# Unidimensional and Multidimensional IRT Modeling with the `mirt` Package

Phil Chalmers

York University

October 15, 2012

YORK U
UNIVERSITY
redefine THE POSSIBLE.

## Introduction

This presentation focuses on unidimensional and multidimensional item response theory (UIRT and MIRT, respectively) models that can be estimated with the `mirt` (Chalmers, 2012) package. In general, I will go over:

- What IRT is, why it exists, and how it relates to other latent variable methods such as factor analysis

- Several types of IRT models and how these can be generalized to more than one dimension

- How to fit UIRT and MIRT models to psychological test data with the `mirt` package

- Useful model comparison techniques, computing latent trait scores and item/person fit statistics, plotting item and test probability curves and information functions, and

- (time permitting) Explore some more advanced methods such as multiple group analysis for detecting DIF, user defined prior parameter distributions and starting values, linear parameter constraints, Wald tests, etc.

YORK U
UNIVERSITY
redefine THE POSSIBLE.

# Classical Test Theory

Classical test theory was largely developed by Spearman, Thurstone, Kuder, Guttman, and Cronbach, as well as a few others. In general to determine the properties of a scale the following aspects were studied (almost entirely by linear regression theory):

1) Estimating the *global* reliability of a test based on how homogeneous the items are with each other ($\alpha$, split-half), and using this to define the *global* standard error of measurement

2) Use the total score of a test as an estimate of ability/'True score' ($X = T + E$) and studying how each individual item relates to this total score

3) Determining the number of linearly related latent factors are manifested in a test (via factor analysis or structure equation modeling), and try to reduce the number of factors down to 1

# Classical Test Theory Problems

- Standard error applies to everyone in the population ($10 \pm 2$, $5 \pm 2$)
- To compare tests to each other forms must be *parallel* (equal item difficulties, same number of items, etc.)
- Individual scores are understood by comparing the person to the group (make total into $z$ or $T$-scores)
- Mixed item formats are difficult to compare (multiple choice vs true-false) and become ambiguous when combined for a total score
- Factor analysis on binary items leads to "difficulty" artifact dimensions
- Change scores cannot be meaningfully compared when initial score levels differ

YORK U
U N I V E R S I T Y
redefine THE POSSIBLE.

## Item Response Theory

- Item response theory (IRT) is a set of latent variable techniques specifically designed to model the interaction between a subject's 'ability' and item level stimuli (difficulty, guessing, etc.)
- Focus is on the pattern of responses rather than on composite variables and linear regression theory, and emphasises how responses can be thought of in probabilistic terms
- Much larger emphases on the error of measurement for *each* test subject rather than a global index of reliability/measurement error
- Widely used in educational and psychological research to study latent variable constructs other than ability (e.g., depression, personality, motivation)

Most common IRT models are still unidimensional, meaning they relate the items to only one latent trait, although multidimensional IRT models are becoming more popular

YORK U
UNIVERSITY
redefine THE POSSIBLE.
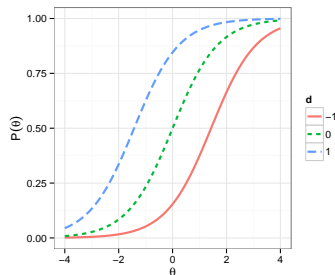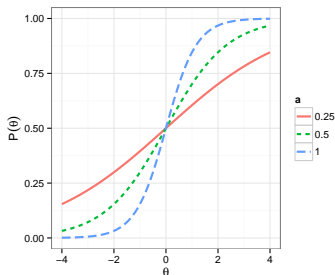
# Unidimensional IRT models (dichotomous)

Traditional IRT models were developed for modeling how a subject's 'ability' ($\theta$) was related to answering a test item correctly ($0 =$ incorrect, $1 =$ correct) given item level proprieties.

$$P(x = 1; \theta, a, d) = \frac{1}{1 + \exp\left(-D(a\theta + d)\right)}$$

This equation represents the 2 parameter logistic model (2PL). The $D$ parameter is a constant used to transform the overall metric to make the model closer to traditional factor analysis, commonly taken to be 1.702.

