Detecting parameter heterogeneity in psychometric models by means of model-based recursive partitioning with psychotree, stablelearner & co.

### Carolin Strobl

jointly with: Achim Zeileis, Florian Wickelmaier, Julia Kopf, Basil Komboz, Michel Philipp, Thomas Rusch, Kurt Hornik, Lennart Schneider, Rudolf Debelak and Mirka Henninger



Detecting parameter heterogeneity

#### CART

Early statistical problems

#### MOB

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

References

・ロト ・ 四ト ・ ヨト ・ ヨト ・ りへぐ

## Classification and regression trees (CART)



Detecting parameter heterogeneity

## Statistical problems of the early algorithms

Breiman et al. (1984), Quinlan (1986)

variable selection bias

pruning vs. early stopping

e.g.: Hothorn, Hornik and Zeileis (2006, *Journal of Computational and Graphical Statistics*), Strobl, Boulesteix and Augustin (2007, *Computational Statistics & Data Analysis*); tutorial paper: Strobl, Malley and Tutz (2009, *Psychological Methods*)

Detecting parameter heterogeneity

#### CAR

Early statistical problems

#### MOE

for BT models

for Rasch models

for non-Rasch models

Can we trust the results? Stability

Effect size stopping

Summary

References

◆□▶ ◆□▶ ◆□▶ ◆□▶ ○□ のQ@

## Model-based recursive partitioning (MOB)

Zeileis, Hothorn and Hornik (2008, JCGS)



Example: Income as a function of age (ALLBUS 2008) Kopf, Augustin and Strobl (2013) Detecting parameter heterogeneity

#### CART

Early statistical problems

#### MOB

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

References

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

# MOB for psychometric models (package psychotree)

### Bradley-Terry models

Strobl, Wickelmaier and Zeileis (2011, *Journal of Educational and Behavioral Statistics*); Eugster, Leisch and Strobl (2014, *Computational Statistics & Data Analysis*); Müller, Strobl et al. (2015, *Journal of Applied Ecology*)

### binary Rasch models

Strobl, Kopf and Zeileis (2015, Psychometrika)

### polytomous Rasch models

Komboz, Strobl, Zeileis (2016, *Educational and Psychological Measurement*)

### 2PL etc. models

Schneider, Strobl, Zeileis and Debelak (2021, Behavior Research Methods); Debelak and Strobl, (2019, Educational and Psychological Measurement) Detecting parameter heterogeneity

### CART

Early statistical problems

#### MOB

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

# MOB for psychometric models (package psychotree)

### Bradley-Terry models

Strobl, Wickelmaier and Zeileis (2011, *Journal of Educational and Behavioral Statistics*); Eugster, Leisch and Strobl (2014, *Computational Statistics & Data Analysis*); Müller, Strobl et al. (2015, *Journal of Applied Ecology*)

### binary Rasch models

Strobl, Kopf and Zeileis (2015, Psychometrika)

### polytomous Rasch models

Komboz, Strobl, Zeileis (2016, *Educational and Psychological Measurement*)

### 2PL etc. models

Schneider, Strobl, Zeileis and Debelak (2021, Behavior Research Methods); Debelak and Strobl, (2019, Educational and Psychological Measurement) Detecting parameter heterogeneity

### CART

Early statistical problems

#### MOE

#### for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

# Example: MOB for BT models – Topmodels

Detecting parameter heterogeneity

**>** sample: n = 192 (96 female and 96 male) raters between the age of 15 and 77

- covariates: gender, age and
  - (q1) Do you know the women on the photos? Do you know the TV show Germany's Next Topmodel?
  - Did you watch the latest season of Germany's Next (q2) Topmodel regularly? (2007)
  - Have you seen the final of the latest season of (a3) Germany's Next Topmodel? Do you know who won the latest season of Germany's Next Topmodel?

where "yes" to one or more parts = overall "yes"

design: forced choice full paired comparison of photos イロト イロト イヨト イヨト ヨー のくべ



for BT models



# Example: MOB for BT models – Beetles



Detecting parameter heterogeneity

## Müller et al. (2015)

different types of habitats attract different types of beetles





# MOB for psychometric models (package psychotree)

### Bradley-Terry models

Strobl, Wickelmaier and Zeileis (2011, *Journal of Educational and Behavioral Statistics*); Eugster, Leisch and Strobl (2014, *Computational Statistics & Data Analysis*); Müller, Strobl et al. (2015, *Journal of Applied Ecology*)

