



The  
University  
Of  
Sheffield.



## *2. INTRODUCING VARIABLES AND UNDERSTANDING THEIR LEVELS OF MEASUREMENT*

Dr Tom Clark & Dr Liam Foster

Department of Sociological Studies | University of Sheffield

# CONTENTS

<i>Recording data numerically</i>	<i>Nominal levels of measurement</i>	<i>Ordinal levels of measurement</i>	<i>Interval (and ratio) levels of measurement</i>
..... page 2	..... page 2	..... page 4	..... page 7
<i>Recoding variables</i>	<i>Identifying levels of measurement</i>	<i>Levels of measurement flow diagram</i>	<i>Why do I have to be able to identify levels of measurement?</i>
..... page 8	..... page 9	..... page 10	..... page 10
<i>Thinking critically about variables</i>	<i>Labelling variables</i>	<i>Research questions, hypotheses, and variables</i>	<i>Rounding up</i>
..... page 10	..... page 11	..... page 14	..... page 15

# 2.

## Introducing and understanding variables

*\*Within qualitative worlds, there are different ways of recording data about our social lives. We may use visual material in the form of pictures, oral testimony in the form of interviews, or thick description in our ethnographies. All of these forms help us to research social action and all have particular advantages and disadvantages.*

### 2.1. Recording data numerically

Equally, not all quantitative data is the same and here are different ways in which we can record quantitative material in order to explore our social worlds and answer our research questions. We may wish to simply count something; we might want to think of measuring an attitude on some sort of scale; or we might even attempt to specify where something lies with respect to an established format like, for instance, the amount of money someone earns in an hour.

Any attempt to measure something quantitatively can be similar, or different, to another attempt to capture data. Indeed, all types of quantitative data can be divided into particular levels of measurement and these levels have different techniques of analysis associated with them. How you summarise data is entirely dependent on the level of data that you have. This is, essentially, why it is so important to be able to recognise the level at which data is being measured.

This workbook will introduce you the different ways in which we can record quantitative material.

By the end of the workbook you should be able to:

- Identify different levels of measurement
- Be able to construct variables at an appropriate level
- Think critically in relation to variable measurement
- Understand the role of a variable with respect to your research aims and objectives, and your hypotheses

### 2.2. Understanding quantitative variables: Identifying levels of measurement

Look at these examples:

Q2.1. Do you smoke cigarettes? e

Yes  No

Q2.2. Do you consider yourself to be a:

non-smoker  light-smoker  
 medium smoker  heavy smoker

Q2.3. How many cigarettes have you smoked in the last seven days?

\_\_\_\_\_ (write here)

Each of these closed questions represents something - a variable - we might be interested in. In this case it is smoking behaviour. However, whilst all of these variables are concerned with smoking, they are recording the information in subtly different ways. These different ways of recording quantitative data are more commonly known as levels of measurement. There are three main levels of measurement: nominal, ordinal, and interval. We shall deal with each in turn.

#### 2.2.1. Nominal levels of measurement

Look again at the first question.

Q2.1. Do you smoke cigarettes? e

Yes  No

Here, the question constructs a category between cigarette smokers on one hand, and non-cigarette smokers on the other. It assumes that all people can be divided into two separate groups.

In this particular case, you are either a cigarette smoker or you are not. There is no middle ground – you are either one or the other. When you are dealing with variables that measure categories such as these, you are dealing with a variable that is operating at a nominal level of measurement. Nominal simply means ‘name’. It is, however, worth noting that nominal variables are also sometimes called ‘categorical’ variables because they measure distinct categories. Confusingly perhaps, whilst all nominal variables are categorical, not all categorical variables are nominal as categorical can also be applied to ordinal variables too – but we’ll deal with those in a minute.

Although this particular smoking variable is a binary category – you are either a smoker or you are not – not all nominal variables are binary. Ethnicity, marital status, nationality, geographical location, and occupation are all commonly used forms of nominal data. Class is also sometimes treated as a nominal variable.

**Here is an example of a variable that is measured at the nominal level: ethnicity. It was included in the 2011 census (ONS, 2012):**

- C. Asian / Asian British
- Indian
  - Pakistani
  - Bangladeshi
  - Chinese
  - Any other background, write in: \_\_\_\_\_
- D. Black / African / Caribbean / Black British
- African
  - Caribbean
  - Any other Black / African / Caribbean background, write in: \_\_\_\_\_
- E. Other ethnic group
- Arab
  - Any other ethnic group, write in: \_\_\_\_\_

**Notice the instruction at the start:**

“  
Choose one section from A to E, then tick one box to best describe your ethnic group or background.  
”

You are either in group A, B, C, D, or E – there is no sliding scale between White, Mixed, Asian, Black, or Other. Once you have decided which category you are in, you now have to ‘tick’ which sub-set of that category you belong to. Again there is no sliding scale between the categories and they are mutually exclusive – they don’t bleed into each other.

