



UNIVERSITY OF GRONINGEN

DATA ANALYSIS AND PROGRAMMING FOR OM

FINAL ASSIGNMENT

Warehouse planning for Gomibo

Authors:
Lucas Krol

Student Numbers:
S4072103

October 26, 2022

Contents

Introduction	2
Data	2
Mathematical model	5
Heuristics	5
Knapsack problem	5
Results	5
Sensitivity analysis	6
Appendix	8
Main file	8
DataFrameConstruction	19
CsvToElastic	22

Introduction

In this report we consider a warehouse planning problem for Gomibo (Belsimpel). Specifically, Gomibo's first warehouse is not large enough anymore to stock all products and therefore Gomibo has to store some products in a second warehouse, which however can not promise next-day delivery. Consequently, products that are stored in the second warehouse will be sold less often. More specifically, products are expected to lose either 20%, 30% or 50% of their sales, for products that belong, respectively, to the class of the 20% most profitable products, the next 30% or the bottom 50%.

Additionally, we have to take into account that Gomibo has products that are correlated with each other. For instance, the Samsung S22 might often be sold together with a protection case for the Samsung S22. If we would then place the protection case in the second warehouse and the Samsung S22 in the first warehouse, combined delivery would take longer and this might lead to customers not ordering the Samsung S22 at all. To prevent this, we add the decision rule that if products are from different classes (based on profits as defined earlier) and are highly correlated (correlation > 0.6) we have to store them in the same warehouse.

Now, for every product we will stock a base-stock level of products. A base-stock level is defined as the average demand plus a mark-up based on the standard deviation. Importantly, we will never store the same product in two different warehouses: this would complicate operations significantly and consequently more likely result in higher costs. In total, the first warehouse has 960 storage boxes, of size 40x40x20 cms. Of this volume, we will assume that 90% can be used efficiently.

We will solve the problem of which product to assign to which warehouse such that the profit loss is minimized in three different ways. Firstly, we use two heuristics that give an upperbound to the total costs for the linear program we will use later. The first heuristic states that you assign products based on the average daily profit that they generate until you run out of storage space. The second heuristic uses the same algorithm, but then based on the ratio of the average daily profit with the number of boxes needed to store the product. In the linear program, we will solve a variation of the Knapsack problem.

Data

We have received data on the sales during a two-year period and the profit margins and dimensions of every product. In total, we have 1263 different products. To get a feeling for this data, we made histograms of the volumes and the average daily profits of the products in figure 1. This shows us that there are quite a lot of products with a low volume, which means you can store a lot of those products in one storage box. These might for example be protection cases, which usually have quite a low volume since it is not packaged in a box like phones are. The average daily profits plot shows us that there are a few products with very high daily profits, but that the epicentre is still at the lower segment. In figure 2 we show the average daily profit per product and showed in red the division where the top 20% ends and the middle 30% begins and where the middle 30% stops and the bottom 50% ends.

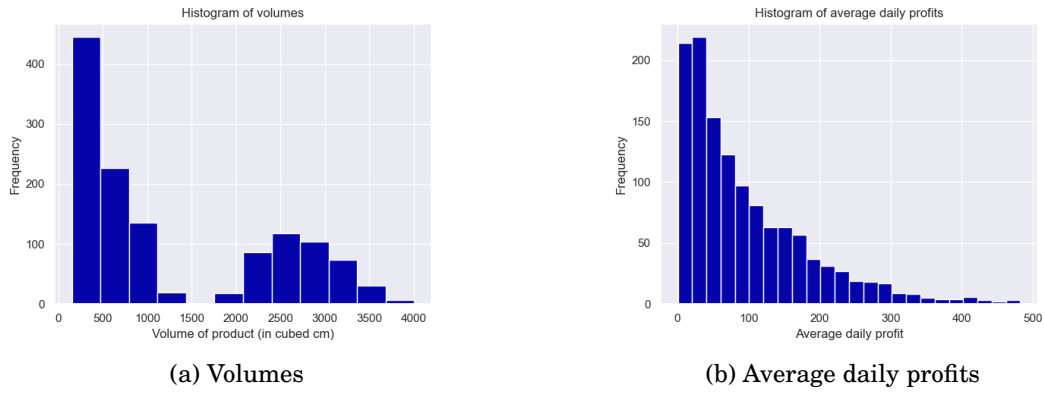


Figure 1: Histograms



Figure 2: Average daily sales per product

In figure 3 we show the number of products in each class and the total daily profit every class generates. This shows us that even though the top 20% has a lot less products, it does generate a lot more profit than the other two classes.

