

# Domača naloga 1

Neza Krzan, Tom Rupnik

## Kazalo

<b>1 Cilji naloge</b>	<b>3</b>
<b>2 Podatki</b>	<b>3</b>
2.1 Urejanje podatkov . . . . .	3
2.2 Analiza povezanosti med spremenljivkami . . . . .	4
2.3 Konstrukcija in analiza Likertovih spremenljivk . . . . .	5
2.4 Povezanost Likertovih spremenljivk . . . . .	7
2.5 Likertovi spremenljivki in tip . . . . .	7
<b>3 Hierarhično razvrščanje</b>	<b>9</b>
3.1 Wardova metoda . . . . .	9
3.2 Minimalna metoda . . . . .	9
3.3 Maksimalna metoda . . . . .	10
3.4 Analiza . . . . .	10
<b>4 Nehierarhično razvrščanje</b>	<b>12</b>
4.1 Razvrščanje K-means . . . . .	12
4.2 GAP statistika . . . . .	12
4.3 Pseudo F (Calinski - Harabasz indeks) . . . . .	13
4.4 Primerjava vrednosti kriterijske funkcije za Wardovo metodo in K-means . . . . .	14
<b>5 Razvrščanje na podlagi modelov</b>	<b>15</b>
5.1 BIC(Bayes Information Criterion) kriterij . . . . .	15
5.2 BIC kriterij na standardiziranih podatkih . . . . .	16
<b>6 Najboljša razvrstitev in predstavitev skupin</b>	<b>18</b>

## Slike

1 Porazdelitve spremenljivk v podatkovnem okviru <b>Swiss banknotes data</b> . . . . .	4
2 Korelacija med spremenljivkami. . . . .	4
3 Korelacija med spremenljivkami in možna povezava med spremenljivkami. . . . .	5
4 Povezanost med kategorično spremenljivko(tip bankovca) in števili spremenljivkami. . . . .	5
5 Porazdelitvi Likertovih spremenljivk. . . . .	6
6 Razsevni diagram dolžine in mere. . . . .	7
7 Povprečja dolžina (levo) in mere (desno) po skupinah spremenljivke tip. . . . .	7
8 Razsevni diagram dolžina in mere po spremenljivki tip. . . . .	8
9 Dendrogrami Wardove metode razvrščanja v skupine. . . . .	9
10 Dendrogrami minimalne metode razvrščanja v skupine. . . . .	10
11 Dendrogrami maksimalne metode razvrščanja v skupine. . . . .	10
12 Povprečja po skupinah za Wardovo metodo. . . . .	11
13 Vrednost Wardove kriterijske funkcije. . . . .	12

14	Vrednost GAP statistike. . . . .	13
15	Vrednost Pseudo F oz. Calinski - Harabasz indeksa. . . . .	13
16	Porazdelitve spremenljivk. . . . .	15
17	BIC kriterij za originalne podatke. . . . .	16
18	BIC kriterij (priorControl) za originalne podatke. . . . .	16
19	BIC kriterij za standardizirane podatke. . . . .	17
20	BIC kriterij (priorControl) za standardizirane podatke. . . . .	17

## Tabele

# 1 Cilji naloge

V nalogi vova poskušala razvrstiti enote v skupine tako, da si bodo enote znotraj skupin čim bolj podobne in enote v različnih skupinah čim bolj različne glede na več spremenljivk.

## 2 Podatki

Uporabila bova podatke *Swiss banknotes data*, ki vsebujejo šest meritev, opravljenih na 100 pravih in 100 ponarejenih starih švicarskih bankovcih za 1000 frankov.

Podatki vsebujejo 7 spremenljivk - 6 številskih in eno opisno. Vsebujejo različne izmerjene dolžine in širine bankovca v milimetrih:

- **length**: dolžina bankovca (na sliki  $x_1$ ),
- **left**: dolžina levega roba (na sliki  $x_2$ ),
- **right**: dolžina desnega roba (na sliki  $x_3$ ),
- **bottom**: dolžina spodnjega roba (na sliki  $x_4$ ) in
- **top**: dolžina zornjega roba (na sliki  $x_5$ ) ter
- **diag**: dolžina diagonale bankovca (na sliki  $x_6$ ).

Opisna spremenljivka **status** pa določa ali je bankovec pravi (**genuine**) ali ponarejen (**counterfeit**). V tabeli imamo torej meritve za 200 različnih bankovcev.

### 2.1 Urejanje podatkov

Imena spremenljivk in vrednosti kategorične spremenljivke sva preimenovala v slovenska imena ter, kot sva že napisala zgoraj, sva podatke skalirala.

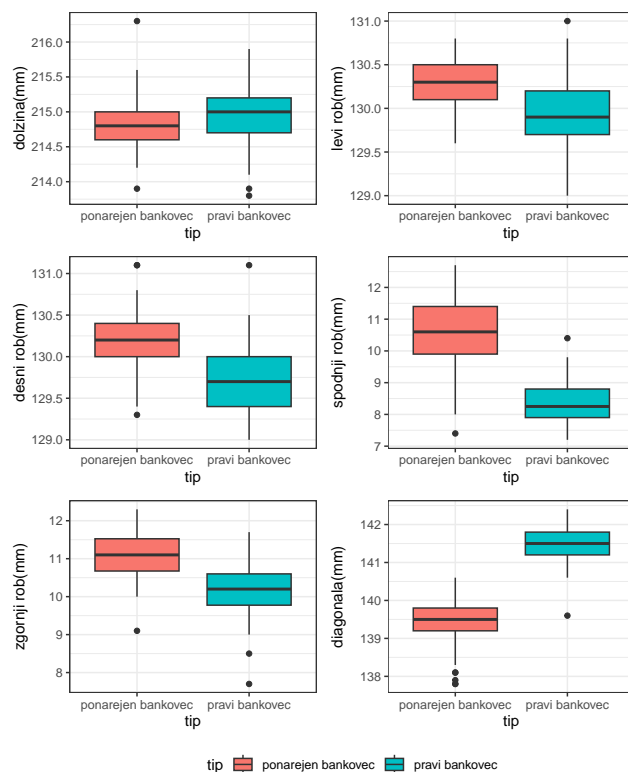
Preimenovane spremenljivke:

- **length**: dolžina,
- **left**: levi.rob,
- **right**: desni.rob,
- **bottom**: spodnji.rob,
- **top**: zgornji.rob,
- **diag**: diagonala in
- **status**: tip, kjer je potem **counterfeit**:ponarejen bankovec in **genuine**:pravi bankovec.

Za lažjo predstavo si pogledajmo opisne statistike številskih spremenljivk, da bomo vedeli s kakšnimi podatki imamo opravka.

Spremenljivke imajo različen razpon vrednosti, zato jih bova, skalirala. Tako bodo imele spremenljivke povprečje 0 in standardni odklon 1. S tem doseževa enakovreden vpliv spremenljivk na razvrstitev. Vidimo pa tudi, da nimamo manjšajočih vrednosti v podatkih.

