Usecase:

Loading the data:

```
#Read in the csv file and convert to a Pandas dataframe
World_Happiness_2015 = pd.read_csv("datafrom_2015.csv ")
World_Happiness_2016 = pd.read_csv("datafrom_2016.csv ")
World_Happiness_2017 = pd.read_csv("datafrom_2017.csv ")
World_Happiness_2018 = pd.read_csv("datafrom_2018.csv ")
World_Happiness_2019 = pd.read_csv("datafrom_2019.csv ")
#World_Happiness_2019 = pd.read_csv("datafrom_2019.csv ")
#World_Happiness = pd.concat([World_Happiness_2015, World_Happiness_2016, World_Happiness_2017, World_Happiness_2018, World_Happiness_2019])
```

merge all datagrams:

```
: # mearge all dataframes

World_Happiness = pd.concat([World_Happiness_2015,World_Happiness_2016,World_Happiness_2017,World_Happiness_2018,World_Happiness_2019])
```

Viewing the dataframe:

Viewing the dataframe

We can get a quick sense of the size of our dataset by using the shape method. This returns a tuple with the number of rows and columns in the dataset.

[117]: World_Happiness

]:	Country	Region	Happiness Rank	Happiness Score	Standard Error	Economy (GDP per Capita)	Family	Health (Life Expectancy)	Freedom	Trust (Government Corruption)	 TrustGovernment.Corruption.	Dystopia.Residua
0	Switzerland	Western Europe	1.0	7.587	0.03411	1.39651	1.34951	0.94143	0.66557	0.41978	 NaN	Naf
1	Iceland	Western Europe	2.0	7.561	0.04884	1.30232	1.40223	0.94784	0.62877	0.14145	 NaN	Naf
2	Denmark	Western Europe	3.0	7.527	0.03328	1.32548	1.36058	0.87464	0.64938	0.48357	 NaN	Naf
3	Norway	Western Europe	4.0	7.522	0.03880	1.45900	1.33095	0.88521	0.66973	0.36503	 NaN	Naf
4	Canada	North America	5.0	7.427	0.03553	1.32629	1.32261	0.90563	0.63297	0.32957	 NaN	Naf
151	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	 NaN	Naf
152	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	 NaN	Naf
153	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	 NaN	Activate Winc
154	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	 NaN	Go to Settings to a

Data Profiling before do consistency processes:

1. Data Profiling:

Data profiling is a comprehensive process of examining the data available in an existing dataset and collecting statistics and information about that data.

```
[119]: World Happiness.info
[119]: <bound method DataFrame.info of
            Switzerland Western Europe
                                                 1.0
                                                               7.587
               Iceland Western Europe
                                                 2.0
                                                               7.561
               Denmark Western Europe
                                                 3.0
                                                               7.527
                Norway
                       Western Europe
                                                 4.0
                                                               7.522
                Canada
                                                 5.0
                                                               7.427
       151
                   NaN
                                  NaN
                                                 NaN
                                                                NaN
       152
                   NaN
                                  NaN
                                                 NaN
                                                                NaN
                                                                NaN
                   NaN
                                  NaN
                                                 NaN
       153
       155
                   NaN
                                  NaN
                                                 NaN
                                                                NaN
            Standard Error Economy (GDP per Capita)
                                                   Family
                  0.03411
                                          1.39651 1.34951
                  0.03328
                                          1.32548
                                                  1.36058
                                          1.45900
                  0.03880
                                                  1.33095
       4
                  0.03553
                                          1.32629
                                                  1.32261
       151
                                              NaN
                      NaN
                                              NaN
                                                      NaN
       153
                      NaN
                                              NaN
                                                      NaN
       154
                      NaN
                                              NaN
                                                      NaN
[121]: World_Happiness.shape
[121]: (782, 30)
[123]: World_Happiness.describe()
                                                                                    Trust
                                          Economy
                                Standard
                                                           Health (Life
            Happiness
                     Happiness
                                                                                                    Dystopia
                                          (GDP per
                                                     Family
                                                                      Freedom
                                                                              (Government
                                                                                         Generosity
                                                                                                            ... Health..Life.Expectancy. Trust..Go
                                                           Expectancy)
                                          Capita)
                                                                                Corruption)
      count 315.000000 315.000000 158.000000 315.000000 470.000000
                                                            315.000000 470.000000
                                                                                315.000000 782.000000 315.000000 ...
                                                                                                                         155.000000
                                0.047885
                                          0.899837
                                                             0.594054
                                                                      0.402828
                                                                                           0.218576
                                                                                                    2.212032 ..
                                                                                                                          0.551341
                                                                                                                          0.237073
        std
            45 538922
                       1 141531
                                0.017146
                                          0.410780
                                                   0.318707
                                                             0.240790
                                                                      0.150356
                                                                                  0.115490
                                                                                           0.122321
                                                                                                    0.558728
                                                                                                                          0.000000
             1.000000
                       2.839000
                                0.018480
                                                             0.000000
                                                                      0.000000
                                                                                  0.000000
                                                                                           0.000000
                                                                                                    0.328580 ...
       min
                                         0.000000
                                                   0.000000
             40.000000
                       4.510000
                                0.037268
                                          0.594900
                                                   0.793000
                                                             0.419645
                                                                       0.297615
                                                                                  0.061315
                                                                                           0.130000
                                                                                                     1.884135 ..
                                                                                                                          0.369866
             79.000000
                       5,286000
                                0.043940
                                          0.973060
                                                   1.025665
                                                             0.640450
                                                                      0.418347
                                                                                  0.106130
                                                                                           0.201982
                                                                                                    2,211260 ...
                                                                                                                          0.606042
       50%
                       6.269000
                                0.052300
                                                                                  0.178610
                                                                                                    2.563470 ...
       max 158.000000
                       7.587000
                                0.136930
                                          1.824270
                                                   1.610574
                                                             1.025250
                                                                      0.669730
                                                                                  0.551910
                                                                                           0.838075
                                                                                                    3.837720 ...
                                                                                                                          0.949492
     8 rows × 27 columns
  [125]:
            World Happiness.columns
  [125]: Index(['Country', 'Region', 'Happiness Rank', 'Happiness Score',
                       'Standard Error', 'Economy (GDP per Capita)', 'Family',
                       'Health (Life Expectancy)', 'Freedom', 'Trust (Government Corruption)',
                       'Generosity', 'Dystopia Residual', 'Lower Confidence Interval',
                       'Upper Confidence Interval', 'Happiness.Rank', 'Happiness.Score',
                       'Whisker.high', 'Whisker.low', 'Economy..GDP.per.Capita.',
                       'Health..Life.Expectancy.', 'Trust..Government.Corruption.',
                       'Dystopia.Residual', 'Overall rank', 'Country or region', 'Score',
                       'GDP per capita', 'Social support', 'Healthy life expectancy',
                       'Freedom to make life choices', 'Perceptions of corruption'],
                     dtype='object')
```

When I merged the data frames, I found that the order of the columns did not match, and their names did not match. Now I will perform some operations to correct these problems and make the data consistent.

1.Reorder the columns:

Dataframs 2015 & 2016:

```
deype- object /
[267]: #reorder dataframe for 2015
      new_or1 = ['Country', 'Happiness Rank', 'Happiness Score',
             'Economy (GDP per Capita)', 'Family',
            'Health (Life Expectancy)', 'Freedom', 'Trust (Government Corruption)',
             'Generosity', 'Dystopia Residual' ,'Region', 'Standard Error']
      World_Happiness_2015 = World_Happiness_2015[new_or1]
      World_Happiness_2015.columns
'Freedom', 'Trust (Government Corruption)', 'Generosity',
             'Dystopia Residual', 'Region', 'Standard Error'],
            dtype='object')
[269]: #reorder dataframe for 2016
      new_or2 = ['Country', 'Happiness Rank', 'Happiness Score',
             'Economy (GDP per Capita)', 'Family', 'Health (Life Expectancy)',
             'Freedom', 'Trust (Government Corruption)', 'Generosity',
             'Dystopia Residual', 'Region', 'Lower Confidence Interval', 'Upper Confidence Interval']
      World_Happiness_2016 = World_Happiness_2016[new_or2]
      World_Happiness_2016.columns
'Freedom', 'Trust (Government Corruption)', 'Generosity',
             'Dystopia Residual', 'Region', 'Lower Confidence Interval',
             'Upper Confidence Interval'],
            dtype='object')
```

Dataframes 2017&2018&2019:

```
[271]: #reorder dataframe for 2017
         new_or3 = ['Country', 'Happiness.Rank', 'Happiness.Score',
               'Economy..GDP.per.Capita.', 'Family', 'Health..Life.Expectancy.',
                'Freedom', 'Trust..Government.Corruption.', 'Generosity',
               'Dystopia.Residual', 'Whisker.high',
                'Whisker.low']
         World_Happiness_2017 = World_Happiness_2017[new_or3]
         World Happiness 2017.columns
 [271]: Index(['Country', 'Happiness.Rank', 'Happiness.Score',
                'Economy..GDP.per.Capita.', 'Family', 'Health..Life.Expectancy.',
                'Freedom', 'Trust..Government.Corruption.', 'Generosity',
                'Dystopia.Residual', 'Whisker.high', 'Whisker.low'],
               dtype='object')
 [273]: #reorder dataframe for 2018 & 2019
         new_or4 = ['Country or region', 'Overall rank', 'Score',
                'GDP per capita', 'Social support', 'Healthy life expectancy',
                'Freedom to make life choices', 'Perceptions of corruption', 'Generosity',
         World_Happiness_2018 = World_Happiness_2018[new_or4]
         World Happiness 2019 = World Happiness 2019 [new or4]
         World_Happiness_2018.columns
         World_Happiness_2019.columns
 [273]: Index(['Country or region', 'Overall rank', 'Score', 'GDP per capita',
                'Social support', 'Healthy life expectancy',
                'Freedom to make life choices', 'Perceptions of corruption',
                'Generosity'],
               dtype='object')
Add Year column for each datafram:
#add a new column for Year
World Happiness 2015['Year'] = 2015
#add a new column for Year
 World Happiness 2016['Year'] = 2016
#add a new column for Year
 World Happiness 2017['Year'] = 2017
 #add a new column for Year
World Happiness 2018['Year'] = 2018
 World Happiness 2019['Year'] = 2019
```

Rename Columns:

```
[275]: #rename columns in dataframs
        #2015
        columns = \{World\_Happiness\_2015.columns[i]: standard\_columns\_name and order[i] \ \ \ for \ i \ \ in \ range(len(standard\_columns\_name and order))\}\}
        World_Happiness_2015.rename(columns=columns, inplace=True)
        columns = {World_Happiness_2016.columns[i]: standard_columns_nameandorder[i] for i in range(len(standard_columns_nameandorder))}
        World_Happiness_2016.rename(columns=columns, inplace=True)
        #2017
        columns = \{World\_Happiness\_2017.columns[i]: standard\_columns\_name and order[i] \ \ \ for \ i \ \ in \ range(len(standard\_columns\_name and order))\}\}
        World_Happiness_2017.rename(columns=columns, inplace=True)
        columns = \{World\_Happiness\_2018.columns[i]: standard\_columns\_name and order[i] \ \ \textbf{for} \ \ i \ \ in \ range(len(standard\_columns\_name and order)))\} \\
        World_Happiness_2018.rename(columns=columns, inplace=True)
        columns = {World_Happiness_2019.columns[i]: standard_columns_nameandorder[i] for i in range(len(standard_columns_nameandorder)))}
        World_Happiness_2019.rename(columns=columns, inplace=True)
```

