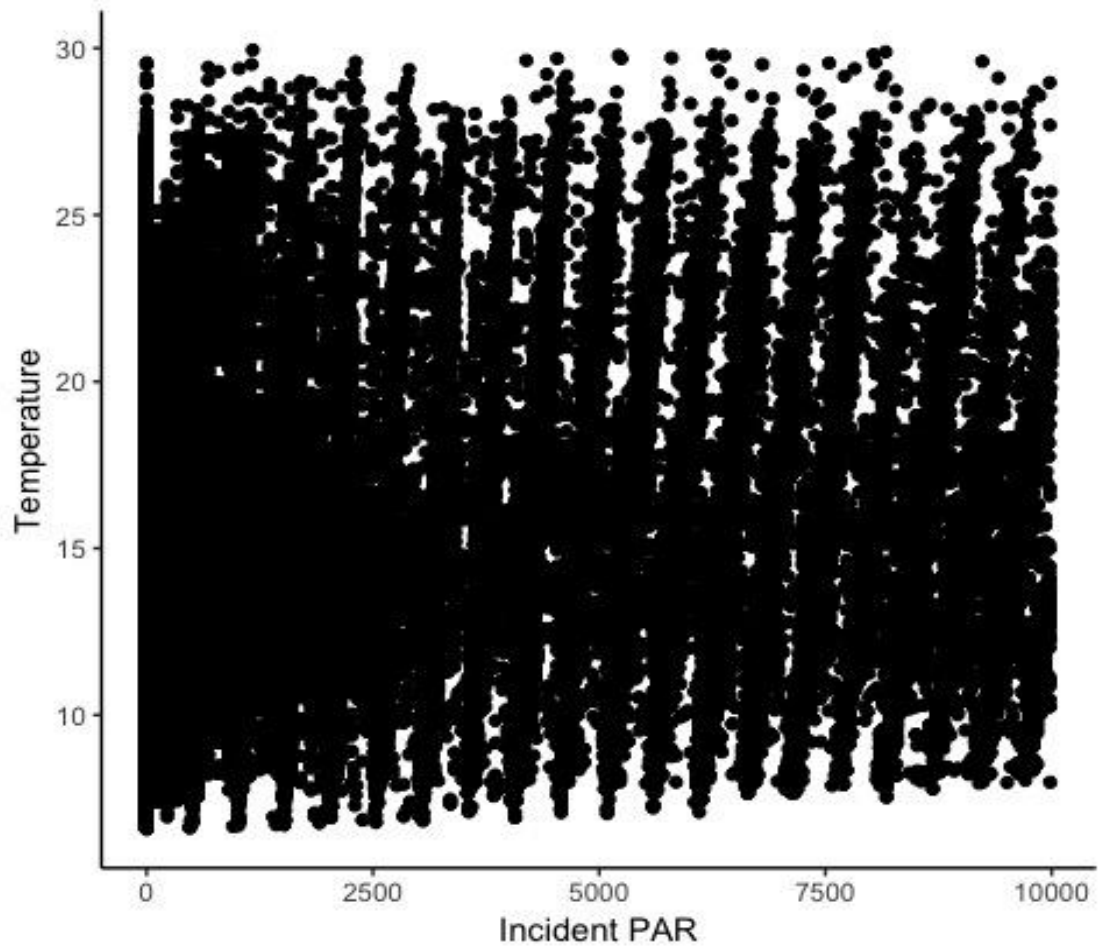

Stat 215A - Week 5

Clustering

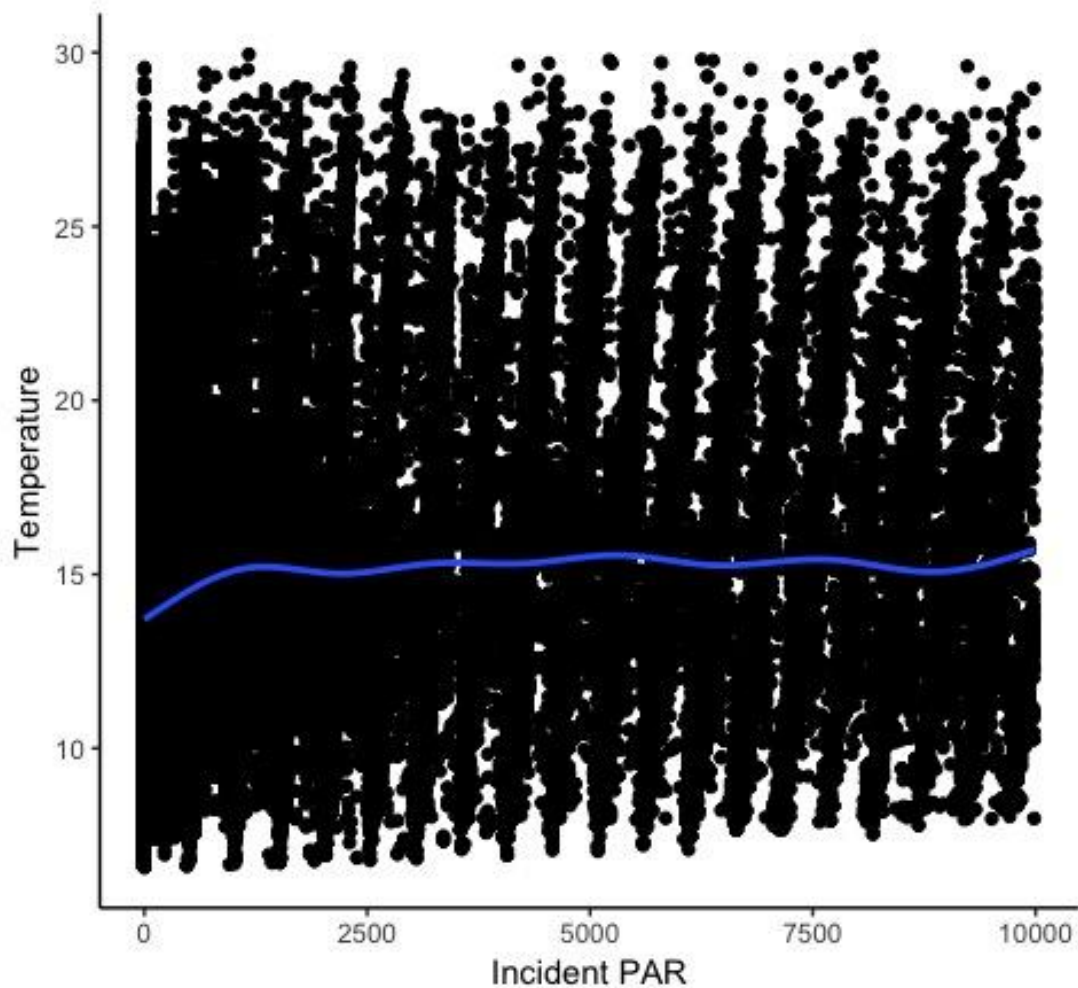
Lab 1 Discussion

Overplotting



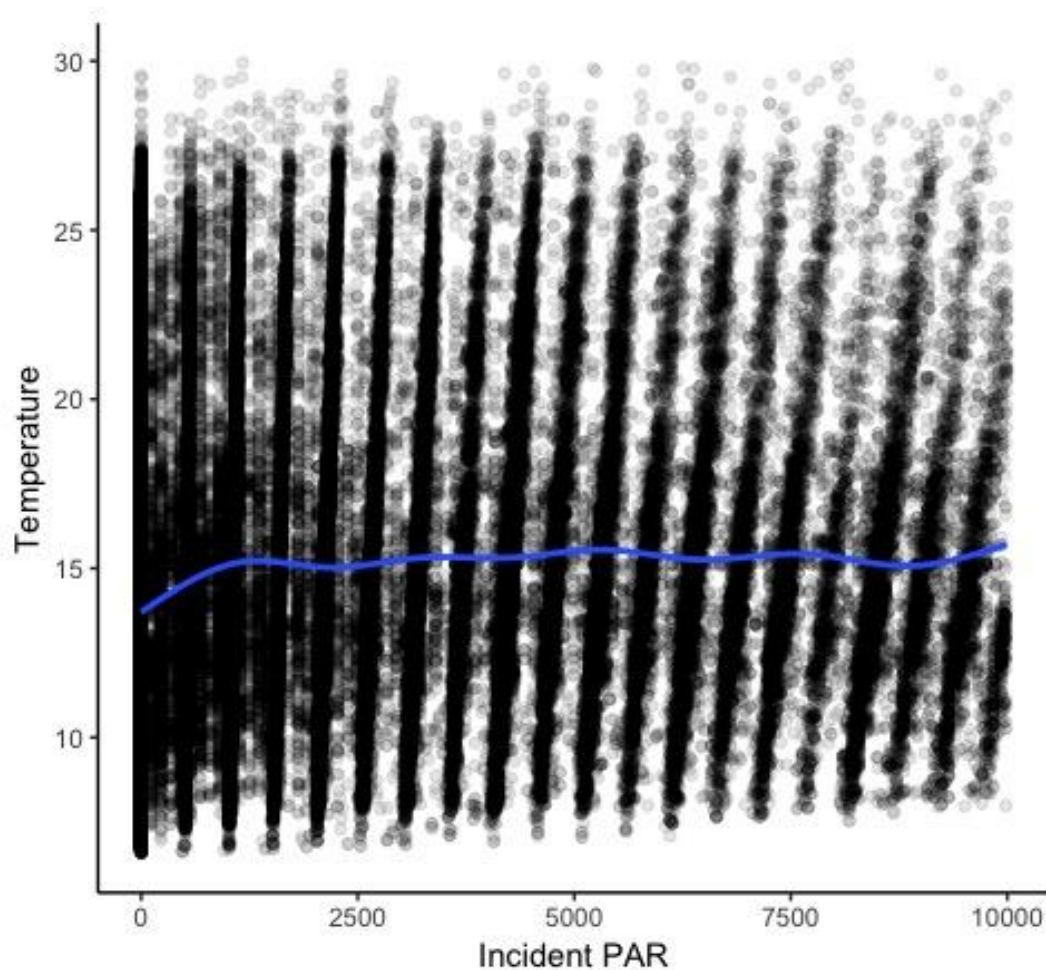
Overplotting

