

```
In [1]: import pandas as pd
import numpy as np
import warnings
warnings.filterwarnings('ignore')

In [2]: df=pd.read_csv('https://raw.githubusercontent.com/dsriccientist/DSData/master/happiness_score_dataset.csv')
df

Out[2]:
   Country      Region  Happiness Rank  Happiness Score  Standard Error  Economy (GDP per Capita)  Family  Health (Life Expectancy)  Freedom  Trust (Government Corruption)  Generosity  Dystopia Residual
0  Switzerland  Western Europe          1         7.587          0.03411          1.39551  1.34951          0.94143  0.66557          0.41978  0.29678  2.51738
1  Iceland      Western Europe          2         7.561          0.04884          1.30232  1.40223          0.94784  0.62877          0.14145  0.43630  2.70201
2  Denmark      Western Europe          3         7.527          0.03328          1.32548  1.36058          0.87464  0.64938          0.48357  0.34139  2.49204
3  Norway        Western Europe          4         7.522          0.03880          1.45900  1.33095          0.88521  0.66973          0.36503  0.34699  2.46531
4  Canada        North America          5         7.427          0.02953          1.32629  1.32261          0.90563  0.63297          0.32957  0.46811  2.45176
...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...
153 Rwanda      Sub-Saharan Africa          154         3.465          0.03464          0.22208  0.77370          0.42864  0.59201          0.55191  0.22628  0.67042
154 Benin        Sub-Saharan Africa          155         3.340          0.03956          0.28605  0.35386          0.31910  0.48450          0.08010  0.18260  1.63328
155 Syria        Middle East and Northern Africa          156         3.006          0.05015          0.66320  0.47489          0.71219  0.15684          0.18906  0.47179  0.32858
156 Burundi      Sub-Saharan Africa          157         2.905          0.08958          0.01530  0.41587          0.22296  0.11850          0.10002  0.19727  1.83302
157 Togo          Sub-Saharan Africa          158         2.839          0.06727          0.20888  0.13995          0.28443  0.36453          0.10731  0.16681  1.56726
158 rows x 12 columns

In [3]: df.isnull().sum()

Out[3]:
Country          0
Region           0
Happiness Rank   0
Happiness Score  0
Standard Error   0
Economy (GDP per Capita)  0
Family           0
Health (Life Expectancy)  0
Freedom          0
Trust (Government Corruption)  0
Generosity       0
Dystopia Residual 0
dtype: int64

In [4]: df.columns

Out[4]:
Index(['Country', 'Region', 'Happiness Rank', 'Happiness Score',
       'Standard Error', 'Economy (GDP per Capita)', 'Family',
       'Health (Life Expectancy)', 'Freedom', 'Trust (Government Corruption)',
       'Generosity', 'Dystopia Residual'],
      dtype='object')

country is always a new categorical data and it also mentioned in problem that happiness score is the target column which depends on the economy,family,health, freedom,trust,generosity,dystopia residual

In [5]: df.dtypes

Out[5]:
Country          object
Region           object
Happiness Rank   int64
Happiness Score  float64
Standard Error   float64
Economy (GDP per Capita)  float64
Family           float64
Health (Life Expectancy)  float64
Freedom          float64
Trust (Government Corruption)  float64
Generosity       float64
Dystopia Residual  float64
dtype: object

In [6]: df=df.drop(columns=['Country','Region'])
df

Out[6]:
   Happiness Rank  Happiness Score  Standard Error  Economy (GDP per Capita)  Family  Health (Life Expectancy)  Freedom  Trust (Government Corruption)  Generosity  Dystopia Residual
0               1         7.587          0.03411          1.39551  1.34951          0.94143  0.66557          0.41978  0.29678  2.51738
1               2         7.561          0.04884          1.30232  1.40223          0.94784  0.62877          0.14145  0.43630  2.70201
2               3         7.527          0.03328          1.32548  1.36058          0.87464  0.64938          0.48357  0.34139  2.49204
3               4         7.522          0.03880          1.45900  1.33095          0.88521  0.66973          0.36503  0.34699  2.46531
4               5         7.427          0.02953          1.30269  1.32261          0.90563  0.63297          0.32957  0.46811  2.45176
...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...
