

Psychology 360: Personality Research

Psychometric Theory – Reliability Theory

William Revelle

Department of Psychology
Northwestern University
Evanston, Illinois USA



NORTHWESTERN
UNIVERSITY

October, 2022

Outline: Part II: Generalizability Theory

Outline of Part II: Generalizability Theory and the IntraClass Correlation

Intraclass correlations

ICC of judges

Kappa

Cohen's kappa

Weighted kappa

Two approaches

Various IRT models

Polytomous items

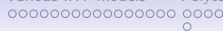
Ordered response categories

Differential Item Functioning

Factor analysis & IRT

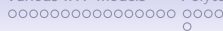
Non-monotone Trace lines

(C) A T



Reliability of judges

- When raters (judges) rate targets, there are multiple sources of variance
 - Between targets
 - Between judges
 - Interaction of judges and targets
- The intraclass correlation is an analysis of variance decomposition of these components
- Different ICC's depending upon what is important to consider
 - Absolute scores: each target gets just one judge, and judges differ
 - Relative scores: each judge rates multiple targets, and the mean for the judge is removed
 - Each judge rates multiple targets, judge and target effects removed



ICC is done by calling anova

```

aov.x <- aov(values ~ subs + ind, data = x.df)
s.aov <- summary(aov.x)
stats <- matrix(unlist(s.aov), ncol = 3, byrow = TRUE)
MSB <- stats[3, 1]
MSW <- (stats[2, 2] + stats[2, 3]) / (stats[1, 2] + stats[1,
  3])
MSJ <- stats[3, 2]
MSE <- stats[3, 3]
ICC1 <- (MSB - MSW) / (MSB + (nj - 1) * MSW)
ICC2 <- (MSB - MSE) / (MSB + (nj - 1) * MSE + nj * (MSJ - MSE) / n
ICC3 <- (MSB - MSE) / (MSB + (nj - 1) * MSE)
ICC12 <- (MSB - MSW) / (MSB)
ICC22 <- (MSB - MSE) / (MSB + (MSJ - MSE) / n.obs)
ICC32 <- (MSB - MSE) / MSB

```

Intraclass Correlations using the ICC function

```
> print(ICC(Ratings), all=TRUE) #get more output than normal
```

```
$results
```

	type	ICC	F	df1	df2	p	lower bound	upper bound
Single_raters_absolute	ICC1	0.32	3.84	5	30	0.01	0.04	0.79
Single_random_raters	ICC2	0.37	10.37	5	25	0.00	0.09	0.80
Single_fixed_raters	ICC3	0.61	10.37	5	25	0.00	0.28	0.91
Average_raters_absolute	ICC1k	0.74	3.84	5	30	0.01	0.21	0.96
Average_random_raters	ICC2k	0.78	10.37	5	25	0.00	0.38	0.96
Average_fixed_raters	ICC3k	0.90	10.37	5	25	0.00	0.70	0.98

```
$summary
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
subs	5	141.667	28.3333	10.366	1.801e-05 ***
ind	5	153.000	30.6000	11.195	9.644e-06 ***
Residuals	25	68.333	2.7333		

```
Signif. codes:  0  ***  0.001  **  0.01  *  0.05  .  0.1  1
```

```
$stats
```

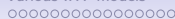
	[,1]	[,2]	[,3]
[1,]	5.000000e+00	5.000000e+00	25.000000
[2,]	1.416667e+02	1.530000e+02	68.333333
[3,]	2.833333e+01	3.060000e+01	2.733333
[4,]	1.036585e+01	1.119512e+01	NA
[5,]	1.800581e-05	9.644359e-06	NA

```
$MSW
```

```
[1] 7.377778
```

```
$Call
```

```
ICC(x = Ratings)
```

Cohen's kappa and weighted kappa

- When considering agreement in diagnostic categories, without numerical values, it is useful to consider the kappa coefficient.
 - Emphasizes matches of ratings
 - Doesn't consider how far off disagreements are.
- Weighted kappa weights the off diagonal distance.
- Diagnostic categories: normal, neurotic, psychotic

Cohen kappa and weighted kappa

