



Stable Learning:

Finding the Common Ground between Causal Inference and Machine Learning

Peng Cui

Tsinghua University

Now AI is stepping into risk-sensitive areas

Healthcare

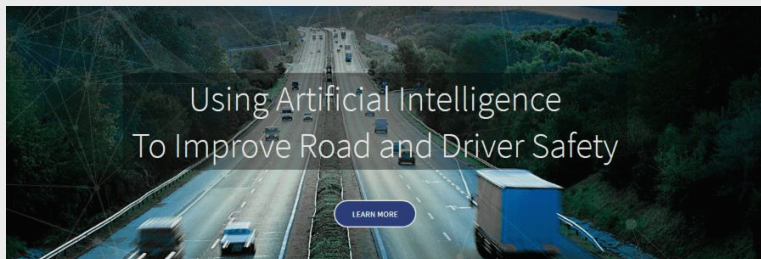


Law



Human

Transportation



Fintech



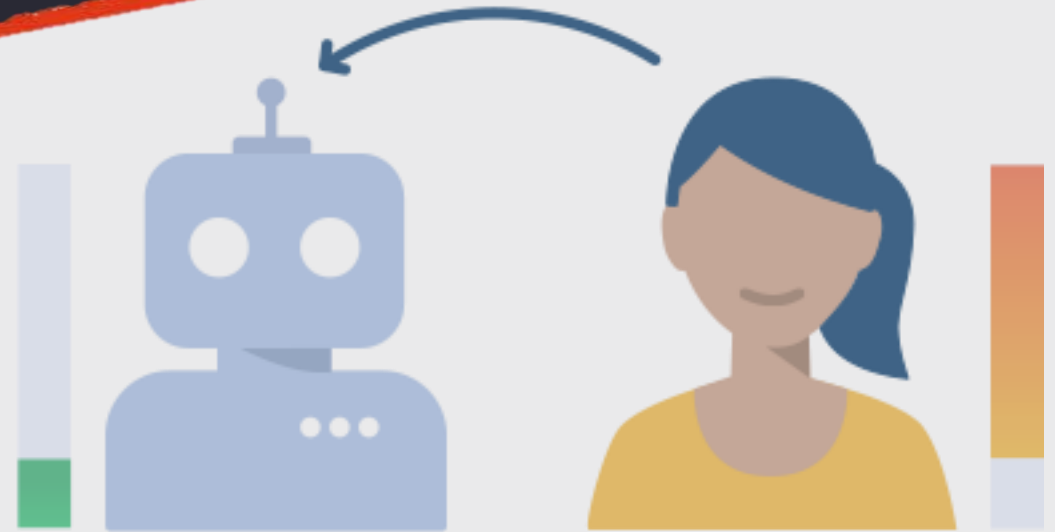
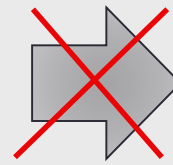
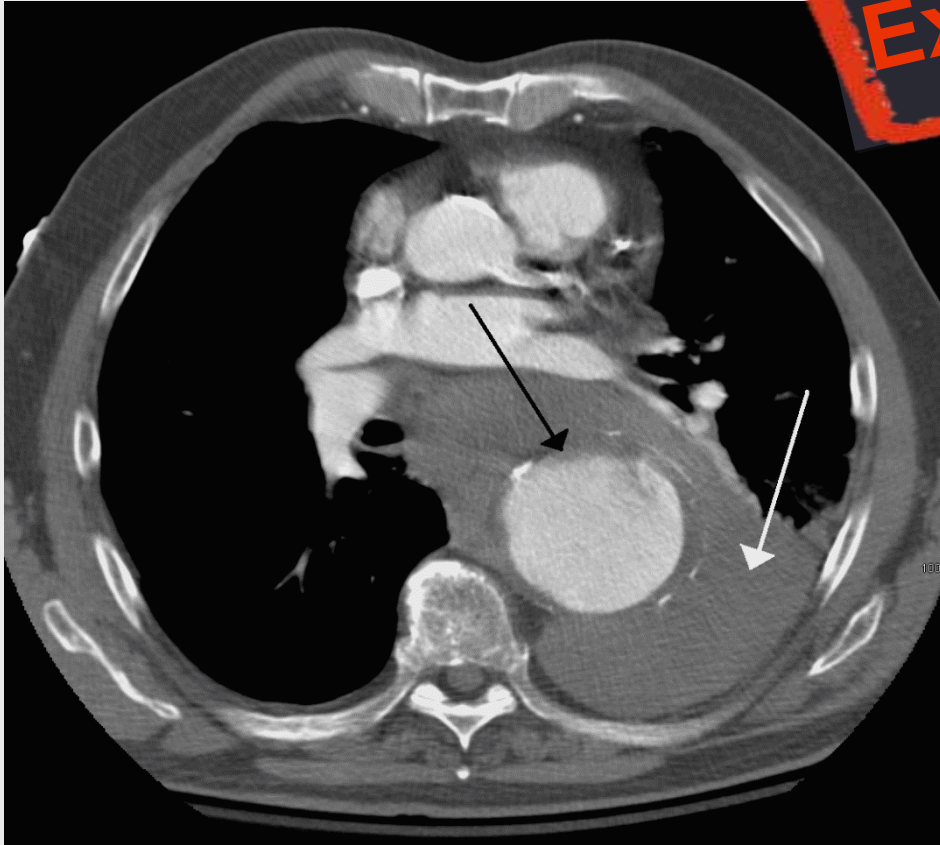
Shifting from *Performance Driven* to *Risk Sensitive*