Merge dataframs again and view dataframe:

[263]: # new merge all data frame agin World_Happiness_v2 = pd.concat([World_Happiness_2015,World_Happiness_2016,World_Happiness_2017,World_Happiness_2018,World_Happiness_20 [265]: World_Happiness_v2

]:		Country	Happiness Rank	Happiness Score	Economy (GDP per Capita)	Family	Health (Life Expectancy)	Freedom	Trust (Government Corruption)	Generosity	Dystopia Residual	Region	Standard Error	Year
	0	Switzerland	1	7.587	1.39651	1.34951	0.94143	0.66557	0.41978	0.29678	2.51738	Western Europe	0.03411	2015
	1	Iceland	2	7.561	1.30232	1.40223	0.94784	0.62877	0.14145	0.43630	2.70201	Western Europe	0.04884	2015
â	2	Denmark	3	7.527	1.32548	1.36058	0.87464	0.64938	0.48357	0.34139	2.49204	Western Europe	0.03328	2015
	3	Norway	4	7.522	1.45900	1.33095	0.88521	0.66973	0.36503	0.34699	2.46531	Western Europe	0.03880	2015
	4	Canada	5	7.427	1.32629	1.32261	0.90563	0.63297	0.32957	0.45811	2.45176	North America	0.03553	2015
1	51	Rwanda	152	3.334	0.35900	0.71100	0.61400	0.55500	0.41100	0.21700	NaN	NaN	NaN	2019
1	52	Tanzania	153	3.231	0.47600	0.88500	0.49900	0.41700	0.14700	0.27600	NaN	NaN	NaN	2019
1	53	Afghanistan	154	3.203	0.35000	0.51700	0.36100	0.00000	0.02500	0.15800	NaN	NaN	NaN	2019
1	54	Central African Republic	155	3.083	0.02600	0.00000	0.10500	0.22500	0.03500	0.23500	NaN	NaN	NaN	2019
		South		2.252										2010

Profiling new dataframe:

Info:

```
[223]: print("Data info for dataframe After doing consistency processes : ",World_Happiness_v2.info)
       Data info for dataframe After doing consistency processes :  <bound method DataFrame.info</pre> of
                         Switzerland
       0
                                                               7.587
                                                                7.561
                             Iceland
       1
       2
                             Denmark
                                                                7.527
                              Norway
                                                                7.522
       4
                              Canada
                                                  5
                                                               7.427
                                                                3.334
       151
                              Rwanda
                                                 152
       152
                            Tanzania
                                                 153
                                                                3.231
                         Afghanistan
                                                 154
                                                                3.203
       154 Central African Republic
                                                 155
                                                                3.083
       155
                         South Sudan
                                                 156
                                                               2.853
            Economy (GDP per Capita) Family Health (Life Expectancy) Freedom \
       0
                             1.39651 1.34951
                                                                0.94143
                                                                        0.66557
       1
                             1.30232 1.40223
                                                                0.94784
                                                                        0.62877
                             1.32548 1.36058
                                                                0.87464 0.64938
       2
                             1.45900
                                     1.33095
                                                                0.88521
                                                                        0.66973
       4
                             1.32629
                                     1.32261
                                                               0.90563
                             0.35900 0.71100
                                                               0.61400 0.55500
       151
                             0.47600 0.88500
                                                                        0.41700
                                                               0.49900
       152
       153
                             0.35000 0.51700
                                                                0.36100
                                                                        0.00000
       154
                             0.02600
                                     0.00000
                                                                0.10500
                                                                        0.22500
       155
                             0.30600
                                     0.57500
                                                               0.29500
                                                                        0.01000
            Trust (Government Corruption) Generosity Dystopia Residual
       0
                                  0.41978
                                              0.29678
                                                                2.51738
                                              0.43630
       2
                                  0.48357
                                              0.34139
                                                                2.49204
       3
                                  0.36503
                                              0.34699
                                                                2.46531
       4
                                  0.32957
                                              0.45811
                                                                2.45176
```

