On Counterfactual Metrics for Social Welfare: Incentives, Ranking, and Information Asymmetry

Serena Wang Joint with Stephen Bates, P. M. Aronow, Michael I. Jordan CIRM, December 14, 2023





Anthony C	Golden Gate Sotheby's	59 \$2	2,059,550.00	
Andrea G	COMPASS	54 \$1 \$1	1,026,346.50 - 1,710,577.50	
University of (California Berke	ley	Subject Score 82.3	
United States Berkel	ey	•	Global Score	
#4 in Best Global Universit	ersities	e	88.7	

BROKERAGE

Real Estate Experts

AGENT NAME

Brett J

and services

10 Top Real Estate Agents in Berkeley by Sales Transactions

SALES IN PRIOR 12

MONTHS*

562

PRICE OF SALES OVER

24 MONTHS* \$1,160,907.75 -

\$1,934,846.25

\$1 235 730 00 -

Enrollment

40,921

The University of California–Berkeley is situated roughly 15 miles from San Francisco in what is known as the Bay Area... Read More $\ensuremath{\text{s}}$

t Charity Navigator	Q k	1	10 Тор	Real Estate Agents in	Berkeley by Sales	Transactions
List of Best Highly Rate	d Charities		AGENT NAME	BROKERAGE	SALES IN PRIOR 12 MONTHS*	PRICE OF SALES OVER 24 MONTHS*
Best Hospital	s Honor Ro)	J	Real Estate Experts	562	\$1,160,907.75 - \$1,934,846.25
Each year, U.S. News ranks h are the top 20 highest rated	nospitals in 15 special hospitals in the U.S.	ties and 20 procedures and conditions. H	bere	Golden Gate Sotheby's	\$ 59	\$1,235,730.00 - \$2 059 550 00
2022-2023 Honor Ro Rankings	M		ea (Hospital Compare	Learning Center	About Of My Hospitals ~
#1 Mayo Clinic Rochester, MN		BEST	•		۹	· · · · · · · · · · · · · · · · · · ·
#2 Cedars-Sinai Medical Center Los Angeles, CA		HOSPITALS	9	Ratings	You Can Trus	st.
#3 NYU Langone Hospitals New York, NY		HONOR ROLL		Cal Hospital Con compare the qu	npare makes it easy to find and ality of hospitals in California.	WESTLAK
SEE FULL RANKINGS LIST »	-	11	Comp	pare Hospitals Search by	Zip Code, City, or Hospital Name	Q C
WISCONSIN Milwaukee , WI		Restant and all		C-Section Honor	Roll Hospitals Announ	ced aral target aimed at
Super-sized Free goods distr	ibution 100		Unive ¹	reducing Cesarean births (C-section	s) for first-time mothers with low-	-risk pregnancies.
	100%		#19 in Best Univ	ersities for Artificial Intelligen	ce	Giobal Score 88.7
			The University of Califo what is known as the Ba	rnia—Berkeley is situated roughly 1 ay Area Read More »	5 miles from San Francisco in	Enrollment 40,921

Metrics often don't match goals

Goodhart (1975) *(Economist)*: "Any observed statistical regularity will tend to collapse once pressure is placed upon it for control purposes."





Campbell (1979) *(Social Psychologist)*: "The more any quantitative social indicator is used for social decision-making, [...] the more apt it will be to distort and corrupt the social processes it is intended to monitor."

Strathern (1997) *(Anthropologist)*: "When a measure becomes a target, it ceases to be a good measure."



Domain

Monetary policy and finance

Social programs

Political processes

Education

Business and management

Behind every algorithm is a metric.

1990: New York



1990: New York



"More severely ill patients experienced dramatically worsened health outcomes" (Dranove et al. 2003).

Dranove, D., Kessler, D., McClellan, M., & Satterthwaite, M. (2003). Is more information better? The effects of "report cards" on health care providers. Journal of Political Economy, 111(3), 555-588.

1990: New York



"More severely ill patients experienced dramatically worsened health outcomes" (Dranove et al. 2003).

Dranove, D., Kessler, D., McClellan, M., & Satterthwaite, M. (2003). Is more information better? The effects of "report cards" on health care providers. Journal of Political Economy, 111(3), 555-588.









Centers for Medicare and

Medicaid Services (CMS)

Patients

Best Hospitals Honor Roll

Each year, U.S. News ranks hospitals in 15 specialties and 20 procedures and conditions. Here are the top 20 highest rated hospitals in the U.S.

Hospital





University of California Berkeley	Subject Score 82.3
#19 in Best Universities for Artificial Intelligence #4 in Best Global Universities	Global Score 88.7
The University of California—Berkeley is situated roughly 15 miles from San Francisco in what is known as the Bay Area Read More »	Enrollment 40,921





Real estate agents

[%] ≥ Zillow

Sales volume

Zillow, Yelp, Google, US News and World, etc.

