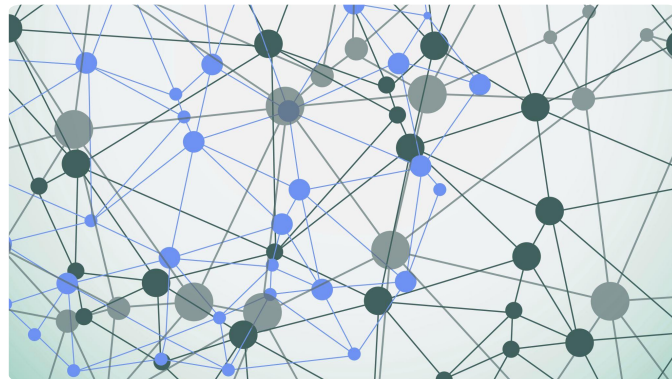


Revisiting Evaluation of Knowledge Base Completion Models

Pouya Pezeshkpour, Yifan Tian, Sameer
Singh



Revisiting Evaluation of KB Completion



Evaluating KG Completion



Shortcomings

1. Semi-Inverse relations
2. Calibration
3. Triple classification robustness



Introducing YAGO3-TC

Revisiting Evaluation of KB Completion



Evaluating KG Completion



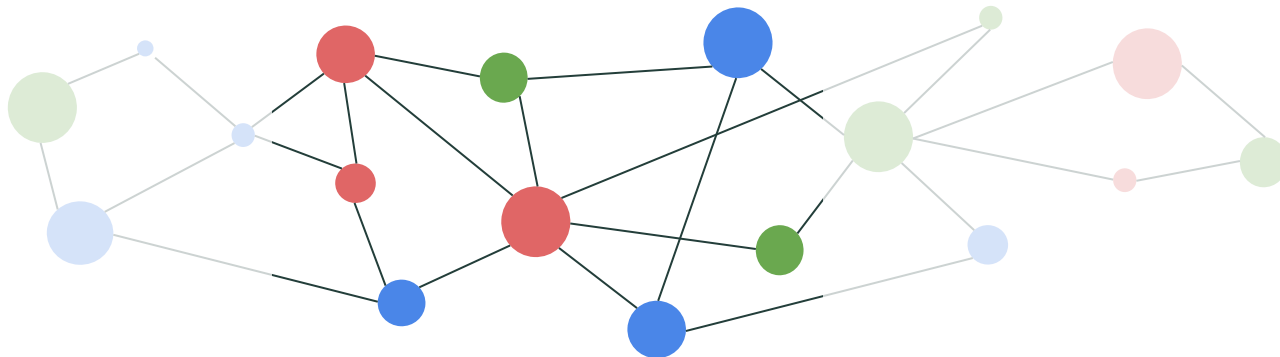
Shortcomings

1. Semi-Inverse relations
2. Calibration
3. Triple classification robustness



Introducing YAGO3-TC

Knowledge Graphs



Unstructured text

Barack Hussein Obama II (/bəˈrɑːk huːˈseɪn oʊˈbɑːmə/ ^[1] listen),^[1] born August 4, 1961) is an American politician and attorney who served as the 44th [president of the United States](#) from 2009 to 2017. A member of the [Democratic Party](#), Barack Obama was the first [African-American](#) president of the United States. He previously served as a [U.S. senator](#) from [Illinois](#) from 2005 to 2008 and an [Illinois state senator](#) from 1997 to 2004.

Structured text

44th President of the United States
<div>In office</div> <div>January 20, 2009 – January 20, 2017</div>
<div>Vice President Joe Biden</div>
<div>Preceded by George W. Bush</div>
<div>Succeeded by Donald Trump</div>
United States Senator from Illinois
<div>In office</div> <div>January 3, 2005 – November 16, 2008</div>
<div>Preceded by Peter Fitzgerald</div>
<div>Succeeded by Roland Burris</div>

Images



Shortcomings of Knowledge Graphs

Because of the way they are created:

