

# Tabular Data Prediction with Heterogeneous Features

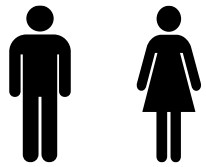
Hedong Yan, CS, HKBU

Supervisor: Prof. Yiuming Cheung

## 2

# Motivation

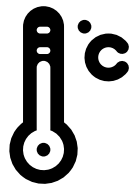
- Fact 1: Tabular data is often heterogenous.
- Fact 2: Neural network often NOT perform well on tabular data.



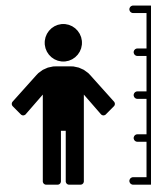
Nominal



Ordinal



Interval



Ratio

**Heterogeneous Features**

---

## Why do tree-based models still outperform deep learning on typical tabular data?

---

Léo Grinsztajn  
Soda, Inria Saclay  
leo.grinsztajn@inria.fr

Edouard Oyallon  
MLIA, Sorbonne University

Gaël Varoquaux  
Soda, Inria Saclay

### Abstract

While deep learning has enabled tremendous progress on text and image datasets, its superiority on tabular data is not clear. We contribute extensive benchmarks of standard and novel deep learning methods as well as tree-based models such as XGBoost and Random Forests, across a large number of datasets and hyperparameter combinations. We define a standard set of 45 datasets from varied domains with clear characteristics of tabular data and a benchmarking methodology accounting for both fitting models and finding good hyperparameters. Results show that tree-based models remain state-of-the-art on medium-sized data (~10K samples) even without accounting for their superior speed. To understand this gap, we conduct an empirical investigation into the differing inductive biases of tree-based models and neural networks. This leads to a series of challenges which should guide researchers aiming to build tabular-specific neural network: 1. be robust to uninformative features, 2. preserve the orientation of the data, and 3. be able to easily learn irregular functions. To stimulate research on tabular architectures, we contribute a standard benchmark and raw data for baselines: every point of a 20 000 compute hours hyperparameter search for each learner.

**Deep learning does NOT perform well**

**Problem: How to address on heterogeneous features effectively?**

# Related works (unsupervised)

One-hot

SEX	Male	Female
Male	1	0
Female	0	1

Ordinal

RATE	RATE
Bad	0
Neutral	1
Good	2

Rank-hot

RATE	Bad	Neutral	Good
Bad	1	0	0
Neutral	1	1	0
Good	1	1	1

Piece-wise linear

$$\begin{array}{c}
 \begin{array}{ccccccc}
 & & x & & & & \\
 | & | & | & | & | & | & \\
 b_0 & b_1 & b_2 & b_3 & b_4 & & \mathbb{R}
 \end{array} \\
 \Downarrow \\
 \text{PLE}(x) = \begin{array}{|c|c|c|c|}
 \hline
 1 & 1 & \frac{x - b_2}{b_3 - b_2} & 0 \\
 \hline
 e_1 & e_2 & e_3 & e_4
 \end{array}
 \end{array}$$

Figure 1. The piecewise linear encoding (PLE) in action, as defined in Equation 4. In the example,  $T = 4$ .

# Related works (supervised)

## Periodic

$$\begin{aligned} f_i(x) &= \text{Periodic}(x) = \text{concat}[\sin(v), \cos(v)], \\ v &= [2\pi c_1 x, \dots, 2\pi c_k x] \end{aligned} \quad (8)$$

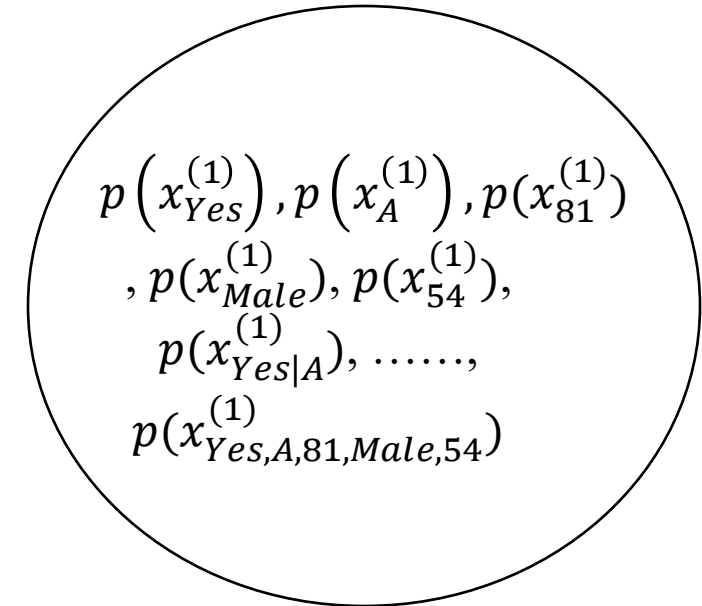
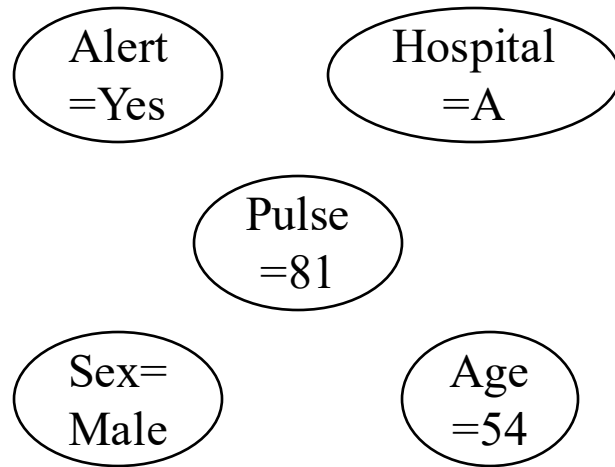
where  $c_i$  are trainable parameters initialized from  $\mathcal{N}(0, \sigma)$ .  $\sigma$  is an important hyperparameter that is tuned using validation sets.

