# CYT-250-Threat Investigation

# Final Project

# Seeing VERIS in Action-Working with VCDB Data

Elaborate by:

Leandro Delgado

Student Number: 114416241

Professor: Tatiana Outkina

**Introduction**

In this lab, I continued working with the VCDB dataset from Verizon, following the exercises outlined in Chapter 8 of the book *Data-Driven Security*. This chapter builds upon the work from Lab 7, where the Python script from Listing 7-2 was already executed. From this point forward, all scripts are provided in R, marking a shift in tools and approach.

As someone new to using the R environment, this lab also served as an opportunity to explore and become comfortable with this powerful yet approachable tool. My previous experience has mostly been with Google Colab, but due to recurring issues with installing libraries and running code there, I decided to switch to RStudio. Although R is not overly complex, I've come to realize that it requires practice and familiarity both of which I'm starting to develop through this lab.

The objective here is to study the chapter thoroughly, execute all the R scripts, and document the process with screenshots, meaningful comments, and personal insights. At each step, I reflect on what I've learned, the challenges I encountered, and how I overcame them, making this not only a technical exercise but also a personal journey of growth in data analysis and cybersecurity.

# Objective: Seeing VERIS in Action. Working with VCDB Data

Sources of information:

**Data-Driven security, Chapter 8. This is in continuation Lab 7 work on VCDB from Verizon. You already run the Python script in Listing 7-2. All other scripts are given in R. Your task is to study the Chapter and run all scripts. Make screenshots and meaningful comments. Explain what you have learned at each step.**
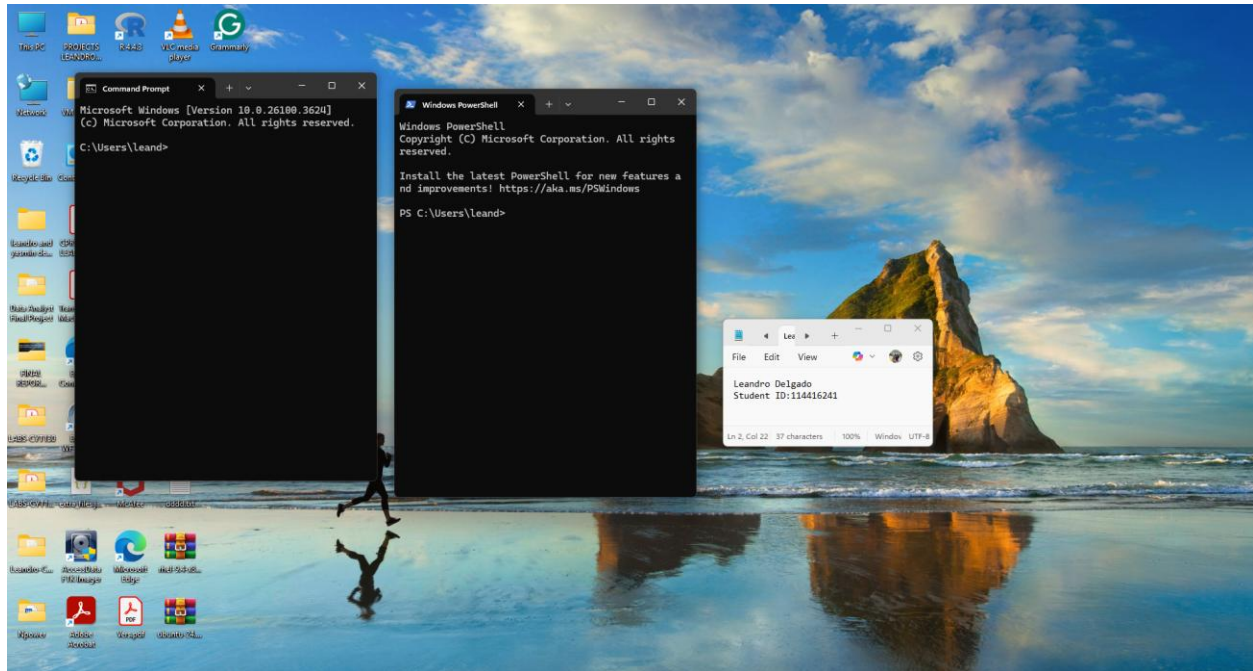


