

APENDICE II - SCRIPT EM PYTHON

TCC Expectativa de Vida

Libraries

Pip Install

```
In [2]: # !pip install inflection
# !pip install pandas
# !pip install sqlalchemy
# !pip install geopy
# !pip install pycountry-convert
# !pip install seaborn
# !pip install plotly
# !pip install pycountry-convert
# !pip install melt
# !pip install sklearn
# !pip install Image
# !pip install missingno
# !pip install boruta
# !pip install xgboost
# !pip install xlrd
# !pip install -U notebook-as-pdf
# !pypeteer-install
```

Import

```
In [259...]:
import pandas as pd
import numpy as np
from geopy.geocoders import Nominatim
from pycountry_convert import country_alpha2_to_continent_code, country_name_to_country_alpha2
import inflection
import sqlite3
from matplotlib import pyplot as plt
from IPython.core.display import HTML
import seaborn as sns
from IPython.display import Image
from scipy.stats import shapiro
from sqlalchemy import create_engine
import missingno as msno
import plotly.express as px
from scipy import stats
from sklearn.impute import KNNImputer
from sklearn.preprocessing import MinMaxScaler
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import RobustScaler
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
from sklearn.linear_model import LinearRegression, Lasso, Ridge
from sklearn.ensemble import RandomForestRegressor
from boruta import BorutaPy
import xgboost as xgb
import random

import warnings
warnings.filterwarnings('ignore')
```

Auxiliary functions

```
In [260...]:
def cross_validation( x_training, kfold, model_name, model, verbose=False ): #funcao que realiza cross validation
    mae_list = []
    mape_list = []
    rmse_list = []
    for k in reversed( range( 1, kfold+1 ) ):
        if verbose:
            print( '\nKFold Number: {}'.format( k ) )
        # start and end date for validation
        validation_start_date = x_training['year'].max() - (k*3)
        validation_end_date = x_training['year'].max() - ((k-1)*3)
        # filtering dataset
        training = x_training[x_training['year'] < validation_start_date]
        validation = x_training[(x_training['year'] >= validation_start_date) & (x_training['year'] <= validation_end_date)]
```

```

# training and validation dataset
# training
xtraining = training.drop( ['year', 'life_expectancy'], axis=1 )
ytraining = training['life_expectancy']

# validation
xvalidation = validation.drop( ['year', 'life_expectancy'], axis=1 )
yvalidation = validation['life_expectancy']

# model
m = model.fit( xtraining, ytraining )

# prediction
yhat = m.predict( xvalidation )

#performance
m_result = ml_error( model_name, np.expm1( yvalidation ), np.expm1( yhat ) )

# store performance of each kfold iteration
mae_list.append( m_result['MAE'] )
mape_list.append( m_result['MAPE'] )
rmse_list.append( m_result['RMSE'] )

return pd.DataFrame( { 'Model Name': model_name,
                      'MAE CV': np.round( np.mean( mae_list ), 2 ).astype( str ) + ' +/- ' + np.round( np.std( mae_list ), 2 ).astype( str ),
                      'MAPE CV': np.round( np.mean( mape_list ), 2 ).astype( str ) + ' +/- ' + np.round( np.std( mape_list ), 2 ).astype( str ),
                      'RMSE CV': np.round( np.mean( rmse_list ), 2 ).astype( str ) + ' +/- ' + np.round( np.std( rmse_list ), 2 ).astype( str ) } )

def mean_percentage_error( y, yhat ):
    '''função que retorna o erro percentual medio'''
    return np.mean( ( y - yhat ) / y )

def mean_absolute_percentage_error( y, yhat ):
    '''função que retorna o percentual do erro medio absoluto'''
    return np.mean( np.abs( ( y - yhat ) / y ) )

def ml_error( model_name, y, yhat ):
    '''função que retorna um dataframe com os indicadores de performance'''
    mae = mean_absolute_error( y, yhat )
    mape = mean_absolute_percentage_error( y, yhat )
    rmse = np.sqrt( mean_squared_error( y, yhat ) )

    return pd.DataFrame( { 'Model Name': model_name,
                           'MAE': mae,
                           'MAPE': mape,
                           'RMSE': rmse }, index=[0] )

def cramer_v( x, y ):
    '''função que retorna a matrix de crammer'''
    cm = pd.crosstab( x, y ).to_numpy()
    n = cm.sum()
    r, k = cm.shape

    chi2 = stats.chi2_contingency( cm )[0]
    chi2corr = max( 0, chi2 - (k-1)*(r-1)/(n-1) )

    kcorr = k - (k-1)**2/(n-1)
    rcorr = r - (r-1)**2/(n-1)
    return np.sqrt( (chi2corr/n) / ( min( kcorr-1, rcorr-1 ) ) )

def jupyter_settings():
    '''função que define os parametros do %notebook'''
    %matplotlib inline
    %pylab inline

    plt.style.use( 'bmh' )
    plt.rcParams[ 'figure.figsize' ] = [ 25, 12 ]
    plt.rcParams[ 'font.size' ] = 24

    display( HTML( '<style>.container { width:75% !important; }</style>' ) )
    pd.options.display.max_columns = None
    pd.options.display.max_rows = None
    pd.set_option( 'display.expand_frame_repr', False )

    sns.set()

    # Supress Scientific Notation

```

```
np.set_printoptions(suppress=True)
pd.set_option('display.float_format', '{:.2f}'.format)
```

Data Acquisition

Dataset Expectativa de vida

In [261...]

```
# **1** Dataframe da OMS e das Nações Unidas com a variavel life-expectation
df_expec=pd.read_csv('../Datasets/Life_Expectancy_Data.csv',parse_dates=[1])
df_expec.sample(5)
```

Out[261...]

	Country	Year	Status	Life expectancy	Adult Mortality	infant deaths	Alcohol	percentage expenditure	Hepatitis B	Measles	BMI	under-five deaths	Polio	expi
348	Botswana	2003-01-01	Developing	46.40	693.00	2	5.51	299.37	9.00	59	31.60	4	96.00	
1241	Iraq	2008-01-01	Developing	69.30	167.00	32	0.17	192.16	66.00	5494	54.50	38	71.00	
156	Azerbaijan	2003-01-01	Developing	67.80	154.00	7	0.55	42.41	51.00	1978	43.60	8	79.00	
1165	Hungary	2004-01-01	Developed	72.90	18.00	1	13.28	146.86	nan	0	58.20	1	99.00	
1530	Lithuania	2007-01-01	Developed	72.00	24.00	0	13.40	1581.51	96.00	0	59.00	0	95.00	

In [262...]

```
print('numero de linhas:{}'.format(df_expec.shape[0]))
print('numero de colunas:{}'.format(df_expec.shape[1]))
```

numero de linhas:2938
numero de colunas:22

Dataset Emission

In [263...]

```
# **2** Dataframe com emissao de dados
df_emi=pd.read_csv('../Datasets/emission data.csv')
df_emi.head()
```

Out[263...]

	Country	1751	1752	1753	1754	1755	1756	1757	1758	1759	1760	1761	1762	1763	1764	1765	1766	1767	1768	1
0	Afghanistan	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	Africa	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	Albania	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	Algeria	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	Americas (other)	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Verificação do nomes em country da tabela emissao

In [264...]

```
countryISO=df_expec['Country'].unique()
emissioncountry=df_emi['Country'].unique()
(set(countryISO) ^ set(emissioncountry))-set(countryISO)
(set(countryISO) ^ set(emissioncountry))-set(emissioncountry)
```

Out[264...]

'Bolivia (Plurinational State of)',
'Brunei Darussalam',
'Cabo Verde',
'Czechia',
"Côte d'Ivoire",
"Democratic People's Republic of Korea",
'Democratic Republic of the Congo',
'Iran (Islamic Republic of)',
"Lao People's Democratic Republic",
'Micronesia (Federated States of)',
'Monaco',
'Republic of Korea',

```
'Republic of Moldova',
'Russian Federation',
'San Marino',
'Syrian Arab Republic',
'The former Yugoslav republic of Macedonia',
'Timor-Leste',
'United Kingdom of Great Britain and Northern Ireland',
'United Republic of Tanzania',
'United States of America',
'Venezuela (Bolivarian Republic of)',
'Viet Nam'}
```

In [265...]

```
countryISO_unmatch=['Bolivia (Plurinational State of)',
'Brunei Darussalam',
'Cabo Verde',
'Czechia',
"Côte d'Ivoire",
"Democratic People's Republic of Korea",
'Democratic Republic of the Congo',
'Iran (Islamic Republic of)',
'Lao People's Democratic Republic',
'Micronesia (Federated States of)',
'Republic of Korea',
'Republic of Moldova',
'Russian Federation',
'Syrian Arab Republic',
'The former Yugoslav republic of Macedonia',
'Timor-Leste',
'United Kingdom of Great Britain and Northern Ireland',
'United Republic of Tanzania',
'United States of America',
'Venezuela (Bolivarian Republic of)',
'Viet Nam']

countryemission_unmatch=['Bolivia','Brunei',
'Cape Verde', 'Czech Republic', "Côte d'Ivoire", 'North Korea', 'Democratic Republic of Congo', 'Iran', 'Laos', 'Micr

dictemi = dict(zip(countryemission_unmatch,countryISO_unmatch))
df_emi['Country']=df_emi['Country'].replace(dictemi)
```

In [266...]

```
# Building DataFrame
df_emi=df_emi[['Country','2000','2001','2002','2003','2004','2005','2006','2007','2008',
'2009','2010','2011','2012','2013','2014','2015']]

df_emi=df_emi.melt(id_vars=['Country']) # empilha os dados com granularidade em country
df_emi.columns=['Country','Year','Emission']

df_emi.head()
```

Out[266...]

	Country	Year	Emission
0	Afghanistan	2000	71717793.00
1	Africa	2000	23640083267.00
2	Albania	2000	196932672.00
3	Algeria	2000	2118624684.00
4	Americas (other)	2000	60974588046.00

Dataset Demographics

In [267...]

```
# *** Dataframe com a população separada em masculina feminina e demografia
df_demo=pd.read_csv('../Datasets/WPP2019_TotalPopulationBySex.csv')
df_demo.head()
```

Out[267...]

	LocID	Location	VarID	Variant	Time	MidPeriod	PopMale	PopFemale	PopTotal	PopDensity
0	4	Afghanistan	2	Medium	1950	1950.50	4099.24	3652.87	7752.12	11.87
1	4	Afghanistan	2	Medium	1951	1951.50	4134.76	3705.39	7840.15	12.01
2	4	Afghanistan	2	Medium	1952	1952.50	4174.45	3761.55	7936.00	12.16
3	4	Afghanistan	2	Medium	1953	1953.50	4218.34	3821.35	8039.68	12.31
4	4	Afghanistan	2	Medium	1954	1954.50	4266.48	3884.83	8151.32	12.49

Verificação do nomes em country da tabela demo

```
In [268...]  
countryISO=df_expec['Country'].unique()  
democountry=df_demo['Location'].unique()  
(set(countryISO) ^ set(democountry))-set(countryISO)  
print((set(countryISO) ^ set(democountry))-set(democountry))  
  
{"Democratic People's Republic of Korea", 'United Kingdom of Great Britain and Northern Ireland', 'Swaziland', 'Micronesia (Federated States of)', 'The former Yugoslav republic of Macedonia'}  
  
In [269...]  
countryISO_unmatch=["Democratic People's Republic of Korea",'Micronesia (Federated States of)', 'The former Yugoslav republic of Macedonia', 'North Macedonia', 'United Kingdom of Great Britain and Northern Ireland', 'Swaziland', 'The former Yugoslav republic of Macedonia']  
countrydemo_unmatch=['Dem. People's Republic of Korea', 'Micronesia (Fed. States of)', 'North Macedonia', 'United Kingdom of Great Britain and Northern Ireland', 'Swaziland', 'The former Yugoslav republic of Macedonia']  
dictdemo = dict(zip(countrydemo_unmatch,countryISO_unmatch))  
df_demo['Location']=df_demo['Location'].replace(dictdemo)  
  
In [270...]  
## Seleção das colunas  
df_demo=df_demo.loc[(df_demo['Time']>=2000 )& (df_demo['Time']<=2015),['Location','Time','PopMale','PopFemale','PopTotal','Density']]  
df_demo.columns=['Country','Year','Pop male','Pop female','Pop total','Density']  
  
In [271...]  
df_demo.head()
```

	Country	Year	Pop male	Pop female	Pop total	Density
50	Afghanistan	2000	10689.51	10090.45	20779.96	31.83
51	Afghanistan	2001	11117.75	10489.24	21606.99	33.10
52	Afghanistan	2002	11642.11	10958.67	22600.77	34.62
53	Afghanistan	2003	12214.63	11466.24	23680.87	36.27
54	Afghanistan	2004	12763.73	11962.96	24726.69	37.87

Merge Datasets

```
In [272...]  
# Transforming variable date  
df_demo['Year']=df_demo['Year'].astype(str)  
df_emi['Year'] = pd.to_datetime(df_emi['Year']).dt.strftime('%Y')  
df_expec['Year'] = pd.to_datetime(df_expec['Year']).dt.strftime('%Y')  
df_demo['Year'] = pd.to_datetime(df_demo['Year']).dt.strftime('%Y')  
  
In [273...]  
# Merge Dataframes  
  
df=pd.merge(df_expec,df_emi,how='left',on=['Country','Year'])  
df=pd.merge(df,df_demo,how='left',on=['Country','Year'])  
df=df.drop(['Population'], axis=1)  
df.head()
```

	Country	Year	Status	Life expectancy	Adult Mortality	infant deaths	Alcohol	percentage expenditure	Hepatitis B	Measles	BMI	under-five deaths	Polio	expenditure
0	Afghanistan	2015	Developing	65.00	263.00	62	0.01	71.28	65.00	1154	19.10	83	6.00	
1	Afghanistan	2014	Developing	59.90	271.00	64	0.01	73.52	62.00	492	18.60	86	58.00	
2	Afghanistan	2013	Developing	59.90	268.00	66	0.01	73.22	64.00	430	18.10	89	62.00	
3	Afghanistan	2012	Developing	59.50	272.00	69	0.01	78.18	67.00	2787	17.60	93	67.00	
4	Afghanistan	2011	Developing	59.20	275.00	71	0.01	7.10	68.00	3013	17.20	97	68.00	

DF2 API geolocation

```
In [274...]  
## Return Latitude  
def geolocate_lat(country):  
    geolocator = Nominatim(user_agent='catuserbot')  
  
    try:  
        # Geolocate the center of the country  
        loc = geolocator.geocode(country)  
        # And return latitude  
        return loc.latitude
```

```

except:
    # Return missing value
    return np.nan

## Return Longitude
def geolocate_long(country):
    geolocator = Nominatim(user_agent='catuserbot')

    try:
        # Geolocate the center of the country
        loc = geolocator.geocode(country)
        # And return longitude
        return loc.longitude

    except:
        # Return missing value
        return np.nan

```

```

In [275...]: # api that get Lat and Long for each country (takes 2 hours to run)
#df['Lat'],df['Long']=df['Country'].apply(lambda x: geolocate_lat(x)),df['Country'].apply(lambda x: geolocate_long(x))

In [276...]: #df.to_csv('Datasets/DF_complete.csv',index=False) #create a file con Lat and Long because the procedure takes 2 hours

In [277...]: df= pd.read_csv('../Datasets/DF_complete.csv')

```

Populating Country code and continent code from function

```

In [278...]: #function to convert to alpha2 country codes

def get_code(col):
    try:
        cn_a2_code = country_name_to_country_alpha2(col)
    except:
        cn_a2_code = 'Unknown'
    try:
        cn_continent = country_alpha2_to_continent_code(cn_a2_code)
    except:
        cn_continent = 'Unknown'
    return (cn_a2_code)

#function to convert to alpha2 country continent

def get_continent(col):
    try:
        cn_a2_code = country_name_to_country_alpha2(col)
    except:
        cn_a2_code = 'Unknown'
    try:
        cn_continent = country_alpha2_to_continent_code(cn_a2_code)
    except:
        cn_continent = 'Unknown'
    return (cn_continent)

In [279...]: df['continent'],df['code']=df['Country'].apply(lambda x: get_continent(x)),df['Country'].apply(lambda x: get_code(x))

In [280...]: aux = df.loc[df['continent']=='Unknown','Country'].unique()
aux

Out[280...]: array(['Bolivia (Plurinational State of)', 'Iran (Islamic Republic of)',
       'Micronesia (Federated States of)', 'Republic of Korea',
       'The former Yugoslav republic of Macedonia', 'Timor-Leste',
       'Venezuela (Bolivarian Republic of)'], dtype=object)

In [281...]: cont=['SA','AS','OC','AS','EU','AS','SA']

In [282...]: for i in range(len(aux)):
    for j in range(len(df)):
        if aux[i] == df.loc[j,'Country'] :
            df.loc[j,'continent']=cont[i]

```

Changing columns name

```
In [283... df.columns

Out[283... Index(['Country', 'Year', 'Status', 'Life expectancy ', 'Adult Mortality',
       'infant deaths', 'Alcohol', 'percentage expenditure', 'Hepatitis B',
       'Measles ', ' BMI ', 'under-five deaths ', 'Polio', 'Total expenditure',
       'Diphtheria ', ' HIV/AIDS', 'GDP', ' thinness 1-19 years',
       ' thinness 5-9 years', 'Income composition of resources', 'Schooling',
       'Emission', 'Pop male', 'Pop female', 'Pop total', 'Density', 'lat',
       'long', 'continent', 'code'],
      dtype='object')

In [284... cols_original=['Country', 'Year', 'Status', 'Life expectancy ', 'Adult Mortality',
       'infant deaths', 'Alcohol', 'percentage expenditure', 'Hepatitis B',
       'Measles ', ' BMI ', 'under-five deaths ', 'Polio', 'Total expenditure',
       'Diphtheria ', ' HIV/AIDS', 'GDP', ' thinness 1-19 years',
       ' thinness 5-9 years', 'Income composition of resources', 'Schooling',
       'Emission', 'Pop male', 'Pop female', 'Pop total', 'Density', 'lat',
       'long', 'continent', 'code']

# snakecase = lambda x: inflection.underscore(x)
# cols = list(map(snakecase, cols_original))
# df.columns=cols

cols=['country', 'year', 'status', 'life_expectancy', 'adult_mortality',
       'infant_deaths', 'alcohol', 'percentage_expenditure', 'hepatitis_b',
       'measles', 'bmi', 'under_five_deaths', 'polio', 'total_expenditure',
       'diphtheria', 'hiv_aids', 'gdp', 'thinness_1_19_years',
       'thinness_5_9_years', 'income_composition_of_resources', 'schooling',
       'emission', 'pop_male', 'pop_female', 'pop_total', 'density', 'lat',
       'long', 'continent', 'code']

df.columns=cols
```

```
In [285... df.head()
```

	country	year	status	life_expectancy	adult_mortality	infant_deaths	alcohol	percentage_expenditure	hepatitis_b	measles
0	Afghanistan	2015	Developing	65.00	263.00	62	0.01	71.28	65.00	1154
1	Afghanistan	2014	Developing	59.90	271.00	64	0.01	73.52	62.00	492
2	Afghanistan	2013	Developing	59.90	268.00	66	0.01	73.22	64.00	430
3	Afghanistan	2012	Developing	59.50	272.00	69	0.01	78.18	67.00	2787
4	Afghanistan	2011	Developing	59.20	275.00	71	0.01	7.10	68.00	3013

Data Preparation

```
In [286... jupyter_settings()

Populating the interactive namespace from numpy and matplotlib
```

```
In [287... df1=df.copy()
```

Creating Features

```
In [288... # female percent
df1['perc_female']=df1['pop_female']/df1['pop_total']

# emission per population tax
df1['emission_pop']=df1['emission']/df1['pop_total']
```

Data Info

```
In [289... print('numero de linhas:{}'.format(df.shape[0]))
print('numero de colunas:{}'.format(df.shape[1]))
```

numero de linhas:2938
numero de colunas:30

In [290...]

