

# Shangai Data: Case study

## Deliverable 1

## Outline

En aquest estudi tractarem la mostra “Shangai Data”, una mostra de 489 individus que han valorat diferents aspectes de les instal·lacions d’un allotjament. Aquests aspectes es poden valorar mitjançant diferents variables quantitatives o qualitatives.

L’estudi es divideix en dues parts: un PCA i un MCA, tots dos enriquits per clustering. Tots els càlculs i els gràfics que es presentaran a continuació s’han realitzat amb l’ajuda del programari R.

## Aplicació del PCA enriquit per clustering

### Inèrcia dels eixos. Quins factors retenim?

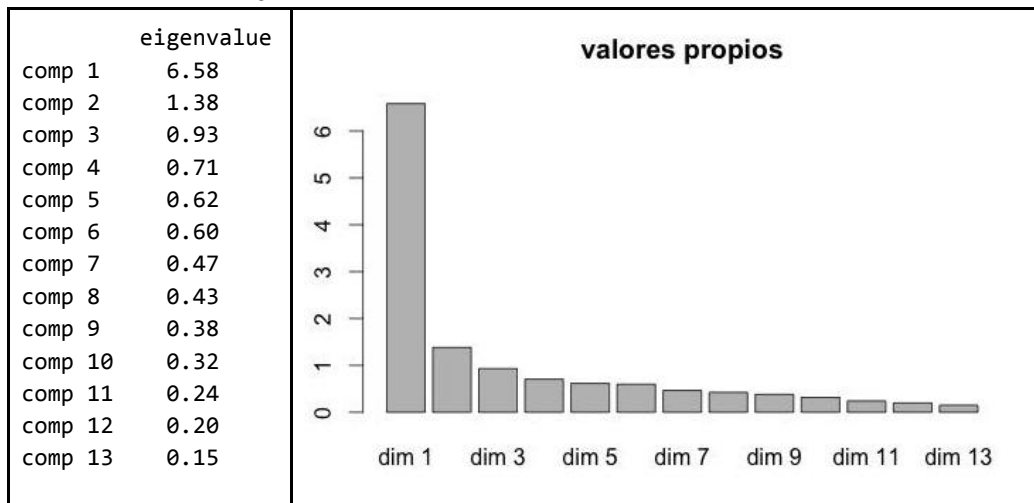


Figura 1. Eigenvalues: valors propis

En un PCA estandarditzat la inèrcia mitjana dels valors propis és 1, per tant només hem de mantenir aquells valors propis superiors a 1. En conseqüència, si observem la Figura 1, podem deduir que centrarem aquest estudi en les dues primeres dimensions. Ara bé, la dimensió 3 té una valor molt proper a 1, les seves dades podrien arribar a ser significatives i aportar informació útil. No descartem utilitzar aquesta tercera dimensió en algun moment.

