

# Psychology 405: Psychometric Theory

## Reliability Theory

William Revelle

Department of Psychology  
Northwestern University  
Evanston, Illinois USA

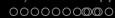


NORTHWESTERN  
UNIVERSITY

April, 2012

## Outline

- 1 Preliminaries
  - Classical test theory
  - Congeneric test theory
- 2 Reliability and internal structure
  - Estimating reliability by split halves
  - Domain Sampling Theory
  - Coefficients based upon the internal structure of a test
  - Problems with  $\alpha$
- 3 Types of reliability
  - Alpha and its alternatives
- 4 Calculating reliabilities
  - Congeneric measures
  - Hierarchical structures
- 5  $2 \neq 1$ 
  - Multiple dimensions - falsely labeled as one
  - Using score.items to find reliabilities of multiple scales
  - Intraclass correlations



## Observed Variables

X

$X_1$

$X_2$

$X_3$

$X_4$

$X_5$

$X_6$

Y

$Y_1$

$Y_2$

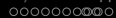
$Y_3$

$Y_4$

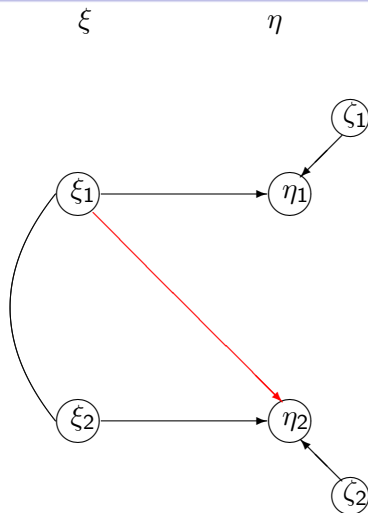
$Y_5$

$Y_6$

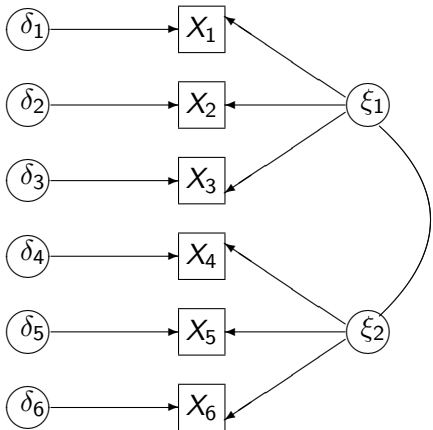




## Theory: A regression model of latent variables

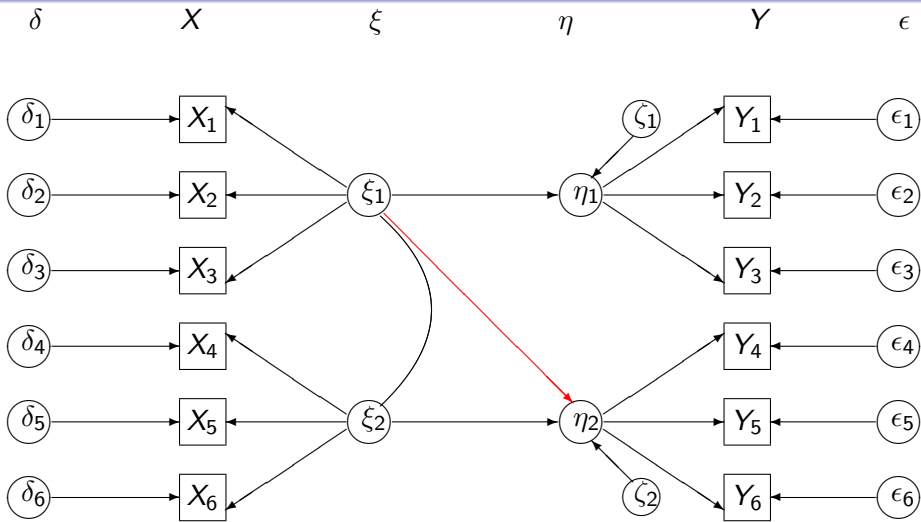


## A measurement model for $X$ – Correlated factors

 $\delta$  $X$  $\xi$ 



## A complete structural model

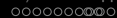






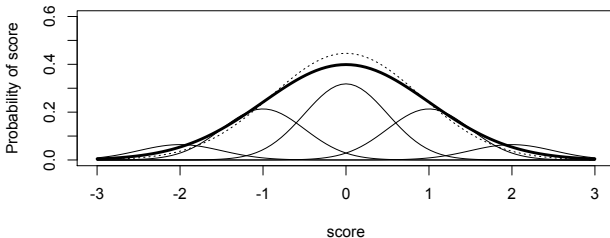
## All data are befuddled with error

*Now, suppose that we wish to ascertain the correspondence between a series of values,  $p$ , and another series,  $q$ . By practical observation we evidently do not obtain the true objective values,  $p$  and  $q$ , but only approximations which we will call  $p'$  and  $q'$ . Obviously,  $p'$  is less closely connected with  $q'$ , than is  $p$  with  $q$ , for the first pair only correspond at all by the intermediation of the second pair; the real correspondence between  $p$  and  $q$ , shortly  $r_{pq}$  has been "attenuated" into  $r_{p'q'}$  (Spearman, 1904, p 90).*

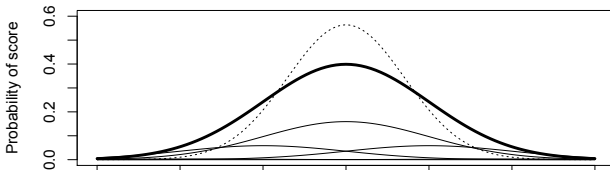


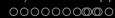
# All data are befuddled by error: Observed Score = True score + Error score

