# Data Science and Global Warming: simple machine learning techniques applied

Yassine BENAZZOU, Charles DOGNIN, Jean-Michel ROUFOSSE

December 23, 2016

## Abstract

We use machine learning techniques to study the evolution of climate change in the world since 1800. The dataset comes from the Kaggle website[1], displaying the evolution since 1750 in degree Celsius of the Land Average Temperature, the Land Average Temperature Uncertainty (the 95% confidence interval around the mean), the Land Max Temperature, the Land Max Temperature Uncertainty, the Land Min Temperature,the Land Min Temperature Uncertainty, the Land And Ocean Average Temperature, and finally the Land And Ocean Average Temperature Uncertainty. The data gathered by Kaggle come from a far-reaching study, The Berkeley Earth Surface Temperature Study which "combines 1.6 billion temperature reports from 16 pre-existing archives"[2]. The second set of data comes from the the US Government's Earth System Research Laboratory, Global Monitoring Division and register the level of CO2 globally on average since 1958.

---

[1] https://www.kaggle.com/berkeleyearth/climate-change-earth-surface-temperature-
[2] http://berkeleyearth.org/data/

# Contents

# 1 Introduction: Why climate change?

The emergence of the Big Data era is an opportunity for scientists to test their theories and to refine their models: *currently, about 25 laboratories across the world support almost 50 climate models, forming the basis of climate projections (predictions) assessed by the Intergovernmental Panel on Climate Change (IPCC), which was established by the United Nations in 1988 and received the 2007 Nobel Peace Prize (shared with former Vice President Al Gore).*[3].But we wanted to see if we could at our level and using simple machine learning models assess climate change. We naturally use simpler machine learning models than those used by climate experts. The Big Data is also an opportunity for the public to realise how serious environmental challenges are. More data indeed imply more graphs, plots that can easily be understood by people without scientific knowledge. In this project we thus also stress the variety and interactivity of the plotting, majorly using a very innovative package called bokeh[3]. We show in the report almost all the graphs we did but in order to fully take advantage of their interactivity, we invite the reader to use our notebook. The next step was to find a significant database. We found one on the website Kaggle. Their are five files in the dataset: the monthly global temperatures since 1750 globally (1), by country (2), by state (3), by major city (4) or by city. In this project, we use the the files (1), (2) and (4). In the last part, we add one variable, the level of CO2 pollution and evaluate its impact on our models.

# 2 Data description and preparation

We began by importing the global temperatures. We import as well the standard data analysis packages Numpy and Pandas, some plotting packages Maplotlib, Seaborn and Bokeh. Finally, the machine learning packages that we use in our models from scikit-learn and XGBoost.

Listing 1: Imports

```python
# Imports

# pandas
import pandas as pd
from pandas import Series,DataFrame

# numpy, matplotlib, seaborn, bokeh
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import bokeh as bk
sns.set_style('whitegrid')
%matplotlib inline

# machine learning sklearn
from sklearn import linear_model
from sklearn import metrics
from sklearn.svm import SVC, LinearSVC
from sklearn.ensemble import RandomForestClassifier
```

[3]Notebook precision: the reader has to run the concerned cells to see the bokeh graphs appear

```
20  from sklearn.ensemble import RandomForestRegressor
21  from sklearn.neighbors import KNeighborsClassifier
22  from sklearn.naive_bayes import GaussianNB
23
24  #machine learning xgboost
25  import xgboost
```

Then we visualise and clean our data. We want to keep only two variables: the Land Average Temperature and the Uncertainty at 95% regarding the temperature. The first problem we faced was the number of NaN values on the dataset. Back in the 18th century, it was probably much more difficult to evaluate the global Land Average Temperature. After this cleaning, using the function *dropna()* we realise that the data set becomes really significant from 1800 onwards. It is consistent with the fact that climate experts generally use temperature data that start in the 19th century[3].

Listing 2: Visualisation, Cleaning

```
1
2   # get global temp csv file as a DataFrame
3   global_temps = pd.read_csv("GlobalTemperatures.csv", sep=',')
4
5   # preview the data
6   global_temps.head()
7
8   #Delete all columns except dates, Land Average Temperatures,
9   #Land Average Temperature Uncertainty
10  global_temps=global_temps.drop(global_temps.columns[3:], axis=1)
11
12  global_temps.head()
13
14  #Delete all empty lines
15  global_temps=global_temps.dropna()
```

We understand that the data will not be easily displayed if we keep the monthly temperatures. We resize the data yearly and build a simpler dataframe.

Listing 3: Resizing

```
1   #Resizing the datasets, we only keep the year averages
2
3   date=global_temps['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(global_temps[y==date]
11      ['LandAverageTemperature'].mean())
12      mean_temps_uncertainty.append(global_temps[y==date]
13      ['LandAverageTemperatureUncertainty'].mean())
14
15  #Defining a function that classifies the average temperatures
16
17  def heat(x):
18
19      if type(x) is str:
```

```
20          return x
21      elif x<=7:
22          return 3
23      elif x>7 and x<=8:
24          return 2
25      else:
26          return 1
27
28  #Building of a new, simpler dataframe
29
30  years=pd.to_numeric(years)
31
32  d1={'Dates': pd.Series(years),'Temperatures':pd.Series(
        mean_temps_global),
33  'Uncertainty':pd.Series(mean_temps_uncertainty)}
34
35  df_global_mean=pd.DataFrame(d1)
36
37  #Adding the heat function
38
39  e=df_global_mean.applymap(heat)
40
41  d2={'Dates': pd.Series(years),'Temperatures':pd.Series(
        mean_temps_global),
42  'Uncertainty':pd.Series(mean_temps_uncertainty),
43      'Heat':e['Temperatures']}
44
45  df_global_mean=pd.DataFrame(d2)
```

# 3   Exploratory Data Analysis

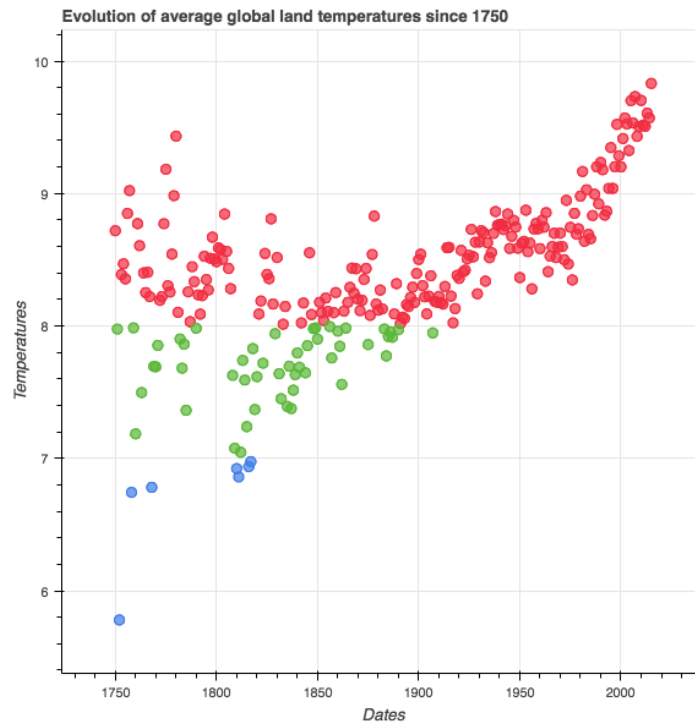We plot the first graph using bokeh:

**Listing 4: Evolution of Average Global Temperatures since 1750**

```
1  #Global Average Temperatures plot using Bokeh
2
3  from bokeh.io import output_notebook
4  from bokeh.charts import Scatter, show
5  from bokeh.plotting import figure
6  output_notebook()
7  p = Scatter(df_global_mean, x='Dates', y='Temperatures',
8  color='Heat', legend=None, background_fill_color="white",
9  title="Evolution of average global land temperatures since 1750"
        )
10  show(p)
```

Figure 1: Evolution of average global land temperatures since 1750



We see a clear path towards increasing temperatures. From 1900 onwards, the average global temperatures variance is decreasing dramatically and from that date, the temperatures are increasing at an ever faster pace. Now we take a look at the uncertainty (in celcius degree) regarding the accuracy of the measured temperatures.