Risks of Today's AI Algorithms

Unexplainable

Explainability

Human in the loop



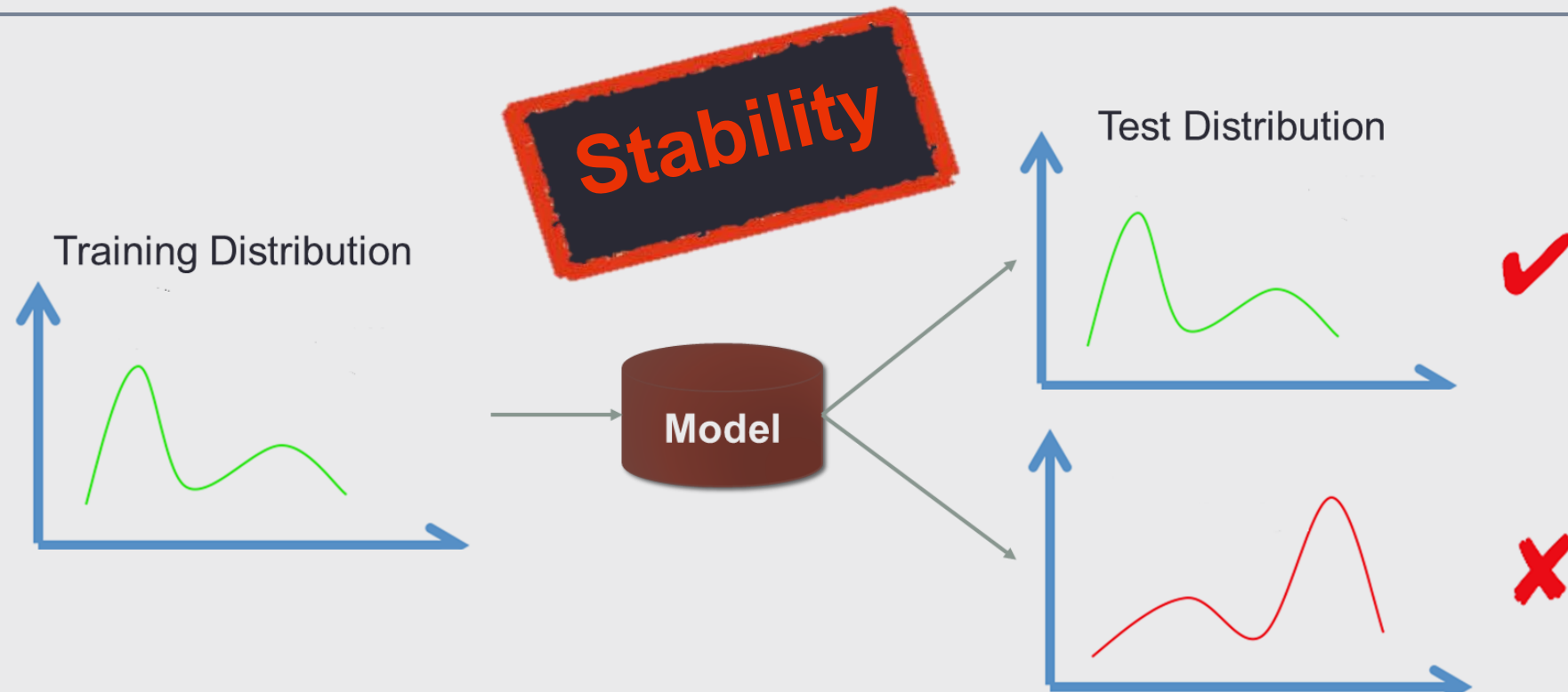
Medical

Military

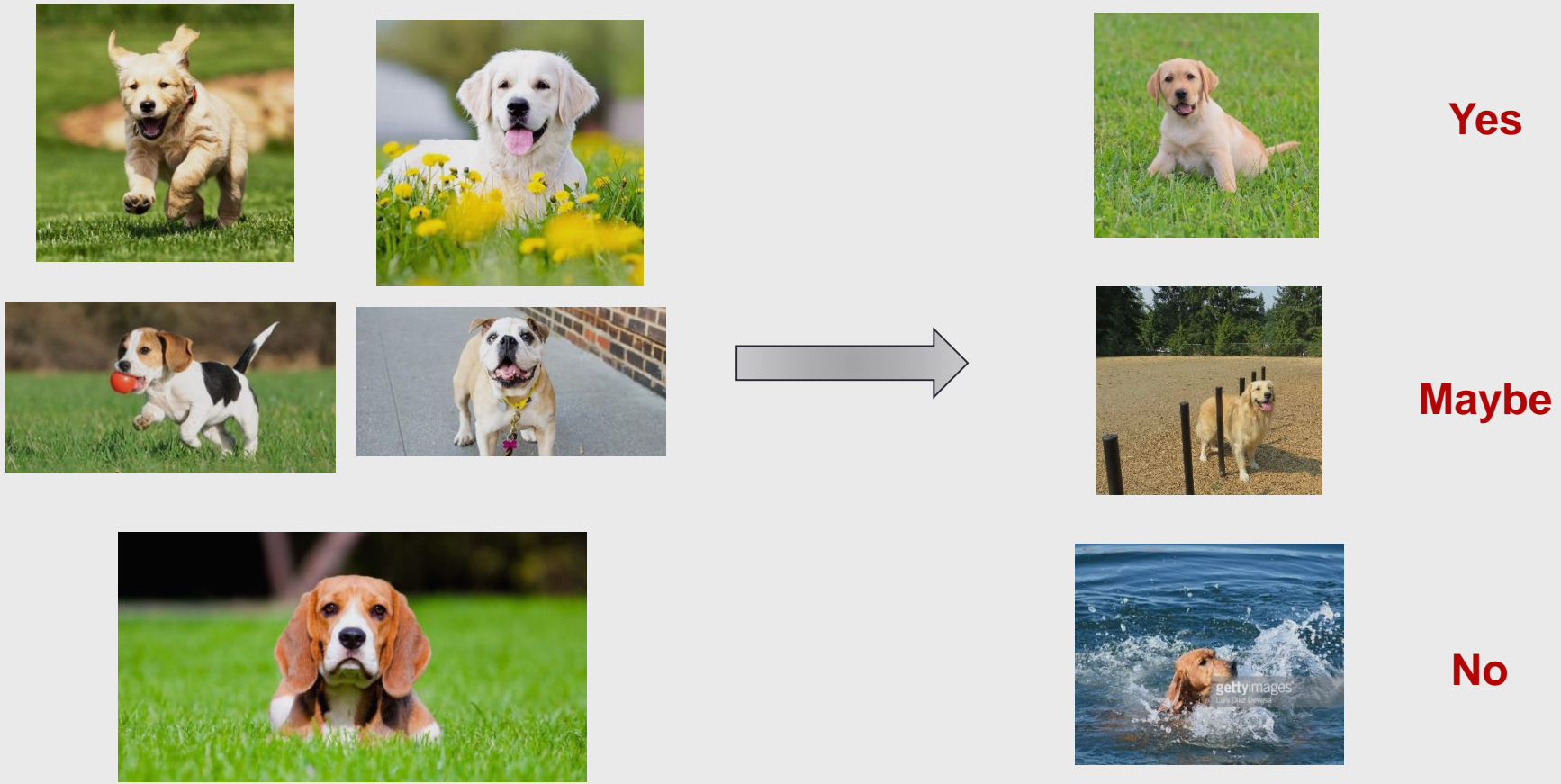
Finance

Risks of Today's AI Algorithms

Most ML methods are developed under I.I.D hypothesis



Risks of Today's AI Algorithms

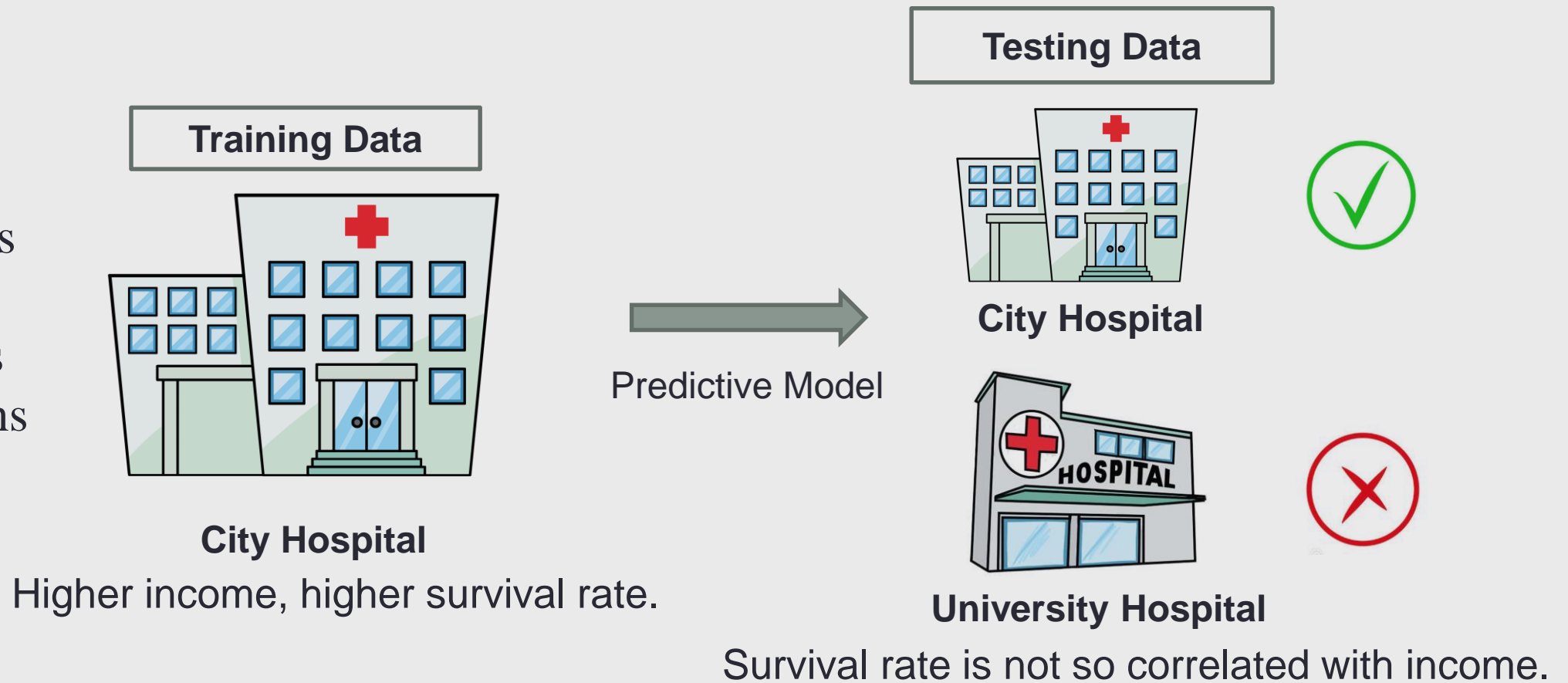


Risks of Today's AI Algorithms

- Cancer survival rate prediction

Features:

- Body status
- **Income**
- Treatments
- Medications



The Current Condition

Explainability

We cannot *understand* AI

Stability

We don't *trust* AI



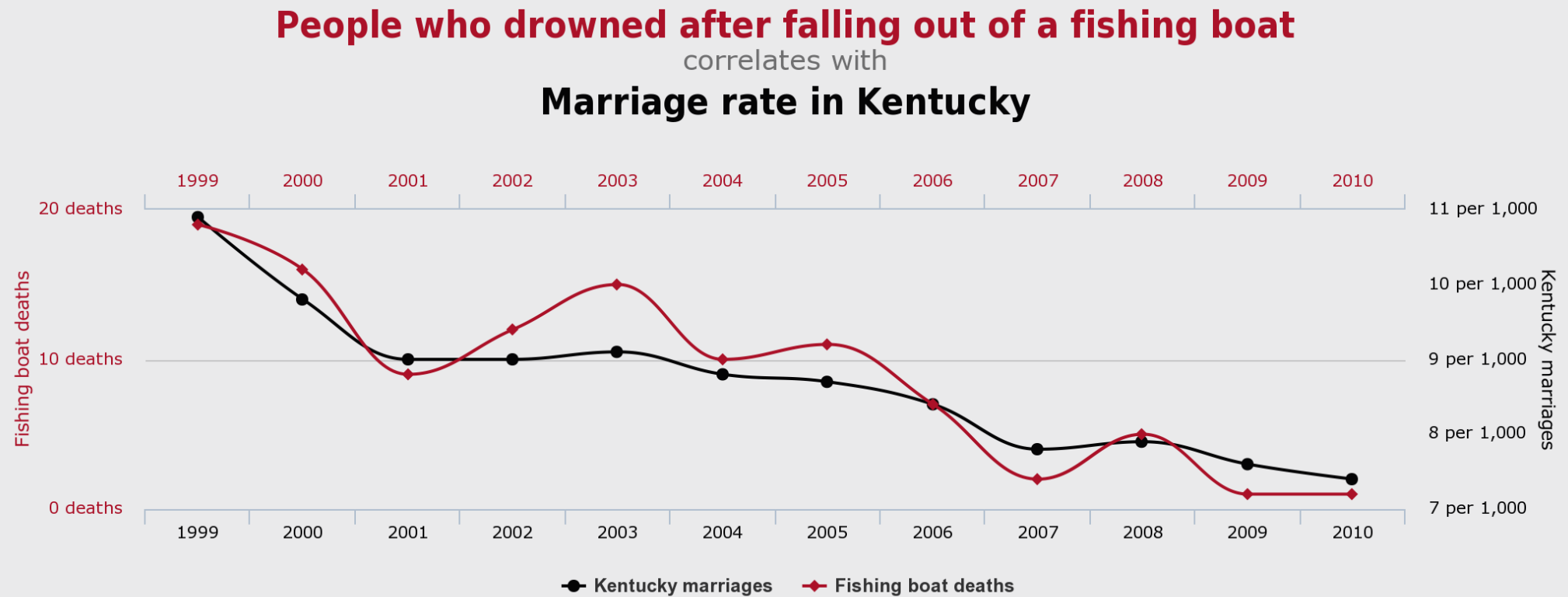
Dilemma

A plausible reason: *Correlation*

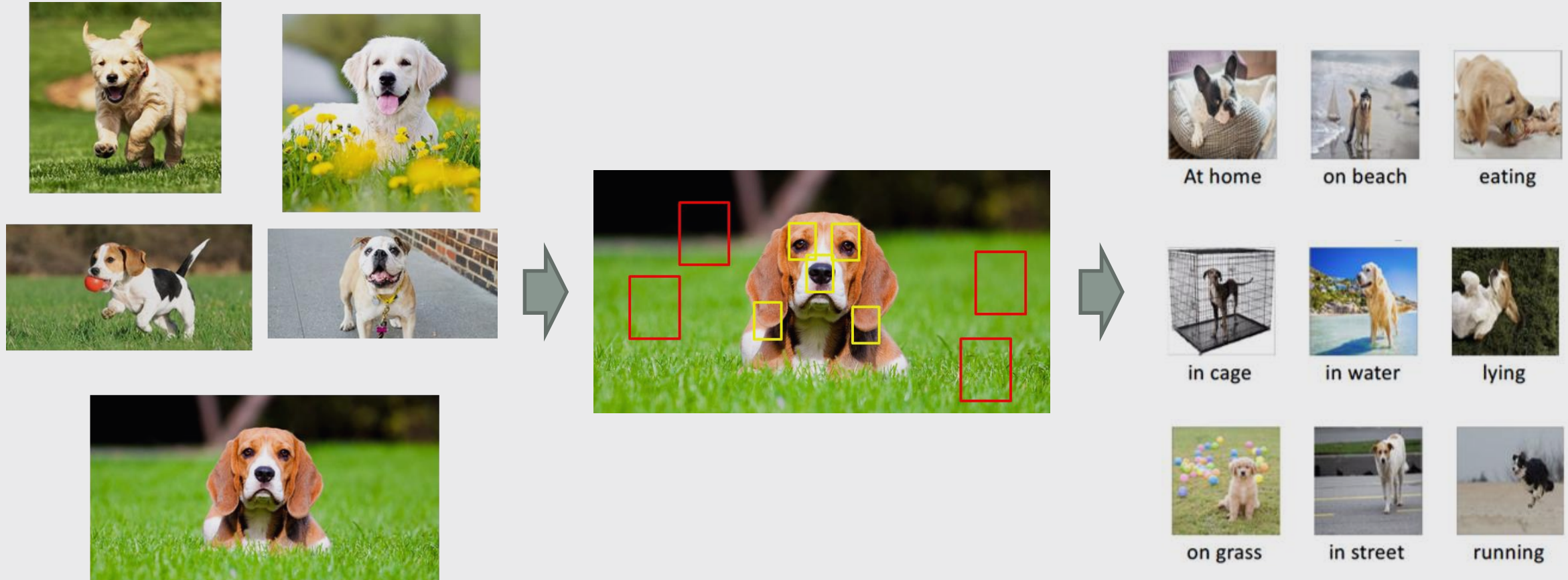
Correlation is the very basics of machine learning.



Correlation is not explainable



Correlation is '*unstable*'

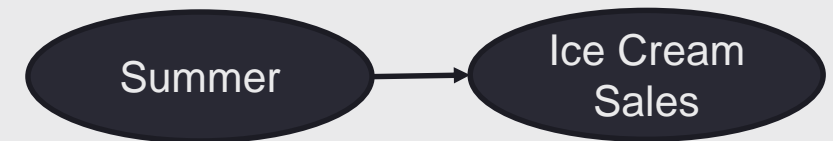
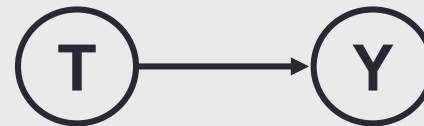


It's not the fault of *correlation*, but the way we use it

• Three sources of correlation:

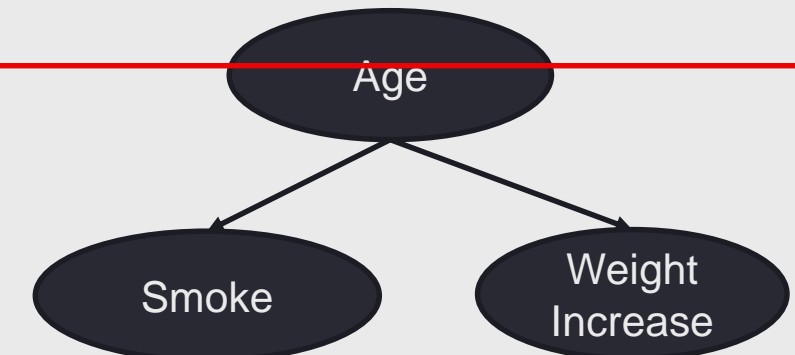
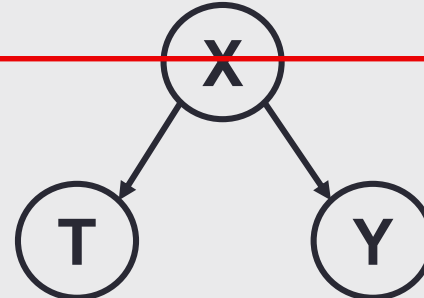
• Causation

- Causal mechanism
- **Stable and explainable**



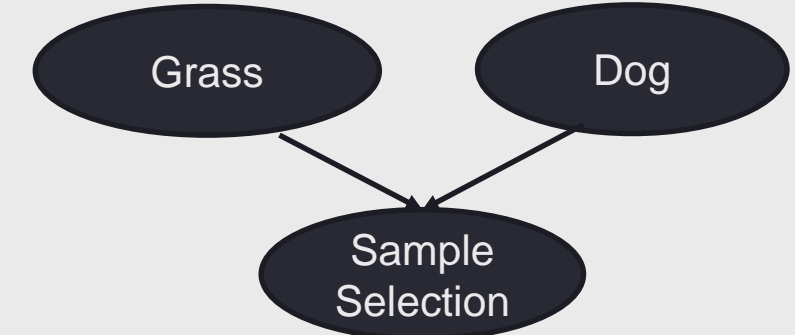
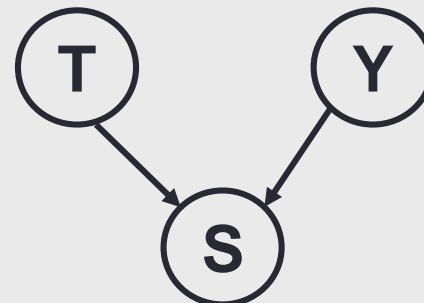
• Confounding

- Ignoring X
- **Spurious Correlation**



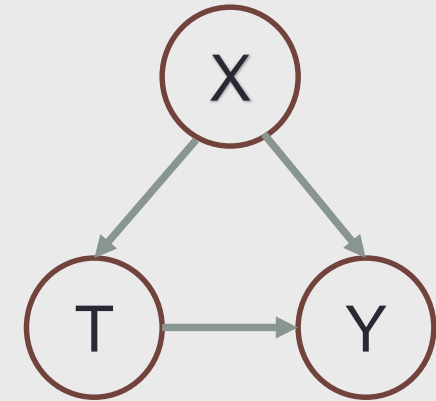
• Sample Selection Bias

- Conditional on S
- **Spurious Correlation**



A Practical Definition of Causality

Definition: T causes Y if and only if
changing T leads to a change in Y,
while keeping everything else constant.

