discriminant analysis)
 - **MDS** (multiple dimensional scaling)
- **Manifold Learning methods**
 - **Isomap** (isometric mapping) [Science 2000]
 - **LLE** (locally linear embedding) [Science 2000]
 - **LE** (Laplacian eigenmaps) [NIPS 2001]
- **Random-walk-based methods**
 - **DeepWalk** [KDD 2014]
 - **Node2vec** [KDD 2016]
- **Deep-learning-based methods**
 - **SDNE** (structural deep network embedding) [KDD 2016]
 - **HNE** (heterogeneous network embedding) [KDD 2015]
- **Others**
 - **LINE** (large-scale information network embedding) [WWW 2015]
 - **GraphSAGE** (sample and aggregate) [NIPS 2017]

Generative Models

- **Generative** graph representation learning model assumes an underlying true connectivity distribution $p_{true}(v|v_c)$ for each vertex v_c
 - The edges can be viewed as observed samples generated by $p_{true}(v|v_c)$
 - Vertex embeddings are learned by maximizing the likelihood of edges
 - E.g., DeepWalk and node2vec



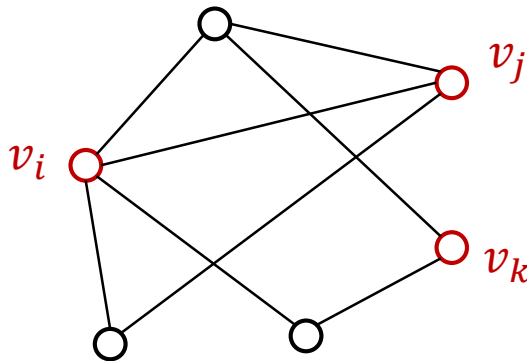
Original graph



$p_{true}(v|v_c)$

Discriminative Models

- **Discriminative** graph representation learning model aims to learn a classifier for predicting edges directly
 - Consider two vertices v_i and v_j jointly as features and predict the probability of an edge existing between them, i.e., $p(\text{edge}|v_i, v_j)$
 - E.g., SDNE and PPNE



$$p(\text{edge}|v_i, v_j) = 0.8$$

$$p(\text{edge}|v_i, v_k) = 0.3$$

.....

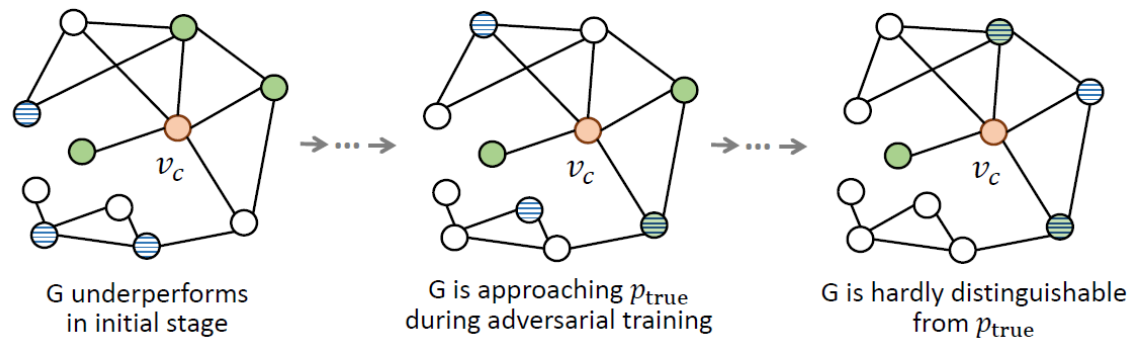
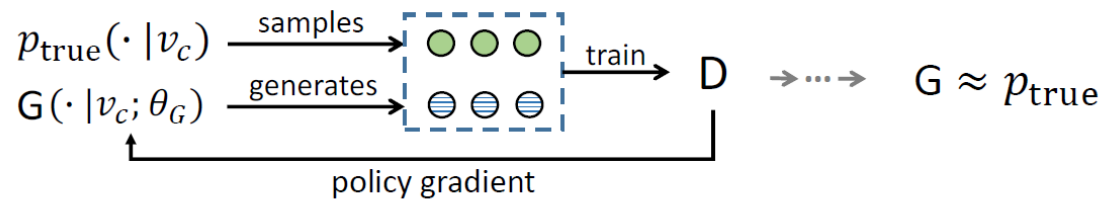
G + D ?

- **GraphGAN**: a framework unifying generative and discriminative models for graph representation learning
- Objective:
 - $G(v | v_c; \theta_G)$: trying to approximate $p_{\text{true}}(v_c)$
 - $D(v, v_c; \theta_D)$: aiming to discriminate the connectivity for the vertex pair (v, v_c)
- The two-player minimax game:

$$\min_{\theta_G} \max_{\theta_D} V(G, D) = \sum_{c=1}^V \left(\mathbb{E}_{v \sim p_{\text{true}}(\cdot | v_c)} [\log D(v, v_c; \theta_D)] + \mathbb{E}_{v \sim G(\cdot | v_c; \theta_G)} [\log (1 - D(v, v_c; \theta_D))] \right)$$

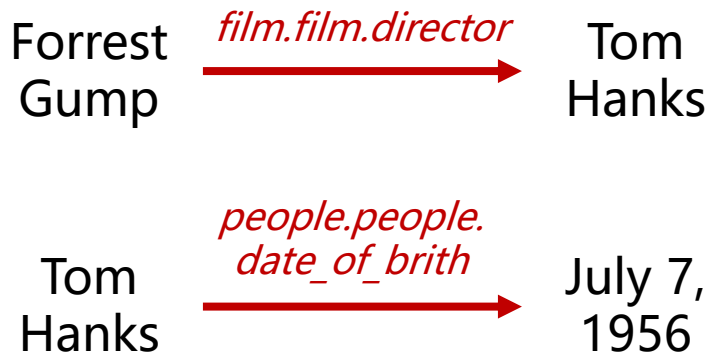
GraphGAN [AAAI 18]

$$\min_{\theta_G} \max_{\theta_D} V(G, D) = \sum_{c=1}^V \left(\mathbb{E}_{v \sim p_{\text{true}}(\cdot | v_c)} [\log D(v, v_c; \theta_D)] + \mathbb{E}_{v \sim G(\cdot | v_c; \theta_G)} [\log (1 - D(v, v_c; \theta_D))] \right)$$

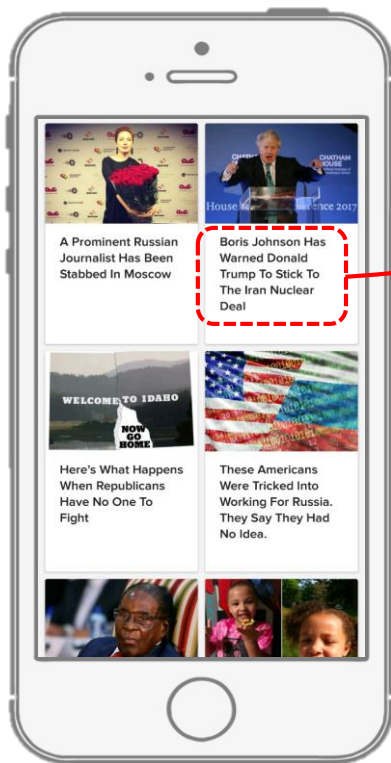


Knowledge Graph

- **Knowledge graph (KG)** is a semantic network in which **nodes** correspond to **entities** and **edges** correspond to **relations**
- A KG usually consists of massive **triples** (head, relation, tail)
- E.g.,

