

```
In [8]: import os
import pandas as pd
```

```
In [9]: import pandas as pd
import numpy as np
import matplotlib as mpl
import matplotlib.pyplot as plt
import seaborn as sns
import datetime, nltk, warnings
import matplotlib.cm as cm
import itertools
from pathlib import Path
from sklearn.preprocessing import StandardScaler
from sklearn.cluster import KMeans
from sklearn.metrics import silhouette_samples, silhouette_score
from sklearn import preprocessing, model_selection, metrics, feature_selection
from sklearn.model_selection import GridSearchCV, learning_curve
from sklearn.svm import SVC
from sklearn.metrics import confusion_matrix
from sklearn import neighbors, linear_model, svm, tree, ensemble

from sklearn.ensemble import AdaBoostClassifier
from sklearn.decomposition import PCA
from IPython.display import display, HTML
import plotly.graph_objs as go
from plotly.offline import init_notebook_mode, iplot
init_notebook_mode(connected=True)
warnings.filterwarnings("ignore")
plt.rcParams["patch.force_edgecolor"] = True
plt.style.use('fivethirtyeight')
mpl.rcParams['patch', edgecolor = 'dimgray', linewidth=1)
%matplotlib inline
```

```
In [10]: os.chdir('D:/My ML Simulations/Model Building')
```

```
In [11]: df_initial = pd.read_csv('data.csv', encoding="ISO-8859-1",
dtype={'CustomerID': str, 'InvoiceID': str})
```

```
In [12]: print('Dataframe dimensions:', df_initial.shape)
```

Dataframe dimensions: (541909, 8)

```
In [13]: df_initial['InvoiceDate'] = pd.to_datetime(df_initial['InvoiceDate'])
```

```
In [14]: tab_info=pd.DataFrame(df_initial.dtypes).T.rename(index={0:'column type'})
tab_info=tab_info.append(pd.DataFrame(df_initial.isnull().sum()).T.rename(index={0:'null values (
tab_info=tab_info.append(pd.DataFrame(df_initial.isnull().sum()/df_initial.shape[0]*100).T.
rename(index={0:'null values (%)'}))
display(tab_info)
```

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
column type	object	object	object	int64	datetime64[ns]	float64	object	object
null values (nb)	0	0	1454	0	0	0	135080	0
null values (%)	0	0	0.268311	0	0	0	24.9267	0

```
In [15]: display(df_initial[:5])
```

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
0	536365	85123A	WHITE HANGING HEART T-LIGHT HOLDER	6	2010-12-01 08:26:00	2.55	17850	United Kingdom
1	536365	71053	WHITE METAL LANTERN	6	2010-12-01 08:26:00	3.39	17850	United Kingdom
2	536365	84406B	CREAM CUPID HEARTS COAT HANGER	8	2010-12-01 08:26:00	2.75	17850	United Kingdom
3	536365	84029G	KNITTED UNION FLAG HOT WATER BOTTLE	6	2010-12-01 08:26:00	3.39	17850	United Kingdom
4	536365	84029E	RED WOOLLY HOTTIE WHITE HEART.	6	2010-12-01 08:26:00	3.39	17850	United Kingdom

```
In [16]: df_initial.dropna(axis = 0, subset = ['CustomerID'], inplace = True)
print('Dataframe dimensions:', df_initial.shape)
```

Dataframe dimensions: (406829, 8)

```
In [17]: tab_info=pd.DataFrame(df_initial.dtypes).T.rename(index={0:'column type'})
tab_info=tab_info.append(pd.DataFrame(df_initial.isnull().sum()).T.rename(index={0:'null values (nb)'}))
tab_info=tab_info.append(pd.DataFrame(df_initial.isnull().sum()/df_initial.shape[0]*100).T.rename(index={0:'null values (%)'}))
display(tab_info)
```

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
column type	object	object	object	int64	datetime64[ns]	float64	object	object
null values (nb)	0	0	0	0	0	0	0	0
null values (%)	0	0	0	0	0	0	0	0

```
In [18]: print('Entrées dupliquées: {}'.format(df_initial.duplicated().sum()))
df_initial.drop_duplicates(inplace = True)
```

Entrées dupliquées: 5225

```
In [19]: temp = df_initial[['CustomerID', 'InvoiceNo', 'Country']].groupby(['CustomerID', 'InvoiceNo', 'Country']).reset_index(drop = False)
temp = temp.reset_index(drop = False)
countries = temp['Country'].value_counts()
print('Display dataframe with countries: {}'.format(len(countries)))
```

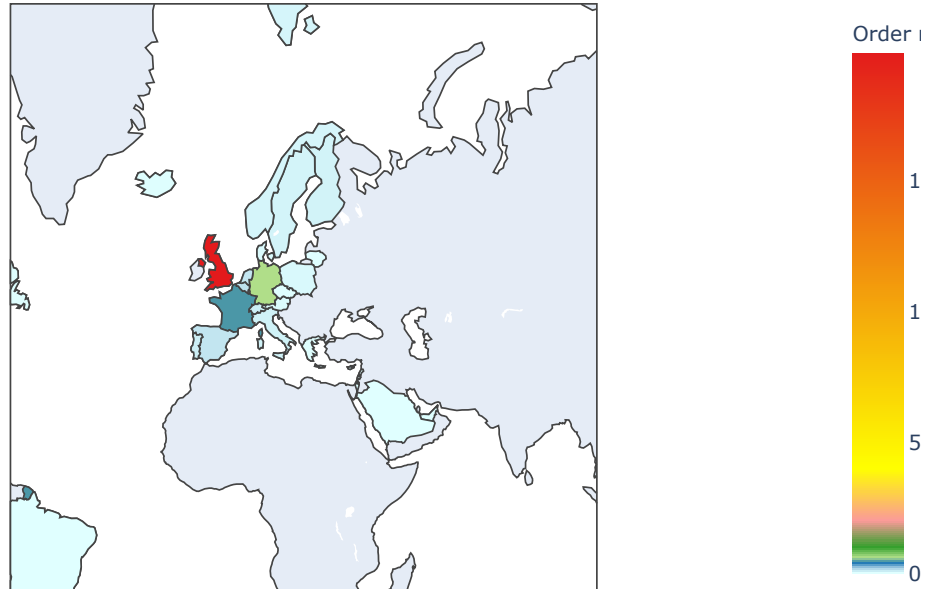
Display dataframe with countries: 37

```
In [20]: data = dict(type='choropleth',
locations = countries.index,
locationmode = 'country names', z = countries,
text = countries.index, colorbar = {'title':'Order nb.'},
colorscale=[[0, 'rgb(224,255,255)'],
[0.01, 'rgb(166,206,227)'], [0.02, 'rgb(31,120,180)'],
[0.03, 'rgb(178,223,138)'], [0.05, 'rgb(51,160,44)'],
[0.10, 'rgb(251,154,153)'], [0.20, 'rgb(255,255,0)'],
[1, 'rgb(227,26,28)']],
reversescale = False)
```

```
In [21]: layout = dict(title='Number of orders per country',
geo = dict(showframe = True, projection={'type':'mercator'}))
```

```
In [22]: choromap = go.Figure(data = [data], layout = layout)
iplot(choromap, validate=False)
```

Number of orders per country



```
In [23]: pd.DataFrame([{'products': len(df_initial['StockCode'].value_counts()),
'transactions': len(df_initial['InvoiceNo'].value_counts()),
'customers': len(df_initial['CustomerID'].value_counts()),
}], columns = ['products', 'transactions', 'customers'], index = ['quantity'])
```

Out[23]:

	products	transactions	customers
quantity	3684	22190	4372

```
In [24]: temp = df_initial.groupby(by=['CustomerID', 'InvoiceNo'], as_index=False)['InvoiceDate'].count()
nb_products_per_basket = temp.rename(columns = {'InvoiceDate': 'Number of products'})
nb_products_per_basket[:10].sort_values('CustomerID')
```

Out[24]:

	CustomerID	InvoiceNo	Number of products
0	12346	541431	1
1	12346	C541433	1
2	12347	537626	31
3	12347	542237	29
4	12347	549222	24
5	12347	556201	18
6	12347	562032	22
7	12347	573511	47
8	12347	581180	11
9	12348	539318	17

```
In [25]: nb_products_per_basket['order_canceled'] = nb_products_per_basket['InvoiceNo'].apply(lambda x:int
display(nb_products_per_basket[:5])
```

	CustomerID	InvoiceNo	Number of products	order_canceled
0	12346	541431	1	0
1	12346	C541433	1	1
2	12347	537626	31	0
3	12347	542237	29	0
4	12347	549222	24	0

```
In [26]: n1 = nb_products_per_basket['order_canceled'].sum()
n2 = nb_products_per_basket.shape[0]
print('Number of orders canceled: {}/{} ({:.2f}%) '.format(n1, n2, n1/n2*100))
```

Number of orders canceled: 3654/22190 (16.47%)

```
In [27]: display(df_initial.sort_values('CustomerID')[:5])
```

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
61619	541431	23166	MEDIUM CERAMIC TOP STORAGE JAR	74215	2011-01-18 10:01:00	1.04	12346	United Kingdom
61624	C541433	23166	MEDIUM CERAMIC TOP STORAGE JAR	-74215	2011-01-18 10:17:00	1.04	12346	United Kingdom
286623	562032	22375	AIRLINE BAG VINTAGE JET SET BROWN	4	2011-08-02 08:48:00	4.25	12347	Iceland
72260	542237	84991	60 TEATIME FAIRY CAKE CASES	24	2011-01-26 14:30:00	0.55	12347	Iceland
14943	537626	22772	PINK DRAWER KNOB ACRYLIC EDWARDIAN	12	2010-12-07 14:57:00	1.25	12347	Iceland

```
In [28]: df_check = df_initial[df_initial['Quantity'] < 0][['CustomerID', 'Quantity',
                                                         'StockCode', 'Description', 'UnitPrice']]

for index, col in df_check.iterrows():
    if df_initial[(df_initial['CustomerID'] == col[0]) & (df_initial['Quantity'] == -col[1])
                  & (df_initial['Description'] == col[2])].shape[0] == 0:
        print(df_check.loc[index])
        print(15*'-'+>+' HYPOTHESIS NOT FULFILLED')
        break
```

```
CustomerID      14527
Quantity        -1
StockCode        D
Description      Discount
UnitPrice        27.5
Name: 141, dtype: object
-----> HYPOTHESIS NOT FULFILLED
```

```
In [29]: df_check = df_initial[(df_initial['Quantity'] < 0) & (df_initial['Description'] != 'Discount')][
        ['CustomerID', 'Quantity', 'StockCode',
         'Description', 'UnitPrice']]

for index, col in df_check.iterrows():
    if df_initial[(df_initial['CustomerID'] == col[0]) & (df_initial['Quantity'] == -col[1])
                  & (df_initial['Description'] == col[2])].shape[0] == 0:
        print(index, df_check.loc[index])
        print(15*'-'+>+' HYPOTHESIS NOT FULFILLED')
        break
```

```
154 CustomerID      15311
Quantity          -1
StockCode          35004C
Description      SET OF 3 COLOURED FLYING DUCKS
UnitPrice          4.65
Name: 154, dtype: object
-----> HYPOTHESIS NOT FULFILLED
```

```

In [31]: df_cleaned = df_initial.copy(deep = True)
df_cleaned['QuantityCanceled'] = 0

entry_to_remove = [] ; doubtful_entry = []

for index, col in df_initial.iterrows():
    if (col['Quantity'] > 0) or col['Description'] == 'Discount': continue
    df_test = df_initial[(df_initial['CustomerID'] == col['CustomerID']) &
                        (df_initial['StockCode'] == col['StockCode']) &
                        (df_initial['InvoiceDate'] < col['InvoiceDate']) &
                        (df_initial['Quantity'] > 0)].copy()

    #
    # Cancellation WITHOUT counterpart
    if (df_test.shape[0] == 0):
        doubtful_entry.append(index)

    #
    # Cancellation WITH a counterpart
    elif (df_test.shape[0] == 1):
        index_order = df_test.index[0]
        df_cleaned.loc[index_order, 'QuantityCanceled'] = -col['Quantity']
        entry_to_remove.append(index)

    #
    # Various counterparts exist in orders: we delete the last one
    elif (df_test.shape[0] > 1):
        df_test.sort_index(axis=0, ascending=False, inplace = True)
        for ind, val in df_test.iterrows():
            if val['Quantity'] < -col['Quantity']: continue
            df_cleaned.loc[ind, 'QuantityCanceled'] = -col['Quantity']
            entry_to_remove.append(index)
            break

