

Data Quality Management

Peer-reviewed author version

VANCAUWENBERGH, Sadia (2019) Data Quality Management. In: Suad Kunosic; Enver Zerem (Ed.). *Scientometrics Recent Advances*, IntechOpen, p. 1-15.

DOI: [10.5772/intechopen.86819](https://doi.org/10.5772/intechopen.86819)

Handle: <http://hdl.handle.net/1942/28871>

Author Queries

[AQ1]	Please provide publisher location for Refs. [7, 8, 19, 20, 22].
[AQ2]	Please provide page range for Ref. [23].
[AQ3]	Please provide volume number and page range for Ref. [24].

02

Data Quality Management

03 *Sadia Vancauwenbergh*04 **Abstract**

05 Data quality is crucial in measuring and analyzing science, technology and
06 innovation adequately, which allows for the proper monitoring of research effi-
07 ciency, productivity and even strategic decision making. In this chapter, the concept
08 of data quality will be defined in terms of the different dimensions that together
09 determine the quality of data. Next, methods will be discussed to measure these
10 dimensions using objective and subjective methods. Specific attention will be paid
11 to the management of data quality through the discussion of critical success fac-
12 tors in operational, managerial and governance processes including training that
13 affect data quality. The chapter will be concluded with a section on data quality
14 improvement, which examines data quality issues and provides roadmaps in order
15 to improve and follow-up on data quality, in order to obtain data that can be used as
16 a reliable source for quantitative and qualitative measurements of research.

17 **Keywords:** data quality, data quality measurement, data quality management,
18 data quality improvement

19 **1. Introduction**

20 Over the past decades, research organizations, administrations and researchers
21 have been collecting data that describe both the input as well as the output side of
22 research. This has resulted in an enormous pile of data on publications, projects, pat-
23 ents, ... researchers and their organizations that are collected within database systems
24 or current research information systems (CRIS). Such data systems are created accord-
25 ing to specific goals and use purposes of individual organizations, which reflects their
26 specific nature and the surrounding context in which they operate. However, over time
27 these data systems, institutions as well as the research ecosystem at large have evolved,
28 thereby potentially threatening the quality of the collected data and the resulting
29 data analyses, particularly if no formal data quality management policy is being
30 implemented. This chapter introduces the readers into the concept of data quality and
31 provides methods to assess and improve data quality, in order to obtain data that can be
32 used as a reliable source for quantitative and qualitative measurements of research.

33 **2. Definition of data quality**

34 In general, data can be considered of high quality if the data is fit to serve a
35 purpose in a given context, for example, in operations, decision making and/or
36 planning [1]. Although this definition of data quality seems to be straightforward,
37 many other definitions exist that differ in terms of the qualitative or quantitative
38 approach towards defining the concept of data quality.

01 **2.1 Qualitative approach**

02 In the qualitative approach, specific attention is drawn to defining data quality
 03 in terms of the different aspects, also termed dimensions. In 1996, Wang and Strong
 04 developed a data quality framework based on a two-stage survey on data quality
 05 aspects important to data consumers, and captured these dimensions in a hierarchi-
 06 cal manner [2]. This model clusters 20 different data quality dimensions into four
 07 major categories: that is, intrinsic, contextual, representational and access data
 08 quality. Although the basis of this model still stands, some minor changes have been
 09 made over the years resulting in the model depicted in **Table 1** [3].

10 In brief, the *intrinsic* category comprises dimensions that define the accuracy of
 11 the data, that is, the extent to which data is certified, error-free, and reliable, as well
 12 as the objectivity of the data based on facts and impartial, and their reputation based
 13 on its sources or content. The *contextual* data quality category comprises dimensions
 14 that must be considered within the context of a specific objective for which one holds
 15 the data, that is, the data should be relevant, up to date, of an appropriate amount,
 16 yet complete, and ready for use for the stated objective. The *representational* category
 17 contains dimensions that reflect how the data are presented within a data system.
 18 Dimensions concerning the format of the data, that is, concise and consistent
 19 representation, as well as their compatibility, their interpretability and whether they
 20 are easy to understand, are considered. The last category is focused on the *accessibil-*
 21 *ity* category that also defines aspects of data quality. Although this category is not
 22 always considered in the literature [4], this is an important aspect of overall data
 23 quality. The related dimensions include the accessibility of the data in terms of their
 24 availability or easily retrievable character, the security measures taken to restrict data
 25 appropriately and the traceability of the data to its source.

Category	DQ dimension
Intrinsic	Accuracy
	Objectivity
	Reputation
Contextual	Completeness
	Appropriate amount
	Value added
	Relevance
	Timeliness
	Actionable
Representational	Interpretable
	Easily understandable
	Consistent
	Concisely represented
	Alignment
Access	Accessibility
	Security
	Traceability

Table 1.
Data quality dimensions.

26

01 These dimensions can also be grouped into an internal and external group of
02 dimensions. The internal group contains the dimensions that can be measured purely
03 in terms of the data, and are generally more objective. Examples of these include the
04 accuracy of the data, which can be examined by calculating a score on the magnitude
05 of errors in the data or the data correctness, which can be measured through the
06 number of errors in the data. On the other hand, the external group of dimensions
07 evaluates how the data are related to their environment, and hence are somewhat
08 more subjective in nature. Examples include the relevancy of data with regards to a
09 stated objective, or their ease of understanding by the consumers of the data.

