
CPU-0008 CRF Level Data De-Identification Notes

And

Database Descriptions

October 31, 2017

CPU0008 Data De-Identification Notes

The Case Report Form (CRF)-level de-identified data sets were prepared by reviewing the CRFs and the data dictionary to identify variables that contained, or could potentially contain, identifying information, such as identifiers, dates and free text fields. These variables were programmatically removed or modified to produce de-identified datasets.

- Patient ID: All patient IDs were transformed to factless key values bearing no relationship to the collected study data. Technical Note: Original patient ID values (prior to transformation) in the AE dataset did not consistently include the leading zero found across all other datasets.
- Site ID: All site IDs were transformed to factless key values bearing no relationship to the collected study data.
- Dates: All dates except date of birth were converted to “days on study”, where the date of randomization = day 0. A negative value for days on study refers to an event that occurred prior to randomization (e.g. consent date, conmed initiated during screening) and a positive value for days on study refers to an event that occurred after randomization (e.g. treatment-phase visit and event dates).
 - Date of Birth: Dates of birth were converted to age values, computed as of randomization date.
 - Note that patients who failed screening will have no computed studyday values in date fields.
 - In cases where the date was stored as 3 distinct columns (month, day, year), these component value were converted to a date value before computing days on study; the resulting variable was named following the naming pattern of the 3 component variable names, but substituting the suffix DT.
 - Partial date sets of variables, e.g. just a Month and Year, may be erased if they refer to potentially identifying milestones, especially in a small patient pool.
 - NOTE: computed days-on-study values, used in place of dates, may disagree with native Study-Day variable values, which often equate date of randomization with day 1, not day 0.
- Initials / Names: Any field designed to record a person’s initials or name was erased.
- Other Identifiers: Sample identifiers and Drug lot numbers were erased.
- Text Fields: Most free text fields long enough to contain dates or narrative were emptied, including
 - Comment fields.
 - Responses derived from “Specify”, “Describe”, “Reason”, or “Other”. For example, other race, other termination reason, reason study drug not given, abnormalities on ECG and lab forms.
 - Narratives regarding interaction with patient, health conditions, treatment progress or follow-up.
 - Medical history and physical exam narrative text.
 - Type of work or job description; narrative descriptions of living situation, illnesses, concerns, drug/alcohol use, feelings/interests.
- Text Field De-identification Exceptions:
 - AE and medical condition terms were *retained* in the database, but detailed descriptions in a separate field providing additional information about adverse events were erased, such as relevant labs, medical history, drug/alcohol use description, probable cause of death.
 - Concomitant medication names and indications were *retained* in the database.
- For a complete listing of the erased text fields by CRF, refer to the “CPU0008 Nulled Values” csv file.

CPU0008 Data De-Identification Notes

See the "CPU0008 Dictionary" Excel workbook file for descriptions of the data, organized into tabbed worksheets for each respective CRF/dataset. The columns that have been nulled for de-identification purposes are indicated by a 'Y' in the Nulled column.