Let us suppose that I’m interested in exploring the relationship between marital status and money-saving behaviour. I would first need to construct a nominal variable that measures marital status. Remember, everyone who answers a question that is measured at the nominal level has to be able answer it – the range of answers needs to cover every possibility.

**T** Try to construct a variable at the nominal level for marital status – make sure that the range of your answers covers every possibility.

**This is how marital status was measured in the British Social Attitudes Survey in 2007:**

**e** Q2.5. Can I just check, which of these applies to you at present?

Please choose the first on the list that applies:

- Married
- In a civil partnership
- Living with a partner
- Separated (after being married)
- Divorced
- Widowed
- Single (never married)
- (Don’t know)
- (Not answered)

This question is actually designed to be read aloud by the interviewer. What is important to recognise, however, is that all possibilities are covered: everybody who responds to the questionnaire will fall into one of these categories.

Now we need to design a variable that will allow us to assess money-saving behaviour: we decide to measure this at the nominal level.

**T** Construct a variable at the nominal level for money-saving behaviour.

**It might look something like this:**

**e** Q2.6. Do you save money on a regular basis?

- Yes
- No

Of course, we could measure money-saving behaviour in a slightly different manner. We might just ask how much people save on a weekly basis. However, such a measurement would not be at the nominal level. We now have all our variables that will enable us to answer our research question.

**Q** However, before we celebrate our newly discovered knowledge by actually conducting a survey with these variables, in which order do you think these variables should be presented within a survey?

Well, the answer does depend on a few things. If we are only interested in measuring these variables and our project is very small, then we don’t really have to worry about it as there are only two questions. However, if this is just one research question of many concerning financial habits then it is often good practice to put the marital status variable with other so-called demographic variables – such as gender, social class, ethnicity, religious persuasion etc, etc.

Demography essentially means the measurement of people. Demographic variables are usually taken to be quite general variables that describe the various categories that a person has and many are nominal in nature<sup>1</sup>. It is often useful to place such variables toward the start of a questionnaire as they are easy to answer and often allow the respondent to ‘get the feel’ of the questions. More substantive variables, that is variables that cover more specific or specialist topics such as the ‘saving status’ one that we constructed above, tend to come later in the questionnaire as they can require a little more thought and can sometimes require some rapport with the respondent. This rapport is often developed by placing the ‘easy questions’ first.

2.2.2. Ordinal levels of measurement

**Look again at the Q2.2. from our smoking example:**

Q2.2. Do you consider yourself to be a:

- Non-smoker
- Light-smoker
- Medium smoker
- Heavy smoker

Like Q1, this question is still concerned with smoking but it is a little more refined as it suggests there is some sort of order to the answer. It is still possible to divide people into non-smokers and smokers, but the question gets a little more specific about how they perceive what type of smoker (or not) they actually are. Unlike Q2.1, there is some middle ground. The world is no longer being reduced to smokers and non-smokers, and instead there is some room for shading between those categories. After all, some people might consider themselves to smoke only in specific circumstances, whilst others might consider themselves to be chain smokers.

<sup>1</sup> It’s worth noting that age, common to demographic measurements, is an outlier in this respect and is not a nominal variable.

However, whilst there is a scale here, it is non specific and there is a lot of room for interpretation: the categories are not neat and well bounded. What 'heavy smoking' means for one person, may be 'light smoking' for another. In this sense, the distances between the points on the scale are not equal or well defined.

Ordinal variables can also be used to group together 'counts' of things. **For instance the following scale, also a smoking variable, is also constructed on an ordinal scale:**

**e** Q.2.2a. How many cigarettes do you smoke a day?  
 0-18    19-30    31-40    40+

Notice how the range between each point is not equal. The first point covers 18 possibilities, but the second only 12, the third just 10 and the final point is unlimited.

**Q** Can you think of any particular reason why the answers are grouped in this manner?

Without any clear rationale for constructing the range of answers like this, data in this format probably would not be useful – the categories are relatively meaningless as they are built upon nothing but guess work. A much better way of constructing this particular question would be to find some precedence within the literature.

