# Writing research protocols: a statistical perspective

## I. M. Wilson & S. Abeyasekera

## January 2006

**The University of Reading**
**Statistical Services Centre**

**Guidance prepared for the DFID**
**Forestry Research Programme**

# Foreword

In carrying a project from inception to completion, one link in the chain is the analysis of data collected. If the material for analysis is inadequately organised or documented, this process may be confused, or done badly. This guide aims to reduce such problems at an early stage.

The material below probably has most to say to relatively inexperienced researchers, such as Ph.D. students taking on particular activities within larger projects. Of course we are well-aware that experienced project leaders already do much of what we are suggesting, insofar as it is relevant to the special circumstances of their project, but we hope that at least some of our points will be helpful to them too.

This, and the parallel document entitled *Writing up research: a statistical perspective*, have been funded by DFID's Forestry Research Programme (Project ZF0183). They form part of a series of guides funded by DFID and written by Statistical Services Centre [SSC] staff since 1998 (see SSC website at http://www.reading.ac.uk/ssc/publications/guides.html ) where more than 20 such guides can be found.

In 2001-2, SSC staff were commissioned to write a review of the *Use of biometrics by DFID RNRRS research projects*, as revealed in Final Technical Reports and accompanying submissions to several RNRRS Programme Managers. These new guides were commissioned by the DFID RNRRS Forestry Research Programme [FRP], following comments about all the projects reviewed.

Our analogy was that Programme Management and research project leaders were primarily concerned with the above-decks aspects of the voyages represented by each project, while the statistician and the data manager are below decks working in the engine room of the ship, and it is from this position that we present our views now.

As outsiders reading DFID Final Technical Reports, we were well aware that we did not have the disciplinary skills and situational knowledge that went into the definition of project aims, and the selection of key results. We concentrated on various aspects of the internal processes of the projects we looked at. We considered primarily statistical and data management issues, including sampling and design aspects. We asked whether the data collected were necessary and sufficient to meet stated objectives or to support claimed conclusions. We asked whether analyses reported were appropriate, and the results fit for purpose.

We concluded that generally too much time and resource had gone into data collection, too little into research design and early-stage planning. In many cases, injudicious organisation and poor documentation of plans for data collection, data management and analysis led to fewer, and poorer quality, outputs than could have been achieved with a little less data, well used. Frequently, there was limited sign of a publicly-available legacy from the projects that covered anything like the full range of their work. Scarcely any of the projects signalled that they had developed a systematic archive of their data.

This guide is a plea for more systematic and more complete recording of planning information, in the form of a protocol, to ensure that the project process – especially the analysis stage – can be efficient and effective, and that the legacy of the project does not languish unusable after the exercise has ended.

We do not set out to challenge the training, judgment and experience of those who develop the processes and select the best methodologies. We recognise and happily acknowledge the importance of achieving and using best-practice standards suited to the central discipline of the research. Our concern is to ensure that project procedures and results are well-recorded, and in particular that they are conveyed in accurate form to data analysts and users.

# Acknowledgements

# Index

## Boxes

# 1.  Introduction

## 1.1  What is a protocol?

We regard the writing of protocols as being an essential component that needs attention throughout a project's lifetime.  The Chambers Dictionary definitions of the word 'protocol' include '*an official or formal account or record*' and '*a factual record of observations, e.g. in scientific experiments*'.  In this guide we take a relatively wide view of what this implies.  Broadly we view a protocol as a document that aims initially to capture the planned aspects of the project, to be later updated as the project progresses so as to reflect what actually happened, and where possible to facilitate the reporting regime.

A project will typically have several protocols, some of which are relevant at the project level and some at the activity level.  These will provide a complete and competent description of all project activities and will represent a true record of the status of the project at any one point in time.  Each protocol must link clearly and effectively to other activities and stages in the project cycle, and often to other parts of an overall research programme.

Although we draw upon experience of the old DFID Renewable Natural Resources Research Strategy (RNRRS) in places, the material that follows is intended to be quite wide in its applicability.   We are thinking about a substantial 'project' with responsibilities for primary collection of non-trivial information, probably involving staff with varied roles, several different sorts of disciplinary background, and probably several nationalities, maybe working in a number of different countries.   We assume (as in the logical framework) that much of the project work is organised into activities, covering facets of the project.   Our comments are not intended to be restricted to RNRRS projects or to those funded by DFID.

## 1.2  Who needs well-organised information?

As statisticians, our primary concern in this guide is with ensuring the right data are collected and with having documented information readily available to undertake timely and effective analysis of the material that the project gathers.

In Boxes 1.2.1, 1.2.2 and 1.2.3 below we outline a few of the *other* uses of the information whose good organisation and documentation we advocate.

Box 3 is intended to cover situations where there may be a role for a specially-constituted body in each host country or region which might be called a 'Project Advisory Committee (PAC)' or an 'external panel'.  It may, for example, provide linkages to a number of organisations over and above the immediate host, provide access to local knowledge, secure support and buy-in from influential parties, give voice to various stakeholders, and help to create a 'market' for project outputs.   We see project documentation as being key to serving the PAC in its role.

There can also be needs for specialised forms of well-documented information.  For example, the project leadership staff who have to appraise job-holders or to recruit new staff will gain from documentation which effectively specifies what each position entails and requires.   This will also help with the induction of, or handover to, new staff or the definition of terms of reference for specialised consultants.

**Box 1.2.1  _Internal_ Uses of Well-Organised Project Documentation**

- Project management and project monitoring tool e.g. looking at timescales, budgets, staff performance, emerging crises

- Guidance to those managing activities within the project

- Systematic creation of records e.g. of decisions taken, e.g. of sources on which current thinking is based, and the development of an 'institutional memory'

- Communication and linking between project teams, establishment of common understandings of objectives or even vocabulary

- Continuity, where staff changes or re-organisation might undermine this

- Orderly handover to local successors after the 'multi-national' project is over

---

**Box 1.2.2  _External_ Uses of Well-Organised Project Documentation**

- Evidence is readily, and quickly, available to external reviewers or evaluators, or when seeking new phases of funding

- Establishing liaison with other entities external to the project is easier, and can be based on better mutual understanding

- Generation of periodic interim reports and a final technical report for a funder is made easier

---

**Box 1.2.3  Project Documentation & a _Project Advisory Committee_**

- If a PAC has any monitoring role, careful records allow them to know, and judge, what is being done

- If the PAC is used extractively to supply inputs e.g. of knowledge, introductions, or advice to the project, the relevance and usefulness of such inputs can be improved where it is clear how they fit in

- If a PAC has a role in facilitating the project, they should have clear enough knowledge of what is under way to ensure project outputs are in suitable forms to achieve acceptance and impact in the context of local policies, institutions and processes

- If a PAC acts to champion the project, they should be able clearly to explain, and advocate for, its legitimate needs so that potential obstacles are removed

- If the PAC helps with dissemination and advocacy of outputs and results, they can create anticipation of, and interest in, the findings and if necessary moderate unreasonable expectations

---

## 1.3 Types of protocols and extent of information needed

We have already referred to the need to have project level and activity level protocols. At the project level there is a need to have an integrative 'project protocol' that will encompass the needs of the project as a whole, list all the activities needed to deliver project outputs, demonstrate how the activities link together to contribute collectively to the overall research goal, and how results will be disseminated and used.

At the project level, there is often also a need to have a 'sampling protocol', and a 'data management' protocol. The former will describe the broad approach to sampling, i.e. how measurement units are selected, giving appropriate justification, it will specify how many units will be sampled and why, and so on. The latter will include procedures for monitoring progress on data collection activities, and will indicate who will be responsible for data management activities, and when and how the data will be archived. Section 5 below expands on these aspects.

Activity level protocols will generally be internal to the project team and will serve to describe the procedures and justification of each component of the activity clearly and completely. In later sections we give more details about what aspects need to be documented within activity protocols.

Clearly, the status of 'the project' makes a great deal of difference to the form and extent of the information that needs to be organised and preserved. If a project exists to serve a one-off and temporary need, less is required than if it is to be the model for replications to be carried out by (a) so far unidentified successors or (b) numerous separate geographical teams.

An example might be the development of a crop forecasting system developed and tested using the expertise of an international organisation, to be handed over to a Ministry of Agriculture for annual use nationally in food-insecure districts.

Where a project has to link together different discipline-based components, there is a need to consider and cope with varied understandings of terminology, methodology, time-scale and approach. Failures to put sufficient thought and effort into integration of disparate discipline-based contributions through developing common documentation can easily lead to poorly synthesised final reports with obviously weak conclusions.

## 1.4 Good documentation to help the data analyst

'The data analysts' are one group for whom good documentation is absolutely essential in a project that collects non-trivial amounts of primary data. The selection of analysis approaches and results has to be based on a clear idea of how uptake and impact are envisaged. Those undertaking data management and analysis certainly need to understand the audience for the results they help to achieve, and the ways in which the outputs may be used.

The project whose professional staff cover a wide range of disciplines, and the lengthy or longitudinal project, have to consider and distinguish a range of types of output. Forethought is required – preferably in consultation with the data analyst – to ensure the appropriate distribution of data collection and analysis resources, in order to meet what may be conflicting requirements posed by different types of output such as those in Box 1.4.1.

*Box 1.4.1  Different Types of Output*

1.   Research projects are increasingly expected to lead to meaningful policy-relevant contributions that will have a poverty impact.   A policy brief usually has limited scope to include details of methods and approaches used to reach the key conclusions.   Yet if successful it may be highly influential, and it is very important that it should be accurate – and defensible if subject to serious scrutiny.

Thus it is usually essential that results are backed up by a body of well-documented technical work that supports the accuracy and reliability of the claims and recommendations made.   A completed formal report and/or a paper submitted for peer-reviewed publication should usually be available by the time of the launch of the policy brief.   The project team must be confident that the data used are fully checked and thoroughly clean, and that the conclusions placed in the public arena will definitely not be materially affected by any errors that could subsequently emerge.   The team must allow the analyst time to ensure this is so.

2.   Relatively long-term research usually involves work on other substantive issues where conclusions are some distance from mature or complete form, but which may require writing-up  (i) to establish interim results, or (ii) to establish what needs to be done and why, in order to move towards final conclusions.

3.   Even by the time a project is completed, the outputs suggested in 1. and 2. above are likely to cover quite highly selected themes, rather than the range of data collected by a large project.   Funders usually regard data collected using public money as constituting a public good and it may be desirable to put in the public domain information outputs covering, as far as the project reasonably can, the range of topics on which the project possesses reliable information that is not otherwise readily available. From the research leaders, intelligent awareness is needed of public sector and public interest groups' likely concerns.   These need to be communicated to the statistical analyst, for consideration in relation to presentation quality, information accuracy and interpretability.

4.   Projects often develop new 'instruments', for example a questionnaire-based asset index or a method of farmer-based data-collection with occasional quality control by project staff.   These may be justified, and made available to others, in papers or reports which are 'methodological' rather than 'substantive' project outputs.   A statistical analyst may have a quite important part to play in the work needed to appraise and validate the methods or instruments being reported.

5.   Where the project needs to coordinate information collection across widely spread sites, or to achieve consistency through time, a fifth class – of documents of record – is likely to be essential to internal users and to others, such as Ph.D. students, who may take up project themes.   One of these documents concerns the sampling selection scheme and details of its implementation.  Other documents of this type can include (a) archived datasets and accompanying metadata, and (b) analysis programs – and in the case of field surveys (c) questionnaires used on each occasion in each country, details of how responses were coded, and composite variables derived, and (d) fieldwork reports. These records are clearly of importance to the work of the analysts, and it is often sensible that they participate from the start in document preparation.

## 1.5   Contents of this guide

In this guide, we assume the project includes at least components where primary data collection takes place: our particular concern is to ensure that data collection and analysis activities are well-organised, purposeful and systematic.   The analyses carried out should be fit-for-purpose i.e. necessary and sufficient to provide good-quality project

---

outputs, *and* the data collected should be fit-for-purpose i.e. necessary and sufficient to provide good-quality analyses.  For example, data collected, but never used in analysis, is a waste of project resource.

In sections above, we pointed out reasons for well-organised documentation, and gave quite an inclusive definition of the term 'protocol.'  We considered how the information may be used internally in a project team, by a Project Advisory Committee, and for other external purposes.

In broad terms, we recognised the need for 'activity protocols' for the activities specified in a logical framework, and an integrative 'project protocol' indicating how activities and their outcomes are drawn together.   Activity protocols are more straightforward and are discussed in section 2.

Activities may be relatively self-contained, but to set them in context effectively, their motivation, scale, timing and outputs need (a) to fit with predecessor, parallel and successor activities, and  (b) to ensure they contribute appropriately to higher levels in the research framework.   The project protocol and these links between one activity and another, or between activities and the whole, are discussed in section 3.

In the text, we consider elements that we believe are essential ingredients of a good protocol.   These concern the 'What?', 'When?' 'Who?', 'Where?' and 'How?' as well as the 'Why?' element of research.   The first five questions are brought into section 2, but the main comments on 'Why?' are in section 4.  The discussion in section 4 is in relation to activities and project together, since it concerns largely the overall logic of the project and the inter-connections of its activities.

In section 5, we concentrate on 'engine-room protocols' for the major exercises needed to ensure data collection projects are competently organised.   At the relatively early project stages where protocols are being conceptualised, the most immediate of these are for sampling and data management.

Included thereafter is a series of appendices giving rather more extensive examples than could be included in the body of the text.   These provide illustrations of a selection of points in the main text.   They are referenced at some relevant points in the text.   To make these examples coherent in themselves, they have to include a certain amount of explanatory material.   In all cases the examples are nevertheless very incomplete, in most cases representing small fractions of much longer documents.   We are very grateful to those who have allowed this limited and selective reproduction from their work.   They are individually acknowledged in the context of the extracts presented.

# 2.   Activity-Level Protocols

## 2.1  Introduction

For simplicity, we look first at the case of a single activity within a project.  This is discussed in a generic way, not specific to any discipline or type of approach to the acquisition of primary data.   The information requirement for a basic factual record at activity level entails including a description of **what** the activity is, **when** and **where** and by **whom** it will be done.  **How** it will be done is discussed briefly in 2.2.4 and 2.2.5 below.   Sub-section 2.3 includes a few comments on **why** questions at activity level, but most discussion of these is in section 4.   Most of what we discuss in section 2 would normally underpin, and maybe be summarised in, the "Methods and Materials" section of a research publication.

'The activity' might be a livelihoods survey in the project operational area, or it might be a series of farmer participatory trials of a new approach to green manuring in conservation farming.   It could be a socio-economic focus group exercise looking at diversification opportunities, or a laboratory study of cell culture of genetically diverse banana planting material.

We visualise that the activity might have its own team of workers, perhaps led by a doctoral student or some other team member responsible to the project leader.   At this stage we assume this activity leader does *not* have the main responsibility for deciding the broader structure of the project beyond his activity, or for how his activity fits into this bigger picture.

## 2.2  A basic factual record

### 2.2.1  Activity-level factual record

The factual component of an activity protocol requires information at activity level, and at various levels below that till we get down to a single 'datum' – for example, a cell entry in a questionnaire corresponding to one respondent's answer to a single simple survey question.

At the activity level, the essential summary is a clear-cut description of what is to be done, when and where and by whom, so that – as a rather extreme example – when the agronomist author is (fatally) run over by a tractor, a competent person brought in as a replacement should be able to reconstruct the agronomic experiment that was to be done, and carry through the activity without any serious doubts about how to proceed.

The author could legitimately have assumed his successor, and most committed reader, would have relevant general skills e.g. Ph.D. in Agronomy and several years experience, but he should also have thought about other readers and what they would need to know to fit the work into their own frameworks – the funder, the project manager, the research assistant, for instance.

So what constitutes a competent description at activity level?

Of course any given discipline and type of study will require records which are fit-for-purpose, and these will differ from case to case.   For example a focus group exercise in which a farmer participated for one afternoon might well require far less information about him than a farmer-managed field experiment where his land, his agricultural management practices, and his self-collected data are the subjective of substantial project investment over a season or more.   For instance, in the example activity protocol in Appendix 2, a substantial baseline survey provided background information about the 6 farmers, and those involved in other activities, while in the example in Appendix 3 only a few basic socio-economic characteristics were measured on each respondent (see A.3.11)

This document is not concerned with specifying exactly what should be recorded by specific types of project, or in a particular type of organisation.   We leave it to readers to supply their own discipline-specific elements of methods, materials and procedures, though some illustrations are given in our sample extracts from protocols in the appendices.    We assume that procedures, including measurement procedures, will follow best practice standards for the relevant discipline. Of course such processes need to be properly validated where they are not well-standardised and proven.   That is a topic for another guide!

Our goal here is to provide a small selection of examples and to list some *generic* elements.   Box 2.2.1 gives some elements of an activity protocol:-

---

**Box 2.2.1  Some Elements of an Activity Protocol**

ACTIVITY TITLE:  *<A title for the activity>*

ACTIVITY LEADERSHIP:  *<Name of scientist, research assistants, technicians etc. responsible for the activity>*

BACKGROUND:  *<Background to the activity and justification>*

*(Comment: This may include reasoning and literature behind the choices of variables. This background provides the link to the project protocol. The latter should state in broad terms what information is needed from the outputs of the activity, but it need not contain details of the approach.   The activity protocol should record justification for decisions about how the activity is designed.  Details of the decisions include the following.)*

OBJECTIVES:  <A clearly specified statement of the objective(s) of the activity>

MATERIALS AND METHODS:

<    (i)   Location(s):  Where is the activity being carried out?

(ii)  …

(iii)  …   >

DATA: <What information is to be collected, how, when and why, data collection procedures and instructions to field or lab. staff involved>

---

Example: If the activity was a crop-based on-farm experimental study, the MATERIALS AND METHODS heading might include:

*(ii) Important dates:*  Start and end dates, planting dates, dates for field staff training, pilot testing, data collection, etc.

*(iii) Study design details:*  Experimental design layout, description of blocking factors, choice of plots within the farm, who manages the trial.   Sample size rationale.

*(iv) Materials to be used:*  What inputs will be used, where from, and how prepared (if relevant), data collection instruments and so on.   See A.2.6 for an example.

*(v) Responsibilities:* Who will be responsible for (a) marking and layout of plots, (b) planting operations, (c) taking measurements, and recording/checking these, (d) supervision.

*(vi) Field Instructions:* How to select samples consistently e.g. soil samples, plants/leaves for disease assessment.   See A.2.6 for an example of this type.