- Given some ability level, $\theta$, the probability of correct endorsement is related to the item easiness ($d$) and it's slope/discrimination ($a$). It may be easier to understand these relationships in the canonical form: $\log(P) \approx a\theta + d$

- This model is tied very closely to factor analysis on tetrachoric correlations, and has an analogous relationship to multiple factor analysis when the number of factors is greater than one (i.e., multidimensional)

YORK
UNIVERSITY
redefine THE POSSIBLE.

# Unidimensional plots (2PL)



Figure : Item response curves when varying the slope and intercept parameters in the 2PL model (not generated from `mirt`)

## Unidimensional IRT models (dichotomous, cont.)

Further generalization of the 2PL model are also possible to accommodate for other psychological phenomenon such as guessing or ceiling effects. For example,
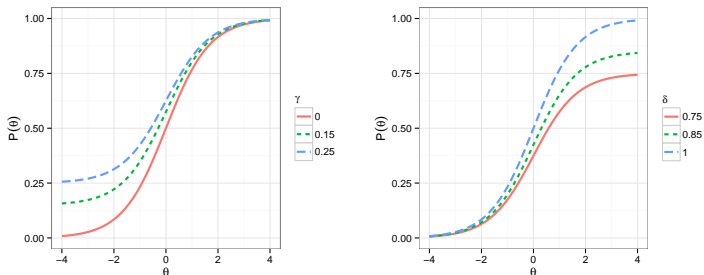
$$P(x = 1; \theta, a, d, \gamma, \delta) = \gamma + \frac{(\delta - \gamma)}{1 + \exp(-1.702(a\theta + d))}$$

This is the (maybe not so popular, but still pretty cool) four parameter logistic model, which when specific constraints are applied reduces to the 3PL, 2PL, 1PL, and Rasch model.

- Given some ability level, $\theta$, the probability of correct endorsement is related to the item easiness ($d$), discrimination ($a$), probability of randomly guessing ($\gamma$), and probability of randomly answering incorrectly ($\delta$).
- For psychological questionnaires the lower and upper bounds often have no rational and are taken to be 0 and 1, respectively (though in clinical instruments they may be justified).

# Unidimensional plots (4PL)



Figure : Item response curves when varying the lower and upper bound parameters in the 4PL model (not generated from `mirt`)
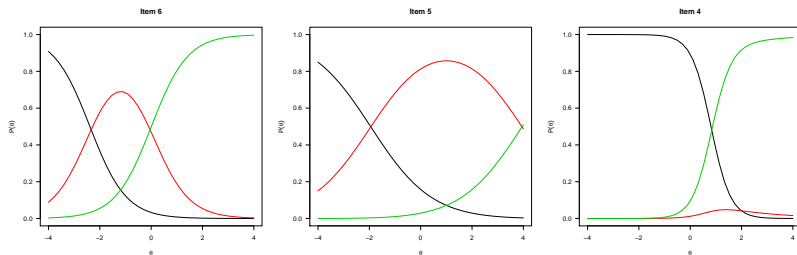
# Unidimensional IRT models (polytomous)

Several different kinds of polytomous item response models exist for ordinal, rating scale, generalized partial credit, and nominal models; all of which extend to the multidimensional case (some of which require some initially counterintuitive parameterizations). Likert scales, for example, are often modeled by ordinal or rating scale models. The ordinal/graded response model can be expressed as:

$$P(x_k = k; \theta, \phi) = P(x \geq k) - P(x \geq k+1)$$

For the generalized partial credit model the $d_k$ values are treated as fixed and ordered values from $0 : (k-1)$.

$$P(x = k; \theta, \psi) = \frac{exp(-1.702[ak_k(a\theta) + d_k])}{\sum_{j=1}^{k} exp(-1.702[ak_k(a\theta) + d_k])}$$

YORK
UNIVERSITY
redefine THE POSSIBLE.