Reliability = .80



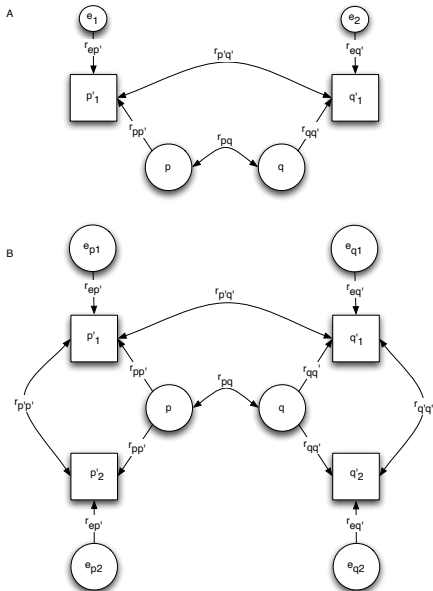
Reliability = .50





## Classical test theory

## Spearman's parallel test theory







# Classical Test Theory

By knowing the correlation between observed score and true score,  $\rho_{xt}$ , and from the definition of linear regression predicted true score,  $\hat{t}$ , for an observed  $x$  may be found from

$$\hat{t} = b_{t.x}x = \frac{\sigma_t^2}{\sigma_x^2}x = \rho_{xt}^2x. \quad (2)$$

All of this is well and good, but to find the correlation we need to know either  $\sigma_t^2$  or  $\sigma_e^2$ . The question becomes how do we find  $\sigma_t^2$  or  $\sigma_e^2$ ?



## Regression effects due to unreliability of measurement

Consider the case of air force instructors evaluating the effects of reward and punishment upon subsequent pilot performance. Instructors observe 100 pilot candidates for their flying skill. At the end of the day they reward the best 50 pilots and punish the worst 50 pilots.

- Day 1
  - Mean of best 50 pilots 1 is 75
  - Mean of worst 50 pilots is 25
- Day 2
  - Mean of best 50 has gone down to 65 ( a loss of 10 points)
  - Mean of worst 50 has gone up to 35 (a gain of 10 points)
- It seems as if reward hurts performance and punishment helps performance.
- If there is no effect of reward and punishment, what is the expected correlation from day 1 to day 2?

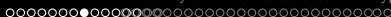


## Correcting for attenuation

*To ascertain the amount of this attenuation, and thereby discover the true correlation, it appears necessary to make two or more independent series of observations of both  $p$  and  $q$ . (Spearman, 1904, p 90)*

Spearman's solution to the problem of estimating the true relationship between two variables,  $p$  and  $q$ , given observed scores  $p'$  and  $q'$  was to introduce two or more additional variables that came to be called *parallel tests*. These were tests that had the same true score for each individual and also had equal error variances. To Spearman (1904b p 90) this required finding "the average correlation between one and another of these independently obtained series of values" to estimate the reliability of each set of measures ( $r_{p'p'}$ ,  $r_{q'q'}$ ), and then to find

$$r_{pq} = \frac{r_{p'q'}}{\sqrt{r_{p'p'}r_{q'q'}}}. \quad (3)$$



## Two parallel tests

The correlation between two parallel tests is the squared correlation of each test with true score and is the percentage of test variance that is true score variance

$$\rho_{xx} = \frac{\sigma_t^2}{\sigma_x^2} = \rho_{xt}^2. \quad (4)$$

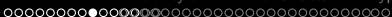
Reliability is the fraction of test variance that is true score variance. Knowing the reliability of measures of p and q allows us to correct the observed correlation between p' and q' for the reliability of measurement and to find the unattenuated correlation between p and q.

$$r_{pq} = \frac{\sigma_{pq}}{\sqrt{\sigma_p^2 \sigma_q^2}} \quad (5)$$

and

$$r_{p'q'} = \frac{\sigma_{p'q'}}{\sqrt{\sigma_{p'}^2 \sigma_{q'}^2}} = \frac{\sigma_{p+e_1' q+e_2'}}{\sqrt{\sigma_{p'}^2 \sigma_{q'}^2}} = \frac{\sigma_{pq}}{\sqrt{\sigma_p^2 \sigma_q^2}} \quad (6)$$





## Modern “Classical Test Theory”

*Reliability* is the correlation between two *parallel tests* where tests are said to be parallel if for every subject, the true scores on each test are the expected scores across an infinite number of tests and thus the same, and the true score variances for each test are the same ( $\sigma_{p'_1}^2 = \sigma_{p'_2}^2 = \sigma_{p'}^2$ ), and the error variances across subjects for each test are the same ( $\sigma_{e'_1}^2 = \sigma_{e'_2}^2 = \sigma_{e'}^2$ ) (see Figure 11), (Lord & Novick, 1968; McDonald, 1999). The correlation between two parallel tests will be

$$\rho_{p'_1 p'_2} = \rho_{p' p'} = \frac{\sigma_{p'_1 p'_2}}{\sqrt{\sigma_{p'_1}^2 \sigma_{p'_2}^2}} = \frac{\sigma_p^2 + \sigma_{pe_1} + \sigma_{pe_2} + \sigma_{e_1 e_2}}{\sigma_{p'}^2} = \frac{\sigma_p^2}{\sigma_{p'}^2}. \quad (7)$$



# Classical Test Theory

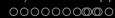
but from Eq 4,

$$\sigma_p^2 = \rho_{p'p'} \sigma_{p'}^2 \quad (8)$$

and thus, by combining equation 5 with 6 and 8 the *unattenuated correlation* between p and q corrected for reliability is Spearman's equation 3

