

Lesson 9: Differential Item Functioning

Patrícia Martinková

Department of Statistical Modelling
Institute of Computer Science, Czech Academy of Sciences

Institute for Research and Development of Education
Faculty of Education, Charles University, Prague

NMST570, December 4, 2018

Outline

- 1 Introduction
- 2 DIF and fairness
- 3 DIF detection methods
- 4 Further Topics
- 5 Conclusion

Review - IRT models

- Item Characteristic Curve (ICC)
 - Item Response Function (IRF)
 - Item Information Function (IIF)
 - Test Information Function (TIF)
 - Likelihood function
 - Parameter estimation: JML, CML, MML, Bayesian approaches
 - Model fit, item fit, person fit
-
- 1PL, 2PL, 3PL, 4PL IRT models
 - Graded Response Model (GRM)
 - Partial Credit Model (PCM)
 - Generalized Partial Credit Model (GPCM)
 - Rating Scale Model (RSM)
 - Nominal Response Model (NRM)

Motivation for differential item functioning (DIF) analysis

Complex validation of Homeostasis Concept Inventory (HCI)

- Males / English as a first language / White and Asian students performed better

Is the test fair?

McFarland et al. Development and Validation of the Homeostasis Concept Inventory. *CBE Life Sciences Education*, vol. 16 no. 2 ar35, 2017. doi [10.1187/cbe.16-10-0305](https://doi.org/10.1187/cbe.16-10-0305)

Motivation: Development and Validation of HCI

Differential Item Functioning (DIF) Analysis

- Analytical method to address item fairness
- Ubiquitous in large-scale assessments development
- Less used in conceptual assessment development
- None of the HCI items exhibited DIF
 - with respect to gender, ethnicity or ELL status

Methods paper: Importance of DIF Analysis

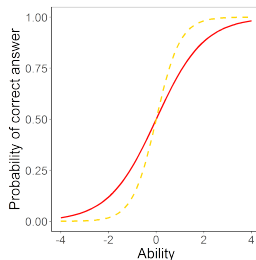
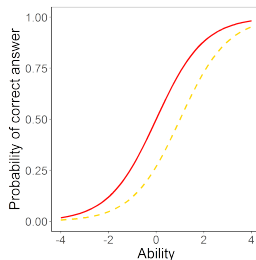
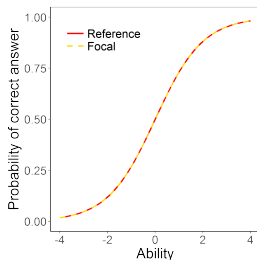
Martinková et al. Checking Equity: Why DIF Analysis should be a Routine Part of Developing Conceptual Assessments. *CBE Life Sciences Education*, 16(2), rm2.
doi [10.1187/cbe.16-10-0307](https://doi.org/10.1187/cbe.16-10-0307)

Differential Item Functioning

Differential Item Functioning (DIF)

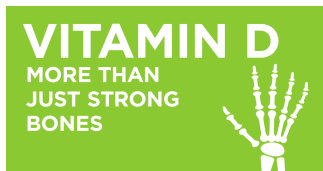
Two subjects with the same underlying ability but from different groups have different probability to answer question correctly

- Two groups referred to as reference and focal (usually minority)
- Two types of DIF - uniform and non-uniform



Example of DIF item

Childhood illnesses (Drabinová & Martinková, 2017)



Deficiency of vitamin D in childhood could cause

- a. rickets
- b. scurvy
- c. dwarfism
- d. mental retardation

Example of DIF item

Tipping example (Martiniello et al., 2012)

Of the following, which is the closest approximation of a 15 percent tip on a restaurant check of \$24.99?

- a. \$2.50
- b. \$3.00
- c. \$3.75
- d. \$4.50

Example of DIF items

- Example: Spelling test (orally administered)
 - spell word girder
- Example (SAT): Runner is to marathon as
 - a. envoy is to embassy
 - b. martyr is to massacre
 - c. oarsman is to regatta
 - d. referee is to tournament
 - e. horse is to stable

Who might have been dissadvantaged?