### Interpretació dels factors i dels eixos

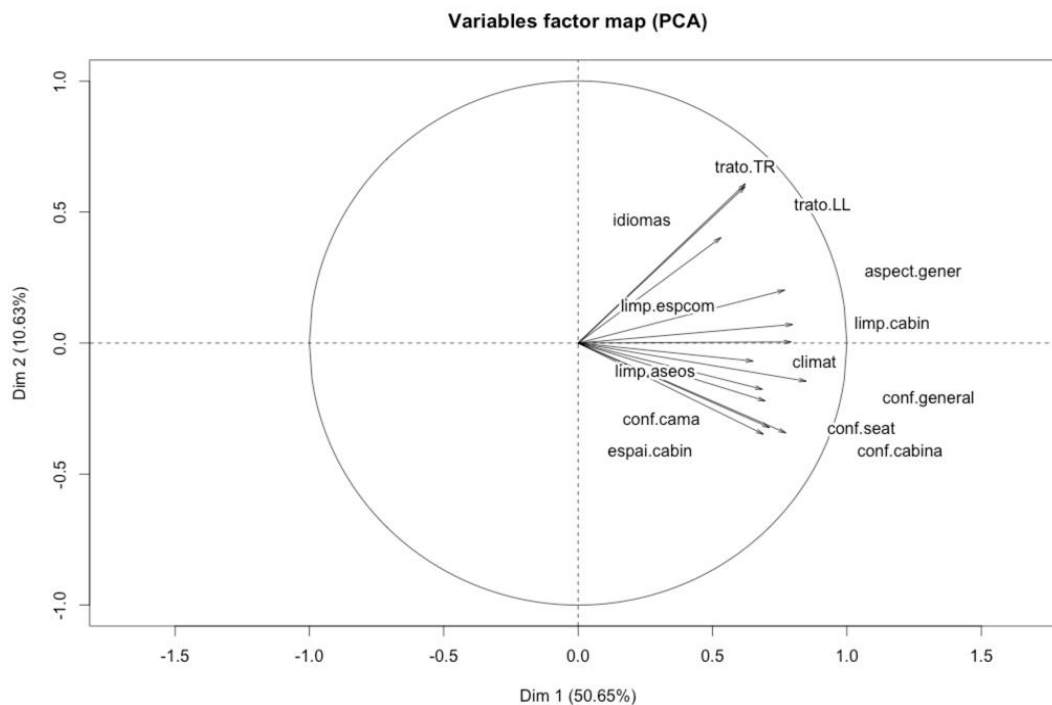


Figura 2. PCA. Mapa factorial de les variables en funció de les dimensions 1 i 2.

Les dimensions són tendències de puntuacions que han fet els diferents individus.

A partir de les figures 2 i 3 podem afirmar que la dimensió 1 en general té bones puntuacions, sobretot respecte el confort i la neteja. En canvi, a les altres dues dimensions podem observar puntuacions negatives entre algunes de les variables. A la dimensió 2 es puntua positivament el tracte però negativament l'espai i el confort a cabina. A la dimensió 3 es puntua positivament el confort i l'espai a cabina i negativament la neteja.

A més, a la Figura 2 podem observar que les variables conf.seat i conf.cabina presenten una correlació alta. Les variables trato.TR i idiomas també.

\$Dim.1\$quanti	correlation	\$Dim.2\$quanti	correlation	\$Dim.3\$quanti	correlation
conf.general	0.8481471	trato.TR	0.6076557	espai.cabin	0.4112225
limp.espcom	0.7985376	trato.LL	0.5966692	conf.cabina	0.2820723
limp.cabin	0.7920013	idiomas	0.4026411	trato.LL	0.2484574
conf.cabina	0.7726858	aspect.gener	0.2021238	trato.TR	0.2363126
aspect.gener	0.7683819	conf.general	-0.1458578	conf.cama	0.1639467
conf.seat	0.7114754	limp.aseos	-0.1766384	climat	0.1075416
conf.cama	0.6958620	conf.cama	-0.2208818	aspect.gener	-0.1481324
espai.cabin	0.6888139	conf.seat	-0.3225243	limp.aseos	-0.3400374
limp.aseos	0.6856033	conf.cabina	-0.3429921	limp.cabin	-0.4291976
climat	0.6511372	espai.cabin	-0.3477319	limp.espcom	-0.4473558
trato.TR	0.6219750				
trato.LL	0.6204638				
idiomas	0.5312163				

Figura 3. Taula amb els valors de correlació de les dimensions 1, 2 i 3.

Tal i com hem havíem dit a l'inici d'aquest estudi, la dimensió 3 ens pot aportar informació rellevant. En aquest cas en relació a les variables de neteja. A la Figura 6 podem observar que la dimensió 3 té una tendència cap a la insatisfacció respecte la neteja en general.

