

Data Mining - Blatt 04

Manuel, Marius

November 30, 2017

1 Abgabe Übung 4

Im Folgenden sind die Plots der einzelnen Konfigurationen der Algorithmen. Leider war es uns nicht möglich das Dataset S2 mit den Algorithmen durchlaufen zu lassen. Wir haben es für mehrere Stunden auf dem Computer laufen lassen, jedoch ohne Erfolg. Auch die Plots wurden trotz vielem Probieren und Ändern der Plot-Funktion nicht wirklich besser. `agnes diana`

Da die Plots nicht wirklich gut sind konnten wir daraus nicht sehr viel über die Datensätze lernen. Jedoch beschreiben wir im folgenden die Hierarchische Clusteranalyse.

1.1 Hierarchische Clusteranalyse

Bei der hierarchischen Clusteranalyse wird mittels distanzbasierten Verfahren eine Clusteranalyse durchgeführt. Ein Cluster besteht dabei aus Objekten, welche zueinander eine höhere Ähnlichkeit haben als zu Objekten eines anderen Clusters.

Die Verfahren der hierarchischen Clusteranalyse können durch die verwendeten Distanz- bzw. Proximitätsmaßen aufgeteilt werden. Ebenso spielt die Berechnungsvorschrift eine Rolle. Anhand derere können zwei wichtige Typen von Verfahren benannt werden.

”Top-down-Verfahren”: Hierbei werden zuerst alle Objekte einem Cluster zugeordnet und das so gebildete Cluster in kleinere Cluster aufgeteilt. Bis nur noch ein Objekt in einem Cluster ist.

”Bottom-up-Verfahren”: Hier wird zuerst jedem Objekt ein Cluster zugeordnet und dann nach und nach die bereits gebildeten Cluster solange zusammengefasst bis alle Objekte zu einem Cluster gehören.

Allgemein gilt, dass gebildete Cluster nicht mehr verändert werden können. Nur die Struktur der Cluster wird somit entweder verfeinert oder vergrößert.¹

