

From Collection Resources to Intelligent Data: Thoughts on the Construction of Intelligent Digital Humanities Platform for Local Historical Documents of Shanghai Jiao Tong University

Qian Yin

Library, Shanghai Jiao Tong University, Shanghai, China

Email address: yqian@lib.sjtu.edu.cn

Xing Zhuoyuan

Library, Shanghai Jiao Tong University, Shanghai, China

Email address: zyxing@lib.sjtu.edu.cn

Shi Xiaohua

Library, Shanghai Jiao Tong University, Shanghai, China

Email address: xhshi@sjtu.edu.cn

Li Yushang

History Department, Shanghai Jiao Tong University, Shanghai, China

Email address: liyushang@sjtu.edu.cn



Copyright © 2019 by Qian Yin, Xing Zhuoyuan, Shi Xiaohua, Li Yushang. This work is made available under the terms of the Creative Commons Attribution 4.0 International License: <http://creativecommons.org/licenses/by/4.0>

Abstract:

Local historical documents originate from daily life of people and belong to special collection resources which are not published publicly. They are valuable assets of universities and libraries. At present, most documents had only finished digitalization or partial datalization work. However, the requirements of deep knowledge mining in documents data, providing visual analysis, and effectively supporting the research of historic humanities scholars had not been fully met.

Taking the local historical documents project of Shanghai Jiao Tong University as an example, using relevant techniques of digital humanities, the in-depth analysis and utilization research of documents data is carried out. On the one hand, the core database of the documents is established based on standardizing metadata cataloguing and establishing metadata association. On the other hand, based on the core database, an intelligent digital humanities system platform is constructed. The platform is to realize full-field retrieval and display of the documents, text analysis, association analysis, statistics, and visual presentation of knowledge. In addition, in the process of using the platform for research, humanities scholars can continuously expand the data dimensions and the relationships between data, achieve intelligent supplementation of documents data and platform self-learning.

The concept of digital humanities has led to a new direction of database construction and platform development. In the exploration and practice of digital humanities, libraries should continue to

widening thinking, improve service and innovation capabilities, and provide better research perspectives, research environments, research support and research experience for humanities scholars.

Key words : Local Historical Documents, Digital Humanities, Intelligent Data, Intelligent Platform Construction

1. Introduction

Local historical documents originate from daily life of specific area people and belong to special collection resources which are not published publicly. For a long time, university history departments and libraries have paid more attention to seek and collect resources than to integrate the resource scientifically, and this will lead to these precious documents cannot be fully and systematically revealed and utilized. At present, most documents had only finished digitalization or partial datalization work. However, the requirements of deep knowledge mining in documents data, providing visual analysis, and effectively supporting research process of historical humanities scholars had not been fully met. In this paper, we introduce the local historical documents project in Shanghai Jiao Tong University (SJTU), take advantage of relevant techniques of digital humanities, and fulfil in-depth analysis intelligent documents data. According to the research needs, we propose a data centralized framework with resources construction, database construction and intelligent platform construction, to provide further research perspective, research environment, research support and research experience for historic humanities scholars.

2. Introduction of Digital Humanities

In era of digital network, academic collections in libraries is more and more similar. The unique academic, historic and valuable special resources have become an important factor for the sustainable special collection strategy in libraries. These are most precious wealth of libraries, and also the embodiment of the advantages and competitiveness of libraries. The competition of library resources in the future will be the competition of special collections with "informal publications" as the core¹.

Library stores a large number of digital special collection documents, but all documents data need to be interpreted by humanities scholars to have professional meaning. In the process of historic humanities research, scholars usually face various difficulties. They have limited ability to obtain available information related to their own research projects from a large number of resources, and are hard to deal with a large-scale of data. The significance of digital humanities support lies in deep mining and intelligent analysis of large scale texts, and the main object of digital humanities research is humanities data resources.

Digital humanities emerged in the 1990s, and accompanied with the changes in research methods of humanities scholars. The application of digital technology enriches the methodological system of humanities research, broadens the horizon of humanities research, changes the traditional humanities research environment, and more importantly, provides novel research methods, tools and platforms for traditional humanities research.