153            154         3.465          0.03464          0.22208  0.77370          0.42864  0.59201          0.55191  0.22628  0.67042
154            155         3.340          0.03956          0.28605  0.35386          0.31910  0.48450          0.08010  0.18260  1.63328
155            156         3.006          0.05015          0.66320  0.47489          0.71219  0.15684          0.18906  0.47179  0.32858
156            157         2.905          0.08958          0.01530  0.41587          0.22296  0.11850          0.10002  0.19727  1.83302
157            158         2.839          0.06727          0.20888  0.13995          0.28443  0.36453          0.10731  0.16681  1.56726
158 rows x 10 columns

since the target column is not dependant much on the categorical columns so they are removed

In [7]: import seaborn as sns
sns.distplot(df['Happiness Rank'],kde=True)

Out[7]:
<AxesSubplot: xlabel='Happiness Rank', ylabel='Density'>


In [8]: sns.distplot(df['Happiness Score'],kde=True)

Out[8]:
<AxesSubplot: xlabel='Happiness Score', ylabel='Density'>


In [9]: sns.distplot(df['Standard Error'],kde=True)

Out[9]:
<AxesSubplot: xlabel='Standard Error', ylabel='Density'>


In [10]: sns.distplot(df['Economy (GDP per Capita)'],kde=True)

Out[10]:
<AxesSubplot: xlabel='Economy (GDP per Capita)', ylabel='Density'>


In [11]: sns.distplot(df['Family'],kde=True)

Out[11]:
<AxesSubplot: xlabel='Family', ylabel='Density'>


In [12]: sns.distplot(df['Health (Life Expectancy)'],kde=True)

Out[12]:
<AxesSubplot: xlabel='Health (Life Expectancy)', ylabel='Density'>


In [13]: sns.distplot(df['Trust (Government Corruption)'],kde=True)

Out[13]:
<AxesSubplot: xlabel='Trust (Government Corruption)', ylabel='Density'>


In [14]: sns.distplot(df['Generosity'],kde=True)

Out[14]:
<AxesSubplot: xlabel='Generosity', ylabel='Density'>


In [15]: sns.distplot(df['Dystopia Residual'],kde=True)

Out[15]:
<AxesSubplot: xlabel='Dystopia Residual', ylabel='Density'>


by seeing the plot we can conclude that the columns are having the skewness

In [16]: df.describe()

Out[16]:
   Happiness Rank  Happiness Score  Standard Error  Economy (GDP per Capita)  Family  Health (Life Expectancy)  Freedom  Trust (Government Corruption)  Generosity  Dystopia Residual
count  158.000000          158.000000          158.000000          158.000000          158.000000          158.000000          158.000000          158.000000          158.000000          158.000000
mean           1.7849371          7.375734          0.147805          0.846137          0.991046          0.630259          0.428635          0.143422          0.372726          2.517380
std           0.4514303          1.146010          0.037146          0.403121          0.277269          0.247078          0.150993          0.120034          0.126665          0.553550
min            1.000000          2.839000          0.018480          0.000000          0.000000          0.000000          0.000000          0.000000          0.000000          0.328580
25%           1.000000          4.520000          0.037208          0.545008          0.896023          0.439185          0.203330          0.061575          0.150553          1.784110
50%           1.784937          7.375734          0.147805          0.846137          0.991046          0.630259          0.428635          0.143422          0.372726          2.517380
75%           2.517380          8.243750          0.052300          1.159445          1.214405          0.811013          0.549062          0.180705          0.398980          2.462415
max            5.000000          9.567000          0.136930          1.899420          1.402230          1.075250          0.669730          0.551910          0.799800          3.602140

In [17]: import matplotlib.pyplot as plt
plt.figure(figsize=(15,8))