```
> cohen
      [,1] [,2] [,3]
[1,] 0.44 0.07 0.09
[2,] 0.05 0.20 0.05
[3,] 0.01 0.03 0.06
> cohen.weights
      [,1] [,2] [,3]
[1,] 0     1     3
[2,] 1     0     6
[3,] 3     6     0
> cohen.kappa(cohen, cohen.weights)
Call: cohen.kappa1(x = x, w = w, n.obs = n.obs, alpha = alpha)
```

Cohen Kappa and Weighted Kappa correlation **coefficients** and confid

	lower	estimate	upper
unweighted kappa	-0.92	0.49	1.9
weighted kappa	-10.04	0.35	10.7

see the other examples in ?cohen.kappa

Outline of Part III: the New Psychometrics

Intraclass correlations

ICC of judges

Kappa

Cohen's kappa

Weighted kappa

Two approaches

Various IRT models

Polytomous items

Ordered response categories

Differential Item Functioning

Factor analysis & IRT

Non-monotone Trace lines

(C) A T

Classical Reliability

1. Classical model of reliability
 - Observed = True + Error
 - Reliability = $1 - \frac{\sigma_{error}^2}{\sigma_{observed}^2}$
 - Reliability = $r_{xx} = r_{x_{domain}}^2$
 - Reliability as correlation of a test with a test just like it
2. Reliability requires variance in observed score
 - As σ_x^2 decreases so will $r_{xx} = 1 - \frac{\sigma_{error}^2}{\sigma_{observed}^2}$
3. Alternate estimates of reliability all share this need for variance
 - 3.1 Internal Consistency
 - 3.2 Alternate Form
 - 3.3 Test-retest
 - 3.4 Between rater
4. Item difficulty is ignored, items assumed to be sampled at random

Estimating the model

The probability of missing an item, q , is just $1 - p(\text{correct})$ and thus the *odds ratio* of being correct for a person with ability, θ_i , on an item with difficulty, δ_j is

$$OR_{ij} = \frac{p}{1-p} = \frac{p}{q} = \frac{\frac{1}{1+e^{\delta_j-\theta_i}}}{1-\frac{1}{1+e^{\delta_j-\theta_i}}} = \frac{\frac{1}{1+e^{\delta_j-\theta_i}}}{\frac{e^{\delta_j-\theta_i}}{1+e^{\delta_j-\theta_i}}} = \frac{1}{e^{\delta_j-\theta_i}} = e^{\theta_i-\delta_j}. \quad (3)$$

That is, the odds ratio will be a exponential function of the difference between a person's ability and the task difficulty. The odds of a particular pattern of rights and wrongs over n items will be the product of n odds ratios

$$OR_{i1} OR_{i2} \dots OR_{in} = \prod_{j=1}^n e^{\theta_i-\delta_j} = e^{n\theta_i} e^{-\sum_{j=1}^n \delta_j}. \quad (4)$$

Difficulty is just a function of probability correct

Similarly, the pattern of the odds of correct and incorrect responses across people for a particular item with difficulty δ_j will be

$$OR_{1j} OR_{2j} \dots OR_{nj} = \frac{P}{Q} = \prod_{i=1}^N e^{\theta_i - \delta_j} = e^{\sum_{i=1}^N (\theta_i) - N\delta_j} \quad (7)$$

and taking logs of both sides leads to

$$\ln \frac{P}{Q} = \sum_{i=1}^N (\theta_i) - N\delta_j. \quad (8)$$

Letting the average ability $\bar{\theta} = 0$ leads to the conclusion that the difficulty of an item for all subjects, δ_j , is the logarithm of Q/P divided by the number of subjects, N ,

$$\delta_j = \frac{\ln \frac{Q}{P}}{N}. \quad (9)$$

Rasch model in words

That is, the estimate of ability (Equation 6) for items with an average difficulty of 0 does not require knowing the difficulty of any particular item, but is just a function of the pattern of corrects and incorrects for a subject across all items.

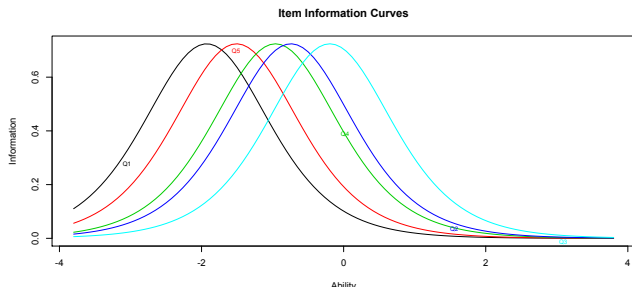
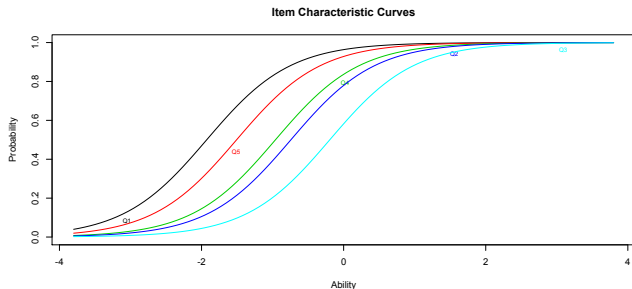
Similarly, the estimate of item difficulty across people ranging in ability, but with an average ability of 0 (Equation 9) is a function of the response pattern of all the subjects on that one item and does not depend upon knowing any one person's ability. The assumptions that average difficulty and average ability are 0 are merely to fix the scales. Replacing the average values with a non-zero value just adds a constant to the estimates.

Rasch as a high jump

The independence of ability from difficulty implied in equations 6 and 9 makes estimation of both values very straightforward. These two equations also have the important implication that the number correct ($n\bar{p}$ for a subject, $N\bar{p}$ for an item) is monotonically, but not linearly related to ability or to difficulty.

That the estimated ability is independent of the pattern of rights and wrongs but just depends upon the total number correct is seen as both a strength and a weakness of the Rasch model. From the perspective of *fundamental measurement*, Rasch scoring provides an additive interval scale: for all people and items, if $\theta_i < \theta_j$ and $\delta_k < \delta_l$ then $p(x|\theta_i, \delta_k) < p(x|\theta_j, \delta_l)$. But this very additivity treats all patterns of scores with the same number correct as equal and ignores potential information in the pattern of responses.

Rasch estimates from ltm



The LSAT example from ltm