## Target Statistic

$$x_k^i = \frac{\sum_{j=1}^n \mathbb{I}_{x_j^i = x_k^i} * y_j + ap}{\sum_{j=1}^n \mathbb{I}_{x_j^i = x_k^i} + a} \quad (1)$$

# Methodology

Idea



**Instance in Heterogeneous Features Space H**

**Instance in Measurement Occurrence Space S**

**Feature can be very heterogeneous. We list all potential measurements of an individual and use the measurement probabilities as new feature space where each axis is a potential measurement.**

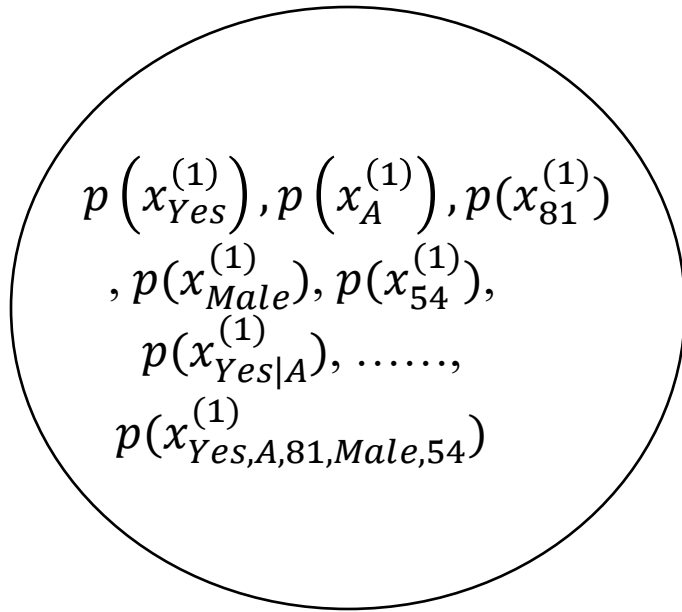
# Methodology

**Challenge 1: the dimensions of Space S is combinational.**

**Solution: We use subspaces S1,..., Sm that were separated by features and concatenate them together.**

**Challenge 2: the number of probabilities is combinational.**

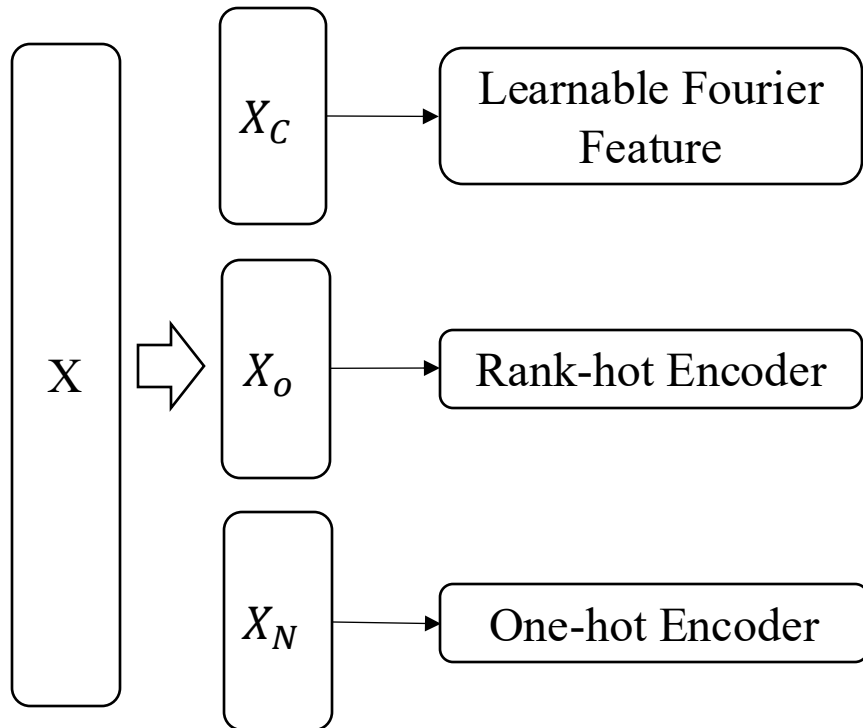
**Solution: We use zero-order probabilities  $p(x_{Yes}^{(1)}), p(x_A^{(1)}), p(x_{81}^{(1)}), p(x_{Male}^{(1)}), p(x_{54}^{(1)})$  as approximation.**



