

# Attention-based Learning for Missing Data Imputation in HoloClean

Richard Wu<sup>1</sup>, Aoqian Zhang<sup>1</sup>, Ihab F. Ilyas<sup>1</sup>  
Theodoros Rekatsinas<sup>2</sup>



UNIVERSITY OF  
**WATERLOO**



**WISCONSIN**  
UNIVERSITY OF WISCONSIN-MADISON

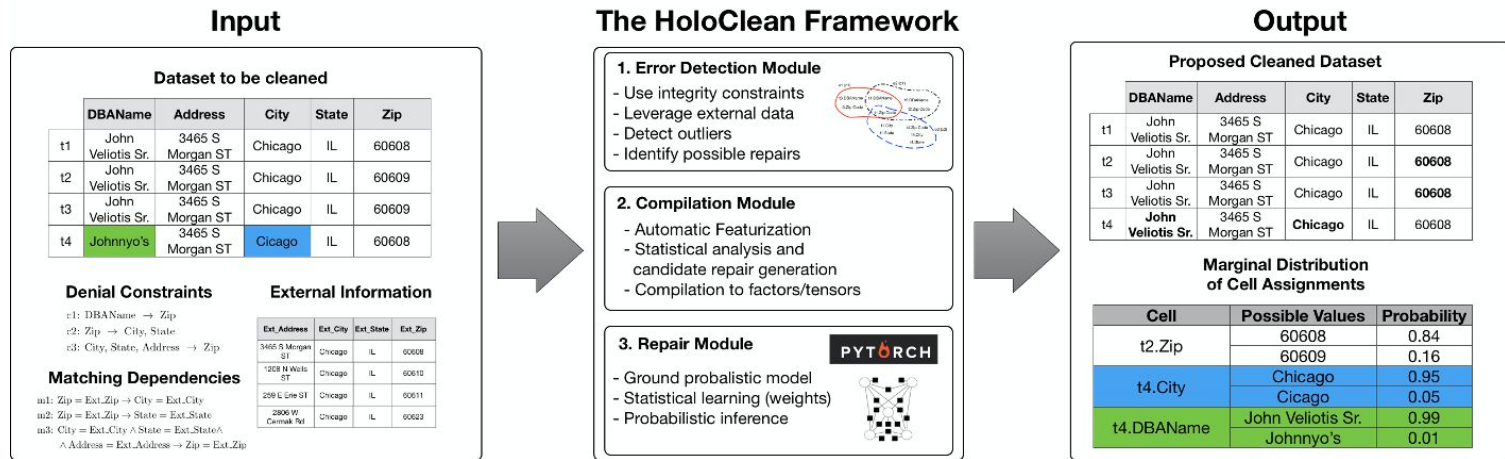
# Problem

- Missing data is a **persistent** problem in many fields
  - Sciences
  - Data mining
  - Finance
- Missing data can reduce downstream statistical power
- Most models require complete data

# Modern ML for Data Cleaning: HoloClean

- Framework for holistic data repairing driven by probabilistic inference
- Unifies **qualitative** (integrity constraints and external sources) with **quantitative** data repairing methods (statistical inference)

Available at [www.holoclean.io](http://www.holoclean.io)



# Missing Values in *Real* Data sets

	First Name	Gender	Start Date	Last Login Time	Salary	Bonus %	Senior Management	Team
0	douglas	Male	8/6/1993	12:42 PM	97308	6.945	True	Marketing
1	thomas	Male	3/31/1996	6:53 AM	61933	4.170	True	NaN
2	maria	Female	4/23/1993	11:17 AM	130590	11.858	False	Finance
3	jerry	Male	3/4/2005	1:00 PM	138705	9.340	True	Finance
4	larry	Male	1/24/1998	4:47 PM	101004	1.389	True	Client Services
5	dennis	Male	4/18/1987	1:35 AM	115163	10.125	False	Legal
6	ruby	Female	8/17/1987	4:20 PM	65476	10.012	True	Product
7	NaN	Female	7/20/2015	10:43 AM	45906	11.598	NaN	Finance
8	angela	Female	11/22/2005	6:29 AM	95570	18.523	True	Engineering
9	frances	Female	8/8/2002	6:51 AM	139852	7.524	True	Business Development
10	louise	Female	8/12/1980	9:01 AM	63241	15.132	True	NaN
11	julie	Female	10/26/1997	3:19 PM	102508	12.637	True	Legal
12	brandon	Male	12/1/1980	1:08 AM	112807	17.492	True	Human Resources

# Challenges

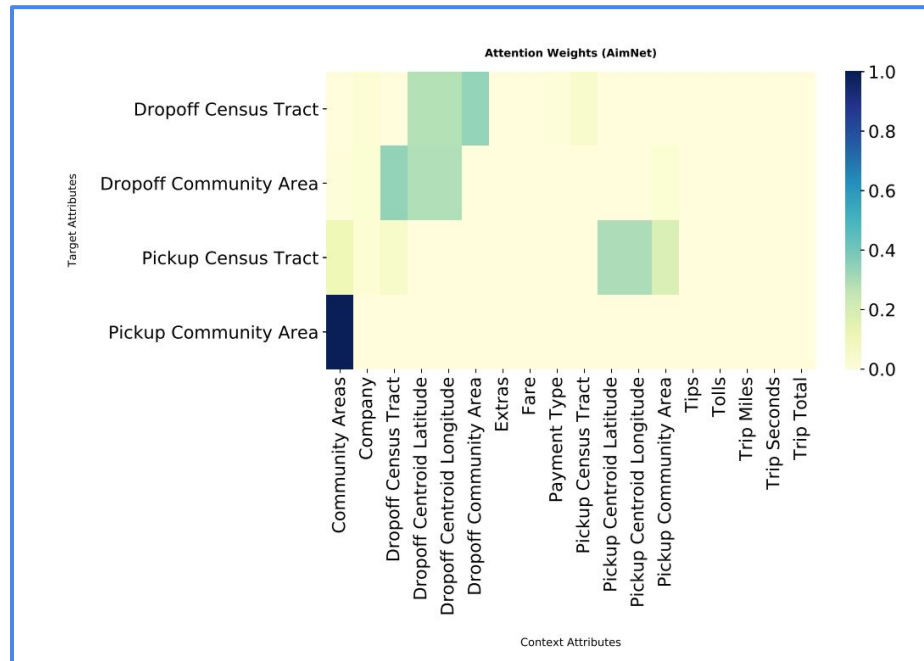
- Values may not be missing completely at random (MCAR/i.i.d.) but **systematically**
- Mixed types (discrete and continuous) introduce **mixed distributions**
- **Drawbacks** of current methods:
  - Heuristic-based (impute **mean/mode**)
  - Requires predefined **rules**
  - **Complex** ML models that are **difficult to train, slow, hard to interpret**

# Contribution

*A **simple** attention architecture that exploits **structure** across attributes*

Our results:

- **>54%** lower run time than baselines
- *Missing at random (MCAR)* : **3%** higher accuracy and **26.7%** reduction in normalized-RMS
- *Systematic*: **43%** higher accuracy and **7.4%** reduction in normalized-RMS



