

Svitlana Chukanova

National University of Kyiv-Mohyla Academy, Research Library

<https://orcid.org/0000-0002-5717-5050>

RESEARCH DATA MANAGEMENT AND DIGITAL CURATION AS A LIBRARY ACTIVITY

Abstract. In the age of information technologies, it is important to recall that research data becomes more significant than mere scholarly publications without data sets. Data sets are treated as texts, they may have DOI, and they can be referred to with the utilization of any bibliographic style. Research data management is a multileveled and convoluted procedure, which can be executed in association with the library or digital curators in the laboratories of research institutions. For advanced analysts, it is urgent to working up a network of individuals to disperse their exploration results and with this reason including the library and data experts.

Keywords: research data management (RDM), research data lifecycle, digital curation, information collections, data librarian

Zarządzanie danymi badawczymi i ich archiwizacja w działalności biblioteki

Streszczenie: W epoce technologii informacyjnych dane badawcze nabrały większego znaczenia niż zwykle publikacje naukowe bez zbiorów danych. Zbiory danych traktowane są jako teksty, mogą posiadać DOI i mogą być cytowane zgodnie ze standardem opisu bibliograficznego. Zarządzanie danymi badawczymi stanowi wieloetapowy i zawły proces, który może być realizowany we współpracy z biblioteką lub jednostkami zajmującymi się zachowaniem zasobów cyfrowych w laboratoriach instytucji badawczych. Dla zaawansowanych analityków ważną sprawą jest praca nad tworzeniem sieci umożliwiającej szeroką eksplorację rezultatów badań i włączenie w ten proces także bibliotek i ekspertów danych.

Słowa kluczowe: zarządzanie danymi badawczymi, cykl życia danych badawczych, archiwizacja danych cyfrowych, zasoby informacji, bibliotekarz danych

Introduction

Today, we observe strong connection between librarianship and information science in the world, which is affecting library practices. As an example, the Library of Glasgow University (Scotland) pays considerable attention to the technical aspects of working with information: creating an online profile of the scholars, depositing thesis and dissertations, tracking information-seeking behavior on the library website, managing research data in accordance with the principles of academic integrity, automating library processes, archiving and digitizing online, etc. The library and IT Center at Glasgow University are combined in entire department which means that library and IT services are provided by one huge team of information specialists. Glasgow University is not the only institution, which takes into consideration research data management (RDM) and data curation, but it definitely has one of the best practices in conducting such processes in the library. We can name the following institutions concerned with this topic in Europe and in the USA such as Edinburgh University¹, University of Glasgow², DataOne³, Florida State University⁴, North Carolina Chapel – Hill University⁵, Ottawa University⁶. Of course, RDM is not just a library practice and it is possible to say that it is more scholarly activity but as research data caution principles are similar to library activities it is logical to provide RDM in libraries and by librarians as well as by research institutions and scholars.

Research data management: An overview

Research data management is now gaining popularity in the world of librarianship and information science. This practice is already widely used by some leading research centers and organizations and is mandatory in some institutions that receive grants for their research projects. The libraries of the universities of

¹ Research Data MANTRA. (n.d.). Retrieved November 3, 2019, from <https://mantra.edina.ac.uk/>

² University of Glasgow – MyGlasgow – Data management support for researchers – RDM at Glasgow. (n.d.). Retrieved November 3, 2019, from <https://www.gla.ac.uk/myglasgow/datamanagement/rdatatglasgow/>

³ DataONE. (n.d.). Retrieved November 3, 2019, from <https://www.dataone.org/>

⁴ Data Management, Florida State University Libraries. (n.d.). Retrieved November 3, 2019, from <https://www.lib.fsu.edu/drs/rdm>

⁵ Research Data Management and Sharing, Coursera. (n.d.). Retrieved November 3, 2019, from <https://www.coursera.org/learn/data-management>

⁶ Library, University of Ottawa. (n.d.). Retrieved November 3, 2019, from <https://biblio.uottawa.ca/en/services/faculty/research-data-management>

Scotland, in Edinburgh and Glasgow, are actively working with their lecturers and academics to provide qualified data management, as funding for universities is directly dependent on their quality of research activity.

