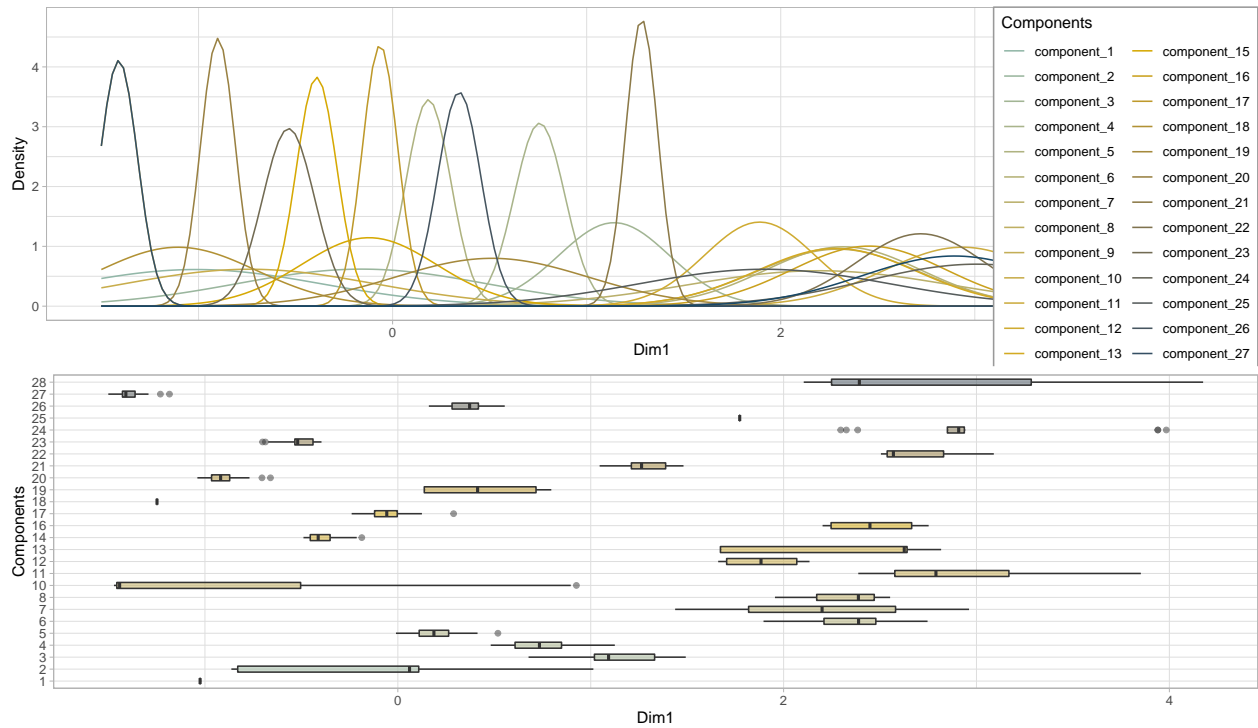


Modell-alapú klaszteranalízis (MKA) R-ben

1. Végezz MKA-t a PTELJ, Pboldog, Pmagány input változókkal, outlier kiszűréssel! Melyik megoldás tűnik a legjobbnak a BIC-grafikon alapján?

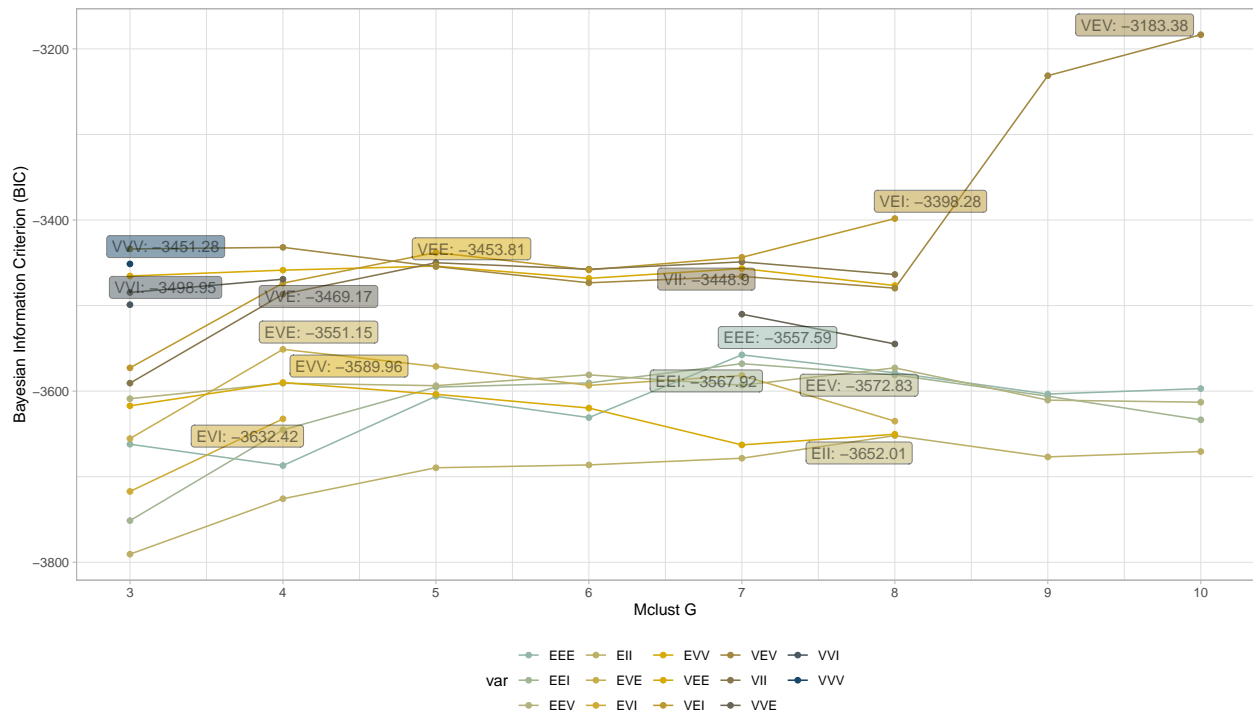
Figure 1 is a line graph showing the Bayesian Information Criterion (BIC) for various clustering methods (EEE, EEI, EEV, EEE, EEI, EEV, EEE, EEI, EEV) across different numbers of clusters (Mclust G). The Y-axis represents BIC, ranging from -4000 to -2000. The X-axis represents Mclust G, ranging from 3 to 30. The graph displays 12 lines, each representing a different clustering method. The lines generally show a downward trend as the number of clusters increases, indicating that the BIC decreases as the model becomes more complex. The legend at the bottom identifies the methods: EEE (green), EEI (yellow), EEV (orange), EEE (blue), EEI (purple), EEV (brown), EEE (pink), EEI (grey), EEV (olive), EEE (light blue), EEI (light yellow), and EEV (light orange).

