# Conformal Meta-learners for Predictive Inference of Individual Treatment Effects

Ahmed M. Alaa

UC Berkeley and UCSF

amalaa@berkeley.edu

Zaid Ahmad

UC Berkeley

zaidahmad@berkeley.edu
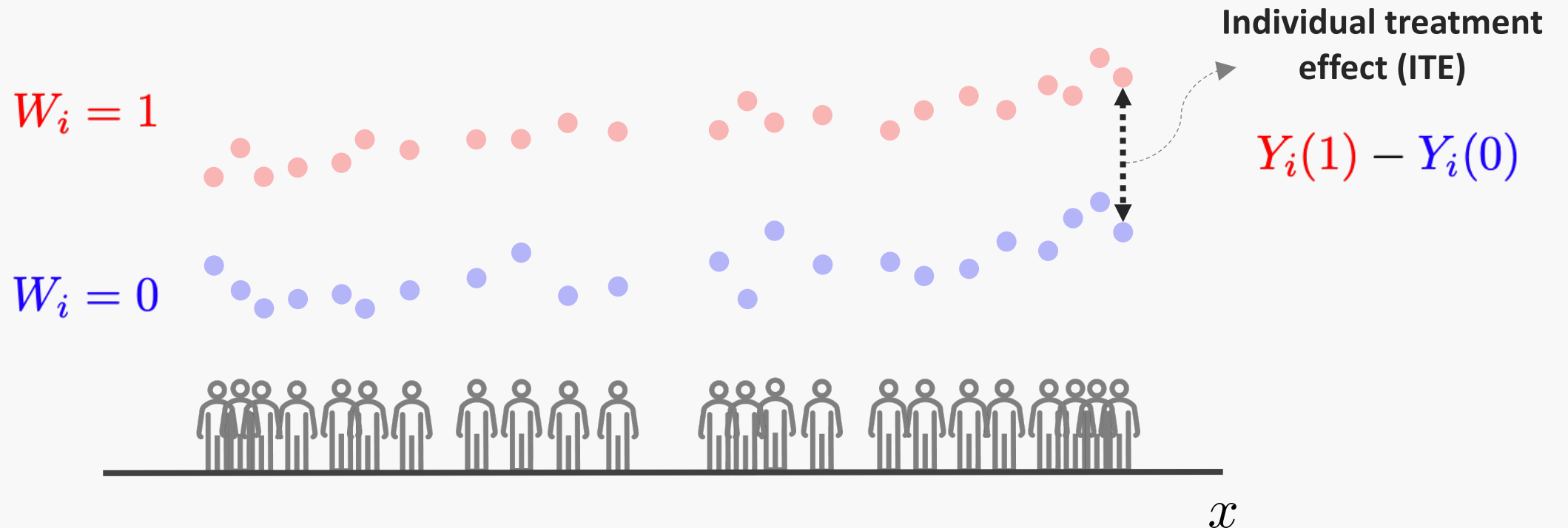
Mark van der Laan

UC Berkeley

laan@berkeley.edu

# Setup: Potential Outcomes Framework

(Neyman 1923; D. Rubin 1974)

- Binary treatment $W_i \in \{0, 1\}$ → two potential outcomes: $Y_i(1)$ and $Y_i(0)$

