

PRA2. Tipologia i cicle de vida de les dades

Autor: Daniel Rodríguez Morente

Juny de 2023

Contents

1	Presentació del projecte i objectiu de l'anàlisi	1
2	Descripció del dataset	2
2.1	Perquè és important i quina pregunta/problema pretén respondre?	2
3	Integració i selecció de les dades d'interès a analitzar. Pot ser el resultat d'addicionar diferents datasets o una subselecció útil de les dades originals, en base a l'objectiu que es vulgui aconseguir.	3
4	Neteja de les dades. Les dades contenen zeros o elements buits?	5
4.1	Les dades contenen zeros o elements buits? Gestiona cadascun d'aquests casos	5
4.2	Identifica i gestiona els valors extrems	6
5	Anàlisi de les dades	7
5.1	Selecció dels grups de dades que es volen analitzar/comparar (p. e., si es volen comparar grups de dades, quins són aquests grups i quins tipus d'anàlisi s'aplicaran?).	7
5.2	Comprovació de la normalitat i homogeneïtat de la variància.	9
5.3	Aplicació de proves estadístiques per comparar els grups de dades. En funció de les dades i de l'objectiu de l'estudi, aplicar proves de contrast d'hipòtesis, correlacions, regressions, etc. Aplicar almenys tres mètodes d'anàlisi diferents.	11
6	Conclusions	16
6.1	A partir dels resultats obtinguts, quines són les conclusions?	16
6.2	Els resultats permeten respondre al problema?	16
7	Taula de contribucions	16

1 Presentació del projecte i objectiu de l'anàlisi

El projecte que es desenvolupa a continuació consisteix en l'estudi de les causes que determinen la probabilitat de patir una cardiopatia, per tal de identificar si posteriors estudis han de tractar de forma separada a homes i dones.

2 Descripció del dataset

2.1 Perquè és important i quina pregunta/problema pretén respondre?

El dataset utilitzat conté informació de diferents indicadors mèdics de persones que han patit o no una cardiopatia.

Les dades han estat publicades per Rashik Rahman sota llicència CC0: Public Domain a [www.kaggle.com](https://www.kaggle.com/rashikrahmanpritom/heart-attack-analysis-prediction-dataset) i es pot accedir a les mateixes a través del següent enllaç: <https://www.kaggle.com/rashikrahmanpritom/heart-attack-analysis-prediction-dataset>

No s'han realitzat modificacions prèvies al conjunt de dades original.

Carreguem el conjunt de dades i fem una revisió del contingut de les diferents variables

```
## 'data.frame':    303 obs. of  14 variables:
## $ age      : int  63 37 41 56 57 57 56 44 52 57 ...
## $ sex      : int  1 1 0 1 0 1 0 1 1 1 ...
## $ cp       : int  3 2 1 1 0 0 1 1 2 2 ...
## $ trtbps   : int  145 130 130 120 120 140 140 120 172 150 ...
## $ chol     : int  233 250 204 236 354 192 294 263 199 168 ...
## $ fbs      : int  1 0 0 0 0 0 0 0 1 0 ...
## $ restecg  : int  0 1 0 1 1 1 0 1 1 1 ...
## $ thalachh : int  150 187 172 178 163 148 153 173 162 174 ...
## $ exng     : int  0 0 0 0 1 0 0 0 0 0 ...
## $ oldpeak  : num  2.3 3.5 1.4 0.8 0.6 0.4 1.3 0 0.5 1.6 ...
## $ slp      : int  0 0 2 2 2 1 1 2 2 2 ...
## $ caa      : int  0 0 0 0 0 0 0 0 0 0 ...
## $ thall    : int  1 2 2 2 2 1 2 3 3 2 ...
## $ output   : int  1 1 1 1 1 1 1 1 1 1 ...
```

Podem observar que es tracta d'un dataset amb 303 observacions i 14 variables, totes elles amb números enters excepte la variable `oldpeak` que conté dades decimals.

Descripció de les variables

- **age**. Edat de la persona
- **sex**. Sexe de la persona (1 = home; 0 = dona)
- **cp**. Mal al tòrax (1: angina de pit amb valor típic; 2: angina de pit amb valor atípic 3: dolor no anginos; 4: assintomàtic)
- **trtbps**. Pressió arterial en repòs (en mm/Hg)
- **chol**. Nivell de colesterol mesurat (en mg/dl)
- **fbs**. Nivell de sucre en sang en dejú (fasting blood sugar > 120 mg/dl) (1: > 120 mg/dl; 0: =< 120 mg/dl)
- **restecg**. resultat de l'electrocardiograma en repòs (0 = normal; 1 = anomalia de l'ona ST-T (inversions de l'ona T i/o elevació o depressió del ST de > 0.05 mV); 2: hipertròfia ventricular esquerra probable o hipertròfia ventricular esquerra segons els criteris d'Estes)
- **thalachh**. Freqüència cardíaca màxima assolida
- **exng**. Angina induïda per exercici físic (1 = sí; 0 = no)
- **oldpeak**. Pics previs
- **slp**. Pendent (0 = sense pendent; 1 = pla; 2 = pendent descendent)

- **caa**. Número de vasos sanguinis principals amb obstrucció (0-4)
- **thall**. Resultats d'una prova d'esforç amb tali (0 = null; 1 = fixed defect; 2 = normal; 3 = reversible defect)
- **output**. Variable objectiu (0 = menys possibilitats de patir una cardiopatia; 1= més possibilitats de patir una cardiopatia)

