

Internet Measurements

CS425- Computer Networks

Vaibhav Nagar (14785)
Email: vaibhavn@iitk.ac.in

November 14, 2016

1. Traffic Measurement

1.1 What is the average packet size, across all traffic in the trace? Describe how you computed this number.

Average Packet Size = Total sum of all bytes(doctets) / Total sum of number of packets(dpks) in each flow

$$\text{Average Packet Size} = 768.18$$

1.2 Plot the Complementary Cumulative Probability Distribution (CCDF) of flow durations (i.e., the finish time minus the start time) and of flow sizes (i.e., number of bytes, and number of packets). First plot each graph with a linear scale on each axis, and then a second time with a logarithmic scale on each axis. What are the main features of the graphs? What artifacts of Netflow and of network protocols could be responsible for these features? Why is it useful to plot on a logarithmic scale?

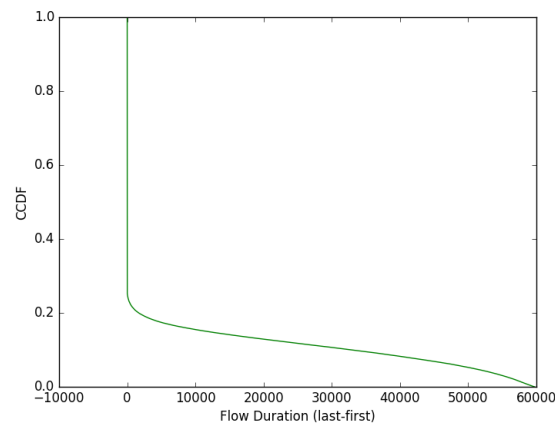


Figure 1: Complementary cumulative distribution of flow durations on linear scale

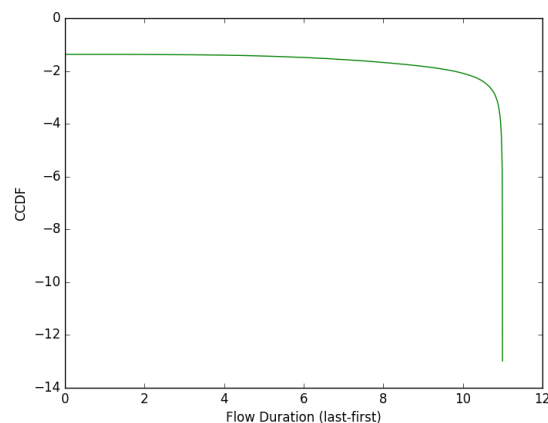


Figure 2: Complementary cumulative distribution of flow durations on logarithmic scale

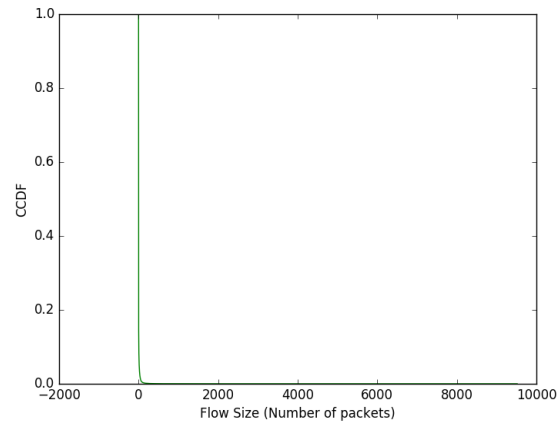


Figure 3: Complementary cumulative distribution of flow size(number of packets) on linear scale

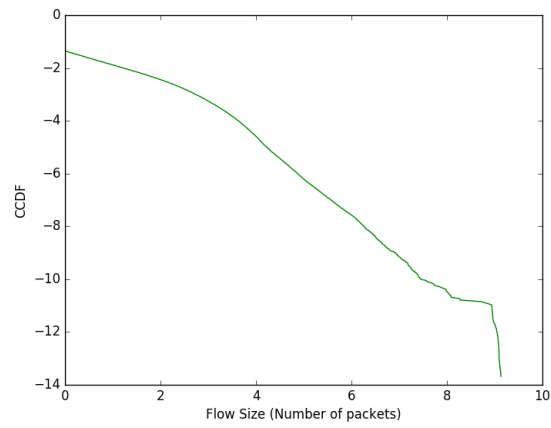


Figure 4: Complementary cumulative distribution of flow size(number of packets) on logarithmic scale

