

# Visualizing Cyber Attacks

```
In [1]: import numpy as np
import pandas as pd
import plotly.express as px
import seaborn as sns

import matplotlib.pyplot as plt
import plotly.graph_objects as go
import matplotlib.ticker as ticker
import plotly.figure_factory as ff

from fbprophet import Prophet
from pandas_datareader import data
from scipy import stats
from sklearn.metrics import mean_absolute_error
from statsmodels.tsa.seasonal import seasonal_decompose

from sklearn.preprocessing import MinMaxScaler

from matplotlib.ticker import FixedFormatter, FixedLocator

import warnings
warnings.filterwarnings('ignore')
```

```
/opt/conda/lib/python3.7/site-packages/geopandas/_compat.py:115: UserWarning: The Shapely GEOS version (3.9.1-CAPI-1.14.2) is incompatible with the GEOS version PyGEOS was compiled with (3.10.1-CAPI-1.16.0). Conversions between both will be slow.
shapely_geos_version, geos_capi_version_string
```

```
In [2]: df = pd.read_csv("../input/aws-honeypot-attack-data/AWS_Honeypot_marx-geo.csv")
```

```
In [3]: month_list = ['Jan.', 'Feb.', 'Mar.', 'Apr.', 'May', 'June', 'July', 'Aug.', 'Sept.', 'Oct.', 'Nov.', 'Dec.']
attack_date = {"year": [], "month": [], "day": [], "time": [], "hour": []}

for d_date in df['datetime']:
    d_date = d_date.split('/')
    mon = d_date[0]
    day = d_date[1]

    year_time = d_date[2].split(' ')
    year = year_time[0]
    time = year_time[1]

    attack_date["year"].append(year)
    attack_date["month"].append(month_list[int(mon) + 1])
    attack_date["day"].append(day)
    attack_date["time"].append(time)
    attack_date["hour"].append(f"{time.split(':')[0]}h")

df["year"] = attack_date["year"]
df["month"] = attack_date["month"]
df["day"] = attack_date["day"]
df["time"] = attack_date["time"]
df["hour"] = attack_date["hour"]
```

In [4]: `df.head(20)`

Out[4]:

	<b>datetime</b>	<b>host</b>	<b>src</b>	<b>proto</b>	<b>type</b>	<b>spt</b>	<b>dpt</b>	<b>srcstr</b>	<b>cc</b>	<b>country</b>	...
<b>0</b>	3/3/13 21:53	groucho-oregon	1032051418	TCP	NaN	6000.0	1433.0	61.131.218.218	CN	China	...
<b>1</b>	3/3/13 21:57	groucho-oregon	1347834426	UDP	NaN	5270.0	5060.0	80.86.82.58	DE	Germany	...
<b>2</b>	3/3/13 21:58	groucho-oregon	2947856490	TCP	NaN	2489.0	1080.0	175.180.184.106	TW	Taiwan	...
<b>3</b>	3/3/13 21:58	groucho-us-east	841842716	UDP	NaN	43235.0	1900.0	50.45.128.28	US	United States	...
<b>4</b>	3/3/13 21:58	groucho-singapore	3587648279	TCP	NaN	56577.0	80.0	213.215.43.23	FR	France	...
<b>5</b>	3/3/13 21:58	groucho-tokyo	3323217250	TCP	NaN	32628.0	2323.0	198.20.69.98	US	United States	...
<b>6</b>	3/3/13 21:59	groucho-oregon	3730416887	TCP	NaN	6000.0	1433.0	222.89.164.247	CN	China	...
<b>7</b>	3/3/13 22:07	groucho-singapore	3738622573	TCP	NaN	6000.0	3306.0	222.214.218.109	CN	China	...
<b>8</b>	3/3/13 22:12	groucho-oregon	3683919430	TCP	NaN	6000.0	1433.0	219.148.38.70	CN	China	...
<b>9</b>	3/3/13 22:14	groucho-singapore	1007884304	TCP	NaN	6000.0	1433.0	60.19.24.16	CN	China	...
<b>10</b>	3/3/13 22:14	groucho-tokyo	3639889826	TCP	NaN	6000.0	1433.0	216.244.79.162	US	United States	...
<b>11</b>	3/3/13 22:20	groucho-oregon	1965603898	TCP	NaN	9907.0	1433.0	117.40.188.58	CN	China	...
<b>12</b>	3/3/13 22:20	groucho-sa	3672981807	TCP	NaN	33367.0	22.0	218.237.65.47	KR	South Korea	...
<b>13</b>	3/3/13 22:20	zeppo-norcal	3672981807	TCP	NaN	33367.0	22.0	218.237.65.47	KR	South Korea	...
<b>14</b>	3/3/13 22:19	groucho-norcal	3672981807	TCP	NaN	33367.0	22.0	218.237.65.47	KR	South Korea	...
<b>15</b>	3/3/13 22:20	groucho-us-east	3672981807	TCP	NaN	33367.0	22.0	218.237.65.47	KR	South Korea	...
<b>16</b>	3/3/13 22:21	groucho-eu	3672981807	TCP	NaN	33367.0	22.0	218.237.65.47	KR	South Korea	...
<b>17</b>	3/3/13 22:21	groucho-sydney	3672981807	TCP	NaN	33367.0	22.0	218.237.65.47	KR	South Korea	...
<b>18</b>	3/3/13 22:23	groucho-oregon	2382398543	TCP	NaN	6000.0	1433.0	142.0.132.79	US	United States	...
<b>19</b>	3/3/13 22:26	groucho-tokyo	1965603898	TCP	NaN	26180.0	1433.0	117.40.188.58	CN	China	...