Add trendline



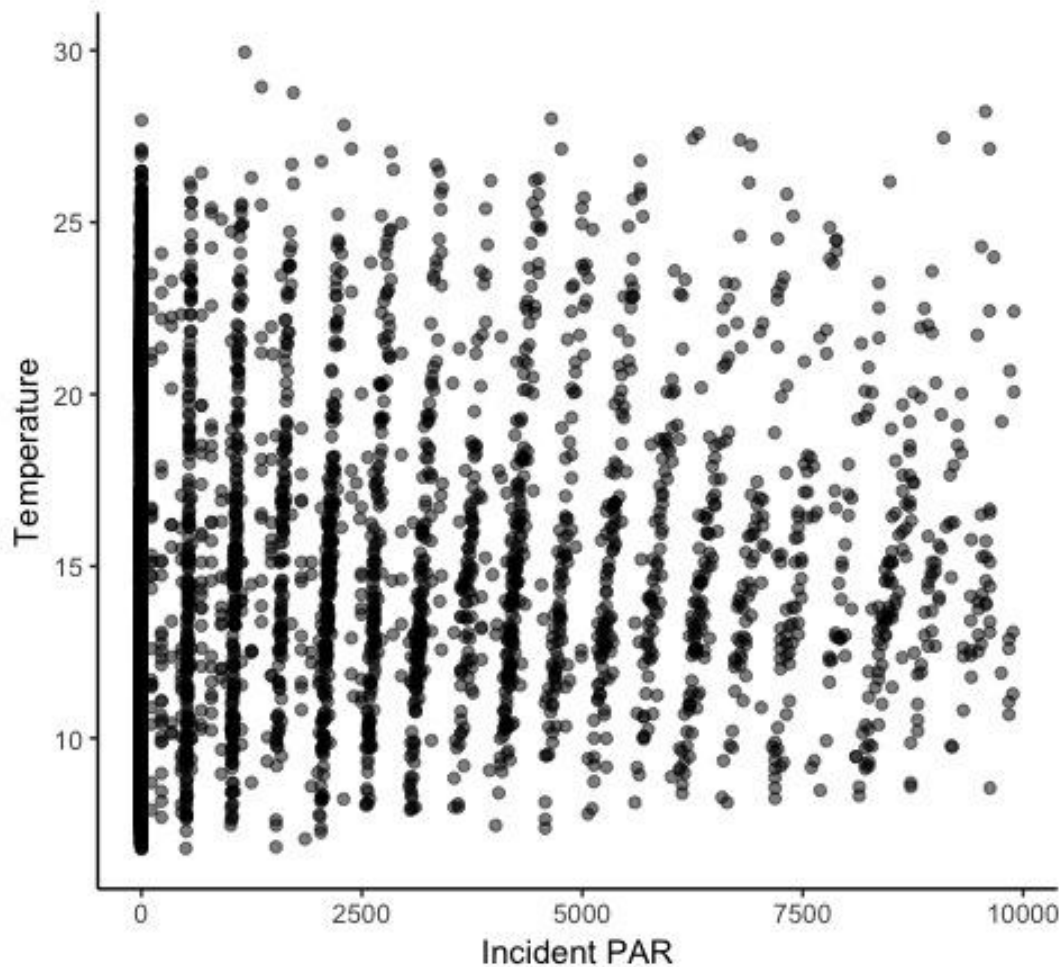
Overplotting

Use transparency



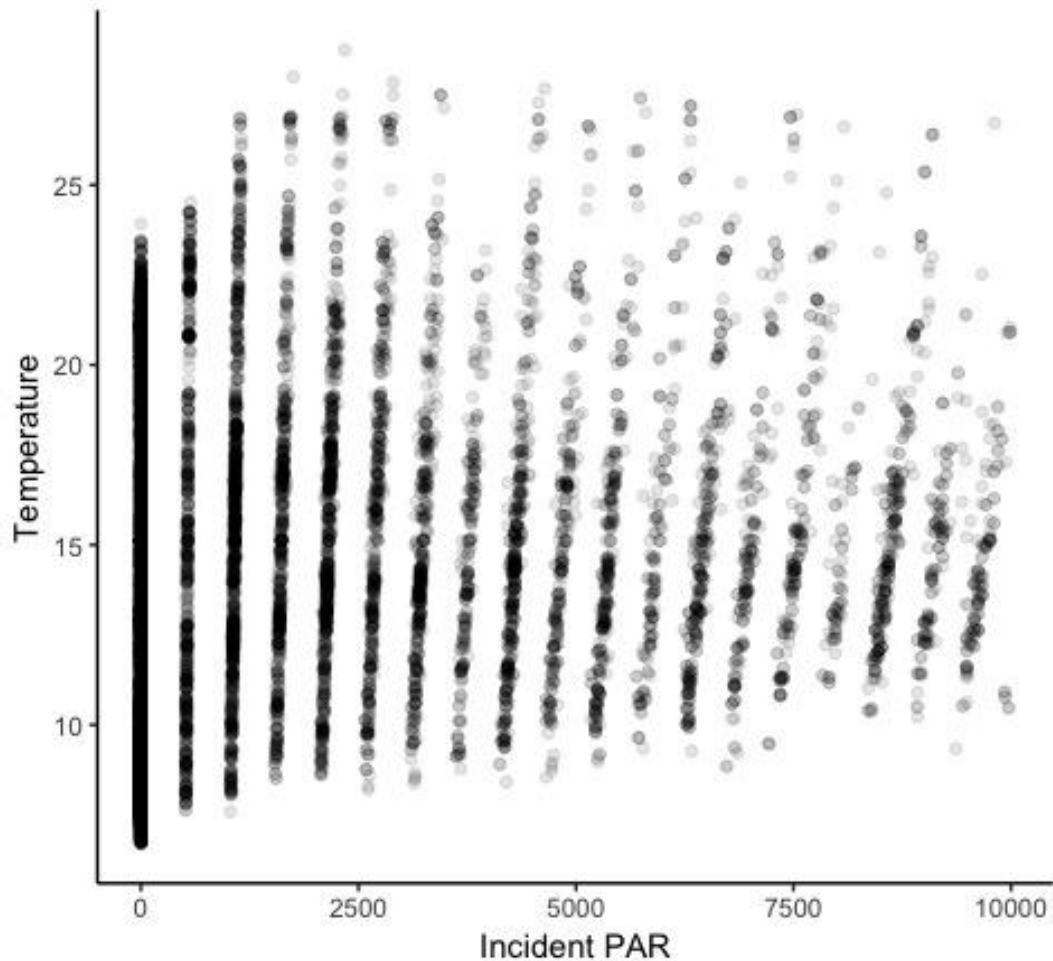
Overplotting

Subsample 10,000
points

