Outline

- Part I. Introduction
- Part II. ML for DI
- Part III. DI for ML
 - Training data creation
 - Data cleaning
- Part IV. Conclusions and research directions

Successful ML requires Data Integration



IM GENET MovieLens



COCO is a large-scale object detection, segmentation, and captioning dataset.

Large collections of **manually curated** training data are necessary for progress in ML.

Noisy data is a bottleneck



What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%

Source: Crowdflower

Cleaning and organizing data comprises 60% of the time spent on an analytics of AI project.

50 Years of Data Cleaning

E.F.Codd

Errors within a source and Understanding relations (installment #7). • across sources FDT - Bulletin of ACM SIGMOD, 7(3):23– Transformation workflows • 28, 1975. and mapping rules; domain-Null-related features of DBs knowledge is crucial **1980s 2000s (Data Repairs)** (Normalization) **1990s 1970s (Nulls) Constraints and Probabilities Integrity Constraints** Normal forms to reduce Warehouses) Dichotomies for consistent query answering redundancy and Minimality-based repairs to integrity obtain consistent instances FDs, MVDs etc. Statistical repairs

Data transforms

Part of ETL

• Anomaly detection

Where are we today?

Machine learning and statistical analysis are becoming more prevalent.

Error detection (Diagnosis)

- Anomaly detection [Chandola et al., ACM CSUR, 2009]
- Bayesian analysis (Data X-Ray) [Wang et al., SIGMOD'15]
- Outlier detection over streams (Macrobase) [Bailis et al., SIMGOD'17]





Where are we today?

Machine learning and statistical analysis are becoming more prevalent.

Data Repairing (Treatment)

• Classical ML (SCARE, ERACER) [Yakout et al., VLDB'11, SIGMOD'13, Mayfield et al., SIGMOD'10]

ST

- Boosting [Krishan et al., 2017]
- Weakly-supervised ML (HoloClean) [Rekatsinas et al., VLDB'17]



				Each cell is a
Address	City	State	Zip	
3465 S Morgan ST	Chicago	IL	60608	Constra
3465 S Morgan ST	. Chicago	IL	60609	c3: City, Stat
3465 S Morgan ST	Chicago	IL	60609	~
3465 S Morgan ST	Cicago	IL	60608	External data i
				Ext_Address Ext_0
				3465 S Morgan





Error Detection: MacroBase [Bailis et al., SIGMOD'17]





Streaming Feature Selection

Setup: Online learning of a classifier (e.g., LR)

Goal: Return top-k discriminative features

Weight-Median Sketch

Sketch of a classifier for fast updates and queries for estimates of each weight and comes with approximation guarantees

[Figure by Kai Sheng Tai]

A data analytics tool that prioritizes attention in large datasets. **Code at: macrobase.stanford.edu**

Data Repairing: BoostClean [Krishnan et al., 2017]



Ensemble learning for error detection and data repairing.

Relies on domain-specific detection and repairing.

Builds upon boosting to identify repairs that will maximize the performance improvement of a downstream classifier.

On-demand cleaning!

Scalable machine learning for data enrichment



loloClean

Code available at: http://www.holoclean.io



Data Repairing: HoloClean [Rekatsinas et al., VLDB'17]



Holistic data cleaning framework: combines a variety of heterogeneous signals (e.g., integrity constraints, external knowledge, quantitative statistics)



Data Repairing: HoloClean [Rekatsinas et al., VLDB'17]

	Address	City	State	Zip
t1	3465 S Morgan ST	Chicago	IL	60608
t2	3465 S Morgan ST	Chicago	IL	60609
t3	3465 S Morgan ST	Chicago	IL.	60609
t4	3465 S Morgan ST	Cicago	IL	60608



Scalable learning and inference: Hard constraints lead to complex and nonscalable models. Novel relaxation to features over individual cells.



Data Repairing: HoloClean [Rekatsinas et al., VLDB'17]



HoloClean is 2x more accurate. Competing methods either do not scale or perform no correct repairs.

HoloClean: our approach combining all signals and using inference Holistic[Chu,2013]: state-of-the-art for constraints & minimality KATARA[Chu,2015]: state-of-the-art for external data SCARE[Yakout,2013]: state-of-the-art ML & qualitative statistics

Probabilistic Unclean Databases [De Sa et al., 2018]

Unclean Database Generation

(A) Schema, Attribute Domain, and Constraint Specification

Tuple ID		Business Listing				Integrity Constraints	
	Tuple Identifiers	Business ID	City	State	Zip Code	PK: Business ID FD: Zip Code → City, State	

(B) The Two-Actor Generation Process **Tuple Identifiers** Business City State Zip Code ID Tuple Constraints t1 Porter Madison WI 53703 t2 Graft Madison WI 53703 Generator Φ t3 EVP Coffee Madison WI 53703 Intentional Data Model ${\cal I}$ Sample of clean intended data] Business City State Zip Code ID t1 Porter Madison WI 53703 Realizer t2 53703 Graft Verona WI t3 EVP Coffee 53703 Madison WI t4 60609 Graft Chicago Dirty data instance J* observed after applying the Realizer

A two-actor noisy channel model for managing erroneous data.

Preprint: A Formal Framework For Probabilistic Unclean Databases

https://arxiv.org/abs/1801.06750

Challenges in Data Cleaning

- Error detection is still a challenge. To what extent is ML useful for error detection? Tuple-scoped approaches seem to be dominating. Is deep learning useful?
- We need a formal framework to describe when automated solutions are possible.
- A major bottleneck is the collection of training data. Can we leverage weak supervision and data augmentation more effectively?
- Limited end-to-end solutions. Data cleaning workloads (mixed relational and statistical workloads) pose unique scalability challenges.

Recipe for Data Cleaning

- Problem definition: Detect and repair data.
- Short answers
 - ML can help partly-automate cleaning. Doma⁻
 expertise is still required.
 - Scalability of ML-based data cleaning methods is a pressing challenge. Exciting systems research!
 - We need more end-to-end systems!

erroneous

_		Each ce	ell is a ra	andom	variable					
te	Zip									
	60608	Constraints introduce								
	60609	correlations c3: City, State, Address \rightarrow Zip								
	60609	~								
	60608	External data introduce evidence								
_		Ext_Address	Ext_City	Ext_State	Ext_Zip					
		3465 S Morgan ST	Chicago	IL	60608					