**Figure 1. 0 Screen**

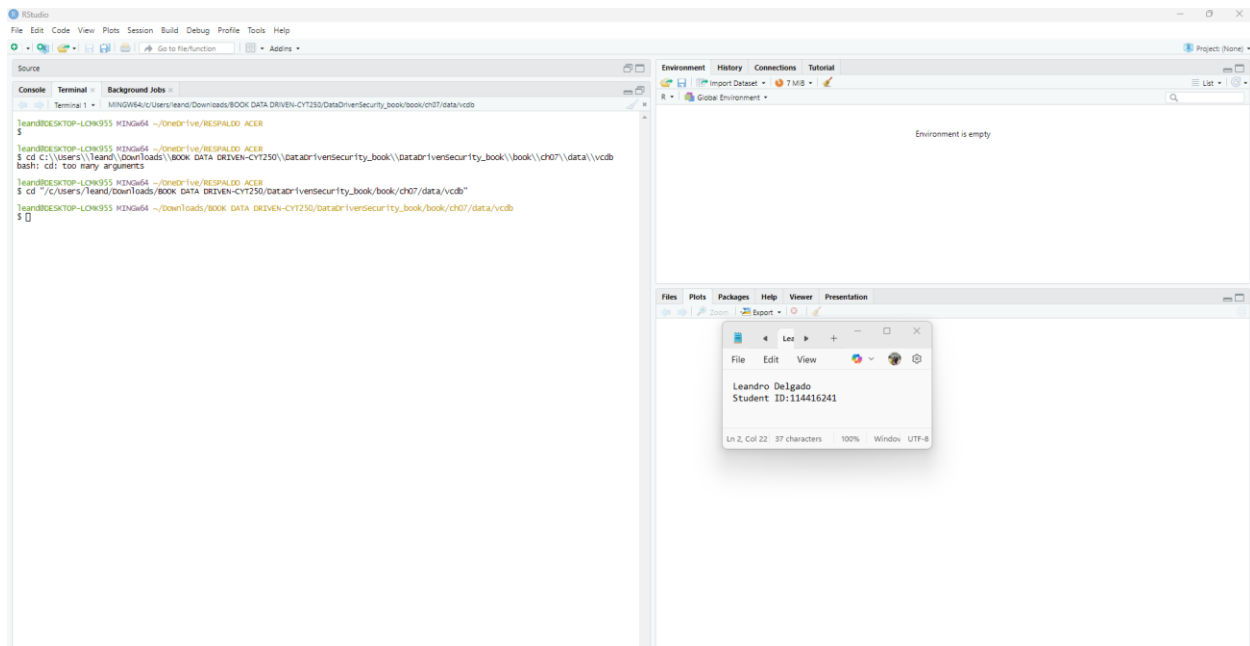**Task 1: Set R environment**

This was my first time using the R environment. My previous experience had been mainly with Google Colab. However, due to several technical issues with installing the right libraries and running the code properly there, I decided to try RStudio instead. It's a new tool for me — not too complex, but I know I need more practice to feel confident using it. I've started to explore how it works and I'm gradually getting more comfortable with it.

**Task 2: Study Chapter 7 material and run Listings 7-3 to 7-11**

**Listing 7-3**

To begin, I configured the working directory so R could easily access the chapter's data files. I then ran a script that checks for essential packages needed for the analysis. If any of them—like ggplot2 or rjson—weren't already installed, the script handled it automatically. This step helped me get the R environment ready for the work ahead without manually managing each dependency.

**Figure 3. Run listing 7-3**



**Figure 4. Result of run Listing 7-3**

**Listing 7-4**

In this step, I loaded a specific JSON file from the VCDB dataset and parsed it using the from JSON() function. Then, I printed the "hacking variety" field, which returned "SQLi", showing that the incident involved a SQL injection. It was interesting to see how structured data from a file can be easily accessed and explored in R.
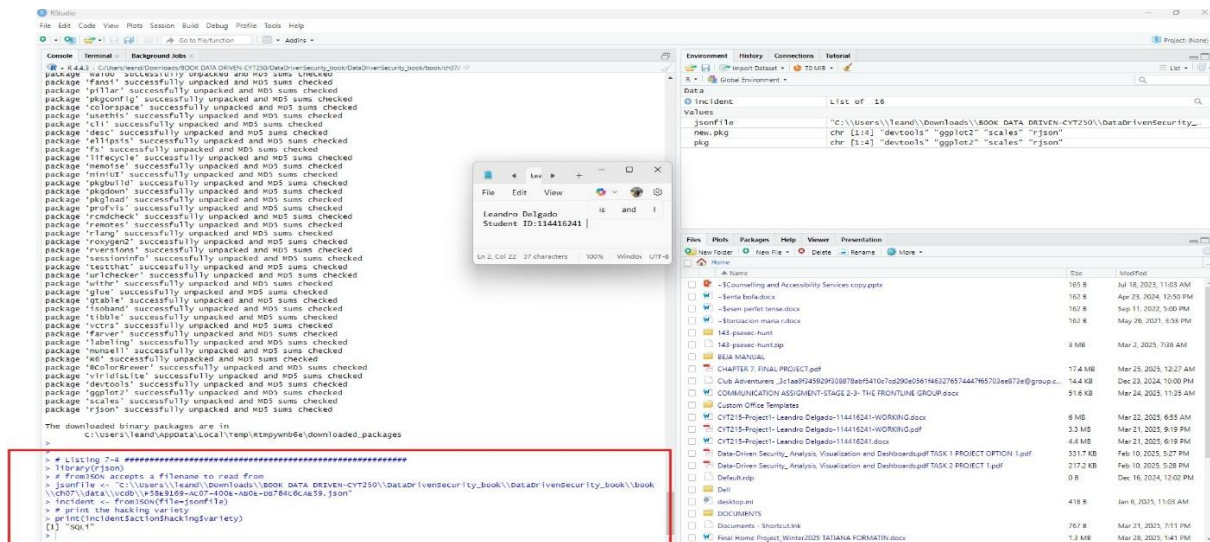


