

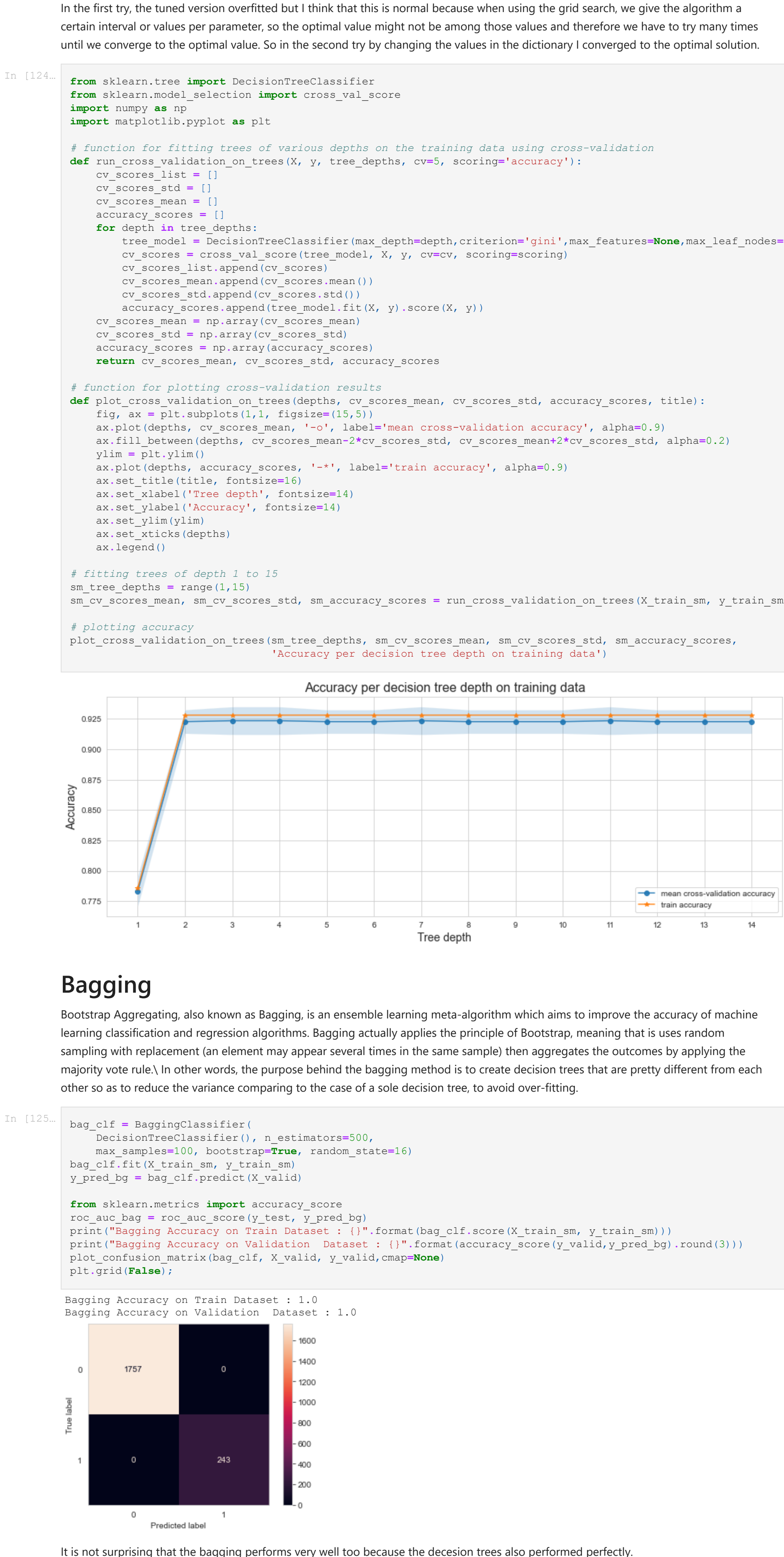

```
Fitting 5 folds for each of 864 candidates, totalling 4320 fits
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 4320 out of 4320 | elapsed: 1.0min finished
GridSearchCV(cv=5, estimator=DecisionTreeClassifier(random_state=16),
             param_grid={'criterion': ['gini'], 'max_depth': [1, 2, 3],
                           'max_features': ['auto', 'log2', 'sqrt', 'none'],
                           'max_leaf_nodes': [1, 3, 5, 10],
                           'min_samples_leaf': [1, 5, 10],
                           'min_weight_fraction_leaf': [0.01, 0.1, 0.4]},
             scorer='accuracy', verbose=1)
```

```
tuning_model.best_params_
{'criterion': 'gini',
 'max_depth': 1,
 'max_features': None,
 'max_leaf_nodes': 5,
 'min_samples_leaf': 1,
 'min_weight_fraction_leaf': 0.01,
 'splitter': 'best',
 'verbose': 1}
```

```
cart_tuned = DecisionTreeClassifier(criterion='gini', max_depth=3, max_features=None, max_leaf_nodes=5, min_samples_leaf=1)
y_cart_tuned = cart_tuned.predict(X_valid)
print("Tuned CART Accuracy on Train Dataset : {}".format(cart_tuned.score(X_train_sm, y_train_sm).round(3)))
print("Tuned CART Accuracy on Validation Dataset : {}".format(cart_tuned.score(X_valid, y_cart_tuned).round(3)))
```

Tuned CART Accuracy on Train Dataset : 1.0
Tuned CART Accuracy on Validation Dataset : 1.0

In the first try, the tuned version overfitted but I think that this is normal because when using the grid search, we give the algorithm a certain interval or values per parameter, so the optimal value might not be among those values and therefore we have to try many times until we converge to the parameter value. So in the second try by changing the values in the dictionary I converged to the optimal solution.



