


▼ Basic imports and function ceil

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from statsmodels.tsa.stattools import adfuller #to see if timeseries is stationary
4 import numpy as np
5 import statsmodels.formula.api as smf
6 import statsmodels.api as sm
7 import datetime
```

```
1 # Fonction pour majorer ou minorer un nombre à l'entier le plus proche
2
3 def my_ceil(predictions):
4     for i in range(len(predictions)):
5         if predictions[i]%1<=0.5:
6             predictions[i] = int(predictions[i])
7         else:
8             predictions[i] = int(predictions[i]) + 1
9     return predictions
```

```
1 import io
2 from google.colab import files
3 uploaded = files.upload()
```



Choose Files No file chosen

Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.


Saving test_input.csv to test_input.csv
Saving train_input.csv to train_input.csv
Saving train_output.csv to train_output.csv

▼ Dataset Exploration

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import numpy as np
5 import datetime
```

```
1 %%capture
2 train_input= pd.read_csv('https://raw.githubusercontent.com/nisrineha/Challenge_data_well_being_at_work/master/datasets/train_input.csv')
3 train_output= pd.read_csv('https://raw.githubusercontent.com/nisrineha/Challenge_data_well_being_at_work/master/datasets/train_output.csv')
4 train= train_input
5 train['Score']= train_output.Score
6 train.head()
7 for i in range(len(train['Date'])):
8     train['Date'][i] = datetime.datetime.strptime(train['Date'][i], '%Y-%m-%d %H:%M:%S')
```

```
1 print(train.shape)
2 print(train.head())
3
```

 (8000, 8)

	ID	Date	Temperature	Humidity	Humex	CO2	Bright	Score
0	0	2017-08-31 23:30:00	22.7	56.0	25.7	534.0	1.0	4.0
1	1	2017-09-01 00:30:00	22.8	55.0	25.7	506.0	1.0	4.0
2	2	2017-09-01 01:30:00	22.9	55.0	25.9	577.0	1.0	4.0
3	3	2017-09-01 02:30:00	23.0	55.0	26.1	630.0	1.0	2.0
4	4	2017-09-01 03:30:00	23.0	55.0	26.1	643.0	1.0	3.0

Verifying if any na value in our dataframe, the result of this query shows us that there is none.

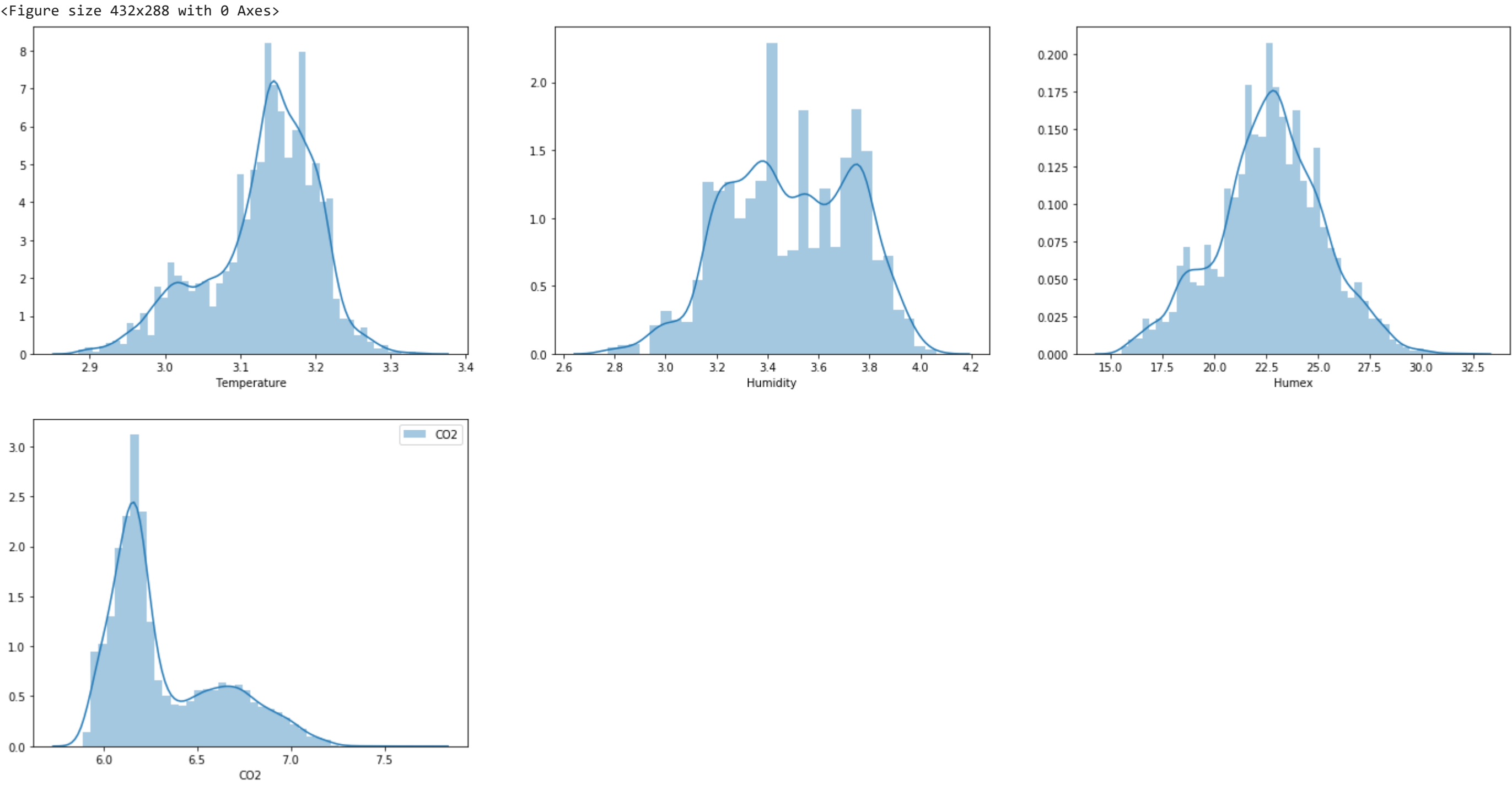
```
1 train.isnull().values.any()
```

 False

▼ Distribution of different features:

```
1 plt.figure("Distribution Plots")
2 fig = plt.figure(figsize = (24,12))
3 plt.subplot(2, 3, 1)
4 sns.distplot(np.log(train.Temperature), label = 'Temperature')
5 plt.subplot(2, 3, 2)
6 sns.distplot(np.log(train.Humidity), label = 'Humidity')
7 plt.subplot(2, 3, 3)
8 sns.distplot(train.Humex, label = 'Humex')
9 plt.subplot(2, 3, 4)
10 sns.distplot(np.log(train.CO2), label = 'CO2')
11 # sns.distplot(well_BE.month, label = 'month')
12 plt.legend()
13 plt.show()
```

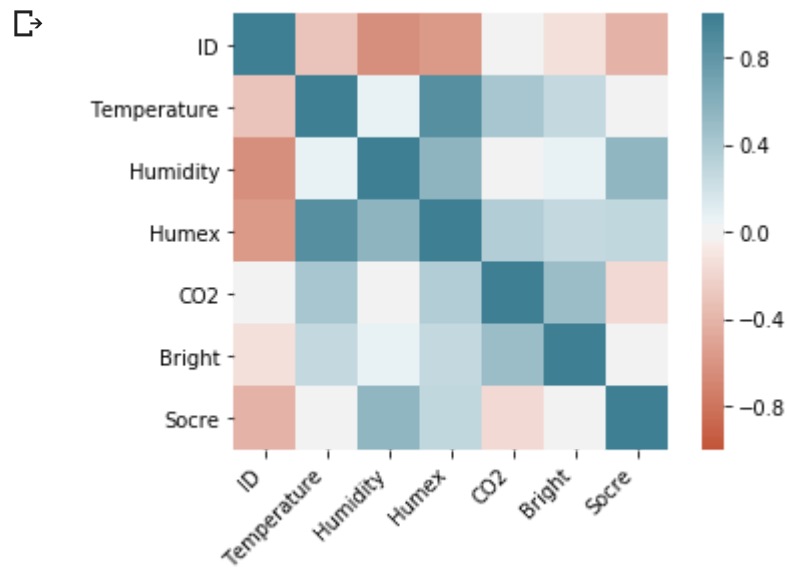




We have the only varibale who has a normal distribution is the Humex variable which represents the air qualtiy