Home buyers/sellers

AGENT NAME	BROKERAGE	SALES IN PRIOR 12 MONTHS*	PRICE OF SALES OVER 24 MONTHS*
Brett J	Real Estate Experts	562	\$1,160,907.75 -
			\$1,934,846.25
Anthony C	Golden Gate Sotheby's	59	\$1,235,730.00 -
			\$2,059,550.00
Andrea G	COMPASS	54	\$1,026,346.50 -
			\$1,710,577.50

Key Questions

Question 1 (Incentives): How do we design quality metrics that lead to better treatment incentives?



Question 1 (Incentives): How do we design quality metrics that lead to better treatment incentives?

Question 2 (Ranking): Do rankings based on those quality metrics behave reasonably?

Key Questions

Question 1 (Incentives): How do we design quality metrics that lead to better treatment incentives?

Question 2 (Ranking): Do rankings based on those quality metrics behave reasonably?

Question 3 (Information Asymmetry): What if the hospitals know more about patients than the agencies?

Related work

"Cream skimming" [Culyer & Newhouse, 2000; Koning & Heinrich, 2013, Glazer & McGuire, 2000]

> Empirical; theoretical

Counterfactual metrics for better incentives in modern ranking systems driven by data/ML.

Culyer, A. J., & Newhouse, J. P. (Eds.). (2000). Handbook of health economics. Elsevier.

Koning, P., & Heinrich, C. J. (2013). Cream-skimming, parking and other intended and unintended effects of high-powered, performance-based contracts. *Journal of Policy Analysis and Management*, 32(3), 461-483. Glazer, J. & McGuire, T. G. (2000). "Optimal risk adjustment in markets with adverse selection: an application to managed care." *American Economic Review* 90(4): 1055-71.

Related work

"Cream skimming" [Culyer & Newhouse, 2000; Koning & Heinrich, 2013, Glazer & McGuire, 2000]



Counterfactual metrics for better incentives in modern ranking systems driven by data/ML.

Culyer, A. J., & Newhouse, J. P. (Eds.). (2000). Handbook of health economics. Elsevier.

Koning, P., & Heinrich, C. J. (2013). Cream-skimming, parking and other intended and unintended effects of high-powered, performance-based contracts. *Journal of Policy Analysis and Management*, 32(3), 461-483. Glazer, J. & McGuire, T. G. (2000). "Optimal risk adjustment in markets with adverse selection: an application to managed care." *American Economic Review* 90(4): 1055-71.

Manski, C. F. (2009). Identification for prediction and decision. Harvard University Press.

Athey, S., & Wager, S. (2021). Policy learning with observational data. Econometrica, 89(1), 133-161.

Kitagawa, T., & Tetenov, A. (2018). Who should be treated? Empirical welfare maximization methods for treatment choice. Econometrica, 86(2), 591-616.

Related work

"Cream skimming" [Culyer & Newhouse, 2000; Koning & Heinrich, 2013, Glazer & McGuire, 2000]



ranking systems driven by data/ML.

Culyer, A. J., & Newhouse, J. P. (Eds.). (2000). Handbook of health economics. Elsevier.

Koning, P., & Heinrich, C. J. (2013). Cream-skimming, parking and other intended and unintended effects of high-powered, performance-based contracts. *Journal of Policy Analysis and Management*, 32(3), 461-483. Glazer, J. & McGuire, T. G. (2000). "Optimal risk adjustment in markets with adverse selection: an application to managed care." *American Economic Review* 90(4): 1055-71.

Manski, C. F. (2009). *Identification for prediction and decision*. Harvard University Press.

Athey, S., & Wager, S. (2021). Policy learning with observational data. Econometrica, 89(1), 133-161.

Kitagawa, T., & Tetenov, A. (2018). Who should be treated? Empirical welfare maximization methods for treatment choice. Econometrica, 86(2), 591-616.

Laffont, J. J., & Martimort, D. (2009). The Theory of Incentives. Princeton University Press.

Notation: potential outcomes (Neyman–Rubin)

Patient covariates: $X_i \in \mathbb{R}^d$

Treatment assignment: $T_i \in \{0, 1\}$

Potential outcomes if untreated or treated: $Y_i(0), Y_i(1)$ Observed outcome: $Y_i = Y_i(T_i)$

Notation: potential outcomes (Neyman–Rubin)

Patient covariates: $X_i \in \mathbb{R}^d$

Treatment assignment: $T_i \in \{0,1\}$

Potential outcomes if untreated or treated: $Y_i(0), Y_i(1)$ Observed outcome: $Y_i = Y_i(T_i)$

Conditional average potential outcomes: $\mu_1(x) = E[Y_i(1)|X_i = x]$ $\mu_0(x) = E[Y_i(0)|X_i = x]$

Principal-agent model



Principal (Health Dept)



Principal-agent model

Decentralized information: Agent knows more about task and effort than principal.



