

DSC250: Advanced Data Mining

Knowledge Graphs (KGs)

Zhiting Hu

Lecture 15, Feb 25, 2025

UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

Outline

- Knowledge graphs
- Presentation
 - Boyang Wang, Zixuan Song: "Inference-Time Intervention: Eliciting Truthful Answers from a Language Model"
 - Benjamin TenWolde, Michael Ko: "Interpreting and Editing Vision-Language Representations to Mitigate Hallucinations"
 - Jay Sawant, Sarthak Kala: "Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding"
 - Dongting Cai, Zhihan Li: "Music Style Analysis among Haydn, Mozart and Beethoven: an Unsupervised Machine Learning Approach"

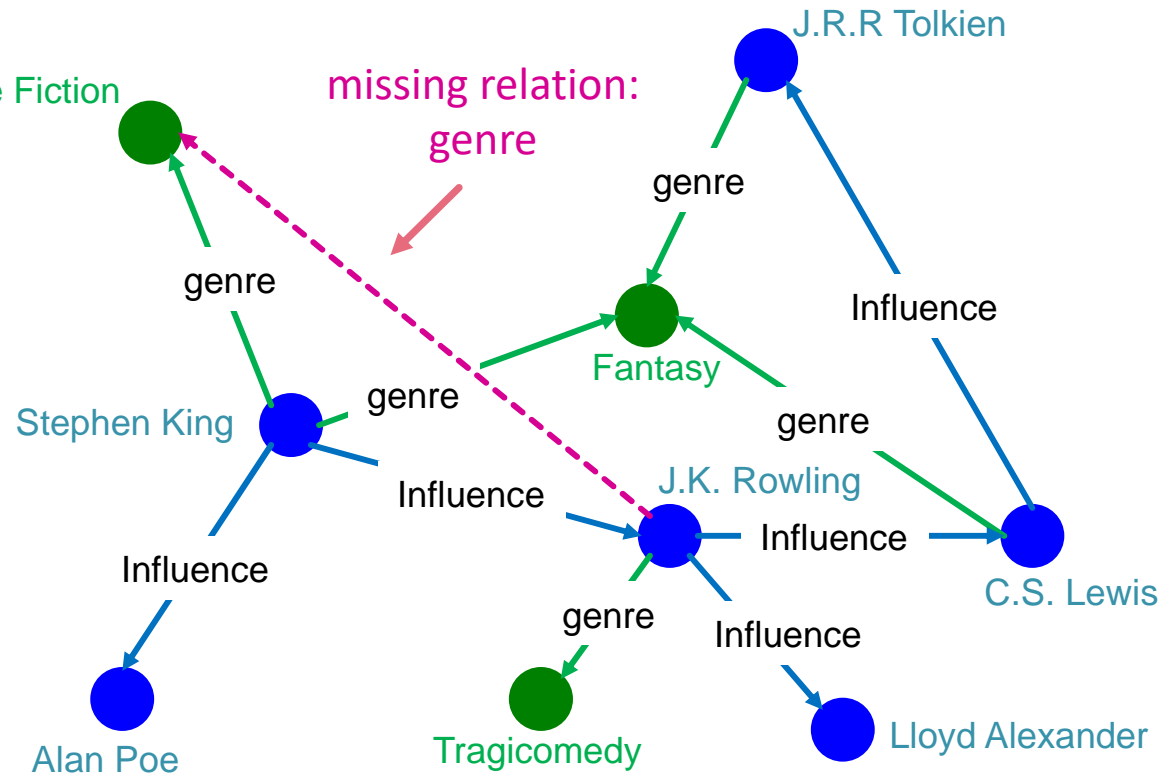
Recap: KG Completion Task

Given an enormous KG, can we complete the KG?

§ For a given (**head**, **relation**), we predict missing **tails**.

§ (Note this is slightly different from link prediction task)

Example task: predict the tail “Science Fiction” for (“J.K. Rowling”, “genre”)



Recap: KG Representation

- i Edges in KG are represented as **triples** (h, r, t)
 - § **head** (h) has **relation** (r) with **tail** (t)
- i **Key Idea:**
 - § Given a triple (h, r, t) , the goal is that the **embedding of (h, r) should be close** to the **embedding of t**
 - § How to embed (h, r) ?
 - § How to define score $s_r(h, t)$?
 - § Score s_r is high if (h, r, t) exists, else s_r is low

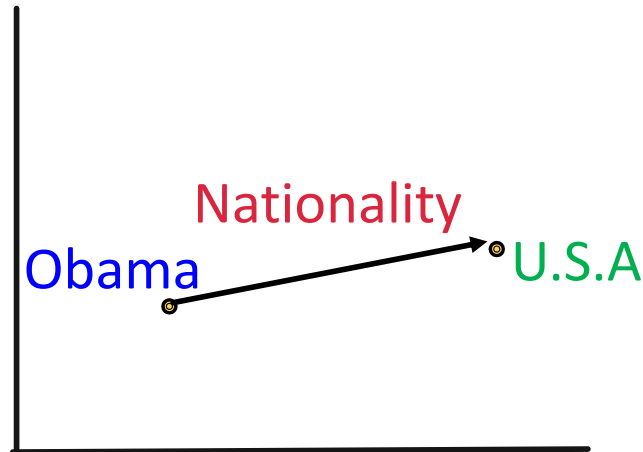
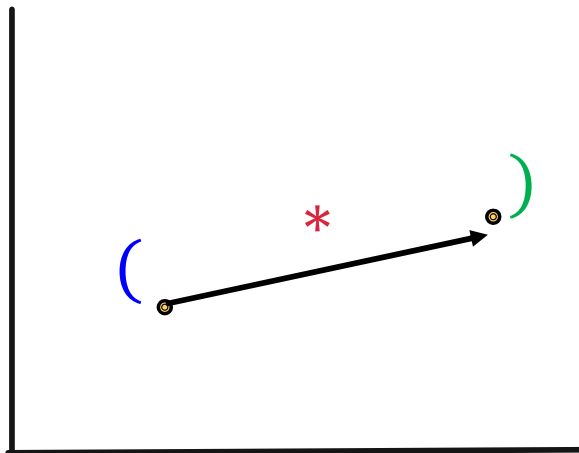
Recap: TransE for KG Completion

i Intuition: Translation

For a triple (h, r, t) , let $\mathbf{h}, \mathbf{r}, \mathbf{t} \in \mathbb{R}^d$ be embedding vectors.

embedding vectors
will appear in
boldface

i **TransE: $\mathbf{h} + \mathbf{r} \approx \mathbf{t}$** if the given link exists else $\mathbf{h} + \mathbf{r} \neq \mathbf{t}$



Four Relationship Patterns

i **Symmetric (Antisymmetric) Relations:**

$$r(h, t) \Rightarrow r(t, h) \quad (r(h, t) \Rightarrow \neg r(t, h)) \quad \forall h, t$$

§ **Example:**

§ Symmetric: Family, Roommate

§ Antisymmetric: Hypernym (a word with a broader meaning: poodle vs. dog)

i **Inverse Relations:**

$$r_{\downarrow}(h, t) \Rightarrow r_{\uparrow}(t, h)$$

§ **Example** : (Advisor, Advisee)

i **Composition (Transitive) Relations:**

$$r_{\downarrow}(x, y) \wedge r_{\downarrow}(y, z) \Rightarrow r_{*}(x, z) \quad \forall x, y, z$$

§ **Example**: My mother's husband is my father.

i **1-to-N relations:**

$r(h, t_{\downarrow}), r(h, t_{\downarrow}), \dots, r(h, t_{+})$ are all True.

