

Student Name : _____ Timothy Chang _____
 Group : _____ SCSB _____
 Date : _____ 22/03/2024 _____

LAB 4: ANALZING NETWORK DATA LOG

You are provided with the data file, in .csv format, in the working directory. Write the program to extract the following informations.

EXERCISE 4A: TOP TALKERS AND LISTENERS

One of the most commonly used function in analyzing data log is finding out the IP address of the hosts that send out large amount of packet and hosts that receive large number of packets, usually know as TOP TALKERS and LISTENERS. Based on the IP address we can obtained the organization who owns the IP address.

List the TOP 5 TALKERS

Rank	IP address	# of packets	Organisation
1	193.62.192.8	3041	European Bioinformatics Institute
2	155.69.160.32	2975	Nanyang Technological University
3	130.14.250.11	2604	National Library of Medicine
4	14.139.196.58	2452	Indian Institute of Technology, Guwahati
5	140.112.9.139	2056	National Taiwan University

TOP 5 LISTENERS

Rank	IP address	# of packets	Organisation
1	103.37.198.100	3841	Agency For Science, Technology and Research
2	137.132.228.15	3715	National University of Singapore
3	202.21.159.244	2446	Republic Polytechnic, Singapore

4	192.101.107.153	2368	Battelle Memorial Institute, Pacific Northwest Division
5	103.21.126.2	2056	Indian Institute of Technology Bombay

EXERCISE 4B: TRANSPORT PROTOCOL

Using the IP protocol type attribute, determine the percentage of TCP and UDP protocol

	Header value	Transport layer protocol	# of packets
1	6	TCP	56064 (82.37%)
2	17	UDP	9462 (13.90%)
3	50	ESP	1698 (2.49%)
4	47	GREs	657 (0.97%)

EXERCISE 4C: APPLICATIONS PROTOCOL

Using the Destination IP port number determine the most frequently used application protocol.
(For finding the service given the port number <https://www.adminsub.net/tcp-udp-port-finder/>)

Rank	Destination IP port number	# of packets	Service
1	443	13423	HTTPS
2	80	2647	HTTP
3	52866	2068	Dynamic / Private ports by IANA, XSAN.
4	45512	1356	Unassigned Ports
5	56152	1341	Dynamic / Private ports by IANA, XSAN

EXERCISE 4D: TRAFFIC

The traffic intensity is an important parameter that a network engineer needs to monitor closely to determine if there is congestion. You would use the IP packet size to calculate the estimated total traffic over the monitored period of 15 seconds. (Assume the sampling rate is 1 in 2048)

Total Traffic (MB)	<u>126519.184</u>
--------------------	-------------------

EXERCISE 4E: ADDITIONAL ANALYSIS

Please append ONE page to provide additional analysis of the data and the insight it provides.
Examples include:

Top 5 communication pairs;

Visualization of communications between different IP hosts;
etc.

Please limit your results within one page (and any additional results that fall beyond one page limit will not be assessed).

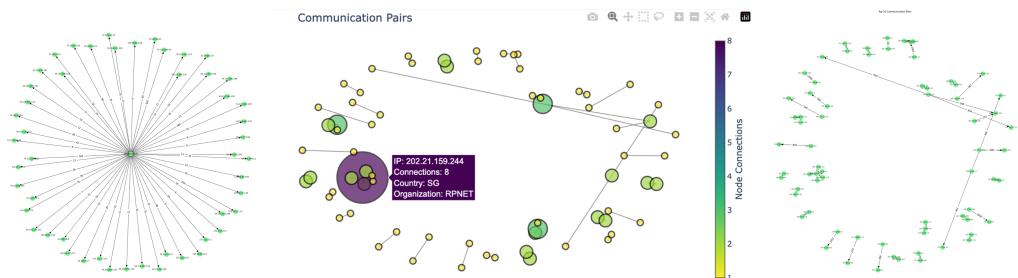
Top 5 Communication pairs

Source: 193.62.192.8 (European Bioinformatics Institute)	Destination: 137.132.228.15 (National University of Singapore)
Packet count: 3041	
Source: 130.14.250.11 (National Library of Medicine)	Destination: 103.37.198.100 (A*STAR)
Packet count: 2599	
Source: 14.139.196.58 (Indian Institute of Technology, Pacific Northwest Division)	Destination: 192.101.107.153 (Battelle Memorial Institute, Pacific Northwest Division)
Packet count: 2368	
Source: 140.112.8.139 (Bombay)	Destination: 103.21.126.2 (Indian Institute of Technology Bombay)
Packet count: 2056	
Source: 137.132.228.15 (National University of Singapore)	Destination: 193.62.192.8 (European Bioinformatics Institute)
Packet count: 1910	

As seen for the results above, there seems to be the most amount of communication between the European Bioinformatics Institute (EBI) and National University of Singapore. There could be several factors that explain this

- Research Collaboration: EBI is a prominent research institute in bioinformatics, while NUS is a leading university with strong research programs in various life science fields. They might be collaborating on research projects that involve exchanging large datasets or frequent communication for analysis.
- Data Sharing: EBI provides a wealth of biological data resources publicly. NUS researchers might be heavily accessing and downloading these datasets for their own research, leading to high traffic towards EBI.
- Cloud Resources: EBI might be using cloud computing services offered by Singapore, with NUS being a potential access point. This could explain high traffic flow if EBI is transferring or processing large datasets on these cloud resources.
- Specific Applications: Both EBI and NUS might be using specialized bioinformatics software or online platforms that rely on heavy data transfer between their servers.