$$W_i = 1$$

$$W_i = 0$$

**Individual treatment effect (ITE)**

$$Y_i(1) - Y_i(0)$$

$$x$$

# Setup: Potential Outcomes Framework

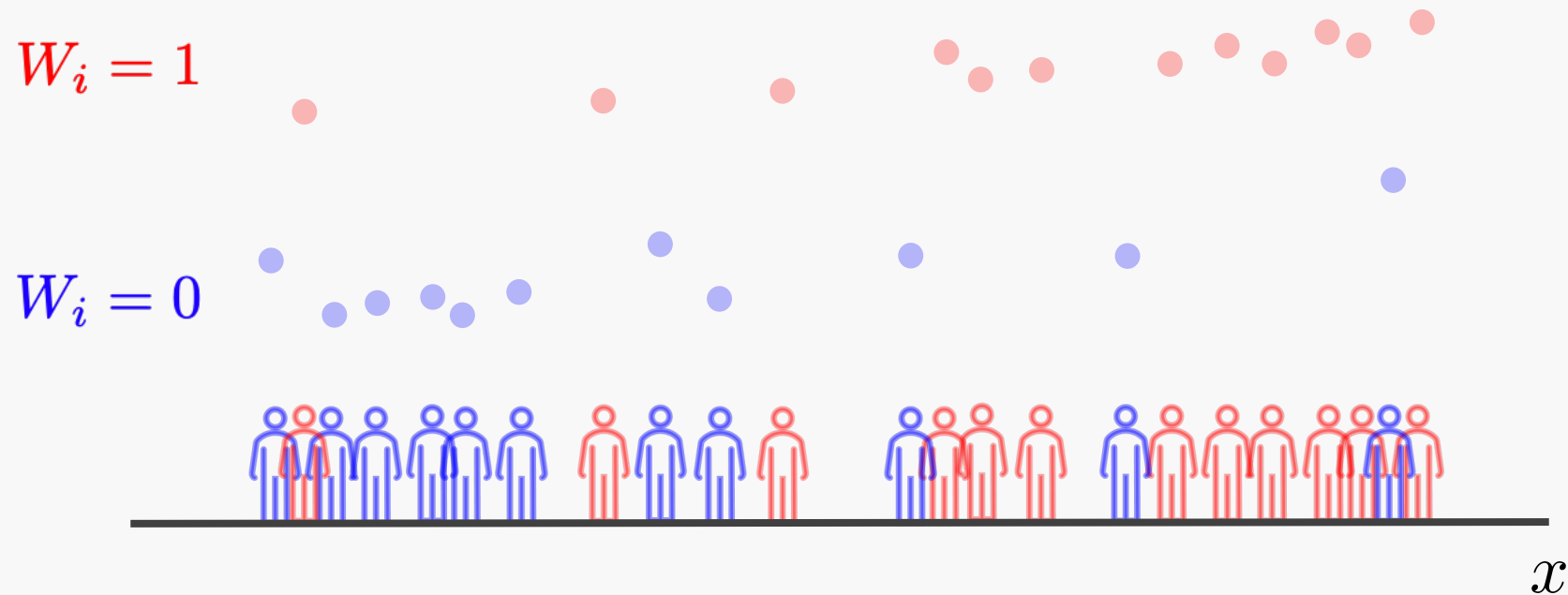- **The fundamental problem of causal inference:** Counterfactuals are not observed!

$$Y_i = W_i Y_i(1) + (1 - W_i) Y_i(0)$$

$W_i = 1$

$W_i = 0$

$x$

# Setup: Potential Outcomes Framework

- **Treatments not randomly assigned:** Propensity score $\rightarrow \pi(x) = P(W = 1 | X = x)$
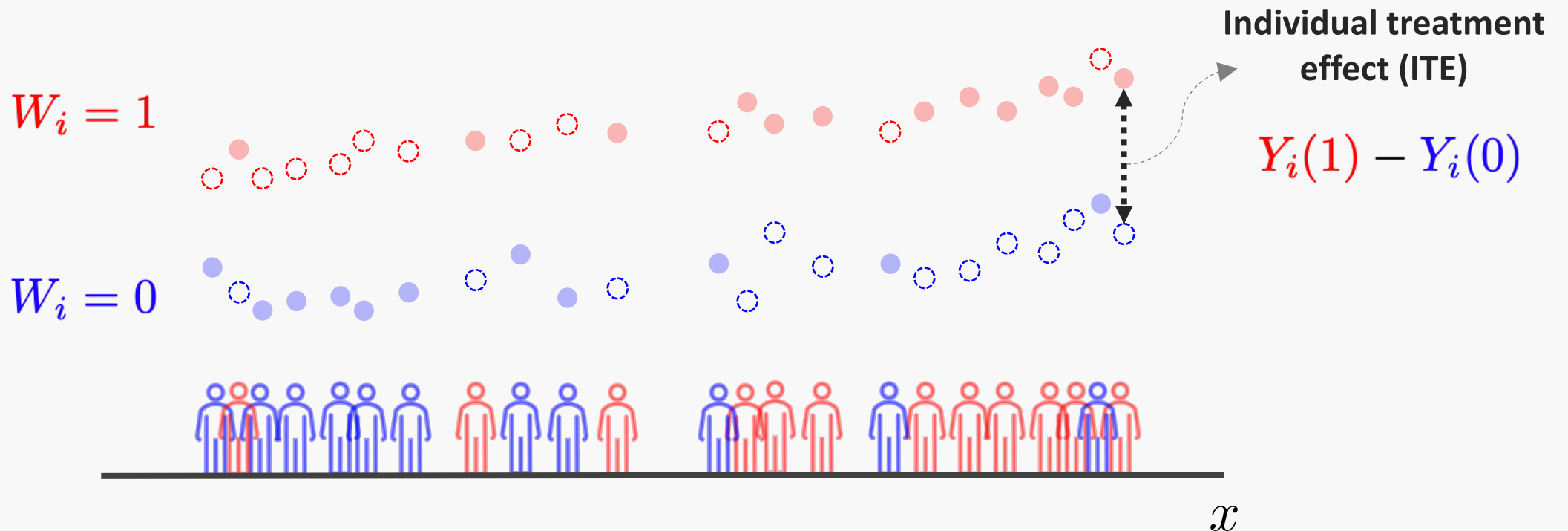
$$Y_i = W_i Y_i(1) + (1 - W_i) Y_i(0)$$