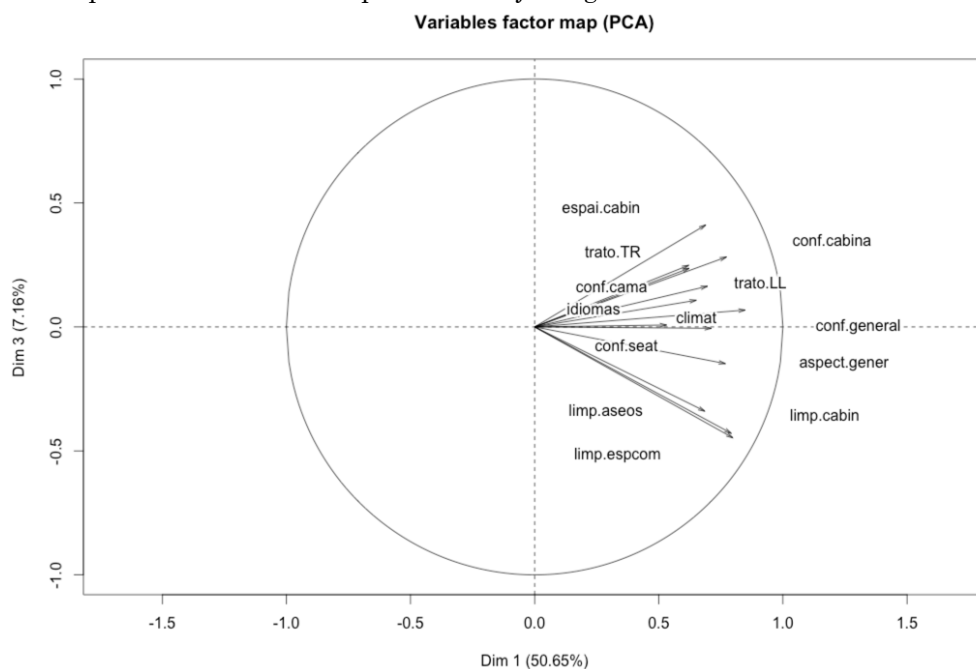
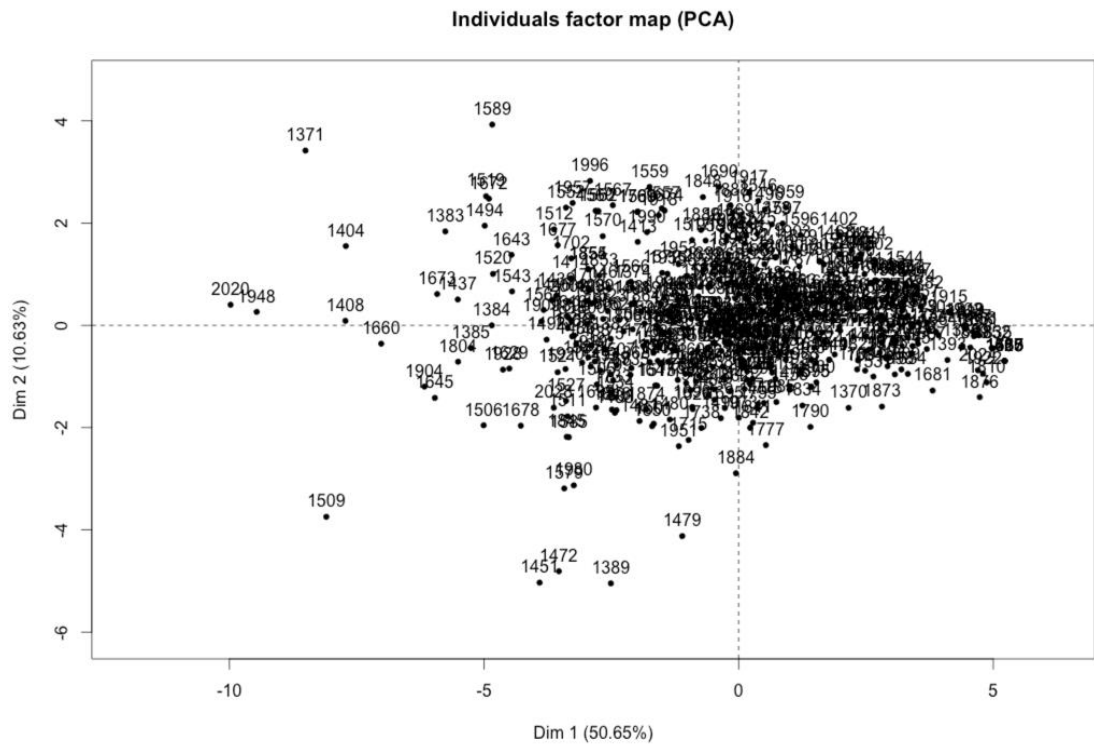