**Q** How might you do this?

A good option is to map our ranges onto ones that are standard within variables of this type. One potential option is to use current understandings of heavy smoking. Whilst there is no universally accepted definition for heavy smoking, there is some acceptance within the medical literature that more than 10 per day is heavy. Some studies have measured 1-9 cigarettes per day as moderate smoking, and less than that as light.

**Using this rationale, a better question might be:**

**e** Q.2.2b. How many cigarettes do you smoke a day?  
 I never smoke    Less than one per day  
 1-9    10+

The possible range of response now corresponds with more standardised definitions of non-smoker, light smoker, moderate smoker, and heavy smoker. However, the values between the scales are still unequal. Hence, this is still an ordinal variable, but now it has a firmer methodological justification.

This is not to say that ordinal scales have to have some theoretical basis. Indeed, other variables that employ an ordinal scale will simply group counts equally as it is an easy way to summarise wide-ranging data. **Many age variables, for instance, will appear in formats like below:**

**e** Q.2.7. Please state your age:  
 0-16    17-29    30-39    40-49  
 50-59    60-69    70-79    80+

Ordinal scales are not, however, only suitable for 'perceptions of behaviour', 'aggregated counts' or other forms of ordered data. Indeed, one of the most popular uses of the ordinal scale is to employ it to measure attitudes. A popular version of this sub-type of ordinal scale is called a 'Likert scale'.

Named after the psychologist who invented it, Rensis Likert, a Likert scale is a series of rating scales on which respondents specify their beliefs, attitudes, or feelings about a particular topic or issue. Typically, Likert scales enable the respondent to numerically express the strength of feeling with respect to a specific statement or question topic. Multiple measurements made on scales are usually bipolar in design in that they are constructed around two polar opposites with a number of points inbetween. However, the number of points on a Likert item does vary and 3 point, 5 point, 7 point, and even 10 point scales are popular. Indeed, even numbered scales are possible if the researcher wants to force a choice in situations where a neutral option is undesirable. Technically, Likert scales are used to summarise a range of responses across a series of variables so that an underlying attitude can be assessed.

However, since its introduction the meaning has loosened and it is now often used to refer to a series of single unconnected items. When used in this singular format, the question is often better referred to as Likert item or a Likert-type variable.

**Look at these examples of Likert items taken from the last wave of the British Household Survey<sup>2</sup> - again, it is a question that is read aloud by an interviewer:**

**e** Q2.8. Please look at the card and tell me how much you agree or disagree with the following statements?

A. It takes too much time and effort to do things that are environmentally friendly:

Strongly agree   Agree   Neither agree nor disagree   Disagree   Strongly disagree  
 1   2   3   4   5

B. Scientists will find a solution to global warming without people having to make big changes to their lifestyle:

Strongly agree   Agree   Neither agree nor disagree   Disagree   Strongly disagree  
 1   2   3   4   5

C. The environment is a low priority for me compared with a lot of other things in my life:

Strongly agree   Agree   Neither agree nor disagree   Disagree   Strongly disagree  
 1   2   3   4   5

D. I am environmentally friendly in most things that I do:

Strongly agree   Agree   Neither agree nor disagree   Disagree   Strongly disagree  
 1   2   3   4   5

In each of these variables, attitude toward an aspect of environmental friendliness is being measured on a five point Likert scale. However, each item is measuring a slightly different aspect of environmental friendliness. The first is measuring the perceived cost of engaging in environmentally friendly behaviour; the second explores the perceived locus of control with respect to the solutions for environmental problems; the third assesses the priority of environmentally friendly behaviour against generalised goals; and, the final measure is examining the perception of self with respect to being environmentally friendly. The items could be used as single measures, or taken collectively to measure general attitude to environmental friendliness.

However, if you were to take these variables as being indicative of an underlying attitude towards environmental friendliness, you also need to note that the weight of the statement is not the same for all four questions: it's not quite as easy as adding up the relative answers and dividing by the number

of observations to get the average. In the first three variables, the emphasis is on agreeing with statements that are broadly negative towards environmentally friendly behaviour. In the fourth, agreement is broadly positive. Indeed, it is important to recognise the direction of measurement: agreement does not always mean the same thing.

If you were to use these variables collectively as a Likert scale, it is also important to recognise that the difference between each point is open to interpretation: what is the distance between 'Strongly Agree' and 'Agree'? The answer will largely depend on individual perception and the points on the scale are not clear and continuous. As a result, they should be taken to be equal in nature: this is exactly why a Likert scale is ordinal in design.