sns.heatmap(df.describe(),annot=True,linewidth=0.1,linewidth=0.1,linestyle='black',font='0.2f')

Out[17]:
<AxesSubplot: >


In [18]: df.corr()

Out[18]:
   Happiness Rank  Happiness Score  Standard Error  Economy (GDP per Capita)  Family  Health (Life Expectancy)  Freedom  Trust (Government Corruption)  Generosity  Dystopia Residual
Happiness Rank    1.000000          -0.992105          0.158518          -0.782267          -0.723844          -0.739513          -0.558886          -0.372715          -0.160142          -0.512199
Happiness Score    0.992105          1.000000          -0.177254          -0.789661          -0.490565          -0.724264          -0.560211          -0.395199          -0.180319          -0.526474
Standard Error    0.158518          -0.177254          1.000000          -0.217651          -0.120778          -0.101678          -0.139773          -0.137825          -0.098439          0.006091
Economy (GDP per Capita)  -0.782267          -0.789666          -0.217651          1.000000          0.645209          0.016478          0.370000          0.307886          -0.010465          0.004589
Family            -0.723844          -0.490565          -0.120778          0.645209          1.000000          0.011104          0.441318          0.250505          0.077513          0.148117
Health (Life Expectancy)  -0.739513          -0.560211          -0.139773          0.016478          0.441318          1.000000          0.360477          0.493524          0.173914          0.082783
Freedom           -0.558886          0.560211          -0.129773          0.370000          0.441518          0.360477          1.000000          0.493524          0.173914          0.082783
Trust (Government Corruption)  -0.372715          -0.395199          -0.136259          -0.136259          -0.299685          -0.146505          -0.493524          1.000000          0.276123          -0.012186
Generosity         -0.160142          -0.180319          -0.098439          -0.010465          0.077513          0.173914          0.173914          0.276123          1.000000          -0.101301
Dystopia Residual  -0.512199          -0.526474          0.006091          0.004589          0.148117          0.082783          0.082783          -0.012186          -0.101301          1.000000

In [19]: plt.figure(figsize=(15,8))
sns.heatmap(df.corr(),annot=True,linewidth=0.1,linewidth=0.1,linestyle='black',font='0.2f')

Out[19]:
<AxesSubplot: >


In [20]: df.skew()

Out[20]:
Happiness Rank    0.988418
Happiness Score    0.977769
Standard Error    0.385439
Economy (GDP per Capita)  -0.315753
Family            -1.085823
Health (Life Expectancy)  -0.795328
Freedom           -0.413482
Trust (Government Corruption)  -1.389463
Generosity        -1.981951
Dystopia Residual  -2.238911
dtype: float64

columns like Standard Error,Family,Health (Life Expectancy),Trust (Government Corruption),Generosity, Dystopia Residual having the skewness

In [21]: df

Out[21]:
   Happiness Rank  Happiness Score  Standard Error  Economy (GDP per Capita)  Family  Health (Life Expectancy)  Freedom  Trust (Government Corruption)  Generosity  Dystopia Residual
0               1         7.587          0.03411          1.39551  1.34951          0.94143  0.66557          0.41978  0.29678  2.51738
1               2         7.561          0.04884          1.30232  1.40223          0.94784  0.62877          0.14145  0.43630  2.70201
2               3         7.527          0.03328          1.32548  1.36058          0.87464  0.64938          0.48357  0.34139  2.49204
3               4         7.522          0.03880          1.45900  1.33095          0.88521  0.66973          0.36503  0.34699  2.46531
4               5         7.427          0.02953          1.32629  1.32261          0.90563  0.63297          0.32957  0.46811  2.45176
...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...