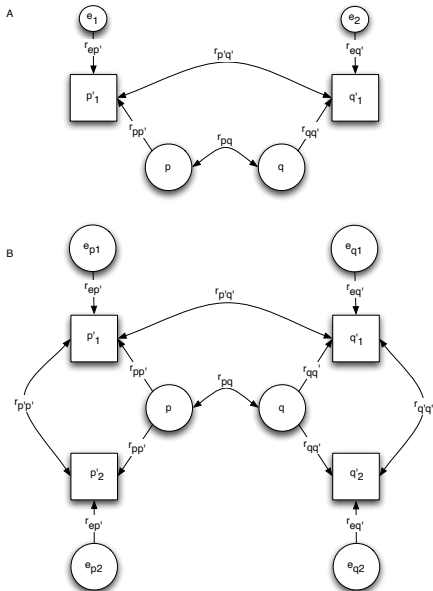
$$r_{pq} = \frac{r_{p'q'}}{\sqrt{r_{p'p'} r_{q'q'}}}. \quad (9)$$

As Spearman recognized, *correcting for attenuation* could show structures that otherwise, because of unreliability, would be hard to detect.



## Classical test theory

## Spearman's parallel test theory



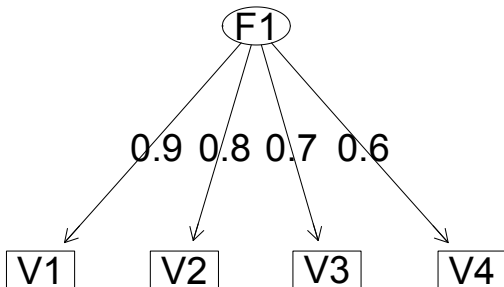


## The problem of parallel tests

Unfortunately, according to this concept of parallel tests, the possibility of one test being far better than the other is ignored. Parallel tests need to be parallel by construction or assumption and the assumption of parallelism may not be tested. With the use of more tests, however, the number of assumptions can be relaxed (for three tests) and actually tested (for four or more tests).

## Four congeneric tests – 1 latent factor

Four congeneric tests



# Observed variables and estimated parameters of a congeneric test

|    | V1       | V2       | V3       | V4      | V1                                      | V2                                      | V3                                      | V4                     |
|----|----------|----------|----------|---------|---|---|---|------------------------|
| V1 | $s_1^2$  |          |          |         | $\lambda_1 \sigma_t^2 + \sigma_{e_1}^2$ |   |   |                        |
| V2 | $s_{12}$ | $s_2^2$  |          |         | $\lambda_1 \lambda_2 \sigma_t^2$        | $\lambda_2 \sigma_t^2 + \sigma_{e_2}^2$ |   |                        |
| V3 | $s_{13}$ | $s_{23}$ | $s_3^2$  |         | $\lambda_1 \lambda_3 \sigma_t^2$        | $\lambda_2 \lambda_3 \sigma_t^2$        | $\lambda_3 \sigma_t^2 + \sigma_{e_3}^2$ |                        |
| V4 | $s_{14}$ | $s_{24}$ | $s_{34}$ | $s_4^2$ | $\lambda_1 \lambda_4 \sigma_t^2$        | $\lambda_2 \lambda_3 \sigma_t^2$        | $\lambda_3 \lambda_4 \sigma_t^2$        | $\lambda_4 \sigma_t^2$ |

## But what if we don't have three or more tests?

Unfortunately, with rare exceptions, we normally are faced with just one test, not two, three or four. How then to estimate the reliability of that one test? Defined as the correlation between a test and a test just like it, reliability would seem to require a second test. The traditional solution when faced with just one test is to consider the internal structure of that test. Letting reliability be the ratio of true score variance to test score variance (Equation 1), or alternatively, 1 - the ratio of error variance to true score variance, the problem becomes one of estimating the amount of error variance in the test. There are a number of solutions to this problem that involve examining the internal structure of the test. These range from considering the correlation between two random parts of the test to examining the structure of the items themselves.



## Split halves

$$\Sigma_{XX'} = \begin{pmatrix} \mathbf{V}_x & \vdots & \mathbf{C}_{xx'} \\ \dots\dots\dots & & \\ \mathbf{C}_{xx'} & \vdots & \mathbf{V}_{x'} \end{pmatrix} \quad (10)$$

and letting  $V_x = \mathbf{1V}_x\mathbf{1}'$  and  $C_{xx'} = \mathbf{1C}_{xx'}\mathbf{1}'$  the correlation between the two tests will be

$$\rho = \frac{C_{xx'}}{\sqrt{V_x V_{x'}}$$

But the variance of a test is simply the sum of the true covariances and the error variances:

$$V_x = \mathbf{1V}_x\mathbf{1}' = \mathbf{1C}_t\mathbf{1}' + \mathbf{1V}_e\mathbf{1}' = V_t + V_e$$



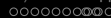
## Split halves

The split half solution estimates reliability based upon the correlation of two random split halves of a test and the implied correlation with another test also made up of two random splits:

$$\Sigma_{XX'} = \left( \begin{array}{cc|cc} \mathbf{V}_{x_1} & \vdots & \mathbf{C}_{x_1x_2} & \mathbf{C}_{x_1x'_1} & \vdots & \mathbf{C}_{x_1x'_2} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ \mathbf{C}_{x_1x_2} & \vdots & \mathbf{V}_{x_2} & \mathbf{C}_{x_2x'_1} & \vdots & \mathbf{C}_{x_2x'_2} \\ \mathbf{C}_{x_1x'_1} & \vdots & \mathbf{C}_{x_2x'_1} & \mathbf{V}_{x'_1} & \vdots & \mathbf{C}_{x'_1x'_2} \\ \mathbf{C}_{x_1x'_2} & \vdots & \mathbf{C}_{x_2x'_2} & \mathbf{C}_{x'_1x'_2} & \vdots & \mathbf{V}_{x'_2} \end{array} \right)$$







## Domain sampling

Other techniques to estimate the reliability of a single test are based on the *domain sampling* model in which tests are seen as being made up of items randomly sampled from a domain of items. Analogous to the notion of estimating characteristics of a population of people by taking a sample of people is the idea of sampling items from a universe of items.