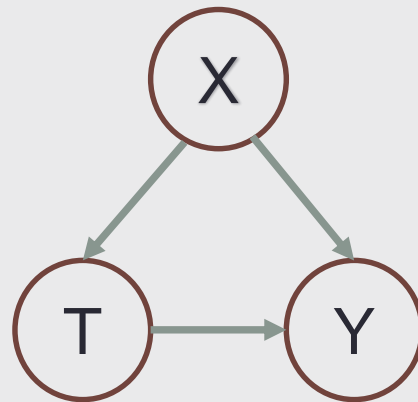


Causal effect is defined as the magnitude by which Y is changed by a unit change in T.

Called the “interventionist” interpretation of causality.

The *benefits* of bringing causality into learning

Causal Framework



T: grass
X: dog nose
Y: label

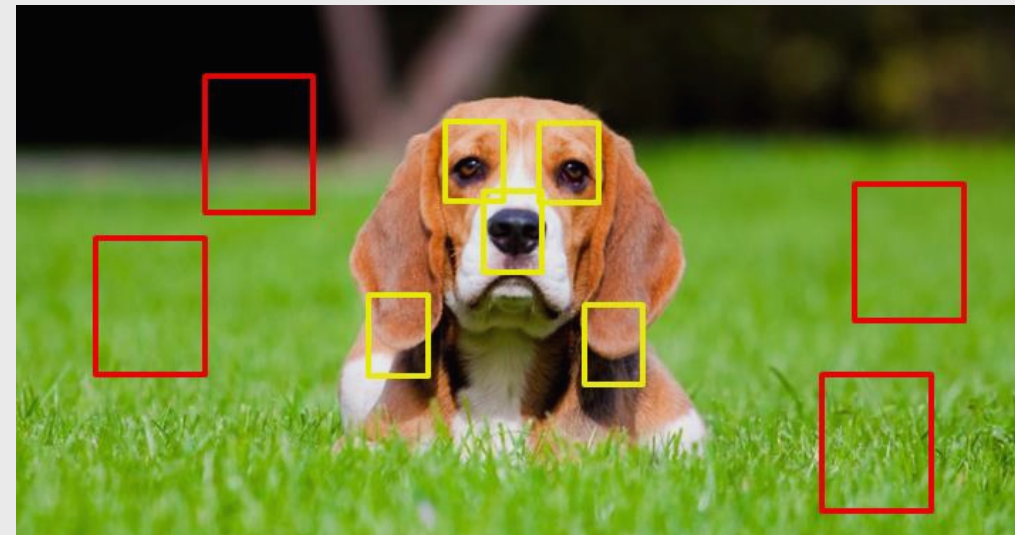


Grass—Label: Strong correlation

Weak causation

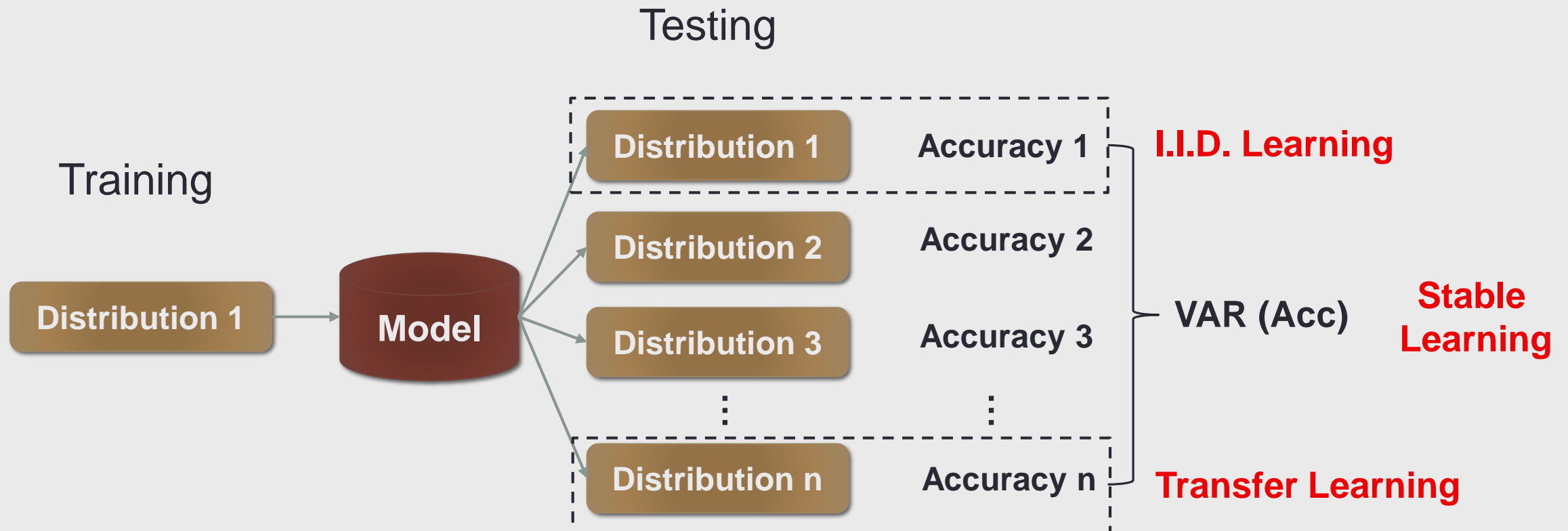
Dog nose—Label: Strong correlation

Strong causation

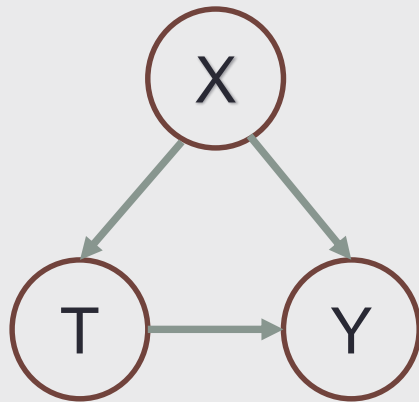


More *Explainable* and More *Stable*

Stable Learning



Revisit Directly Balancing for causal inference



Typical Causal Framework

Directly Confounder Balancing

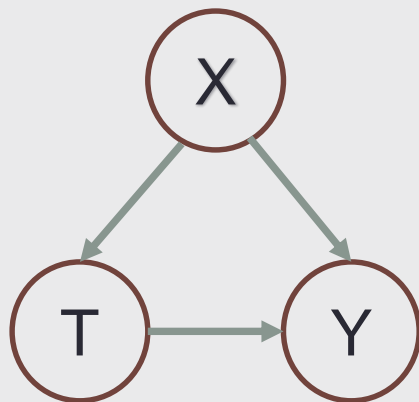
Given a feature T

Assign different weights to samples so that the samples with T and the samples without T have similar distributions in X

Calculate the difference of Y distribution in treated and controlled groups. (correlation between T and Y)

Sample reweighting can make a variable independent of other variables.