Principal (Health Dept)

Principal-agent model

Decentralized information: Agent knows more about task and effort than principal.

Decentralized interests: Principal sets contract that determines agent's reward.



Question 1 (Incentives): How do we design quality metrics that lead to better hospital incentives?



Principal (Health Dept)

Principal (Health Dept)

Observes

Covariates $\mathbf{X} = X_1, ..., X_n$ Treatments $\mathbf{T} = T_1, ..., T_n$ Outcomes $\mathbf{Y} = Y_1, ..., Y_n$ (Possibly missing if $T_i = 0$)

Principal (Health Dept)

Observes

Covariates
$$\mathbf{X} = X_1, ..., X_n$$

Treatments $\mathbf{T} = T_1, ..., T_n$
Outcomes $\mathbf{Y} = Y_1, ..., Y_n$
(Possibly missing if $T_i = 0$)

Chooses

Reward function $w(\mathbf{Y}, \mathbf{T}, \mathbf{X})$

Principal (Health Dept)

Observes

Covariates
$$\mathbf{X} = X_1, ..., X_n$$

Treatments $\mathbf{T} = T_1, ..., T_n$
Outcomes $\mathbf{Y} = Y_1, ..., Y_n$
(Possibly missing if $T_i = 0$)

Chooses

Reward function $w(\mathbf{Y},\mathbf{T},\mathbf{X})$

Agent (Hospital)

Observes

Covariates
$$\mathbf{X} = X_1, ..., X_n$$

Principal (Health Dept)

Observes

Covariates
$$\mathbf{X} = X_1, ..., X_n$$

Treatments $\mathbf{T} = T_1, ..., T_n$

Outcomes
$$\mathbf{Y} = Y_1, ..., Y_n$$

(Possibly missing if $T_i = 0$)

Chooses

Reward function $w(\mathbf{Y},\mathbf{T},\mathbf{X})$

Agent (Hospital) Observes Covariates $\mathbf{X} = X_1, ..., X_n$ Knows Conditional average potential outcomes $\mu_0(X_i), \mu_1(X_i)$

Principal (Health Dept)

Observes

Covariates
$$\mathbf{X} = X_1, ..., X_n$$

Treatments $\mathbf{T} = T_1, ..., T_n$

Outcomes $\mathbf{Y} = Y_1, ..., Y_n$ (Possibly missing if $T_i = 0$)

Chooses

Reward function $w(\mathbf{Y},\mathbf{T},\mathbf{X})$

Agent (Hospital) Observes Covariates $\mathbf{X} = X_1, ..., X_n$ Knows Conditional average potential outcomes $\mu_0(X_i), \mu_1(X_i)$

Chooses

Treatment rule

$$\pi(x) = P(T_i = 1 | X_i = x)$$

Game

1. Principal (Health Dept) chooses reward function w



Principal (Health Dept)

Game

- 1. Principal (Health Dept) chooses reward function w
- 2. Agent (Hospital) best responds with treatment rule π^w



Defining welfare and regret

Principal's goal: design a reward to maximize welfare effect (Manski, 2009):

$$V(\pi) = E[Y_i(T_i^{\pi}) - Y_i(0)]$$

Outcomes under treatment rule

Outcomes if no one gets treated

Defining welfare and regret

Principal's goal: design a reward to maximize welfare effect (Manski, 2009):

$$V(\pi) = E[Y_i(T_i^{\pi}) - Y_i(0)]$$

Outcomes under treatment rule

Outcomes if no one gets treated

Regret:

$$R(\pi^{w}) = \max_{\pi} V(\pi) - \underbrace{V(\pi^{w})}_{\text{Welfare effect of policy}} - \underbrace{V(\pi^{w})}_{\text{Welfare effect of policy}}$$
Welfare effect of policy optimized under reward function w .