¹https://de.wikipedia.org/wiki/Hierarchische_Clusteranalyse

Figure 1: Dataset Seeds

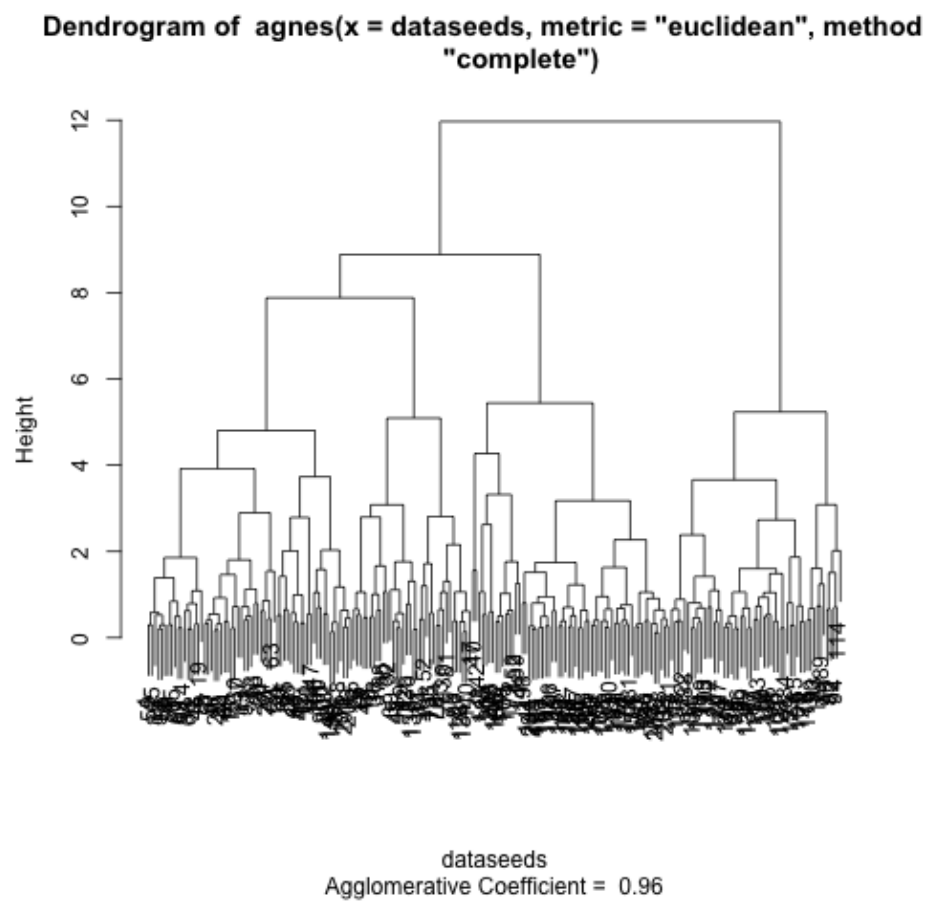
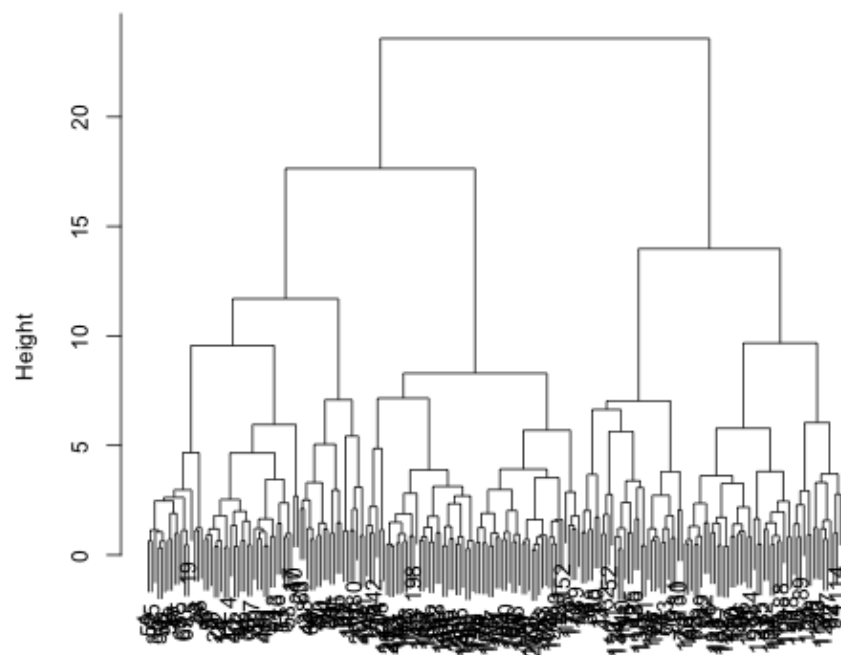