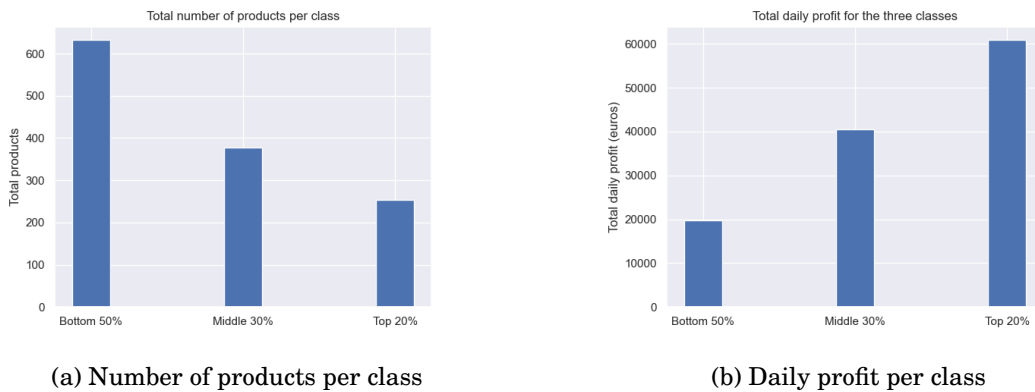


Figure 3: Comparison of classes

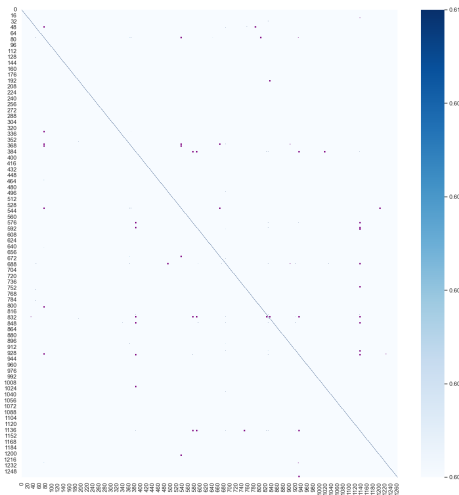


Figure 4: Heat map of the correlation matrix

In figure 4 we can see the correlations higher than 0.6 marked for the products. In here, we marked the "normal" correlations in blue and the correlations between products that should be kept in one warehouse in purple. Obviously, on the diagonal we have one big blue line: a product is 100% correlated with itself. On aggregate, we see that we have no big clumps of correlations, so the correlation requirement will not be too restrictive.

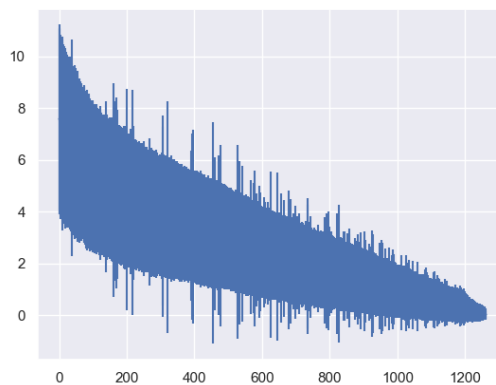


Figure 5: Error bars per product

Lastly, in figure 5 we see the error bars per product, with as mean the average number of orders per day and the length of the error bars is based on the standard deviation. This plot gives a small insight on how many products you will have to stock as base-stock level for a certain product, although the standard deviation is multiplied with a scalar bigger than 1 in the base-stock formula. Hence, these error bars give smaller base-stock levels than the actual ones.

Mathematical model

Heuristics

In the heuristics we prioritize products based on a certain measure. We will then place the products with the highest measure in the boxes, until all boxes are full. Importantly, if we place a product in the first warehouse, we also need to place its correlated products in the same warehouse. Obviously, if a product has a higher measure but there are not enough boxes available for the product (and its correlated products) we will still try to put a product with a lower measure in the first warehouse, as long as there are boxes available.

Specifically, in heuristic 1 the measure will be based on the average daily profit that a product generates, whereas in heuristic 2 the measure will be based on the ratio of the average daily profit to the boxes needed to store the product.

Knapsack problem

Let x_j denote if the product with product id j is stored in the **second** warehouse, $x_j \in \{0, 1\}$.

Let p_j denote the number of boxes needed to store a base-stock level of product id j .

Let c_j denote the average daily costs of storing product j in the second warehouse.

Then, we want to solve the following (integer) linear problem:

$$\begin{aligned} \min \quad & c^T x \\ \text{s.t.} \quad & p^T(\iota_n - x) \leq 960 \\ & x_i(1 - x_j) = 0 \quad \{i, j\} \in S \quad \forall S \\ & x_j \in \{0, 1\} \quad j = 1, \dots, n \end{aligned}$$

where the sets S contain two product ids, one of a higher class and one of a lower class, which are correlated in demand with a correlation of at least 0.6. Furthermore, ι_n is a (column) vector of ones of length n , where n is the number of products we have. Then, since x is binary and denotes the products stored in the second warehouse, $\iota_n - x$ yields the products stored in the first warehouse.

Results

In table 1 we report the costs of the different policies. We see that heuristic one does not perform very well, which makes sense: after all if a product generates a lot of profit but does need a lot of storage boxes, heuristic 1 will assign it all the boxes whereas in heuristic 2 or in the knapsack problem the product likely won't be assigned, leaving more room for the products that generate more profit per storage box. Heuristic 2 however performs very good, having less than a 2% difference in costs with the (for this set of assumptions) optimal solution given by the knapsack problem. However, even a 2% difference might mean a lot when we start talking about revenues of millions of dollars, as we are in the context of this report. The main reason why the knapsack problem performs slightly better is likely since it takes into account the fact that you do not merely need to account for the storage boxes needed for your own product, but also for the correlated

products. E.g., if we have a product that only needs one storage box, it will likely have a high ratio and therefore be preferred by heuristic 2, but if it has a correlated product that needs 6 storage boxes to be stored, the actual ratio would be a lot smaller. This, in turn, will result in higher costs for heuristic 2.