**Figure 5. Result listing 7-4**

**Listing 7-5**

In this step, I used the devtools package to install veris directly from GitHub. Since this package isn't available on CRAN, installing it this way is necessary. The process required downloading several dependencies, and I was also reminded to install Rtools, which is needed to build packages from source. This step helped me understand how to bring in external R packages that aren't part of the standard repositories.
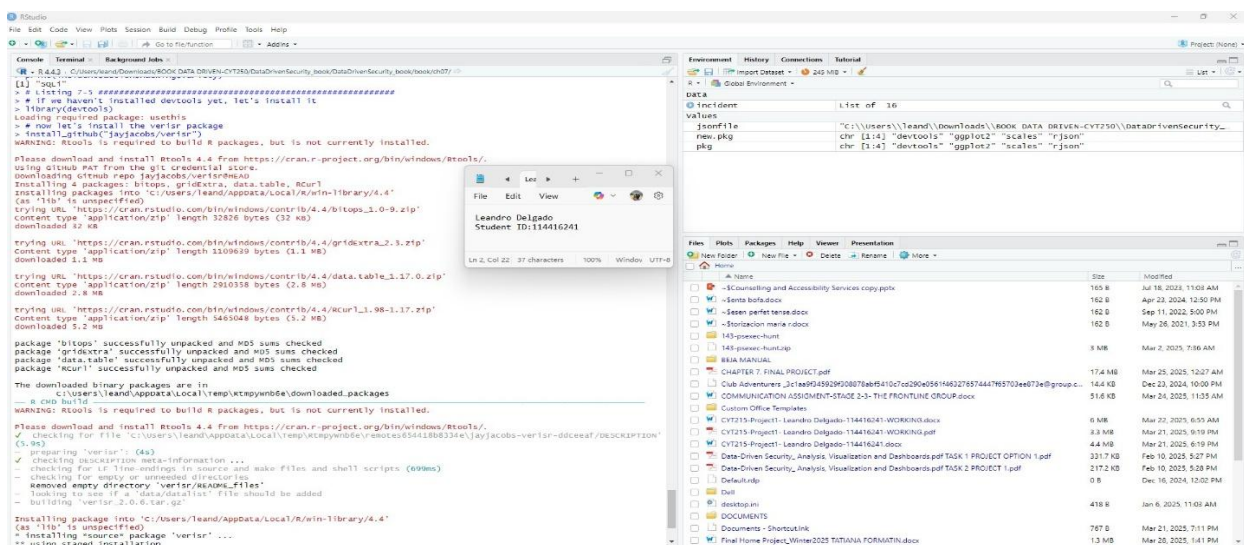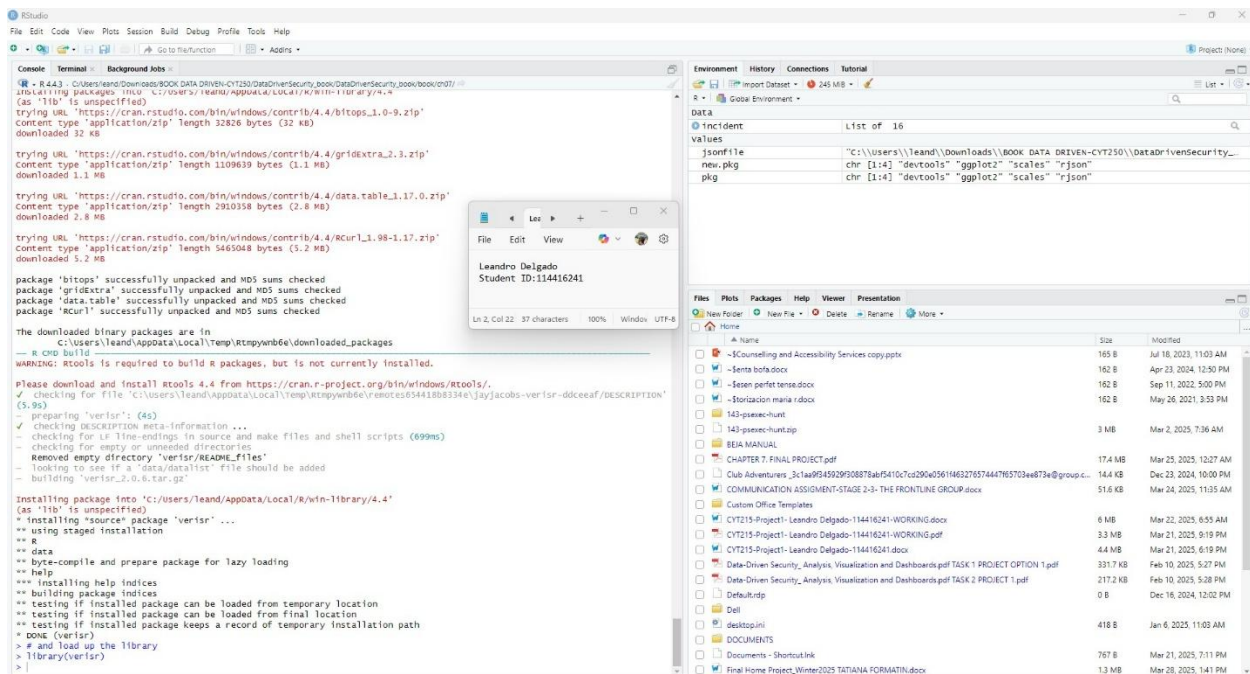