Kaplan² tries to represent Big Data research in digital humanities as a structured research field, and intend to draw a map for Big Data digital humanities as big cultural datasets, digital culture and digital experiences. Manovich³ presents a number of core concepts from data science that

are relevant to digital art history and the use of quantitative methods to study any cultural artifacts or processes in general.

Digital Humanities promotes the paradigm change of humanities research⁴. At the same time, processing, experience and achievements of humanities research on special collection resources can in turn promote the improvement of digital technology.

3. Collection and Collation of Local Historical Documents Resources

With increasing number of discoveries and collections of local historical documents, thousands of confused documents will bring more difficulties for users to use them. Scholars can utilize these resources better only through orderly organization. Effective use must be based on the premise of good documents collation. With abundant resources preservation, collation and organizational experience, libraries have become an indispensable force in the process of collecting and collating local historical documents.

3.1 Introduction to the overall situation of local historical documents

Local historical documents, also known as folk historical documents, folk documents and so on, mainly come from the daily life of local people. All documents are found and obtained by means of folk collection. These documents are not published or reorganized. This kind of characteristic is close to the nature of archives, they are the words and other forms of materials produced in the course of people's daily activities. Their main forms include contract documents, litigation documents, village regulations, account books, diaries, letters, singing books, scripts, religious ritual books, prescriptions, daily miscellaneous books, etc. which cover a wide range of social, economic, political and cultural fields of folk life⁵. Scholars can get a glimpse of folk historical memory and restore the rich and colourful life of civil society.

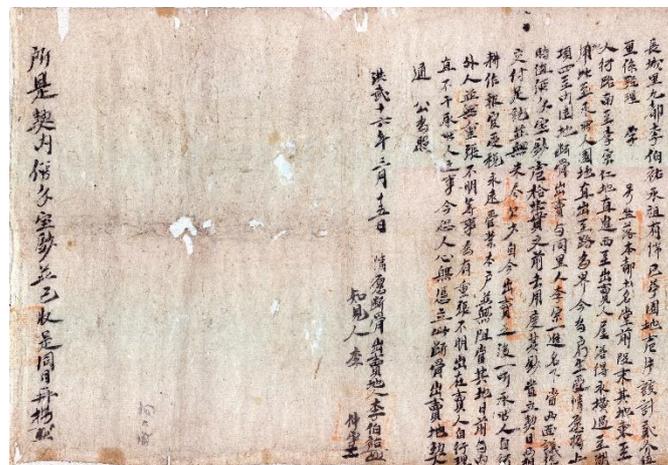


Figure 1 An Example of Local Historical Documents

Shanghai Jiao tong University began to engage in local historical documents collation and research in 2007. After 2012, the collected documents were handed over to the library for preservation, restoration and collation. By 2018, more than 350,000 pieces of local historical documents have been collected, and these documents mainly from Zhejiang, Anhui, Fujian, Jiangxi, etc. And the collection have covered 36% of these four provinces' counties⁶. At present, they are the most systematic documents group that could reflecting the traditional society in Southeast China. According to the characteristics of documents from different sources, libraries adopt different measures of collation and preservation.

3.2 Digital Construction of Documents

In recent years, “The trend of synchronization between historical research and digital era is becoming more and more obvious”⁷. Digital resources have become an important basis for academic exchanges and research. As a basic work to provide raw resources for historical research, libraries have carried out many large-scale digital processing projects of local historical documents. With transforming the physical form into electronic form to store and use, libraries are responsible for long-term preservation, and accelerate the popularization of humanistic knowledge, and provide better support for the research of humanities scholars.

We divides the digitalization of local historical documents into two levels: long-term preservation level and data service level. The main difference between the two is that long-term preservation level will use higher scanning resolution, while data service level will take into account the speed of access and server pressure to adopt lower scanning resolution. According to the characteristics of different forms of documents to develop the corresponding technical standards and formats. There is no technical obstacle in digitization itself, but it requires more effective on-site management. If documents damaged seriously, they need to be repaired first and then scanned⁸.