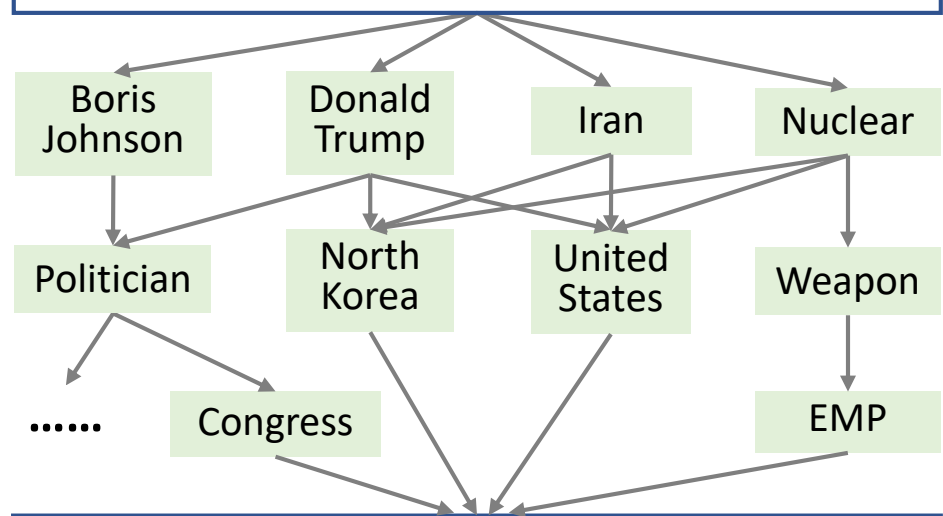


Why Using KG in RS ?



*News the user
have read*

Boris Johnson Has Warned **Donald Trump**
To Stick To The **Iran Nuclear** Deal



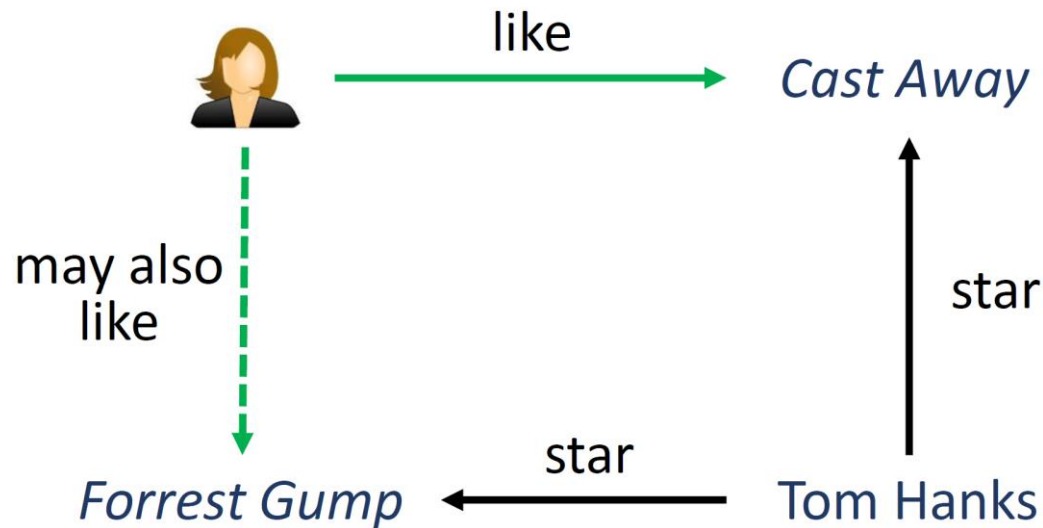
*News the user
may also like*

North Korean EMP Attack Would Cause Mass
U.S. Starvation, Says **Congressional** Report

There is more...

Advantages of using KG in RS

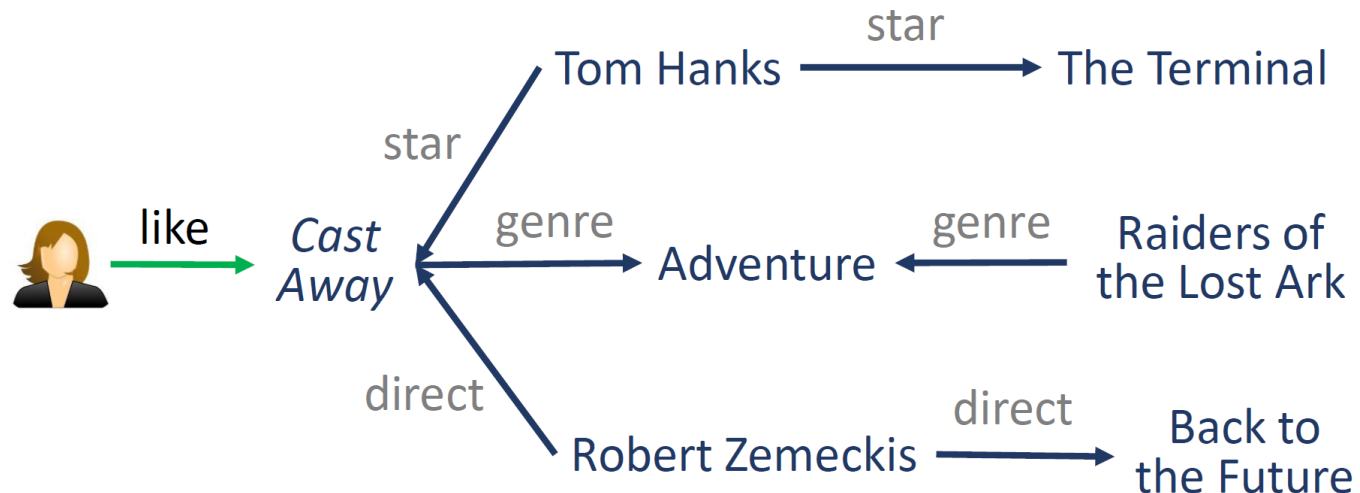
- **Precision**



There is more...

Advantages of using KG in RS

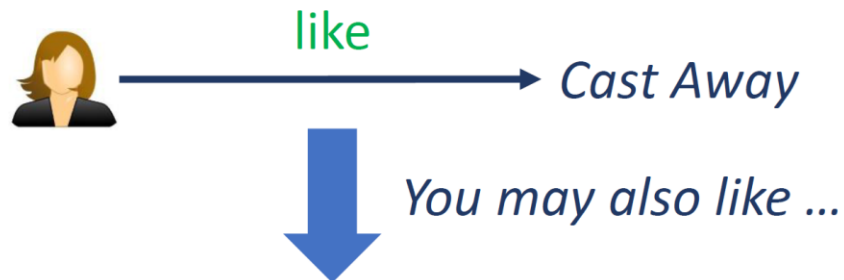
- **Diversity**



There is more...

