

Data Management and Digital Archaeology

Andreas Angourakis

Tim Klingenberg

08 October, 2024

Table of contents

Course overview	5
Course schedule/ <i>Kursplan</i> :	5
Evaluation / <i>Kursbewertung</i>	5
 I Getting started	 6
1 Data and research data management	7
1.1 Introduction to Research Data	7
1.1.1 What is data? What is it for?	7
1.1.2 Archaeological Data: Particularities	7
1.1.3 Open Science	7
1.1.4 FAIR Principles	8
1.1.5 Open Source	8
1.1.6 Reproducibility	8
1.1.7 The lay of the land: Organizations, Institutions, and Where to Find Tools and Datasets	8
1.2 Introduction to Programming for Research	10
1.2.1 What is Programming?	10
1.2.2 Importance of Learning Programming	10
1.2.3 Overview of Common Programming Languages	10
1.2.4 Concept of Research Software: Tools and Scripts	11
1.2.5 Where to learn (and keep learning)	11
1.3 Hands-on Practice	12
 2 Git and GitHub	 13
2.1 Version-control and Git	13
2.2 GitHub	13
2.3 Hands-on Practice	14
 II Programming in R	 15
3 Introduction to R	16
3.1 Preparation	16
3.2 R syntax and workflow	16

3.3	Basic Data Structures in R	16
3.4	Data Manipulation in R	16
3.5	Data Visualization	17
3.6	(EXTRA)Interactive Visualizations	17
3.7	Hands-on Practice	17
4	Best practices in programming	21
4.1	Code Organisation	21
4.2	Writing Clean and Readable Code	21
4.3	Writing Efficient and Scalable Code	22
4.4	(EXTRA)Testing and Validation	22
4.5	(EXTRA)Error Handling and Debugging	22
4.6	(EXTRA)Code Reusability and Sharing	23
4.7	Hands-on Practice	23
III	Data Science in R	25
5	Count data and seriation	26
5.1	Introduction to Count Data in Archaeology	26
5.2	Basic Statistical Methods for Count Data	26
5.3	Introduction to Seriation	27
5.4	Using the <code>tesselle</code> Package for Seriation	27
5.5	Hands-on Practice	27
6	Compositional data	29
6.1	Introduction to Compositional Data in Archaeology	29
6.2	Basic Concepts in Compositional Data Analysis	29
6.3	Exploratory Data Analysis	30
6.4	Hands-on Practice	30
IV	Databases	31
7	Databases (I)	32
7.1	Introduction to Databases in Archaeology	32
7.2	Relational Database Concepts	32
7.3	Introduction to SQL (Structured Query Language)	33
7.4	Database Tools for Archaeology	34
7.5	Hands-on Practice	35
8	Databases (II)	36
8.1	Using GIS Databases for Archaeology	36
8.2	Hands-on Practice	36

V	GIS	38
9	GIS (I)	39
10	GIS (II)	40
	References	41

Course overview

040468 Datenmanagement und digitale Archäologie (ÜB, unterrichtet in Englisch) Übungen
(3 CP)

Time slot / *Zeitfenster*: Fr 14-16 Uhr c.t.

Place / *Ort*: Raum 2

Course instructors / *Kursleiter*: Andreas Angourakis / Tim Klingenberg

Support / *Unterstützung*: Thomas Rose

Course schedule/*Kursplan*:

	Date	Block	Topic
1	2024.10.18	Getting started	Data and research data management
2	2024.10.25	Getting started	Git and GitHub
3	2024.11.08	Programming in R	Introduction to R
4	2024.11.15	Programming in R	Best practices in programming
5	2024.11.22	Data Science in R	Data Science Workflow
6	2024.11.29	Data Science in R	Count data and seriation
7	2024.12.06	Data Science in R	Compositional data
8	2024.12.13	Databases	Databases (I)
9	2024.12.20	Databases	Databases (II)
10	2025.01.10	GIS	GIS (I)
11	2025.01.17	GIS	GIS (II)

Evaluation / *Kursbewertung*

(Attendance and final examination)

Part I

Getting started

1 Data and research data management

1.1 Introduction to Research Data

1.1.1 What is data? What is it for?

Data refers to factual information, often in quantitative or qualitative form, used as a basis for reasoning, discussion, or calculation. It serves as the foundational element for analysis and supports decision-making across diverse fields. In research, data helps to generate new insights, verify hypotheses, and contribute to the overall body of knowledge in a specific area.

