

K -means og hierarkisk clustering

Lecture 2

Torben Tvedebrink
tvede@math.aau.dk

Institut for Matematiske Fag



AALBORG UNIVERSITY
DENMARK

Self-study opgaver

Lotte, Kenneth og Nikolaj



K-means

1 Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

The floor is yours!

Supervised vs. unsupervised learning



Supervised learning er data mining problemer hvor vi kender responsen Y og en række kovariater X . Vi ønsker at være i stand til at lave *inferens* eller *prediktion*, således vores information X giver hhv. indsigt i en sammenhæng eller lav prediktions fejl.

Unsupervised learning optræder når vi kun har X og ingen klar respons Y . I visse tilfælde kan hver vektor i X betragtes som respons, men dette er lidt misvisende. Vi kan således være interesserede i at gruppere vores data i klustre, hvor medlemmerne af hvert kluster er mere ens end elementer fra øvrige klustre.

Semi-supervised learning er en kombination hvor vi kender Y for m observationer men ikke for de $n - m$ resterende observationer (ikke umiddelbart en del af dette kursus).

K -means

Self-study opgaver

2 Supervised vs. unsupervised learning

K -means

Eksempel

Hierarkiske klustre

Dendrogram

Opgaver

C4.5 vs. K -means



K -means

Self-study opgaver

3 Supervised vs. unsupervised learning

K -means

Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

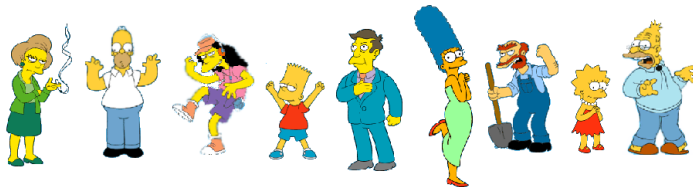
Learning tree (som C4.5) knytter sig til *supervised learning*, hvor hver observation i træningsdatasættet har en klasse (fx. syg/rask, ja/nej, ...).

I konstruktionen af træet søges attributter som opdeler klasser bedst muligt.

Cluster analysis (på dansk: klyngeanalyse) forsøger at gruppere observationer som er ens i samme cluster (så tæt på hinanden som muligt), og observationer der er forskellige i andre clustre (så langt fra hinanden som muligt).

Denne opdeling foregår oftes som *unsupervised learning*, idet klassen for observationerne er ukendt.

Hvad er en naturlig gruppering?



K-means

Self-study opgaver

4 Supervised vs.
unsupervised
learning

K-means

Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

Hvad er en naturlig gruppering?

K-means

Self-study opgaver

4 Supervised vs. unsupervised learning

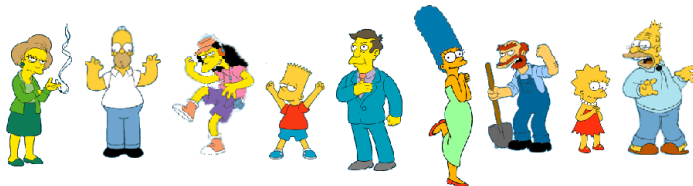
K-means

Eksempel

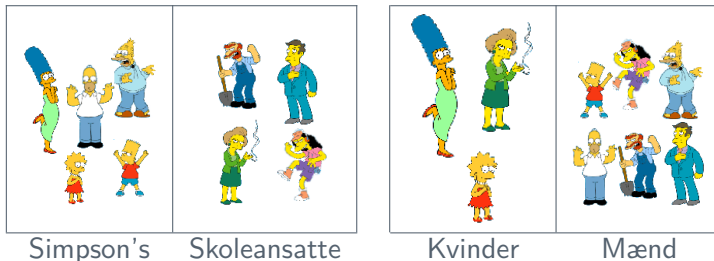
Hierarkiske clustre

Dendrogram

Opgaver



Clustering er subjektiv



Hvad er forskellighed?



K-means

Self-study opgaver

5 Supervised vs.
unsupervised
learning

K-means

Eksempel

Hierarkiske clustre

Dendrogram

Opgaver



Forskellighed (dissimilarity)



K-means

Self-study opgaver

6 Supervised vs.
unsupervised
learning

K-means

Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

Til at afgøre hvilke objekter der ens - og hvilke der er forskellige, behøves et mål for “(dis)similarity”.

Lad \mathbf{x}_i være den i 'te række/observation. En dissimilarity afstand, D , skal opfylde:

- ▶ $D(\mathbf{x}_i, \mathbf{x}_j) = D(\mathbf{x}_j, \mathbf{x}_i)$
- ▶ $D(\mathbf{x}_i, \mathbf{x}_i) = 0$
- ▶ $D(\mathbf{x}_i, \mathbf{x}_j) = 0$ hvis og kun hvis $\mathbf{x}_i = \mathbf{x}_j$.
- ▶ $D(\mathbf{x}_i, \mathbf{x}_j) \leq D(\mathbf{x}_i, \mathbf{x}_k) + D(\mathbf{x}_k, \mathbf{x}_j)$

Bemærk at trekantsuligheden gør D til en metrik - denne betingelse behøver ikke at være opfyldt for alle clustering algoritmer.

Beregne parwise “dissimilarities” i R



K-means

Self-study opgaver

7

Supervised vs. unsupervised learning

K-means

Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

I R kan vi beregne “dissimilarities” mellem objekter/observationer i en data.frame meget enkelt vha. **dist**-funktionen:

```
> data <- data.frame(x1 = c(0,1,3), x2=c(0,1,4))  
> (res <- dist(data, method="euclidian", diag=TRUE, upper=TRUE))
```

```
      1      2      3  
1 0.000000 1.414214 5.000000  
2 1.414214 0.000000 3.605551  
3 5.000000 3.605551 0.000000
```

```
> class(res)  
[1] "dist"
```

R-pakken `proxy` udvider `dist`-funktionen med en masse forskellige afstandsmål.

