



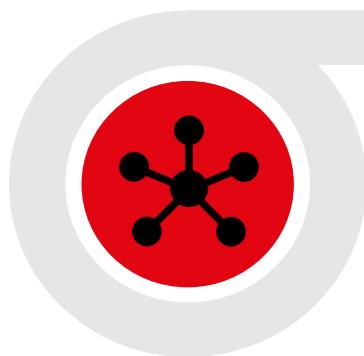
Disentangling Treatment Assignment Bias in Counterfactual Outcome Prediction and Biomarker Identification

First Authors: Michael Vollenweider¹, Manuel Schürch^{1,2,4} [michavol@ethz.ch, manuel_schurch@dfci.harvard.edu]

Co-Authors: Chiara Rohrer¹, Gabriele Gut^{2,3}, Michael Krauthammer^{2,3}, Andreas Wicki^{2,3}

¹ETHZ: Federal Institute of Technology Zurich; ²University of Zurich; ³University Hospital Zurich; ⁴Harvard University

Outline



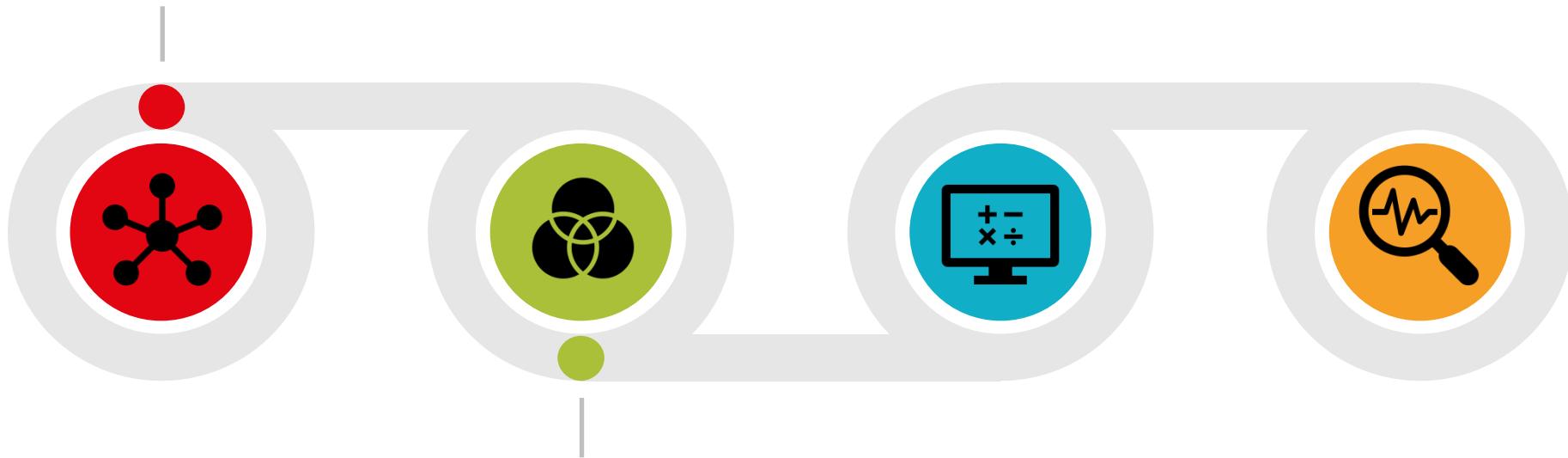
Outline

1. Treatment Effect Prediction &
Selection Bias



Outline

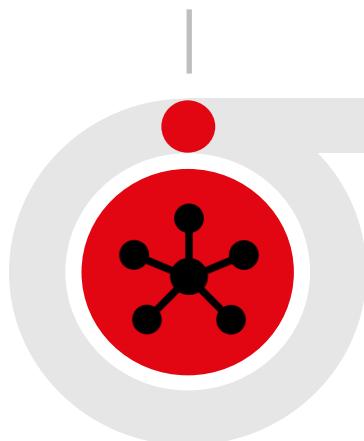
1. Treatment Effect Prediction &
Selection Bias



