

Selection Bias in Observational Data: Opt-In vs Randomized Experiments

Justin S. Eloriaga

Emory University

ECON 521

- A fitness app launches **streak reminders** to increase weekly workouts.
- Two rollout strategies:
 1. **Observational (opt-in)**: users choose to enable reminders.
 2. **Randomized (RCT)**: app randomly assigns reminders to 50% of users.
- Claim: If highly motivated users are more likely to opt in, the observational comparison *overstates* the causal effect.

Potential Outcomes Notation

For user i :

- $Y(1)_i$: weekly workouts *if* treated (has reminders).
- $Y(0)_i$: weekly workouts *if not* treated.
- $\Delta_i \equiv Y(1)_i - Y(0)_i$: individual causal effect.
- $D_i \in \{0, 1\}$: opt-in indicator (1 = chose reminders).
- $T_i \in \{0, 1\}$: randomized assignment (1 = assigned reminders).
- v_i : **unobserved motivation** (higher $v_i \Rightarrow$ more likely to opt in and to work out more).

Observed outcome under each regime:

$$Y_i^{\text{obs}} = \begin{cases} D_i Y(1)_i + (1 - D_i) Y(0)_i & \text{(observational)} \\ T_i Y(1)_i + (1 - T_i) Y(0)_i & \text{(randomized)} \end{cases}$$

Selection-Distorted Observational Difference (SDO):

$$\text{SDO} \equiv \mathbb{E}[Y^{\text{obs}} \mid D = 1] - \mathbb{E}[Y^{\text{obs}} \mid D = 0] = \mathbb{E}[Y(1) \mid D = 1] - \mathbb{E}[Y(0) \mid D = 0].$$

Selection-Distorted Observational Difference (SDO):

$$\text{SDO} \equiv \mathbb{E}[Y^{\text{obs}} \mid D = 1] - \mathbb{E}[Y^{\text{obs}} \mid D = 0] = \mathbb{E}[Y(1) \mid D = 1] - \mathbb{E}[Y(0) \mid D = 0].$$

Average Treatment Effect (ATE, via RCT):

$$\text{ATE} \equiv \mathbb{E}[Y^{\text{obs}} \mid T = 1] - \mathbb{E}[Y^{\text{obs}} \mid T = 0] = \mathbb{E}[Y(1) - Y(0)] \quad \text{since } T \perp \{Y(0), Y(1)\}.$$

Why SDO is Biased (Decomposition)

$$\begin{aligned}\text{SDO} &= \mathbb{E}[Y(1) \mid D = 1] - \mathbb{E}[Y(0) \mid D = 0] \\ &= \underbrace{\mathbb{E}[Y(1) - Y(0) \mid D = 1]}_{\text{effect among opt-ins}} + \underbrace{\left(\mathbb{E}[Y(0) \mid D = 1] - \mathbb{E}[Y(0) \mid D = 0] \right)}_{\text{selection term}}.\end{aligned}$$

- If motivated users (v high) both opt in ($D = 1$) and have higher baselines $Y(0)$, then the **selection term** > 0 .
- Result: $\text{SDO} > \text{ATE}$ on average (upward bias).

- Heterogeneous effects: $\Delta_i \sim \mathcal{N}(0.5, 0.2^2)$.
- Draw (e_{0i}, e_{1i}, v_i) from a trivariate normal with positive correlations.
- Potential outcomes: $Y(0)_i = e_{0i}$, $Y(1)_i = e_{1i} + \Delta_i$.
- **Opt-in rule:** $D_i = \mathbf{1}\{v_i > 0\}$.
- **RCT rule:** $T_i \sim \text{Bernoulli}(0.5)$.

Python: Imports & Seed

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 np.random.seed(7)
```


Python: One Large Run + Bias Decomposition

```
1 mu = np.array([0.0, 0.0, 0.0])
2 Sigma = np.array([[1.0, 0.6, 0.4],
3                   [0.6, 1.0, 0.4],
4                   [0.4, 0.4, 1.0]])
5 n = 100_000
6
7 delta = np.random.normal(loc=0.5, scale=0.2, size=n)
8 E = np.random.multivariate_normal(mu, Sigma, size=n)
9 e0, e1, v = E[:,0], E[:,1], E[:,2]
10
11 y0 = e0
12 y1 = e1 + delta
13
14 D = (v > 0).astype(int)
15 T = np.random.binomial(1, 0.5, size=n)
```

Python: One Large Run + Bias Decomposition

```
1 y_obs_obs = D*y1 + (1-D)*y0
2 y_obs_rct = T*y1 + (1-T)*y0
3
4 SDO = y_obs_obs[D==1].mean() - y_obs_obs[D==0].mean()
5 ATE = y_obs_rct[T==1].mean() - y_obs_rct[T==0].mean()
6 true_ATE = delta.mean()
7
8 effect_optins = (y1 - y0)[D==1].mean()
9 selection_term = y0[D==1].mean() - y0[D==0].mean()
10
11 print(f"SDO: {SDO:.3f}, ATE: {ATE:.3f}, true : {true_ATE:.3f}")
12 print("Decomp:",
13       f"E[Y(1)-Y(0) | D=1]={effect_optins:.3f}",
14       f"Sel={selection_term:.3f}")
```

Python: Repeat Experiments (10k Sims)