Fx. behøver vi ikke blot at være interesserede i en Euklidiske afstand, men også *Manhattan* afstanden som er

$$D(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|_1 = \sum_{k=1}^p |x_{ik} - x_{jk}|.$$

*[Prøv det for **data** – forrige slide]*

K-means

Self-study opgaver

8

Supervised vs.
unsupervised
learning

K-means

Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

R-pakken proxy udvider dist-funktionen med en masse forskellige afstandsmål.

```
library(proxy)
> summary(pr_DB)
* Similarity measures:
Braun-Blanquet, Chi-squared, correlation, cosine, Cramer,
Dice, eJaccard, Fager, Faith, fJaccard, Gower, Hamman,
Jaccard, Kulczynski1, Kulczynski2, Michael, Mountford, Mozley,
Ochiai, Pearson, Phi, Phi-squared, Russel, simple matching,
Simpson, Stiles, Tanimoto, Tschuprow, Yule, Yule2
```

```
* Distance measures:
Bhattacharyya, Bray, Canberra, Chord, divergence, Euclidean,
Geodesic, Hellinger, Kullback, Levenshtein, Mahalanobis,
Manhattan, Minkowski, Podani, Soergel, supremum, Wave,
Whittaker
```

K-means

Self-study opgaver

8

Supervised vs.
unsupervised
learning

K-means

Eksempel

Hierarkiske clustre

Dendrogram

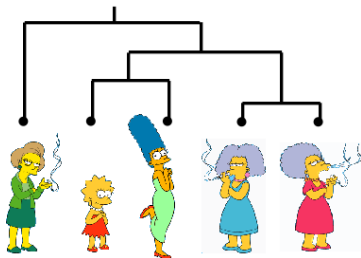
Opgaver

To typer af clustering

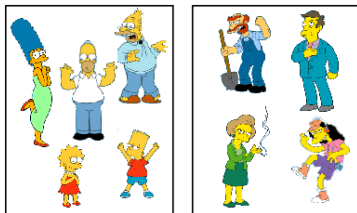
Der findes overordnet to typer af clustering metoder:
hierarkisk og **partitionering**.

Forskellen er at hierarkisk clustering opdeler clustre i underclustre, mens partitionering opdeler data i disjunkte clustre.

Hierarkisk



Partitionering



K-means

Self-study opgaver

9 Supervised vs. unsupervised learning

K-means

Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

To typer af clustering



Der findes overordnet to typer af clustering metoder:
hierarkisk og **partitionering**.

► Hierarkisk

- *Agglomerative* Bottom-up metode som starter med hver observation i eget cluster, og sammensmelter mest *ens* clustere (forskellige mål for *ens*)
- *Divisive* Starter med alle observationer i et cluster som derefter rekursivt opdeles i de to mest *forskellige* enheder (forskellige mål for *forskel*)

► Partitionering

- Laver disjunkte clustre hvor hvert cluster intern er så ens som muligt, og så forskellig fra de øvrige som muligt (fx. *K*-means).

K-means

Self-study opgaver

9 Supervised vs. unsupervised learning

K-means

Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

K-means



K -means algoritmen er meget simpel og beror på en opdeling af data i K clustre hvor *within* kvadratsummen ønskes så lille som mulig:

$$\sum_{i=1}^n \min_{j=1, \dots, k} \|\mathbf{x}_i - \mathbf{c}_j\|_2^2$$

K -means

Self-study opgaver

Supervised vs.
unsupervised
learning

10

K -means

Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

33

K-means algoritmen er meget simpel og beror på en opdeling af data i K clustre hvor *within* kvadratsummen ønskes så lille som mulig:

$$\sum_{i=1}^n \min_{j=1, \dots, k} \|\mathbf{x}_i - \mathbf{c}_j\|_2^2$$

Algoritme:

- 0.1 Vælg antal clustre K
- 0.2 Vælg tilfældigt K cluster repræsentanter tilfældigt (fx. K data punkter)
 1. Hver observation, \mathbf{x}_i , allokeres til nærmeste center, \mathbf{c}_j .
 2. Opdater $\mathbf{c}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} \mathbf{x}_{i_j}$, hvor $\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_{n_j}}$ tilhører kluster j .
 3. Gentag 1-2 indtil konvergens (ingen ændring i allokering).

Self-study opgaver

Supervised vs. unsupervised learning

10 K-means

Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

Antal clusterer?

K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

11

K-means

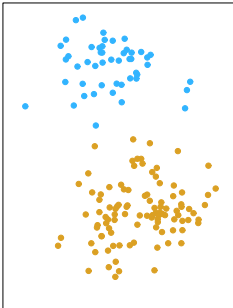
Eksempel

Hierarkiske clusterer

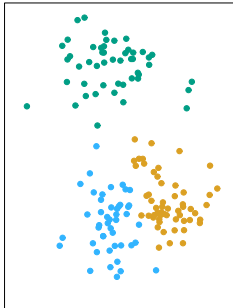
Dendrogram

Opgaver

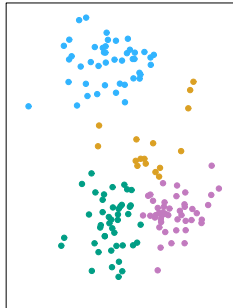
K=2



K=3



K=4



Problemer med K -means



K -means

Self-study opgaver

Supervised vs.
unsupervised
learning

12 K -means

Eksempel

Hierarkiske clustre

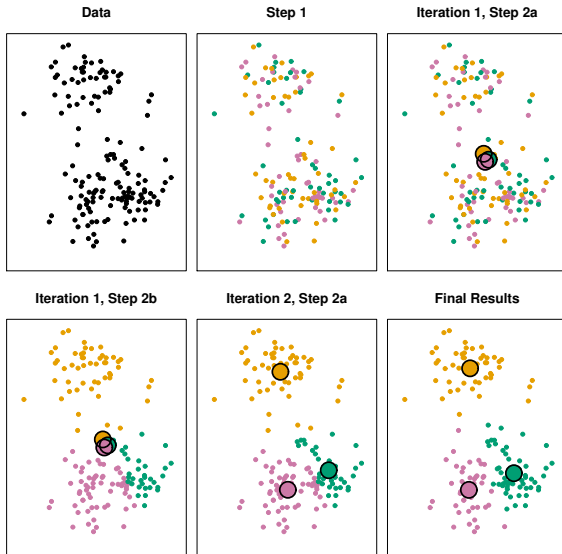
Dendrogram

Opgaver

- ▶ Skal beslutte antal clustre, K , før analyse
- ▶ Følsom overfor start værdier (se simulations eksempel)
- ▶ Følsom over for outliers
- ▶ ...

