

*Partial
Transportability* for
**Domain
Generalization**



Kasra Jalaldoust*
Columbia University



Alexis Bellot*
Google DeepMind



Elias Bareinboim
Columbia University

00

Preface

The Research Landscape

Guided by Empirical Evaluations – *Data*



“There is increasing concern that in modern research, false findings may be the vast majority of published research claims” – **John P. A. Ioannidis**

“Well-designed experiments demonstrate the effect only for those specific conditions which the experimental and control group have in common, e.g., same geographical region, historical moment, orientation of the stars, orientation in the magnetic field, etc.” – **Donald T. Campbell**

“A learner that makes no a priori assumptions regarding the identity of the target concept has no rational basis for classifying any unseen instances”
– **Tom Mitchell**

Training and Deploying AI

Predicting outcomes for patients (i.e. *test domain*) that may **differ** in multiple aspects from the training domain.

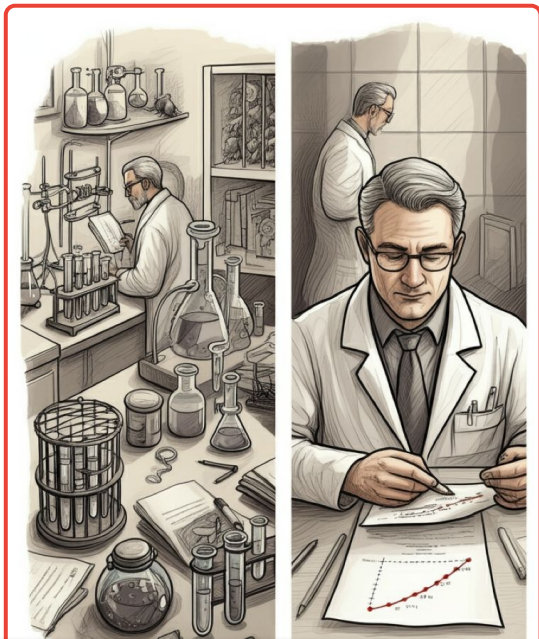


Learning from data collected under particular circumstances (i.e. *training domain*, lab, RCT)

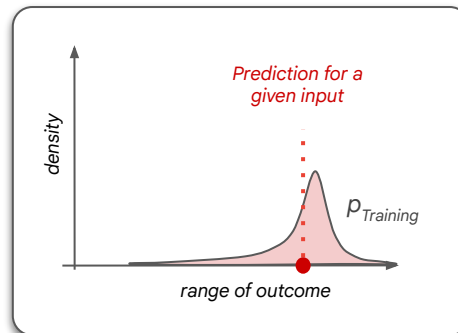


Training and Deploying AI

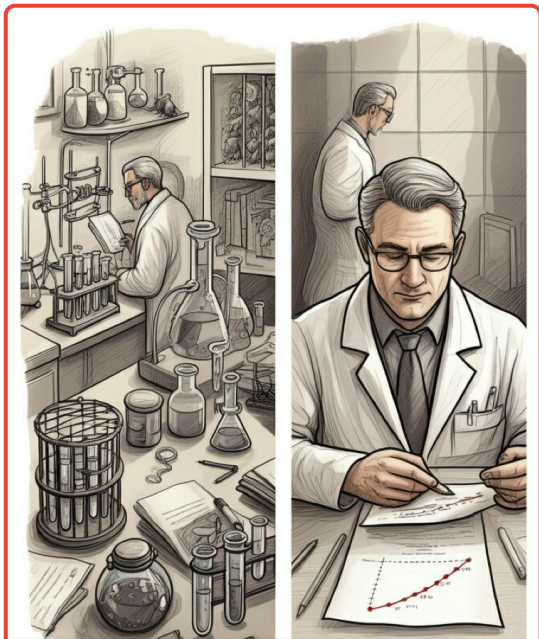
Predicting outcomes for patients (i.e. *test domain*) that may **differ** in multiple aspects from the training domain.



Learning from data collected under particular circumstances (i.e. *training domain*, lab, RCT)



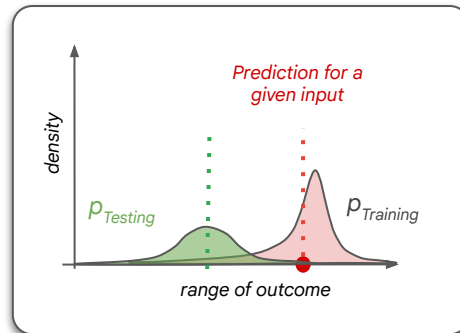
Training and Deploying AI



Learning from data collected under particular circumstances (i.e. training domain, lab, RCT)



Predicting outcomes for patients (i.e. test domain) that may **differ** in multiple aspects from the training domain.



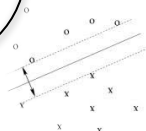
Without an understanding of the relationship between domains, **extrapolation** or **generalization** is impossible.

Paradigms in the Data Sciences

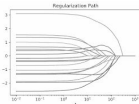
Statistical Learning



Empirical Risk Minimization



Regularization



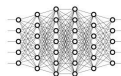
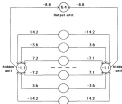
Causal Inference



$$Y_x(u)$$



Backpropagation



Representation learning



Deep Learning

Extrapolation



Interactions between treatment and context



External validity of experiments



Transportability Theory



Mechanisms and generalization

Data-driven
1970s - Present

Model-based
1950s - Present

Paradigms in the Data Sciences

Statistical Learning

Empirical Risk Minimization

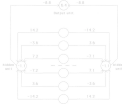
Regularization

Causal Inference

$Y_x(u)$

We have adopted many of the techniques from *data-driven* accounts of ML to solve problems that require the *framing* of *model-based* accounts.

Backpropagation



Representation learning



Deep Learning

External validity of experiments

Interactions between treatment and context



Mechanisms and generalization

Data-driven

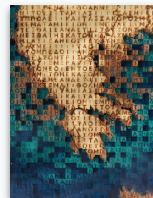
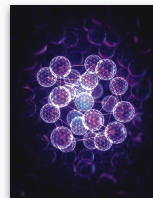
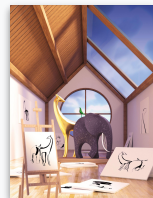
1970s - Present

Model-based

1950s - Present

Clarifying a Subtle Point...

