

```
1 ### md
2 # Verisk - computer purchase problem
3
4 *by Jolanta Śliwa*
5 ### md
6 ## Problem Statement
7 Suppose you're trying to help a company determine
  which computers to purchase.
8 ### Data - utilization data by employee:
9 The company has been able to pull utilization data by
  employee that classifies users into 3 bins,
  depending on how much they use their computer in
  their work:
10 * Low usage - spends a lot of time in meetings,
   checking email, doing people management
11 * Average usage - requires some compute power, with
   balanced mix of heads down/technical work along with
   a
12 good amount of meetings/email writing
13 * High usage - power user, relies heavily on computer
   performance
14
15 ###
16 import pandas as pd
17
18
19 utilization = pd.read_csv(
20     "https://raw.githubusercontent.com/shubhamkalra27
   /dsep-2020/main/datasets/util_b_emp.csv"
21 )
22 ###
23 utilization.info()
24 ###
25 utilization.head()
26 ### md
27 Checking types of utilization
28 ###
29 utilization["utilization_bin"].unique()
30 ### md
31 average usage is stored as medium
32 ### md
```

```
33 ## Data - survey
34 Additionally, they've surveyed employees to collect
    the relative importance of the following variables
    describing a
35 computer's performance:
36 * Memory
37 * Processing
38 * Storage
39 * Price inverse - this metric was given to you by the
    company as you can see in the dataset, with the
    directive that
40 price inverse being fixed at a 25% weight in the
    purchase decision
41 ###
42 survey = pd.read_csv(
43     "https://raw.githubusercontent.com/shubhamkalra27
        /dsep-2020/main/datasets/survey_emp.csv"
44 )
45 ###
46 survey.info()
47 ###
48 survey.head()
49 ### md
50 Checking whether we have 100% in every column
51 ###
52 for i, row in survey.iterrows():
53     if sum(row[1:]) != 1:
54         print("problem")
55         break
56 print("ok")
57 ### md
58 ## Data - computers
59 Lastly, the company is looking to purchase a maximum
    of 3 different computer models, and have compiled the
    following
60 list scoring their memory, processing, storage, and
    relative price. Each dimension is scored from 0-10,
    with 10 being the best.
61 ###
62 computers = pd.read_csv(
63     "https://raw.githubusercontent.com/shubhamkalra27
```

```
63 /dsep-2020/main/datasets/vendor_options.csv"
64 )
65 #%%
66 computers.info()
67 #%%
68 computers.head()
69 #%% md
70 Checking "real" range of scores:
71 #%%
72 print(computers.max())
73 #%%
74 print(computers.min())
75 #%% md
76 ## Task
77 **Given this information, provide the company with a
  recommendation on which computers to purchase.**
78 #%% md
79 List of parameters:
80 #%%
81 parameters = computers.columns[1:]
82 print(parameters)
83 #%% md
84 It will be more convenient for me to store all
  employees related data in one DataFrame instead of
  two.
85
86 Merging survey and utilization into employees column
  :
87 #%%
88 employees = utilization.merge(survey, left_on="
  employee_id", right_on="employee_id")
89 #%%
90 employees.info()
91 #%%
92 employees.head()
93 #%% md
94 Normally (if we would like to make some predictions
  ) it would be better to store "object" (nominal data
  - categorical data) using one-hot-encoding but in
  this case it is more convenient for me to leave it
  like that.
```

```

95 ### md
96 Under we will see outputs of different metrics. For
   each there will be solution provided by:
97 * Simulated Annealing
98 * Naive algorithm
99 ### md
100 Results given by algorithms: id of computer in
    DataFrame - not the one in colum computer_id
101 ###
102 from problem import (
103     ProblemNothing,
104     ProblemMax,
105     ProblemScale,
106     ProblemMaxHalf,
107     ProblemScaleHalf,
108 )
109 ### md
110 * ProblemNothing - metric without using utilization
    info
111 * ProblemMax - metric where every group have a
    different max score (3, 7, 10)
112 * ProblemMaxHalf - like above with different values
    (5, 7.5, 10)
113 * ProblemScale - metric where we scale computers
    scores by multiplying them (3/1, 3/2, 3/3)
114 * ProblemScaleHalf - like above with different
    values (4/2, 4/3, 4/4)
115 ###
116 from simulated_annealing import SimulatedAnnealing,
    SimulatedAnnealingConfig
117 from naive_solution import Naive
118 ### md
119 ### No utilization value
120 ###
121 prob_n = ProblemNothing(computers, employees)
122 ###
123 annealing_n = SimulatedAnnealing(
    SimulatedAnnealingConfig(), prob_n)
124
125 annealing_n.solve()
126 ###

```

```
127 naive_n = Naive(prob_n)
128
129 naive_n.solve()
130 ### md
131 As we can see simulated annealing sometimes returns
    solutions that aren't optimal
132 Let's try one more time
133 ###
134 annealing_n = SimulatedAnnealing(
    SimulatedAnnealingConfig(), prob_n)
135
136 annealing_n.solve()
137 ### md
138 This time we manage to get "the best" solution for
    this metric
139
140 Let's see how it looks like:
141 ###
142 prob_n.calculate_state_cost([10, 1, 8])
143 ###
144 prob_n.calculate_state_cost([8, 2, 1])
145 ### md
146 So actually there are at least 2 optimal solutions
147 ### md
148 ### Problem MAX
149 #### max: (3, 7, 10)
150 ###
151 prob_max = ProblemMax(computers, employees)
152 ###
153 annealing_max = SimulatedAnnealing(
    SimulatedAnnealingConfig(), prob_max)
154
155 result_max = annealing_max.solve()
156 ###
157 top_max = prob_max.get_wanted_computers(result_max)
158 ###
159 print(top_max.keys())
160 ###
161 naive_max = Naive(prob_max)
162
163 result_max_n = naive_max.solve()
```

```
164 #%%
165 print(prob_max.get_wanted_computers(result_max_n).
