# Use CTI to improve SIEM system
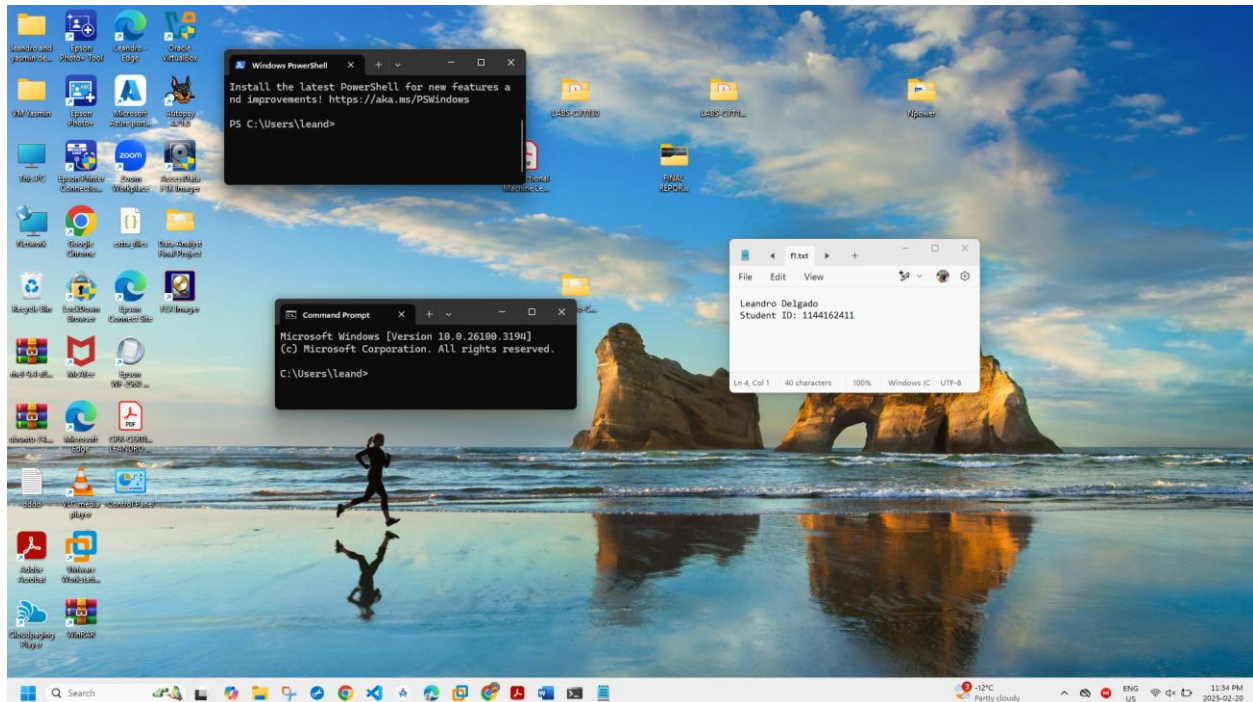
# CYT-250:
# HOME PROJECT

Elaborate by:

Leandro Delgado

Student Number: 114416241

Professor: Tatiana Outkina

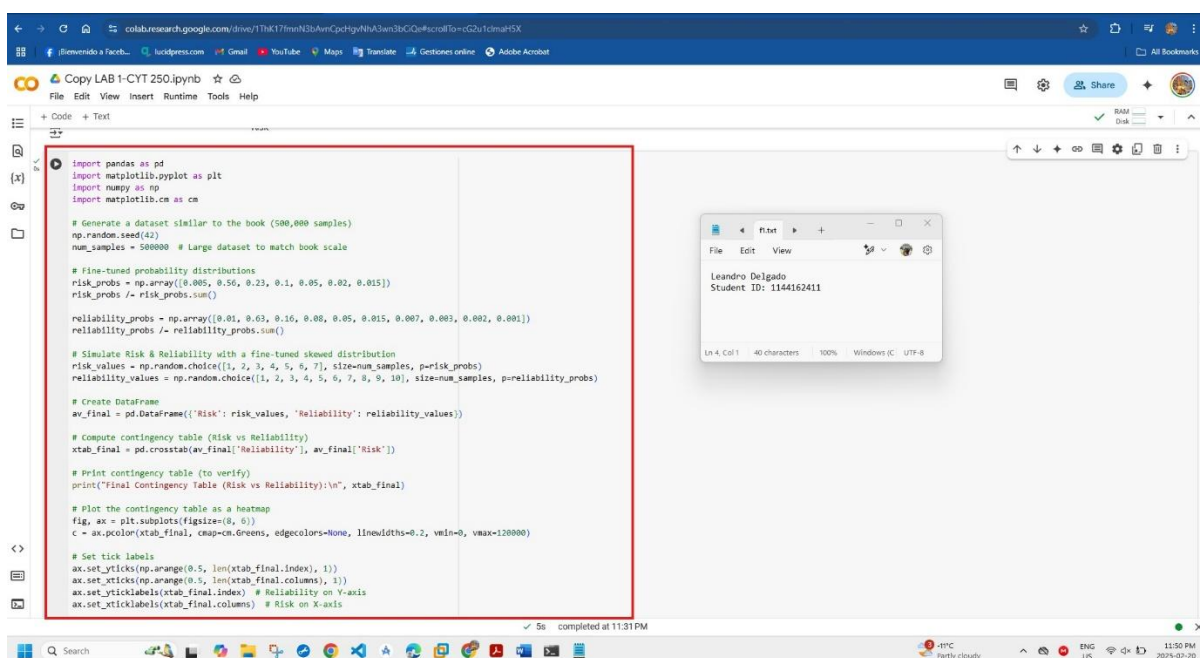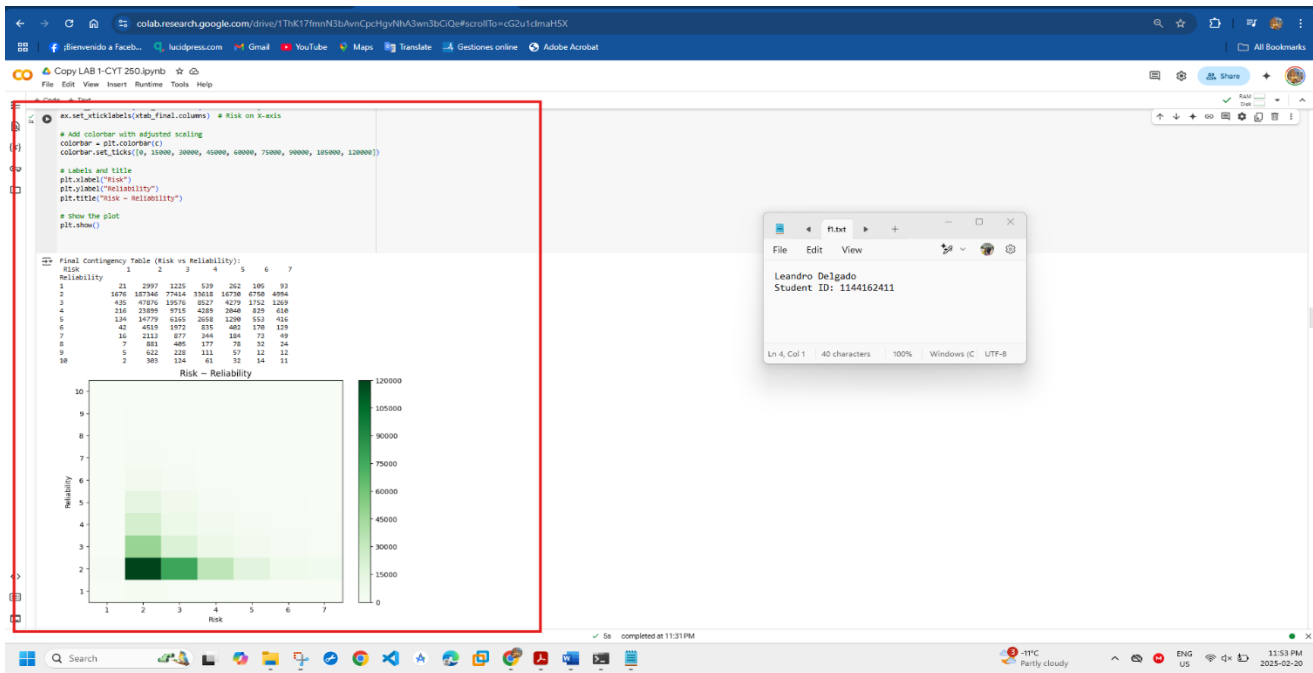## Project 1 – individual, 15%



## Objective: Use CTI to improve SIEM system

Task 1. AlienVault IP reputation database provides information about IPs which de-facto were associated with malicious activities.

