

Pràctica (Part 1): Selecció del conjunt de dades

Mireia Arnal Avalo

Visualització de dades · 06/05/2024

Índex

1	Justificació	2
2	Rellevància	2
3	Complexitat	2
4	Originalitat	4
5	Diccionari	9

La temàtica d'estudi és la transformació del [Nivell educatiu a Barcelona](#) al llarg dels anys segons el districte, barri, edat o sexe del participant. Per tal de fer-ho, s'ha seleccionat 3 datasets del catàleg Open Data BCN (tercer nivell) en un rang de 10 anys entre cadascun:

- 2003_bcn.csv
- 2013_bcn.csv
- 2023_bcn.csv

1 Justificació

L'educació, des d'una perspectiva social, és un dels pilars fonamentals per al desenvolupament individual i col·lectiu d'una societat. L'anàlisi de la seva evolució és essencial per a comprendre i revelar disparitats i desigualtats que necessiten abordar-se per a la creació d'un entorn més inclusiu per a tothom.

En aquest sentit, Barcelona ofereix factors únics i rics en termes d'origen ètnic, cultural i socioeconòmic. És una ciutat cosmopolita i dinàmica, amb canvis polítics i socials significatius en les últimes dècades. Tot plegat, crea una gran varietat de situacions i reptes que poden ser interessants d'explorar per a la posterior presa de decisions polítiques, la planificació urbana i la implementació d'intervencions educatives de qualitat.

Quant al motiu de selecció d'aquest conjunt de dades, ha estat àmpliament influenciat per causes personals. Pretén, llavors, ser un treball per enriquir el coneixement i fer una posada en pràctica de les competències apreses al llarg de l'assignatura.

2 Rellevància

Tot i que les dades no són necessàriament actuals, abasten un període de temps valuós que permet investigar tendències i canvis. També cal destacar que és un registre històric amb una freqüència d'actualització anual, de manera que cada setembre s'afegeix les dades de l'any actual. Com que som a maig, encara no hi ha les de 2024 i s'empra les de l'any anterior.

D'altra banda, i com s'ha mencionat anteriorment, analitzar el nivell educatiu de la població té una rellevància significativa en el context social i educatiu, proporcionant informació valuosa sobre els reptes i les oportunitats de millora a futur. Pot ser d'especial interès per a figures polítiques, educadors, famílies, investigadors i altres agents interessats.

A més a més, es dona una especial importància a la perspectiva de gènere perquè sigui una anàlisi completa. La inclusió de dades desglossades per sexe permet identificar possibles disparitats en el nivell d'estudis entre homes i dones.

3 Complexitat

```
# Carreguem el conjunt de dades
bcn03 <- read.csv('2003_bcn.csv', row.names = NULL, stringsAsFactors = TRUE)
bcn13 <- read.csv('2013_bcn.csv', row.names = NULL, stringsAsFactors = TRUE)
bcn23 <- read.csv('2023_bcn.csv', row.names = NULL, stringsAsFactors = TRUE)
```

Primerament, l'arxiu 2003_bcn.csv:

```
# Mostrem l'estructura del conjunt de dades
str(bcn03)
```

```
## 'data.frame': 13071 obs. of 9 variables:
## $ Data_Referencia: Factor w/ 1 level "2003-01-01": 1 1 1 1 1 1 1 1 1 ...
## $ Codi_Districte : int 1 1 1 1 1 1 1 1 1 ...
## $ Nom_Districte : Factor w/ 10 levels "Ciutat Vella",...: 1 1 1 1 1 1 1 1 1 ...
## $ Codi_Barri : int 1 1 1 1 1 1 1 1 1 ...
## $ Nom_Barri : Factor w/ 73 levels "Baró de Viver",...: 23 23 23 23 23 23 23 23 23 ...
## $ Valor : Factor w/ 723 levels "...", "10", "100",...: 398 376 528 48 559 107 646 98 611 71 ..
## $ NIV_EDUCA_esta : int 1 1 1 1 1 1 1 1 1 ...
## $ EDAT_Q : int 3 3 4 4 5 5 6 6 7 7 ...
## $ SEXE : int 1 2 1 2 1 2 1 2 1 2 ...
```

Presenta un total de **13071 observacions** i **9 variables**, les quals són del tipus:

- Quantitatives:
 - Integer: Codi_Districte, Codi_Barri, NIV_EDUCA_esta, EDAT_Q i SEXE.
- Qualitatives:
 - Factor: Data_Referencia, Nom_Districte, Nom_Barri i Valor.

