

Customer Dataset Segregation and Analysis

Building a Customer analysis and segregation using Support Vector Machine Classifier (SVC) and visualization involves several key steps, including defining analysis objectives, collecting campaign data, and processing and cleaning the data.

Here's a step-by-step guide on how to get started:

1. Define Analysis Objectives:

Start by clearly defining the objectives of your customer segregation and analysis. What specific insights are you trying to gain from the data?

The objective of the design analysis for “Customer Segregation and Analysis” is to assess the effectiveness, efficiency, and scalability of the proposed system in partitioning customers into distinct segments based on relevant criteria, and subsequently evaluating the analytical capabilities to derive actionable insights and tailored strategies for each segment. This analysis aims to ensure that the design optimally addresses the business objectives, accommodates diverse customer profiles, and provides a robust framework for ongoing analysis and adaptation.

2.Implementing a dataset

Kaggle dataset

Dataset Link: <https://www.kaggle.com/datasets/akram24/mall-customers>

CustomerID	Genre	Age	Annual Income (k\$)	S
1	Female	18	15	1
0001	Male	19	15	39
0002	Male	21	15	81
0003	Female	20	16	6
0004	Female	23	16	77
0005	Female	31	17	40
0006	Female	22	17	76

3. Collect Customer Data:

To collect data for your analysis, you'll need access to the data source. This may involve reaching out to relevant public, commodities, retail shops, or partners. Ensure you have permission to use and analyse the data.

```
In [1]: #imports necessary libraries to do basic things on the dataset
import pandas as pd
import numpy as np

import seaborn as sns
import matplotlib.pyplot as plt

print('Successfully imported')
```

Successfully imported

4. Process and Clean Data:

Data processing and cleaning are critical to ensuring the quality and accuracy of your analysis. Here are some steps to follow:

- a. Data Integration:
 - If your customer data is stored in different formats or sources, integrate them into a single dataset. This may involve using ETL (Extract, Transform, Load) tools.
- b. Data Cleaning:
 - Handle missing data: Identify and deal with missing values, either by imputation or removal.
 - Remove duplicates: Eliminate duplicate records.
 - Data format standardisation: Ensure that date formats, units of Measurement, and naming conventions are consistent.

```
In [3]: #Check the dataset for missing data
if data.isnull().sum().sum() == 0 :
    print ('There is no missing data in our dataset')
else:
    print('There is {} missing data in our dataset '.format(data.isnull().sum().sum()))

There is 1892 missing data in our dataset
```

```
In [4]: #Check our missing data from which columns and how many unique features they have.
frame = pd.concat([data.isnull().sum(), data.nunique(), data.dtypes], axis = 1, sort= False)
frame
```

c. Data Transformation:

- Convert data types: Ensure that data types are appropriate for analysis (e.g., dates as date objects, numbers as numeric types).
- Create calculated fields: Generate new variables if needed.
- Aggregation: Summarise data as needed, e.g., daily, weekly, or monthly aggregates.

d. Data Quality Check:

- Check for outliers and anomalies that may affect analysis.
- Validate data against the defined analysis objectives.

5. Classification of customers:

In this part, the objective will be to adjust a classifier that will classify consumers in the different client categories that were established in the previous section. The objective is to make this classification possible at the first visit. To fulfill this objective, I will test several classifiers implemented in scikit-learn. First, in order to simplify their use, I define a class that allows to interface several of the functionalities common to these different classifiers:

```
[64]: class Class_Fit(object):
    def __init__(self, clf, params=None):
        if params:
            self.clf = clf(**params)
        else:
            self.clf = clf()

    def train(self, x_train, y_train):
        self.clf.fit(x_train, y_train)

    def predict(self, x):
        return self.clf.predict(x)

    def grid_search(self, parameters, Kfold):
        self.grid = GridSearchCV(estimator = self.clf, param_grid = parameters,
        cv = Kfold)

    def grid_fit(self, X, Y):
        self.grid.fit(X, Y)

    def grid_predict(self, X, Y):
        self.predictions = self.grid.predict(X)
        print("Precision: {:.2f} % ".format(100*metrics.accuracy_score(Y, self.
        predictions)))
```

Since the goal is to define the class to which a client belongs and this, as soon as its first visit, I only keep the variables that describe the content of the basket, and do not take into account the variables related to the frequency of visits or variations of the basket price over time:

```
[65]: columns = ['mean', 'categ_0', 'categ_1', 'categ_2', 'categ_3', 'categ_4' ]
X = selected_customers[columns]
Y = selected_customers['cluster']
```

Finally, I split the dataset in train and test sets:

```
[66]: X_train, X_test, Y_train, Y_test = model_selection.train_test_split(X, Y,
train_size = 0.8)
```

5.1 Support Vector Machine Classifier (SVC)

The first classifier I use is the SVC classifier. In order to use it, I create an instance of the `Class_Fit` class and then call `grid_search()`. When calling this method, I provide as parameters:

- the hyperparameters for which I will seek an optimal value - the number of folds to be used for cross-validation

```
svc = Class_Fit(clf = svm.LinearSVC)
svc.grid_search(parameters = [{'C': np.logspace(-2, 2, 10)}], Kfold = 5)
```

Once this instance is created, I adjust the classifier to the training data:

```
svc.grid_fit(X = X_train, Y = Y_train)
```

then I can test the quality of the prediction with respect to the test data:

```
svc.grid_predict(X_test, Y_test)
```

Precision: 80.75

Confusion matrix The accuracy of the results seems to be correct. Nevertheless, let us remember that when the different classes were defined, there was an imbalance in size between the classes obtained. In particular, one class contains around 40% of the clients. It is therefore interesting to look at how the predictions and real values compare to the breasts of the different classes. This is the subject of the confusion matrices and to represent them, I use the code of the [sklearn documentation](#):

