

# Data Warehousing and Data Mining

## Complex Engineering Problem

### Import important Libraries and Load dataset

The first step in dealing with any data science problem is to import the important data handling libraries i.e., NumPy and Pandas and then load the respective data set

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
```

```
In [2]: terrorism_dataset = pd.read_csv("globalterrorismdb_0718dist.csv",encoding = "ISO-8859-1")
```

### Examine Dataset

After loading the dataset, we must examine to know the total number of rows and columns and the amount of NaN values associated with in each column

```
In [4]: terrorism_dataset
```

```
Out[4]:
```

	eventid	iy	imonth	iday	approxdate	extended	resolution	country	country_txt	region	...	addnotes	scite1	scite2	scite3
0	197000000001	1970	7	2	NaN	0	NaN	58	Dominican Republic	2	...	NaN	NaN	NaN	NaN
1	197000000002	1970	0	0	NaN	0	NaN	130	Mexico	1	...	NaN	NaN	NaN	NaN
2	197001000001	1970	1	0	NaN	0	NaN	160	Philippines	5	...	NaN	NaN	NaN	NaN
3	197001000002	1970	1	0	NaN	0	NaN	78	Greece	8	...	NaN	NaN	NaN	NaN
4	197001000003	1970	1	0	NaN	0	NaN	101	Japan	4	...	NaN	NaN	NaN	NaN
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
181686	201712310022	2017	12	31	NaN	0	NaN	182	Somalia	11	...	NaN	"Somalia: Al-Shabaab Militants Attack Army Che...	"Highlights: Somalia Daily Media Highlights 2 ...	"Highlights: Somalia Daily Media Highlights 1 ...
181687	201712310029	2017	12	31	NaN	0	NaN	200	Syria	10	...	NaN	"Putin's 'victory' in Syria has turned into a ...	"Two Russian soldiers killed at Hmeymim base i...	"Two Russian servicemen killed in Syria mortar...
181688	201712310030	2017	12	31	NaN	0	NaN	160	Philippines	5	...	NaN	"Maguindanao clashes trap tribe members," Phil...	NaN	NaN
181689	201712310031	2017	12	31	NaN	0	NaN	92	India	6	...	NaN	"Trader escapes grenade attack in Imphal," Bus...	NaN	NaN
181690	201712310032	2017	12	31	NaN	0	NaN	160	Philippines	5	...	NaN	"Security tightened in Cotabato following IED ...	"Security tightened in Cotabato City," Manila ...	NaN

181691 rows x 135 columns

## Data Cleaning

Data cleaning or Data cleansing is very important from the perspective of building intelligent automated systems. Data cleansing is a preprocessing step that improves the data validity, accuracy, completeness, consistency, and uniformity. It is essential for building reliable machine learning models that can produce good results. Otherwise, no matter how good the model is, its results cannot be trusted. In short, data cleaning means fixing bad data in your data set. Bad data could be:

1. Empty cells
2. Data in wrong format
3. Wrong data
4. Duplicates

So initially we apply a loop to filter out the columns whose 1/3<sup>rd</sup> values are NaN

```
td2=terrorism_dataset
for i in terrorism_dataset.columns:
    x =terrorism_dataset[f"{i}"].isna().sum()
    if(x>50000):
        td2 = td2.drop(f"{i}",axis=1)
td2
```

Secondly, we remove repetitive columns from our data set i.e. In our data set we have one column with title country and other with country\_txt both provide same information however one is numeric and other is alphabetic so we will get rid of one

```
td2.drop("attacktype1_txt",axis=1,inplace=True)
td2.drop("targtype1_txt",axis=1,inplace=True)
td2.drop("targsubtype1_txt",axis=1,inplace=True)
td2.drop("natlty1_txt",axis=1,inplace=True)
td2.drop("region_txt",axis=1,inplace=True)
td2.drop("country_txt",axis=1,inplace=True)
td2.drop("weaptype1_txt",axis=1,inplace=True)
td2.drop("weapsubtype1_txt",axis=1,inplace=True)
```

Thirdly we will remove columns that are not giving much information regarding our defined task after analysis of dataset

