

What Can Conformal Inference Offer to Statistics?

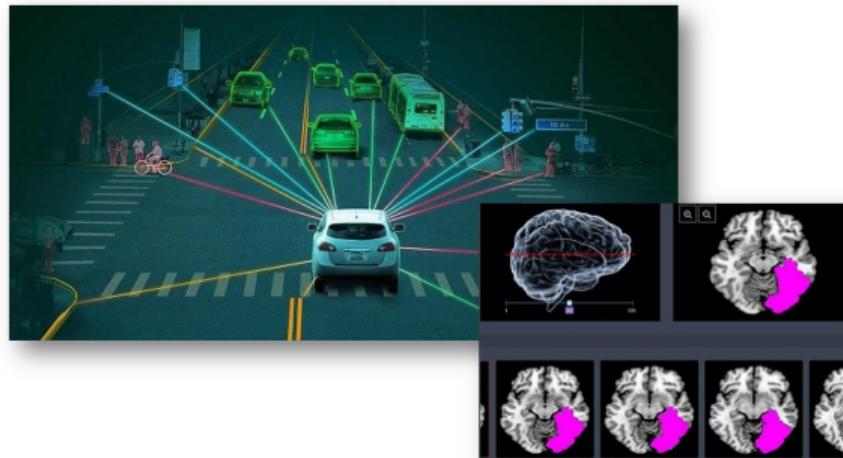
Lihua Lei

Department of Statistics, Stanford University

Job Talk, 2021-2022

ML in critical applications

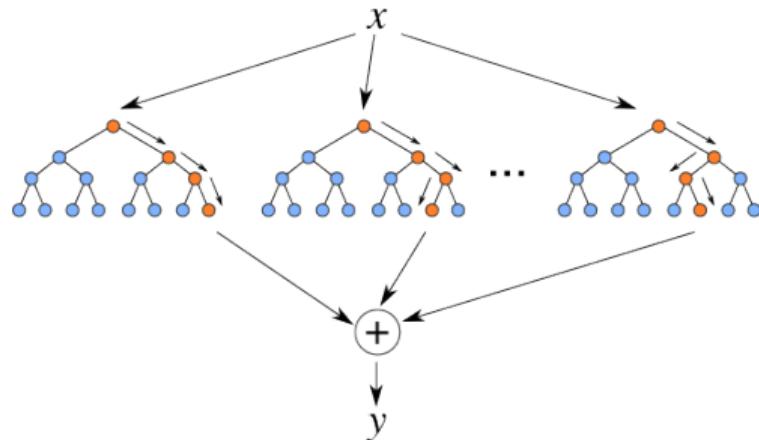
ML tools make potentially high-stakes decisions: self-driving cars, disease diagnosis, ...



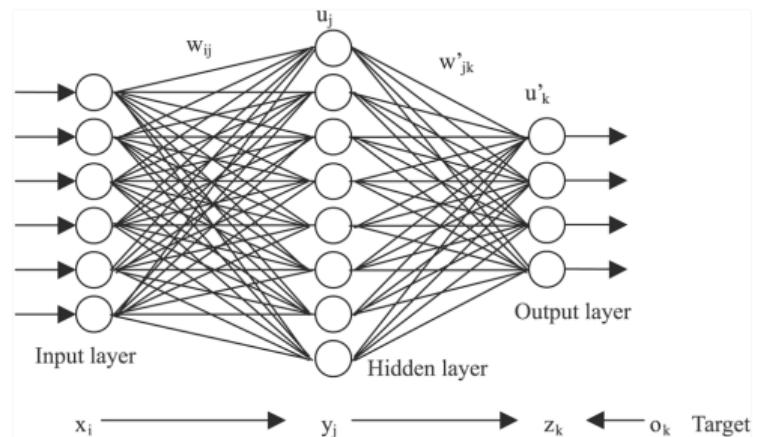
Can we have reliable uncertainty quantification (confidence) in these predictions?

Today's predictive algorithms

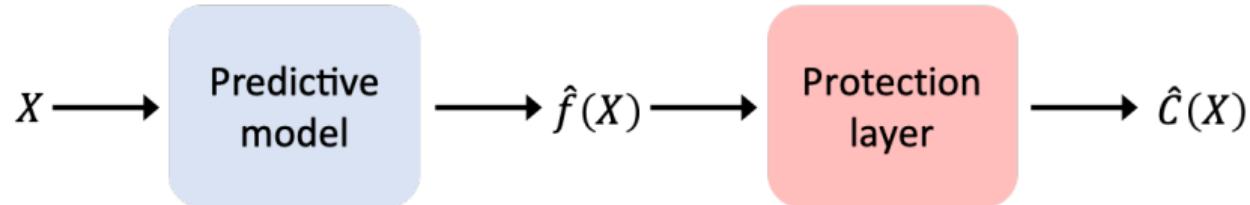
random forests, gradient boosting



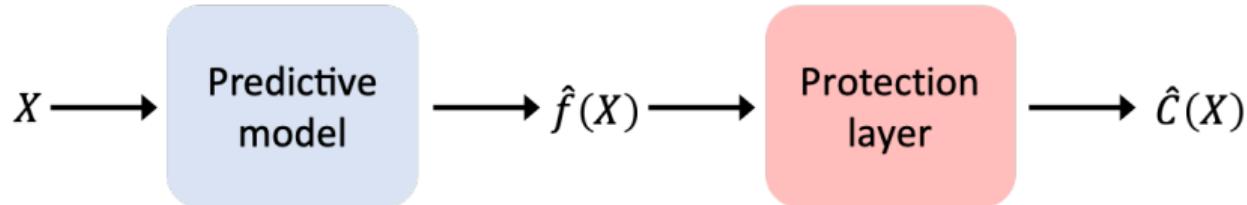
neural networks



Previous work on conformal inference



Previous work on conformal inference



- ▶ i.i.d. training samples (X_i, Y_i) , $i = 1, \dots, n$
- ▶ Test point $(X, Y = ?)$ from the same distribution
- ▶ Conformal inference Vovk et al. '99, Papadopoulos et al. '12, Lei et al. '18, Barber et al. '19, Romano et al. '19

Constructs predictive interval $\hat{C}(x)$ with $\mathbb{P}(Y \in \hat{C}(X)) \geq 90\%$

- ▶ Holds in finite samples for any distribution of (X, Y) and any predictive algorithm \hat{f}

Statistical inference is more complicated than predicting “the seen”

Statistical inference is more complicated than predicting “the seen”

- ▶ Counterfactual
 - ▶ Causal inference, offline policy evaluation, algorithmic fairness, ...
 - ▶ What would have been one's response had one taken the treatment
 - ▶ Observable for those in a particular “treatment arm”

Statistical inference is more complicated than predicting “the seen”

- ▶ Counterfactual
 - ▶ Causal inference, offline policy evaluation, algorithmic fairness, ...
 - ▶ What would have been one's response had one taken the treatment
 - ▶ Observable for those in a particular “treatment arm”
- ▶ Time-to-event (survival) outcome
 - ▶ Survival analysis, industrial life testing, economics, ...
 - ▶ Censored by study termination, loss to follow-up, ...
 - ▶ Observable for those whose event (e.g., death) has occurred

Statistical inference is more complicated than predicting “the seen”

- ▶ Counterfactual
 - ▶ Causal inference, offline policy evaluation, algorithmic fairness, ...
 - ▶ What would have been one's response had one taken the treatment
 - ▶ Observable for those in a particular “treatment arm”
- ▶ Time-to-event (survival) outcome
 - ▶ Survival analysis, industrial life testing, economics, ...
 - ▶ Censored by study termination, loss to follow-up, ...
 - ▶ Observable for those whose event (e.g., death) has occurred

Goal: construct calibrated prediction intervals for partially observed outcomes

Part I: conformalized counterfactual prediction



Emmanuel Candès

Inference of counterfactuals? Potential outcomes

Neyman '23, Rubin '74

Assumptions

- stable unit treatment values (SUTVA)

| science table | | | | | |
|-----------------|-------|-------|----------|----------|--------------------|
| Unit | x_i | T_i | $Y_i(1)$ | $Y_i(0)$ | Y_i^{obs} |
| Treatment Group | | | | | |
| 1 | ✓ | 1 | ✓ | ✗ | $Y_1(1)$ |
| 2 | ✓ | 1 | ✓ | ✗ | $Y_2(1)$ |
| 3 | ✓ | 1 | ✓ | ✗ | $Y_3(1)$ |
| 4 | ✓ | 1 | ✓ | ✗ | $Y_4(1)$ |
| 5 | ✓ | 1 | ✓ | ✗ | $Y_5(1)$ |
| Control Group | | | | | |
| 6 | ✓ | 0 | ✗ | ✓ | $Y_6(0)$ |
| 7 | ✓ | 0 | ✗ | ✓ | $Y_7(0)$ |
| 8 | ✓ | 0 | ✗ | ✓ | $Y_8(0)$ |
| 9 | ✓ | 0 | ✗ | ✓ | $Y_9(0)$ |
| 10 | ✓ | 0 | ✗ | ✓ | $Y_{10}(0)$ |

Inference of counterfactuals? Potential outcomes

Neyman '23, Rubin '74

Assumptions

- ▶ stable unit treatment values (SUTVA)
- ▶ super population (i.i.d.)
- ▶ unconfoundedness ($Y(1), Y(0)$) $\perp\!\!\!\perp T | X$

| science table | | | | | |
|-----------------|-------|-------|----------|----------|--------------------|
| Unit | x_i | T_i | $Y_i(1)$ | $Y_i(0)$ | Y_i^{obs} |
| Treatment Group | | | | | |
| 1 | ✓ | 1 | ✓ | ✗ | $Y_1(1)$ |
| 2 | ✓ | 1 | ✓ | ✗ | $Y_2(1)$ |
| 3 | ✓ | 1 | ✓ | ✗ | $Y_3(1)$ |
| 4 | ✓ | 1 | ✓ | ✗ | $Y_4(1)$ |
| 5 | ✓ | 1 | ✓ | ✗ | $Y_5(1)$ |
| Control Group | | | | | |
| 6 | ✓ | 0 | ✗ | ✓ | $Y_6(0)$ |
| 7 | ✓ | 0 | ✗ | ✓ | $Y_7(0)$ |
| 8 | ✓ | 0 | ✗ | ✓ | $Y_8(0)$ |
| 9 | ✓ | 0 | ✗ | ✓ | $Y_9(0)$ |
| 10 | ✓ | 0 | ✗ | ✓ | $Y_{10}(0)$ |

Inference of counterfactuals? Potential outcomes

Neyman '23, Rubin '74

Assumptions

- ▶ stable unit treatment values (SUTVA)
- ▶ super population (i.i.d.)
- ▶ unconfoundedness ($Y(1), Y(0)$) $\perp\!\!\!\perp T | X$

Goal: find interval estimate $\hat{C}_1(X)$, s.t.,

$$\mathbb{P}(Y(1) \in \hat{C}_1(X) | T = 0) \geq 90\%$$

| science table | | | | | |
|-----------------|-------|-------|----------|----------|--------------------|
| Unit | x_i | T_i | $Y_i(1)$ | $Y_i(0)$ | Y_i^{obs} |
| Treatment Group | | | | | |
| 1 | ✓ | 1 | ✓ | ✗ | $Y_1(1)$ |
| 2 | ✓ | 1 | ✓ | ✗ | $Y_2(1)$ |
| 3 | ✓ | 1 | ✓ | ✗ | $Y_3(1)$ |
| 4 | ✓ | 1 | ✓ | ✗ | $Y_4(1)$ |
| 5 | ✓ | 1 | ✓ | ✗ | $Y_5(1)$ |
| Control Group | | | | | |
| 6 | ✓ | 0 | ✗ | ✓ | $Y_6(0)$ |
| 7 | ✓ | 0 | ✗ | ✓ | $Y_7(0)$ |
| 8 | ✓ | 0 | ✗ | ✓ | $Y_8(0)$ |
| 9 | ✓ | 0 | ✗ | ✓ | $Y_9(0)$ |
| 10 | ✓ | 0 | ✗ | ✓ | $Y_{10}(0)$ |

Counterfactual inference

Assign treatment by a coin toss for each subject based on the **propensity score** $e(x)$