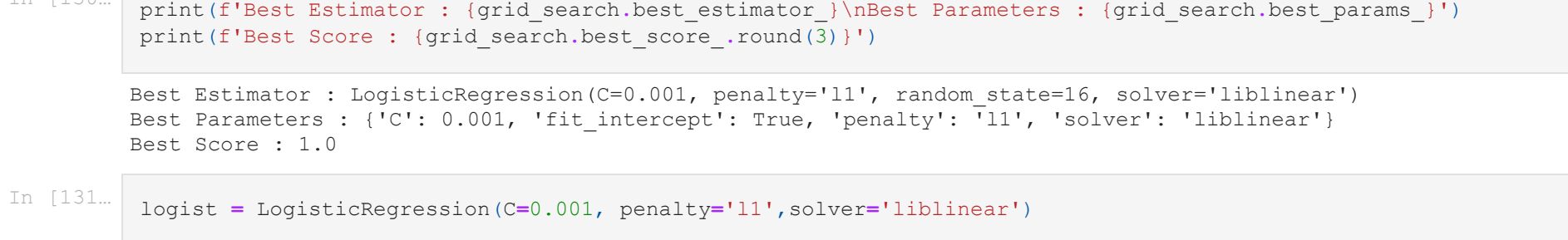
Bagging

Bootstrap Aggregating, also known as Bagging, is an ensemble learning meta-algorithm which aims to improve the accuracy of machine learning classification and regression algorithms. Bagging actually applies the principle of Bootstrap, meaning that it uses random sampling with replacement (an element may appear several times in the same sample) then aggregates the outcomes by applying the majority vote rule. In other words, the purpose behind the bagging method is to create decision trees that are pretty different from each other so as to reduce the variance, comparing to the case of a single decision tree, to avoid over-fitting.

```
bag_clf = BaggingClassifier(
    DecisionTreeClassifier(n_estimators=500,
                          max_samples=100, bootstrap=True, random_state=16),
    bag_clf_fit(X_train_sm, y_train_sm),
    y_pred_bg = bag_clf.predict(X_valid)

from sklearn.metrics import accuracy_score
roc_auc_bag = roc_auc_score(y_test, y_pred_bg)
print("Bagging Accuracy on Train Dataset : {}".format(bag_clf.score(X_train_sm, y_train_sm)))
print("Bagging Accuracy on Validation Dataset : {}".format(accuracy_score(y_valid, y_pred_bg).round(3)))
plot_confusion_matrix(bag_clf, X_valid, y_valid, cmap=None)
plt.grid(False)
```

Bagging Accuracy on Train Dataset : 1.0
Bagging Accuracy on Validation Dataset : 1.0



It is not surprising that the bagging performs very well too because the decision trees also performed perfectly.

Hyperparameter Tuning

Logistic Regression (GridSearchCV)

We'll be optimizing 4 different hyperparameters of our logistic regression model :

- **penalty**
- **C**
- **fit_intercept**
- **solver**

```
# Feature scaling
from sklearn.preprocessing import StandardScaler
sc = StandardScaler().fit(X_train_sm)
X_train_sc = sc.transform(X_train_sm)
X_valid_sc = sc.transform(X_valid) # same training scale for the validation
```

```
log_model = LogisticRegression(max_iter=1000,
                               param_grid={'C': (0.001, 0.02, 0.01, 0.03, 0.1),
                                             'penalty': ('l1', 'l2'),
                                             'fit_intercept': (True, False),
                                             'solver': ('liblinear', 'saga')})
```

```
grid_search = GridSearchCV(estimator=log_model, param_grid=param_grid, scoring='accuracy', n_jobs=-1, cv=5, verbose=1)
grid_search.fit(X_train_sc, y_train_sc) # scaled
```

Fitting 5 folds for each of 40 candidates, totalling 600 fits

| [Parallel(n_jobs=-1)]: | Using backend | Using backend | Using backend |
|------------------------|---------------------|---------------|---------------|
| [Parallel(n_jobs=-1)]: | Done 34 tasks | elapsed: 3.6s | |
| [Parallel(n_jobs=-1)]: | Done 346 tasks | elapsed: 6.0s | |
| [Parallel(n_jobs=-1)]: | Done 600 out of 600 | elapsed: 9.8s | finished |

```
GridSearchCV(cv=5, estimator=LogisticRegression(random_state=16), n_jobs=-1,
             param_grid={'C': (0.001, 0.02, 0.01, 0.03, 0.1),
                           'fit_intercept': (True, False),
                           'penalty': ('l1', 'l2'),
                           'solver': ('liblinear', 'saga')},
             scoring='accuracy', verbose=1)
```

```
print(f"Best Estimator : {grid_search.best_estimator_}\nBest Parameters : {grid_search.best_params_}")
print(f"Best Score : {grid_search.best_score_.round(3)}")
```

Best Estimator : LogisticRegression(C=0.001, penalty='l1', solver='liblinear')
Best Parameters : {'C': 0.001, 'fit_intercept': True, 'penalty': 'l1', 'solver': 'liblinear'}
Best Score : 1.0

```
logist = LogisticRegression(C=0.001, penalty='l1', solver='liblinear')
```

```
# Confusion matrix
from sklearn.metrics import confusion_matrix, accuracy_score
logist_fit(X_train_sc, y_train_sc)
y_pred_logist = logist.predict(X_valid_sc)
print(f"Confusion matrix(y_valid, y_pred_logist)\nValidation Accuracy : {accuracy_score(y_valid, y_pred_logist).round(3)}")
```

```
[[1757  0]
 [  0 243]]
Validation Accuracy : 1.0
```

```
print("Tuned Logistic Regression Accuracy on Train Dataset : {}".format(logist.score(X_train_sc, y_train_sc)))
print("Tuned Logistic Regression Accuracy on Validation Dataset : {}".format(accuracy_score(y_valid, y_pred_logist).round(3)))
# Visualizing on the test set or valid conditionner des infos de la train, il faut analyser le train etc
```

Tuned Logistic Regression Accuracy on Train Dataset : 1.0
Tuned Logistic Regression Accuracy on Validation Dataset : 1.0

Tuning Random Forest

```
RandomSearchCV
random_model = RandomForestClassifier(random_state=16)

n_estimators = [int(x) for x in np.linspace(start = 1, stop = 20, num = 20)] # number of trees in the random forest
max_features = ['auto', 'sqrt'] # number of features in consideration at every split
min_depth = [int(x) for x in np.linspace(10, 100, num = 12)] # maximum number of levels allowed in each decision tree
min_samples_split = [2, 5, 10] # minimum sample number to split a node
min_samples_leaf = [1, 2, 5, 10] # minimum sample number that can be stored in a leaf node
bootstrap = [True, False] # whether to use bootstrap samples

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

```
rfo_random = RandomizedSearchCV(estimator = random_model, param_distributions = random_grid,
                               n_iter = 100, cv = 30, verbose=1, random_state=35, n_jobs = -1)
rfo_random.fit(X_train_sm, y_train_sm)
print(f"Best Estimator : {rfo_random.best_estimator_}")
```

Fitting 50 folds for each of 100 candidates, totalling 3000 fits

| [Parallel(n_jobs=-1)]: | Using backend | Using backend | Using backend |
|------------------------|-----------------------|-----------------|---------------|
| [Parallel(n_jobs=-1)]: | Done 34 tasks | elapsed: 2.1s | |
| [Parallel(n_jobs=-1)]: | Done 184 tasks | elapsed: 6.0s | |
| [Parallel(n_jobs=-1)]: | Done 522 tasks | elapsed: 22.1s | |
| [Parallel(n_jobs=-1)]: | Done 960 tasks | elapsed: 35.7s | |
| [Parallel(n_jobs=-1)]: | Done 1860 tasks | elapsed: 1.1min | |
| [Parallel(n_jobs=-1)]: | Done 2800 tasks | elapsed: 1.7min | |
| [Parallel(n_jobs=-1)]: | Done 3000 out of 3000 | elapsed: 1.8min | finished |

