

## Praktikum 8

### Klaster- ja peakomponentanalüüs

- ☞ Salvestage arvutisse Eesti puude andmestik *R*-i failina:  
[http://www.eau.ee/~ktanel/DK\\_0007/puud.rda](http://www.eau.ee/~ktanel/DK_0007/puud.rda)
- ☞ Avage *R* ja käivitage lisamoodul *Rcmdr*.
- ☞ Võtke *R*-is kasutusele salvestatud puude andmestik
  - *R Commander*is käsud *Data -> Load data set ...*
  - või skripti aknas käsk *load*, näiteks  
`load("C:/Documents and Settings/Tanel/DK_0007/puud.rda")`
- ☞ **puud1 = na.omit(puud)** # tekitamiseks uut puuduvate väärtusteta andmestikku  
 (*R Commander* oskab oma menüükäskudes küll puuduvad väärtused välja jätta, aga nende tulemuste edasine skriptiaknas töötlemine võib viia veateadeten; klasteranalüüsi üksnes skriptiaknasse sisestatud baas-*R*-i käskude abil teostamine eeldabki puuduvate väärtusteta andmestikku)
- ☞ **attach(puud1)** # edasiste käskude lihtsama esituse huvides
- ☞ Puude andmestikus vastab üks rida ühele puule, veergudes on vastavalt:
  - veerus nimega 'A' puu vanused aastates;
  - veerus 'D' puu diameeter sentimeetrites;
  - veerus 'H' puu kõrgus meetrites;
  - veerus 'ARENGUKL' puu arenguklass (A – lage, N – noorendikud, L – latimets (noorendikust järgmine), K – keskealised, V – valmiv, Y – küps, S – selgusetu, – puuduv väärtus);
  - veerus 'PE' puu liik (HB – haab, KS – kask, KU – kuusk, LH - lehis, LM – sanglepp, LV – hall lepp, MA – mänd, RE – remmelgas, SA – saar, TA – tamm);
  - veerus 'KKT' kasvukohatüüp (AN – angervaksa, JM – mustika, ...)<sup>1</sup>;
  - veerus 'H100' baaskõrgus (prognoositav kõrgus saja aasta vanuses).

#### 1. Klasteranalüüs

Kaks peamist klasterdamisega seotud funktsiooni *R*-is on

`hclust` (hierarhiline klasterdamine) ja  
`kmeans` (k-keskmiste meetod).

*R Commander*is on mõlemad meetodid leitavad menüüdest

*Statistics -> Dimensional analysis -> Cluster analysis ->*

Nagu *R Commander* puhul tavaline, on menüüdest valitud analüüside tulemusel kirjutatav programm pisut teistsugune – *R Commander* kasutab vähe erinevat süntaksit, samas töötavad *R Commander* skripti aknas ka tavalisele *R*-ile omased käsud `hclust` ja `kmeans`.

Kuna hierarhiline klasterdamine on märksa töömahukam ja suurte andmestike korral ebaülevaatlikum võrreldes k-keskmiste meetodiga, kasutatakse esimest enamasti väikeste ja teist suurte andmestike korral.

<sup>1</sup> Kasvukohatüüp (ka metsakasvukohatüüp) on mullastikult ja taimestikult ühtlane metsaala. Nimetus tuleb enamesineva taime järgi (näiteks kasvab naadi kasvukohas väga palju naate, jänesekapsa kasvukohatüübis jänesekapsaid jne) – <http://et.wikipedia.org/wiki/Kasvukohatüüp>. Täpsemalt vt näiteks <http://www.hot.ee/sinumets3/sinumets03-08.pdf>.

## 2. Püüdke jagada üksikud puud klastritesse kasutades andmeid puude vanuse ja diameetri kohta.

Kuna puid on palju (puuduvate väärtusteta andmestikus „puud1“  $n=25436$ ), siis on mõistlik kasutada k-keskmiste meetodit ja jagada võiks puud näiteks nelja klastrisse.

# 1. võimalus – käsuna skriptiaknas:

```
tree_kmean1 = kmeans(cbind(A, D), 4)
```

#2. võimalus – R Commander'i menüüdest:

*Statistics -> Dimensional analysis -> Cluster analysis -> kmeans cluster analysis ...*

Eelkõige on oluline määrata soovitud klastrite arv, ülejäänud parameetrid (nn algseemnete arv ja maksimaalne iteratsioonide arv) on seotud iteratiivse klasterdamisalgoritmi tööga ja võivad enamusel juhtudel jääda muutmata (sama käsku skriptiaknas käivitades võib vastavad argumendid – `iter.max` ja `num.seeds` – üldse ära jätta).

Valikuga '*Print cluster summary*' tellitakse olulisem klasterdamise tulemuste kohta käiv info,

valiku '*Bi-plot clusters*'

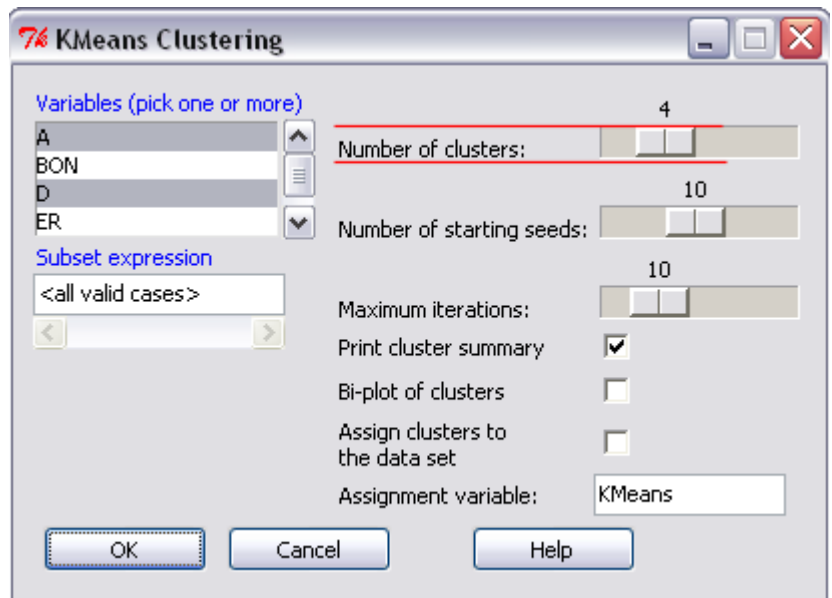
tulemusena teostatakse ette antud tunnustega peakomponentide analüüs ja joonistatakse kahe esimese peakomponendi põhjal hajuvusdiagramm, kus iga objekti märgib klastrite number, millesse vastav objekt on määratud – mõningatel juhtudel võib taoline graafik anda täpsema ettekujutuse sellest, milliste omadustega objektid mingisse klastrisse kuuluvad (et üksnes kahe tunnuse alusel teostatud klasteranalüüsi tulemuste illustreerimiseks on ka lihtsamaid võimalusi, on antud näites see valik tühjaks jäetud),

