

# Task 4: Exploratory Data Analysis on Dataset - Terrorism

Author: Mitul D Shinde

```
##
In [1]:
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

In [2]:
data = pd.read_csv('gt.csv')

C:\Users\Admin\anaconda3\lib\site-packages\Python\cores\interactiveshell.py:3165: DtypeWarning: Columns (4,6,3,1,3,61,62,63,76,79,90,92,94,96,114,115,121) have mixed types.Specify dtype option on import or set low memory=False.
has_escaped = await self.run_ast_nodes(code_ast.body, cell_name,

In [3]:
data.head()

Out[3]:
   eventid  year  month  iday  approximate  extended  resolution  country  country_txt  region  ...  addnotes  scite1  scite2  scite3  db
0  197000000001  1970  7  2  NaN  0  NaN  58  Dominican Republic  2  ...  NaN  NaN  NaN  NaN  NaN
1  197000000002  1970  0  0  NaN  0  NaN  130  Mexico  1  ...  NaN  NaN  NaN  NaN  NaN
2  197001000001  1970  1  0  NaN  0  NaN  78  Philippines  5  ...  NaN  NaN  NaN  NaN  NaN
3  197001000002  1970  1  0  NaN  0  NaN  160  Greece  8  ...  NaN  NaN  NaN  NaN  NaN
4  197001000003  1970  1  0  NaN  0  NaN  101  Japan  4  ...  NaN  NaN  NaN  NaN  NaN

5 rows x 135 columns

In [4]:
data.tail()

Out[4]:
   eventid  year  month  iday  approximate  extended  resolution  country  country_txt  region  ...  addnotes  scite1
181686  201712310022  2017  12  31  NaN  0  NaN  182  Somalia  11  ...  NaN
181687  201712310023  2017  12  31  NaN  0  NaN  200  Syria  10  ...  NaN
181688  201712310030  2017  12  31  NaN  0  NaN  160  Philippines  5  ...  NaN
181689  201712310031  2017  12  31  NaN  0  NaN  92  India  6  ...  NaN
181690  201712310032  2017  12  31  NaN  0  NaN  160  Philippines  5  ...  NaN

5 rows x 135 columns

In [5]:
data.shape

Out[5]:
(181691, 135)

In [6]:
data.describe()

Out[6]:
   eventid  year  month  iday  extended  country  region  latitude  longitude
mean  2002705e+11  2002.638997  6.46277  15.05644  0.04346  7.16098  177135.00000  -53.15463  -4.58695e+02
std  1325957e+09  13.259440  3.38833  8.814045  0.20063  112.414635  2.93408  18.569242  2.047790e+01
min  1970000e+11  1970.000000  0.000000  0.000000  0.000000  4.000000  1.000000  -53.15463  -8.618590e+07
25%  1991022e+11  1991.000000  4.000000  8.000000  0.000000  78.000000  5.000000  11.510046  4.54540e+00
50%  2009102e+11  2009.000000  6.000000  15.000000  0.000000  98.000000  6.000000  31.467463  4.344651e+01
75%  2014081e+11  2014.000000  9.000000  23.000000  0.000000  160.000000  10.000000  34.685087  6.871034e+01
max  2017123e+11  2017.000000  12.000000  31.000000  1.000000  1004.000000  12.000000  74.633553  1.793667e+02

8 rows x 77 columns

In [7]:
data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 181691 entries, 0 to 181690
Columns: 135 entries, eventid to related
dtypes: float(41), int64(42), object(58)
memory usage: 187.1 MB

In [8]:
data.columns

Out[8]:
Index(['eventid', 'year', 'month', 'iday', 'approximate', 'extended',
      'resolution', 'country', 'country_txt', 'region',
      'addnotes', 'scite1', 'scite2', 'scite3', 'dsource', 'INT_LOG',
      'INT_IDBO', 'INT_MISC', 'INT_ANY', 'related'],
      dtype='object', length=135)

In [9]:
data.value_counts()

Out[9]:
bound method DataFrame.value_counts of
   eventid  year  month  iday  approximate  extended  resoluti
0  197000000001  1970  7  2  NaN  0  NaN
1  197000000002  1970  0  0  NaN  0  NaN
2  197001000001  1970  1  0  NaN  0  NaN
3  197001000002  1970  1  0  NaN  0  NaN
4  197001000003  1970  1  0  NaN  0  NaN
181686  201712310022  2017  12  31  NaN  0  NaN
181687  201712310023  2017  12  31  NaN  0  NaN
181688  201712310030  2017  12  31  NaN  0  NaN
181689  201712310031  2017  12  31  NaN  0  NaN
181690  201712310032  2017  12  31  NaN  0  NaN

country  country_txt  region  \
0  58  Dominican Republic  2  ...  addnotes  \
1  130  Mexico  1  ...  NaN
2  160  Philippines  5  ...  NaN
3  78  Greece  8  ...  NaN
4  101  Japan  4  ...  NaN
...  ...  ...  ...