Consider a test meant to assess English vocabulary. A person's vocabulary could be defined as the number of words in an unabridged dictionary that he or she recognizes. But since the total set of possible words can exceed 500,000, it is clearly not feasible to ask someone all of these words. Rather, consider a test of  $k$  words sampled from the larger domain of  $n$  words. What is the correlation of this test with the domain? That is, what is the correlation across subjects of test scores with their domain scores.?



## Correlation of an item with the domain

First consider the correlation of a single (randomly chosen) item with the domain. Let the domain score for an individual be  $D_i$  and the score on a particular item,  $j$ , be  $X_{ij}$ . For ease of calculation, convert both of these to deviation scores.  $d_i = D_i - \bar{D}$  and  $x_{ij} = X_{ij} - \bar{X}_j$ . Then

$$r_{x_j d} = \frac{\text{COV}_{x_j d}}{\sqrt{\sigma_{x_j}^2 \sigma_d^2}}$$

Now, because the domain is just the sum of all the items, the domain variance  $\sigma_d^2$  is just the sum of all the item variances and all the item covariances

$$\sigma_d^2 = \sum_{j=1}^n \sum_{k=1}^n \text{COV}_{x_{jk}} = \sum_{j=1}^n \sigma_{x_j}^2 + \sum_{j=1}^n \sum_{k \neq j} \text{COV}_{x_{jk}}$$

## Correlation of an item with the domain

Then letting  $\bar{c} = \frac{\sum_{j=1}^{j=n} \sum_{k \neq j} cov_{x_j k}}{n(n-1)}$  be the average covariance and

$\bar{v} = \frac{\sum_{j=1}^{j=n} \sigma_{x_j}^2}{n}$  the average item variance, the correlation of a randomly chosen item with the domain is

$$r_{x_j d} = \frac{\bar{v} + (n-1)\bar{c}}{\sqrt{\bar{v}(n\bar{v} + n(n-1)\bar{c})}} = \frac{\bar{v} + (n-1)\bar{c}}{\sqrt{n\bar{v}(\bar{v} + (n-1)\bar{c})}}.$$

Squaring this to find the squared correlation with the domain and factoring out the common elements leads to

$$r_{x_j d}^2 = \frac{(\bar{v} + (n-1)\bar{c})}{n\bar{v}}.$$

and then taking the limit as the size of the domain gets large is

$$\lim_{n \rightarrow \infty} r_{x_j d}^2 = \frac{\bar{c}}{\bar{v}}. \tag{13}$$

That is, the squared correlation of an average item with the domain is the ratio of the average interitem covariance to the average item variance. Compare the correlation of a test with true



## Domain sampling – correlation of an item with the domain

$$\lim_{n \rightarrow \infty} r_{x_j d}^2 = \frac{\bar{c}}{\bar{v}}. \quad (14)$$

That is, the squared correlation of an average item with the domain is the ratio of the average interitem covariance to the average item variance. Compare the correlation of a test with true score (Eq 4) with the correlation of an item to the domain score (Eq 14). Although identical in form, the former makes assumptions about true score and error, the latter merely describes the domain as a large set of similar items.



## Correlation of a test with the domain

A similar analysis can be done for a test of length  $k$  with a large domain of  $n$  items. A  $k$ -item test will have total variance,  $V_k$ , equal to the sum of the  $k$  item variances and the  $k(k-1)$  item covariances:

$$V_k = \sum_{i=1}^k v_i + \sum_{i=1}^k \sum_{j \neq i}^k c_{ij} = k\bar{v} + k(k-1)\bar{c}.$$

The correlation with the domain will be

$$r_{kd} = \frac{\text{cov}_k d}{\sqrt{V_k V_d}} = \frac{k\bar{v} + k(n-1)\bar{c}}{\sqrt{(k\bar{v} + k(k-1)\bar{c})(n\bar{v} + n(n-1)\bar{c})}} = \frac{k(\bar{v} + (n-1)\bar{c})}{\sqrt{nk(\bar{v} + (k-1)\bar{c})(\bar{v} + (n-1)\bar{c})}}$$

## Correlation of a test with the domain

Then the squared correlation of a  $k$  item test with the  $n$  item domain is

$$r_{kd}^2 = \frac{k(\bar{v} + (n-1)\bar{c})}{n(\bar{v} + (k-1)\bar{c})}$$

and the limit as  $n$  gets very large becomes

$$\lim_{n \rightarrow \infty} r_{kd}^2 = \frac{k\bar{c}}{\bar{v} + (k-1)\bar{c}}. \quad (15)$$

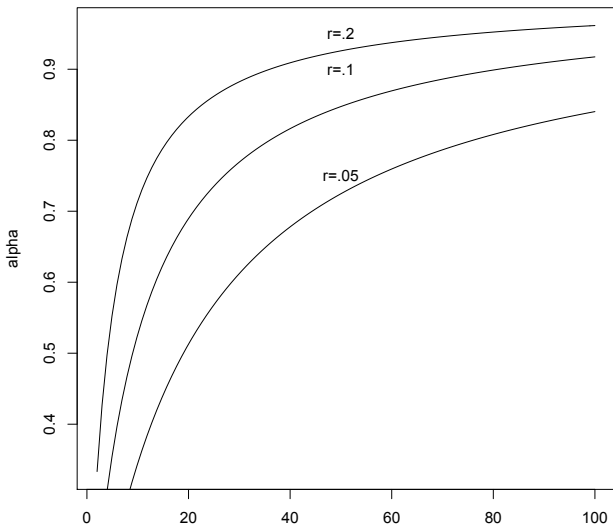
## Coefficient $\alpha$

Find the correlation of a test with a test just like it based upon the internal structure of the first test. Basically, we are just estimating the error variance of the individual items.