20 rows × 21 columns

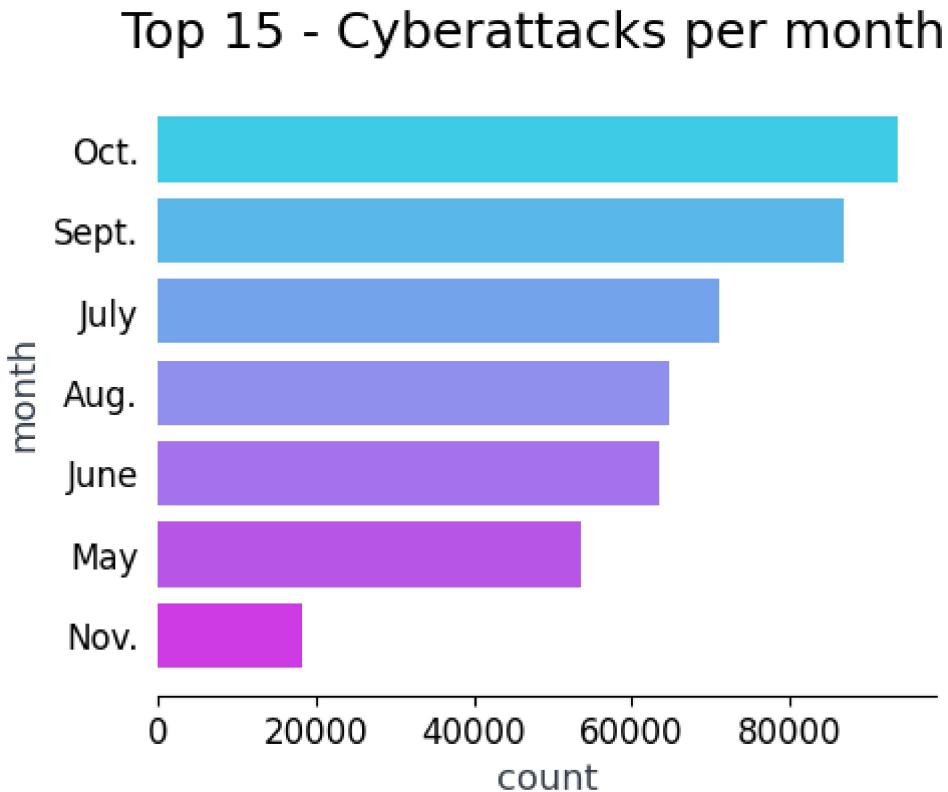
```
In [5]: def bar_plot_data(df: dict, field: str, subtitle: str, figsize=(5, 4), top_filter=15):
    fig, ax1 = plt.subplots(figsize=figsize, dpi=100)

    for spline in ['top', 'right', 'left']:
        ax1.spines[spline].set_visible(False)

    df_filter = df[field].value_counts().rename_axis(field).reset_index(name='counts')
    if top_filter:
        df_filter = df_filter.head(top_filter)
    sns.barplot(data=df_filter, palette='cool', x='counts', y=field)
    ax1.tick_params(axis='both', which='both', labelsize=12, bottom=True, left=False)
    ax1.set_xlabel('count', fontsize=13, color = '#333F4B')
    ax1.set_ylabel(f'{field}', fontsize=13, color = '#333F4B')

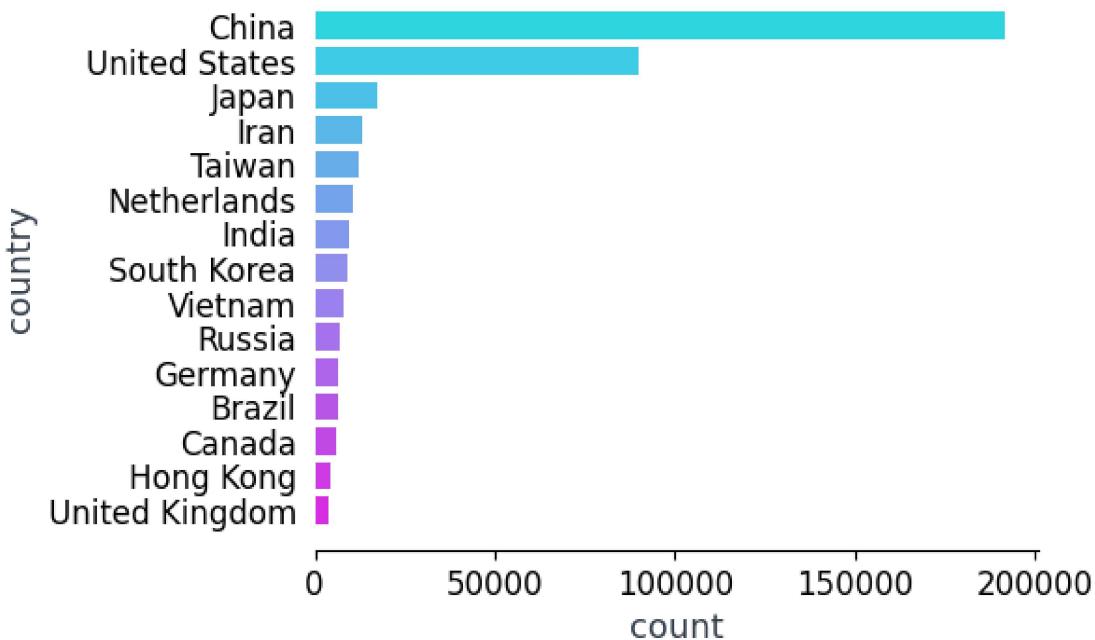
    plt.plot(color="white", lw=3)
    fig.suptitle(subtitle, fontsize=18)
    plt.show()
```

```
In [6]: bar_plot_data(df, 'month', 'Top 15 - Cyberattacks per month')
```