<sup>2</sup> See [http://www.iser.essex.ac.uk/survey/bhps/documentation/pdf\\_versions/questionnaires/bhpsw18q.pdf](http://www.iser.essex.ac.uk/survey/bhps/documentation/pdf_versions/questionnaires/bhpsw18q.pdf).

It is worth noting, however, that there is some discussion amongst statisticians and social researchers concerning whether Likert scales should be treated at the ordinal or interval level. It is actually common practice, particularly in areas of psychometrics, to treat Likert-based attitude measures at the interval level – particularly when multiple items are used in conjunction with one another. This is an important issue to recognise because the appropriate descriptive and inferential techniques used to analyse quantitative data differ between ordinal and interval variables. If the wrong statistical technique is used, the researcher increases the chance of coming to the wrong conclusion about the significance (or otherwise) of their findings. So, until you are able to clearly justify why you want to use a Likert scale at the interval level, it is best to err on the side of caution and treat these measures at the ordinal level.

Look at the examples from the British Household Survey above. Let us suppose that I want to measure the extent to which respondents think that people who haven't paid any taxes should not receive any benefits – construct an ordinal variable that would allow you to do this.

**This is how it was measured in the British Social Attitudes Survey in 2007 (BSA, 2008):**

The government raises money through taxation to pay for social benefits like the state pension, unemployment benefits and sickness benefits. How much do you agree or disagree with each of the statements below about social benefits like these?

Q.2.9. People who haven't paid taxes should not be entitled to any benefits:

- A. Agree Strongly
- B. Agree
- C. Neither agree nor disagree
- D. Disagree
- E. Disagree strongly
- F. Can't choose

As you can see, sometimes Likert-type ordinal scales employ the use of short vignettes, prompts, or cues to contextualise the issue at hand. This method is particularly appropriate to employ where there might be some misunderstanding about what is being asked of the respondent.

Likert scales also need not be numbered within a survey or questionnaire. In the set of questions at Q8, for example, numbers are included; in Q9, they are not. However, it should always be possible to put the points on a numerical scale if you so wish later on. In the case of this particular variable, it would be relatively easy to alter the letters so they reflected a numerical scale of 1 to 5. The final category 'Can't Choose' would not be on the scale, and would instead serve as a count of those people who were unable to offer an answer.

**2.2.3. Interval (and ratio) levels of measurement**

**Let us go back to the third question – it is yet more specific than the first two smoking variable we saw.**

Q.2.3. How many cigarettes have you smoked in the last seven days? \_\_\_\_\_ (write here)

This is an example of an interval variable. An interval variable is similar to an ordinal variable, except that the intervals between the values of the interval variable are equally spaced and the gaps between each point are clear, consistent, and continuous. In the above example, the difference between one cigarette and two cigarettes is the same as it is between five and six – the gap between each measure is always the same. To repeat, where there is no room for individual interpretation between each point on a scale, and the distance between points are equal, then it is likely that you are dealing with an interval measure.

Let us suppose that I'm interested in the amount of hours young people work. Construct a variable at the interval level that will help you to do this.

**This is how it was done in the British Household Panel Survey for Young People (Wave 17).**

Q2.10. How many hours paid work did you do last week? If you have more than one job please write in the total hours you worked at all of them.

Write in hours: \_\_\_\_\_

However, not all interval variables are the same. A ratio variable has all the properties of an interval variable, but also has a clear definition of 0. Age, for instance, has a logical 0 point and is a ratio variable. Indeed, like the example above, most interval level variables within social statistics are actually also ratio variables. In fact, I can't actually think of one that you are likely to come across in social research that isn't, but many textbooks will still mention it. However, temperature is a good example of an interval variable that isn't a ratio variable – the 0 on the Celsius temperature scale is arbitrary. Although it is based on the freezing point of water, the 0 point could just as easily be based on the freezing point of alcohol. In any case, you are unlikely to use this in social research.

Strictly speaking, interval measures are not that common in social research as there are not many social phenomena that take interval form. Items like income are usually the exception rather than the rule. Where they are used, they tend to be constructed from multiple measures, as is the case with psychometric intelligence tests. As previously stated, some researchers do occasionally use ordinal level data as if it were interval level data, but this should only be done with extreme caution.

**2.2.4. Recoding variables**

Although variables measured at the interval level are often considered the best form of data - mainly because they are the most refined - sometimes we may want to recode interval data into ordinal data, or even nominal data.

Can you think of a reason why would we want to recode a variable?