Terminology: Reference group (R), Focal group (F)

DIF as multidimensionality problem

DIF as multidimensionality problem:

- Existence of another dimension tested on the particular item besides the primary latent variable

What is the primary and the secondary latent variable tested in mentioned examples?

DIF and item fairness

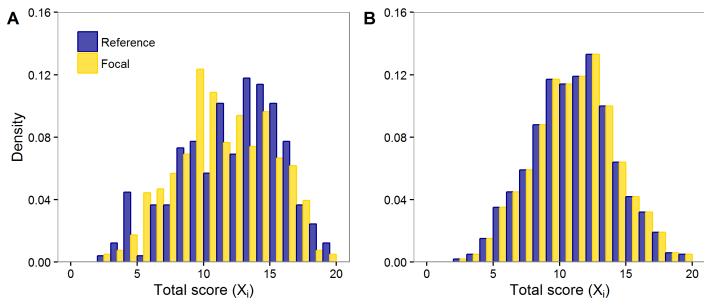
DIF items are **potentially** unfair

- Content experts must decide on item fairness
- Secondary latent trait causing DIF
 - Unrelated to content being tested
 - DIF item is considered unfair
 - Item should be reworded or removed
 - Example: Tipping
 - Related to content being tested
 - DIF item is not considered unfair
 - Item can inform teaching
 - Example: Item on childhood illnesses
as part of Czech Medical School Admission Test in Biology

DIF vs. Difference in total scores

Comparing total scores only can lead to incorrect conclusions about item/test fairness:

- Case study 1: Homeostasis Concept Inventory
 - Significant difference between males and females in total score (Fig A)
- Case study 2: Simulated dataset based on GMAT
 - Identical distributions of total score (Fig B)

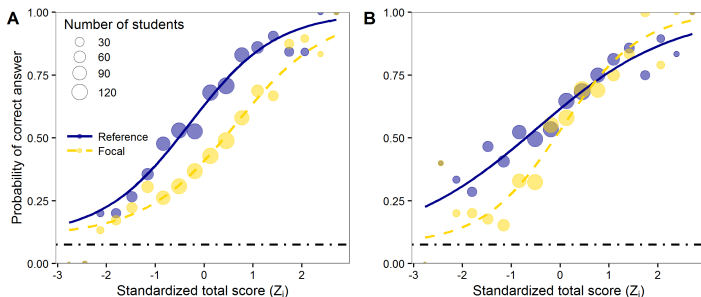


Martinková et al. (2017)

DIF vs. Difference in total scores (cont.)

Comparing total scores only can lead to incorrect conclusions about item/test fairness:

- Case study 1: No HCI item detected as DIF
- Case study 2: DIF detected in two items of simulated dataset
 - Item 1 exhibits uniform DIF (Fig A)
 - Item 2 exhibits non-uniform DIF (Fig B)



Martinková et al. (2017)

DIF detection methods

- Based on **total score**

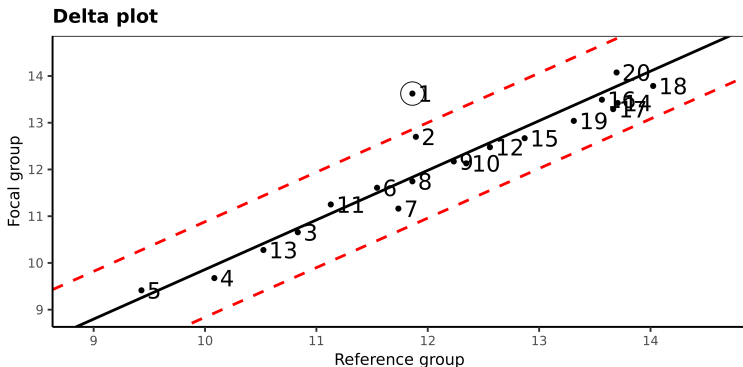
- Mantel-Haenszel test
 - + simple, easily implemented
 - cannot detect non-uniform DIF
 - doesn't account for possibility of guessing/inattention
- Logistic regression
 - + simple, easily implemented, detects both forms of DIF
 - doesn't account for possibility of guessing/inattention

