

BROOKE ANDERSON, MICHAEL LYONS, MERCEDES GONZALEZ-  
JUARRERO, MARCELA HENAO-TAMAYO, AND GREGORY ROBERT-  
SON

# IMPROVING THE REPRODUCIBIL- ITY OF EXPERIMENTAL DATA RECORDING AND PRE-PROCESSING



# Contents

<i>1</i>	<i>Overview</i>	<i>5</i>
<i>1.1</i>	<i>License</i>	<i>7</i>
<i>2</i>	<i>Experimental Data Recording</i>	<i>9</i>
<i>2.1</i>	<i>Separating data recording and analysis</i>	<i>9</i>
<i>2.2</i>	<i>Principles and power of structured data formats</i>	<i>9</i>
<i>2.3</i>	<i>The ‘tidy’ data format</i>	<i>10</i>
<i>2.4</i>	<i>Designing templates for ‘tidy’ data collection</i>	<i>10</i>
<i>3</i>	<i>Experimental Data Pre-Processing</i>	<i>11</i>
<i>4</i>	<i>Bibliography</i>	<i>13</i>



# I

## Overview

The recent NIH-Wide Strategic Plan (U.S. Department of Health and Human Services, National Institutes of Health, 2016) describes an integrative view of biology and human health that includes translational medicine, team science, and the importance of capitalizing on an exponentially growing and increasingly complex data ecosystem (U.S. Department of Health and Human Services, National Institutes of Health, 2018). Underlying this view is the need to use, share, and re-use biomedical data generated from widely varying experimental systems and researchers. Basic sources of biomedical data range from relatively small sets of measurements, such as animal body weights and bacterial cell counts that may be recorded by hand, to thousands or millions of instrument-generated data points from various imaging, -omic, and flow cytometry experiments. In either case, there is a generally common workflow that proceeds from measurement to data recording, pre-processing, analysis, and interpretation. However, in practice the distinct actions of data recording, data pre-processing, and data analysis are often merged or combined as a single entity by the researcher using commercial or open source spreadsheets, or as part of an often proprietary experimental measurement system / software combination (Figure 1.1), resulting in key failure points for reproducibility at the stages of data recording and pre-processing.

It is widely known and discussed among data scientists, mathematical modelers, and statisticians (Broman and Woo, 2018; Krishnan et al., 2016) that there is frequently a need to discard, transform, and reformat various elements of the data shared with them by laboratory-based researchers, and that data is often shared in an unstructured format, increasing the risks of introducing errors through reformatting before applying more advanced computational methods. Instead, a critical need for reproducibility is for the transparent and clear sharing across research teams of: (1) raw data, directly from hand-recording or directly output from experimental equipment; (2) data that has been pre-processed as necessary (e.g., gating for flow cytometry data, feature identification for metabolomics data), saved in a consistent, structured format, and (3) a clear and repeatable description of how the pre-processed data was

