

Item Characteristic Curves generated from common CTT Item Statistics

Diego Figueiras¹ & John T. Kulas¹

¹ Montclair State University

Author Note

Add complete departmental affiliations for each author here. Each new line herein must be indented, like this line.

Enter author note here.

The authors made the following contributions. John T. Kulas: .

Correspondence concerning this article should be addressed to Diego Figueiras, Dickson Hall 226. E-mail: figueirasd1@montclair.edu

Item Characteristic Curves generated from common CTT Item Statistics

Introduction

Item characteristic curves are referenced by psychometricians as visual indicators of important attributes of assessment items - most commonly *difficulty* and/or *discrimination*. Within these visual presentations the x-axis ranges along “trait” levels (by convention annotated with the greek θ), whereas the y-axis displays probabilities of responding to the item within a given response category. In the context of true tests, the response categories are binary, and the y-axis probability refers to the likelihood of a “correct” response. From this orientation, each item within a test can be represented with a visualized “curve.” By looking at it, we can know the likelihood with which respondents of any trait level would answer any item correctly. If the curve is leaning towards the lower end of the trait level, this indicates that it is easy to answer the item correctly. On the contrary, if the curve is leaning towards the higher end of the trait level, this indicates that the item is difficult. If the curve is steep, this indicates high discrimination among respondents; if it is flat, it indicates no discrimination.

Psychometricians who examine ICCs usually do it using Item Response Theory and Rasch models to get the parameters necessary to plot the curves. In a 2PL model, these would be item difficulty and item discrimination. Item difficulty is the necessary trait level for a respondent to have a 50/50 chance to answer the item correctly. Item discrimination is the degree to which an item can differentiate among individuals with low and high levels of the trait. From a Classical Test Theory (CTT) frame of thinking, the difficulty of an item is determined by looking at the p-values of the items, while discrimination is determined by checking the Cronbach alpha and the corrected item total correlations. Psychometricians who look at these CTT parameters don’t typically use them to plot ICCs. There is no reason for them not to, since ICCs based on CTT parameters could provide information as valuable as those based on IRT or Rasch without the need of being familiar with these

models and with how to compute the necessary estimates. Fan states in summary that IRT and CTT "... framework produce very similar item and person statistics" (p.379).

There is research that shows that there is little difference between the parameters of both frameworks. @hambleton1993comparison concluded that "no study provides enough empirical evidence on the extent of disparity between the two frameworks and the superiority of IRT over CTT despite the theoretical differences."

Fan (1998) conducted a study to empirically test the differences between the two frameworks. According to him, "The findings here simply show that the two measurement frameworks produced very similar item and person statistics both in terms of the comparability of item and person statistics between the two frameworks and in terms of the degree of invariance of item statistics from the two competing measurement frameworks." In his study, Fan (1998) looked at the correlations between ability estimates and item difficulty in CTT and all three IRT models. These correlations were very high, between high .80 and low .90. As of item discrimination, correlations were moderate to high, with only a few being very low.

He also looked at the item invariance for all models. In theory, the major advantage of IRT models over CTT is that the latter has a circular dependency between the item and person statistics, while IRT has no such dependency, which means that the item parameters don't depend on the sample and the person parameters don't depend on the set of items. This property of invariance is very important, since item estimates can be used regardless of the sample you are giving the test or assessment to. An item will always have the same level of difficulty regardless of who is responding, for example.

What Fan (1998) got on his study, however, shows empirical evidence against this supposed advantage of IRT against CTT. The CTT item difficulty and discrimination degrees of invariance were highly correlated with those of IRT, indicating that they were highly comparable.

Lord (2012) described a function that approximates the relationship between IRT parameters and the CTT discrimination index of an item-test biserial correlation:

$$a_i \cong \frac{r_i}{\sqrt{1 - r_i^2}}$$

This formula wasn't intended for practical purposes but rather to assist in the conceptual comprehension of the discrimination parameter in IRT for people who were more familiar with CTT procedures. In an effort to move from the conceptual to a practical application, Kulas et al. (2017) proposed a modification that minimized the average residual (either a_i or r_i , where r_i is the *corrected* item-total *point-biserial* correlation).