Focusing on **optimal training prediction** is worthwhile.



There is social, economic, and scientific **benefit** in these pursuits.

But neglecting the context in which models are trained is a **limiting perspective**.



Lab

vs.



Real use case

01

A Causal Perspective

Transportability Problem

Depends on invariances across domains

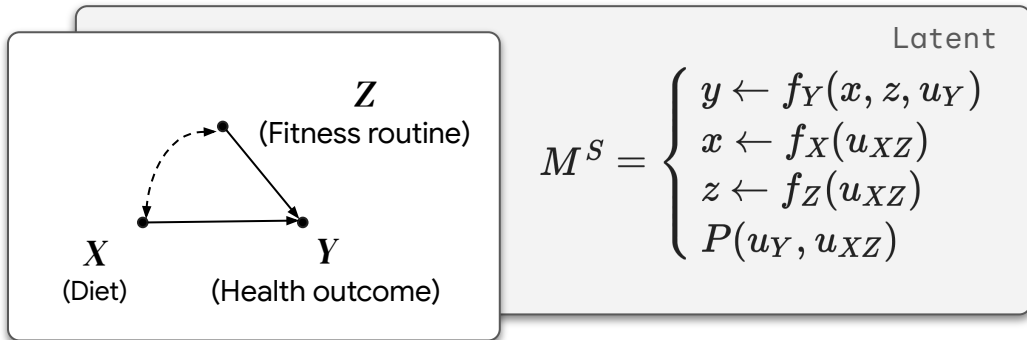
Q: What is the error of a given prediction function in a **target** domain ... given data from a **source** domain?

A: Sometimes can be computed.

$$\mathbb{E}_{P^T} [\underbrace{\mathcal{L}(y, h(x))}_{\text{Loss/Error}}] \quad \text{given} \quad \underbrace{P^S(x, y, \dots)}_{\text{Source distribution}} + \dots$$

Target distribution
Prediction function

Example



Transportability Problem

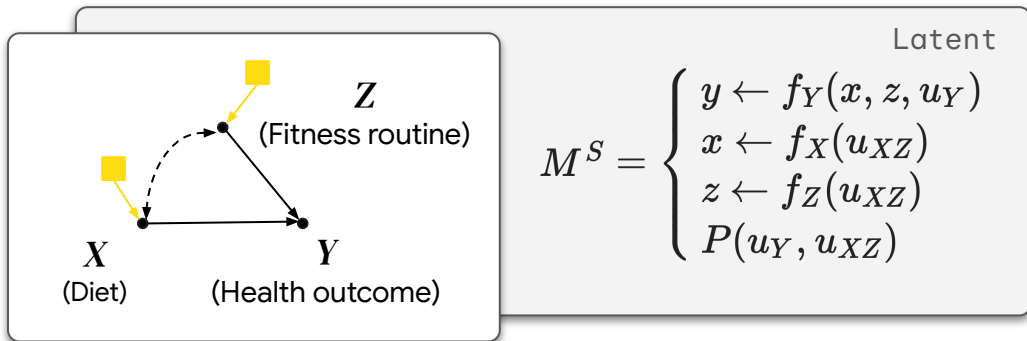
Depends on invariances across domains

Q: What is the error of a given prediction function in a **target** domain ... given data from a **source** domain?

A: Sometimes can be computed.

$$\mathbb{E}_{P^T} [\underbrace{\mathcal{L}}_{\text{Loss / Error}}(\underbrace{y, h(x)}_{\text{Prediction function}})] \quad \text{given} \quad \underbrace{P^S}_{\text{Source distribution}}(x, y, \dots) + \dots$$

Example



Differences across domains captured in Selection diagram.

Transportability Problem

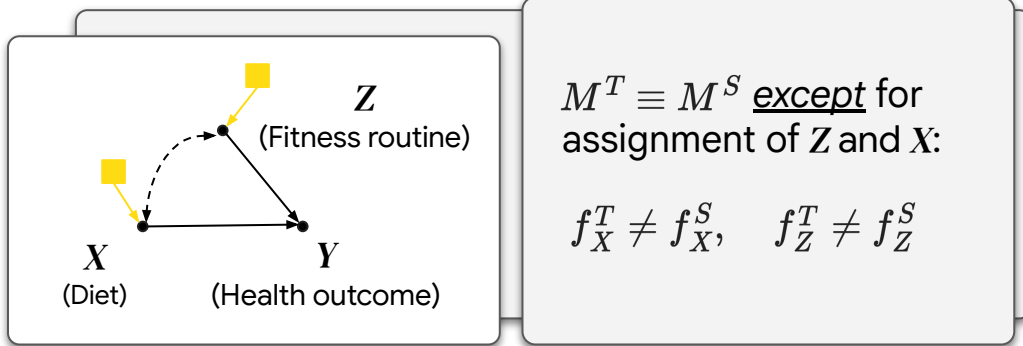
Depends on invariances across domains

Q: What is the error of a given prediction function in a **target** domain ... given data from a **source** domain?

A: Sometimes can be computed.

$$\mathbb{E}_{P^T} [\underbrace{\mathcal{L}}_{\text{Loss / Error}}(\underbrace{y, h(x)}_{\text{Prediction function}})] \quad \text{given} \quad \underbrace{P^S(x, y, \dots)}_{\text{Source distribution}} + \dots$$

Example



Differences across domains captured in Selection diagram

Transportability Problem

Depends on invariances across domains

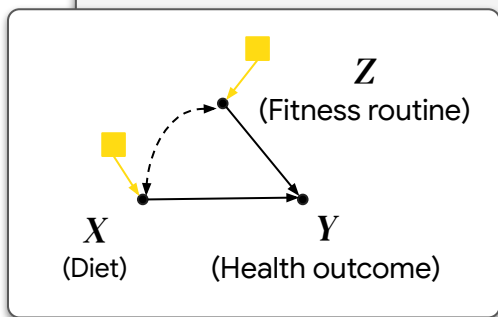
Q: What is the error of a given prediction function in a **target** domain ... given data from a **source** domain?