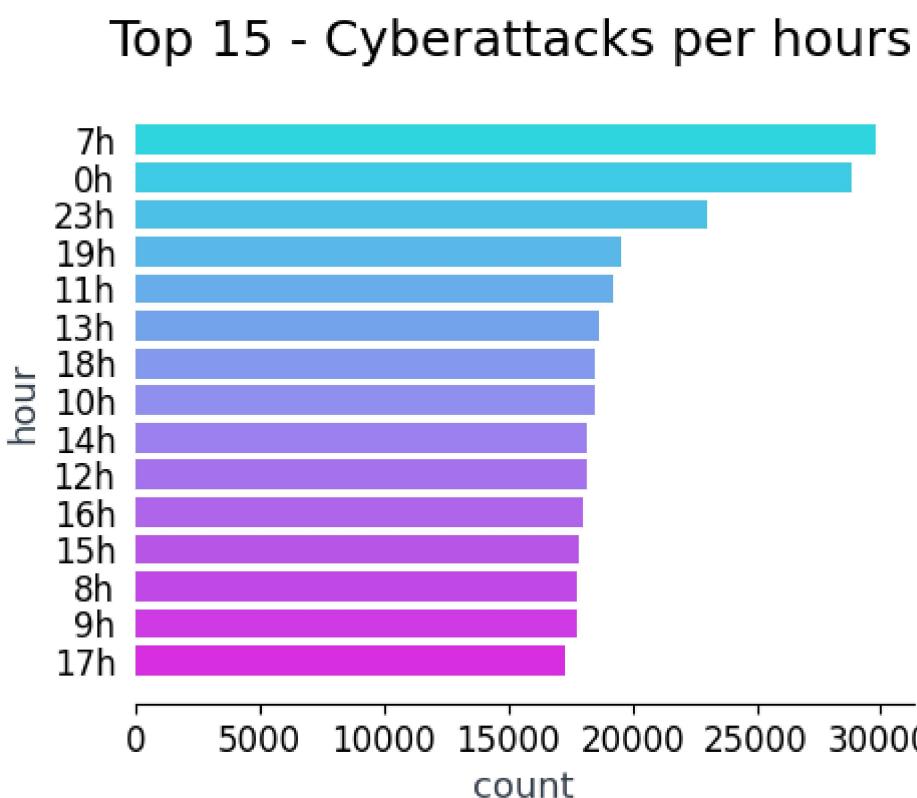


```
In [21]: bar_plot_data(df, 'country', 'Top 15 - Cyberattacks per country and regions')
```

## Top 15 - Cyberattacks per country and regions



```
In [8]: bar_plot_data(df, 'hour', 'Top 15 - Cyberattacks per hours')
```



```
In [9]: df['srcstr'].value_counts().rename_axis('ipaddress').reset_index(name='counts')
```

Out[9]:

	ipaddress	counts
0	175.146.199.252	18472
1	2.186.189.218	11116
2	203.178.148.19	7086
3	128.9.168.98	7046
4	129.82.138.44	6772
...	...	...
69597	113.161.81.118	1
69598	116.252.182.193	1
69599	113.81.65.41	1
69600	122.144.11.167	1
69601	175.201.224.104	1

69602 rows × 2 columns

In [10]:

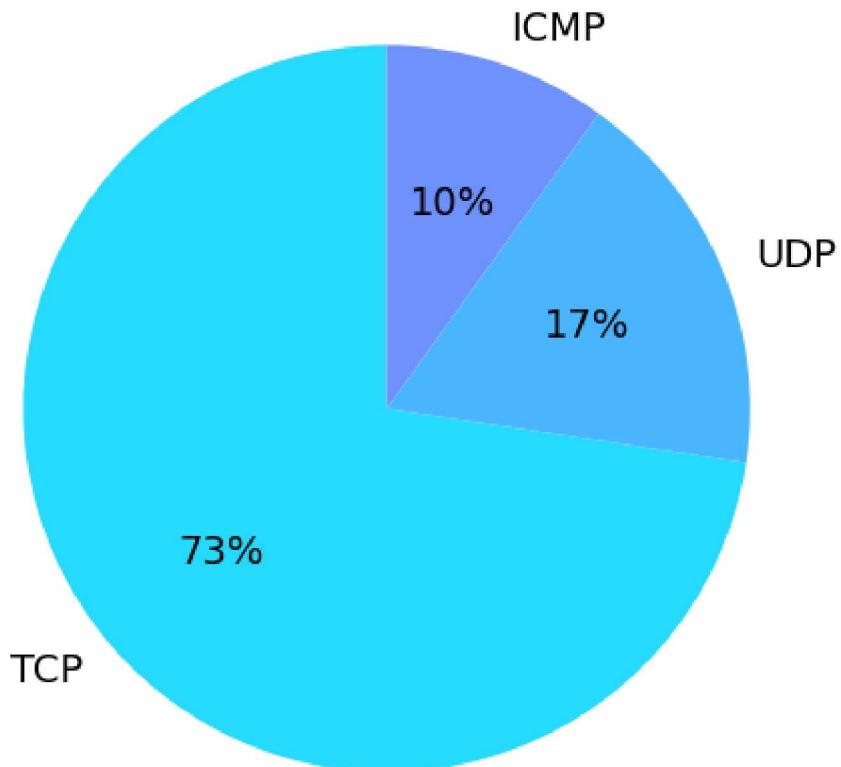
```
def get_pie_plot_count(df: dict, field: str, subtitle: str):
    fig, ax1 = plt.subplots(figsize=(10, 6), dpi=100)
    colors = sns.color_palette('cool')
    df_filter = df[field].value_counts().rename_axis(field).reset_index(name='counts')
    ax1.tick_params(axis='both', which='both', labelsize=12, bottom=True, left=False)
    plt.pie(df_filter['counts'], colors=colors, labels=df_filter[field], autopct = '%.1f%%',
            startangle=90, textprops={'fontsize': 14})

    plt.plot(color="white", lw=3)
    fig.suptitle(subtitle, fontsize=18)
    plt.show()
```

In [11]:

```
get_pie_plot_count(df, 'proto', 'Cyberattacks per communication protocol')
```

## Cyberattacks per communication protocol



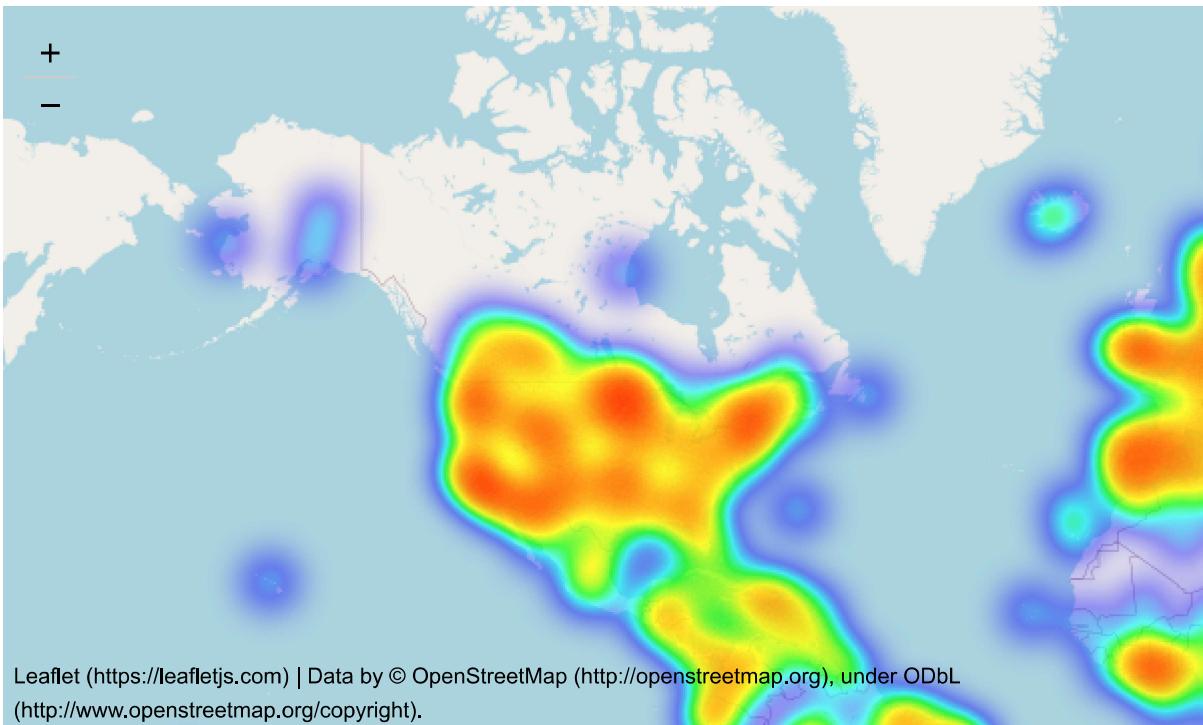
```
In [12]: attack= df[['latitude','longitude']]
attack.latitude.fillna(0, inplace = True)
attack.longitude.fillna(0, inplace = True)

World =folium.Map(location=[0,0],zoom_start=2)
HeatMap(data=attack, radius=16).add_to(World)

print('Top cyberattacks by country')
World
```