# How does AimNet improve on the MVI problem?

Key idea:  
**Exploit the structure in data**

model that learns **schema-level relationships** between attributes



**dot product attention**

# Architecture overview

## (1) Model mixed data

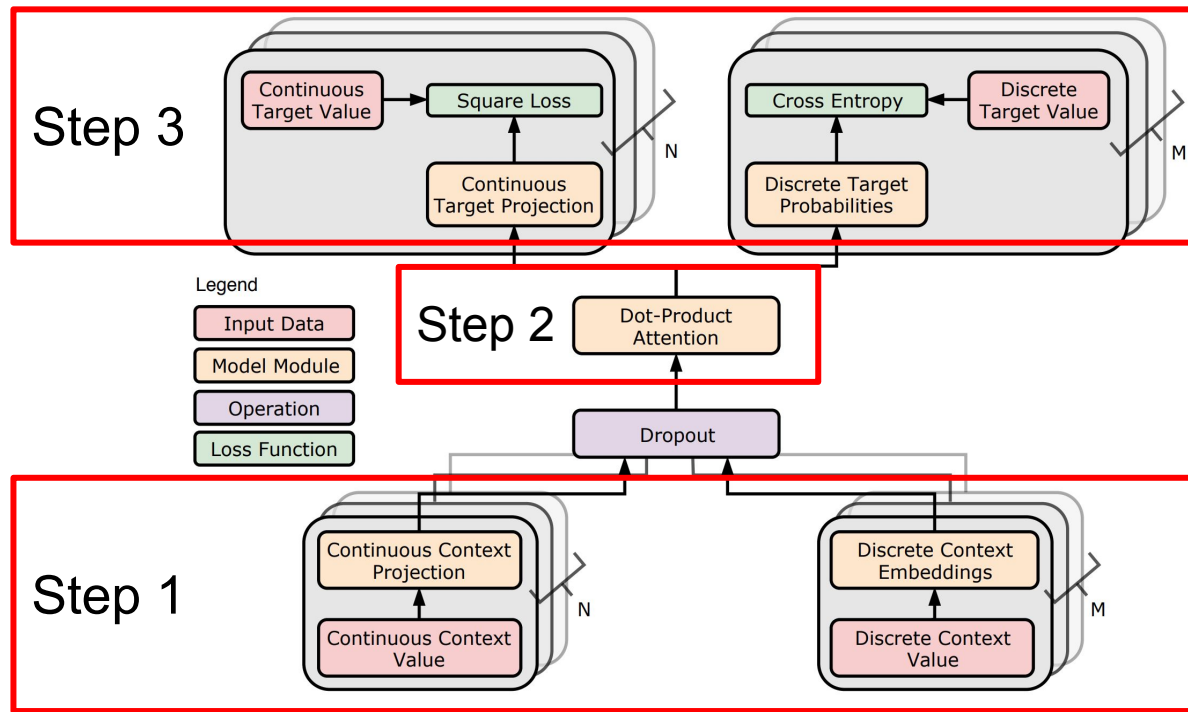
- Encode w/ non-linear layers (continuous)
- Embedding lookup (discrete)

## (2) Identify relevant context

- Attention helps identify schema-level importance

## (3) Prediction

- Inverse of encoding (continuous)
- Softmax over possible values (discrete)

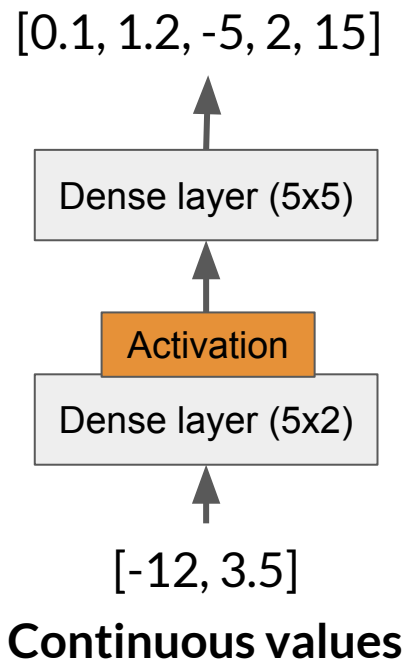


Learned via self-supervision: mask and predict observed values

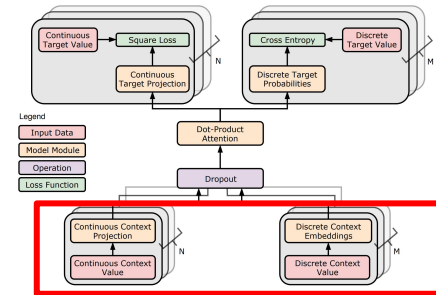
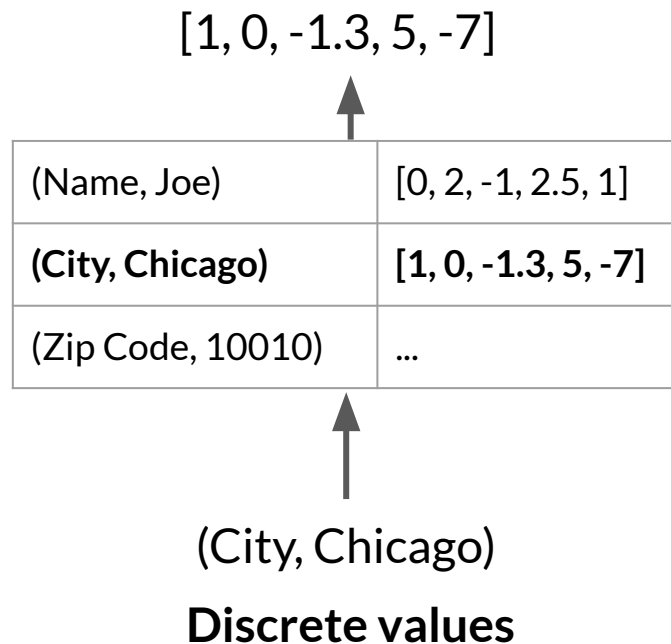


# How do we encode mixed types?

Convert **context values** to **vector embeddings**.