Simulations identified systematic slope and inflection differences across item with differing item difficulty values, so the formula was further changed to include the following modifiers This revised formula is used in the current presentation:

$$\hat{a}_i \cong [(.51 + .02z_g + .3z_g^2)r] + [(.57 - .009z_g + .19z_g^2)\frac{e^r - e^{-r}}{e - e^r}]$$

Where g is the absolute deviation from 50% responding an item correctly and 50% responding incorrectly (e.g., a “p-value” of .5). Z_g is the standard normal deviation associated with g . The transformation of the standard p-value was recommended in order to scale this index along an interval-level metric more directly analogous to the IRT b parameter. Figure XX visualizes the re-specifications of Lord's formula at p-values (difficulty) of .5, .3 (or .7), and .1 (or .9).

As we can see, the higher the corrected item-total correlations, the higher the estimated IRT a-parameter (discrimination). Also, as the p-values (difficulty) deviates from 0, the relationship between the estimated IRT a-parameter and the corrected item-total correlations becomes stronger.

Practitioners and researchers that don't use IRT or Rasch models and instead opt to

follow a CTT philosophy would benefit from having ICCs that use CTT statistics. This study intends to show evidence of the overlapping nature of CTT and IRT parameters when it comes to plotting ICCs.

Study 1 - Visual of discrimination relationship

The purpose of study 1 is to look at the visualizations resulting from Kulas et al. (2017) formula on simulated data. We hypothesize that the relationship between the estimated IRT a-parameter and the corrected item-total correlations will be stronger as the latter deviates from 0, which would mean that the item has more discrimination.

Procedure and methods

We simulated data using Han (2007) software. Our sample was 10,000 observations, with a mean of 0 and a standard deviation of 1. The number of items were 100, with response categories of either correct or incorrect (1 and 0). The mean for the a parameter was 2, and the standard deviation 0.8. The mean for parameter b was 0 and the standard deviation 0.5.

Results

Study 2 - Item Characteristic Curves comparisons.

The purpose of study 2 is to simulate a lot of test data and then generate ICCs based on the IRT model and then we compare that to our CTT estimates.

Procedure and materials

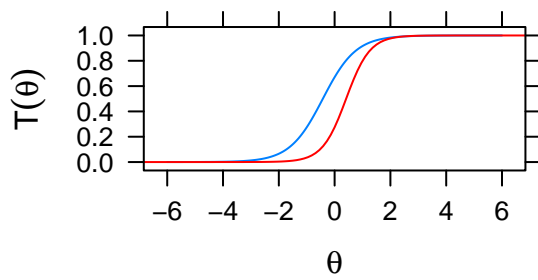
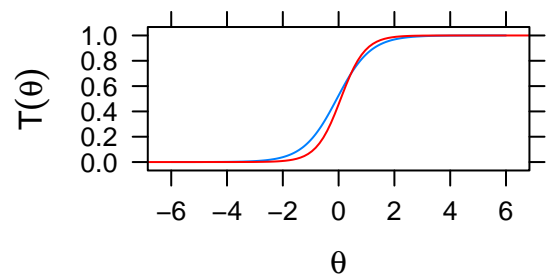
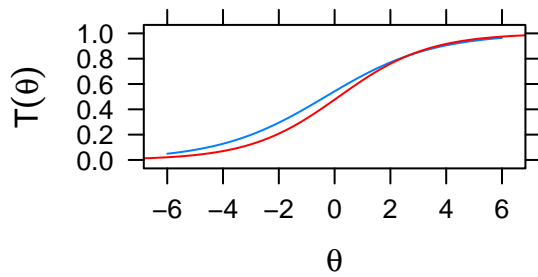
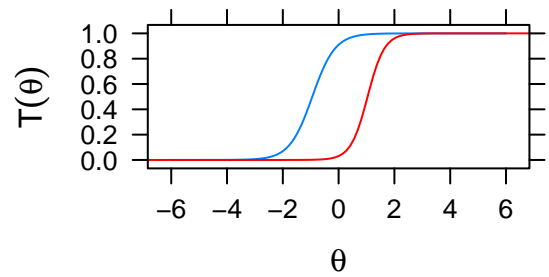
The same simulated data as in study 1 was used. The mirt package was used to compute the IRT statistics. The blue curves were plotted using 2PL IRT parameters (a

and b), while the red curves were plotted using CTT parameters (p-values and corrected item-total correlations, modifying them with Kulas et al. (2017) formulas).