Listing 5: Evolution of the Uncertainty regarding the Average Global Temperatures since 1750

```
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,
13 title='A clear decreasing of uncertainty')
14
15 p.line(x, y, line_dash="4 4", line_width=1, color='gray')
16
```
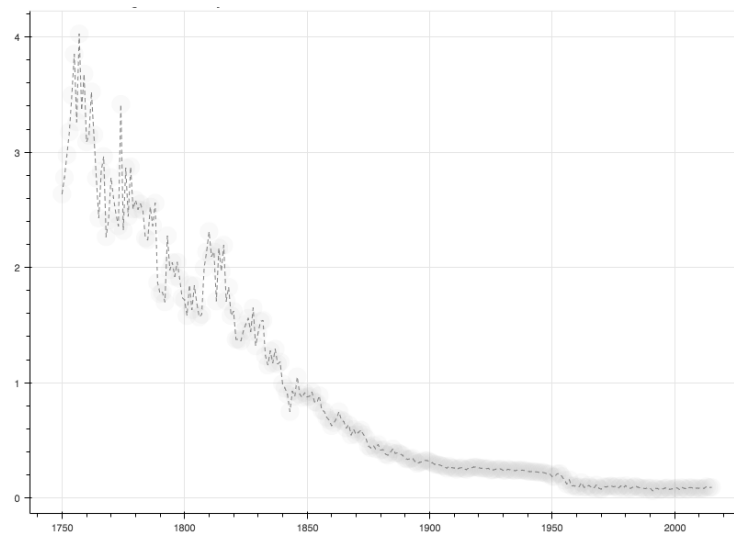
```
17  # Add a circle, that is visible only when selected
18  source = ColumnDataSource({'x': x, 'y': y})
19  invisible_circle = Circle(x='x', y='y', fill_color='gray',
20  fill_alpha=0.05, line_color=None, size=20)
21  visible_circle = Circle(x='x', y='y', fill_color='firebrick',
22  fill_alpha=0.5,
23  line_color=None, size=20)
24  cr = p.add_glyph(source, invisible_circle,
25  selection_glyph=visible_circle, nonselection_glyph=
        invisible_circle)
26
27  # Add a hover tool, that selects the circle
28  code = "source.set('selected', cb_data['index']);"
29  callback = CustomJS(args={'source': source}, code=code)
30  p.add_tools(HoverTool(tooltips=None, callback=callback,
31  renderers=[cr], mode='hline'))
32
33  show(p)
```

Figure 2: Evolution of uncertainty regarding the average global land temperatures since 1750



The uncertainty is clearly decreasing as well and is near 0 for the past few years. Global warming is a reality, not a fiction. To be able to fully appreciate the interactivity of the graph, we encourage the reader to take a look at the notebook.

After doing the same manipulation as those done for the average, we explore the evolution of the extremes and we plot them.

Listing 6: Evolution of the Min and Max Temperatures

```
1  #Loading the Min and Max data
2
3  global_temps2 = pd.read_csv("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]
22     ['LandAverageTemperature'].mean())
23     max_temps_1800.append(global_temps2[y==date]
24     ['LandMaxTemperature'].mean())
25     min_temps_1800.append(global_temps2[y==date]
26     ['LandMinTemperature'].mean())
27
28 #Building of a new, simpler dataframe
29
30 list_of_years_1800=pd.to_numeric(list_of_years_1800)
31
32 d1={'Average Temperature':pd.Series(mean_temps_1800),
33 'Dates': pd.Series(list_of_years_1800),
34 'Min Temperature':pd.Series(min_temps_1800),
35 'Max Temperature':pd.Series(max_temps_1800)}
36
37 global_temps_1800=pd.DataFrame(d1)
38
39 global_temps_1800=global_temps_1800[['Dates','Average
       Temperature',
40 'Min Temperature','Max Temperature']]
41
42 from bokeh.plotting import figure, show
43 from bokeh.models import ColumnDataSource, Circle, HoverTool,
       CustomJS
44 from bokeh.io import output_notebook
45
46 output_notebook()
47
48 x = global_temps_1800['Dates']
49 y = global_temps_1800['Average Temperature']
50 z = global_temps_1800['Min Temperature']
51 t = global_temps_1800['Max Temperature']
52
53 p = figure(plot_width=600, plot_height=600)
54
55
56 p.line(x, y, legend="Average Temperature", line_color="green",
       line_width=3)
57
58 p.line(x,z, legend="Min Temperature", line_color="blue",
       line_width=3)
59
60 p.line(x, t, legend="Max Temperature", line_color="red",
       line_width=3)
61
62 # change just some things about the x-grid
```
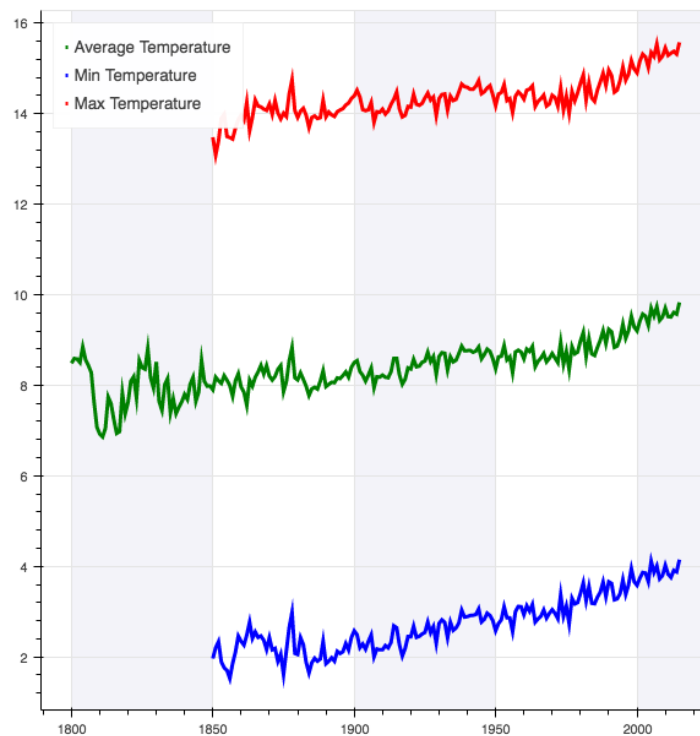
```
63 | p.xgrid.band_fill_alpha = 0.05
64 | p.xgrid.band_fill_color = "navy"
65 |
66 | p.legend.location = "top_left"
67 | p.legend.glyph_width = 2
68 | p.legend.label_width=2
69 |
70 | show(p)
```

Figure 3: Comparison Extreme and Average Temperatures since 1800



We notice the lack of data regarding the extremes between 1800 and 1850. The upward trend is confirmed.

The result of our first Exploratory Data Analysis helps us understand the Average Temperature trend globally over the past few centuries. We then train and test several machine learning models on the dataset.

# 4 Predictive Modelling

After having explored the data visualy. We try several machine learning models. To evaluate their accuracy, we define a metric that measures the fit of the model to the data. We call it Mean Euclidean Distance (MED), it is defined as:

**Definition 1.**

$$\mathbb{E}[\|\beta - \hat{\beta}\|]$$

*With $\beta$ being the true value of the temperature, $\hat{\beta}$ the predicted value and $\|.\|$ being the Euclidean Norm.*

We use this metric for almost every machine learning model we devised. The only exceptions being the Linear and Polynomial Regression for which we use more classical metrics like the mean squared error or the explained variance score. Furthermore, the visual aspect of regression is sufficient to see if it fits or not the data.

## 4.1 Linear Regression

We use first the simpler continuous machine learning model: the linear regression. The output variable is the temperature, the input variables are the dates. We evaluate the relationships between time and temperature. First, we train the model on the whole dataset. To test the accuracy of the model, we use graphs and classical regression metrics.

Listing 7: Linear Regression

```
1   #Regression on all the data (without prediction) => See if a
        regression can be a good estimation
2
3   import matplotlib.pyplot as plt
4   import numpy as np
5   from sklearn import linear_model
6
7
8   #Define the input and output
9   temp=global_temps_1800['Average Temperature']
10  time=global_temps_1800['Dates']
11
12
13  # The data sets
14  dates_X = time.iloc[0:,]
15
16  # The targets
17  temperature_y = temp.iloc[0:,]
18
19  # Create linear regression object
20  reg = linear_model.LinearRegression()
21
22  # Train the model using the data sets
23  reg.fit(dates_X.to_frame(),temperature_y.to_frame())
24
25  # The estimation
26  output_reg=reg.predict(dates_X.to_frame())
27  output_simple=np.copy(output_reg)
28  output=pd.DataFrame(index=pd.DataFrame(output_simple).index.
        values)
```

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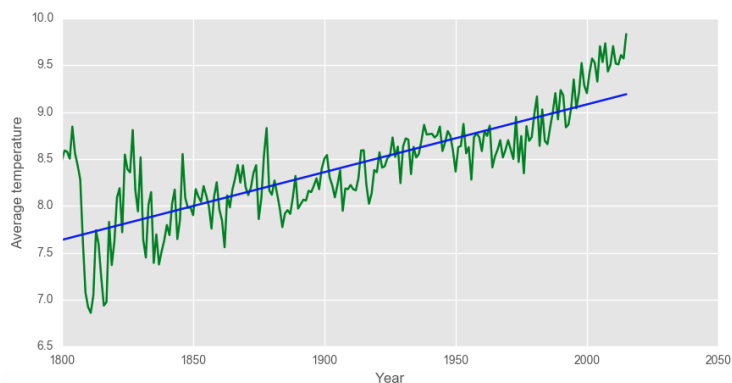
We obtain the following results:

<div align="center">

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

</div>

And the following graph:

<div align="center">

Figure 4: Linear Regression 1

</div>



As shown by the graphe, the regression seems to give accurate estimates of our real temperatures. Then we try to use it as a predictive model, separating the train and test databases. The training database spans from 1800 to 1949 and the test database from 1950 to 2015.