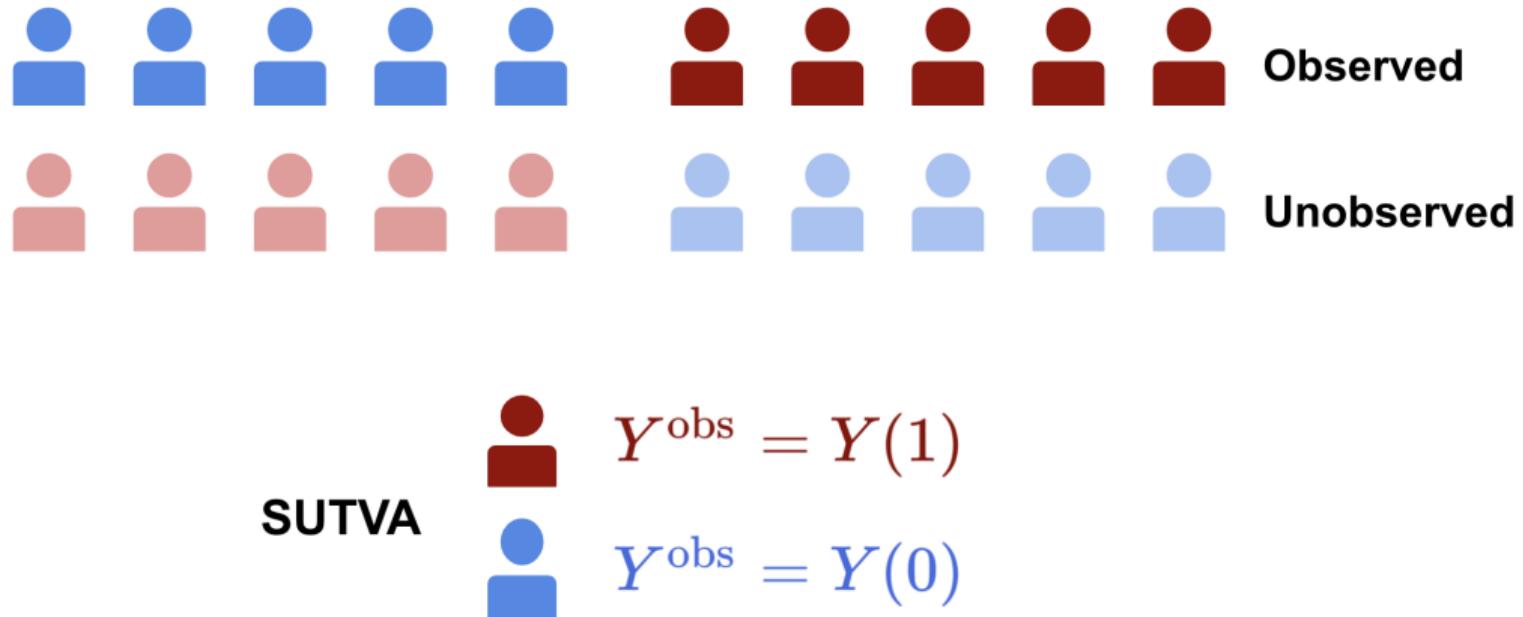


$$\mathbb{P}(\text{treated} \mid X = x) = e(x)$$

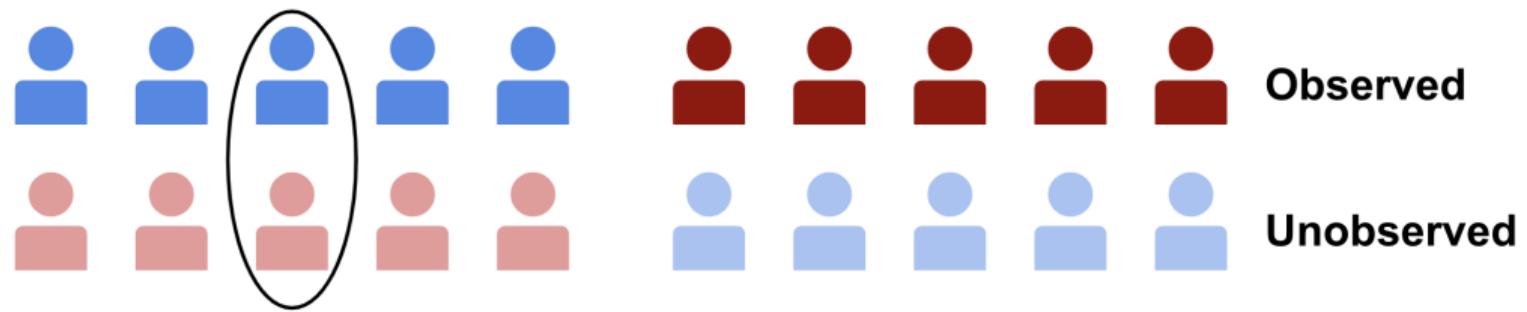
$$\mathbb{P}(\text{control} \mid X = x) = 1 - e(x)$$

Counterfactual inference

Each subject has potential outcomes ($Y(1)$, $Y(0)$) and the observed outcome Y^{obs}

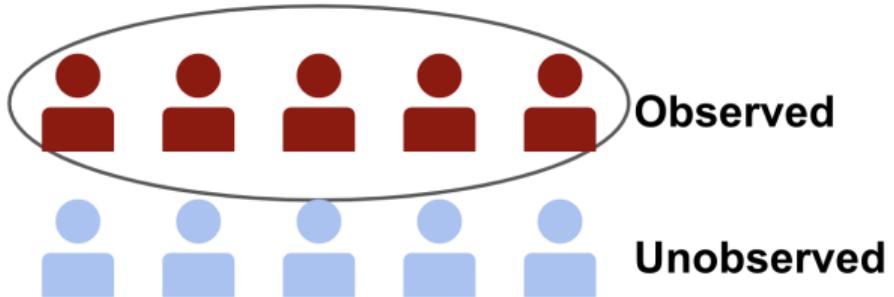
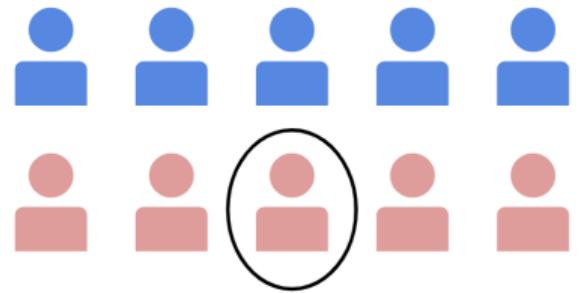


Counterfactual inference



How to infer $Y(1)$ of ?

Counterfactual inference

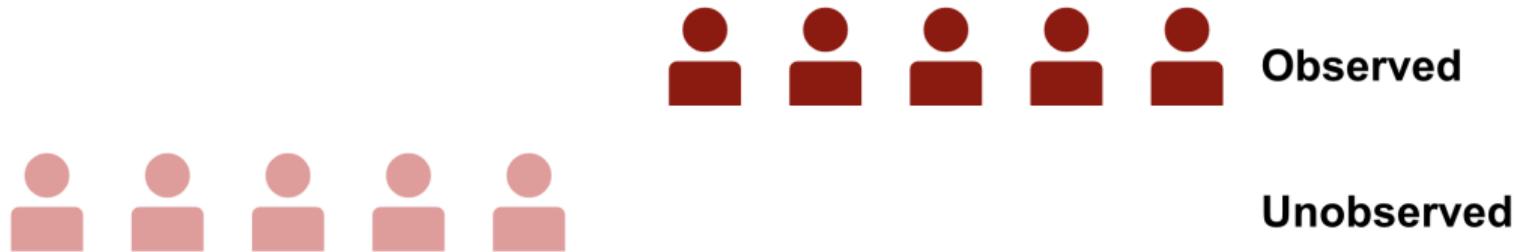


Use observed treated units



Covariate shift under unconfoundedness $Y(1) \perp\!\!\!\perp T | X$

$$P_{X|T=1} \times P_{Y(1)|X}$$

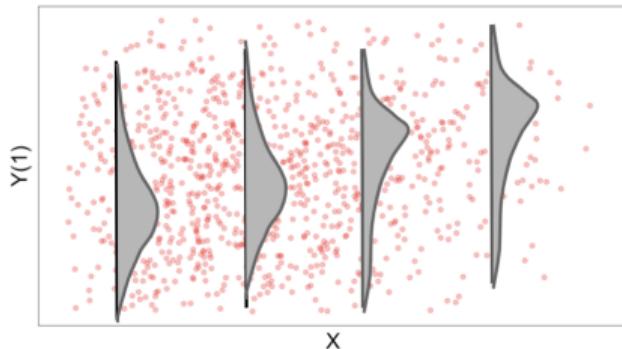


$$P_{X|T=0} \times P_{Y(1)|X}$$

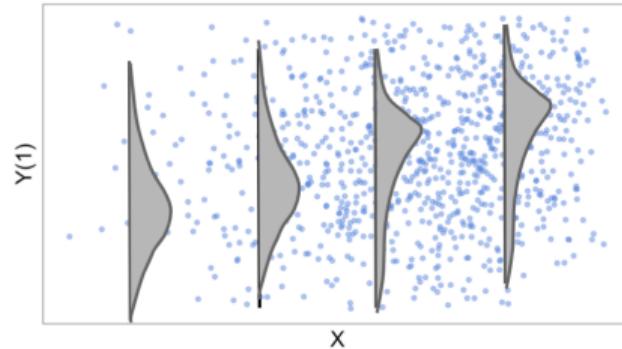
Distribution mismatch! Covariate shift

Covariate shift under unconfoundedness $Y(1) \perp\!\!\!\perp T | X$

Real world (treated units)



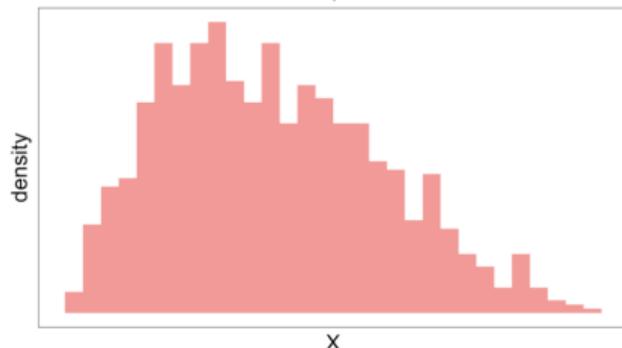
Counterfactual world



$$P_{Y(1)|X}$$



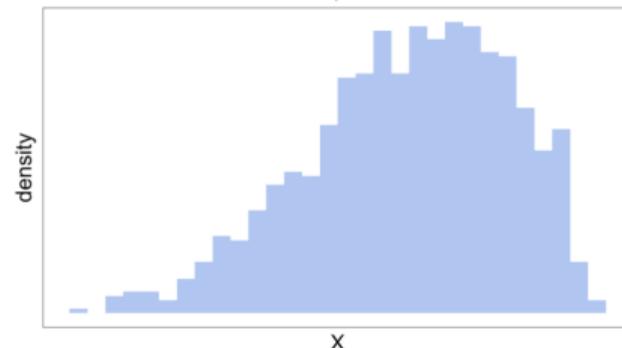
$$P_{X|T=1}$$



$$P_X$$



$$P_{X|T=0}$$



The counterfactual inference problem and covariate shift

Use i.i.d. samples (observed treated units) from $P_{X|T=1} \times P_{Y(1)|X}$ to construct $\hat{C}_1(X)$ with

$$\mathbb{P}(Y(1) \in \hat{C}_1(X)) \geq 90\% \quad \text{under } P_{X|T=0} \times P_{Y(1)|X}$$

Covariate shift $w(x) \triangleq \frac{dP_{X|T=0}}{dP_{X|T=1}}(x) \propto \frac{1 - e(x)}{e(x)}$

Conformal inference under covariate shift

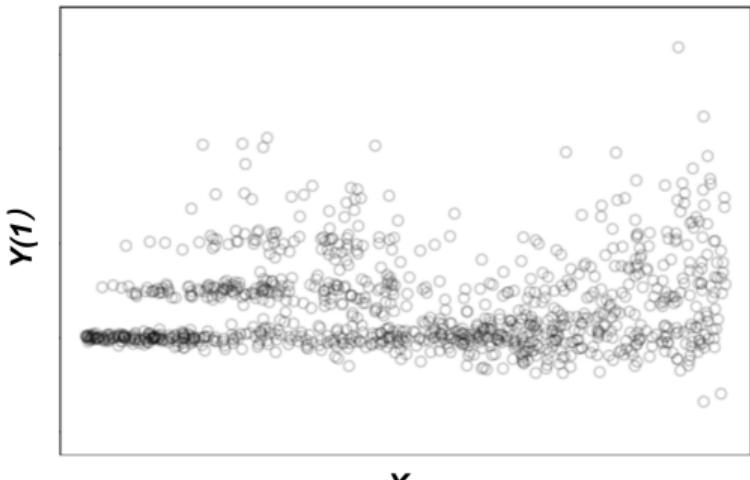
Weighted Split Conformalized Quantile Regression (CQR)

Tibshirani, Barber, Candès, Ramdas ('19); Romano, Patterson, Candès ('19)

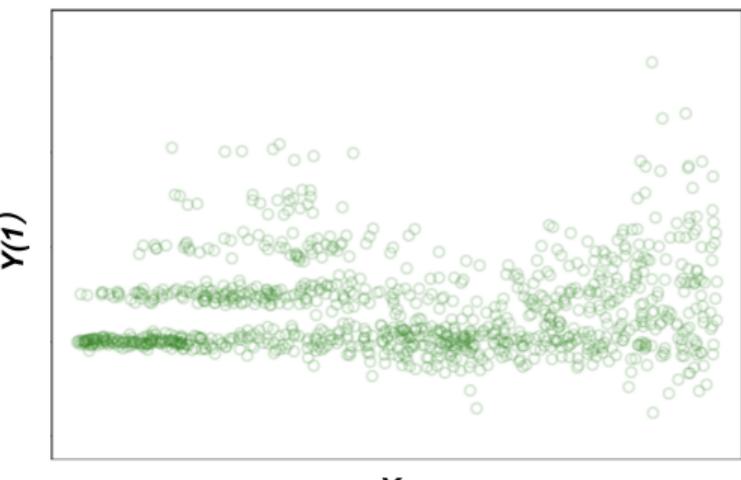
$$(X_i, Y_i) \stackrel{i.i.d.}{\sim} P_{\textcolor{red}{X}} \times P_{Y|X} \implies \mathbb{P}_{(X, Y) \sim \textcolor{blue}{Q}_X \times P_{Y|X}}(Y \in \hat{C}(X)) \geq 90\%$$