Results

We used R [Version 4.1.1; R Core Team (2020)] and the R-packages *ape* [R-ape], *dplyr* [Version 1.0.7; Wickham et al. (2021)], *DT* [Version 0.19; Xie et al. (2021)], *forcats* [Version 0.5.1; Wickham (2021a)], *formattable* (Ren & Russell, 2021), *geiger* [Version 2.0.7; Alfaro et al. (2009); Eastman et al. (2011); Slater et al. (2012); Harmon et al. (2008); Pennell et al. (2014)], *ggplot2* [Version 3.3.5; Wickham (2016)], *gridExtra* [Version 2.3; Auguie (2017)], *irtplay* [Version 1.6.2; Lim (2020)], *jpeg* [Version 0.1.9; Urbanek (2021)], *knitr* [Version 1.34; Xie (2015)], *lattice* [Version 0.20.44; Sarkar (2008); Sarkar and Andrews (2019)], *latticeExtra* [Version 0.6.29; Sarkar and Andrews (2019)], *markdown* [Version 1.1; Allaire et al. (2019); Xie et al. (2018); Xie et al. (2020)], *mirt* [Version 1.34; Chalmers (2012)], *officer* (Gohel, 2021), *papaja* [Version 0.1.0.9997; Aust and Barth (2020)], *pdfutils* [Version 3.0.1; Ooms (2021)], *psych* [Version 2.1.9; Revelle (2021)], *purrr* [Version 0.3.4; Henry and Wickham (2020)], *readr* [Version 2.0.1; Wickham and Hester (2021)], *readxl* [Version 1.3.1; Wickham and Bryan (2019)], *reticulate* [Version 1.22; Ushey et al. (2021)], *rmarkdown* [Version 2.11; Xie et al. (2018); Xie et al. (2020)], *shiny* [Version 1.7.0; Chang et al. (2021)], *stringr* [Version 1.4.0; Wickham (2019)], *tibble* [Version 3.1.4; Müller and Wickham (2021)], *tidyr* [Version 1.1.3; Wickham (2021b)], *tidyverse* [Version 1.3.1; Wickham et al. (2019)], and *tinytex* [Version 0.33; Xie (2019)] for all our analyses. The area between ICC's was calculated between CTT-derived and IRT-derived ICCs. The average difference for all 100 curves was 0.35.

Item Characteristic Curves

**Medium area between curves (0.36)****Low area between curves (0.03)****Low area between curves (0.09)****Big area between curves (0.81)**

Results

Discussion

References

- Alfaro, M., Santini, F., Brock, C., Alamillo, H., Dornburg, A., Rabosky, D., Carnevale, G., & Harmon, L. (2009). Nine exceptional radiations plus high turnover explain species diversity in jawed vertebrates. *Proceedings of the National Academy of Sciences of the United States of America*, 106, 13410–13414.
- Allaire, J., Horner, J., Xie, Y., Marti, V., & Porte, N. (2019). *Markdown: Render markdown with the c library 'sundown'*.
<https://CRAN.R-project.org/package=markdown>
- Auguie, B. (2017). *gridExtra: Miscellaneous functions for "grid" graphics*.
<https://CRAN.R-project.org/package=gridExtra>
- Aust, F., & Barth, M. (2020). *papaja: Create APA manuscripts with R Markdown*.
<https://github.com/crsh/papaja>
- Chalmers, R. P. (2012). mirt: A multidimensional item response theory package for the R environment. *Journal of Statistical Software*, 48(6), 1–29.
<https://doi.org/10.18637/jss.v048.i06>
- Chang, W., Cheng, J., Allaire, J., Sievert, C., Schloerke, B., Xie, Y., Allen, J., McPherson, J., Dipert, A., & Borges, B. (2021). *Shiny: Web application framework for r*. <https://CRAN.R-project.org/package=shiny>
- Eastman, J., Alfaro, M., Joyce, P., Hipp, A., & Harmon, L. (2011). A novel comparative method for identifying shifts in the rate of character evolution on trees. *Evolution*, 65, 3578–3589.
- Fan, X. (1998). Item response theory and classical test theory: An empirical comparison of their item/person statistics. *Educational and Psychological Measurement*, 58(3), 357–381.