Com que R (i el mateix dataset) atorga un tipus de variables que no coincideixen amb les reals, es realitzaran els següents canvis:

- Data_Referencia passarà a integer i només es deixarà l'any corresponent al nom del fitxer (2003 per 2003_bcn.csv, 2013 per 2013_bcn.csv i 2023 per 2023_bcn.csv).
- Valor passarà a integer.
- NIV_EDUCA_esta passarà a factor i s'afegirà, en text, el nivell d'estudis.
- SEXE seguirà sent integer, però passarà de ser 1 i 2 a 0 i 1, respectivament.

A banda d'aquestes millores, també es realitzarà un canvi en el nom de les variables perquè hi hagi homogeneïtat en l'estil, s'afegirà una nova variable (id) i es combinarà els tres fitxers en un.

Després hi ha l'arxiu 2013_bcn.csv:

```
# Mostrem l'estructura del conjunt de dades
str(bcn13)
```

```
## 'data.frame': 12915 obs. of 9 variables:
## $ Data_Referencia: Factor w/ 1 level "2013-01-01": 1 1 1 1 1 1 1 1 1 ...
## $ Codi_Districte : int 1 1 1 1 1 1 1 1 1 ...
## $ Nom_Districte : Factor w/ 10 levels "Ciutat Vella",...: 1 1 1 1 1 1 1 1 1 ...
## $ Codi_Barri : int 1 1 1 1 1 1 1 1 1 ...
## $ Nom_Barri : Factor w/ 73 levels "Baró de Viver",...: 23 23 23 23 23 23 23 23 23 ...
## $ Valor : Factor w/ 743 levels "...", "10", "100",...: 617 158 390 735 236 21 225 668 400 532
## $ NIV_EDUCA_esta : int 1 1 1 1 1 1 1 1 1 ...
## $ EDAT_Q : int 3 3 4 4 5 5 6 6 7 7 ...
## $ SEXE : int 1 2 1 2 1 2 1 2 1 2 ...
```

Presenta un total de **12915 observacions** i **9 variables**, les quals són les mateixes que en el cas anterior. Per acabar, el fitxer `2023_bcn.csv`:

```
# Mostrem l'estructura del conjunt de dades
str(bcn23)

## 'data.frame':    11562 obs. of  9 variables:
## $ Data_Referencia: Factor w/ 1 level "2023-01-01": 1 1 1 1 1 1 1 1 1 1 ...
## $ Codi_Districte : int  1 1 1 1 1 1 1 1 1 1 ...
## $ Nom_Districte  : Factor w/ 10 levels "Ciutat Vella",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ Codi_Barri     : int  1 1 1 1 1 1 1 1 1 1 ...
## $ Nom_Barri      : Factor w/ 73 levels "Baró de Viver",...: 23 23 23 23 23 23 23 23 23 23 ...
## $ Valor          : Factor w/ 802 levels "...", "10", "100",...: 1 1 1 1 2 1 590 1 1 1 ...
## $ NIV_EDUCA_esta : int  1 1 1 1 1 1 1 1 1 1 ...
## $ EDAT_Q         : int  4 5 5 6 6 7 7 8 8 9 ...
## $ SEXE           : int  2 1 2 1 2 1 2 1 2 1 ...
```

Presenta un total de **11562 observacions** i **9 variables**, les quals són les mateixes que en el primer cas.

Amb més d'11000 registres per conjunt de dades, el volum és suficient per proporcionar una mostra representativa de la realitat i dur a terme anàlisis detallades.

Tenim un total de 9 variables simples i estructurades. Per una banda, les dades categòriques del joc de dades descriuen aspectes ordinals (nivell d'estudis) i nominals (nom del districte o barri). Per l'altra, les dades numèriques descriuen quantitats, sempre senceres (edat o valor). També s'inclou un altre tipus de dades, que és la binària (sexe) i la data (any).

