



# Scorecards, Benchmarking, and the Search for Unusual Hospitals, Communities, and Cops

Greg Ridgeway

Department of Criminology

Department of Statistics and Data Science

# Scorecards are Popular

## Texas Education Scorecard

### Accountability Rating

C

HOUSTON HEIGHTS CHARTER SCHOOL earned a C (70-79) for acceptable performance by serving many students well but needs to provide additional academic support to many more students.

State accountability ratings are based on three domains: Student Achievement, School Progress, and Closing the Gaps. The graph below provides summary results for HOUSTON HEIGHTS CHARTER SCHOOL. Scores are scaled from 0 to 100 to align with letter grades.

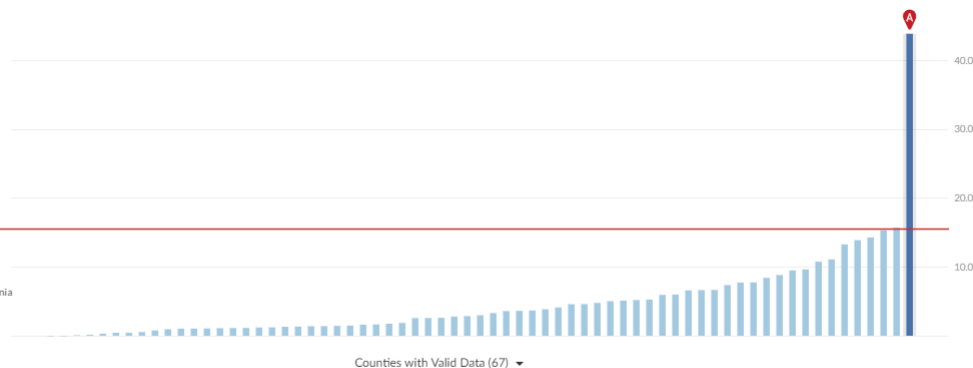
Student Achievement

School Progress

Closing the Gaps

15.52%  
of Cases

Failure to Pay  
Low Monetary  
Bail  
Statewide in Pennsylvania



In Philadelphia County, PA, of the cases in which defendants failed to pay monetary bail, 43.94% had \$500 bail or less.

## California Children's Well-Being

All topics Health Education Child Welfare Early Childhood

Los Angeles %'s for all races  
Viewing 1-6 of 6

Children in the child welfare system who exited to permanency within one year	35%
Children in the child welfare system who had been in one placement after 24 months in care	39%
Adolescents in the child welfare system who were placed in family-like settings	74%
	65%
	78%
	47%

## Measures for Justice

Have these scorecards addressed fundamental differences?

# Scorecards Generate Reputational Concern for Local Government

## Accountability Rating

C

HOUSTON HEIGHTS CHARTER SCHOOL earned a C (70-79) for acceptable performance by serving many students well but needs to provide additional academic support to many more students.

State acco  
Student A  
below pro  
Scores are

All topics Health Education Child Welfare Early Childhood

Los Angeles %s for all races  
Viewing 1–6 of 6

Children in the child welfare system who exited to permanency within one year	35%
Children in the child welfare system who had been in one placement after 24 months in care	39%
Adolescents in the child welfare system who were placed in family-like settings	74%
Children in the child welfare system who had a timely dental exam	65%
Children in the child welfare system who had a timely medical exam	78%

15.52%  
of Cases

Failure to Pay  
Low Monetary  
Bail  
Statewide in Pennsylvania

Counties with Valid Data (67)

In Philadelphia County, PA, of the cases in which defendants failed to pay monetary bail, 43.94% had \$500 bail or less.

Focus on issues that

- are sensitive for the government
- garner exposure through public attention
- force governments to prioritize the issues measured in the scorecard

# Four benchmarking applications

- Which officers stop black pedestrians at an unusual rate?
- Which communities are particularly dissatisfied with the police?
- Which counties contribute most to racial disparities in incarceration sentences?
- Which hospitals have...
  - excessive opioid prescriptions?
  - unusually high mortality rates?

# Four benchmarking applications

- Which officers stop black pedestrians at an unusual rate?
- Which communities are particularly dissatisfied with the police?
- Which counties contribute most to racial disparities in incarceration sentences?
- Which hospitals have...
  - excessive opioid prescriptions?
  - unusually high mortality rates?

G. Ridgeway and J.M. MacDonald (2009). “Doubly Robust Internal Benchmarking and False Discovery Rates for Detecting Racial Bias in Police Stops,” *Journal of the American Statistical Association* 104:661–668

# Is an Officer Who Stops 86% Black Pedestrians Unusual?

Stop Characteristic	Example Officer (%) n = 392	
% black pedestrians stopped	86%	

- Combine three statistical techniques to answer this question
  - Propensity score weighting
  - Doubly robust estimation
  - False discovery rate

# We Know a Lot About the Environment of this Officer's Stops

Stop Characteristic		Example Officer (%) n = 392	
% black pedestrians stopped		86%	
Month	January	3	
	February	4	
	March	8	
Day of the week	Monday	13	
	Tuesday	11	
	Wednesday	14	
Time of day	(4-6 p.m.)	9	
	(6-8 p.m.)	8	
	(8-10 p.m.)	23	
	(10 p.m. -12 a.m.)	17	
Patrol borough	Brooklyn North	100	
Precinct	B	98	
	C	1	
Outside		96	
In uniform	Yes	99	
Radio run	Yes	1	

# We Also Know the Exact Location of This Officer's Stops



**Example Officer**



# Idea: Reweight Stops Made By Other Officers to Resemble This Officer's Stops



Example Officer

- Align their distributions  
 $f(\mathbf{x}|t = 1) = w(\mathbf{x})f(\mathbf{x}|t = 0)$

Example officer

Other officers

- Solving for  $w(\mathbf{x})$  yields the propensity score weight

$$w(\mathbf{x}) \propto \frac{P(t = 1|\mathbf{x})}{1 - P(t = 1|\mathbf{x})}$$

- Estimate  $P(t = 1|\mathbf{x})$  using boosted logistic regression as implemented in `gbm`

# Reweightings Aligns the Distribution of Stop Locations



**Example Officer**



**Matched Stops**

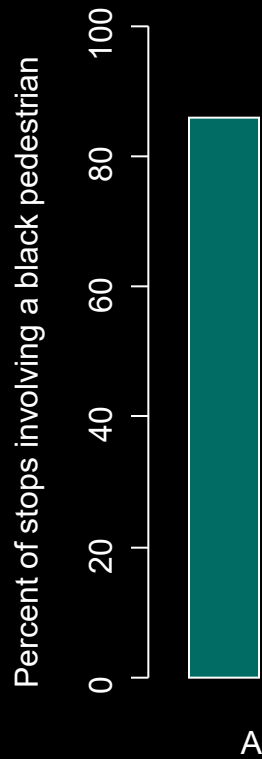
# Reweightings Also Aligns the Distribution of All Other Stop Features

Stop Characteristic		Example Officer (%) n = 392	Internal Benchmark (%) ESS = 3,676
% black pedestrians stopped		86%	
Month	January	3	3
	February	4	4
	March	8	9
Day of the week	Monday	13	13
	Tuesday	11	10
	Wednesday	14	15
Time of day	(4-6 p.m.)	9	10
	(6-8 p.m.)	8	8
	(8-10 p.m.)	23	23
	(10 p.m. -12 a.m.)	17	17
Patrol borough	Brooklyn North	100	100
Precinct	B	98	98
	C	1	1
Outside		96	94
In uniform	Yes	99	97
Radio run	Yes	1	3

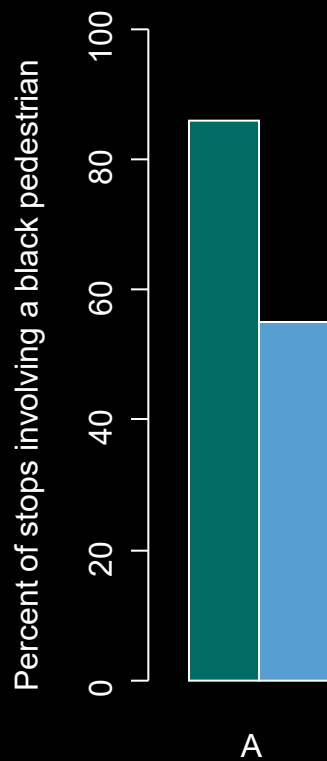
# Colleagues at the Same Time, Place, and Context Stop 55% Black Pedestrians

Stop Characteristic		Example Officer (%) n = 392	Internal Benchmark (%) ESS = 3,676
% black pedestrians stopped		86%	55%
Month	January	3	3
	February	4	4
	March	8	9
Day of the week	Monday	13	13
	Tuesday	11	10
	Wednesday	14	15
Time of day	(4-6 p.m.)	9	10
	(6-8 p.m.)	8	8
	(8-10 p.m.)	23	23
	(10 p.m. -12 a.m.)	17	17
Patrol borough	Brooklyn North	100	100
Precinct	B	98	98
	C	1	1
Outside		96	94
In uniform	Yes	99	97
Radio run	Yes	1	3

# 86% of the Officer's Stops Were Black...



# ...Compared with 55% for the Benchmark



- Doubly robust benchmark estimate obtainable from weighted logistic regression

$$\ell(\beta) = \sum_{i=1}^n w_i \left( y_i (\delta t_i + \beta' \mathbf{x}_i) - \log(1 + e^{\delta t_i + \beta' \mathbf{x}_i}) \right)$$

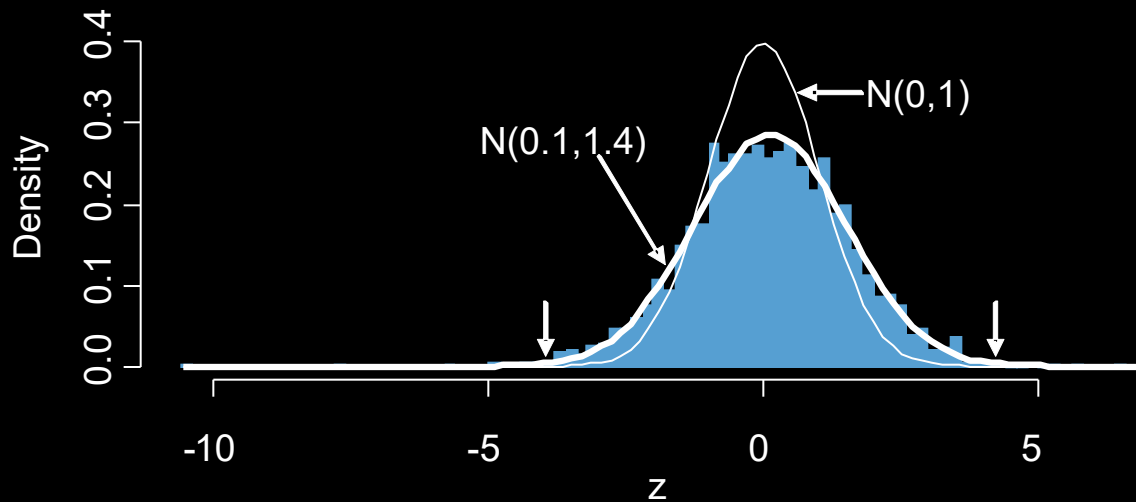
- Disparity computed as

$$\hat{\theta}_A^{DR} = \frac{1}{\sum t_i} \sum_{i=1}^n t_i \left( \underbrace{\frac{1}{1 + \exp(-\delta - \beta' \mathbf{x}_i)}}_{\text{Predicted probability stopped pedestrian is black for example officer}} - \underbrace{\frac{1}{1 + \exp(-\beta' \mathbf{x}_i)}}_{\text{Predicted probability stopped pedestrian is black for other officers}} \right)$$