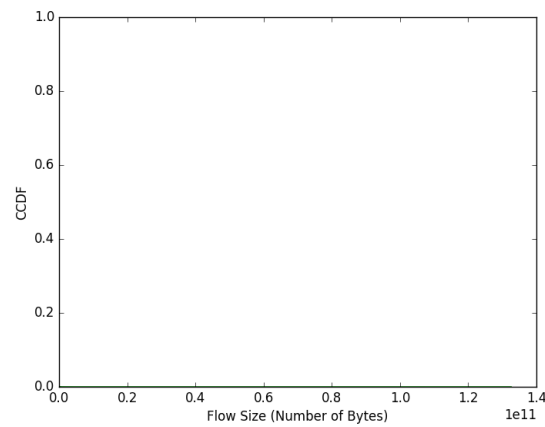


Figure 5: Complementary cumulative distribution of flow size(number of bytes) on linear scale

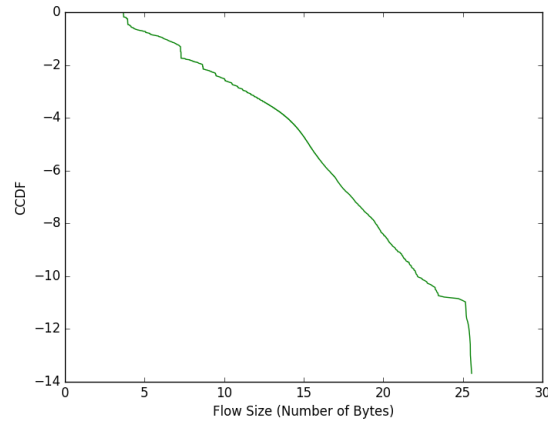


Figure 6: Complementary cumulative distribution of flow size(number of bytes) on logarithmic scale

- Graph of CCDF of flow durations shows that there are huge number of flows whose duration is zero and then after $x = 0$ it gradually decreases. The logarithmic graph of ccdf of flow duration shows clearly that after e^{10} duration there is significant fall.
- CCDF graph of flow size (both number of packets and number of bytes) doesn't reveal much as linear scale is too large and large number of packets are concentrated at smaller range. So the graph rapidly falls to 0 just after $x = 0$.
- Their logarithmic graphs clearly shows the packet density falls after certain point. The logarithmic graph of ccdf of flow size in bytes start with $y = 0$ upto almost $x = 4$ which is due to minimum acceptable size of packet is at least 8 bytes (UDP) and 20 bytes(TCP) and thus drastically decreases after $x = 25$ due to maximum acceptable size of packets.
- Plotting in logarithmic scale shrinks the x-axis. In linear scale when some values are too large and most of the data are concentrated in small region then the graph doesn't render much info. Logarithmic graphs perfectly responds to skewness towards large values and shows percent change or multiplicative factors.

1.3 Summarize the traffic by which TCP/UDP port numbers are used. Create two tables, listing the top-ten port numbers by sender traffic volume (i.e., by source port number) and by receiver traffic volume (i.e., by destination port number), including the percentage of traffic (by bytes) they contribute. Where possible, explain what applications are likely responsible for this traffic. Explain any significant differences between the results for sender vs. receiver port numbers

Sender traffic at port 80 is maximum as expected. Since port 80 is well-known and registered port for HTTP protocol and thus used by global servers. Port 22 and 443 also contributes for high traffic value as port 22 is used for SSH and 443 for HTTPS globally.

Port Numbers	Percentage of traffic (in bytes)
80	43.94
33001	7.36
1935	3.66
22	2.17
443	1.72
55000	1.62
388	1.34
16402	0.76
20	0.671
873	0.58

Table 1: Top-ten port numbers by *sender* traffic volume

Port Numbers	Percentage of traffic (in bytes)
33002	4.02
80	2.98
49385	2.09
62269	1.23
443	0.77
43132	0.76
16402	0.74
22	0.65
5500	0.64
57493	0.35

Table 2: Top-ten port numbers by *receiver* traffic volume

Servers send large number of bytes to receivers whose port may not necessarily be the well-known or registered port, when requested by clients. Thus percentage of receiver traffic is not as much high as senders.