### binary Rasch models

Strobl, Kopf and Zeileis (2015, Psychometrika)

### polytomous Rasch models

Komboz, Strobl, Zeileis (2016, *Educational and Psychological Measurement*)

### 2PL etc. models

Schneider, Strobl, Zeileis and Debelak (2021, Behavior Research Methods); Debelak and Strobl, (2019, Educational and Psychological Measurement) Detecting parameter heterogeneity

### CART

Early statistical problems

#### MOE

#### for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

# MOB for psychometric models (package psychotree)

### Bradley-Terry models

Strobl, Wickelmaier and Zeileis (2011, *Journal of Educational and Behavioral Statistics*); Eugster, Leisch and Strobl (2014, *Computational Statistics & Data Analysis*); Müller, Strobl et al. (2015, *Journal of Applied Ecology*)

### binary Rasch models

Strobl, Kopf and Zeileis (2015, Psychometrika)

### polytomous Rasch models

Komboz, Strobl, Zeileis (2016, *Educational and Psychological Measurement*)

### 2PL etc. models

Schneider, Strobl, Zeileis and Debelak (2021, Behavior Research Methods); Debelak and Strobl, (2019, Educational and Psychological Measurement) Detecting parameter heterogeneity

### CART

Early statistical problems

#### MOE

for BT models

#### for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

# MOB for binary Rasch models

 $+\,$  identifies previously unknown groups with DIF



Strobl, Kopf and Zeileis (2015, Psychometrika)



Detecting parameter heterogeneity

for Rasch models

# Selection of splitting variables

identify groups of people with different parameters by means of tests for parameter instability (score tests, LM tests):

individual contributions to the score-funktion

$$\psi(y_i, \theta) = \frac{\partial \Psi(y_i, \theta)}{\partial \theta}$$

 $\blacktriangleright$  cumulated over all values of covariate  $\ell$ 





Detecting parameter heterogeneity

#### CART

Early statistical problems

#### MOE

for BT models

#### for Rasch models

for non-Rasch models

#### Can we trust the results?

Stabili

Effect size stopping

Summary

## Selection of cutpoints

maximize partitioned log-likelihood (tree)



Detecting parameter heterogeneity

#### CAR

Early statistical problems

#### MOE

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

References

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● のへで

# Stopping criteria

### $\blacktriangleright$ p value > 0.05 tree

		Node 1	Node 2	Node 3	
Age	Statistic	41.237	48.448	28.924	
	p value	0.171	0.018*	0.593	
Gender	Statistic	41.479	—	—	
	p value	0.006*	—	—	
Motivation	Statistic	112.368	94.680	84.078	
	p value	0.290	0.740	0.432	

minimal node-size < 20</p>

Detecting parameter heterogeneity

#### CAR

Early statistical problems

#### мов

for BT models

#### for Rasch models

for non-Rasch models

Can we trust the

results

Stabilit

Effect size stopping

Summary

References

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

## Properties

as shown by Strobl et al. (2015):

- no alpha inflation
  - correct distributions for optimally selected statistics, closed testing procedure, Bonferroni adjustment
  - even in the presence of ability differences due to CML
- higher power than the LR test to detect DIF
  - in numeric variables (where the true cutpoint is unknown/misspecified)
  - for other non-standard groups (formed, e.g., by interaction or u-shaped patterns)

Detecting parameter heterogeneity

#### CART

Early statistical problems

#### MOE

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary



Detecting parameter heterogeneity

#### CAR

Early statistical problems

#### MOE

for BT models

#### for Rasch models

for non-Rasch models

#### Can we trust the results?

Stabilit

Effect size stopping

Summary

References

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● のへで



Nr. 4: Where is Hessen? (indicate location on a map)

Nr. 5: What is the capital of Rheinland-Pfalz? (Mainz)



#### CART

Early statistical problems

MOE

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stability

Effect size stopping

Summary

References

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○ 三 ○ ○ ○

# MOB for psychometric models (package psychotree)

### Bradley-Terry models

Strobl, Wickelmaier and Zeileis (2011, *Journal of Educational and Behavioral Statistics*); Eugster, Leisch and Strobl (2014, *Computational Statistics & Data Analysis*); Müller, Strobl et al. (2015, *Journal of Applied Ecology*)