Global Balancing: making all variables independent



Typical Causal Framework

Analogy of A/B Testing

Given **ANY** feature T

Assign different weights to samples so that the samples with T and the samples without T have similar distributions in X

Calculate the difference of Y distribution in treated and controlled groups. (correlation between T and Y)

If all variables are independent after sample reweighting,
Correlation = Causality

Theoretical Guarantee

PROPOSITION 3.3. If $0 < \hat{P}(\mathbf{X}_i = x) < 1$ for all x , where $\hat{P}(\mathbf{X}_i = x) = \frac{1}{n} \sum_i \mathbb{I}(\mathbf{X}_i = x)$, *there exists a solution W^* satisfies equation (4) equals 0 and variables in \mathbf{X} are independent after balancing by W^* .*

$$\sum_{j=1}^p \left\| \frac{\mathbf{X}_{:,j}^T \cdot (W \odot \mathbf{X}_{:,j})}{W^T \cdot \mathbf{X}_{:,j}} - \frac{\mathbf{X}_{:,j}^T \cdot (W \odot (1 - \mathbf{X}_{:,j}))}{W^T \cdot (1 - \mathbf{X}_{:,j})} \right\|_2^2, \quad (4)$$



0

b

PROOF. Since $\|\cdot\| \geq 0$, Eq. (8) can be simplified to $\forall j, \forall k \neq j$

$$\lim_{n \rightarrow \infty} \left(\frac{\sum_{t: \mathbf{X}_{t,k}=1, \mathbf{X}_{t,j}=1} W_t}{\sum_{t: \mathbf{X}_{t,j}=1} W_t} - \frac{\sum_{t: \mathbf{X}_{t,k}=1, \mathbf{X}_{t,j}=0} W_t}{\sum_{t: \mathbf{X}_{t,j}=0} W_t} \right) = 0$$

with probability 1. For W^* , from Lemma 3.1, $0 < P(\mathbf{X}_i = x) < 1$, $\forall x, \forall i, t = 1$ or 0 ,

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t: \mathbf{X}_{t,j}=t} W_t^* &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{x: x_j=t} \sum_{t: \mathbf{X}_t=x} W_t^* \\ &= \lim_{n \rightarrow \infty} \sum_{x: x_j=t} \frac{1}{n} \sum_{t: \mathbf{X}_t=x} \frac{1}{P(\mathbf{X}_t=x)} \\ &= \lim_{n \rightarrow \infty} \sum_{x: x_j=t} P(\mathbf{X}_t=x) \cdot \frac{1}{P(\mathbf{X}_t=x)} = 2^{p-1} \end{aligned}$$

with probability 1 (Law of Large Number). Since features are binary,

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t: \mathbf{X}_{t,k}=1, \mathbf{X}_{t,j}=1} W_t^* &= 2^{p-2} \\ \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t: \mathbf{X}_{t,j}=0} W_t^* &= 2^{p-1}, \quad \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t: \mathbf{X}_{t,k}=1, \mathbf{X}_{t,j}=0} W_t^* = 2^{p-2} \end{aligned}$$

and therefore, we have following equation with probability 1:

$$\lim_{n \rightarrow \infty} \left(\frac{\mathbf{X}_{:,k}^T (W^* \odot \mathbf{X}_{:,j})}{W^{*T} \mathbf{X}_{:,j}} - \frac{\mathbf{X}_{:,k}^T (W^* \odot (1 - \mathbf{X}_{:,j}))}{W^{*T} (1 - \mathbf{X}_{:,j})} \right) = \frac{2^{p-2}}{2^{p-1}} - \frac{2^{p-2}}{2^{p-1}} = 0.$$

□

Causal Regularizer

Set feature j as treatment variable

$$\sum_{j=1}^p \left\| \frac{X_{-j}^T \cdot (W \odot I_j)}{W^T \cdot I_j} - \frac{X_{-j}^T \cdot (W \odot (1 - I_j))}{W^T \cdot (1 - I_j)} \right\|_2^2,$$

All features
excluding
treatment j

Sample
Weights

Indicator of
treatment
status

Causally Regularized Logistic Regression

$$\begin{aligned}
 \min \quad & \sum_{i=1}^n W_i \cdot \log(1 + \exp((1 - 2Y_i) \cdot (x_i \beta))), \\
 \text{s.t.} \quad & \sum_{j=1}^p \left\| \frac{X_{-j}^T \cdot (W \odot I_j)}{W^T \cdot I_j} - \frac{X_{-j}^T \cdot (W \odot (1 - I_j))}{W^T \cdot (1 - I_j)} \right\|_2^2 \leq \lambda_1, \\
 & W \geq 0, \quad \|W\|_2^2 \leq \lambda_2, \quad \|\beta\|_2^2 \leq \lambda_3, \quad \|\beta\|_1 \leq \lambda_4, \\
 & (\sum_{k=1}^n W_k - 1)^2 \leq \lambda_5,
 \end{aligned}$$

Sample
reweighted
logistic loss

Causal
Contribution

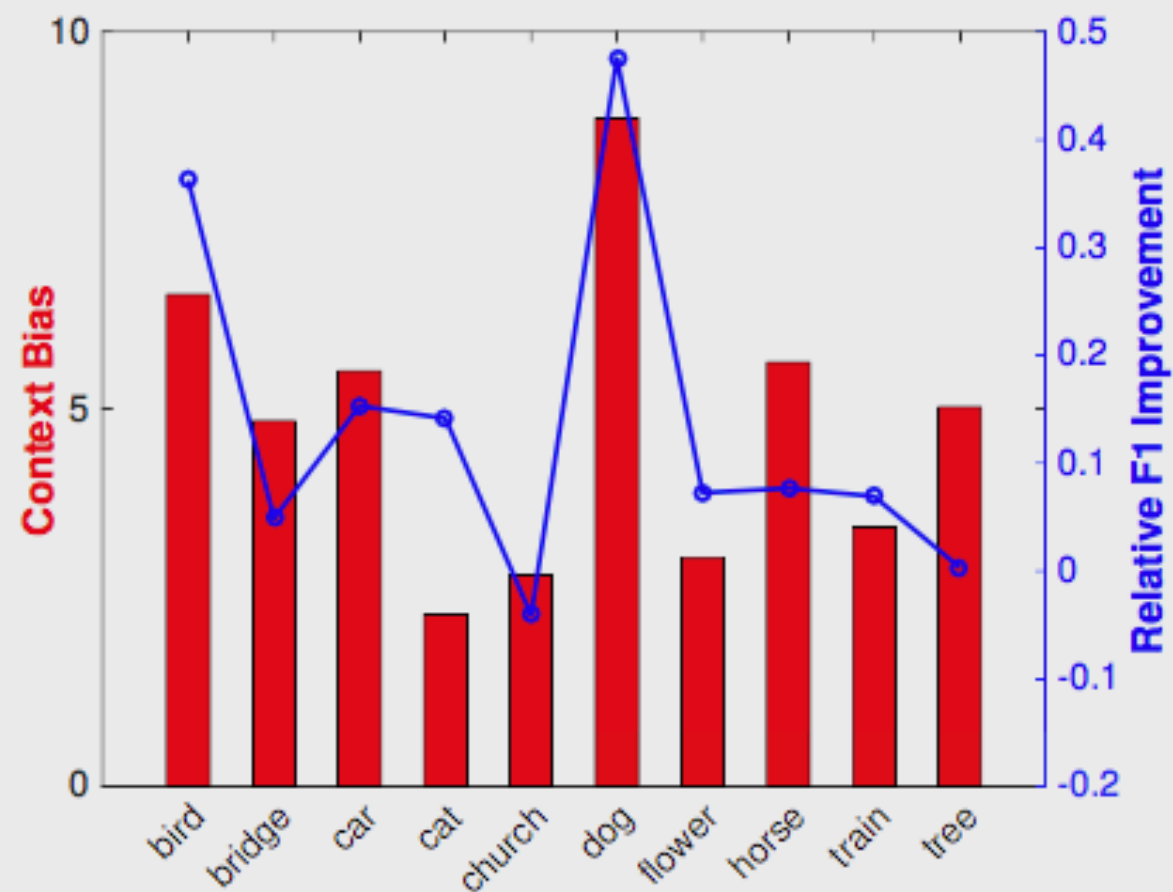
NICO - Non-I.I.D. Image Dataset with Contexts

- Data size of each class in NICO
 - Sample size: thousands for each class
 - Each superclass: 10,000 images
 - Sufficient for some basic neural networks (CNN)
- Samples with contexts in NICO

<i>Animal</i>	DATA SIZE	<i>Vehicle</i>	DATA SIZE
BEAR	1609	AIRPLANE	930
BIRD	1590	BICYCLE	1639
CAT	1479	BOAT	2156
COW	1192	BUS	1009
DOG	1624	CAR	1026
ELEPHANT	1178	HELICOPTER	1351
HORSE	1258	MOTORCYCLE	1542
MONKEY	1117	TRAIN	750
RAT	846	TRUCK	1000
SHEEP	918		



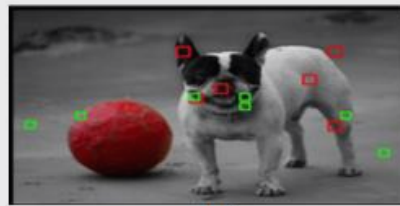
Experimental Result - insights



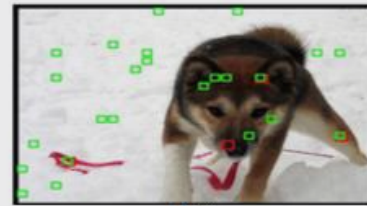
Experimental Result - insights



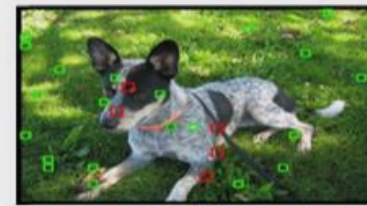
(a)



