

Climate Change :

Importing the Necessary Libraries

```
1 import pandas as pd
2 from pandas import Series,DataFrame
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import bokeh as bk
7 sns.set_style('whitegrid')
8 %matplotlib inline
9 from sklearn import linear_model
10 from sklearn import metrics
11 from sklearn.svm import SVC, LinearSVC
12 from sklearn.ensemble import RandomForestClassifier
13 from sklearn.ensemble import RandomForestRegressor
14 from sklearn.neighbors import KNeighborsClassifier
15 from sklearn.naive_bayes import GaussianNB
16 import xgboost
```

Reading the Data using Pandas

```
1 # get global temp csv file as a DataFrame
2 df = pd.read_csv("C://Users//user//Downloads//GlobalTemperatures.csv", sep=',')
3 df.head()
```

	dt	LandAverageTemperature	LandAverageTemperatureUncertainty	LandMaxTemperature	LandMaxTemperatureUncertainty	LandMinTemperature	LandM
0	1750-01-01	3.034	3.574	NaN	NaN	NaN	NaN
1	1750-02-01	3.083	3.702	NaN	NaN	NaN	NaN
2	1750-03-01	5.626	3.076	NaN	NaN	NaN	NaN
3	1750-04-01	8.490	2.451	NaN	NaN	NaN	NaN
4	1750-05-01	11.573	2.072	NaN	NaN	NaN	NaN

Data has Multiple Null values , So we are Dropping Few Columns

```
1 # Delete all columns except dates, Land Average Temperatures, Land Average Temperature Uncertainty
2 df=df.drop(df.columns[3:], axis=1)
3
4 #Delete all empty lines
5 df=df.dropna()
6
```

Resizing the Datasets and only keeps Averages

```

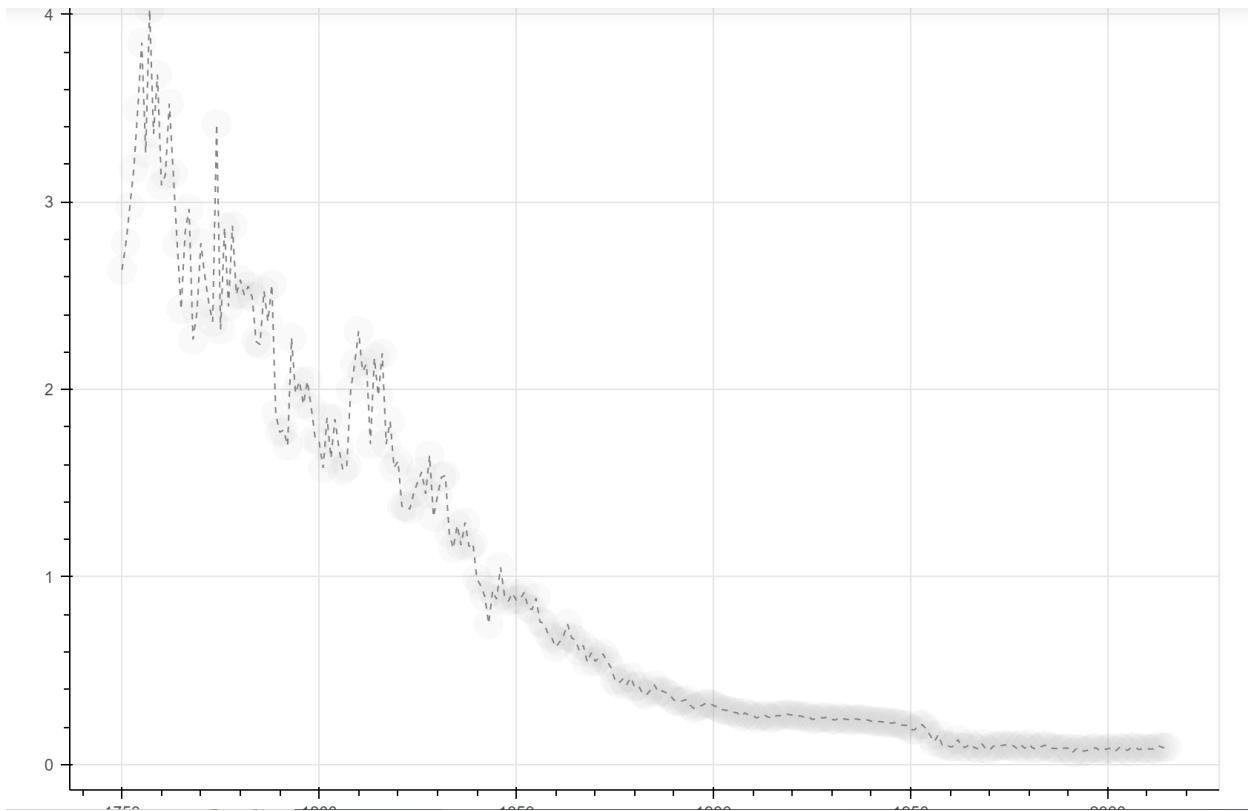
1 #Resizing the datasets, we only keep the year averages
2
3 date=df['dt'].apply(lambda x: x[:4])
4 years = np.unique(date)
5 mean_temps_global=[]
6 mean_temps_uncertainty=[]
7
8 for y in years:
9
10     mean_temps_global.append(df[y==date]['LandAverageTemperature'].mean())
11     mean_temps_uncertainty.append(df[y==date]['LandAverageTemperatureUncertainty'].mean())
12
13 #Defining a function that classifies the average temperatures
14
15 def heat(x):
16
17     if type(x) is str:
18         return x
19     elif x<=7:
20         return 3
21     elif x>7 and x<=8:
22         return 2
23     else:
24         return 1
25
26 #Building of a new, simpler dataframe
27
28 years=pd.to_numeric(years)
29
30 d1={'Dates': pd.Series(years), 'Temperatures':pd.Series(mean_temps_global), 'Uncertainty':pd.Series(mean_temps_uncertainty)}
31 df_global_mean=pd.DataFrame(d1)
32
33 #Adding the heat function
34
35 e=df_global_mean.applymap(heat)
36
37 d2={'Dates': pd.Series(years), 'Temperatures':pd.Series(mean_temps_global), 'Uncertainty':pd.Series(mean_temps_uncertainty),
38      'Heat':e['Temperatures']}
39
40 df_global_mean=pd.DataFrame(d2)
41
```

## Plotting the Uncertainty

```

1 #Now, we plot the 95% Uncertainty regarding the temperature evolution
2
3 from bokeh.plotting import figure, show
4 from bokeh.models import ColumnDataSource, Circle, HoverTool, CustomJS
5 from bokeh.io import output_notebook
6
7 output_notebook()
8
9 (x, y) = (df_global_mean['Dates'], df_global_mean['Uncertainty'])
10
11 # Basic plot setup
12 p = figure(width=800, height=600, toolbar_location=None, title='A clear decreasing of uncertainty')
13
14 p.line(x, y, line_dash=[4, 4], line_width=1, color='gray')
15
16
17 # Add a circle, that is visible only when selected
18 source = ColumnDataSource({'x': x, 'y': y})
19 invisible_circle = Circle(x='x', y='y', fill_color='gray', fill_alpha=0.05, line_color=None, size=20)
20 visible_circle = Circle(x='x', y='y', fill_color='firebrick', fill_alpha=0.5, line_color=None, size=20)
21 cr = p.add_glyph(source, invisible_circle, selection_glyph=visible_circle, nonselection_glyph=invisible_circle)
22
23 # Add a hover tool, that selects the circle
24 code = "source.set('selected', cb_data['index']);"
25 callback = CustomJS(args={'source': source}, code=code)
26 p.add_tools(HoverTool(tooltips=None, callback=callback, renderers=[cr], mode='hline'))
27
28 show(p)

