# CSC-691 — Data Mining
# Assignment 4

Esteban Murillo Burford

Thursday 24ᵗʰ October, 2019

## Summary

For this assignment, everything related to the *k-Means* algorithm was implemented successfully. This time, a consensus has been reached between the two different versions of the algorithm (the one self-implemented and the one from the *sklearn* library), meaning that the generated clusters are almost always the same for both. However, it has been noted that in a small amount of scenarios, the elements belonging to the clusters vary slightly due to the nature of the algorithm. Either way, the results are satisfactory. Refer to Figure 1 for details about the program's output.

Figure 1: Final results for both *k-Means* algorithms (custom and *sklearn*)

# *k-Means* comparison

As it can be seen in Figure 2, the results are very similar for both implementations of the *k-Means* algorithm. All scores are within are reasonable range, save for $k = 5$ & $k = 6$, which differ a little bit between both algorithms. To see the actual values used for plotting, refer to Figure 1.
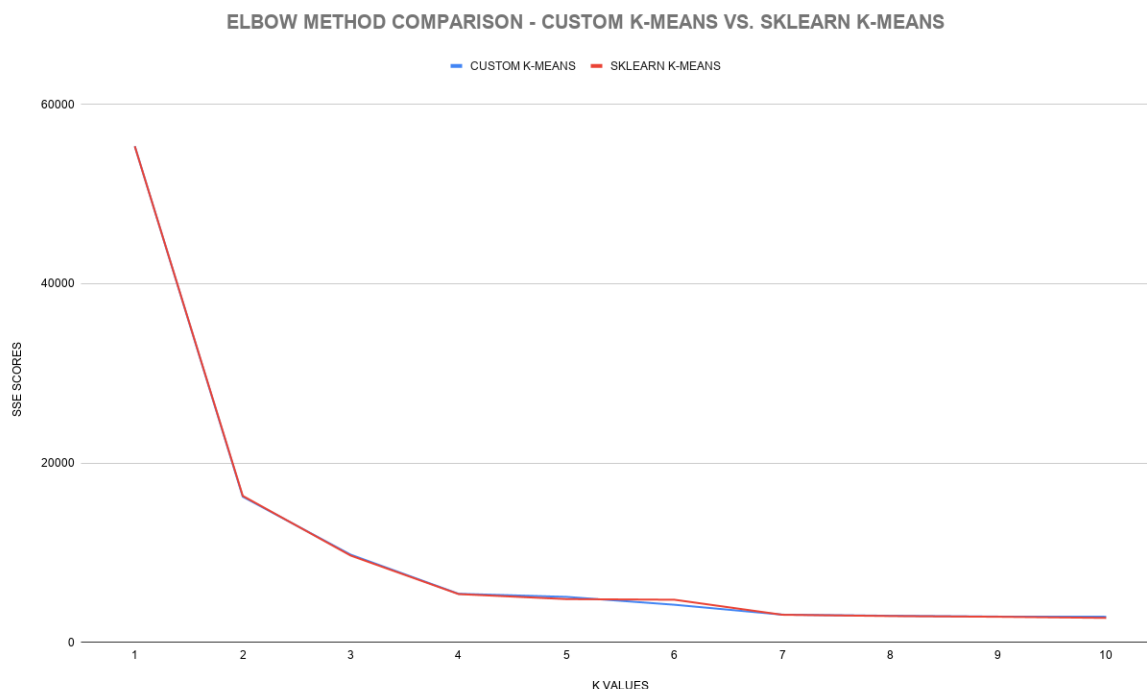


Figure 2: *k-Means* elbow method comparison — self-implemented vs. *sklearn*

# Notes

- To modify the behavior of the program, change values in *config.py* accordingly

- Remember to run code using the **version 3 of the *Python* interpreter**

# References

[1] *Beyond the k-Means – the Right k.* Retrieved on October 24th from
https://www.edupristine.com/blog/beyond-k-means