

Computer network Intrusion Detection:

Name	ID
Hossam el-din ahmed	6721
Adham el Kazzaz	6713
Saeed abokhatwa	6829

## Intro:

The KDD-CUP-1999 dataset is a famous benchmark dataset used in the field of network intrusion detection. It was created for a competition organized by the University of California, Irvine (UCI), and the International Conference on Knowledge Discovery and Data Mining (KDD). The goal of the competition was to encourage the development of effective and efficient methods for detecting network intrusions in real-time by analyzing network traffic data. The dataset consists of a set of network traffic data captured over nine weeks from an environment that simulated a typical US Air Force local area network (LAN), subjected to various types of attacks. The dataset has become a standard benchmark for evaluating the performance of intrusion detection systems (IDSs), and it is often used to test and compare the effectiveness of various machine learning algorithms and techniques for detecting different types of network attacks.

(Attack vs normal)

```
[23]: print(f'percentage of unique data is {(len(training_df)-len(training_df.drop_duplicates(
percentage of unique data is 70.53040255373759 %
```

```
[24]: def remove_duplicates(df):
# count the duplicate rows and add this count as a new column
duplicates_count = df.groupby(list(df.columns)).size().reset_index().rename(columns=
df_with_counts = pd.merge(df, duplicates_count, on=list(df.columns), how='left')
return df_with_counts.drop_duplicates()
```

Since 30% of the data are duplicated , removing the duplicates will improve the effectiveness of clustering but also just removing is not the optimal solution , in this case I introduced a new feature to each instance named “duplicate count” which state the number of times the instance was repeated throughout the dataset.

This will keep the information of the rows that was removed.

```
: training_df['Duplicate_count'].value_counts()
1      126209
2       8251
3       3625
4       2165
5       1375
...
148         1
127         1
153         1
111         1
59          1
Name: Duplicate_count, Length: 147, dtype: int64
```

```
[30]: training_df.corr()['target_encoded'].sort_values()
```

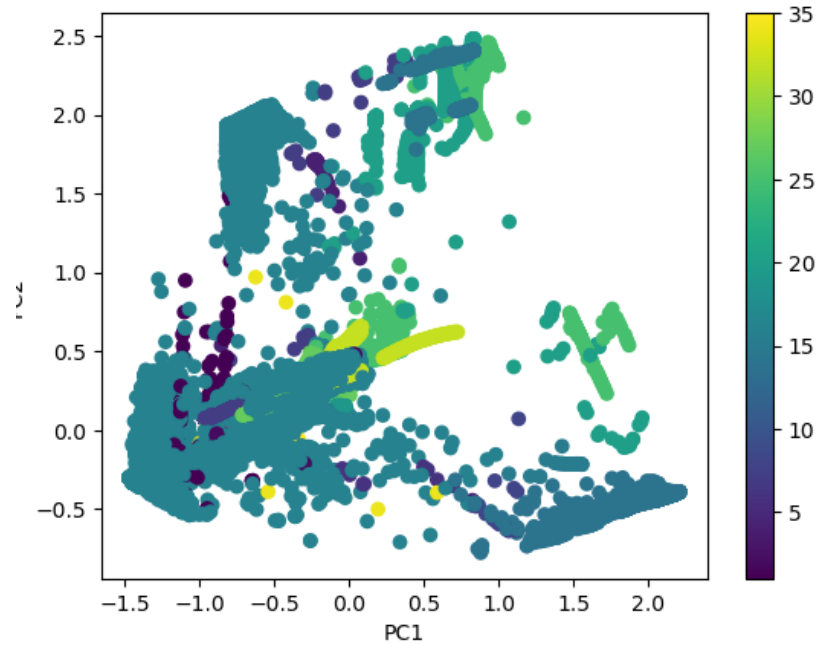
```
[30...] dst_host_serror_rate    -0.326148
      srv_serror_rate      -0.325529
      dst_host_srv_serror_rate -0.325256
      flag_S0              -0.325138
      serror_rate          -0.323583
      ...
      wrong_fragment       0.469979
      target_encoded       1.000000
      num_outbound_cmds    NaN
      is_host_login        NaN
      service_icmp         NaN
      Name: target_encoded, Length: 121, dtype: float64
```

num\_outbound\_cmds NaN is\_host\_login NaN so dropping them will affect nothing

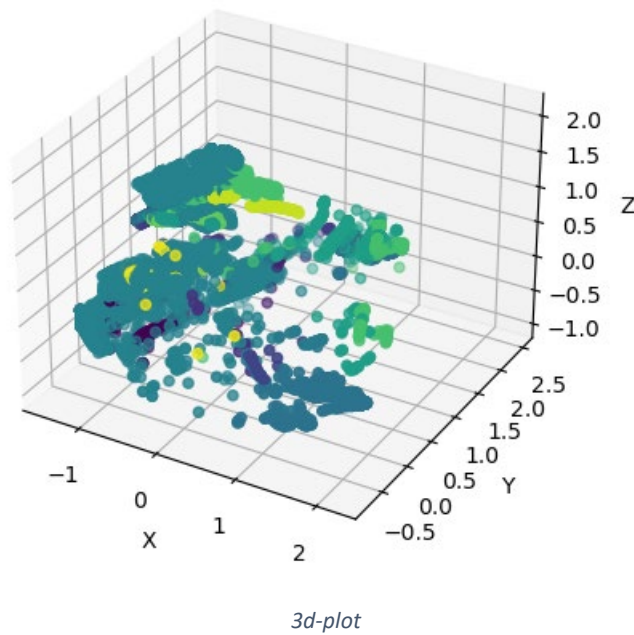
```
[31]: training_df.drop(inplace=True, columns=['num_outbound_cmds', 'is_host_login', 'service_icmp'])
      testing_df.drop(inplace=True, columns=['num_outbound_cmds', 'is_host_login', 'service_icmp'])
```

dropped the feature who has no correlation with the target label, this will save us some time and will provide better results.