10 **2.2 Quantitative approach**

11 In the quantitative approach, data quality has been defined by J. M. Juran as the
12 fitness of the data to serve a purpose in a given context, that is, in operations, deci-
13 sion making and/or planning as perceived by its users [1]. This concept is denoted
14 as ‘fitness for use’ and is based on Juran’s five principles: that is, who uses the data,
15 how are the data used, is there a danger for human safety, what are the economic
16 resources of the producers and users of the data and what are the characteristics
17 taken into account by users when determining the fitness for use. This definition
18 is widely accepted in both academic and industrial settings. However, in practice
19 the fitness for use is a rather subjective measure as this highly depends on the users’
20 judgement over the degree of conformity of the data to their intended use.

21 For example, consider the score of a student on an exam. If scores are rounded
22 to integers, this can potentially influence the final grade that a student receives.
23 Therefore, the rounding procedure might be accurate enough for the professors, but
24 by rounding numbers, the students might miss out on obtaining a final grade and
25 thus might be not accurate enough from the perspective of the student.

26 On the other hand, it might well be that not all uses of the data are known,
27 neither its potential future use purposes. Therefore, DQ might be hard to evaluate
28 using this definition.

29 Some definitions of data quality use the notion of zero defects, which aims to
30 reduce defects by motivating people to prevent making mistakes by developing a
31 constant, conscious desire to do the job right from the first time [5]. This zero-defect
32 concept has been incorporated by P. Crosby in its *Absolutes of Quality Management*
33 [6]. According to Crosby’s *Absolutes*, data quality should conform to its requirements
34 and prevention should be used as a manner to guarantee zero defects, which sets the
35 performance standard. Consequently, data quality can be measured as the price of
36 nonconformance. Although this zero-defect concept is not widely used in the data
37 quality literature, it does emphasize again the necessity to measure data quality.

38 **3. Measuring data quality**

39 Based on the definitions of data quality, several DQ measurement methods have
40 been developed, that can generally be divided into objective and subjective meth-
41 ods. While objective methods tend to evaluate data quality rather from the perspec-
42 tive of the data producer based on hard criteria, subjective methods rather take the
43 user’s perspectives and beliefs into account.

44 **3.1 Objective DQ measurement methods**

45 Measurements of data quality are generally intended to assess the dimensions
46 of data quality as defined in the previous section. As a first step, a framework must

01 be set up with the indicators that one wants to assess. Next, a proper reference for
02 verification of the data within the data systems must be determined.

03 Ideally, the data are compared using real world data, which allows for validation
04 and, if required immediate corrective actions. This method is termed *data audit-*
05 *ing* and is the only way of measuring the quality level of dimensions like accuracy,
06 completeness. Furthermore, by going through the data itself, one can discover
07 data quality issues that were unexpected and therefore are of great value for taking
08 corrective measures to improve data quality. However, data auditing comes at a high
09 cost as it is very time consuming and the need of experts in the respective field is
10 required. Furthermore, data auditing can be also very labor-intensive and requires
11 that data controllers have access to the actual data.

12 For example, consider the metadata of publications that are contained in pub-
13 lication databases. If a data controller validates the content of the metadata fields
14 with the metadata as indicated on the publications, inaccuracies can be detected.
15 These can contain expected flaws like spelling errors but can also provide valuable
16 information on unexpected errors that also might be highly relevant in the context
17 of bibliometric analyses.

18 If the conditions for data auditing are not met, data controllers can *use rule-*
19 *based checking* in order to determine data quality. This method heavily relies on
20 business rules that are drafted based upon the domain knowledge and experience
21 that the data controllers have with regards to the data. Consequently, these rules
22 can only check for flaws that were anticipated by the data controllers. However,
23 rule-based checking also offers important advantages, especially as they can be
24 automated after conversion to validation rules, which allows for the identifica-
25 tion of the errors (or possibly correct outliers!) via data mining techniques.
26 Nevertheless, the presumed errors still need to be corrected, which remains
27 labor-intensive.

28 3.2 Subjective DQ measurement methods

29 Some dimensions, however, cannot be measured objectively because of their
30 intrinsic properties. For example, the dimension relevancy pertains to the extent to
31 which data is applicable and helpful for the stated objective. Obviously, this dimen-
32 sion can only be evaluated using the *perception of the users*. Although this results
33 in a subjective scoring, user evaluations are the only way to measure dimensions
34 that describe external data quality attributes. Internal data quality dimensions on
35 the contrary are preferably measured using objective DQ measurement methods as
36 described above.

37 Regardless of which methodology is chosen to measure data quality, it is always
38 important to provide information about the measurement method and parameters
39 in addition to the dimension under evaluation, in order that the measurement
40 results can be interpreted correctly by everyone. Furthermore, although a lot of
41 attention always goes to correcting errors, it is important to stress that eliminating
42 the root cause should always be the ultimate goal [7].

43 4. Data quality management

44 4.1 Data quality frameworks

45 As data are extremely valuable resources in today's society, a plethora of data
46 quality management frameworks have been published in the last decades that all
47 strive to preserve the quality of data and to make it accessible for future use. The

01 most popular models are listed below, however more DQM frameworks can be
02 found throughout the literature that show slight differences.