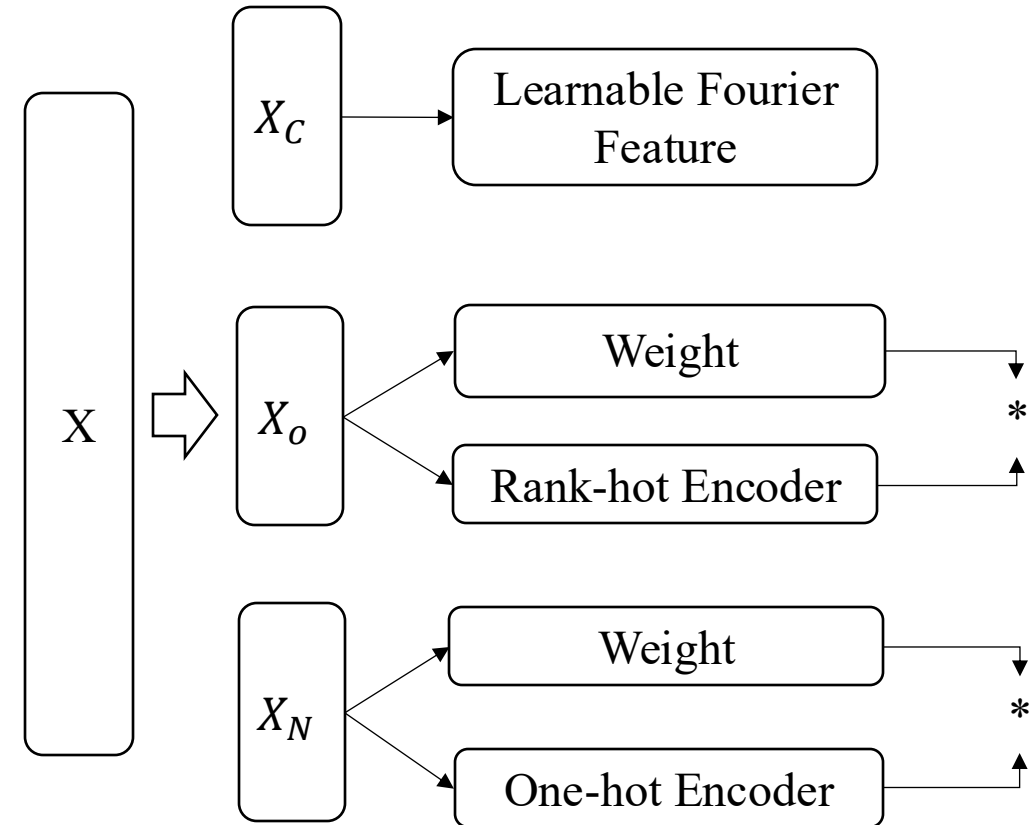
**Measurement Occurrence Space S**

# Methodology

## Encoder of A Trivial MLP



## Encoder of HetMLP

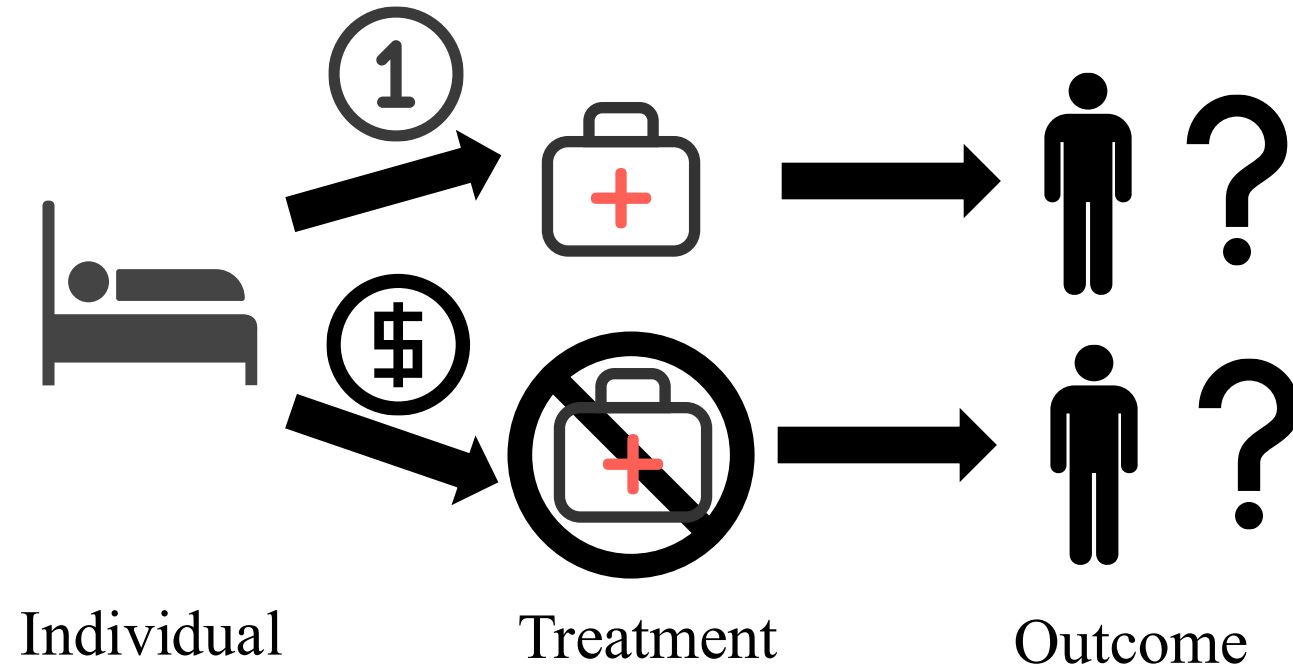


**Finally, we get a weighted encoder for heterogeneous features where weight is occurrence's inverse probability. It can better deal with rare event.**



# Experiment

## Outcome Prediction Task For RCT



### ➤ Task Goal

Predict Outcome in Randomized Control Trial

### ➤ Meaning

We can compare the predicted outcome difference between different treatment for an individual to decide whether a patient should accept the treatment.

### ➤ Metric

Mean Average Precision for Classification  
Mean Squared Error for Regression

**Assumption:** for an individual, if a unbiased model can predict its factual outcome in RCT better, then it can predict its counterfactual outcome in RCT better.

# Experiment

## Our Gathered Dataset

TABLE IV: Heterogeneous datasets for outcome prediction

Dataset	Instance	Outcome	Treatment
Safety and Preliminary Efficacy of Intranasal Insulin for Cognitive Impairment in Parkinson Disease and Multiple System Atrophy	16	Parkinson disease	Intranasal insulin
<a href="https://physionet.org/content/inipdmsa/1.0/">https://physionet.org/content/inipdmsa/1.0/</a>			
Tai Chi, Physiological Complexity, and Healthy Aging - Gait	60	Gait and EMG data	Tai Chi
<a href="https://physionet.org/content/taichidb/1.0.2/">https://physionet.org/content/taichidb/1.0.2/</a>			
ECG Effects of Dofetilide, Moxifloxacin, Dofetilide+Mexiletine, Dofetilide+Lidocaine and Moxifloxacin+Diltiazem	22	ECG	Dofetilide, Moxifloxacin, Dofetilide+Mexiletine, Dofetilide+Lidocaine and Moxifloxacin+Diltiazem
<a href="https://physionet.org/content/ecgdmml/1.0.0/">https://physionet.org/content/ecgdmml/1.0.0/</a>			
ECG Effects of Ranolazine, Dofetilide, Verapamil, and Quinidine	22	ECG	Ranolazine, Dofetilide, Verapamil, and Quinidine
<a href="https://physionet.org/content/ecgrdvq/1.0.0/">https://physionet.org/content/ecgrdvq/1.0.0/</a>			
