

Article

Not peer-reviewed version

Teaching Data Science with Literate Programming Tools

[Marcus Birkenkrahe](#)*

Posted Date: 27 July 2023

doi: 10.20944/preprints202307.1847.v1

Keywords: Data Science; Literate Programming; Teaching, Emacs; Org-mode; IDE; Case Study



Preprints.org is a free multidiscipline platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Article

Teaching Data Science with Literate Programming Tools

Marcus Birkenkrahe ^{1,2}

¹ Lyon College, 2300 Highland Road, Batesville, AR 72501; birkenkrahe@lyon.edu

² Berlin School of Economics and Law, Badensche Str. 52, 10825 Berlin, Germany (on leave of absence)

Abstract: This paper presents a case study on using Emacs and Org-mode for literate programming in undergraduate data science courses. Over three academic terms, the author mandated these tools across courses in R, Python, C++, SQL, and more. Onboarding relied on simplified Emacs tutorials and starter configurations. Students gained proficiency after initial practice. Live coding sessions demonstrated the flexible instruction enabled by literate notebooks. Assignments and projects required documentation alongside functional code. Student feedback showed enthusiasm for learning a versatile IDE, despite some frustration with the learning curve. Skilled students highlighted efficiency gains in a unified environment. However, uneven adoption of documentation practices pointed to a need for better incorporation into grading. Additionally, some students found Emacs unintuitive, desiring more accessible options. This highlights a need to match tools to skills levels, potentially starting novices with graphical IDEs before introducing Emacs. Key takeaways include literate programming aids comprehension but requires rigorous onboarding and reinforcement; Emacs excels for advanced workflows but has a steep initial curve. With proper support, these tools show promise for data science education.

Keywords: data science; literate programming; teaching; Emacs; Org-mode; IDE; case study

1. Introduction

The author began teaching data science at a small liberal arts college in 2021, after a long career of teaching business informatics courses at a German business school. The COVID-19 pandemic prompted him to look for different ways of working with students in the classroom.

A couple of years before, the author had returned to an old friend from his days as a graduate student of physics, GNU Emacs, which is tersely described as an "extensible, customizable text editor" [1]. For programmers, GNU Emacs is an early tool for "literate programming," which assembles documentation, code, and output in one text document that can be converted to either source code for compilation or a document for printing [2]. For humans, literate programs are easier to understand, debug, and maintain.

In production data science, literate programming is the norm and not the exception, thanks to interactive notebooks that were first popularized by IPython and the Jupyter project [3]. Today, every development platform and IDE offers notebooks, such as Google Colaboratory, RStudio, Kaggle notebooks, and dozens of others. This is partly due to the interdisciplinary character of data science, which relies on programming, mathematics, and domain knowledge and skills, and which is used by practitioners with diverse backgrounds who are often not trained in computer science, such as biologists, psychologists, sociologists, and medical professionals.

The path to learning data science is highly fragmented. It involves a choice of programming languages (e.g., Python vs. R), data visualization techniques, mathematical and process modeling, and knowledge of computing infrastructure [4].

When teaching data science, a core problem is to give enough of an overview of all these aspects to enable students to solve real-world problems. In his classes, the author employs a triadic progression of didactic concepts: instruction, interaction, and immersion [5]. As a teacher, he had

always used the ancient GNU Emacs editor (created in 1975, first launched in 1985, and first used by me in 1991) and the more modern Org-mode package (created by Carsten Dominik in 2003), while the students used a variety of different platforms, such as GitHub for lecture materials, Google Colaboratory for coding along, and DataCamp and Canvas for assignments and tests.

Then, at EmacsConf in November 2021, the author watched Daniel German explain how he uses Org-mode to prepare and present teaching materials to his students when teaching programming in a variety of languages [6]. This emboldened the author to embark on an ambitious plan: to mandate Emacs and Org-mode as the central platform for interaction, instruction, and immersion in all his courses.

This was ambitious because Emacs is said to have a steep learning curve and to be useful mainly to a small group of devoted developers and professionals [7]. However, the Internet has experienced somewhat of an Emacs renaissance, with several popular YouTube channels and blogs featuring Emacs as a viable alternative to commercial products like VS Code by Microsoft or RStudio by Posit and making it more accessible [8–10].

After the author started using Emacs in class in spring 2022, students liked it, and so I kept going until the end of the spring term of 2023.

In this paper, the author will present and reflect on the choice of using Emacs and Org-mode as a mandatory literate programming tool for teaching data science in a variety of undergraduate courses for the programming languages R, SQL, SQLite, C/C++, bash, databases, data visualization, machine learning, and operating systems over three academic terms at a small college in rural Arkansas.

1.1. The theory and practice of data science teaching

Data science as a recent field of interdisciplinary study and practice is too young to have a well-established "best practice" for teaching [11]. However, there is basic agreement about what students need to learn in each of its disciplines: at the center of most courses and textbooks stands a workflow that begins with data and ends with storytelling. Between these two, data need to be stored, managed, transformed, modeled, and analyzed [12]. Figure 1 shows a shortened version of this flow with the four essential phases of any data science project - from cleaning data for analysis over modeling and visualization to presenting insights. This data science workflow leads to different job titles like "data engineer", "data analyst" or "machine learning scientist", who oversee different parts of the pipeline while "data scientists" can take charge of the whole pipeline.

Until a couple of years ago, the focus of data science education was on graduate level programs for people with undergraduate degrees in computer science, software engineering, mathematics, and statistics, or in a specific data domain like biology, psychology, or business [13]. As the demand for data science graduates increased, more and more undergraduate programs were springing up.

Mastering the data science workflow requires a lot of diverse skills. As an example: the data science major at Lyon College requires core competence in computer science, mathematics, and statistics, two data science specialization courses (e.g., visualization and machine learning), two domain specific specializations (social science/humanities, business/economics, or science). It includes training in R, Python, SQL, C/C++ or Java, and foundations of digital logic, database design and operating systems alongside shell scripting and version control with the Git version control software [14].

The teaching response to this ambitious set of skills has been to create infrastructures that integrate different tools so that students could focus on the data science task at hand - e.g., import and clean a dataset, get an overview of the dataset structure, and create some explorative visualizations.

Training platforms like DataCamp have perfected this approach, offering smooth sailing through different topics in only a few hours while hiding many of the tricky bits - like finding and loading suitable software packages, managing files and processes, mastering the interface between graphics and operating system etc. Integrated Development Environments (IDEs) like RStudio, VS Code, Google Colaboratory, Spyder, or Anaconda have gone a similar path. The same methodological

attitude is behind the recent trend to “code intelligence” - automatic comment generation and code completion enabled by transformer-based models [15], and behind some of the most popular textbooks on the market [16,17].

The problem with these integrated infrastructures is that they do not represent the systemic structural messiness of the real world, and they do not train students specifically in meeting real world requirements, especially transparency and reproducibility [18]. This would, for example, include having to struggle with setting up, managing, and debugging a multi-part work environment on one's own, and making it work with remotely with other developers, who may or may not share that setup.