Top cyberattacks by country

Out[12]:



Leaflet (<https://leafletjs.com>) | Data by © OpenStreetMap (<http://openstreetmap.org>), under ODbL (<http://www.openstreetmap.org/copyright>).

In [13]:

```
def get_histplot_central_tendency(df: dict, fields: list):
    colors = sns.color_palette('cool')
    for field in fields:
        f, (ax1) = plt.subplots(figsize=(10, 5), dpi=100)
        v_dist_1 = df[field].values

        for spline in ['top', 'right', 'left']:
            ax1.spines[spline].set_visible(False)

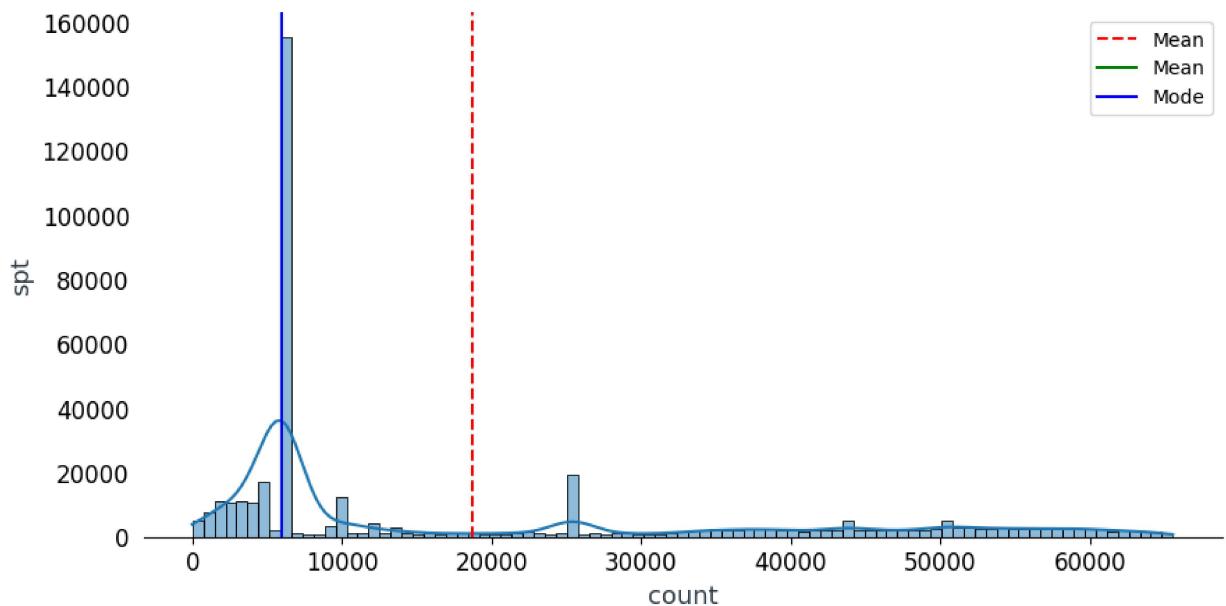
        sns.histplot(v_dist_1, ax=ax1, palette=colors, kde=True)

        mean=df[field].mean()
        median=df[field].median()
        mode=df[field].mode().values[0]

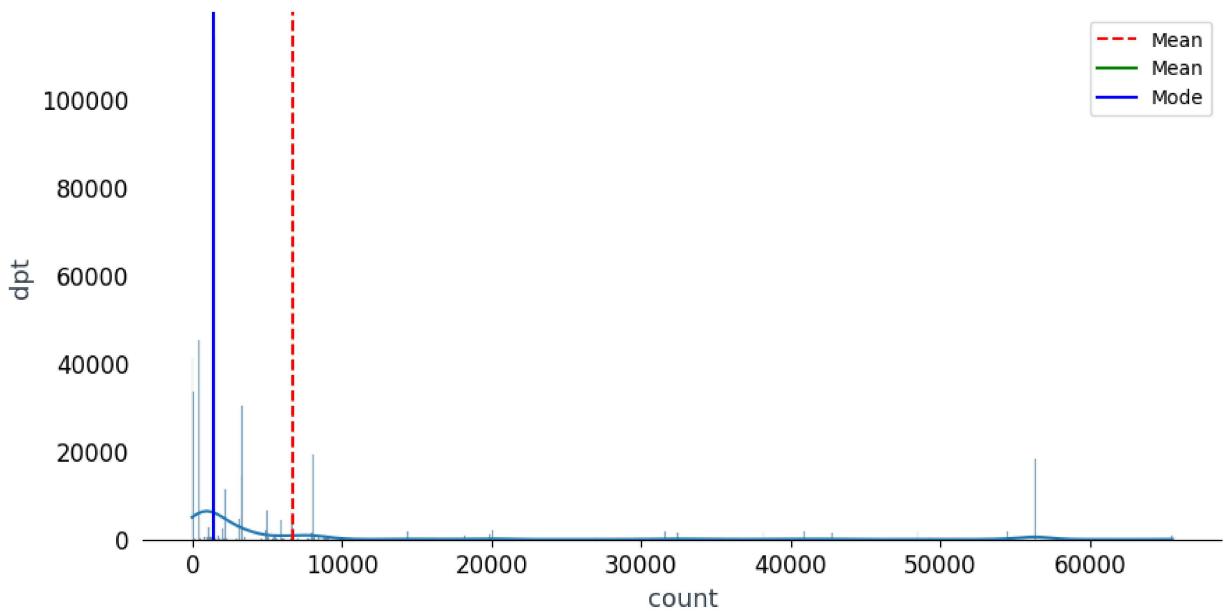
        ax1.axvline(mean, color='r', linestyle='--', label="Mean")
        ax1.axvline(median, color='g', linestyle='-', label="Median")
        ax1.axvline(mode, color='b', linestyle='-', label="Mode")
        ax1.tick_params(axis='both', which='both', labelsize=12, bottom=True, left=False)
        ax1.set_xlabel(f'count', fontsize=13, color = '#333F4B')
        ax1.set_ylabel(f'{field}', fontsize=13, color = '#333F4B')
        ax1.legend()
        plt.grid(False)
        plt.plot(color="white", lw=3)
        f.suptitle(f"Representation Histogram for {field}", fontsize=18)
```