The example of RDM training activities provided by the Library of Edinburgh University is presented by online and classroom-based activities⁷. Online activities include MOOC developed jointly with North Carolina University (Research Data Management and Sharing) and free non-credit course MANTRA (Management Training) which has video and text materials explaining main RDM principles and approaches of different scholars.

Classroom-based training activities include workshops oriented towards developing and improving RDM skills and scholars' profile. Classroom workshops are dedicated to the following topics: benefits of RDM for scholars, RDM planning activities and tools, ethical issues of involving vulnerable and personal data, Open Refine tool usage for cleaning the data, SPSS tools for data, how to assess risk of data disclosure, quality and quantity assurance, effective use of the data in research, data visualizing tools and techniques. Library users may request additional training activities on specific topics concerning their research.

The second example of RDM activities provided by the library is Glasgow University Library. The library provides research support and depositing data at University repository Enlighten. RDM activities are performed according to Research Data Management Policy of the University accepted in 2015⁸. RDM support services at Glasgow University offer to researchers a variety of activities and include general RDM training, RDM plans writing consultations and review, templates for RDM plans, support of compiling grant applications, depositing data assistance.

In order to understand why such successful research and education institutions of the world are concerned with RDM, we should regard and learn what RDM is and focus on its importance to the modern-day researchers. We should mention that the reputation of the research institution is closely connected to its effectiveness. In order to show the effective use of their resources and providing new knowledge, the researchers should remember about managing their data. Research data management is one of the best practices for making a good public image of the research institution.

⁷ *Research data training and skills*, The University of Edinburgh. (n.d.). Retrieved January 18, 2020, from <https://www.ed.ac.uk/information-services/research-support/research-data-service/training>

⁸ *Good Management of Research Data Policy*. Retrieved January 18, 2020, from http://www.gla.ac.uk/media/media_406078_en.pdf

To begin with, we need to indicate what research data is. To answer this question, we are going to address American researcher Christine Borgman. According to her, data is considered as the main research unit, which represents factual interpreted information in forms of observation, calculation, records, experimental data, and digital data and exists in such types as texts, numbers, and multimedia, software, specified by discipline, specified by tools⁹.

Data also exists in different formats, has numerous versions, which affects naming conventions during data curation. With this purpose librarians use different software open codes or subscribe in order to arrange these elements in a proper data set. Usage of naming conventions helps in organizing the data and provides better data preservation and access.

Research data undergoes several stages of lifecycle, which may vary from one research institution to another. In order to understand the basic principles of the research data lifecycle, we shall regard one of the simple examples developed by the Library of Ottawa University¹⁰. This example contains the main elements of the research data lifecycle and explains what each stage means in order to proceed to the next level and, thus, pushing the research forward.

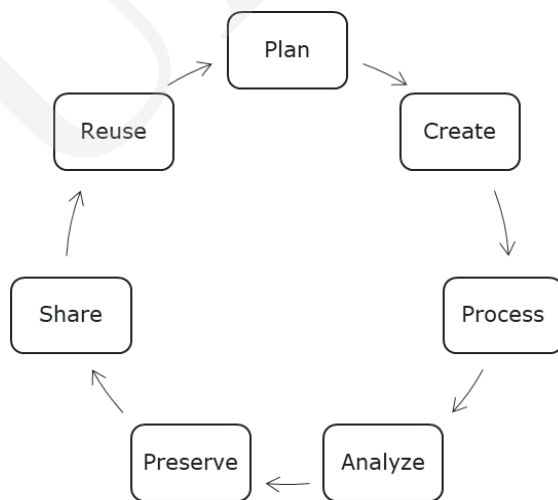


Figure 1. Research Data Lifecycle developed by the Library of Ottawa University.