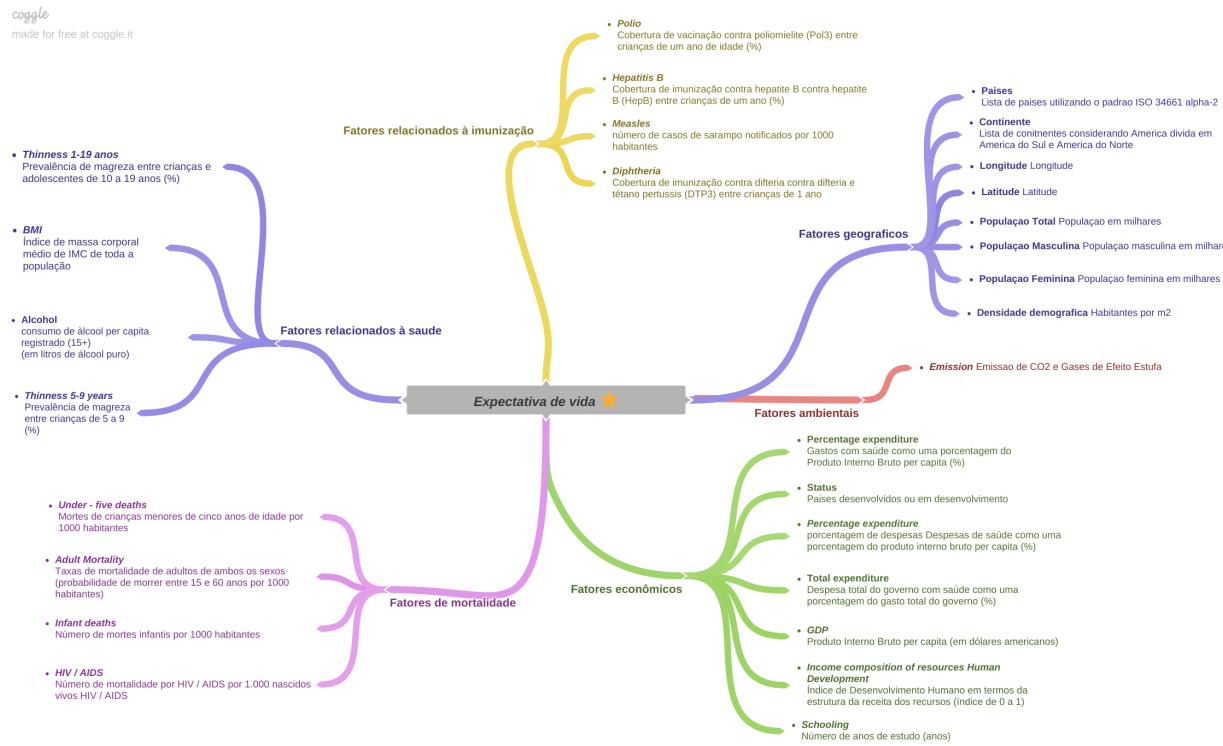
```
df1.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2938 entries, 0 to 2937
Data columns (total 32 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   country          2938 non-null    object  
 1   year              2938 non-null    int64  
 2   status             2938 non-null    object  
 3   life_expectancy   2928 non-null    float64 
 4   adult_mortality   2928 non-null    float64 
 5   infant_deaths     2938 non-null    int64  
 6   alcohol            2744 non-null    float64 
 7   percentage_expenditure  2938 non-null    float64 
 8   hepatitis_b        2385 non-null    float64 
 9   measles            2938 non-null    int64  
 10  bmi                2904 non-null    float64 
 11  under_five_deaths  2938 non-null    int64  
 12  polio               2919 non-null    float64 
 13  total_expenditure  2712 non-null    float64 
 14  diphtheria         2919 non-null    float64 
 15  hiv_aids           2938 non-null    float64 
 16  gdp                2490 non-null    float64 
 17  thinness_1_19_years 2904 non-null    float64 
 18  thinness_5_9_years  2904 non-null    float64 
 19  income_composition_of_resources 2771 non-null    float64 
 20  schooling           2775 non-null    float64 
 21  emission            2936 non-null    float64 
 22  pop_male            2912 non-null    float64 
 23  pop_female          2912 non-null    float64 
 24  pop_total            2922 non-null    float64 
 25  density              2922 non-null    float64 
 26  lat                 2922 non-null    float64 
 27  long                2920 non-null    float64 
 28  continent            2938 non-null    object  
 29  code                2938 non-null    object  
 30  perc_female          2912 non-null    float64 
 31  emission_pop         2920 non-null    float64 
dtypes: float64(24), int64(4), object(4)
memory usage: 734.6+ KB
```

In [291...]

```
Image('../img/Expectativa_de_vida_star.png')
```

Out[291...]



Descriptive Statistics

```
In [292...]: num_attributes = df1.select_dtypes( include=['int64', 'float64'] )
cat_attributes = df1.select_dtypes( exclude=['int64', 'float64', 'datetime64[ns]'] )
```

Numerical Atributes

```
# Central Tendency - mean, median
ct1 = pd.DataFrame( num_attributes.apply( np.mean ) ).T
ct2 = pd.DataFrame( num_attributes.apply( np.median ) ).T

# dispersion - std, min, max, range, skew, kurtosis
d1 = pd.DataFrame( num_attributes.apply( np.std ) ).T
d2 = pd.DataFrame( num_attributes.apply( min ) ).T
d3 = pd.DataFrame( num_attributes.apply( max ) ).T
d4 = pd.DataFrame( num_attributes.apply( lambda x: x.max() - x.min() ) ).T
d5 = pd.DataFrame( num_attributes.apply( lambda x: x.skew() ) ).T
d6 = pd.DataFrame( num_attributes.apply( lambda x: x.kurtosis() ) ).T

# concatenar
m = pd.concat( [d2, d3, d4, ct1, ct2, d1, d5, d6] ).reset_index()
m.columns = ['attributes', 'min', 'max', 'range', 'mean', 'median', 'std', 'skew', 'kurtosis']
m
```

Out[293...]	attributes	min	max	range	mean	median	std	skew	kurtosis
0	year	2000.00	2015.00	15.00	2007.52	2008.00	4.61	-0.01	-1.21
1	life_expectancy	36.30	89.00	52.70	69.22	nan	9.52	-0.64	-0.23
2	adult_mortality	1.00	723.00	722.00	164.80	nan	124.27	1.17	1.75
3	infant_deaths	0.00	1800.00	1800.00	30.30	3.00	117.91	9.79	116.04
4	alcohol	0.01	17.87	17.86	4.60	nan	4.05	0.59	-0.80
5	percentage_expenditure	0.00	19479.91	19479.91	738.25	64.91	1987.58	4.65	26.57
6	hepatitis_b	1.00	99.00	98.00	80.94	nan	25.06	-1.93	2.77
7	measles	0.00	212183.00	212183.00	2419.59	17.00	11465.32	9.44	114.86
8	bmi	1.00	87.30	86.30	38.32	nan	20.04	-0.22	-1.29
9	under_five_deaths	0.00	2500.00	2500.00	42.04	4.00	160.42	9.50	109.75

TCC_Expectativa de vida

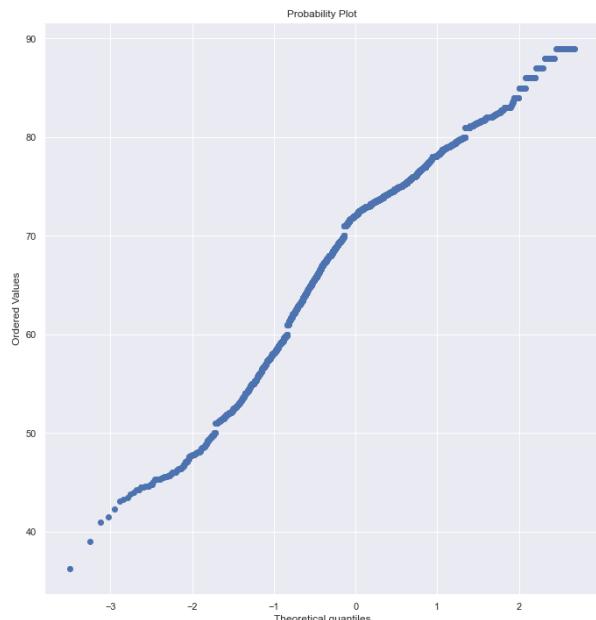
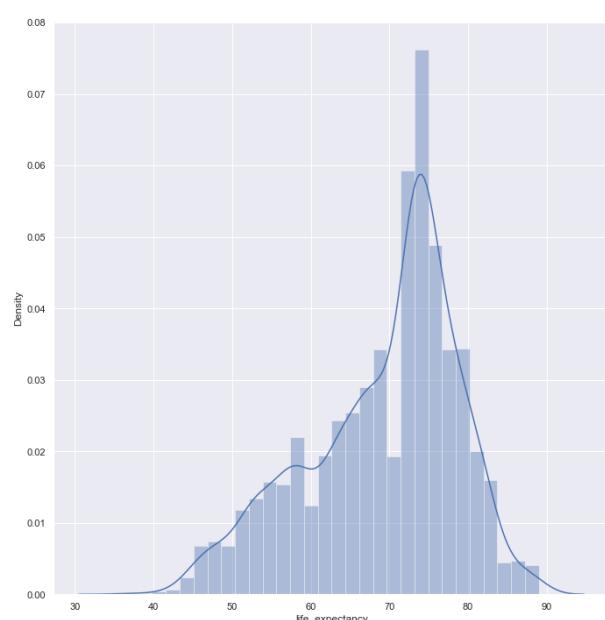
	attributes	min	max	range	mean	median	std	skew	kurtosis
10	polio	3.00	99.00	96.00	82.55	nan	23.42	-2.10	3.78
11	total_expenditure	0.37	17.60	17.23	5.94	nan	2.50	0.62	1.16
12	diphtheria	2.00	99.00	97.00	82.32	nan	23.71	-2.07	3.56
13	hiv_aids	0.10	50.60	50.50	1.74	0.10	5.08	5.40	34.89
14	gdp	1.68	119172.74	119171.06	7483.16	nan	14267.30	3.21	12.33
15	thinness_1_19_years	0.10	27.70	27.60	4.84	nan	4.42	1.71	3.97
16	thinness_5_9_years	0.10	28.60	28.50	4.87	nan	4.51	1.78	4.36
17	income_composition_of_resources	0.00	0.95	0.95	0.63	nan	0.21	-1.14	1.39
18	schooling	0.00	20.70	20.70	11.99	nan	3.36	-0.60	0.89
19	emission	0.00	389000000000.00	389000000000.00	6553095034.87	nan	29181208055.13	9.45	103.59
20	pop_male	35.71	722508.01	722472.30	18584.31	nan	70463.98	8.34	73.22
21	pop_female	40.30	684339.86	684299.57	18285.66	nan	66363.16	8.19	71.49
22	pop_total	1.61	1406847.87	1406846.26	36743.90	nan	136593.19	8.28	72.65
23	density	1.54	24764.43	24762.89	177.21	nan	701.84	19.56	560.18
24	lat	-41.50	64.98	106.48	18.42	nan	24.39	-0.21	-0.65
25	long	-175.20	179.16	354.36	18.45	nan	65.55	-0.07	0.26
26	perc_female	0.23	0.54	0.31	0.50	nan	0.03	-5.33	37.14
27	emission_pop	0.00	1313605.50	1313605.50	183512.06	nan	257204.18	1.98	3.94



In [294...]

```
# Response variable
plt.subplot( 1, 2, 1 )
sns.distplot(df1['life_expectancy']);

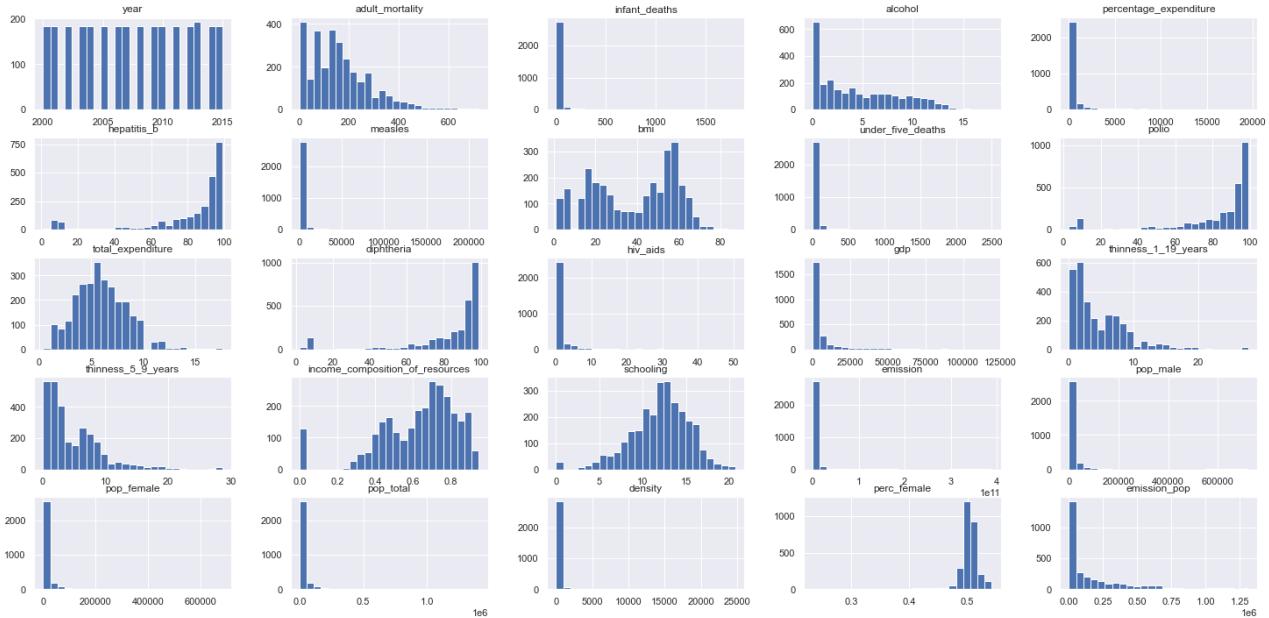
plt.subplot( 1, 2, 2 )
stats.probplot(df1['life_expectancy'], dist='norm', plot=pylab)
pylab.show()
```



In [295...]

```
cols_drop = ['lat', 'long', 'life_expectancy']
num_attributes1=num_attributes.drop(cols_drop, axis=1 )
num_attributes1.hist( bins=25 );
```

TCC_Expectativa de vida



```
In [296...]
# plt.subplot( 1, 6, 1)
# sns.boxplot(x = num_attributes1['emission'])

# plt.subplot( 1, 6, 2)
# sns.boxplot(x = num_attributes1['gdp'])

# plt.subplot( 1, 6, 3)
# sns.boxplot(x = num_attributes1['measles'])

# plt.subplot( 1, 6, 4)
# sns.boxplot(x = num_attributes1['under_five_deaths'])

# plt.subplot( 1, 6, 5)
# sns.boxplot(x = df1['perc_female'])

# plt.subplot( 1, 6, 6)
# sns.boxplot(x = df1['pop_female'])

# plt.subplot( 1, 6, 6)
# sns.boxplot(x = df1['pop_male'])
```

Categorical Atributes

```
In [297...]
cat_attributes.apply( lambda x: x.unique().shape[0] )
```

```
Out[297...]
country      193
status        2
continent     6
code         188
dtype: int64
```

Missing values and Bias data

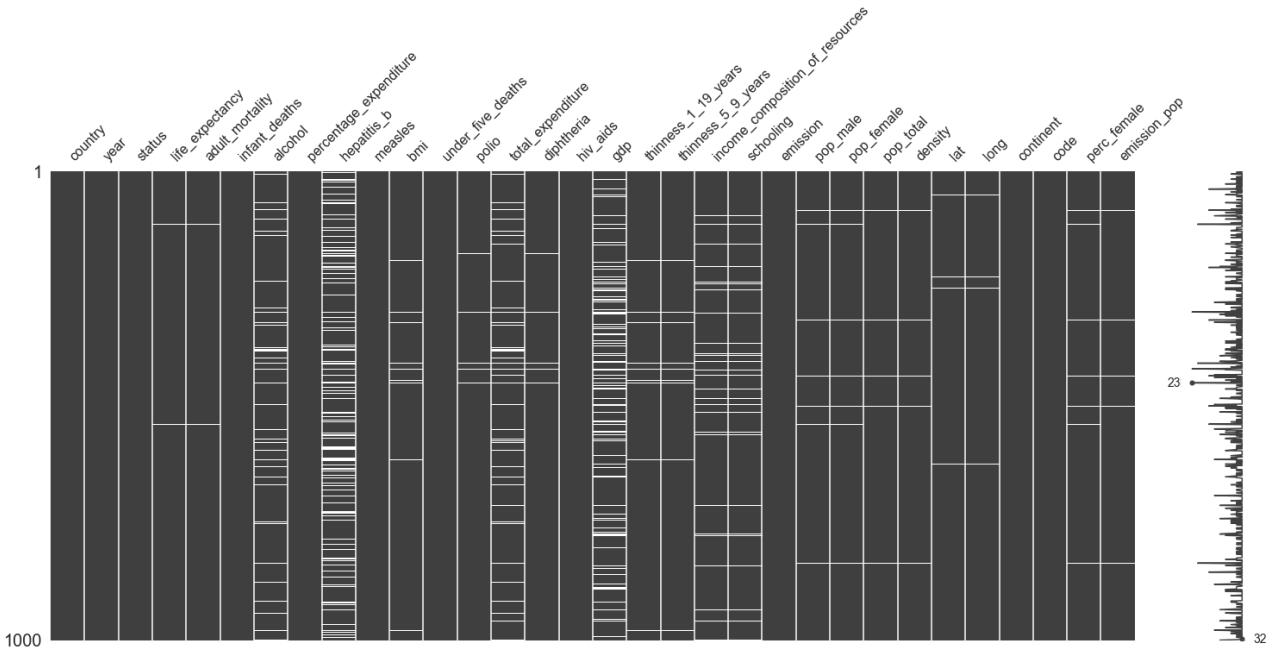
```
In [298...]
# duplicate
df1.drop_duplicates()

print('not exist data duplicated')

print('numero de linhas:{}' .format(df1.shape[0]))
print('numero de colunas:{}' .format(df1.shape[1]))
```

not exist data duplicated
numero de linhas:2938
numero de colunas:32

```
In [299...]
%matplotlib inline
sns.set()
sns.set_style("whitegrid")
sns.set_color_codes("muted")
sns.set_context("notebook")
```



```
In [300...]: aux=df1.isna().sum().sort_values(ascending=False)
aux1=df1.isna().sum().sort_values(ascending=False)/df1.shape[0]*100

na=pd.concat([aux,aux1],axis=1)
na.columns=['NaN', 'NaN %']
na
```

Out[300...]