```
data(bock)
> ord <- order(colMeans(lsat6), decreasing=TRUE)
> lsat6.sorted <- lsat6[,ord]
> describe(lsat6.sorted)
> Tau <- round(-qnorm(colMeans(lsat6.sorted)), 2) #tau = estimates of threshold
> rasch(lsat6.sorted, constraint=cbind(ncol(lsat6.sorted)+1, 1.702))
```

	var	n	mean	sd	median	trimmed	mad	min	max	range	skew	kurtosis	se
Q1	1	1000	0.92	0.27	1	1.00	0	0	1	1	-3.20	8.22	0.01
Q5	2	1000	0.87	0.34	1	0.96	0	0	1	1	-2.20	2.83	0.01
Q4	3	1000	0.76	0.43	1	0.83	0	0	1	1	-1.24	-0.48	0.01
Q2	4	1000	0.71	0.45	1	0.76	0	0	1	1	-0.92	-1.16	0.01
Q3	5	1000	0.55	0.50	1	0.57	0	0	1	1	-0.21	-1.96	0.02

```
> Tau
  Q1    Q5    Q4    Q2    Q3
-1.43 -1.13 -0.72 -0.55 -0.13
```

Call:

```
rasch(data = lsat6.sorted, constraint = cbind(ncol(lsat6.sorted) +
  1, 1.702))
```

Coefficients:

Dffclt.Q1	Dffclt.Q5	Dffclt.Q4	Dffclt.Q2	Dffclt.Q3	Dscrnm
-1.927	-1.507	-0.960	-0.742	-0.195	1.702

Item information

When forming a test and evaluating the items within a test, the most useful items are the ones that give the most information about a person's score. In classic test theory, *item information* is the reciprocal of the squared *standard error* for the item or for a one factor test, the ratio of the item communality to its uniqueness:

$$I_j = \frac{1}{\sigma_{e_j}^2} = \frac{h_j^2}{1 - h_j^2}.$$

When estimating ability using IRT, the information for an item is a function of the first derivative of the likelihood function and is maximized at the inflection point of the *icc*.

Estimating item information

The information function for an item is

$$I(f, x_j) = \frac{[P'_j(f)]^2}{P_j(f)Q_j(f)} \quad (10)$$

For the 1PL model, P' , the first derivative of the probability function $P_j(f) = \frac{1}{1+e^{\delta-\theta}}$ is

$$P' = \frac{e^{\delta-\theta}}{(1 + e^{\delta-\theta})^2} \quad (11)$$

which is just $P_j Q_j$ and thus the information for an item is

$$I_j = P_j Q_j. \quad (12)$$

That is, information is maximized when the probability of getting an item correct is the same as getting it wrong, or, in other words, the best estimate for an item's difficulty is that value where half of the subjects pass the item.