Listing 8: Linear Regression

```
1   import matplotlib.pyplot as plt
2   import numpy as np
3   from sklearn import linear_model
4
5   #Define the input and output
6   temp=global_temps_1800['Average Temperature']
7   time=global_temps_1800['Dates']
8
9   # Split the data into training/testing sets
10  dates_X_train = time.iloc[0:150,]
11  dates_X_test = time.iloc[150:,]
```

11

```
12
13  # Split the targets into training/testing sets
14  temperature_y_train = temp.iloc[0:150,]
15  temperature_y_test = temp.iloc[150:,]
16
17  # Create linear regression object
18  reg = linear_model.LinearRegression()
19
20  # Train the model using the training sets
21  reg.fit(dates_X_train.to_frame(),temperature_y_train.to_frame())
22
23  # The coefficients
24  print('Coefficients: \n', reg.coef_)
25  # The intercept
26  print('Intercept: \n', reg.intercept_)
27
28  # The prediction
29  output_simple=reg.predict(dates_X_test.to_frame())
30
31  # The mean squared error
32  print("Mean squared error: %.2f"
33         % np.mean((reg.predict(dates_X_test.to_frame())
34         - temperature_y_test.to_frame()) ** 2))
35
36  # Explained variance score: 1 is perfect prediction
37  print('Variance score: %.2f' % reg.score(dates_X_test.to_frame()
           ,
38   temperature_y_test.to_frame()))
39
40  # Linear regression curve on the predicted values
41
42  fig, axes = plt.subplots(figsize=(12,10))
43  axes.set_ylabel('Average temperature')
44  axes.set_xlabel('Year')
45  plt.plot(dates_X_test, temperature_y_test, label='Annual Mean
           temperature',
46  color='g')
47  plt.plot(dates_X_test,output_simple,label='Forecast Temperature
           by Regression',
48  color='b')
```

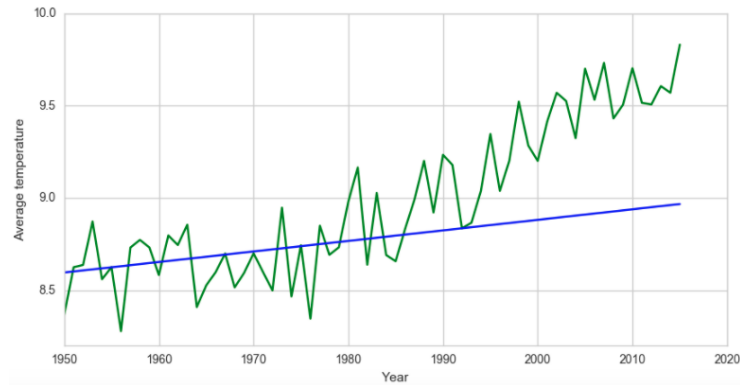The metrics of the output are:

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

Figure 5: Linear Regression 2



We deduce from the metrics and the graph that the linear regression is not relevant when it comes to testing. Whereas the parameters of the model were accurately fitting the evolution of the temperatures when they took their information from the whole dataset, they are not relevant anymore when they are built from a part of the dataset. Therefore, the linear regression is not an efficient predictive model.

## 4.2 Polynomial Regression

"*Historically, the standard way to extend linear regression to settings in which the relationship between the predictors and the response is nonlinear has been to replace the standard linear model*"[4]. It seems to be the case for time and temperature. We thus try to implement a Polynomial Regression with $\deg(\mathbb{P}) \leq 3$. We realise the same experience as regarding the linear regression. First, we train the model on the whole dataset. To evaluate the fit, we use once again a graph and traditional regression metrics.

Listing 9: Polynomial Regression 2

```
1   #Polynomial Regression t^2 + t^3
2
3   #Regression on all the data (without prediction) => See if a
        regression can be a good estimation
4
5   #Define the input and output
6   time=pd.DataFrame()
7   time["Dates"]=global_temps_1800['Dates']
8   time["Dates^2"]=global_temps_1800['Dates']*global_temps_1800['
        Dates']
9   time["Dates^3"]=global_temps_1800['Dates']*global_temps_1800['
        Dates']*global_temps_1800['Dates']
10  temp=global_temps_1800['Average Temperature']
11
12
13  # The data sets
14  dates_X = time.iloc[0:,0]
15
```

---

[4][2]

13

```
16   # The targets
17   temperature_y = temp.iloc[0:,]
18
19   # Create linear regression object
20   reg = linear_model.LinearRegression()
21
22   # Train the model using the training sets
23   reg.fit(time.iloc[0:,],temperature_y.to_frame())
24
25   # The prediction
26   output_reg=reg.predict(time.iloc[0:,])
27   output_multiple=np.copy(output_reg)
28   output=pd.DataFrame(index=pd.DataFrame(output_simple).index.
         values)
29   # The coefficients
30   print('Coefficients: \n', reg.coef_)
31   # The intercept
32   print('Intercept: \n', reg.intercept_)
33
34   # The mean squared error
35   print("Mean squared error: %.2f" % np.mean((output_reg -
         temperature_y.to_frame()) ** 2))
36
37   # Explained variance score: 1 is perfect prediction
38   print('Variance score: %.2f' % reg.score(time.iloc[0:,],
         temperature_y.to_frame()))
39
40   # Euclidian Distance
41   output['Predicted Temperature']=pd.DataFrame(output_multiple)
42   output["Average Temperature"]=temp
```

The results are the following ones:
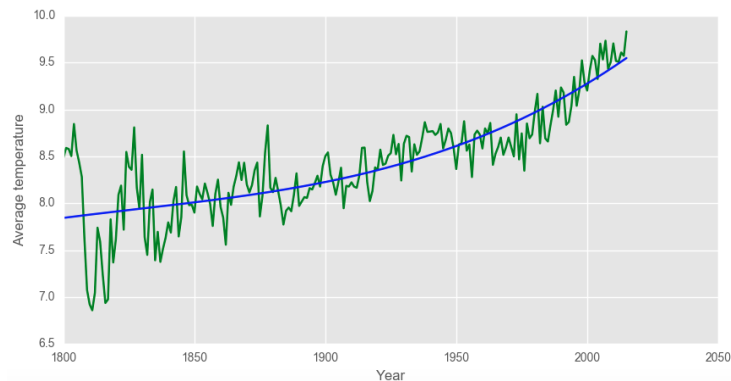
Coefficients:[1.54583760e+00 -8.43972200e-04 1.53912499e-07]
Intercept:[-937.81103859]
Mean squared error: 0.10
Variance score: 0.69

And the graph:

Figure 6: Linear Regression 2



The model seems to fit the data even better than the first linear regression. We study then as before the predictive power of it. Once again, the train

database spans from 1800 to 1949 and the test database from 1950 to 2015.
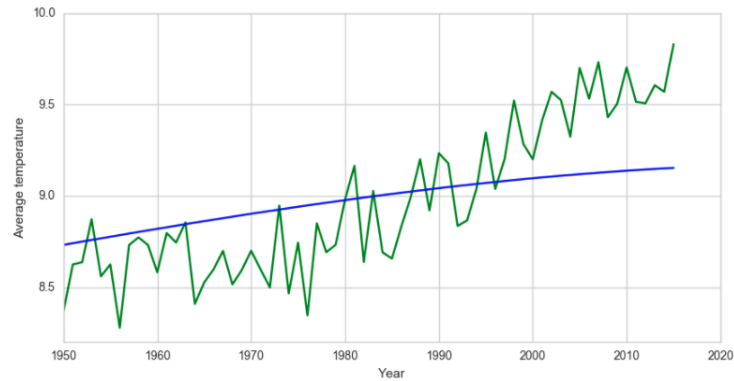
**Listing 10: Polynomial Regression 2**

```python
#Polynomial Regression t^2 + t^3

#Define the input and output
temp=global_temps_1800['Average Temperature']
time=pd.DataFrame()
time["Dates"]=global_temps_1800['Dates']
time["Dates^2"]=global_temps_1800['Dates']*global_temps_1800['Dates']
time["Dates^3"]=global_temps_1800['Dates']
*global_temps_1800['Dates']*global_temps_1800['Dates']

# Split the data into training/testing sets
dates_X_train = time.iloc[0:150,0]
dates_X_test = time.iloc[150:,0]

# Split the targets into training/testing sets
temperature_y_train = temp.iloc[0:150,]
temperature_y_test = temp.iloc[150:,]

# Create linear regression object
reg = linear_model.LinearRegression()

# Train the model using the training sets
reg.fit(time.iloc[0:150,],temperature_y_train.to_frame())

# The coefficients
print('Coefficients: \n', reg.coef_)
# The intercept
print('Intercept: \n', reg.intercept_)

# The prediction
output_multiple=reg.predict(time.iloc[150:,])

# The mean squared error
print("Mean squared error: %.2f"
      % np.mean((reg.predict(time.iloc[150:,])
      - temperature_y_test.to_frame()) ** 2))

# Explained variance score: 1 is perfect prediction
print('Variance score: %.2f' % reg.score(time.iloc[150:,]
, temperature_y_test.to_frame()))

# Multiple Linear Regression: t,t^2, t^3:

fig, axes = plt.subplots(figsize=(12,10))
axes.set_ylabel('Average temperature')
axes.set_xlabel('Year')
# Plot
plt.plot(dates_X_test, temperature_y_test,
label='Annual Mean temperature', color='g')
plt.plot(dates_X_test,output_multiple,label='Forecast Reg',color
     ='b')
```

The results are:

Coefficients:[-3.37849536e+00 1.75645409e-03 -3.03542634e-07]
Intercept:[2168.61282322]

15

Mean squared error: 0.10
Variance score: 0.42

Figure 7: Polynomial Regression



The variance score is significantly higher than the 2nd linear regression's one. Still, we cannot consider it as a relevant predictive model for our data. We turn to the correlation matrix to see if we can learn more about our variables.
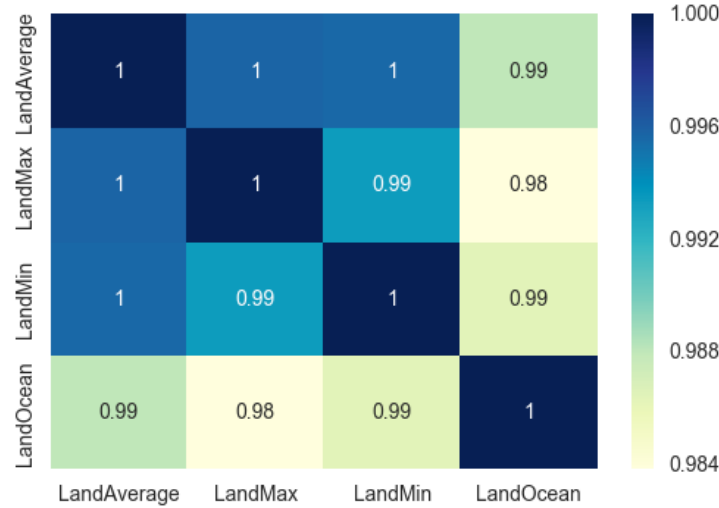
## 4.3 Correlation Matrix

We calculate the correlations between some of our variables. We only look at all the temperatures variables, we disregard here the uncertainty.

