

DATA STRUCTURES FOR EARLY DETECTION OF BLIGHT PLAGUE ON COFFEE PLANTATIONS

Simón Flórez Silva
Universidad Eafit
Colombia
sflorezs1@eafit.edu.co

Adrián Alberto Gutiérrez Leal
Universidad Eafit
Colombia
aagutierrl@eafit.edu.co

Mauricio Toro
Universidad Eafit
Colombia
mtorobe@eafit.edu.co

ABSTRACT

In the following paper, we address the blight plague on coffee plantations, a problem that affects Colombian agriculture and therefore, its economy. In addition, lack of automation makes it harder to make an early diagnosis and, of course, makes it harder to control. As it was previously mentioned, the importance of this problem relapses on its consequences over the coffee exportation, from which thousands of Colombian families depend and, of course, the economy of the country as a whole. Some have already tried to treat decision problems such as this by implementing algorithms such as ID3 and C4.5 (which is a slightly better version of the latter) for decision trees.

KEYWORDS

• Mathematics of computing~Trees • Mathematics of computing~Graph algorithms • Information systems~Data access methods • Human-centered computing~Displays and imagers • Theory of computation~Recursive functions • Theory of computation~Data structures design and analysis • Hardware~Sensors and actuators • Hardware~Design modules and hierarchy • Software and its engineering~Interpreters • Software and its engineering~Virtual machines • Software and its engineering~Software performance • Software and its engineering~Object oriented languages • Software and its engineering~Abstract data types • Software and its engineering~Data types and structures • Software and its engineering~Classes and objects • Software and its engineering~Recursion

1. INTRODUCTION

In the current times, it becomes evident that automation can make some difficult processes a lot easier. That is the case for the early diagnosis and control of the blight plague on coffee plantations, a problem that has been harming productivity and incomes of Colombian coffee production during the last years. In this paper, the team will approach different algorithmic solutions and implement decision trees in order to detect this fungus as soon as possible.

2. PROBLEM

The main problem, from a software perspective, is to use Data Structures for the design of an algorithm based on

decision trees to predict if a caturra coffee batch is infected with the blight plague. It's worth noting that it has an immense impact on Colombian economy, as coffee consists in one of Colombia's main exports and therefore, many people depend on it to sustain themselves.

3. RELATED WORK

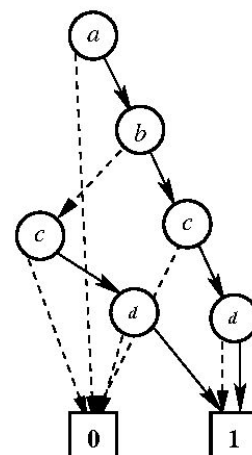
3.1 Fuzzy matching (Similar data matching)

According to Megter, Fuzzy matching is a technique used in computer-assisted translation as a special case of record linkage. It works with matches that may be less than 100% perfect when finding correspondences between segments of a text and entries in a database of previous translations. It usually operates at sentence-level segments, but some translation technology allows matching at a phrasal level. It is used when the translator is working with translation memory.

A solution example could be the n-gram algorithm, which is a contiguous sequence of n items from a given sequence of text or speech. The items can be phonemes, syllables, letters, words or base pairs according to the application. An n-gram model is a type of probabilistic language model for predicting the next item in such a sequence in the form of a $(n - 1)$ -order Markov model.

Example for the n-gram algorithm

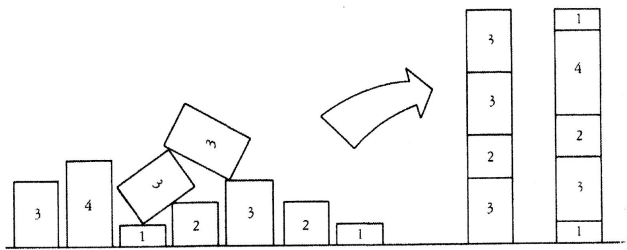
Sequences	Unordered 4-grams	Combinations
aabc	abc	abc
cabc	abc	abc
abcd	abcd	abcd
aacd	acd	acd



3.2 Partition Problem

Following Naumenko, F., Partition Problem consists of deciding whether a given set of positive integers with count of N can be partitioned into K subsets such that the sum of the numbers in each subset is equal. In actual practice, the exact equality of the sums is usually unattainable. A common example of its application would be the distribution of a set of independent works, each with different execution time, in several concurrent threads. If each thread runs on a separate core, we would like to be confident in the uniform load of each core in the CPU node.