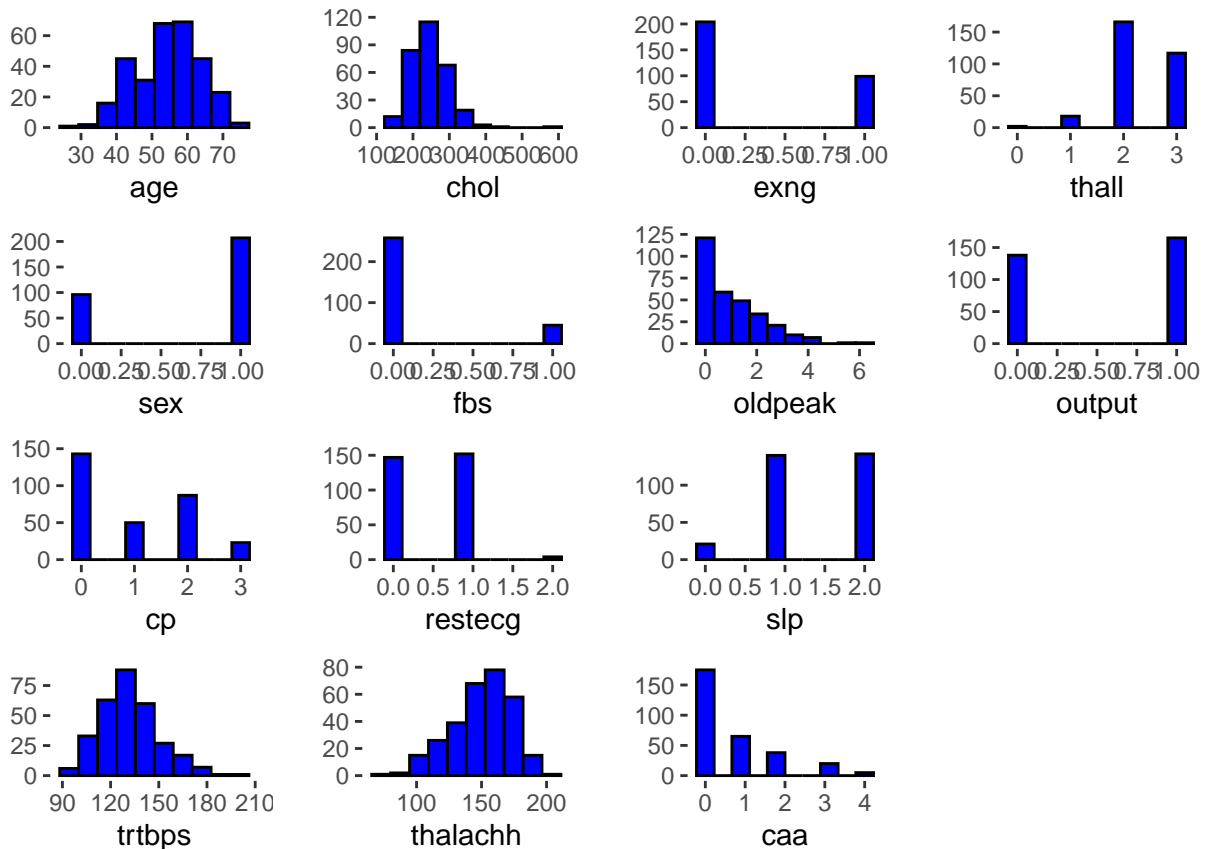
3 Integració i selecció de les dades d'interès a analitzar. Pot ser el resultat d'addicionar diferents datasets o una subselecció útil de les dades originals, en base a l'objectiu que es vulgui aconseguir.

Revisem la distribució de les diferents variables

```
##      age      sex      cp      trtbps
## Min.   :29.00  Min.   :0.0000  Min.   :0.000  Min.   : 94.0
## 1st Qu.:47.50  1st Qu.:0.0000  1st Qu.:0.000  1st Qu.:120.0
## Median :55.00  Median :1.0000  Median :1.000  Median :130.0
## Mean   :54.37  Mean   :0.6832  Mean   :0.967  Mean   :131.6
## 3rd Qu.:61.00  3rd Qu.:1.0000  3rd Qu.:2.000  3rd Qu.:140.0
## Max.   :77.00  Max.   :1.0000  Max.   :3.000  Max.   :200.0
##      chol      fbs      restecg      thalachh
## Min.   :126.0  Min.   :0.0000  Min.   :0.0000  Min.   : 71.0
## 1st Qu.:211.0  1st Qu.:0.0000  1st Qu.:0.0000  1st Qu.:133.5
## Median :240.0  Median :0.0000  Median :1.0000  Median :153.0
## Mean   :246.3  Mean   :0.1485  Mean   :0.5281  Mean   :149.6
## 3rd Qu.:274.5  3rd Qu.:0.0000  3rd Qu.:1.0000  3rd Qu.:166.0
## Max.   :564.0  Max.   :1.0000  Max.   :2.0000  Max.   :202.0
##      exng      oldpeak      slp      caa
## Min.   :0.0000  Min.   :0.00  Min.   :0.000  Min.   :0.0000
## 1st Qu.:0.0000  1st Qu.:0.00  1st Qu.:1.000  1st Qu.:0.0000
## Median :0.0000  Median :0.80  Median :1.000  Median :0.0000
## Mean   :0.3267  Mean   :1.04  Mean   :1.399  Mean   :0.7294
## 3rd Qu.:1.0000  3rd Qu.:1.60  3rd Qu.:2.000  3rd Qu.:1.0000
## Max.   :1.0000  Max.   :6.20  Max.   :2.000  Max.   :4.0000
##      thall      output
## Min.   :0.000  Min.   :0.0000
## 1st Qu.:2.000  1st Qu.:0.0000
## Median :2.000  Median :1.0000
## Mean   :2.314  Mean   :0.5446
## 3rd Qu.:3.000  3rd Qu.:1.0000
## Max.   :3.000  Max.   :1.0000
```

A priori, no observem que hi hagi valor perduts, però més endavant farem una comprovació adhoc.