Advantages of using KG in RS

- **Explainability**

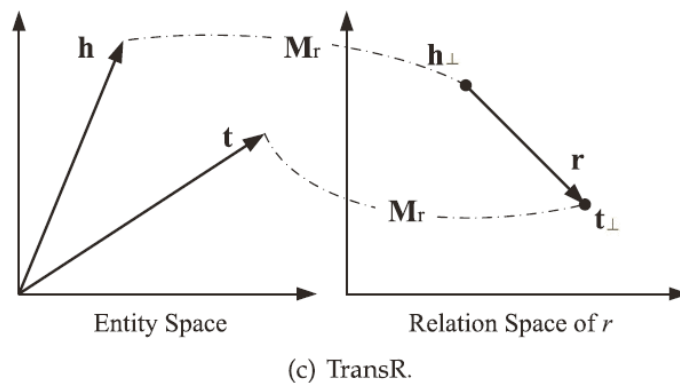
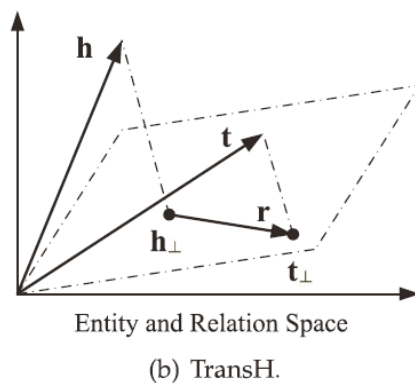
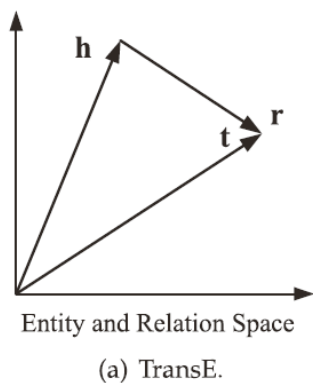


The Terminal, as they share the same **star**;
Raiders of the Lost Ark, as they share the same **genre**;
Back to the Future, as they share the same **director**;

.....

Knowledge Graph Embedding

- **Knowledge graph embedding (KGE)** aims to learn a low-dimensional representation vector for each entity and relation in a KG
- Translational distance models (TransX)
 - TransE: $f_r(h, t) = \|\mathbf{h} + \mathbf{r} - \mathbf{t}\|_2^2$
 - TransH: $f_r(h, t) = \|\mathbf{h}_\perp + \mathbf{r} - \mathbf{t}_\perp\|_2^2$, where $\mathbf{h}_\perp = \mathbf{h} - \mathbf{w}_r^T \mathbf{h} \mathbf{w}_r$ and $\mathbf{t}_\perp = \mathbf{t} - \mathbf{w}_r^T \mathbf{t} \mathbf{w}_r$
 - TransR: $f_r(h, t) = \|\mathbf{h}_r + \mathbf{r} - \mathbf{t}_r\|_2^2$, where $\mathbf{h}_r = \mathbf{h} \mathbf{M}_r$ and $\mathbf{t}_r = \mathbf{t} \mathbf{M}_r$



RS + NRL ?

- If we treat RS and NRL as two tasks, there are three possible ways of combining RS and NRL together...
 - **One-by-one learning**
 - **Alternate learning**
 - **Joint learning**

One-by-one learning KG \xrightarrow{KGE} Entity embeddings
Relation embeddings \xrightarrow{apply} RS \xrightarrow{learn} User representation
Item representation

Joint learning KG } \xrightarrow{KGE} Entity and relation embeddings
RS } \xrightarrow{learn} User and item representation

Alternate learning KG \xrightarrow{KGE} Entity and relation embeddings
RS \xrightarrow{learn} User and item representation

One-by-one Learning: DKN [WWW 2018]

Knowledge distillation

Trump praises **Las Vegas** medical team
Apple CEO Tim Cook: iPhone 8 and **Apple Watch Series 3** are sold out in some places
EU Spain: Juncker does not want **Catalonian** independence
.....

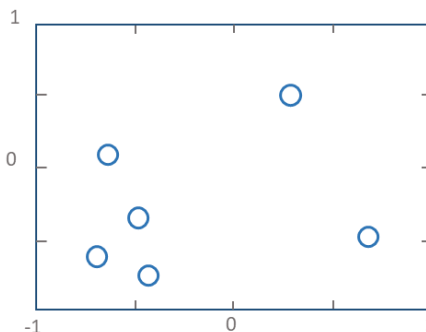
Entity
linking

Donald Trump: Donald Trump is the 45th president ...
Las Vegas: Las Vegas is the 28th-most populated city ...
Apple Inc.: Apple Inc. is an American multinational ...
CEO: A chief executive officer is the position of the ...
Tim Cook: Timothy Cook is an American business ...
iPhone 8: iPhone 8 is smartphone designed, ...
.....

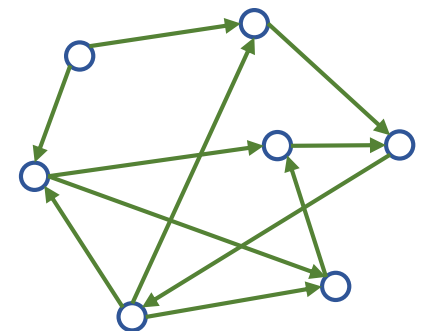
Knowledge subgraph
construction

Donald Trump: (0.32, 0.48)
Las Vegas: (0.71, -0.49)
Apple Inc.: (-0.48, -0.41)
CEO: (-0.57, 0.06)
Tim Cook: (-0.61, -0.59)
iPhone 8: (-0.46, -0.75)