```

```

In [32]: print("entry_to_remove: {}".format(len(entry_to_remove)))
print("doubtfull_entry: {}".format(len(doubtfull_entry)))

```

```

entry_to_remove: 7521
doubtfull_entry: 1226

```

```

In [33]: df_cleaned.drop(entry_to_remove, axis = 0, inplace = True)
df_cleaned.drop(doubtfull_entry, axis = 0, inplace = True)
remaining_entries = df_cleaned[(df_cleaned['Quantity'] < 0) & (df_cleaned['StockCode'] != 'D')]
print("nb of entries to delete: {}".format(remaining_entries.shape[0]))
remaining_entries[:5]

```

```

nb of entries to delete: 48

```

Out[33]:

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	QuantityCanceled
77598	C542742	84535B	FAIRY CAKES NOTEBOOK A6 SIZE	-94	2011-01-31 16:26:00	0.65	15358	United Kingdom	0
90444	C544038	22784	LANTERN CREAM GAZEBO	-4	2011-02-15 11:32:00	4.95	14659	United Kingdom	0
111968	C545852	22464	HANGING METAL HEART LANTERN	-5	2011-03-07 13:49:00	1.65	14048	United Kingdom	0
116064	C546191	47566B	TEA TIME PARTY BUNTING	-35	2011-03-10 10:57:00	0.70	16422	United Kingdom	0
132642	C547675	22263	FELT EGG COSY LADYBIRD	-49	2011-03-24 14:07:00	0.66	17754	United Kingdom	0

```
In [34]: df_cleaned[(df_cleaned['CustomerID'] == 14048) & (df_cleaned['StockCode'] == '22464')]
```

Out[34]:

InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	QuantityCanceled
-----------	-----------	-------------	----------	-------------	-----------	------------	---------	------------------

```
In [35]: list_special_codes = df_cleaned[df_cleaned['StockCode'].str.contains('^[a-zA-Z]+', regex=True)][  
list_special_codes
```

Out[35]: array(['POST', 'D', 'C2', 'M', 'BANK CHARGES', 'PADS', 'DOT'],
dtype=object)

```
In [36]: list_special_codes:  
{:<15} -> {:<30}".format(code, df_cleaned[df_cleaned['StockCode'] == code]['Description'].unique())
```

POST	-> POSTAGE
D	-> Discount
C2	-> CARRIAGE
M	-> Manual
BANK CHARGES	-> Bank Charges
PADS	-> PADS TO MATCH ALL CUSHIONS
DOT	-> DOTCOM POSTAGE

```
In [37]: df_cleaned['TotalPrice'] = df_cleaned['UnitPrice'] * (df_cleaned['Quantity'] - df_cleaned['QuantityCanceled'])  
df_cleaned.sort_values('CustomerID')[:5]
```

Out[37]:

	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	QuantityCanceled	1
61619	541431	23166	MEDIUM CERAMIC TOP STORAGE JAR	74215	2011-01-18 10:01:00	1.04	12346	United Kingdom	74215	
148288	549222	22375	AIRLINE BAG VINTAGE JET SET BROWN	4	2011-04-07 10:43:00	4.25	12347	Iceland	0	
428971	573511	22698	PINK REGENCY TEACUP AND SAUCER	12	2011-10-31 12:25:00	2.95	12347	Iceland	0	
428970	573511	47559B	TEA TIME OVEN GLOVE	10	2011-10-31 12:25:00	1.25	12347	Iceland	0	
428969	573511	47567B	TEA TIME KITCHEN APRON	6	2011-10-31 12:25:00	5.95	12347	Iceland	0	

```
In [38]: #
# somme des achats / utilisateur & commande
temp = df_cleaned.groupby(by=['CustomerID', 'InvoiceNo'], as_index=False)['TotalPrice'].sum()
basket_price = temp.rename(columns = {'TotalPrice': 'Basket Price'})
#
# date de la commande
df_cleaned['InvoiceDate_int'] = df_cleaned['InvoiceDate'].astype('int64')
temp = df_cleaned.groupby(by=['CustomerID', 'InvoiceNo'], as_index=False)['InvoiceDate_int'].mean()
df_cleaned.drop('InvoiceDate_int', axis = 1, inplace = True)
basket_price.loc[:, 'InvoiceDate'] = pd.to_datetime(temp['InvoiceDate_int'])
#
# selection des entrées significatives:
basket_price = basket_price[basket_price['Basket Price'] > 0]
basket_price.sort_values('CustomerID')[:6]
```

Out[38]:

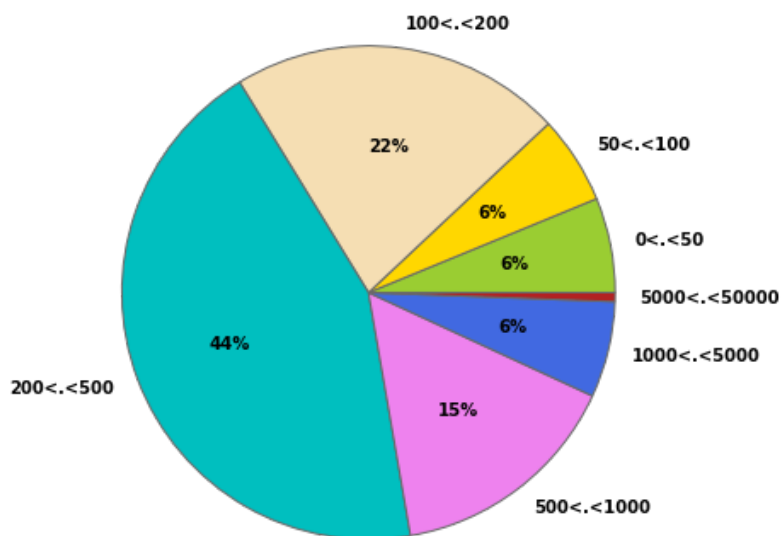
	CustomerID	InvoiceNo	Basket Price	InvoiceDate
1	12347	537626	711.79	2010-12-07 14:57:00.000001024
2	12347	542237	475.39	2011-01-26 14:29:59.999999744
3	12347	549222	636.25	2011-04-07 10:42:59.999999232
4	12347	556201	382.52	2011-06-09 13:01:00.000000256
5	12347	562032	584.91	2011-08-02 08:48:00.000000000
6	12347	573511	1294.32	2011-10-31 12:25:00.000001280


```

In [39]: #
# Décompte des achats
price_range = [0, 50, 100, 200, 500, 1000, 5000, 50000]
count_price = []
for i, price in enumerate(price_range):
    if i == 0: continue
    val = basket_price[(basket_price['Basket Price'] < price) &
                      (basket_price['Basket Price'] > price_range[i-1])]['Basket Price'].count()
    count_price.append(val)
#
# Représentation du nombre d'achats / montant
plt.rc('font', weight='bold')
f, ax = plt.subplots(figsize=(11, 6))
colors = ['yellowgreen', 'gold', 'wheat', 'c', 'violet', 'royalblue', 'firebrick']
labels = [ '{<.<}'.format(price_range[i-1], s) for i,s in enumerate(price_range) if i != 0]
sizes = count_price
explode = [0.0 if sizes[i] < 100 else 0.0 for i in range(len(sizes))]
ax.pie(sizes, explode = explode, labels=labels, colors = colors,
      autopct = lambda x: '{:1.0f}%'.format(x) if x > 1 else '',
      shadow = False, startangle=0)
ax.axis('equal')
f.text(0.5, 1.01, "Répartition des montants des commandes", ha='center', fontsize = 18);

```

Répartition des montants des commandes



```

In [40]: is_noun = lambda pos: pos[:2] == 'NN'

def keywords_inventory(dataframe, colonne = 'Description'):
    stemmer = nltk.stem.SnowballStemmer("english")
    keywords_roots = dict() # collect the words / root
    keywords_select = dict() # association: root <-> keyword
    category_keys = []
    count_keywords = dict()
    icount = 0
    for s in dataframe[colonne]:
        if pd.isnull(s): continue
        lines = s.lower()
        tokenized = nltk.word_tokenize(lines)
        nouns = [word for (word, pos) in nltk.pos_tag(tokenized) if is_noun(pos)]

        for t in nouns:
            t = t.lower() ; racine = stemmer.stem(t)
            if racine in keywords_roots:
                keywords_roots[racine].add(t)
                count_keywords[racine] += 1
            else:
                keywords_roots[racine] = {t}
                count_keywords[racine] = 1

    for s in keywords_roots.keys():
        if len(keywords_roots[s]) > 1:
            min_length = 1000
            for k in keywords_roots[s]:
                if len(k) < min_length:
                    clef = k ; min_length = len(k)
            category_keys.append(clef)
            keywords_select[s] = clef
        else:
            category_keys.append(list(keywords_roots[s])[0])
            keywords_select[s] = list(keywords_roots[s])[0]

    print("Nb of keywords in variable '{}': {}".format(colonne, len(category_keys)))
    return category_keys, keywords_roots, keywords_select, count_keywords

```

```

In [41]: df_produits = pd.DataFrame(df_initial['Description'].unique()).rename(columns = {0: 'Description'})

```

```

In [43]: keywords, keywords_roots, keywords_select, count_keywords = keywords_inventory(df_produits)