1.4 Aggregate the traffic volumes based on the source IP prefix. What fraction of the total traffic comes from the most popular (by number of bytes) 0.1% of source IP prefixes? The most popular 1% of source IP prefixes? The most popular 10% of source IP prefixes? Some flows will have a source mask length of 0. Report the fraction of traffic (by bytes) that has a source mask of 0, and then exclude this traffic from the rest of the analysis. That is, report the top 0.1%, 1%, and 10% of source prefixes that have positive mask lengths.

1. Fraction of the total traffic comes from the most popular (by number of bytes):

- 0.1% of source IP prefixes = 0.58944
- 1% of source IP prefixes = 0.82253
- 10% of source IP prefixes = 0.98383

2. The fraction of traffic (by bytes) that has a source mask of 0 = 0.43259

3. Now after excluding the traffic that has source mask 0, fraction of the total traffic comes from the most popular (by number of bytes):

- 0.1% of source IP prefixes = 0.27005
- 1% of source IP prefixes = 0.45622
- 10% of source IP prefixes = 0.56383

1.5 Princeton has the 128.112.0.0/16 address block. What fraction of the traffic (by bytes and by packets) in the trace is sent by Princeton? To Princeton?

1. Fraction of the traffic in the trace is sent by Princeton:

- In terms of Bytes = 0.00701368
- In terms of Packets = 0.02191742

2. Fraction of the traffic in the trace is sent to Princeton:

- In terms of Bytes = 0.010148
- In terms of Packets = 0.014689

2. BGP Measurement

2.1 How many BGP updates per minute does the session handle, on average? Count each announcement or withdrawal for a single prefix as a BGP update, even if multiple prefixes appear in a single update message or a single prefix has multiple announcement or withdrawal messages in a short period of time. Describe how you computed this result.

Average BGP updates per minute is calculated by counting the number of BGP update messages in the text file for all intervals of a session. Update messages for IPv6 are excluded and all prefixes are counted in a single update message. Then updates per minute is averaged by dividing total count by number of minutes in a session.

- **Session 1:** RouteViews on January 3, 2014 between 12pm and 2pm from Sydney, Australia

$$\text{Average BGP updates per minute} = 623.66$$

- **Session 2:** RouteViews on February 3, 2014 between 12pm and 2pm from Sydney, Australia

$$\text{Average BGP updates per minute} = 1020.21$$

- **Session 3:** RouteViews on March 3, 2014 between 12pm and 2pm from Sydney, Australia

$$\text{Average BGP updates per minute} = 1959.05$$

2.2 What fraction of IP prefixes experience no update messages? (Count each prefix equally, independently of what fraction of address space they cover or whether one prefix is contained inside another.)

- **Session 1:** RouteViews on January 3, 2014 between 12pm and 2pm from Sydney, Australia
 - Fraction of IP prefixes experience no update messages: 0.95736
- **Session 2:** RouteViews on February 3, 2014 between 12pm and 2pm from Sydney, Australia
 - Fraction of IP prefixes experience no update messages: 0.93444
- **Session 3:** RouteViews on March 3, 2014 between 12pm and 2pm from Sydney, Australia
 - Fraction of IP prefixes experience no update messages: 0.88125

2.3 What prefix (or prefixes) experiences the most updates, and how frequent are they?

Prefixes that experiences the most updates are:

- **Session 1:** RouteViews on January 3, 2014 between 12pm and 2pm from Sydney, Australia
 - 121.52.149.0/24 = 1020 updates in two hour session
 - 121.52.144.0/24 = 1020 updates in two hour session
 - 121.52.145.0/24 = 1020 updates in two hour session
 - 121.52.150.0/24 = 1020 updates in two hour session
- **Session 2:** RouteViews on February 3, 2014 between 12pm and 2pm from Sydney, Australia
 - 89.221.206.0/24 = 712 updates in two hour session
- **Session 3:** RouteViews on March 3, 2014 between 12pm and 2pm from Sydney, Australia
 - 109.161.64.0/20 = 728 updates in two hour session

2.4 What fraction of all update messages come from the most unstable 0.1% of prefixes? The most unstable 1% of prefixes? The most unstable 10% of prefixes?