2PL and 2PN models

$$p(\text{correct}_{ij}|\theta_i, \alpha_j, \delta_j) = \frac{1}{1 + e^{\alpha_j(\delta_j - \theta_i)}} \quad (14)$$

while in the *two parameter normal ogive (2PN)* model this is

$$p(\text{correct}|\theta, \alpha_j, \delta) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\alpha(\theta - \delta)} e^{-\frac{u^2}{2}} du \quad (15)$$

where $u = \alpha(\theta - \delta)$.

The information function for a two parameter model reflects the item discrimination parameter, α ,

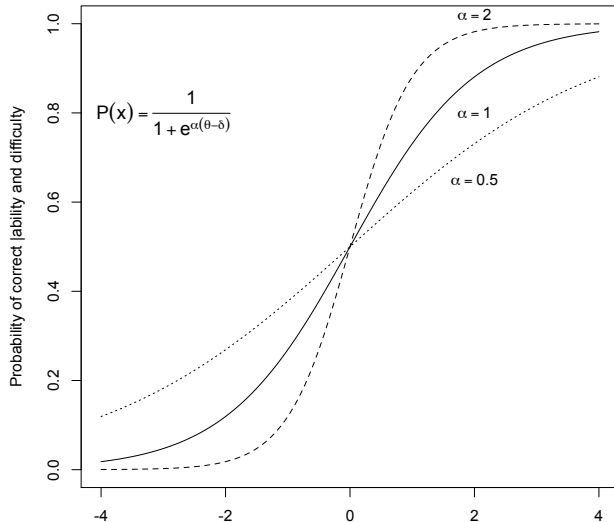
$$I_j = \alpha^2 P_j Q_j \quad (16)$$

which, for a 2PL model is

$$I_j = \alpha_j^2 P_j Q_j = \frac{\alpha_j^2}{(1 + e^{\alpha_j(\delta_j - \theta_j)})^2} \quad (17)$$

The problem of non-parallel trace lines

2PL models differing in their discrimination parameter



Parameter explosion – better fit but at what cost

The 3 parameter model adds a guessing parameter.

$$p(\text{correct}_{ij}|\theta_i, \alpha_j, \delta_j, \gamma_j) = \gamma_j + \frac{1 - \gamma_j}{1 + e^{\alpha_j(\delta_j - \theta_i)}} \quad (18)$$

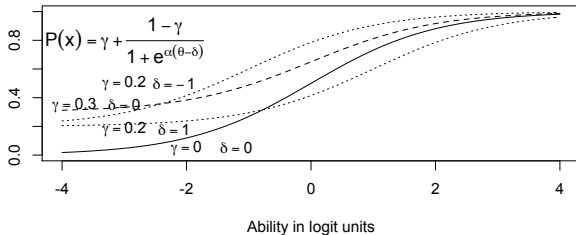
And the four parameter model adds an asymptotic parameter

$$P(x|\theta_i, \alpha, \delta_j, \gamma_j, \zeta_j) = \gamma_j + \frac{\zeta_j - \gamma_j}{1 + e^{\alpha_j(\delta_j - \theta_i)}}. \quad (19)$$

frame

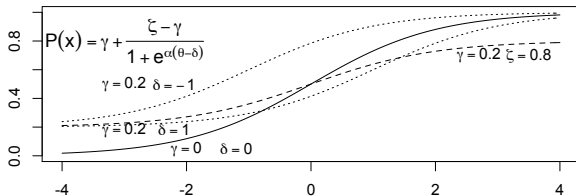
Probability of correct | ability and difficulty

3PL models differing in guessing and difficulty



probability of correct | ability and difficulty

4PL items differing in guessing, difficulty and asymptote



Personality items with monotone trace lines

A typical personality item might ask “How much do you enjoy a lively party” with a five point response scale ranging from “1: not at all” to “5: a great deal” with a neutral category at 3. An alternative response scale for this kind of item is to not have a neutral category but rather have an even number of responses. Thus a six point scale could range from “1: very inaccurate” to “6: very accurate” with no neutral category