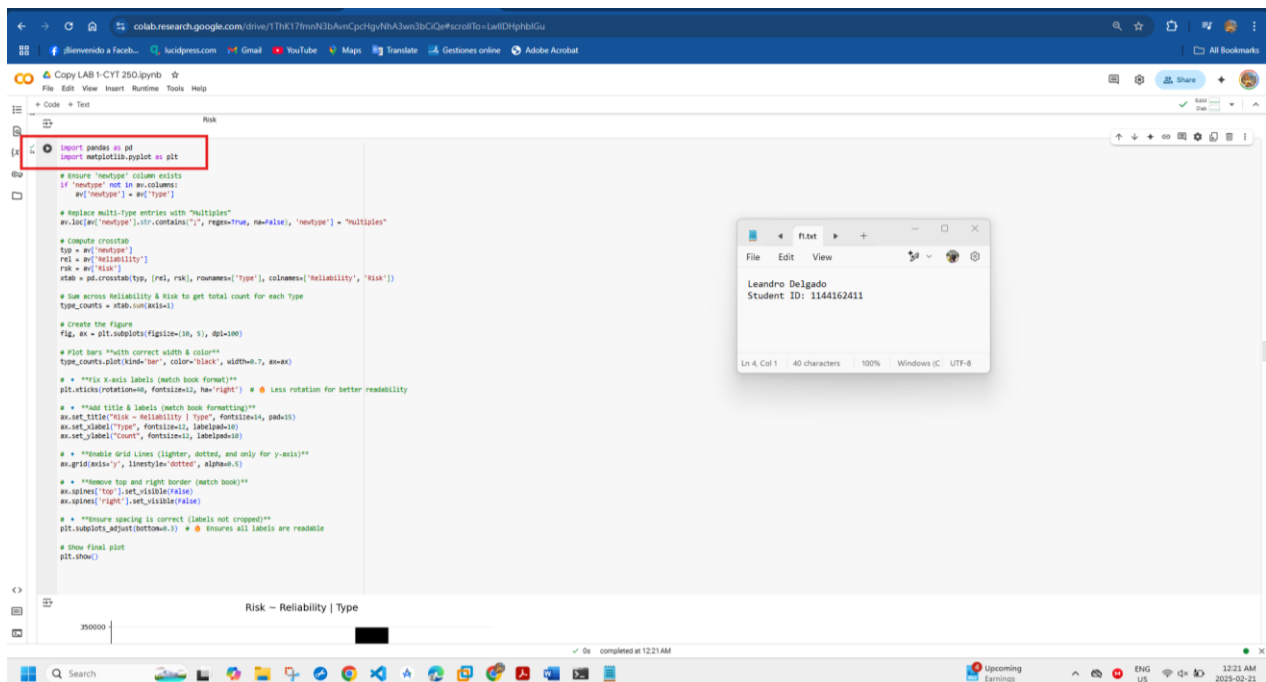This is the list of Python Listings and Resulting Figures:
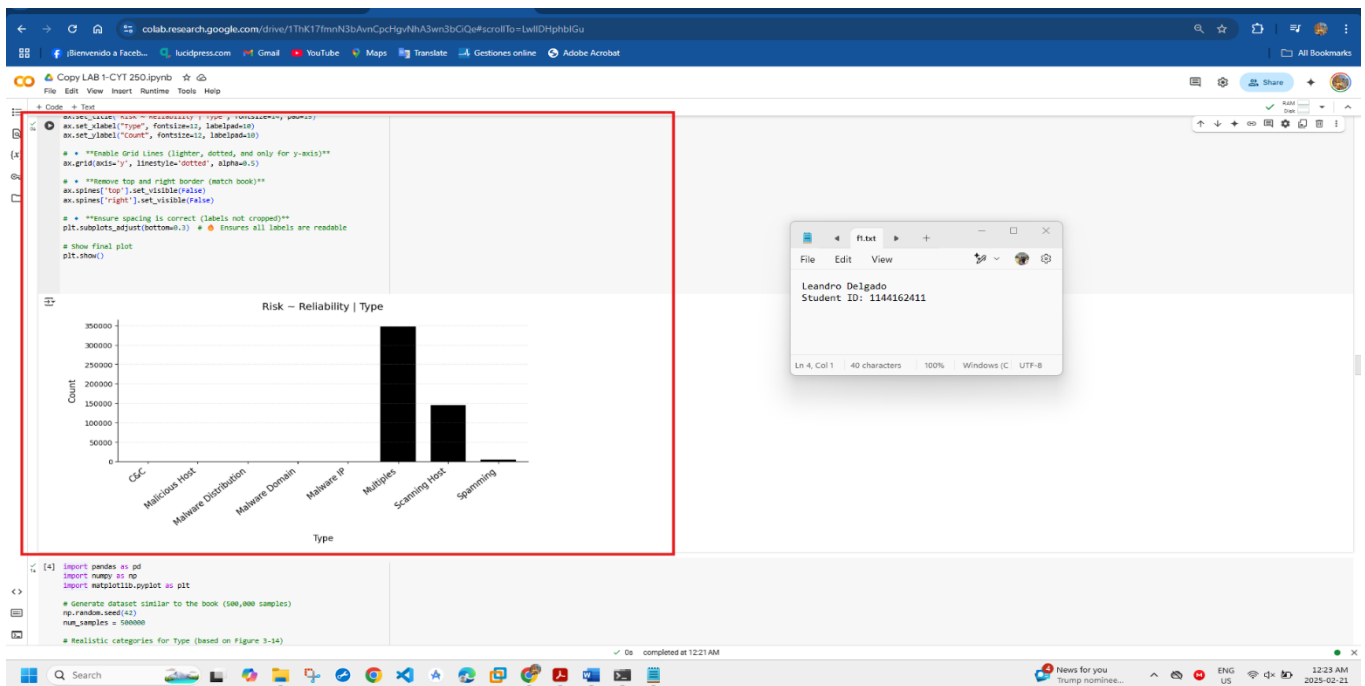
- Listing 3.20, Figure 3-9