Fem una representació de les diferents variables per tal de facilitar la revisió prèvia del dataset



Hi ha més informació d'homes que de dones i el número de registres amb output igual a 1 és lleugerament superior al valor 0.

Observant les gràfiques, veiem que hi ha quatre variables que podrien tenir una distribució similar a una normal (age, chol, trtbps i thalachh), la qual cosa revisarem en un punt posterior de l'estudi.

Modifiquem els valors de la variable sex per facilitar la seva interpretació

Revisem si tenim registres amb idèntics valors a totes les variables per tal de identificar si tenim registres duplicats

```
## [1] 302 14
```

Comprovem que hi ha 302 registres diferents, per la qual cosa, donat el nivell d'especificitat de les dades, considerem que hi ha un registre repetit.

Eliminem el registre repetit i conservem la resta donat que tenim un número de registres perfectament gestionable i, per tant no és necessari plantejar agrupacions que facilitin l'ús del dataset

Donat que l'estudi està dirigit a identificar diferències entre homes i dones, ens interessa comprovar quina informació tenim per cada grup

```
## [1] "Distribució entre homes i dones en valors absoluts:"
```

```
##
##      0    1 Sum
## Dona  24   72  96
## Home 114   92 206
## Sum  138  164 302
```

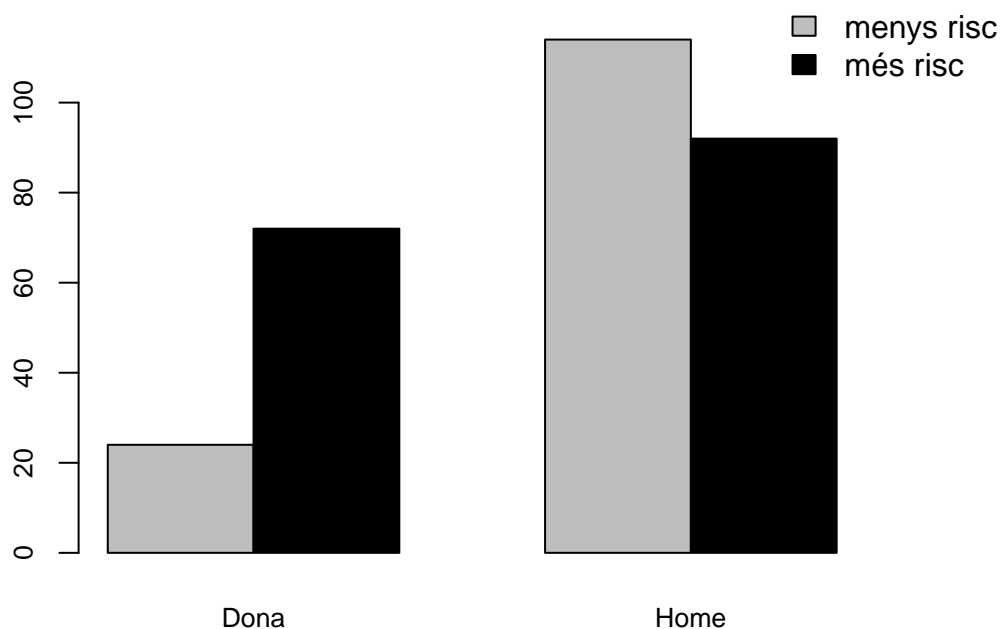
```
## [1] "Pes relatiu de cada sexe dins el valor de la variable output:"
```

```
##  
##           0    1  
## Dona 0.17 0.44  
## Home 0.83 0.56
```

```
## [1] "Pes relatiu de la variable output dins de cada sexe:"
```

```
##  
##           0    1  
## Dona 0.25 0.75  
## Home 0.55 0.45
```

Hi ha 207 homes i 96 dones i la distribució del camp output dins de cada grup és diferent, tenint més pes el valor 1 en dones que en homes.



4 Neteja de les dades. Les dades contenen zeros o elements buits?

4.1 Les dades contenen zeros o elements buits? Gestiona cadascun d'aquests casos

Tot i que amb al resum del dataset no apareixien valors perduts, fem una revisió adhoc de valors perduts i en blanc

```
##      age      sex      cp      trtbps      chol      fbs      restecg      thalachh
##      0        0        0        0        0        0        0        0
##      exng    oldpeak    slp      caa      thall      output
##      0        0        0        0        0        0

##      age      sex      cp      trtbps      chol      fbs      restecg      thalachh
##      0        0        0        0        0        0        0        0
##      exng    oldpeak    slp      caa      thall      output
##      0        0        0        0        0        0
```