1.1.2 Archaeological Data: Particularities

In archaeology, data encompasses various types such as field notes, artefacts, images, geospatial coordinates, and more. By systematically collecting, organizing, and analysing this data, researchers can reconstruct past human behaviours, understand environmental contexts, and explore cultural practices.

Archaeological data is unique due to its diverse formats and the complexity of its collection. It often includes material remains like pottery, bones, tools, and structures, as well as environmental data like pollen samples and soil types. Since archaeological data often comes from excavation sites, it is usually non-renewable; once excavated, a site cannot be restored to its original state.

Archaeologists rely heavily on careful documentation to preserve as much information as possible for future study. Furthermore, the context in which artefacts are found is crucial, as it helps interpret their use, significance, and the broader cultural setting.

1.1.3 Open Science

Open Science promotes transparency and collaboration in research by making methodologies, data, and findings accessible to the public and other researchers. In archaeology, this can involve sharing excavation reports, datasets, and analysis methods to facilitate wider understanding and scrutiny of findings.

1.1.4 FAIR Principles

FAIR stands for Findable, Accessible, Interoperable, and Reusable. These principles aim to enhance the usefulness of digital assets by ensuring they can be easily located, understood, and utilized by others. Applying FAIR principles in archaeology involves creating well-documented, open-access datasets that other researchers can readily use.

1.1.5 Open Source

Open Source refers to software and tools that are freely available for anyone to use, modify, and distribute. For archaeologists, open-source tools can offer affordable solutions for data analysis, visualization, and data management, promoting a more inclusive research environment.

1.1.6 Reproducibility

Reproducibility refers to the ability to replicate the results of a study using the original author's assets or following their methodology. It ensures that findings can be independently verified under similar conditions, reinforcing the reliability of scientific research (National Academies of Sciences et al. 2019).

In scientific research, reproducibility is crucial for confirming the validity of experimental findings and building upon existing knowledge. It enhances transparency and trustworthiness in scientific practices, promoting better peer review and collaboration “GRN · German Reproducibility Network” (n.d.).

Most research in archaeology does not necessarily involve controlled experiments and archaeological survey and excavation are destructive, thus unrepeatable. However, the reproducibility of data collection, processing and analysis is not a trivial concern.

Making research reproducible must be considered as a spectrum of practices, in which researchers should thrive for doing better, despite the challenges (Marwick 2017).

1.1.7 The lay of the land: Organizations, Institutions, and Where to Find Tools and Datasets

Several organizations and platforms support Open Science and the use of FAIR and Open Source principles, providing archaeologists with access to valuable tools and datasets:

- Research Organizations:

- [OpenAIRE](#): OpenAIRE AMKE is a non-profit organization with a mission to promote open scholarship and improve discoverability, accessibility, shareability, reusability, reproducibility, and monitoring of data-driven research results, globally (Iatropoulou n.d.).
 - [Go FAIR Initiative](#): Provides guidelines on implementing FAIR principles, with resources tailored for researchers in various fields, including archaeology (“FAIR Principles” n.d.).
 - [The Turing Way](#): collaborative writing of an online handbook (Community 2022).
- Data Repositories and Tools:
 - [Zenodo](#): A repository where researchers can deposit datasets, software, and publications.
 - [Open Science Foundation](#): OSF is a free, open platform to support your research and enable collaboration.
 - [Archaeology Data Service \(ADS\)](#): An archive for archaeological data from the UK, which provides access to various datasets, including excavation reports and geospatial data.
 - [Open Context](#): A platform offering access to archaeological data from various global sources, adhering to FAIR principles [noauthor_open_nodate-1].
 - [GitHub](#): GitHub is a developer platform that allows developers to create, store, manage and share their code through the use of Git software and a series of additional automated services.
 - Open Source Tools:
 - [Markdown](#): Markdown is a lightweight markup language designed for creating formatted text using a plain-text editor (“Markdown Guide” n.d.).
 - [R](#) and [Python](#): Programming languages with extensive libraries for data analysis, which are widely used in archaeology for statistical analysis and modeling (R Core Team 2024, noauthor_welcome_2024). See more about these and other programming languages below.
 - [RStudio](#): RStudio IDE is an integrated development environment for R, a programming language for statistical computing and graphics.
 - [Quarto](#): An open-source scientific and technical publishing system.
 - [Git](#): Git is a distributed version control system designed for tracking changes in source code during software development.
 - [Zotero](#): Zotero is a free, easy-to-use tool to help collect, organize, annotate, cite, and share metadata on references.