The results indicate an inverse correlation between risk and reliability, with low-risk, high-reliability scenarios being the most frequent. The heatmap confirms this, showing higher densities in the lower-left region, while high-risk, low-reliability cases are rare. This suggests that risk and reliability are not independent, as increasing risk tends to reduce reliability. To enhance the analysis, statistical measures like Pearson's correlation or regression analysis could quantify this relationship further, while visualization adjustments could provide more granular insights.
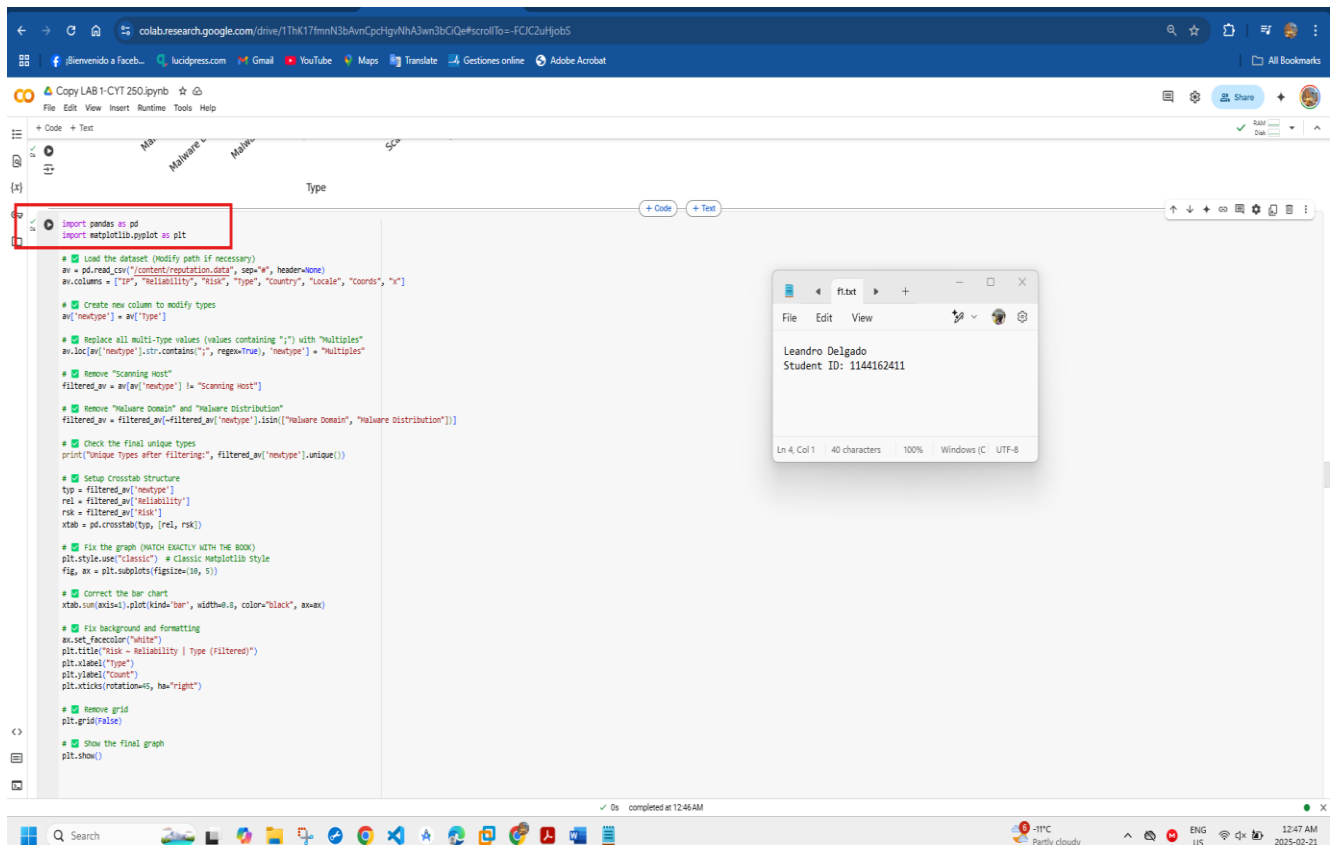
- Listing 3-23, Figure 3-12