- There are many missing facts

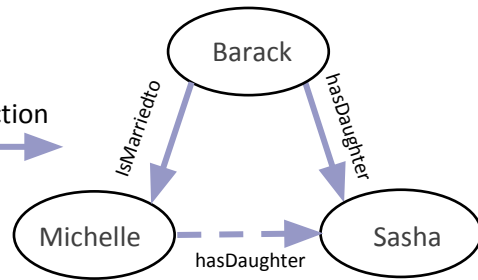
➡ KG Completion

<Barack, IsMarriedto, Michelle>

<Barack, hasDaughter, Sasha>

< Michelle, hasDaughter, ?>

Link Prediction



- The factuality of non-existent links is unknown (open-world assumption)

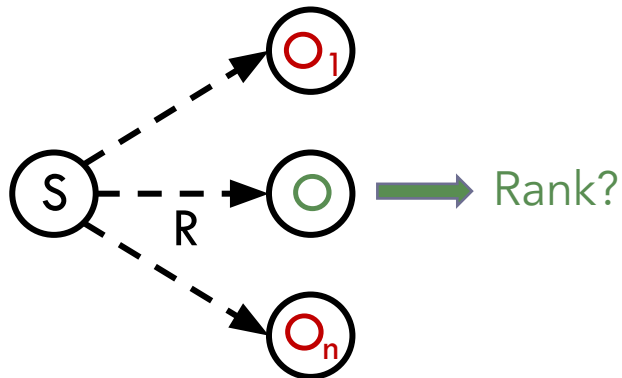
➡ Use of ranking for evaluation

But, real-world application mostly care about an information being True or False and not the ranking

Knowledge Graph Completion Evaluation

- Two evaluation approach for scoring target triple $\langle S, R, O \rangle$:

1. Ranking Metrics:



Doesn't correspond
to real-world use
case in many
instances

2. Triple Classification:



Learning thresholds τ_R by randomly choosing
negative samples for validation data

Revisiting Evaluation of KB Completion



Overview of Knowledge
Graph Completion



Shortcomings

1. Semi-Inverse relations
2. Calibration
3. Triple classification robustness

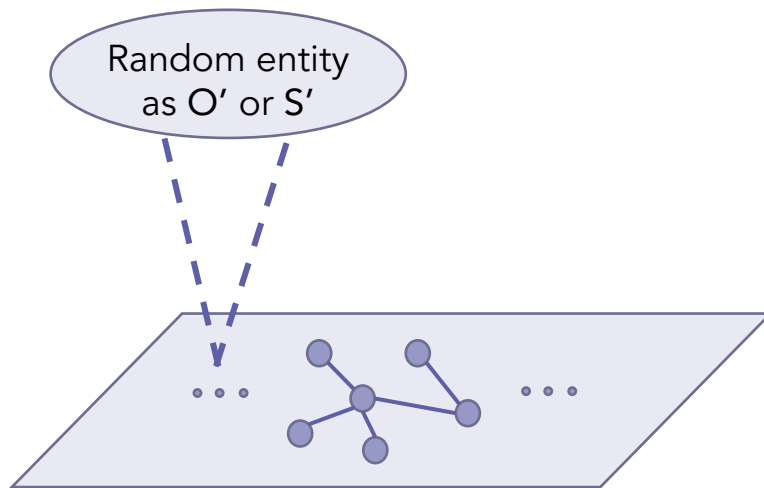


Introducing YAGO3-TC

Negative Sampling

- We consider 3 different negative sampling for target triple $\langle S, R, O \rangle$:

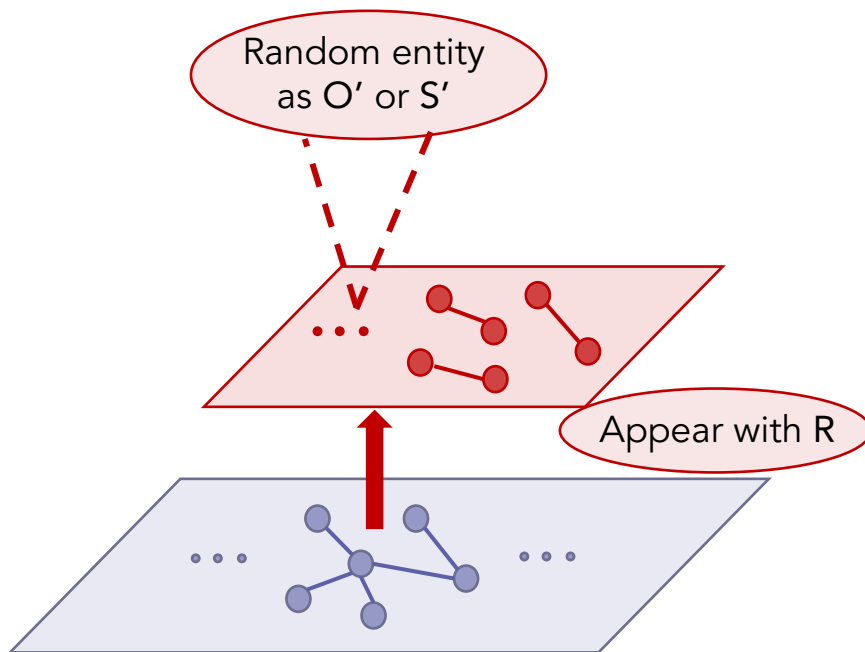
1. Random Sampling:



