# Debiasing Evaluations
# That Are Biased by Evaluations

Jingyan Wang

Carnegie Mellon/Georgia Tech

# Motivation 1: teaching evaluation

- Students are asked to rate instructors' teaching effectiveness
- Correlation between ratings vs. teaching quality can be negative
  [Carrell & West, 2008; Braga et al., 2014; Boring et al., 2016]
- Highly biased by grading leniency:

  *"…the **effects of grades** on teacher–course evaluations are both **substantively and statistically important**…"*

  [Johnson, 2003]
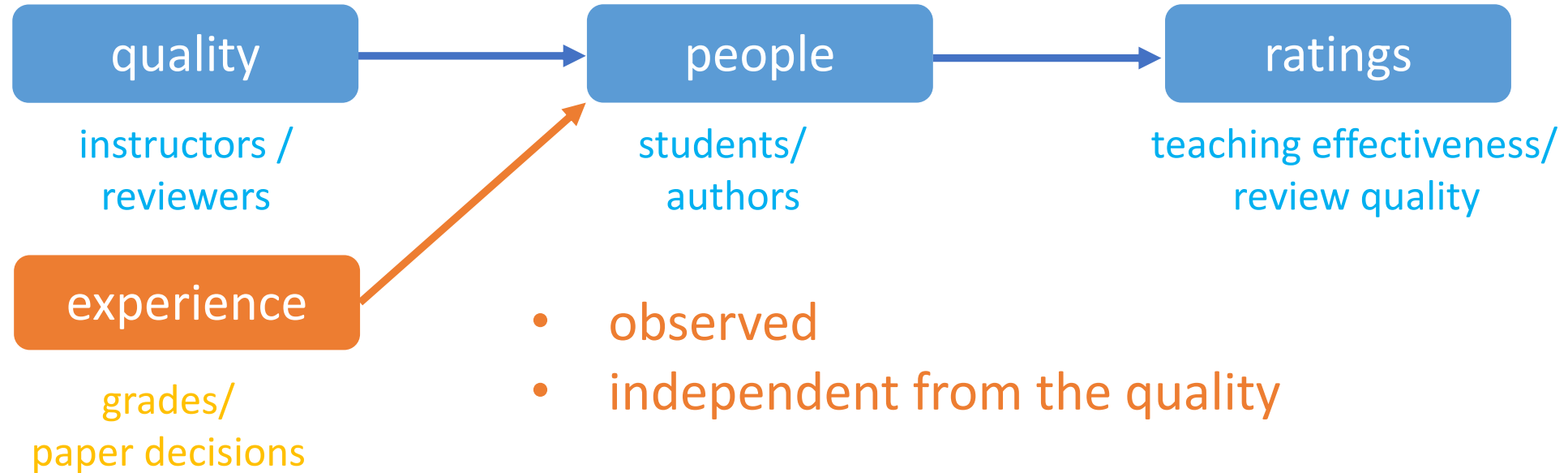
# Motivation 2: peer review

- Authors are asked to rate the reviews they receive
- Highly biased by positiveness of reviews: [Weber et al., 2002; Papagiannaki, 2007; Khosla, 2013]