Figure 2: Dataset Seeds

Dendrogram of `agnes(x = dataseeds, metric = "manhattan", method "complete")`



dataseeds
Agglomerative Coefficient = 0.96

Figure 3: Detaset Seeds

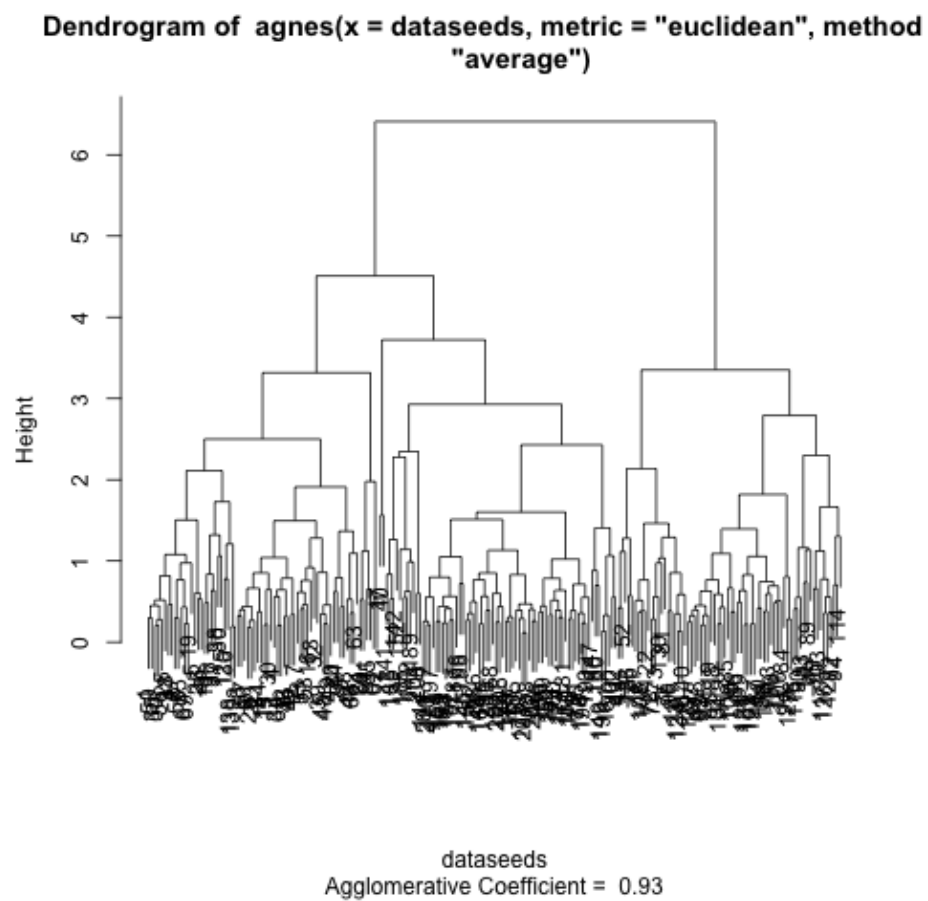


