

Two aspects: • Reliability • Validity

Reliability (Reproducibility, Precision)

Extent to which obtain the same answer/value/score if object/subject is measured repeatedly under similar situations...

Some ways to quantify Reliability:

- For one subject: average variation of individual measurements around their mean... either the square root of the average of squared deviations) i.e. standard deviation (SD); or the average absolute deviation, which will usually be quite close to the SD. Could also use range or other measures such as Inter Quartile Range
- For one subject: average variation or SD as % of the mean of the measurements for that subject...called the {within-subject} Coefficient of Variation (CV) if calculate it as $[SD/Mean] \times 100$.
- For several subjects: : average the the CV's calculated for the different subjects; if CV's are highly variable, may want to give some sense of this using the range or other measure of spread of the CV's.

Unfortunately, CV gives no sense of how well the measurements of different subjects (ss) segregate from each other

How about

$$\frac{SD \text{ of within-ss measurements}}{SD \text{ of between-ss measurements}} \quad ?? \text{ see last item below*}$$

- Correlation (Pearson or Spearman) if 2 assessments of each ss. ??

- Using correlation between scores on random halves of a test, can estimate how 'reproducible' the full test is (helpful if cannot repeat the test)
- If the measurement in question concerns a population (eg the percentage of smokers among Canadian adults) and if it is measured (estimated) using a statistic: e.g. the proportion in a random sample of 1000 adults, it is possible from statistical laws concerning averages to quantify the reliability of the statistic without having to actually perform repeated measurements (samples). For simple random sampling, the formula

$$SE[\text{average}] = \frac{SD[\text{individuals}]}{\sqrt{\text{number of individuals measured}}}$$

allows us to quantify the reliability indirectly. If we didn't know this formula, we could also arrive at an answer by various re-sampling methods applied to the individuals in the sample at hand -- again without resorting to observing any additional individuals.

- * Some function of Variance of Within-ss measurements and Variance of Between-ss values? ? Estimate these COMPONENTS OF VARIANCE USING Analysis of Variance (ANOVA)

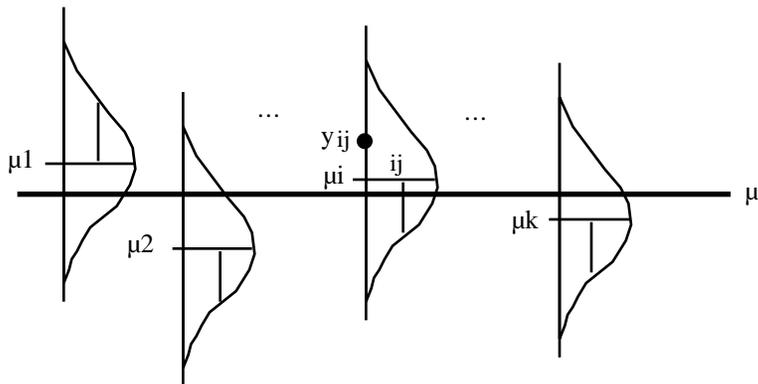
First, a General Orientation to ANOVA and its primary use, namely testing differences between μ 's of k (≥ 2) different groups.

E.g. 1-way ANOVA:

DATA:

	Group					
	1	2	.	i	.	k
Subject						
1	y ₁₁
2
.
j	.	.	.	y _{ij}	.	.
.
n	y _{kn}
Mean	\bar{y}_1	\bar{y}_2	.	\bar{y}_i	.	\bar{y}_k
Variance	s^2_1	s^2_2	.	.	.	s^2_k

MODEL



refers to the variation (SD) of all possible individuals in a group; It is an (unknowable) parameter; it can only be ESTIMATED.

Or, in symbols...

$$y_{ij} = \mu_i + e_{ij} = \mu + (\mu_i - \mu) + e_{ij}$$

DE-COMPOSITION OF OBSERVED (EMPIRICAL) VARIATION

$$(\bar{y}_{ij} - \bar{y})^2 = (\bar{y}_i - \bar{y})^2 + (\bar{y}_{ij} - \bar{y}_i)^2$$