However, electronic documents will not enhance the use value, except for their ease of dissemination. Digital humanities is not a simple humanities digitalization, which contains a large number of new topics to be solved by humanities scholars and digital technicians. Turning resources into intelligent data, so that they can integrate into the scientific research process and support scientific research and innovation. This is common goal of the history departments and the libraries.

4. Intelligent Data Construction of Local Historical Documents

The rapid development of big data brings out an important concept: smart data. In the field of digital humanities, smart data can be understood as a kind of information that is meaningful to humanities scholars⁹. Zeng et al.¹⁰ put forwards how the digital humanities have embodied Smart Data and Big Data concepts and approaches, which demonstrate an emerging and significant change in terms of methodology.

Smart data has strong semantic expression ability and association ability. Zeng¹¹ introduces a number of semantic enrichment methods and efforts that can be applied to LAMs (libraries, archives, and museums) data at various levels, aiming to support deeper and wider exploration and use of LAM data in DH research. It shifts the focus of big data from “big” to the essence of data at the knowledge level, which can fully represent the semantic attributes and characteristics of data resources. By effectively organizing and extracting the data of local historical documents, using digital technology and combining with the knowledge of humanities scholars themselves, the data could match the needs of humanities research intelligently, and make the researcher’s judgement, decision-making and behaviour more wisely. Smart data has played and will continue to play a huge role in the field of digital humanities.

4.1 Self-built Metadata Scheme

Metadata structure determines the way in which documents are retrieved and used. Good metadata construction will transform the use of local historical documents from “reading” to “analysis”. In construction of the local historical documents database, how to establish a reasonable metadata structure and clarify internal relations of different documents, such as time sequence, geographical distribution and interpersonal network, should be the issues that humanities scholars and technicians need to consider.

Taking contract documents in local historical documents as an example, library tried to combine the knowledge of archival science and library science, used special metadata design methods to extract resource characteristics and user needs, and designed a set of metadata specifications applicable to contract documents.

We investigated the demand of historical humanities scholars for resources, and paid more attention to correlation between resources, people and families. We analysed the document resources, extracted the resource attributes, and summarized three major attribute modules: External Physical Characteristics, Content Features and Identity Recognition Features. Based on reusing of DC, 18 metadata elements is defined, including 4 custom elements¹², as shown in the following table.

Table 1 Local Historical Documents Elements Set

Attributes	ID	Element	Reuse	Element Description
Content Features	1	Title	dc:title	title of local historical documents
	2	Character	custom element	all the important persons and their roles in the original text
	3	Family	custom element	administrative divisions and family information
	4	Cause	dc:description	events or acts recorded in document, such as litigation, trading, tax-paid behaviour, etc.
	5	Geographic Information	dc:description	geographical information in document content
	6	Area Code	dc:spatial	code information of a certain region obtained by querying the code book
	7	Document Date	dc:data	the date (in Chinese year number) that the document was generated
	8	Gregorian Date	dc:data	Gregorian calendar year, corresponding to the Chinese year number
	9	Object	custom element	the object of the transaction, land, houses, goods, rights, etc.
	10	Amount	custom element	amount due to transfer of property rights
Physical Features	11	Number of Page	dc:extent	quantity of documents
	12	Size	dc:extent	document size
	13	Material	dc:format	materials of documents
	14	Note	dc:description	other important information about the physical form of local documents
Identity Recognition Features	15	Type	dc:type	types of local documents
	16	Identifier	dc:identifier	the serial number naturally generated when the document is recorded
	17	File Number	dc:identifier	unique number for each document
	18	Language	dc:language	language information

The ultimate goal of self-built metadata is to construct a local historical documents data service platform. After collating and cataloguing, we will form a relatively systematic and complete collection of resources, which is the basic data resources to support research and writing of humanities history.

4.2 Realizing Multiple Associations of Resource Data

Through metadata, descriptive text can be transformed into analysable data, and data use will be more inclusive and flexible. And metadata can implement contextual association, people and events correlation analysis, and relevance analysis with other documents. In order to clarify the internal relations of different documents, such as time series, geographical layout, event correlation and interpersonal network, multiple related elements are set up in the specification of describing metadata.