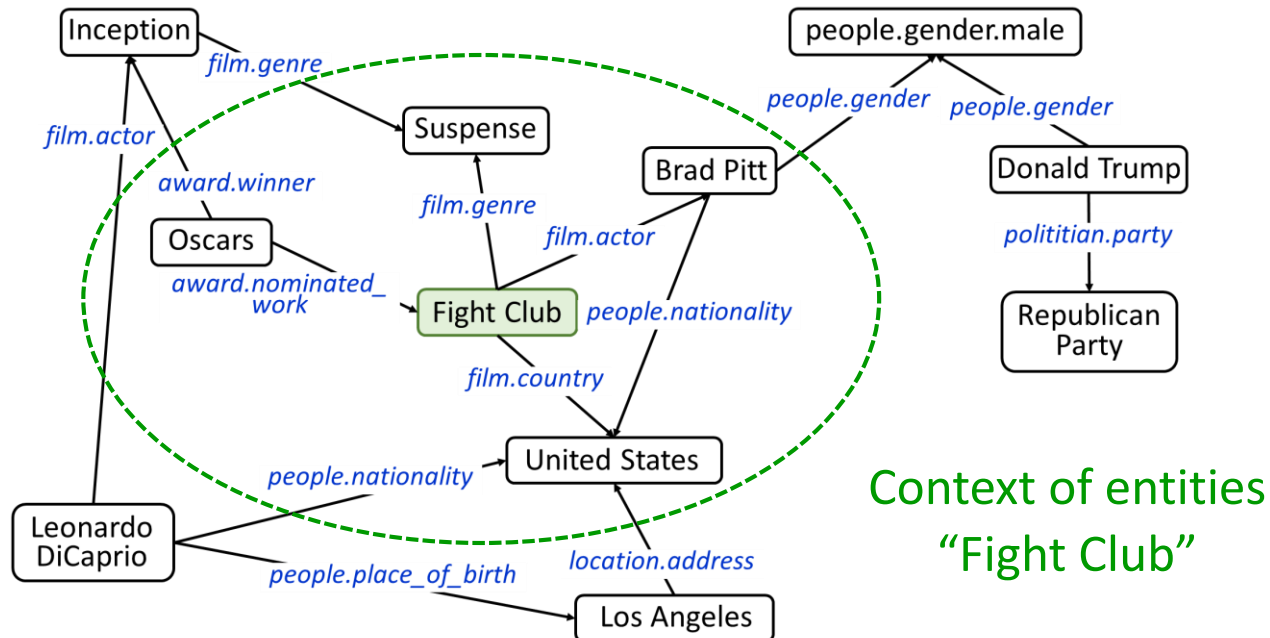
Entity
embedding



Knowledge
graph
embedding



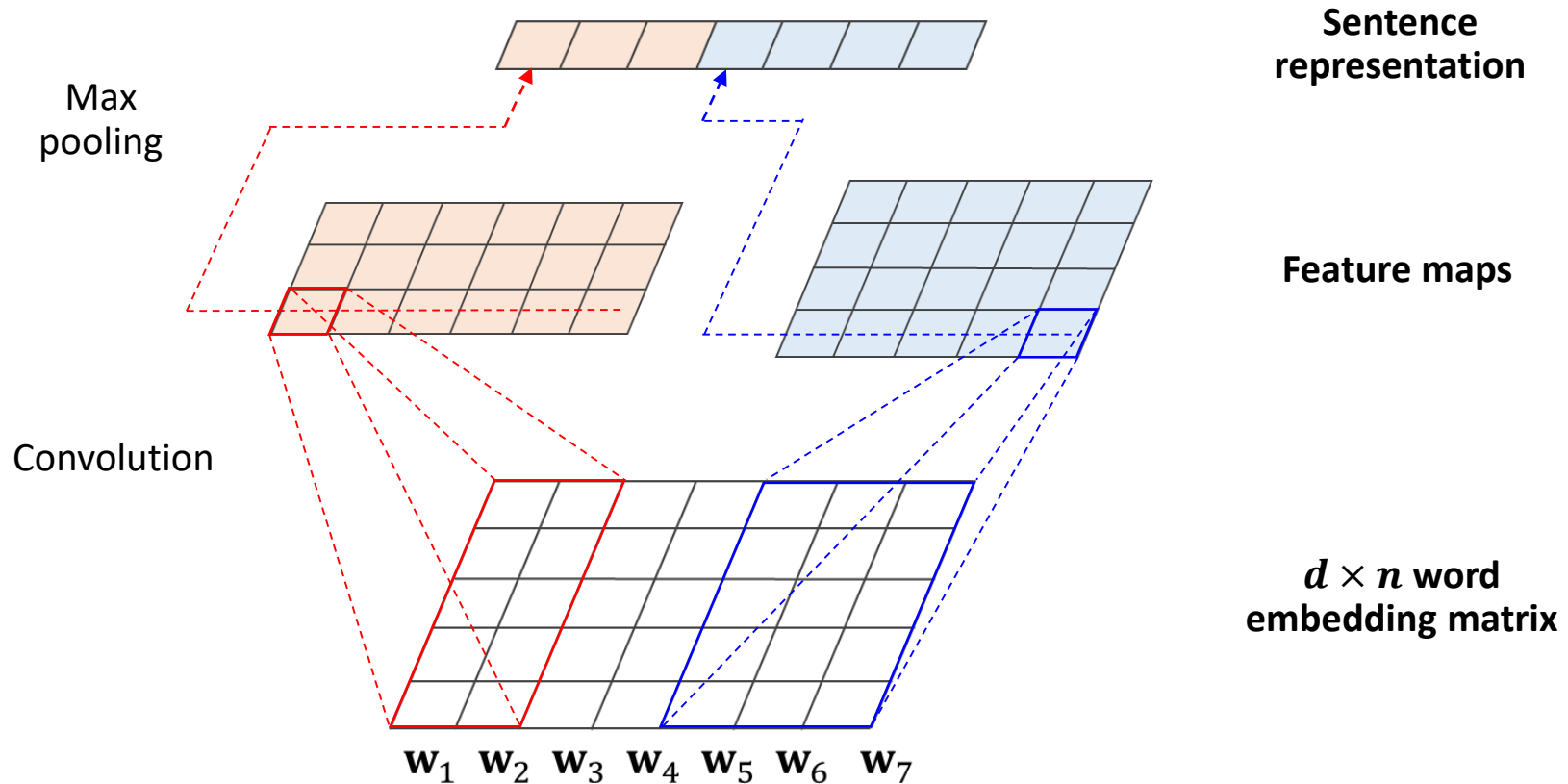
Context embedding



$$\bar{e} = \frac{1}{|context(e)|} \sum_{e_i \in context(e)} e_i$$

DKN

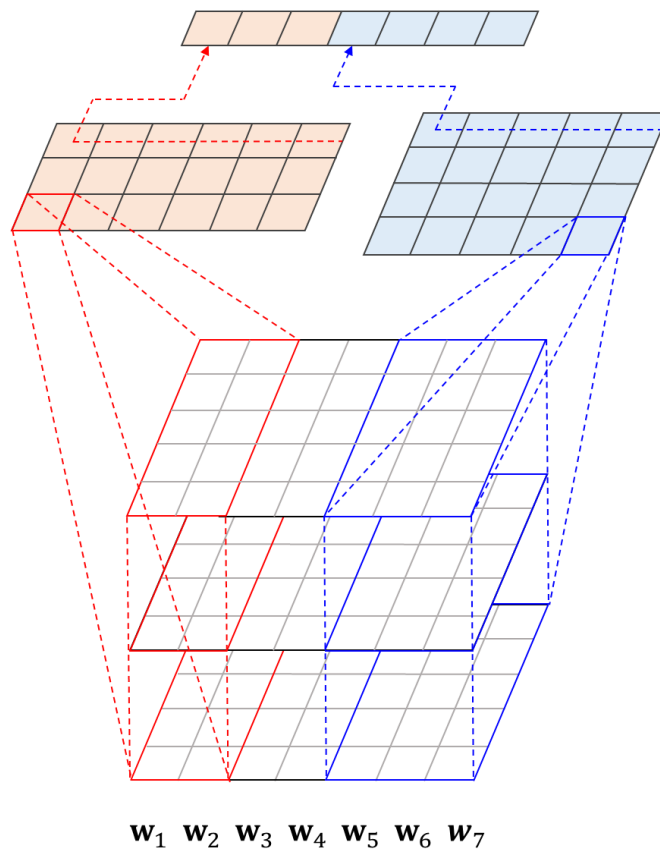
Kim CNN



$w_{1:n} = [Donald Trump praises Las Vegas medical team]$

Sentence

Knowledge-aware CNN



pooling

CNN layer

$d \times n$ transformed
context embeddings

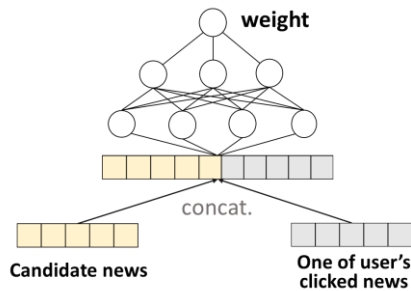
$d \times n$ transformed
entity embeddings

$d \times n$ word
embeddings

multiple
channels

DKN

Attention Network



Attention net:

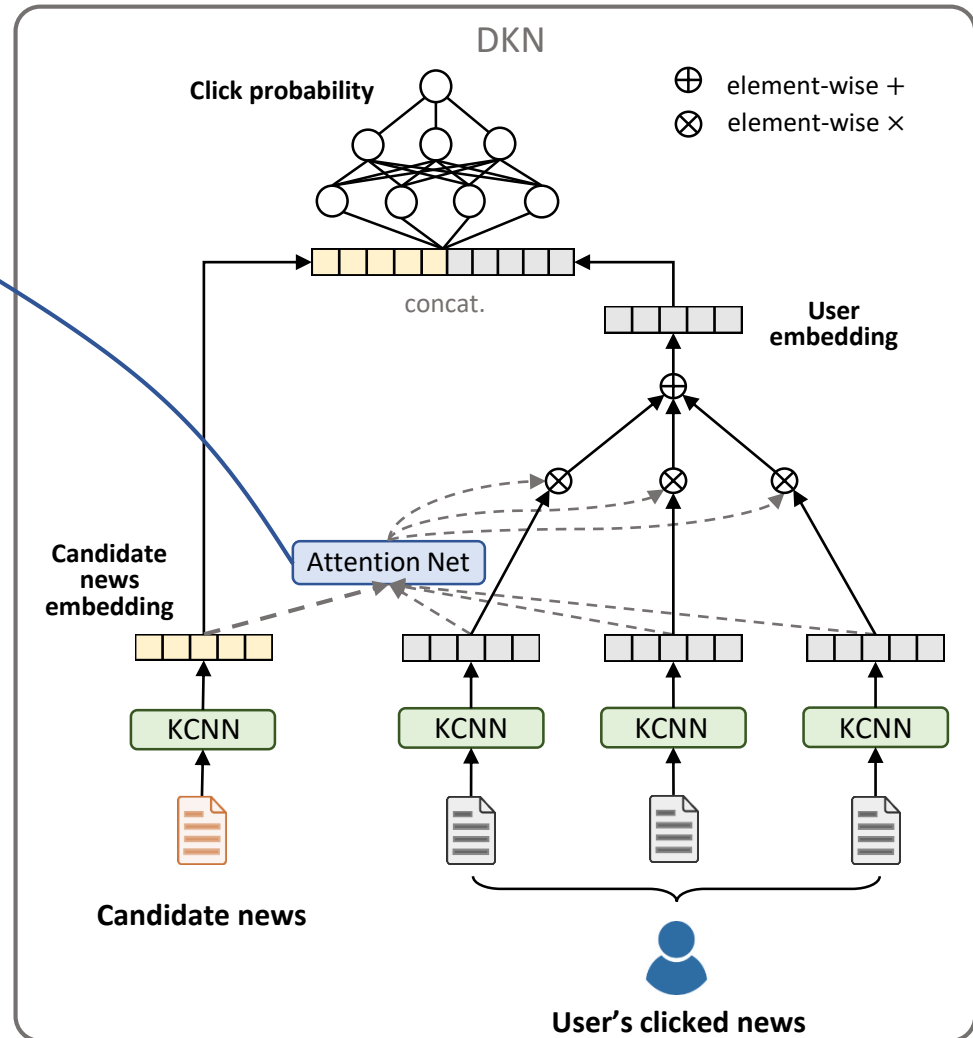
$$s_{t_k^i, t_j} = \text{softmax}(\mathcal{H}(e(t_k^i), e(t_j))) = \frac{\exp(\mathcal{H}(e(t_k^i), e(t_j)))}{\sum_{k=1}^{N_i} \exp(\mathcal{H}(e(t_k^i), e(t_j)))}$$

User interest extraction:

$$e(i) = \sum_{k=1}^{N_i} s_{t_k^i, t_j} e(t_k^i).$$

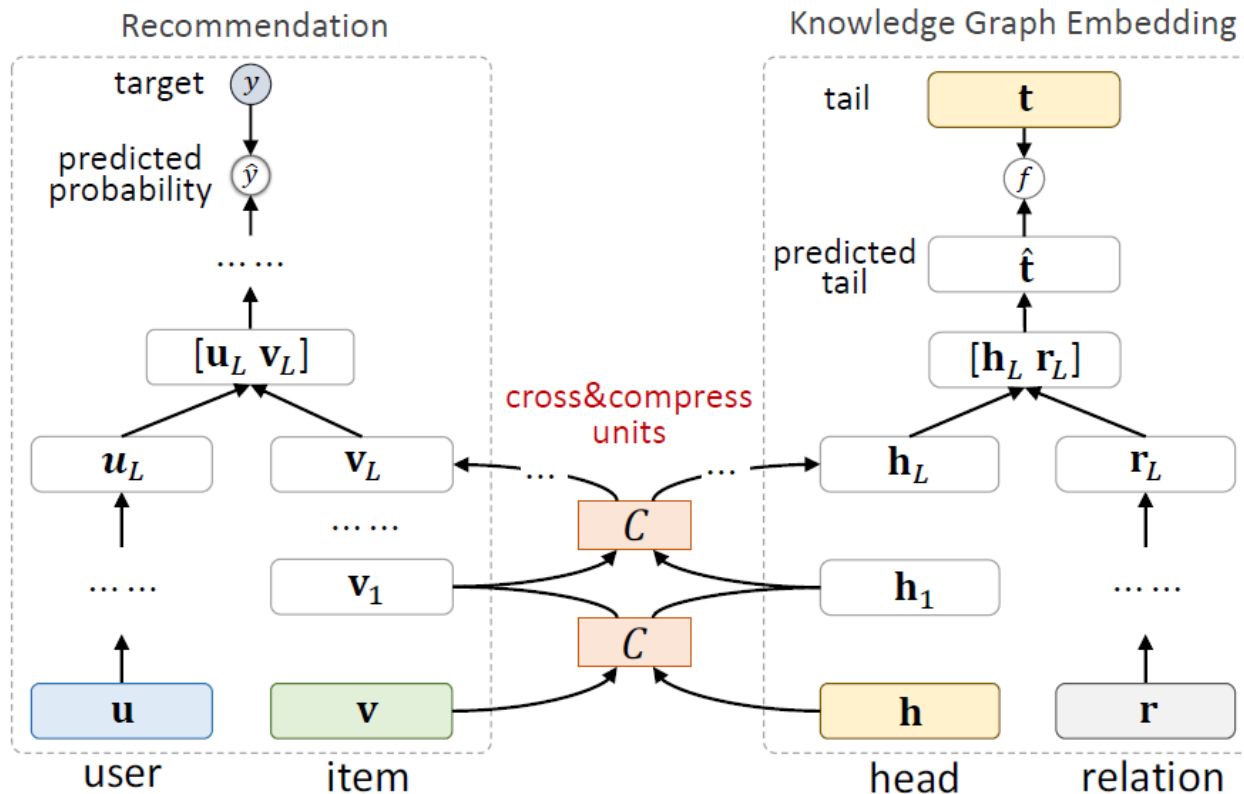
CTR prediction:

$$p_{i, t_j} = \mathcal{G}(e(i), e(t_j))$$

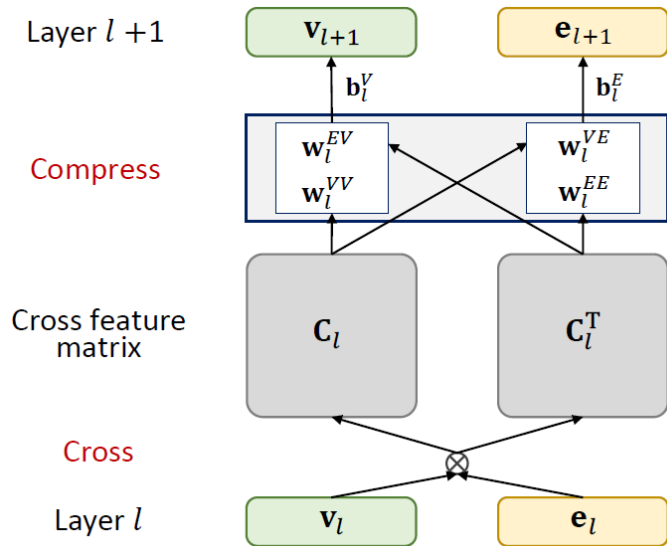


Alternate Learning: MKR [NIPS in sub.]