In [14]: `get_histplot_central_tendency(df, ['spt', 'dpt'])`

### Representation Histogram for spt



### Representation Histogram for dpt

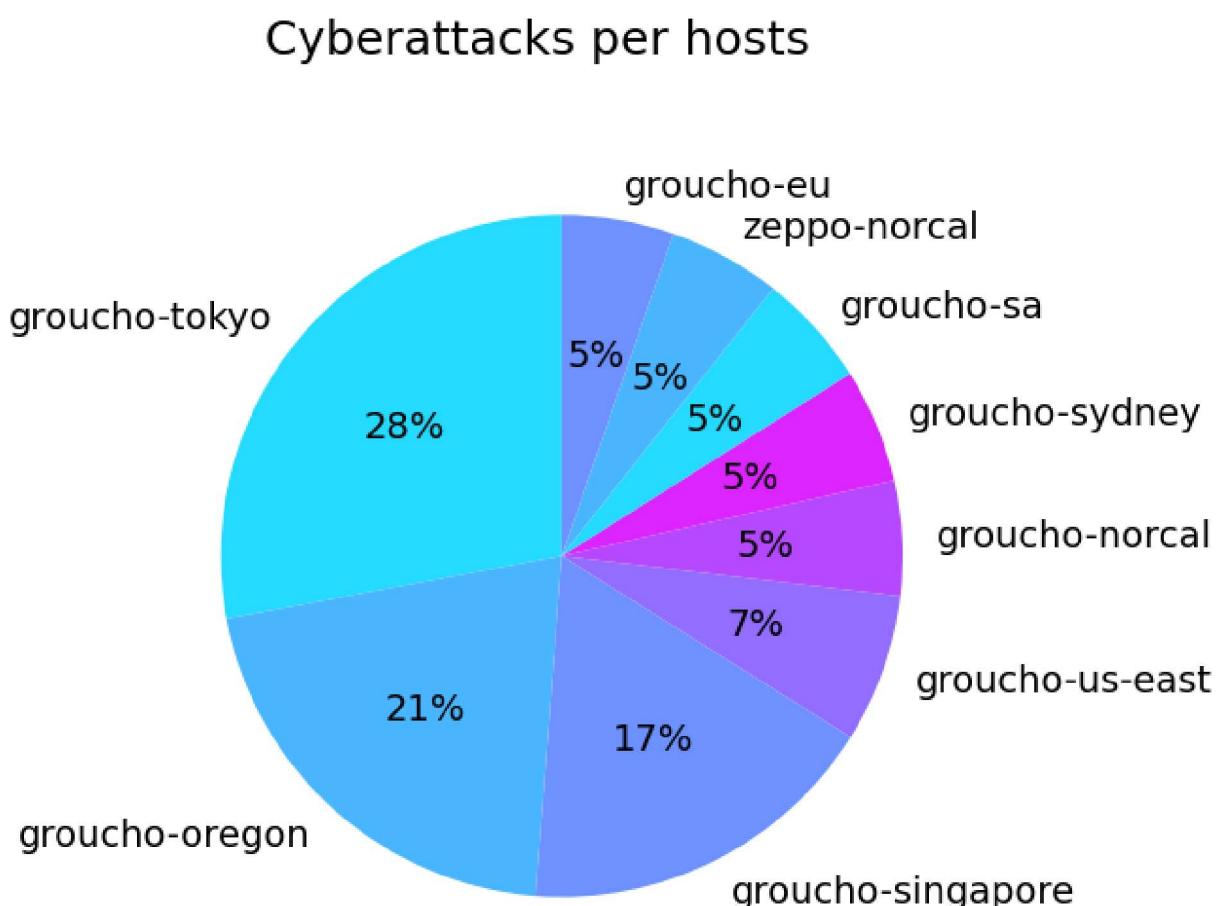


```
In [15]: df['type'].value_counts().rename_axis('id_type_attack').reset_index(name='counts')
```

```
Out[15]: id_type_attack  counts
```