```

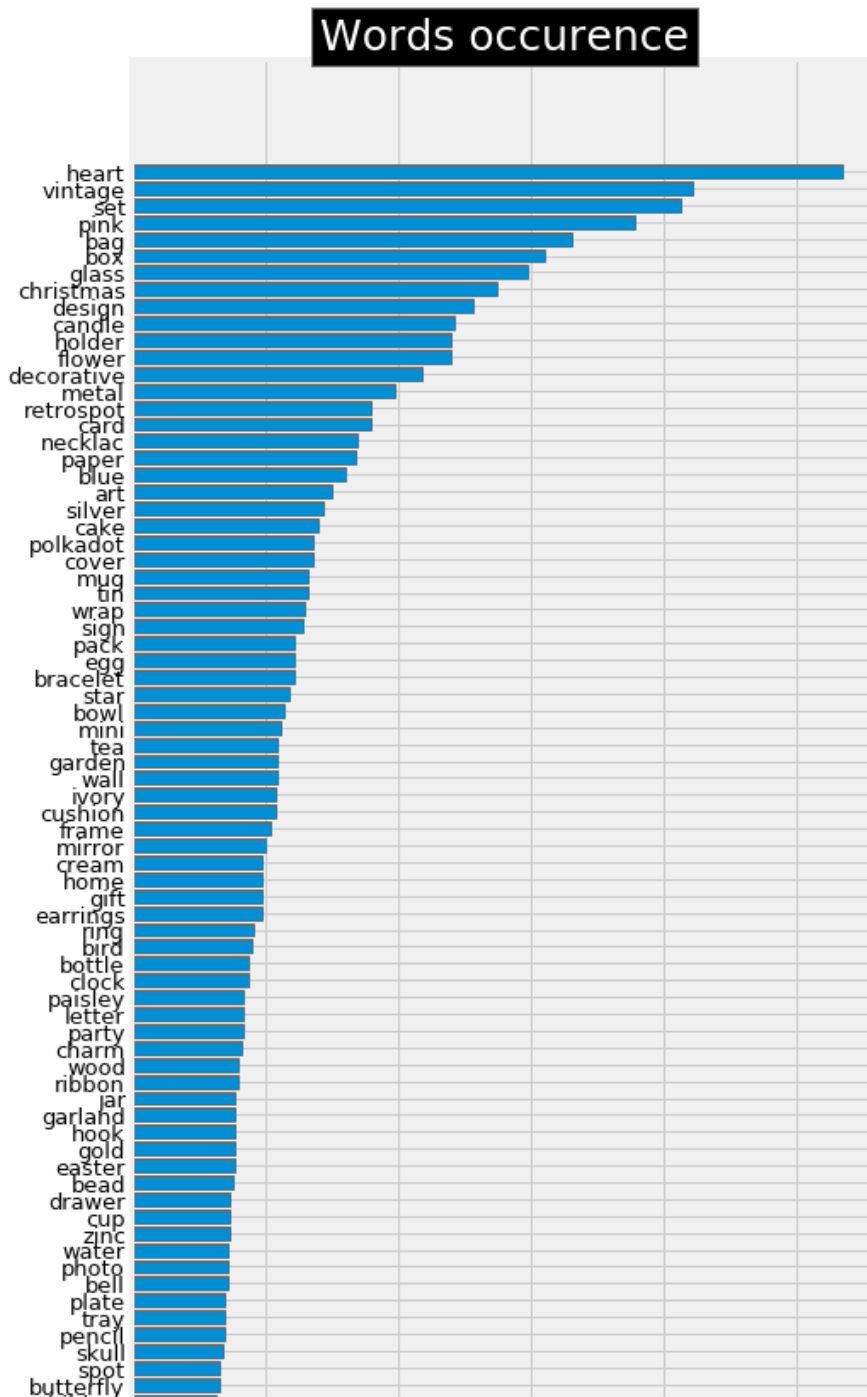
Nb of keywords in variable 'Description': 1484

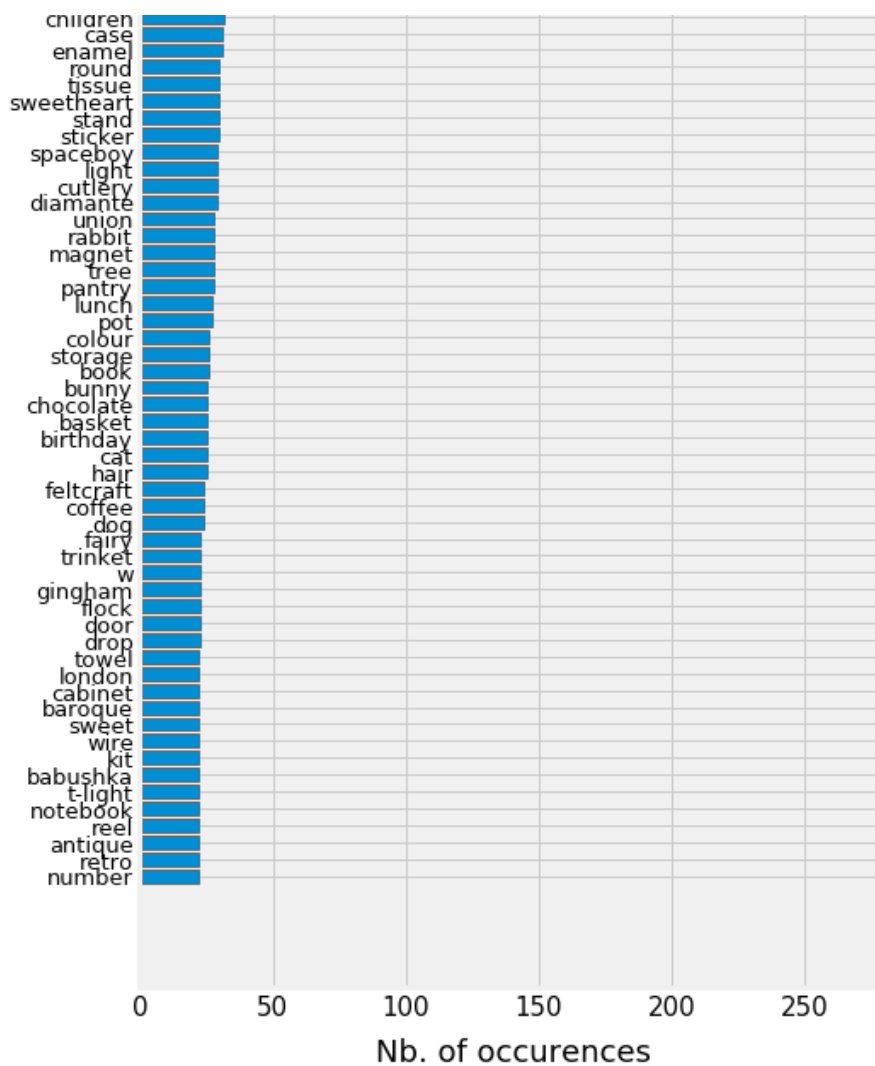
```

In [44]: list_products = []
for k,v in count_keywords.items():
    list_products.append([keywords_select[k],v])
list_products.sort(key = lambda x:x[1], reverse = True)

```

```
In [45]: liste = sorted(list_products, key = lambda x:x[1], reverse = True)
#
plt.rc('font', weight='normal')
fig, ax = plt.subplots(figsize=(7, 25))
y_axis = [i[1] for i in liste[:125]]
x_axis = [k for k,i in enumerate(liste[:125])]
x_label = [i[0] for i in liste[:125]]
plt.xticks(fontsize = 15)
plt.yticks(fontsize = 13)
plt.yticks(x_axis, x_label)
plt.xlabel("Nb. of occurrences", fontsize = 18, labelpad = 10)
ax.barh(x_axis, y_axis, align = 'center')
ax = plt.gca()
ax.invert_yaxis()
#
plt.title("Words occurrence",bbox={'facecolor':'k', 'pad':5}, color='w',fontsize = 25)
plt.show()
```





```
In [46]: list_products = []
for k,v in count_keywords.items():
    word = keywords_select[k]
    if word in ['pink', 'blue', 'tag', 'green', 'orange']: continue
    if len(word) < 3 or v < 13: continue
    if ('+' in word) or ('/' in word): continue
    list_products.append([word, v])
#
list_products.sort(key = lambda x:x[1], reverse = True)
print('mots conservés:', len(list_products))
```

mots conservés: 193

```
In [47]: liste_products = df_cleaned['Description'].unique()
X = pd.DataFrame()
for key, occurrence in list_products:
    X.loc[:, key] = list(map(lambda x:int(key.upper() in x), liste_products))
```

```
In [48]: threshold = [0, 1, 2, 3, 5, 10]
label_col = []
for i in range(len(threshold)):
    if i == len(threshold)-1:
        col = '}>{}'.format(threshold[i])
    else:
        col = '{}<.<{}'.format(threshold[i], threshold[i+1])
    label_col.append(col)
    X.loc[:, col] = 0

for i, prod in enumerate(liste_produits):
    prix = df_cleaned[ df_cleaned['Description'] == prod]['UnitPrice'].mean()
    j = 0
    while prix > threshold[j]:
        j+=1
        if j == len(threshold): break
    X.loc[i, label_col[j-1]] = 1
```

```
In [49]: matrix = X.as_matrix()
for n_clusters in range(3,10):
    kmeans = KMeans(init='k-means++', n_clusters = n_clusters, n_init=30)
    kmeans.fit(matrix)
    clusters = kmeans.predict(matrix)
    silhouette_avg = silhouette_score(matrix, clusters)
    print("For n_clusters =", n_clusters, "The average silhouette_score is :", silhouette_avg)
```

```
For n_clusters = 3 The average silhouette_score is : 0.10071681758064248
For n_clusters = 4 The average silhouette_score is : 0.12609893747265383
For n_clusters = 5 The average silhouette_score is : 0.12666226252150753
For n_clusters = 6 The average silhouette_score is : 0.14544497724586775
For n_clusters = 7 The average silhouette_score is : 0.15071006331939618
For n_clusters = 8 The average silhouette_score is : 0.15051164969896913
For n_clusters = 9 The average silhouette_score is : 0.1607323479959968
```

```
In [50]: n_clusters = 5
silhouette_avg = -1
while silhouette_avg < 0.145:
    kmeans = KMeans(init='k-means++', n_clusters = n_clusters, n_init=30)
    kmeans.fit(matrix)
    clusters = kmeans.predict(matrix)
    silhouette_avg = silhouette_score(matrix, clusters)

    #km = kmodes.KModes(n_clusters = n_clusters, init='Huang', n_init=2, verbose=0)
    #clusters = km.fit_predict(matrix)
    #silhouette_avg = silhouette_score(matrix, clusters)
    print("For n_clusters =", n_clusters, "The average silhouette_score is :", silhouette_avg)
```