Fraction of all update messages come from the most unstable:

Session	0.1% of prefixes	1% of prefixes	10% of prefixes
Session-1	0.3876	0.8443	0.8504
Session-2	0.2791	0.7026	0.8800
Session-3	0.0718	0.2278	0.5727

Table 3: Fraction of all update messages come from unstable prefixes

2.5 Briefly summarize your results and what you learned about BGP stability from them.

After looking at the fraction of prefixes which is close to 1, this indicates that very less number of prefixes announce or withdraw. Still the average updates per minute is quite high in all three sessions which indicates BGP instability. Thus due to large number of updates per minute BGP packets consumes significant amount of bandwidth and CPU resources on routers and disrupt the delivery of data traffic.

3. Python Scripts (Python-3.5.2)

3.1 Plot CCDF of traffic flow for Q1.1 and Q1.2

```
1 import pandas as pd
2 import sys
3 import numpy as np
4 import matplotlib.pyplot as plt
5
6
7 def plot_ccdf(dataset, xstring, ystring, img, fig_num, log_scale):
8     dataset.sort()
9     xvals = np.array(list(set(dataset)))
10    xvals.sort()
11
12    counts = []
13    c = 0
14    d_prev = dataset[0]
15
16    for d in dataset:
17        if d_prev != d:
18            d_prev = d
19            counts = counts + [c]
20            c = 1
21        else:
22            c += 1
23    counts = counts + [c]
24
25    assert len(counts) == len(xvals)
26
27    cdf = np.cumsum(counts)
28
29    assert len(dataset) == cdf[len(cdf)-1]
30
31    yvals = 1 - cdf/float(len(dataset))
32
33    if log_scale:
34        xvals = np.log(xvals)
35        yvals = np.log(yvals)
36    else:
37        xvals = np.append([xvals[0]-1], xvals)
38        yvals = np.append([1], yvals)
39        xvals = np.append([xvals[0]-1], xvals)
40        yvals = np.append([1], yvals)
41
42
43    assert len(xvals) == len(yvals)
44
45    plt.figure(fig_num)
46    plt.plot(xvals, yvals, color="green")
47    plt.xlabel(xstring)
48    plt.ylabel(ystring)
```



```

49 # plt.show()
50 plt.savefig(img)
51
52
53 try:
54     data = pd.read_csv("ft-v05.2010-09-29.235501+0000.csv")
55 except:
56     print("Unable to read csv file")
57     sys.exit(1)
58
59 # extract number of packets in the flow <dpkts>
60 dpkts = data.dpkts
61
62 # extract bytes in the flow <doctets>
63 doctets = data.doctets
64
65 assert len(dpkts) == len(doctets)
66
67 packet_list = [int(pk) for pk in dpkts]
68 octets_list = [int(oc) for oc in doctets]
69
70 avg = sum(octets_list)/sum(packet_list)
71
72 print("Total Size: ", sum(octets_list))
73 print("Number of Samples: ", len(doctets))
74 print("Average Packet Size: ", avg)
75
76 finish = data['last']
77 start = data['first']
78
79 assert len(finish) == len(start)
80
81 durations = [f-s for f,s in zip(finish, start)]
82
83 plot_ccdf(durations, "Flow Duration (last-first)", "CCDF", "traffic_q1.2.1.png", 1, False)
84 plot_ccdf(durations, "Flow Duration (last-first)", "CCDF", "traffic_q1.2.1_log.png", 2,
85           True)
86
87 plot_ccdf(packet_list, "Flow Size (Number of packets)", "CCDF", "traffic_q1.2.2.png", 3,
88           False)
89 plot_ccdf(packet_list, "Flow Size (Number of packets)", "CCDF", "traffic_q1.2.2_log.png",
90           4, True)
91
92 plot_ccdf(octets_list, "Flow Size (Number of Bytes)", "CCDF", "traffic_q1.2.3.png", 5,
93           False)
94 plot_ccdf(octets_list, "Flow Size (Number of Bytes)", "CCDF", "traffic_q1.2.3_log.png", 6,
95           True)