	NaN	NaN %
hepatitis_b	553	18.82
gdp	448	15.25
total_expenditure	226	7.69
alcohol	194	6.60
income_composition_of_resources	167	5.68
schooling	163	5.55
thinness_1_19_years	34	1.16
bmi	34	1.16
thinness_5_9_years	34	1.16
perc_female	26	0.88
pop_female	26	0.88
pop_male	26	0.88
polio	19	0.65
diphtheria	19	0.65
emission_pop	18	0.61
long	18	0.61
pop_total	16	0.54
density	16	0.54
lat	16	0.54
life_expectancy	10	0.34
adult_mortality	10	0.34
emission	2	0.07
status	0	0.00
year	0	0.00
infant_deaths	0	0.00
hiv_aids	0	0.00
percentage_expenditure	0	0.00

	NaN	NaN %
measles	0	0.00
under_five_deaths	0	0.00
continent	0	0.00
code	0	0.00
country	0	0.00

Bias gdp

```
In [301... df1=df1.drop(['gdp'], axis=1) # drop original column

df_gdp=pd.read_excel('../Datasets/GDPPC-USD-countries.xls',header=2)# New set Dataframe with standard country name
df_gdp=df_gdp[['Country',2000,2001,2002,2003,2004,2005,2006,2007,2008,2009,2010, 2011, 2012, 2013, 2014, 2015]]
df_gdp=df_gdp.melt(id_vars=['Country']) # stacks the data with granularity in the country
df_gdp.columns=['country','year','gdp']
```



```
In [302... # Checking the country names

countryISO=df['country'].unique()
gdpcountry=df_gdp['country'].unique()

(set(countryISO) ^ set(gdpcountry))-set(countryISO)

(set(countryISO) ^ set(gdpcountry))-set(gdpcountry)
```



```
Out[302... {"Democratic People's Republic of Korea",
'Democratic Republic of the Congo',
'Lao People's Democratic Republic',
'Micronesia (Federated States of)',
'Niue',
'Saint Vincent and the Grenadines',
'Swaziland',
'The former Yugoslav republic of Macedonia',
'United Kingdom of Great Britain and Northern Ireland',
'United Republic of Tanzania',
'United States of America'}
```



```
In [303... countryISO_unmatch=["Democratic People's Republic of Korea",
'Democratic Republic of the Congo',
'Lao People's Democratic Republic',
'Micronesia (Federated States of)',
'Saint Vincent and the Grenadines',
'The former Yugoslav republic of Macedonia',
'United Kingdom of Great Britain and Northern Ireland',
'United Republic of Tanzania',
'United States of America']

countrygdp_unmatch=[

'D.P.R. of Korea',
'D.R. of the Congo',
'Lao People's DR',
'Micronesia (FS of)',
'St. Vincent and the Grenadines',
'North Macedonia',
'United Kingdom',
'U.R. of Tanzania: Mainland',
'United States',
]

dicgdp = dict(zip(countrygdp_unmatch,countryISO_unmatch))
df_gdp['country']=df_gdp['country'].replace(dicgdp)

df_gdp['year']=df_gdp['year'].astype(str)
df_gdp['year'] = pd.to_datetime( df_gdp['year'] ).dt.strftime('%Y')
df1['year']=df1['year'].astype(str)
df1['year'] = pd.to_datetime( df1['year'] ).dt.strftime('%Y')

df1=pd.merge(df1,df_gdp,how='left',on=['country','year'])

df1.head()
```

	country	year	status	life_expectancy	adult_mortality	infant_deaths	alcohol	percentage_expenditure	hepatitis_b	measles
0	Afghanistan	2015	Developing	65.00	263.00	62	0.01	71.28	65.00	1154
1	Afghanistan	2014	Developing	59.90	271.00	64	0.01	73.52	62.00	492
2	Afghanistan	2013	Developing	59.90	268.00	66	0.01	73.22	64.00	430
3	Afghanistan	2012	Developing	59.50	272.00	69	0.01	78.18	67.00	2787
4	Afghanistan	2011	Developing	59.20	275.00	71	0.01	7.10	68.00	3013

◀ ▶

Bias percentage _expenditure

```
In [304... #drop column bias in percentage _expenditure derived from Total_expenditure
df1=df1.drop('percentage_expenditure', axis=1)
```

NaN com algoritmo KNN

```
In [305... knn_imputer = KNNImputer(n_neighbors=5, weights="uniform", metric='nan_euclidean')
df1['hepatitis_b'] = knn_imputer.fit_transform(df1[['hepatitis_b']])
df1['alcohol'] = knn_imputer.fit_transform(df1[['alcohol']])
df1['total_expenditure'] = knn_imputer.fit_transform(df1[['total_expenditure']])
df1['income_composition_of_resources'] = knn_imputer.fit_transform(df1[['income_composition_of_resources']])
df1['schooling'] = knn_imputer.fit_transform(df1[['schooling']])
```

NaN lat long

```
In [306... null_cols_code=df1.loc[df1['long'].isnull(),'country'].unique()
null_cols_code
```

```
Out[306... array(['Brunei Darussalam', 'Sierra Leone',
       'The former Yugoslav republic of Macedonia'], dtype=object)
```

```
In [307... lat=[-16.2837065,32.4207423,41.6137143]
lon=[-63.5493965,53.6830157,21.743258]

for i in range(len(null_cols_code)):
    for j in range(len(df1)):
        if null_cols_code[i] == df1.loc[j,'country'] :
            df1.loc[j,'lat']=lat[i]
            df1.loc[j,'long']=lon[i]
```

Drop NaN

```
In [308... df1.isna().sum().sort_values(ascending=False)
```

thinness_1_19_years	34
thinness_5_9_years	34
bmi	34
gdp	33
pop_female	26
pop_male	26
perc_female	26
diphtheria	19
polio	19
emission_pop	18
density	16
pop_total	16
life_expectancy	10
adult_mortality	10
emission	2
alcohol	0
hepatitis_b	0
status	0
infant_deaths	0
year	0
measles	0
continent	0
under_five_deaths	0
total_expenditure	0
long	0
hiv_aids	0
code	0
income_composition_of_resources	0
schooling	0

```
lat          0
country      0
dtype: int64
```

```
In [309... # drop Nan's
df1= df1.dropna(subset=['gdp','polio','bmi','life_expectancy'])
```

```
In [310... df1.isna().sum().sort_values(ascending=False)
```

```
Out[310... gdp          0
hiv_aids      0
year          0
status         0
life_expectancy 0
adult_mortality 0
infant_deaths 0
alcohol        0
hepatitis_b    0
measles        0
bmi            0
under_five_deaths 0
polio          0
total_expenditure 0
diphtheria     0
thinness_1_19_years 0
emission_pop   0
thinness_5_9_years 0
income_composition_of_resources 0
schooling       0
emission        0
pop_male        0
pop_female      0
pop_total       0
density         0
lat             0
long            0
continent       0
code            0
perc_female     0
country         0
dtype: int64
```

```
In [311... print('novo numero de linhas:{}'.format(df1.shape[0]))
print('novo numero de colunas:{}'.format(df1.shape[1]))
```

```
novo numero de linhas:2872
novo numero de colunas:31
```

Exploratory Data Analysis

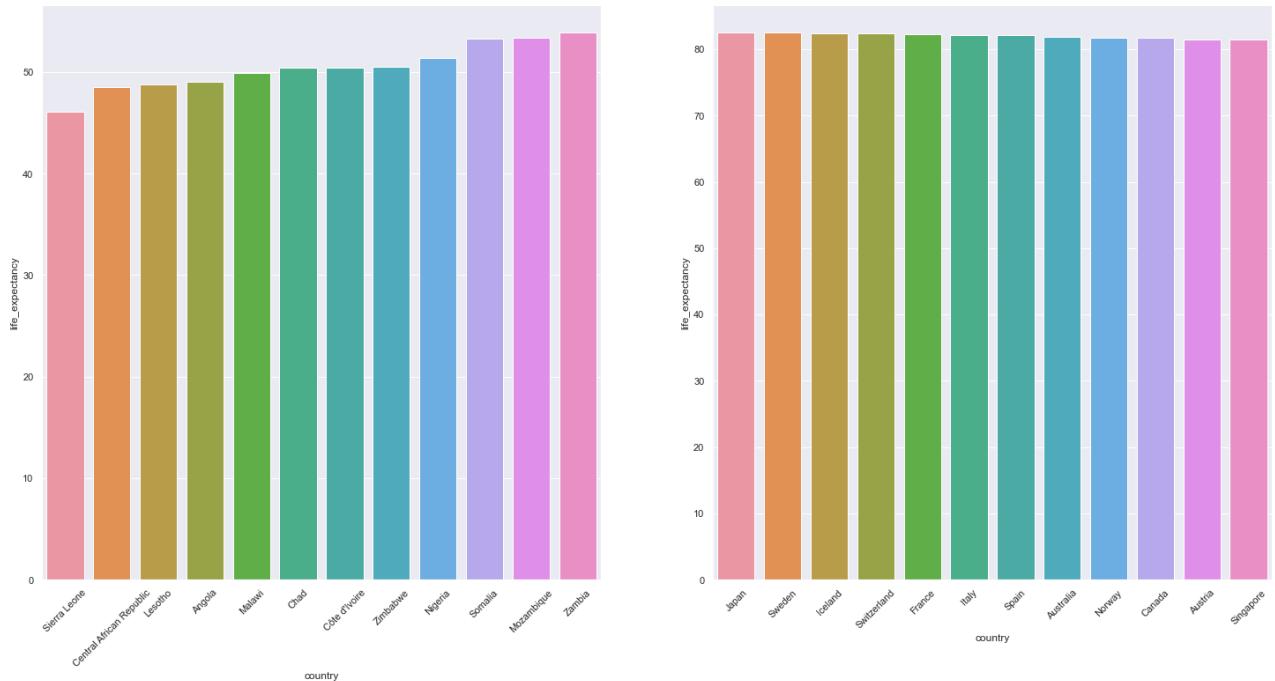
```
In [312... jupyter_settings()
```

```
Populating the interactive namespace from numpy and matplotlib
```

```
In [313... df2=df1.copy()
```

```
In [314... ## continent
plt.subplot( 1, 2, 2 )
a=df2[['country','life_expectancy']].groupby('country').mean().sort_values(ascending=False,by='life_expectancy').re
a=a.head(12)
sns.barplot(x='country',y='life_expectancy',data=a);
plt.xticks(rotation=45);

plt.subplot( 1, 2, 1 )
b=df2[['country','life_expectancy']].groupby('country').mean().sort_values(ascending=True,by='life_expectancy').re
b=b.head(12)
sns.barplot(x='country',y='life_expectancy',data=b);
plt.xticks(rotation=45);
```

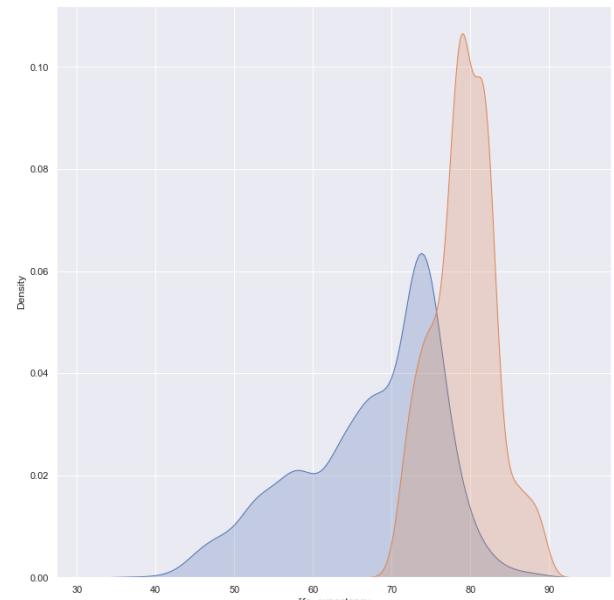
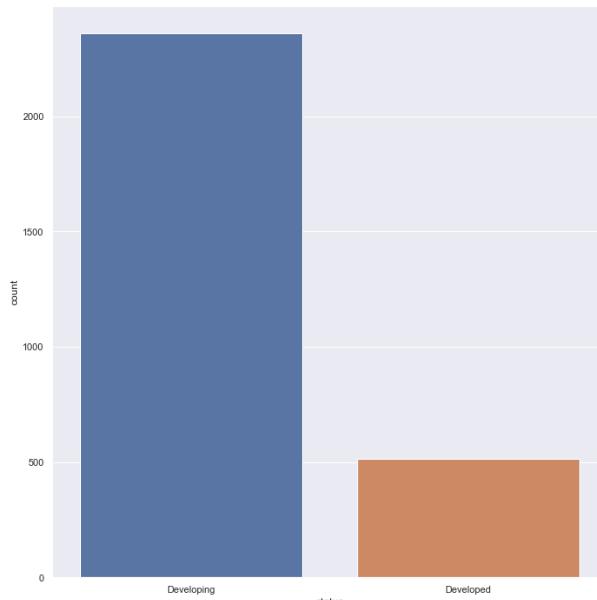


In [315...]

```
## Status

plt.subplot( 1, 2, 1 )
sns.countplot( df1['status'] )

plt.subplot( 1, 2, 2 )
sns.kdeplot( df1[df1['status'] =='Developing']['life_expectancy'], label='status', shade=True )
sns.kdeplot( df1[df1['status'] =='Developed']['life_expectancy'], label='status', shade=True );
```



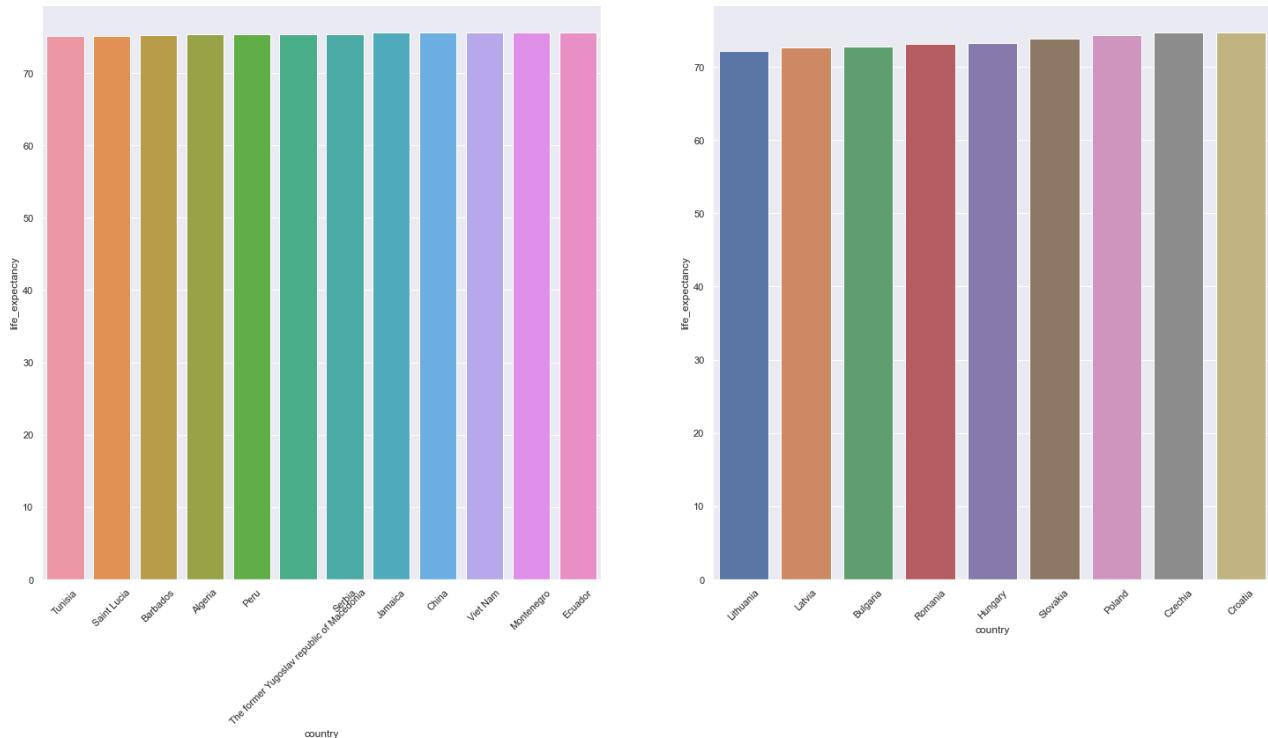
In [316...]

```
df_statusb=df2[df2['status']=='Developing']
df_statusb=df_statusb[df_statusb['life_expectancy']>75]

df_statususa=df2[df2['status']=='Developed']
df_statususa=df_statususa[df_statususa['life_expectancy']<75]

plt.subplot( 1, 2, 1 )
b=df_statusb[['country','life_expectancy']].groupby('country').mean().sort_values(ascending=True,by='life_expectancy')
b=b.head(12)
sns.barplot(x='country',y='life_expectancy',data=b);
plt.xticks(rotation=45);

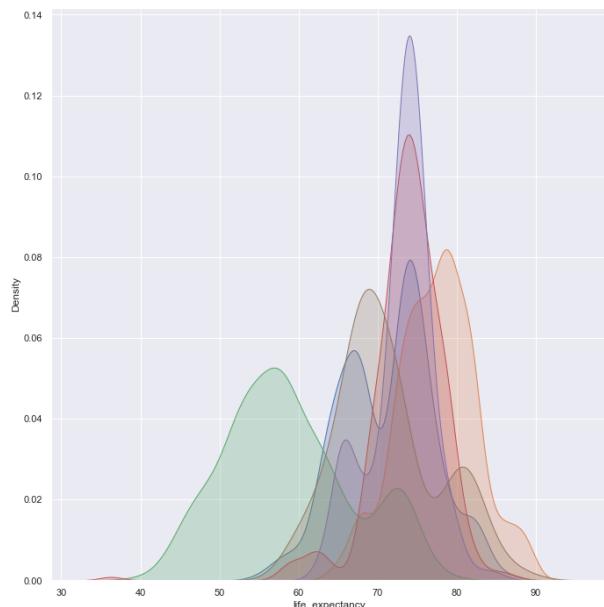
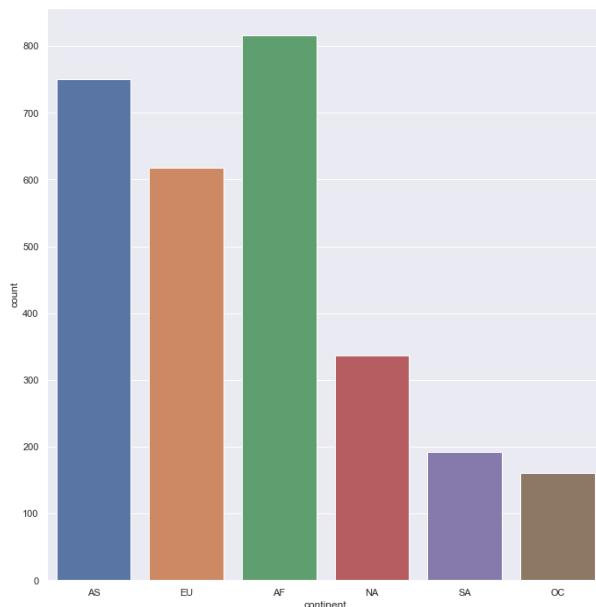
plt.subplot( 1, 2, 2 )
a=df_statususa[['country','life_expectancy']].groupby('country').mean().sort_values(ascending=True,by='life_expectancy')
a=a.head(12)
sns.barplot(x='country',y='life_expectancy',data=a);
plt.xticks(rotation=45);
```