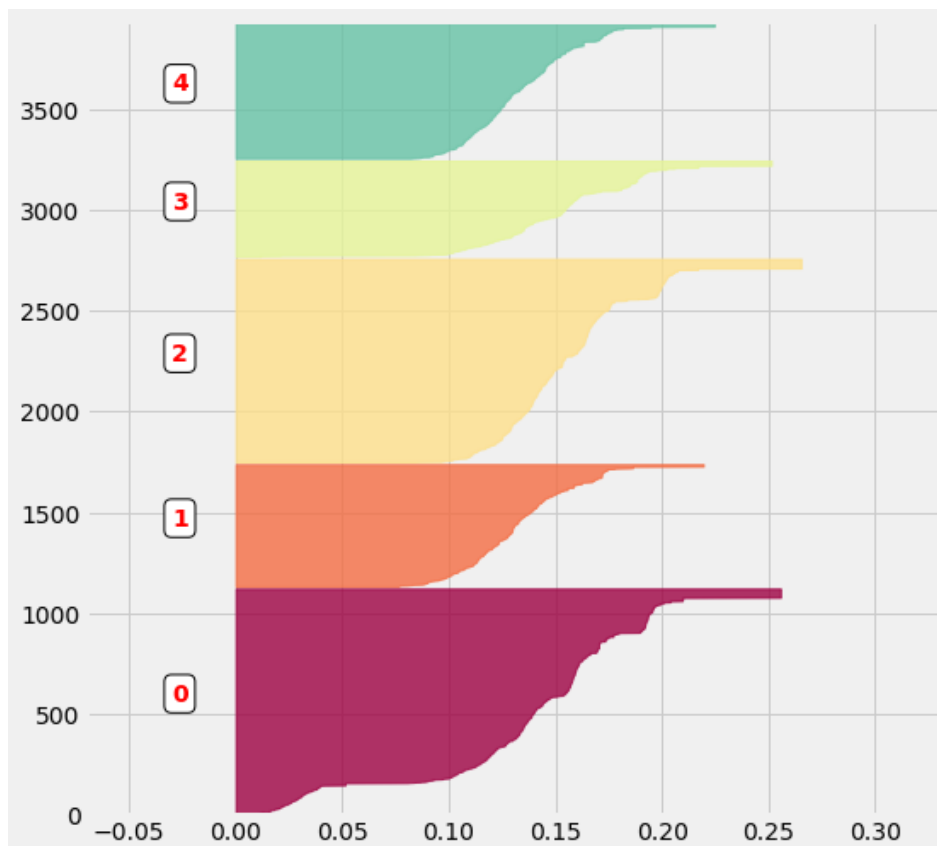
```
For n_clusters = 5 The average silhouette_score is : 0.13087995415351875
For n_clusters = 5 The average silhouette_score is : 0.1454649521746832
```

```
In [51]: pd.Series(clusters).value_counts()
```

```
Out[51]: 0    1114
         2    1009
         4     673
         1     606
         3     476
dtype: int64
```

```
In [52]: def graph_component_silhouette(n_clusters, lim_x, mat_size, sample_silhouette_values, clusters):
plt.rcParams["patch.force_edgecolor"] = True
plt.style.use('fivethirtyeight')
mpl.rc('patch', edgecolor = 'dimgray', linewidth=1)
#
fig, ax1 = plt.subplots(1, 1)
fig.set_size_inches(8, 8)
ax1.set_xlim([lim_x[0], lim_x[1]])
ax1.set_ylim([0, mat_size + (n_clusters + 1) * 10])
y_lower = 10
for i in range(n_clusters):
#
# Aggregate the silhouette scores for samples belonging to cluster i, and sort them
ith_cluster_silhouette_values = sample_silhouette_values[clusters == i]
ith_cluster_silhouette_values.sort()
size_cluster_i = ith_cluster_silhouette_values.shape[0]
y_upper = y_lower + size_cluster_i
cmap = cm.get_cmap("Spectral")
color = cmap(float(i) / n_clusters)
ax1.fill_betweenx(np.arange(y_lower, y_upper), 0, ith_cluster_silhouette_values,
facecolor=color, edgecolor=color, alpha=0.8)
#
# Label the silhouette plots with their cluster numbers at the middle
ax1.text(-0.03, y_lower + 0.5 * size_cluster_i, str(i), color = 'red', fontweight = 'bold',
bbox=dict(facecolor='white', edgecolor='black', boxstyle='round', pad=0.3))
#
# Compute the new y_lower for next plot
y_lower = y_upper + 10
```

```
In [53]: # define individual silhouette scores
sample_silhouette_values = silhouette_samples(matrix, clusters)
#
# and do the graph
graph_component_silhouette(n_clusters, [-0.07, 0.33], len(X), sample_silhouette_values, clusters)
```



```
In [54]: liste = pd.DataFrame(liste_produits)
liste_words = [word for (word, occurrence) in list_products]

occurrence = [dict() for _ in range(n_clusters)]

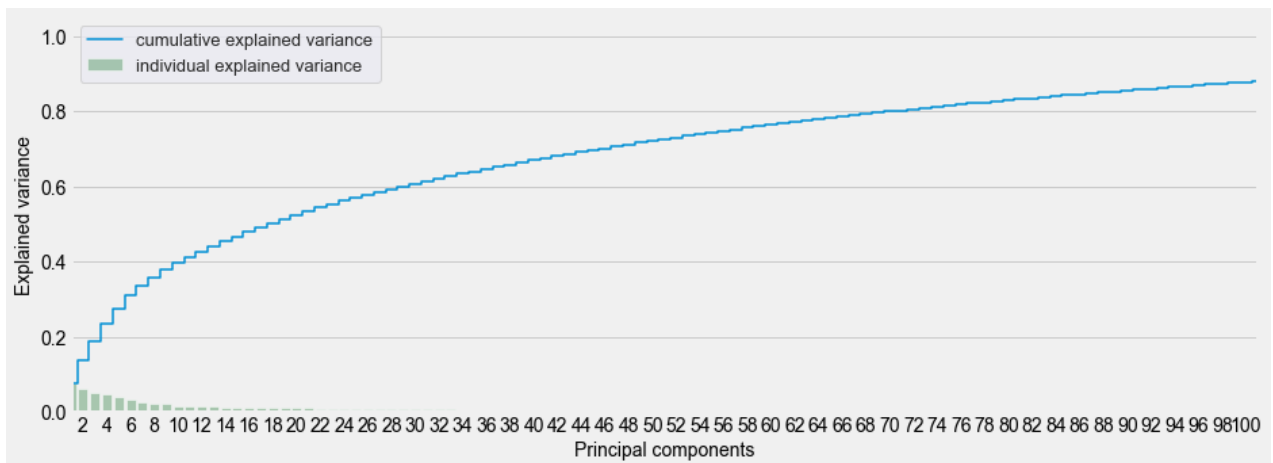
for i in range(n_clusters):
    liste_cluster = liste.loc[clusters == i]
    for word in liste_words:
        if word in ['art', 'set', 'heart', 'pink', 'blue', 'tag']: continue
        occurrence[i][word] = sum(liste_cluster.loc[:, 0].str.contains(word.upper()))
```

```
In [56]: pca = PCA()
pca.fit(matrix)
pca_samples = pca.transform(matrix)
```

```
In [57]: fig, ax = plt.subplots(figsize=(14, 5))
sns.set(font_scale=1)
plt.step(range(matrix.shape[1]), pca.explained_variance_ratio_.cumsum(), where='mid',
        label='cumulative explained variance')
sns.barplot(np.arange(1, matrix.shape[1]+1), pca.explained_variance_ratio_, alpha=0.5, color = 'g',
        label='individual explained variance')
plt.xlim(0, 100)

ax.set_xticklabels([s if int(s.get_text())%2 == 0 else '' for s in ax.get_xticklabels()])

plt.ylabel('Explained variance', fontsize = 14)
plt.xlabel('Principal components', fontsize = 14)
plt.legend(loc='upper left', fontsize = 13);
```



```
In [58]: pca = PCA(n_components=50)
matrix_9D = pca.fit_transform(matrix)
mat = pd.DataFrame(matrix_9D)
mat['cluster'] = pd.Series(clusters)
```

```

In [59]: import matplotlib.patches as mpatches

sns.set_style("white")
sns.set_context("notebook", font_scale=1, rc={"lines.linewidth": 2.5})

LABEL_COLOR_MAP = {0:'r', 1:'gold', 2:'b', 3:'k', 4:'c', 5:'g'}
label_color = [LABEL_COLOR_MAP[l] for l in mat['cluster']]

fig = plt.figure(figsize = (15,8))
increment = 0
for ix in range(4):
    for iy in range(ix+1, 4):
        increment += 1
        ax = fig.add_subplot(2,3,increment)
        ax.scatter(mat[ix], mat[iy], c= label_color, alpha=0.4)
        plt.ylabel('PCA {}'.format(iy+1), fontsize = 12)
        plt.xlabel('PCA {}'.format(ix+1), fontsize = 12)
        ax.yaxis.grid(color='lightgray', linestyle=':')
        ax.xaxis.grid(color='lightgray', linestyle=':')
        ax.spines['right'].set_visible(False)
        ax.spines['top'].set_visible(False)

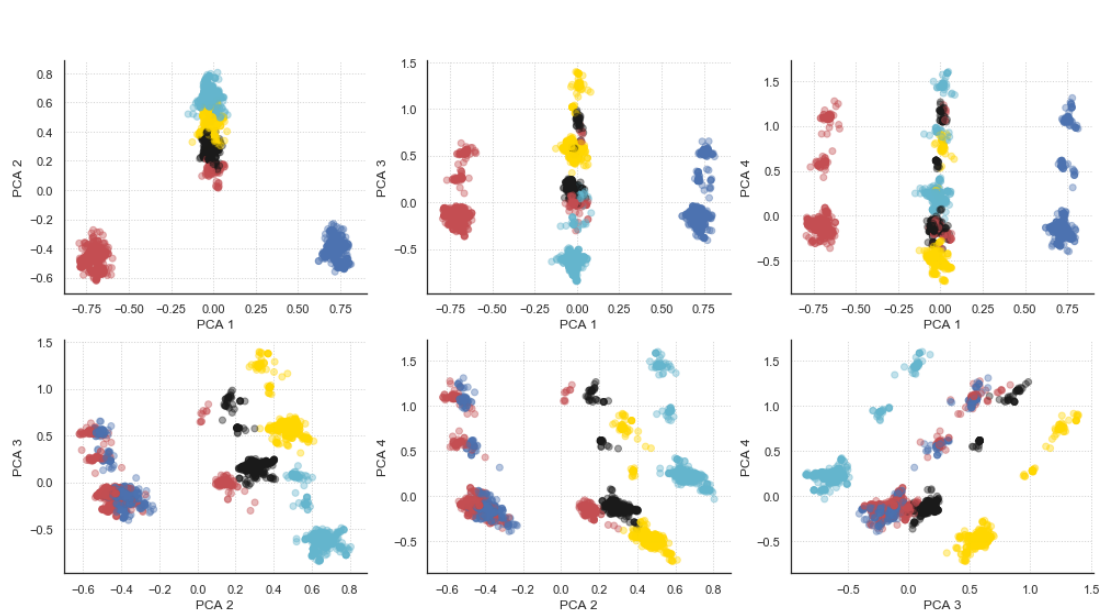
        if increment == 9: break
    if increment == 9: break

#
# I set the Legend: abbreviation -> airline name
comp_handler = []
for i in range(5):
    comp_handler.append(mpatches.Patch(color = LABEL_COLOR_MAP[i], label = i))

plt.legend(handles=comp_handler, bbox_to_anchor=(1.1, 0.97),
          title='Cluster', facecolor = 'lightgrey',
          shadow = True, frameon = True, framealpha = 1,
          fontsize = 13, bbox_transform = plt.gcf().transFigure)

plt.show()

