

Kingdom of Saudi Arabia
Ministry of Education
King Saud University
College of Business
Administration
Department of Management
Information Systems

المملكة العربية السعودية
وزارة التعليم
جامعة الملك سعود
كلية إدارة الأعمال
قسم نظم المعلومات الإدارية



Exploring the Depths of Data Science: A Comprehensive Journey with Rust, Docker, and
Apache Arrow

Salman Almuzaini

Data Odyssey: A Personal Journey through Analysis and Code

King Saud University

November, 2023

Table of Contents

N.1 Day 1: Rust, APIs, and Containerization Challenges	5
N.1.1 Setting Up with Cargo.toml	5
N.1.2 The Build: Rust's Compilation Dance	5
N.1.3 The Runtime: A Tryst with Errors	5
N.1.4 Docker: The Container Conundrum	5
N.1.5 Project Scope and Reflections	6
N.2 Day 2: Dockerization and Execution Woes	6
N.2.1 Docker Build: A Smooth Sail	6
N.2.2 Execution Hurdles: The Binary Conundrum	6
N.2.3 Troubleshooting: A Deep Dive	6
N.2.4 Reflections and Learnings	7
N.3 Day 3: Rust Error Handling and Data Management	7
N.3.1 The Error Conundrum: Rust's Type System at Play	7
N.3.2 Data Export Trials: CSV Format Handling	7
N.3.3 User Interaction and Support	8
N.3.4 Reflections on Progress	8
N.4 Day 4: Transitioning Environments and Refining Data Handling	8
N.4.1 Workflow Shift: From Docker to Local Development	8
N.4.2 Project Restructuring: Modularization for Clarity	9
N.4.3 CSV Parsing: Overcoming API Limitations with Local Data	9
N.4.4 Data Science with Rust: Evaluating Against Python	9
N.4.5 Insights on Language Proficiency and Prospective Pathways	9
N.5 Day 5: Overcoming Compilation Errors and Strategically Pivoting to Structured Data	10
N.5.1 Compilation Errors: A Steep Learning Curve	10
N.5.2 Embracing Apache Arrow for Scalability	10
N.5.3 Overcoming Data Complexity with Python Scripting	10
N.5.4 The Shift to a Cleaned Dataset: 'organization-100.csv'	10
N.5.5 Reflections on Adaptability and Tool Selection	10
N.6 Day 6: Refinement and Error Resolution in CSV to Arrow Conversion	10
N.6.1 Type Handling and Schema Validation	11
N.6.2 Debugging Data Parsing Issues	11
N.6.3 Progress and Reflections	11
N.7 Day 7: Structuring for Documentation and Compliance	11
N.7.1 Emphasizing Task Management	12
N.7.2 Commitment to Legal Acknowledgement	12
N.7.3 Day's Reflection: Beyond Coding	12
N.8 Day 8: Refinement and Augmentation of Error Handling Mechanisms	12
N.8.1 Enhancing Error Management: Introduction of error_handling.rs	12
N.8.2 The Iterative Process: Debugging and Problem Resolution	12
N.8.3 Data Interaction: Developing read_csv.rs	13
N.8.4 Reflections on Progress and Mastery	13
N.9 Day 9: Integrating Citations into Documentation	13
N.9.1 Documentation Refinement: The Role of Citations	13
N.9.2 The refs.bib File: A Repository of Citations	13
N.9.3 Embedding Citations in Typst	13
N.9.4 Reflections on Academic Rigor and Integrity	14

N.10 Day 10: Navigating Between LaTeX and Typst Formats	14
N.10.1 Deliberating Documentation: LaTeX vs. Typst	14
N.10.2 Experimenting with LaTeX: A Brief Foray	14
N.10.3 The Reversion: A Return to Typst	14
N.10.4 Reflections: Simplicity vs. Sophistication	14
N.11 Day 11: Elevating Documentation with a New Template	14
N.11.1 Embracing a Structured Approach with a New Template	15
N.11.2 Strengthening Academic Integrity with Citations	15
N.11.3 Introducing a Summary and Table of Contents for Easy Navigation	15
N.11.4 Personalizing the Footer for a Unique Touch	15
N.11.5 Implementing a Logo-Background and Rectifying Citation Issues	15
N.11.6 Reflecting on Day 11: A Day of Refinement and Enhancement	15
N.12 Day 12: Overcoming Data Analysis Challenges and Enhancing Functionality	15
N.12.1 Resolving the UTF-8 Error: Insightful Debugging	15
N.12.2 Identifying the Root Cause: Misinterpreted Data Format	16
N.12.3 Implementing Basic Data Analysis Functions	16
N.12.4 Reflections: The Importance of Detail-Oriented Debugging	16
N.13 Day 13: Expanding Data Processing and Tackling New Challenges	16
N.13.1 Transitioning to a Directory of CSV Files	16
N.13.2 Dynamic Schema Challenges and Organizational Shifts	16
N.13.3 Encountering and Resolving Data Format Errors	16
N.13.4 Debugging Rust Code: The Borrow Checker	17
N.13.5 Reflections: Scaling Up Data Analysis	17
Bibliography	17

Summary

This document is a detailed record of my foray into data science with Rust, capturing the essence of the journey from initial setup using Cargo to tackling various challenges. It narrates the progression through stages that involved integration with APIs, meticulous error management, and navigating Docker for environmental stability, emphasizing the intricacies of cross-compiling and static linking.