- 03 • DAMA DMBOK's Data governance model [8]
- 04 • EWSolutions' EIM Maturity Model [9]
- 05 • Oracle's Data Quality Management Process [10]

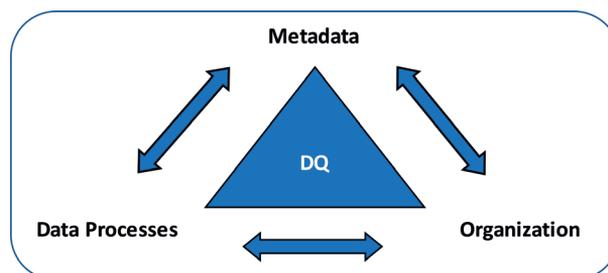
06 All frameworks are basically centered around three basic elements, that is,
07 the metadata associated with the data, the processes involved in the registration,
08 organization and (re)use of the data, and the organizational context in relation to
09 the data (**Figure 1**). The quality of each individual element, as well as the interplay
10 in between them, ultimately determines the quality and thus the true value of an
11 organization's data heritage. Ideally, an organization uses metadata standards that
12 are understandable throughout the organization and aligned with the organiza-
13 tion's processes, business strategies and goals. Rather than describing all popular
14 frameworks, we will describe critical success factors that are useful for developing
15 effective DQ management strategies, and that can be found in all DQ frameworks.

16 4.2 Critical success factors

17 Critical success factors, also termed CSFs, have been defined by Milosevic and
18 Patanakul as '*characteristics, conditions, or variables that can have a significant impact*
19 *on the success of i.e., a company or a project when properly sustained, maintained, or*
20 *managed*' [11]. In 2014, Baskarada described 11 CSFs in the field of information
21 quality management that provide valuable means for developing effective DQ
22 management strategies [12]. These CSFs can be clustered into four major groups,
23 that is, training, governance, management and operational processes, that have
24 inter-dependencies with each other.

25 4.2.1 Operational processes

26 The first group of critical success factors deals with the operational processes
27 involved in the collection, storage, analysis and security of the data, which are all
28 highly interdependent. As data is a valuable good, its quality should be managed
29 throughout its entire lifecycle. In practice this comes down to taking measures that
30 maximize, whenever possible, the **automated capture** of data in **real-time**, directly
31 from its **original source**. This minimizes the risk of errors introduced by manual
32 data entry, which can result in typo's, inaccuracies, missing values, erroneous data



33 **Figure 1.**
The cornerstone of data quality frameworks.

01 due to misinterpretations, multiple copies of the same data entry. Such errors have
02 been identified in almost all existing research and innovation databases, but have
03 a significant impact on the resulting scientometric analyses. Suppose a highly cited
04 paper is included in the Web of Science with typos in the author's name. This can
05 erroneously lead to the omission of this paper in the bibliometric analyses per-
06 formed on this author, which on its turn can have a major impact on this researcher
07 career perspectives in terms of chances of success in obtaining grants, promotion.

08 In addition, these errors can be due to a lack of the use of common **standards**
09 for the concepts contained within the databases and a uniform interpretation
10 thereof by both information providers as well consumers throughout the entire
11 organization. Nevertheless, such standards are available, that is, the Common
12 European Research Information Format (CERIF) is a well-known standard for
13 exchanging research information created by the EuroCRIS organization and is
14 widely used throughout Europe [13], the CASRAI dictionary is a standard created
15 by the organization on Consortia Advancing Standards in Research Administration
16 Information (CASRAI) and was created in Canada [14]. Although both communi-
17 ties work closely together to align the concepts and meanings described in the
18 standards, some differences remain which might cause difficulties in exchanging
19 information in between CRIS systems. Furthermore, the inclusion of a standard in
20 the information model of a data system does not safeguard that all data providers
21 use the standard similarly, nor that the data users grasp the information as intended.
22 Next to using standards for aligning the concepts and meanings of research-related
23 data, the formats of the data fields should be standardized as well. A well-known
24 example here includes the various formats in which a (publication) date is recorded.
25 By means of standardizing this format in a data system, important gains can be
26 obtained in terms of ease of interpretation of the data, leading to more accurate
27 analyses. However as described above, efforts should also be made to clarify what
28 the concept of (publication) date means. For instance, it could point to the creation
29 date, submission date, the published online date, the publication date for in print
30 papers, the date on which the material was made available.

31 Furthermore, when storing research-related data, it is highly recommended to
32 provide **traceability** to the raw data, which ensures that the data quality can always
33 be controlled. Most bibliometric databases, including the Web of Science and
34 Scopus, comply to this rule by providing a link to the journal article. Research data
35 repositories mostly refer to the creator of the datasets involved. However, over time,
36 researchers can switch positions and thus institutions and as the data are stored
37 in institutional repositories, it would be more meaningful to refer to the research
38 institution in question. In addition, **versioning** should be included when storing
39 research data, as this can be very helpful to understand and potentially (re)use data.
40 Although this is frequently observed in research data repositories, bibliometric and
41 patent databases usually do not show version control. Finally, **back-up** and **data**
42 **recovery** processes should be ensured when storing research-related data, which is
43 mostly realized via back-up servers at various physical places.

44 The access to research information should be managed using an **information**
45 **security management** plan in order to safeguard the intellectual property rights
46 of the researchers that created the information, including their respective institu-
47 tions. Although large data repositories on bibliometric, innovation and research
48 data control accessibility rights, researchers themselves do not always closely follow
49 the measures taken to control access. Particularly when it comes down to research
50 data that may contain sensitive data [15], strict follow-up of information security
51 measures is needed as emphasized by the EU Regulation 2016/679, also known as
52 the General Data Protection Regulation (GDPR) that protects natural persons with
53 regards to the processing of personal data and on the free movement of such data.