**Figure 6. Result listing 7-5**

**Figure 7. Result Listing 7-5**

## Listing 7-6

In this step, I used the veris package to read in a collection of VCDB incident files from a directory. The data was loaded into a vcdb object, which contains 1,643 records and over 2,400 variables. This gives me a structured dataset ready for analysis, with detailed information on cyber incidents.
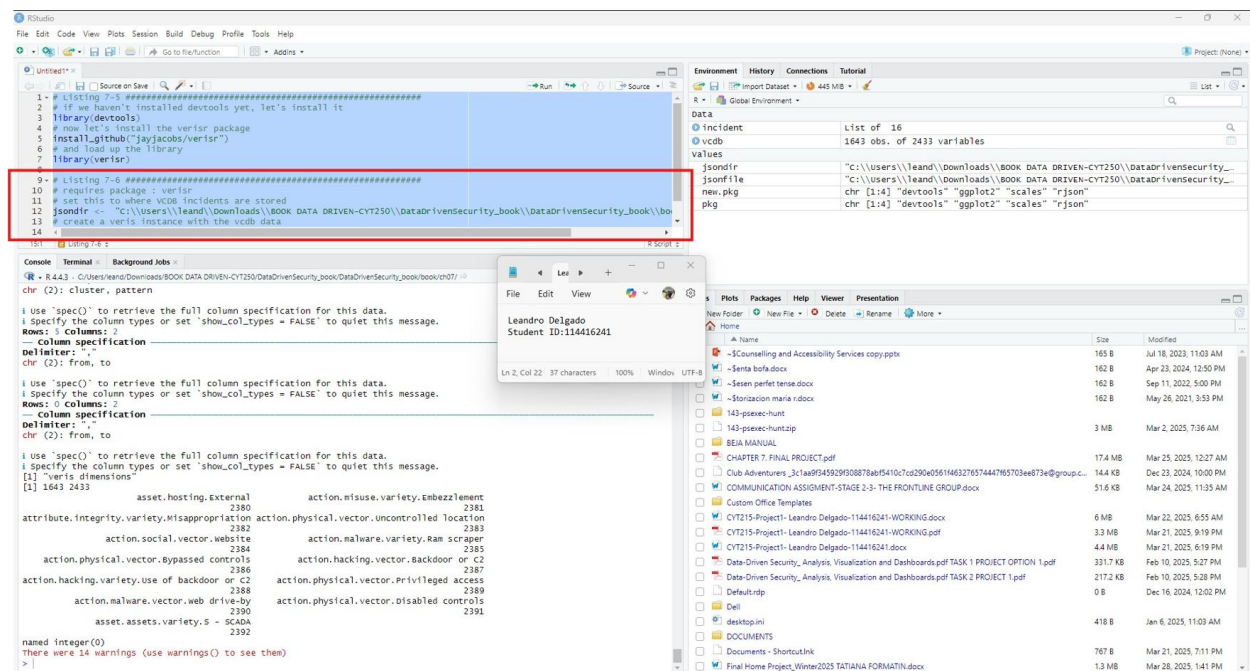


**Figure 8. Result listing 7-6**

## Listing 7-7

In this step, I used the summary() function on the vcdb object to get an overview of the dataset. It returned useful insights, such as the number of incidents involving external actors (955), and the most common impact being on confidentiality (1,604). Although the message confirms that the dataset contains 1,643 incidents, the actual records are not shown—only a high-level summary is provided. To explore individual incidents, I would need to use functions like head(vcdb) or View(vcdb).
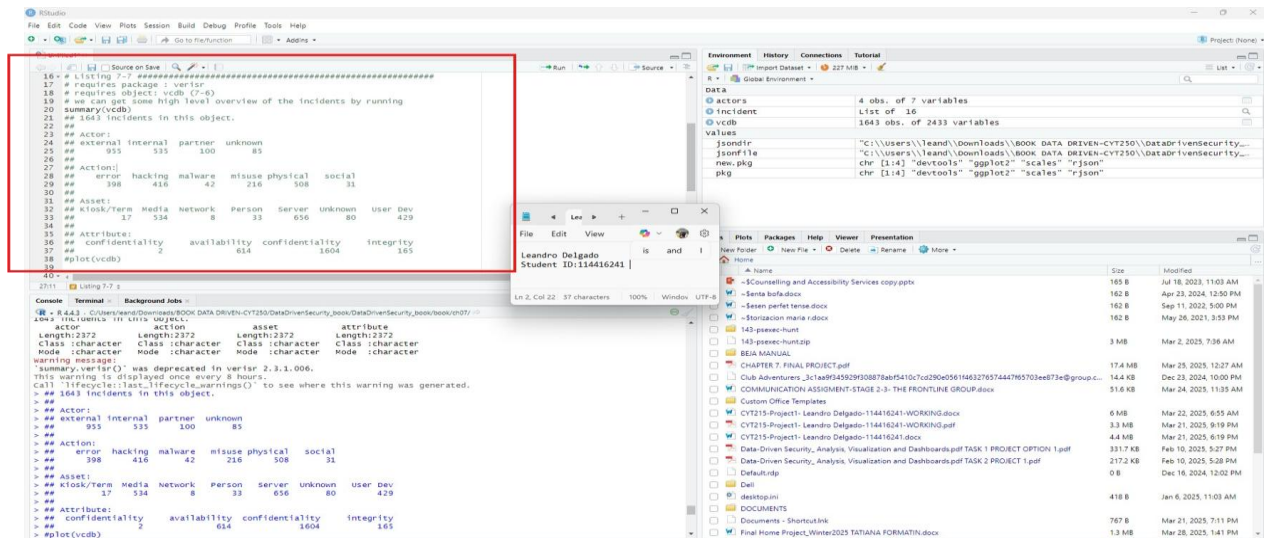


