# Development of Education Curriculum in the Data Science Area for a Liberal Arts University

Zhihua Zhang, Toshiyuki Yamamoto and Koji Nakajima

# Development of Education Curriculum in the Data Science Area for a Liberal Arts University

Zhihua ZHANG, Toshiyuki YAMAMOTO and Koji NAKAJIMA

Kansai University of International Studies, Kobe, Hyogo 650-0006, Japan
z-zhang@kuins.ac.jp

**Abstract.** Data science has emerged as a field that will revolutionize science and industry. The development of human resources for data science has become an urgent issue in every aspect of the digitizing society. However, a curriculum to meet the needs in such a digitizing society is not available to higher education in Japan, especially in the realm of liberal arts. Such liberal arts students do not have enough basic math education such as statistics before entering the university. In response to the required situation of the approved program for Mathematics, Data Science, and AI Smart Higher education, we proposed a conceptual curriculum model for the data science education program, which systematically incorporates the knowledge module of data science while remedying the weakness in the basic math skills and barriers to be considered in the process of learning data science concepts. The goal of this paper is to propose an integrated curriculum based on the conceptual model for the faculty members in a small-sized private liberal arts university, where students lack basic math skills, IT skills, and the basic knowledge of data science. Issues consist of the curriculum on knowledge area and subjects, the implementation approach of data science education courses, and the fusion of data science with expertise education are discussed. A sample course will be showcased at the end.

**Keywords:** Conceptual Curriculum Model, Curriculum Development, Data Science Education, Liberal Arts University, Stage-wised Refinement Model.

## 1    Introduction

With the advancement of advanced information technology, work across nearly all domains is becoming more data-driven in society, where many various kinds of data are generated and relatively easily available, it is required to utilize these data to create new value. Digital transformation is also being promoted speedily in all industries so that new digital technologies can be used to develop new business models. In various fields of social, industrial, and business situations, problem-solving based on the existence of big data is emphasized. Therefore, there is an urgent need to develop human resources who have mathematical thinking ability and data analysis/utilization ability and who can create value and solve problems based on this, in addition to specialized education, in each field of humanities or science. Now, as more data and

ways of analyzing them become available, more aspects of the economy, society, and daily life will become dependent on data. It is imperative that educators, administrators, and students begin today to consider how to best prepare for and keep pace with this data-driven era of tomorrow [1].

In response to the required situation of MDASH, which is an approved program for Mathematics, Data Science and AI Smart Higher education, accredited education program promoted by the MEXT of Japan, universities nationwide are actively developing their own mathematical, data science, and AI education curriculums [2]. On the other hand, a literacy level and applied level model curriculum formulated by the Inter-University Consortium for Mathematics and Data Science Education is limited to about 4 credits each. Here, to secure the opportunity for all university students to take courses, it has no choice but to enroll in common basic education, even when considering securing teachers in charge of classes [3].

In general, data scientists need to acquire knowledge of mathematics and statistics, IT, especially programming skills, and knowledge of specialized fields, and the range are wide. In addition, there are systematic patterns due to the remarkable difference in difficulty, but the ones that are suitable for liberal arts college students are not organized. Whether it is possible to realize a systematic curriculum model that comprehensively considers "knowledge module", "implementation approach", and "barriers to be considered" in the lesson design of Data Science Education at a private university of liberal arts. This is an academic question in this research.
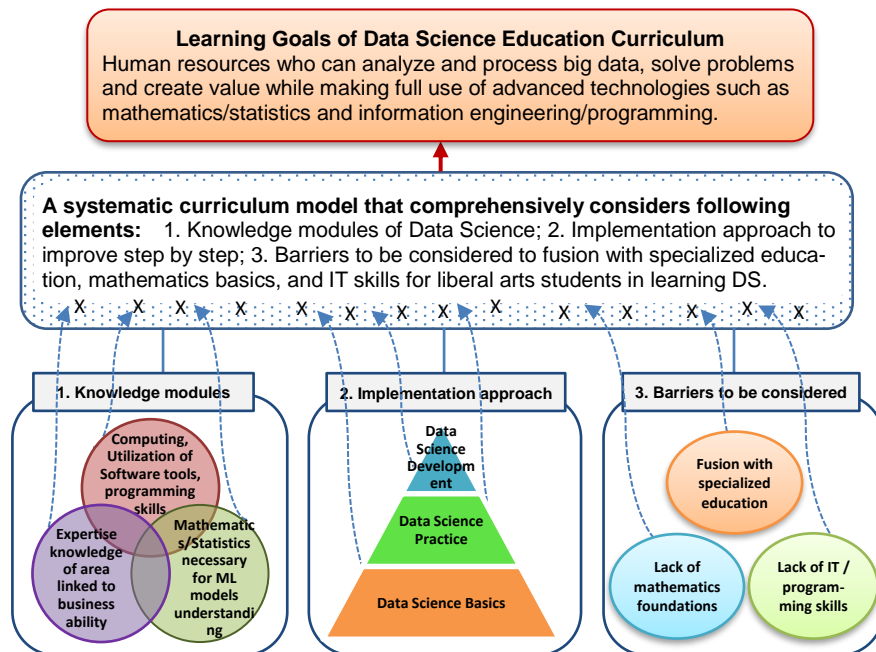