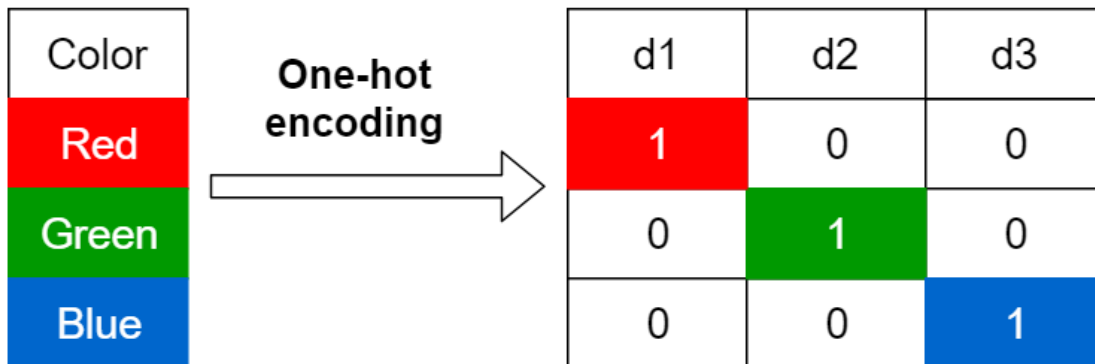
## Visualization of the data using PCA



It's obvious here the result of any clustering algorithm will be bad since there is overlapping points of different classes and more likely to cluster different classes together.



Encoding chosen -> OneHotEncoding:



One hot encoding is a popular technique used in machine learning to represent categorical data as numerical values. It works by converting each category value into a binary vector that has a dimension equal to the total number of categories. In this vector, only one bit is set to 1, representing the category value, while all the other bits are set to 0. One hot encoding is used because many machine learning algorithms require numerical input data, and categorical data cannot be used directly. By converting categorical data into numerical values using one hot encoding, machine learning models can effectively process and analyze categorical data.

## Kmeans clustering:

On training data:

K	Precision	Recall	F-measure	Entropy
7	[0.99721913236929 93, 0.629620106704 3774, 0.9830784022 375666, 0.96094532 3452834, 0.9364118 092354277, 0.99901 63594947195, 0.820 5407614247769]	[0.3460054033191 818, 0.1890196835 198765, 0.6402791 693232535, 0.0925 8973369355461, 0. 169004463065852 98, 0.37238517946 73871, 0.13717096 274706259]	0.40183054443829 01	0.35811831459 179383
15	[0.99947932937623 66, 0.998149935649 9357, 0.3540229885 0574715, 0.9879336 349924586, 1.0, 0.86 78068410462777, 0. 7734056987788331, 0.983205701721131 2, 0.8933764135702 746, 0.88039623333 74098, 0.540816326 5306123, 0.9864864 864864865, 0.96622 61380323054, 0.961 3303947104789, 0.9 061054579093432]	[0.1092768011658 621, 0.7183905827 865689, 0.9609984 399375975, 0.0372 8709354221696, 0. 288607796702796 24, 0.04910510975 4986796, 0.019468 986246470536, 0.2 952796247381364 6, 0.012592221513 799072, 0.0819632 9356043355, 0.004 223973039438928 5, 0.053192458329 53821, 0.04494944 894799162, 0.0925 8973369355461, 0. 189019683519876 5]	0.23376920248258 506	0.19716632018 485683
23	[0.98619499568593 62, 0.858895705521 4724, 0.9988998604 701084, 0.99943470 88750707, 0.987244 8979591837, 0.4, 1.0 , 0.88045471213788 04, 0.678111587982 8327, 0.8679775280 898876, 0.97555791 7109458, 0.9843665 318714303, 0.96108 2910321489, 0.9381 255686988171, 1.0, 0.973248973248973 3, 0.6763990267639]	[0.2212291647691 0465, 0.006375808 361417251, 0.7183 905827865689, 0.0 603880134802805 35, 0.00440613899 2622279, 0.000751 4345568813189, 0. 037093542216959 65, 0.08200883504 87294, 0.98595943 83775351, 0.74278 84615384616, 0.03 135531469168412, 0.26883140541032 88, 0.00646689133]	0.23539282598618 855	0.18424488025 534586