Algoritmens basale skridt

ISLR, Figure 10.6, p. 389



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

13

Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

33

Torben Tvedebrink
tvede@math.aau.dk

Simuleret data



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

14 Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

Simulér data fra fire bivariat normal fordelinger med hvert sit centrum, μ_i , og varians, $\sigma_i^2 I_2$.

Som i bogen (pp. 27-29) vælges

$$\mu_1 = \begin{pmatrix} -3 \\ -3 \end{pmatrix} \quad \mu_2 = \begin{pmatrix} 3 \\ -3 \end{pmatrix} \quad \mu_3 = \begin{pmatrix} -1 \\ 2 \end{pmatrix} \quad \mu_4 = \begin{pmatrix} 1 \\ 2 \end{pmatrix},$$

$$\sigma_1^2 = 0.0625 \quad \sigma_2^2 = 1 \quad \sigma_3^2 = 1 \quad \sigma_4^2 = 1$$

$$n_1 = 200 \quad n_2 = 200 \quad n_3 = 150 \quad n_4 = 150.$$

Simuleret data



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

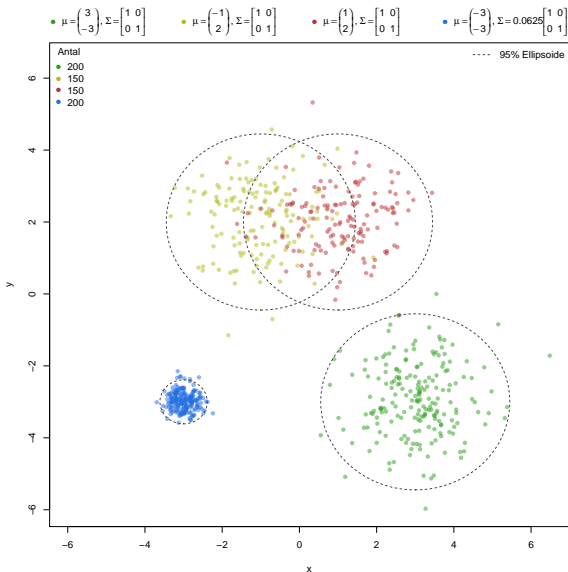
K-means

14 Eksempel

Hierarkiske clustre

Dendrogram

Opgaver



K-means på simulationerne

Begyndelsesværdier 1



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

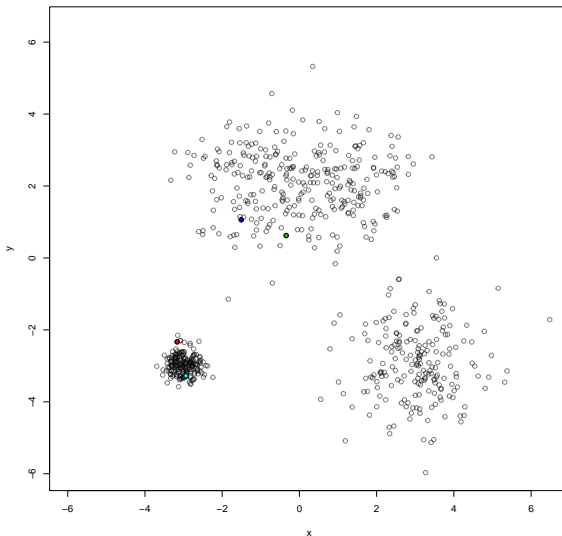
K-means

15 Eksempel

Hierarkiske clustre

Dendrogram

Opgaver



K-means på simulationerne

Begyndelsesværdier 1



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

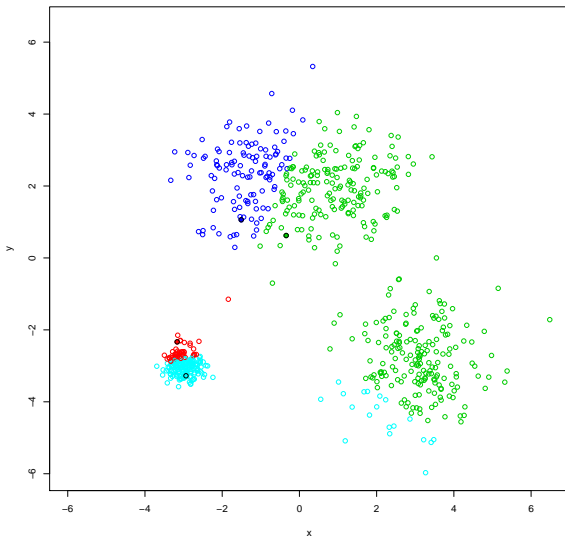
K-means

15 Eksempel

Hierarkiske clustre

Dendrogram

Opgaver



K-means på simulationerne

Begyndelsesværdier 1



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