There is an easy and well-known algorithmic approach for this problem, that could actually be solved by children, for this, the algorithm iterates through the numbers in descending order, assigning each of them to whichever subset has the smaller sum. This approach has a running time of $O(n \log n)$.

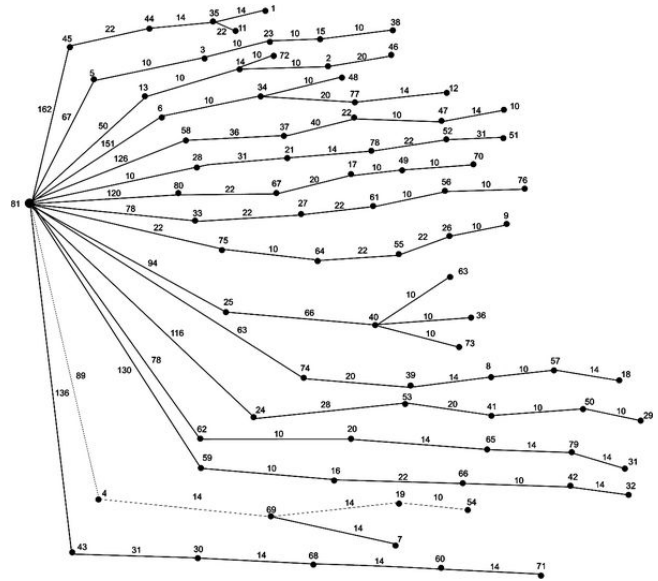


3.3 Capacitated minimum spanning tree

Capacitated minimum spanning tree is a minimal cost spanning tree of a graph that has a designated root node " r " and satisfies the capacity constraint " c ". The capacity constraint ensures that all subtrees (maximal subgraphs connected to the root by a single edge) incident on the root node r have no more than c nodes. If the tree nodes have weights, then the capacity constraint may be interpreted as follows: the sum of weights in any subtree should be no greater than c . The edges connecting the subgraphs to the root node are called gates.

For its solution, there is the Esau-Williams heuristic, which finds suboptimal CMST that are very close to the exact solutions, but on average EW produces better results than many other heuristics. Initially, all nodes are connected to the root r (star graph) and the network's cost is ; each of these edges is a gate. At each iteration, it seeks the closest neighbor a_i for every node in $G - r$ and evaluate the tradeoff function: $t(a_i) = g_i - c_{ij}$. Then, it looks for the greatest $t(a_i)$ among the positive tradeoffs and, if the resulting subtree does not violate the capacity constraints, removes the gate g_i connecting the i -th subtree to a_i by an

edge c_{ij} . We repeat the iterations until we can not make any further improvements to the tree.

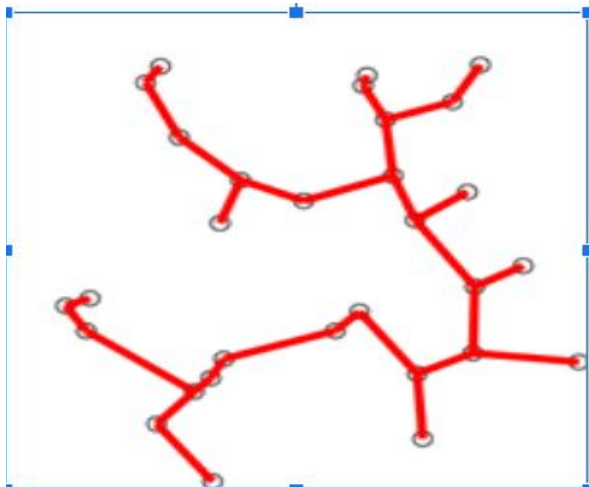


Example for the Esau-Williams Heuristic

3.4 Minimum spanning tree

A minimum spanning tree (MST) is a subset of the edges of a connected, edge-weighted undirected graph that connects all the vertices together, without any cycles and with the minimum possible total edge weight. That is, it is a spanning tree whose sum of edge weights is as small as possible. More generally, any edge-weighted undirected graph (not necessarily connected) has a minimum spanning forest, which is a union of the minimum spanning trees for its connected components.