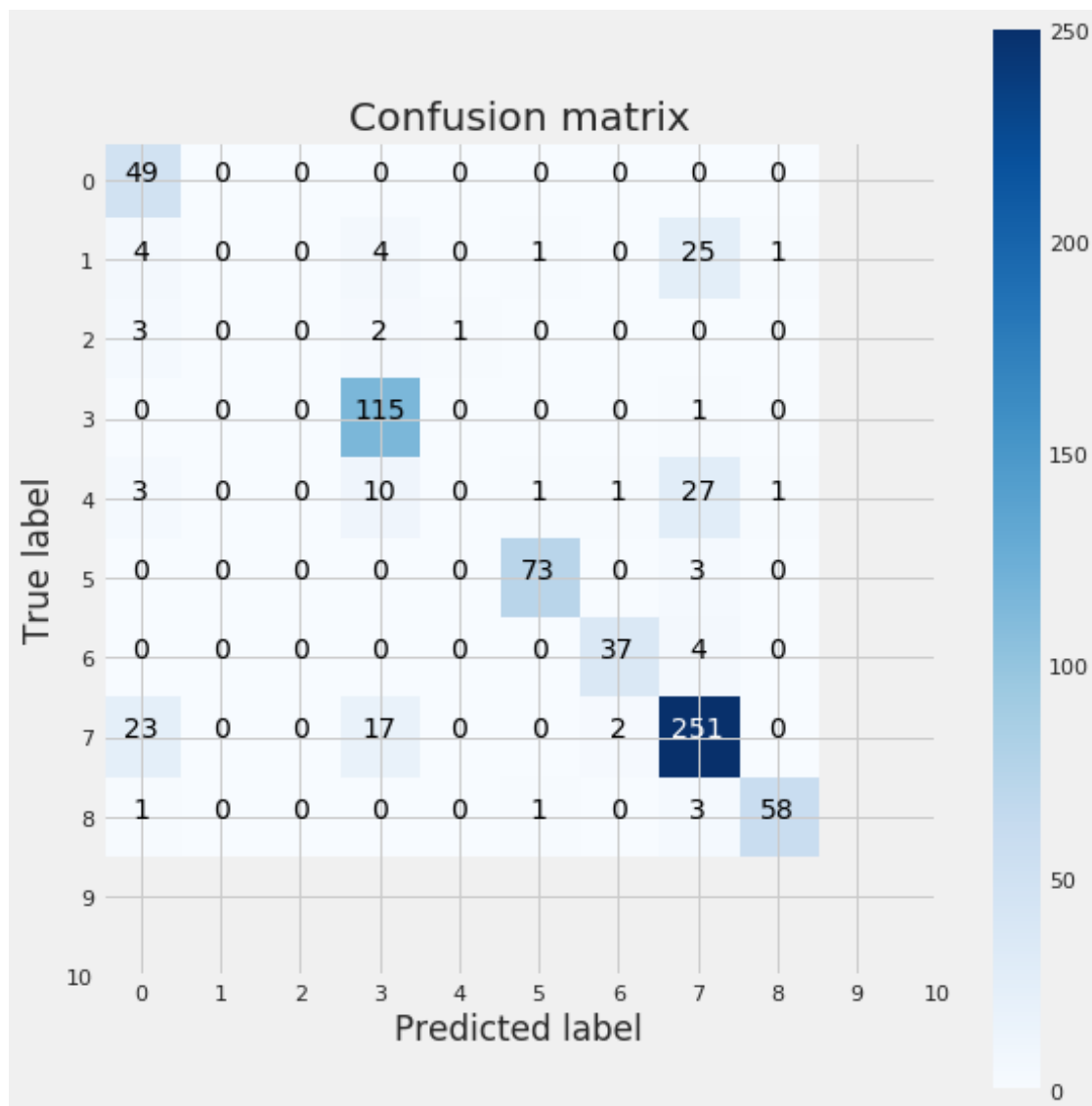
```
[70]: def plot_confusion_matrix(cm, classes, normalize=False, title='Confusion_
      ↪matrix', cmap=plt.cm.Blues):
      if normalize:
          cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
          print("Normalized confusion matrix")
      else:
          print('Confusion matrix, without normalization')
      #-----
      plt.imshow(cm, interpolation='nearest', cmap=cmap)
      plt.title(title)
      plt.colorbar()
      tick_marks = np.arange(len(classes))
      plt.xticks(tick_marks, classes, rotation=0)
      plt.yticks(tick_marks, classes)
      #-----
```

```
      fmt = '.2f' if normalize else 'd'
      thresh = cm.max() / 2.
      for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
          plt.text(j, i, format(cm[i, j], fmt),
                  horizontalalignment="center",
                  color="white" if cm[i, j] > thresh else "black")
      #-----
      plt.tight_layout()
      plt.ylabel('True label')
      plt.xlabel('Predicted label')
```

from which I create the following representation:

```
[71]: class_names = [i for i in range(11)]
      cnf_matrix = confusion_matrix(Y_test, svc.predictions)
      np.set_printoptions(precision=2)
      plt.figure(figsize = (8,8))
      plot_confusion_matrix(cnf_matrix, classes=class_names, normalize = False,
      ↪title='Confusion matrix')
```

