Data Quality Mining in Electronic News Paper

S. Brintha Rajakumari*

Department of CSE, Bharath Institute of Science and Technology, Chennai, Tamil Nadu, India; brintha.ramesh@gmail.com

Abstract

This paper presents the study about the attributes of quality representation of data and a case study about how effective, the data representation has been made with "Science & Technology" column of 'The Hindu' daily news paper web portal.

Keywords: Data Quality, Intrinsic DQ, Representational DQ, Web Portal

1. Introduction

Data quality is a new research area that represents one of the biggest challenges for data mining. It plays an important role in the efficiency and effectiveness of organizations and businesses. Data quality refers to the accuracy and completeness of the data, also measured by the structure and consistency. Web portals need to offer Data quality that represents a common interest for data consumers and portal providers. A web portal or public portal has lot of information from multiple sources¹ that meets user requirements². Web portals should be userfriendly and there is a need for users to ensure that the data obtained are right for their needs.

2. Classification of Data Quality

Data Quality is classified into four categories, Intrinsic DQ, Accessibility DQ, Contextual DQ and Representational DQ. Each category has many dimensions like Accuracy, Completeness, Consistency, Timeliness, etc.³ (Table 1). Accuracy of data is the degree to which data correctly reflects the real world object or an event being described. An example of data Accuracy is the bank balance in the customer's account is the real value customer deserves from the Bank. Completeness of data is the extent to which the expected attributes of data are provided. For example, a customer data is considered as complete if all customers' addresses, contact details and other information are available. Consistency of data means that data across the enterprise should be in synchronized with each other or the absence of data conflicts. An example of data inconsistency is a credit card is cancelled, and inactive, but still the card billing status shows due. The timeliness of data is extremely important which depends on user expectation. Quality of data in the web portals can be analyzed using the survey method. The survey has been made with the web users who are regular to use the online "The Hindu" web portal.

The scope of the study in this paper includes only the intrinsic and representational data quality categories of "Science & Technology" column of 'The Hindu' web

Table 1. DQ Categories and dimensions

-	0
DQ Category	DQ Dimensions
Intrinsic DQ	Accuracy, Timeliness/Currency.
Accessibility DQ	Accessibility, Access Security
Contextual DQ	Relevancy,Value-added, Completeness
Representational DQ	Content coverage/Amount of data, Consistent Representation/ Writing Style, Interactivity, Layout, Multimedia Presentation, Navigation Quality,Organization, Achieves/Documentation.

portal. Table 2 shows the Data quality, its dimensions and its definitions.

3. Quality Analysis

Data Quality (DQ) is often defined as "fitness for use" i.e., the ability of a collection of data to meet user requirements^{4, 5}.

 Table 2.
 Definitions for the DQ Dimensions

Category	Dimensions	Definitions
	Accuracy	Ensure data are the correct and valid values.
Intrinsic		The news is up to date.
	Timeliness or	Information in the
	Currency	articles is useful to our work or life.
		The website
	Content	includes appropriate
	coverage	information and features.
	Consistant	The pages of the portal should be consistent
	Consistent	and apply it to all the
	or Writing	nages in the portal
	style	Alternatively, try not to
		use more than two or
		three styles.
		Easy to effectively
	Interactivity	retrieve specific
		information on the site.
		The art of the overall
	Layout	design of a page, such
Representational		as arrangement of
		graphics and text
	Multimedia	Use of audio and video
	presentation	content.
	Navigation	Links to other websites
	quality	or between pages
		The information
		presented on the
	Organization	pages of the portal
		should be organized
		visual characteristics
		such as size of letters
		images, colours, data
		grouping etc.
		Storage and provision
	Archives	of past articles or past
		newspapers.

This definition and the current view of assessing DQ, involve understanding DQ from the users point of view⁶. Newspapers can provide online versions, that are not mirror images of print versions, instead offer something extra such as interactive features or information that could not fit in print version⁷. There are number of newspapers available on internet some with general information and some papers are complete with archives. The Hindu newspaper is one among the complete newspaper available on the internet via the web portal http://www. thehindu.com/.The online web portal of this paper consists of many columns which covers various information every day. But the case study in this paper has analyzed the data qualities like Intrinsic DQ, and Representational DQ in the 'S&T' (Science & Technology) column alone.

The "S&T" Column of the portal includes several sub columns like Agriculture, Energy & Environment, Gadgets, Internet, Science and Technology. The survey has been done by feedback analysis using statistical tool. A questionnaire has been framed and the feedback has been collected from the undergraduate and postgraduate Students, Research scholars, Academicians of various disciplines and web users who go through this portal in a regular basis.