2. Information-Theoretic
Modelling of Bias

Outline

1. Treatment Effect Prediction &
Selection Bias



3. Simulations
& Experiments

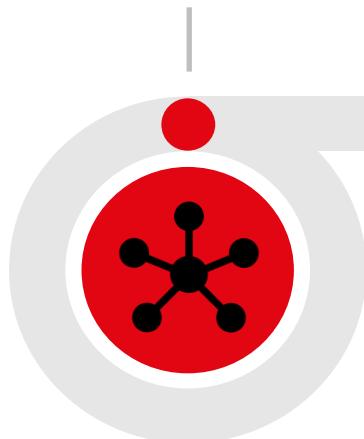


2. Information-Theoretic
Modelling of Bias

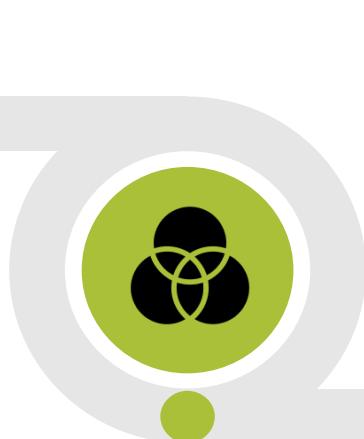


Outline

1. Treatment Effect Prediction &
Selection Bias



2. Information-Theoretic
Modelling of Bias



3. Simulations
& Experiments



4. Results



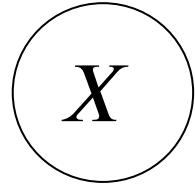


Treatment Effect Model



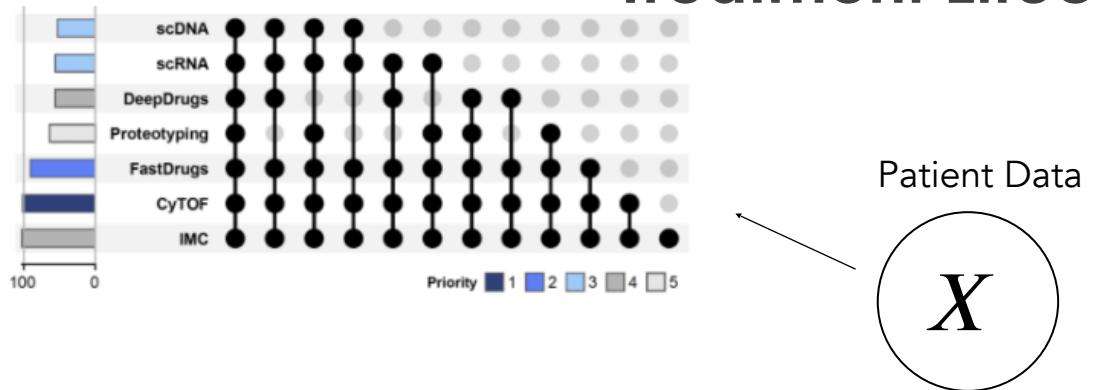
Treatment Effect Model

Patient Data



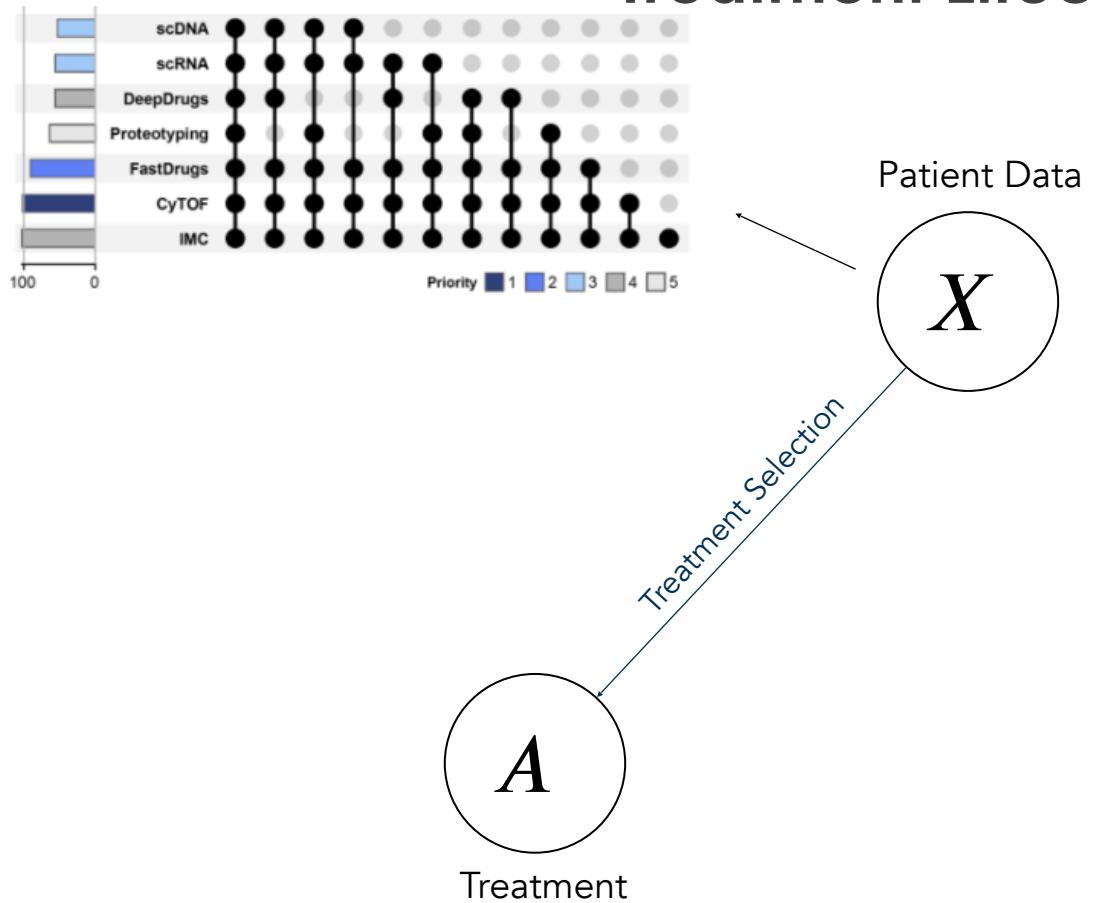


Treatment Effect Model



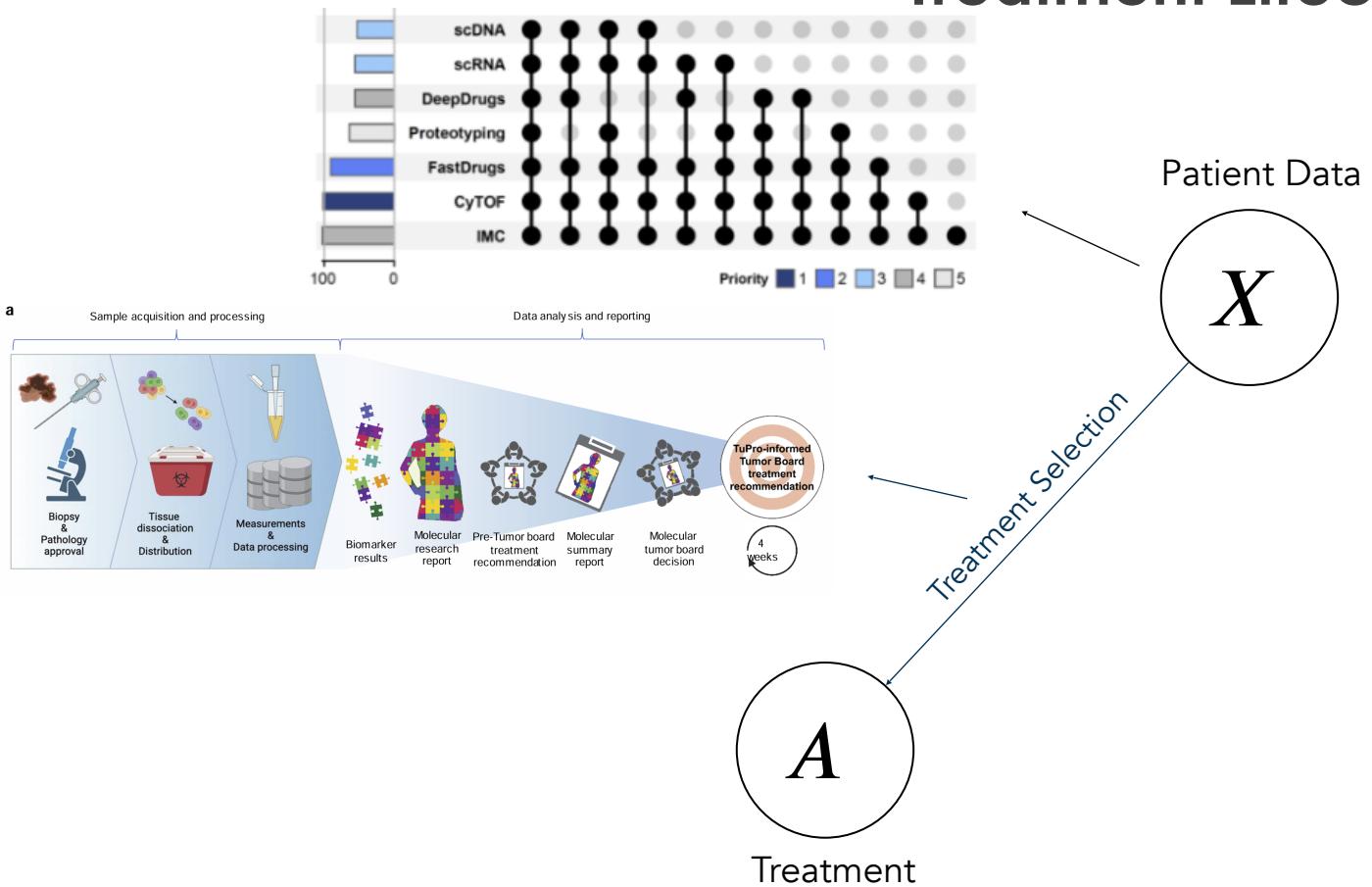


Treatment Effect Model



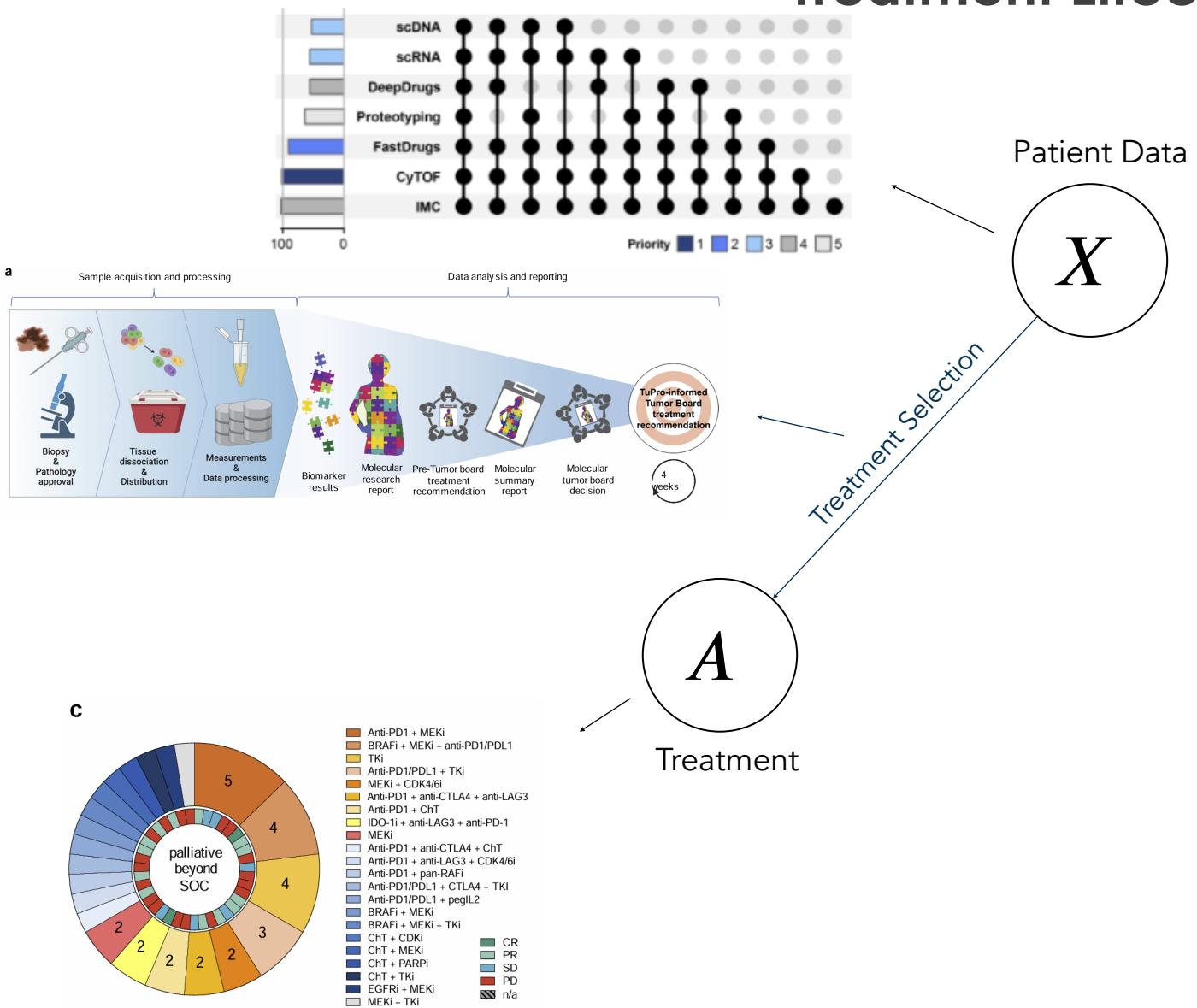


Treatment Effect Model



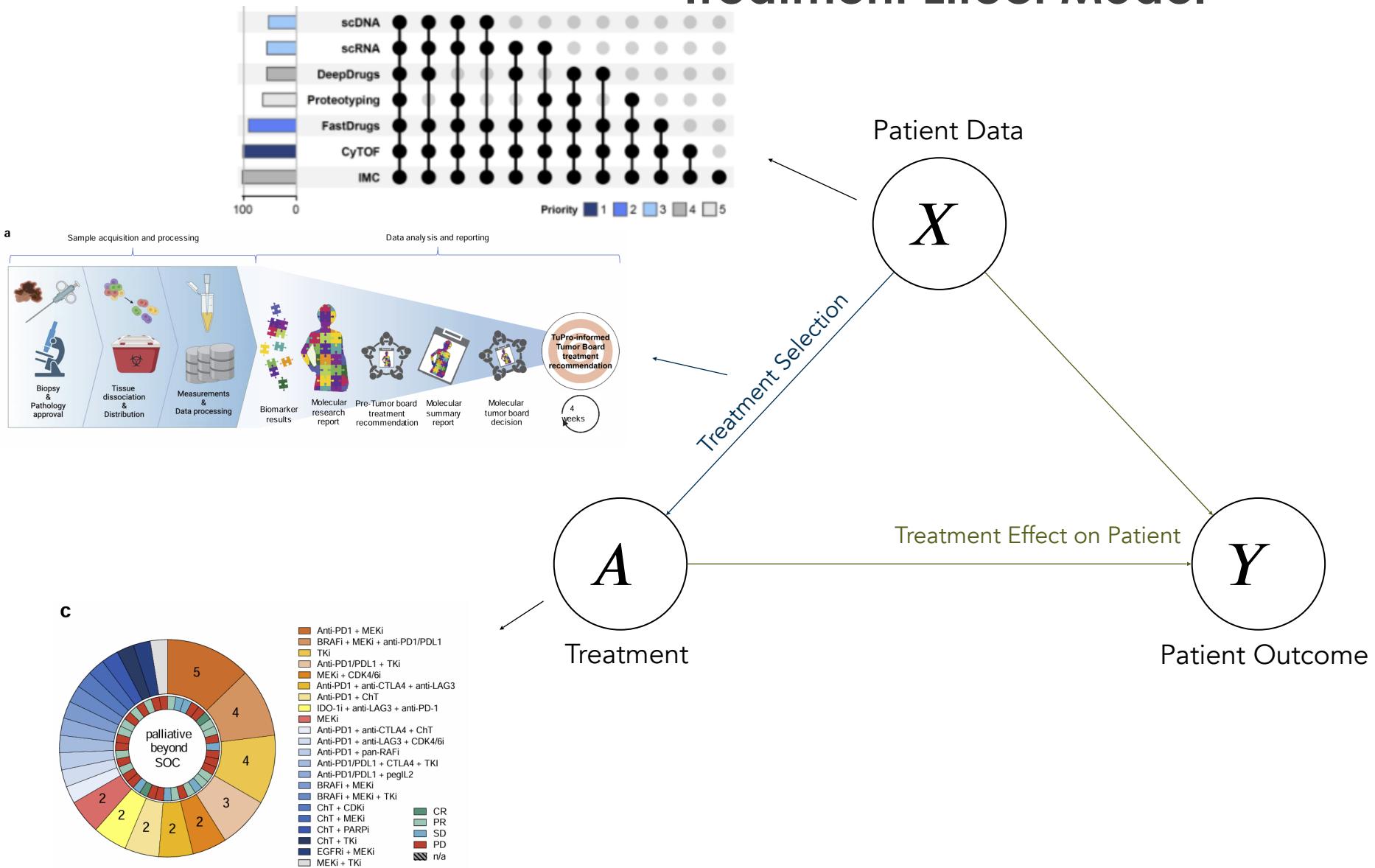


Treatment Effect Model



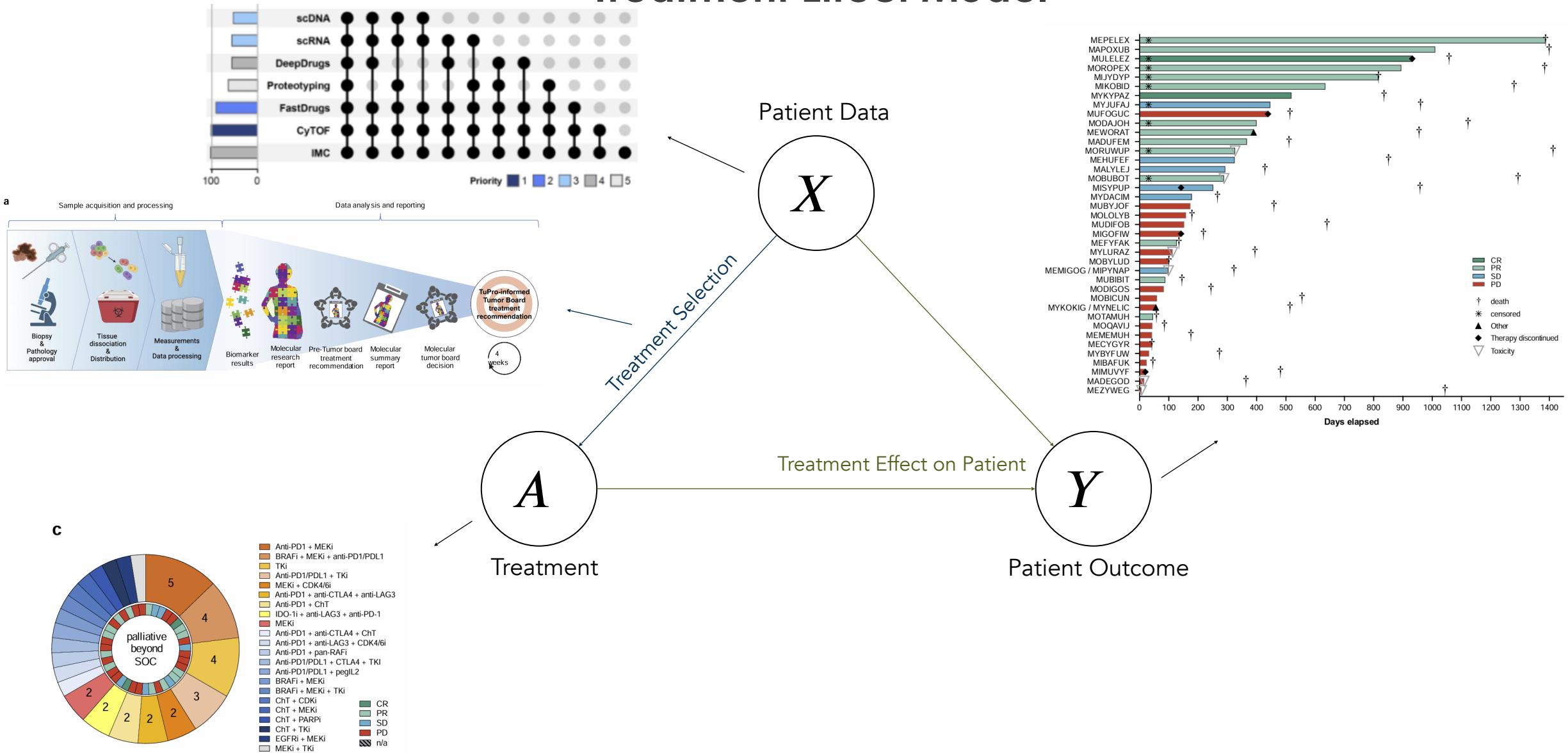


Treatment Effect Model





Treatment Effect Model





CATE Prediction

Potential Outcomes Framework by Rubin (2005):

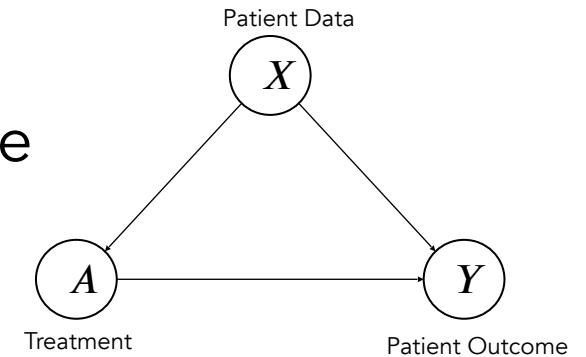


CATE Prediction

Potential Outcomes Framework by Rubin (2005):

Conditional Expected Potential Outcome

$$\mu_a(x) = E[Y_i^a \mid X_i = x]$$



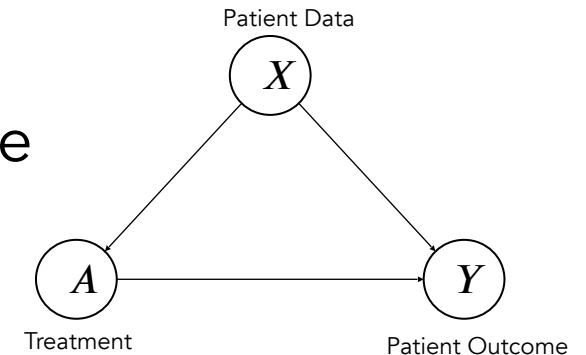


CATE Prediction

Potential Outcomes Framework by Rubin (2005):

Conditional Expected Potential Outcome

$$\mu_a(x) = E[Y_i^a | X_i = x]$$

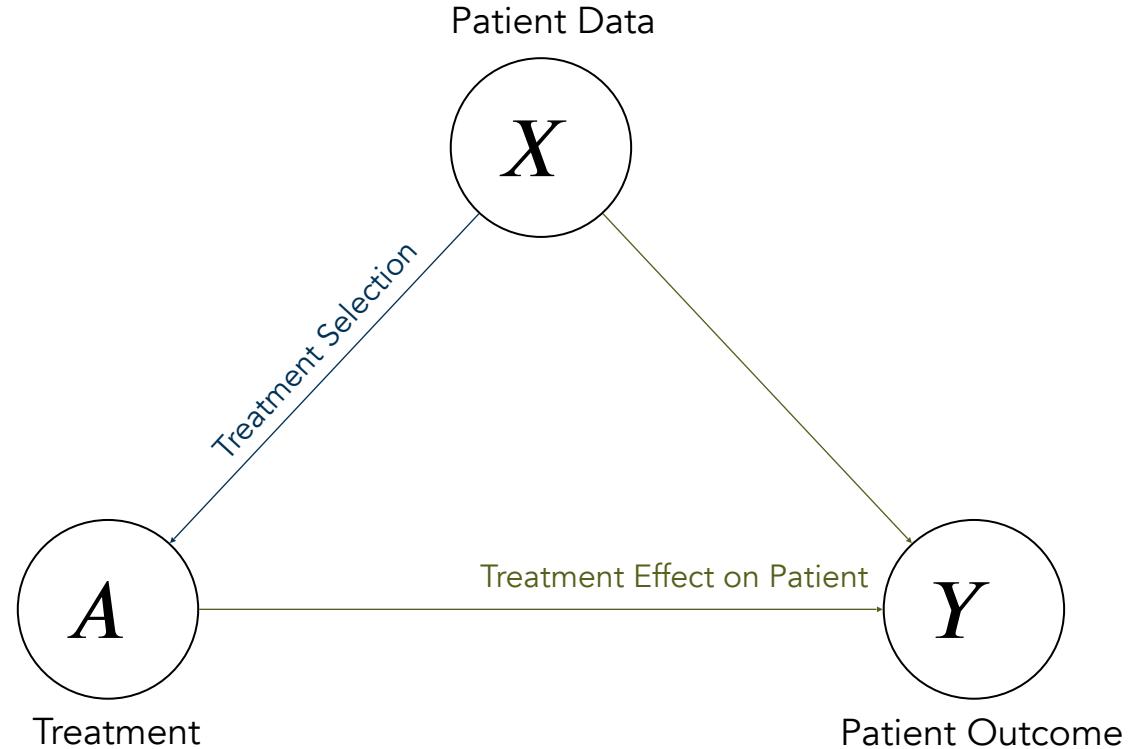


Conditional Average Treatment Effect - CATE

$$\tau(x) = E[Y_i^1 - Y_i^0 | X_i = x] = \mu_1(x) - \mu_0(x)$$



Treatment Effect Model

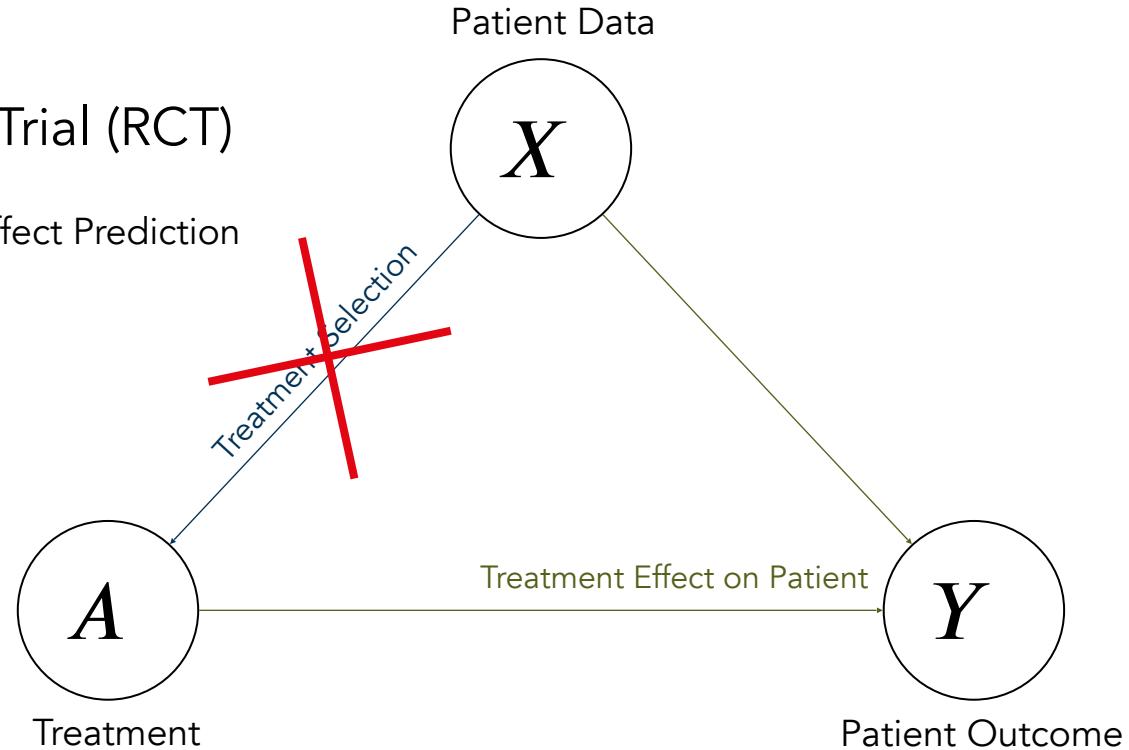




Treatment Effect Model

Randomized Controlled Trial (RCT)