Then, for this problem, one of its possible solutions could be Prim's (also known as Jarník's) algorithm, a greedy algorithm that finds a minimum spanning tree for a weighted undirected graph, by finding a subset of the edges that forms a tree that includes every vertex, where the total weight of all the edges in the tree is minimized. The algorithm operates by building this tree one vertex at a time, from an arbitrary starting vertex, at each step adding the cheapest possible connection from the tree to another vertex.



Representation of Prim's Algorithm

4. Decision Tree of Coffee Data

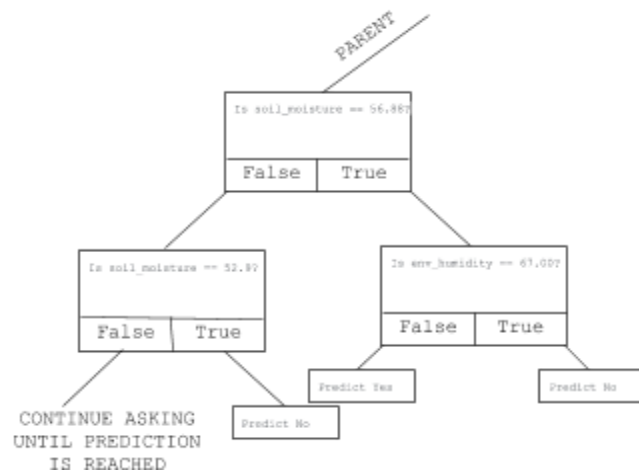
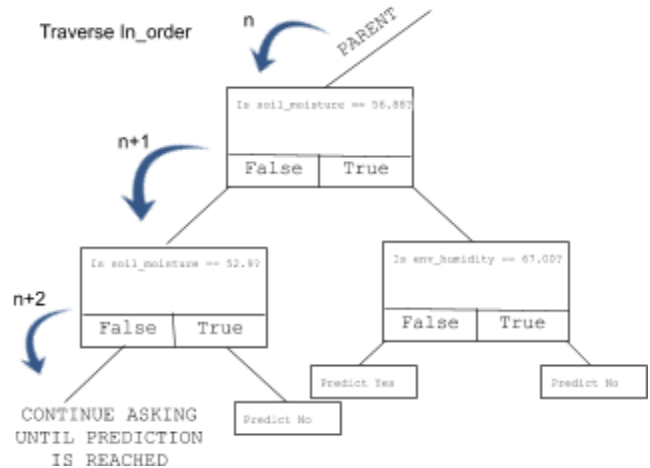


Figure 1: Tree of Coffee Data. Each node represents a type of data which will generate a hierarchy based on the type and value of the data to determine as quick as possible if the plants are sick.

4.1 Operations of the data structure

The traverse method, “traverses” each node of the Tree and counts them:

Figure 2: Traverse



The find method performs the task of searching each node for a specific input data:

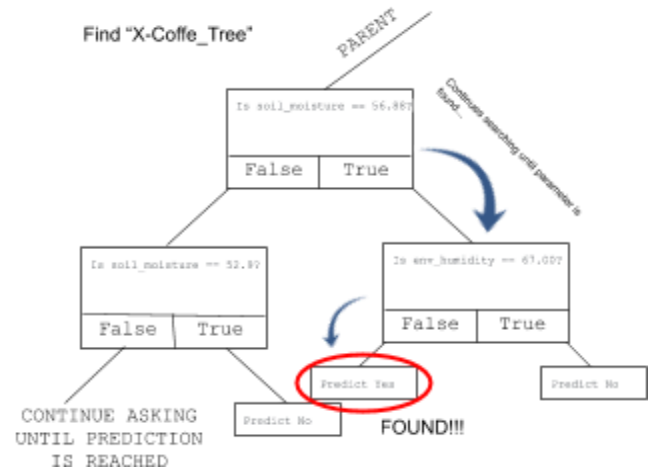


Figure 3: Find

As for the remove method, it searches for a node and deletes it based on a specific input data.