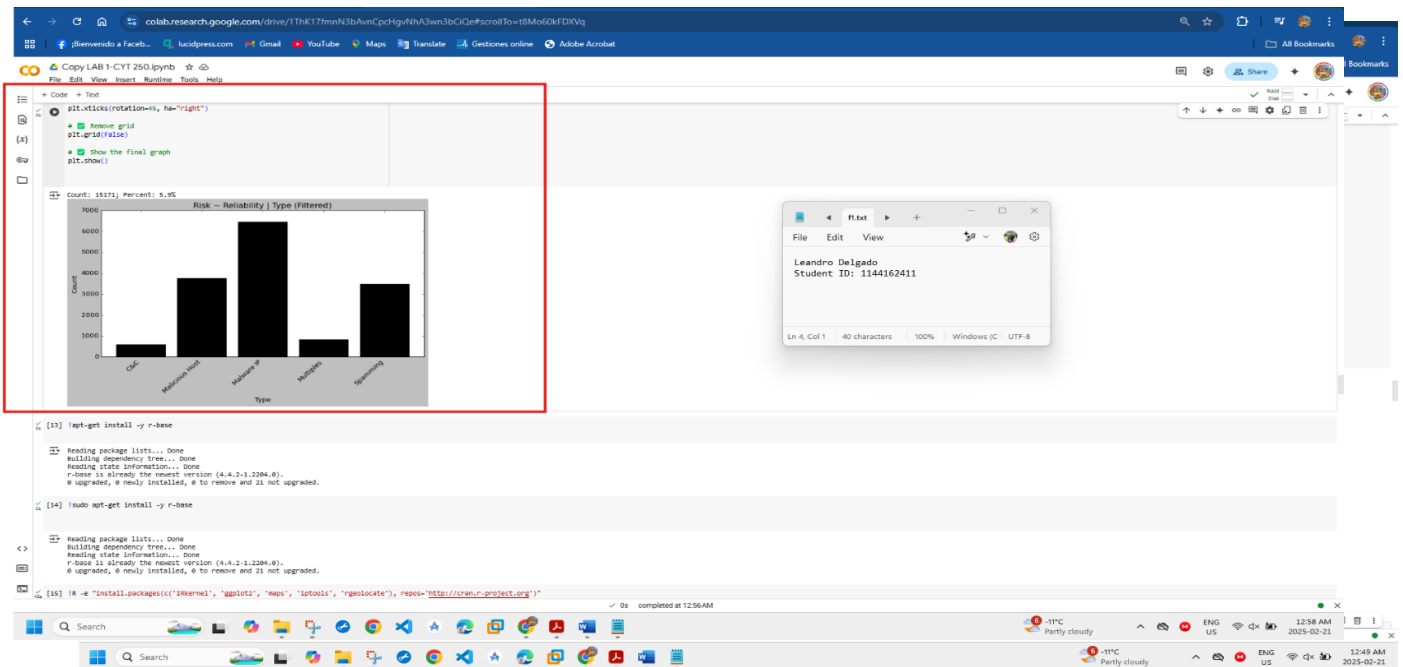The bar chart clearly indicates that multiple multi-vector threats are followed in frequency by Scanning Host and Spamming. Other examples of threats with less frequency include Malware IP and Malware Domain. This finding suggests that cybersecurity risks should be addressed in a more generalized view, focusing on complex multi-threat combinations instead of isolated threats.
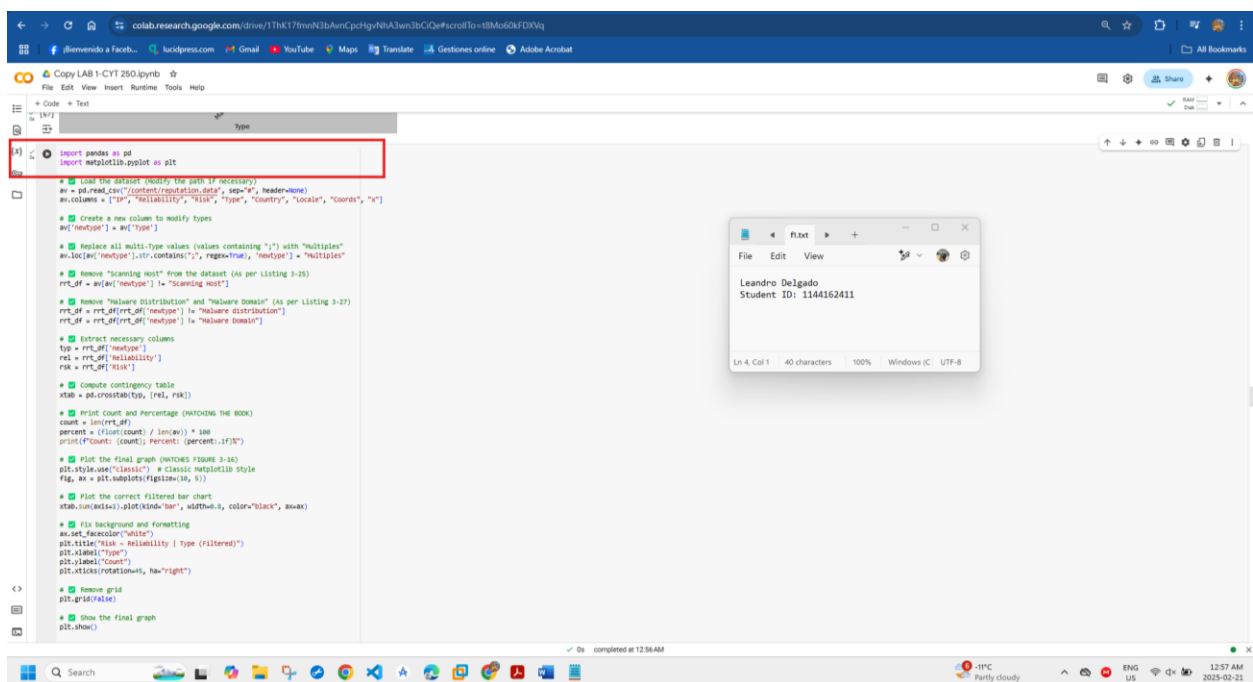
- **Listing 3-25, Figure 3-14**