- [QGIS](#): A free, open-source GIS tool useful for mapping and spatial analysis in archaeology (“Spatial Without Compromise · QGIS Web Site” n.d.).

By utilizing these resources, archaeologists can ensure that their research aligns with Open Science, FAIR, and Open Source principles, ultimately enhancing the transparency, accessibility, and longevity of their work.

1.2 Introduction to Programming for Research

1.2.1 What is Programming?

Programming is the process of designing and implementing instructions that a computer can follow to perform specific tasks. It involves writing code in various programming languages that communicate with the computer’s hardware and software to solve problems, process data, and automate repetitive tasks. At its core, programming is about creating a sequence of steps, known as algorithms, which help achieve a desired outcome (“What Is Programming? And How To Get Started” 2024).

1.2.2 Importance of Learning Programming

For researchers, learning programming offers several significant advantages:

- *Efficiency and Automation*: Programming can help automate data collection, processing, and analysis, saving time and reducing human error. It also enables researchers to handle large datasets and complex calculations with ease.
- *Reproducibility*: Writing scripts to perform analysis allows other researchers to replicate experiments, thus ensuring results can be verified and reproduced, a fundamental aspect of scientific research.
- *Access to Powerful Tools*: With programming skills, researchers can access a wide range of tools for data visualization, statistical analysis, machine learning, and simulation. These tools can enhance the scope and quality of research projects.

1.2.3 Overview of Common Programming Languages

Several programming languages are popular in research due to their specific features and libraries (“Introduction to Programming Languages” 2018):

- *Python*: Known for its readability and versatility, Python is widely used for data analysis, machine learning, and automation. It has extensive libraries such as NumPy, pandas, and matplotlib, which are particularly useful for data-intensive research.

- *R*: A language specifically developed for statistical computing and graphics, R is preferred in data science, bioinformatics, and fields requiring extensive statistical analysis. The Comprehensive R Archive Network (CRAN) offers thousands of packages that can handle a range of analytical tasks.
- *MATLAB*: Commonly used in engineering and scientific research, MATLAB excels at numerical computing, simulation, and algorithm development. It is particularly popular in fields like physics, engineering, and finance.
- *JavaScript*: While primarily a web development language, JavaScript is also used in research for developing interactive data visualizations and web-based applications.

1.2.4 Concept of Research Software: Tools and Scripts

Research software comprises tools, libraries, and scripts that aid in conducting and managing research activities. It can range from simple scripts for data cleaning and preprocessing to complex software packages for statistical analysis and simulation. Although any software used in research could be considered as research software, advance training focus on programming or scripting skills, not on graphical user interface operations (e.g., creating a plot in R or Python rather than in Microsoft Excel).

- *Tools*: Research software tools like [Jupyter Notebooks](#), [RStudio](#), [PyCharm](#), [Visual Studio Code](#), etc., offer integrated development environments (IDEs) tailored to a specific range of languages, enhancing productivity and facilitating code sharing and version control.
- *Scripts*: Scripts are text files encoding programs written to perform specific tasks. In research, scripts are often used for data cleaning, model training, and result visualization. Scripts can be used even to generate and format an entire dataset (have a peek on how the [course schedule has been built with R code](#)). These scripts can be shared among researchers to ensure that analyses are reproducible and consistent.

Learning programming for research allows scientists to leverage these tools and scripts effectively, making research more efficient, reproducible, and insightful.