Az alábbi ábrán az egyes keverék-komponensek sűrűségét, illetve box-whiskers eloszlásait látjuk.



2. ábra MKA sűrűség és eloszlási eredmények G=28 között

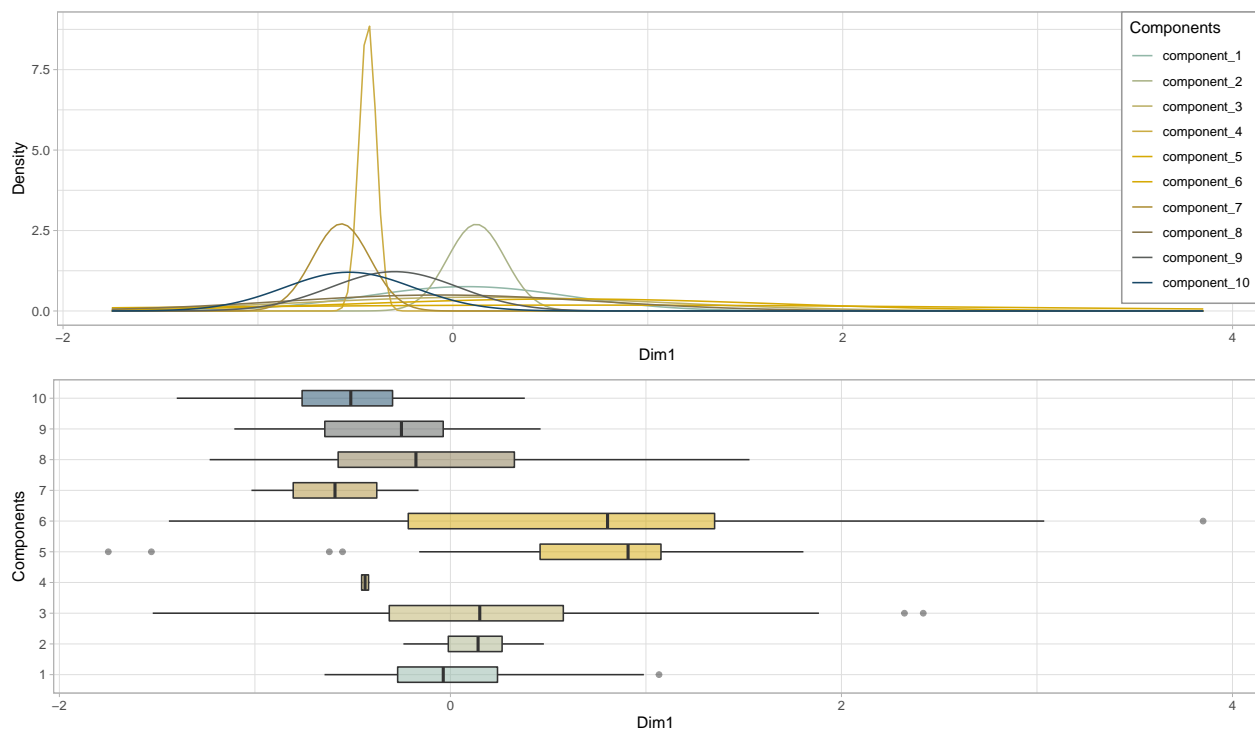
Mindezekon túl azonban, az értelmezhetőség végett az elemzést a G=3 és G=10 tartományra is megismételttem.