CAST RR Interval Sub-Study Database	734	Cardiac arrhythmia suppression	Encainide, flecainide, moricizine (antiarrhythmic drugs) or a placebo
<a href="https://physionet.org/content/crisdb/1.0.0/">https://physionet.org/content/crisdb/1.0.0/</a>			
Randomized trial of AKI alerts in hospitalized patients	6030	Acute Kidney Injury	Electronic AKI alert versus usual care
<a href="https://datadryad.org/stash/dataset/doi:10.5061%2Fdryad.59zw3r27n">https://datadryad.org/stash/dataset/doi:10.5061%2Fdryad.59zw3r27n</a>			
Telerehabilitation program for COVID-19 survivors (TERECO) - Randomized controlled trial	120	Exercise capacity, lower-limb muscle strength (LMS), pulmonary function, health-related quality of life (HRQOL), and dyspnoea	Telerehabilitation program for COVID-19 survivors
<a href="https://datadryad.org/stash/dataset/doi:10.5061%2Fdryad.59zw3r27n">https://datadryad.org/stash/dataset/doi:10.5061%2Fdryad.59zw3r27n</a>			
Bicycling comfort video experiment	15289	Bicycle rating	Video Type
<a href="https://datadryad.org/stash/dataset/doi:10.25338%2FB8KG77">https://datadryad.org/stash/dataset/doi:10.25338%2FB8KG77</a>			
Megafon uplift competition	1.5 million	User conversion	Exposure
<a href="https://ods.ai/tracks/df21-megafon/competitions/megafon-df21-comp/data">https://ods.ai/tracks/df21-megafon/competitions/megafon-df21-comp/data</a>			
Infant Health and Development Program	1090	Cognitive development, Behavior problems, Health status	Home visits, attendance at a special child development center
<a href="https://www.icpsr.umich.edu/web/HMCA/studies/9795">https://www.icpsr.umich.edu/web/HMCA/studies/9795</a>			
National Supported Work Evaluation Study	6600	effects of the Supported Work Program	Offered a job in supported work
<a href="https://www.icpsr.umich.edu/web/ICPSR/studies/7865">https://www.icpsr.umich.edu/web/ICPSR/studies/7865</a>			
CPAP Pressure and Flow Data from a Local Trial of 30 Adults at the University of Canterbury	30	Breathing	Continuous positive airway pressure
<a href="https://physionet.org/content/cpap-data-canterbury/1.0.1/">https://physionet.org/content/cpap-data-canterbury/1.0.1/</a>			