1.2.5 Where to learn (and keep learning)

- [Datacamp](#) (“Learn R, Python & Data Science Online” n.d.)
- [Udemy](#) (“Online Courses - Learn Anything, On Your Schedule | Udemy” n.d.)
- [Coursera](#) (“Coursera | Degrees, Certificates, & Free Online Courses” n.d.)
- Courses and workshops by high education institutions (e.g., [RUB Research School](#))

1.3 Hands-on Practice

- Task 1: register yourself as user at:
 - [GitHub](#)
 - [Zenodo](#)
 - [Zotero](#) (optional)
- Task 2: installations:
 - [Git](#)
 - [GitHub Desktop](#)
 - [R](#)
 - [RStudio](#)
 - [QGIS](#)
- Task 3: explore datasets at [Open Context](#)
- Task 4: have a go with writing in markdown at [Markdown Live Preview](#).

2 Git and GitHub

2.1 Version-control and Git

- Version control: general concept and its usefulness
- Git software (installation check, but preferably done before)

2.2 GitHub

- GitHub: creating a user
- Introduction to workflows: local and on web browser
- Introduction to Markdown (GitHub-flavoured)
- Best practices and conventional files:
 - Folders:
 - * `source` or `src` (source code)
 - * `doc` or `docs` (documentation)
 - Version tags
 - Files:
 - * `README.md`: Provides an overview of the project, including what it does, how to set it up, and how to contribute. A few sections examples are:
 - General description
 - Authors and/or contributors
 - Acknowledgments

- Funding
 - Installation or use instructions
 - Contributing
 - * LICENSE: Specifies the terms under which the code can be used, modified, and distributed.
 - * CITATION.cff: human- and machine-readable citation information for software (and datasets). See example [here](#).
 - * .gitignore: Lists files and directories that Git should ignore, such as build outputs and temporary files.
 - * CHANGELOG.md: Logs a chronological record of all notable changes made to the project, often following conventions like Conventional Commits.
 - * References.bib: one or more files containing references cited within the markdown files of the repository.
- Releases and GitHub-Zenodo connection

2.3 Hands-on Practice

- Task 1: create your own repository in GitHub. Add files, change files, commit (online). Edit your README file and create other conventional files.
- Task 2: fork our repository, clone, modify, commit, pull request back to ours (online)
- Task 3: Set up a local workflow and clone your fork. Modify a file and commit. Push into your fork in GitHub. Create a pull request to merge it with ours.
- Task 4: Setting up the GitHub-Zenodo connection. Publishing your repository and updating it with its DOI.

Part II

Programming in R

3 Introduction to R

3.1 Preparation

- Installing R and RStudio.
- Overview of RStudio interface.

3.2 R syntax and workflow

- Basic R syntax: variables and data types.
- Writing and executing R scripts.
- Arithmetic operations, logical operations in R.
- Algorithm structures: `if`, `else`, `while`, `function`
- Packages. Mention the most important packages used later, including `tidyverse` and `tesselle`.

3.3 Basic Data Structures in R

- Vectors, matrices, data frames, lists.
- Basic operations on data structures (indexing, subsetting, adding/removing elements).

3.4 Data Manipulation in R

- Importing data: reading data from CSV files, using canonical datasets (`iris`, `archdata::DartPoints`).

- Basic data manipulation using base R (`subset`, `merge`, `apply` functions).
- Introduction to `dplyr` package for data manipulation (filtering, selecting, mutating data).

3.5 Data Visualization

- Introduction to plots: histograms, bar plots, scatter plots.
- Creating plots in R with base R graphics.
- Creating multiple plot figures with `layout`.
- Creating plots in `ggplot2`.
- Creating multiple plot figures with `gridExtra::grid.arrange`.
- Base R graphics and `ggplot2`: comparative
- Saving plots: open and closing graphic devices in R.

3.6 (EXTRA)Interactive Visualizations

- Introduction to creating interactive visualizations.
- Example: Building an interactive plot using `plotly` and `knitr`.

3.7 Hands-on Practice

- Create a new project in RStudio, placing it at the root directory of your own repository (cloned local branch).
- Create a `data.frame` named “stone_tools_data” directly in R with the following characteristics (based on Carlson 2017, p. 26):
 - Set of six stone tools with inventory number
 - Recording of dimensions (length, breadth, thickness), material type, and whether the material is local or non-local.