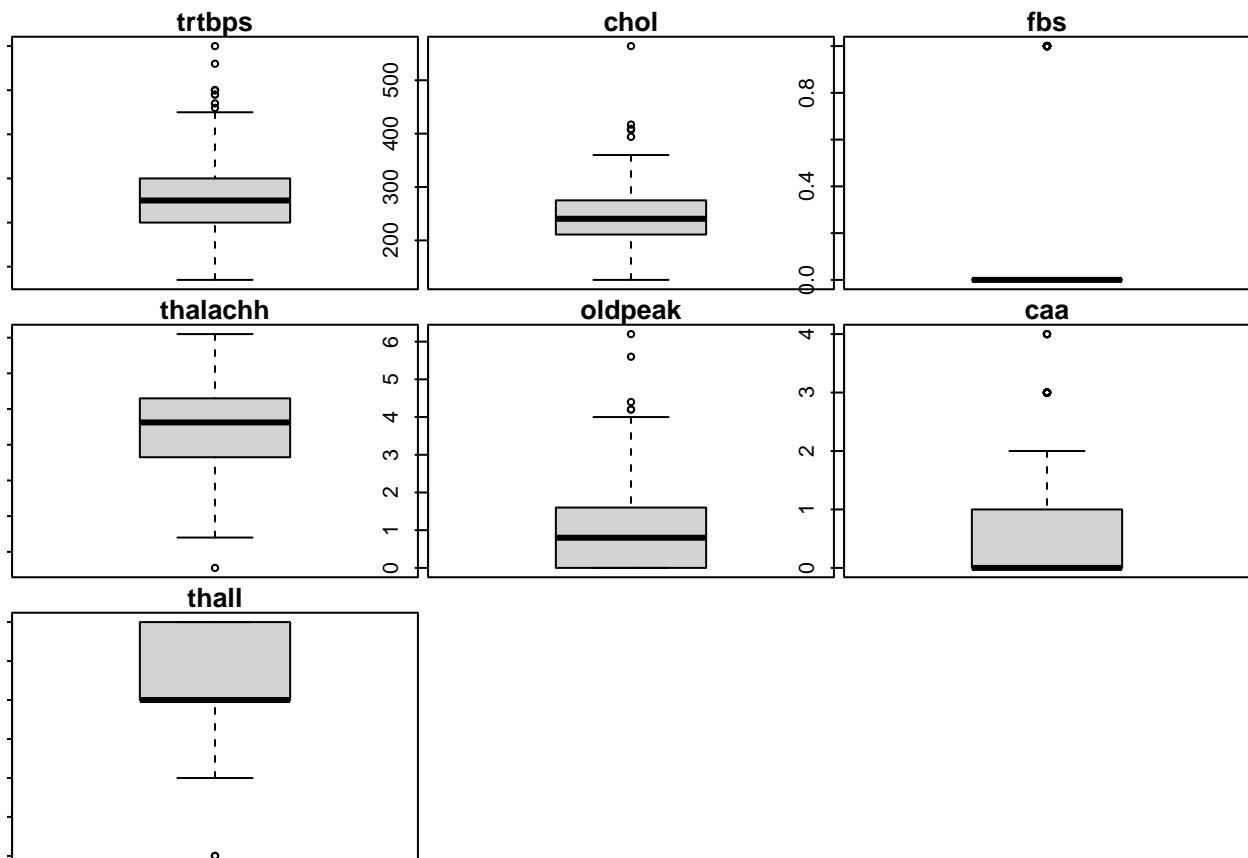
Es confirma que no tenim valors perduts, per tant no em de fer cap modificació al dataset.

4.2 Identifica i gestiona els valors extrems

Fem una revisió de la possible existència de valors extrems en el conjunt de les dades

```
## [1] "trtbps" "chol" "fbs" "thalachh" "oldpeak" "caa" "thall"
```

Veiem que tenim 7 variables que a priori contenen valors extrems



Per les variables fbs, caa i thall, tot i el resultat gràfic, podem descartar el fet que hi hagi valors extrems donat que es tracta de variables discretes i els valors observats estan dins de les categories considerades.

Apliquem el criteri de les dues desviacions estàndard per tal de identificar si mantenim els valors originals

```
## [1] "Número d'outliers considerant les dades de forma global: 49"
```

Donat que l'estudi es basa en la comparació entre homes i dones, fem una valoració dels valors extrems per separat, per tal d'evitar que els valors d'un sexe amaguin informació rellevant a l'altra

```
## [1] "trtbps"      "chol"      "fbs"      "thalachh" "exng"      "oldpeak"   "caa"
## [8] "thall"

## [1] "trtbps"      "fbs"      "thalachh" "oldpeak"   "caa"      "thall"
```

Així com en l'estudi conjunt trobavem valors extrems a les variables trtbps, chol, thalachh i oldpeak, quan estudiem els sexes per separat varien aquests resultats. Per les dones no hi ha variació en quant a variables, donat que la variable exng en ser discreta no la podem considerar, i pels homes no hi haurà valors extrems per la variable chol.

Revisem els registres que contenen valors extrems considerant dues desviacions estàndar i comptem quants registres es veuen afectats per outliers, considerant les dades de forma global i separant per sexes

```
## [1] "Número d'outliers considerant les dades de forma global: 49"

## [1] "Número d'outliers en dones considerant les dades de forma global: 19"

## [1] "Número d'outliers en homes considerant les dades de forma global: 30"

## [1] "Número d'outliers en dones: 17"

## [1] "Número d'outliers en homes: 25"
```

Podem extreure una primera conclusió sobre la importància de tractar les dades per separat donat que es redueix el número d'outliers. Tractar les dades conjuntament implicaria fer un tractament de les dades errònia i descartar registres o imputar valors de forma equivocada, a banda de que ens facilita una primera informació al respecte dels diferents valors observats en funció de si es tracta de dones o d'homes.

Per tal de decidir si realment els valors trobats són erronis cal tenir un coneixement ampli del tipus de dades i de si els valors que estem identificant són relament erronis. Per altra banda, amb una mostra de 302 registres si 41 tenen dades errònies, hauríem de considerar que hi ha hagut massa errors en la recolecció de les dades i la mostra no és gaire útil. Per tant, per continuar amb l'estudi, considerarem que els valors detectats són correctes i continuarem l'anàlisi sense imputar nous valors, tenint present que aquest és un exercici teòric i que en un cas real hauríem de consultar amb els experts per tal de validar quin és el tractament correcte.