Assignment policy	Heuristic 1	Heuristic 2	Knapsack problem
Daily costs	€ 16626.62	€ 10033.7	€ 9884.26

Table 1: Average daily costs of allocation

Sensitivity analysis

In addition to our regular results, we also present a sensitivity analysis. Our first sensitivity analysis is on the number of boxes that the first warehouse has. We are primarily interested in what the changes of first warehouse sizes would have for impact on the average daily profit loss. Of course, having a bigger first warehouse will always result in a lower profit loss, but from figure 6 we can see that the rate at which the costs decrease gets smaller for larger percentages. From figure 6 we can conclude that adding for instance 5% warehouse space would be beneficial as long as that comes with a corresponding operations costs smaller than approximately 600 euros per day. From figure 7 we can conclude that the results of heuristic 2 stay very good for all different percentages and are therefore quite robust, whereas heuristic 1 is better for very few boxes in warehouse 1 or when almost all products can be stored in the first warehouse. Hence, heuristic 1 performs badly overall.

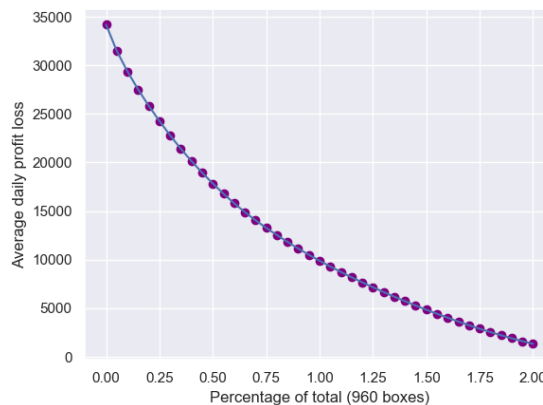


Figure 6: Daily costs for different warehouse sizes (0 - 1920 boxes)

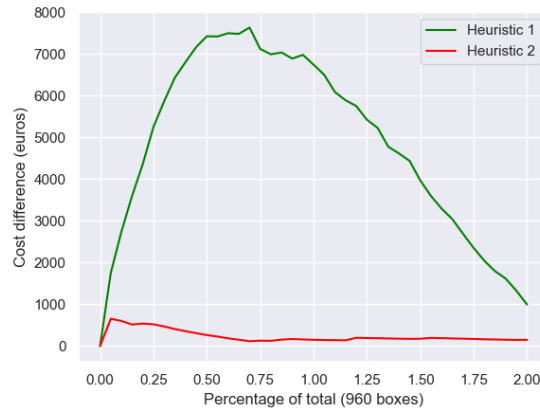


Figure 7: Difference in daily costs between Knapsack solution and heuristics

Our second sensitivity analysis is on the daily profit losses: we investigate what the effect of having different average daily profit losses per product is on the total profit loss. In figure 8 we plot this for different percentages of the usual amount and it might not come as a surprise that the costs simply linearly increase with the percentage. The allocation of which products to store in the first warehouse will after all stay the same and it will only result in different absolute costs, but the costs will stay the same relatively.