Literate programming aligns well with the goals of reproducible, transparent, and open data science education. The literate programming methodology promoted through Emacs and Org-Mode allows students to unify theory, code, and output within a single coherent text document using only one single application while still being in full control of work environment and its customizations.

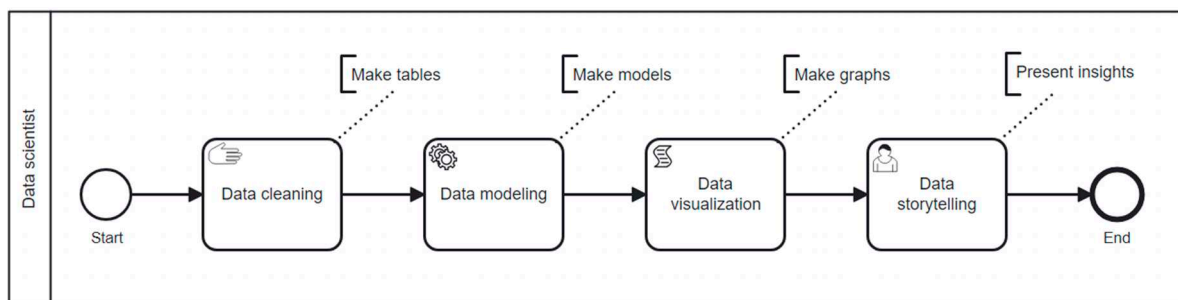


Figure 1. Simplified data science workflow.

2.1. The rationale for Emacs as an IDE and its learning curve

Emacs, as an Integrated Development Environment (IDE), provides a unique set of advantages, particularly within the realm of data science education [19]. Its high degree of customizability, extensive language support, and open-source nature underscore its utility. However, it's important to acknowledge that Emacs also introduces certain complexities that can make it challenging for newcomers to learn.

Emacs' customizability is a double-edged sword. On the one hand, it empowers users to tailor their programming environment to their specific needs - an ideal feature for a data science course that requires different programming tools and workflows. Yet, the flip side of this flexibility is that Emacs lacks a definitive 'out-of-the-box' setup. For beginners, this means they must invest significant effort into learning how to configure Emacs effectively. It's not just about learning a new tool; it's about shaping that tool to fit your needs.

Additionally, while Emacs' support for a wide range of programming languages adds to its versatility, it also necessitates the user to have a deeper understanding of these languages and their integration into the Emacs ecosystem. Navigating these integrations can prove challenging for novices, who may need to grapple not just with the intricacies of the languages themselves, but also with how to effectively set up and use Emacs to code in these languages.

Furthermore, while being open source offers advantages in terms of reliability and flexibility, it also introduces another layer of complexity. Users often must sift through an abundance of community-generated resources, deciphering what's relevant and reliable, which can be daunting and time-consuming for beginners.

Lastly, Emacs' extensibility, through packages such as Emacs Speaks Statistics (ESS) [20], designed for statistical programming and data analysis, solidifies its standing as a versatile tool for data science tasks. This extensibility, however, comes with a learning curve. The sheer number of available packages can be overwhelming to new users, requiring them to understand the utility and application of each one.

When using Emacs in class as a teacher, and especially when mandating its use by students, it is important to reign in expectations and provide a pre-configured environment. Figure 2 shows the minimal Emacs configuration file provided to the students. It enables them to:

1. Run code in Emacs in C/C++, R, SQL, SQLite, R, Python, and bash.
2. Update Emacs packages from a central repository.
3. Create code blocks easily with skeleton commands.
4. Auto-load the ESS package for using R in Emacs.
5. Disable toolbar and graphical menu bars (not being able to use graphical menus discourages use of the mouse and helps relying on the keyboard as the only and faster way to get things done).

```
(require 'ob-sqlite)
(require 'ob-sql)
(require 'python)
(require 'ob-emacs-lisp)
(require 'ob-R)
(require 'ob-C)
(require 'ob-shell)
(require 'ob-python)
(org-babel-do-load-languages
 'org-babel-load-languages
 '((R . t)
  (sql . t)
  (python . t)
  (emacs-lisp . t)
  (C . t)))
(setq org-confirm-babel-evaluate nil
      org-src-fontify-natively t
      org-src-tab-acts-natively t)
(require 'org-tempo)
(require 'package)
(add-to-list 'package-archives
 '("melpa-stable" . "https://stable.melpa.org/packages/"))
(global-set-key (kbd "<f6>") 'org-display-inline-images)
(global-set-key (kbd "<f7>") 'org-remove-inline-images)
(add-to-list 'load-path "~/emacs.d/elpa/ess-20221121.1627")
(load "ess-autoloads")
(tool-bar-mode -1)
(menu-bar-mode -1)
1-DD-\----F1 .emacs Top (28,19) (ELisp/d ivy ElDoc) -----
```

Figure 2. Emacs configuration file .emacs.

3.1. The rationale for Org-mode as a literate programming tool

Org-mode is a structured plain text format with notebook-like live code execution that offers an ideal platform for literate programming, a methodology that intermingles code, documentation, and output within a singular document [21]. Conceived by Donald Knuth in 1984, this practice promotes the creation of programs that are not just functional, but are also easy to understand, debug, and maintain. In the context of data science education, this form of programming plays a key role in unifying theory and practice, enabling students to visualize the results of their code parallel to the theoretical constructs being explained.

For instance, suppose a data science student is working on a machine learning project to predict housing prices. Using Org-mode, they can describe the theoretical concepts of their chosen regression model, then input the corresponding code, and finally display the generated outputs, all in one unified document. This seamless presentation not only makes it easier for the student to comprehend

the link between theory and implementation, but it also enhances the readability for others who might review or collaborate on the project.

Org-mode also bolsters the creation of reproducible research documents, a cornerstone of modern data science. Org-mode documents, integrating code, results, and narrative text, are ideal for assignments, projects, and collaborative work. For example, a student could use Org-mode to write a report on a data cleaning project. The report could contain blocks of R code for handling missing data, interspersed with explanations of why certain strategies were chosen. The output from the code (like summary statistics or visualizations) can be included directly below the code blocks, creating a comprehensive, easily understandable narrative.

Furthermore, Org-mode's capability to export these documents to various formats, including HTML, PDF, and LaTeX, eases the process of sharing work. A research team could collaborate on a data analysis project in Org-mode, then export the project as a PDF to share with their client, or as an HTML page to publish on the web. The LaTeX export option allows for creating formal academic articles, replete with features such as bibliographies and index creation.

Figure 3 illustrates the purpose of a literate program: the Org-mode file (left) contains the data story and the code with rich meta data. This file can be tangled into source code for compilation, and woven into a document that includes text, code, and output.