Figura 4. PCA. Mapa factorial de les variables a les dimensions 1 i 3.

Punts de vista dels individus

A partir de la Figura 5 podem observar que la gran majoria d'individus de la mostra estan més o menys centrats amb una lleugera tendència cap a satisfets. Els que estan concentrats a la part dreta del gràfic són els més satisfets en tots els aspectes. En canvi, els que estan dispersos per la part esquerra del gràfic són els menys satisfets, moltes vegades per aspectes diferents, és a dir, que han puntuat malament determinats aspectes.



## Distància entre els nivells de satisfacció

La figura 7 ens mostra que les diferents categories de la variable satisfacció formen una línia gairebé recta seguint l'eix horitzontal. A més, estan ordenades de menys a més satisfeta. Per tant, podem deduir que les respostes han estat coherents.

Podem observar una gran diferència entre la variable “nada” i la variable “poco”. Aquest fet ens indica que són variables molt diferenciades i que els individus que es troben a “nada” realment estan molt insatisfets. Ara bé, per contra trobem que les categories “poco” i “regular” estan una mica més juntes. Si estiguessin una mica més unides podríem arribar a dir que descriuen el mateix nivell satisfacció tot i ser categories diferents.

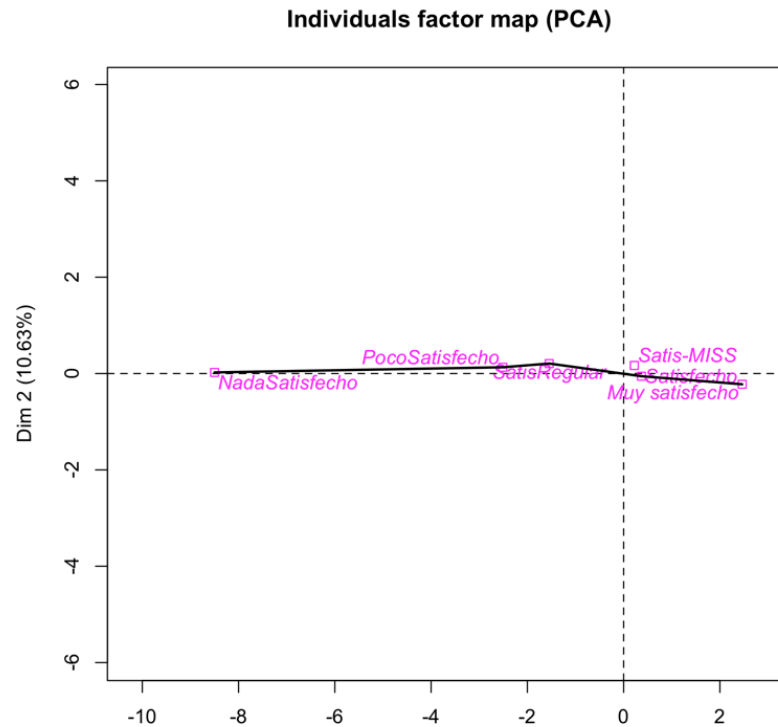


Figura 7. PCA. Mapa factorial de les categories de satisfacció a les dimensions 1 i 2.