Multi-task learning for Knowledge graph enhanced Recommendation (MKR)



Cross&compress unit



$$\mathbf{C}_l = \mathbf{v}_l \mathbf{e}_l^T = \begin{bmatrix} v_l^{(1)} e_l^{(1)} & \cdots & v_l^{(1)} e_l^{(d)} \\ \vdots & & \vdots \\ v_l^{(d)} e_l^{(1)} & \cdots & v_l^{(d)} e_l^{(d)} \end{bmatrix}$$

$$\begin{aligned} \mathbf{v}_{l+1} &= \mathbf{C}_l \mathbf{w}_l^{VV} + \mathbf{C}_l^T \mathbf{w}_l^{EV} + \mathbf{b}_l^V = \mathbf{v}_l \mathbf{e}_l^T \mathbf{w}_l^{VV} + \mathbf{e}_l \mathbf{v}_l^T \mathbf{w}_l^{EV} + \mathbf{b}_l^V \\ \mathbf{e}_{l+1} &= \mathbf{C}_l \mathbf{w}_l^{VE} + \mathbf{C}_l^T \mathbf{w}_l^{EE} + \mathbf{b}_l^E = \mathbf{v}_l \mathbf{e}_l^T \mathbf{w}_l^{VE} + \mathbf{e}_l \mathbf{v}_l^T \mathbf{w}_l^{EE} + \mathbf{b}_l^E \end{aligned}$$

Theoretical Analysis

- Polynomial approximation

Theorem 1 Denote the input of item and entity in MKR network as $\mathbf{v} = [v_1 \cdots v_d]^\top$ and $\mathbf{e} = [e_1 \cdots e_d]^\top$, respectively. Then the cross terms about \mathbf{v} and \mathbf{e} in $\|\mathbf{v}_L\|_1$ and $\|\mathbf{e}_L\|_1$ (the $L1$ -norm of \mathbf{v}_L and \mathbf{e}_L) with maximal degree is $k_{\alpha,\beta} v_1^{\alpha_1} \cdots v_d^{\alpha_d} e_1^{\beta_1} \cdots e_d^{\beta_d}$, where $k_{\alpha,\beta} \in \mathbb{R}$, $\alpha_i, \beta_i \in \mathbb{N}$ for $i \in \{1, \dots, d\}$, $\alpha_1 + \cdots + \alpha_d = 2^{L-1}$ and $\beta_1 + \cdots + \beta_d = 2^{L-1}$ ($L \geq 1, \mathbf{v}_0 = \mathbf{v}, \mathbf{e}_0 = \mathbf{e}$).

sufficient approximation ability

Theoretical Analysis

- **Generalization**

- **Factorization Machines**

Proposition 1 *The L1-norm of \mathbf{v}_1 and \mathbf{e}_1 can be written as the following form:*

$$\|\mathbf{v}_1\|_1 \text{ (or } \|\mathbf{e}_1\|_1) = \left| b + \sum_{i=1}^d \sum_{j=1}^d \langle w_i, w_j \rangle v_i e_j \right|,$$

where $\langle w_i, w_j \rangle = w_i + w_j$.

- **Deep&Cross Network**

Proposition 2 *In the formula of \mathbf{v}_{l+1} in Eq. (2), if we restrict \mathbf{w}_l^{VV} in the first term to satisfy $\mathbf{e}_l^\top \mathbf{w}_l^{VV} = 1$ and restrict \mathbf{e}_l in the second term to be \mathbf{e}_0 (and impose similar restrictions on \mathbf{e}_{l+1}), the cross&compress unit is then conceptually equivalent to DCN layer in the sense of multi-task learning:*

$$\mathbf{v}_{l+1} = \mathbf{e}_0 \mathbf{v}_l^\top \mathbf{w}_l^{EV} + \mathbf{v}_l + \mathbf{b}_l^V, \quad \mathbf{e}_{l+1} = \mathbf{v}_0 \mathbf{e}_l^\top \mathbf{w}_l^{VE} + \mathbf{e}_l + \mathbf{b}_l^E.$$

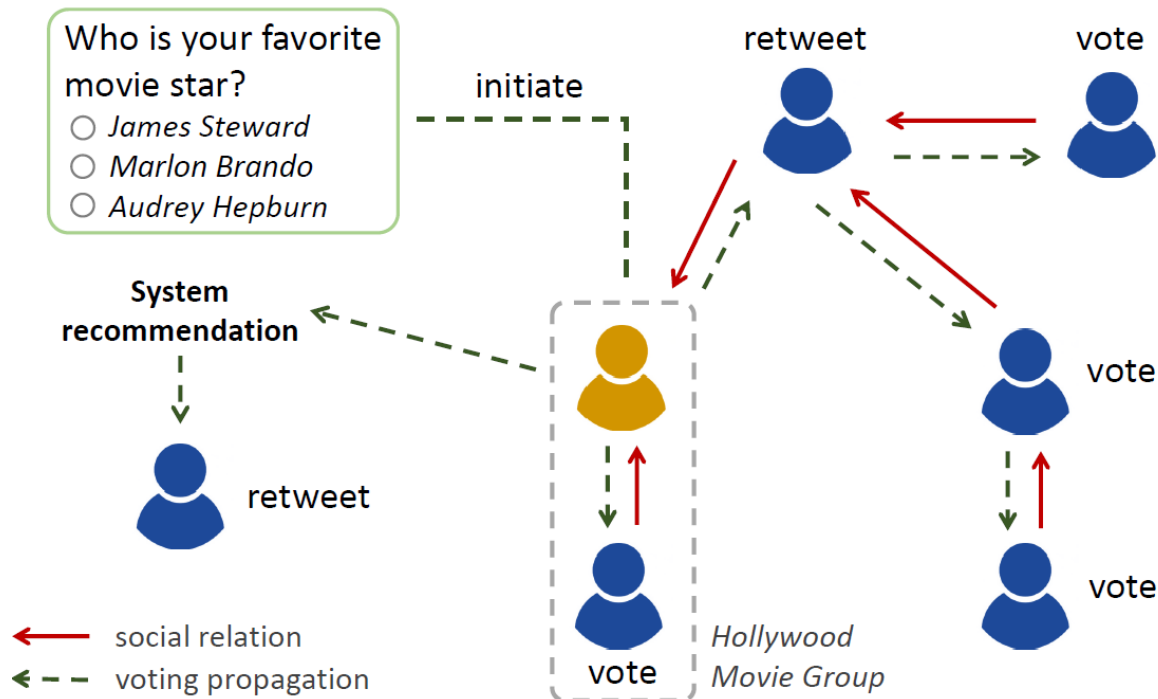
- **Cross-stitch Network**

Proposition 3 *If we omit all biases in Eq. (2), the cross&compress unit can be written as*

$$\begin{bmatrix} \mathbf{v}_{l+1} \\ \mathbf{e}_{l+1} \end{bmatrix} = \begin{bmatrix} \mathbf{e}_l^\top \mathbf{w}_l^{VV} & \mathbf{v}_l^\top \mathbf{w}_l^{EV} \\ \mathbf{e}_l^\top \mathbf{w}_l^{VE} & \mathbf{v}_l^\top \mathbf{w}_l^{EE} \end{bmatrix} \begin{bmatrix} \mathbf{v}_l \\ \mathbf{e}_l \end{bmatrix}.$$