181686  182  Somalia  11  ...  NaN
181687  200  Syria  10  ...  NaN
181688  160  Philippines  5  ...  NaN
181689  92  India  6  ...  NaN
181690  160  Philippines  5  ...  NaN

scite1  \
0  NaN
1  NaN
2  NaN
3  NaN
4  NaN

...
181686  "Somalia: Al-Shabab Militants Attack Army Camp..."
181687  "Putin's 'Victory' in Syria has turned into ..."
181688  "Maguindanao clashes trap tribe members," Phil...
181689  "Trader escapes grenade attack in Imphal," Bus...
181690  "Security tightened in Cotabato following IED ..."

scite2  \
0  NaN
1  NaN
2  NaN
3  NaN
4  NaN

...
181686  "Highlights: Somalia Daily Media Highlights 2 ..."
181687  "Two Russian soldiers killed at Rameyn base ..."
181688  "Maguindanao clashes trap tribe members," Phil...
181689  "Security tightened in Cotabato City," Manila ...

scite3  \
0  NaN
1  NaN
2  NaN
3  NaN
4  NaN

...
181686  "Highlights: Somalia Daily Media Highlights 1 ..."
181687  "Two Russian servicemen killed in Syria mortar..."
181688  "Maguindanao clashes trap tribe members," Phil...
181689  "Security tightened in Cotabato City," Manila ...

dsource  INT_LOG  INT_IDBO  INT_MISC  INT_ANY  related
0  PGIS  0  0  0  0  NaN
1  PGIS  -9  -9  1  1  NaN
2  PGIS  -9  -9  1  1  NaN
3  PGIS  -9  -9  1  1  NaN
4  PGIS  -9  -9  1  1  NaN

...
181686  START Primary Collection  0  0  0  0  NaN
181687  START Primary Collection  -9  -9  1  1  NaN
181688  START Primary Collection  0  0  0  0  NaN
181689  START Primary Collection  -9  -9  0  -9  NaN
181690  START Primary Collection  -9  -9  0  -9  NaN

(181691 rows x 135 columns)

In [10]:
data.isnull().sum()

Out[10]:
eventid      0
year          0
month         0
iday          0
approximate  172452
INT_LOG      0
INT_IDBO     0
INT_MISC     0
INT_ANY      0
related      156653
Length: 135, dtype: int64

In [11]:
data.isnull().any()

Out[11]:
eventid      False
year          False
month         False
iday          True
approximate   True

INT_LOG      False
INT_IDBO     False
INT_MISC     False
INT_ANY      False
related      True
Length: 135, dtype: bool

In [12]:
data.count()

Out[12]:
eventid      181691
year         181691
month        181691
iday         181691
approximate   9239
INT_LOG      181691
INT_IDBO     181691
INT_MISC     181691
INT_ANY      181691
related      25038
Length: 135, dtype: int64

In [13]:
data.country_txt.value_counts()

Out[13]:
Iraq          2436
Pakistan      1468
Afghanistan   1273
India         1190
Colombia      876
...