## Síntesi del PCA a través del clustering

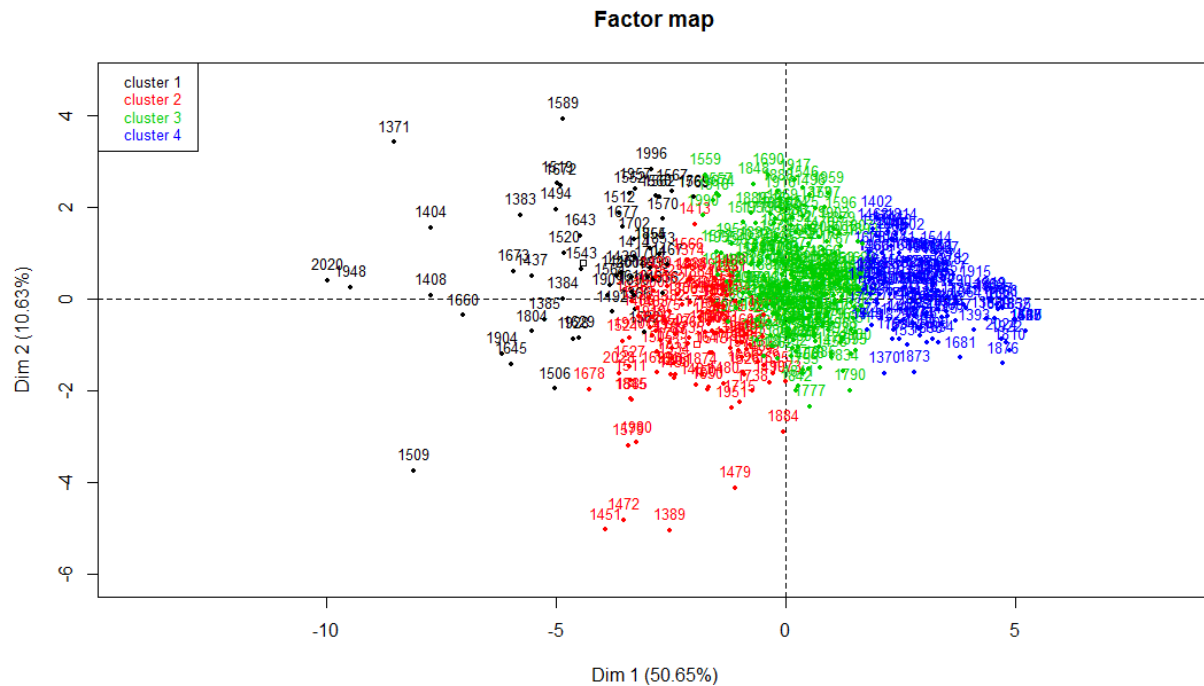


Figura 8. PCA. Mapa factorial dels individus en funció de les dimensions 1 i 2, classificats per clustering.

Tenim 4 clústers amb el següent nombre de components: clúster 1 (97), clúster 2 (141), clúster 3 (121), clúster 4 (130).

Cluster 1		
	Mean in category	Overall mean
trato.LL	1.998922	8.207039
conf.seat	4.158652	6.194617
trato.TR	6.679624	8.307851
conf.cama	4.882927	6.821978
conf.cabina	4.597938	6.628337
conf.general	4.896907	7.018443
mean.SAT	5.126571	7.166966

Al clúster 1 trobem el grup dels insatsifets (categories nada i poc). Ho han puntuat tot per sota de la mitjana. Estan especialment insatsifets amb el tracte i el confort.

Cluster 2		
	Mean in category	Overall mean
espai.cabin	5.801418	5.465164
mean.SAT	6.938419	7.166966

El clúster 2 també està format majoritàriament per gent insatsifeta però amb puntuacions no tant per sota de la mitjana com al clúster 1. A més, han puntuat positivament l'espai de la cabina.

Cluster 3		
	Mean in category	Overall mean
trato.TR	9.250478	8.307851
trato.LL	9.133942	8.207039
idiomas	8.687140	7.959233
aspect.gener	8.395084	7.978355
espai.cabin	4.896406	5.465164
conf.seat	5.588385	6.194617
conf.cabina	6.104366	6.628337

El clúster 3 està format per gent majoritàriament satisfeta, sobretot amb el tracte, els idiomes i l'aspecte general. Ara bé, han puntuat malament l'espai de la cabina i la comoditat.