- Data per object:
 - * LN15:
 - Length: 18
 - Breadth: 9
 - Thickness: 3
 - Material type: chert
 - Material provenance: local
 - * LN17:
 - Length: 14
 - Breadth: 7
 - Thickness: 2
 - Material type: chert
 - Material provenance: local
 - * LN18:
 - Length: 21
 - Breadth: 10
 - Thickness: 3
 - Material type: obsidian
 - Material provenance: local
 - * LN21:
 - Length: 14
 - Breadth: 7
 - Thickness: 3

- Material type: chert
- Material provenance: non-local
- * LN23:
 - Length: 17
 - Breadth: 8
 - Thickness: 3
 - Material type: obsidian
 - Material provenance: local
- * LN24:
 - Length: 16
 - Breadth: 8
 - Thickness: 2
 - Material type: obsidian
 - Material provenance: non-local
- Check that the data and data types are coherent with the specifications. Save it as a CSV file and load it back as a new R object (e.g. “stone_tools_data2”). Compare.
- Create a plot showing the counts of objects made of chert and obsidian. Save it as a PNG file.
- Create a new variable (“type_and_provenance”) that combines type and provenance and create a plot showing the counts in each category. Save it as a PNG file.
- Create a single figure displaying the variable distribution of the three dimensions measured. Save it as both a PNG and a SVG file.
- Create a plot displaying the relationship between length and breadth. Save it as a PNG file.

- Create a plot displaying the relationship between length and breadth, this time marking (point type, colour) objects by their “type_and_provenance”. Save it as both a PNG and a EPS file.
- (EXTRA) Create a figure to help explore the question: Do stone tools of different material and provenance tend to be of different size?
- Commit all changes and push to the remote using RStudio.
- Q&A and troubleshooting.

4 Best practices in programming

4.1 Code Organisation

- **Modular Programming**
 - Importance of modularity: breaking down code into functions and modules.
 - Example: Creating in-script custom functions.
 - Example: Creating and importing custom R scripts.
- **Code Structuring**
 - Structuring a data science project: folder organization, separating code, data, and outputs.
 - Example: Setting up a basic project structure in R and RStudio.

4.2 Writing Clean and Readable Code

- **Naming Conventions**
 - Using meaningful and consistent names for variables, functions, and files.
 - Example: Best practices in naming conventions in R ([tidyverse style guide](#)).
- **Commenting and Documentation**
 - Importance of comments and inline documentation.
 - Example: Writing using `roxygen2` in R for documenting functions.
- **Avoiding Magic Numbers and Hardcoding**

- Use of constants and configuration files.
- Example: Using constants in R.

4.3 Writing Efficient and Scalable Code

- **Vectorization**
 - Avoiding loops by using vectorized operations for efficiency.
 - Example: Implementing vectorized operations in R (base R, `dplyr`).
- **Memory Management**
 - Managing memory usage and avoiding memory leaks.
 - Example: Best practices for handling large datasets in R (using `data.table`).

4.4 (EXTRA)Testing and Validation

- **Writing Unit Tests**
 - Importance of testing: ensuring code correctness.
 - Example: Writing basic unit tests in Python (`unittest` or `pytest`) and R (`testthat`).
- **Data Validation**
 - Validating data inputs and outputs, ensuring data integrity.
 - Example: Implementing data validation checks in data processing scripts.

4.5 (EXTRA)Error Handling and Debugging

- **Error Handling Techniques**
 - Using `tryCatch` in R.

- Writing meaningful error messages.
- Example: Implementing error handling in a data processing script.
- **Debugging Tools**
 - Introduction to debugging tools: `browser` in R.
 - Example: Walkthrough of a debugging session in R.

4.6 (EXTRA)Code Reusability and Sharing

- **Creating Reusable Code**
 - Writing functions and libraries for reuse across projects.
 - Example: Creating a simple R package.
- **Sharing Code**
 - Sharing code with others: publishing packages, sharing notebooks.
 - Example: Publishing an R package on CRAN/GitHub.