In [317...]

```
## continent
plt.subplot( 1, 2, 1 )
sns.countplot( df1['continent'] )

plt.subplot( 1, 2, 2 )
sns.kdeplot( df1[df1['continent'] == 'AS']['life_expectancy'], label='continent', shade=True )
sns.kdeplot( df1[df1['continent'] == 'EU']['life_expectancy'], label='continent', shade=True )
sns.kdeplot( df1[df1['continent'] == 'AF']['life_expectancy'], label='continent', shade=True )
sns.kdeplot( df1[df1['continent'] == 'NA']['life_expectancy'], label='continent', shade=True )
sns.kdeplot( df1[df1['continent'] == 'SA']['life_expectancy'], label='continent', shade=True )
sns.kdeplot( df1[df1['continent'] == 'OC']['life_expectancy'], label='continent', shade=True );
```



H1 - Países densamente povoados ou altamente populosos tendem a ter menor expectativa de vida?

In [318...]

```
map=df2[['continent','status','country','lat','long','pop_total','density','life_expectancy']].groupby(['continent'])
map.head()
```

Out[318...]

	continent	country	status	lat	long	pop_total	density	life_expectancy
0	AF	Algeria	Developing	28.00	3.00	34821.37	14.62	73.62
1	AF	Angola	Developing	-11.88	17.57	21622.86	17.34	49.02

	continent	country	status	lat	long	pop_total	density	life_expectancy
2	AF	Benin	Developing	9.53	2.26	8628.81	76.52	57.57
3	AF	Botswana	Developing	-23.17	24.59	1888.51	3.33	56.05
4	AF	Burkina Faso	Developing	12.08	-1.69	14618.29	53.43	55.64

In [319...]

```
fig=px.scatter_mapbox(map,
                      hover_name='country',
                      hover_data=["life_expectancy", "density",'status'],
                      lat='lat',
                      lon='long',
                      size='density',
                      color='life_expectancy',
                      color_continuous_scale=px.colors.cyclical.IceFire_r,
                      size_max=60,
                      zoom=1)

fig.update_layout(mapbox_style='open-street-map')
fig.update_layout(height=400, margin={'r':0,'l':0,'t':0,'b':0})
fig.show()
```

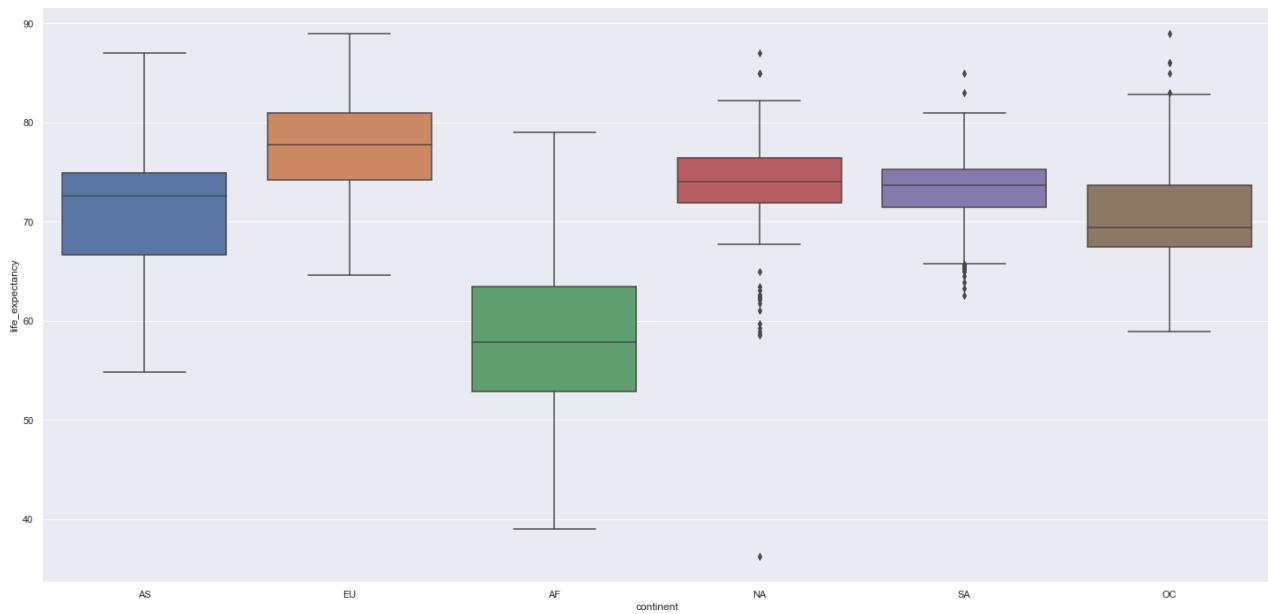
In [320...]

```
fig=px.scatter_mapbox(map,
                      hover_name='country',
                      hover_data=["life_expectancy", "density",'status'],
                      lat='lat',
                      lon='long',
                      size='pop_total',
                      color='life_expectancy',
                      color_continuous_scale=px.colors.cyclical.IceFire_r,
                      size_max=60,
                      zoom=1)

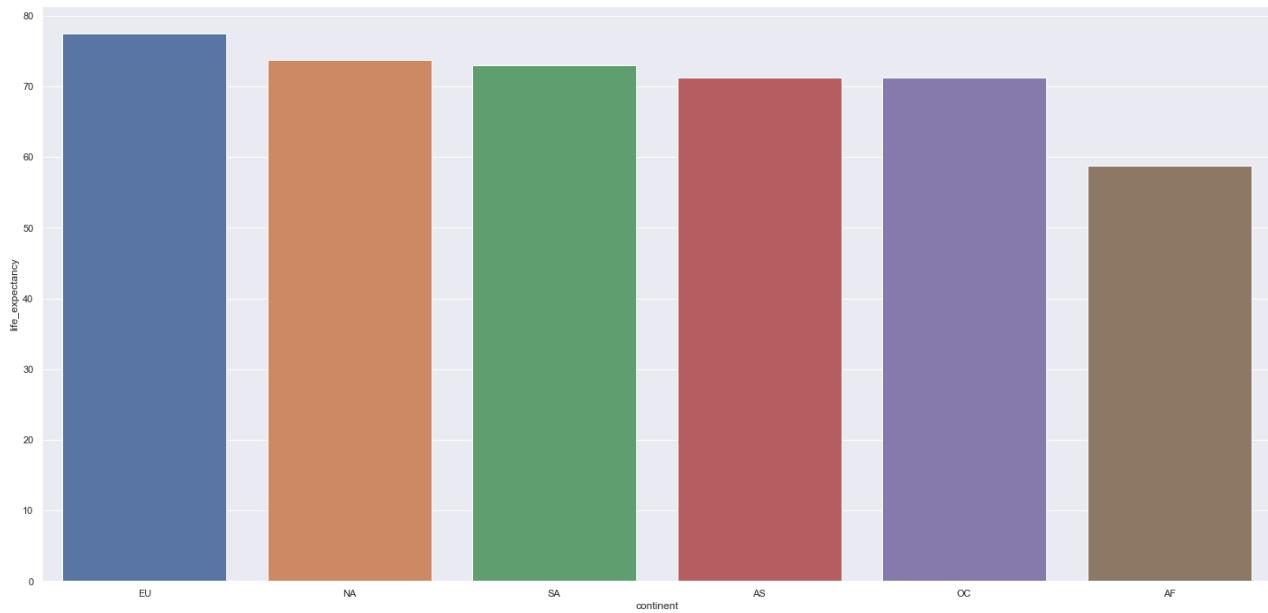
fig.update_layout(mapbox_style='open-street-map')
fig.update_layout(height=400, margin={'r':0,'l':0,'t':0,'b':0})
fig.show()
```

H2- Como é a variação da expectativa de vida dentro dos continentes

```
In [321... aux = df2[['continent','country','life_expectancy','year']]  
#Life_expectancy per continent  
sns.boxplot( x='continent', y='life_expectancy',data=df2);
```

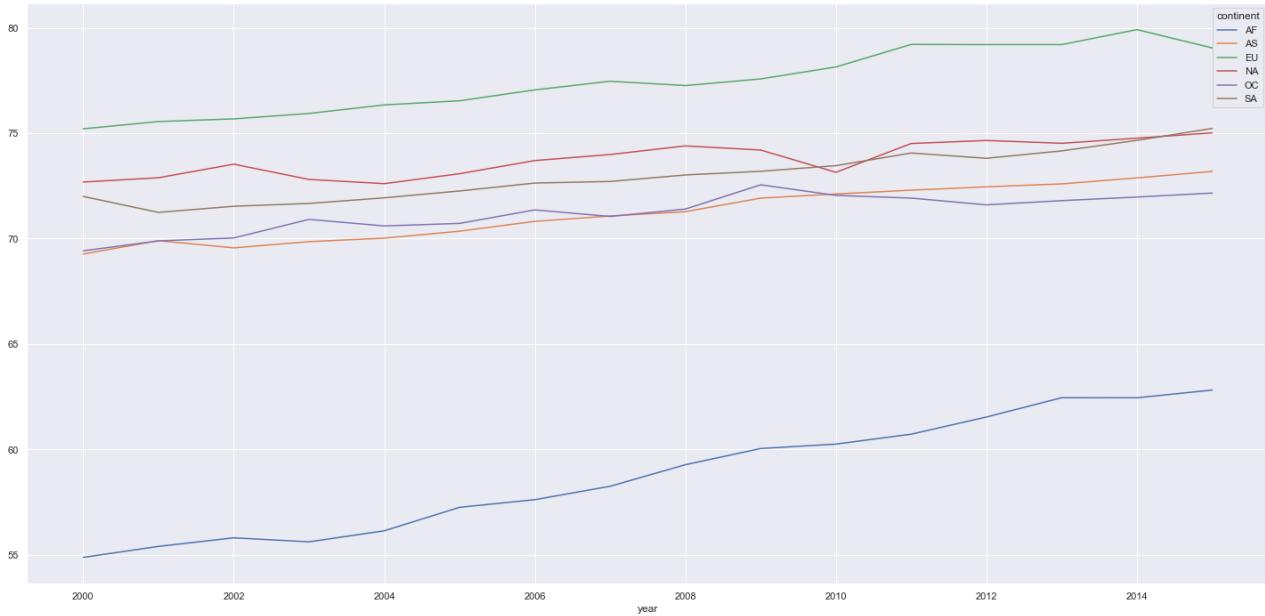


```
In [322... ## continent  
a=df2[['continent','life_expectancy']].groupby('continent').mean().sort_values(ascending=False,by='life_expectancy'  
a=a.head(12)  
sns.barplot(x='continent',y='life_expectancy',data=a);
```



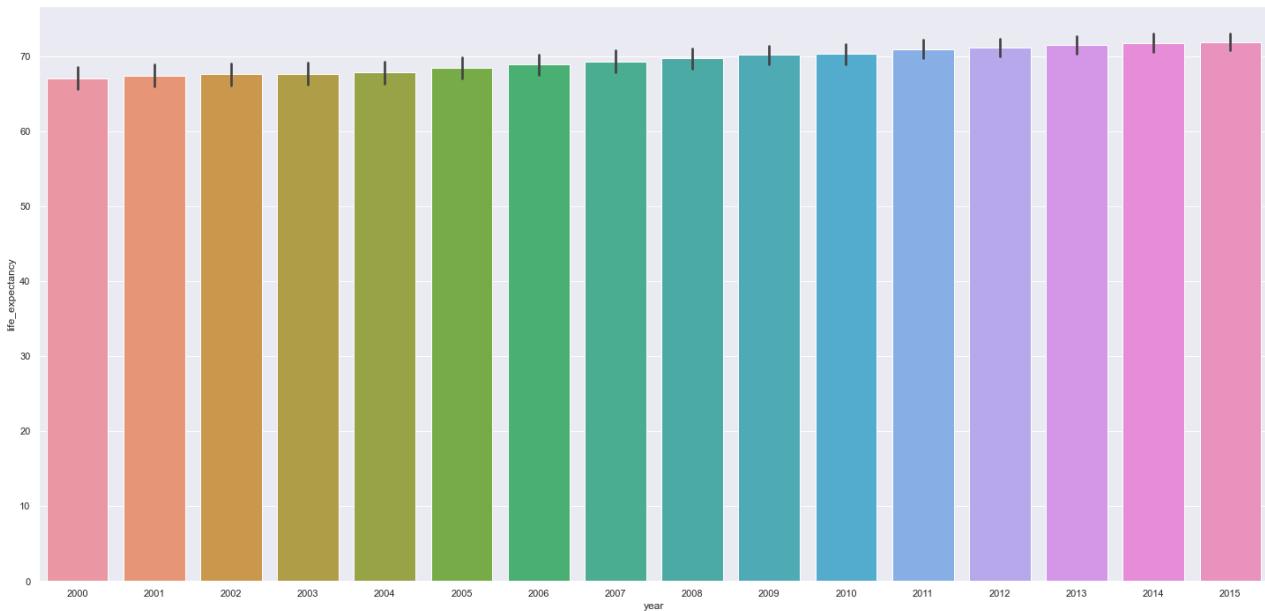
H3 - Como foi a evolução da expectativa de vida ao longo dos anos

```
In [323... #Life_expectancy mean per year  
aux1 = df2[['year','continent','life_expectancy']].groupby(['continent','year']).mean().reset_index()  
aux1.pivot(index='year',columns='continent',values='life_expectancy').plot();
```



In [324]:

```
#Life_expectancy mean per year
aux1 = df2[['year','country','life_expectancy']].groupby(['country','year']).mean().reset_index()
sns.barplot(x='year',y='life_expectancy',data=aux1);
```



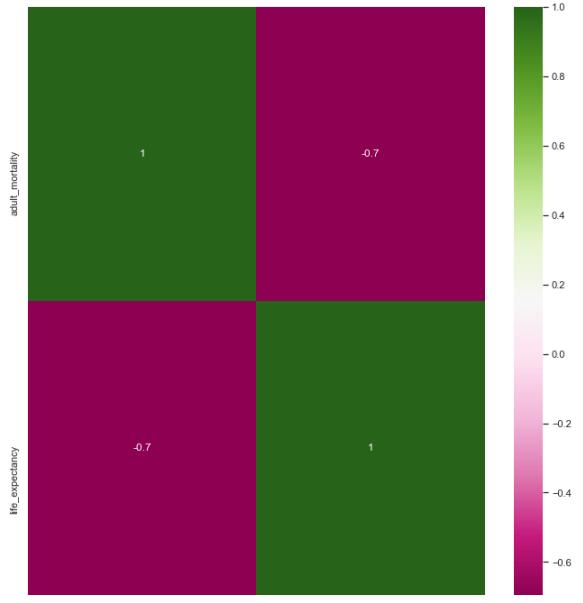
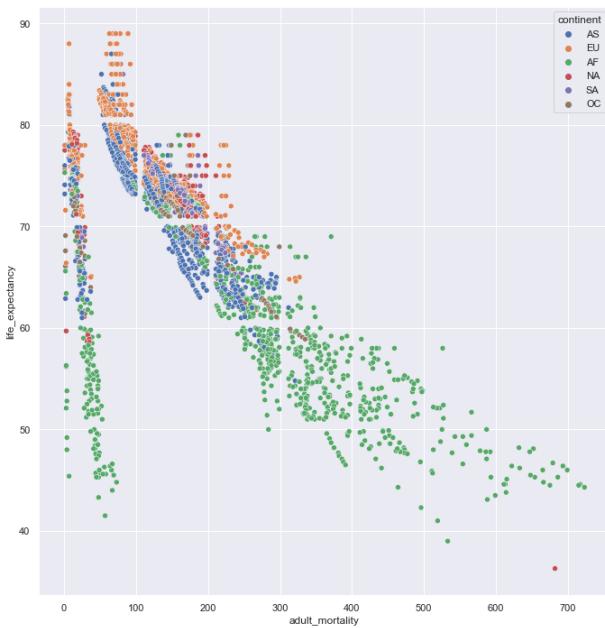
H4 - Como as taxas de mortalidade de bebês e adultos afetam a expectativa de vida?