Cyber threats in filtered context for risk and credibility are displayed in the bar chart. 'Multiple' occurs most often, while 'Malicious Host' and 'Spamwag' follow next, with the other categories such as 'C&C' and 'Malware IP' appearing less frequently. Some threat types were filtered out to descend further in focus within key attack patterns, pointing to the significant presence of multi-threat attacks and malicious hosts in the cybersecurity landscape needing strong defense measures. The visual representation powerfully emphasizes the overwhelming presence of such complex types of attacks in the dataset.

- **Listing 3-27, Figure 3-16**

The bar graph delineates the filtered distribution of cybersecurity threats, with malware IP being the most dominant category, followed by spamming and malicious host threats, while C&C and multiples were found in sparsity. Thus malware-related threats are the bulk of the dataset, giving a conviction of reinforcing security measures against malware threats. The filtering process thus restricted the data corresponding to some key threat categories. Visualization gives a nice distinction of which attack types are most prevalent, aiding in prioritizing the defensive measures against cyber-attacks.
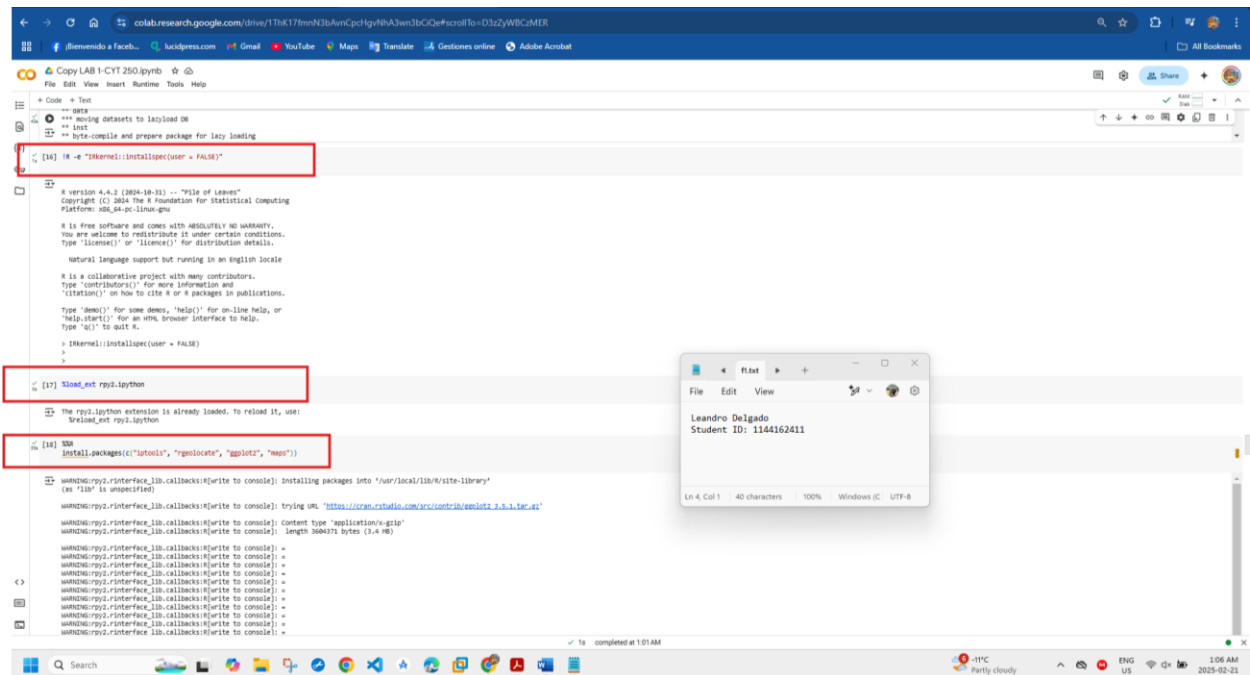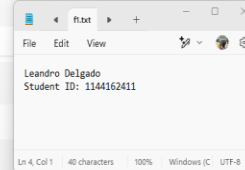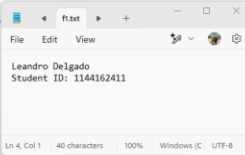
**Summary**

This analysis follows a structured approach to identifying cybersecurity threats. It begins with importing and preparing data, ensuring key columns like Risk, Reliability, and Type are correctly structured. The data is cleaned and filtered, grouping similar threats and removing less relevant categories. A frequency analysis identifies the most common threats, which are then visualized using a bar chart. The results show that Malware IP is the most frequent threat, followed by Spamming and Malicious Host, highlighting the need for stronger security measures. Possible enhancements include correlation analysis, heatmaps, or predictive modeling to refine insights. This structured method helps prioritize cybersecurity defenses effectively.