Weighted CQR

Randomly split $(X_i, Y_i^{\text{obs}})_{T_i=1}$ into two folds



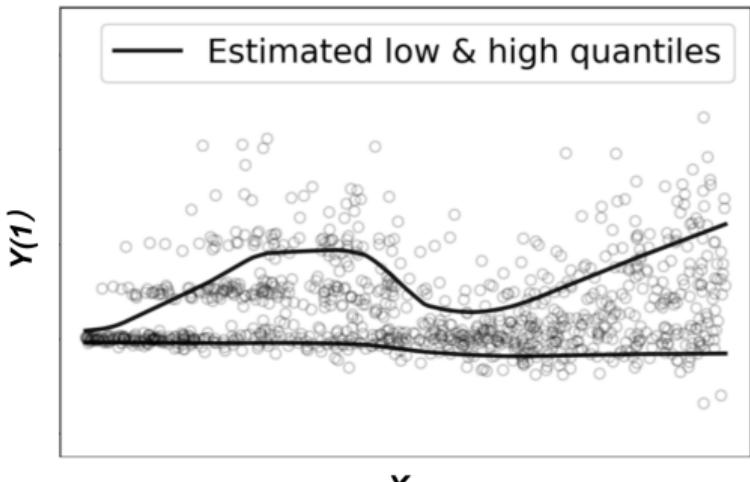
Proper training set



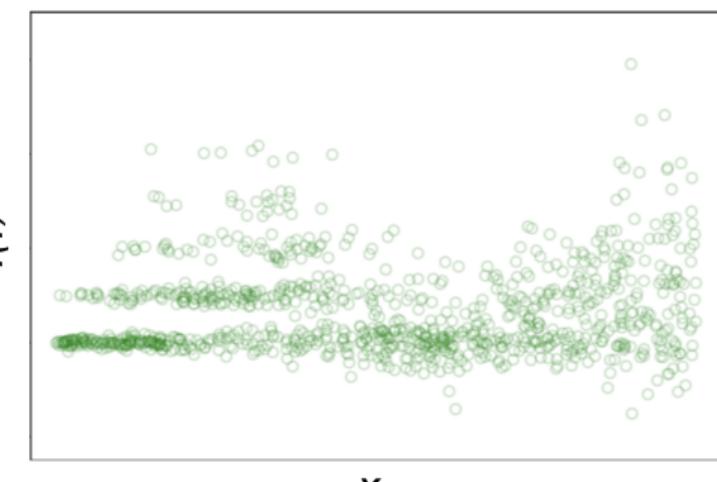
Calibration set

Weighted CQR

Fit 5 & 95%-th quantiles of $Y(1) | X$ on training fold

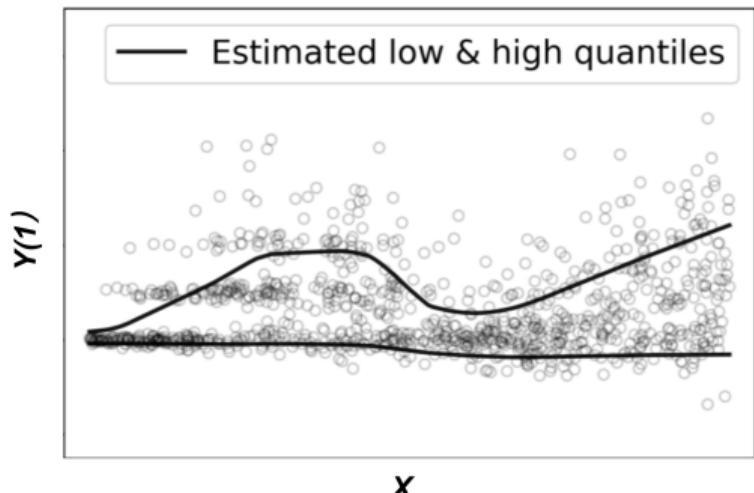


Apply quantile regression

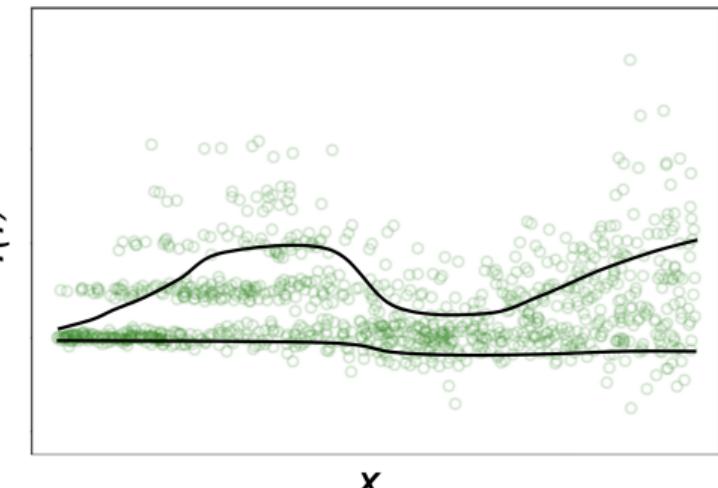


Calibration set

Weighted CQR



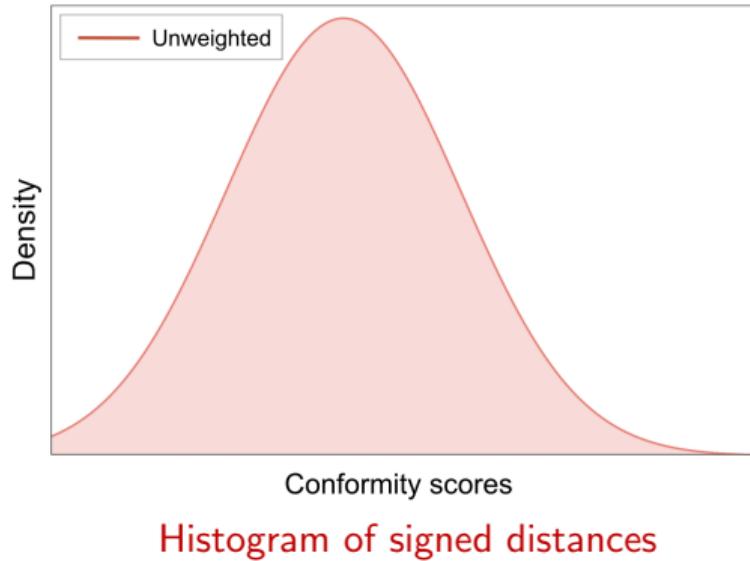
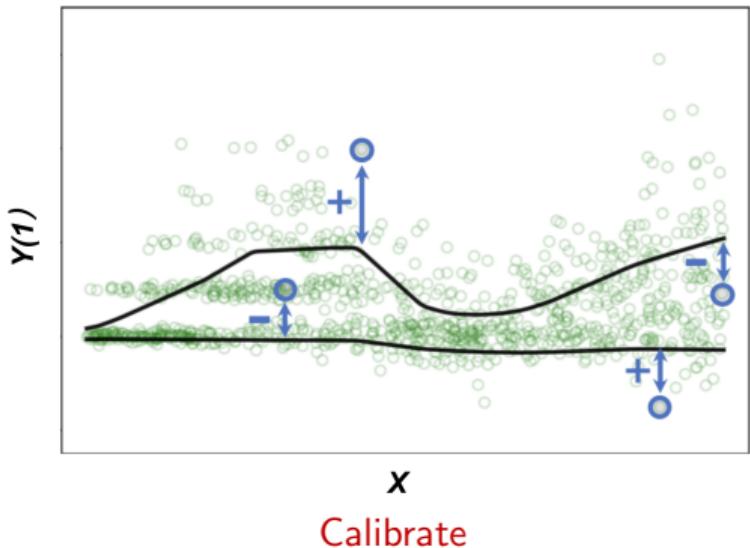
Apply quantile regression



Calibration set

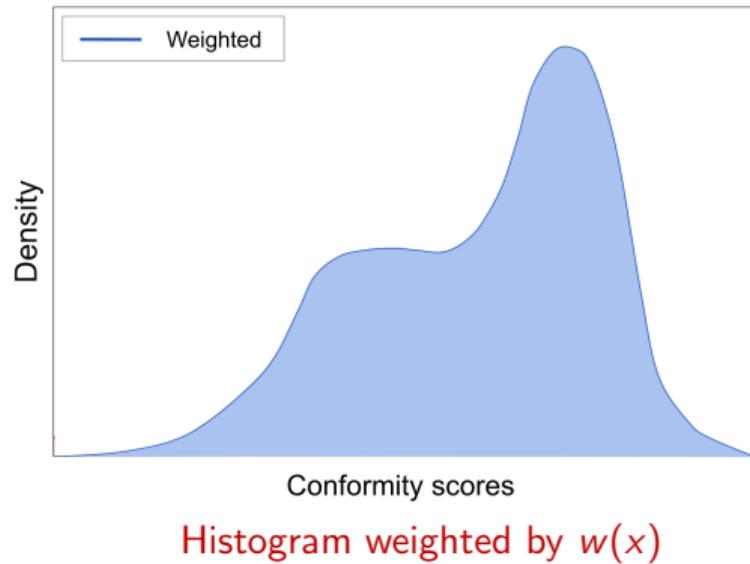
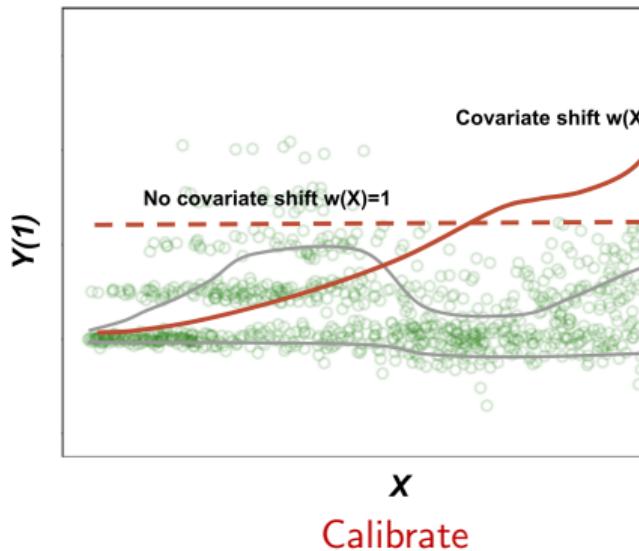
Weighted CQR

Signed distance: $V_i \triangleq \max\{\hat{q}_{0.05}(X_i) - Y_i(1), Y_i(1) - \hat{q}_{0.95}(X_i)\}$



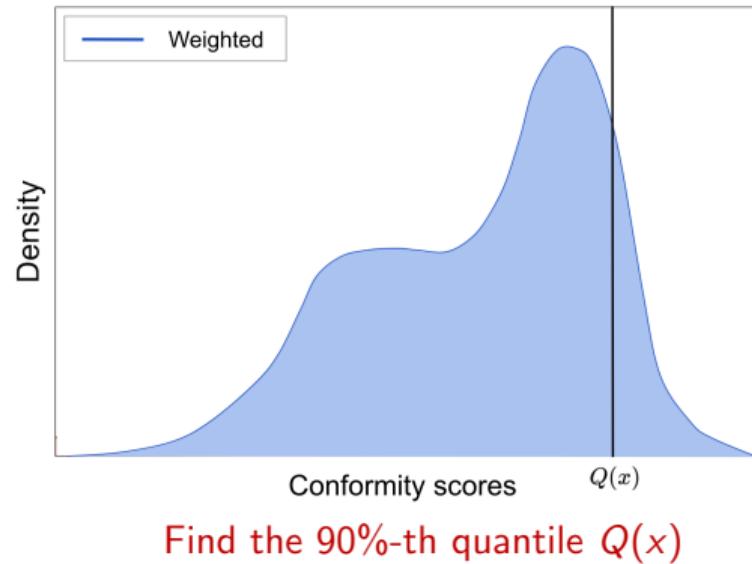
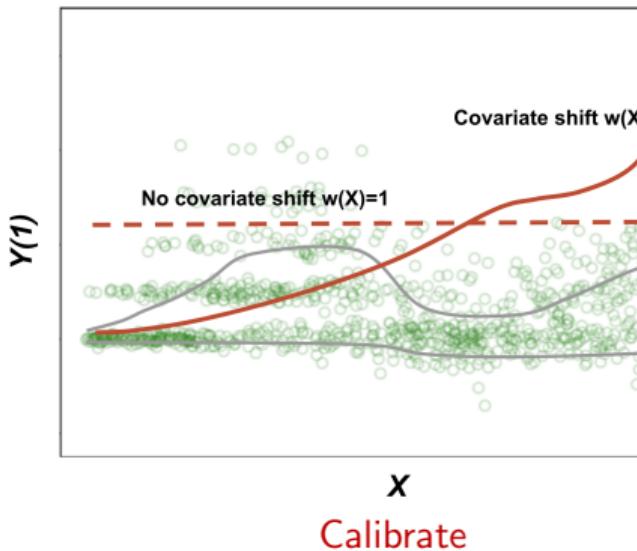
Weighted CQR