```
1 nsample = 1000
2 nsim = 10_000
3 sdo_vals, ate_vals, true_deltas = [], [], []
4
5 for _ in range(nsim):
6     delta = np.random.normal(0.5, 0.2, nsample)
7     E = np.random.multivariate_normal(mu, Sigma, size=nsample)
8     e0, e1, v = E[:,0], E[:,1], E[:,2]
9     y0 = e0
10    y1 = e1 + delta
11
12    D = (v > 0).astype(int)
13    T = np.random.binomial(1, 0.5, size=nsample)
14
15    y_obs_obs = D*y1 + (1-D)*y0
16    y_obs_rct = T*y1 + (1-T)*y0
```

Python: Repeat Experiments (10k Sims)

```
1
2     sdo_vals.append(y_obs_obs[D==1].mean() - y_obs_obs[D==0].mean())
3     ate_vals.append(y_obs_rct[T==1].mean() - y_obs_rct[T==0].mean())
4     true_deltas.append(delta.mean())
5
6 sdo_vals = np.array(sdo_vals)
7 ate_vals = np.array(ate_vals)
8 true_deltas = np.array(true_deltas)
9
10 print(f"Mean SDO: {sdo_vals.mean():.3f}")
11 print(f"Mean ATE: {ate_vals.mean():.3f}")
12 print(f"Mean true : {true_deltas.mean():.3f}")
```

Python: Plot Distributions

```
1 # SDO histogram
2 fig, ax = plt.subplots(figsize=(6,4))
3 ax.hist(sdo_vals, bins=50)
4 ax.axvline(true_deltas.mean(), linestyle='--')
5 ax.set_title("Distribution of SDO (observational opt-in)")
6 ax.set_xlabel("Estimated effect"); ax.set_ylabel("Frequency")
7 plt.show()
8
9 # ATE histogram
10 fig, ax = plt.subplots(figsize=(6,4))
11 ax.hist(ate_vals, bins=50)
12 ax.axvline(true_deltas.mean(), linestyle='--')
13 ax.set_title("Distribution of ATE (randomized)")
14 ax.set_xlabel("Estimated effect"); ax.set_ylabel("Frequency")
15 plt.show()
```

Figure Placeholders

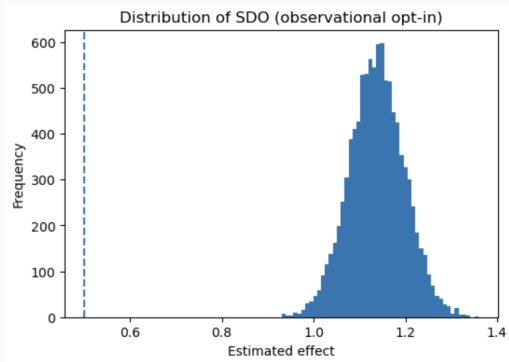


Figure 1: Paste SDO histogram. Dashed line = mean true effect.

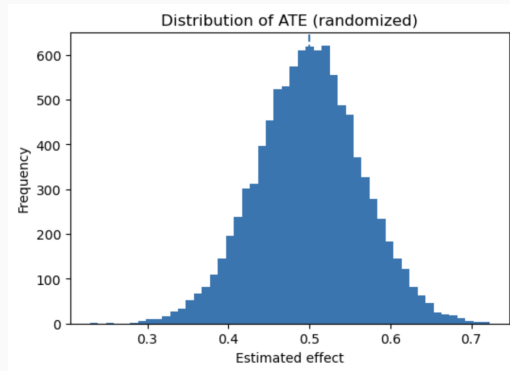


Figure 2: Paste ATE histogram. Dashed line = mean true effect.

- **Observational (opt-in):** D correlates with unobserved v , which also raises $Y(0)$.
- $\Rightarrow \mathbb{E}[Y(0) \mid D = 1] > \mathbb{E}[Y(0) \mid D = 0] \Rightarrow$ positive **selection term**.
- **Randomization:** $T \perp \{Y(0), Y(1)\} \Rightarrow$ treated vs control difference recovers $\mathbb{E}[Y(1) - Y(0)]$.

- $Y(1)$ and $Y(0)$ are the **counterfactual** outcomes; you only ever observe one.
- **SDO** mixes treatment effect with selection on unobservables.
- **RCT ATE** centers on the true causal effect under random assignment.
- If RCT is infeasible, consider strong quasi-experimental designs (IV, RD, DID with good parallel-trends, matching with rich covariates, etc.).