Predicted probability  
stopped pedestrian is  
black for example officer

Predicted probability  
stopped pedestrian is  
black for other officers

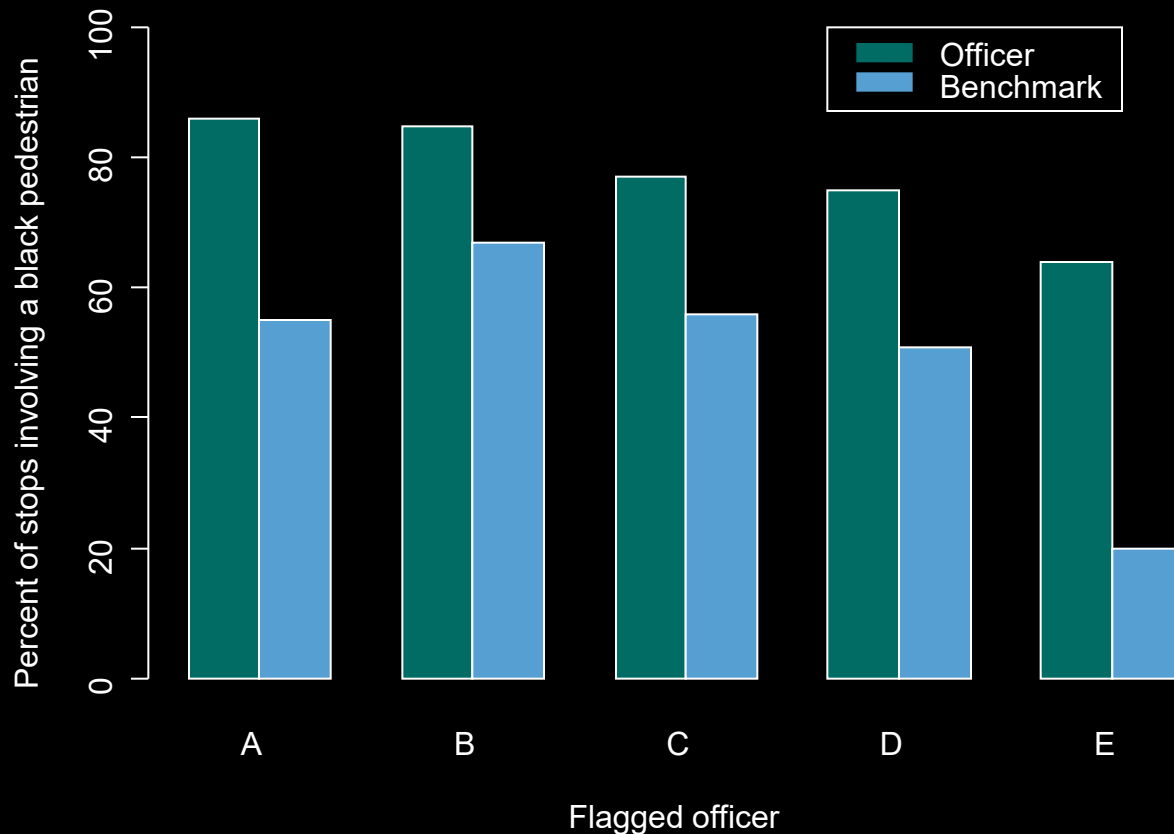
# Repeat for Nearly 3,000 NYPD Officers Actively Involved in Stops



- $$P(\text{problem}|z) = 1 - \frac{f(z|\text{no problem})f(\text{no problem})}{f(z)}$$

$$\geq 1 - \frac{f_0(z)}{f(z)}$$
- Right tail consists of 5 officers with “problem officer” probabilities in excess of 50%
- Standard cutoff of  $z > 2.0$  flags 242 officers, 90% of which have  $\text{fdr}$  estimated to be greater than 0.999

# Analysis in NYPD Flagged Five Officers





# Benchmarking

1. Apply propensity score weights so benchmark is based on activities in similar context
2. Compute z-statistic from a propensity score weighted regression
3. Repeat for all units, customizing benchmark for each
4. Compute false discovery rate based on empirical distribution of z-statistics

# Four benchmarking applications

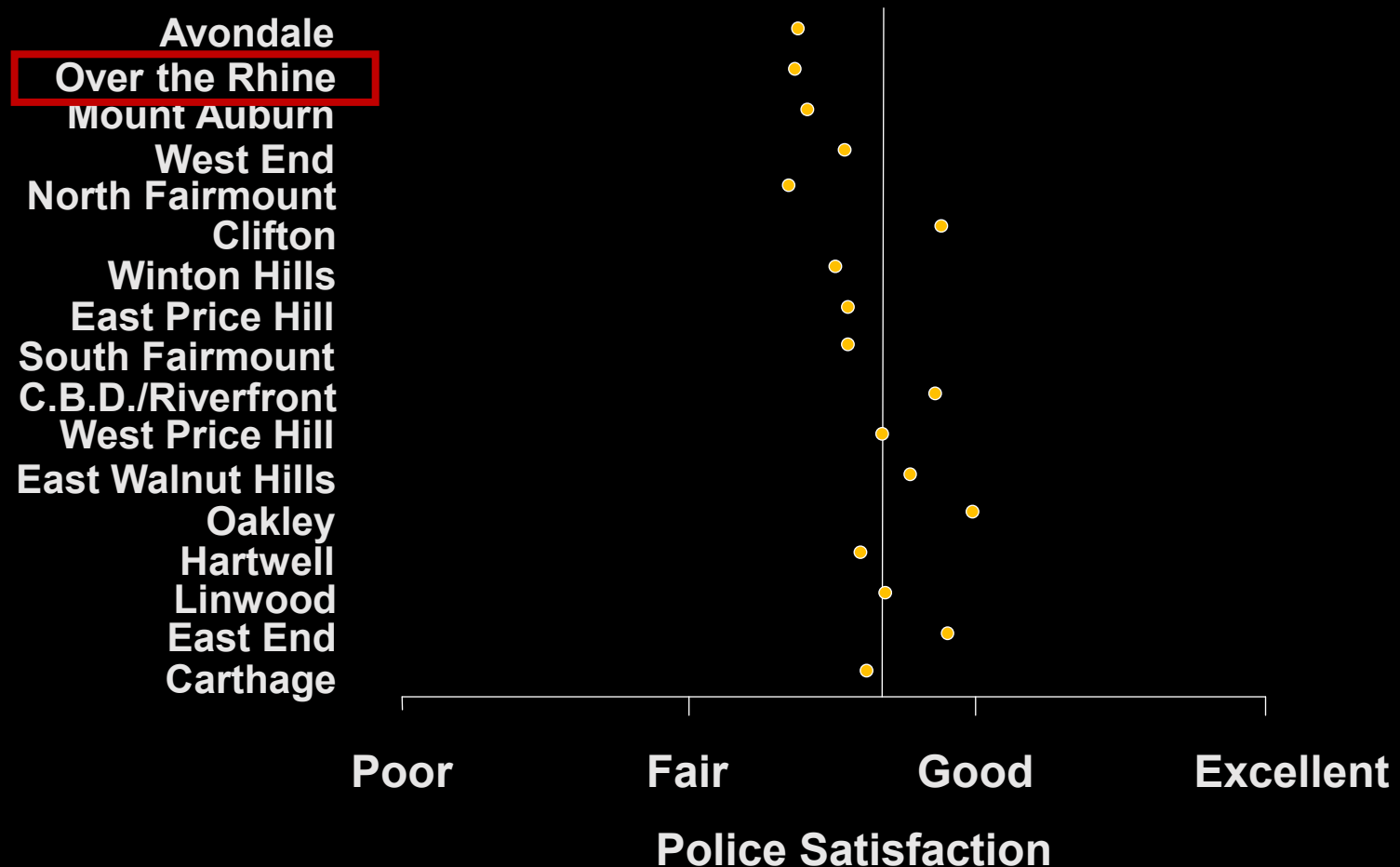
- Which officers stop black pedestrians at an unusual rate?
- Which communities are particularly dissatisfied with the police?
- Which counties contribute most to racial disparities in incarceration sentences?
- Which hospitals have...
  - excessive opioid prescriptions?
  - unusually high mortality rates?