```



## Loading the Min and Max Data

```

2
3 global_temps2 = pd.read_csv("C://Users//user//Downloads//GlobalTemperatures.csv", sep=',')
4
5 # Starting in 1800 when data is more significant (less NaN)
6
7 list_of_years_1800=years[50:]
8
9 # Compute Average, Max, Min Temperature by Year since 1800
10
11 date=global_temps2['dt'].apply(lambda x: x[:4])
12 list_of_years_1800=np.unique(date)
13 list_of_years_1800=list_of_years_1800[50:]
14
15 mean_temps_1800=[]
16 min_temps_1800=[]
17 max_temps_1800=[]
18
19 for y in list_of_years_1800:
20
21     mean_temps_1800.append(global_temps2[y==date]['LandAverageTemperature'].mean())
22     max_temps_1800.append(global_temps2[y==date]['LandMaxTemperature'].mean())
23     min_temps_1800.append(global_temps2[y==date]['LandMinTemperature'].mean())
24

```

```

25 #Building of a new, simpler dataframe
26
27 list_of_years_1800=pd.to_numeric(list_of_years_1800)
28
29 d1={'Average Temperature':pd.Series(mean_temps_1800), 'Dates': pd.Series(list_of_years_1800),
30      'Min Temperature':pd.Series(min_temps_1800), 'Max Temperature':pd.Series(max_temps_1800)}
31
32 global_temps_1800=pd.DataFrame(d1)
33
34 global_temps_1800=global_temps_1800[['Dates','Average Temperature','Min Temperature','Max Temperature']]
35

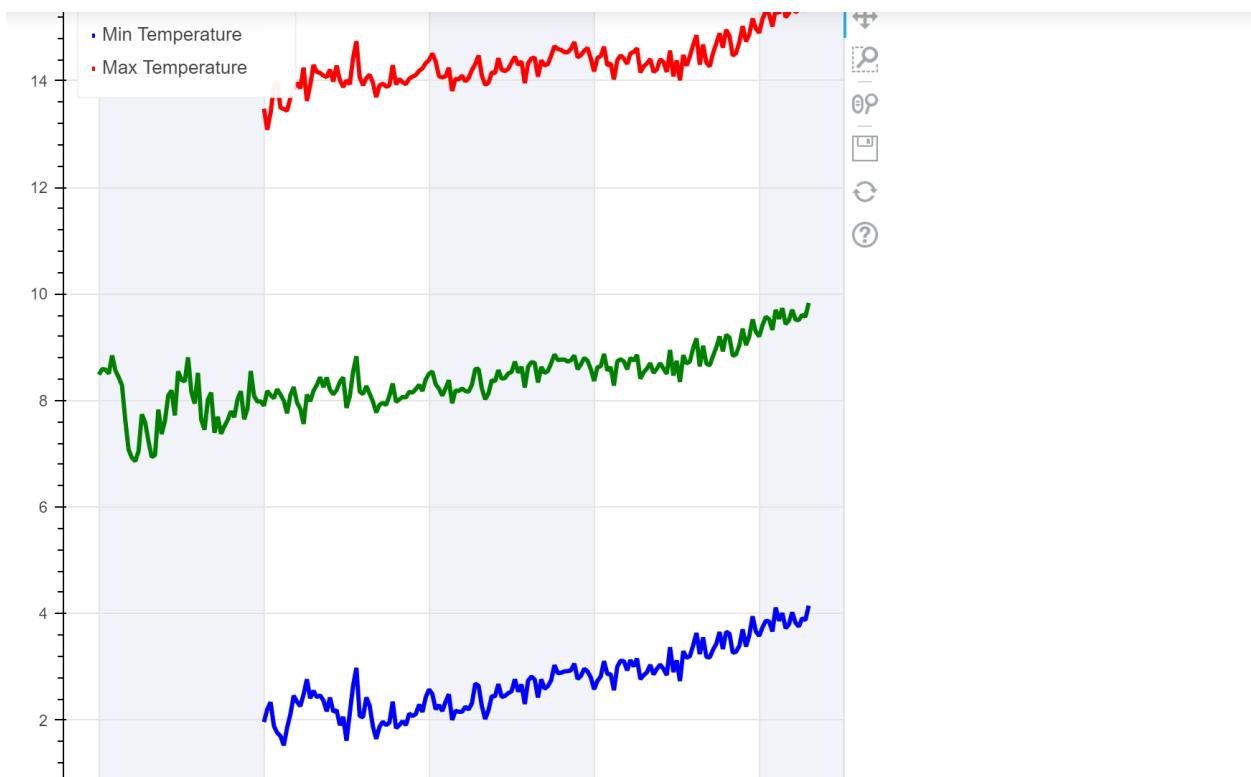
```

## Plotting the Tempurature

```

1 from bokeh.plotting import figure, show
2 from bokeh.models import ColumnDataSource, Circle, HoverTool, CustomJS
3 from bokeh.io import output_notebook
4
5 output_notebook()
6
7 x = global_temps_1800['Dates']
8 y = global_temps_1800['Average Temperature']
9 z = global_temps_1800['Min Temperature']
10 t = global_temps_1800['Max Temperature']
11
12 p = figure(plot_width=600, plot_height=600)
13
14
15 p.line(x, y, legend="Average Temperature", line_color="green", line_width=3)
16
17 p.line(x,z, legend="Min Temperature", line_color="blue", line_width=3)
18
19 p.line(x, t, legend="Max Temperature", line_color="red", line_width=3)
20
21 # change just some things about the x-grid
22 p.xgrid.band_fill_alpha = 0.05
23 p.xgrid.band_fill_color = "navy"
24
25 p.legend.location = "top_left"
26 p.legend.glyph_width = 2
27 p.legend.label_width=2
28
29
30 show(p)

```



## Linear Regression

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 from sklearn import linear_model
4
5
6 temp=global_temps_1800['Average Temperature']
7 time=global_temps_1800['Dates']
8
9
10 # The data sets
11 dates_X = time.iloc[0:,]
12
13 # The targets
14 temperature_y = temp.iloc[0:,]
15
16 # Create linear regression object
17 reg = linear_model.LinearRegression()
18
19 # Train the model using the data sets
20 reg.fit(dates_X.to_frame(),temperature_y.to_frame())
21
22 # The estimation
23 output_reg=reg.predict(dates_X.to_frame())
24 output_simple=np.copy(output_reg)
25 output=pd.DataFrame(index=pd.DataFrame(output_simple).index.values)
```