153            154         3.465          0.03464          0.22208  0.77370          0.42864  0.59201          0.55191  0.22628  0.67042
154            155         3.340          0.03956          0.28605  0.35386          0.31910  0.48450          0.08010  0.18260  1.63328
155            156         3.006          0.05015          0.66320  0.47489          0.71219  0.15684          0.18906  0.47179  0.32858
156            157         2.905          0.08958          0.01530  0.41587          0.22296  0.11850          0.10002  0.19727  1.83302
157            158         2.839          0.06727          0.20888  0.13995          0.28443  0.36453          0.10731  0.16681  1.56726
158 rows x 10 columns

In [22]: columns_list=df.columns.values
ncol=10
mrow=15
plt.figure(figsize=(ncol,ncol))
for i in range(0,int((columns_list[1]))):
    plt.subplot(mrow,ncol,i+1)
    sns.boxplot(df[columns_list[i]],color='red',orient='h')
    plt.tight_layout()



by seeing the box plot we can conclude that the data set having the outliers

In [23]: from scipy.stats import skew
zmap=skew(df)
print(df.shape)
print(z.shape)
threshold=3
print(np.where(z>3))

(158, 10)
(158, 9)
(array([ 27, 48, 64, 115, 128, 147, 153, 155, 157], dtype=int64), array([ 2, 2, 2, 8, 4, 7, 9, 4], dtype=int64))

In [24]: df_new=df[(z<3).all(axis=1)]
print(df_new.shape)

(149, 10)

all the outliers are removed

In [25]: data_loss=(158-149)/158*100
data_loss

5.6962055164557

In [26]: x=df_new.iloc[:,1:]
y=df_new.iloc[:,0]
print(x)
print(y)

Happiness Rank  Happiness Score  Standard Error  Economy (GDP per Capita)  Family  Health (Life Expectancy)  Freedom  Trust (Government Corruption)  Generosity  Dystopia Residual
0               1         7.587          0.03411          1.39551  1.34951          0.94143  0.66557          0.41978  0.29678  2.51738
1               2         7.561          0.04884          1.30232  1.40223          0.94784  0.62877          0.14145  0.43630  2.70201
2               3         7.527          0.03328          1.32548  1.36058          0.87464  0.64938          0.48357  0.34139  2.49204
3               4         7.522          0.03880          1.45900  1.33095          0.88521  0.66973          0.36503  0.34699  2.46531
4               5         7.427          0.02953          1.32629  1.32261          0.90563  0.63297          0.32957  0.46811  2.45176
...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...
153            154         3.465          0.03464          0.22208  0.77370          0.42864  0.59201          0.55191  0.22628  0.67042
154            155         3.340          0.03956          0.28605  0.35386          0.31910  0.48450          0.08010  0.18260  1.63328
155            156         3.006          0.05015          0.66320  0.47489          0.71219  0.15684          0.18906  0.47179  0.32858
156            157         2.905          0.08958          0.01530  0.41587          0.22296  0.11850          0.10002  0.19727  1.83302
157            158         2.839          0.06727          0.20888  0.13995          0.28443  0.36453          0.10731  0.16681  1.56726
158 rows x 9 columns
0      2.51738
1      2.70201
2      2.49204
3      2.46531
4      2.45176
...  ...
153      0.67042
154      1.63328
155      0.32858
156      1.83302
157      1.56726
158      1.56726

In [27]: from sklearn.preprocessing import MinMaxScaler
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import train_test_split

In [28]: mms=fit_transform(x)
array([[ 0.055578,  1.
,  0.38824838, ...,  0.99131621,  0.96175317,
0.85988121],
[0.01546685,  0.99552352,  0.61212633, ...,  0.91594943,  0.58524637,
0.89989852],
[0.02998488,  0.98965724,  0.36372291, ...,  0.95783832,  0.99616915,
0.78139333],
[0.09984889,  0.17673459,  0.31363499, ...,  0.2637885,  0.43182928,
0.76778853],
[0.0984892,  0.11644652,  0.42697649, ...,  0.64511577,  0.39188377,
0.46898873],
[ 0.
,  0.
