# Data Quality Checklist

**Objective**: Evaluate all data received for completeness, validity, consistency, accuracy, and interpretability in a timely manner so that:

- Basic and broad understanding of the data can be established as early as possible
- Problems with the data can be communicated to the client as early as possible, allowing sufficient flexibility in timeline should the data issues be severe enough to require re-extraction
- The client expectation can be managed effectively from the beginning

Key Things to Remember:
- Assume that all data are wrong until proven otherwise. Question everything.
- Do not try to interpret the data at this stage. This is a very fact-based exercise.
- No records are to be removed or contents altered with very few exceptions, until the final frequencies are confirmed with the client and the disposition of such data mutually agreed to. Absolutely no interpolation at this stage.
- Always consider the downstream impact of how you manage the data. Have justification for everything you do.

| From Receipt[1]: | At the Minimum[2]: |
|---|---|
| First 24 Hours:<br><br>Check *really* obvious facts | <ul><li>Verify files and sources against data specifications in the project design document</li><li>Verify that the files can be uncompressed and opened without errors</li><li>Verify presence of data dictionary for each file where needed<ul><li>Note that the layouts are absolutely required ONLY if the file has no column headers; however, the contents of the field may not be discernable from the column name alone</li></ul></li><li>Anything else that are really obvious</li><li>Acknowledge receipt of the data and communicate any issues to client via email:<ul><li>Confirm files received by name</li><li>List any missing files (check against the documented data requirements)</li><li>List any missing data dictionaries</li><li>State any issues with uncompressing and opening files</li><li>State any other really obvious issues</li><li>Otherwise state that a data quality analysis output will be delivered within 7 days barring critical data issues and will contact if any questions or issues come up in the meantime</li></ul></li></ul> |

---

[1] Of the respective file. These are milestones that should be reached at the end of the stated time, and it is expected that some work for the subsequent milestones is done in parallel as appropriate.

[2] Depending on the data or project, there may be other items that may need to be added to this list.

| First 48 Hours: Check obvious facts | <ul><li>Verify format of each file</li><li>Create a ditto/xxd output (or equivalent) for each file<ul><li>This is necessary especially for identifying non-displayable characters that impact your data load</li></ul></li><li>Visually inspect the data<ul><li>Unless the files are small, create head and tail samples for each file</li></ul></li><li>Verify layout of each file: do not assume that the layout is correct<ul><li>Delimited files: visually verify the field order against the raw data (not just the headers; also look at the contents), then verify that the length and the formats are correct</li><li>Fixed-column files: visually verify the positions and formats against raw data</li></ul></li><li>Verify record count of each file, using the line count or record count function that is native to the operating system<ul><li>Do not use SAS/R to do this: Software tools have their own interpretation of what a "line" is, and must isolate facts from software interpretation</li></ul></li><li>Verify fields required against the documented data requirements</li><li>Start to load raw files<ul><li>As the files are loaded, check that the loaded record count matches the record count obtained above</li><li>Determine the best informat based on the layout, your verification of the layout, the messages SAS provides in the log, your understanding of the data, and common sense</li><li>Do not use PROC IMPORT<ul><li>Assumes that file is clean, consistent, and relatively small—in practice this is rarely the case</li></ul></li><li>Verify the width and the format of the field, even if they are given in a layout or data dictionary<ul><li>They are not always correct</li><li>SAS could have its own interpretation which can vary from the definition native to the data source</li><li>For delimited records, use %varlen SAS macro/program</li><li>Considerations: character vs. numeric, number of decimal places, implied vs. explicit decimals, packed decimals, other non-displayable meta-characters, etc.</li></ul></li><li>Additional things to consider:<ul><li>If the same/similar fields are found in multiple files, should they all be read in with the same format/length?</li><li>How will these fields impact each other across different files?</li><li>What is the implication of reading each field in the format you have chosen?</li><li>Why do I need to load this line/file/variable/etc. this way?</li></ul></li></ul></li><li>Note any issues in data load</li><li>Check for presence of correct merge/unique key in each file<ul><li>Based on the descriptions and the initial looks of the fields included, can tables be merged (theoretically)?</li></ul></li><li>Anything else that are obvious</li><li>Communicate any issues to client via email</li></ul> |

