

## **Part A. EpiData Entry**

### **Part A: Quality-assured data capture with EpiData Manager and EpiData EntryClient**

- Exercise 1 A data documentation sheet for a simple questionnaire
- Exercise 2 Create a basic data entry form
- Exercise 3 Create a value-label pair from external data
- Exercise 4 Create a composite identifier
- Exercise 5 Data entry and validation
- Exercise 6 Upgrading an EpiData 3.1 REC/CHK file pair to an EPX file
- Exercise 7 Relational database

#### **Acknowledgments:**

We thank Ajay M V Kumar who had made valuable suggestions to improve the structure and flow of argumentation of the preceding version (using EpiData Entry 3.1) of Part A which we partially incorporated into this revised version.

**Exercise 1: A data documentation sheet for a simple questionnaire**

The first step in the process is to prepare a plan for data entry. This plan is called the **data documentation sheet**. This should not be confused with *data collection form* or *case report form* which is the proforma used for collecting the data from study participants or extracted out of the program records. The data documentation sheet is a codebook containing the details of all the variables (like field names, field labels, field type, field length, possible field values, and field value labels) to be entered.

*Note: Please do not worry if you do not yet understand all the technical terms at this point in time. Be assured that you will appreciate these as you go along.*

But let us proceed step by step and say that we have the following questionnaire:

Laboratory serial number: \_\_\_\_

Date specimen received (dd/mm/yyyy): \_\_\_\_/\_\_\_\_/\_\_\_\_

Sex: \_\_\_\_

Age in years: \_\_\_\_

Reason for examination: \_\_\_\_

Result of specimen 1: \_\_\_\_

Result of specimen 2: \_\_\_\_

Result of specimen 3: \_\_\_\_

This might present a typical simple questionnaire as used by an interviewer. Often such questionnaires are first completed on a paper *Case Report Form*. The above is actually an excerpt from the original Tuberculosis Laboratory Register proposed by The Union (note that the current version has been slightly changed):

Tuberculosis Programme

Form 2

**Tuberculosis laboratory register**

Year \_\_\_\_\_

Lab Serial No.	Date specimen received	Name	Sex M/F	Age	Name of referring facility	Address - patient for diagnosis	Reason for examination*		Results of specimen			Only for SS+ for diagnosis: TB Number or treatment centre**	Remarks
							Diagnosis (tick)	Month of follow up	1	2	3		

We will use this register as the basis for this course. For the time being, you plan to write a short and concise electronic data capture form, retaining only variables that are easy to capture and are likely to be useful for the analysis. *Please note this as a first principle in being efficient – capture only those variables which you will use for analysis!*

Each of the questions can be conceived of as a variable and the answer to the question as the value that the variable takes for a particular individual. **Variables** are also referred to as **'Fields'** in EpiData – both refer to the same and will be used interchangeably. We will give each variable a unique name. A completed entered data form for one study subject is called a **'Record'**. A set of such records is called a **'data file'**. The data file thus contains several records and each record contains information about one individual with respect to several variables. We are going

to describe each variable with respect to several attributes in the data documentation sheet. Let us now understand some terminology we are going to use.

- **Field name:** This is the name of the variable and in EpiData, there are certain rules to be followed in arriving at this name. We will come to these rules in a short while.
- **Field Label:** This is the descriptive name for the variable and contains a more detailed description than the variable name / field name can convey.
- **Field Type:** This describes the type of the variable – text, numeric or date being the major types.
- **Field length:** This describes the number of characters that a value can take.
- **Field values:** This describes the possible values that a variable can take.
- **Value labels:** These are descriptive names for the values. For categorical variables which are numerically coded, it is always useful to label them so that it is easier to read and understand what each of the codes mean.

“**Labels**” are also called “**metadata**” or “**data about data**”. They play a key role in data files. We may have entered a value “9” for a given field, but this number remains meaningless without specifying for what this value stands. It is important to get acquainted with these terms and understand them clearly since we will be using them frequently. We will be using several examples later in this chapter to clarify these terms.

**Field name:** There are some software-specific rules in naming a variable.

First, a Field name has to be **single word** that has **not more than ten characters**. This means that you cannot use a space in the name as a space makes it more than a word. Also, you cannot use any special characters like comma, semicolon, full-stop or underscore, and the first character should not be a number. Thus, do not start the variable name with a number. It cannot be ‘1v’, but it can be ‘v1’.

Note that some other analysis software may accept only a field length of eight characters. If you later plan to export your EpiData files for analysis to such a software package and you had used the full field length of ten, then your field names get truncated.

Second, use a name which is **intuitive** to understand what it means instead of generic field names like v1, v2 and so on.