	id_type_attack	counts
0	8.0	38597
1	3.0	4251
2	11.0	1156
3	0.0	536
4	13.0	142
5	5.0	127
6	12.0	2

```
In [16]: get_pie_plot_count(df, 'host', 'Cyberattacks per hosts')
```



```
In [17]: df_filter = df[(df.host == "groucho-oregon") & (df.month == "May") & (df.hour == "20h")]  
df_filter
```

Out[17]:

	datetime	host	src	proto	type	spt	dpt	srcstr	cc	country	.
1437	3/4/13 20:00	groucho-oregon	2036715561	TCP	NaN	27984.0	22.0	121.101.208.41	CN	China	
1441	3/4/13 20:05	groucho-oregon	3736745024	TCP	NaN	6000.0	1433.0	222.186.52.64	CN	China	
1442	3/4/13 20:06	groucho-oregon	3658020275	TCP	NaN	6000.0	1433.0	218.8.245.179	CN	China	
1444	3/4/13 20:08	groucho-oregon	3630406730	TCP	NaN	6000.0	1433.0	216.99.156.74	US	United States	
1445	3/4/13 20:09	groucho-oregon	1019141837	TCP	NaN	6000.0	1433.0	60.190.222.205	CN	China	
...	...	...	...	...	...	...	...	...	...	...	...
53246	3/31/13 20:46	groucho-oregon	1995742862	TCP	NaN	6000.0	1433.0	118.244.158.142	CN	China	
53248	3/31/13 20:49	groucho-oregon	3736734735	TCP	NaN	6000.0	1433.0	222.186.12.15	CN	China	
53249	3/31/13 20:49	groucho-oregon	3689613428	TCP	NaN	6000.0	1433.0	219.235.8.116	CN	China	
53258	3/31/13 20:52	groucho-oregon	3683824005	TCP	NaN	6000.0	1433.0	219.146.177.133	CN	China	
53265	3/31/13 20:57	groucho-oregon	1995742035	TCP	NaN	6000.0	3306.0	118.244.155.83	CN	China	