In [325]:

```
aux2 = df2[['adult_mortality','life_expectancy','continent']]

plt.subplot( 1, 2, 1 )
sns.scatterplot(x='adult_mortality',y='life_expectancy',hue='continent',data=df2)

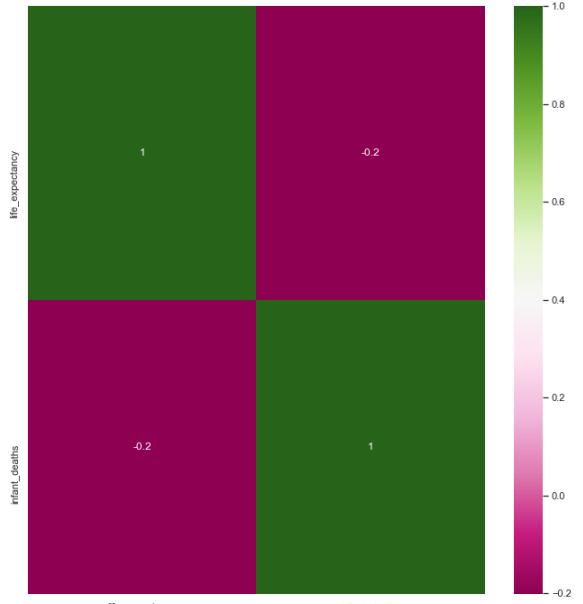
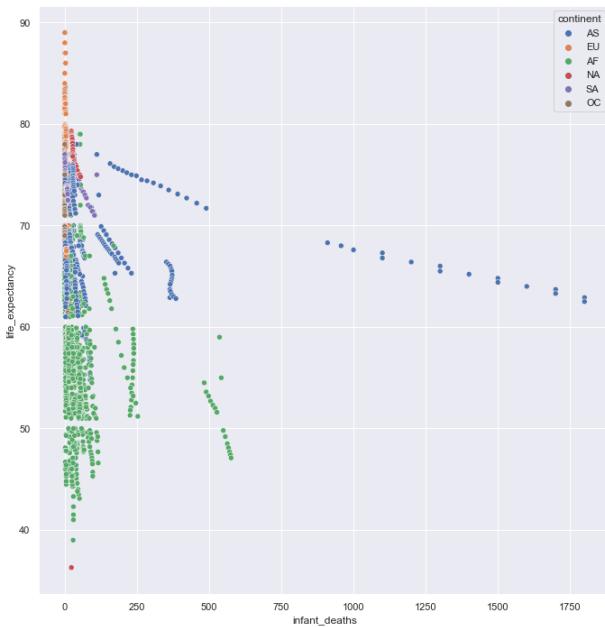
plt.subplot( 1, 2, 2 )
sns.heatmap(aux2.corr(method='pearson'),annot=True,cmap="PiYG");
```



```
In [326]: aux3 = df2[['life_expectancy','infant_deaths','continent']]
```

```
plt.subplot( 1, 2, 1 )
sns.scatterplot(x='infant_deaths',y='life_expectancy',hue='continent', data=df2);

plt.subplot( 1, 2, 2 )
sns.heatmap(aux3.corr(method='pearson'), annot=True, cmap="PiYG");
```

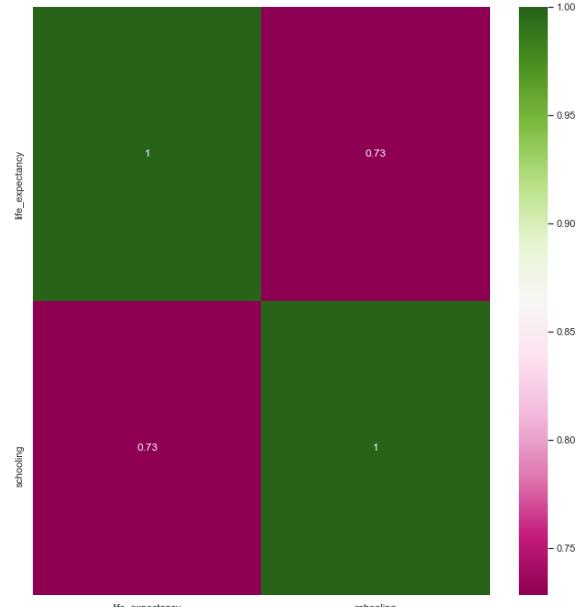
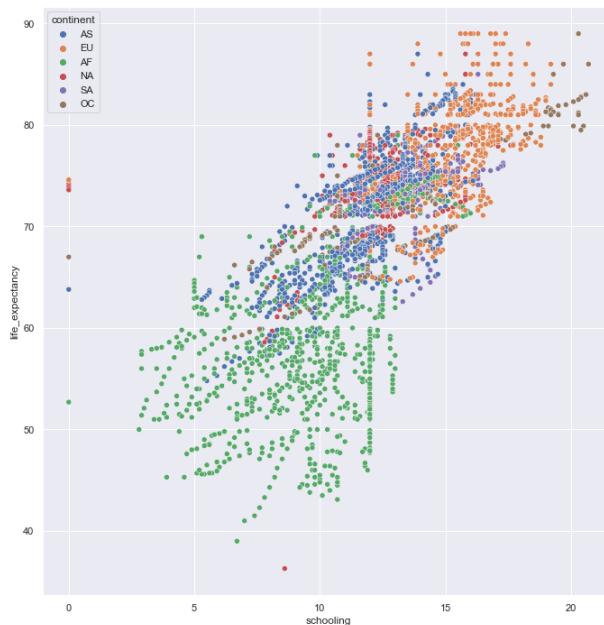


H5 - Qual é o impacto da escolaridade na expectativa de vida?

```
In [327]: aux4 = df2[['life_expectancy','schooling','continent']]
```

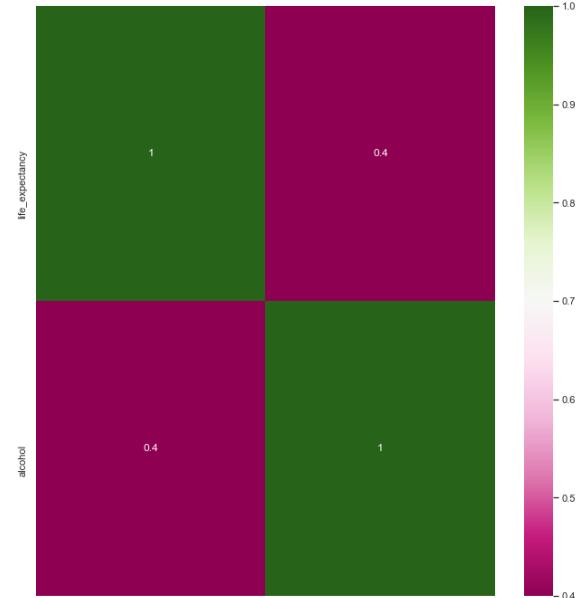
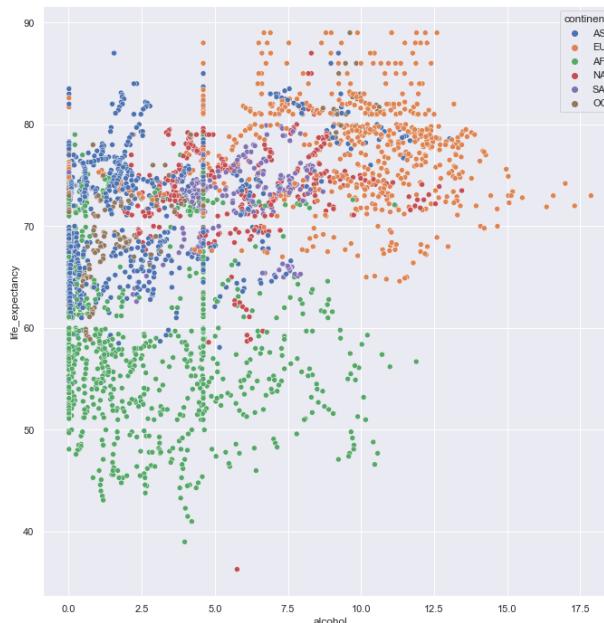
```
plt.subplot( 1, 2, 1 )
sns.scatterplot(x='schooling',y='life_expectancy',hue='continent', data=df2);

plt.subplot( 1, 2, 2 )
sns.heatmap(aux4.corr(method='pearson'), annot=True, cmap="PiYG");
```



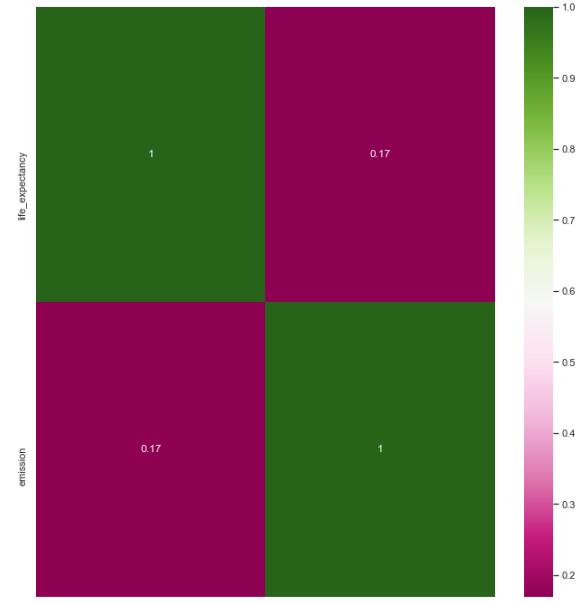
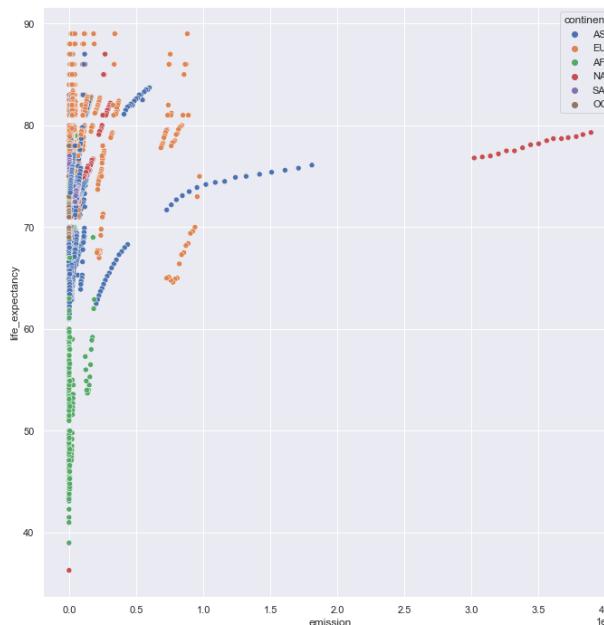
H6 - A expectativa de vida tem uma relação positiva ou negativa com o consumo de álcool?

```
In [328]: aux5 = df2[['life_expectancy','alcohol','continent']]
plt.subplot( 1, 2, 1 )
sns.scatterplot(x='alcohol',y='life_expectancy',hue='continent', data=df2);
plt.subplot( 1, 2, 2 )
sns.heatmap(aux5.corr(method='pearson'),annot=True,cmap="PiYG");
```



H7 - Paises mais inquinados representam uma expectativa de vida menor?

```
In [329]: aux6 = df2[['life_expectancy','emission','continent']]
plt.subplot( 1, 2, 1 )
sns.scatterplot(x='emission',y='life_expectancy',hue='continent', data=df2);
plt.subplot( 1, 2, 2 )
sns.heatmap(aux6.corr(method='pearson'),annot=True,cmap="PiYG");
```



H8 - Qual é o impacto da cobertura de imunização na expectativa de vida?

```
In [330]: cols_drop = ['adult_mortality', 'infant_deaths',  
    'alcohol', 'percentage_expenditure', 'bmi',  
    'under_five_deaths', 'total_expenditure', 'gdp', 'pop_total', 'thinness_1_19_years',  
    'thinness_5_9_years', 'income_composition_of_resources', 'schooling',  
    'emission', 'lat', 'long', 'year', 'emission_pop', 'perc_female', 'density', 'pop_male', 'pop_female']  
  
num_attributes2 = num_attributes.drop(cols_drop, axis=1)  
correlation = num_attributes2.corr( method='pearson' )  
sns.heatmap( correlation, annot=True );
```



```
In [331]: aux7 = df2[['life_expectancy', 'measles', 'continent']]  
aux8 = df2[['life_expectancy', 'polio', 'continent']]  
aux9 = df2[['life_expectancy', 'hepatitis_b', 'continent']]  
aux0 = df2[['life_expectancy', 'diphtheria', 'continent']]  
  
plt.subplot( 4, 2, 1 )  
sns.scatterplot(x='measles',y='life_expectancy',hue='continent', data=df2);  
  
plt.subplot( 4, 2, 2 )  
sns.heatmap(aux7.corr(method='pearson'), annot=True, cmap="PiYG");  
  
plt.subplot( 4, 2, 3 )  
sns.scatterplot(x='polio',y='life_expectancy',hue='continent', data=df2);
```

```

plt.subplot( 4, 2, 4 )
sns.heatmap(aux8.corr(method='pearson'), annot=True, cmap="PiYG");

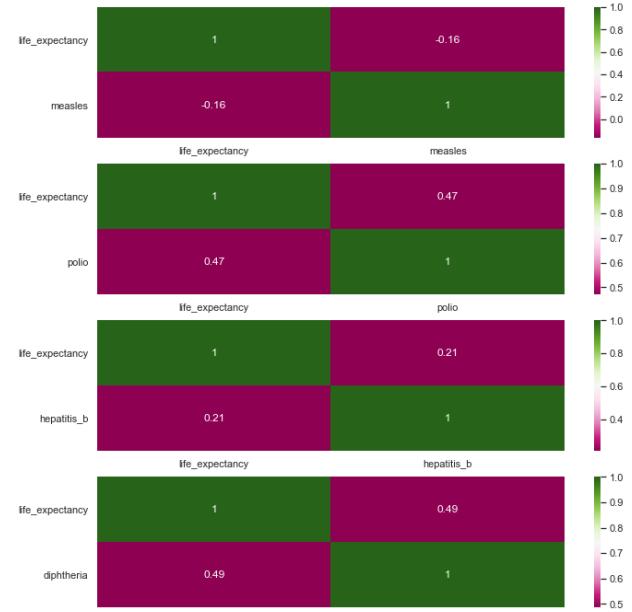
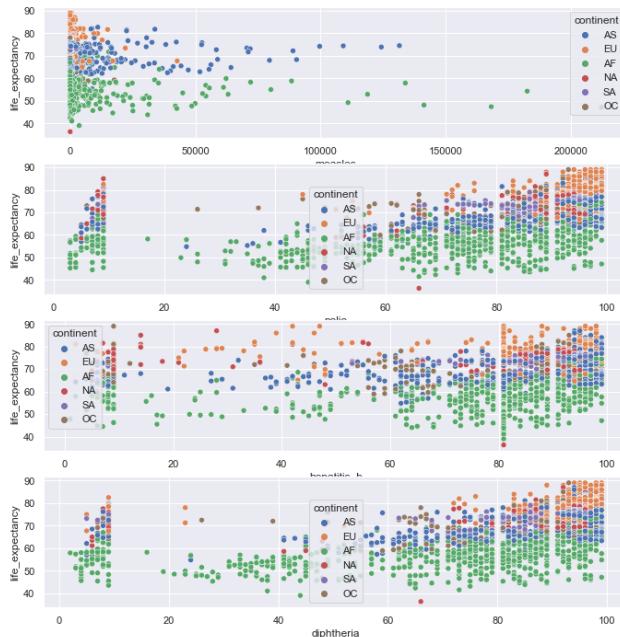
plt.subplot( 4, 2, 5 )
sns.scatterplot(x='hepatitis_b',y='life_expectancy',hue='continent', data=df2);

plt.subplot( 4, 2, 6 )
sns.heatmap(aux9.corr(method='pearson'), annot=True, cmap="PiYG");

plt.subplot( 4, 2, 7 )
sns.scatterplot(x='diphtheria',y='life_expectancy',hue='continent', data=df2);

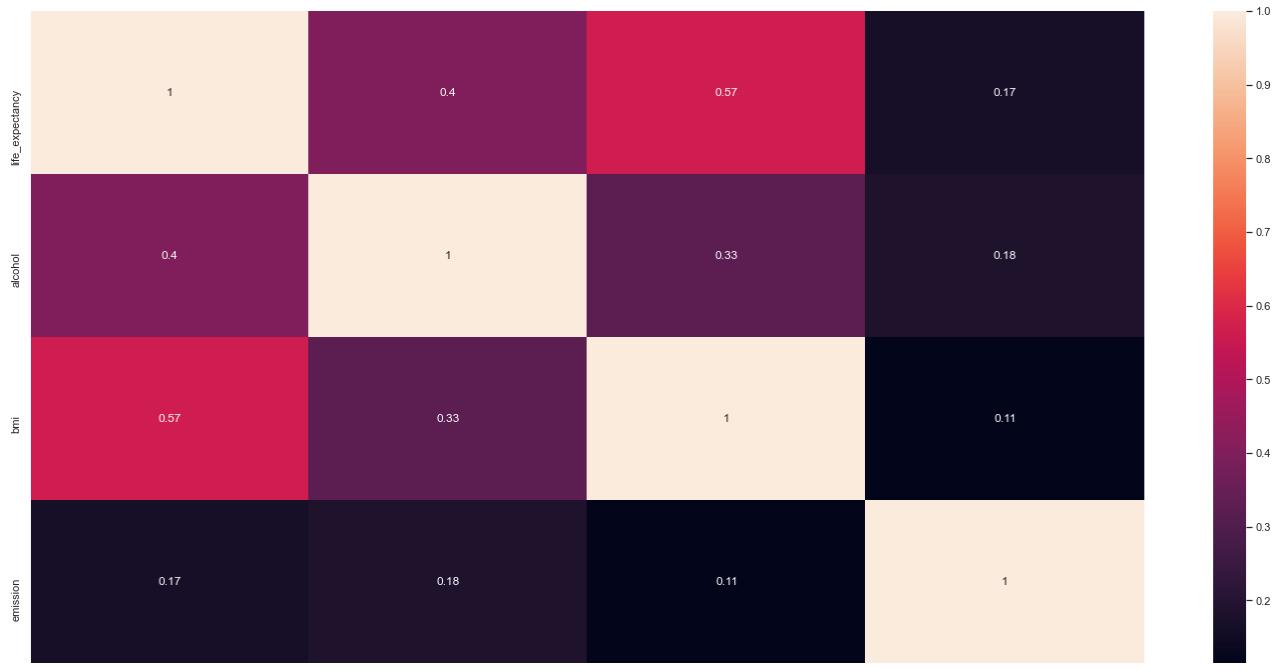
plt.subplot( 4, 2, 8 )
sns.heatmap(aux0.corr(method='pearson'), annot=True, cmap="PiYG");