▼ Matrice de correlation

```
1 corr = train.corr()
2 ax = sns.heatmap(
3     corr,
4     vmin=-1, vmax=1, center=0,
5     cmap=sns.diverging_palette(20, 220, n=200),
6     square=True
7 )
8 ax.set_xticklabels(
9     ax.get_xticklabels(),
10    rotation=45,
11    horizontalalignment='right'
12 );
```



```
1 corr
```

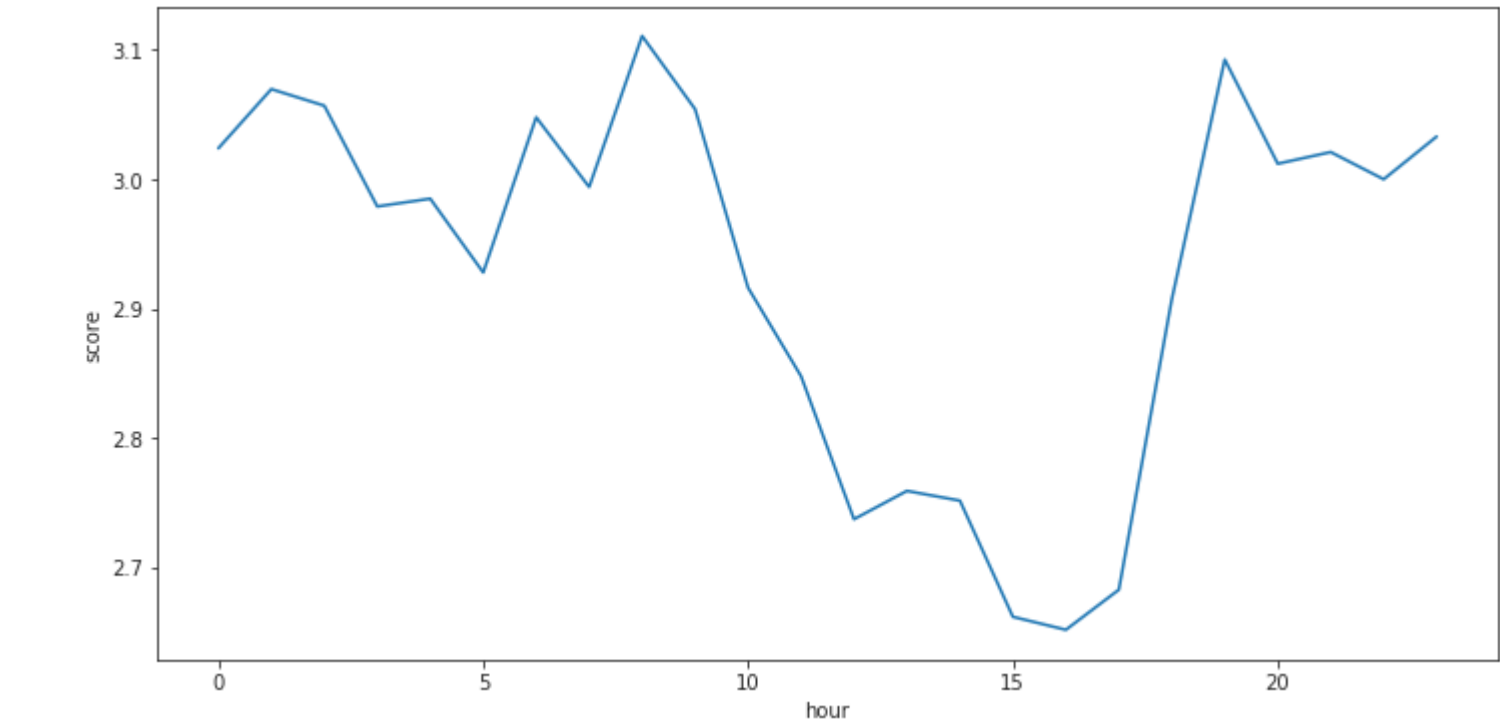
	ID	Temperature	Humidity	Humex	CO2	Bright	Socre
ID	1.000000	-0.306515	-0.630485	-0.571340	0.006029	-0.131696	-0.416685
Temperature	-0.306515	1.000000	0.071968	0.861334	0.426742	0.275884	0.028900
Humidity	-0.630485	0.071968	1.000000	0.566265	0.015807	0.075624	0.541904
Humex	-0.571340	0.861334	0.566265	1.000000	0.361522	0.268284	0.287406
CO2	0.006029	0.426742	0.015807	0.361522	1.000000	0.486382	-0.172919
Bright	-0.131696	0.275884	0.075624	0.268284	0.486382	1.000000	-0.024523
Socre	-0.416685	0.028900	0.541904	0.287406	-0.172919	-0.024523	1.000000

Again, there is only one variable **Humidty** which is hihgly correlated with our traget variable **score**

▼ Relation between our target variable and the variable Date

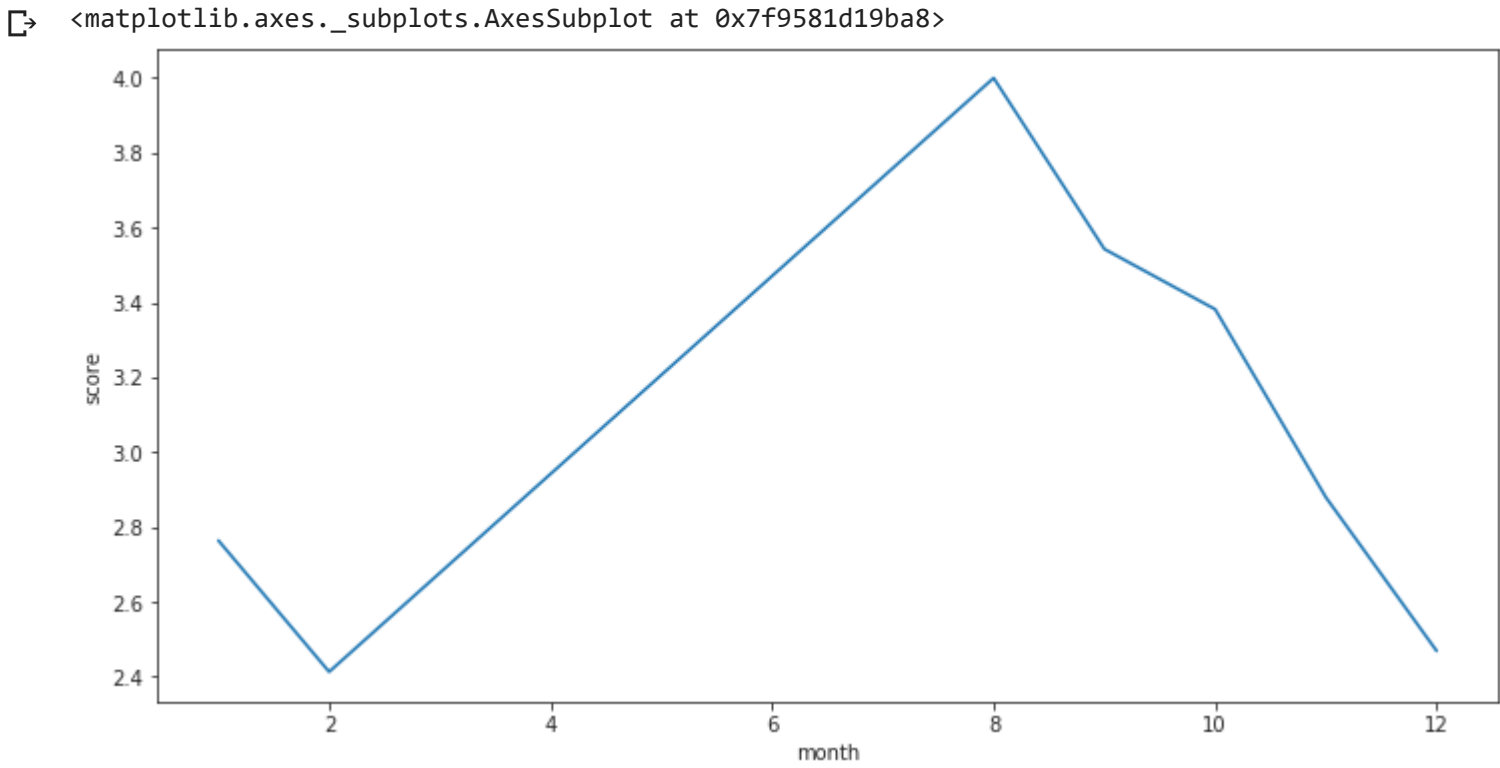
```
1 series = pd.Series(np.array(train.Score), index=train.Date)
2 groupHour = series.groupby(series.index.hour).mean()
3 groupHour = pd.DataFrame({'hour':groupHour.index, 'score':groupHour.values})
4
5 # plt.plot(groupHour.hour, groupHour.score)
6 fig, ax = plt.subplots(1,1, figsize = (12,6))
7 #sns.barplot(x = 'hour', y = 'score', data = groupHour)
8
9 sns.lineplot(x = 'hour', y = 'score', data = groupHour)
10 sns.set_style("ticks", {"xtick.major.size": 16, "ytick.major.size": 8})
```