§ **Example**: # is "StudentsOf"

Antisymmetric Relations in TransE

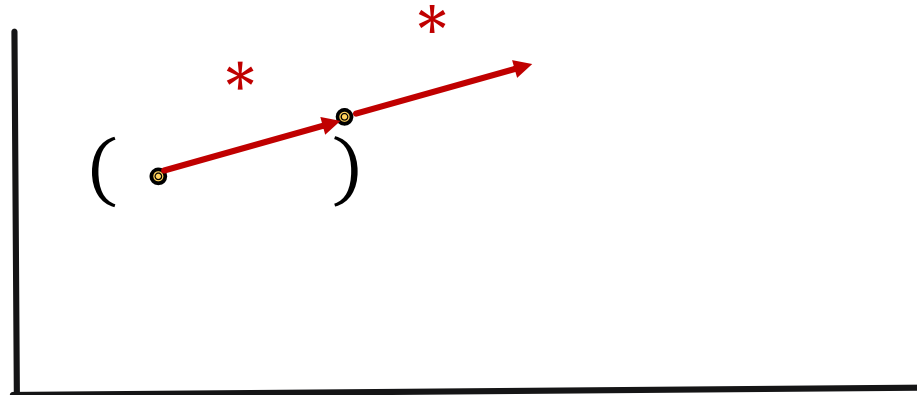
- Antisymmetric Relations:

$$h \text{ } \phi \Rightarrow \neg (\phi h) \quad \square \forall h, \phi$$

§ **Example:** Hypernym (a word with a broader meaning: poodle vs. dog)

- TransE can model antisymmetric relations \ddot{u}

§ $(+ * =)$, but $) + * \neq ($



Inverse Relations in TransE

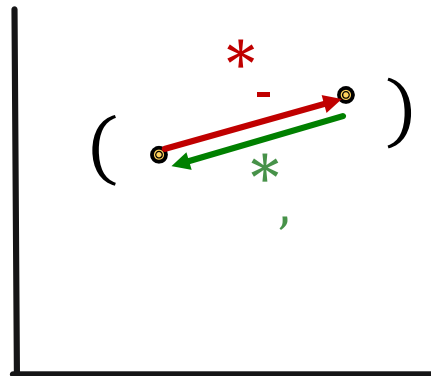
- **Inverse Relations:**

$$r_{\#}(h, \phi) \Rightarrow r_{\$}(\phi, h)$$

§ **Example** : (Advisor, Advisee)

- **TransE** can model inverse relations \ddot{u}

§ $(+ *, =)$, we can set $*) = -*($



Composition in TransE

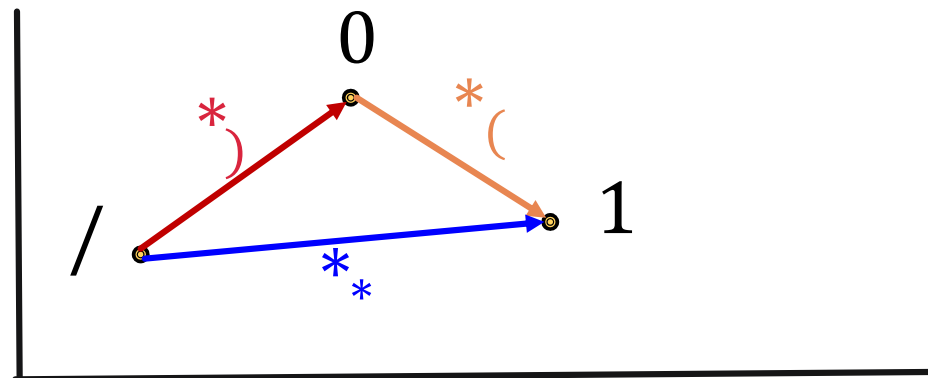
- Composition (Transitive) Relations:**

$$R_{\$} (7, 8) \wedge R_{\#} (8, :) \Rightarrow R_{\%} (7, :) \quad \forall 7, 8, :$$

§ **Example:** My mother's husband is my father.

- TransE** can model composition relations

$$R_{\%} = R_{\$} + R_{\#}$$



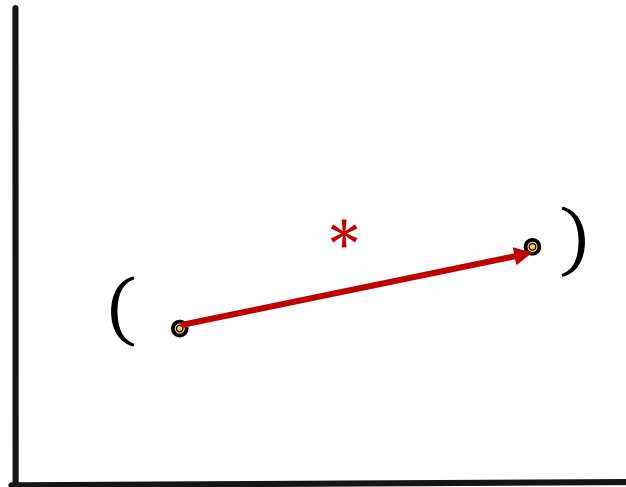
Limitations of TransE: Symmetric Relations

i Symmetric Relations:

$$r(h, t) \Rightarrow r(t, h) \quad \forall h, t$$

§ Example: Family, Roommate

i TransE cannot model symmetric relations \hat{U} only if $(r = 0, t = h)$



For all h, t that satisfy $r(h, t)$, $r(t, h)$ is also True, which means $\|\mathbf{h} + \mathbf{r} - \mathbf{t}\| = 0$ and $\|\mathbf{t} + \mathbf{r} - \mathbf{h}\| = 0$. Then $\mathbf{r} = 0$ and $\mathbf{h} = \mathbf{t}$, however h and t are two different entities and should be mapped to different locations.

Limitations of TransE: 1-to-N Relations

i 1-to-N Relations:

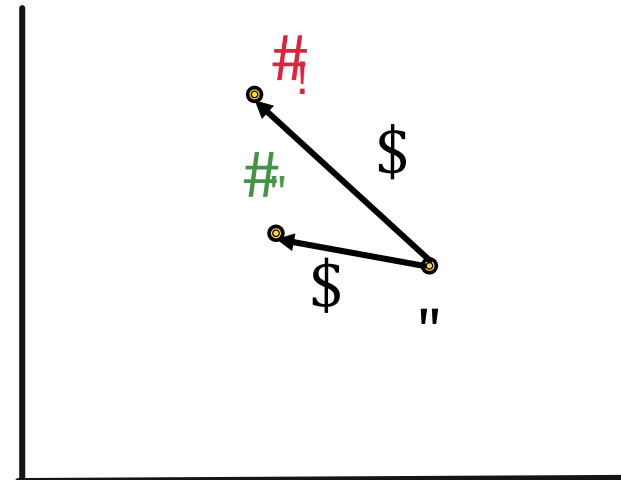
§ **Example:** (h, r, t_j) and (h, r, t_c) both exist in the knowledge graph, e.g., r is “StudentsOf”

i TransE cannot model 1-to-N relations \hat{u}

§ t_j and t_c will map to the same vector, although they are different entities

i $t_j = h + r = t_c$

i $t_j \neq t_c$ **contradictory!**



KG Completion Methods

Model	Score	Embedding	Sym.	Antisym.	Inv.	Compos.	1-to-N
TransE	$-\ \mathbf{h} + \mathbf{r} - \mathbf{t}\ $	$\mathbf{h}, \mathbf{t}, \mathbf{r} \in \mathbb{R}^l$	\hat{u}	\ddot{u}	\ddot{u}	\ddot{u}	\hat{u}
TransR	$-\ M_{\mathbf{r}} \mathbf{h} + \mathbf{r} - M_{\mathbf{t}} \mathbf{t}\ $	$\mathbf{h}, \mathbf{t} \in \mathbb{R}^l,$ $\mathbf{r} \in \mathbb{R}^{\#},$ $M_{\mathbf{r}}, M_{\mathbf{t}} \in \mathbb{R}^{\# \times l}$	\ddot{u}	\ddot{u}	\ddot{u}	\ddot{u}	\ddot{u}
DistMult	$\langle \mathbf{h}, \mathbf{r}, \mathbf{t} \rangle$	$\mathbf{h}, \mathbf{t}, \mathbf{r} \in \mathbb{R}^l$	\ddot{u}	\hat{u}	\hat{u}	\hat{u}	\ddot{u}
Complex	$\text{Re}(\langle \mathbf{h}, \mathbf{r}, \bar{\mathbf{t}} \rangle)$	$\mathbf{h}, \mathbf{t}, \mathbf{r} \in \mathbb{C}^l$	\ddot{u}	\ddot{u}	\ddot{u}	\hat{u}	\ddot{u}
RotateE	$-\ \mathbf{h} \circ \mathbf{r} - \mathbf{t}\ $	$\mathbf{h}, \mathbf{t}, \mathbf{r} \in \mathbb{C}^l$	\ddot{u}	\ddot{u}	\ddot{u}	\ddot{u}	\ddot{u}

Outline

- Overview
- Knowledge Graph Completion (Link Prediction)
- Reasoning on Knowledge Graphs

Reasoning over KGs

- i **Goal:**

- § How to perform multi-hop reasoning over KGs?