G. Ridgeway and J.M. MacDonald (2014). "A Method for Internal Benchmarking of Criminal Justice System Performance," *Crime & Delinquency* 60(1):145-162

# In Which Neighborhoods Are Police Underperforming?

- Cincinnati Police Department sponsored a citywide survey of citizens
  - Citizen satisfaction with the police
  - Perceptions of racially discriminatory police practices
  - Whether residents felt that they had personally experienced racial profiling
- 6,000 residents in Cincinnati selected via random-digit dialing and list-assisted sampling methods
  - Stratified to cover 45 defined Cincinnati neighborhoods
  - Respondents were 18 years or older

# Neighborhoods Varied on Police Satisfaction



# Respondents Differ on Key Features Associated with Police Satisfaction

Respondent characteristics	Respondents from Over-the-Rhine (N=146)	Respondents from other neighborhoods (N=5,671)	
Less than HS	21	10	

# Respondents Differ on Key Features Associated with Police Satisfaction

Respondent characteristics	Respondents from Over-the-Rhine (N=146)	Respondents from other neighborhoods (N=5,671)
Less than HS	21	10
College degree+	23	33
Black	66	42
White	30	53
\$20,000 or less	47	25
\$100,000 or more	6	11
Employed (%)	60	58
Married (%)	15	38
Male (%)	43	36
Age 22-29	16	8
Age 65+	13	25
Homeowner (%)	20	60
Children at home (%)	40	31

# Constructed Benchmark Matches Neighborhoods on These Features

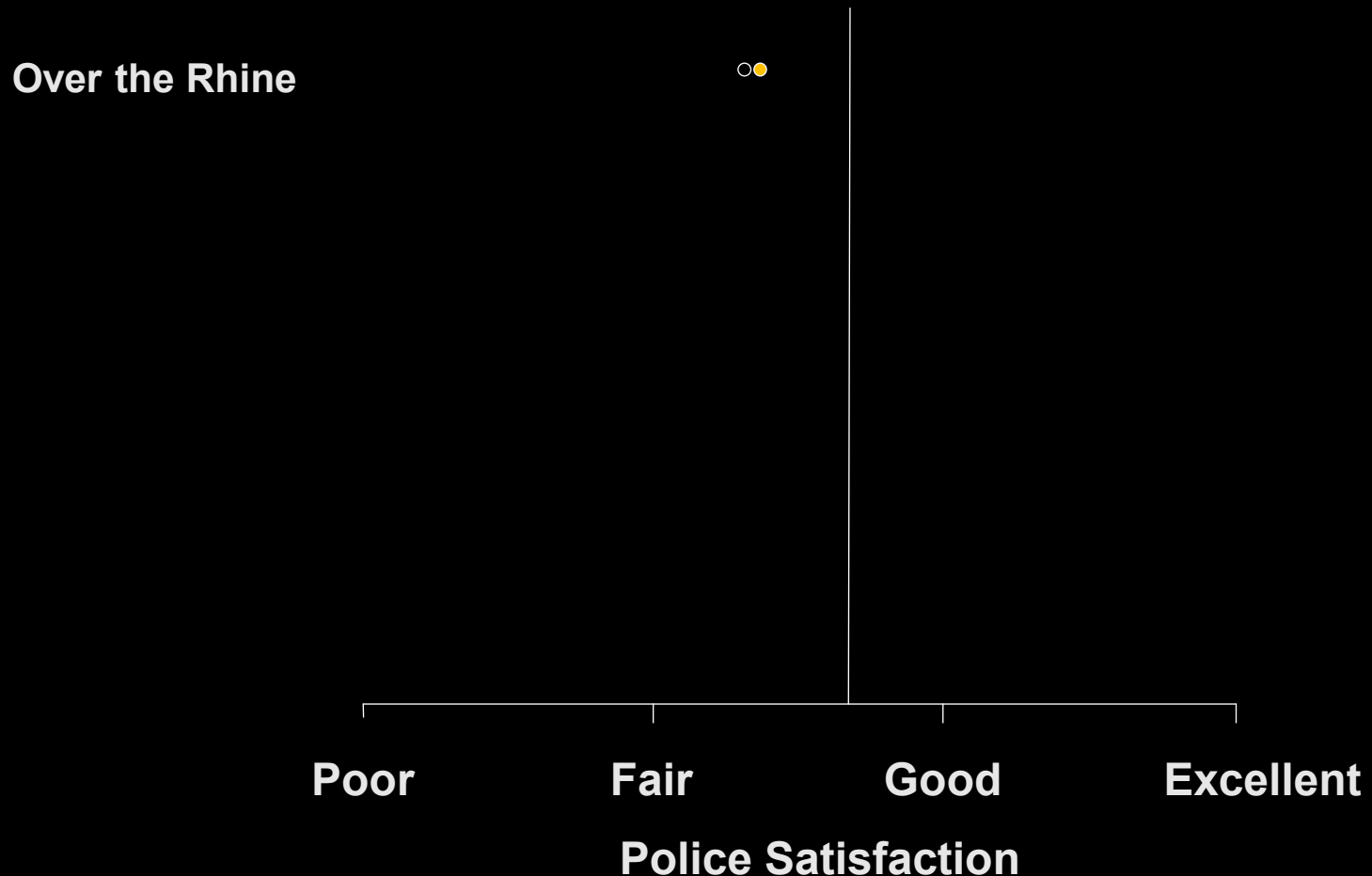
Respondent characteristics	Respondents from Over-the-Rhine (N=146)	Respondents from other neighborhoods (N=5,671)	Weighted respondents from other neighborhoods (N=422)
Less than HS	21	10	21
College degree+	23	33	22
Black	66	42	65
White	30	53	32
\$20,000 or less	47	25	45
\$100,000 or more	6	11	5
Employed (%)	60	58	58
Married (%)	15	38	16
Male (%)	43	36	42
Age 22-29	16	8	17
Age 65+	13	25	13
Homeowner (%)	20	60	21
Children at home (%)	40	31	38

# Constructed Benchmark Matches Neighborhoods on These Features

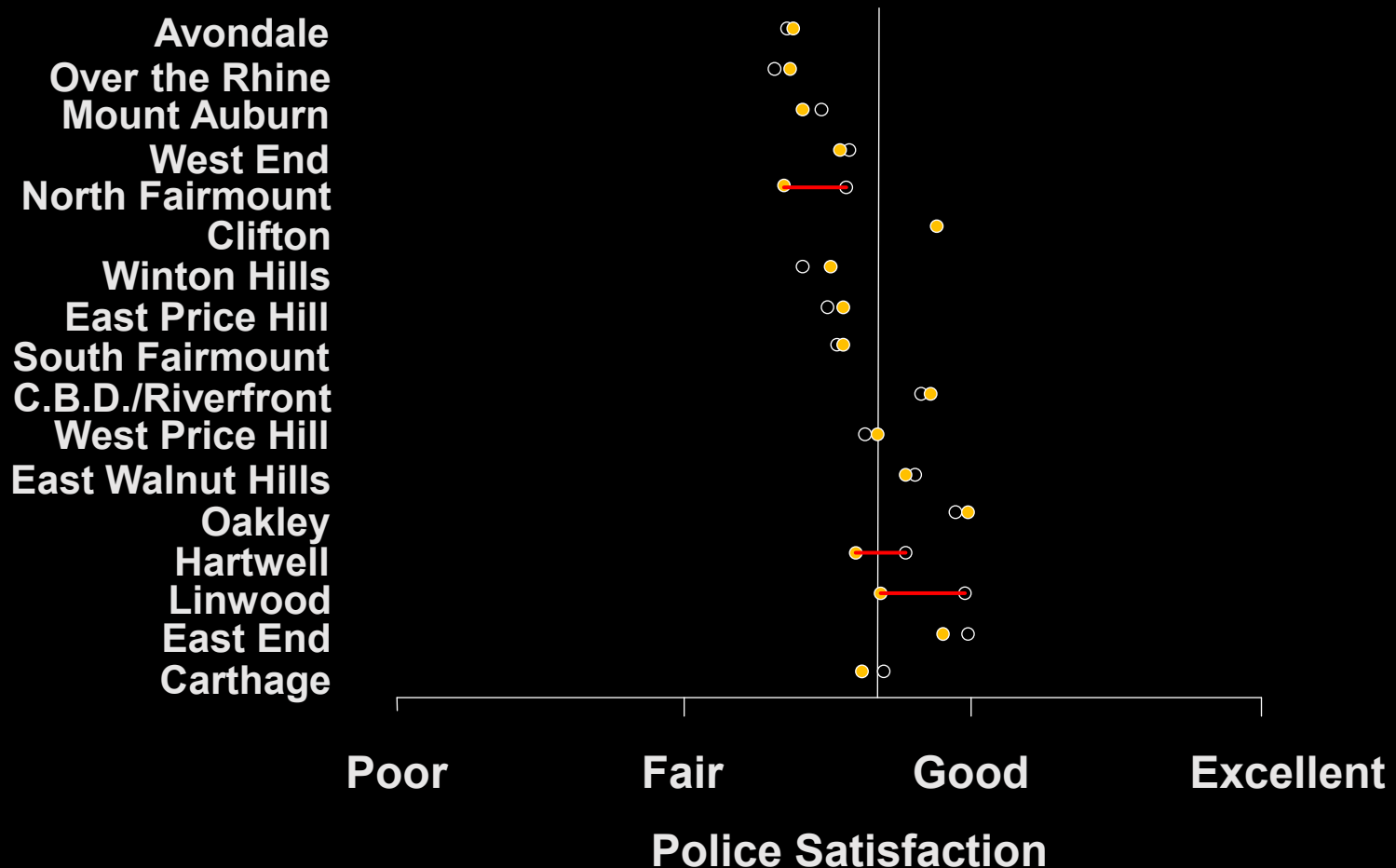
Respondent characteristics	Respondents from Over-the-Rhine (N=146)	Respondents from other neighborhoods (N=5,671)	Weighted respondents from other neighborhoods (N=422)
Less than HS	21	10	21
College degree+	23	33	22
Black	66	42	65
White	30	53	32
\$20,000 or less	47	25	45
\$100,000 or more	6	11	5
Employed (%)	60	58	58
Married (%)	15	38	16
Male (%)	43	36	42
Age 22-29	16	8	17
Age 65+	13	25	13
Homeowner (%)	20	60	21
Children at home (%)	40	31	38



# Police Satisfaction in Over the Rhine is Close to Expectation



# Few Neighborhoods Differ from Benchmarks



# Four benchmarking applications

- Which officers stop black pedestrians at an unusual rate?
- Which communities are particularly dissatisfied with the police?
- Which counties contribute most to racial disparities in incarceration sentences?
- Which hospitals have...
  - unusually high mortality and readmission rates?
  - excessive opioid prescriptions?