Poglejmo si še porazdelitve spremenljivk.

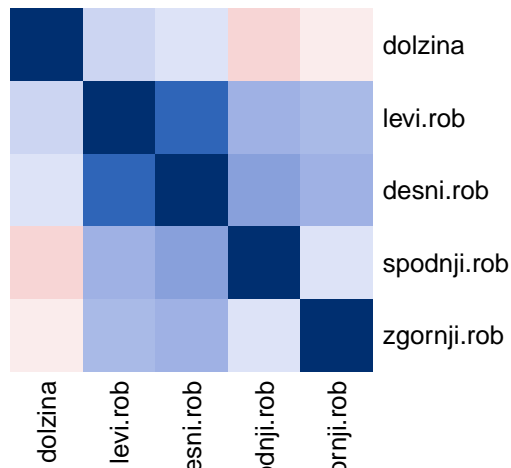


Slika 1: Porazdelitve spremenljivk v podatkovnem okviru `Swiss banknotes data`.

Opazna je razlika med pravimi bankovci in ponarejenimi pri vseh spremenljivkah.

Za razvrščanje bova uporabljala samo številske spremenljivke, in sicer `dolzina`, `levi.rob`, `desni.rob`, `spodnji.rob`, `zgornji.rob`; za analizo pa spremenljivki `tip` in `diagonala`. Ker je `diagonala` edina številska spremenljivka pri analizi, le ta ne bo skalirana.

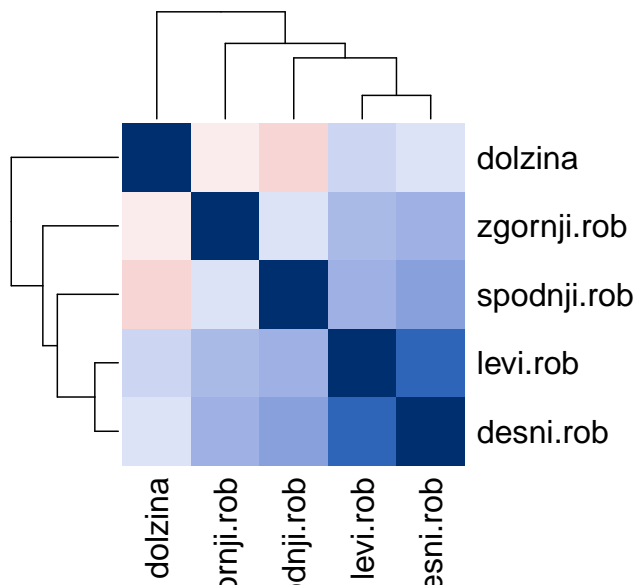
## 2.2 Analiza povezanosti med spremenljivkami



Slika 2: Korelacija med spremenljivkami.

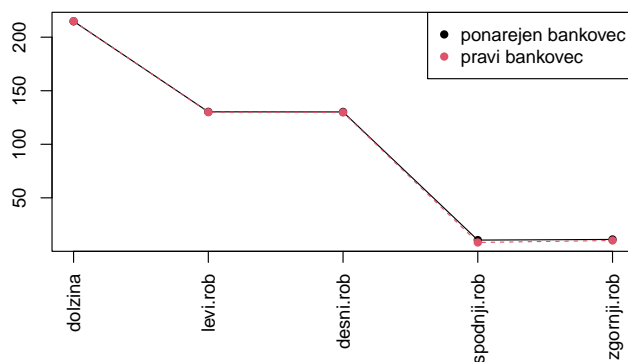
Iz grafa vidimo, da sta levi in desni rob med seboj močno povezana. Prav tako imata zgornji in spodnji rob z njima srednje močno korelacijo. Dolzina ima z vsemi zelo sibko korelacijo iz česar bi lahko zaključili, da tvori svojo skupino.

To je še bolje razvidno iz spodnjega grafa.



Slika 3: Korelacija med spremenljivkami in možna povezava med spremenljivkami.

Poglejmo si povezanost med kategorično spremenljivko tip bankovca in številskimi spremenljivkami za razvrščanje.



Slika 4: Povezanost med kategorično spremenljivko (tip bankovca) in številskimi spremenljivkami.

Torej sklepava lahko, da spremenljivka **dolzina** tvori svoj sklop, ostale spremenljivke pa svojega iz zgornje analize in grafov.

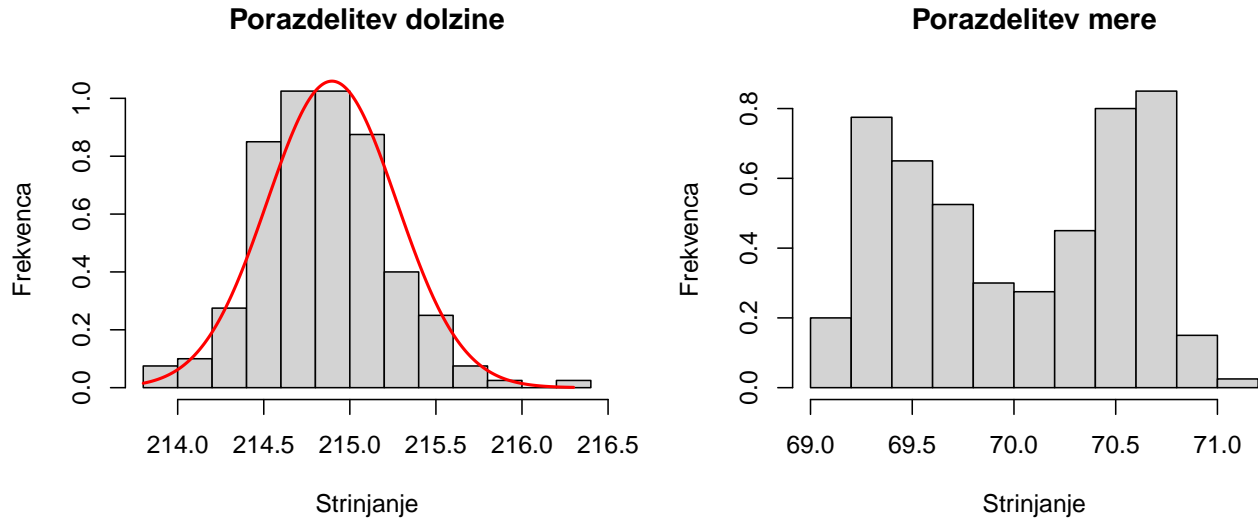
## 2.3 Konstrukcija in analiza Likertovih spremenljivk

Sedaj bomo za spremenljivke izdelali dve Likertovi lestvici, torej razdelila sva jih v dve skupini, tako da znotraj skupine spremenljivke merijo približno enako stvar. Za vsako skupino sva ustvarili eno Likertovo lestvico, tako da za vsako opazovano enoto izračunamo povprečje izbranih spremenljivk.

- Prvi sklop predstavlja dolžino bankovca in vsebuje **dolzina**.

- Druga sklop predstavlja mere in vsebuje `levi.rob`, `desni.rob`, `spodnji.rob`, `zgornji.rob`.