(1) Resource Physical Association

Using the file number element as the unique identity of each local historical document. The file number is fixed before and after digitalization.

(2) Chinese-Western Calendar Association

The local historical document uses the Chinese year number to record the date, and it is hard for users to clearly understand local area specific time. Setting a Gregorian date corresponds to it, which could provide a perspective to observe the changes of local economy and society on the time axis. At same time, the events recorded in documents can be placed in the domestic and international position of the same period for comparison and calculation.

(3) Family Association

Local historical documents have accumulated many generations of documents of some rural families. There are close connections between different families, and also inheritance. By setting the family element to record the administrative division and family information, and try to gather the documents of the same region and the same family together.

(4) Geographical Association

The combination of geographical element and family element can more clearly reflect the geographical relationship between the documents, helping researchers to further explore the geographic information of the resources¹³.

(5) Character Association

Most local historical documents record the civil behaviour of individuals, families, and organizations. Therefore, the element of the character is set. And through the statistics of the identity information of the character, the family, social and economic relations of the relevant persons can be related.

5. Intelligent Digital Humanities Platform Construction

5.1 Consideration on Platform Construction

Caroline et al.¹⁴ argue some urgent questions, with the recent turn towards what has come to be called 'platformisation', that is the construction of a single digital system that acts as a technical monopoly within a particular sector, and it is certainly the case that the implications of machine learning infrastructures and their black-boxed techniques for sorting, classification and ordering large amounts of data. Cornelius et al.¹⁵ compared two academic networking platforms, HASTAC and Hypotheses, to show the distinct ways in which they serve specific communities in the Digital Humanities (DH) in different national and disciplinary contexts.

By investigating several digital humanities platforms and sorting out the needs of historic humanities scholars, from the perspective of application, the intelligent digital humanities platform should have the following characteristics:

(1) Retrieve and Discover Data Resources

Helping humanities scholars to obtain the information they need from a large number of local historical documents resources is the basic function of the platform. The system can provide the retrieval of all the metadata field and feed the retrieval results back to the scholars.

(2) Text Analysis and Statistics Based on Data Mining

Local historical documents have a lot of content to be revealed and reused, but they are all hidden inside the entity and need to be deeply explored and analysed. Based on the multiple associations of data, time span distribution, geographical distribution, and family distribution of the documents can be statistically analysed, and this information can also be used to perform association analysis and reveal the potential value of the documents.

(3) Preserve and Manage Scientific Research Data

In the process of research, humanities scholars will produce a lot of research data. These data are various in form, complex in variety and numerous in number. They can be properly preserved and revealed, which can form a complete research context of scholars and their teams. Humanities scholars can set version numbers for their research outcomes based on time nodes to form a comparable and traceable outcome set.

(4) Analyse and Display Research Data

Traditional data analysis and visualization tools need to use more professional knowledge or require higher learning costs, develop an easy-to-use tool set in the platform, help researchers analyse the research data and sort out the research results. And present the results in a more intuitive and visual way for analysis and discussion by the research team members, which in turn promotes the development of the entire research process.

(5) Open Data Service Support

Without data interoperation, DH platform will lack the driving force of development, and the value of data will not be fully utilized. It is an effective way to improve the vitality of resource by sorting out the underlying data to obtain data catalogues and providing data support services at database level, data interface level or data application level. At the same time, researchers can also add and modify data under the data review process of the platform to ensure the accuracy and completeness of the data.

(6) Application of Artificial Intelligence

The application of artificial intelligence in the field of digital humanities platform can be tried from the following two aspects: Intelligent recommendation of information resources through user behaviour data analysis; Resource learning algorithm with supervised feedback. Scholars mark resources in their research process, which can more accurately index the attributes and labels of resources, and also can expand the dimension of data and the association between data.

(7) Community Communication

Users are accustomed to sharing and exchanging information in a community-based manner. All kinds of resource service platforms have certain community functions. For example, book reviews, ratings, labels, etc. of academic resources. The support of community communication

can enhance the communication between users of the platform and enhance the utilization and activity of the system.