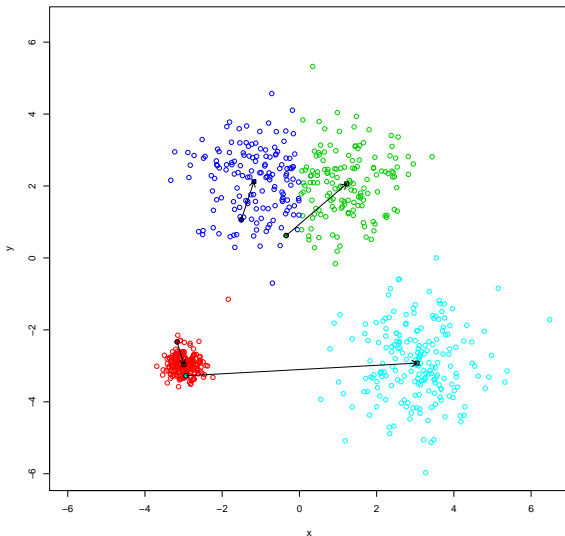
K-means

15 Eksempel

Hierarkiske clustre

Dendrogram

Opgaver



K-means på simulationerne

Begyndelsesværdier 2



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

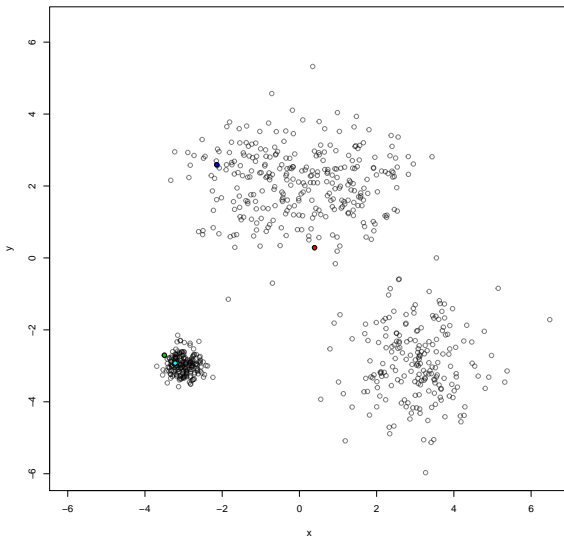
K-means

15 Eksempel

Hierarkiske clustre

Dendrogram

Opgaver



K-means på simulationerne

Begyndelsesværdier 2



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

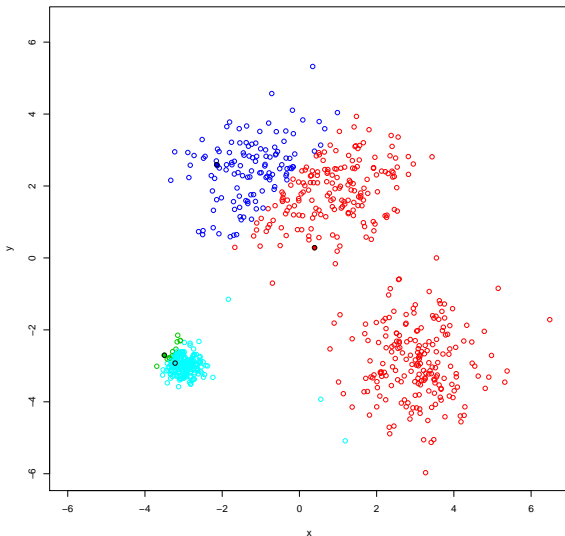
K-means

15 Eksempel

Hierarkiske clustre

Dendrogram

Opgaver



K-means på simulationerne

Begyndelsesværdier 2



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

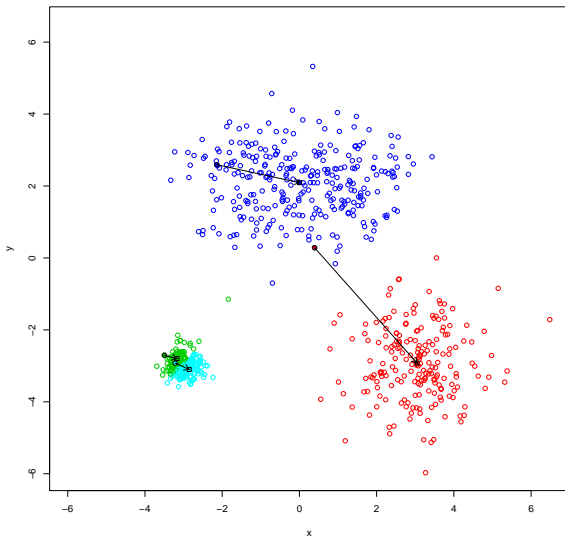
K-means

15 Eksempel

Hierarkiske clustre

Dendrogram

Opgaver



K-means på simulationerne

Begyndelsesværdier 3



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

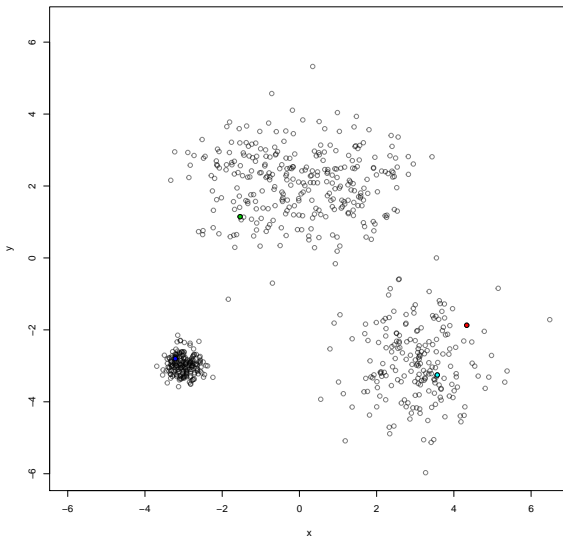
K-means

15 Eksempel

Hierarkiske clustre

Dendrogram

Opgaver



K-means på simulationerne

Begyndelsesværdier 3



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

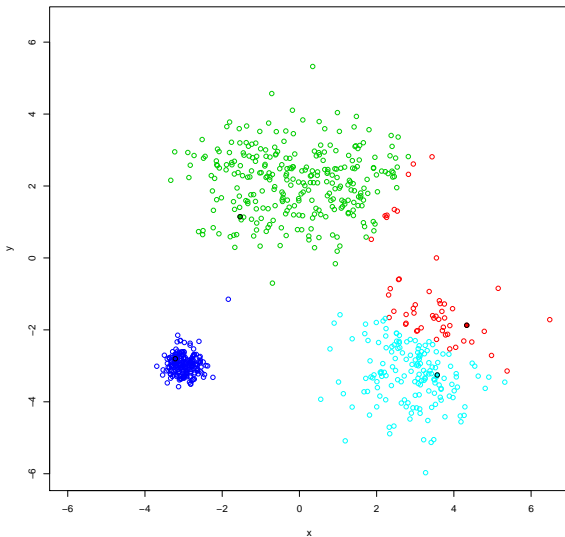
K-means

15 Eksempel

Hierarkiske clustre

Dendrogram

Opgaver



K-means på simulationerne

Begyndelsesværdier 3



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

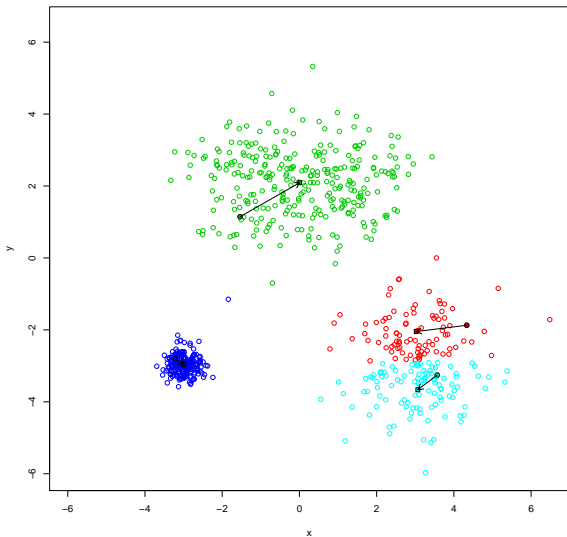
K-means

15 Eksempel

Hierarkiske clustre

Dendrogram

Opgaver



K-means er tilgængelig i R uden yderligere pakker ved at bruge `kmeans`. Se på dokumentationen for `kmeans` for at gøre jer bekendte med argumenterne.

