# Causal Data Science

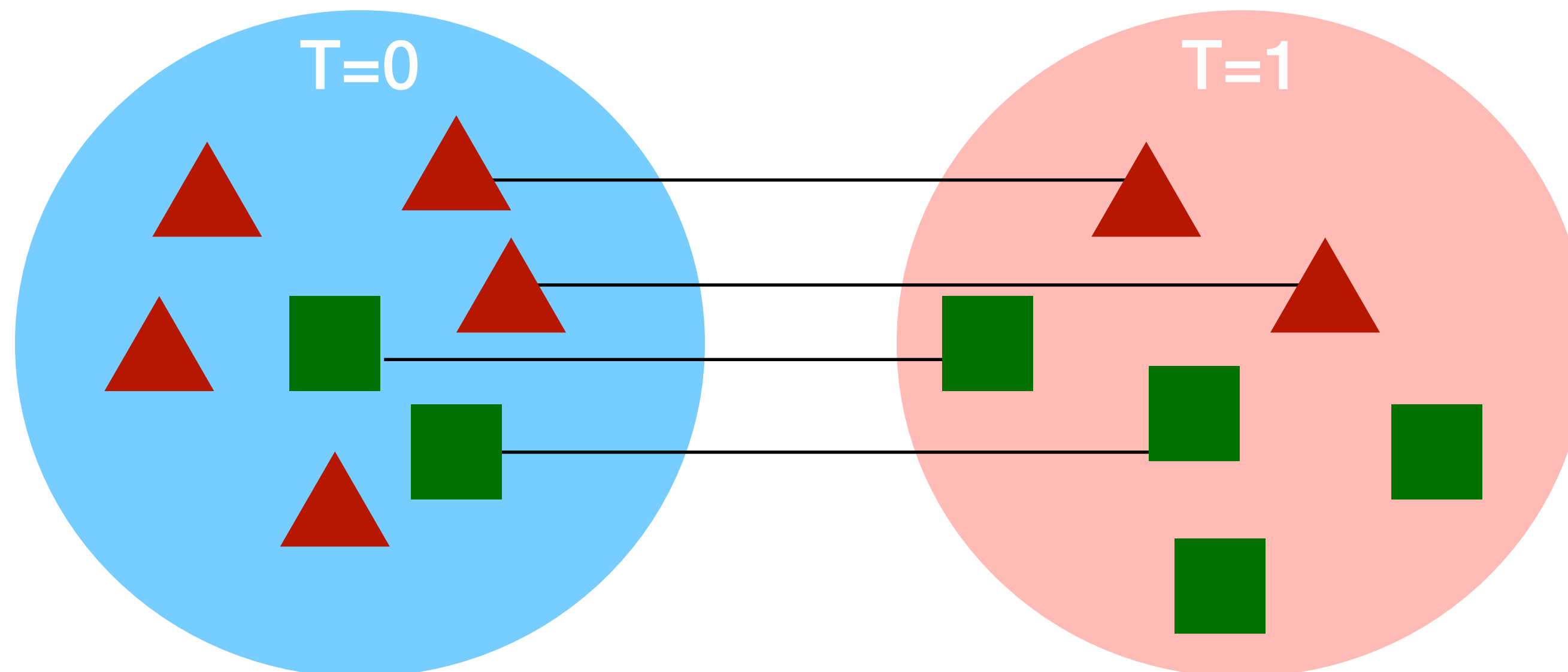**Lecture 8.1: Estimation methods 2**
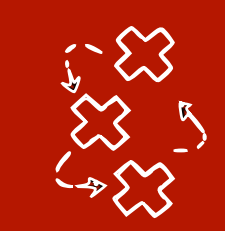
Lecturer: Sara Magliacane