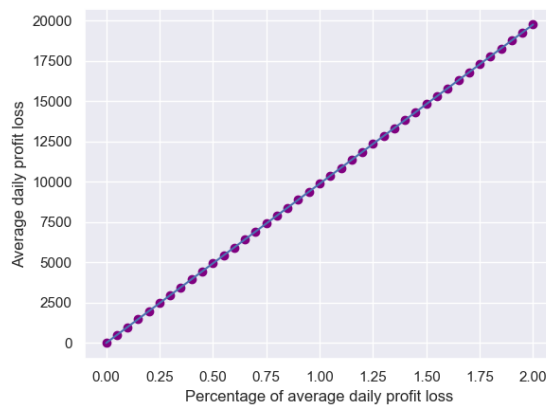


Figure 8: Daily costs for different average daily profit losses

Appendix

Main file

```
1  import pandas as pd
2  from elasticsearch import Elasticsearch
3  import matplotlib.pyplot as plt
4  import numpy as np
5  from scipy.stats import norm # Mainly used for the z-statistics
6  import seaborn as sns; sns.set_theme() # For the heatmap
7  from gurobipy import Model, GRB
8  from matplotlib.patches import Rectangle
9
10
11 pd.options.mode.chained_assignment = None # Ignores chained warnings
12 import urllib3
13 urllib3.disable_warnings(urllib3.exceptions.InsecureRequestWarning)
14 # Ignores unverified https request warning
15
16 import warnings
17 warnings.filterwarnings("ignore", category=DeprecationWarning)
18 # Ignores DeprecationWarnings
19 ## Questions:
20
21
22 ## Loading the csv file for sales into the elastic database
23 es = Elasticsearch( # Specifies the elasticsearch "location" to store the data
24     hosts="https://elastic:pancakes@localhost:9200/",
25     verify_certs=False
26 )
27
28 ## The next lines are commented, as we only want to run them once.
29 ## After all, if we have made the database there is no need to make it again
30
31
32 # Creates the elastic index in which we will put the data
33 settings = {
34     'settings': {
35         "number_of_shards" : 3
36     }
37     'mappings': {
38         'properties': {
39             'pickup_datetime': {'type': 'date',
40                                 "format": "yyyy-MM-dd HH:mm:ss" },
41             'dropoff_datetime': {'type': 'date',
42                                 "format": "yyyy-MM-dd HH:mm:ss" },
```

```
43         'pickup_location': {'type': 'geo_point' },
44         'dropoff_location': {'type': 'geo_point'}
45     }
46 }
47 }
48 es.indices.create(index="sales", body=settings)
49 print("index created")
50
51 import csv_to_elastic # Importing the python functions that allow us to,
52 # quickly, put the csv data into our elastic database
53 csv_to_elastic.ingest_csv_file_into_elastic_index(basepath + "/sales.csv",
54 es, "sales", 5000)
55
56
57
58 ## We load in the pandas dataframes (made in DataFrameConstruction)
59 dfProducts = pd.read_pickle("dfProducts.pkl")
60 dfOrderPerDay = pd.read_pickle("dfOrderPerDay.pkl")
61 ndays = dfProducts.shape[0]
62
63 # Here we construct an errorbar plot, with as mean the average orders and
64 # error the standard deviation
65 dfSorted = dfProducts.sort_values(by = ["AvgOrders"], ascending=False)
66 # Sorts the dataframe descending by Avg Orders
67 plt.errorbar(x = range(0, len(dfSorted["AvgOrders"])),
68             y = dfSorted["AvgOrders"], xerr = 0, yerr = dfSorted["StdOrders"])
69 plt.savefig('Errorbar orders.png')
70 plt.show()
71
72 # Here we make a histogram of the average daily profit generated by each product
73 dfProducts["AvgDailyProfit"] = dfProducts["AvgOrders"] * dfProducts["Margins"]
74 plt.hist(x=dfProducts["AvgDailyProfit"], bins='auto', color='#0504aa')
75 plt.xlabel('Average daily profit')
76 plt.ylabel('Frequency')
77 plt.title('Histogram of average daily profits')
78 plt.savefig('Histogram average daily profits.png')
79 plt.show()
80
81 # Here we make a histogram of the volume of each product
82 dfProducts["Volume"] = dfProducts["Length"] * dfProducts["Height"] \
83                        * dfProducts["Width"]
84 plt.hist(x=dfProducts["Volume"], bins='auto', color='#0504aa')
85 plt.xlabel('Volume of product (in cubed cm)')
86 plt.ylabel('Frequency')
87 plt.title('Histogram of volumes')
88 plt.savefig('Histogram volumes.png')
```

```
89 plt.show()
90
91 # We divide the products up in three different classes,
92 # the top 20% profit generating products, the following 30% and the bottom 50%
93 dfProfitSort = dfProducts.sort_values(by=["AvgDailyProfit"])
94 nrow = dfProfitSort.shape[0]
95 bottom50 = dfProfitSort[0:(round(0.5*nrow))]
96 middle30 = dfProfitSort[(round(0.5*nrow)):(round(0.8*nrow))]
97 top20 = dfProfitSort[(round(0.8*nrow)):nrow]
98
99 # We will compute the total profit and the number of products per class and
100 # display them in a barplot
101 fig = plt.figure()
102 totprofits = [sum(bottom50["AvgDailyProfit"]), sum(middle30["AvgDailyProfit"]),
103               sum(top20["AvgDailyProfit"])]
104 plt.bar(["Bottom 50%", "Middle 30%", "Top 20%"], totprofits, color = 'b',
105         width = 0.25)
106 plt.ylabel("Total daily profit (euros)")
107 plt.title("Total daily profit for the three classes")
108 plt.savefig('Total daily profit per class.png')
109 plt.show()
110
111 fig = plt.figure()
112 totsalses = [len(bottom50["AvgOrders"]), len(middle30["AvgOrders"]),
113              len(top20["AvgOrders"])]
114 plt.bar(["Bottom 50%", "Middle 30%", "Top 20%"], totsalses, color = 'b'
115         , width = 0.25)
116 plt.ylabel("Total products")
117 plt.title("Total number of products per class")
118 plt.savefig('Barplot products.png')
119 plt.show()
120
121 # Making a barplot that shows the daily average profits, marking the
122 # products on the boundary of the product classes
123 fig = plt.figure()
124 colours = ["blue"] * nrow
125 max_profit = dfProfitSort["AvgDailyProfit"].iloc[-1]
126 plt.bar(range(0, nrow), dfProfitSort["AvgDailyProfit"][:, -1],
127         color=colours, width = 0.2)
128 plt.plot([round(0.5*nrow), round(0.5*nrow)], [0,max_profit],
129         color='r', linestyle='-', linewidth=2)
130 plt.plot([round(0.2*nrow), round(0.2*nrow)], [0,max_profit],
131         color='r', linestyle='-', linewidth=2)
132 plt.ylabel("Average daily sales (in euros)")
133 plt.title("Average daily sales")
134 plt.savefig('Barplot daily sales.png')
```

```
135 plt.show()
136
137 # Calculating the mean and standard deviation for the demand over the
138 # replenishment period
139 avg_list = 7 * dfProfitSort["AvgOrders"]
140 std_list = np.sqrt(7) * dfProfitSort["StdOrders"]
141
142 # Calculating the base-stock levels
143 quantile = [norm.ppf(0.9), norm.ppf(0.95), norm.ppf(0.99)]
144 # We do everything times 7 (sqrt(7)) because it is based on weekly data
145 bottom50["Basestock"] = round(bottom50["AvgOrders"] * 7 +
146                                np.sqrt(7) * bottom50["StdOrders"] * quantile[0])
147 middle30["Basestock"] = round(middle30["AvgOrders"] * 7 +
148                                np.sqrt(7) * middle30["StdOrders"] * quantile[1])
149 top20["Basestock"] = round(top20["AvgOrders"] * 7 +
150                              np.sqrt(7) * top20["StdOrders"] * quantile[2])
151
152 volume_box = 40 * 40 * 20 # in cm^3
153 bottom50["NrProdPerBox"] = (0.9 * volume_box) // bottom50["Volume"]
154 # Notice the floor division and that we only use 90% efficiently
155 middle30["NrProdPerBox"] = (0.9 * volume_box) // middle30["Volume"]
156 top20["NrProdPerBox"] = (0.9 * volume_box) // top20["Volume"]
157
158 # Calculating the correlation matrix:
159 corr_basis = dfOrderPerDay.drop(columns="Day")
160 corr_matrix = np.corrcoef(corr_basis, rowvar = False)
161 plt.figure(figsize=(16, 16), dpi=100)
162 g = sns.heatmap(corr_matrix, vmin = 0.6, vmax = 0.61, cmap="Blues", cbar = True)
163 plt.savefig('heat_map.png')
164 plt.show()
165
166 # Finding the product couples
167 # Firstly, we set the diagonal elements to 1 (as they should equal,
168 # but sometimes they are 0.9999 right now)
169 for i in range(len(corr_matrix[0])):
170     corr_matrix[i][i] = 1
171