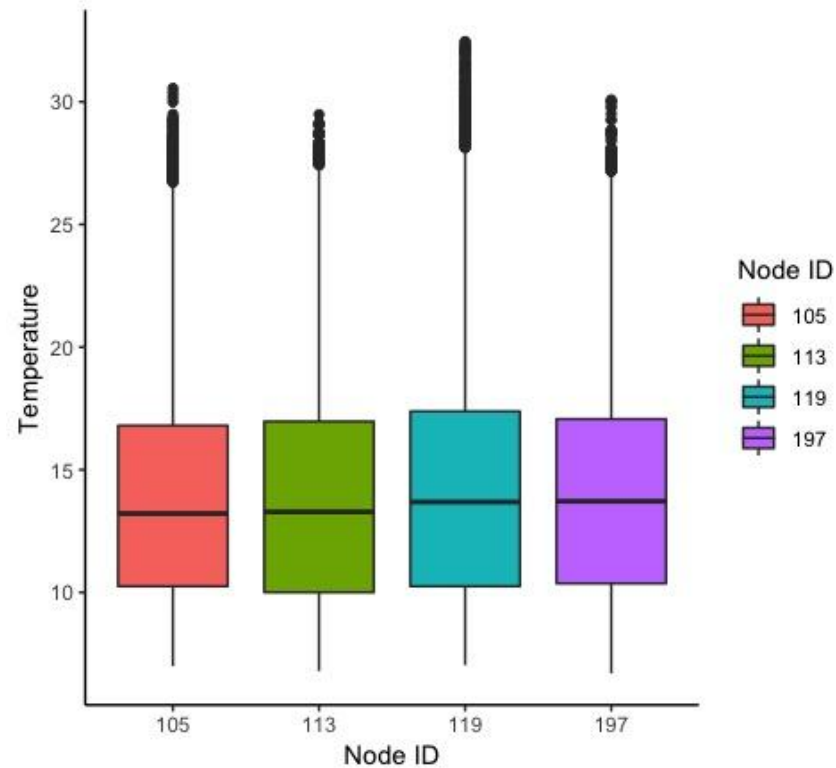


Overplotting

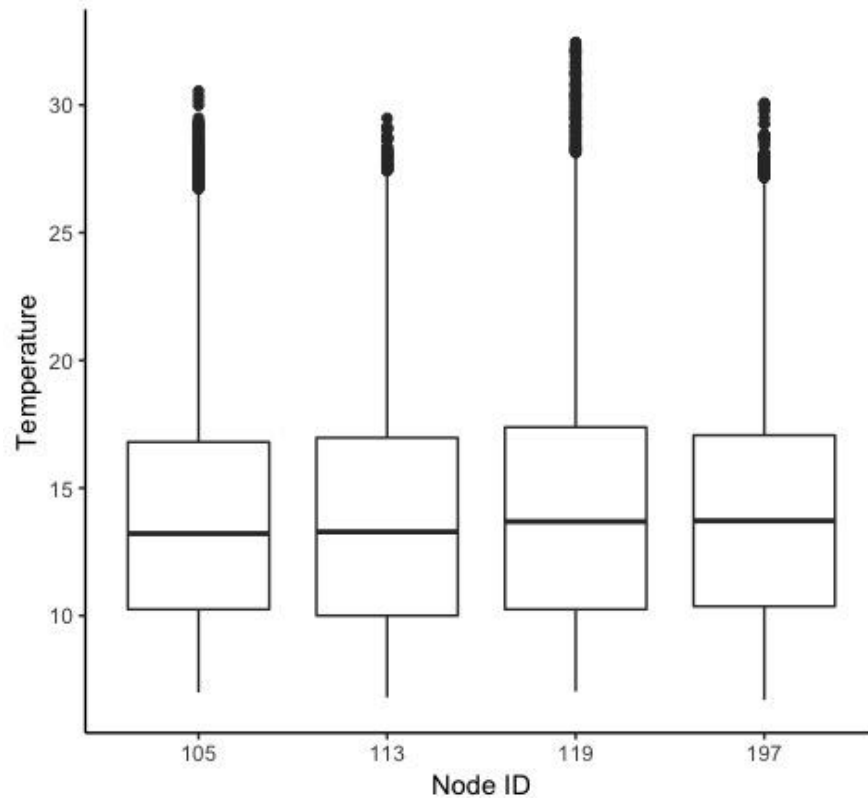
Subsample a single node



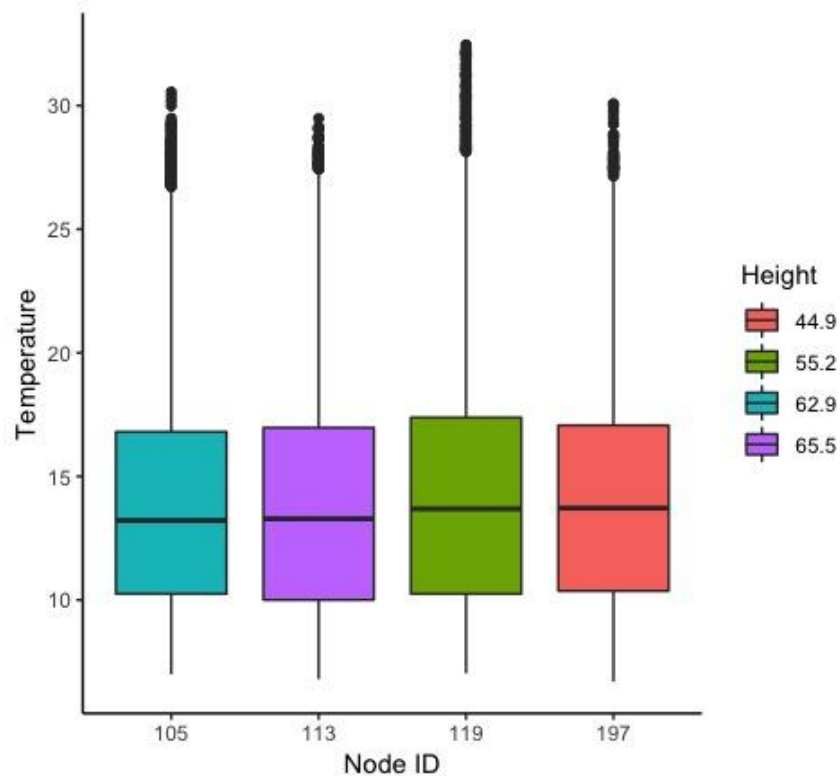
Using color appropriately



Using color appropriately

