

Levels of measurement in psychological research:

Psychology is a science. As such it generally involves objective measurement of the phenomena being studied, whatever these might be. However, not all measurements are the same. There are important distinctions between different kinds of measurements that you need to be aware of, because how you measure things affects what kinds of statistical test you can use on your data. We covered levels of measurement at the start of the autumn term, but that may have seemed rather abstract at the time; now that you have completed a couple of practicals, and have had some experience of encountering different levels of measurement, it's a good time to revisit the topic.

1. "Numbers" which are really names - the "nominal" scale of measurement:

Sometimes all you can do is place people into categories and record the frequency with which each category occurs. In this situation, you might use *numbers* as *names* for the categories. The examples I used in one of my autumn lectures were the numbers on footballers' jerseys, and house numbers. These are not "real" numbers, and you cannot do any arithmetic with them other than count how many instances of each category occur.

This can sometimes be confusing, especially when using SPSS, which requires you to use numbers in order to code participants on various attributes. For example, in order to tell SPSS about the gender of your participants, you might use "1" to stand for "male" and "2" to stand for "female", like this:

subject name	gender
1	1
2	1
3	1
4	2
5	2
6	2

You can certainly *try* to perform various mathematical operations on these data, because SPSS will unintelligently go along with your demands. However the results will be quite meaningless. Here, the mean of the subject names is 3.5 and the mean gender is 1.5. Neither of these makes any sense, because it is impossible to combine names or genders in this way. If I go to "variable view" in SPSS, and use the "value label" option, I can make SPSS show names as words instead of code numbers. The absurdity of trying to do arithmetic on these values is now even more obvious: what's the average of three "males" and three "females", or the average of "Bob", "Bill", "Eric", etc?

subject name	gender
Bob	male
Bill	male
Eric	male
Cynthia	female
Ethel	female
Doris	female

The Wason task - an example of nominal data:

In the Wason card sorting practical, the data were measurements on a nominal scale of measurement: all we did was record whether or not each participant got the right answer to the two problems that they attempted (the formal logic problem and its "concrete" counterpart). We coded the data using "1" for "right" and "2" for "wrong".

All you can do with nominal data is use the Chi-Square test to see if there are any significant differences in the frequencies with which the various categories occur. This is what we did with the Wason task: we merely counted up how many people passed or failed each reasoning task, and then looked to see if these frequencies differed from what we would expect to have obtained by chance. It makes no sense to calculate means and standard deviations for frequency data; a frequency is a frequency, pure and simple. As a result, any graphs would simply show the frequency with which each category occurred, with no error bars.

2. Measurements using proper numbers - the "ordinal", "interval" and "ratio" scales:

Deciding whether or not you have nominal (frequency) data is usually fairly straightforward. Think about the data provided by each participant: if all you know is that the participant falls into one of a number of categories, then you have data on a nominal level of measurement. If you have one or more scores from each participant, then it is clear that you do not have data on a nominal scale. However what you then need to do is to decide whether your data are measured on an ordinal, interval or ratio scale. This is sometimes tricky to decide. It comes down to two issues:

- (a) are there **equal intervals** between the various points on your measuring scale?
- (b) does the measuring scale have a **true zero point**, as opposed to an arbitrary one?

If data are measured on an **ordinal** scale, then (as the name implies!) they can be placed in some kind of *order*. Examples of ordinal scales might be: "small", "medium", "big"; "very tired", "quite tired", "awake", "very awake"; and "very happy", "happy", "neutral", "unhappy", "very unhappy". However, the points on an ordinal scale are not necessarily equally spaced. You can't do anything other than arrange the values in order of magnitude (amount of whatever it is you are trying to measure, such as size, alertness or mood, in the case of the scales just mentioned) is pretty much all that you can do with them.

The classic example of an ordinal scale is sporting performance. If you are told who comes first, second and third in a horse race, you know that the horse who came first was faster than the horse who came second, who in turn was faster than the horse who came third. Thus you can rank the horses in order of "speed". However, if this is all your data consist of, you don't know anything more about the horses' performance: it might be that the first horse beat the second one by a few seconds, and that the second horse beat the third one by a minute. Or it might be that the first horse beat the other two by minutes, and the second and third horses had very similar times. An ordinal scale of "first", "second" and "third" contains no information about the distances between these points on the scale.

In contrast to ordinal scales, if data are measured on an **interval** or **ratio** scale, the distances between the various points on the scale are equivalent across the whole range of

measurements. The distinction between interval and ratio scales is rather subtle: a ratio scale has a true zero point, whereas the interval scale does not. If there is a zero value on an interval scale, it is merely an arbitrary point on the scale that is regarded as "zero" by definition.

