

Algorithmic Methods of Data Mining - Assignment 2

Maksad Donayorov

November 10, 2018

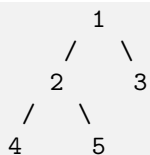
1 Problem. *Clustering given a hierarchy:*

1.1 *k-means on trees can be solved optimally in polynomial time*

To prove the k-means on tree can be solved optimally in polynomial time let's start by stating a pseudocode that solves this problem and then discuss how it works:

```
1 optimal_clusters = []
2 for node in TraversalPostorder(tree):
3     if node does not have children:
4         tuple = (1, 0, [node.index], [node.value])
5         optimal_clusters.append(tuple)
6     if node has children:
7         children = node.left_child + node.right_child
8         numb_of_clusters = children.clusters
9         node_self = (
10             numb_of_clusters,
11             calc_self_var,
12             [children.indices],
13             [children.elements]
14         )
15         tuple = (node_self, children)
16         optimal_clusters.append(tuple)
17         tuple.filter_the_optimal
18         tuple.trim_elements
```

Assume that we have a tree like this:



Then the *TraversalPostorder(tree)* will run as:

```
4 5 2 3 1
```

This allows us to travel through the tree only once. Based on the pseudocode we should not need to refer to the children once we have computed the needed parameters. The main idea in

here is that we create a *tuple* of the needed parameters once we are visiting a node or a leaf, store all of the parameters and move to the next. For instance, at line 4:

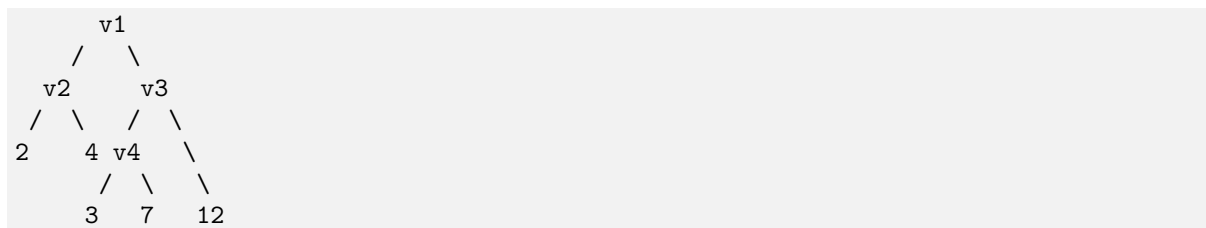
```
tuple = (1, 0, [node.index], [node.value])
```

We are creating the tuple out of a leaf. The first parameter of that *tuple* is an indicator of clusters and since a leaf can only have a cluster of one element, its default value is 1.

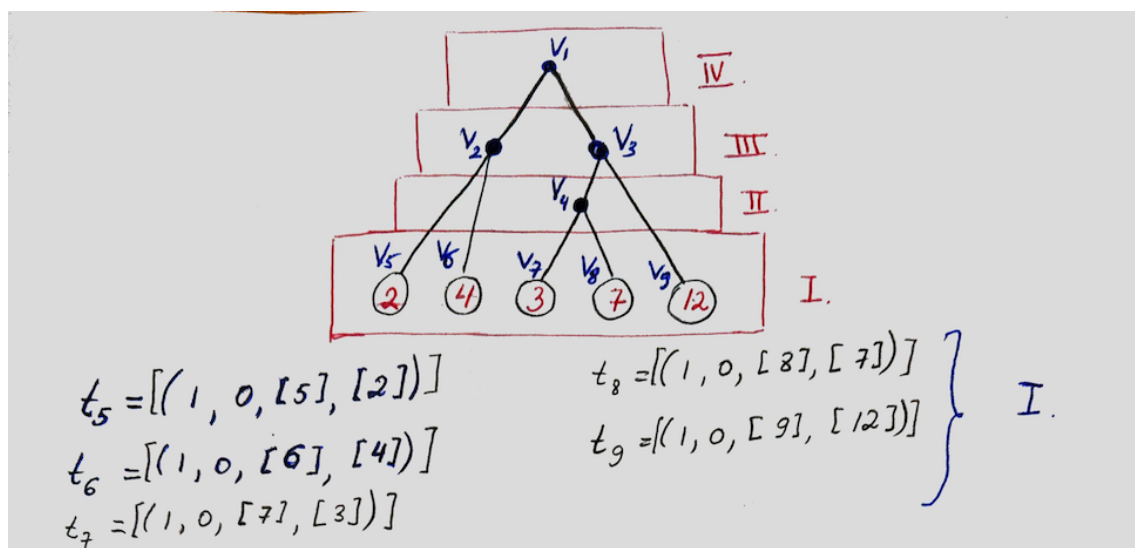
The second parameter of that tuple is 0, that indicates the variance at that point. Again, since there is no more than one element at the leaf the variance will be 0.