⁹ Borgman, C. (2010). *Scholarship in the Digital Age: Information, Infrastructure, and the Internet* (pp. 120–122). Cambridge: MIT Press.

¹⁰ *What is research data management?*, Library. University of Ottawa. (n.d.). Retrieved November 3, 2019, from <https://biblio.uottawa.ca/en/services/faculty/research-data-management/what-research-data-management>

Every research starts with planning, especially when it is a part of the grant project. As Joyce M. Ray¹¹ in the Introduction to Research Data Management states, sharing research data returns a profit to investments made into the science, especially when the research may be very crucial to the public health and wellbeing. Ray also indicates that it is necessary to avoid duplication of research data, thus, economizing time for providing further research and not duplicating work¹².

In order to provide a plan for RDM, some researchers may use online instruments like DMP Online or DMP Tool and compile the 2-page document, which can be altered during the research performance. Planning also includes the sources review and informed consent from the participants of the medical, educational, sociological research, etc. After planning, the researcher creates data sets while gathering the raw data during experiments, interviews or surveys, etc. At the stage of creation, it is worth to assign metadata to the packages. The raw data must undergo some processing: cleansing, renaming, versioning, anonymizing, validation, etc.

When data is ready for analysis, some Data Science tools may be used in order to perform operations in gathering important information from the data, which will support the hypothesis of the research. At this stage, some articles may be produced, and it is important to cite the data as it would serve as a basis for research results. In order to cite any data, it should be kept appropriately thus preserved on research institution servers, institutional repositories or in case of absence of the last – on open repositories for data, which have the seal of approval. The data sets and packages are shared with the help of the repositories, as well as stored in them. The accessibility of data also enables the second usage of data by the researcher him/herself or by other researchers interested in similar topics.

As one can observe the lifecycle is very simple and easy to understand. The library can help researchers at several stages of research data processing. The library can help with data curation, renaming, versification, anonymization, etc. The library may help in sharing, citing, preserving data in the same way as it helps with text repositories of the institution. With this purpose (providing data curation for the researchers), librarians need to know how to use some tools in order to curate the data and to consult researchers on how they can manage their research data sets.

¹¹ Ray, J.M. (2014). *Research Data Management. Practical Strategies for Information Professionals* (p. 12). West Lafayette: Purdue University Press.

¹² *Ibidem*.

Planning tools

The most popular tools for building up an RDM plan that librarians can suggest to researchers in order to comply with grant organization requirements are DMPTool and DMP online. In order to create an RDM plan, the researcher should register an account on one of these platforms and follow the suggestions of the system while creating the document. DMPTool has a dashboard, which opens after registering an account¹³. The researcher can choose a template, which suits to requirements of providing research grants. At the stage of creating a plan, the system asks a curator or researcher about the project, research organization and granting institution. After selecting and filling in these points the system opens a template, which the researcher must fulfill. This template is composed of several elements where the researcher describes project details by indicating its title, funding, abstract, etc. After the first stage, there is a plan overview which helps the researcher to check whether the plan fits the standards of selected funding organization and after reading instruction the researcher fills in the plan itself by answering the questions concerning the points of collecting and preserving data, ethics, metadata, and other important aspects. When the document is ready, it can be shared and downloaded.

DMP online platform¹⁴ has similar principles of creating plans. In order to create an RDM plan, the research must register for an account and undergo a similar procedure of filling in the information about research project details and funding information.

The role of the librarian at this stage is either to consult a researcher on how to write the plan or to write the part that concerns the details on metadata creation, preservation and sharing, thus, contributing to the research as a curator. In the research data management plan, the researcher will indicate who is responsible for data curation: either the librarian or the researcher him/herself.

Creation

As we mentioned before, the researcher creates data sets while performing surveys, carrying out some experiments and observations, thus, the library has a little impact on this process as this must be the researcher's responsibility to choose what data will be necessary to testify the main hypothesis or theory. The only possible help which can the librarian suggests to the researcher at this stage is to assign metadata in order to describe selected data sets for further preservation, sharing, and reuse.