- i **Reasoning over Knowledge Graphs**

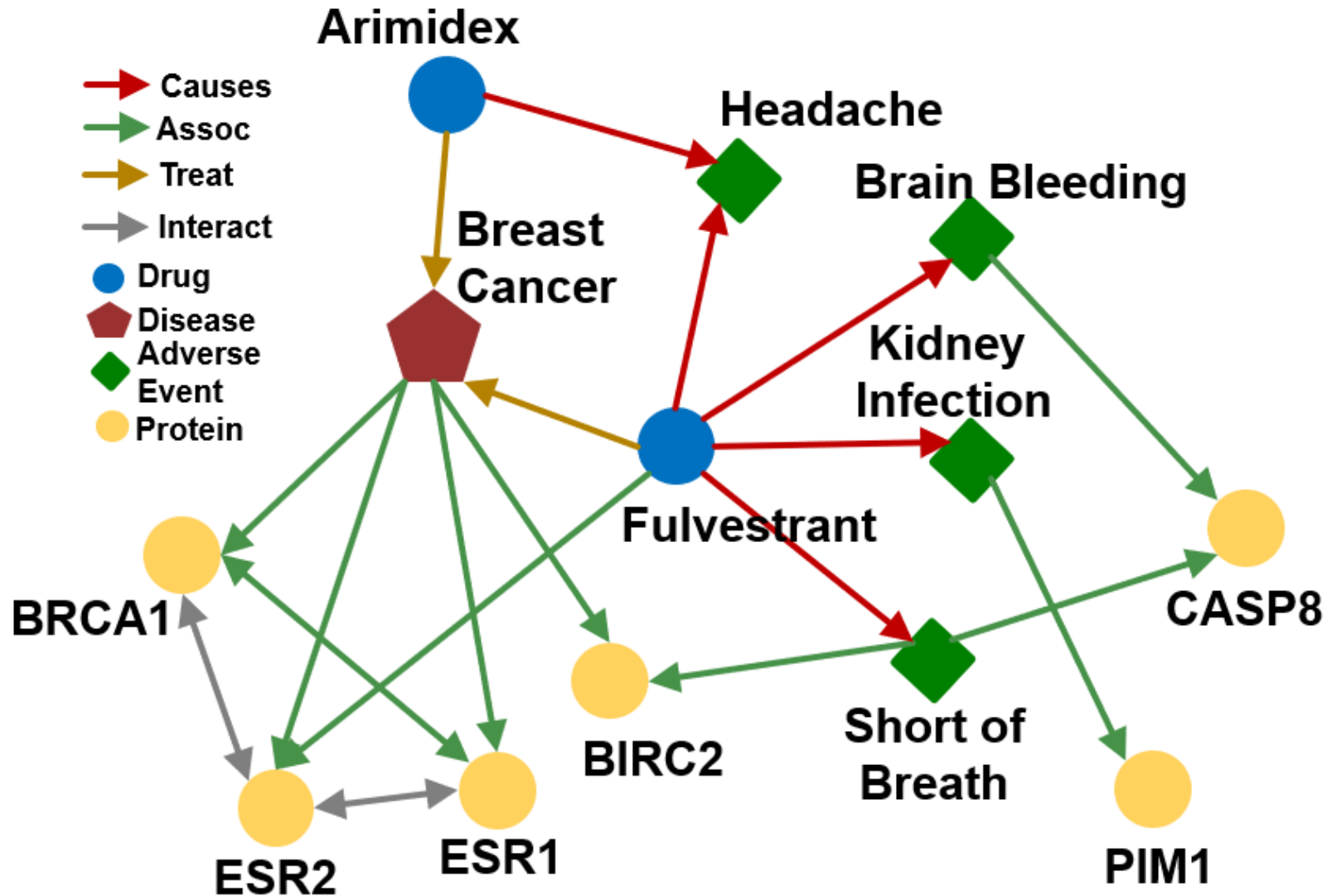
- § Answering multi-hop queries

- § Path Queries

- § Conjunctive Queries

- § Query2Box

Example KG: Biomedicine



Predictive Queries on KG

Can we do multi-hop reasoning, i.e., **answer complex queries on an incomplete, massive KG?**

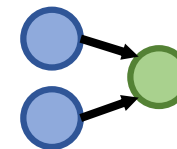
Query Types	Examples: Natural Language Question, Query
One-hop Queries	What adverse event is caused by Fulvestrant? (e:Fulvestrant, (r:Causes))
Path Queries	What protein is associated with the adverse event caused by Fulvestrant? (e:Fulvestrant, (r:Causes, r:Assoc))
Conjunctive Queries	What is the drug that treats breast cancer and caused headache? ((e:BreastCancer, (r:TreatedBy)), (e:Migraine, (r:CausedBy)))



One-hop Queries



Path Queries

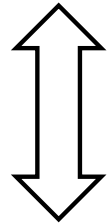


Conjunctive Queries

Predictive One-hop Queries

i We can formulate knowledge graph completion problems as answering one-hop queries.

i **KG completion:** Is link (h , r , t) in the KG?



i **One-hop query:** Is t an answer to query (h , r)?

§ **For example:** What side effects are caused by drug Fulvestrant?

Path Queries

! Generalize one-hop queries to path queries by **adding more relations on the path**.

! An n -hop path query q can be represented by

$$q = (v_!, (r_1, \dots, r_n))$$

§ $v_!$ is an “anchor” entity,

§ Let answers to " q " in graph G be denoted by $[[q]]_G$.

Path Queries

- Generalize one-hop queries to path queries by **adding more relations on the path**.

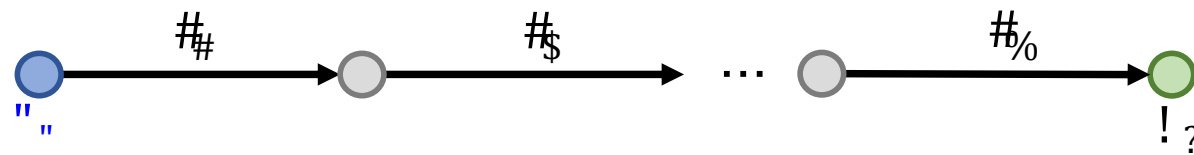
- An n -hop path query q can be represented by

$$q = (v_1, (r_1, \dots, r_n))$$

§ v_1 is an “anchor” entity,

§ Let answers to q in graph G be denoted by $[[q]]_G$.