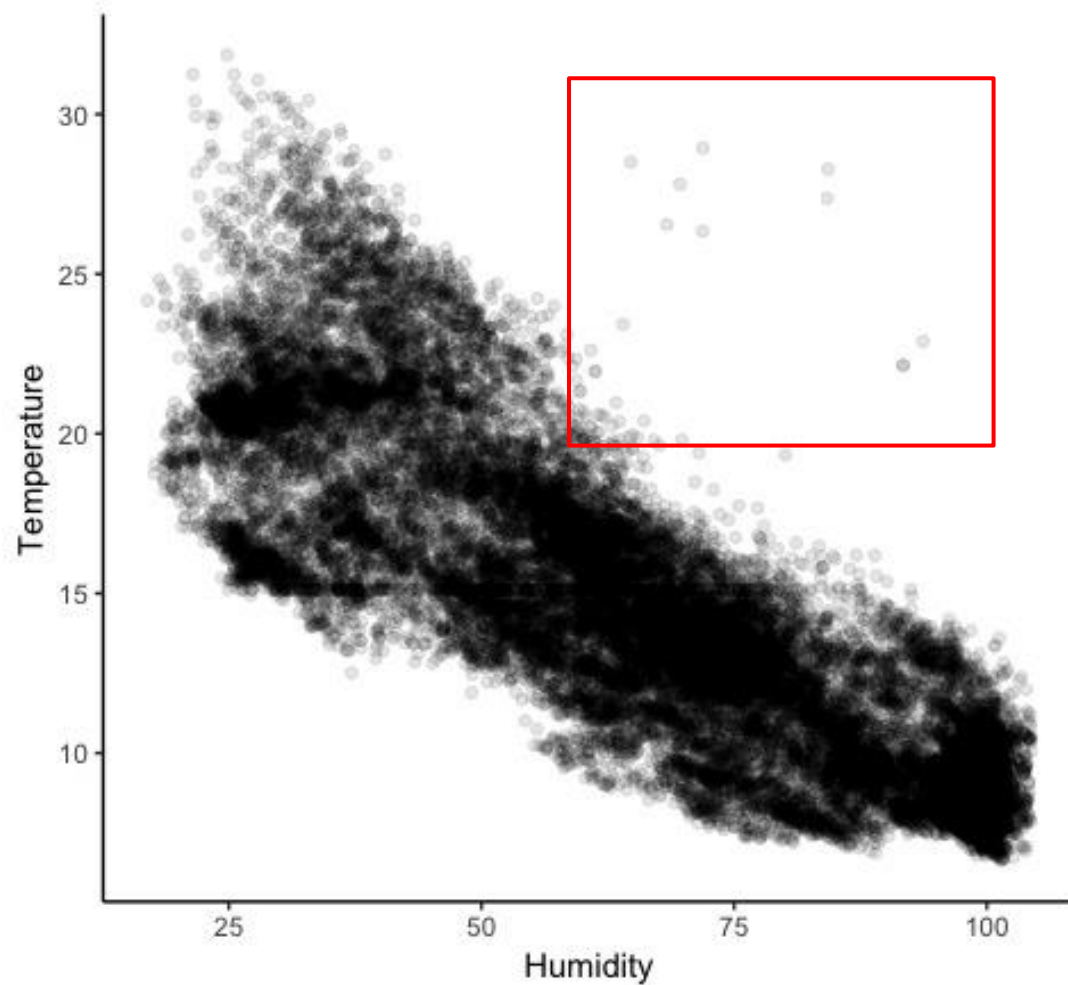


Using color appropriately



Data cleaning

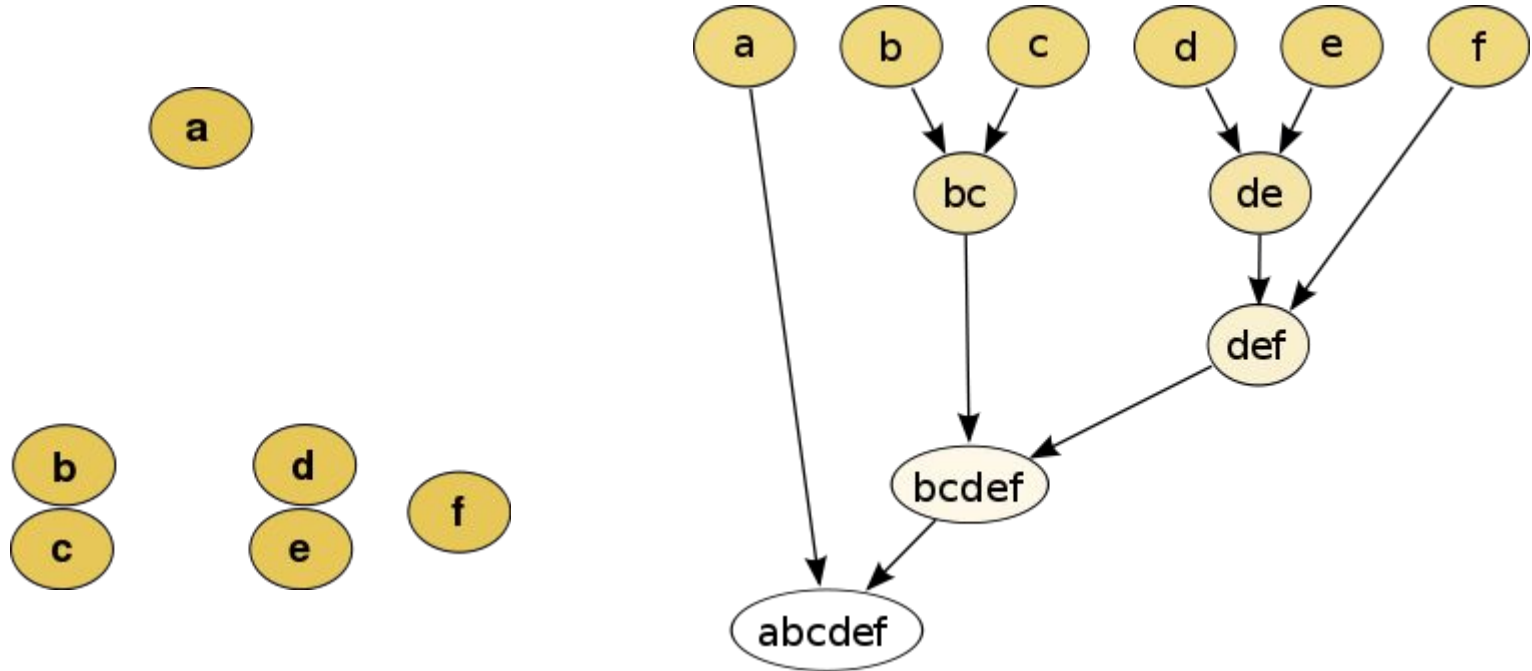
Problematic to only clean by apply bounds to the variables.