```
Best Estimator : RandomForestClassifier(bootstrap=False, max_depth=10, min_samples_leaf=2, min_samples_split=12, n_estimators=14, n_jobs=-1,
                                       random_state=16)
```

```
# Confusion matrix
from sklearn.metrics import confusion_matrix, accuracy_score
rfo_opt_fit(X_train_sm, y_train_sm)
y_pred_rfo = rfo_opt.predict(X_valid)
print(f"Validation Accuracy : {accuracy_score(y_valid, y_pred_rfo).round(3)}")
plot_confusion_matrix(rfo_opt, X_valid, y_valid, cmap=None)
plt.grid(False)
```

Validation Accuracy Score: 1.0



```
print(f"Train Accuracy : {rfo_opt.score(X_train_sm, y_train_sm).round(3)}")
rfprediction=rfo_opt.predict(X_valid)
print(f"Validation Accuracy : {accuracy_score(y_valid, rfprediction).round(3)}")
print(confusion_matrix(y_valid, rfprediction))
```

```
Train Accuracy : 1.0
Validation Accuracy : 1.0
[[1757  0]
 [  0 243]]
```

```
GridSearchCV
parameters = {
    'n_estimators': [1, 5, 10, 15],
    'criterion': ['gini', 'entropy'],
    'max_features': ['sqrt', 'log2'],
}
est = GridSearchCV(estimator=RandomForestClassifier(),
                   param_grid=parameters, cv=10)
clf = GridSearchCV(est, parameters, cv=10)
clf.fit(X_train_sm, y_train_sm)
```

```
GridSearchCV(cv=10, estimator=RandomForestClassifier(),
             param_grid={'criterion': ['gini', 'entropy'],
                           'max_features': ['sqrt', 'log2'],
                           'n_estimators': [1, 5, 10, 15]})
```

To measure impurity (ie. heterogeneity), we seek, among all the possible splits, the one that best separates the classes, namely the one that minimizes impurity so as to get the most homogeneous nodes. Among the criteria used to measure this impurity are the Gini index and entropy.

```
clf.best_params_
{'criterion': 'entropy', 'max_features': 'sqrt', 'n_estimators': 15}
```

```
rfo_optg = RandomForestClassifier(criterion='gini', max_features='sqrt', n_estimators=15)
rfo_optg_fit(X_train_sm, y_train_sm)
y_pred_rfo = rfo_optg.predict(X_valid)
```

```
print(f"(confusion_matrix(y_valid, y_pred_rfo))\nValidation Accuracy : {accuracy_score(y_valid, y_pred_rfo).round(3)}")
[[1757  0]
 [  0 243]]
Validation Accuracy : 1.0
```

```
# Grid search graphs
fp, ax = plt.subplots(1, 2, figsize=(15, 5))
n_estimators = [1, 5, 10, 15, 20, 30, 40, 50, 60, 70, 80, 90, 100]
ax = grid_search.cv_results_['mean_test_score'].groupby(n_estimators).mean()
plt.grid(False)
```



I would like to remind that the purpose behind using tuning methods in this data set is just a matter of practice, since the initial results were very satisfying I did not have to tune the models.

Tuning K-NN (Optuna : Bayesian Optimization)

K-Nearest Neighbors is a supervised learning algorithm that makes predictions on qualitative or quantitative variables. This algorithm is made up of several steps. First, we must choose the number K of neighbors, the default value of K is 5. We must then take the K-nearest neighbors of the new data point based on the Euclidean distance, the Hamming distance, the Minkowski distance or the Manhattan distance or any other distance that can be considered. The next step is to count the number of data points in each category and then assign the new point to the category that contains the most neighbors, hence the name K-Nearest Neighbors. This classification method does not require the model to be fitted and is easy to implement. Being a non-parametric algorithm, K-NN has no assumptions to consider, unlike parametric models such as linear regression. It should also be noted that K-NN requires homogeneous features and it could be that the algorithm works slowly, is not suitable for imbalanced data and becomes problematic when there are missing values.

```
def optimize(trial):
    # Instantiate scaler
    scalers = trial.suggest_categorical("scalers", ['minmax', 'standard', 'robust'])

    # Make a pipeline
    pipeline = make_pipeline(scaler, knn)

    # Cross-validate the features reduced by dimensionality reduction methods
    kfold = StratifiedKFold(n_splits=10)
    score = cross_val_score(pipeline, X_train_sm, y_train_sm, scoring='accuracy', cv=kfold)
    score = score.mean()
    return score

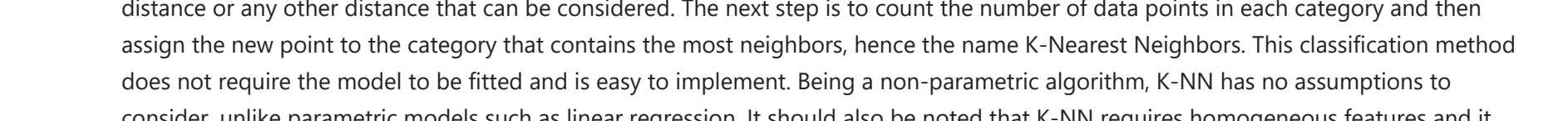
# Optuna
sampler = TPESampler(seed=42) # create a seed for the sampler for reproducibility
study = optuna.create_study(direction='maximize', sampler=sampler)
study.optimize(objective, n_trials=100)
```

```
print("Best trial out of 100 is:")
study.best_trial
```

Best trial out of 100 is :

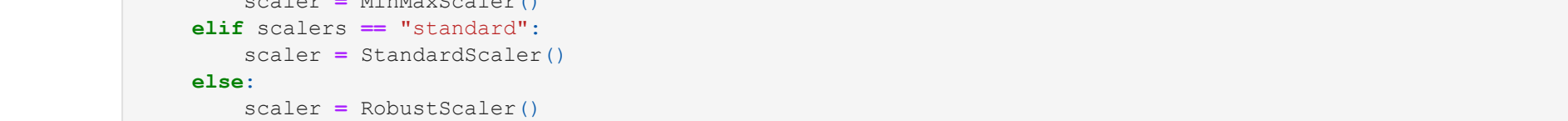
```
Frontier(number=10, values=[0], datetime_start=datetime.datetime(2022, 2, 3, 17, 23, 29, 194988), datetime_end=datetime.datetime(2022, 2, 3, 17, 23, 36, 453737), params={'scalers': 'minmax', 'n_neighbors': 11, 'weights': 'distance', 'metric': 'euclidean'}, distributions={'scalers': 'categorical', 'n_neighbors': 'int', 'weights': 'categorical', 'metric': 'categorical', 'distance': 'categorical'}, categorical_distribution_choices=('uniform', 'standard', 'distance'), metric=('euclidean', 'manhattan', 'euclidean'), manha
tran=('minmax', 'uniform'), user_attrs={}, system_attrs={}, intermediate_values={}, trial_id=10, state='TrialState.COMPLETE', value=None)
```

```
optuna.visualization.plot_optimization_history(study)
```



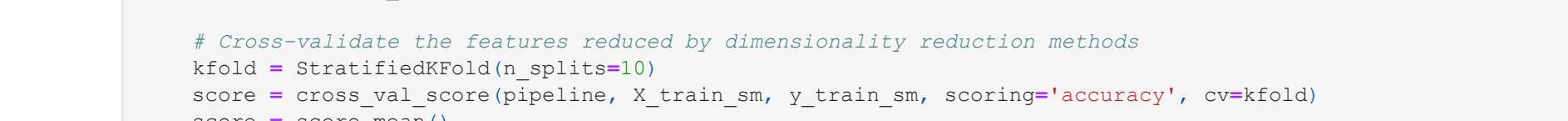
The objective value became almost stable after 10 trials.