### binary Rasch models

Strobl, Kopf and Zeileis (2015, Psychometrika)

### polytomous Rasch models

Komboz, Strobl, Zeileis (2016, *Educational and Psychological Measurement*)

### 2PL etc. models

Schneider, Strobl, Zeileis and Debelak (2021, Behavior Research Methods); Debelak and Strobl, (2019, Educational and Psychological Measurement) Detecting parameter heterogeneity

### CART

Early statistical problems

#### MOE

for BT models

#### for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

# MOB for psychometric models (package psychotree)

### Bradley-Terry models

Strobl, Wickelmaier and Zeileis (2011, *Journal of Educational and Behavioral Statistics*); Eugster, Leisch and Strobl (2014, *Computational Statistics & Data Analysis*); Müller, Strobl et al. (2015, *Journal of Applied Ecology*)

### binary Rasch models

Strobl, Kopf and Zeileis (2015, Psychometrika)

### polytomous Rasch models

Komboz, Strobl, Zeileis (2016, *Educational and Psychological Measurement*)

### 2PL etc. models

Schneider, Strobl, Zeileis and Debelak (2021, Behavior Research Methods); Debelak and Strobl, (2019, Educational and Psychological Measurement) Detecting parameter heterogeneity

### CART

Early statistical problems

#### MOE

for BT models

#### for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

## MOB for polytomous Rasch models

"Effect plots", inspired by Van der Linden and Hambleton (1997) and Fox and Hong (2009)



Latent trait  $\theta$ 

・ロト ・ 四ト ・ ヨト ・ ヨト ・ りへぐ

Detecting

parameter heterogeneity

# Example: MOB for PCM - Verbal Agression

· Komboz, Strobl and Zeileis (2017, Psychological Measurement)



parameter heterogeneity

Detecting

for Rasch models

Can we trust the results?

Stability

Effect size stopping

Summary

# MOB for psychometric models (package psychotree)

### Bradley-Terry models

Strobl, Wickelmaier and Zeileis (2011, *Journal of Educational and Behavioral Statistics*); Eugster, Leisch and Strobl (2014, *Computational Statistics & Data Analysis*); Müller, Strobl et al. (2015, *Journal of Applied Ecology*)

### binary Rasch models

Strobl, Kopf and Zeileis (2015, Psychometrika)

### polytomous Rasch models

Komboz, Strobl, Zeileis (2016, *Educational and Psychological Measurement*)

### 2PL etc. models

Schneider, Strobl, Zeileis and Debelak (2021, Behavior Research Methods); Debelak and Strobl, (2019, Educational and Psychological Measurement) Detecting parameter heterogeneity

### CART

Early statistical problems

#### MOE

for BT models

#### for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

# MOB for psychometric models (package psychotree)

### Bradley-Terry models

Strobl, Wickelmaier and Zeileis (2011, *Journal of Educational and Behavioral Statistics*); Eugster, Leisch and Strobl (2014, *Computational Statistics & Data Analysis*); Müller, Strobl et al. (2015, *Journal of Applied Ecology*)

### binary Rasch models

Strobl, Kopf and Zeileis (2015, Psychometrika)

### polytomous Rasch models

Komboz, Strobl, Zeileis (2016, *Educational and Psychological Measurement*)

### 2PL etc. models

Schneider, Strobl, Zeileis and Debelak (2021, Behavior Research Methods); Debelak and Strobl, (2019, Educational and Psychological Measurement) Detecting parameter heterogeneity