```

td2.drop("INT_LOG",axis=1,inplace=True)
td2.drop("INT_IDEO",axis=1,inplace=True)
td2.drop("INT_MISC",axis=1,inplace=True)
td2.drop("INT_ANY",axis=1,inplace=True)
td2.drop("dbsource",axis=1,inplace=True)
td2.drop("corp1",axis=1,inplace=True)
td2.drop("ishostkid",axis=1,inplace=True)
td2.dropna(subset=["doubtterr"],axis=0 , inplace=True)
td2.dropna(subset=["multiple"],axis=0 , inplace=True)
td2.dropna(subset=["specificity"],axis=0 , inplace=True)

td2.drop("guncertain1",axis=1,inplace=True)
td2.drop("longitude",axis=1,inplace=True)
td2.drop("latitude",axis=1,inplace=True)
td2.drop("targsubtype1",axis=1,inplace=True)
td2.drop("weapsubtype1",axis=1,inplace=True)

```

Our data consist upon more than 180000 rows hence getting rid of few row that hold NaN values will not affect the data set much hence we will remove such rows too

```

td2.dropna(subset=["provstate"],axis=0 , inplace=True)
td2.dropna(subset=["city"],axis=0 , inplace=True)

```

Lastly, we will fill NaN values of some important columns calculating their mean or mode or hardcode some value based on our analysis

```

td2["nwound"].fillna(3,inplace=True)
td2["nkill"].fillna(2,inplace=True)
td2.natlty1.fillna(0,inplace=True)
td2.target1.fillna("Unknown",inplace=True)

```

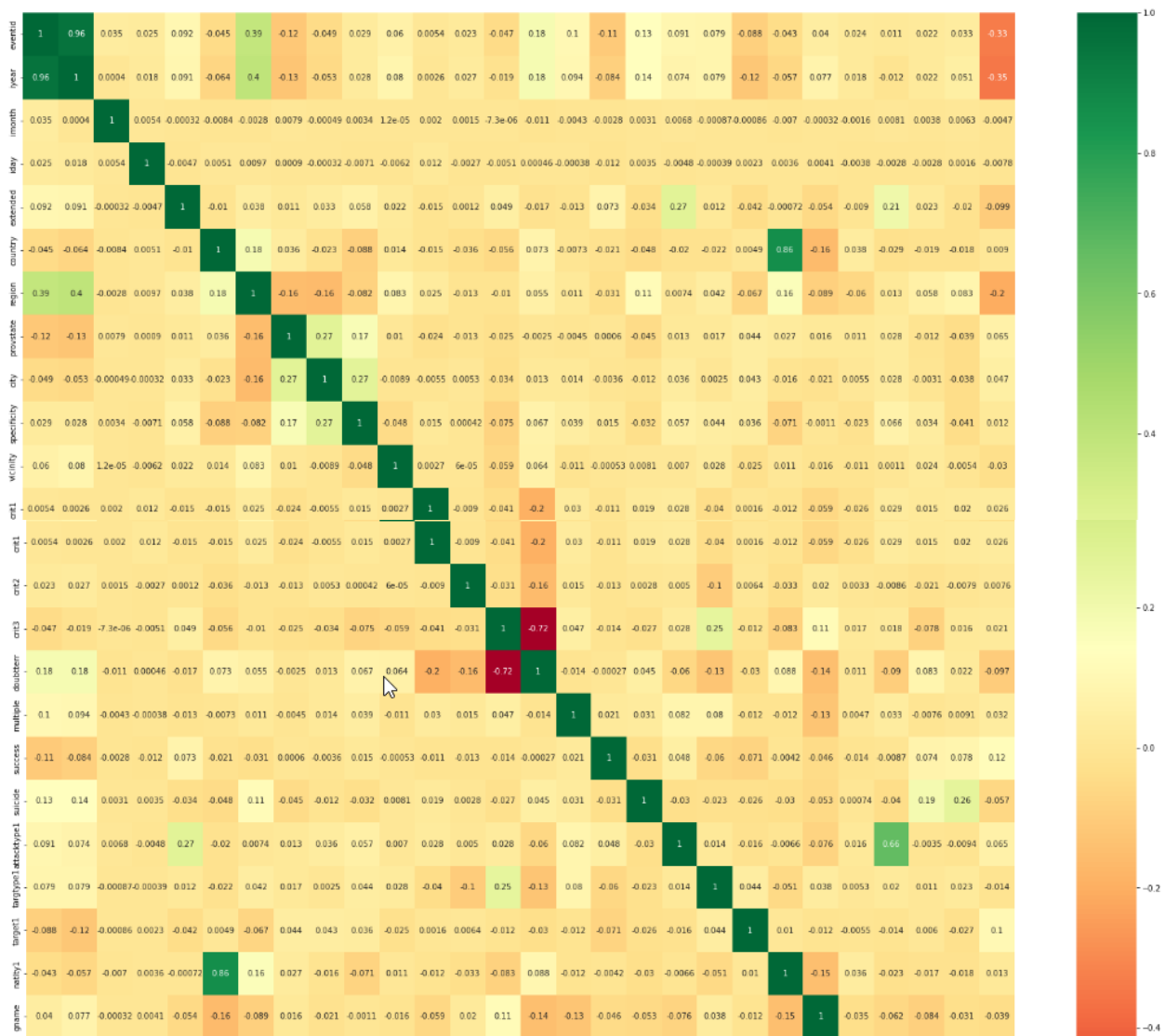
## Apply Label Encoder

Label Encoding refers to converting the labels into a numeric form so as to convert them into the machine-readable form. Machine learning algorithms can then decide in a better way how those labels must be operated. It is an important pre-processing step for the structured dataset in supervised learning.