G. Ridgeway, R. Moyer, and S. Bushway (2020). "Sentencing Scorecards: Reducing Racial Disparities in Prison Sentences at Their Source," *Criminology & Public Policy* 19(4):1113-1138

# New York's Permanent Commission on Sentencing Sought to Identify Sources of Racial Disparities

- Created in 2010, the Commission was charged with evaluating sentencing laws and practices and recommending reforms to improve sentencing policy
- Commission decided to move forward with new statewide analytical efforts on racial disparity
- Data for this analysis comes from the New York Division of Criminal Justice Services (DCJS)
  - 584,299 felony cases
  - January 1, 2000 and December 31, 2014
  - Detailed information about defendant, their criminal history, and case features

# Match Defendants on Detailed Case and Defendant Features

[illegible]

# Match Defendants on Detailed Case and Defendant Features

[illegible]

# Match Defendants on Detailed Case and Defendant Features

[illegible]

# Match Defendants on Detailed Case and Defendant Features

[illegible]



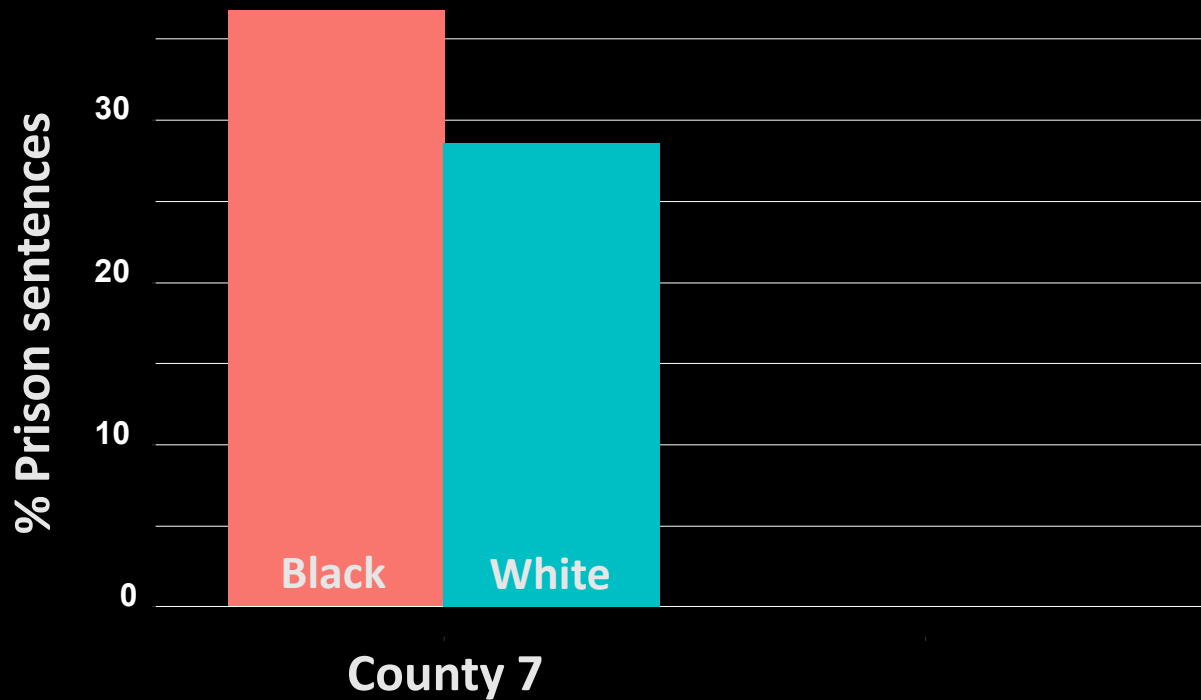
# Match Defendants on Detailed Case and Defendant Features

Case/defendant feature	Within County 7		Outside County 7	
	Black	White	Black	White
	n < 2,000	ESS = 1,354	ESS = 19,402	ESS = 29,977
Age at arrest (average)	30.4	30.4	30.3	30.3
Male (%)	81.0	81.1	82.1	82.3
No prior felony arrests (%)	49.5	52.0	47.2	48.1
Prior arrests (average count)				
Felonies	1.4	1.2	1.7	1.4
Drugs	0.4	0.4	0.6	0.5
Firearms	0.07	0.05	0.08	0.05
Violent crimes	0.4	0.3	0.5	0.4
Prior convictions (average count)				
Weapons	0.1	0.1	0.1	0.1
Violent crimes	0.1	0.1	0.1	0.1
Specific top charge (%)				
PL 120.05(02)	3.8	3.6	3.8	3.8
PL 140.25(02)	7.1	7.2	7.1	7.0
PL 155.30	5.6	6.0	5.7	5.9
PL 220.39(01)	6.5	5.6	6.0	6.6

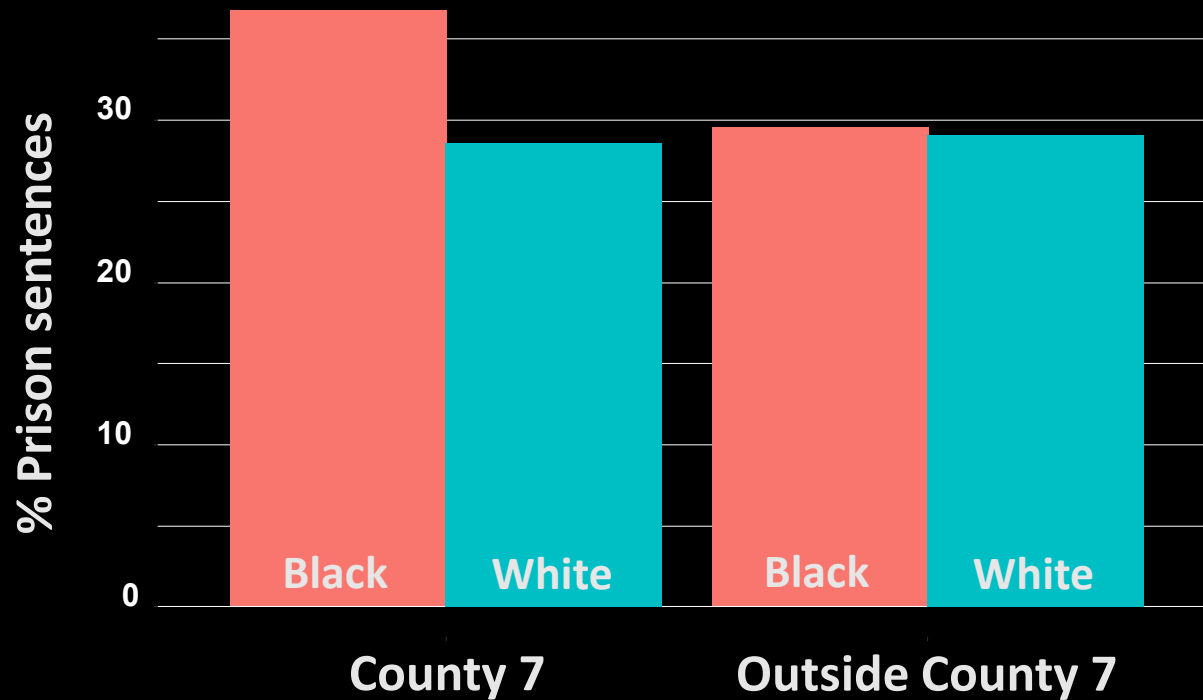
# Match Defendants on Detailed Case and Defendant Features

Case/defendant feature	Within County 7		Outside County 7	
	Black n < 2,000	White ESS = 1,354	Black ESS = 19,402	White ESS = 29,977
Age at arrest (average)	30.4	30.4	30.3	30.3
Male (%)	81.0	81.1	82.1	82.3
No prior felony arrests (%)	49.5	52.0	47.2	48.1
Prior arrests (average count)				
Felonies	1.4	1.2	1.7	1.4
Drugs	0.4	0.4	0.6	0.5
Firearms	0.07	0.05	0.08	0.05
Violent crimes	0.4	0.3	0.5	0.4
Prior convictions (average count)				
Weapons	0.1	0.1	0.1	0.1
Violent crimes	0.1	0.1	0.1	0.1
Specific top charge (%)				
PL 120.05(02)	3.8	3.6	3.8	3.8
PL 140.25(02)	7.1	7.2	7.1	7.0
PL 155.30	5.6	6.0	5.7	5.9
PL 220.39(01)	6.5	5.6	6.0	6.6
General top charge features (%)				
Violent crime	20.5	20.0	21.5	20.6
Class D felony	39.7	40.0	39.1	38.7
Firearm Related	5.0	4.2	4.8	4.7
Disposition month: June (%)	11.4	9.4	8.1	8.4

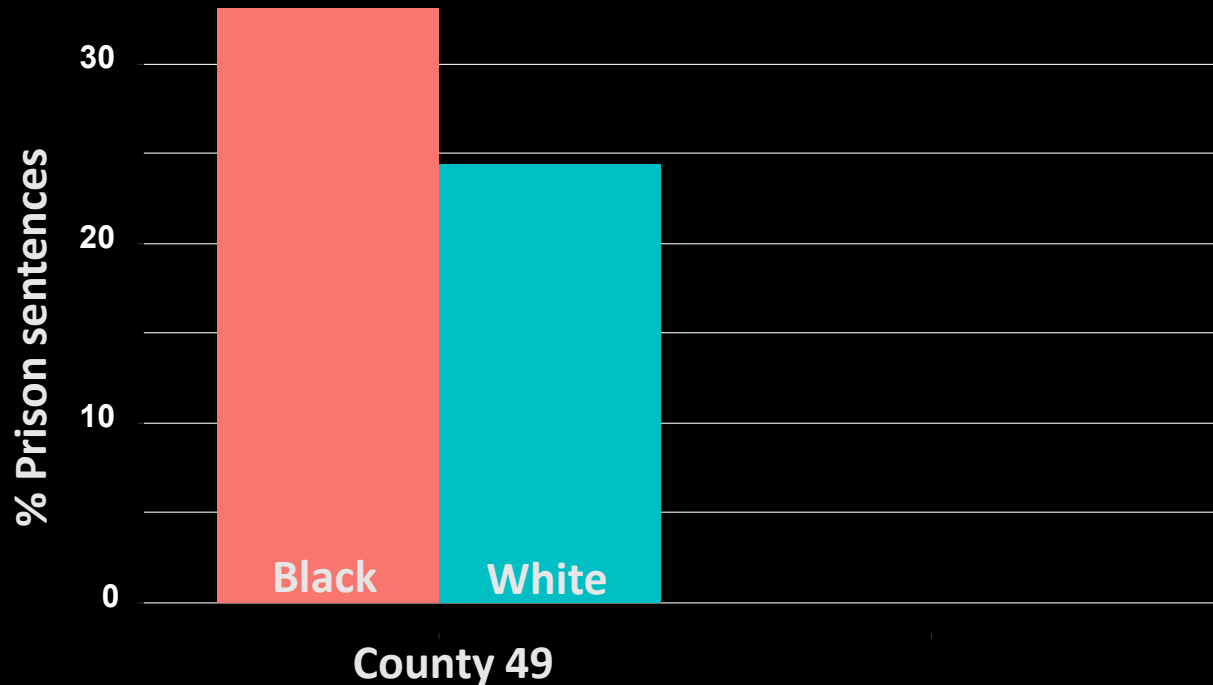
# County 7 Incarcerates Black Defendants at Higher Rates