→ Gold Standard for Treatment Effect Prediction

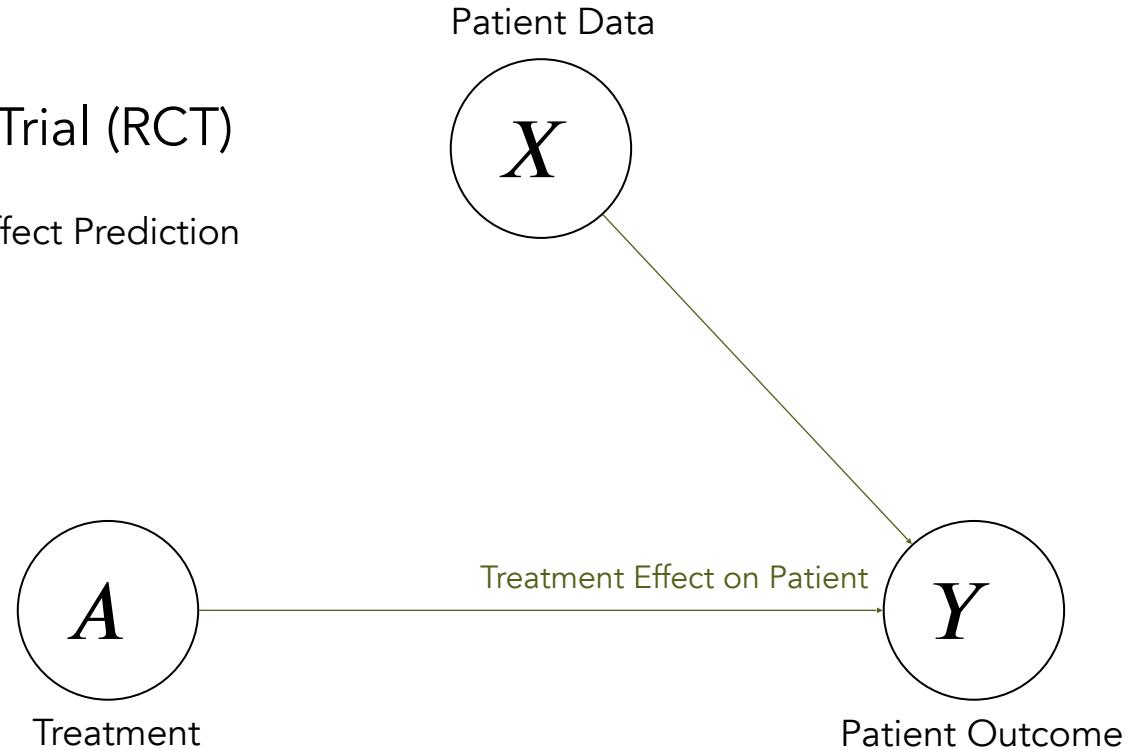




Treatment Effect Model

Randomized Controlled Trial (RCT)

→ Gold Standard for Treatment Effect Prediction



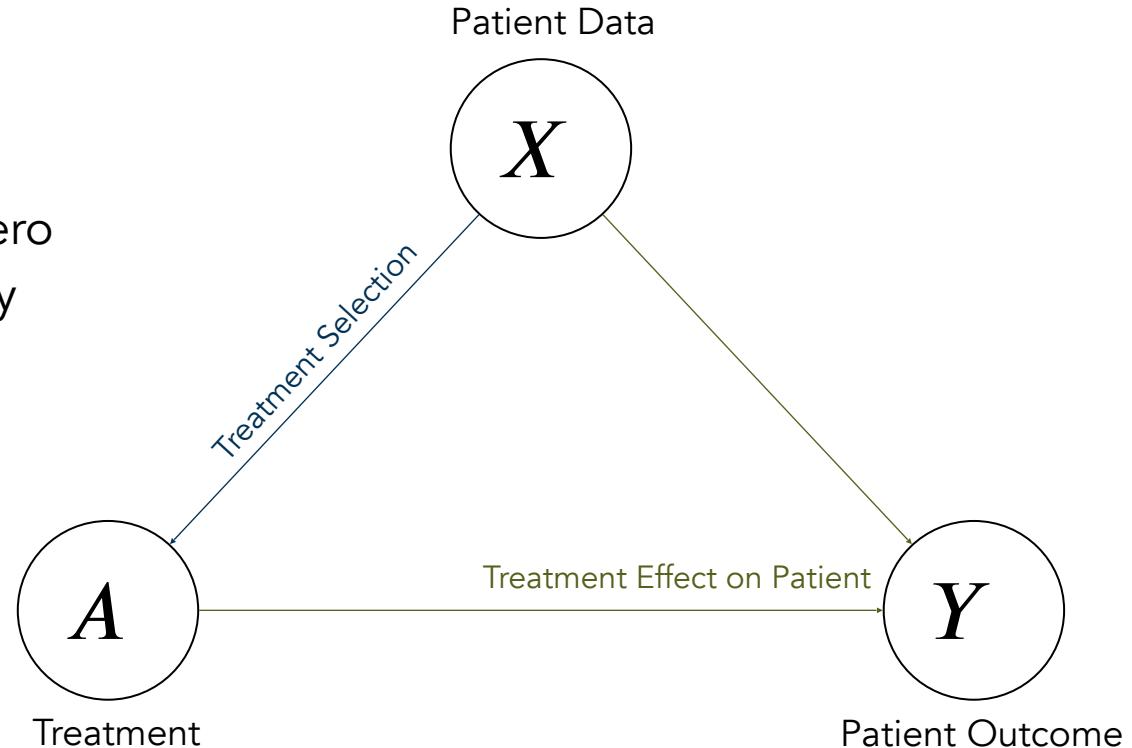


Treatment Effect Model

Overlap Assumption:

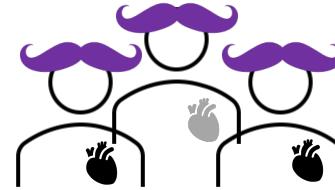
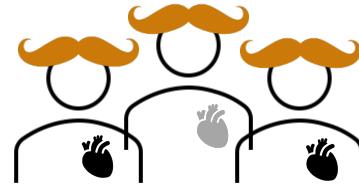
"Each patient has a non-zero probability of receiving any treatment"

(Rosenbaum and Rubin, 1983)



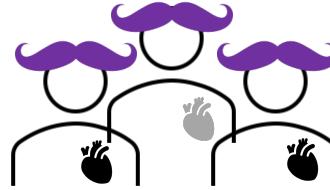
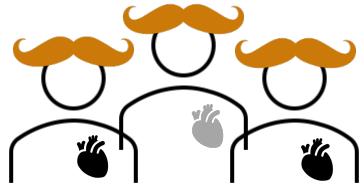


Types of Bias





Types of Bias



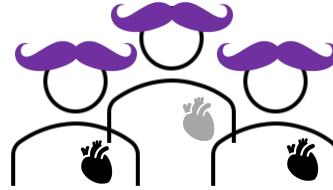
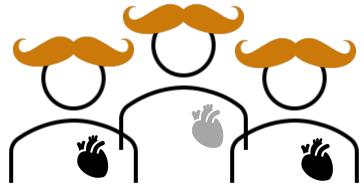
Treatment 1

Treatment selection
based on hair color.

Treatment 2



Types of Bias



Treatment 1

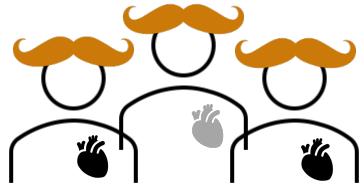
Treatment 2

Treatment selection
based on hair color.

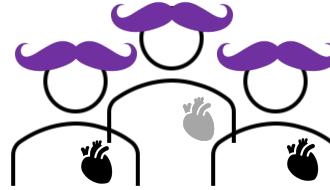
Overlap Assumption not Fulfilled!



Types of Bias



Treatment 1



Treatment 2

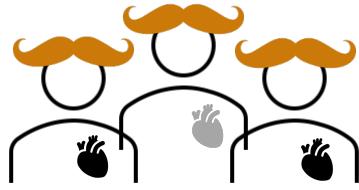
Treatment selection
based on hair color.

Overlap Assumption not Fulfilled!

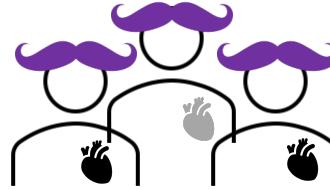
But does it matter?



Types of Bias



Treatment 1



Treatment 2

Treatment selection
based on hair color.

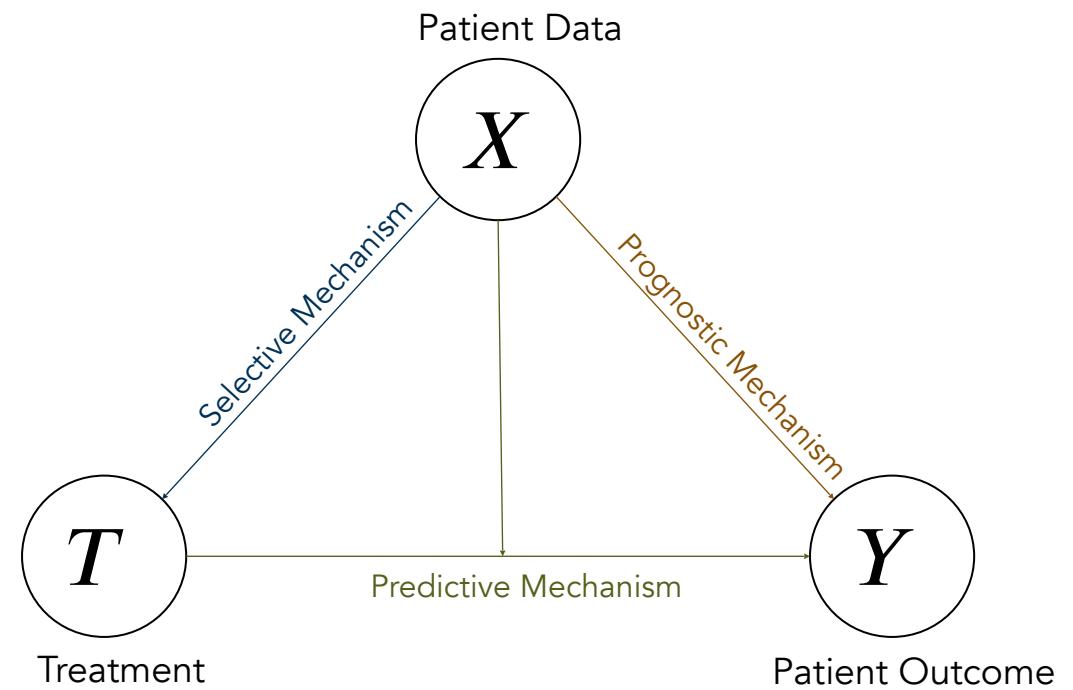
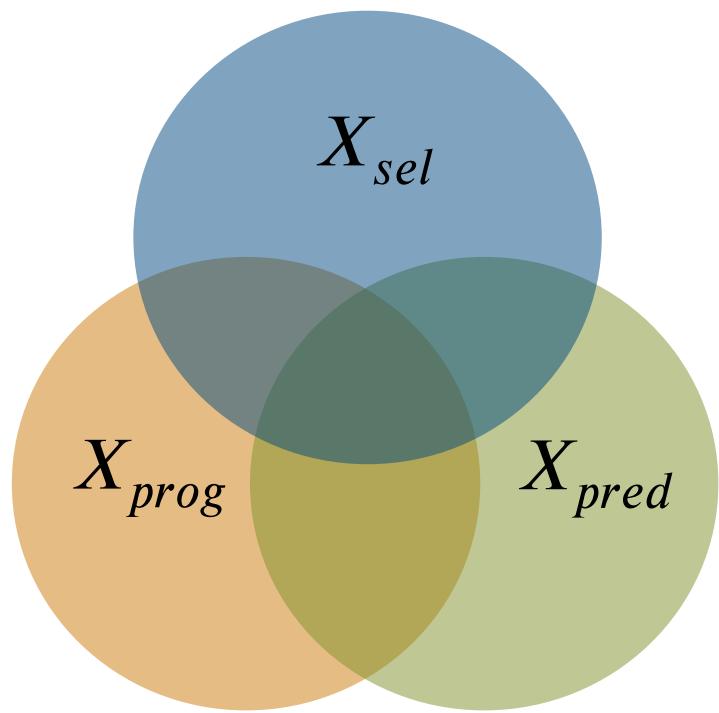
Overlap Assumption not Fulfilled!

But does it matter?

With regard to 'relevant' features,
the study is still randomized.



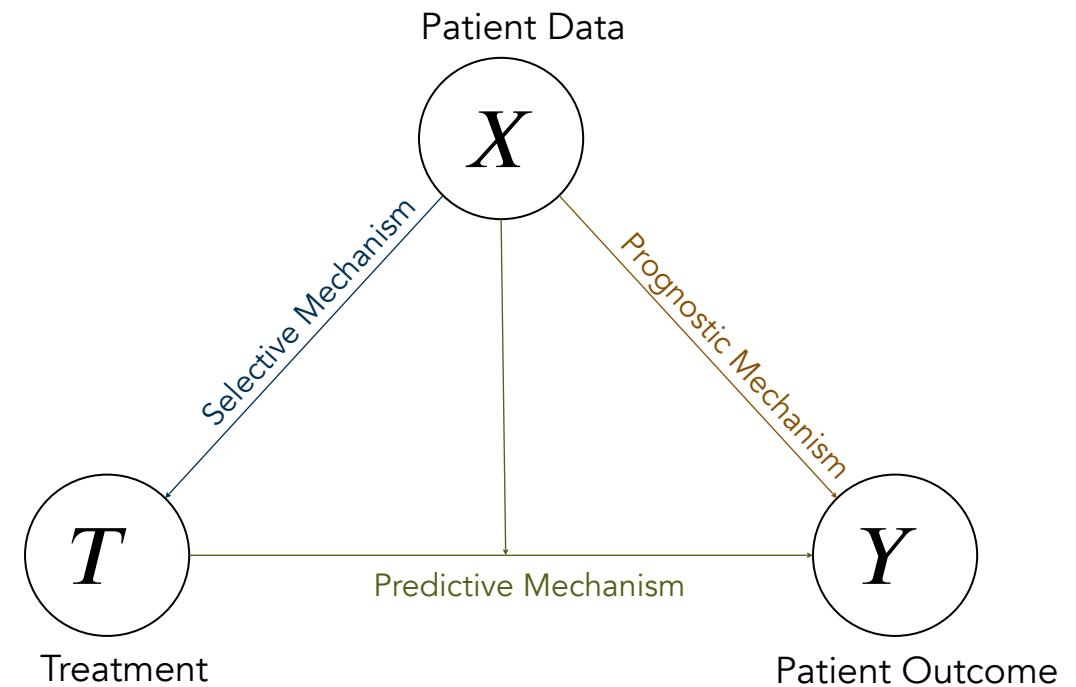
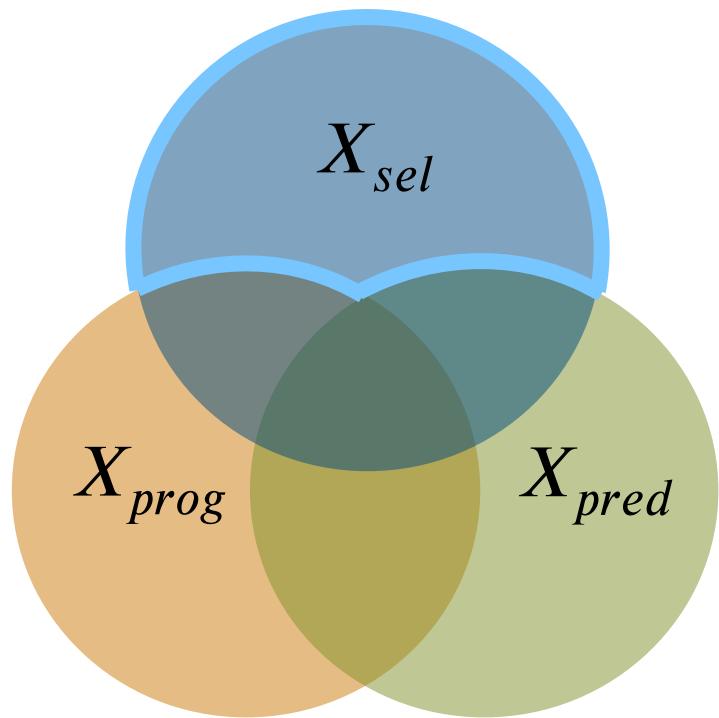
Irrelevant Features





Irrelevant Features

Features that do not have any impact on the patient outcome but still introduce treatment selection bias





Data Model

How can we formalize the
type and strength of
treatment selection bias?