North Korea  1
New Hebrides 1
International 1
South Vietnam 1
El Salvador 1
Name: country_txt, Length: 205, dtype: int64

In [14]:
c = data.country_txt.value_counts()
print(c)

Out[14]:
eventid  year  month  iday  extended  country  \
0  1.000000  0.999996  0.002706  0.018354  0.091761  -0.135039
year  0.999996  1.000000  0.000139  0.018254  0.091754  -0.135023
month  0.002706  0.000139  1.000000  0.005497  0.000468  -0.006305
iday  0.018354  0.018254  0.005497  1.000000  0.000468  -0.006305
extended  0.091761  0.091754  0.000468  0.000468  1.000000  -0.020466
...
released -0.181612 -0.181566 -0.011535 0.007168 0.192155 -0.044331
INT_LOG -0.143600 -0.143601 -0.002302 -0.001540 0.071768 0.067564
INT_IDBO -0.133252 -0.133253 -0.002034 -0.001621 0.075147 0.067564
INT_MISC -0.077852 -0.077847 -0.002554 -0.002027 0.027335 0.027281
INT_ANY -0.175605 -0.175596 -0.006336 -0.001199 0.080767 0.153118

region  latitude  longitude  specificity  ...  ransomom  \
0  0.401371  0.166886  0.003907  0.030641  ...  -0.009990
year  0.401384  0.166933  0.003917  0.030626  ...  -0.009984
month  -0.002999 -0.015978 -0.003880  0.030621  ...  -0.009984
iday  0.009710  0.003423 -0.002285 -0.006991  ...  0.012755
extended  0.038389 -0.024749  0.000523  0.057897  ...  -0.008114
...
released -0.149511  0.002790 -0.017745 -0.030631  ...  0.054571
INT_LOG -0.082584 -0.099727  0.002272  0.073022  ...  0.035821
INT_IDBO -0.071917 -0.094470  0.002268  0.071333  ...  0.039053
INT_MISC  0.043139  0.097652  0.000371 -0.019197  ...  0.023815
INT_ANY -0.047900 -0.041530  0.002497  0.061389  ...  0.028054

ransomom  ransompaid  ransompaidus  hostkidoutcome  released  \
0 -0.018001 -0.014094 -0.163422  0.256113 -0.181612
year -0.018216 -0.014038 -0.163373  0.256092 -0.181556
month 0.046989  0.058878  0.003146  0.012295 -0.001765
iday 0.010502  0.003148  0.000581 -0.006706  0.001765
extended 0.028177  0.001966  0.009367  0.232933 -0.192155
...
released 0.034843  0.049322  0.016832 -0.555478  1.000000
INT_LOG 0.031079  0.007029 -0.045504 -0.015442  0.039388
INT_IDBO 0.041983  0.013162 -0.039944 -0.016234  0.040947
INT_MISC 0.125162  0.037227  0.129274 -0.119776  0.085055
INT_ANY 0.053484  0.007275  0.056438 -0.061946  0.064759

INT_LOG  INT_IDBO  INT_MISC  INT_ANY
0 -0.143600 -0.133252 -0.077852 -0.175605
year -0.143601 -0.133253 -0.077847 -0.175596
month -0.002302 -0.002034 -0.002554 -0.001199
iday -0.001540 -0.001621 -0.002027 -0.001336
extended 0.071768  0.075147  0.027335  0.080767
...
released 0.039388  0.040947  0.085055  0.064759
INT_LOG 1.000000  0.996211  0.052537  0.891051
INT_IDBO 0.996211  1.000000  0.082014  0.893811
INT_MISC 0.052537  0.082014  1.000000  0.252193
INT_ANY 0.891051  0.893811  0.252193  1.000000

(77 rows x 77 columns)

In [15]:
c = data.country_txt.value_counts()
print(c)

Out[15]:
eventid  year  month  iday  extended  country  \
0  1.000000  0.999996  0.002706  0.018354  0.091761  -0.135039
year  0.999996  1.000000  0.000139  0.018254  0.091754  -0.135023
month  0.002706  0.000139  1.000000  0.005497  0.000468  -0.006305
iday  0.018354  0.018254  0.005497  1.000000  0.000468  -0.006305
extended  0.091761  0.091754  0.000468  0.000468  1.000000  -0.020466
...