4.7 Hands-on Practice

- **Refactoring Code (30min)**
 - First attempt: Refactoring a sample script to follow best practices (clean code, modularity, documentation).
 - Second attempt: try using a Large Language Model (LLM) to refactor.
- **Collaborative Exercise (40min)**
 - Simulating a collaborative workflow with Git: making and reviewing pull requests.
 - Groups of two or three

- Re-use one of the repositories created in GitHub (Session 2) and populated by R code and output files (Session 3).
 - Mutual reviews and change suggestions.
 - Discussion, accepting/rejecting changes, and merge decision
- **Open discussion (10min)**
 - Addressing common challenges in applying best practices to real-world projects.

Part III

Data Science in R

5 Count data and seriation

5.1 Introduction to Count Data in Archaeology

- **Understanding Count Data**
 - Definition and significance of count data in archaeological contexts (e.g., artifact counts, feature frequencies).
 - Example: Overview of typical archaeological datasets involving count data.
- **Challenges in Analysing Count Data**
 - Issues with skewness, overdispersion, and zero inflation.
 - Example: Common problems encountered in archaeological count data analysis.

5.2 Basic Statistical Methods for Count Data

- **Poisson and Negative Binomial Distributions**
 - Introduction to Poisson distribution and its application to count data.
 - Example: Fitting a Poisson model using `glm` in R.
 - Introduction to the Negative Binomial distribution for overdispersed data.
 - Example: Fitting a Negative Binomial model using `MASS::glm.nb`.
- **Goodness-of-Fit Testing**
 - Assessing the fit of count data models.
 - Example: Performing a chi-square goodness-of-fit test in R.

5.3 Introduction to Seriation

- **What is Seriation?**
 - Overview of seriation techniques and their importance in archaeology for ordering artefacts or sites chronologically.
 - Example: Historical applications of seriation in archaeology.
- **Seriation Techniques**
 - Introduction to different seriation methods (e.g., frequency seriation, contextual seriation).
 - Example: Basic seriation using traditional methods.

5.4 Using the `tesselle` Package for Seriation

- **Introduction to `tesselle`**
 - Overview of the `tesselle` package and its tools for seriation.
 - Example: Installation and loading of `tesselle`.
- **Practical Seriation in R**
 - Performing seriation.
 - Example: Applying seriation to an archaeological dataset (e.g., pottery styles, stratigraphic data).
- **Visualizing Seriation Results**
 - Creating visual representations of seriation results.
 - Example: Plotting seriation outputs.

5.5 Hands-on Practice

- **Case Study: Seriation of Archaeological Artefacts**

- Step-by-step walkthrough of a seriation analysis using count data.
 - Example: Seriation of pottery fragments or lithic tools using `tesselle`.
- **Q&A and Troubleshooting**
 - Addressing common issues in count data analysis and seriation in archaeological contexts.

6 Compositional data

6.1 Introduction to Compositional Data in Archaeology

- **Understanding Compositional Data**
 - Definition and examples of compositional data in archaeology (e.g., proportions of different materials, chemical compositions).
 - Example: Typical archaeological datasets that include compositional data.
- **Challenges in Analysing Compositional Data**
 - Issues with relative proportions and the “closed” nature of compositional data.
 - Example: Limitations of traditional statistical methods on compositional data.

6.2 Basic Concepts in Compositional Data Analysis

- **Overview of methods in Multivariate statistics**
- **Log-Ratio Transformations**
 - Introduction to log-ratio transformations (CLR, ILR, ALR) for compositional data.
 - Example: Applying a centred log-ratio (CLR) transformation in R (e.g. `nexus::transform_clr`).
- **Visualizing Compositional Data**
 - Techniques for visualizing compositional data (ternary plots, bar charts).
 - Example: Creating a ternary plot using the `isopleuros::ternary_plot` and `ggtern`.

6.3 Exploratory Data Analysis

- **Principal Component Analysis (PCA)**
 - Introduction to PCA tailored for compositional data.
 - Example: Performing PCA on compositional data using `tesselle::nexus::pca`.
 - Biplots and screeplots
 - Example: Plotting PCA results as a biplot and add visualisation aids using `tesselle::dimensio` functions.
- **Cluster Analysis**
 - Overview of clustering techniques for exploring groupings in compositional data.
 - Example: Applying hierarchical clustering on transformed compositional data using `tesselle`.
 - `ggplot2`: dendrograms with `ggraph`
- **Correspondence Analysis**
 - Correspondence Analysis for compositional data.
 - Example: Performing Correspondence Analysis using `tesselle::ca`.