Poglejmo si porazdelitvi novo ustvarjenih spremenljivk Likertove lestvice tj. `dolzina` in `mere`.



Slika 5: Porazdelitvi Likertovih spremenljivk.

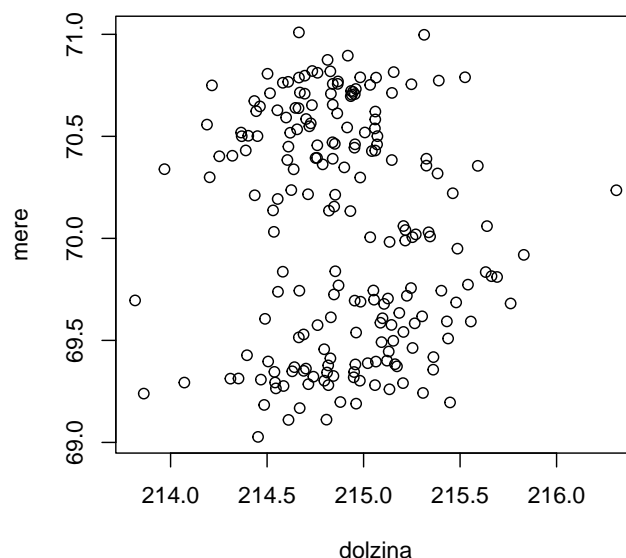
Porazdelitvi `dolzina` in `mere` si nista podobni. Povprečna nagnjenost k `dolzini` je 214.9 s standardnim odklonom 0.38 in je dokaj simetrična. Povprečna nagnjenost k `meri` je 70.04, medtem ko je standardni odklon enak 0.56 in njena porazdelitev je precej nedefinirana.

Da preverimo domnevo o enakosti povprečij `dolzina` in `mere` uporabimo t-test za odvisna vzorca z ničelno domnevo, da je razlika povprečij enaka 0. Ta vrne vrednost  $p < 0.001$ , torej lahko ničelno domnevo pri stopnji značilnosti 0.05 zavrnemo. 95% interval zaupanja za razliko povprečij med `dolzino` in `mero` je  $[-144.96, -144.76]$ .

Pogledala bova še ali sta ustvarjeni spremenljivki Likertove lestvice povezani med seboj in ali sta povezani z ostalimi spremenljivkami.

## 2.4 Povezanost Likertovih spremenljivk

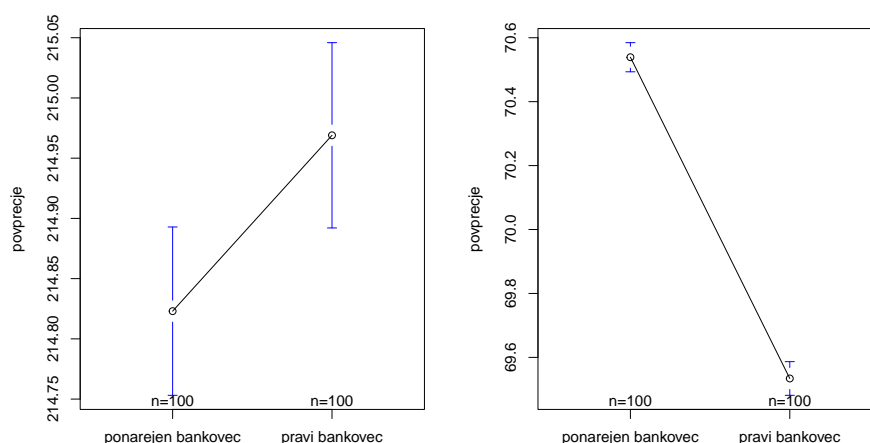
Grafično si oglejmo razsevni diagram dolzine in mer.



Slika 6: Razsevni diagram dolzine in mer.

Da sva preverila ničelno povezanost, sva uporabila korelacijski test za odvisna vzorca na osnovi Pearsonovega koeficienta korelacije, ki vrne kot oceno za Pearsonov korelacijski koeficient vrednost -0.0795, ampak vrednost 0 pa pade v 95% IZ ter vrednost  $p = 0.26 > 0.05$ , torej ne moreva pri stopnji značilnosti 5% trditi, da povezanost obstaja na populaciji.

## 2.5 Likertovi spremenljivki in tip



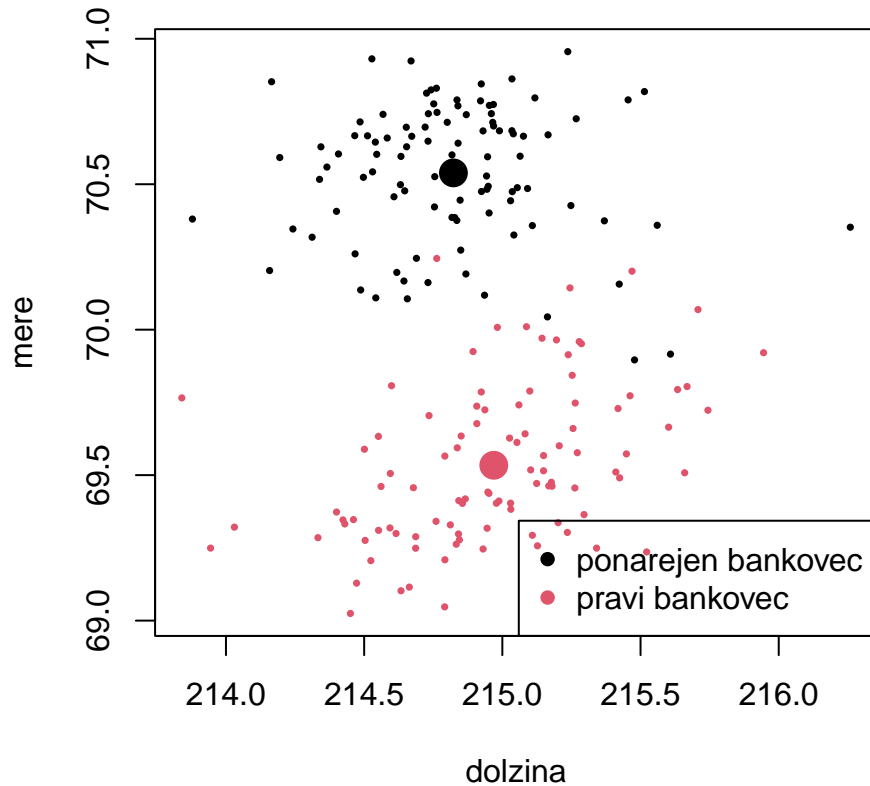
Slika 7: Povprečja dolzina (levo) in mere (desno) po skupinah spremenljivke tip.

Tukaj sva uporabljala Welchov test za primerjavo povprečij, ki pokaže, da se skupini tipa bankovcev pri stopnji značilnosti 5% v povprečju statistično značilno razlikujeta pri obeh Likertovi spremenljivki, ker je  $p < 0.05$ .