It's often easier to see patterns in our data when we use 'lumpier' categories - hence sometimes it is a good idea to recode interval/ratio level variables

into ordinal or even categorical ones. For instance, although 'Age in years' is a ratio variable, it is often recoded into an ordinal variable with six categories as follows:

- 0 - 15 = 1
- 16 - 30 = 2
- 31 - 45 = 3
- 46 - 60 = 4
- 61 - 75 = 5
- over 75 = 6

There are many ways of recording the age variable and interval level data can even be recoded into a nominal level data. For instance, if you were interested in differences in attitude towards engaging with environmentally friendly behaviour (see Q8) between those of working age and those of retirement age it would be a relatively straight-forward job to aggregate all those between 16 and 64, and those who were 65+.

Recoding variables often helps you to see patterns and trends in your data more easily because you are interpreting a smaller number of categories. In some cases, particularly where you are trying to employ inferential statistics, you may even find yourself having to recode a variable because the data fails to satisfy an assumption of the techniques you are attempting to employ. However, do this with care. Recoded data is less and less refined and you inevitably lose the sensitivity of your data to identify smaller differences and similarities.

2.3. Identifying levels of measurement

You should now be able to identify the levels of measurement for any particular variable.

**T** Try to identify the level of measurement for each of these variables taken from the British Social Attitudes Self Completion Questionnaire 2006.

**e** Q.2.11. How many journeys of less than two miles do you make by car in a typical week?  
Please write in: \_\_\_\_\_

**e** Q.2.12. How much do you agree or disagree with each of these statements?

A. A lot of false benefit claims are a result of confusion rather than dishonesty:

Strongly agree **1**    Agree **2**    Neither agree nor disagree **3**    Disagree **4**    Strongly disagree **5**

B. The reason that some people on benefit cheat the system is that they don't get enough to live on:

Strongly agree **1**    Agree **2**    Neither agree nor disagree **3**    Disagree **4**    Strongly disagree **5**

**e** Q.2.13. Which is it more important for the Government to do? PLEASE SELECT ONE ANSWER ONLY

To get people to claim benefits to which they are entitled  
OR  
To stop people claiming benefits to which they are not entitled  
OR  
Can't choose

**e** Q2.14a. Consider this situation: a person in work takes on an extra weekend job and is paid in cash. He does not declare it for tax and so is £500 in pocket. Do you feel this is right or wrong?

A. Not wrong     B. A bit wrong     C. Wrong     D. Seriously wrong     E. Can't choose

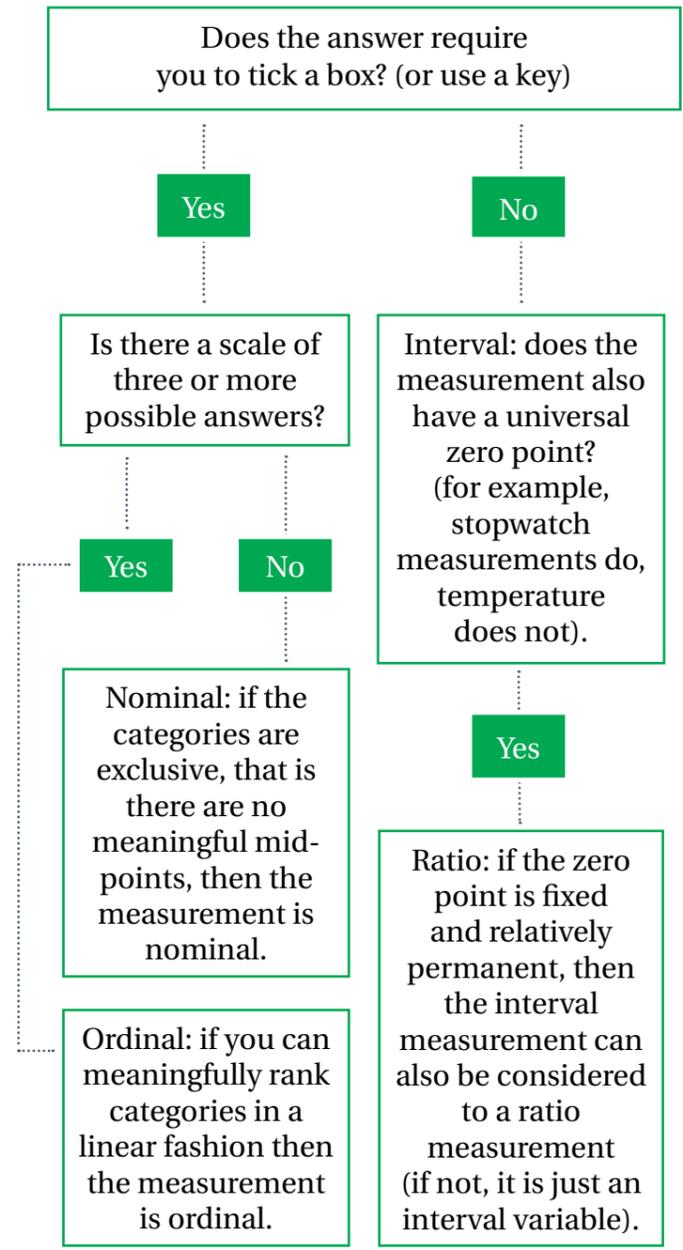
Q2.14b. Consider this situation: and how likely do you think it is that you would do this if you found yourself in this situation?