**Figure 9.Result Listing 7-7**

## Listing 7-8

This step uses the getenum() function to extract and count the different types of actors involved in the VCDB incidents. The results show that most incidents were caused by external actors (955), followed by internal (535), partners (100), and unknown sources (85). The output is stored as a data frame named actors, making it easy to view and use in further analysis.
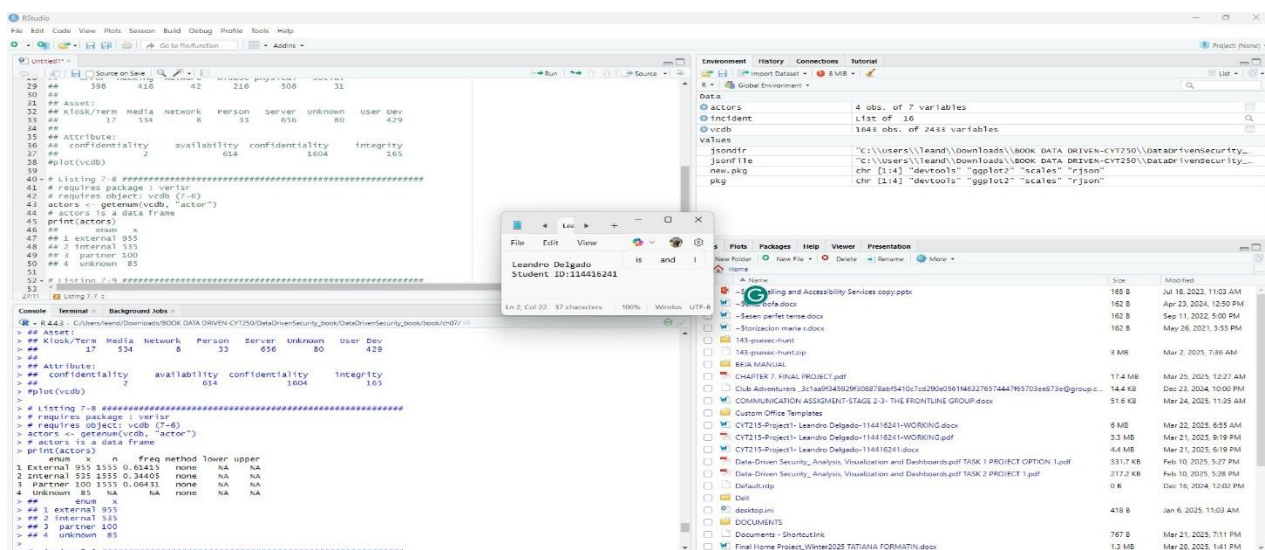


**Figure 10. Result 7-8**

## Listing 7-9

In this step, I ran an updated version of the previous command, this time adding arguments to include both the total count of incidents (n) and the frequency (freq) for each actor category. The results show that external actors were responsible for over 58% of the incidents, followed by internal actors at 32%. This gave me a clearer picture of the relative impact of each actor type within the dataset.
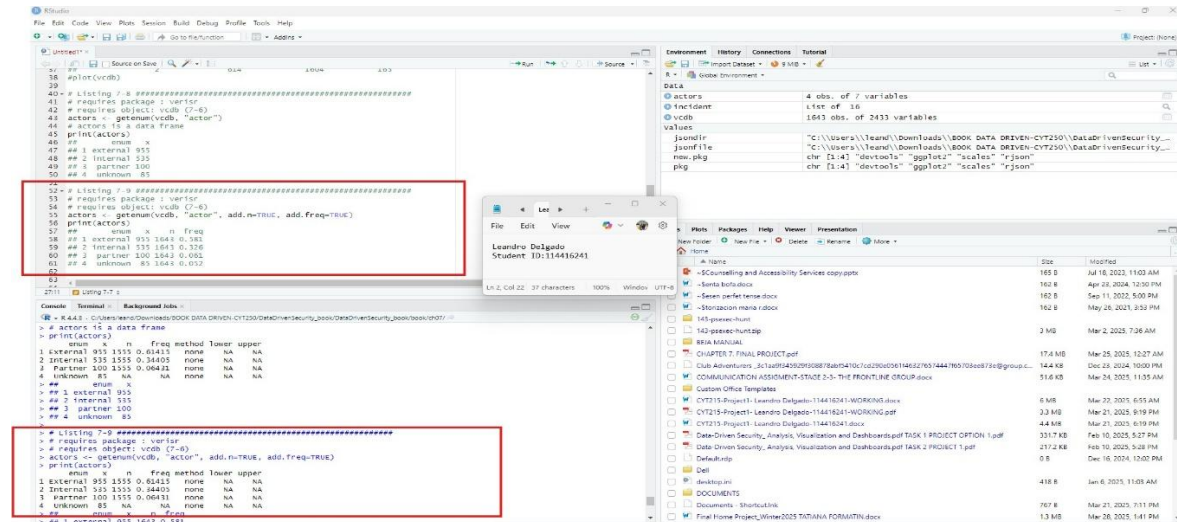


**Figure 11. Result listing 7-9**

## Listing 7-10

In this step, I created a custom function called verisplot() to generate bar charts from specific VCDB fields using ggplot2. I used it to visualize the different types of data affected under the attribute confidentiality.data.variety field. The chart showed that medical and personal data were the most compromised. This helped me understand how to combine data aggregation and visualization in R for clearer insights.
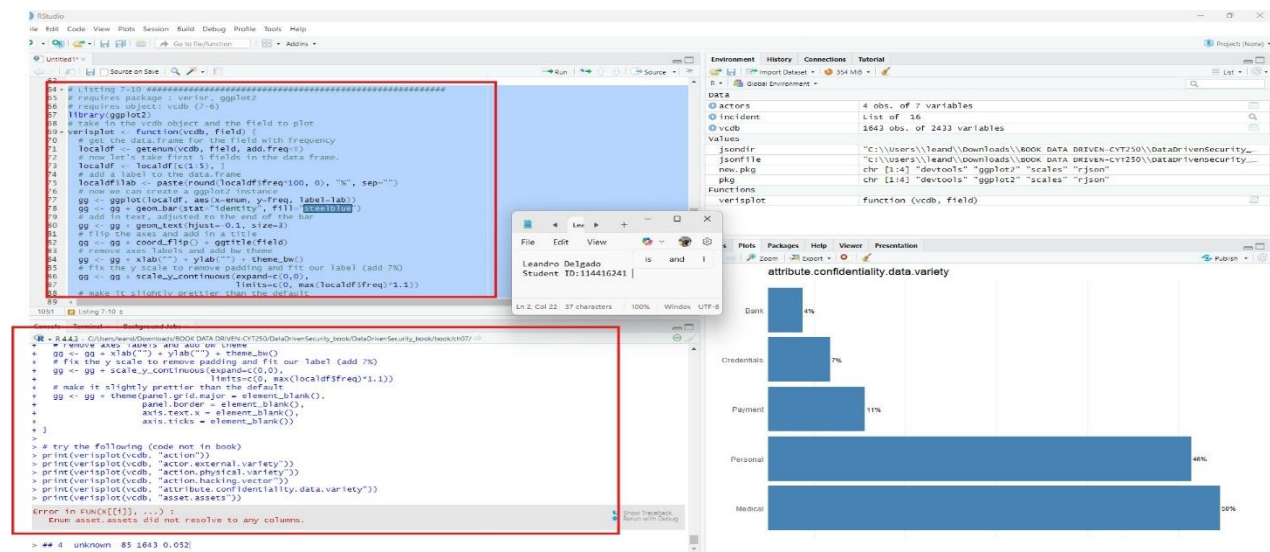