Ocena za povprečje pri dolzini za ponarejen tip bankovca je 214.82 in za pravi bankovec 214.97 pri 95%

intervalu zaupanja  $[-0.25, -0.043]$ . Ocena za povprečje pri mere pa je za ponarejen tip bankovca 70.54 in za pravi bankovec 69.53 pri 95% intervalu zaupanja  $[0.94, 1.07]$ .

Nariševa še razsevni diagram spremenljivke **tip**, kjer za osi vzameva Likertovi lestvici, točke pa predstavljajo posamezne enote. Dodamo tudi dve večji točki, ki predstavljata povprečji po skupinah obeh Likertovih spremenljivk.



Slika 8: Razsevni diagram dolzina in mere po spremenljivki tip.

Iz grafa opazimo razlike med tipom bankovca, ampak obstajajo enote, ki se malenkost pomešajo.

Torej spremenljivke sva razdelila v dva sklopa glede na povezanost in tako ustvarila dve Likertovi spremenljivki. Prva spremenljivka predstavlja dolžino, druga pa ostale mere. S 95% zaupanjem lahko trdimo, da dolžina in mere nista statistično značilno povezani. Glede na tip bankovca se skupini statistično razlikujeta pri obeh Likertovi spremenljivki.



### 3 Hierarhično razvrščanje

Pri hierarhičnem razvrščanju začnemo s tem, da je vsaka enota v svoji skupini. Potem pa se na vsakem koraku, glede na izračunane matrike različnosti, v kateri so razdalje med pari skupin, združujejo skupine, ki so si najbližje. Nato se izračunajo različnosti novih združenih skupin od ostalih, kar se nadaljuje dokler niso vse enote v eni skupini. Dobra lastnost hierarhičnega razvrščanja je, da uporabniku ni potrebno vnaprej določiti števila skupin.

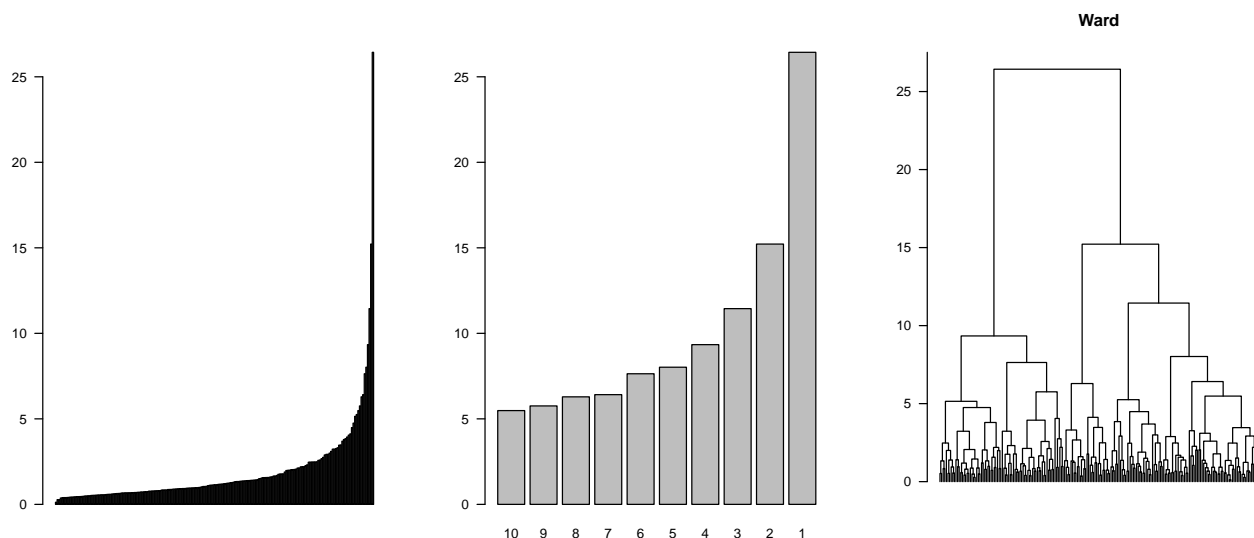
Kot mero različnosti bova uporabila evklidsko razdaljo.

Torej za razvrščanje uporabljava spremenljivke `dolzina`, `levi.rob`, `desni.rob`, `spodnji.rob` in `zgornji.rob` ter primerjala bova tri različne metode in sicer, Wardovo metodo, minimalno metoda (single linkage) in maksimalno metoda (complete linkage).

Število skupin lahko določimo na podlagi dendograma, ki grafično prikazuje potek združevanja v skupine. Število skupin pa določimo tako na podlagi vidnejšega zmanjšanja razdalj skupinami.

#### 3.1 Wardova metoda

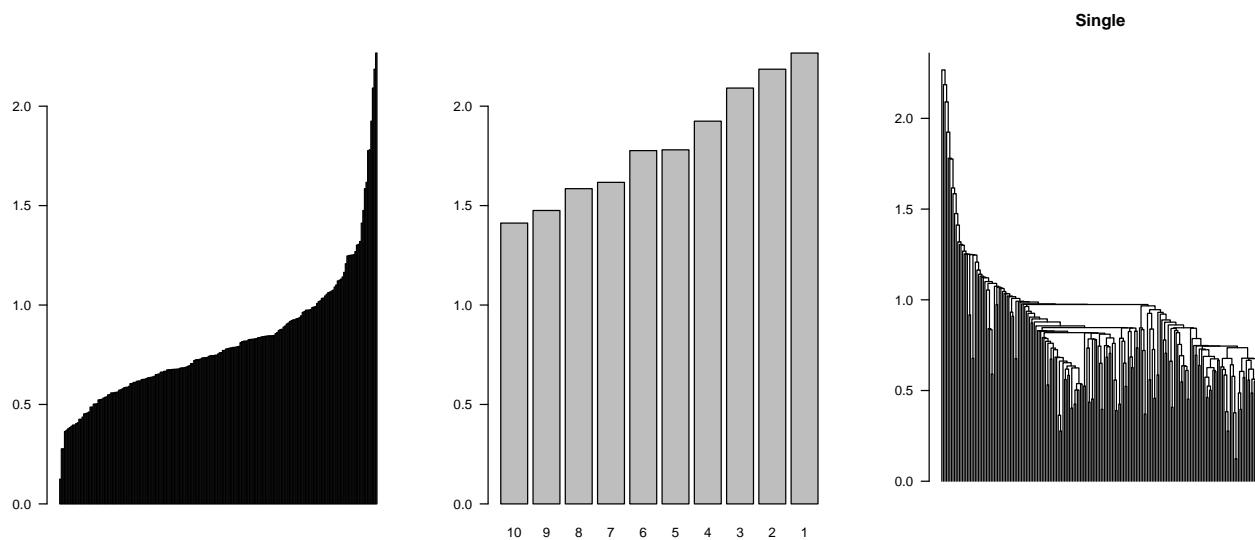
Wardova metoda je primerna za eliptične skupine.