kolmas lisavalik võimaldab lisada klastrinumbrid algseesse andmebaasi eraldi veeruna (mille nime saab ka eraldi määrata).

Tulemuseks on programm kujul

```
.cluster <- kMeans(model.matrix(~-1 + A + D, puud1), centers=4, iter.max=10, num.seeds=10)
.cluster$size # Cluster Sizes
.cluster$centers # Cluster Centroids
.cluster$withinss # Within Cluster Sum of Squares
.cluster$tot.withinss # Total Within Sum of Squares
.cluster$betweenss # Between Cluster Sum of Squares
remove(.cluster)
```

} Need käsud on valiku '*Print cluster summary*' tulemus.



Kuna *R Commander*

- 1) ei võimalda omistada teostatud klasteranalüüsi tulemustele oma nime ja
- 2) kustutab analüüsi tulemustest moodustatud muutujad peale tellimisaknas valitud analüüside teostamist (automaatselt genereeritav käsk `remove(.cluster)`),

siis tuleb täiendavate analüüside teostamiseks *R Commanderi* poolt genereeritud programm uuesti käivitada, muutes igaks juhuks ka nime, mille all analüüsi tulemusi hoitakse (et mõni *R Commanderi* abil edaspidi teostatav klasteranalüüs tulemusi üle ei kirjutaks):

```
tree_kmean1_RC <- KMeans(model.matrix(~-1 + A + D, puud1), centers=4)
```

Järgnevate käskude puhul on ükskõik, kas kasutada käsuga `kmeans` teostatud klasteranalüüsi tulemusi (muutuja `tree_kmean1`) või *R Commanderi* käsuga `KMeans` teostatud klasteranalüüsi tulemusi (muutuja `tree_kmean1_RC`).

☞ Info, mis klastrisse mingi puu kuulub, on kirjas muutujas `tree_kmean1$cluster`. Näiteks 100 esimese puu klastrid on välja trükitavad käsuga

```
tree_kmean1$cluster[1:100]
```

☞ Klastrite suurused:

```
table(tree_kmean1$cluster)
```

Sama, mis eelnev käsk, üksnes ilma klastrite numbriteta (sama käsk sisaldub ka *R Commanderi* väljastatavas klasteranalüüsi kokkuvõttes)

```
tree_kmean1$size
```

☞ Klastrite keskpunktid (ka see käsk sisaldub *R Commanderi* klasteranalüüsi kokkuvõttes)

```
tree_kmean1$centers
```

```
> tree_kmean1$centers
      A      D
1  54.55353 16.513966
2 119.40579 26.728858
3  80.99565 22.630178
4  24.09864  6.334685
```

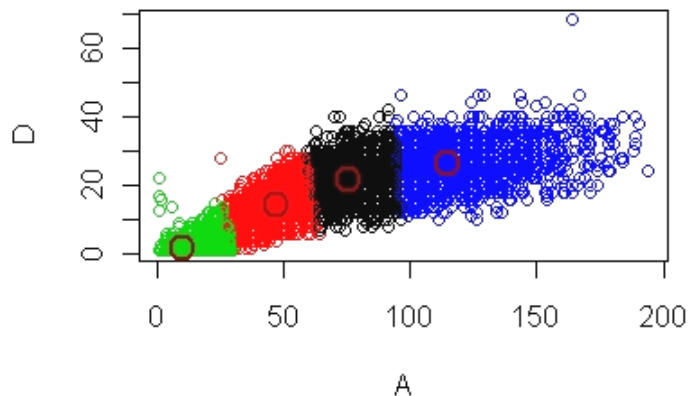
↙                      ↙                      ↙

Klastri number                      Klastrisse kuuluvate puude keskmine vanus                      Klastrisse kuuluvate puude keskmine diameeter

☞ Klasterdades vaatluseid (=puid) vaid kahe tunnuse alusel, on toimuvat võimalik üsna hästi illustreerida graafiliselt (siin esimese käsk joonistab puude vanuse ja diameetri hajuvusdiagrammi, kus erinevatesse klastritesse kuuluvad puud on tähistatud erineva värviga, teine käsk lisab klastrite keskpunktid ringidena):

```
plot(A, D, col=tree_kmean1$cluster)
points(tree_kmean1$centers, cex=2, col="red4", bg=1:4, pch=21, lwd=2)
```

Tulemus:



Veendumaks, et nii käsuga `kmeans` kui ka *R Commanderi* käsuga `KMeans` teostatud klasteranalüüsi tulemused on samad, võite mõlemate kohta tellida joonised ühele lehele:

```
par(mfrow=c(2,1))
plot(A, D, col=tree_kmean1$cluster)
points(tree_kmean1$centers, cex=2, col="red4", bg=1:4, pch=21, lwd=2)

plot(A, D, col=tree_kmean1_RC$cluster)
points(tree_kmean1_RC$centers, cex=2, col="red4", bg=1:4, pch=21, lwd=2)
par(mfrow=c(1,1))
```

☞ Miks tundub joonist vaadates, et klastrid on tehtud puu vanuse (tunnus 'A') järgi?