A: Sometimes can be computed.

$$\mathbb{E}_{P^T} [\underbrace{\mathcal{L}(y, h(x))}_{\text{Loss / Error}}] \quad \text{given} \quad P^S(x, y, \dots) + \dots$$

Target distribution
Prediction function
Source distribution

Example



Compute error of $h(x)$ with access to $P^S(y, x, z), P^T(x, z)$

$$\mathbb{E}_{P^T} [\mathcal{L}(y, h(x))] = \sum_{x,y} \mathcal{L}(y, h(x)) P^T(x, y) \quad \text{by definition}$$

Transportability Problem

Depends on invariances across domains

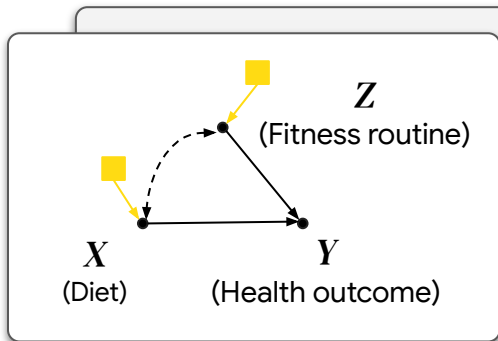
Q: What is the error of a given prediction function in a **target** domain ... given data from a **source** domain?

A: Sometimes can be computed.

$$\mathbb{E}_{P^T} [\underbrace{\mathcal{L}(y, h(x))}_{\text{Loss / Error}}] \quad \text{given} \quad \underbrace{P^S(x, y, \dots)}_{\text{Source distribution}} + \dots$$

Target distribution
Prediction function

Example



Compute error of $h(x)$ with access to $P^S(y, x, z), P^T(x, z)$

$$\begin{aligned} \mathbb{E}_{P^T} [\mathcal{L}(y, h(x))] &= \sum_{x,y} \mathcal{L}(y, h(x)) P^T(x, y) && \text{by definition} \\ &= \sum_{x,y,z} \mathcal{L}(y, h(x)) P^T(y | x, z) P^T(x, z) && \text{by marginal.} \end{aligned}$$

Differences across domains captured in Selection diagram

Transportability Problem

Depends on invariances across domains

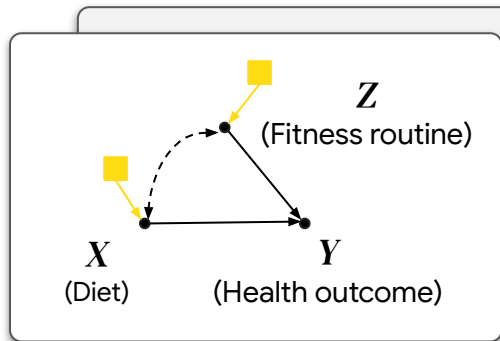
Q: What is the error of a given prediction function in a **target** domain ... given data from a **source** domain?

A: Sometimes can be computed.

$$\mathbb{E}_{P^T} [\mathcal{L}(y, h(x))] \quad \text{given} \quad P^S(x, y, \dots) + \dots$$

Target distribution
Loss / Error
Prediction function
Source distribution

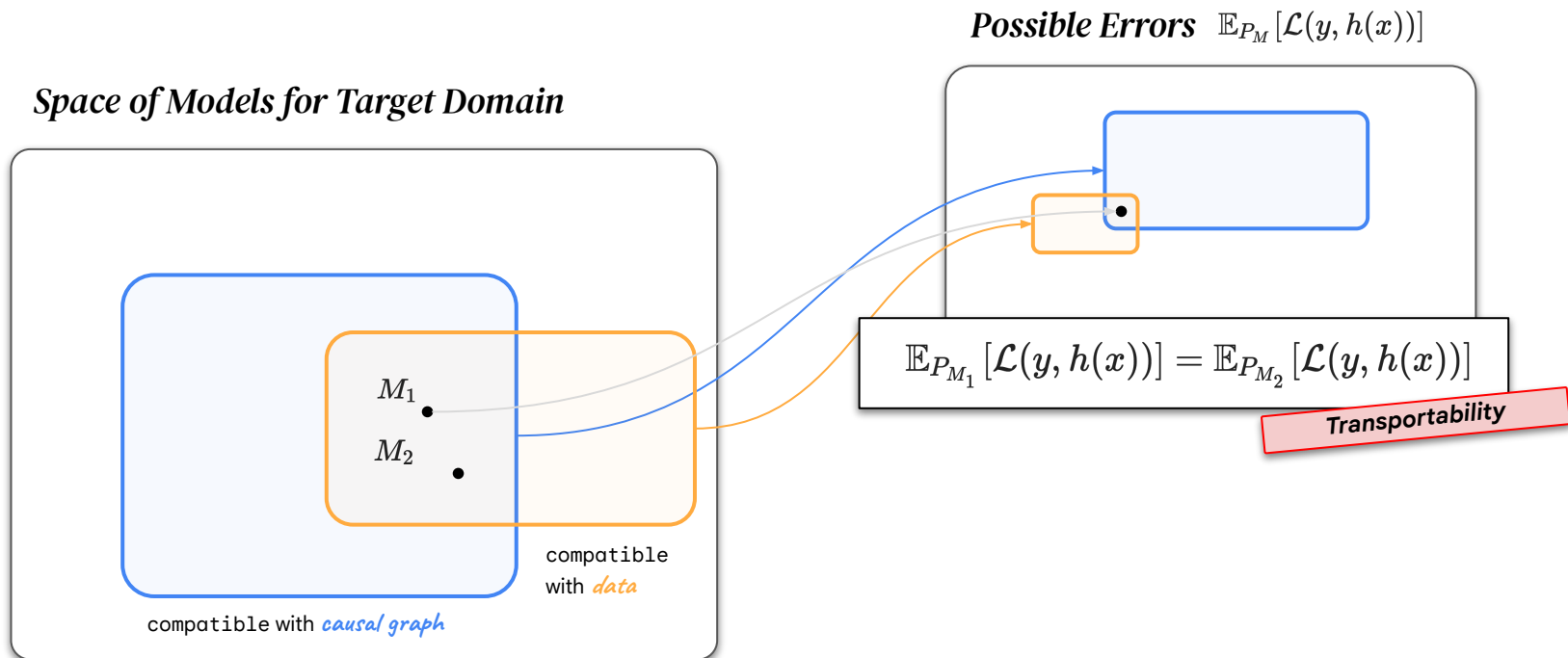
Example



Compute error of $h(x)$ with access to $P^S(y, x, z), P^T(x, z)$

$$\begin{aligned} \mathbb{E}_{P^T} [\mathcal{L}(y, h(x))] &= \sum_{x,y} \mathcal{L}(y, h(x)) P^T(x, y) && \text{by definition} \\ &= \sum_{x,y,z} \mathcal{L}(y, h(x)) P^T(y | x, z) P^T(x, z) && \text{by marginal.} \\ &= \sum_{x,y,z} \mathcal{L}(y, h(x)) P^S(y | x, z) P^T(x, z) && \text{by invariance} \end{aligned}$$

Thinking in terms of Models ...