# Last class: Exact matching (simplified)

- **Usually for ATT,** sometimes for ATE

- **Intution:** find the most similar couple of units in terms of covariates $\mathbf{X}$, such that one is in the treatment and the other in the control group
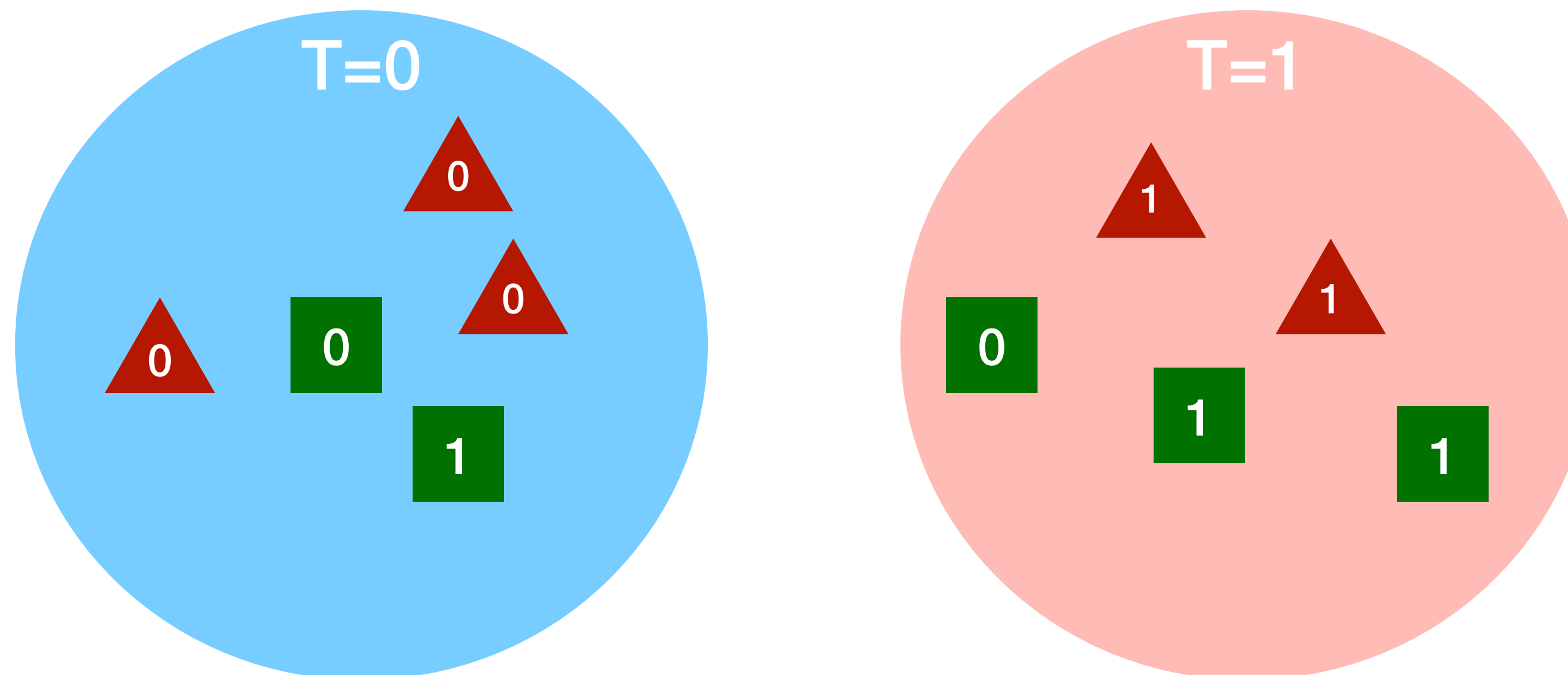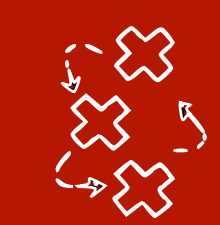
# Advanced: Exact matching (slightly less simplified)

- **Usually for ATT,** with M multiple matches

$$\hat{\text{ATT}} = \frac{1}{n_t} \sum_{i=1}^{n_t} (Y_i - \frac{1}{M} \sum_{j=1}^{M} Y_{m_j(i)})$$

$Y_{m_j(i)}$  **match j for i**

# Advanced: Exact matching (slightly less simplified)

- **Usually for ATT,** with M multiple matches

$$\hat{\text{ATT}} = \frac{1}{n_t} \sum_{i=1}^{n_t} (Y_i - \frac{1}{M} \sum_{j=1}^{M} Y_{m_j(i)})$$