En definitiva, l'existència de múltiples variables i la possibilitat de relacionar-les entre si aporta complexitat i la possibilitat d'una investigació en profunditat.

4 Originalitat

Per tal de donar un enfocament nou o una perspectiva complementària a les dades educatives, es podria afegir dades econòmiques com són la [Renda tributària neta mitjana per llar](#). No obstant això, com que aquestes agafen el rang d'entre l'any 2015 i 2021, no es poden afegir al projecte.

Llavors només queda enriquir el conjunt de dades existent mitjançant la transformació de variables existents per la generació de noves mètriques o indicadors.

Primer, però, es farà els canvis mencionats en l'anterior apartat i es fusionarà tot en un únic dataset.

```
# Deixem l'any corresponent al fitxer per la variable 'Data_Referencia'
bcn03$Data_Referencia <- 2003
bcn13$Data_Referencia <- 2013
bcn23$Data_Referencia <- 2023

# Combinem els conjunts de dades en un de sol
bcn_educacio <- rbind(bcn03, bcn13, bcn23)

# Creem una nova variable 'id' com a identificador únic
```

```
bcn_educacio$id <- 1:nrow(bcn_educacio)
```

```
# Guardem el dataset combinat en un nou fitxer CSV  
write.csv(bcn_educacio, 'bcn_educacio.csv', row.names = FALSE)
```

```
# Mostrem l'estructura del conjunt de dades  
str(bcn_educacio)
```

```
## 'data.frame': 37548 obs. of 10 variables:  
## $ Data_Referencia: num 2003 2003 2003 2003 2003 ...  
## $ Codi_Districte : int 1 1 1 1 1 1 1 1 1 1 ...  
## $ Nom_Districte : Factor w/ 10 levels "Ciutat Vella",...: 1 1 1 1 1 1 1 1 1 1 ...  
## $ Codi_Barri : int 1 1 1 1 1 1 1 1 1 1 ...  
## $ Nom_Barri : Factor w/ 73 levels "Baró de Viver",...: 23 23 23 23 23 23 23 23 23 23 ...  
## $ Valor : Factor w/ 992 levels "...", "10", "100",...: 398 376 528 48 559 107 646 98 611 71 ..  
## $ NIV_EDUCA_esta : int 1 1 1 1 1 1 1 1 1 1 ...  
## $ EDAT_Q : int 3 3 4 4 5 5 6 6 7 7 ...  
## $ SEXE : int 1 2 1 2 1 2 1 2 1 2 ...  
## $ id : int 1 2 3 4 5 6 7 8 9 10 ...
```