Slika 9: Dendogrami Wardove metode razvrščanja v skupine.

#### 3.2 Minimalna metoda

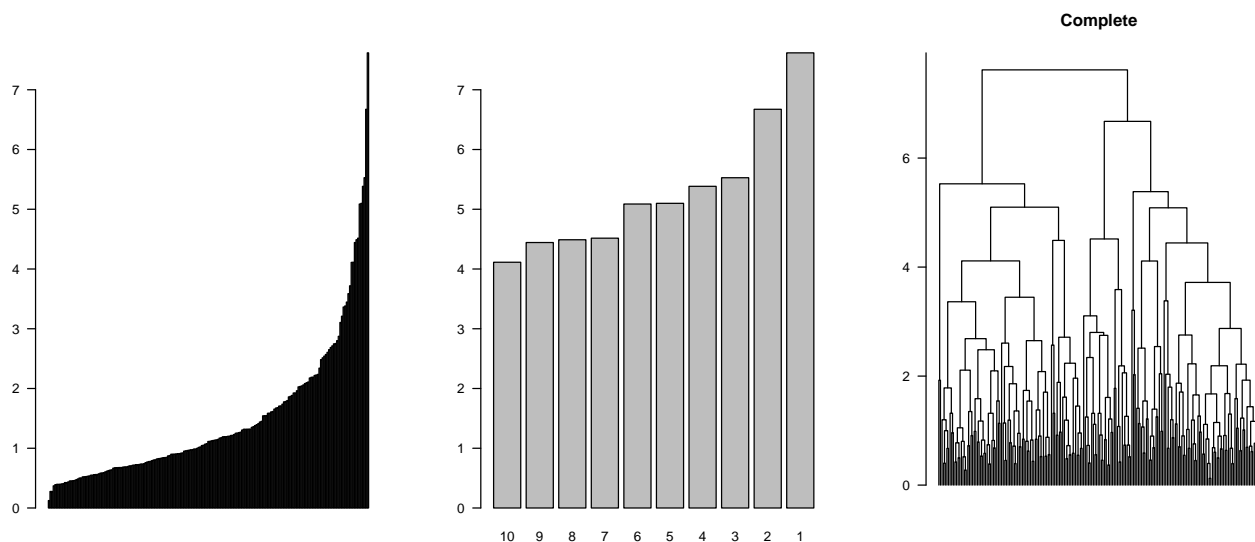
Minimalna metoda (enojna povezanost - single linkage) je primerna za dolge in neeliptične skupine, ki so jasno ločene med seboj. Kadar skupine med seboj niso jasno ločene pri minimalni metodi pride do problema veriženja. Na takem dendogramu ne moremo določiti števila skupin in zato rečemo, da je skupina zgolj ena.



Slika 10: Dendogrami minimalne metode razvrščanja v skupine.

### 3.3 Maksimalna metoda

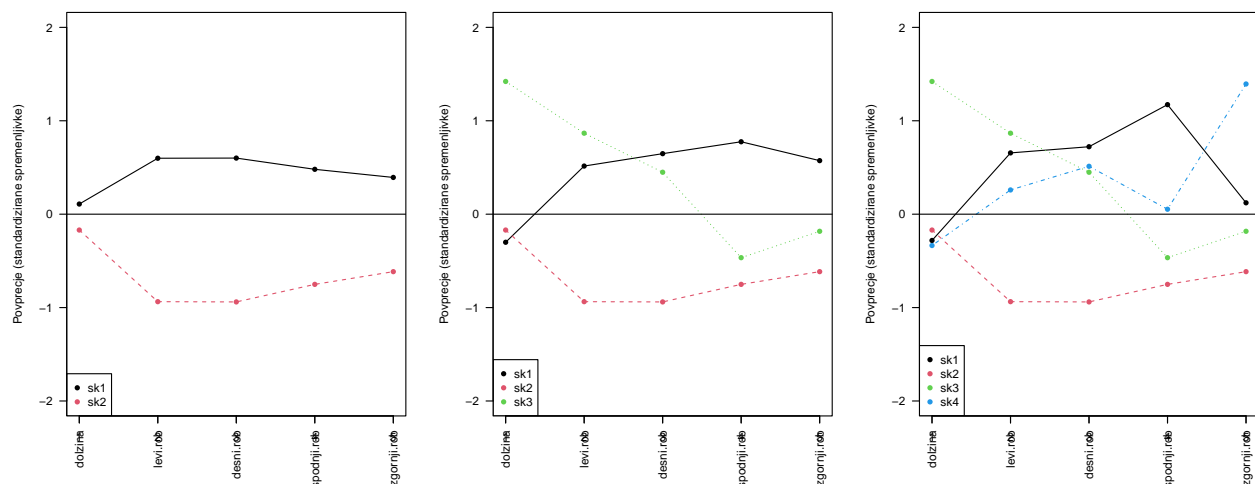
Maksimalna metoda (polna povezanost - complete linkage) pa je primerna za okrogle skupine.



Slika 11: Dendogrami maksimalne metode razvrščanja v skupine.

### 3.4 Analiza

Glede na izgled grafov (razvrstitve) sva se odločila, da je najbolj primerna razvrstitev po Wardovi metodi. Pri ostalih dveh metodah so različnosti dokaj majhne (ni tako izrazitih skokov v višini). Grafe bomo narisali za 2, 3 in 4 skupine, saj so tu razlike bolj izrazite.



Slika 12: Povprečja po skupinah za Wardow metodo.

Če si pogledamo skupino 2 na vseh treh grafih, vidimo, da zavzema podpovprečne vrednosti. Ravno obratno vidimo pri skupini 1, ki na prvem grafu zavzema nadpovprečne vrednosti, na drugih dveh pa zavzema podpovprečne vrednosti samo pri dolžini bankovca. Skupina 3 pa je v nekaterih primerih nadpovprečna v nekaterih pa podpovprečna (spodnji.rob, zgornji.rob). Pri zadnjem grafu se skupina 4 pri spremenljivki *dolzina* približa povprečju zelo dobro, pri vseh ostalih spremenljivkah je nadpovprečna in pri zadnji močno podpovprečna.