Visualisation



The leftmost figure is a visual representation of network communication patterns, focusing on the top 50 communication pairs initiated by a specific IP address. The directed graph shows the direction of communication flow, node sizes represent communication volume, and edge labels indicate the number of packets exchanged between each pair. This visualization helps to identify communication hubs, high-traffic pathways, and potential bottlenecks within the network. The rightmost figure is similar to the leftmost one, but instead focuses on displaying the communication pairs of the dataset, revealing a broader view of network interactions. The middle figure builds on the analysis of the two, but is done in a more informative and interactive way. It allows for quick identification of high-activity nodes and exploration of associated details on hover.

Republic Polytechnic seems to have the most amount of connections, especially to sites like Microsoft and Facebook. Further investigation into a specific IP address, '155.69.160.32', revealed it to be part of our institution's network (NTU-3) managed by APNIC. Analyzing this node's most frequent destination IP addresses identified the top four as belonging to Google LLC. While initial assumptions pointed towards these connections originating primarily from Google Search, a deeper examination of the data revealed a more nuanced picture.

Three out of the four Google-related IP addresses predominantly utilized ports 80 (HTTP) and 443 (HTTPS), suggesting web browsing activity. However, the remaining IP address employed port 19305, which is commonly associated with Google Talk. This suggests that a portion of the communication with Google may involve both web browsing and potentially the use of Google Talk software.

In conclusion, the analysis identified active communication hubs and Google as a major destination for web traffic, with a possibility of Google Talk usage from NTU-3.

EXERCISE 4F: SOFTWARE CODE

Please also submit your code to the NTULearn lab site.

```
● ● ●  
1 #Installing dependencies  
2  
3 %pip install pandas  
4 %pip install ipwhois  
5 %pip install seaborn  
6 %pip install networkx  
7 %pip install plotly
```



```
● ● ●  
1 # Package imports  
2  
3 # Import pandas library for data manipulation (analysis and wrangling)  
4 import pandas as pd  
5 # Import requests library to potentially download data from web sources  
6 import requests  
7 # Import math library for mathematical operations (e.g., ratios, percentages)  
8 import math  
9 # Import IPWhois class from ipwhois library to lookup IP address information  
10 from ipwhois import IPWhois  
11 from ipwhois.exceptions import IPDefinedError  
12 # Ignore warnings  
13 import warnings  
14 warnings.filterwarnings('ignore')
```



```
● ● ●  
1 # This code snippet prepares a pandas DataFrame for network traffic analysis from a CSV file.  
2  
3 # 1. Define meaningful column names based on Table 1 in the lab manual. These names correspond to the fields in the traffic data.  
4 # 2. Read the CSV file 'Data_3.csv' using pandas. Assign temporary placeholder names (a-u) for columns during this process.  
5 # 3. Select all columns except the last one. This might be done because the last column could be empty or contain irrelevant data.  
6 # 4. Assign the defined column names (col_names) to the DataFrame. This makes the data more readable and easier to understand.  
7 # 5. Display the DataFrame containing the network traffic data with proper column names.  
8  
9  
10 col_names = [  
11     "type", "sfloew_agent_address", "inputPort", "outputPort", "src_MAC",  
12     "dst_MAC", "etherent_type", "in_vlan", "out_vlan", "src_IP", "dst_IP",  
13     "IP_protocol", "ip_tos", "ip_ttl", "src_transport_port", "dst_transport_port",  
14     "tcp_flags", "packet_size", "IP_size", "sampling_rate"  
15 ]  
16 df = pd.read_csv('Data_3.csv', names = list('abcdefghijklmopqrstuvwxyz'))  
17 df = df.iloc[:, :-1]  
18 df.columns = col_names  
20 df
```