Listing 11: Correlation Matrix

```python
#Correlation Matrix for the following features :
#LandAverageTemperature, LandMaxTemperature,
#LandMinTemperature,LandAndOceanAverageTemperature

global_temperatures=pd.read_csv("GlobalTemperatures.csv", sep=',')
global_temperatures['dt'] = pd.to_datetime(global_temperatures['dt'])
global_temperatures=global_temperatures.drop(
['LandAverageTemperatureUncertainty',"
    LandMaxTemperatureUncertainty","
    LandMinTemperatureUncertainty","
    LandAndOceanAverageTemperatureUncertainty"],axis=1)
global_temperatures=global_temperatures.dropna()
corr=global_temperatures.corr()

# Heatmap for the previous correlation matrix
sns.heatmap(corr,xticklabels=['LandAverage','LandMax','LandMin',
    'LandOcean'], yticklabels=['LandAverage','LandMax',
'LandMin','LandOcean'], annot=True,cmap="YlGnBu")
```

Figure 8: Correlation Matrix



The variables are very linked. Only the variable *LandandOceanTemperature* is slightly less correlated with than the others between them. We turn to another predictive model, the Random Forest Regressor.

## 4.4   Random Forest

We use the Random Forest Regressor proposed by the scikit-learn package. We tested different values regarding the number of trees and we realised that $n = 10$ was the most efficient. We decide to choose a larger training database than for the previous regressions. We regress the average temperatures on the landandocean, min and max temperatures, one month before. Otherwise, the output (average) would be based on inputs (min, max, landandoceanaverage) of the same year, which would biaised the result. The train database spans from 1800 to 1980 and the test database from 1981 onwards. It allows for a finer tuning of our model. To test the fit of the model, we use the metric defined above, what we named the Mean Euclidean Distance (MED).

Listing 12: First Random Forest

```
1
2   # Train
3   train_df_copy=pd.DataFrame.copy(global_temperatures[
        global_temperatures['dt'].
4   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
9   # To predict each average land temperature we use the
10  #previous years' temperatures (ocean, min, max)
11  for j in (train_df.index):
12      for column in train_df.columns[1:-1]:
13          train_df.loc[j,column]=train_df_copy.loc[j-1,column]
```

17

```
14  forest = RandomForestRegressor(n_estimators=10)
15  train_df=train_df.drop(['dt'],axis=1)
16  forest = forest.fit(train_df.drop(["LandAverageTemperature"],
        axis=1),
17  train_df["LandAverageTemperature"])
18
19  # Test
20  test_df_copy=pd.DataFrame.copy(global_temperatures[
        global_temperatures['dt'].
21  map(lambda x: x.year)>=1980])
22  test_df=pd.DataFrame.copy(test_df_copy)
23  test_df.drop(test_df.index[:1], inplace=True)
24  test_df_copy.drop(test_df.index[-1], inplace=True)
25
26  # To predict each average land temperature we use the
27  #previous year's temperatures (landandocean, min , max)
28  for j in (test_df.index):
29      for column in test_df.columns[1:-1]:
30          test_df.loc[j,column]=test_df_copy.loc[j-1,column]
31  test_df=test_df.drop(['dt'],axis=1)
32  output = forest.predict(test_df.drop(["LandAverageTemperature"],
        axis=1))
33
34  #Adding the predicted temperature to our output DF
35  test_df['Predicted_temp']=output
36
37  # Distance computing
38  MED=abs(test_df['LandAverageTemperature']-test_df['
        Predicted_temp']).mean()
```

We found a MED of 0.21886310904872408. We store this value and will compare it to the values obtained with other models.

## 4.5 XGBoost

After having used the Random Forest Regressor we use the eXtreme Gradient Boosting Regressor, XGBoostRegressor. The training and test sets, the output and input variables are the same as before. We decide to test the model with the default parameters and to analyse the outcome.

Listing 13: First XGBoost

```
1
2   # Packages
3   import xgboost
4   from xgboost import XGBRegressor
5
6   # Train
7   train_df_copy_2=pd.DataFrame.copy(global_temperatures[
        global_temperatures['dt'].
8   map(lambda x: x.year)<1980])
9   train_df_2=pd.DataFrame.copy(train_df_copy_2)
10  train_df_2.drop(train_df_2.index[:1], inplace=True)
11  train_df_copy_2.drop(train_df_2.index[-1], inplace=True)
12
13  # To predict each average land temperature we use the
14  #previous year's temperatures (landandocean, min , max)
15  for j in (train_df_2.index):
16      for column in train_df_2.columns[1:-1]:
17          train_df_2.loc[j,column]=train_df_copy_2.loc[j-1,column]
```

```
18  forest2 = XGBRegressor()
19  train_df_2=train_df_2.drop(['dt'],axis=1)
20  forest2 = forest2.fit(train_df_2.drop(["LandAverageTemperature"
        ],
21  axis=1),train_df_2["LandAverageTemperature"])
22
23  # Test
24  test_df_copy_2=pd.DataFrame.copy(global_temperatures[
        global_temperatures['dt'].
25  map(lambda x: x.year)>=1980])
26  test_df_2=pd.DataFrame.copy(test_df_copy_2)
27  test_df_2.drop(test_df_2.index[:1], inplace=True)
28  test_df_copy_2.drop(test_df_2.index[-1], inplace=True)
29
30  # To predict each average land temperature we use the previous
31  #year's temperatures (landandocean, min , max)
32  for j in (test_df_2.index):
33      for column in test_df_2.columns[1:-1]:
34          test_df_2.loc[j,column]=test_df_copy_2.loc[j-1,column]
35  test_df_2=test_df_2.drop(['dt'],axis=1)
36  output = forest2.predict(test_df_2.drop(["LandAverageTemperature
        "],axis=1))
37
38  # Adding the predicted temperature to our output DF
39  test_df_2['Predicted_temp']=output
40
41  # Distance computing
42  MED=abs(test_df_2['LandAverageTemperature']-test_df_2['
        Predicted_temp']).mean()
```

We obtain a MED of 0.20515603448730613. The model is thus more accurate than the Random Forest Regressor.

## 4.6   K-Nearest Neighbors

We turn to the K-Nearest Neighbors machine learning method. After having tested a large amount of different values, we find that we minimize the MED with $k = 2$. The training and test sets, the output and input variables are the same as before.

**Listing 14: K-Nearest Neighbors**

```
1
2   ## Packages
3   from sklearn.neighbors import KNeighborsRegressor
4
5   # Train
6   neigh = KNeighborsRegressor(n_neighbors=2)
7   train_df_copy_3=pd.DataFrame.copy(global_temperatures[
        global_temperatures['dt'].map(lambda x: x.year)<1980])
8   train_df_3=pd.DataFrame.copy(train_df_copy_3)
9   train_df_3.drop(train_df_3.index[:1], inplace=True)
10  train_df_copy_3.drop(train_df_3.index[-1], inplace=True)
11
12  # To predict each average land temperature we use the
13  #previous year's temperatures (landandocean, min , max)
14  for j in (train_df_3.index):
15      for column in train_df_3.columns[1:-1]:
16          train_df_3.loc[j,column]=train_df_copy_3.loc[j-1,column]
17  train_df_3=train_df_3.drop(['dt'],axis=1)
```

```
18  neighbors=neigh.fit(train_df_3.drop(["LandAverageTemperature"],
        axis=1),
19  train_df_3["LandAverageTemperature"])
20
21  # Test
22  test_df_copy_3=pd.DataFrame.copy(global_temperatures[
        global_temperatures['dt'].
23  map(lambda x: x.year)>=1980])
24  test_df_3=pd.DataFrame.copy(test_df_copy_3)
25  test_df_3.drop(test_df_3.index[:1], inplace=True)
26  test_df_copy_3.drop(test_df_3.index[-1], inplace=True)
27
28  # To predict each average land temperature we use the
29  #previous year's temperatures (landandocean, min , max)
30  for j in (test_df_3.index):
31      for column in test_df_3.columns[1:-1]:
32          test_df_3.loc[j,column]=test_df_copy_3.loc[j-1,column]
33  test_df_3=test_df_3.drop(['dt'],axis=1)
34  neigh_output = neighbors.predict(test_df_3.drop(["
        LandAverageTemperature"],axis=1))
35
36  # Adding the predicted temperature to our output DF
37  test_df_3['Predicted_temp']=neigh_output
38
39  MED=abs(test_df_3['LandAverageTemperature']-test_df_3['
        Predicted_temp']).mean()
```

We obtain a MED of 0.23760904872389807, which is higher than both Random Forest and XGBoost. We now want to see if we can find hyperparameters that would give us a lower MED for our most efficient model, the XGB Regressor. We use two hyperparameter optimization techniques: the grid-search technique, with the function $GridSearchCV$ proposed by scikit-learn and a manual minimization.