Figure 4: Detaset Seeds

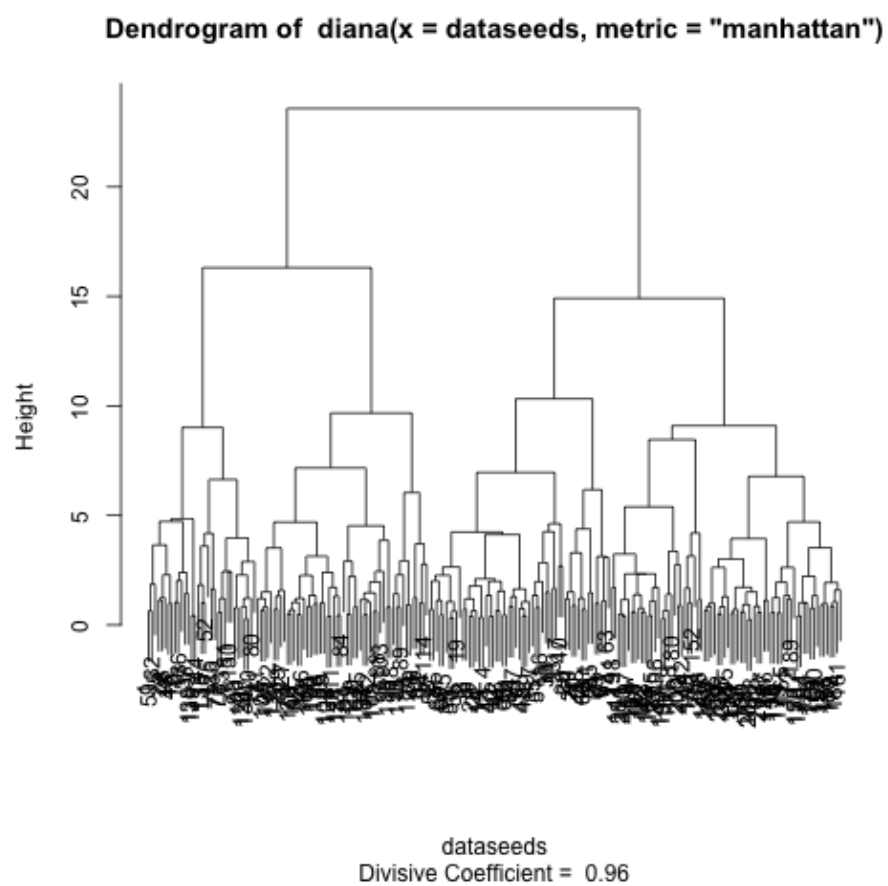


Figure 5: Detaset Seeds

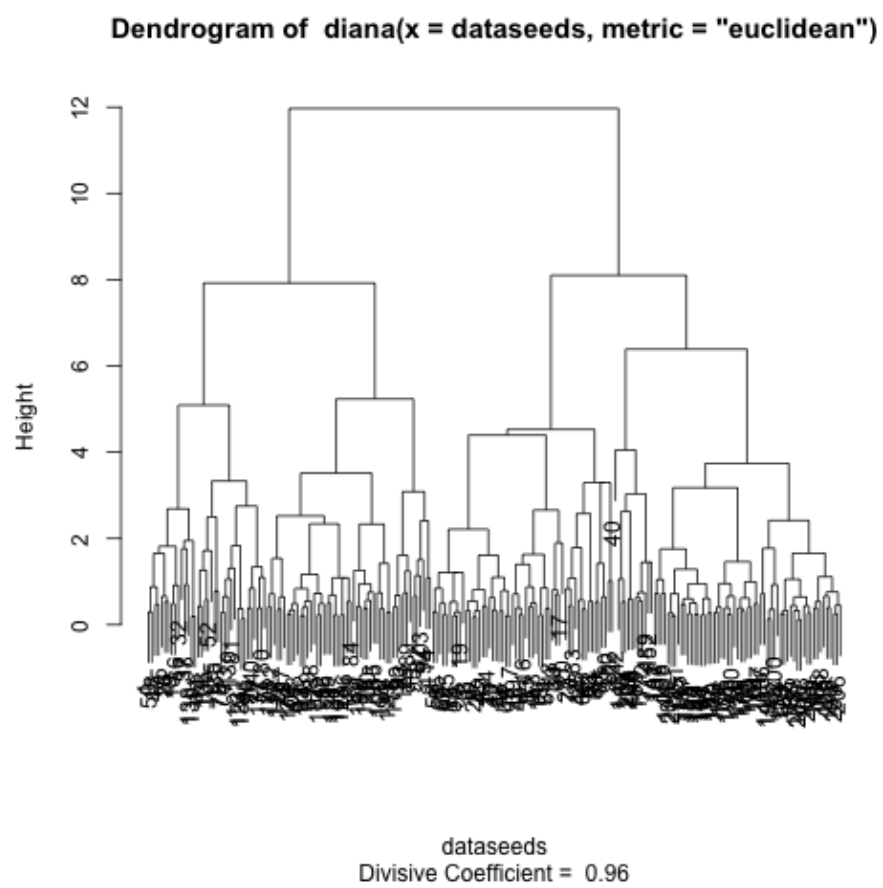
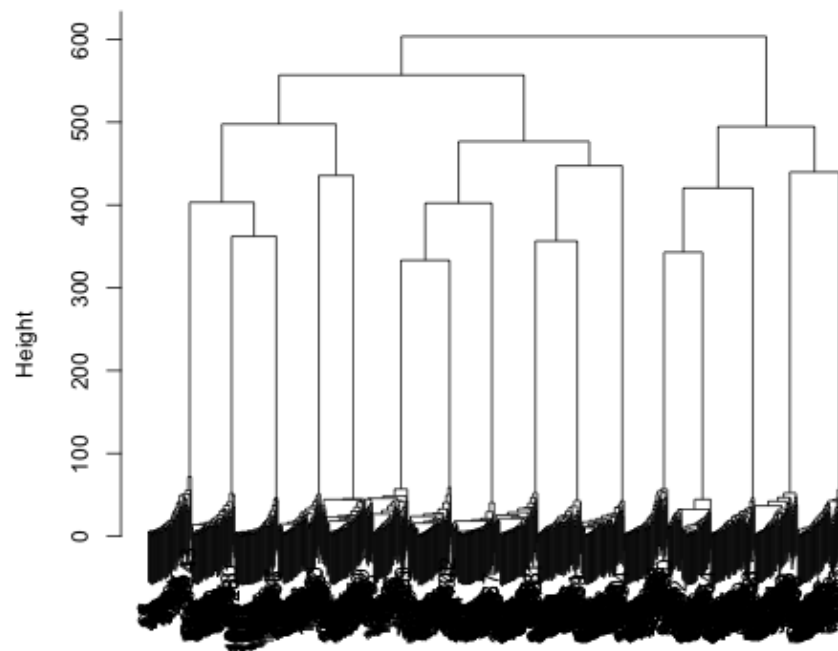