TOTAL Sum of Squares = BETWEEN Groups + WITHIN Group Sum of Squares

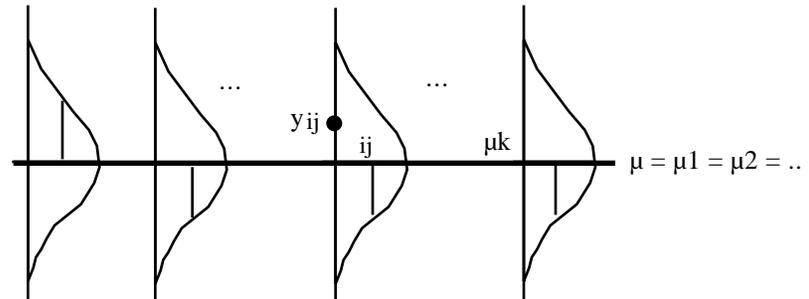
ANOVA TABLE

	Sum of Squares	Degrees of Freedom	Mean Square	F Ratio	P-Value
SOURCE	SS	df	MS (= SS / df)	$\frac{MS_{BETWEEN}}{MS_{WITHIN}}$	Prob(>F)
BETWEEN	xx.x	k-1	xx.x	x.xx	0.xx
WITHIN	xx.x	k(n-1)	xx.x		

LOGIC FOR F-TEST (Ratio of variances) as a test of

$$H_0: \mu_1 = \mu_2 = \dots = \mu_i = \dots = \mu_k$$

UNDER H0



Means, based on samples of n , should vary around μ with a variance of $\frac{2}{n}$

Thus, if H_0 is true, and we calculate the empirical variance of the k different \bar{y}_i 's, it should give us an unbiased estimate of $\frac{2}{n}$

i.e. $\frac{[\bar{y}_i - \bar{y}]^2}{k-1}$ is an unbiased estimate of $\frac{2}{n}$

i.e. $\frac{n}{k-1} \frac{[\bar{y}_i - \bar{y}]^2}{k-1}$ is an unbiased estimate of 2

i.e. $\frac{[\bar{y}_i - \bar{y}]^2}{k-1} = MS_{\text{BETWEEN}}$ is an unbiased estimate of 2

Whether or not H_0 is true, the empirical variance of the n (within-group) values

y_{i1} to y_{in} i.e. $\frac{[\bar{y}_{ij} - \bar{y}_i]^2}{n-1}$ should give us an unbiased estimate of 2

i.e. $s^2_i = \frac{[\bar{y}_{ij} - \bar{y}_i]^2}{n-1}$ is an unbiased estimate of 2

so the average of the k different estimates,

$$\frac{1}{k} s^2_i = \frac{1}{k} \frac{[\bar{y}_{ij} - \bar{y}_i]^2}{n-1}$$

is also an unbiased estimate of 2

i.e. $\frac{[\bar{y}_{ij} - \bar{y}_i]^2}{k[n-1]} = MS_{\text{WITHIN}}$ is an unbiased estimate of 2

THUS, under H_0 , both MS_{BETWEEN} and MS_{WITHIN} are unbiased estimates of estimates of 2 and so their ratio should, apart from sampling variability, be 1. IF however, H_0 is not true, MS_{BETWEEN} will tend to be larger than MS_{WITHIN} , since it contains an extra contribution that is proportional to how far the μ 's are from each other.

In this "non-null" case, the MS_{BETWEEN} is an unbiased estimate of

$$2 + \frac{n[\mu_i - \bar{\mu}]^2}{k-1}$$

and so we expect that, apart from sampling variability, the ratio $\frac{MS_{\text{BETWEEN}}}{MS_{\text{WITHIN}}}$ should be greater than 1. The tabulated values of the F distribution (tabulated under the assumption that the numerator and denominator of the ratio are both estimates of the same quantity) can thus be used to assess how extreme the observed F ratio is and to assess the evidence against the H_0 that the μ 's are equal.

How ANOVA can be used to estimate Components of Variance used in quantifying Reliability.

The basic ANOVA calculations are the same, but the MODEL underlying them is different. First, in the more common use of ANOVA just described, the groups can be thought of as all the levels of the factor of interest. The number of levels is necessarily finite. The groups might be the two genders, all of the age groups, the 4 blood groups, etc. Moreover, when you publish the results, you explicitly identify the groups.

When we come to study subjects, and ask "How big is the intra-subject variation compared with the inter-subject variation, we will for budget reasons only study a sample of all the possible subjects of interest. We can still number them 1 to k , and we can make n measurements on each subject, so the basic layout of the data doesn't change. All we do is replace the word 'Group' by 'Subject' and speak of BETWEEN-SUBJECT and WITHIN-SUBJECT variation. So the data layout is...