```




```
In [60]: corres = dict()
for key, val in zip(liste_produits, clusters):
    corres[key] = val
#
df_cleaned['categ_product'] = df_cleaned.loc[:, 'Description'].map(corres)
```

```
In [61]: for i in range(5):
    col = 'categ_{}'.format(i)
    df_temp = df_cleaned[df_cleaned['categ_product'] == i]
    price_temp = df_temp['UnitPrice'] * (df_temp['Quantity'] - df_temp['QuantityCanceled'])
    price_temp = price_temp.apply(lambda x: x if x > 0 else 0)
    df_cleaned.loc[:, col] = price_temp
    df_cleaned[col].fillna(0, inplace = True)
#
df_cleaned[['InvoiceNo', 'Description', 'categ_product', 'categ_0', 'categ_1', 'categ_2', 'categ_3', 'categ_4']]
```

Out[61]:

	InvoiceNo	Description	categ_product	categ_0	categ_1	categ_2	categ_3	categ_4
0	536365	WHITE HANGING HEART T-LIGHT HOLDER	4	0.0	0.00	0.0	0.0	15.3
1	536365	WHITE METAL LANTERN	1	0.0	20.34	0.0	0.0	0.0
2	536365	CREAM CUPID HEARTS COAT HANGER	1	0.0	22.00	0.0	0.0	0.0
3	536365	KNITTED UNION FLAG HOT WATER BOTTLE	1	0.0	20.34	0.0	0.0	0.0
4	536365	RED WOOLLY HOTTIE WHITE HEART.	1	0.0	20.34	0.0	0.0	0.0

```
In [62]: #
# somme des achats / utilisateur & commande
temp = df_cleaned.groupby(by=['CustomerID', 'InvoiceNo'], as_index=False)['TotalPrice'].sum()
basket_price = temp.rename(columns = {'TotalPrice': 'Basket Price'})
#
# pourcentage du prix de la commande / categorie de produit
for i in range(5):
    col = 'categ_{}'.format(i)
    temp = df_cleaned.groupby(by=['CustomerID', 'InvoiceNo'], as_index=False)[col].sum()
    basket_price.loc[:, col] = temp
#
# date de la commande
df_cleaned['InvoiceDate_int'] = df_cleaned['InvoiceDate'].astype('int64')
temp = df_cleaned.groupby(by=['CustomerID', 'InvoiceNo'], as_index=False)['InvoiceDate_int'].mean()
df_cleaned.drop('InvoiceDate_int', axis = 1, inplace = True)
basket_price.loc[:, 'InvoiceDate'] = pd.to_datetime(temp['InvoiceDate_int'])
#
# selection des entrées significatives:
basket_price = basket_price[basket_price['Basket Price'] > 0]
basket_price.sort_values('CustomerID', ascending = True)[:5]
```

Out[62]:

	CustomerID	InvoiceNo	Basket Price	categ_0	categ_1	categ_2	categ_3	categ_4	InvoiceDate
1	12347	537626	711.79	23.40	293.35	187.2	124.44	83.40	2010-12-07 14:57:00.000001024
2	12347	542237	475.39	122.59	169.20	130.5	0.00	53.10	2011-01-26 14:29:59.999999744
3	12347	549222	636.25	119.25	115.00	330.9	0.00	71.10	2011-04-07 10:42:59.999999232
4	12347	556201	382.52	41.40	168.76	74.4	19.90	78.06	2011-06-09 13:01:00.000000256
5	12347	562032	584.91	99.55	158.16	109.7	97.80	119.70	2011-08-02 08:48:00.000000000

```
In [63]: print(basket_price['InvoiceDate'].min(), '->', basket_price['InvoiceDate'].max())
```

2010-12-01 08:26:00 -> 2011-12-09 12:50:00

```
In [64]: #
# nb de visites et stats sur Le montant du panier / utilisateurs
transactions_per_user=basket_price.groupby(by=['CustomerID'])['Basket Price'].agg(['count','min',
for i in range(5):
    col = 'categ_{}'.format(i)
    transactions_per_user.loc[:,col] = basket_price.groupby(by=['CustomerID'])[col].sum() /\
        transactions_per_user['sum']*100

transactions_per_user.reset_index(drop = False, inplace = True)
basket_price.groupby(by=['CustomerID'])['categ_0'].sum()
transactions_per_user.sort_values('CustomerID', ascending = True)[:5]
```

Out[64]:

	CustomerID	count	min	max	mean	sum	categ_0	categ_1	categ_2	categ_3	categ_4
0	12347	7	224.82	1294.32	615.714286	4310.00	15.674478	29.540371	26.375870	7.604176	20.805104
1	12348	4	227.44	892.80	449.310000	1797.24	58.046783	0.000000	41.953217	0.000000	0.000000
2	12349	1	1757.55	1757.55	1757.550000	1757.55	30.145373	10.713778	26.506216	20.389178	12.245455
3	12350	1	334.40	334.40	334.400000	334.40	23.654306	0.000000	48.444976	0.000000	27.900718
4	12352	7	144.35	840.30	340.815714	2385.71	50.930331	14.601523	15.705178	14.691643	4.071325

```
In [65]: last_date = basket_price['InvoiceDate'].max().date()

first_registration = pd.DataFrame(basket_price.groupby(by=['CustomerID'])['InvoiceDate'].min())
last_purchase      = pd.DataFrame(basket_price.groupby(by=['CustomerID'])['InvoiceDate'].max())

test = first_registration.applymap(lambda x:(last_date - x.date()).days)
test2 = last_purchase.applymap(lambda x:(last_date - x.date()).days)

transactions_per_user.loc[:, 'LastPurchase'] = test2.reset_index(drop = False)['InvoiceDate']
transactions_per_user.loc[:, 'FirstPurchase'] = test.reset_index(drop = False)['InvoiceDate']

transactions_per_user[:5]
```

Out[65]:

	CustomerID	count	min	max	mean	sum	categ_0	categ_1	categ_2	categ_3	categ_4
0	12347	7	224.82	1294.32	615.714286	4310.00	15.674478	29.540371	26.375870	7.604176	20.805104
1	12348	4	227.44	892.80	449.310000	1797.24	58.046783	0.000000	41.953217	0.000000	0.000000
2	12349	1	1757.55	1757.55	1757.550000	1757.55	30.145373	10.713778	26.506216	20.389178	12.245455
3	12350	1	334.40	334.40	334.400000	334.40	23.654306	0.000000	48.444976	0.000000	27.900718
4	12352	7	144.35	840.30	340.815714	2385.71	50.930331	14.601523	15.705178	14.691643	4.071325

```
In [66]: n1 = transactions_per_user[transactions_per_user['count'] == 1].shape[0]
n2 = transactions_per_user.shape[0]
print("nb. de clients avec achat unique: {:<2}/{:<5} ({:<2.2f}%)".format(n1,n2,n1/n2*100))
```

nb. de clients avec achat unique: 1489/4327 (34.41%)

```
In [67]: list_cols = ['count','min','max','mean','categ_0','categ_1','categ_2','categ_3','categ_4']
#
selected_customers = transactions_per_user.copy(deep = True)
matrix = selected_customers[list_cols].as_matrix()
```

```
In [68]: scaler = StandardScaler()
scaler.fit(matrix)
print('variables mean values: \n' + 90*'- ' + '\n' , scaler.mean_)
scaled_matrix = scaler.transform(matrix)
```

variables mean values:

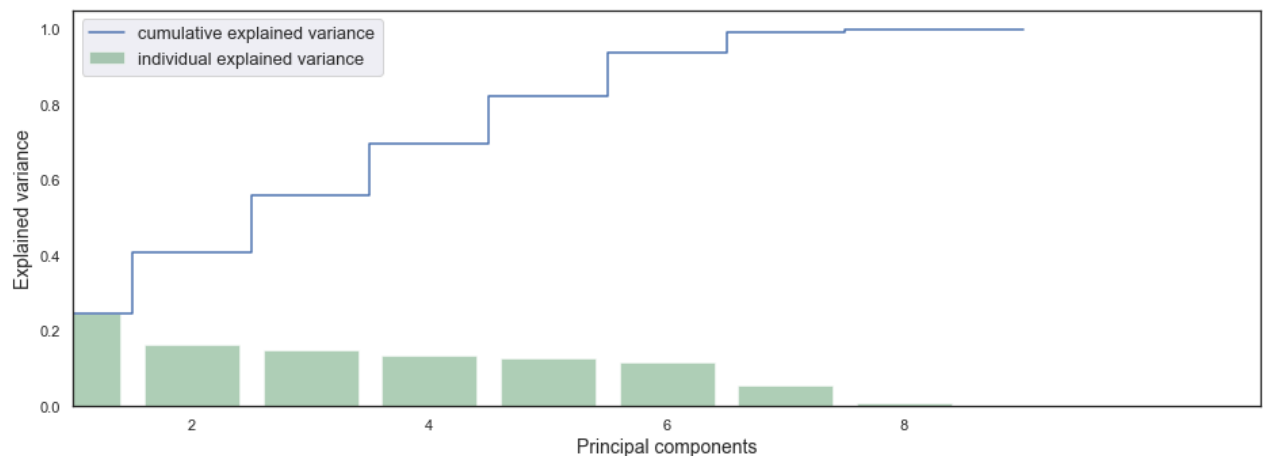
```
-----
[  4.25190663 241.38253571 578.87676959 372.25705475 21.29203891
 16.97264402 25.30761412 14.80638872 21.62911535]
```

```
In [69]: pca = PCA()
pca.fit(scaled_matrix)
pca_samples = pca.transform(scaled_matrix)
```

```
In [70]: fig, ax = plt.subplots(figsize=(14, 5))
sns.set(font_scale=1)
plt.step(range(matrix.shape[1]), pca.explained_variance_ratio_.cumsum(), where='mid',
        label='cumulative explained variance')
sns.barplot(np.arange(1,matrix.shape[1]+1), pca.explained_variance_ratio_, alpha=0.5, color = 'g'
        label='individual explained variance')
plt.xlim(0, 10)

ax.set_xticklabels([s if int(s.get_text())%2 == 0 else '' for s in ax.get_xticklabels()])

plt.ylabel('Explained variance', fontsize = 14)
plt.xlabel('Principal components', fontsize = 14)
plt.legend(loc='best', fontsize = 13);
```



```
In [71]: n_clusters = 11
kmeans = KMeans(init='k-means++', n_clusters = n_clusters, n_init=100)
kmeans.fit(scaled_matrix)
clusters_clients = kmeans.predict(scaled_matrix)
silhouette_avg = silhouette_score(scaled_matrix, clusters_clients)
print('score de silhouette: {:.3f}'.format(silhouette_avg))
```

score de silhouette: 0.219

```
In [72]: pd.DataFrame(pd.Series(clusters_clients).value_counts(), columns = ['nb. de clients']).T
```