Negative Sampling

- We consider 3 different negative sampling for target triple $\langle S, R, O \rangle$:

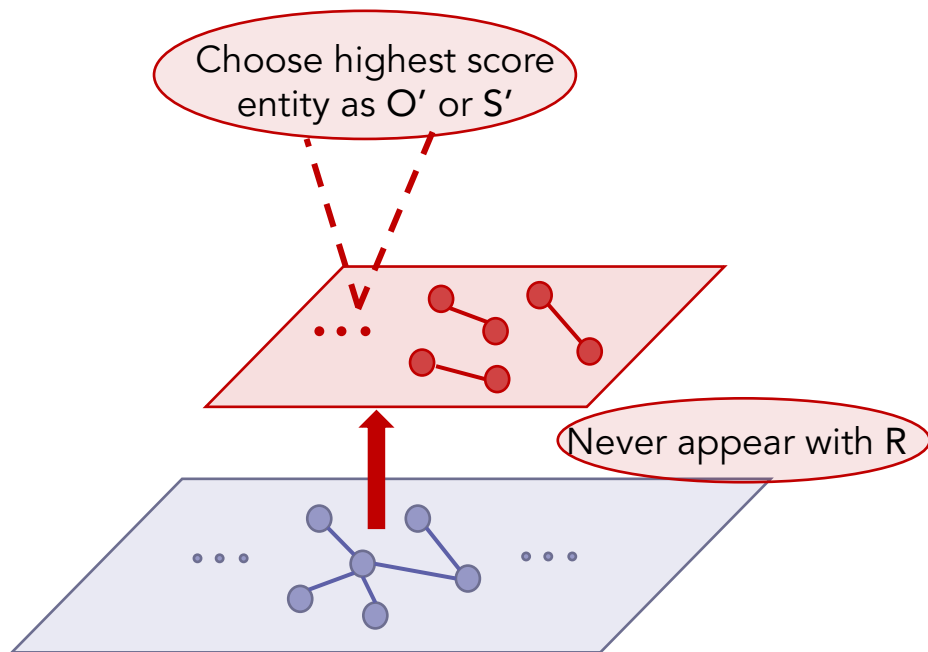
2. Constraint Sampling:



Negative Sampling

- We consider 3 different negative sampling for target triple $\langle S, R, O \rangle$:

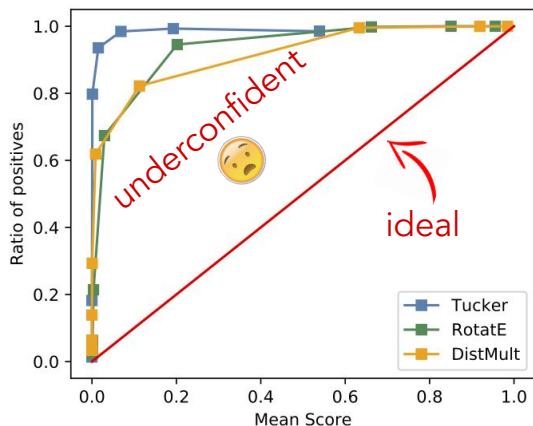
3. Careful Sampling:



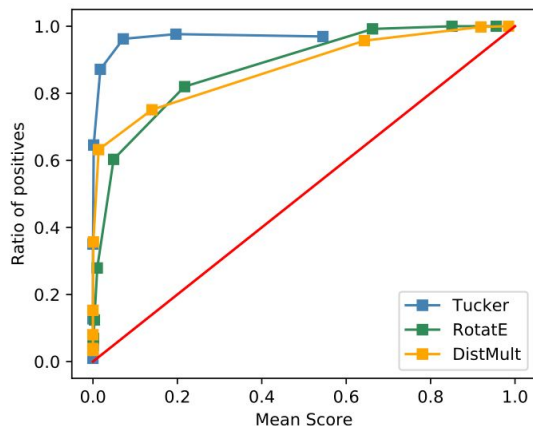
Calibration

- Calibration study is not well defined

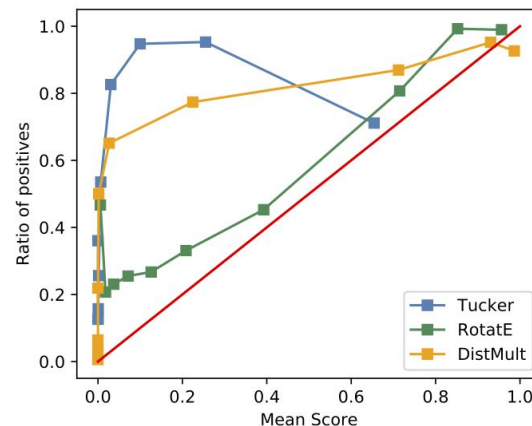
1. Random sampling



2. Constrained sampling



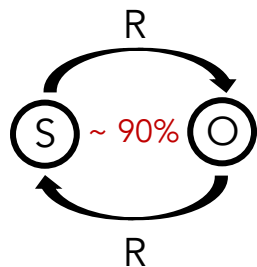
3. Careful sampling



Extremely different conclusions
for different negative samplings

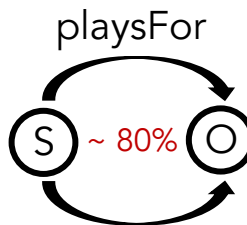
Semi-Inverse Relations

WN18RR

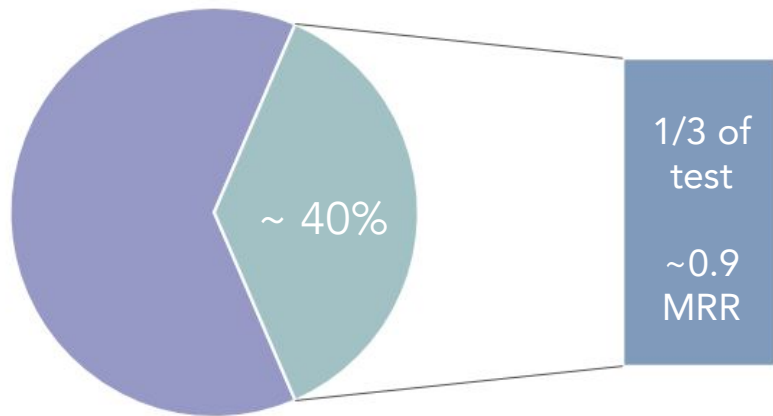


Ranking metrics do not
reflect reasoning power

YAGO3-10

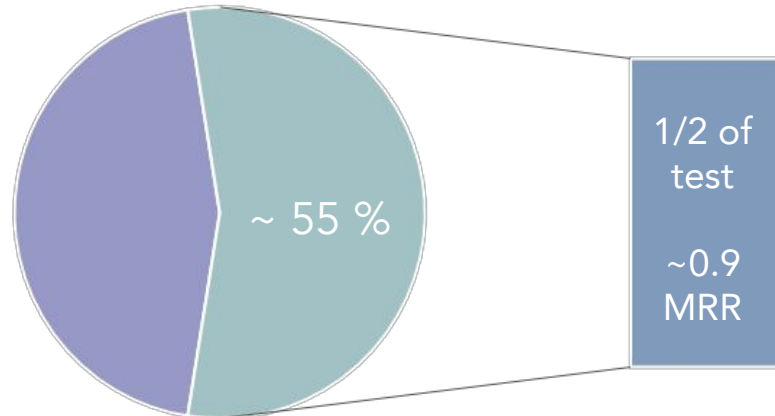


Train data



■ Other R ■ Semi-Inverse

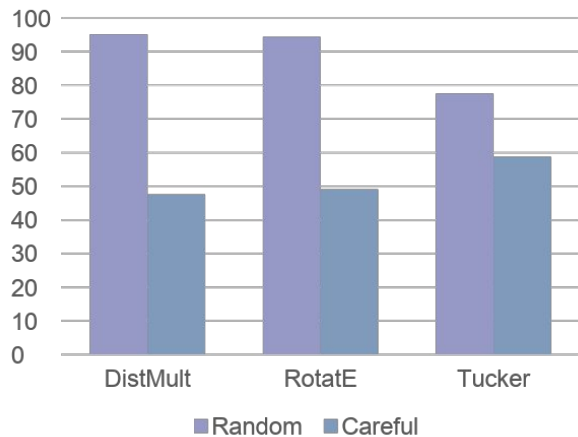
Train data



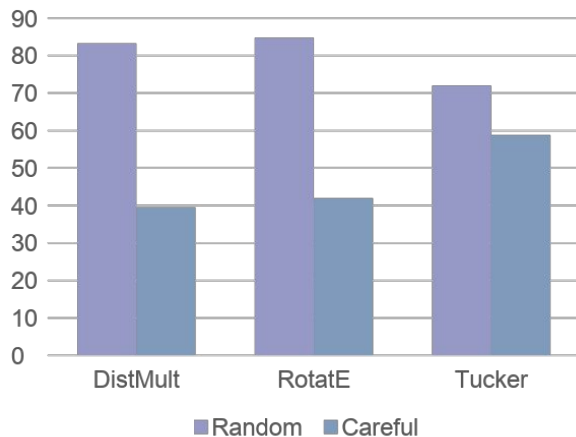
■ Other R ■ Semi-Inverse

Triple Classification Robustness

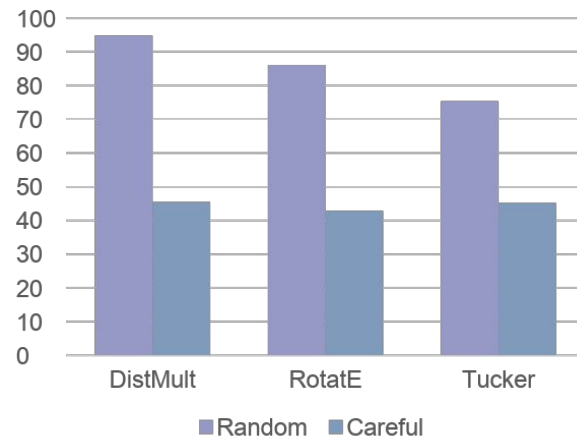
