

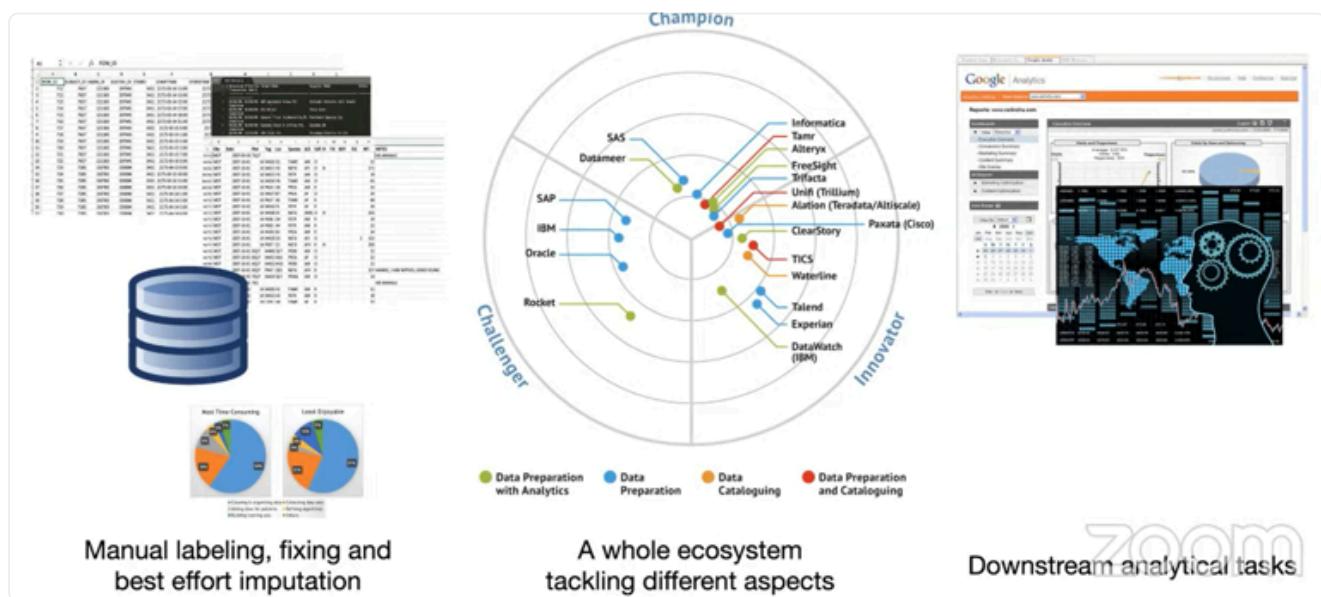
GROUP > SEMINAR & TALK > STANFORD MLSYS SEMINAR

Episode 18

Structure is all you need

- Software 2.0 for data quality management
- Theodoros Rekatsinas | UW-Madison

The Notorious data quality problem



Data quality management task

Error detection tasks

- Tuple (sample) validation

- Cell-value validation

Data repairs

- Missing data imputation
- Data repairs (value replacement)

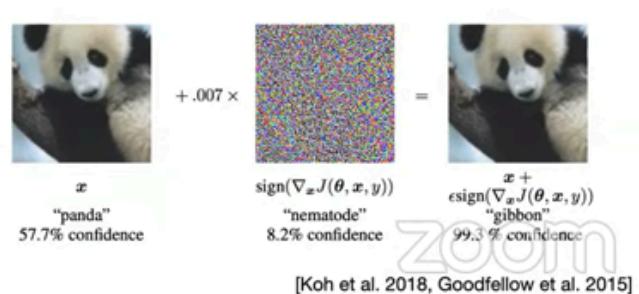
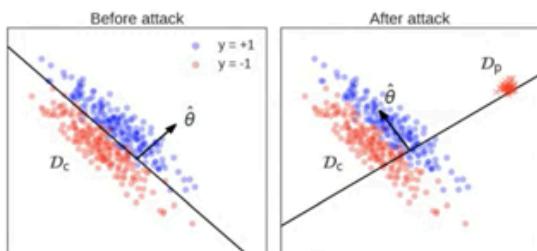
Address	City	State	Zip	Property Tax Rate
640 W Washington Ave	Madison	WI	53703	1.76
652 W Washington Ave	Madison	WI	53703	1.76
2707 W E Washington Ave	Madison	WI	53704	1.05
641 W Washington Ave	adison, WI	53703	53703	1.65

Invalid tuple
 Erroneous value
 Data repair
 Value imputation

- Push to ML model?

ML models are sensitive to low-quality data

- In the training stage, corrupted data can bias the ML models
- Data poisoning techniques destroy a model by adding a small fraction of adversarially crafted points
- In the inference stage, errors in the inference queries can result in wrong predictions
- Adversarial attacks add noise to flip the prediction



Goal: Streamline data quality management

- **Many projects:** mainly HoloClean (Inductiv) and recently Picket and Marius
- **Question:** Can we automate data quality management tasks?
- **Focus:** data validation tasks (i.e., data cleaning, out-of-distribution detection)
- **Automation opportunity:** adopt probabilistic semantics and rely on self-supervised ML approaches to learn how clean data looks like and use that to drive automation

Reasoning about **structured context** is key for automated data validation

- Marius: graphs, heterogeneous-structure data

Example: data validation for mean estimation

The simplest analytical model: compute the mean \mathbf{X} of the following sample

X1	X2	X3	X4	X5
2.39	1.34	1.05	0.01	0.35
NaN	-1.3	1.92	NaN	1.14
0.19	NaN	1.76	1.90	NaN
0.99	-2.16	3.15	1.72	-3.15

Coordinate-wise mean: $\mathbf{X} = (1.19, -0.706, 1.97, 1.21, -0.553)$

Mean-after filtering: $\mathbf{X} = (1.69, -0.41, 2.1, 0.865, -1.4)$

- Discrepancies between two estimate

X1	X2	X3	X4	X5
2.39	1.34	1.05	0.01	0.35
0.62	-1.3	1.92	NaN	1.14
0.19	-1.57	1.76	1.90	NaN
0.99	-2.16	3.15	1.72	-3.15

Improved estimate $\mathbf{X} = (1.048, -0.923, 1.97, 1.21, -0.553)$; why?