Out[72]:

	1	8	0	9	5	2	4	3	7	6	10
nb. de clients	1904	600	421	370	313	303	197	183	17	12	7

```
In [73]: pca = PCA(n_components=6)
matrix_3D = pca.fit_transform(scaled_matrix)
mat = pd.DataFrame(matrix_3D)
mat['cluster'] = pd.Series(clusters_clients)
```

```

In [74]: import matplotlib.patches as mpatches

sns.set_style("white")
sns.set_context("notebook", font_scale=1, rc={"lines.linewidth": 2.5})

LABEL_COLOR_MAP = {0:'r', 1:'tan', 2:'b', 3:'k', 4:'c', 5:'g', 6:'deeppink', 7:'skyblue', 8:'dark
                  10:'yellow', 11:'tomato', 12:'seagreen'}
label_color = [LABEL_COLOR_MAP[l] for l in mat['cluster']]

fig = plt.figure(figsize = (12,10))
increment = 0
for ix in range(6):
    for iy in range(ix+1, 6):
        increment += 1
        ax = fig.add_subplot(4,3,increment)
        ax.scatter(mat[ix], mat[iy], c= label_color, alpha=0.5)
        plt.ylabel('PCA {}'.format(iy+1), fontsize = 12)
        plt.xlabel('PCA {}'.format(ix+1), fontsize = 12)
        ax.yaxis.grid(color='lightgray', linestyle=':')
        ax.xaxis.grid(color='lightgray', linestyle=':')
        ax.spines['right'].set_visible(False)
        ax.spines['top'].set_visible(False)

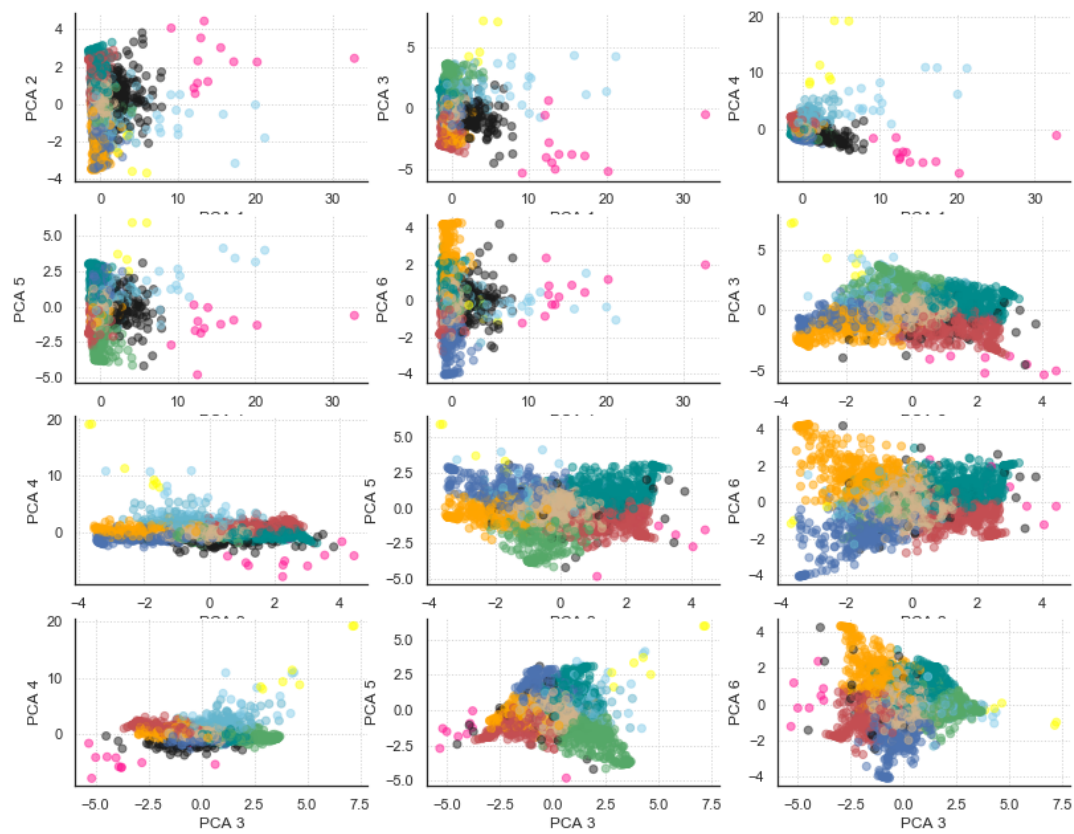
        if increment == 12: break
    if increment == 12: break

#
# I set the legend: abbreviation -> airline name
comp_handler = []
for i in range(n_clusters):
    comp_handler.append(mpatches.Patch(color = LABEL_COLOR_MAP[i], label = i))

plt.legend(handles=comp_handler, bbox_to_anchor=(1.1, 0.9),
           title='Cluster', facecolor = 'lightgrey',
           shadow = True, frameon = True, framealpha = 1,
           fontsize = 13, bbox_transform = plt.gcf().transFigure)

plt.tight_layout()

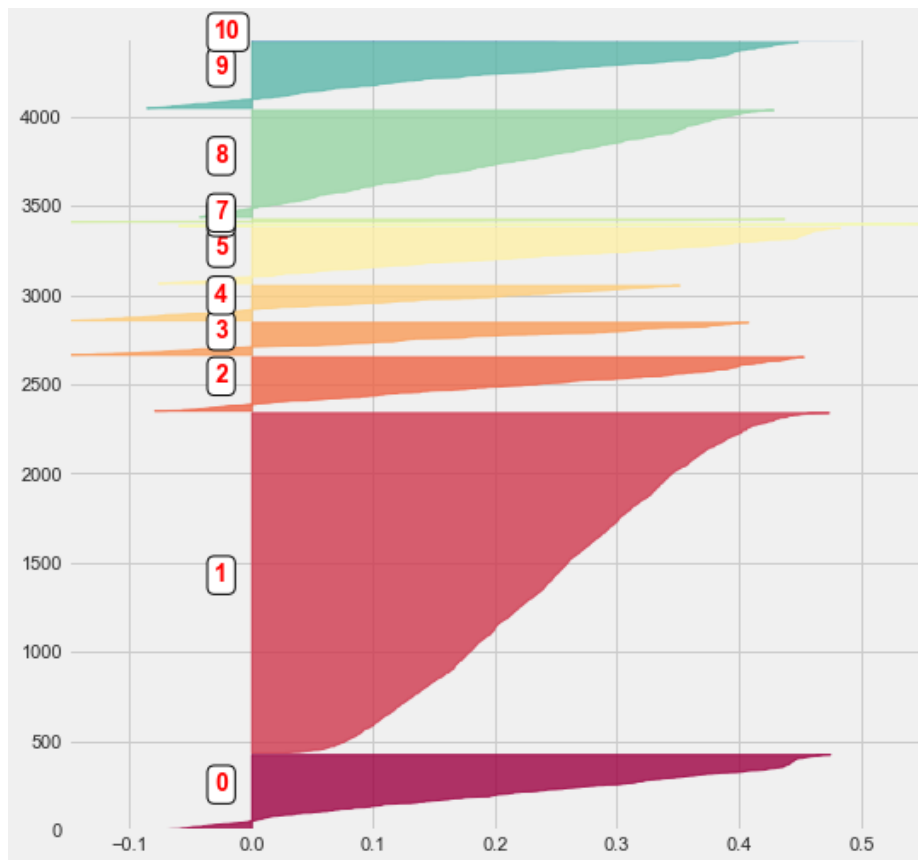
```



```
In [75]: silhouette_values = silhouette_samples(scaled_matrix, clusters_clients)

# Individual silhouette scores
silhouette_values = silhouette_samples(scaled_matrix, clusters_clients)

# Plot
plot_silhouette(n_clusters, [-0.15, 0.55], len(scaled_matrix), sample_silhouette_values, clusters_clients)
```



```
In [76]: selected_customers.loc[:, 'cluster'] = clusters_clients
```

```
In [77]: merged_df = pd.DataFrame()
for i in range(n_clusters):
    test = pd.DataFrame(selected_customers[selected_customers['cluster'] == i].mean())
    test = test.T.set_index('cluster', drop = True)
    test['size'] = selected_customers[selected_customers['cluster'] == i].shape[0]
    merged_df = pd.concat([merged_df, test])

#
merged_df.drop('CustomerID', axis = 1, inplace = True)
print('number of customers:', merged_df['size'].sum())

merged_df = merged_df.sort_values('sum')
```

number of customers: 4327

```
In [78]: liste_index = []
for i in range(5):
    column = 'categ_{}'.format(i)
    liste_index.append(merged_df[merged_df[column] > 45].index.values[0])
#
liste_index_reordered = liste_index
liste_index_reordered += [ s for s in merged_df.index if s not in liste_index]
#
merged_df = merged_df.reindex(index = liste_index_reordered)
merged_df = merged_df.reset_index(drop = False)
display(merged_df[['cluster', 'count', 'min', 'max', 'mean', 'sum', 'categ_0',
                    'categ_1', 'categ_2', 'categ_3', 'categ_4', 'size']])
```

	cluster	count	min	max	mean	sum	categ_0	categ_1	categ_2	categ_3
0	0.0	2.235154	201.864157	331.190071	252.998048	617.173990	60.174179	6.955456	15.421611	7.701101
1	2.0	2.656766	196.731188	358.192871	268.073645	867.094125	9.887365	53.785835	13.421984	9.563101
2	8.0	2.740000	206.865183	350.671400	273.144140	770.938285	17.112146	8.474484	54.890504	5.374101
3	9.0	2.948649	190.611595	335.951622	257.235481	789.610568	12.021964	12.804910	11.344624	51.600101
4	5.0	3.006390	197.446805	405.058019	295.163570	974.089169	9.770183	7.502067	13.644960	6.013101
5	1.0	3.684874	205.505663	490.470489	333.269256	1258.022449	19.907995	18.151790	24.920991	14.484101
6	3.0	2.027322	1028.421475	1501.049678	1250.289513	2699.014268	20.264346	18.925267	26.016947	13.750101
7	6.0	1.500000	4075.741667	6245.503333	5073.737778	7931.028333	29.221945	15.135140	18.218655	12.841101
8	4.0	21.076142	68.539746	1331.796447	473.278136	9857.763401	18.468181	17.378025	23.336639	15.741101
9	10.0	127.285714	10.585714	2248.087143	381.948905	49672.060000	21.329216	14.286531	25.288514	17.040101
10	7.0	31.058824	85.434118	13750.221765	3025.061266	90558.188235	16.268942	18.540635	20.021898	15.226101


```

In [79]: def _scale_data(data, ranges):
    (x1, x2) = ranges[0]
    d = data[0]
    return [(d - y1) / (y2 - y1) * (x2 - x1) + x1 for d, (y1, y2) in zip(data, ranges)]

class RadarChart():
    def __init__(self, fig, location, sizes, variables, ranges, n_ordinate_levels = 6):

        angles = np.arange(0, 360, 360./len(variables))

        ix, iy = location[:]; size_x, size_y = sizes[:]

        axes = [fig.add_axes([ix, iy, size_x, size_y], polar = True,
            label = "axes{}".format(i)) for i in range(len(variables))]

        _, text = axes[0].set_thetagrids(angles, labels = variables)

        for txt, angle in zip(text, angles):
            if angle > -1 and angle < 181:
                txt.set_rotation(angle - 90)
            else:
                txt.set_rotation(angle - 270)

        for ax in axes[1:]:
            ax.patch.set_visible(False)
            ax.xaxis.set_visible(False)
            ax.grid("off")

        for i, ax in enumerate(axes):
            grid = np.linspace(*ranges[i], num = n_ordinate_levels)
            grid_label = [""] + ["{:0f}.".format(x) for x in grid[1:-1]]
            ax.set_rgrids(grid, labels = grid_label, angle = angles[i])
            ax.set_ylim(*ranges[i])

        self.angle = np.deg2rad(np.r_[angles, angles[0]])
        self.ranges = ranges
        self.ax = axes[0]

    def plot(self, data, *args, **kw):
        sdata = _scale_data(data, self.ranges)
        self.ax.plot(self.angle, np.r_[sdata, sdata[0]], *args, **kw)

    def fill(self, data, *args, **kw):
        sdata = _scale_data(data, self.ranges)
        self.ax.fill(self.angle, np.r_[sdata, sdata[0]], *args, **kw)

    def legend(self, *args, **kw):
        self.ax.legend(*args, **kw)

    def title(self, title, *args, **kw):
        self.ax.text(0.9, 1, title, transform = self.ax.transAxes, *args, **kw)