**Fig. 1.** Academic questions about the systematic curriculum model.

Here we proposed a conceptual diagram (model) of curriculum design for our university's data science education program (see Fig. 1).

In this conceptual model, firstly, we define the learning content of data science as the relevant "knowledge modules"; secondly, we define the educational "implementation approach" as the learning stage that is gradually refined from simple to complex; finally, the "barriers to be considered" for liberal arts students, such as the characteristics of weak basics of mathematical knowledge, insufficient basic information skills and programming skills, and how to integrate data science education into the education of their major field knowledge in teaching and other considerations, to achieve a systematic data science teaching plan for liberal arts college students. These are the academic questions about the systematic curriculum model we defined.

## 2 Knowledge Area and Subjects of Data Science Education Curriculum

### 2.1 Knowledge Area of Data Science Education

Data science is a growing field of study that combines discipline expertise, programming skills, and knowledge of mathematics and statistics to derive meaningful insights from the data. Moreover, data science experts should apply machine learning algorithms to numbers, texts, images, videos, audio, and more to build artificial intelligence (AI) systems that can perform tasks that normally require human intelligence. Because it has become possible to handle big data efficiently, data science is creating many new values not only in business but also in various situations. For instance, in business sites, it is the ability to imagine business models and consumer behavior from data. In the task of aggregating and visualizing data, you may face the challenge of how to interpret the aggregated results for the real site business field. Students should have the ability to organically combine and utilize the data science knowledge they have learned to complete these tasks.

However, despite being used as a general term covering multiple disciplines, data science is difficult to be accurately defined. The famous one is Drew Conway's Venn diagram (2013), which consists of the knowledge of substantive expertise, mathematics and statistics, and computer and programming technology. The Japan Data Scientist Association talks about the qualities of the data science skillset required as a data scientist with the following three powers. 1) Business problem solving: Ability to organize and solve business problems after understanding the background of the problem; 2) Data science: Information science such as information processing, artificial intelligence, and statistics. And the ability to understand and use the wisdom of the system; 3) Data engineering: Ability to use, implement, and operate data science in a meaningful way.

On the other hand, in the report of the National Academies of Sciences, Engineering, and Medicine for the emerging discipline of data science at the undergraduate level [1], the following points were made: new majors and minors will initially combine ingredients from existing courses, in areas such as computer science, statistics,

business analytics, information technology, optimization, applied mathematics, and numerical computing.

The knowledge module of data science education subjects constructed in this research includes such learning content: computing, utilization of software tools and programming skills on data analyzing, modeling, visualization, and presentation; the foundations of mathematics/statistics necessary for machine learning models understanding; and the expertise of area linked to real site business ability; in the lesson design of data science education, this allows for that choice and combination, making the construction of the syllabus more flexible.

The subjects of the data science education curriculum performed at our university are mainly related to three knowledge areas. The area I is about basic statistics and its utilization; Area II is about ICT and programming skills on data science; Area III is for the basics of big data, and AI evolving subjects on problem-solving content.

