

Anonymisation of Clinical Trial Datasets

1. Introduction

Providing access to data in ways that allows further research while maintaining the privacy and confidentiality of research participants is critical. There are also privacy laws and regulatory guidance which need to be followed (for example guidance from European data protection regulators and Code of Federal Regulations - Title 45: Public Welfare, Subtitle A §164.514). Publications in this area which provide guidance^{1,2}.

This document describes the approach taken by a number of the study sponsors to prepare data for sharing with other researchers in a way that:

- Minimises risks to the privacy and confidentiality of research participants.
- Ensures compliance with data privacy legal requirements.

Other study sponsors achieve these objectives using other approaches (see the members page on vivli.org).

2. General Approach

Access is provided to anonymised data. Anonymisation involves:

- Removing personally identifiable information (PII) from the dataset.** This includes recoding identifiers (by replacing the original code number with a new code number), removing free text verbatim terms, Replacing date of birth with year of birth or age and replacing all dates relating to individual subjects with dummy dates or replacing them with a study day.
- Destroying the link (code key) between the dataset that is provided and the original dataset.** Some Data Protection Authorities in Europe suggest that the data can only be considered anonymised if personal information is removed (or redacted) and the subject code number cannot be linked to a research participant. Therefore, research participants' identification code numbers are anonymised by destroying the code key that was used to generate the new code number from the original (i.e. destroying the link between the two code numbers).

¹ Hrynaszkiewicz I, Norton ML, *et al.* Preparing raw clinical data for publication: guidance for journal editors, authors, and peer reviewers. *BMJ* 2010; **340**: c181.

² De-identification of Clinical Trials Data Demystified. *Jack Shostak, Duke Clinical Research Institute (DCRI), Durham, NC* <http://www.lexjansen.com/pharmasug/2006/publichealthresearch/pr02.pdf>

3. Removing personally identifiable information (PII) from the dataset

The 18 identifiers (as defined by HIPAA –see [Code of Federal Regulations - Title 45: Public Welfare, Subtitle A §164.514](#)) are removed from the datasets (and related documentation). In addition any other PII that may be present is removed.

This involves removing:

- any names and initials,
- (or recoding) kit numbers and device numbers
- geographic information such as place of work.
- some sponsors also remove socioeconomic data such as occupation, income or education. Household and family composition. Multiple pregnancies.

In addition the following steps are undertaken:

- Recoding identifiers (or code numbers).
- Removing free text verbatim terms.
- Replacing date of birth with year of birth or age at randomisation. Ages above 89 which are aggregated into a single category of “90 or older”. (This is a specific HIPAA requirement).
- Replacing all original dates relating to individual subjects with randomly generated offsets which are then applied to create ‘dummy dates’ or deleting these dates and replacing them with a study day. (see below)
- Reviewing and removing other PII

These steps are described in further detail below.

3.1 Recoding Identifiers (or code numbers)

The following identifiers (code numbers) are re-coded and the code key that was used to generate the new code number from the original code number is destroyed (as described in section 5):

- The investigator identifier (or code number) is re-coded or set to blank for each investigator. The investigator name is set to “blank” or dropped from the dataset (see Appendix 1 & 2).
 - A new subject identifier (or code number) for each research participant.
 - Some sponsors also re-code the centre identification number.
 - Some sponsors aggregate patients from centres with less than 10 patients into a single centre.
- The same new identifiers (or code numbers) are used across all datasets applicable to a single study e.g. raw dataset, analysis-ready dataset. This includes (where applicable) PK datasets, genetic datasets etc.
 - Extension studies use the same new identifiers (or code numbers) as used for the initial study to enable individual subject data to remain linked. This also applies to long term follow-up studies where separate reports are published. This is achieved by repeating the data anonymisation process for the initial study data at the same time as the extension/follow up data.

3.2 Removing Free Text Verbatim Terms

Information in a descriptive free text verbatim term may compromise a subject’s anonymity.

- Free text verbatim terms are set to “blank” or dropped from the dataset including:
 - Adverse Events
 - Medications
 - Other e.g. Medical History
 - Other specific verbatim free text

Certain free text fields may be retained if they do not contain PII and removal of these fields may impact the scientific value of the dataset (e.g. medical history that has not been coded).

- All dictionary coded terms with decode and/or verbatim terms that use a pre-specified list are retained.

3.3 Replacing Date of Birth

Information relating to a research participant’s date of birth and identification of specific ages above 89 may compromise anonymity.

- Date of birth is replaced with the year of birth or age at randomisation with the exception of ages above 89 which are aggregated into a single category of “90 or older”

3.4 Replacing all Original Dates relating to a Research Participant

Study sponsors use one of two methods as described below.

3.4.1 Dummy Date Method

Specific dates (other than year) directly related to a research participant may compromise a research participant’s anonymity.

All dates are replaced: A random offset is generated for each research participant and applied to all dates for that research participant. All original dates are replaced with the new dummy dates so that the relative times for each research participant are retained.

Example: If the original reference date was 01APR2008 and the date of death was 01MAY2008, a random offset is generated (in this case 91 days). Dummy dates are then calculated using this offset of 91 days.

	Original Date	New Date	
Reference date	01APR2008	01JUL2008	Apply offset = 91 days
Date of Death	01May2008	31Jul2008	Apply offset=91 days
Relative Time of death	30 days	30 days	

3.4.2. Study Day Method

All dates are removed from the datasets. The Study Day is calculated for each observation with days relative to a reference date. In order of priority the reference date is defined as the date of first study treatment, date of randomisation or date of consent. For example if a patient is randomised, but does not take the study treatment (i.e. the date of first treatment is missing), the date of randomisation will be used as the reference date to calculate the study day for any assessments recorded.