Clustering

K-means, spectral clustering, hierarchical clustering

Hierarchical clustering



Hierarchical clustering

Pros:

- ❑ Does not require choosing number of clusters
- ❑ Not sensitive to choice of distance metric
- ❑ Good when underlying data has a hierarchical structure

Cons:

- ❑ Slower than k-means

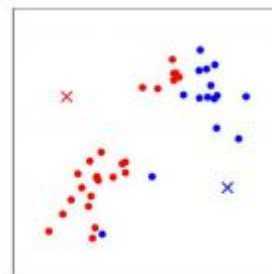
K - means



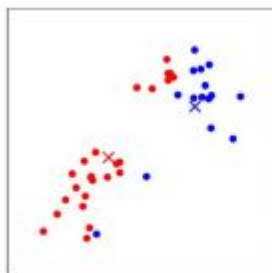
(a)



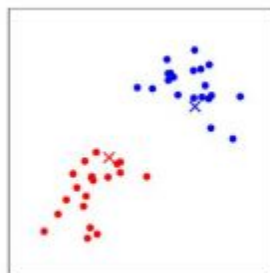
(b)



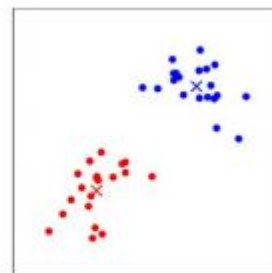
(c)



(d)



(e)



(f)

K-means: pros and cons

Pros:

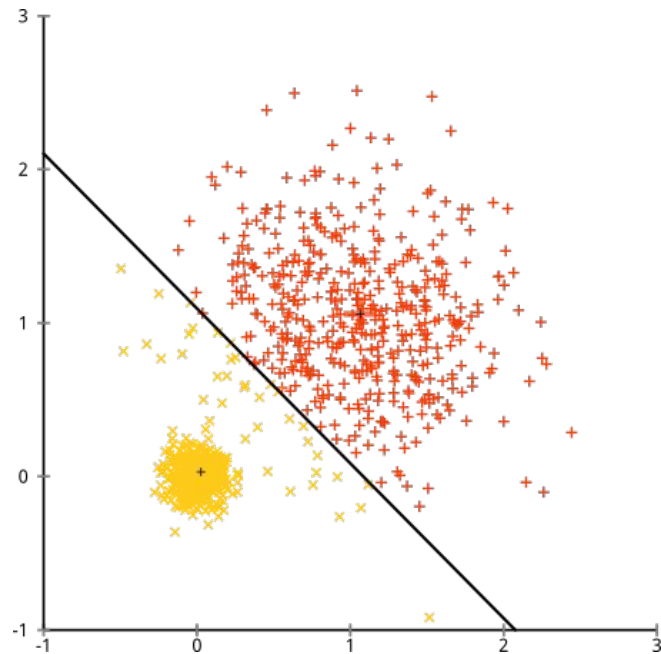
- ❑ Linear complexity: $O(n)$

Cons:

- ❑ Have to manually choose the number of clusters
- ❑ Randomly choosing starting points for clusters
- ❑ Assumes the variance of each variable is spherical
- ❑ Assumes all variables have the same variance

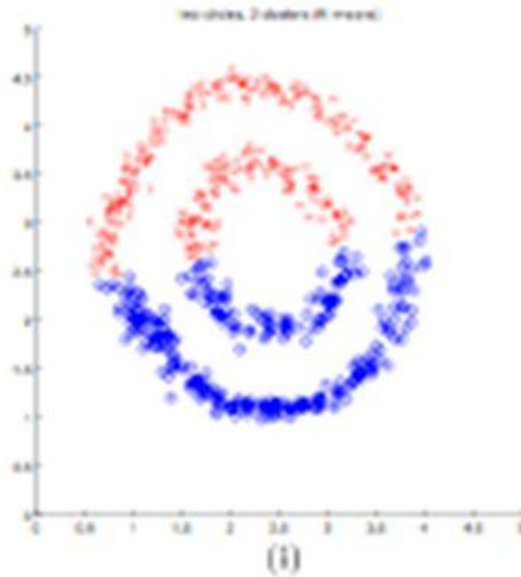
<https://towardsdatascience.com/the-5-clustering-algorithms-data-scientists-need-to-know-a36d136ef68>

K-means: different variances

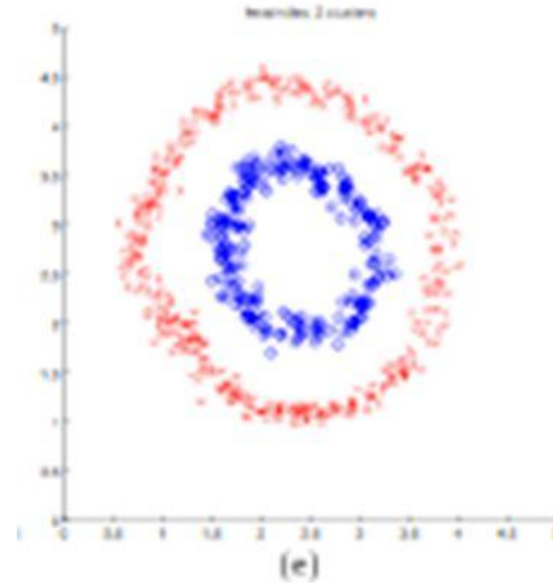


<https://stats.stackexchange.com/questions/133656/how-to-understand-the-drawbacks-of-k-means>