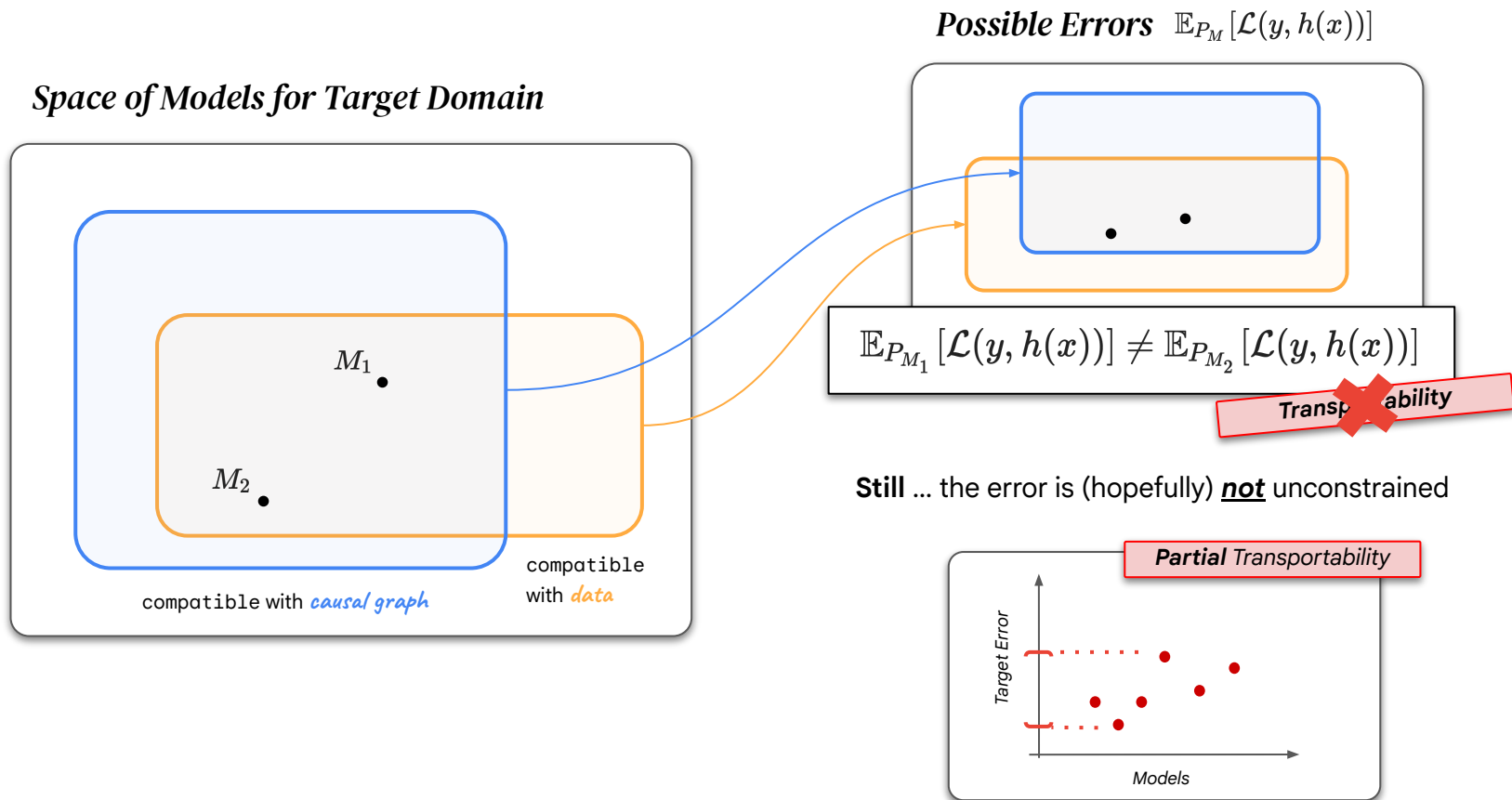


Causal Inference and The Data-Fusion Problem. E. Bareinboim, J. Pearl. PNAS-16.