The section on **DATA** might list types of data to be collected (e.g. socio-economic, labour use, climate data, disease assessments etc), how they relate to overall activity objectives, measurements to be made within each data type (e.g. agronomic performance for maize measured by seed weight, number of cobs, stand count, plant height).

Example: If the activity was a qualitative consultative exercise based on focus groups, the **MATERIALS AND METHODS** heading might include:

*(ii) Important dates:*  Activity schedule for field operations,

*(iii) Study design details*:  Selection of sites, rules and procedures for recruiting focus group members.   Rationale for sample sizes.

*(iv) Field Instructions*:  Details of procedures and participatory tools to be used with each focus group.   See Appendix 6 for an example.

*(v) Materials to be used:*  De-briefing sheets, pens, blank cards, masking tape, flip chart sheets, picture cards.   See A.3.11 for an example for an example of a de-briefing document to be completed after each focus group, by the group facilitator.

*(vi) Who is responsible for* (a) facilitation; (b) note keeping; (c) assembling respondents, etc.

The section on  **DATA** might list background information to be noted down as to the focus group respondents who are to be recruited (e.g. gender, education level, reading ability, land ownership); and a checklist of topics to be covered e.g. assessment of current information channels, evaluation of information sources and media.

If for example, the activity was a laboratory study, there would probably be extensive attention paid to experimental procedures.  See for example the illustration in Appendix 5.

*2.2.2   What the data manager and analyst need at activity level*

Statisticians involved in analysis or data management are sure to need access to metadata i.e. *the detailed description of* information at activity level and below.   One necessary, but not sufficient,  strand of this is an accurate general (activity-level) description of the dataset e.g. (i) that the data are from the 1999-2000 experiment on enhanced plant nutrition as a management option for banana leaf spot diseases conducted at Luwero benchmark site as part of the IPM project, (ii) the list of farmers who participated initially, those who dropped out, why and when etc.   A summary of such information should also be a part of the computerised data file, so that it can be understood without reference to the protocol

*2.2.3   Datum level factual record*

As well as this, however, it is likely that the most onerous demand posed to a research team by a statistical analyst looking at a large dataset from a project will be at the most detailed level – the individual datum level.   The statistician will want full answers to questions such as, 'What does the entry '6' in this cell represent?'   The *generic* questions that need to be answered by a competent basic factual record, at the activity level – and at the datum level - concern:-

<div align="center">What?   Where?   When?   Who?   How?</div>

These questions are answered at varying levels of detail according to circumstances.  In the case of the statistician asking what '6' represents, (s)he might need to find out :-

**What?**        This is the yield of cowpea from the plot the farmer planted using the bag of seed provided in the Food Security Pack

**Where?**      It came from the farm of Mr. S. Moses (link to database file of farmer details)

---

**When?** This was a yield at harvest in April 2002

**Who?** Measurement recorded by farmer (Mr. S. Moses)

**How?** Harvested crop measured in buckets (nominal 50 litres)

The analysis of such data is a link in the chain connecting the study design to the creation of outputs, and clearly the whole process is jeopardised if there are communication failures at this stage. Ensuring that all the relevant information is properly preserved and conveyed at this ultimate level of detail is very important. Thus the protocol for a data collection activity must include a system for careful documentation of what every item of the data represents. See for example A.2.6.

Relevant links in the data files also need to be maintained. For example analysing how yields in 2003 compared to those in 2001, or analysing farmers' changes in planted area between these dates, may require identifying and using data from farmers who were in the sample on both dates, and matching the records as coming from the same farmer.

*2.2.4 Documenting How? – Data collection procedure*

It is very common in data analysis to want to know *exactly* what is represented (at <u>variable</u> <u>level</u>) by one of the item in the dataset, especially when something apparently obvious starts to reveal peculiarities e.g. '90% of your farmers are citing maize yields of under 8 tonnes per hectare, but there are a few impossible values of over 100 tonnes per hectare: why is that?'

It is strongly advocated that fieldwork and data entry instructions are written down in a data management protocol (see section 5.3 below), documented, and used as the basis of quality checks. Often, when questions arise about strange values found in datasets, the explanation turns out to represent:-

(a) simplistic assumptions at the time of data collection instrument design, or

(b) carelessness in transmitting instructions.

To avoid *simplistic assumptions* as in (a) above, it is often invaluable to have additional scrutiny of what seems to its author like a very clear, well-thought-out, and obvious data-gathering instrument or process:

(i) a scientific peer may add intellectual breadth, or may constructively challenge the assumptions, or the usefulness of the results;

(ii) a hard-working and experienced data analyst will often foresee some errors that will arise from a draft data collection schedule;

(iii) someone with supervisory experience of somewhat similar work may foresee what will go wrong at the field or laboratory level when data are being generated.

Inadequate data often results from under-prepared, inadequately-tested instruments. For example, those who design questionnaires, and those who administer them, should be thoroughly conversant with practical features and local conventions, as in the following.

<u>Example</u> of data collection: in a survey of maize yields, the questionnaire designer should be well-aware of (1) the full range of measures used by respondents anywhere in the area covered e.g. yield measured in terms of 90 Kg. or 50 Kg. bags, oxcarts, buckets and so forth, (2) the full range of forms in which product may stored or sold, (3) different fates of the crop - for example in a hungry season some maize may never reach the harvest receptacle to be measured, being plucked and eaten green by farm family or other dependents, or stolen, or sold in small quantities. Enumerators cannot be expected to fill in logical gaps left by more highly-paid survey designers.

Even if tiresome, avoiding *careless transmission of instructions*, as in (ii) above, is important because all other efforts in a study may turn out to be wasted if the stages leading up to the recording of each datum have been done badly. If field recording was inadequate the presumed 'data' may be uninterpretable.

If the field data recording conventions are not fully understood by the data entry personnel, good results from the field can quickly become an incomprehensible muddle when the data are computerised as the following shows.

Example of data entry: 'Helpful' but inexperienced data entry staff were presented with 'unit of measurement' data as 50 [i.e. 50 Kg. bags of maize] and 'yield' as 3 [bags] meaning total yield was 3 x 50 Kg. i.e. 150 Kg. in total. One or more of the clerks misunderstood and entered 50 and then 'helpfully' computed 'yield' of 3 x 50 = 150, but this then entered the analyst's data set as representing a yield of 150 bags, not Kg., an implausibly huge yield from a small seed pack.

---

**Box 2.2.2   Some Activity Protocol Procedures**

An activity protocol needs to include painstaking and detailed procedures for:-

♦   establishing how data collection and computerisation should be carried out,

♦   meticulous training of the staff doing these 'mundane' jobs,

♦   pilot-testing all aspects of the process,

♦   close, intelligent and attentive supervision to ensure there are no ambiguities, misunderstandings or opportunities for slovenly practice.

♦   and, where necessary, retraining or replacing errant staff.

---

The above focuses on key pieces of documentation needed for the statistical analysis with which we are most concerned, particularly copies of data collection instruments, such as recording sheets for experimental studies, questionnaires for surveys, and field manuals telling field staff in detail how to record particular responses. The issue of units of measurements used by farmers (Mr. Moses' 50 litre bucket) is an example of where careful attention to detail is essential.

One other key component that needs documentation by the data analyst in consultation with research scientists is a *data analysis protocol*. This is elaborated in the next sub-section.

### 2.2.5   Documenting How?  Data Analysis Protocols

When planning a particular activity, it is important to look forward to the stage of data analysis with the aim of ensuring that the data to be collected are appropriately structured, and sample sizes sufficient, to enable an analysis that will meet the research objective(s) to which the activity is addressed. Not only will this facilitate and speed up the data analysis when data collection, computerisation and checking is complete, but thought in advance as to how the data will be used can also ensure that every datum is necessary in achieving the research objective(s).

We therefore advocate preparing a data analysis protocol, giving broad expectations of how the data analysis will be done. Some key elements of a data analysis plan are given below but specific details under some of these headings will be very dependent on the nature of the activity concerned. For guidance on data analysis procedures see SSC (2001 a, b, c).

---

*Analysis objectives and variables for analysis*

The data analysis protocol must first relate to clearly defined research questions, identified after a careful consideration of the main research objective(s) set out for the activity.  The next step is to determine the specific variables, needed to address these questions.   See Appendix 7 which illustrates this process for a study involving participatory assessments.

*Software to be used*

The software to be used for the data analysis should be specified, ensuring that it is capable of doing the analysis that is wanted.   For example, initial simple exploration of the data may be done using MS-Excel, but statistical analysis involving hypothesis testing or statistical modelling procedures may use a statistical package such as Genstat, SAS, or SPSS.

*Organising the data for analysis*

Consider and specify the *sampling unit* for each component of the analysis, and the underlying *sampling structure*.  This involves ensuring that the overall sampling plan in the sampling protocol meets the needs of the specific analyses envisaged.   One important element of this is ensuring that the sample sizes will be adequate to generate useful results.   This is particularly relevant to hierarchical data sets (locations, farms within locations, plots within farms) where the sampling units differ according to the level at which they occur in the hierarchy.  Some summary tables of the sampling structure would also be valuable at this stage to guide the analysis at a later stage. See tables in A.3.6 for example.

Document the procedures for calculation of derived variables, e.g. adjustment of maize yields for moisture content, synthesis of an asset index as a measure of wealth status of a household.  Where there are hierarchical data structures, document how variables may be aggregated to a higher level in the hierarchy for any analysis relating to sampling units at the higher level.  For example, aggregating quantitative variables may be in terms of a mean or total, while aggregating a binary variable e.g. yes/no answers, may be in terms of the number of "yes" answers.

Specify the format in which the data should be organised to enable an analysis.  For example, if data is collected on several sampling occasions and data corresponding to each occasion is kept in a separate sheet in an Excel workbook, then for a combined analysis across all sampling occasions, it would be necessary to have all the data in a single worksheet with an additional column to indicate the sampling occasion.

*Exploratory data analysis (EDA)*

Document, and allow time for, the data exploration work that may be done before data analysis aimed at addressing the research questions.  This is important since it gives a good feel for the data and more importantly, it serves to identify errors that may still be present in the data.

For example the documentation may specify the following procedures:

(i) Do simple **descriptive summaries** for all variables selected for analysis.  This will include frequency tables for categorical variables, and summary statistics (number of cases, mean, maximum, minimum, standard deviation) for the quantitative variables.  The maximum and minimum values in particular are useful indicators of possible data errors.

(ii) Do **frequency tables of variables** having only a few likely values (e.g. number of conflicts in the community in the past year).  In cases where very few observations result for particular values, it may suggest re-coding the variable into one with fewer values.

(iii) If the main variable for analysis is quantitative (e.g. days since a forestry co-management scheme came into operation), produce **tables giving counts and mean**

**values** of this variable for each value(or 'level') of factors relating to sampling or treatment structures.

(iv) In all of the above summaries, check that the number of observations (overall and in sub-group categories) are as expected according to the data structure.

(v) **Graphs and charts** are also valuable.  Consider producing:

(a) *Box plots*, to compare groups of data and highlight outliers;

(b) *Scatterplots* between quantitative measurements using separate colours or symbols for each level of factors relating to sampling or treatment structures;

(c) *Bar charts, multiple bar charts, graphs in time order*, can also be considered for some appropriate variables and will be useful for observing trends across selected grouping variables.

*Analysis plan for testing the research hypothesis*

Specify the actual statistical analysis procedure to be used to answer specific research questions.

Example:  Suppose interest is in determining factors contributing to sustainable fisheries and their successful co-management.

▪ The first essential is to specify the main response(s) for analysis, i.e. the outcome variables.   These might be catch per unit area (CPUA) measured in tonnes $km^{-2}$ and catch per unit effort (CPUE – a proxy for resource abundance or biomass) measured in tonnes per fisher per year.

▪ Also essential is to set out the set of explanatory variables, e.g. water body type (seasonal, perennial, both); ecosystem type (rivers, beels, lakes, reefs, others); number of fishing villages; number of fishers of all types.

▪ Then specify the method of analysis to be used, e.g. analysis of variance (ANOVA) using a general linear model.  The ANOVA structure may also be specified.

Example:  For a survey, the analysis plan may list two-way tables that are of *a priori* interest, giving table titles and row and column headings.  Such tables may also indicate data summary and presentation formats that would help to address the objectives of the activity.  For example consider a study involving separate focus group discussions with men, women, boys and girls groups in each of 20 villages.   Suppose one item in the discussion involved determining

which diseases are perceived as being most prevalent in the sampled villages, and whether these views differ across men, women, boys and girls.

For this study component, the following may form elements of an analysis plan:

▪ a one-way table of the frequency of mention of each disease across all 80 focus groups

▪ a two-way table of perceived diseases (rows) tabulated by type of focus group (columns), with table cells containing both frequencies i.e. number of villages, and column percentages.

▪ For those diseases perceived as being present by most groups (identified by results from  above, a bar chart of the percent of groups perceiving each disease according to group type.

Dummy tables would be produced for the former outputs, and the form of a multiple bar chart for the last.

*2.2.6 Other documents forming part of the basic factual record*

Of course there would normally be a range of other documents in the factual record of a substantial project activity. These may be relatively unlikely to be used to inform data analysis, yet remain important as a record of the activity, for example timelines and Gantt charts, activity analyses, work records, or financial files. Thinking further, in good time, about the user community for results, and about their uptake and dissemination, is elaborated in 'Scaling-up and communication: Guidelines for enhancing the developmental impact of natural resources research', DFID-NRSP, 2002. We avoid duplicating this important complementary advice here.

## 2.3 Answering 'Why?' questions at Activity Level

Looking at the other questions above, we could take the rationale for our one activity as 'given'. We could ignore questions as to *why* a particular study was done at all, *why* it was done by using a household survey rather than focus groups, *why* it included respondent Msimba, or *why* he was asked in a particular way about his knowledge of the functions of lime as a soil remediant.

However, analysis and the planning of outputs, will not be purposeful or well-directed unless we have clear answers to 'Why?' questions at several different levels of detail, like the above. Even when accepting the activity rationale, design and general decision-making as being fixed by other people responsible for the larger project, the analyst needs to understand what outputs the project designers need or expect from the activity.

Activities need to tie up in terms of the overall integration of the larger project and the other activities to which it links, and the analysis should be based on an understanding of how the particular activity contributes to project objectives, including how effectively its data will tie up with data from other components.

For example if the results from analysis are supposed to generalise to a wider domain than that covered by the sample under the activity, great detail about a small number of sampled communities may be misplaced.

# 3.  Project-Level Protocols

## 3.1  Introduction

This section 3 is relatively short because this guideline document is largely concerned with project aspects that have some relation to data analysis and outputs, and much of this work is at activity level.   This section focuses on the proper involvement and responsibilities of a data manager and/or analyst at the level of project organisation.

### 3.1.1  Background to the project

Of course the selection of goal and objectives takes place at a higher level to which data analytic activities are subsidiary.  In most cases, the initial project proposal and associated documents will contain much background information that will accompany and inform any additional protocols:-

---

**Box 3.1.1   Potential Background Information**

♦  an analysis of the project setting, e.g. relevant policies, institutions and processes as at the time the proposal was submitted, and e.g. assessment of knowledge, attitudes and practices of intended beneficiaries;

♦  a review of existing information with identification of knowledge gaps;

♦  broad statements about goal and purpose, as well as justification of the project in terms of pre-existing demand or perceptions of need for the advances that the project sets out to deliver;

♦  a brief summary, in the logical framework, that should clarify the overall objectives;

♦  results of any consultations about stakeholders' perceived needs for information, and

♦  indications of its intended uses and required forms.

---

### 3.1.2  Initial planning

Before detailed design of data-collection instruments, the inception phase of a project will involve  looking:-

---

**Box 3.1.2   Views during Planning**

♦  *forward* and undertaking the early-stage development of a workplan leading from clear and more detailed specification of objectives up to activities, analyses and outputs;

♦  *outward*  e.g. to farmers, to dissemination and uptake activities, seeking impact and arranging handover to users of the outputs;

♦  *upward* e.g. to senior civil servants, to (i) find means to ensure policy relevance, and to develop the potential markets for the results, and the creation of interest and expectation as to the project outputs, and (ii) analyse institutional and organisational impediments, and how to tackle these.

♦  *downward* when the time comes to define activities in detail.

---

### 3.1.3   Involving the data manager/analyst

Like others in the engine-room of the project, the data manager and analyst will need guidance through parts of the above, insofar as the relevant sections contribute to the analyst's understanding of the project objectives and priorities in relation to their own work.   The need is clear with regard to (a), but for example, appreciation of (c) above may lead the analyst to focus on presentation-style tables of results, not just on more detailed and turgid outputs.

At the research planning stage, it is obviously very important to focus on the project's most challenging main substantive themes, and ways to approach them.   It is tempting to focus exclusively on these issues, but organisational matters should not be overlooked.   The availability of suitable personnel needs to be checked and the allocation of responsibilities needs to be documented.   Herein this is seen as part of the broadly-defined 'project protocol'.

*Section 5 below* concerns the development of other project-level protocol documents - a *Project Sampling Protocol, Project Data Management Protocol*, and one or more *Analysis Protocols*.   Suitable people need to begin work on these at an early stage, if necessary also acquiring additional training, or seeking some consultant inputs.

It is tempting for a project manager to retain this role personally, because it provides oversight of much that happens along the way, but it should be borne in mind that it can involve a great deal of distractingly time-consuming, and detailed, work.   Development of clear terms of reference for a subordinate may be better in the longer-term.

While similar issues arise for various aspects of a project, our example in the next sub-section concerns duties that may in some cases be written into the terms of reference of a data manager and/or analyst.

### 3.1.4   Day-to-day management of support staff

Ordinary project management tools include work schedules, and Gantt charts which show interdependencies between activities, such as the knock-on effects of delays to predecessor activities.   In projects which extensively use full-time technicians or field staff, shared across several activities, one of the most important uses of such tools may be to avoid double-booking of technician time, or long periods when these staff are unoccupied.

There should be clear lines of responsibility for managing shared support staff.   Where major components of their workloads are for fieldwork, data collection and/or data entry or data management, the project data manager/analyst may reasonably be designated as responsible for their training, supervision, and/or scheduling their time allocation on specific tasks.

## 3.2   Detailed Project Planning

The general planning and conceptual development issues mentioned in the preceding section are not the responsibility of a specialist data manager or analyst.   These people can have more to offer when the project looks *downward* and finalises the list of activities and how they fit together.   Appendix 1 provides an illustration: see A.1.14 and A.1.15.