Incentive problems with the status quo

Reward function 1: Status quo mortality rate, "Average Treated Outcome"

$$w_{\text{ATO}}(\mathbf{X}, \mathbf{T}^{\pi}, \mathbf{Y}) = \frac{\sum_{i=1}^{n} Y_i T_i^{\pi}}{\sum_{i=1}^{n} T_i^{\pi}}$$

Incentive problems with the status quo

Reward function 1: Status quo mortality rate, "Average Treated Outcome"

$$w_{\text{ATO}}(\mathbf{X}, \mathbf{T}^{\pi}, \mathbf{Y}) = \frac{\sum_{i=1}^{n} Y_i T_i^{\pi}}{\sum_{i=1}^{n} T_i^{\pi}}$$

Agent's best response: Treat only the patients with the highest expected treated outcome.

$$\pi^{w_{\text{ATO}}}(x) = \mathbb{1}(x \in \arg \max_{x} \mu_{1}(x) \text{ and } \mu_{1}(x) > 0)$$

Proposition 1: Regret is unbounded!

Change 1: reward benefit from treatment

Reward function 2: "Average Treatment Effect on the Treated" (ATT)

$$w_{\text{ATT}}(\mathbf{Y}, \mathbf{T}^{\pi}, \mathbf{X}) = \frac{\sum_{i=1}^{n} (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi}}{\sum_{i=1}^{n} T_i^{\pi}}$$

Estimated untreated counterfactual
Change 2: reward total effect

Reward function 3 (proposed): Aligned with welfare, "Total Treatment Effect"

$$w_{\mathrm{TT}}(\mathbf{Y}, \mathbf{T}^{\pi}, \mathbf{X}) = \sum_{i=1}^{n} (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi}$$
$$\boxed{\sum_{i=1}^{n} T_i^{\pi}}$$

Change 2: No denominator ("total effect" instead of "average effect")

Change 2: reward total effect

Reward function 3 (proposed): Aligned with welfare, "Total Treatment Effect"

$$w_{\mathrm{TT}}(\mathbf{Y}, \mathbf{T}^{\pi}, \mathbf{X}) = \sum_{i=1}^{n} (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi}$$

Agent's best response: Treat all patients with a positive treatment effect.

$$\pi^{w_{\mathrm{TT}}}(x) = \mathbb{1}(\tau(x) > 0)$$

Proposition 3: Regret is zero as long as $\hat{\mu}_0(x)$ is unbiased

Change 2: reward total effect

Reward function 3 (proposed): Aligned with welfare, "Total Treatment Effect"

$$w_{\mathrm{TT}}(\mathbf{Y}, \mathbf{T}^{\pi}, \mathbf{X}) = \sum_{i=1}^{n} (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi}$$

Agent's best response: Treat all patients with a positive treatment effect.

$$\pi^{w_{\mathrm{TT}}}(x) = \mathbb{1}(\tau(x) > 0)$$

Proposition 3: Regret is **zero as long as** $\hat{\mu}_0(x)$ **is unbiased**

Causal - inference problem

Question 1 (Incentives): How do we design quality metrics that lead to better hospital incentives?

Key result: "total treatment effect" reward function

$$w_{\mathrm{TT}}(\mathbf{Y}, \mathbf{T}^{\pi}, \mathbf{X}) = \sum_{i=1}^{n} (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi}$$

Reward based on agent treated data

Question 1 (Incentives): How do we design quality metrics that lead to better hospital incentives?

Key result: "total treatment effect" reward function



Question 2 (Ranking): Do rankings based on those quality metrics behave reasonably?



Basic ranking desiderata: better hospitals should be ranked higher.

Basic ranking desiderata: better hospitals should be ranked higher.

$$w_{\mathbf{T}\mathbf{T}}(\mathbf{Y}, \mathbf{T}^{\pi}, \mathbf{X}) = \sum_{i=1}^{n} (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi}$$

Better hospitals should be ranked higher (formal)

Definition 1 (Uniform Rank Preservation): If **hospital 1** is uniformly better than **hospital 2** for all x, then $w_1 > w_2$.



Better hospitals should be ranked higher (formal)

Definition 2 (Relative Rank Preservation): If **hospital 1** is better than **hospital 2** on average for a reference population, then $w_1 > w_2$.



Basic ranking desiderata: better hospitals should be ranked higher.

$$w_{\mathbf{T}\mathbf{T}}(\mathbf{Y}, \mathbf{T}^{\pi}, \mathbf{X}) = \sum_{i=1}^{n} (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi}$$

Basic ranking desiderata: better hospitals should be ranked higher.

$$w_{\mathrm{TT}}(\mathbf{Y}, \mathbf{T}^{\pi}, \mathbf{X}) = \sum_{i=1}^{n} (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi}$$

Basic ranking desiderata: better hospitals should be ranked higher.

$$w_{\mathrm{TT}}(\mathbf{Y}, \mathbf{T}^{\pi}, \mathbf{X}) = \sum_{i=1}^{n} (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi}$$

Problems with this reward function when comparing hospitals j and k:

• Advantages larger hospitals ($n_k > n_j$)

Basic ranking desiderata: better hospitals should be ranked higher.

$$w_{\mathrm{TT}}(\mathbf{Y}, \mathbf{T}^{\pi}, \mathbf{X}) = \sum_{i=1}^{n} (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi}$$

Problems with this reward function when comparing hospitals j and k:

- Advantages larger hospitals ($n_k > n_j$)
- Different hospitals serve diverse patient populations (Kim et al. 2022) Advantages "easier" population $X^{(k)}$ vs. $X^{(j)}$.