### CART

Early statistical problems

#### MOE

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

### Schneider, Strobl, Zeileis and Debelak (2021, Beh. Res. Meth.)

Model	Response type	Estimation	Model reference	Score-based tests published in
Rasch Model	dichotomous	CML	Rasch ( <u>1960</u> )	Strobl et al., ( <u>2015</u> )
1PL	dichotomous	MML	Birnbaum ( <u>1968</u> )	
2PL	dichotomous	MML	Birnbaum ( <u>1968</u> )	Debelak and Strobl (2019)
3PL	dichotomous	MML	Birnbaum ( <u>1968</u> )	Debelak and Strobl (2019)
3PLu	dichotomous	MML	Barton and Lord (1981)	
4PL	dichotomous	MML	Barton and Lord (1981)	
ideal point model	dichotomous	MML	Maydeu-Olivares et al., (2006)	
rating scale model	polytomous	CML	Andrich (1978)	Komboz et al., ( <u>2018</u> )
partial credit model	polytomous	CML	Masters ( <u>1982</u> )	Komboz et al., (2018)
generalized partial credit model	polytomous	MML	Muraki ( <u>1992</u> )	
graded response model	polytomous	MML	Samejima ( <u>1969</u> )	
nominal response model	polytomous	MML	Bock ( <u>1972</u> )	
generalized graded unfolding model	polytomous	MML	Roberts et al., (2000)	
monotonic polynomial model	polytomous	MML	Falk and Cai (2016)	



Detecting

parameter heterogeneity

#### ◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ○臣○

### Schneider, Strobl, Zeileis and Debelak (2021, Beh. Res. Meth.)

lodel	Response type	Estimation	Model reference	Score-based tests published in
asch Model	dichotomous	CML	Rasch ( <u>1960</u> )	Strobl et al., (2015)
PL	dichotomous	MML	Birnbaum ( <u>1968</u> )	
PL	dichotomous	MML	Birnbaum ( <u>1968</u> )	Debelak and Strobl (2019)
PL	dichotomous	MML	Birnbaum ( <u>1968</u> )	Debelak and Strobl (2019)
PLu	dichotomous	MML	Barton and Lord (1981)	
PL	dichotomous	MML	Barton and Lord (1981)	
leal point model	dichotomous	MML	Maydeu-Olivares et al., (2006)	
ating scale model	polytomous	CML	Andrich (1978)	Komboz et al., ( <u>2018</u> )
artial credit model	polytomous	CML	Masters ( <u>1982</u> )	Komboz et al., ( <u>2018</u> )
eneralized partial credit model	polytomous	MML	Muraki ( <u>1992</u> )	
raded response model	polytomous	MML	Samejima ( <u>1969</u> )	
ominal response model	polytomous	MML	Bock ( <u>1972</u> )	
eneralized graded unfolding model	polytomous	MML	Roberts et al., ( <u>2000</u> )	
nonotonic polynomial model	polytomous	MML	Falk and Cai (2016)	

Detecting parameter heterogeneity

#### ▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへぐ

m

### Schneider, Strobl, Zeileis and Debelak (2021, Beh. Res. Meth.)

lodel	Response type	Estimation	Model reference	Score-based tests published in
asch Model	dichotomous	CML	Rasch ( <u>1960</u> )	Strobl et al., ( <u>2015</u> )
PL	dichotomous	MML	Birnbaum ( <u>1968</u> )	
PL	dichotomous	MML	Birnbaum ( <u>1968</u> )	Debelak and Strobl (2019)
PL	dichotomous	MML	Birnbaum ( <u>1968</u> )	Debelak and Strobl (2019)
PLu	dichotomous	MML	Barton and Lord (1981)	
PL	dichotomous	MML	Barton and Lord (1981)	
leal point model	dichotomous	MML	Maydeu-Olivares et al., (2006)	
ating scale model	polytomous	CML	Andrich ( <u>1978</u> )	Komboz et al., ( <u>2018</u> )
artial credit model	polytomous	CML	Masters (1982)	Komboz et al., ( <u>2018</u> )
eneralized partial credit model	polytomous	MML	Muraki ( <u>1992</u> )	
raded response model	polytomous	MML	Samejima ( <u>1969</u> )	
ominal response model	polytomous	MML	Bock ( <u>1972</u> )	
eneralized graded unrolaing model	polytomous	MML	Roberts et al., ( <u>2000</u> )	
nonotonic polynomial model	polytomous	MML	Falk and Cai (2016)	

