

Ancianos

Ferran Garcia

2024-01-25

Contents

1	Introducción	2
1.1	Contexto	2
1.2	Objetivos	2
1.3	Estructura	2
2	Metodología	3
2.1	Técnicas	3
2.2	Herramientas	3
3	Reducción dimensional	4
3.1	Selección de variables y adecuación de los datos	4
3.2	Número, rotación y puntuaciones factoriales	7
4	Clústers	11
4.1	Método jerárquico	11
4.2	Método de k-medias	11
5	Annexo	12
5.1	Variables excluidas	12
5.2	Gráficos	12
6	Bibliografía	14

1 Introducción

1.1 Contexto

El envejecimiento demográfico emerge como un fenómeno geográfico crítico que podría acarrear serias implicaciones macroeconómicas estructurales en nuestra sociedad. En consecuencia, resulta imperativo abordar de manera inmediata políticas orientadas a mejorar la calidad de vida de la población mayor, fundamentadas en estudios que delinee su situación espacial y social, particularmente en aquellos entornos urbanos donde su presencia es más significativa.

1.2 Objetivos

En el actual escenario, nos enfrentamos al desafío de localizar áreas urbanas en el municipio de Sevilla que compartan una estructura demográfica similar, especialmente en lo que respecta a la población mayor. Para abordar este objetivo, emplearemos distintos métodos de análisis estadístico multivariante y espacial. La meta es desarrollar servicios sociales y de asistencia que tomen en cuenta las características específicas de los ancianos y su distribución geográfica en la ciudad. En este contexto, nuestro enfoque consiste en identificar segmentos homogéneos de la población anciana en áreas urbanas. La ejecución de este proceso nos permitirá diseñar servicios sociales y de asistencia que se adapten a las necesidades particulares de los ancianos, considerando tanto su tipología como su ubicación en la ciudad.

1.3 Estructura

2 Metodología

2.1 Técnicas

2.2 Herramientas

3 Reducción dimensional

Disponemos de muchas variables, eso motiva las técnicas de reducción dimensional. Usaremos análisis factorial de componentes principales, la idea es sintetizar todas las medidas disponibles en variables latentes. Lo haremos aprovechando la correlación que comparten. Para ello seleccionamos las variables que incluiremos a partir de un breve análisis exploratorio. Luego comprobaremos si es pertinente usar los procedimientos en cuestión en el conjunto de datos resultante. Entonces realizaremos varios ajustes con diferente número de factores hasta alcanzar la descomposición más satisfactoria en términos de representabilidad de las variables iniciales. Finalmente aplicaremos técnicas de rotación con tal de lograr factores ortogonales que sean interpretables. Con éstos extraeremos las puntuaciones factoriales que más adelante usaremos para crear clústers.

3.1 Selección de variables y adecuación de los datos

La selección de variables sirve para subsanar potenciales problemas y afinar en el cálculo de factores. Es importante tanto para que se pueda realizar el análisis factorial como para que éste sea eficiente y se ajuste un modelo robusto con interpretación relativamente simple evitando sobreajuste.

Por ejemplo, si nos fijamos en las variables demográficas veremos que se incluyen múltiples medidas de población quedando algunas determinadas completamente por el resto. En cuanto a la *situación laboral* tenemos población activa, ocupada, inactiva y parados. Entonces el total queda determinado por otra variable, ya sea la población total o u subgrupo de ésta y por tanto son linealmente dependientes. Por tanto la matriz de correlaciones no será invertible y en consecuencia no se podrá realizar el análisis factorial. Lo mismo ocurre con algunas variables del estado de las viviendas. En el anexo se muestra una lista completa de las variables que han sido excluidas por este motivo.

Por otro lado es importante que las variables seleccionadas estén correlacionadas entre sí. Para comprobarlo representamos gráficamente la correlación entre las variables seleccionadas y recogemos en una tabla estadísticos relevantes para la selección de variables como:

- Estimación inicial de comunalidades¹
- Suma de correlaciones absolutas
- Número de variables no correlacionadas²

Las variables que se muestran en la tabla quedan excluidas del análisis factorial.