3. ábra MKA BIC eredmények G=3 és G=10 között

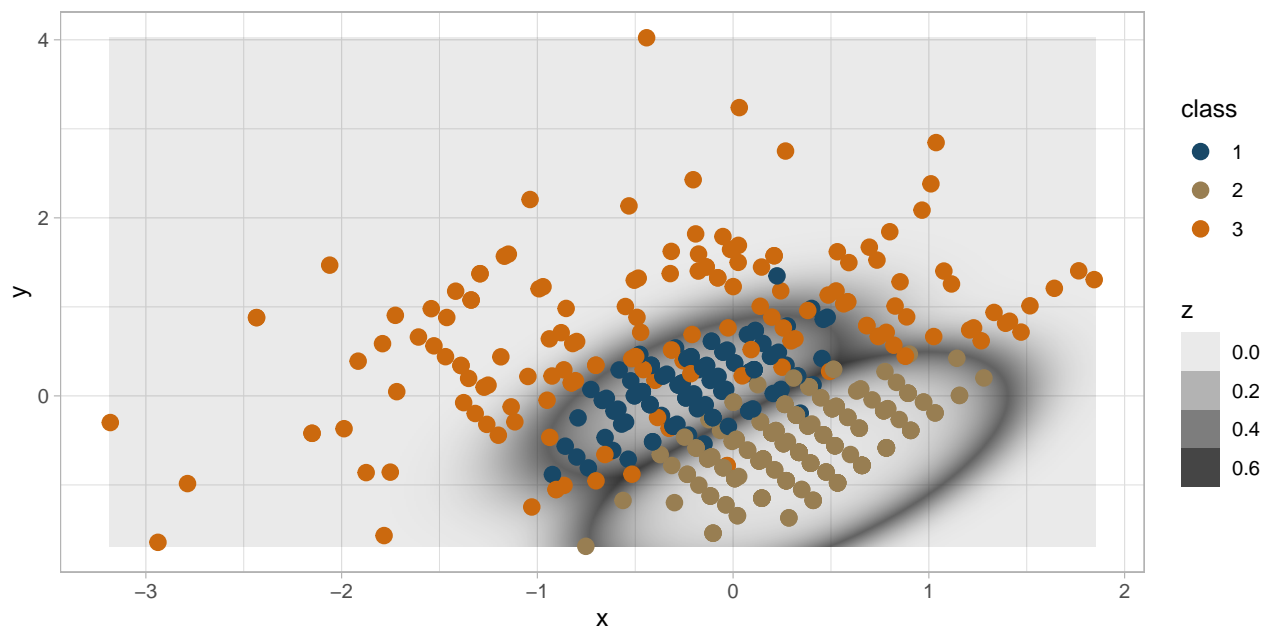
Ezen elemzés nyomán azt látjuk, a VEV eloszlástípus bizonyul a legjobbnak, mely a fenti, kiterjesztett elemzésnél is a legjobbak között volt, így az képződött struktúra értelmezhetőségének reményében a G=10, VEV struktúrát fogadom el.

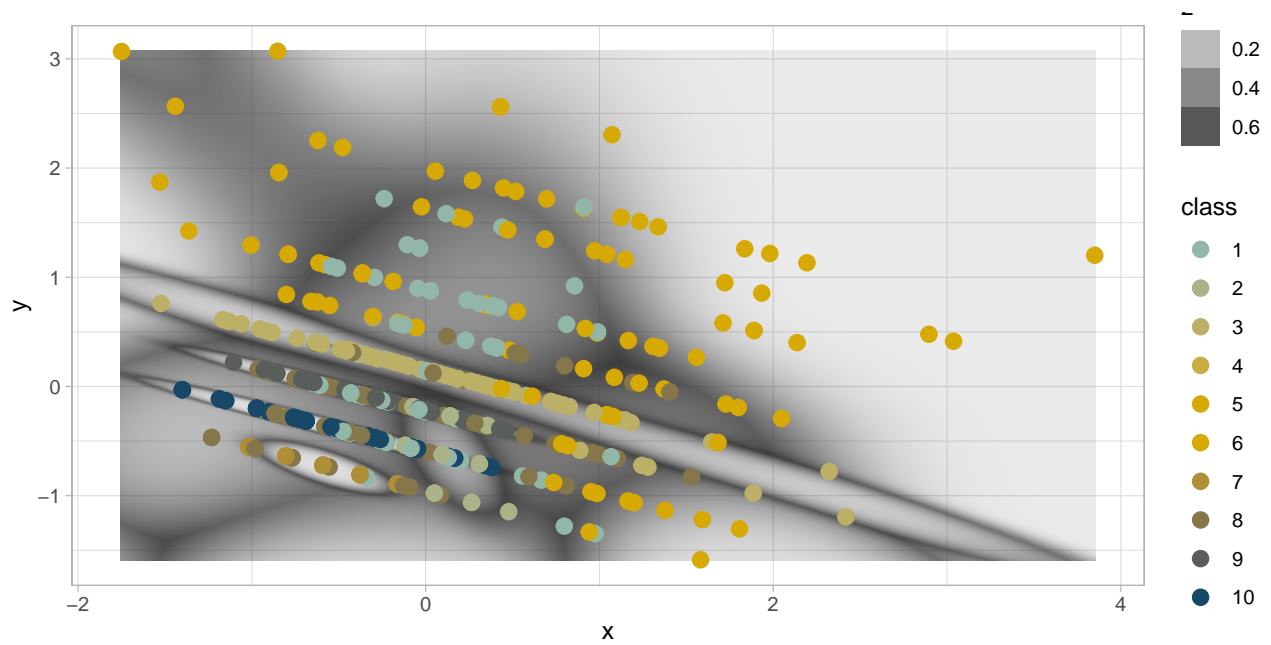
Ha az keverékkomponensek sűrűségét tekintjük (alábbi, 4. ábra)



4. ábra MKA sűrűség és eloszlási eredmények G=10 között

A keverékkomponensek sűrűsödés vizsgálata három jobban elkülönülő struktúrát jelez, így mindez felveti hipotézisként, vajon mennyire értelmezhetőbb egy G=3 struktúra az információvesztés ellenére a G=10-zel szemben.