Throughout this exploration, data played a central role, with CSV files such as “Table 2.4.4U. Price Indexes for Personal Consumption Expenditures by Type of Product” and “Table 2.6. Personal Income and Its Disposition, Monthly” serving as fundamental resources. These files, derived from the Bureau of Economic Analysis (BEA), were crucial in analyzing economic indicators within the Rust environment.

A notable point in the journey was the division of the “Personal Consumption Expenditures Price Index.xlsx” file into individual CSV sheets such as “Contents.csv”, “Table 1.csv”, through “Table 7.csv”, a process facilitated by a Python script sourced from GitHub. This script, `getsheets.py`, was instrumental in transforming the multifaceted Excel document into a structured array of CSV files, thereby simplifying the dataset for subsequent analysis.

The narrative also highlights the strategic shift to a more manageable dataset, ‘organization-100.csv’, for effective analysis, and the utilization of ‘iris.csv’ for exploring data processing capabilities within Rust. Additionally, the Rust Jupyter notebook, sourced from the Rust user’s forum and adapted from online tutorials, provided an interactive platform for executing and testing code snippets, further enriching the learning experience.

The document reflects on the importance of adaptability in learning and the precision required for robust error handling in Rust. It is a testament to the technical milestones achieved and an introspection of the strategic decision-making and problem-solving approaches vital for data science within the Rust ecosystem.

N.1 Day 1: Rust, APIs, and Containerization Challenges

Embarking on a journey to master Rust for data science, I set my sights on leveraging APIs for dynamic data retrieval. The KAPSARC (King Abdullah Petroleum Studies and Research Center) API [1] promised a wealth of energy data, but tapping into this resource proved to be a complex task.

N.1.1 Setting Up with Cargo.toml

The first order of business was configuring the Cargo.toml to declare the dependencies necessary for making HTTP requests (request) and handling JSON data (serde). After overcoming an initial hiccup where request wasn't recognized, a corrected Cargo.toml looked like this:

```
[dependencies]
request = "0.11"
serde = "1.0"
serde_json = "1.0"
```

This setup was my gateway to the world of Rust dependencies.

N.1.2 The Build: Rust's Compilation Dance

Executing cargo build set off a cascade of downloads and compilations, a testament to the rich ecosystem Rust offers. Watching libraries like cfg-if and percent-encoding get ready for action was my first taste of victory.

N.1.3 The Runtime: A Tryst with Errors

However, my celebration was short-lived as cargo run threw a stark runtime error — a missing file. It underscored the importance of crafting precise error handling in Rust. Adjusting the error handling logic, I braced myself for the next attempt.

But the hurdles continued. The API's JSON response was missing the total_count field, expected by my parsing logic. Rust's Option type came to the rescue, allowing me to handle the unpredictability of real-world data:

```
#[derive(Deserialize)]
struct ApiResponse {
    total_count: Option<i32>, // Now optional, no missing field error
    // Other fields as required
}
```

This adjustment was a key learning moment in data flexibility and error resilience.

N.1.4 Docker: The Container Conundrum

As the project progressed, I decided to encapsulate the development environment using Docker, with Alpine Linux as the base image for its minimalistic footprint. This decision led to a confrontation with a new beast: building Rust applications for musl libc instead of the default glibc. This choice, aimed at creating a more secure and portable executable, introduced me to the challenges of cross-compiling Rust applications.

The openssl-sys crate, necessary for HTTPS requests, became the center of my struggle as it clashed with the musl-based environment. The errors were obscure, leading to a deep dive into Alpine's package manager and Rust's compilation flags. The learning curve was steep, as I navigated through the intricacies of static linking and system dependencies.

N.1.5 Project Scope and Reflections

Amidst the technical battles, the scope of my project began to shift. I had to reassess my goals and adapt my approach to align with my growing understanding of Rust's capabilities and limitations. The initial intention to use Jupyter notebooks [2] for an interactive Rust experience was shelved in favor of focusing on the API integration aspect.

Day 1 was not just about code; it was about adaptability, problem-solving, and continuous learning. It reinforced the iterative nature of programming, where each error leads to a deeper insight, and every decision shapes the trajectory of the project.

In the rearview mirror, the day was a confluence of technical challenges, strategic pivots, and the inevitable realization that learning Rust was as much about understanding its ecosystem as it was about the language constructs themselves. It was a day of laying the groundwork for a journey that was only just beginning.

N.2 Day 2: Dockerization and Execution Woes

With the foundations set on Day 1, Day 2's mission was to containerize the Rust environment using Docker to ensure a consistent development and execution environment. This was done with the aim of creating a lightweight and secure setup that could potentially be deployed anywhere with ease.