(b)



(c)



(d)



(e)



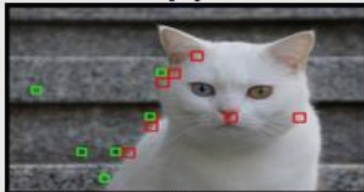
(f)



(g)



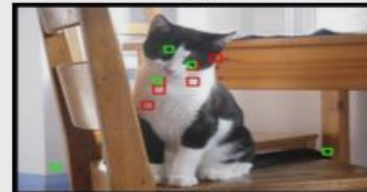
(h)



(i)



(j)



(k)



(l)



(m)



(n)



(o)



(p)

Stable Learning with *Continuous* Variables

Variable Decorrelation by Sample Reweighting:

$$\min_W \sum_{j=1}^p \left\| \mathbb{E}[\mathbf{X}_{:,j}^T \Sigma_W \mathbf{X}_{:,-j}] - \mathbb{E}[\mathbf{X}_{:,j}^T W] \mathbb{E}[\mathbf{X}_{:,-j}^T W] \right\|_2^2$$

Decorrelated Weighted Regression:

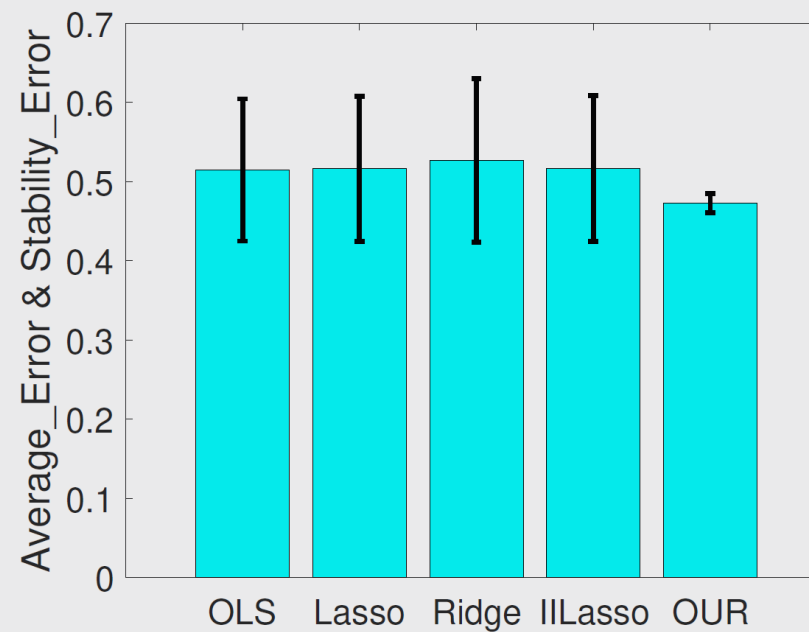
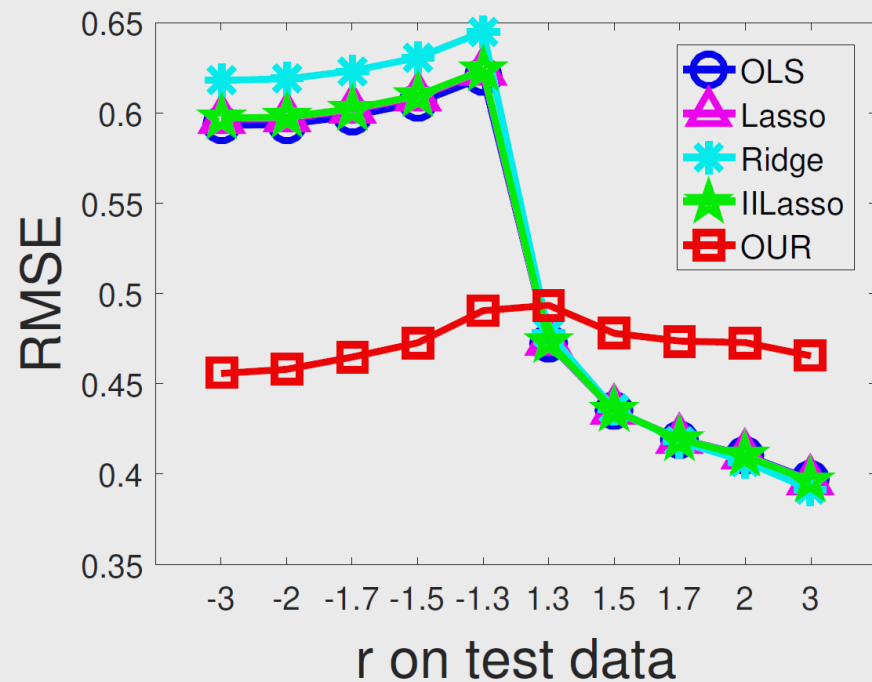
$$\min_{W, \beta} \sum_{i=1}^n W_i \cdot (Y_i - \mathbf{X}_i \beta)^2 \quad (12)$$

$$s.t. \quad \sum_{j=1}^p \left\| \mathbf{X}_{:,j}^T \Sigma_W \mathbf{X}_{:,-j} / n - \mathbf{X}_{:,j}^T W / n \cdot \mathbf{X}_{:,-j}^T W / n \right\|_2^2 < \lambda_2$$

$$|\beta|_1 < \lambda_1, \quad \frac{1}{n} \sum_{i=1}^n W_i^2 < \lambda_3,$$

$$\left(\frac{1}{n} \sum_{i=1}^n W_i - 1 \right)^2 < \lambda_4, \quad W \succeq 0,$$

Stable Learning with *Continuous* Variables



Stable Learning with *Differentiated* Variables

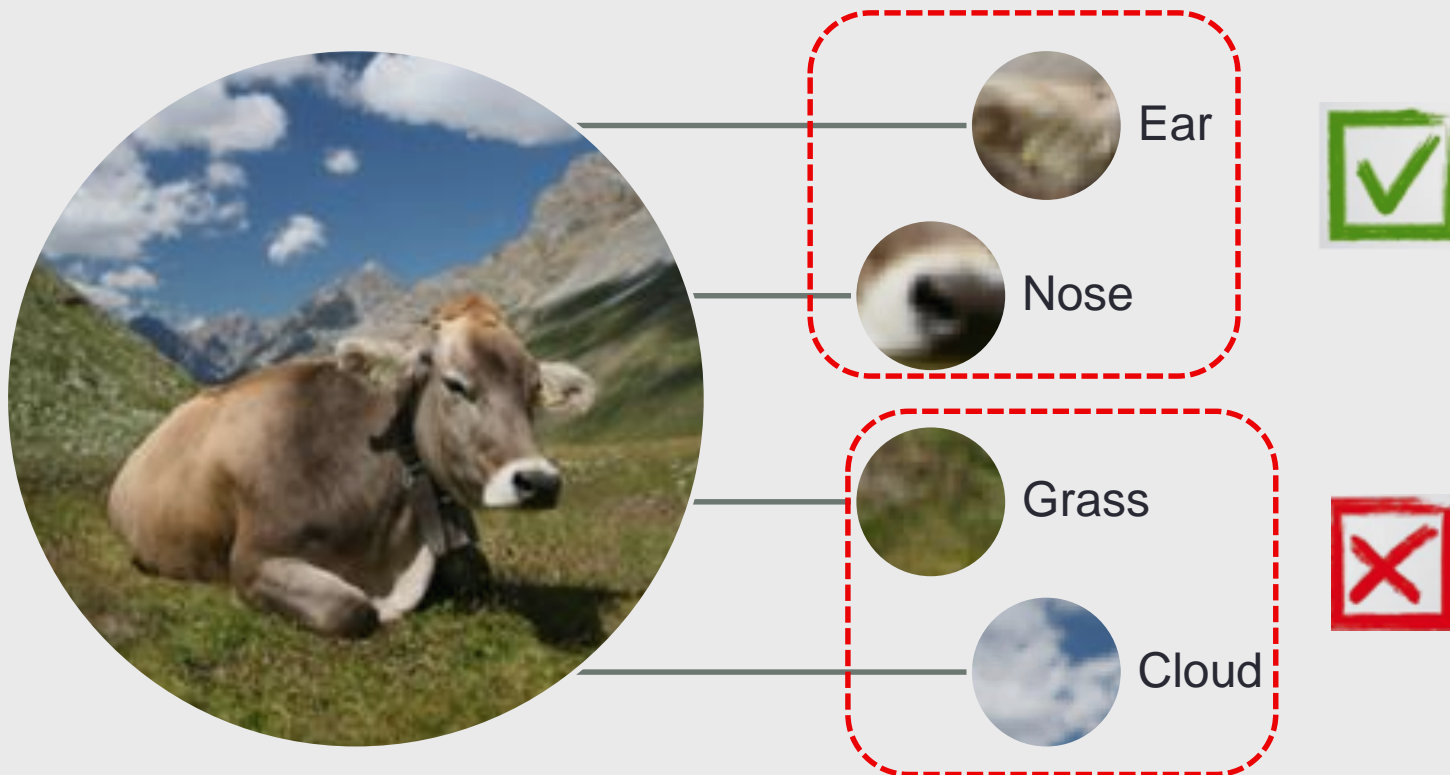
- More detailed analysis:

$$\begin{aligned}\hat{\beta}_{VOLS} &= \beta_V + \left(\frac{1}{n} \sum_{i=1}^n \mathbf{v}_i^T \mathbf{v}_i \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{v}_i^T g(\mathbf{s}_i) \right) \\ &\quad + \left(\frac{1}{n} \sum_{i=1}^n \mathbf{v}_i^T \mathbf{v}_i \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{v}_i^T \mathbf{s}_i \right) (\beta_S - \hat{\beta}_{SOLS}) \\ \hat{\beta}_{SOLS} &= \beta_S + \left(\frac{1}{n} \sum_{i=1}^n \mathbf{s}_i^T \mathbf{s}_i \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{s}_i^T g(\mathbf{s}_i) \right) \\ &\quad + \left(\frac{1}{n} \sum_{i=1}^n \mathbf{s}_i^T \mathbf{s}_i \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{s}_i^T \mathbf{v}_i \right) (\beta_V - \hat{\beta}_{VOLS})\end{aligned}$$

- We can focus on only the **spurious part** of correlation
- But how?
- **Leveraging the abundant sources of unlabeled data!**

Stable Learning with *Differentiated* Variables

ASSUMPTION 3. The variables $\mathbf{X} = \{X_1, X_2, \dots, X_p\}$ could be partitioned into k distinct groups $\mathbf{G}_1, \mathbf{G}_2, \dots, \mathbf{G}_k$. For $\forall i, j, i \neq j$ and $X_i, X_j \in \mathbf{G}_l, l \in \{1, 2, \dots, k\}$, we have $P_{X_i X_j}^e = P_{X_i X_j}$.



Clustering?

Stable Learning with *Differentiated* Variables

- Feature Partition by Stable Correlation Clustering
 - Define the dissimilarity of two variables:

$$Dis(X_i, X_j) = \sqrt{\frac{1}{M-1} \sum_{l=1}^M \left(Corr(X_i^l, X_j^l) - Ave_Corr(X_i, X_j) \right)^2},$$

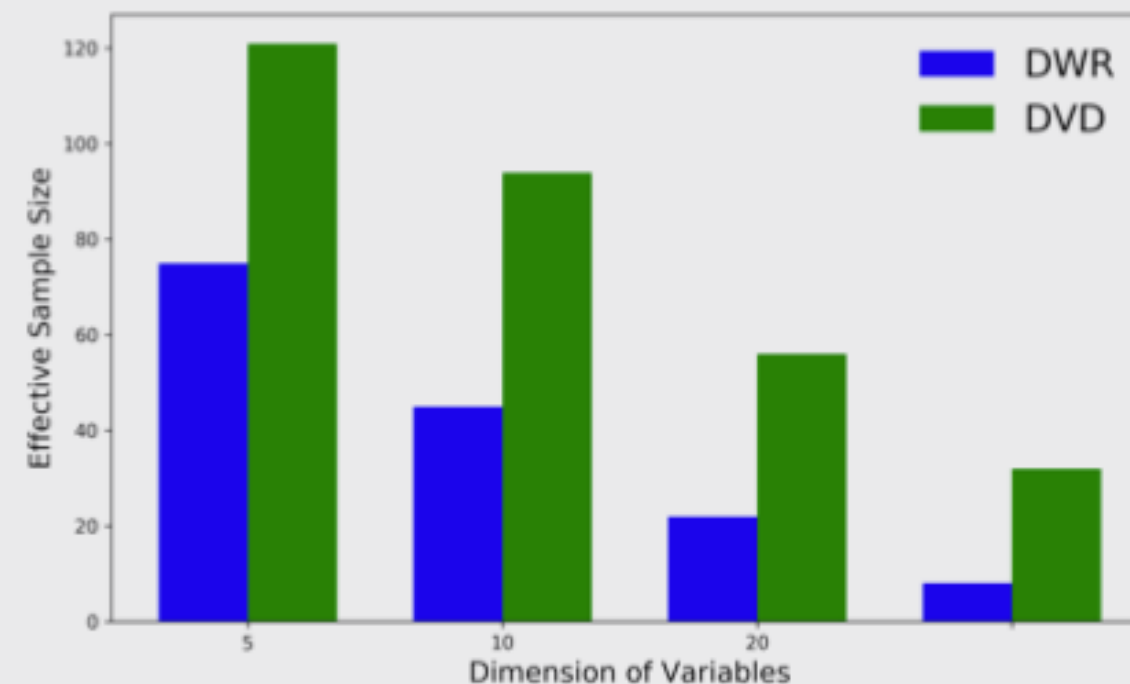
- Remove the correlation between variables via sample reweighting:

$$\min_W \sum_{i \neq j} \mathbb{I}(i, j) \left\| (\mathbf{X}_{:,i}^T \Sigma_W \mathbf{X}_{:,j} / n - \mathbf{X}_{:,i}^T W / n \cdot \mathbf{X}_{:,j}^T W / n) \right\|_2^2$$

$$\text{s.t. } \frac{1}{n} \sum_{i=1}^n W_i^2 < \gamma_1, \quad \left(\frac{1}{n} \sum_{i=1}^n W_i - 1 \right)^2 < \gamma_2, \quad W \geq 0$$

Stable Learning with *Differentiated* Variables

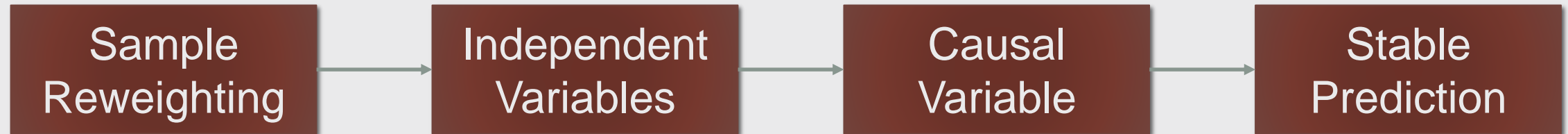
Scenario 1: varying sample size n						
n, p_{v_b}, r	$n = 120, p_{v_b} = p * 0.2, r = 1.9$			$n = 160, p_{v_b} = p * 0.2, r = 1.9$		
Methods	β_Error	Average_Error	Stability_Error	β_Error	Average_Error	Stability_Error
OLS	1.988	0.470	0.087	1.870	0.489	0.105
Lasso	2.021	0.476	0.092	1.905	0.494	0.110
IIIasso	2.035	0.475	0.094	1.920	0.498	0.113
DWR	2.012	0.545	0.099	1.991	0.502	0.076
Our	1.892	0.469	0.040	1.741	0.489	0.050
Scenario 2: varying number of unstable variables p_{v_b}						
n, p_{v_b}, r	$n = 200, p_{v_b} = p * 0.2, r = 1.9$			$n = 200, p_{v_b} = p * 0.3, r = 1.9$		
Methods	β_Error	Average_Error	Stability_Error	β_Error	Average_Error	Stability_Error
OLS	1.839	0.522	0.121	2.128	0.563	0.179
Lasso	1.876	0.529	0.129	2.176	0.571	0.186
IIIasso	1.894	0.538	0.149	2.196	0.575	0.191
DWR	1.656	0.485	0.081	1.881	0.469	0.092
Our	1.369	0.476	0.042	1.641	0.460	0.064
Scenario 3: varying bias rate r on training data						
n, p_{v_b}, r	$n = 200, p_{v_b} = p * 0.2, r = 1.6$			$n = 200, p_{v_b} = p * 0.2, r = 1.8$		
Methods	β_Error	Average_Error	Stability_Error	β_Error	Average_Error	Stability_Error
OLS	1.296	0.452	0.064	1.780	0.510	0.117
Lasso	1.321	0.455	0.067	1.812	0.516	0.123
IIIasso	1.339	0.457	0.070	1.829	0.519	0.125
DWR	1.153	0.457	0.033	1.262	0.458	0.035
Our	1.236	0.463	0.021	1.236	0.450	0.023



Effective Sample Size

From *Causal* problem to *Learning* problem

- Previous logic:



- More direct logic:



Interpretation from Statistical Learning perspective

- Consider the linear regression with misspecification bias

$$y = x^\top \bar{\beta}_{1:p} + \bar{\beta}_0 + b(x) + \epsilon$$

Goes to infinity when perfect collinearity exists!