172 # We can then find the values between 0.6 and 1 (exclusive), those are the
173 # relevant correlations over the threshold value
174 corr_list = [] # List of indices for which we have a high enough correlation
175 for i in range(len(corr_matrix[0])):
176     temp_list = (corr_matrix[i] > 0.6) & (corr_matrix[i] < 1)
177     corr_list.append([i for i, x in enumerate(temp_list) if x])
178
179 prod_couples = []
180 # In this loop we find the combination of top20 products with lower products
```

```
181 for i in top20["product_id"]:
182     if len(corr_list[i]) > 0: # If there is an element in corr_list
183         for j in range(0, len(corr_list[i])): # Loop over it
184             if (sum(top20["product_id"] == corr_list[i][j]) == 0): # Not in top 20%
185                 prod_couples.append([i, corr_list[i][j]])
186 # In this loop we find the combination of middle30 products with lower products
187 for i in middle30["product_id"]:
188     if len(corr_list[i]) > 0: # If there is an element in corr_list
189         for j in range(0, len(corr_list[i])): # Loop over it
190             if (sum(middle30["product_id"] == corr_list[i][j]) == 0): # Not in middle 30%
191                 prod_couples.append([i, corr_list[i][j]])
192
193 print("Number of products couples between classes is", len(prod_couples))
194
195 # Highlighting these cells in the heatmap
196 plt.figure(figsize=(16, 16), dpi=100)
197 g = sns.heatmap(corr_matrix, vmin = 0.6, vmax = 0.61, cmap="Blues", cbar = True)
198 for i in prod_couples:
199     g.add_patch(Rectangle((i[0], i[1]), 1, 1, fill=False, edgecolor='purple',
200                          lw=2))
201 plt.savefig('heat_map_marked.png')
202 plt.show()
203
204 ##### Optimization #####
205 # Calculating the average daily profit loss per product, which is the demand
206 # loss percentage times the average daily profit
207 bottom50["AvgProfitLoss"] = 0.5 * bottom50["AvgDailyProfit"]
208 middle30["AvgProfitLoss"] = 0.3 * middle30["AvgDailyProfit"]
209 top20["AvgProfitLoss"] = 0.2 * top20["AvgDailyProfit"]
210
211 # The dataframe that will be used in all optimization procedures
212 dfMerged = pd.concat([top20, middle30, bottom50], axis = 0) # Merging the dataframes
213 dfMerged["BoxesNeeded"] = -(dfMerged["Basestock"] // -dfMerged["NrProdPerBox"])
214 dfMerged = dfMerged.sort_values(by = ["product_id"])
215
216 dfMerged["Ratio"] = dfMerged["AvgProfitLoss"] / dfMerged["BoxesNeeded"]
217 dfMerged = dfMerged.sort_values(by = ["product_id"])
218
219 ## Heuristic 1:
220 # In this heuristic, we will place a base-stock level of products in boxes
221 # for the product with the highest AvgProfitLoss, then fill boxes with the
222 # second highest AvgProfitLoss etc. until all boxes are full
223 # Flaws:
224 # 1. Assumes that the product_ids follow on each other (so 1,2,3 not 1,2,5 etc.)
225 # 2. Background: Correlations calculated between classes based on Profit, not
226 # products that on average have the highest Profit Loss
```

```
227
228 # We use ceiling division here, but since this is not defined with a simple
229 # operator in python, we use the floor division operator in a smart way
230
231 def heuristic1(dataframe, boxes):
232     boxes_available = boxes
233     dfPop = dataframe.copy() # No reference!
234     dfPop = dfPop.sort_values(by=["AvgProfitLoss"], ascending=False)
235     # dfPop = dfMerged.copy() # No reference!
236
237     prod_ids_first = [] # This will contain the products that are stored in
238                       # the first warehouse
239
240     i = 0
241     while i < dfPop.shape[0]: # dfPop will contain the products that are not in
242                               # the first warehouse yet
243         prod_id = dfPop.iloc[i, 0]
244         corr_prods = []
245
246         # In this loop we check if we need to add other products to the first
247         # warehouse too
248         for j in range(0, len(prod_couples)):
249             if prod_id == prod_couples[j][0]: # First element is highest sale value
250                 corr_prods.append(prod_couples[j][1] - 1)
251         if len(corr_prods) > 0:
252             corr_prods.append(prod_id - 1)
253             row_indices = corr_prods
254         else:
255             row_indices = [prod_id - 1]
256
257         tot_size_needed = np.sum(dataframe.iloc[row_indices, -2])
258         # Total size needed to fit the product and its correlated products
259
260         if (tot_size_needed <= boxes_available): # If we have enough room available
261             # We can assign the products to the boxes
262             boxes_available -= tot_size_needed
263
264         for k in range(0, len(row_indices)): # Loop over the row_indices
265             # Remove the indices from the dataframe such that we do not
266             # assign the product to a box twice
267             dfPop.drop(dfPop[(dfPop['product_id'] == row_indices[k] + 1)].index,
268                       inplace=True)
269             prod_ids_first.append(row_indices[k] + 1)
270
271         i -= 1 # Since we removed a row at place i, what was row i + 1 is now
272             # row i and therefore we want i to stay the same. Since we add one later
273             # we place -1 here. Notice that the rows after row i that were removed
```

```
273         # (those from the correlated products), are irrelevant for our choice of i
274
275         # Small if statement to make sure we do not run the loop longer than necessary
276         if (boxes_available == 0):
277             break
278
279         i += 1
280     return prod_ids_first
281
282     # Now we also find the products that are stored in the second warehouse:
283     prod_ids_first = heuristic1(dfMerged, 960)
284     prod_ids_first.sort()
285     prod_ids_second = [x for x in range(1, nrow + 1) if x not in prod_ids_first]
286
287     # The actual loss made with this policy:
288     total_loss = 0
289     for i in range(0, len(prod_ids_second)):
290         total_loss += dfMerged.iloc[prod_ids_second[i] - 1, -3]
291     # Can do this since dfMerged is sorted on product id
292     print("The total daily cost of implementing Heuristic 1 is: €", round(total_loss, 2),
293           sep = "")
294
295
296
297     ## Heuristic 2:
298     # In this heuristic, we will place a base-stock level of products in boxes
299     # for the product with the highest ratio AvgProfitLoss / BoxesNeeded, then
300     # fill boxes with the second highest ratio etc. until all boxes are full
301     # Flaws:
302     # 1. Assumes that the product_ids follow on each other (so 1,2,3 not 1,2,5 etc.)
303     # 2. Does not consider the ratio of the correlated products
304     # 3. Background: Correlations calculated between classes based on Profit, not
305     # products that on average have the highest Profit Loss
306     def heuristic2(dataframe, boxes):
307         boxes_available = boxes
308         dfPop = dataframe.copy() # No reference!
309         dfPop = dfPop.sort_values(by = ["Ratio"], ascending = False) # No reference!
310
311         prod_ids_first = [] # This will contain the products that are stored in
312         # the first warehouse
313         i = 0
314         while i < dfPop.shape[0]: # dfPop will contain the products that are not in
315             # the first warehouse yet
316             prod_id = dfPop.iloc[i, 0]
317             corr_prods = []
318
```

```