Cluster 4		
	Mean in category	Overall mean
mean.SAT	8.834925	7.166966

El clúster 4 està format per gent satisfeta o molt satisfeta. Ho han puntuat tot per sobre de la mitjana.



## Aplicació del MCA enriquit per clustering

### Inèrcia dels eixos. Quins factors retenim?

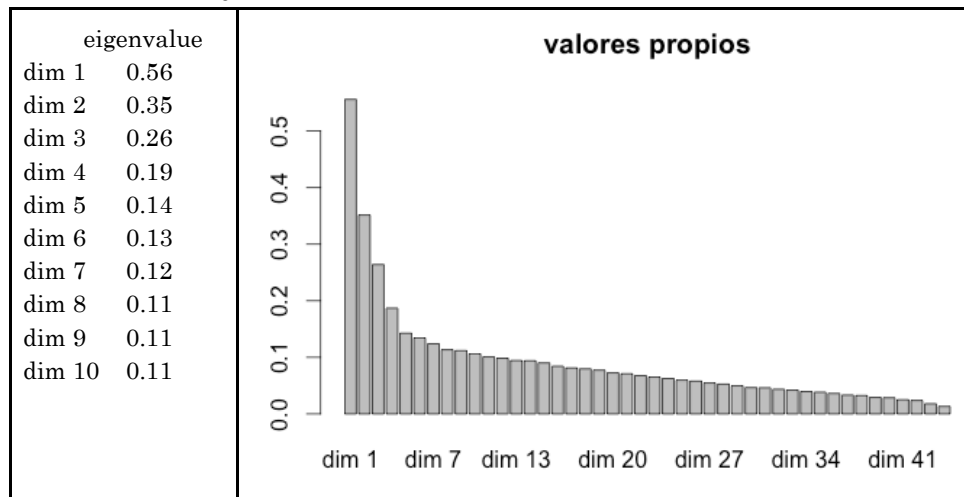


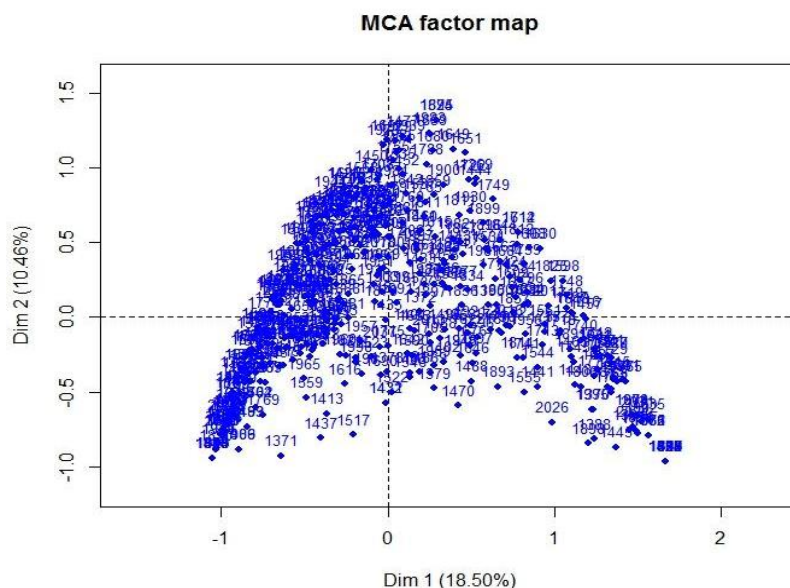
Figura 9. Eigenvalues: valors propis

A partir de la Figura 9 i de la norma del colze, tindrem en compte les 4 primeres dimensions durant aquesta part de l'estudi, el MCA. A partir de la quarta dimensió la diferència entre els diferents valors és molt petita i per tant aportaran informació molt menys rellevant que les anteriors.