```
# most important parameters
optuna.visualization.plot_param_importances(study)
```



We notice that the most important parameters here are the scalers (which we are going to apply in the next step) and the number of neighbors.

```
optuna.visualization.plot_parallel_coordinate(study)
```



The parallel coordinates plot allows to compare the feature of several individual observations on a set of numeric variables. Each vertical bar represents a parameter and has its own scale.

```
# Have a look at the best parameters for the tuned model
print("Best parameters after tuning using mean accuracy:")
study.best_params
```

Best parameters after tuning using mean accuracy:

```
{'scalers': 'minmax',
 'n_neighbors': 11,
 'weights': 'distance',
 'metric': 'euclidean'}
```

```
minmax = MinMaxScaler().fit(X_train_sm)
X_train_minmax = minmax.transform(X_train_sm)
X_valid_minmax = minmax.transform(X_valid)
X_test_minmax = minmax.transform(X_test)
```

```
knn = KNeighborsClassifier(n_neighbors=11, weights='distance', metric='euclidean')
knn.fit(X_train_minmax, y_train_sm)

print(f"Train Accuracy : {knn.score(X_train_minmax, y_train_sm).round(3)}")
prediction=knn.predict(X_valid_minmax)
print(f"Validation Accuracy : {accuracy_score(y_valid, prediction).round(3)}")
print(confusion_matrix(y_valid, prediction))
print(classification_report(y_valid, prediction))
```

```
Train Accuracy : 1.0
Validation Accuracy : 1.0
[[1757  0]
 [  0 243]]
precision    recall  f1-score   support

0     1.00     1.00     1.00     1757
1     1.00     1.00     1.00     243

accuracy avg 1.00 1.00 1.00 2000
macro avg 1.00 1.00 1.00 2000
weighted avg 1.00 1.00 1.00 2000
```

- **Precision** : Precision, also known as Positive Prediction Value (PPV), shows the proportion of true positive predictions, i.e. how many predictions were truly positive out of the total positive predicted values. Precision is good to use when there is an important proportion of FP, which should be minimized.

$$Precision = \frac{TP}{TP + FP}$$

- **Recall** : Also known as True Positive Rate (TPR) or Sensitivity, Recall is a classification error metric which shows how many predictions were correctly predicted out of the total positive actual values. This metric is preferred when there is a large proportion of FN, which should be minimized.

$$Recall = \frac{TP}{TP + FN}$$

- **F-Score** : F-Score or F-Measure combines both precision and recall into a unique measure that captures both properties. The general F-Beta formula is :

$$F_\beta = (1 + \beta^2) \frac{Precision * Recall}{(\beta^2 * Precision) + Recall}$$

Selecting the beta value depends on FP and FN. If FP and FN are equally important, then $\beta = 1$ and thus the formula becomes :

$$F_1 = 2 * \frac{Precision * Recall}{Precision + Recall}$$

which is basically the harmonic mean of precision and recall, in this case it is called F_1 -Score. If FP are more important than FN (i.e. more weight on precision), the beta will be reduced to a range between 1 and 1, a $\beta = 0.5$ for instance refers to an $F_{0.5}$ -Score. Finally, if FN have a higher impact than FP (i.e. more weight on recall), beta value will be higher than 1, a $\beta = 2$ for instance refers to an F_2 -Score and so on.

```
# KNN Permutation Feature Importance (on the validation set)
# n = permutation importance(knn, X_train_max, y_valid,
#                           n_repeats=30,
#                           random_state=0)

for i in range(1, n_repeats+1):
    if r.importances_mean[i] - 2 * r.importances_std[i] > 0:
        print(f"X_train_max column {i} is :")
        f = r.importances_mean[i].3f
        f += - r.importances_std[i].3f"
```

```
isNewBuilt 0.127 +/- 0.006
Pool      0.123 +/- 0.005
numPrevOw 0.118 +/- 0.005
numPrevOw 0.002 +/- 0.001
attico    0.001 +/- 0.001
```

The permutation feature importance is the decrease in a model score when a single feature value is randomly shuffled. After training the model we measure the importance of the features. For instance, if we look at the isNewBuilt column, we permute this column randomly and we measure the decrease in the metric (accuracy for instance) compared to the non-permuted case. This decrease can be interpreted as the column importance. In this example I chose 30 iterations and the decrease in performance is calculated in each iteration. For the isNewBuilt column for instance, the 0.127 represents the mean by which the performance of the model decreased on the 30 iterations and 0.006 is the standard deviation.

```
# Testing the bagging method on the test set
y_pred_bag = bag_clf.predict(X_test)
print(f"Bagging classifier Test Accuracy : {accuracy_score(y_test, y_pred_bag).round(3)}")

Bagging classifier Test Accuracy : 1.0
```

2. Regression

It turned out that the target which we chose for a classification gave perfect predictions. Now, we will quickly work on a regression case on the price variable. As we saw previously, this variable has an approximately symmetric distribution and is platykurtic.

```
df.head(1)
```

| | squareMeters | numberOfRooms | floors | cityPartRange | numPrevOwners | made | isNewBuilt | hasStormProtector | basement | attic | garage | ha |
|---|--------------|---------------|--------|---------------|---------------|------|------------|-------------------|----------|-------|--------|-----|
| 0 | 75523 | 3 | 23 | 3 | 8 | 2005 | 0 | 1 | 4313 | 9005 | 956 | |
| 1 | 80771 | 39 | 98 | 6 | 6 | 2015 | 1 | 0 | 3653 | 2436 | 128 | |
| 2 | 55712 | 58 | 19 | 8 | 23 | 8 | 2021 | 0 | 0 | 2937 | 8852 | 135 |
| 3 | 73216 | 47 | 6 | 10 | 4 | 2012 | 0 | 1 | 659 | 7141 | 359 | |
| 4 | 30429 | 19 | 90 | 3 | 7 | 1990 | 1 | 0 | 8435 | 2429 | 292 | |

```
# Train and Test split
X_req = df.drop('price', axis=1)
y_req = df['price']
X_train_req, X_test_req, y_train_req, y_test_req = train_test_split(X_req, y_req, train_size=0.7, random_state=16)
print(f"X_train shape : {X_train_req.shape} X_test shape : {X_test_req.shape}")
print(f"y_train shape : {y_train_req.shape} y_test shape : {y_test_req.shape}")
```