¹³ DMPTool. (n.d.). Retrieved November 3, 2019, from <https://dmptool.org/>

¹⁴ DMPonline. (n.d.). Retrieved November 3, 2019, from <https://dmponline.dcc.ac.uk/>

Process

After the data received, it is necessary to process it. During processing, the researcher cleans, renames, anonymizes data, and creates data packages. One of the tools, which can be suggested by the librarian or digital curator for this purpose is Data Package Creator from Frictionless Data¹⁵. This tool is equipped with an easy-to-understand guide that enables better data package creation. It has a detailed description of how to assign metadata to your data package, which is very important in RDM and for the long-term preservation of the data.

Analysis

Different types of data can be analyzed depending on the sphere where it is used. Several tools may be applied in this process: Jupiter Notebook, Zeppelin Notebook, Watson Studio, and other tools peculiar to Data Science¹⁶. Data packages become at this stage a kind of evidence on which the hypothesis and theories are based, thus, data serves as the basis for numerous publications and must be cited as any other source of information. At this stage, the librarian can suggest the best ways of citing data by means of bibliographical description.

Preservation and sharing (data repositories)

At this stage, the researcher should care about ethics, in particular, academic integrity and research ethics, security (where to keep and store the data either on a personal USB drive or on the research institution server, repository, etc.) In order to provide security to the research data packages, it is always better to deposit data packages to thematic or institutional data repositories. In her book *Digital Curation for Libraries and Archives*, Stacy T. Kowalczyk regards among other examples of research data lifecycles the one developed by the Library of Congress. This lifecycle has traditional and digital versions, and – what is important – the digital version

¹⁵ *Well packaged datasets*. (n.d.). Retrieved November 3, 2019, from <https://frictionlessdata.io/field-guide/well-packaged-datasets/>

¹⁶ Open Source tools for Data Science, Coursera. (n.d.). Retrieved November 3, 2019, from <https://ru.coursera.org/learn/open-source-tools-for-data-science>

indicates that after each stage of lifecycle, i.e. creation, distribution, library, and archival collection, and long-term access, there is a stage of preservation¹⁷.

According to Kowalczyk, the following scheme can be presented as a linear, one direction model which represents moving forward.

The model of traditional preservation can be depicted in the form of a timeline. Thus, the model of traditional preservation is rather linear, and the action of preservation takes place once after the collection is compiled and is in the stage of preparation for long-term access.



Figure 2. Model of traditional preservation according to Kowalczyk.

Unlike the traditional model, the model of digital preservation includes the stage of preservation action frequently because it deals with other types of resources which require updates.

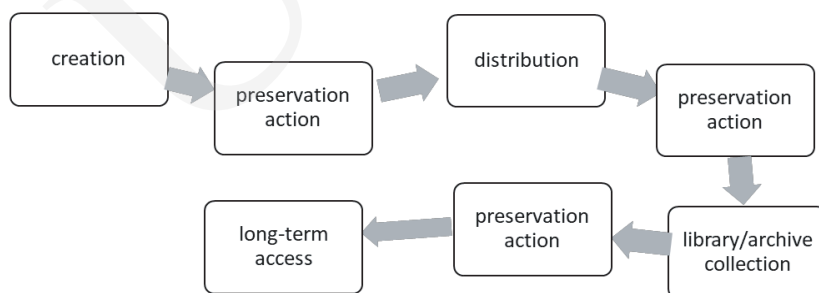


Figure 3. Model of digital preservation according to Kowalczyk.

One can see that in the model of digital preservation, the frequency of the preservation stage is larger than in the traditional model. The traditional model is used for preserving physical objects, while the digital model is designed for digital resources of information. As digital information changes rapidly, it is important to preserve these objects frequently which is reflected in the model.