type	sflow_agent_address	inputPort	outputPort	src_MAC	dst_MAC	ether_type	in_vlan	out_vlan	src_IP	dst_IP	ip_protocol	ip_tos	ip_ttl	src_transport_port	dst_transport_port	tcp_flags	packet_size	ip_size	sampling_rate	
0	FLOW	203.30.38.251	137	200	d404ff55fd4d	8071fc76001	0x0800	919	280	130.246.176.22	140.115.32.81	6	0x00	50	81216	23505	0x10	1518	1500	2048
1	FLOW	203.30.38.251	129	193	609cb9851b00	031466b23cf	0x0800	11	919	195.69.160.33	84.233.188.128	6	0x00	96	23159	80	0x10	74	62	2048
2	FLOW	203.30.38.251	137	200	d404ff55fd4d	8071fc76001	0x0800	919	280	130.246.176.53	140.115.32.83	6	0x00	50	50641	20729	0x10	1518	1500	2048
3	FLOW	203.30.38.251	129	193	609cb9851b00	02688bd9cf	0x0800	11	919	195.69.160.32	54.168.174.79	17	0x00	120	54241	28610	0x10	118	94	2048
4	FLOW	203.30.38.251	130	199	00398cd97cf	1448dcf9a7df	0x0800	919	600	137.182.228.16	193.62.193.8	6	0x00	56	53923	34262	0x10	70	52	2048
69365	FLOW	203.30.38.251	258	199	20447fcfb0f	cce4857044	0x0800	537	600	207.241.228.157	210.48.222.9	6	0x00	56	443	57434	0x10	1522	1500	2048
69366	FLOW	203.30.38.251	131	193	00a7422339e	031466b23cf	0x0800	43	919	192.123.13.136	216.68.203.234	6	0x00	121	4920	443	0x10	1442	1420	2048
69367	FLOW	203.30.38.251	130	193	00239cd97cf	1448dcf9a7df	0x0800	919	600	137.182.228.16	193.62.192.8	6	0x00	56	34153	3123	0x10	82	64	2048
69368	FLOW	203.30.38.251	129	193	609cb9851b00	031466b23cf	0x0800	11	919	195.69.196.9	74.125.56.6	17	0x00	58	56221	60786	0x10	1267	1245	2048
69369	FLOW	203.30.38.251	137	200	d404ff55fd4d	8071fc76001	0x0800	919	280	14.139.196.58	192.101.107.153	6	0x00	57	34625	41211	0x10	58	40	2048

Exercise 4A

```
● ● ●  
1 # Identify the top 5 talkers (senders) and listeners (receivers) in the network traffic data  
2 # Talkers are those who send the most traffic, listeners receive the most.  
3 # Get top 5 IP addresses based on source IP count. Get top 5 IP addresses based on destination IP count  
4 # Define a function to retrieve the organization associated with an IP address.  
5 # This function utilizes the IPWhois library to query WHOIS information.  
6 # Create an IPWhois object. Perform a RDAP lookup (WHOIS protocol for IP info).  
7 # Extract and return the 'asn_description' field  
8  
9  
10 talkers = df['src_IP'].value_counts().head(5)  
11 listeners = df['dst_IP'].value_counts().head(5)  
12  
13 def org(ip_address):  
14     ipwhois = IPWhois(ip_address)  
15     result = ipwhois.lookup_rdap()  
16     return result['asn_description']
```

```
● ● ●  
1 # Print results for Top 5 Talkers (rank, IP, packets, organization)  
2 # Lookup organization for this IP and print results  
3 # Print results for Top 5 Listeners (rank, IP, packets, organization)  
4 # Lookup organization for this IP and print results  
5  
6 # Retrieve the top 5 source IP addresses and their respective packet counts  
7 top_source_ips = df['src_IP'].value_counts().head(5)  
8  
9 # Retrieve the top 5 destination IP addresses and their respective packet counts  
10 top_dest_ips = df['dst_IP'].value_counts().head(5)  
11  
12 # Function to retrieve the organization associated with an IP address  
13 def get_organization(ip_address):  
14     try:  
15         ipwhois = IPWhois(ip_address)  
16         result = ipwhois.lookup_rdap()  
17         return result['asn_description']  
18     except:  
19         return 'Unknown'  
20  
21 # Print top 5 talkers with their packet counts and organizations  
22 print("Top 5 Talkers")  
23 print("Rank\tIP Address\t\t# of Packets\tOrganization")  
24 for index, (ip_address, packets) in enumerate(top_source_ips.items(), start=1):  
25     organization = get_organization(ip_address)  
26     print(f"{index}\t{ip_address}\t{packets}\t{organization}")  
27  
28 # Print top 5 listeners with their packet counts and organizations  
29 print("\nTop 5 Listeners")  
30 print("Rank\tIP Address\t\t# of Packets\tOrganization")  
31 for index, (ip_address, packets) in enumerate(top_dest_ips.items(), start=1):  
32     organization = get_organization(ip_address)  
33     print(f"{index}\t{ip_address}\t{packets}\t{organization}")
```

Top 5 Talkers			
Rank	IP Address	# of Packets	Organization
1	193.62.192.8	3041	JANET Jisc Services Limited, GB
2	155.69.160.32	2975	NTU-AS-AP Nanyang Technological University, SG
3	130.14.250.11	2604	NLM-GW, US
4	14.139.196.58	2452	NKN-EDGE-NW NKN EDGE Network, IN
5	140.112.8.139	2056	Unknown

Top 5 Listeners			
Rank	IP Address	# of Packets	Organization
1	103.37.198.100	3841	A-STAR-AS-AP A-STAR, SG
2	137.132.228.15	3715	NUS-AS-AP NUS Information Technology, SG
3	202.21.159.244	2446	REPUBLICPOLYTECHNIC-AS Republic Polytechnic. Multihoming AS Singapore, SG
4	192.101.107.153	2368	ESNET-AS, US
5	103.21.126.2	2056	IITB-IN Powai, IN