Confusion matrix, without normalization



Learning curve A typical way to test the quality of a fit is to draw a learning curve. In particular, this type of curves allow to detect possible drawbacks in the model, linked for example to over- or under-fitting. This also shows to which extent the mode could benefit from a larger data sample. In order to draw this curve, I use the [scikit-learn documentation code again](#)

```
[72]: def plot_learning_curve(estimator, title, X, y, ylim=None, cv=None,
                                n_jobs=-1, train_sizes=np.linspace(.1, 1.0, 10)):
    """Generate a simple plot of the test and training learning curve"""
    plt.figure()
    plt.title(title)
    if ylim is not None:
```

```

plt.ylim(*ylim)
plt.xlabel("Training examples")
plt.ylabel("Score")
train_sizes, train_scores, test_scores = learning_curve(
    estimator, X, y, cv=cv, n_jobs=n_jobs, train_sizes=train_sizes)
train_scores_mean = np.mean(train_scores, axis=1)
train_scores_std = np.std(train_scores, axis=1)
test_scores_mean = np.mean(test_scores, axis=1)
test_scores_std = np.std(test_scores, axis=1)
plt.grid()

plt.fill_between(train_sizes, train_scores_mean - train_scores_std,
                 train_scores_mean + train_scores_std, alpha=0.1, color="r")
plt.fill_between(train_sizes, test_scores_mean - test_scores_std,
                 test_scores_mean + test_scores_std, alpha=0.1, color="g")
plt.plot(train_sizes, train_scores_mean, 'o-', color="r", label="Training score")
plt.plot(train_sizes, test_scores_mean, 'o-', color="g", label="Cross-validation score")

plt.legend(loc="best")
return plt

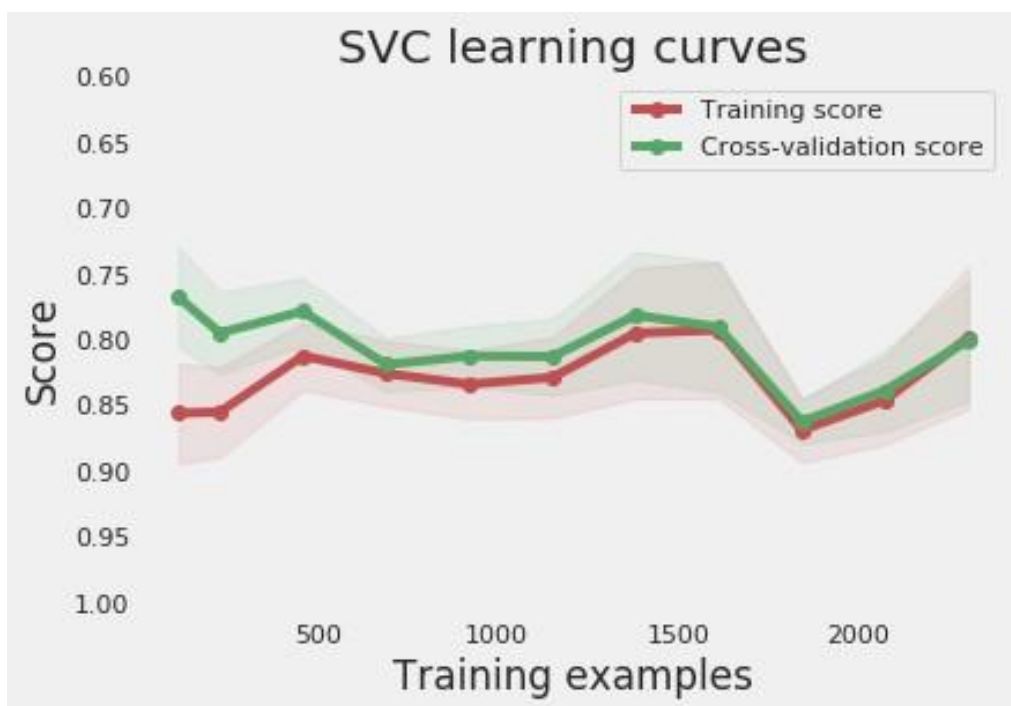
```

from which I represent the learning curve of the SVC classifier:

```

[73]: g = plot_learning_curve(svc.grid.best_estimator_,
                             "SVC learning curves", X_train, Y_train, ylim = [1.01,
                             0.6],
                             cv = 5, train_sizes = [0.05, 0.1, 0.2, 0.3, 0.4, 0.5,
                             0.6, 0.7, 0.8, 0.9, 1])

```



On this curve, we can see that the train and cross-validation curves converge towards the same limit when the sample size increases. This is typical of modeling with low variance and proves that the model does not suffer from overfitting. Also, we can see

that the accuracy of the training curve is correct which is synonymous of a low bias. Hence the model does not underfit the data.

5.2 Logistic Regression

I now consider the logistic regression classifier. As before, I create an instance of the `Class_Fit` class, adjust the model on the training data and see how the predictions compare to the real values:

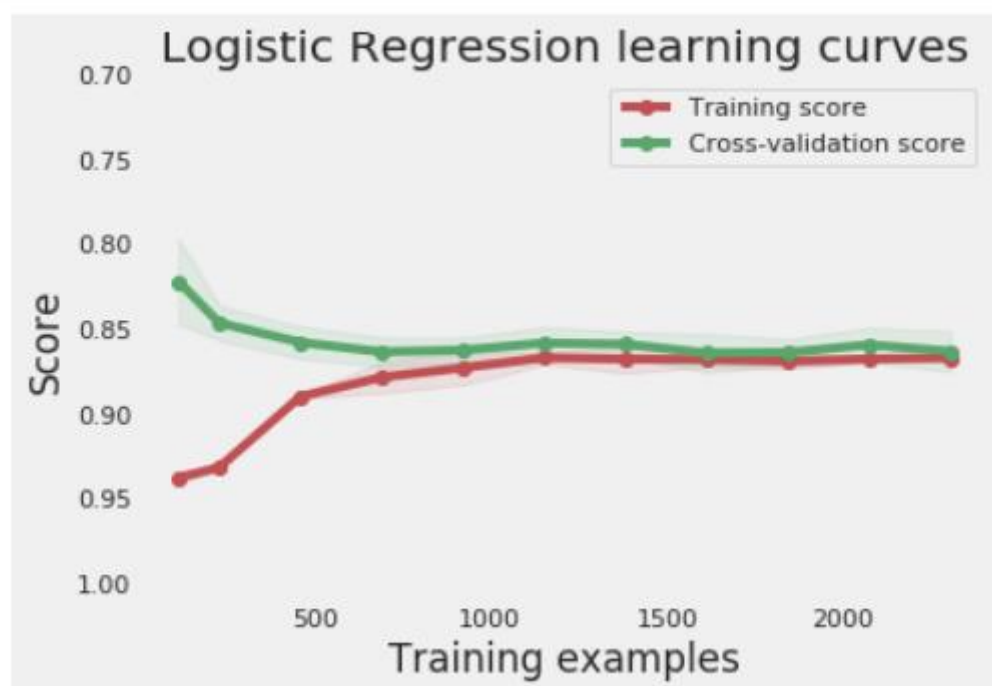
```
[74]: lr = Class_Fit(clf = linear_model.LogisticRegression)
lr.grid_search(parameters = [{'C': np.logspace(-2, 2, 20)}], Kfold = 5)
lr.grid_fit(X = X_train, Y = Y_train)
lr.grid_predict(X_test, Y_test)
```

Precision: 86.29 %

Then, I plot the learning curve to have a feeling of the quality of the model:

```
[75]: g = plot_learning_curve(lr.grid.best_estimator_, "Logistic Regression learning_
    ↪curves", X_train, Y_train,
    ylim = [1.01, 0.7], cv = 5,
```

```
    ↪0.8, 0.9, 1])
    train_sizes = [0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1])
```