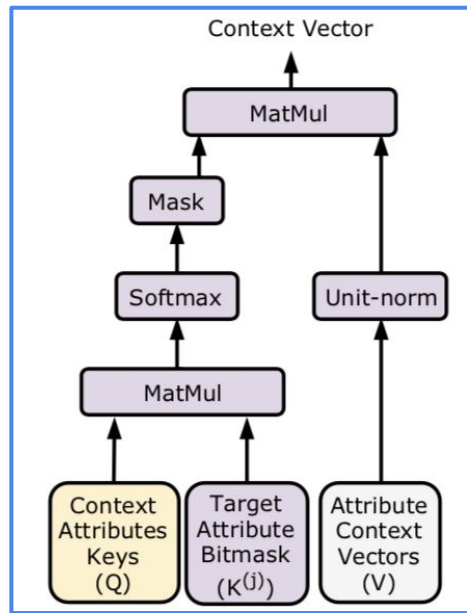
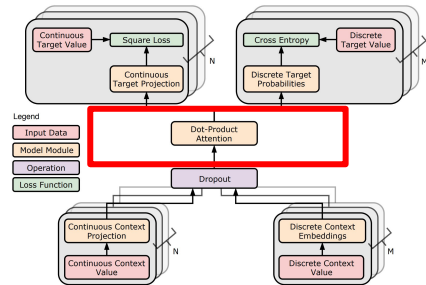
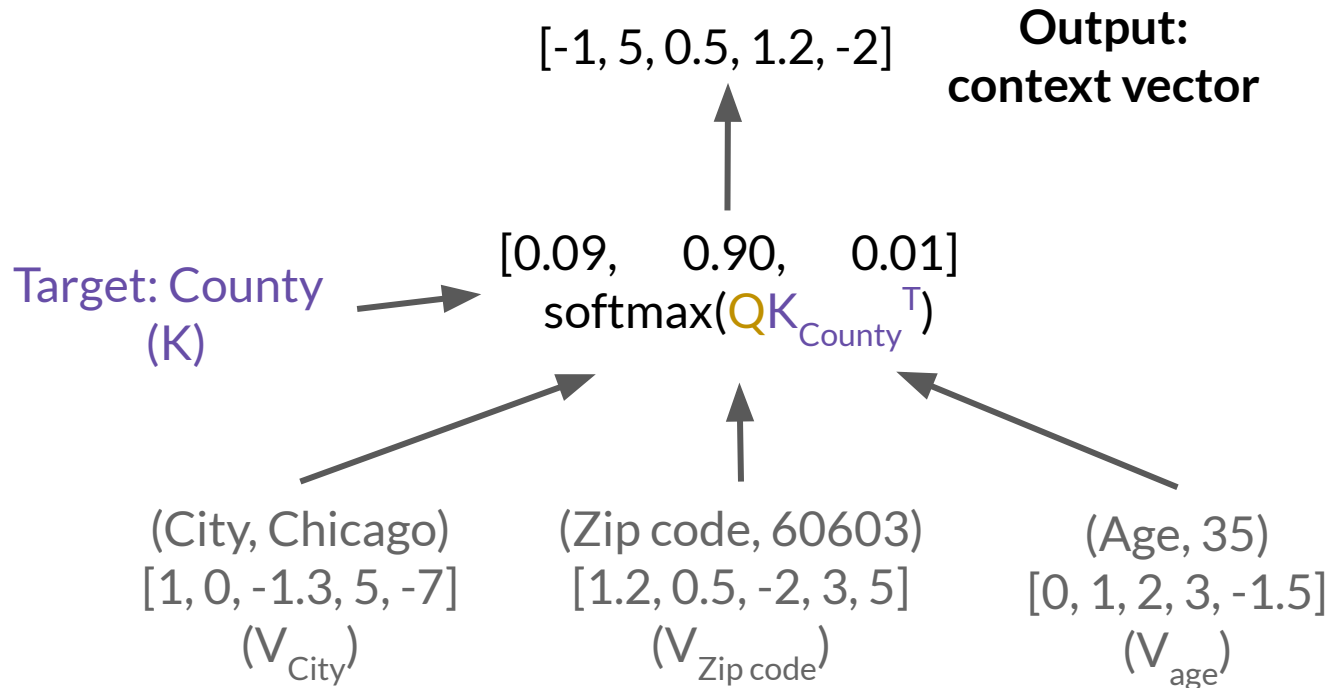


**Output:**  
embeddings



# Attention layer

Attention where Q/K are derived from **attributes** rather than values



# Prediction

Input: context vector

$[-1, 5, 0.5, 1.2, -2]$

Dense layer (5x5)

Activation

Dense layer (1x5)

Output: 100600

Salary (continuous)

matmul

County A:  $[0, 100, 0, 0, 0]^T$

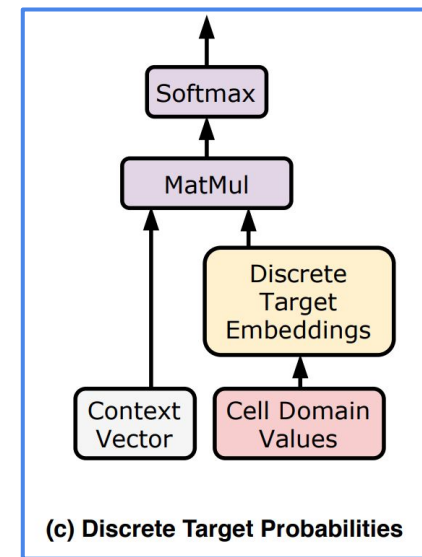
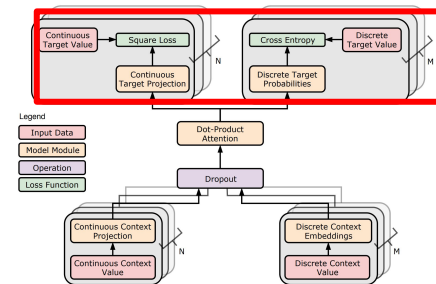
County B:  $[0, 0, 0, 0, 50]^T$

softmax

$[0.99, 0.01]$

Output: County A

County (discrete)



# Questions

- Can AimNet impute missing completely at random (**MCAR/i.i.d.**) values?
- Does AimNet's emphasis on structure help it with **systematic bias** in missing values?
- Can we **interpret** the structure that AimNet learns in the data?

# Experimental setup

- 14 real data sets
- Missing types
  - MCAR/i.i.d.
  - Systematic
- Evaluation
  - Accuracy (discrete)
  - normalized-RMS (continuous)

Data Set	$ r $	# Continuous Attributes	# Discrete Attributes
Tic-Tac-Toe	958	0	10
Hospital	1000	2	14
Mammogram	831	1	5
Thoracic	470	3	14
Contraceptive	1473	2	8
Solar Flare	1066	3	10
NYPD	32399	4	13
Credit	653	6	10
Australian	691	6	9
Chicago	400k	11	7
Balance	625	4	1
Eye EEG	14976	14	1
Phase	9628	4	0
CASP	45730	10	0

Mostly discrete



$$NRMS_j = (\sum_i^{n_j} ((y_i - \hat{y}_i)^2) / (n_j \cdot \sigma(\vec{y})^2))^{-1/2}$$

Mostly continuous

- Training: self-supervised learning where targets = observable values

# Experiment results

- >54% lower run time than baselines
- *Missing at random (MCAR)* : 3% higher accuracy and **26.7%** reduction in normalized-RMS
- *Systematic*: **43%** higher accuracy and **7.4%** reduction in normalized-RMS

**Attention** identifies **structure between attributes** that helps it deal with systematic bias in missing values

# MCAR (20%)

AimNet outperforms on both **discrete** and **continuous** attributes on almost all data sets