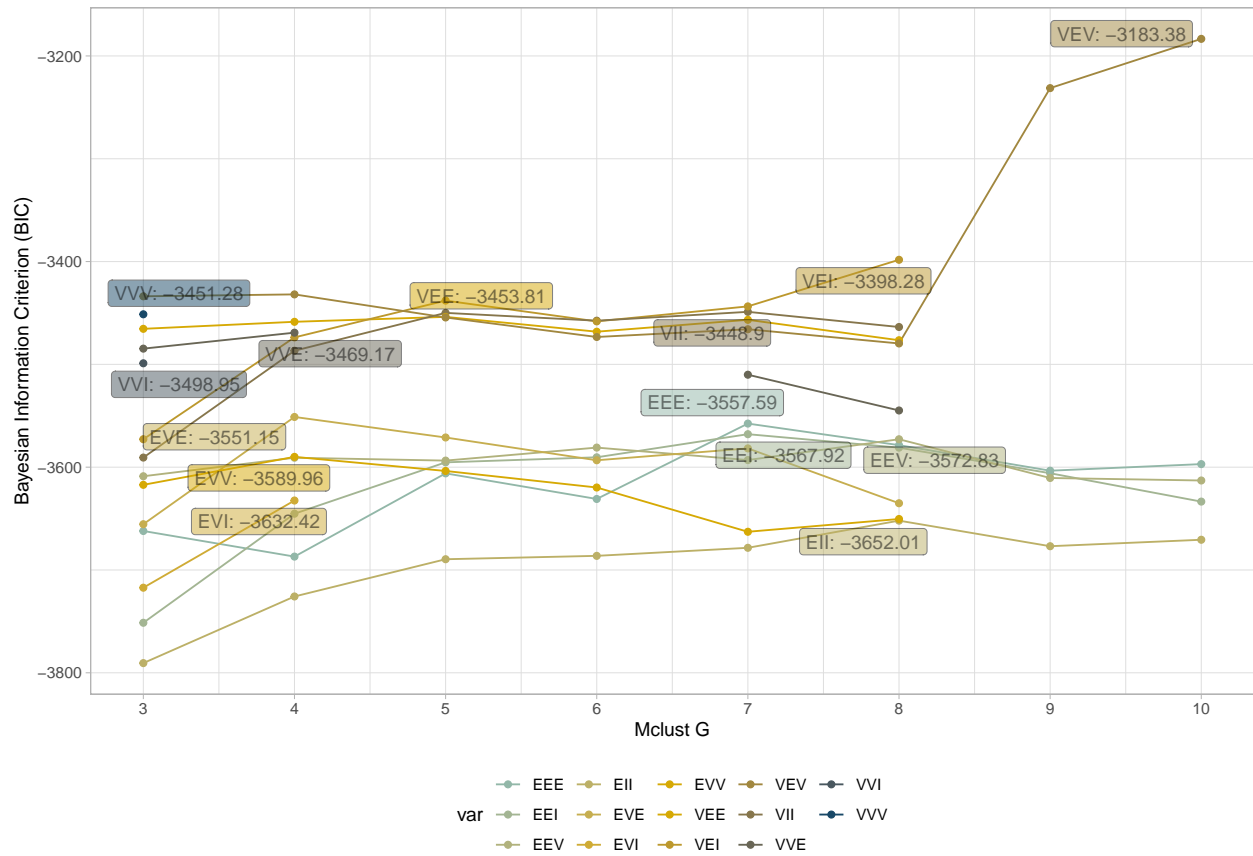




5.ábra A $G=3$ és a $G=10$ struktúrák összevetése az adatok első két főkomponense mentén képzett két dimenziós síkban; a klaszterhatárok bizonytalanságát a háttér árnyalata (z -paraméter) jelzi.

Noha értelmezhetőbb struktúrát kapunk $G=3$ értékkel, ez jelentősen rosszabb BIC struktúra mint a $G=10$.

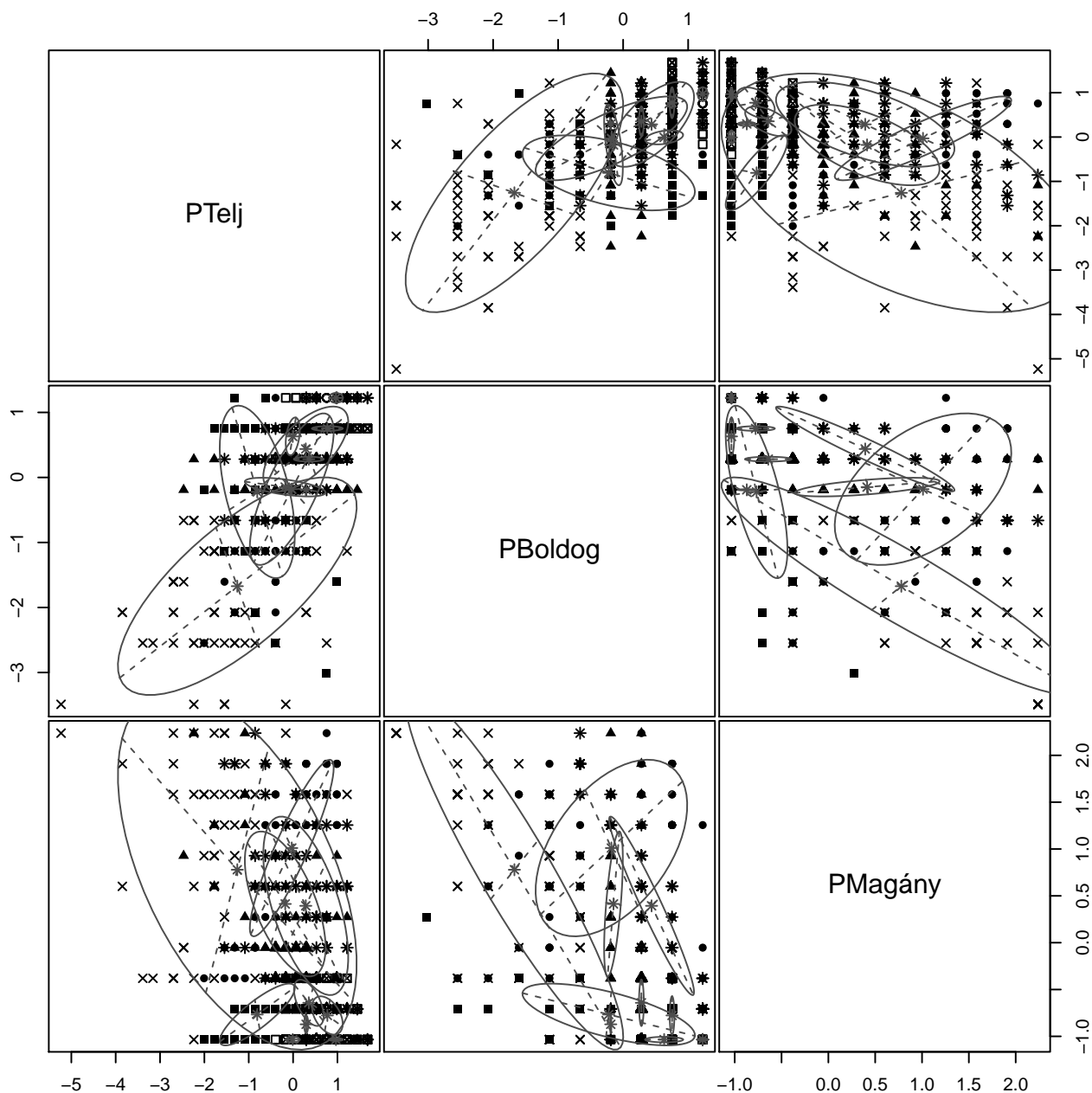
2. Készítsd el az 1. feladat BIC-grafikonját $k = 3$ és 10 között!