The third element is an array of indices, we need that the very end and I will explain this later. The fourth parameter is the an array of element that that leaf holds. This will be needed once we are in the parent node and want to compute the variance.

After we do the computation at the first leaf we move to the next leaf and eventually to the parent where we compute the parameters (later in the example these parameters are noted as tuples). This sounds very complicated and indeed it is, as there isn't a better approach I came up with. Let's look into more concrete example hope to make things more clear. Let's assume that we have this tree with given leaves:



In this tree we have $n = 5$ and the leaves are 2, 4, 3, 7, 12. We have 4 "layers" and we can compute the variance, number of clusters and number of vertices. The first thing we do is to compute the level I as that's the easiest (please note that we are not following at this point the postorder's steps). This is how it looks like:



Here we create tuples for each leaf. For example, the:

$t_5 = (1, 0, [5], [2])$

means that this is the vertex 5, that has one "cluster", it's index is 5 and the elements that it holds are [2]. With that representation we can move to the second and third "levels".

It is a bit tricky to calculate the tuples for nodes. That's because the nodes have 2 combinations:

1. When we are clustering the elements in the node (let's denote it by t_{self}).
2. When we are clustering the elements in the children nodes/leaves ($t_{left} + t_{right}$).

Let's now calculate the parameters (t_2) for *vertex*₂:

$$\begin{aligned}
 t_2 &= [t_{2self}, t_{2left} + t_{2right}] = [t_{2self}, t_5 + t_6] \\
 &= [t_{2self}, (1, 0, [5], [2]) + (1, 0, [6], [4])] \\
 &= [t_{2self}, (2, 0, [5, 6], [2, 4])] \\
 t_{2self} &= (1, t_{2var}, [2], [2, 4])
 \end{aligned}$$

Now that we have calculated the left and the right children of t_2 , we need to find the variance at t_2 that will be used for t_{self} :

to calculate the variance of t_{2var} we can use the formula

$$\underbrace{\sum_i x_i^2}_{\text{sumSquare}} - \frac{(\sum_i x_i)^2}{N} \rightarrow \text{sum}^2$$

N \hookrightarrow number of elements

$$t_{2var} = \text{sumSquare} - \frac{\text{sum}^2}{N} = 20 - \frac{36}{2} = \frac{40 - 36}{2} = 2$$

$$\begin{aligned}
 \text{sumSquare} &= 2^2 + 4^2 = 20 \\
 \text{sum}^2 &= (2+4)^2 = 36
 \end{aligned}$$

$$\textcircled{t_2} = [(1, 2, [2], [2, 4]), (2, 0, [5, 6], [2, 4])]$$

As you might have noticed we are not using the formula that is given in the problem description to calculate the variance (link to the formula). Instead, we are using a more efficient formula that relies on number of elements, sum of the elements and "square sum" of those elements.

Now we can move to the next vertex and certainly the same logic applies:

$$\begin{aligned}
t_4 &= [(t_{4self}), (t_{4left} + t_{4right})] = [t_{4self}, t_7 + t_8] \\
&= [t_{4self}, (1, 0, [7], [3]) + (1, 0, [8], [7])] \\
&= [t_{4self}, (2, 0, [7, 8], [3, 7])] \\
t_{4self} &= (1, t_{4var}, [4], [3, 7]) \\
t_{4var} &= \text{sumsquare} - \frac{\text{sum}^2}{N} = 58 - \frac{100}{2} = \frac{116 - 100}{2} = 8 \\
\textcircled{t_4} &= [(1, 8, [4], [3, 7]), (2, 0, [7, 8], [3, 7])]
\end{aligned}$$