The assumption is that the more sociable one is, the higher the response alternative chosen. The probability of endorsing a 1 will increase monotonically the less sociable one is, the probability of endorsing a 5 will increase monotonically the more sociable one is.

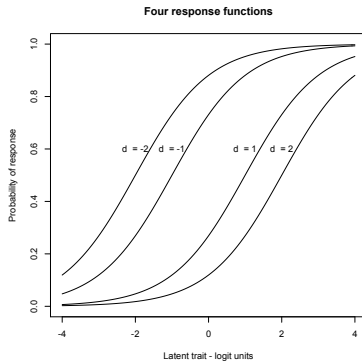
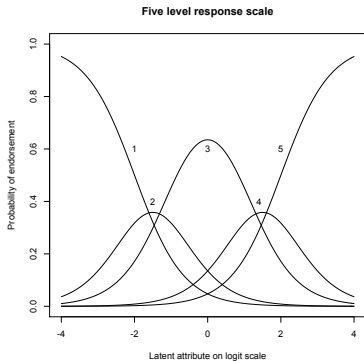
Threshold models

For the 1PL or 2PL logistic model the probability of endorsing the k^{th} response is a function of ability, item thresholds, and the discrimination parameter and is

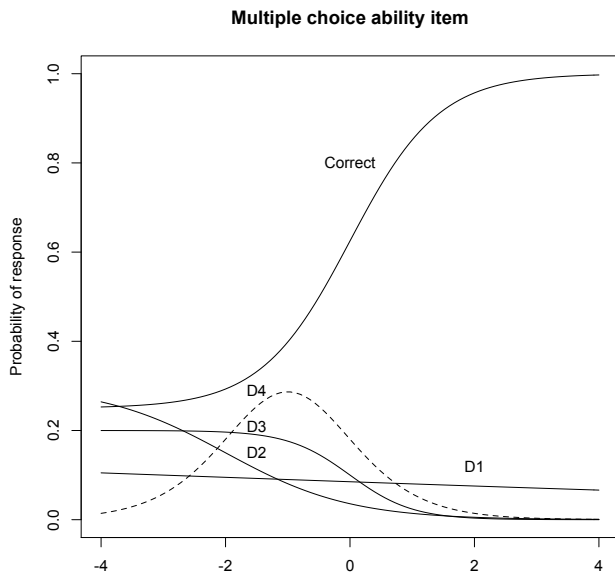
$$P(r = k | \theta_i, \delta_k, \delta_{k-1}, \alpha_k) = P(r | \theta_i, \delta_{k-1}, \alpha_k) - P(r | \theta_i, \delta_k, \alpha_k) = \frac{1}{1 + e^{\alpha_k(\delta_{k-1} - \theta_i)}} - \frac{1}{1 + e^{\alpha_k(\delta_k - \theta_i)}} \quad (20)$$

where all b_k are set to $b_k = 1$ in the 1PL Rasch case.

Responses to a multiple choice polytomous item



Differences in the response shape of multiple choice items



FA and IRT

If the correlations of all of the items reflect one underlying latent variable, then factor analysis of the matrix of tetrachoric correlations should allow for the identification of the regression slopes (α) of the items on the latent variable. These regressions are, of course just the factor loadings. Item difficulty, δ_j and item discrimination, α_j may be found from factor analysis of the tetrachoric correlations where λ_j is just the factor loading on the first factor and τ_j is the normal threshold reported by the tetrachoric function (McDonald, 1999; Lord & Novick, 1968; Takane & de Leeuw, 1987).

$$\delta_j = \frac{D\tau}{\sqrt{1 - \lambda_j^2}}, \quad \alpha_j = \frac{\lambda_j}{\sqrt{1 - \lambda_j^2}} \quad (21)$$

where D is a scaling factor used when converting to the parameterization of *logistic* model and is 1.702 in that case and 1 in the case of the normal ogive model.

FA and IRT

IRT parameters from FA

$$\delta_j = \frac{D\tau}{\sqrt{1 - \lambda_j^2}}, \quad \alpha_j = \frac{\lambda_j}{\sqrt{1 - \lambda_j^2}} \quad (22)$$

FA parameters from IRT

$$\lambda_j = \frac{\alpha_j}{\sqrt{1 + \alpha_j^2}}, \quad \tau_j = \frac{\delta_j}{\sqrt{1 + \alpha_j^2}}.$$

the irt.fa function