5.2 Platform Architecture Design

Based on the above considerations, the architecture of the intelligent digital humanities system platform is shown in Figure 2.

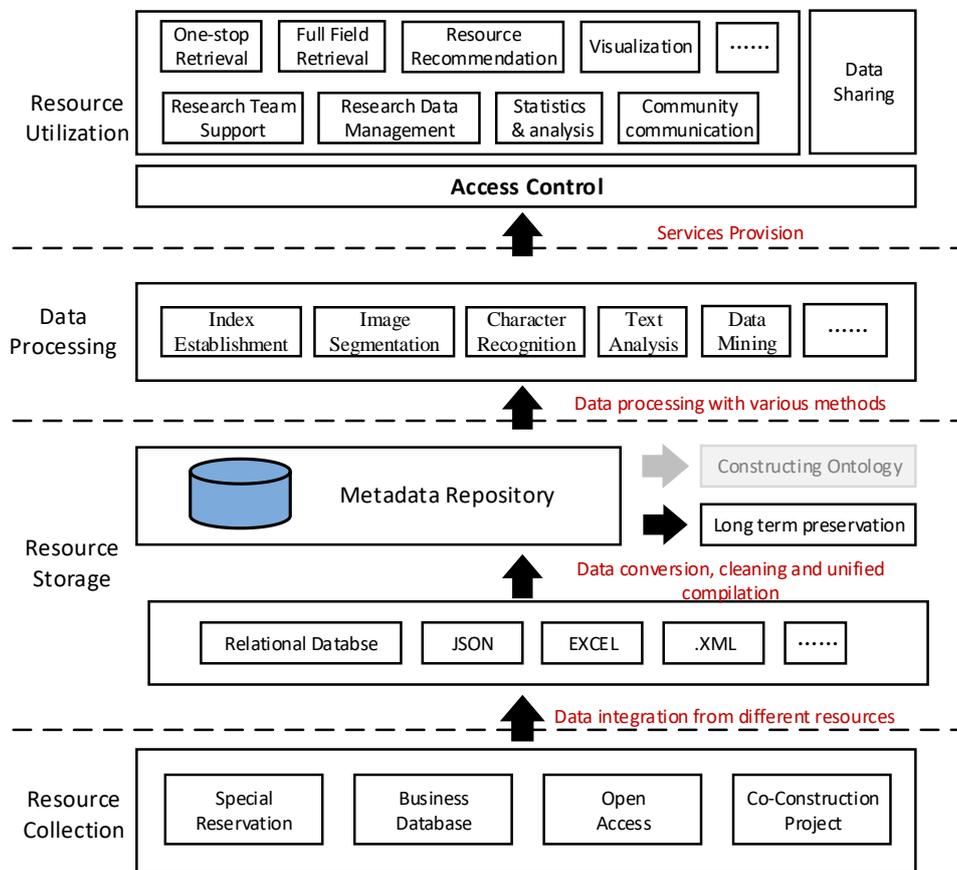


Figure 2 Platform Architecture

(1) Resource Collection

There are four types of data sources : Library Special Reservation, Business Database, Open Access Data and Co-Construction Project Data.

(2) Resource Storage

Resource data need to be translated, removed duplicates, cleaned and combined because of the different data source, different ways of obtaining data, different metadata coding and different resource description. According to our metadata scheme of local historical documents, the metadata is coded uniformly and stored in the metadata repository. These data will be preserved for a long time. In the future, we can identify metadata entities, build knowledge ontology, and publish linked data to realize the purpose of automatic association, reuse and sharing with external resources.

(3) Data Processing

Based on the metadata repository, several data processing and intelligence analysis technology can be used. For example, index establishment, image segmentation, character recognition, text analysis and data mining and so on. Many hided information will be show to the higher level after data processing step.

(4) Resource Utilization

The basis for data utilization is access control, and any data call should be made under access control protection. On this basis, the system can provide various types of services for humanities scholars, such as one-stop search service, full-field search and acquisition services, and provide recommended resources. Another important aspect of data utilization is data sharing. Through data interface, the data is opened to third-party systems for use, improving system interoperability.