## 4.7 Minimization of the MED

### 4.7.1 Manual Minimization

We first attempt to find the hyperparameters of the XGBoost that minimize the MED manually. We decide to restrain the possible values of the hyperparameters for the following reasons: if max_depth is too large, the model tends to overfit, regarding the learning_rate we tested several parameters and decided to choose a specific efficient range. As for the n_estimators, we choose not to go beyond 389 because once again, we tested different values.

Listing 15: Manual Minimization

```
1
2
3   #Definition of the function we want to minimize
4
5   def MED_XGB(n,maxi,learn,
6       X_train, y_train, X_test,y_test,
7       MED_limit):
8
9       '''
10      n         : n_estimators  -> Number of boosted trees to fit
11      maxi      : max_depth     -> Maximum tree depth for base
          learners
```

```python
12          learn      : learning_rate -> Boosting learning rate
13          ...
14          MED_limit : -> threshold for minimum MED required
15          '''
16
17          xgb_model = xgboost.XGBRegressor(n_estimators = n ,
                learning_rate = learn, max_depth = maxi )
18          model_fit = xgb_model.fit(X_train , y_train )
19          y_predict = model_fit.predict(X_test)
20          MED=abs(y_test-y_predict).mean()
21
22          if MED<MED_limit:
23              return MED , n, learn, maxi
24          else:
25              return 0
26
27  #### Data Sets ####
28
29  #Train
30  train_df_copy=pd.DataFrame.copy(global_temperatures[
        global_temperatures['dt'].map(lambda x: x.year)<1980])
31  train_df=pd.DataFrame.copy(train_df_copy)
32  train_df.drop(train_df.index[:1], inplace=True)
33  train_df_copy.drop(train_df.index[-1], inplace=True)
34  for j in (train_df.index):
35      for column in train_df.columns[1:-1]:
36          train_df.loc[j,column]=train_df_copy.loc[j-1,column]
37
38  train_df=train_df.drop(['dt'],axis=1)
39
40  #Test
41  test_df_copy=pd.DataFrame.copy(global_temperatures[
        global_temperatures['dt'].map(lambda x: x.year)>=1980])
42  test_df=pd.DataFrame.copy(test_df_copy)
43  test_df.drop(test_df.index[:1], inplace=True)
44  test_df_copy.drop(test_df.index[-1], inplace=True)
45  for j in (test_df.index):
46      for column in test_df.columns[1:-1]:
47          test_df.loc[j,column]=test_df_copy.loc[j-1,column]
48
49  test_df=test_df.drop(['dt'],axis=1)
50
51
52  X_train = train_df.drop("LandAverageTemperature", axis = 1)
53  y_train = train_df['LandAverageTemperature']
54
55  X_test = test_df.drop("LandAverageTemperature", axis = 1)
56  y_test = test_df['LandAverageTemperature']
57
58  #### Hyperparameters ####
59
60  max_depth = [2+i for i in range(5)]
61  learning_rate = [0.01*(i+1) for i in range(5)]
62  n_estimators = [370+i for i in range(20)]
63  MED_limit =  0.21
64  l=[]
65
66  #### Very long loop depending on parameters ####
67
68  for n in n_estimators:
69      for learn in learning_rate:
70          for maxi in max_depth:
```

```
71              a = MED_XGB(n,maxi,learn,X_train,y_train,X_test,
                    y_test,MED_limit)
72              if a != 0:
73                  l.append(a)
74                  l.sort()
75
76  l[:5] # 5 Best MED Results
```

We obtain the following results:

$$[(0.2024272751996246, 385, 0.04, 2),$$
$$(0.2024369978097088, 386, 0.04, 2),$$
$$(0.20244117096460612, 387, 0.04, 2),$$
$$(0.20250207667782244, 384, 0.04, 2),$$
$$(0.20251076595633877, 383, 0.04, 2)]$$

Those results show us that there are hyperparameters that allow us to beat our previos XGBoost result. We want to compare this result with a Grid Search method.

Listing 16: Minimization using Grid Search

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

We obtain the following results:

$$(0.9921251007122502, \text{'learning\_rate'}: 0.02, \text{'max\_depth'}: 2, \text{'n\_estimators'}: 381)$$

The first number being a score, the others, the hyperparameters that minimize the MED. We observe that these hyperparamaters give us an output MED of 0.21173377635362803 which is higher than the XGBoost taking the default hyperparameters and also higher than the MED we obtain with a manual minimization. We thus decide to take the first result of our manual optimization. Taking these hyperparameters to run the XGBoost, we evaluate the importance of the different features in the XGBoost Regressor prediction.

Listing 17: Model Features

```
1
2   #Using the results of our minimization
3   maxi = 2
```
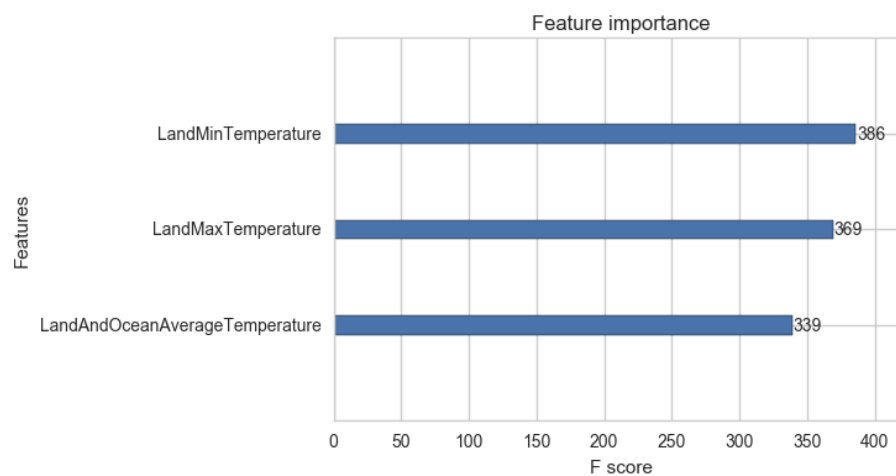
```
 4  learn = 0.04
 5  n = 385
 6
 7  xgb_model = xgboost.XGBRegressor(max_depth = maxi,
 8                                   learning_rate = learn,
 9                                   n_estimators = n)
10
11  model_fit = xgb_model.fit(X_train,y_train)
12
13  #Ranking the features'importance
14  plot_importance(model_fit)
```

Figure 9: Features Importance



We obtain that the most discriminant feature is the LandMinTemperature, closely followed by the LandMaxTemperature. We then plot the most efficient tree of the XGBoost model.
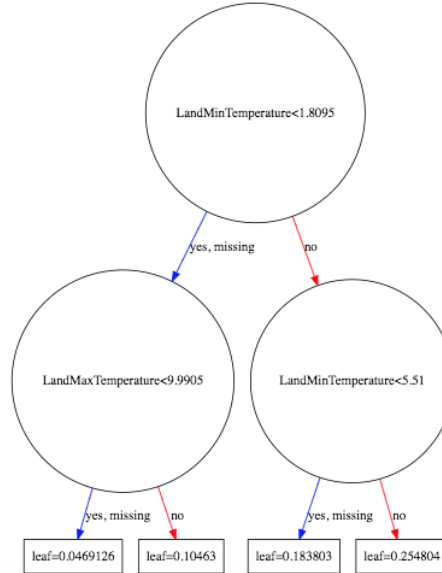
Listing 18: Decision Tree

```
 1  #Plotting an exemple of node and leafs
 2  to_graphviz(model_fit)
```

23

Figure 10: Decision Tree



Naturally, we observe that the most discriminatory feature LandMinTemperature appears at the top of the tree.

We seem to have found hyperparameters that allow us to minimize significantly our MED, as well as identified the most discriminatory features of the model. Now, we evaluate the impact of the introduction of a new variable, the CO2 level into our analysis.

## 4.8 Impact of a new variable on our models

We decide to add a CO2 variable to our dataset. The CO2 data[5] spans from 1958 to 2016. We drop the 2016 year because our initial database stops in 2015. We also decide to assess the impact on our two most efficient models to date: Random Forest Regressor and XGBoost Regressor. Here we use directly the optimization techniques seen before, in order to find the most efficient hyperparameters. Before, we import the new data and merge them with the initial dataset.

Listing 19: Import CO2 variable

```
1  # Adding the CO2 level as a a feature
2
3  CO2=pd.read_csv("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)
```

[5]http://data.okfn.org/data/core/

24

```
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()
```

We can now use our predictive models on the new dataset. This time, we decide to directly use the hyperparameters optimization techniques. The set on which we are minimizing was chosen after having done several trials.