Exercise 4B



```
1 # Analyze distribution of protocols used in the network traffic data
2
3 protocols = df["IP_protocol"].value_counts()
4
5 protocols
```

IP_protocol

```
6      56064
17     9462
50    1698
0     1261
47     657
41     104
1      74
381    45
58      4
103    1
Name: count, dtype: int64
```



```
1 # Calculate packet counts and percentages for TCP and UDP protocols
2 # Calculate TCP and UDP packet percentages relative to total packets
3
4 tcp_count = protocols[6] # TCP = 6
5 udp_count = protocols[17] # UDP = 17
6
7 tcp_percentage = tcp_count / len(df.index) * 100
8 udp_percentage = udp_count / len(df.index) * 100
9
10 print("Transport layer protocol\t# of packets")
11 print(f"TCP\t\t\t{tcp_count} ({tcp_percentage:.2f}%)")
12 print(f"UDP\t\t\t{udp_count} ({udp_percentage:.2f}%)")
```

Transport layer protocol	# of packets
TCP	56064 (80.82%)
UDP	9462 (13.64%)

Exercise 4C

```
● ● ●  
1 # Function to lookup the service name associated with a port number  
2 # ports description according to this website - https://www.webopedia.com/reference/well-known-tcp-port-numbers/  
3 # and this website - https://www.adminsub.net/  
4  
5 def checkPort(port):  
6     if port == 443:  
7         service = "HTTPS"  
8     elif port == 80:  
9         service = "HTTP"  
10    elif port == 52866:  
11        service = "Dynamic and/or Private Ports by IANA, Xsan. Xsan Filesystem Access by Apple Inc."  
12    elif port == 45512:  
13        service = "Unassigned"  
14    elif port == 56152:  
15        service = "Dynamic and/or Private Ports by IANA, Xsan. Xsan Filesystem Access by Apple Inc."  
16    else:  
17        service = "Unknown"  
18    return service
```

```
● ● ●  
1 # Print table header for Top 5 Destination Ports  
2 # Iterate through the top 5 destination ports  
3 # Lookup service name using the checkPort function  
4 # Print results in a formatted table-like manner  
5  
6  
7 print("Rank\tDestination IP port number\t# of packets\tService")  
8 for i, (port, count) in enumerate(top_ports.items(), start=1):  
9     service = checkPort(port)  
10    print(f"{i}\t{port}\t{count}\t{service}")
```

Rank	Destination IP port number	# of packets	Service
1	443	13423	HTTPS
2	80	2647	HTTP
3	52866	2068	Dynamic and/or Private Ports by IANA, Xsan. Xsan Filesystem Access by Apple Inc.
4	45512	1356	Unassigned
5	56152	1341	Dynamic and/or Private Ports by IANA, Xsan. Xsan Filesystem Access by Apple Inc.

Exercise 4D

```
● ● ●  
1 # Calculate total traffic size in Megabytes (MB)  
2  
3 total_packet_size = df["IP_size"].sum() * 2048 / (1024**2) # 1024**2 is 1MB  
4 print(f"Total traffic: {total_packet_size: .3f} MB")
```

```
Total traffic: 126519.184 MB
```

Exercise 4E

```
● ● ●  
1 # Analyze traffic data initiated from source IP address '155.69.160.32'  
2 # Filter for this source IP and get the top destinations it sent to  
3 # Print results  
4  
5 # 155.69.160.32 belongs to NTU-3, provided by APNIC  
6  
7 df[df['src_IP'] == '155.69.160.32']['dst_IP'].value_counts().head(5)
```

```
dst_IP  
74.125.200.95      520  
172.217.27.14      454  
64.233.188.128    348  
74.125.200.127    161  
137.222.0.78       80  
Name: count, dtype: int64
```

```
● ● ●  
1 # Create a new DataFrame (df_new) containing traffic data from a specific source IP address  
2  
3 df_new = df[df['src_IP'] == '155.69.160.32']
```

```
dst_transport_port  
443      520  
Name: count, dtype: int64
```

```
● ● ●  
1 # Analyze destination ports used in traffic sent to '172.217.27.14'  
2  
3 df_new[df_new['dst_IP'] == '172.217.27.14']['dst_transport_port'].value_counts()
```

```
dst_transport_port  
443      451  
80       3  
Name: count, dtype: int64
```

```
● ● ●  
1 # Analyze destination ports used in traffic sent to '64.233.188.128'  
2  
3 df_new[df_new['dst_IP'] == '64.233.188.128']['dst_transport_port'].value_counts()
```

```
dst_transport_port  
443      253  
80       95  
Name: count, dtype: int64
```

```
● ● ●  
1 # Analyze destination ports used in traffic sent to '74.125.200.127'  
2  
3 df_new[df_new['dst_IP'] == '74.125.200.127']['dst_transport_port'].value_counts()
```

```
dst_transport_port  
19305    161  
Name: count, dtype: int64
```

```
● ● ●  
1 # Analyze destination ports used in traffic sent to '137.222.0.78'  
2  
3 df_new[df_new['dst_IP'] == '137.222.0.78']['dst_transport_port'].value_counts()
```

```
dst_transport_port  
443      80  
Name: count, dtype: int64
```

```
● ● ●  
1 # Analyze traffic destinations for source IP address '155.69.160.32'  
2  
3 # 1. Filter data to include only traffic initiated by source IP '155.69.160.32'  
4 # 2. Group the filtered data by destination IP address  
5 # 3. Count the number of packets sent to each destination IP address within each group created in step 2  
6 # 4. Convert the grouped results (destination & counts) into a DataFrame and rename the automatically generated column name to "count" for clarity  
7 # 5. Sort the DataFrame by the "count" column in descending order. This prioritizes destinations that received the most traffic  
8 # Now 'grouped_df' is a DataFrame containing the destination IP addresses and the corresponding number of packets sent to each destination from '155.69.160.32'  
9  
10 grouped_df = df[df['src_IP'] == '155.69.160.32'].groupby(by='dst_IP').size().reset_index(names="count").sort_values("count", ascending=False)  
11 # [['dst_IP', 'inputPort', 'outputPort', 'src_transport_port', 'dst_transport_port']]  
12 grouped_df
```