$Y_{m_j(i)}$ **match j for i**



$$M = 1$$

$$\hat{\text{ATT}} = \frac{1}{5} \sum_{i=1}^{5} (Y_i - Y_{m(i)})$$

$$\hat{\text{ATT}} = \frac{1}{5}[1 + 1 + 0 + 0] = \frac{2}{5}$$

# Advanced: Exact matching (slightly less simplified)

- **Usually for ATT,** with M multiple matches
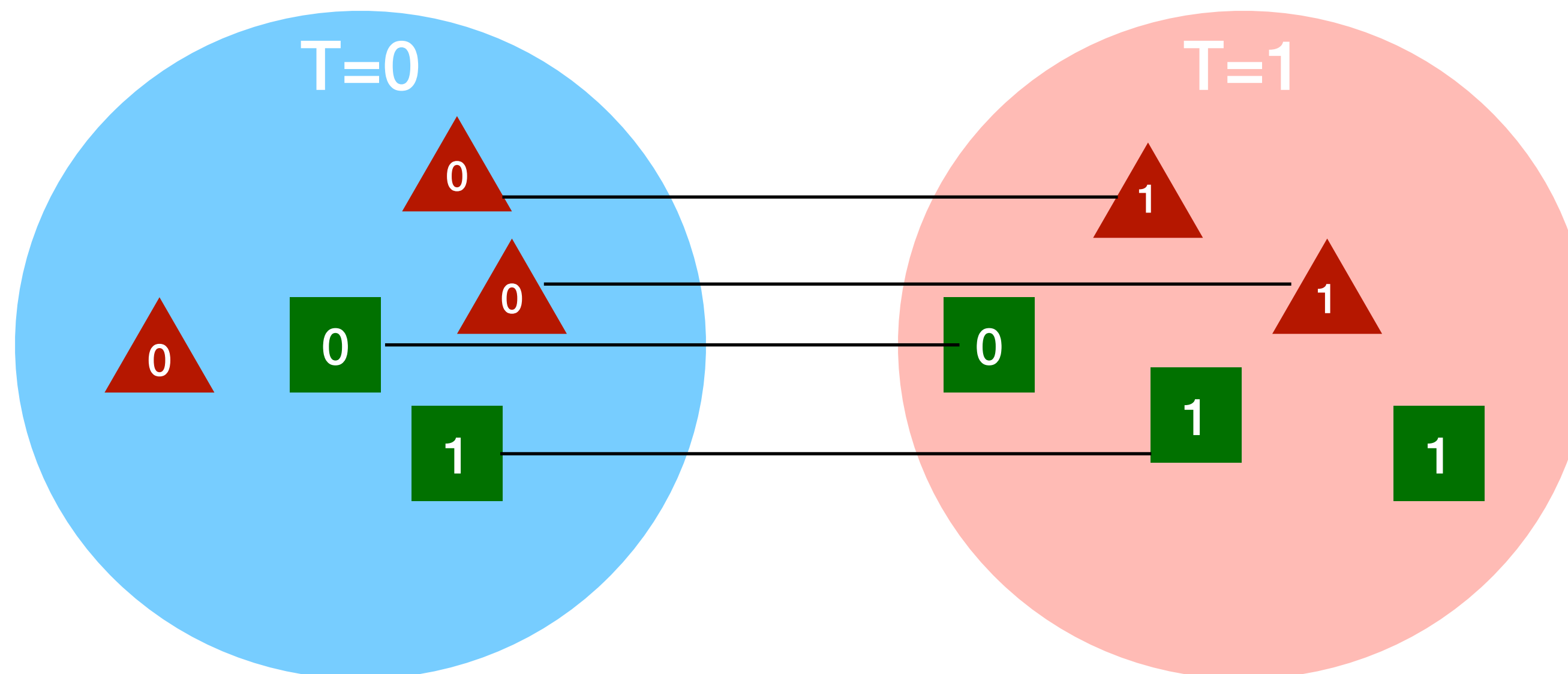
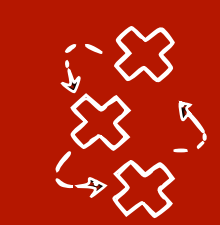$$\hat{ATT} = \frac{1}{n_t} \sum_{i=1}^{n_t} (Y_i - \frac{1}{M} \sum_{j=1}^{M} Y_{m_j(i)})$$

$Y_{m_j(i)}$  **match j for i**

$$M = 2$$



$$\hat{ATT} = \frac{1}{5} \sum_{i=1}^{5} (Y_i - \frac{1}{2} \sum_{j=1}^{M} Y_{m_j(i)})$$

$$\hat{ATT} = \frac{1}{5}[1 + 1 - \frac{1}{2} + \frac{2}{2}] = \frac{1}{5} \cdot \frac{5}{2} = \frac{1}{2}$$

# Advanced: Exact matching (slightly less simplified)

- **ATE** with M multiple matches (e.g. M=2, can be random):

$$\hat{\text{ATE}} = \frac{1}{n_t + n_c} [ \sum_{i=1}^{n_t} (Y_i - \frac{1}{M} \sum_{j=1}^{M} Y_{m_j(i)}) + \sum_{j=1}^{n_c} (\frac{1}{M} \sum_{i=1}^{M} Y_{m_i(j)} - Y_j)]$$
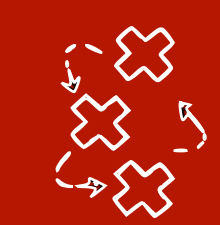


$$\hat{\text{ATE}} = \frac{1}{10} [\frac{5}{2} + 3 + \frac{1}{2} - \frac{1}{2}] = \frac{11}{20}$$