released -0.181612 -0.181566 -0.011535 0.007168 0.192155 -0.044331
INT_LOG -0.143600 -0.143601 -0.002302 -0.001540 0.071768 0.067564
INT_IDBO -0.133252 -0.133253 -0.002034 -0.001621 0.075147 0.067564
INT_MISC -0.077852 -0.077847 -0.002554 -0.002027 0.027335 0.027281
INT_ANY -0.175605 -0.175596 -0.006336 -0.001199 0.080767 0.153118

region  latitude  longitude  specificity  ...  ransomom  \
0  0.401371  0.166886  0.003907  0.030641  ...  -0.009990
year  0.401384  0.166933  0.003917  0.030626  ...  -0.009984
month  -0.002999 -0.015978 -0.003880  0.030621  ...  -0.009984
iday  0.009710  0.003423 -0.002285 -0.006991  ...  0.012755
extended  0.038389 -0.024749  0.000523  0.057897  ...  -0.008114
...
released -0.149511  0.002790 -0.017745 -0.030631  ...  0.054571
INT_LOG -0.082584 -0.099727  0.002272  0.073022  ...  0.035821
INT_IDBO -0.071917 -0.094470  0.002268  0.071333  ...  0.039053
INT_MISC 0.043139  0.097652  0.000371 -0.019197  ...  0.023815
INT_ANY -0.047900 -0.041530  0.002497  0.061389  ...  0.028054

ransomom  ransompaid  ransompaidus  hostkidoutcome  released  \
0 -0.018001 -0.014094 -0.163422  0.256113 -0.181612
year -0.018216 -0.014038 -0.163373  0.256092 -0.181556
month 0.046989  0.058878  0.003146  0.012295 -0.001765
iday 0.010502  0.003148  0.000581 -0.006706  0.001765
extended 0.028177  0.001966  0.009367  0.232933 -0.192155
...
released 0.034843  0.049322  0.016832 -0.555478  1.000000
INT_LOG 0.031079  0.007029 -0.045504 -0.015442  0.039388
INT_IDBO 0.041983  0.013162 -0.039944 -0.016234  0.040947
INT_MISC 0.125162  0.037227  0.129274 -0.119776  0.085055
INT_ANY 0.053484  0.007275  0.056438 -0.061946  0.064759

INT_LOG  INT_IDBO  INT_MISC  INT_ANY
0 -0.143600 -0.133252 -0.077852 -0.175605
year -0.143601 -0.133253 -0.077847 -0.175596
month -0.002302 -0.002034 -0.002554 -0.001199
iday -0.001540 -0.001621 -0.002027 -0.001336
extended 0.071768  0.075147  0.027335  0.080767
...