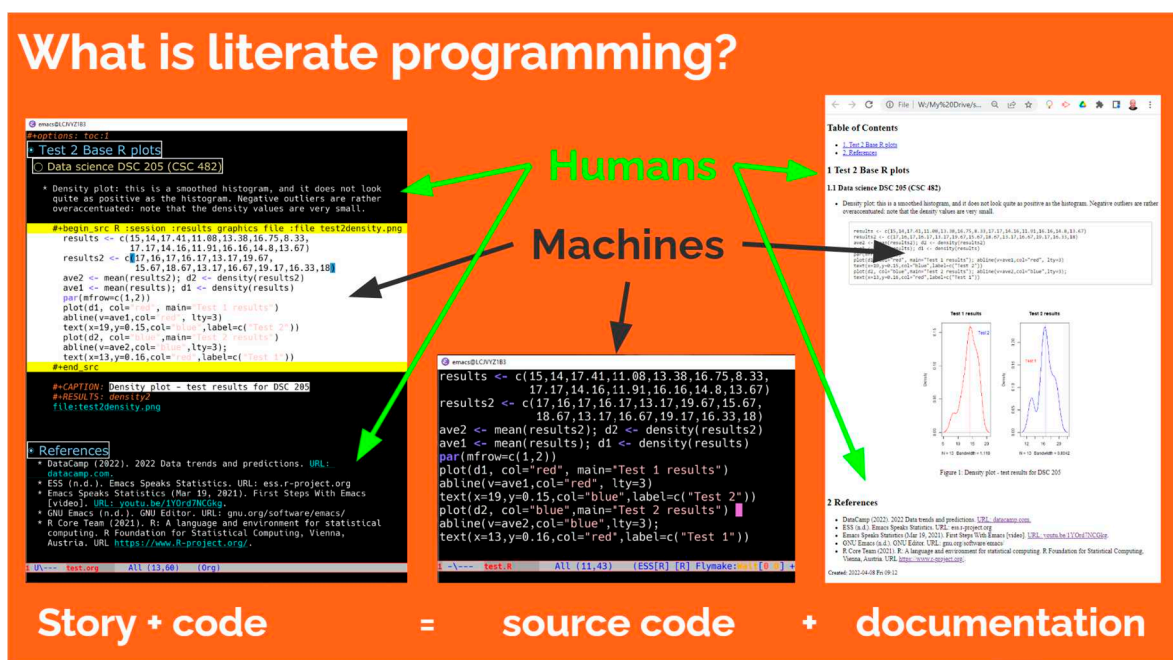


Figure 3. Literate programming is story + code = source code + documentation.

2. Methodology

This study is a design-based research (DBR) study, which seeks to improve educational practices by iteratively designing, implementing, and evaluating interventions [22,23]. The goal of this study is to develop evidence-based design principles for the use of literate programming tools in a series of linked courses.

The study was conducted in a real high-education setting, and the author of the paper was involved in all aspects of the study, from instigation to design to implementation to evaluation. The author was a participant-observer in the study, which means that they were both a participant in the courses and an observer of the students' learning process.

The evidence for the study was gathered through insights rather than statistical significance. This means that the author looked for patterns and trends in the students' learning, rather than trying to prove a specific hypothesis. The insights were gathered using systemic action research methods,

which are a type of qualitative research that focuses on understanding the complex interactions between people and their environment.

Here are some of the key features of the DBR methodology that are used in this study:

- The study is *iterative*, meaning that the design of the intervention was refined based on the findings of the evaluation.
- The study is situated in a *real-world setting*, which allows for the findings to be more generalizable.
- The study is *participatory*, meaning that the author was involved in all aspects of the study.
- The study is *qualitative*, meaning that the evidence is gathered through insights rather than statistical significance.

3. Case Study: Teaching data science with Emacs and Org-mode

This case study details the use of Emacs and Org-mode as central teaching tools across multiple undergraduate data science courses by the author. The following sections describe the specific courses involved, the student participants, how instruction was tailored to address knowledge gaps, the Emacs packages and configurations used, onboarding students onto these platforms, in-class teaching methods leveraging literate programming, assignment and project implementation, participant feedback and evaluations, and key challenges encountered along with lessons learned. Through this multi-term case study spanning a variety of classes and languages, the goal was to assess the potential benefits and limitations of mandating literate programming with Emacs/Org-mode in data science education.

3.1. Courses and participant profiles

Over the course of three academic terms (spring 2022, fall 2022, and spring 2023), the author used Emacs and Org-mode as the primary tools for teaching a variety of undergraduate courses. The class sizes varied between 6 and 28 participants. The participants included undergraduate students from different levels of their degrees, including first-year students, sophomores, juniors, and seniors. The students were from different majors, including computer science, mathematics, and engineering. None of the participants had had any contact with Emacs or Org-mode before.

Table 1 shows all courses involved in this case study with the language mainly used in the course, the course level and denomination (CSC = Computer Science Course, DSC = Data Science Course), and the number of participants in each course.

The material for all these courses is freely available in GitHub repositories at github.com/birkenkrahe (distributed under a GPL-3.0 copyleft license).

Table 1. Courses involved in the case study (CSC: CompSci, DSC: Data Science).

COURSE NAME (MAIN LANGUAGE)	LEVEL	WHEN	PARTICIPANTS
Intro to programming in C++ (C)	CSC 100	Spring 22/23	13/13
		Summer 22	6
Intro to data science (R)	DSC 105	Fall 22	13
Intro to advanced data science (R)	DSC 205	Spring 23	13
Digital humanities - text mining (R)	CSC 105	Spring 23	6
Database theory & applications (SQLite)	CSC 330	Spring 22	28
Data visualization (R)	DSC 302	Fall 22	15
Machine learning (R)	DSC 305	Spring 23	20
Operating systems (bash)	CSC 420	Spring 22	22
Applied math in data science (R)	DSC 482/MTH 445	Fall 22	20

3.2. Addressing specific student knowledge gaps

In several courses, the students had to answer an entry survey before the first session to help the author establish a baseline of what they already knew. The survey results revealed that most students did not know how to open a command line terminal, and that most students knew at most one programming language. These findings helped to tailor instructions to the specific needs of the students and help them develop a better understanding of their operating system and the computing infrastructure.

Specifically, the author addressed the following student needs to improve and enable basic computing infrastructure knowledge:

- **Opening and using the command line terminal (CLI).** A tutorial was created that showed students how to open the command line terminal and how to navigate the file system from the command line. Participants were also provided with exercises to practice these skills.
- **Exploring the file system from the command line.** The author showed the students how to use the command line to navigate the file system, including how to create, delete, and rename files and directories. He also showed them how to use the command line to search for files and directories.
- **Explaining and practicing CLI compilation** or, equivalently, running scripts using an interpreter. The author explained the difference between compilation and interpretation and showed students how to compile and run programs from the command line. He also provided them with exercises to practice these skills.
- **Using the shell and creating shell scripts in Emacs.** The author showed students how to use the shell to interact with the operating system, and how to create small shell scripts. He also showed them how to use Emacs to edit shell scripts.
- **Options for editing, executing, and debugging programs.** For instance, when it came to the statistical programming language R, the author introduced them to several alternatives. These alternatives included running R through a console in a terminal, utilizing the GUI provided by base R, accessing R online through platforms like Google Colaboratory and replit.com, utilizing DataCamp workspace with Jupyter Lab, using the RStudio IDE, and integrating R within Emacs, including running it in the background while executing an Org-mode code block.