## Unidimensional plots (polytomous)



Figure : Probability curves for ordinal (left), generalized partial credit (middle), and nominal (right) response models

# Item and test information

Item and test information are very important concepts in IRT and form the building blocks of more advanced applications such as computerized adaptive testing (CAT). The information in a test depends on the items used **as well as the ability of the subject**, and is inversely related to *reliability*. IRT advances the concept of reliability by treating it as a function of the $\theta$ values

- For example, easy items and tests tend to tell us very little about individuals in the upper end of the $\theta$ distribution ($\theta_{Einstein}$ v.s. $\theta_{Hawking}$) but can tell us something about lower ability subjects (whether $\theta_{Larry} < \theta_{Curly} < \theta_{Moe}$).

- Formally this information function (dependent on $\theta$) is defined as:

$$I(\theta) = \sum_{k=1} \left( \frac{(\partial P/\partial \theta)^2}{P} - \partial^2 P/\partial \theta \right)$$

- *Test information* is simply the sum over each item information function $T(\theta) = \sum_{i=1} I_i(\theta)$. CAT applications often stop when the information reaches a pre-specified tolerance (since $SE(\theta) = \sqrt{T(\theta)^{-1}}$). These ideas also readily generalize to multiple latent traits

# Ability estimation

Three algorithms are typically used to obtain estimates of latent trait values and their standard errors:

1) Maximum likelihood (ML) – Maximize likelihood vector w.r.t. $\theta$ directly with iterative methods. Doesn't allow for all/none patterns

2) Maximum a posteriori (MAP) – Given a prior (typically [multivariate] normal) maximize the posterior distribution. Requires iterative methods for each response pattern but works for all patterns

3) Expected a posteriori (EAP) – Similar to MAP but is not iterative and often a consequence of the estimation process (mean estimate rather than mode). Most often used method
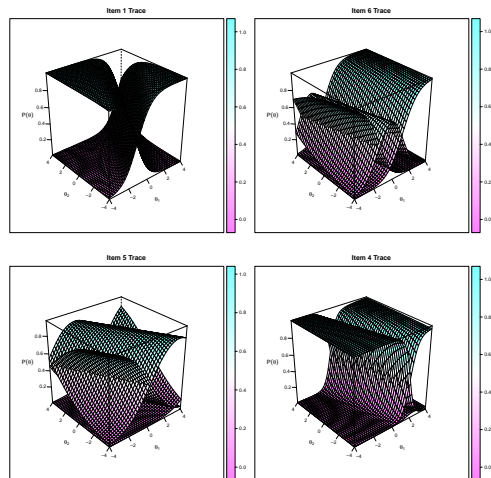
# Multidimensional IRT models

Multidimensional IRT models replace the single $\theta$ and $a$ values with vectors $\boldsymbol{\theta}$ and $\mathbf{a}$, respectively. This is analogous to the transition from zero-order regression to multiple regression (expect that the predictors are latent and non-linear).

$$P(x = 1; \boldsymbol{\theta}, \mathbf{a}, d, \gamma, \delta) = \gamma + \frac{(\delta - \gamma)}{1 + \exp\left[-1.702(\mathbf{a}'\boldsymbol{\theta} + d)\right]}.$$

This model has a very intimate relationship to nonlinear factor analysis when $\gamma = 0$ and $\delta = 1$, (since $log(P) \approx \mathbf{a}'\boldsymbol{\theta} + d$) and is often called a 'compensatory' model for the relationships between latent trait scores.

- Similar relationships exists for the generalized partial credit, graded, and nominal models, but other special types of models that don't follow these trends (e.g., partially compensatory, polynomial/exponential related traits) are also possible.

YORK U
UNIVERSITY
redefine THE POSSIBLE.

# Multidimensional plots



Figure : Probability curves for multidimensional 2PL and ordinal (top), generalized partial credit and nominal models (bottom)

# Model estimation

IRT item parameters are typically estimated by maximizing the observed likelihood

$$L(\mathbf{\Psi}; \mathbf{X}) = \prod_{i=1}^{N} \left[ \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} L_\ell(\mathbf{x}; \mathbf{\Psi}, \boldsymbol{\theta}) g(\boldsymbol{\theta}) d\boldsymbol{\theta} \right].$$

- Maximizing the above equation directly quickly becomes infeasible due to the number of parameters estimated
- Instead an EM algorithm is often employed to capitalize on a more manageable complete-data likelihood (creating artificial tables of number of participants with given response patterns)
- Effectively this approach lessens the problem of maximizing all the parameters at each iteration, but the integrals must still be evaluated

YORK U
UNIVERSITY
redefine THE POSSIBLE.