6. ábra MKA sűrűség és eloszlási eredmények $G=10$ között

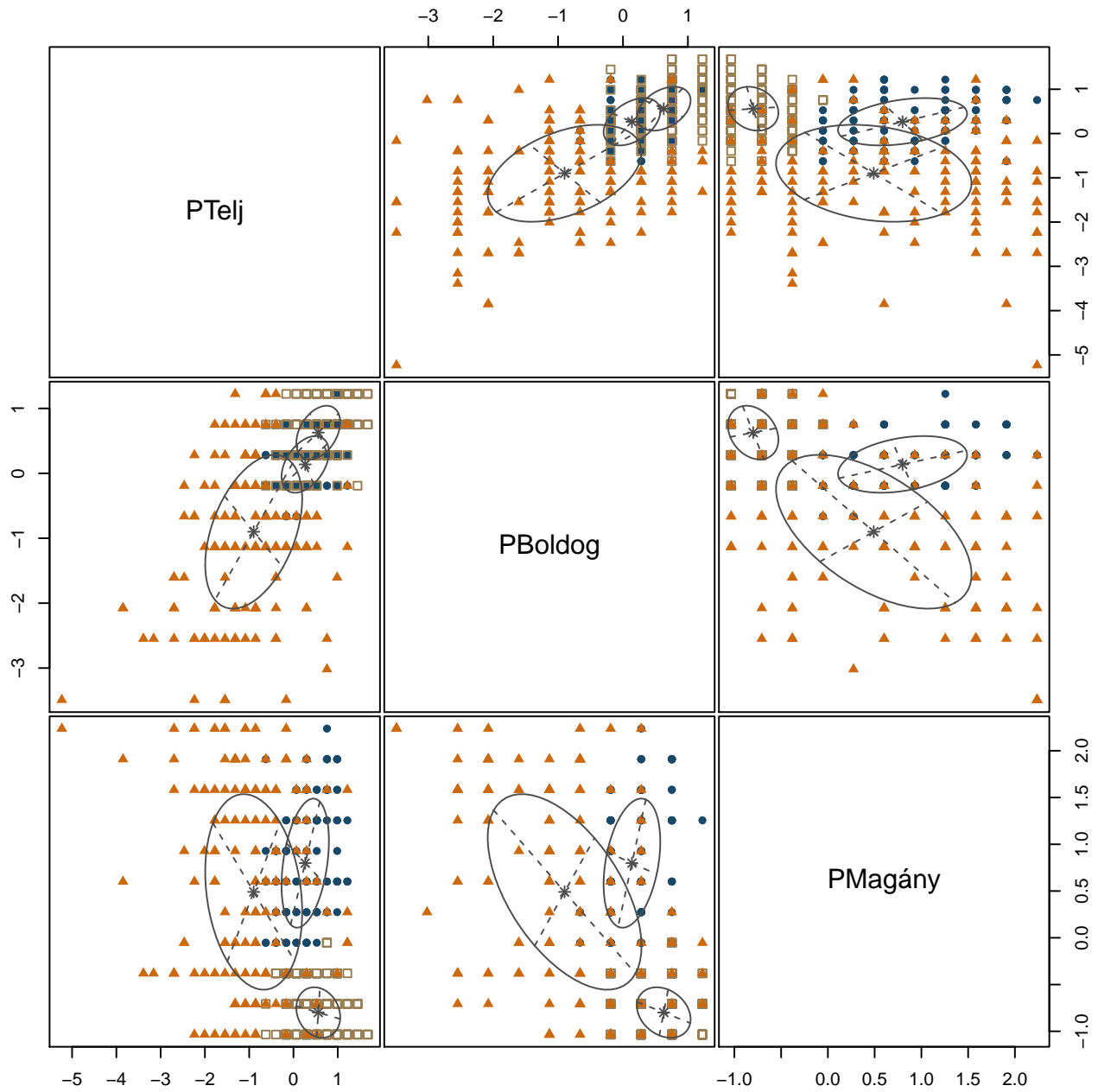
BIC ábráimon a szövegdozok az adott keveréktípus maximális értékénél állnak. Ennek értelmében $G=10$ -nél a legjobbnak tűnő eloszlástípus a VEV, melynek BIC értéke -3183.38.