3.3. Emacs version and packages used

The author's own Emacs configuration, though not strictly speaking minimal, did not require deep understanding of Emacs or Emacs-Lisp (the language used to configure and program the Emacs editor).

He used the vanilla GNU Emacs editor for Windows 10 (version 28.2 at the time of writing). Because it is easy to customize Emacs, a multitude of customization exist. In line with the general GNU philosophy, it is possible to adapt the Emacs to pretty much any workflow, aesthetic preference or keyboard and language setting [1].

In those courses where R was the primary programming language, I used *Org-mode* [21,24] in connection with the *ESS* package [20]. For version control, I used the *magit* package, integrated with GitHub [25].

For courses teaching C/C++, SQLite or bash, Org-mode with the built-in *Babel* extension for multi-lingual code was sufficient to run code [24]. Unlike all other interactive notebook environments, Emacs allows code blocks of different languages in one and the same file.

For classroom presentations, the author used the Emacs *org-present* package to render Org-mode files in presentation format, and the *modus-themes* package as the general Emacs theme. To make code clearer, I used the *rainbow-delimiters* package, and *org-appear* to hide emphasis markers used e.g., to highlight code [26–29].

Figure 4 shows an Emacs Org-mode buffer with a "Hello World" code block and corresponding output in (from the top): Python, SQLite, C, R and bash. This works only, of course, if the respective

languages are installed and if the environment is initialized properly so that Emacs can find the interpreters or compilers needed to execute the code.

```

#+begin_src python :results output
  print('Hello world')
#+end_src

#+RESULTS:
: Hello world

#+begin_src sqlite :db test.db :results output
  SELECT 'Hello world'
#+end_src

#+RESULTS:
: "Hello world"

#+begin_src C :results output :main yes :includes <stdio.h>
  printf("Hello world");
#+end_src

#+RESULTS:
: Hello world

#+begin_src R :results output
  print('Hello world')
#+end_src

#+RESULTS:
: [1] "Hello world"

#+begin_src sh :results output
  echo "Hello world"
#+end_src

#+RESULTS:
: Hello world

1 -\**~ babel.org All (30,20) (Org org-ai ivy)
Code block evaluation complete (took 0.2s).

```

Figure 4. Emacs Org-mode buffer with "hello world" program in 5 languages..

3.4. Tools used across all courses

When the author taught data science at a German business school, he used twelve different tools [30]. When he began to use Emacs for preparation, lecture, and practice, he could reduce the number of tools drastically:

GitHub. GitHub was the central repository for all course materials except quizzes, tests, and grade data. Git was fully integrated with Emacs, allowing me to work in different locations while maintaining a central, up to date material collection for the students. Students received minimal instruction in using GitHub (registration was not mandatory).

DataCamp. I used the DataCamp online classroom in all my courses for home assignments: students were required to complete lessons of relevant DataCamp courses (e.g., "Introduction to R", or "Understanding machine learning") on a regular basis.

Learning Management System. Participant activities in the LMS included weekly online quizzes and tests and home programming assignments whose solutions were submitted in the LMS.

The students could follow their personal progress at any time using the built-in, up-to-date gradebook. Both Schoology (Spring 2022) and Canvas (Fall 2022, Spring 2023) were used in this way.

Zoom. Zoom was used to share my screen with the students so that they could discern details e.g., when coding along, and to record each session for later viewing.

DataCamp, the Learning Management System (LMS) and Zoom remained decoupled from the Emacs-oriented workflow and were required to be accessed outside of class.

3.5. Teaching Emacs and Org-mode onboarding

A key part of the study was teaching students how to use Emacs and Org-mode, since most were unfamiliar with these tools. To facilitate onboarding, the author developed a simplified Emacs tutorial focused on the basics alongside a brief cheat sheet of the most common commands, including:

- Navigation and modes
- Managing files and buffers
- Customizing the interface
- Keyboard shortcuts

This hands-on tutorial was delivered interactively during class. Students followed along on their computers as the author demonstrated Emacs features followed by practice applying the skills through simple editing exercises.

Additionally, the author provided sample initialization and configuration files for Emacs and Org-mode. Students could use these as a starting point to tailor the tools to their needs. Throughout the term, troubleshooting help during classes and office hours were offered.

No separate Org-mode tutorial was used. Instead, Org-mode skills like adding metadata and using code blocks were taught through integrated examples during regular class programming exercises.

Though the students may not be conscious of markup methods, most had encountered the effects of separating content and layout instructions, e.g., in an HTML file. Hence, the way Org-mode meta data were used to control output and appearance, was not entirely foreign to them.

This combination of hands-on practice, custom configurations, and integrated learning was used to help participants quickly gain proficiency. By the second week of classes, most students were able to competently use Emacs and Org-mode for their assignments. Frequent reinforcement of skills also contributed to students' learning.

Figure 5 shows the top of the tutorial (as a Markdown file on GitHub) with the table of contents.