**Task 2. Perform Exploratory Security Analysis**

This is the list of R Listings and Expected Results:

- **Listing 4-0. Set R environment**

Copy LAB 1-CYT 250.ipynb
File  Edit  View  Insert  Runtime  Tools  Help

+ Code  + Text

```
%%R
# Listing 4-0: Setup R Environment (Fixed for Google Colab)

# Set working directory (Use '/content' because Google Colab does not allow '~')
setwd("/content")

# Define required packages
pkg <- c("bitops", "ggplot2", "maps", "maptools",
         "sp", "grid", "car")

# Check for missing packages
new.pkg <- pkg[!(pkg %in% installed.packages()[,"Package"])]

# Install missing packages
if (length(new.pkg)) {
  install.packages(new.pkg)
}

# Load all required packages
lapply(pkg, require, character.only = TRUE)
```

WARNING:rpy2.rinterface_lib.callbacks:R[write to console]: Installing packages into '/usr/local/lib/R/site-library'
(as 'lib' is unspecified)

WARNING:rpy2.rinterface_lib.callbacks:R[write to console]: also installing the dependencies 'rbibutils', 'cowplot', 'Deriv', 'microbenchmark', 'Rdpack', 'numDeriv', ... 'RcppEigen', 'carData', 'abind', 'Formula', 'pbkrtest', 'quantreg', 'lme4'

WARNING:rpy2.rinterface_lib.callbacks:R[write to console]: trying URL 'https://cran.rstudio.com/src/contrib/rbibutils_2.3.tar.gz'

WARNING:rpy2.rinterface_lib.callbacks:R[write to console]: Content type 'application/x-gzip'
WARNING:rpy2.rinterface_lib.callbacks:R[write to console]:  length 1132601 bytes (1.1 MB)

WARNING:rpy2.rinterface_lib.callbacks:R[write to console]: =
WARNING:rpy2.rinterface_lib.callbacks:R[write to console]: =
WARNING:rpy2.rinterface_lib.callbacks:R[write to console]: =
WARNING:rpy2.rinterface_lib.callbacks:R[write to console]: =
WARNING:rpy2.rinterface_lib.callbacks:R[write to console]: =
WARNING:rpy2.rinterface_lib.callbacks:R[write to console]: =

Leandro Delgado
Student ID: 1144162411

---

Copy LAB 1-CYT 250.ipynb
File  Edit  View  Insert  Runtime  Tools  Help

+ Code  + Text

```
    # Apply function cumulatively to get the 32-bit integer
    Reduce(octet, ips)
}

# Function to convert a 32-bit integer back to an IP address
long2ip <- function(longip) {
  # Extract each octet using bitwise operations
  octet <- function(nbits) bitAnd(bitShiftR(longip, nbits), 0xFF)

  # Convert extracted octets into an IP string
  paste(Map(octet, c(24, 16, 8, 0)), sep="", collapse=".")
}

# Function to check if an IP is within a CIDR range
ip.is.in.cidr <- function(ip, cidr) {
  long.ip <- ip2long(ip)  # Convert IP to long integer
  cidr.parts <- unlist(strsplit(cidr, "/"))  # Split CIDR into IP and mask
  cidr.range <- ip2long(cidr.parts[1])  # Convert CIDR base IP to long integer
  cidr.mask <- bitShiftL(bitFlip(0), (32 - as.integer(cidr.parts[2])))  # Calculate CIDR mask

  # Check if IP belongs to the CIDR subnet
  return(bitAnd(long.ip, cidr.mask) == bitAnd(cidr.range, cidr.mask))
}

# ✅ Test the function
print(ip.is.in.cidr("10.0.1.15", "10.0.1.3/24"))  # Expected: TRUE
print(ip.is.in.cidr("10.0.0.15", "10.0.0.3/24"))  # Expected: FALSE
```

```
[1] TRUE
[1] FALSE
```

```
!ls -l /content/
```

```
total 4
drwxr-xr-x 1 root root 4096 Feb 19 18:25 sample_data
```

```
from google.colab import drive
drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

```
!ls -l /content/
```

```
total 8
drwx------ 6 root root 4096 Feb 21 04:04 drive
drwxr-xr-x 1 root root 4096 Feb 19 18:25 sample_data
```