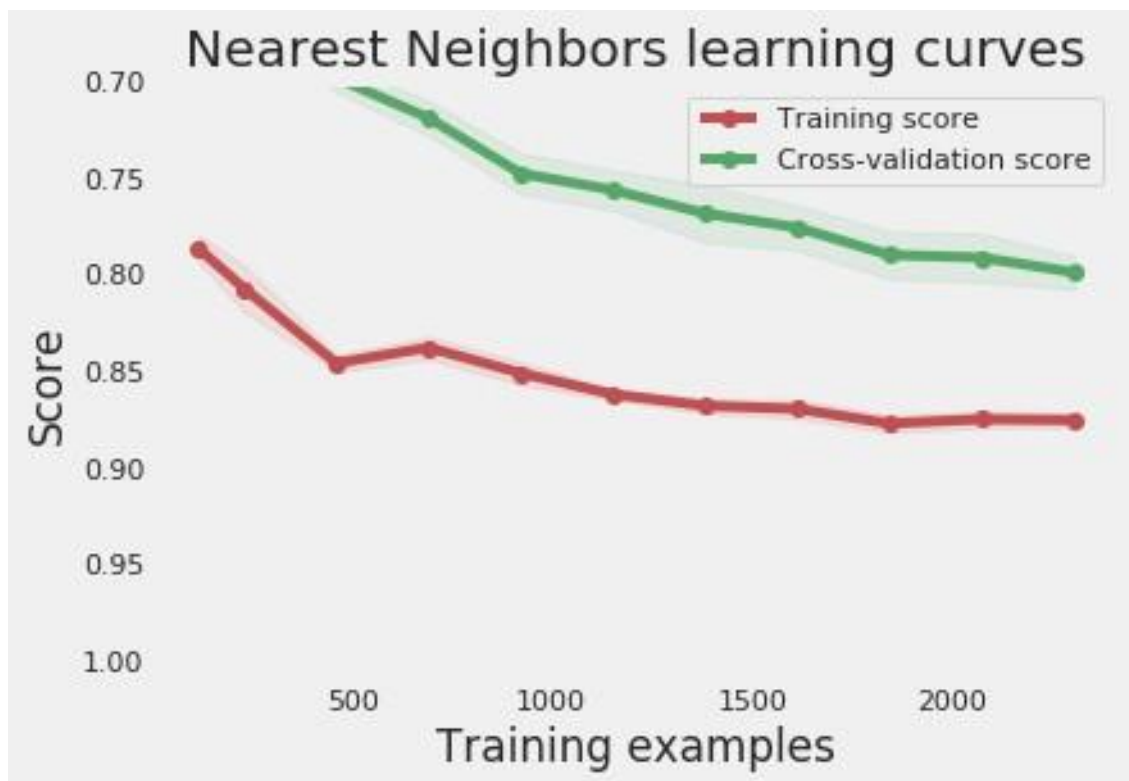


1.5.3 5.3 k-Nearest Neighbors

```
[76]: knn = Class_Fit(clf = neighbors.KNeighborsClassifier)
      knn.grid_search(parameters = [{'n_neighbors': np.arange(1,50,1)}], Kfold = 5)
      knn.grid_fit(X = X_train, Y = Y_train)
      knn.grid_predict(X_test, Y_test)
```

Precision: 79.78 %

```
[77]: g = plot_learning_curve(knn.grid.best_estimator_, "Nearest Neighbors learning_
      ↪curves", X_train, Y_train,
      ylim = [1.01, 0.7], cv = 5,
      train_sizes = [0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7,
      ↪0.8, 0.9, 1])
```

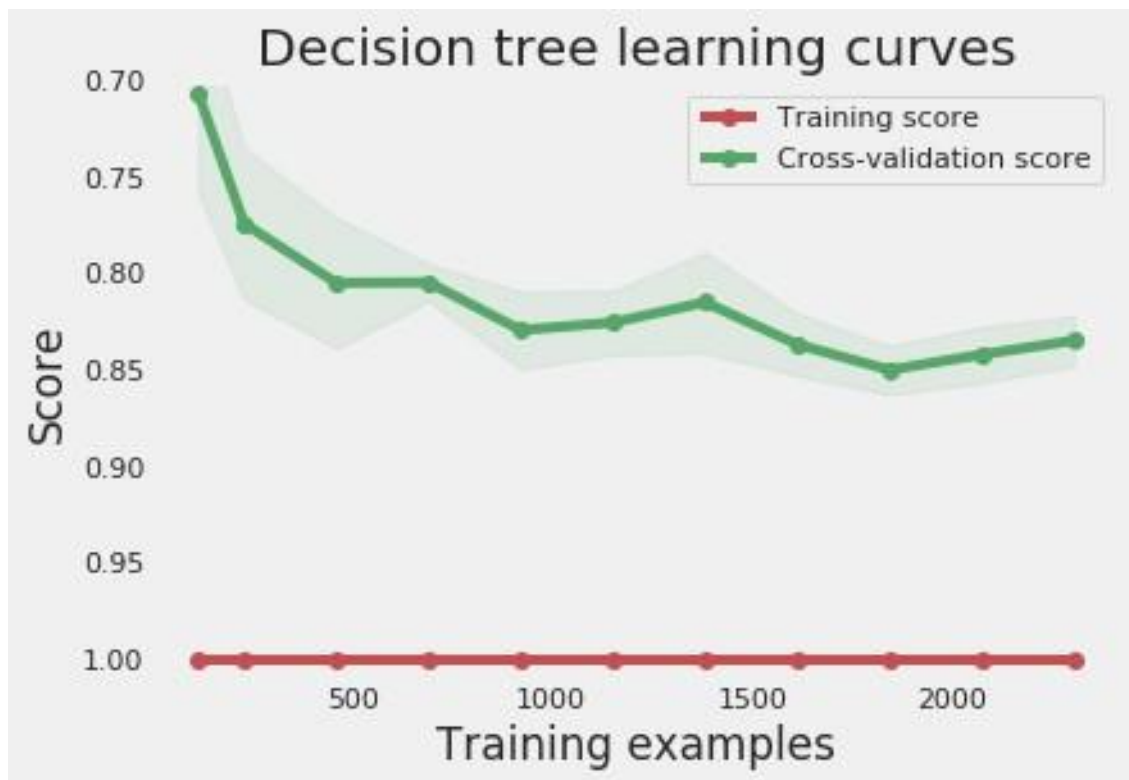


5.4 Decision Tree

```
[78]: tr = Class_Fit(clf = tree.DecisionTreeClassifier)
      tr.grid_search(parameters = [{'criterion' : ['entropy', 'gini'], 'max_features' :
      ↪['sqrt', 'log2']}], Kfold = 5)
      tr.grid_fit(X = X_train, Y = Y_train)
      tr.grid_predict(X_test, Y_test)
```