Data Model

Definition 1 (Z-Bias) *The Z-bias for any treatment assignment $A^\pi \sim \pi(X)$ and random variable Z is defined as*

$$B_Z^\pi = \frac{\mathbb{I}(A^\pi; Z)}{\mathbb{H}[A^\pi]},$$

which measures the degree of bias of treatment assignment A^π with respect to Z .

❖ **Observable Bias B_X^π**

Quantifies how much information all patient characteristics provide about treatment selection. This bias can also be seen to quantify the degree of violation of the overlap assumption in causal inference.

❖ **Outcome Biases B_{Y0}^π and B_{Y1}^π**

Describe to what extent the potential outcomes of the treatments $a = 0$ and $a = 1$ determine the treatment decision.

❖ **Treatment Effect Bias B_{Y1-Y0}^π**

Quantifies how much the treatment decision is biased with respect to the true individualized treatment effect.

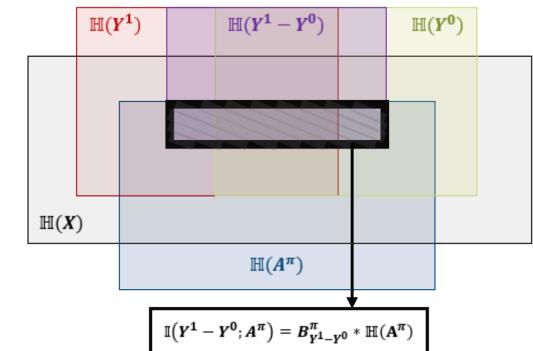
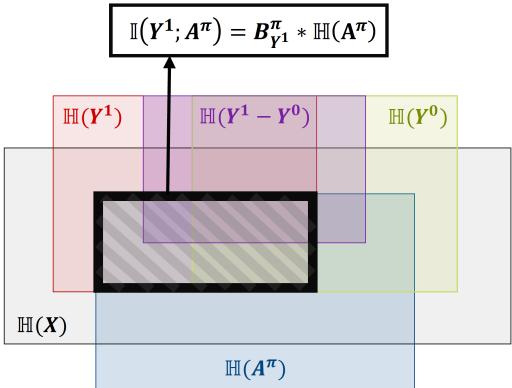


Data Model

Definition 1 (Z-Bias) *The Z-bias for any treatment assignment $A^\pi \sim \pi(X)$ and random variable Z is defined as*

$$B_Z^\pi = \frac{\mathbb{I}(A^\pi; Z)}{\mathbb{H}[A^\pi]},$$

which measures the degree of bias of treatment assignment A^π with respect to Z .



❖ **Observable Bias B_X^π**

Quantifies how much information all patient characteristics provide about treatment selection. This bias can also be seen to quantify the degree of violation of the overlap assumption in causal inference.

❖ **Outcome Biases $B_{Y^0}^\pi$ and $B_{Y^1}^\pi$**

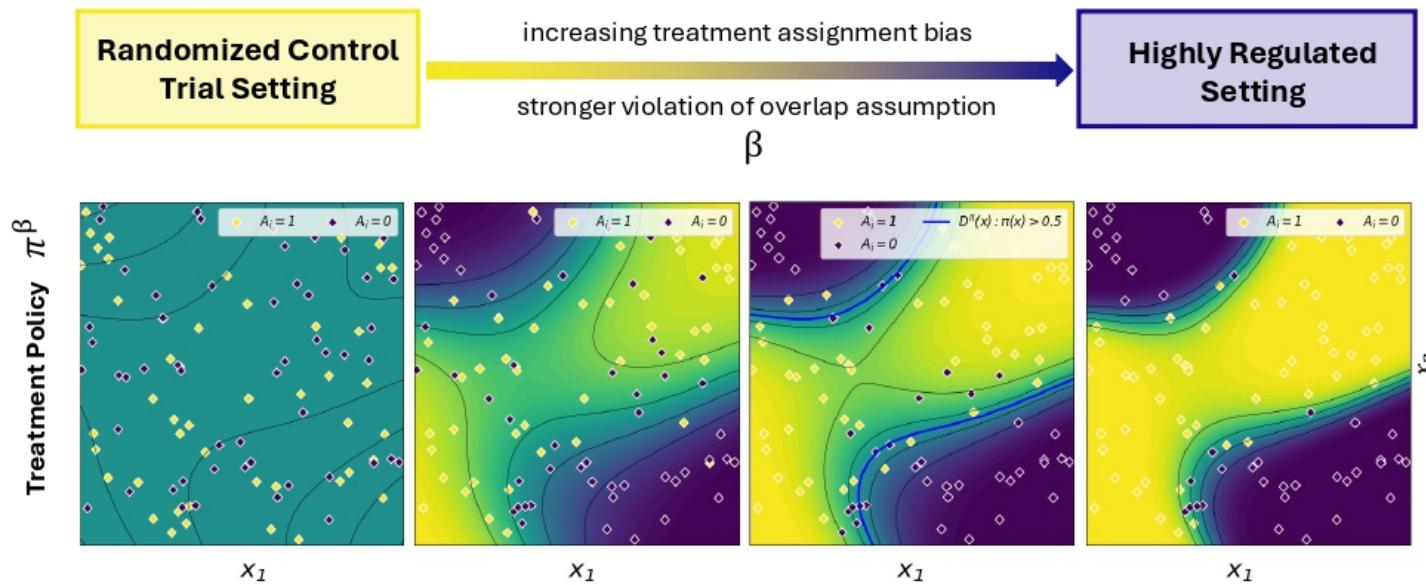
Describe to what extent the potential outcomes of the treatments $a = 0$ and $a = 1$ determine the treatment decision.

❖ **Treatment Effect Bias $B_{Y^1 - Y^0}^\pi$**

Quantifies how much the treatment decision is biased with respect to the true individualized treatment effect.



Policy Simulation



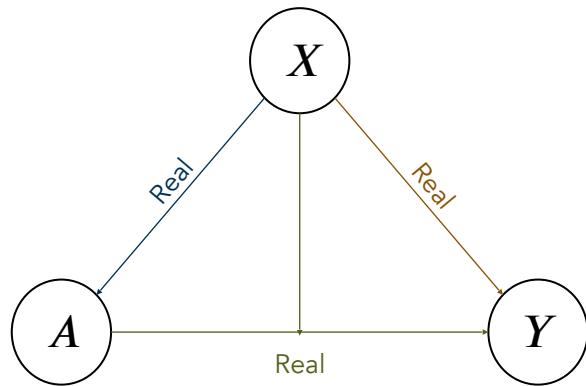
Definition 4 (Z-Policy) For a given random variable Z and parameter β , the Z-policy is defined as

$$\pi_Z^\beta(x) := \sigma(\beta Z(x)), \quad \sigma(x) := \frac{e^x}{1+e^x}$$



Simulations

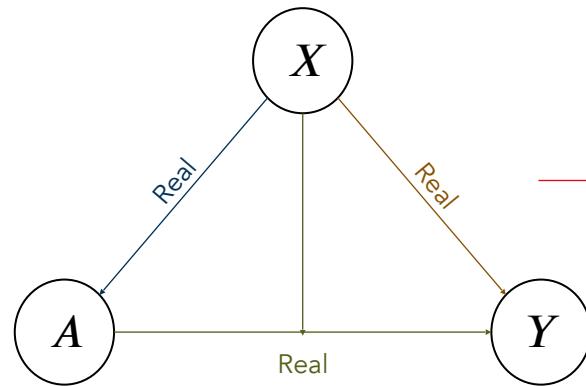
No Simulation
(Clinical Data)





Simulations

No Simulation
(Clinical Data)

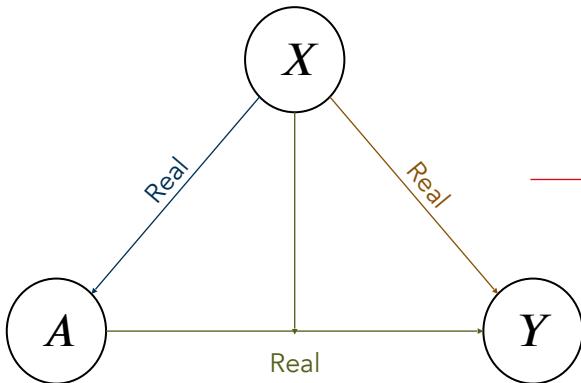


No ground-truth effect available!



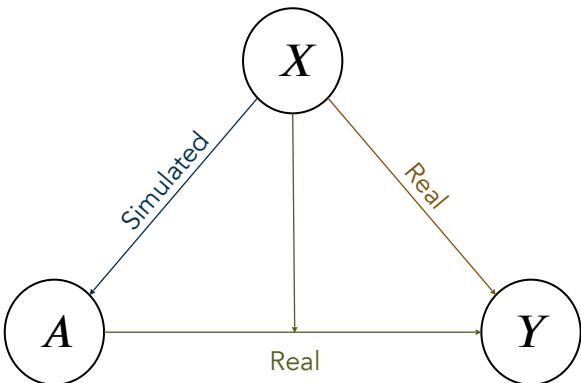
Simulations

No Simulation
(Clinical Data)



No ground-truth effect available!

A-Simulation

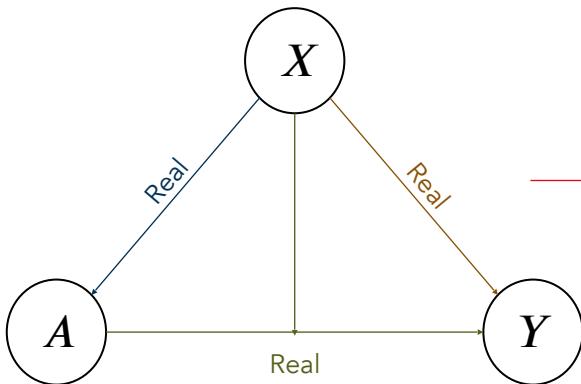


Less Realistic
More Evaluation



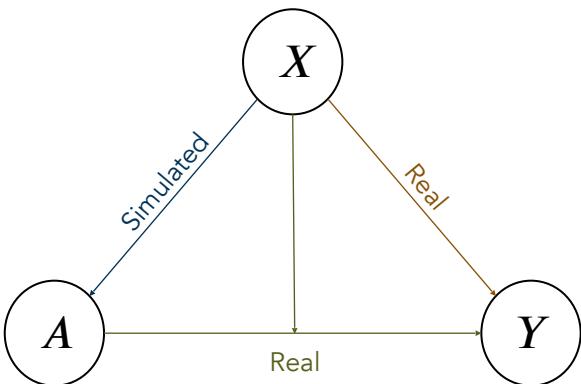
Simulations

No Simulation
(Clinical Data)



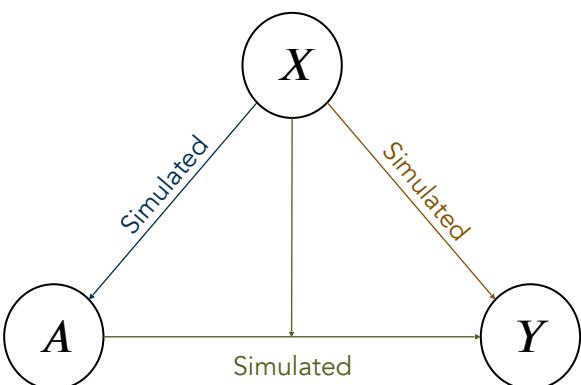
No ground-truth effect available!

A-Simulation



Less Realistic
More Evaluation

AY-Simulation

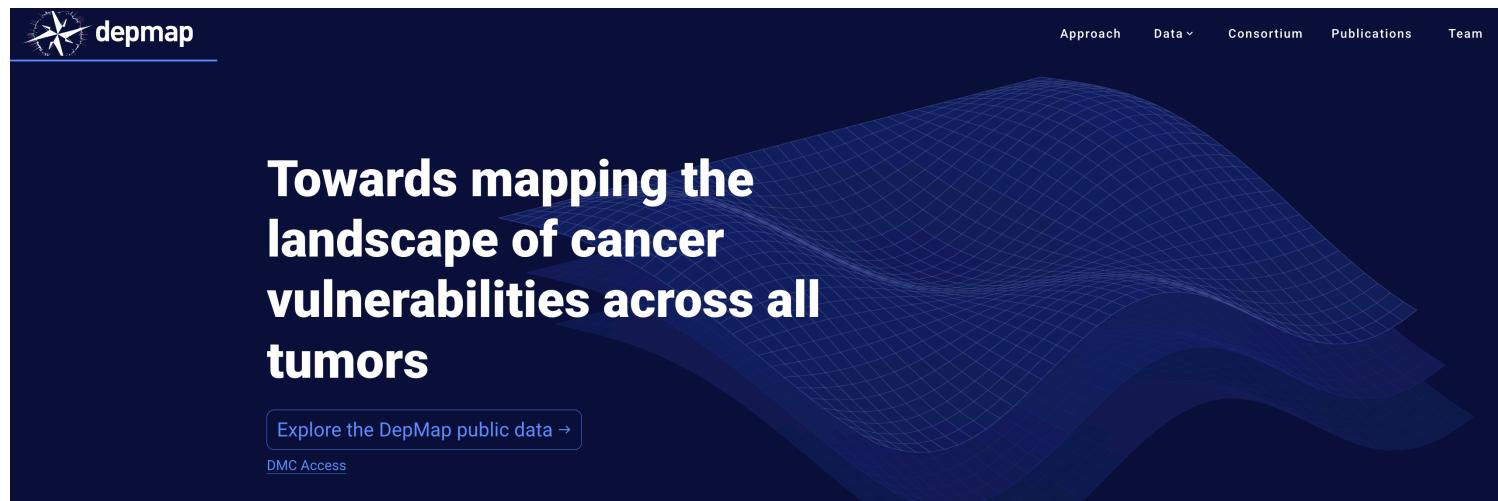
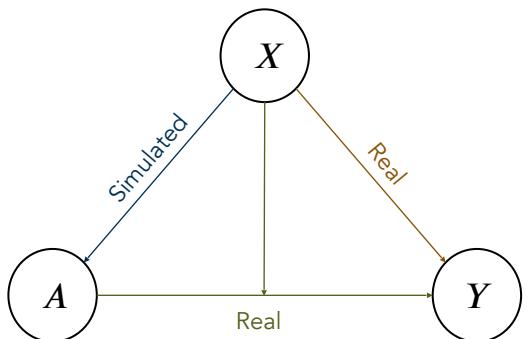




A-Simulation

'patients' ... ~1000 cell lines as

'treatments' ... Drugs or CRISPR knock-outs as 'treatments'