```



H9 - A expectativa de vida tem correlação positiva ou negativa com hábitos alimentares, estilo de vida, exercícios, fumo, bebida alcoólica etc.

```
In [332]: num_attributes1 = df2[['life_expectancy', 'alcohol', 'bmi', 'emission']]
correlation = num_attributes1.corr( method='pearson' )
sns.heatmap( correlation, annot=True );
```



```
In [ ]:
```

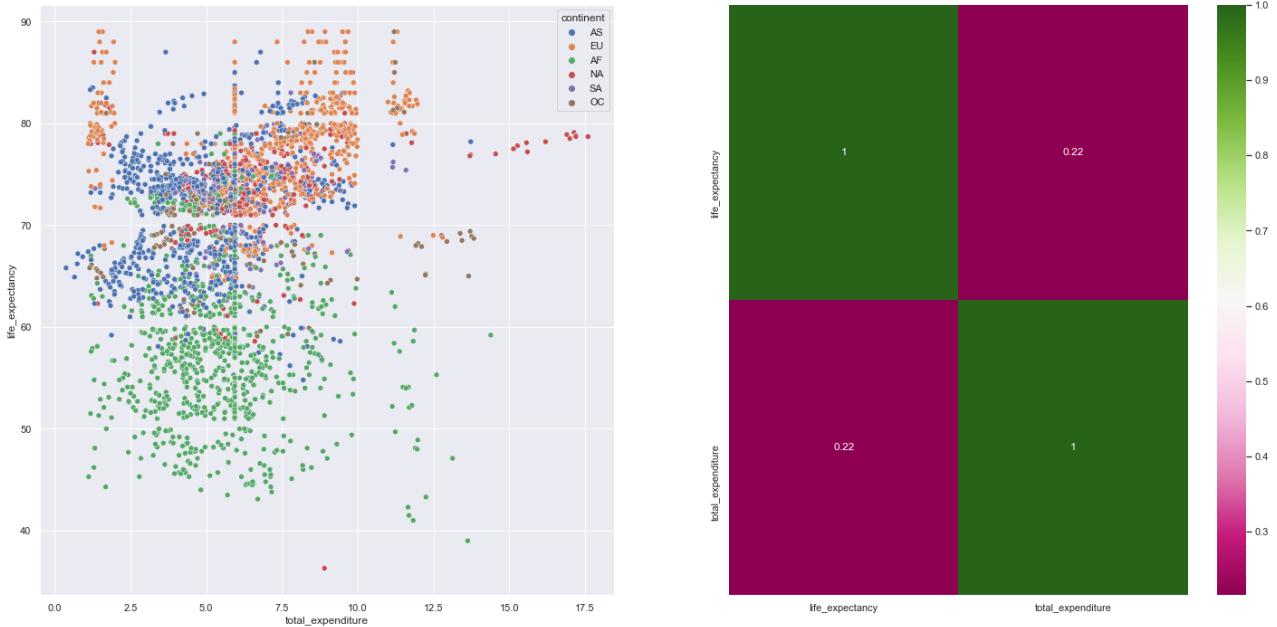
H10 - Qual a relação da expectativa de vida com a despesa de saúde?

In [333...]

```
auxc = df2[['life_expectancy','total_expenditure','continent']]

plt.subplot( 1, 2, 1 )
sns.scatterplot(x='total_expenditure',y='life_expectancy',hue='continent', data=auxc);

plt.subplot( 1, 2, 2 )
sns.heatmap(auxc.corr(method='pearson'), annot=True, cmap="PiYG");
```



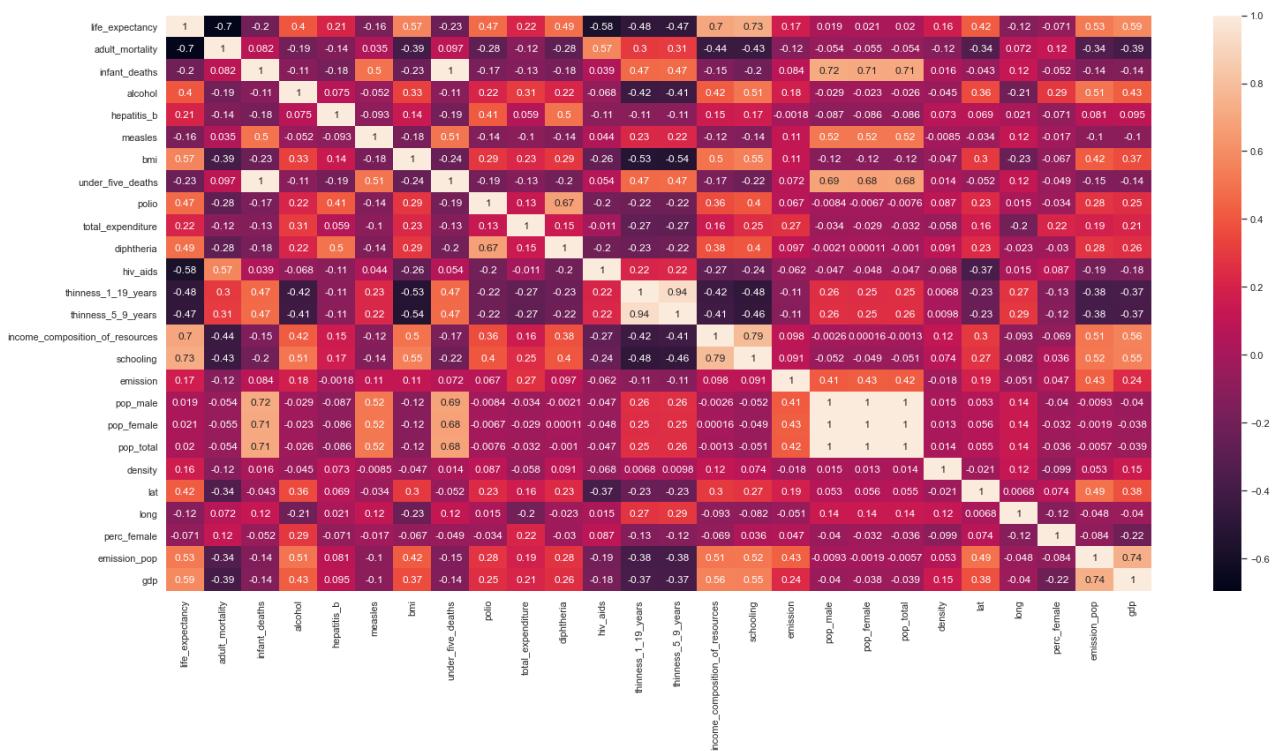
Numerical Heatmap

In [334...]

```
df3=df2.copy()
num_attributes4 = df3.select_dtypes( include=['int64', 'float64'] )
cat_attributes4 = df3.select_dtypes( exclude=['int64', 'float64', 'datetime64[ns]' ] )
```

In [335...]

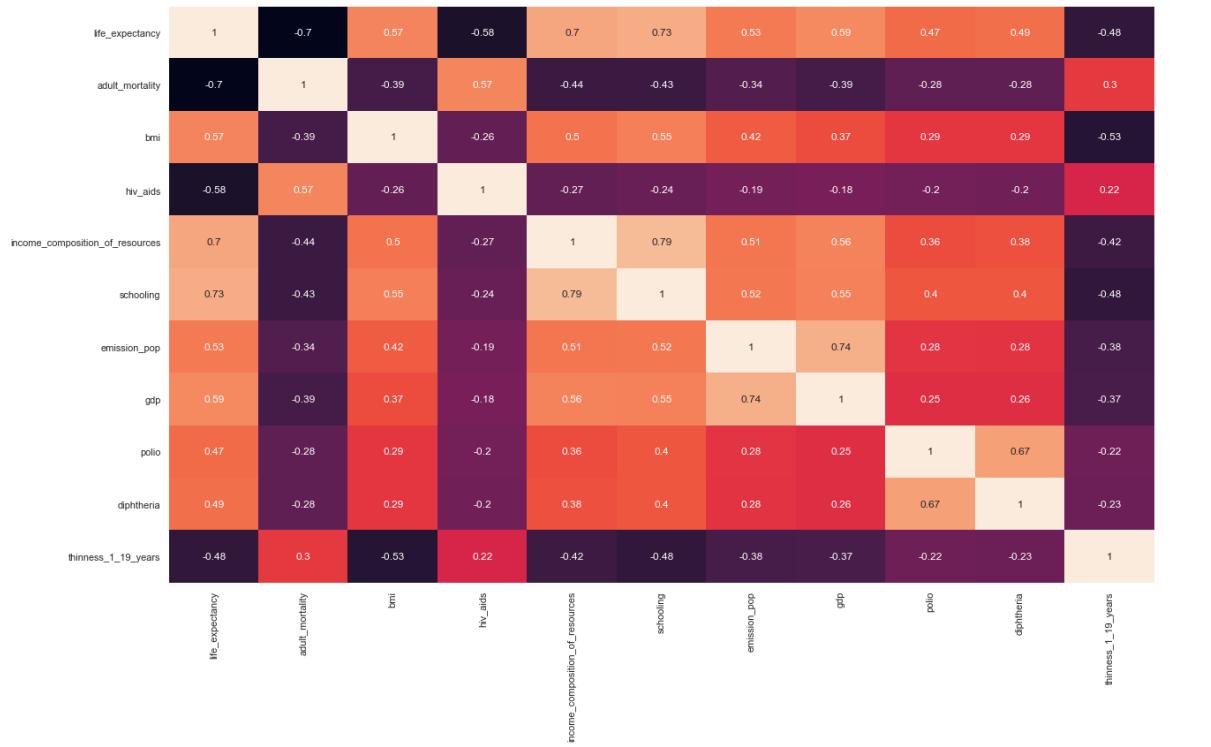
```
correlation = num_attributes4.corr( method='pearson' )
sns.heatmap( correlation, annot=True );
```



In [336...]

```
num_attributes5=num_attributes4[['life_expectancy','adult_mortality','bmi','hiv_aids','income_composition_of_resou
correlation = num_attributes5.corr( method='pearson' )
```

```
sns.heatmap( correlation, annot=True );
```

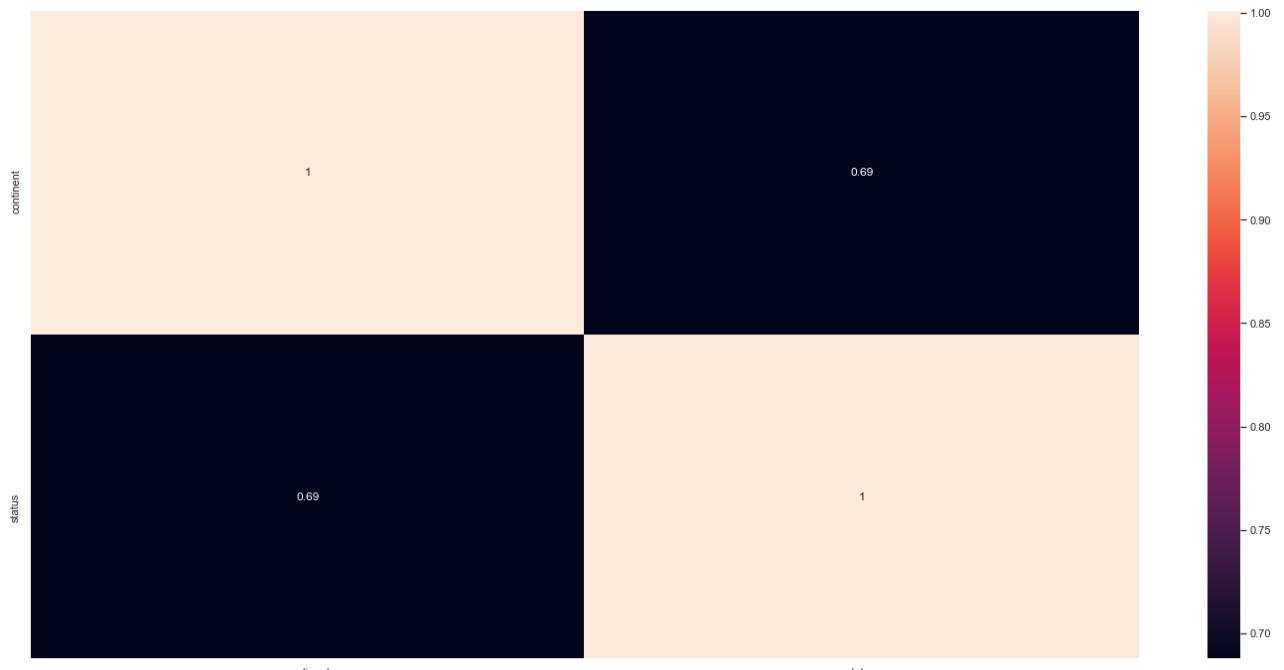


Categorical Heatmap

```
In [337...]
# # # Calculate cramer V
a = df2.select_dtypes( include='object' )

a1 = cramer_v( a['continent'], a['continent'] )
a2 = cramer_v( a['continent'], a['status'] )
a3 = cramer_v( a['status'], a['continent'] )
a4 = cramer_v( a['status'], a['status'] )

# # Final dataset
d = pd.DataFrame( {'continent': [a1, a2], 'status': [a3, a4]} )
d = d.set_index( d.columns )
sns.heatmap( d, annot=True );
```



Preprocessing Data

In [338]: `jupyter_settings()`

Populating the interactive namespace from numpy and matplotlib

In [339]: `df3=df2.copy()`

In [340]: `df3.info()`

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2872 entries, 0 to 2937
Data columns (total 31 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   country          2872 non-null    object  
 1   year              2872 non-null    object  
 2   status             2872 non-null    object  
 3   life_expectancy   2872 non-null    float64 
 4   adult_mortality   2872 non-null    float64 
 5   infant_deaths     2872 non-null    int64   
 6   alcohol            2872 non-null    float64 
 7   hepatitis_b       2872 non-null    float64 
 8   measles            2872 non-null    int64   
 9   bmi                2872 non-null    float64 
 10  under_five_deaths 2872 non-null    int64   
 11  polio              2872 non-null    float64 
 12  total_expenditure 2872 non-null    float64 
 13  diphtheria         2872 non-null    float64 
 14  hiv_aids          2872 non-null    float64 
 15  thinness_1_19_years 2872 non-null    float64 
 16  thinness_5_9_years 2872 non-null    float64 
 17  income_composition_of_resources 2872 non-null    float64 
 18  schooling          2872 non-null    float64 
 19  emission            2872 non-null    float64 
 20  pop_male            2872 non-null    float64 
 21  pop_female          2872 non-null    float64 
 22  pop_total            2872 non-null    float64 
 23  density              2872 non-null    float64 
 24  lat                 2872 non-null    float64 
 25  long                2872 non-null    float64 
 26  continent            2872 non-null    object  
 27  code                2872 non-null    object  
 28  perc_female          2872 non-null    float64 
 29  emission_pop         2872 non-null    float64 
 30  gdp                 2872 non-null    float64 
dtypes: float64(23), int64(3), object(5)
memory usage: 798.0+ KB
```

In [341]: `rs = RobustScaler()
mms = MinMaxScaler()
le=LabelEncoder()`

Robust Scaler

Variaveis a serem rescaladas: (problema con range)

- emission
- gdp
- infant_deaths
- percentage_expenditure
- measles
- under_five_deaths
- pop_male
- pop_female
- pop_total
- density
- perc_female
- emission_pop

In [342]: `df3['emission'] = rs.fit_transform(df3[['emission']].values)
df3['gdp'] = rs.fit_transform(df3[['gdp']].values)
df3['infant_deaths'] = rs.fit_transform(df3[['infant_deaths']].values)
df3['measles'] = rs.fit_transform(df3[['measles']].values)
df3['under_five_deaths'] = rs.fit_transform(df3[['under_five_deaths']].values)
df3['pop_male'] = rs.fit_transform(df3[['pop_male']].values)`

```
df3['pop_female'] = rs.fit_transform( df3[['pop_female']].values )
df3['pop_total'] = rs.fit_transform( df3[['pop_total']].values )
df3['density'] = rs.fit_transform( df3[['density']].values )
df3['perc_female'] = rs.fit_transform( df3[['perc_female']].values )
df3['emission_pop'] = rs.fit_transform( df3[['emission_pop']].values )
```

Transformation

Order Encoder

- year