From the graph before which represents the score in function of day hours, the confort subjectif drop at 9:00 Am and raise again at 19:00, maybe people are working at this hours and they are less comfortable.

```
1 series = pd.Series(np.array(train.Score), index=train.Date)
2 groupHour = series.groupby(series.index.month).mean()
3 groupHour = pd.DataFrame({'month':groupHour.index, 'score':groupHour.values})
4
5 fig, ax = plt.subplots(1,1, figsize = (12,6))
6 #sns.barplot(x = 'month', y = 'score', data = groupHour)
7 sns.lineplot(x = 'month', y = 'score', data = groupHour)
```

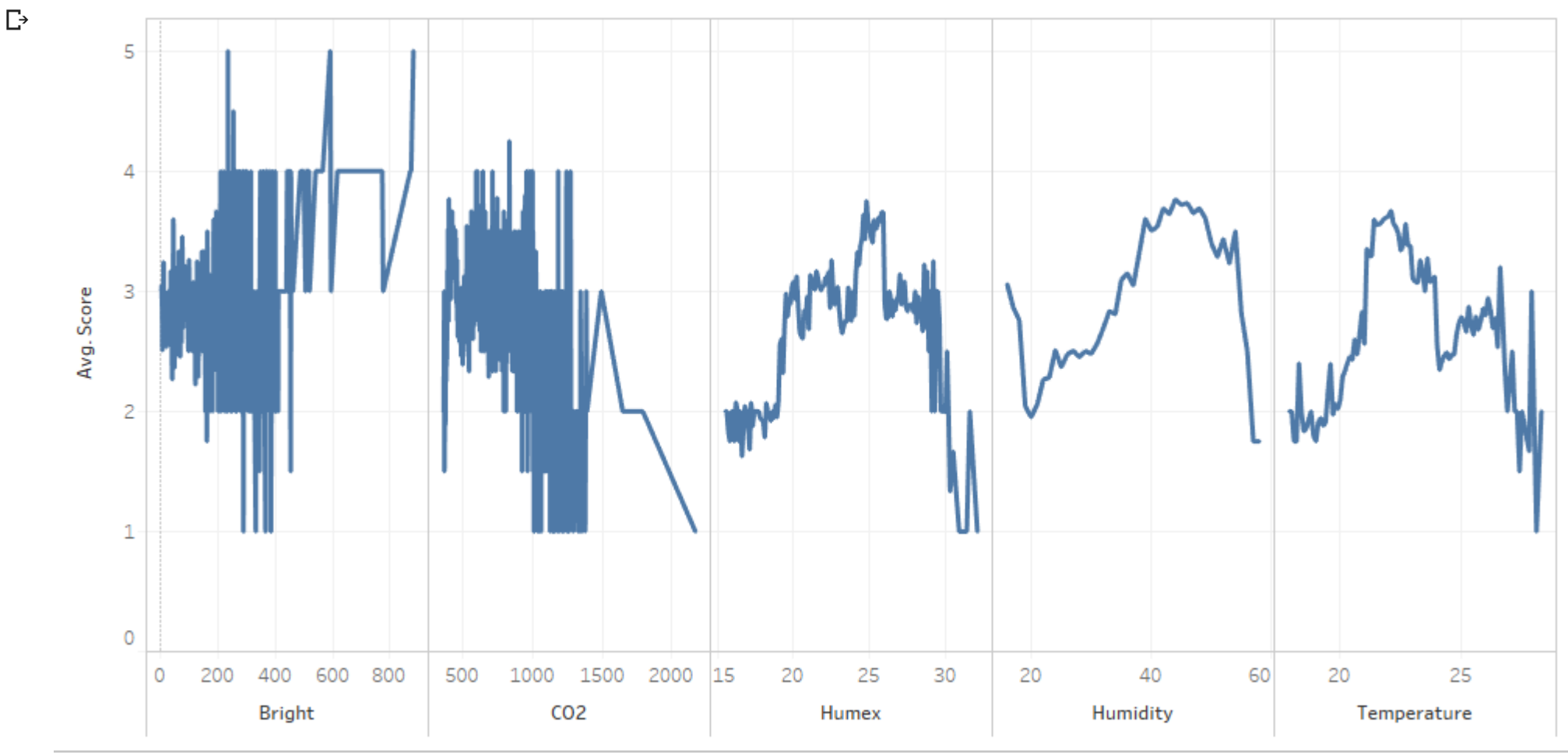


We did a variation of the score in function of months, but it is insignificant, since we have only few months in our dataset

▼ Tableau software

In order to have a better idea about our dataset we used the software Tableau

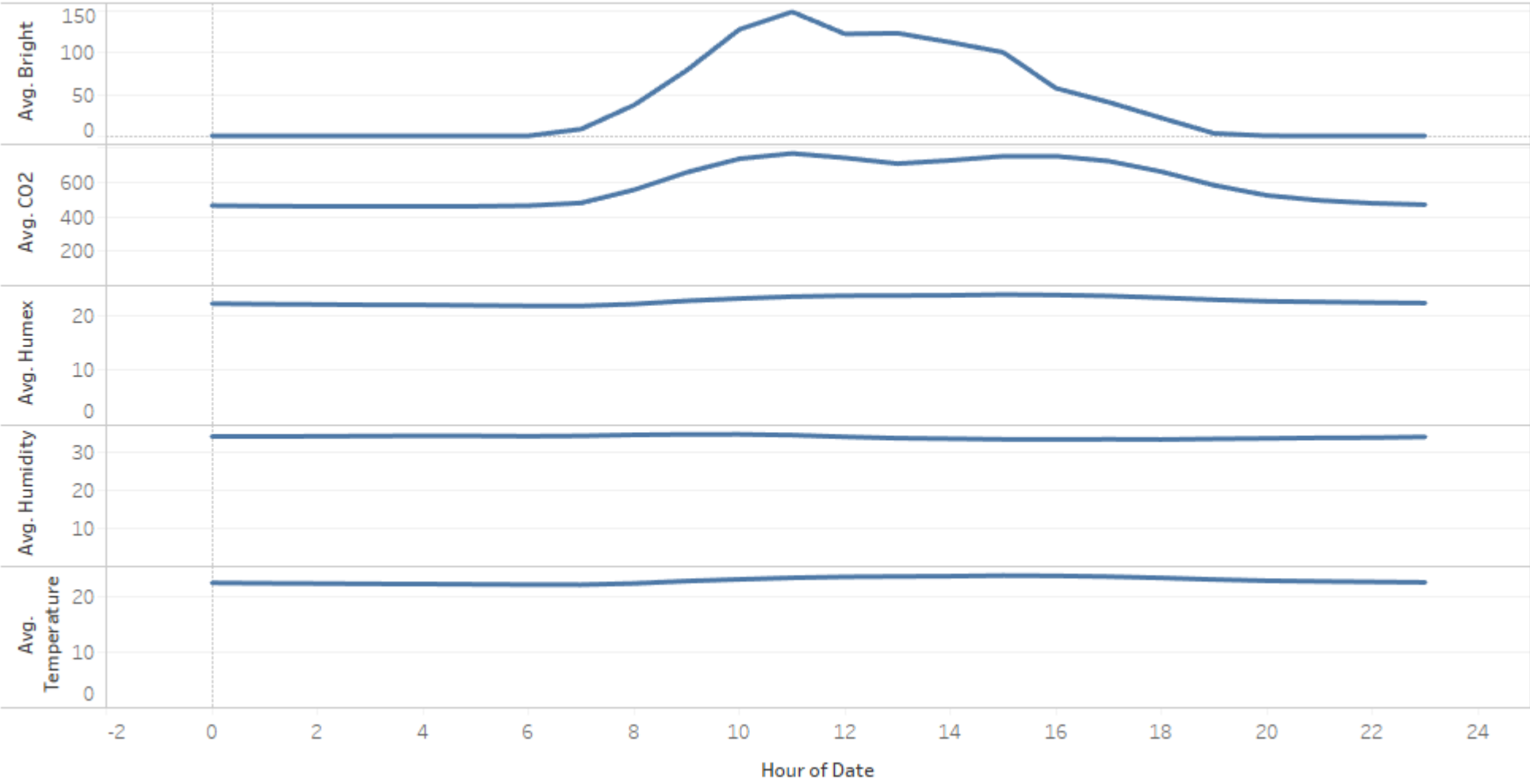
```
1 from IPython.core.display import Image, display
2 display(Image('/content/score(autresvariables).PNG'))
```