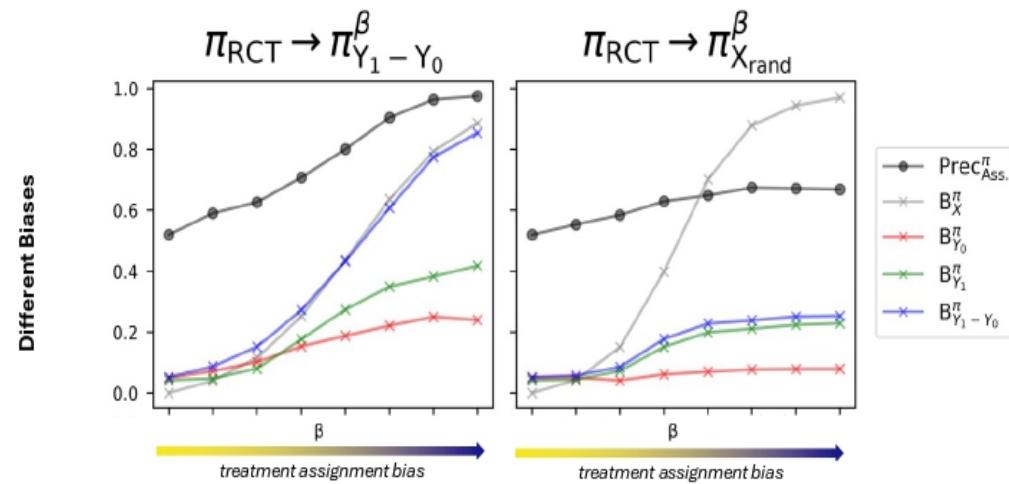




Results – Drug Screen

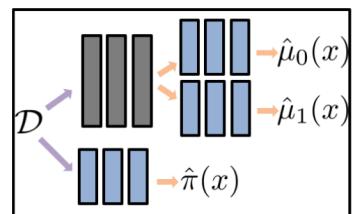
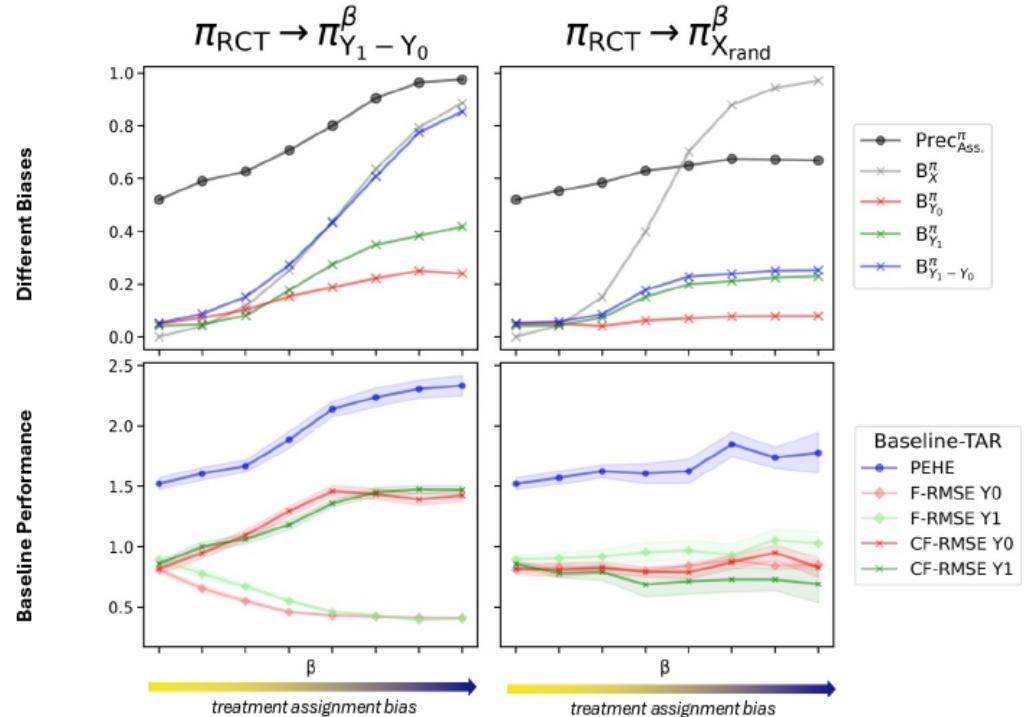
Definition 4 (Z-Policy) For a given random variable Z and parameter β , the Z-policy is defined as

$$\pi_Z^\beta(x) := \sigma(\beta Z(x)), \quad \sigma(x) := \frac{e^x}{1+e^x}$$





Results – Drug Screen

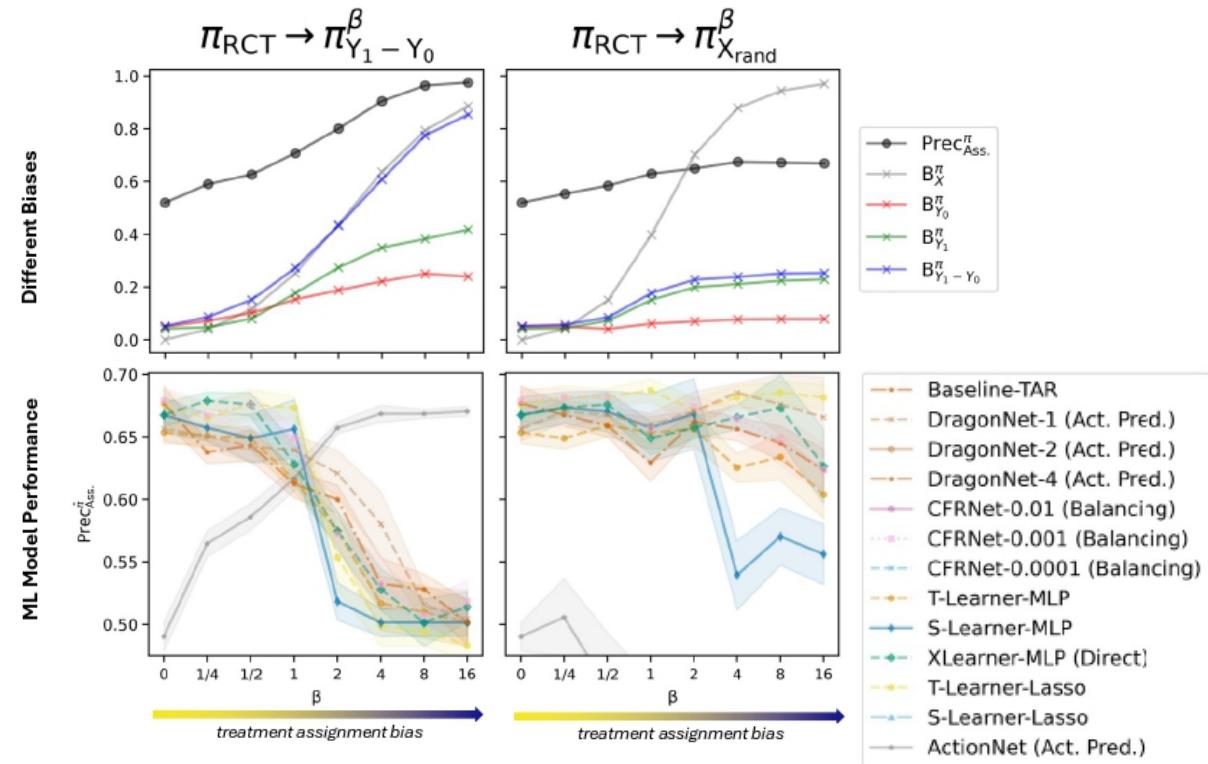


TARNet (SNet-1)

$$PEHE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{\tau}(x_i) - \tau(x_i))^2}, \quad \tau(x) = E[Y_i^1 - Y_i^0 | X_i = x]$$



Results – Drug Screen





Conclusion

❖ Formalizing Treatment Biases

Modeled various types of treatment assignment biases using mutual information.

❖ Realistic Data Benchmarks

Introduced novel biological datasets for robust evaluations.

❖ Bias Impact Analysis

Provide simulation framework for how different biases affect counterfactual predictions and biomarker identification.

❖ Overlap Assumption

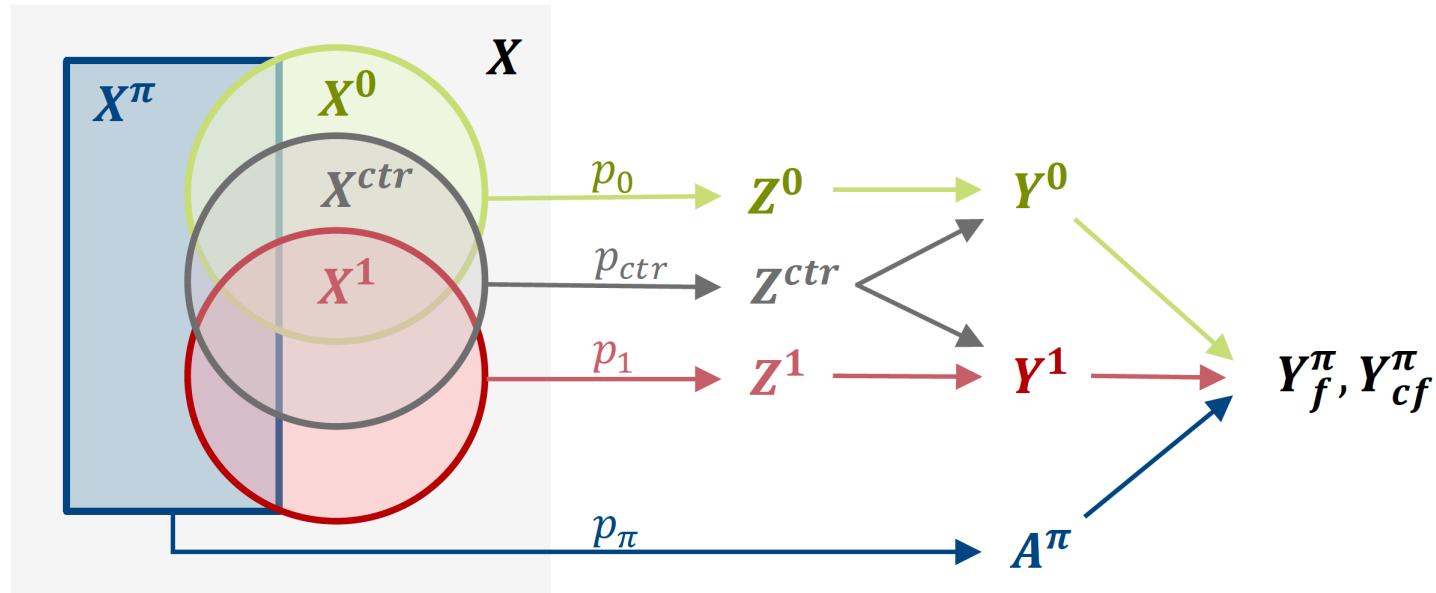
Demonstrated that the violation of the overlap assumption is not necessarily detrimental.

❖ Model Development

Highlighted the need for bias-aware ML models in precision medicine applications.



Data Model



$$A^\pi \sim \pi(X^\pi)$$

$$Z^a \sim p_a(X^a)$$

$$Z^{ctr} \sim p_{ctr}(X^{ctr})$$

$$Y^a := Z^{ctr} + Z^a$$

$$Y_f^\pi := A^\pi * Y^1 + (1 - A^\pi) * Y^0$$

$$Y_{cf}^\pi := A^\pi * Y^0 + (1 - A^\pi) * Y^1$$



Data Model

Proposition 2 *Under the non-confounding assumption, we have*

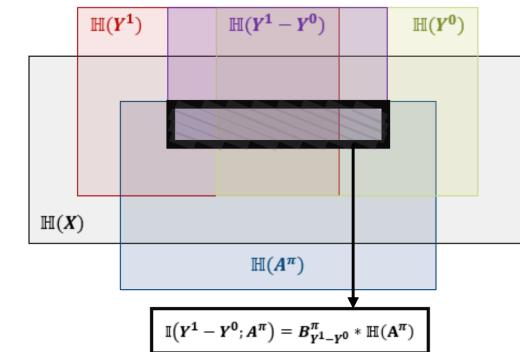
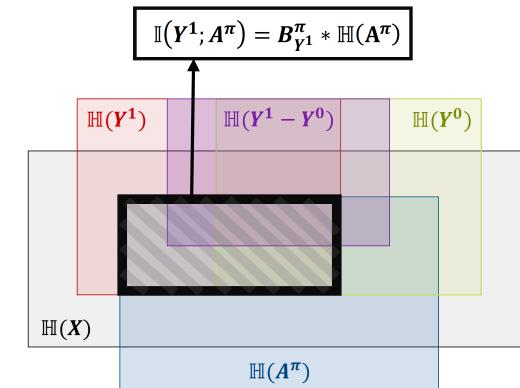
$$B_X^\pi \geq B_{Y^0, Y^1}^\pi \geq B_{Y^0}^\pi, B_{Y^1}^\pi, B_{Y^1 - Y^0}^\pi.$$

Proposition 3 *Z-bias relates to the overlap assumption (OA) as follows*

$B_X^\pi = 1 \implies OA \text{ is violated.}$

Definition 4 (Z-Policy) For a given random variable Z and parameter β , the Z -policy is defined as

$$\pi_Z^\beta(x) := \sigma(\beta Z(x)), \quad \sigma(x) := \frac{e^x}{1+e^x}$$





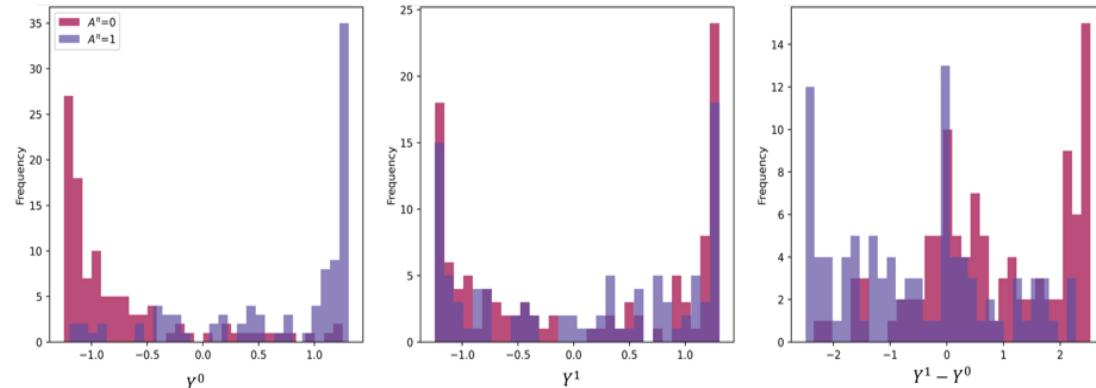
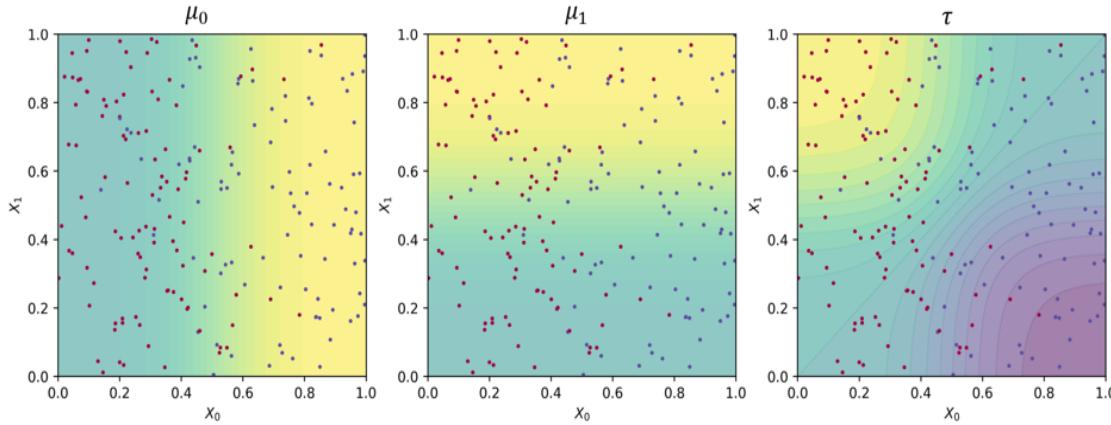
Datasets

PRISM repurposing drug screens. ([DepMap and Kocak, 2024](#)), performed on 877 cell lines according to [Corsello et al. \(2020\)](#) record the effect of hundreds of drugs on various cancer cell lines. The goal is to identify new therapeutic uses for existing drugs through a multiplexed approach. For our use case we focused on the measurements for two drugs, imatinib and az-628. Imatinib is a well-known tyrosine kinase inhibitor used mainly for chronic myeloid leukemia (CML) and gastrointestinal stromal tumors (GISTs). Az-628 is a selective inhibitor of RAF kinase, a key player in the MAPK/ERK signaling pathway, which is often involved in cancer development. As covariates we use RNA transcriptomics data, retaining the 200 most correlated features as covariates.

CRISPR knock-out screens. ([DepMap, Broad, 2021](#)) were conducted on 1067 cell lines following the methodology described by [Behan et al. \(2019\)](#), and targeted a myriad of genes of interest. They provide critical insights into gene dependencies by knocking out specific genes to observe their effects on cancer cell viability. As the two treatment options in our setting, we used the knock-out of genes EGFR (Epidermal Growth Factor Receptor) and KRAS (Kirsten Rat Sarcoma Viral Oncogene Homolog), which are crucial in cell signaling pathways related to cancer, with mutations or overexpression frequently observed in various cancers. Again, the most correlated 200 RNA transcriptomics features were used as covariates for our analysis. We refer to these datasets as A-CRISPR and A-DRUG.