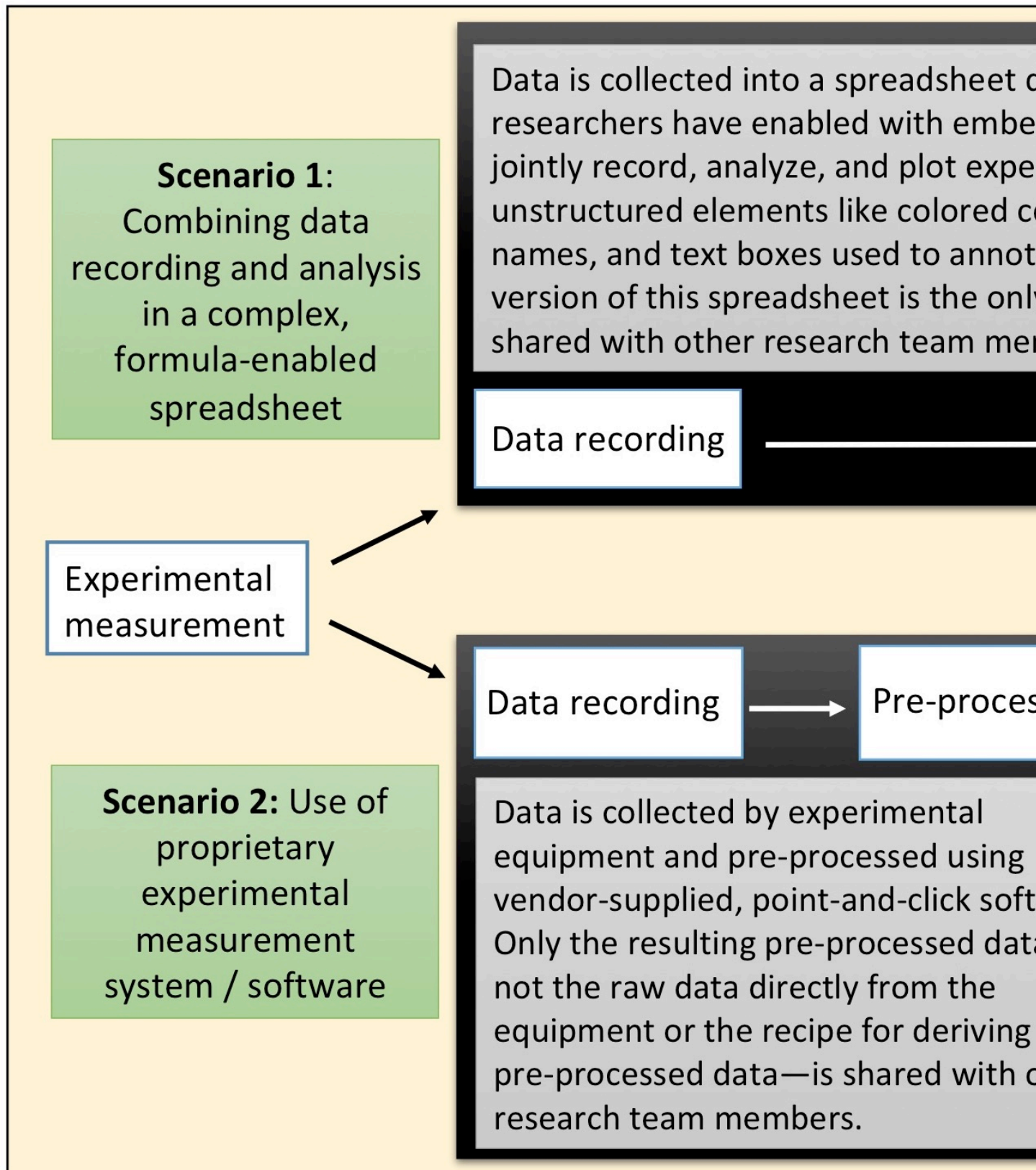


Figure 1.1: Two scenarios where 'black boxes' of non-transparent, non-reproducible data handling exist in research data workflows at the stages of data recording and pre-processing. These create potential points of failure for reproducible research. Red arrows indicate where data is passed to other research team members, including statisticians / data analysts, often within complex or unstructured spreadsheet files.

generated from the raw data (Broman and Woo, 2018; Ellis and Leek, 2018).

To enhance data reproducibility, it is critical to create a clear separation among data recording, data pre-processing, and data analysis—breaking up commonly existing “black boxes” in data handling across the research process. Such a rigorous demarcation requires some change in the conventional understanding and use of spreadsheets and a recognition by biomedical researchers that recent advances in computer programming languages, especially the R programming language, provide user-friendly and accessible tools and concepts that can be used to extend a transparent and reproducible data workflow to the steps of data recording and pre-processing. Among our team, we have found that there are many common existing practices—including use of spreadsheets with embedded formulas that concurrently record and analyze experimental data, problematic management of project files, and reliance on proprietary, vendor-supplied point-and-click software for data pre-processing—that can interfere with the transparency, reproducibility, and efficiency of laboratory-based biomedical research projects, problems that have also been identified by others as key barriers to research reproducibility (Broman and Woo, 2018; Bryan, 2018; Ellis and Leek, 2018, Marwick et al. (2018)). In these training modules, we have chosen topics that tackle barriers to reproducibility that have straightforward, easy-to-teach solutions, but which are still very common in biomedical laboratory-based research programs.