```
from sklearn.preprocessing import LabelEncoder
label_encoder = LabelEncoder()
encode_result = td2.apply(label_encoder.fit_transform)
```

## Feature Selection

Correlation states how the features are related to each other or the target variable. Correlation can be positive (increase in one value of feature increases the value of the target variable) or negative (increase in one value of feature decreases the value of the target variable) Heatmap makes it easy to identify which features are most related to the target variable.



## Apply Machine Learning Model

Firstly, we define x with set of features that will be input for our model and y with the output feature, in our case we are supposed to predict the attack type therefore we will put that in y and some other selected features in x

```
x=encode_result[["iyear","country","region","city","provstate","gname","target1","weaptype1"]]  
y=encode_result["attacktype1"]
```

Secondly, we Apply train\_test\_split to x and y The train-test split is used to estimate the performance of machine learning algorithms that are applicable for prediction-based Algorithms/Applications. This method is a fast and easy procedure to perform such that we can compare our own machine learning model results to machine results.

```
from sklearn.model_selection import train_test_split  
np.random.seed(42)  
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.2)
```

Lastly, we setup ML model and train it and see the score to know its accuracy

```
In [13]: # Setup model  
from sklearn.ensemble import RandomForestClassifier  
model = RandomForestClassifier()  
  
model.fit(x_train, y_train)  
model.score(x_test, y_test)
```

```
Out[13]: 0.8342642260686833
```

## Confusion Matrix

A confusion matrix is a table that is used to define the performance of a classification algorithm. A confusion matrix visualizes and summarizes the performance of a classification algorithm.

```

from sklearn.metrics import confusion_matrix
rclassifier_pred = model.predict(x_test)
a=confusion_matrix(y_test,rclassifier_pred)
a1=a.flatten()
x=a1[0:4]
print(x)

```

```
[1871 1100  565    2]
```

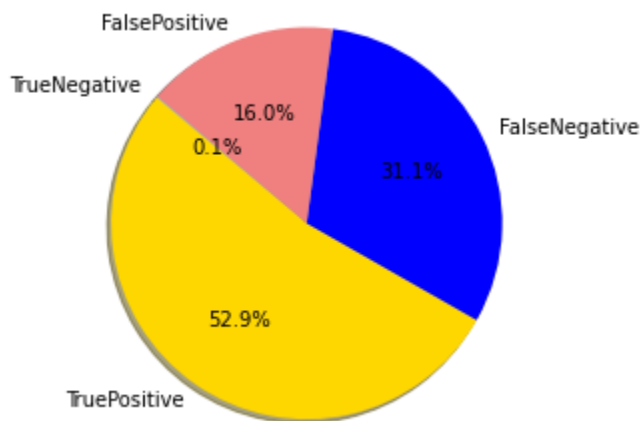
```

labels = 'TruePositive', 'FalseNegative', 'FalsePositive', 'TrueNegative'
colors = ['gold', 'blue', 'lightcoral', 'lightskyblue']