Weighted dist.: $\sum_{i=1}^n p_i(x) \delta_{V_i} + p_\infty(x) \delta_\infty$ where $p_i(x) = w(X_i) / (\sum_{i=1}^n w(X_i) + w(x))$



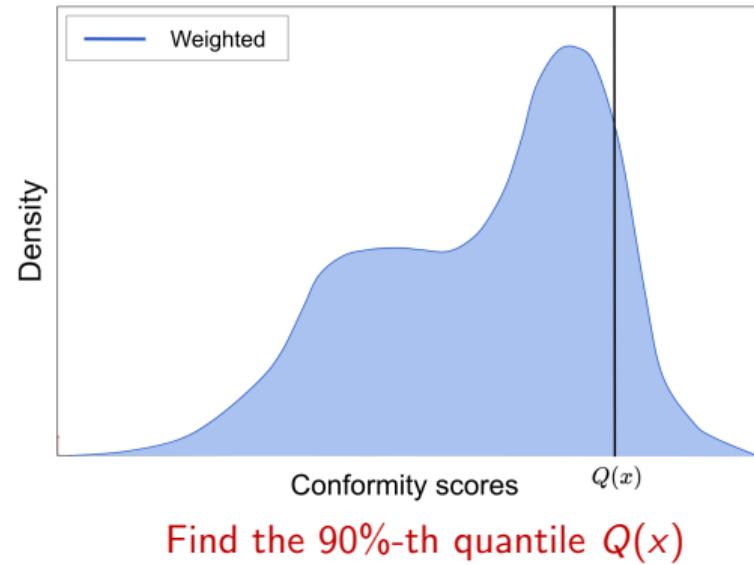
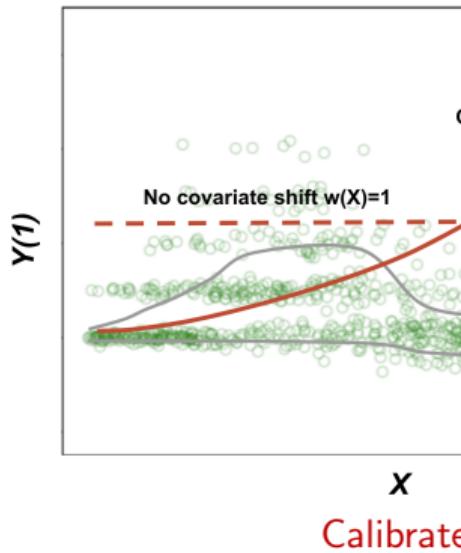
Weighted CQR

Cutoff: $Q(x) \triangleq \text{Quantile} \left(90\%, \sum_{i=1}^n p_i(x) \delta_{V_i} + p_\infty(x) \delta_\infty \right)$



Weighted CQR

$$\text{Interval: } \hat{C}_1(x) = [\hat{q}_{0.05}(x) - Q(x), \hat{q}_{0.95}(x) + Q(x)]$$



Near-exact counterfactual inference in finite samples

Theorem (L. and Candès, '20, for randomized experiments)

Set $w(x) = (1 - e(x))/e(x)$ ($e(x)$ known) in weighted split-CQR. Then

$$90\% \leq \mathbb{P}(Y(1) \in \hat{C}_1(X) \mid T = 0) \leq 90\% + c/n$$

- ▶ Lower bound holds without extra assumption
- ▶ Upper bound holds if V_i 's are a.s. distinct & overlap holds, and c only depends on the overlap

- ✓ Any conditional distribution $P_{Y(1)|X}$
- ✓ Any sample size
- ✓ Any procedure to fit conditional quantiles
- ▶ Reconcile Bayesians and frequentists (if $Q(x) \geq 0$)



Approximate counterfactual inference

Theorem (informal, L. and Candès, 2020, for observational studies)

Let $\hat{e}(x)$ be an estimate of $e(x)$. Set $w(x) = (1 - \hat{e}(x))/\hat{e}(x)$ in weighted split-CQR. Then

$$\mathbb{P}(Y(1) \in \hat{C}_1(X) \mid T = 0) \approx 90\%$$

if (1) $\hat{e}(x) \approx e(x)$ **OR** (2) $\hat{q}_{0.05/0.95}(x) \approx q_{0.05/0.95}(x)$.

Similar to the **double robustness** for ATE

Adaptivity to good outcome modelling

Theorem (informal, L. and Candès, 2020, for observational studies)

If $\hat{q}_{0.05/0.95}(x) \approx q_{0.05/0.95}(x)$,

$\hat{C}_1(x) \approx [q_{0.05}(x), q_{0.95}(x)]$ (*oracle counterfactual interval*),

and $\mathbb{P}(Y(1) \in \hat{C}_1(X) \mid T = 0, X) \approx 90\%$ with high probability (*conditional coverage!*)

- ▶ $[q_{0.05}(x), q_{0.95}(x)]$ is the optimal interval for symmetric unimodal conditional distribution
- ▶ Continue to hold if $(0.05, 0.95) \rightarrow (\beta, \beta + 1 - \alpha)$, e.g., $(0.01, 0.91)$ Romano and Sesia, '21
- ▶ Good outcome modelling \implies good intervals and conditional coverage!
- ▶ Robustness (marginal coverage) + adaptivity (efficiency and conditional coverage)

Technical conditions

Theorem (L. and Candès, '20)

Assume one of the following holds:

- (1) $\mathbb{E} |1/\hat{e}(X) - 1/e(X)| = o(1);$
- (2) $\mathbb{P}(Y(1) = y \mid X = x)$ uniformly bounded away from 0 and ∞ and there exists $\delta > 0$

$$\mathbb{E} [1/\hat{e}(X)^{1+\delta}] = O(1), \quad \mathbb{E} [H(X)/\hat{e}(X)], \mathbb{E} [H(X)/e(X)] = o(1),$$

$$\text{where } H(x) = \max\{|\hat{q}_{0.05}(x) - q_{0.05}(x)|, |\hat{q}_{0.95}(x) - q_{0.95}(x)|\}.$$

Then

$$\mathbb{P}(Y(1) \in \hat{C}_1(X) \mid T = 0) \geq 90\% - o(1).$$

Furthermore, if (2) holds, then

$$\mathbb{P}(Y(1) \in \hat{C}_1(X) \mid T = 0, X) \geq 90\% - o_{\mathbb{P}}(1).$$

From counterfactuals to individual treatment effects (ITE)



$$(X_i, Y_i^{\text{obs}}) \stackrel{i.i.d.}{\sim} P_{X|T=0} \times P_{Y(0)|X}$$



$$(X_i, Y_i^{\text{obs}}) \stackrel{i.i.d.}{\sim} P_{X|T=1} \times P_{Y(1)|X}$$



$$X \sim Q_X$$

L. and Candès, '20

Prediction interval for $\text{ITE} = Y(1) - Y(0)$ (not CATE = $\mathbb{E}[\text{ITE} | X]$)

$$\mathbb{P}_{X \sim Q_X}(\text{ITE} \in \hat{C}_{\text{ITE}}(X)) \geq 90\%$$

Our R package cfcausal (github.com/lihualei71/cfcausal)

The screenshot shows the GitHub page for the cfcausal package. At the top, there's a navigation bar with tabs for 'cfcausal' (selected), '0.2.0', 'Home', 'Reference', and 'Articles'. Below the navigation, the title 'cfcausal' is displayed in a large font. A subtitle 'An R package for conformal inference of counterfactuals and individual treatment effects' follows. The 'Overview' section contains a detailed description of the package's purpose and features, mentioning weighted conformal inference-based procedures for counterfactuals and individual treatment effects. It also highlights the split conformal inference and cross-validation+. The 'Developers' section lists Lihua Lei as the Maintainer. On the right side, there are links for 'License' (with options for 'Full license' and 'MIT + file LICENSE'), 'Citation' (with a link to 'Citing cfcausal'), and 'Install' (with a code block for installing the package). The 'Install' section includes the following R code:

```
if (!require("devtools")){
  install.packages("devtools")
}
devtools::install_github("lihualei71/cfcausal")
```

Summary for conformalized counterfactual inference

Conformal inference of counterfactuals is reliable

- ▶ Randomized experiments: **near-exact** coverage in finite samples with any black-box
- ▶ Observational studies: **doubly robust** guarantees of coverage