# Last class:  Propensity score matching (PSM)

- **Assumptions**: binary treatment $T$, $\mathbf{X}$ is valid adjustment set

- **Propensity score:** the probability of getting assigned the treatment

$$e(x) \quad \pi(x) := P(T = 1 \mid \mathbf{X} = x)$$

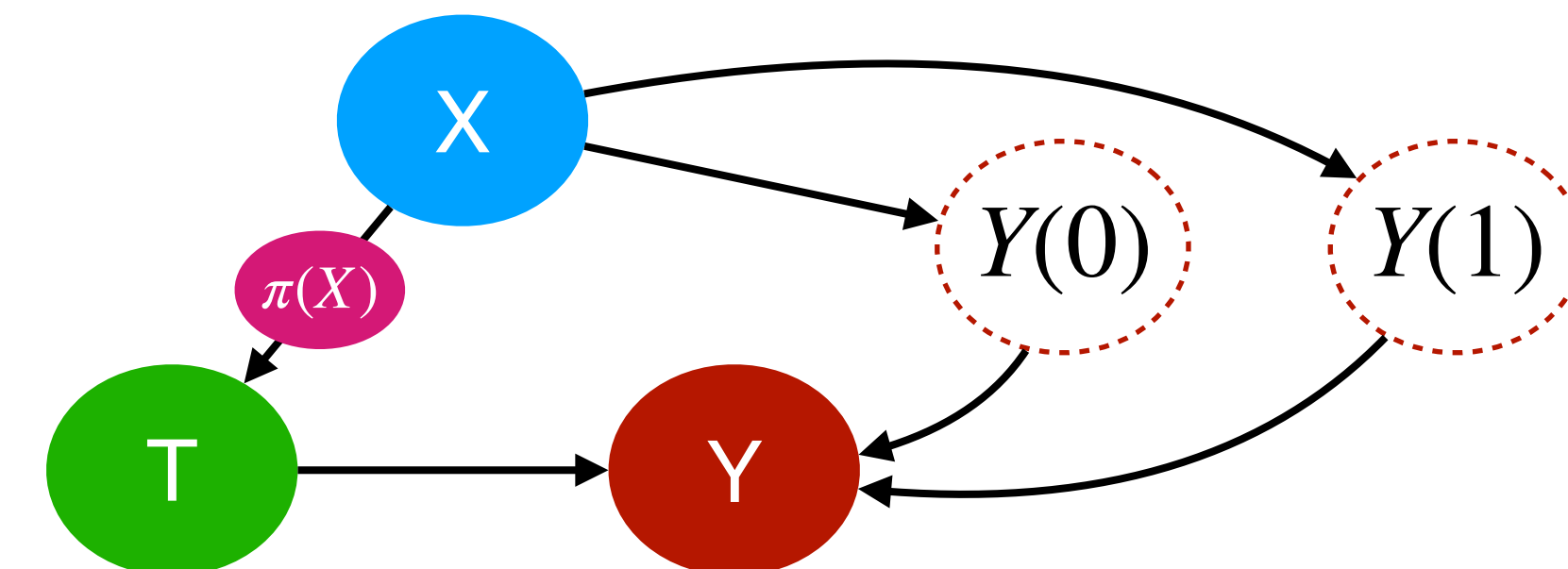**Conditional ignorability**

- We can show that $T \perp\!\!\!\perp \mathbf{X} \mid \pi(\mathbf{X})$ and that if $Y(0), Y(1) \perp\!\!\!\perp T \mid \mathbf{X}$ then

$$Y(0), Y(1) \perp\!\!\!\perp T \mid \pi(\mathbf{X})$$

e.g. with **logistic regression**

- We can estimate $\pi$ from data and use it to match

- If $\mathbf{X}$ has a lot of covariates, it is easier to match since it's a single number

# Estimation method: Inverse probability weighting (IPW)

- **Idea:** rather than match **(and discard some samples)**, reweight (downweight or upweight) samples

- **Inverse probability (of treatment) weighting:** weight by inverse of probability of treatment **received**:
  - For treated $T = 1$: weight by the inverse of $\pi = P(T = 1 | \mathbf{X})$
  - For untreated $T = 0$: weight by the inverse of $1 - \pi = P(T = 0 | \mathbf{X})$

# Inverse probability weighting (IPW) - derivation

- We can estimate the average causal effect/**average treatment effect**

$$\text{ATE} = \mathbb{E}[Y(t=1) - Y(t=0)] = \mathbb{E}[Y|\text{do}(T=1)] - \mathbb{E}[Y|\text{do}(T=0)]$$