$W_i = 1$

$W_i = 0$

$x$

# Problem: Valid Predictive Inference on ITEs

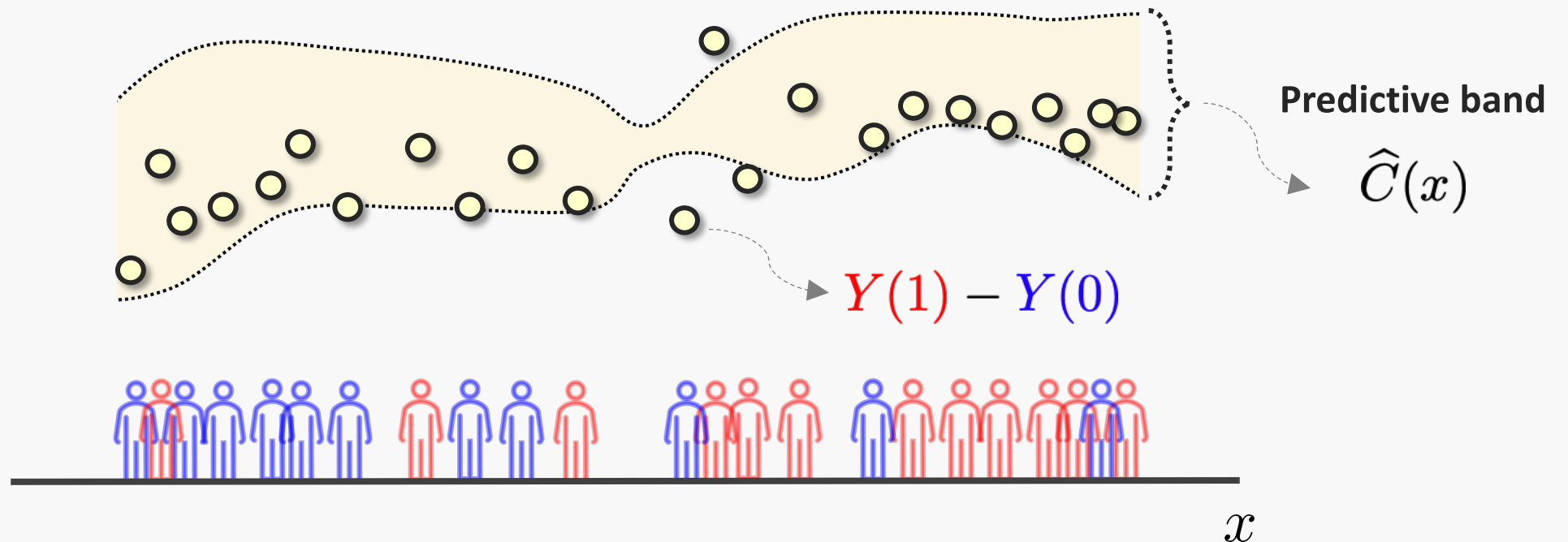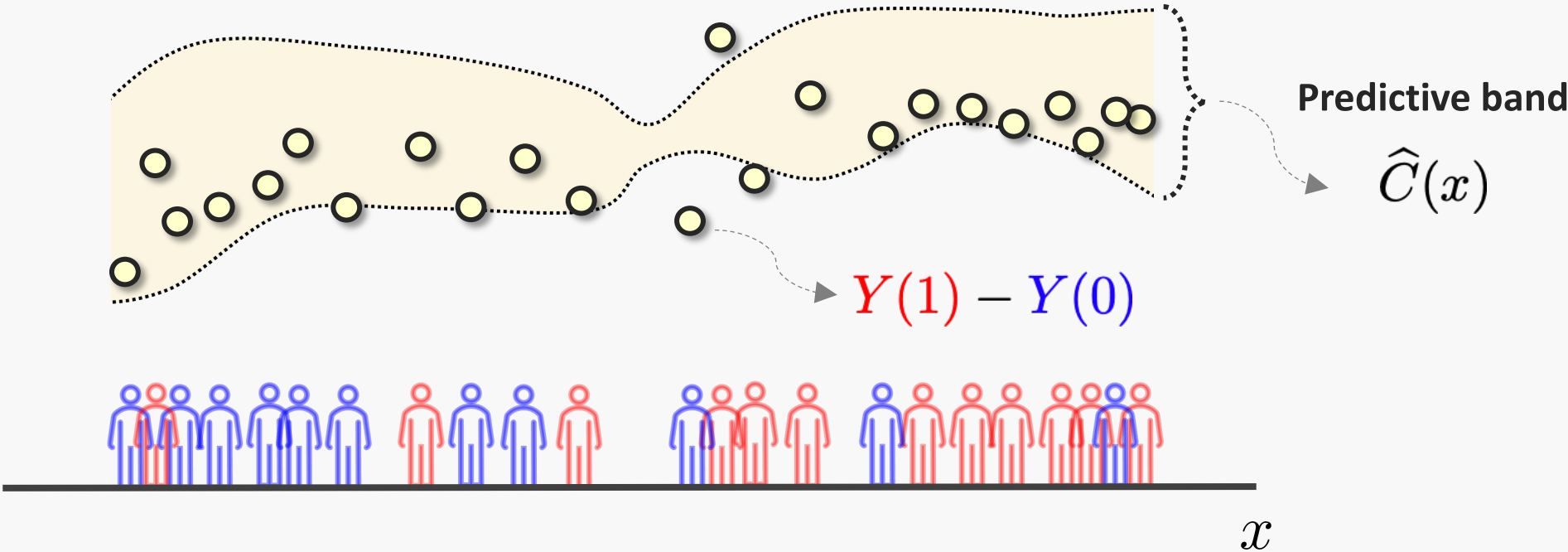● **Predictive Inference:** Construct predictive intervals $\widehat{C}(X_{n+1})$ that cover ITEs

$$P(Y_{n+1}(1) - Y_{n+1}(0) \in \widehat{C}(X_{n+1})) \geq 1 - \alpha$$



$W_i = 1$

$W_i = 0$

Individual treatment effect (ITE)

$Y_i(1) - Y_i(0)$

$x$

# Problem: Valid Predictive Inference on ITEs

- **Predictive Inference:** Construct predictive intervals $\widehat{C}(X_{n+1})$ that cover ITEs

$$P(\textcolor{red}{Y_{n+1}(1)} - \textcolor{blue}{Y_{n+1}(0)} \in \widehat{C}(X_{n+1})) \geq 1 - \alpha$$