*"Satisfaction [of the author with the review] had a **strong, positive association with acceptance of the manuscript** for publication... Quality of the review of the manuscript was not associated with author satisfaction."*
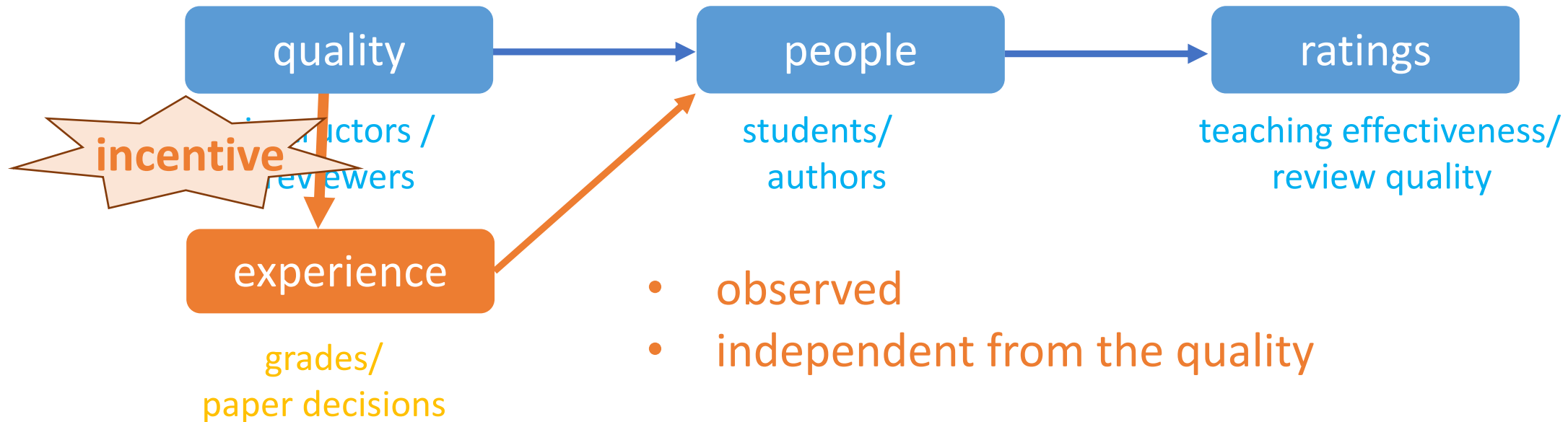
[Weber et al., 2002]

# High-level problem

| quality | → | people | → | ratings |
|---------|---|--------|---|---------|

instructors / reviewers

experience

students/ authors

teaching effectiveness/ review quality

grades/ paper decisions

- observed
- independent from the quality

Unfair for rigorous and strict instructors

**This work: correct experience-induced bias**

# Incentives



**quality** → **people** → **ratings**

instructors / reviewers

students/ authors

teaching effectiveness/ review quality

**incentive**

**experience**

grades/ paper decisions

- observed
- independent from the quality

Introduce incentives for inflating grades, reducing content, "teaching to test" etc.

[Carrell & West, 2008; Braga et al., 2014]

*"... instructors can often **double their odds** of receiving high evaluations from students simply by awarding A's rather than B's or C's."* [Johnson, 2003]

**This work: Correcting experience-induced bias reduces such incentives.** 4

# Problem formulation

- $n$ courses to evaluate: unknown true quality $x_i^*$ for $i \in [n]$
- $d$ students per course
- Student $j \in [d]$ in course $i \in [n]$ gives ratings:

$$y_{ij} = x_i^* + \text{bias} + \text{noise}$$

- **Noise:** iid zero-mean normal
- **Bias:** marginally distributed as normal

<span style="color:blue">The observed experience gives structural information about the bias</span>
- <span style="color:blue">Higher grades → better ratings</span>

# Problem formulation

Example 1:  total ordering of grades

$n = 2,\ \ d = 3$

90　　　85　　　60

Course 1 $(x_1^*)$

Course 2 $(x_2^*)$

95　　　80　　　70

Bias:　　$b_{95} \geq b_{90} \geq b_{85} \geq b_{80} \geq b_{70} \geq b_{60}$

# Problem formulation

Example 2:   partial ordering of grades

$n = 2, \quad d = 6$

A　　　　　　B　　　　　C

Course 1 $(x_1^*)$

Course 2 $(x_2^*)$

Bias:　　　$b_A \geq b_B$　　　$b_B \geq b_C$

Ratings:　$Y = x1^T + B + \text{noise}$

Goal:　　estimate $x^*$ (given $Y$ and ordering)