- $\mathbf{X}$ is a valid adjustment set for the causal effect of T on Y, so:

$$P(Y=y|\text{do}(T=1)) = \sum_{\mathbf{x}} P(Y=y|\mathbf{X}=\mathbf{x}, T=1)P(\mathbf{X}=\mathbf{x})$$
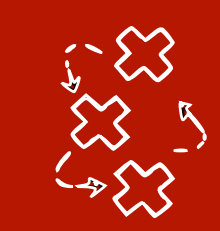
# Inverse probability weighting (IPW) - derivation

- We can estimate the average causal effect/**average treatment effect**

$$\text{ATE} = \mathbb{E}[Y(t = 1) - Y(t = 0)] = \mathbb{E}[Y \,|\, \text{do}(T = 1)] - \mathbb{E}[Y \,|\, \text{do}(T = 0)]$$

- $\mathbf{X}$ is a valid adjustment set for the causal effect of T on Y, so:

$$P(Y = y \,|\, \text{do}(T = t)) = \sum_{\mathbf{x}} P(Y = y \,|\, \mathbf{X} = \mathbf{x}, T = t) P(\mathbf{X} = \mathbf{x})$$

$$\mathbb{E}[Y \,|\, \text{do}(T = t)] = \sum_{y} y \sum_{\mathbf{x}} P(Y = y \,|\, \mathbf{X} = \mathbf{x}, T = t) P(\mathbf{X} = \mathbf{x})$$
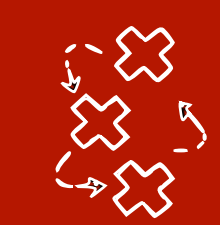
# Inverse probability weighting (IPW) - derivation

$$\mathbb{E}[Y \mid \mathrm{do}(T = t)] = \sum_y \sum_{\mathbf{x}} y \cdot P(Y = y \mid \mathbf{X} = \mathbf{x}, T = t) P(\mathbf{X} = \mathbf{x})$$

Assuming $P(T = t \mid \mathbf{X} = \mathbf{x}) \neq 0$

$$= \sum_y \sum_{\mathbf{x}} y \cdot P(Y = y \mid \mathbf{X} = \mathbf{x}, T = t) P(\mathbf{X} = \mathbf{x}) \frac{P(T = t \mid \mathbf{X} = \mathbf{x})}{P(T = t \mid \mathbf{X} = \mathbf{x})}$$

$$= P(Y = y, \mathbf{X} = \mathbf{x}, T = t)$$
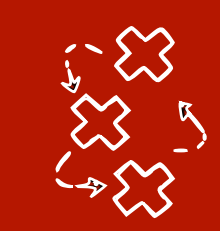
$$= \sum_y \sum_{\mathbf{x}} y \cdot P(Y = y \mid \mathbf{X} = \mathbf{x}, T = t) P(\mathbf{X} = \mathbf{x}) \frac{P(T = t \mid \mathbf{X} = \mathbf{x})}{P(T = t \mid \mathbf{X} = \mathbf{x})}$$

$$= \sum_y \sum_{\mathbf{x}} \frac{y \cdot P(Y = y, \mathbf{X} = \mathbf{x}, T = t)}{P(T = t \mid \mathbf{X} = \mathbf{x})} \qquad \pi \text{ for t} = 1, \ (1 - \pi) \text{ for t} = 0$$

11

# Estimation method: Inverse probability weighting (IPW)

- **Inverse probability (of treatment) weighting:** weight by inverse of probability of treatment **received**:

  - For treated $T = 1$: weight by the inverse of $\pi = P(T = 1 | \mathbf{X})$

  - For untreated $T = 0$: weight by the inverse of $1 - \pi = P(T = 0 | \mathbf{X})$

$$\hat{\mathbb{E}}(Y(t = 1)) = \frac{1}{n} \sum_{i=1}^{n} Y_i \cdot 1\{T = 1\} \cdot \frac{1}{P(T = 1 | X_i)} \qquad \pi$$

$$\hat{\mathbb{E}}(Y(t = 0)) = \frac{1}{n} \sum_{i=1}^{n} Y_i \cdot 1\{T = 0\} \cdot \frac{1}{P(T = 0 | X_i)} \qquad (1 - \pi)$$

# IPW Example

$$X$$
$$\swarrow \qquad \searrow$$
$$T \longrightarrow Y$$

Population:

| | $X=0$ | $X=1$ |
|---|---|---|
| $T=0$ | 1 | 9 |
| $T=1$ | 4 | 1 |

$$P(T=1 \mid X=1) = 0.1$$
$$P(T=1 \mid X=0) = 0.8$$
$$P(T=0 \mid X=1) = 0.9$$
$$P(T=0 \mid X=0) = 0.2$$