Detecting parameter heterogeneity

#### ▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへぐ

### Schneider, Strobl, Zeileis and Debelak (2021, Beh. Res. Meth.)

lodel	Response type	Estimation	Model reference	Score-based tests published in
asch Model	dichotomous	CML	Rasch ( <u>1960</u> )	Strobl et al., (2015)
PL	dichotomous	MML	Birnbaum ( <u>1968</u> )	
PL	dichotomous	MML	Birnbaum ( <u>1968</u> )	Debelak and Strobl (2019)
PL	dichotomous	MML	Birnbaum ( <u>1968</u> )	Debelak and Strobl (2019)
PLu	dichotomous	MML	Barton and Lord (1981)	
PL	dichotomous	MML	Barton and Lord (1981)	
leal point model	dichotomous	MML	Maydeu-Olivares et al., (2006)	
ating scale model	polytomous	CML	Andrich (1978)	Komboz et al., ( <u>2018</u> )
artial credit model	polytomous	CML	Masters (1982)	Komboz et al., ( <u>2018</u> )
eneralized partial credit model	polytomous	MML	Muraki ( <u>1992</u> )	
raded response model	polytomous	MML	Samejima ( <u>1969</u> )	
ominal response model	polytomous	MML	Bock ( <u>1972</u> )	
eneralized graded unrolaing model	polytomous	MML	Roberts et al., ( <u>2000</u> )	
nonotonic polynomial model	polytomous	MML	Falk and Cai (2016)	

 $MML \Rightarrow$  need to model true group differences in means of person parameter distributions (impact)

Detecting parameter heterogeneity



## Can we trust the results?

Detecting parameter heterogeneity

#### CAR

Early statistical problems

#### MOE

for BT models

for Rasch models

for non-Rasch models

# Can we trust the results?

Stabilit

Effect size stopping

Summary

References

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

# MOB for psychometric models (package psychotree)

### Bradley-Terry models

Strobl, Wickelmaier and Zeileis (2011, *Journal of Educational and Behavioral Statistics*); Eugster, Leisch and Strobl (2014, *Computational Statistics & Data Analysis*); Müller, Strobl et al. (2015, *Journal of Applied Ecology*)

### binary Rasch models

Strobl, Kopf and Zeileis (2015, Psychometrika)

### polytomous Rasch models

Komboz, Strobl, Zeileis (2016, *Educational and Psychological Measurement*)

### 2PL etc. models

Schneider, Strobl, Zeileis and Debelak (2021, Behavior Research Methods); Debelak and Strobl, (2019, Educational and Psychological Measurement) Detecting parameter heterogeneity

### CART

Early statistical problems

#### MOE

for BT models

for Rasch models

for non-Rasch models

# Can we trust the results?

Stabilit

Effect size stopping

Summary

# MOB for psychometric models (package psychotree)

### Bradley-Terry models

Strobl, Wickelmaier and Zeileis (2011, *Journal of Educational and Behavioral Statistics*); Eugster, Leisch and Strobl (2014, *Computational Statistics & Data Analysis*); Müller, Strobl et al. (2015, *Journal of Applied Ecology*)

### binary Rasch models

Strobl, Kopf and Zeileis (2015, Psychometrika)

### polytomous Rasch models

Komboz, Strobl, Zeileis (2016, *Educational and Psychological Measurement*)

### 2PL etc. models

Schneider, Strobl, Zeileis and Debelak (2021, Behavior Research Methods); Debelak and Strobl, (2019, Educational and Psychological Measurement) Detecting parameter heterogeneity

### CART

Early statistical problems

#### MOE

for BT models

for Rasch models

for non-Rasch models

# Can we trust the results?

Stabilit

Effect size stopping

Summary

#### Gende p < 0.001 Female, Missing Mak Occupation Age p < 0.001 n < 0.001Stude Pubil 18 in Education. Inominand Student Part Time Student Retired + Activity p = 0.01n < 0.00 Retired Others Not in Employment -<u>10</u>-Age p = 0.02< 24 Can we trust the Node 11 (n = 359) 1.21 Node 12 (n = 397) 1.21 Node 5 (n = 919) Node 8 (n = 311) Node 13 (n = 623) Node 4 (n = 1047) Node 6 (n = 664) 1.21 1.21 1.21 1.21 1.21 12 4 Sunsprary 12 4 6 89 12 4 6 8.9 12 4 6 89 12 4 6 89 12 4 6 8.0 12 4 6 8 9

References

Detecting

parameter heterogeneity

now: Natural Science items, tree for 5000 students

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● のへで

## Can we trust the results?

1. stability:

if we drew another random sample from the data, would the tree still look the same?