### 4.8.1 Random Forest

Listing 20: Random Forest Optimization

```
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']
20
21   #Definition of the Minimization Function
22
23   def MED_RF(n,maxi,
24           X_train, y_train, X_test,y_test,
25           MED_limit):
26
27   #Note that there is no learning_rate for the random forest
         regressor
28
29       '''
30       n        : n_estimators  -> Number of boosted trees to fit
31       maxi     : max_depth     -> Maximum tree depth for base
             learners
32       ...
33       MED_limit : -> threshold for minimum MED required
34       '''
35
36       forest = RandomForestRegressor(n_estimators=n,max_depth =
             maxi)
37       model_fit = forest.fit(X_train , y_train )
```

25

```
38    y_predict = model_fit.predict(X_test)
39    MED=abs(y_test-y_predict).mean()
40
41    if MED<MED_limit:
42        return MED, n, maxi
43    else:
44        return 0
45
46 #Hyperparameters
47
48 max_depth = [8+i for i in range(7)]
49 n_estimators = [i+1 for i in range(100)]
50 MED_limit = 0.2
51 l=[]
52
53 #### Minimization (Very long loop depending on parameters) ####
54
55 for n in n_estimators:
56     for maxi in max_depth:
57         a = MED_RF(n,maxi,X_train,y_train,X_test,y_test,
58                 MED_limit)
58         if a != 0:
59             l.append(a)
60             l.sort()
61
62 l[:5] # 5 Best MED Results
63
64 # GridSearch (very long)
65
66 clf = GridSearchCV(forest,
67                     {'max_depth': max_depth,
68                      'n_estimators': n_estimators}, verbose=1)
69
70 clf.fit(X_train,y_train)
71 clf.best_score_ , clf.best_params_
```

We obtain the following output for the manual minimization:

$$[(0.12805185110461015, 23, 13),$$
$$(0.13384931231904434, 36, 9),$$
$$(0.1344451730051446, 92, 13),$$
$$(0.1360461221964427, 25, 10),$$
$$(0.1361788863109051, 20, 13)]$$

And for the GridSearch we obtain:

$$(0.9994277108488541, \text{'max\_depth': } 14, \text{'n\_estimators': } 54)$$

Note that we are minimizing only two hyperparameters. Contrary to the XG-Boost, Random Forest does not have a learning_rate hyperparameter. We test the parameters obtained with the Grid Search.

Listing 21: Random Forest-Grid Search Testing

```
1 #Random Forest (GridSearch MED Test)
2
3 # To predict each average land temperature we use the previous
      year's temperatures (landandocean, min , max) and the CO2
      level of the current year
```

```
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=14,n_estimators=54)
8   train_df=train_df.drop(['dt',"Date"],axis=1)
9   forest = forest.fit(train_df.drop(["LandAverageTemperature"],
        axis=1),train_df["LandAverageTemperature"])

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))

17  # Distance computing
18  test_df['Predicted_temp']=output
19  MED=abs(test_df['LandAverageTemperature']-test_df['
        Predicted_temp']).mean()
20  MED
```

The result is: 0.15592940620434984. The GridSearch is again less efficient than the manual minimization, whose top result was 0.12805185110461015. Without the CO2 variable, the best Random Forest MED was 0.21886310904872408. The introduction of the CO2 variable decreases substantially our MED, which means that the model, here the Random Forest Regressor, better fits the new, richer data set. We turn to the XGB Regressor.

### 4.8.2  XGBoost

Listing 22: XGBoost Optimization

```
1   #Sets
2
3   # Train
4   train_df_copy_2=pd.DataFrame.copy(temperatures[temperatures['dt'
        ].map(lambda x: x.year)<1980])
5   train_df_2=pd.DataFrame.copy(train_df_copy_2)
6   train_df_2.drop(train_df_2.index[:1], inplace=True)
7   train_df_copy_2.drop(train_df_2.index[-1], inplace=True)
8
9   # Test
10  test_df_copy_2=pd.DataFrame.copy(temperatures[temperatures['dt'
        ].map(lambda x: x.year)>=1980])
11  test_df_2=pd.DataFrame.copy(test_df_copy_2)
12  test_df_2.drop(test_df_2.index[:1], inplace=True)
13  test_df_copy_2.drop(test_df_2.index[-1], inplace=True)
14
15  #Hyperparameters
16  max_depth = [2+i for i in range(4)]
17  learning_rate = [0.01+0.01*i for i in range(5)]
18  n_estimators = [450+i*10 for i in range(10)]
19  MED_limit = 2
20  l=[]
21
22  #### Minimization (Very long loop depending on parameters) ####
23
24  for n in n_estimators:
25      for learn in learning_rate:
```

```python
26            for maxi in max_depth:
27                a = MED_XGB(n,maxi,learn,X_train,y_train,X_test,
                       y_test,MED_limit)
28                if a != 0:
29                    l.append(a)
30                    l.sort()
31
32  l[:5] # 5 Best MED Results
33
34  # GridSearch (very long)
35
36  clf = GridSearchCV(forest2,
37                     {'max_depth': max_depth,
38                      'n_estimators': n_estimators,
39                      'learning_rate': learning_rate}, verbose=1)
40
41  clf.fit(X_train,y_train)
42  clf.best_score_ , clf.best_params_
43
44  #XGBoost (GridSearch MED Test)
45
46  # To predict each average land temperature we use the previous
         year's temperatures (landandocean, min, max) and the CO2
         level of the current year
47
48  for j in (train_df_2.index):
49      for column in train_df_2.columns[1:-2]:
50          train_df_2.loc[j,column]=train_df_copy_2.loc[j-1,column]
51  forest2 = XGBRegressor(learning_rate=0.01, max_depth=4,
         n_estimators=450)
52  forest2 = forest2.fit(train_df_2.drop(["LandAverageTemperature",
         "dt","Date"],axis=1),train_df_2["LandAverageTemperature"])
53
54
55  for j in (test_df_2.index):
56      for column in test_df_2.columns[1:-2]:
57          test_df_2.loc[j,column]=test_df_copy_2.loc[j-1,column]
58  output = forest2.predict(test_df_2.drop(["LandAverageTemperature
         ","dt","Date"],axis=1))
59
60  # Distance computing
61  test_df_2['Predicted_temp']=output
62  MED=abs(test_df_2['LandAverageTemperature']-test_df_2['
         Predicted_temp']).mean()
63  MED
```

The results are:

$$[(0.2016687891311823, 460, 0.04, 2),$$
$$(0.20183603119020402, 470, 0.04, 2),$$
$$(0.20186279160130058, 450, 0.04, 2),$$
$$(0.20253749801830465, 480, 0.04, 2),$$
$$(0.20266475269025555, 490, 0.04, 2)]$$

The results or GridSearch Optimization are:

$$(0.9227900204652463, \text{'learning\_rate'}: 0.01, \text{'max\_depth'}: 4, \text{'n\_estimators'}: 450)$$

And the MED XGBoost result:

$$0.3102404183345716$$

We observe that the XGBoost MED results, 0.2016687891311823 at best are sensibly lower than the Random Forest's ones, 0.12805185110461015 at best. Remember that we find above that without the new variable, the XGB Regressor is more efficient than the Random Forest. We explain this new counter-intuitive result by two facts. First, the training database is far smaller than before, because the new variable does not have values before 1958. Second, our computing power is limited and we believe that if we had been able to minimize our function on a bigger set, we would have found more promising results for the XGBRegressor. We conclude that the introduction of the new variable has a significant impact on the MED regarding the Random Forest Regressor. As for the XGB Regressor, the impact on the MED is not significant (0.2016687891311823 here against 0.2024272751996246 without the new variable). We deduce that The Random Forest Regressor better fits our new dataset.

# 5 Country Clustering

In this section, we study the existing similarities in terms of climate path between countries throughout the world. We clean and prepare the dataframe to adapt it to the clustering.

**Listing 23: Cleaning and Preparing the Dataframe**

```python
country_temps = pd.read_csv("GlobalLandTemperaturesByCountry.csv
    ", sep=',')

#Drop Uncertainty for the clustering

country_temps = country_temps.drop("
    AverageTemperatureUncertainty", axis = 1)

# Delete all empty lines, change time, and put time as Index

country_temps['dt'] = pd.to_datetime(country_temps['dt'])

country_temps.tail()

#Setting the Countries as Index, time as Columns and Average
    Temperature as the main Data
pivot = country_temps.pivot("Country","dt","AverageTemperature")

pivot = pivot.drop(pivot.columns[[i for i in range(1274)
    ]+[3238]], axis = 1)

"""
    Data before 1850 is deleted because

    *30%  of Database is NaN if dt starts at 1773
    *18%  of Database is NaN if dt starts at 1800
    *5.2% of Database is NaN if dt starts at 1850

    Numbers obtained doing :
        -> np.count_nonzero(~np.isnan(pivot))

    Plus, the last month (9/2013) doesn't have any data

"""

pivot.head()

#Getting rid of all countries with no data (Only one country
    eliminated)

pivot = pivot.dropna(how='all')

#We decide to replace all the missing values by the median
    temperature of the concerned country. The mean would have
    been too influenced by extreme values.

from sklearn.preprocessing import Imputer

imp = Imputer(strategy="median", axis = 1)
# Beware: here axis = 1 means rows -> http://scikit-learn.org/
    stable/modules/generated/sklearn.preprocessing.Imputer.html
```

```
46  #Imputer changes type of DataFrame.
47
48  trans_pivot = imp.fit_transform(pivot)
49  nonan_pivot = pd.DataFrame(trans_pivot)
50
51  nonan_pivot.columns = pivot.columns
52  nonan_pivot.index = pivot.index
53
54  nonan_pivot.tail()
55
56  print(np.count_nonzero(~np.isnan(nonan_pivot)))
57
58  #No missing data, and let's now call the DataFrame : clusterdf
59  clusterdf = nonan_pivot
```

We define our Clustering function, using a scikit-learn package and plot it.