Shape:

```
]: print("The shape of dataframe After doing consistency processes : ",World_Happiness_v2.shape)
```

The shape of dataframe After doing consistency processes: (782, 17)

Description:

```
[229]: print("The Description of dataframe After doing consistency processes: ")
World_Happiness_v2.describe()

The Description of dataframe After doing consistency processes:
```

9]:		Happiness Rank	Happiness Score	Economy (GDP per Capita)	Family	Health (Life Expectancy)	Freedom	Trust (Government Corruption)	Generosity	Dystopia Residual	Standard Error	Lower Confidence Interval	Upper Confidence Interval	Dystopia.Re
	count	782.000000	782.000000	782.000000	782.000000	782.000000	782.000000	781.000000	782.000000	315.000000	158.000000	157.000000	157.000000	155.C
	mean	78.698210	5.379018	0.916047	1.078392	0.612416	0.411091	0.125436	0.218576	2.212032	0.047885	5.282395	5.481975	1.8
	std	45.182384	1.127456	0.407340	0.329548	0.248309	0.152880	0.105816	0.122321	0.558728	0.017146	1.148043	1.136493	0.5
	min	1.000000	2.693000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.328580	0.018480	2.732000	3.078000	0.3
	25%	40.000000	4.509750	0.606500	0.869363	0.440183	0.309768	0.054000	0.130000	1.884135	0.037268	4.327000	4.465000	1.5
	50%	79.000000	5.322000	0.982205	1.124735	0.647310	0.431000	0.091000	0.201982	2.211260	0.043940	5.237000	5.419000	1.8
	75%	118.000000	6.189500	1.236187	1.327250	0.808000	0.531000	0.156030	0.278832	2.563470	0.052300	6.154000	6.434000	2.1
	max	158.000000	7.769000	2.096000	1.644000	1.141000	0.724000	0.551910	0.838075	3.837720	0.136930	7.460000	7.669000	3.1

Columns:

Data Quality Checks:

Data Quality Checks

Data quality checks involve the process of ensuring that the data is accurate, complete, consistent, relevant, and reliable.

Here are typical steps involved in checking data quality:

1. Reliability:

Evaluate the data's source and collection process to determine its trustworthiness.

[133]: #In the kaggle page mentioned, the data source is Creative Commons Organization

2. Timeliness:

Ensure the data is up-to-date and reflective of the current situation or the period of interest for the analysis.

[136]: #Data from 2015 to 2019

Consistency:

Check the data type:

3. Consistency:

Confirm that the data is consistent within the dataset and across multiple data sources. For exam

[273]: World_Happiness_v2.dtypes [273]: Country object int64 Happiness Rank float64 Happiness Score Economy (GDP per Capita) float64 Family float64 Health (Life Expectancy) float64 Freedom float64 Trust (Government Corruption) float64 float64 Generosity float64 Dystopia Residual object Region float64 Standard Error int64 Lower Confidence Interval float64 Upper Confidence Interval float64 Dystopia.Residual float64 float64 Whisker.high Whisker.low float64 dtype: object

Some countries have different names depending on the data set.

```
'Count': country_value_counts.values
       table = tabulate(country_counts_df, headers='keys', tablefmt='grid')
       print(table)
        154 | North Cyprus
                                       3 |
       | 155 | Comoros
                                       3 |
       | 156 | Belize
                                       3 |
                                       2 |
       157 | Northern Cyprus
       | 158 | Suriname
                                       2 |
       | 159 | Swaziland
                                       2 |
       | 160 | Puerto Rico
                                       1 |
       | 161 | Somaliland Region
                                                                                                                   Activate Wind
       | 162 | Oman
                                      1 |
```

4. Relevance:

Next Day working:

4. Relevance:

```
[]: #| I will delete some columns because some columns are not present in all the dataframes.

#These are the columns that are not present in all dataframes and also do not help me [Dystopia Residual, Region, Standard Error, Lower Confidence Interval #Upper Confidence Interval, Dystopia. Residual, Whisker.high, Whisker.low]

[521]: World_Happiness_v2.drop(['Dystopia Residual', 'Region','Lower Confidence Interval','Standard Error', 'Upper Confidence Interval', 'Dystopia.Residual', 'Whisker.high', 'Whisker.low'], axis=1, inplace=True)

[523]: World_Happiness_v2.shape

[523]: (782, 10)
```

5. Uniqueness:

The data has zero duplicate.



6.Completeness:

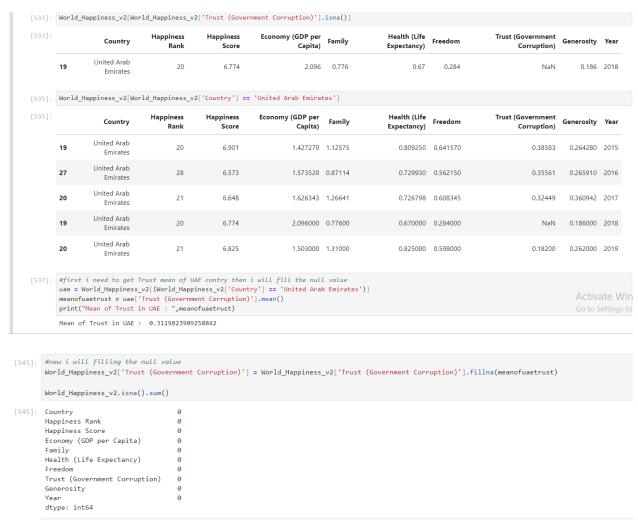
One null value in "Trust (Government Corruption)" column.



Data Cleaning

Handling missing values:

Here the missing value was the trust in the government for the UAE in 2018. I took the mean of the trust in the government for all years for the UAE and then compensated for it in the missing value.



Correcting errors

Change the value name:

```
#Here I will change the names to one name
World_Happiness_v2['Country'].replace('Hong Kong S.A.R., China', 'Hong Kong', inplace=True)
World_Happiness_v2['Country'].replace('Somaliland region', 'Somalia', inplace=True)
World_Happiness_v2['Country'].replace('Taiwan Province of China', 'Taiwan', inplace=True)
World_Happiness_v2['Country'].replace('North Macedonia', 'Macedonia', inplace=True)
World_Happiness_v2['Country'].replace('Trinidad & Tobago', "Trinidad and Tobago", inplace=True)
World_Happiness_v2['Country'].replace('Northern Cyprus', "North Cyprus", inplace=True)
```

After make change:

There are countries that were mentioned only a few times and this is because they have less than 5 Count.

#here the row that have country name Northern Cyprus is 5 that because I changed North Cyprus to Northern Cyprus

[1560]:	#There	I will review e are countries the row that I _Happiness_v2_\	☆ ©	↑ ↓ ±	₽						
[1560]:		Country	Happiness Rank	Happiness Score	Economy (GDP per Capita)	Family	Health (Life Expectancy)	Freedom	Trust (Government Corruption)	Generosity	Year
	65	Northern Cyprus	66	5.695	1.208060	1.070080	0.923560	0.490270	0.142800	0.261690	2015
	61	Northern Cyprus	62	5.771	1.311410	0.818260	0.841420	0.435960	0.165780	0.263220	2016
	60	Northern Cyprus	61	5.810	1.346911	1.186303	0.834647	0.471204	0.155353	0.266846	2017
	57	Northern Cyprus	58	5.835	1.229000	1.211000	0.909000	0.495000	0.154000	0.179000	2018
	63	Northern Cyprus	64	5.718	1.263000	1.252000	1.042000	0.417000	0.162000	0.191000 Activ	ate V