

Ranking vs. Classifying: Measuring Knowledge Base Completion Quality

Marina Speranskaya, Martin Schmitt, Benjamin Roth

CIS, LMU Munich

22.6.2020 – 24.6.2020



Motivation

Rank = 3,
not bad!

Thomas Lennon – has profession – ?

top-3

e4 musician 0.9

e7 artist 0.8

e3 film producer 0.77

e2 screenwriter 0.5

...

e1 New Zealand 0.04

e6 US dollar 0.01

With what triples
do I actually
update my KB?

Top-3, top-5?

How does it change
the **KB quality**?

TP = {e3}
FP = {e4, e7}
FN = {e2}

➔ Precision & recall

Motivation

Thomas Lennon – has profession

What other possible scenarios are there?

Reflects the closed-world problematic

top-3

e4	musician	0.9
e7	artist	0.8
e3	film producer	0.77
e2	screenwriter	0.5

●	0.8
●	0.72
●	0.4
●	0.3

○	0.8
○	0.5
○	0.3
○	0.2

...

e1	New Zealand	0.04
e6	US dollar	0.01

○	0.04
○	0.01

Paris – has profession - ?

→ Precision & recall

Motivation

Thomas Lennon – has profession

How to account for different number of answers?

top-3

e4	musician	0.9
e7	artist	0.8
e3	film producer	0.77
e2	screenwriter	0.5

> 0.7

●	0.8
●	0.72
●	0.4
●	0.3

○	0.8
○	0.5
○	0.3
○	0.2

...

e1	New Zealand	0.04
e6	US dollar	0.01

○	0.04
○	0.01

Despite good ranking the classification is not optimal yet..

➡ Precision & recall
➡ Score-based threshold

Motivation

Thomas Lennon – has profession

How to account for different number of answers?

e4	musician	0.9
e7	artist	0.8
e3	film producer	0.77
e2	screenwriter	0.5

> 0.7

●	0.8	0.81
●	0.72	0.75
●	0.4	0.73
●	0.3	0.71

○	0.8	0.5
○	0.5	0.5
○	0.3	0.3
○	0.2	0.3

...

e1	New Zealand	0.04
e6	US dollar	0.01

○	0.04	0.03
○	0.01	0.02

Despite good ranking the classification is not optimal yet..