5.3 Platform Function Module Design

The functions of the platform are divided and modularized from humanist dimension, library dimension and digital humanities technology dimension, as shown in Figure 3.

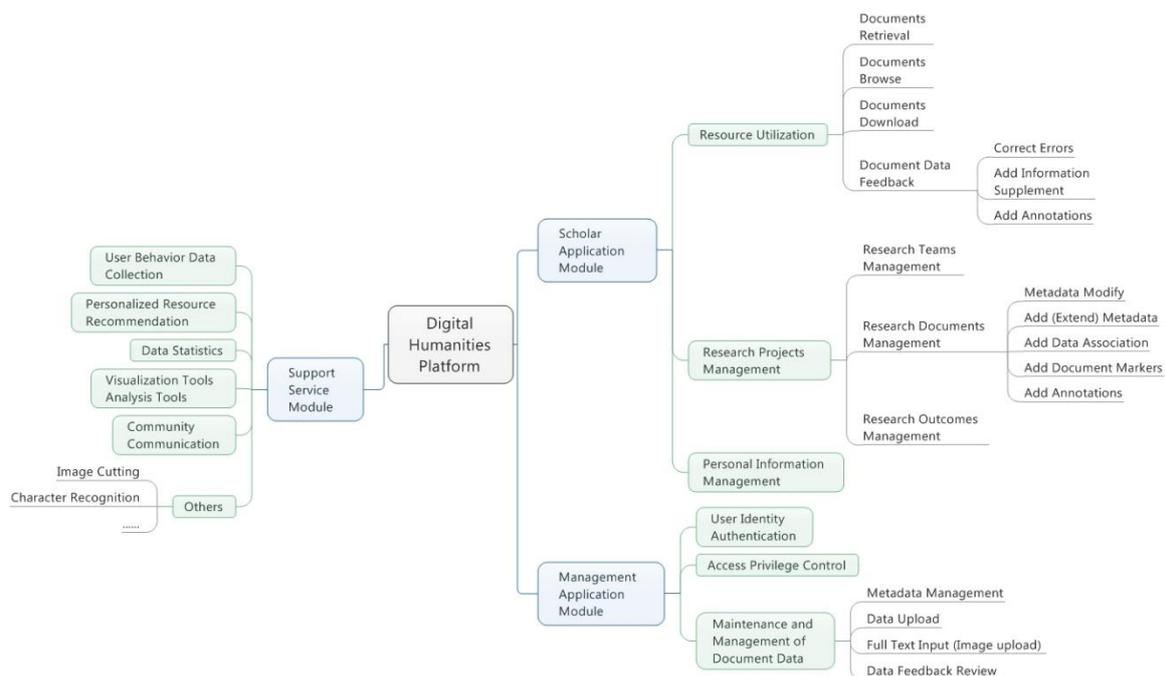


Figure 3 Function Module Design of Platform

(1) Scholar Application Module

It is a functional module based on the operation of humanities scholars. There are three main functions as follow.

a. Resource Utilization

Scholars can search, view and download the documents, and can also store the documents in their research document management module for centralized manage and use. When scholars find that the documents information is missing, wrong, or can provide relevant annotations, they can submit feedback to the system to help the system improve the quality and utilization of the documents.

b. Research Projects Management

Scholars can set up their own research projects on the platform and manage them. They can research in the name of individuals or teams. They can manage the research documents, such as modify the metadata, add data associations, add tags to data, and so on. These contents can be selected to be visible only to themselves, visible to team members, or visible to public.

Scholars can also manage the research outcomes of the project, and they can choose to disclose or not disclose these outcomes.

c. Personal Information Management

Scholars can manage their personal information.

(2) Management Application Module

This module is set up to manage whole platform. These functions mainly include user identity authentication, determine whether the user has the right to use the platform, and set access restrictions on the document resources. For example, only allow intra-school access or limit IP segment access. Manage and maintain the core data of the platform, such as metadata management, data update and upload, full-text input, review the data feedback submitted by scholars, and update the corresponding content of the database after confirmation.

