

```
In [39]: import pandas as pd  
import numpy as np
```

```
In [2]: pip install ucimlrepo
```

```
Collecting ucimlrepo  
  Downloading ucimlrepo-0.0.6-py3-none-any.whl (8.0 kB)  
Installing collected packages: ucimlrepo  
Successfully installed ucimlrepo-0.0.6  
Note: you may need to restart the kernel to use updated packages.
```

```
In [8]: from ucimlrepo import fetch_ucirepo  
import pandas as pd  
import numpy as np  
import matplotlib as plt  
  
# fetch dataset  
rt_iot2022 = fetch_ucirepo(id=942)  
  
# data (as pandas dataframes)  
X = rt_iot2022.data.features  
y = rt_iot2022.data.targets  
  
# metadata  
print(rt_iot2022.metadata)  
  
# variable information  
print(rt_iot2022.variables)
```

```
{'uci_id': 942, 'name': 'RT-IoT2022 ', 'repository_url': 'https://archive.ics.uci.edu/dataset/942/rt-iot2022', 'data_url': 'https://archive.ics.uci.edu/static/public/942/data.csv', 'abstract': 'The RT-IoT2022, a proprietary dataset derived from a real-time IoT infrastructure, is introduced as a comprehensive resource integrating a diverse range of IoT devices and sophisticated network attack methodologies. This dataset encompasses both normal and adversarial network behaviours, providing a general representation of real-world scenarios.\nIncorporating data from IoT devices such as Thing Speak-LED, Wipro-Bulb, and MQTT-Temp, as well as simulated attack scenarios involving Brute-Force SSH attacks, DDoS attacks using Hping and Slowloris, and Nmap patterns, RT-IoT2022 offers a detailed perspective on the complex nature of network traffic. The bidirectional attributes of network traffic are meticulously captured using the Zeek network monitoring tool and the Flowmeter plugin. Researchers can leverage the RT-IoT 2022 dataset to advance the capabilities of Intrusion Detection Systems (IDS), fostering the development of robust and adaptive security solutions for real-time IoT networks. ', 'area': 'Engineering', 'tasks': ['Classification', 'Regression', 'Clustering'], 'characteristics': ['Tabular', 'Sequential', 'Multivariate'], 'num_instances': 123117, 'num_features': 83, 'feature_types': ['Real', 'Categorical'], 'demographics': [], 'target_col': ['Attack_type'], 'index_col': ['id'], 'has_missing_values': 'no', 'missing_values_symbol': None, 'year_of_dataset_creation': 2023, 'last_updated': 'Fri Mar 08 2024', 'dataset_doi': '10.24432/C5P338', 'creators': ['B. S.', 'Rohini Nagapadma'], 'intro_paper': {'title': 'Quantized autoencoder (QAE) intrusion detection system for anomaly detection in resource-constrained IoT devices using RT-IoT2022 dataset', 'authors': 'B. S. Sharmila, Rohini Nagapadma', 'published_in': 'Cybersecurity', 'year': 2023, 'url': 'https://www.semanticscholar.org/paper/753f6ede01b4acaa325e302c38f1e0c1ade74f5b', 'doi': None}, 'additional_info': {'summary': None, 'purpose': None, 'funded_by': None, 'instances_represent': None, 'recommended_data_splits': None, 'sensitive_data': None, 'preprocessing_description': None, 'variable_info': 'Column Details:\nid.orig_p\nid.resp_p\nnproto\nservice\nnflow_duration\nnfwd_pkts_tot\nnbwd_pkts_tot\nnfwd_data_pkts_tot\nnbwd_data_pkts_tot\nnfwd_pkts_per_sec\nnbwd_pkts_per_sec\nnflow_pkts_per_sec\nndown_up_ratio\nnfwd_header_size_tot\nnfwd_header_size_min\nnfwd_header_size_max\nnbwd_header_size_tot\nnbwd_header_size_min\nnbwd_header_size_max\nnflow_FIN_flag_count\nnflow_SYN_flag_count\nnflow_RST_flag_count\nnfwd_PSH_flag_count\nnbwd_PSH_flag_count\nnflow_ACK_flag_count\nnfwd_URG_flag_count\nnbwd_URG_flag_count\nnflow_CWR_flag_count\nnflow_ECE_flag_count\nnfwd_pkts_payload.min\nnfwd_pkts_payload.max\nnfwd_pkts_payload.tot\nnfwd_pkts_payload.avg\nnfwd_pkts_payload.std\nnbwd_pkts_payload.min\nnbwd_pkts_payload.max\nnbwd_pkts_payload.tot\nnbwd_pkts_payload.avg\nnbwd_pkts_payload.std\nnflow_pkts_payload.min\nnflow_pkts_payload.max\nnflow_pkts_payload.tot\nnflow_pkts_payload.avg\nnflow_pkts_payload.std\nnfwd_iat.min\nnfwd_iat.max\nnfwd_iat.tot\nnfwd_iat.avg\nnfwd_iat.std\nnbwd_iat.min\nnbwd_iat.max\nnbwd_iat.tot\nnbwd_iat.avg\nnbwd_iat.std\nnflow_iat.min\nnflow_iat.max\nnflow_iat.tot\nnflow_iat.avg\nnflow_iat.std\nnpayload_bytes_per_second\nnfwd_subflow_pkts\nnbwd_subflow_pkts\nnfwd_subflow_bytes\nnbwd_subflow_bytes\nnfwd_bulk_bytes\nnbwd_bulk_bytes\nnfwd_bulk_packets\nnbwd_bulk_packets\nnfwd_bulk_rate\nnbwd_bulk_rate\nnactive.min\nnactive.max\nnactive.tot\nnactive.avg\nnactive.std\nnidle.min\nnidle.max\nnidle.tot\nnidle.avg\nnidle.std\nnfwd_init_window_size\nnbwd_init_window_size\nnfwd_last_window_size\nnAttack_type', 'citation': None}}
```