FB15K-237



WN18RR



YAGO3-10



Careful negative sampling results in a dramatic drop

Results are around 90 %

Revisiting Evaluation of KB Completion



Overview of Knowledge
Graph Completion



Shortcomings
of evaluation
metrics

1. Semi-Inverse relations
2. Calibration
3. Triple classification robustness



Introducing YAGO3-TC



YAGO3-TC Dataset

What is Our goal?

- Create a benchmark that align with real-world application
- Properly differentiate between models
- Capture reasoning powers

What are existing challenges?

- The knowledge graphs are not complete
- There so many non existent links
- Identifying the factuality of missing information is hard

YAGO3-TC Creation

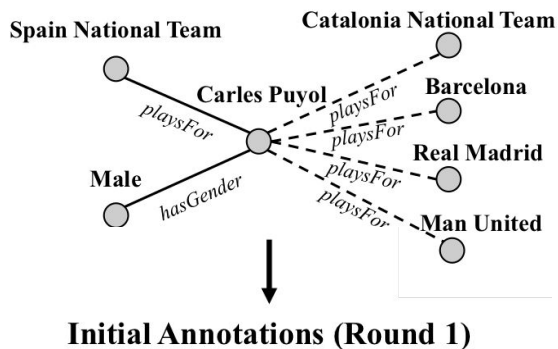
1. Randomly choose a subset of YAGO3-10 test