Often data entry and analysis staff should have a hand in development and validation of study instruments such as de-briefing documents or data recording sheets, in field-worker training so that results are properly recorded, in pilot-testing and assessment of fieldwork(er) quality.   They then need to have meaningful involvement, and a voice, in the project research planning from the very start.

Of course, when activities are listed, each individually needs a working definition of its objectives and outputs, timescale, staff time and budget allocation.   In some form each will have a conceptual framework and research questions, scientific rules, and a set of

---

enabling actions such as training, equipment purchase, and mechanisms to motivate or incentivise respondents.

Activities *inter-relate* in that some need to run simultaneously or to follow in successional order.   They may be competing for, or may be able to share, field staff time or other resources.   Recording these managerial aspects in the project's documentation is acknowledged to be important, but is not developed here, being seen as outside the scope of these guidelines.

In relation to project-level planning and management, the data manager or analyst is likely to be concerned with technical aspects of the *linkage* and the *balance* between activities, and to be particularly keen that these are negotiated, sorted out, written down in the protocol, and accepted by all concerned.   The data manager/analyst is also likely to be concerned with those common plans and resources which cut across activities and relate directly to their own performance i.e. sampling plan, data management system, and analysis strategy.

## 3.3   Linkage of Activities

Despite the desirability of having project work conveniently divided up into separate activities it is important to recognise that ensuring the coherency of project work often involves looking at how these link up.   These may be very coarsely divided into more 'intellectual' and more 'practical' linkages:-

---

**Box 3.3.1   'Intellectual' Linkage of Activities**

♦ in relation to *objectives* – more than one activity may contribute to an overall objective,   and sometimes a portmanteau survey is expected to serve more than one objective.  Either of  these eventualities needs to be recognised and clear ground rules established to avoid gaps, duplication, competition or recriminations;

♦ in relation to *cross-activity analysis* – plot-level data from an on-farm trial may need to be aggregated to farm level for linking to farmer information, or householders' individual  responses aggregated to community level to be linked to focus group findings.  Many past projects have under-performed in terms of exploitation of the value of linking datasets in this way.   See for example A.1.15 for a limited example of doing this.

♦ Activities may also be linked in relation to *inter-disciplinary work* – this may be an objective in itself, or may be planned because it can strengthen two or more activities.   Of course activities rooted in disparate disciplines such as anthropology and economics have different requirements and different traditions and may make quite different demands in terms of numbers of subjects and time spent with each subject.   Notwithstanding that, it can be possible for studies led from two such disparate stand-points to be mutually reinforcing, for example one contributing greater breadth and generalisability, the other greater depth of understanding or better means to communicate results to beneficiaries.   Achieving this means overcoming any barriers to inter-disciplinary partnership, which may require establishment of a documented, formal agreement to do so.   Given that, there is also a technical element in terms of selection of samples – and linking up of data, often across the qualitative-quantitative divide.

---

---

**Box 3.3.2  'Practical' Linkage of Activities**

♦ in relation to *sampling* – for example it makes sense to ensure that the same beneficiary individuals, households or communities are included in both (i) a socio-economic survey or one that explores beneficiaries' knowledge, attitudes and practices, and (ii) a survey of off-take of common property resources or an on-farm experiment.  *See section 5.2 below.*

♦ in relation to *shared identification systems* – common conventions are essential for items that may become key fields linking together different sets of data.  This is also desirable at simpler levels e.g. helping fieldworkers by always coding Yes = 1, No = 0.

♦ in relation to *scheduling* – where the analysis of one extensive set of results has to precede work on another successor activity, delays or deficiencies in the first are prone to result in undue pressure on data entry and analysis personnel, because of pressures to start the later activity on time; where extensive material arrives at uncontrolled times it is difficult to plan data entry workloads, and their supervision.

♦ In relation to *sequencing* – where one action or activity has to take place after, or in the light of, another.   For example an in-depth qualitative study may involve sampling its small number of communities or households in the light of information from an earlier socio-economic survey.   Data analysis work associated with an activity is often delayed by over-running of fieldwork, of data entry or of responses to queries arising about implausible data.

♦ in relation to *data-sharing* – understandings need to be established that forestall any temptation for activity leaders to treat data from 'their' activity as a private resource, and to prevent its integration with other parts of the project's resources, by concealment or delay in its proper organisation.

---

To be able to think constructively about which of these linkages are important, and to build a project protocol that makes the right links and codifies what to do about them, it can be very helpful to have a structured plan of a set of linked activities, as in the following example.

**Box 3.3.3   Specimen Flow chart of Project Activities**



The above is reproduced with permission from the final technical report on a DFID Natural Resources Systems Programme project "Investigation of Livelihood Strategies and Resource Use Patterns in Floodplain Production Systems in Bangladesh" by Julian Barr.

## 3.4  Balance of Activities

After plans have been drafted, the project leadership has to examine critically and in detail, the relevance, quality, quantity and timeliness of outputs that can be expected from each activity team.  One part of the review process, which the data manager/analyst may lead, is to look at proposed data collection plans and to appraise whether their costs and contents are necessary and sufficient for intended outputs, and/or are disproportionate to the importance – to key project targets and research questions – of the likely results.

Clearly some activities are by their nature relatively resource-intensive, but trade-offs between them have to be negotiated in terms of the breadth, depth and importance of their outputs.   There is often a temptation to collect too much data, leaving inadequate time for data management, analysis and reporting.   Negotiated, agreed and documented limits on time and cost are particularly important for activities whose over-running threatens other important components.

## 3.5  Summary list of Project Protocol Headings

We give below a checklist of items that may be important components in a project level protocol.   Some of these are illustrated in the extract in Appendix 1.   It should be noted this list is not claimed or intended to be definitive.

**3.5.1  PROJECT TITLE:**  *<Title of the project as in project proposal>*


**3.5.2  PROJECT LEADER:**  *<Named scientist from the institute responsible for delivering project outputs>*


**3.5.3  PROJECT MANAGER:**  *<Lead scientist for the project in the collaborating institute>*


**3.5.4  RESEARCH PARTNERS:**  *<Name and organisation of each collaborator>*


**3.5.5  PROJECT FUNDING**:  *<Source of funding or name of funder>*


**3.5.6  START AND END DATES:**  *<Project start and end dates>*


**3.5.7  PROJECT PURPOSE:**  *<Purpose as stated in the project proposal>*


**3.5.8  PROJECT JUSTIFICATION:**  *<Brief outline of demand for the research - extracted from project proposal>*


**3.5.9  PROJECT OBJECTIVES:**  *<Specification of what the project expects to achieve>*


**3.5.10  SPECIFIC OBJECTIVES WITH JUSTIFICATION:**  *<Objectives defined in more precise detail, and explained>*

**3.5.11  RESEARCH ACTIVITIES LINKED TO SPECIFIC PROJECT OBJECTIVES:**


**3.5.12  LINKAGES BETWEEN THE RESEARCH ACTIVITIES:** *(Comment: As in section 3.3 above this may well include a diagrammatic representation, and indications of how to handle various relevant forms of linkage, where these are of importance)*


**3.5.13  RESEARCH ACTIVITIES LINKED THROUGH DATA ANALYSIS:**


**3.5.14  PROCEDURE FOR IMPLEMENTING EACH STUDY ACTIVITY:** *(Comment:  Although the activity protocols will give full details of each activity, it is useful in the Project Protocol to briefly indicate when, who, how (in broad terms) and expected date of completion. This could be in the form of a table with rows listing the activities and columns representing the when, who, how and completion date.)*


**3.5.15 SAMPLING PROTOCOL:** *(Comment:  For example, method of selecting farmers for the baseline survey, for on-farm studies and for farmer assessment of cultivars with respect to post-harvest utilisation; addressing sample size issues, and demonstrating how sampling activities link together)*


**3.5.16  DATA MANAGEMENT PROTOCOL:**  *(Comment: Elements needed here include: Identifying persons responsible for data management (if there is no designated data manager), data entry, supervising data collection, entry and validation procedures, having a strategy for data entry and checking, organising the data, archiving, keeping back-up files, maintaining the master copy of the data, ensuring there is an audit trail to track changes made to data files, keeping recording sheets in a safe place, etc.)*


**3.5.17  LIST OF DOCUMENTS RELATING TO THE PROJECT:** *(Comment : This will be updated over the duration of the project.   It may include planning meeting minutes, workshop reports, consultant reports, progress reports, and technical documents.   It should of course be associated with filing system(s) where definitive copies are kept, and should form one resource from which to select non-ephemeral documents for a project archive.)*


**3.5.18  PLANS FOR DISSEMINATION:**  *<These are likely to have been specified in the project proposal, but it is beneficial here to have either a reference to the proposal or a brief outline of what is intended.>*


**3.5.19  LIST OF PUBLICATIONS, CONFERENCE PAPERS, AND OTHER TECHNICAL ARTICLES:** *(Comment : This will be updated over the duration of the project.)*

---

# 4.   Answering 'Why?' questions

Often an otherwise competent basic factual record does not adequately explore 'Why?' questions, and those are our main focus in this section.   The general justification of a project will normally explain broad structure and general rationale, but where the project has been commissioned in response to a 'Call for Proposals' there may be limited inclination to examine in written detail just *why* a specific problem should have been put at the core of the research, or a specific approach favoured.

However, there is some basis for arguing that outputs will be easier to structure, more coherent and have greater impact, if project leadership have successfully reflected, concurred, and prepared documentation on, their interpretations of 'Why?' and have done so at appropriate levels of detail.   It is not enough to have only a plausible short answer at the title/goal level of the project.   Engine-room workers on a project, including data managers and statistical analysts, take many decisions and make many compromises: to do so constructively needs the guidance of well-explained objectives developed in more detail.

## 4.1   Project or Activity?

If an activity is distinguished as such it should have an identifiable rationale, that can stand alone, for why the activity is to be undertaken and why it should take a specific form.   Usually, though, this is closely dependent on the thinking at project level, parts of which need to be thought through in more detail to justify the chosen structure of an activity.

In the following we move between project and activity or do not distinguish them where the argument is more or less similar in both cases.

## 4.2   Background Reasons for Project Decisions

### 4.2.1   Demand

A justification in terms of demand may include identifying knowledge gaps, and a broad judgment as to why these are important and how the gaps might best be filled by the project or activity.   In many instances this will include a demonstration of demand for some output.   The situation analysis describing the background should then provide some explanation of how the demand serves to justify the proposed work.   A well-understood demand should help to frame the outputs that will adequately meet this need: who wants to know, what do they want to know, what will they do with the information, how much detail and accuracy do they need?

Of course demand may arise from various quarters.   It is common for stakeholder consultations to be used to assess demands in very broad form, and to delimit what is feasible within the project's timespan and intellectual and financial scope.

### 4.2.2   Choices and their justification

Often the lines of action taken in a reported project are self-evidently reasonable and interesting, and the results provide information that was not previously available. Yet there is limited evidence that the project leadership ever weighed up other lines of action or attempted to justify their design choices, for instance why a questionnaire survey rather than participatory discussion was used in an early-stage scoping study.   There is a need to lay out clearly the process by which it is determined that a particular approach is not only scientifically or intellectually important and valuable, but also is more important than other potentially reasonable approaches.

Thinking through, and writing down, some justification of decisions is desirable at various levels.   A.1.9 exemplifies in brief form a summary project justification, and A.1.10 some justification of the objectives of that project.

Though presented as referring to fundamental research design – choices of different types of study – the above argument applies at a more detailed level, within an activity, to noting down a justification for decisions about the 'shape and size' of a chosen study. In an experimental activity, a design involving 4 replicates of a 2 x 2 x 2 design maybe should have competed with one involving 2 replicates of a 3 x 3 x 2 design. What trade-offs of information did the planners consider in deciding on the former?

At a detailed level, it is often desirable to have a justification document to accompany an individual research instrument: for instance this may include explanation of (a) why a wealth index (based on households' reported ownership of various goods) is preferred to an assessment of income, consumption, or expenditure, and (b) why a particular set of contributors to the wealth index was selected for inclusion – maybe based on previous reportedly successful and well-validated use by earlier studies.

This exercise of justifying the research design decisions is not primarily undertaken to inform the analysts' work, but it should be undertaken with analysis in mind. For example, if a lengthy questionnaire means a small sample size, there is a trade-off between depth and breadth of the information collected, and those setting up studies should satisfy themselves, before spending much money, that (a) they have considered alternatives of all sorts, (b) their choices are the best available, and that (c) the chosen route should lead to worthwhile results justifying the time and money involved.

Often such decisions are phased and are not finalised in one step, in the case of a field study instrument not until after a process of pilot-testing of drafts. It is worth noting that when a pilot-test leads to any substantial modification, the revised instrument should be (a) pilot-tested again, and (b) checked to ensure it still provides a coherent basis for analysis and outputs.

## 4.3  Getting down to Detailed Planning

### 4.3.1  Consensus building

In terms of clarifying minds there are reasons for going into some detail at the planning stage. A team brought together from different backgrounds to carry out a project or an activity are very likely to share the same broad understandings, but at the start to have

(i)      different pieces of experience and special knowledge to share,

(ii)     different detailed understanding of what is meant by vaguely stated general objectives or abstract terminology,

(iii)    different beliefs about the nature of the outputs that will emerge and are required, of what constitutes adequate or good evidence, and of what urgency there is to reach the activity milestones and end-points.

Even if detailed consideration of these issues is time-consuming, fractious and frustrating, it is important to develop a shared understanding at a detailed level (levels 3 or 4 in section 4.3.2 below), within the activity team and between its members and the rest of the project. If the data manager and analyst are party to these discussions, they are likely to want at least an outline analysis plan (see 2.2.5 above) to be examined and agreed, so that they can assure themselves they can meet the needs of the main specialisms and sectoral interests represented e.g. the anthropologist, the biologist, and the marketing specialist in a project on forest products, e.g. the demographers, the social researchers and the advocacy-oriented NGO, in a project about causes and effects of distress migration.

### 4.3.2  Information at different levels

It may be useful to consider four levels at which discussions need to reach a satisfactory state. These are broadly similar to the *goal*, *purpose*, *output* and *activity* rows of a logical framework.

---

**Box 4.3.1   Levels at which Information is Conceptualised**

1.  *'Intellectual Universe'*  At the most general level, a broad group of interested and informed people can participate somewhat knowledgeably in discussion of the concepts and problems under consideration e.g. stimulation of trade in forest products or e.g. distress migration.   This is the level where stakeholder workshops are likely to make an input to the thinking of the project leadership.

2.  *'Conceptual Framework'*  At the second level, the project team has to develop what is often referred to as a conceptual framework for the activity, an agreed extract from, and systematic synthesis of, what the literature and the informed public have to say about the subject-matter of the activity.   It is usually at this level that cause and effect are sketched and connections set out between larger (e.g. state government) scale and smaller (e.g. impact on community) scale entities.   This framework has to set boundaries on what will be covered in the activity, establish sequences and linkages, and provide the basis for planning and decision-making such as that above.   The conceptual framework for an activity is usually a more detailed version of the relevant part of a project conceptual framework.   A case-study example of an extract from a project conceptual framework is presented in sub-section 4.4 below.

3.  *'Output Level'*  A third level is that where the team's concept set is restricted and made more explicit, so that the research hypotheses and other driving forces can be phrased in an explicit form when defining activity data collection and in writing up outputs.  For example, terms like 'child care' or 'access to information' need to be given more explicit content at this level, so that we can talk in an informative and discriminating way about variant forms and degrees of the refined concept.

4.  *'Determinand Level'*  The fourth level is where we define the contents of the research instrument – the questions in the survey questionnaire, or the analytical procedures and measurements in the laboratory, the field recording regimes of the on-farm trial, or the themes on the focus group topic list.   For instance in a questionnaire, this is the level where we ensure that something very specific can be unambiguously elicited from respondents.

The instrument used at 'determinand level' must provide 'atoms' of data that are combined in analysis into 'molecules' of information about synthesised variables.   This returns us to level 3 above, and these synthesised variables are used in the write-up of outputs.   The 'molecules' will often represent something too complicated to be collected directly e.g. several survey questions may contribute to a measure of 'quality of access to agricultural extension advice'; e.g. the Young Lives project (http://www.younglives.org.uk/) has made use, within a longer questionnaire, of a standard inventory of 20 questions addressed to the caregivers of the index children, the score from which is a measure referred to as '*caregiver's mental health'*.

---

The leadership tasked with setting the general targets for activities in a project are likely to operate at the 'higher' levels in this sequence (1 and 2 above), but general instructions and decisions at the higher levels often fail to translate into effective processes lower down (3 and 4 above).   The designer of the data recording sheet, the data manager and the statistical or other analyst have to deal with data at the lowest level of the following structure, but for outcomes to be satisfactory, the leadership need to be engaged in overseeing what happens at all levels.

What goes wrong?   Objective setting is often recorded in retrievable form only at a very general level, for example at the activity purpose level, 2 above.   This is an appropriate starting level, perhaps, but is insufficient to guide the progress of detailed 'below decks' work, e.g. where the data manager is setting up links between files in a database of study data.   So for example 'improvement in food security' needs amplification to level

---

3.   Are we talking about national (maybe district) food security, or household (maybe individual) food security?   What we mean by 'food security' has to be clear to make the research question meaningful.   Which set of households, and associated sampling issues, have to be clear to make the generalisability of the conclusions meaningful. What kind of conclusions about food security do we expect the activity to yield?   Given clarity at that level, what information can we collect (at level 4) from child carers in households about their food habits, food availability, informal networks, coping strategies or sufferings?

## 4.4   Conceptual Framework – a Case Study

The DFID-funded Young Lives Project (see http://www.younglives.org.uk/) looks at the life trajectories of a large cohort of 'millennium children' up to adulthood around 2015. One part of one version of its conceptual framework is illustrated below.

Here the left hand box is concerned with policies, institutions and processes operating at higher than household level: macro and meso level interventions affecting e.g. districts and communities.   The 'messy meso'-level is omitted from the diagram below.   The right hand box is concerned with selected outcomes for children, say by age 15.   Taking each one of these, the various determinants are discussed and main 'causal linkages' are proposed, initially at a fairly conceptual level.   For example the external determinant 'Mother's Educational Status' (part of 'Family Attributes') may affect the moderator 'Mother's Perception of the Benefits of Education' (part of 'Intra-family Environment') and through that mechanism the 'Child's Educational Outcome'.   Of course the same determinant may affect 'Mother's Earning Potential' and this may affect 'Ability to Pay for Education (or bear the opportunity costs)' and by this additional route affect the same outcome.   The number of influence pathways can become huge very quickly, e.g. 'Mother's Educational Status' may well affect other outcomes such the child's health outcome.    Influence pathways for individual variables are not shown here.