**Predictive band**

$\widehat{C}(x)$

$\textcolor{red}{Y(1)} - \textcolor{blue}{Y(0)}$

$x$

# Concept 1: Pseudo-outcome Regression (Meta-learners)

- **Pseudo-outcomes:** Transformations of $(X, W, Y)$ that preserve conditional effects

$$E[\phi(X, W, Y)|X = x] = E[Y(1) - Y(0)|X = x]$$

**Inverse Propensity Weighting**

$$\phi_{\text{IPW}} = \frac{W - \pi(X)}{\pi(X)(1 - \pi(X))} Y$$
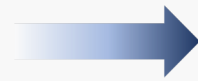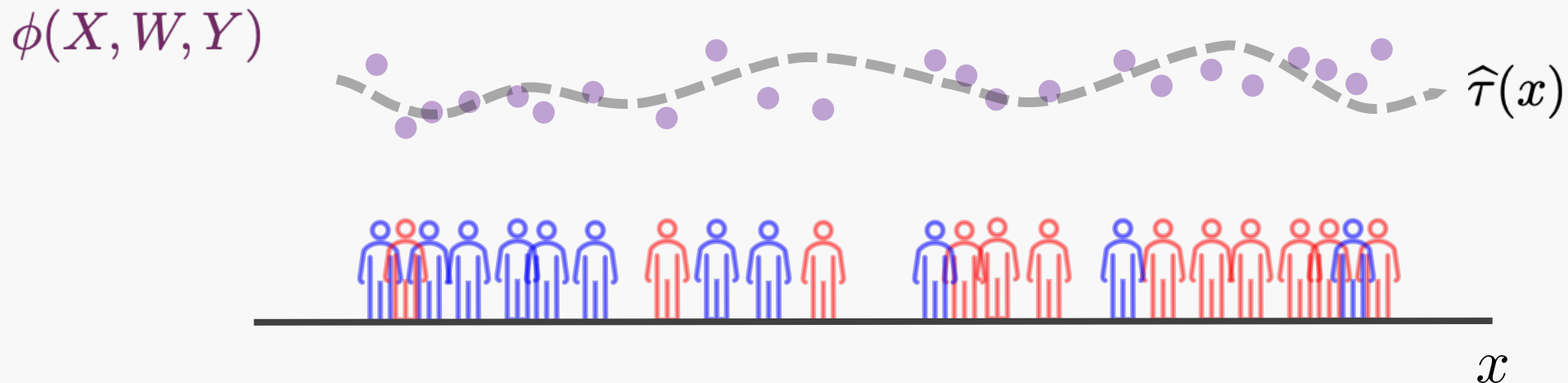
$\phi(X, W, Y)$

$x$

# Concept 1: Pseudo-outcome Regression (Meta-learners)

- **Pseudo-outcomes:** Transformations of $(X, W, Y)$ that preserve conditional effects

$$E[\phi(X, W, Y)|X = x] = E[Y(1) - Y(0)|X = x]$$

$$\tau = E[Y(1) - Y(0)] \quad \Longrightarrow \quad \widehat{\tau} = \frac{1}{n}\sum_i \phi(X_i, W_i, Y_i)$$

$\phi(X, W, Y)$

$x$

# Concept 1: Pseudo-outcome Regression (Meta-learners)

- **Pseudo-outcomes:** Transformations of $(X, W, Y)$ that preserve conditional effects

$$E[\phi(X, W, Y)|X = x] = E[Y(1) - Y(0)|X = x]$$

$$\tau(x) = E[Y(1) - Y(0)|X = x] \longrightarrow \widehat{\tau}(x): \text{ Regress } \phi(X, W, Y) \text{ on } X$$

$\phi(X, W, Y)$

$\widehat{\tau}(x)$

$x$

# Concept 2: Conformal Prediction

- **Conformal Prediction:** A general approach for **post-hoc** predictive inference (V. Vovk 2012)