DATA:

	Subject					
	1	2	.	i	.	k
Measurement						
1	y_{11}
2
j	.	.	.	y_{ij}	.	.
.
n	y_{kn}
Mean	\bar{y}_1	\bar{y}_2	.	\bar{y}_i	.	\bar{y}_k
Variance	s^2_1	s^2_2	.	.	.	s^2_k

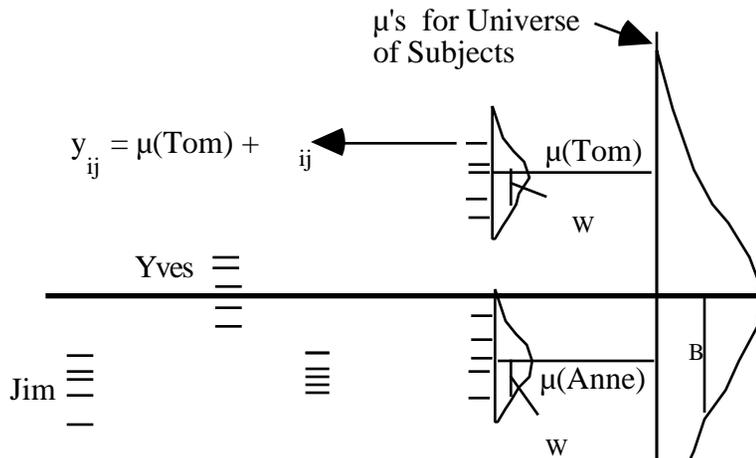
MODEL

The model is different. There is no interest in the specific subjects. Unlike the critical labels "male" and "female", or "smokers", "nonsmokers" and "exsmokers" to identify groups of interest, we certainly are not going to identify subjects as Yves, Claire, Jean, Anne, Tom, Jim, and Harry in the publication, and nobody would be fussed if in the dataset we used arbitrary subject identifiers to keep track of which measurements were made on whom. we wouldn't even care if the research assistant lost the identities of the subjects -- as long as we know that the correct measurements go with the correct subject!

The "Random Effects" Model uses 2 stages:

- (1) random sample of subjects, each with his/her own μ
- (2) For each subject, series of random variations around his/her μ

Notice the diagram has considerable 'segregation' of the measurements on different individuals. There is no point in TESTING for (inter-subject) differences in the μ 's. The task is rather to estimate the relative magnitudes of the two variance components σ_B^2 and σ_W^2 .



B refers to the SD of the universe of μ 's ; It is an unknowable parameter and can only be ESTIMATED

w refers to the variation (SD) of all possible measurements on a subject
It is an (unknowable) parameter; it can only be ESTIMATED.

Or, in symbols...

$$y_{ij} = \mu_i + e_{ij} = \mu + (\mu_i - \mu) + e_{ij}$$

$$= \mu + \mu_i + e_{ij}$$

$$\mu_i \sim N(0, \sigma_B^2)$$

$$e_{ij} \sim N(0, \sigma_W^2)$$

DE-COMPOSITION OF OBSERVED (EMPIRICAL) VARIATION

$$(\bar{y}_{ij} - \bar{y})^2 = (\bar{y}_i - \bar{y})^2 + (\bar{y}_{ij} - \bar{y}_i)^2$$

TOTAL Sum of Squares = BETWEEN Subjects Sum of Squares + WITHIN Subjects Sum of Squares

ANOVA TABLE (Note absence of F and P-value Columns)

SOURCE	Sum of Squares (SS)	Degrees of Freedom (df)	Mean Square (MS) (= SS / df)	What the Mean Square is an estimate of*
BETWEEN Subjects	xx.x	k-1	xx.x	$\sigma_W^2 + n \sigma_B^2$
WITHIN Subjects	xx.x	k(n-1)	xx.x	σ_W^2

ACTUAL ESTIMATION OF 2 Variance Components

$MS_{BETWEEN}$ is an unbiased estimate of $\sigma_W^2 + n \sigma_B^2$

MS_{WITHIN} is an unbiased estimate of σ_W^2

By subtraction...

$MS_{BETWEEN} - MS_{WITHIN}$ is an unbiased estimate of $n \sigma_B^2$

$\frac{MS_{BETWEEN} - MS_{WITHIN}}{n}$ is an unbiased estimate of σ_B^2

This is the **definitional** formula; the **computational** formula may be different.

* Pardon my ending with a preposition, but I find it difficult to say otherwise. These parameter combinations are also called the "Expected Mean Squares". They are the long-run expectations of the MS statistics. As Winston Churchill would say, "For the sake of clarity, this one time this wording is something up which you would put".

