## ▾ 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
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_wel
3   train_output= pd.read_csv('https://raw.githubusercontent.com/nisrineha/Challenge_data_we
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
```

↱

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

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

↱

## ▾ 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 a

## 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
```

⤷

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 reprensents the score in function of day hours, the confort subjectif dro

maybe people are working at this hours and they are less confortable.

```
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 mc

# ▾ 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

## 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)
```

## 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()
```

```
1    dummy_weekday = pd.get_dummies(well_B['weekday'])
2    dummy_weekday.columns = ['lundi','mardi','mercredi','jeudi','vendredi','samedi','dimanch
3
4    # dummy_weekday.rename(columns = {'0':'lundi','1':'mardi','2':'mercredi','3':'jeudi','4'
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','
 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()
```

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

## ▾ Test Data

### ▸ Import test data

↳ *1 cellule masquée*

### ▸ Add weekdays to columns as dummies

↳ *3 cellules masquées*

### ▸ Add hours to columns as dummy variables

↳ *1 cellule masquée*

## ▾ 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=
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 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

> ↳ 10 cellules masquées

## XGboost

> ↳ 3 cellules masquées

# 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_bein
2    train = pd.read_csv('https://raw.githubusercontent.com/nisrineha/Challenge_data_well_bei
3    test_with_ID = pd.read_csv('https://raw.githubusercontent.com/nisrineha/Challenge_data_w
4    test_with_date= pd.read_csv('https://raw.githubusercontent.com/nisrineha/Challenge_data_
```

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

⤷

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=
```

# ▾ 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
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, GradientBoostin
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
```

⤷

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

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],
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_)
```

⤷

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

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

```
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)
```

⤷

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

▸ Random forest model

> ↳ *2 cellules masquées*

▸ K nearest neighbor

↳ 2 cellules masquées

▶ Gradient Boosting Classification

↳ 3 cellules masquées

## ▼ 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'))
```

```
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)
```