## Unfortunately . . .

Every new $\theta$ estimated requires a new integral to be evaluated in the observed likelihood.

- The difficult task is to evaluate the likelihood numerically, which requires integration by quadrature (e.g., Gauss-Hermite) or simulation methods
- Quadrature techniques often become intractable as the dimensions increase since the number of quadratures required increases exponentially
- Bayesian methods have been used to circumvent this integration problem at the cost of longer estimation times and often high computation demand

# Estimation (cont.)

An alternative approach is to capitalize on the complete-data likelihood function directly

$$L(\mathbf{\Psi}; \mathbf{X}, \boldsymbol{\theta}) = \prod_{i=1}^{N} L_\ell(\mathbf{x}_i; \mathbf{\Psi}, \boldsymbol{\theta}_i) g(\boldsymbol{\theta}_i; \mu, \mathbf{\Sigma}).$$

- What is required here is that we obtain 'known' values for $\boldsymbol{\theta}$ and maximize this function instead
- The Metropolis-Hastings Robbins-Monro (MH-RM) algorithm works well in this situation and is surprisingly fast and accurate
- MH sampler to obtain $\boldsymbol{\theta}$ values, treat values as 'known' and update parameters using standard numerical optimization methods (e.g., Newton-Raphson), and use Robbins-Monro method help remove the sampling error borne from the MH draws

### mirt package tip

I recommend using the MH-RM over the EM when the number of dimensions in the model becomes higher than 3–4

## mirt package

# mirt package

# Why the `mirt` package?

1) Multidimensional IRT functions in `R` offered limited features, were slow, and sometimes computationally demanding (e.g., `ltm`, `MCMCpack`)

2) Wanted an open source version of `TESTFACT` and `POLYFACT` which would easily integrate with useful R packages (e.g., `plink`, `GPArotation`)

3) Also wanted to utilize the MH-RM algorithm (Cai, 2010) for higher dimensional and confirmatory IRT models (analogous to confirmatory factor analysis in SEM)

4) Wanted to fit more general item response models (e.g., nominal, generalized partial credit, partially compensatory, polynomial related traits, etc.)

5) For multiple group estimation, which is important for testing the bias in testing instruments. Existed in proprietary software (even then, only in a select few) but couldn't work for MIRT models

## Functions

The mirt package consists of 4 estimation functions: mirt(), bfactor(), confmirt(), and multipleGroup(). All of these function can be used to model any mixture of dichotomous and polytomous items.

- mirt() uses a fixed quadrature estimation method (Bock & Aitkin, 1981) for obtaining ML parameter estimates with the EM algorithm. The syntax used is similar to the standard factor analysis routines in R, but also allows for confmirt.model() defined objects to be passed for confirmatory IRT models

- bfactor() uses dimension reduction algorithm for confirmatory bi-factor models described by Gibbons et al. (2007). These have the benefit of remaining computationally efficient and accurate regardless of the number of specific factors defined

- confmirt() uses the MH-RM algorithm for exploratory and confirmatory IRT models, which may also include non-compensatory item types and polynomial factor relationships

- multipleGroup() uses the MH-RM or EM algorithm to perform multiple group estimation useful for testing the invariance of parameters between potentially heterogeneous groups

# Functions (cont.)

Some useful generic functions which work on the returned estimated objects:

- coef() and summary() – extract unstandardized and standardized (i.e., factor loadings) coefficients, respectively
- plot() – two- and three-dimensional probability and information plots for item bundles
- anova() – comparison between nested models with $\chi^2$, AIC, BIC, etc.
- residuals() and fitted() – linear dependence or pattern based residuals
- itemplot() – plots individual item response curves
- fscores() – compute EAP, MAP, or ML factor scores
- itemfit() – $Z$, $\chi^2$, infit, and outfit statistics to judge item fit
- personfit() – $Z$, infit, and outfit for detecting person misfit