5 Anàlisi de les dades

5.1 Selecció dels grups de dades que es volen analitzar/comparar (p. e., si es volen comparar grups de dades, quins són aquests grups i quins tipus d'anàlisi s'aplicaran?).

Com ja s'ha comentat, l'objectiu de l'estudi és la comparació de les dades existents entre homes i dones per tal de valorar si posteriors estudis han de considerar les dades per separat.

Ja hem vist que el tractament diferenciat ens permet detectar valors extrems diferents en el cas d'homes i de dones.

Tot i que es tracta d'una tècnica utilitzada per reduir la dimensionalitat de les dades, farem una anàlisi PCA per tal de valorar si les variables tenen la mateixa importància pels dos grups i d'aquesta manera refermar la idea de fer un tractament diferenciat

```
## [1] "ACP dones:"
```

```
## Importance of components:
```

```
##          PC1      PC2      PC3      PC4      PC5      PC6      PC7
## Standard deviation    1.9353 1.2488 1.1511 1.05246 1.00397 0.95263 0.91461
## Proportion of Variance 0.2881 0.1200 0.1019 0.08521 0.07753 0.06981 0.06435
## Cumulative Proportion 0.2881 0.4081 0.5100 0.59522 0.67276 0.74257 0.80691
##          PC8      PC9      PC10     PC11     PC12     PC13
## Standard deviation    0.87141 0.69795 0.63454 0.55813 0.54281 0.5048
## Proportion of Variance 0.05841 0.03747 0.03097 0.02396 0.02266 0.0196
## Cumulative Proportion 0.86532 0.90280 0.93377 0.95773 0.98040 1.0000
```

```
## [1] "ACP homes:"
```

```
## Importance of components:
```

```
##          PC1      PC2      PC3      PC4      PC5      PC6      PC7
## Standard deviation    1.7966 1.2412 1.12198 1.0268 0.99126 0.94694 0.91272
## Proportion of Variance 0.2483 0.1185 0.09683 0.0811 0.07558 0.06898 0.06408
## Cumulative Proportion 0.2483 0.3668 0.46364 0.5447 0.62032 0.68929 0.75338
##          PC8      PC9      PC10     PC11     PC12     PC13
## Standard deviation    0.87956 0.80828 0.76272 0.67198 0.64432 0.57508
## Proportion of Variance 0.05951 0.05026 0.04475 0.03474 0.03193 0.02544
## Cumulative Proportion 0.81289 0.86314 0.90789 0.94263 0.97456 1.00000
```

S'aprecia una diferència entre el resultat pels homes i per les dones, sent força significativa la diferència de la primera component principal, amb un valor de quasi un 4% més per les dones.

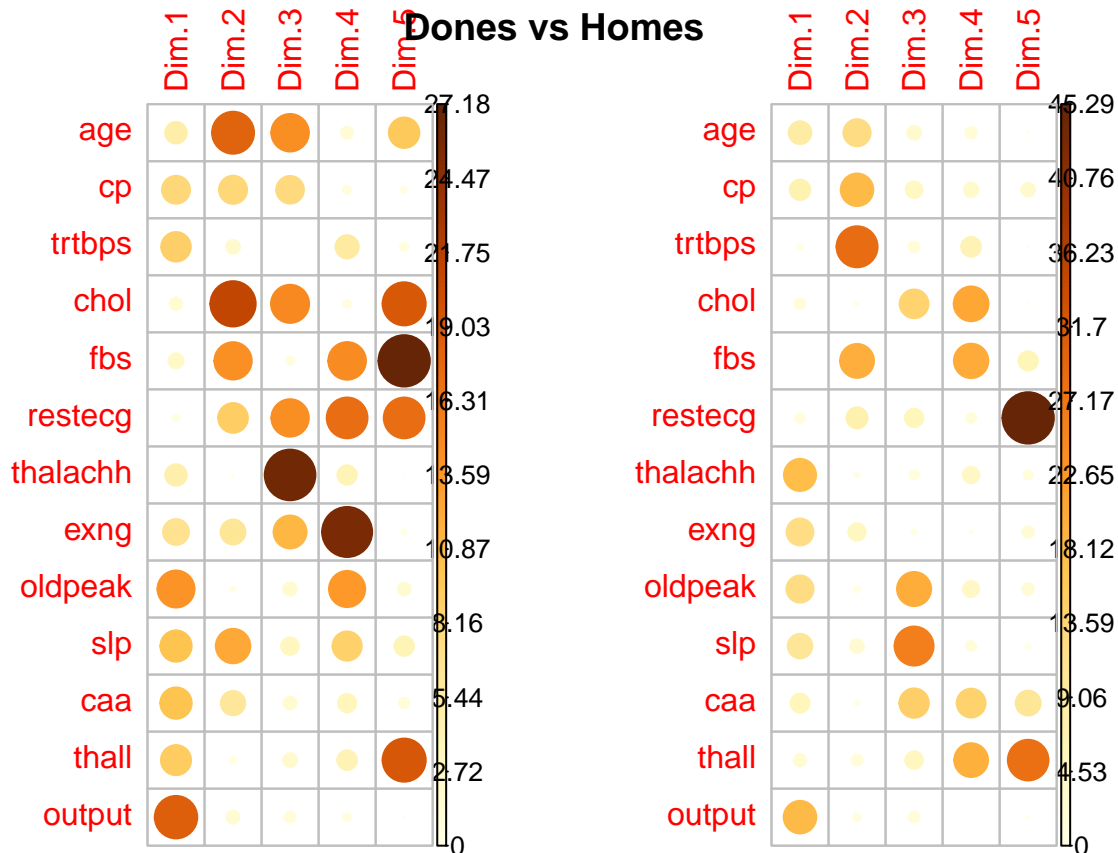
En quant a la interpretació del resultat, per les dones les dues primeres components principals expliquen el 40,81% de la variància, mentre que pels homes acumulen 36,68%, però la participació de les diferents components principals està força repartida, necessitant fins 9 per les dones i 10 pels homes per arribar al 90%.