# Plot
plt.pie(x, labels=labels, colors=colors,
autopct='%1.1f%%', shadow=True, startangle=140)

plt.axis('equal')
plt.show()

```



## Spatial-Temporal Analysis Using Clustering

### Feature scaling using MinMaxScaler.

Feature scaling through standardization can be an important preprocessing step for many machine learning algorithms. Standardization involves rescaling the features such that they have the properties of a standard normal distribution with a mean of zero and a standard deviation of one.

The MinMaxScaler is a type of scaler that scales the minimum and maximum values to be 0 and 1 respectively. While the StandardScaler scales all values between min and max so that they fall within a range from min to max.

```
from sklearn.preprocessing import normalize
from sklearn.preprocessing import MinMaxScaler
scalar = MinMaxScaler()

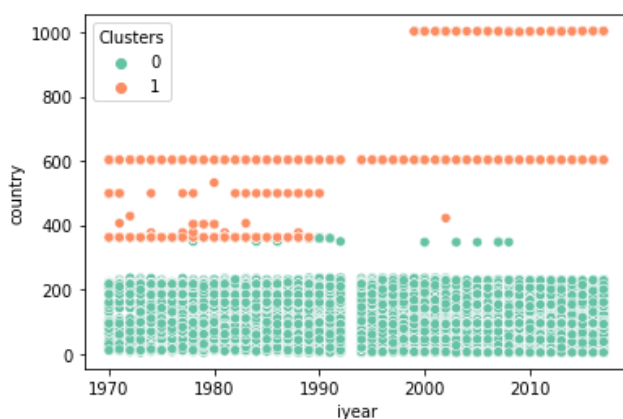
scale = scalar.fit_transform(td2[['iyear', 'country']])
df_scale = pd.DataFrame(scale, columns = ['iyear', 'country']);
df_scale.head(5)
```

## K-Mean Algorithm

K-means Algorithm is an Iterative algorithm that divides a group of n datasets into k subgroups or clusters based on the similarity and their mean distance from the centroid of that subgroup formed.

```
: from sklearn.cluster import KMeans
import sklearn.cluster as cluster
km=KMeans(n_clusters=2)
y_predicted = km.fit_predict(td2[['iyear', 'region', 'attacktype1', 'country', 'imonth', 'iday', 'nkill', 'nwound']])
td2['Clusters'] = km.labels_
sns.scatterplot(x="iyear", y="country", hue = 'Clusters',
data=td2, palette='Set2')
```

```
: <AxesSubplot:xlabel='iyear', ylabel='country'>
```



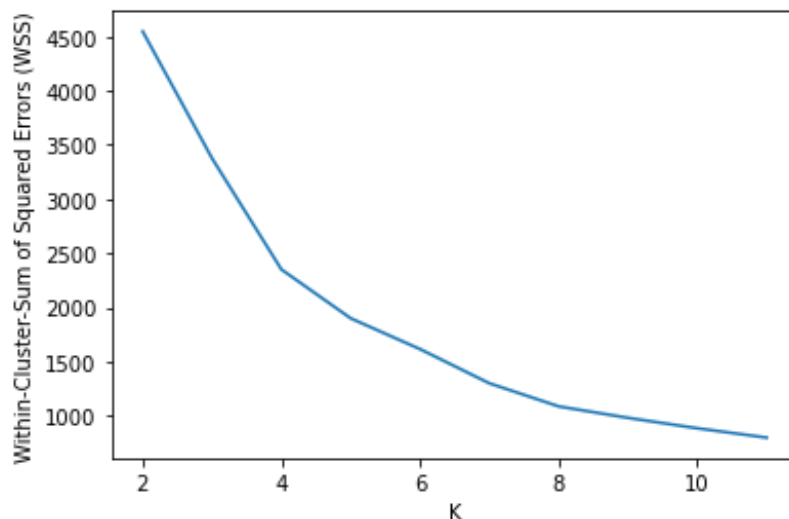
## Optimum number of Clusters in K Means

Imagine we set  $k$  to its maximum value  $n$  (where  $n$  is number of samples) each sample will form its own cluster meaning sum of squared distances equals zero. Below is a plot of sum of squared distances for  $k$  in the range specified above. If the plot looks like an arm, then the elbow on the arm is optimal  $k$ .