The questionnaire has been framed with 5 to 6 questions for each dimension. The web user has to enter their rating percentage values in the specified columns.

Likewise more than 80 feedback forms collected and calculated the average of each dimension. Table 3 shows the part of the attribute questionnaire.

4. Intrinsic Quality

The Intrinsic DQ specifies the basic qualities of data like accuracy and timeliness. Accuracy ensure data are correct and valid values, Timeliness refers to the information is up to date and the articles are useful to our work or life. Chart 1 represents the Intrinsic DQ in which the accuracy is 80% and the timeliness is 90%. On an average, the intrinsic quality of data, that's accuracy and timeliness is measured as 85% from the feedback collected.

5. Representational Data Quality

The Representational DQ specifies the way in which the data are presented or made available in the web portal.

Table 3. Attribute Questionnaires

Theouracy: Ensure data are the correct and variat	51			
Questions:	Low	Medium	High	Very High
The presence of advice for agriculture growth and advancement from concerned experts.				
The information about Energy & Environment.				
The arrival of new gadgets.				
The information of new web addresses and contact details.				
The innovative research projects explained clearly.				

1. Accuracy: Ensure data are the correct and valid values.

2. Timeliness or Currency: The news is up to date. Information in the articles is useful to our work or life.

Questions:	Low	Medium	High	Very High
The information about agriculture are up to date.				
The information about current energy & Environment.				
The issues about various gadgets are given in right time.				
The new research opportunities reach the user in time.				
The latest news in science and technology are provided in time.				

The writing style in S&T column changes periodically.



Chart 1. Intrinsic data quality.

The representational DQ includes content coverage, writing style, interactivity, layout, multimedia presentation, navigation, organization and archive. These factors help the online web portal to present their information in a most effective manner to the wide user. Chart 2 represents the Representational DQ in which the data representational quality has been observed through various factors.



Chart 2. Representational DQ.

From the chart 2 it is observed that the navigation of data is very high as 86%, and the Layout, organization and archive of the presentation of data are high and found to be 85%, 84 % and 85% with a very small difference of 1% among them from the feedback collected.

Content coverage %	Consistent Representation%	Interactivity %	Layout %	Multimedia Presentation %	Navigation %	Organization %	Archive %
65	80	70	85	45	86	84	85

 Table 5.
 Representational DQ Percentages

Content coverage and interactivity are found to be 65% and 70%. Multimedia presentation is found to be a medium value of 45%.

6. Conclusion

Understanding content and consumer preferences is unique, rather than asking consumers to describe what kind of news and information they want and how they should be covered, this study measured online newspaper content and measured consumer reaction. The study on the "S&T" column of "The Hindu" web portal have shown just the amount of presence of Intrinsic and Representational Data qualities which is quantified by their Data quality dimensions as previously mentioned in the data classifications section. Through quantifying the data quality dimensions, the study has been made with the exact presence of intrinsic and representational data qualities. This paper has made a sample study to quantify the Data qualities through their dimensions, so that importance can be given to areas in which a poor quantifying measure is shown. Future study can lead to all the columns of the paper, identification of lacking data quality in the portal, suggestions to improve the data quality can also be included.

7. References

- 1. Pernici, B. and Scannapieco, M.,2002. Data Quality in Web Information Systems. Proceeding of the 21st International Conference on Conceptual Modeling, pp: 397-413.
- 2. Caro, C. Calero, I. Caballero, and M. Piattini. Defining a Data Quality Model for Web Portals. in WISE2006, The 7th International Conference on Web Information Systems Engineering. 2006. Wuhan, China: Springer LNCS 4255. p. 363-374.
- 3. M. Angelica Caro,Coral Calero, Ismael Caballero, Mario Piattini., Data Quality In Web Applications: A State Of The Art ,IADIS International Conference on WWW/Internet 2005, pp 364-368.
- 4. C. Cappiello, C. Francalanci, and B. Pernici., Data quality assessment from the user's perspective in International Workshop on Information Quality in Information Systems, (IQIS2004). 2004. Paris, Francia: ACM. p. 68-73.
- 5. D. Strong, Y. Lee, and R. Wang, Data Quality in Context. Communications of the ACM, 1997. Vol. 40, N° 5: p. 103 -110.
- 6. S.A. Knight and J.M. Burn, Developing a Framework for Assessing Information Quality on the World Wide Web. Informing Science Journal, 2005. 8: p. 159-172.
- Chyi, H.I. & Lasorsa D., Access, Use and Preferences for Online Newspapers. Newspaper Research Journal, 1999, 20(4), 2-13.