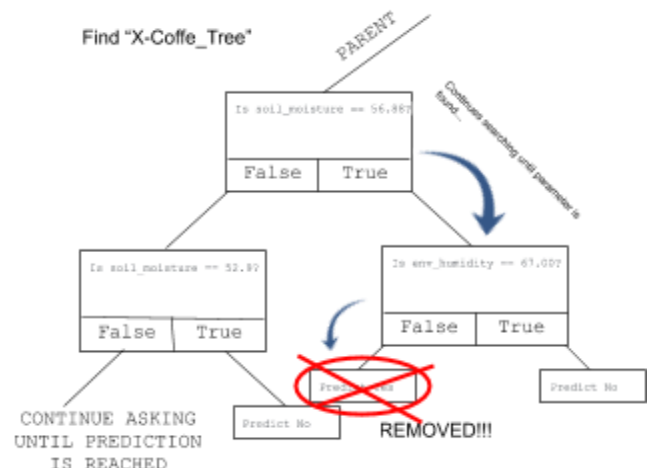


Figure 4: Remove

4.2 Design criteria of the data structure

The design of this data structure was based on the implementation of the Decision Tree data structure, from which branches are generated and so, create a hierarchy of the data collected from a coffee tree, using the importance of the type of data and whether or not its values are outside the reviewed parameters to decide as fast as possible if the tree is infected or if it's healthy. For that matter, one of the best ways to decide is by implementing a decision tree, for it considers every outcome based on the data it has and allows to take the best possible decision for a problem for which solution would be otherwise uncertain.

It's worth to add that these Trees are also pretty easy to implement, due to them being based on simple math and provide a visual representation of the data, which makes even easier the decision taking process.

4.3 Complexity analysis

The complexity for each operation can be seen in the next table:

Method	Complexity
Remove	$O(n)$
Traverse (Post-Order)	$O(n)$
Traverse (Pre-Order)	$O(n)$
Find	$O(n)$
Create Tree	$O(v \cdot \log n)$

Table 1: Table to report complexity analysis

4.4 Execution time

	Data set test 1	Data set test 2	Data set test 3
Full Algorithm (Create Tree)	35 sec	32 sec	38 sec
Remove	0.097 ms	0.073 ms	0.088 ms
Traverse	0.077 ms	0.053 ms	0.074 ms
Find	0.094 ms	0.035 ms	0.108 ms

Table 2: Execution time of the operations of the data structure for each data set.

4.5 Memory used

	Data set 1	Data set 2	Data set 3
Memory Consumption	10 MB	9.9 MB	10.7 MB

Table 3: Memory used for each data set.

4.6 Result analysis

For the data structure we designed, the worst case scenario of its runtime would be $O(v \cdot \log n)$, still it works relatively fast considering the amount of data that it has to manage. Some functions such as removing, finding or traversing take as little as $O(n)$ to fulfill their respective tasks, being this, of course, an approximation of the worst case scenario, forgiven a rather lucky value, it could take $O(1)$ to run, which is a very optimistic, even ideal and yet possible runtime for these algorithmic solutions.

Values during execution	
Task	Performance
Memory Space	10 MB
Create Tree	32 sec. - 38 sec. 35 sec. in average
Remove	0.073 ms. - 0.097 ms. 0.085 ms. in average
Traverse	0.053 ms. - 0.077 ms. 0.065 ms. in average
Find	0.035 ms. - 0.108 ms. 0.072 ms. in average

Table 4: Analysis of the results

5. Final Data Structure designed

Decision Tree of Coffee Data

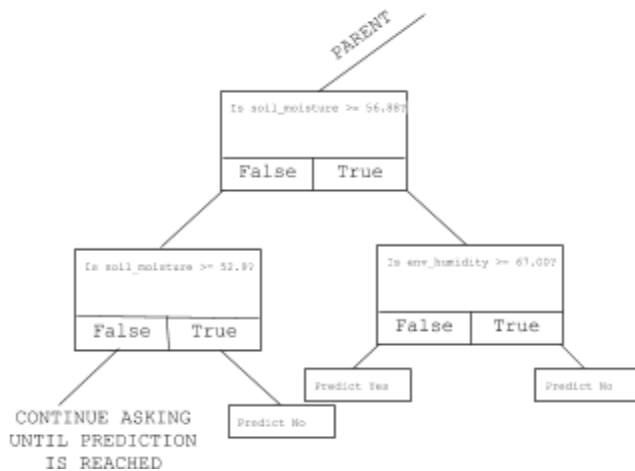


Figure 1: Tree of Coffee Data. Each node represents a type of data which will generate a hierarchy based on the type and value of the data to determine as quick as possible if the plants are sick.