**Figure 12. Result listing 7-10**

**Listing 7-10**

After defining the verisplot() function, I used it to generate several additional charts by specifying different fields from the VCDB dataset. One of the most interesting results came from the action.hacking.vector field, where 93% of the incidents involved web applications, highlighting a major attack surface. I also attempted to visualize asset. assets, but received an error indicating that this field could not be resolved. This likely means the column doesn't exist in the current version of the dataset or is named differently. It was a helpful reminder to verify field names before plotting.
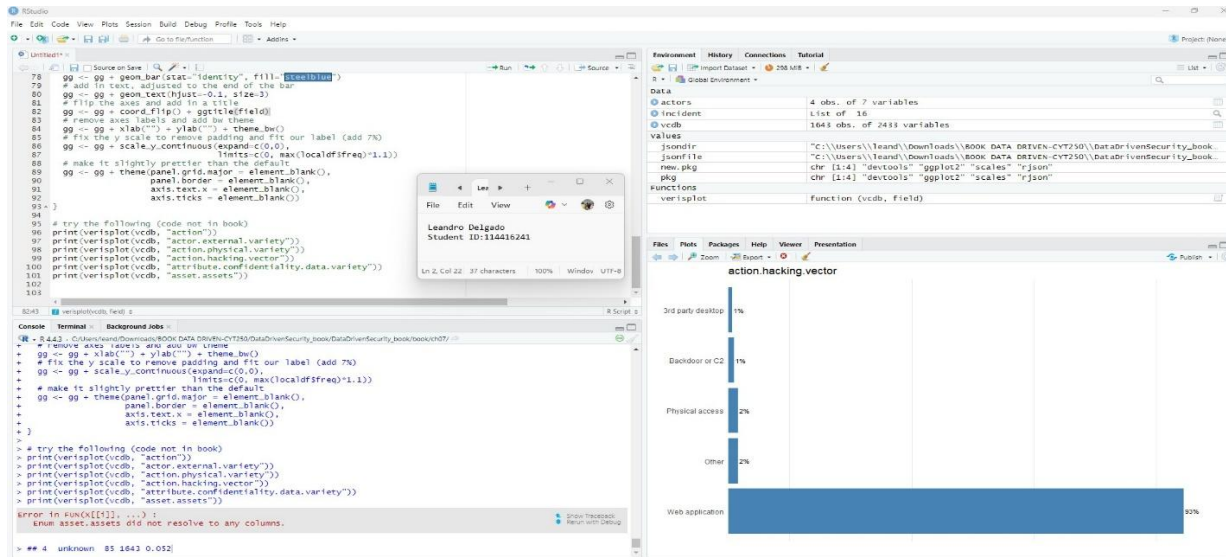


Figure 13. Result run listing 7-10

Using the verisplot() function, I explored the action.physical.variety field. The resulting bar chart shows that theft was the most frequent physical action, accounting for 95% of all such incidents. This highlights how physical security threats like stolen devices remain a major concern in cybersecurity incidents.
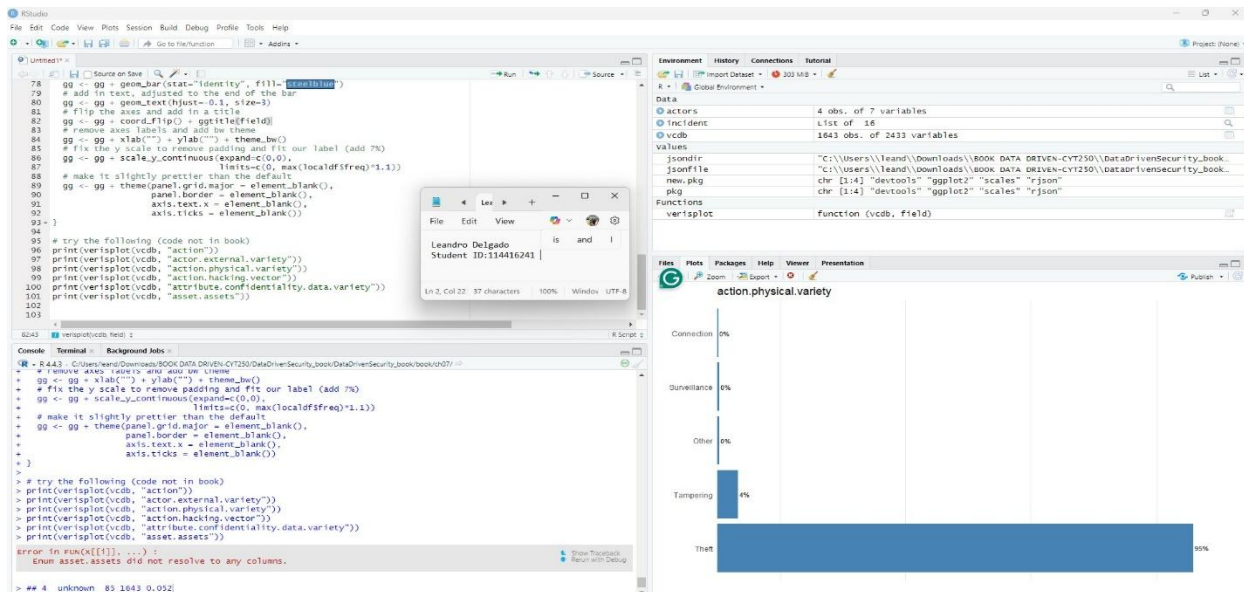