## 2.2 Subjects Related to the Knowledge Area

Based on the knowledge area classification of the above data science education courses, we have designed the corresponding related subjects, which are detailed below.

**Area I, Basic Statistics and Its Utilization:** Area I consists of the three subjects shown in Table 1. It is mainly about the foundations of mathematics/statistics skills.

**Table 1.** "Data Science Education Curriculum" the Area I Subjects

| Subject name | Contents | Credit |
|---|---|---|
| Basic Statistics A | A subject related to basic knowledge for reading official statistics, simple research reports, and fieldwork dissertations. Educational content includes how to read descriptive statistical data such as simple tabulation, frequency distribution, representative value, dispersal degree, cross-tabulation, and correlation coefficient; how to read graphs, and how to calculate and create them. | 2 |
| Basic Statistics B | A subject on the basic knowledge of inference statistics necessary for compiling and analyzing statistical data. Educational content includes basics of probability and probability distribution, population and sample, basic statistics and its properties, test/estimation theory and its application, and basic regression analysis. | 2 |
| Introduction to Research | This is a course aimed at acquiring basic knowledge and skills in qualitative and quantitative research, tabulation, and analysis. Specific survey methods, observation surveys, interview surveys, and questionnaire surveys will be taken up, and the basics of the survey will be acquired in an exercise format. | 1 |

**Area II, ICT and Programming Skills:** Area II consists of the three subjects shown in Table 2. It is mainly about Area II is about ICT and data science basics on utilization of software tools and programming skills.

**Table 2.** "Data Science Education Curriculum" the Area II Subjects

| Subject name | Contents | Credit |
|---|---|---|
| ICT Literacy | The purpose is to familiarize yourself with how to use data analysis using Excel or R using BYOD and to give lectures and exercises on the basics of data analysis and basic techniques of data visualization. | 2 |
| Utilization of ICT A | The purpose is to familiarize yourself with how to use Python and R language and programming technology using a personal computer and to give lectures and exercises on the basics of analysis using actual data and basic techniques of data visualization. | 2 |
| Data Science | Learn the concepts and methods of mathematics, data science, and AI behind big data and AI technology, and examples of their use. You will also learn the basics of handling data, points to keep in mind about data, and information security. In the lecture, we will understand the role of DS through lectures inviting people who are active in the front lines of society and research. | 2 |

**Area III, Basics of Big Data and AI Evolving Subjects on Problem-Solving:** Area III consists of the three subjects shown in Table 3. Area III is for the programming practice techniques and the basics of IoT, big data, and AI evolving subjects in the basics of problem-solving content.

**Table 3.** "Data Science Education Curriculum" the Area III Subjects

| Subject name | Contents | Credit |
|---|---|---|
| Data Science Theory | For beginners in programming, the R language is easier to understand. For this reason, R is useful in statistical processing that aims to analyze data and explain something based on the results. In this course, students will acquire practical DS skills through R language programming techniques, involvement with data science, data analysis, and visualization of actual problems. | 2 |
| Data Science Practice Exercise | Python is a simple, readable, and versatile programming language. Python has abundant libraries in all fields such as statistical analysis of data, AI application/machine learning, and IoT data utilization, and is used in a wide range of application fields. The purpose of this course is to learn the basic knowledge of Python, deepen the understanding of algorithms through Python programming, and learn how to utilize typical APIs and services. | 2 |
| Basics of Artificial Intelligence | Learn basic technologies of information system engineering such as algorithm/data structure, information security, information communication network, and artificial intelligence related to the characteristics of big data and the utilization of AI. Introductory level engineering knowledge subjects. Content created by external organizations (MOOC, etc.) may be used. | 2 |

# 3 Implementation Approach of the Data Science Education Courses

Using data science and artificial intelligence technology to solve practical problems in the business field is an extremely complex process. It concerns variously skilled data science professionals with diverse technical and business backgrounds. Also, it is multiple tasks across the data science life cycle. According to the data analysis model (CRISP-DM), a standard cross-industry data analysis process, including:

1) Understanding of business: Understand the business situation, problems and set project goals. It is important to accurately grasp the business situation as a numerical value from the perspective of the marketer and select the theme.
2) Understanding of data: Examine whether the data is available. Investigate data items, quantities, quality, etc. Often, it includes outliers and missing values. Make the data available.
3) Data preparation: As a pre-processing for mining, the usable data is shaped into data suitable for analysis. Missing value processing, data type maintenance, normalization, sampling, etc.
4) Modeling: Model creation or method selection. Correlation analysis, regression analysis, market basket analysis, cluster analysis, genetic algorithm (GA), decision tree, etc.
5) Evaluation: Evaluate from a business perspective whether the model is sufficient to achieve well-defined business goals.
6) Implementation: Make concrete plans to apply the results of data mining to the specific business field and take concrete actions to achieve the set goals.

Many companies have had considerable difficulty in developing human resources with these abilities. Most of the time, the first introduction to Data Science was successful, but it lacked the ability to have a holistic view of where to start in subsequent practice and advanced problem-solving. For students to have such abilities, it is essential to have a PBL-type lesson using actual data after understanding the basics of statistics and mathematics, an introduction to Data Science, basic languages and programming techniques, etc.

In order to carry out Data Science education smoothly, and let the data science talent training course succeed, we proposed a Stage-wised Refinement Model in the Data Science education program according to the defined educational "implementation approach" in the conceptual model 5).

As shown in Fig. 2, it consists of four stages from bottom to top as pillars that support the educational curriculum of data science at our university. The difference from the conceptual model is the addition of "Stage 0: Invitation to data science". This is a result of consideration given to the basics of mathematics and the lack of ICT skills at liberal arts universities.

The other three stages are Stage 1: Data science basics, Stage 2: Data science practice, and Stage 3: Data science development, respectively, which will be described in detail below.
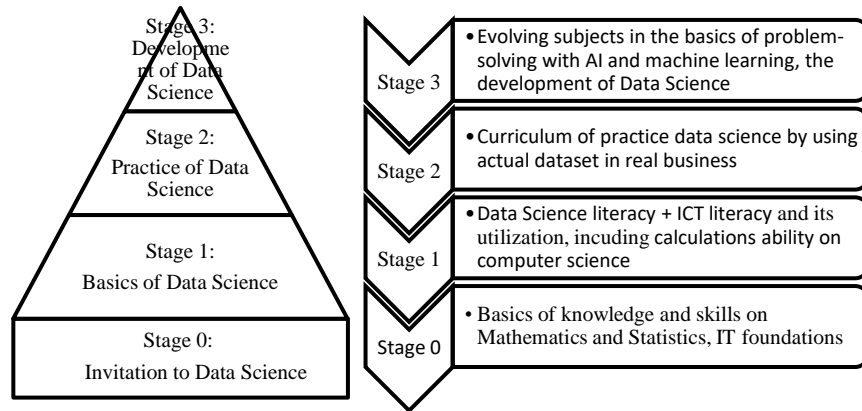
**Fig. 2.** Concept of Stage-wised Refinement Model in data science education program.

**Stage 1: Basics of Data Science**, the contents include the data Science literacy corresponding points of introduction in the Consortium model curriculum, in addition, ICT literacy skills including calculation's ability on computer science.

**Stage 2: Practice of Data Science**, the contents include practice data science by using actual datasets in real business or simulation with a specific domain application. Learn how to handle data and databases for liberal arts students and analyze data using statistical analysis software R or entry of programming language Python.

**Stage 3: Development of Data Science**, the contents include evolving subjects in the basics of problem-solving with AI and machine learning. Learn how to perform machine learning using actual data by using Python or R programming language, AI tools, etc. If it is possible, take it to PBL education that practices problem-solving and value creation.

**Stage 0: Invitation to Data Science**, the contents include the basics of knowledge and skills on Mathematics and Statistics which is aimed to clear the hurdles for liberal arts students in data science study.

For liberal arts students, the foundations of mathematics (including statistics) that are indispensable for Data Science in order to study data science literacy level in the future, to take future developmental curriculum, and to acquire problem-solving ability beyond literacy level. It is necessary to have a minimum of strength on the foundations of mathematics, as well as IT (including programming) skills.

