

A Comparative Analysis on the Effects of Cognitive Tools in Data Inquiry Cultivation

Hui ZHANG*, Bian WU, Yi-Ling HU & Yu-Jie XU

Department of Education Information Technology, East China normal University, China

*51214108043@stu.ecnu.edu.cn

Abstract: The advent of the big data era has elevated the importance of understanding, analyzing, and mining data, necessitating data science education with cognitive tools aiding data inquiry. However, limited appraisals exist, especially regarding diverse tools' impact on data inquiry cognitive aspects. This study addresses this void by examining cognitive tool impact on data inquiry. The study collected discourse and questionnaire data generated during online collaborative tasks. Epistemic network analysis and difference testing unveiled cognitive pattern disparities, learning motivation, cognitive load, and self-efficacy variations between groups. Block-based group displayed robust cognitive connections in data understanding and preparation. Text-based group focused on modeling and optimization iterations. Motivation and load didn't differ significantly, yet block-based group showed higher self-efficacy. Study concludes by acknowledging limitations and suggesting future research directions.

Keywords: Data inquiry, cognitive tools, CSCL, human-computer interaction

1. Introduction

Data inquiry is the process of extracting meaningful and valuable information from data, including the identification of understandable patterns and relationships, as well as the construction of representative data models. In recent years, many researchers have provided various perspectives on the connotative structure of data inquiry. Most of them perceive data inquiry as a collection of data skills, encompassing elements such as data collection, pre-processing, analysis, model building, and evaluation (Donoho 2017).

Various cognitive tools have emerged to support data practice and learning. These tools can be categorized into three types: block-based, text-based programming, and menu-selected (Bart et al. 2020). Block-based tools, such as RapidMiner and IBM SPSS Modeler, offer encapsulated modular systems and user-friendly graphical interfaces, enabling quick prototyping and validation of predictive models. They simplify the entire analysis process and provide instant gratification. On the other hand, text-based programming tools like RStudio and Python offer greater flexibility and scalability. They allow users to explore algorithmic details, create customized program algorithms, and offer system openness.

Block-based cognitive tools offer a potential resolution to the educational challenge of data inquiry. Their algorithmic intricacies are concealed, enabling learners to concentrate on configuring and adapting predefined algorithms and parameters. By seamlessly connecting modules through drag-and-drop actions, learners swiftly execute data inquiry, with a shallow learning curve fostering confidence and problem-solving belief. Modular encapsulation delegates computational tasks to cognitive tools, enabling learners to focus on problem solving, expanding mental capacity and reducing cognitive load (Javadpour 2022). This phenomenon is well-established in computational thinking development. Additionally, graphical cognitive tools deliver instant interactive feedback via module interconnections, prompting immediate data-inquiry knowledge reflection, potentially heightening motivation to learn.

However, there have been no studies exploring the differences in the effects of block-based and traditional text-based cognitive tools on learners' data inquiry abilities. At the same time, researchers have pointed out that future studies should explore whether certain data science practices are influenced by tool affordance (Jiang and Kahn 2020). In order to better

design teaching and learning based on block-based data science cognitive tools, clarify how cognitive tools affect learning processes and knowledge construction, this study proposes the following research questions:

- 1) What are the differences in the frequency distribution of data inquiry cognitive elements between the block-based and the text-based tools in online collaboration activities?
- 2) What are the differences in the cognitive patterns of data inquiry between the block-based and the text-based tools in online collaboration activities?
- 3) Are there any differences in learning motivation, cognitive load, and self-efficacy between the block-based and the text-based tools?

2. Method

2.1 Research Context and Participants

This empirical study was conducted during the Fall Semester of 2022-2023 at a comprehensive university in eastern China, focusing on the "Educational Data Mining" course. Designed for senior college students pursuing careers in education, the study participants were selected from two classes within the same university department, totaling 44 individuals - 10 males and 34 females. All participants had not received prior data mining instruction. Collaborative learning was facilitated in groups of 3-4 learners. One class (16 students, 5 groups, average age 21) used text-based cognitive tools, while the other class (28 students, 9 groups, average age 23) employed block-based cognitive tools.