Finite-sample validity

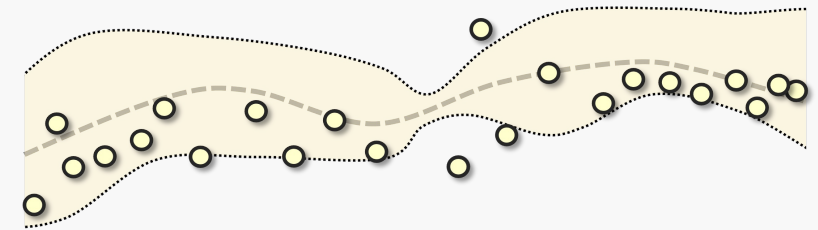Model-free

Distribution-free

# Concept 2: Conformal Prediction

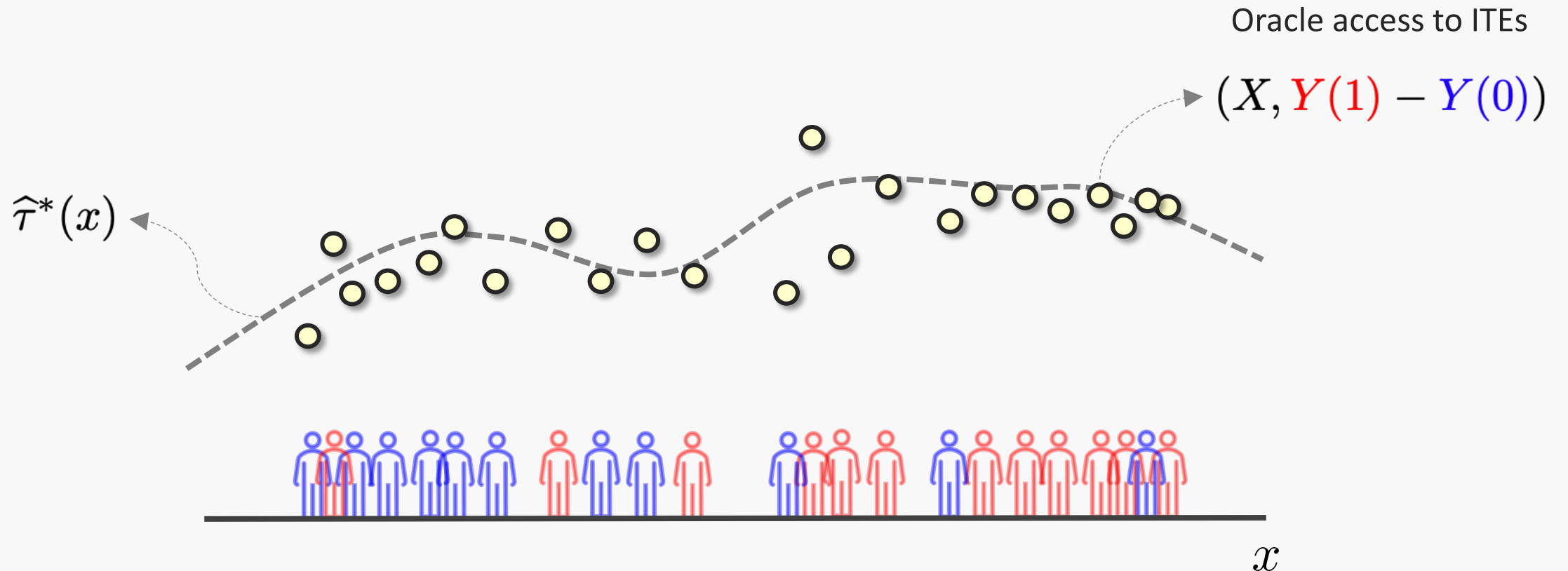- **Conformal Prediction:** A general approach for **post-hoc** predictive inference (V. Vovk 2012)