Joint Learning: JTS-MF [CIKM 2017]

Joint Topic-Semantic-aware Matrix Factorization (JTS-MF)



Joint Topic-Semantic-aware Matrix Factorization (JTS-MF)

$$\begin{aligned} L = & \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^M I'_{i,j} \left(R_{i,j} - \mathbf{Q}_i \mathbf{P}_j^\top \right)^2 + \frac{\alpha}{2} \sum_{i=1}^N \left\| \mathbf{Q}_i - \sum_{k \in \mathcal{F}_i^+} \hat{S}_{i,k} \mathbf{Q}_k \right\|_2^2 \\ & + \frac{\beta}{2} \sum_{i=1}^N \left\| \mathbf{Q}_i - \sum_{k \in \mathcal{G}_i} \hat{G}_{i,k} \mathbf{Q}_k \right\|_2^2 + \frac{\gamma}{2} \sum_{j=1}^M \left\| \mathbf{P}_j - \sum_{t \in \mathcal{V}_j} \hat{T}_{j,t} \mathbf{P}_t \right\|_2^2 + \frac{\lambda}{2} \left(\|\mathbf{Q}\|_F^2 + \|\mathbf{P}\|_F^2 \right) \end{aligned}$$

Q: user feature matrix

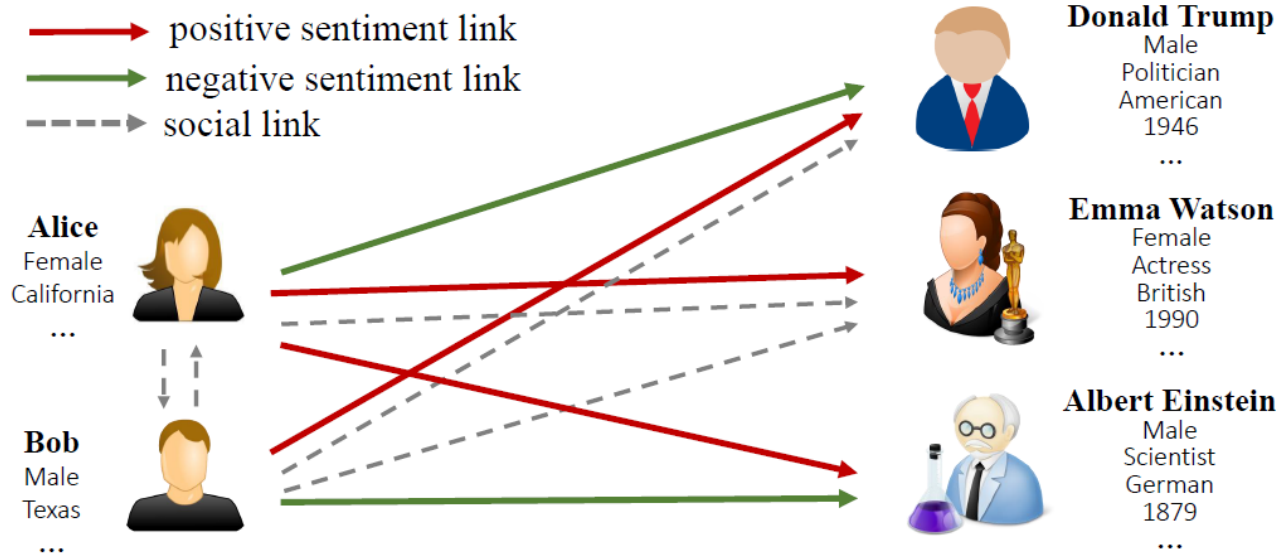
P: item feature matrix

$R_{i,j}$: rating

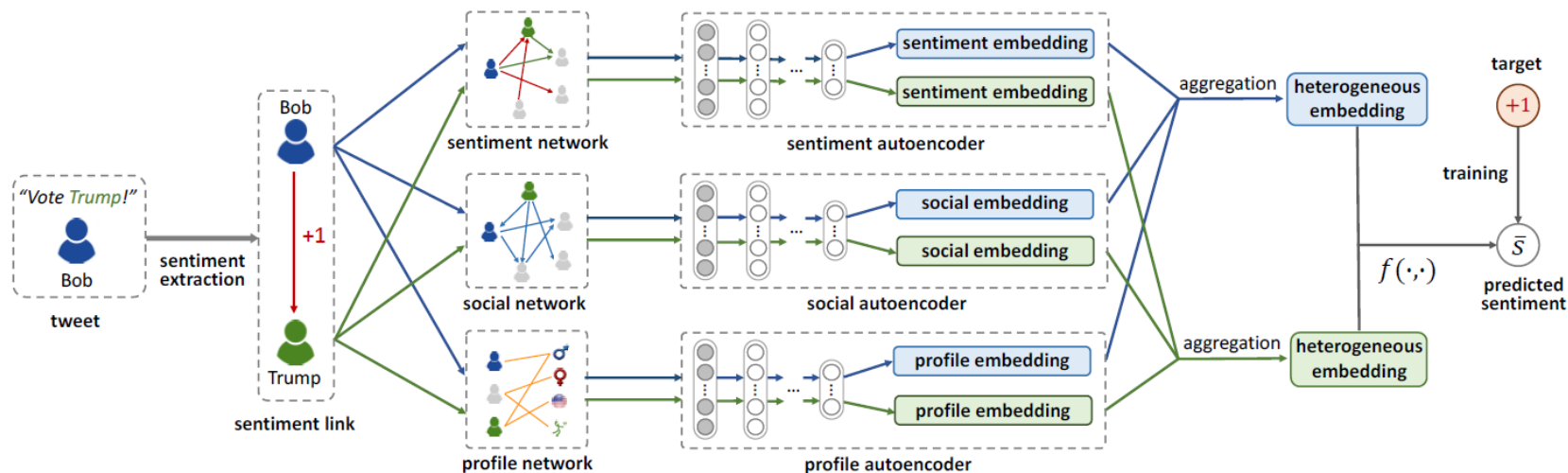
$\hat{S}, \hat{G}, \hat{T}$: similarity coefficient

Joint Learning: SHINE [WSDM 2018]

Signed Heterogeneous Information Network Embedding (SHINE)