```
26 # The coefficients
27 print('Coefficients: \n', reg.coef_)
28 # The intercept
29 print('Intercept: \n', reg.intercept_)
30
31
32 # The mean squared error
33 print("Mean squared error: %.2f" % np.mean((reg.predict(dates_X.to_frame()) - temperature_y.to_frame()) ** 2))
34
35 # Explained variance score: 1 is perfect prediction
36 print('Variance score: %.2f' % reg.score(dates_X.to_frame(), temperature_y.to_frame()))
37
```

Coefficients:  
[[0.00721859]]  
Intercept:  
[-5.35578585]  
Mean squared error: 0.12  
Variance score: 0.64

### Regression Curve on Predicted Values

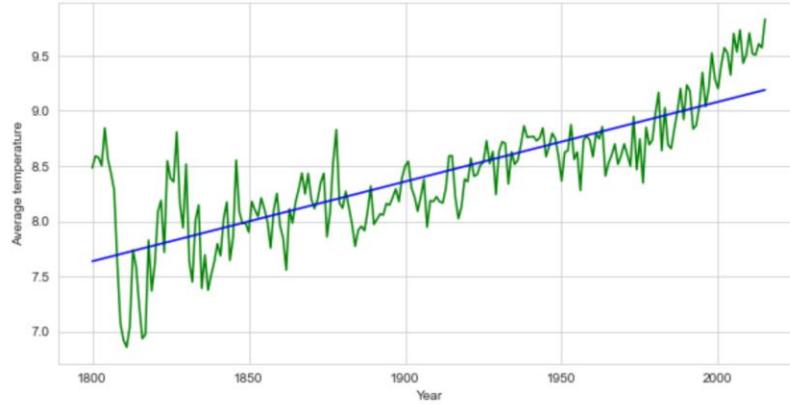
```
1 # Linear regression cuve on the predicted values
2
3 # Simple Linear Regression: t : Regression done on the values between 1800 and 1950 and testes between 1950 and 2016
4 fig, axes = plt.subplots(figsize=(10,5))
5 axes.set_ylabel('Average temperature')
6 axes.set_xlabel('Year')
7 plt.plot(dates_X, temperature_y, label='Annual Mean temperature', color='g')
8 plt.plot(dates_X,output_simple,label='Forecast Temperature by Regression',color='b')
```

```

1 # Linear regression curve on the predicted values
2
3 # Simple Linear Regression: t : Regression done on the values between 1800 and 1950 and testes between 1950 and 2016
4 fig, axes = plt.subplots(figsize=(10,5))
5 axes.set_ylabel('Average temperature')
6 axes.set_xlabel('Year')
7 plt.plot(dates_X, temperature_y, label='Annual Mean temperature', color='g')
8 plt.plot(dates_X,output_simple,label='Forecast Temperature by Regression',color='b')

```

[<matplotlib.lines.Line2D at 0x1ca24f527c0>]



## LINEAR MODEL

```

4 from sklearn import linear_model
5
6 #Define the input and output
7 temp=global_temps_1800['Average Temperature']
8 time=global_temps_1800['Dates']
9
10 # Split the data into training/testing sets
11 dates_X_train = time.iloc[0:150,:]
12 dates_X_test = time.iloc[150:,:]
13
14 # Split the targets into training/testing sets
15 temperature_y_train = temp.iloc[0:150,:]
16 temperature_y_test = temp.iloc[150:,:]
17
18 # Create linear regression object
19 reg = linear_model.LinearRegression()
20
21 # Train the model using the training sets
22 reg.fit(dates_X_train.to_frame(),temperature_y_train.to_frame())
23
24 # The coefficients
25 print('Coefficients: \n', reg.coef_)
26 # The intercept
27 print('Intercept: \n', reg.intercept_)

```

```

29 # The prediction
30 output_simple=reg.predict(dates_X_test.to_frame())
31 # The mean squared error
32 print("Mean squared error: %.2f"
33     % np.mean((reg.predict(dates_X_test.to_frame()) - temperature_y_test.to_frame()) ** 2))
34 # Explained variance score: 1 is perfect prediction
35 print('Variance score: %.2f' % reg.score(dates_X_test.to_frame(), temperature_y_test.to_frame()))
36

Coefficients:
[[0.00570766]]
Intercept:
[-2.53305848]
Mean squared error: 0.14
Variance score: 0.16

```

## Simple Linear Regression

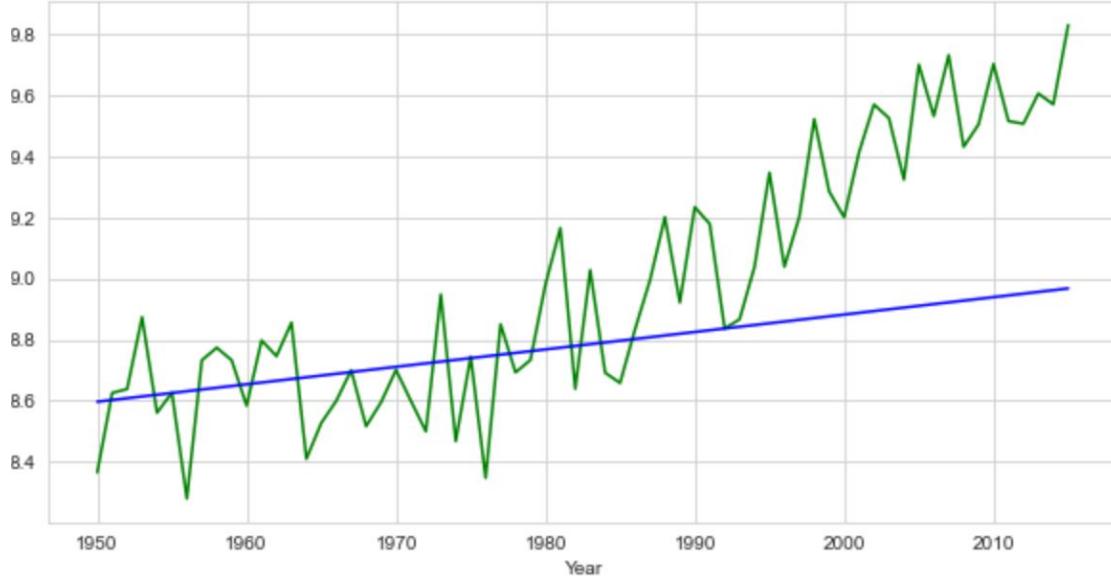
```

1 # Simple Linear Regression: t : Regression done on the values between 1800 and 1950 and testes between 1950 and 2016
2 fig, axes = plt.subplots(figsize=(10,5))
3 axes.set_ylabel('Average temperature')
4 axes.set_xlabel('Year')
5 plt.plot(dates_X_test, temperature_y_test, label='Annual Mean temperature', color='g')
6 plt.plot(dates_X_test,output_simple,label='Forecast Temperature by Regression',color='b')
7

[<matplotlib.lines.Line2D at 0x1ca26013ca0>]

```

[`<matplotlib.lines.Line2D at 0x1ca26013ca0>`]



## Polynomial Regression

## Polynomial Regression