01 Although the GDPR regulation only applies to personal data *in se*, it nicely under-
02 pins some elements present in information security management plans.

03 These information security management plans indeed not only entail the acces-
04 sibility rights of individuals, including user authentication and a regular update of
05 their access rights, but also include the secure storage, archival, transmission, and
06 if required, destruction of the information. In case of research data on natural per-
07 sons, this can be achieved via pseudonymization, for example, through encryption,
08 or via anonymization of the research information residing in data systems or on
09 data carriers. Obviously, when transmitting research information, the proper legal
10 agreements should be put in place, for example, non-disclosure agreements with
11 third parties are well-known examples used to secure research information. Finally,
12 information security management plans should also contain audit trails in order to
13 constantly monitor and adjust the security of research-related information.

14 4.2.2 Management processes

15 A second group of CSFs encompasses the managerial processes that are imposed
16 on these operational processes, and which are primarily aimed at the alignment of
17 the data quality with the organization's goals with regards to the data and the result-
18 ing data analyses. Consider for example, the information requirement of a univer-
19 sity that wants to monitor the research funds obtained via researchers. In order to
20 answer this question, the concepts of research funds and researchers should be clear
21 and uniform between information providers and users. Although this might seem
22 straightforward, it could well be that the interpretation of 'researcher' is different in
23 between stakeholders, that is, while some might include PhD students, other might
24 omit this group. Furthermore, it could well be that the university does not have a
25 specific label for clustering funds as belonging to the 'research' category, or that the
26 information is only partially provided by the researchers. These examples clearly
27 illustrate that the lack of management of operational DQ processes, has a devastat-
28 ing effect on the data analyses and the conclusions based thereon.

29 Managerial processes of data quality essentially focus on four sequential pro-
30 cesses, that is, the determination of the information quality requirements, the
31 assessment of the risks associated with DQ issues, the assessment or monitoring of
32 DQ and the continuous improvement of the related DQ processes [16]. First, the
33 **information quality requirements** should be determined of the collected data,
34 considering all stakeholders. Next, a conceptual information model should be
35 drafted using high-level data constructs, generally described in non-technical terms
36 in order to be understandable by executives and managers. This model should then
37 be translated into a logical data model that uses entities, attributes and relationships
38 that are customized towards the organization's use of the data, in terms of the orga-
39 nization's terminology, semantics as well as the prevailing business rules. Finally,
40 the logical model should be transferred to developers that can derive a physical data
41 model in line with this logical model including validation rules, based upon the
42 business rules, that are useful for automating data quality control. Obviously, the
43 constructed models must consider the importance of the data within the organiza-
44 tion. For example, certain data will be more important than others, and poor DQ
45 of those data might have a larger negative impact in terms of loss of reputation,
46 financial loss. of the organization. The explicit **management of these DQ risks** is a
47 must as a manner to guarantee data quality. As stated by Baskarada '*using gut feeling*
48 *will result in inefficiency and an ineffective use of resources*' [16].

49 Next, a framework of key performance DQ indicators needs to be set up in line
50 with the organization's goals, in order to **assess the DQ performance**. This assess-
51 ment must be performed on a regular basis in order to allow for the **continuous**

01 **improvement of data quality** in terms of analyzing the root cause of the errors as
02 well as cleansing erroneous data.

03 The application of such DQ managerial processes has already been implemented
04 to some extent in CRIS systems that contain research information. For example,
05 the Flanders Research Information Space, also termed FRIS, is a research informa-
06 tion portal sustained by the Department of Economy, Science and Innovation in
07 Flanders, Belgium that collects research information from a wide range of Flemish
08 stakeholders in the research field, that is, research universities, higher education
09 colleges, strategic research centers and research institutions (www.researchportal.be) [17]. Underlying the FRIS architecture, a conceptual metamodel was developed
10 in order to model all concepts, attributes and relationships that are contained within
11 FRIS. This conceptual model is based on the CERIF standard, but customized to
12 the Flemish context. In addition, in line with the use purposes of this CRIS system,
13 business rules were drafted to safeguard the quality of the contained information.
14 These business rules were translated to validation rules that are used for the auto-
15 mated quality control of the research information received. If non-compliances to
16 these rules are detected, the research information is rejected, and the information
17 providers receive a notification thereby allowing for immediate data cleansing.
18 Furthermore, the Flemish government also performs manual quality checks on a
19 regular basis in order to validate the research information contained as validation
20 rules in general are not well suited for detecting unpredicted errors. Such errors
21 generally provide valuable input for root cause analyses that can identify important
22 underlying problems which can be caused by human, process, organizational or
23 technological factors.
24

25 4.2.3 Governance process

26 A third group of CSFs encompasses the governance processes associated with
27 DQ management. These processes can be largely summarized as the **commitment**
28 **of an organization's top management** to set DQ management as a priority and to
29 stimulate a culture change throughout the entire organization in this respect. In the
30 field of information governance, Gartner Research defined information governance
31 as *'the specification of decision rights and an accountability framework to encourage*
32 *desirable behavior in the valuation, creation, storage, use, archival and deletion of*
33 *information'* [18]. In practice, information governance basically comes down to
34 allocating budget and resources to the process of DQ management by defining roles
35 and responsibilities, making agreements on related concepts, terms and associated
36 DQ processes, including the monitoring, control and improvement thereof. The
37 FRIS-system as indicated above has included data governance in order to ensure
38 proper DQ management [17].

