# Visual grids for managing data completeness in clinical research datasets

Robert R. Kelley [a], William A. Mattingly [a,*], Timothy L. Wiemken [a], Mohammad Khan [a], Daniel Coats [a], Daniel Curran [a], Julia H. Chariker [b], Julio Ramirez [a]

[a] Division of Infectious Diseases, Department of Medicine, University of Louisville, Louisville, KY, USA
[b] Department of Psychological and Brain Sciences, University of Louisville, Louisville, KY, USA

## ARTICLE INFO

## ABSTRACT

Missing data arise in clinical research datasets for reasons ranging from incomplete electronic health records to incorrect trial data collection. This has an adverse effect on analysis performed with the data, but it can also affect the management of a clinical trial itself. We propose two graphical visualization schemes to aid in managing the completeness of a clinical research dataset: the binary completeness grid (BCG) for single patient observation, and the gradient completeness grid (GCG) for an entire dataset. We use these tools to manage three clinical trials. Two are ongoing observational trials, while the other is a cohort study that is complete. The completeness grids revealed unexpected patterns in our data and enabled us to identify records that should have been purged and identify missing follow-up data from sets of observations thought to be complete. Binary and gradient completeness grids provide a rapid, convenient way to visualize missing data in clinical datasets.

© 2015 Elsevier Inc. All rights reserved.

## 1. Introduction

The initial product of a clinical trial is the dataset collected during its execution. Missing data in such a dataset complicates the statistical analysis process. In fact, missing data may preclude certain types of analysis because there is not enough data to analyze; in the extreme case, the entire dataset may not be useful. However, discovering that data are incomplete after the data collection period is often not sufficient, particularly in the case of a study in which follow ups at various time points are collected. Therefore, managing missing data during the execution of a clinical trial is a critical aspect of the clinical trial process.

Data for medical research studies, such as retrospective chart reviews or prospective randomized clinical trials (RCT), are typically collected using paper or electronic case report forms (CRFs) [1]. An investigator, or a designee, examines the electronic health record (EHR) or paper chart for data required by the study, and then transfers these data to a CRF by hand or by retyping it into an electronic system. In the case of a study in which follow ups are required, the process usually involves the investigator contacting the subject by

telephone or through a follow up clinical visit to collect further data on their condition and progress.

Study data collection policies and practices and differences in electronic health records systems (EHRs) have the most effect on the quality of data collected for a clinical trial. First, if data are missing or not captured in the EHR, it will also be missing in the CRF and subsequently the study database. Moreover, data required for a clinical trial is often stored as unstructured free text in an EHR; data collection personnel are required to locate these data, which are often located in clinician's notes, which may necessitate a clinician's experience to successfully parse. Furthermore, in a facility still using paper charts, potential legibility problems may interfere with transcribing data to the CRF that, in turn, leads to missing data in the study database.

Outside of the EHR, typographical errors and misunderstanding of data entry policy by those entering data may also lead to data quality issues. Several preventative measures can be used to address these problems including: documenting the data collection process (e.g. manual of operations), training data collection staff on proper data collection procedures, performing double-data entry, fostering frequent communication between investigators and data collectors regarding potential data issues, and using monitoring reports to track the state of the dataset [2].

The data quality literature frequently categorizes missing data as "completeness", a characteristic of the Contextual Information

* Corresponding author at: Division of Infectious Diseases, Department of Medicine, University of Louisville, 501 East Broadway, Suite 140B, Louisville, KY 40202, USA. Fax: +1 502 852 1555.