```
#Finding optimum value of K
K=range(2,12)
wss = []
for k in K:
    kmeans=cluster.KMeans(n_clusters=k)
    kmeans=kmeans.fit(df_scale)
    wss_iter = kmeans.inertia_
    wss.append(wss_iter)

import matplotlib.pyplot as plt
#Plotting the graph
plt.xlabel('K')
plt.ylabel('Within-Cluster-Sum of Squared Errors (WSS)')
plt.plot(K,wss)
```

[<matplotlib.lines.Line2D at 0x17c869fc7c0>]

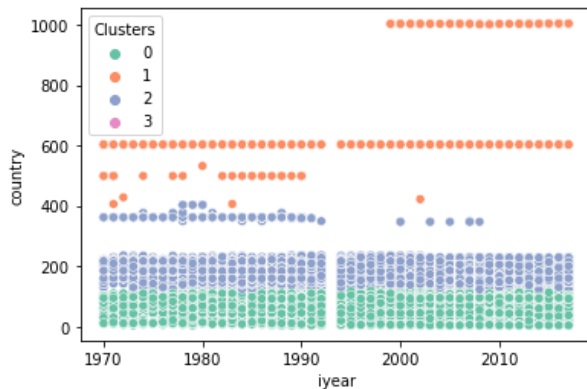




## Applying Optimum value of K for Clustering