- Filling some of the values, and leaving the others

X1	X2	X3	X4	X5
2.39	1.34	1.05	0.01	0.35
0.62	-1.3	1.92	NaN	1.14
0.19	-1.57	1.76	1.90	NaN
0.99	-2.16	3.15	1.72	-3.15

Improved estimate $\mathbf{X} = (1.048, -0.923, 1.97, 1.21, -0.553)$; why?

Exploit the fact that $X_1 = X_2 + X_3$, and X_4, X_5 are independent

- If know this dependency in advanced, then we could make better estimate

Structure-aware data cleaning is necessary

- Two-step meta-algorithm for robust mean estimation:
 - (1) Recover: use the dependencies across coordinates of the data (i.e., the structure) to recover the values of corrupted samples
 - (2) Estimate: After fixing corruptions, perform standard statistical estimation
- Information theoretically optimal error-estimation in the presence of structure:

Structure	Entire-sample corruption (epsilon-samples corrupted)	Coordinate-level Corruption (alpha entries corrupted)
No Structure	$\Theta(\epsilon)$	$\Theta(\alpha n)$
Linear Structure ($x = Az$)	$\Theta(\epsilon)$	$\Theta(\alpha n/m_A)$

 n : number of coordinates ϵ : % of corrupted samples α : % of corrupted entries m_A : number of corruptions to reduce the row space by one

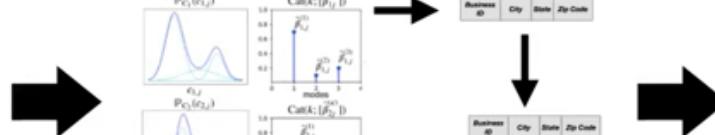
Ref: Structure-aware robust mean estimation [Liu et al., 2020]
<https://arxiv.org/abs/2002.04137>



Heterogeneous types of structure

Domain constraints		External knowledge	Dataset statistics																																						
<table border="1"> <thead> <tr> <th>Address</th> <th>City</th> <th>State</th> <th>Zip</th> <th>Property Tax Rate</th> </tr> </thead> <tbody> <tr> <td>640 W Washington Ave</td> <td>Madison</td> <td>WI</td> <td>53703</td> <td>1.78</td> </tr> <tr> <td>652 W Washington Ave</td> <td>Madison</td> <td>WI</td> <td>53703</td> <td>1.78</td> </tr> <tr> <td>2707 W E Washington Ave</td> <td>Madison</td> <td>WI</td> <td>53704</td> <td>1.65</td> </tr> <tr> <td>641 W Washington Ave</td> <td>Madison, WI</td> <td>WI</td> <td>53703</td> <td>1.65</td> </tr> </tbody> </table> <p><i>Integrity Constraints</i> (e.g., functional dependencies) Zip → State</p> <p><i>Schema Constraints</i> City: Text State: Char(2) Zip: ^[0-9]{5}</p>	Address	City	State	Zip	Property Tax Rate	640 W Washington Ave	Madison	WI	53703	1.78	652 W Washington Ave	Madison	WI	53703	1.78	2707 W E Washington Ave	Madison	WI	53704	1.65	641 W Washington Ave	Madison, WI	WI	53703	1.65	<p>US Street Address API</p> <p>HTTP Request: URL Composition Proper URL construction is required for all API request https://us-street.api.smartyatstreet.com/</p>  <p>External catalogues and knowledge bases</p>	<table border="1"> <thead> <tr> <th>State</th> <th>Zip</th> </tr> </thead> <tbody> <tr> <td>WI</td> <td>53703</td> </tr> </tbody> </table> <p>$P(WI \mid 53703) = 0.98$</p>  <p>Data redundancy allows us to recover statistical dependencies</p>	State	Zip	WI	53703										
Address	City	State	Zip	Property Tax Rate																																					
640 W Washington Ave	Madison	WI	53703	1.78																																					
652 W Washington Ave	Madison	WI	53703	1.78																																					
2707 W E Washington Ave	Madison	WI	53704	1.65																																					
641 W Washington Ave	Madison, WI	WI	53703	1.65																																					
State	Zip																																								
WI	53703																																								
WI	53703																																								
WI	53703																																								
WI	53703																																								
WI	53703																																								
WI	53703																																								

Contextual ML for automated data quality

Context	Generative Models	Data Quality Tasks
Domain Constraints External Knowledge Dataset statistics	 <p>Learn a model of how clean data is generated and how errors are introduced</p>	<p>Tuple (sample) validation Cell-value validation Data repairs (imputation)</p> <p>Data quality ops are inference queries</p>

- Unsupervised manner?
 - Can I run this inference queries fast

Probabilistic Unclean Databases (ICDT'19) Theoretical Noisy-channel Framework for Structured Data	Structure-aware Recovery (UAI'19) Inference over Structured Noisy Data with near-optimal Guarantees
FDX (SIGMOD'20) Contextual Profiling and Structure Discovery	HoloClean (VLDB'17; SIGMOD'19; MLSys'20) From Probabilistic Data Cleaning to Attention-based Data Cleaning
	Picket (Under Submission 2021) Self-supervised Transformers for Data Validation in ML pipelines