- Gohel, D. (2021). *Officer: Manipulation of microsoft word and PowerPoint documents*. <https://CRAN.R-project.org/package=officer>
- Han, K. (2007). WinGen3: Windows software that generates IRT parameters and item responses [computer program]. *Amherst, MA: Center for Educational Assessment, University of Massachusetts Amherst*.
- Harmon, L., Weir, J., Brock, C., Glor, R., & Challenger, W. (2008). GEIGER: Investigating evolutionary radiations. *Bioinformatics*, *24*, 129–131.
- Henry, L., & Wickham, H. (2020). *Purrr: Functional programming tools*. <https://CRAN.R-project.org/package=purrr>
- Kulas, J. T., Smith, J. A., & Xu, H. (2017). Approximate functional relationship between IRT and CTT item discrimination indices: A simulation, validation, and practical extension of lord’s (1980) formula. *Journal of Applied Measurement*, *18*(4), 393–407.
- Lim, H. (2020). *Irtplay: Unidimensional item response theory modeling*. <https://CRAN.R-project.org/package=irtplay>
- Lord, F. M. (2012). *Applications of item response theory to practical testing problems*. Routledge.
- Müller, K., & Wickham, H. (2021). *Tibble: Simple data frames*. <https://CRAN.R-project.org/package=tibble>
- Ooms, J. (2021). *Pdftools: Text extraction, rendering and converting of PDF documents*. <https://CRAN.R-project.org/package=pdfutils>
- Pennell, M., Eastman, J., Slater, G., Brown, J., Uyeda, J., Fitzjohn, R., Alfaro, M., & Harmon, L. (2014). Geiger v2.0: An expanded suite of methods for fitting macroevolutionary models to phylogenetic trees. *Bioinformatics*, *30*, 2216–2218.
- R Core Team. (2020). *R: A language and environment for statistical computing*. R

- Foundation for Statistical Computing. <https://www.R-project.org/>
- Ren, K., & Russell, K. (2021). *Formattable: Create 'formattable' data structures*.
<https://CRAN.R-project.org/package=formattable>
- Revelle, W. (2021). *Psych: Procedures for psychological, psychometric, and personality research*. Northwestern University.
<https://CRAN.R-project.org/package=psych>
- Sarkar, D. (2008). *Lattice: Multivariate data visualization with r*. Springer.
<http://lmdvr.r-forge.r-project.org>
- Sarkar, D., & Andrews, F. (2019). *latticeExtra: Extra graphical utilities based on lattice*. <https://CRAN.R-project.org/package=latticeExtra>
- Slater, G., Harmon, L., Wegmann, D., Joyce, P., Revell, L., & Alfaro, M. (2012). Fitting models of continuous trait evolution to incompletely sampled comparative data using approximate bayesian computation. *Evolution*, 66, 752–762.
- Urbanek, S. (2021). *Jpeg: Read and write JPEG images*.
<https://CRAN.R-project.org/package=jpeg>
- Ushey, K., Allaire, J., & Tang, Y. (2021). *Reticulate: Interface to 'python'*.
<https://CRAN.R-project.org/package=reticulate>
- Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>
- Wickham, H. (2019). *Stringr: Simple, consistent wrappers for common string operations*. <https://CRAN.R-project.org/package=stringr>
- Wickham, H. (2021a). *Forcats: Tools for working with categorical variables (factors)*. <https://CRAN.R-project.org/package=forcats>

Wickham, H. (2021b). *Tidyr: Tidy messy data*.

<https://CRAN.R-project.org/package=tidyr>

Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R.,

Grolemund, G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T. L.,

Miller, E., Bache, S. M., Müller, K., Ooms, J., Robinson, D., Seidel, D. P.,

Spinu, V., ... Yutani, H. (2019). Welcome to the tidyverse. *Journal of Open*

Source Software, 4(43), 1686. <https://doi.org/10.21105/joss.01686>

Wickham, H., & Bryan, J. (2019). *Readxl: Read excel files*.