```
km=KMeans(n_clusters=4)
y_predicted = km.fit_predict(td2[['iyear','region','attacktype1','country','imonth','iday','nkill','nwound']])
y_predicted
td2['Clusters'] = km.labels_
sns.scatterplot(x="iyear", y="country",hue = 'Clusters',
data=td2,palette='Set2')
```

<AxesSubplot:xlabel='iyear', ylabel='country'>



## Description of Graph

The Graph show clustering based on Space and time with time along x axis and country code along y axis the list pf countries along with their codes has been attached below.

4 = Afghanistan	22 = Belize	41 = Central African Republic
5 = Albania	23 = Benin	42 = Chad
6 = Algeria	24 = Bermuda*	43 = Chile
7 = Andorra	25 = Bhutan	44 = China
8 = Angola	26 = Bolivia	45 = Colombia
10 = Antigua and Barbuda	28 = Bosnia-Herzegovina	46 = Comoros
11 = Argentina	29 = Botswana	47 = Republic of the Congo
12 = Armenia	30 = Brazil	49 = Costa Rica
14 = Australia	31 = Brunei	50 = Croatia
15 = Austria	32 = Bulgaria	51 = Cuba
16 = Azerbaijan	33 = Burkina Faso	53 = Cyprus
17 = Bahamas	34 = Burundi	54 = Czech Republic
18 = Bahrain	35 = Belarus	55 = Denmark
19 = Bangladesh	36 = Cambodia	56 = Djibouti
20 = Barbados	37 = Cameroon	57 = Dominica
21 = Belgium	38 = Canada	58 = Dominican Republic

59 = Ecuador	106 = Kuwait	157 = Papua New Guinea		
60 = Egypt	107 = Kyrgyzstan	158 = Paraguay		
61 = El Salvador	108 = Laos	159 = Peru		
62 = Equatorial Guinea	109 = Latvia	160 = Philippines		
63 = Eritrea	110 = Lebanon	161 = Poland		
64 = Estonia	111 = Lesotho	162 = Portugal		
65 = Ethiopia	112 = Liberia	163 = Puerto Rico*		
66 = Falkland Islands	113 = Libya	164 = Qatar		
67 = Fiji	114 = Liechtenstein*	166 = Romania		
68 = Finland	115 = Lithuania	167 = Russia		
69 = France	116 = Luxembourg	168 = Rwanda		
70 = French Guiana	117 = Macau	169 = Saba (Netherlands		
71 = French Polynesia	118 = Macedonia	Antilles)*		
72 = Gabon	119 = Madagascar	173 = Saudi Arabia	209 = Turkey	229 = Democratic Republic of the
73 = Gambia	120 = Malawi	174 = Senegal	210 = Turkmenistan	Congo
74 = Georgia	121 = Malaysia	175 = Serbia-Montenegro	212 = Tuvalu*	230 = Zambia
75 = Germany	122 = Maldives	176 = Seychelles	213 = Uganda	231 = Zimbabwe
76 = Ghana	123 = Mali	177 = Sierra Leone	214 = Ukraine	233 = Northern Ireland*
78 = Greece	124 = Malta	178 = Singapore	215 = United Arab Emirates	235 = Yugoslavia
79 = Greenland*	125 = Man, Isle of*	179 = Slovak Republic	216 = Great Britain*	236 = Czechoslovakia
80 = Grenada	126 = Marshall Islands*	180 = Slovenia	217 = United States	238 = Corsica*
81 = Guadeloupe	127 = Martinique	181 = Solomon Islands	218 = Uruguay	334 = Asian*
83 = Guatemala	128 = Mauritania	182 = Somalia	219 = Uzbekistan	347 = East Timor
84 = Guinea	129 = Mauritius	183 = South Africa	220 = Vanuatu	349 = Western Sahara
85 = Guinea-Bissau	130 = Mexico	184 = South Korea	221 = Vatican City	351 = Commonwealth of
86 = Guyana	132 = Moldova	185 = Spain	222 = Venezuela	Independent States*
87 = Haiti	134 = Mongolia*	186 = Sri Lanka	223 = Vietnam	359 = Soviet Union
88 = Honduras	136 = Morocco	189 = St. Kitts and Nevis	225 = Virgin Islands (U.S.)*	362 = West Germany (FRG)
89 = Hong Kong	137 = Mozambique	190 = St. Lucia	226 = Wallis and Futuna	377 = North Yemen
90 = Hungary	138 = Myanmar	192 = St. Martin*	228 = Yemen	403 = Rhodesia
91 = Iceland	139 = Namibia	195 = Sudan		
92 = India	141 = Nepal	196 = Suriname		
93 = Indonesia	142 = Netherlands	197 = Swaziland		
94 = Iran	143 = New Caledonia	198 = Sweden		
95 = Iraq	144 = New Zealand	199 = Switzerland		
96 = Ireland	145 = Nicaragua	200 = Syria		
97 = Israel	146 = Niger	201 = Taiwan		
98 = Italy	147 = Nigeria	202 = Tajikistan		
99 = Ivory Coast	149 = North Korea	203 = Tanzania		
100 = Jamaica	151 = Norway	204 = Togo		
101 = Japan	152 = Oman*	205 = Thailand		
102 = Jordan	153 = Pakistan	206 = Tonga*		
103 = Kazakhstan	155 = West Bank and Gaza Strip	207 = Trinidad and Tobago		
104 = Kenya	156 = Panama	208 = Tunisia		
				406 = South Yemen
				422 = International
				428 = South Vietnam
				499 = East Germany (GDR)
				520 = Sinhalese*
				532 = New Hebrides
				603 = United Kingdom
				604 = Zaire
				605 = People's Republic of the
				Congo
				999 = Multinational*
				1001 = Serbia
				1002 = Montenegro
				1003 = Kosovo
				1004 = South Sudan

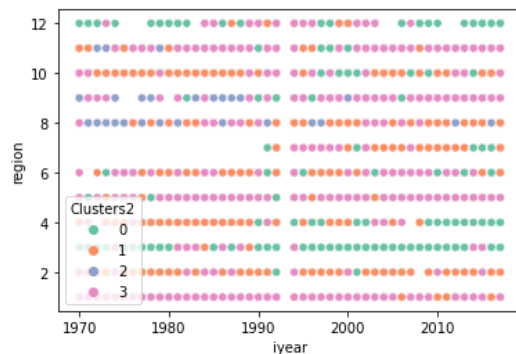
20

## More clustering on optimal k values for spatial temporal analysis