	dst_IP	count
246	74.125.200.95	520
36	172.217.27.14	454
205	64.233.188.128	348
229	74.125.200.127	161
17	137.222.0.78	80
...
119	52.74.58.153	1
121	52.76.136.147	1
122	52.76.201.57	1
124	52.76.235.25	1
260	89.106.36.184	1

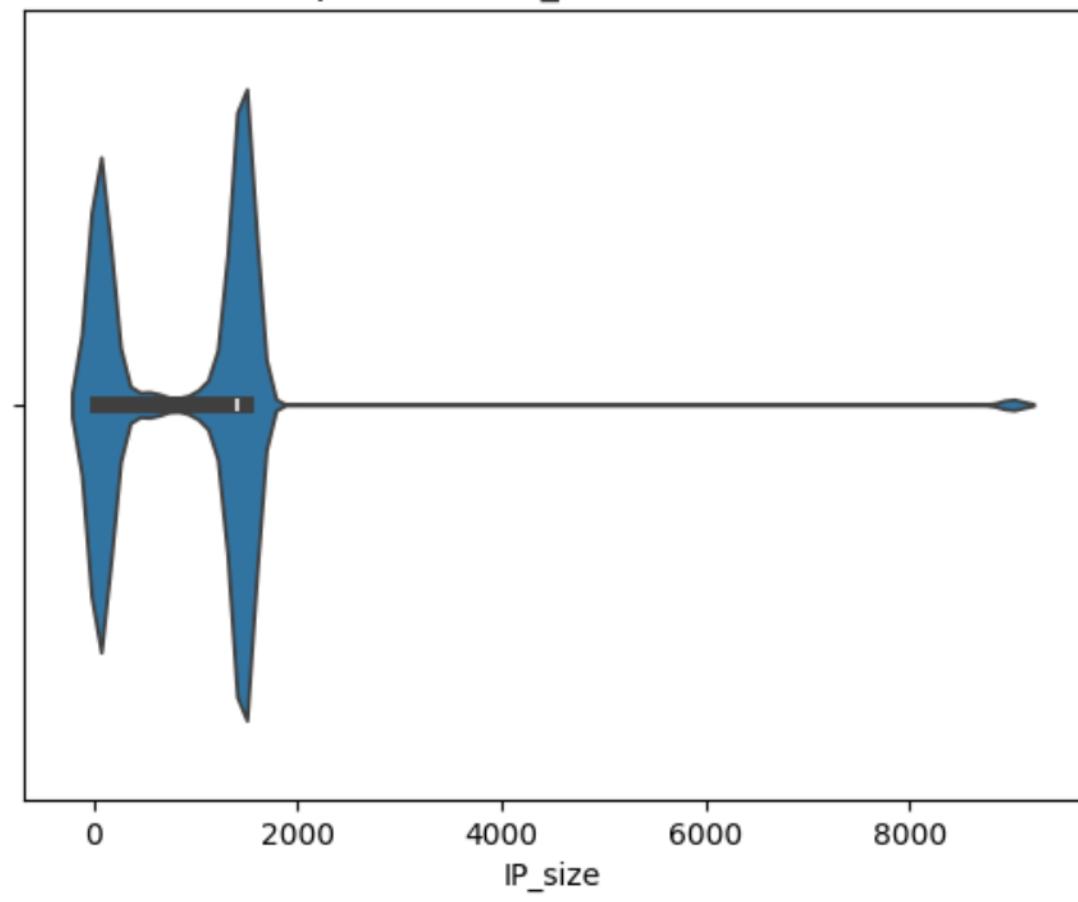


```
1 # Check out the size of the IP_size  
2  
3 df.IP_size.describe()
```

```
count    69370.000000  
mean      933.801672  
std       983.516558  
min       0.000000  
25%      52.000000  
50%     1400.000000  
75%     1500.000000  
max      9000.000000  
Name: IP_size, dtype: float64
```

```
● ● ●  
1 # Import libraries for data visualization  
2 # Create violin plot to visualize IP packet size distribution  
3 # Customize plot title  
4 # Display the plot  
5  
6 import seaborn as sns  
7 import matplotlib.pyplot as plt  
8  
9 sns.violinplot(data = df.IP_size, orient='h')  
10 plt.title("Violin plot of the IP_size column values")  
11 plt.show()
```