N.2.1 Docker Build: A Smooth Sail

The day began with a sense of optimism as I initiated the Docker build process. The Dockerfile was designed to use Alpine Linux as the base image, keeping the overall footprint small. The build sequence was crafted carefully to ensure that all Rust tooling and dependencies were in place:

The curl command was used to test the API connectivity within the container. The jq tool was installed to allow for command-line JSON processing. A new Rust project RustDataVoyager was created within the container. The project source code was copied into the container. A final cargo build command was invoked to compile the Rust application. The Docker image build log showed no errors, and the tag dataquest-rust-docker:latest was successfully applied to the image. It was a moment of accomplishment as the image was now ready to be instantiated into a container.

N.2.2 Execution Hurdles: The Binary Conundrum

However, my excitement was tempered when I attempted to run the newly created Rust application inside the Docker container. The command cargo run failed, with the error message indicating that the target/debug/rust_data_voyager binary could not be executed due to it not being found. This was perplexing as the binary should have been created during the build process.

I was dropped into an Alpine Linux shell prompt inside the container to further investigate the issue. Commands like uname -a confirmed the Linux environment, while a listing of the RustDataVoyager directory showed all the expected files in place, yet the binary remained elusive.

N.2.3 Troubleshooting: A Deep Dive

The troubleshooting process was meticulous:

I reviewed the build output during the Docker image build to ensure there were no hidden errors. I verified that the Dockerfile and build process were supposed to compile the Rust application and generate the binary. I checked the target/debug directory within the source code to look for the binary. I made sure that the Dockerfile set the working directory correctly to /usr/src/RustDataVoyager. Each step was taken with a careful approach to uncover the root cause of the missing binary. The process involved a lot of back-and-forth, examining build logs, and verifying directory paths.

N.2.4 Reflections and Learnings

Day 2 brought with it a stark reminder of the complexities of containerization and cross-environment development. While Docker promised a uniform environment, it also introduced new variables that required attention to detail and a deep understanding of both the Rust ecosystem and Docker's mechanisms.

Despite the setbacks, the day was rich with learning opportunities. It underscored the importance of verifying every step of the development pipeline, from writing code to deploying executables. It was a day that tested patience and perseverance but ultimately contributed to a growing repository of knowledge and skills.

N.3 Day 3: Rust Error Handling and Data Management

On Day 3, I confronted a critical juncture in my Rust learning journey. The initial idea of leveraging a live API for data science tasks was met with unexpected challenges. The API proved to be unreliable, or perhaps it was my inexperience with its intricacies that led to inconsistent results. In response, I pivoted my approach:

- **Reliability:** I opted for downloaded CSV files, which offered a more stable and reliable data source than the unpredictable API.
- **Simplicity:** This shift allowed me to focus my efforts on data processing with Rust, without the added complexity of API management.
- **Control:** Static CSV files provided a controlled environment for a more straightforward debugging and learning process.

Acknowledging the steep learning curve with the API, I redirected my project to work with CSV files, setting the stage for a more focused and fruitful exploration of data science with Rust.

N.3.1 The Error Conundrum: Rust's Type System at Play

The day's work began with an error handling issue that Rust's compiler brought to light. My attempt to create a file and write to it within a function that returns a `Result<(), reqwest::Error>` resulted in compilation errors. The compiler was unsparing in its feedback:

```
let mut file = File::create(&file_path)?;  
file.write_all(&bytes)?;
```

The `?` operator, which I had used to propagate errors up the call stack, was designed to work with compatible error types. However, the errors arising from file operations (`std::io::Error`) were not automatically convertible to the expected `reqwest::Error`. This was Rust's type system enforcing strict error type consistency.

I learned that I had to map the error to the appropriate type or change the function's return type to accommodate multiple error types using an enum or a boxed error type. This was a deep dive into Rust's powerful, yet sometimes complex, error handling model.

N.3.2 Data Export Trials: CSV Format Handling

The day also saw me exploring data exportation. I had been interacting with the KAPSARC API [1], trying to export datasets. After some trial and error, I decided to attempt exporting a dataset in CSV format. Using the curl command, I fetched the data and saved it to "main.csv". The command looked like this:

```
curl -X 'GET' 'https://datasource.kapsarc.org/api/explore/v2.1/catalog/datasets/
the-renewable-energy-policy-paradox/exports/csv' -H 'accept: application/json;
charset=utf-8' -o main.csv
```

The CSV file, once opened, presented a new challenge — it contained metadata headers that needed parsing to access the actual dataset content:

```
recordid;_record_id;record_timestamp;_record_timestamp;...
```

I learned that dealing with real-world data often involves an additional layer of processing to separate useful information from metadata. This was an important lesson in data cleaning and preparation, essential skills for any data scientist.

N.3.3 User Interaction and Support

As part of the learning experience, I also engaged in discussions about how to handle the CSV data. The headers in the CSV seemed to pertain to metadata rather than the data itself. I advised on opening the file in a text editor or spreadsheet program to inspect the content and find the actual dataset data.

This interaction was a reminder of the importance of good communication skills in programming — being able to clearly explain and discuss issues is as crucial as solving them.

N.3.4 Reflections on Progress

The third day brought with it a recognition of the multifaceted nature of programming. It was not just about writing code; it was also about understanding the ecosystem in which the code operates and the data it manipulates. Error handling, data processing, and user support — each aspect was a piece of the larger puzzle.