```
km=KMeans(n_clusters=4)
y_predicted = km.fit_predict(td2[['iyear','region','attacktype1','country','imonth','iday']])
td2['Clusters2'] = km.labels_
sns.scatterplot(x="iyear", y="region", hue = 'Clusters2',
data=td2,palette='Set2')
```

```
<AxesSubplot:xlabel='iyear', ylabel='region'>
```

C:\Users\Dell\sample\_project\_1\env\lib\site-packages\IPython\core\pylabtools.py:151: UserWarning  
t" can be slow with large amounts of data.  
fig.canvas.print\_figure(bytes\_io, \*\*kw)



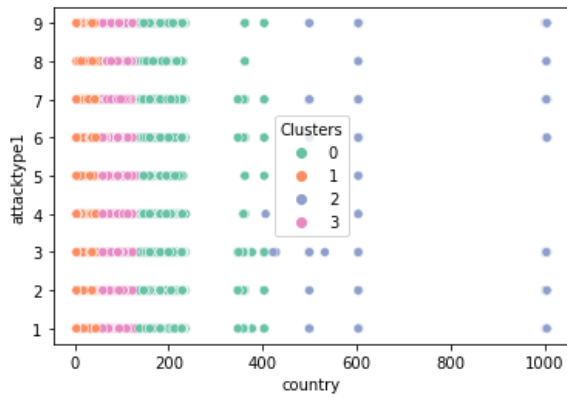
```

km=KMeans(n_clusters=4)
y_predicted = km.fit_predict(td2[['iyear', 'region', 'attacktype1', 'country', 'imonth', 'iday', 'nkill', 'nwound']])
td2['Clusters'] = km.labels_
sns.scatterplot(x="country", y="attacktype1", hue = 'Clusters',
data=td2,palette='Set2')

```

<AxesSubplot:xlabel='country', ylabel='attacktype1'>

C:\Users\Dell\sample\_project\_1\env\lib\site-packages\IPython\core\pylabtools.py:151: UserWarning: Creating legend t" can be slow with large amounts of data.  
fig.canvas.print\_figure(bytes\_io, \*\*kw)

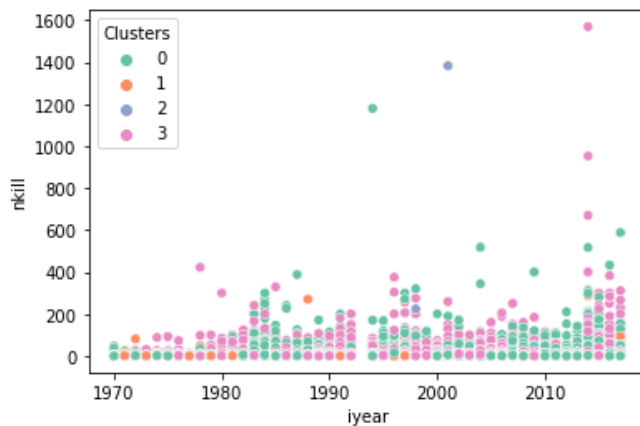


```

: km=KMeans(n_clusters=4)
y_predicted = km.fit_predict(td2[['iyear', 'region', 'attacktype1', 'country', 'imonth', 'iday', 'nkill', 'nwound']])
td2['Clusters'] = km.labels_
sns.scatterplot(x="iyear", y="nkill", hue = 'Clusters',
data=td2,palette='Set2')

```

: <AxesSubplot:xlabel='iyear', ylabel='nkill'>



```

km=KMeans(n_clusters=4)
y_predicted = km.fit_predict(td2[['iyear','region','attacktype1','country','imonth','iday','nkill','nwound']])
td2['Clusters'] = km.labels_
sns.scatterplot(x="region", y="nkill", hue = 'Clusters',
data=td2,palette='Set2')

```

<AxesSubplot:xlabel='region', ylabel='nkill'>