For eksempel, man kan køre K-means med ti tilfældige sæt af start værdier ved at angive `nstart = 10`.

Vi kan tilgå en totale variabilitet og variabiliteten forklaret ved "modellen" fra outputet fra `kmeans`:

`"totss" "withinss" "tot.withinss" "betweenss"`

Sums of squares

Dekomponere variationen i data



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

17 Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

I statistik ønsker vi at forklare så meget af variabiliteten ved systematiske komponenter. I clustering er de systematiske komponenter cluster tilhørsforholdet.

Sums of squares

Dekomponere variationen i data: SS_{TOT} , SS_W og SS_B



Lad x_i være en skalar (fx. reelle tal – argumentet holder også for vektorer) så har vi:

$$\bar{x} = n^{-1} \sum_{i=1}^n x_i \quad \text{og} \quad \bar{x}_k = |C_k|^{-1} \sum_{i \in C_k} x_i$$

Den *totale sum of squares* og *within sums of squares* er da givet ved

$$SS_{TOT} = \sum_{i=1}^n (x_i - \bar{x})^2 \quad \text{og} \quad SS_{W_k} = \sum_{i \in C_k} (x_i - \bar{x}_k)^2$$

Ydermere, *between sums of squares* er defineret som

$$SS_W = \sum_{k=1}^K SS_{W_k} \quad \text{and} \quad SS_B = SS - SS_W$$

K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

17

Eksempel

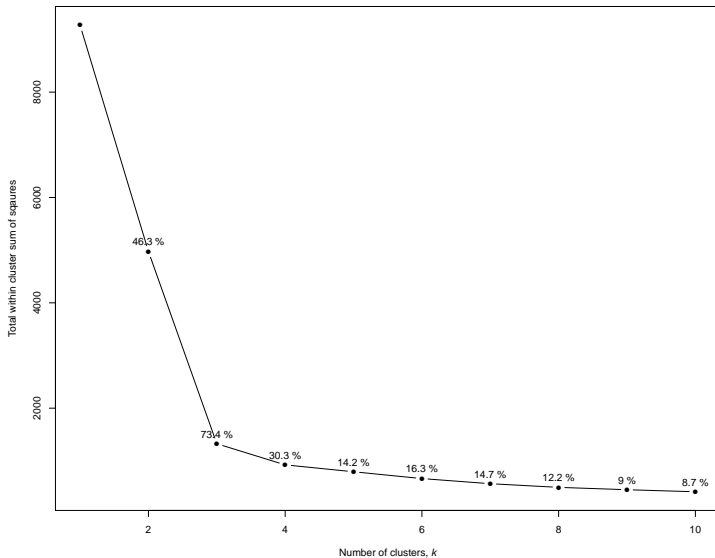
Hierarkiske clustre

Dendrogram

Opgaver

Sums of squares

Simulations eksempel fortsat



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

17

Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

33

Torben Tvedebrink
tvede@math.aau.dk

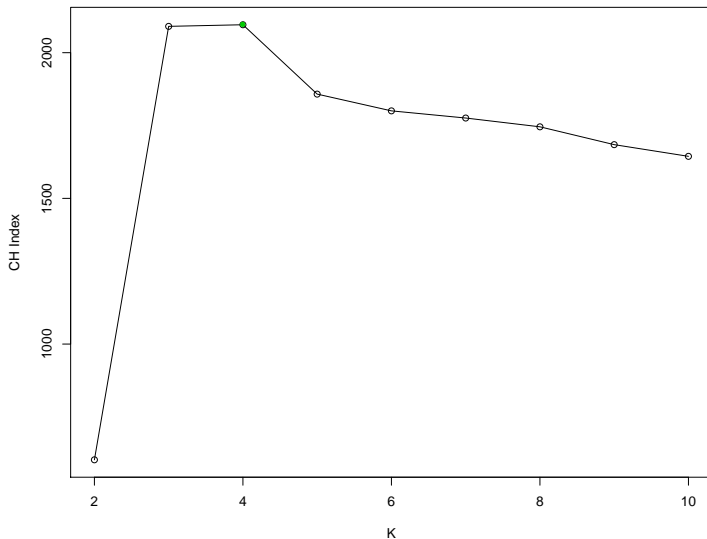
Til at vælge antal clustre foreslog Calinski and Harabasz (1974), “A dendrite method for cluster analysis” dette index, som ønskes størst muligt:

$$CH(K) = \frac{B(K)/(K-1)}{W(K)/(n-K)} = \frac{SS_B/(K-1)}{SS_W/(n-K)}$$

hvor $B(K)$ og $W(K)$ er hhv. between-sum-of-squares og within-sum-of-squares for K clustre.

CH Index

Simulations eksempel fortsat



K -means

Self-study opgaver

Supervised vs.
unsupervised
learning

K -means

18 Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

Kør modellen flere gange

SS_W , ISLR, Figure 10.7, p. 390



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

19 Eksempel

Hierarkiske clustre

Dendrogram

Opgaver



K-means alternativer

K-mediods



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

20 Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

I stedet for at bruge gennemsnittet som centre, c_j , kan medianen bruges (K -mediods, tilgængeligt via [pam](#) i [cluster](#)-pakken).

K -means alternativer

X -means

I stedet for at bruge gennemsnittet som centre, c_j , kan medianen bruges (K -mediods, tilgængeligt via **pam** i **cluster**-pakken).

En anden metode, X -means, består i at bruge et Informations Kriteire til at bestemme hvor mange klustre data understøtter, dvs. estimere K på baggrund af hvor godt data fitter forskellige modeller. Scoren/kvadratsummen straffes for flere klustre/parametre vha. BIC ved $\frac{q}{2} \log n$, hvor q er antal parametre i modellen og n er antal observationer.