¹⁷ Kowalczyk, S. (2018). *Digital Curation for Libraries and Archives* (pp. 24–38). Santa Barbara: Libraries Unlimited.

How to find Data Repository

In case when the institution does not possess its own repository, the researcher can upload the materials to a thematic data repository. In order to find one, it can be useful to use the tool re3data.org¹⁸. During the selection of the repository, the researcher will encounter 27 filters and numerous subfilters according to which it is easier to find out what repository will suit the research data for particular studies. The main filters are the following:

1. Subjects,
2. Content Types,
3. Countries.¹⁹

The variety of subjects and subtopics is represented by 4 main sectors:

- Humanities and Social Sciences,
- Life Sciences,
- Natural Sciences,
- Engineering Sciences.

Each of these sectors is divided between a variety of topics and subtopics and each subtopic offers some number of repositories for the research data. This tool is very useful for the researchers and for the librarians and data curators during their training and consultations for the researchers.

Reuse (Data citation principles)

As was stated before, data act as a basis for the research output. In order to demonstrate these data packages, each researcher needs to cite the data properly in the articles and conference proceedings.

As is the case for any other material, it is no exception for data to be cited. The data citation is ruled by certain principles. The initiative group FORCE 11 (The Future of Research Communications and e-Scholarship) developed these principles. Joint Declaration of Data Citation Principles underlines the following important statements²⁰:

1. Importance,
2. Credit and attribution,
3. Evidence,

¹⁸ Home, re3data.org. (n.d.). Retrieved November 3, 2019, from <https://www.re3data.org/>

¹⁹ *Ibidem*.

²⁰ *Joint Declaration of Data Citation Principles – FINAL*, FORCE11. (n.d.). Retrieved November 3, 2019, from <https://www.force11.org/datacitationprinciples>

4. Unique identification,
5. Access,
6. Persistence,
7. Specificity and Verifiability,
8. Interoperability and Flexibility.

The importance principle means that data should be considered as important as any other research output (e.g. research paper, etc.). It must be cited as it serves as a shred of evidence, which supports the theory or hypothesis. In order to provide better data citation, unique identifiers such as DOI would serve the best part. These data citations and identifiers provide access to the data packages and research results, so the community and scientific world can test the hypothesis or benefit from the research output. Identifiers also provide persistence, which means that data is findable and reusable while it exists, thus, it is flexible and interoperable because other scholars may use it.

These data citation principles act as certain rules or codes of conduct for supporting Academic Integrity principles in the research world thus, the researchers would not be afraid to share their data because of some misconduct or lack of attribution.

FAIR principles

In open knowledge society, it is important to take care of the quality of the data which we use in the researches as we should remember that the data is the main research unit and serves as evidence for the research output. FAIR principles is an acronym that is deciphered as: Findable, Accessible, Interoperable, and Reusable²¹.

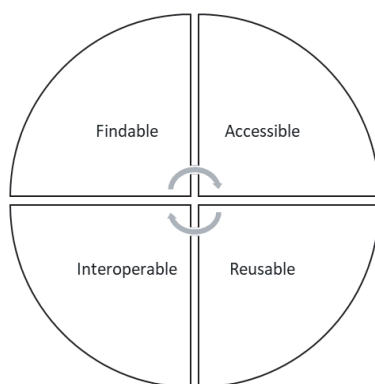


Figure 4. FAIR principles according to LIBER.

²¹ *FAIR Principles*, GO FAIR. (n.d.). Retrieved November 3, 2019, from <https://www.go-fair.org/fair-principles/>

According to LIBER (Association of European Research Libraries), these principles may be explained in the following way:

- The findable principle implies the existence of DOI or other links or identifiers, metadata descriptions necessary for finding data via search mechanisms or data repositories.
- The accessible principle means that data are kept securely in a data repository with the seal of approval and can be read and understood both by humans and machines.
- The reusable principle provides necessary licenses that are clear and easy to comply with in order to use these data packages.
- The principle of interoperability means that the language of metadata is understandable and widespread across the scholarly community²².