Example....

Measurement	Subject					Variance = 0.614
	Tom	Anne	Yves	Jean	Claire	
1	4.8	5.5	5.1	6.4	5.8	4.5
2	4.7	5.2	4.9	6.2	6.3	4.1
3	4.9	5.2	5.3	6.6	5.6	4.0
Mean	4.8	5.3	5.1	6.4	5.9	4.2
Variance	0.01	0.03	0.04	0.04	0.13	0.07

ANOVA TABLE (Check... I did it by hand!)

SOURCE	Sum of Squares	Degrees of Freedom	Mean Square	What the Mean Square is an estimate of... *
	SS	df	MS (= SS /df)	
BETWEEN Subjects	9.205	5	1.841	$\sigma^2_W + n \sigma^2_B$
WITHIN Subjects	0.640	12	0.053	σ^2_W
TOTAL	9.845	17		

ESTIMATES OF VARIANCE COMPONENTS

$MS_{WITHIN} = 0.053$ is an unbiased estimate of σ^2_W

$\frac{1.841 - 0.053}{3} = 0.596$ is an unbiased estimate of σ^2_B

1-Way ANOVA Calculations performed by SAS; Components estimated manually

```
PROC GLM in SAS ==> estimating components 'by hand'
```

DATA a; INPUT Subject Value; LINES;

1 4.8
1 4.7

...
6 4.5

```
proc glm; class subject; model value=subject / ss3;
random subject ;
```

See worked example using earsize data.

If unequal numbers of measurements per subject, see formula in A&B or Fleiss

Estimating Components of Variance using "Black Box"

```
PROC VARCOMP; class subject ; model Value = Subject ;
See worked example following...
```

2 measurements (in mm) of earsize of 8 subjects by each of 4 observers

subject	1				2				3				4			
obsr	1	2	3	4	1	2	3	4	1	2	3	4	1	2	3	4
1st	67	65	65	64	74	74	74	72	67	68	66	65	65	65	65	65
2nd	67	66	66	66	74	73	71	73	68	67	68	67	64	65	65	64

subject	5				6				7				8			
obsr	1	2	3	4	1	2	3	4	1	2	3	4	1	2	3	4
1st	65	62	62	61	59	56	55	53	60	62	60	59	66	65	65	63
2nd	61	62	60	61	57	57	57	53	60	65	60	58	66	65	65	65

INTRA-OBSERVER VARIATION (e.g. observer #1)

e.g. observer #1

```
PROC GLM in SAS ==> estimating components 'by hand'
```

```
INPUT subject rater occasion earsize; if observer=1;
The data set has 16 obsns & 4 variables.
```

```
proc glm; class subject; model earsize=subject / ss3;
random subject ;
```

General Linear Models Procedure: Class Level Information

Class Levels Values

SUBJECT 8 1 2 3 4 5 6 7 8 ; # of obsns. in data set = 16

Dependent Variable: EARSIZE

Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	7	341.00	48.71	35.43	0.0001
Error	8	11.00	1.38		
Corrected Total	15	352.00			

R-Square C.V. Root MSE EARSIZE Mean
0.968750 1.80 1.17260 65.0

Source	DF	Type III SS	Mean Square	F Value	Pr > F
SUBJECT	7	341.00	48.71	35.43	0.0001

Source Type III Expected Mean Square
SUBJECT Var(Error) + 2 Var(SUBJECT)

Var(Error) + 2 Var(SUBJECT) = 48.71
 $\frac{Var(Error)}{2} = 1.38$
 2 Var(SUBJECT) = 47.33
 Var(SUBJECT) = 47.33 / 2 = 23.67

Estimating Variance components using PROC VARCOMP in SAS

```
proc varcomp; class subject ; model earsize = subject ;
Variance Components Estimation Procedure: Class Level Information
```

```
Class      Levels      Values
```

```
SUBJECT    8      1 2 3 4 5 6 7 8 ; # obsns in data set = 16
```

```
MIVQUE(0) Variance Component Estimation Procedure
```

Variance Component	Estimate
	EARSIZE
Var(SUBJECT)	23.67
Var(Error)	1.38

• **ICC** (Fleiss § 1.3)

$$\text{ICC} = \frac{\text{Var}(\text{SUBJECT})}{\text{Var}(\text{SUBJECT}) + \text{Var}(\text{Error})} = \frac{23.67}{23.67 + 1.38} = 0.94$$

1-sided 95% Confidence Interval (see Fleiss p 12)

df for F in CI: (8-1)= 7 and 8

so from Tables of F distribution with 7 & 8 df, F = 3.5

So lower limit of CI for ICC is

$$= \frac{35.43 - 3.5}{35.43 + (2 - 1) \cdot 3.5} = \underline{0.82}$$

EXERCISE: Carry out the estimation procedure for one of the other 3 observers.