Kim, H., Mahmood, A., Hammarlund, N. E., & Chang, C. F. (2022). Hospital value-based payment programs and disparity in the United States: A review of current evidence and future perspectives. *Frontiers in Public Health*, 10, 882715.

$$w_k(\mathbf{Y}, \mathbf{T}, \mathbf{X}) = \sum_{i=1}^{n_k} \left(Y_i - \hat{\mu}_0(X_i^{(k)}) \right) T_i^{\pi} \left(\underbrace{\frac{1}{n_k} \frac{p_0(X_i^{(k)})}{p_k(X_i^{(k)})}}_{\uparrow} \right)$$

Reweighting term

$$w_k(\mathbf{Y}, \mathbf{T}, \mathbf{X}) = \sum_{i=1}^{n_k} \left(Y_i - \hat{\mu}_0(X_i^{(k)}) \right) T_i^{\pi} \left(\frac{1}{n_k} \frac{p_0(X_i^{(k)})}{p_k(X_i^{(k)})} \right)$$

density for reference population X_0

$$w_{k}(\mathbf{Y}, \mathbf{T}, \mathbf{X}) = \sum_{i=1}^{n_{k}} \left(Y_{i} - \hat{\mu}_{0}(X_{i}^{(k)}) \right) T_{i}^{\pi} \left(\frac{1}{n_{k}} \frac{p_{0}(X_{i}^{(k)})}{p_{k}(X_{i}^{(k)})} \right)$$

density for hospital population $X^{(k)}$

(See paper for Radon-Nikodym derivative form of the likelihood ratio.)

density for reference population X_0

 $X^{(k)}$

$$w_{k}(\mathbf{Y}, \mathbf{T}, \mathbf{X}) = \sum_{i=1}^{n_{k}} \left(Y_{i} - \hat{\mu}_{0}(X_{i}^{(k)}) \right) T_{i}^{\pi} \left(\frac{1}{n_{k}} \frac{p_{0}(X_{i}^{(k)})}{p_{k}(X_{i}^{(k)})} \right)$$

density for hospital population

Theorem 1 (Ranking Desiderata Satisfied): Hospitals will be ranked higher if they are (i) uniformly better at treating all patients, or (ii) better on average for reference population X_0 .

density for reference population X_0

$$w_k(\mathbf{Y}, \mathbf{T}, \mathbf{X}) = \sum_{i=1}^{n_k} \left(Y_i - \hat{\mu}_0(X_i^{(k)}) \right) T_i^{\pi} \left(\frac{1}{n_k} \frac{p_0(X_i^{(k)})}{p_k(X_i^{(k)})} \right)$$

density for hospital population $X^{(k)}$

Theorem 2 (Incentive Alignment): Reweighting preserves incentive-alignment.

General incentive-aligned form

Any positive function; policymaker's choice!

$$w_{\mathrm{TT}}^g(\mathbf{X}, \mathbf{T}^{\pi}, \mathbf{Y}) = \sum_{i=1}^n (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi} \underline{g(X_i)}$$

Theorem 2 (Incentive Alignment): Reweighting preserves incentive-alignment.

Accounting for policy desiderata

Any positive function; policymaker's choice!

$$w_{\mathrm{TT}}^{g}(\mathbf{X}, \mathbf{T}^{\pi}, \mathbf{Y}) = \sum_{i=1}^{n} (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi} \underline{g(X_i)}$$

Theorem 2 (Incentive Alignment): Reweighting preserves incentive-alignment.

Key result: this functional form decouples ranking desiderata from incentive alignment! A designer can optimize g for any ranking policy desiderata.

Accounting for policy desiderata

Any positive function; policymaker's choice!

$$w_{\mathrm{TT}}^g(\mathbf{X}, \mathbf{T}^{\pi}, \mathbf{Y}) = \sum_{i=1}^n (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi} g(X_i) -$$

Theorem 2 (Incentive Alignment): Reweighting preserves incentive-alignment.

Key result: this functional form decouples ranking desiderata from incentive alignment! A designer can optimize g for any ranking policy desiderata.