```
X_train.shape : (7000, 16)
X_test.shape : (3000, 16)
y_train.shape : (7000,)
y_test.shape : (3000,)
```

```
# Feature scaling
from sklearn.preprocessing import StandardScaler
sc = StandardScaler().fit(X_train)
X_train = sc.transform(X_train)
X_test = sc.transform(X_test)
```

Multiple linear regression

```
from sklearn.linear_model import LinearRegression
regressor = LinearRegression()
regressor.fit(X_train, y_train)
```

```
LinearRegression()

# convert to array
X_test = X_test.to_numpy()
y_test = y_test.to_numpy()
```

```
# with scaling
from sklearn.metrics import r2_score
r2_test = r2_score(y_test, regressor.predict(X_test))
np.set_printoptions(precision=2)
print(np.concatenate((y_pred.reshape(len(y_pred), 1), y_test.reshape(len(y_test), 1)), 1))
print(f"r2 score : {r2_test}")
rmseLR = rmse_cv(regressor).mean().round(2)
print(f"RMSE Linear Regression : {rmseLR}")
```

```
[[2500073.67 2499963.1 ]
 [7159590.52 7160303.2 ]
 [207905.43 2039464.6 ]
 [ 96833.72  82916.5 ]
 [ 1927.39  1927.39 ]
 RMSE Linear Regression = 1898.03
```

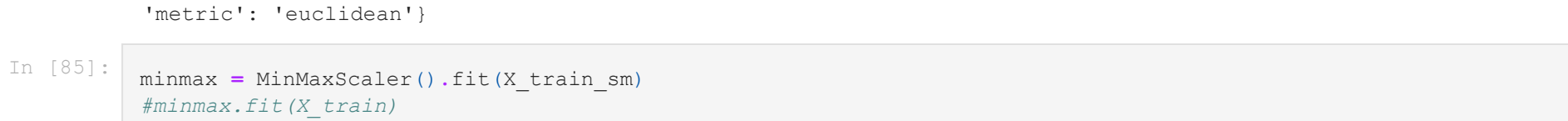
RIDGE & LASSO

```
from sklearn.linear_model import Ridge, RidgeCV, ElasticNet, LassoCV, LassoLarsCV, Lasso
from sklearn.model_selection import cross_val_score, GridSearchCV

rsmse = np.sqrt(-cross_val_score(model, X_train, y_train, scoring='neg_mean_squared_error', cv=5))
print(rsmse)

RidgeCV(rsmse)

# Ridge
ridge = Ridge(alpha=0.01, 0.3, 1, 3, 5)
cv_ridge = cross_val_score(ridge, X_train, y_train, scoring='neg_mean_squared_error', cv=5)
cv_ridge = pd.Series(cv_ridge, index=alphas)
cv_ridge.plot()
plt.xlabel('alpha')
plt.ylabel('rmse')
```



```
cv_ridge.min()
# For the Ridge regression we get a rmse of 1898
```

```
1898.1808797738165
```

```
# optimal alpha
parameters = {'alpha': (0.01, 1, 5, 10, 20),
             'fit_intercept': (True, False),
             'normalize': (False, True)}
ridge_regressor = RidgeCV(cv=5, parameters=parameters, scoring='neg_mean_squared_error', cv=5)
ridge_regressor.fit(X_train, y_train)
print(ridge_regressor.best_params_)
print(ridge_regressor.best_score_)
print(f"RIDGE RMSE on the training set : {rmse_cv(ridge_regressor).mean().round(2)}")
```

```
{'alpha': 0.01, 'fit_intercept': True, 'normalize': False}
RIDGE RMSE on the training set : 1898.03
```

LASSO

Also known as L1 Regularization, Lasso shrinks the less important features' coefficient to 0. Lasso is computationally inefficient on non-sparse cases. Ridge and Lasso can be both used for feature selection in case of collinearity to avoid overfitting. Lasso provides the posterior mode based on a prior belief that the coefficients are Laplace distributed with mean zero. So the estimates are biased towards zero.

$$\sum_{i=1}^n (y_i - \sum_{j=1}^p x_{ij} \beta_j)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

```
model_lasso = LassoCV(alphas = [1, 0.1, 0.001, 0.0005]).fit(X_train, y_train)
```

```
lasso = Lasso()
parameters = {'alpha': (20, 25, 30, 35, 40), # alpha always > 0
             'fit_intercept': (True, False),
             'normalize': (True, False)}
lasso_regressor = GridSearchCV(lasso, parameters, scoring='neg_mean_squared_error', cv=5)
lasso_regressor.fit(X_train, y_train)
print(lasso_regressor.best_params_)
print(lasso_regressor.best_score_.round(2))
```

```
{'alpha': 28, 'fit_intercept': True, 'normalize': False}
{'alpha': 28, 'fit_intercept': True, 'normalize': False}
LASSO RMSE on the training set : 1902.52
```

Ridge's RMSE was lower than Lasso's, we will use Ridge for the test set.

```
# Test
from sklearn.linear_model import Ridge
ridge_reg = Ridge(alpha=1, solver='cholesky', fit_intercept=True, normalize=False, random_state=16)
ridge_reg.fit(X_train, y_train)
y_pred_ridge = ridge_reg.predict(X_test)
print(np.concatenate((y_pred.reshape(len(y_pred), 1), y_test.reshape(len(y_test), 1)), 1))
```

```
[[2500440.27 2499963.1 ]
 [7159590.52 7160303.2 ]
 [207905.43 2039464.6 ]
 [ 96833.72  82916.5 ]
 [ 1927.39  1927.39 ]
 RMSE Linear Regression = 1898.03
```

Elastic Net

There exists a combination of Ridge and Lasso called Elastic Net.

$$Ridge = SSE + \lambda_1 ||m||^2 \rightarrow L2norm$$
$$Lasso = SSE + \lambda_2 ||m||_1 \rightarrow L1norm$$
$$ElasticNet = SSE + \lambda_1 ||m||^2 + \lambda_2 ||m||_1$$

The idea is to find the best combination of λ_1 and λ_2 in order to minimize the formula of elastic net.

Conclusion

We conclude, I have tried, during this exercise to use different machine learning techniques after oversampling the data. But tuning them was meaningless since they were initially cross-correlated. I usually use a package called PyCaret at first, to have a quick overview of the data and the algorithms, then I start doing this procedure manually. As we have seen, the predictions were perfect on the three sets (training, validation and test) but please notice that before arriving to this final version, I have tried many different targets (Pool, Price and Price categories) but I decided to work with this one because it was imbalanced and I wanted to work on a classification method and practice some sampling techniques as well. Then I tried to predict the price using some regression models (Linear, Ridge and Lasso), tune them and compare their RMSE.

PART II - Covid19 (Clustering) : by VALAP Sophia

We decided to use COVID-19 database to implement some techniques :

- Clustering with GMM algorithm
- BIC and AIC criteria
- PCA

1) Pre-processing