Part II: conformalized survival analysis

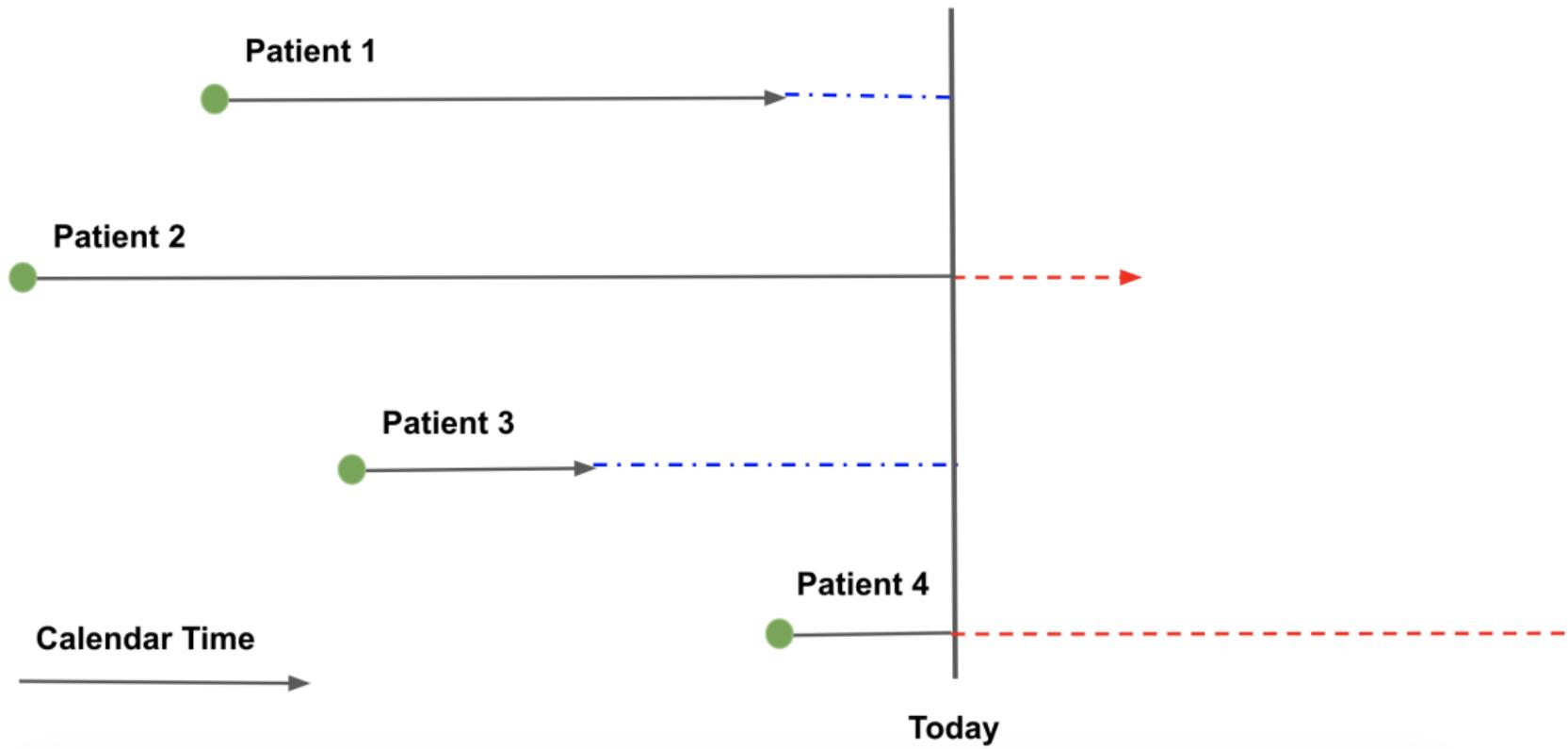


Zhimei Ren

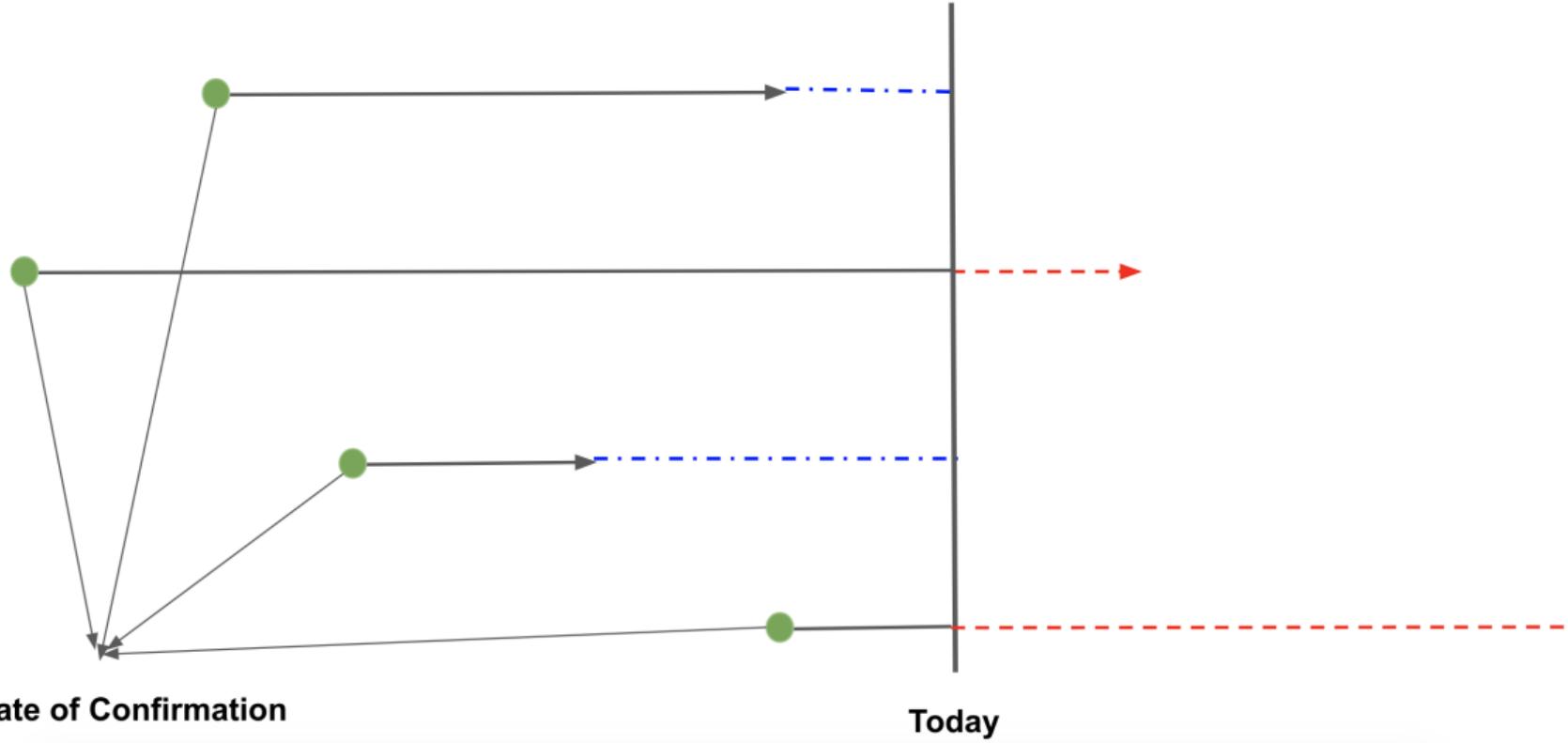


Emmanuel Candès

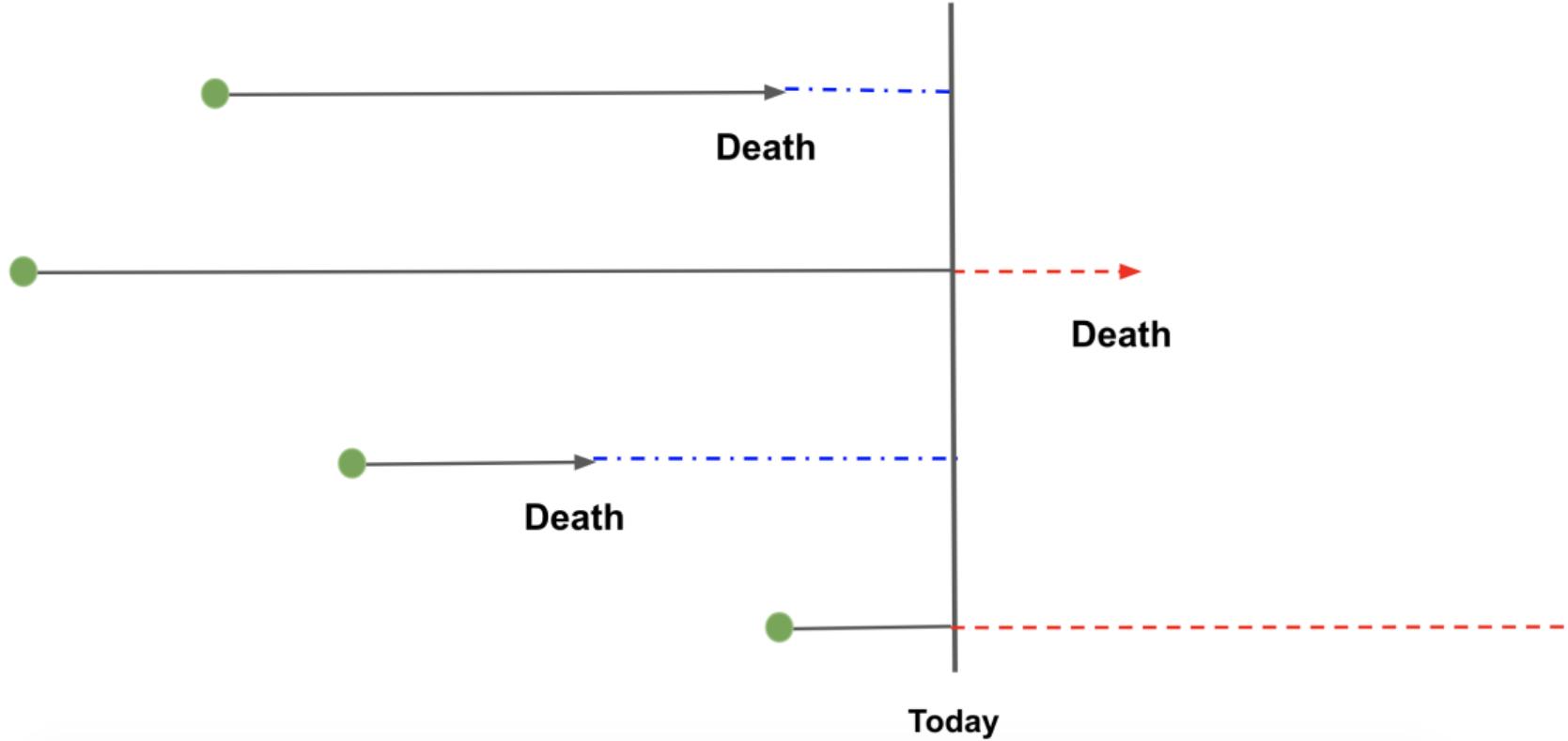
Right Censored Data: Type-I Censoring



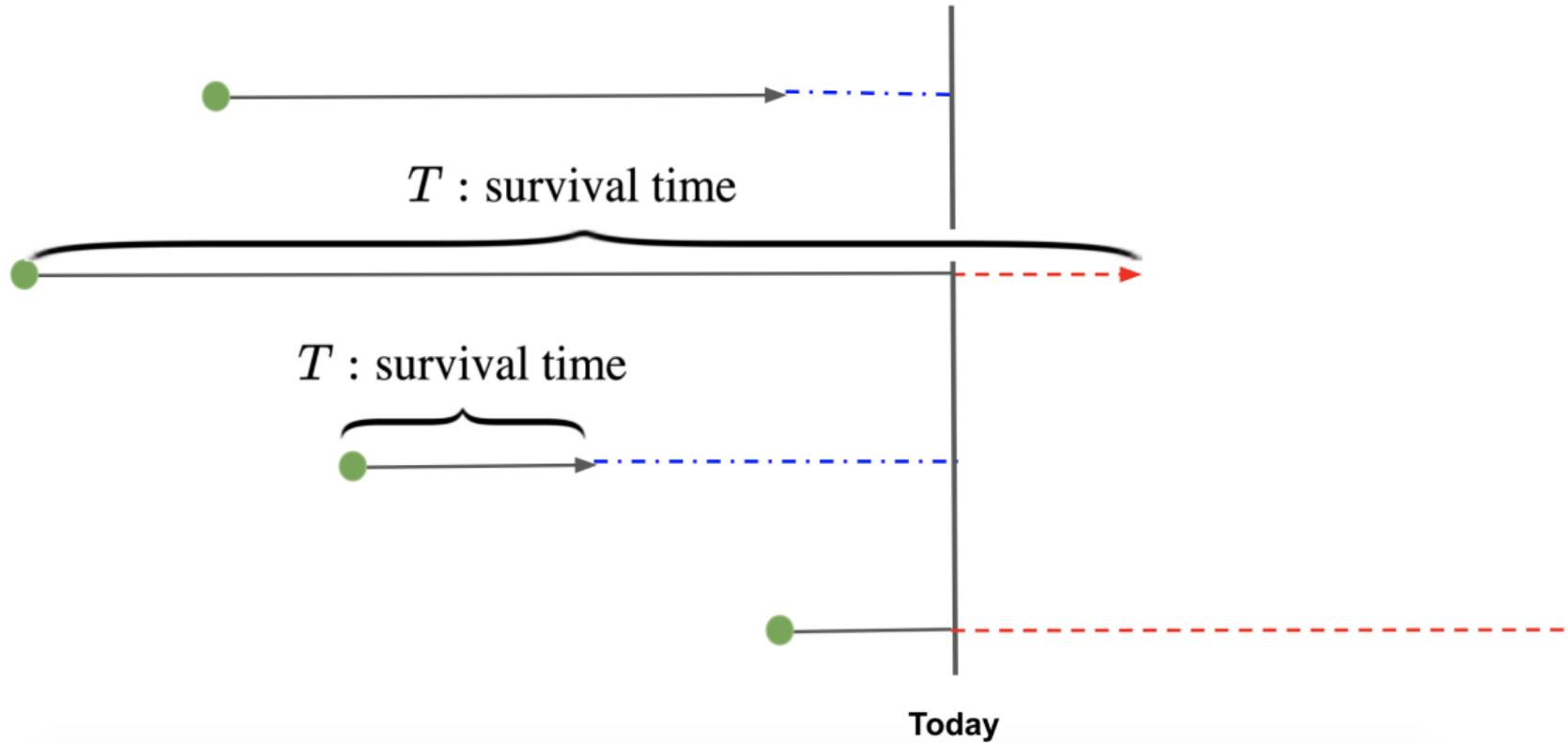
Right Censored Data: Type-I Censoring



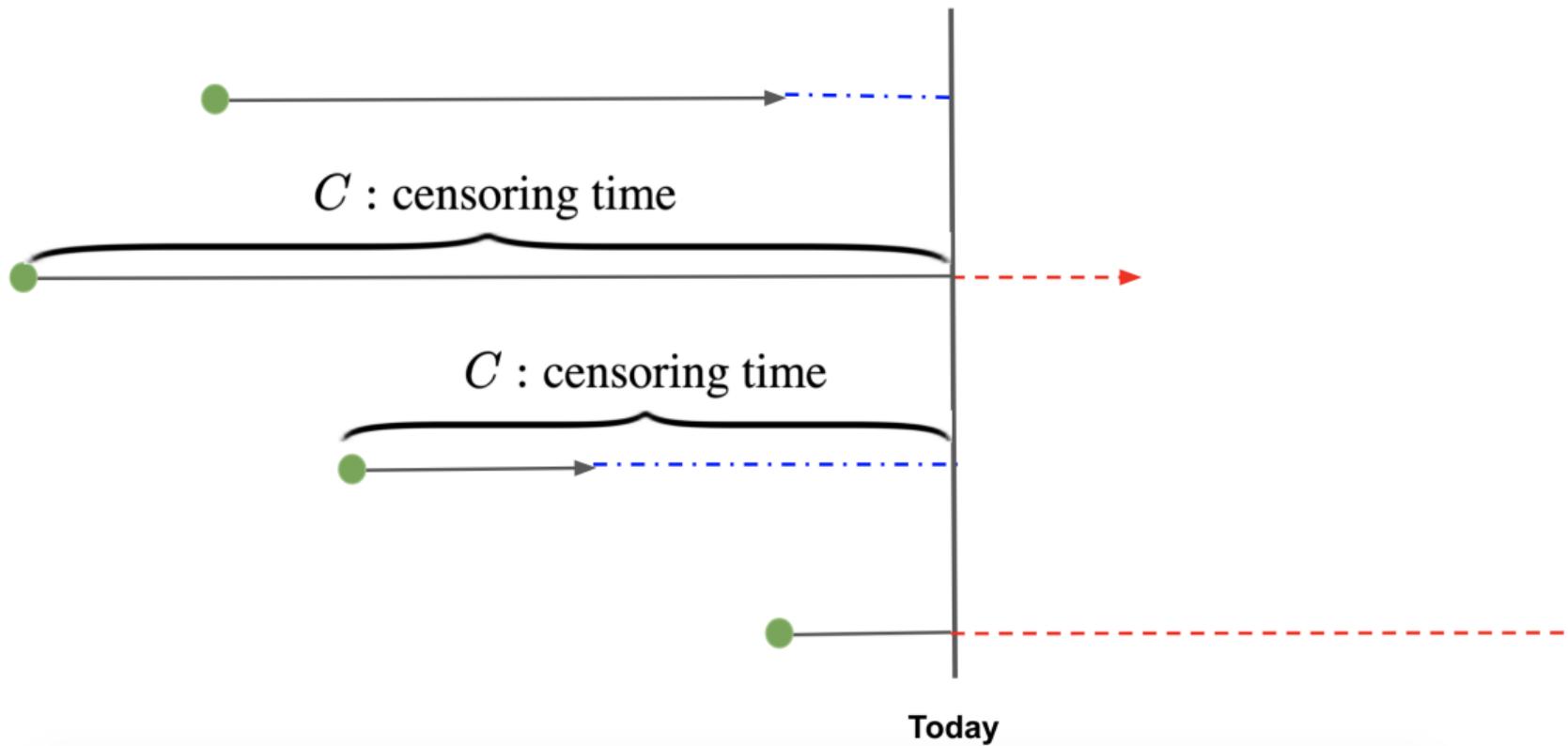
Right Censored Data: Type-I Censoring



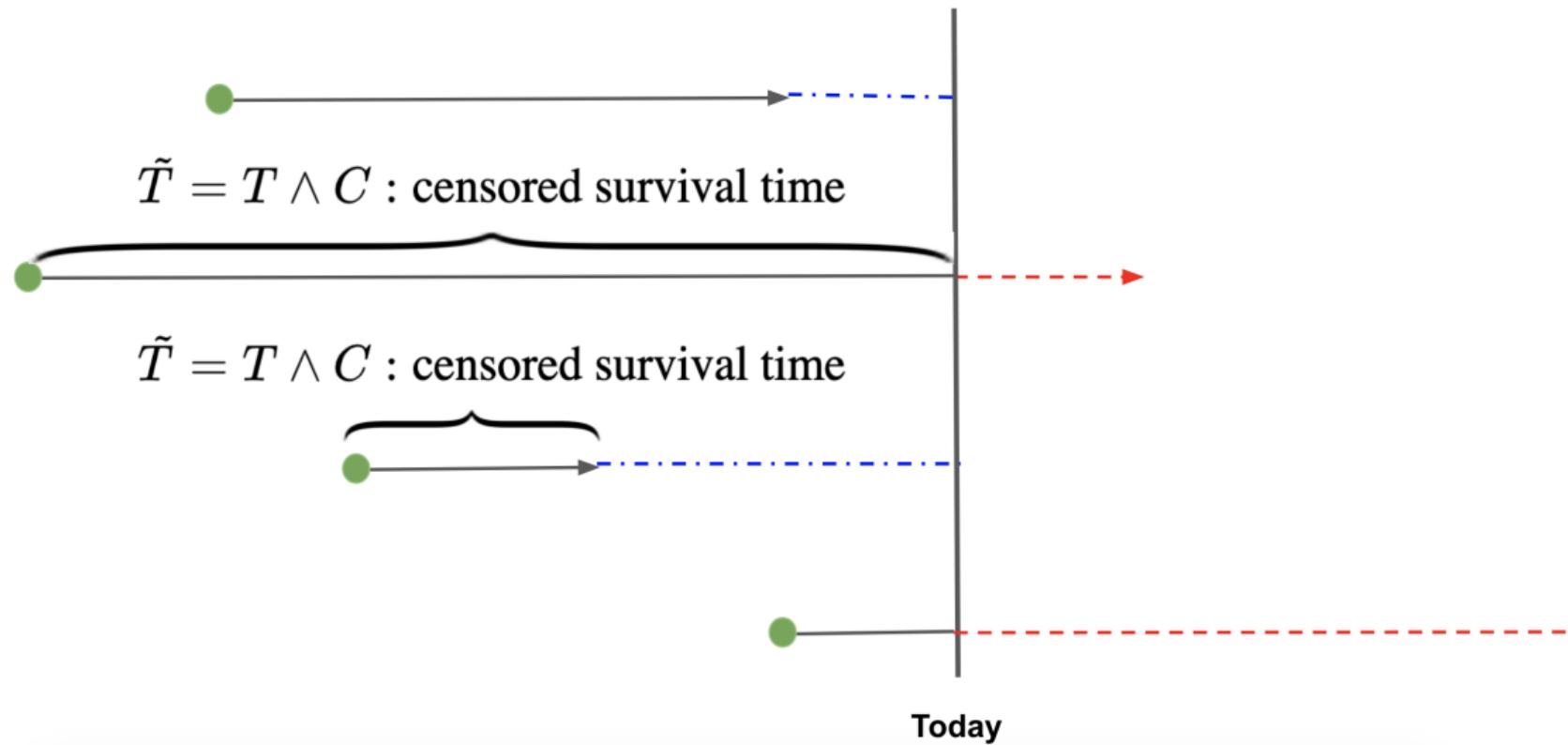
Right Censored Data: Type-I Censoring



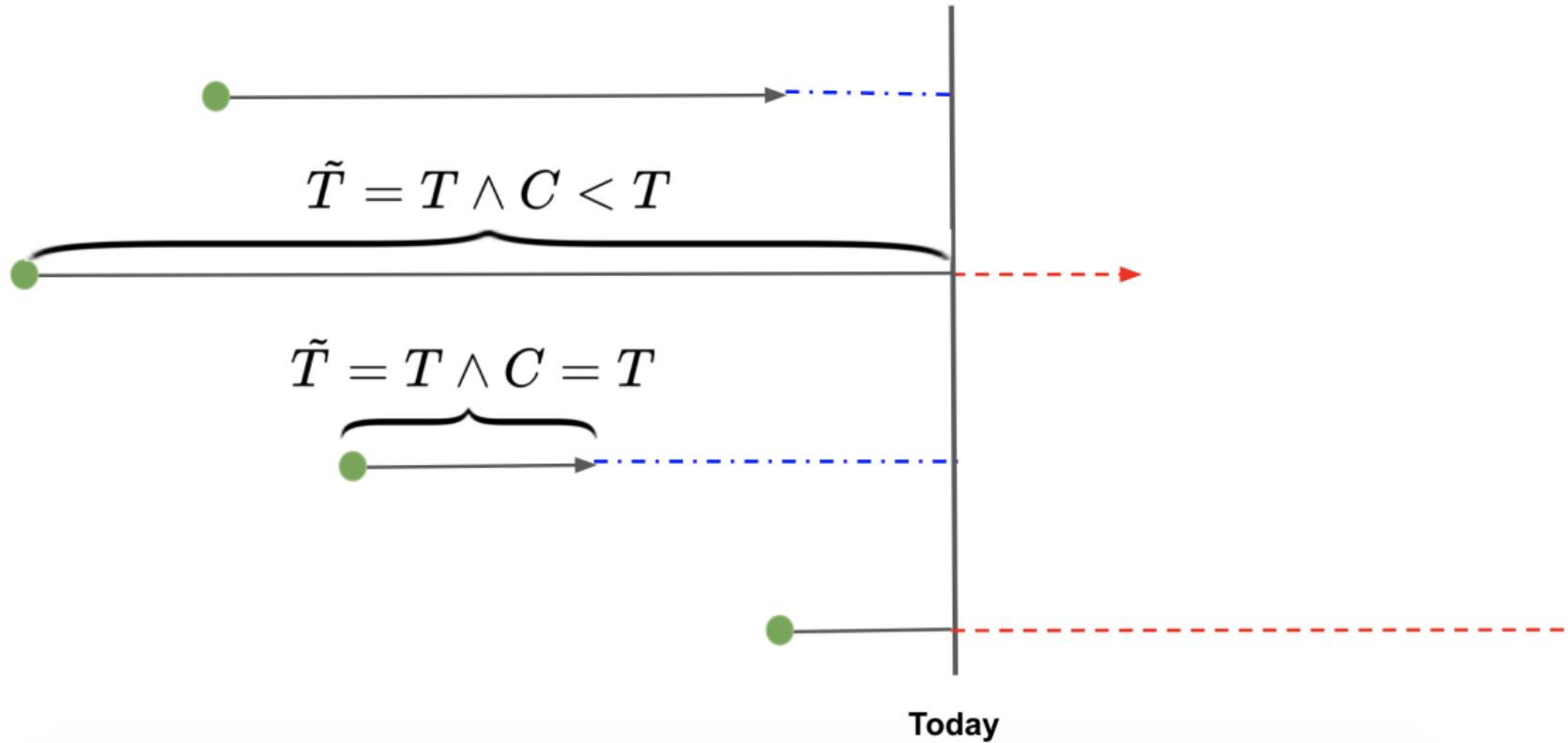
Right Censored Data: Type-I Censoring