# Concept 2: Conformal Prediction

- **Step 1:** Train a machine learning model $\widehat{\tau}^*(x)$ using $\{(X_i, Y_i(1) - Y_i(0))\}_i$
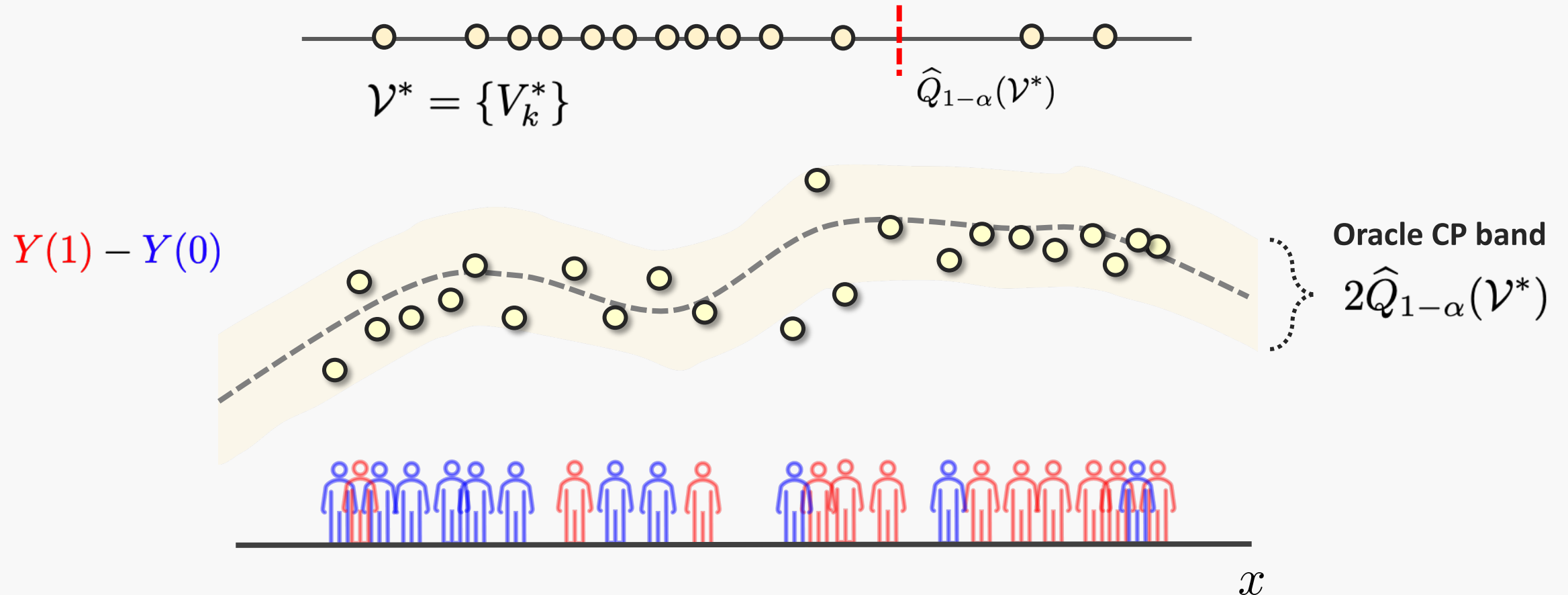


Oracle access to ITEs

$(X, Y(1) - Y(0))$

$\widehat{\tau}^*(x)$

$x$

# Concept 2: Conformal Prediction

- **Step 2:** Evaluate **conformity scores** on a held-out calibration set

$$V_k^*(\widehat{\tau}) = V(\widehat{\tau}(X_k), \textcolor{red}{Y_k(1)} - \textcolor{blue}{Y_k(0)})$$

# Concept 2: Conformal Prediction

- **Step 3:** Construct a predictive interval using the empirical quantile of conformity scores



$$\mathcal{V}^* = \{V_k^*\}$$

$$\widehat{Q}_{1-\alpha}(\mathcal{V}^*)$$

$$Y(1) - Y(0)$$

**Oracle CP band**

$$2\widehat{Q}_{1-\alpha}(\mathcal{V}^*)$$

$$x$$

# Concept 2: Conformal Prediction

- **Step 3:** Construct a predictive interval using the empirical quantile of conformity scores

$$P(Y_{n+1}(1) - Y_{n+1}(0) \in \widehat{C}^*(X_{n+1})) \geq 1 - \alpha$$

If calibration and test data are exchangeable



$Y(1) - Y(0)$

**Oracle CP band**

$2\widehat{Q}_{1-\alpha}(\mathcal{V}_c^*)$

$x$

# Method: Conformal Meta-learners

- **Key Idea:** Apply CP to pseudo-outcomes instead of unobserved ITEs!

$$V_{\varphi,k}(\widehat{\tau}) = V(\widehat{\tau}(X_k), \phi(X_k, W_k, Y_k))$$

# Method: Conformal Meta-learners

- **Key Idea:** Apply CP to pseudo-outcomes instead of unobserved ITEs.

$$P(\phi(X_{n+1}, W_{n+1}, Y_{n+1}) \in \widehat{C}_\varphi(X_{n+1})) \geq 1 - \alpha$$



$\phi(X, W, Y)$

**Pseudo-outcome CP band**

$2\widehat{Q}_{1-\alpha}(\mathcal{V}_\varphi)$

$x$

# Method: Validity of Meta-learners via Stochastic Ordering Theory

- Under what conditions are predictive intervals for pseudo-outcomes valid for ITEs?