The ultimate purpose of data science is to provide an academic basis for making rational decisions in the business field, etc., based on the analyzed data. Thus, data science is closely related to the basics of mathematics and IT skills an inseparable field. That is why we add the "Stage 0: Invitation to data science" to the Stage-wised Refinement Model shown in Fig.2.

## 4 Barriers to be Considered on the Implementation Approach

### 4.1 Problems of Lack of Mathematics Basics and Information Technology Skills

In the framework of Conway, as mentioned above, the data science education curriculum consists of the knowledge of substantive expertise, mathematics/statistics, computer and programming technology, etc. First, one starting point is to assume a curriculum organization that conforms to this. But in this process, it is hard to say that the student's learning history of mathematics and information subjects in high school was fully considered.

At a private liberal arts university such as our university, where most of the students are not good at mathematics and IT basics, even if literacy level mathematics and data science education are conducted, the basic ability of students' mathematics and information will be a big challenge. In some cases, you may also need a remedial education on mathematics and IT foundations in junior high school and high school.

To solve this problem, we are newly constructing basic mathematics content for data science education for college students in the liberal arts, which will be developed by on-demand learning. In other words, we decided to develop learning content that utilizes Python programming and corresponds to the middle and high school math knowledge modules that are indispensable for data science learning.

As an effect of this effort, it is expected to be foundational research of the Data Science education curriculum model, which is a field under development in the practice of Data Science education. Remedial education of mathematics basics and IT basics necessary for implementing the Data Science education curriculum at the university can be smoothly realized, and it is expected to contribute to the implementation of educational contents suitable for the actual situation of students and the guarantee of Data Science education quality. From the perspectives of teachers and students, there are the following merits.

From the teacher's point of view:

・On-campus resources in data science education and programming education for liberal arts students will be available.

・Teachers in the field of education can provide learning opportunities with the teaching materials of this project for the problems of students' foundations of mathematics and ICT literacy.

From the student's point of view:

・Overcome the consciousness that students are not good at mathematics before taking data science.

・It will be an opportunity to improve students' IT literacy and programming skills and they will be able to keep up with the learning of data science.

・Become able to understand the basic mathematics that was lacking in high school.

・Become able to understand the basics of IT literacy that should be learned in high school.

・Students will be able to fully utilize their BYOD as an IT stationery by learning.
・The mindset of data science will be the cornerstone of lifelong learning.

### 4.2 Problem Of the Fusion of Data Science Education with Substantive Expertise Education

Another problem about the practice of data science education is how to deal with the relationship of data science education and expertise education in each major. Discussions held above are under the policy of promoting the concreteness of the curriculum as a pillar in the framework of Conway. It turns out that at least "expertise" is lacking in the latter. However, considering that data science education is positioned as a university-wide education that goes beyond that framework, apart from expertise education, it is natural to think that "expertise knowledge" is educated in the specialized field of each faculty to which it belongs.

But can we become data scientists in a specialized field without expertise and fusion with Data Science? The answer to this question is clearly no, and it can be said that education on thinking and processes that lead to value creation by comprehensively utilizing the knowledge and background cultivated in each field is indispensable.

Almost common in many faculties, there is a need for human resources who can analyze data scientifically with a deep understanding of their major field data, and for that purpose, education is aimed at acquiring the ability to utilize data based on the "expertise knowledge". Regarding the introductory education in the proposed systematic curriculum model, first, data science centered on the utilization of big data with the aim of understanding the image of data science that utilizes the three elements of mathematics/statistics, information technology, and specialized knowledge. However, we have prepared "Data Science" and "Utilization of ICT A" as the entry subjects to give the lecture on what is brought about in solving social problems, including case studies if necessary.

We think that the most important background to acquire in practicing data science is not just the skill to master the software tools and the data processing method, but the imagination and creativity to envision value creation on a real site. In general, data acquisition is costly, so it is important to emphasize what can be done by using data and clearly show the usefulness of acquiring and accumulating data.