- 3% in accuracy
- 26.7% in NRMS

HCQ	XGB	MIDAS	GAIN	MF	MICE
HoloClean with quantization	XGBoost	Denoising Autoencoder	GAN	Random Forest	Linear regression with multiple iterations

data set	AimNet	Accuracy on discrete attributes (ACC $\pm$ std)					
		HCQ	XGB	MIDAS	GAIN	MF	MICE
Tic-Tac-Toe	<b>0.61 <math>\pm</math> 0.01</b>	0.53 $\pm$ 0.01	0.57 $\pm$ 0.02	0.46 $\pm$ 0.01	0.32 $\pm$ 0.01	0.52 $\pm$ 0.01	0.58 $\pm$ 0.02
Hospital	<b>0.99 <math>\pm</math> 0.0</b>	<b>0.99 <math>\pm</math> 0.0</b>	0.97 $\pm$ 0.01	0.24 $\pm$ 0.0	0.13 $\pm$ 0.01	<b>0.99 <math>\pm</math> 0.0</b>	0.82 $\pm$ 0.01
Mammogram	<b>0.75 <math>\pm</math> 0.01</b>	<b>0.74 <math>\pm</math> 0.02</b>	<b>0.74 <math>\pm</math> 0.02</b>	<b>0.74 <math>\pm</math> 0.01</b>	0.35 $\pm$ 0.02	0.68 $\pm$ 0.02	0.64 $\pm$ 0.02
Thoracic	<b>0.86 <math>\pm</math> 0.01</b>	0.84 $\pm$ 0.01	<b>0.85 <math>\pm</math> 0.01</b>	0.84 $\pm$ 0.01	0.59 $\pm$ 0.09	<b>0.86 <math>\pm</math> 0.01</b>	0.38 $\pm$ 0.4
Contraceptive	<b>0.65 <math>\pm</math> 0.01</b>	<b>0.64 <math>\pm</math> 0.01</b>	0.63 $\pm$ 0.01	0.63 $\pm$ 0.01	0.42 $\pm$ 0.01	0.63 $\pm$ 0.02	0.57 $\pm$ 0.01
Solar Flare	<b>0.78 <math>\pm</math> 0.02</b>	<b>0.77 <math>\pm</math> 0.02</b>	0.76 $\pm$ 0.01	0.65 $\pm$ 0.01	0.48 $\pm$ 0.02	0.76 $\pm$ 0.02	0.67 $\pm$ 0.02
NYPD	0.92 $\pm$ 0.0	0.89 $\pm$ 0.0	<b>0.93 <math>\pm</math> 0.0</b>	0.79 $\pm$ 0.01	0.14 $\pm$ 0.01	0.92 $\pm$ 0.0	0.72 $\pm$ 0.0
Credit	<b>0.76 <math>\pm</math> 0.01</b>	0.73 $\pm$ 0.02	0.75 $\pm$ 0.01	0.61 $\pm$ 0.02	0.4 $\pm$ 0.01	<b>0.76 <math>\pm</math> 0.01</b>	0.68 $\pm$ 0.01
Australian	<b>0.72 <math>\pm</math> 0.02</b>	0.69 $\pm$ 0.02	0.71 $\pm$ 0.02	0.61 $\pm$ 0.03	0.45 $\pm$ 0.01	<b>0.73 <math>\pm</math> 0.01</b>	0.63 $\pm$ 0.02
Balance	<b>0.79 <math>\pm</math> 0.04</b>	<b>0.78 <math>\pm</math> 0.04</b>	0.75 $\pm$ 0.05	0.68 $\pm$ 0.08	0.5 $\pm$ 0.05	0.69 $\pm$ 0.05	0.72 $\pm$ 0.05
Eye EEG	0.71 $\pm$ 0.01	0.63 $\pm$ 0.01	0.81 $\pm$ 0.01	0.55 $\pm$ 0.01	0.54 $\pm$ 0.01	<b>0.87 <math>\pm</math> 0.01</b>	0.54 $\pm$ 0.01
data set	AimNet	NRMS on continuous attributes (NRMS $\pm$ std)					
		HCQ	XGB	MIDAS	GAIN	MF	MICE
Hospital	<b>0.72 <math>\pm</math> 0.06</b>	1.4 $\pm$ 0.36	0.87 $\pm$ 0.07	611.12 $\pm$ 129.15	1.19 $\pm$ 0.02	0.86 $\pm$ 0.07	1.13 $\pm$ 0.13
Mammogram	<b>0.91 <math>\pm</math> 0.04</b>	1.03 $\pm$ 0.08	0.96 $\pm$ 0.06	1.12 $\pm$ 0.05	1.0 $\pm$ 0.11	0.99 $\pm$ 0.05	1.27 $\pm$ 0.11
Thoracic	<b>1.1 <math>\pm</math> 0.41</b>	1.78 $\pm$ 2.33	1.76 $\pm$ 2.45	1.71 $\pm$ 1.17	<b>1.5 <math>\pm</math> 0.95</b>	1.66 $\pm$ 1.61	3.62 $\pm$ 4.82
Contraceptive	<b>0.84 <math>\pm</math> 0.02</b>	1.06 $\pm$ 0.05	0.87 $\pm$ 0.04	1.09 $\pm$ 0.03	1.13 $\pm$ 0.05	0.88 $\pm$ 0.04	1.14 $\pm$ 0.06
Solar Flare	<b>0.94 <math>\pm</math> 0.15</b>	<b>0.94 <math>\pm</math> 0.13</b>	1.21 $\pm$ 0.53	17698.23 $\pm$ 8959.74	<b>0.96 <math>\pm</math> 0.16</b>	<b>0.96 <math>\pm</math> 0.2</b>	1.22 $\pm$ 0.34
NYPD	0.15 $\pm$ 0.01	1.28 $\pm$ 0.53	0.14 $\pm$ 0.0	0.62 $\pm$ 0.04	3.19 $\pm$ 0.41	<b>0.1 <math>\pm</math> 0.0</b>	0.37 $\pm$ 0.01
Credit	<b>0.94 <math>\pm</math> 0.03</b>	1.84 $\pm$ 0.83	1.09 $\pm$ 0.15	1.29 $\pm$ 0.24	1.18 $\pm$ 0.09	1.04 $\pm$ 0.13	1.84 $\pm$ 0.83
Australian	<b>0.94 <math>\pm</math> 0.03</b>	2.47 $\pm$ 2.0	1.09 $\pm$ 0.12	1.22 $\pm$ 0.13	1.24 $\pm$ 0.18	1.07 $\pm$ 0.24	1.58 $\pm$ 0.28
Balance	<b>0.92 <math>\pm</math> 0.02</b>	1.37 $\pm$ 0.08	1.0 $\pm$ 0.03	1.02 $\pm$ 0.02	1.03 $\pm$ 0.03	1.08 $\pm$ 0.05	1.26 $\pm$ 0.07
Eye EEG	0.4 $\pm$ 0.0	0.94 $\pm$ 0.34	0.39 $\pm$ 0.0	0.84 $\pm$ 0.01	0.65 $\pm$ 0.04	<b>0.35 <math>\pm</math> 0.0</b>	0.62 $\pm$ 0.0
Phase	<b>0.45 <math>\pm</math> 0.01</b>	0.54 $\pm$ 0.03	<b>0.45 <math>\pm</math> 0.01</b>	0.95 $\pm$ 0.01	0.76 $\pm$ 0.14	0.5 $\pm$ 0.01	0.63 $\pm$ 0.01
CASP	0.45 $\pm$ 0.02	1.45 $\pm$ 0.14	0.43 $\pm$ 0.02	0.82 $\pm$ 0.01	0.72 $\pm$ 0.04	<b>0.41 <math>\pm</math> 0.02</b>	0.64 $\pm$ 0.03