Precision: 83.24 %

```
[79]: g = plot_learning_curve(tr.grid.best_estimator_, "Decision tree learning_
      ↪curves", X_train, Y_train,
      ylim = [1.01, 0.7], cv = 5,
      train_sizes = [0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7,
      ↪0.8, 0.9, 1])
```

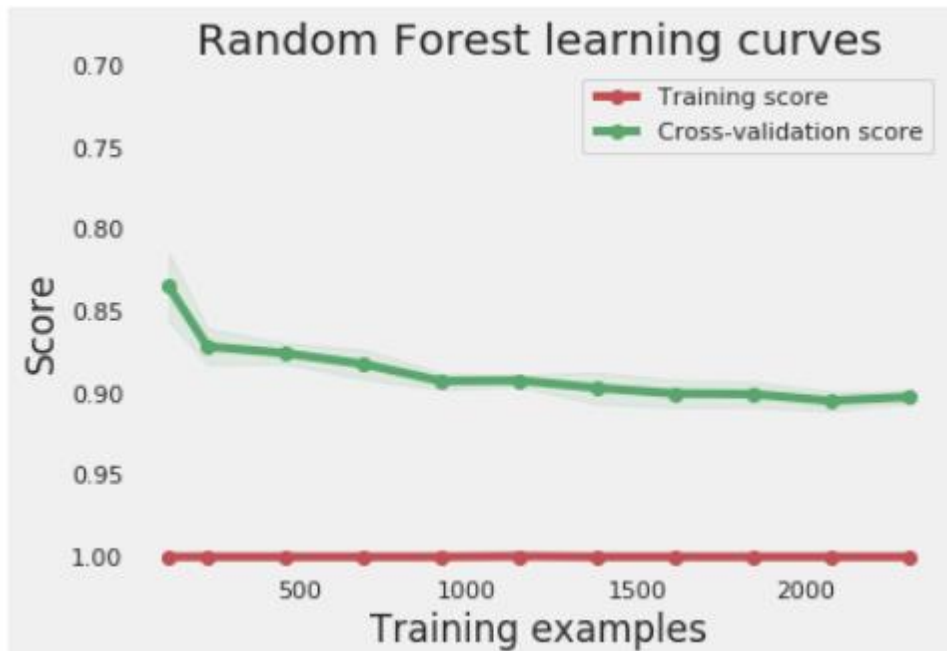


5.5 Random Forest

```
[80]: rf = Class_Fit(clf = ensemble.RandomForestClassifier)
      param_grid = {'criterion' : ['entropy', 'gini'], 'n_estimators' : [20, 40, 60, 80, 100],
                    'max_features' : ['sqrt', 'log2']}
      rf.grid_search(parameters = param_grid, Kfold = 5)
      rf.grid_fit(X = X_train, Y = Y_train)
      rf.grid_predict(X_test, Y_test)
```

Precision: 89.61 %

```
[81]: g = plot_learning_curve(rf.grid.best_estimator_, "Random Forest learning curves", X_train, Y_train,
                              ylim = [1.01, 0.7], cv = 5,
                              train_sizes = [0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1])
```

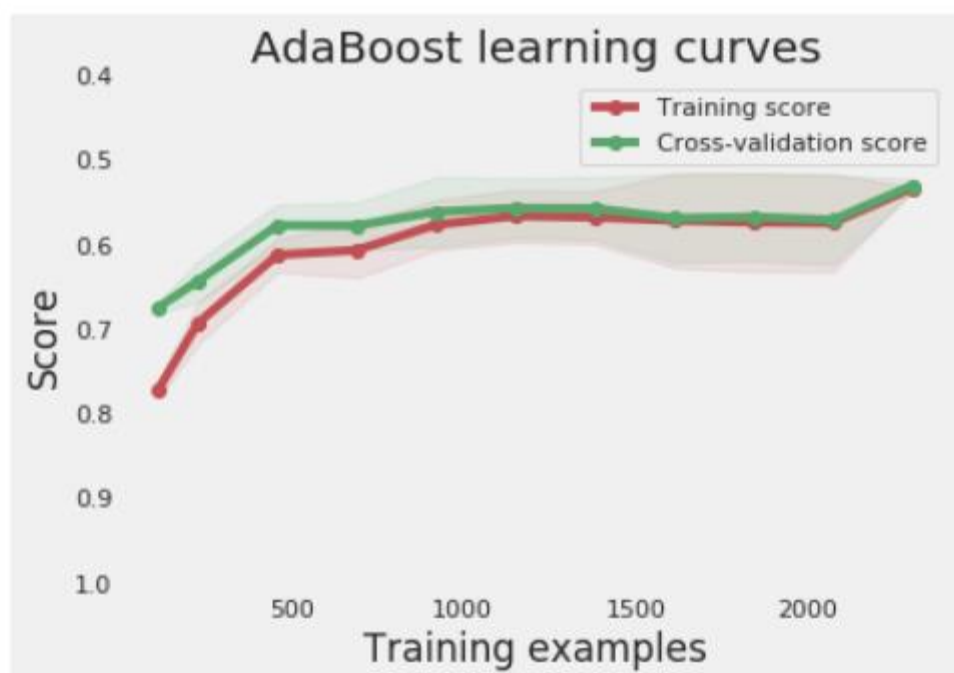


1.5.4 5.6 AdaBoost Classifier

```
[82]: ada = Class_Fit(clf = AdaBoostClassifier)
      param_grid = {'n_estimators' : [10, 20, 30, 40, 50, 60, 70, 80, 90, 100]}
      ada.grid_search(parameters = param_grid, Kfold = 5)
      ada.grid_fit(X = X_train, Y = Y_train)
      ada.grid_predict(X_test, Y_test)
```