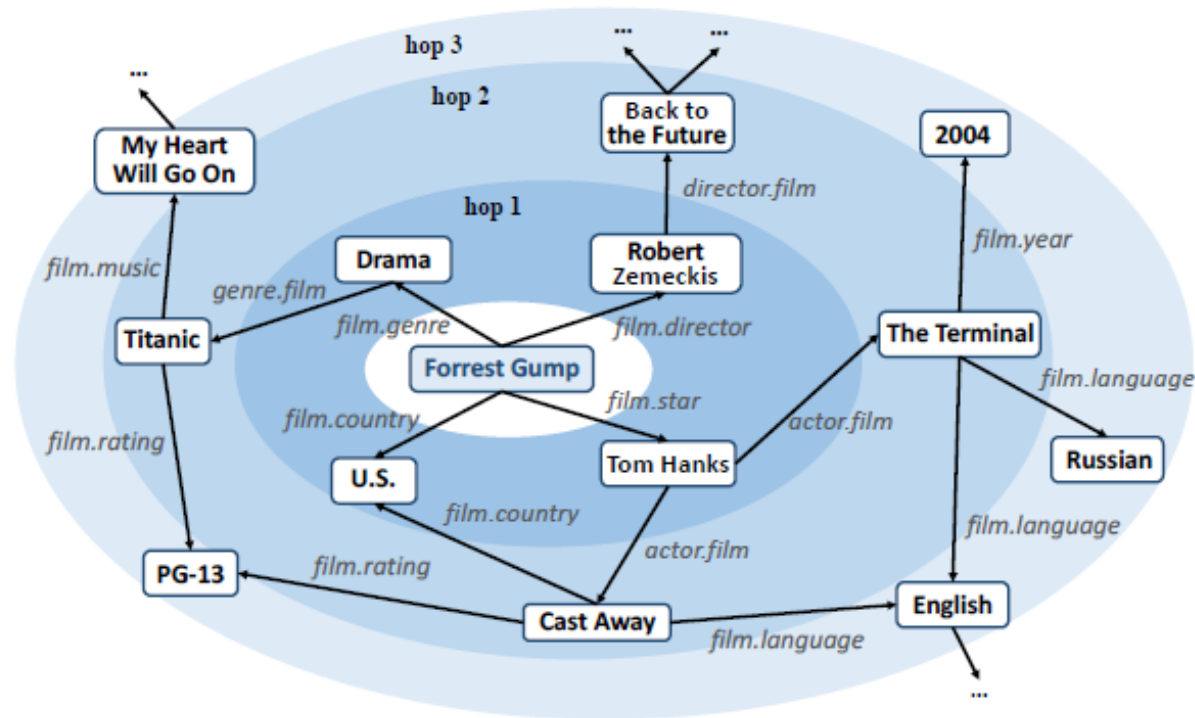
Signed Heterogeneous Information Network Embedding (SHINE)



$$\begin{aligned} \mathcal{L} = & \sum_{i \in V} \|(\mathbf{x}_i - \mathbf{x}'_i) \odot \mathbf{l}_i\|_2^2 + \lambda_1 \sum_{i \in V} \|(\mathbf{y}_i - \mathbf{y}'_i) \odot \mathbf{m}_i\|_2^2 \\ & + \lambda_2 \sum_{i \in V} \|(\mathbf{z}_i - \mathbf{z}'_i) \odot \mathbf{n}_i\|_2^2 + \lambda_3 \sum_{s_{ij} = \pm 1} (f(\mathbf{e}_i, \mathbf{e}_j) - s_{ij})^2 \\ & + \lambda_4 \mathcal{L}_{reg}, \end{aligned}$$

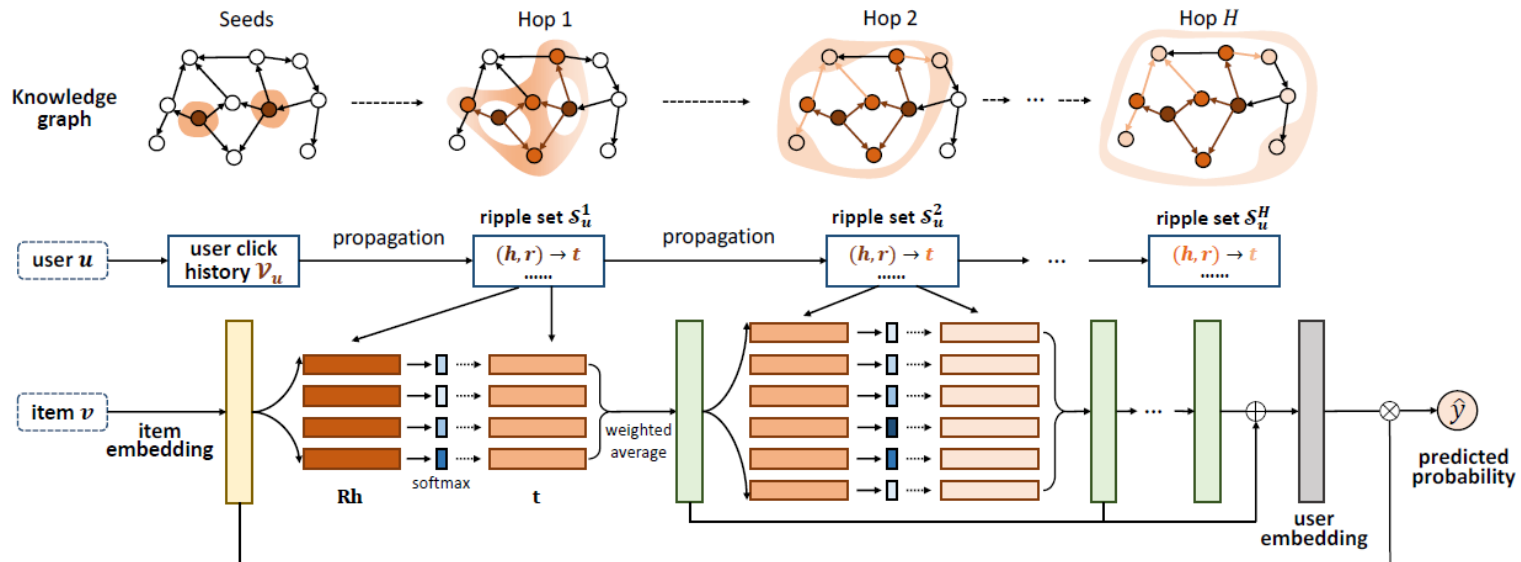
Joint Learning: RippleNet [CIKM in sub.]

Ripple Network



RippleNet

Ripple Network



$$\begin{aligned}
 \min \mathcal{L} &= -\log(p(Y|\Theta, \mathcal{G}) \cdot p(\mathcal{G}|\Theta) \cdot p(\Theta)) \\
 &= \sum_{(u,v) \in Y} -\left(y_{uv} \log \sigma(\mathbf{u}^T \mathbf{v}) + (1 - y_{uv}) \log (1 - \sigma(\mathbf{u}^T \mathbf{v}))\right) \\
 &\quad + \frac{\lambda_2}{2} \sum_{r \in \mathcal{R}} \|\mathbf{I}_r - \mathbf{E}^T \mathbf{R} \mathbf{E}\|_2^2 + \frac{\lambda_1}{2} \left(\|\mathbf{V}\|_2^2 + \|\mathbf{E}\|_2^2 + \sum_{r \in \mathcal{R}} \|\mathbf{R}\|_2^2 \right)
 \end{aligned}$$

Discussion

In general...

- **Efficiency: one-by-one > alternate > joint**
 - The updating frequency of KG is far less than RS...
 - All embeddings need to be learned from scratch in joint learning...
- **Performance: joint > alternate > one-by-one**
 - Joint learning methods are end-to-end...
 - Entity and relation embeddings are learned in advance in one-by-one learning, lacking the supervision from RS...

Summary

- **Recommender systems** are key technique in user-oriented web services
- **Network-structured data** are common in RS
 - User-item interaction, social network, knowledge graph
- **Network representation learning**: a popular technique of processing network-structured data
 - GraphGAN
 - Knowledge graph embedding
- RS + NRL:
 - **One-by-one learning**
 - **Alternate learning**
 - **Joint learning**

Q & A

Thanks!