Conformity scores evaluated
on pseudo-outcome

Oracle Conformity scores
evaluated on true ITEs

$$V_{\varphi,k}(\widehat{\tau}) = V(\widehat{\tau}(X_k), \phi(X_k, W_k, Y_k))$$

$$V_k^*(\widehat{\tau}) = V(\widehat{\tau}(X_k), Y_k(1) - Y_k(0))$$



$\mathcal{V}_\varphi = \{V_{\varphi,k}\}$

$\widehat{Q}_{1-\alpha}(\mathcal{V}_\varphi)$

$\mathcal{V}_c^* = \{V_k^*\}$

$\widehat{Q}_{1-\alpha}(\mathcal{V}_c^*)$

$\phi(X, W, Y)$

Pseudo-outcome
CP band

$2\widehat{Q}_{1-\alpha}(\mathcal{V}_\varphi)$

$Y(1) - Y(0)$

Oracle CP band

$2\widehat{Q}_{1-\alpha}(\mathcal{V}_c^*)$

$x$

$x$

# Method: Validity of Meta-learners via Stochastic Ordering Theory

- **Sufficient condition for validity:** First-order stochastic dominance!

$$V_\varphi(\widehat{\tau}) \succeq V^*(\widehat{\tau})$$

PDF

$$\mathcal{V}_\varphi = \{\mathcal{V}_{\varphi,k}\}$$

$$\widehat{Q}_{1-\alpha}(\mathcal{V}_\varphi)$$

PDF

$$\mathcal{V}_c^* = \{V_k^*\}$$

$$\widehat{Q}_{1-\alpha}(\mathcal{V}_c^*)$$

CDF

$$Q_\varphi(\alpha) \geq Q^*(\alpha), \ \forall \alpha \in [0,1]$$

$$F_{V^*}$$

$$F_{V_\varphi}$$

$$\alpha$$

$$Q^*(\alpha) \ Q_\varphi(\alpha)$$

Conformity score

# Method: Validity of Meta-learners via Stochastic Ordering Theory

- Unified analysis of validity of meta-learners = stochastic orders of $V_\varphi(\widehat{\tau})$ and $V^*(\widehat{\tau})$

| Meta-learner | Pseudo-outcome |
|---|---|
| X-learner | $\phi_{\mathrm{x}} = W(Y - \widehat{\mu}_0(X)) + (1 - W)(\widehat{\mu}_1(X) - Y)$ |
| IPW-learner<br>Inverse propensity weighted | $\phi_{\mathrm{IPW}} = \frac{W - \pi(X)}{\pi(X)(1 - \pi(X))} Y$ |
| DR-learner<br>Doubly-robust learner | $\phi_{\mathrm{DR}} = \frac{W - \pi(X)}{\pi(X)(1 - \pi(X))}(Y - \widehat{\mu}_W(X)) + (\widehat{\mu}_1(X) - \widehat{\mu}_0(X))$ |

# Method: Validity of Meta-learners via Stochastic Ordering Theory

- Commonly-used meta-learners guarantee model-/distribution-free stochastic orders!

| Meta-learner | Stochastic orders of Conformity Scores |
|---|---|
| X-learner | No distribution-free stochastic order! |
| IPW-learner <br> Inverse propensity weighted | $V^* \succeq_{(2)} V_{\mathrm{IPW}}$ |
| DR-learner <br> Doubly-robust learner | $V^* \succeq_{(2)} V_{\mathrm{DR}}$ |

# Results and Takeaway

- **TL;DR: Conformal meta-learners** = valid predictive inference + accurate point predictions



(a) Empirical assessment of stochastic orders

(b) Coverage, efficiency and RMSE for **Setup A** (top) and **Setup B** (bottom)

(c) Performance at different levels of target coverage

# Poster Session 2

## Tuesday Dec. 12

### 5:15 pm – 7:15 pm