# Adaptive Data Augmentation for Supervised Learning over Missing Data

#### Tongyu Liu, Ju Fan, Yinqing Luo, Xiaoyong Du



#### Nan Tang



#### **Guoliang Li**



#### Missing Data is Everywhere

- A A A A
- The raw data we collect from the real world always contain missing values.



- How to handle the missing values?
  - Dropping the records containing missingness.
  - Using Imputation methods to complete data.
    - It is hard to obtain the ground-truth via value imputation

## Noise Shift in Source and Target Datasets

 Separately impute source and target data might cause an even bigger divergence on data distributions.



If we cannot repair them to be correct, can we repair them to be "similar"?

#### Adaptive Data Augmentation



#### Source Data

age	cholesterol	glucose	smoke	alcohol	cardio
35	2	NA	1	0	no
50	3	3	0	0	yes
65	2	3	1	1	yes

<sup>(2)</sup> Source Data Adaptation

Target Mask Generation Target Data (1)

> Target Pattern

#### ③ Simple Imputation

age	cholesterol	glucose	smoke	alcohol	cardio
25	1	NA	0	1	?
37	NA	3	0	0	?
70	3	2	1	1	?

age	cholesterol	glucose	smoke	alcohol	cardio
35	2	2	1	0	no
50	3	3	0	0	yes
65	2	3	1	1	yes

#### (4) Retrain the model





age	cholesterol	glucose	smoke	alcohol	cardio
25	1	2	0	1	?
37	2	3	0	0	?
70	3	2	1	1	?



#### **Target Mask Generation**



We utilize Conditional GAN to learn the p(mask|observed data)  $\mathcal{L}_m(D_m, G_m) = \mathbb{E}_{(\mathbf{x}_t, \mathbf{m}_t) \sim \mathcal{U}_t}[D_m(\mathbf{m}_t, \mathbf{x}_t)]$ 

 $-\mathbb{E}_{\boldsymbol{\sigma}\sim p(\boldsymbol{\sigma}),(\boldsymbol{x}_{t},\boldsymbol{m}_{t})\sim\mathcal{U}_{t}}[D_{m}(G_{m}(\boldsymbol{\sigma},\boldsymbol{x}_{t}),\boldsymbol{x}_{t})]$ 

Target Data

age	cholesterol	glucose	smoke	alcohol
25	1	NA	0	1
37	NA	3	0	0
70	3	2	1	1



Mask Matrix

age	cholesterol	glucose	smoke	alcohol
1	1	0	1	1
1	0	1	1	1
1	1	1	1	1

#### Source Data Adaptation





#### Adaptive Data Augmentation



Finally, we can transform a source data to the target data by using the mask generator  $G_m$  and data generator  $G_x$ .

8

### Experiment

• Datasets

	Dataset	Area	#Rec	<b>#N</b>	#C	#L	Label Distribution
ſ	Ipums	Social	16329	16	43	7	65:9:11:5:6:3:1
Real Missing	Okcupid	Social	50789	3	13	3	7:2:1
l	Welfare	Financial	20309	5	0	2	1.0 : 1.2
Synthetic ∫	EyeState	Health	14977	14	0	2	1.0 : 1.3
Missing ]	Adult	Social	13567	6	8	2	1.0 : 3.0

- Baselines:
  - MICE: Multiple Imputation with Chained Equations.
  - MISF: MissForest.
  - GAIN: Generative Adversarial Imputation Nets.



#### 9

#### Experiment

- Evaluation Framework
  - Evaluation Task:
    - Classification.
  - Metrics:
    - F1 score of model.
  - For Imputation Methods:
    - Train the model on **imputed source data**, then apply it to the **imputed target data**.
  - For DAGAN:
    - Train the model on **adapted source data**, then directly apply it to the **target data**.



## Evaluation on Real-World Datasets.





\* MISF fails on the Ipums dataset because its R-based implementation cannot handle categorical attributes with more than 53 categories

# **Evaluating Effect of Data Adaptation**





Repairing source and target to be "similar" is beneficial for ML training.

Adult Dataset

# **Robustness Evaluation**

- A A A A
- We consider three different settings of missing value injection.
  - MCAR (Missing Completely at Random)
  - MAR (Missing at Random)
  - MNAR (Missing Not at Random)



DAGAN is robust for different missing data patterns.

# Conclusion

- We propose DAGAN to adapt the source data to better serve the prediction on the unseen target data.
- DAGAN performs well over different missing patterns.

Missing	Missing Dataset-Setting		MISF	MICE	DAGAN
	EyeState-Overlap	-	+		++
MNAR	Adult-Overlap	+	-		++
WINAR	EyeState-NoOverlap	-		+	++
	Adult-NoOverlap	-		+	++
	EyeState-Overlap		+	-	++
MAR	Adult-Overlap		-	+	++
MIAR	EyeState-NoOverlap	-		+	++
	Adult-NoOverlap		-	++	+
	EyeState-Overlap		++	-	+
MCAR	Adult-Overlap	-		+	++
MCAK	EyeState-NoOverlap		++	+	-
	Adult-NoOverlap		++	-	+

Average F1 across missing rates under different datasets and missing patterns



# Thank you!

Github repository: https://github.com/ruc-datalab/dagan

#### Multi-ADA





15