$$\alpha = r_{xx} = \frac{\sigma_t^2}{\sigma_x^2} = \frac{k^2 \frac{\sigma_x^2 - \sum \sigma_i^2}{k(k-1)}}{\sigma_x^2} = \frac{k}{k-1} \frac{\sigma_x^2 - \sum \sigma_i^2}{\sigma_x^2} \quad (16)$$

# Alpha varies by the number of items and the inter item correlation

Alpha varies by  $r$  and number of items



○○○○○○○○○○○○○○○○○○○○○○●○○○○○○○○○○○○

○○○○○

○○○○○○○○○○

Coefficients based upon the internal structure of a test

# Find alpha using the alpha function

```
> alpha(bfi[16:20])
```

Reliability analysis

Call: alpha(x = bfi[16:20])

| raw_alpha | std.alpha | G6(smc) | average_r | mean | sd  |
|-----------|-----------|---------|-----------|------|-----|
| 0.81      | 0.81      | 0.8     | 0.46      | 15   | 5.8 |

Reliability if an item is dropped:

|    | raw_alpha | std.alpha | G6(smc) | average_r |
|----|-----------|-----------|---------|-----------|
| N1 | 0.75      | 0.75      | 0.70    | 0.42      |
| N2 | 0.76      | 0.76      | 0.71    | 0.44      |
| N3 | 0.75      | 0.76      | 0.74    | 0.44      |
| N4 | 0.79      | 0.79      | 0.76    | 0.48      |
| N5 | 0.81      | 0.81      | 0.79    | 0.51      |

Item statistics

|    | n   | r    | r.cor | mean | sd  |
|----|-----|------|-------|------|-----|
| N1 | 990 | 0.81 | 0.78  | 2.8  | 1.5 |
| N2 | 990 | 0.79 | 0.75  | 3.5  | 1.5 |
| N3 | 997 | 0.79 | 0.72  | 3.2  | 1.5 |
| N4 | 996 | 0.71 | 0.60  | 3.1  | 1.5 |
| N5 | 992 | 0.67 | 0.52  | 2.9  | 1.6 |

○○○○○○○○○○○○○○○○○○○○○○○○●○○○○○○○○○○

○○○○○

○○○○○○○○○○

Coefficients based upon the internal structure of a test

## What if items differ in their direction?

```
> alpha(bfi[6:10], check.keys=FALSE)
```

Reliability analysis

```
Call: alpha(x = bfi[6:10], check.keys = FALSE)
```

```
raw_alpha std.alpha G6(smc) average_r mean sd
-0.28      -0.22      0.13      -0.038  3.8 0.58
```

Reliability if an item is dropped:

```
raw_alpha std.alpha G6(smc) average_r
C1      -0.430      -0.472      -0.020      -0.0871
C2      -0.367      -0.423      -0.017      -0.0803
C3      -0.263      -0.295       0.094      -0.0604
C4      -0.022       0.123       0.283       0.0338
C5      -0.028       0.022       0.242       0.0057
```

Item statistics

```
      n      r r.cor  r.drop mean sd
C1 2779 0.56  0.51  0.0354  4.5 1.2
C2 2776 0.54  0.51 -0.0076  4.4 1.3
C3 2780 0.48  0.27 -0.0655  4.3 1.3
C4 2774 0.20 -0.34 -0.2122  2.6 1.4
C5 2784 0.29 -0.19 -0.1875  3.3 1.6
```

oooooooooooooooooooooooooooo●oooooooooooo

ooooo

ooooooooooooo

Coefficients based upon the internal structure of a test

## But what if some items are reversed keyed?

```
alpha(bfi[6:10])
```

Reliability analysis

```
Call: alpha(x = bfi[6:10])
```

```
raw_alpha std.alpha G6(smc) average_r mean sd
      0.73      0.73      0.69      0.35 3.8 0.58
```

Reliability if an item is dropped:

```
raw_alpha std.alpha G6(smc) average_r
C1      0.69      0.70      0.64      0.36
C2      0.67      0.67      0.62      0.34
C3      0.69      0.69      0.64      0.36
C4-     0.65      0.66      0.60      0.33
C5-     0.69      0.69      0.63      0.36
```

Item statistics

```
      n      r r.cor r.drop mean sd
C1 2779 0.67 0.54 0.45 4.5 1.2
C2 2776 0.71 0.60 0.50 4.4 1.3
C3 2780 0.67 0.54 0.46 4.3 1.3
C4- 2774 0.73 0.64 0.55 2.6 1.4
C5- 2784 0.68 0.57 0.48 3.3 1.6
```

```
Warning message: In alpha(bfi[6:10]) :
```

Some items were negatively correlated with total scale and were automatically



## Guttman's alternative estimates of reliability

Reliability is amount of test variance that is not error variance. But what is the error variance?