Please pay attention to the summation of left and right children that have more than one tuple:

calculating t_3 follows the same logic:

$$\begin{aligned}
t_3 &= [t_{3self}, t_4 + t_9] \\
t_4 + t_9 &= \begin{matrix} (1, 8, [4], [3, 7]) \\ (2, 0, [7, 8], [3, 7]) \end{matrix} \rightarrow (1, 0, [9], [12]) \\
&= [(2, 8, [4, 9], [3, 7, 12]), (3, 0, [7, 8, 9], [3, 7, 12])] \\
t_{3self} &= (1, t_{3var}, [3], [3, 7, 12]) \\
t_{3var} &= \text{sumsquare} - \frac{\text{sum}^2}{N} = 202 - \frac{22^2}{3} = 40,6 \\
\textcircled{t_3} &= [(1, 40,6, [3], [3, 7, 12]), (2, 8, [4, 9], [3, 7, 12]), (3, 0, [7, 8, 9], [3, 7, 12])]
\end{aligned}$$

Up to this point we have computed all the parameters (tuples) for all the children. Now, it's time to calculate them for $vertex_1$.

Now we have to compute the last level - 1:

$$t_1 = [t_{self}, t_2 + t_3]$$

$$t_2 + t_3 = \begin{matrix} (1, 2, [2], [2, 4]) \\ (2, 0, [5, 6], [2, 4]) \end{matrix} \begin{matrix} \rightarrow (1, 40.6, [3], [3, 7, 12]) \\ \rightarrow (2, 8, [4, 9], [3, 7, 12]) \\ \rightarrow (3, 0, [7, 8, 9], [3, 7, 12]) \end{matrix}$$

$$\checkmark (2, 42.6, [2, 3], [2, 3, 4, 7, 12])$$

$$= \checkmark (3, 10, [2, 4, 9], [2, 3, 4, 7, 12])$$

$$\checkmark (4, 2, [2, 7, 8, 9], [2, 3, 4, 7, 12])$$

$$\times (3, 40.6, [3, 5, 6], [2, 3, 4, 7, 12])$$

$$\times (4, 8, [4, 5, 6, 9], [2, 3, 4, 7, 12])$$

$$\checkmark (5, 0, [5, 6, 7, 8, 9], [2, 3, 4, 7, 12])$$

we don't take these two since their variance is too big

$$t_{self} = (1, t_{var}, [1], [2, 3, 4, 7, 12])$$

$$t_{var} = \text{SumSquare} - \frac{\text{sum}^2}{N} = (2^2 + 3^2 + 4^2 + 7^2 + 12^2) - \frac{2+3+4+7+12}{5} = \frac{110 - 784}{5} = 65.2$$

A very important point in the above calculation is how we choose some of the tuples. As you can see, those tuples that we choose are marked with the blue tick and those which are not chosen marked by red cross. This implies that we have 2 clusters with the same value of k and we only choose the optimal one, which is a tuple with the smaller variance. In the pseudocode that's done on line 16.

The final version of the t_1 looks like this:

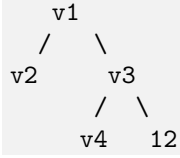
$$t_1 = \begin{bmatrix} (1, 65.2, [1], [2, 3, 4, 7, 12]) \\ (2, 42.6, [2, 3], [2, 3, 4, 7, 12]) \\ (3, 10, [2, 4, 9], [2, 3, 4, 7, 12]) \\ (4, 2, [2, 7, 8, 9], [2, 3, 4, 7, 12]) \\ (5, 0, [5, 6, 7, 8, 9], [2, 3, 4, 7, 12]) \end{bmatrix}$$

The final result that our algorithm returns is a list of the tuples with all the possible combination of clusters where $1 \leq k \leq 5$. Looking at this table we can immoderately tell what is the optimal cluster, what kind of vertices it will have and what is its variance (the elements are that each tuple holds are trimmed at line 17 of the pseudocode):