	name	role	type	demographic	description	units	\
0	id.orig_p	Feature	Integer	None	None	None	
1	id.resp_p	Feature	Integer	None	None	None	
2	proto	Feature	Categorical	None	None	None	
3	service	Feature	Continuous	None	None	None	
4	flow_duration	Feature	Continuous	None	None	None	
..	
80	fwd_init_window_size	Feature	Integer	None	None	None	
81	bwd_init_window_size	Feature	Integer	None	None	None	
82	fwd_last_window_size	Feature	Integer	None	None	None	
83	Attack_type	Target	Categorical	None	None	None	
84	id	ID	Integer	None	None	None	

missing_values

0 no

```

1          no
2          no
3          no
4          no
..         ...
80         no
81         no
82         no
83         no
84         no

```

[85 rows x 7 columns]

```

-----
NameError                                Traceback (most recent call last)
Cell In[8], line 19
     16 # variable information
     17 print(rt_iot2022.variables)
--> 19 pd.concat([x,y])

NameError: name 'x' is not defined

```

In [10]: `pd.concat([X,y])` *#merge the two dataframes (X,y) into one dataframe*

Out[10]:

	id.orig_p	id.resp_p	proto	service	flow_duration	fwd_pkts_tot	bwd_pkts_tot	fwd_data_pkts
0	38667.0	1883.0	tcp	mqtt	32.011598	9.0	5.0	
1	51143.0	1883.0	tcp	mqtt	31.883584	9.0	5.0	
2	44761.0	1883.0	tcp	mqtt	32.124053	9.0	5.0	
3	60893.0	1883.0	tcp	mqtt	31.961063	9.0	5.0	
4	51087.0	1883.0	tcp	mqtt	31.902362	9.0	5.0	
...
123112	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
123113	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
123114	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
123115	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
123116	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

246234 rows × 84 columns

In [30]: `df = pd.concat([X,y], axis=1)` *#concatinated the NaN values column wise (axis=1)*
`df`

Out[30]:

	id.orig_p	id.resp_p	proto	service	flow_duration	fwd_pkts_tot	bwd_pkts_tot	fwd_data_pkts
0	38667	1883	tcp	mqtt	32.011598	9	5	
1	51143	1883	tcp	mqtt	31.883584	9	5	
2	44761	1883	tcp	mqtt	32.124053	9	5	
3	60893	1883	tcp	mqtt	31.961063	9	5	
4	51087	1883	tcp	mqtt	31.902362	9	5	
...
123112	59247	63331	tcp	-	0.000006	1	1	
123113	59247	64623	tcp	-	0.000007	1	1	
123114	59247	64680	tcp	-	0.000006	1	1	
123115	59247	65000	tcp	-	0.000006	1	1	
123116	59247	65129	tcp	-	0.000006	1	1	