Ara queda fer els canvis per cada variable. Es mira que cada columna sigui una variable, cada fila sigui una observació i cada cel·la només tingui un valor.

```
# Convertim la variable 'Data_Referencia' de double a integer  
bcn_educacio$Data_Referencia <- as.integer(bcn_educacio$Data_Referencia)  
  
# Canviem el nom de la variable 'Data_Referencia' a 'Any'  
names(bcn_educacio)[names(bcn_educacio) == 'Data_Referencia'] <- 'Any'  
  
# Convertim la variable 'Valor' de factor a integer  
bcn_educacio$Valor <- as.integer(bcn_educacio$Valor)  
  
# Reemplaçem el valor '..' per NA a la variable 'Valor'  
bcn_educacio$Valor[bcn_educacio$Valor == '..'] <- NA  
  
# Canviem els valors de la variable 'NIV_EDUCA_esta' i la convertim d'int a factor  
bcn_educacio$NIV_EDUCA_esta <- factor(bcn_educacio$NIV_EDUCA_esta,  
                                     levels = c(1, 2, 3, 4, 5),  
                                     labels = c('Sense estudis', 'Estudis primaris o EGB',  
                                                'ESO o CFGB', 'Batxillerat, BUP, COU o CFGM',  
                                                'Estudis universitaris o CFGS'))  
  
# Canviem el nom de la variable 'NIV_EDUCA_esta' a 'Nivell_Educatiu'  
names(bcn_educacio)[names(bcn_educacio) == 'NIV_EDUCA_esta'] <- 'Nivell_Educatiu'  
  
# Canviem el nom de la variable 'EDAT_Q' a 'Edat'  
names(bcn_educacio)[names(bcn_educacio) == 'EDAT_Q'] <- 'Edat'  
  
# Reemplaçem el valor '1' i '2' per '0' i '1', respectivament, a la variable 'SEXE'  
bcn_educacio$SEXE <- ifelse(bcn_educacio$SEXE == 1, 0,  
                            ifelse(bcn_educacio$SEXE == 2, 1, bcn_educacio$SEXE))  
  
# Canviem el nom de la variable 'SEXE' a 'Sexe'
```

```
names(bcn_educacio)[names(bcn_educacio) == 'SEXE'] <- 'Sexe'
```

```
# Convertim la variable 'Sexe' de numeric a integer
```

```
bcn_educacio$Sexe <- as.integer(bcn_educacio$Sexe)
```

```
# Capturem la variable 'id', l'eliminem i l'afegim a la primera posició
```

```
id_column <- bcn_educacio$id
```

```
bcn_educacio <- bcn_educacio[, -which(names(bcn_educacio) == 'id')]
```

```
bcn_educacio <- cbind(id = id_column, bcn_educacio)
```

```
# Mostrem l'estructura del conjunt de dades
```

```
str(bcn_educacio)
```

```
## 'data.frame': 37548 obs. of 10 variables:
## $ id : int 1 2 3 4 5 6 7 8 9 10 ...
## $ Any : int 2003 2003 2003 2003 2003 2003 2003 2003 2003 2003 ...
## $ Codi_Districte : int 1 1 1 1 1 1 1 1 1 1 ...
## $ Nom_Districte : Factor w/ 10 levels "Ciutat Vella",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ Codi_Barri : int 1 1 1 1 1 1 1 1 1 1 ...
## $ Nom_Barri : Factor w/ 73 levels "Baró de Viver",...: 23 23 23 23 23 23 23 23 23 23 ...
## $ Valor : int 398 376 528 48 559 107 646 98 611 71 ...
## $ Nivell_Educatiu: Factor w/ 5 levels "Sense estudis",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ Edat : int 3 3 4 4 5 5 6 6 7 7 ...
## $ Sexe : int 0 1 0 1 0 1 0 1 0 1 ...
```

Després d'haver solucionat els errors de format i duplicats, a continuació es mirarà si hi ha valors atípics (*outliers*) o nuls (*missing*).

```
# Realitzem una estadística descriptiva
```

```
summary(bcn_educacio)
```

```
##      id      Any      Codi_Districte      Nom_Districte
## Min.   : 1      Min.   :2003      Min.   : 1.000      Nou Barris      :6235
## 1st Qu.: 9388      1st Qu.:2003      1st Qu.: 3.000      Horta-Guinardó:5482
## Median :18775      Median :2013      Median : 7.000      Sant Martí     :5184
## Mean   :18775      Mean   :2013      Mean   : 6.159      Sants-Montjuïc:4143
## 3rd Qu.:28161      3rd Qu.:2023      3rd Qu.: 8.000      Sant Andreu    :3500
## Max.   :37548      Max.   :2023      Max.   :10.000      Eixample       :3377
##                                     (Other)      :9627
##      Codi_Barri      Nom_Barri      Valor
## Min.   : 1.00      la Nova Esquerra de l'Eixample : 579      Min.   : 1.0
## 1st Qu.:18.00      el Raval                        : 573      1st Qu.: 38.0
## Median :36.00      les Corts                      : 566      Median :168.0
## Mean   :36.35      Sant Gervasi - la Bonanova     : 564      Mean   :250.5
## 3rd Qu.:55.00      l'Antiga Esquerra de l'Eixample: 563      3rd Qu.:450.0
## Max.   :73.00      la Sagrada Família            : 563      Max.   :992.0
##                                     (Other)      :34140
##      Nivell_Educatiu      Edat      Sexe
## Sense estudis           :6185      Min.   : 3.00      Min.   :0.0000
## Estudis primaris o EGB   :7578      1st Qu.: 7.00      1st Qu.:0.0000
## ESO o CFGB               :7253      Median :11.00      Median :0.0000
## Batxillerat, BUP, COU o CFGM:7114      Mean   :10.88      Mean   :0.4924
```

```
## Estudis universitaris o CFGS:6816    3rd Qu.:15.00    3rd Qu.:1.0000
## NA's                               :2602    Max.      :20.00    Max.      :1.0000
##
```