Precision: 54.57 %

```
[83]: g = plot_learning_curve(ada.grid.best_estimator_, "AdaBoost learning curves",
      X_train, Y_train,
      ylim = [1.01, 0.4], cv = 5,
      train_sizes = [0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7,
      0.8, 0.9, 1])
```

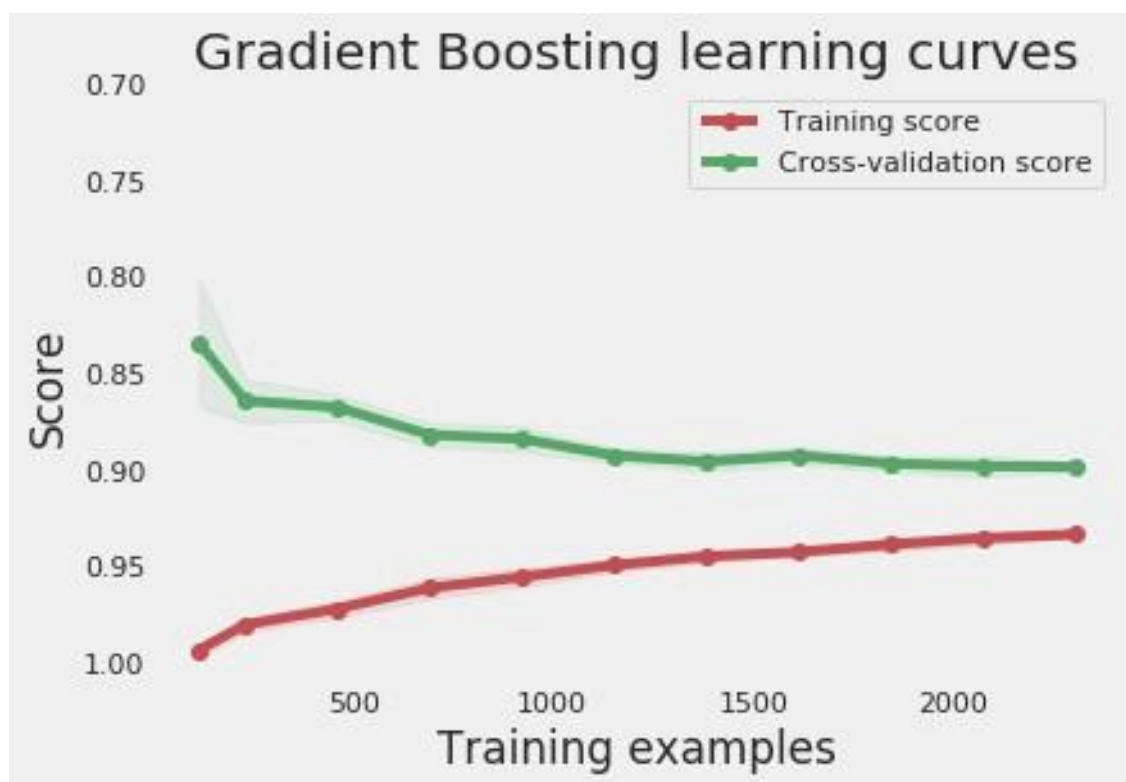


5.7 Gradient Boosting Classifier

```
[84]: gb = Class_Fit(clf = ensemble.GradientBoostingClassifier)
      param_grid = {'n_estimators' : [10, 20, 30, 40, 50, 60, 70, 80, 90, 100]}
      gb.grid_search(parameters = param_grid, Kfold = 5)
      gb.grid_fit(X = X_train, Y = Y_train)
      gb.grid_predict(X_test, Y_test)
```

Precision: 89.47 %

```
[85]: g = plot_learning_curve(gb.grid.best_estimator_, "Gradient Boosting learning_
      ↪curves", X_train, Y_train,
      ylim = [1.01, 0.7], cv = 5,
      train_sizes = [0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7,
      ↪0.8, 0.9, 1])
```



5.8 Let's vote !

Finally, the results of the different classifiers presented in the previous sections can be combined to improve the classification model. This can be achieved by selecting the customer category as the one indicated by the majority of classifiers. To do this, I use the `VotingClassifier` method of the `sklearn` package. As a first step, I adjust the parameters of the various classifiers using the *best* parameters previously found:

```
rf_best = ensemble.RandomForestClassifier(**rf.grid.best_params_)
gb_best = ensemble.GradientBoostingClassifier(**gb.grid.best_params_)
svc_best = svm.LinearSVC(**svc.grid.best_params_)
tr_best = tree.DecisionTreeClassifier(**tr.grid.best_params_)
knn_best = neighbors.KNeighborsClassifier(**knn.grid.best_params_)
lr_best = linear_model.LogisticRegression(**lr.grid.best_params_)
```

Then, I define a classifier that merges the results of the various classifiers:

```
votingC = ensemble.VotingClassifier(estimators=[('rf', rf_best), ('gb', gb_best),  
                                              ('knn', knn_best)],  
                                   voting='soft')
```

and train it:

```
votingC = votingC.fit(X_train, Y_train)
```

Finally, we can create a prediction for this model:

```
predictions = votingC.predict(X_test)  
print("Precision: {:.2f} %".format(100*metrics.accuracy_score(Y_test,   
                                                                predictions)))
```

Note that when defining the `votingC` classifier, I only used a sub-sample of the whole set of classifiers defined above and only retained the *Random Forest*, the *k-Nearest Neighbors* and the *Gradient Boosting* classifiers. In practice, this choice has been done with respect to the performance of the classification carried out in the next section.