Com a criteri de selecció considerarem les components amb una variància superior a 1

```
## [1] 3.7455634 1.5595211 1.3251408 1.1076699 1.0079545 0.9075065 0.8365109
## [8] 0.7593494 0.4871390 0.4026351 0.3115098 0.2946376 0.2548619
```

```
## [1] 3.2278984 1.5405512 1.2588347 1.0542459 0.9826044 0.8966912 0.8330578
## [8] 0.7736301 0.6533156 0.5817416 0.4515623 0.4151547 0.3307122
```

Per tant, considerarem les 5 primeres per les dones i les 4 primeres pels homes. De totes maneres, revisarem la importància de cada variable a les 5 primeres components principals diferenciant per sexes



De forma visual es pot apreciar que hi ha diferències en quant a la importància de cada variable en la contribució a les 5 primeres components principals, per tant podem considerar que a l'hora de seleccionar les variables a estudiar serà important diferenciar entre els dos sexes.

5.2 Comprovació de la normalitat i homogeneïtat de la variància.

Tenim cinc variables contínues: age, trtbps, chol, thalachh i oldpeak a les quals aplicarem el test de Shapiro per comprovar la seva **normalitat** a les dades en conjunt

```
##
## Shapiro-Wilk normality test
##
## data:  dades$age
## W = 0.98664, p-value = 0.006745

##
## Shapiro-Wilk normality test
##
## data:  dades$trtbps
## W = 0.96573, p-value = 1.419e-06

##
## Shapiro-Wilk normality test
##
## data:  dades$chol
## W = 0.94658, p-value = 5.196e-09
```

```
##
## Shapiro-Wilk normality test
##
## data:  dades$thalachh
## W = 0.97679, p-value = 8.268e-05
```

```
##
## Shapiro-Wilk normality test
##
## data:  dades$oldpeak
## W = 0.84522, p-value < 2.2e-16
```

En els cinc casos, observant el valor de p, podem dir que, segons el test Shapiro, es rebutja la hipòtesi nul · la i, per tant, no es distribueixen com una normal.

Apliquem el test de Kolmogorov_Smirnov per tal de valorar si els resultats coincideixen

```
##
## One-sample Kolmogorov-Smirnov test
##
## data:  dades$age
## D = 0.075788, p-value = 0.06228
## alternative hypothesis: two-sided
```

```
##
## One-sample Kolmogorov-Smirnov test
##
## data:  dades$trtbps
## D = 0.10258, p-value = 0.003475
## alternative hypothesis: two-sided
```

```
##
## One-sample Kolmogorov-Smirnov test
##
## data:  dades$chol
## D = 0.055822, p-value = 0.3035
## alternative hypothesis: two-sided
```

```
##
## One-sample Kolmogorov-Smirnov test
##
## data:  dades$thalachh
## D = 0.070819, p-value = 0.09669
## alternative hypothesis: two-sided
```

```
##
## One-sample Kolmogorov-Smirnov test
##
## data:  dades$oldpeak
## D = 0.18458, p-value = 2.313e-09
## alternative hypothesis: two-sided
```

Obtenim resultats contradictoris per age, chol i thalachh, per tant, sent conservadors, considerarem que cap de les cinc variables es distribueix segons una normal

Comprovem l'**homoscedasticitat** per les quatre variables

Aplicuem el test de fligner sobre les cinc variables, al considerar que no segueixen una distribució normal

```
##
## Fligner-Killeen test of homogeneity of variances
##
## data:  trtbps by sex
## Fligner-Killeen:med chi-squared = 0.93812, df = 1, p-value = 0.3328

##
## Fligner-Killeen test of homogeneity of variances
##
## data:  oldpeak by sex
## Fligner-Killeen:med chi-squared = 8.372, df = 1, p-value = 0.00381

##
## Fligner-Killeen test of homogeneity of variances
##
## data:  age by sex
## Fligner-Killeen:med chi-squared = 0.68171, df = 1, p-value = 0.409

##
## Fligner-Killeen test of homogeneity of variances
##
## data:  chol by sex
## Fligner-Killeen:med chi-squared = 9.2927, df = 1, p-value = 0.002301

##
## Fligner-Killeen test of homogeneity of variances
##
## data:  thalachh by sex
## Fligner-Killeen:med chi-squared = 5.3763, df = 1, p-value = 0.02041
```

Segons els resultats obtingut podem concloure que oldpeak, chol i thalachh tenen variàncies estadísticament diferents per cada sexe mentre que age i trtbps tenen variàncies estadísticament iguals.

5.3 Aplicació de proves estadístiques per comparar els grups de dades. En funció de les dades i de l'objectiu de l'estudi, aplicar proves de contrast d'hipòtesis, correlacions, regressions, etc. Aplicar almenys tres mètodes d'anàlisi diferents.

En base a l'estudi de components principals realitzats prèviament, considerarem les següents variables: oldpeak, age i chol, per la incidència que tenen en el grup de dones, i thalachh, trtbps i fbs, per la incidència que tenen en el grup d'homes. A més utilitzarem la variable output.

Donat que hem considerat que les variables contínues no segueixen una distribució normal, utilitzarem tests no paramètrics per aquestes variables.