Detecting parameter heterogeneity

#### CAR

Early statistical problems

#### MOE

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

References

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへで

#### Gende p < 0.001 Female, Missing Mak Occupation Age p < 0.001 n < 0.001Stude Pubil 18 in Education. Inominand Student Part Time Student Retired + Activity p = 0.01n < 0.00 Retired Others Not in Employment -<u>10</u>-Age p = 0.02< 24 Can we trust the Node 11 (n = 359) 1.21 Node 12 (n = 397) 1.21 Node 5 (n = 919) Node 8 (n = 311) Node 13 (n = 623) Node 4 (n = 1047) Node 6 (n = 664) 1.21 1.21 1.21 1.21 1.21 12 4 Sunsprary 12 4 6 89 12 4 6 8.9 12 4 6 89 12 4 6 89 12 4 6 8.0 12 4 6 8 9

References

Detecting

parameter heterogeneity

now: Natural Science items, tree for 5000 students

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● のへで

Which cutpoints were selected in trees for 125 re-samples?



(variable Age appears more than once in many trees  $\Rightarrow$ more than 125 times in total)

Philipp, Zeileis, Strobl (2016, COMPSTAT Proceedings)

Detecting parameter heterogeneity

#### CART

Early statistical problems

#### MOB

or BT models

for Rasch models

for non-Rasch models

### Can we trust the results?

Stability

Effect size stopping

#### Summary



Which variables were selected in trees for 125 re-samples?



Detecting parameter heterogeneity

Stability

Structural vs. semantic similarity of tree results



# Philipp, Rusch, Hornik and Strobl (2018, JCGS)

Detecting parameter heterogeneity

CART

Early statistical problems

MOB

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stability

Effect size stopping

Summary

References

ъ

Structural vs. semantic similarity of tree results



Philipp, Rusch, Hornik and Strobl (2018, JCGS)

Detecting parameter heterogeneity

#### CART

Early statistical problems

#### MOB

for BT models

for Rasch models

for non-Rasch models

### Can we trust the results?

Stability

Effect size stopping

Summary

References

э

Structural vs. semantic similarity of tree results



Detecting parameter heterogeneity

#### CART

Early statistical problems

#### MOB

for BT models

for Rasch models

for non-Rasch models

### Can we trust the results?

Stability

Effect size stopping

Summary

References

 $x_1$ 

Are trees always instable?

Detecting parameter heterogeneity

#### CAR

Early statistical problems

#### MOE

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stability

Effect size stopping

Summary

References

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = ● ● ●

Are trees always instable?



given enough observations, tree shows high stability for step-function (orange results), but not for smooth function (blue results)

Detecting parameter heterogeneity

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへで

Are trees always instable?



given enough observations, tree shows high stability for step-function (orange results), but not for smooth function (blue results)

### stability = f(learner $\times$ dgp)

▲□▶ ▲□▶ ▲ 臣▶ ▲ 臣▶ 三臣 - のへで

Detecting

parameter heterogeneity

## Can we trust the results?

1. stability:

if we drew another random sample from the data, would the tree still look the same?

Detecting parameter heterogeneity

#### CAR

Early statistical problems

#### MOE

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

#### Stability

Effect size stopping

Summary

References

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへで

## Can we trust the results?

1. stability:

if we drew another random sample from the data, would the tree still look the same?

2. power vs. effect size:

the group differences are significant – but are they practically relevant?

Detecting parameter heterogeneity

#### CAR

Early statistical problems

#### MOE

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stability

Effect size stopping

Summary

References

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ のQ@

#### Gende p < 0.001 Female, Missing Mak Occupation Age p < 0.001 n < 0.001Stude Pubil 18 in Education. Unemployed. Student Part Time Student Retired + Activity p = 0.01n < 0.001 Retired Others Not in Employment -<u>10</u>-Age p = 0.02< 24 Node 11 (n = 359) 1.21 Node 12 (n = 397) Node 5 (n = 919) Node 8 (n = 311) Node 13 (n = 623) Node 4 (n = 1047) Node 6 (n = 664) 1.21 1.21 1.21 1.21 n 1.21 stopping 12 4 Sunsprary 12 4 6 89 12 4 6 8.9 12 4 6 89 12 4 6 89 12 4 6 8.9 12 4 6 8 9