# Proposed estimator
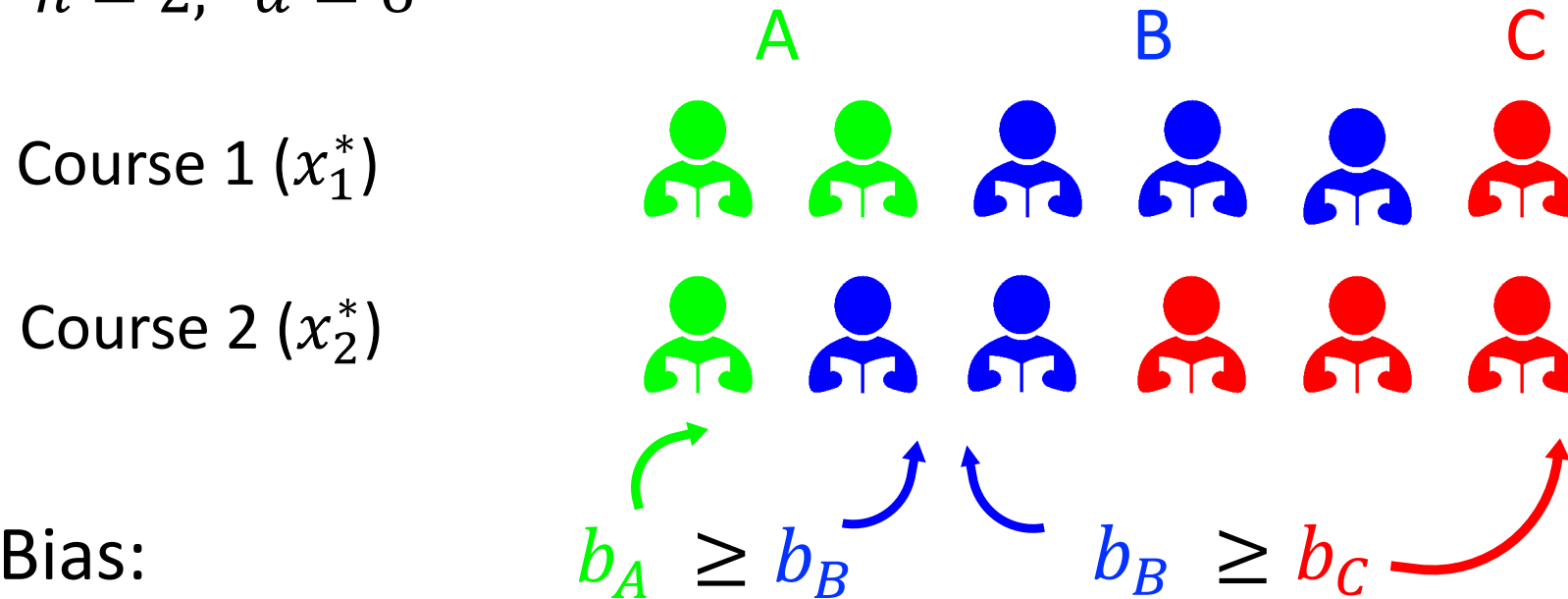
$$\hat{x}^{(\lambda)} \in \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} \quad \underset{B \text{ obeys ordering}}{\min} \quad \|Y - x1^T - B\|_F^2 + \lambda\|B\|_F^2$$

Difference between raw ratings $y$ vs. experience-corrected ratings $x + b$

Regularization on magnitude of $b$

- Analyze two extremal cases: $\lambda = 0$ and $\lambda = \infty$
- Choose $\lambda$ based on the data

# Extremal case 1: $\lambda = 0$

$$\hat{x}^{(\lambda)} \in \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} \quad \underset{\substack{B \text{ obeys} \\ \text{ordering}}}{\min} \quad \|Y - x1^T - B\|_F^2 + \lambda\|\cancel{B}\|_F^2$$

- No regularization means we "explain" the ratings as much as possible by $B$

- Closed-form solution

# Extremal case 1: $\lambda = 0$

$$\hat{x}^{(\lambda)} \in \underset{x \in \mathbb{R}^n}{\mathrm{argmin}} \quad \underset{B \text{ obeys} \atop \text{ordering}}{\min} \quad \|Y - x1^T - B\|_F^2 + \lambda\|B\|_F^2$$

- No regularization means we "explain" the ratings as much as possible by $B$

- Closed-form solution

- Works well when there is no/little noise

**Theorem 1 (informal).** Our estimator (with $\lambda = 0$) is consistent when there is no noise.

- Sample mean is not consistent

# Extremal case 2: $\lambda \to \infty$

$$\hat{x}^{(\lambda)} \in \underset{x \in \mathbb{R}^n}{\mathrm{argmin}} \quad \underset{B \text{ obeys ordering}}{\min} \quad \|Y - x1^T - B\|_F^2 + \lambda\|B\|_F^2$$

- $B \approx 0$
- $\hat{x}^{(\infty)} \approx \underset{x \in \mathbb{R}^n}{\mathrm{argmin}} \|Y - x1^T\|_F^2 = $ taking sample mean
- Formally, define $\hat{x}^{(\infty)} = \underset{\lambda \to \infty}{\lim} \hat{x}^{(\lambda)}$

**Theorem 2.** $\hat{x}^{(\infty)}$ is equivalent to taking the sample mean.

- Our class of estimators includes one of the most commonly-used methods
- Minimax optimal when there is no bias. [Wainwright 2019]

# Choosing $\lambda$

- $\lambda = 0$ and $\lambda = \infty$ work well respectively when there is no noise and no bias.