Listing 24: Clustering

```
1   def Clustering(df,x,y,orientation,size):
2
3       """
4
5       df           : DataFrame
6       x,y          : plt dimensions
7       orientation : top, bottom, right or left
8       size         : leaf_font_size
9
10      """
11      from sklearn.cluster import AgglomerativeClustering
12      from scipy.cluster.hierarchy import dendrogram
13      import matplotlib.pyplot as plt
14      plt.style.use('ggplot')
15
16      ward = AgglomerativeClustering(linkage='ward',
              compute_full_tree=True).fit(df)
17      dendro = [ ]
18      for a,b in ward.children_:
19          dendro.append([a,b,float(len(dendro)+1),len(dendro)+1])
20
21      fig = plt.figure( figsize=(x,y) )
22      ax = fig.add_subplot(1,1,1)
23      r = dendrogram(
24          dendro,
25          color_threshold=None,
26          labels=df.index,
27          distance_sort = True,
28          show_leaf_counts=False,
29          leaf_font_size = size,
30          ax=ax,
31          orientation = str(orientation)
32      )
33
34  def Clustering_truncated(df,x,y,orientation,p,size):
35
36      """
37
38      df           : DataFrame
39      x,y          : plt dimensions
40      orientation : top, bottom, right or left
41      p            : Only the last p merged clusters are shown
42      size         : leaf_font_size
```
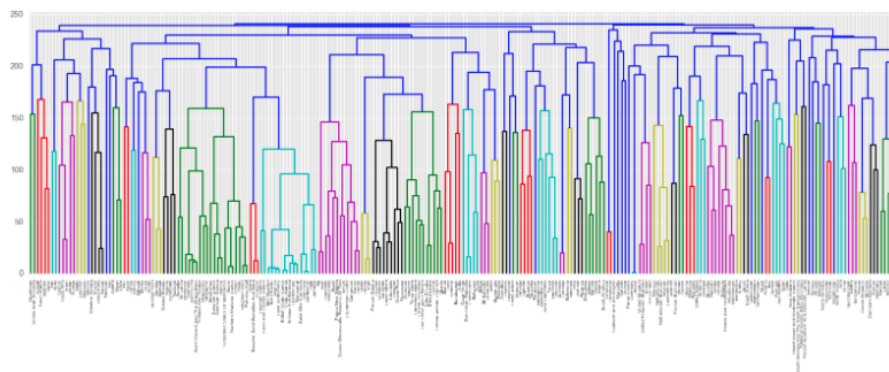
```
43
44        """
45        from sklearn.cluster import AgglomerativeClustering
46        from scipy.cluster.hierarchy import dendrogram
47        import matplotlib.pyplot as plt
48        plt.style.use('ggplot')
49
50        ward = AgglomerativeClustering(linkage='ward',
              compute_full_tree=True).fit(df)
51        dendro = [ ]
52        for a,b in ward.children_:
53            dendro.append([a,b,float(len(dendro)+1),len(dendro)+1])
54
55        fig = plt.figure( figsize=(x,y) )
56        ax = fig.add_subplot(1,1,1)
57        r = dendrogram(
58            dendro,
59            color_threshold=None,
60            labels=df.index,
61            show_leaf_counts=False,
62            truncate_mode = 'lastp',  # show only the last p merged
                  clusters
63            p = p,  # show only the last p merged clusters
64            leaf_font_size = size,
65            show_contracted = True,  # to get a distribution
                  impression in truncated branches
66            ax=ax,
67            orientation = str(orientation)
68
69        )
70        return r["ivl"]
71
72  #We can now plot the clustering
73  Clustering(clusterdf,20,6,"top",6)
```

We obtain the following graph:

Figure 11: First Clustering



We enlarge the graph and put it vertically. Given the important size of the graph, we decide not to display it here. The reader can see it in the notebook. As shown by the previous graph, many countries do not share any similarities with others. The similarity is measured by a certain distance. The more important
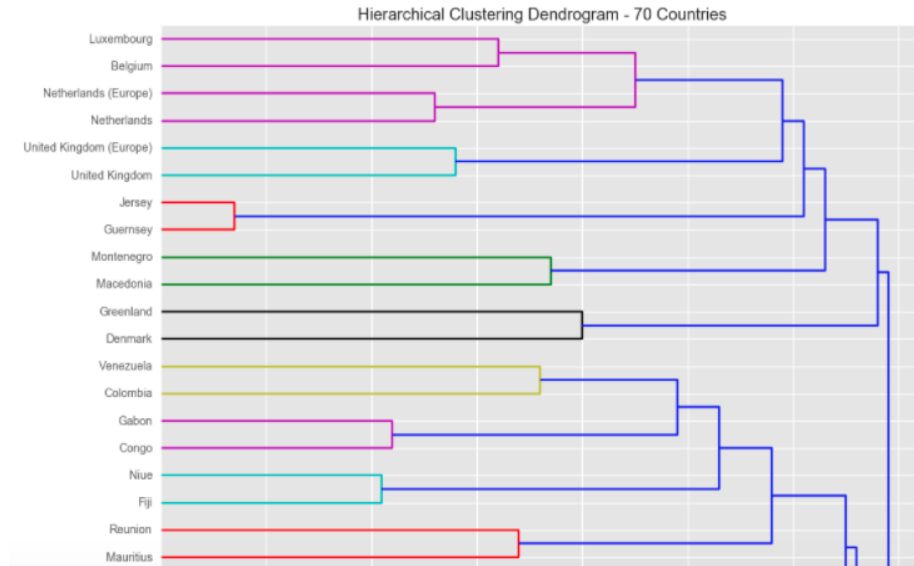
the distance, the less relevant the merging of the countries into clusters is. We thus create a second clustering which gathers only the countries that closely share climate similarities.

```python
#We remove from the cluster the most interesting countries in
    terms of shared climate similarities

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

#We define a function that only keeps the interesting #countries

def delete_symbol(l,sym):
    for i in range(l.count(sym)):
        l.remove(sym)
    return l, len(l)

cntry_to_del = delete_symbol(Cluster_70,'')
cntry_to_del[1]

#We deleted 172 countries
clusterdf_new = clusterdf.drop(cntry_to_del[0], axis = 0)

#We can now plot the new dendrogram
Clustering(clusterdf_new,12,30,"right",10)
plt.title('Hierarchical Clustering Dendrogram - 70 Countries')
plt.xlabel('Distance')
plt.ylabel('Countries')
```

We display the top part of the graph.

Figure 12: Second Clustering



Hierarchical Clustering Dendrogram - 70 Countries

We cannot see it here but the distance axis that spanned from 0 to 250 in the previous graphs stops here at 70. We thus have only the country, whose temperature trends are very similar.

# 6 Zoom in: Paris Temperature Peak?

We use the file GlobalLandTemperaturesByMajorCity.csv and focus on Paris. Once again, we decide to keep only the years from 1800 onwards given the numerous NaN we have before. We design an interactive heating-map using the bokeh package, displaying the temperature in Paris, every month since 1800, classified by colors. The reader can hover over the map on the notebook.