The day concluded with a sense of accomplishment and a deeper appreciation for Rust's capability to enforce robustness through its strict type system. The journey through Rust was proving to be as much about learning the language as it was about understanding the principles of software development and data science.

N.4 Day 4: Transitioning Environments and Refining Data Handling

The fourth day of my data science journey with Rust was one of significant transitions and deep technical exploration. Embracing substantial changes to my development setup, I tackled the intricacies of data manipulation and began to reshape my project's structure for enhanced maintainability and scalability.

N.4.1 Workflow Shift: From Docker to Local Development

A pivotal shift occurred in my workflow as I moved away from Docker and Alpine Linux containers. I opted to run Rust locally on my Mac, which provided immediate benefits:

- **Local Environment:** The integrated development process and immediate feedback loop improved my efficiency and speed.
- **Simplicity:** I removed the overhead of Docker management, streamlining my focus onto Rust itself.
- **Performance:** By avoiding the potential limitations of containerized setups, my Mac's resources were better utilized.

This transition represented a growing understanding of my development tools and a commitment to optimizing the learning environment.

N.4.2 Project Restructuring: Modularization for Clarity

To address the increasing complexity of my project, I reorganized my codebase into several distinct modules:

- `concatenation.rs`: Handled the merging of disparate data sources.
- `arrow_converter.rs`: Encapsulated the logic for converting CSV data into the Apache Arrow format.
- `data_analysis.rs`: Contained functionalities for data analysis and insights extraction.
- `data_validation.rs`: Focused on ensuring data integrity and accuracy.

The introduction of a `lib.rs` file allowed these modules to be orchestrated cohesively, with `main.rs` now acting as the entry point that called upon the library for specific functionalities.

N.4.3 CSV Parsing: Overcoming API Limitations with Local Data

Faced with the challenges of an unresponsive API, I turned to downloading CSV files for local data analysis. The datasets included:

1. Table 2.4.4U. Price Indexes for Personal Consumption Expenditures by Type of Product.csv [3]
2. Table 2.4.5U. Personal Consumption Expenditures by Type of Product.csv [4]
3. Table 2.4.6U. Real Personal Consumption Expenditures by Type of Product, Chained Dollars.csv [5]
4. Table 2.6. Personal Income and Its Disposition, Monthly.csv [6]
5. Table 2.8.7. Percent Change From Preceding Period in Prices for Personal Consumption Expenditures by Major Type of Product, Monthly.csv [7]
6. Personal Consumption Expenditures Price Index.xlsx [8]

During the process, I encountered a schema mismatch error:

```
Error: InvalidArgumentError("number of columns(1) must match number of fields(47) in schema")
```

This led to the enhancement of the `create_arrow_schema` function to more accurately mirror the structure of the CSV files. Subsequently, I adopted the Apache Arrow format for its robust scalability and suitability for simulating complex data environments.

The decision to use Apache Arrow's columnar storage was pivotal in replicating real-world data workflows, reinforcing my commitment to utilizing appropriate tools for authentic data science endeavors.

N.4.4 Data Science with Rust: Evaluating Against Python

I invested time in evaluating Rust's data manipulation capabilities against Python's, tackling operations like CSV parsing, handling missing values, and adjusting data types. Rust provided granular control, contrasting with Python's user-friendly approach.

N.4.5 Insights on Language Proficiency and Prospective Pathways

Day 4 marked a significant milestone in my exploration of Rust, where I juxtaposed performance with usability. The insights gained not only deepened my understanding of Rust's strengths in data manipulation but also set the stage for future project hurdles.

N.5 Day 5: Overcoming Compilation Errors and Strategically Pivoting to Structured Data

Day 5 of my Rust journey was a blend of technical deep dives and strategic shifts. The day was split between navigating through compilation errors with Apache Arrow and making a pivotal transition to structured data for more efficient analysis.

N.5.1 Compilation Errors: A Steep Learning Curve

The initial part of the day was spent on addressing the complexities of Rust's type system and the Apache Arrow library, which were integral to my goal of validating and converting CSV data to the Apache Arrow format.

Type Mismatch Error: Faced with an `InvalidArgumentError`, I revised my code to dynamically align the Arrow schema with the CSV structure, as the number of columns in my CSV did not match the schema's fields.

Incorrect Method Calls: After encountering an error due to a non-existent method on the `FileReader` struct, I delved into the Arrow crate's documentation to correct my usage, which was a rich learning moment.

Incorrect Error Propagation: The misuse of the `?` operator led me to a better understanding of Rust's error handling patterns, refining my approach to iterators and non-`Result` types.

N.5.2 Embracing Apache Arrow for Scalability

Despite the initial setbacks, I remained committed to utilizing Apache Arrow, appreciating its scalability and the efficiency it offers for handling large datasets—key for real-world data science tasks.

N.5.3 Overcoming Data Complexity with Python Scripting

Acknowledging the complexity of the previous CSV files, I utilized a Python script [9] to break down a comprehensive Excel file into multiple CSV files, leading to a cleaner and more manageable dataset.