2.2 Research Process

The experimental procedure is illustrated in Figure 1. The study spanned 18 weeks, comprising a 14-week instructional phase dedicated to learners' acquisition of data inquiry skills. Both classes were taught exactly the same by a veteran instructor with over a decade of teaching experience. Subsequently, participants undertook data analysis projects during the final 4 weeks. Upon course completion, all students were required to complete questionnaires assessing cognitive load, learning motivation, self-efficacy, and satisfaction. The questionnaire completion process typically lasted around 10 minutes.

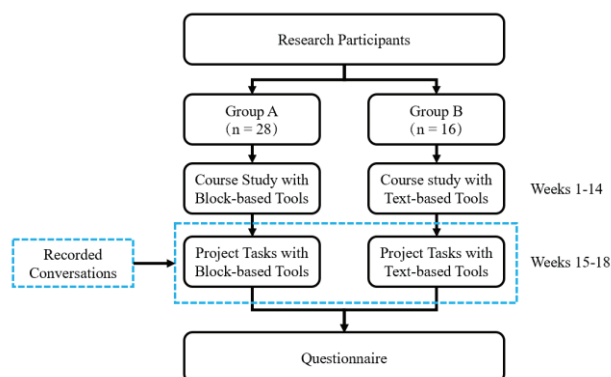


Figure 1. The flowchart of research experiment.

The two student cohorts employed distinct cognitive tools – block-based and text-based – for their data inquiry tasks. The block-based group utilized the RapidMiner platform, a robust data mining tool characterized by its visual workflow interface, facilitating intuitive design, execution, and assessment of diverse data mining tasks. Conversely, the text-based group leveraged the RStudio platform, an interactive development environment based on the R programming language. This platform encompasses a variety of functions, including a code editor and data viewer, facilitating efficient organization and management of code and data.

2.3 Data Collection

In the project's final 4 weeks, each group conducted 5 online discussions lasting 1 to 3 hours, recorded for analysis. The study employed adapted percentage scored questionnaires for learning motivation (6 items, Cronbach's $\alpha = 0.848$), cognitive load (2 items, Cronbach's $\alpha = 0.836$), and self-efficacy (3 items, Cronbach's $\alpha = 0.843$). Learning motivation drew from Lau & Lee (2008), measuring knowledge mastery and grades on a five-point scale. Cognitive load, adapted from Leppink et al. (2013), assessed internal and external load using a nine-point scale. Self-efficacy, from Tsai et al. (2020), focused on data inquiry self-efficacy with a five-point Likert scale.

2.4 Code Framework

A coding framework was applied to analyze cognitive elements in data inquiry for frequency and patterns. Adapted from Wirth & Hipp (2000), a recognized data mining model, the initial framework had 6 codes. Grounded theory guided manual analysis of pre-experiment conversation data, refining the framework. The final version contained 7 codes: planning, understanding, preparation, modeling, evaluation, technical inquiries, and tangential topics. For code breakdown, meanings, and case examples, consult Table 1.

To ensure coding framework reliability, two coders independently processed a randomly selected 15% of session data, defining content analysis units as uninterrupted student presentations. Results showed substantial agreement, Kappa score at 0.784. Coders discussed discrepancies to enhance shared understanding and scheme interpretation. Remaining data (8,615 sentences) was then independently coded, resulting in 5,986 coded data points.