---

**Box 4.4.1   Specimen Summary Conceptual Framework**

| MACRO- POLICY | EXTERNAL DETERMINANTS | MEDIATING DETERMINANTS | OUTCOMES |
|---|---|---|---|
| Government spending<br><br>Education policy<br><br>Child labour policy<br><br>Trade regulations | Child attributes<br>Family attributes<br>Physical environment<br>Socio-economic assets<br>Social capital of the household<br>Economic shocks<br>Community attributes | Child care<br>Education<br>Work<br>Leisure/play<br>Security<br>Health care<br>Intra-family environment<br>Social capital of the child<br>Migration | Nutritional status<br>Physical health<br>Mental health<br>Life skills (*numeracy and literacy*)<br>Development stage for age<br>Child's perceptions of well-being |

---

**Box 4.4.2   Using a Conceptual Framework**

A useful conceptual model will:-

♦ evolve only after much discussion and will then

♦ sort and systematise the concepts, maybe into columns similar to those defined by the boxes in the Young Lives version above.

♦ group related themes together as 'rows' (rather untidy ones, usually) e.g. one row might concern income-earning opportunities and factors affecting these.

♦ simplify the picture by agreeing only to include main themes, and ones which are researchable by the kinds of approaches which the project will use.   For example, after serious heart-searching, the first phase of Young Lives (YL) dropped all consideration of child abuse issues because it could not get high-quality information using the survey tools planned.

♦ help to organise the concept descriptions to be at similar 'levels' e.g. figure 4.4.1 is brief because the concepts expressed all represent quite substantial intellectual areas to define and research

When the model gradually evolves to a reasonably stable form, there should be an emerging consensus amongst the research team about:-

(i) which subset of themes they can tackle e.g. YL omits household income, consumption and expenditure, because these are difficult and time-consuming to do well

(ii) what research methods are feasible e.g. YL uses formal questionnaires, sentinel site sampling, qualitative thematic projects etc

(iii) what exactly the people involved understand by the (often vague) conceptual terms in the conceptual model.

Properly developed and used, the conceptual framework can help to focus thought, trim off marginally relevant issues, and help to ensure everyone involved understands broad concepts in exactly the same way. At a pragmatic level, it helps to organise the drafting of a survey questionnaire and to ensure sets of questions capture the necessary aspects of the higher-level concepts in the framework, and illuminate the linkages between them.

With a clearly-stated conceptual framework, the survey questionnaire or other research instruments can be explained clearly to fieldworkers and other project staff, and a justification of the coverage can be developed.   For example YL can justify measuring 'Mother's Educational Status' when the child is aged 1, maybe 15 years before relating it to Child's Educational Outcome, in order to relate it immediately to the 1-year-old child's quality of care and physical health.   In this way we can check that all the questions included in the survey will in fact be used in the analysis we plan to do.

# 5.  Engine-room Protocols – Sampling & Data Management

## 5.1  Introduction

In 5.2 and 5.3 below, we address two of the pre-occupations of data managers and analysts, namely the need to prepare a 'Sampling Protocol' and a 'Data Management Protocol'.   These are areas where we have often criticised the authors of project reports for scrappy reporting, and apparently incomplete conceptualisation of what is needed.

For this reason, we have maybe erred in sub-section 5.2 in the direction of focusing too much on describing how and why we think sample size issues need to be tackled, rather than just stating what needs to be included in a protocol.   The reader is referred to sub-sections A.3.6 and A.4.6 for limited forms of sampling protocol.   Fuller examples can be seen in Wilson and Huttly (2003) and Bojanic et al. (2003).   Other issues related to sampling structures rather than sample sizes are discussed in SSC (2000b).

## 5.2  Sample Size

One general principle that applies at either the project or the activity level is that in the planning stage, we work backwards from what we need to have achieved, so as to define what can and should be done now, and this has extremely important consequences for sampling and sample sizes, amongst other aspects of research design.   As explained in earlier sections, there needs to be a well-developed sense of why the project is collecting particular data, and how it will advance the research agenda in terms of substantive results.   This *must* be the general basis on which intelligent and constructive sample size and sample structure decisions are made after thoughtful debate and consensus building.

### 5.2.1  Writing Up Processes of Sample Size Determination

It is important to think about sample size where primary data collection involves substantial numbers of 'cases' or 'units' or 'sites'.   If the sample size is too big answers will emerge, but at the cost of wasting a lot of resources on collecting more data than is actually needed for the purpose in hand.    If the sample size is too small and there is no further opportunity to augment it, nothing definite will emerge, and once again data collection results will have been wasted.   Very often, accurate sample size determination is difficult, maybe even impossible, at the early stages of a research process, because there is insufficient information.   This is often a good reason for sequencing, or phasing, the contributory studies in such a way that early work informs future sample size decision-making.

Sample size decisions deserve to be thought through, discussed and written down as part of each relevant activity protocol, because the consequences can be very important in terms of resources expended, time taken, and quality of information resulting.

Usually preferred sample sizes have to pared down to reduce costs, so at project level it is important to consider the competing demands of different activities, to think through the risks involved, and of course to cut down on over-ambitious activities if the resources are truly over-stretched.  It is also desirable to think out at project level the benefits of linking the sampled units between various different project activities, with a view to analyses that draw on data from more than one activity.   This requires clear decision-making at project level to ensure activity teams link up properly and work towards common ends.

Writing up sampling plans in a clear way is not the most popular part of protocol-writing. It involves justifying processes of thought, parts of which are often only intelligent guesses.   Usually it is important to accept that there is no right answer.   There is usually no magic formula, for the sample size in a multipurpose survey.    However, there are good reasons to make the effort.   A clear account of the decisions made and

the procedures subsequently followed is key to the data manager's organisation of the resulting data, the analyst's understanding of what the data mean, and the author's ability to justify his conclusions.   The information is often crucial to follow-on activities, and successor researchers.

### 5.2.2   Statistical formulae or structured use of common sense?

There is an extensive technical statistical literature that provides the basis of deciding on sample numbers, but only in certain rather specific situations.   Most of the more accessible literature helps with sample size determination where the task is to estimate a single well-defined quantitative measure with specified accuracy.   It should be used if possible, but of course some objectives are more diffuse in character.   For instance, in the case of a survey, it might be predicted that there would be a need for certain basic tables so as to look at patterns of information across the responses – this would be decided in the light of understanding users' information needs, the project objectives, and the intended policy and research audiences for results.

All too frequently, there is misuse of general 'statistical' formulae for sample size. Usually 'the formula' is based on the assumption that a single-stage, unstratified random sample of 'people' are asked a single Yes/No question where our only interest is in the proportion of people who answer Yes.   In certain limited contexts this is justifiable.   For instance, a survey sampling strategy was developed by the Expanded Programme of Immunisation (EPI) of WHO.    They looked at samples of children in clusters of households and ascertained whether each child had been immunised, 'Yes' or 'No'. They wanted **only** to estimate the overall immunisation rate for the sampled population, and agreed they would get reasonable accuracy from a standard type of sample – 7 children in each of 30 clusters sampled at random from the population.   This is very different from the example below:-

Example: a survey example – working back from outputs to sample sizes.   In a survey of household heads, one set of desired results is stratified by region, restricted to households whose primary livelihood source is own agricultural production, subdivided by gender of household head, and is cross-tabulated using 'Main Tenure Type' by 'Education level':-

| Male head of household | | | | ← Tenure → | Female head of household | | | |
|---|---|---|---|---|---|---|---|---|
| Own | Share out | Share in | Rent | ↓ Education ↓ | Own | Share out | Share in | Rent |
| | | | | None | | | | |
| | | | | Some Primary | | | | |
| | | | | Complete Primary | | | | |
| | | | | Some Secondary | | | | |

For this set of results to be meaningful, enough households of this livelihood type are needed in each region so that the pattern of results in the tables will be interpretable and useful to readers.   That requires 'sufficiently big' numbers of observations that the patterns of relative frequencies in the 'popular' cells have settled down and are well-established.   Then the data are fit-for-purpose.   This is a complex requirement on a particular segment of the survey, and that it has been satisfied is a qualitative judgment related to the use that may be made of the data.

For livelihoods surveys and others with complicated sets of responses, the formula-based approach is simply wrong, and means the researchers are evading their responsibility to think properly about the situation. For a wide-ranging, holistic, multi-purpose livelihood study, any such approach using a 'magic total number' formula risks ignoring factors such as the following.

---

**Box 5.2.1   Some Features Invalidating Over-simplified Sample Size Determinations**

♦ stratification of the population within that sample size, e.g. we may know there are distinct subpopulations who behave or think differently

♦ exclusions from reported results because some respondents don't qualify under the restrictions imposed on the results

♦ refusals or unavailable respondents and other equivalent forms of loss of anticipated data, e.g. experimental livestock lost to wild predators or stolen, e.g. instrument breakdowns

♦ unpredicted choices by collaborators e.g. in farmer-managed trials

♦ the fact that most questions are not of the Yes/No type e.g. several categories of Tenure or Education

♦ the fact that much analysis does not look at one question at a time e.g. cross-tabulation as in the survey example above

♦ the hierarchical structures within which 'people' live e.g. household, village, district in a rural setting – this important theme is expanded a little in sub-section 5.2.3 below.

---

Before data collection, project or activity leadership should think ahead.

---

**Box 5.2.2   Sample size – Think Ahead to Outputs**

♦ think ahead to the main output results e.g. tables as in the example above

♦ try to decide what the most important output results will be

♦ try out what are thought to be low, medium and high sample sizes

♦ try to predict how the results will look (e.g. how data in tables will be spread over cells)

♦ ask if reasonable and useful results can be expected with that sample size

---

### 5.2.3   Observation hierarchies

This sub-section deals with one specific form of structure that needs to be considered carefully, incorporated in the study design, and documented for careful use in analysis.

The definition of objectives, research questions, and sampling plans has to take account of the hierarchical structures in the population studied and the variability arising at each level. These might be regions, agro-ecological zones, and forest fragments within the zones. In a study of chemical concentrations in leaves of tea plants, one might sample clones, fields planted with the same clone, bushes within fields, then branches as the levels in the hierarchy. In a study focussing on factors affecting the food security of rural people in Ethiopia, we may need to consider district administrations at woreda

level, peasant associations at village level, then households, and perhaps target individuals within them.

In particular it is important to think out the sample sizes needed at each of the levels of the hierarchy. See for example SSC (2000 b), pp. 8 –11.

The present authors have criticised projects in the past which have sampled very few first-stage units (regions, clones or woredas above) and very many ultimate units (forest fragments, branches, or individuals respectively) within each one. A study with only two first-stage units has a sample of size two as the basis for any generalisation to other units at that level. This is almost always very weak, logically and statistically, unless there has been a very careful preliminary study to demonstrate that the small sample is typical and capable of representing other areas!

Crucially, this demonstration has to show the first-stage units are typical in every way that might be important to the conclusions. The woredas should have 'typical' population density, climate, ethnic mix, access to schools and hospitals, governing body composition and approach, access to input and output markets etc. Random sampling of first-stage units cannot guarantee that any small sample meets the criterion of 'typicality': the properties of random sampling are beneficial only when sample sizes are large.

Studies, qualitative or quantitative, that entail substantial data collection from settings where there are hierarchies, need to ensure:-

---

***Box 5.2.3   Sample Size – Accommodating Hierarchies***

♦   Adequate numbers of first-stage units are selected in a way that utilises available information about their differences

♦   Reasonable numbers of second and later stage units are selected

♦   Usually objective, maybe random, is planned for selection of ultimate sampling units, these generally being too numerous to be 'well-known' to the consumers of project outputs

♦   Careful decision-making as to what hypotheses and objectives, what analyses, and what types of conclusion, are expected at each hierarchical level – and documentation of these as part of the protocol.

---

## 5.3   Data Management Protocols

A *Data Management Protocol* is another 'engine-room protocol' that we advocate as being necessary at the project planning stage to focus attention on strategies for data computerisation, checking, organisation, analysis and archiving. Once data has begun to be computerised, the protocol should include data status reports so that progress, and file versions, can be tracked. This is an area where minor-seeming deficiencies of procedure – where some quality, relationship or meaning is lost from the data – can lead to the dissipation of a disproportionate amount of the analytical potential of a project.

SSC has wide-ranging experience of being drawn in to help with various stages of data-cleaning, organisation and analysis, and several negative generalisations or dire warnings are possible from this experience, for example the following.

(i)  Where project staff have put off data entry on the basis that things will be easier later, the difficulties of doing so sensibly generally increase rapidly as the delay gets longer.

(ii)  Where data entry plans have been insufficiently detailed, data analysis has usually been heavily damaged, especially if weakly organised data have been left unused for some time.

---

(iii)  In cases where data-sharing and ownership issues have not been resolved from the start, a claim to 'ownership' or 'privacy' is a very convenient way for a collaborator to cover up the fact that they have not done what they promised, that their team messed up the data collection, or that their organisation of research materials is in a  hopeless mess.

(iv)  Analysis expertise, bought in to help out at a late stage, will be much less productive when the expensive analysts spend a great deal of their time expensively chasing up non-ignorable errors in the data, and writing off chunks of analysis because it is no longer possible to get to an adequately documented or reasonably clean dataset!

(v)  A skilled, professional data manager is likely to be very widely employable, and projects should make every effort to ensure that they are not vulnerable to losing one – by ensuring they have back-up systems and shared knowledge, and by treating data management staff considerately!

In terms of realising the potential of projects that gather substantial data, its effective management is laborious and may seem unexciting, but it is a crucial element. Highlighted below are key elements that should be included in a data management protocol so as to avoid problems like the above.

### 5.3.1   Documenting responsibilities

There is a series of tasks associated with data-related activities.   These include (a) preparing data collection forms, (b) conduct and supervision of data collection, (c) checking data after collection, (d) designing data entry forms with appropriate validation checks, (e) entering the data and checking it once computerised, (f) managing the data, including archiving datasets and associated metadata.

It is important that the allocation of responsibilities for these tasks should be agreed and recorded, along with activity time charts.   Doing this should facilitate the tracking of progress on the activities that need to precede data analysis.

### 5.3.2   Data quality control

Any large-scale data-collection exercise, whatever its source, needs a consistent system to ensure a very high level of accuracy and consistency throughout the 'data chain' from first measurement or elicitation to the use and archiving of the completed set of data. There should be a documented strategy for this process.  Box 5.3.1 below gives an illustration.  The issues involved are discussed in greater detail in SSC (1998) and SSC (2000a).

### 5.3.3   Managing and organising the data

There must be consideration given and agreement amongst project partners as to whether a proper database system is to be used for managing the data or whether a simple spreadsheet package like MS-Excel will suffice.  Where a project has substantial information collection activities in a specific geographic area, it is likely the information will have too complex a structure for it to be sensible to organise the data using spreadsheets.   Some type of database is likely to be preferable, allowing the linkage of data files of different types to be preserved using unique identifiers e.g. the GPS reference of the household.   Decisions about the software, the file design, access permissions and the like, and about training of users, need to be documented and followed up with buy-in from project and activity leaders, with implementation responsibility assigned to someone with appropriate background.

*Box 5.3.1  Some Components of a Strategy for Data Quality Control*

- Software for data entry, validation, management and archiving

- Procedure for preparing data recording sheets (who, when, how) with data entry and field procedures in mind

- Pilot testing process (where, by who, how) and follow-on actions

- Training of field data collectors, including field instructionswritten out in detail in a field manual, to ensure no information is lost or misunderstood at the time of data collection

- Procedure for checking data collection forms for completeness and accuracy when returned from the field and follow-on action to deal with queries

- Setting validation procedures on the computer system used for data entry

- Computer entry of pilot data (checking suitability of formats for data entry) and modification to data entry formats if required.

- Procedures for checking computer entries, e.g. should double data entry* be done, or manual checks.

- Exploratory data analysis by data manager or analyst for logical checks and to ensure graphical and simple data summaries follow expected patterns (see also 2.2.5)

- Queries referred back to data collectors

- Procedures for storing of raw(unedited) data files (where, who, when, etc) and filing of paper copies of the data collection instruments

- Procedures for backing-up files and updating the master copy of the data.


*In double data entry, two completely separate individuals or teams enter the same body of data, then effectively  one set is subtracted from the other item by item.   Every non-zero value, representing different entries for the same field, is then followed up and the correct value ascertained and entered to replace the queries data item.

### 5.3.4   Ensuring documentation of the meta-data

Even a sophisticated database can be useless unless there are accessible and comprehensible explanations of what exactly the records represent!   The emphasis in section 2.2.3 on the need to record details relating to each datum in terms of answering the what, where, when who and how components is intended to point out that when we come to analyse data from some time back, the individuals who collected it may have left, or forgotten the conventions they were then using.   Without detailed records properly organised at the time, the data will have no value.   We usually need the database and other documentation, carefully annotated, kept together, and still interpretable, to make use of the results.

### 5.3.5   Naming and organising files

To those who naturally have well organised minds this section may sound trivial, but our experience is that little thought is given at the project planning stage to a discussion about (a) conventions to be used in naming files and (b) a system for organising data and other files on a computer.  Both of these depend very much on the particular project and the way in which the activities are organised.   Addressing these issues at the start

makes archiving and retrieval of the information easier as the project progresses and after.

Example: to cite possible conventions for naming files, we quote from the chapter on data management in Lawson-McDowall et al (2001), the following convention used for farm-trial activities in the Farming Systems Integrated Pest Management Project in Malawi.

"A standard notation should be used for data file names.  This makes it easy to find data files when they are needed for analysis purposes.  The first two characters of the filename can be used to identify the trial (e.g. ST for *Striga* trial), the second pair to identify the year (97, 98, 99), the next two characters to identify the crop (e.g. MZ=maize, PP=pigeonpeas) and the following character to identify the type of data (H for harvest data, D for damage data), while the last character is reserved for a version number (1, 2, etc.).  For example, the filename ST99PPH2.xls would contain the second version of data on pigeonpea harvest yields from the 1999 *Striga* trial."

As an example of file organisation on the computer, we illustrate in Figure 5.3.1 below the varied conventions adopted by three CPP-funded projects in the National Banana Programme in Uganda, namely an IPM project entitled 'Integrated management of banana diseases in Uganda', a BSV project entitled 'Epidemiology, vector studies, and control of Banana Streak Virus', and a project on 'Integrated management of the banana weevil in Uganda'.   Discussions among research team members within each project indicated different requirements as shown in Figure 5.3.1.  For example, in the IPM project, the main categorisation was by location first, then by type of study, while in the BSV project it was by type of scientific study and then by location.  The convention used must be acceptable by, and intuitive to, all relevant project personnel.