Violin plot of the IP_size column values



```

● ● ●

1 # Identify top 5 source-destination IP pairs with most packets
2 # Analyze and print results for the top 5 traffic flows
3 # Attempt to lookup organization for source IP using external API
4 # Attempt to lookup organization for destination IP using external API
5 # Print results in a formatted table-like manner
6
7
8 top_freq = df.groupby(["src_IP", "dst_IP"])\n9         .size()\n10        .sort_values(ascending=False)\n11        .nlargest(5)
12
13 for src_ip, dst_ip in top_freq.index:
14     count = top_freq[(src_ip, dst_ip)]
15
16     res_src = requests.get(f"http://ip-api.com/json/{src_ip}")
17     org_src = res_src.json()["org"] if res_src.status_code == 200 else "No response from server"
18
19     res_dst = requests.get(f"http://ip-api.com/json/{dst_ip}")
20     org_dst = res_dst.json()["org"] if res_dst.status_code == 200 else "No response from server"
21
22     print(f"Source: {src_ip}:15 ({org_src}):35\tDestination: {dst_ip}:15 ({org_dst}):35\tPacket count: {count}")

```

```

Source: 193.62.192.8 (European Bioinformatics Institute) Destination: 137.132.228.15 (National University of Singapore)
Packet count: 3041
Source: 130.14.250.11 (National Library of Medicine) Destination: 103.37.198.100 (A*STAR)
Packet count: 2599
Source: 14.139.196.58 (Indian Institute of Technology) Destination: 192.101.107.153 (Battelle Memorial Institute,
Pacific Northwest Division) Packet count: 2368
Source: 140.112.8.139 () Destination: 103.21.126.2 (Indian Institute of Technology
Bombay) Packet count: 2056
Source: 137.132.228.15 (National University of Singapore) Destination: 193.62.192.8 (European Bioinformatics Institute)
Packet count: 1910

```

```

● ● ●

1 # We can see these pairs were the top 5 pairs that have the highest amount of communication (packets)
2 # Now to focus on the top 50 communication pairs and create a visualization of communications between different IP hosts
3
4 communication_pairs

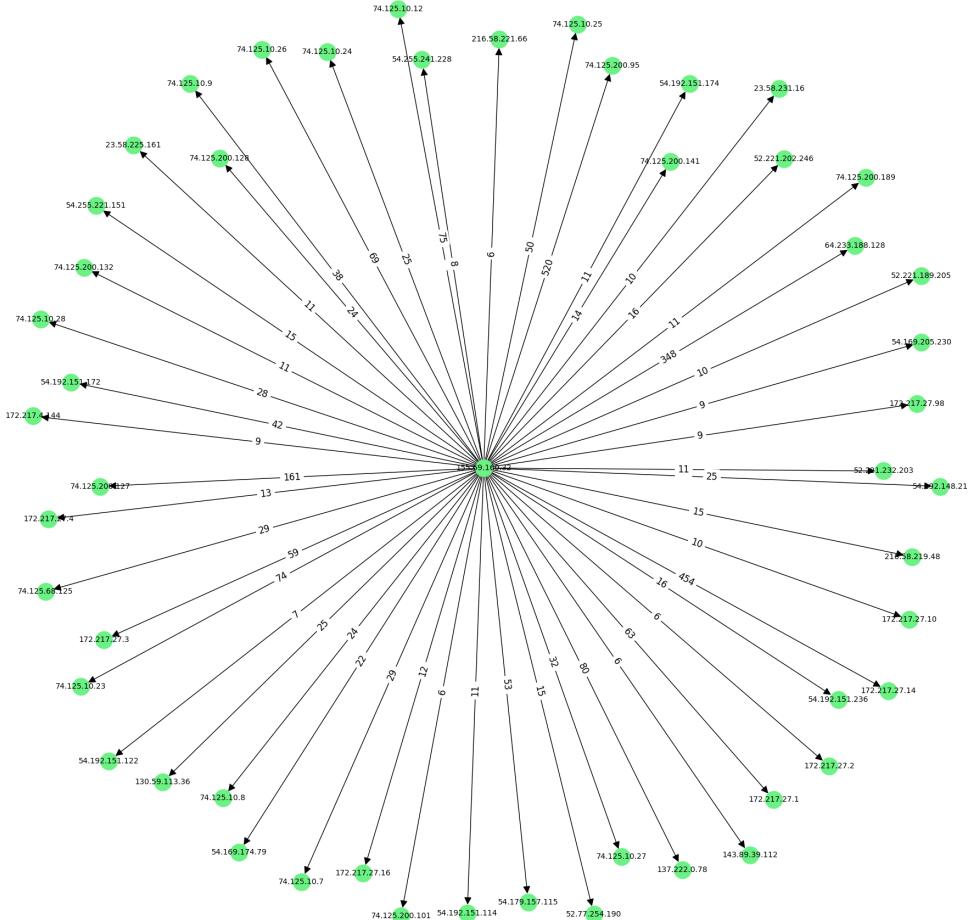
```