	902, 0.96622613803 23054, 0.969678469 6784697, 0.7160804 020100503, 0.93383 5446658404, 0.8561 151079136691, 0.99 75594874923734]	8008926, 0.011738 318608252118, 0.0 486269241278805, 0.09982694234447 581, 0.8540706605 222734, 0.0449494 4894799162, 0.091 95291393284447, 0.00973449312323 5268, 0.189019683 5198765, 0.043355 49685763731, 0.01 861508334092358 2]		
31	[0.98789388575674 14, 0.940289910722 8568, 0.7191011235 955056, 0.95060483 87096774, 0.995296 3311382879, 0.5107 913669064749, 0.96 92902176123652, 1. 0, 1.0, 0.9834399698 908544, 0.56480505 79557429, 0.985939 0518007423, 0.9976 359338061466, 0.99 9767576990122, 0.9 25249169435216, 0. 9204892966360856, 0.998847926267281 1, 0.9998167155425 219, 0.99668141592 92036, 0.914201183 4319527, 0.3455710 9557109555, 0.9991 620404315492, 0.57 71495877502945, 0. 8022284122562674, 0.986290101453249 2, 0.9746972318339 1, 0.9882210812443 37, 0.983193277310 9243, 0.9706717123 935666, 0.86063829 78723405, 0.999435 3472614342]	[0.2573549503597 7773, 0.189019683 5198765, 0.769230 7692307693, 0.018 19760710150521, 0.01204572365424 902, 0.0008083614 172511158, 0.0919 7221150135083, 0. 021700519172966 573, 0.0269036342 10766008, 0.02974 997722925585, 0.6 002239641657335, 0.21475088805902 176, 0.0336323891 06476, 0.33203396 37205712, 0.00634 1652245195373, 0. 003426996994261 7726, 0.009871117 58812278, 0.06210 72046634484, 0.01 025822023863739 8, 0.017590399854 267236, 0.9251170 046801872, 0.3681 5901196449247, 0. 016736496948720 283, 0.3225083986 56215, 0.04095318 3350031876, 0.102 62774387466983, 0.03725293742599 508, 0.0053283541 30612989, 0.01168 1391747882321, 0. 009210766007833	0.18042175838862 057	0.16280182471 318053

		136, 0.0604563257 12724295]		
45	[0.99981597350018 4, 1.0, 0.4667093469 910371, 1.0, 0.97062 13417379559, 0.949 4845360824742, 0.8 333333333333334, 1.0, 0.99763593380 61466, 1.0, 1.0, 0.989 2413961733603, 1.0 , 1.0, 0.77171717171 71717, 0.961116650 0498505, 1.0, 0.9992 589455854329, 0.96 54978962131837, 0. 9758874748827864, 0.976666666666666 7, 1.0, 0.5382352941 17647, 0.983624454 1484717, 0.5887445 887445888, 0.99956 98924731182, 0.982 2784810126582, 0.9 165919282511211, 0.999642474079370 8, 0.6587677725118 484, 0.97805343511 45038, 1.0, 0.687640 4494382022, 0.9997 59239195859, 0.993 4904996481351, 0.9 414141414141414, 0.999707602339181 3, 0.9041666666666 667, 1.0, 1.0, 0.99324 32432432432, 0.859 4386600271616, 0.8 950437317784257, 0.999351070733290 1, 0.9962292609351 433]	[0.0618567264778 213, 0.0433441114 8556335, 0.008299 936241916385, 0.0 149489935331086 62, 0.16640293323 041297, 0.0177730 6059436511, 5.692 686036979689e-0 5, 0.113443847344 93124, 0.03363238 9106476, 0.024410 237726568904, 0.0 188541761544767 27, 0.18424947627 28846, 0.63465783 66445916, 0.03571 5912196010564, 0. 017396848529009 93, 0.01097549867 929684, 0.0524410 2377265689, 0.182 14974913160942, 0.03918845067856 8174, 0.016588487 111758813, 0.0226 16750289463528, 0.05279814743342 339, 0.0041670461 79069132, 0.01025 8220238637398, 0. 001548410602058 4752, 0.026459604 69988159, 0.00441 7524364696238, 0. 011635850259586 484, 0.0318335003 1879042, 0.854070 6605222734, 0.011 670006375808362, 0.06335959559158 394, 1.0, 0.3205326 128907758, 0.0642 931961016486, 0.0 106111667729301 4, 0.197935160169 8186, 0.008375144 731763798, 0.0136 39675744603334, 0.00169818602856 04014, 0.00669459	0.14809281862049 387	0.12690482531 99128



		8779488113, 0.043 230257764823754, 0.73798076923076 93, 0.02971825549 9807025, 0.015040 076509700337]		
--	--	----------------------------------------------------------------------------------------------------------------------	--	--

On testing data:

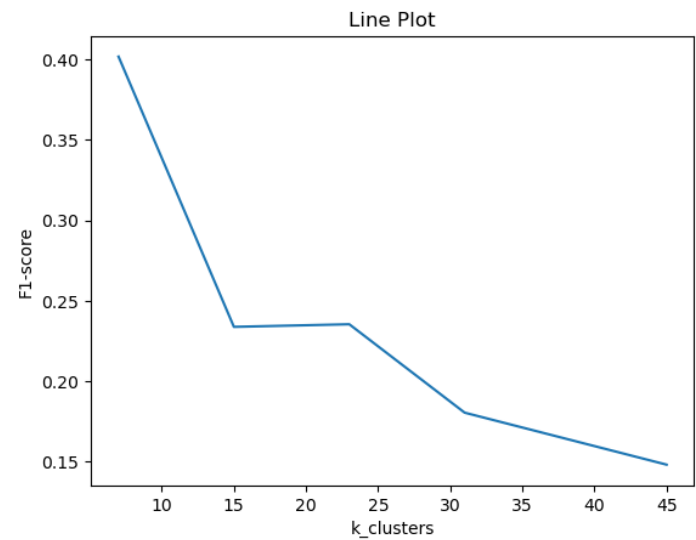
K	Precision	Recall	F-measure	Entropy
7	[0.97211538461538 46, 0.880115644073 2413, 0.9697504035 763069, 0.58023320 37756802, 0.609375 , 0.998046875, 0.648 8982501620221]	[0.1491737163092 6618, 0.673765492 8192013, 0.814956 27491495, 0.05139 681290576431, 0.0 993049902949095 2, 0.125663977965 76825, 0.08358900 507169244]	0.36343807229313 91	0.67633615415 69572
15	[0.90852094613850 1, 0.9837713882518 963, 0.67970479704 79705, 0.995791245 7912458, 1.0, 0.6419 408812046848, 0.89 52583156404812, 0. 912260793315108, 0.339301700984780 7, 0.8678925719816 114, 0.69006622516 55629, 0.502577319 5876289, 0.6007130 124777184, 0.58510 63829787234, 0.895 56238154369]	[0.0665372654603 1348, 0.274296675 19181583, 0.98397 43589743589, 0.04 938117003735938, 0.37042138876714 04, 0.02402270782 459875, 0.0264020 20328512094, 0.36 91273767036086, 0.00791017051739 6113, 0.074864859 22401019, 0.51996 00798403193, 0.24 55919395465995, 0.39186046511627 91, 0.05139681290 576431, 0.6739622 270312807]	0.33527367622261 11	0.54124981855 54266
23	[1.0, 0.79661016949 15254, 0.997317116 7948489, 0.7805325 987144169, 0.85103 01109350238, 0.322 42424242424245, 0. 001, 0.86789257198 16114, 0.855955678 6703602, 0.4101796 407185629, 0.99900 34877927254, 0.928	[0.2813641391689 1034, 0.001961889 2576127563, 0.274 24749163879597, 0.03548097593554 985, 0.6763224181 360201, 0.3575268 817204301, 0.001, 0.07486485922401 019, 0.9903846153 846154, 0.1306005	0.32303140987276 55	0.44373689644 74166

	8025889967637, 0.3 1752178121974833, 0.762376237623762 4, 0.9873501997336 884, 1.0, 0.55109489 0510949, 0.6017857 142857143, 0.89898 9898989899, 0.5225 67703109328, 0.919 2909420533136, 0.6 719846841097639, 1.0]	719733079, 0.0418 46680441633796, 0.33544132072715 127, 0.4408602150 5376344, 0.146806 48236415633, 0.03 095193371318849, 0.13161354955857 49, 0.97419354838 70968, 0.39186046 51162791, 0.04815 0698406452884, 0. 519960079840319 3, 0.673372024395 0423, 0.021977333 91772588, 0.02469 058501867969]		
31	[0.91784192522354 91, 0.924562706827 8517, 0.5398550724 637681, 0.8, 0.66666 6666666666, 0.350 46113306982873, 0. 8924339106654512, 0.987465181058495 8, 1.0, 0.001, 0.92054 79452054794, 1.0, 0. 6923076923076923, 1.0, 0.37792397660 818716, 0.69189907 03851262, 0.866612 7728375101, 1.0, 1.0 , 0.94830508474576 27, 0.678864824495 8925, 0.9983759642 71214, 0.632450331 1258278, 0.001, 0.59 47521865889213, 1. 0, 1.0, 0.9095477386 934674, 0.75227272 72727273, 0.801020 4081632653, 0.8424 182358771061]	[0.3192035564460 5847, 0.673322840 8420224, 0.284080 0762631077, 0.95, 0.11764705882352 941, 0.3575268817 204301, 0.0481506 98406452884, 0.01 479765408135579, 0.01619602195646 2754, 0.001, 0.3353 2934131736525, 0. 272619122158913 04, 0.00093920230 41763195, 0.15301 003344481606, 0.3 970814132104455, 0.51996007984031 93, 0.02237388600 1711436, 0.054682 44526537683, 0.01 465155594515058 5, 0.023354830630 517815, 0.9711538 461538461, 0.1209 4235687586072, 0. 019931960010853 005, 0.001, 0.25692 695214105793, 0.1 349738066912946 3, 0.049735979796 71488, 0.00377768 03790203076, 0.38 488372093023254, 0.00327677248345	0.25575791551987 453	0.42899907916 190977