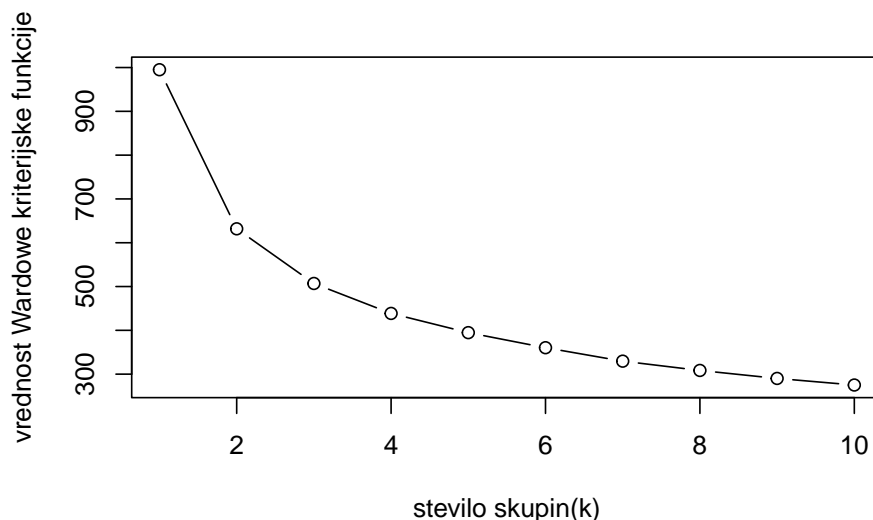
## 4 Nehierarhično razvrščanje

### 4.1 Razvrščanje K-means

K-means je metoda voditeljev oz. nehierarhičnega razvrščanja. Voditelji so “predstavniki skupin”, vsaka enota pa pripada skupini, kateremu voditelju je najbližje (razdalja je evklidska) oz. mu je najbolj podobna; voditelj predstavlja povprečje skupine. Spremenljivke pri metodi *k-means* morajo biti vsaj intervalne.

Tukaj pri tej metodi mora biti število skupin podano vnaprej, kar je morda slaba lastnost in se glede tega razlikuje od npr. Wardove metode. Na začetku določimo voditelje, potem pa na vsakem koraku vsako enoto priredimo voditelju oz. skupini, kateremu je najbližja glede na evklidsko razdaljo. Na vsakem koraku se izračunajo novi voditelji kot povprečja skupin. Postopek se zaključi, ko so novi voditelji enaki starim.

Izberemo tisto razvrstitev, ki ima najmanjšo vrednost Wardove kriterijske funkcije, za katero vemo, da pada z naraščanjem števila skupin. Torej za optimalno število skupin ponavadi vzamemo tisto vrednost, kjer se zgori t.i. “koleno” funkcije. Če to “koleno” ni jasno razvidno, lahko sklepamo, da skupine niso jasno ločene. Postopek običajno večkrat ponovimo, saj za različne začetne voditelje lahko dobimo različne rešitve, torej razvrstitve v skupine.

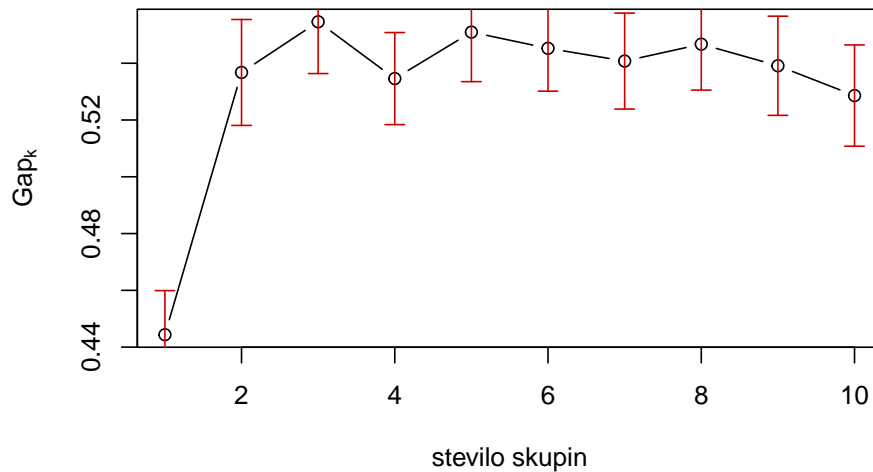


Slika 13: Vrednost Wardove kriterijske funkcije.

Sprememba naklona funkcije izgleda največja pri 2 ali 4 skupinah oziroma je tam “koleno” najbolj razvidno.

### 4.2 GAP statistika

Pri določevanju števila skupin si lahko pomagamo tudi z GAP statistiko, kjer iščemo skupine, ki so podatki bolj homogeni, kot kjer ni skupin. Gre za primerjavo razdalj znotraj skupin z razdaljami na podatkih brez skupin. Izberemo pa tisto najmanjše število skupin  $k$ , kjer je vrednost  $GAP(k)$  statistike vsaj tolikšna kot  $GAP(k+1) - SE(GAP(k+1))$ ;  $SE$  je standardna napaka GAP statistike.

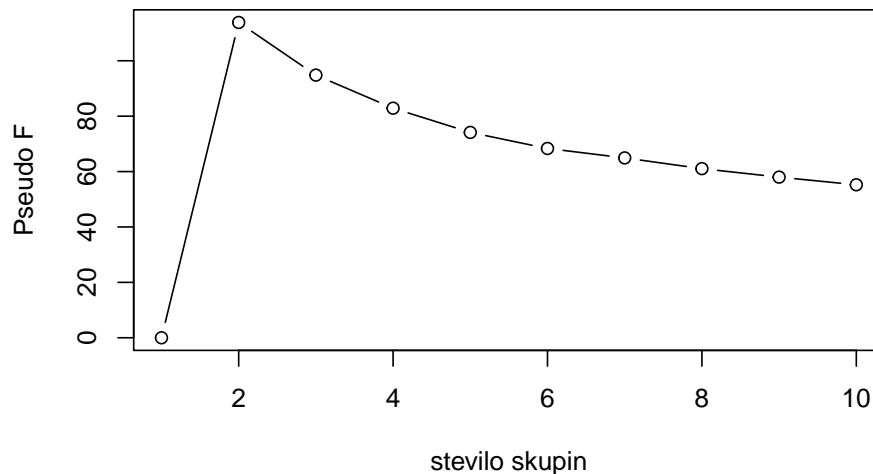


Slika 14: Vrednost GAP statistike.

Na podlagi grafičnega prikaza vrednosti GAP statistike pri različnem številu skupin se odločimo za **4** skupine, saj tam doseže najvišjo točko.

### 4.3 Pseudo F (Calinski - Harabasz indeks)

Uporabimo pa lahko tudi indeks Calinski-Harabasz, ki ocenjuje razmerje med razpršenostjo znotraj skupin in razpršenostjo med skupinami. Uporabljamo ga za oceno primernosti števila skupin v metodi gručenja (angl. *clustering*). Višje vrednosti indeksa Calinski-Harabasz označujejo boljše gručenje, pri čemer optimalno število skupin običajno doseže maksimum tega indeksa.



Slika 15: Vrednost Pseudo F oz. Calinski - Harabasz indeksa.

Tukaj je maksimum dosežen pri **2** skupinah.

Torej, če povzameva celotno analizo, bi, glede na posamezen graf, izbrala

- WSS: sprememba naklona izgleda največja pri **4** skupinah,
- Pseudo F: maksimum doseže pri **2** skupinah,
- gap statistika: najvišjo točko preden začne padati doseže pri **4** skupinah.

Na podlagi zgornjih analiz in ugotovitev pri hierarhičnem razvrščanju, kjer smo se odločali med 2, 3 ali 4 skupinami, bi se tu določili raje za **4** skupine, kot za 2, saj težimo k večjemu številu skupin kot je 2.