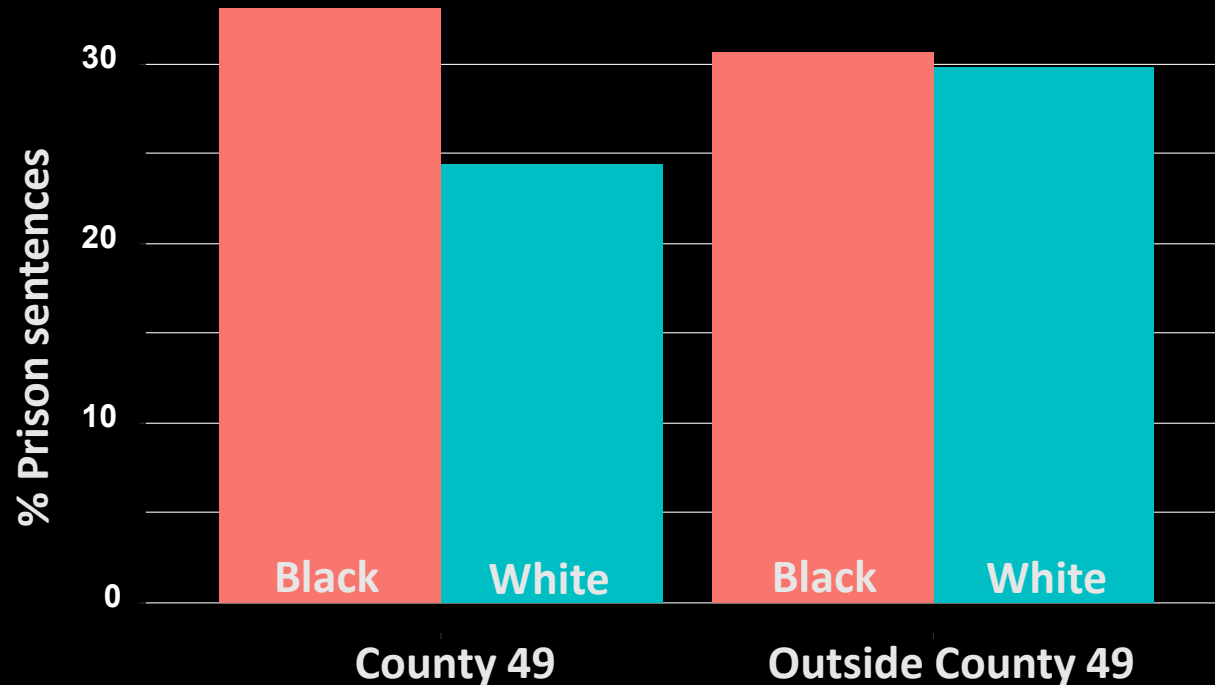
# County 7 Incarcerates Black Defendants at Higher Rates



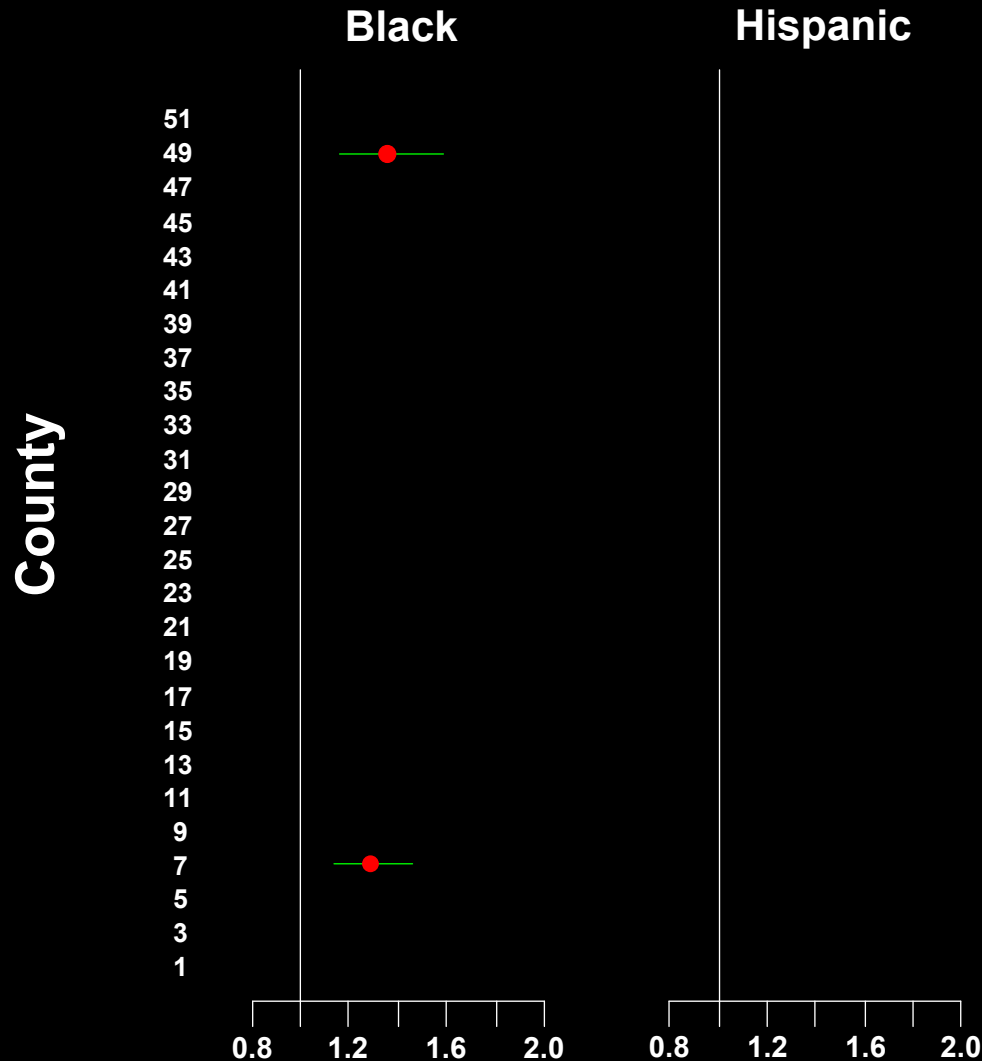
# County 49 Incarcerates Black Defendants at Higher Rates



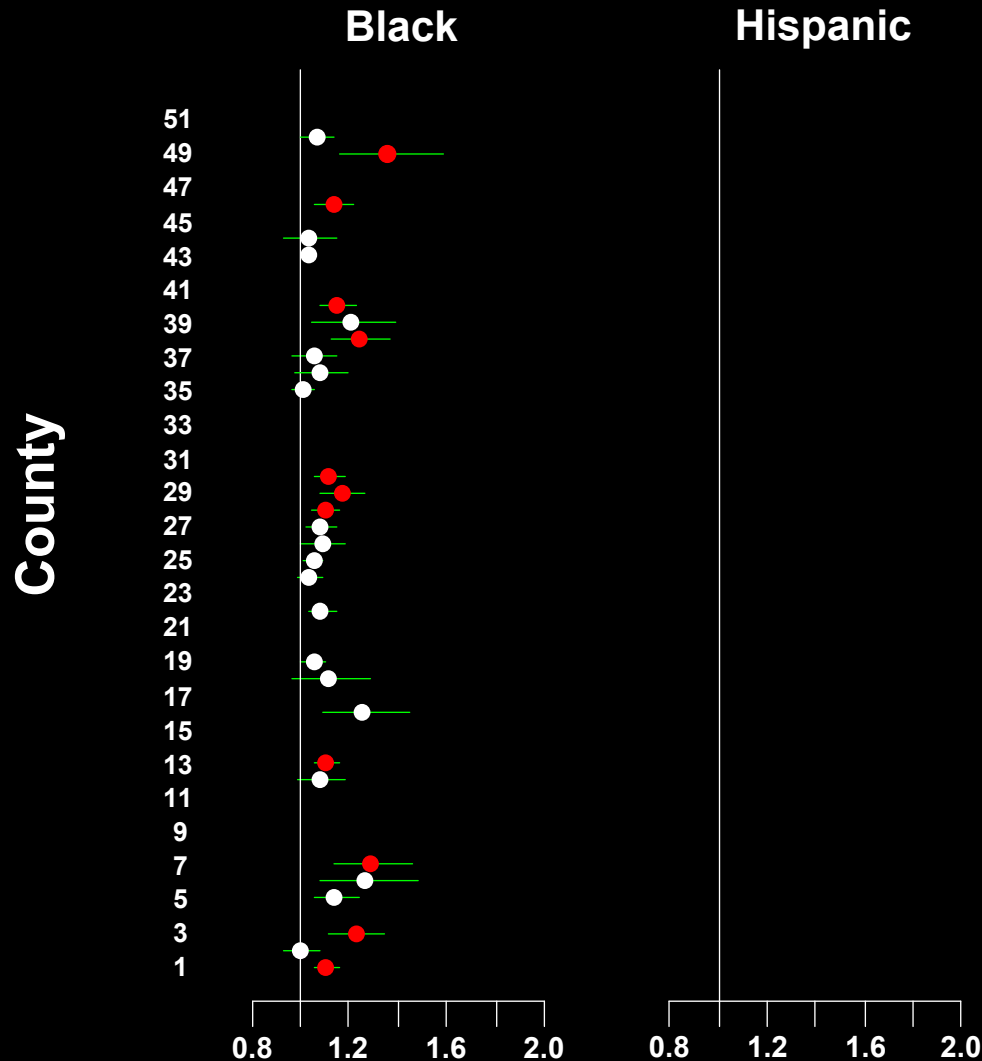
# County 49 Incarcerates Black Defendants at Higher Rates