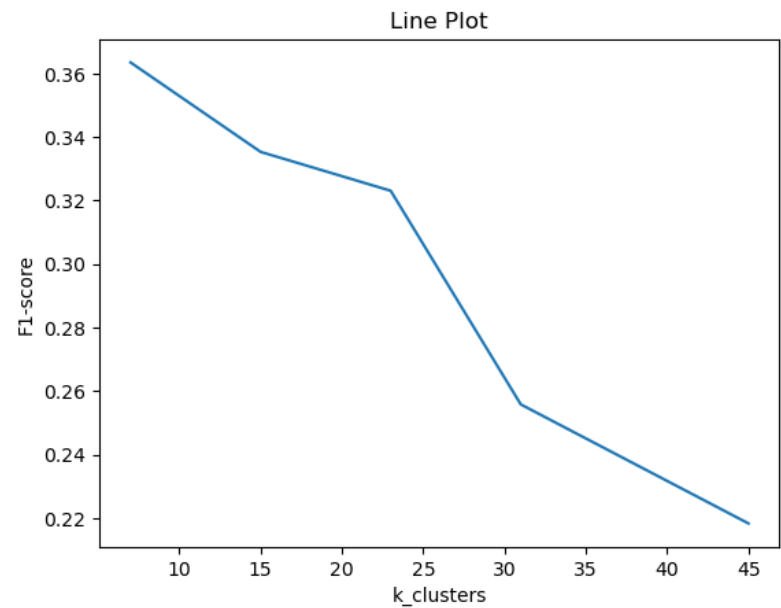
		9604, 0.035480975 93554985]		
45	[1.0, 1.0, 0.79591836 73469388, 1.0, 0.983 3884153765873, 0.8 172043010752689, 0.6, 1.0, 0.69230769 23076923, 1.0, 0.987 3417721518988, 0.9 351612038939001, 0.690671031096563 , 0.56521739130434 78, 0.925592804578 9044, 0.3443223443 223443, 1.0, 0.99834 16252072969, 0.984 13344182262, 0.621 6798277099784, 0.7 536186280679673, 1.0, 0.69518716577 54011, 0.992987377 2791024, 0.8579088 471849866, 0.65147 67932489451, 0.540 4040404040404, 0.9 640479360852197, 0.960762331838565 1, 0.5067114093959 731, 0.80266666666 66666, 1.0, 0.709016 3934426229, 1.0, 0.9 71926647045506, 0. 6434108527131783, 1.0, 0.34822601839 684625, 0.99765807 96252927, 1.0, 1.0, 0. 6741790083708951, 0.670658682634730 5, 0.9967105263157 895, 1.0]	[0.0546407029407 4677, 0.056498236 38678438, 1.0, 0.00 851543422453196 5, 0.556118433995 6718, 0.95, 6.26134 8694508798e-05, 0 .143697952538976 9, 0.000939202304 1763195, 0.031223 25882328387, 0.01 302360528457829 9, 0.218541940600 67206, 0.49069767 44186046, 0.01637 27959697733, 0.02 362615574061319 5, 0.631720430107 5269, 0.072840356 47945233, 0.05921 6997835923664, 0. 050487341640055 94, 0.01807442656 4815393, 0.117794 60948258902, 0.02 759197324414715 6, 0.518962075848 3033, 0.014776782 91904076, 0.80604 53400503779, 0.01 611253730720263 8, 0.002233214367 7081377, 0.015110 72151608123, 0.01 788658610398013, 0.97419354838709 68, 0.35, 0.0770772 0242940329, 0.003 610711080500073, 0.13530395435766 28, 0.08959989981 842088, 0.0791229 7426120114, 0.079 52980523313004, 0.35618279569892 475, 0.0088911151 46202493, 0.00098 36710603974032, 2.08711623150293 24e-05, 0.0218521 06943835704, 0.64	0.21827929427567 405	0.34659061464 387547

		36781609195402, 0.01490261656502 0657, 0.021142487 425124704]		
--	--	------------------------------------------------------------------------	--	--

Training graph:



Testing graph:



## Testing K\_means on iris dataset:

presicion :

[0.8333333333333334, 1.0, 0.8076923076923077]

recall :

[0.8, 1.0, 0.84]

F\_1 score :

0.8799519807923168

the conditional entropy of T given clustering C

[0.6500224216483541, 0.0, 0.7062740891876007]

entropy :

0.4528488591791749

## Using sklearn k-means:

presicion :

[1.0, 0.9230769230769231, 0.7704918032786885]

recall :

[1.0, 0.72, 0.94]

F\_1 score :

0.8852785369639302

the conditional entropy of T given clustering C

[0.0, 0.39124356362925566, 0.7771529943226336]

entropy :

0.4177655442348108

Normalized Cut (Spectral Clustering):

*Data preprocessing:*

The KDD Cup 1999 was loaded into the notebook in a pandas dataframe, with number of rows equal to the number of entries in the dataset (4,898,431), and number of columns equal to the number of features. The features' names were retrieved from the dataset in order to name the corresponding columns after them.

```
]: training_data
```