442 rows × 21 columns



```
In [18]: df_filter['proto'].value_counts()
```

```
Out[18]: TCP      392
          UDP      35
          ICMP     15
          Name: proto, dtype: int64
```

```
In [19]: import networkx as nx

G = nx.Graph()
G.add_node('groucho-oregon', color='red')

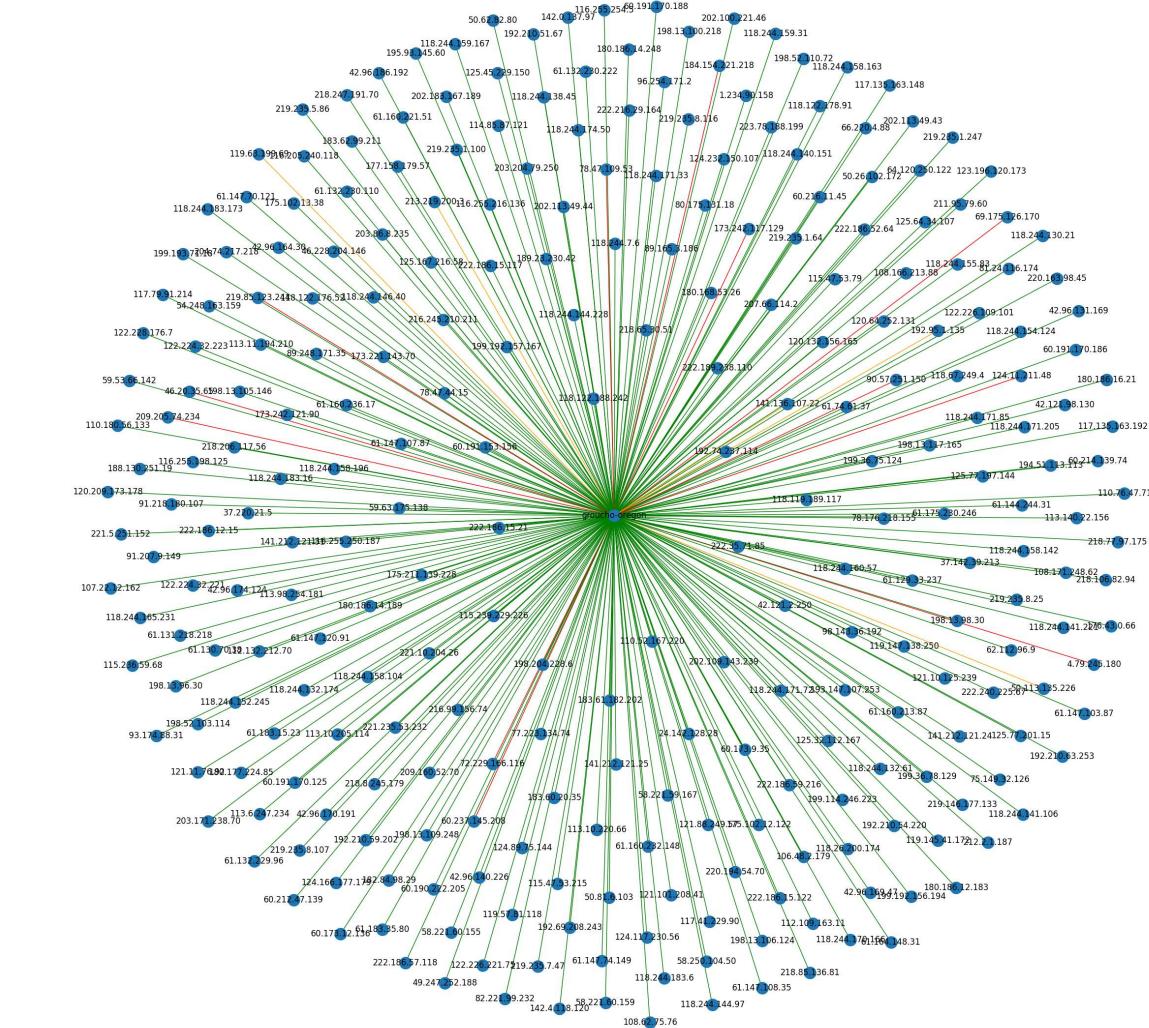
ip_listx = []
edge_color = []
for ip, proto in zip(df_filter['srcstr'], df_filter['proto']):
    if ip not in ip_listx:
        if 'TCP' == proto:
            edge_color.append('green')
        if 'UDP' == proto:
            edge_color.append('red')
        if 'ICMP' == proto:
            edge_color.append('orange')
    G.add_node(ip)
    ip_listx.append(ip)
```

```

for ip in ip_listx:
    G.add_edge('groucho-oregon', ip)
fig = plt.figure(1, figsize=(25, 25), dpi=100)
print('TPC == green\nUDP == red\nICMP == orange')
nx.draw(G, with_labels=True, font_weight='normal', edge_color=edge_color, arrowsize=30)

TPC == green
UDP == red
ICMP == orange

```



```
In [20]: from matplotlib.pyplot import figure, text
```

```

Gx = nx.Graph()
Gx.add_node('groucho-oregon')

country_list = []
ip_list = []
edge_color = []
for country, ip, proto in zip(df_filter['cc'], df_filter['srcstr'], df_filter['proto']):
    if country != 'CN' and country != 'US':
        if country not in country_list:

```

```

Gx.add_node(country)
country_list.append(country)

if ip not in ip_list:
    if 'TCP' == proto:
        edge_color.append('green')
    if 'UDP' == proto:
        edge_color.append('red')
    if 'ICMP' == proto:
        edge_color.append('orange')

Gx.add_node(ip)
Gx.add_edge(country, ip)
ip_list.append(ip)

for country in country_list:
    Gx.add_edge('groucho-oregon', country)

options = {
    "node_color": "black",
    "node_size": 50,
    "linewidths": 0,
    "width": 0.1,
}

pos = nx.spring_layout(G)
d = dict(G.degree)
fig = plt.figure(figsize=(26, 26), dpi=100)
nx.draw(Gx, alpha=0.6, node_size=2050, with_labels=True, font_weight='normal',
        edge_color=edge_color, arrowsize=30, arrowstyle='fancy')

for node, (x, y) in pos.items():
    text(x, y, node, fontsize=12)

```