A. Very likely     B. Fairly likely     C. Not very likely     D. Not at all likely     E. Can't choose

- Q2.11. is an interval variable
- Q2.12A. and Q2.12B. are ordinal variables
- Q2.13 is a nominal variable
- Q2.14A and Q2.14B are ordinal variables
- Q2.15 is an ordinal variable

2.3.1. Levels of measurement flow diagram

If you are still having trouble understanding which level of measurement is which, the following flow chart should help you to identify what particular levels of measurement your variables are.



2.3.2. Why do I have to be able to identify levels of measurement?

You need to be able to tell the difference between a screwdriver, a hammer and a spanner if you are to fix a broken pipe or build a wall. Different tools always do different jobs. The same is true with statistics; you need to be able to identify and use the right level of measurement because the appropriate descriptive and inferential statistical methods that you will use to analyse your data will differ for nominal, ordinal and interval variables. If you use the wrong tool to fix your pipe, the chances are that you won't be able to fix your pipe and you will come to the conclusion that you need a new one. But this would be the wrong conclusion to make. Equally, if you use the wrong statistical technique the likelihood is that you will come to the wrong conclusion about your findings. Your ability to investigate your hypothesis, and subsequently achieve your research aim, will be severely limited.

For instance, it would not make sense to calculate an average ethnicity. This is because there is no intrinsic ordering of the levels of the categories – there is no mid-point. We could say which particular ethnic group is most or least common, for example, or we could work out the proportion of people in each particular group, but we could not talk about a meaningful average.

We'll deal with these statistical techniques in more detail later, but for now we will concentrate on another reason that it is important to identify and understand levels of measurement. We need to be able to think critically about what a variable is actually measuring – we need to have an idea of how to assess its reliability and its validity.

2.4. Thinking critically about variables

2.4.1. Validity and reliability

Broadly speaking, the validity of a variable refers to the ability of the variable to measure what it suggests it is measuring. On the other hand, the reliability of a variable refers to its stability and its consistency. Thinking critically about the reliability and validity of a variable is crucial in assessing how good our answer to our research question is, and whether we have actually achieved our research aims. Think about it: if a variable is not measuring 'the thing' it

is supposed to then it is not a good measure of that 'thing'. Similarly, if our measurement of that 'thing' varies wildly each time you measure it, then we will struggle to say anything meaningful about it because by the time we have measured it, the answer will have changed. Therefore validity and reliability are crucial to understanding the usefulness of a variable.

Of course, we can't exhaust the range of potential problems with all variables here - there are simply too many variables to do this - but we will demonstrate how you can think critically about the validity and reliability of variables, and why it is important to do so.

2.4.2. Labelling variables

One way of thinking critically about the validity of a variable is to be very specific about what a variable is doing when we label it. The more specific we are about what a variable is actually doing, the less likely we are to go beyond what that variable is actually telling us. **Let's go back to the original smoking examples.**

Q2.1. Do you smoke cigarettes?  
 Yes  No

Q2.2. Do you consider yourself to be a:  
 Non-smoker  Light-smoker  
 Medium smoker  Heavy smoker

Q2.3. How many cigarettes have you smoked in the last seven days?  
 \_\_\_\_\_ (write here)

The temptation might be to label these variables: smoking status; level of smoking; amount of smoking. However, these labels are not necessarily valid.

The first thing to note is that they are all self-reported measures. That is, they deal with perceived smoking behaviour not the actual behaviour. Perceived smoking behaviour is not always a reliable predictor of actual smoking behaviour. Indeed, there is a good deal of evidence to suggest that risky behaviours are commonly under-reported in survey questions of this type.