Vastus – puud paigutati klastritesse eelkõige vanuse alusel põhjusel, et puude vanust märkivad arvud (eelkõige nende varieeruvus) on suuremad, kui diameetrit märkivad arvud. Soovides, et klasterdamisalgoritm peaks puu diameetrit sama oluliseks kui vanust, tuleb klasterdamiseks kasutatavad tunnused eelnevalt standardiseerida (lahutada kõigist väärtustest keskmine ja jagada standardhälbega).

Tunnuseid saab standardiseerida käsu `scale` abil:

```
tree_kmean2 = kmeans(scale(cbind(A, D)), 4)
```

Sama käsku saab kasutada ka muutmaks *R Commanderi* käsku skriptiaknas:

```
tree_kmean2_RC <- KMeans(model.matrix(~-1 + scale(A) + scale(D), puud1), centers=4)
```

Veel ühe alternatiivina võib standardiseeritud tunnused lisada uute tunnustena andmestikku, mida saab *R Commanderis* tellida menüüdest

*Data -> Manage variables in active dataset -> Standardize variables ...*

ja rakendada seejärel klasteranalüüsi juba standardiseeritud tunnustele.

Tulemuste illustreerimiseks:

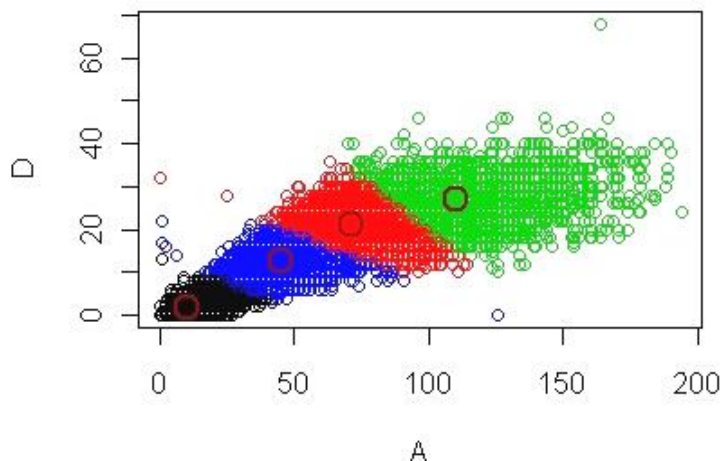
```
plot(A, D, col=tree_kmean2$cluster)
```

Klastrite keskpunktide joonisele kandmiseks tuleb aga teha täiendavaid arvutusi, sest klasteranalüüsiga on leitud klastrite keskpunktid standardiseeritud tunnuste jaoks.

```
meanA = tree_kmean2$centers[,1]*sd(A) + mean(A)
```

```
meanD = tree_kmean2$centers[,2]*sd(D) + mean(D)
```

```
points(meanA, meanD, cex=2, col="red4", bg=1:4, pch=21, lwd=2)
```



Nagu näha, on klasterdamisprogramm peale standardiseerimist käsitlenud puude vanust ja diameetrit samaväärsetena (klastrite piirid sõltuvad mõlemast tunnusest).

☞ Vaatame veelkord, kuidas puud on klasterdunud – sedakorda sõltuvalt arenguklassist.

```
table(tree_kmean2$cluster, puud1$ARENGUKL)
```

```
> table(tree_kmean2$cluster, puud1$ARENGUKL)
```

	A	K	L	N	S	V	Y
1	0	794	0	1	3	624	2152
2	0	5990	1482	0	1	388	241
3	1	159	2527	2566	3	110	41
4	0	4928	1	0	0	1635	1777

Nagu näha, kuuluvad arenguklassid N ja osaliselt ka L ühte klastrisse (pisitillukesed ja noorukesed puud); klastrisse 2 kuuluvatest puudest on enamus arenguklassidest K ja L (küpsust saavutavad puud); klastrisse 1 kuuluvate puude hulgas on märgatavalt palju arenguklassi Y kuuluvaid puud, märkimisväärselt ka arenguklassi V puud ja natuke ka arenguklassi K puud (suured ja vanad puud); 4. klastrisse kuuluvad arenguklasside K ja V enamus, veidi ka arenguklassi Y puud – samuti suured ja vanad, kenasti raieküpsed puud.

### 3. Püüame klasterdada puuliike, kasutades hierarhilist klasterdamist.

Puuliikide klasterdamiseks peame esmalt moodustama andmetabeli, kus objektideks on puuliigid. Lisaks peame otsustama, millised tunnused sellesse andmetabelisse kaasata. Võtame nendeks tunnusteks puuliigi keskmise diameetri 20 ja 50 aasta vanuses ning diameetri juurdekasvu vanuses 60-80.

	D50	D20	D60_80
Seega tahame jõuda andmetabelini kujul	HB 19.49209	6.098114	5.535394
	KS 15.44849	4.799597	4.449109
	KU 15.59229	4.680480	5.265263
	LH 20.27949	6.810392	4.777650
	LM 16.56126	6.799664	4.510453
	LV 16.62557	6.259825	0.000000
	MA 13.73487	4.201077	4.890281
	RE 22.64630	6.816267	0.000000
	SA 16.02317	4.926619	4.814877
	TA 17.33763	5.815752	5.491468

Kuidas sellist andmetabelit koostada? Kuna meil alguses andmetabelis selliseid tunnuseid pole, tuleb nende väärtused prognoosida olemasolevate tunnuste baasil. Näiteks võime iga puuliigi tarvis püüda prognoosida puu diameetrit sõltuvalt tema vanusest 4. järku polünoomi abil (see peaks olema piisavalt täpne võimaldamiseks modelleerida ka mittelineaarset juurdekasvu), leitud mudeli alusel saab juba omakorda prognoosida kõigi vajalike tunnuste väärtused.

☞ Kirjeldatu teostamiseks on üks variant teha arvutused iga puuliigi jaoks eraldi:

```
model_KU = lm(D~A+I(A^2)+I(A^3)+I(A^4), data=puud1[PE=="KU",])
predict(model_KU, data.frame(A=c(20,50,60,80)))
```

Järgnevalt tuleks tulemused välja kirjutada ja jätkata uue puuliigiga.

☞ Alternatiivne variant on lasta kõik teha arvutil.