```
1 time=pd.DataFrame()
2 time["Dates"] = global_temps_1800['Dates']
3 time["Dates^2"] = global_temps_1800['Dates']*global_temps_1800['Dates']
4 time["Dates^3"] = global_temps_1800['Dates']*global_temps_1800['Dates']*global_temps_1800['Dates']
5 temp=global_temps_1800['Average Temperature']
6
7
8 # The data sets
9 dates_X = time.iloc[0:,0]
10
11 # The targets
12 temperature_y = temp.iloc[0:,]
13
14 # Create Linear regression object
15 reg = linear_model.LinearRegression()
16
17 # Train the model using the training sets
18 reg.fit(time.iloc[0:,],temperature_y.to_frame())
19
20 # The prediction
21 output_reg=reg.predict(time.iloc[0:,])
22 output_multiple=np.copy(output_reg)
23 output=pd.DataFrame(index=pd.DataFrame(output_simple).index.values)
24 # The coefficients
25 print('Coefficients: \n', reg.coef_)
26 # The intercept
27 print('Intercept: \n', reg.intercept_)
28
29 # The mean squared error
30 print("Mean squared error: %.2f" % np.mean((output_reg - temperature_y.to_frame()) ** 2))
31
32 # Explained variance score: 1 is perfect prediction
33 print('Variance score: %.2f' % reg.score(time.iloc[0:,], temperature_y.to_frame()))
34
35 # Euclidian Distance
36 output[ 'Predicted Temperature']=pd.DataFrame(output_multiple)
37 output[ "Average Temperature"]=temp
```

```
Coefficients:
 [[ 1.54583760e+00 -8.43972200e-04  1.53912499e-07]]
Intercept:
 [-937.81103858]
Mean squared error: 0.10
Variance score: 0.69
```

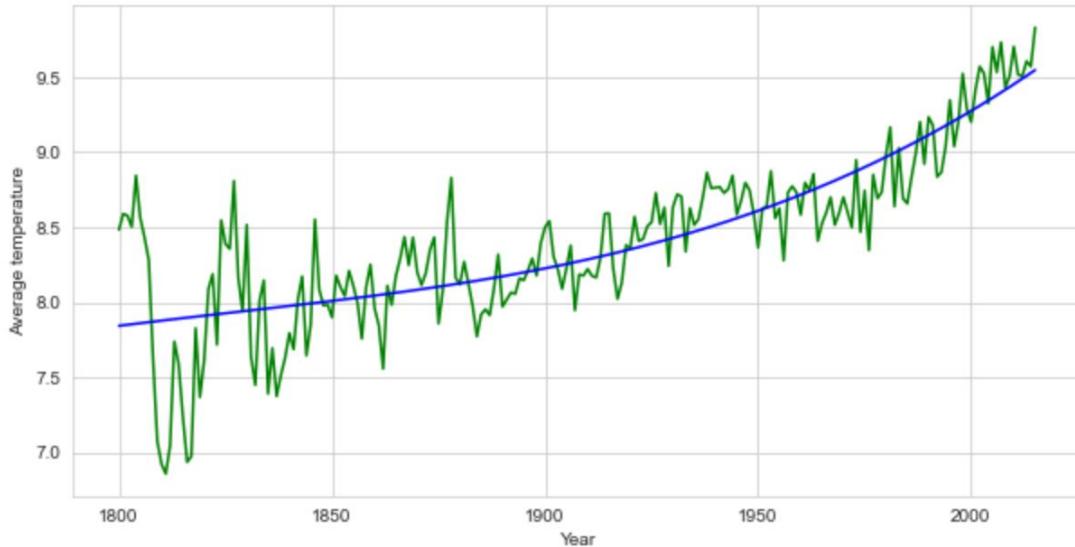
It has High Variance Score

```
1 #Regression done on the values between 1800 and 1950 and testes between 1950 and 2016
2 fig, axes = plt.subplots(figsize=(10,5))
3 axes.set_ylabel('Average temperature')
4 axes.set_xlabel('Year')
5 plt.plot(dates_X, temperature_y, label='Annual Mean temperature', color='g')
6 plt.plot(dates_X,output_multiple,label='Forecast Reg',color='b')
```

```
|: [<matplotlib.lines.Line2D at 0x1ca260a02b0>]
```



```
|: [<matplotlib.lines.Line2D at 0x1ca260a02b0>]
```



```
1 temp=global_temps_1800['Average Temperature']
2 time=pd.DataFrame()
3 time["Dates"]=global_temps_1800['Dates']
4 time["Dates^2"]=global_temps_1800['Dates']*global_temps_1800['Dates']
5 time["Dates^3"]=global_temps_1800['Dates']*global_temps_1800['Dates']*global_temps_1800['Dates']
6
7 # Split the data into training/testing sets
8 dates_X_train = time.iloc[0:150,0]
9 dates_X_test = time.iloc[150:,0]
10
11 # Split the targets into training/testing sets
12 temperature_y_train = temp.iloc[0:150,:]
13 temperature_y_test = temp.iloc[150:,:]
14
15 # Create Linear regression object
16 reg = linear_model.LinearRegression()
17
18 # Train the model using the training sets
19 reg.fit(time.iloc[0:150,:],temperature_y_train.to_frame())
20
21 # The coefficients
22 print('Coefficients: \n', reg.coef_)
23 # The intercept
24 print('Intercept: \n', reg.intercept_)
25
26 # The prediction
27 output_multiple=reg.predict(time.iloc[150:,:])
```

```

29 |     print("Mean squared error: %.2f"
30 |         % np.mean((reg.predict(time.iloc[150:,:]) - temperature_y_test.to_frame()) ** 2))
31 | # Explained variance score: 1 is perfect prediction
32 | print('Variance score: %.2f' % reg.score(time.iloc[150:,:], temperature_y_test.to_frame()))

```

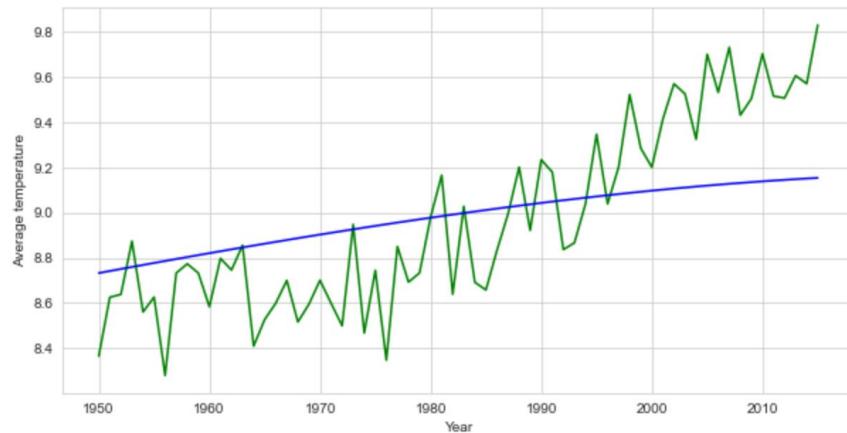
Coefficients:  
 $[-3.37849536e+00 \quad 1.75645409e-03 \quad -3.03542634e-07]$   
Intercept:  
 $[2168.61282323]$   
Mean squared error: 0.10  
Variance score: 0.42