Example design choice: Selecting a reference population X_0 :



Accounting for policy desiderata

Any positive function; policymaker's choice!

$$w_{\mathrm{TT}}^g(\mathbf{X}, \mathbf{T}^{\pi}, \mathbf{Y}) = \sum_{i=1}^n (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi} \underline{g(X_i)}$$

Theorem 2 (Incentive Alignment): Reweighting preserves incentive-alignment.

Key result: this functional form **decouples** ranking desiderata from incentive alignment! A designer can optimize g for any ranking **policy desiderata**.

Example design choice: Selecting a reference population X_0 :







Tailor to age/gender/race/income/etc.

Question 1 (Incentives): How do we design quality metrics that lead to better hospital incentives?

Question 2 (Ranking): Are rankings based on quality metrics useful for patients?

Question 1 (Incentives): How do we design quality metrics that lead to better hospital incentives?

Question 2 (Ranking): Are rankings based on quality metrics useful for patients?

$$w_{\mathrm{TT}}^{g}(\mathbf{X}, \mathbf{T}^{\pi}, \mathbf{Y}) = \sum_{i=1}^{n} (Y_{i} - \hat{\mu}_{0}(X_{i})) T_{i}^{\pi} g(X_{i})$$

Question 1 (Incentives): How do we design quality metrics that lead to better hospital incentives?

Question 2 (Ranking): Are rankings based on quality metrics useful for patients?



Question 1 (Incentives): How do we design quality metrics that lead to better hospital incentives?

Question 2 (Ranking): Are rankings based on quality metrics useful for patients?



Question 1 (Incentives): How do we design quality metrics that lead to better hospital incentives?

Question 2 (Ranking): Are rankings based on quality metrics useful for patients?



Question 3 (Information Asymmetry): What if the hospitals know more about patients than the agencies?

"Providers may be able to improve their ranking by selecting patients on the basis of characteristics that are **unobservable to the analysts but predictive of good outcomes**" (Dranove et al., 2003).

Principal (Health Dept)

Agent (Hospital)

Principal (Health Dept)

Observes

Covariates $\mathbf{X} = X_1, ..., X_n$ Treatments $\mathbf{T} = T_1, ..., T_n$ Outcomes $\mathbf{Y} = Y_1, ..., Y_n$ Agent (Hospital)

Principal (Health Dept)

Observes

Covariates $\mathbf{X} = X_1, ..., X_n$ Treatments $\mathbf{T} = T_1, ..., T_n$ Outcomes $\mathbf{Y} = Y_1, ..., Y_n$

Chooses

Reward function $w(\mathbf{Y},\mathbf{T},\mathbf{X})$

Agent (Hospital)

Principal (Health Dept)

Observes

Covariates
$$\mathbf{X} = X_1, ..., X_n$$

Treatments $\mathbf{T} = T_1, ..., T_n$
Outcomes $\mathbf{Y} = Y_1, ..., Y_n$

Chooses

Reward function $w(\mathbf{Y},\mathbf{T},\mathbf{X})$

Agent (Hospital)

Observes

Covariates
$$\mathbf{X} = X_1, ..., X_n$$

 $\mathbf{U} = U_1, ..., U_n$

Principal (Health Dept)

Observes

Covariates
$$\mathbf{X} = X_1, ..., X_n$$

Treatments $\mathbf{T} = T_1, ..., T_n$
Outcomes $\mathbf{Y} = Y_1, ..., Y_n$

Chooses

Reward function $w(\mathbf{Y},\mathbf{T},\mathbf{X})$

Agent (Hospital)

Observes

Covariates
$$\mathbf{X} = X_1, ..., X_n$$

 $\mathbf{U} = U_1, ..., U_n$

Knows

Expected potential outcomes $E[Y_i(1)|X_i, U_i] = E[Y_i(0)|X_i, U_i]$

Principal (Health Dept)

Observes

Covariates
$$\mathbf{X} = X_1, ..., X_n$$

Treatments $\mathbf{T} = T_1, ..., T_n$
Outcomes $\mathbf{Y} = Y_1, ..., Y_n$

Chooses

Reward function $w(\mathbf{Y},\mathbf{T},\mathbf{X})$

Agent (Hospital)

Observes

Covariates
$$\mathbf{X} = X_1, ..., X_n$$

 $\mathbf{U} = U_1, ..., U_n$

Knows

Expected potential outcomes $E[Y_i(1)|X_i, U_i] = E[Y_i(0)|X_i, U_i]$

Chooses

Treatment rule

$$\pi(x, u) = P(T_i = 1 | X_i = x, U_i = u)$$
$$w_{\mathrm{TT}}(\mathbf{Y}, \mathbf{T}^{\pi}, \mathbf{X}) = \sum_{i=1}^{n} (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi}$$

$$w_{\mathrm{TT}}(\mathbf{Y}, \mathbf{T}^{\pi}, \mathbf{X}) = \sum_{i=1}^{n} (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi}$$