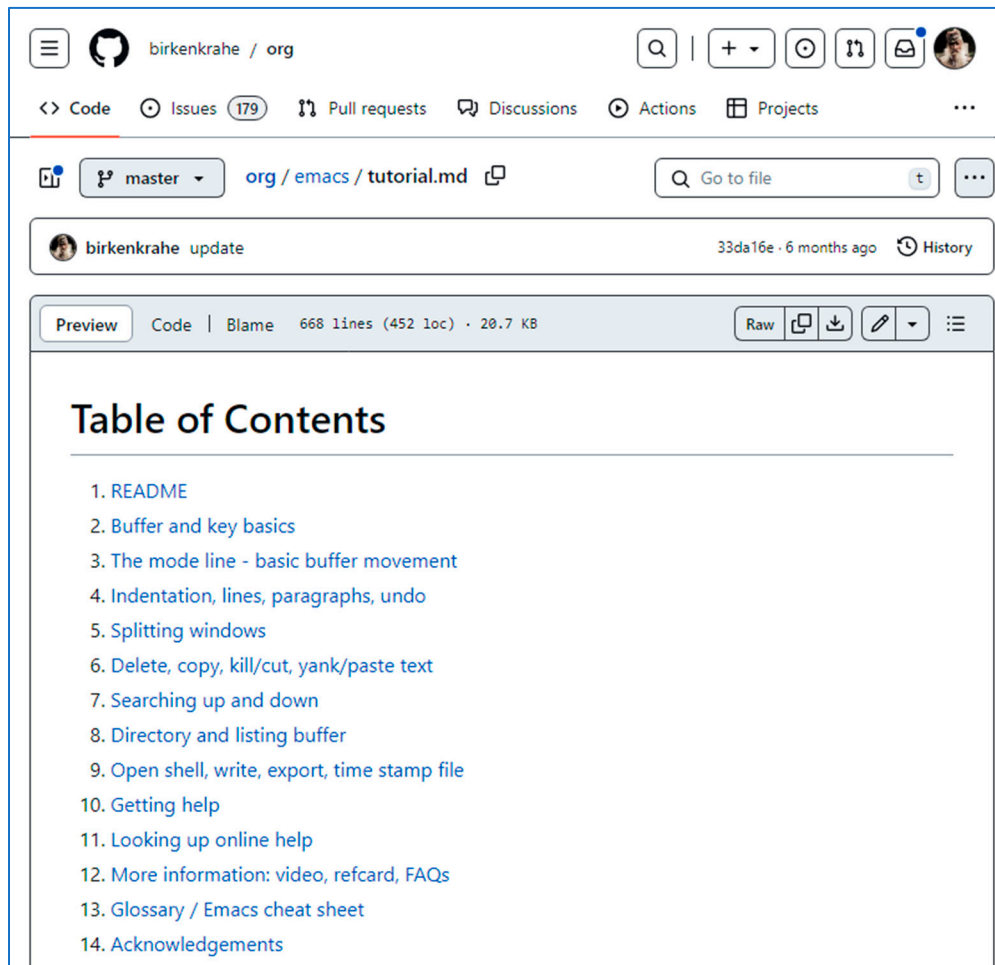


Figure 5. Table of contents of the Emacs tutorial on GitHub.

3.6. In-class Instruction

The course participants were instructed directly in Emacs using Org-mode files only. The author's Emacs screen showing an Org-mode file (aka notebook) was shared with the students (and recorded) via Zoom.

He would demonstrate practical skills and the students would use their own Emacs Org-mode files to either repeat them or solve simple exercises in class. Depending on the complexity of the material, the author would prepare practice files or ask the students to create their own Org-mode files from scratch.

Completed Org-mode files from class coding sessions could be uploaded to the LMS for credit and would regularly be checked and commented upon inside the students' notebooks.

Delivering live code presentations with Org-mode provided flexibility: the author could evaluate code blocks on the fly, modify examples and expand explanations using the literate programming approach. This interactivity would have been impossible to achieve using static, slide-based programs, and harder to do in an IDE where code and documentation (apart from program comments) are not in the same location.

Figure 6 shows the top of an Emacs Org-mode practice file for DSC 205 ("Introduction to advanced data science") used by the students in class in parallel to my lecture on user-defined functions in R. An R source code block was prepared where the students could put their solution. The file header includes file meta data (`#+TITLE`, `#+AUTHOR`, `#+SUBTITLE`), layout meta data (`#+STARTUP`, `#+OPTIONS`) and code block meta data - header arguments for an R code block with output to the screen whose code runs in an R session buffer named `*R*`. The `:noweb yes` argument means that noweb-style chunk substitution is enabled (though this feature was not used in this example) [31].



```

9_functions_practice.org - GNU Emacs at LCJ/VYZ1B3
#+TITLE:Writing Your Own Functions - PRACTICE FILE
#+AUTHOR: [yourname]
#+SUBTITLE:DSC205 Introduction to Advanced Data Science
#+STARTUP:overview hideblocks indent
#+OPTIONS: toc:nil num:nil ^:nil
#+PROPERTY: header-args:R :exports both :results output :session *R* :noweb yes
● README

This practice file accompanies the lecture on functions.

● TODO Identify and pledge yourself

1) In Emacs, replace the placeholder [yourname] at the top of this...
2) Go with the cursor on the headline and hange the TODO label to DONE
   by entering S-<right> ("Shift + right-arrow").

● TODO Example: hello, world!

- The function arguments are not workspace objects. Check that:...

- Modify hello_world - create a new function hello that takes a
  name as an argument and prints it to the screen:
  1) define a function named hello
  2) hello should have one argument, name
  3) return the name together with "Hello," using paste
  4) call the function with your name as the (string) argument
  5) check if name is in the list of user-defined objects using any
#+begin_src R
█
#+end_src

● TODO Example: Fibonacci sequence generator...
1 -(Unix)**- 9_functions_practice.org Top (37,0) Git:main (Org org-ai Ind i

```

Figure 6. student practice Org-mode file for classroom use (DSC 205).

3.6. Assessment

All courses were graded based on different, equally weighted categories shown in Table 2 as published in the syllabus:

- **Multiple-choice tests** (of no less than 10 questions per week) were administered weekly covering the week's topics and DataCamp assignments and reviewed at the start of every week.
- The **final exam** was a representative selection of those questions from the tests, in which the participants had performed worst.
- For **home assignments**, **In-class assignments** and **projects** see below.

Table 2. Grading system published in the syllabus document.

REQUIREMENT	UNITS	PPU	TOTAL	% of TOTAL
Final exam	1	100	100	20.
Home assignments	10	10	100	20.
Class assignments	10	10	100	20.
Project sprint reviews	5	20	100	20.
Multiple-choice tests	10	10	100	20.

TOTAL	500	100.
-------	-----	------

3.6.1. Assignments

There were three types of coding assignments:

1. Assignments consisting of completing a DataCamp lesson.
2. Home coding assignments (Org-mode notebooks).
3. In-class assignments (Org-mode notebooks).

The DataCamp lessons (for R and SQL) used a modified, customized platform consisting of a series of problems often preceded by an instructional video. The platform consisted of text, a code editor, and a console.

For home and in-class assignments, the students used Emacs as their text editor for writing and debugging code in Org-files. The answers to programming exercises had to be submitted as Org-mode files complete with documentation, code, sample input, sample output and references.

Figure 7 shows a programming assignment for the course "Introduction to programming in C++" including the creation of a literate program both in Org-mode and in "tangled" form (i.e., as C source code).

1. Write a program that prompts the user to enter a telephone number in the form `(xxx) xxx-xxxx`, and then displays the number in the form `xxx.xxx.xxxx`.
2. Example input/output of the first program, `phone1.c`:


```
Enter phone number [(xxx) xxx-xxxx]: (870) 456-7890
You entered: 870.456.7890
```
3. Write another program that asks for the input format in the form `xxx\xxx\xxxx`, and then displays the number in the form `(xxx)xxx-xxx`.
4. Example input/output of the second program, `phone2.c`:


```
Enter phone number [xxx\xxx\xxxx]: 870\456\7890
You entered: (870) 456-7890
```
5. Submit one Emacs Org-mode file `phone.org` with both programs in it as code blocks that can be tangled as `phone1.c` and `phone2.c`, respectively.
6. The header information of your Org-mode file should look like this:


```
#+TITLE: Phone number conversion
#+AUTHOR: [your name]
#+HONOR: pledged
#+PROPERTY: header-args:C :main yes :includes <stdio.h> :results output :tangle yes
```
7. Tip: some characters, like `\` are protected because they are part of the file `PATH`. If you want to use them, you have to "escape" them with an extra `\`, like the newline character `\n`. So to print (or to scan) the character `\`, you use `\\`.
8. [Here is a short video](#) (9 min) that explains in detail how to get started with this exercise in Emacs + Org-mode + C.