N.5.4 The Shift to a Cleaned Dataset: 'organization-100.csv'

I pivoted to the 'organization-100.csv' [10] dataset, which was not only cleaner but also free from the encoding issues of the 'iris.csv' [11] file, thereby enhancing my ability to focus on data processing and analysis within Rust.

N.5.5 Reflections on Adaptability and Tool Selection

The day was marked by the importance of adaptability in data science. The strategic shift to a more suitable dataset, combined with the use of cross-language scripting for data preparation, underscored the iterative nature of programming and the adaptive approach necessary for managing data effectively.

Day 5 was both challenging and rewarding, a testament to the iterative nature of software development. Each error message provided a learning opportunity, and the strategic pivot to structured data emphasized the project's adaptability and growth.

N.6 Day 6: Refinement and Error Resolution in CSV to Arrow Conversion

Day 6 of my Rust for data science journey was dedicated to refining the CSV to Apache Arrow conversion process, with a focus on resolving compilation errors and improving data validation.

By Day 6, I was fully immersed in refining the CSV to Apache Arrow conversion process. With the project scope change on Day 3 already established, I could dedicate my efforts to overcoming the technical complexities of Rust and the Arrow library.

- **Error Handling:** I tackled type mismatches and data parsing errors head-on, each one offering new insights into Rust's strict type safety and error handling.
- **Data Conversion:** The conversion process from CSV to Apache Arrow format became smoother as I applied the lessons learned from previous days, appreciating the efficiency and scalability of the Arrow format for big data tasks.
- **Technical Insight:** Day 6 was marked by a series of breakthroughs in understanding and utilizing Rust's powerful features for data processing.

Reflecting on the decision made on Day 3, the benefits of working with CSV files were clear. It provided a more predictable and manageable environment for mastering Rust's data science tools, setting a strong foundation for future developments.

N.6.1 Type Handling and Schema Validation

I encountered a type mismatch error that prevented the successful compilation of my code. The issue was with the handling of an Arc object, a reference-counted pointer used in Rust to allow multiple ownership of data, which is especially useful for complex data structures like Apache Arrow schemas. The error message was clear:

```
Error: InvalidArgumentError("number of columns(1) must match number of fields(47) in schema")
```

This pointed to a mismatch between the expected schema and the CSV file's actual structure. To resolve this, I updated the `validate_csv_input` and `convert_csv_to_arrow` functions to ensure they handled the Arc correctly, providing a consistent and accurate schema for the data conversion process.

N.6.2 Debugging Data Parsing Issues

Another challenge was a `ParseFloatError`, indicating that the CSV data contained fields that could not be parsed into floating-point numbers as expected by the schema. This error highlighted the importance of robust data validation and error handling in data processing pipelines:

```
Error: ParseFloatError { kind: Invalid }
```

To address this, I had to carefully inspect the CSV file to ensure that all float fields were in the correct format and did not contain any invalid characters or formats. Adding more granular error handling in the `validate_record` function became necessary to gracefully handle such parsing errors and to provide detailed debugging information, such as the line number and field content that caused the issue.

N.6.3 Progress and Reflections

Despite these challenges, by the end of Day 6, I had made significant strides in my understanding of Rust's error handling and data validation mechanisms. I learned the importance of precise type handling when dealing with complex data structures and the necessity of thorough data validation to ensure the integrity of the conversion process.

Moreover, the day's work underscored the value of Apache Arrow in handling large datasets efficiently, validating my choice to use it for scalability and performance in data science tasks.

N.7 Day 7: Structuring for Documentation and Compliance

On the seventh day of my data science journey with Rust, I began with the administrative task of setting up a proper documentation and task management system. This was crucial as my project grew in complexity and required a more structured approach to keep track of the various components and dependencies involved.

N.7.1 Emphasizing Task Management

I started the day by creating a `TODO.md` file, which would act as a centralized checklist for my ongoing and upcoming tasks. This file is vital for maintaining a clear overview of my objectives and ensuring that nothing important slips through the cracks as the complexity of the project increases.

The primary item on my to-do list was to document my learning journey, a task that would not only solidify my understanding but also serve as a guide for future reference. This documentation was to be carried out using Typst, chosen for its combination of simplicity and power, ideal for a project where clarity is paramount.

TODO

- [] Document my learning journey using Typst.

N.7.2 Commitment to Legal Acknowledgement

Recognizing the significance of the open-source crates that empowered my project, I decided to add another task aimed at creating a `NOTICE.md` file. This file would acknowledge the contributions of the various crates and their authors, as well as ensuring that I adhere to the licensing requirements associated with their use.

Copy code

- [] Add a `NOTICE.md` file for crediting all crates used in the project, with their respective licenses.

N.7.3 Day's Reflection: Beyond Coding

Day 7 might not have been characterized by coding breakthroughs, but it was a day that highlighted the importance of the often overlooked aspects of software development: documentation and compliance. By addressing these areas, I was laying down a solid foundation for managing the project's growth and complexity.

The lessons from this day were clear—successful software development is multifaceted, with coding being just one aspect among many that include project management, documentation, and legal considerations. The groundwork laid on this day was crucial for a smooth journey ahead.