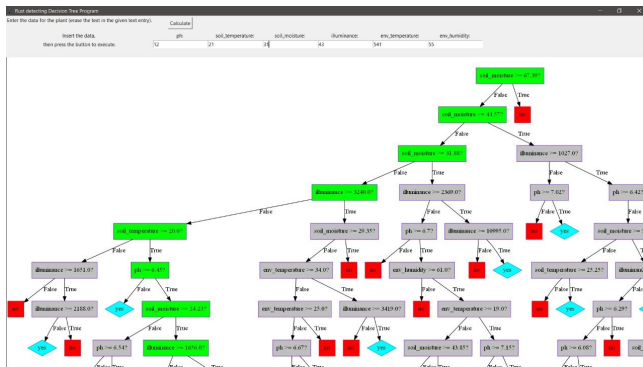


Figure 1.2 Actual footage of the application.

5.1 Operations of the data structure

The print method, shows each node of the Tree as it generates them:

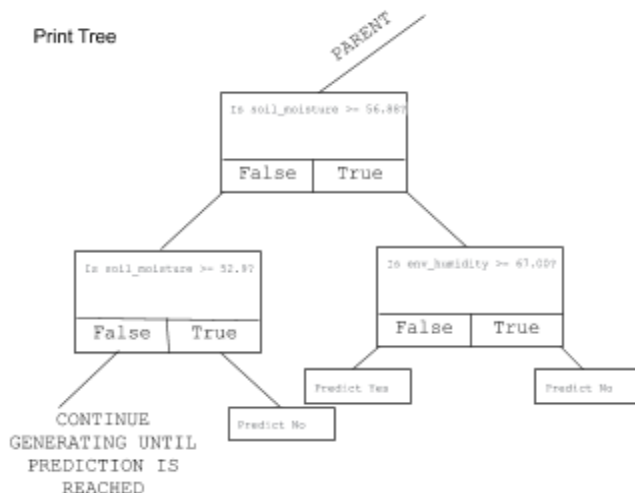


Figure 2: Print

Does the same as the find method, performs the task of searching each node for a specific input data:

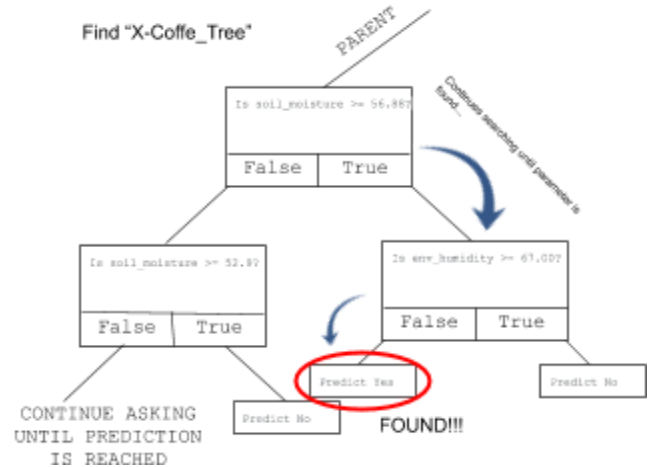


Figure 3: Classify

5.2 Design criteria of the data structure

The design of this data structure was based on the implementation of the Decision Tree data structure, from which branches are generated and so, create a hierarchy of the data collected from a coffee tree to decide as fast as possible if the tree is infected. For that matter, one of the best ways to decide is by implementing a decision tree, for it considers every outcome based on the data it has and allows to take the best possible decision for a problem for which solution would be otherwise uncertain.

It's worth to add that these Trees are also pretty easy to implement, due to them being based on simple math and provide a visual representation of the data, which makes even easier the decision taking process and most of their standard methods take $O(n)$ as worst case to execute.

For the final product, it was greatly taken into account ways to improve its performance in time and, most importantly, the accuracy of its prediction by implementing new algorithms, like C4.5.

As for the graphic interface, it was designed using python libraries (tkinter and graphviz) so it was clear and easy to use for everyone, as long as they are familiar with the data it uses as input.

5.3 Complexity analysis

The complexity for each new operation can be seen in the next table:

Method	Complexity
Print Tree	$O(2^n)$
Classify	$O(\log_2 n)$
Train (with C4.5)	$O(m \cdot n^2)$

Table 5: Table to report complexity analysis.