Therefore, we designed exactly what kind of data is likely to be obtained in what form in the area targeted for problem-solving, and what kind of value can be created by utilizing it. Subjects at the stage of "data science development" in the Stage-wised Refinement Model, need to be done in connection with specialized knowledge education. The essence of data science education is to develop the ability to systematically explain, but the key is how to deepen the experience. And it is considered effective to incorporate PBL (Project Based Learning), which deals with practical issues in collaboration with the real site, into the curriculum, although it is commonplace.

In addition, subjects such as computer science, business analysis, information technology, applied mathematics, and mathematical calculation are also important in the common curriculum. Further discussion is needed on these issues in the future.

# 5 Course Structure and Qualification Certification

From the following year, we will start a university-wide Data Science Education Program. The course structure of this education program is designed based on the proposed systematic curriculum model as a case study, which is necessary to comprehensively develop the subjects on computer and information system knowledge, statistics-related skills, domain knowledge of specific fields, problem-solving experience, critical thinking skills, and communication and presentation skills.

The practice work performed mainly considers two aspects. The first one is to systematize the relationships between the subjects while standardizing the contents of the subjects that are currently being implemented.

The second one is going to award a "Data Science Education Curriculum Certificate" by acquiring the prescribed credits from the subject group (see Fig. 3), for the purpose of encouraging students to actively participate in data science education programs.
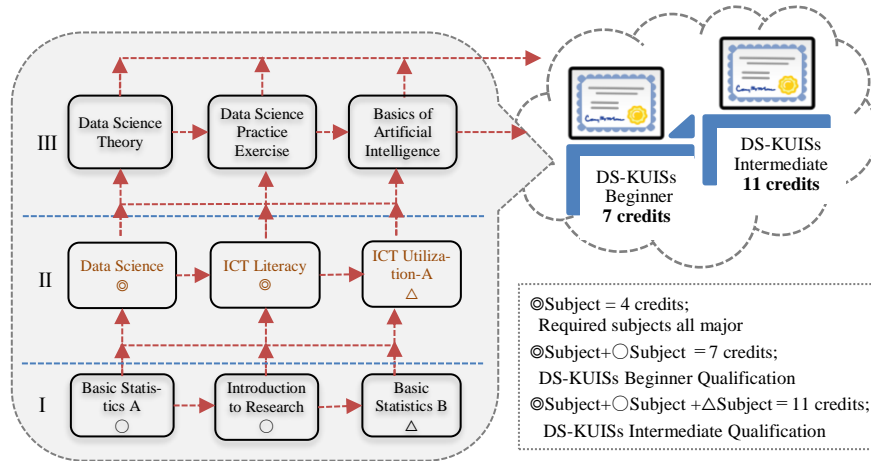


**Fig. 3.** Course structure and qualification certification at KUISs.

Each subject in the subject group I, II, III as shown in Fig. 3 corresponds to the three category-level subjects described in the knowledge area of Section 2. This is the Course Structure and Qualification Certification of the Data Science Education Program at KUISs. Here, the symbol ◎ means compulsory subjects, the symbol ◯ means elective subjects, and the symbol △ means recommended subjects. In this program, students acquire the target certification program by taking a combination of three category-level courses. If they got a certain number of credits, they will be qualified for Data Science.

Curriculum subjects including Area I and Area II are called DS-KUISs beginner, while Area III is reached are called DS-KUISs Intermediate. If a student has completed the prescribed 4 subjects from it with 7 credits, he/she will be reached the "Data Science Education Curriculum Beginner's level". If for DS-KUISs Intermediate,

he/she should complete the prescribed 6 subjects with 11 credits and will be recognized as the "Data Science Education Curriculum Intermediate level".

## 6 Conclusion

As an umbrella term, Data Science includes a broad range of theories, algorithms, methodologies, and software tools that help us to use datasets to understand and solve problems in the real world. The knowledge area of data science also widely includes mathematics/statistics, programming skills, data management techniques, and in some cases, machine learning models and algorithms widely used in the sciences, engineering, humanities, education, medicine, and business. In addition, due to the rapid development of digital transformation in society, there is an urgent need to develop human resources who have mathematical thinking ability and data analysis/utilization ability and who can create value and solve problems based on this, in addition to specialized education in each field, regardless of whether it is humanities or science.

