# Improving the Peer-Review Process with Model-Based Estimates of Inter-Rater Reliability and Detection of Rating Bias:

From Teacher Selection to Journal Submissions and Grant Applications

Patricia Martinkova<sup>1,2</sup> Dan Goldhaber<sup>3</sup> Elena Erosheva<sup>3</sup> & Carole Lee<sup>3</sup>

<sup>1</sup>Institute of Computer Science, Czech Academy of Sciences <sup>2</sup>College of Education, Charles University, Prague <sup>3</sup>University of Washington, Seattle, USA

PEERE International Conference, 7-9 March 2018, Rome

## Acknowledgements

- Fulbright-Masaryk Fellowship (University of Washington)
- IES grants R305C130030, R305A060018
- Czech Science Foundation grant GJ15-15856Y
- PEERE, TD Cost Action TD1306 (STSM to University of Washington, Nov 2017)
- Data processing support by Malcolm Wolf and Adela Drabinova

## Outline

- Introduction: Peer-Review in Teacher Hiring
- 2 Rating bias
- Model Based Inter-Rater Reliability
- Implications for Other Types of Peer-Review

Onclusion

## Introduction: Teacher Selection Process



Applicants to classroom job openings in Spokane Public Schools during years (2008/09 - 2012/13)

## Introduction: Ratings as Source of Error

#### 54-Pt Screening Rubric:

- Certificate and Education
- Training
- Experience
- Classroom Management
- Flexibility
- Instructional Skills
- Interpersonal Skills
- Cultural Competency
- Preferred Qualifications
- (Quality of Recom. Letters)

DATE:	SCREENER:
Job # / Position Title:	
APPLICANT NAME:	
BATING	1.0
SCREENING CRITERIA 3-6 Story of	vidence to support life as on orns of strongth
3 - 4 Solights 2 - 1 Solights	any neutronics to support that an area of density of the sport
CERTIFICATE AND	
EDUCATION	(has anyones of anote of and), are from our farmer of peopleg: electron
Washington State Certificate Yes / No	
Required Endorsement Yes / No	
Rating (1 - 6) 4	
TRAINING	Look for quality, depth and level of candidates unlikinged entring extension to be position.
Rating (1 - 6) 4	
EXPERIENCE	plane degras to which appendix appoint the production of sectors, for plan discussion of parts: A trappointy combiner could be were highly
Rating (1 - 6) 4	
CLASSROOM MANAGEMENT	Table to particle replacements in manufact biologies. This may not more part and many the planetal and device applicationly double large much an education concorrect only devices prompt, the chapter reserve and prevalues mersures lowering contributes (non-parentice), and required appropriately.
Rating (1 - 6) 4	
FLEXIBILITY	Their insight addression, another, reacting strends, shader, handly at dates, or consulty append with beam one compto and providents, accounted transfer a randor of antipartons, effectively some sample address polys.
Rating (1 - 6) 4	
INSTRUCTIONAL SKILLS	1.200 groups representation in support of the relation was a plane supported, instantic matter or material relation of a analysis approaches, another and adjusts, una calcinally responsive manyor appropriate or app, hadgeoind a senseled downing of materia.
Rating (1 - 6) 4	
INTERPERSONAL SKILLS	Develation and managine effective working relationships with alternar and, students, annexis (perchans, and some
Rating (1 - 6) 4	
CULTURAL COMPETENCY	card for specific references to necessifie destinging for heading and neartaining a relationship with next analysis share family. This may not be explicitly mentioned, but the following neuropics affer some evidence of estimat
individual and caltural differences (nore, religion, second origination, gandar, abilities, oricle-constantic status, etc.) and regular implementation of a trust promoting inclusion.	industrial repetitive language about making and dimities a being that all ability our achieve as high levels, nor of complex remainteners interview powerses, specific interviewing charage of proceedings of the anguate of anisotration which are able reporting control and any powerses and the second of the control and the second of
Rating (1 - 6) 4	
PREFERRED QUALIFICATIONS AS INDICATED ON POSTING	
Rating (1 - 6) 4	
LETTERS OF RECOMMENDATION	(cold for correct inters of recommendation from next the most mean supervised). There were should reduct its quality and recomp of the recommendation as well or the orders of the https://doi.org/10.1006/j.com/pro- certiset/approx/10.1007/j.