Query Plan of q :

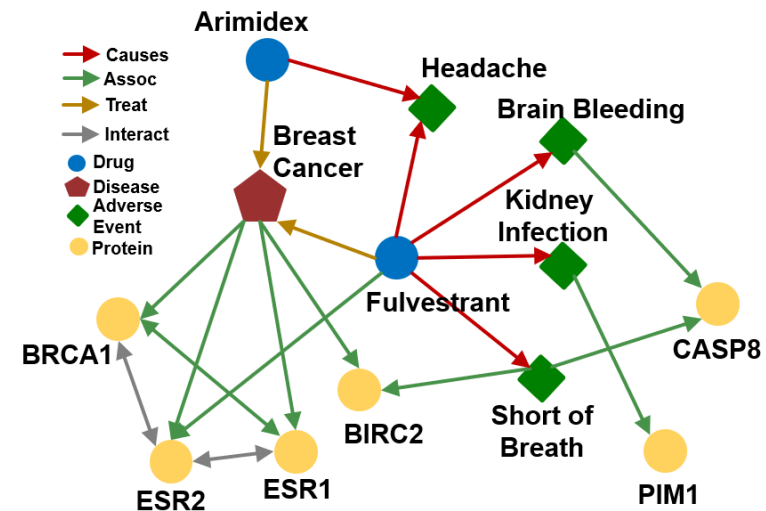
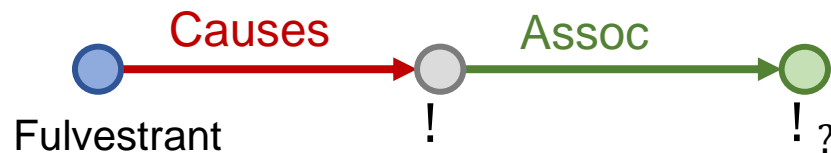


Query plan of path queries is a chain.

Path Queries

Question: “What proteins are *associated* with adverse events *caused* by *Fulvestrant*?”

- ! is e:Fulvestrant
- (\$_1, \$_2) is (r:Causes, r:Assoc)
- Query: (e:Fulvestrant, (r:Causes, r:Assoc))

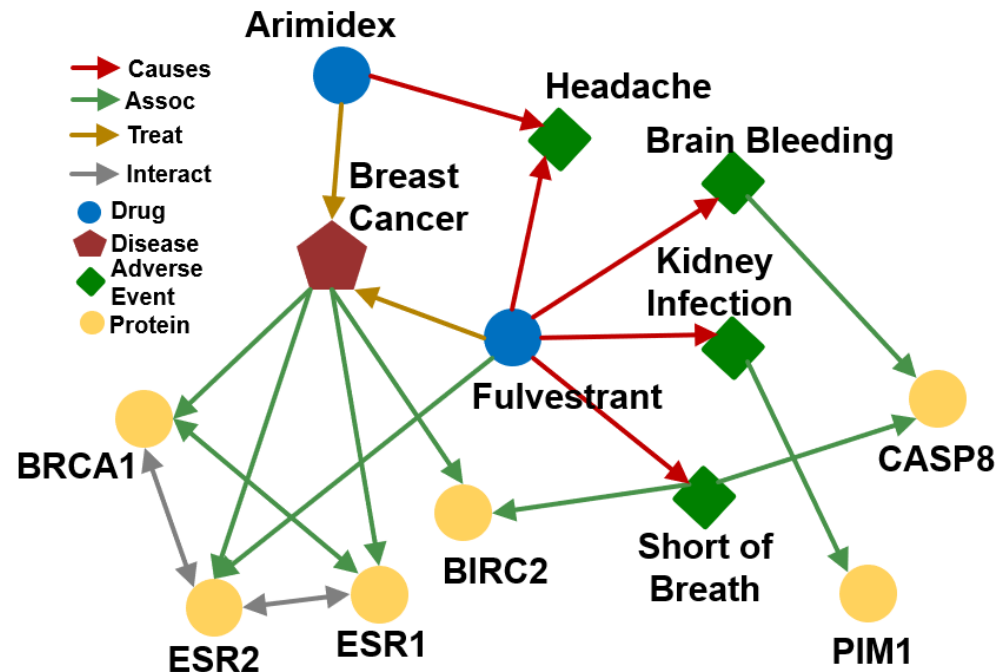


Path Queries

Question: “What proteins are *associated* with adverse events *caused* by *Fulvestrant*?”

Query: (e:Fulvestrant, (r:Causes, r:Assoc))

Given a KG, how to answer a path query?



Traversing Knowledge Graphs

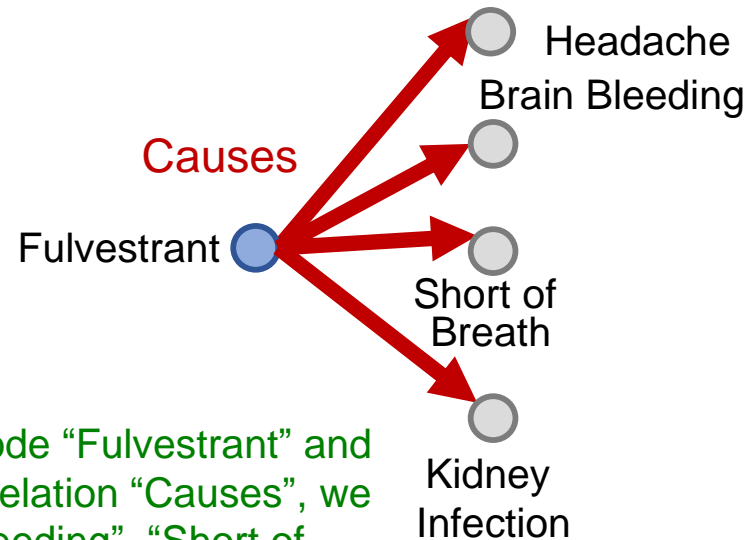
- i We answer path queries by traversing the KG:
“What proteins are *associated* with adverse events *caused* by *Fulvestrant*?”
- i Query: (e:Fulvestrant, (r:Causes, r:Assoc))

Fulvestrant 

Start from the
anchor node
(Fulvestrant).

Traversing Knowledge Graphs

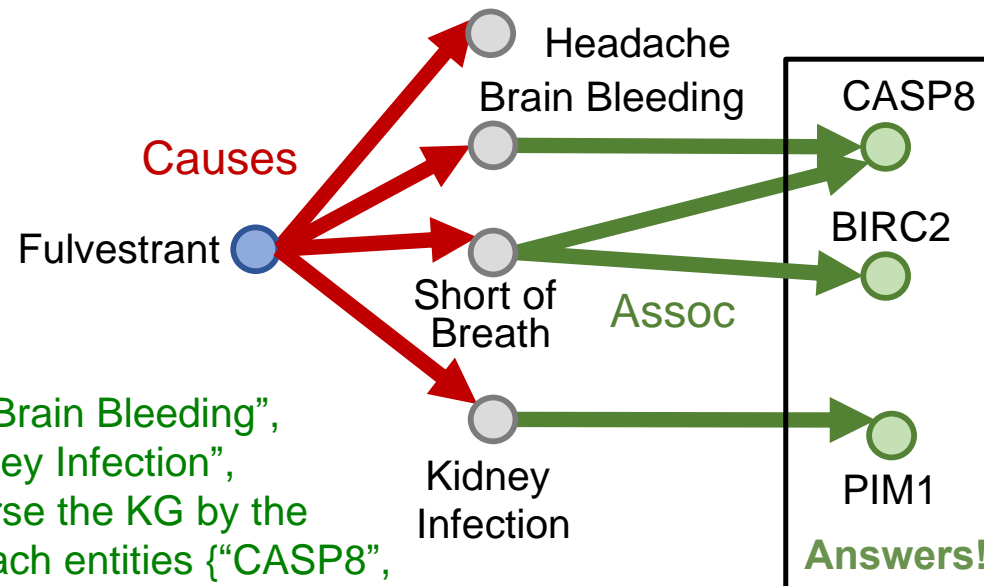
- ! We answer path queries by traversing the KG:
“What proteins are *associated* with adverse events *caused* by *Fulvestrant*?”
- ! Query: (e:Fulvestrant, (r:Causes, r:Assoc))



Start from the anchor node “Fulvestrant” and traverse the KG by the relation “Causes”, we reach entities {“Brain Bleeding”, “Short of Breath”, “Kidney Infection”, “Headache”}.

Traversing Knowledge Graphs

- ! We answer path queries by traversing the KG:
“What proteins are *associated* with adverse events *caused* by *Fulvestrant*?”
- ! Query: (e:Fulvestrant, (r:Causes, r:Assoc))



Start from the nodes {“Brain Bleeding”, “Short of Breath”, “Kidney Infection”, “Headache”} and traverse the KG by the relation “Assoc”, we reach entities {“CASP8”, “BIRC2”, “PIM1”}. These are the answers.