Table 1: Correlaciones de las variables seleccionadas, se muestran las 15 con menor comunalidad

	Comunalidad	Suma	p.val
CON1HABI	0.5205	6.88	16
AP	0.4817	9.65	28
DEL81AL90	0.3475	4.44	7
CON5OMAS	0.2745	4.26	8
C	0.2607	7.85	25
EA	0.2160	6.15	27
TRAEVEN	0.2101	4.25	2
AFAM	0.1894	5.05	7
COOP	0.1453	2.61	1

¹Correlación múltiple al cuadrado

²Test de correlación de Pearson, nivel de significación del 5%

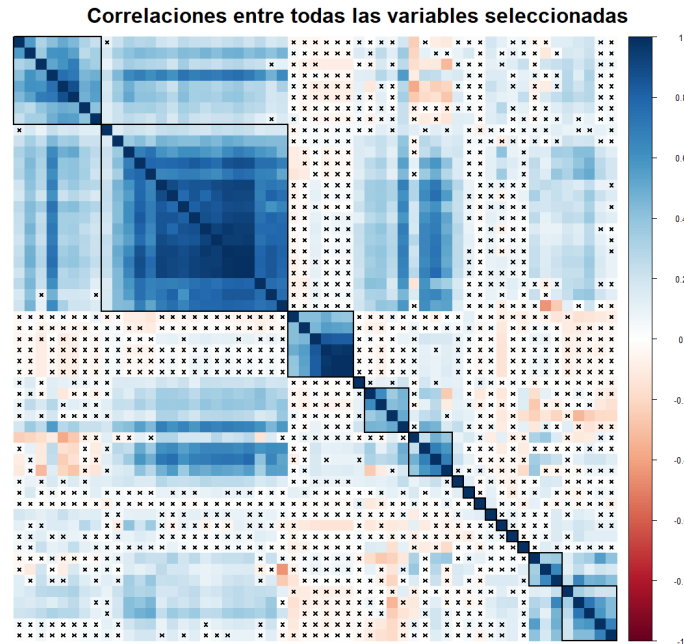


Figure 1: Correlaciones entre todas las variables.

Veamos ahora si el conjunto de datos seleccionados resulta apropiado para el análisis factorial o si debemos seleccionar variables con un criterio más estricto. Para ello realizamos varios tests:

Esfericidad de Barlett:

- Hipótesis nula: La matriz de correlaciones es la identidad (no existe ninguna correlación)
- Bajo H_0 el estadístico de contraste sigue una χ^2 , asintóticamente
- Disponemos de más de 400 observaciones, se sostiene la suposición asintótica

```
cortest.bartlett(R1, n = nrow(X1))["p.value"] # podemos rechazar  $H_0$ 
```

```
## [1] 0
```

Kaiser, Meyer, Olkin:

- Medida de adecuación de la muestra para un análisis factorial
- Compara la suma de cuadrados de la “imagen” de la matriz de correlaciones con la original
- Está acotada entre 0 y 1, según Kaiser:
 - Valores mayores a 0,9 son maravillosos
 - Valores mayores a 0,8 son meritorios
 - Valores mayores a 0,7 son medios
 - Valores mayores a 0,6 son mediocres
 - Valores mayores a 0,5 son miserables
 - Valores menores a 0,5 son inaceptables

```
KMO(R1)[[1]] # según Kaiser, la adecuación de la muestra es meritoria
```

```
## [1] 0.8141468
```

Cabe remarcar que tras decidir que el análisis factorial una práctica habitual es estandarizar las variables con con propósito de comparabilidad entre ellas. Sin embargo tras realizar pruebas con y sin estandarizar hemos obtenido mejores resultados sin estandarizar. Teniendo en cuenta que la extracción de factores nos servirá para formar clústers, que es nuestro principal objetivo consideramos más importante conseguir un buen ajuste que obtener comparabilidad en una etapa intermedia del trabajo.