39 4.2.4 Training

40 Although an organization might have all operational, managerial and gover-
41 nance processes perfectly in place, a complete implementation of DQ management
42 also requires the investment in training throughout the organization. A first and
43 foremost important goal is to inform people on the importance of qualitative data
44 to the organization. Secondly, people should receive training via training programs,
45 course series, mentorships. on the rules as set out in the operational, managerial
46 and governance processes in order to ensure a systematic implementation of DQ
47 throughout the entire organization. Finally, a continuous follow-up is also needed
48 which allows for swift adjustments in case of unpredicted errors, adjustment of
49 business rules, etc.

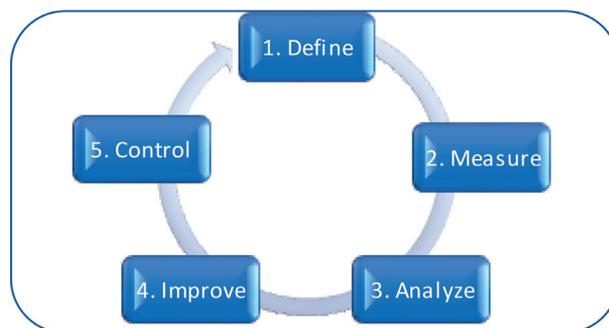
01 5. Data quality improvement

02 In order to safeguard the continuous monitoring of data quality and the adop-
03 tion of measures to improve data quality, a DQ improvement workflow needs to
04 be established. This workflow essentially comprises a repetitive workflow of five
05 consecutive phases, that is, the definition, measurement, analyze, improvement
06 and control phase as depicted in **Figure 2**. A best practice is to formalize this
07 data quality improvement process, in terms of properly documenting all related
08 processes and activities in each phase, as this allows for the tracking of progress
09 throughout the entire DQ improvement workflow.

10 5.1 Definition of DQ project

11 The DQ improvement workflow starts with defining the scope of the DQ
12 improvement project. This includes the selection of a dataset relevant to a specific
13 business goal, and the determination of the data attributes required. When collect-
14 ing this information, it is very important to discuss the meaning of the metadata
15 required with all stakeholders in order to be able to identify any discrepancies in
16 interpretation of the required data attributes versus the meaning of the existing
17 metadata, as this prevents erroneous data collection, analysis and interpretation.
18 All obtained information should be documented using domain modeling techniques
19 that include information on the data and the associated operations on the data [19].
20 Examples of such techniques include Business Process Model Notation (BPMN)
21 diagrams [20], data flow diagrams, of which the resulting information should be
22 contained in data governance tools together with the accompanying semantics. In
23 addition, data quality dimensions important to the specific use purposes of the data
24 should be determined, and if possible, these are preferably defined in a measurable
25 manner which facilitates further steps in the DQ improvement process.

26 For example, consider the use of bibliometric data as part of a researcher's evalu-
27 ation in the context of career-wise promotions. In order to provide an adequate,
28 qualitative data-analysis, a clear framework should be defined by an organization's
29 management comprising what should be evaluated, that is, which publications
30 (books, journals.), validation criteria (peer reviewed, group author.) are to be
31 used as well as the accompanying processes. This information should be discussed
32 with all stakeholders, that is, researchers, librarians, data analysts and IT-staff in
33 order to harmonize the data flow, the accompanying semantics, procedures and
34 models in accordance with the management's goals. Next, the *As Is* situation should
35 be evaluated with regards to these intentions and according to the relevant data
36 dimensions. In bibliometric analyses, accuracy, completeness, timeliness, relevance,



37 **Figure 2.**
38 *Data quality improvement workflow.*

01 accessibility, traceability of the data are all relevant dimensions, of which the
 02 accurate and complete collection and analysis of a researcher's published works are
 03 the foremost ones.

04 **5.2 Measurement DQ**

05 In order to determine the quality level of the current data in relation to the
 06 organization's objectives, the quality dimensions need to be expressed in a measur-
 07 able manner. While the internal dimensions can be scored in a quantitative manner
 08 by means of expressing the errors in the data set in terms of magnitude, number of
 09 errors or missing records., the external dimensions are measured in a qualitative
 10 manner based on the context of the data's use purposes. Independent of the dimen-
 11 sion under analysis, measurements must always be relevant for the purpose for
 12 which the data will be used and according to the task's requirements. Although in
 13 most cases, common sense will be used to identify task requirements, in other cases
 14 specific techniques like sensitivity analysis might be used which allows for identify-
 15 ing critical factors and errors in data models [21, 22]. Furthermore, data profiling
 16 is another technique frequently used in DQ assessment as a method to discover the
 17 true content, structure and quality of data by means of rule-based checking [23].
 18 Obviously, this technique does not find all inaccurate data, as it can only identify
 19 violations to the predefined rules, and hence expected errors. For instance, data
 20 profiling can identify invalid data values (i.e., using column property analysis),
 21 invalid data combinations (i.e., through structure analysis), inaccurate data (i.e.,
 22 through value rule analysis). Importantly, data profiling also provides metrics on
 23 the data inaccuracies in a dataset, that is, the number of violations, the frequency
 24 of invalid data values, etc. Such metrics can be useful as a means to communicate to
 25 stakeholders on the (in)accuracy of a data set, and the follow-up of the progression
 26 in subsequent DQ improvement programs.