However, KGs are incomplete

- | Answering queries seems easy: Just traverse the graph.
- | **But**

However, KGs are incomplete

- | **Answering queries seems easy: Just traverse the graph.**
- | **But KGs are incomplete and unknown:**
 - § Many relations between entities are missing or are incomplete
 - § For example, we lack all the biomedical knowledge
 - § Enumerating all the facts takes non-trivial time and cost, we cannot hope that KGs will ever be fully complete
- | **Due to KG incompleteness, one is not able to identify all the answer entities**

Can KG Completion Help?

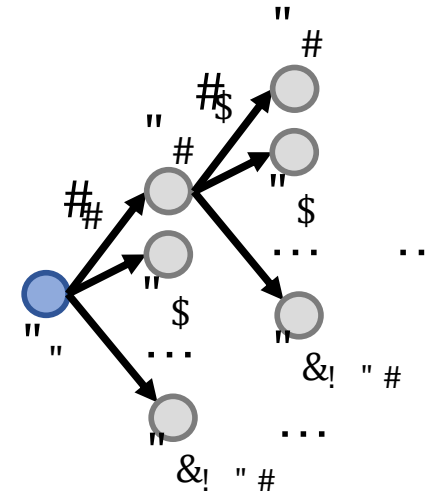
Can we first do KG completion and then traverse the completed (probabilistic) KG?

Can KG Completion Help?

Can we first do KG completion and then traverse the completed (probabilistic) KG?

- ! **No!** The “completed” KG is a **dense graph!**
 - § Most (h, r, t) triples (edge on KG) will have some non-zero probability.

- ! Time complexity of traversing a dense KG is exponential as a function of the path length, $\sim (|V|)^L$



Task: Predictive Queries

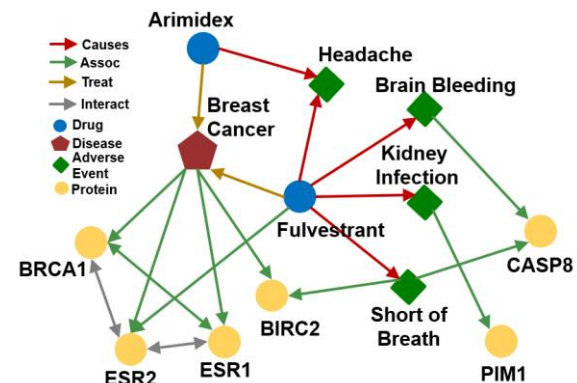
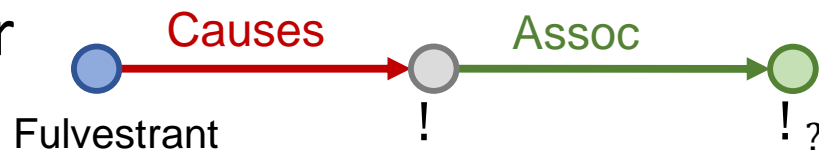
- i We need a way to answer path-based queries over an incomplete knowledge graph.
- i We want our approach to implicitly impute and account for the incomplete KG.

Task: Predictive Queries

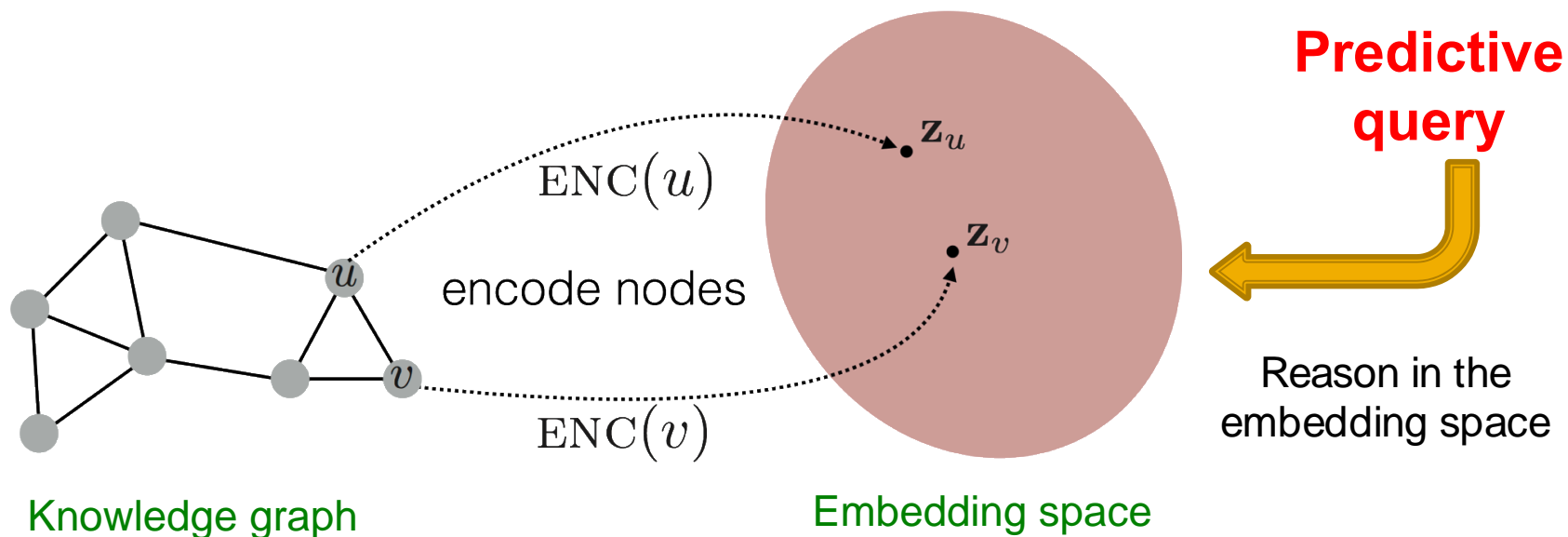
- i We need a way to answer path-based queries over an incomplete knowledge graph.
- i We want our approach to implicitly impute and account for the incomplete KG.
- i **Task: Predictive queries**

§ Want to be able to answer arbitrary queries while implicitly imputing for the missing information

§ **Generalization of the link prediction task**



A General Idea



Map queries into embedding space. **Learn to reason in that space**

- i Embed query into a single **point** in the Euclidean space: answer nodes are close to the query.

[[Embedding Logical Queries on Knowledge Graphs](#). Hamilton, et al., NeurIPS 2018]

[[Query2box: Reasoning over Knowledge Graphs in Vector Space Using Box Embeddings](#). Ren, et al., ICLR 2020]

Traversing KG in Vector Space

- | **Key idea: Embed queries!**

- § Generalize **TransE** to multi-hop reasoning.

Traversing KG in Vector Space

i Key idea: Embed queries!

§ Generalize **TransE** to multi-hop reasoning.

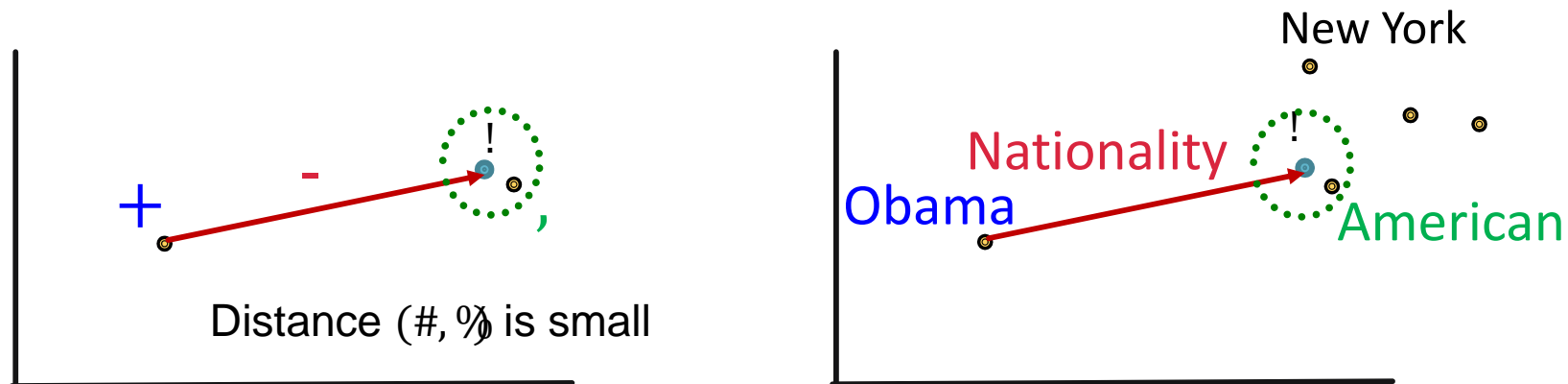
§ **Recap: TransE**: Translate ! to " using # with score function $f(h, r) = -\|h + r - t\|$.