6. Testing predictions

In the previous section, a few classifiers were trained in order to categorize customers. Until that point, the whole analysis was based on the data of the first 10 months. In this section, I test the model the last two months of the dataset, that has been stored in the `set_test` dataframe:

```
[90]: basket_price = set_test.copy(deep = True)
```

In a first step, I regroup reformattes these data according to the same procedure as used on the training set. However, I am correcting the data to take into account the difference in time between the two datasets and weights the variables `** count **` and `** sum **` to obtain an equivalence with the training set:

```
[91]: transactions_per_user=basket_price.groupby(by=['CustomerID'])['Basket Price'].  
      <agg(['count', 'min', 'max', 'mean', 'sum'])  
      for i in range(5):  
          col = 'categ_{}'.format(i)  
          transactions_per_user.loc[:,col] = basket_price.  
          <groupby(by=['CustomerID'])[col].sum() /\br/>                                              transactions_per_user['sum']*100  
  
      transactions_per_user.reset_index(drop = False, inplace = True)  
      basket_price.groupby(by=['CustomerID'])['categ_0'].sum()  
  
      #-----  
      # Correcting time range  
      transactions_per_user['count'] = 5 * transactions_per_user['count']  
      transactions_per_user['sum'] = transactions_per_user['count'] *  
      <transactions_per_user['mean']  
  
      transactions_per_user.sort_values('CustomerID', ascending = True)[:5]
```

CustomerID		count	min	...	categ_2	categ_3	categ_4
0	12347	10	224.82	...	12.696657	37.379043	5.634767
1	12349	5	1757.55	...	4.513101	37.877728	20.389178
2	12352	5	311.73	...	6.672441	34.398358	17.290604
3	12356	5	58.35	...	0.000000	100.000000	0.000000
4	12357	5	6207.67	...	5.089832	22.895547	25.189000

Then, I convert the dataframe into a matrix and retain only variables that define the category to which consumers belong. At this level, I recall the method of normalization that had been used on the training set:

```
[92]: list_cols = [
    → ['count', 'min', 'max', 'mean', 'categ_0', 'categ_1', 'categ_2', 'categ_3', 'categ_4']
    # -----
    matrix_test = transactions_per_user[list_cols].as_matrix()
    scaled_test_matrix = scaler.transform(matrix_test)
```

Each line in this matrix contains a consumer's buying habits. At this stage, it is a question of using these habits in order to define the category to which the consumer belongs. These categories

have been established in Section 4. ** At this stage, it is important to bear in mind that this step does not correspond to the classification stage itself. **Here, we prepare the test data by defining the category to which the customers belong. However, this definition uses data obtained over a period of 2 months (via the variables count , min , max ** and ** sum **).** The classifier defined in Section 5 uses a more restricted set of variables that will be defined from the first purchase of a client.

Here it is a question of using the available data over a period of two months and using this data to define the category to which the customers belong. Then, the classifier can be tested by comparing

its predictions with these categories. In order to define the category to which the clients belong, I recall the instance of the kmeans method used in section 4. The predict method of this instance calculates the distance of the consumers from the centroids of the 11 client classes and the smallest

distance will define the belonging to the different categories:

```
[93]: Y = kmeans.predict(scaled_test_matrix)
```

Finally, in order to prepare the execution of the classifier, it is sufficient to select the variables on which it acts:

```
[94]: columns = ['mean', 'categ_0', 'categ_1', 'categ_2', 'categ_3', 'categ_4']
X = transactions_per_user[columns]
```

It remains only to examine the predictions of the different classifiers that have been trained in section 5:

```
[95]: classifiers = [(svc, 'Support Vector Machine'),
                    (lr, 'Logostic Regression'),
                    (knn, 'k-Nearest Neighbors'),
```

```

        (tr, 'Decision Tree'),
        (rf, 'Random Forest'),
        (gb, 'Gradient Boosting')]
#-----
for clf, label in classifiers:
    print(30*'_ ', '\n{}'.format(label))
    clf.grid_predict(X, Y)

```

Support Vector Machine

Precision: 65.93 %

Logostic Regression

Precision: 71.34 %

k-Nearest Neighbors

Precision: 67.58 %

Decision Tree

Precision: 71.38 %

Random Forest

Precision: 75.38 %

Gradient Boosting

Precision: 75.23 %

Finally, as anticipated in Section 5.8, it is possible to improve the quality of the classifier by combining their respective predictions. At this level, I chose to mix *Random Forest*, *Gradient Boosting* and *k-Nearest Neighbors* predictions because this leads to a slight improvement in predictions:

```

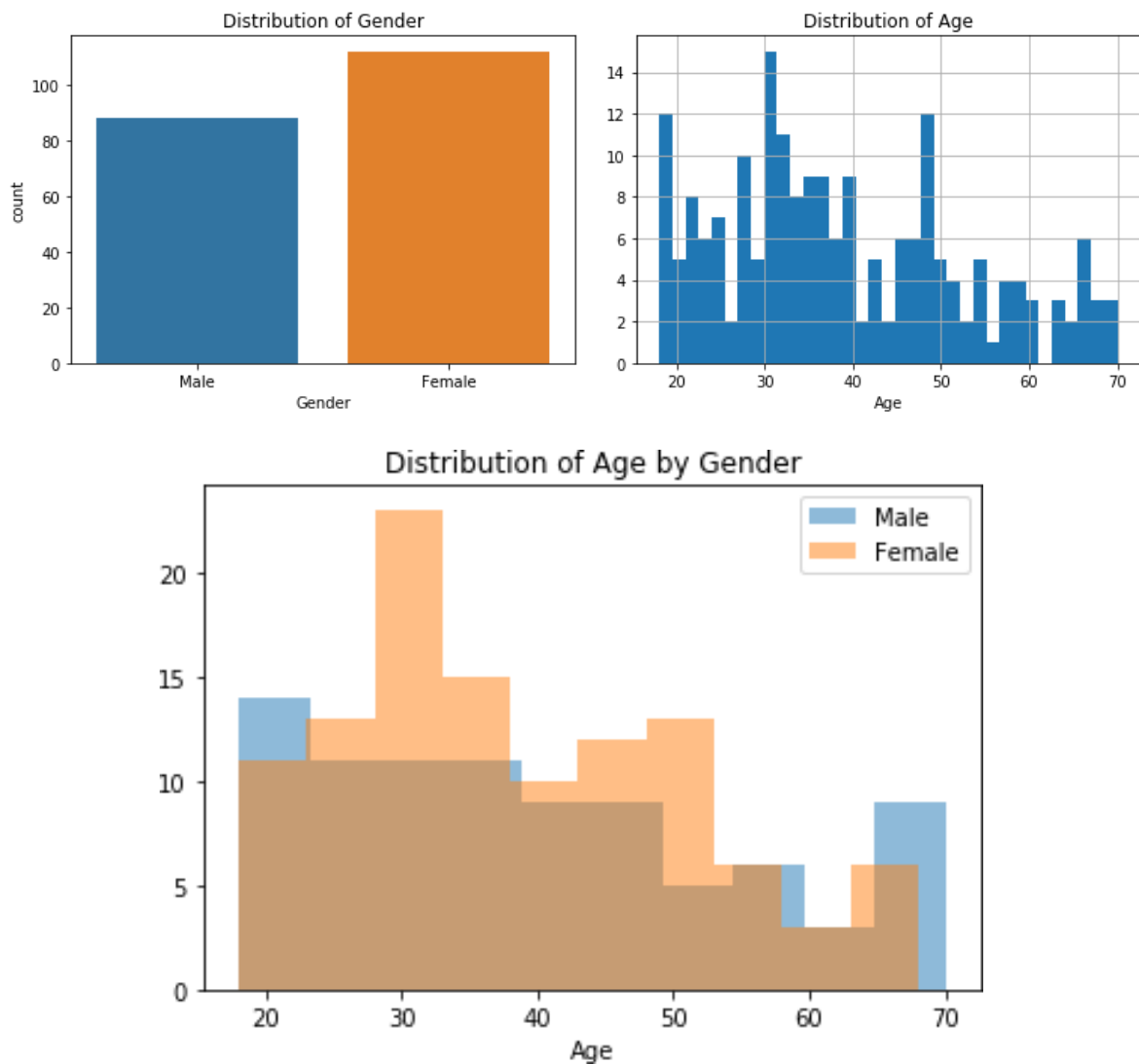
[96]: predictions = votingC.predict(X)
      print("Precision: {:.2f} % ".format(100*metrics.accuracy_score(Y, predictions)))

      Precision: 75.46 %