123117 rows × 84 columns

In [15]: `df.dtypes` *#we want to change the datatypes of proto, service, and Attack_type to numer*

Out[15]:

```
id.orig_p          int64
id.resp_p          int64
proto              object
service            object
flow_duration      float64
...
idle.std           float64
fwd_init_window_size int64
bwd_init_window_size int64
fwd_last_window_size int64
Attack_type        object
Length: 84, dtype: object
```

In [35]: `df['proto'].unique()` *#we can see that these are the objects we want to replace as numer*

Out[35]: `array(['tcp', 'udp', 'icmp'], dtype=object)`

In [36]: `df['service'].unique()` *#we can see that these are the objects we want to replace as nu*

Out[36]: `array(['mqtt', '-', 'http', 'dns', 'ntp', 'ssl', 'dhcp', 'irc', 'ssh', 'radius'], dtype=object)`

In [37]: `df['Attack_type'].unique()` *#we can see that these are the objects we want to replace a*

Out[37]: `array(['MQTT_Publish', 'Thing_Speak', 'Wipro_bulb', 'ARP_poisoning', 'DDOS_Slowloris', 'DOS_SYN_Hping', 'Metasploit_Brute_Force_SSH', 'NMAP_FIN_SCAN', 'NMAP_OS_DETECTION', 'NMAP_TCP_scan', 'NMAP_UDP_SCAN', 'NMAP_XMAS_TREE_SCAN'], dtype=object)`

In []:

In [16]: `df.describe()`

Out[16]:

	id.orig_p	id.resp_p	flow_duration	fwd_pkts_tot	bwd_pkts_tot	fwd_data_pkts_tot
count	123117.000000	123117.000000	123117.000000	123117.000000	123117.000000	123117.000000
mean	34639.258738	1014.305092	3.809566	2.268826	1.909509	1.471218
std	19070.620354	5256.371994	130.005408	22.336565	33.018311	19.635196
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	17702.000000	21.000000	0.000001	1.000000	1.000000	1.000000
50%	37221.000000	21.000000	0.000004	1.000000	1.000000	1.000000
75%	50971.000000	21.000000	0.000005	1.000000	1.000000	1.000000
max	65535.000000	65389.000000	21728.335580	4345.000000	10112.000000	4345.000000

8 rows × 81 columns

In [17]: `df.head()`

Out[17]:

	id.orig_p	id.resp_p	proto	service	flow_duration	fwd_pkts_tot	bwd_pkts_tot	fwd_data_pkts_tot
0	38667	1883	tcp	mqtt	32.011598	9	5	3
1	51143	1883	tcp	mqtt	31.883584	9	5	3
2	44761	1883	tcp	mqtt	32.124053	9	5	3
3	60893	1883	tcp	mqtt	31.961063	9	5	3
4	51087	1883	tcp	mqtt	31.902362	9	5	3

5 rows × 84 columns

In [18]: `df.tail()`

Out[18]:

	id.orig_p	id.resp_p	proto	service	flow_duration	fwd_pkts_tot	bwd_pkts_tot	fwd_data_pkts_tot
123112	59247	63331	tcp	-	0.000006	1	1	
123113	59247	64623	tcp	-	0.000007	1	1	
123114	59247	64680	tcp	-	0.000006	1	1	
123115	59247	65000	tcp	-	0.000006	1	1	
123116	59247	65129	tcp	-	0.000006	1	1	

5 rows × 84 columns

Conclusion

In this quiz, the dataset came with missing values (NaN), what I did was concatenate the two dataframes (X, y) into a single dataframe column wise. Furthermore, after concatenating the two

separate dataframes, we now check the dtypes of every column. We were tasked to change the dtypes of a column that is labeled as object. To do that it's either we create a function on our own or use the lambda function. However, I wasn't able to convert the said columns into a numeric datatype. I should practice more of my OOP so that next quiz I would complete all of my tasks