```
# Import database
data = pd.read_csv('Covid19Cases.csv', sep = ',',
                  data.head()
```


| | ID | Country | TotalCases | TotalDeaths | TotalRecovered | ActiveCases | TotalCasesPerMillion | TotalDeathsPerMillion | TotalTests | TotalTestsPerMillion |
|----------|----|---------|------------|-------------|----------------|-------------|----------------------|-----------------------|------------|----------------------|
| In (38): | 0 | 1 | Yemen | 9369 | 17780 | 5921 | 1670 | 306 | 58.0 | 26253.0 |
| | 1 | 2 | Vietnam | 838662 | 20555.0 | 782199 | 36908 | 8529 | 209.0 | 42517091.0 |
| | 2 | 3 | Uzbekistan | 1786637 | 1271.0 | 1742139 | 3153 | 5241 | 37.0 | 1377915.0 |
| | 3 | 4 | UAE | 737890 | 2114.0 | 731295 | 4481 | 73477 | 211.0 | 87246490.0 |
| | 4 | 5 | Turkey | 7444552 | 66180.0 | 6893476 | 48496 | 87077 | 774.0 | 89847975.0 |
| | | | | | | | | | | 1050 |

```
In (60): print(data.shape)
data.info()
```

```
(49, 11)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 49 entries, 0 to 48
Data columns (total 11 columns):
 #   Column              Non-Null Count  Dtype
---  -
 0   ID                  49 non-null    int64
 1   Country             49 non-null    object
 2   TotalCases          49 non-null    int64
 3   TotalDeaths         48 non-null    float64
 4   TotalRecovered      49 non-null    int64
 5   ActiveCases         49 non-null    int64
 6   TotalCasesPerMillion 49 non-null    float64
 7   TotalDeathsPerMillion 48 non-null    float64
 8   TotalTests          49 non-null    int64
 9   TotalTestsPerMillion 48 non-null    float64
10   TotalPopulation     49 non-null    int64
dtypes: float64(4), int64(6), object(1)
memory usage: 4.3+ KB
```

```
In (279): # Copy
X = data.copy()

# Remove ID and Country features which will not be usefull for the clustering
X.drop(['ID','Country'], axis=1, inplace=True)
X.drop(['TotalCases'], axis = 1, inplace= True)

# Missing values check
X.isnull().sum(axis=0)
```

```
Out(279): TotalCases      0
TotalDeaths     1
TotalRecovered   0
ActiveCases     0
TotalCasesPerMillion  0
TotalDeathsPerMillion  1
TotalTests       1
TotalTestsPerMillion  1
TotalPopulation   0
dtype: int64
```

There is a missing value for each of the variables TotalDeaths, TotalDeathsPerMillion, TotalTests and TotalTestsPerMillion.

There are several ways to handle missing values:

- Replace this value with the mean or median of our dataset
- Use the fillna function which offers 4 methods of filling missing values:
 - pad/fill which consist in replacing the missing value by the one preceding it in the dataset
 - backfill/bfill which replace the missing value with the next one

Since our dataset has only a small number of observations, we did not want to replace this value by the mean or the median. We also did not want to use the fillna method, which does not seem to be adapted to our case. Indeed, replacing for example the number of deaths for a country by the number of deaths of another country (which does not have the same number of inhabitants) does not seem appropriate to us. Therefore, we have decided to delete these lines

```
In (280): X = X.dropna()

# Convert as integer
X = X.astype(int)
```

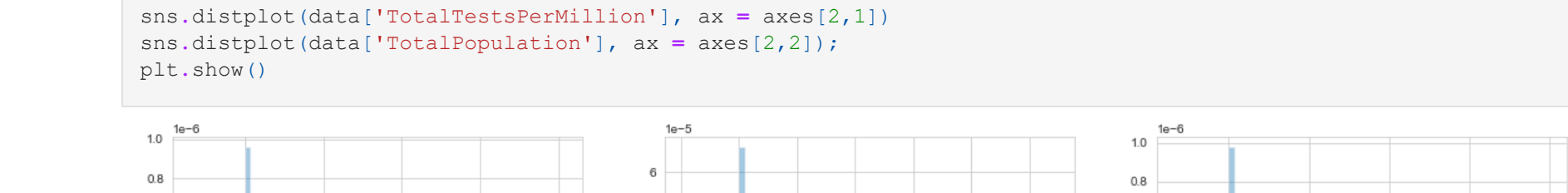
2) Data analysis

```
In (63): # Distribution of each variable
fig, axes = plt.subplots(3, 3, figsize=(20, 10))
sns.distplot(data['TotalCases'], ax = axes[0,0]);
sns.distplot(data['TotalDeaths'], ax = axes[0,1]);
sns.distplot(data['TotalRecovered'], ax = axes[0,2]);
sns.distplot(data['ActiveCases'], ax = axes[1,0]);
sns.distplot(data['TotalCasesPerMillion'], ax = axes[1,1]);
sns.distplot(data['TotalDeathsPerMillion'], ax = axes[1,2]);
sns.distplot(data['TotalTests'], ax = axes[2,0]);
sns.distplot(data['TotalTestsPerMillion'], ax = axes[2,1]);
sns.distplot(data['TotalPopulation'], ax = axes[2,2]);
plt.show()
```



We can see that most of the variables follow a Gaussian distribution except for the variables TotalCasesPerMillion, TotalTestsPerMillion and TotalDeathsPerMillion

```
In (64): # Correlation
corrMatrix = X.corr()
plt.figure(figsize=(12,6))
sns.heatmap(corrMatrix, annot = True);
```



3) Clustering - GMM

Summary:

- 1. Use of raw data without dimension reduction
- 1. Dimension reduction with PCA - 2D
- 1. Dimension reduction with PCA - 3D

1. Use of raw data without dimension reduction

```
In (281): from sklearn.preprocessing import StandardScaler, normalize
from sklearn.mixture import GaussianMixture
from sklearn.metrics import silhouette_score
import itertools
```

```
In (282): # Standardize Data
scaler = StandardScaler()
scaled_X = scaler.fit_transform(X)
```

```
# Normalizing the Data
normalized_X = normalize(scaled_X)
normalized_X = pd.DataFrame(normalized_X)
normalized_X.head()
```