released 0.039388  0.040947  0.085055  0.064759
INT_LOG 1.000000  0.996211  0.052537  0.891051
INT_IDBO 0.996211  1.000000  0.082014  0.893811
INT_MISC 0.052537  0.082014  1.000000  0.252193
INT_ANY 0.891051  0.893811  0.252193  1.000000

(77 rows x 77 columns)

Visualization of Data

In [16]:
data.hist(figsize=(20,20))
plt.show()

In [17]:
countries_with_most_terrorism = data.country_txt.value_counts().head(10)

Index(['Iraq', 'Pakistan', 'Afghanistan', 'India', 'Colombia', 'Philippines',
      'Peru', 'El Salvador', 'United Kingdom', 'Turkey'],
      dtype='object')

In [18]:
countries = list(countries_with_most_terrorism.index)
print("Countries with most terrorism are as follows:", countries)

Countries with most terrorism are as follows:
['Iraq', 'Pakistan', 'Afghanistan', 'India', 'Colombia', 'Philippines', 'Peru', 'El Salvador', 'United Kingdom', 'Turkey']

In [19]:
fig, ax = plt.subplots(figsize=(14,5))
ax.bar(countries_with_most_terrorism.index, countries_with_most_terrorism.values, color='r')
plt.title("Countries with most terrorist")
plt.xlabel('Country')
plt.ylabel('Values')
plt.show()

Countries with most terrorist

In [20]:
sns.displot(data, x='country', kind='kde', fill='tree', palette='colorblind', color='g')
sns.set_style('darkgrid')

In [21]:
sns.distplot(data[data['country']])

C:\Users\Admin\anaconda3\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function) or 'kdeplot' (an axes-level function for histograms).
warnings.warn(msg, FutureWarning)

Out[21]:
In [22]:
data['year'].value_counts().head(10)

Out[22]:
2014    16903
2015    14965
2016    13587
2017    12515
2018    10900
2019     8522
2020     5076
2021     5071
2022     4826
2023     4809
Name: year, dtype: int64

In [23]:
years = list(data['year'].value_counts().head(10).index)
print("Years which had most terrorist attacks in decreasing order : \n", years)

Years which had most terrorist attacks in decreasing order :
[2014, 2015, 2016, 2013, 2017, 2012, 2011, 1993, 2010, 2008]

In [24]:
year = data['year'].value_counts().head(5)
plt.figure(figsize=(12,6))
ax.bar(year.index, year.values, color='yellow')
plt.title("Years with most terrorist attacks")
plt.xlabel("year")
plt.ylabel("Values")
plt.show()

years with most terrorist attacks

In [25]:
sns.displot(data, x='year', kind='hist', palette='colorblind', color='b')
plt.style.use('ggplot')

In [26]:
sns.displot(data, x='year', kind='kde', fill='tree', palette='colorblind', color='b')
sns.axes_grid1.FacetGrid at 9617963d400>

Out[26]:
In [27]:
sns.boxplot(x = data['year'])

Out[27]:
In [28]:
sns.scatterplot(y=data['year'], x=data['country'])

Figure size 144x144 with 0 Axes>

Out[28]:
In [29]:
year = data['year'].unique()
years_count = data['year'].value_counts().sort_index()
sns.barplot(x=year, y = years_count, palette = 'colorblind')
plt.xticks(rotation=50)
plt.ylabel("Number of attacks each year")
plt.title("Attack in years")
plt.show()

attack in years

In [30]:
fig=plt.subplots(figsize=(18,9))
sns.countplot(data['year'], order=data['year'].value_counts().index, palette='colorblind')
plt.xticks(rotation=80)
plt.ylabel("count")
plt.title("Attack per year")
plt.show()

attack per year

C:\Users\Admin\anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as keyword arg: x = 'year'. From version 0.12, the only valid positional argument will be "data", and passing other arguments without an explicit keyword will result in an error or misinterpretation.