,  0.9478978, ...,  0.12274438,  0.48456604,
0.49662251]])

In [29]: for i in range(0,1000):
x_train,x_test,y_train,y_test=train_test_split(x,y,random_state=i,test_size=0.2)
rf=fit(x_train,y_train)
pred_train=rf.predict(x_train)
pred_test=rf.predict(x_test)
if round(r2_score(y_train,pred_train)*100,1)!=round(r2_score(y_test,pred_test)*100,1):
print('at random state '+str(i), 'the model performs very well')
print('r2_score= ',r2_score(y_train,pred_train))
print('training accuracy score = ',r2_score(y_train,pred_train)*100)
print('testing accuracy score = ',r2_score(y_test,pred_test)*100)

at random state 24 the model performs very well
training accuracy score = 96.884313184872
testing accuracy score = 96.8120561362182
at random state 95 the model performs very well
training accuracy score = 96.8872545877154
testing accuracy score = 96.8892808918105
at random state 104 the model performs very well
training accuracy score = 96.94566629193
testing accuracy score = 96.9142033830383
at random state 180 the model performs very well
training accuracy score = 96.9327255984836
testing accuracy score = 96.955643189981
at random state 356 the model performs very well
training accuracy score = 96.9395464527536
testing accuracy score = 96.9380998981389
at random state 359 the model performs very well
training accuracy score = 96.9359909818015
testing accuracy score = 96.9435599286122
at random state 445 the model performs very well
training accuracy score = 96.90002674757341
testing accuracy score = 96.9162591299665
at random state 526 the model performs very well
training accuracy score = 96.87523598853625
testing accuracy score = 96.86459455788598
at random state 522 the model performs very well
training accuracy score = 96.9066114967749
testing accuracy score = 96.9374831353879
at random state 596 the model performs very well
training accuracy score = 96.9523863867684
testing accuracy score = 96.9702897777282
at random state 639 the model performs very well
training accuracy score = 96.9096538728459
testing accuracy score = 96.9505909548495
at random state 671 the model performs very well
training accuracy score = 96.9592888678891
testing accuracy score = 96.9565779358123
at random state 673 the model performs very well
training accuracy score = 96.86280318862175
testing accuracy score = 96.8707028787814
at random state 703 the model performs very well
training accuracy score = 96.84054035724191
testing accuracy score = 96.823521930584
at random state 751 the model performs very well
training accuracy score = 96.81577901197028
testing accuracy score = 96.7681548738571
at random state 752 the model performs very well
training accuracy score = 96.898297458324
testing accuracy score = 96.9013893971527
at random state 786 the model performs very well
training accuracy score = 96.9425458918667
testing accuracy score = 96.9237255984836
at random state 791 the model performs very well
training accuracy score = 96.9484131848722
testing accuracy score = 96.9120561362182
at random state 864 the model performs very well
training accuracy score = 96.9245844583731
testing accuracy score = 96.9120561362182
at random state 864 the model performs very well
training accuracy score = 96.9179998493655
testing accuracy score = 96.93564513161213
at random state 865 the model performs very well
training accuracy score = 96.95888495215653
testing accuracy score = 96.96095164455126

In [30]: x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,random_state=869)
rf=fit(x_train,y_train)
pred_train=rf.predict(x_train)
pred_test=rf.predict(x_test)
print('r2_score= ',r2_score(y_train,pred_train))
print('training accuracy score = ',r2_score(y_train,pred_train)*100)
print('testing accuracy score = ',r2_score(y_test,pred_test)*100)

r2_score= 96.86995164455126

In [31]: from sklearn.ensemble import RandomForestRegressor
rf=RandomForestRegressor()
rf.fit(x_train,y_train)
rf.predict(x_train,y_train)
pred_train=rf.predict(x_train)
pred_test=rf.predict(x_test)
print('r2_score= ',r2_score(y_train,pred_train))
print('training accuracy score = ',r2_score(y_train,pred_train)*100)
print('testing accuracy score = ',r2_score(y_test,pred_test)*100)

r2_score= 67.13346236727979
cross_val_score = 88.77892682652584

linear regressor have the good R2 score with Good cv_score value so model is performing well and giving the good results for linear regression (R2_score=96.96 and CV_score=88.14)
```