	src_IP	dst_IP	count
0	-	-	1
1	0		90
2	10.1.26.28	192.203.230.10	1
3	10.11.14.202	160.36.2.94	1
4	10.11.7.211	198.202.121.142	1
...
6584	95.108.213.253	123.136.68.137	1
6585	95.57.73.150	120.124.170.136	1
6586	95.80.49.94	163.22.18.74	1
6587	fe80:0000:0000:0000:72e4:22ff:fe69:fb68	fe80:0000:0000:0000:224e:710b:c5cf:1b0f	1
6588	fe80:0000:0000:0000:b2a8:6e03:ca76:6716	fe80:0000:0000:0000:ce4e:24ff:fe9c:6314	1

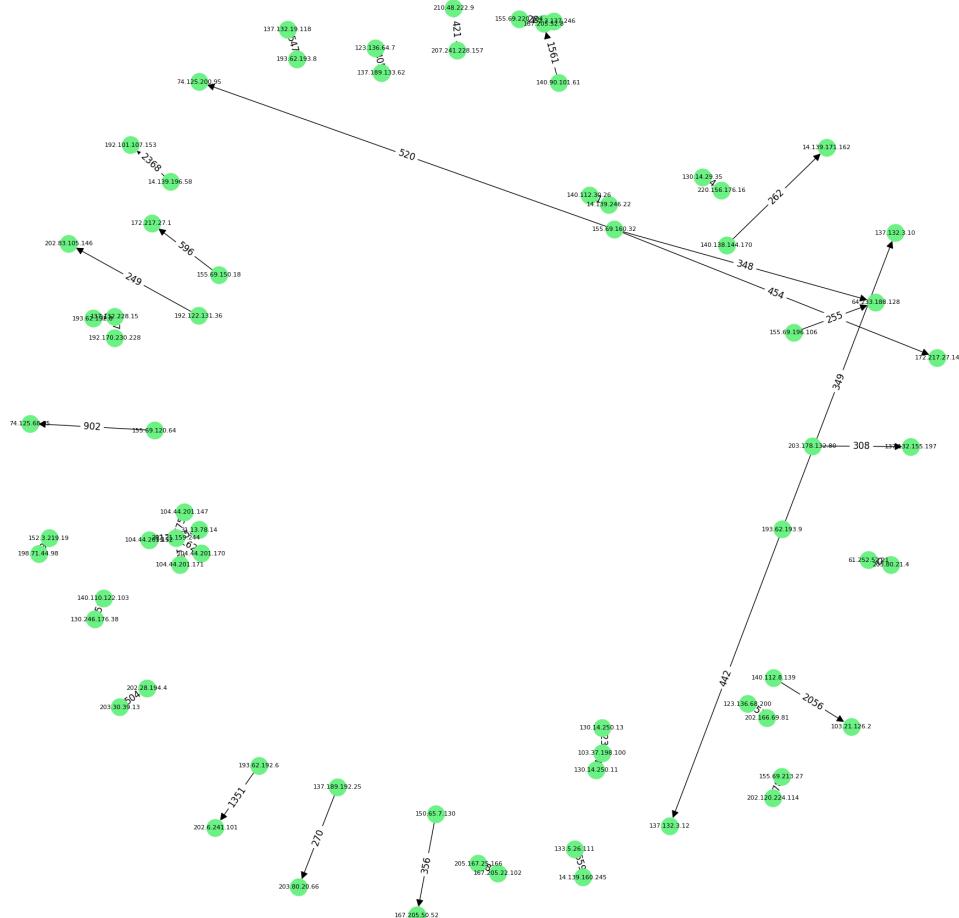
```

1 # This code investigates communication patterns for a specific source IP address ('155.69.160.32') by identifying its top 50 communication partners (destination IPs) based on packet count.
2 # It then constructs a directed graph where nodes represent IP addresses and directed edges depict communication flow.
3 # The graph visualization highlights frequently contacted destinations and the relative volume of traffic exchanged with each, providing insights into network traffic patterns.
4
5 import networkx as nx
6
7 communication_pairs = df[df['src_IP'] == '155.69.160.32'].groupby(['src_IP', 'dst_IP']).size().reset_index(name='count')
8 top_50_pairs = communication_pairs.sort_values('count', ascending=False).head(50)
9
10 # Create a new column in the top_50_pairs DataFrame
11 top_50_pairs['src_IP_dst_IP'] = list(zip(top_50_pairs['src_IP'], top_50_pairs['dst_IP']))
12 top_pairs = top_50_pairs[['src_IP_dst_IP']].to_list()
13
14 # Filter the data to include only the top 50 communication pairs
15 filtered_df = df[df.apply(lambda row: (row['src_IP'], row['dst_IP']) in top_pairs, axis=1)]
16
17 # Create a directed graph from the filtered dataframe
18 G_filtered = nx.from_pandas_edgelist(filtered_df, "src_IP", "dst_IP", create_using=nx.DiGraph())
19
20 # Draw the filtered graph with node labels
21 pos_filtered = nx.spring_layout(G_filtered, seed=42)
22 plt.figure(figsize=(20, 20))
23 nx.draw(G_filtered, pos_filtered, with_labels=True, node_size=500, node_color="lightgreen", arrowsize=20, font_size=10)
24
25 # Add edge labels with packet count
26 edge_labels_filtered = {(src, dst): filtered_df[(filtered_df['src_IP'] == src) & (filtered_df['dst_IP'] == dst)].shape[0] for src, dst in G_filtered.edges()}
27 nx.draw_networkx_edge_labels(G_filtered, pos_filtered, edge_labels=edge_labels_filtered, font_size=12)
28
29 # plt.title("Top 50 Communication Pairs")
30 plt.show()