### Interpretació dels factors

Després d'imputar valors als nuls (valors NA de la mostra) tenint en compte els 2 eixos suggerits per l'estimació de R i obtenir una nova taula de valors, discretitzem les nostres dades a partir dels punts de tall deduïts dels quartils de la mitjana de satisfacció.

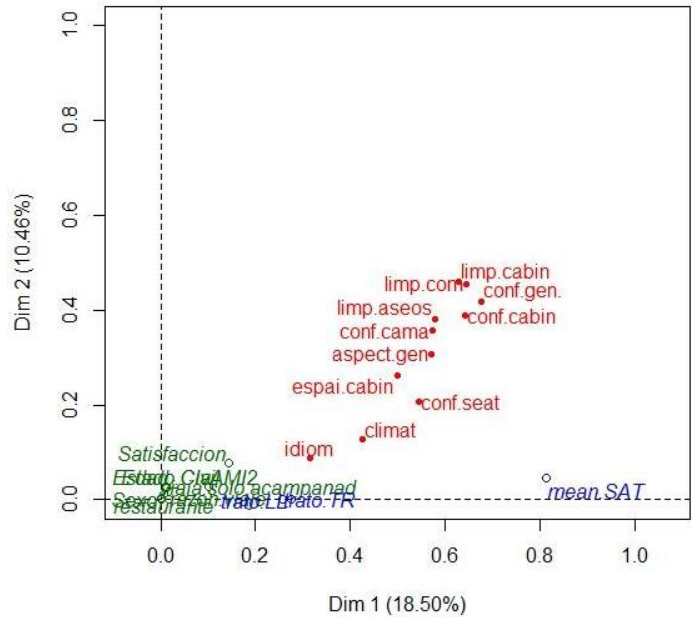
Ara podem començar a analitzar el nostre estudi MCA. Per començar, podem observar a la Figura 10 observar que el núvol de punts-individus té una repartició bastant regular sobre el pla factorial. Això vol dir que reflecteix tendències compartides per grups grans d'individus.



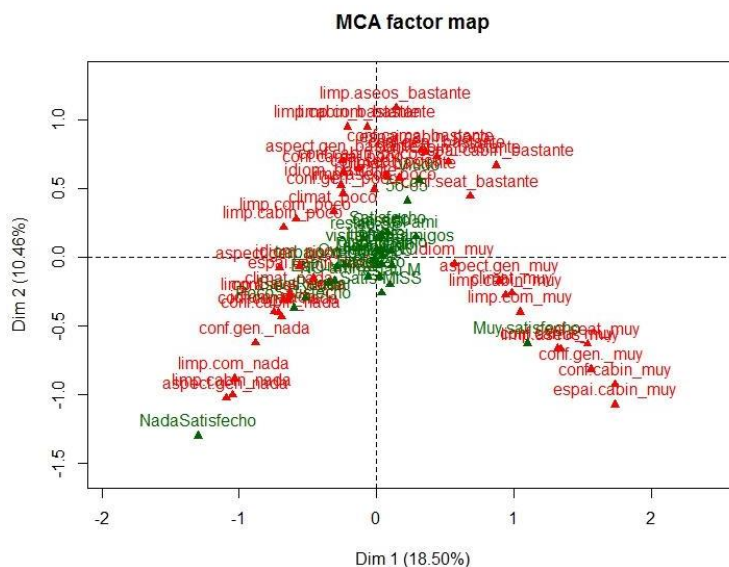
També podem observar que es produeix l'efecte Gutmann. Aquest efecte indica que hi ha un fort factor d'escala entre les categories de les variables. Es manifesta per la forma que pren el núvol de punts, gairebé parabòlica. Això és bo ja que ens diu que reflecteix tendències generals i no específiques.

Figura 10. MCA. Mapa factorial dels individus en funció de les dimensions 1 i 2.