```

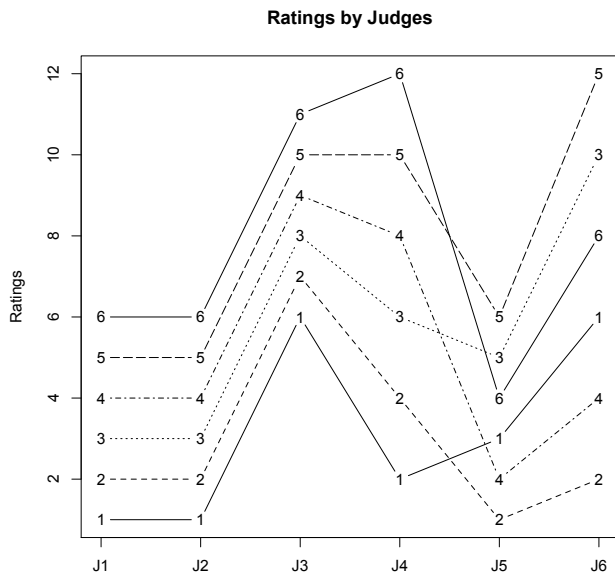
> set.seed(17)
> items <- sim.npn(9,1000,low=-2.5,high=2.5)$items
> p.fa <- irt.fa(items)
  
```

Summary information **by factor** and item

```

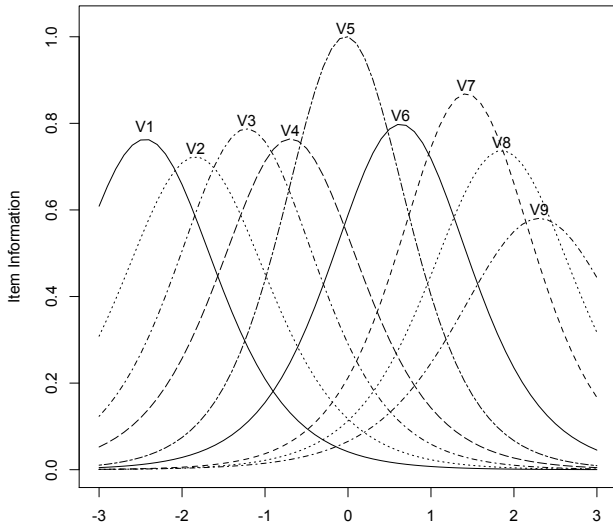
Factor = 1
      -3    -2    -1     0     1     2     3
V1      0.61 0.66 0.21 0.04 0.01 0.00 0.00
V2      0.31 0.71 0.45 0.12 0.02 0.00 0.00
V3      0.12 0.51 0.76 0.29 0.06 0.01 0.00
V4      0.05 0.26 0.71 0.54 0.14 0.03 0.00
V5      0.01 0.07 0.44 1.00 0.40 0.07 0.01
V6      0.00 0.03 0.16 0.59 0.72 0.24 0.05
V7      0.00 0.01 0.04 0.21 0.74 0.66 0.17
V8      0.00 0.00 0.02 0.11 0.45 0.73 0.32
V9      0.00 0.00 0.01 0.07 0.25 0.55 0.44
Test Info 1.11 2.25 2.80 2.97 2.79 2.28 0.99
SEM      0.95 0.67 0.60 0.58 0.60 0.66 1.01
Reliability 0.10 0.55 0.64 0.66 0.64 0.56 -0.01
  
```

Item Characteristic Curves from FA

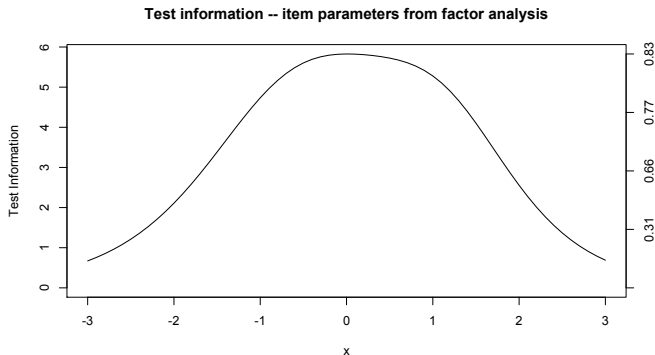


Item information from FA

Item information from factor analysis



Test Information Curve



Comparing three ways of estimating the parameters