Spectral clustering



K-means



Spectral Clustering

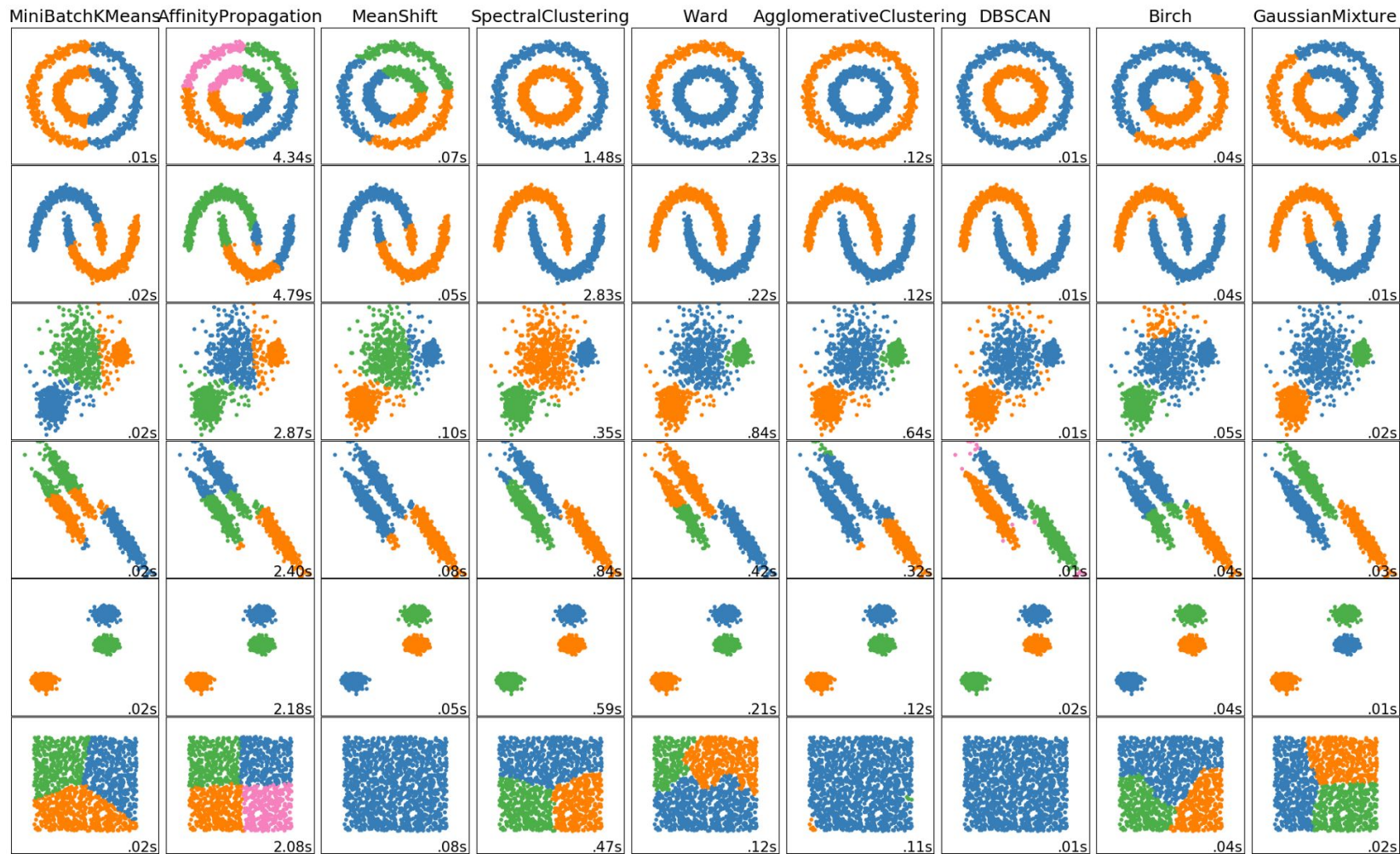
Spectral clustering

Pros:

- ❑ Good for clustering data with more complex shapes
- ❑ Reasonably fast if the data is sparse
- ❑ Produces good results

Cons:

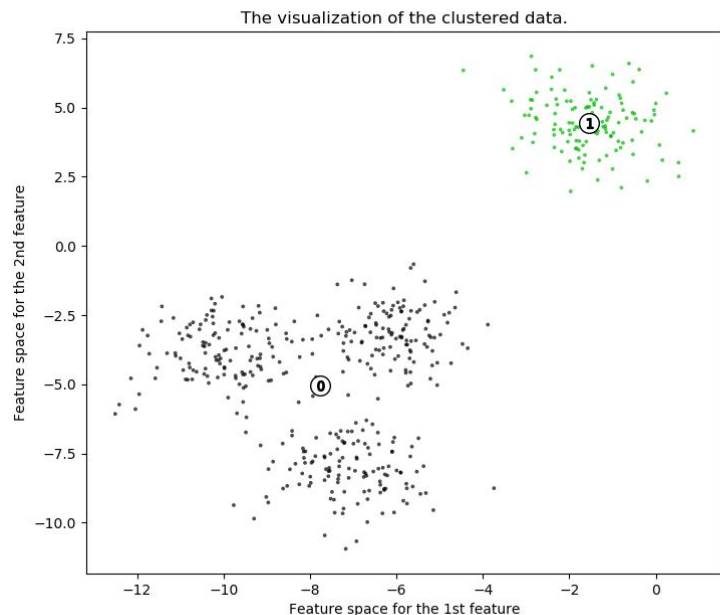
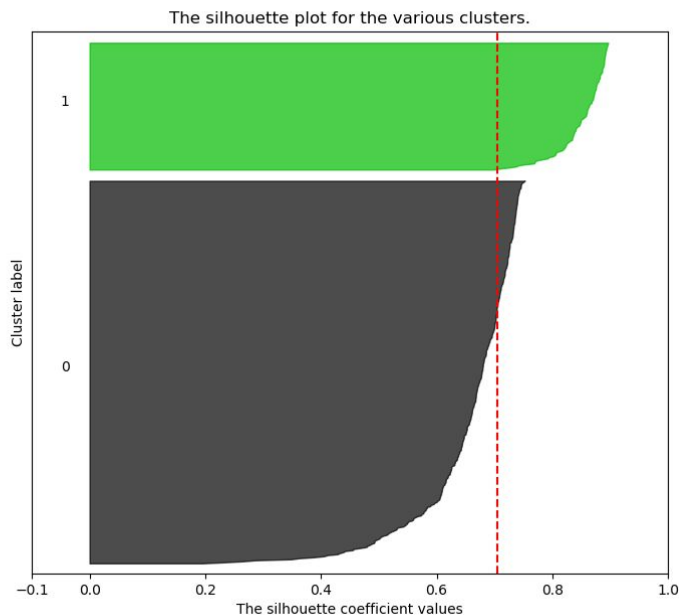
- ❑ Have to manually choose the number of clusters
- ❑ Slow in large datasets