- Most of those tabular datasets are heterogeneous features.
- More detail can be seen at: <https://github.com/herdonyan/RandomizedTrialDataset>

# Experiment

## Alert2AKI Dataset

<b>Intervention</b>	AKI Alert or Not
<b>Main-outcome</b>	AKI Progression in 14 Days
<b>Pre-treatment</b>	EHR Records
<b>Patients Num</b>	6030 in 5 Hospitals (5082/948)

SCALE	NUM
Nominal	9
Ordinal	19
Interval	3
Ratio	20

## PR-AUC (5 Random Splits)

HetMLP	Trivial MLP	MLP	Random
<b>.2117<math>\pm</math>.0009</b>	.2087 $\pm$ .0164	.2009 $\pm$ .0329	.1568 $\pm$ .0089

Our HetMLP got **1.43%** performance up compared with Trivial MLP.

# Will the patients be benefited from the alert?

## Splitting 1

Metrics	Num
Patients Num	3536
Benefited: AKI=1 $\rightarrow$ AKI=0	15
Harmful: AKI=0 $\rightarrow$ AKI=1	14

## Splitting 2

Metrics	Num
Patients Num	3552
Benefited: AKI=1 $\rightarrow$ AKI=0	8
Harmful: AKI=0 $\rightarrow$ AKI=1	2

## Splitting 3

Metrics	Num
Patients Num	3504
Benefited: AKI=1 $\rightarrow$ AKI=0	26
Harmful: AKI=0 $\rightarrow$ AKI=1	4

## Splitting 4

Metrics	Num
Patients Num	3536
Benefited: AKI=1 $\rightarrow$ AKI=0	9
Harmful: AKI=0 $\rightarrow$ AKI=1	9

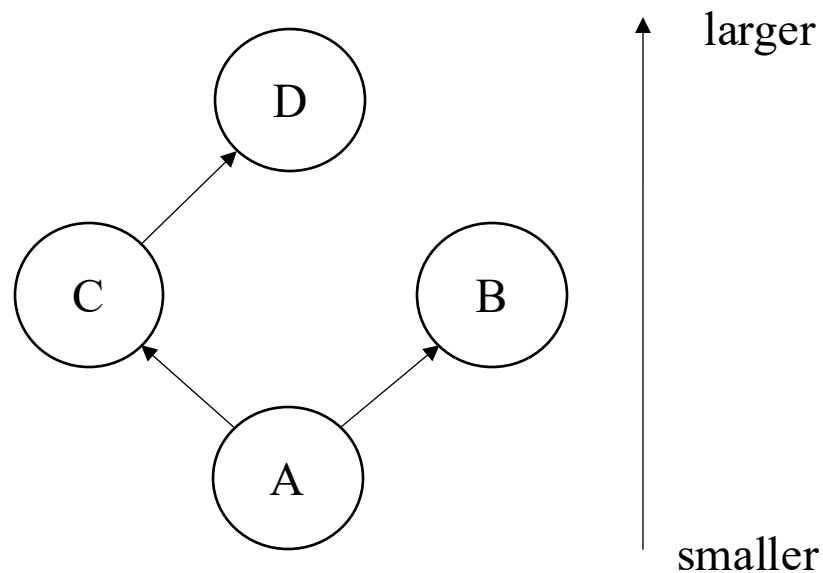
# Futural Plan

1. Add more models and datasets for further detailed comparison in experiment
2. Consistency constrain
3. Extend to time-series data

# Motivation: Evaluate Individual Treatment Effect

Fundamental Problem: counterfactual is unknown

Methodology:



M1,M2对factual的预测是无偏的，对counterfactual的预测是无偏的，所有预测误差服从高斯分布，其中M1和M2对于counterfactual的预测方差相比factual的预测方差较大，则在左图假设下，在factual上具有较低MSE的模型以接近1的概率具有较低的ITE误差，因此，可以通过Factual上的MSE来评估ITE。