Bias term with bound $b(x) \leq \delta$

- By accurately estimating $\bar{\beta}$ with the property that $b(x)$ is uniformly small for all x , we can achieve stable learning.
- However, the estimation error caused by misspecification term can be as bad as $\|\hat{\beta} - \bar{\beta}\|_2 \leq 2(\delta/\gamma) + \delta$, where γ^2 is the smallest eigenvalue of centered covariance matrix.

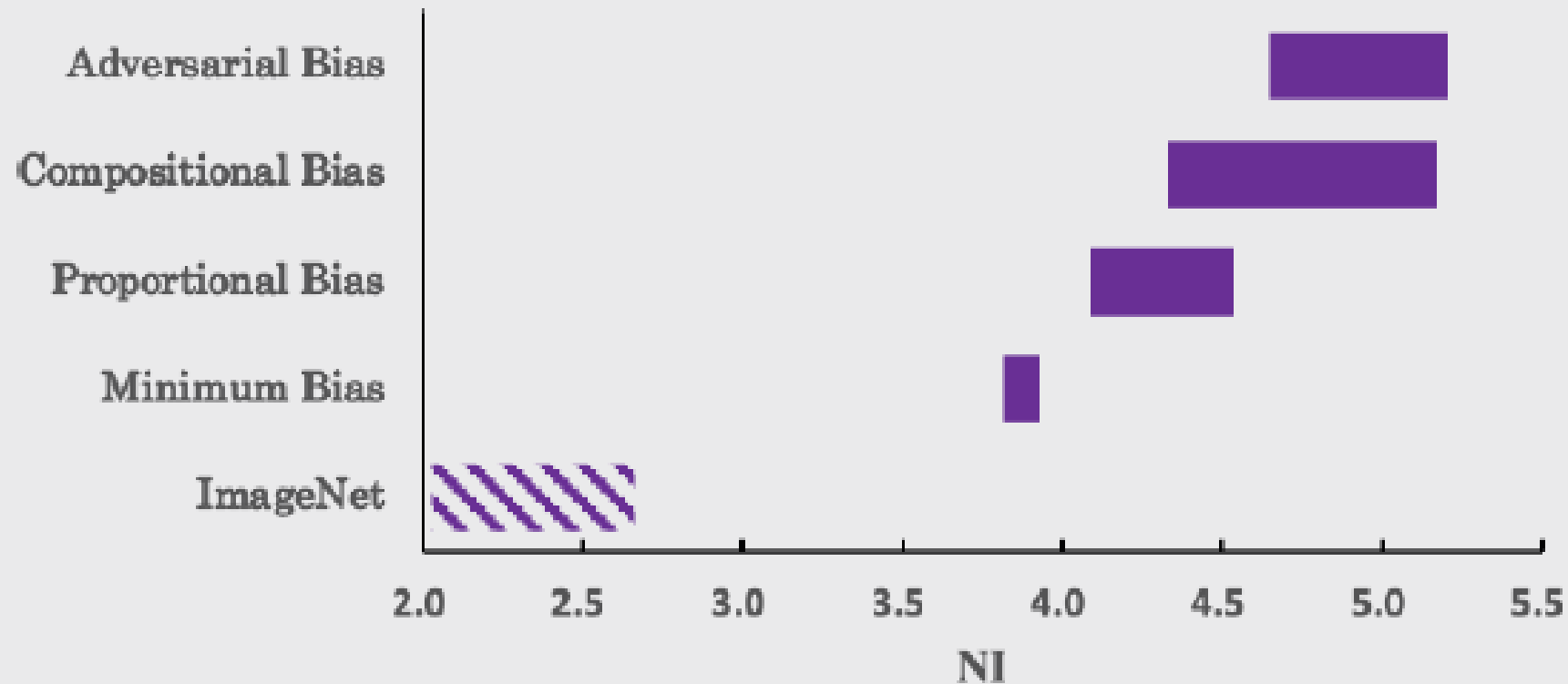
NICO - Non-I.I.D. Image Dataset with Contexts

- Data size of each class in NICO
 - Sample size: thousands for each class
 - Each superclass: 10,000 images
 - Sufficient for some basic neural networks (CNN)
- Samples with contexts in NICO

<i>Animal</i>	DATA SIZE	<i>Vehicle</i>	DATA SIZE
BEAR	1609	AIRPLANE	930
BIRD	1590	BICYCLE	1639
CAT	1479	BOAT	2156
COW	1192	BUS	1009
DOG	1624	CAR	1026
ELEPHANT	1178	HELICOPTER	1351
HORSE	1258	MOTORCYCLE	1542
MONKEY	1117	TRAIN	750
RAT	846	TRUCK	1000
SHEEP	918		



NICO - Non-I.I.D. Image Dataset with Contexts



<http://nico.thumedia lab.com/>

Conclusions

- Why can't the current AI generalize well to unknown environments?

Know What, but don't know Why

知其 然 ， 但不知其 所以然

Correlation

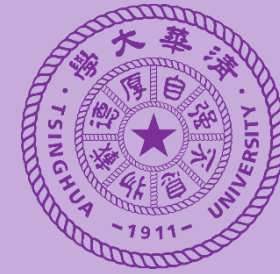
Causality

Stable Learning: Try to promote the convergence of causal inference and machine learning.

Reference

- Hao Zou, Peng Cui, Bo Li, Zheyang Shen, Jianxin Ma, Hongxia Yang, Yue He. Counterfactual Prediction for Bundle Treatments. **NeurIPS**, 2020.
- Zheyang Shen, Peng Cui, Jiashuo Liu, Tong Zhang, Bo Li and Zhitang Chen. Stable Learning via Differentiated Variable Decorrelation. **KDD**, 2020.
- Yue He, Peng Cui, Jianxin Ma, Zou Hao, Xiaowei Wang, Hongxia Yang and Philip S. Yu. Learning Stable Graphs from Multiple Environments with Selection Bias. **KDD**, 2020.
- Renzhe Xu, Peng Cui, Kun Kuang, Bo Li, Linjun Zhou, Zheyang Shen and Wei Cui. Algorithmic Decision Making with Conditional Fairness. **KDD**, 2020.
- Yue He, Zheyang Shen, Peng Cui. Towards Non-I.I.D. Image Classification: A Dataset and Baselines. **Pattern Recognition**, 2020.
- Zheyang Shen, Peng Cui, Tong Zhang. Stable Learning via Sample Reweighting. **AAAI**, 2020.
- Kun Kuang, Ruoxuan Xiong, Peng Cui, Susan Athey, Bo Li. Stable Prediction with Model Misspecification and Agnostic Distribution Shift. **AAAI**, 2020.
- Hao Zou, Kun Kuang, Boqi Chen, Peng Cui, Peixuan Chen. Focused Context Balancing for Robust Offline Policy Evaluation. **KDD**, 2019.
- Jianxin Ma, Peng Cui, Kun Kuang, Xin Wang, Wenwu Zhu. Disentangled Graph Convolutional Networks. **ICML**, 2019.
- Kun Kuang, Peng Cui, Susan Athey, Ruoxuan Li, Bo Li. Stable Prediction across Unknown Environments. **KDD**, 2018.
- Zheyang Shen, Peng Cui, Kun Kuang, Bo Li. Causally Regularized Learning on Data with Agnostic Bias. **ACM Multimedia**, 2018.
- Kun Kuang, Peng Cui, Bo Li, Shiqiang Yang. Estimating Treatment Effect in the Wild via Differentiated Confounder Balancing. **KDD**, 2017.
- Kun Kuang, Peng Cui, Bo Li, Shiqiang Yang. Treatment Effect Estimation with Data-Driven Variable Decomposition. **AAAI**, 2017.

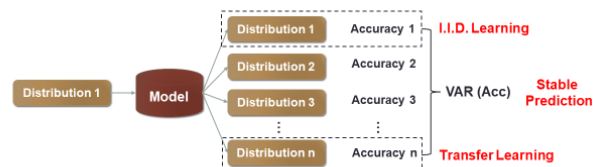
Thanks!



Peng Cui
 cuip@tsinghua.edu.cn
<http://pengcui.thumedia lab.com>

Research Problems

- Comes down to the Model

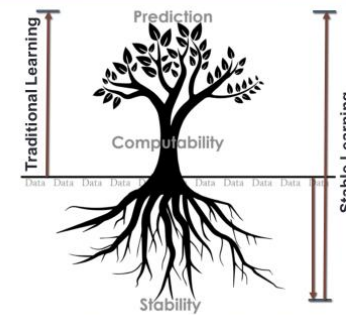


Stable Learning

Prediction
Performance

Learning Process

True Model



Bin Yu (2016), Three Principles of Data Science: predictability, computability, stability