```
set.seed(17)
items <- sim.npn(9,1000,low=-2.5,high=2.5)$items
p.fa <- irt.fa(items)$coefficients[1:2]
p.ltm <- ltm(items~z1)$coefficients
p.ra <- rasch(items, constraint = cbind(ncol(items) + 1, 1))$coeff
a <- seq(-2.5,2.5,5/8)
p.df <- data.frame(a,p.fa,p.ltm,p.ra)
round(p.df,2)
```

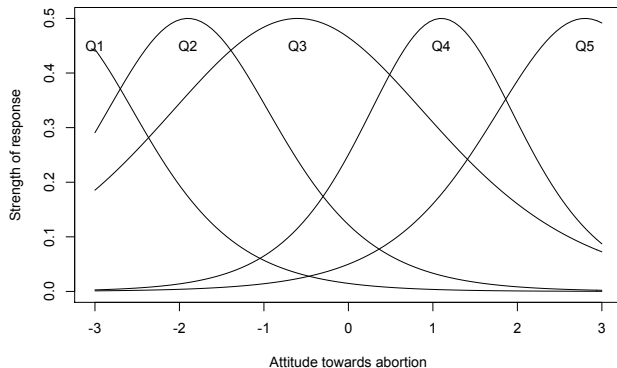
	a	Difficulty	Discrimination	X. Intercept.	z1	beta.i	be
Item 1	-2.50	-2.45	1.03	5.42	2.61	3.64	
1							
Item 2	-1.88	-1.84	1.00	3.35	1.88	2.70	
1							
Item 3	-1.25	-1.22	1.04	2.09	1.77	1.73	
1							
Item 4	-0.62	-0.69	1.03	1.17	1.71	0.98	
1							
Item 5	0.00	-0.03	1.18	0.04	1.94	0.03	
1							
Item 6	0.62	0.63	1.05	-1.05	1.68	-0.88	
1							
Item 7	1.25	1.43	1.10	-2.47	1.90	-1.97	
1							

Attitudes might not have monotone trace lines

1. *Abortion is unacceptable under any circumstances.*
2. *Even if one believes that there may be some exceptions, abortions is still generally wrong.*
3. *There are some clear situations where abortion should be legal, but it should not be permitted in all situations.*
4. *Although abortion on demand seems quite extreme, I generally favor a woman's right to choose.*
5. *Abortion should be legal under any circumstances.*

Ideal point models of attitude

Attitudes reflect an unfolding (ideal point) model



- Lord, F. M. & Novick, M. R. (1968). *Statistical theories of mental test scores*. The Addison-Wesley series in behavioral science: quantitative methods. Reading, Mass.: Addison-Wesley Pub. Co.
- McDonald, R. P. (1999). *Test theory: A unified treatment*. Mahwah, N.J.: L. Erlbaum Associates.
- Rasch, G. (1960). *Probabilistic models for some intelligence and attainment tests*. Chicago: reprinted in 1980 by The University of Chicago Press /Paedagogike Institut, Copenhagen.
- Takane, Y. & de Leeuw, J. (1987). On the relationship between item response theory and factor analysis of discretized variables. *Psychometrika*, 52, 393–408. 10.1007/BF02294363.