T=0

T=1

# IPW Example

$$X$$

$$T \longrightarrow Y$$

Population:

$P(T=1 \mid X=1) = 0.1$

$P(T=1 \mid X=0) = 0.8$

$P(T=0 \mid X=1) = 0.9$

$P(T=0 \mid X=0) = 0.2$

**Reweight by** $\dfrac{1}{P(T_i \mid X_i)}$

T=0

T=1

|       | X=0   | X=1   |
|-------|-------|-------|
| T=0   | 1/0.2 | 9/0.9 |
| T=1   | 4/0.8 | 1/0.1 |

|       | X=0 | X=1 |
|-------|-----|-----|
| T=0   | 5   | 10  |
| T=1   | 5   | 10  |

$T \perp\!\!\!\perp \mathbf{X}$ **in pseudo-population**

# Estimating (conditional) ATEs -S/X learners

- We can estimate the average causal effect/**average treatment effect**

$$\text{ATE} = \mathbb{E}[Y(t=1) - Y(t=0)] = \mathbb{E}_{\text{X}}[\mathbb{E}[Y(t=1)\,|\,\text{X}] - \mathbb{E}[Y(t=0)\,|\,\text{X}]]$$

**Outcome model**        $\hat{\mu}(1,\mathbf{X})$        $\hat{\mu}(0,\mathbf{X})$

**We still assume $\mathbf{X}$ is a valid adjustment set!**

# Estimating (conditional) ATEs -S/X learners

- We can estimate the average causal effect/**average treatment effect**

$$\text{ATE} = \mathbb{E}[Y(t=1) - Y(t=0)] = \mathbb{E}_X[\mathbb{E}[Y(t=1) \mid X] - \mathbb{E}[Y(t=0) \mid X]]$$

$$\hat{\mu}(1, \mathbf{X}) \qquad \hat{\mu}(0, \mathbf{X})$$

$$\hat{\text{ATE}} = \frac{1}{n} \sum_{i=1}^{n} \hat{\mu}(1, \mathbf{x}_i) - \hat{\mu}(0, \mathbf{x}_i)$$

**We still assume $\mathbf{X}$ is a valid adjustment set!**

16

# Estimating (conditional) ATEs -S/X learners

- We can estimate the average causal effect/**average treatment effect**

$$\text{ATE} = \mathbb{E}[Y(t=1) - Y(t=0)] = \mathbb{E}_{\text{X}}[\mathbb{E}[Y(t=1)\,|\,\text{X}] - \mathbb{E}[Y(t=0)\,|\,\text{X}]]$$

$$\hat{\mu}(1,\mathbf{X}) \qquad\qquad \hat{\mu}(0,\mathbf{X})$$

$$\hat{\text{ATE}} = \frac{1}{n}\sum_{i=1}^{n}\hat{\mu}(1,\mathbf{x}_i) - \hat{\mu}(0,\mathbf{x}_i)$$

**We still assume $\mathbf{X}$ is a valid adjustment set!**

- We can also estimate the **conditional average treatment effect:**

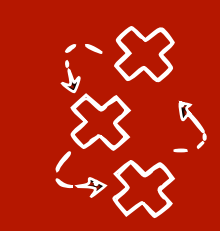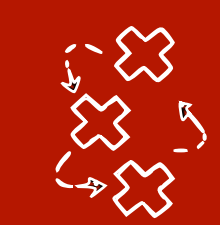$$\text{CATE}(\text{w}) = \mathbb{E}[Y(t=1) - Y(t=0)\,|\,\text{W}=\text{w}]$$

# Estimating (conditional) ATEs -S/X learners

- We can estimate the average causal effect/**average treatment effect**

$$\text{ATE} = \mathbb{E}[Y(t=1) - Y(t=0)] = \mathbb{E}_{\text{X}}[\underbrace{\mathbb{E}[Y(t=1) \mid \text{X}]}_{\hat{\mu}(1,\mathbf{X})} - \underbrace{\mathbb{E}[Y(t=0) \mid \text{X}]}_{\hat{\mu}(0,\mathbf{X})}]$$

$$\hat{\text{ATE}} = \frac{1}{n} \sum_{i=1}^{n} \hat{\mu}(1,\mathbf{x}_i) - \hat{\mu}(0,\mathbf{x}_i)$$