ERT STESCREENINGFORM.X.:

### Data structure

- 3474 filled forms
- 1090 applicants
- 137 raters
- 54 job locations (schools)

#### **Applicant status**

- Internal applicant (2322 forms)
  - was previously employed as a teacher in the district or
  - had completed their student teaching in the district
- External applicant (1152 forms)
- 51 applicants external for some and internal for other ratings

### Data structure

- 3474 filled forms
- 1090 applicants
- 137 raters
- 54 job locations (schools)

#### **Applicant status**

- Internal applicant (2322 forms)
  - was previously employed as a teacher in the district or
  - had completed their student teaching in the district
- External applicant (1152 forms)
- 51 applicants external for some and internal for other ratings

1. Introduction 2. Rating bias 3. Model-based Inter-Rater Reliability 4. Implications for peer-review 5. Conclusion

## Ratings of a single applicant



Applicants ranked by averaged total score

1. Introduction 2. Rating bias 3. Model-based Inter-Rater Reliability 4. Implications for peer-review 5. Conclusion

## Ratings of two applicants



Applicants ranked by averaged total score

## Ratings of all applicants



Applicants ranked by averaged total score

## Ratings of all applicants by Internal/External Status



Applicants ranked by averaged total score

# Rating distributions



#### • About 3pt higher ratings for internal applicants

## Rating distributions



- Higher ratings for internal applicants across all subcomponents
- More skewed distribution for internal applicants

## Testing for bias with respect to applicant status

Model controlling for quality measures, accounting for data structure

$$Y_{ijk} = \mu + \omega_i \beta_0 + \mathbf{X}_i \mathbf{\beta} + \mathbf{A}_i + \mathbf{B}_j + \mathbf{S}_k + \mathbf{A} \mathbf{S}_{ik} + \mathbf{e}_{ijk}$$

#### Applicant internal/external status ω<sub>i</sub>

- Applicant quality measures X<sub>i</sub> (e.g. experience, licensure test scores, teacher value added estimates)
- Applicant latent quality  $A_i \sim N(0, \sigma_A^2)$
- Rater severity/leniency  $B_j \sim N(0, \sigma_B^2)$
- School severity/leniency  $S_k \sim N(0, \sigma_S^2)$
- Applicant-school matching effect (interaction)  $AS_{ik} \sim N(0, \sigma_{AS}^2)$

## Testing for bias with respect to applicant status

Model controlling for quality measures, accounting for data structure

$$Y_{ijk} = \mu + \omega_i \beta_0 + \mathbf{X}_i \mathbf{\beta} + \mathbf{A}_i + \mathbf{B}_j + \mathbf{S}_k + \mathbf{A} \mathbf{S}_{ik} + \mathbf{e}_{ijk}$$

- Applicant internal/external status  $\omega_i$
- Applicant quality measures X<sub>i</sub> (e.g. experience, licensure test scores, teacher value added estimates)
- Applicant latent quality A<sub>i</sub> ~ N(0, σ<sup>2</sup><sub>A</sub>)
- Rater severity/leniency  $B_j \sim N(0, \sigma_B^2)$
- School severity/leniency  $S_k \sim N(0, \sigma_S^2)$
- Applicant-school matching effect (interaction)  $AS_{ik} \sim N(0, \sigma_{AS}^2)$

## Testing for bias with respect to applicant status

Model controlling for quality measures, accounting for data structure

$$Y_{ijk} = \mu + \omega_i \beta_0 + \mathbf{X}_i \mathbf{\beta} + \mathbf{A}_i + \mathbf{B}_j + \mathbf{S}_k + \mathbf{A} \mathbf{S}_{ik} + \mathbf{e}_{ijk}$$