(3) Support Service Module

This module serves the digital humanities platform itself and assist users to make full use of platform functions. These functions include: collection and analysis of user behaviour data, personalized resource recommendation, statistics of various documents data, providing analysis and visualization tools, and effective analysis and visualize scholar's research data and research results, community communication and some other features.

The function of the digital humanities platform is not only the storage and retrieval of data, the main purpose is to make the resource truly usable by humanities scholars, and become useful data, fresh data, and intelligent data. And on this basis, it provides a better research environment and research support for humanities scholar, helping them to reorganize knowledge, discover problems, bring new research perspectives, and provide decision-making basis for the future work.

6. Conclusion

In this paper, we expound how to collect, organize and utilize special collection resources, such as local historical documents in SJTU, by digital humanities thoughts, and propose an intelligent digital humanities data framework to effectively support humanities research. We hope that these useful thoughts will provide some inspiration and reference for institutions and individuals who use the same type of special collection resources. At the same time, in the road of exploration and practice in digital humanities, libraries should not only do important job in protection, development and open utilization of special collection resources, but also constantly expand research ideas, enhance service and innovation capabilities with various data innovation efforts, and seize opportunities to enhance academic and social status of libraries.

Acknowledgments

This work was supported by the Social Science Planning of Shanghai (Grant No. 2018BTQ002)

References

- 1 Wanguo,L.,Ying,H.. A Summary of the Academic Seminar in Digital Resource Construction and Knowledge Service in 2017[J]. Journal of academic library,2018,36(01):12-17. (in Chinese)
- 2 Kaplan, F.: A map for big data research in digital humanities. *Frontiers in digital humanities* 2, 1 (2015)
- 3 Manovich, L.: Data science and digital art history. *International Journal for Digital Art History* (1) (2015)
- 4 Ling,Z.Review of the Library and Digital Humanities Forum[J]. Journal of academic library,2018,36(2):5-10. (in Chinese)
- 5 Zhengman,Z. Folk Historical Documents and Cultural Inheritance Research [J]. *Southeast Academic*, 2004(S1) : 293-296.(in Chinese)
- 6 Fang,L.,Jin,C.,Xin,W. Planning and Practice for the Historical Documents Digitalization in recent Llibrary holding of Shanghai JiaoTong University [J],*Journal of academic library*,2015,33(02):77-83. (in Chinese)
- 7 Hang,L. Digitization of historical documents calls for active involvement of scholars [N]. *Chinese Social Science Today*,2014-04-02(A02). (in Chinese)
- 8 Baoguo,Z.,Xiaohua,S.,Xin,W. Research of Quality Control in the Process of the Large-scale Book Digitization [J].*Research on Library Science*,2017(04):51-55. (in Chinese)
- 9 Kobielius J. The evolution of big data to smart data[EB/OL] · Keynote, Smart Data Online 2016 · (2016-07-13)[2019-03-08] · Video available at: <http://www.dataversity.net/big-data-smart-data-big-drivers-smart-decisionmaking/>
- 10 Zeng, M.L.: Semantic enrichment for enhancing lam data and supporting digital humanities. review article. *El profesional de la informaci_on* 28(1) (2019)
- 11 Puschmann, C., Bastos, M.: How digital are the digital humanities? an analysis of two scholarly blogging platforms. *PloS one* 10(2), e0115035 (2015)
- 12 Jie,Z.,Fang,L.,Meng,T. The Design and Use of Metadata Application Profile for Local Historical Property Contracts[J].*Library and Information Service*,2017,61(08):106-111. (in Chinese)
- 13 Xin,W., Jie,Z., Meng,T. Study of Pathway of Information Organization and Display in Deed Document From Geographic Aspect[J]. *New Century Library*,2018(04):55-59. (in Chinese)
- 14 Bassett, C., Berry, D.M., Fazi, M.B., Pay, J., Roberts, B.: Critical digital humanities and machine-learning. In: *DH* (2017)
- 15 Manovich, L.: Data science and digital art history. *International Journal for Digital Art History* (1) (2015)