Confounding bias may or may not be present in auxiliary data: $E[\hat{\mu}_0(x)] \neq \mu_0(x)$

$$w_{\mathrm{TT}}(\mathbf{Y}, \mathbf{T}^{\pi}, \mathbf{X}) = \sum_{i=1}^{n} (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi}$$

Confounding bias may or may not be present in auxiliary data: $E[\hat{\mu}_0(x)] \neq \mu_0(x)$ Unmeasured heterogeneity: Only measuring $\mu_0(x)$, but agent chooses $\pi(x, u)$.

$$w_{\mathrm{TT}}(\mathbf{Y}, \mathbf{T}^{\pi}, \mathbf{X}) = \sum_{i=1}^{n} (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi}$$
Confounding bias may or may
not be present in auxiliary data:
$$E[\hat{\mu}_0(x)] \neq \mu_0(x)$$
Unmeasured heterogeneity:
Only measuring $\mu_0(x)$,
but agent chooses $\pi(x, u)$.

Key result: When there is information asymmetry, unconfoundedness is not enough! There is still regret, even if $\hat{\mu}_0(x)$ is unbiased.

Assumption (Bounded Unmeasured Heterogeneity): the effect of U on the untreated potential outcome given X is bounded:

$$E[|\mu_0(X_i) - \mu_0(X_i, U_i)|] \le \gamma_{\text{marg}}$$

Assumption (Bounded Unmeasured Heterogeneity): the effect of U on the untreated potential outcome given X is bounded:

$$E[|\mu_0(X_i) - \mu_0(X_i, U_i)|] \le \gamma_{\text{marg}}$$

Upper bound: the regret under reward $w_{\rm TT}$ is upper bounded by the degree of unmeasured heterogeneity: $R(\pi^{w_{\rm TT}}) \leq 2\gamma_{\rm marg}$

Assumption (Bounded Unmeasured Heterogeneity): the effect of U on the untreated potential outcome given X is bounded:

$$E[|\mu_0(X_i) - \mu_0(X_i, U_i)|] \le \gamma_{\text{marg}}$$

Upper bound: the regret under reward $w_{\rm TT}$ is upper bounded by the degree of unmeasured heterogeneity: $R(\pi^{w_{\rm TT}}) \leq 2\gamma_{\rm marg}$

Lower bound: the regret is lower bounded by $\gamma_{
m marg}$.

Question 1 (Incentives): How do we design quality metrics that lead to better hospital incentives?

Question 2 (Ranking): Do rankings based on those quality metrics behave reasonably?

Question 3 (Information Asymmetry): What if the hospitals know more about patients than the agencies?

Question 1 (Incentives): How do we design quality metrics that lead to better hospital incentives?

Key results: Reward total treatment effect by estimating the untreated counterfactual.

Question 2 (Ranking): Do rankings based on those quality metrics behave reasonably?

Question 3 (Information Asymmetry): What if the hospitals know more about patients than the agencies?

Question 1 (Incentives): How do we design quality metrics that lead to better hospital incentives?

Key results: Reward total treatment effect by estimating the untreated counterfactual.

Question 2 (Ranking): Do rankings based on those quality

metrics behave reasonably?

Key results: Not necessarily, but reweighting preserves incentive alignment.

Question 3 (Information Asymmetry): What if the hospitals know more about patients than the agencies?

Question 1 (Incentives): How do we design quality metrics that lead to better hospital incentives?

Key results: Reward total treatment effect by estimating the untreated counterfactual.

Question 2 (Ranking): Do rankings based on those quality

metrics behave reasonably?

Key results: Not necessarily, but reweighting preserves incentive alignment.

Question 3 (Information Asymmetry): What if the hospitals know more about patients than the agencies?

Key results: Regret exists even without confounding. Success depends on the degree of unmeasured heterogeneity.

Can **robust policy learning** help mitigate problems of information asymmetry?



Can **robust policy learning** help mitigate problems of information asymmetry?





Incentives for information sharing: who has information on how to improve metrics, and when would they share it?

Model extensions:

- Cost of treatment/resource constraints
- Patient decisions
- Proxies for health risk
- Competition between hospitals
- Competition between ranking platforms
- Truthfulness
- Hospital investment in improvement
- Multiple metrics
- ...etc.

Behind every algorithm is a metric.

Many thanks!

serenalwang@berkeley.edu serenalwang.com

Appendix

Example data sources for estimating $\hat{\mu}_0(x)$

Source 1: Auxiliary/historical untreated data

Requirements: Ignorability: $Y_i(0), Y_i(1) \perp T_i | X_i, P(T_i = 1 | X_i = x) < 1$ **Advantages:** Can come from untreated observational data. **Challenges:** Confounding; distribution shift.

