

Rasch and IRT Models: A Case Study

Item response models may be used to model the responses of subjects to a number of questions or test items. An item response model with one parameter for item difficulty is known as a Rasch model. Georg Rasch (1901-1980), a Danish statistician, gave an axiomatic derivation of the model in the 1960s. We will be using a conditional (fixed-effects) logit model to illustrate the model, however, Rasch's derivation used a different approach, but one that turns out to be equivalent to the fixed-effects logit. Rasch models are one of the dominant models for binary items (e.g., success/failure on test items) in psychometrics.

Suppose n students are taking a test consisting of k true/false questions. In the Rasch model the log odds of subject i giving a correct response to item j may be modeled using a one-parameter logistic model

$$\begin{aligned} y_{ij} &\sim \text{Ber}(p_{ij}), \\ \text{logit}(p_{ij}) &= \alpha_i - \delta_j, \quad 1 \leq i \leq n, 1 \leq j \leq k, \end{aligned} \tag{1}$$

where p_{ij} is the probability of subject i answering the item j correctly, α_i represent the *ability* of the subject i , and δ_j represent the *difficulty* of the item j . Observed y_{ij} is a scoring – has value 1 if the answer is correct, and 0 if the answer is wrong.

In terms of probability

$$P(y_{ij} = 1) = \text{logit}^{-1}(\alpha_i - \delta_j) = \frac{\exp\{\alpha_i - \delta_j\}}{1 + \exp\{\alpha_i - \delta_j\}}.$$



Figure 1: Georg William Rasch (1901-1980), Danish Statistician.

The failure of an item to fit the model can be traced to two main sources. One is that the model is too simple. It takes account of only one item characteristic – item easiness. This model assumes that all items have the same discrimination, and that the effect of guessing is negligible. Parameters for discrimination and guessing can be included in a more general model. However, their inclusion makes the application of the model to actual measurement involved.

The model (1) is not identifiable because a common constant can be added to all the abilities α_i and all the difficulties δ_j , and the predictions of the model will remain the same,

$$\alpha_i - \delta_j = \alpha_i + c - (\delta_j + c) = \alpha'_i - \delta'_j.$$

The probabilities depend only on the difference of the ability and difficulty parameters, but not on their individual locations.



Figure 2: Test

From the standpoint of classical logistic regression, this nonidentifiability is a simple case of collinearity and can be resolved by several ways. For example, by constraining the parameters (i) setting $\alpha_1 = 0$ (that is, using the person #1 as a baseline), (ii) by setting $\delta_1 = 0$ (so that the first item is the comparison point), (iii) constraining $\sum_i \alpha_i$ to 0, or (iv) constraining $\sum_j \delta_j$ to 0.

A common Bayesian model for (1) assigns normal priors to the ability and difficulty parameters:

$$\begin{aligned}\alpha_i &\sim \mathcal{N}(\mu_\alpha, \tau_\alpha), \quad i = 1, \dots, n \\ \delta_j &\sim \mathcal{N}(\mu_\delta, \tau_\delta), \quad j = 1, \dots, k.\end{aligned}$$

The priors for these parameters are assigned hyperpriors and estimated conditional on the data. This is also referred to as a partial pooling or hierarchical approach (remember the concept of “borrowing strength” in hierarchical models). The model is nonidentifiable for the reasons mentioned above: this time, it is μ_α and μ_δ that are not identifiable, because a constant can be added to each without changing the predictions. Bayesian analysis can proceed with or without the model parameters being identified, since identification is a property of a likelihood. Priors do not really “solve” identification problems, except in the degenerate case of a point-mass spike priors (i.e., parameter restrictions by any other means), however, they could place the estimators in the range where classical restricted estimators would fall. The simplest way to identify (in this Bayesian-vague way) this hierarchical model is set μ_α to 0 (or to set μ_δ to 0, but probably not both due to their relationship in the likelihood).

```
model Rasch
{
for (i in 1:n)
{
for (j in 1:k)
{
y[i,j] ~ dbern(p[i,j])
logit(p[i,j]) <- alpha[i] - delta[j]
# or p[i,j] <- (exp(alpha[i] - delta[j]))/(1+exp(alpha[i] - delta[j]))
}

```

[illegible]

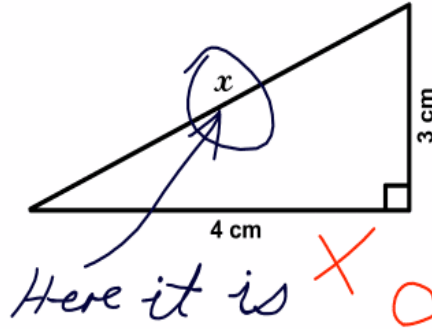
delta[30]	-0.4043	0.2141	0.008461	-0.8185	-0.4013	0.01447	501	2500
delta[31]	0.3747	0.215	0.008428	-0.03354	0.3738	0.7909	501	2500
delta[32]	-0.8829	0.2142	0.008372	-1.32	-0.8804	-0.4746	501	2500
delta[33]	0.7065	0.2145	0.007964	0.3084	0.702	1.142	501	2500