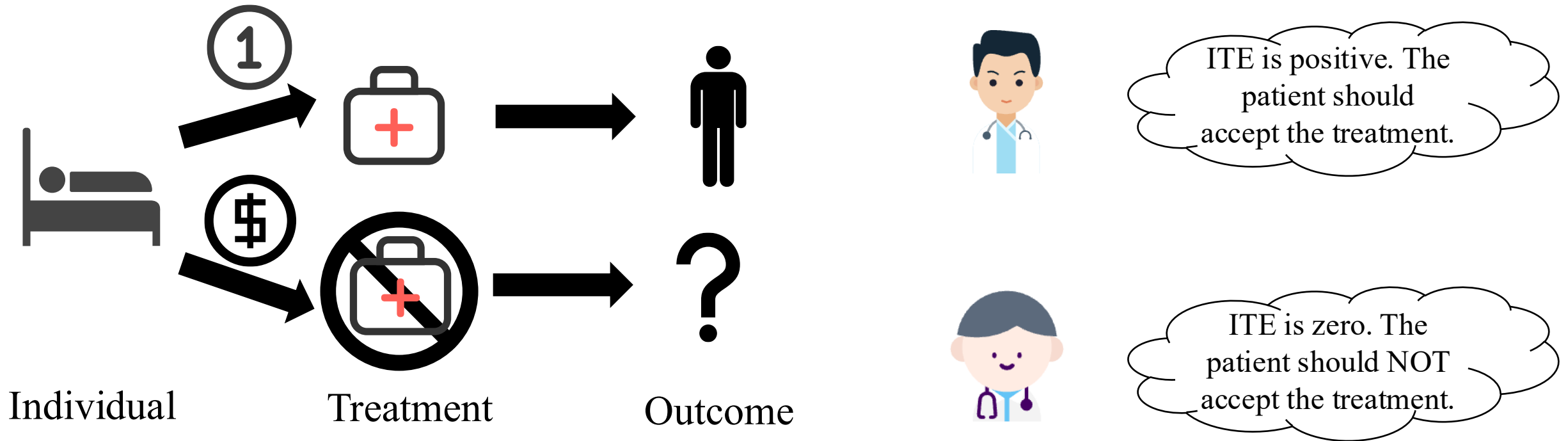
A = MAE of M1 for factual

B = MAE of M1 for counterfactual

C = MAE for M2 for factual

D = MAE for M2 for counterfactual

# Individual Treatment Effect Evaluation



How to evaluate causal models for ITE estimation task?

Fundamental Problem: only one clinical ending can be measured

# Problem Formulation

- Population Evaluation: Given  $M_1: (A, X) \rightarrow Y$ ,  $M_2: (A, X) \rightarrow Y$ , test dataset  $(A, X, Y, F=1)$  where  $A$  is randomized, compare  $MSE \left( M(1, X) - M(0, X) - (Y(1) - Y(0)) \right)$  of  $M_1$  and  $M_2$  where only one of  $Y_i(1)$ ,  $Y_i(0)$  is given.
- Individual Evaluation: Given  $M_1: (A, X) \rightarrow Y$ ,  $M_2: (A, X) \rightarrow Y$ , individual  $(1, x, y, F=1)$  where  $A$  is randomized, compare  $|M(1, x) - M(0, x) - (y - Y(0))|$  of  $M_1$  and  $M_2$  where  $Y(0)$  is not given.

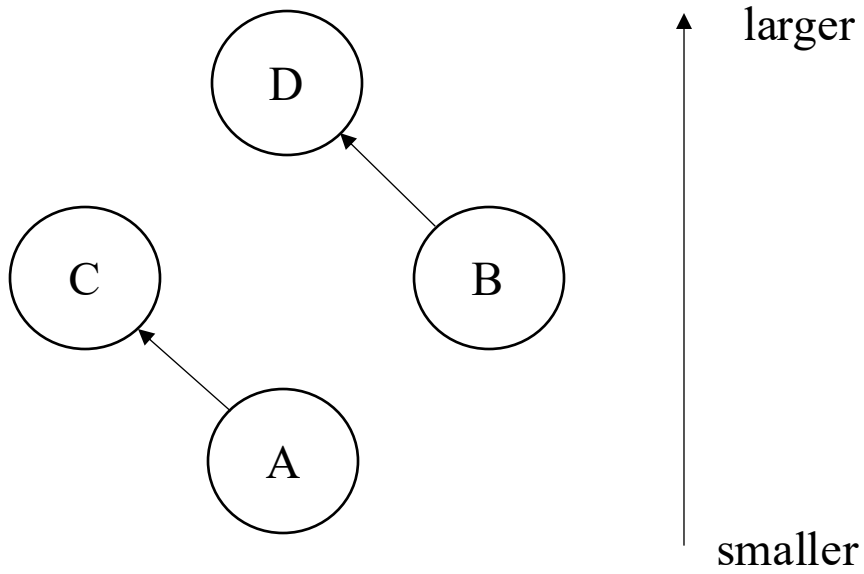


# Methodology

- $M_1, M_2$  is unbiased on factual dataset and counterfactual dataset
- Factual individual (Y can be measured) and counterfactual individual (Y can not be measured) with same randomized treatment A follows identical distribution
  - $MSE(M_1, [F]) < MSE(M_2, [F]) \rightarrow MSE(M_1, [F, CF]) < MSE(M_2, [F, CF])$   
with high probability ( $\geq 95\%$ ) as  $n$  increase  $n \geq \frac{16}{((\frac{MSE(M_2)}{MSE(M_1)})^2 - 1)^2}$ ,  $n = 77$  if  
ratio is 1.1
- So, we can use MSE on randomized factual data to evaluate ITE