warnings.warn(msg, FutureWarning)

In [31]:
data['month'].value_counts().head(10)

Out[31]:
5    16975
7    16268
8    15800
9    15563
6    15359
3    12579
4    11512
1    14936
11   14906
9    14160
Name: month, dtype: int64

In [32]:
month_with_most_terrorism = data['month'].value_counts().head(10)
sns.barplot(month_with_most_terrorism.index, month_with_most_terrorism.values, color='r')
print("Month recording most terrorism:\n", month_with_most_terrorism)

Month recording most terrorism:
[5, 7, 8, 10, 6, 3, 4, 1, 11, 9]

In [33]:
fig, ax = plt.subplots(figsize=(14,5))
ax.bar(month_with_most_terrorism.index, month_with_most_terrorism.values, color='r')
plt.title("Months with most terrorism")
plt.xlabel('month no')
plt.ylabel('Values')
plt.show()

months with most terrorism

In [34]:
data['target_type_txt'].value_counts().head(10)

Out[34]:
Private Citizens & Property    43511
Military                      27984
Police                       24506
Government (General)         21283
Business                     20669
Transportation               19358
Utilities                    18799
Unknown                      18779
Religious Figures/Institutions 4440
International Institution     4322
Name: target_type_txt, dtype: int64

In [35]:
targets_with_max_attacks = data['target_type_txt'].value_counts().head(10)
print("Targets with max attacks:\n", targets)

Targets with max attacks:
['Private Citizens & Property', 'Military', 'Police', 'Government (General)', 'Business', 'Transportation', 'Utilities', 'Unknown', 'Religious Figures/Institutions', 'Educational Institution']

In [36]:
target = data['target_type_txt'].value_counts().head(10)
fig, ax = plt.subplots(figsize=(16,5))
ax.bar(target.index, target.values, color='purple')
plt.xticks(rotation=90)
plt.title("Type of targets")
plt.xlabel('targets')
plt.ylabel('Values')
plt.show()

type of targets

In [37]:
data['attack_type_txt'].value_counts().head(10)

Out[37]:
Bombing/Explosion    88255
Armed Assault       42669
Assassination       13512
Government Taking (Kidnapping) 11158
Facility/Infrastructure Attack 10356
Unknown             7176
Unarmed Assault     1015
Hostage Taking (Barricade Incident) 991