#### 4.4 Primerjava vrednosti kriterijske funkcije za Wardovo metodo in K-means

Primerjala sva tudi vrednosti kriterijskih funkcij za Wardovo metodo in metodo K-means, ker sta podobno oziroma delujeta na isti princip. Je pa metoda K-means boljša, ker išče lokalne minimume, za razliko do Wardove, ki deluje hierarhično in vedno poda enak rezultat. Ocenjujeva sva pa po principu, da ima boljša razvrstitev manjšo vrednost karaketristične funkcije. Pomembno pa je tudi to, da so podatki standardizirani, saj drugače med seboj ne bi bilo primerljivo.

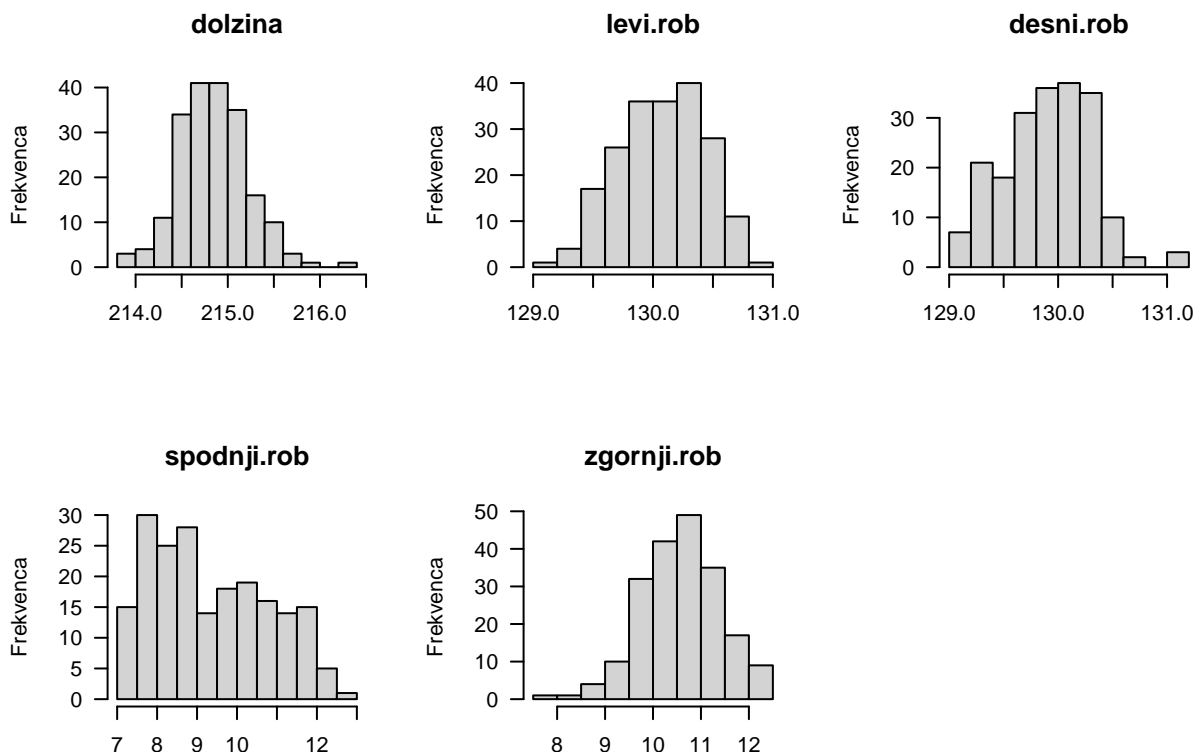
Vidimo da ima v vseh primerih (z izjemo prvega kjer sta enaka) K-means manjšo vrednost, kar si tudi želimo. Primerjavo razvrstitev bomo naredili na številu skupin  $k = 4$ .

Največje elemente imamo na diagonali kontingenčne tabele, tudi te vrednosti niso ekstremno velike (npr. 400). Za izvendiagonalne elemente si želimo, da bi bili čim manjši oziroma zelo blizu 0, kar pa po večini so, ali pa so celo kar enaki 0 (iztopa le ena vrednost - 17).

Poglejmo si še Randow indeks, ki predstavlja delež parov enot, ki so si v obeh razbitjih usklajeni - v obeh razbitjih v isti skupini ali pa v obeh razbitjih v različnih skupinah. Pogledala sva si popravljen Randow indeks, zaradi boljše primerljivosti. Enak je 0.644, kar je blizu 0,5, torej gre skoraj za neko srednjo podobnost, sicer pa večji kot je, boljše je - vrednost 1 pomeni identični razbitji, vrednost 0 pa, da sta si razbitji tako podobni po slučaju.

## 5 Razvrščanje na podlagi modelov

Tukaj predpostavimo, da so podatki generirani iz multivariatnih normalnih porazdelitev z različnimi parametri oziroma komponentami; vsaka skupina ima svojo multivariatno normalno porazdelitev. Skupina je večja po volumnu, če ima večjo variabilnost, omejimo pa se z domnevami oziroma predhodnim znanjem, kakšne naj bi ti skupine bile. Zato si pogledimo porazdelitve spremenljivk ne glede na tip bankovca.



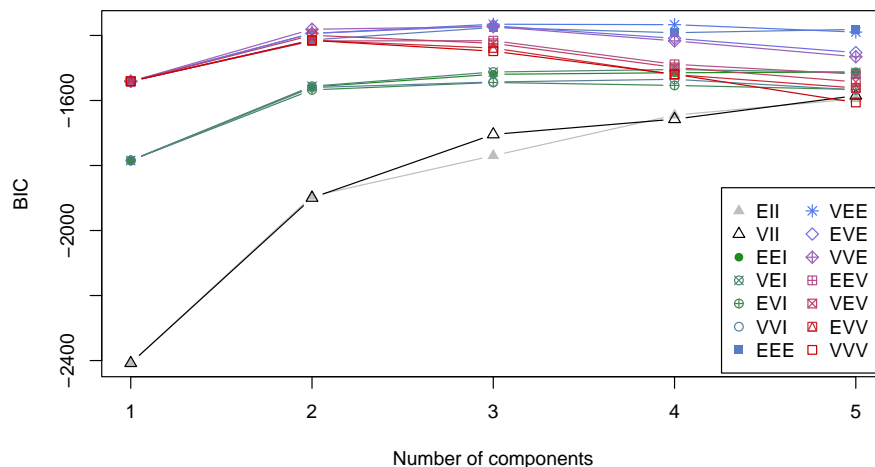
Slika 16: Porazdelitve spremenljivk.

Spremenljivka **spodnji.rob** in tudi **desni.rob** nista porazdeljeni po normalni porazdelitvi, zato ne moremo trditi, da je zadoščen ta pogoj. Ostale so porazdeljene po normalni, nekatere asimetrične v desno (npr. spremenljivka **dolzina**) in nekatere v levo (npr. spremenljivka **levi.rob**).