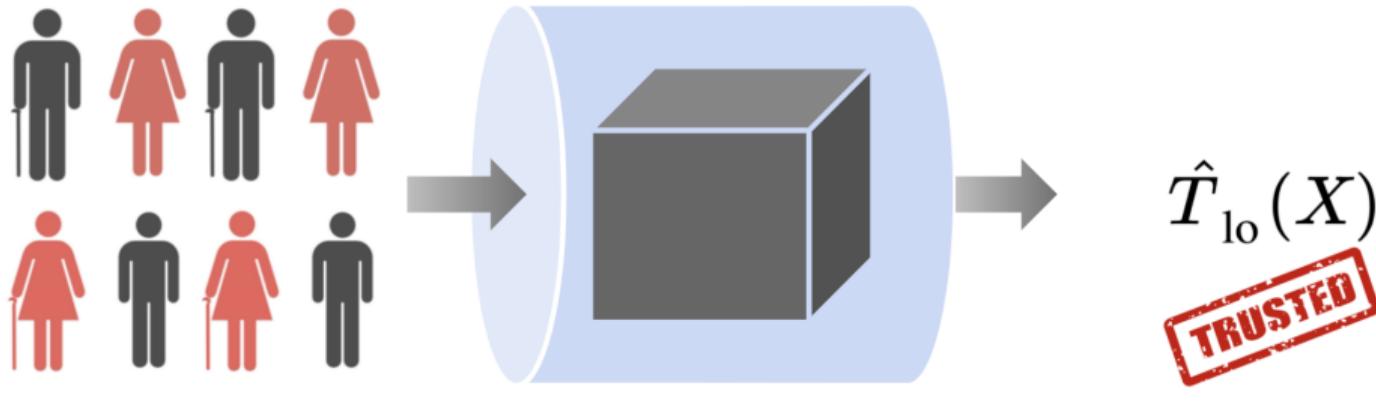
Right Censored Data: Type-I Censoring



Right Censored Data: Type-I Censoring



A reliable predictive system for survival times



Find lower predictive bound $\hat{T}_{\text{lo}}(X)$, s.t. $\mathbb{P}(T \geq \hat{T}_{\text{lo}}(X)) \geq 90\%$

First thought: survival times as counterfactuals?

- ▶ Event indicator $\Delta = I(T < C)$:

$$\tilde{T} = \begin{cases} T & \text{if } \Delta = 1 \\ C & \text{if } \Delta = 0 \end{cases} .$$

- ▶ Treat T as a “potential outcome” under the “treatment” $\Delta = 1$?

First thought: survival times as counterfactuals?