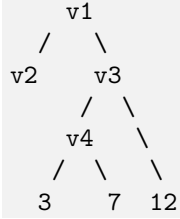
```
[
  (1, 65.2, [1],      ),
  (2, 42.6, [2, 3],   ),
  (3, 10,   [2, 4, 9], ),
  (4, 2,   [2, 7, 8, 9], ),
  (5, 0,   [5, 6, 7, 8, 9])
]
```

For example, if $k = 3$, we can say that the optimal cluster will have variance of 10 with vertices

indices of [2, 4, 9]:



Or if $k = 4$, then the optimal cluster will have a variance of 2 with vertices [2, 7, 8, 9]:



1.2 correctness of this algorithm

As we have seen the algorithm returns only the optimal "tuples" with the optimal number of clusters. Therefore, it will always filter the tuples that have a bigger variance (see line 17 of the pseudocode) before moving to the next step. For example, suppose that we have $k = 3$ and the list of tuples looks like this:

```
(3, 10, [2, 4, 9], [2, 3, 4, 7, 12])
(3, 40.6, [3, 5, 6], [2, 3, 4, 7, 12])
```

Then the algorithm will return the t with the smallest variance, which is in this case 10:

```
(3, 10, [2, 4, 9], [2, 3, 4, 7, 12])
```

For detailed explanation of the correctness of this algorithm, please look at the at **section 1.1**.

1.3 complexity of the algorithm

There are 3 main component to this algorithm that we have to take into consideration when calculate the complexity of it.

1. The *Postorder Traversal* algorithm that we are using, which has complexity of $\mathcal{O}(n)$.
2. Number of computation of tuples t .
3. The variance for each t .

Let's start with the variance. We are using a formula that has a complexity of $4n$:

$$\sum_i x_i^2 - \frac{(\sum_i x_i)^2}{N}$$

We are using that formula for each tuple. The number of times t is calculated equals to $2n - 1$. n times for the leaves and $n - 1$ times for the parent nodes.

So in total, t and *variance* with the *Postorder Traversal algorithm* are calculated $n(4n^2 - 4n) = 4n^3 - 4n^2$ times.

Consequently, the complexity of this algorithm is $\mathcal{O}(n^3)$.

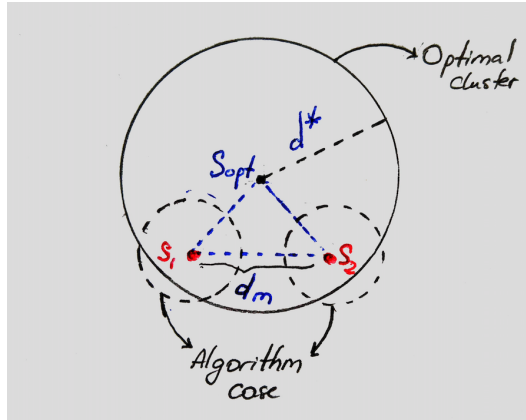
2 Clustering a data stream

2.1 Prove of SF algorithm produces a clustering that has at most k cluster centers.

We can prove that the Streaming-Furthest algorithm produces a clustering that has at most k cluster centers, by contradiction. Let's assume that the algorithm produces more clusters than the optimal, which can be denoted by $p > k$. Where k is the number of clusters that was produced by the optimal and p is the number cluster that were produced by the algorithm. In other words:

$$\begin{aligned} C^* &= \{C_1, C_2, C_3 \dots C_k\} - \text{optimal case} \\ Ac &= \{Ac_1, Ac_2, Ac_3 \dots Ac_p\} - \text{algorithm case} \\ &\text{where } p > k \end{aligned}$$

In this case there will be at least two centers: $\{s_1, s_2\} \in Ac$, that lie inside one of the cluster found in the optimal C^* . The representation and depiction of this statement will look like this:



From here we can make our analogy that:

$$d(s_1, s_2) \leq d(s_1, d^*) + d(s_2, d^*) \leq 2d^*$$

From the above inequality we can see that $d(s_1, s_2) \leq 2d^*$ and based on the algorithm this cannot happen. Because, in the algorithm the cluster centers are assigned with the condition that $d_m > 2d^*$. Consequently, we prove by contradiction that the Streaming-Furthest algorithm produces a clustering that has at most k clusters.

2.2 Prove of SF algorithm is still a factor-2 approximation.

2.3 Modification of the SF algorithm.