- Applicant internal/external status  $\omega_i$
- Applicant quality measures X<sub>i</sub> (e.g. experience, licensure test scores, teacher value added estimates)
- Applicant latent quality  $A_i \sim N(0, \sigma_A^2)$
- Rater severity/leniency  $B_j \sim N(0, \sigma_B^2)$
- School severity/leniency  $S_k \sim N(0, \sigma_S^2)$
- Applicant-school matching effect (interaction)  $AS_{ik} \sim N(0, \sigma_{AS}^2)$

	Model A Internal Only N = 3474	Model B Experience Only N = 3473	Model C WESTB N = 1411	Model D1 VA Math Only N = 303	Model D2 VA Read Only N = 336	Model D Both VA N = 267
Intercept Internal Experience WESTB	36.03 (0.48)*** 3.08 (0.31)***	35.57 (0.50)*** 3.16 (0.31)*** 0.11 (0.03)	36.23 (0.60)*** 2.84 (0.50)***	37.34 (1.32)*** 3.97 (1.29)**	36.96 (1.11)*** 4.15 (1.11)***	36.74 (1.37)*** 4.80 (1.35)***
Writing Reading Math						
Value Added Math Reading						5.62 (2.46)* -3.10 (3.04)

Notes:

- Models include random effects of applicant, rater, school and applicant-school interaction.
- Experience in years
- WESTB scores on state licensure tests.
- Value Added teacher value added estimates based on changes of student perfomrance on achievement tests.

	Model A Internal Only N = 3474	$\begin{array}{l} \mbox{Model B} \\ \mbox{Experience Only} \\ \mbox{N} = 3473 \end{array}$	Model C WESTB N = 1411	Model D1 VA Math Only N = 303	Model D2 VA Read Only N = 336	Model D Both VA N = 267
Intercept Internal Experience WESTB	36.03 (0.48)*** 3.08 (0.31)***	35.57 (0.50)*** 3.16 (0.31)*** 0.11 (0.03)	36.23 (0.60)*** 2.84 (0.50)***	37.34 (1.32)*** 3.97 (1.29)**	36.96 (1.11)*** 4.15 (1.11)***	36.74 (1.37)*** 4.80 (1.35)***
Writing Reading Math						
Math Reading						5.62 (2.46)* -3.10 (3.04)

Notes:

- Models include random effects of applicant, rater, school and applicant-school interaction.
- Experience in years
- WESTB scores on state licensure tests.
- Value Added teacher value added estimates based on changes of student perfomrance on achievement tests.

	Model A Internal Only N = 3474	$\begin{array}{l} \mbox{Model B} \\ \mbox{Experience Only} \\ \mbox{N} = 3473 \end{array}$	$\begin{array}{l} \mbox{Model C} \\ \mbox{WESTB} \\ \mbox{N} = 1411 \end{array}$	Model D1 VA Math Only N = 303	Model D2 VA Read Only N = 336	Model D Both VA N = 267
Intercept Internal Experience WESTB	36.03 (0.48)*** 3.08 (0.31)***	35.57 (0.50)*** 3.16 (0.31)*** 0.11 (0.03)	36.23 (0.60)*** 2.84 (0.50)***	37.34 (1.32)*** 3.97 (1.29)**	36.96 (1.11)*** 4.15 (1.11)***	36.74 (1.37)*** 4.80 (1.35)***
Writing Reading Math			0.11 (0.35) 0.40 (0.33) 0.09 (0.27)			
Math Reading						5.62 (2.46)* -3.10 (3.04)

Notes:

Models include random effects of applicant, rater, school and applicant-school interaction.

Experience in years

WESTB - scores on state licensure tests.