**Listing 26: Zoom in, Paris**

```python
#Import of Paris'data

import pandas as pd
df=pd.read_csv('GlobalLandTemperaturesByMajorCity.csv')
df_Paris=df[df.City == 'Paris']
df_Paris=df_Paris.dropna()
df_Paris=df_Paris.drop(df_Paris.columns[2:], axis=1)

#Starts at 1800

df_Paris=df_Paris.iloc[602:]

#Building of a new Dataframe adapted to the map

date=df_Paris['dt'].apply(lambda x: x[:4])
years = np.unique(date)
mean_temps_Paris=[]

for y in years:

    mean_temps_Paris.append(df_Paris[y==date]['
        AverageTemperature'].mean())

jan=[]
for i in range(0,2562,12):
        jan.append(df_Paris.iloc[i]['AverageTemperature'])
feb=[]
for i in range(1,2562,12):
        feb.append(df_Paris.iloc[i]['AverageTemperature'])
mar=[]
for i in range(2,2562,12):
        mar.append(df_Paris.iloc[i]['AverageTemperature'])
apr=[]
for i in range(3,2562,12):
        apr.append(df_Paris.iloc[i]['AverageTemperature'])
may=[]
for i in range(4,2562,12):
        may.append(df_Paris.iloc[i]['AverageTemperature'])
jun=[]
for i in range(5,2562,12):
        jun.append(df_Paris.iloc[i]['AverageTemperature'])
jul=[]
for i in range(6,2562,12):
        jul.append(df_Paris.iloc[i]['AverageTemperature'])
aug=[]
for i in range(7,2562,12):
        aug.append(df_Paris.iloc[i]['AverageTemperature'])
sep=[]
for i in range(8,2562,12):
```

```python
           sep.append(df_Paris.iloc[i]['AverageTemperature'])
octb=[]
for i in range(9,2562,12):
        octb.append(df_Paris.iloc[i]['AverageTemperature'])
nov=[]
for i in range(10,2562,12):
        nov.append(df_Paris.iloc[i]['AverageTemperature'])
dec=[]
for i in range(11,2562,12):
        dec.append(df_Paris.iloc[i]['AverageTemperature'])

#New Dataframe

df_Paris=pd.DataFrame({'Annual':pd.Series(mean_temps_Paris),
'Year':years,'Jan':pd.Series(jan),
                          'Feb':pd.Series(feb),'Mar':pd.Series(mar)
                          ,
                          'Apr':pd.Series(apr),'May':pd.Series(may)
                          ,
                          'Jun':pd.Series(jun),'Jul':pd.Series(jul),
                          'Aug':pd.Series(aug),'Sep':pd.Series(sep),
                          'Oct':pd.Series(octb),
                          'Nov':pd.Series(nov),'Dec':pd.Series(dec)
                            })

df_Paris=df_Paris[['Year','Jan','Feb','Mar','Apr','May','Jun','
        Jul','Aug',
'Sep','Oct','Nov','Dec','Annual']]
df_Paris=df_Paris.dropna()

#Interactive map with monthly average temperatures in Paris
        since 1800

from collections import OrderedDict
import pandas as pd
import numpy as np
from bokeh.io import output_notebook, show
from bokeh.plotting import (ColumnDataSource, figure)
from bokeh.models import HoverTool

output_notebook()

# Read in the data with pandas. Convert the year column to
        string
df_Paris['Year'] = [str(x) for x in df_Paris['Year']]
years = list(df_Paris['Year'])
months = ["Jan","Feb","Mar","Apr","May","Jun","Jul","Aug","Sep",
        "Oct",
"Nov","Dec"]
df_Paris = df_Paris.set_index('Year')

# this is the colormap from the original plot

colors = ["#75968f", "#75968f", "#a5bab7", "#c9d9d3", "#e2e2e2",
        "#dfccce",
        "#ddb7b1", "#cc7878", "#933b41", "#550b1d","#75968f","
                #75968f"]

# Set up the data for plotting. We will need to have values for
        every
# pair of year/month names. Map the Temperature to a color.
```

```python
102  def mapping_to_colors(x):
103      if x<=0:
104          return 1
105      elif x>=0 and x<4:
106          return 2
107      elif x>=4 and x<8:
108          return 3
109      elif x>=8 and x<12:
110          return 4
111      elif x>=12 and x<14:
112          return 5
113      elif x>=14 and x<19:
114          return 6
115      elif x>=19 and x<22:
116          return 7
117      else:
118          return 8
119
120  month = []
121  year = []
122  color = []
123  Temperature = []
124  for y in years:
125      for m in months:
126          month.append(m)
127          year.append(y)
128          Monthly_Temperature = df_Paris[m][y]
129          Temperature.append(Monthly_Temperature)
130          color.append(colors[mapping_to_colors(
131              Monthly_Temperature)])
131
132  source = ColumnDataSource(
133      data=dict(month=month, year=year, color=color, Temperature=
134          Temperature)
134  )
135
136  TOOLS = "resize,hover,save,pan,box_zoom,wheel_zoom"
137
138  #Display parameters
139  p = figure(title="Paris Temperature (1800 - 2012)",
140      x_range=years, y_range=list(reversed(months)),
141      x_axis_location="above", plot_width=3000, plot_height=500,
142      toolbar_location="left", tools=TOOLS)
143
144  p.rect("year", "month", 1, 1, source=source,
145      color="color", line_color=None)
146
147  p.grid.grid_line_color = None
148  p.axis.axis_line_color = None
149  p.axis.major_tick_line_color = None
150  p.axis.major_label_text_font_size = "5pt"
151  p.axis.major_label_standoff = 0
152  p.xaxis.major_label_orientation = np.pi/3
153
154  #Hovers parameters
155  hover = p.select(dict(type=HoverTool))
156  hover.tooltips = OrderedDict([
157      ('date', '@month @year'),
158      ('Temperature', '@Temperature'),])
159
160  show(p)
```

Figure 13: Paris Heating Map Beginning



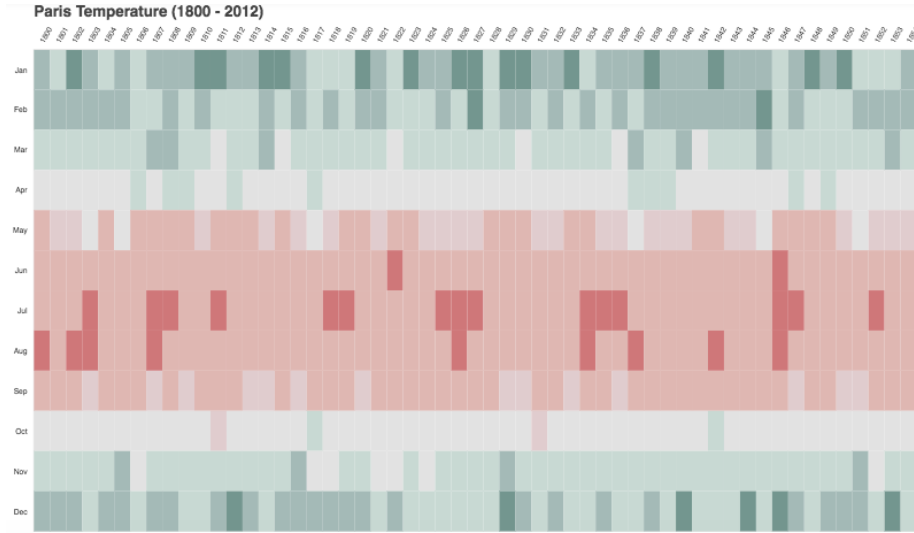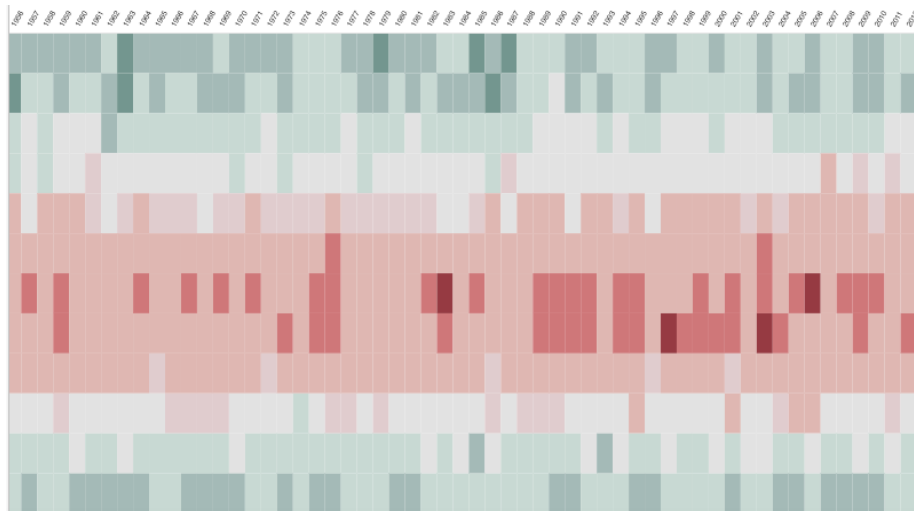**Paris Temperature (1800 - 2012)**

Figure 14: Paris Heating Map Last seventy years



We clearly see that the temperatures are rising in the French capital, with a significant acceleration starting in the 1980's. The center of the map, corresponding to the summer registers more and more extreme values. The borders of the map correponding to the winter, spring and autumn register less and less extremely cold values and becomes more homogeneous. This zoom on Paris echoes the recent pollution peak the city experienced. It is now clear that Paris' temperatures are growing and pollution peaks can only worsen the phenomenon.

# 7 General Conclusion

Following our data exploratory section, we tried out several machine learning models and tested their accuracy. Then we clustered the countries by temperature and finally we plotted an interactive heatmap of Paris. The efficient models are according to the MED, the Random Forest Regressor and the XGBoost Regressor regarding the initial dataset. When adding a new variable, we observe that the Random Forest Regressor produces better results. We regret the very slow process of minimization of the MED, whether it is through the GridSearch or via manual optimization. A more powerful computing power would have certainly helped us obtain finer results. There are no reason why the model would have a limited lifespan, except if the global warming becomes boiling and make human disappear. And we find no reason why the model would lose efficiency with time. Besides, our models can perfectly be trained with higher datasets.

In order to spare costs, we could only take the yearly temperature data. Our models could be used and executed by artificially intelligent robots if we automatise the data scraping and the structuration of the dataframes. It is also possible to detect the failing data updates with a simple findnan function scanning through the new data.

Regarding the business opportunity surrounding our models. We think an application would not spark interest because of the numerous weather related app already existing. We could nonetheless sell this model to industries that do not trust the already existing weather-predicting models and want to have another opinion. Agriculture that heavily relies on the climate could take advantage of our predictive models. Hedge Funds or Banks that trade on complex weather-derivatives or simply on commodities-related derivatives could also take advantage of these predictive models. An Excel extension could be proposed.

With the CO2 variable included, our models become more accurate. The MED is around 0.1 which is very small. In order to further reduce it, we have to add more variables that have relationships with the evolution of the temperatures. The level of human related pollutions generally, the deforestation, the ocean levels could enrich our models and reduce the errors. Nowadays, climate scientists actually use those kind of models, including many more variables.[3]

# List of Figures

# References

[1] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*, volume 1. Springer series in statistics Springer, Berlin, 2001.

[2] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An introduction to statistical learning*, volume 6. Springer, 2013.

[3] Claire Monteleoni, Gavin A Schmidt, and Scott McQuade. Climate informatics: accelerating discovering in climate science with machine learning. *Computing in Science & Engineering*, 15(5):32–40, 2013.

[4] Sebastian Raschka. *Python Machine Learning*. Packt Publishing Ltd, 2015.

[5] Willi Richert. *Building machine learning systems with python*. Packt Publishing Ltd, 2013.