Natural Science items, tree for 5000 students

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

Detecting

parameter heterogeneity



Natural Science items, tree for 500 students

Detecting parameter heterogeneity

#### CART

Early statistical problems

#### MOE

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

References

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ● のへで

Henninger, Debelak and Strobl (2022, Edu. and Psy. Measurement)

additional stopping criterion: Mantel-Haenszel odds ratio = effect size measure for DIF

Educational Testing Service (ETS) classification criteria:

- A = negligible DIF
- $\mathsf{B} = \mathsf{medium} \ \mathsf{DIF}$
- $\mathsf{C} = \mathsf{large} \; \mathsf{DIF}$





MOB

or BT models

or Rasch models

for non-Rasch models

Can we trust the results?

Stability

Effect size stopping

Summary

References

◆□▶ ◆□▶ ◆□▶ ◆□▶ ○□ のQ@

### tree stopped with Mantel-Haenszel criterion



Detecting

parameter heterogeneity

### tree stopped with Mantel-Haenszel criterion



Detecting parameter heterogeneity

- What I have talked about: (package: psychotree) model-based recursive partitioning can
  - detect groups of observations with different parameters
  - need not be specified a priori (data-driven, exploratory)

Detecting parameter heterogeneity

#### CART

Early statistical problems

#### MOB

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

References

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

- What I have talked about: (package: psychotree) model-based recursive partitioning can
  - detect groups of observations with different parameters
  - need not be specified a priori (data-driven, exploratory)
  - can be used to detect DIF and DSF

#### CART

Early statistical problems

#### MOB

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

References

▲□▶ ▲圖▶ ▲匡▶ ▲匡▶ ― 匡 - のへで

- What I have talked about: (package: psychotree) model-based recursive partitioning can
  - detect groups of observations with different parameters
  - need not be specified a priori (data-driven, exploratory)
  - can be used to detect DIF and DSF
  - can we trust the results?
    - stability assessment (package: stablelearner)
    - effect-size-based stopping (currently on github.com/mirka-henninger)

Detecting parameter heterogeneity

#### CART

Early statistical problems

#### MOB

or BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

References

◆□▶ ◆□▶ ◆□▶ ◆□▶ ○□ のQ@

- What I have talked about: (package: psychotree) model-based recursive partitioning can
  - detect groups of observations with different parameters
  - need not be specified a priori (data-driven, exploratory)
  - can be used to detect DIF and DSF
  - can we trust the results?
    - stability assessment (package: stablelearner)
    - effect-size-based stopping (currently on github.com/mirka-henninger)
- What I haven't talked about:
  - anchoring: making parameters comparable between groups (= end nodes, possibly many) in the presence of DIF and DSF

Detecting parameter heterogeneity

#### CART

Early statistical problems

#### MOB

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

## Key References

Henninger M, Debelak R, Strobl C (2023). A New Stopping Criterion for Rasch Trees Based on the Mantel–Haenszel Effect Size Measure for Differential Item Functioning. *Educational and Psychological Measurement*, **83**(1), 181–212.

Strobl C, Kopf J, Zeileis A (2015). Rasch Trees: A New Method for Detecting Differential Item Functioning in the Rasch Model. *Psychometrika*, **80**(2), 289–316.

Strobl C, Malley J, Tutz G (2009). An Introduction to Recursive Partitioning: Rationale, Application and Characteristics of Classification and Regression Trees, Bagging and Random Forests. *Psychological Methods*, **14**(4), 323–348.

Philipp M, Rusch T, Hornik K, Strobl C (2018). Measuring the Stability of Results From Supervised Statistical Learning. *Journal of Computational and Graphical Statistics*, **27**(4), 685–700. Detecting parameter heterogeneity

#### CART

Early statistical problems

#### MOE

for BT models

or Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

# Stability of cutpoints



Detecting parameter heterogeneity

#### CART

Early statistical problems

#### MOE

for BT models

for Rasch models

for non-Rasch models

Can we trust the results?

Stabilit

Effect size stopping

Summary

References

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへの

Which cutpoints were selected in trees for 125 re-samples?



Detecting parameter heterogeneity

・ロト ・ 同ト ・ ヨト ・ ヨー・ つへぐ

### tree stopped with Mantel-Haenszel criterion



Detecting parameter heterogeneity