Et liiga väikese esinemissagedusega puude puhul ei pruugi prognoosid eriti usaldusväärsed tulla, jätame analüüsist välja liigid, mille kohta on vähem kui 10 mõõtmist.

```
species = names(table(PE)[table(PE)>10])
```

Järgnevalt programmeerime tsükli üle kõigi puuliikide (mille arv on tähistatud  $n$ -ga), moodustame esmalt tühjad vektorid kõigi tunnuste tarvis ja täidame need rida reall konkreetsele puuliigile prognoositud väärtustega (konkreetselt puuliiki ja vastavat rida tunnuste vektoreis märgib muutuja  $i$ ). Vältimaks ebareaalset prognoose puuliikidele, mille puhul 80-aastaseid puid polegi, leiame täiendavalt ka maksimaalse vanuse iga puuliigi tarvis ja võrdsustame diameetri juurdekasvu vanuses 60-80 nendel liikidel, kelle maksimaalne vanus  $<80$ , nulliga.

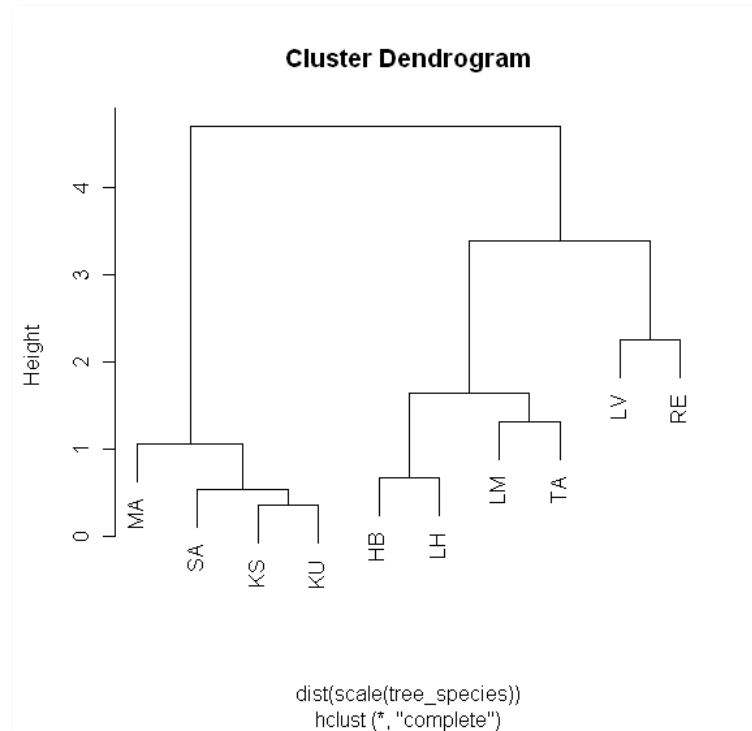
```
n = length(species)
D50=rep(NA,n); D20=rep(NA,n); D60_80=rep(NA,n); AMAX=rep(NA,n)
for (i in 1:n) {
  m1 = lm(D~A+I(A^2)+I(A^3)+I(A^4), data=puud1[PE==species[i],])
  D50[i] = predict(m1, data.frame(A=50))
  D20[i] = predict(m1, data.frame(A=20))
  D60_80[i] = predict(m1, data.frame(A=80)) - predict(m1,
data.frame(A=60))
  AMAX[i] = max(A[PE==species[i]])
}
D60_80[AMAX<80]=0
```

Koondame moodustatud vektorid üheks andmestikuks:

```
tree_species = data.frame(D50, D20, D60_80)
rownames(tree_species)=species
```

☞ Teostame uue andmestikuga hierarhilise klasteranalüüsi ja tellime tulemusena dendrogrammi (joonis, mis esitab klasterduse tulemused puu kujul).

```
tree_hclust1 = hclust(dist(scale(tree_species)))
plot(tree_hclust1)
```



Kui soovime puuliigid jagada kolme klastrisse, saame seda teha käsuga `cutree`:

```
tree_hclust1_3 = cutree(tree_hclust1, 3)
tree_hclust1_3
```

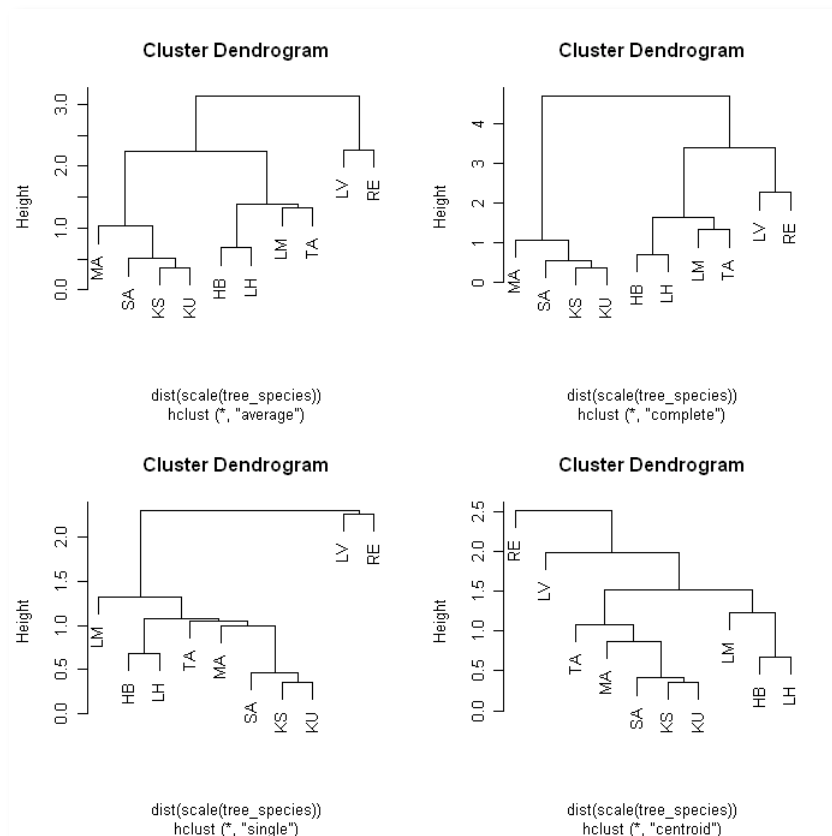
```
> tree_hclust1_3
HB KS KU LH LM LV MA RE SA TA
1 2 2 1 1 3 2 3 2 1
```

Esimesse klastrisse kuuluvad puuliigid haab, sanglepp, lehis ja tamm; teise klastrisse kask, kuusk mänd ja saar; kolmandasse klastrisse remmelgas ja hall lepp.