§ Another way to interpret this is that:

§ **Query embedding**: $q = h + r$

§ Goal: **query embedding** q is close to the **answer embedding** t

$$f(t) = -\|q - t\|$$

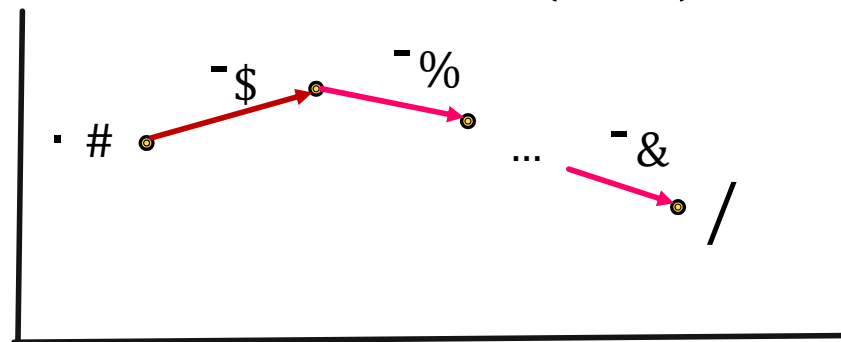


Traversing KG in Vector Space

! Key idea: Embed queries!

§ Generalize TransE to multi-hop reasoning.

Given a path query $q = (0, (1), \dots, 1^*)$,



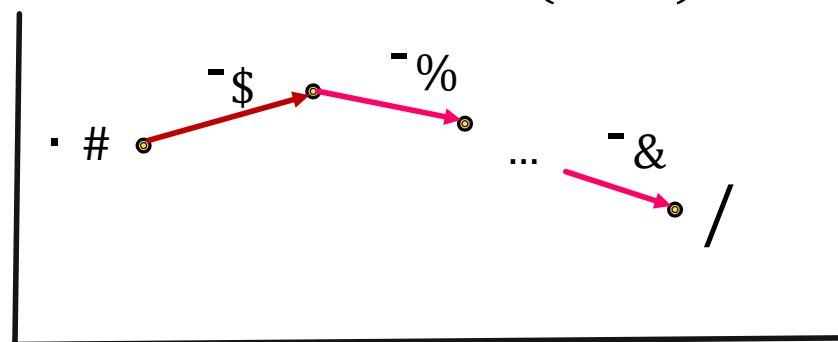
$$2 = 3_1 + 5_1 + \dots + 5_{\#}$$

Traversing KG in Vector Space

- Key idea: Embed queries!

§ Generalize TransE to multi-hop reasoning.

Given a path query $q = (0, (1), \dots, 1^*)$,



$$2 = 3_1 + 5_1 + \dots + 5_{\#}$$

- The embedding process **only involves vector addition**, **independent of # entities** in the KG!

Traversing KG in Vector Space

Embed path queries in vector space.

- | **Question:** “What proteins are *associated* with adverse events *caused* by *Fulvestrant*?”
- | **Query:** (e:Fulvestrant, (r:Causes , r:Assoc))

Follow the query plan:

Query Plan

Embedding Process

Fulvestrant ●

Fulvestrant ●

Traversing KG in Vector Space

Embed path queries in vector space.

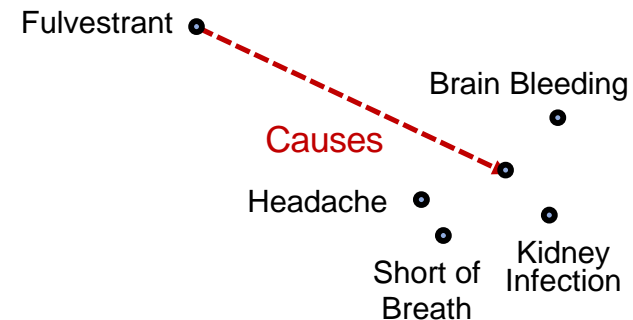
- Question: “What proteins are *associated* with adverse events *caused* by *Fulvestrant*?”
- Query: (e:Fulvestrant, (r:Causes , r:Assoc))

Follow the query plan:

Query Plan



Embedding Process



Traversing KG in Vector Space

Embed path queries in vector space.

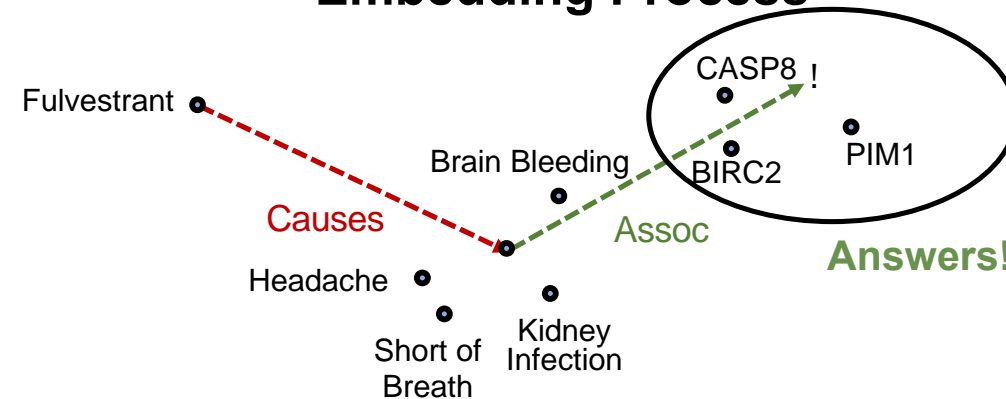
- Question: “What proteins are *associated* with adverse events *caused* by *Fulvestrant*?”
- Query: (e:Fulvestrant, (r:Causes , r:Assoc))

Follow the query plan:

Query Plan



Embedding Process



Traversing KG in Vector Space

Insights:

- i We can train **TransE** to optimize knowledge graph completion objective
- i Since **TransE** can naturally handle **compositional relations**, it can handle path queries by translating in the latent space **for multiple hops using addition of relation embeddings.**

Recommender System (RecSys)

Slides adapted from:

- Y. Sun, CS 247: Advanced Data Mining
- Jure Leskovec, Stanford CS224W: Machine Learning with Graphs

Recommender System Examples

You may also like



Jack & Jones
JAMIE - Polo shirt - orange
£21.00
Free delivery & returns

ALTERNATIVE PRODUCTS

Beko Washing Machine
Code: WMB81431LW
£269.99

Zanussi Washing Machine
Code: ZWH6130P
£269.99

Blomberg Washing Machine
Code: WNF6221
£299.99

You may also like



★★★★☆ (109)



★★★★☆ (53)



★★★★☆ (33)

Related hotels...



Hotel 41
★★★★★ 1,170 Reviews
London, England

Show Prices

Read Commented Recommended



Germany Just Rejected The
Idea That The European Bailout
Fund Would Buy Spanish Debt



There Is Almost No Gold In The
Olympic Gold Medal



MOST POPULAR

RECOMMENDED

How to Break NRA's Grip on Politics: Michael R. Bloomberg +

Growth in U.S. Slows as Consumers Restrain Spending +

Why Recommender Systems?