Therefore, for us to be able to say that reported

measures accurately reflect the 'actual' number of cigarettes smoked, three assumptions would need to be met. Firstly, that people can recognise and interpret their behaviour meaningfully and consistently; secondly, that those interpretations can be meaningfully and consistently mapped on to the measures we have built; and thirdly, that respondents will choose to represent their interpretations as accurately as is possible. Can we say this with confidence?

It's actually fairly easy to question all these assumptions. For the first assumption - do people recognise and interpret their behaviour meaningfully and consistently - how many people do you know who smoke when they drink alcohol and still consider themselves to be non-smokers? What about people who only smoke cannabis? What about people who seem to be in a constant cycle of 'trying to give up' and 'giving up' - how long does it take before you can actually class yourself as actually given up?

For the second assumption - can the interpretations be mapped meaningfully and consistently on to the measures we have built - what actually constitutes a medium smoker? If we asked the question again two weeks later, would they respond the same? For that matter what counts as a cigarette?

For the final assumption - do respondents choose to represent their interpretations as accurately as is possible - are people really going to admit to their smoking levels or are they likely to under-estimate it? Besides, don't smoking levels vary from week to week? This final difficulty is basically a problem of reliability - can we really infer accurate levels of smoking from one measure at one particular point in time? This might seem to be a little pedantic, and perceived levels of behaviour can sometimes be good predictors of actual behaviours, however, measures that deal with social behaviour are often much more prone to problems of reliability and validity than their 'natural' counterparts - a red blood cell does not react to the fact that it is being counted but a person might.

Think about the difference between these two statements:

**T** 57% of University of Sheffield students smoke less than 10 cigarettes a week.

57% of University of Sheffield students in the present sample reported that they smoked less than 10 cigarettes a week.

The difference is subtle, but the inclusion of 'reported' and 'the present sample' in the second case is important. Accurate naming of variables can help you to avoid making spurious statements of fact with impressive numbers to match - this is all too easily done when reporting statistical data. Whenever you see statements such as these two examples, you need to be asking the questions: 'how is that variable being measured?' and 'does that variable allow us to come to that conclusion securely?'

As a result of all of this, it is probably better to name our three variables: self-reported smoking status; perceived level of smoking; estimated level of cigarettes smoked per week.

However, it is not just the labelling of variables that you need to be aware of when working with variables. Having a critical understanding of any potential problems of measuring that variable is also crucial. Let's look again at the ethnicity question that has been proposed for the 2011 census on the next page.

Let us return to the validity assumptions we suggested above:

- Is our sample likely to recognise and interpret their behaviour meaningfully and consistently?
- Can those interpretations be meaningfully and consistently mapped on to the measures we have built?
- Will the respondents choose to represent their interpretations as accurately possible?

**Q** Take each assumption and think about it in relation to the ethnicity question that was asked in the UK Census of 2011: Can you identify any problems?

Q2.4. What is your ethnic group?

Choose one section from A to E, then tick one box to best describe your ethnic group or background

- A. White
- English / Welsh / Scottish / Northern Irish / British
- Irish
- Gypsy or Irish Traveller
- Any other background, write in: \_\_\_\_\_

B. Mixed / Multiple ethnic groups

- White and Black Caribbean
- White and Black African
- White and Asian
- Any other Mixed / multiple ethnic background, write in: \_\_\_\_\_

C. Asian / Asian British

- Indian
- Pakistani
- Bangladeshi
- Chinese
- Any other background, write in: \_\_\_\_\_

D. Black / African / Caribbean / Black British

- African
- Caribbean
- Any other Black / African / Caribbean background, write in: \_\_\_\_\_

E. Other ethnic group

- Arab
- Any other ethnic group, write in: \_\_\_\_\_

Ethnicity is not a fixed or natural category: it is a social one that needs to be understood in relation to its context - it is certainly not as rigid as this question might suggest.

In answer to the first question, the categories presume that ethnicity is largely related to a stable and linear line of ancestral origin. That is certainly part of the answer, but ethnicity is much richer - and more complicated - than this. Look again at category C "Asian / British Asian". Here we have a tick box for the category 'Indian'. As you may know, India is a huge country, with many distinct regions that have very distinct traditions. Similarly, immigration from India to Britain has been on-going since at least the 18th century and has historically been comprised of a diverse range of people including princes, servants, students, seamen, doctors, and performers. Whilst some are relatively recent arrivals, some have lived in the UK for over four generations. Further still, whilst recent immigration has been comprised of economic migrants who came to work in the burgeoning manufacturing trade, many are actually professionally qualified political immigrants

who came via East Africa in the 1960's and 1970's.