☞ Järgnevalt võite püüda muuta klasterdamismeetodit ja jälgida, kas ja kuidas tulemused muutuvad.

```
tree_hclust1 = hclust(dist(scale(tree_species)), method="ave")
tree_hclust2 = hclust(dist(scale(tree_species)), method="complete")
tree_hclust3 = hclust(dist(scale(tree_species)), method="single")
tree_hclust4 = hclust(dist(scale(tree_species)), method="centroid")

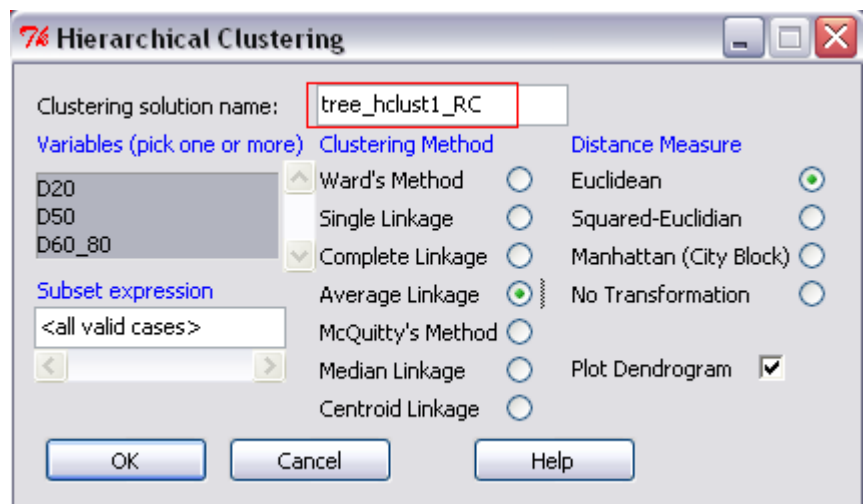
par(mfrow=c(2,2))
plot(tree_hclust1); plot(tree_hclust2); plot(tree_hclust3); plot(tree_hclust4)
par(mfrow=c(1,1))
```



ᄁ) Analooget ülesande võib lahendada *R Commanderi* abil, kusjuures erinevalt k-keskmiste meetodist saab hierarhilise klasterdamise korral omistada tulemusi sisaldavale muutujale ka nime, mille alusel saab konkreetse analüüsi kohta hiljem lisaanalüüsi tellida.

#### Statistics

- > Dimensional analysis
- > Cluster analysis
- > Hierarchical cluster analysis ...



Hoolimata sellest, et hierarhilisel klasterdamisel kasutab *R Commander* funktsiooni `hclust` nii nagu ka *R*-i baasversioon, võivad tulemused siiski erineda tulla. Võrdle näiteks:

```
tree_hclust1 = hclust(dist(scale(tree_species)), method="ave")
tree_hclust1_RC = hclust(dist(model.matrix(~-1 + D20+D50+D60_80,
tree_species)), method="average")
```

```
par(mfrow=c(2,1))
plot(tree_hclust1); plot(tree_hclust1_RC)
par(mfrow=c(1,1))
```



Põhjus on selles, et *R Commander* ei kasuta vaikimisi tunnuste standardiseerimist ega võimalda seda teha ka analüüsi tellimise käigus. Resultaadina mõjutavad suuremate väärtustega ja enam varieeruvad tuunused klasterdamise tulemusi enam ...

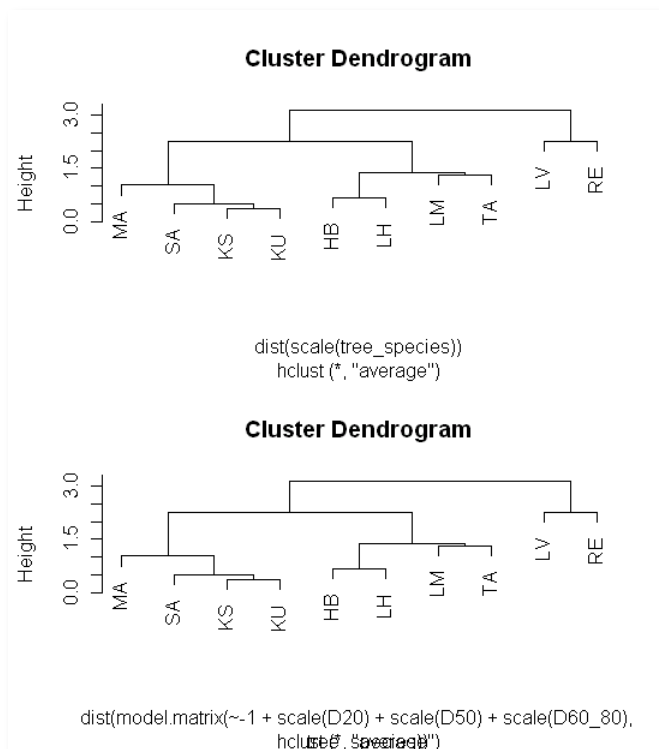
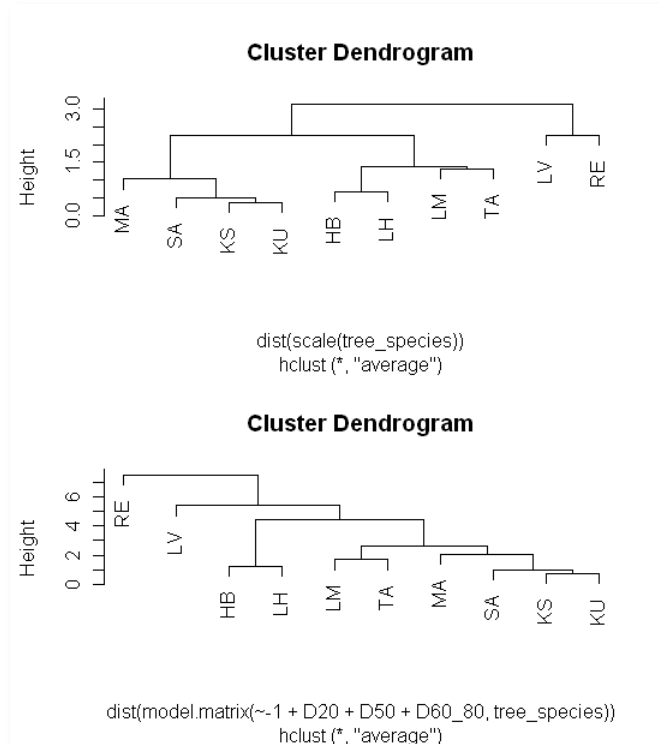
Lahenduseks on kas tunnuste eelnev standardiseerimine või siis *R Commanderi* poolt väljastatud programmi käsitsi täiendamine.