K -means

Self-study opgaver

Supervised vs.
unsupervised
learning

K -means

20 Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

K-means alternativer

Soft clustering



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

20

Eksempel

Hierarkiske clustre

Dendrogram

Opgaver

I stedet for at bruge gennemsnittet som centre, c_j , kan medianen bruges (K -mediods, tilgængeligt via `pam` i `cluster`-pakken).

I modsætning til K -means som udfører *hard clustering* kan man benytte *soft clustering*. Her bliver hvert datapunkt tilskrevet et cluster med en given sandsynlighed – hvor *hard clustering* angiver “sandsynligheder” svarende til 0 og 1.

Vi vender tilbage til dette i lektion 5. En god pakke til dette er `mclust` – med en udførlig vignette:

```
library(mclust)
vignette("mclust")
```

- *Agglomerative*

Bottom-up metode: Starter med n “clustre” og sammensmelter *nærmeste* clustre i hvert trin

- *Divisive*

Stop-down metode: Starter med ét “cluster” og adskiller hvert cluster indtil n ‘blade’ er nået.

K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

Eksempel

21 Hierarkiske clustre

Dendrogram

Opgaver

Hierarkiske clustre



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

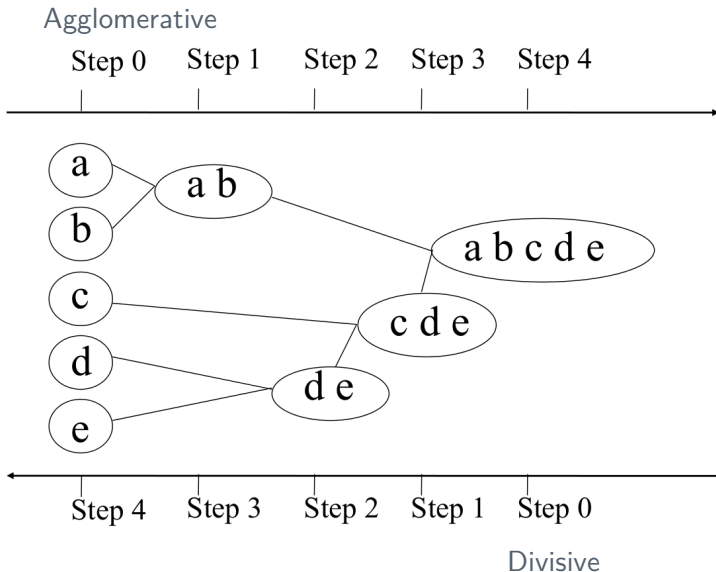
Eksempel

21

Hierarkiske clustre

Dendrogram

Opgaver



33

Eksempel

Complete link



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

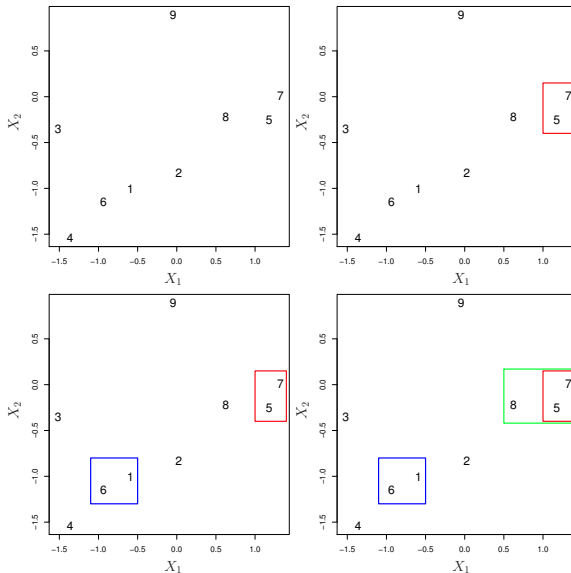
K-means

Eksempel

22 Hierarkiske clustre

Dendrogram

Opgaver



Dissimilarity mellem grupper

Linkages



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

Eksempel

23 Hierarkiske clustre

Dendrogram

Opgaver

- ▶ **Single linkage** (nærmeste nabo) - afstanden mellem de to nærmeste elementer fra forskellige clustre

$$d_{\text{Single linkage}}(C_k, C_{k'}) = \min_{i \in C_k, i' \in C_{k'}} d_{ii'}.$$

- ▶ **Complete linkage** (fjerneste nabo) - afstanden mellem de to fjerneste elementer fra forskellige clustre

$$d_{\text{Complete linkage}}(C_k, C_{k'}) = \max_{i \in C_k, i' \in C_{k'}} d_{ii'}.$$

- ▶ **Average linkage** - den gennemsnitlige afstand i mellem alle par af punkter i hvert cluster

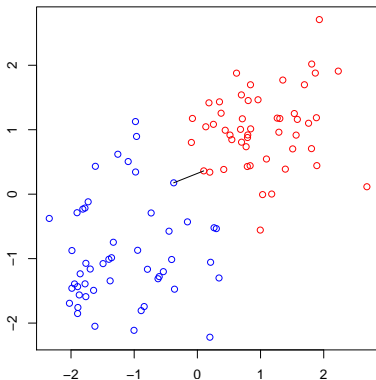
$$d_{\text{Average linkage}}(C_k, C_{k'}) = \frac{1}{|C_k| \cdot |C_{k'}|} \sum_{i \in C_k, i' \in C_{k'}} d_{ii'}.$$

- ▶ **Ward's method** - minimerer SS_W når to clustre sammensmelter.
- ▶ ...

Single link

Single linkage (nærmeste nabo) - afstanden mellem de to nærmeste elementer fra forskellige clustre

$$d_{\text{Single linkage}}(C_k, C_{k'}) = \min_{i \in C_k, i' \in C_{k'}} d_{ij'}.$$



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

Eksempel

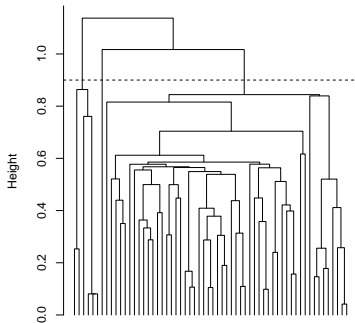
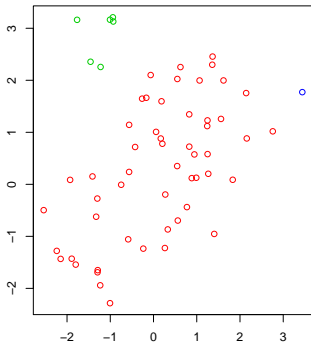
24 Hierarkiske clustre