Figure 14. Run listing 7-10

Using the verisplot() function on the actor.external.variety field, I was able to see which types of external actors were most involved in security incidents. The chart revealed that activists were responsible for the majority (54%) of incidents, followed by unaffiliated individuals (25%) and organized crime (10%). This helped me better understand the diversity of external threats represented in the dataset.
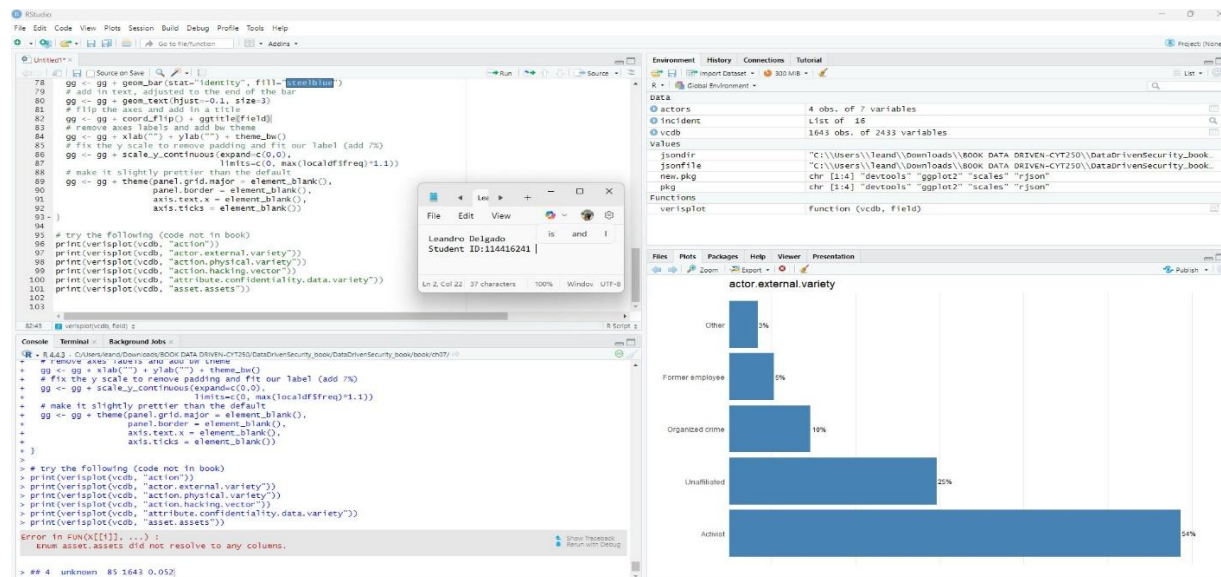


**Figure 15.Result Listing 7-10**

Using the verisplot() function on the action field, I generated a bar chart showing the types of actions involved in the incidents. The data revealed that physical actions (33%), hacking (27%), and errors (26%) were the most common, followed by misuse and malware. This helped me understand the broader categories of how breaches occurred across the dataset.
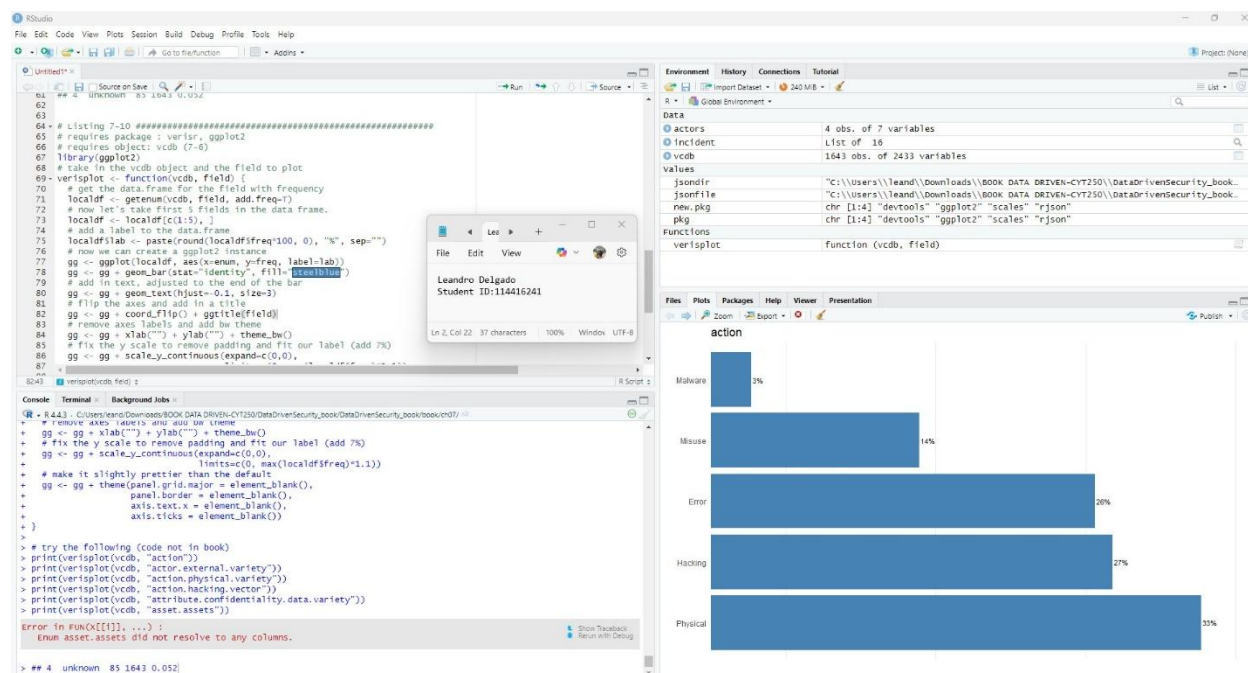