County 7's Relative Risk = 1.3  
County 49's Relative Risk = 1.4

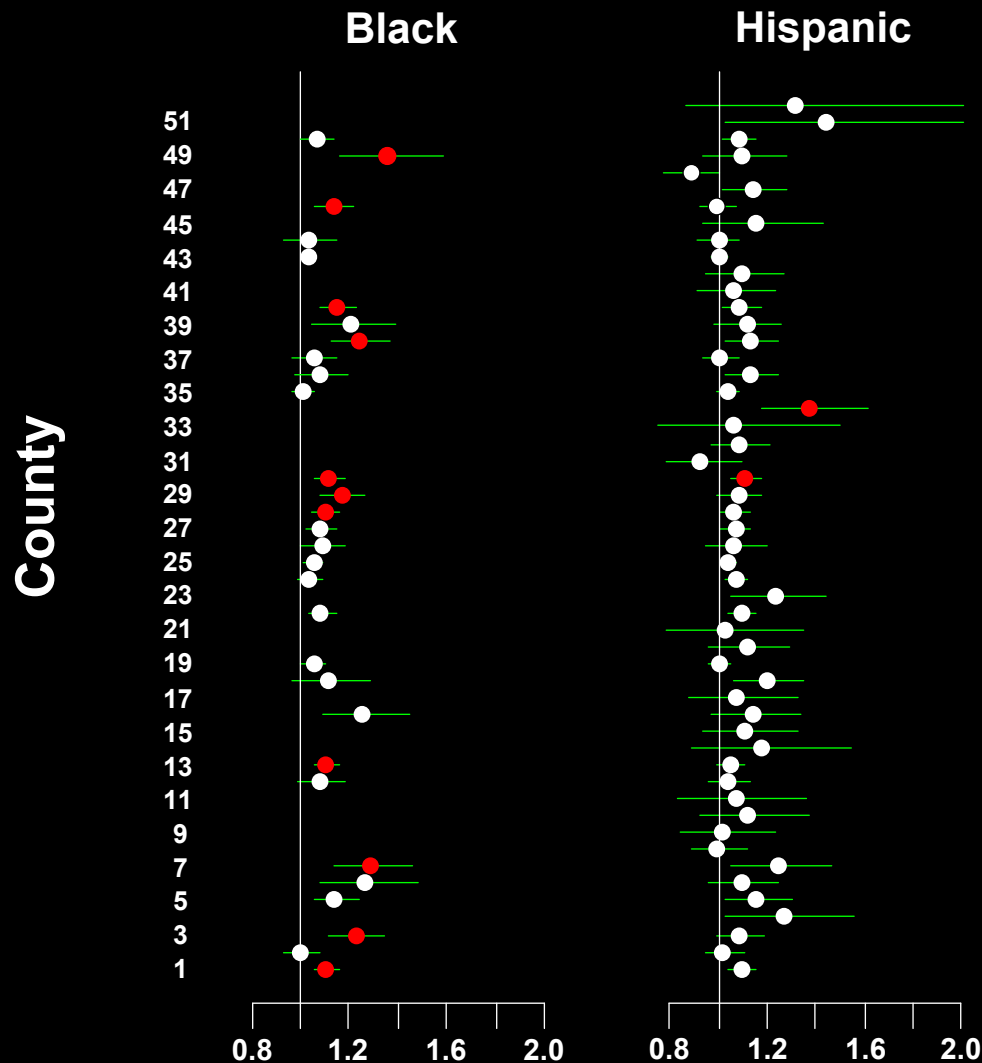


# 11 Counties Had Racial Disparities in Sentencing Black Defendants





# Two Counties Had Racial Disparities in Sentencing Hispanic Defendants



# Scorecard Ultimately Failed, Buckling Under Political Pressure

- Presented to the Commission starting in 2017, masking county identifiers, to get agreement on the approach
- Opposition emerged from some commissioners after identifying outlier counties
- DAs from the counties identified were fiercely opposed to the entire approach
- The Permanent Commission on Sentencing was dissolved during the summer of 2018
- June 2020, Chief Judge DiFiore announced Jeh Johnson would conduct a review of race in the New York courts
  - “...analysis on case outcomes is critically important to identifying the points at which racial disparities exist”

# Four benchmarking applications

- Which officers stop black pedestrians at an unusual rate?
- Which communities are particularly dissatisfied with the police?
- Which counties contribute most to racial disparities in incarceration sentences?
- Which hospitals have...
  - unusually high mortality and readmission rates?
  - excessive opioid prescriptions?

G. Ridgeway, M. Nørgaard, T.B. Rasmussen, W.D. Finkle, L. Pedersen, H.E. Bøtker, and H.T. Sørensen (2019). "Benchmarking Danish Hospitals on Mortality and Readmission Rates After Cardiovascular Admission," *Clinical Epidemiology* 11:67-80

# Compare Performance of 26 Danish Hospitals

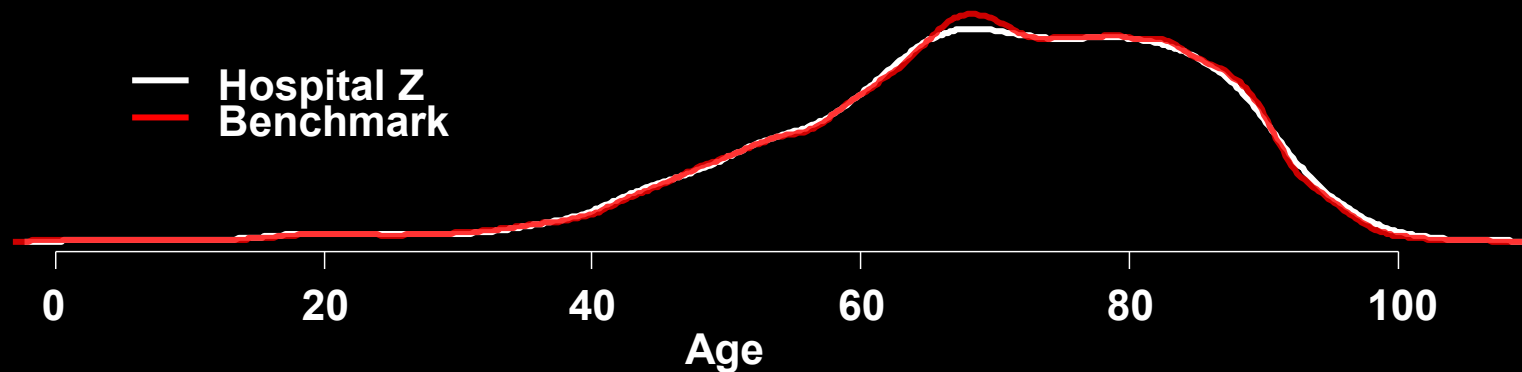
- Data from all Danish hospitals
  - 331,513 patients
  - Danish National Patient Registry and the Danish National Health Service Prescription Database
  - discharged with a primary cardiovascular diagnosis
  - from one of 26 Danish hospitals during 2011-2015
- Main outcome measures
  - 30-day post-admission mortality rates
  - 30-day post-discharge readmission rates
- Patient features
  - age, sex
  - primary discharge diagnosis
  - diagnosis history
  - medications
  - previous cardiac procedures
  - comorbidities

# Benchmark Patients at Other Hospitals Resemble Hospital Z's Patients

	Hospital Z	Benchmark patients	All other hospitals
<b>Age, average</b>	69.9	69.9	<b>68.6</b>
<b>Male, %</b>	55.7	55.2	<b>57.4</b>

# Distributions Match, Not Only Means

	Hospital Z	Benchmark patients	All other hospitals
Age, average	69.9	69.9	68.6
Male, %	55.7	55.2	57.4



# Patients Match on 105 Discharge Diagnoses

	Hospital Z	Benchmark patients	All other hospitals
<b>Myocardial infarction (any)</b>	8.8	8.9	<b>10.5</b>

# Patients Match on 105 Discharge Diagnoses

	Hospital Z	Benchmark patients	All other hospitals
<b>Myocardial infarction (any)</b>	8.8	8.9	<b>10.5</b>
<b>STEMI</b>	0.5	0.5	<b>3.1</b>
<b>Unstable angina</b>	4.2	4.2	<b>2.4</b>
<b>Stable coronary artery disease</b>	15.7	15.7	<b>11.4</b>
<b>Arterial hypertension</b>	8.2	8.2	<b>5.4</b>
<b>Atrial fibrillation or flutter</b>	27.7	27.9	<b>23.8</b>
<b>Ischemic stroke</b>	4.7	4.7	<b>11.4</b>
...			



# Patients Match on 5-year Cardiovascular Diagnosis History

	Hospital Z	Benchmark patients	All other hospitals
<b>Myocardial infarction</b>	9.3	9.1	9.4
<b>Heart Failure</b>	11.2	11.6	16.0
<b>Arterial hypertension</b>	27.8	28.3	33.0
<b>Valvular heart disease</b>	5.0	5.2	8.1
<b>Stroke (any)</b>	7.0	7.1	8.3
...			

# Patients Match on Current Cardiovascular Medication

	Hospital Z	Benchmark patients	All other hospitals
Current use of prescribed cardiovascular medications			
<b>Betablockers</b>	44.3	44.1	<b>40.9</b>
<b>Diuretics</b>	44.2	43.6	<b>37.5</b>

# Patients Match on Procedures

	Hospital Z	Benchmark patients	All other hospitals
Current use of prescribed cardiovascular medications			
<b>Betablockers</b>	44.3	44.1	<b>40.9</b>
<b>Diuretics</b>	44.2	43.6	<b>37.5</b>
Previous cardiac procedures			
<b>Implantable cardiac defibrillator</b>	1.4	1.4	<b>2.0</b>
<b>Aortic valve surgery</b>	1.1	1.0	<b>1.6</b>