The classic illustration of interval and ratio scales is temperature. Both the Centigrade and Fahrenheit temperature scales are interval scales. In both cases there are zero points, but these do not represent a true absence of temperature - they are merely arbitrary points on the scale. On both of these scales, it is quite possible to have temperatures below zero. In contrast, the Kelvin temperature scale is a ratio scale: zero degrees on this scale is defined as a complete absence of heat. Ratio scales have a true zero point, marking a total absence of the attribute being measured. The existence of this zero point means that you can make additional statements about the relationships between different points on a ratio scale.

To illustrate this, consider the example of temperature again. On all three scales, the points on the scale are equally spaced wherever you happen to be on the scale. Therefore on all three scales, it is possible to say that the temperature has increased by one degree, or decreased by two degrees, etc. A degree of temperature is a constant amount, and so a change from 21 to 22 degrees is the same amount of change in temperature as a change from 3 to 4 degrees. However, with the two interval scales (Fahrenheit and Centigrade), the absence of a true zero point makes it impossible to make ratio statements such as "it is twice as hot today as it was yesterday". You *can* make statements like this with a ratio scale, because on a ratio scale the zero point is a true zero (a total absence of the property being measured) and not just another point on the scale.

Centigrade	100	90	80	70	60	50	40	30	20	10	0	-10	-20	-30	-40	-50	-60 ..	-273
Fahrenheit	212	194	176	158	140	122	104	86	68	50	32	14	-4	-22	-40	-58	-76 ..	-460
Kelvin	373	363	353	343	333	323	313	303	293	283	273	263	253	243	233	223	213 ..	0

Examples of ratio scales in psychology are things such as reaction time, and individual scores such as "number of items correctly recalled" or "number of errors". With these kinds of measures, it is valid to make statements about ratios, such as "Fred was twice as fast as Dorothy", or "Fred made half as many mistakes as Cynthia". The statement "Fred got no items correct" is also valid, because there is a true zero on a "number correct" scale, representing a complete absence of correct responses.

Examples of interval scales include most IQ tests. There is no true zero point on an IQ test, so although I can say that "my IQ is 70 points higher than yours", I cannot say that "I have an IQ of 140 and you have an IQ of 70, so therefore I am twice as intelligent as you".

In practice, you don't need to worry too much about the difference between interval and ratio scales, because that won't affect your choice of statistical test. A simple way to choose between them is to think of whether a score of zero on your scale represents a complete absence of the thing being measured. If it does, you have a ratio scale; if not, you have an interval scale.

You *do* need to be able to appreciate the difference between these scales and an ordinal scale, because parametric statistical tests require data to be measurements on either an interval or ratio scale (i.e. they should not be used on ordinal data).

The nature of the data produced by the "Risky Shift" practical:

Sometimes it can be quite complicated to work out precisely what kind of scale you have. We can illustrate this with the data from the "Risky Shift" experiment. Here's one of the questions in the questionnaire that was used for this practical.

Imagine that you are advising Mr A. Look at the probabilities or odds of the new company proving financially sound, and choose the lowest acceptable probability in order for Mr A to take the new job.

1. Please tick the <u>lowest</u> probability that you would consider acceptable to make it worthwhile for <u>Mr A</u> to take the new job.	RISK SCORE
<input type="checkbox"/> The chances are 1 in 10 that the company will prove financially sound	<input type="checkbox"/> → 1
<input type="checkbox"/> The chances are 3 in 10 that the company will prove financially sound	<input type="checkbox"/> → 3
<input type="checkbox"/> The chances are 5 in 10 that the company will prove financially sound	<input type="checkbox"/> → 5
<input type="checkbox"/> The chances are 7 in 10 that the company will prove financially sound	<input type="checkbox"/> → 7
<input type="checkbox"/> The chances are 9 in 10 that the company will prove financially sound	<input type="checkbox"/> → 9
<input type="checkbox"/> Tick here if you think <u>Mr A</u> should <u>not</u> take the new job, no matter what the probabilities.	<input type="checkbox"/> → 11

What sort of data are these, and why?

1. In the pre- and post-discussion conditions, we are asking each person to give us a number of scores - one score for each question. We therefore have quantitative data, with each participant giving us a set of numerical scores (as opposed to each participant merely falling into a category, such as "risky" versus "not risky").

2. For each question, we have six options. We code these as 1 to 11. The options represent different levels of risk (from 1, which represents acceptance of the highest level of risk, through to the final option, coded as 11, which represents accepting no risk whatsoever). Thus we have a scale of some kind.