	duration	protocol_type	service	flag	src_bytes	dst_bytes	land	wrong_fragment	urgent	hot	...	dst_host_srv_count	dst_host_same_srv_rate	dst_host_diff_srv_rate	dst_host_same_src_port_rate	dst_host_srv_diff_...
0	0	tcp	http	SF	215	45076	0	0	0	0	...	0	0.0	0.0	0.00	
1	0	tcp	http	SF	162	4528	0	0	0	0	...	1	1.0	0.0	1.00	
2	0	tcp	http	SF	236	1228	0	0	0	0	...	2	1.0	0.0	0.50	
3	0	tcp	http	SF	233	2032	0	0	0	0	...	3	1.0	0.0	0.33	
4	0	tcp	http	SF	239	486	0	0	0	0	...	4	1.0	0.0	0.25	
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
498426	0	tcp	http	SF	212	2288	0	0	0	0	...	255	1.0	0.0	0.33	
498427	0	tcp	http	SF	219	236	0	0	0	0	...	255	1.0	0.0	0.25	
498428	0	tcp	http	SF	218	3610	0	0	0	0	...	255	1.0	0.0	0.20	
498429	0	tcp	http	SF	219	1234	0	0	0	0	...	255	1.0	0.0	0.17	
498430	0	tcp	http	SF	219	1098	0	0	0	0	...	255	1.0	0.0	0.14	

488431 rows × 42 columns

To be able to operate on the data, we had to encode 4 categorical features into numerical features. As for the protocol\_type, service and flag features, one hot encoding was applied. As for the target feature, label encoding was applied.

```
icmp      2833545
tcp       1870598
udp        194288
Name: protocol_type, dtype: int64
3
ecr_i      2811660
private    1100831
http       623091
smtp       96554
other      72653
...
tftp_u          3
harvest         2
aol             2
http_8001       2
http_2784       1
Name: service, Length: 70, dtype: int64
70
SF      3744328
S0      869829
REJ     268874
RSTR     8094
RSTO     5344
SH      1040
S1       532
S2       161
RSTOS0    122
OTH       57
S3        50
Name: flag, dtype: int64
11
```

The original protocol\_type, service and flag features were dropped, as they were replaced with the above one hot encoded features. One hot encoding was favored to avoid any unintentional biasing of the model calculations, as all one hot encoded features only take a 0 or 1 value. As for the categorical values of the target feature, they were replaced with the label encoded numerical ones. Label encoding was used on this feature to avoid an unnecessary increase in dimensionality. The original dataframe had 42 features, after dropping 3 features, and adding 84 features, the encoded dataframe has a total of 123 features.

encoded\_data = encode(training\_data)  
encoded\_data

	duration	src_bytes	dst_bytes	land	wrong_fragment	urgent	hot	num_failed_logins	logged_in	num_compromised	...	flag_REJ	flag_RSTO	flag_RSTOS0	flag_RSTR	flag_S0	flag_S1	flag_S2	flag_S3	flag_SF
0	0	215	45076	0	0	0	0	0	1	0	...	0	0	0	0	0	0	0	0	1
1	0	162	4528	0	0	0	0	0	1	0	...	0	0	0	0	0	0	0	0	1
2	0	236	1228	0	0	0	0	0	1	0	...	0	0	0	0	0	0	0	0	1
3	0	233	2032	0	0	0	0	0	1	0	...	0	0	0	0	0	0	0	0	1
4	0	239	486	0	0	0	0	0	1	0	...	0	0	0	0	0	0	0	0	1
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
488426	0	212	2288	0	0	0	0	0	1	0	...	0	0	0	0	0	0	0	0	1
488427	0	219	236	0	0	0	0	0	1	0	...	0	0	0	0	0	0	0	0	1
488428	0	218	3610	0	0	0	0	0	1	0	...	0	0	0	0	0	0	0	0	1
488429	0	219	1234	0	0	0	0	0	1	0	...	0	0	0	0	0	0	0	0	1
488430	0	219	1098	0	0	0	0	0	1	0	...	0	0	0	0	0	0	0	0	1

4893431 rows × 123 columns

Analysis of the target feature:

The target feature which contains the labels for each row of data, we are informed that our data is classified into a total of 23 classes, but not all classes are equally represented as shown in the figure. According to the paper "Reducing Redundancy in the KDD CUP 1999 Data Set" by Ali et al. (2011), the percentage of redundancies in the original KDD Cup 1999 dataset was found to be approximately 70%. This high percentage of redundancies can have a negative impact on the performance of machine learning models trained on the dataset. Although redundancies are important to consider, we had to remove them, and work on a small portion of the remaining data to be able to complete the clustering task.

smurf.	2807886
neptune.	1072017
normal.	972781
satan.	15892
ipsweep.	12481
portsweep.	10413
nmap.	2316
back.	2203
warezclient.	1020
teardrop.	979
pod.	264
guess_passwd.	53
buffer_overflow.	30
land.	21
warezmaster.	20
imap.	12
rootkit.	10
loadmodule.	9
ftp_write.	8
multihop.	7
phf.	4
perl.	3
spy.	2
Name: target, dtype: int64	