**Figure 16.Result Listing 7-10**

**Listing 7-11**

In this step, I created a heatmap to explore the relationship between different types of actions and the assets they impacted. I grouped the data, counted the most frequent combinations, and used ggplot2 to visualize it. The darker the color, the more frequently that action-asset pair appeared. This helped me quickly identify which actions were most often associated with specific types of compromised assets across the dataset.
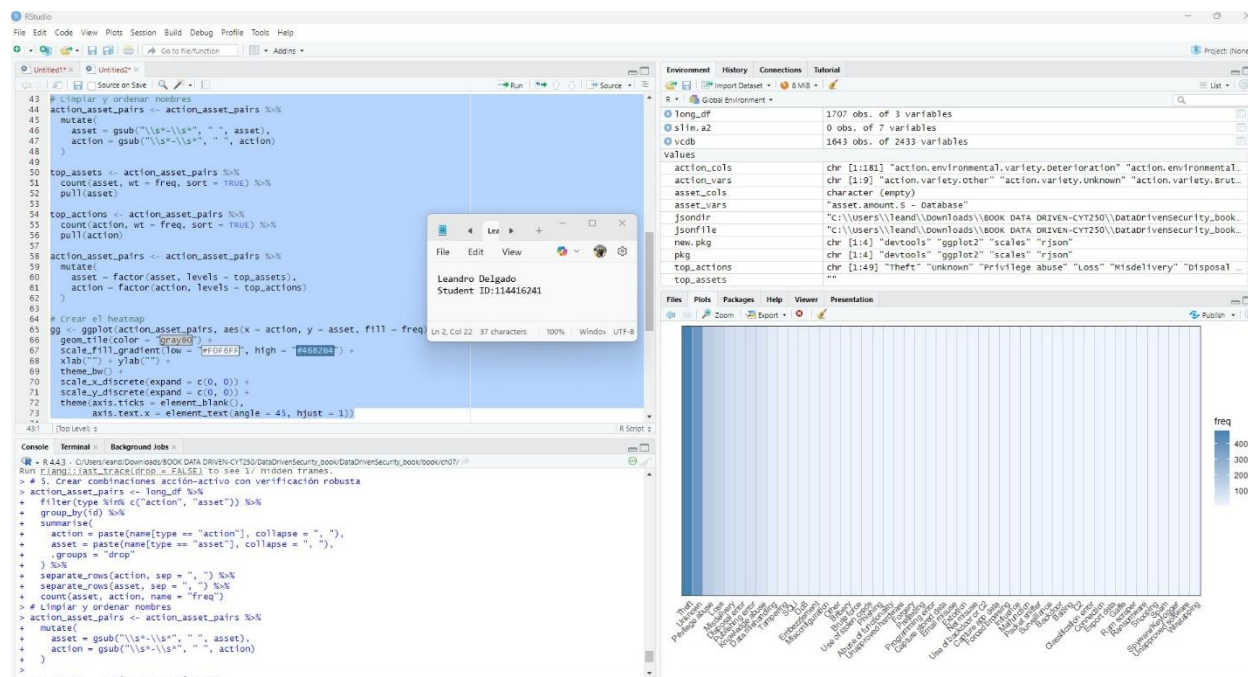


*Figure 17. Result listing 7-11*

I created a refined heatmap to visualize the relationships between different types of assets and the types of actions that affected them. After filtering out unknown and irrelevant categories, I used geom_tile() in ggplot2 to plot a matrix where darker tiles represent higher frequencies of incidents. This allowed me to quickly identify patterns, such as how laptops are most often involved in physical attacks or how web applications are commonly linked to hacking incidents. This final heatmap provided a more complete and insightful view of how different types of cyber actions target specific assets.

**Summary of the Final Project**

Reading through Chapter 7 of Data-Driven Security has given me real-life on hands practice analyzing incident data in real life through R. It introduced me to many technical challenges such as package dependencies, path errors, and undefined objects, all of which have turned out to be my greatest opportunities to sharpen my troubleshooting and data-wrangling skills.

At the end of this project, I was able to import the JSON formatted VCDB data, explore its structure, do some descriptive summaries, and make visible clear views through ggplot2; some of these were frequency bar charts and heatmaps showing patterns between actors, actions, and affected assets.

Such practical learning has made it easy for me to understand VERIS data structures and how to work with and visualize complicated datasets in R. It has prepared me for future cybersecurity data analysis tasks where structured data, incident classification, and insight generation are key.