3.2 Número, rotación y puntuaciones factoriales

Ahora nuestra intención es identificar una estructura latente dentro del conjunto de datos. Se focaliza el interés en los factores capaces de explicar una proporción significativa de la variabilidad presente en los datos. Para determinar cuántos factores debemos retener, empleamos dos criterios. En primer lugar, se puede dibujar un scree plot para evaluar los *eigenvalues* de los factores estimados. Además, aplicamos la Regla de Kaiser, que sugiere retener aquellos factores cuyos *eigenvalues* superen la unidad. Otra estrategia consiste en seleccionar tantos factores como sean necesarios para explicar alrededor de un 70% de la variabilidad de los datos.

Regla de Kaiser:

En el ámbito del álgebra una matriz (como la de correlaciones) se puede interpretar como una transformación lineal. Los *eigenvectors* asociados a una transformación son aquellos vectores cuya dirección es invariable a la misma y los *eigenvalues* aparejados son la medida en la que la dirección se alarga o mengua en magnitud. Bajo ciertas condiciones de regularidad estos vectores forman una base sobre la que se puede descomponer la transformación en cuestión. La idea de Kaiser es tomar tantos factores como autovectores de “alarguen” (su importancia crezca) en el proceso de descomposición.³

```
autov_1 = eigen(R1)[["values"]]; sort(autov_1, decreasing = T) %>% round(digits = 2)
```

```
## [1] 16.41  5.18  3.86  3.52  2.08  1.82  1.30  1.01  0.96  0.90  0.85  0.75
## [13]  0.71  0.66  0.62  0.59  0.55  0.49  0.39  0.36  0.33  0.30  0.26  0.23
## [25]  0.22  0.21  0.21  0.20  0.18  0.16  0.13  0.12  0.11  0.10  0.07  0.05
## [37]  0.04  0.04  0.02  0.01  0.01  0.00  0.00  0.00  0.00  0.00  0.00
```

```
sum(autov_1 > 1) # según el criterio de Kaiser, usaremos 8 factores
```

```
## [1] 8
```

En este caso Kaiser recomienda tomar 8 factores. No obstante los últimos están muy cerca de 1, es algo que debemos tener en cuenta para la selección del número de factores.

Proporción de la varianza explicada:

Para emplear el método de la varianza explicada por factores necesitamos realizar un análisis factorial como tal. De manera que ajustamos un modelo con un número arbitrario de factores y comprobamos cuántos son necesarios

```
fit_1 = factanal(na.omit(X1[, -1]), factors = 10, lower = .01)
fit_1[["loadings"]]
```

```
##
## Loadings:
##          Factor1 Factor2 Factor3 Factor4 Factor5 Factor6 Factor7 Factor8
## TPOBDCHO    0.940   0.275         0.139
## 65-69AÑOS    0.925   0.205                 0.123
## 70-74AÑOS    0.883   0.235         0.159
## 75-79AÑOS    0.797   0.299         0.301           0.132
## 80-85AÑOS    0.648   0.303         0.376           0.292
## MÁS85AÑOS    0.538   0.431         0.298           0.334
## POBDV        0.881   0.254                 0.175           0.122
```

³Para más información sobre los *eigenvectors* y *eigenvalues* consultar este enlace