Viimane tähendab funktsiooni `scale` lisamist kõigi tunnuste ette:

```
tree_hclust1_RCsc = hclust(dist(model.matrix(~-1 + D20 + D50 + D60_80, tree_species))
  + scale(D50) + scale(D60_80), tree_species), method="average")
```

Nüüd peaksid tulemused tulema juba analoogsed. Soovi korral vaata järgi:

```
par(mfrow=c(2,1))
plot(tree_hclust1); plot(tree_hclust1_RCsc)
par(mfrow=c(1,1))
```



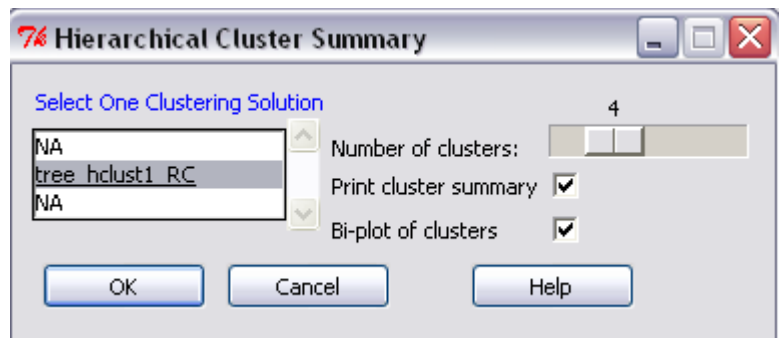
☞ Lisaks võimaldab *R Commander* valida menüüdest mõningaid lisaanalüüse hierarhilise klasterdamise tulemuste alusel. Halb on see, et aktsepteerib *R Commander* vaid enese poolt teostatud analüüse, skripti aknas jooksvatuid klasteranalüüside kohta menüüdest midagi lisaks tellida ei saa. Küll aga on võimalik menüüdest tellitud lisaanalüüside programmides asendada muutujate nimesid ja selliselt muudetud programmid skripti aknas uuesti käivitada ...

Näiteks *R Commanderiga* teostatud klasteranalüüsi tulemuste `tree_hclust1_RC` kohta lisaanalüüside tellimiseks tuleb menüüst

*Statistics -> Dimensional analysis -> Cluster analysis*

valida käsk

*Summarize Hierarchical clustering ...*



Jagades puuliigid nelja klastrisse, saame lisavaliku '*Print cluster summary*' tulemusena puliikide arvud ja keskmised klasterdamise aluseks olnud tunnuste väärtused klastrite kaupa:

```
> summary(as.factor(cutree(tree_hclust1_RC, k = 4))) # Cluster Sizes
 1 2 3 4
2 6 1 1
```

ja

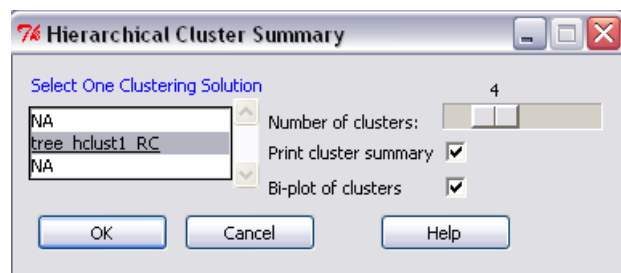
```
> by(model.matrix(~-1 + D20 + D50 + D60_80, tree_species), as.factor(cutree(tree_hclust1_RC, k = 4)), mean)
INDICES: 1
      D20      D50      D60_80
6.634692 19.889509  5.157483
-----
INDICES: 2
      D20      D50      D60_80
5.195316 15.852870  4.755391
-----
INDICES: 3
      D20      D50      D60_80
6.378851 16.665620  0.000000
-----
INDICES: 4
      D20      D50      D60_80
6.673248 22.431491  0.000000
```

NB! Need tulemused vastavad *R Commanderi* poolt teostatud standardiseerimata tunnustega analüüsile (klastrisse 1 kuulub ainult kaks puuliiki, kuigi peaks kuuluma neli liiki – vrdl jooniseid eelmisel leheküljel)!

Standardiseeritud tunnuste alusel moodustatud klastrite tarvis tuleb eelnevalt toodud väljatrükkis näha olevais programmides asendada *R Commanderiga* konstrueeritud mudeli nimi '`tree_hclust1_RC`' asendada ise skriptiaknast käivitatud mudeliga '`tree_hclust1`' (või '`tree_hclust1_RCsc`').