Figure 6: Dataset Dim032

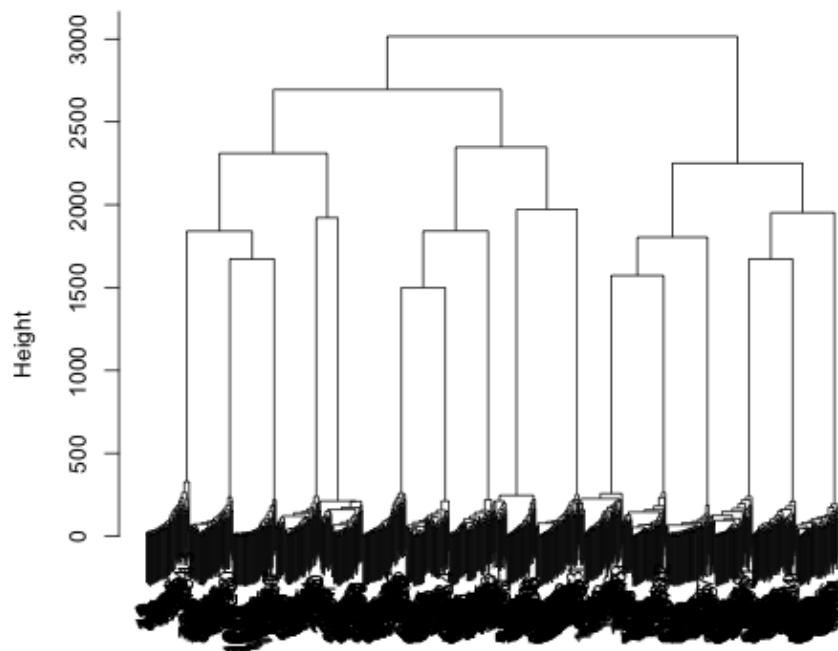
Dendrogram of `agnes(x = datadim032, metric = "euclidean", method "complete")`