# Chicago taxi data set

- Benchmark in TFX data validation pipeline
- Pickup/dropoff info, fare, company
- Naturally-occurring missing values w/ ground truth
- Systematic bias between companies

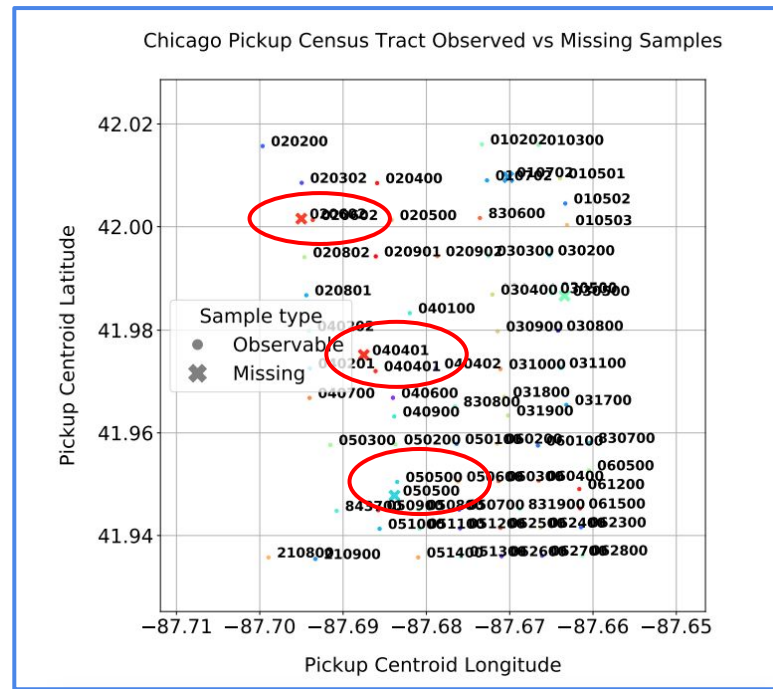
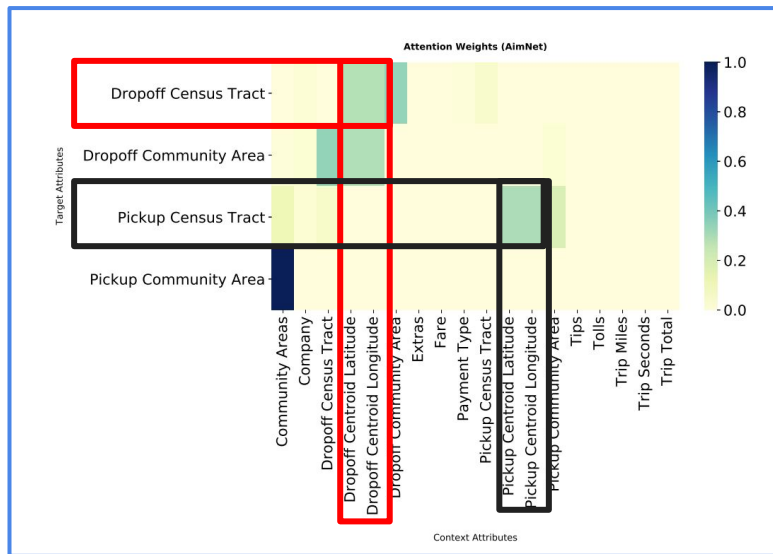
	Company	Pickup Census Tract	Pickup Centroid Latitude	Pickup Centroid Longitude
2366	Chicago Medallion Leasing INC	nan	41.975171	-87.687516
78445	Dispatch Taxi Affiliation	nan	41.975171	-87.687516
57109	Taxi Affiliation Services	17031040401	41.972036	-87.686100

All within "17031040401" census tract



# Chicago taxi: naturally-occurring missing data

- Values are missing systematically (not i.i.d.)
- **Attention** learns relationship between **Census Tract** and **Latitude/Longitude**



# Chicago taxi results

AimNet outperforms baselines by a **huge margin**

- Accuracy: **73%** vs 27% (XGB)
- Run time: **53 mins.** vs 124 mins (HoloClean w/ Quantization)

Accuracy on discrete attributes for the Chicago data set						
AimNet	HCQ	XGB	MIDAS	GAIN	MF	MICE
<b><math>0.73 \pm 0.01</math></b>	$0.07 \pm 0.0$	$0.27 \pm 0.0$	$0.09 \pm 0.01$	$0.01 \pm 0.01$	$0.3 \pm 0.0$	—
Run time (minutes) for the Chicago data set						
AimNet	HCQ	XGB	MIDAS	GAIN	MF	MICE
<b>53</b>	124	5350	176	186	7439	—

# What if we inject systematic errors into other real data sets?

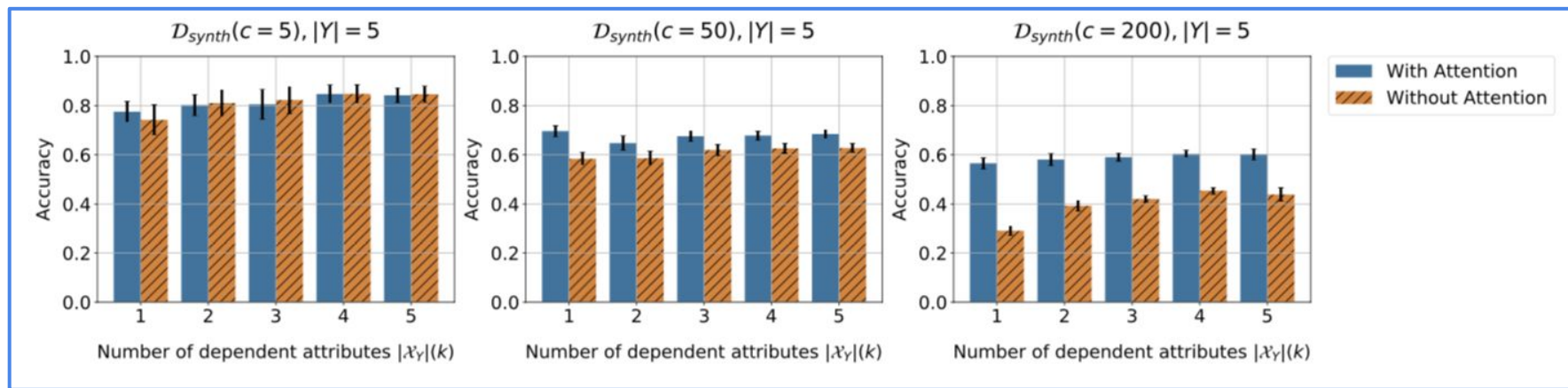
AimNet still outperforms baselines in almost all cases