$$r_{xx} = \frac{V_x - V_e}{V_x} = 1 - \frac{V_e}{V_x}. \quad (17)$$

$$\lambda_1 = 1 - \frac{\text{tr}(\mathbf{V}_x)}{V_x} = \frac{V_x - \text{tr}(\mathbf{V}_x)}{V_x}. \quad (18)$$

$$\lambda_2 = \lambda_1 + \frac{\sqrt{\frac{n}{n-1} C_2}}{V_x} = \frac{V_x - \text{tr}(\mathbf{V}_x) + \sqrt{\frac{n}{n-1} C_2}}{V_x}. \quad (19)$$

$$\lambda_3 = \lambda_1 + \frac{\frac{V_x - \text{tr}(\mathbf{V}_x)}{n(n-1)}}{V_x} = \frac{n\lambda_1}{n-1} = \frac{n}{n-1} \left(1 - \frac{\text{tr}(\mathbf{V}_x)}{V_x}\right) = \frac{n}{n-1} \frac{V_x - \text{tr}(\mathbf{V}_x)}{V_x} = \alpha \quad (20)$$

$$\lambda_4 = 2 \left(1 - \frac{V_{X_a} + V_{X_b}}{V_x}\right) = \frac{4c_{ab}}{V_x} = \frac{4c_{ab}}{V_{X_a} + V_{X_b} + 2c_{ab} V_{X_a} V_{X_b}}. \quad (21)$$

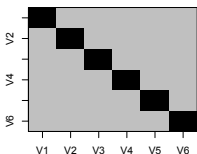
$$\lambda_6 = 1 - \frac{\sum e_j^2}{V_x} = 1 - \frac{\sum (1 - r_{smc}^2)}{V_x} \quad (22)$$



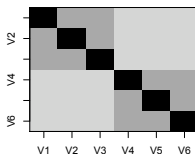
Problems with  $\alpha$

## Four different correlation matrices, one value of $\alpha$

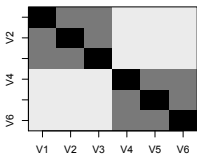
S1: no group factors



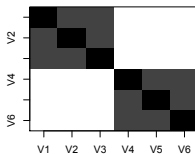
S2: large g, small group factors



S3: small g, large group factors



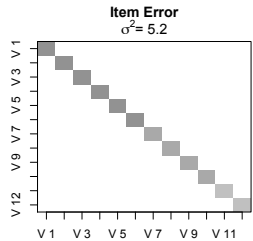
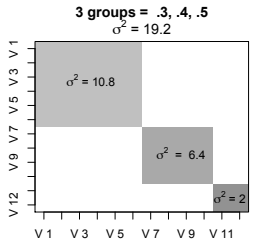
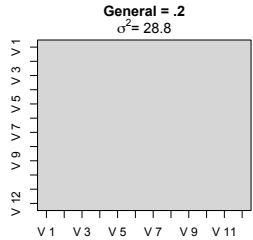
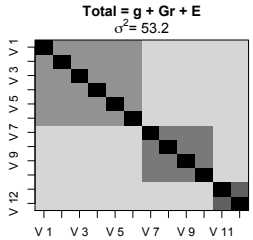
S4: no g but large group factors



- 1 The problem of group factors
- 2 If no groups, or many groups,  $\alpha$  is ok

Problems with  $\alpha$

# Decomposing a test into general, Group, and Error variance



- 1 Decompose total variance into general, group, specific, and error
- 2  $\alpha < \text{total}$
- 3  $\alpha > \text{general}$

## Two additional alternatives to $\alpha$ : $\omega_{\text{hierarchical}}$ and $\omega_{\text{total}}$

If a test is made up of a general, a set of group factors, and specific as well as error:

$$\mathbf{x} = \mathbf{c}\mathbf{g} + \mathbf{A}\mathbf{f} + \mathbf{D}\mathbf{s} + \mathbf{e} \quad (23)$$

then the communality of item $_j$ , based upon general as well as group factors,

$$h_j^2 = c_j^2 + \sum f_{ij}^2 \quad (24)$$

and the unique variance for the item

$$u_j^2 = \sigma_j^2(1 - h_j^2) \quad (25)$$

may be used to estimate the test reliability.

$$\omega_t = \frac{\mathbf{1}\mathbf{c}\mathbf{c}'\mathbf{1}' + \mathbf{1}\mathbf{A}\mathbf{A}'\mathbf{1}'}{V_x} = 1 - \frac{\sum(1 - h_j^2)}{V_x} = 1 - \frac{\sum u^2}{V_x} \quad (26)$$

## McDonald (1999) introduced two different forms for $\omega$

$$\omega_t = \frac{\mathbf{1cc}'\mathbf{1}' + \mathbf{1AA}'\mathbf{1}'}{V_x} = 1 - \frac{\sum(1 - h_j^2)}{V_x} = 1 - \frac{\sum u^2}{V_x} \quad (27)$$

and

$$\omega_h = \frac{\mathbf{1cc}'\mathbf{1}}{V_x} = \frac{(\sum \Lambda_i)^2}{\sum \sum R_{ij}} \quad (28)$$

These may both be found by factoring the correlation matrix and finding the g and group factor loadings using the omega function.

# Using omega on the Thurstone data set to find alternative reliability estimates

```
> lower.mat(Thurstone)
> omega(Thurstone)
```

|                 | Sntnc | Vcblr | Snt.C | Frs.L | 4.L.W | Sffxs | Ltt.S | Pdgrs | Ltt.G |
|-----------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| Sentences       | 1.00  |       |       |       |       |       |       |       |       |
| Vocabulary      | 0.83  | 1.00  |       |       |       |       |       |       |       |
| Sent.Completion | 0.78  | 0.78  | 1.00  |       |       |       |       |       |       |
| First.Letters   | 0.44  | 0.49  | 0.46  | 1.00  |       |       |       |       |       |
| 4.Letter.Words  | 0.43  | 0.46  | 0.42  | 0.67  | 1.00  |       |       |       |       |
| Suffixes        | 0.45  | 0.49  | 0.44  | 0.59  | 0.54  | 1.00  |       |       |       |
| Letter.Series   | 0.45  | 0.43  | 0.40  | 0.38  | 0.40  | 0.29  | 1.00  |       |       |
| Pedigrees       | 0.54  | 0.54  | 0.53  | 0.35  | 0.37  | 0.32  | 0.56  | 1.00  |       |
| Letter.Group    | 0.38  | 0.36  | 0.36  | 0.42  | 0.45  | 0.32  | 0.60  | 0.45  | 1.00  |