```

3.2 Fractions for Q1.3 and Q1.4 and Q1.5

```

1 import pandas as pd
2 import sys
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import operator
6
7 def top_ten(port_list, octet_list, proto_list):
8     _traffic = {}
9
10    for port in list(set(port_list)):
11        _traffic[int(port)] = 0
12
13    print(len(_traffic))
14    total_volume = 0

```

```

15 for port, vol, proto in zip(port_list, octet_list, proto_list):
16     if (int(proto) == 6) or (int(proto) == 17):
17         _traffic[int(port)] += int(vol)
18     total_volume += int(vol)
19
20 for key in _traffic:
21     _traffic[key] = (_traffic[key]/total_volume)*100
22
23 return sorted(_traffic.items(), key=operator.itemgetter(1), reverse=True)[:10]
24
25
26 def ip_traffic_fraction(ipaddr, octet_list, mask_list, masked, fractions):
27     ip_traffic = {}
28     for ip in list(set(ipaddr)):
29         ip_traffic[ip] = 0
30
31     total_traffic = 0
32     for ip, octets, mask in zip(ipaddr, octet_list, mask_list):
33         if masked:
34             if mask != 0:
35                 ip_traffic[ip] += octets
36             else:
37                 ip_traffic[ip] += octets
38             total_traffic += octets
39
40     for f in fractions:
41         f = f/100
42         a = dict(sorted(ip_traffic.items(), key=operator.itemgetter(1), reverse=True)[:int(len
43             (ip_traffic)*f)])
44         print("Fraction of the total traffic comes from the most popular ", f*100," percent of
45             source IP prefixes: ", sum(a.values())/total_traffic)
46     return
47
48 def ip_traffic_zero_mask(octet_list, mask_list):
49     res = 0
50     tot = 0
51     for oc, mask in zip(octet_list, mask_list):
52         if mask == 0:
53             res += oc
54             tot += oc
55     print(tot)
56     print(res)
57     print("Fraction of traffic (by bytes) that has a source mask of 0", res/tot)
58
59 def specific_ip_traffic(ipaddr, _list, ip_string, addr_block):
60     total_traffic = 0
61     res = 0
62     for ip, octets in zip(ipaddr, _list):
63         if ip_string in ip[:addr_block]:
64             res += octets
65             total_traffic += octets
66     print(total_traffic)
67     print(res)
68
69     return res/total_traffic
70
71
72 try:
73     data = pd.read_csv("ft-v05.2010-09-29.235501+0000.csv")
74 except:
75     print("Unable to read csv file")

```

```

76     sys.exit(1)
77
78 sender = top_ten(data.srcport, data.doctets, data.prot)
79 reciever = top_ten(data.dstport, data.doctets, data.prot)
80 print("Top Ten Senders: ")
81 print(sender)
82 print("Top Ten Recievers: ")
83 print(reciever)
84
85 ip_traffic_fraction(data.srcaddr, data.doctets, data.src_mask, False, [0.1, 1, 10])
86
87 ip_traffic_zero_mask(data.doctets, data.src_mask)
88
89 print("Non zero mask traffic:")
90
91 ip_traffic_fraction(data.srcaddr, data.doctets, data.src_mask, True, [0.1, 1, 10])
92
93 p1 = specific_ip_traffic(data.srcaddr, data.doctets, "128.112", 16)
94 p2 = specific_ip_traffic(data.srcaddr, data.dppts, "128.112", 16)
95 p3 = specific_ip_traffic(data.dstaddr, data.doctets, "128.112", 16)
96 p4 = specific_ip_traffic(data.dstaddr, data.dppts, "128.112", 16)
97
98 print("fraction of the traffic (by bytes) in the trace is sent by Princeton: ", p1)
99 print("fraction of the traffic (by packets) in the trace is sent by Princeton: ", p3)
100 print("fraction of the traffic (by bytes) in the trace is to by Princeton: ", p2)
101 print("fraction of the traffic (by packets) in the trace is to by Princeton: ", p4)