## Possible MIRT models

From the `mirt` documentation:

### itemtype

type of items to be modeled, declared as a vector for each item or a single value which will be repeated globally. The NULL default assumes that the items follow a graded or 2PL structure, however they may be changed to the following: 'Rasch', '1PL', '2PL', '3PL', '3PLu', '4PL', 'graded', 'grsm', 'gpcm', 'nominal', 'mcm', 'PC2PL', and 'PC3PL', for the Rasch/partial credit, 1 and 2 parameter logistic, 3 parameter logistic (lower asymptote and upper), 4 parameter logistic, graded response model, rating scale graded response model, generalized partial credit model, nominal model, multiple choice model, and 2-3PL partially compensatory model, respectively

See `?mirt` for more details.

## Running example

To demonstrate some of the features in mirt I've constructed a simple dataset of 6 items consisting of 2PL, ordinal, gpcm, and nominal item models with an orthogonal bi-factor structure (one general factor that affects all items + specific item factors that form a Thurstonian 'simple structure'). This dataset was used to generate the previous figures as well and came from the mirt function simdata().

```
> cat(itemtype)

## 2PL 2PL 2PL nominal gpcm graded

> head(dat)

##      Item_1 Item_2 Item_3 Item_4 Item_5 Item_6
## [1,]      0      0      0      1      1      2
## [2,]      1      0      0      1      1      2
## [3,]      0      1      1      1      1      3
## [4,]      0      0      0      1      1      2
## [5,]      0      1      1      1      1      3
## [6,]      1      1      1      3      1      2
```

## mirt() estimation

```
> # one factor
> mixedmod <- mirt(dat, 1, itemtype = itemtype)
> # two factor (exploratory)
> mixedmod2 <- mirt(dat, 2, itemtype = itemtype)

> mixedmod

##
## Call:
## mirt(data = dat, model = 1, itemtype = itemtype)
##
## Full-information item factor analysis with 1 factors
## Converged in 16 iterations with 40 quadrature.
## Log-likelihood = -14060
## AIC = 28151
## BIC = 28252
## G^2 = 265.3, df = 100, p = 0
## TLI = 0.899, RMSEA = 0.02
```

## Estimation times

| Subroutine | 2-factor | 3-factor | 4-factor |
|---|---|---|---|
| `mirt()` | 4.2 | 9.2 | 128.8 |
| `ltm()` | 1353.1 | — | — |
| `TESTFACT` | 9.6 | 175.3 | 946.3 |
| `confmirt()` | 117.5 | 172.9 | 202.1 |
| `MCMCirtKd()` | 2150.7 | 2368.6 | 2479.5 |

Table : Estimation times in seconds for three factor population model. See Chalmers (2012) for more detail.
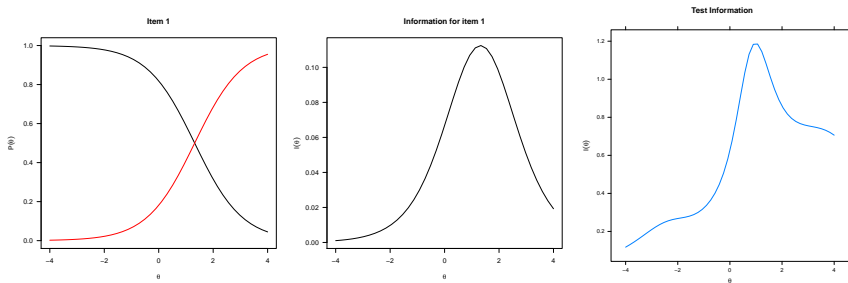
# summary()

```
> summary(mixedmod2, rotate = "oblimin", suppress = 0.3)

##
## Rotation:   oblimin
##
## Rotated factor loadings:
##
##           F_1     F_2    h2
## Item_1     NA -0.730 0.510
## Item_2     NA -0.408 0.177
## Item_3     NA -0.533 0.374
## Item_4 0.630     NA 0.450
## Item_5 0.568     NA 0.244
## Item_6 0.593     NA 0.414
##
## Rotated SS loadings:  1.084 1.007
##
## Factor correlations:
##
##         F_1    F_2
## F_1  1.000 -0.668
## F_2 -0.668  1.000
```