Omega

Call: omega(m = Thurstone)

Alpha: 0.89

G.6: 0.91

Omega Hierarchical: 0.74

Omega H asymptotic: 0.79

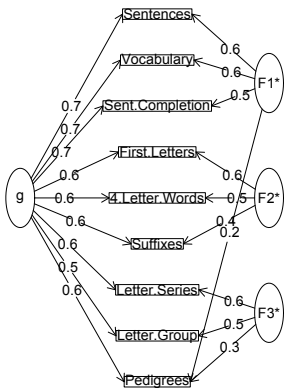
Omega Total 0.93



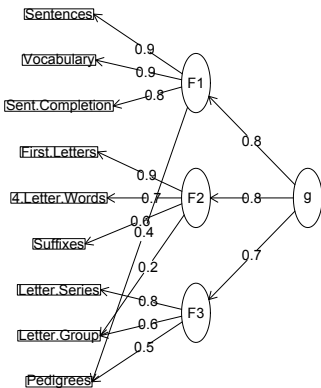
Problems with  $\alpha$

## Two ways of showing a general factor

Omega



Omega



oo●ooo

ooooo

ooooooooooooo

Problems with  $\alpha$

## omega function does a Schmid Leiman transformation

```
> omega(Thurstone,sl=FALSE)
```

Omega

```
Call: omega(m = Thurstone, sl = FALSE)
```

```
Alpha:                0.89
```

```
G.6:                  0.91
```

```
Omega Hierarchical:  0.74
```

```
Omega H asymptotic:  0.79
```

```
Omega Total          0.93
```

```
Schmid Leiman Factor loadings greater than 0.2
```

|                 | g    | F1*  | F2*  | F3*  | h2   | u2   | p2   |
|-----------------|------|------|------|------|------|------|------|
| Sentences       | 0.71 | 0.57 |      |      | 0.82 | 0.18 | 0.61 |
| Vocabulary      | 0.73 | 0.55 |      |      | 0.84 | 0.16 | 0.63 |
| Sent.Completion | 0.68 | 0.52 |      |      | 0.73 | 0.27 | 0.63 |
| First.Letters   | 0.65 |      | 0.56 |      | 0.73 | 0.27 | 0.57 |
| 4.Letter.Words  | 0.62 |      | 0.49 |      | 0.63 | 0.37 | 0.61 |
| Suffixes        | 0.56 |      | 0.41 |      | 0.50 | 0.50 | 0.63 |
| Letter.Series   | 0.59 |      |      | 0.61 | 0.72 | 0.28 | 0.48 |
| Pedigrees       | 0.58 | 0.23 |      | 0.34 | 0.50 | 0.50 | 0.66 |
| Letter.Group    | 0.54 |      |      | 0.46 | 0.53 | 0.47 | 0.56 |

```
With eigenvalues of:
```

| g    | F1*  | F2*  | F3*  |
|------|------|------|------|
| 3.58 | 0.96 | 0.74 | 0.71 |





## Types of reliability

- Internal consistency
    - $\alpha$
    - $\omega_{hierarchical}$
    - $\omega_{total}$
    - $\beta$
  - Intraclass
  - Agreement
  - Test-retest, alternate form
  - Generalizability
- Internal consistency
    - alpha,  
score.items
    - omega
    - iclust
  - icc
  - wkappa,  
cohen.kappa
  - cor
  - aov



## Alpha and its alternatives

- Reliability =  $\frac{\sigma_t^2}{\sigma_x^2} = 1 - \frac{\sigma_e^2}{\sigma_x^2}$
- If there is another test, then  $\sigma_t = \sigma_{t_1 t_2}$  (covariance of test  $X_1$  with test  $X_2 = C_{xx}$ )
- But, if there is only one test, we can *estimate*  $\sigma_t^2$  based upon the observed covariances within test 1
- How do we find  $\sigma_e^2$  ?
- The worst case, (Guttman case 1) all of an item's variance is error and thus the error variance of a test  $X$  with variance-covariance  $C_x$ 
  - $C_x = \sigma_e^2 = \text{diag}(C_x)$
  - $\lambda_1 = \frac{C_x - \text{diag}(C_x)}{C_x}$
- A better case (Guttman case 3,  $\alpha$ ) is that that the average covariance between the items on the test is the same as the average true score variance for each item.
  - $C_x = \sigma_e^2 = \text{diag}(C_x)$
  - $\lambda_3 = \alpha = \lambda_1 * \frac{n}{n-1} = \frac{(C_x - \text{diag}(C_x)) * n / (n-1)}{C_x}$







oo

oo●oo

oooooooooooo

Congeneric measures

## Find $\alpha$ and related stats for the simulated data

```
> alpha(v4$observed)
```

Reliability analysis

Call: alpha(x = v4\$observed)

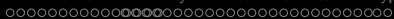
| raw_alpha | std.alpha | G6(smc) | average_r | mean   | sd   |
|-----------|-----------|---------|-----------|--------|------|
| 0.71      | 0.72      | 0.67    | 0.39      | -0.036 | 0.72 |

Reliability if an item is dropped:

|    | raw_alpha | std.alpha | G6(smc) | average_r |
|----|-----------|-----------|---------|-----------|
| V1 | 0.59      | 0.60      | 0.50    | 0.33      |
| V2 | 0.63      | 0.64      | 0.55    | 0.37      |
| V3 | 0.65      | 0.66      | 0.59    | 0.40      |
| V4 | 0.72      | 0.72      | 0.64    | 0.46      |