**Challenge:** don't know the amount of bias vs. noise 😔

**Idea:** carefully design a cross-validation algorithm to choose $\lambda$ 🙂

# Algorithm (sketch)

1. **Split** data to $(Y_{\text{train}}, Y_{\text{val}})$ in a "balanced" way

$Y_{\text{train}}$    $Y_{\text{val}}$

# Algorithm (sketch)

1. **Split** data to $(Y_{\text{train}}, Y_{\text{val}})$ in a "balanced" way
2. **Compute** validation error for each $\lambda$

$Y_{\text{train}}$    $Y_{\text{val}}$

estimator $(\lambda)$

error
$$\left\| Y_{\text{val}} - \hat{x}_{\text{train}} 1^T - \hat{B}_{\text{train}} \right\|_{\text{val}}^2$$

$(\hat{x}_{\text{train}}, \hat{B}_{\text{train}})$

**Challenge:** different bias on different individuals 😔

# Algorithm (sketch)

1. **Split** data to $(Y_{\text{train}}, Y_{\text{val}})$ in a "balanced" way
2. **Compute** validation error for each $\lambda$



$$Y_{\text{train}} \quad Y_{\text{val}}$$

estimator $(\lambda)$

error

$$\left\| Y_{\text{val}} - \hat{x}_{\text{train}} 1^T - \hat{B}_{\text{val}} \right\|_{\text{val}}^2$$

$(\hat{x}_{\text{train}}, \hat{B}_{\text{train}})$    ordering    $\hat{B}_{\text{val}}$

Interpolate $\hat{B}_{\text{val}}$ using $(\hat{B}_{\text{train}}, \text{ordering})$

**Challenge:** different bias on different individuals 😔

**Idea:** interpolate train bias → val bias 🙂

13

# Algorithm (sketch)

$Y_{\text{train}}$     $Y_{\text{val}}$

estimator $(\lambda)$

error

$$\left\| Y_{\text{val}} - \hat{x}_{\text{train}} 1^T - \hat{B}_{\text{val}} \right\|_{\text{val}}^2$$

$(\hat{x}_{\text{train}}, \hat{B}_{\text{train}})$    ordering    $\hat{B}_{\text{val}}$

Interpolate $\hat{B}_{\text{val}}$ using $(\hat{B}_{\text{train}}, \text{ordering})$

$\hat{B}_{\text{train}}$   $-1$        $0$      $2$   $3$

......                                 ...... increasing

$\hat{B}_{\text{val}}$       ?   ?       ?

14

# Algorithm (sketch)

$$Y_{\text{train}} \qquad Y_{\text{val}}$$

estimator $(\lambda)$

error

$$\left\| Y_{\text{val}} - \hat{x}_{\text{train}} 1^T - \hat{B}_{\text{val}} \right\|_{\text{val}}^2$$

$(\hat{x}_{\text{train}}, \hat{B}_{\text{train}})$

ordering

$\hat{B}_{\text{val}}$

Interpolate $\hat{B}_{\text{val}}$ using $(\hat{B}_{\text{train}}, \text{ordering})$

$\hat{B}_{\text{train}}$  $-1$ \qquad 0 \qquad 2 \qquad 3

...... \qquad ...... increasing

$\hat{B}_{\text{val}}$ \quad $-1$ \quad ? \qquad ?

