Outline

- Part I. Introduction
- Part II. ML for DI
 - ML for entity linkage
 - ML for data extraction
 - ML for data fusion
 - ML for schema alignment
- Part III. DI for ML
- Part IV. Conclusions and research direction



What is Entity Linkage?

• Definition: Partition a given set R of records, such that each partition corresponds to a distinct real-world entity.

SEE RANK

Are they the same entity?





Actress | Music Department | Soundtrack

Anahi was born in Mexico. She's had roles in Tu y Yo, in which she played a 17 year old girl while she was 13, and Vivo Por Elena, in which she played Talita, a naive and innocent teenager. Anahi lives with her mother and sister name Marychelo. She hopes to become a fashion designer one day, and is currently pursuing a career in singing. See full bio »

Born: May 14, 1982 in Mexico City, Distrito Federal, Mexico

More at IMDbPro » **Contact Info:** View manager



Three Steps in Entity Linkage

• **Blocking**: efficiently create small blocks of



Three Steps in Entity Linkage

• Pairwise matching: compare all record



Three Steps in Entity Linkage

• **Clustering**: group records into entities



50 Years of Entity Linkage

Rule-based and stats-based

 Blocking: e.g., same name Matching: e.g., avg similarity of attribute values Clustering: e.g., transitive closure, etc 		 Supervised learning Random forest for matching F-msr: >95% w. ~1M labels Active learning for blocking & matching F-msr: 80%-98% w. ~1000 labels 	
•	~2000 (Early ML)		2018 (Deep ML)
1969 (Pre-ML)	 ~2015 (ML) Sup / Unsup learning Matching: Decision tree, SVM F-msr: 70%-90% w. 500 labels Clustering: Correlation clustering, Markov clustering 		 Deep learning Deep learning Entity embedding

Rule-Based Solution

Rule-based and stats-based

- Blocking: e.g., same name
- Matching: e.g., avg similarity of attribute values
- Clustering: e.g., transitive closure, etc.

1969 (Pre-ML)

- [Fellegi and Sunter, 1969]
 - Match: $sim(r, r') > \boldsymbol{\theta}_h$
 - Unmatch: $sim(r, r') < \boldsymbol{\theta}_1$
 - Possible match:

$$\boldsymbol{\theta}_{l} < sim(r, r') < \boldsymbol{\theta}_{h}$$

Early ML Models

• [Köpcke et al, VLDB'10]

~2000 (Early ML)

Sup / Unsup learning

- Matching: Decision tree, SVM
 F-msr: 70%-90% w. 500 labels
- Clustering: Correlation clustering, Markov clustering



Supervised learning

- Random forest for matching F-msr: >95% w. ~1M labels
- AL for blocking & matching F-msr: 80%-98% w. ~1000 labels

```
~2015 (ML)
```

- Features: attribute similarity measured in various ways. E.g.,
 - string sim: Jaccard, Levenshtein
 - number sim: absolute diff, relative diff
- ML models on Freebase vs. IMDb
 - Logistic regression: Prec=0.99, Rec=0.6
 - Random forest: Prec=0.99, Rec=0.99

Supervised learning

- Random forest for matching F-msr: >95% w. ~1M labels
- AL for blocking & matching F-msr: 80%-98% w. ~1000

labels

~2015 (ML)

- Expt 1. IMDb vs. Freebase
 - Logistic regression: Prec=0.99, Rec=0.6
 - Random forest: Prec=0.99, Rec=0.99



Supervised learning

- Random forest for matching F-msr: >95% w. ~1M labels
- AL for blocking & matching F-msr: 80%-98% w. ~1000 labels

~2015 (ML)

- Features: attribute similarity measured in various ways. E.g.,
 - name sim: Jaccard, Levenshtein
 - age sim: absolute diff, relative diff
- ML models on Freebase vs. IMDb
 - Logistic regression: Prec=0.99, Rec=0.6
 - Random forest: Prec=0.99, Rec=0.99
 - XGBoost: marginally better, but sensitive to hyper-parameters

Supervised learning

- Random forest for matching F-msr: >95% w. ~1M labels
- AL for blocking & matching F-msr: 80%-98% w. ~1000 labels

labels

~2015 (ML)

- Expt 2. IMDb vs. Amazon movies
 - 200K labels, ~150 features
 - Random forest: Prec=0.98, Rec=0.95



State-of-the-Art ML Models [Das et al., SIGMOD'17]

Supervised learning

- Random forest for matching F-msr: >95% w. ~1M labels
- AL for blocking & matching F-msr: 80%-98% w. ~1000 labels

~2015 (ML)

 Falcon: apply active learning both for blocking and for matching; ~1000 labels

Dataset	Accuracy (%)		(%)	Cost
	P	R	F_1	(# Questions)
Products	90.9	74.5	81.9	57.6(960)
Songs	96.0	99.3	97.6	\$54.0(900)
Citations	92.0	98.5	95.2	65.5(1087)

Supervised learning

~2015 (ML)

- Random forest for matching F-msr: >95% w. ~1M labels
- AL for blocking & matching F-msr: 80%-98% w. ~1000 labels

Apply active learning to minimize #labels



For 99% precision and recall, active learning reduces #labels by 2 orders of magnitude

Reaching prec=99% and rec=~99% requires 1.5M labels



Deep Learning Models [Mudgal et al., SIGMOD'18]

Check-out at poster session on Wednesday! Code at: deepmatcher.ml

2018 (Deep ML)

Deep learning

- Deep learning
- Entity embedding

- Bi-RNN w. attention
- Similar performance for structured data;
 Significant improvement on texts and dirty data





Deep Learning Models [Trivedi et al., ACL'18]

• LinkNBed: Generate embeddings for entities as in knowledge embedding



Deep Learning Models [Trivedi et al., ACL'18]

- LinkNBed: Generate embeddings for entities as in knowledge embedding
- Performance better than previous knowledge embedding methods, but not comparable to random forest
- Enable linking different types of entities

2018 (Deep ML)

Deep learning

- Deep learning
- Entity embedding

Challenges in Applying ML on EL

- How can we obtain abundant training data for many types, many sources, and dynamically evolving data??
- From two sources to multiple sources



Challenges in Applying ML on EL

- How can we obtain abundant training data for many types, many sources, and dynamically evolving data??
- From one entity type to multiple types



Challenges in Applying ML on EL

- How can we obtain abundant training data for many types, many sources, and dynamically evolving data??
- From static data to dynamic data



Recipe for Entity Linkage

- Problem definition: Link references to the same entity
- Short answers
 - **RF w. attributesimilarity features**
- Production Ready
- DL to handle texts and noises
- End-to-end solution is future work

