

Weekly Journal

Leila Uy

June 21, 2021

1 Work Update

This journal is going to be short because I have focused a lot of my time on the Literature Review.

1.1 Code Base

I spent some time on the code base for a serial k-means clustering algorithm in an R notebook. Since our last meeting on Friday, I have not touched the code because I wanted to finish my Literature Review first. Jishnu gave me an interesting article which provided pseudocode that will be useful when we implement our algorithm.

Figure 1: Pseudocode for a serial k-means clustering from [1]

```
Result:  $C = \{C_1, \dots, C_k\}$ ,  $c_j, j \in 1, \dots, k$ 
1  $c_j = \text{random } x_i \text{ in } \mathcal{X}, \quad j = 1, \dots, k, \quad c_j \neq c_i \quad \forall i \neq j;$ 
2 do
3    $C_j = \emptyset, \quad j = 1, \dots, k;$ 
4   foreach  $x_i \in \mathcal{X}$  do
5      $j = \text{argmin} D(c_j, x_i);$ 
6      $C_j = C_j \cup x_i$ 
7   end
8   foreach  $c_i \in C$  do
9      $c_i = \frac{1}{|C_j|} \sum_{x_i \in C_j} x_i;$ 
10  end
11 while convergence;
```

Figure 2: Pseudocode for a parallel k-means clustering from [1]

```
Result:  $C = \{C_1, \dots, C_k\}$ ,  $c_j, j \in 1, \dots, k$ 
1 if threadID = 0 then
2    $c_j = \text{random } x_i \text{ in } \mathcal{X}, \quad j = 1, \dots, k, \quad c_j \neq c_i \quad \forall i \neq j;$ 
3 end
4 synchronize threads;
5 do
6   foreach  $x_i \in \mathcal{X}_{\text{threadID}}$  do
7      $l_i = \text{argmin} D(c_j, x_i);$ 
8   end
9   synchronize threads;
10  if threadID = 0 then
11    foreach  $x_i \in \mathcal{X}$  do
12       $c_l = c_l + x_i;$ 
13       $m_{l_i} = m_{l_i} + 1;$ 
14    end
15    foreach  $c_j \in C$  do
16       $c_j = \frac{1}{m_j} c_j;$ 
17    end
18    if convergence then
19      signal threads to terminate;
20    end
21  end
22 while convergence;
```

Currently the code uses for loops in R which I want to change to apply because it makes for more readable code and since apply does the for loops in C, it is slightly faster. I am also going to terminate my single processor instance on EC2 because there is constant charges for the EBS storage even when the instance is stopped and I am going to move up to a 4 processor instance (m4.xlarge).

For the assignment of each row of the data frame (pixel of the raster) to a cluster, I am thinking of using a format of dict{clusternumber: [list of pixels]} because it is easier to visualize with n number of

variables. I think the important thing while implementing the clustering algorithm is to make sure that we allow for the flexibility of n number of variables and k number of clusters.

2 Literature Review

Rather than looking for new articles for my Literature Review, I looked at the citations of my old articles to expand upon certain points. For example, an old article I had from a previous week [3] talked about k-means clustering in great lengths and even went on to explain the different extensions of k-means clustering. I know from the article the summary of the different extensions and what makes them unique like fuzzy c and kmeans++ but I wanted to go back to find the articles and skim through the papers. This means that I am able to cite them if I need to and I know if the algorithms are suitable for us to implement in the future [2, 4, 5].

I am finished my Literature Review and Diljot, Jishnu, Valerie and I are sending our Literature Reviews to each other for editing on Tuesday.

References

- [1] Salvatore Cuomo, Vincenzo De Angelis, Gennaro Farina, Livia Marcellino, and Gerardo Toraldo. A gpu-accelerated parallel k-means algorithm. *Computers & electrical engineering*, 75:262–274, 2019.
- [2] James Churchill Dunn. A fuzzy relative of the isodata process and its use in detecting compact well-separated clusters. *Journal of cybernetics*, 3(3):32–57, 1973.
- [3] Anil Jain. Data clustering: 50 years beyond k-means. *Pattern recognition letters*, 31(8):651–666, 2010.
- [4] Dan Pelleg and Andrew Moore. Accelerating exact k-means algorithms with geometric reasoning. In *Proceedings of the Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 277–281. Association for Computing Machinery, 1999.
- [5] Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural computation*, 10(5):1299–1319, 1998.