<https://CRAN.R-project.org/package=readxl>

Wickham, H., François, R., Henry, L., & Müller, K. (2021). *Dplyr: A grammar of*

data manipulation. <https://CRAN.R-project.org/package=dplyr>

Wickham, H., & Hester, J. (2021). *Readr: Read rectangular text data*.

<https://CRAN.R-project.org/package=readr>

Xie, Y. (2015). *Dynamic documents with R and knitr* (2nd ed.). Chapman;

Hall/CRC. <https://yihui.org/knitr/>

Xie, Y. (2019). TinyTeX: A lightweight, cross-platform, and easy-to-maintain

LaTeX distribution based on TeX live. *TUGboat*, 1, 30–32.

<http://tug.org/TUGboat/Contents/contents40-1.html>

Xie, Y., Allaire, J. J., & Grolemund, G. (2018). *R markdown: The definitive guide*.

Chapman; Hall/CRC. <https://bookdown.org/yihui/rmarkdown>

Xie, Y., Cheng, J., & Tan, X. (2021). *DT: A wrapper of the JavaScript library*

'DataTables'. <https://CRAN.R-project.org/package=DT>

Xie, Y., Dervieux, C., & Riederer, E. (2020). *R markdown cookbook*. Chapman;

Hall/CRC. <https://bookdown.org/yihui/rmarkdown-cookbook>

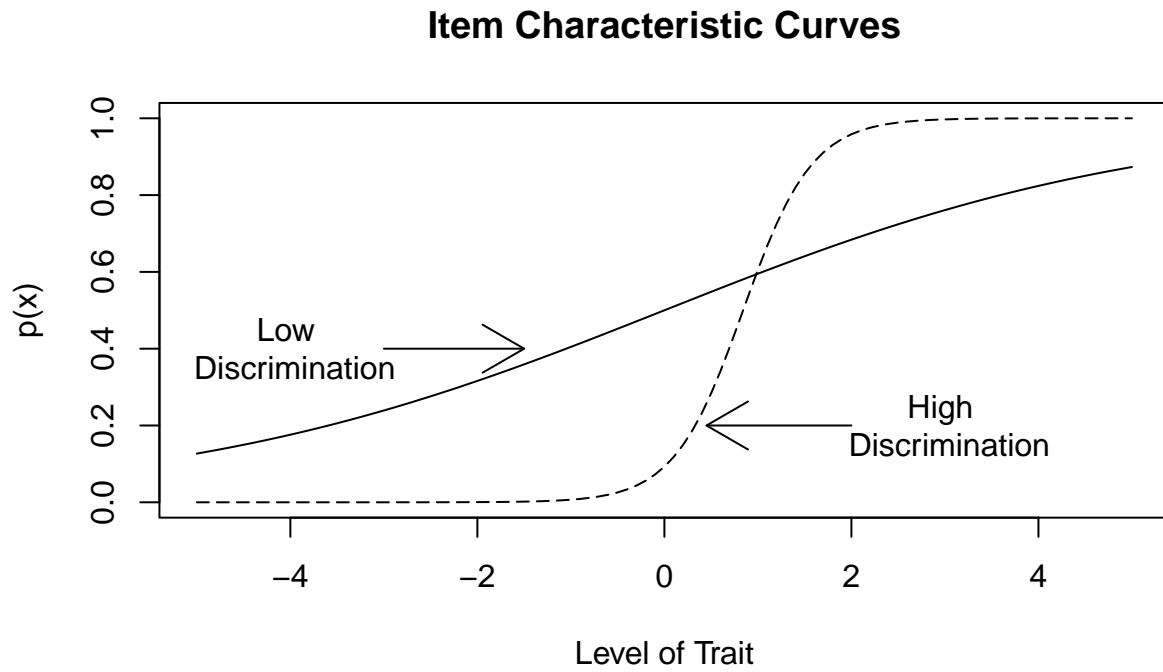


Figure 1. Item characteristic curves primarily reflecting differences in discrimination.