---

**Box 5.3.2 – Specimen Approaches to Structuring Project Files**



---

*5.3.6   Backing-up data files, data archiving and data security*

The project team needs to document an agreed strategy for regularly backing-up files, and establish procedures for data archiving and ensuring data security.   Questions such as "how often should back-ups be taken?" will need to be addressed.   Where data security is an issue, the data custodian will need a list of those who have authority to access the data.   Exactly when data archiving activities begin and how the archive is maintained need to be discussed, agreed and signed up to by those responsible.   A documented strategy is also needed as to who will have copies of the archive once the project has been completed.

# 6.   Last Words

The patient reader who has been through the above text is very likely to have felt that some of the statistical recommendations above pose too heavy a demand on the project that (s)he has in mind.   The projects that SSC staff get involved in helping, or reviewing, are generally those with particularly complex data collection and analysis demands, and there is no doubt that the authors' experience is biased towards these relatively complex cases.

A brief summary follows for those who want something simpler to take away from this guide:-

---

**Box 6.1   Summary**

The processes involved in a data-related activity are sometimes described as the 'data-chain', which indicates how each step in the sequence is linked to those before and after. Each link needs to be sound in terms of being carefully thought-out, quality-controlled and properly documented.   The links need to be properly fastened together, with no losses between them of thinking about objectives, or of meaning.  The larger project is likely to involve a number of inter-linked data chains.

Generally a competent data manager and/or analyst should take a hand from an early stage in the establishment and checking of such data-chains.

---

# Appendices.  Examples of Parts of Protocols

This appendix gives several skeleton examples of activity protocols and one example of an integrative project protocol.  These are drawn from our own experiences of involvement with RNRRS projects but have been deliberately left incomplete for brevity.  They do not attempt to provide full answers to all the requirements we advocate in sections 3 and 4, but serve to demonstrate aspects of what might be expected in activity protocols of different types, i.e. a survey, an experimental study, a lab-based study, and a consultative study with focus groups.  In the main we elaborate on components that fall within the *Materials and Methods* section of an activity protocol.

Further examples of activity-level protocols can be found at SSC (2001d).  Three of these appear within Case Study 7 in our series of Case Studies of Good Statistical Practice, drawn from one Plant Sciences Research Programme project and two Natural Resources Systems Programme (NRSP) projects.  A more comprehensive example[1] is provided in an excellent document, available at

http://quin.unep-wcmc.org/forest/ntfp/outputs.cfm

which describes methodological procedures undertaken by team members[2] of a Forestry Research Programme project entitled "Commercialisation of non-timber forest products: factors influencing success".

This document first puts the project in context by giving background details, recommendation domains for conclusions resulting from the project, the project's research hypotheses, and planned outputs.  It then has a substantial methodology section which describes the mix of qualitative and quantitative approaches employed, full details of sampling procedures used in selection of non-timber forest products, communities, parts of the communities, focus groups (for PRA exercises), and selection of interviewees for household questionnaires inside and outside the community, as well as procedures for data management and data quality.  The data analysis procedures to be undertaken with data collected from each of three principal data sources[3] are also described fully, including statements of the research hypothesis and how they will be investigated.  The final section demonstrates how the results will be integrated across the different studies.

The appendices we provide below cover a range of different studies and topics as follows.

Appendix A1.  Components of an integrative "Project Protocol"

Appendix A2.  An activity protocol for an on-farm experimental study

Appendix A3.  An activity protocol for a participatory study

Appendix A4.  An activity protocol for a survey study

Appendix A5.  An activity protocol for a laboratory study

Appendix A6.  Field instructions for participatory assessments

Appendix A7.  Identifying research hypothesis and variables in an Analysis Plan.

We hope these appendices will give the reader a better understanding of some protocol component contents that would be helpful to the research team, project reviewers and others.  They emphasise elements that are particularly helpful to a data analyst in providing a better service to the research team in achieving its objectives.

---

[1]  Practical Tools for Researching Successful NTFP Commercialization: A Methods Manual (2006) by Elaine Marshall, Jonathan Rushton and Kate Schreckenberg.

[2]  Team members were:  Kate Schreckenberg (ODI), Elaine Marshall (Project Leader, UNEP-WCMC), Adrian Newton (University of Bournemough), Jonathan Rushton (CEVEP, Bolivia), Dirk Willem te Velde (ODI).  Their permission to use their material is gratefully acknowledged.

[3] Data sources:  Community Reports, Market Reports and Questionnaires.

# A.1 Components of an integrative "Project Protocol"

The example below is drawn from a DFID Crop Protection Programme integrated pest management (IPM) project based at Kawanda Agricultural Research Institute in Uganda. The project was led by [4]CABI Bioscience and the School of Agriculture, Policy and Development of the University of Reading, with collaborating partners from the Uganda National Banana Research Programme (UNBRP).

Only an outline is given of some of the protocol elements. The main aim is to demonstrate how each objective was justified and how activities link to the objectives and to each other to achieve the project's goal. This example also illustrate some of the points made in sections 3.2, 3.3, 3.5 and 4.2.2.

---

### A.1.1 PROJECT TITLE: Integrated management of banana diseases in Uganda

### A.1.2 PROJECT LEADER: *<Named scientist from the institute responsible for delivering project outputs>*

### A.1.3 PROJECT MANAGER: *<Lead scientist for the project in the collaborating institute>*

### A.1.4 RESEARCH PARTNERS: *<Name and organisation of each collaborator>*

### A.1.5 PROJECT FUNDING: DFID Crop Protection Programme

### A.1.6 START AND END DATES: January 2000 to June 2003

### A.1.7 PROJECT PURPOSE: Promotion of strategies to reduce the impact of pests in herbaceous crops in Forest Agriculture systems, for the benefit of poor people.

### A.1.8 PROJECT JUSTIFICATION:

(**Comment:** *Only a summary is presented below. Further details are available in the project proposal and in reports of two planning meetings held during the first six months of the project.*)

Banana is the most important single crop for food and income security in Uganda. Yet over the last 44 years there has been a steady but marked decline in production of bananas in Uganda. While the area of land under bananas (c. 1.5 million hectares) is double that of 1956, banana production in traditional producing areas of central and eastern Uganda has severely declined. The decline in these areas has been reflected by a shift in production from central regions in particular, such as Luwero, to western Uganda. But even in these relatively productive regions, there has been a gradual decline with yields currently at only 17 tons/ha/year respectively (compared with 60 tons/ha/year attainable on research stations).

---

[4] Project leaders were Mike Rutherford (CABI Bioscience) and Simon Gowen (Univ. Reading), with collaborating partners W.Tushemereirwe (UNBRP) and Cliff Gold (IITA). Their permission to use this material is gratefully acknowledged.

Baseline research conducted by the UNBRP throughout the banana growing areas has identified and prioritised a number of key constraints to production, including declining soil fertility, a complex of pests and diseases, post harvest problems, socio-economic constraints and low genetic diversity.  This project was aimed at addressing pest and disease problems using an integrated pest management approach.


### A.1.9  PROJECT OBJECTIVES:

To validate cultural management technologies with potential for enhancing banana plant health and productivity and alleviating losses due to banana diseases and pests, and to identify and facilitate uptake of practices which are most effective and acceptable to farmers as part of an IPM approach.


### A.1.10  SPECIFIC OBJECTIVES WITH JUSTIFICATION: *(Comment: Objectives below are restricted to those which lead to activities in Luwero district.  This is to minimise the level of detail given in this example)*

(a) To provide information on banana production prior to the implementation of interventions through a baseline survey.

**WHY?** (i) The information will be used as benchmark against which changes (impacts) brought about, as direct effects of the intervention(s) will be assessed;  (ii) the baseline information will enable intended beneficiaries to express their expectations, thus providing a basis for their assessment and involvement of the planned interventions.

(b) To determine the effect of enhanced plant nutrition, e.g. application of organic materials (like cow dung, compost, manure) at least once a year  (20 – 25 kg manure per mat) as a management strategy for control of banana leaf spots (*Black Sigatoka*), and associated cost benefits.

**WHY?**  Because previous studies have indicated that better nutrition speeded up growth whereas the damage caused by the disease remained the same, and thus better nutrition counteracted the impact of the disease.

The research would be done with cultivars Kisansa, Mpologoma, Mbwazirume, Namaliga, Atwalira.  **WHY?**  Because although they are the popular local high yielding East African Highland banana cultivars, they are highly susceptible to pests and diseases.

(c) To evaluate exotic banana cultivars *FHIA 25, PITA 8, PITA 14, PITA 17 and SABA* for their agronomic performance and their resistance to pests and diseases.

**WHY?**  Because on-station germplasm evaluation activities have shown that these cultivars have resistance to weevils, Black Sigatoka and nematodes, and that they have promising post harvest characteristics.  However they still require evaluation with farmers under farmer conditions.

(d) To promote released exotic banana cultivars, namely *FHIA 01, FHIA 17, FHIA 23 and KM5* among the wider farming population.

**WHY?**  Banana farmers in Luwero are generally unaware of the benefits of the released exotic cultivars.  Procedures are needed to promote these cultivars amongst the wider population of banana farmers in Central Uganda.

(e) To determine how cultivars under (c) and (d) above are rated with respect to their agronomic performance and their resistance to pests and diseases when grown under conditions of poor management (no inputs, only weeding) and good management (using mulch and manure).   **WHY?**  In Luwero, banana productivity is very low and has been declining over several years.  It is important to establish whether any of the exotic cultivars, bred primarily for their disease and pest resistance, would also be able to withstand poor management.  It also serves to inform farmers about the benefits of good management.

---

(f) To determine farmers' acceptance of exotic cultivars under (c) and (d) above.

**WHY?**  It is clear that high agronomic performance and resistance to pests and diseases alone will not convince farmers to grow the exotic cultivars unless they are suitable for their intended use(s), e.g. for food, juice, beer, sale, etc.  Evaluating the cultivars for their post-harvest utilisation is therefore important.

(g) To assess the economic benefits of improved technologies investigated through above activities.

**WHY?**  Unless there is an economic benefit, promoting the technologies has little value.


**A.1.11  RESEARCH ACTIVITIES AND LINKAGE TO SPECIFIC PROJECT OBJECTIVES:** *(**Comment:** Each activity will have its own protocol with a named activity leader.  Here we only list the activities and show how they link to objectives).*

1.  Baseline socio-economic survey.  Its objective is as stated in A.1.10(a).

2.  On-farm trial with East African Highland bananas, using enhanced plant nutrition as a management option for the control of Black Sigatoka.  Its objective is as stated in A.1.10(b).

3.  On-farm trial with exotic banana cultivars *FHIA 25, PITA 8, PITA 14, PITA 17 and SABA* grown under good and poor management.  Information from this trial contributes to objectives in A.1.10(c) and (e).

4.  On-farm trial with exotic banana cultivars *FHIA 01, FHIA 17, FHIA 23 and KMA.* Information from this trial contributes to objectives in A.1.10(d), (e) and (g).  Objective (d) is addressed by dissemination of suckers to the wider population of banana farmers by requesting trial farmers to give out two free suckers for every tissue cultured plant received, and encouraging them to give out further suckers of the exotic cultivars to assist dissemination of suckers to the larger population of banana farmers.

A.  Evaluation by farmers of improved exotic banana cultivars *FHIA 25, PITA 8, PITA 14, PITA 17 and SABA*, and *FHIA 01, FHIA 17, FHIA 23 and KM5,* based on farmers' own post-harvest criteria.  Its objective is as stated in A.1.10(f).

6.  Survey of labour and other inputs used by farmers during trial management.  It contributes to objective in A.1.10(g).


**A.1.12  LINKAGES BETWEEN THE RESEARCH ACTIVITIES:**

*(**Comment:**  Although project proposals often set milestones for major components of the project, there is little effort at demonstrating links between proposed activities.  Figure A.1 shows an attempt at this by setting out the research activities in approximate time sequence.)*

## Figure A.1 Links across research activities and in time sequence

&lt;start&gt;

```
Planning baseline socio-economic survey with stakeholders
                        │
                        ▼
        Implementing the baseline socio-economic survey
                        │
                        ▼
          Planning and establishing on-farm trials
```

| Setting up evaluation trial | Setting up of enhanced plant nutrition trial | Setting up of promotion trial |
| --- | --- | --- |

Sucker distribution by promotion trial farmers

Data collection, computerisation, and validation

Survey of labour and other farmer inputs during trial management

Analysing data from each trial

Assessment of adoption of released cultivars

| Assessment with respect to agronomic performance and resistance to Black Sigatoka | Cultivar evaluation on the basis of agronomic performance and resistance to pests and diseases. | Assessment of the effect of management practices for each cultivar |
| --- | --- | --- |

Planning the activity and evaluation by farmers of improved technologies based on farmers' criteria

Integrative data analysis

Identification of technologies for promotion

&lt;end of data related activities&gt;

---

### A.1.13 RESEARCH ACTIVITIES LINKED THROUGH DATA ANALYSIS:

(**Comment:** *It is common in many research projects to aim for a holistic approach by bringing together different subject specialists, e.g. pathologists, soil scientists, socio-economists. However it is extremely rare to see an integrative analysis done as shown in the penultimate box in Figure A.1. After individual pieces of data analyses from the separate study activities, an overall integrative analysis serves to link different activities together to provide a more complete picture of the research situation being investigated. This is particularly relevant in a project that aims at an integrated pest management approach. We demonstrate below ways in which this may be done.*)

We will suppose for example, that the individual analysis of on-farm trial data showed that one or more of the exotic cultivars performed well under farmer conditions, it would be important to address the question *Are these results consistent across all the farmers included in the trial*? If not, results from the baseline socio-economic survey can be used to investigate, using statistical modelling procedures, whether the farmers' socio-economic characteristics, or other farm-level characteristics (e.g. soil conditions, farmer inputs, etc.) can be used to explain why. This type of analysis would allow appropriate farmer recommendation domains to be set up that will allow different technologies to be promoted to different types of farmers. For example, if the farmer is extremely poor and cannot afford mulch or manure, promoting a cultivar that performs well only under good management would be of little benefit. A blanket recommendation of a set of technologies to all farmers in the target population would not be advisable in such a situation.

Another relevant question which requires links between data from different activities is:

*Is good agronomic performance of one or more cultivars supported by farmers' preferences for adoption of these cultivars?*

There is a need here to link the data from the farmer assessment of improved cultivars according to their post-harvest utilization with the agronomic and pest/disease assessment data.

Questions such as these allows more information to be extracted from the data. Anticipating such questions by appropriate documentation helps the data analyst to plan the data collection and computerisation activities in such a way that allows the data files for the socio-economic survey, the agronomic assessments, the pest and disease assessments and the farmer acceptability results to be combined through well-thought out key identifiers (as advocated in the third and last bullet of Box 3.3.1).

### A.1.14 PROCEDURE FOR IMPLEMENTING EACH STUDY ACTIVITY: (**Comment:** *Although the activity protocols will give full details of each activity, it is useful in the Project Protocol to briefly indicate when, who, how (in broad terms) and expected date of completion. This could be in the form of a table with rows listing the activities and columns representing the when, who, how and expected completion date.*)

### A.1.15 SAMPLING PROTOCOL: *<Method of selecting farmers for the baseline survey, for on-farm studies and for farmer assessment of cultivars with respect to post-harvest utilisation; addressing sample size issues, and demonstrating how sampling activities link together>*

### A.1.16 DATA MANAGEMENT PROTOCOL: (**Comment:** *Elements needed here include: Identifying persons responsible for data management (if there is no designated data manager), data entry, supervising data collection, entry and validation procedures, having a strategy for data entry and checking, organising the data, archiving, keeping back-up files, maintaining the master copy of the data, ensuring there is an audit trail to track changes made to data files, keeping recording sheets in a safe place, etc.*)

**A.1.17 LIST OF DOCUMENTS RELATING TO THE PROJECT:** (*Comment:* *This will be updated over the duration of the project and will include planning meeting minutes, workshop reports, progress reports, short technical documents, etc. These will help in checking items for inclusion in the project archive.*)

**A.1.18 PLANS FOR DISSEMINATION:** (*Comment:* *This is likely to be specified in the project proposal but either a reference to the proposal or a brief outline of what is intended is beneficial here as shown below.*)

⇒ Training and active participation of potential beneficiaries, e.g. scientists, extension services, farmers in project trials.

⇒ Annual workshops to discuss research progress.

⇒ Preparation and dissemination of quarterly, annual and final technical reports.

⇒ Publications in banana-oriented journals and peer-reviewed journals.

⇒ Presentations at national and international meetings

⇒ Interview broadcast in African by Wren Media

⇒ Publications of articles in local newspapers.

**A.1.19 LIST OF PUBLICATIONS, CONFERENCE PAPERS, AND OTHER TECHNICAL ARTICLES:**
<*This list has to be updated as the project progresses*>

*END OF APPENDIX 1*

# A.2 An activity protocol for an on-farm experimental study

This example is drawn from a DFID bilateral farming systems integrated pest management project[5] in Malawi. We emphasise below requirements within the *Materials and Methods* and *Data Management* sections of an activity protocol, with a view to illustrating points made in sections 2.2.1, 2.2.3 and 5.3. Other details are provided briefly to provide context.

---

### A.2.1 ACTIVITY TITLE: Technologies for management of *Striga asiatica* in Blantyre Shire Highlands

### A.2.2 ACTIVITY LEADER(S): <*Name of scientist responsible for managing the research activity, and names of others involved in the activity (e.g. research assistants, technicians, field data collectors, etc>*

### A.2.3 BACKGROUND *(Why activity is being done):* Informal consultations with farmers at the start of the project identified whitegrubs, termites and witchweed (*Striga asiatica)* as major constraints to maize production. *Striga asiatica* is parasitic on the host plant and plants that are attacked often wilt, leading to stunting and failure to produce seed. Traditionally control was by avoiding infested fields and long fallowing, some hand-pulling, and use of crop-rotation. In Blantyre Shire highlands in southern Malawi, crop rotation and fallowing are not possible due to severe land pressure, leading to low soil fertility which favours *Striga* infestation. Use of fertilizer and farmyard manure have been found to have positive effects on the yields of *Striga*-infested maize, while trap-cropping is a viable alternative option to rotation and fallowing. The latter involves the use of a crop which stimulates Striga seed germination but does not support the *Striga* plants since no root attachment occurs. There was thus a need to explore the use of green manure and trap-crops as management options for control of *Striga*.