- Value for the customer
 - Find things that are interesting
 - Narrow down the set of choices
 - Help me explore the space of options
 - Discover new things
 - Entertainment
 - ...
- Value for the provider
 - Additional and probably unique personalized service for the customer
 - Increase trust and customer loyalty
 - Increase sales, click through rates, conversion etc.
 - Opportunities for promotion, persuasion
 - Obtain more knowledge about customers
 - ...

Formulating the RecSys problem (I): Matrix Completion

Users	Movie1	Movie2	Movie3	Movie4	Movie5	Movie6	...
User1	?	?	4	?	1	?	...
User2	2	5	2	?	?	2	...
User3	?	?	5	3	2	4	...
User4	1	?	?	4	?	?	...
User5	2	3	?	?	?	?	...
...

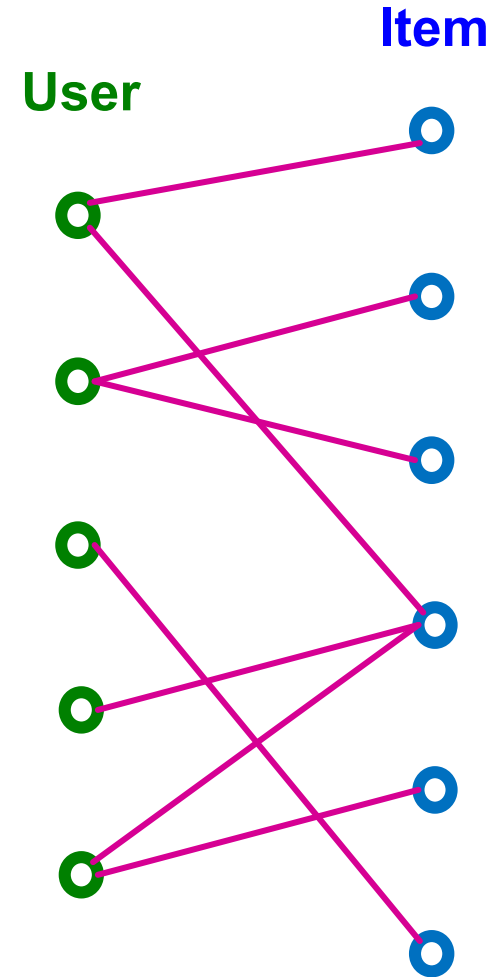
Explicit Feedback vs. Implicit Feedback

- **Explicit Feedback**
 - Know the ratings
- **Implicit Feedback**
 - only know whether user and item has interacted
 - Like (1) vs. unknown (0)

		Items				...
						
Users	Alice	1	1	0	0	
	Bob	0	0	1	1	
	Corey	1	0	1	0	
	...					

Formulating the RecSys problem (II): Link Prediction

- i Recommender system can be naturally modeled as a **bipartite graph**
 - § A graph with two node types: **users** and **items**.
 - § **Edges** connect users and items
 - § Indicates user-item interaction (e.g., click, purchase, review etc.)
 - § Often associated with timestamp (timing of the interaction).



Formulating the RecSys problem (II): Link Prediction

i Given

§ Past user-item interactions

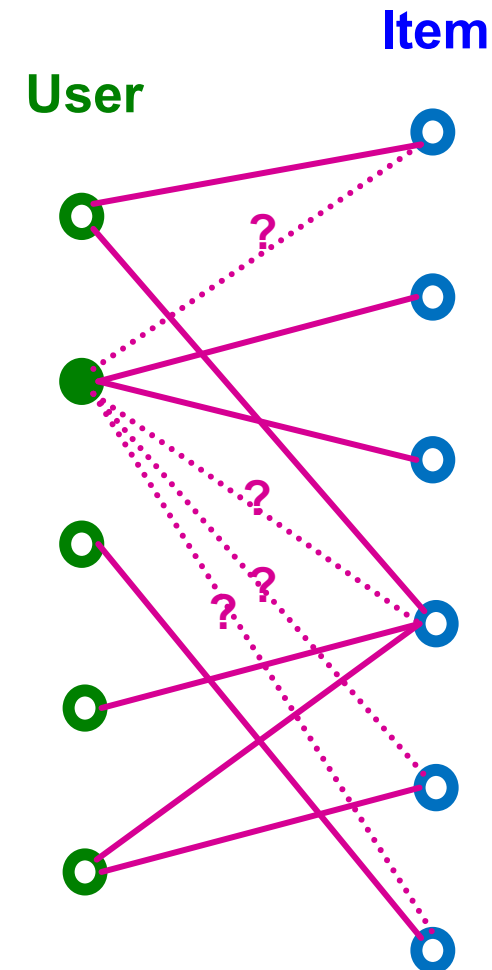
i Task

§ Predict new items each user will interact in the future.

§ Can be cast as **link prediction** problem.

§ Predict new user-item interaction edges given the past edges.

§ For $! \in \#$, $\% \in \&$, we need to get a real-valued **score** $' (!, \%)$.



Modern Recommender System

- ! **Problem:** Cannot evaluate $r(u, i)$ for every user u – item i pair.

Example $r(u, i)$:
 $r(u, i) = \sum_j u_j \cdot i_j$

Modern Recommender System

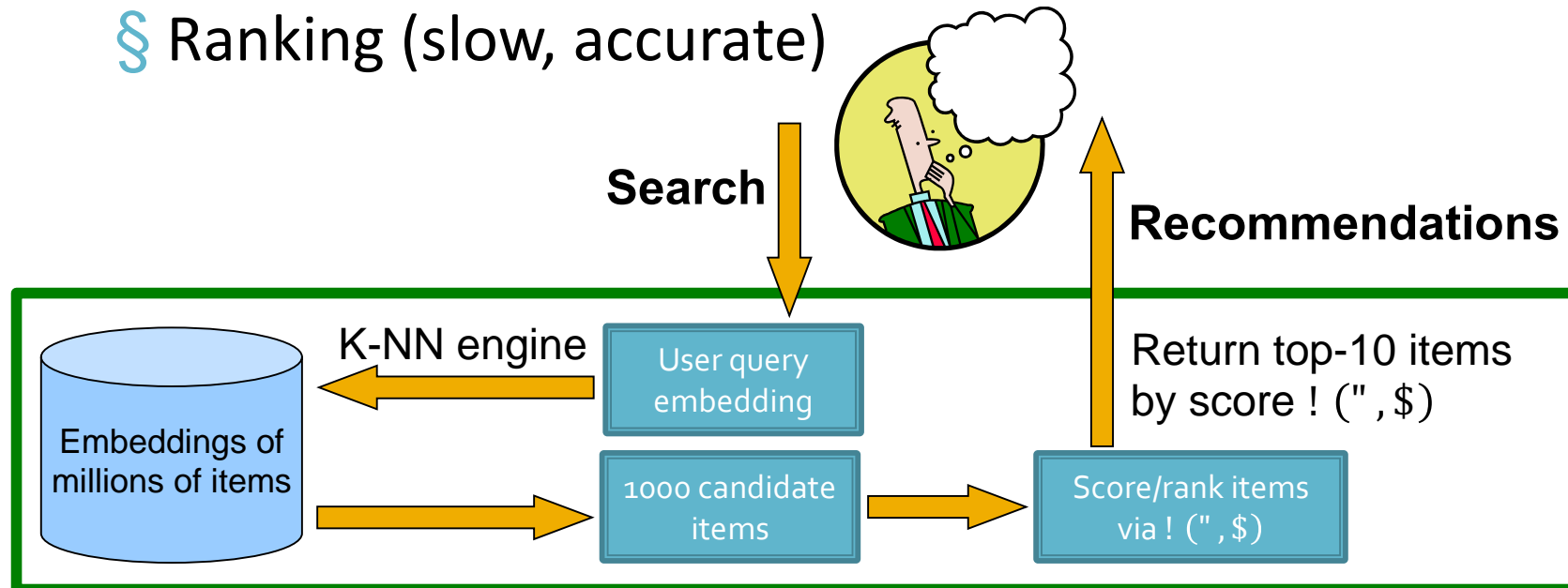
❗ **Problem:** Cannot evaluate $! (" , \$)$ for every user " – item \$ pair.