data set	Attribute	AimNet	Accuracy on discrete attributes (ACC $\pm$ std)					
			HCQ	XGB	MIDAS	GAIN	MF	MICE
Balance	class	<b><math>0.83 \pm 0.07</math></b>	$0.4 \pm 0.37$	$0.48 \pm 0.32$	$0.7 \pm 0.16$	$0.5 \pm 0.12$	$0.46 \pm 0.34$	<b><math>0.78 \pm 0.15</math></b>
	ADDR_PCT_CD	<b><math>0.67 \pm 0.03</math></b>	$0.19 \pm 0.05$	$0.41 \pm 0.05$	$0.13 \pm 0.01$	$0.04 \pm 0.04$	$0.59 \pm 0.03$	$0.23 \pm 0.01$
NYPD	BORO_NM	<b><math>0.92 \pm 0.07</math></b>	<b><math>0.85 \pm 0.09</math></b>	$0.58 \pm 0.18$	$0.78 \pm 0.04$	$0.23 \pm 0.03$	$0.84 \pm 0.09$	$0.58 \pm 0.09$
	PATROL_BORO	<b><math>0.83 \pm 0.07</math></b>	$0.69 \pm 0.03$	$0.57 \pm 0.17$	$0.6 \pm 0.06$	$0.13 \pm 0.02$	$0.72 \pm 0.1$	$0.58 \pm 0.07$
data set	Attribute	AimNet	NRMS on continuous attributes (NRMS $\pm$ std)					
			HCQ	XGB	MIDAS	GAIN	MF	MICE
Phase	A	<b><math>0.75 \pm 0.13</math></b>	$1.26 \pm 0.46$	<b><math>0.81 \pm 0.18</math></b>	$1.45 \pm 0.5$	$2.04 \pm 1.73$	$0.94 \pm 0.18$	$1.41 \pm 0.66$
	B	<b><math>0.77 \pm 0.07</math></b>	$0.99 \pm 0.27$	<b><math>0.83 \pm 0.16</math></b>	$1.18 \pm 0.33$	$1.54 \pm 0.54$	$0.98 \pm 0.13$	$1.32 \pm 0.4$
	C	<b><math>0.79 \pm 0.11</math></b>	$1.25 \pm 0.26$	<b><math>0.81 \pm 0.13</math></b>	$1.44 \pm 0.25$	$1.33 \pm 0.34$	$1.09 \pm 0.23$	$1.29 \pm 0.1$
	D	<b><math>0.62 \pm 0.09</math></b>	$1.1 \pm 0.24$	<b><math>0.65 \pm 0.09</math></b>	$1.47 \pm 0.56$	$1.57 \pm 0.63$	$0.85 \pm 0.17$	$1.03 \pm 0.13$
	Trip Total	$0.82 \pm 0.21$	$2.38 \pm 1.16$	<b><math>0.48 \pm 0.32</math></b>	$3.84 \pm 0.7$	$1.87 \pm 0.7$	<b><math>0.71 \pm 0.57</math></b>	—
Chicago	Fare	<b><math>1.35 \pm 0.76</math></b>	$5.18 \pm 1.94$	<b><math>1.16 \pm 0.73</math></b>	$12.73 \pm 7.93$	$58.09 \pm 36.03$	$2.78 \pm 3.0$	—
	Tips	<b><math>0.52 \pm 0.01</math></b>	$1.29 \pm 0.09$	<b><math>0.5 \pm 0.12</math></b>	$1.59 \pm 0.49$	$11.64 \pm 6.76$	$0.8 \pm 0.28$	—

# Does the attention layer actually help?

As the domain size increases, attention leads to better performance

- Learns schema-level dependencies



5 classes

50 classes

200 classes

# Architecture summary

- **Encode:** learns projections for **continuous** and embeddings for **discrete** data
- **Structure:** new variation of attention to learn structural dependencies between **attributes**
- **Prediction:** mixed-type prediction using projections (continuous) and softmax classification (discrete)

# Conclusion

- A **simple attention-based architecture** modestly outperforms existing methods on i.i.d. missing values
- AimNet outperforms state of the art in the presence of **systematically** missing values by a **large margin**
- **Attention** mechanism learns **structural properties of the data** which improves MVI with systematic bias

# Appendix

# Hyperparameter Sensitivity

Accuracy on discrete attributes for the NYPD data set				
	dropout	max domain size	embedding size	AimNet
base	0.25	50	64	0.921
(dropout rate)	0.0			0.918
	0.5			0.920
(max domain size)		10		0.917
		100		0.921
(embedding size)			16	0.920
			32	0.921
			128	<b>0.922</b>
			256	<b>0.920</b>
NRMS on continuous attributes for the NYPD data set				
	dropout rate	max domain size	embedding size	AimNet
base	0.0	50	64	<b>0.150</b>
(dropout rate)	0.25			0.281
	0.5			0.509
(embedding size)			16	0.159
			32	<b>0.150</b>
			128	0.153
			256	<b>0.144</b>



# Multi-task and Single-task

Data Set	Accuracy on discrete attributes (ACC $\pm$ std)	
	Single	Multi-Task
Tic-Tac-Toc	0.61 $\pm$ 0.01	0.61 $\pm$ 0.01
Hospital	0.99 $\pm$ 0.0	0.99 $\pm$ 0.0
Mammogram	0.75 $\pm$ 0.01	0.75 $\pm$ 0.01
Thoracic	0.86 $\pm$ 0.01	0.86 $\pm$ 0.01
Contraceptive	0.65 $\pm$ 0.01	0.65 $\pm$ 0.01
Solar Flare	0.78 $\pm$ 0.02	0.78 $\pm$ 0.01
NYPD	0.92 $\pm$ 0.0	0.92 $\pm$ 0.0
Credit	0.76 $\pm$ 0.01	0.76 $\pm$ 0.01
Australian	0.72 $\pm$ 0.02	0.72 $\pm$ 0.01
Chicago	0.73 $\pm$ 0.01	0.7 $\pm$ 0.03
Balance	0.79 $\pm$ 0.04	0.79 $\pm$ 0.04
Eye EEG	0.71 $\pm$ 0.01	0.69 $\pm$ 0.01
Data Set	NRMS on continuous attributes (NRMS $\pm$ std)	
	Single	Multi-Task
Hospital	0.72 $\pm$ 0.06	0.77 $\pm$ 0.06
Mammogram	0.91 $\pm$ 0.04	0.93 $\pm$ 0.03
Thoracic	1.1 $\pm$ 0.41	1.3 $\pm$ 0.83
Contraceptive	0.84 $\pm$ 0.02	0.84 $\pm$ 0.02
Solar Flare	0.94 $\pm$ 0.15	0.87 $\pm$ 0.16
NYPD	0.15 $\pm$ 0.01	0.15 $\pm$ 0.01
Credit	0.94 $\pm$ 0.03	0.94 $\pm$ 0.03
Australian	0.94 $\pm$ 0.03	0.93 $\pm$ 0.02
Eye EEG	0.4 $\pm$ 0.0	0.44 $\pm$ 0.0
Phase	0.45 $\pm$ 0.01	0.45 $\pm$ 0.01
CASP	0.45 $\pm$ 0.02	0.48 $\pm$ 0.01

Data Set	Run time (seconds)	
	Single	Multi-Task
Tic-Tac-Toc	9	38
Hospital	18	148
Mammogram	9	39
Thoracic	11	69
Contraceptive	15	102
Solar Flare	15	131
NYPD	378	6320
Credit	15	198
Australian	16	145
Chicago	3180	462756
Balance	11	12
Eye EEG	188	5306
Phase	66	250
CASP	648	5800