Risqueing           859
Name: attack_type_txt, dtype: int64
```



```
In [38]: attack_types=data.attack_type.value_counts().head(10)
attack_type = list(attack_type.index)
print("types of attacks:\n",attack_type)
```

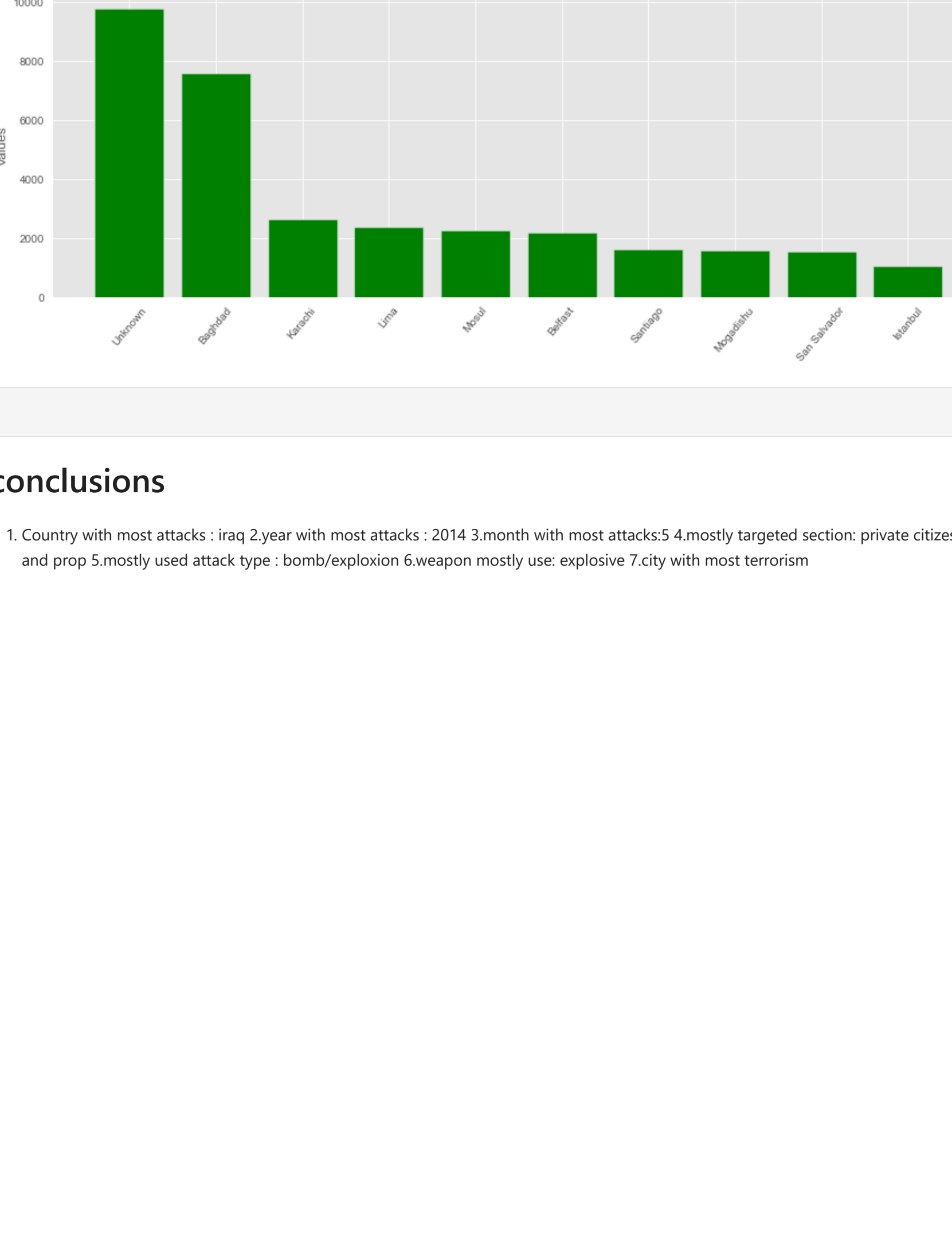
types of attacks:  
'Bombing/Explosion', 'Armed Assault', 'Assassination', 'Hostage Taking (Kidnapping)', 'Facility/Infrastructure Attack', 'Unknown', 'Unarmed Assault', 'Hostage Taking (Barricade Incident)', 'Rijacking']

```
In [39]: fig, ax = plt.subplots(figsize=(14,5))
ax.bar(attack_types.index,attack_types.values,color='r')
plt.xticks(rotation=90)
plt.title("type of attacks")
#plt.xlabel('attacks')
plt.ylabel('Values')
plt.show()
```



```
In [40]: weapon_type=data['weapntype'].value_counts().head(10)
weapon = list(weapon_type.index)
print("types of weapons used the most :\n",weapon)
fig, ax = plt.subplots(figsize=(14,5))
ax.bar(weapon_type.index,weapon_type.values,color='green')
plt.xticks(rotation=90)
plt.title("type of weapons used")
#plt.xlabel('weapons')
plt.ylabel('Values')
plt.show()
```

types of weapons used the most is:  
'Explosives', 'Firearms', 'Unknown', 'Incendiary', 'Melee', 'Chemical', 'Sabotage Equipment', 'Vehicle (not to include vehicle-borne explosives, i.e., car or truck bombs)', 'Other', 'Biological']



```
In [41]: cities_with_most_terrorism = data.city.value_counts().head(10)
cities = list(cities_with_most_terrorism.index)
print("cities with most terrorism are:\n",cities)
fig, ax = plt.subplots(figsize=(15,5))
ax.bar(cities_with_most_terrorism.index,cities_with_most_terrorism.values,color='green')
plt.xticks(rotation=90)
plt.title("cities with most terrorist attacks")
#plt.xlabel('cities')
plt.ylabel('Values')
plt.show()
```

cities with most terrorism are:  
'Unknown', 'Baghdad', 'Karachi', 'Lima', 'Mosul', 'Belfast', 'Santiago', 'Mogadishu', 'San Salvador', 'Istanbul']



```
In [ ]:
```

## conclusions

1. Country with most attacks : Iraq 2.year with most attacks : 2014 3.month with most attacks: 4.mostly targeted section: private citizens and prop 5.mostly used attack type : bomb/explosion 6.weapon mostly use: explosive 7.city with most terrorism