3. Készítsd el az 1. feladat legjobb BIC megoldásának classification ábráját!



7. ábra A G=10 megoldás klasszifikációs ábrája

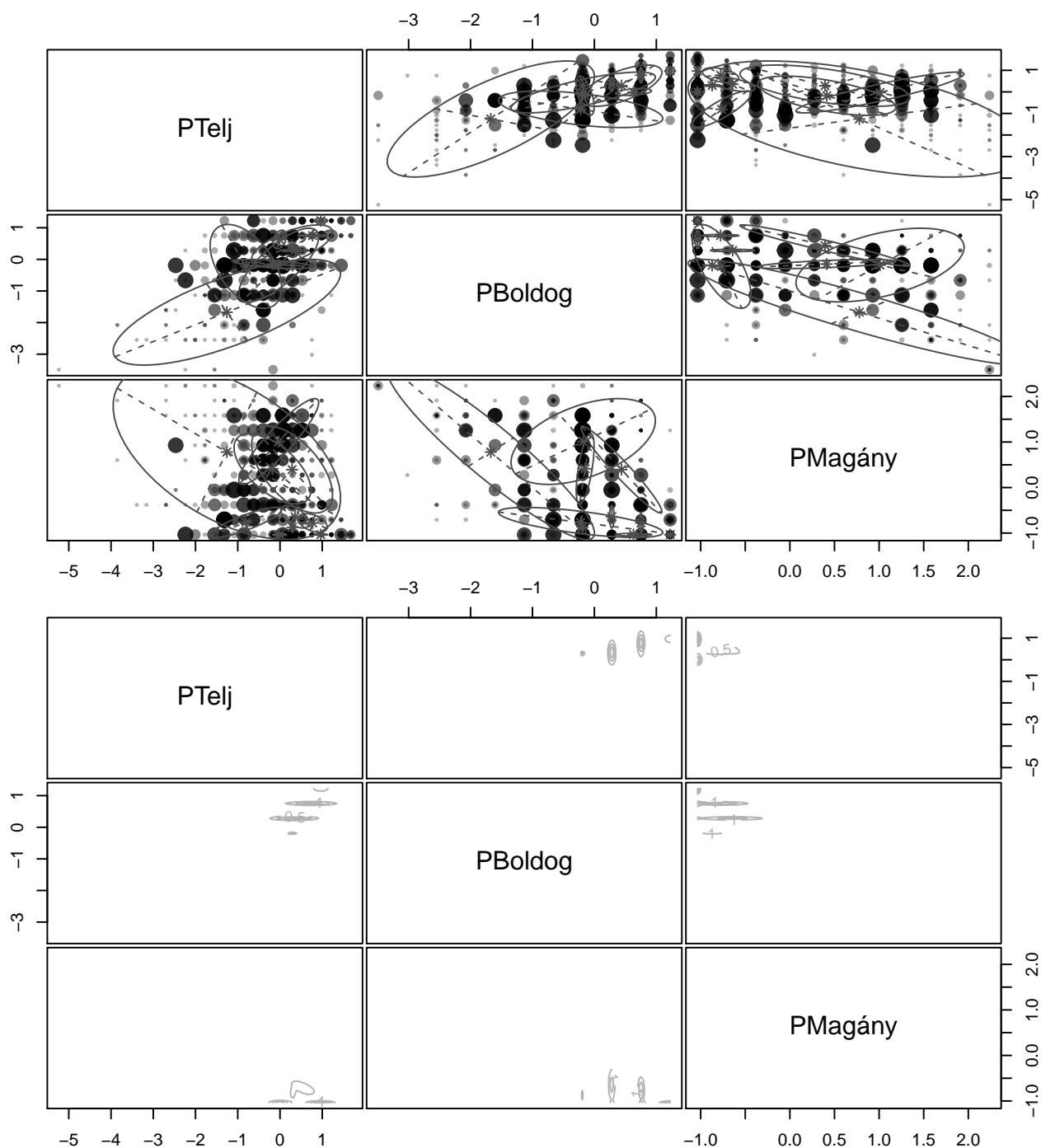
Jól látható három nagyobb klaszter elkülönülése, és több kisebb, részben átfedő struktúra is. Ez ismét felveti a kérdést, mennyiben értelmezhetőbb egy G=3 struktúra. Különösen a Boldogság és Magányosság szeleteiben láthatóak néhány személyt magukban foglaló apró klaszterek.