```

1 | fig, axes = plt.subplots(figsize=(10,5))
2 | axes.set_ylabel('Average temperature')
3 | axes.set_xlabel('Year')
4 |
5 | plt.plot(dates_X_test, temperature_y_test, label='Annual Mean temperature', color='g')
6 | plt.plot(dates_X_test, output_multiple, label='Forecast Reg', color='b')
7 |

```

: [`<matplotlib.lines.Line2D at 0x1ca26102250>`]

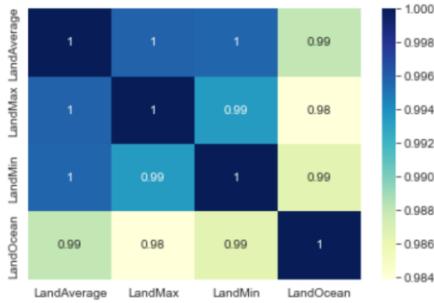


## Plotting the Correlation Matrix

## Correlation Matrix

```
[]: #Correlation Matrix for the following features : LandAverageTemperature, LandMaxTemperature, LandMinTemperature, LandAnd
1 global_temperatures=pd.read_csv("C://Users//user//Downloads//GlobalTemperatures.csv", sep=',')
2 global_temperatures['dt'] = pd.to_datetime(global_temperatures['dt'])
3 global_temperatures=global_temperatures.drop(['LandAverageTemperatureUncertainty','LandMaxTemperatureUncertainty','Land
4 global_temperatures=global_temperatures.dropna()
5 corr=global_temperatures.corr()
6 sns.heatmap(corr,xticklabels=['LandAverage','LandMax','LandMin','LandOcean'], yticklabels=['LandAverage','LandMax','Land
7
8
```

t[26]: <AxesSubplot:>



## Random Forest

### Train Data

```
[1] # Train
2 train_df_copy=pd.DataFrame.copy(global_temperatures[global_temperatures['dt'].map(lambda x: x.year)<1980])
3 train_df=pd.DataFrame.copy(train_df_copy)
4 train_df.drop(train_df.index[:1], inplace=True)
5 train_df_copy.drop(train_df.index[-1], inplace=True)
6
7 # To predict each average Land temperature we use the previous year's temperatures (Landandocean, min, max)
8 for j in (train_df.index):
9     for column in train_df.columns[1:-1]:
10         train_df.loc[j,column]=train_df_copy.loc[j-1,column]
11 forest = RandomForestRegressor(n_estimators=10)
12 train_df=train_df.drop(['dt'],axis=1)
13 forest = forest.fit(train_df.drop(["LandAverageTemperature"],axis=1),train_df["LandAverageTemperature"])
```

### Test Data

```
[1] # Test
2 test_df_copy=pd.DataFrame.copy(global_temperatures[global_temperatures['dt'].map(lambda x: x.year)>=1980])
3 test_df=pd.DataFrame.copy(test_df_copy)
4 test_df.drop(test_df.index[:1], inplace=True)
5 test_df_copy.drop(test_df.index[-1], inplace=True)
6
7 # To predict each average Land temperature we use the previous year's temperatures (Landandocean, min , max)
8 for j in (test_df.index):
9     for column in test_df.columns[1:-1]:
10         test_df.loc[j,column]=test_df_copy.loc[j-1,column]
11 test_df=test_df.drop(['dt'],axis=1)
12 output = forest.predict(test_df.drop(["LandAverageTemperature"],axis=1))
```

```

1 #Adding the predicted temperature to our output DF
2 test_df['Predicted_temp']=output

1 # Distance computing
2 MED=abs(test_df['LandAverageTemperature']-test_df['Predicted_temp']).mean()

```

```
1 MED
```

0.24371299303944335

IT has MED 0.24 score

XGB REGRESSOR And Train

## XGBoost

```

# Packages
import xgboost
from xgboost import XGBRegressor

# Train
train_df_copy_2=pd.DataFrame.copy(global_temperatures[global_temperatures['dt'].map(lambda x: x.year)<1980])
train_df_2=pd.DataFrame.copy(train_df_copy_2)
train_df_2.drop(train_df_2.index[:1], inplace=True)
train_df_copy_2.drop(train_df_2.index[-1], inplace=True)

# To predict each average Land temperature we use the previous year's temperatures (Landandocean, min , max)
for j in (train_df_2.index):
    for column in train_df_2.columns[1:-1]:
        train_df_2.loc[j,column]=train_df_copy_2.loc[j-1,column]
forest2 = XGBRegressor()
train_df_2=train_df_2.drop(['dt'],axis=1)
forest2 = forest2.fit(train_df_2.drop(["LandAverageTemperature"],axis=1),train_df_2["LandAverageTemperature"])

```

Test

```

# Test
test_df_copy_2=pd.DataFrame.copy(global_temperatures[global_temperatures['dt'].map(lambda x: x.year)>=1980])
test_df_2=pd.DataFrame.copy(test_df_copy_2)
test_df_2.drop(test_df_2.index[:1], inplace=True)
test_df_copy_2.drop(test_df_2.index[-1], inplace=True)

# To predict each average Land temperature we use the previous year's temperatures (Landandocean, min , max)
for j in (test_df_2.index):
    for column in test_df_2.columns[1:-1]:
        test_df_2.loc[j,column]=test_df_copy_2.loc[j-1,column]
test_df_2=test_df_2.drop(['dt'],axis=1)
output = forest2.predict(test_df_2.drop(["LandAverageTemperature"],axis=1))

```

MED on XGB

```

▶ 1 test_df_2['Predicted_temp']=output

▶ 1 # Distance computing
  2 MED=abs(test_df_2['LandAverageTemperature']-test_df_2['Predicted_temp']).mean()

▶ 1 MED
8]: 0.24781643837399955

```

## KNN and Train Data

### Nearest neighbors

```

: ▶ 1 from sklearn.neighbors import KNeighborsRegressor

: ▶ 1 # Train
  2 neigh = KNeighborsRegressor(n_neighbors=2)
  3 train_df_copy_3=pd.DataFrame.copy(global_temperatures[global_temperatures['dt'].map(lambda x: x.year)<1980])
  4 train_df_3=pd.DataFrame.copy(train_df_copy_3)
  5 train_df_3.drop(train_df_3.index[:1], inplace=True)
  6 train_df_copy_3.drop(train_df_3.index[-1], inplace=True)
  7
  8 # To predict each average Land temperature we use the previous year's temperatures (Landandocean, min , max)
  9 for j in (train_df_3.index):
 10     for column in train_df_3.columns[1:-1]:
 11         train_df_3.loc[j,column]=train_df_copy_3.loc[j-1,column]
 12 train_df_3=train_df_3.drop(['dt'],axis=1)
 13 neighbors=neigh.fit(train_df_3.drop(["LandAverageTemperature"],axis=1),train_df_3["LandAverageTemperature"])

```

## Test Data