```

7. Create Visualizations and Reports:

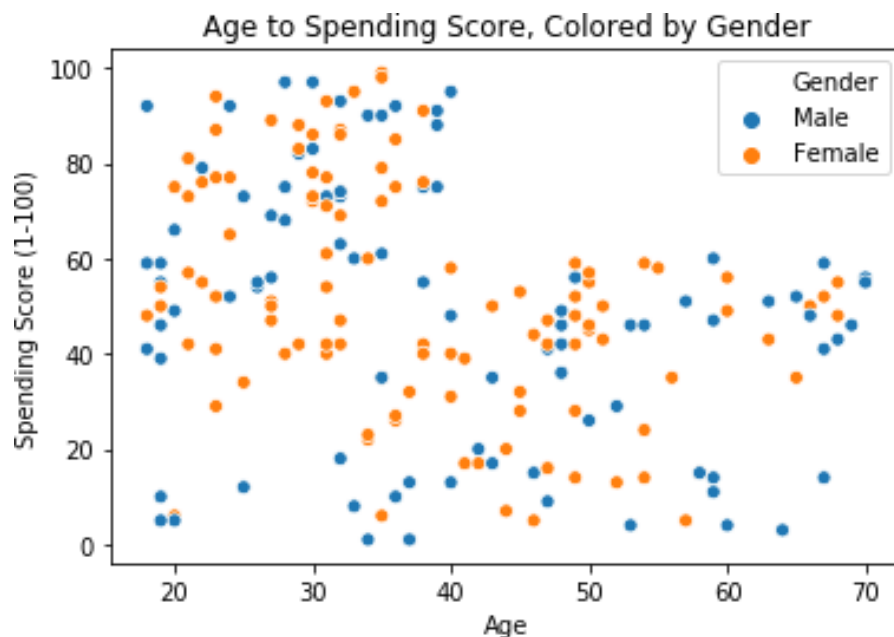
With your data, you can start building visualizations and reports to address your analysis objectives. You can use various chart types, tables, and graphs to present the data effectively.



8. Analyze the Data and Share Insights:

Analyze the data to draw insights and conclusions related to your dataset's effectiveness. Use features with to perform statistical analysis or apply business intelligence techniques to discover patterns and trends.

Finally, share the insights and findings from your analysis with relevant stakeholders, managers, product manager, executives, or other decision-makers. Use the reports and visualizations created in IBM Cognos for this purpose.



9. Iterate and Refine:

Data analysis is an iterative process. If your analysis reveals areas for Improvement in the analysis, work with stakeholders to refine strategies and potentially run new programs to improve the efficiency.

Remember to document your analysis process, including the data sources, data cleaning steps, and the rationale behind your analysis choices. This documentation will be valuable for future reference and for ensuring the transparency and reproducibility of your work.

11. Conclusion

The work described in this notebook is based on a database providing details on purchases made on an E-commerce platform over a period of one year. Each entry in the dataset describes the purchase of a product, by a particular customer and at a given date.

In total, approximately 4000 clients appear in the database. Given the available information, I decided to develop a classifier that allows to anticipate the type of purchase that a customer will make, as well as the number of visits that he will make during a year, and this from its first visit to the E-commerce site.

The first stage of this work consisted in describing the different products sold by the site, which was the subject of a first classification. There, I grouped the different products into 5 main categories of goods.

In a second step, I performed a classification of the customers by analyzing their consumption habits over a period of 10 months. I have classified clients into 11 major categories based on the type of products they usually buy, the number of visits they make and the amount they spent during the 10 months.

Once these categories established, I finally trained several classifiers whose objective is to be able to classify consumers in one of these 11 categories and this from their first purchase. For this, the classifier is based on 5 variables which are:

mean: amount of the basket of the current purchase

categ_ N with $N \in [0 : 4]$: percentage spent in product category with index N

Finally, the quality of the predictions of the different classifiers was tested over the last two months of the dataset. The data were then processed in two steps: first, all the data was considered (over the 2 months) to define the category to which each client belongs, and then, the classifier predictions were compared with this category assignment. I then found that 75% of clients are awarded the right classes. The performance of the classifier therefore seems correct given the potential shortcomings of the current model. In particular, a bias that has not been dealt with concerns the seasonality of purchases and the fact that purchasing habits will potentially depend on the time of year (for example, Christmas).

In practice, this seasonal effect may cause the categories defined over a 10-month period to be quite different from those extrapolated from the last two months. In order to correct such bias, it would be beneficial to have data that would cover a longer period of time.