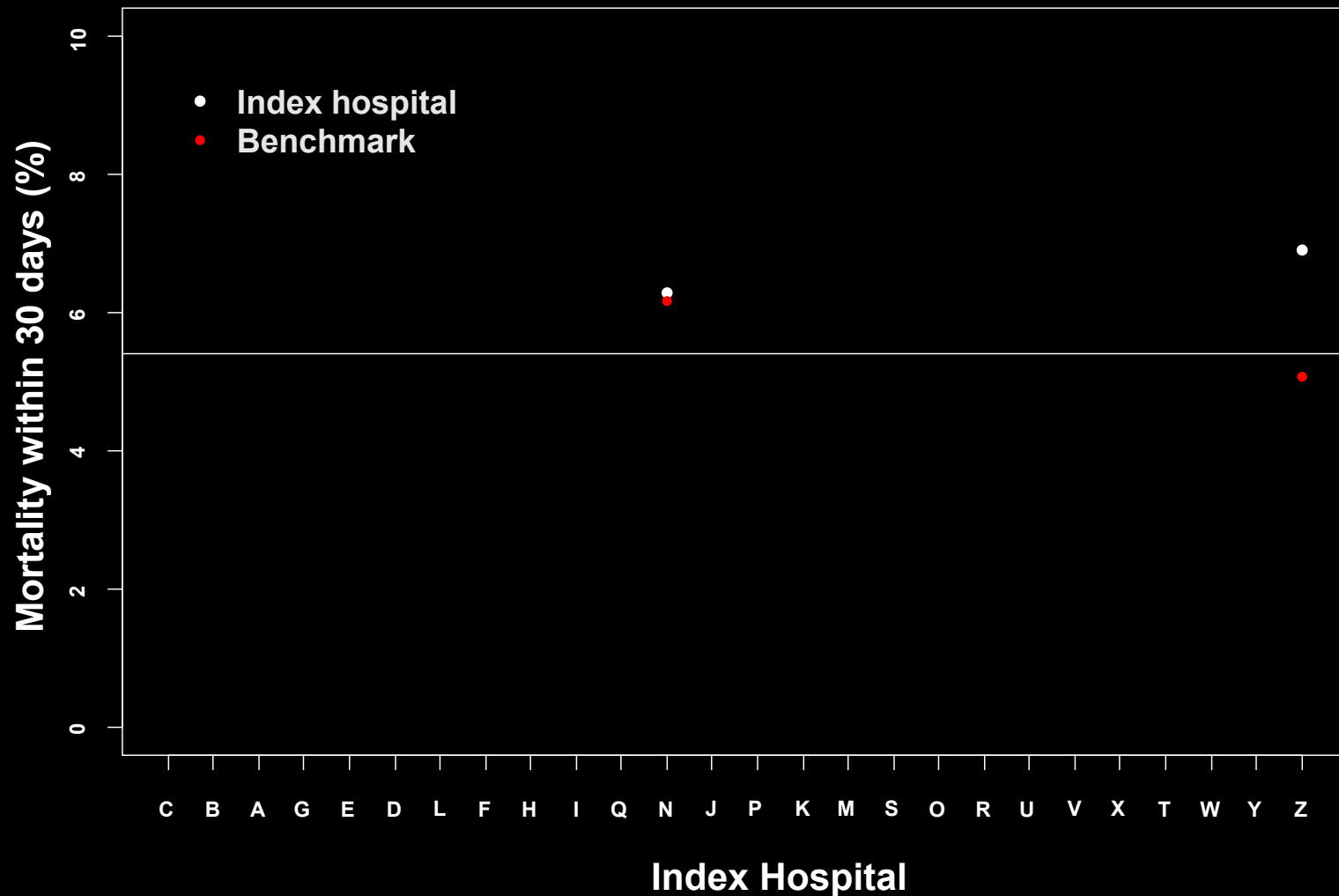
# Patients Match on Comorbidities

	Hospital Z	Benchmark patients	All other hospitals
Current use of prescribed cardiovascular medications			
Betablockers	44.3	44.1	40.9
Diuretics	44.2	43.6	37.5
Previous cardiac procedures			
Implantable cardiac defibrillator	1.4	1.4	2.0
Aortic valve surgery	1.1	1.0	1.6
Selected comorbidity diagnosis history			
Diabetes	11.1	11.4	12.9
Liver disease	0.8	0.8	1.5

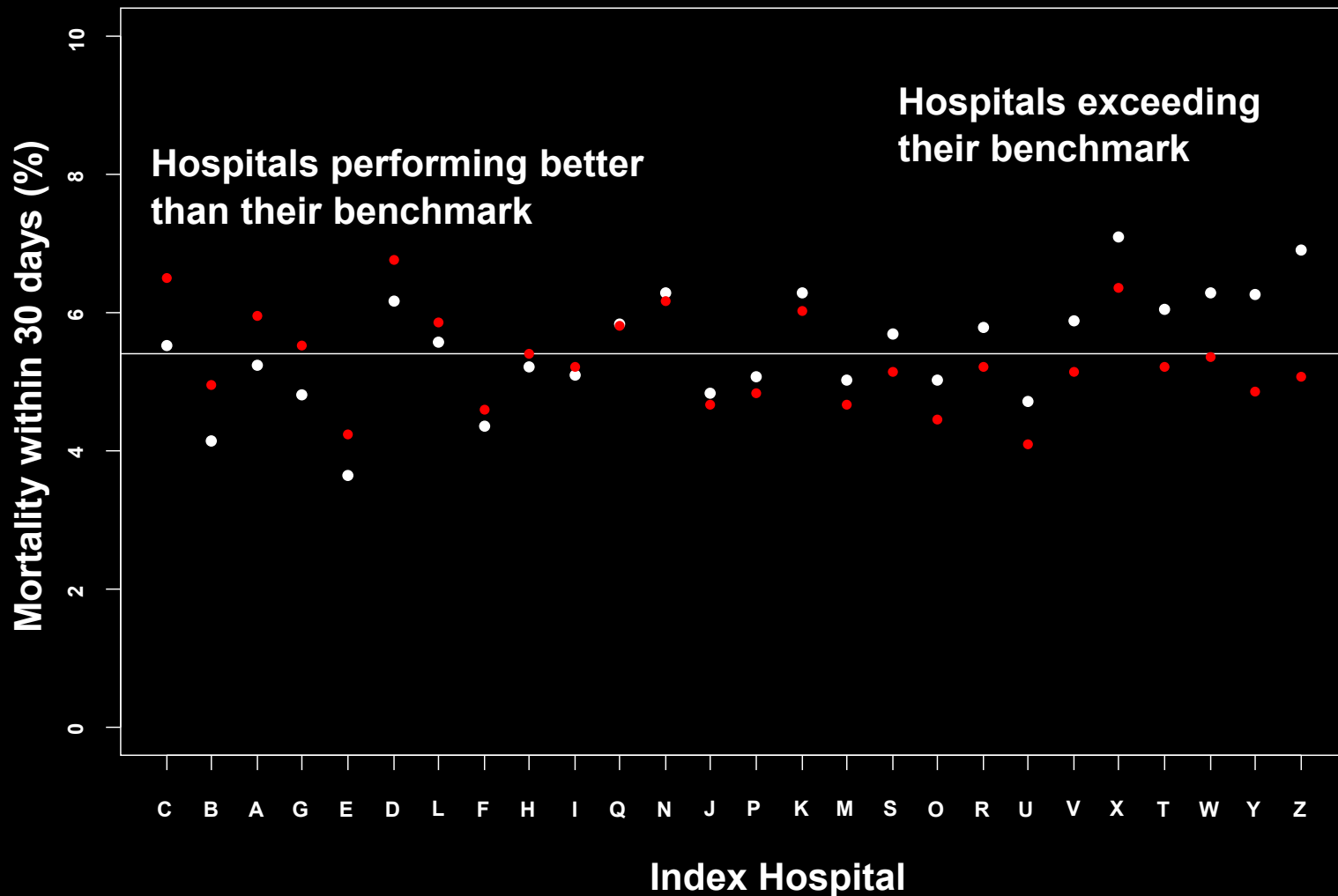
# Patients Match on Other Prescribed Medication

	Hospital Z	Benchmark patients	All other hospitals
Current use of prescribed cardiovascular medications			
<b>Betablockers</b>	44.3	44.1	<b>40.9</b>
<b>Diuretics</b>	44.2	43.6	<b>37.5</b>
Previous cardiac procedures			
<b>Implantable cardiac defibrillator</b>	1.4	1.4	<b>2.0</b>
<b>Aortic valve surgery</b>	1.1	1.0	<b>1.6</b>
Selected comorbidity diagnosis history			
<b>Diabetes</b>	11.1	11.4	<b>12.9</b>
<b>Liver disease</b>	0.8	0.8	<b>1.5</b>
Current use of selected prescribed other medications			
<b>Antidepressants</b>	9.8	9.4	<b>7.9</b>
<b>Antidiabetics</b>	13.5	13.8	<b>14.1</b>

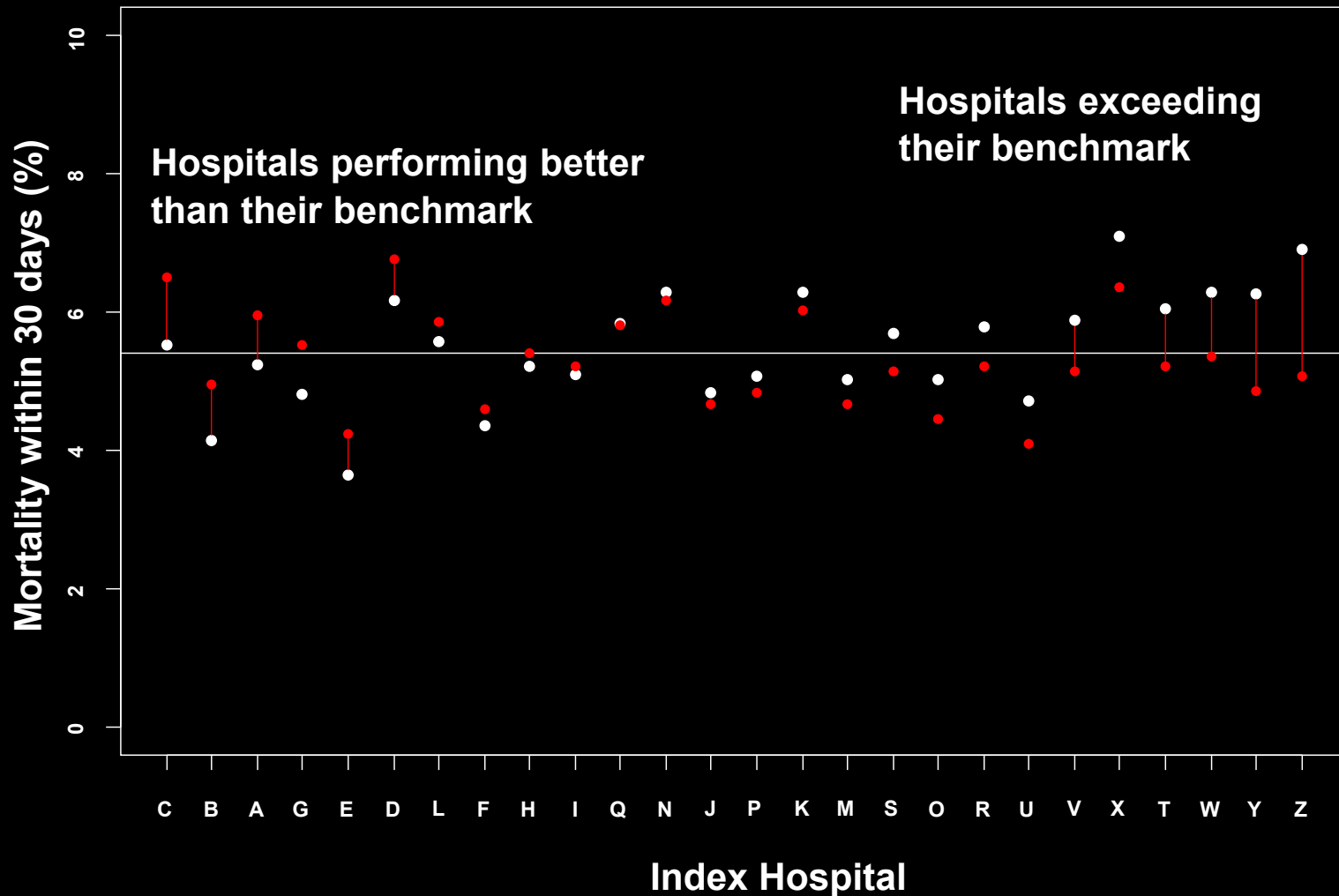
# Compare Every Hospital to Its Customized Benchmark