As we can see, the huge role of metadata (data about data – thus, description) is observed through the whole range of principles. The role of librarians in supporting FAIR principles lies in promotion, consultation, curation, and providing guidance on depositing in correspondence with FAIR principles.

On the GO FAIR platform, there is a detailed guide on how to adjust data packages in accordance with FAIR principles. The “FAIRification of the data”²³ implies several stages:

1. We receive non-FAIR data;
2. We analyze it from the perspective of the content type and structure;
3. We define the semantic model (description of links between entities in the data set);
4. We make data linkable via Link technologies with the use of the semantic model;
5. We assign a license for effective reuse of the data;
6. Define rich metadata for the data set to be read by humans and machines;
7. Deploy (publish) FAIR data resource.

As with any other information resource, the data sets should be made accessible and clear for understanding. As we may find any record in the e-catalog of any library, we can find any data set on data repository nowadays. This became possible with implementing of FAIR principles, Data Citation principles, development of

²² *Implementing FAIR Data Principles: The Role of Libraries*. (2017). Retrieved November 3, 2019, from <https://libereurope.eu/wp-content/uploads/2017/12/LIBER-FAIR-Data.pdf>

²³ *FAIRification Process*, GO FAIR. (n.d.). Retrieved November 3, 2019, from <https://www.go-fair.org/fair-principles/fairification-process/>

data repositories, and effective data curation practices developed for new library professionals – data stewards or data curators. Robin Rice and John Southall in their work *The Data Librarian's Handbook* state that Data Librarianship is like IT data support and is closely connected to cultural norms of research funding in academia. This kind of library practice started first with technical support but now relates to scholarly communication principles in general.²⁴ These library scientists also pay attention to the role of librarians as advisors on consent agreements when dealing with sensitive data. They consider assistance in creating forms of consent an opportunity for the librarian to be involved in the research process, thus providing feedback. The consent forms have different types so, it is better for the librarian to gather wide-spread templates as with an RDM plan. It may be of assistance in further consultations. Consent agreement is integral part of research data packages when sensitive data involved, so it is the document which will be kept along with the data set in the repository and the librarian is obliged to check it for compliance with the standards for these kind of researches.²⁵

Conclusion

Research Data Management is a relatively new practice for modern-day librarianship in the world. This practice implies the possession of strong competencies in IT, research, librarianship, information literacy, academic integrity, and Open Access principles. With the development of data and technologies for its processing there is a need for qualified specialists to manage these data sets and curate them. Thus, data curators or data librarians became new and demanded professions in information society, especially as different research grant programs have requirements connected to Research Data Management Plans.

The task of data librarian is like the task of traditional librarian – curate information resources, consult customers (researchers) on how to use them. Information age implies a huge amount of data which we may call “Big Data” and these types of information can be analyzed only by means of Data Science and special software, but the data librarians can consult how to curate those data sets which are selected by researchers as necessary for their research projects and need long-term preservation.

With the development of the scholarly world we observe information resources development. There is a need for newly created professionals that would curate

²⁴ Rice, R., & Southall, J. (2016). *The Data Librarian's Handbook* (pp. 15–17). London: Facet Publishing.

²⁵ *Ibidem*, pp. 123–125.

these resources thus, data sets. Data librarians help researchers to comply with FAIR principles by making data sets findable, accessible, interoperable, and reusable. These principles pay a considerable amount of attention to the metadata which means that the librarian's help is necessary for such purposes. As data sets must be findable and accessible it is very important to involve information professionals – librarians in the description process of data packages and sets because library specialists will assign keywords and fill in metadata fields in the data repository in a more professional way than ordinary repository users.