### A.2.4 OBJECTIVES: To investigate the effects of trap crops and green manures with and without nitrogen fertiliser in a maize-pigeonpea cropping system at Striga-infested sites.

### A.2.5 START AND END DATES (When activity is taking place): Two seasons, 1997/98 and 1998/99.

### A.2.6 MATERIALS AND METHODS (Where, when, how and why):

Location: Farmers' fields in Matapwata Extension Planning Area (EPA) in Blantyre Shire Highlands in southern Malawi.

Important dates associated with the trial: *Start and end dates, planting dates, dates for field training and data collection, etc.*

Selection of farmers: The project team selected two EPAs, namely Chiradzulu and Matapwata, for conduct of project activities on the basis of several criteria including

---

[5] Project's TC Team Leader based in Malawi was Mark Ritchie (then at the Natural Resources Institute, Chatham, U.K.). His permission to use this material is gratefully acknowledged.

serious pest problems identified through reconnaissance surveys topography, rainfall, cropping patterns, population density and recommendations from extension officials. Two villages were then selected from each EPA. Since this activity was a small component of larger trials to investigate technologies for control of whitegrubs and termites in maize-based cropping systems, farmers for this activity were chosen from one EPA, i.e. Matapwata. They were selected purposively to meet the Project's socio-economic objective of targeting poorer smallholders. Farmers selected were those who had fields with high Striga infestation so as to maximise the chance of demonstrating the effects of the Striga management technologies. The trial included 6 farmers.

Selection of fields/plots within the farm:  Four main plots were used in each farm, each split into two sub-plots. Two of the farmers had enough land to include two sets of four main plots, while one farmer was able to accommodate three sets of main plots. With 6 farmers involved in the trial, this led to a total of 10 "blocks" with 4 main plots within each "block".

Treatments:  2 x 3 factorial treatment structure, the fertiliser factor at 2 levels (with or without fertiliser) and the legume treatment factor with three levels (Tephrosia, cowpea or neither), as given below.

Input materials and planting method:

⇒ Fertiliser used was 50kg N (CAN) per hectare, dolloped to both sides of the maize plant at sowing, no top dressing.

⇒ In *Tephrosia* plots, the *Tephrosia* (4 seeds/station) was planted on one side of the ridge at a spacing of 45 am between maize and pigeonpea.

⇒ In cowpea plots, the cowpea (3 seeds/station) was planted between maize and pigeonpea on top of ridge.

⇒ Maize variety was MH18.  Pigeonpea intercrop was the local variety.


Experimental Design: To give farmers the opportunity to compare legume treatments and observe fertiliser effects, a split-plot design was used with each "block" consisting of 4 main plots (10.8 x A.4 m), each divided into two (one with and one without fertiliser) to give a total of eight sub- or split-plots of A.4 (6 maize planting stations) x A.4 m (6 ridges).  Of the 4 main plots, one had *Tephrosia,* one had cowpea.  In 1997/98 2 plots were left as controls with no legume, but in 1998/99, one control plot was planted with *Crotalaria*.  The allocation of "treatments" to the main plots and split-plots is shown in Figure A.2.

Trial management:  Plots were managed by the researcher to ensure that weeding was carried out at the same time on all plots.  This was essential to ensure *Striga* emergence was comparable between plots.  At flowering the heads of *Tephrosia* were removed and thrown into the furrow to prevent nitrogen being concentrated in seed production.


### A.2.7 DATA:

What measurements and when taken:

⇒ Time of first observed emergence of Striga .

⇒ Number of emerged *Striga* stems collected weekly at six sampling occasions.

⇒ Number of *Striga* plants found dead without flowering, collected weekly at six sampling occasions.

⇒ Number of *Striga* plants that flowered weekly at six sampling occasions.

⇒ Fortnightly stand counts of maize, pigeonpeas, *Tephrosia* and cowpea; cause of death of any of the plants.

⇒ Yields of maize, pigeonpea, cowpeas and *Tephrosia* seed and biomass of cowpea and *Tephrosia* from treatment net plots at the appropriate harvest times.

⇒ How measured?

⇒ *Striga* counts on each sampling occasion was measured in each of three quadrats of 0.9 m x 0.9m, formed by enclosing area between four maize stations with quadrats placed between non-contiguous groups of maize plants within net plot.

⇒ Maize yield data collected as total grain weight and usable grain (kg) and adjusted for moisture content.

⇒ *Tephrosia* wet leaf biomass and wet wood biomass(kg)

⇒ Cowpea pod weight and seed weight (kg)

## Figure A.2 Plan of the allocation of treatments to split-plots in fields of 6 farmers

| Farmer | Block No | Plot 1 | | Plot 2 | | Plot 3 | | Plot 4 | |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | $f_1 t_0$ | $f_0 t_0$ | $f_1 t_1$ | $f_0 t_1$ | $f_1 t_2$ | $f_0 t_2$ | $f_1 t_0$ | $f_0 t_0$ |
| 2 | 2 | $f_1 t_0$ | $f_0 t_0$ | $f_1 t_1$ | $f_0 t_1$ | $f_1 t_2$ | $f_0 t_2$ | $f_1 t_0$ | $f_0 t_0$ |
| 3 | 3 | $f_1 t_0$ | $f_0 t_0$ | $f_1 t_1$ | $f_0 t_1$ | $f_1 t_2$ | $f_0 t_2$ | $f_1 t_0$ | $f_0 t_0$ |
| 4A | 4 | $f_1 t_0$ | $f_0 t_0$ | $f_1 t_1$ | $f_0 t_1$ | $f_1 t_2$ | $f_0 t_2$ | $f_1 t_0$ | $f_0 t_0$ |
| 4B | 5 | $f_1 t_0$ | $f_0 t_0$ | $f_1 t_1$ | $f_0 t_1$ | $f_1 t_2$ | $f_0 t_2$ | $f_1 t_0$ | $f_0 t_0$ |
| 5A | 6 | $f_1 t_0$ | $f_0 t_0$ | $f_1 t_1$ | $f_0 t_1$ | $f_1 t_2$ | $f_0 t_2$ | $f_1 t_0$ | $f_0 t_0$ |
| 5B | 7 | $f_1 t_0$ | $f_0 t_0$ | $f_1 t_1$ | $f_0 t_1$ | $f_1 t_2$ | $f_0 t_2$ | $f_1 t_0$ | $f_0 t_0$ |
| 6A | 8 | $f_1 t_0$ | $f_0 t_0$ | $f_1 t_1$ | $f_0 t_1$ | $f_1 t_2$ | $f_0 t_2$ | $f_1 t_0$ | $f_0 t_0$ |
| 6B | 9 | $f_1 t_0$ | $f_0 t_0$ | $f_1 t_1$ | $f_0 t_1$ | $f_1 t_2$ | $f_0 t_2$ | $f_1 t_0$ | $f_0 t_0$ |
| 6C | 10 | $f_1 t_0$ | $f_0 t_0$ | $f_1 t_1$ | $f_0 t_1$ | $f_1 t_2$ | $f_0 t_2$ | $f_1 t_0$ | $f_0 t_0$ |

Notation:
$f_0$ represents no fertilizer,
$f_1$ represents 50Kg N/ha dolloped fertilizer

$t_0$ represents no *Tephrosia* or cowpeas
$t_1$ represents *Tephrosia*
$t_2$ represents cowpea

**A.2.8 DATA MANAGEMENT:** *Below is an outline of some of the key elements taken to ensure an effective data management system. Full details can be found as Case Study No. 6 amongst SSC's "Case Studies of Good Biometric Practices" found at www.reading.ac.uk/ssc/workareas/development/case_studies.html*

Field training and allocation of responsibilities:

⇒ Training course to field assistants in the 1998/99 season.

⇒ Schedule of agreed responsibilities for data entry and quality control.

⇒ Pilot testing of data recording sheets.

⇒ Training in MS-Excel and procedures for backing-up data files to data entry personnel.

⇒ Assigning senior staff member to carry out regular checks on data recording sheets and the entered data.

<u>Organising the data for analysis and archiving:</u>

⇒  Separate sheets within the same MS-Excel workbook used for data collected on different sampling occasions.  This enables each worksheet to be simple and facilitates completing the data entry immediately after data collection.  It also allows consistency checks between sampling occasions.

⇒  A master copy of data maintained and data errors found during data validation and data analysis corrected in this master copy.

⇒  Meta-data included with the data within each data file.

⇒  Standard notation used for naming files.

⇒  Back-up files kept up-to-date throughout the activity.

⇒  A system set up for archiving the data.


**A.2.9  DATA FILE NAMES:**  *<List of all data files, data recording format files (e.g. questionnaires), program files, reporting documents, planning documents.>*


**A.2.10  DATA ANALYSIS PLAN:**  *<Identification of the specific objectives of the analysis, listing variables to be used, noting steps needed to organise the data into the right format for analysis and an indication of the type of approach to be undertaken during the data analysis and the software to be used>*


*END OF APPENDIX 2*

# A.3  An activity protocol for a participatory study

This example is drawn from a project funded by DFID's Natural Resources Systems Programme (NRSP) and led by Rothamsted Research Agriculture and Environment Division, collaborating with the International and Rural Development Department (IRDD) at the University of Reading[6] and a Communications Consultant.  The extract is primarily aimed at demonstrating sampling needs of a protocol and to show the type of format (see A.3.11) that can be used to capture much of the information gathered during a participatory study.  This example also illustrates points made in the main text in sections 2.2.1, 2.2.5 and 5.1.

---

### A.3.1  ACTIVITY TITLE:  Knowledge and Information Systems in the North-West regions of Bangladesh

### A.2.2  ACTIVITY LEADER(S):  <*Name of scientist responsible for managing the research activity, and names of others involved in the activity (e.g. research assistants, technicians, field data collectors, etc)*>

### A.3.3  BACKGROUND (Why activity is being done):  The project *Strengthened Rural Services for Improved Livelihoods in Bangladesh* was set up following a need identified in a previous NRSP project, namely *Feasibility of Integrated Crop Management in Bangladesh*, for a decision-support system capable of strengthening farmers' access to information on ICM-related technologies.  A primary aim of the project was to develop and promote efficient systems for the provision of rural services to the poor.  One way to fulfil this aim in part is to enlist a local NGO to identify some specific approaches to improving their information provision and explore the impact of this intervention.  A necessary pre-cursor to this process is the gathering of information on the current status of knowledge and information flows in areas targeted for project activities.  This protocol describes one activity carried out in the north-west region of Bangladesh to collect information on the current situation.

### A.3.4  OBJECTIVES:  To establish an understanding of instruments and mechanisms by which people obtain information from available sources, the perceptions that different client groups have of the quality of information itself and the reasons for choosing different information sources.

### A.3.5  START AND END DATES:  April 2003 to September 2003.

### A.3.6  MATERIALS AND METHODS (Where, when, how and why):

Location:  North-west region of Bangladesh since this is one of the three areas where the DFID/IRRI PETRRA (Poverty Elimination Through Rice-Research Assistance) project is operating.  However Char lands in the far east of the area are excluded where poor peoples' livelihoods are not mainly dependent on crop production and where no PETRRA sub-projects are situated.

---

---

<u>Important dates associated with the activity</u>:  Testing of field methodology in April and November 2002; field work in January 2003; data computerisation, checking and analysis, and reporting: February – September 2003.

<u>The field methodology:</u>

RDRS's mode of facilitating people's development was through the formation of people's "organisations", first as "primary groups", graduating in time to "secondary groups" when they demonstrate their capability for managing group activities by themselves. Participatory discussions with such self-help groups seemed the appropriate field methodology to use rather than a household survey.  Short field visits therefore took place in April 2002 to three working areas of RDRS in the north-west to explore the feasibility of some common PRA tools for gathering information about information flows. A draft methodology for the field work, based on conducting several participatory exercises, was developed and pre-tested in November 2002.  (***Comment:***  *The actual process was more detailed, involving several iterations of discussions of the methodology by the project field team and the external collaborators, consultations with the local NGO partners, documenting the process, further field testing, discussions and modifications to the procedures, and re-testing in the field before finalising the exact checklist of questions and how exactly this information was to be elicited during participatory assessments.*)

<u>Sampling Method:</u>  (***Comment:***  *Below is an outline of some of the key elements in the sampling process.  Further details and a full justification for the procedures undertaken can be found in the activity sampling protocol found at* [*www.reading.ac.uk/ssc/workareas/development/case_studies.html*](http://www.reading.ac.uk/ssc/workareas/development/case_studies.html) *as Case Study No. 9 in SSC's series of Case Studies of Good Statistical Practice.*)

*(a) Identifying the target population*:  In anticipation of Phase 2 of project activities, which would involve RDRS trying out one or more new interventions to improve their information provision, the survey sample was restricted to RDRS's target population, i.e. those with no more than 1.5 acres of cultivable land.  The project excluded the landless since its focus was on uses of information to make decisions about resource allocation, innovation and investment in crop production, and in particular integrated crop management.

*(b) Specifying the sampling unit:*  A focus group of 12-16 persons comprised the sampling unit.  Information flows were expected to be different for RDRS formed groups compared to "control" groups of persons from a community having no involvement with RDRS.  To maximize the chance of observing an impact with interventions put in place by RDRS in phase 2, it was decided to sample two recently formed (less than 1 year old) RDRS groups for every sampled control group.

*(c) Sample size:*  It was decided to conduct a total of 30 focus group discussions, with a reasonable gender balance, on the grounds that this was feasible within time and resource constraints.  It was regarded as being sufficient to enable a meaningful analysis to compare results from male groups with those from female groups, and between RDRS groups and control groups.  During the field work however, the team also felt a practical need to carry out two more participatory exercises, one with an Adibashi[7] female group and the other with a poor Hindu community, bringing the total number of discussion groups to 32.

*(d) Method of sample selection:*

⇒ Five upazilas (sub-districts) were chosen purposively from the 29 upazilas in the north-west of Bangladesh.  The purposive selection of upazilas (i) ensured crop-based agriculture was a primary livelihood, (ii) provided a reasonable coverage of the north-west and (iii) represented the two major agro-ecological zones in the area, (iv)

---

[7] Indigenous disadvantaged people suppressed by neighbouring communities.

considered level of remoteness and distance from RDRS headquarters, and (v) considered whether the upazila included recently formed RDRS groups.

⇒ Within the chosen upazilas, unions (next recognised administrative unit) were selected as far as possibly at random whilst ensuring they were non-neighbours and included newly-formed RDRS groups.

⇒ The selection of a group (RDRS or control) within a union was made after discussion with RDRS staff and visiting the village.

Tables 1 and 2 show the distribution of the number of groups across the upazilas covered during field work and the total number of people participating in the discussions.

*Table A.3.1  Distribution of groups across zones, upazilas and gender*

| Zone | Upazila | Number of participating groups | | | | |
| | | RDRS Male | RDRS Female | Control Male | Control Female | Total |
|---|---|---|---|---|---|---|
| Tista-Korotoa Flood Plain | Aditmari | 3 | 1 | 1 | 1 | 6 |
| | Kaliganj | 2 | 2 | 1 | 1 | 6 |
| Old-Himalayan Piedmont zone | Panchagarh | | 4 | 1 | 1 | 6 |
| | Pirganj | 1 | 3 | 1 | 2 | 7 |
| | Thakurgaon | | 4 | 3 | | 7 |
| | Totals | 6 | 14 | 7 | 5 | 32 |

*Table A.3.2  Participating numbers across zones, upazilas and gender*

| Zone | Upazila | Number of people participating in discussions | | | | |
| | | RDRS Male | RDRS Female | Control Male | Control Female | Total |
|---|---|---|---|---|---|---|
| Tista-Korotoa Flood Plain | Aditmari | 58 | 17 | 15 | 15 | 105 |
| | Kaliganj | 39 | 42 | 19 | 21 | 121 |
| Old-Himalayan Piedmont zone | Panchagarh | | 68 | 15 | 15 | 98 |
| | Pirganj | 17 | 46 | 15 | 28 | 106 |
| | Thakurgaon | | 71 | 45 | | 116 |
| | Totals | 114 | 244 | 109 | 79 | 546 |

*(e) Preparing the field manual:*  Discussions following the piloting of field procedures led to the development of a field manual to enable a systematic approach to collecting the necessary information, while preserving the flexibility associated with conducting PRA exercises.  Elements in this field manual included the following components (a copy of the field manual would normally be a part of the activity protocol):

⇒ Schedule of visits to the field to cover 32 focus group discussions (FGDs)

⇒ Introductory remarks to be made on visit to the field

⇒ Background information to be collected from individual respondents present at the FGD

⇒ Field procedure and resources needed for assessing current knowledge, gaps, information sources & preference for different media to receive information, including details of participatory tools to be utilised.


**A.3.7 DATA:** *(**Comment:** If information from a participatory exercise are to be analysed in a meaningful way to produce generalisable conclusions, then field facilitators need to record the information from the field exercise in a systematic way, directly after completion of each focus group discussion (FGD), using a "De-Briefing Document".  It aims to capture the key pieces of information in a systematic way across all FGDs conducted.  An example is provided at the end of this section.  Further information concerning the focus groups discussions may still reside in the field notebooks which should be retained to add depth to reports of activity findings.)*


**A.3.8 DATA MANAGEMENT:** *<Description of how data will be computerised, organised and managed and plans for data analysis procedures, together with lists of data file names and other documentation>*


**A.3.9 DATA FILE NAMES:** *<List of all data files, data recording format files (e.g. questionnaires), program files, reporting documents, planning documents>*


**A.3.10 DATA ANALYSIS PLAN:** *<Identification of the specific objectives of the analysis, listing variables to be used, noting steps needed to organise the data into the right format for analysis and an indication of the type of approach to be undertaken during the data analysis and the software to be used>*

### A.3.11 EXAMPLE OF A DE-BRIEFING DOCUMENT:

**Strengthened Rural Services for Improved Livelihoods in Bangladesh**

**(NRSP Project R8083)**

**KIS PARTICIPATORY EVALUATION IN THE NORTH-WEST OF BANGLADESH**

**DE-BRIEFING DOCUMENT**

Upazila:_____     Identification Number[8]     | | | | |

Name of Group: _____

Village:_____     Union: _____     Date: _____

Facilitator: _____     Co-Facilitator: _____

---

## 1. BACKGROUND INFORMATION CONCERNING THE FOCUS GROUP

(Circle appropriate answer where relevant)

Group Gender:          1=Male;          2=Female

Group Type:          1=RDRS;          2=Control

If RDRS,          (a) Date of joining: _____     (mm/dd/yy)

(b) Which RDRS Extension Officer is in charge of Group: Crops (✓)

| EO Crops | EO Livestock | EO Fisheries | EO Soc.Dev.&Ed. | Other |
|----------|--------------|--------------|-----------------|-------|
|          |              |              |                 |       |

Has anyone in group been involved with external research/development activities?   1=Yes;     2 = No.