Considerem que els grups dades separats per sexe són independents per la qual cosa aplicarem el test Mann-Whitney per comparar les distribucions

```
##
## Wilcoxon rank sum test with continuity correction
##
## data: oldpeak by sex
## W = 8638, p-value = 0.07187
## alternative hypothesis: true location shift is not equal to 0

##
## Wilcoxon rank sum test with continuity correction
##
## data: age by sex
## W = 11064, p-value = 0.09585
## alternative hypothesis: true location shift is not equal to 0

##
## Wilcoxon rank sum test with continuity correction
##
## data: chol by sex
## W = 11710, p-value = 0.009943
## alternative hypothesis: true location shift is not equal to 0

##
## Wilcoxon rank sum test with continuity correction
##
## data: thalachh by sex
## W = 10420, p-value = 0.4519
## alternative hypothesis: true location shift is not equal to 0

##
## Wilcoxon rank sum test with continuity correction
##
## data: trtbps by sex
## W = 10552, p-value = 0.3464
## alternative hypothesis: true location shift is not equal to 0
```

Donat que obtenim una p-valor inferior al 0,05 per chol podem concloure que aquesta variable té diferències estadísticament significatives per cada sexe, mentre que a les altres quatre variables contínues no.

Apliquem el test chi-quadrat a la variable discreta fbs

```
##
## Pearson's Chi-squared test with Yates' continuity correction
##
## data: table(dades[, c("sex", "fbs")])
## X-squared = 0.39221, df = 1, p-value = 0.5311
```

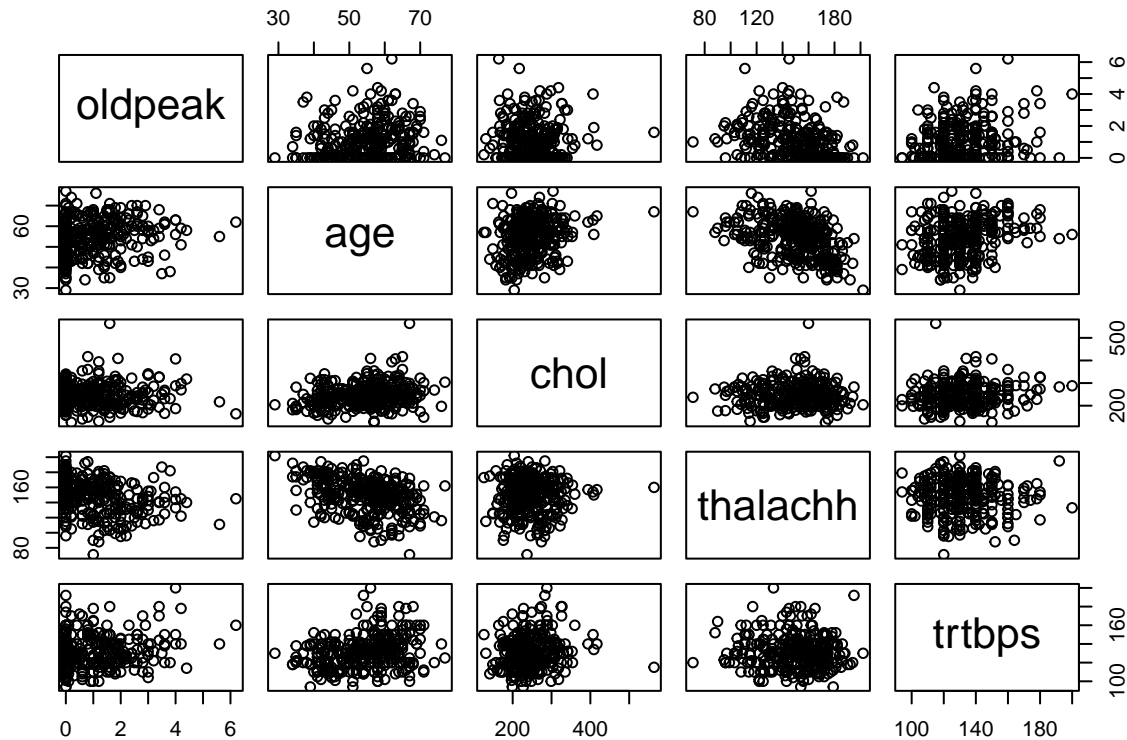
Podem concloure que no hi ha diferències significatives en la variable fbs per sexes.

Fem el mateix exercici amb la variable output

```
##
## Pearson's Chi-squared test with Yates' continuity correction
##
## data: table(dades[, c("sex", "output")])
## X-squared = 23.084, df = 1, p-value = 1.551e-06
```

En aquest cas sí s'aprecien diferències significatives, la qual cosa és especialment rellevant pel nostre estudi, sempre i quan considerem que la mostra és suficientment representativa de la població.

Observem les relacions entre les diferents variables contínues per tal valorar visualment si hi pot haver una correlació entre elles que ens permeti generar un model de regressió lineal.



A priori, sembla que hi podria haver una certa relació entre age i chol i també entre age i thalachh.

Estudiem el grau de correlació entre les variables age i chol

```
##
## Call:
## lm(formula = age ~ chol, data = dades)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -23.8809  -6.5405   0.4127   6.3624  23.3727
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 45.490546   2.486992  18.291  < 2e-16 ***
## chol         0.036227   0.009875   3.669 0.000288 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.866 on 300 degrees of freedom
## Multiple R-squared:  0.04294,    Adjusted R-squared:  0.03975
## F-statistic: 13.46 on 1 and 300 DF,  p-value: 0.0002884
```

Fem el mateix estudi per les variables age i thalachh

```
##
## Call:
## lm(formula = age ~ thalachh, data = dades)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -22.4760  -6.6945   0.5537   6.3865  24.5203
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  77.77379     3.17004   24.534 < 2e-16 ***
## thalachh    -0.15614     0.02095   -7.452 9.86e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8.325 on 300 degrees of freedom
## Multiple R-squared:  0.1562, Adjusted R-squared:  0.1534
## F-statistic: 55.54 on 1 and 300 DF,  p-value: 9.858e-13
```

Com podem observar en els valors obtingut pel coeficient de determinació, la correlació entre els dos parells de variables és feble.