```

1 # Test
2 test_df_copy_3=pd.DataFrame.copy(global_temperatures[global_temperatures['dt'].map(lambda x: x.year)>=1980])
3 test_df_3=pd.DataFrame.copy(test_df_copy_3)
4 test_df_3.drop(test_df_3.index[:1], inplace=True)
5 test_df_copy_3.drop(test_df_3.index[-1], inplace=True)
6
7 # To predict each average Land temperature we use the previous year's temperatures (Landandocean, min , max)
8 for j in (test_df_3.index):
9     for column in test_df_3.columns[1:-1]:
10         test_df_3.loc[j,column]=test_df_copy_3.loc[j-1,column]
11 test_df_3=test_df_3.drop(['dt'],axis=1)
12 neigh_output = neighbors.predict(test_df_3.drop(["LandAverageTemperature"],axis=1))

```

```

1 # Adding the predicted temperature to our output DF
2 test_df_3['Predicted_temp']=neigh_output

1 MED=abs(test_df_3['LandAverageTemperature']-test_df_3['Predicted_temp']).mean()

1 MED

```

0.23760904872389793

Knn Has low MED

Minimizing the MED

```

1 #Definition of the function we want to minimize
2
3 def MED_XGB(n,maxi,learn,
4             X_train, y_train, X_test,y_test,
5             MED_limit):
6
7
8
9     xgb_model = xgboost.XGBRegressor(n_estimators = n , learning_rate = learn, max_depth = maxi )
10    model_fit = xgb_model.fit(X_train , y_train )
11    y_predict = model_fit.predict(X_test)
12    MED=abs(y_test-y_predict).mean()
13
14    if MED<MED_limit:
15        return MED , n, learn, maxi
16    else:
17        return 0

```

Separating the Train and Test

```

1 ##### Data Sets #####
2
3 #Train
4 train_df_copy=pd.DataFrame.copy(global_temperatures[global_temperatures['dt'].map(lambda x: x.year)<1980])
5 train_df=pd.DataFrame.copy(train_df_copy)
6 train_df.drop(train_df.index[:1], inplace=True)
7 train_df_copy.drop(train_df.index[-1], inplace=True)
8 for j in (train_df.index):
9     for column in train_df.columns[1:-1]:
10         train_df.loc[j,column]=train_df_copy.loc[j-1,column]
11
12 train_df=train_df.drop(['dt'],axis=1)
13
14 #Test
15 test_df_copy=pd.DataFrame.copy(global_temperatures[global_temperatures['dt'].map(lambda x: x.year)>=1980])
16 test_df=pd.DataFrame.copy(test_df_copy)
17 test_df.drop(test_df.index[:1], inplace=True)
18 test_df_copy.drop(test_df.index[-1], inplace=True)
19 for j in (test_df.index):
20     for column in test_df.columns[1:-1]:
21         test_df.loc[j,column]=test_df_copy.loc[j-1,column]
22
23 test_df=test_df.drop(['dt'],axis=1)
24

```

```

24
25
26 X_train = train_df.drop("LandAverageTemperature", axis = 1)
27 y_train = train_df[ 'LandAverageTemperature' ]
28
29 X_test = test_df.drop("LandAverageTemperature", axis = 1)
30 y_test = test_df[ 'LandAverageTemperature' ]

```

```
1 ##### HyperParameters #####

```

Hyperparameters and their Values

```

1 ##### HyperParameters #####
2
3 max_depth = [2+i for i in range(5)]
4 learning_rate = [0.01*(i+1) for i in range(5)]
5 n_estimators = [370+i for i in range(20)]
6 MED_limit = 0.21
7 l=[]

1 ##### Very Long Loop depending on parameters #####
2
3 for n in n_estimators:
4     for learn in learning_rate:
5         for maxi in max_depth:
6             a = MED_XGB(n,maxi,learn,X_train,y_train,X_test,y_test,MED_limit)
7             if a != 0:
8                 l.append(a)
9                 l.sort()
10
11 l[:5] # 5 Best MED Results

```

Best Results and Their Values

```

[(0.20242728570995633, 385, 0.04, 2),
 (0.2024370215962491, 386, 0.04, 2),
 (0.20244119143209427, 387, 0.04, 2),
 (0.20250209825166132, 384, 0.04, 2),
 (0.20251078753017762, 383, 0.04, 2)]

```

Applying Grid Search

```

1 ##### GridSearch (Take some time to compute) #####
2
3 from xgboost import plot_importance,to_graphviz, plot_tree
4 from sklearn.model_selection import GridSearchCV
5
6 xgb_model = xgboost.XGBRegressor()
7
8 clf = GridSearchCV(xgb_model,
9                     {'max_depth': max_depth,
10                      'n_estimators': n_estimators,
11                      'learning_rate': learning_rate}, verbose=1)
12
13 clf.fit(X_train,y_train)
14 clf.best_score_, clf.best_params_

```

Fitting 5 folds for each of 500 candidates, totalling 2500 fits

```
9]: (0.9925820953088234,
{'learning_rate': 0.02, 'max_depth': 2, 'n_estimators': 370})
```

It has got best score of 99.25

## Feature Importances

### Features importance

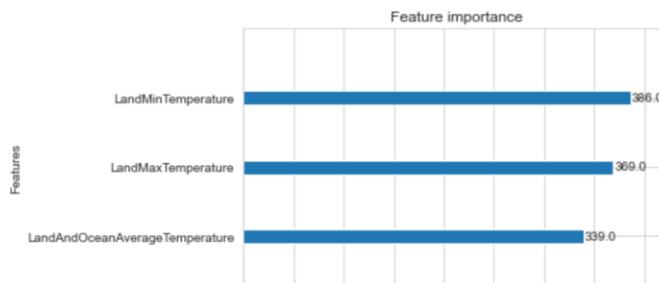
```

: ┌ 1 maxi = 2
  2 learn = 0.04
  3 n = 385
  4
  5 xgb_model = xgboost.XGBRegressor(max_depth = maxi,
  6                                     learning_rate = learn,
  7                                     n_estimators = n)
  8
  9 model_fit = xgb_model.fit(X_train,y_train)

```

```
: ┌ 1 #Ranking the features' importance
  2 plot_importance(model_fit)
```

```
[61]: <AxesSubplot:title={'center':'Feature importance'}, xlabel='F score', ylabel='Features'>
```



Adding CO2 as a Feature

```

1 # Adding the CO2 Level as a feature
2
3 CO2=pd.read_csv("C://Users//user//Downloads//co2-mm-mlo.csv",sep=",")
4
5 # We drop the year 2016 since we don't have the temperatures of 2016
6 CO2 = CO2.drop(CO2.index[-10:])
7 CO2["Date"]=CO2["Date"]+"-01"
8 CO2["Date"]=pd.to_datetime(CO2["Date"])
9 CO2=CO2.drop(["Decimal Date","Interpolated","Trend","Number of Days"],axis=1)
10
11 # We merge the temperatures data with the CO2 data and we keep
12 #only the months where we have the average CO2 Levels
13 temperatures=pd.merge(left=global_temperatures,right=CO2, left_on='dt', right_on='Date', how='outer')
14 temp=pd.DataFrame.copy(temperatures[temperatures['dt'].map(lambda x: x.year)<1980])
15 temp=temp.dropna()

```

### Random Forest