Figure 7. Assignment for an introductory course in C (CSC 100).

3.6.2. Projects

In addition to regular assignments, students in almost all courses had to complete a team-based term project using an adapted agile methodology called Scrum, with the author as the product owner (cp. Hidalgo, 2019). Teams of 2-3 students could pick a project topic from a provided list or choose their own idea relevant to the course material.

Throughout the term, teams presented their progress during four sprint review meetings, held approximately every 3-4 weeks. This allowed them to receive ongoing feedback and guidance.

The final project deliverable was a fully documented interactive report in the form of an Emacs Org-mode notebook. At the end of the term, each team also gave a presentation demonstrating their work to the entire class.

The author was able to provide input and assessment during the sprint reviews in addition to evaluating the final report. This helped simulate a real-world collaborative environment.

In spring 2022, no projects took place, and multiple students provided feedback asking to bring them back in future classes. The agile, team-based projects with regular feedback proved engaging and beneficial based on student responses.

The success of this adapted Scrum approach to team projects in computer science education was described in more detail in prior work [32]. The projects help reinforce practical application of concepts learned in class through hands-on collaborative problem solving.

The use of literate programming and interactive notebooks for the project reports provided unique benefits compared to traditional code-only deliverables. By having to interleave documentation, references, and outputs alongside functional code, students gained skill at communicating their technical work for an audience in a reproducible format. The requirement to attend to code quality, explanations, and results throughout the project lifecycle, not just at the end, improved depth of analysis and understanding. Teams had to consciously connect their design decisions and discoveries to academic sources.

Overall, the literate programming approach enhanced student learning and resulted in higher quality project work, validated by both my assessments and student feedback. The interactive element added unique value that static reports lack in demonstrating capabilities and engaging viewers.

Table 3 lists a selection of topics from projects in different courses that were completed and presented by participant teams:

Table 3. Selected term project presentation titles.

PROJECT TOPIC	COURSE	TERM
<i>Regression models on Twitter impressions</i>	DSC302	SP23
<i>Reactivity in R Shiny Dashboards</i>	DSC302	SP23
<i>Exploring data science salaries with ggvis</i>	DSC302	SP23
<i>Legality of AI art</i>	DSC305	SP23
<i>Langrangian Polynomial for regression</i>	DSC305	FA22
<i>ChatGPT in higher education</i>	CSC105	SP23
<i>An introduction to RcppArmadillo</i>	MTH445	FA22
<i>Spin rate in Baseball</i>	MTH445	FA22
<i>Introduction to Power BI</i>	DSC482	FA22

3.7. Participant evaluations

All courses were subject to anonymous participant evaluations in the anthology Course Evaluations survey "CourseEval" software. The survey report consisted of general 5-point Likert scale questions on course quality and instructor quality, followed by three open questions:

1. What are the best features of this course?
2. Do you have any suggestions for improvement of this course?
3. Do you have any additional comments on this course and/or the instructor?

The response rate was high with an average participation of 82.86% and detailed feedback on the use of Emacs in class. Tables 4 and 5 contains the answers that relate to the use of Emacs and/or the use of interactive Org-mode notebooks in the first two open questions. The negative answers show that the anonymous evaluation worked in the sense that students were not reluctant to share criticism. The answers also show that the overwhelming number of students, who answered, were happy with Emacs and Org-mode as literate programming tools.

Table 4. Student answers to “What were the best features of this course?”.

NO	WHAT ARE THE BEST FEATURES OF THIS COURSE?
1	"The use of Emacs. Learning such a powerful tool will, I feel, truly help me excel after college."
2	"The new possibilities and features that Emacs brings to not only my programming classes but all of them."
3	"The notebook[s] through emacs were very helpful."
4	"Emacs and working with .org files"
5	"Emacs was my favorite part about this course it was a challenging program to use but worth learning."
6	"The interactive notebooks. While the lecture is going on, being able to use what you are learning at the same time helps to remember the information later on."
7	"I like how we use Emacs to take notes and code along with the professor, so that we can practice the skills we learn in class."
8	"The use of emacs and the terminal over a GUI, while difficult at first, is more representative of tools we may use in the future and are very important skills to learn for competent use of computers. The introduction to emacs specifically is quite daunting, but it's useful for those that accept it."
9	"Being able to be hands on with the programs in class."
10	"The interactive notebooks that allow students to use what they are learning to use immediately."
11	"I enjoy the use of emacs as a literate programming environment."
12	"The constant practices we do during the semester."
13	"Getting to learn first-hand with pre-planned activities"
14	"Learning to use Emacs and GitHub."
15	"The emacs notebooks that allow you to use what you are learning in real time."
17	"I like the use of emacs because it is a versatile tool that might help us later in our careers."
18	"I like how interactive the course is."
19	"Being able to code alongside [the teacher]."
20	"Learning how to create data and move through my computer without touching my mouse."
21	"The instructor walks through the coding programs and makes sure you understand. Also stops if you are stuck."
22	"Using emacs as the text editor also helped me develop coding skills and get more familiar with managing files, switching between buffers, and other skills that will be fundamental in any career in the Computer Science area that I follow after college."

Table 5. Student answers to “Do you have any suggestions for improvement of the course?”

NO	DO YOU HAVE ANY SUGGESTIONS FOR IMPROVEMENT OF THE COURSE?
1	"I really did not like using Emacs. I felt like it complicated things more than it helped, and it took a couple weeks at the beginning of the semester to learn that we could have been learning C or C++."
2	"It would have been helpful to spend more time on explaining how EMACS works. Even now, I am not confident about my EMACS skills."
3	"Having a cheat sheet for emacs that focuses on the most common commands used for assignments."
4	"Emacs is a hassle to work in; since we have so many home assignments in R it would have been very useful to have spent class time downloading it, I missed several home assignments early in the year because I could not figure out how to download it."
5	"Emacs is a pain to use."
6	"Try and give the students more challeng[ing] notebooks to complete on their own."
7	"The installation of Emacs on our PCs was a little stressful on my end."
8	"eMacs is extremely clunky for me personally, and it crashed at least once every time I tried to work on things outside of class. I'd rather just use a normal compiler at that point. In addition to having to submit the sprint reviews as org-mode files, didn't feel great. I'd rather just submit them as a word document or even a text document. I know several others have also complained about eMacs to me personally."
9	"Emacs is a convoluted text editor. It would be nice if we used something else that was more streamlined."