N.8 Day 8: Refinement and Augmentation of Error Handling Mechanisms

The eighth day of immersing myself in the Rust landscape was marked by an expansion of the project's error handling capabilities and the integration of new functionalities to enhance data interaction.

N.8.1 Enhancing Error Management: Introduction of `error_handling.rs`

With the increasing complexity of the project, it became imperative to consolidate and refine the error handling mechanisms. To this end, I introduced a new module, `error_handling.rs`, dedicated exclusively to managing the myriad of potential errors in a more structured and coherent manner. This module's purpose was to encapsulate the various error types and their handling logic, thereby promoting code cleanliness and maintainability.

N.8.2 The Iterative Process: Debugging and Problem Resolution

The bulk of the day revolved around the iterative process of debugging. I encountered an array of bugs and errors that demanded attention—each one offering a puzzle to be solved. This process was not just about fixing issues but also about understanding the deeper workings of the Rust language and the third-party libraries it interacts with. Through this rigorous exercise, my proficiency with the language's nuances continued to strengthen.

N.8.3 Data Interaction: Developing read_csv.rs

In parallel with refining error handling, I developed another module, read_csv.rs. This new addition to the project was tasked with reading and displaying CSV data directly within the console. By outputting the contents of the CSV to the screen, I could immediately verify the correctness of data ingestion and parsing processes—a crucial step for ensuring the integrity of subsequent data manipulation and analysis.

N.8.4 Reflections on Progress and Mastery

Day 8 was a testament to the evolving nature of software development. The day's efforts were directed not only toward problem-solving but also toward improving the project's robustness through better error handling and data interaction. It underscored the importance of resilience and patience, as each challenge surmounted paved the way for a deeper understanding and mastery of the Rust programming language.

As I concluded the day's tasks, it was evident that each line of code refined, and each bug squashed contributed to a more stable and reliable data science toolset within the Rust ecosystem. The journey continued to be as much about the destination as it was about the learnings gleaned along the way.

N.9 Day 9: Integrating Citations into Documentation

The ninth day of my Rust journey was dedicated to enhancing the documentation of my project by integrating citations. This was a crucial step in acknowledging the sources of the datasets and tools that have been instrumental in my data science exploration.

N.9.1 Documentation Refinement: The Role of Citations

Understanding the importance of proper citation for the credibility and reproducibility of scientific work, I set out to include references to all the datasets and tools I had utilized. This meticulous process involved creating a structured references file in the BibTeX format and embedding citations within the project documentation.

N.9.2 The refs.bib File: A Repository of Citations

I created a refs.bib file to store all the bibliographic references in a structured format suitable for LaTeX documents. This BibTeX file contained detailed entries for each citation, ensuring that all sources were accurately and consistently documented. Here's a glimpse of what the entries looked like:

```
@online{bea244u,  
  author = {{Bureau of Economic Analysis}},  
  title = {Table 2.4.4U. Price Indexes for Personal Consumption Expenditures by  
Type of Product},  
  year = {2020/2022},  
  url = {https://apps.bea.gov/iTable/?reqid=19&step=3&isuri=1&select_all_years=0&  
nipa_table_list=2013&series=m&first_year=2020&last_year=2022&scale=-99&categories=  
underlying&thetable=},  
  urldate = {2023-11-07}  
}
```

N.9.3 Embedding Citations in Typst

With the refs.bib file in place, I proceeded to embed citations within my Typst documentation. Using the citation keys, I could reference the datasets and tools inline, providing readers with the context and sources of the data I analyzed. For example:

... the data obtained from the Bureau of Economic Analysis [[@bea244u](#)] was pivotal in ...

N.9.4 Reflections on Academic Rigor and Integrity

Day 9 was not characterized by coding breakthroughs but by the reinforcement of academic rigor and integrity in my project. By integrating citations, I was not only paying homage to the creators of the tools and compilers of the data but also ensuring that my documentation met the high standards of scientific communication.

It was a day that emphasized that data science is not just about analysis and algorithms; it's also about the responsible presentation and dissemination of information. The addition of citations was a small but significant step toward upholding these values in my learning journey.

N.10 Day 10: Navigating Between LaTeX and Typst Formats

The tenth day of my Rust journey found me grappling with documentation formats, as I oscillated between LaTeX and Typst. This exploration was driven by a quest to find the most effective and efficient way to document my data science project.

N.10.1 Deliberating Documentation: LaTeX vs. Typst

Initially, I was inclined to transition my documentation from Typst to LaTeX, attracted by LaTeX's maturity, extensive array of templates, and comprehensive package ecosystem. This move seemed promising, given LaTeX's longstanding reputation in academic and technical writing.

N.10.2 Experimenting with LaTeX: A Brief Foray

Embracing the LaTeX format, I began converting my existing Typst documentation. The process involved adapting to LaTeX's syntax, exploring its diverse templates, and leveraging its robust package support. However, this journey was not without its challenges. The complexity of LaTeX's syntax, compared to the more straightforward nature of Typst, soon became apparent.