The graph shows us that there is no apprent relation between our traget and features

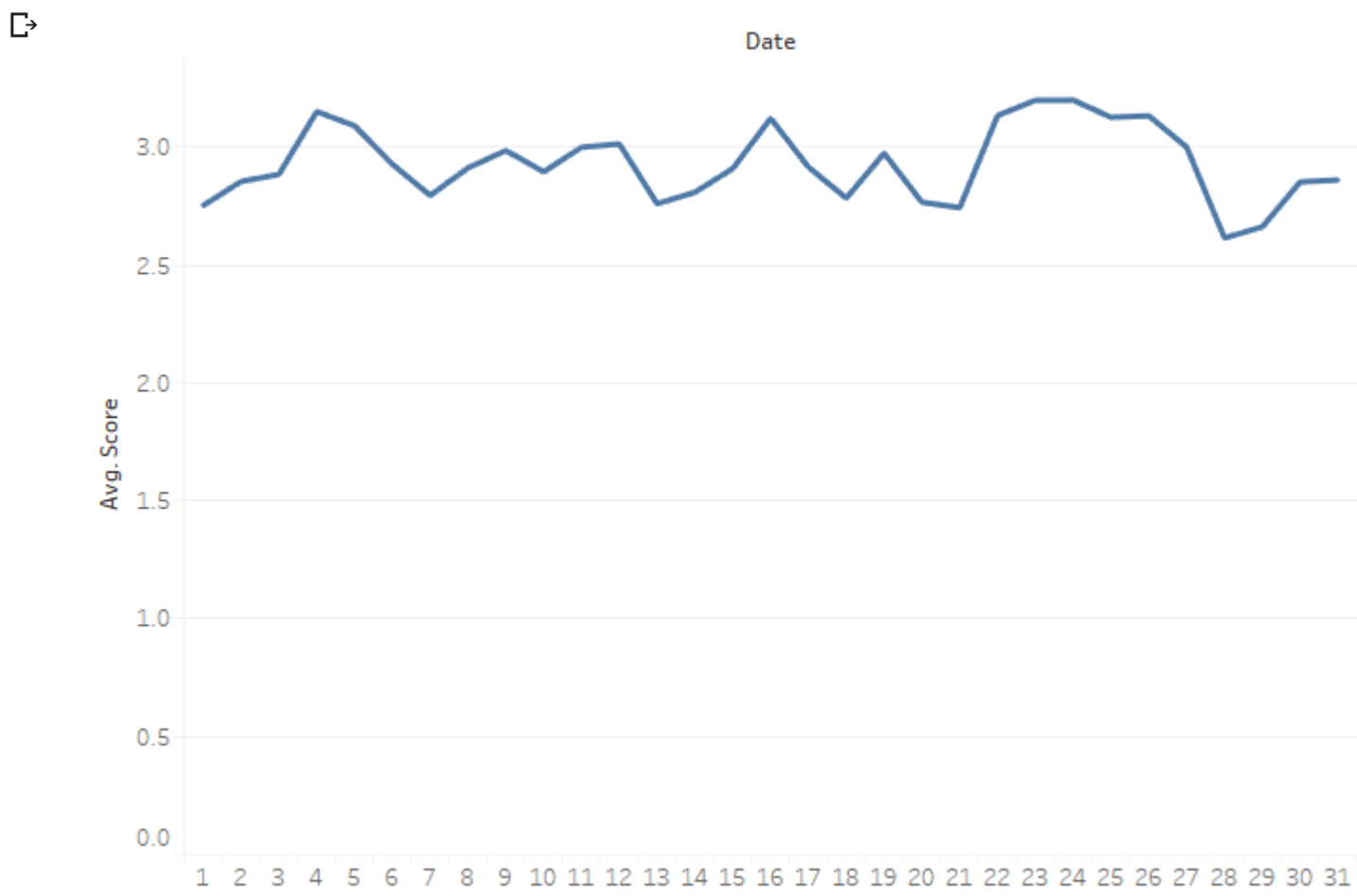
```
1 display(Image('/content/variable(hours).PNG'))
2
```





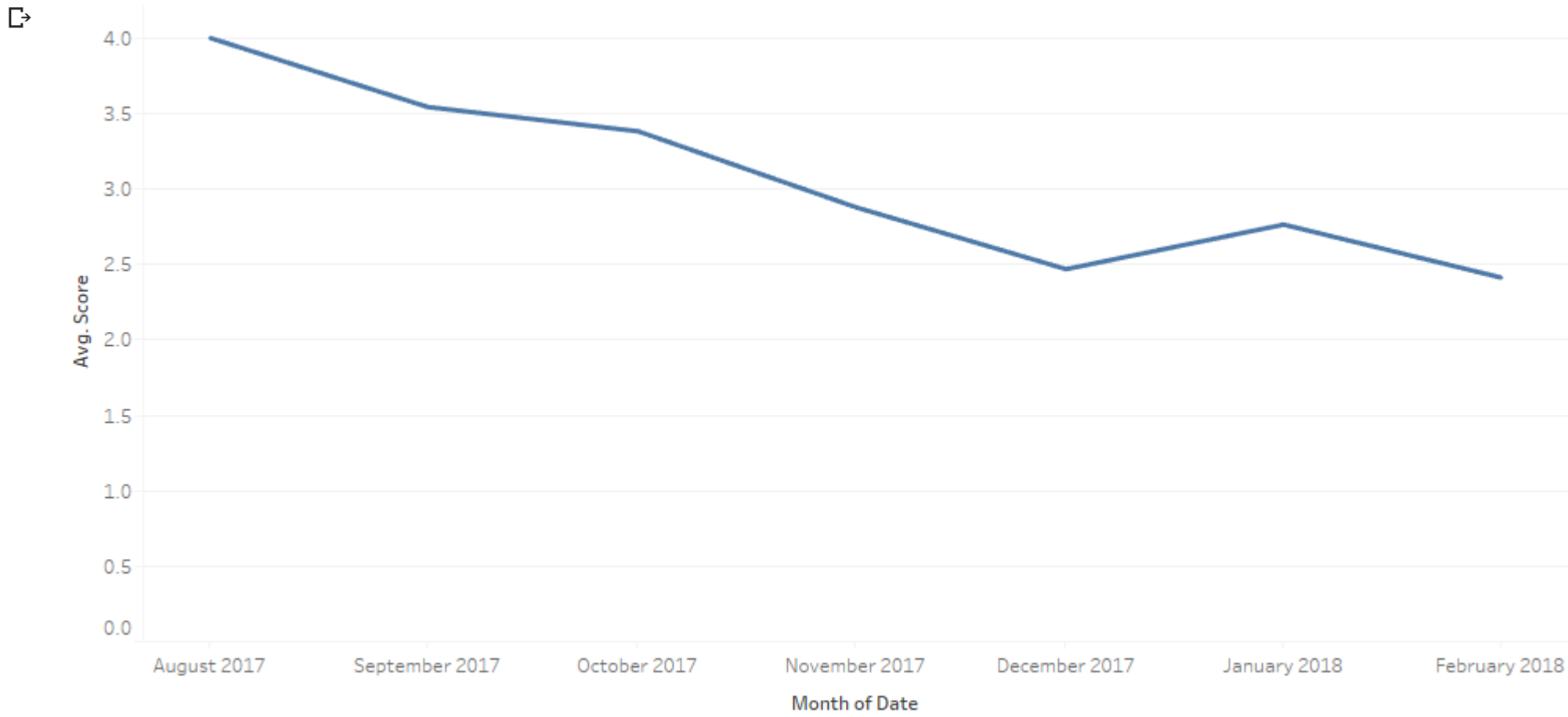
The graph shows us that except the bright and CO2, the other features are constant during the day.

```
1 display(Image('/content/score(day).PNG'))
2
```



this graph shows us that there is no apparent seasonality concerning the score and the time.

```
1 display(Image('/content/score(month).PNG'))
```



This graph shows us that the score is dropping from August to Decemeber and maybe of the season from summer to winter.