```

```

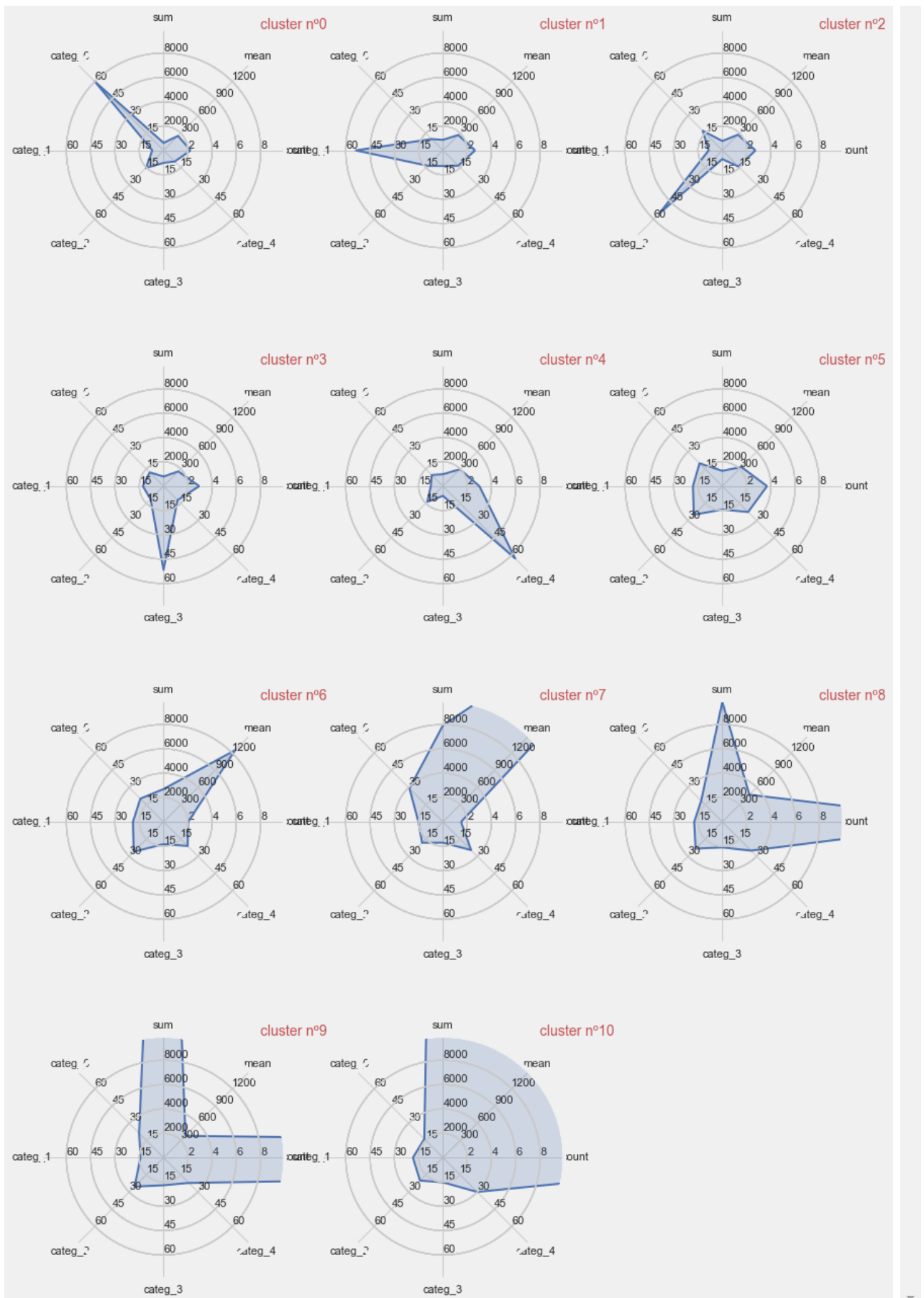
In [80]: fig = plt.figure(figsize=(10,12))

attributes = ['count', 'mean', 'sum', 'categ_0', 'categ_1', 'categ_2', 'categ_3', 'categ_4']
ranges = [[0.01, 10], [0.01, 1500], [0.01, 10000], [0.01, 75], [0.01, 75], [0.01, 75], [0.01, 75]]
index = [0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]

n_groups = n_clusters ; i_cols = 3
i_rows = n_groups//i_cols
size_x, size_y = (1/i_cols), (1/i_rows)

for ind in range(n_clusters):
    ix = ind%3 ; iy = i_rows - ind//3
    pos_x = ix*(size_x + 0.05) ; pos_y = iy*(size_y + 0.05)
    location = [pos_x, pos_y] ; sizes = [size_x, size_y]
    #
    data = np.array(merged_df.loc[index[ind], attributes])
    radar = RadarChart(fig, location, sizes, attributes, ranges)
    radar.plot(data, color = 'b', linewidth=2.0)
    radar.fill(data, alpha = 0.2, color = 'b')
    radar.title(title = 'cluster n°{}'.format(index[ind]), color = 'r')
    ind += 1

```



```
In [81]: 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)))
```

```
In [82]: columns = ['mean', 'categ_0', 'categ_1', 'categ_2', 'categ_3', 'categ_4' ]
X = selected_customers[columns]
Y = selected_customers['cluster']
```

```
In [83]: X_train, X_test, Y_train, Y_test = model_selection.train_test_split(X, Y, train_size = 0.8)
```

```
In [84]: svc = Class_Fit(clf = svm.LinearSVC)
svc.grid_search(parameters = [{ 'C': np.logspace(-2,2,10)}], Kfold = 5)
```

```
In [85]: svc.grid_fit(X = X_train, Y = Y_train)
```

```
In [86]: svc.grid_predict(X_test, Y_test)
```

Precision: 79.33 %

```

In [87]: 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)

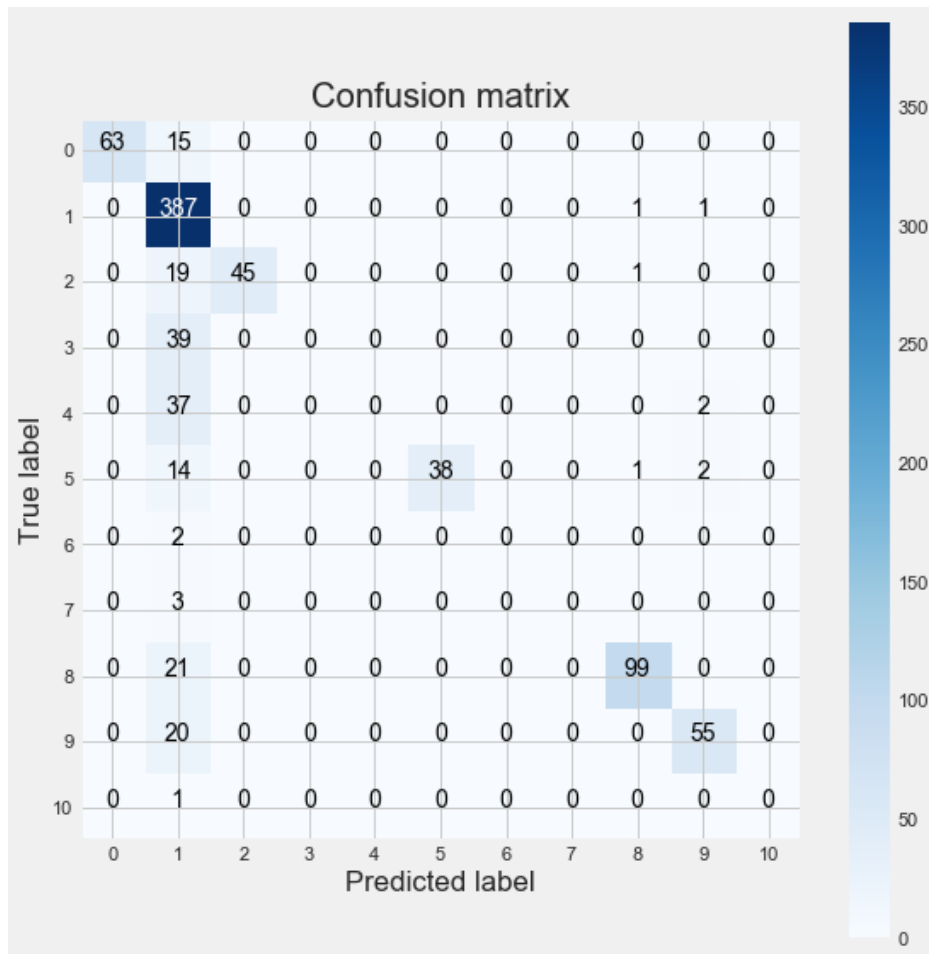
    #
    fmt = '.2f' if normalize else 'd'
    thresh = cm.max() / 2.
    for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
        plt.text(j, i, format(cm[i, j], fmt),
                 horizontalalignment="center",
                 color="white" if cm[i, j] > thresh else "black")

    #
    plt.tight_layout()
    plt.ylabel('True label')
    plt.xlabel('Predicted label')