➔ Precision & recall

➔ Score-based threshold

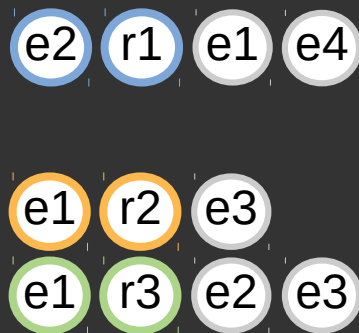
➔ Score calibration

FB14k-QAQ: Query Answering Quality

FB15k-237



FB14k-QAQ



➡ Precision & recall

➡ Query-based

➡ Score-based threshold

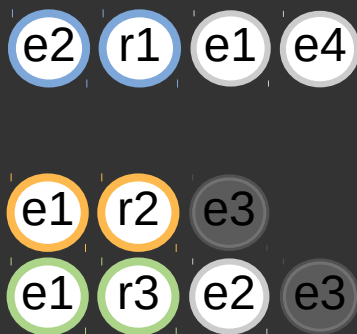
➡ Score calibration

FB14k-QAQ: Query Answering Quality

FB15k-237



FB14k-QAQ



new evaluation set

e2, r1, {e1, e4}

e1, r2, {}

e1, r3, {e2}

+ nonsensical

Make inaccessible
by removing

➡ Precision & recall

➡ Query-based

➡ Score-based threshold

➡ Score calibration

Results

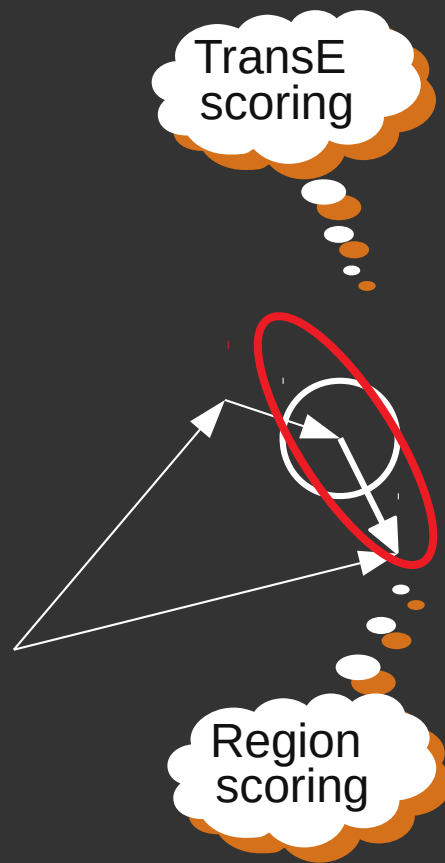
Model (d)	MRR		F1
ConvE 128	.321	1 3	.134
ConvE 64	.263	7 2	.157
Complex 128	.293	2 7	.021
Complex 64	.293	3 8	.009
TransE 128	.293	4 6	.108
TransE 64	.283	5 5	.111
DistMult 64	.266	6 1	.159
DistMult 128	.221	8 4	.133

Complete > nonsensical > inaccessible

Relative order of the models changes

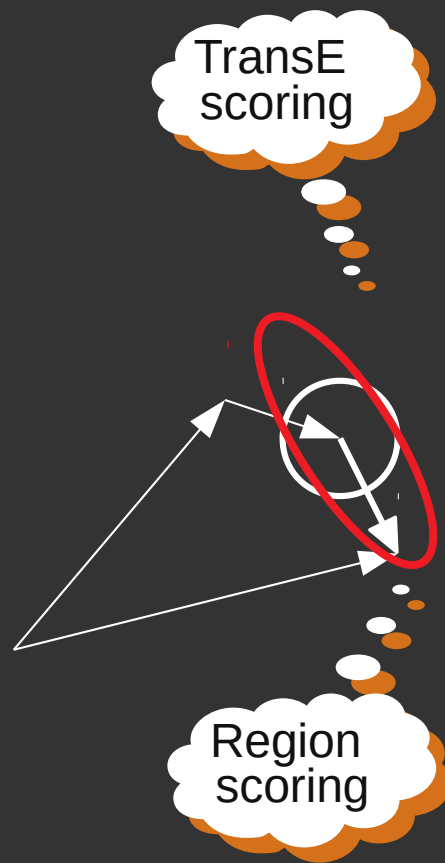
Results

Model (d)	MRR	F1
ConvE 128	.321	.134
ConvE 64	.263	.157
ComplEx 128	.293	.021
ComplEx 64	.293	.009
TransE 128	.293	.108
TransE 64	.283	.111
DistMult 64	.266	.159
DistMult 128	.221	.133



Results

Model (d)	MRR	F1
ConvE 128	.321	.134
ConvE 64	.263	.157
Complex 128	.293	.021
Complex 64	.293	.009
TransE 128	.293	.108
TransE 64	.283	.111
DistMult 64	.266	.159
DistMult 128	.221	.133
Region 64	.330	.146



Contributions

1. A classification-based evaluation setting that provides **interpretable performance metrics**
2. A new benchmark FB14k-QAQ reflecting a **comprehensive set of KBC scenarios**
3. A simple, yet effective strategy to improve TransE's **ability to calibrate scores**

Available at: https://github.com/marina-sp/classification_lp

Thank you for your attention!

This work was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) - RO 5127/2-1, and by the BMBF as part of the project MLWin (01IS18050).

Results

C – complete
F – nonsensical
I – incomplete

Different complexity
of the sets

Model (<i>d</i>)	MRR	F ₁ global threshold				F ₁ multiple thresholds			
		full	C	C ∪ F	T	full	C	C ∪ F	T
ConvE 128	.321	1 3 .134	.272	.211	.105	.204	.317	.286	.150
ConvE 64	.263	7 2 .157	.307	.261	.108	.189	.312	.280	.135
ComplEx 128	.293	2 7 .021	.169	.017	.042	.190	.296	.261	.143
ComplEx 64	.293	3 8 .009	.157	.005	.057	.181	.282	.245	.143
TransE 128	.293	4 6 .108	.150	.106	.108	.159	.172	.168	.154
TransE 64	.283	5 5 .111	.164	.112	.110	.161	.176	.175	.154
DistMult 64	.266	6 1 .159	.273	.226	.129	.184	.256	.239	.148
DistMult 128	.221	8 4 .133	.275	.194	.138	.206	.194	.138	

Relative order of
the models changes