Donat que hem considerat sis variables com a més significatives per l'estudi, calcularem la correlació existent entre les cinc que són contínues i output per valorar el nivell de correlació existent

```
##              oldpeak      age      chol      thalachh      trtbps      output
## oldpeak    1.00000000  0.2636254  0.03956479 -0.43049461  0.15680732 -0.4196306
## age        0.26362540  1.00000000  0.18890292 -0.39345342  0.28970501 -0.2348453
## chol       0.03956479  0.1889029  1.00000000 -0.04036747  0.13021023 -0.1170065
## thalachh   -0.43049461 -0.3934534 -0.04036747  1.00000000 -0.04269948  0.4263680
## trtbps     0.15680732  0.2897050  0.13021023 -0.04269948  1.00000000 -0.1234777
## output    -0.41963058 -0.2348453 -0.11700649  0.42636798 -0.12347766  1.0000000
```

Podem veure als resultats que la correlació és baixa per oldpeak i thalachh, sent inversa per la primera. Pel cas d'age, chol i trtbps tenim una correlació molt baixa.

Fem ara el mateix càlcul diferenciant homes i dones

```
##              oldpeak      age      chol      thalachh      trtbps      output
## oldpeak    1.00000000  0.2209536  0.09984333 -0.3385933  0.2864445 -0.3864213
## age        0.22095360  1.00000000  0.23125909 -0.3707971  0.3538976 -0.2059031
## chol       0.09984333  0.2312591  1.00000000 -0.0748373  0.2090334 -0.1510627
## thalachh   -0.33859334 -0.3707971 -0.07483730  1.00000000 -0.1458489  0.2683794
## trtbps     0.28644453  0.3538976  0.20903340 -0.1458489  1.00000000 -0.3455143
## output    -0.38642132 -0.2059031 -0.15106267  0.2683794 -0.3455143  1.0000000

##              oldpeak      age      chol      thalachh      trtbps      output
## oldpeak    1.00000000  0.3041437  0.03695629 -0.46452683  0.10252637 -0.40835321
## age        0.30414368  1.00000000  0.14288329 -0.41083513  0.25639198 -0.29539569
## chol       0.03695629  0.1428833  1.00000000 -0.04345756  0.06853068 -0.17464040
## thalachh   -0.46452683 -0.4108351 -0.04345756  1.00000000 -0.01676926  0.49344417
## trtbps     0.10252637  0.2563920  0.06853068 -0.01676926  1.00000000 -0.05500951
## output    -0.40835321 -0.2953957 -0.17464040  0.49344417 -0.05500951  1.00000000
```

Podem veure com els resultats varien. En el cas de les dones, respecte a output, les variacions més significatives són la disminució de la correlació per thalachh i l'augment per trtbps, sense arribar a ser alta. Pel cas del homes, la correlació respecte a output per trtbps passa a ser quasi zero, mentre que puja la correlació de les altres variables.

Per analitzar si les correlacions són significatives, farem un test d'Spearman per oldpeak i trtbps, per les dones, i per oldpeak i thalachh pels homes

```
##
## Spearman's rank correlation rho
##
## data:  dades_dones$oldpeak and dades_dones$output
## S = 204414, p-value = 0.0001008
## alternative hypothesis: true rho is not equal to 0
## sample estimates:
##      rho
## -0.3864213

##
## Spearman's rank correlation rho
##
## data:  dades_dones$trtbps and dades_dones$output
## S = 198383, p-value = 0.0005651
## alternative hypothesis: true rho is not equal to 0
## sample estimates:
##      rho
## -0.3455143

##
## Spearman's rank correlation rho
##
## data:  dades_homes$oldpeak and dades_homes$output
## S = 2051879, p-value = 1.108e-09
## alternative hypothesis: true rho is not equal to 0
## sample estimates:
##      rho
## -0.4083532

##
## Spearman's rank correlation rho
##
## data:  dades_homes$thalachh and dades_homes$output
## S = 738019, p-value = 4.856e-14
## alternative hypothesis: true rho is not equal to 0
## sample estimates:
##      rho
## 0.4934442
```

Podem apreciar com en cap cas són les correlacions són significativament diferents de zero.

6 Conclusions

6.1 A partir dels resultats obtinguts, quines són les conclusions?

L'estudi de les diferents varibles permet determinar que és necessari diferenciar entre homes i dones per tal d'estudiar les probabilitats de patir cardiopaties donats que tant les distribucions de les variables com la importància de les mateixes per cada sexe és diferent i no fer-ho pot implicar descartar dades de forma errònia, com hem vist en l'estudi d'outliers.

Per altra banda, considerem que el consell d'un expert per tal de valorar si els valors identificats com a outliers ho són o no, pot suposar que el resultat obtingut sigui diferent i, per tant, modificar les conclusions.

6.2 Els resultats permeten respondre al problema?

Els resultats són determinants i responen clarament el problema plantejat.

7 Taula de contribucions

Contribucions	Signatura
Investigació prèvia	Daniel Rodríguez Morente
Redacció de les respostes	Daniel Rodríguez Morente
Desenvolupament	Daniel Rodríguez Morente
Participació al vídeo	Daniel Rodríguez Morente