5.4 Execution time

	Data set test 1	Data set test 2	Data set test 3
Train	5.9 s	8.53 s	6.9 s
Print Tree	0.0014 s	0.0068 s	0.0036 s
Classify	0.076 ms	0.037 ms	0.93 ms

Table 6: Execution time of the operations of the data structure for each data set.

5.5 Memory used

	Data set 1	Data set 2	Data set 3
Memory Consumption	44.8 MB	48.5 MB	48.2 MB

Table 7: Memory used for each data set.

5.6 Result analysis

For the data structure we designed, the worst case scenario of its runtime would be $O(m \cdot n^2)$, while the function classify takes a decent $O(\log_2 n)$ to fulfill its respective task. Regardless, printing the tree can take as long as $O(2^n)$, which can be justified by the amount of data that it needs to manage. Still, with our dataset, it takes about 50MB of space and about 6 seconds of time in average. .

Values during execution	
Task	Performance
Memory Space	50 MB
Train	5.9 sec. - 8.53 sec. 6.9 sec. in average
Classify	0.037 ms. - 0.093 ms. 0.076 ms. in average
Print Tree	0.014 sec. - 0.068 sec.

	0.036 sec. in average
--	-----------------------

Table 8: Analysis of the results

6. CONCLUSIONS

In conclusion, our addressed problem, which was to predict the blight plague on coffee plantations was solvable thanks to predictive models on data structures and different implementations for each of them, that we came across during the investigation process, from which we had to decide the best one for our purpose, decision trees.

With our last design, a pretty decent accuracy was finally reached, along with a graphic interface that makes the program easier to use. The execution time and consumption of memory were optimized as well as it shows a graphic that leads to the prediction with every individual data that the tree used, from its root down to the answer “leaf”.

The reason that led us to our final solution was that almost every early attempt had a high consumption of time and little accuracy, that directed us to improve in those failing aspects, without losing other important ones from our sight. Furthermore, some methods were now useless for the purposes of this project, so they were eliminated.

6.1 Future work

For the future, our team would like to improve the accuracy of the prediction which, although it is pretty high, would be better if it was 100% accurate. Also, to implement a better, more instinctive and comfortable user interface to use our application.

ACKNOWLEDGEMENTS

This research was partially supported by institutions such as the Colombian Government and Alcaldía de Medellín, who invest in our studies.

We mainly thank for assistance in the investigation for accurate and faster solutions, and also for expressing new ideas for their implementation to students from the same course, such as Julián Bueno, Armando Ríos and Juan Diego Mejía, whose comments and critics help to improve our solution.

REFERENCES

1. Fischer, G. and Nakakoji, K. Amplifying designers' creativity with domainoriented design environments. in Dartnall, T. ed. *Artificial Intelligence and Creativity: An Interdisciplinary Approach*, Kluwer Academic Publishers, Dordrecht, 1994, 343-364.
2. Megter: Fuzzy Matching Algorithms To Help Data Scientists Match Similar Data, 2016. Retrieved August 10, 2019, from Data Science Central: <https://www.datasciencecentral.com/profiles/blogs/fuzzy-matching-algorithms-to-help-data-scientists-match-similar>
3. Naumenko, F: Fast and Practically 'perfect' Partition Problem Solution, 2018. Retrieved in August 10, 2019, from Code Project: <https://www.codeproject.com/Articles/1265125/Fast-and-Practically-perfect-Partition-Problem-Sol>
4. Kurai R., Minato S., Zeugmann T.: *Unordered N-gram Representation Based on Zero-suppressed BDDs for Text Mining and Classification*, 2007. Retrieved in August 17, 2019, from: <http://bit.ly/2Zk4I2x>
5. Wikipedia contributors. "Capacitated minimum spanning tree." *Wikipedia, The Free Encyclopedia*. Wikipedia, The Free Encyclopedia, 9 Feb. 2019. Web. Retrieved in August 10, 2019, form: <http://bit.ly/2KDfLLY>
6. Wikipedia contributors. "Minimum spanning tree." *Wikipedia, The Free Encyclopedia*. Wikipedia, The Free Encyclopedia, 11 Aug. 2019. Web. 17 Aug. 2019. Retrieved in August 10, 2019, form: <http://bit.ly/2Hev1N0>
7. Kritikos, M. & Ioannou, G. *J Oper Res Soc* (2017) 68: 1223. <https://doi.org/10.1057/s41274-016-014>