```
In [343...]: a=['2000','2001','2002','2003','2004','2005','2006','2007','2008','2009','2010','2011','2012','2013','2014','2015']
b=np.arange(1,17,1)
year_dict = dict(zip(a, b))
df3['year']=df3['year'].map(year_dict)
df3.head()
```

Out[343...]:

	country	year	status	life_expectancy	adult_mortality	infant_deaths	alcohol	hepatitis_b	measles	bmi	under_five_deaths
0	Afghanistan	16	Developing	65.00	263.00	2.81	0.01	65.00	3.17	19.10	3.04
1	Afghanistan	15	Developing	59.90	271.00	2.90	0.01	62.00	1.32	18.60	3.15
2	Afghanistan	14	Developing	59.90	268.00	3.00	0.01	64.00	1.15	18.10	3.27
3	Afghanistan	13	Developing	59.50	272.00	3.14	0.01	67.00	7.72	17.60	3.42
4	Afghanistan	12	Developing	59.20	275.00	3.24	0.01	68.00	8.35	17.20	3.58

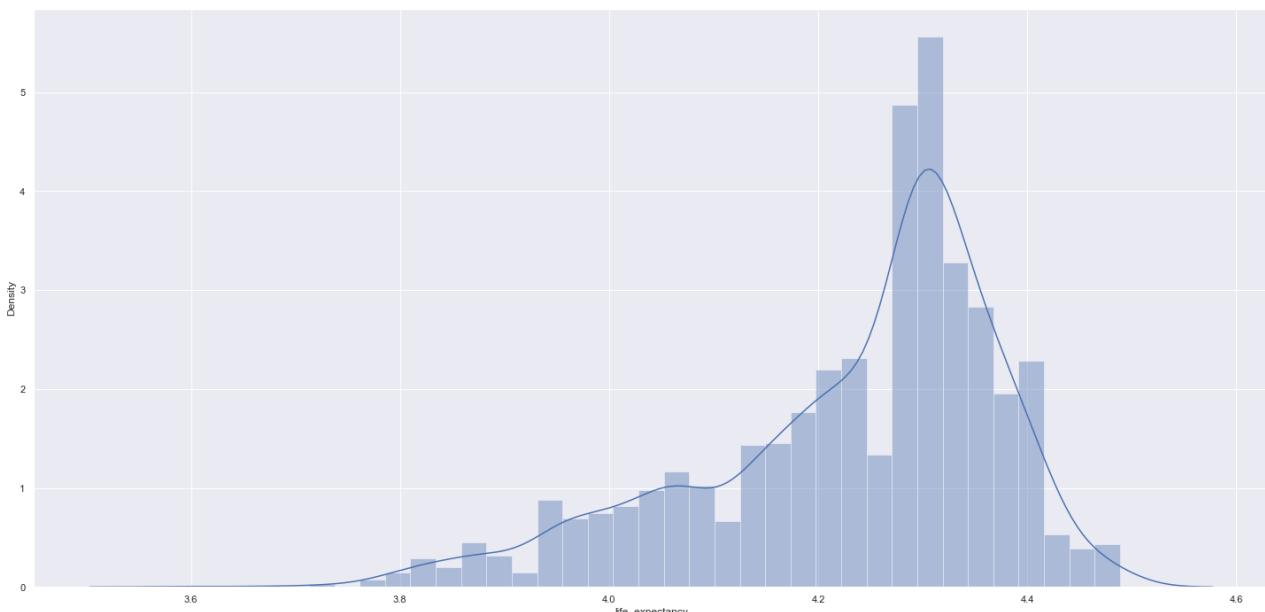
Label Encoder

- status
- continent
- code
- country

```
In [344...]: # Label Encoder
df3['status']=le.fit_transform(df3['status'])
df3['continent']=le.fit_transform(df3['continent'])
df3['country']=le.fit_transform(df3['country'])
df3['code']=le.fit_transform(df3['code'])
```

Response variable Transformation

```
In [345...]: df3['life_expectancy']=np.log(df3['life_expectancy'])
sns.distplot(df3['life_expectancy']);
```



Features Selection

```
In [346... df4=df3.copy()
df4.head()
```

```
Out[346... country year status life_expectancy adult_mortality infant_deaths alcohol hepatitis_b measles bmi under_five_deaths polio
0 0 16 1 4.17 263.00 2.81 0.01 65.00 3.17 19.10 3.04 6.00
1 0 15 1 4.09 271.00 2.90 0.01 62.00 1.32 18.60 3.15 58.00
2 0 14 1 4.09 268.00 3.00 0.01 64.00 1.15 18.10 3.27 62.00
3 0 13 1 4.09 272.00 3.14 0.01 67.00 7.72 17.60 3.42 67.00
4 0 12 1 4.08 275.00 3.24 0.01 68.00 8.35 17.20 3.58 68.00
```

```
In [347... # # # Remocao de atributos irrelevantes
cols_drop = ['code','lat','long','pop_total','perc_female']
df4=df4.drop(cols_drop, axis=1)
```

Split dataframe into training and test dataset

```
# training dataset
X_train = df4[df4['year'] < 14]
y_train = X_train['life_expectancy']

# test dataset
X_test = df4[df4['year'] >= 14]
y_test = X_test['life_expectancy']

print( 'Training Min Date: {}'.format( X_train['year'].min() ) )
print( 'Training Max Date: {}'.format( X_train['year'].max() ) )

print( '\nTest Min Date: {}'.format( X_test['year'].min() ) )
print( 'Test Max Date: {}'.format( X_test['year'].max() ) )
```

Training Min Date: 1
 Training Max Date: 13

Test Min Date: 14
 Test Max Date: 16

```
In [349... X_train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2332 entries, 3 to 2937
Data columns (total 26 columns):
 #   Column           Non-Null Count Dtype
 ---  -- 
 0   country          2332 non-null  int32
 1   year              2332 non-null  int64
 2   status             2332 non-null  int32
 3   life_expectancy    2332 non-null  float64
 4   adult_mortality    2332 non-null  float64
 5   infant_deaths     2332 non-null  float64
 6   alcohol            2332 non-null  float64
 7   hepatitis_b        2332 non-null  float64
 8   measles            2332 non-null  float64
 9   bmi                 2332 non-null  float64
 10  under_five_deaths  2332 non-null  float64
 11  polio               2332 non-null  float64
 12  total_expenditure  2332 non-null  float64
 13  diphtheria         2332 non-null  float64
 14  hiv_aids           2332 non-null  float64
 15  thinness_1_19_years 2332 non-null  float64
 16  thinness_5_9_years   2332 non-null  float64
 17  income_composition_of_resources 2332 non-null  float64
 18  schooling           2332 non-null  float64
 19  emission             2332 non-null  float64
 20  pop_male             2332 non-null  float64
 21  pop_female           2332 non-null  float64
 22  density               2332 non-null  float64
 23  continent             2332 non-null  int32
 24  emission_pop         2332 non-null  float64
 25  gdp                  2332 non-null  float64
```

dtypes: float64(22), int32(3), int64(1)
memory usage: 464.6 KB

Subset Select- Boruta

In [350...]

```
# ## training and test dataset for Boruta
X_train_n = X_train.drop(['life_expectancy', 'year'], axis=1).values #retorna uma matrix
y_train_n = y_train.values.ravel()

# ## define RandomForestRegressor
rf = RandomForestRegressor( n_jobs=-1 ) # cria as arvores em paralelo

# # define Boruta
boruta = BorutaPy( rf, n_estimators='auto', verbose=2, random_state=15 ).fit( X_train_n, y_train_n )# devo passar
```

```
Iteration: 1 / 100
Confirmed: 0
Tentative: 24
Rejected: 0
Iteration: 2 / 100
Confirmed: 0
Tentative: 24
Rejected: 0
Iteration: 3 / 100
Confirmed: 0
Tentative: 24
Rejected: 0
Iteration: 4 / 100
Confirmed: 0
Tentative: 24
Rejected: 0
Iteration: 5 / 100
Confirmed: 0
Tentative: 24
Rejected: 0
Iteration: 6 / 100
Confirmed: 0
Tentative: 24
Rejected: 0
Iteration: 7 / 100
Confirmed: 0
Tentative: 24
Rejected: 0
Iteration: 8 / 100
Confirmed: 21
Tentative: 3
Rejected: 0
Iteration: 9 / 100
Confirmed: 21
Tentative: 2
Rejected: 1
Iteration: 10 / 100
Confirmed: 21
Tentative: 2
Rejected: 1
Iteration: 11 / 100
Confirmed: 21
Tentative: 2
Rejected: 1
Iteration: 12 / 100
Confirmed: 21
Tentative: 2
Rejected: 1
Iteration: 13 / 100
Confirmed: 21
Tentative: 2
Rejected: 1
Iteration: 14 / 100
Confirmed: 21
Tentative: 2
Rejected: 1
Iteration: 15 / 100
Confirmed: 21
Tentative: 2
Rejected: 1
Iteration: 16 / 100
Confirmed: 22
Tentative: 1
Rejected: 1
Iteration: 17 / 100
Confirmed: 22
```

```
Tentative:      1
Rejected:      1
Iteration:    18 / 100
Confirmed:     22
Tentative:      1
Rejected:      1
Iteration:    19 / 100
Confirmed:     23
Tentative:      0
Rejected:      1
```

BorutaPy finished running.

```
Iteration:    20 / 100
Confirmed:     23
Tentative:      0
Rejected:      1
```

```
In [351...]: cols_selected = boruta.support_.tolist()

## best features
X_train_fs = X_train.drop( [ 'life_expectancy','year'], axis=1 )
cols_selected_boruta = X_train_fs.iloc[:, cols_selected].columns.to_list()

# ## not selected boruta
cols_not_selected_boruta = list( np.setdiff1d( X_train_fs.columns, cols_selected_boruta ) )
```

```
In [352...]: cols_selected_boruta
```

```
Out[352...]: ['country',
'adult_mortality',
'infant_deaths',
'alcohol',
'hepatitis_b',
'measles',
'bmi',
'under_five_deaths',
'polio',
'total_expenditure',
'diphtheria',
'hiv_aids',
'thinness_1_19_years',
'thinness_5_9_years',
'income_composition_of_resources',
'schooling',
'emission',
'pop_male',
'pop_female',
'density',
'continent',
'emission_pop',
'gdp']
```

```
In [353...]: cols_not_selected_boruta
```

```
Out[353...]: ['status']
```

```
In [354...]: features_selection= [
'country',
'adult_mortality',
'infant_deaths',
'alcohol',
'hepatitis_b',
'measles',
'bmi',
'under_five_deaths',
'polio',
'total_expenditure',
'diphtheria',
'hiv_aids',
'thinness_1_19_years',
'thinness_5_9_years',
'income_composition_of_resources',
'schooling',
'emission',
'pop_male',
'pop_female',
'density',
```

```

    'continent',
    'emission_pop',
    'gdp'
]

# columns to add
feat_to_add = ['life_expectancy', 'year']

features_selection_full = features_selection.copy()
features_selection_full.extend( feat_to_add )

```

In [355...]: features_selection_full

Out[355...]:

```

['country',
 'adult_mortality',
 'infant_deaths',
 'alcohol',
 'hepatitis_b',
 'measles',
 'bmi',
 'under_five_deaths',
 'polio',
 'total_expenditure',
 'diphtheria',
 'hiv_aids',
 'thinness_1_19_years',
 'thinness_5_9_years',
 'income_composition_of_resources',
 'schooling',
 'emission',
 'pop_male',
 'pop_female',
 'density',
 'continent',
 'emission_pop',
 'gdp',
 'life_expectancy',
 'year']

```

Machine Learning Modelling

In [356...]:

```

x_train = X_train[ features_selection ]
x_test = X_test[ features_selection ]

# Time Series Data Preparation
x_training = X_train[ features_selection_full ]

```

Linear regression model

In [357...]:

```

# model
lr = LinearRegression().fit( x_train, y_train )

# prediction
yhat_lr = lr.predict( x_test )

```

In [358...]:

```

print( 'Score Train: {}'.format( lr.score(x_train,y_train) ) )
print( 'Score Test: {}'.format( lr.score(x_test,y_test) ) )

# performance
lr_result = ml_error( 'Linear Regression', np.expm1( y_test ), np.expm1( yhat_lr ) )
lr_result

```

Score Train: 0.8394402986660525
Score Test: 0.8061552252295672

Out[358...]:

	Model Name	MAE	MAPE	RMSE
0	Linear Regression	2.83	0.04	3.71

In [359...]:

```

print("Descrição do modelo: ")
s = ["{0}: {1:0.5f}".format(a, v) for a, v in zip(features_selection_full, lr.coef_)]
print("w: {} b: {:.5f}".format(s, lr.intercept_))
print("Número de atributos usados: {}".format(np.sum(lr.coef_ != 0)))

```

Descrição do modelo:
w: ['country: 0.00008', 'adult_mortality: -0.00025', 'infant_deaths: 0.02676', 'alcohol: 0.00094', 'hepatitis_b: -

```
0.00022', 'measles: -0.00013', 'bmi: 0.00064', 'under_five_deaths: -0.02633', 'polio: 0.00044', 'total_expenditure: -0.00057', 'diphtheria: 0.00061', 'hiv_aids: -0.00884', 'thinness_1_19_years: 0.00015', 'thinness_5_9_years: 0.00139', 'income_composition_of_resources: 0.05992', 'schooling: 0.01000', 'emission: -0.00008', 'pop_male: 0.00366', 'pop_female: -0.00223', 'density: 0.00204', 'continent: 0.01292', 'emission_pop: 0.00044', 'gdp: 0.01424'] b: 3.98768
```

Número de atributos usados: 23

Linear Regression Model - Cross Validation

```
In [360... lr_result_cv = cross_validation( x_training, 3, 'Linear Regression', lr, verbose=False ) lr_result_cv
```

	Model Name	MAE CV	MAPE CV	RMSE CV
0	Linear Regression	3.17 +/- 0.206	0.05 +/- 0.004	4.12 +/- 0.173

Linear Regression Regularized Model - Lasso

```
In [361... # model lrr = Lasso( alpha=0.01 ).fit( x_train, y_train ) # prediction yhat_lrr = lrr.predict( x_test )
```

```
In [362... print( 'Score Train: {}'.format( lrr.score(x_train,y_train) ) ) print( 'Score Test: {}'.format( lrr.score(x_test,y_test) ) ) # performance lrr_result = ml_error( 'Linear Regression - Lasso', np.expm1( y_test ), np.expm1( yhat_lrr ) ) lrr_result
```

Score Train: 0.8236920836289436
Score Test: 0.8007282381425641

	Model Name	MAE	MAPE	RMSE
0	Linear Regression - Lasso	2.81	0.04	3.71

```
In [363... print("Descrição do modelo: ") s = ["{}: {:.0f}{}".format(a, v) for a, v in zip(features_selection_full, lrr.coef_)] print("w: {} b: {:.5f}{}".format(s, lrr.intercept_)) print("Número de atributos usados: {}".format(np.sum(lrr.coef_ != 0)))
```

Descrição do modelo:
w: ['country: 0.00006', 'adult_mortality: -0.00030', 'infant_deaths: -0.00000', 'alcohol: 0.00042', 'hepatitis_b: -0.00026', 'measles: -0.00028', 'bmi: 0.00085', 'under_five_deaths: -0.00170', 'polio: 0.00049', 'total_expenditure: -0.00000', 'diphtheria: 0.00075', 'hiv_aids: -0.00870', 'thinness_1_19_years: -0.00000', 'thinness_5_9_years: 0.00000', 'income_composition_of_resources: 0.00000', 'schooling: 0.01324', 'emission: 0.00000', 'pop_male: 0.00244', 'pop_female: 0.00000', 'density: 0.00178', 'continent: 0.00508', 'emission_pop: 0.00000', 'gdp: 0.00728'] b: 3.99654
Número de atributos usados: 15

Linear Regression Regularized Model - Cross Validation

```
In [364... lrr_result_cv = cross_validation( x_training, 3, 'Linear Regression - Lasso', lrr, verbose=False ) lrr_result_cv
```

	Model Name	MAE CV	MAPE CV	RMSE CV
0	Linear Regression - Lasso	3.04 +/- 0.163	0.05 +/- 0.003	4.02 +/- 0.199

Linear Regression Regularized Model - Ridge

```
In [365... # model lrri = Ridge( alpha=0.9 ).fit( x_train, y_train ) # prediction yhat_lrri = lrri.predict( x_test )
```

```
In [366... print( 'Score Train: {}'.format( lrri.score(x_train,y_train) ) ) print( 'Score Test: {}'.format( lrri.score(x_test,y_test) ) ) # performance
```

```
lrri_result = ml_error( 'Linear Regression - Ridge', np.expm1( y_test ), np.expm1( yhat_lrri ) )
lrri_result
```

Score Train: 0.8394387798012231
Score Test: 0.806042258039995

Out[366...]

	Model Name	MAE	MAPE	RMSE
0	Linear Regression - Ridge	2.83	0.04	3.71

In [367...]

```
print("Descrição do modelo: ")
s = ["{0}: {1:0.5f} ".format(a, v) for a, v in zip(features_selection_full, lrri.coef_)]
print("w: {} b: {:.5f}".format(s, lrri.intercept_))
print("Número de atributos usados: {}".format(np.sum(lrri.coef_ != 0)))
```

Descrição do modelo:
w: ['country: 0.00008', 'adult_mortality: -0.00025', 'infant_deaths: 0.02665', 'alcohol: 0.00093', 'hepatitis_b: -0.00022', 'measles: -0.00013', 'bmi: 0.00064', 'under_five_deaths: -0.02622', 'polio: 0.00044', 'total_expenditure: -0.00057', 'diphtheria: 0.00061', 'hiv_aids: -0.00884', 'thinness_1_19_years: 0.00014', 'thinness_5_9_years: 0.00139', 'income_composition_of_resources: 0.05852', 'schooling: 0.01005', 'emission: -0.00009', 'pop_male: 0.00350', 'pop_female: -0.00204', 'density: 0.00204', 'continent: 0.01291', 'emission_pop: 0.00046', 'gdp: 0.01426'] b: 3.98782

Número de atributos usados: 23

Linear Regression Regularized Model - Cross Validation

In [368...]

```
lrri_result_cv = cross_validation( x_training, 3, 'Linear Regression - Ridge', lrri, verbose=False )
lrri_result_cv
```

Out[368...]

	Model Name	MAE CV	MAPE CV	RMSE CV
0	Linear Regression - Ridge	3.16 +/- 0.201	0.05 +/- 0.004	4.12 +/- 0.166

Random Forest Regressor

In [369...]

```
RandomForestRegressor().get_params()
```

Out[369...]

```
{'bootstrap': True,
 'ccp_alpha': 0.0,
 'criterion': 'mse',
 'max_depth': None,
 'max_features': 'auto',
 'max_leaf_nodes': None,
 'max_samples': None,
 'min_impurity_decrease': 0.0,
 'min_impurity_split': None,
 'min_samples_leaf': 1,
 'min_samples_split': 2,
 'min_weight_fraction_leaf': 0.0,
 'n_estimators': 100,
 'n_jobs': None,
 'oob_score': False,
 'random_state': None,
 'verbose': 0,
 'warm_start': False}
```

In [370...]