# Algorithm (sketch)

$Y_{\text{train}}$     $Y_{\text{val}}$

estimator $(\lambda)$

error

$$\left\| Y_{\text{val}} - \hat{x}_{\text{train}} 1^T - \widehat{B}_{\text{val}} \right\|_{\text{val}}^2$$

ordering

$(\hat{x}_{\text{train}}, \widehat{B}_{\text{train}})$     $\widehat{B}_{\text{val}}$

Interpolate $\widehat{B}_{\text{val}}$ using $(\widehat{B}_{\text{train}}, \text{ordering})$

$\widehat{B}_{\text{train}}$   $-1$      $0$     $2$   $3$

......      ...... increasing

$\widehat{B}_{\text{val}}$    $-1$   $0$     $?$

14

# Algorithm (sketch)

$$Y_{\text{train}}$$   $$Y_{\text{val}}$$

estimator $(\lambda)$

error

$$\left\| Y_{\text{val}} - \hat{x}_{\text{train}} 1^T - \hat{B}_{\text{val}} \right\|_{\text{val}}^2$$

ordering

$(\hat{x}_{\text{train}}, \hat{B}_{\text{train}})$   $\hat{B}_{\text{val}}$

Interpolate $\hat{B}_{\text{val}}$ using $(\hat{B}_{\text{train}}, \text{ordering})$

$\hat{B}_{\text{train}}$   $-1$   $0$   $2$   $3$

...... increasing

$\hat{B}_{\text{val}}$   $-1$   $0$   $\dfrac{0+2}{2} = 1$

14

# Algorithm (sketch)

$$\|Y_{\text{val}} - \hat{x}_{\text{train}} 1^T - \hat{B}_{\text{val}}\|^2_{\text{val}}$$

error

estimator $(\lambda)$

$(\hat{x}_{\text{train}}, \hat{B}_{\text{train}})$

ordering

$\hat{B}_{\text{val}}$

Interpolate $\hat{B}_{\text{val}}$ using $(\hat{B}_{\text{train}}, \text{ordering})$

$\hat{B}_{\text{train}}$ $-1$     $0$    $2$   $3$

......                          ...... increasing

$\hat{B}_{\text{val}}$        $-1$   $0$     $\dfrac{0+2}{2} = 1$

14

# Theoretical guarantees

**Theorem 3 (informal).** In cases of common partial orderings,
- when there is **no noise**, we have
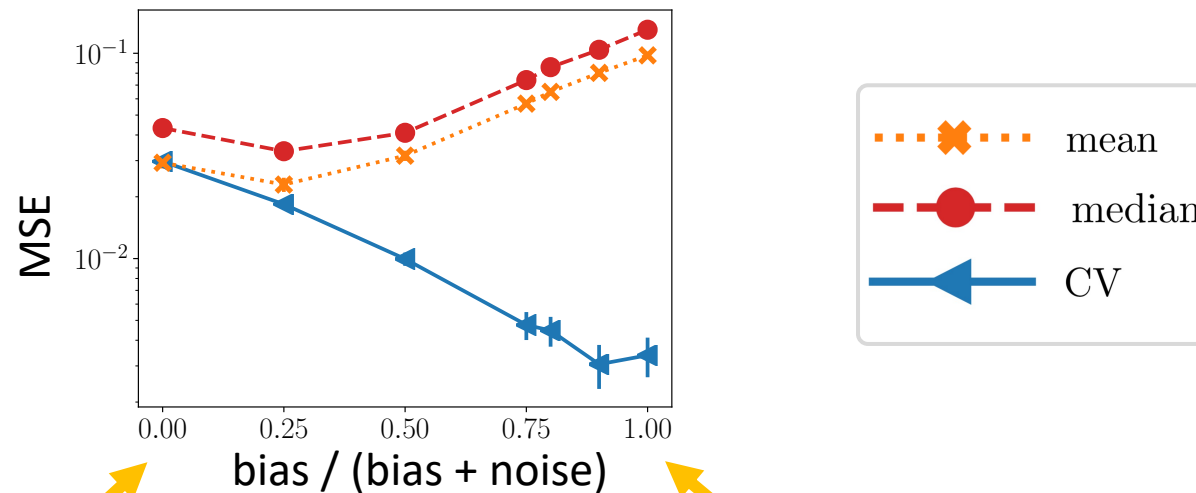$$\hat{x}_{CV} \rightarrow \hat{x}^{(0)};$$

- when there is **no bias**, we have
$$\hat{x}_{CV} \rightarrow \hat{x}^{(\infty)}.$$

Our cross-validation successfully recovers the two extremal cases.