```
Out(282): 0 -0.195696 -0.195218 -0.192462 -0.274204 -0.591646 -0.501332 -0.206787 -0.383017 -0.143258
1 -0.147685 -0.049149 -0.149296 -0.063428 -0.744290 -0.484299 0.131356 -0.379068 -0.001020
2 -0.182550 -0.207766 -0.178810 -0.275163 -0.548826 -0.548710 -0.207207 -0.387985 -0.141600
3 -0.042954 -0.076244 -0.040108 -0.101352 0.143551 -0.124011 0.156586 0.956247 -0.073995
4 0.218123 0.114287 0.204657 0.915197 0.169043 0.119445 0.129427 -0.014854 -0.008752
```

Use of the GMM algorithm with the BIC selection criterion

The GMM algorithm contains two main parameters:

- **n_components**: number of clusters
- **covariance_type**: type of covariance

We will use the **BIC** criterion to select the best model. It corresponds to the lowest BIC score. So, using a for loop, we will first go through our model with different possible values for the number of clusters and we will indicate the best model thanks to the BIC criterion.

Another criterion can be used to determine the performance of clustering. It is the **Silhouette Score**. For each value, we calculate the distance between our sample and the closest sample which belongs to the other class. This is called the silhouette coefficient. The silhouette score returns the average silhouette coefficient over all samples with the following properties:

- 1 is the highest value
- -1 is the lowest value
- 0 indicates overlapping clusters
- Negative values generally indicate that a sample has been assigned to the wrong group, as a different group is more similar

```
In (283): # BIC equals infinity (we want to minimize)
lowest_bic = np.infty
# List that contains the different BIC score
bic = []
# List that contains the silhouette score
s = []

n_components_range = range(2, 7)

# For loop on the number of components/clusters of the model
for n_components in n_components_range:
    # GMM algorithm
    gmm = mixture.GaussianMixture(n_components=n_components, random_state=0)
    gmm.fit(normalized_X)
    # Silhouette score calculation
    score = silhouette_score(normalized_X, gmm.predict(normalized_X)).round(3)
    # Adding the silhouette score to the list s
    s.append(score)
    # Adding the BIC score to the list bic
    bic.append(gmm.bic(normalized_X).round(2))
    # If the criterion bic found is lower than lowest_bic, then it becomes the new lowest_bic
    if bic[-1] < lowest_bic:
        lowest_bic = bic[-1]
        best_gmm = gmm

bic = np.array(bic)
clf = best_gmm
# List that contains the best gmm model silhouette score
silhouette = silhouette_score(normalized_X, best_gmm.predict(normalized_X)).round(3)
print("Silhouette score for the selected model \n", clf, "is", silhouette, " for a BIC score of ", min(bic))
```

Silhouette score for the selected model GaussianMixture(n_components=6, random_state=0) is 0.417 for a BIC score of -745.32

```
In (284): plt.figure(figsize=(15,3))

ax1 = plt.subplot(1, 2, 1)
ax2 = plt.subplot(1, 2, 2)

ax1.title.set_text('Critere BIC')
ax1.plot(n_components_range, bic)

ax2.title.set_text('Silhouette score')
ax2.plot(n_components_range, s);
```



We observe here that the number of clusters that minimizes the BIC criterion is 6. This is also the highest silhouette score.

In the following, we will integrate the type of covariance in the search for the best model.

The type of covariance will allow us to visualize the clusters in different shapes

```
In (285): # Formatting the algorithm as a function
def algo_gmm(X):
    lowest_bic = np.infty
    bic = []
    s = []
    n_components_range = range(2, 7)
    cv_types = ["spherical", "tied", "diag", "full"]
    for cv_type in cv_types:
        # GMM algorithm
        for n_components in n_components_range:
            gmm = mixture.GaussianMixture(n_components=n_components, covariance_type=cv_type)
            gmm.fit(X)
            score = silhouette_score(X, gmm.predict(X)).round(3)
            s.append(score)
            bic.append(gmm.bic(X).round(2))
            if bic[-1] < lowest_bic:
                lowest_bic = bic[-1]
                best_gmm = gmm

    bic = np.array(bic)
    clf = best_gmm
    silhouette = silhouette_score(X, best_gmm.predict(X)).round(3)
    return clf, bic, silhouette, n_components_range
```

```
In (286): for i in range(1,11):
    best_gmm, min_bic, silhouette, n_components_range = algo_gmm(normalized_X)
    print("***** Round ", i, " : Model : ", best_gmm, " BIC : ", min(bic), " Silhouette score : ", silhouette)
```

```
***** Round 1 : Model : GaussianMixture(n_components=2) BIC : -745.32 Silhouette score : 0.417
***** Round 2 : Model : GaussianMixture(n_components=3) BIC : -745.32 Silhouette score : 0.4
***** Round 3 : Model : GaussianMixture(n_components=4) BIC : -745.32 Silhouette score : 0.438
***** Round 4 : Model : GaussianMixture(n_components=5) BIC : -745.32 Silhouette score : 0.431
***** Round 5 : Model : GaussianMixture(n_components=6) BIC : -745.32 Silhouette score : 0.438
***** Round 6 : Model : GaussianMixture(n_components=7) BIC : -745.32 Silhouette score : 0.378
***** Round 7 : Model : GaussianMixture(n_components=8) BIC : -745.32 Silhouette score : 0.287
***** Round 8 : Model : GaussianMixture(n_components=9) BIC : -745.32 Silhouette score : 0.417
***** Round 9 : Model : GaussianMixture(n_components=10) BIC : -745.32 Silhouette score : 0.37
***** Round 10 : Model : GaussianMixture(n_components=11) BIC : -745.32 Silhouette score : 0.39
```

We can also visualize with which classes our data belong

```
best_gmm.predict(normalized_X)
```

```
Out(287): array([5, 5, 3, 1, 5, 2, 5, 5, 0, 3, 5, 5, 4, 2, 4, 2, 3, 5, 5, 4, 4,
1, 4, 5, 5, 4, 4, 4, 5, 4, 0, 1, 2, 1, 3, 4, 3, 2, 5, 5, 5, 4,
0, 4, 5], dtype=int64)
```

We will create a dataframe of size (1,9) that will contain the average of each column of our dataframe. This will allow us to predict which class the average of our data would belong to

```
In (288): # Mean of each column
X_mean = X.mean(axis = 0)

# Creation and filling of the dataset
X_new = pd.DataFrame(0, index = [0], columns = list(X_mean.index))
i = 0
for col in X_new.columns:
    X_new[col] = X_mean.values[i]
    i = i + 1
X_new
```

```
Out(288): TotalCases TotalRecovered ActiveCases TotalCasesPerMillion TotalDeathsPerMillion TotalTests TotalTestsPerMillion
0 1.64131e+06 24235255319 1.573923e+06 43154.12766 45698.957447 464.787234 3.015736e+07 1.198093e+06
```

```
In (289): # Standardize Data
scaler = StandardScaler()
scaled_X_new = scaler.fit_transform(X_new)
```

```
# Normalizing data
normalized_X_new = normalize(scaled_X_new)
normalized_X_new = pd.DataFrame(normalized_X_new)
```

```
In (290): # Prediction on the new dataframe
best_gmm.predict(normalized_X_new)
```

```
Out(290): array([4], dtype=int64)
```

The average of our data would belong to the class 4). The predict_proba function displays the probability that our point belongs to each class. This function allows us not only to assign a point to a cluster but also to check if the point could not have belonged to a different cluster. If we want, we can assign a different status to these points.