We have proposed a conceptual curriculum model for the data science education program, which systematically considered the knowledge module related to data science learning, the implementation approach as an educational method. But the discussions are held under the policy of promoting the concreteness of the curriculum as a pillar in the framework of Conway. As a case study, this paper discussed the development of a curriculum for data science education at a private liberal arts university. However, this is a result of the fact that data science education is positioned as university-wide education that transcends that framework, apart from specialized education.

To ensure that the educational program runs smoothly, it is expected to contribute to the implementation of educational content suitable for the actual situation of students of liberal arts for the guarantee of data science education quality. That is the discussion on clear the barriers of liberal arts students and the problem of the fusion of data science education with substantive expertise education. The essence of data science education is to develop the ability to systematically explain, but the key is how to deepen the experience. This is considered effective to incorporate PBL-based learning, which deals with practical issues in collaboration with real business site problems.

Finally, we describe the data science education program starting from the following year as a case study, the course structure and the qualification certification of the data science education program at KUISs. As a future research topic, we will improve the proposed model, and let it be able to withstand the construction of a curriculum for minors in combination with existing undergraduate courses.

## Acknowledgment

# References

1. National Academies of Sciences, Engineering, and Medicine 2018. Data Science for Undergraduates: Opportunities and Options. Washington, DC: The National Academies Press, (2018).
2. MEXT, "Mathematical / Data Science / AI Education Program (Literacy Level) Requirements", Open call for participants briefing materials, (2021), https://www.mext.go.jp/content/20210305_mext_senmon01-000012801_1.pdf, last accessed 2021/04/20.
3. The Japan Inter-University Consortium for Mathematics and Data Science Education, Mathematics / Data Science / AI (literacy level) Model curriculum, (2020).
4. Zhihua Zhang: The development of new information technology and necessity of data science education, Bulletin of Sanyo Women's College, No.41, 1-20 (2020).
5. Zhihua Zhang: A Study of Digital Transformation Human Resources Development and Data Science Education Programs at a Private Liberal Arts University, Bulletin of Kansai University of International Studies Research Series, March, (2022).
6. Cabinet Office Japan: AI Strategy 2019 -AI for all people, industry, regions, and governments-, Integrated Innovation Strategy Promotion Council Decision, (June 2019), https://www.kantei.go.jp/jp/singi/ai_senryaku/pdf/aistratagy2019.pdf, last accessed 2021/12/11.
7. Peter K., Carlie I., Erick B., Pieter den H., Farhan C., Afraz J., Shubhangi V., Gartner R.: Magic Quadrant for Data Science and Machine Learning Platforms, (2021), https://content.dataiku.com/gartner-mq-2021, last accessed 2022/02/24.
8. Jeanne R., Cynthia B.: the Digital Challenge: How to Transform Your Business in the Midst of Crisis, MIT Industrial Liaison Program Webinar Series, (2020).
9. http://www.mi.u-tokyo.ac.jp/consortium/pdf/model_literacy.pdf, last accessed 2020/03/20.
10. The recommendations of the Science Council of Japan, Training of Human Resources for the Big Data Era, (2014), https://www.scj.go.jp/ja/info/kohyo/pdf/kohyo-22-t198-2.pdf.
11. Yumi M. & Kazue T.: Issues in data science education at liberal arts and private universities, Bulletin of Edogawa University No. 31 (March 2021) pp. 249-255.
12. Jeanne R., Cynthia B.: the Digital Challenge: How to Transform Your Business amid Crisis, MIT Industrial Liaison Program Webinar Series, (2020).
13. Curriculum Guidelines for Undergraduate Programs in Statistical Science, (2014), https://www.amstat.org/asa/files/pdfs/EDU-guidelines2014-11-15.pdf.
14. UCI Data Science Initiative "What is Data Science?", http://datascience.uci.edu/about/, Last accessed 2020/11/ 25.
15. Alan David Fekete: Variations in Data Science Curriculum: A View From Computing Education, Harvard Data Science Review, Issue 3.2, Spring (2021).