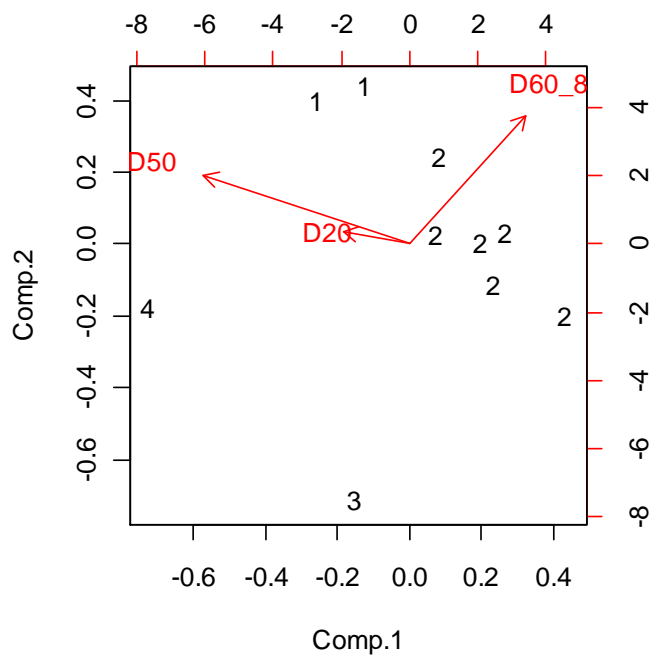
☞ Lisavaliku 'Bi-plot clusters' tulemusena teostab *R Commander* klasteranalüüsi aluseks olnud tunnustega

**peakomponentanalüüsi** ja väljastab kahe esimese peakomponendi laadungite – so peakomponentide ja algsete tunnuste vaheliste korrelatsioonide – joonise, millel on erinevate numbritega märgitud erinevatesse klastritesse kuuluvate puuliikide paigutus peakomponentide mõistes.



Esimene peakomponent on negatiivselt seotud diameetriga 20 ja 50 aasta vanuses ning positiivselt juurdekasvuga vanuses 60-80, omades seega väikeseid väärtusi juhul, kui puuliik kasvab kiiresti, saavutades maksimumi lähedase läbimõõdu juba 50.-ks eluaastaks (juurdekasv vanuses 60-80 on väike), ning suuri väärtusi juhul, kui algne kasv on suhteliselt aeglane, aga kasv vanuses 60-80 märgatav.

Teine peakomponent väljendab üldist kasvukiirust, olles positiivselt korreleeritud kõigi tunnustega.



Klastrite paigutus näitab seda, et 3. ja 4. klastrisse kuuluvad liigid, mis on suhteliselt kiirekasvulised ja ei ela enamasti kaua – hall lepp ja remmelgas – vt `tree_hclust1_RC` kohta joonistatud dendrogrammi või siis uuri liikide klastritesse kuuluvust käsuga

```
cutree(tree_hclust1_RC, k=4);
```

klastrisse 2 kuuluvad liigid kasvavad vanuses 60-80 kõige rohkem ja klastrisse 1 kuuluvad liigid võib liigitada keskmise kasvukiirusega liikide hulka.

NB! Ka toodud peakomponentanalüüsi tulemused on *R Commander* leidnud standardiseerimata tunnuste alusel, mistap võib tunnuse 'D20' nõrgem seos leitud peakomponentidega (lühem nool ülaltoodud joonisel) peegeldada hoopis seda, et tunnuse 'D20' mõõteskaala on väiksem ja roll peakomponentide moodustamisel seeläbi ka väiksem ...  
Pisut täpsemalt juba järgmises punktis.

4. Rakendame peakomponentanalüüsi puuliikidele arvutatud karakteristikutele 'D20', 'D50' ja 'D60\_80'.

Peakomponentanalüüsi rakendamiseks võib kasutada nii *R Commanderi* vastavat menüükäsku kui ka trükkida vajalik käsk otse skriptiaknasse.

# 1. võimalus – käsuna skriptiaknas:

```
tree_PC <- princomp(~D20+D50+D60_80, cor=TRUE, data=tree_species)
```

Lisavalik `cor=TRUE` käseb *R*-l teostada peakomponentide analüüs tunnuste korrelatsioonimaatriksi alusel, mis seeläbi elimineerib peakomponentanalüüsi tulemuste potentsiaalse sõltuvuse argumentide erinevast mõõteskaalast.

☒) Käsuga

```
summary(tree_PC)
```

on tellitav esmane kokkuvõtte analüüsi tulemustest:

```
> summary(tree_PC) # proportions of variance
Importance of components:
              Comp.1   Comp.2   Comp.3
Standard deviation  1.4086273 0.8683579 0.5115893
Proportion of Variance 0.6614103 0.2513485 0.0872412
Cumulative Proportion 0.6614103 0.9127588 1.0000000
```

Esimene peakomponent kirjeldab ära 66,1% algsete tunnuste koguvarieeruvusest, teine komponent 25,1% ja kolmas komponent 8,7%. Kokku kirjeldavad kaks esimest peakomponenti ära 91,3% kolme analüüsi kaasatud tunnuse koguvarieeruvusest.

☒) Käsuga

```
tree_PC$sd^2
```

on tellitavad erinevate peakomponentide dispersioonid (korrelatsioonimaatriksi omaväärtused):

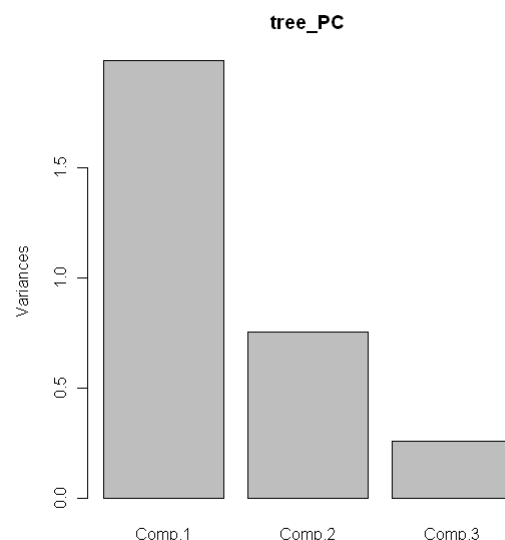
```
> tree_PC$sd^2 # component variances
  Comp.1   Comp.2   Comp.3
1.9842310 0.7540454 0.2617236
```

Ühest suurem on vaid esimese peakomponenti dispersioon, mistap sisaldab vaid esimene peakomponent enam infot kui üksiktunnused.