# MCAR (40% missing) results

Data Set	Accuracy on discrete attributes (ACC $\pm$ std)						
	AimNet	HCQ	XGB	MIDAS	GAIN	MF	MICE
Tic-Tac-Toe	<b>0.53 <math>\pm</math> 0.01</b>	0.5 $\pm$ 0.01	0.52 $\pm$ 0.01	0.44 $\pm$ 0.02	0.35 $\pm$ 0.01	0.5 $\pm$ 0.01	0.46 $\pm$ 0.01
Hospital	<b>0.95 <math>\pm</math> 0.0</b>	<b>0.95 <math>\pm</math> 0.0</b>	0.91 $\pm$ 0.01	0.24 $\pm$ 0.01	0.14 $\pm$ 0.02	0.94 $\pm$ 0.01	0.7 $\pm$ 0.01
Mammogram	<b>0.73 <math>\pm</math> 0.02</b>	<b>0.72 <math>\pm</math> 0.01</b>	<b>0.72 <math>\pm</math> 0.02</b>	0.71 $\pm$ 0.02	0.35 $\pm$ 0.01	0.66 $\pm$ 0.02	0.63 $\pm$ 0.02
Thoracic	<b>0.85 <math>\pm</math> 0.01</b>	0.84 $\pm$ 0.01	0.84 $\pm$ 0.01	0.83 $\pm$ 0.02	0.52 $\pm$ 0.15	<b>0.85 <math>\pm</math> 0.01</b>	0.75 $\pm$ 0.03
Contraceptive	<b>0.63 <math>\pm</math> 0.01</b>	<b>0.63 <math>\pm</math> 0.01</b>	0.62 $\pm$ 0.01	0.62 $\pm$ 0.01	0.43 $\pm$ 0.01	0.62 $\pm$ 0.01	0.55 $\pm$ 0.01
Solar Flare	<b>0.76 <math>\pm</math> 0.01</b>	0.75 $\pm$ 0.01	0.75 $\pm$ 0.01	0.66 $\pm$ 0.01	0.46 $\pm$ 0.02	0.74 $\pm$ 0.01	0.65 $\pm$ 0.01
NYPD	0.87 $\pm$ 0.0	0.85 $\pm$ 0.0	<b>0.88 <math>\pm</math> 0.0</b>	0.75 $\pm$ 0.0	0.15 $\pm$ 0.01	<b>0.88 <math>\pm</math> 0.0</b>	0.58 $\pm$ 0.0
Credit	<b>0.73 <math>\pm</math> 0.01</b>	0.7 $\pm$ 0.01	<b>0.73 <math>\pm</math> 0.01</b>	0.6 $\pm$ 0.01	0.39 $\pm$ 0.01	<b>0.73 <math>\pm</math> 0.01</b>	0.63 $\pm$ 0.01
Australian	<b>0.7 <math>\pm</math> 0.01</b>	0.66 $\pm$ 0.01	0.68 $\pm$ 0.01	0.6 $\pm$ 0.01	0.46 $\pm$ 0.01	0.69 $\pm$ 0.01	0.59 $\pm$ 0.01
Balance	<b>0.73 <math>\pm</math> 0.03</b>	<b>0.72 <math>\pm</math> 0.03</b>	<b>0.71 <math>\pm</math> 0.03</b>	0.64 $\pm$ 0.03	0.45 $\pm$ 0.05	0.63 $\pm$ 0.05	0.64 $\pm$ 0.04
Eye EEG	0.67 $\pm$ 0.01	0.62 $\pm$ 0.01	0.73 $\pm$ 0.01	0.55 $\pm$ 0.01	0.52 $\pm$ 0.03	<b>0.78 <math>\pm</math> 0.01</b>	0.53 $\pm$ 0.01
Data Set	NRMS on continuous attributes (NRMS $\pm$ std)						
	AimNet	HCQ	XGB	MIDAS	GAIN	MF	MICE
Hospital	<b>0.81 <math>\pm</math> 0.04</b>	1.1 $\pm$ 0.08	0.92 $\pm$ 0.07	440.63 $\pm$ 61.35	2.26 $\pm$ 1.18	0.89 $\pm$ 0.04	1.23 $\pm$ 0.09
Mammogram	<b>0.92 <math>\pm</math> 0.02</b>	1.02 $\pm$ 0.04	0.98 $\pm$ 0.05	1.12 $\pm$ 0.08	1.05 $\pm$ 0.06	1.01 $\pm$ 0.03	1.25 $\pm$ 0.07
Thoracic	<b>0.94 <math>\pm</math> 0.01</b>	1.09 $\pm$ 0.05	1.03 $\pm$ 0.11	5.64 $\pm$ 7.16	1.23 $\pm$ 0.22	0.99 $\pm$ 0.06	1.32 $\pm$ 0.12
Contraceptive	<b>0.9 <math>\pm</math> 0.02</b>	1.12 $\pm$ 0.04	0.94 $\pm$ 0.02	1.11 $\pm$ 0.02	1.17 $\pm$ 0.05	0.99 $\pm$ 0.02	1.23 $\pm$ 0.06
Solar Flare	<b>0.93 <math>\pm</math> 0.09</b>	<b>0.98 <math>\pm</math> 0.09</b>	<b>1.0 <math>\pm</math> 0.1</b>	10772.7 $\pm$ 4057.37	<b>1.0 <math>\pm</math> 0.09</b>	1.04 $\pm$ 0.11	1.16 $\pm$ 0.07
NYPD	0.32 $\pm$ 0.01	0.44 $\pm$ 0.13	0.28 $\pm$ 0.0	0.69 $\pm$ 0.03	3.63 $\pm$ 0.17	<b>0.22 <math>\pm</math> 0.01</b>	0.62 $\pm$ 0.01
Credit	<b>0.97 <math>\pm</math> 0.03</b>	1.24 $\pm$ 0.03	1.26 $\pm$ 0.41	1.15 $\pm$ 0.07	1.2 $\pm$ 0.08	1.12 $\pm$ 0.18	1.34 $\pm$ 0.11
Australian	<b>0.96 <math>\pm</math> 0.02</b>	1.23 $\pm$ 0.03	1.19 $\pm$ 0.2	1.14 $\pm$ 0.12	1.27 $\pm$ 0.16	1.07 $\pm$ 0.13	1.6 $\pm$ 0.7
Eye EEG	0.48 $\pm$ 0.0	0.71 $\pm$ 0.03	0.47 $\pm$ 0.0	0.91 $\pm$ 0.01	1.0 $\pm$ 0.28	<b>0.44 <math>\pm</math> 0.0</b>	0.67 $\pm$ 0.01
Phase	<b>0.52 <math>\pm</math> 0.01</b>	0.58 $\pm$ 0.0	0.53 $\pm$ 0.01	0.97 $\pm$ 0.01	1.14 $\pm$ 0.26	0.58 $\pm$ 0.01	0.73 $\pm$ 0.01
CASP	0.5 $\pm$ 0.01	1.5 $\pm$ 0.26	<b>0.49 <math>\pm</math> 0.01</b>	0.88 $\pm$ 0.01	0.83 $\pm$ 0.09	<b>0.48 <math>\pm</math> 0.01</b>	0.73 $\pm$ 0.03