If YES, give details (e.g. PETRRA, FLE, etc):

_____

_____

---

[8] For field work in the North West, the ID number should start with the letters NW, followed by sequential numbers 01, 02, etc.

Information concerning each member of the group.      For yes/no answers, yes = ✓   and no = ×

| Name | Education level (CLASS) | Reading ability? (✓/×) | Own cultivable land? (✓/×) | Rent land? (✓/×) | Income source (or occupation) | |
|---|---|---|---|---|---|---|
| | | | | | Main source | Secondary source[9] |
| 1. | | | | | | |
| 2. | | | | | | |
| 3. | | | | | | |
| 4. | | | | | | |
| etc | | | | | | |
| etc | | | | | | |
| etc | | | | | | |
| etc | | | | | | |
| 14. | | | | | | |
| 15. | | | | | | |
| 16. | | | | | | |

## 2. *INFORMATION NEEDS*

Sketch diagram of main topic (crop-focused) and all the information needs (themes) corresponding to this, as identified by the group.

Show the allocation of 100 seeds into ALL themes.  Select for further discussion only 5 of these at most, i.e. the most important ones.

*In the actual document, nearly a page of space was left in this box for capturing the information*

Key points raised during the discussion (e.g. reasons for allocating a higher number of seeds to one particular theme)

*About ¼ of a page left in actual document to record notes.*

---

[9]  This last column may be left blank if there is only one occupation

## 3. INFORMATION GAPS AND SOURCES

| Themes | Most important types of information needed within theme (i.e. STRANDS) | Have you ever received this information? | No. of votes | If YES, from whom, or from where did you get this information? If NO, from whom, or from where would you *expect* to get this information? If unknown, then write "don't know". |
|---|---|---|---|---|
| 1. | 1. | Yes | | |
| | | No | | |
| | 2. | Yes | | |
| | | No | | |
| | 3. | Yes | | |
| | | No | | |
| | 4. | Yes | | |
| | | No | | |
| | 5. | Yes | | |
| | | No | | |

Additional key points resulting from the discussions.  Make special note if any use their own knowledge about (say) their own methods of experimentation.


## 4.  INFORMATION CHANNELS

Ask the participants to explain how they are getting information from different sources, i.e. means / media.  Ask them to show on the ground by drawing a flow chart in two steps.   Step one is how they receive the information from different sources including means / media.   Step two is how they are disseminating that information to the next lower level or neighbouring farmers.   Use colour cards to draw the flow chart and make it visible.  Sketch the flow chart below.

*In actual document space of about 2/3 of a page was left in this box*

**NOTES:** *1/3 of a page was left for notes in the actual document.*

## 5. *EVALUATION OF INFORMATION CHANNELS (MEDIA)*

Construct a matrix with all identified information channels from the flow chart.  Carry out pairwise ranking and record below.  (Add rows / columns as necessary)

| Information media/means | 1. | 2. | 3. | 4. | 5. | 6. |
|---|---|---|---|---|---|---|
| 1. | x | | | | | |
| 2. | | x | | | | |
| 3. | | | x | | | |
| 4. | | | | x | | |
| 5. | | | | | x | |
| 6. | | | | | | x |
| TOTALS | | | | | | |

**NOTES:** *In the actual document, remainder of the page was left blank here.*

## 6. *EVALUATION OF SOURCES*

Use pocket chart to construct a matrix of STRANDS and SOURCES (as many as have been identified) and facilitate participants to score individually using seeds (max 5 seeds per cell).

Add rows/columns as necessary.

| STRANDS | SOURCES | | | | |
|---|---|---|---|---|---|
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |
| | | | | | |

| TOTALS: | | | | | |
|---|---|---|---|---|---|

Any further comments about the group or key points emerging from the group discussion.

*In the actual document, half a page was left for these notes.*

### *END OF APPENDIX 3*

# A.4  An activity protocol for a survey study

The example below comes from a DFID Forestry Research Programme funded project on "The promotion and development of African rattans: an ecological and socio-economic approach"[10], with a research focus in Cameroon, Nigeria and Ghana.  In contrast to typical RNRRS projects where project activities are listed very briefly in the project proposal, here the activities where primary data collection is expected are described in considerable detail, addressing in particular, sampling needs of three major activities, namely participatory rural appraisals, intensive on-farm studies and extensive socio-economic surveys.  We focus below on the last of these, emphasising the project's work on elaborating more fully on the sampling aspects of their survey work.  This provides an illustration of the sampling protocol requirements described in section 5.2.

---

### A.4.1  ACTIVITY TITLE:  Study of present patterns of rattan usage and sales, and their implications for livelihood

### A.4.2  ACTIVITY LEADER(S):  < *Name of scientist responsible for managing the research activity, and names of others involved in the activity (e.g. research assistants, technicians, field data collectors, etc>*

### A.4.3  BACKGROUND (Why activity is being done):  Activity contributes to the first project output, i.e. to evaluate the socio-economic acceptability of different approaches to sustainable rattan cultivation and enrichment planting for different categories of low-income farmers.

### A.4.4  OBJECTIVES:

The aim was to gain a more comprehensive and socially differentiated view of the significance of rattan for rural livelihoods within the regional economies of the socio-economic and ecological zones where the intensive studies are being undertaken.

### A.4.5  START AND END DATES:  Cameroon: August 2000 to May 2003; Nigeria: September 2000 to May 2003; Ghana: February 2001 to May 2003.

It is important to note that the survey sites were visited periodically (every 3-5 months) during this time.

### A.4.6  MATERIALS AND METHODS (Where, when, how and why):

Location:  This activity (and others within the project) are restricted to the major rattan supply areas which lie within the humid forest areas in each target country.  Within these areas, the following administrative regions were selected:

⇒  Southwest Province, Cameroon

⇒  Cross River State, Nigeria

⇒  Western Region, Ghana

---

[10] Project Leader:  Terry Sunderland, University College London.  Activity Leaders:  Ruth Malleson, Socio-Economic Consultant and Phil Burnham, Socio-Economic Adviser.  Country Partners: Forest Research Institute of Ghana; Limbe Botanic Garden, Cameroon; Living Earth Foundation and Ekuri Initiative, Nigeria.

The above administrative regions were selected for the following reasons:

⇒ All regions include areas where rattans are particularly abundant. As the project memorandum makes clear, "the areas where rattans flourish are rather patchy, even in forested areas, and do not correspond to administrative or other census units that could be used to define the population universe to be sampled". Rattan also flourishes in other southern provinces of Cameroon, however in Ghana and Nigeria forest degradation is much more widespread. Cross River State, Nigeria and Western Region, Ghana still contain forested area where rattans are still present in considerable densities, whilst in other regions rattan is far less common.

⇒ All the selected regions contain people who are involved in rattan-related activities and/or use items made with rattan cane in every day life.

⇒ All can be divided into areas that may be referred to as "zones" with contrasting access to markets and forest resources.

⇒ Collaborating institutions in Cameroon and Nigeria are based in these regions and carry out research and development activities here, so logistically it was convenient to limit the research to these areas. In the case of Ghana the collaborating institution is not based in the Western Region but carries out rattan and other non-timber forest product-related research work in this region.

Target population: This consists of low-income farmers and small-scale urban craft producers in the chosen regions. Information on the social characteristics and wealth status of sampled households allowed the target population to be defined more precisely.

Sampling methodology: Several further stages of sampling took place within each administrative region. Each stage is described below to answer the "how" and "why" questions and to indicate the sample size.

⇒ *Selecting 1st stage units:* Three different types of zone were identified in each administrative region, i.e. cross-border zones, remote zones and on-road zones. These were chosen purposively to cover differences in accessibility to local and cross-border markets, forest resources, roads and international borders. The "zones" were not recognised administrative units, but the settlements within each zone have similar socio-economic characteristics. Descriptions of the characteristics of the different zones appear in country reports.

⇒ *Selecting 2nd stage units:* Details of the settlements (≅ villages) chosen in each target country are available in country reports. Settlements were purposively selected for the socio-economic surveys within each of the above zones, using the following criteria:

   ▪ Proximity to areas where rattan resources are abundant . Because areas where rattans grow are patchy even within the selected zones, only a limited number of settlements within each zone have access to the rattan resource.

   ▪ Activities of collaborating institutions in each of the target countries. Settlements chosen are located within the area where collaborating institutions operate for logistical reasons and for convenience. Some of the communities selected for survey work are involved in forest conservation and development activities with collaborating local institutions or other conservation NGOs.

   ▪ Socio-economic characteristics of rural settlements within these zones tend to be fairly similar. So although the non-random selection of study settlements may introduce some bias into sampling, the settlements chosen are fairly typical of other settlements within the same zone.

   ▪ The number of settlements per zone varied between 1 and 5, depending on the size of the settlements sampled.

⇒ *Selecting 3ʳᵈ stage units:* Using key informants within each settlement, a mapping exercise was carried out to identify each household in the settlement. A wealth ranking exercise, together with the preliminary analysis of household census data, then allowed the households to be divided into 2 groups: those relatively wealthy and those relatively poor. This information was used to select a sub-sample of households to which the multi-round and rattan consumption questionnaires were administered. The wealth ranking also served to yield descriptive information on livelihoods pursued by households belonging to different wealth groups. The number of households selected for inclusion in the sample was based on the requirement that 120 households should be selected from each zone, spread over the chosen settlements in that zone. This was undertaken using systematic random sampling. The exact procedure is detailed in the project's fieldwork and survey manual.

Important dates associated with the activity:  *<Dates associated with testing of field methodology, field work; data computerisation, checking and analysis, and reporting – NOT REPRODUCED HERE>.*

Field operations:  Field procedures to be undertaken are fully described in a Fieldwork and Survey Manual, and form a part of this protocol.

The survey questionnaires:  *<The activity protocol should include (say as appendices) all questionnaires used in the survey  – NOT REPRODUCED HERE >*

**A.4.7 DATA:** Four survey questionnaires formed the data collection instruments – a household census questionnaire, a multi-round income questionnaire, a short rattan questionnaire and a long rattan questionnaire. In the project's fieldwork and survey manual, the purpose of each questionnaire is explicitly stated, and each question in the questionnaires justified with appropriate instructions to the field staff about how to elicit the information required.

**A.4.8 DATA MANAGEMENT:** *<Description of how data will be computerised, organised and managed and plans for data analysis procedures, together with lists of data file names and other documentation – NOT REPRODUCED HERE >*

**A.4.9 DATA FILE NAMES:** *<List of all data files, data recording format files (e.g. questionnaires), program files, reporting documents, planning documents – NOT REPRODUCED HERE >*

**A.4.10 DATA ANALYSIS PLAN:** *<Identification of the specific objectives of the analysis, listing variables to be used, noting steps needed to organise the data into the right format for analysis and an indication of the type of approach to be undertaken during the data analysis and the software to be used – NOT REPRODUCED HERE >*

*END OF APPENDIX 4*

# A.5 An activity protocol for a laboratory study

This example is drawn from a DFID Crop-Protection Programme research project on integrated management of the banana weevil, based at Kawanda Agricultural Research Institute (KARI) in Uganda. The project was led by [11]the School of Agriculture, Policy and Development of the University of Reading, with collaborating partners from the Uganda National Banana Research Programme (UNBRP). This protocol is aimed at demonstrating documentation required in relation to the preparation of experimental material in a lab-based study, and procedures relating to the experimental design. Although the first of these is not directly needed by the data analyst, this information is nevertheless an important part of the protocol because it provides the details required for the study to be repeated at a future time. This is particular necessary in studies such as the one described below where the prepared material is key to likely study recommendations.

---

### A.5.1 ACTIVITY TITLE: Effect of soil amendments in the delivery of Beauveria Bassiana for the control of the banana weevil

### A.5.2 ACTIVITY LEADER(S): *< Name of scientist responsible for managing the research activity, and names of others involved in the activity (e.g. research assistants, technicians, field data collectors, etc>*

### A.5.3 BACKGROUND (Why activity is being done):

Recent studies indicate that the entomopathogenic fungus, *Beauveria bassiana* has a high potential as a biological control agent for the banana weevil in Africa. However, the biotic and abiotic factors that may influence the efficacy and persistence of this fungus under field conditions are not yet fully evaluated. Therefore, this activity aims at evaluating the efficacy and persistence of various *B. bassiana* formulations under laboratory conditions and evaluating the most effective formulation for use in subsequent field experiments.

### A.5.4 OBJECTIVES:

To study the infectivity of different *B. bassiana* formulations to the banana weevil (*Cosmopolites sordidus*). More specifically to determine the amount of conidia produced from different B. bassiana substrates and to evaluate the infectivity of different B. bassiana formulations against the banana weevil under laboratory conditions.

### A.5.5 START AND END DATES: *<Dates refer to the start and end dates of the activity as a whole>*

### A.5.6 MATERIALS AND METHODS (Where, when, how and why):

Location: Banana nematology/weevil laboratory, Kawanda

---

[11] Project leaders were Simon Gowen (Univ. Reading), and Caroline Nankinga Kukiriza (Banana Research Programme). Activity Leader: Evarist Magara, author of this protocol. Permission from these persons to use this material is gratefully acknowledged.

---

Source of materials:

| Material | Source |
|---|---|
| Cracked maize and maize bran | Kawempe maize mill |
| Bagasse | Lugazi sugar works |
| "Machicha" | Kawanda malwa (local brew) joint |
| Cotton husks | Kawempe ginnery |
| *B. bassiana* inoculum | lab. reserved conidia, continuously recultured, and kept in the fridge at $4^0$C. |
| Banana weevils | parent stock collected from Masaka District, then reared in metallic drums in a shade outside the laboratory |
| Spent yeast | Uganda Breweries |
| Sucrose (sugar) | Purchased from retail shops |
| Clay and loam soils | KARI swamp and field respectively |

Preparation of experimental materials and data collection

*(a) B. bassiana spore (conidia) counts*

METHOD (Ref: LUBILOSA)

⇒ 1g of fungal substrate weighed into a test tube

⇒ Mix with 100 ml of distilled water, then add 2 drops of liquid soap

⇒ Let the solution settle for about 10 minutes

⇒ Shake and mix thoroughly

⇒ Measure out 1 ml and mix it with 9 mls of distilled water (= $10^{-1}$ dilution)

⇒ Using a dropper to introduce one drop into the counting chamber

⇒ Count the spores in the 5 diagonal big squares, in the 2 grids

⇒ Finally use the formula $C = A \times 5 \times 10^4$, where C is the concentration of spores/ml in the diluted

⇒ quantity and A is the average spore counts from the 2 grids

⇒ The concentration of spores in the original solution before dilution:

⇒ $S = C \times 10^n$ where n is the number of dilutions, i.e. $S = A \times 5 \times 10^4 \times 10^n$.

*(b) Banana weevil rearing*

The initial batch of banana weevils were trapped from KARI and farmers' banana plantations in Masaka using split pseudo stem traps.  The weevils were reared in metallic drums on fresh banana corms under a shade as described by Nankinga (1999).  The adult weevils were introduced to pared banana corms to oviposit eggs for seven days and thereafter the banana corms were maintained in metallic drums for 60 days to allow development of eggs to adults.  The drums were covered with papyrus mats to avoid desiccation.

*(c) B. bassiana culturing and formulation*

One strain of the fungus, code G41, known to have high pathenogenicity to *C. sordidus*, superior growth and sporulation was used.  It was cultured in KARI insect pathology laboratory on the substrates under evaluation; cracked maize, maize bran, ''machicha'', cotton husks, bagasse, cotton husks + maize bran, maize bran +bagasse and bagasse + spent yeast.  The substrates were cultured following the modified diphasic method described by Nankinga (1999).  Where substrate mixtures were made, this was done to the ratio of 1:1 by volume.

''Machicha'' is spent millet and yeast residue obtained after a local potent gin (''malwa'') has been extracted.  This was collected from the local drinking places, washed, dried and used for culturing the fungus.  The amount of conidia produced in each gram substrate was determined using the improved Neubuer Hemacytometer counting chamber (0.100mm deep), as described in the section on spore counts[12].

*(d) Formulations for Laboratory bioassays*

The formulations evaluated were *B. bassiana* grown on cracked maize seed, maize bran and ''machicha'', applied alone or formulated with loam soil or clay soil.  The formulations were chosen depending on their conidia yields.  The loam soil was collected from the banana field at KARI with the physical characteristics of estimated levels of sand (52%), silts (28-50%), clay (7-28%), and high water holding capacity (23%).  The clay soil used was the grey type, mined from water logged swamps, with particle size of approximately 0.002 mm.  Thus, eleven (11) *B. bassiana* formulations were evaluated under laboratory conditions and these are;

▪ Maize bran alone, maize bran + loam soil, maize bran + clay soil

▪ "Machicha" ("bussa") alone, "machicha" + loam soil, "machicha" + clay soil

▪ Cracked maize alone, cracked maize + loam soil, cracked maze + clay soil

▪ Loam soil alone or clay soil alone with nothing added.

The *B. bassiana* grown on cracked maize substrate was used as the standard. 1g of this substrate was mixed with 1g of the sterile formulation (1:1 ratio).  The 2g was then weighed into plastic petri-dishes and replicated 3 times.  The amount of conidia in each treatment was standardized to the same level as in cracked maize.  The amounts of the other substrates used depended on the amount of conidia per gram determined. They were also in the ratio of 1:1 per formulation.

Key dates associated with the trial:

*(a) B .bassiana* culturing and conidia counts:  29/11/01 - 20/05/02.

*(b)* Laboratory tests for the different *B. bassiana*  formulations:   24/05 - 24/06/02.

Experimental treatments

*(a) No. of substrates* = 8 (for objective 2i above); these are; Cracked maize, Maize bran, "Machicha", Cotton husks, Bagasse, Cotton husks + maize bran, Maize bran + bagasse, Bagasse + spent yeast.

*(b) No. of formulations* =11 (for objective 2ii above) and these are; Clay soil alone, Loam soil alone,  Cracked maize + clay, Maize bran + clay, "Machicha" + clay, Cracked maize alone, Maize bran alone, "Machicha" alone, Maize bran + loam, Cracked maize + loam, "Machicha" + loam soil.