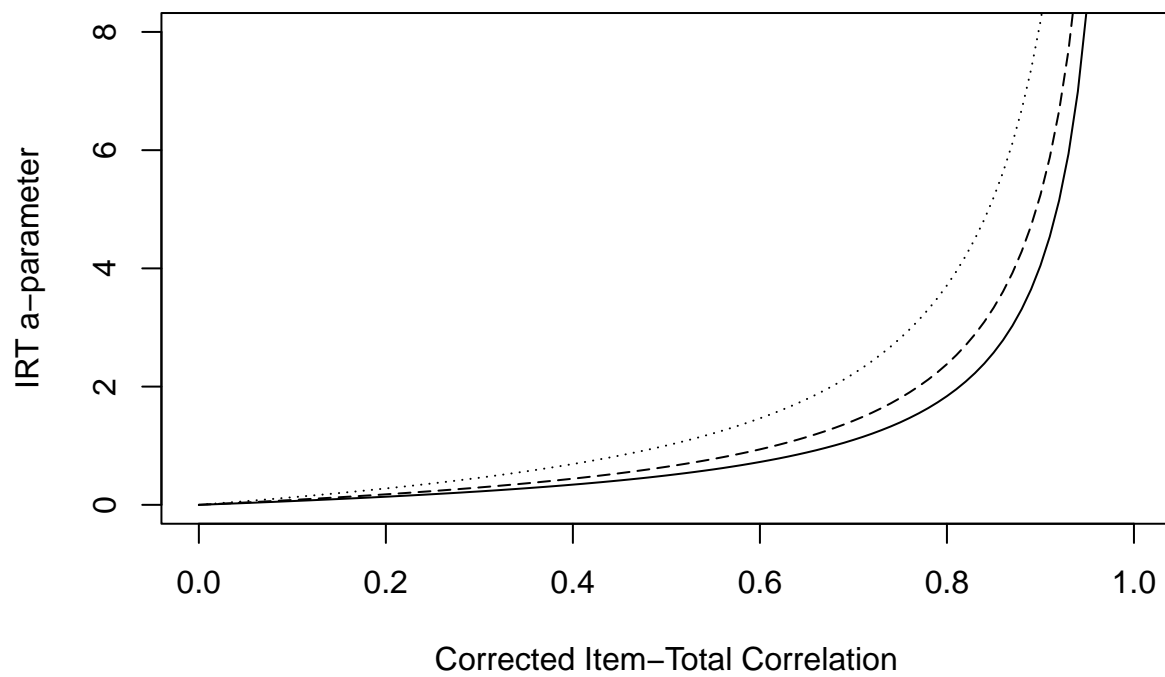


Figure 2. Empricially-derived functional relationship between the IRT a parameter and the CTT corrected-item total correlation as a function of item difficulty (p-value; solid = .5, dashed = .3/.7, dotted = .1/.9).