**Thanks!**

# Methodology



## Assumption

M1, M2 is unbiased on factual data

M1, M2 is unbiased on counterfactual data

$$A < C \rightarrow A+B < C+D$$

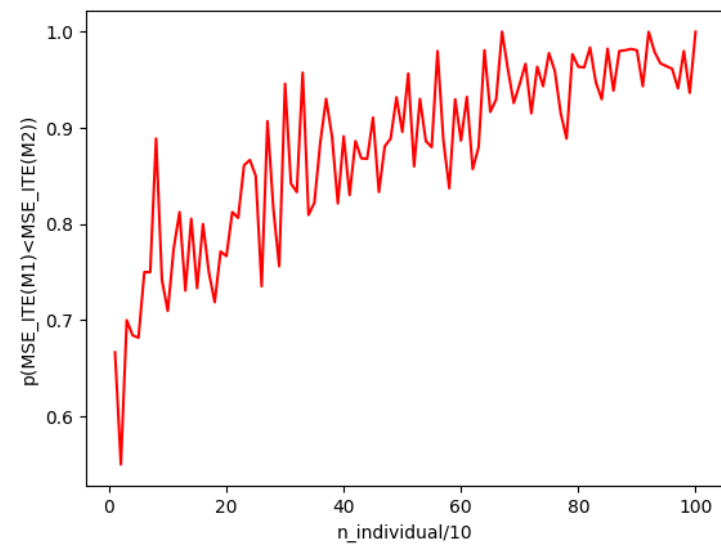
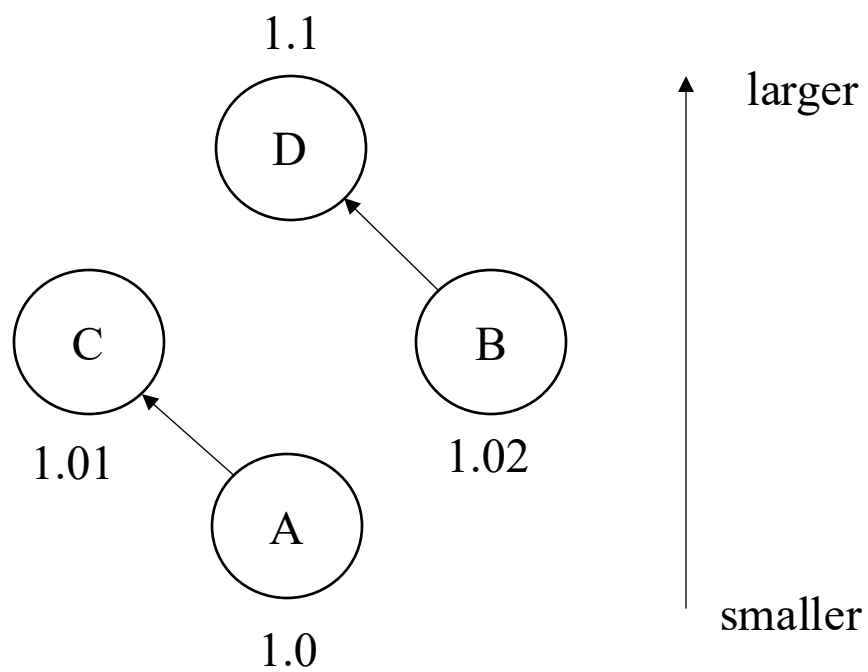
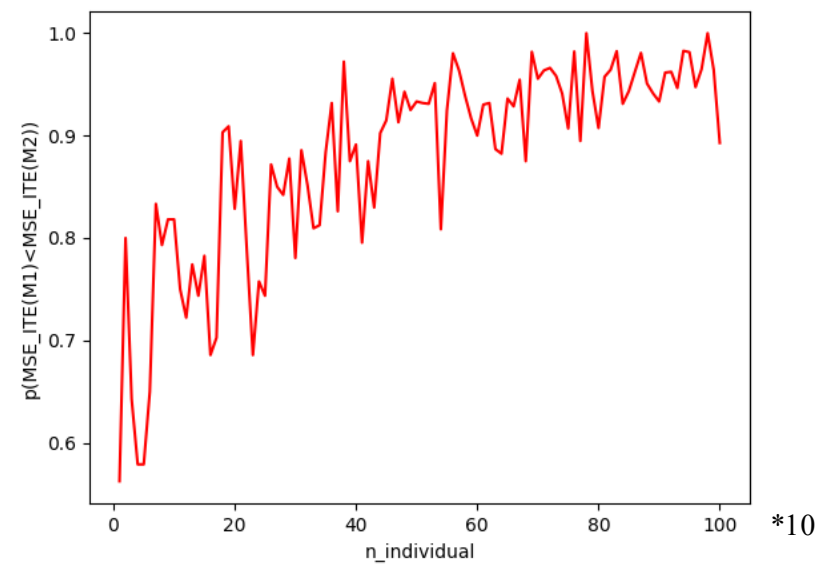
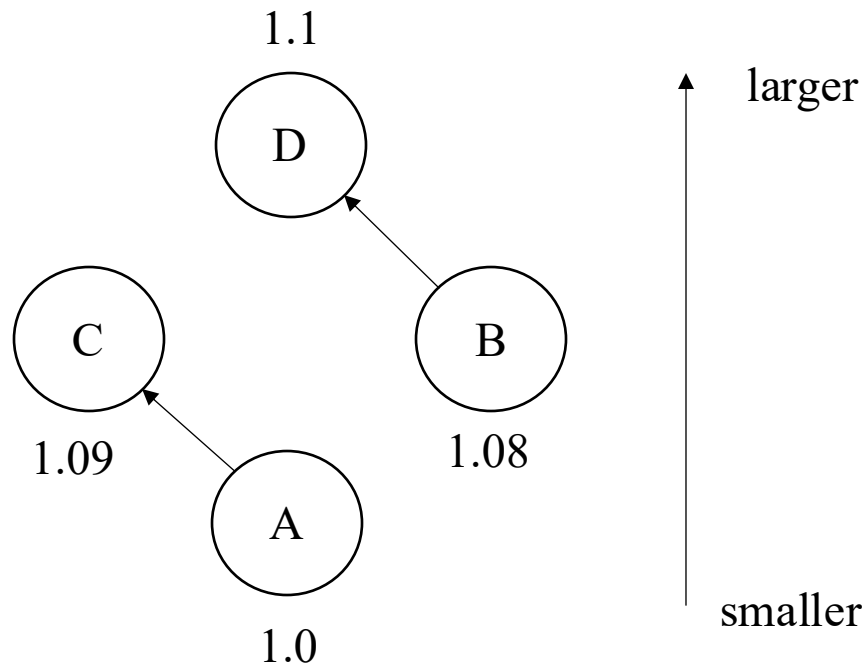
Errors are both sampling from Gaussian