3. What kind of data are these? Think about the different options. We can consider them to be measurements on a quantitative scale, since they represent different levels of riskiness. So we can rule out the possibility that these are categorical/nominal data. They are measurements on a scale, and hence either ordinal, interval or ratio measurements.

At first sight, you might think that this is a *ratio* scale, especially if we used 0 rather than 11 to represent refusal to take any risk whatsoever. You could certainly argue that the final option in the question above does represent a true absence of "riskiness". However, these are *not* ratio data. In order for data to be regarded as ratio data, the scale has to have equal gaps between different points on the scale, as well as a true zero point. On our scale, the different points do not necessarily represent equal intervals of "riskiness", even though the probabilities attached to the options might lead you to think that they do.

Always think about what the participant is being asked to do, and hence what scores really represent. All we can say in this case is that someone who picks a "chance of 1 in 10" is prepared to be more risky than someone who picks a "chance of 3 in 10" or "5 in 10". While the options in each question might give the *impression* that these ratings are equally spaced along a scale, we cannot know this for certain. We don't have enough information about the underlying

psychological construct of "riskiness" to know whether the options are really as equally spaced as they might appear from the wording of the questions.

Also, because there is an option of not taking any risk at all (in the case of the example above: choosing not to take the job at all), we cannot say that there are equal gaps between the riskiness scores. In fact the final option is different in nature from the others. Not taking the job (which gets a score of 11), is not a logical next step after taking the job if the chances are 9 in 10 (which gets a score of 9). Taking the job if there was a chance of 10 in 10 (i.e. success was a certainty) would have been the next logical step. Because of all these reasons, this scale is *not* a ratio, *nor* an interval scale. Therefore it is best to treat these scores as *ordinal* data: they are ratings of "riskiness" on a 6-point scale.

4. What kind of descriptive statistics can we do on these data? We ask you to produce a mean riskiness score for each participant. You might think this is odd, given what I've just said about the points on the riskiness scale not being equally spaced. However, there is a subtle difference that needs to be appreciated here. It is quite valid to use means as descriptions of the ratings themselves: in effect, all we are saying is that "this is typically the option that participants chose on this question". If the average for the pre-discussion questionnaire was 5 and the average for the post-discussion questionnaire was 10, it would be valid to say that average riskiness had diminished, because the means do reflect some kind of increase in riskiness. (Remember that this questionnaire is "reverse coded", so that *high* scores represent *low* riskiness).. However it would *not* be valid to say that riskiness had "halved": all we really know is that higher ratings are less risky than lower ones, but we don't know precisely how our scale is related to the underlying psychological construct of riskiness.

So - the main message to take home here is "beware of scales that make things look more precise than they really are". The use of "acceptable probabilities" in these questions makes them appear like proper numbers, but really they are just ordinal data.

The nature of the data produced by the "Memory Perspectives" practical:

You'll be relieved to know that the data for the final lab report, on "memory perspectives", are far more straightforward! Here we are recording the number of items correctly recalled by each participant. This is most definitely a *ratio* scale of measurement, with equal intervals and a true zero point (a participant could, in theory, remember no items at all from the passage; and it is meaningful to make statements like "participants in one group recalled twice as many items as participants in the other group").

What kinds of statistical tests can you perform on ordinal, interval and ratio data?

Statistical tests can be divided into two kinds: parametric tests (which make certain assumptions about the nature of the data on which you are performing the test) and nonparametric tests (which don't make those assumptions).

The three assumptions that need to be met in order for you to perform a parametric test are:

- (a) the data should be roughly normally distributed;
- (b) the data should show homogeneity of variance (the spread of scores in the different conditions of the study should be roughly similar);
- (c) the data should be measured on an interval or ratio scale.

For the purposes of the exam, we have a simple policy: you can perform a parametric test ONLY if all three of these requirements are met. In other words, we say that you should not perform a parametric test on ordinal data (such as ranks), and in the exam we would mark any attempt to do so as being incorrect. However, an important part of the university experience is learning to be tolerant of ambiguity, and so you should be aware that this is a grey area in practice. Not all researchers and statisticians think that it's a problem to use parametric tests on ranked data, so you may well come across published research that uses parametric tests on ordinal data such as personality measures, attitude scale data, Likert scale scores, etc.