Silhouette plots

Silhouette plots

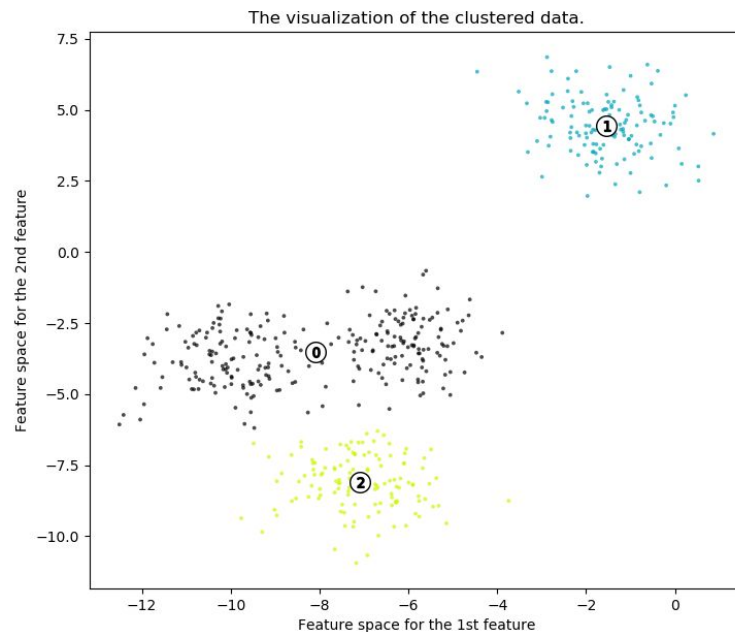
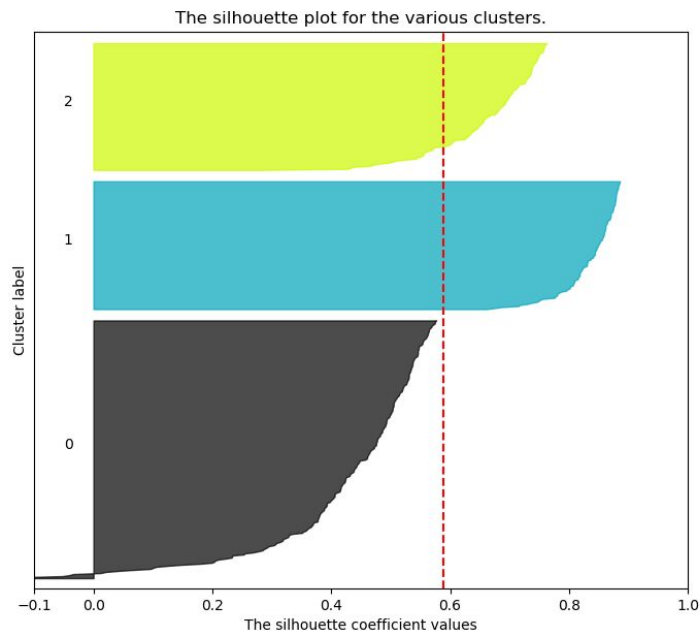
Silhouette analysis for KMeans clustering on sample data with $n_clusters = 2$



http://scikit-learn.org/stable/auto_examples/cluster/plot_kmeans_silhouette_analysis.html

Silhouette plots

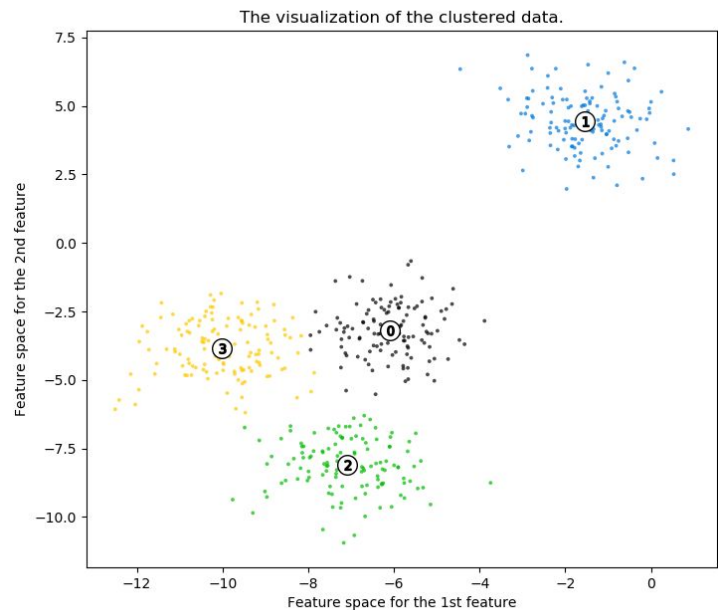
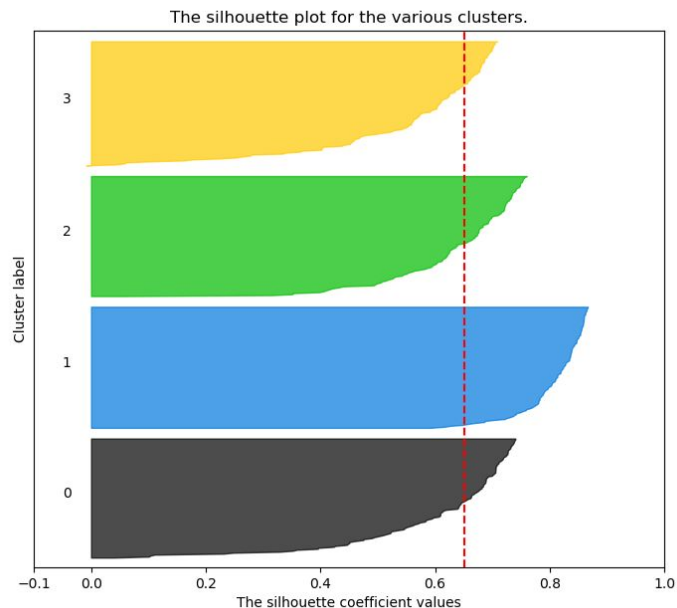
Silhouette analysis for KMeans clustering on sample data with $n_clusters = 3$



http://scikit-learn.org/stable/auto_examples/cluster/plot_kmeans_silhouette_analysis.html

Silhouette plots

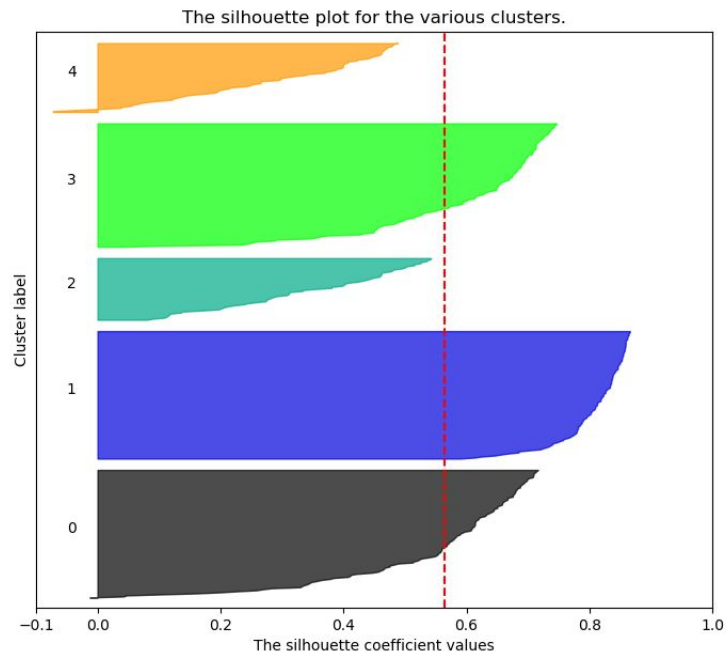
Silhouette analysis for KMeans clustering on sample data with $n_clusters = 4$



http://scikit-learn.org/stable/auto_examples/cluster/plot_kmeans_silhouette_analysis.html