A = Testing MSE of M1 for factual

B = Testing MSE of M1 for counterfactual

C = Testing MSE for M2 for factual

D = Testing MSE for M2 for counterfactual



## Conclusion

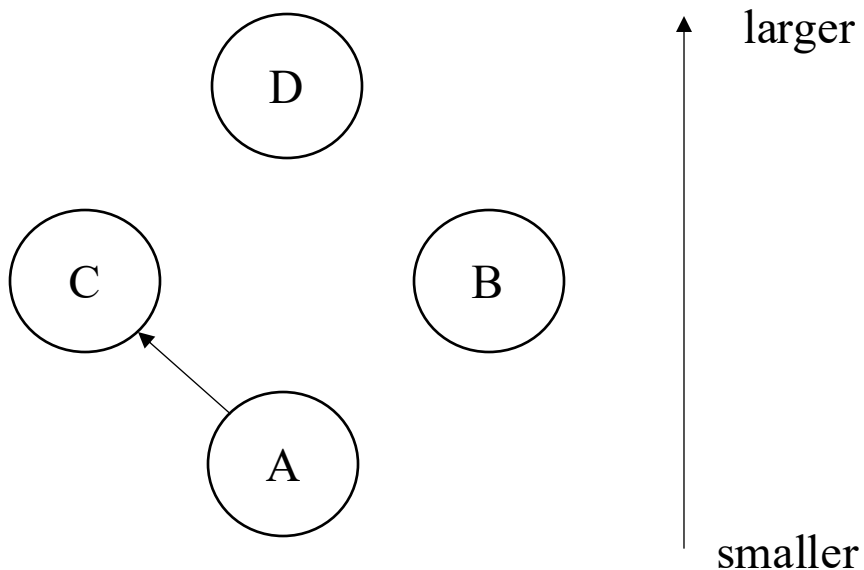
容易证明

当 $n$ 趋向无穷，  
 $\text{MAE}(\text{ITE\_M1}) < \text{MAE}(\text{ITE\_M2})$ 的概率是1，  
收敛速率和 $\sqrt{(A+B)/(C+D)}$ 线性相关

因此，评估模型时不需要知道反事实输出，只需要计算事实数据上的预测的MSE误差对A和C进行ite估计进行评估

- $N \geq 16 / ((C+D)/(A+B))$
- $N \geq \frac{16}{\frac{C+D^2}{A+B} - 1^2}$ , 95%置信度
- 如果比值为1.1，仅需要77个样本，比值越大需要的样本越少

# 进一步解释和弱化假设



A = Testing MSE of M1 for factual  
B = Testing MSE of M1 for counterfactual  
C = Testing MSE for M2 for factual  
D = Testing MSE for M2 for counterfactual

$$A < C \rightarrow A+B < C+D$$

- 对于任意模型，接受治疗和不治疗的人数相同，对于测试集中全部个体估计  $\text{treatment}=1$  时的结局时，假设治疗组个体和不治疗组个体的误差为同分布
- 这是因为
  - 1、test 中的  $\text{treatment}$  是被随机分配的，因此对于测试中  $\text{treatment}$  相同的个体以及反事实个体 ( $A=1, X, Y$ ) 可以认为在同一个分布
  - 2、test 的  $\text{factual}$  数据 ( $A, X, Y$ ) 是训练时没见过的，而对应的  $\text{counterfactual}$  数据 ( $!A, X, Y$ ) 也是训练时没见过的，所以他们的均方误差相同

因此，对于测试集， $A < C$  推出  $A+B < C+D$  是可靠的假设