27 In our bibliometric example, the accuracy and completeness of the bibliometric
 28 records for a given author, collected in a university's database system should be
 29 compared to a publication list provided by the author. By manually auditing the
 30 registered data found within the database system, one could indeed record the com-
 31 pleteness of information. Furthermore, the accuracy can be tested using a manual
 32 auditing procedure. This allows for the identification of spelling errors, erroneous
 33 exchange of an author's last versus first name, etc. In addition, manual auditing
 34 also allows for identification of rather unexpected data entries, like changes in the
 35 author's first or last name over time. The latter example of a DQ inaccuracy, can
 36 however not be detected through data profiling as rule-based checking is unable to
 37 test for unexpected errors. Nevertheless, data profiling has an important role in DQ
 38 measurement as it allows for automated and thus efficient screening of DQ.

39 **5.3 Analyzing DQ issues**

40 Once DQ inaccuracies have been detected, these should be analyzed in order to
 41 screen for the potential existence of (groups of) common underlying root causes.
 42 For example, author names can have various problems like misspelling, last names
 43 mistaken for first names, etc. The grouping of such errors that show similar pat-
 44 terns, also called error cluster analysis, allows for the identification of common
 45 causes and is often more efficient in terms of time and resources as compared to
 46 handling all inaccuracies in a stand-alone way. In addition, a data event analysis can
 47 be performed which evaluates the time points when data are created and updated in
 48 order to facilitate the identification of the root causes of problems. For example, the
 49 manual entry of author names in a database system might result in misspelling, the

01 lack of automated verification in the recording process, the lack of domain specific
02 knowledge of the persons responsible for recording the data, ... might affect the
03 occurrence of DQ inaccuracies.

04 Commonly used techniques to identify root causes include the auditing of the
05 data, the surveying of the user perceptions and the evaluation of the data process.
06 The identified causes can then be depicted in cause and effect diagrams, also termed
07 Ishikawa or fishbone diagrams [24]. These diagrams cluster causes together in
08 groups which is instrumental in identifying, classifying and prioritizing the impact
09 of root causes to a problem. In our example root cause analysis could result in the
10 identification of the field 'author name', as a string datatype, that is, completed
11 according to the data provider's interpretation and accuracy. Because the datatype is
12 set as a string, multiple inaccuracies can occur during the registration process.

13 **5.4 DQ improvement trajectories**

14 In the next phase, the focus resides on finding solutions to eliminate the root
15 cause of the problem. These solutions, also termed remedies, are in fact changes
16 to data systems or processes in order to prevent data inaccuracies from happen-
17 ing including the swift detection upon their occurrence. While some solutions
18 might be oriented towards improving the data registration, others might focus
19 on the implementation of validation rules or periodic data profiling. In addition,
20 re-engineering of associated data processes and even training of the data provider
21 and user community on data quality aspects, should be considered. Data cleansing
22 might be applied as well, however this mostly is not a solution to eliminate the root
23 cause itself.

24 Although solutions might be found using common sense, in most cases more
25 efforts are needed. A frequently used method encompasses the organization
26 of topic-oriented brainstorm sessions in the presence of all stakeholders. This
27 approach has the benefit to tackle the problem for multiple viewpoints and at the
28 same time enables a higher engagement of the stakeholders. Importantly, all rel-
29 evant solutions to the problem should be listed and effects of the proposed solutions
30 should be investigated carefully. In general, continuous, short-term improvements
31 are to be preferred as these might result in quick wins which can result in additional
32 business benefits (as DQ improvement is mostly not a goal in itself).

33 In our example many solutions can be found that focus on improving the correct
34 registration of the author name. However, if an author ID would be registered and
35 coupled to an author name, the specific focus on registering the name perfectly
36 in a wide variety of bibliometric sources diminishes. Although this seems an easy
37 solution at first glance, this strategy also includes the re-engineering of business
38 processes, that is, the authentication of research publications by an author using
39 its author ID. In order to investigate the effect of this proposed solution, one could
40 investigate the number of publications that can be attributed to a group of authors
41 that has registered and authenticated their research publications versus a group of
42 authors that have no author ID (i.e., the control group) in an experimental setting.
43 By measuring the DQ of both groups in terms of accuracy and completeness, one
44 can see the effect of the proposed solution.

45 **5.5 DQ control and follow-up**

46 Based on all DQ solutions tested, the most appropriate solution(s) should be
47 selected for implementation. It is important to note here that the success of imple-
48 mentation is dependent on the guidance foreseen to all stakeholders. In essence,
49 this comes down to providing information on the solution and its effectuation on

01 all (related) business processes to everybody involved. In addition, business rules,
 02 definitions, roles and responsibilities must be defined in consultation with all
 03 stakeholders.

04 Obviously, a close monitoring is needed in order to follow-up on the effective-
 05 ness of the implemented DQ solution in the real-world setting as a means to validate
 06 the (positive) impact of the proposed DQ solution. At the same time, it allows for
 07 the detection of unexpected errors that were unanticipated in the experimental test
 08 phase, and the swift adoption of corrective measure in case required. Specific moni-
 09 toring tools that can be used here include control charts, also known as Shewhart
 10 charts, cause and effect diagrams, check sheets, histograms, Pareto charts, scatter
 11 diagrams, ... [25].