```

1 #Sets
2
3 # Train
4 train_df_copy=pd.DataFrame.copy(temp)
5 train_df=pd.DataFrame.copy(train_df_copy)
6 train_df.drop(train_df.index[:1], inplace=True)
7 train_df_copy.drop(train_df.index[-1], inplace=True)
8
9 X_train = train_df.drop(["LandAverageTemperature","dt","Date"], axis = 1)
10 y_train = train_df['LandAverageTemperature']
11
12 # Test
13 test_df_copy=pd.DataFrame.copy(temperatures[temperatures['dt'].map(lambda x: x.year)>=1980])
14 test_df=pd.DataFrame.copy(test_df_copy)
15 test_df.drop(test_df.index[:1], inplace=True)
16 test_df_copy.drop(test_df.index[-1], inplace=True)
17
18 X_test = test_df.drop(["LandAverageTemperature","dt","Date"],axis = 1 )
19 y_test = test_df['LandAverageTemperature']

```

### Minimizing the Function

```

1 #Definition of the Minimization Function
2
3 def MED_RF(n,maxi,
4             X_train, y_train, X_test,y_test,
5             MED_limit):
6
7 #Note that there is no Learning_rate for the random forest regressor
8
9
10
11     forest = RandomForestRegressor(n_estimators=n,max_depth = maxi)
12     model_fit = forest.fit(X_train , y_train )
13     y_predict = model_fit.predict(X_test)
14     MED=abs(y_test-y_predict).mean()
15
16     if MED<MED_limit:
17         return MED, n, maxi
18     else:
19         return 0

```

Printing the top 5 results

```
1 #Hyperparameters
2
3 max_depth = [8+i for i in range(7)]
4 n_estimators = [i+1 for i in range(100)]
5 MED_limit = 0.2
6 l=[]

1 ##### Minimization (Very Long Loop depending on parameters) #####
2
3 for n in n_estimators:
4     for maxi in max_depth:
5         a = MED_RF(n,maxi,X_train,y_train,X_test,y_test,MED_limit)
6         if a != 0:
7             l.append(a)
8         l.sort()
9
10 l[:5] # 5 Best MED Results
```

: [(0.13500554266563564, 9, 12),  
(0.13506879608146466, 30, 10),  
(0.13692821121428825, 23, 9),  
(0.13729237068377267, 17, 12),  
(0.13780974477958347, 63, 14)]

Grid Search on Random Forest with added Features

```
1 # GridSearch (very Long)
2
3 clf = GridSearchCV(forest,
4                     {'max_depth': max_depth,
5                      'n_estimators': n_estimators}, verbose=1)
6
7 clf.fit(X_train,y_train)
8 clf.best_score_, clf.best_params_
```

Fitting 5 folds for each of 700 candidates, totalling 3500 fits

: (0.9995378623714621, {'max\_depth': 14, 'n\_estimators': 60})

```

1 #Random Forest (GridSearch MED Test)
2
3 # To predict each average Land temperature we use the previous year's temperatures (Landandocean, min
4 for j in (train_df.index):
5     for column in train_df.columns[1:-2]:
6         train_df.loc[j,column]=train_df_copy.loc[j-1,column]
7 forest = RandomForestRegressor(max_depth=10,n_estimators=462)
8 train_df=train_df.drop(['dt','Date'],axis=1)
9 forest = forest.fit(train_df.drop(["LandAverageTemperature"],axis=1),train_df["LandAverageTemperature"]
10
11 for j in (test_df.index):
12     for column in test_df.columns[1:-2]:
13         test_df.loc[j,column]=test_df_copy.loc[j-1,column]
14 test_df=test_df.drop(['dt','Date'],axis=1)
15 output = forest.predict(test_df.drop(["LandAverageTemperature"],axis=1))
16
17 # Distance computing
18 test_df['Predicted_temp']=output
19 MED=abs(test_df['LandAverageTemperature']-test_df['Predicted_temp']).mean()
20 MED
21

```

73]: 0.15355030831921537

## Clustering

### Clustering

#### Preparing the DataFrame

```

1 country_temps = pd.read_csv("C://Users//user//Downloads//GlobalLandTemperaturesByCountry.csv", sep=',')
2
3 #Drop Uncertainty for the clustering
4 country_temps = country_temps.drop("AverageTemperatureUncertainty", axis = 1)
5
6 # Delete all empty lines, change time, and put time as Index
7 country_temps['dt'] = pd.to_datetime(country_temps['dt'])
8 country_temps.tail()
9
10 #Setting the Countries as Index, time as Columns and Average Temperature as the main Data
11 pivot = country_temps.pivot("Country","dt","AverageTemperature")
12 pivot = pivot.drop(pivot.columns[[i for i in range(1274)]+[3238]], axis = 1)
13 pivot.head()
14
15 #Getting rid of all countries with no data
16
17 pivot = pivot.dropna(how='all') #aucune donnée pour le pays
18
19 #Only one country eliminated

```

#### Clustering on Country Wise

```

1 from sklearn.impute import SimpleImputer
2
3 imp = SimpleImputer(strategy="median")
4 # Beware: here axis = 1 means rows -> http://scikit-Learn.org/stable/modules/generated/skLearn.preprocessing.Imputer
5
6 #Imputer changes type of DataFrame.
7
8 trans_pivot = imp.fit_transform(pivot)
9 nonan_pivot = pd.DataFrame(trans_pivot)
10
11 nonan_pivot.columns = pivot.columns
12 nonan_pivot.index = pivot.index
13
14 nonan_pivot.tail()
15
16 print(np.count_nonzero(~np.isnan(nanon_pivot)))
17
18 #No missing data, and let's now call the DataFrame : clusterdf
19
20 clusterdf = nonan_pivot

```

### Clustering Code

```

1 def Clustering(df,x,y,orientation,size):
2
3     from sklearn.cluster import AgglomerativeClustering
4     from scipy.cluster.hierarchy import dendrogram
5     import matplotlib.pyplot as plt
6     plt.style.use('ggplot')
7
8     ward = AgglomerativeClustering(linkage='ward', compute_full_tree=True).fit(df)
9     dendro = [ ]
10    for a,b in ward.children_:
11        dendro.append([a,b,float(len(dendro)+1),len(dendro)+1])
12
13    fig = plt.figure( figsize=(x,y) )
14    ax = fig.add_subplot(1,1,1)
15    r = dendrogram(
16        dendro,
17        color_threshold=None,
18        labels=df.index,
19        distance_sort = True,
20        show_leaf_counts=False,
21        leaf_font_size = size,
22        ax=ax,
23        orientation = str(orientation)
24    )

```