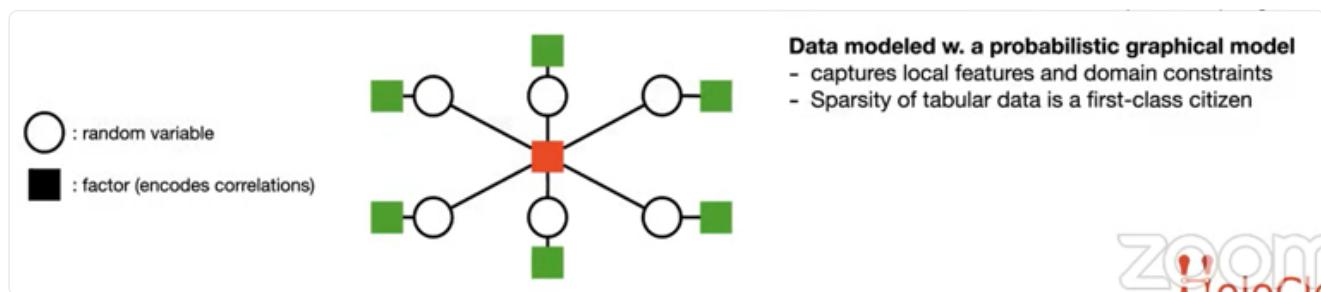
HoloClean: Probabilistic Data Repairs

Each cell is a random variable

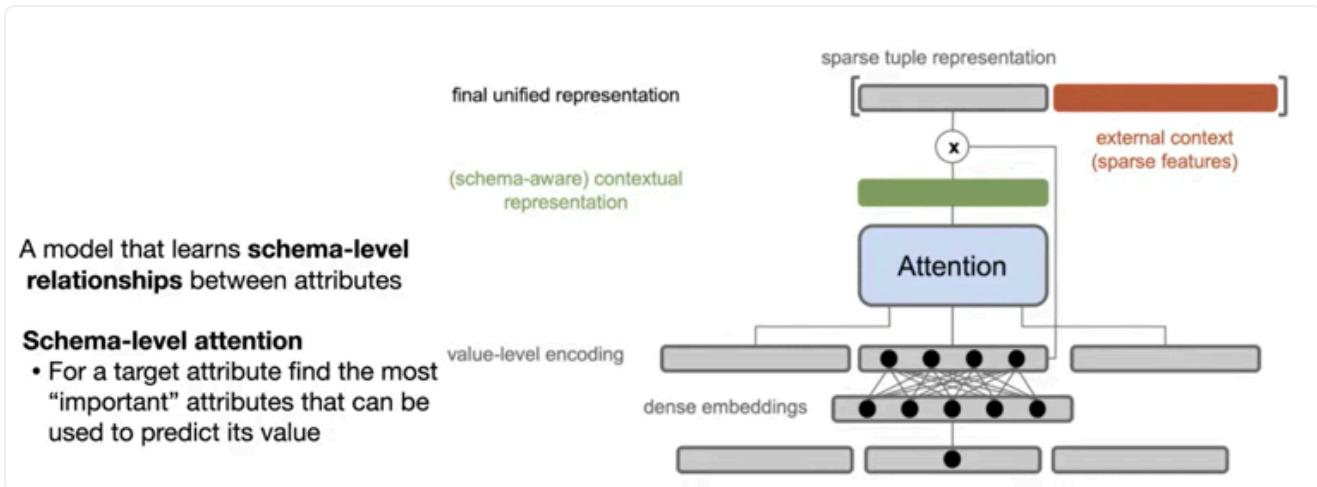
Value co-occurrence capture statistics

Trip ID	Taxi ID	Pickup Census Tract	Pickup Centroid Latitude	Pickup Centroid Longitude
0f5f37514df104...	4864601f96d54f...	17031010400	42.004764559	-87.659122427
c0768e8fa92fbb...	4864601f96d54f...	17031010400	42.004764559	-87.659122427
c5877ebeb9513...	4864601f96d54f...	17031010400	42.004764559	-87.659122427
c0de690f3cd8e...	4979769b1242f...	17031010501	42.009412547	-87.663958214
c15b88fd6dab7...	4979769b1242f...	17031010501	42.009412547	-87.663958214
68f082244ea6fe...	4979769b1242f...	17031010501	42.009412547	-87.663958214

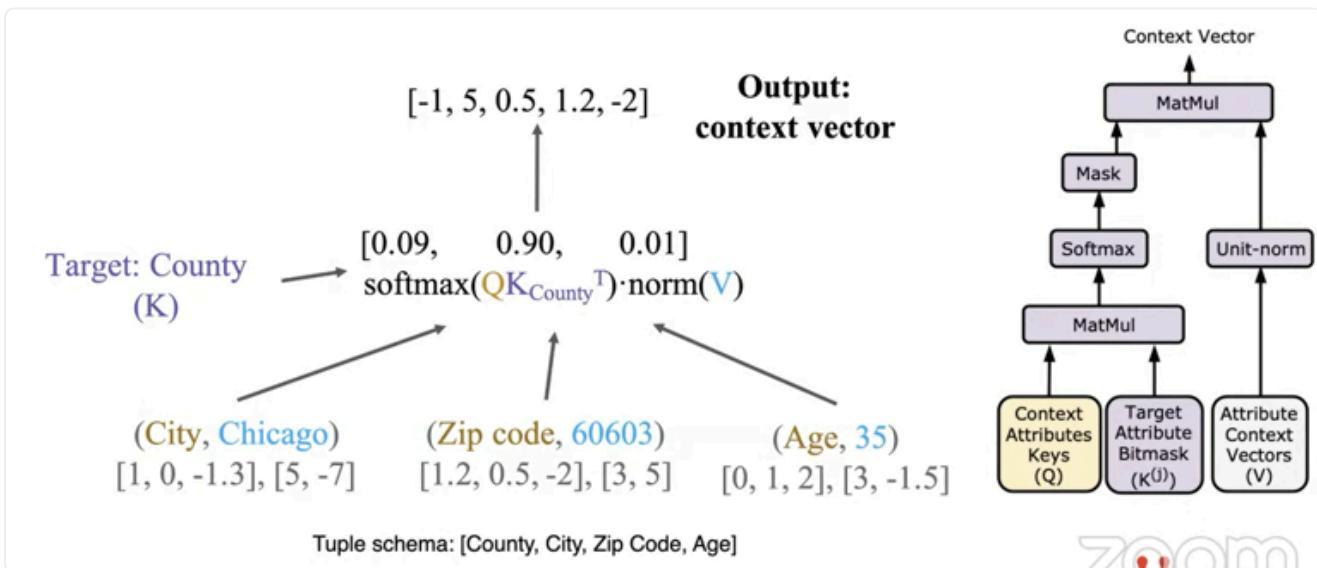
Domain constraints describe correlations
Pickup Lat. \wedge Pickup Long. \rightarrow Pickup Tract



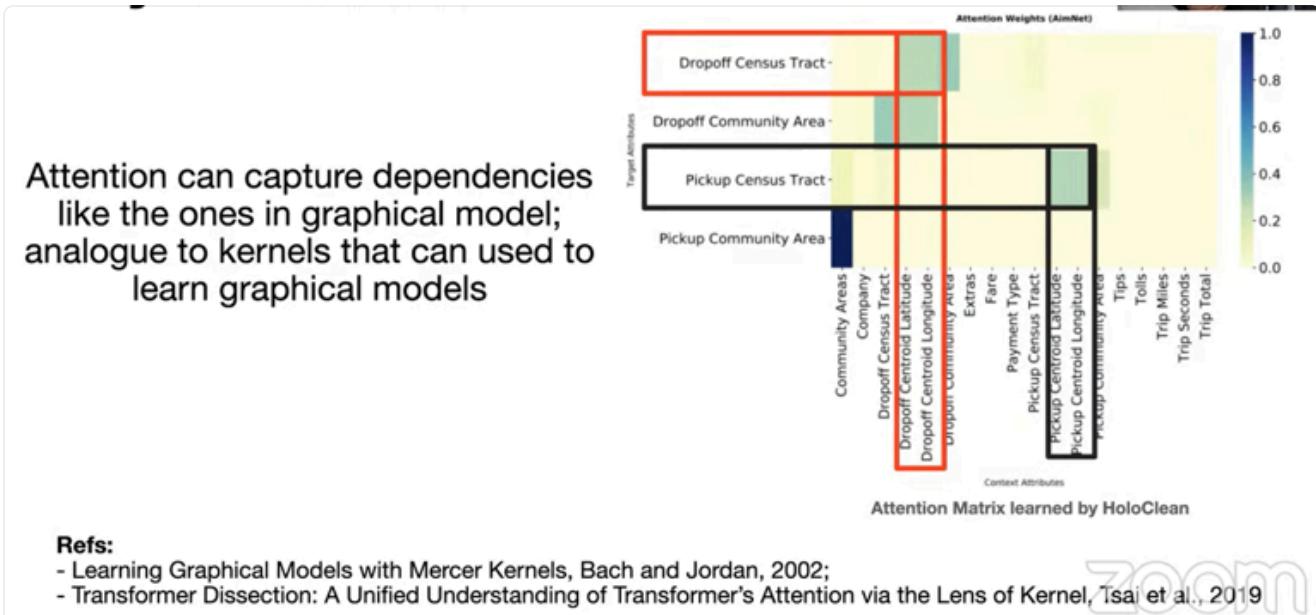
("Take two")



Schema-level Attention



Why Attention?



Naturally-occurring missing data

Chicago taxi data

- Benchmark in Google's TFX data validation
- Pickup/drop-off info, fare, company
- Naturally-occurring missing values w/ ground truth
- Systematic bias between companies

	Company	Pickup Census Tract	Pickup Centroid Latitude	Pickup Centroid Longitude
2366	Chicago Medallion Leasing INC	nan	41.975171	-87.687516
78445	Dispatch Taxi Affiliation	nan	41.975171	-87.687516
57109	Taxi Affiliation Services	17031040401	41.972036	-87.686100

All within "17031040401" census tract

- 40+ F1 points improvements over next-best method
- 10x faster

AimNet (MLSys'20)	HCQ	XGB	MIDAS	GAIN	MF	MICE
HoloClean with Attention	HoloClean with quantization	XGBoost	Denoising Autoencoder	GAN	Random Forest	Linear regression

Accuracy on discrete attributes for the Chicago data set						
AimNet	HCQ	XGB	MIDAS	GAIN	MF	MICE
0.73 ± 0.01	0.07 ± 0.0	0.27 ± 0.0	0.09 ± 0.01	0.01 ± 0.01	0.3 ± 0.0	—

Run time (minutes) for the Chicago data set						
AimNet	HCQ	XGB	MIDAS	GAIN	MF	MICE
53	124	5350	176	186	7439	—

24

Use case: Data Categorization

Problem: Market research data (purchase transactions) missing the canonical ID (SKU) for the purchased item

Previous approach:

- Manual annotation
- Labor-intensive process taking months

Automating with HoloClean:

- Accuracy > 97% in less than 15 hours
- Wrong human-labeled data identified

k	Accuracy	Avg. Confidence
1	96.8%	97.22%
2	99.4%	95.2%
3	99.8%	94.8%

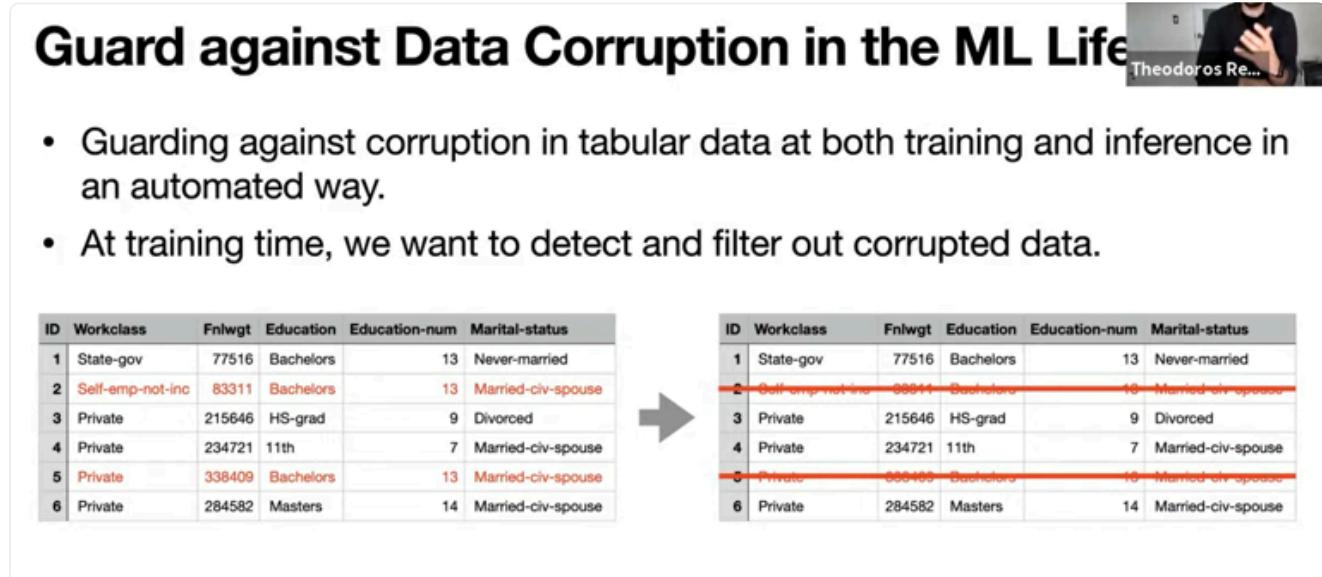
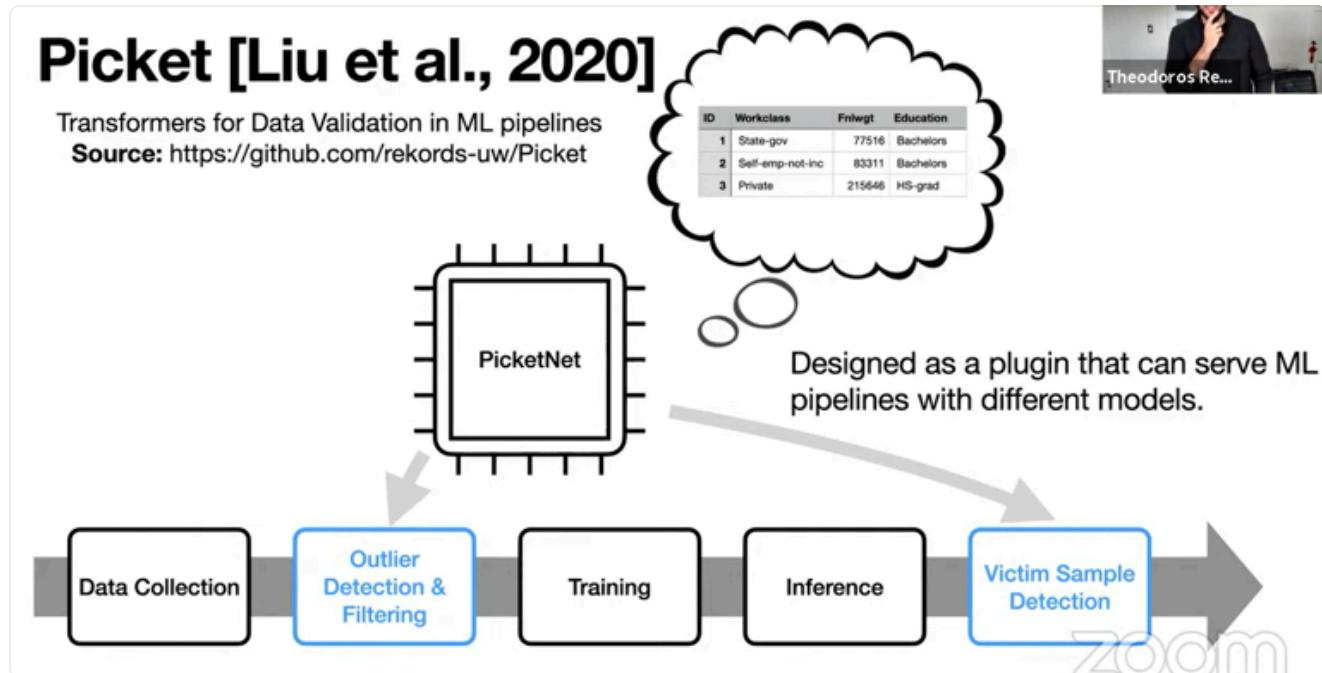
Error Type	Sampled
Ground truth is incorrect	1128 (71.7%)
Prediction is Incorrect	333 (21.1%)
Uncertain	112 (7.1%)



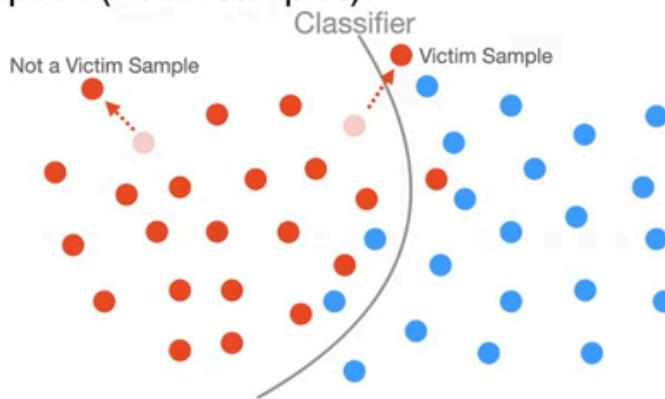
Other use cases

- Error detection in demographic data used for policy decision
- KIP tracking
- Imputation of numerical data for industrial machinery monitoring

Picket: self-supervised transformers for data validation in ML pipelines



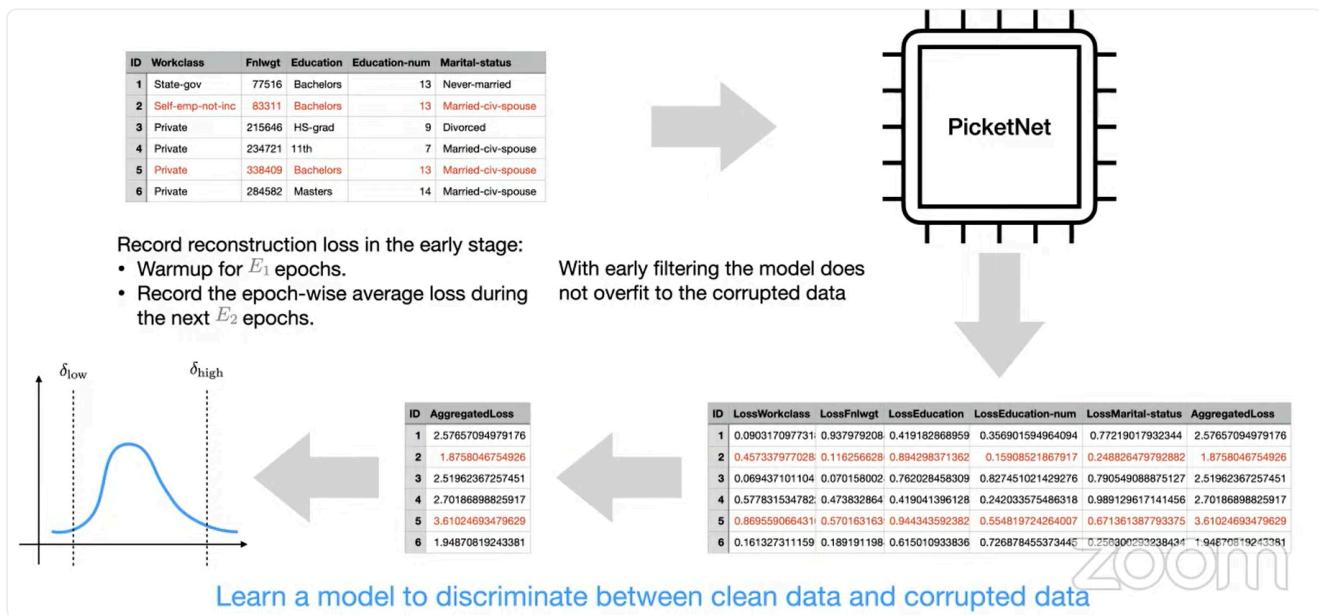
- Guarding against corruption in tabular data at both training and inference in an automated way.
- At training time, we want to detect and filter out corrupted data.



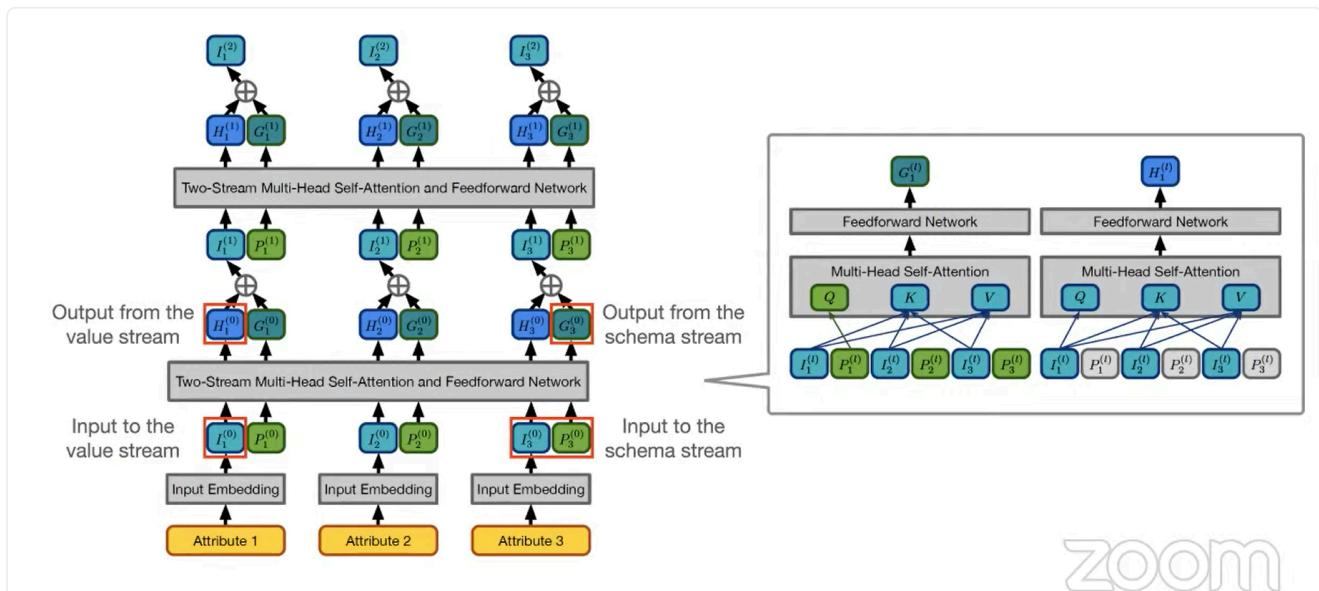
zoom

Loss-based Outlier Detection and Filtering

1. Go back to the idea of learning a model to capture the clean data, and use this model in decisions
 - a. PicketNet: transformer
 - b. Outlier detection problem