319     # In this loop we check if we need to add other products to the first
320     # warehouse too
321     for j in range(0, len(prod_couples)):
322         if prod_id == prod_couples[j][0]: # First element is highest sale value
323             corr_prods.append(prod_couples[j][1] - 1)
324     if len(corr_prods) > 0:
325         corr_prods.append(prod_id - 1)
326         row_indices = corr_prods
327     else:
328         row_indices = [prod_id - 1]
329
330     tot_size_needed = np.sum(dataframe.iloc[row_indices, -2])
331     # Total size needed to fit the product and its correlated products
332     # Notice the -2 now since we added a column to dfMerged
333
334
335     if (tot_size_needed <= boxes_available): # If we have enough room
336         # available we can assign the products to the boxes
337         boxes_available -= tot_size_needed
338
339         for k in range(0, len(row_indices)): # Loop over the row_indices
340             # Remove the indices from the dataframe such that we do not assign
341             # the product to a box twice
342             dfPop.drop(dfPop[(dfPop['product_id'] == row_indices[k] + 1)].index,
343                       inplace=True)
344             prod_ids_first.append(row_indices[k] + 1)
345
346         i -= 1 # Since we removed a row at place i, what was row i + 1 is now
347         # row i and therefore we want i to stay the same. Since we add one later
348         # we place -1 here. Notice that the rows after row i that were removed
349         # (those from the correlated products), are irrelevant for our choice of i
350
351     # Small if statement to make sure we do not run the loop longer than necessary
352     if (boxes_available == 0):
353         break
354
355     i += 1
356
357     return prod_ids_first
358
359 prod_ids_first = heuristic2(dfMerged, 960)
360 prod_ids_first.sort()
361 # Now we also find the products that are stored in the second warehouse:
362 prod_ids_second = [x for x in range(1, nrow + 1) if x not in prod_ids_first]
363
364 # The actual loss made with this policy:
```



```
365 total_loss = 0
366 for i in range(0, len(prod_ids_second)):
367     total_loss += dfMerged.iloc[prod_ids_second[i] - 1, -3]
368     # Can do this since dfMerged is sorted on product id
369     print("The total daily cost of implementing Heuristic 2 is: €",
370           round(total_loss, 2), sep = "")
371
372 ## Knapsack problem:
373 def Knapsack(dataframe, boxes):
374     m = Model("Knapsack problem")
375     x = m.addVars(nrow, vtype=GRB.BINARY)
376     c = dataframe["AvgProfitLoss"] # Costs of adding product to second warehouse
377     p = dataframe["BoxesNeeded"] # Boxes needed per product
378
379     # The constraint that makes sure no more than 960 boxes are used
380     m.addConstr(sum(p[j] * (1 - x[j]) for j in range(len(x))) <= boxes)
381
382     # The constraints that make sure correlated products are in the same warehouse
383     for i,j in prod_couples:
384         m.addConstr(x[i] * (1 - x[j]) == 0)
385
386     m.setObjective(sum(c[j] * x[j] for j in range(len(x))), GRB.MINIMIZE)
387     m.setParam('OutputFlag', 0)
388     m.optimize()
389
390     return m.X
391
392 prod_ids_second = Knapsack(dfMerged, 960) # Recall our definition specified
393 # the x-variables as the products stored in the second (!) warehouse
394 prod_ids_second = [i + 1 for i in range(0, len(prod_ids_second))
395                   if prod_ids_second[i] > 0]
396 prod_ids_first = [x for x in range(1, nrow + 1) if x not in prod_ids_second]
397
398 total_loss = 0
399 for i in range(0, len(prod_ids_second)):
400     total_loss += dfMerged.iloc[prod_ids_second[i] - 1, -3]
401     print("The total daily cost of implementing the Knapsack solution is: €",
402           round(total_loss, 2), sep = "")
403     print("The products stored in the first warehouse are:", prod_ids_first)
404
405
406 #### Sensitivity Analysis ####
407 # On the available boxes
408 dfSensCost = pd.DataFrame({"Percentage": [], "Heuristic 1": [],
409                           "Heuristic 2": [], "Knapsack": []})
410 for i in range(0, 41):
```

```
411     print(i)
412     percentage = i * 0.05
413     boxes_available = 960 * percentage
414     h1_first = heuristic1(dfMerged, boxes_available)
415     h2_first = heuristic2(dfMerged, boxes_available)
416     k_second = Knapsack(dfMerged, boxes_available)
417     k_second = [i + 1 for i in range(0, len(k_second)) if k_second[i] > 0]
418     h1_second = [x for x in range(1, nrow + 1) if x not in h1_first]
419     h2_second = [x for x in range(1, nrow + 1) if x not in h2_first]
420     h1_costs = 0
421     h2_costs = 0
422     k_costs = 0
423     for i in range(0, len(k_second)):
424         k_costs += dfMerged.iloc[k_second[i] - 1, -3]
425     for i in range(0, len(h1_second)):
426         h1_costs += dfMerged.iloc[h1_second[i] - 1, -3]
427     for i in range(0, len(h2_second)):
428         h2_costs += dfMerged.iloc[h2_second[i] - 1, -3]
429     cost_frame = pd.DataFrame({"Percentage": [percentage],
430                               "Heuristic 1": [h1_costs],
431                               "Heuristic 2": [h2_costs],
432                               "Knapsack": [k_costs]})
433     dfSensCost = pd.concat([dfSensCost, cost_frame], ignore_index=True, axis = 0)
434
435     print(dfSensCost)
436
437     plt.plot(dfSensCost.iloc[:, 0], dfSensCost.iloc[:, -1])
438     plt.xlabel("Percentage of total (960 boxes)")
439     plt.ylabel("Average daily profit loss")
440
441     for i in range(0, 41):
442         plt.scatter(dfSensCost.iloc[i, 0], dfSensCost.iloc[i, -1], color = "purple")
443     plt.savefig('Costs_knapsack.png')
444     plt.show()
445
446     dfSensCost["K-h1"] = dfSensCost["Heuristic 1"] - dfSensCost["Knapsack"]
447     dfSensCost["K-h2"] = dfSensCost["Heuristic 2"] - dfSensCost["Knapsack"]
448
449     plt.plot(dfSensCost.iloc[:, 0], dfSensCost.iloc[:, -2], color = "green",
450             label = "Heuristic 1")
451     plt.plot(dfSensCost.iloc[:, 0], dfSensCost.iloc[:, -1], color = "red",
452             label = "Heuristic 2")
453     plt.xlabel("Percentage of total (960 boxes)")
454     plt.ylabel("Cost difference (euros)")
455     plt.legend(loc = "upper right")
456     plt.savefig('Cost_difference.png')
```

```
457 plt.show()
458
459 # On the average daily profit loss
460 dfSensCost2 = pd.DataFrame({"Percentage": [], "Knapsack": []})
461 for i in range(0, 41):
462     print(i)
463     percentage = i * 0.05
464     dfCopy = dfMerged.copy()
465     dfCopy["AvgProfitLoss"] = percentage * dfCopy["AvgProfitLoss"]
466     k_second = Knapsack(dfCopy, 960)
467     k_second = [i + 1 for i in range(0, len(k_second)) if k_second[i] > 0]
468     k_costs = 0
469     for i in range(0, len(k_second)):
470         k_costs += dfCopy.iloc[k_second[i] - 1, -3]
471     cost_frame = pd.DataFrame({"Percentage": [percentage], "Knapsack": [k_costs]})
472     dfSensCost2 = pd.concat([dfSensCost2, cost_frame], ignore_index=True, axis = 0)
473
474 plt.plot(dfSensCost2.iloc[:, 0], dfSensCost2.iloc[:, 1])
475 plt.xlabel("Percentage of average daily profit loss")
476 plt.ylabel("Average daily profit loss")
477
478 print(dfSensCost2)
479 for i in range(0, 41):
480     plt.scatter(dfSensCost2.iloc[i, 0], dfSensCost2.iloc[i, 1], color = "purple")
481 plt.savefig('Costs_knapsack_2.png')
482 plt.show()
```