Example $! (" , \$)$:
 $! (" , \$) = \& \cdot \&$

❗ **Solution:** 2-stage process:

§ Candidate generation (cheap, fast)

§ Ranking (slow, accurate)



Top-K Recommendation

- i For each user, we recommend % items.
 - § For recommendation to be effective, * needs to be much smaller than the total number of items (up to billions)
 - § + is typically in the order of 10—100.

Top-K Recommendation

- i For each user, we recommend % items.
 - § For recommendation to be effective, * needs to be much smaller than the total number of items (up to billions)
 - § + is typically in the order of 10—100.
- i The goal is to include as many **positive items** as possible in the top-% recommended items.
 - § **Positive items = Items that the user will interact with in the future.**
- i **Evaluation metric:** Recall@% (defined next)

Evaluation Metric: Recall@K

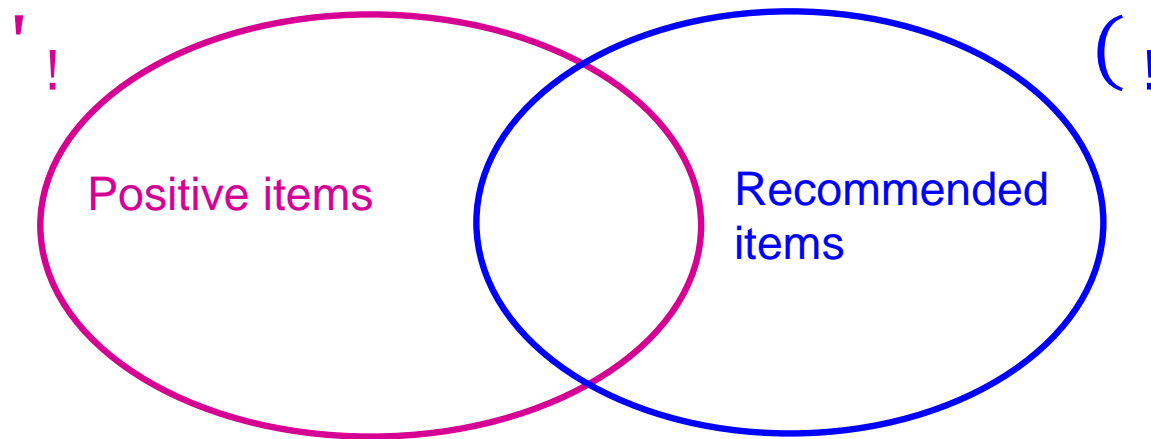
! For each user u ,

§ Let P_u be a set of positive items the user will interact in the future.

§ Let R_u be a set of items recommended by the model.

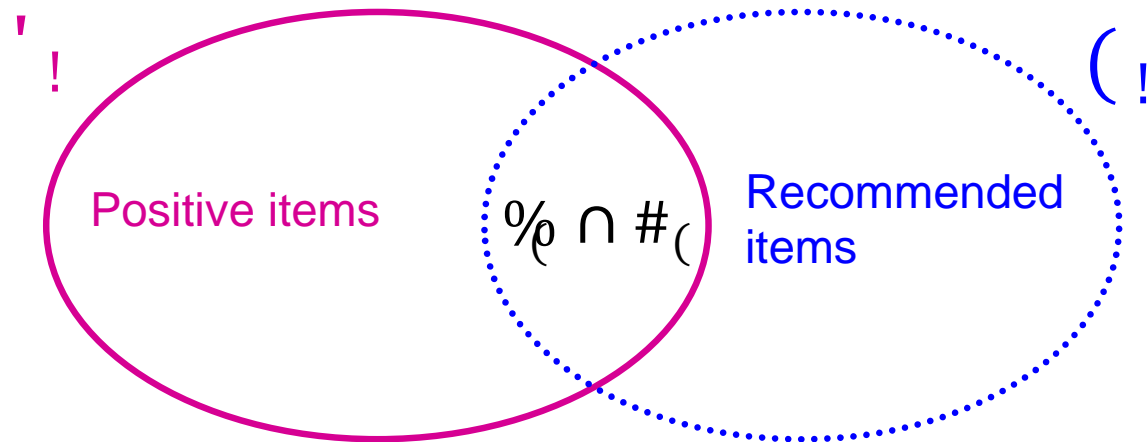
§ In top- K recommendation, $|R_u| = K$.

§ Items that the user has already interacted are excluded.



Evaluation Metric: Recall@K

- ! **Recall@K** for user u is $\frac{|P_u \cap R_u|}{|R_u|}$.
- § Higher value indicates more positive items are recommended in top- K for user u .



- ! The final Recall@K is computed by averaging the recall values across all users.

Methods

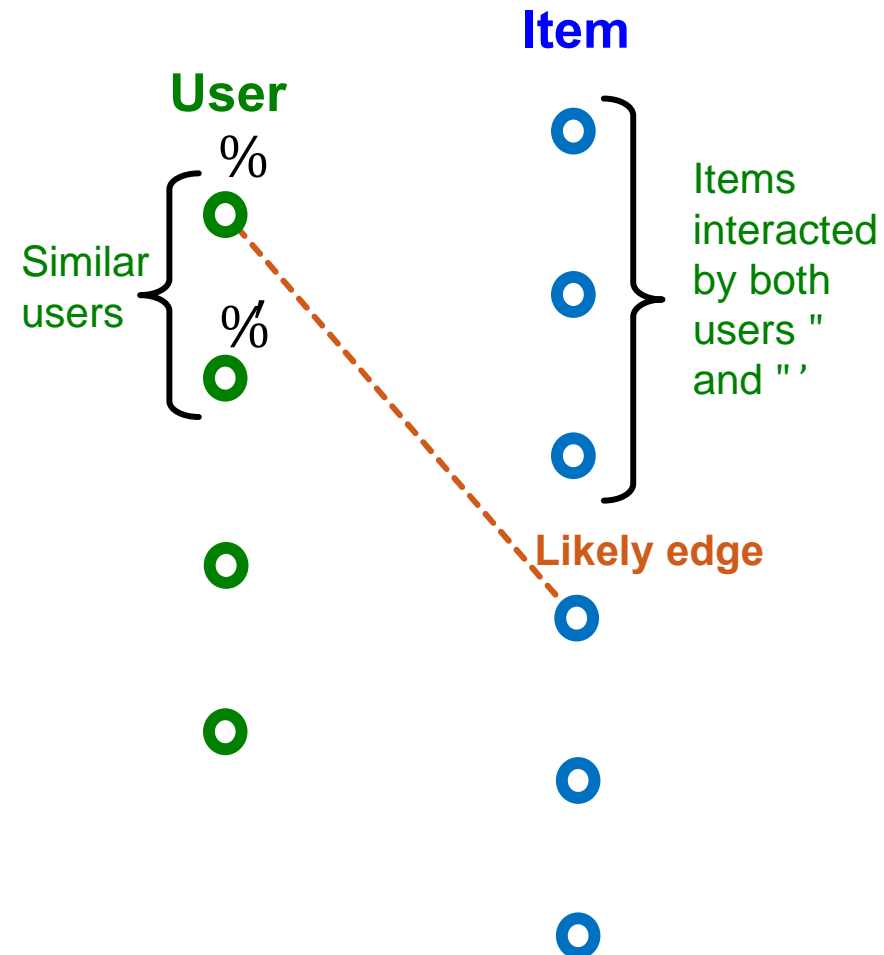
- Collaborative filtering
- Content-based recommendation
- Hybrid methods

Methods

- Collaborative filtering
- Content-based recommendation
- Hybrid methods

Collaborative Filtering (CF)

- i **Underlying idea:**
 - § **Collaborative filtering**
 - § Recommend items for a user by **collecting preferences of many other similar users.**
 - § **Similar users tend to prefer similar items.**
- i **Key question: How to capture similarity between users/items?**



Questions?