Dendrogram

Opgaver

Single link

Eksempel



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

Eksempel

24

Hierarkiske clustre

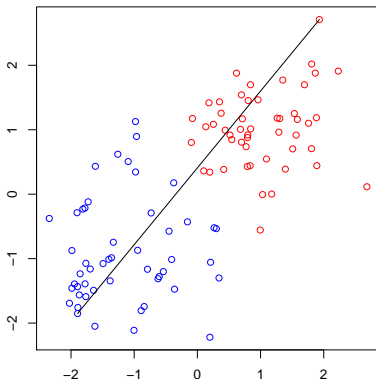
Dendrogram

Opgaver

33

Complete linkage (fjerneste nabo) - afstanden mellem de to fjerneste elementer fra forskellige clustre

$$d_{\text{Complete linkage}}(C_k, C_{k'}) = \max_{i \in C_k, i' \in C_{k'}} d_{ii'}.$$



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

Eksempel

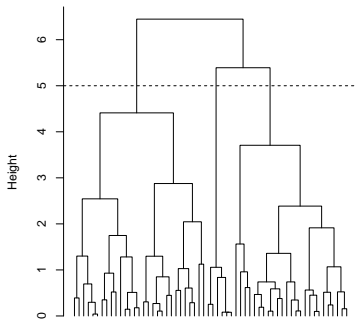
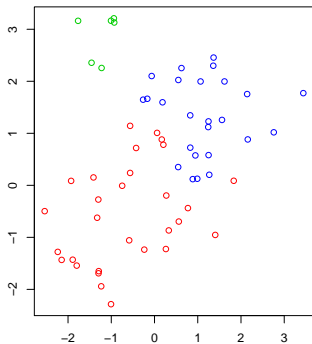
25 Hierarkiske clustre

Dendrogram

Opgaver

Complete link

Eksempel



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

Eksempel

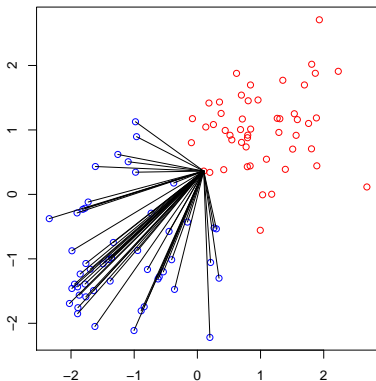
25 Hierarkiske clustre

Dendrogram

Opgaver

Average linkage - den gennemsnitlige afstand i mellem alle par af punkter i hvert cluster

$$d_{\text{Average linkage}}(C_k, C_{k'}) = \frac{1}{|C_k| \cdot |C_{k'}|} \sum_{i \in C_k, i' \in C_{k'}} d_{ii'}.$$



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

Eksempel

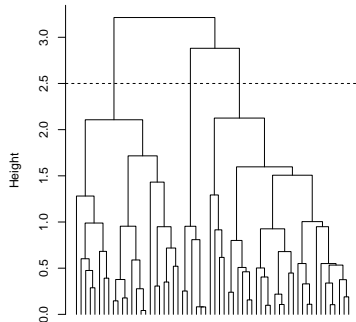
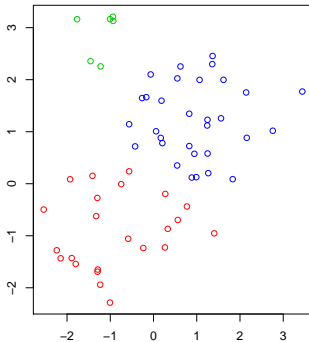
26 Hierarkiske clustre

Dendrogram

Opgaver

Average link

Eksempel



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

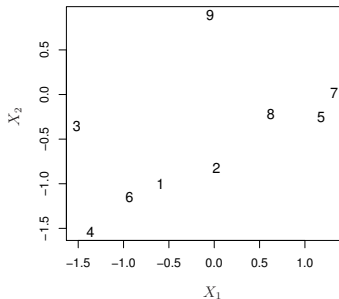
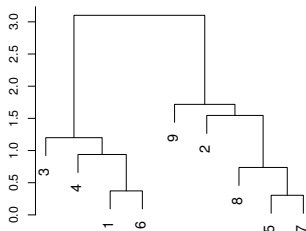
Eksempel

26 Hierarkiske clustre

Dendrogram

Opgaver

Dendrogrammet



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

Eksempel

Hierarkiske clustre

27 Dendrogram

Opgaver

Hvert af de foregående linkages giver anledning til forskellige clustre (nogen minder naturligvis om hinanden for simple eksempler).

Alt afhængigt af formålet med cluster analysen vil forskellige mål være at foretrække.

K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

Eksempel

Hierarkiske clustre

28

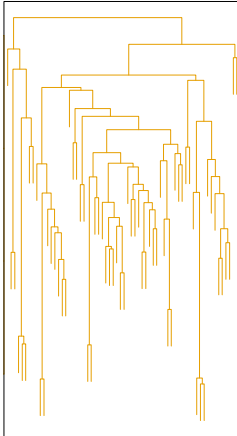
Dendrogram

Opgaver

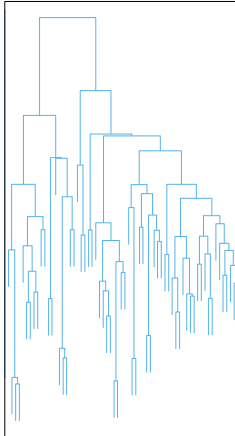
33

Linkages

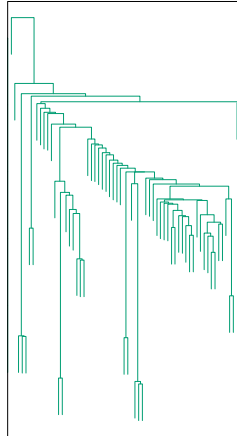
Average Linkage



Complete Linkage



Single Linkage



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

Eksempel

Hierarkiske clustre

28

Dendrogram

Opgaver

Partitionering af hierarkiske clusterer

Forskellige cut højder



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

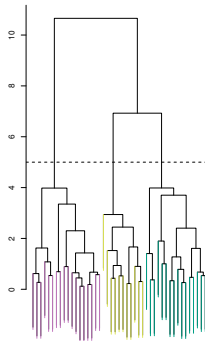
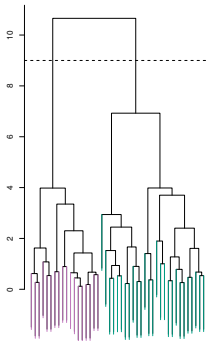
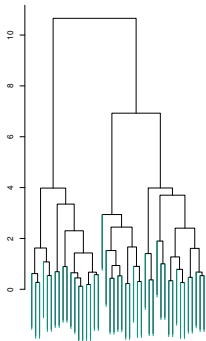
Eksempel