```

```
In [88]: 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

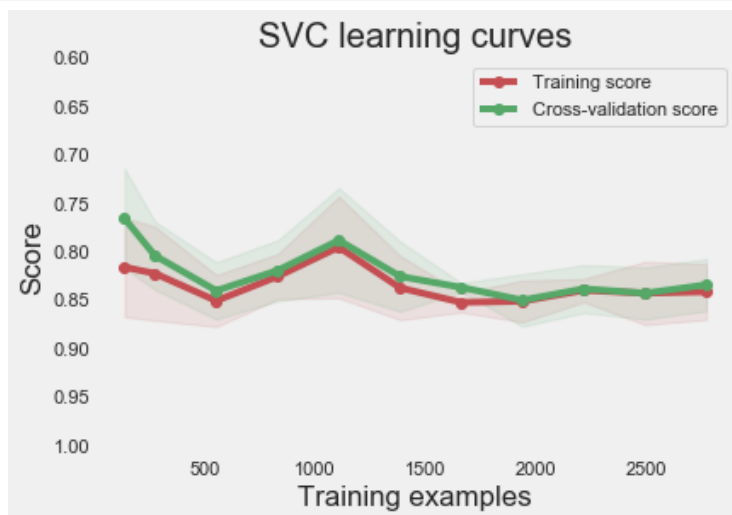


```
In [89]: def plot_learning_curve(estimator, title, X, y, ylim=None, cv=None,
                                n_jobs=-1, train_sizes=np.linspace(.1, 1.0, 10)):
    """Generate a simple plot of the test and training learning curve"""
    plt.figure()
    plt.title(title)
    if ylim is not None:
        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
```

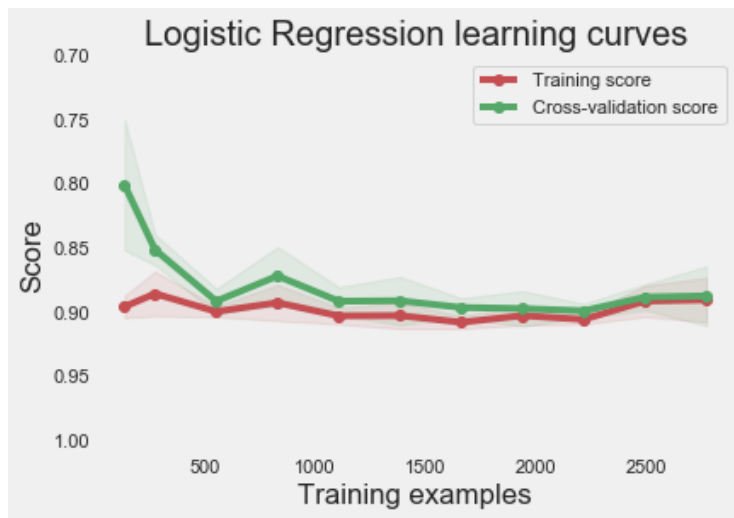
```
In [90]: g = plot_learning_curve(svc.grid.best_estimator_,
                                "SVC learning curves", X_train, Y_train, ylim = [1.01, 0.6],
                                cv = 5, train_sizes = [0.05, 0.1, 0.2, 0.3, 0.4, 0.5,
                                                         0.6, 0.7, 0.8, 0.9, 1])
```



```
In [91]: 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: 90.65 %

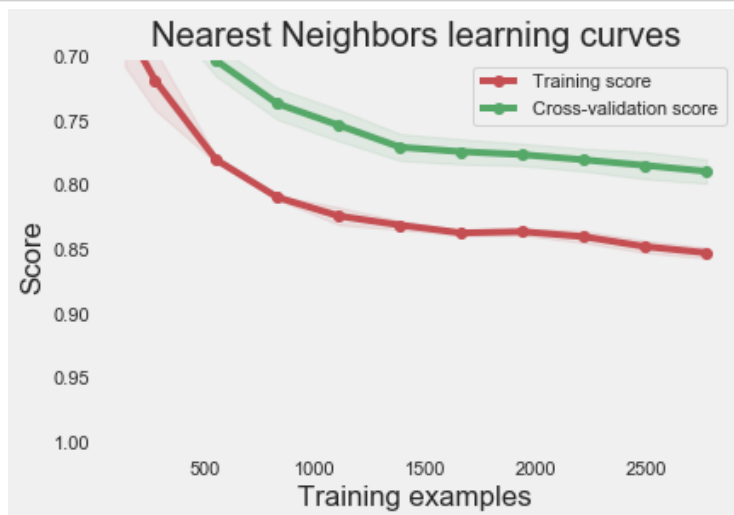
```
In [94]: g = plot_learning_curve(lr.grid.best_estimator_, "Logistic Regression learning curves", X_train,
                                ylim = [1.01, 0.7], cv = 5,
                                train_sizes = [0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1])
```



```
In [95]: 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: 80.72 %

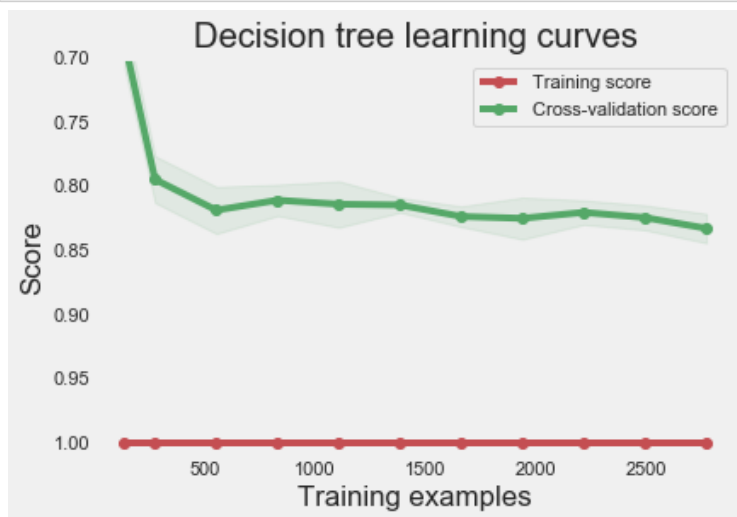
```
In [96]: g = plot_learning_curve(knn.grid.best_estimator_, "Nearest Neighbors learning curves", X_train, Y
                                ylim = [1.01, 0.7], cv = 5,
                                train_sizes = [0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1])
```




```
In [93]: tr = Class_Fit(clf = tree.DecisionTreeClassifier)
tr.grid_search(parameters = [{'criterion' : ['entropy', 'gini'], 'max_features' : ['sqrt', 'log2']}
tr.grid_fit(X = X_train, Y = Y_train)
tr.grid_predict(X_test, Y_test)
```

Precision: 81.99 %

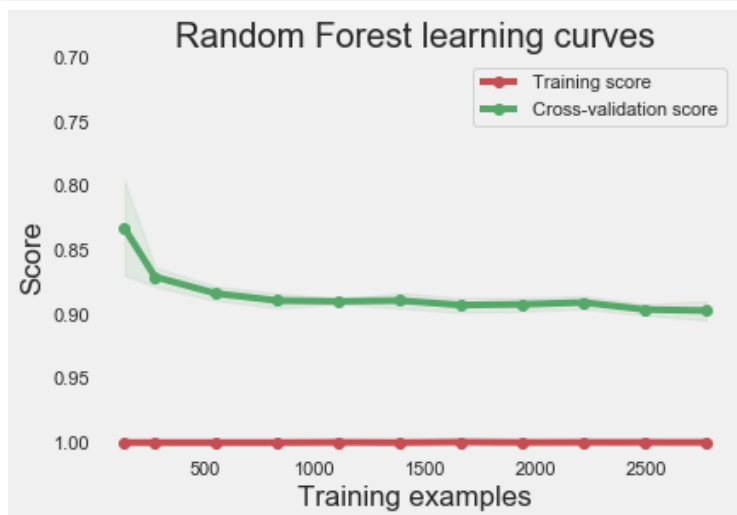
```
In [97]: g = plot_learning_curve(tr.grid.best_estimator_, "Decision tree learning curves", X_train, Y_train,
                                ylim = [1.01, 0.7], cv = 5,
                                train_sizes = [0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1])
```



```
In [98]: 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: 90.53 %

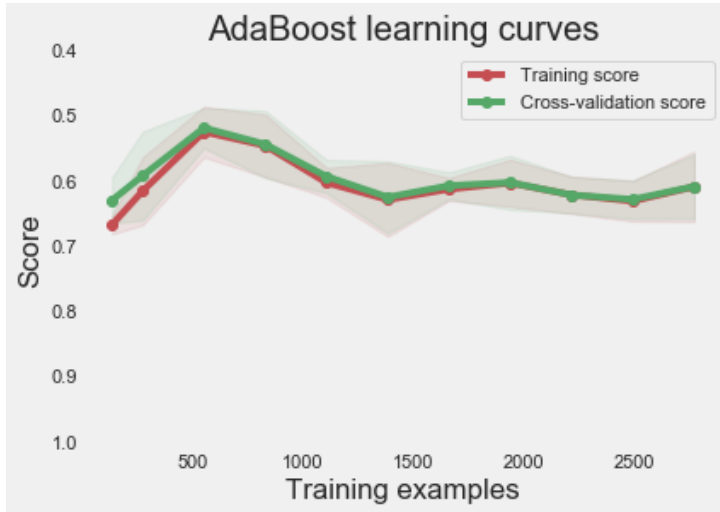
```
In [99]: g = plot_learning_curve(rf.grid.best_estimator_, "Random Forest learning curves", X_train, Y_train,
                                ylim = [1.01, 0.7], cv = 5,
                                train_sizes = [0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1])
```



```
In [100]: 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: 65.94 %

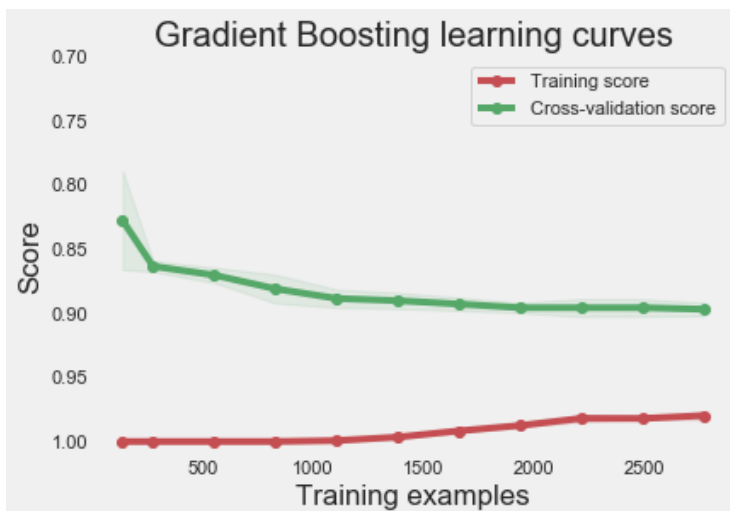
```
In [101]: g = plot_learning_curve(ada.grid.best_estimator_, "AdaBoost learning curves", X_train, Y_train,
                                ylim = [1.01, 0.4], cv = 5,
                                train_sizes = [0.05, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1])
```



```
In [102]: 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: 90.42 %

```
In [103]: 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])
```



*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.