A la Figura 11 podem observar la distribució de la puntuació de les diferents variables en funció de la satisfacció. Com més a la dreta es troben els punts, més disperses estan les puntuacions d'aquella variable i més associació tenen amb la satisfacció. Com més a l'esquerra es troben els punts, més comprimides estan les puntuacions d'aquella variable i menys associació tenen amb la satisfacció. No s'observa relació entre les variables categòriques i les suplementàries.



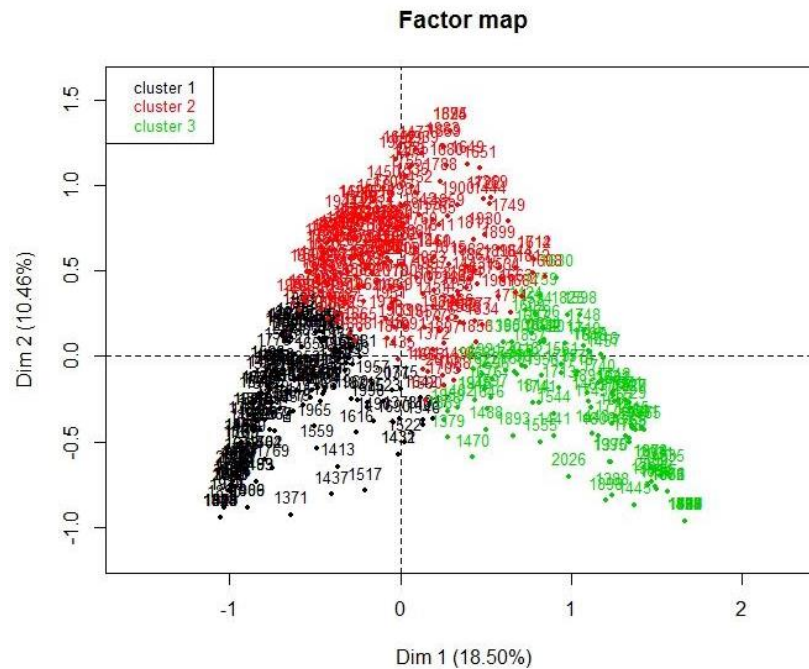
*Figura 11. MCA. Mapa factorial de la distribució de les variables quantitatives i qualitatives en funció de les dimensions 1 i 2.*



*Figura 12. MCA. Mapa factorial de la distribució de les categories de les variables quantitatives i les categories de la variable de satisfacció en funció de les dimensions 1 i 2.*

A la Figura 12 podem observar que la distribució de les categories de les variables quantitatives i les categories de la variable de satisfacció tenen formes parabòliques similars. Per tant, tornem a tenir un indicador de que es produeix l'efecte Gutmann.

## Síntesi del MCA a través dels clusters



*Figura 10. MCA. Mapa factorial dels individus en funció de les dimensions 1 i 2, classificats per clustering.*

Tenim 3 clústers amb els següents nombres d'individus: clúster 1 (183), clúster 2 (194) i clúster 3 (112).

### Descripció dels clústers

El clúster 1 agrupa els individus insatsifets. Han puntuat negativament per sota de la mitjana global, especialment el tracte.

El clúster 2 agrupa individus que han puntuat per sobre de la mitjana global.

El clúster 3 agrupa individus que han puntuat positivament per sobre de la mitjana global, especialment el tracte.

## Conclusions globals

Una primera aproximació a les dades de la mostra fa pensar que tots els individus estan raonablement satisfets amb l'experiència (mitjana de satisfacció de 7.167). Ara bé, un estudi més profund ha revelat que existeixen tendències en les puntuacions. Els individus puntuen determinats aspectes positiva o negativament en funció de les seves característiques personals com l'edat o la raó per la que viatgen. A més, hem pogut observar que si una tendència (dimensió) és molt forta, influeix molt en la resta i moltes vegades les difumina.

Deixant de banda l'estudi en sí, la realització d'aquest treball ens ha permès fer un ús més extens de R del que havíem fet fins ara i aprendre una mica més del seu potencial, per exemple, amb l'ús de paquets addicionals com FactoMineR i MissMDA.