

Student Name : Toh Kok Soon

Group : SCSI

Date : 15/10/2025

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	13.107.4.50	5960	MSFT
2	130.14.250.7	4034	NLM-ETHER
3	155.69.160.38	3866	NTUNET1
4	171.67.77.19	2656	NETBLK-SUNET
5	155.69.199.255	2587	NTUNET1

TOP 5 LISTENERS

Rank	IP address	# of packets	Organisation
1	137.132.228.33	5908	NUSNET
2	192.122.131.36	4662	A-STAR-AS-AP
3	202.51.247.133	4288	NUSGP
4	137.132.228.29	4022	NUSNET
5	103.37.198.100	3741	A-STAR-AS-AP

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	137707
2	17	UDP	36852
3			

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	43208	HTTPS
2	80	11018	HTTP
3	50930	2450	Dynamic Ports
4	15000	2103	Dynamic Ports
5	8160	1354	Dynamic Ports

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)	20.258 Mb
--------------------	-----------

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).

EXERCISE 4F: SOFTWARE CODE

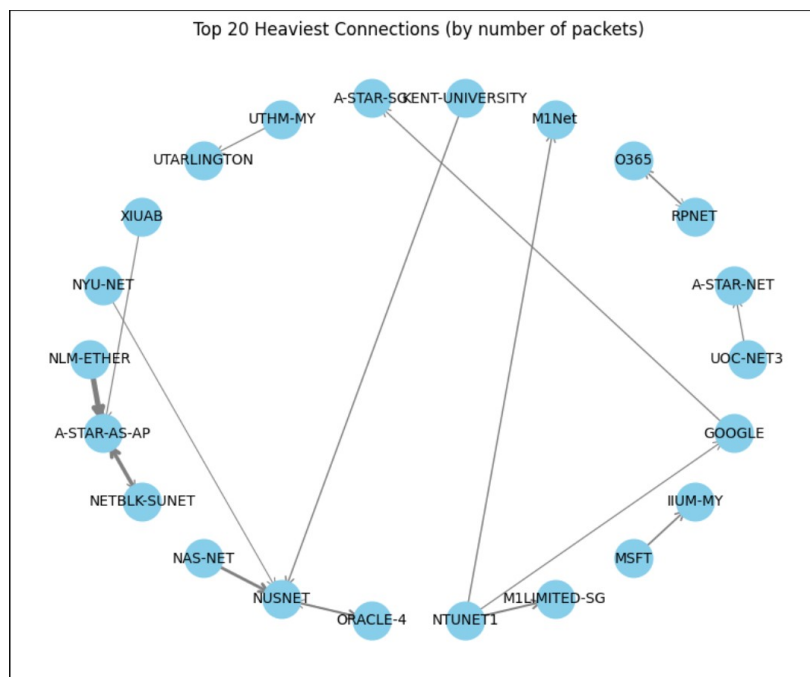
Please also submit your code to the NTU Learn lab site.

Additional Analysis

Top 5 Connection Pairs

	src_IP	dst_IP	Number of Packets	From	To
0	130.14.250.7	103.37.198.100	3739	NLM-ETHER	A-STAR-AS-AP
1	171.67.77.19	192.122.131.36	2656	NETBLK-SUNET	A-STAR-AS-AP
2	129.99.230.54	137.132.22.74	2097	NAS-NET	NUSNET
3	137.132.228.42	137.131.17.212	1553	NUSNET	ORACLE-4
4	155.69.252.133	138.75.242.36	1475	NTUNET1	M1LIMITED-SG

Visualisation of Top 20 Heaviest Connections (by number of Packets)



Observations from graph

1. The graph shows most connection are mostly unidirectional
2. There's no obvious central hub
3. The graph is very sparse with some of the connection being only between the pair
4. A-STAR-AS-AP is the most heavy in traffic by number of packets which can be seen by the thickness of the connection
5. Most of these nodes seem to belong to universities, software services as well as telcoms

Insights

The observed network is likely to be a network of research organisation since it consists mostly of universities and research organisations. The many isolated pairs and dominance of unidirectional traffic suggest the observed network connections are mostly peer to peer. Furthermore, no obvious central hub could be seen from the graph, therefore reinforcing the idea that this is a mostly peer to peer network. In a client-server network, we would expect heavy traffic on one of the nodes in terms of both number of packets as well as nuber of connections.

Learnings

1. We can gain deeper insights into network through visualisation like network graphs
2. Whois lookup is slow and is a bottleneck when it comes to the analysis of the network log, most operation can be done straight from the IP address so we should only do lookups when necessary such as in the final visualisations, this can be further mitigated through caching in a dictionary map to avoid re-lookups.