	Model A Internal Only N = 3474	$\begin{array}{l} \mbox{Model B} \\ \mbox{Experience Only} \\ \mbox{N} = 3473 \end{array}$	$\begin{array}{l} Model \ C \\ WESTB \\ N = 1411 \end{array}$	$\begin{array}{l} \mbox{Model D1} \\ \mbox{VA Math Only} \\ \mbox{N} = 303 \end{array}$	Model D2 VA Read Only N = 336	Model D Both VA N = 267
Intercept Internal Experience WESTB Writing Reading Math Value Added	36.03 (0.48)*** 3.08 (0.31)***	35.57 (0.50)*** 3.16 (0.31)*** 0.11 (0.03)	36.23 (0.60)*** 2.84 (0.50)*** 0.11 (0.35) 0.40 (0.33) 0.09 (0.27)	37.34 (1.32)*** 3.97 (1.29)**	36.96 (1.11)*** 4.15 (1.11)***	36.74 (1.37)*** 4.80 (1.35)***
Math Reading				3.90 (2.00)		5.62 (2.46)* -3.10 (3.04)

Notes:

Models include random effects of applicant, rater, school and applicant-school interaction.

Experience in years

WESTB - scores on state licensure tests.

	Model A Internal Only N = 3474	$\begin{array}{l} \mbox{Model B} \\ \mbox{Experience Only} \\ \mbox{N} = 3473 \end{array}$	$\begin{array}{l} \mbox{Model C} \\ \mbox{WESTB} \\ \mbox{N} = 1411 \end{array}$	$\begin{array}{l} \mbox{Model D1} \\ \mbox{VA Math Only} \\ \mbox{N} = 303 \end{array}$	$\begin{array}{l} \mbox{Model D2} \\ \mbox{VA Read Only} \\ \mbox{N} = 336 \end{array}$	Model D Both VA N = 267
Intercept Internal Experience WESTB Writing Reading Math	36.03 (0.48)*** 3.08 (0.31)***	35.57 (0.50)*** 3.16 (0.31)*** 0.11 (0.03)	36.23 (0.60)*** 2.84 (0.50)*** 0.11 (0.35) 0.40 (0.33) 0.09 (0.27)	37.34 (1.32)*** 3.97 (1.29)**	36.96 (1.11)*** 4.15 (1.11)***	36.74 (1.37)*** 4.80 (1.35)***
Math Reading				3.90 (2.00)	3.29 (2.27)	5.62 (2.46)* -3.10 (3.04)

Notes:

Models include random effects of applicant, rater, school and applicant-school interaction.

Experience in years

WESTB - scores on state licensure tests.

	Model A Internal Only N = 3474	$\begin{array}{l} \mbox{Model B} \\ \mbox{Experience Only} \\ \mbox{N} = 3473 \end{array}$	$\begin{array}{l} Model \ C \\ WESTB \\ N = 1411 \end{array}$	$\begin{array}{l} \mbox{Model D1} \\ \mbox{VA Math Only} \\ \mbox{N} = 303 \end{array}$	Model D2 VA Read Only N = 336	Model D Both VA N = 267
Intercept	36.03 (0.48)***	35.57 (0.50)***	36.23 (0.60)***	37.34 (1.32)***	36.96 (1.11)***	36.74 (1.37)***
Internal	3.08 (0.31)***	3.16 (0.31)***	2.84 (0.50)***	3.97 (1.29)**	4.15 (1.11)***	4.80 (1.35)***
Experience		0.11 (0.03)				
WESTB						
Writing			0.11 (0.35)			
Reading			0.40 (0.33)			
Math			0.09 (0.27)			
Value Added						
Math				3.90 (2.00)		5.62 (2.46)*
Reading					3.29 (2.27)	-3.10 (3.04)

Notes:

Models include random effects of applicant, rater, school and applicant-school interaction.

Experience in years

WESTB - scores on state licensure tests.

## Inter-Rater Reliability (Model 1)

$$Y_{ij} = \mu + A_i + B_j + e_{ij}$$

- applicant true quality  $A_i \sim N(0, \sigma_A^2)$ ,
- rater leniency  $B_j \sim N(0, \sigma_B^2)$ ,
- error  $e_{ij} \sim N(0, \sigma_e^2)$

Inter-Rater Reliability:

$$R = \operatorname{cor}(Y_{ij}, Y_{ij'}) = \operatorname{ICC} = \frac{\sigma_A^2}{\sigma_Y^2} = \frac{\sigma_A^2}{\sigma_A^2 + \sigma_B^2 + \sigma_e^2}$$