```



Top 50 Communication Pairs



```

1 # This code defines a function get_ip_info that retrieves information about an IP address using the IPWhois service.
2 # It first creates an IPWhois object and attempts an RDAP lookup to get details. If successful, it returns a dictionary with the retrieved information.
3 # The function handles errors like invalid IP addresses and other exceptions during the lookup process, printing informative messages and returning None to indicate failure.
4
5
6 # code for getting ip information using IPWhois
7 def get_ip_info(ip):
8     try:
9         ipwhois = IPWhois(ip)
10        result = ipwhois.lookup_rdap(depth=1)
11        return result
12    except IPDefinedError:
13        return None
14    # except http error as well
15    except Exception as e:
16        print(e)
17        return None

```

```

1 # This code visualizes a network communication graph using Plotly, enhancing it with interactive features and IP details.
2 # It constructs edge and node traces based on a filtered graph, maps coordinates for proper layout, and incorporates IP information for hover text and node coloring.
3 # The plot visually maps communication patterns with edge connections, interactively reveals IP details and connections upon hovering, and highlights frequently connected nodes through color-coding and sizing.
4
5
6 import plotly.graph_objs as go
7
8 # Initialize edge trace
9 edge_trace = go.Scatter(
10    x=[],
11    y=[],
12    Line=dict(width=1, color="gray"),
13    hoverinfo="none",
14    mode="lines",
15 )
16
17 # Add edge coordinates to the edge trace
18 for src, dst in G.filtered.edges():
19    pos_filtered[src]
20    x1, y1 = pos_filtered[src]
21    edge_trace["x"] += (x0, x1, None)
22    edge_trace["y"] += (y0, y1, None)
23
24 # Initialize node trace
25 node_trace = go.Scatter(
26    x=[],
27    y=[],
28    text=[],
29    mode="markers+text",
30    hoverinfo="text",
31    hovertext=[],
32    marker=dict(
33        symbol="circle",
34        showscale=True,
35        colorscale="Viridis",
36        reversescale=True,
37        color=[],
38        size=[],
39        colorbar=dict(thickness=15, title="Node Connections", xanchor="left", titleside="right"),
40        line=dict(width=2, color="black"),
41    ),
42    textposition="bottom center",
43 )
44
45 # Add node coordinates, labels, and hover text to the node trace
46 for node in G.filtered.nodes():
47    x, y = pos_filtered[node]
48    num_connections = G.filtered.degree[node]
49    node_trace["x"] += tuple([x])
50    node_trace["y"] += tuple([y])
51    node_trace["text"] += tuple([""])
52    ip_info = get_ip_info(node)
53    if ip_info:
54        Country = ip_info["network"]["country"]
55        org = ip_info["network"]["name"]
56        hover_text = f"IP: {node}<br>Connections: {num_connections}<br>Country: {Country}<br>Organization: {org}"
57    else:
58        hover_text = f"IP: {node}<br>Connections: {num_connections}"
59    node_trace["marker"]["color"] += tuple([num_connections])
60    node_trace["marker"]["size"] += tuple([num_connections * 10])
61
62 # Create a Plotly figure with the edge and node traces
63 fig = go.Figure(
64    data=[edge_trace, node_trace],
65    layout=Layout(
66        title="Communication Pairs",
67        showlegend=False,
68        hovermode="closest",
69        margin=Margin(b=20, l=5, r=5, t=40),
70        xaxis=dict(showgrid=False, zeroline=False, showticklabels=False),
71        yaxis=dict(showgrid=False, zeroline=False, showticklabels=False),
72        plot_bgcolor="white",
73    ),
74 )
75
76 # Show the figure
77 fig.show()

```

Communication Pairs