- ▶ Event indicator $\Delta = I(T < C)$:

$$\tilde{T} = \begin{cases} T & \text{if } \Delta = 1 \\ C & \text{if } \Delta = 0 \end{cases}.$$

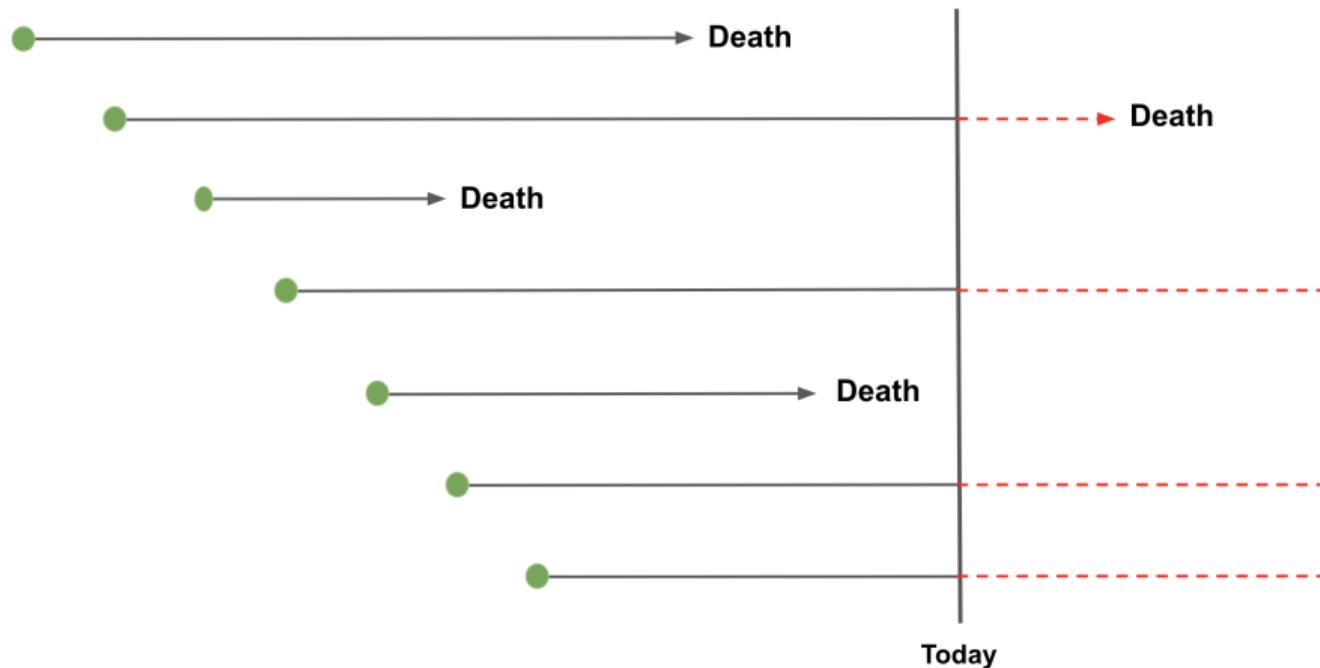
- ▶ Treat T as a “potential outcome” under the “treatment” $\Delta = 1$?
- ▶ **Invalid** because “unconfoundedness” does not hold:

$$(T, C) \not\perp\!\!\!\perp I(T < C) \mid X$$

- ▶ $(X_i, T_i)_{\Delta_i=1}$ has shifts in both the covariate distribution and conditional distribution

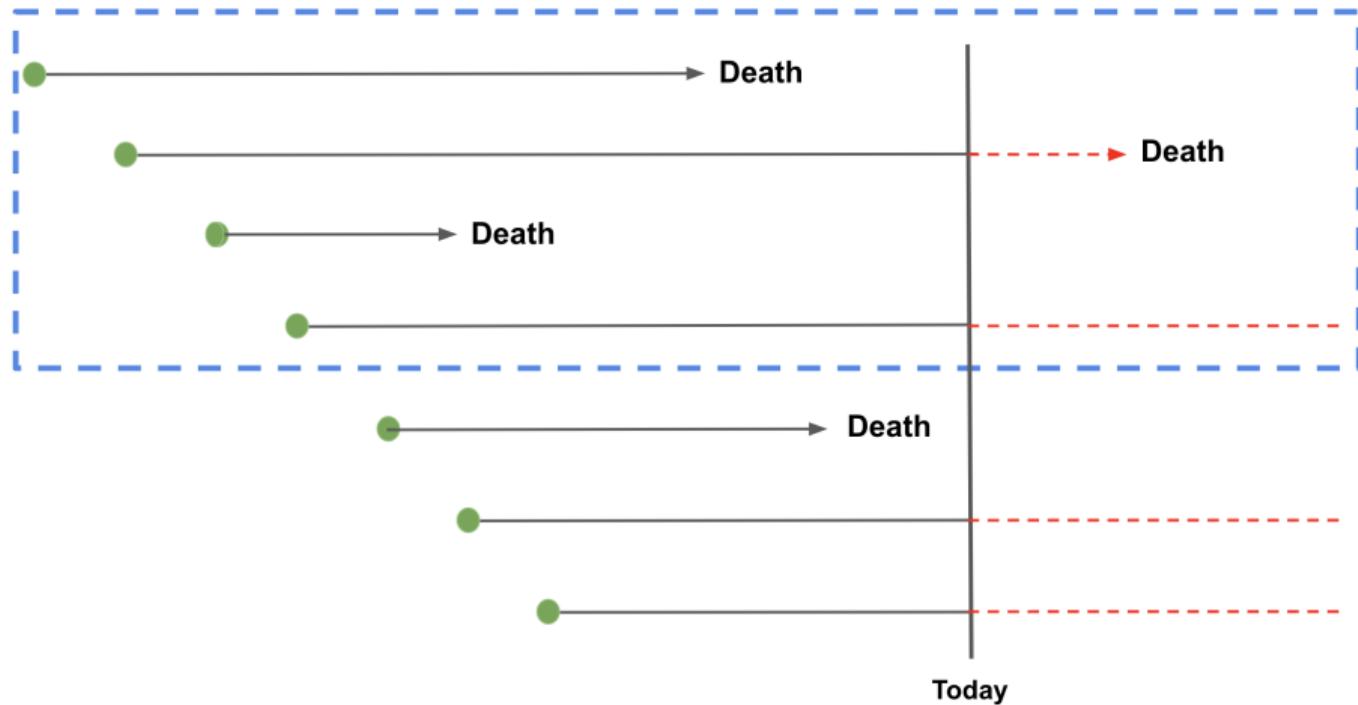
Conformalized survival analysis

Order the units by censoring times C_i



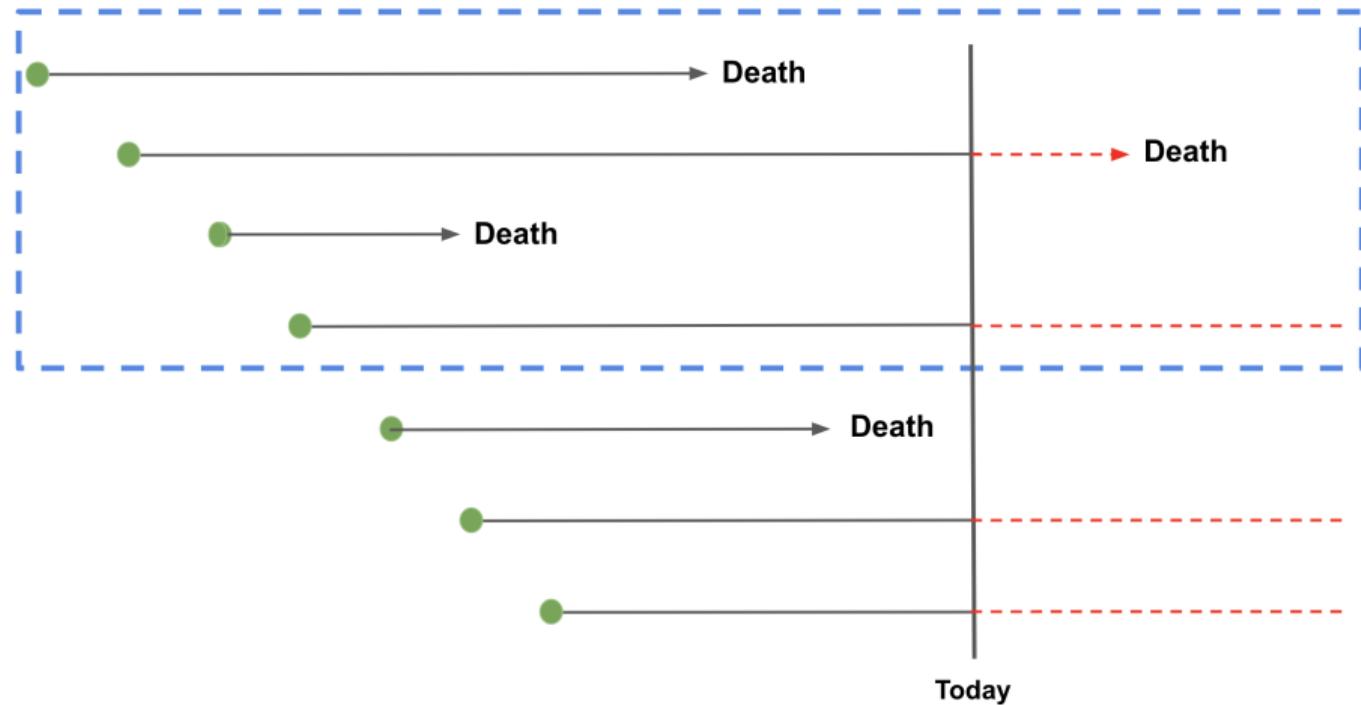
Conformalized survival analysis

Study population: $C_i \geq c_0$ (c_0 chosen via data splitting)



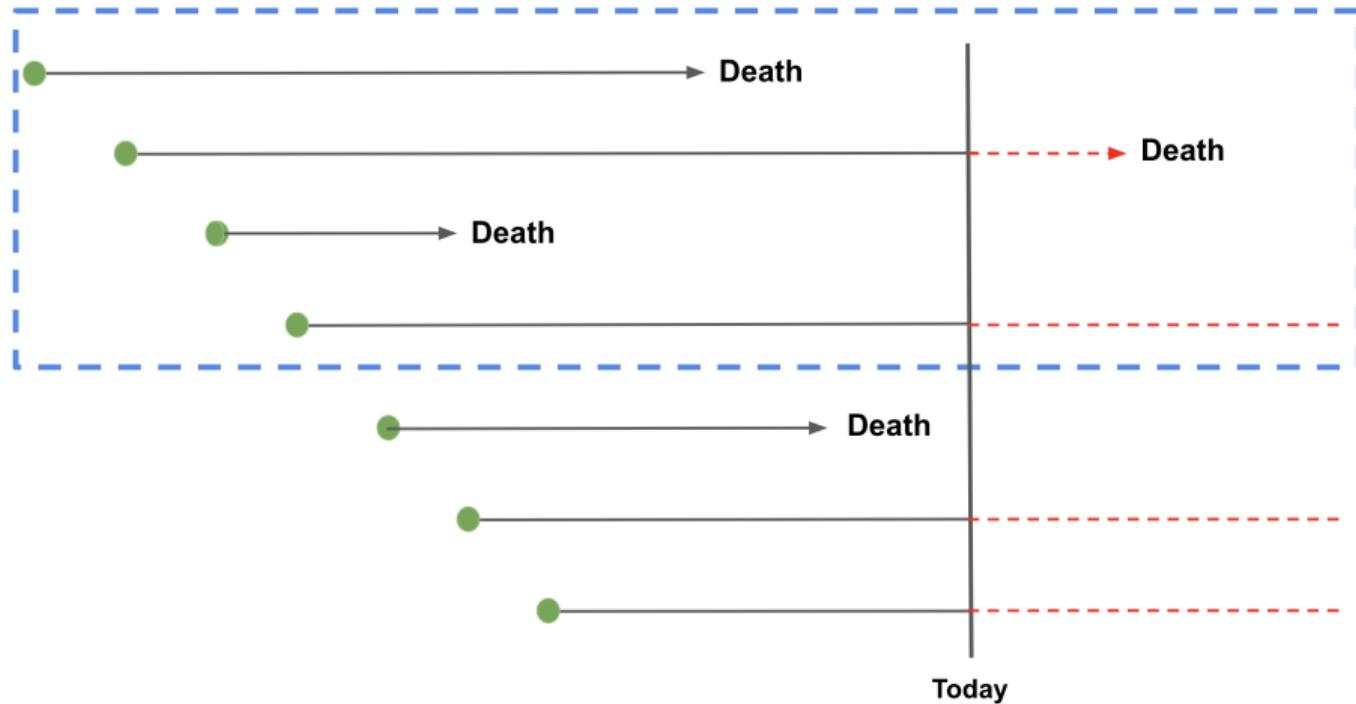
Conformalized survival analysis

On this population ($C \geq c_0$), the **surrogate outcome** $T \wedge c_0 = \tilde{T} \wedge c_0$ is always observable



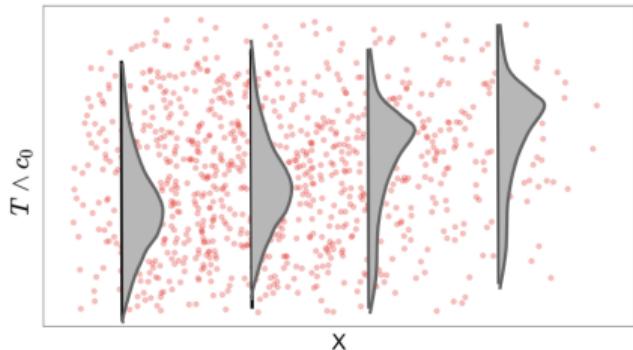
Conformalized survival analysis

$$\mathbb{P}(T \wedge c_0 \geq \hat{T}_{\text{lo}}(X)) \geq 90\% \implies \mathbb{P}(T \geq \hat{T}_{\text{lo}}(X)) \geq 90\%$$



Covariate shift under conditionally independent censoring $T \perp\!\!\!\perp C | X$