### *train\_test\_split:*

We used the sklearn function to get a training set containing 0.004 of the data (X\_train), and the target feature was used as the (Y\_train). However, we still had to remove redundancies in them. To be able to do this, and ensure that both their rows still correspond, they were combined in a single dataframe, and re-split into Xtrain\_unique and Ytrain\_unique. This leaves as with 7974 rows of data on which we apply the clustering.

```
#Combine to remove corresponding duplicates
combined = pd.concat([X_train, y_train], axis=1)
combined = combined.drop_duplicates()
combined
```

	duration	src_bytes	dst_bytes	land	wrong_fragment	urgent	hot	num_failed_logins	logged_in	num_compromised	...	flag_RSTO	flag_RSTO50	flag_RSTR	flag_S0	flag_S1	flag_S2	flag_S3	flag_SF	flag_SH	
2451599	0	1032	0	0		0	0	0	0	0	0 ...	0	0	0	0	0	0	0	0	1	0
999149	0	1032	0	0		0	0	0	0	0	0 ...	0	0	0	0	0	0	0	0	1	0
4477733	0	145	0	0		0	0	0	0	0	0 ...	0	0	0	0	0	0	0	0	1	0
654692	0	0	0	0		0	0	0	0	0	0 ...	0	0	0	1	0	0	0	0	0	0
3418860	0	45	115	0		0	0	0	0	0	0 ...	0	0	0	0	0	0	0	0	1	0
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
374455	0	415	10439	0		0	0	0	0	1	0 ...	0	0	0	0	0	0	0	0	1	0
3651134	0	0	0	0		0	0	0	0	0	0 ...	0	0	0	1	0	0	0	0	0	0
636658	0	0	0	0		0	0	0	0	0	0 ...	0	0	0	1	0	0	0	0	0	0
3872637	0	0	0	0		0	0	0	0	0	0 ...	0	0	0	1	0	0	0	0	0	0
3547162	0	0	0	0		0	0	0	0	0	0 ...	0	0	0	1	0	0	0	0	0	0

7974 rows × 123 columns

### *Setting the # of clusters to the # of existing classes in the remaining data:*

The remaining data consists of only 11 classes of the original 23. However, classes are still unequally represented, as the most represented class is represented almost equally to the sum of all other remaining classes.

9	3938
11	3758
18	123
17	57
5	47
15	27
0	9
10	6
20	4
21	4
14	1

## Normalized Cut Algorithm:

---

**Algorithm 16.1:** Spectral Clustering Algorithm

---

**SPECTRAL CLUSTERING ( $\mathbf{D}, k$ ):**

- 1 Compute the similarity matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$
  - 2 **if** *ratio cut* **then**  $\mathbf{B} \leftarrow \mathbf{L}$
  - 3 **else if** *normalized cut* **then**  $\mathbf{B} \leftarrow \mathbf{L}^s$  or  $\mathbf{L}^a$
  - 4 Solve  $\mathbf{B}\mathbf{u}_i = \lambda_i \mathbf{u}_i$  for  $i = n, \dots, n - k + 1$ , where  $\lambda_n \leq \lambda_{n-1} \leq \dots \leq \lambda_{n-k+1}$
  - 5  $\mathbf{U} \leftarrow (\mathbf{u}_n \quad \mathbf{u}_{n-1} \quad \dots \quad \mathbf{u}_{n-k+1})$
  - 6  $\mathbf{Y} \leftarrow$  normalize rows of  $\mathbf{U}$  using Eq. (16.23)
  - 7  $\mathcal{C} \leftarrow \{C_1, \dots, C_k\}$  via K-means on  $\mathbf{Y}$
- 

This pseudocode from the textbook was followed to apply spectral clustering:

The similarity matrix was computed using sklearn's `rbf_kernel` which computes the Gaussian kernel. The degree matrix was then computed, and the Laplacian matrix was then computed as the difference between the similarity and degree matrix. To apply normalized cut, we computed the asymmetric Laplacian matrix using the following equation:

$$\mathbf{L}^a = \Delta^{-1} \mathbf{L}$$

As this is an asymmetric matrix, eig was used to compute the eigenvalues and eigenvectors  $\mathbf{U}$ . Only the real components of  $\mathbf{U}$  were kept, and they were then normalized to return the unit vector  $\mathbf{Y}$ , on which Kmeans will be applied.

## Evaluation:

precision:

recall:

[0.9583741429970617,	[0.9939055358049771,
0.9518779342723005,	0.21580627993613624,
0.9905660377358491,	0.02794039382650346,
0.7651515151515151,	0.08062799361362427,
0.9578256794751641,	0.271953166577967,
0.9915966386554622,	0.03139968068121341,
0.9558011049723757,	0.04603512506652475,
0.9415384615384615,	0.08142629058009579,
0.8939051918735892,	0.10537519957424162,
0.9739130434782609,	0.05960617349654071,
0.9590643274853801]	0.043640234167110166]

F1 score = 0.23949682543756703

Conditional entropy: 0.3464060044294866

We can see that except for the precision, the evaluation scores are affected due to the low quality and noisiness of the data upon which spectral clustering was applied. After all, a model is only as good as the data it trains on.

### *Further analysis:*

We asked the spectral clustering algorithm to cluster the data into different numbers of clusters (7, 9, 13, 15).

For the numbers of clusters that are smaller than than the number of remaining classes (7, 9) in the data, an improved score of F1 and conditional entropy were recorded.

```
K = 7
F1 score = 0.35974169308547094
Conditional entropy: 0.4205720755700158
K = 9
F1 score = 0.29010667036097765
Conditional entropy: 0.3551763509243745
```

The opposite was recorded as the numbers of clusters exceeded the number of remaining classes (13, 15).

```
K = 13
F1 score = 0.20732136012072405
Conditional entropy: 0.32136218488940527
K = 15
F1 score = 0.18119829877015414
Conditional entropy: 0.30070067845164533
```

# DBSCAN

## I)Introduction)

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is a popular unsupervised machine learning algorithm used for clustering data points in a dataset. It is designed to group together data points that are close to each other based on their density, while identifying outliers or noise points that do not belong to any cluster.