8. ábra A G=3 megoldás klasszifikációs ábrája

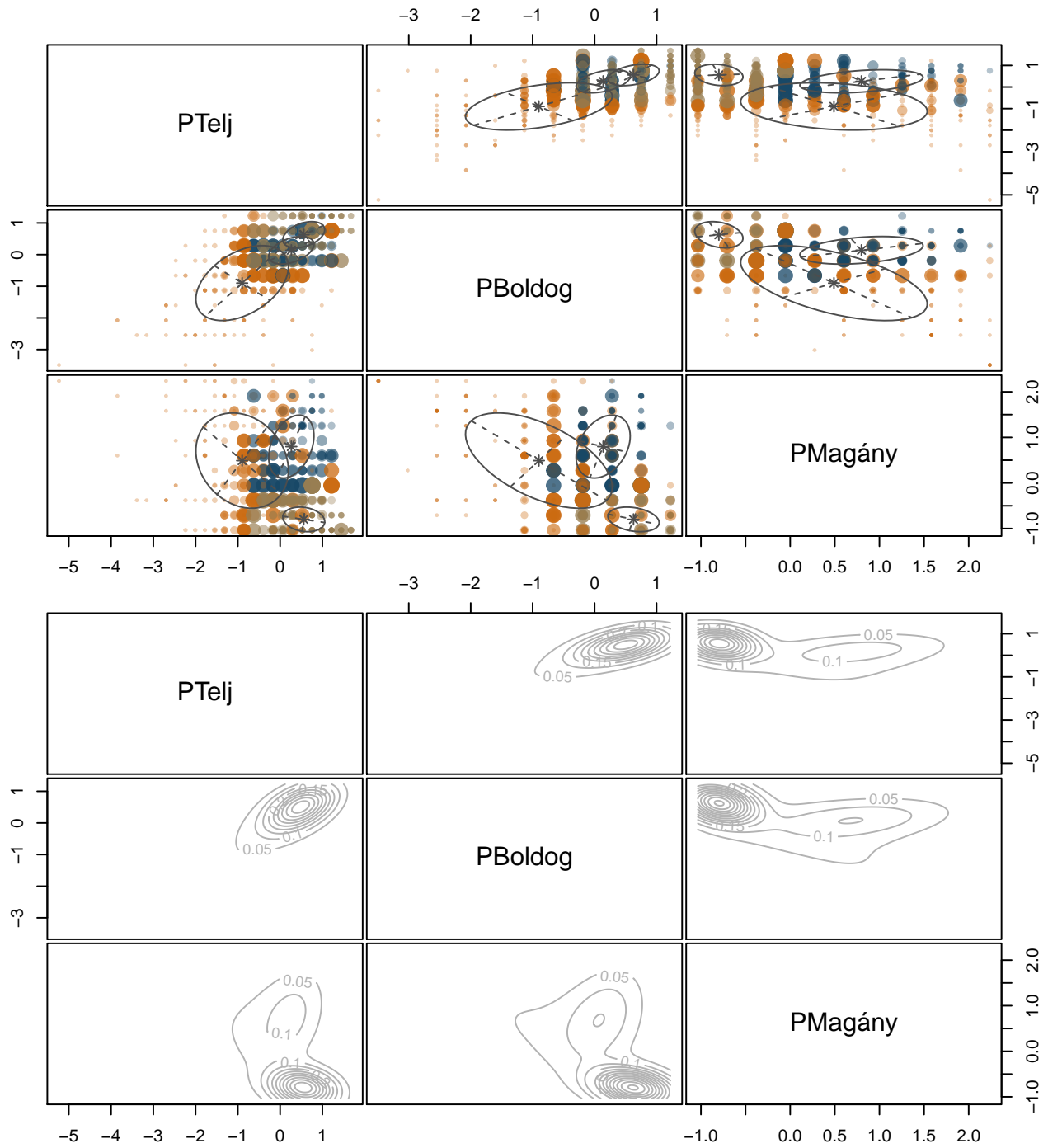
A BIC-értékben való csökkenés ellenére egy jobban értelmezhető, kevésbé átfedő, kevésbé redundáns megoldást kapunk.