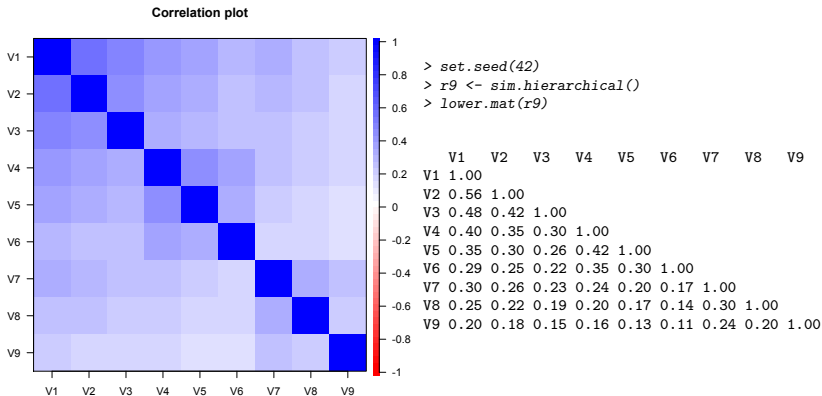
Item statistics

|    | n   | r    | r.cor | r.drop | mean   | sd   |
|----|-----|------|-------|--------|--------|------|
| V1 | 200 | 0.80 | 0.72  | 0.60   | -0.015 | 0.93 |
| V2 | 200 | 0.76 | 0.64  | 0.53   | -0.060 | 0.98 |
| V3 | 200 | 0.73 | 0.59  | 0.50   | -0.119 | 0.92 |
| V4 | 200 | 0.66 | 0.46  | 0.40   | 0.049  | 1.09 |



## A hierarchical structure

`cor.plot(r9)`



## $\alpha$ of the 9 hierarchical variables

```
> alpha(r9)
```

```
Reliability analysis
```

```
Call: alpha(x = r9)
```

| raw_alpha | std.alpha | G6(smc) | average_r |
|-----------|-----------|---------|-----------|
| 0.76      | 0.76      | 0.76    | 0.26      |

```
Reliability if an item is dropped:
```

|    | raw_alpha | std.alpha | G6(smc) | average_r |
|----|-----------|-----------|---------|-----------|
| V1 | 0.71      | 0.71      | 0.70    | 0.24      |
| V2 | 0.72      | 0.72      | 0.71    | 0.25      |
| V3 | 0.74      | 0.74      | 0.73    | 0.26      |
| V4 | 0.73      | 0.73      | 0.72    | 0.25      |
| V5 | 0.74      | 0.74      | 0.73    | 0.26      |
| V6 | 0.75      | 0.75      | 0.74    | 0.27      |
| V7 | 0.75      | 0.75      | 0.74    | 0.27      |
| V8 | 0.76      | 0.76      | 0.75    | 0.28      |
| V9 | 0.77      | 0.77      | 0.76    | 0.29      |

```
Item statistics
```

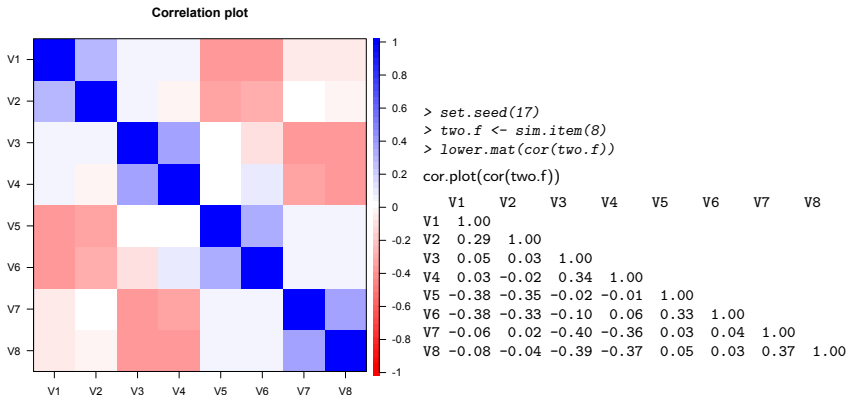
|    | r    | r.cor |
|----|------|-------|
| V1 | 0.72 | 0.71  |
| V2 | 0.73 | 0.72  |





Multiple dimensions - falsely labeled as one

## An example of two different scales confused as one







## Score as two different scales

First, make up a keys matrix to specify which items should be scored, and in which way

```
> keys <- make.keys(nvars=8,keys.list=list(one=c(1,2,-5,-6),two=c(3,4,-7,-8)))  
> keys  
      one two  
[1,]  1  0  
[2,]  1  0  
[3,]  0  1  
[4,]  0  1  
[5,] -1  0  
[6,] -1  0  
[7,]  0 -1  
[8,]  0 -1
```

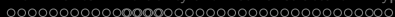
Using score.items to find reliabilities of multiple scales

Now score the two scales and find  $\alpha$  and other reliability estimates

```

> score.items(keys,two.f)
Call: score.items(keys = keys, items = two.f)
(Unstandardized) Alpha:
  one two
alpha 0.68 0.7
Average item correlation:
  one two
average.r 0.34 0.37
Guttman 6* reliability:
  one two
Lambda.6 0.62 0.64
Scale intercorrelations corrected for attenuation
raw correlations below the diagonal, alpha on the diagonal
corrected correlations above the diagonal:
  one two
one 0.68 0.08
two 0.06 0.70
Item by scale correlations:
corrected for item overlap and scale reliability
  one two
V1 0.57 0.09
V2 0.52 0.01
V3 0.09 0.59
V4 -0.02 0.56
V5 -0.58 -0.05
V6 -0.57 -0.05
V7 -0.05 -0.58
V8 -0.09 -0.59

```



## 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.obs)
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