 $\bullet~\mathrm{R} \in [0,1],$  low values mean a lot of measurement error

## Inter-Rater Reliability (Model 1)

$$Y_{ij} = \mu + A_i + B_j + e_{ij}$$

- applicant true quality  $A_i \sim N(0, \sigma_A^2)$ ,
- rater leniency  $B_j \sim N(0, \sigma_B^2)$ ,
- error  $e_{ij} \sim N(0, \sigma_e^2)$

#### Inter-Rater Reliability:

$$R = \operatorname{cor}(Y_{ij}, Y_{ij'}) = \operatorname{ICC} = \frac{\sigma_A^2}{\sigma_Y^2} = \frac{\sigma_A^2}{\sigma_A^2 + \sigma_B^2 + \sigma_e^2}$$

•  $R \in [0,1]$ , low values mean a lot of measurement error

## Inter-Rater Reliability (Model 1)

$$Y_{ij} = \mu + A_i + B_j + e_{ij}$$

- applicant true quality  $A_i \sim N(0, \sigma_A^2)$ ,
- rater leniency  $B_j \sim N(0, \sigma_B^2)$ ,
- error  $e_{ij} \sim \mathrm{N}(\mathbf{0}, \sigma_e^2)$

#### Inter-Rater Reliability:

$$R = \operatorname{cor}(Y_{ij}, Y_{ij'}) = \operatorname{ICC} = \frac{\sigma_A^2}{\sigma_Y^2} = \frac{\sigma_A^2}{\sigma_A^2 + \sigma_B^2 + \sigma_e^2}$$

•  $\mathrm{R} \in [0,1],$  low values mean a lot of measurement error

## Within-School IRR (Model 2)

$$Y_{ijk} = \mu + A_i + B_j + S_k + AS_{ik} + e_{ijk}$$

- School leniencyl  $S_k \sim N(0, \sigma_S^2)$
- Applicant-school matching effect (interaction)  $AS_{ik} \sim N(0, \sigma_{AS}^2)$

Within-school IRR:

$$R = \operatorname{cor}(Y_{ijk}, Y_{ij'k}) = \frac{\sigma_A^2 + \sigma_S^2 + \sigma_{AS}^2}{\sigma_A^2 + \sigma_B^2 + \sigma_S^2 + \sigma_{AS}^2 + \sigma_e^2}$$

## Within-School IRR (Model 2)

$$Y_{ijk} = \mu + A_i + B_j + \frac{S_k}{S_k} + \frac{AS_{ik}}{S_{ik}} + e_{ijk}$$

- School leniencyl  $S_k \sim N(0, \sigma_S^2)$
- Applicant-school matching effect (interaction)  $AS_{ik} \sim N(0, \sigma_{AS}^2)$

Within-school IRR:

$$R = \operatorname{cor}(Y_{ijk}, Y_{ij'k}) = \frac{\sigma_A^2 + \sigma_S^2 + \sigma_{AS}^2}{\sigma_A^2 + \sigma_B^2 + \sigma_S^2 + \sigma_{AS}^2 + \sigma_e^2}$$

#### • Q: Does IRR differ in ratings of internal vs. external applicants?

- Model 3: Variance components may vary by group
  - e.g. Rater variance may higher when rating external applicants

$$Y_{ijk} = \mu + \omega_i \beta_0 + (1 - \omega_i) A_{0i} + \omega_i A_{1i} + (1 - \omega_i) B_{0j} + \omega_i B_{1j} + (1 - \omega_i) S_{0k} + \omega_i S_{1k} + A S_{ik} + e_{ijk}$$

•  $\omega_i = 1$  for internal and 0 for external applicants •  $A_{0i} \sim N(0, \sigma_{A0}^2)$  and  $A_{1i} \sim N(0, \sigma_{A1}^2)$ •  $B_{0j} \sim N(0, \sigma_{B0}^2)$  and  $B_{1j} \sim N(0, \sigma_{B1}^2)$ •  $S_{0k} \sim N(0, \sigma_{S0}^2)$  and  $S_{1k} \sim N(0, \sigma_{S1}^2)$ 