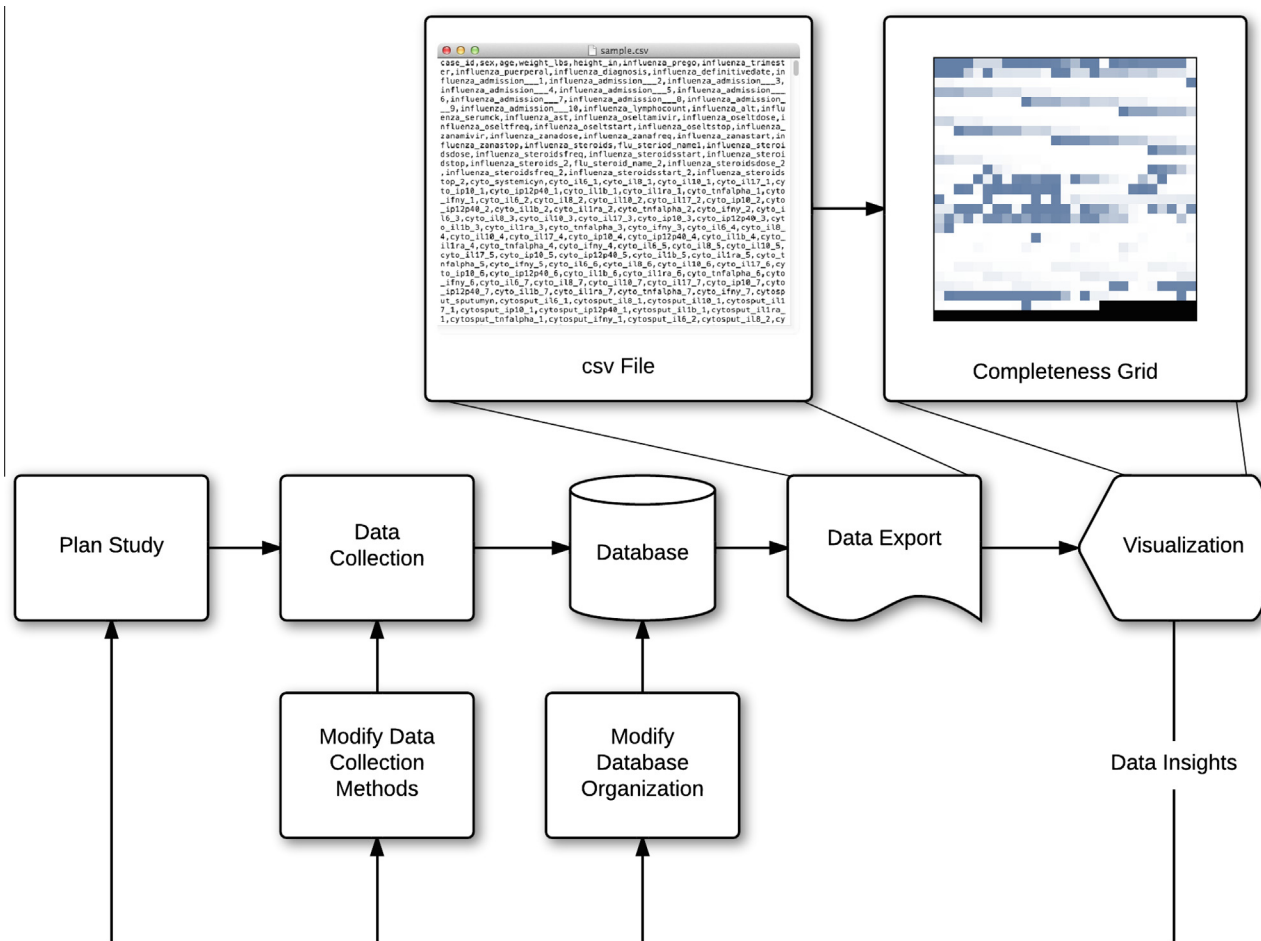E-mail address: wamatt02@louisville.edu (W.A. Mattingly).

**Fig. 1.** Clinical study workflow using completeness grid for visualization. After data has been collected, completeness grids can be created which inform the processes of plan study, data collection, and database organization.

Quality dimension [3–6]. Completeness is specifically defined as "the extent to which data are not missing" [7]. The statistical analysis literature is replete with warnings about the effects of missing data on research results, as well as approaches to mitigate the problems it causes [8–11]. The two most significant issues with missing data are that they can introduce bias into the statistical results and they reduce the power of the study. Common techniques for handling missing data include complete case analysis, analysis of observations available, using data imputation to "fill-in" data with statistical methods, and intent to treat [12].

Missing data can be more effectively recognized using data visualization techniques. Data visualization is a means of representing data in a graphical form allowing the recognition of details and patterns that might otherwise remain obscure. Ways to represent data while taking into account missing or incomplete datasets is an active area of research in data visualization and data quality. Many visualization methods for missing data depend on knowing the uncertainty associated with data [13]. Different types of visual attributes, like color hue and texture can be used in a dataset visualization to draw attention to uncertainty values [14,15]. However, in the case of many sparse datasets, missing data is the norm, and estimation and imputation should be avoided in favor of alerting the viewer as to which data fields are empty [16].

There are also many statistical software packages for visualizing and interacting with missing data. Missing Are Now Equally Treated (MANET), XGobi, GGobi and Mondrian, were all designed to allow users to explore and visualize data and missing data [17–19]. The R programming language contains many packages and functions for data analysis and visualization, and the VIM package has been designed specifically to facilitate the analysis and interactive exploration of missing data [20].

One drawback of previous missing data visualizations is their interdependence with complex data analysis. Many stakeholders in the clinical data collection process are interested in quickly and easily viewing the completeness of a dataset, but do not have the background to engage in complex analysis. To facilitate this we have developed two visualization methods, the binary completeness grid and the gradient completeness grid, tools for clinical trial data managers to monitor and evaluate missing data and effectively visualize large two-dimensional datasets like those collected for clinical trials.

Fig. 1 shows a typical workflow for using completeness grids. As data is collected for a study and entered into a database, the entire dataset or subsets can be exported in an easy-to-process format such as a comma-separated value (csv) file. Completeness grids can then be generated providing a visual representation of the completeness of the dataset. Using this information, data monitors can investigate unexpected patterns and resolve potential problems with the database or data collection methodology.

## 2. Materials and methods

### 2.1. Completeness grids

Completeness grids account for the visualization needs of individual users in a clinical study. Investigators are generally involved