*(c) No. of replicates per formulation* = 3; for each experiment.

*(d) No. of weevils per replicate* = 10 of mixed sex   (1:1 ratio).

Experimental design:

Completely randomised design (CRD), since the laboratory area used was uniform.  First the treatments were allocated to petri-dishes at random.  An area measuring 1x1m was marked on the laboratory bench.  The positions for placement of petri-dishes were marked on the bench and each petri-dish randomised to marked positions, using a table of random numbers.

---

[12] The full protocol included the corresponding details

---

### A.5.7 DATA:

Measurements:

(a) amount of conidia per unit gram of substrate.

(b) weevil mortality in the different formulations.

The numbers of dead weevils were recorded at different time points i.e. by observing the weevils after every 5 days for mortality, over a 30-day period.  Any dead weevils were removed, and put into a moist chamber and observed for any *B. bassiana* fungal growth.


### A.5.8 DATA MANAGEMENT:  *<Description of how data will be computerised, organised and managed and plans for data analysis procedures, together with lists of data file names and other documentation.>*


### A.5.9 DATA FILE NAMES:  *<List of all data files, data recording format files, program files, reporting documents, planning documents.>*


### A.5.10 DATA ANALYSIS PLAN:  *<Identification of the specific objectives of the analysis, listing variables to be used, noting steps needed to organise the data into the right format for analysis and an indication of the type of approach to be undertaken during the data analysis and the software to be used>*


*END OF APPENDIX 5*

# A.6 Field instructions for participatory assessments

The example below, drawn from a DFID Crop Protection Programme research project[13] on promotion of pest management strategies in Uganda, is aimed at illustrating the level of detail needed in the component of the protocol which describes details of the data collection procedure. In projects involving extensive data collection activities, this would be a part of a separate Field Instructions Manual, but in other cases, as below, it would be sufficient to include such details in the *Materials and Methods* section of the protocol. This example is referred to in the main text section 2.2.1. It provides brief background details but omits other components included in the complete activity protocol.

---

### A.6.1 ACTIVITY TITLE: Farmer perceptions of technologies for banana pest and disease management in Luwero district, Uganda.

### A.6.2 BACKGROUND (Why activity is being done):

In a previous project (see Appendix 1) on-farm trials were established to explore a range of technologies for the management of pests and diseases in banana. The technologies included adopting good management practices and the use of improved cultivars. A subsequent project was set up to promote promising technologies. One promotional activity concerned promoting cultivars that performed well on farmers' fields. As a preliminary to promoting these cultivars, the research team took the view that it was important to seek farmers' views on criteria they would use when choosing cultivars to grow on their farms, and based on these criteria, to determine which cultivars were most preferred by banana farmers. The activity below relates to the first of these.

### A.6.3 OBJECTIVES:

To determine criteria that banana farmers would use in choosing cultivars to grow on their farms.

### A.6.4 MATERIALS AND METHODS (Addressing the "how"):

Field data collection procedure

*(a) Materials needed:* Flip chart, many cards, flip chart pens, a bag of bean seeds, an empty bag (or container) a stapler, a notebook. In the field, need a suitable spot for having the discussion.

*(b) Preliminaries:* After introductions, record the name, age, educational status, of farmers and their village. Also ask them individually if they use banana as a main crop for their livelihoods and whether they use banana as a food crop, a cash crop or as both a food and cash crop. Explain the purpose of the meeting, mentioning in particular that promotional activities in three districts are intended, and that the researchers need to know which cultivars should be promoted. This decision needs to be based on farmers' perceptions of which cultivars they like, and this in turn requires (as a first step) identifying the criteria that farmers would consider before growing the cultivars in their own fields.

---

[13] Project Leader: Mike Rutherford (CABI Bioscience) and Simon Gowen (School of Agriculture, University of Reading) with collaborating partners W. Tushemereirwe and Caroline Nankinga Kukiriza (Banana Research Programme, Uganda), Lawrence Kenyon (NRI, Chatham) and Tim Wheeler (University of Reading).

*(c) Identifying criteria:*  Invite the group to tell you about the different criteria they would consider in selecting a cultivar to grow on their farm.  Write out their suggestions on separate cards.  As each criterion is mentioned, the cards on which the criterion is written can be placed (visibly) for farmers to see, e.g. on the ground.  Make sure there are no further suggestions, particularly from those who have thus far been quiet.

Now show them the (pre-prepared) flip-chart (or cards) list of cultivars (Table 2) and ask whether they would like to add any more criteria to the list if asked to evaluate these cultivars.  Ask them which cultivars are familiar to them, and taking each in turn, ask them about additional criteria, if any, they would like to add to their previous list.

Suggestions will be written on new cards, all cards being visible to the whole group. (Note: If there are illiterate farmers in the group, there will be a need to draw some symbols or pictures to indicate what the different criteria are).

Now consider the criteria set as a whole.  Ask which of the listed criteria would be the most important to consider if they (the farmers) were selecting a cultivar to grow on their farm.  <u>It is important at this stage to note down reasons the farmers give for regarding a particular criterion as important</u>.  Include such comments under section 3 of the recording schedule shown in Appendix 1.  As the discussion progresses, re-arrange the cards so that the important ones are at the top, and the ordering is according to farmers' perception of importance.  (Least important criteria would be at the lowest end of the list of suggestions).

*(d) Secret voting:*  Hand out 5 bean seeds from your bag of seeds to each farmer.  Take the top criterion and tell them that you would like them to allocate more seeds to this criterion if they think its very important, and few seeds if it is not too important.  Pass an empty bag round the group and ask them to put into the bag, either 1, 2, 3, 4 or 5 seeds, according to their views.  When the bag gets back to you, empty contents of the bag next to the card for the corresponding criterion.  Count the seeds and note down the total number of seeds, either on another card placed next to the criterion card or on the same card which named the criterion.  Pass the (empty) bag round again to collect the remaining seeds.  Then again give out 5 seeds per person and carry out the voting for the second criterion and so on.

It would be good to have one pilot run of this first to ensure that the group are clear about what they have to do.  Then repeat for the criterion concerned and also for each of the other remaining criteria.

Once all the results on the cards are visible on the ground, point out to the group the criteria that appear to be the most important (i.e. ones having the highest totals), and check they are happy with the differences in importance.  At this time (and indeed during the whole field exercise) it would be very desirable that the second facilitator notes down comments made by the farmers in a notebook.

If during this exercise participants wish to discuss criteria in relation to the cultivars you showed them, then there should be flexibility to do so.

*(e) Leaving the group:*  Use appropriate ways to take your leave from the group after thanking them for their time.

*(f) Post-FGD work:*  After each focus group discussion (FGD), record the information from the flip-charts onto the de-briefing document.  Once all the FGDs have been completed, use pre-prepared computer formats to computerise the data.

*END OF APPENDIX 6*

# A.7 Identifying research hypotheses and variables in an analysis plan

The example below, drawn from a DFID Natural Resources Systems Programme (NRSP) [14]Project in Bangladesh, provides an illustration to section 2.2.5 of the main text where the importance of setting up research hypotheses and associated variables during the preparation of an analysis plan was emphasised.  Unlike in Appendices 1 to 5, we provide just a little background to the project and then move directly to the main point we wish to illustrate.   Full details may be found in Sultana *et al* (2007).

In the activity we consider below, a new participatory planning approach, aimed at facilitating and developing an appropriate system for community based fishery management (CBFM), is compared with other approaches to CBFM.  To achieve this, the researchers first identified a series of clearly stated research hypotheses, and spent a considerable period of time (several days) to consider very carefully the way in which specific variables could be collected and used in an analysis to address each hypothesis. We illustrate below how a research hypothesis matrix was constructed to help this process.

---

### A.7.1  ACTIVITY TITLE : Evaluating the effectiveness of Participatory Action Plan Development (PAPD)

### A.7.2  BACKGROUND (Why activity is being done):

This activity builds on an earlier consensus building study which attempted to assess the impacts of PAPD immediately after the workshop process through interview surveys to assess social capital.  The approach could however be criticised for not assessing organisational and institutional change, or changes in community action in resources management, or impacts on the livelihoods of stakeholders.  This activity will address this gap based on comparative assessments of the PAPD approach with other forms of CBFM activities.

### A.7.3  OBJECTIVES:

To evaluate the effectiveness of the PAPD method compared to other forms of approach towards sustainable co-management of fishery resources.

### A.7.4  DATA ANALYSIS PLAN: (*Comment:  Only a component of the data analysis plan is shown below*).

We describe here in outline the steps undertaken to address the research objective through the construction of a structured hypothesis matrix as shown in Figure 7.4.1.

First, it was necessary to identify a series of research hypotheses that would support the claim that the PAPD planning process was indeed more effective than other approaches to community based fishery management.  The first column of Figure 7.4.1 lists all the hypotheses, grouped according to the type of hypothesis.  In the next step, the researchers identified one or more responses that would allow each research hypothesis to be tested.  These responses are shown in the second column of Figure 7.4.1.

The third column shows in detail how each response variable can be measured and the data source.  Some of these (e.g. item **a2** in 3$^{rd}$ column of hypothesis iii) are still in need of further clarification as to how exactly they are calculated, and this was done in the full data analysis plan.  The final column identifies variables that may confound the comparison of sites where PAPD was done and sites where it was not.  The data analysis can adjust for such confounders and the listing assists the data analysis process to a considerable degree.

---

[14] This project is a DFID NRSP programme development project led by Parvin Sultana, an independent consultant in Bangladesh.  Her permission to use this example is gratefully acknowledged.

---

## Table 1. Research hypotheses, data sources and confounding variables

| Achievable research hypothesis | Main response variable(s) to address hypothesis | Variables (with data file names) contributing to main response (& data source[15]) | Confounding factors and variables relating to data structure (and source) |
|---|---|---|---|
| *Community Based Organisation (CBO) development* | | | |
| i. PAPD results in faster setting up of community based organisations (CBOs) | Number of days from start of CBFM activities to the formation of a water body management committee (CBO). | Date of NGO signing contract to undertake CBFM work and date of first forming the CBO - from Quarterly Monitoring Report (QMR). | Waterbody type |
| ii. PAPD results in more active CBOs. | (a) Average number of CBO meetings per month since start of CBO. (b) % attendance at CBO meetings. (c) Number of awareness raising activities with organisations outside the CBO. (d) Number of conflicts resolved by CBO. | (a) and (b): Number of CBO meetings and % attended are obtained from CBO QMR. (c) from FGD-MC. (d) Resolved number of internal and external conflicts from IMF (only where conflicts did take place) | (a) None (b) None (c) Waterbody type, because floodplains have less fisheries activities during dry season (d) Size of waterbody; type of waterbody. |
| iii. PAPD results in the formation of CBOs that are more holistic, and where poor are better represented. | (a) Number of categories of stakeholders involved in the CBO. (b) Proportion of poor fishers and landless in CBO. | (a) Number of different stakeholder types (fishers, farmers, landless, official, other), on scale of 1-5. (b)Number in membership of fishers and landless, and total number of members. All above from IMF. | (a) & (b) Waterbody type and size (For some sites, estimates but not accurate information on area were available from CBFM-2). |
| *Social capital* | | | |
| iv. PAPD results in greater social cohesion | (a) Measure of the degree of change in social cohesion in community. | (a) From FGD, an average over participants' −5 to +5 measure of the change in social cohesion since start of CBFM. Obtained from FGD-MC and FGD-Fsh. | For both (a) and (b): (i) waterbody type. (ii) number of other development activities in the area. (iii) number of categories of other uses of fishery. (iv) Date since CBO formation (QMR). (v) proportion of fishers in the waterbody catchment area who fish for an income (census). (vi) Proportion of "better-offs" in community (from census). |
| *Sustainability of fishery* | | | |
| v. PAPD results in greater community awareness and concern for collective sustainability and security actions. | Measure of community interests in sustainability and security of the fishery as judged by benefits to self and family, short-term benefits to community , long-term benefits to community. | Sum of scores given by FGD-MC for importance of benefits they list. Sum of scores given by FGD-Fsh for importance of benefits they list. Six summaries will result (MC-Fsh vs 3 groups of benefits). | (i) Waterbody type. (ii) Total number of fishers - census. (iii) No. of categories of other uses of fishery. |

[15] QMR – Quarterly Monitoring Report;   IMF – Institutional-Monitoring Form; FGD-MC – Focus Group Discussions with Management Committee; FGD-Fsh – Focus Group Discussions with fishers and traders

| Achievable research hypothesis | Main response variable(s) to address hypothesis | Variables (with data file names) contributing to main response (& data source[15]) | Confounding factors and variables relating to data structure (and source) |
|---|---|---|---|
| **Collective action** | | | |
| vi. PAPD results in faster uptake of community actions for NR management. | (a) Number of days between action date (date first key fishery management action was implemented) and start of CBFM activities (fielding of staff). (b) Number of days between action date and CBO formation. | (a) From QMR, date CBFM started) and date first action implemented. (b) From QMR, date of first formation of CBO and date first action implemented. | Waterbody type for both (a) and (b). |
| vii. PAPD results in more community/ collective actions for NR management. | (a) Number of actions planned and not implemented. (b) Number of actions implemented. | Cumulative numbers from QMR for these variables. | Number of conflicts from IMF. Number of categories of other uses of the fishery. Number of other development activities in area. |
| viii. PAPD results in community actions with greater compliance. | (a) Number of rules in place out of total relevant. (b) Number of rule breaking incidents. (c) % of community reported know of rules. (d) Number of conflicts. | All from IMF. (a) number of rules (and number of actions planned (ticked) on page 1 of QMR). (b), (c) and (d) similarly obtained. | Waterbody type, and waterbody size. % fishers in CBO – QMR. Number of categories of other uses of the fishery. Number of other development activities in area. |
| **Livelihood outcomes and linkages** | | | |
| ix. PAPD results in community actions involving wider coverage of communities that perceive benefits. | Number of stakeholder categories that may benefit (or have benefited) from CBFM. | FGD – number of categories benefited, yes/no if reported that key stakeholders (poor, fishers, poor women) benefited. | Number of different stakeholders in FGD. |
| x. PAPD results in better links with local government | (a) Whether stakeholders get support from government and form of support. (b) Attitude and understanding of CBFM in local government. (c) Number of links fishers have with outside groups (government or otherwise). | (a) FGD – types and numbers of government support incidents – MC group only. (b) FGD – Union Parishad[16] and Upazila[17] attitude rating; assessment by MC members of officials understanding – MC group only. (c) FGD – number of links with local government organisations and others. | None that can be realistically collected. |
| **Time /transaction costs** | | | |
| xi. PAPD actions requires greater time input from participant communities. | Number of person days involved in CBFM activities in general. | Transaction cost assessment from FGD for (a) CBO leaders (as a group) and (b) for general resource users. | (a) Number of CBO members. (b) None. |

*END OF APPENDIX 7*

---

[16] Union Parishad – lowest level of government, an elected council covering about 10 villages
[17] Upazila – lowest administrative level of government based on officers of different agencies covering about 10 unions

# References

DFID-NRSP (2002) *Scaling-up and communication: Guidelines for enhancing the developmental impact of natural resources research*.   DFID Natural Resources Systems Programme, 8 pp.

Lawson-McDowall, J.M., Abeyasekera, S., Mwale, B., Ritchie, J.M., Orr, A. and Chanza, C. (2001)  *IPM for Smallholders: A Researcher's Casebook from Malawi.*   Natural Resources Institute, Chatham.   [Note: an updated version of the chapter on Data Management is available as Case Study 6 in SSC (2001d) below]

Lindsey, J.K. (1999) *Revealing Statistical Principles*.   Arnold,  ISBN 0 340 74120 1

Patel, B.K. Muir-Leresche, K. Coe, R. and Hainsworth, S.D. (2004) *The Green Book: a guide to effective graduate research in African agriculture, environment and rural development.*   African Crop Science Society[18].   ISBN (CD version) 9970-866-00-1

SSC (1998) *Data Management Guidelines for Experimental Project*. Statistical Guidelines Series supporting DFID Natural Resources Projects, Statistical Services Centre, The University of Reading, UK.  www.reading.ac.uk/ssc/publications/guides.html

SSC (2000a) *Disciplined use of Spreadsheets for Data Entry*. Statistical Guidelines Series supporting DFID Natural Resources Projects, Statistical Services Centre, The University of Reading, UK.  www.reading.ac.uk/ssc/publications/guides.html

SSC (2000b) *Some Basic Ideas of Sampling*. Statistical Guidelines Series supporting DFID Natural Resources Projects, Statistical Services Centre, The University of Reading, UK.  www.reading.ac.uk/ssc/publications/guides.html

SSC (2001a) *Approaches to the Analysis of Survey Data*. Statistical Guidelines Series supporting DFID Natural Resources Projects, Statistical Services Centre, The University of Reading, UK.  www.reading.ac.uk/ssc/publications/guides.html

SSC (2001b) *Modern Approaches to the Analysis of Experimental Data*.  Statistical Guidelines Series supporting DFID Natural Resources Projects, Statistical Services Centre, The University of Reading, UK.  www.reading.ac.uk/ssc/publications/guides.html

SSC (2001c) *Modern Methods of Analysis*. Statistical Guidelines Series supporting DFID Natural Resources Projects, Statistical Services Centre, The University of Reading, UK.  www.reading.ac.uk/ssc/publications/guides.html

SSC (2001d) *Case Studies of Good Statistical Practice*.  The University of Reading Statistical Services Centre Guideline Series for DFID, available at http://www.rdg.ac.uk/ssc/workareas/development/case_studies.html

Stern, R.D., Coe, R., Allan, E.F., and Dale, I.C. (eds) (2004). *Good Statistical Practice for Natural Resources Research*.  CAB International, Wallingford, UK. 388 pp.  ISBN 0-85199-722-8

Sultana, P., Abeyasekera, S. Thompson, P. (2007)  Methodological rigour in assessing participatory development.  *Agricultural Systems,* 94, pp. 220-230.

Van Belle, G. (2002) *Statistical Rules of Thumb*.  Wiley, ISBN 0-471-40227-3

*Wilson, I.M. and Huttly, S.R.A. (2003) Young Lives : A Case Study of Sample Design for Longitudinal Research.   Young Lives Working Paper 10,* ISBN 1-904427-11-1

---

[18] African Crop Science Society, Faculty of Agriculture and Forestry, Makerere University, P.O. Box 7062, Kampala, Uganda.