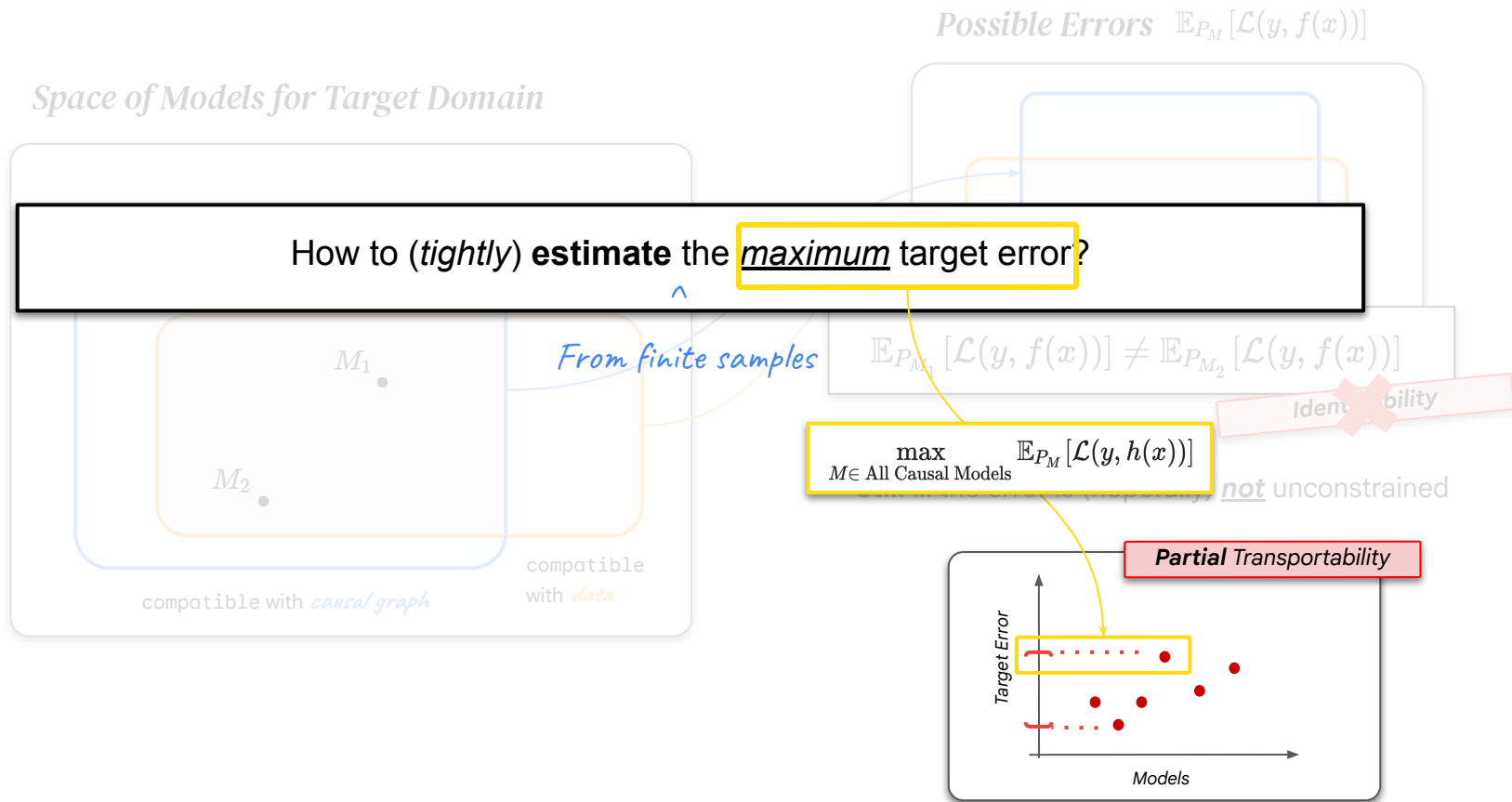
From Statistical Transportability to Estimating the Effect of Stochastic Interventions. J. Correa, E. Bareinboim. IJCAI-19.

General Identifiability with Arbitrary Surrogate Experiments. S. Lee, J. Correa, E. Bareinboim. UAI-19.

Thinking in terms of Models ...



Thinking in terms of Models ...



02

Evaluating (Maximum) Target Error

Possible **if** observed variables are discrete:

$$P^T(Y = y, \mathbf{X} = \mathbf{x})$$

*Non-parametric regime, only assume
data and graph!*

All discrete probabilities and constraints (selection diagrams) induced by SCMs

... may be **equivalently** generated by a special subset of SCMs.