```

3.3 Average BGP updates for Q2.1

Change session value accordingly.

```

1 import pandas as pd
2 import sys
3 import numpy as np
4
5 session = "1"
6
7 intervals = [1200, 1215, 1230, 1245, 1300, 1315, 1330, 1345]
8
9 updates = []
10
11 for interval in intervals:
12     filename = "./updates-asst4/updates.20140" + session + "03." + str(interval) + ".txt"
13
14     try:
15         print(filename)
16         cols = ["A", "B", "C", "D", "E", "F", "G", "H", "I", "J", "K", "L", "M", "N", "O"]
17         data = pd.read_csv(filename, sep="|", header=None, names=cols)
18     except:
19         #raise
20         print("Unable to read file")
21         sys.exit(1)
22
23     count = 0
24     for aw, entry in zip(data["C"], data["F"]): # index "F" is the sixth column
25         if aw == "A" or aw == "W":
26             ent = entry.split() # multiple prefixes in one update entry
27             for e in ent:
28                 if ':' not in e: # not ipv6
29                     count += 1
30     updates += [count]
31
32 print(updates)
33 print(sum(updates))

```

```
34 print("Average BGP updates per minute: ", sum(updates)/120)
```

3.4 Fractions for Q2.2, Q2.3 and Q2.4

Change session value accordingly.

```
1 import pandas as pd
2 import sys
3 import numpy as np
4 import operator
5
6 try:
7     cols = ["A", "B", "C", "D", "E", "F", "G", "H", "I", "J", "K", "L", "M", "N", "O"]
8     rib = pd.read_csv("./rib/rib.20140103.1200.txt" , sep="|", header=None, names=cols)
9 except:
10     raise
11     print("Unable to read file")
12     sys.exit(1)
13
14 rib_data = list(set(rib["F"]))
15
16 session = "3"
17
18 intervals = [1200, 1215, 1230, 1245, 1300, 1315, 1330, 1345]
19
20 updates_data = []
21 update_messages = 0
22
23 for interval in intervals:
24     filename = "./updates-asst4/updates.20140" + session + "03." + str(interval) + ".txt"
25
26     try:
27         print(filename)
28         cols = ["A", "B", "C", "D", "E", "F", "G", "H", "I", "J", "K", "L", "M", "N", "O"]
29         data = pd.read_csv(filename , sep="|", header=None, names=cols)
30     except:
31         #raise
32         print("Unable to read file")
33         sys.exit(1)
34
35     count = 0
36     update_messages += len(data["F"])
37     for entry in data["F"]:
38         # index "F" is the sixth column
39         ent = entry.split()
40         # multiple prefixes in one update entry
41         for e in ent:
42             if ':' not in e:
43                 # not ipv6
44                 updates_data += [e]
45
46 print("This may take some time...")
47
48 ip_updates = {}
49 count = 0
50 for prefix in rib_data:
51     if ":" not in prefix:
52         ip_updates[prefix] = updates_data.count(prefix)
53         if ip_updates[prefix] == 0:
54             count += 1
55
56 print("count: ", count)
57 print("rib size: ", len(rib["F"]))
58 print("rib_data size(non-repititive): ", len(rib_data))
59 print("Total Update Messages: ", update_messages)
```

```

58 print("Fraction of IP prefixes experience no update messages: ", count/len(rib_data))
59
60 b = dict(sorted(ip_updates.items(), key=operator.itemgetter(1), reverse=True)[:10])
61 print(b)
62
63 fractions = [0.1, 1, 10]
64 for f in fractions:
65     f = f/100
66     a = dict(sorted(ip_updates.items(), key=operator.itemgetter(1), reverse=True)[:int(len(
        ip_updates)*f)])
67     print("Sum of update messages for fraction ", f*100, ": ", sum(a.values()))
68     print("Fraction of all update messages come from the most popular ", f*100, " percent of
        prefixes: ", sum(a.values())/update_messages)

```