2 Parameter IRT Model (2PLM)

The Rasch model can be generalized to 2-parameter item-response model, known as 2PLM, by allowing the slope of the logistic regression to vary by item:

$$\begin{aligned} y_{ij} &\sim \text{Ber}(p_{ij}), \\ \text{logit}(p_{ij}) &= \gamma_j(\alpha_i - \delta_j), \quad 1 \leq i \leq n, 1 \leq j \leq k, \end{aligned} \quad (2)$$

3. Find x.



Ocular Trauma - by Wade Clarke ©2005

Figure 3: Example of an y_{ij}

In the model (3), parameter γ_j is called the *discrimination* of item j : if $\gamma_j = 0$, then the item does not discriminate at all, and $P(y_{ij} = 1) = 0.5$ for any person, whereas high values of γ_j correspond to a strong relation between ability and the probability of getting a correct response. Negative values of γ_j correspond to items where low-ability persons do better. Such items typically represent mistakes in the construction of the test since test designers generally try to create questions with a high positive discrimination value.

The addition of the discrimination parameter induces a new invariance problem. Model (3) has a so called multiplicative aliasing. This arises when multiplying the γ_j by a constant and dividing the $\alpha_i - \delta_j$ by that same constant.

This indeterminacy can be resolved constraining the α_i 's to have mean 0 and standard deviation 1 or, in Bayesian context, by giving the α_i 's a fixed prior distribution, for example $\mathcal{N}(0, 1)$.

model 2-parameter IR

```
{
for (i in 1:n)
{
for (j in 1:k)
{
y[i,j] ~ dbern(p[i,j])
logit(p[i,j]) <- gamma[j] * ( alpha[i] - delta[j] )
}
```

```

    }
    alpha[i] ~ dnorm(0,1)
  }
  for (j in 1:k)
  {
    delta[j] ~ dnorm(mu.delta, tau.delta)
    gamma[j] ~ dnorm(mu.gamma, tau.gamma)
  }
mu.gamma ~ dnorm(0, 0.001)
tau.gamma ~ dnorm(0, 0.001) I(0,)
mu.delta ~ dnorm(0, 0.001)
tau.delta ~ dnorm(0, 0.001) I(0,)
var.gamma <- 1/tau.gamma
var.delta <- 1/tau.delta
}

```

3 Parameter IRT Model (3PLM)

The 3LPM model is given by

$$\begin{aligned}
 y_{ij} &\sim \mathcal{Ber}(p_{ij}), \\
 p_{ij} &= c_j + (1 - c_j) \frac{\exp\{\gamma_j(\alpha_i - \delta_j)\}}{1 + \exp\{\gamma_j(\alpha_i - \delta_j)\}}, \quad 1 \leq i \leq n, 1 \leq j \leq k.
 \end{aligned} \tag{3}$$

The c parameter is commonly called the *guessing* parameter or *lower asymptote* parameter, because it indicates the probability of responding positively for examinees having very low α .

The more general IRT model is a diffuse polytomous model where each question has m categories of response, with $m > 2$. In this case the probability related to each category depends on item and person parameter as in the other case and also on the category thresholds.

Further Reading

There are many papers in fields different than education where Rasch-type models are appropriate: manufacturing and industry (conforming/non-conforming products), medicine and health systems (healthy/non-healthy), psychology (conscientiousness or cognitive ability/opposite), genetics, etc.

Some references are:

Hahn et al. (2005). Cross-Cultural Evaluation of Health Status Using Item Response Theory: FACT-B Comparisons. *Eval Health Prof.*, **28**, 233–259.

Strong, D. R., Breen, R., & Lejuez, C. W. (2004). Using item response theory to examine gambling affinity as an underlying vulnerability across a continuum of gambling involvement. *Personality and Individual Differences*, **36**, 1515–1529.

Leplege A, Ecosse E. and the WHOQOL Rasch Project Scientific Committee, (2000) Methodological issues in using the Rasch model to develop a quality of life index: the analysis of four WHOQOL-100 data sets (Argentina, France, Hong Kong,, UK). *Journal of Applied Measurement*, **1** (4) 389–418.

Li H and Hong F. (2001) Cluster-Rasch models for microarray gene expression data. *Genome Biol.*, **2**(8).