6.4 Hands-on Practice

- **Case Study: Analysis of Compositional Data from Archaeological Sites**
 - Step-by-step walkthrough of an exploratory analysis of compositional data.
 - Example: Analysing chemical compositions of ceramics or soils using `tesselle`.
- **Q&A and Troubleshooting**
 - Addressing challenges in visualizing and analyzing compositional data in archaeological research.

Part IV

Databases

7 Databases (I)

7.1 Introduction to Databases in Archaeology

- **What is a Database?**
 - Definition and importance of databases in archaeology (data storage, retrieval, and management).
 - Example: Common uses of databases for managing archaeological data (e.g., artifacts, excavation records).
- **Types of Databases**
 - Overview of relational vs. non-relational databases.
 - Example: Advantages of relational databases for structured archaeological data.

7.2 Relational Database Concepts

- **Basic Concepts**
 - Introduction to tables, records, fields, primary and foreign keys.
 - Example: Organizing excavation data in relational tables (e.g., site locations, stratigraphy, finds).
- **Normalization**
 - Introduction to database normalization (avoiding redundancy, ensuring data integrity).
 - Example: Structuring artifact data into normalized tables (e.g., artifact type, material, condition).

- **Entity-Relationship (ER) Models**

- Creating ER models to visualize relationships between archaeological datasets.
- Example: Designing an ER diagram for an archaeological database (e.g., site, context, finds).

7.3 Introduction to SQL (Structured Query Language)

- **Basic SQL Commands**

- Overview of SQL syntax and common commands (e.g., `CREATE`, `INSERT`, `SELECT`, `UPDATE`).
- Example: Creating tables for archaeological data and inserting records.

- **Creating a Database with SQL**

- Step-by-step process of building an archaeological database using SQL.
- Example: Creating an artifact catalogue with fields such as artifact ID, type, material, and context.

- **Indexing and Keys**

- Explanation of primary keys, foreign keys, and indexing for optimizing database performance.
- Example: Defining keys to link excavation sites to finds in different tables.

- **Basic Querying with SQL**

- Writing basic queries to retrieve data from a database.
- Example: Using `SELECT` statements to extract information about finds or excavation layers.

- **Filtering and Sorting Data**

- Using `WHERE`, `ORDER BY`, and `GROUP BY` clauses to filter and sort data.

- Example: Querying artifacts by material type or sorting excavation records by date.
- **Joining Tables**
 - Using JOIN statements to combine related tables (e.g., linking artifact data with stratigraphy).
 - Example: Writing queries to retrieve artifacts based on their excavation context.

7.4 Database Tools for Archaeology

- **Accessing Databases from R**
 - Introduction to R packages (DBI, RSQLite, RPostgres) for interacting with databases.
 - Example: Connecting to an SQLite or PostgreSQL database from R for archaeological data analysis.
- **SQLite for small projects**
 - Introduction to SQLite as a lightweight database solution for archaeological projects.
 - Example: Creating a simple SQLite database for site data using R (RSQLite, e.g. <https://cran.r-project.org/web/packages/RSQLite/vignettes/RSQLite.html>) or Python.
- **PostgreSQL for larger projects**
 - Overview of PostgreSQL for larger archaeological datasets.
 - Example: Setting up a PostgreSQL database in R for managing excavation records.
- **Querying Databases in R**
 - Writing SQL queries in R and retrieving data for further analysis.
 - Example: Querying an excavation database from R and visualizing the results with ggplot2.

- **Data Wrangling and Cleaning**

- Using R to clean and manipulate database data for analysis.
- Example: Using dplyr in R to filter and transform queried data.

7.5 Hands-on Practice

- **Building a Small Archaeological Database**

- Step-by-step walkthrough of creating a database schema for a hypothetical excavation project (using RSQLite).
- Example: Creating tables for stratigraphy, finds, and context.