The algorithm works by defining a neighborhood around each data point, based on a specified radius. The density of each point is then calculated based on the number of neighboring points within that radius. A point is considered a "core point" if it has at least a minimum number of neighboring points within that radius.

Starting with a core point, the algorithm expands the cluster by iteratively adding neighboring points that also meet the density requirement, until there are no more qualifying points to add. This process is repeated until all points have been assigned to a cluster, or marked as noise points.

The algorithm is flexible in that it can handle clusters of arbitrary shapes, and does not require the number of clusters to be specified beforehand. However, it does require the setting of two parameters: the radius of the neighborhood (epsilon) and the minimum number of neighboring points required for a point to be considered a core point (min\_samples).

Overall, DBSCAN is a powerful clustering algorithm that is widely used in a variety of applications such as image segmentation, anomaly detection, and customer segmentation in marketing.

Before using the algorithm, some data analysis was done on the data and the result was:

---

- a. Reduced the dimensions of the dataset from 42 to 30
  - b. Dropping all the redundancy from the dataset
  - c. Using only 0.005 from the data
-

## Pseudocode:

---

**Algorithm 15.1: Density-based Clustering Algorithm**

---

**DBSCAN ( $\mathbf{D}, \epsilon, minpts$ ):**

```
1  $Core \leftarrow \emptyset$ 
2 foreach  $\mathbf{x}_i \in \mathbf{D}$  do // Find the core points
3   Compute  $N_\epsilon(\mathbf{x}_i)$ 
4    $id(\mathbf{x}_i) \leftarrow \emptyset$  // cluster id for  $\mathbf{x}_i$ 
5   if  $N_\epsilon(\mathbf{x}_i) \geq minpts$  then  $Core \leftarrow Core \cup \{\mathbf{x}_i\}$ 
6  $k \leftarrow 0$  // cluster id
7 foreach  $\mathbf{x}_i \in Core$ , such that  $id(\mathbf{x}_i) = \emptyset$  do
8    $k \leftarrow k + 1$ 
9    $id(\mathbf{x}_i) \leftarrow k$  // assign  $\mathbf{x}_i$  to cluster id  $k$ 
10  DENSITYCONNECTED ( $\mathbf{x}_i, k$ )
11  $\mathcal{C} \leftarrow \{C_i\}_{i=1}^k$ , where  $C_i \leftarrow \{\mathbf{x} \in \mathbf{D} \mid id(\mathbf{x}) = i\}$ 
12  $Noise \leftarrow \{\mathbf{x} \in \mathbf{D} \mid id(\mathbf{x}) = \emptyset\}$ 
13  $Border \leftarrow \mathbf{D} \setminus \{Core \cup Noise\}$ 
14 return  $\mathcal{C}, Core, Border, Noise$ 
```

**DENSITYCONNECTED ( $\mathbf{x}, k$ ):**

```
15 foreach  $\mathbf{y} \in N_\epsilon(\mathbf{x})$  do
16    $id(\mathbf{y}) \leftarrow k$  // assign  $\mathbf{y}$  to cluster id  $k$ 
17   if  $\mathbf{y} \in Core$  then DENSITYCONNECTED ( $\mathbf{y}, k$ )
```

---

As seen in the pseudocode, We have 2 hyperparameters to tune which are epsilon and minpts.

We chose minpts according to the following equation:

$$minpts = D + 1$$

Where D is the number of dimensions in our dataset.

Now we need to get the epsilon value KNN Algorithm from Sklearn was used was a value of K equal to the number of features and we got 0.5 for epsilon.

Using those 2 hyperparameters we got the following results:

```
Precision= [0.8114334470989761, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0]
Recall= [0.40832975525976817, 0.765142150803461, 0.734375, 0.5021468441391155, 0.18149979398434282, 0.04229282954057535, 0.0098754830399313, 0.6438356
164383562, 0.01051953628166595, 0.012666380420781451, 0.00622584800343495, 0.00794332331472735]
F_Measure= 0.3492410127614332
Entropy= 0.25905530128112064
```

---

References:

- 1- [http://www.ccs.neu.edu/home/vip/teach/DMcourse/2\\_cluster\\_EM\\_mixed/notes\\_slides/revisitofrevisitDBSCAN.pdf](http://www.ccs.neu.edu/home/vip/teach/DMcourse/2_cluster_EM_mixed/notes_slides/revisitofrevisitDBSCAN.pdf)
- 2- <https://link.springer.com/article/10.1023%2FA%3A1009745219419>
- 3- <https://www.datanovia.com/en/lessons/dbscan-density-based-clustering-essentials/>
- 4- <https://towardsdatascience.com/machine-learning-clustering-dbscan-determine-the-optimal-value-for-epsilon-eps-python-example-3100091cfbc>
- 5- [https://dataminingbook.info/book\\_html/chap15/book.html](https://dataminingbook.info/book_html/chap15/book.html)