**We assume $\mathbf{X} \cup W$ is a valid adjustment set!**

- We can also estimate the **conditional average treatment effect:**

$$\text{CATE(w)} = \mathbb{E}[Y(t=1) - Y(t=0) \mid \text{W} = \text{w}]$$

$$= \mathbb{E}_{\text{X}}[\underbrace{\mathbb{E}[Y(t=1) \mid \text{X}, \text{W} = \text{w}]}_{\hat{\mu}(1,\mathbf{x}_i, w)} - \underbrace{\mathbb{E}[Y(t=0) \mid \text{X}, \text{W} = \text{w}]}_{\hat{\mu}(0,\mathbf{x}_i, w)}]$$

# S-learners [Küntzel et al 2019]

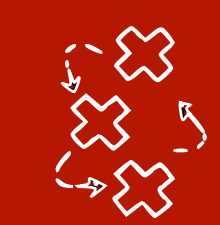- We learn a single model to predict the both potential outcomes $Y_i(0), Y_i(1)$

$$\hat{\text{ATE}} = \frac{1}{n} \sum_{i=1}^{n} \hat{\mu}(1,\mathbf{x}_i) - \hat{\mu}(0,\mathbf{x}_i)$$

$$\hat{\text{CATE}}(w) = \frac{1}{n_w} \sum_{i=1}^{n} 1(W = w)[\hat{\mu}(1,\mathbf{x}_i, w) - \hat{\mu}(0,\mathbf{x}_i, w)]$$

- **Issue:** for high-dimensional $\mathbf{X}$, S-learners can ignore the treatment

# X-learners [Küntzel et al 2019]

1. Learn two separate models $\hat{\mu}_1(\mathbf{x}_i)$ (only treated) and $\hat{\mu}_0(\mathbf{x}_i)$ (only control)

2. We impute the treatment effect per unit (*individual treatment effect*)

**Treatment group**

$Y(0)$ **Control group**

$$\hat{\tau}_{i,1} = Y_i - \hat{\mu}_0(\mathbf{x}_i)$$

$$\hat{\tau}_{i,0} = \hat{\mu}_1(\mathbf{x}_i) - Y_i$$

**Estimated from control**    **Estimated from treated**

| Unit | Y(0) | Y(1) | T | X |
|------|------|------|---|---|
| 1 | ? | 1 | 1 | 1 |
| 2 | 1 | ? | 0 | 1 |
| 3 | ? | 0 | 1 | 0 |
| 4 | 0 | ? | 0 | 0 |
| 5 | ? | 1 | 1 | 1 |
| 6 | ? | 0 | 1 | 0 |

# X-learners [Küntzel et al 2019]

1. Learn two separate models $\hat{\mu}_1(\mathbf{x}_i)$ (only treated) and $\hat{\mu}_0(\mathbf{x}_i)$ (only control)

2. We impute the treatment effect per unit (*individual treatment effect*)

**Treatment group**            $Y(0)$  **Control group**

$$\hat{\tau}_{i,1} = Y_i - \boxed{\hat{\mu}_0(\mathbf{x}_i)}$$        $$\hat{\tau}_{i,0} = \boxed{\hat{\mu}_1(\mathbf{x}_i)} - Y_i$$

**Estimated from control**        **Estimated from treated**

| Unit | Y(0) | Y(1) | T | X |
|------|------|------|---|---|
| 1 | ? | 1 | 1 | 1 |
| 2 | 1 | ? | 0 | 1 |
| 3 | ? | 0 | 1 | 0 |
| 4 | 0 | ? | 0 | 0 |
| 5 | ? | 1 | 1 | 1 |
| 6 | ? | 0 | 1 | 0 |

| Unit | Y(0) | Y(1) | T | X |
|------|------|------|---|---|
| 1 | 1 | 1 | 1 | 1 |
| 2 | 1 | 0 | 0 | 1 |
| 3 | 1 | 0 | 1 | 0 |
| 4 | 0 | 0 | 0 | 0 |
| 5 | 1 | 1 | 1 | 1 |
| 6 | 1 | 0 | 1 | 0 |

# X-learners [Küntzel et al 2019]

1. Learn two separate models $\hat{\mu}_1(\mathbf{x}_i)$ (only treated) and $\hat{\mu}_0(\mathbf{x}_i)$ (only control)

2. We impute the treatment effect per unit (*individual treatment effect)*

$Y(0)$

**Treatment group**        **Control group**

$$\hat{\tau}_{i,1} = Y_i - \boxed{\hat{\mu}_0(\mathbf{x}_i)} \qquad\qquad \hat{\tau}_{i,0} = \boxed{\hat{\mu}_1(\mathbf{x}_i)} - Y_i$$