12 With regards to the author disambiguation example described, it will be
 13 required to install business process that allow for the coupling of a unique author
 14 ID with corresponding research publications. This includes the close cooperation
 15 of the authors, research administrators, data analysts and data system/IT-staff on
 16 the definitions, business rules and responsibilities of each stakeholder. For instance,
 17 it might well be that authors are obliged to enter a unique author ID in a database
 18 system in fixed format, rather than a free text field. A business rule could be that for
 19 each author, an author ID of a given type (i.e., ORCID, Researcher ID, Scopus ID,
 20 Research Gate ID.) should be kept in a data system, which translates to a value of a
 21 given format, that is, an integer, in terms of a derived validation rule. This author ID
 22 field might be used to search large bibliometric databases such as Web of Science,
 23 Scopus, ... for publications that might be coupled to this author ID, which could be
 24 added to the bibliometric profile of a researcher. Furthermore, publications might
 25 also be retrieved using an author name search that are not yet coupled to this author
 26 ID. Therefore, an authentication step is required here in which the author has a
 27 critical responsibility to validate these publications. Research administrators and
 28 data analysts should be informed on the process of authentication in order to use
 29 the information in a correct manner. Although this might seem a perfect solution,
 30 the reality demonstrates that a continuous follow-up is required as practice demon-
 31 strates that authors sometimes use several author IDs of the same type. Therefore,
 32 a corrective action could be to adapt the business rules in order to allow for only
 33 one author ID of a give type within the data system as well as the notification to the
 34 author to take corrective measures in this respect and the follow-up thereof.

35 It is clear from the example described above, that data quality improvement is
 36 a process that requires continuous monitoring due to internal and external fac-
 37 tors that might effectuate data quality and its related processes. Therefore, the
 38 systematic and continuous retaking of the DQ improvement workflow will be the
 39 only manner to constantly have qualitative data instrumental for high quality data
 40 analyses.

41 **6. Conclusion**

42 Research organizations worldwide are using data on research input and out-
 43 put, that is, publications, patents, research data nowadays for a wide variety of
 44 use purposes, such as evaluation, reporting and visualization of a researcher' or
 45 research organization's expertise. This places high demands on the quality of the
 46 data gathered for these purposes, which have—in most cases—largely outgrown the
 47 initial intentions when the data systems were constructed. Moreover, the research
 48 world has evolved in a global, dynamic in which research data are increasingly
 49 being used in order to monitor the efficiency of research processes, the research
 50 productivity and even strategic decision making. In order to safeguard correct data

05 analysis, research-related data must be assessed on all relevant quality dimensions,
06 and inaccuracies must be addressed using data quality improvement trajectories as
07 discussed in this chapter. The integration of a data quality management policy, is
08 the only way to ensure the fitness for use of research-related data for various appli-
09 cations and business processes across the research world as the impact of inaccurate
10 data can have tremendous effects on a researcher's or research organization's future
11 prospects.

12 **Acknowledgements**

13 This work is carried out for the Expertise Centre for Research and Development
14 Monitoring (ECOOM) in Flanders, which is supported by the Department of
15 Economy, Science and Innovation, Flanders.

16 **A. Abbreviations**

17	BPMN	Business Process Model Notation
18	CASRAI	Consortia Advancing Standards in Research Administration 19 Information
20	CERIF	Common European Research Information Format
21	CRIS	current research information systems
22	FRIS	Flanders Research Information Space
23	DQ	data quality
24	DQM	data quality management

01 **Author details**

02 **Sadia Vancauwenbergh**
03 **ECOOM-Hasselt and Hasselt University, Hasselt, Belgium**

04 ***Address all correspondence to: vancauwenbergh@uhasselt.be**

IntechOpen

© 2019 The Author(s). Licensee IntechOpen. This chapter is distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/3.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited. 

References

01 [1] Juran JM, Blanton Godfrey A. Juran’s 44
 02 Quality Handbook. 7th ed. Europe: 45
 03 McGraw-Hill Education; 2016. p. 992.
 04 ISBN-10: 9781259643613

05 [2] Wang RY, Strong DM. Beyond 46
 06 accuracy: What data quality means to 47
 07 data consumers. Journal of Management 48
 08 Information Systems. 1996;12(4):5-33. 49
 09 DOI: 10.1080/07421222.1996.11518099 50
 51

10 [3] Moges H-T, Dejaeger K, Lemahieu W, 52
 11 Baesens B. A total data quality 53
 12 management for credit risk: New 54
 13 insights and challenges. International 55
 14 Journal of Information Quality. 56
 15 2012;3(1):1-27. DOI: 10.1504/ 57
 16 IJIQ.2012.050036

17 [4] Culnan M. The dimensions of 58
 18 accessibility to online information: 59
 19 Implications for implementing office 60
 20 information systems. ACM Transactions 61
 21 on Office Information Systems. 62
 22 1984;2(2):141-150. DOI: 10.1145/521.523 63

23 [5] Halpin JF. Zero Defects: A New 64
 24 Dimension in Quality Assurance. 65
 25 New York City: McGraw-Hill; 1966. 66
 26 p. 228. OCLC 567983091 67

27 [6] Crosby PB. 8: Quality Improvement 68
 28 Program. Quality Is Free: The Art of 69
 29 Making Quality Certain. New York City: 70
 30 McGraw-Hill; 1979. pp. 127-139. ISBN 71
 31 9780070145122. OCLC 3843884

AQ1 32 [7] Redman TC. Data Quality: The Field 72
 33 Guide. Digital Press; 2001. p. 256. ISBN- 73
 34 10 1555582516 74
 75
 76
 77
 78

35 [8] DAMA International. The DAMA 79
 36 Guide to the Data Management Body of 80
 37 Knowledge (DAMA-DMBOK). Technics 81
 38 Publications, LLC; 2009. p. 430. ISBN- 82
 39 10: 0977140083 83