3.7. Participant experiences

In addition to the anonymous evaluations, in Spring 2022, I asked students to comment on their experience with Emacs in class. The answers were positive throughout. Here are four exemplary statements (all seniors at the time; student names used with permission):

- "It was the first time I had ever used an editing software, and it will definitely have an impact on how I take notes/write code in the future." (Hunter Perkins)
- "I find [Emacs] very helpful. The ability to use it for multiple programming languages alone makes it powerful and worth learning." (Victor Noppe)
- "The ability to have just one app to code in all of these different languages with minimal setup is a breath of fresh air [...] a very useful tool despite its learning curve." (Jacob Wolfrom)
- "I learn best with examples and by doing, so when we started doing the Org mode notebooks in class, I really started to learn to program." (Spencer Rhoden)

4. Discussion and limitations

This chapter provides critical examination of the study results, including challenges faced and lessons derived, along with analysis of limitations. The first section details key difficulties encountered in the implementation of Emacs and Org-mode for teaching data science, as well as best practices identified that may inform future pedagogical approaches. The next section discusses the role of traditional literate programming in the era of low-code platforms and AI coding assistants, arguing its continued relevance despite technological shifts. Finally, the limitations of the present case study are discussed, including potential biases, the inherent complexity of the tools assessed, and restrictions of the evaluation methods employed. While valuable insights emerge, further

research is needed to address these limitations and gain a more comprehensive understanding of utilizing literate programming tools in data science education.

4.1. Challenges and lessons learned

Some key challenges from mandating Emacs as the primary, literate programming data science platform were:

- Steep learning curve of Emacs was an initial hurdle for many students. Some skipped the intro tutorial which slowed coding progress.
- Uneven student adoption of documentation practices in Org-mode. Some focused just on code rather than full literate programming.
- Most students did not utilize references in Org-mode assignments, despite academic importance.
- A portion of students found Emacs clunky or frustrating to use.
- Technical issues like Emacs crashing outside of class caused problems for some students.

The lessons learned included:

- Need to mandate and follow-up on Emacs tutorial completion early on.
- Documentation skills must be incorporated into grading rubrics.
- Instructor should model best practices in class materials and demos.
- Highlight exemplary student work to showcase capabilities of tools.
- Match tools to student coding skill levels. Emacs may be better for advanced courses after all.
- Ongoing reinforcement needed to overcome reluctance and ensure adoption.
- Overall, while powerful, Emacs and Org-mode had a steep initial learning curve. Adoption was mixed, though skilled students demonstrated the tools' potential.

Despite the challenges, using Emacs and Org-mode also had notable benefits for both teaching and learning. The students who took time to master the tools were able to work efficiently across multiple programming languages within a single environment. Emacs provided a unified framework for coding in languages like Python, R, C++ and more.

Additionally, the literate programming approach enabled through Org-mode increased student comprehension. Integrating code, results, and documentation into one document provided an interactive learning experience that reinforced concepts.

From the instructor perspective, Emacs and Org-mode facilitated flexible delivery of materials. Live coding demonstrations could be performed seamlessly, with the ability to evaluate and modify code blocks on-the-fly during lectures.

Overall, students who embraced Emacs and Org-mode highlighted the tools' versatility. They appreciated learning an environment common in industry. The instructor also benefited from streamlined lesson planning and presentation. With proper support, these tools demonstrate strong potential for enhancing data science education.

4.2. Literate programming in the age of low code and AI assistants

The advent of low code platforms and AI-powered coding assistants has transformed the programming landscape [33,34]. Tools like GitHub Copilot and services like Bubble.io abstract away much complexity, allowing those with little traditional coding knowledge to build applications through simple drag-and-drop interfaces and plain language commands.

In light of this shift towards égalisation and automation of coding, what role remains for traditional literate programming using tools like Emacs and Org-mode? At first glance, these may appear antiquated compared to cutting-edge AI technologies. However, literate programming still provides unique benefits that warrant its continued teaching: Most low code platforms and AI assistants take a functional approach, focused narrowly on generating executable code. Documentation is secondary. This limits developers' ability to understand and explain the inner workings of their programs.

In contrast, literate programming interweaves human-readable documentation with code in one cohesive document. This comprehensive encapsulation of process makes programs easier to comprehend, debug, maintain and share.

The hands-on nature of literate programming also builds deeper developer skills. Working directly in environments like Emacs engenders a firmer grasp of coding techniques and problem solving. AI assistants can expedite work, but relying solely on their suggestions deprives developers of learning opportunities.

Lastly, literate programming remains essential for many scientific and academic contexts where reproducibility and transparency are paramount. Computational essays and other reproducible documents require human-authored narrative combined with code.

Rather than displacing traditional coding knowledge, low code and AI can complement skills gained through literate programming. Just as calculator use strengthens rather than atrophies math skills, leveraging automation can free developers to focus on higher value tasks, informed by a deeper understanding of code.

In data science education, foundations in literate programming continue to provide lasting benefit. Emacs and Org-Mode teach vital concepts applicable across multiple programming languages and paradigms. This establishes a basis for effectively utilizing AI coding tools.

The landscape has changed, but literate programming's unique merits persist. It remains an essential component of a well-rounded data science education, laying groundwork for both automation-assisted workflows and reproducible research. The approaches synergize: literate programming supplies understanding, while AI accelerates application.

4.3. Limitations of the case study

While this study provides valuable insights into the use of literate programming tools in teaching data science, it is not without its limitations. One of the primary limitations is the potential bias in the sample selection. The study was conducted in a specific educational setting, and the participants were students who had chosen to study data science. Their motivation and interest in the subject might be higher than average, which could have influenced the results. Therefore, the findings might not be generalizable to all students or educational settings.

Another limitation is the inherent complexity and steep learning curve associated with literate programming tools like Emacs and Org-mode. While these tools offer numerous benefits, they require a significant investment of time and effort to master. This could potentially deter students who are new to programming or those who are looking for quicker, more straightforward solutions. The study did not fully explore the potential impact of this learning curve on student motivation and engagement.

Furthermore, the study focused primarily on the benefits of using literate programming tools in teaching data science, without fully exploring the potential drawbacks or challenges. For instance, issues related to software compatibility, technical glitches, or difficulties in understanding the syntax of these tools were not thoroughly examined. These factors could potentially affect the effectiveness of literate programming tools in a teaching environment.

Lastly, the study did not consider the impact of other pedagogical strategies or teaching aids that could be used in conjunction with literate programming tools. The use of these tools does not exist in a vacuum, and their effectiveness could be influenced by a variety of other factors, such as the quality of instruction, the use of supplementary materials, or the level of student-teacher interaction.

In conclusion, while this study provides a starting point for understanding the potential of literate programming tools in teaching data science, further research is needed to address these limitations and provide a more comprehensive picture of this complex issue.

5. Conclusions and outlook

The experiment of utilizing Emacs and Org-mode to teach various data science courses yielded mixed results. There were certainly benefits in terms of the versatility and power of these tools for skilled students. However, the initial learning curve proved prohibitive for some students.

The key advantage of this approach was the immersive coding environment enabled by Emacs and Org-mode. Students who embraced the tools were able to work efficiently across multiple languages like Python, R, and C++ within a single IDE. The literate programming methodology also enhanced their understanding by integrating code, results, and documentation.

However, this experience highlighted that Emacs and Org-mode may not be the optimal choice for all data science students, especially those new to programming. Despite best efforts at onboarding through tutorials and examples, some students struggled to become proficient. This hampered their learning during practical exercises.