However, this ethnicity question ignores all of the potential difference within groups and largely lumps everyone together in a series of homogenous groups. Indeed, it is not hard to question how this measure can adequately represent all the potential difference in an increasingly globalised world.

In answer to our second and third assumptions, we could also question whether people will accurately represent this information. Many excluded groups - which includes some ethnic groups - have a particularly high rate of non-response in the census. In fact, a pilot ethnicity question that was administered in the 1979 census met with such hostility, the question was removed from the 1981 census. Although the question met much less resistance in 1991, some people still refused to answer it, or simply couldn't. Think about the various mixed ethnicities we can create. Imagine someone born in Britain from a first generation Indian mother (by way of Uganda) and a second generation father of mixed White and Black Caribbean heritage - where would they fit?

Of course, not all data is of the self-reported type, and this is not to say self-reported measures are useless - far from it - but any variable is a pragmatic solution to difficult problems of measurement. We have to measure social phenomena somehow, and whilst these measurements might not be perfect they are often the best we can manage given the circumstances. The question for ethnicity was, in part, designed to evaluate the 1978 Race Relations Act and it does allow us to see general patterns of inequality across the country. It is a problematic variable, but it is also a potentially useful one if we recognise its limitations. We are not trying to put you off using variables like these in your research, but you need to be aware of any potential issues/limitations with variables you choose and the possible impact that this could have on your research

Indeed, a failure to recognise the problems of any variable can result in you going beyond what is actually being suggested by the data that you get from your variables. Hence we need to think critically about any variable that we choose to use. To paraphrase Deming (1986), a famous statistician: the better we understand the limitations of an inference...the more useful becomes the inference. This starts with a clear understanding of what your variables are actually measuring - and the devil is often in the detail.

2.4.3. *Research questions, hypotheses, and variables: Working examples*

Variables and research questions and hypotheses are a lot like the cogs in a machine. The cogs allow the machine to work and the machine gives the cogs purpose - but without each other they are quite limited. Similarly, variables are not too much use without a hypothesis and a hypothesis is not much use without a research question. Therefore we need to make them work in tandem with each other.

Suppose that I'm interested in the relationship between gender and body image - this is my research aim. More specifically, I want to discover whether women are more likely to be dissatisfied with their body image than men - this is my research question.

In order to investigate this, and subsequently achieve my research aim, I'm making the assumption that gender is related to body image - or not: this is, effectively, my hypothesis. So now I need to construct two variables to help to assess it. Firstly, I need to build a variable that will measure the gender of the respondent. Secondly, I need to construct a variable that will allow me to measure how satisfied the respondent is with their body.

Identify the level of measurement required for these variables and attempt to construct them.

**As gender is a categorical variable, I can use a nominal level of measurement. So:**

Q2.16. Please state your gender: e

Male  Female

Having looked at some similar studies in the literature, I find that there is methodological precedence for a measure of body image satisfaction that uses a five point Likert scale. Hence I am measuring my body satisfaction variable at the ordinal level. **So:**

Q2.17. How satisfied are you with the way you look: e

A. Very dissatisfied  
 B. Dissatisfied  
 C. Neither satisfied or dissatisfied  
 D. Satisfied  
 E. Very satisfied

Look again at the process we have worked through in this example. I begin with the original idea - my research aim; from that idea I then develop a hypothesis; using this hypothesis I can then identify the variables that I need to measure; from these variables I can then identify the level of measurement that are required by these variables; finally, I construct a question with respective answers.

2.5. *Rounding up...*

This workbook has introduced you to the different ways in which we can record quantitative material, and how you can use variables to help you to answer your research questions and hypotheses.

You should now be able to:

- Identify different levels of measurement
- Be able to construct variables at level appropriate to your research aims
- Understand issues of reliability and validity and think critically about them in relation to variable measurement
- Understand the role of a variable to answer your research questions and hypothesis

By now, you should be able to generate a research project, write a rationale and accompanying aims, questions and hypotheses, and design some survey questions that are appropriate to your study. However, you're not quite ready to go yet. Before you begin collecting data, you need to be able to have a good understanding of analysis and how you will make sense of your data. Indeed, the presentation of the research process looks to be a linear one. Before you carry out any data collection, you need to understand how you will analyse your data. The next book in this series will introduce you to some of the main issues of analysis that you will need to answer your research questions.

This workbook by Tom Clark and Liam Foster is licensed under a Creative Commons Attribution Non Commercial - ShareAlike 4.0 International License.

Contains public sector information licensed under the Open Government Licence v2.0.  
Crown Copyright.