DataFrameConstruction

```
1  import statistics
2  import pandas as pd
3  from elasticsearch import Elasticsearch
4  from pandas import json_normalize
5
6  import urllib3
7  urllib3.disable_warnings(urllib3.exceptions.InsecureRequestWarning)
8  # Ignores unverified https request warning
9
10 import warnings
11 warnings.filterwarnings("ignore", category=DeprecationWarning)
12 # Ignores DeprecationWarnings
13
14 es = Elasticsearch( # Specifies the elasticsearch "location" to store the data
15     hosts="https://elastic:pancakes@localhost:9200/",
16     verify_certs=False
17 )
18
19 ## Loading in margins and dimensions from csv to pandas dataframe
20 basepath = "C:/Users/lucas/OneDrive/Documenten/RUG/Msc Year 1/Trimester 1" \
21     "/Data analysis and programming for OM/Assignment"
22 margins = pd.read_csv(basepath + "/margins.csv")
23 dimensions = pd.read_csv(basepath + "/dimensions.csv")
24
25 search_body = { # search_body that aggregates on product_id,
26     # so giving the demand per product over the whole period
27     "size": 0,
28     "aggs": {
29         "PerProduct": {
30             "terms": { # Here we specify that we aggregate on days
31                 "field": "product_id",
32                 "size": 10000,
33                 "order": {
34                     "_key": "asc"
35                 }
36             }
37         }
38     }
39 }
40
41 result = es.search(index="sales", body=search_body)
42 dfProducts = json_normalize(result['aggregations']['PerProduct']['buckets'])
43 dfProducts.rename(columns={'key': 'product_id'}, inplace=True)
44
```

```
45 search_body = { # search_body that aggregates on days, so giving the
46                 # demand per day over all products
47     "size": 0,
48     "aggs": {
49         "PerDay": {
50             "terms": { # Here we specify that we aggregate on days
51                 "field": "day",
52                 "size": 10000,
53                 "order": {
54                     "_key": "asc"
55                 }
56             }
57         }
58     }
59 }

60
61 result = es.search(index="sales", body=search_body)
62 days = result['aggregations']['PerDay']['buckets']
63 ndays = len(days)
64 df = pd.DataFrame(list(range(1, ndays + 1)), columns=["Day"])
65 # Dataframe with the days (1 to 730 in this case). We specify this one, as we
66 # want to merge our dataframes such that if we have no demand on say day 5,
67 # we get a zero there instead of ignoring that row
68
69 for index, row in dfProducts.iterrows():
70     search_body = {
71         "size": 10000, # Needs to be larger than number of product orders
72         'query': { # Here we find the orders with product number
73                 # row["product_id"]
74                 'term': {
75                     'product_id': row["product_id"]
76                 }
77         },
78         "aggs": { # Here we find the orders per day for said product
79             "OrdersPerDay": {
80                 "terms": {
81                     "field": "day",
82                     "size": 1000,
83                     "order": {
84                         "_key": "asc"
85                     }
86                 }
87             }
88         }
89     }
90     result = es.search(index="sales", body=search_body)
```

```
91
92     secondFrame = pd.DataFrame(result["aggregations"]["OrdersPerDay"]["buckets"])
93     # The orders per day
94     secondFrame = secondFrame.rename(columns={"key": "Day", "doc_count": "ID" +
95                                             str(row["product_id"])}))
96
97     df = df.merge(secondFrame, how="left", on=["Day"]).fillna(0)
98     # Constructs a dataframe with Orders per day, putting
99     # 0 at the places where we had no orders.
100
101     Prod_list_mean = []
102     Prod_list_std = []
103     for i in range(1, df.shape[1]):
104         Prod_list_mean.append(statistics.mean(list(df["ID" + str(i)])))
105         Prod_list_std.append(statistics.stdev(list(df["ID" + str(i)])))
106
107     dfProducts["AvgOrders"] = Prod_list_mean
108     dfProducts["StdOrders"] = Prod_list_std
109     dfProducts["Margins"] = margins["margin"]
110     dfProducts["Length"] = dimensions["length"]
111     dfProducts["Width"] = dimensions["width"]
112     dfProducts["Height"] = dimensions["height"]
113
114     dfProducts.to_pickle("dfProducts.pkl")
115     df.to_pickle("dfOrderPerDay.pkl")
```