- **Inserting and Managing Data**

- Practical examples of adding and managing archaeological data.
- Example: Inserting excavation records into the database using SQL.

- **Q&A and Troubleshooting**

- Addressing challenges in database design and SQL queries.

8 Databases (II)

8.1 Using GIS Databases for Archaeology

- **Introduction to Spatial Databases**
 - Overview of spatial databases and their use in archaeology (e.g., storing geospatial excavation data).
 - Example: Introduction to PostGIS for spatial data management in PostgreSQL.
- **Linking GIS Data to Databases**
 - Using databases to manage spatial data for archaeological analysis.
 - Example: Storing site coordinates and linking them to excavation records in a spatial database.
- **Querying Spatial Data**
 - Writing spatial queries to retrieve geospatial data from a database.
 - Example: Querying sites within a specific radius or extracting artifact distributions across layers.

8.2 Hands-on Practice

- **Querying and Analyzing Archaeological Data**
 - Walkthrough of writing SQL queries to extract meaningful insights from archaeological data.
 - Example: Retrieving and analyzing artifact distributions based on context and stratigraphy.

- **Combining Database and Spatial Data**
 - Practical examples of integrating database queries with GIS data for archaeological site analysis.
 - Example: Using a database to visualize the spatial distribution of artifacts across an excavation site.
- **Q&A and Troubleshooting**
 - Addressing challenges in querying and analyzing archaeological databases.

Part V

GIS

9 GIS (I)

(Basic concepts, installation and setting up, loading files)

Cool archive video about ARCInfo (1988): <https://youtu.be/7xqNyUOIRCsi?si=FulmlUVzaThGE9BU>

10 GIS (II)

(get a map image displaying data from shape and raster files)

References

- Baker, Monya. 2016. “1,500 Scientists Lift the Lid on Reproducibility.” *Nature* 533 (7604): 452–54. <https://doi.org/10.1038/533452a>.
- Community, The Turing Way. 2022. “The Turing Way: A Handbook for Reproducible, Ethical and Collaborative Research.” Zenodo. <https://doi.org/10.5281/ZENODO.3233853>.
- “Coursera | Degrees, Certificates, & Free Online Courses.” n.d. Accessed October 8, 2024. <https://www.coursera.org/>.
- “FAIR Principles.” n.d. *GO FAIR*. Accessed October 7, 2024. <https://www.go-fair.org/fair-principles/>.
- “GRN · German Reproducibility Network.” n.d. Accessed October 7, 2024. <https://reproducibilitynetwork.de/>.
- Iatropoulou, Katerina. n.d. “OpenAIRE.” *OpenAIRE*. Accessed October 7, 2024. <https://www.openaire.eu/>.
- “Introduction to Programming Languages.” 2018. *GeeksforGeeks*. <https://www.geeksforgeeks.org/introduction-to-programming-languages/>.
- “Learn R, Python & Data Science Online.” n.d. Accessed October 8, 2024. <https://www.datacamp.com>.
- “Markdown Guide.” n.d. Accessed October 7, 2024. <https://www.markdownguide.org/>.
- Marwick, Ben. 2017. “Open Science in Archaeology,” January. <https://doi.org/10.17605/OSF.IO/3D6XX>.
- National Academies of Sciences, Engineering, Policy and Global Affairs, Engineering Committee on Science, Board on Research Data Information, Division on Engineering and Physical and Sciences, Committee on Applied and Theoretical Statistics, Board on Mathematical Sciences Analytics, et al. 2019. “Understanding Reproducibility and Replicability.” In *Reproducibility and Replicability in Science*. National Academies Press (US). <https://www.ncbi.nlm.nih.gov/books/NBK547546/>.
- “Online Courses - Learn Anything, On Your Schedule | Udemy.” n.d. Accessed October 8, 2024. <https://www.udemy.com/>.
- R Core Team. 2024. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- “Spatial Without Compromise · QGIS Web Site.” n.d. Accessed October 7, 2024. <https://qgis.org/>.
- “What Is Programming? And How To Get Started.” 2024. *Coursera*. <https://www.coursera.org/articles/what-is-programming>.