# MCAR (60% missing) results

Data Set	Accuracy on discrete attributes (ACC $\pm$ std)						
	AimNet	HCQ	XGB	MIDAS	GAIN	MF	MICE
Tic-Tac-Toe	<b>0.48 <math>\pm</math> 0.01</b>	<b>0.48 <math>\pm</math> 0.0</b>	0.47 $\pm$ 0.01	0.43 $\pm$ 0.01	0.39 $\pm$ 0.01	0.44 $\pm$ 0.01	0.4 $\pm$ 0.01
Hospital	<b>0.86 <math>\pm</math> 0.01</b>	<b>0.86 <math>\pm</math> 0.01</b>	0.68 $\pm$ 0.01	0.24 $\pm$ 0.0	0.11 $\pm$ 0.01	0.79 $\pm$ 0.01	0.37 $\pm$ 0.01
Mammogram	<b>0.69 <math>\pm</math> 0.01</b>	<b>0.69 <math>\pm</math> 0.01</b>	<b>0.68 <math>\pm</math> 0.01</b>	<b>0.68 <math>\pm</math> 0.01</b>	0.34 $\pm$ 0.02	0.62 $\pm$ 0.02	0.58 $\pm$ 0.02
Thoracic	<b>0.85 <math>\pm</math> 0.01</b>	0.84 $\pm$ 0.01	0.83 $\pm$ 0.0	0.84 $\pm$ 0.01	0.51 $\pm$ 0.13	0.84 $\pm$ 0.01	0.72 $\pm$ 0.04
Contraceptive	<b>0.62 <math>\pm</math> 0.01</b>	<b>0.62 <math>\pm</math> 0.01</b>	0.61 $\pm$ 0.01	0.61 $\pm$ 0.01	0.42 $\pm$ 0.02	0.6 $\pm$ 0.01	0.53 $\pm$ 0.01
Solar Flare	<b>0.72 <math>\pm</math> 0.01</b>	<b>0.72 <math>\pm</math> 0.01</b>	0.71 $\pm$ 0.01	0.66 $\pm$ 0.01	0.45 $\pm$ 0.03	0.7 $\pm$ 0.01	0.61 $\pm$ 0.01
NYPD	0.77 $\pm$ 0.0	0.76 $\pm$ 0.0	<b>0.79 <math>\pm</math> 0.0</b>	0.67 $\pm$ 0.0	0.15 $\pm$ 0.0	0.78 $\pm$ 0.0	0.45 $\pm$ 0.0
Credit	<b>0.69 <math>\pm</math> 0.01</b>	0.66 $\pm$ 0.01	0.68 $\pm$ 0.01	0.6 $\pm$ 0.01	0.38 $\pm$ 0.02	0.68 $\pm$ 0.01	0.57 $\pm$ 0.01
Australian	<b>0.67 <math>\pm</math> 0.01</b>	0.65 $\pm$ 0.01	0.65 $\pm$ 0.01	0.59 $\pm$ 0.01	0.45 $\pm$ 0.02	0.65 $\pm$ 0.01	0.54 $\pm$ 0.02
Balance	<b>0.67 <math>\pm</math> 0.03</b>	<b>0.64 <math>\pm</math> 0.04</b>	0.63 $\pm$ 0.03	0.51 $\pm$ 0.06	0.46 $\pm$ 0.04	0.54 $\pm$ 0.03	0.55 $\pm$ 0.04
Eye EEG	0.63 $\pm$ 0.01	0.6 $\pm$ 0.01	0.66 $\pm$ 0.01	0.54 $\pm$ 0.01	0.52 $\pm$ 0.03	<b>0.67 <math>\pm</math> 0.01</b>	0.52 $\pm$ 0.01
Data Set	NRMS on continuous attributes (NRMS $\pm$ std)						
	AimNet	HCQ	XGB	MIDAS	GAIN	MF	MICE
Hospital	<b>0.9 <math>\pm</math> 0.04</b>	1.16 $\pm$ 0.12	1.01 $\pm$ 0.08	139.97 $\pm$ 24.15	3.79 $\pm$ 0.36	0.95 $\pm$ 0.05	1.27 $\pm$ 0.09
Mammogram	<b>0.94 <math>\pm</math> 0.02</b>	1.05 $\pm$ 0.06	1.0 $\pm$ 0.05	1.13 $\pm$ 0.06	1.1 $\pm$ 0.11	1.01 $\pm$ 0.04	1.28 $\pm$ 0.08
Thoracic	<b>0.99 <math>\pm</math> 0.02</b>	1.13 $\pm$ 0.05	1.17 $\pm$ 0.11	3.54 $\pm$ 5.28	1.18 $\pm$ 0.09	1.07 $\pm$ 0.04	1.43 $\pm$ 0.2
Contraceptive	<b>0.94 <math>\pm</math> 0.01</b>	1.15 $\pm$ 0.04	1.01 $\pm$ 0.02	1.12 $\pm$ 0.02	1.31 $\pm$ 0.14	1.14 $\pm$ 0.04	1.29 $\pm$ 0.05
Solar Flare	<b>0.98 <math>\pm</math> 0.06</b>	<b>1.01 <math>\pm</math> 0.06</b>	1.05 $\pm$ 0.09	4454.64 $\pm$ 1610.59	1.07 $\pm$ 0.17	1.13 $\pm$ 0.14	1.24 $\pm$ 0.19
NYPD	0.56 $\pm$ 0.0	0.58 $\pm$ 0.04	0.45 $\pm$ 0.0	0.8 $\pm$ 0.01	3.51 $\pm$ 0.14	<b>0.42 <math>\pm</math> 0.01</b>	0.98 $\pm$ 0.0
Credit	<b>0.99 <math>\pm</math> 0.01</b>	1.33 $\pm$ 0.19	1.23 $\pm$ 0.25	1.14 $\pm$ 0.11	1.24 $\pm$ 0.11	1.13 $\pm$ 0.13	1.43 $\pm$ 0.26
Australian	<b>0.98 <math>\pm</math> 0.01</b>	1.37 $\pm$ 0.26	1.29 $\pm$ 0.42	1.1 $\pm$ 0.04	1.25 $\pm$ 0.12	1.13 $\pm$ 0.19	1.38 $\pm$ 0.09
Eye EEG	0.59 $\pm$ 0.0	0.82 $\pm$ 0.04	<b>0.57 <math>\pm</math> 0.0</b>	0.97 $\pm$ 0.01	1.61 $\pm$ 0.24	<b>0.57 <math>\pm</math> 0.0</b>	0.79 $\pm$ 0.0
Phase	<b>0.64 <math>\pm</math> 0.0</b>	0.71 $\pm$ 0.01	0.65 $\pm$ 0.0	1.0 $\pm$ 0.0	1.45 $\pm$ 0.51	0.71 $\pm$ 0.01	0.91 $\pm$ 0.01
CASP	<b>0.58 <math>\pm</math> 0.01</b>	2.06 $\pm$ 0.51	0.59 $\pm$ 0.01	0.94 $\pm$ 0.01	1.2 $\pm$ 0.22	0.62 $\pm$ 0.0	0.88 $\pm$ 0.03

# Census Tracts form Voronoi-like cells