N.10.3 The Reversion: A Return to Typst

After delving into the depths of LaTeX, I made a strategic decision to revert to Typst. This decision was influenced by several factors, chief among them the simplicity and ease of use offered by Typst's syntax. Despite the allure of LaTeX's advanced features, the user-friendly nature of Typst ultimately held more appeal for my project's needs.

N.10.4 Reflections: Simplicity vs. Sophistication

Day 10 was a day of reflection and decision-making. It underscored the importance of choosing the right tools for the job, balancing the sophistication and capabilities of LaTeX against the simplicity and intuitiveness of Typst. This experience highlighted that while advanced features are valuable, the accessibility and ease of use are equally crucial for efficient documentation.

In conclusion, Day 10 was a journey of exploration between two worlds - the established realm of LaTeX and the emerging domain of Typst. It was a reminder that in the realm of data science, the choice of documentation tools is as vital as the analytical tools themselves, shaping the way we present and communicate our findings.

N.11 Day 11: Elevating Documentation with a New Template

On Day 11 of my Rust journey, I delved deeper into the realm of documentation, embracing a new report template to enhance the presentation of my data science project.

N.11.1 Embracing a Structured Approach with a New Template

I discovered an innovative report template from LostPy's Typst Templates [12], a tool that promised to streamline my documentation process. This template offered a structured, professional layout, enhancing the overall appearance and organization of my project.

N.11.2 Strengthening Academic Integrity with Citations

To uphold the academic integrity of my work, I integrated a citation feature into my documentation. This addition was crucial for acknowledging sources and reinforcing the credibility of my research, a key aspect of any scientific or technical endeavor.

N.11.3 Introducing a Summary and Table of Contents for Easy Navigation

To provide a concise overview and facilitate navigation, I added a summary and a table of contents to the report. The summary offered a quick glimpse into the key aspects of the project, while the table of contents allowed for easy and organized access to different sections.

N.11.4 Personalizing the Footer for a Unique Touch

Aiming for a more personalized touch, I customized the footer font, enabling me to experiment with different styles and choose one that best complemented the report's aesthetics and readability.

N.11.5 Implementing a Logo-Background and Rectifying Citation Issues

I took steps to further personalize the document by implementing a logo-background and addressing citation-related errors, ensuring both visual appeal and accuracy in referencing.

N.11.6 Reflecting on Day 11: A Day of Refinement and Enhancement

Day 11 was a pivotal moment in my Rust journey, dedicated to refining and elevating my documentation tools. The introduction of a new template, along with features like citation support, a summary, a table of contents, and footer customization, marked a significant enhancement in the professionalism and comprehensiveness of my project's presentation.

In summary, Day 11 was characterized by building upon the foundational work of previous days, focusing on details that elevate the effectiveness of documentation. This day underscored the ongoing evolution of my journey in Rust and data science, highlighting the importance of continuous improvement and adaptability.

N.12 Day 12: Overcoming Data Analysis Challenges and Enhancing Functionality

On the twelfth day of my Rust journey, I encountered and overcame a significant challenge in my data analysis process, leading to a deeper understanding of the tools and data formats involved.

N.12.1 Resolving the UTF-8 Error: Insightful Debugging

While working with my data set, I faced an error:

```
Csv(Error(Utf8 { pos: Some(Position { byte: 0, line: 1, record: 0 }), err:
Utf8Error { field: 0, valid_up_to: 8 } })))
```

To diagnose and fix this issue, I took a methodical approach, adding more print statements in `data_cleaning.rs` and `data_analysis.rs`. These additional logs helped me trace the data flow more effectively.

Initially, after reading the entire CSV file in `read_csv.rs`, I printed it out again post-cleaning in `data_cleaning.rs`. In `data_analysis.rs`, I inserted detailed print statements to track the analysis process, such as:

```
println!("Starting data analysis for file: {}", input_file);
println!("File opened successfully.");
println!("Processing record batch.");
```

N.12.2 Identifying the Root Cause: Misinterpreted Data Format

Through this detailed logging, I discovered that I was mistakenly processing an Arrow file as if it were a CSV file. As you might know, an Arrow file is binary, not a text file. This revelation was crucial, as it allowed me to correctly interpret and process the data.

N.12.3 Implementing Basic Data Analysis Functions

After rectifying the file format issue, I proceeded to add basic data analysis functions. These functions provided valuable insights into the dataset, which I outputted with print statements:

```
println!("Average number of employees: {}", average_employees);
println!("Organizations per country:");
For each country and its count in the dataset, println!("{}", country, count);
println!("Average founding year: {}", average_founded_year);
```

N.12.4 Reflections: The Importance of Detail-Oriented Debugging

Day 12 was a day marked by a crucial learning experience in data analysis and debugging. It highlighted the importance of being detail-oriented and methodical in troubleshooting. This experience not only helped me overcome the immediate error but also reinforced my understanding of data formats and the importance of correct data processing in Rust.

In conclusion, Day 12 was an enlightening chapter in my Rust journey, teaching me valuable lessons about the intricacies of data analysis and the significance of careful debugging and data format awareness.

N.13 Day 13: Expanding Data Processing and Tackling New Challenges

On the thirteenth day of my Rust journey, I embarked on a significant expansion of my data processing capabilities, transitioning from handling a single CSV file to managing an entire directory of diverse data sets.