Source 2: Agent's untreated units

Requirements: The agent's treatment policy depends only on X_i . **Advantages:** No distribution shift. **Challenges:** Unable to guarantee positivity.

Central threat: Confounding

Change 1: reward benefit from treatment

Reward function 2: "Average Treatment Effect on the Treated" (ATT)

$$w_{\text{ATT}}(\mathbf{Y}, \mathbf{T}^{\pi}, \mathbf{X}) = \frac{\sum_{i=1}^{n} (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi}}{\sum_{i=1}^{n} T_i^{\pi}}$$
Estimated untreated counterfactual

Agent's best response: Treat only the patients with the highest expected treatment effect.

$$\pi^{w_{\text{ATT}}}(x) = \mathbb{1}(x \in \operatorname{arg\,max}_x \tau(x) \text{ and } \tau(x) > 0)$$

Proposition 2: Regret is at most the max utility: $R(\pi^{w_{\text{ATT}}}) \leq \max_{\pi \in \Pi} V(\pi)$

Today



The US government spent **\$1.3 billion** on quality metric development over the last 10 years (Wadhera et al. 2020).

Wadhera, R. K., Figueroa, J. F., Maddox, K. E. J., Rosenbaum, L. S., Kazi, D. S., & Yeh, R. W. (2020). Quality measure development and associated spending by the Centers for Medicare & Medicaid Services. JAMA, 323(16), 1614-1616.



Merit-Based Incentive Payment System (MIPS)







Centers for Medicare and Medicaid Services (CMS)

Patients

Merit-Based Incentive Payment System (MIPS)



Best Hospitals Honor Roll

Each year, U.S. News ranks hospitals in 15 specialties and 20 procedures and conditions. Here are the top 20 highest rated hospitals in the U.S.

2022-2023 Honor Roll Rankings #1 Mayo Clinic Rochester, MN

- #2 Cedars-Sinai Medical Center Los Angeles, CA
- #3 NYU Langone Hospitals New York, NY

SEE FULL RANKINGS LIST »



Accounting for policy desiderata

Any positive function; policymaker's choice!

$$w_{\mathrm{TT}}^g(\mathbf{X}, \mathbf{T}^{\pi}, \mathbf{Y}) = \sum_{i=1}^n (Y_i - \hat{\mu}_0(X_i)) T_i^{\pi} g(X_i) -$$

Theorem 2 (Incentive Alignment): Reweighting preserves incentive-alignment.

Key result: this functional form decouples ranking desiderata from incentive alignment! A designer can optimize g for any ranking policy desiderata.

Example design choice: Boost for public vs. private institutions.





Assumption (Bounded Unmeasured Heterogeneity): the effect of U on the untreated potential outcome given X is bounded:

$$E[|\mu_0(X_i) - \mu_0(X_i, U_i)|] \le \gamma_{\text{marg}}$$

Theorem (Bounded Regret): the regret under reward $w_{\rm TT}$ is upper bounded by the degree of unmeasured heterogeneity: $R(\pi^{w_{\rm TT}}) \leq 2\gamma_{\rm marg}$

Low γ : X= smoking, U= drinking "The magnitude of risk related to smoking is far larger than any ostensible benefit related to moderate drinking" (Mukamal 2006).

Mukamal, K. J. (2006). The effects of smoking and drinking on cardiovascular disease and risk factors. Alcohol Research & Health, 29(3), 199.

Assumption (Bounded Unmeasured Heterogeneity): the effect of U on the untreated potential outcome given X is bounded:

$$E[|\mu_0(X_i) - \mu_0(X_i, U_i)|] \le \gamma_{\text{marg}}$$

Theorem (Bounded Regret): the regret under reward $w_{\rm TT}$ is upper bounded by the degree of unmeasured heterogeneity: $R(\pi^{w_{\rm TT}}) \leq 2\gamma_{\rm marg}$

Low γ : X= smoking, U= drinking "The magnitude of risk related to smoking is far larger than any ostensible benefit related to moderate drinking" (Mukamal 2006). **High** γ : X = sex hormones, U = diabetes

"Cardiovascular risks associated with diabetes also appear to be higher in women" (Rodgers et al. 2019).

Rodgers, J. L. et al. (2019). Cardiovascular risks associated with gender and aging. *Journal of Cardiovascular Development and Disease*, 6(2), 19. Mukamal, K. J. (2006). The effects of smoking and drinking on cardiovascular disease and risk factors. *Alcohol Research & Health*, 29(3), 199.