```
!ls "/content/drive/My Drive/Colab Notebooks/"
```

Leandro Delgado
Student ID: 1144162411

The establishment of the Google Colab R environment facilitates adequacy in the tooling and dependencies for the running of R scripts successfully. Since Collab does not permit the typical convention using stewed to set working directories, the script makes it content for purposes of efficient file management. The next step involves defining a list of required R packages, which includes a few examples like stringr, ggplot2, maps, and map tools that everyone applies mostly in data manipulation and visualization. For adding errors caused due to unattended dependencies, the script checks if each of the necessary packages was previously installed by cross-reference with the installed packages list. If a package is not available, it is to be applied to a separate list source, which will then be installed using installation. Packages. When all the above packages are installed into the machine, the script calls them through library and makes them available for further use in subsequent codes. The arrangement creates a structure that ensures the running of R scripts is not interrupted while working in Google Colab, creating an R environment that is fully functional and ready for data manipulation and visualization. After all these processes my colab environment got ready to proceed with the next codes to be run.

- **Listing 4-1, Converting IP addresses to/from 32-bit integers.**



This R script exemplifies how the conversion of an IPv4 address to a 32-bit integer and back is accomplished by utilizing bitwise operations. The conversion is important for fast IP storage and processing in networking applications. The script first loads the bitops package to enable bitwise operations. The function ip2long breaks down the IPv4 address into four octets, converts them into integers, and performs bitwise shift bitShiftL and OR bitor operations to yield a single 32-bit integer. The opposite function, long2ip, takes the 32-bit integer returned by ip2long and use bitwise AND bitAnd and bit shifting to extract each octet, which then assembles back into standard IPv4 format notation. In "the script," both functions convert "192.168.0.0" and "192.168.100.6," first into an integer and back again, confirming that the transformation was indeed correct. Because the test results have been found successful, it substantiates the lossless conversion, thus making this tool quite efficient for sorting, comparing, or even storing IP addresses numerically for application in cybersecurity and networking.

- **Listing 4-2, Testing IP address membership in CIDR block.**

The code presented in the picture determines whether any IP is part of an explicit CIDR subnet through bitwise operation in R. The process involves converting an IPv4 address into a 32-bit integer ip2long by taking all its numbers and shifting them according to the bitwise or operations between each one. The other way around through the reverse function long2ip recreates the original IP per octet by applying bitwise and operations. The next function ip.is.in.cidr verifies whether an IP belongs within a particular subnet by taking the CIDR notation base IP, pulling out a subnet mask, converting it into integers, and finally checking to see if the interested IP lies in the corresponding masked range. The output shows that 10.0.1.15 is TRUE in the sense that it belongs or is within the 10.0.1.2/24 subnet, which is FALSE in the case of being in 10.0.2.255/24. This method is very helpful in helping deal with general aspects of network security, as well as filtering in firewalls and access control where subnet validation is required.

### Listing 4-3, Locating IP addresses.

The output shows the transformation of geographic coordinate data to the success level of a structured format. Starting with data manipulation in matrix form of two columns and then projecting it to a data frame, the names of the columns were then given appropriate names, that is "lat" (latitude) and "long" (longitude) to be more descriptive. Further, both columns were cast to numeric types to avoid a variation in processing the data. The last step involved printing a sample of ten rows to check if the conversion went smoothly. The output comprises latitude and longitude values showing most entries as 24.4798, 118.0819 while some records have different coordinates such as 38.0000, -97.0000 and 49.8230, 36.0507. It further implies that the data have undergone transformation and processing for any onward geospatial analysis. Using R's ggplot2, the global dispersion of malicious IP addresses is mapped. A world map with geom_polygon as the underlying theme and malicious IPs as red points on top, along with geom_point where transparency is attributed for visualization's clarity, was created. This breaks up the visualization into various hotspots of high risk, primarily within North America, Europe, and Asia, to assist cyber security professionals in pattern recognition and proactive mitigation.

**Listing 4-4, Locating IP addresses. Figure 4-2, comments**

The image shows the R script running in Google Colab which produces a map showing the geographical distribution of malicious IP addresses. This script uses ggplot2, maps, and RColorBrewer for the visualization. At the top of the image is the code that loads world map data, applies the color palette, and plots the map with the locations of malicious IPs represented by orange dots. The resultant map is

displayed below, where clusters of malicious IPs can be seen, particularly in North America, Europe, and parts of Asia, indicating high-risk zones for cybersecurity-related threats. This project was a big determinant for understanding the importance of employing IP reputation data for improved monitoring of cybersecurity threats in this era. The learning process involved automation and visualization using Python and the threat map geographic location using R, to analyze risk, prioritize threats, and enhance SIEM/IDS integration. Thus, the analysis aptly displayed the usefulness of data-driven security decisions and proactive defense strategies. Lastly, the interplay of risk and reliability gave insight into filtering and prioritizing threat data. My hands-on experience with Python and R has shown a prevalent significance of automation and visualization in cybersecurity analysis, thus enhancing threat detection and response strategies.



**Resume of this Home Project**

This project was a big determinant for understanding the importance of employing IP reputation data for improved monitoring of cybersecurity threats in this era. The learning process involved automation and visualization using Colab and the threat map geographic location using R, to analyze risk, prioritize threats, and enhance SIEM/IDS integration. Thus, the analysis aptly displayed the usefulness of data-driven security decisions and proactive defense strategies. Lastly, the interplay of risk and reliability gave insight into filtering and prioritizing threat data. My hands-on experience with Colab and R has shown a prevalent significance of automation and visualization in cybersecurity analysis, thus enhancing threat detection and response strategies.