N.13.1 Transitioning to a Directory of CSV Files

The day's task involved changing the focus from a single file, `organization-100.csv`, to a directory named `Datasets/csv-files`. This directory housed multiple CSV files, each of which required cleaning, validation, and eventual conversion into Arrow format, stored in `Datasets/arrow-files`.

N.13.2 Dynamic Schema Challenges and Organizational Shifts

A key challenge I faced was creating a dynamic Arrow schema capable of accommodating a variety of CSV files with different column structures. I also realized that my conversion logic from CSV to Arrow was mistakenly placed in `data_validation.rs`, while it should have been in `arrow_converter.rs`. Rectifying this helped streamline the process.

N.13.3 Encountering and Resolving Data Format Errors

During this process, I encountered a specific error:

```
Error: Error(UnequalLengths { pos: Some(Position { byte: 246, line: 3, record:
3 }), expected_len: 1, len: 47 })
```


N.13.4 Debugging Rust Code: The Borrow Checker

```
error[E0716]: temporary value dropped while borrowed
--> src/data_cleaning.rs:57:77
|
57 |         let mut csv_reader =
ReaderBuilder::new().has_headers(true).from_reader(data_lines.join("\n").as_bytes());
|
~~~~~ - temporary value is freed at the end of this
statement
|
|
| creates a temporary value which is freed while still in use
```

N.13.5 Reflections: Scaling Up Data Analysis

This day taught me the intricacies of managing larger datasets and the importance of flexible schema design in Rust. It was a step up in complexity, highlighting the differences between handling individual files and dealing with a collection of diverse datasets.

In conclusion, Day 13 marked a significant leap in my data processing capabilities. It was a day filled with learning and problem-solving, demonstrating the evolving nature of my Rust journey and the continuous adaptation required to handle increasingly complex data science tasks.

1957

Bibliography

- [1] K. A. P. Studies and R. Center, “KAPSARC's API Portal”. [Online]. Available: <https://datasource.kapsarc.org/api/explore/v2.1/console>
- [2] ShahinRostami, “Rust Jupyter Notebook”. Accessed: Nov. 07, 2023. [Online]. Available: <https://users.rust-lang.org/t/getting-started-with-rust-jupyter-notebooks/39004>
- [3] Bureau of Economic Analysis, “Table 2.4.4U. Price Indexes for Personal Consumption Expenditures by Type of Product”. Accessed: Nov. 07, 2023. [Online]. Available: https://apps.bea.gov/iTable/?reqid=19&step=3&isuri=1&select_all_years=0&nipa_table_list=2013&series=m&first_year=2020&last_year=2022&scale=-99&categories=underlying&thetable=
- [4] Bureau of Economic Analysis, “Table 2.4.5U. Personal Consumption Expenditures by Type of Product”. Accessed: Nov. 07, 2023. [Online]. Available: <https://apps.bea.gov/>

- iTable/?reqid=19&step=3&isuri=1&select_all_years=0&nipa_table_list=2014&series=m&first_year=2020&last_year=2022&scale=-99&categories=underlying&thetable=
- [5] Bureau of Economic Analysis, “Table 2.4.6U. Real Personal Consumption Expenditures by Type of Product, Chained Dollars”. Accessed: Nov. 07, 2023. [Online]. Available: https://apps.bea.gov/iTable/?reqid=19&step=3&isuri=1&select_all_years=0&nipa_table_list=2014&series=m&first_year=2020&last_year=2022&scale=-99&categories=underlying&thetable=
 - [6] Bureau of Economic Analysis, “Table 2.6. Personal Income and Its Disposition, Monthly”. Accessed: Nov. 07, 2023. [Online]. Available: <https://apps.bea.gov/iTable/?reqid=19&step=3&isuri=1&1921=survey&1903=76>
 - [7] Bureau of Economic Analysis, “Table 2.8.7. Percent Change From Preceding Period in Prices for Personal Consumption Expenditures by Major Type of Product, Monthly”. Accessed: Nov. 07, 2023. [Online]. Available: <https://apps.bea.gov/iTable/?reqid=19&step=3&isuri=1&1921=survey&1903=84>
 - [8] Bureau of Economic Analysis, “Personal Consumption Expenditures Price Index”. [Online]. Available: <https://www.bea.gov/data/income-saving/personal-income>
 - [9] S. Ming, “getsheets.py”. Accessed: Nov. 07, 2023. [Online]. Available: <https://gist.github.com/scottming/99c09685360376d4cac2de7c891e8050/9670acd5044d11871188e04b8123263ad132f512>
 - [10] Datablist, “organization-100.csv”. Accessed: Nov. 07, 2023. [Online]. Available: <https://www.datablist.com/learn/csv/download-sample-csv-files>
 - [11] J. Shin, “iris.csv”. Accessed: Nov. 07, 2023. [Online]. Available: <https://gist.github.com/netj/8836201>
 - [12] LostPy, “Report Template”. [Online]. Available: https://gitlab.com/LostPy/lostpy-typst-templates/-/blob/master/templates/report-template/report-template.typ?ref_type=heads