Silhouette plots

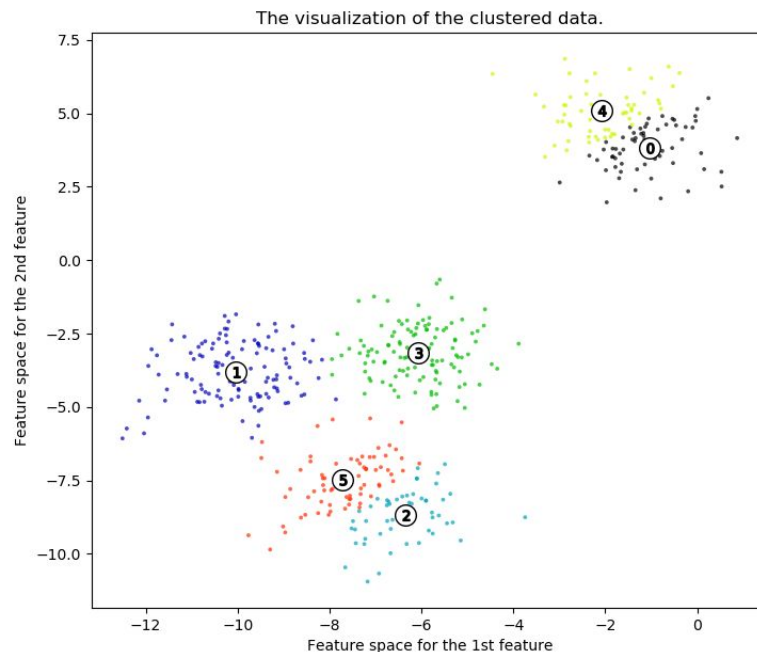
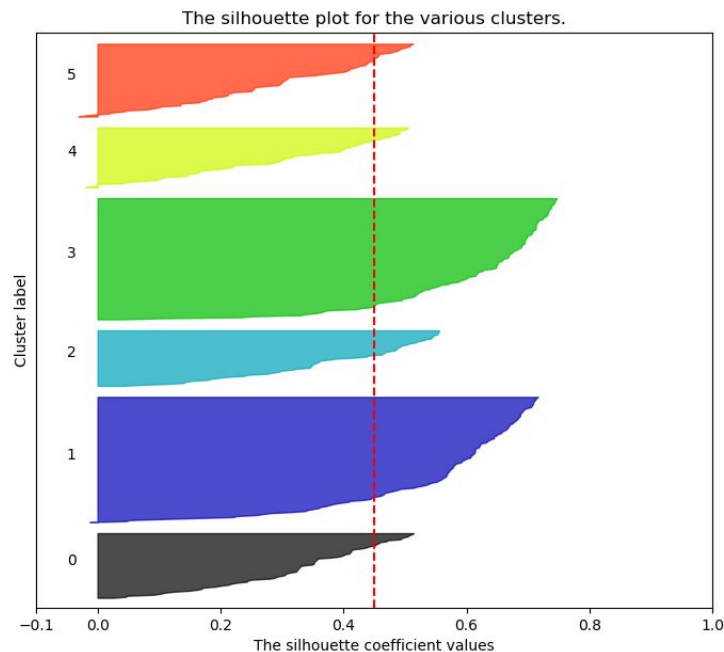
Silhouette analysis for KMeans clustering on sample data with $n_clusters = 5$



http://scikit-learn.org/stable/auto_examples/cluster/plot_kmeans_silhouette_analysis.html

Silhouette plots

Silhouette analysis for KMeans clustering on sample data with $n_clusters = 6$



http://scikit-learn.org/stable/auto_examples/cluster/plot_kmeans_silhouette_analysis.html

Silhouette plots (pick k = 2)

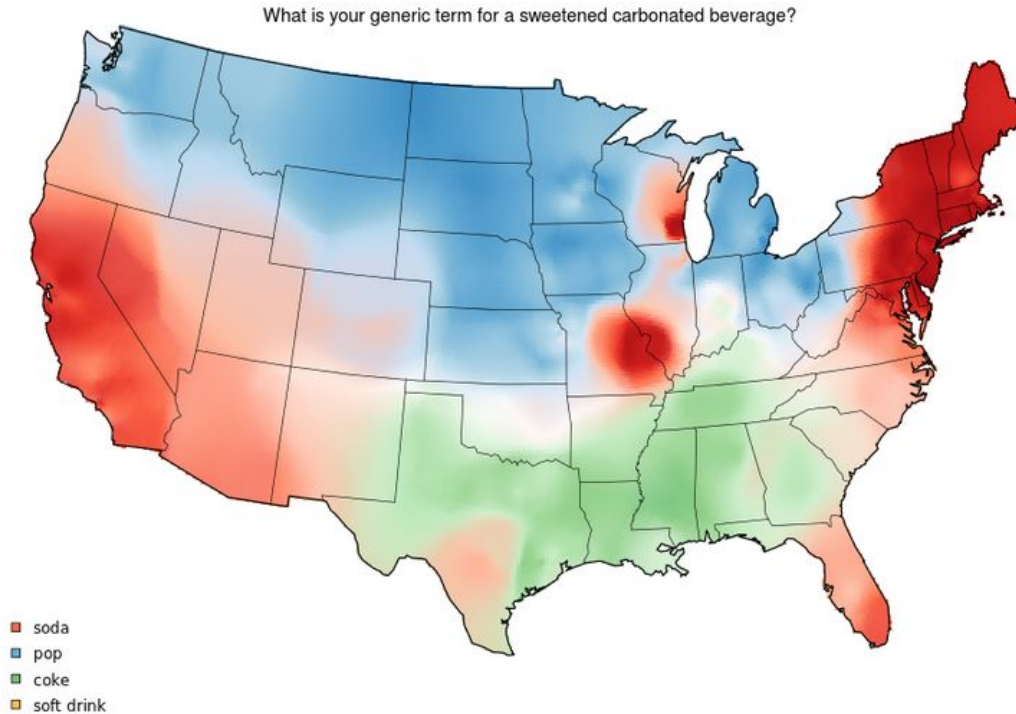
```
For n_clusters = 2 The average silhouette_score is : 0.7049787496083261
For n_clusters = 3 The average silhouette_score is : 0.5882004012129721
For n_clusters = 4 The average silhouette_score is : 0.6505186632729437
For n_clusters = 5 The average silhouette_score is : 0.56376469026194
For n_clusters = 6 The average silhouette_score is : 0.4504666294372765
```

Clustering exercises

1. Load wine.csv (14 characteristics of 178 wines from three different cultivars)
2. Plot the wines in the space defined by the first two principal components. Color each wine by cultivar
3. Run k-means with 3 cluster centers using all variables (except cultivar). Color each point in step 2 by cluster
4. Run k-means using the first two principal components only. Color each point in 2 by cluster. Compare the clustering results from (3) and (4)
5. Re-run steps 3 and 4 each four times. Do the results change?
6. Re-run steps 3-5 with 10 clusters. Compare silhouette plots.

Lab 2 Introduction

Lab 2



Joshua Katz, Department of Statistics, NC State University

<https://www.businessinsider.com/22-maps-that-show-the-deepest-linguistic-conflicts-in-america-2013-6#ok-this-one-is-crazy-everyone-pronounces-pecan-pie-differently-10>