No s'observen valors atípics que no concorden amb la resta de les dades. Sembla, per tant, que és un joc de dades net. El que sí que hi ha són dades no registrades (NA) en la variable Nivell_Educatiu.

```
# Eliminem les dades amb valors nuls
bcn_educacio <- bcn_educacio[!is.na(bcn_educacio$Nivell_Educatiu),]

# Mostrem l'estructura del conjunt de dades
str(bcn_educacio)
```

```
## 'data.frame':    34946 obs. of  10 variables:
## $ id             : int  1 2 3 4 5 6 7 8 9 10 ...
## $ Any            : int  2003 2003 2003 2003 2003 2003 2003 2003 2003 2003 ...
## $ Codi_Districte : int  1 1 1 1 1 1 1 1 1 1 ...
## $ Nom_Districte  : Factor w/ 10 levels "Ciutat Vella",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ Codi_Barri     : int  1 1 1 1 1 1 1 1 1 1 ...
## $ Nom_Barri      : Factor w/ 73 levels "Baró de Viver",...: 23 23 23 23 23 23 23 23 23 23 ...
## $ Valor          : int  398 376 528 48 559 107 646 98 611 71 ...
## $ Nivell_Educatiu: Factor w/ 5 levels "Sense estudis",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ Edat           : int  3 3 4 4 5 5 6 6 7 7 ...
## $ Sexe           : int  0 1 0 1 0 1 0 1 0 1 ...
```

Ara tenim **34946 observacions**. Com s'ha dit abans, es vol enriquir el joc de dades. Per això, es barrejaran les variables existents per generar nous indicadors. S'estudia les proporcions, sobre 100, de persones per cada combinació de districte i nivell educatiu segons el sexe.

```
# Creem la taula de contingència
taula_contingencia <- table(bcn_educacio$Nom_Districte, bcn_educacio$Nivell_Educatiu, bcn_educacio$Sexe)

# Calculem les proporcions per files
proporcions_per_districte <- prop.table(taula_contingencia, margin = 1) * 100
proporcions_per_districte
```

```
## , , = 0
##
##
##          Sense estudis Estudis primaris o EGB ESO o CFGB
## Ciutat Vella          9.616343          10.662681  10.413553
## Eixample              8.892508          10.553746  10.456026
## Gràcia                8.709809          10.826211  10.582011
## Horta-Guinardó        8.758272          11.132736  10.568315
## Les Corts             8.423181          10.849057  10.646900
## Nou Barris            9.719149          11.353191  10.553191
## Sant Andreu           9.522370          11.094317  10.429262
## Sant Martí            8.957081          11.071947  10.470661
## Sants-Montjuïc        9.796238          10.971787  10.370951
## Sarrià-Sant Gervasi   8.288227          10.791610  10.723951
##
##          Batxillerat, BUP, COU o CFGM Estudis universitaris o CFGS
## Ciutat Vella          10.164425          10.014948
```

##	Eixample	10.325733	10.293160
##	Gràcia	10.419210	10.215710
##	Horta-Guinardó	10.198521	10.140132
##	Les Corts	10.579515	10.309973
##	Nou Barris	9.770213	9.174468
##	Sant Andreu	9.975816	9.340992
##	Sant Martí	10.221854	10.035248
##	Sants-Montjuïc	9.822362	9.587252
##	Sarrià-Sant Gervasi	10.690122	10.588633
##			
##	, , = 1		
##			
##			
##		Sense estudis	Estudis primaris o EGB
##	Ciutat Vella	9.018435	10.264076
##	Eixample	8.403909	10.390879
##	Gràcia	7.977208	10.500611
##	Horta-Guinardó	8.213313	10.918645
##	Les Corts	8.018868	10.309973
##	Nou Barris	8.987234	11.012766
##	Sant Andreu	9.159613	10.852479
##	Sant Martí	8.480199	10.574331
##	Sants-Montjuïc	9.273772	10.762800
##	Sarrià-Sant Gervasi	7.476319	10.453315
##			
##		Batxillerat, BUP, COU o CFGM	Estudis universitaris o CFGS
##	Ciutat Vella	10.064773	9.865471
##	Eixample	10.423453	10.065147
##	Gràcia	10.378510	10.012210
##	Horta-Guinardó	10.198521	9.692487
##	Les Corts	10.309973	10.377358
##	Nou Barris	10.110638	8.919149
##	Sant Andreu	10.157195	9.250302
##	Sant Martí	10.201120	9.724238
##	Sants-Montjuïc	10.005225	9.221526
##	Sarrià-Sant Gervasi	10.419486	10.317997