PicketNet: two-stream transformer for tabular data



- Benefits:

Value stream: flexibility

Reading List

Schema stream: regularization



Experimental Highlights: Poisoning Attacks

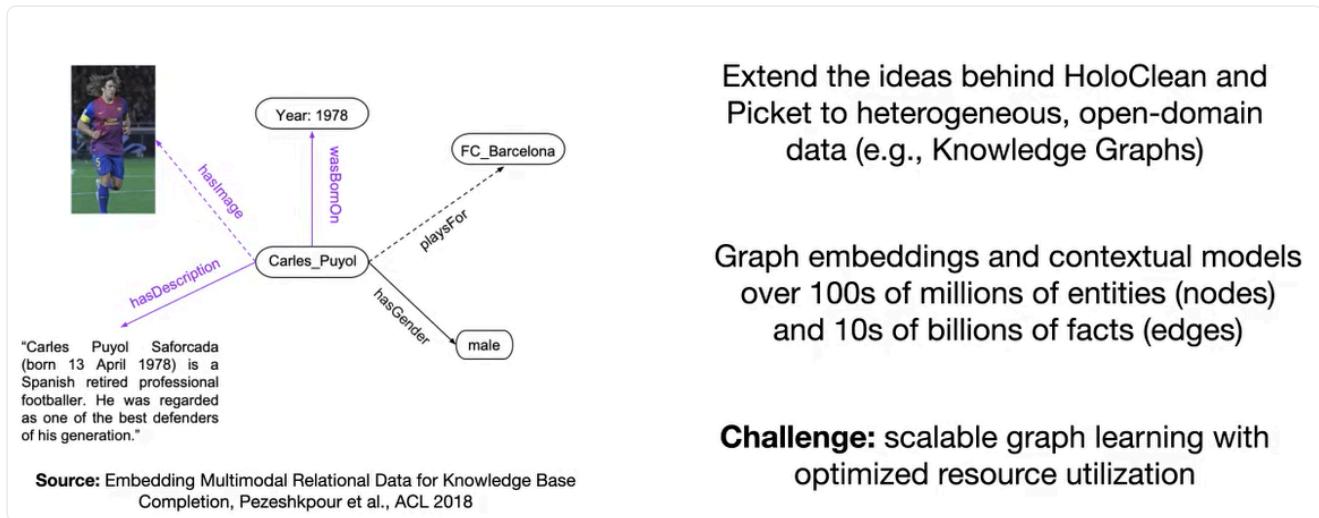
- Aim to destroy

Dataset	Downstream Model	Clean training set				Corrupted training set with no filtering	
		IF	OCSVM	RVAE	Picket	CL	NF
Wine	LR	0.7261	0.6976	0.7051	0.7312	0.7349	0.6745
	SVM	0.7286	0.6933	0.7082	0.7310	0.7386	0.6727
	NN	0.7210	0.6894	0.7035	0.7320	0.7365	0.6722
HTRU2	LR	0.8884	0.9015	0.8811	0.9067	0.9396	0.8799
	SVM	0.8884	0.8979	0.8887	0.9232	0.9424	0.8832
	NN	0.8671	0.8707	0.8643	0.9000	0.9280	0.8646