Type of data:	Permissible descriptive statistics:	Permissible inferential statistics:
Nominal	Counts (frequencies). Statements like "more people chose coffee than tea as their preferred drink".	Chi-Square
Ordinal	Median, mode (mean, though arguable). Statements like "people liked coffee more than tea". (But we don't know by how much).	Nonparametric tests (e.g. Wilcoxon, Mann-Whitney, Friedmans, Kruskal-Wallis)
Interval	Median, mode, mean. Statements like "on the 'Beverage Appreciation Scale', people gave coffee a higher score than tea".	Parametric tests (e.g. t-tests, ANOVA)
Ratio	Median, mode, mean. Statements like "people drank twice as many cups of coffee than they did tea". (This is a ratio statement).	same as interval

The case of Likert Scales:

A popular measuring tool in psychology is the Likert Scale. This usually consists of a statement plus a rating scale that goes in apparently equal increments from an extreme negative to an extreme positive opinion. For example:

"Cats are evil little monsters".

Strongly disagree	Disagree	Neither agree nor disagree	Agree	Strongly agree

A participant who gives a rating of "strongly agree" clearly feels more deeply that cats are evil than does someone who gives a rating of "disagree". Sometimes the verbal labels are replaced by numbers, such as 1-5, where 1 would be "strongly disagree" and 5 would be "strongly agree".

What kind of data are produced by Likert Scales? It's clearly not a ratio scale, as there is no true zero point, so at least we can exclude that option. However, are these ordinal or interval data?

Some people would argue that Likert Scales produce interval data, others that they are really ordinal data¹. The central issue is whether or not the increments on the scale are truly equally spaced. At first sight they appear to be, especially if numbers are used to represent the different points on the scale, so that the scale runs from 1-5, 1-9, -3 to +3 or whatever. However, if you think about the psychological property that you are trying to measure with this scale, it's clear that in fact it is an ordinal scale.

If I give a rating of 5 and you give a rating of 5, we know that we both strongly agree with the statement, but we have no way of knowing for certain whether we really do have similar depths of antipathy to cats. We might both be using the same verbal label to represent different levels of feeling. Replacing the verbal label with a number from 1 to 5 makes these data look like an interval scale because the *numbers* are equally spaced - but we cannot know whether the *psychological property underlying the responses* is also equally spaced. How can we be certain that "amount of depth of feeling about cats" falls on such a linear scale?

We have no way at all of knowing whether the differences between the different points on the scale are truly equivalent, as they must be in order for it to be regarded as an interval scale. Is the difference in cat-hatred between me and someone who gives a rating of "agree" really the same as the difference in cat-hatred between that person and someone who gives a rating of "neither agree nor disagree"? We cannot do anything more than place people in order of magnitude of cat loathing, on the basis of their responses to this item. Therefore this is an ordinal scale.

¹ These include the people involved in this course, which is why the practical slides sometimes refer to Likert scales as interval data while the lecture slides always refer to them as ordinal data!

What statistics can I do with Likert Scale data?

If you accept the argument above, then Likert Scale data are not suitable for parametric tests which require the data to be measured on an interval or ratio scale. However, as mentioned earlier, this is a grey area: in practice, researchers often do perform parametric tests on them.

However, I reiterate: for the purposes of this course, all rating scale data (Likert scales included) are to be treated as ordinal data and hence only analysed with nonparametric statistical tests.

What about descriptive statistics? Does it make sense to summarise these data with means and standard deviations? The answer is a qualified yes! Any means and standard deviations obtained from rating data (whether from the "Risky Shift" data or from a Likert Scale) are perfectly valid as *descriptions of participants' behaviour*, i.e. how participants responded when faced with a question and asked to pick a response. So it is fine to say something like "the mean rating chosen was 4.6, with a standard deviation of 1.2". This tells us that "typical" performance was to pick a rating of "agree" or similar, although there was some spread around this choice.

However, what this actually means in the context of the *underlying psychological construct* of "attitudes to cat morality" is a different question. Our data tell us that most people think cats are rather evil, but as we cannot know for certain that everyone who gave a particular rating really did feel exactly the same way about cats, we would have to be cautious in interpreting these data. For example, if we had three groups (cat lovers, people who were quite indifferent to cats, and cat haters) and they gave us different mean ratings (say "1", "2" and "5" respectively) we could say that the mean ratings for the groups differ, and that the three groups differ in their attitudes to cats. We would not be able to say much more than that.

In short, the use of Likert Scales raises an important point: in psychology, you need to distinguish between measuring people's behaviour, and interpreting what those measurements actually represent in psychological terms. This problem isn't unique to Likert Scales. For example, suppose you measure the reaction times of young people and old people and find a difference between them. The difference itself is real, but what gives rise to it may be much less clear-cut. It might stem from cognitive decline in the elderly participants, or the use of different strategies between the two groups, or some combination of the two. When you perform a study, always think carefully about what it is that you are *really* measuring.

Thanks to Linda Tip and Sarah Laurence for their contributions to the arguments in this handout.