*Powerful because we can search
over NCMs with gradient-descent*

$$\{f_V : V \in \mathbf{V}\}$$

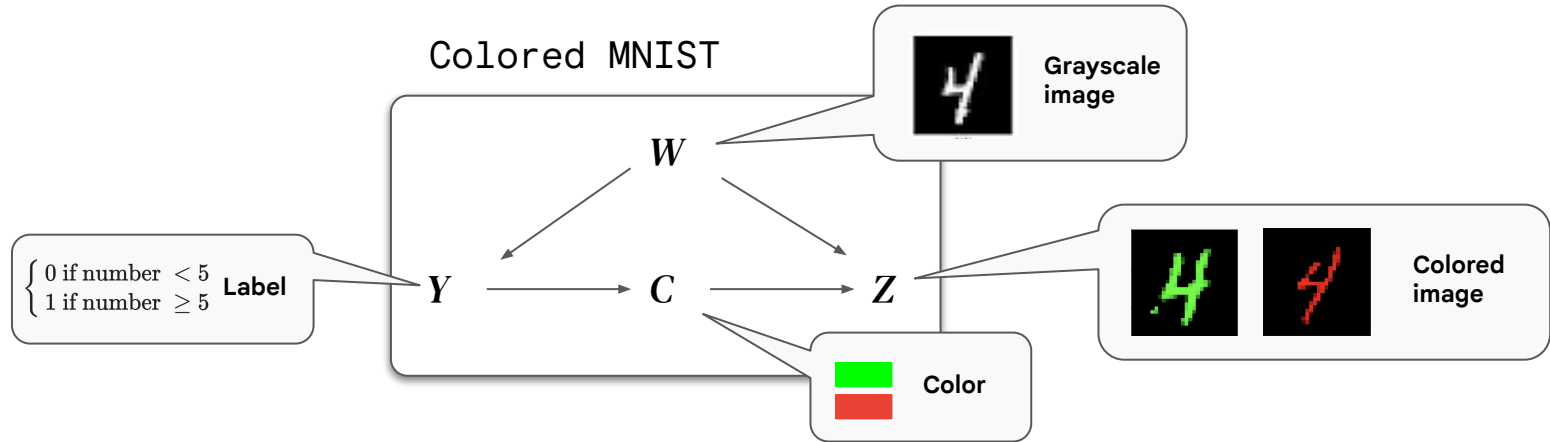
NN

$$\{P(U_V) : U_V \in \mathbf{U}\}$$

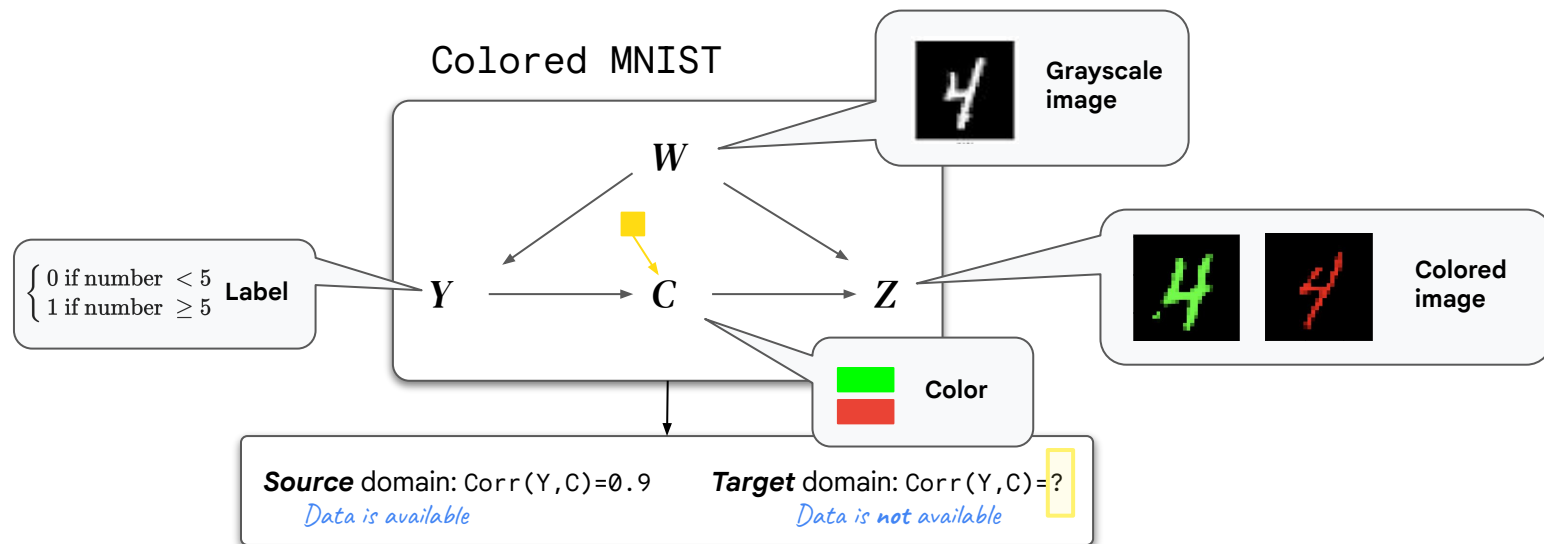
Uniform

$$\max_{M \in \text{Neural Causal Models}} \mathbb{E}_{P_M} [\mathcal{L}(y, h(x))] = \max_{M \in \text{All Causal Models}} \mathbb{E}_{P_M} [\mathcal{L}(y, h(x))]$$

Example of how this works ...



Example of how this works ...



What is the **worst target error** of $h(z)$?

$$M^S = \begin{cases} y \leftarrow \text{NN}_Y(w, u_Y) \\ c \leftarrow \text{NN}_C^S(y, u_C) \\ \dots \end{cases} \quad M^T = \begin{cases} y \leftarrow \text{NN}_Y(w, u_Y) \\ c \leftarrow \text{NN}_C^T(y, u_C) \\ \dots \end{cases}$$

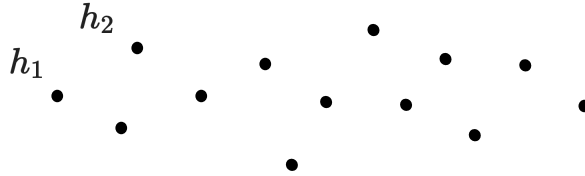
Constrained Optimization Problem:

- $\max_{M \in \text{NCMs}} \mathbb{E}_{P_M} [\mathcal{L}(y, h(x))]$ *Max Target Error*
- $-\lambda \cdot d(P_M, P_{\text{obs}})$ *Match Source Data*

Parameterize according to graph



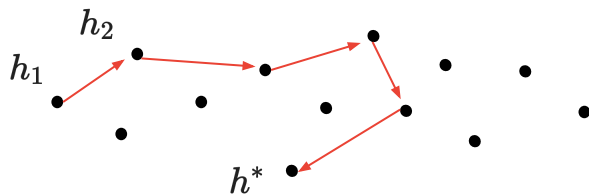
I know the maximum target error of any classifier!



But ... which classifier has the **lowest** maximum target error ?



I know the maximum target error of any classifier!



But ... which classifier has the **lowest** maximum target error ?

Optimization / Search

$$\underbrace{\arg \min_h}_{\text{Optimization / Search}} \underbrace{\max_{M \in \text{All Causal Models}} \mathbb{E}_{P_M} [\mathcal{L}(y, h(x))]}_{\text{Evaluation Max Target Error (Neural-TR)}}$$

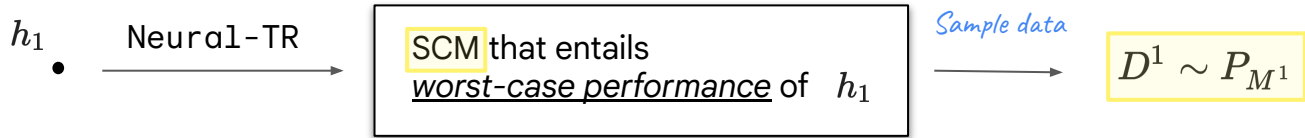
Evaluation Max Target Error (Neural-TR)

03

Optimizing for best
(Maximum) Target Error



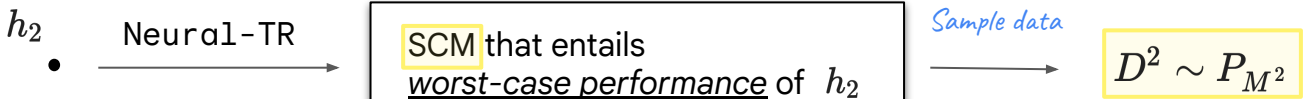
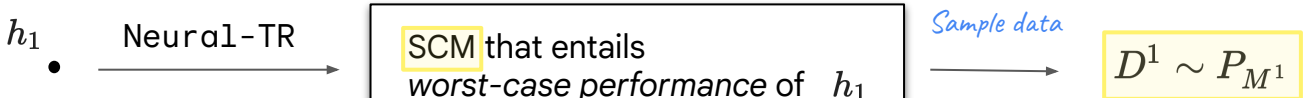
I can **generate adversarial data** for any prediction function !



• $h_2 \leftarrow \arg \min_h \mathbb{E}_{D_1} [\mathcal{L}(y, h(x))]$



I can **generate adversarial data** for any prediction function !



• $h_3 \leftarrow \arg \min_h \max_{D \in \mathbb{D}^*} \mathbb{E}_D [\mathcal{L}(y, h(x))] \quad \mathbb{D}^* = D_1 \cup D_2$

Best worst-case predictor
across collection of datasets



I can generate adversarial data for any prediction function !