Conclusion

After exploring of our data, in the next section we will preprocess the data for modelisation.


Preprocessing phase

Train data

Import train data

```
1 wellB_in = pd.read_csv(io.BytesIO(uploaded['train_input.csv']))
2 wellB_out = pd.read_csv(io.BytesIO(uploaded['train_output.csv']))
3 well_B = wellB_in
.
```

```
4
5 # Convert Date column from String to Date
6 for i in range(len(well_B['Date'])):
7     well_B['Date'][i] = datetime.datetime.strptime(well_B['Date'][i], '%Y-%m-%d %H:%M:%S')
8
9 train_ID = well_B.ID
10 train_Date = well_B.Date
11
12 well_B = well_B.drop(['ID','Date'], axis = 1)
```


 /usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:7: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy.

import sys

▼ Add weekdays to columns as dummies

```
1 well_B['weekday'] = train_Date
2
3 # Add weekday column to data
4 for i in range(len(well_B.weekday)):
5     well_B.weekday[i] = well_B.weekday[i].weekday()
```

 /usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:5: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy.


"""

```
1 dummy_weekday = pd.get_dummies(well_B['weekday'])
2 dummy_weekday.columns = ['lundi','mardi','mercredi','jeudi','vendredi','samedi','dimanche']
3
4 # dummy_weekday.rename(columns = {'0':'lundi','1':'mardi','2':'mercredi','3':'jeudi','4':'vendredi','5':'samedi','6':'dimanche'})
5 # dummy_weekday.head()
6
7 well_B = pd.concat([well_B,dummy_weekday], axis=1)

1 # Drop Weekday column with categorical values
2 well_B = well_B.drop('weekday', axis = 1)
```

▼ Add hours to columns as dummy variables

```
1 well_B['hour'] = train_Date
2
3 for i in range(len(well_B.hour)):
4     well_B.hour[i] = well_B.hour[i].hour
5
6 dummy_hour = pd.get_dummies(well_B['hour'])
7 dummy_hour.columns = ['midnight','AM1','AM2','AM3','AM4','AM5','AM6','AM7','AM8','AM9','AM10','AM11','midday','PM1','PM2','PM3','PM4','PM5','PM6','PM7','PM8','PM9','PM10','PM11']
8
9 well_B = pd.concat([well_B,dummy_hour], axis = 1)
10
11 well_B = well_B.drop('hour', axis = 1)
12
13
14 well_B.head()
```

 /usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:4: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy.

after removing the cwd from sys.path.


	Temperature	Humidity	Humex	CO2	Bright	lundi	mardi	mercredi	jeudi	vendredi	samedi	dimanche	midnight	AM1	AM2	AM3	AM4	AM5	AM6	AM7	AM8	AM9	AM10	AM11	midday	PM1	PM2
0	22.7	56.0	25.7	534.0	1.0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
1	22.8	55.0	25.7	506.0	1.0	0	0	0	0	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	
2	22.9	55.0	25.9	577.0	1.0	0	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	
3	23.0	55.0	26.1	630.0	1.0	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	
4	23.0	55.0	26.1	643.0	1.0	0	0	0	0	1	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	

```
1 # Add objective variable to data
2 well_B['Score'] = wellB_out.Score
```

▼ Test Data

▼ Import test data

```
1 well_Btest = pd.read_csv(io.BytesIO(uploaded['test_input.csv']))
2
3 # Convert Date column from String to Date
4 for i in range(len(well_Btest['Date'])):
5     well_Btest['Date'][i] = datetime.datetime.strptime(well_Btest['Date'][i], '%Y-%m-%d %H:%M:%S')
6
7 test_ID = well_Btest.ID
8 test_Date = well_Btest.Date
9
10 well_Btest = well_Btest.drop(['ID','Date'], axis = 1)
11 # well_Btest_ = well_Btest
```

 /usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:5: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy.


"""

▼ Add weekdays to columns as dummies

```
1 well_Btest['weekday'] = test_Date
2
3 # Add weekday column to data
4 for i in range(len(well_Btest.weekday)):
```



```
5 well_Btest.weekday[i] = well_Btest.weekday[i].weekday()
```

 /usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:5: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame


See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy.

```
1 dummy_weekday = pd.get_dummies(well_Btest['weekday'])
2 dummy_weekday.columns = ['lundi', 'mardi', 'mercredi', 'jeudi', 'vendredi', 'samedi', 'dimanche']
3
4 # dummy_weekday.rename(columns = {'0':'lundi', '1':'mardi', '2':'mercredi', '3':'jeudi', '4':'vendredi', '5':'samedi', '6':'dimanche'})
5 # dummy_weekday.head()
6
7 well_Btest = pd.concat([well_Btest,dummy_weekday], axis=1)

1 # Drop Weekday column with categorical values
2 well_Btest = well_Btest.drop('weekday', axis = 1)
```

▼ Add hours to columns as dummy variables

```
1 well_Btest['hour'] = test_Date
2
3 for i in range(len(well_Btest.hour)):
4     well_Btest.hour[i] = well_Btest.hour[i].hour
5
6 dummy_hour = pd.get_dummies(well_Btest['hour'])
7 dummy_hour.columns = ['midnight', 'AM1', 'AM2', 'AM3', 'AM4', 'AM5', 'AM6', 'AM7', 'AM8', 'AM9', 'AM10', 'AM11', 'midday', 'PM1', 'PM2', 'PM3', 'PM4', 'PM5', 'PM6', 'PM7', 'PM8', 'PM9', 'PM10', 'PM11']
8
9 well_Btest = pd.concat([well_Btest,dummy_hour], axis = 1)
10
11 well_Btest = well_Btest.drop('hour', axis = 1)
12
13 well_Btest.head()
```

 /usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:4: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: http://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy.
after removing the cwd from sys.path.

	Temperature	Humidity	Humex	CO2	Bright	lundi	mardi	mercredi	jeudi	vendredi	samedi	dimanche	midnight	AM1	AM2	AM3	AM4	AM5	AM6	AM7	AM8	AM9	AM10	AM11	midday	PM1	PM2
0	20.1	26.0	17.9	377.0	1.0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	20.2	26.0	18.0	374.0	1.0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	20.1	26.0	17.9	379.0	1.0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	20.1	26.0	17.9	380.0	1.0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	20.0	26.0	17.8	379.0	1.0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

▼ K-fold Cross validation: regression and classification

▼ Imports and functions

```
1 from sklearn.model_selection import KFold # import KFold
2 from sklearn.metrics import accuracy_score #To calculate accuracy
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.neighbors import KNeighborsClassifier
5 from sklearn.preprocessing import StandardScaler
```

```
1 kf = KFold(n_splits = 7)
```

```
1 def cross_validation_XG_classifier(model):
2     accuracy = []
3     for train_index, test_index in kf.split(well_B):
4         # print("TRAIN:", train_index, "TEST:", test_index)
5         #train part
6         X_train = np.array(well_B.drop('Score', axis = 1))[train_index]
7         y_train = np.array(well_B.Score)[train_index]
8         #test part
9         X_test = np.array(well_B.drop('Score', axis = 1))[test_index]
10        y_test = np.array(well_B.Score)[test_index]
11
12        # Define the scaler
13        # scaler = StandardScaler().fit(X_train)
14
15        # X_train = scaler.transform(X_train)
16        # X_test = scaler.transform(X_test)
17
18        # rgr_time = RandomForestClassifier(n_estimators = 400, random_state = 0, max_depth=20)
19        model.fit(X_train, y_train)
20
21        y_pred = model.predict(X_test)
22
23        y_pred = my_ceil(np.array(y_pred))
24
25        accuracy.append(accuracy_score(y_pred, y_test))
26        # print("The score of the " + str(cpt) + " is " + str(accuracy_score(y_pred, y_test)))
27        # cpt = cpt+1
28
29    return np.average(accuracy)
30
31 # Converts floats to integers for classification
32 def my_ceil(predictions):
33     for i in range(len(predictions)):
34
35         if predictions[i]%1 <= +0.5:
36             predictions[i] = int(predictions[i])
37         else:
38             predictions[i] = int(predictions[i]) + 1
39    return predictions
40
41
42 def cross_validation_Lregressor(formula):
43     accuracy = []
44     for i in range(10):
45         # Split the data into training and testing sets
46         # Create a Linear Regression model
47         # Fit the model to the training data
48         # Predict the scores for the testing data
49         # Calculate the accuracy
50         # Print the accuracy
51         # Increment the counter
52     return accuracy
```

```
44     for train_index, test_index in Kf.split(well_B):
45         # print("TRAIN:", train_index, "TEST:", test_index)
46         #train part
47         train = well_B.iloc[train_index]
48         #test part
49         test = well_B.iloc[test_index]
50
51         # Define the scaler
52         # scaler = StandardScaler().fit(X_train)
53
54         # X_train = scaler.transform(X_train)
55         # X_test = scaler.transform(X_test)
56
57         model_lr = formula
58         result_lr = smf.ols(model_lr, data = train).fit()
59
60         y_pred = np.array(result_lr.predict(test))
61
62         y_pred = my_ceil(y_pred)
63
64         accuracy.append(accuracy_score(y_pred, test.Score))
65         # print("The score of the " + str(cpt) + " is " + str(accuracy_score(y_pred, y_test)))
66         # cpt = cpt+1
67
68     return np.average(accuracy)
69
```