**Estimated from control**       **Estimated from treated**

3. Learn two separate models $\hat{\tau}_1(\mathbf{x}_i)$ (only treated) and $\hat{\tau}_0(\mathbf{x}_i)$ (only control)

4. The final estimator is a weighted average where $g(\mathbf{x}) : \mathscr{X} \to [0,1]$

$$\hat{\tau}(\mathbf{x}) = g(\mathbf{x}_i)\hat{\tau}_1(\mathbf{x}_i) + (1 - g(\mathbf{x}_i))\hat{\tau}_0(\mathbf{x}_i)$$
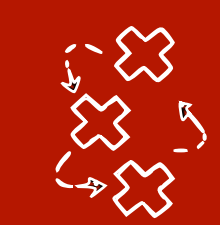
# Issues: Inverse probability weighting (IPW)

- **Inverse probability (of treatment) weighting:** weight by inverse of estimated probability of treatment **received**:

  - For treated $T = 1$: weight by the inverse of $\hat{\pi}(X_i)$

  - For untreated $T = 0$: weight by the inverse of $1 - \hat{\pi}(X_i)$

$$\hat{\mathbb{E}}(Y(t = 1)) = \frac{1}{n} \sum_{i=1}^{n} Y_i \cdot T_i \cdot \frac{1}{\hat{\pi}(X_i)}$$

**We estimate $\hat{\pi}$ e.g. with logistic regression**

**What if the estimated $\hat{\pi}(X_i)$ is biased?**

$$\hat{\mathbb{E}}(Y(t = 0)) = \frac{1}{n} \sum_{i=1}^{n} Y_i \cdot (1 - T_i) \cdot \frac{1}{1 - \hat{\pi}(X_i)}$$

# Advanced: Augmented Inverse probability weighting (AIPW)

- We assume we can estimate in an **unbiased way at least one of the two:**

    1. Propensity scores $\hat{\pi}(\mathbf{x}_i)$

    **-> we say this is a doubly robust method**

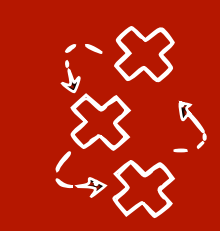    2. S-learner (outcome model) $\hat{\mu}(t_i, \mathbf{x}_i) \approx y_i$

Then: $$\hat{\text{ATE}}_{S-learn} = \frac{1}{n} \sum_{i=1}^{n} \hat{\mu}(1, \mathbf{x}_i) - \hat{\mu}(0, \mathbf{x}_i)$$

**Adjustment on the residuals of the S-learner with IPW**

$$\hat{\text{Adj}}_{S-learm} = \frac{1}{n} \sum_{i=1}^{n} \frac{T_i}{\hat{\pi}(x_i)} (Y_i - \hat{\mu}(1, x_i)) - \frac{1 - T_i}{1 - \hat{\pi}(x_i)} (Y_i - \hat{\mu}(0, x_i))$$

$$\hat{\text{ATE}}_{AIPW} = \hat{\text{ATE}}_{S-learn} + \hat{\text{Adj}}_{S-learn}$$

**This is unbiased if either propensity score or S-learner are unbiased**

# Advanced: Missing data (very briefly)

- **Typical approaches in practice (depending on the assumptions):**
  - Remove all samples with a missing feature (listwise deletion), or
  - Ignore the problem and use the non-missing features of all samples, or
  - Impute (predict) the missing values

# Advanced: Missing data (very briefly)

- **Typical approaches in practice (depending on the assumptions):**
  - Remove all samples with a missing feature (listwise deletion), or
  - Ignore the problem and use the non-missing features of all samples, or
  - Impute (predict) the missing values

- **Typical assumptions (see talk by Karthika Mohan):**
  - $R_X$ is an indicator variable that is 0 if $X$ is missing and 1 otherwise
  - Missing completely at random (MCAR): $R_X \perp\!\!\!\perp \mathbf{X_V}$ (all variables)
  - Missing at random (MAR): $R_X \perp\!\!\!\perp X \,|\, \mathbf{X_V} \backslash \{X\}$
  - Missing not at random (MNAR) - anything else

# Advanced: Missing at random (MAR)

- Missing completely at random (MCAR): coin toss, quite unrealistic

- **Missing at random (MAR):** missing at random given the completely observed (not missing) variables
  - Similar to **ignorability/unconfoundedness**
  - Imputation with EM
  - Multiple imputation (Rubin 1987) - impute m datasets, analyse, combine
    - (Augmented) IPW can be used to analyse/estimate ATE of each dataset

  - See https://scikit-learn.org/stable/modules/impute.html#impute, http://juliejosse.com/wp-content/uploads/2018/07/LectureNotesMissing.html