CsvToElastic

```
1 from elasticsearch import Elasticsearch, helpers
2 import uuid
3 import pandas as pd
4 import json
5
6
7 def convert_csv_file_to_bufferized_json_lines_list(csv_file_name, buffer_size):
8     df = pd.read_csv(csv_file_name)
9     json_list = json.loads(df.to_json(orient='records'))
10
11     buffer_list = []
12     total_nr_docs = len(json_list)
13     for i in range(0, total_nr_docs, buffer_size):
14         buffer_list.append(json_list[i: i + buffer_size])
15     return buffer_list
16
17 def bulk_json(json_buffer, _index):
18     for doc in json_buffer:
19         # use a `yield` generator so that the data
20         # isn't loaded into memory
21         if '{"index"' not in doc:
22             yield {
23                 "_index": _index,
24                 "_id": uuid.uuid4(),
25                 "_source": doc
26             }
27
28 def ingest_csv_file_into_elastic_index(csv_file_name,
29                                       elastic_client: Elasticsearch, index_name, buffer_size=5000):
30     chunks = convert_csv_file_to_bufferized_json_lines_list(csv_file_name,
31                                                             buffer_size=buffer_size)
32     for i, buffer in zip(range(len(chunks)), chunks):
33         try:
34             # print("__TEST: first document from buffer: ", buffer[0])
35             response = helpers.bulk(elastic_client, bulk_json(json_buffer=buffer,
36                                                             _index=index_name))
37             print("bulk_json() RESPONSE for chunk:", i, response)
38         except Exception as e:
39             print("\nERROR:", e)
```