Iterate this process ... converge to ...

$$\arg \min_h \max_{M \in \text{All Causal Models}} \mathbb{E}_{P_M} [\mathcal{L}(y, h(x))]$$

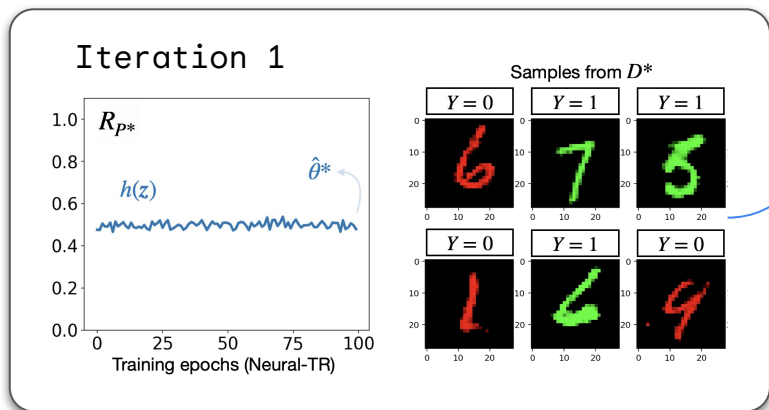
worst-case performance of h_2

$D_1 \cup D_2$

$$h_3 \leftarrow \arg \min_h \max_{D \in \mathbb{D}^*} \mathbb{E}_D [\mathcal{L}(y, h(x))] \quad \mathbb{D}^* = D_1 \cup D_2$$

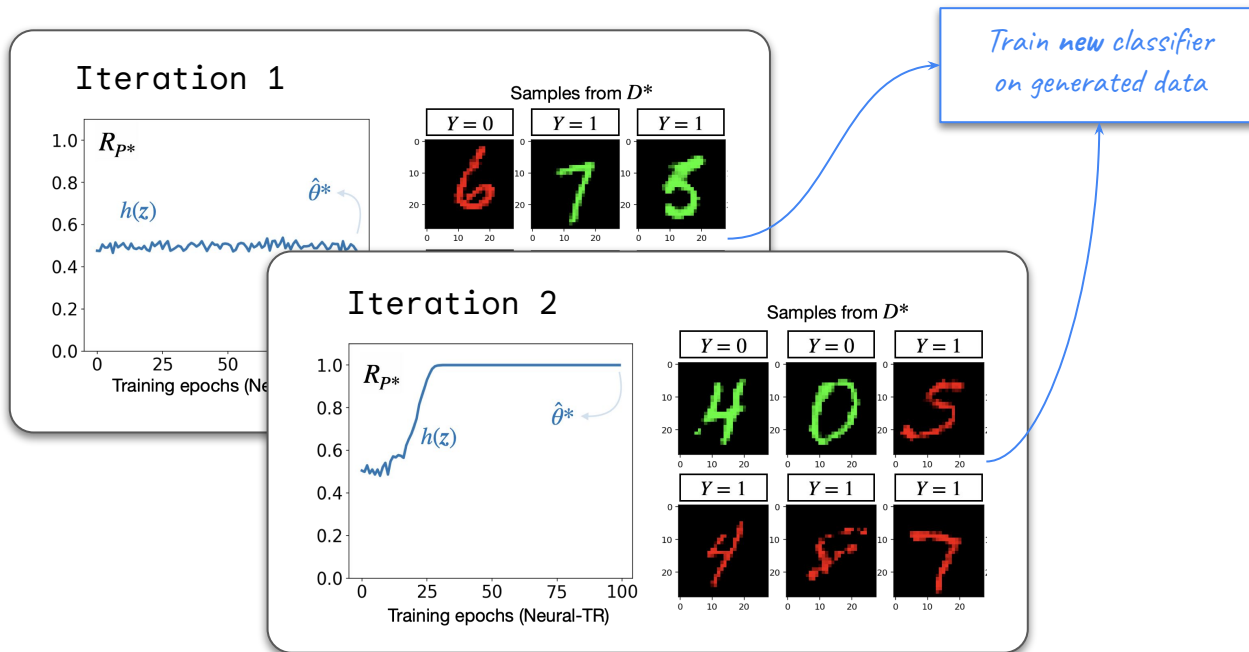
Best worst-case predictor across collection of datasets

What is the best (worst-case) classifier for CMNIST?