- Q: Does IRR differ in ratings of internal vs. external applicants?
- Model 3: Variance components may vary by group
  - e.g. Rater variance may higher when rating external applicants

$$\begin{aligned} Y_{ijk} &= \mu + \omega_i \beta_0 + (1 - \omega_i) A_{0i} + \omega_i A_{1i} \\ &+ (1 - \omega_i) B_{0j} + \omega_i B_{1j} \\ &+ (1 - \omega_i) S_{0k} + \omega_i S_{1k} \\ &+ A S_{ik} + e_{ijk} \end{aligned}$$

•  $\omega_i = 1$  for internal and 0 for external applicants •  $A_{0i} \sim N(0, \sigma_{A0}^2)$  and  $A_{1i} \sim N(0, \sigma_{A1}^2)$ •  $B_{0j} \sim N(0, \sigma_{B0}^2)$  and  $B_{1j} \sim N(0, \sigma_{B1}^2)$ •  $S_{0k} \sim N(0, \sigma_{S0}^2)$  and  $S_{1k} \sim N(0, \sigma_{S1}^2)$ 

- Q: Does IRR differ in ratings of internal vs. external applicants?
- Model 3: Variance components may vary by group
  - e.g. Rater variance may higher when rating external applicants

$$\begin{aligned} Y_{ijk} &= \mu + \omega_i \beta_0 + (1 - \omega_i) A_{0i} + \omega_i A_{1i} \\ &+ (1 - \omega_i) B_{0j} + \omega_i B_{1j} \\ &+ (1 - \omega_i) S_{0k} + \omega_i S_{1k} \\ &+ A S_{ik} + e_{ijk} \end{aligned}$$

- $\omega_i = 1$  for internal and 0 for external applicants
- $A_{0i} \sim N(0, \sigma_{A0}^2)$  and  $A_{1i} \sim N(0, \sigma_{A1}^2)$
- $B_{0j} \sim N(0, \sigma_{B0}^2)$  and  $B_{1j} \sim N(0, \sigma_{B1}^2)$
- $S_{0k} \sim N(0, \sigma_{50}^2)$  and  $S_{1k} \sim N(0, \sigma_{51}^2)$

#### Within-school IRR:

• For internal applicant :

$$R_{1} = \operatorname{cor}(Y_{ijk}, Y_{ij'k}) = \frac{\sigma_{A1}^{2} + \sigma_{S1}^{2} + \sigma_{AS}^{2}}{\sigma_{A1}^{2} + \sigma_{B1}^{2} + \sigma_{S1}^{2} + \sigma_{AS}^{2} + \sigma_{e}^{2}}$$

• For external applicant:

$$R_0 = \operatorname{cor}(Y_{ijk}, Y_{ij'k}) = \frac{\sigma_{A0}^2 + \sigma_{50}^2 + \sigma_{A5}^2}{\sigma_{A0}^2 + \sigma_{B0}^2 + \sigma_{50}^2 + \sigma_{A5}^2 + \sigma_{e}^2}$$

#### Within-school IRR:

• For internal applicant :

$$R_{1} = \operatorname{cor}(Y_{ijk}, Y_{ij'k}) = \frac{\sigma_{A1}^{2} + \sigma_{S1}^{2} + \sigma_{AS}^{2}}{\sigma_{A1}^{2} + \sigma_{B1}^{2} + \sigma_{S1}^{2} + \sigma_{AS}^{2} + \sigma_{e}^{2}}$$

• For external applicant:

$$R_0 = \operatorname{cor}(Y_{ijk}, Y_{ij'k}) = \frac{\sigma_{A0}^2 + \sigma_{50}^2 + \sigma_{AS}^2}{\sigma_{A0}^2 + \sigma_{B0}^2 + \sigma_{50}^2 + \sigma_{AS}^2 + \sigma_{e}^2}$$



High applicant-school variability

- Lower applicant variability for external applicants
- Higher rater variability for external applicants
- Lower inter-rater reliability for external applicants



#### • High applicant-school variability

- Lower applicant variability for external applicants
- Higher rater variability for external applicants
- Lower inter-rater reliability for external applicants