## plot() and itemplot()

```
> itemplot(mixedmod, item = 1)
> itemplot(mixedmod, item = 1, type = "info")
> plot(mixedmod)
```

## fscores()

EAP, MAP, and ML factor scores available for all estimated objects.

```
> tabscores <- fscores(mixedmod)

##
## Method:  EAP
##
## Empirical Reliability:
##      F1
## 0.6169

> head(tabscores)

##      Item_1 Item_2 Item_3 Item_4 Item_5 Item_6 Freq      F1  SE_F1
## [1,]      0      0      0      1      1      1   22 -1.9745 0.7306
## [2,]      0      0      0      1      1      2   45 -1.3384 0.6882
## [3,]      0      0      0      1      1      3    6 -0.7796 0.6886
## [4,]      0      0      0      1      2      1   21 -1.5224 0.7314
## [5,]      0      0      0      1      2      2   88 -0.9499 0.6678
## [6,]      0      0      0      1      2      3   49 -0.4041 0.6459
```

# residuals()

```
> residuals(mixedmod2)

## LD matrix (lower triangle) and standardized values:

##         Item_1 Item_2 Item_3 Item_4 Item_5 Item_6
## Item_1      NA  0.007  0.006  0.012  0.012  0.006
## Item_2 -0.189     NA  0.013  0.035  0.013  0.011
## Item_3  0.129  0.664     NA  0.012  0.013  0.016
## Item_4 -0.615  4.935  0.574     NA  0.016  0.015
## Item_5 -0.618  0.682  0.662 -0.968     NA  0.012
## Item_6  0.169  0.472 -1.019  0.949  0.599     NA

> head(residuals(mixedmod, restype = "exp"))

##   Item_1 Item_2 Item_3 Item_4 Item_5 Item_6 Freq   exp    res
## 1      0      0      0      1      1      1   22 17.02  1.237
## 2      0      0      0      1      1      2   45 40.90  0.684
## 3      0      0      0      1      1      3    6 11.36 -1.574
## 4      0      0      0      1      2      1   21 20.06  0.240
## 5      0      0      0      1      2      2   88 80.38  0.909
## 6      0      0      0      1      2      3   49 35.64  2.283
```

# `itemfit()` and `personfit()`

Values for detecting peculiar response patterns (e.g., someone answers all the hard questions right but easy ones wrong). Same for items, but could also also calculate a $\chi^2$ test and plot the fitted values.

```
> pfit <- personfit(mixedmod)
> print(pfit[1:3, ])

##   Item_1 Item_2 Item_3 Item_4 Item_5 Item_6  outfit  infit      Zh
## 1      0      0      0      1      1      2  6.5873 3.0444 -0.1761
## 2      1      0      0      1      1      2 11.8257 4.5174 -1.8120
## 3      0      1      1      1      1      3  0.8106 0.6045  1.0104


> ifit <- itemfit(mixedmod, X2 = TRUE)
> print(ifit[1:3, ])

##      item outfit infit      Zh df      X2
## 1 Item_1   1152 2.024  4.8399 18   86.72
## 2 Item_2   3818 6.372  0.2018 18   54.26
## 3 Item_3   1363 2.137 13.8587 18  326.51
```
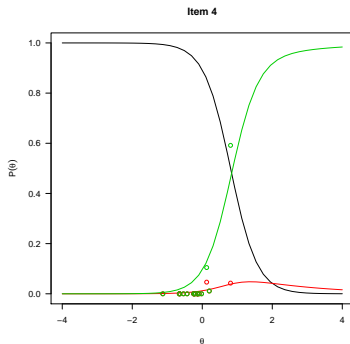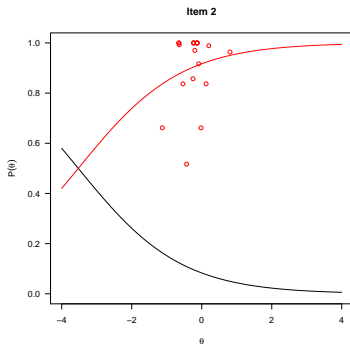
YORK U
UNIVERSITY
redefine THE POSSIBLE.