Example If the original reference date was 01JAN2008 and the date of death was 01MAY2008, the date of death would be 122 expressed as Study Days.

	Original Date	Reference Date	Study Day
Date of Death	01May2008	01Jan2008	122

3.5 Reviewing and Removing Other PII

- Other data elements that contain PII are removed. For example:
 - Information from variable names e.g. lab names may contain location information
 - Investigator comments may be used to identify a subject
 - Genetic data that would enable a direct trace back to an individual subject

Appendix 1: Illustrates non-real examples of how these steps are applied.

4. Review and Quality Control

A final review of the HIPAA 18 identifiers is made to determine if further removal is required. Quality Control checks and documentation (QC record) is conducted for the processing of the data and supportive metadata documentation.

5. Destroying the link (key code) between the dataset that is provided and the original dataset

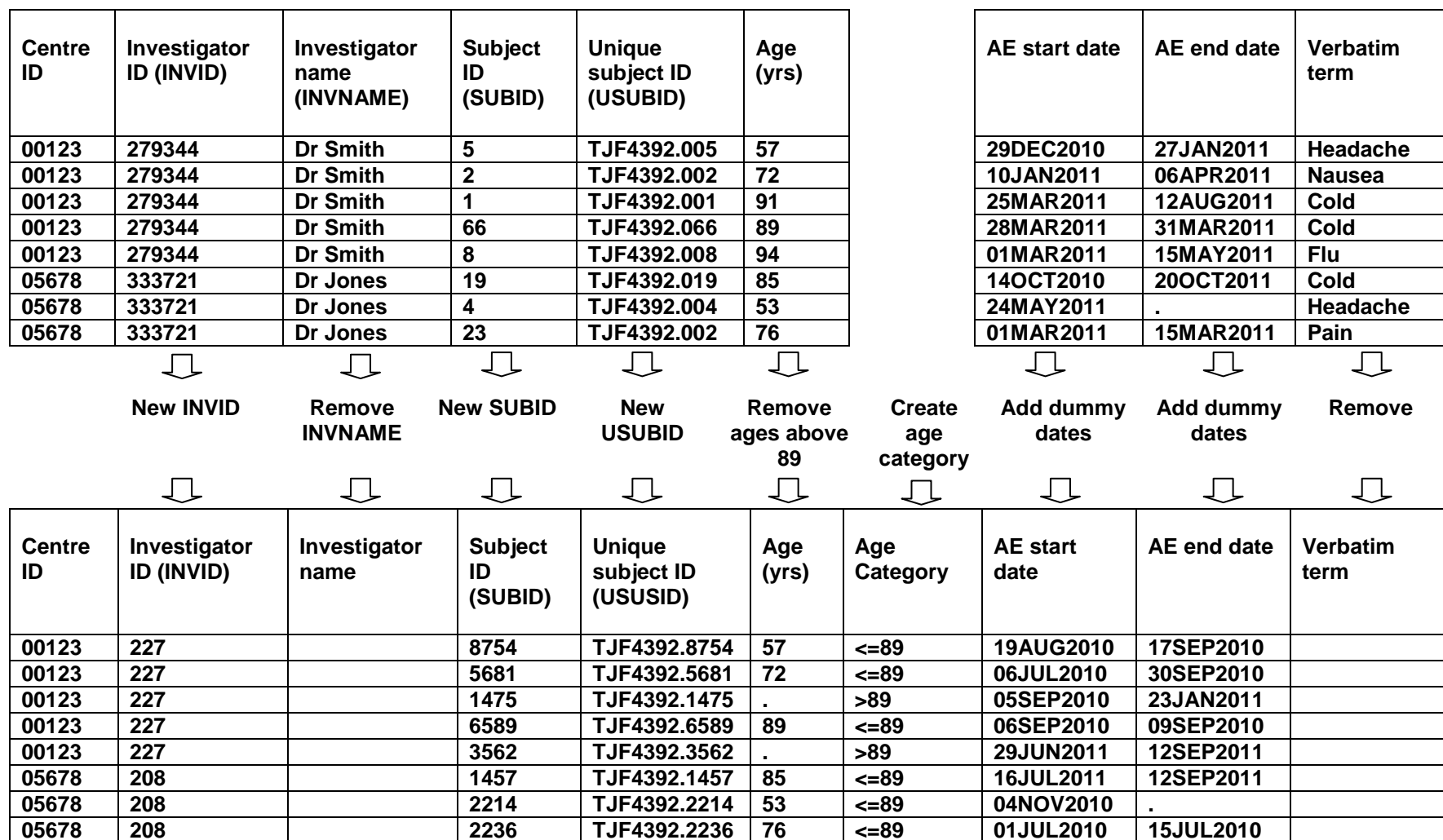
Research participants' identification code numbers are anonymised by replacing the original code number with a new code number (as described in 3.1) and destroying the code key that was used to generate the new code number from the original (i.e. destroying the link between the two code numbers).

The following specific items are discarded:

- Any transactional copies of anonymised datasets
- De-identification tables (links for original variable and new anonymised variable)
- Any QC output datasets
- Any Log or LST files
- The seed utilised for random number generation

The anonymised datasets are stored in a separate secure location to the original coded datasets.

Appendix 1: A non-real example illustrating removal of personally identifiable information using the dummy date method



Appendix 2: A non-real example illustrating removal of personally identifiable information using the study day method and aggregation of small centres