40 [9] EWSolutions, Foundations 84
 41 of Enterprise Data Management 85
 42 [Internet]. 2013. Available from: 86
 43 [https://www.ewsolutions.com/](https://www.ewsolutions.com/foundations-enterprise-data-management/)

[10] Oracle, Oracle Warehouse Builder 46
 Users Guide 10g Release 2(10.2.0.2) 47
 [Internet]. 2009. Available from: [https://](https://docs.oracle.com/cd/B31080_01/doc/owb.102/b28223/concept_data_quality.htm) 48
[docs.oracle.com/cd/B31080_01/doc/](https://docs.oracle.com/cd/B31080_01/doc/owb.102/b28223/concept_data_quality.htm) 49
[owb.102/b28223/concept_data_quality.](https://docs.oracle.com/cd/B31080_01/doc/owb.102/b28223/concept_data_quality.htm) 50
[htm](https://docs.oracle.com/cd/B31080_01/doc/owb.102/b28223/concept_data_quality.htm) [Accessed: 18 April 2019] 51

[11] Milosevic D, Patanakul P. 52
 Standardized project management 53
 may increase development projects 54
 success. International Journal of Project 55
 Management. 2005;23:181-192. DOI: 56
 10.1016/j.ijproman.2004.11.002 57

[12] Baskarada S, Koronios A. A 58
 critical success factor framework for 59
 information quality management. 60
 Information Systems Management. 61
 2014;31(4):276-295. DOI: 62
 10.1080/10580530.2014.958023 63

[13] EuroCRIS, CERIF: Main features of 64
 CERIF. Available from: [https://www.](https://www.eurocris.org/cerif/main-features-cerif) 65
[eurocris.org/cerif/main-features-cerif](https://www.eurocris.org/cerif/main-features-cerif) 66
 [Accessed: 18 April 2019] 67

[14] CASRAI: CASRAI dictionary. 68
 Available from: [https://dictionary.casrai.](https://dictionary.casrai.org/Main_Page) 69
[org/Main_Page](https://dictionary.casrai.org/Main_Page) [Accessed: 19 April 70
 2019] 71

[15] European Commission: Sensitive 72
 data. Available from: [https://](https://ec.europa.eu/info/law/law-topic/data-protection/reform/rules-business-and-organisations/legal-grounds-processing-data/sensitive-data_en) 73
[ec.europa.eu/info/law/law-topic/](https://ec.europa.eu/info/law/law-topic/data-protection/reform/rules-business-and-organisations/legal-grounds-processing-data/sensitive-data_en) 74
[data-protection/reform/rules-business-](https://ec.europa.eu/info/law/law-topic/data-protection/reform/rules-business-and-organisations/legal-grounds-processing-data/sensitive-data_en) 75
[and-organisations/legal-grounds-](https://ec.europa.eu/info/law/law-topic/data-protection/reform/rules-business-and-organisations/legal-grounds-processing-data/sensitive-data_en) 76
[processing-data/sensitive-data_en](https://ec.europa.eu/info/law/law-topic/data-protection/reform/rules-business-and-organisations/legal-grounds-processing-data/sensitive-data_en) 77
 [Accessed: 19 April 2019] 78

[16] European Commission: Data 79
 protection in the EU. Available from: 80
[https://ec.europa.eu/info/law/law-topic/](https://ec.europa.eu/info/law/law-topic/data-protection/data-protection-eu_en) 81
[data-protection/data-protection-eu_en](https://ec.europa.eu/info/law/law-topic/data-protection/data-protection-eu_en) 82
 [Accessed: 19 April 2019] 83

[17] Vancauwenbergh S, De Leenheer P, 84
 Van Grootel G. On research information 85
 and classification governance in an 86

- 01 inter-organizational context: The
02 Flanders Research Information Space.
03 Scientometrics. 2016;**108**(1):425-439.
04 DOI: 10.1007/s11192-016-1912-7
- 05 [18] Bell T, Logan D, Friedman T. Key
06 issues for establishin information
07 governance policies, processes and
08 organization. Gartner Research; 2008
- 09 [19] Fowler M. Patterns of Enterprise
10 Application Architecture. Addison
11 Wesley; 2003. p. 116. ISBN-10:
12 0321127420
- 13 [20] White SA. Process Modeling
14 Notations and Workflow Patterns. IBM
15 Corporation; 2006
- 16 [21] Saltelli A. Sensitivity analysis
17 for importance assessment. Risk
18 Analysis. 2002;**22**(3):1-12. DOI:
19 10.1111/0272-4332.00040
- 20 [22] Saltelli A, Ratto M, Andres T,
21 Campolongo F, Cariboni J, Gatelli D,
22 et al. Global Sensitivity Analysis: The
23 Primer. John Wiley & Sons; 2008. ISBN-
24 10: 0470059974
- AQ2** 25 [23] Woodall P, Oberhofer M, Borek A.
26 A classification of data quality
27 assessment and improvement methods.
28 International Journal of Information
29 Quality. 2014;**3**(4). DOI: 10.1504/
30 ijiq.2014.068656
- AQ3** 31 [24] Ishikawa K. Guide to quality control.
32 Asian Productivity Organization. 1976.
33 ISBN: 92-833-1036-5
- 34 [25] Tague NR. The Quality Toolbox.
35 Milwaukee, Wisconsin: American
36 Society for Quality; 2005. p. 15. ISBN-
37 10: 0873896394