☒) Käsk

```
plot(tree_PC)
```

esitab peakomponentide varieeruvuse graafiliselt:



☒) Käsk

```
unclass(loadings(tree_PC))
```

annab tulemuseks peakomponentide laadungid (korrelatsioonid peakomponentide ja algsete tunnuste vahel):

```
> unclass(loadings(tree_PC)) # component loadings
      Comp.1  Comp.2  Comp.3
D20    -0.6260724 0.3580368 0.69270700
D50    -0.6378554 0.2758336 -0.71906628
D60_80 0.4485240 0.8920345 -0.05568388
```

Sarnaselt klasteranalüüsi lõpus leituga on esimene peakomponent negatiivselt seotud diameetriga 20 ja 50 aasta vanuses ning positiivselt juurdekasvuga vanuses 60-80, omades seega väikeseid väärtusi juhul, kui puuliik kasvab kiiresti, saavutades maksimumi lähedase läbimõõdu juba 50.-ks eluaastaks (juurdekasv vanuses 60-80 on väike), ning suuri väärtusi juhul, kui algne kasv on suhteliselt aeglane, aga kasv vanuses 60-80 märgatav.

Teine peakomponent peegeldab üldist ja pisut hilisemat kasvukiirust, olles positiivselt korreleeritud kõigi tunnustega, aga eelkõige juurdekasvuga vanuses 60-80 aastat.

⌘) Käsk

```
biplot(tree_PC)
```

esitab samal joonisel nii viimatised korrelatsioonid kui ka, vastavalt oma peakomponentide väärtustele, üksikud vaatlused (antud juhul puuliigid):

Joonise alumine ja vasakpoolne telg vastavad korrelatsioonikordajatele ning nende alusel on joonisele kantud algsete tunnuste ja kahe esimese peakomponendi vahelist seost illustreerivad nooled;

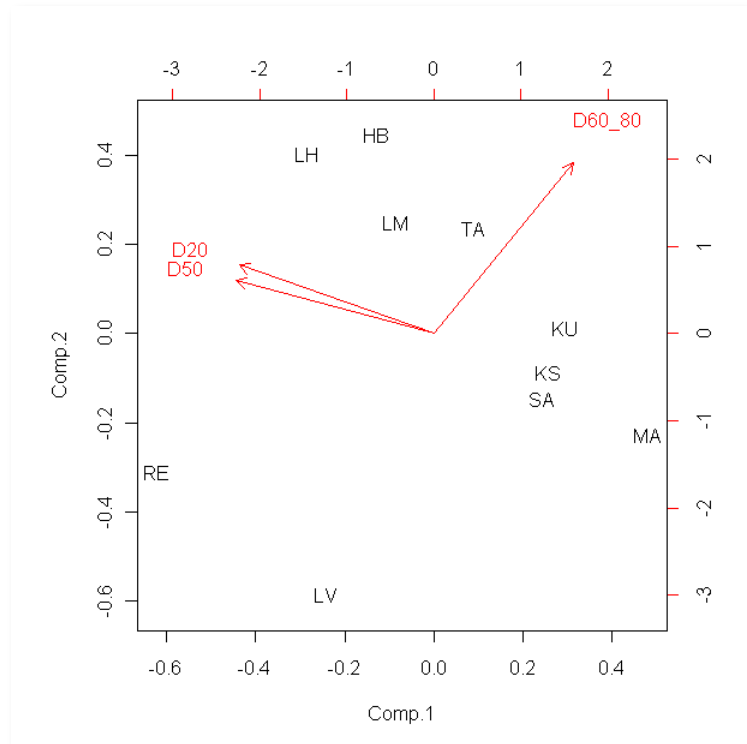
joonise parempoolne ja ülemine telg vastavad aga peakomponentide väärtustele ja nende alusel on joonisele kantud kõik andmebaasi vaatlused (antud juhul puuliigid, kusjuures vaikumisi kasutatavaks tähiseks joonisel on vastava rea nimi).

Järeldused üksikute puuliikide kasvamise kohta on samuti osaliselt toodud juba klasteranalüüsi lõpus. Halli lepa (LV) ja remmelga (RE) nii esimese kui ka teise peakomponendi väärtused on

negatiivsed, viidates sellele, et tegu on vanuses 20 ja 50 aastat kiirelt, aga vanuses 60-80 aastat aeglaselt (või enam üldse mitte) kasvavate liikidega (vastavalt peakomponentide ja algsete tunnuste vahelistele korrelatsioonidele on esimese peakomponendi väärtused seda väiksemad, mida suuremad on tunnuste 'D20' ja 'D50' väärtused, ning teise peakomponendi väärtused seda väiksemad, mida väiksemad on tunnuse 'D60\_80' väärtused).

Kuusk (KU), kask (KS), saar (SA) ja mänd (MA) kasvavad vanuses 20 ja 50 aastat suhteliselt aeglaselt (nende liikide esimese peakomponendi väärtused on positiivsed) ega hiilga jõudsa juurdekasvuga ka vanuses 60-80 aastat (teise peakomponendi väärtused on negatiivsed).

Haab (HB), lehis (LH) ja sanglepp (LM) kasvavad alguses suhteliselt kiiresti (esimese peakomponendi väärtused on negatiivsed) ja säilitavad oma diameetri suhteliselt kiire suurenemise ka vanuses 60-80 (teise peakomponendi väärtused on positiivsed).



☞) Käskudega

```
tree_species$PC1 <- tree_PC $scores[,1]
```

```
tree_species$PC2 <- tree_PC $scores[,2]
```

on vastavalt esimese ja teise peakomponendi väärtused (skoorid; antud juhul nimede 'PC1' ja 'PC2' all) lisatavad ka andmetabelile nende võimaliku edasise analüüsimise tarvis.

# 2. võimalus – *R Commander*'i menüüdest on peakomponentide analüüs tellitav järgmiselt:  
*Statistics -> Dimensional analysis -> Principal-components analysis ...*

Kuigi käsu süntaks on natuke erinev eelnevalt skriptiaknasse sisestatud, on tulemused identsed. St, et erinevalt hierarhilise klasterdamise nõ lisanalüüsis kasutab *R Commander* eraldi peakomponentanalüüsi tehes siiski algsete tunnuste vahelist korrelatsioonimaatriksit, mistap ei mängi tunnuste pisut erinev mõõteskaala rolli.