# Experiment

- Indiana University Bloomington

- 10 sessions of a course

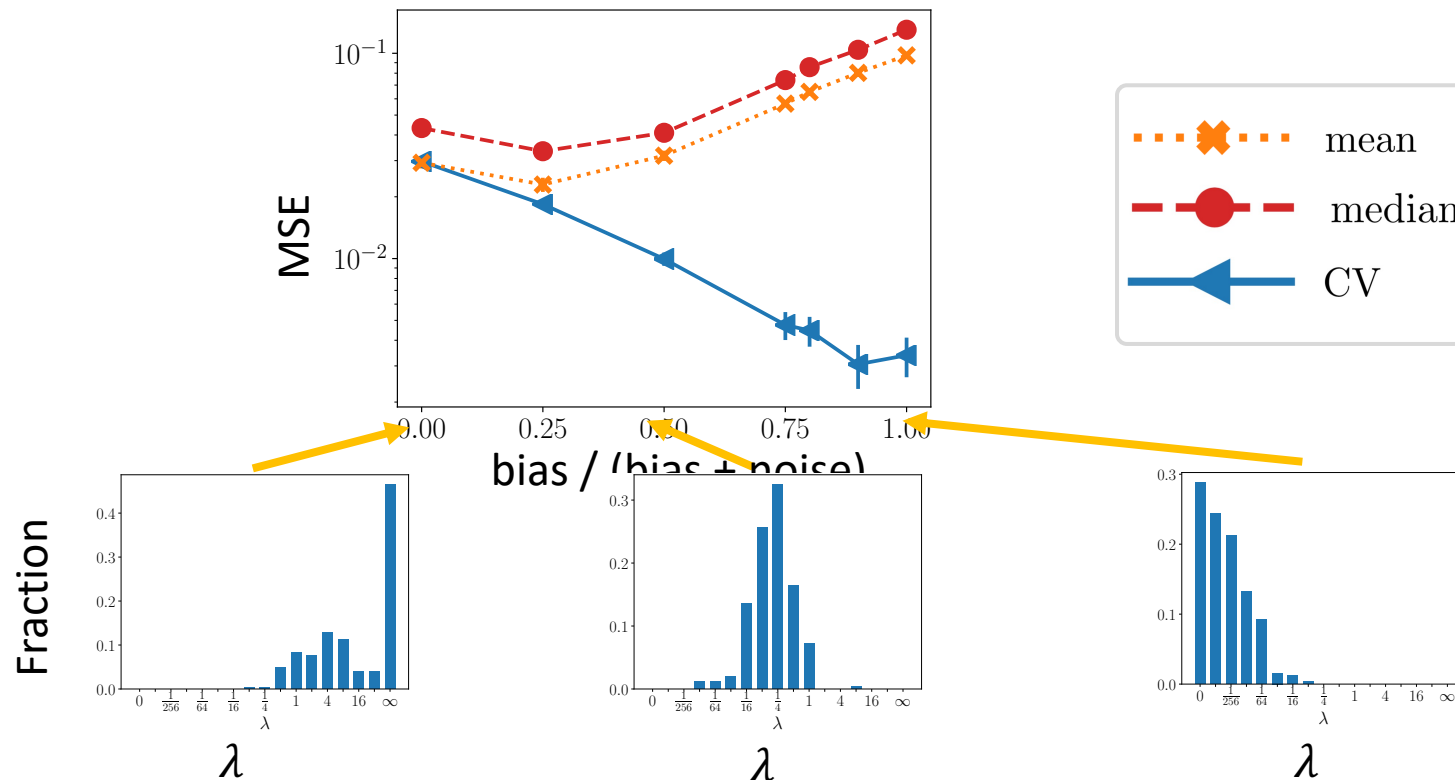- Simulate bias and noise using real grading statistics



**no bias:**

all estimators work well

**lots of bias:**

our estimator significantly better than {mean, median}

16

# Experiment

- Indiana University Bloomington
- 10 sessions of a course
- Simulate bias and noise using real grading statistics

# Take-aways

- Use an ordering constraint to model experience-induced bias, without making restrictive assumptions

- Design a novel CV algorithm to tease out bias vs noise

# Future work

- Sharp statistical bounds on error rates / sample complexity + when there is both bias and noise

- Combining with a game-theoretic approach to design mechanisms

Thanks :)

jingyanw@cmu.edu