▼ Regression

▼ Linear regression

```
1 import statsmodels.formula.api as smf
2 import statsmodels.api as sm
```

▼ First model

```
1 model_lr = 'Score ~ Temperature + Humidity + Humex + CO2 + Bright'
2 result_lr = smf.ols(model_lr, data = well_B).fit()
3 print(result_lr.summary())
```



OLS Regression Results

=====

Dep. Variable:	Score	R-squared:	0.380
Model:	OLS	Adj. R-squared:	0.380
Method:	Least Squares	F-statistic:	979.6
Date:	Thu, 19 Dec 2019	Prob (F-statistic):	0.00
Time:	18:09:19	Log-Likelihood:	-8541.0
No. Observations:	8000	AIC:	1.709e+04
Df Residuals:	7994	BIC:	1.714e+04
Df Model:	5		
Covariance Type:	nonrobust		

=====

	coef	std err	t	P> t	[0.025	0.975]
-----	-----	-----	-----	-----	-----	-----
Intercept	-18.7846	0.790	-23.779	0.000	-20.333	-17.236
Temperature	2.0279	0.080	25.209	0.000	1.870	2.186
Humidity	0.2986	0.010	30.564	0.000	0.279	0.318
Humex	-1.5151	0.061	-24.712	0.000	-1.635	-1.395
CO2	-0.0010	4.75e-05	-20.754	0.000	-0.001	-0.001
Bright	0.0004	0.000	3.234	0.001	0.000	0.001

=====

Omnibus:	10.461	Durbin-Watson:	1.163
Prob(Omnibus):	0.005	Jarque-Bera (JB):	10.717
Skew:	-0.070	Prob(JB):	0.00471
Kurtosis:	3.112	Cond. No.	6.32e+04


=====

Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 6.32e+04. This might indicate that there are strong multicollinearity or other numerical problems.

```
1 cross_validation_Lregressor(model_lr)
```



0.5186297859877849

▼ Second model

```
1 model_lr_1 = 'Score ~ Temperature + Humidity + Humex + CO2 + Bright + lundi + mardi + mercredi + jeudi + vendredi + samedi + dimanche'
2 result_lr_1 = smf.ols(model_lr_1, data = well_B).fit()
3 print(result_lr_1.summary())
```



OLS Regression Results						
=====						
Dep. Variable:	Score	R-squared:	0.386			
Model:	OLS	Adj. R-squared:	0.385			
Method:	Least Squares	F-statistic:	456.1			
Date:	Thu, 19 Dec 2019	Prob (F-statistic):	0.00			
Time:	18:09:27	Log-likelihood:	-8503.0			
No. Observations:	8000	AIC:	1.703e+04			
Df Residuals:	7988	BIC:	1.711e+04			
Df Model:	11					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]

Intercept	-16.2137	0.691	-23.456	0.000	-17.569	-14.859
Temperature	2.0151	0.080	25.098	0.000	1.858	2.172
Humidity	0.2984	0.010	30.612	0.000	0.279	0.317
Humex	-1.5116	0.061	-24.714	0.000	-1.631	-1.392
CO2	-0.0010	5.11e-05	-20.223	0.000	-0.001	-0.001
Bright	0.0003	0.000	2.857	0.004	0.000	0.001
lundi	-2.3511	0.101	-23.205	0.000	-2.550	-2.152
mardi	-2.3638	0.102	-23.267	0.000	-2.563	-2.165
mercredi	-2.2852	0.101	-22.581	0.000	-2.484	-2.087
jeudi	-2.2261	0.101	-22.046	0.000	-2.424	-2.028
vendredi	-2.2103	0.101	-21.947	0.000	-2.408	-2.013
samedi	-2.3632	0.100	-23.712	0.000	-2.559	-2.168
dimanche	-2.4141	0.099	-24.280	0.000	-2.609	-2.219
=====						
Omnibus:	11.224	Durbin-Watson:	1.174			
Prob(Omnibus):	0.004	Jarque-Bera (JB):	11.583			
Skew:	-0.071	Prob(JB):	0.00305			
Kurtosis:	3.121	Cond. No.	1.25e+18			
=====						

Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The smallest eigenvalue is 1.99e-27. This might indicate that there are strong multicollinearity problems or that the design matrix is singular.

1 cross_validation_Lregressor(model_lr_1)

 0.5176292543116874

▼ Third model

```
1 allcolumns = '+'.join(well_B.columns[:len(well_B.columns)-1])
2 allcolumns
3 model_lr_2 = 'Score ~'+ allcolumns
4 model_lr_2
5 result_lr_2 = smf.ols(model_lr_2, data = well_B).fit()
6 print(result_lr_2.summary())
```

OLS Regression Results						
=====						
Dep. Variable:	Score	R-squared:	0.393			
Model:	OLS	Adj. R-squared:	0.391			
Method:	Least Squares	F-statistic:	151.9			
Date:	Thu, 19 Dec 2019	Prob (F-statistic):	0.00			
Time:	18:22:00	Log-Likelihood:	-8452.9			
No. Observations:	8000	AIC:	1.698e+04			
Df Residuals:	7965	BIC:	1.722e+04			
Df Model:	34					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]

Intercept	-15.5752	0.665	-23.410	0.000	-16.879	-14.271
Temperature	2.0095	0.080	25.025	0.000	1.852	2.167
Humidity	0.2977	0.010	30.458	0.000	0.279	0.317
Humex	-1.5074	0.061	-24.603	0.000	-1.628	-1.387
CO2	-0.0011	5.94e-05	-18.148	0.000	-0.001	-0.001
Bright	0.0004	0.000	2.615	0.009	9.34e-05	0.001
lundi	-2.2590	0.098	-23.131	0.000	-2.450	-2.068
mardi	-2.2672	0.098	-23.161	0.000	-2.459	-2.075
mercredi	-2.1918	0.098	-22.478	0.000	-2.383	-2.001
jeudi	-2.1334	0.097	-21.931	0.000	-2.324	-1.943
vendredi	-2.1180	0.097	-21.827	0.000	-2.308	-1.928
samedi	-2.2777	0.096	-23.638	0.000	-2.467	-2.089
dimanche	-2.3282	0.096	-24.211	0.000	-2.517	-2.140
midnight	-0.6755	0.047	-14.256	0.000	-0.768	-0.583
AM1	-0.6304	0.048	-13.258	0.000	-0.724	-0.537
AM2	-0.6536	0.048	-13.746	0.000	-0.747	-0.560
AM3	-0.7399	0.048	-15.499	0.000	-0.834	-0.646
AM4	-0.7341	0.048	-15.395	0.000	-0.828	-0.641
AM5	-0.7869	0.048	-16.510	0.000	-0.880	-0.693
AM6	-0.6627	0.048	-13.885	0.000	-0.756	-0.569
AM7	-0.7093	0.048	-14.903	0.000	-0.803	-0.616
AM8	-0.5489	0.047	-11.606	0.000	-0.642	-0.456
AM9	-0.5267	0.048	-11.027	0.000	-0.620	-0.433
AM10	-0.5867	0.049	-11.989	0.000	-0.683	-0.491
AM11	-0.6063	0.050	-12.225	0.000	-0.704	-0.509
midday	-0.7021	0.048	-14.506	0.000	-0.797	-0.607
PM1	-0.7001	0.048	-14.522	0.000	-0.795	-0.606
PM2	-0.6717	0.048	-14.028	0.000	-0.766	-0.578
PM3	-0.7147	0.048	-14.998	0.000	-0.808	-0.621
PM4	-0.7004	0.047	-14.880	0.000	-0.793	-0.608
PM5	-0.7102	0.047	-15.104	0.000	-0.802	-0.618
PM6	-0.5424	0.047	-11.616	0.000	-0.634	-0.451
PM7	-0.4362	0.046	-9.392	0.000	-0.527	-0.345
PM8	-0.5937	0.047	-12.663	0.000	-0.686	-0.502
PM9	-0.6194	0.047	-13.197	0.000	-0.711	-0.527
PM10	-0.6652	0.047	-14.116	0.000	-0.758	-0.573
PM11	-0.6582	0.047	-13.884	0.000	-0.751	-0.565
=====						
Omnibus:	13.195	Durbin-Watson:	1.185			
Prob(Omnibus):	0.001	Jarque-Bera (JB):	13.656			
Skew:	-0.078	Prob(JB):	0.00108			
Kurtosis:	3.129	Cond. No.	2.99e+18			
=====						

Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The smallest eigenvalue is 3.49e-28. This might indicate that there are strong multicollinearity problems or that the design matrix is singular.

1 cross_validation_Lregressor(model_lr_2)

 0.5175061304727452

▼ XGboost

```
1 import xgboost

1 model_xg = xgboost.XGBRegressor(objective ='reg:squarederror', colsample_bytree = 0.3, learning_rate = 0.05,
2                               alpha = 10,base_score= 0.4, booster='gbtree', n_estimators = 150)

1 cross_validation_XG_classifier(model_xg)

0.7156269433045913
```

▼ Classification


```
1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from google.colab import files

1 def my_ceil(predictions):
2     for i in range(len(predictions)):
3         if predictions[i]%1<=0.5:
4             predictions[i] = int(predictions[i])
5         else:
6             predictions[i] = int(predictions[i]) + 1
7     return predictions
8
9 #Export function
10 def export( data_test, predictions):
11     result_ = pd.DataFrame({'ID': data_test.ID, 'Score': my_ceil(predictions)})
12     result_.to_csv('results_.csv', index = False)
13
14     files.download('results_.csv')
15
```

▼ Upload preprocessed dataset train and test

```
1 test = pd.read_csv('https://raw.githubusercontent.com/nisrineha/Challenge_data_well_being_at_work/master/datasets/data_test_processed.csv')
2 train = pd.read_csv('https://raw.githubusercontent.com/nisrineha/Challenge_data_well_being_at_work/master/datasets/data_train_processed.csv')
3 test_with_ID = pd.read_csv('https://raw.githubusercontent.com/nisrineha/Challenge_data_well_being_at_work/master/datasets/test_input.csv')
4 test_with_date= pd.read_csv('https://raw.githubusercontent.com/nisrineha/Challenge_data_well_being_at_work/master/datasets/train_input.csv')
```

```
1 test_with_date.head()
2
```



	ID	Date	Temperature	Humidity	Humex	C02	Bright
0	0	2017-08-31 23:30:00	22.7	56.0	25.7	534.0	1.0
1	1	2017-09-01 00:30:00	22.8	55.0	25.7	506.0	1.0
2	2	2017-09-01 01:30:00	22.9	55.0	25.9	577.0	1.0
3	3	2017-09-01 02:30:00	23.0	55.0	26.1	630.0	1.0
4	4	2017-09-01 03:30:00	23.0	55.0	26.1	643.0	1.0

Creation of features and the target variable Score

```
1 y= train['Score']
2 train1= train
3 train1= train1.drop('Score', axis= 1)
4 X= train1
5 train1.head()
```

Spliting of the data

```
1 X_train, X_test, y_train, y_test= train_test_split(X, y , test_size= 0.2, random_state= 42)
```

▼ Pipeline method

Implement of pipeline method using different transformer: numeric and categorial

```
1 from sklearn.pipeline import Pipeline
2 from sklearn.impute import SimpleImputer
3 from sklearn.preprocessing import StandardScaler, OneHotEncoder
4
5 #SimpleImputer fill any missing values
6 #Scaler numeric transformer
7
8 numeric_transformer = Pipeline(steps=[
9     ('imputer', SimpleImputer(strategy='median')),
10    ('scaler', StandardScaler())])
11
12 #One hot encoder to transform categorial values into integers.
13
14 categorical_transformer = Pipeline(steps=[
15     ('imputer', SimpleImputer(strategy='constant', fill_value='missing')),
16     ('onehot', OneHotEncoder(handle_unknown='ignore'))])
```

Transform the categorical features and numeric on train dataset and test

```
1 #Select les columns numeric
2 #Select les columns categoric
3
4
5 integer_features = list(X.columns[X.dtypes == 'int64'])
6 continuous_features = list(X.columns[X.dtypes == 'float64'])
7 categorical_features = list(X.columns[X.dtypes == 'object'])
8 numeric_features = integer_features + continuous_features
9
10
```

```
11 from sklearn.compose import ColumnTransformer
12 preprocessor = ColumnTransformer(
13     transformers=[
14         ('num', numeric_transformer, numeric_features),
15         ('cat', categorical_transformer, categorical_features)])
16
17 integer_features_test = list(test.columns[test.dtypes == 'int64'])
18 continuous_features_test = list(test.columns[test.dtypes == 'float64'])
19 categorical_features_test = list(test.columns[test.dtypes == 'object'])
20 numeric_features = integer_features + continuous_features
21
22
23 from sklearn.compose import ColumnTransformer
24 preprocessor = ColumnTransformer(
25     transformers=[
26         ('num', numeric_transformer, numeric_features),
27         ('cat', categorical_transformer, categorical_features)])
```

▼ Model selection

In this section, we chose different classifier from sklearn, to get the best classifier for our dataset, we calculate the score thanks to the splitting of the dataset that we did before.

```
1 from sklearn.metrics import accuracy_score, log_loss
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.svm import SVC, LinearSVC, NuSVC
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier, GradientBoostingClassifier
6 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
7 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis
8 classifiers = [
9     KNeighborsClassifier(3),
10    SVC(kernel="rbf", C=0.025, probability=True),
11    DecisionTreeClassifier(),
12    RandomForestClassifier(),
13    AdaBoostClassifier(),
14    GradientBoostingClassifier(),
15    LinearDiscriminantAnalysis(),
16    QuadraticDiscriminantAnalysis()
17 ]
18 pipes= []
19 for classifier in classifiers:
20     pipe = Pipeline(steps=[('preprocessor', preprocessor),
21                            ('classifier', classifier)])
22     pipe.fit(X_train, y_train)
23     pipes.append(pipe)
24     print(classifier)
25     print("model score: %.3f" % pipe.score(X_test, y_test))
26
```

```
➤ KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                      metric_params=None, n_jobs=None, n_neighbors=3, p=2,
                      weights='uniform')

model score: 0.587
SVC(C=0.025, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
    kernel='rbf', max_iter=-1, probability=True, random_state=None,
    shrinking=True, tol=0.001, verbose=False)
model score: 0.439
DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=None,
                      max_features=None, max_leaf_nodes=None,
                      min_impurity_decrease=0.0, min_impurity_split=None,
                      min_samples_leaf=1, min_samples_split=2,
                      min_weight_fraction_leaf=0.0, presort=False,
                      random_state=None, splitter='best')

model score: 0.601
RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                      max_depth=None, max_features='auto', max_leaf_nodes=None,
                      min_impurity_decrease=0.0, min_impurity_split=None,
                      min_samples_leaf=1, min_samples_split=2,
                      min_weight_fraction_leaf=0.0, n_estimators=10,
                      n_jobs=None, oob_score=False, random_state=None,
                      verbose=0, warm_start=False)

model score: 0.677
/usr/local/lib/python3.6/dist-packages/sklearn/ensemble/forest.py:245: FutureWarning: The default value of n_estimators will change from 10 in version 0.20 to 100 in 0.22.
  "10 in version 0.20 to 100 in 0.22.", FutureWarning)
AdaBoostClassifier(algorithm='SAMME.R', base_estimator=None, learning_rate=1.0,
                  n_estimators=50, random_state=None)
model score: 0.565
GradientBoostingClassifier(criterion='friedman_mse', init=None,
                          learning_rate=0.1, loss='deviance', max_depth=3,
                          max_features=None, max_leaf_nodes=None,
                          min_impurity_decrease=0.0, min_impurity_split=None,
                          min_samples_leaf=1, min_samples_split=2,
                          min_weight_fraction_leaf=0.0, n_estimators=100,
                          n_iter_no_change=None, presort='auto',
                          random_state=None, subsample=1.0, tol=0.0001,
                          validation_fraction=0.1, verbose=0,
                          warm_start=False)

model score: 0.750
LinearDiscriminantAnalysis(n_components=None, priors=None, shrinkage=None,
                          solver='svd', store_covariance=False, tol=0.0001)

model score: 0.562
QuadraticDiscriminantAnalysis(priors=None, reg_param=0.0,
                              store_covariance=False, tol=0.0001)

model score: 0.138
/usr/local/lib/python3.6/dist-packages/sklearn/discriminant_analysis.py:388: UserWarning: Variables are collinear.
  warnings.warn("Variables are collinear.")
/usr/local/lib/python3.6/dist-packages/sklearn/discriminant_analysis.py:693: UserWarning: Variables are collinear
  warnings.warn("Variables are collinear")
```