## TMUN	0.835		0.289		-0.160		
## TCOM	0.819	0.337		0.102	0.116		
## TESP	0.594	0.707					
## TEXT	0.165	0.510				0.207	
## ANALF	0.258	-0.313			0.249	0.109	0.861
## SE	0.828	-0.332	0.124		0.198		0.107
## 1G	0.786	0.191	0.116	0.201	-0.114		-0.259
## 3G	0.124	0.955			-0.122		
## INACTIVOS	0.960	0.188	0.129				
## AP		0.581					
## IND	0.423	0.300				-0.127	
## C	0.242				0.162		
## EMPEMP	0.120	0.492					-0.124
## TRAFIJO	0.353	0.557			-0.146		
## VPRINC	0.949	0.249	0.119				
## VNOPRIN	0.181	0.166					0.147
## COLEC		0.160	-0.113	0.125		0.961	
## CON1VIV			0.232		0.950		0.168
## CON2VIV		0.110	0.648		0.518		
## CON3VIV	0.828	0.173			-0.507		
## MENOS30			0.681				
## D30A60	0.692	-0.365		-0.580	-0.102		
## D60A90	0.325			0.871			
## D90A121		0.674	0.115	0.165	0.158		
## CON2HABI	0.219		0.436	-0.306	0.192		0.202
## CON4HABI	0.756	-0.262					
## CON5HABI	0.362			0.570			
## CON6HABI	0.140	0.913	0.109		0.134		
## CON1OCUP	0.878	0.144	0.188				0.148
## CON2OCUP	0.930	0.232					
## CON4OCUP	0.143	0.266	0.296				
## ANT1941	0.101	0.141	0.867				
## DEL61AL80	0.316		-0.374	0.331			0.102
## PROP	0.662	0.259	-0.272	0.323	0.129		
## ALQ	0.220	0.212	0.779				
## AGUAVIV			0.995				
## ELECTRIC			0.995				
## NOREFRIGER			0.972				
## NOGAS	0.110		0.829				
##	Factor9	Factor10					
## TPOBDCHO							
## 65-69AÑOS							
## 70-74AÑOS							
## 75-79AÑOS							
## 80-85AÑOS							
## MÁS85AÑOS							
## POBDV		0.144					
## TMUN							
## TCOM							
## TESP							
## TEXT							
## ANALF							
## SE	0.354						
## 1G	-0.441						


```

## 3G
## INACTIVOS
## AP
## IND
## C
## EMPEMP
## TRAFIJO
## VPRINC
## VNOPRIN
## COLEC
## CON1VIV
## CON2VIV
## CON3VIV
## MENOS30
## D30A60          0.125
## D60A90
## D90A121
## CON2HABI        0.108
## CON4HABI
## CON5HABI
## CON6HABI
## CON1OCUP        0.348
## CON2OCUP       -0.237
## CON4OCUP
## ANT1941
## DEL61AL80    0.219    0.142
## PROP
## ALQ
## AGUAVIV
## ELECTRIC
## NOREFRIGER
## NOGAS
##
##
##          Factor1 Factor2 Factor3 Factor4 Factor5 Factor6 Factor7 Factor8
## SS loadings    14.344   5.620   3.688   3.424   1.906   1.842   1.311   1.096
## Proportion Var   0.312   0.122   0.080   0.074   0.041   0.040   0.029   0.024
## Cumulative Var   0.312   0.434   0.514   0.589   0.630   0.670   0.699   0.722
##
##          Factor9 Factor10
## SS loadings     0.430   0.289
## Proportion Var   0.009   0.006
## Cumulative Var   0.732   0.738

```

Esta salida muestra por un lado las cargas factoriales de cada variable inicial y por el otro la proporción de varianza explicada por cada uno de los factores. Nos interesa la parte final. Nos fijamos en que se explica alrededor del 70% de la varianza con 7 u 8 factores. También notamos que los primeros factores recogen sustancialmente más varianza que el resto.

Scree plot:

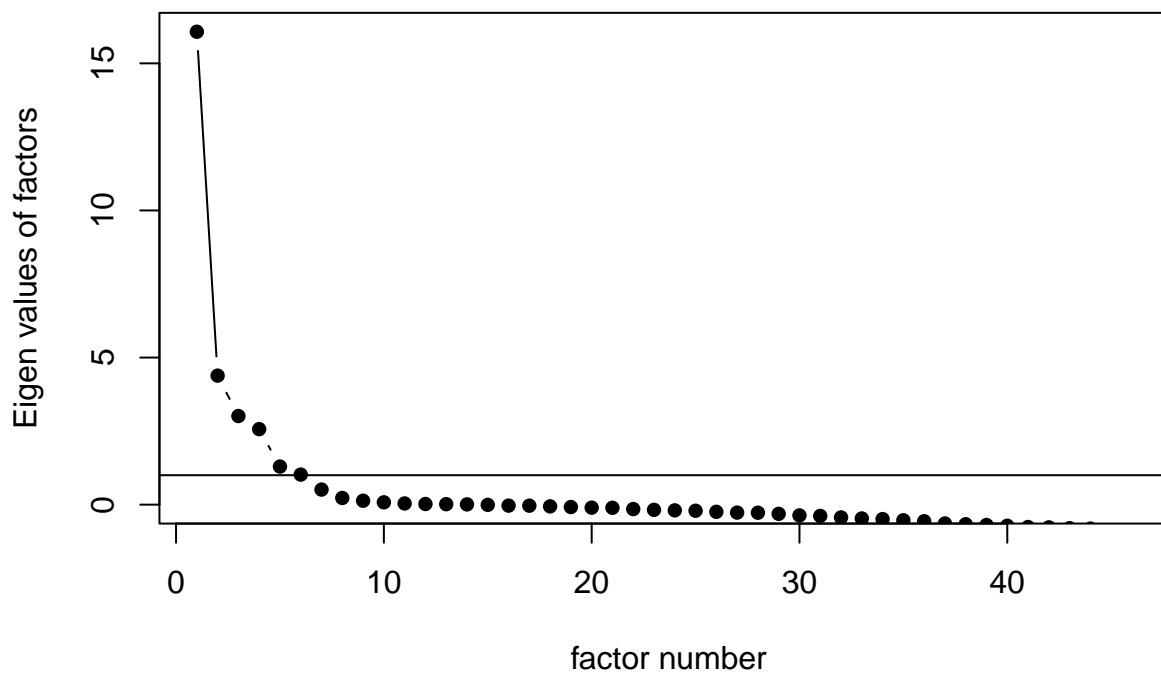
En este caso interpretamos gráficamente los *eigenvalues* asociados a factores con el número que usamos. Notamos que los cuatro primeros son relativamente más altos a los dos siguientes, que resultan ser los últimos que superan el umbral de 1. En este caso el gráfico sugiere usar alrededor de 6 factores.

```
scree(X1[,-1], pc = FALSE)
```

```
## Warning in fa.stats(r = r, f = f, phi = phi, n.obs = n.obs, np.obs = np.obs, :  
## The estimated weights for the factor scores are probably incorrect. Try a  
## different factor score estimation method.
```

```
## Warning in fac(r = r, nfactors = nfactors, n.obs = n.obs, rotate = rotate, : An  
## ultra-Heywood case was detected. Examine the results carefully
```

Scree plot



Interpretación en conjunto:

Los diversos métodos indican una cantidad de factores similar: Entre 6 y 8. Puesto que utilizar 7 factores concuerda con los resultados de los diferentes criterios y está en medio de las sugerencias de cada uno usaremos 7 factores para reducir la dimensionalidad de los datos.

```
fit_2 = factanal(na.omit(X1[,-1]), factors = 7, lower = .007, nstart)
```

breve interpretación de los factores

4 Clústers

4.1 Método jerárquico

4.2 Método de k-medias

5 Anexo

5.1 Variables excluidas

5.1.1 Variables demográficas

- Categoría de edad: 60-64 años
- Sexo: Población masculina
- Procedencia: Provincial
- Nivel de estudios: Segundo grado
- Situación laboral: Activos
- Situación de actividad: Ocupados
- Categoría profesional: Servicios
- Posición profesional: Empleados
- Temporalidad profesional: Otro
- Residencia: Alojamiento

5.1.2 Variables de la vivienda

- Tamaño en m²: +120
- Número de habitaciones: 3
- Número de ocupantes: 3
- Año de construcción: 41-60
- Régimen legal: Otro