Hierarkiske clusterer

29

Dendrogram

Opgaver



Partitionering af hierarkiske clustere

Ved tre clustre



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

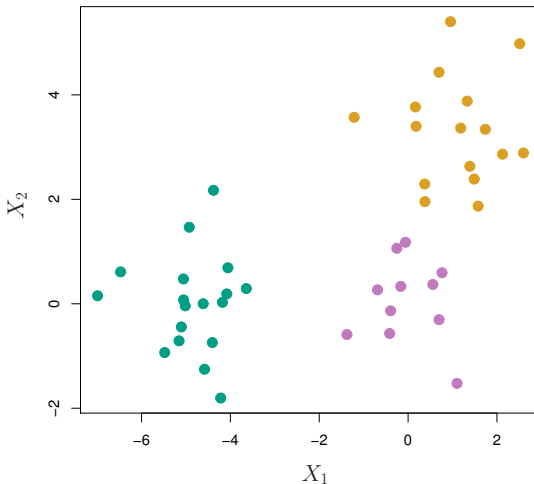
Eksempel

Hierarkiske clustre

29

Dendrogram

Opgaver



I R kan man lave hierarkisk clustering vha. `hclust` funktionen i R.

`hclust` tager følgende vigtige argumenter: `d` og `method`, hvor `d` er et `dist`-objekt fx. fra `dist`-funktionen. Det andet argument, `method`, refererer til linkage metoden, fx. "complete", "single", "average", "ward.D2"

Mange andre hierarkiske clustering metoder er tilgængelige via `cluster` pakken i R. Ved at bruge funktionerne fra denne pakke, se listen af funktioner og datasæts ved `help(package="cluster")`.

Eksempel



K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

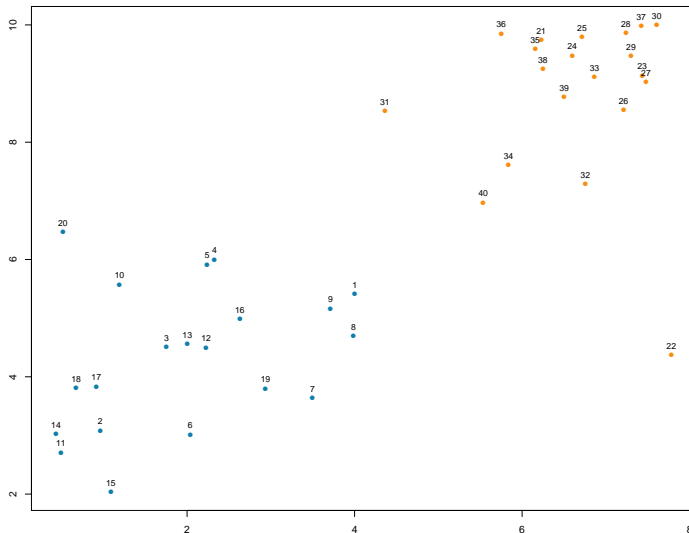
K-means

Eksempel

Hierarkiske clustre

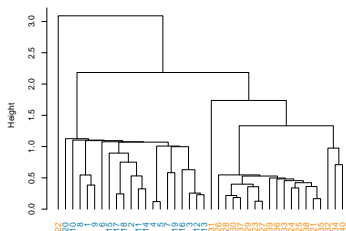
31 Dendrogram

Opgaver

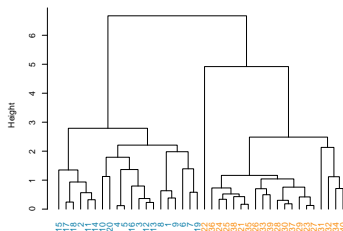


33

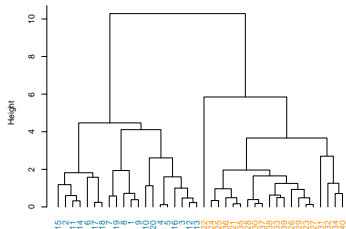
Eksempel



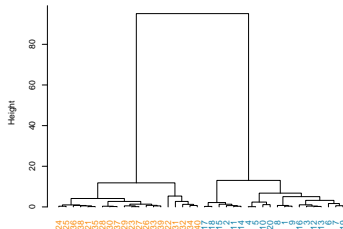
Single link



Average link



Complete link



Ward's link

K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

Eksempel

Hierarkiske clustre

31

Dendrogram

Opgaver

33

Plot dendogrammet



Efter vi har udført hierarkisk clustering og gemt objektet som et R objekt, fx.:

```
hc <- hclust(dist(USArrests), "ave")  
plot(hc)  
plot(hc, hang = -1)
```

vi plotter altså ved at bruge **plot**-kommandoen. I R virker mange funktioner på forskellig vis afhængigt af typen af objektet.

Man kan “overskære” træet i en bestemt højde eller antal undertræer vha. **cutree**-kommandoen, ved hhv. **h** eller **k** argumenterne.

K-means

Self-study opgaver

Supervised vs.
unsupervised
learning

K-means

Eksempel

Hierarkiske clustre

32

Dendrogram

Opgaver

33

- ▶ Løs opgaver 1–9 fra afsnit 2.7 i Wu og Kumar.
- ▶ Antag at data er *række*-standardiseret. Vis at $\|\mathbf{x}_i - \mathbf{x}_j\|_2^2 = 2(1 - r_{ij})$, hvor r_{ij} er korrelationen mellem række i og j .
Vis ligeledes resultatet numerisk for datasættet USArrests (fx. ved et plot).
- ▶ Gennemgå analyserne i MASS4-cluster.pdf (moodle) – og suppler hvor nødvendigt for at forstå de anvendte metoder.