4. Készítsd el az 1. feladat legjobb BIC megoldásának uncertainty és density ábráját!



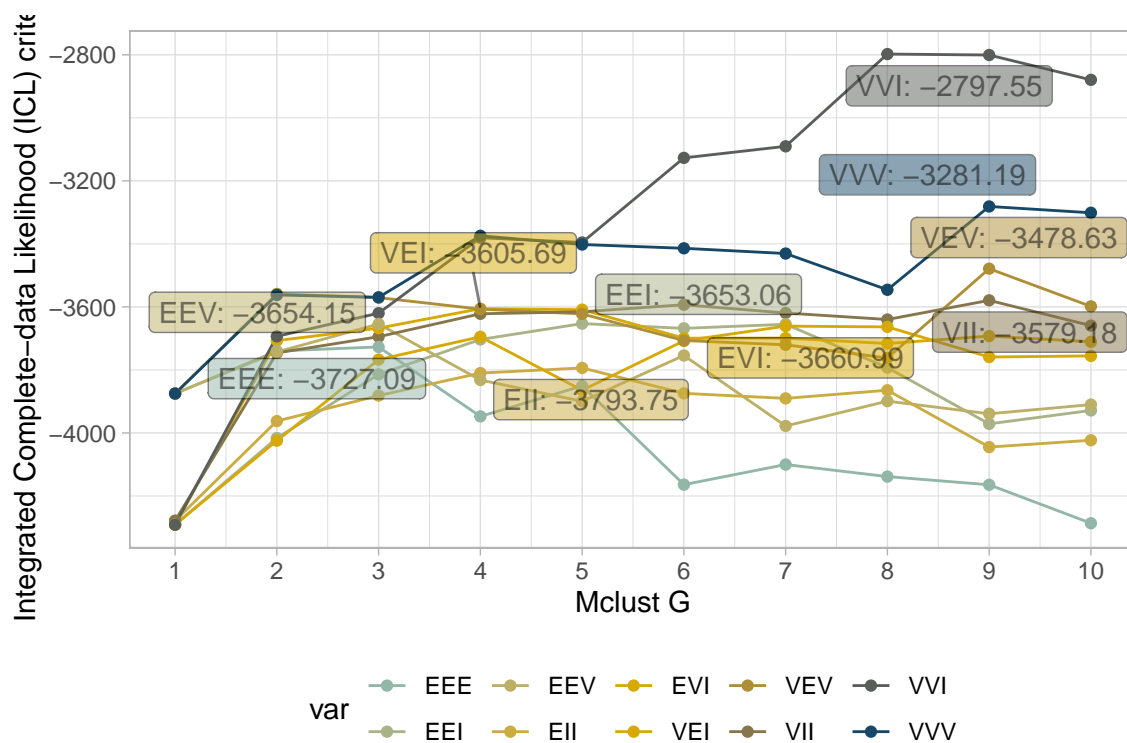
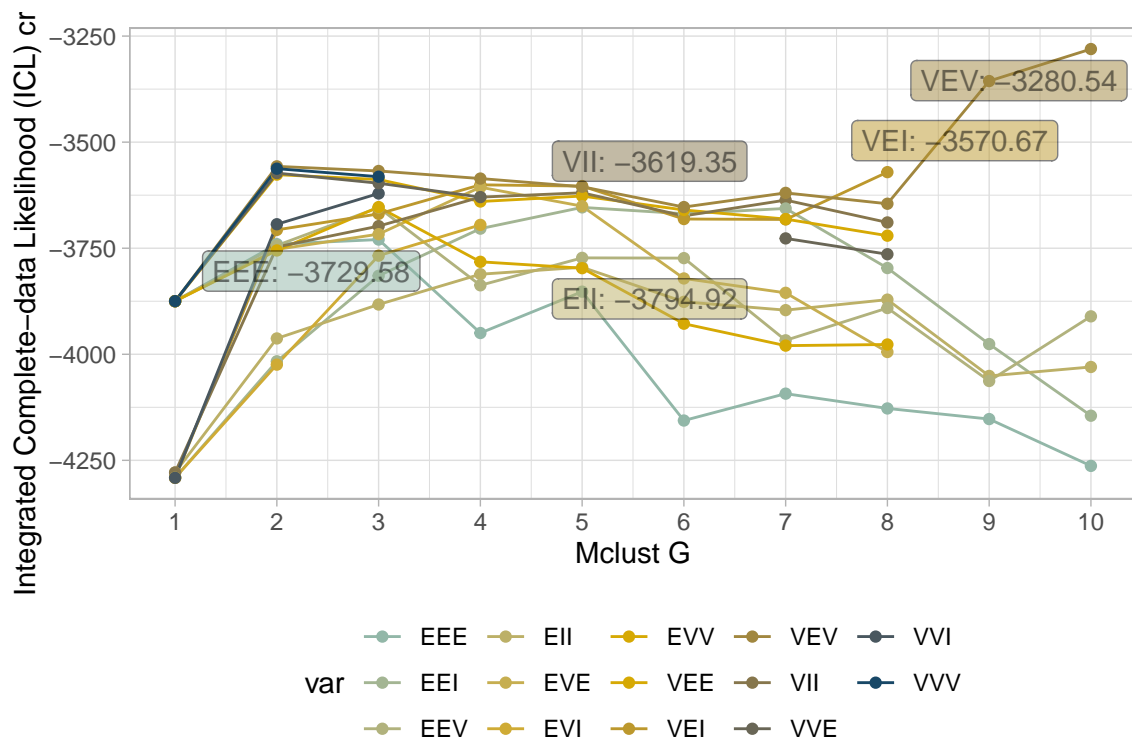
9. ábra A G=10 megoldás “uncertainty” és sűrűsödés ábrája

G=10 esetében a klaszterek átfednek, több esetben is a klaszterbe tartozás bizonytalansága emelkedett.



10. ábra A G=3 megoldás “uncertainty” és sűrűsödés ábrája

5. Készítsd el az 1. feladat ICL-grafikonját $k = 1$ és 9 között! Ugyanaz a modell tűnik a legjobbnak, mint a BIC-grafikon alapján?



11. **ábra** A $G=1$ és $G=10$ közti megoldások ICL ábrái (a felső a prior opció nélkül, az alsó pedig ezzel kiegészített lefutás)

A *priorControl()* hangolási opció nélkül, a felső ábra szerint hasonlóképpen a VEV keveréktípus a legjobb az ICL információs kritérium szeriunt is, akár $G=9$, akár $G=10$ opciót tekintjük is.

A finomhangolással együtt azonban (alsó ábra) a VVI eloszlástípus mutat kedvezőbb illeszkedést. Ezt úgy értelmezem, hogy méretükben, tengelyeik hosszában eltérő, de a főtengelyekkel és egymás tengelyirányultságában egyező eloszlásokat modellez a legjobb ICL-lel leírható modell.

A VVI típus G=8 struktúra esetén tetőzik, így ezt fogadom el a legjobb megoldásnak, mely az alacsonyabb komponensszám miatt az értelmezhetőségben is kedvezőbb.

6. Mentsd el a legjobb BIC-megoldást, tedd át ROPstatba és számítsd ki a Validálás modullal a főbb QC mutatókat! Hasonlítsd össze a kapott értékeket a 8. óra 1. feladatában kapott QC-értékekkel!

Az alábbi táblázatban mutatom be a validálás eredményét.

A

EESS%	Pontbisz	XBmod	Sil.eh.	HCátlag	CLdelta	GDI24	HCmin-HCmax	G/K	Type
51.52	0.362	0.513	0.631	0.974	0.754	0.447	0.33-2.22	3	MKA
54.99	0.189	-0.795	0.181	0.915	0.555	0.039	0.09-2.37	9	MKA
58.49	0.176	-0.776	0.216	0.842	0.544	0.041	0.02-2.25	10	MKA
78.72	0.355	0.456	0.598	0.434	0.893	0.381	0.24-1.03	8	KKA
76.83	0.370	0.534	0.610	0.471	0.895	0.484	0.24-1.03	7	KKA
74.70	0.369	0.491	0.614	0.513	0.888	0.345	0.24-1.44	6	KKA
80.26	0.351	0.496	0.590	0.404	0.898	0.381	0.24-1.03	9	KKA
81.48	0.342	0.473	0.592	0.379	0.896	0.342	0.18-1.03	10	KKA
79.50	0.233	-0.877	0.136	0.438	0.789	0.011	0.00-2.21	28	MKA
1. táb	lázat** Ade	kvációs m	utatók öss	zevetése					