```
In (291): # Use of predict_proba function
predict_proba = best_gmm.predict_proba(normalized_X)
predict_proba = pd.DataFrame(predict_proba)
```

```
# Adding max column to get the max proba of each row and check if our values are satisfying
predict_proba['Max_proba'] = predict_proba.max(axis = 1)
predict_proba[predict_proba['Max_proba'] < 0.5]
```

```
Out(291): 0 1 2 3 4 5 Max_proba
```

No proba max below 0.5. The results are satisfying \ In the following, we will use dimension reduction via PCA to visualize our clusters

2. Dimension reduction with PCA - 2D

```
In (292): from sklearn.decomposition import PCA

# PCA dimension reduction
pca = PCA(n_components = 2)
X_reduced = pca.fit_transform(normalized_X)
X_reduced = pd.DataFrame(X_reduced)
X_reduced.columns = ['P1', 'P2']
X_reduced.head(2)
```

```
Out(292): P1 P2
0 -0.725313 -0.121381
1 -0.747473 0.415930
```

We observe that when the random_state is different from 0, the results are different. To make sure to have an optimal result, we will run the algorithm with several different values

```
In (296): for i in range(1,11):
    best_gmm, bic, silhouette, n_components_range = algo_gmm(X_reduced)
    print("***** Round ", i, " : Model : ", best_gmm, " BIC : ", min(bic), " Silhouette score : ", silhouette)
```

```
***** Round 1 : Model : GaussianMixture(n_components=2) BIC : 116.85 Silhouette score : 0.404
***** Round 2 : Model : GaussianMixture(n_components=3) BIC : 116.85 Silhouette score : 0.404
***** Round 3 : Model : GaussianMixture(n_components=4) BIC : 116.85 Silhouette score : 0.404
***** Round 4 : Model : GaussianMixture(n_components=5) BIC : 116.85 Silhouette score : 0.404
***** Round 5 : Model : GaussianMixture(n_components=6) BIC : 116.85 Silhouette score : 0.404
***** Round 6 : Model : GaussianMixture(n_components=7) BIC : 116.85 Silhouette score : 0.404
***** Round 7 : Model : GaussianMixture(n_components=8) BIC : 116.85 Silhouette score : 0.404
***** Round 8 : Model : GaussianMixture(n_components=9) BIC : 116.85 Silhouette score : 0.404
***** Round 9 : Model : GaussianMixture(n_components=10) BIC : 116.85 Silhouette score : 0.404
***** Round 10 : Model : GaussianMixture(n_components=11) BIC : 116.85 Silhouette score : 0.404
```

We do not have any specification for the covariance type of the model. Therefore, it is the default covariance which is of type full

```
In (297): # Model
gmm = GaussianMixture(n_components = 2)
gmm.fit(X_reduced)
cluster = gmm.predict(X_reduced)
```

```
# Visualizing the clustering
plt.scatter(X_reduced['P1'], X_reduced['P2'],
           c = GaussianMixture(n_components=2).fit_predict(X_reduced), cmap=plt.cm.winter, alpha = 0.6)
plt.show()
```



- **ICL criterion**

We can observe that the classes are not homogeneous. There is a criterion that would have allowed us to obtain more homogeneous classes: the ICL criterion. It would be the most optimal criterion for model selection in a clustering case. It is a BIC criterion with a penalty when the classes are not homogeneous enough. Due to the lack of implementation of the function in the python libraries, we will not use this criterion. The models will continue to be selected with the BIC criterion.

- **AIC criterion**

Another very popular model selection criterion is the AIC criterion. **Warning! It is not recommended to use the AIC criterion when you have a clustering problem.** The purpose of this criterion is to try to do density estimation. It will propose as many classes as possible to improve the fit to the data. Its goal is not to find clusters. It has not been developed in a Gaussian setting. We will see this with the graph below:

```
In (298): n_components = np.arange(2, 21)
models = [GaussianMixture(n, covariance_type='full', random_state=0).fit(X_reduced)
          for n in n_components]

plt.plot(n_components, (n.bic(X_reduced) for n in models), label='BIC')
plt.plot(n_components, (n.aic(X_reduced) for n in models), label='AIC')
plt.legend(loc='best')
plt.xlabel('n_components');
```



As we can see, the AIC criterion suggests taking the largest number of clusters while the BIC criterion finds its minimum at 2 clusters. The more the number of iterations is large, the more the AIC criterion moves away from the BIC criterion until it completely moves away from it.

3. Dimension reduction with PCA - 3D

```
In (302): # Dimension reduction
pca = PCA(n_components = 3)
X_principal = pca.fit_transform(normalized_X)
X_principal = pd.DataFrame(X_principal)
X_principal.columns = ['P1', 'P2', 'P3']
```

```
# GMM
for i in range(1,11):
    best_gmm, bic, silhouette, n_components_range = algo_gmm(X_principal)
    print("***** Round ", i, " : Model : ", best_gmm, " BIC : ", min(bic), " Silhouette score : ", silhouette)
```

```
***** 1 éme tour : Modèle : GaussianMixture(covariance_type='diag', n_components=2) BIC : 161.8 Silhouette score : 0.287
***** 2 éme tour : Modèle : GaussianMixture(n_components=6) BIC : 157.24 Silhouette score : 0.377
***** 3 éme tour : Modèle : GaussianMixture(covariance_type='diag', n_components=2) BIC : 161.8 Silhouette score : 0.287
***** 4 éme tour : Modèle : GaussianMixture(covariance_type='diag', n_components=2) BIC : 161.8 Silhouette score : 0.287
***** 5 éme tour : Modèle : GaussianMixture(covariance_type='diag', n_components=2) BIC : 161.8 Silhouette score : 0.287
***** 6 éme tour : Modèle : GaussianMixture(covariance_type='diag', n_components=2) BIC : 161.8 Silhouette score : 0.287
***** 7 éme tour : Modèle : GaussianMixture(covariance_type='diag', n_components=2) BIC : 161.8 Silhouette score : 0.287
***** 8 éme tour : Modèle : GaussianMixture(covariance_type='diag', n_components=2) BIC : 161.8 Silhouette score : 0.287
***** 9 éme tour : Modèle : GaussianMixture(covariance_type='diag', n_components=2) BIC : 161.8 Silhouette score : 0.287
***** 10 éme tour : Modèle : GaussianMixture(n_components=6) BIC : 160.38 Silhouette score : 0.465
```

We choose the GMM with 6 components because it has the best silhouette score

```
In (303): # 3D plot result
from mpl_toolkits import mplot3d

fig = plt.figure()
ax = plt.axes(projection='3d')
ax.scatter3D(X_principal['P1'], X_principal['P2'], X_principal['P3'],
            c = GaussianMixture(n_components = 6).fit_predict(X_principal))
plt.show()
```



The average of our data would belong to the class 4). The predict_proba function displays the probability that our point belongs to each class. This function allows us not only to assign a point to a cluster but also to check if the point could not have belonged to a different cluster. If we want, we can assign a different status to these points.