Tukaj ocenimo število skupin in parametre za vsako skupino ter kateri skupini posamezna enota pripada. V najinem primeru, kjer je predpostavka o multivariatni normalni porazdelitvi kršena, se simulacija ne izkaže za optimalno. Razvrstitev se dela na originalnih podatkih oz. nestandardiziranih podatkih, ker s tem omogočimo različno velikost skupin.

### 5.1 BIC (Bayes Information Criterion) kriterij

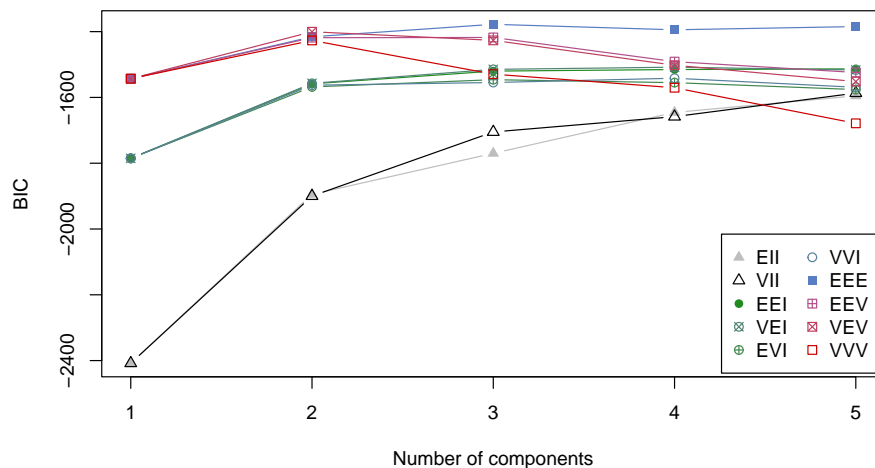
Naredimo torej razvrstitev na originalnih, nestandardiziranih podatkih, kjer funkcija sama izbere naprimernejši model.



Slika 17: BIC kriterij za originalne podatke.

Na podlagi BIC kriterija (Bayesian Information Criterion), ki zavzame vrednost -592.64 izberemo model VEE s tremi skupinami, kar pomeni, da gre za elipsoidne(angl. *ellipsoidal*) skupine, ki so različno velike, različnih oblik in enako usmerjene.

Kriterij temelji na "Bayesovski" statistiki, zato lahko določimo tudi apriorne verjetnosti(torej neko naše predhodno znanje oziroma prepričanja).



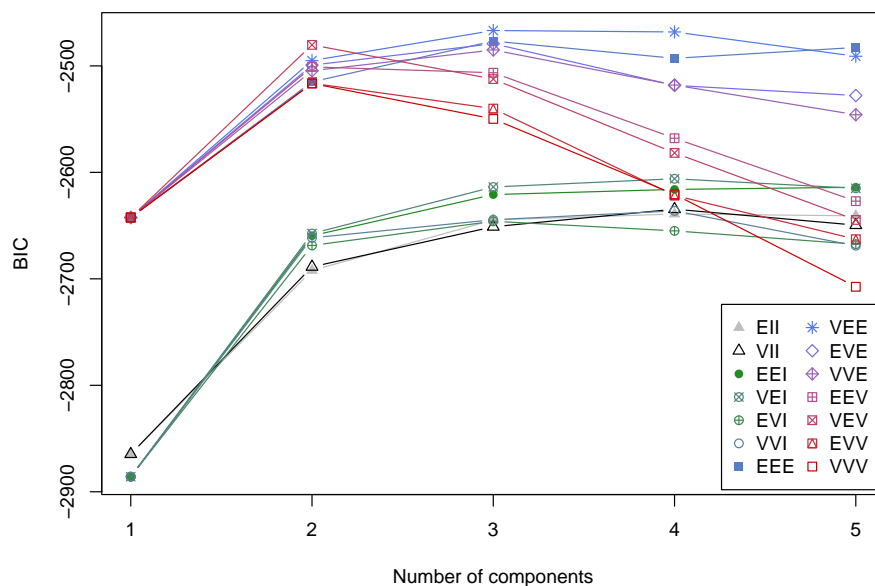
Slika 18: BIC kriterij (priorControl) za originalne podatke.

Na podlagi BIC kriterija z uporabljenim argumentom o apriornih verjetnostih se odločimo za model EEE s tremi skupinami, kar pomeni, da gre za različno velike skupine, različnih oblik in enake usmerjenosti.

## 5.2 BIC kriterij na standariziranih podatkih

Poglejmo si še, iz radovednosti, kako je z oceno modela na standariziranih podatkih, ampak vrednosti BIC kriterija niso primerljive med standariziranimi in nestandariziranimi podatki.

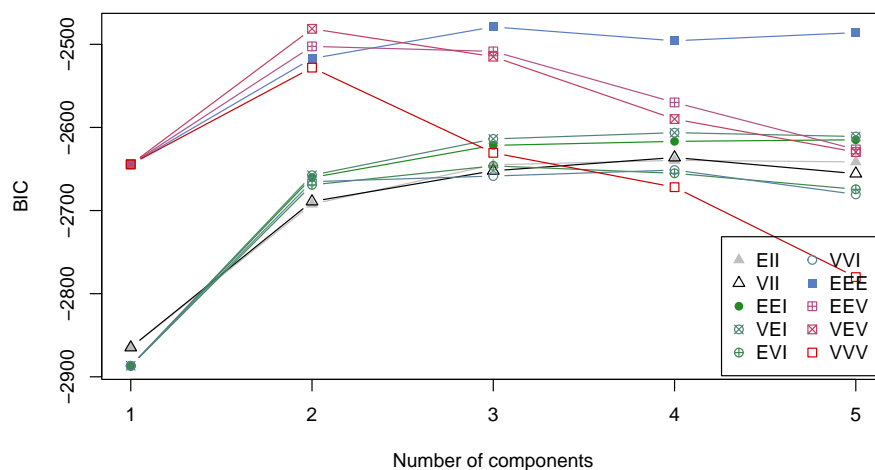




Slika 19: BIC kriterij za standardizirane podatke.

Na podlagi BIC kriterija (Bayesian Information Criterion), ki zavzame vrednost -1143.31 ponovno izberemo model VVE s tremi skupinami.

Tudi tukaj lahko primerjamo z vključitvijo apriornih verjetnosti.



Slika 20: BIC kriterij (priorControl) za standardizirane podatke.

Tudi tukaj se na podlagi BIC kriterija z uporabljenim argumentom o apriornih verjetnostih se odločimo za model EEE s tremi skupinami, kar pomeni, da gre za različno velike skupine, različnih oblik in enake usmerjenosti.

Na podlagi vseh štirih kriterijev se zaradi enostavnosti odločimo za model VEE s tremi skupinami - torej različno velike skupine, različnih oblik in enakih usmerjenosti.

## 6 Najboljša razvrstitev in predstavitev skupin

Tukaj naju pa zanima kako podobne so si naše razvrstitve, ki sva jih v prejšnjih poglavjih izbrala na podlagi različnih modelov.