Table 1. *Collaborative Conversation Data Coding Framework*

Code	Meaning	Cases
Plan	Develop plans to achieve the goals. Review models and processes, issues that arise, and identify areas for improvement in subsequent work.	For the next step, we may have to continue to extract the features. I'm missing a descriptive statistic for each type of question, such as the total value of the mcss type and its average value, and time.
Understand	Understand project goals and requirements, understand and infer data context and types, hidden trends, anomalies, patterns in data, explore data distribution, correlations and trends.	One is to determine whether he answered the question or not, and the other is to calculate the threshold value. These are the two problems that need to be solved in order to go ahead and make predictions based on these.
Preparation	Fixing data problems, converting data formats, merging multiple data sets, extracting new features from data, transforming features, selecting some data rows and features.	We need to standardize the test dataset. I roughly sieve a few features, you guys look at. We can count the type of operation. We can also add other columns, such as how many times to enter, how many times to exit, how many times to draw.
Modeling	Find, search, introduce, and filter algorithmic models that can be used, build models using cognitive tools, understand the meaning of model parameters, and change parameters to improve model performance.	Can you use generalized linear regression, but we're not a continuum here. There should be a template from when you wrote the logistic regression before. Take it and change it. I adjusted the tree model, because the previous 1500 may overfitting, so I set it to 1000, may be a little better.
Evaluation	Select evaluation metrics to assess model performance, calculate selected metrics using software manipulation or code programming, and understand	I see logistic regression can convert probabilistic output to labeled output, it can also output precision confusion matrix, exact value. How to use kappa for logistic regression? AUC is still 0.3536, Kappa is still 0, very low.

	what the evaluation results represent.	If kappa is equal to 0, it means that the two judgment results are caused by chance.
Technology problem	Only address technical issues that arise during interaction with RStudio and RapidMiner.	No matter how I access it, I cannot access its core content. So, I'm stuck here, stuck for a long time, stuck me for almost four or five hours.
Irrelevant topic	Distractions, or project related emotional expressions.	R language is too disgusting. I'm infected COVID-19 too.

2.5 Data Analysis

To address research question one, chi-square tests examined differences in frequency distribution across the seven dimensions between the two groups. For research question two, ENA method generated cognitive process patterns. Unit of analysis was cognitive tool type and group number, with section size set at 4. For research question three, descriptive statistics were initially performed. Data normality was assessed using Shapiro-Wilk, and homoscedasticity with Levene's method. Normally distributed data underwent independent sample t-tests, while non-normally distributed data were analyzed using Mann-Whitney U tests.

3. Results

3.1 What Are the Differences in the Frequency Distribution of Data Inquiry Cognitive Elements Between the Block-Based and the Text-Based Tools?

A total of 3062 codes were generated in the Block-based group, and 3148 codes were generated in the text-based group. Specifically, the block-based group had a significantly higher cognitive frequency than the text-based group in planning ($\chi^2 = 11.239$, $p < 0.01$) and understand ($\chi^2 = 164.512$, $p < 0.01$). Code frequencies for modeling ($\chi^2 = 102.044$, $p < 0.01$), evaluation ($\chi^2 = 88.256$, $p < 0.01$), and irrelevant topics ($\chi^2 = 83.195$, $p < 0.01$) were significantly higher in the text-based group. See Table 2 for more detailed results.

Table 2. *The Frequency Distribution of Cognitive Elements of Data Inquiry*

Code	block-based		text-based		χ^2	p
	(F)	(%)	(F)	(%)		
Plan	274	8.95%	190	6.04%	11.239	0.001
Understand	652	21.29%	259	8.23%	164.512	0
Preparation	1239	40.46%	1160	36.85%	0.034	0.866
Modeling	268	8.75%	514	16.33%	102.044	0
Evaluation	273	8.92%	499	15.85%	88.256	0
Technology problem	147	4.80%	117	3.72%	1.92	0.169
Irrelevant problem	209	6.83%	409	12.99%	83.195	0
Total	3062	100%	3148	100%		

3.2 What Are the Differences in the Cognitive Patterns of Data Inquiry Between the Block-Based and the Text-Based Tools?

The connection lines between the nodes represent the relationship between the data inquiry elements, and the thickness of the connection line reflects the number of times the two nodes appear simultaneously in the students' utterances. By observing the ENA diagrams and comparison plot of two groups, as shown in Figure 2, we found that the block-based group's plan-understand, understand-preparation are highly correlated. In the text-based group, evaluation frequently co-occurs with preparation, modeling, and technical problem. In addition to this, preparation is also more highly relevant to modeling and technical problem.