En resum, per aportar originalitat al conjunt de dades, s'ha considerat la creació i exploració de les proporcions d'habitants per cada districte segons el sexe i el nivell educatiu.

Es veuen unes lleugeres diferències entre el tipus d'estudis per cada sexe i districte. Pel cas de les dones (0), la menor proporció sense estudis és al Districte Sarrià-Sant Gervasi (8.29), mentre que la major és a Sants-Montjuïc (9.80). Quant a la proporció més gran amb estudis universitaris o CFGS, torna a ser a Sarrià-Sant Gervasi (10.59), i la menor és a Nou Barris (9.17).

Pel que fa als homes (1), la menor proporció sense estudis és al Districte Sarrià-Sant Gervasi (7.48), mentre que la major és a Sants-Montjuïc (9.27). Quant a la proporció més gran amb estudis universitaris o CFGS, és les Corts (10.38), i la menor és a Nou Barris (8.92).

Amb tot això es pot evidenciar que el districte Sarrià-Sant Gervasi, amb alt valor adquisitiu, presenta el nombre més gros d'individus amb estudis superiors, sent les dones les que tenen un major percentatge en estudis universitaris o CFGS per qualsevol dels deu districtes, exceptuant les Corts.

En la pròxima part de la pràctica, es realitzaran les primeres visualitzacions d'aquestes dades, ja que mai s'ha realitzat un estudi. S'ha pogut enriquir el conjunt de dades ja existent a partir d'una actualització i ajustament de les variables.

5 Diccionari

El plantejament de les qüestions que es respondran en la visualització de la segona part tindran en compte els punts anteriors per assegurar una perspectiva completa i significativa de les dades. Això implica assegurar que les qüestions són adequades i donen coneixement valuós.

Així doncs, es farà una investigació exhaustiva del nombre d'habitants per cada districte segons el nivell d'estudis i el sexe. A més, es compararà segons siguin registres del 2003, 2013 o 2023, amb l'objectiu de veure un canvi significatiu en les conductes de la població al llarg del temps.

Com el conjunt de dades no ha estat plantejat en altres visualitzacions o projectes, hi ha llibertat total per anar provant diferents tipus de gràfics sense cap mena de problema per repetir estudis.

```
# Guardem el dataset definitiu en un nou fitxer CSV
write.csv(bcn_educacio, 'bcn_educacio_v2.csv', row.names = FALSE)
```

- Variables quantitatives: `id`, `Any`, `Codi_Districte`, `Codi_Barri`, `Valor`, `Edat` i `Sexe`.
- Variables qualitatives: `Nom_Districte`, `Nom_Barri` i `Nivell_Educatiu`.

Creem un diccionari de dades:

- `id`: identificador únic per a cada registre.

FETS A ESTUDIAR

- `Valor`: nombre de persones.
- `Nivell_Educatiu`: nivell educatiu (5 nivells).
- `Edat`: edat (17 nivells).
 - 3: 15-19 anys.
 - 4: 20-24 anys.
 - ...
 - 19: 95-99 anys.
 - 20: > 100 anys.
- `Sexe`: sexe (2 nivells).
 - 0: Dona.
 - 1: Home.

DIMENSIÓ GEOGRÀFICA

- `Codi_Districte`: codi del districte (de l'1 al 10).
- `Nom_Districte`: nom del districte (10 nivells).
- `Codi_Barri`: codi del barri (de l'1 al 73).
- `Nom_Barri`: nom del barri (73 nivells).

DIMENSIÓ TEMPORAL

- `Any`: any (2003, 2013 o 2023).