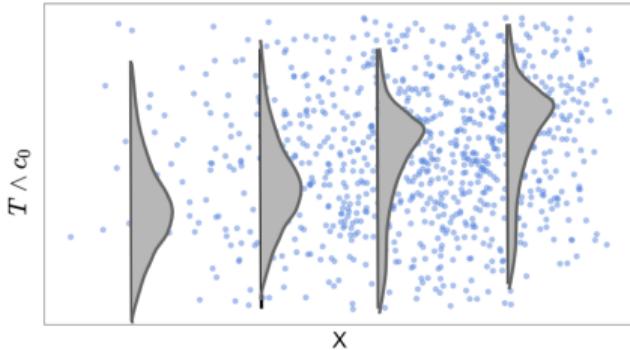
Subpopulation with $C \geq c_0$



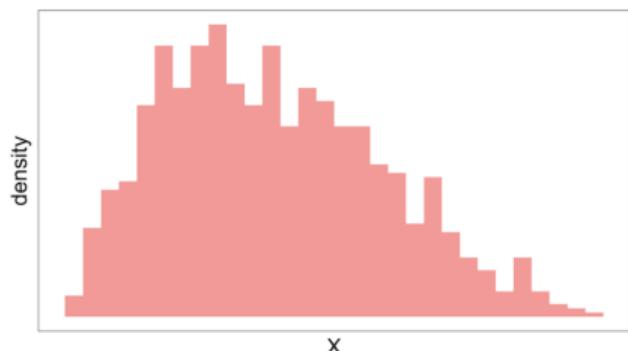
$$P_{T \wedge c_0 | X}$$



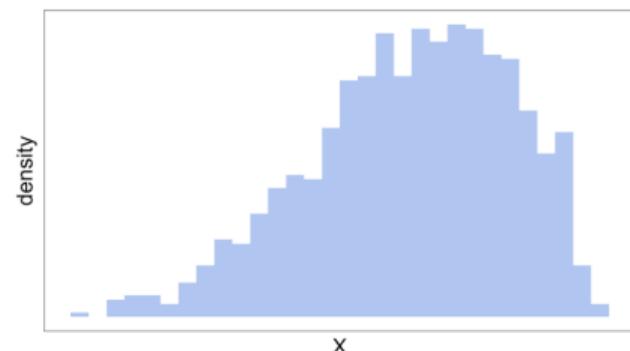
Target population



$$P_X$$



$$P_X$$



Conformalized survival analysis

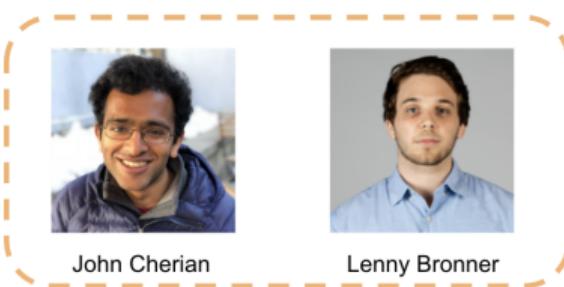
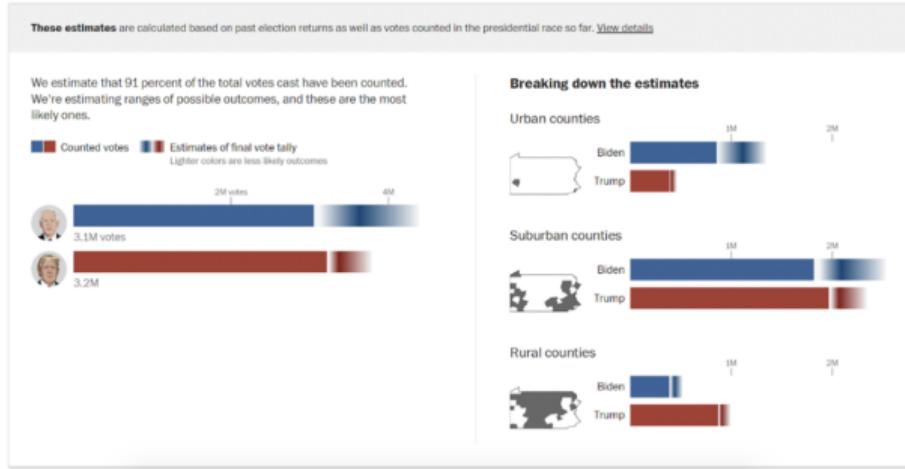
- ▶ Assumptions:
 - ▶ (T_i, C_i, X_i) are i.i.d.
 - ▶ Conditionally independent censoring ($T \perp\!\!\!\perp C | X$)
 - ▶ Type-I censoring (also useful beyond this setting)
- ▶ Finite-sample validity: $\hat{T}_{lo}(X)$ is valid if $P(C | X)$ is known
- ▶ Double robustness: $\hat{T}_{lo}(X)$ is approximately valid if $P(C | X)$ or $P(T | X)$ is estimated well

Part III: what else?

Election night model: prediction intervals for aggregated outcomes

The Washington Post Election Night Model

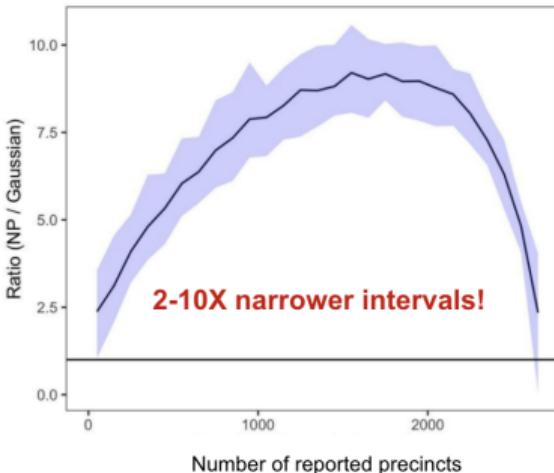
Where the vote could end up



Georgia Senate Runoff Election

Parametric conformal inference

Ratio of average length

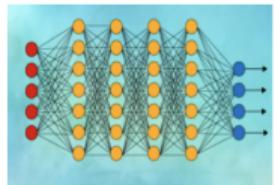
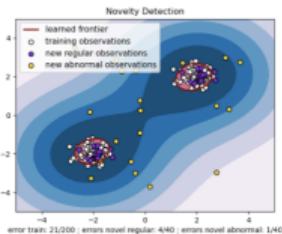
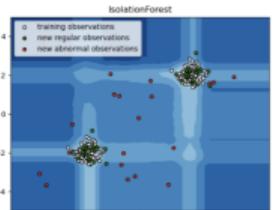
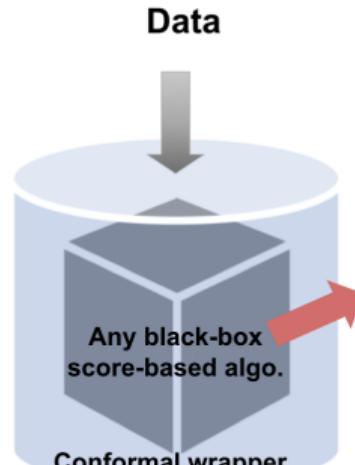
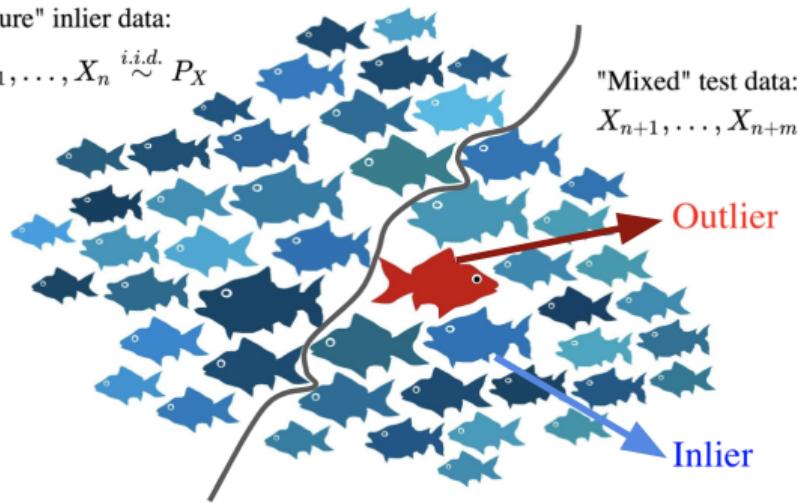


$$\mathbb{P} \left(\underbrace{\sum_{i:\text{precincts}} \text{Votes}_i}_{\text{aggregated outcome}} \in \hat{C}(\{X_i\}_{\text{precincts}}) \right) \geq 90\%$$

Outlier detection with conformal p-values

"Pure" inlier data:

$$X_1, \dots, X_n \stackrel{i.i.d.}{\sim} P_X$$



Stephen Bates



Emmanuel Candès

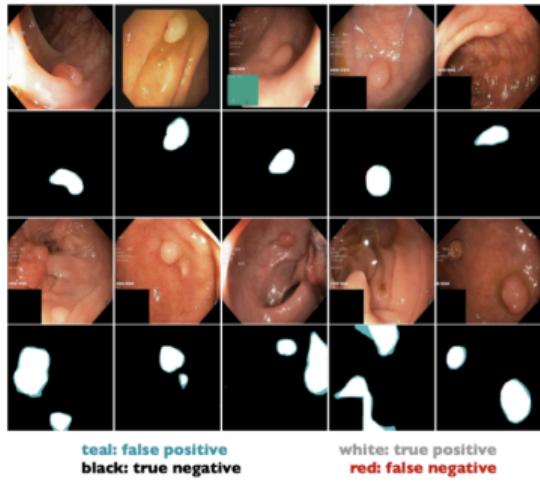


Yaniv Romano



Matteo Sesia

Risk calibrated prediction



Tumor/Polyp detection



black: correct label
red: raw output
teal: our output



Hierarchical classification



Anastasios
Angelopoulos

Stephen Bates

Emmanuel Candès

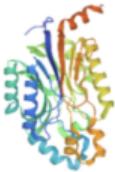
Michael Jordan

Jitendra Malik

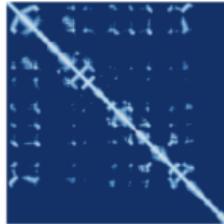
$$\mathbb{P} \left(\underbrace{\mathbb{E}[L(Y, \hat{Y})]}_{\text{risk}} \leq \gamma \right) \geq 90\%$$

Risk calibrated prediction

3D rendering of T0995



prediction (Å)



length of lower interval (Å)



length of upper interval (Å)



Protein structure prediction



Object detection

$$\mathbb{P} \left(\underbrace{\mathbb{E}[L(Y, \hat{Y})]}_{\text{risk}} \leq \gamma \right) \geq 90\%$$



Anastasios
Angelopoulos



Stephen Bates



Emmanuel Candès



Michael Jordan

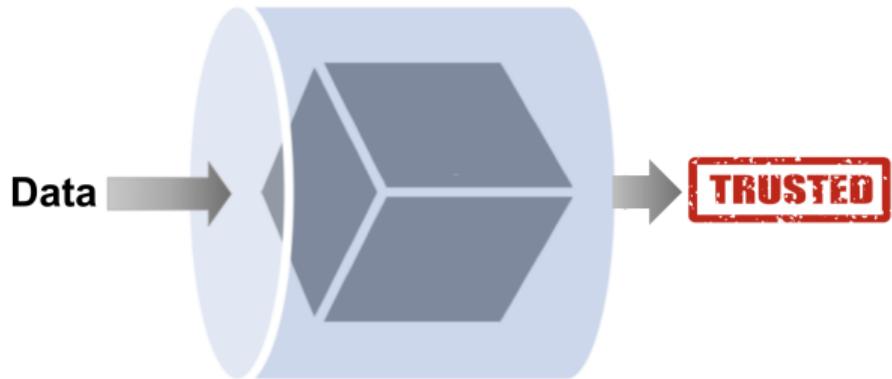


Jitendra Malik

What can conformal inference offer to statistics?

What can conformal inference offer to statistics? A LOT!

- ▶ Causal inference
- ▶ Time-to-event analysis
- ▶ Election night model
- ▶ Outlier detection
- ▶ Risk calibration
- ▶ ...



Valid inference under partial/complete misspecification!

Long-term career goal: principles for inference under misspecification

- ▶ **Network:** (hierarchical) clustering under misspecified SBMs

w/ Tianxi Li, Sharmo Bhattacharyya, Purna Sarkar, Peter Bickel, Liza Levina, Xiaodong Li, Xingmei Lou

- ▶ **Multiple testing:** FDR control with side information

w/ Will Fithian, Aaditya Ramdas, Chiao-Yu Yang, Nhat Ho, Yixiang Luo

- ▶ **Causal inference:**

- ▶ **Randomized experiments:** model-free regression adjustment

w/ Peng Ding

- ▶ **Observational studies:** distribution-free assessment of overlap

w/ Alex D'Amour, Peng Ding, Avi Feller, Jas Sekhon

- ▶ **High-d inference:** finite-sample valid test for linear models with exchangeable errors

w/ Peter Bickel

- ▶ **Econometrics:** panel data analysis under heterogeneous treatment effects

w/ Dmitry Arkhangelsky, Guido Imbens, Xiaoman Luo

All models are wrong, but some are (hopefully) useful

All models are wrong, but we can make them safe and useful!

All models are wrong, but we can make them safe and useful!

Thank you!

Check out my CV and other works on my website!

lihualei71.github.io