# Empirical plot

```
> itemfit(mixedmod, empirical.plot = 2)
> itemfit(mixedmod, empirical.plot = 4)
```

## bfactor() estimation

```
> # specify where the specific factor load
> sp <- c(1, 1, 1, 2, 2, 2)
> bfactor.mod <- bfactor(dat, sp, itemtype, SE = TRUE)
> coef(bfactor.mod)

## $Item_1
##          a1    a2 a3      d  g  u
## pars 0.758 0.489  0 -1.011  0  1
## SE   0.056 0.049 NA  0.054 NA NA
##
## $Item_2
##          a1    a2 a3     d  g  u
## pars 0.419 0.246  0 1.450  0  1
## SE   0.058 0.057 NA 0.065 NA NA
##
## $Item_3
##          a1    a2 a3     d  g  u
## pars 0.735 0.459  0 0.015  0  1
## SE   0.048 0.042 NA 0.038 NA NA
##
## $Item_4
##          a1 a2    a3 ak0  ak1 ak2 d0     d1    d2
## pars 0.827  0 0.446   0 1.601   2  0 -1.137 1.513
## SE   0.031 NA 0.022  NA 0.120  NA NA  0.125 0.053
##
```

# confmirt() estimation

This estimation function requires that a structural model be defined (models can also be passed to mirt(), however confmirt() can be more accurate and faster in higher dimensions)

```
> model <- confmirt.model()
+     G = 1-6
+    S1 = 1-3
+    S2 = 4-6

> conf.mod <- confmirt(dat, model, itemtype = itemtype, verbose = FALSE)
> anova(mixedmod, conf.mod)

##
## Chi-squared difference:
##
## X2 = 83.75, df = 3, p = 0
## AIC difference = 77.75
## BIC difference = 58.87
```

# Advanced features

### Advanced features

Does anybody have the time? Or the patience? Or, preferably, *both*?

## Customizing values and estimation

I've centered several methods for constraints, starting/fixed values, prior distributions, etc., on the idea of returning a `values` index to see how `mirt` codes the parameters. The data frame returned can then be modified and input back into the function, or users can observe what the parameter numbers are and apply linear constraints or prior parameter distributions.

```
> values <- mirt(dat, model, itemtype, pars = "values")
> head(values)

##    group   item name parnum  value   est
## 1    all Item_1   a1      1  0.500  TRUE
## 2    all Item_1   a2      2  0.500  TRUE
## 3    all Item_1   a3      3  0.000 FALSE
## 4    all Item_1    d      4 -1.059  TRUE
## 5    all Item_1    g      5  0.000 FALSE
## 6    all Item_1    u      6  1.000 FALSE

> # change start value
> values[1, 5] <- 1
> newmod <- mirt(dat, model, itemtype, pars = values)
```

## Constraints and prior distributions

Once the parameter index has been obtained users can use this information to impose equality constraints or give prior distributions to help control unstable parameters.

```
> #set first two slopes equal
> constrmod <- mirt(dat, model, itemtype,
+     constrain = list(c(1,7)))
>
> #normal prior on first intercept (N ~ (0,2))
> priormod <- mirt(dat, model, itemtype,
+     parprior = list(c(4, 'norm', 0, 2)))
```

YORK U
UNIVERSITY
redefine THE POSSIBLE.

# Multiple group estimation

Multiple group analysis (MGA) takes into account empirical grouping clusters that are thought to behave differently to the response data. For instance, items may be more difficult for one group or another, may have unequal slopes, etc., and these play a key role in determining the 'fairness' of a test.