Third, it may be a good practice to keep the field names in **lower case**. While EpiData is not case-sensitive, some other software is. Important among such are e.g. R and Stata® which are what we call “case-sensitive”. In the latter two, a field name of ‘sex’ (lower case), ‘SEX’ (Upper case), and ‘Sex’ (mixed case) are understood as different variables. If you later plan to export your EpiData files for analysis to such a software package, it may make life unnecessarily complicated if you have been inconsistent without a defined rule. Hence, the recommendation to keep it uniformly, “lower case”.

The following words ‘date’, ‘month’, and ‘year’ are functions in EpiData and are reserved names. Hence they cannot be used as variable names.

**Field label:** This is the full description of the variable and can be more than a word.

**Field Type:** There are different types of entry fields for the variables (we will follow the EpiData notation and call them “Fields”):

- **Text fields:** These fields take letters or numbers or a combination of these as possible values, like PETER, KOCH1882, giraffe, 45677 etc. You can type anything on the keyboard into this field. If you enter a number into such a field, it is accepted but you will not be able to make any calculation with it. These fields are also sometimes designated as character or alphanumeric fields, or most simply “**string**” (denoted by **S**) fields as they take any string of characters.
- **Numeric fields:** These are numbers. The numbers might be integers (denoted by **I**) like 885, 33, 1235 or real numbers like 3.4, 6.88, and 66.5 (also called **floats** and denoted by ‘**F**’). You can make calculations with numeric fields.
- **Date fields:** (denoted by ‘**D**’): In different countries, different ways of writing dates are used and this can be confusing for people from another culture. Some write *5 March 2005*, others *March 5 2005*, and again others *2005 March 5*. EpiData lets you choose the type of date you wish to take. We made our choice for European dates in the Preferences as we will be using European dates in this course, i.e. dates of the format *5 March 2005* or symbolized with DD/MM/YYYY.
- One other type of variables is called “**logic**” or “**Boolean**” variables. This is sometimes used in food-borne outbreak investigations. There, answers to questions on food items eaten is limited to “yes” and “no” and “missing”. There is no need for using this field type and we discourage its use as it might pose problems in analysis. The alternative is a numeric field with a label block.

While you are asked to limit the length of the field name, you have much more flexibility with the length of the value a field can take, but we will try to use it as efficiently as possible, that is we will limit the value length to the minimum needed.

## Data Documentation Sheet

It is good practice to write what we call a **data documentation sheet** before you make your actual data entry form in EpiData Manager. As mentioned earlier, EpiData refers to this as **Codebook**.

In the past, fields like `sex` were commonly made text field with values “F” or “M” denoting Females and Males. It is efficient as it uses only a length of 1 to remain unambiguous (well, as long as the language is English or French, at least!). Things would get less efficient, if we would have to code treatment outcome with the possible values “cured”, “completed”, “died”, “failed”, “lost from follow-up”, “transferred out”, and “outcome not recorded”. Numeric coding is much simpler as there are up to ten possible values with a field length of just 1:

- 1 Cured
- 2 Treatment completed
- 3 Died of any cause
- 4 Failed bacteriologically
- 5 Lost from follow-up
- 6 Transferred out
- 9 Outcome not recorded

Later, in the analysis, you will also realize that it is very convenient to apply numeric selection criteria when you want to select a subset of data and undertake analysis only on the subset. Of course, a prerequisite is that the link between the numeric value and the text label is unambiguously clear. The role of labels is of enormous importance, also called meta-data or

“data about data” as mentioned above. We are going to make full use of numeric coding of field values and using explicit text as value labels.

Let us now go through a few examples of the variables (from the tuberculosis laboratory register example) and describe the various attributes of the variables in the data documentation sheet. As you go through, you will note that preparation of data documentation sheet requires thinking and knowledge of study data.

This is how we would write such a data documentation sheet:

Field name	Question [Field label]	Field type	Field length	Field values	Value labels	Notes [Comment]
serno	Laboratory serial number *	I	4	1-9000, 9001, 9002,		Serial number starting with 1 each year Enter 9001, 9002,... if serial number is <i>not unique</i> or <i>missing</i> , and write a data entry note
regdate	Registration date	D	10	01/01/2000-31/12/2005		Range of legal registration dates
sex	Examinee's sex	I	1	1 2 9	Female sex Male sex Sex not recorded	

\* **Note:** Often, it will be preferable to make the identifier a text field. If it is a number, as in this case here with the laboratory serial number, precautions must be taken to distinguish e.g. “0001” from “1”, requiring that the numeric value is entered into one field. We can make use of the possibility in EpiData Manager to add leading zeros where appropriate.

**Task:**

- o Complete the data documentation sheet for all fields in the questionnaire. Note that you should always define a value if no answer was provided to a question.*
- o Think of the most efficient ways to code reason for examination and results of microscopic examination*