INTERPRETING YOUR GRE SCORES

(Blurb from Educational Testing Service)

Your test score is an estimate, not a complete and perfect measure, of your knowledge and ability in the area tested. In fact, if you had taken a different edition of the test that contained different questions but covered the same content, it is likely that your score would have been slightly different. The only way to obtain perfect assessment of your knowledge and ability in the area tested would be for you to take all possible test editions that could ever be constructed. Then assuming that your ability and knowledge did not change, the average score on all those editions, referred to as your "true score," would be a perfect measure of your knowledge and ability in the content areas covered by the test. Therefore, scores are estimates and not perfect measures of a person's knowledge and ability. Statistical indices that address the imprecision of scores in terms of standard error of measurement and reliability are discussed in the next two sections.

STANDARD ERROR OF MEASUREMENT

The difference between a person's true and obtained scores is referred to as "error of measurement."* The error of measurement for an individual person cannot be known because a person's true score can never be known. The average size of these errors, however, can be estimated for a group of examinees by the statistic called the "standard error of measurement for individual scores:" The standard error of measurement for individual scores is expressed in score points. About 95 percent of examinees will have test scores that fall within two standard errors of measurement of their true scores. For example, the standard error of measurement of the GRE Psychology Test is about 23 points. Therefore, about 95 percent of examinees obtain scores in Psychology that are within 46 points of their true scores. About 5 percent of examinees, however, obtain scores that are more than 46 points higher or lower than their true scores.

Errors of measurement also affect any comparison of the scores of two examinees. Small differences in scores may be due to measurement error and not to true differences in the abilities of the examinees. The statistic "standard error of measurement of score differences" incorporates the error of measurement in each examinee's score being compared. This statistic is about 1.4 times as large as the standard error of measurement for the individual scores themselves. Approximately 95 percent of the differences between the obtained scores of examinees who have the same true score will be less than two times the standard error of measurement of score differences. Fine distinctions should not be made when comparing the scores of two or more examinees.

RELIABILITY

The reliability of a test is an estimate of the degree to which the relative position of examinees' scores would change if the test had been administered under somewhat different conditions (for example, examinees were tested with a different test edition).

Reliability is represented by a statistical coefficient that is affected by errors of measurement. Generally, the smaller the errors of measurement in a test, the higher the reliability. Reliability coefficients may range from 0 to 1, with 1 indicating a perfectly reliable test (i.e., no measurement error) and zero reliability indicating a test that yields completely inconsistent scores. Statistical methods are used to estimate the reliability of the test from the data provided by a single test administration. Average reliabilities of the three scores on the General Test and of the total scores on the Subject Tests range from .88 to .96 on recent editions. Average reliabilities of subscores on recent editions of the Subject Test range from .82 to .90.

Data regarding standard errors of measurement and reliability of individual GRE tests may be found in the leaflet *Interpreting Your GRE General and Subject Test Scores*, which will be sent to you with your GRE Report of Scores.

VALIDITY

The validity of a test—the extent to which it measures what it is intended to measure—can be assessed in several ways. One way of addressing validity is to delineate the relevant skills and areas of knowledge for a test, and then, when building each edition of the test, make sure items are included for each area. This is usually referred to as content validity. A committee of ETS specialists defines the content of the General Test, which measures the content skills needed for graduate study. For Subject Tests, ETS specialists work with professors in that subject to define test content. In the assessment of content validity, content representativeness studies are performed to ensure that relevant content is covered by items in the test edition.

Another way to evaluate the validity of a test is to assess how well test scores forecast some criterion, such as success in grade school. This is referred to as predictive validity. Indicators of success in graduate school may include measures such as graduate school grades, attainment of a graduate degree, faculty ratings, and departmental examinations. The most commonly used measure of success in assessing the predictive validity of the GRE tests is graduate first year grade point average. Reports on content representativeness and predictive validity studies of GRE tests may be obtained through the GRE Program office.

* The term "error of measurement" does not mean that someone has made a mistake in constructing or scoring the test. It means only that a test is an imperfect measure of the ability or knowledge being tested.