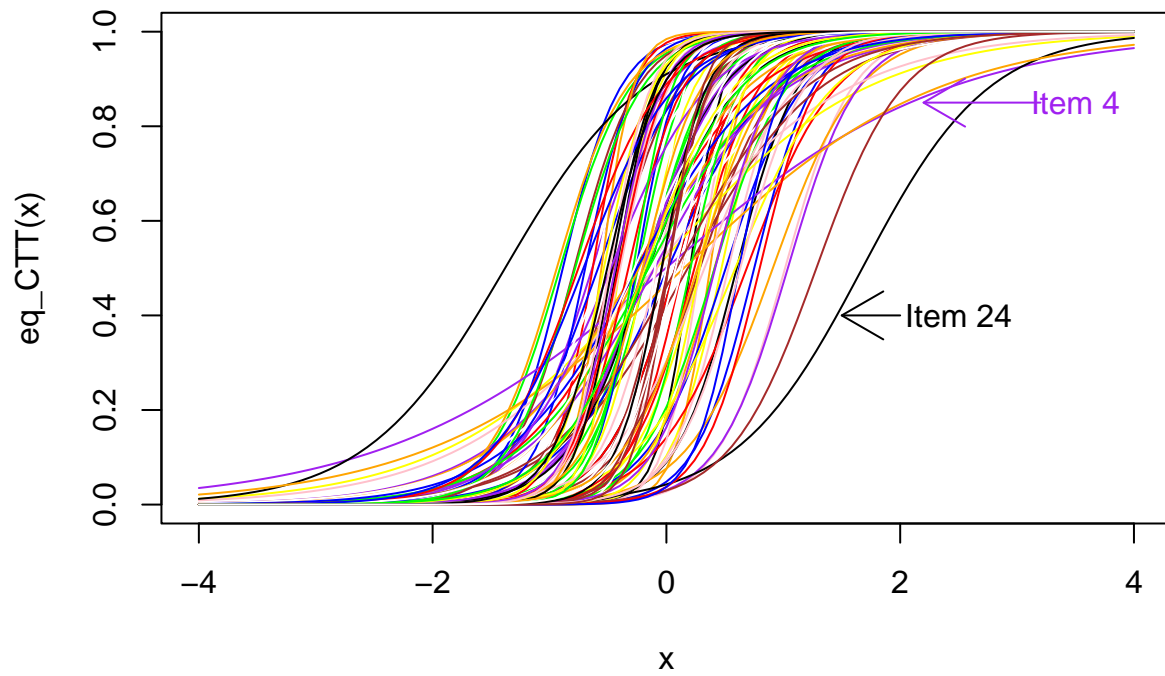


Figure 3. ICCs derived from only CTT parameters (with two noteworthy ICCs annotated).

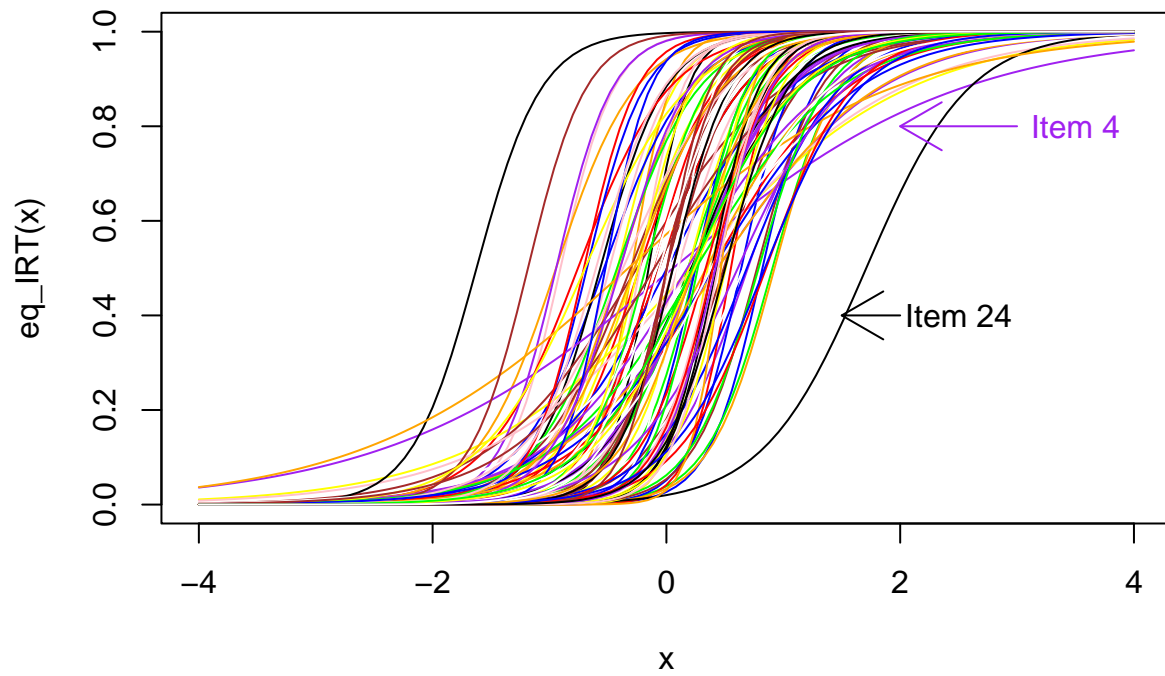


Figure 4. Typical ICCs derived from IRT parameters (same noteworthy items annotated).