- Two extremes of MGA are that all the parameters are equal across groups (equivalent to fitting any of the previous methods to all the data while ignoring group membership), or that all groups are completely independent (equivalent to sub-setting the data by group and estimating independent models)
- MGA becomes useful when models lie somewhere in the middle of these extremes, where we seek for a simpler model than strict independence while being mindful of population differences

YORK U
UNIVERSITY
redefine THE POSSIBLE.

# Multiple group estimation (cont.)

The `multipleGroup()` function begins at the strict independence end of MGA. Although it's entirely possible to declare values manually I've included a few common across group constraints such as `slopes`, `intercepts`, `free_means`, etc., that can be passed to an optional `invariance` input.

```
> #strictly independent model
> levels(group)

## [1] "D1" "D2"

> # model can also be a confmirt.model() object
> mg1 <- multipleGroup(dat, model = 1, group = group,
+      method = 'EM', verbose = FALSE)
```

YORK U
UNIVERSITY
redefine THE POSSIBLE.

# Multiple group estimation (cont.)

Equal slopes across groups (Wald test may be useful here too). *Note*: can use previously estimated models to give the current model free parameters better starting values.

```
> mg2 <- multipleGroup(dat, model = 1, group = group,
+       prev.mod = mg1, invariance = 'slopes', method = 'EM',
+       verbose = FALSE)
> anova(mg2, mg1)

##
## Chi-squared difference:
##
## X2 = 9.43, df = 6, p = 0.1508
## AIC difference = -2.57
## BIC difference = -40.33

> #models not sig diff, equal slopes accross
> #groups probably kool
```
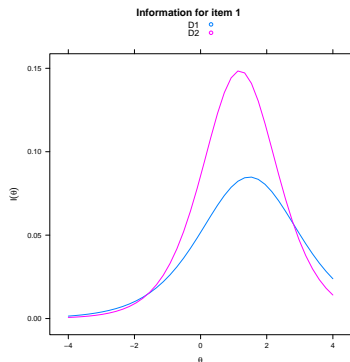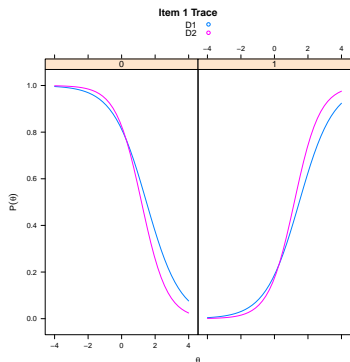
# Multiple group estimation itemplots

Superimposed item trace and information plots with each group. Also available for polytomous and two factor IRT models.

```
> itemplot(mg1, item = 1)
> itemplot(mg1, item = 1, type = "info")
```

# Future developments

This package is geared towards making complex IRT modeling accessible to those who may (or may not) be proficient with R, while still giving front end users the flexibility to explore particular models that they are comfortable with.

In the future I plan to add support for the following features:

- Parallel processing for Monte Carlo methods
- Explanatory IRT models
- Multilevel IRT models
- Imputation methods for obtaining $G^2$ and related statistics when data are missing
- Limited-information fit statistics and bootstrapped goodness of fit measures

# References

Bock, R. D., & Aitkin, M. (1981). Marginal maximum likelihood estimation of item parameters: Application of an EM algorithm. *Psychometrika*, *46*(4), 443–459.

Cai, L. (2010). High-dimensional exploratory item factor analysis by a Metropolis-Hastings Robbins-Monro algorithm. *Psychometrika*, *75*(1), 33–57. doi: 10.1007/S11336--009--9136--X

Chalmers, R. P. (2012). mirt: A multidimensional item response theory package for the R environment. *Journal of Statistical Software*, *48*(6), 1–29. Retrieved from http://www.jstatsoft.org/v48/i06

Gibbons, R. D., Darrell, R. B., Hedeker, D., Weiss, D. J., Segawa, E., Bhaumik, D. K., . . . Stover, A. (2007). Full-information item bifactor analysis of graded response data. *Applied Psychological Measurement*, *31*(1), 4–19.

YORK U
UNIVERSITY
redefine THE POSSIBLE.