- Based on **latent ability**

- Item Response Theory models (non-linear mixed effect models)
 - + detects both forms of DIF, accounts for possibility of guessing/inattention
 - more complex, computationally demanding

Delta plot

- Angoff & Ford (1973)
- compares proportions of correct answers
- displays non-linear transformation of proportions (using quantiles)
- detection threshold
 - fixed to 1.5
 - normal approximation (Magis & Facon, 2012).



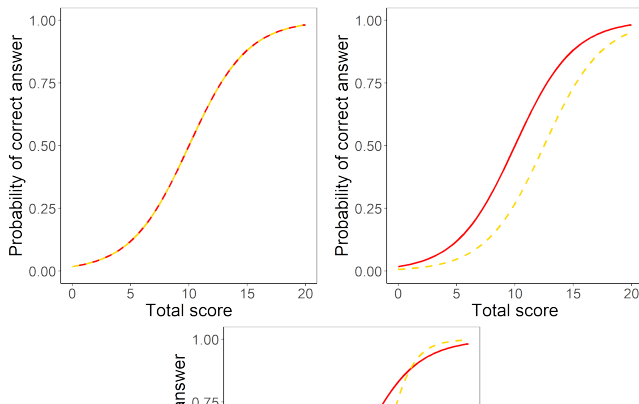
Mantel-Haenszel test

- Test of independence of two binary variables: item score and group membership.
- X^2 test, but incorporating also ability score
- Looking at contingency tables **for each level of total score**, adding up

Logistic regression for DIF detection

$$P(Y_{ij} = 1|X_i, G_i) = \frac{e^{\beta_{0j} + \beta_{1j}X_i}}{1 + e^{\beta_{0j} + \beta_{1j}X_i}} \frac{e^{\beta_{0j} + \beta_{1j}X_i + \beta_{2j}G_i}}{1 + e^{\beta_{0j} + \beta_{1j}X_i + \beta_{2j}G_i}} \frac{e^{\beta_{0j} + \beta_{1j}X_i + \beta_{2j}G_i + \beta_{3j}X_i}}{1 + e^{\beta_{0j} + \beta_{1j}X_i + \beta_{2j}G_i + \beta_{3j}X_i}}$$

= probability of correct answer of student i to item j
 X_i total score, G_i group



Further Topics

Further DIF detection methods:

- Non-linear regression
- SIBTEST
- IRT-based methods

Further issues in DIF detection:

- Correction for multiple comparisons
- Purification
- DIF Effect size

Conclusion

DIF/DDF analysis should be used routinely in test development

- to check for fairness with respect to groups
- to inform teaching

DIF detection methods

- Delta-Plot
- Mantel-Haenszel test
- Logistic regression
- Further (NLR, SIBTEST, IRT/based methods)

Thank you for your attention!

www.cs.cas.cz/martinkova

References

- McFarland, Price, Wenderoth, Martinková, Cliff, Michael, Modell and Wright (2017). Development and Validation of the Homeostasis Concept Inventory. *CBE Life Sciences Education*, vol. 16 no. 2 ar35.
[doi 10.1187/cbe.16-10-0305](https://doi.org/10.1187/cbe.16-10-0305)
- Martinková, Drabinová, Liaw, Sanders, McFarland & Price (2017). Checking Equity: Why DIF Analysis should be a Routine Part of Developing Conceptual Assessments. *CBE-Life Sciences Education*, 16(2), rm2.
[doi 10.1187/cbe.16-10-0307](https://doi.org/10.1187/cbe.16-10-0307)
- Drabinová & Martinková (2017). Detection of Differential Item Functioning with Non-Linear Regression: Non-IRT Approach Accounting for Guessing. *Journal of Educational Measurement*, 54(4), pp. 498-517, 2017.
dx.doi.org/10.1111/jedm.12158