      keys())
166 #%% md
167 ##### max: (5, 7.5, 10)
168 #%%
169 prob_max_half = ProblemMaxHalf(computers, employees)
170 #%%
171 annealing_max_half = SimulatedAnnealing(
      SimulatedAnnealingConfig(), prob_max_half)
172
173 result_max_half = annealing_max_half.solve()
174 #%%
175 print(prob_max_half.get_wanted_computers(
      result_max_half).keys())
176 #%%
177 naive_max_half = Naive(prob_max_half)
178
179 result_max_half_n = naive_max_half.solve()
180 #%%
181 print(prob_max_half.get_wanted_computers(
      result_max_half_n).keys())
182 #%% md
183 ### Problem Scale
184 ##### scale: (3/1, 3/2, 3/3)
185 #%%
186 prob_scale = ProblemScale(computers, employees)
187 #%%
188 annealing_scale = SimulatedAnnealing(
      SimulatedAnnealingConfig(), prob_scale)
189
190 result_scale = annealing_scale.solve()
191 #%%
192 print(prob_scale.get_wanted_computers(result_scale).
      keys())
193 #%%
194 naive_scale = Naive(prob_scale)
195
196 result_scale_n = naive_scale.solve()
197 #%%
198 print(prob_scale.get_wanted_computers(result_scale_n
```

```
198 ).keys())
199 ### md
200 #### scale: (4/2, 4/3, 4/4)
201 ###
202 prob_scale_half = ProblemScaleHalf(computers,
    employees)
203 ###
204 annealing_scale_half = SimulatedAnnealing(
    SimulatedAnnealingConfig(), prob_scale_half)
205
206 result_scale_half = annealing_scale_half.solve()
207 ###
208 print(prob_scale_half.get_wanted_computers(
    result_scale_half).keys())
209 ###
210 naive_scale_half = Naive(prob_scale_half)
211
212 result_scale_half_n = naive_scale_half.solve()
213 ###
214 print(prob_scale_half.get_wanted_computers(
    result_scale_half_n).keys())
215 ### md
216 There is still a question which 3 computers are the
    best for that company?
217
218 We can see according to all presented metrics
    computer in row nr 1 is always in top three
219
220 Others can vary depending on metric we use. We can
    choose one of the above method but there are only 3
    candidates for 2 positions.
221 nr five wasn't choosen only by problem max
222 ###
223 print(result_max_half)
224 ###
225 print(prob_max_half.calculate_state_cost(
    result_max_half))
226 ###
227 print(prob_max_half.calculate_state_cost([1, 3, 5]))
228 ###
229 print(prob_max_half.calculate_state_cost([1, 5, 4]))
```

```
230 ### md
231 For that metric there is not a big difference
    between computer 5 and 4
232 nr 3 seems to better option to leave
233 ### md
234 On the other hand we see that nr 4 also wasn't
    chosen only by on metric: scale half
235 So it seems that we have final three
236 ###
237 print(result_scale_half)
238 ###
239 print(prob_scale_half.calculate_state_cost(
    result_scale_half))
240 ###
241 print(prob_scale_half.calculate_state_cost([1, 5, 4
    ]))
242 ###
243 print(prob_scale_half.calculate_state_cost([1, 4, 3
    ]))
244 ### md
245 Computer nr 4 seems to be less valuable and changing
    it with 3 or 5 give almost the same results
246 ### md
247 Just in case: Let's check nr 3 - wasn't chosen by 2
    metrics
248 ###
249 print(result_max)
250 ###
251 print(prob_max.calculate_state_cost(result_max))
252 ###
253 print(prob_max.calculate_state_cost([3, 5, 1]))
254 ###
255 print(prob_max.calculate_state_cost([3, 4, 1]))
256 ###
257 print(result_scale)
258 ###
259 print(prob_scale.calculate_state_cost(result_scale))
260 ###
261 print(prob_scale.calculate_state_cost([3, 5, 1]))
262 ###
263 print(prob_scale.calculate_state_cost([3, 4, 1]))
```



```
264 #%% md
265 As we can see using above metrics and adding
    computer 3 instead of any other (4 or 5) results in
    a bigger i other cases decrease of cost value
266 #%% md
267 Under we I displayed top three computers:
268 #%%
269 computers.iloc[[1, 4, 5]]
270 #%% md
271 seems like we ended up with rather ballanced final
    state
272 #%%
273 computers.iloc[[1, 4, 5]]["computer_id"]
274 #%% md
275 And here we have computer 3 as an addiction
276 #%%
277 computers.iloc[[3]]
278 #%% md
279 it is simmilar to computer nr 3 but with a slightly
    better price but worse memory
280 #%%
281
```