- High applicant-school variability
- Lower applicant variability for external applicants
- Higher rater variability for external applicants
- Lower inter-rater reliability for external applicants



- High applicant-school variability
- Lower applicant variability for external applicants
- Higher rater variability for external applicants
- Lower inter-rater reliability for external applicants



- High applicant-school variability
- Lower applicant variability for external applicants
- Higher rater variability for external applicants
- Lower inter-rater reliability for external applicants

1. Introduction 2. Rating bias 3. Model-based Inter-Rater Reliability 4. Implications for peer-review 5. Conclusion

## IRR for Internal and External Applicants (Model 3)

- IRR is estimated simultaneously for both groups within Model 3
- Bootstrapped confidence intervals



• Significant difference in IRR between Internal and External applicants



## Conclusion for Teacher Hiring Data

#### • Rating is school-specific

- Accounting for applicant-school matching in the model is important
- Significantly lower ratings of external applicants confirmed
  - Accounting for previous experience and licensure scores
  - Accounting for subsequent teacher value added
- Singificantly lower inter-rater reliability when rating external applicants
  - Similar variance decomposition in stratified data
  - Our approach allows for testing differences in variance terms and in IRR by group

## Conclusion for Teacher Hiring Data

- Rating is school-specific
  - Accounting for applicant-school matching in the model is important
- Significantly lower ratings of external applicants confirmed
  - Accounting for previous experience and licensure scores
  - Accounting for subsequent *teacher value added*
- Singificantly lower inter-rater reliability when rating external applicants
  - Similar variance decomposition in stratified data
  - Our approach allows for testing differences in variance terms and in IRR by group

## Conclusion for Teacher Hiring Data

- Rating is school-specific
  - Accounting for applicant-school matching in the model is important
- Significantly lower ratings of external applicants confirmed
  - Accounting for previous experience and licensure scores
  - Accounting for subsequent *teacher value added*
- Singificantly lower inter-rater reliability when rating external applicants
  - Similar variance decomposition in stratified data
  - Our approach allows for testing differences in variance terms and in IRR by group

### Implications for Peer-Review in other areas

Model-based IRR is applicable to testing differences w/ respect to:

- assessee status (experienced, matching gender etc.)
  - more likely to matter in fellowships or grant reviews
  - not expected to matter in double-blind review
- other grouping variable
  - reviewer type (experience, research field)
  - journal, journal type
  - grant panel, grant type
  - etc.

### Implications for Peer-Review in other areas

Model-based IRR is applicable to testing differences w/ respect to:

- assessee status (experienced, matching gender etc.)
  - more likely to matter in fellowships or grant reviews
  - not expected to matter in double-blind review
- other grouping variable
  - reviewer type (experience, research field)
  - journal, journal type
  - grant panel, grant type
  - etc.

- Significantly lower ratings and lower IRR showed for external applicants to teacher hiring positions
- Model-based approach allows to
  - account for data structure (applicant-school matching etc.)
  - test for difference in IRR between groups
- Method is aaplicable to grant or journal peer-review

# Thank you for your attention!

http://www.cs.cas.cz/martinkova/

- Significantly lower ratings and lower IRR showed for external applicants to teacher hiring positions
- Model-based approach allows to
  - account for data structure (applicant-school matching etc.)
  - test for difference in IRR between groups
- Method is aaplicable to grant or journal peer-review

# Thank you for your attention!

http://www.cs.cas.cz/martinkova/

- Significantly lower ratings and lower IRR showed for external applicants to teacher hiring positions
- Model-based approach allows to
  - account for data structure (applicant-school matching etc.)
  - test for difference in IRR between groups
- Method is aaplicable to grant or journal peer-review

# Thank you for your attention!

http://www.cs.cas.cz/martinkova/

- Significantly lower ratings and lower IRR showed for external applicants to teacher hiring positions
- Model-based approach allows to
  - account for data structure (applicant-school matching etc.)
  - test for difference in IRR between groups
- Method is aaplicable to grant or journal peer-review

# Thank you for your attention!

http://www.cs.cas.cz/martinkova/