## *1.1 License*

This book is licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License, while all code in the book is under the MIT license.

Click on the **Next** button (or navigate using the links at the top of the page) to continue.





## 2

# Experimental Data Recording

### 2.1 *Separating data recording and analysis*

Many biomedical laboratories use spreadsheets, with embedded formulas, to both record and analyze experimental data. This practice impedes transparency and reproducibility of both data recording and data analysis. In this module, we will describe this common practice and will outline alternative approaches that separate the steps of data recording and data analysis.

**Objectives.** After this module, the trainee will be able to:

- Explain the difference between data recording and data analysis
- Understand why collecting data on spreadsheets with embedded formulas impedes reproducibility
- List alternative approaches to improve reproducibility

#### 2.1.1 *Data recording versus data analysis*

#### 2.1.2 *Common practices combining recording and analysis*

#### 2.1.3 *Hazards of combining recording and analysis*

#### 2.1.4 *Approaches to separate recording and analysis*

#### 2.1.5 *Discussion questions*

### 2.2 *Principles and power of structured data formats*

The format in which experimental data is recorded can have a large influence on how easy and likely it is to implement reproducibility tools in later stages of the research workflow. Recording data in a “structured” format brings many benefits. In this module, we will explain what makes a dataset “structured” and why this format is a powerful tool for reproducible research.

**Objectives.** After this module, the trainee will be able to:

- List the characteristics of a structured data format
- Describe benefits for research transparency and reproducibility
- Outline other benefits of using a structured format when recording data

### 2.2.1 *Characteristics of a structured data format*

### 2.2.2 *Benefits of a structured data format*

### 2.2.3 *Applied exercise*

## 2.3 *The ‘tidy’ data format*

The “tidy” data format is an implementation of a structured data format popular among statisticians and data scientists. By consistently using this data format, researchers can combine simple, generalizable tools to perform complex tasks in data processing, analysis, and visualization. We will explain what characteristics determine if a dataset is “tidy” and how use of the “tidy” implementation of a structure data format can improve the ease and efficiency of “Team Science”.

**Objectives.** After this module, the trainee will be able to:

- List characteristics defining the “tidy” structured data format
- Explain the difference between the a structured data format (general concept) and the ‘tidy’ data format (one popular implementation)

### 2.3.1 *The “tidy” data format*

### 2.3.2 *The “tidy” data format as a structured data format*

### 2.3.3 *Practice quiz*

## 2.4 *Designing templates for ‘tidy’ data collection*

### 2.4.1 *Subsection 1*

### 2.4.2 *Subsection 2*

3

*Experimental Data Pre-Processing*



## 4

### *Bibliography*

Broman, K. W. and Woo, K. H. (2018). Data organization in spreadsheets. *The American Statistician*, 72(1):2–10.

Bryan, J. (2018). Excuse me, do you have a moment to talk about version control? *The American Statistician*, 72(1):20–27.

Ellis, S. E. and Leek, J. T. (2018). How to share data for collaboration. *The American Statistician*, 72(1):53–57.

Krishnan, S., Haas, D., Franklin, M. J., and Wu, E. (2016). Towards reliable interactive data cleaning: A user survey and recommendations. In *Proceedings of the Workshop on Human-In-the-Loop Data Analytics*, page 9. ACM.

Marwick, B., Boettiger, C., and Mullen, L. (2018). Packaging data analytical work reproducibly using R (and friends). *The American Statistician*, 72(1):80–88.

U.S. Department of Health and Human Services, National Institutes of Health (2016). NIH-Wide Strategic Plan, Fiscal Years 2016-2020: Turning Discovery Into Health. Accessed: 2018-06-24.

U.S. Department of Health and Human Services, National Institutes of Health (2018). NIH Strategic Plan for Data Science. Accessed: 2018-06-24.