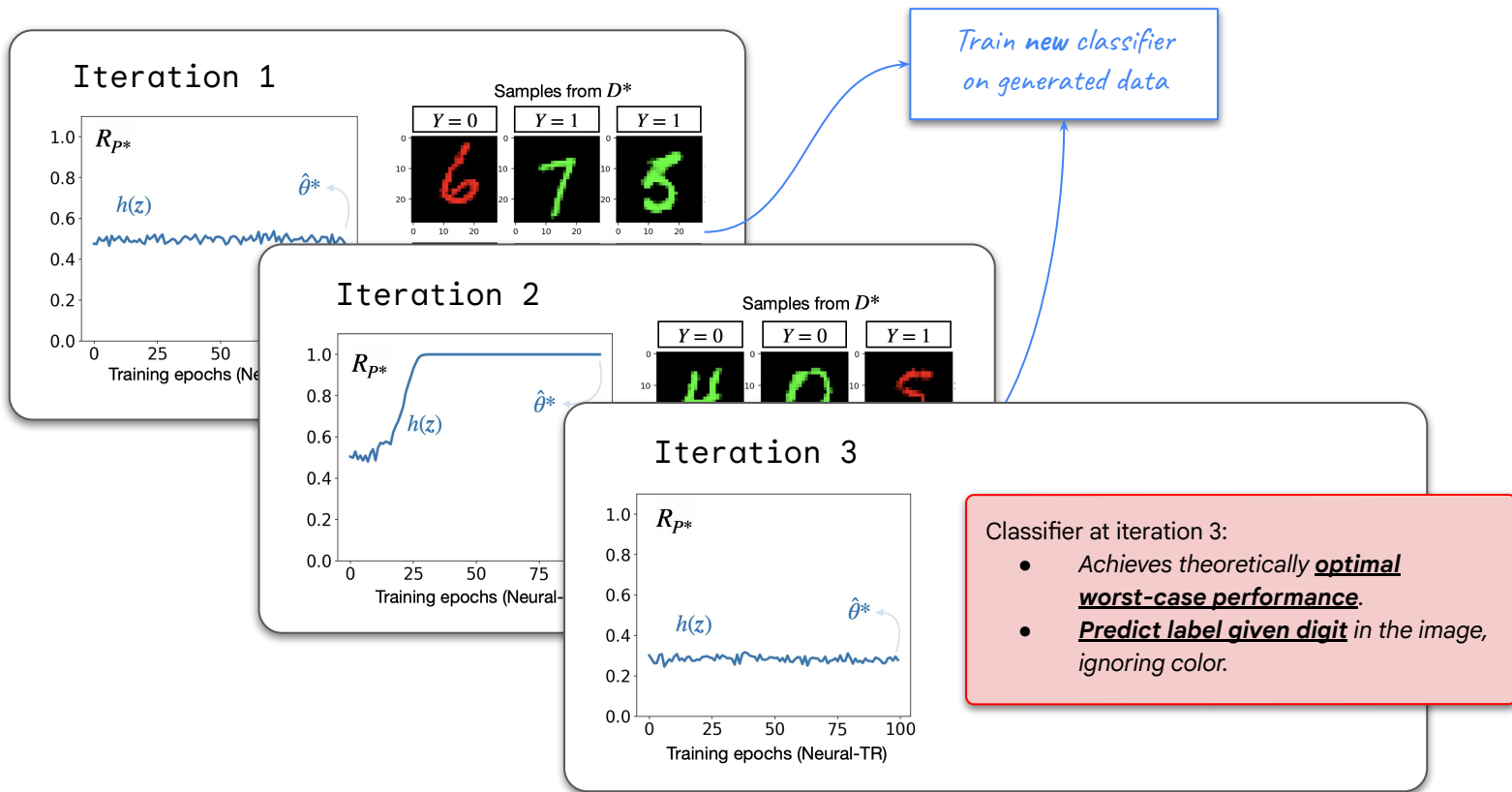


Train new classifier
on generated data

What is the best (worst-case) classifier for CMNIST?

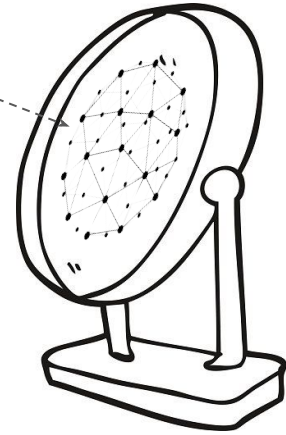
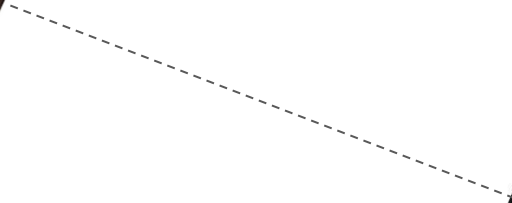


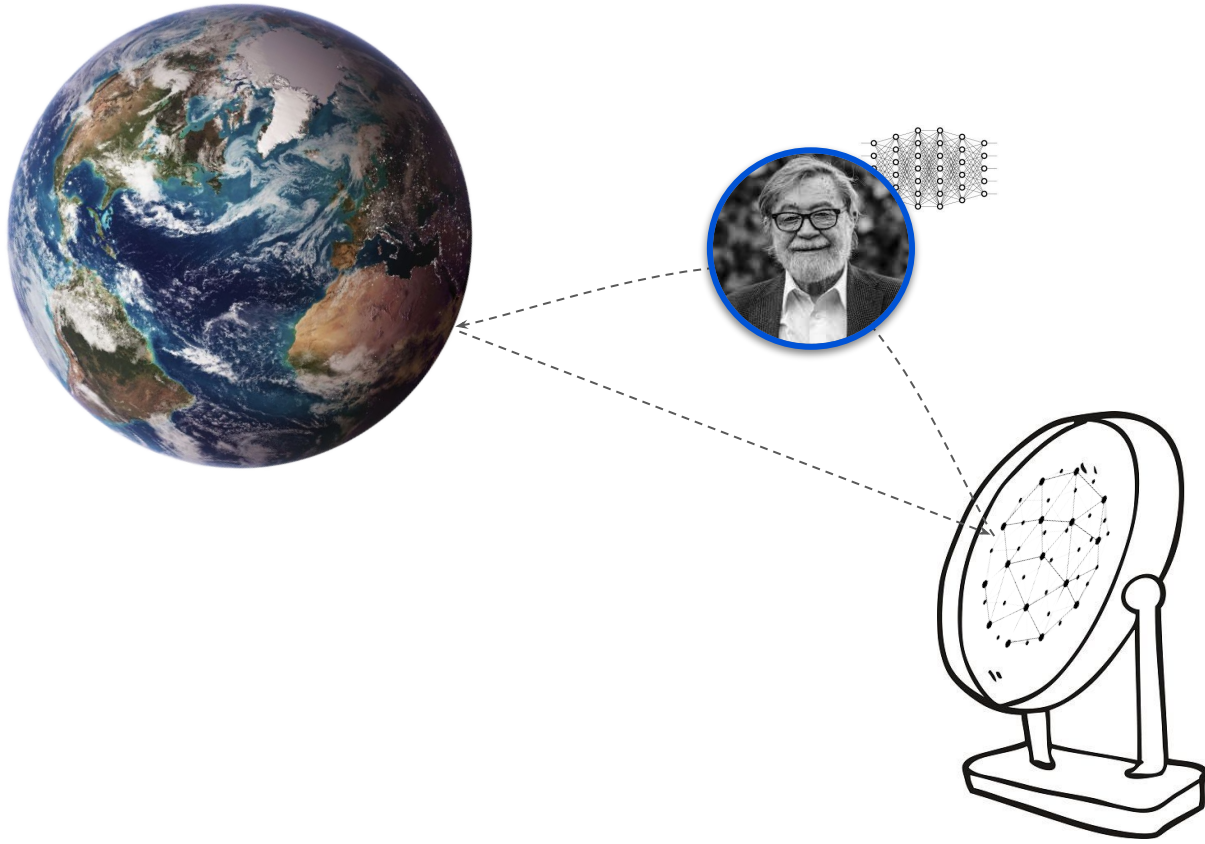
What is the best (worst-case) classifier for CMNIST?



04

Conclusion





Questions?