```

def Clustering_truncated(df,x,y,orientation,p,size):

    from sklearn.cluster import AgglomerativeClustering
    from scipy.cluster.hierarchy import dendrogram
    import matplotlib.pyplot as plt
    plt.style.use('ggplot')

    ward = AgglomerativeClustering(linkage='ward', compute_full_tree=True).fit(df)
    dendro = [ ]
    for a,b in ward.children_:
        dendro.append([a,b,float(len(dendro)+1),len(dendro)+1])

    fig = plt.figure( figsize=(x,y) )
    ax = fig.add_subplot(1,1,1)
    r = dendrogram(
        dendro,
        color_threshold=None,
        labels=df.index,
        show_leaf_counts=False,
        truncate_mode = 'lastp', # show only the last p merged clusters
        p = p, # show only the last p merged clusters
        leaf_font_size = size,
        show_contracted = True, # to get a distribution impression in truncated branches
        ax=ax,
        orientation = str(orientation)

    )
    return r["ivl"]

```



## Hierarchical Clustering

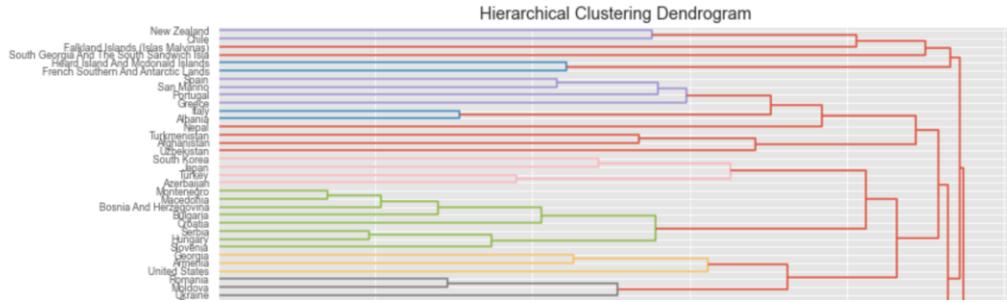
```

1 Clustering(clusterdf,12,30,"right",8.5)
2 plt.title('Hierarchical Clustering Dendrogram')
3 plt.ylabel('Countries')
4 plt.xlabel('Distance')
5
6 #Now, even though scrolling is necessary, this vertical dendrogram is much clearer

```

C:\Users\user\anaconda3\anaconda\_install\lib\site-packages\sklearn\utils\validation.py:1688: FutureWarning: Feature names only support names that are all strings. Got feature names with dtypes: ['Timestamp']. An error will be raised in 1.2.  
warnings.warn(

84]: Text(0.5, 0, 'Distance')



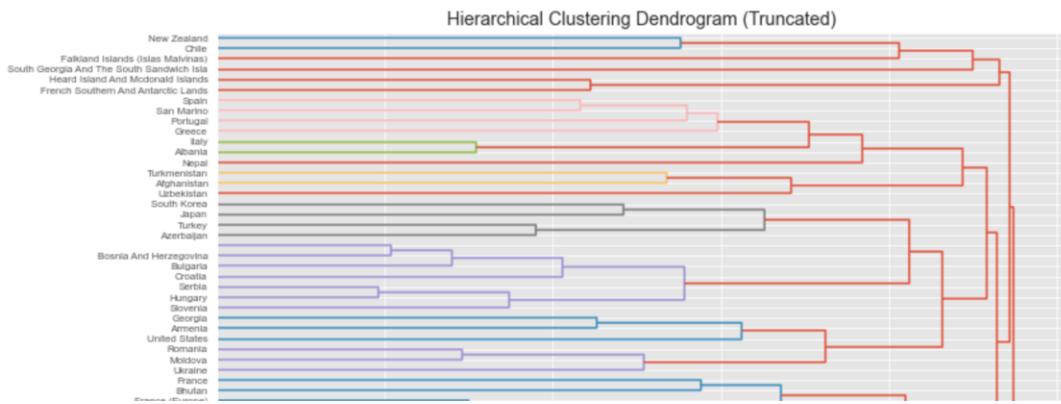
```

1 Cluster_70 = Clustering_truncated(clusterdf,12,30,"right",200,8)
2 plt.title('Hierarchical Clustering Dendrogram (Truncated)')
3 plt.xlabel('Distance')
4 plt.ylabel('Countries')
5

```

C:\Users\user\anaconda3\anaconda\_install\lib\site-packages\sklearn\utils\validation.py:1688: FutureWarning: Feature names only support names that are all strings. Got feature names with dtypes: ['Timestamp']. An error will be raised in 1.2.  
warnings.warn(

85]: Text(0, 0.5, 'Countries')



```
In [86]: 1 def delete_symbol(l,sym):
2     for i in range(l.count(sym)):
3         l.remove(sym)
4     return l, len(l)

In [87]: 1 cntry_to_del = delete_symbol(Cluster_70,'')
2 cntry_to_del
Kuwait ,
'Saudi Arabia',
'Bahrain',
'Qatar',
'United Arab Emirates',
'Mali',
'Mauritania',
'Aruba',
'Belize',
'Bahamas',
'Cuba',
'Brazil',
'Baker Island',
'Malaysia',
'Singapore',
'Costa Rica',
'Panama',
```

## Country wise and Date wise Dataframe

```
] 1 clusterdf_new = clusterdf.drop(cntry_to_del[0], axis = 0)
2 clusterdf_new.head()

t[88]:
      dt 1850- 1850- 1850- 1850- 1850- 1850- 1850- 1850- 1850- 1850- ... 2012- 2012- 2013- 2013- 2013- 2013-
          01-01 02-01 03-01 04-01 05-01 06-01 07-01 08-01 09-01 10-01 ...
          11-01 12-01 01-01 02-01 03-01 04-01 05-01 06-01

  Country
American Samoa 13.960 15.370 18.443 20.359 21.059 23.902 24.512 24.570 23.448 22.030 ... 27.963 27.316 27.241 27.793 27.683 27.824 28.045 27.650
Anguilla 24.075 24.985 25.064 24.859 25.734 27.073 27.922 27.949 27.696 26.701 ... 28.203 26.476 25.919 25.899 26.187 26.679 27.667 28.414
Antigua And Barbuda 23.930 24.845 24.984 24.879 25.711 26.988 27.714 27.672 27.346 26.351 ... 27.890 26.329 25.787 25.802 26.089 26.601 27.706 28.322
Australia 13.960 15.370 18.443 20.359 21.059 23.902 24.512 24.570 23.448 22.030 ... 26.935 28.417 29.861 28.311 26.590 23.305 19.377 16.018
Barbados 24.249 25.140 25.271 25.304 25.918 27.006 27.312 27.109 26.968 26.298 ... 27.922 26.744 26.088 26.182 26.504 26.806 27.849 28.020

5 rows × 1964 columns
```

```
1 Clustering(clusterdf_new,12,30,"right",10)
2 plt.title('Hierarchical Clustering Dendrogram - 70 Countries')
3 plt.xlabel('Distance')
4 plt.ylabel('Countries')
5
```

C:\Users\user\anaconda3\anaconda\_install\lib\site-packages\sklearn\utils\validation.py:1688: FutureWarning: Feature names only support names that are all strings. Got feature names with dtypes: ['Timestamp']. An error will be raised in 1.2.  
warnings.warn(

```
:9]: Text(0, 0.5, 'Countries')
```