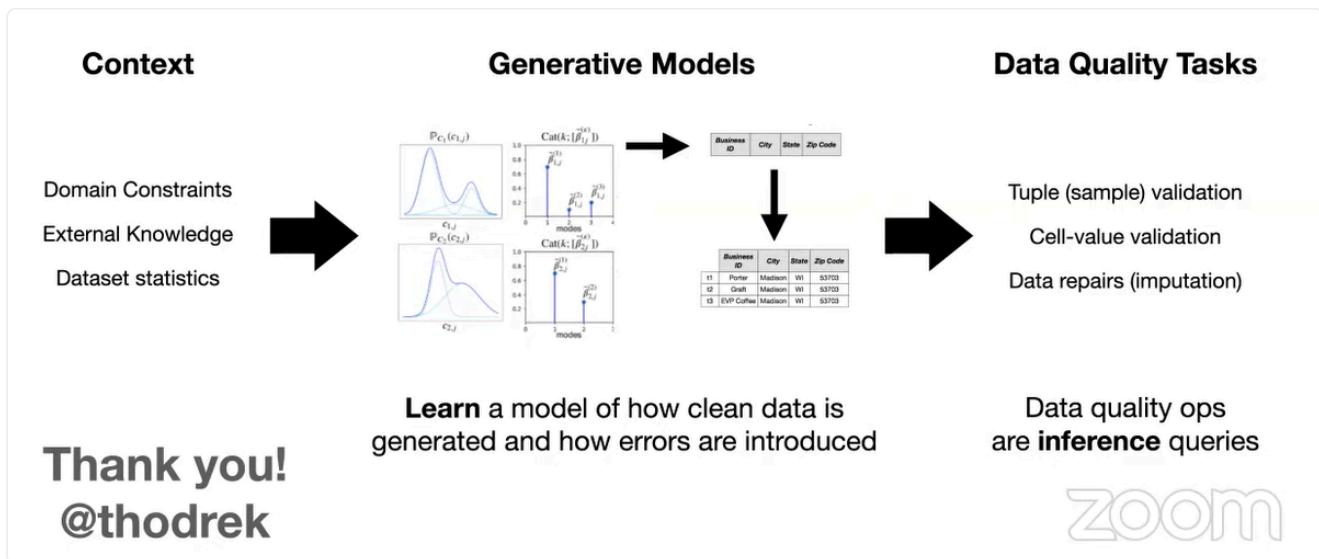
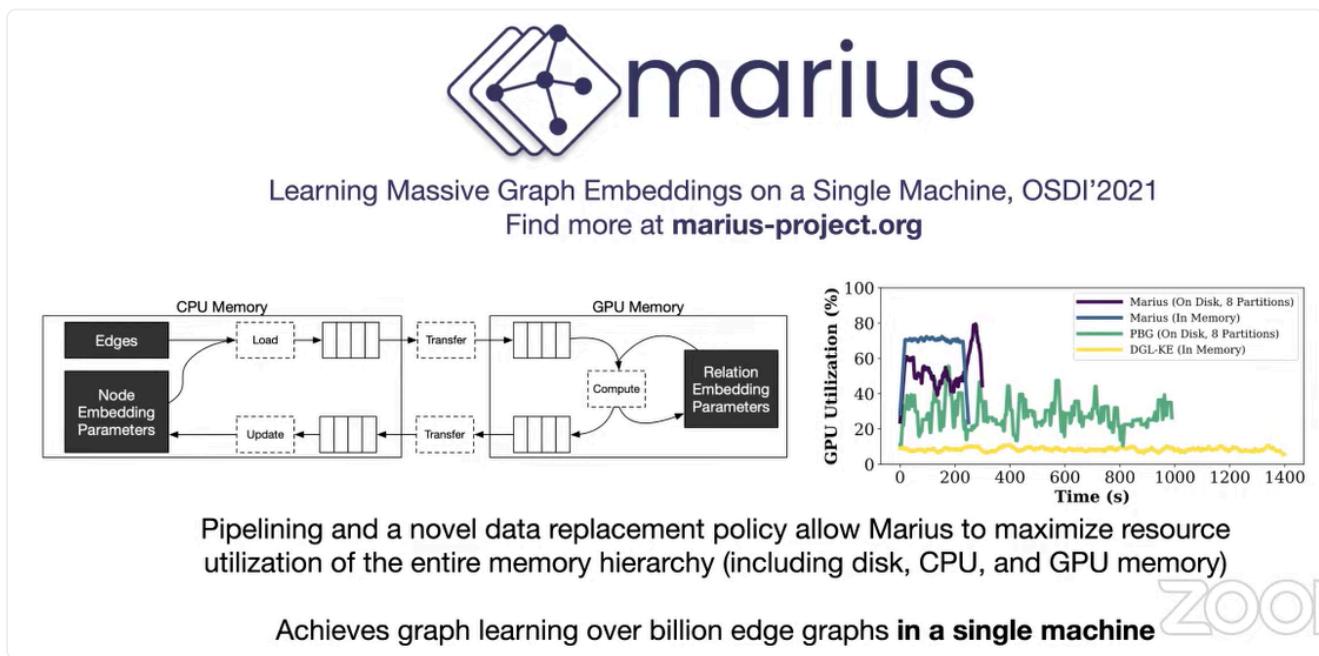
Downstream Accuracy

With Picket, the downstream accuracy is much better than no filtering, and gets close to the clean training set.

Contextual ML for automated data quality ops



Scalable no-code graph learning



- Systematic variation of the data?
 - Types of noise: real-world data, no assumption about the noise.
Types of noise can be random, or systematic error (integrating things across different sources)
 - Random
 - Systematic: repeated instances of the noise, if we condition, then it's not random
 - Adversarial noise: are aware of the downstream task, and go and attack that system
- Holoclean
 - Attention: handle this gracefully, pick up the type of strong bias
- Picket:
 - Worst-possible case (Adversarial)
 - Not overfitting to examples
- Distinguish between out-of-distribution or systematic change
 - Solution: two streams (scheme, value)
 - Value: robust towards the case. Kernel structure that operates on this level.
 - Profiling mechanism
- Heterogeneous data types
 - Tabula data: higher level constraints
 - Encode them as functional dependencies in DB
 - or pick them up through attention mechanism
 - Back in the day: user specification
 - Structure learning over the data
 - Exactly the attention matrix, in a faster and cheaper way
 - Heterogenous

- Same mechanism can potentially hold for a graph
- Running structural learning type of profiling
 - Identify homogenous area as a pre-processing step
 - And preprocessing ...
 - Filtering away and keep
- Doing this heavily in Holoclean
- Monitor data and see if something is going on in the data pipeline
 - Reliable data
 - How that setting is different? Or some of the goals change?
 - Reconstruction
 - Signal and context
 - With high confidence, then it should be outlier
 - Information about the likelihood

Goal:

- applying the rules at scale (ETL, ...)
- Start-up and companies: specific problem
 - Identify duplicates in records
 - Infer rules to prepare and standardize
 - AI --> platform for error detection and fixing
- Nobody is targeting
 - Automating this
 - Position: reasoning about noisy structured data

Challenges and what that issue looks like:

- What is the model doing? Aspects

- Attention: interpretable (know the semantics of the attributes, put more weights or less)
- Allow people not immediately accept them. Have confidence over the prediction of the model.
 - Accept the one that makes sense
- Also, allow users to introduce business /external features that would allow you to introduce

Holoclean

- Accept logic rules, convert them to features
- Support matching functions

Which one should I trust?

- Ensemble (weighted vote)
- In real cases, people believe their rules...

Previous
Episode 17

Next
Index

Last updated 3 years ago