Additionally, not all students adopted the discipline of documentation through Org-mode despite its importance for reproducible analysis. The motivation to fully utilize these tools was uneven.

Moving forward, the author plans to transition his teaching to tools that offer greater accessibility. For R and Python in particular, Jupyter Notebooks and Google Colaboratory seem better suited for getting novice students started with coding. These platforms remove friction while still supporting literacy programming.

For lower-level languages like C++, beginner-friendly, cross-platform IDEs may provide a gentler initial ramp. The graphical user interface is less intimidating but also delivers fewer insights.

While Emacs and Org-mode remain excellent options for advanced data science work, they may be better suited as secondary tools introduced after students have built core coding skills. The author still plans to cover these tools to some extent given their prevalence in the field. But they will no longer be mandatory for all coursework.

This evolution in his teaching reflects lessons learned about matching tools to students' skills levels.

Finding the optimal set of platforms is an ongoing process as new technologies emerge. But the lessons from this experiment will guide future decisions, with the goal of choosing tools best suited to each course and level. The priorities are supporting student learning and building practical data science skills.

Acknowledgments: The author would like to thank the students at Lyon College who participated in his courses and provided valuable feedback on the use of Emacs and Org-mode. Their willingness to learn these new tools was essential to this study. He also wishes to express his gratitude to the faculty and staff at Lyon College for their support during his visiting professorship from 2021-2023. The opportunity to teach a variety of data science courses and experiment with innovative approaches was invaluable. Finally, the author is grateful to the Berlin School of Economics and Law for granting him academic leave to pursue this visiting position. The time spent teaching and researching in the United States has expanded his perspectives and capabilities as an educator. The author received helpful assistance in preparing this paper from writing tools utilized ethically and appropriately.

References

1. Stallman, R.; Steele, G. GNU Emacs Manual; Version 28.2; Free Software Foundation, Inc.: Boston, MA, USA, 2022.
2. Knuth, D. Literate Programming. *Comput. J.* 1984, 27, 97–111.
3. Kluyver, T.; Ragan-Kelley, B.; Pérez, F.; Granger, B.; Bussonnier, M.; Frederic, J.; Kelley, K.; Hamrick, J.; Grout, J.; Corlay, S. et al. Jupyter Notebooks – a publishing format for reproducible computational workflows. In *Positioning and Power in Academic Publishing: Players, Agents and Agendas*; Loizides, F., Schmidt, B., Eds.; IOS Press: Netherlands, 2016; pp 87-90.
4. Adhikari, A.; DeNero, J.; Jordan, M.I. Interleaving Computational and Inferential Thinking: Data Science for Undergraduates at Berkeley. *Harv. Data Sci. Rev.* 2021, 3, 2.
5. Birkenkrahe, M. Building Graduate-Level, Gamified xMOOCs In Moodle. In *Proceedings of the EADTU – The Online, Open and Flexible Higher Education Conference, Hagen, Germany, 29-30 October 2015*; pp. 57-73.

6. German, D. Teaching with Org-mode. In Proceedings of the EmacsConf 2021, Online, 2 October 2021.
7. Johnson, T. Emacs as a Tool for Modern Science: The use of open-source tools to improve scientific workflows. *J. Johnson Matthey Technol. Rev.* 2022, 66, 122-129.
8. Prevos, P. Analyzing data science code with R and Emacs. *Opensource.com* 2020.
9. Taylor, D. Emacs: The Best Text Editor! YouTube, 2022.
10. Wilson, D. Emacs from Scratch. YouTube Playlist, 2022.
11. Davenport, T.H.; Patil, D.J. Data Scientist: The Sexiest Job of the 21st Century. *Harv. Bus. Rev.* 2012, 90, 70–76.
12. Wing, J.M. Computational Thinking's Influence on Research and Education for All. *Ital. J. Educ. Technol.* 2017, 25, 7–14.
13. Irizarry, R. Introduction to Data Science; CRC Press: Boca Raton, FL, USA, 2020.
14. Lyon College. Data Science Major. Available online: lyon.edu/data-science-major (accessed on 24 July 2023).
15. Li, Y.; Tarlow, D.; Brockschmidt, M.; Zemel, R. CodeBERT: A Pre-Trained Model for Programming and Natural Languages. *arXiv* 2022, arXiv:2002.08155.
16. VanderPlas, J. Python Data Science Handbook, 2nd ed.; O'Reilly Media, Inc.: Sebastopol, CA, USA, 2023.
17. Wickham, H.; Grolemund, G. R for Data Science, 2nd ed.; O'Reilly Media, Inc.: Sebastopol, CA, USA, 2023.
18. Davenport, T.H.; Patil, D.J. Is Data Scientist Still the Sexiest Job of the 21st Century? *Harv. Bus. Rev.* 2022.
19. Giorgi, F.M.; Ceraolo, C.; Mercatelli, D. The R Language: An Engine for Bioinformatics and Data Science. *Life* 2022, 12, 648.
20. Rossini, A.J. ; et al. Emacs Speaks Statistics (ESS) add-on package, Version 2023.
21. Guerry, B. ; et al. Org Mode, Version 2023.
22. Wang, F.; Hannafin, M.J. Design-based Research and Technology-enhancing Learning Environments. *Educ. Technol. Res. Dev.* 2005, 53, 5-23.
23. Hevner, A.R.; Park, S.T.; Ram, S. Design science in IS Research. *MIS Q.* 2004, 28, 77-105.
24. Schulte, E.; Davison, D.; Dominik, C. A Multi-Language Computing Environment for Literate Programming and Reproducible Research. *J. Stat. Softw.* 2012, 46, 1-24.
25. Lemmer-Weber, M. Using Programming Environments for Academic Research and Writing. *Archeomatica* 2021, 13, 30-31.
26. Lister, R. org-present-mode, Version 2014.
27. Stavrou, P. Modus Themes for GNU Emacs, Version 2023.
28. Rayman, J.; Linithien, F. rainbow-delimiters, Version 2021.
29. Istleyeva, A. org-appear, Version 2022.
30. Birkenkrahe, B. Teaching Data Science in a Synchronous Online Introductory Course at a Business School – A Case Study. In Proceedings of the Innovations in Learning and Technology for the Workplace and Higher Education, Online, 2021; pp. 17-29.
31. Ramsey, N. Literate programming simplified. *IEEE Softw.* 1994, 11, 97:105.
32. Gunnoltz, J.; Birkenkrahe, M. Bit-wise learning: a didactical–methodological design and evaluation concept for a mobile learning pilot project in junior enterprises. *Int. J. Innov. Educ.* 2013, 2, 70-85.
33. Porter, L.; Zingaro, D. Learn AI-Assisted Python Programming with GitHub Copilot and ChatGPT; Manning: Shelter Island, NY, USA, 2023.
34. Greengard, S. AI Rewrites Coding. *Commun. ACM* 2023, 66, 12-14.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.