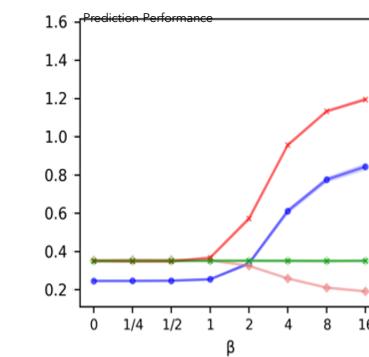
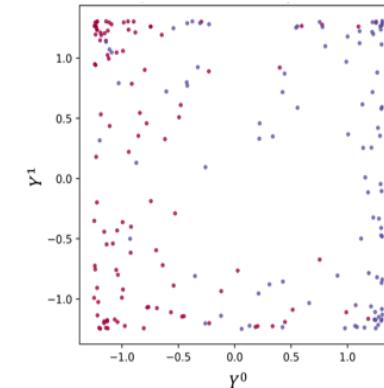
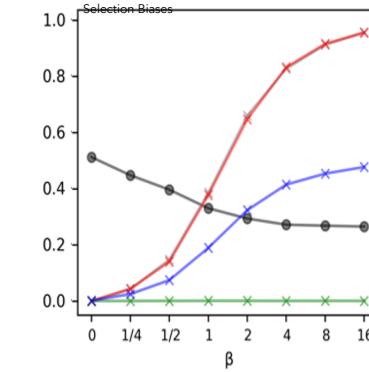
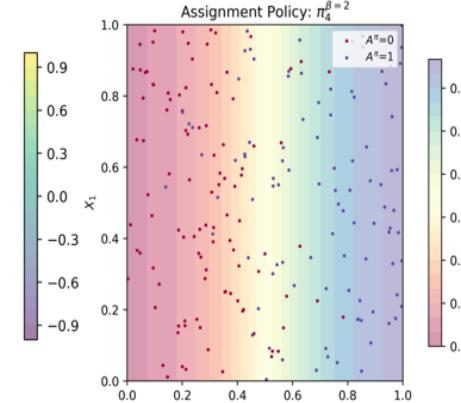
Toy Example



$$Y^a = f_{nl}(f_i^a(x_0, x_1))$$

$$\begin{aligned} f_4^0(x_0, x_1) &:= x_0 \\ f_4^1(x_0, x_1) &:= x_1 \\ \pi_4^\beta(x_0, x_1) &\sim x_0 \end{aligned}$$

$$f_{nl}(\cdot) = \frac{1}{1 + \exp(-10 * (-0.5))}$$



Legend for Selection Biases:

- $\text{Prec}^\pi_{\text{Ass.}}$
- B_X^π
- $B_{Y_0}^\pi$
- $B_{Y_1}^\pi$
- $B_{Y_1 - Y_0}^\pi$

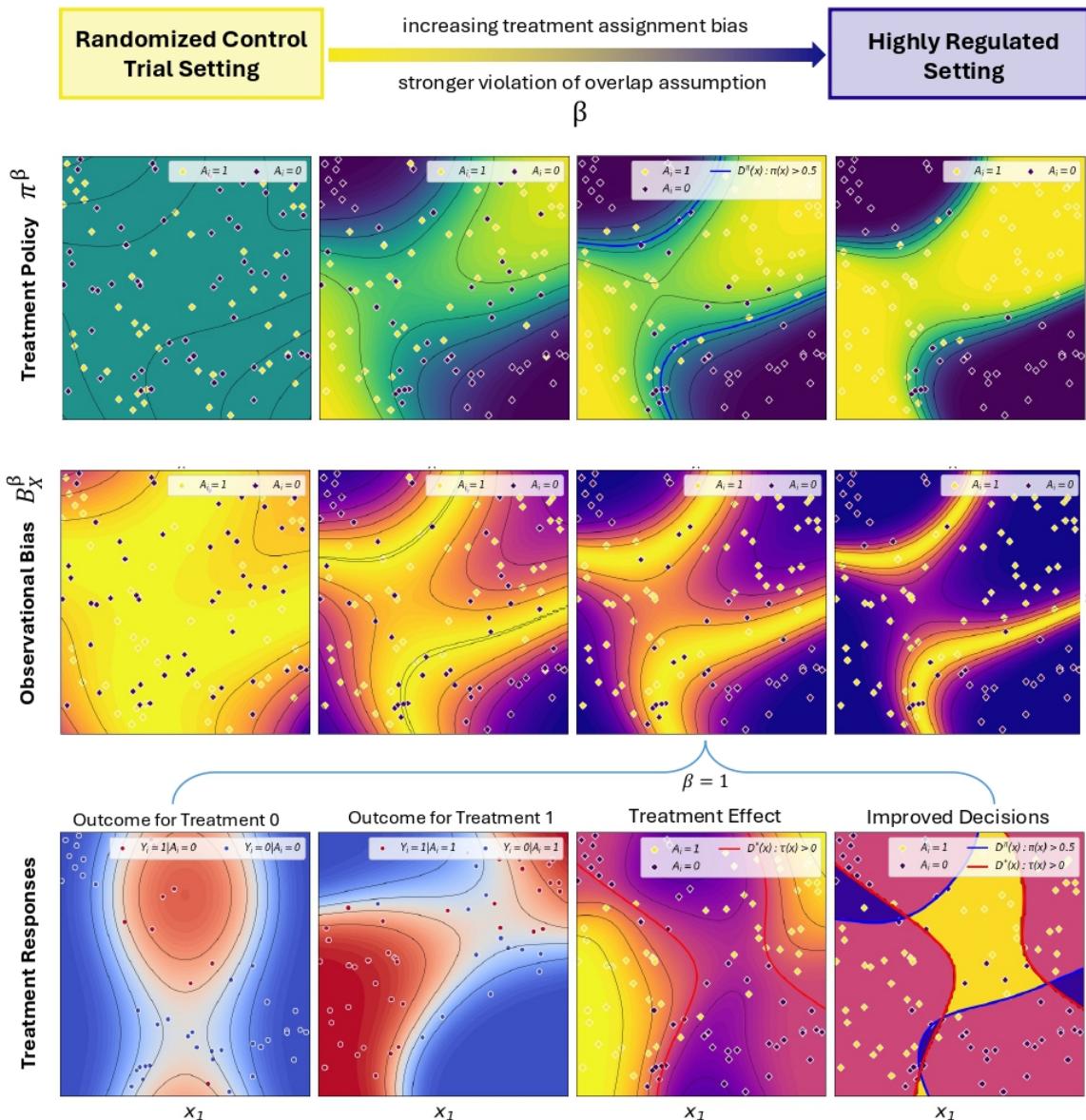
Legend for Prediction Performance:

T-Learner-Lasso

- PEHE
- F-RMSE Y0
- F-RMSE Y1
- CF-RMSE Y0
- CF-RMSE Y1



Toy Example





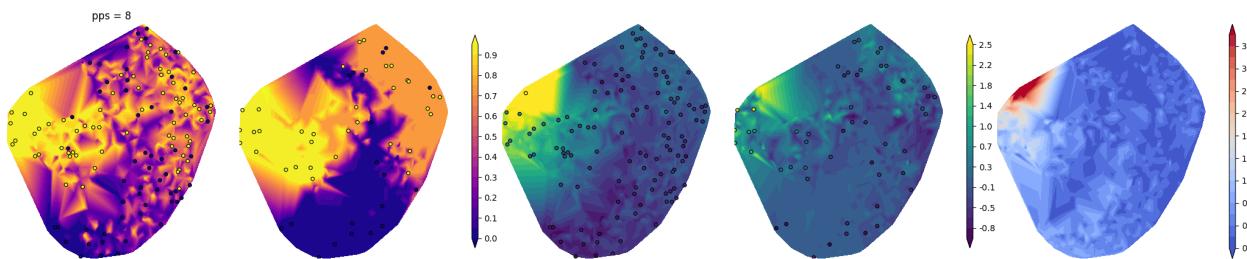
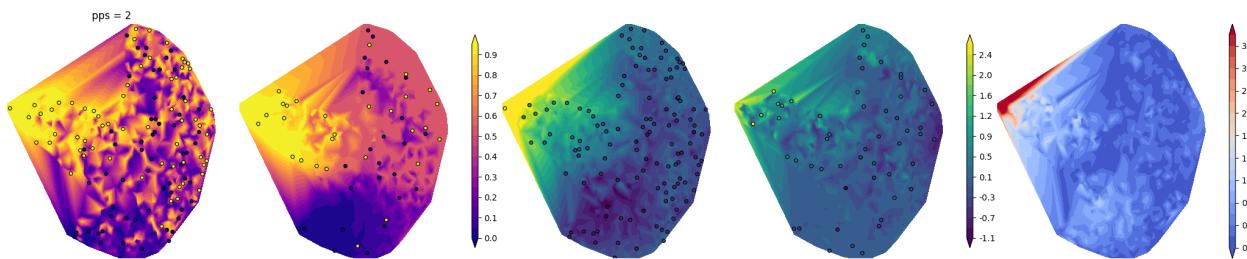
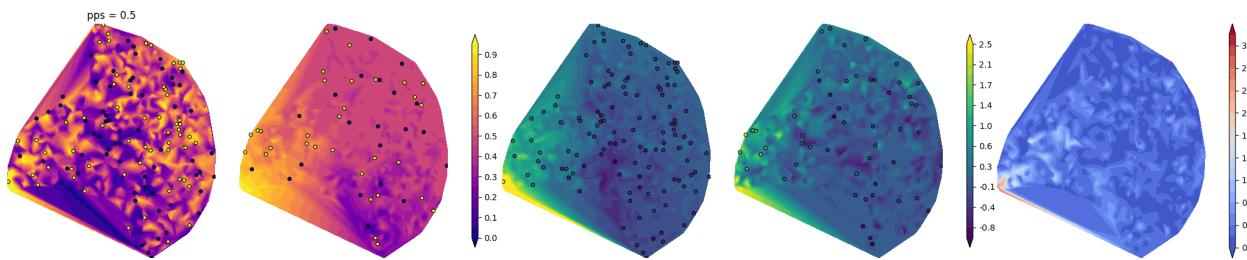
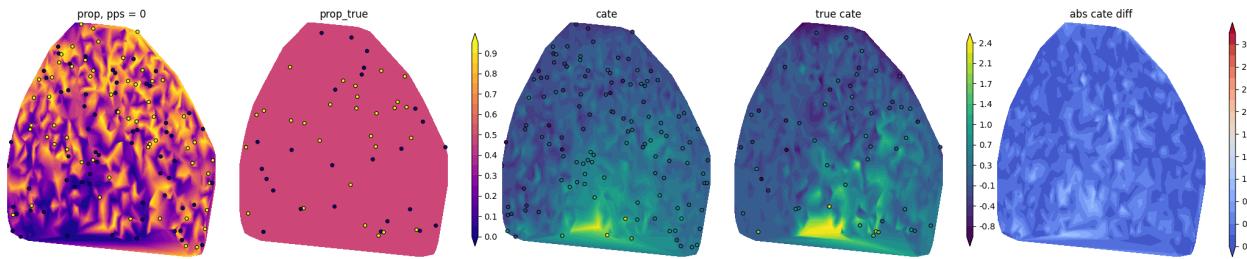
Simulation Settings

Setting	Overall Usefulness based on Mike's Intuition (0-10)	Overall Closeness to Reality (0-10)	Real Patient Data	Real Outcome Mechanism	Real Counterfactual Outcomes	Counterfactual Outcome and PEHE Evaluation	Selective Feature Importance Evaluation	Prognostic/Predictive Feature Importance Evaluation	Sample Size Effect	Propensity Scale Effect, T Imbalance	Y Imbalance, #Num Dims #Treatments, #Imp. Feat.
Toy, TY-simulated	2 (debugging)	0	✗	✗	✗	✓	✓	✓	✓	✓	✓
In-vivo	3 (realistic, but no counterfactual evaluation)	10	✓	✓	✗	✗	✗	✗	✗	✗	✗
Real X, TY-simulated (linear)	4 (all evaluation types, but too simple)	3	✓	✗	✗	✓	✓	✓	✗	✓	✓
Synthetic X, TY-simulated (linear)	4 (all evaluation types, but very unrealistic)	2	✗	✗	✗	✓	✓	✓	✓	✓	✓
Real X, TY-simulated (Y learned via ML)	7 (more realistic, but too few samples)	6	✓	~	✗	✓	✓	✓	✗	✓	✓
Synthetic X, TY-simulated (Y learned via ML)	8 (more samples, all evaluation, but synthetic data)	5	✗	~	✗	✓	✓	✓	✓	✓	✓
Real X, TY-simulated (biological sim.)	8 (more realistic, domain knowledge based, but too few samples)	7	✓	~	~	✓	✓	✓	✗	✓	✓
Synthetic X, TY-simulated (biological sim.)	8 (more samples, all evaluation, domain knowledge based, but synthetic data)	5	✗	~	~	✓	✓	✓	✓	✓	✓
In-vitro	9 (realistic, counterfactual, but no propensity analysis)	10	✓	✓	✓	✓	✗	✗	✗	✗	✗
In-vitro, T-simulated	9 (realistic, counterfactual, propensity analysis, but few samples)	9	✓	✓	✓	✓	✓	✗	✗	✓	✗
In-vitro, T-simulated, Large	10 (realistic, counterfactual, propensity analysis, sample size analysis)	9	✓	✓	✓	✓	✓	✗	✓	✓	✗



Simulation Settings

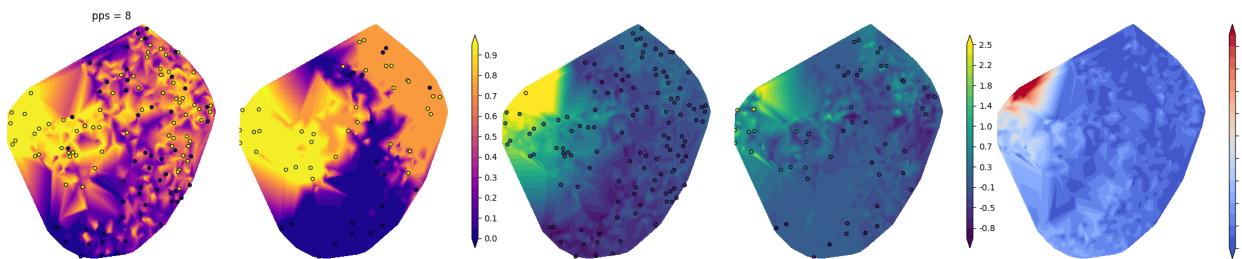
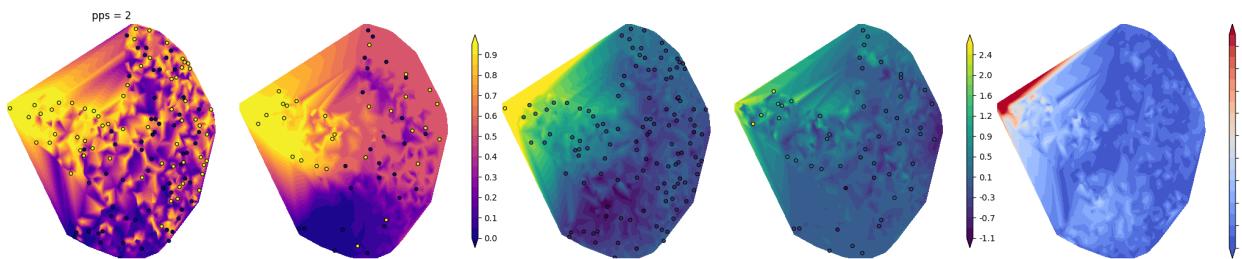
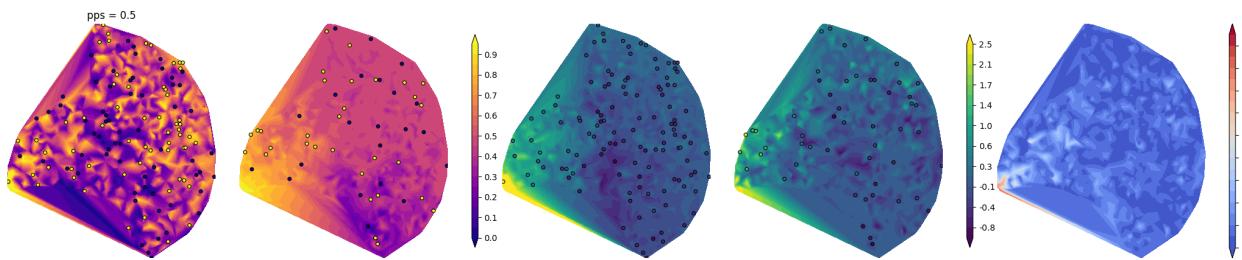
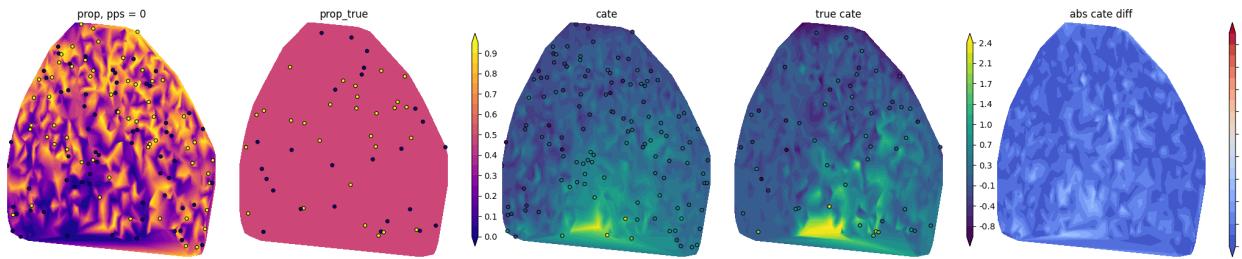
Sample Size Sensitivity, N = 175