Summarizing, we may say that the role of the librarian in the process of research data management is crucial in terms of providing effective scholarship activity. Librarians act as consultants while providing training activities on RDM principles, special tools for different stages of research data lifecycle, and RDM planning. Librarians act as data curators when they are involved in RDM process and provide data services by depositing it to a repository or cleaning and visualizing data sets. Research support is inevitable part of library work in modern-day research libraries and institutions and RDM is a process which in its essence is like traditional library processing of information resources. Data amount is growing which means that its management must be performed as effectively and efficiently as possible. In the variety of the data it is important to select the most valuable information which is necessary for the research. If finding the data for a project is a task of researcher, then to keep and help in providing access to it is a task of library professional.

References

- Borgman, C. (2010). *Scholarship in the Digital Age: Information, Infrastructure, and the Internet* (pp. 120–122). Cambridge, MA: MIT Press.
- Data Management, Florida State University Libraries. (n.d.). Retrieved November 3, 2019, from <https://www.lib.fsu.edu/drs/rdm>
- DataONE. (n.d.). Retrieved November 3, 2019, from <https://www.dataone.org/>
- DMPonline. (n.d.). Retrieved November 3, 2019, from <https://dmponline.dcc.ac.uk/>
- DMPTool. (n.d.). Retrieved November 3, 2019, from <https://dmptool.org/>
- FAIR Principles, GO FAIR. (n.d.). Retrieved November 3, 2019, from <https://www.go-fair.org/fair-principles/>
- FAIRification Process, GO FAIR. (n.d.). Retrieved November 3, 2019, from <https://www.go-fair.org/fair-principles/fairification-process/>
- Good Management of Research Data Policy. Retrieved January 18, 2020, from http://www.gla.ac.uk/media/media_406078_en.pdf
- Home, re3data.org. (n.d.). Retrieved November 3, 2019, from <https://www.re3data.org/>
- Implementing FAIR Data Principles: The Role of Libraries. (2017). Retrieved November 3, 2019, from <https://libereurope.eu/wp-content/uploads/2017/12/LIBER-FAIR-Data.pdf>

- Joint Declaration of Data Citation Principles – FINAL*, FORCE11. (n.d.). Retrieved November 3, 2019, from <https://www.force11.org/datacitationprinciples>
- Kowalczyk, S. (2018). *Digital Curation for Libraries and Archives* (p. 38). Santa Barbara, CA: Libraries Unlimited.
- Library. University of Ottawa. (n.d.). Retrieved November 3, 2019, from <https://biblio.uottawa.ca/en/services/faculty/research-data-management>
- Open Source tools for Data Science, Coursera. (n.d.). Retrieved November 3, 2019, from <https://ru.coursera.org/learn/open-source-tools-for-data-science>
- Ray, J. M. (2014). *Research data management : practical strategies for information professionals* (p. 12) . West Lafayette: Purdue University Press.
- Research Data Management and Sharing | Coursera. (n.d.). Retrieved November 3, 2019, from <https://www.coursera.org/learn/data-management>
- Research Data MANTRA. (n.d.). Retrieved November 3, 2019, from <https://mantra.edina.ac.uk/>
- Research data training and skills*, The University of Edinburgh. (n.d.). Retrieved January 18, 2020, from <https://www.ed.ac.uk/information-services/research-support/research-data-service/training>
- Rice, R., & Southall, J. (2016). *The Data Librarian's Handbook* (pp. 15–17). London: Facet Publishing. DOI: <https://doi.org/10.29085/9781783301836>.
- University of Glasgow – MyGlasgow – Data management support for researchers – RDM at Glasgow. (n.d.). Retrieved November 3, 2019, from <https://www.gla.ac.uk/myglasgow/datamanagement/rdmatglasgow/>
- Well packaged datasets*. (n.d.). Retrieved November 3, 2019, from <https://frictionlessdata.io/field-guide/well-packaged-datasets/>
- What is research data management?*, Library. University of Ottawa. (n.d.). Retrieved November 3, 2019, from <https://biblio.uottawa.ca/en/services/faculty/research-data-management/what-research-data-management>