datadim032
Agglomerative Coefficient = 0.98

Figure 7: Dataset Dim032

Dendrogram of `agnes(x = datadim032, metric = "manhattan", method = "complete")`



datadim032
Agglomerative Coefficient = 0.98

Figure 8: Dataset Dim032

Dendrogram of `agnes(x = datadim032, metric = "euclidean", method "average")`

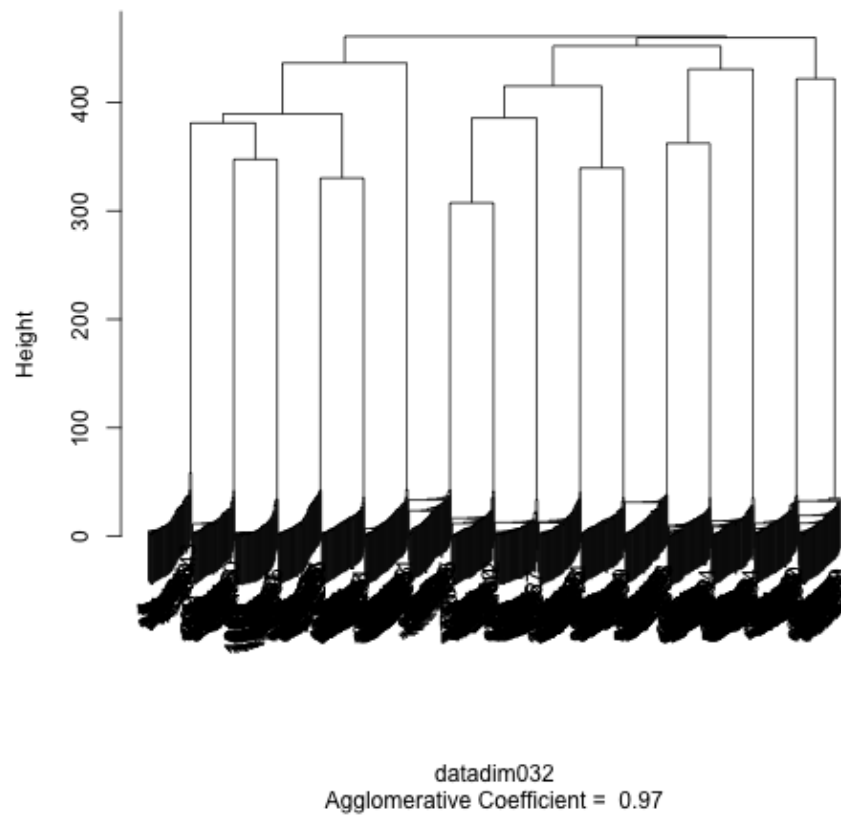


Figure 9: Detaset Dim032

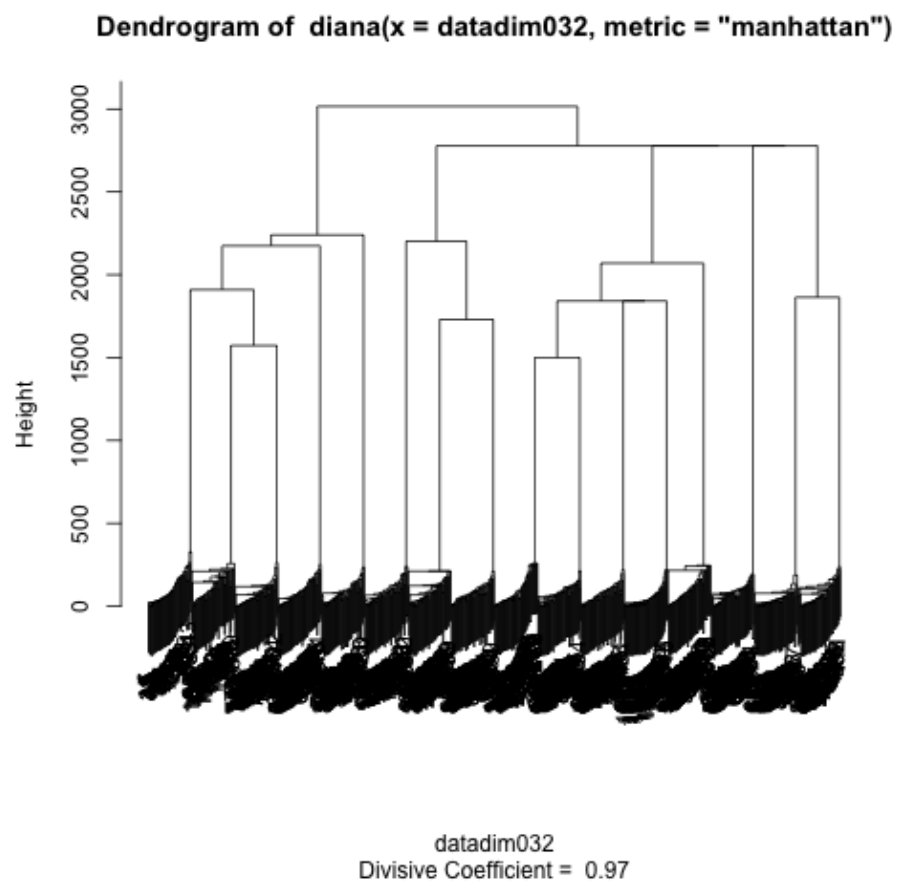


Figure 10: Detaset Dim032