Simulation Settings

Sample Size Sensitivity, N = 175





A-Drug Results

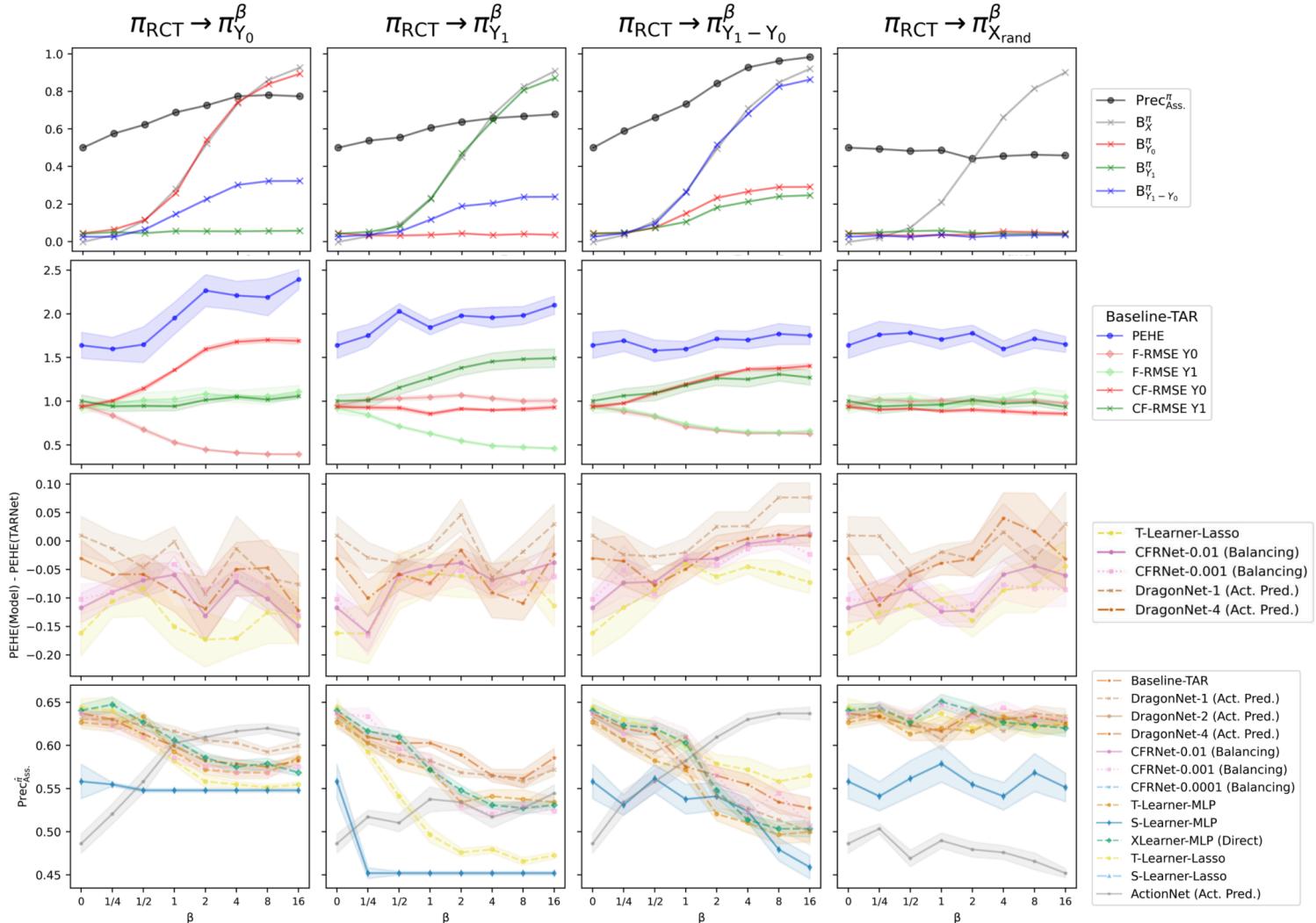
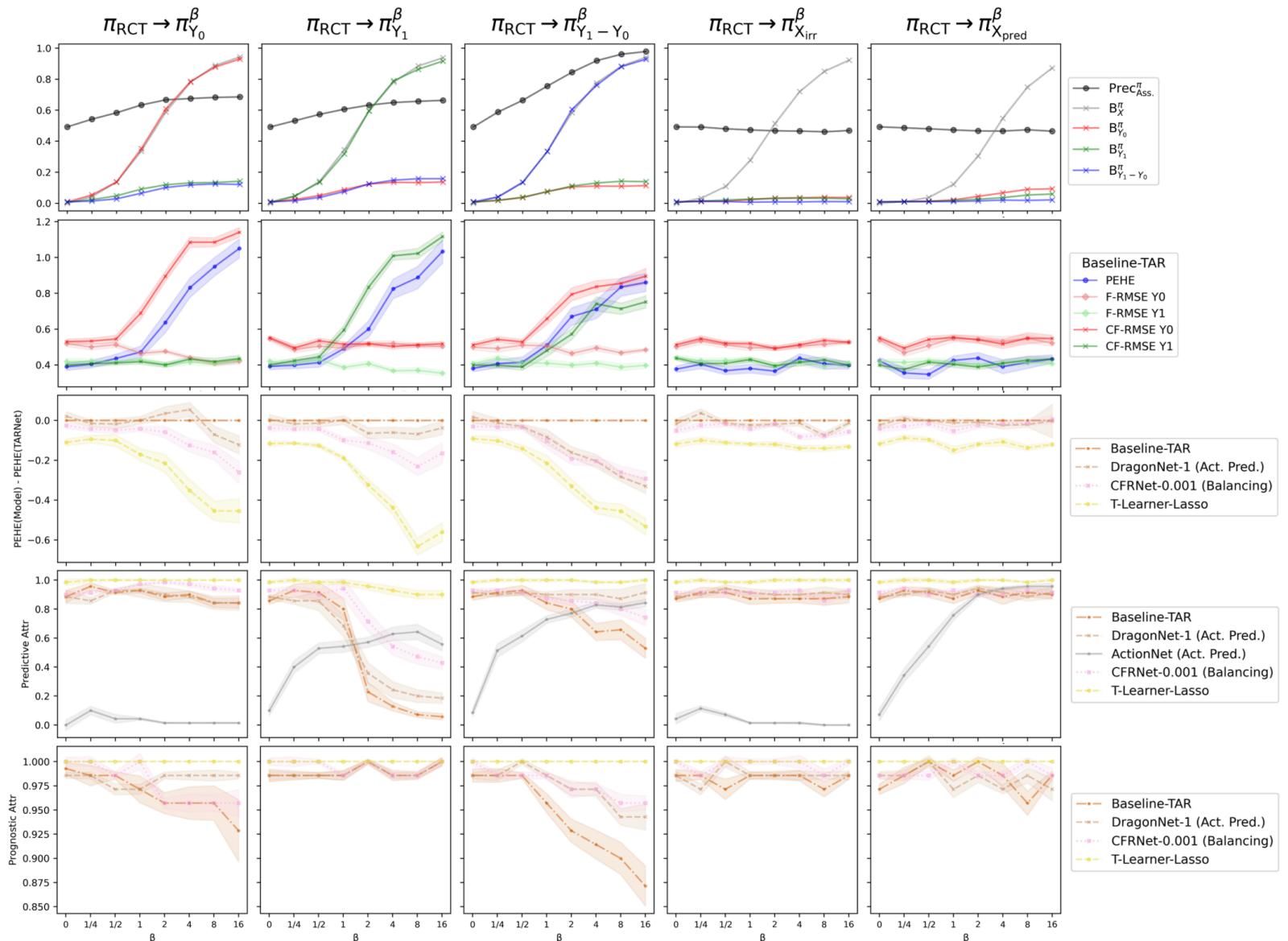


Figure 9: Results for A-Drug.



AY-TCGA Results





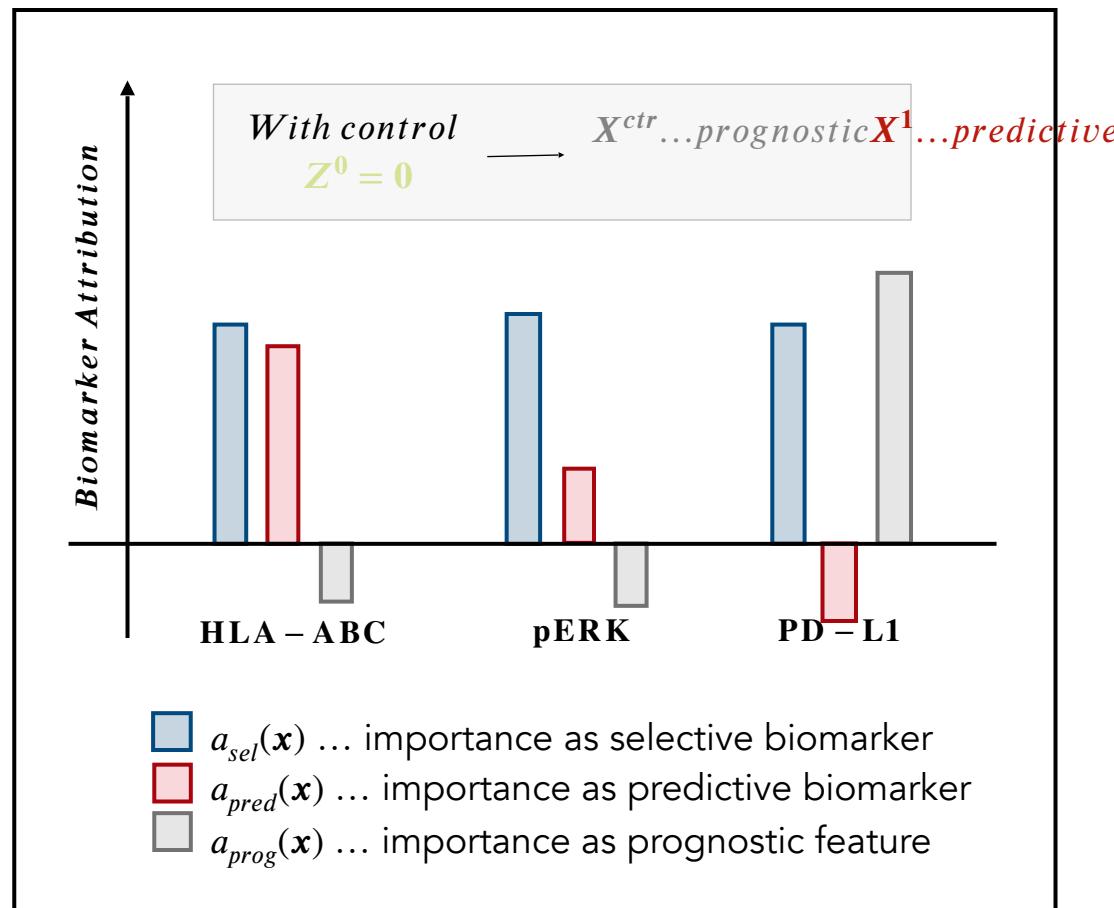
Simulation Settings

With control \longrightarrow $X^{ctr} \dots$ prognostic $\textcolor{red}{X^1} \dots$ predictive
 $Z^0 = 0$

- $a_{sel}(x)$... importance as selective biomarker
- $a_{pred}(x)$... importance as predictive biomarker
- $a_{prog}(x)$... importance as prognostic feature