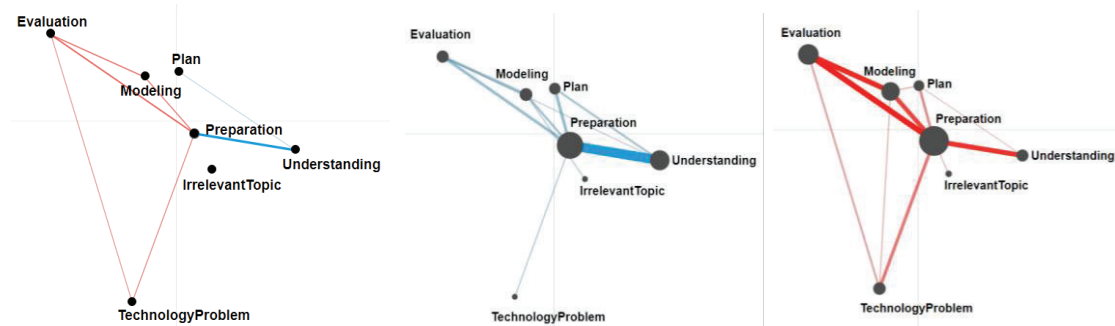


Figure 2. Comparison plot of block-based group and text-based group (left), ENA network of block-based group (middle), ENA network of text-based group (right).

3.3 Are There Any Differences in Learning Motivation, Cognitive Load, and Self-Efficacy Between the Block-Based and the Text-Based Tools?

The Block-based group exhibited a mean learning motivation of 76.000, while the text-based group had 72.667. T-test results ($t = 0.814$, $p = 0.427 > 0.05$) indicated no significant difference between them. Comparing cognitive load, the block-based group averaged 81.482, and the text-based, 83.704. Independent sample t-test ($t = -0.787$, $p = 0.438 > 0.05$) revealed both groups had higher cognitive load, but insignificantly different. However, self-efficacy showed significance, $U = 66.000$, $p = 0.047 < 0.05$, with block-based group at 72.889, and text-based group at 63.111 mean. See Table 3 for detailed results.

Table 3. Difference Test of Learning Motivation, Cognitive Load and Self-Efficacy

	Block-based		Text-based		T/U	p-value
	Mean	SD	Mean	SD		
Learning motivation	76.000	14.971	72.667	7.473	0.814 (T)	0.427
Cognitive load	81.482	7.987	83.704	7.753	-0.787 (T)	0.438
Self-efficacy	72.889	7.753	63.111	14.224	66.000 (U)	0.047

4. Discussion and Conclusion

This study extensively examined how distinct cognitive tools influence students' data inquiry skills by analyzing data inquiry patterns within block-based and text-based groups. Both tool types exhibit distinct merits in enhancing data comprehension and modeling.

4.1 Block-Based Cognitive Tools Make It Easier for Beginners to Enter Data Inquiry

The block-based group exhibited significantly higher cognitive frequencies for planning and understanding than the text-based group. Furthermore, ENA network connections, particularly between understand and preparation, were markedly stronger in the block-based group. Upon comprehending task objectives and data meanings, the block-based group strategized converting complex datasets into a streamlined format using the pivot module operator in RapidMiner. This transformation enhanced data comprehension, revealing trends like incomplete question responses, which the original dataset obscured. As understanding deepened, the group refined their inquiry plans. This iterative process, underpinned by solid data contextualization, facilitates pattern identification (Wilkerson, Lanouette, and Shareff 2021), emphasizing data inquiry's importance. Conversely, the text-based group encountered coding challenges, causing disruption between planning, understanding, and preparation. This resulted in decreased self-efficacy. Similar findings were noted by Price and Barnes (2015) in computational thinking comparisons. Novices, facing difficulties, may cease autonomous learning after few attempts (Tawfik, Payne, and Olney 2022).