From the results above, we observe that the best classifiers are **GradientBoostingClassifier** and **RandomForestClassifier**

we chose GradientBoostingClassifier for submission, we had our best score which is **0,6990**

```
1 x = test
2 y_pred= pipes[-3].predict(x)
3 export(test_with_ID, y_pred)
```

▼ Using pipeline in GridSearch

```
1 param_grid = {
2     'classifier__n_estimators': [ 200, 300, 400, 500]
```

```
2     'classifier__n_estimators': [100, 200, 300, 400],
3     'classifier__max_features': ['auto', 'sqrt', 'log2'],
4     'classifier__max_depth' : [10, 20, 25, 30],
5     'classifier__criterion' :['gini', 'entropy']}
6 from sklearn.model_selection import GridSearchCV
7 CV = GridSearchCV(rf, param_grid, n_jobs= 1)
8
9 CV.fit(X_train, y_train)
10 print(CV.best_params_)
11 print(CV.best_score_)

🔗 /usr/local/lib/python3.6/dist-packages/sklearn/model_selection/_split.py:1978: FutureWarning: The default value of cv will change from 3 to 5 in version 0.22. Specify it explicitly to silence warnings.warn(CV_WARNING, FutureWarning)
{'classifier__criterion': 'entropy', 'classifier__max_depth': 10, 'classifier__max_features': 'auto', 'classifier__n_estimators': 400}
0.71359375
```

we had from our grid search our best max depth and number of estimators for our randomforset model.

```
1 #Fitting the classifier
2 from sklearn.ensemble import RandomForestClassifier
3 rf = Pipeline(steps=[('preprocessor', preprocessor),
4                       ('classifier', RandomForestClassifier(n_estimators= 400, max_depth=20))])

1 rf.fit(X_train, y_train )
2
3 pipe = Pipeline(steps=[('preprocessor', preprocessor),
4                       ('classifier',RandomForestClassifier(n_estimators= 400, max_depth=20) )])
5 pipe.fit(X_train, y_train)
6 print(classifier)
7 print("model score: %.3f" % pipe.score(X_test, y_test))
8 x = test
9 y_pred= pipe.predict(x)
10 export(test_with_ID, y_pred)

🔗 QuadraticDiscriminantAnalysis(priors=None, reg_param=0.0,
                                store_covariance=False, tol=0.0001)
model score: 0.734
```

We submitted our results, but we had a score less than the one when we implemented **GradientBoostingClassifier**

▼ Random forest model

```
1 rf_time = RandomForestClassifier(n_estimators = 400, random_state = 0, max_depth=20)

1 cross_validation_XG_classifier(rf_time)

👤 0.706623909314313
```

▼ K nearest neighbor

```
1 neigh_time = KNeighborsClassifier(n_neighbors=5)

1 cross_validation_XG_classifier(neigh_time)

👤 0.48610790989129854
```

▼ Gradient Boosting Classification

```
1 from sklearn.ensemble import GradientBoostingClassifier

1 gbC_time = GradientBoostingClassifier()

1 cross_validation_XG_classifier(gbC_time)

👤 0.6976108065300944
```

▼ Basic deep learning model


```
1 # Import `Sequential` from `keras.models`
2 from keras.models import Sequential
3
4 # Import `Dense` from `keras.layers`
5 from keras.layers import Dense
6
7 # Initialize the constructor
8 model = Sequential()
9
10 # Add an input layer
11 model.add(Dense(12, activation='softmax', input_shape=(36,)))
12
13 # Add one hidden layer
14 model.add(Dense(12, activation='relu'))
15
16 # # Add one hidden layer
17 # model.add(Dense(12, activation='relu'))
18
19 # Add an output layer
20 model.add(Dense(output_dim = 5, activation = 'softmax'))

👤 /usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:19: UserWarning: Update your `Dense` call to the Keras 2 API: `Dense(activation="relu", units=5)`

1 def dummies_categ(y):
2     cat = []
3     for i in range(len(y)):
4         ind = np.argmax(y_pred[i]) + 1
5         cat.append(ind)
6     return cat

1 accuracy = []
2 for train_index, test_index in kf.split(well_B):
3     # print("TRAIN:", train_index, "TEST:", test_index)
4     #train part
5     X_train = np.array(well_B.drop('Score', axis = 1))[train_index]
```

```
6     y_train = np.array(well_B.Score)[train_index]
7
8     y_train = pd.get_dummies(y_train)
9     #test part
10    X_test = np.array(well_B.drop('Score', axis = 1))[test_index]
11    y_test = np.array(well_B.Score)[test_index]
12
13    y_test = pd.get_dummies(y_test)
14
15
16    # Define the scaler
17    scaler = StandardScaler().fit(X_train)
18
19    X_train = scaler.transform(X_train)
20    X_test = scaler.transform(X_test)
21
22    model.compile(loss='binary_crossentropy',
23                  optimizer='adam',
24                  metrics=['accuracy'])
25
26    model.fit(X_train, y_train,epochs=2, batch_size=1, verbose=1)
27
28    y_pred = model.predict(X_test)
29
30    y_pred = dummies_categ(y_pred)
31    y_test = np.array(well_B.Score)[test_index]
32    # print(y_pred)
33    accuracy.append(accuracy_score(y_pred, y_test))
34    # print("The score of the " + str(cpt) + " is " + str(accuracy_score(y_pred, y_test)))
35    # cpt = cpt+1
36
37 np.average(accuracy)
```



Epoch 1/2
6857/6857 [=====] - 12s 2ms/step - loss: 0.5317 - acc: 0.8609
Epoch 2/2
6857/6857 [=====] - 11s 2ms/step - loss: 0.5344 - acc: 0.8626
Epoch 1/2
6857/6857 [=====] - 12s 2ms/step - loss: 0.5495 - acc: 0.8596
Epoch 2/2
6857/6857 [=====] - 10s 2ms/step - loss: 0.5380 - acc: 0.8617
Epoch 1/2
6857/6857 [=====] - 12s 2ms/step - loss: 0.5503 - acc: 0.8602
Epoch 2/2
6857/6857 [=====] - 11s 2ms/step - loss: 0.5311 - acc: 0.8629
Epoch 1/2
6857/6857 [=====] - 12s 2ms/step - loss: 0.5359 - acc: 0.8615
Epoch 2/2
6857/6857 [=====] - 11s 2ms/step - loss: 0.5310 - acc: 0.8646
Epoch 1/2
6857/6857 [=====] - 13s 2ms/step - loss: 0.5457 - acc: 0.8645
Epoch 2/2
6857/6857 [=====] - 11s 2ms/step - loss: 0.5314 - acc: 0.8665
Epoch 1/2
6857/6857 [=====] - 13s 2ms/step - loss: 0.5260 - acc: 0.8682
Epoch 2/2
6857/6857 [=====] - 11s 2ms/step - loss: 0.5188 - acc: 0.8701
Epoch 1/2
6858/6858 [=====] - 13s 2ms/step - loss: 0.5262 - acc: 0.8665
Epoch 2/2
6858/6858 [=====] - 11s 2ms/step - loss: 0.5191 - acc: 0.8684
0.6390032025331335