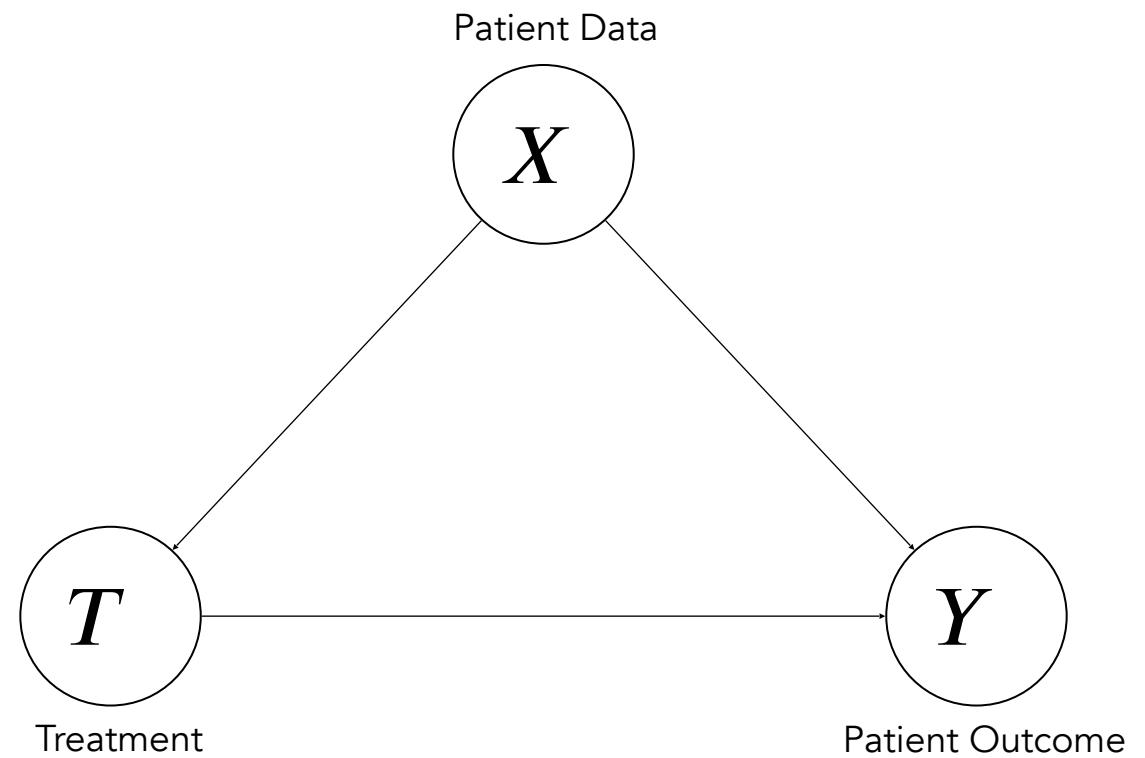


Simulation Settings





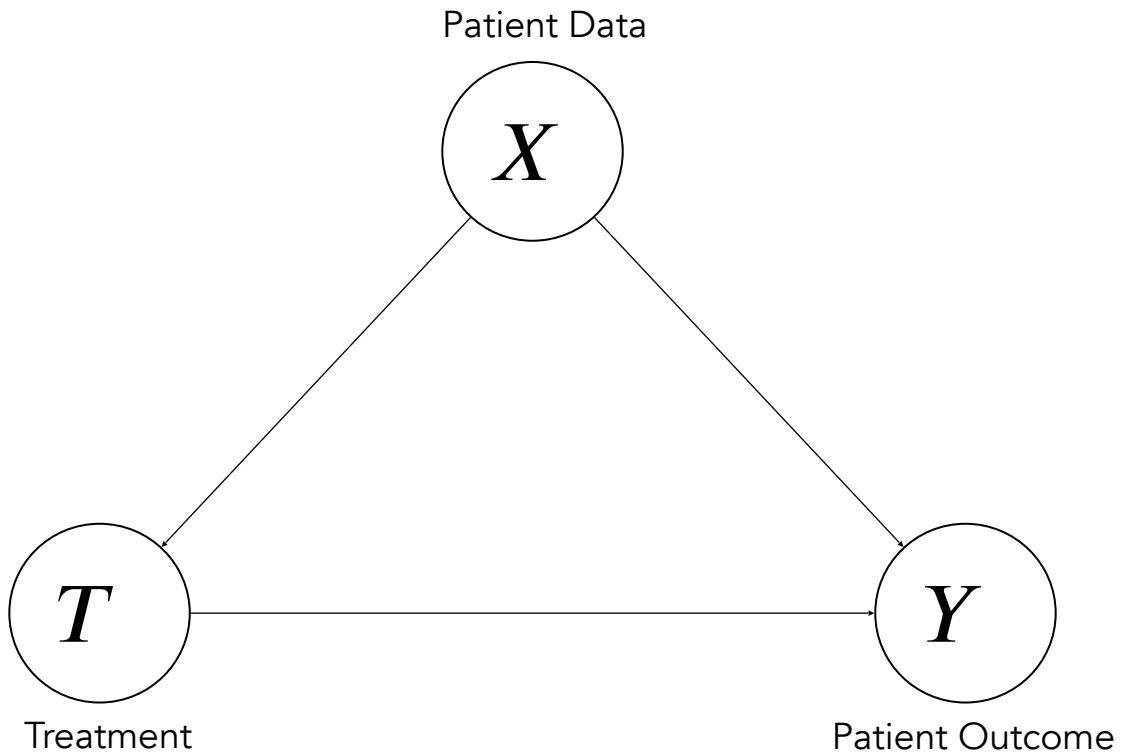
Counterfactuals





Counterfactuals

$$P(Y|X, T)$$

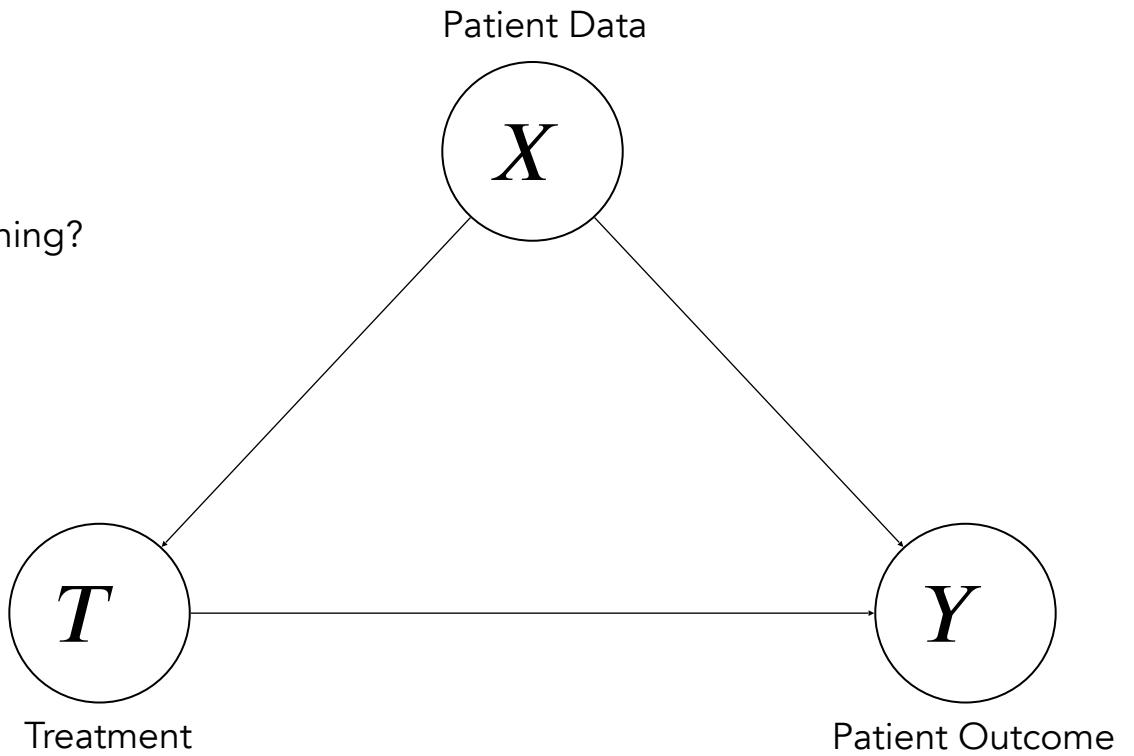




Counterfactuals

$$P(Y|X, T)$$

→ What is happening?

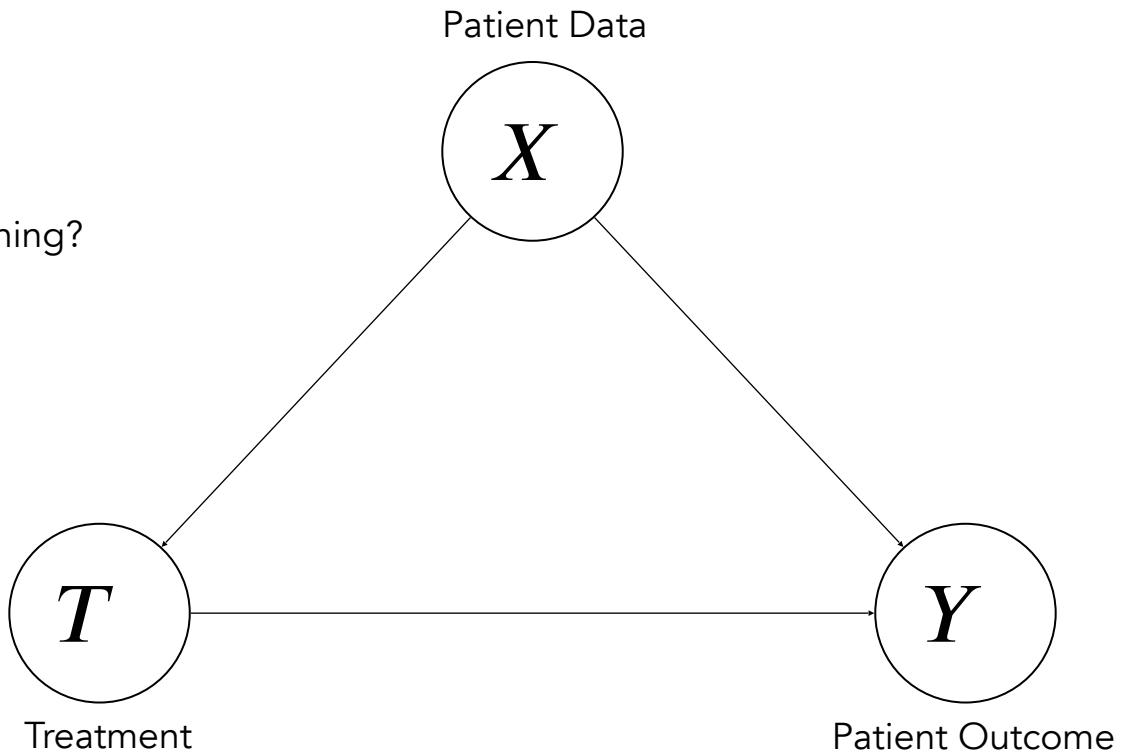




Counterfactuals

$$P(Y|X, T)$$

→ What is happening?





Counterfactuals

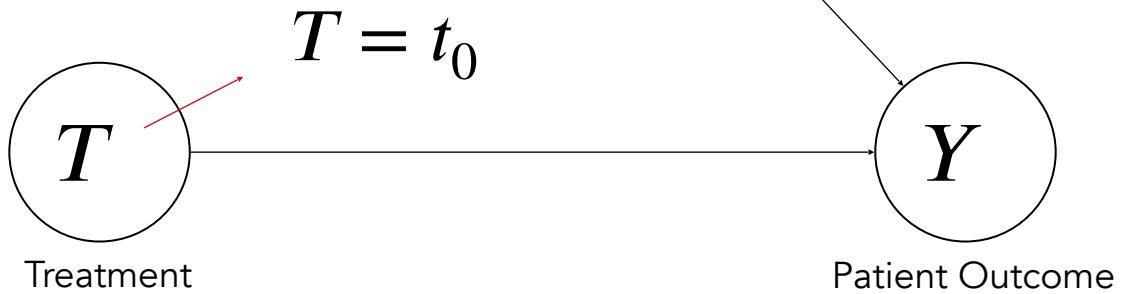
$$P(Y|X, T)$$

→ What is happening?

→ What happens if?

Patient Data

X





Counterfactuals

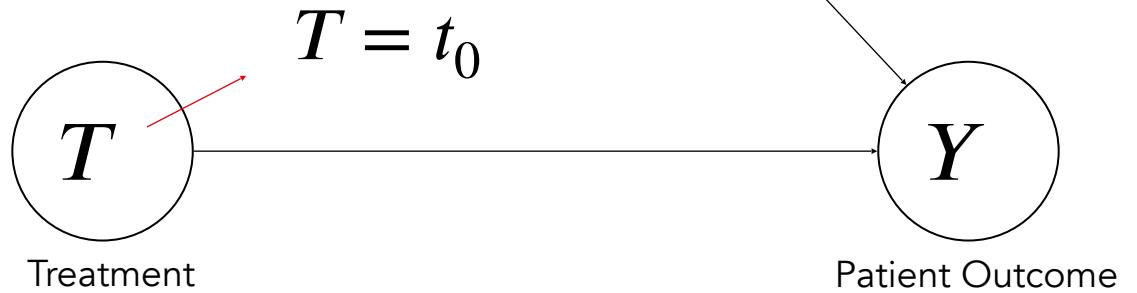
$$P(Y|X, T)$$

→ What is happening?

→ What happens if?

Patient Data

X





Counterfactuals

$$P(Y|X, T)$$

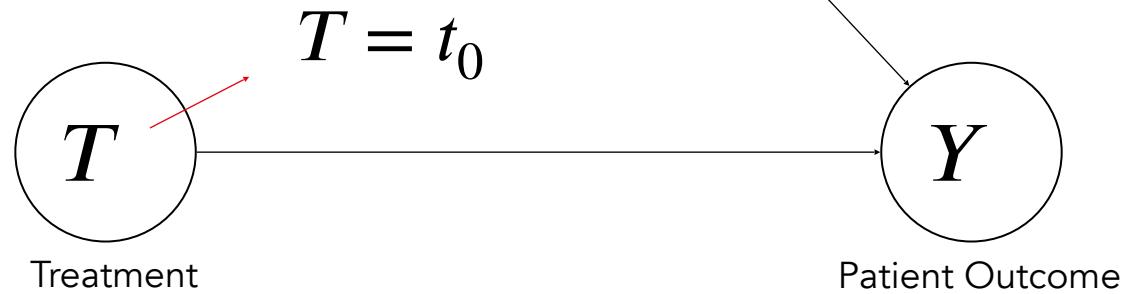
→ What is happening?

→ What happens if?

→ What would have happened if?

Patient Data

X





Counterfactuals

$$P(Y|X, T)$$

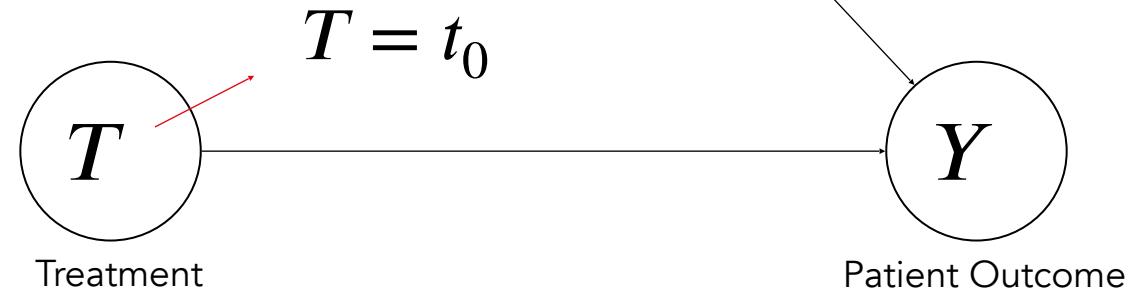
→ What is happening?

→ What happens if?

→ What would have happened if?

Patient Data

X



$$T = t_0$$

Treatment

Y

Patient Outcome

→ This is what we want to learn from data!



Clinical Biomarkers

Predictive

A biomarker used to identify individuals who are more likely than similar individuals without the biomarker to experience a favorable or unfavorable effect from exposure to a medical product or an environmental agent.

Prognostic

A biomarker used to identify likelihood of a clinical event, disease recurrence or progression in patients who have the disease or medical condition of interest.



Clinical Biomarkers

Predictive

A biomarker used to identify individuals who are more likely than similar individuals without the biomarker to experience a favorable or unfavorable effect from exposure to a medical product or an environmental agent.

Prognostic

A biomarker used to identify likelihood of a clinical event, disease recurrence or progression in patients who have the disease or medical condition of interest.



Clinical Biomarkers

Susceptibility/Risk

A biomarker that indicates the potential for developing a disease or medical condition in an individual who does not currently have clinically apparent disease or the medical condition.

Diagnostic

A biomarker used to detect or confirm presence of a disease or condition of interest or to identify individuals with a subtype of the disease.

Monitoring

A biomarker measured repeatedly for assessing status of a disease or medical condition or for evidence of exposure to (or effect of) a medical product or an environmental agent.

Safety

A biomarker measured before or after an exposure to a medical product or an environmental agent to indicate the likelihood, presence, or extent of toxicity as an adverse effect.

Prognostic

A biomarker used to identify likelihood of a clinical event, disease recurrence or progression in patients who have the disease or medical condition of interest.

Predictive

A biomarker used to identify individuals who are more likely than similar individuals without the biomarker to experience a favorable or unfavorable effect from exposure to a medical product or an environmental agent.

Response

A biomarker used to show that a biological response, potentially beneficial or harmful, has occurred in an individual who has been exposed to a medical product or an environmental agent.



Clinical Biomarkers

Prognostic

A biomarker used to identify likelihood of a clinical event, disease recurrence or progression in patients who have the disease or medical condition of interest.

Predictive

A biomarker used to identify individuals who are more likely than similar individuals without the biomarker to experience a favorable or unfavorable effect from exposure to a medical product or an environmental agent.



Clinical Biomarkers

Prognostic

A biomarker used to identify likelihood of a clinical event, disease recurrence or progression in patients who have the disease or medical condition of interest.

Predictive

A biomarker used to identify individuals who are more likely than similar individuals without the biomarker to experience a favorable or unfavorable effect from exposure to a medical product or an environmental agent.



True vs. Clinical



True vs. Clinical

Clinical Predictive Biomarker → "Selective Feature"

A biomarker used to identify individuals who are more likely than similar individuals without the biomarker to experience a favorable or unfavorable effect from exposure to a medical

→ How should the patient be treated?



True vs. Clinical

Clinical Predictive Biomarker → "Selective Feature"

A biomarker **used to identify** individuals who are more likely than similar individuals without the biomarker to experience a favorable or unfavorable effect from exposure to a medical

→ **How should the patient be treated?**

Clinical Prognostic Biomarker → "Selective Feature"

A biomarker **used to identify** likelihood of a clinical event, disease recurrence or progression in patients who have the disease or medical condition of interest.

→ **Should the patient be treated?**



True vs. Clinical

Clinical Predictive Biomarker → "Selective Feature"

A biomarker **used to identify** individuals who are more likely than similar individuals without the biomarker to experience a favorable or unfavorable effect from exposure to a medical

→ **How should the patient be treated?**

Clinical Prognostic Biomarker → "Selective Feature"

A biomarker **used to identify** likelihood of a clinical event, disease recurrence or progression in patients who have the disease or medical condition of interest.

→ **Should the patient be treated?**

True Predictive Biomarker → "Predictive Feature"

A biomarker **used to identify determining** individuals who are more likely than similar individuals without the biomarker to experience a favorable or unfavorable effect from exposure

→ **What happens when the patient is treated?**



True vs. Clinical

Clinical Predictive Biomarker → "Selective Feature"

A biomarker **used to identify** individuals who are more likely than similar individuals without the biomarker to experience a favorable or unfavorable effect from exposure to a medical

→ **How should the patient be treated?**

Clinical Prognostic Biomarker → "Selective Feature"

A biomarker **used to identify** likelihood of a clinical event, disease recurrence or progression in patients who have the disease or medical condition of interest.

→ **Should the patient be treated?**

True Predictive Biomarker → "Predictive Feature"

A biomarker **used to identify** **determining** individuals who are more likely than similar individuals without the biomarker to experience a favorable or unfavorable effect from exposure

→ **What happens when the patient is treated?**

True Prognostic Biomarker → "Prognostic Feature"

A biomarker **used to identify** **determining** the likelihood of a clinical event, disease recurrence or progression in patients who have the disease or medical condition of interest.

→ **What happens when the patient remains untreated?**