4.2 Text-Based Group Used External Resources for In-Depth Modeling

Text-based group exhibited significantly higher frequencies of modeling and evaluation than block-based group. In-depth analysis revealed text-based group's model selection and parameter challenges, necessitating external resource usage (blogs, papers). This process enhanced their model understanding and internal models, aiding subsequent feature screening goals, as supported by Oh & Oh (2011). In contrast, block-based group relied on tool-provided modules within its interface, missing external learning opportunities. Although attempting more models, they lacked direction and understanding. Instructors must ensure resource exploration beyond tool for less reliance on blind attempts.

Cognitive load and learning motivation didn't significantly differ. However, working memory constraints and simultaneous module identification hinder effective attention allocation to modeling and optimization. Text-based group efficiently treats multiline text codes, reducing element interactivity (Chen, Kalyuga, and Sweller 2017) and cognitive load. Encapsulating multiple modules into one may ease resource consumption.

Acknowledgements

This paper was supported by the East China Normal University graduate student international conference special fund.

References

- Bart, Austin Cory et al. 2020. "Design and Evaluation of a Block-Based Environment with a Data Science Context." *IEEE Transactions on Emerging Topics in Computing* 8(1): 182–92.
- Chen, Ouhaio, Slava Kalyuga, and John Sweller. 2017. "The Expertise Reversal Effect Is a Variant of the More General Element Interactivity Effect." *Educational Psychology Review* 29(2): 393–405.
- Donoho, David. 2017. "50 Years of Data Science." *Journal of Computational and Graphical Statistics* 26(4): 745–66.
- Javadpour, Leili. 2022. "Using RapidMiner for Executing Queries and Visualization in a Traditional Database Course." *Journal of Education for Business* 97(4): 247–52.
- Jiang, Shiyao, and Jennifer Kahn. 2020. "Data Wrangling Practices and Collaborative Interactions with Aggregated Data." *International Journal of Computer-Supported Collaborative Learning* 15(3): 257–81.
- Lau, Kit-Ling., and John C. K. Lee. 2008. "Validation of a Chinese Achievement Goal Orientation Questionnaire." *British Journal of Educational Psychology* 78(2): 331–53.
- Leppink, Jimmie et al. 2013. "Development of an Instrument for Measuring Different Types of Cognitive Load." *Behavior Research Methods* 45(4): 1058–72.
- Oh, Phil Seok, and Sung Jin Oh. 2011. "What Teachers of Science Need to Know about Models: An Overview." *International Journal of Science Education* 33(8): 1109–30.
- Price, Thomas W., and Tiffany Barnes. 2015. "Comparing Textual and Block Interfaces in a Novice Programming Environment." In *Proceedings of the Eleventh Annual International Conference on International Computing Education Research*, Omaha Nebraska USA: ACM, 91–99. <https://dl.acm.org/doi/10.1145/2787622.2787712> (October 12, 2022).
- Tawfik, Andrew A., Linda Payne, and Andrew M. Olney. 2022. "Scaffolding Computational Thinking Through Block Coding: A Learner Experience Design Study." *Technology, Knowledge and Learning*. <https://link.springer.com/10.1007/s10758-022-09636-4> (May 19, 2023).
- Tsai, Chia-Lin, Moon-Heum Cho, Rose Marra, and Demei Shen. 2020. "The Self-Efficacy Questionnaire for Online Learning (SeQoL)." *Distance Education* 41(4): 472–89.
- Wilkerson, Michelle Hoda, Kathryn Lanouette, and Rebecca L. Shareff. 2021. "Exploring Variability during Data Preparation: A Way to Connect Data, Chance, and Context When Working with Complex Public Datasets." *Mathematical Thinking and Learning*: 1–19.
- Wirth, Rüdiger, and Jochen Hipp. 2000. "CRISP-DM: Towards a Standard Process Model for Data Mining." In *Proceedings of the 4th International Conference on the Practical Applications of Knowledge Discovery and Data Mining*, Manchester, 29–39.