Msight
Analytics · Consulting

| From Receipt[1]: | At the Minimum[2]: |
|---|---|
| First 72 Hours | <ul><li>Complete the first round of data load</li><li>Check for exact duplicates<ul><li>Do not remove them unless why they happen is known and it is confirmed appropriate to remove</li></ul></li><li>Identify unique key(s) that uniquely identify a record<ul><li>The purpose of the exercise is to understand the record level (i.e. what each record represents), as it often is NOT what you understand it to be</li><li>Start with your understanding of which fields/concepts would reasonably represent a record and their hierarchy</li><li>Use common sense. Just because a combination of fields uniquely identifies a record does not mean that the field is an appropriate ID variable.</li></ul></li><li>Quantify the uniqueness of each ID (or ID-like) fields</li><li>Investigate any "duplicates"<ul><li>NOTE: Nowhere does it say "remove duplicates"</li><li>Try to answer the question: what is different among the "duplicates"?</li></ul></li><li>Check for "mergeability"<ul><li>Do merge keys actually merge and how well do they merge?</li><li>This is further to the previous day where we simply answered whether it is theoretically possible)</li></ul></li><li>Note any issues in data load</li><li>Check for general reasonableness of the sample received:<ul><li>Volume</li><li>Target event rate or population average of the target metric, as applicable for the project (preliminary, or some other related metric)</li><li>Subpopulation count</li><li>Etc. This will vary by project.</li></ul></li><li>Communicate any issues to client via email</li></ul> |

| From Receipt[1]: | At the Minimum[2]: |
|---|---|
| First 7 Days | <ul><li>Finalize data load<ul><li>At the minimum, your SAS log should be free of "error", "warning", "not initialized" (check both uppercase and lowercase)—none of these are acceptable.</li></ul></li><li>Full frequency of every raw variable<ul><li>Avoid shortcuts; do not assume a priori grouping is appropriate since it often hides problems. Instead, do a frequency on raw values with limited number of observations until you get a feel for the data, then apply groupings as necessary.</li><li>Work iteratively—this is a critical opportunity to see raw values of every variable with the naked eye.</li></ul></li><li>An "eight-number summary" of every numeric variable for which this makes sense (#Missing, minimum, 1st percentile, 1st quartile, median, 3rd quartile, 99th quartile, maximum)<ul><li>Note any potential default values</li><li>Consider whether any negative values make sense</li></ul></li><li>Further check against the data requirements</li><li>Verify that the values in the distributions make sense</li><li>Verify that definitions are found for all coded values<ul><li>Use the frequency above and label all coded values as they are found in the data dictionary and/or other relevant documents</li><li>Identify which values are missing descriptions, are inconsistent, etc.</li><li>Include all comments in the frequency output for client to evaluate (see below)</li></ul></li><li>Verify that the distribution of key variables are reasonable</li><li>Verify consistency of fields within and across files (for example, similar fields from different data sources, class vs. subclass consistency, etc.) to the extent possible</li><li>Deliver a full report of the initial data check with comments and questions identified (an Excel file is generally acceptable especially at this stage):<ul><li>File names with description and record count for each</li><li>Duplicate and merge check statistics</li><li>Frequency and distributions with variable definition and value descriptions indicated</li><li>All questions about the data</li><li>Comments on anything incomplete, unusual, unexpected, unknown, not understood, not explained, etc.</li><li>Comments on your initial insights about the data—validate your understanding of the data</li><li>Request an ETA for the answers to outstanding questions</li></ul></li><li>Sign-off on the data if no uncorrectable error present; develop plan with timeline if uncorrectable error present or if the correctability is in question (e.g. requires client response to questions)<ul><li>If the final sign-off on the file in question depends on the receipt of other files, then sign off conditionally and develop plan with timeline</li></ul></li></ul> |

| From Receipt[1]: | At the Minimum[2]: |
|---|---|
| Questions to ask yourself prior to Day 7: | • Do I have the description for 100% of the fields?<br>• Do I have all of the definitions?<br>• Do I understand how the files relate to each other?<br>• Do I understand everything in the data? Can I explain everything in the data?<br>• Do the values make sense?<br>• Are there any strange or unexpected values?<br>• Are there any variables that should not be used from the practical perspective?<br>• Am I ready to sign off on the data? (Once you sign off on the data, all fatal issues are your issues, not the client's!) |
| Questions to ask the client at Day 7: | • Any response to comments and observations (not just questions)?<br>• Are there any variables that do not make sense?<br>• Are there any variables that should not be used and why? |