5.2 Gráficos

5.2.1 Correlaciones

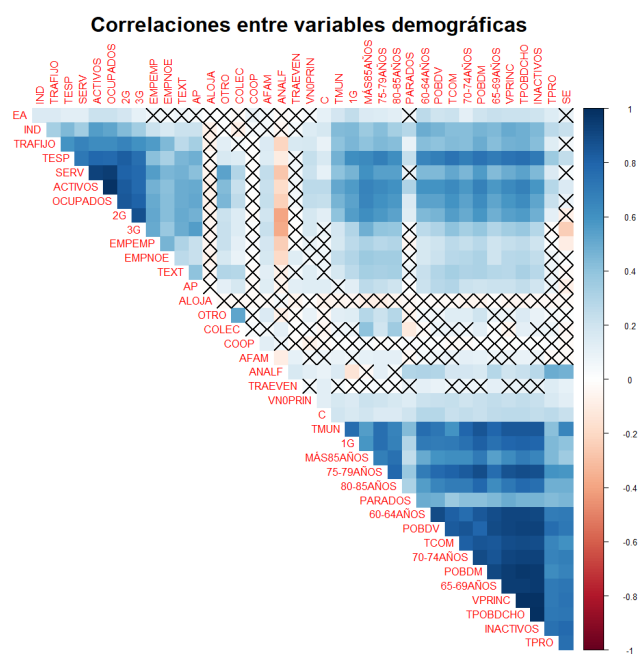


Figure 2: Correlaciones entre variables demográficas.

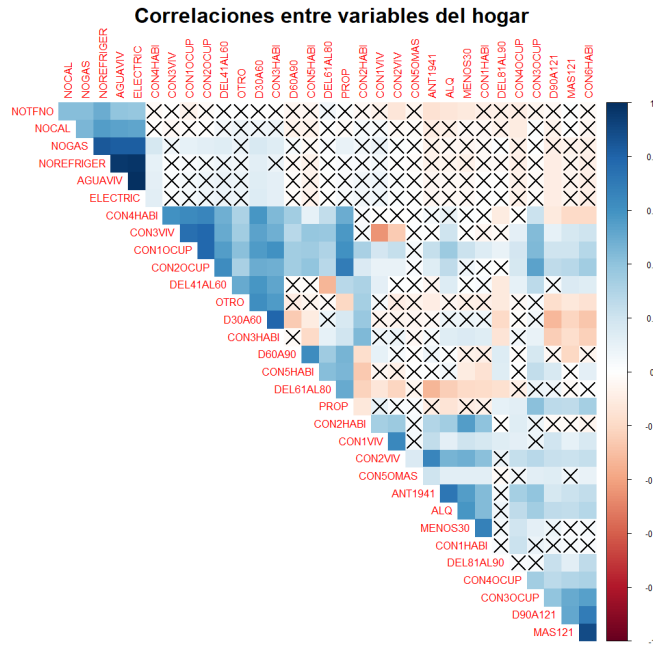


Figure 3: Correlaciones entre variables de las viviendas.

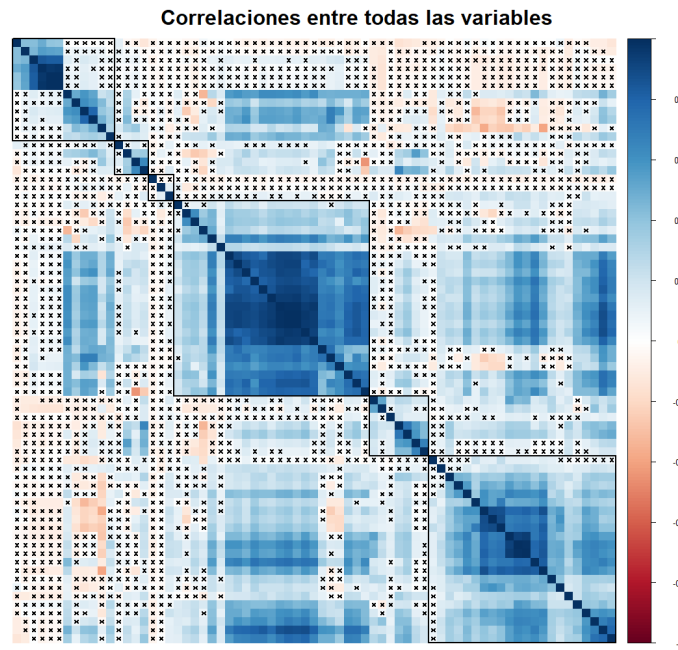


Figure 4: Correlaciones entre todas las variables.

6 Bibliografía

- Datos vectoriales
- Análisis clúster I
- Análisis clúster II
- Análisis clúster III
- ClustGeo
- ClusterR