# Compare Every Hospital to Its Customized Benchmark

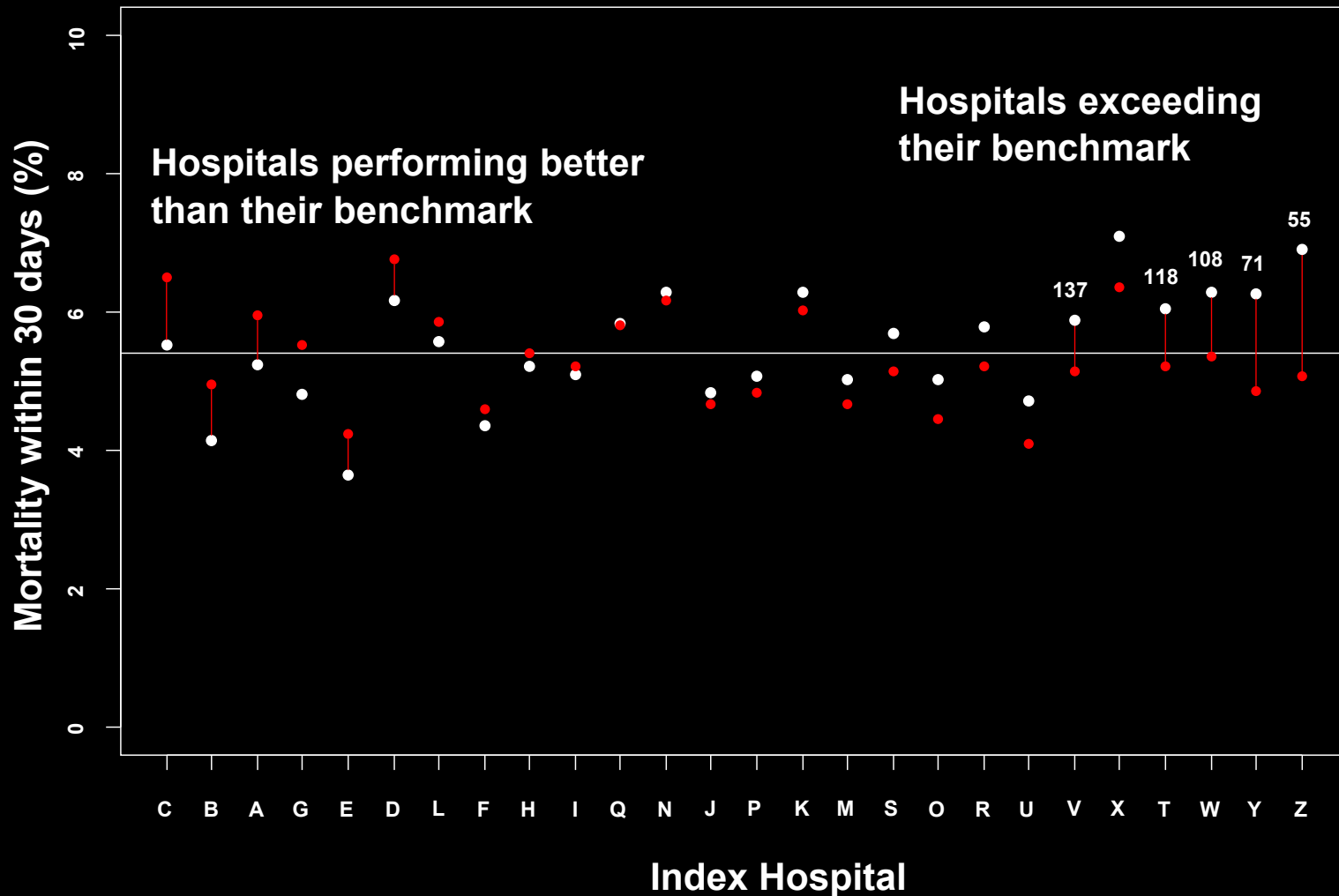


# False Discovery Rate Below 5% for Five Hospitals Exceeding Benchmark

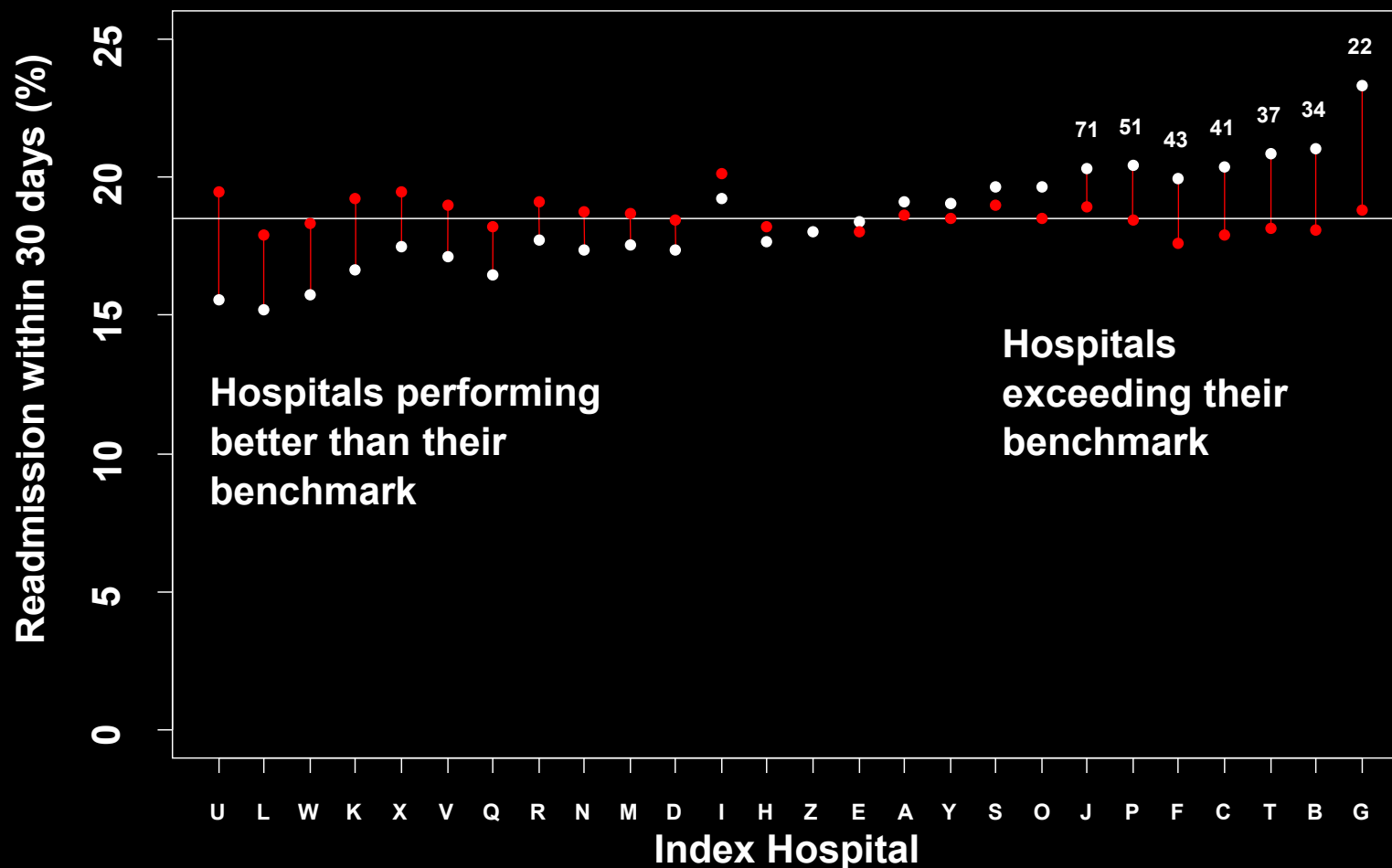




# Number Needed to Harm is Low at Hospital Z



# Hospital T Also Has High 30-day Readmission Rates



# Broad Applicability in Creating Hospital Scorecards

	Hospital X	Benchmark	All Patients
30-day readmission	15.8%	11.7%	7.1%

Consider 287 hospitals

- MarketScan Medicaid Multi-State Database
- Admissions between January 2012-September 2014

# Broad Applicability in Creating Hospital Scorecards

	Hospital X	Benchmark	All Patients
30-day readmission	15.8%	11.7%	7.1%
Oxygen expense (90-day)	\$12.63	\$5.30	\$2.97
Oxygen prescribed (per 100)	9.7	9.7	6.2

# Broad Applicability in Creating Hospital Scorecards

	Hospital X	Benchmark	All Patients
30-day readmission	15.8%	11.7%	7.1%
Oxygen expense (90-day)	\$12.63	\$5.30	\$2.97
Oxygen prescribed (per 100)	9.7	9.7	6.2
Oxycodone supply (30-day)	5.7	5.0	2.5
Oxycodone supply (90-day)	12.3	11.4	5.2
Opiate supply (30-day)	10.1	12.1	6.5
Opiate supply (90-day)	23.5	29.0	14.1
Any opiate prescribed	49.2%	57.0%	42.2%

# Traditional Regression Approach Flags Hospital X on Several Outcomes

	Hospital X	Benchmark	All Patients
30-day readmission	15.8%	11.7%	7.1%
Oxygen expense (90-day)	\$12.63	\$5.30	\$2.97
Oxygen prescribed (per 100)	9.7	9.7	6.2
Oxycodone supply (30-day)	5.7	5.0	2.5
Oxycodone supply (90-day)	12.3	11.4	5.2
Opiate supply (30-day)	10.1	12.1	6.5
Opiate supply (90-day)	23.5	29.0	14.1
Any opiate prescribed	49.2%	57.0%	42.2%

# But the False Discovery Rate is Low Only for Oxygen Expense

	Hospital X	Benchmark	FDR
30-day readmission	15.8%	11.7%	1.00
Oxygen expense (90-day)	\$12.63	\$5.30	0.06
Oxygen prescribed (per 100)	9.7	9.7	1.00
Oxycodone supply (30-day)	5.7	5.0	1.00
Oxycodone supply (90-day)	12.3	11.4	1.00
Opiate supply (30-day)	10.1	12.1	0.39
Opiate supply (90-day)	23.5	29.0	0.49
Any opiate prescribed	49.2%	57.0%	0.27

# Identify Hospitals with Unusual Opioid Prescription Patterns

ID	Hospital	Benchmark	Hospital # Patients	Benchmark # Patients	False Discovery Rate
	Rate of prescription per 100 discharges				
XP	62.1	51.8	642	28,104	0.01
XD	36.6	31.8	4,270	28,744	0.01
XH	63.6	36.1	228	3,827	0.01
ZA	61.4	46.7	526	5,366	0.01
	Days supply 30 days post-discharge				
XD	3.1	2.5	4,270	28,744	0.08
XP	13.4	9.4	642	28,104	0.14
XH	8.5	4.6	228	3,827	0.14



# Broad Applicability of General Approach

- Justice
  - Racial profiling
  - Police performance
  - Sentencing disparities
  - Judicial decision making
- Healthcare
  - Mortality
  - Expenses
  - Prescription practice
- Education?
- Transportation?



# Scorecards, Benchmarking, and the Search for Unusual Hospitals, Communities, and Cops

Greg Ridgeway

Department of Criminology

Department of Statistics and Data Science