2. Identify top scoring triples from the models that are unknown to be true

3. Filter triples

4. Crowdsourcing pipeline

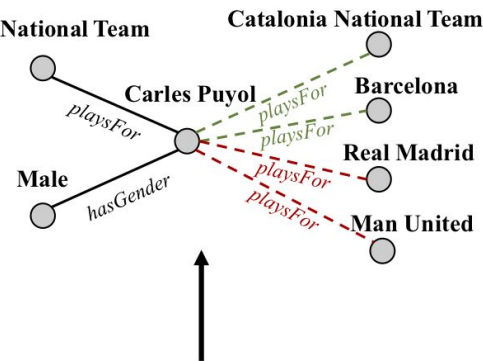
Crowdsourcing Pipeline



Carles Puyol plays for?

<input type="checkbox"/> Catalonia National Team	✓	✓	×
<input type="checkbox"/> Barcelona	✓	✓	✓
<input type="checkbox"/> Real Madrid	✓	×	×
<input type="checkbox"/> Man United	×	×	×

✓ Catalonia National Team
✓ Barcelona
× Man United
? Real Madrid



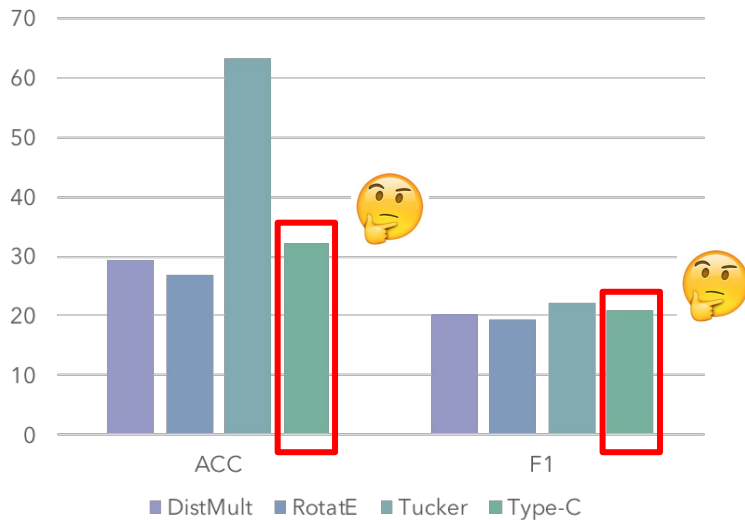
Carles Puyol plays for?

<input type="checkbox"/> Real Madrid	×	×
--------------------------------------	---	---

- ~ 30 K triples
- ~ 10% positives

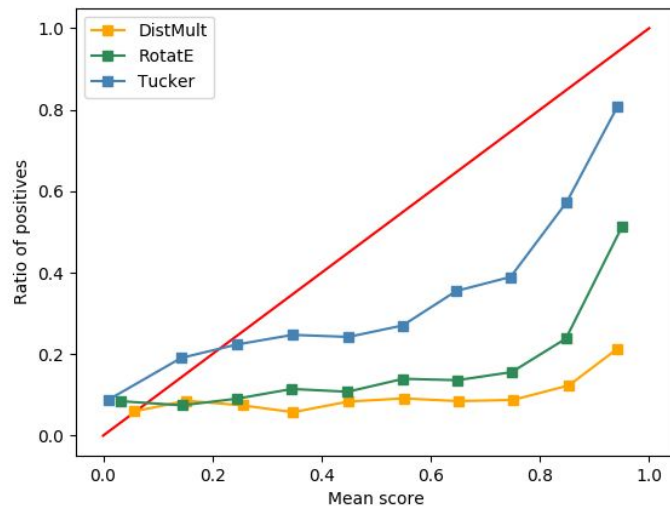
Evaluation Using YAGO3-TC

- Triple classification



- SOTA models perform poorly
- Huge drop in accuracy

- Calibration



- Reverse order of models
- Overconfident

Discussion

- Ranking metrics are not very trustworthy
- Triple classification is not robust
- Real-world adoption of KG needs better evaluation techniques
- YAGO3-TC is the first step toward this goal

We propose a web-hosted evaluation platform to update YAGO3-TC using new KG completion models

Thank You!

Website (code, data and leaderboard):

pouyapez.github.io/yago3-tc/

Contact me: **pezeshkp@uci.edu**