```
# model
rf = RandomForestRegressor(bootstrap= True,
                           min_samples_leaf= 1,
                           min_samples_split= 2,
                           n_estimators=100,
                           n_jobs=-1,
                           ).fit( x_train, y_train )

# prediction
yhat_rf = rf.predict( x_test )
```

In [371...]

```
print( 'Score Train: {}'.format(rf.score(x_train,y_train) ) )
print( 'Score Test: {}'.format( rf.score(x_test,y_test) ) )

# performance
rf_result = ml_error( 'Random Forest Regressor', np.expm1( y_test ), np.expm1( yhat_rf ) )
rf_result
```

Score Train: 0.9956753889550254
Score Test: 0.9211404340414775

Out[371...]

	Model Name	MAE	MAPE	RMSE
0	Random Forest Regressor	1.45	0.02	2.29

Random Forest- Cross Validation

```
In [372...]: rf_result_cv = cross_validation( x_training, 3, 'Random Forest Regressor', rf, verbose=False )
rf_result_cv
```

	Model Name	MAE CV	MAPE CV	RMSE CV
0	Random Forest Regressor	1.63 +/- 0.093	0.03 +/- 0.001	2.49 +/- 0.124

XGBoost Regressor

```
In [373...]: xgb.XGBRegressor().get_params()
```

```
{'objective': 'reg:squarederror',
 'base_score': None,
 'booster': None,
 'colsample_bylevel': None,
 'colsample_bynode': None,
 'colsample_bytree': None,
 'enable_categorical': False,
 'gamma': None,
 'gpu_id': None,
 'importance_type': None,
 'interaction_constraints': None,
 'learning_rate': None,
 'max_delta_step': None,
 'max_depth': None,
 'min_child_weight': None,
 'missing': nan,
 'monotone_constraints': None,
 'n_estimators': 100,
 'n_jobs': None,
 'num_parallel_tree': None,
 'predictor': None,
 'random_state': None,
 'reg_alpha': None,
 'reg_lambda': None,
 'scale_pos_weight': None,
 'subsample': None,
 'tree_method': None,
 'validate_parameters': None,
 'verbosity': None}
```

```
In [374...]: # model
model_xgb = xgb.XGBRegressor( objective='reg:squarederror',
                               n_estimators=100,
                               max_depth=5,
                               subsample=0.7
                               ).fit( x_train, y_train )

# prediction
yhat_xgb = model_xgb.predict( x_test )

# performance
xgb_result = ml_error( 'XGBoost Regressor', np.expm1( y_test ), np.expm1( yhat_xgb ) )
xgb_result
```

	Model Name	MAE	MAPE	RMSE
0	XGBoost Regressor	1.66	0.02	2.44

XGBoost - Cross Validation

```
In [375...]: xgb_result_cv = cross_validation( x_training, 3, 'XGBoost Regressor', model_xgb, verbose=True )
xgb_result_cv
```

KFold Number: 3

KFold Number: 2

KFold Number: 1

Out[375...]

Model Name	MAE CV	MAPE CV	RMSE CV
0 XGBoost Regressor	1.75 +/- 0.023	0.03 +/- 0.001	2.57 +/- 0.033

Performance Metrics

```
In [376...]: modelling_result = pd.concat( [lr_result, lrr_result, lrri_result, rf_result, xgb_result] )
modelling_result.sort_values( 'RMSE' )
```

	Model Name	MAE	MAPE	RMSE
0	Random Forest Regressor	1.45	0.02	2.29
0	XGBoost Regressor	1.66	0.02	2.44
0	Linear Regression	2.83	0.04	3.71
0	Linear Regression - Ridge	2.83	0.04	3.71
0	Linear Regression - Lasso	2.81	0.04	3.71

Real Performance - Cross Validation

```
In [377...]: modelling_result_cv = pd.concat( [lr_result_cv, lrr_result_cv, lrri_result_cv, rf_result_cv, xgb_result_cv] )
modelling_result_cv.sort_values( 'RMSE CV' )
```

	Model Name	MAE CV	MAPE CV	RMSE CV
0	Random Forest Regressor	1.63 +/- 0.093	0.03 +/- 0.001	2.49 +/- 0.124
0	XGBoost Regressor	1.75 +/- 0.023	0.03 +/- 0.001	2.57 +/- 0.033
0	Linear Regression - Lasso	3.04 +/- 0.163	0.05 +/- 0.003	4.02 +/- 0.199
0	Linear Regression - Ridge	3.16 +/- 0.201	0.05 +/- 0.004	4.12 +/- 0.166
0	Linear Regression	3.17 +/- 0.206	0.05 +/- 0.004	4.12 +/- 0.173

Fine Tuning

Random search

```
In [378...]: # Search the best hyperparameters using Random Search
## Cria o random grid

n_estimators = [int(x) for x in np.linspace(start = 200, stop = 2000, num = 100)]

max_features = ['auto', 'sqrt']

max_depth = [int(x) for x in np.linspace(10, 110, num = 20)]
max_depth.append(None)

min_samples_split = [2, 5, 10]

min_samples_leaf = [1, 2, 4]

bootstrap = [True, False]

random_grid = {'n_estimators': n_estimators,
               'max_features': max_features,
               'max_depth': max_depth,
               'min_samples_split': min_samples_split,
               'min_samples_leaf': min_samples_leaf,
               'bootstrap': bootstrap}
```

```
In [379...]: final_result = pd.DataFrame()

n_interactions = 15

for i in range( n_interactions ):
    # choose values for parameters randomly
    hp = { k: np.random.choice( v, 1 )[0] for k, v in random_grid.items() }
    print( hp )

    # model
```

```

model_rf = RandomForestRegressor(
    n_estimators=hp['n_estimators'],
    max_depth=hp['max_depth'],
    max_features=hp['max_features'],
    min_samples_split=hp['min_samples_split'],
    min_samples_leaf=hp['min_samples_leaf'],
    bootstrap=hp['bootstrap']
)

# performance
result = cross_validation( x_training, 3, 'model_RandomForest', model_rf, verbose=False )
final_result = pd.concat( [final_result, result] )

final_result

```

{'n_estimators': 1581, 'max_features': 'auto', 'max_depth': 25, 'min_samples_split': 2, 'min_samples_leaf': 4, 'bootstrap': True}
{'n_estimators': 345, 'max_features': 'auto', 'max_depth': 25, 'min_samples_split': 10, 'min_samples_leaf': 1, 'bootstrap': True}
{'n_estimators': 781, 'max_features': 'auto', 'max_depth': 83, 'min_samples_split': 2, 'min_samples_leaf': 1, 'bootstrap': False}
{'n_estimators': 709, 'max_features': 'auto', 'max_depth': 15, 'min_samples_split': 5, 'min_samples_leaf': 2, 'bootstrap': False}
{'n_estimators': 1581, 'max_features': 'sqrt', 'max_depth': 36, 'min_samples_split': 2, 'min_samples_leaf': 2, 'bootstrap': False}
{'n_estimators': 1163, 'max_features': 'sqrt', 'max_depth': 94, 'min_samples_split': 5, 'min_samples_leaf': 4, 'bootstrap': False}
{'n_estimators': 1418, 'max_features': 'sqrt', 'max_depth': 62, 'min_samples_split': 10, 'min_samples_leaf': 4, 'bootstrap': False}
{'n_estimators': 1781, 'max_features': 'sqrt', 'max_depth': 73, 'min_samples_split': 5, 'min_samples_leaf': 2, 'bootstrap': True}
{'n_estimators': 1127, 'max_features': 'sqrt', 'max_depth': 10, 'min_samples_split': 2, 'min_samples_leaf': 4, 'bootstrap': True}
{'n_estimators': 1309, 'max_features': 'sqrt', 'max_depth': 110, 'min_samples_split': 2, 'min_samples_leaf': 4, 'bootstrap': False}
{'n_estimators': 418, 'max_features': 'auto', 'max_depth': None, 'min_samples_split': 2, 'min_samples_leaf': 2, 'bootstrap': False}
{'n_estimators': 854, 'max_features': 'auto', 'max_depth': 94, 'min_samples_split': 2, 'min_samples_leaf': 4, 'bootstrap': True}
{'n_estimators': 618, 'max_features': 'sqrt', 'max_depth': 15, 'min_samples_split': 2, 'min_samples_leaf': 4, 'bootstrap': True}
{'n_estimators': 1727, 'max_features': 'auto', 'max_depth': 20, 'min_samples_split': 5, 'min_samples_leaf': 4, 'bootstrap': True}
{'n_estimators': 200, 'max_features': 'auto', 'max_depth': None, 'min_samples_split': 2, 'min_samples_leaf': 4, 'bootstrap': True}

Out[379...]

	Model Name	MAE CV	MAPE CV	RMSE CV
0	model_RandomForest	1.68 +/- 0.067	0.03 +/- 0.001	2.54 +/- 0.106
0	model_RandomForest	1.69 +/- 0.068	0.03 +/- 0.001	2.55 +/- 0.111
0	model_RandomForest	2.04 +/- 0.041	0.03 +/- 0.001	3.06 +/- 0.093
0	model_RandomForest	2.08 +/- 0.014	0.03 +/- 0.0	3.07 +/- 0.041
0	model_RandomForest	1.52 +/- 0.023	0.02 +/- 0.001	2.32 +/- 0.062
0	model_RandomForest	1.58 +/- 0.019	0.02 +/- 0.001	2.38 +/- 0.045
0	model_RandomForest	1.6 +/- 0.009	0.02 +/- 0.001	2.39 +/- 0.043
0	model_RandomForest	1.59 +/- 0.019	0.02 +/- 0.001	2.39 +/- 0.029
0	model_RandomForest	1.72 +/- 0.026	0.03 +/- 0.001	2.52 +/- 0.02
0	model_RandomForest	1.58 +/- 0.014	0.02 +/- 0.001	2.38 +/- 0.044
0	model_RandomForest	2.07 +/- 0.008	0.03 +/- 0.0	3.06 +/- 0.035
0	model_RandomForest	1.68 +/- 0.065	0.03 +/- 0.001	2.54 +/- 0.104
0	model_RandomForest	1.69 +/- 0.048	0.03 +/- 0.002	2.49 +/- 0.013
0	model_RandomForest	1.68 +/- 0.076	0.03 +/- 0.001	2.54 +/- 0.116
0	model_RandomForest	1.69 +/- 0.068	0.03 +/- 0.001	2.54 +/- 0.109

Final Model

Final Model Random Forest

In [380...]

```
# model
rf = RandomForestRegressor(n_estimators= 1345,
                           max_features= 'sqrt',
```

```

        max_depth= 110,
        min_samples_split= 2,
        min_samples_leaf= 1,
        bootstrap= False
    ).fit( x_train, y_train )

# prediction
yhat_rf = rf.predict( x_test )

print( 'Score Test: {}'.format( rf.score(x_test,y_test) ) )

# performance
rf_result = ml_error( 'Random Forest Regressor', np.expm1( y_test ), np.expm1( yhat_rf ) )
rf_result

```

Score Test: 0.9354477721610472

Out[380...]

	Model Name	MAE	MAPE	RMSE
0	Random Forest Regressor	1.34	0.02	2.11

Translation and interpretation of the error

In [381...]

```

df5 = X_test[ features_selection_full ]
# rescale
df5['life_expectancy'] = np.expm1( df5['life_expectancy'] )
df5['predictions'] = np.expm1( yhat_rf )

```

Performance

In [382...]

```

b=['2000','2001','2002','2003','2004','2005','2006','2007','2008','2009','2010','2011','2012','2013','2014','2015']
a=np.arange(1,17,1)
year_dict = dict(zip(a, b))
df5[['year']] = df5[['year']].map(year_dict)
df5[['country','year','life_expectancy','predictions']].head()

```

Out[382...]

	country	year	life_expectancy	predictions
0	0	2015	64.00	60.46
1	0	2014	58.90	59.09
2	0	2013	58.90	59.16
16	1	2015	76.80	75.29
17	1	2014	76.50	75.43

In [383...]

```

# sum of predictions
df51 = df5[['country', 'predictions','life_expectancy']].groupby( 'country' ).mean().reset_index()

# # MAE and MAPE
df5_aux1 = df5[['country', 'life_expectancy', 'predictions']].groupby( 'country' ).apply( lambda x: mean_absolute_
df5_aux2 = df5[['country', 'life_expectancy', 'predictions']].groupby( 'country' ).apply( lambda x: mean_absolute_

# Merge
df5_aux3 = pd.merge( df5_aux1, df5_aux2, how='inner', on='country' )
df52 = pd.merge( df51, df5_aux3, how='inner', on='country' )

# # Scenarios
df52['worst_scenario'] = df52['predictions'] - df52['MAE']
df52['best_scenario'] = df52['predictions'] + df52['MAE']

# # "#Country with worst MAPE"
df52 = df52[['country', 'life_expectancy','predictions', 'worst_scenario', 'best_scenario', 'MAE', 'MAPE']]
df52.sort_values( 'MAPE', ascending=False ).head(5)

```

Out[383...]

	country	life_expectancy	predictions	worst_scenario	best_scenario	MAE	MAPE
143	143	50.03	55.72	50.03	61.41	5.69	0.11
3	3	50.73	55.47	50.73	60.22	4.74	0.09
50	50	77.67	71.35	65.04	77.67	6.31	0.08
162	162	70.10	75.69	70.10	81.27	5.59	0.08
166	166	62.60	57.96	53.32	62.60	4.64	0.07

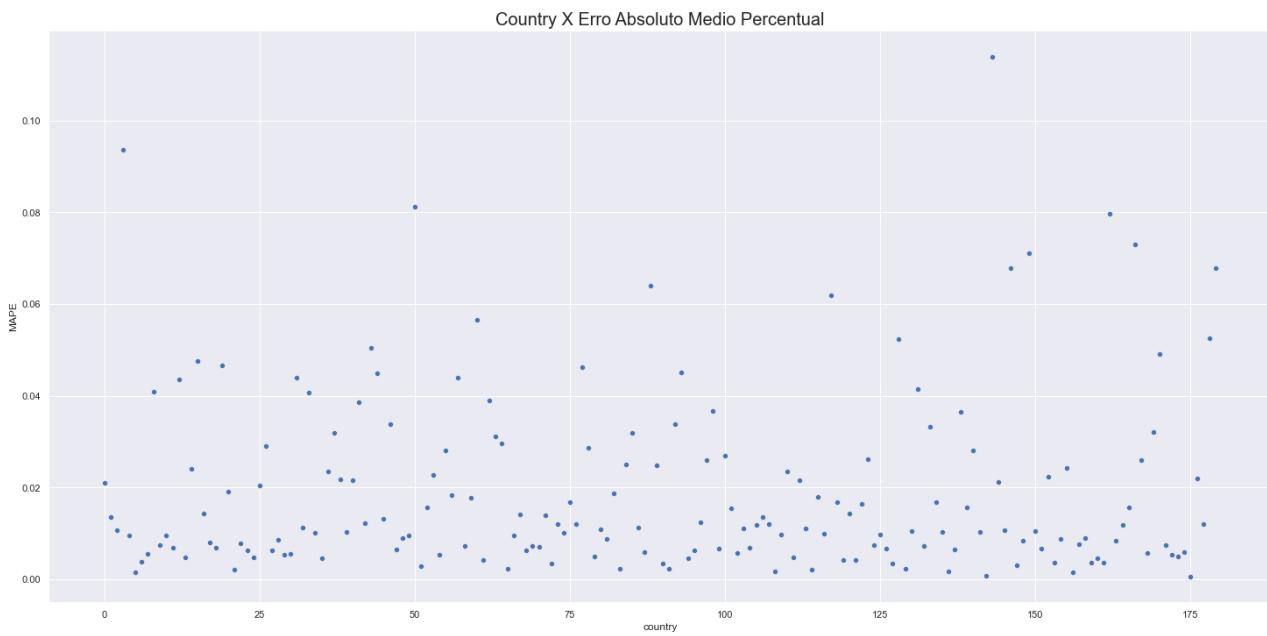
03/02/22, 00:38

TCC_Expectativa de vida

```
In [384...]: a=df2[['country','life_expectancy']].groupby('country').mean().reset_index()
a[a.index==143]
```

```
Out[384...]: country life_expectancy
143 Sierra Leone    46.11
```

```
In [385...]: sns.scatterplot(x='country',y='MAPE', data=df52)
plt.title('Country X Erro Absoluto Medio Percentual', fontsize=20);
```

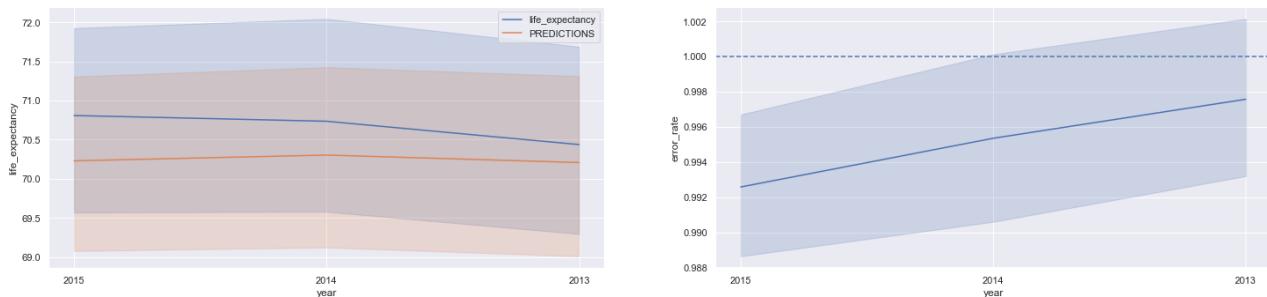


Machine Learning Performance

```
In [386...]: df5['error'] = df5['life_expectancy'] - df5['predictions']
df5['error_rate'] = df5['predictions'] / df5['life_expectancy']
```

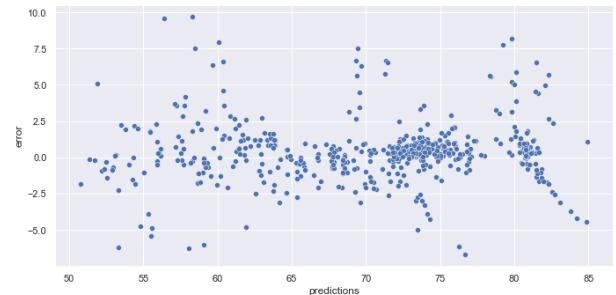
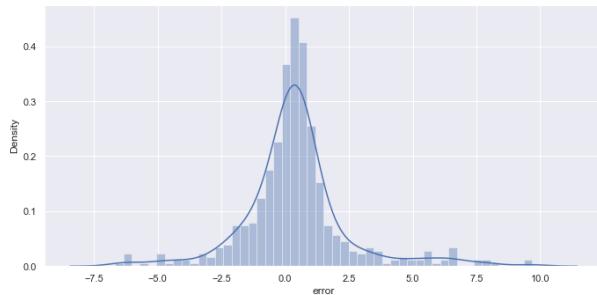
```
In [387...]: plt.subplot( 2, 2, 1 )
sns.lineplot( x='year', y='life_expectancy', data=df5, label='life_expectancy' )
sns.lineplot( x='year', y='predictions', data=df5, label='PREDICTIONS' );

plt.subplot( 2, 2, 2 )
sns.lineplot( x='year', y='error_rate', data=df5)
plt.axhline( 1, linestyle='--');
```



```
In [388...]: plt.subplot( 2, 2, 3 )
sns.distplot( df5['error'] )

plt.subplot( 2, 2, 4 )
sns.scatterplot( df5['predictions'], df5['error'] );
```



In []:

In []:

In []: