

Séance 2.3: Exploration des données

Visseho Adjiwanou, PhD.

30 May 2021

Enquête sociale générale, 1996

- Il s'agit du CROP Socio-Cultural Survey de 1996
- Dans cette partie, nous allons apprendre à :
 - Sélectionner les **variables**
 - Sélectionner les observations
 - Réorganiser les données
 - Créer de nouvelles variables avec des fonctions de variables existantes (`mutate()`)
 - Recoder des variables existantes
 - Calculer des statistiques univariées

Dressons la table

```
# Effacer votre environnement

rm(list = ls())

# Installer les package dont vous avez besoin

#install.packages("tidyverse")
#install.packages("summarytools")
#install.packages("tinytex")

# Charger les packages - Étape fondamentales

library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.0 --
## v ggplot2 3.3.3      v purrr  0.3.4
## v tibble  3.1.2      v dplyr  1.0.6
## v tidyr   1.1.3      v stringr 1.4.0
## v readr   1.4.0      v forcats 0.4.0

## Warning: package 'ggplot2' was built under R version 3.6.2
## Warning: package 'tibble' was built under R version 3.6.2
## Warning: package 'tidyr' was built under R version 3.6.2
## Warning: package 'readr' was built under R version 3.6.2
## Warning: package 'purrr' was built under R version 3.6.2
```

```
## Warning: package 'dplyr' was built under R version 3.6.2
## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag() masks stats::lag()
library(summarytools)

## Registered S3 method overwritten by 'pryr':
##   method      from
##   print.bytes Rcpp

## Warning in system2("/usr/bin/otool", c("-L", shQuote(DSO)), stdout = TRUE):
## running command ''/usr/bin/otool' -L '/Library/Frameworks/R.framework/Resources/
## library/tcltk/libs//tcltk.so'' had status 1

##
## Attaching package: 'summarytools'

## The following object is masked from 'package:tibble':
##
##   view
```

Téléchargement de la base de données

```
crsc96 <- read_csv("../Données/cora-crsc1996-E-1996_F1.csv")

##
## -- Column specification -----
## cols(
##   .default = col_double()
## )
## i Use `spec()` for the full column specifications.
```

Sélectionnons les données qui nous intéressent

q1 : - I hate being bossed around: I must feel that I have total control over all the different areas of my life - **Je déteste être patronisé: je dois sentir que j'ai un contrôle total sur tous les différents domaines de ma vie**

q2: - An unmarried girl of 18 should not have sexual relations - *Une fille non mariée de 18 ans ne devrait pas avoir de relations sexuelles*

q3: - The best way to get something from someone is by putting your foot down - *La meilleure façon d'obtenir quelque chose de quelqu'un est de mettre le pied à terre (dialoguer)*

q4: - In a household where both partners are working, is not right for the wife to earn more than the husband - *Dans un ménage où les deux partenaires travaillent, il n'est pas normal que la femme gagne plus que le mari*

q44: - Overpopulation in third world countries doesn't really affect our country - *La surpopulation dans les pays du tiers monde n'affecte pas vraiment notre pays*

q95: - An extramarital affair from time to time is not that serious - *Une liaison extraconjugale de temps en temps n'est pas si grave*

q96: - I would like to have a religious service at my funeral - *J'aimerais avoir un service religieux à mes funérailles*

```
crsc96_small <-
  crsc96 %>%
    select(sexq, region, age, ageq, q1, q2, q3, q4, q44, q95, q96)

crsc96_small

## # A tibble: 2,859 x 11
##   sexq region  age  ageq   q1   q2   q3   q4  q44  q95  q96
##   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1     2     9   33     3     1     5     5     5     4     5     1
## 2     2     9   34     3     2     5     4     5     5     5     5
## 3     2     9   56     4     2     2     4     5     5     4     5
## 4     1     9   69     5     1     4     2     4     5     5     2
## 5     1     9   43     3     4     4     4     5     5     4     4
## 6     2     9   28     2     4     5     4     5     5     5     2
## 7     1     9   27     2     2     4     2     4     4     5     4
## 8     1     9   51     4     1     4     4     5     5     4     2
## 9     1     9   41     3     1     5     5     5     4     4     5
## 10    1     9   39     3     4     2     5     5     4     5     1
## # ... with 2,849 more rows

crsc96_small <-
  crsc96_small %>%
    mutate(age4 = case_when(
      age < 20 ~ "adolescent",
      age >= 20 & age < 34 ~ "jeune",
      age >= 35 & age < 59 ~ "adulte",
      age >= 60 ~ "ainé"
    ))
```

Statistiques univariées

Statistiques univariées

Les objectifs de la statistique descriptive sont de : - définir le ou les groupes étudiées (population ou échantillon) - définir le codage des observations - définir la présentation des données : numérique et/ou graphique - réduire les données à quelques indicateurs statistiques synthétiques

Distribution de fréquences et de pourcentage

Utilisation de base R

Utilisation de tidyverse

```
nombre_sexe <-
  crsc96_small %>%
    count(sexe = sexq)

nombre_age4 <-
```

```
crsc96_small %>%
count(age = age4)
```

Calculer des proportions

```
proportion1 <-
  crsc96_small %>%
  count(sexe = sexq, age = age4) %>%
  mutate(proportion = n / (sum(n)))
proportion1
```

```
## # A tibble: 10 x 4
##   sexe age      n proportion
##   <dbl> <chr>   <int>     <dbl>
## 1     1 adolescent  137     0.0479
## 2     1 adulte    551     0.193
## 3     1 ainé     231     0.0808
## 4     1 jeune    386     0.135
## 5     1 <NA>      56     0.0196
## 6     2 adolescent  140     0.0490
## 7     2 adulte    617     0.216
## 8     2 ainé     295     0.103
## 9     2 jeune    394     0.138
## 10    2 <NA>      52     0.0182
```

proportion

```
proportion2 <-
  crsc96_small %>%
  group_by(sexq) %>%
  count(age = age4) %>%
  mutate(proportion = n / (sum(n)))
proportion2
```

```
## # A tibble: 10 x 4
## # Groups:   sexq [2]
##   sexq age      n proportion
##   <dbl> <chr>   <int>     <dbl>
## 1     1 adolescent  137     0.101
## 2     1 adulte    551     0.405
## 3     1 ainé     231     0.170
## 4     1 jeune    386     0.284
## 5     1 <NA>      56     0.0411
## 6     2 adolescent  140     0.0935
## 7     2 adulte    617     0.412
## 8     2 ainé     295     0.197
## 9     2 jeune    394     0.263
## 10    2 <NA>      52     0.0347
```

Avec le package Summarytools

```
freq(crsc96$sexq)
```

```
## Frequencies
## crsc96$sexq
## Type: Numeric
##
##          Freq  % Valid  % Valid Cum.  % Total  % Total Cum.
## -----
##          1  1361    47.60      47.60    47.60    47.60
##          2  1498    52.40     100.00    52.40   100.00
##         <NA>    0      0.00      0.00    0.00   100.00
##        Total 2859   100.00     100.00   100.00   100.00
```

```
freq(crsc96$q1)
```

```
## Frequencies
## crsc96$q1
## Type: Numeric
##
##          Freq  % Valid  % Valid Cum.  % Total  % Total Cum.
## -----
##          1  1065    37.25      37.25    37.25    37.25
##          2  1410    49.32      86.57    49.32    86.57
##          3    9     0.31      86.88     0.31    86.88
##          4   326    11.40      98.29    11.40    98.29
##          5    49     1.71     100.00     1.71   100.00
##         <NA>    0      0.00      0.00    0.00   100.00
##        Total 2859   100.00     100.00   100.00   100.00
```

```
freq(crsc96$region)
```

```
## Frequencies
## crsc96$region
## Type: Numeric
##
##          Freq  % Valid  % Valid Cum.  % Total  % Total Cum.
## -----
##          0   218     7.63      7.63     7.63     7.63
##          1   270     9.44     17.07     9.44    17.07
##          2   564    19.73     36.80    19.73    36.80
##          3   531    18.57     55.37    18.57    55.37
##          4   211     7.38     62.75     7.38    62.75
##          5   351    12.28     75.03    12.28    75.03
##          6   124     4.34     79.36     4.34    79.36
##          7   117     4.09     83.46     4.09    83.46
##          8   240     8.39     91.85     8.39    91.85
##          9   233     8.15    100.00     8.15   100.00
##         <NA>    0      0.00      0.00    0.00   100.00
##        Total 2859   100.00     100.00   100.00   100.00
```

```
freq(crsc96$q44)
```

```
## Frequencies
## crsc96$q44
```

```
## Type: Numeric
##
##          Freq  % Valid  % Valid Cum.  % Total  % Total Cum.
## -----
##          1    118     4.13         4.13     4.13         4.13
##          2    414    14.48        18.61    14.48        18.61
##          3     18     0.63        19.24     0.63        19.24
##          4   1293    45.23        64.46    45.23        64.46
##          5   1016    35.54       100.00    35.54       100.00
##         <NA>     0         0.00         0.00       100.00
##        Total  2859   100.00       100.00   100.00       100.00
```

```
freq(crsc96$q95)
```

```
## Frequencies
## crsc96$q95
## Type: Numeric
##
##          Freq  % Valid  % Valid Cum.  % Total  % Total Cum.
## -----
##          1     66     2.31         2.31     2.31         2.31
##          2    240     8.39        10.70     8.39        10.70
##          3     22     0.77        11.47     0.77        11.47
##          4    605    21.16        32.63    21.16        32.63
##          5   1926    67.37       100.00    67.37       100.00
##         <NA>     0         0.00         0.00       100.00
##        Total  2859   100.00       100.00   100.00       100.00
```

Paramètres de tendances centrales et de dispersion

Avec base R

La commande `summary` nous donne une première indication sur l'ensemble des variables de notre base de données. Il faut prêter attention aux variables manquantes.

```
summary(crsc96_small)
```

```
##      sexq      region      age      ageq      q1
## Min.   :1.000  Min.   :0.000  Min.   :15.00  Min.   :1.000  Min.   :1.00
## 1st Qu.:1.000  1st Qu.:2.000  1st Qu.:28.00  1st Qu.:2.000  1st Qu.:1.00
## Median :2.000  Median :3.000  Median :39.00  Median :3.000  Median :2.00
## Mean   :1.524  Mean   :3.907  Mean   :41.45  Mean   :3.226  Mean   :1.91
## 3rd Qu.:2.000  3rd Qu.:5.000  3rd Qu.:54.00  3rd Qu.:4.000  3rd Qu.:2.00
## Max.   :2.000  Max.   :9.000  Max.   :99.00  Max.   :5.000  Max.   :5.00
##      q2      q3      q4      q44      q95
## Min.   :1.00  Min.   :1.000  Min.   :1.000  Min.   :1.000  Min.   :1.000
## 1st Qu.:2.00  1st Qu.:2.000  1st Qu.:4.000  1st Qu.:4.000  1st Qu.:4.000
## Median :4.00  Median :4.000  Median :5.000  Median :4.000  Median :5.000
## Mean   :3.28  Mean   :3.685  Mean   :4.524  Mean   :3.936  Mean   :4.429
## 3rd Qu.:5.00  3rd Qu.:5.000  3rd Qu.:5.000  3rd Qu.:5.000  3rd Qu.:5.000
## Max.   :5.00  Max.   :5.000  Max.   :5.000  Max.   :5.000  Max.   :5.000
##      q96      age4
## Min.   :1.0  Length:2859
## 1st Qu.:1.0  Class :character
```

```
## Median :2.0   Mode  :character
## Mean   :2.3
## 3rd Qu.:4.0
## Max.   :5.0
```

L'inconvénient, c'est que c'est mal présenté, et ce ne sont pas l'ensemble des variables de notre base de données qui nous concernent. Les informations sur les variables nominales ne sont pas fournies.

Paramètres de position (Base R)

```
age_moyen <- mean(crsc96_small$age)
age_moyen
```

```
## [1] 41.45261
```

```
age_mediane <- median(crsc96_small$age)
age_mediane
```

```
## [1] 39
```

Cette approche n'est pas la bonne car elle nous demande beaucoup de coding (avec la création de plusieurs objets)

Paramètres de position - tidyverse

La fonction `summarise` permet de calculer l'ensemble des indicateurs dont nous avons besoin. Dans toute étude, il est important de résumer l'information contenue dans les variables pour se faire une première idée.

```
age_position <-
  crsc96_small %>%
    summarise(age_moyen = mean(age),
              age_mediane = median(age),
              age_Q1 = quantile(age, prob = 0.25),
              age_Q3 = quantile(age, prob = 0.75),
              age_min = min(age))
age_position
```

```
## # A tibble: 1 x 5
##   age_moyen age_mediane age_Q1 age_Q3 age_min
##   <dbl>      <dbl> <dbl> <dbl> <dbl>
## 1    41.5         39    28    54    15
```

Statistiques univariées: Mode

Il n'y a aucune fonction qui permet de calculer directement le mode. Alors, il faut la créer soi-même.

```
getmode <- function(v) {
  uniqv <- unique(v)
  uniqv[which.max(tabulate(match(v, uniqv)))]
}
```

```
age_position <-
  crsc96_small %>%
```

```
summarise(age_moyen = mean(age),
          age_median = median(age),
          age_Q1 = quantile(age, prob = 0.25),
          age_Q3 = quantile(age, prob = 0.75),
          age_mode = getmode(age))
```

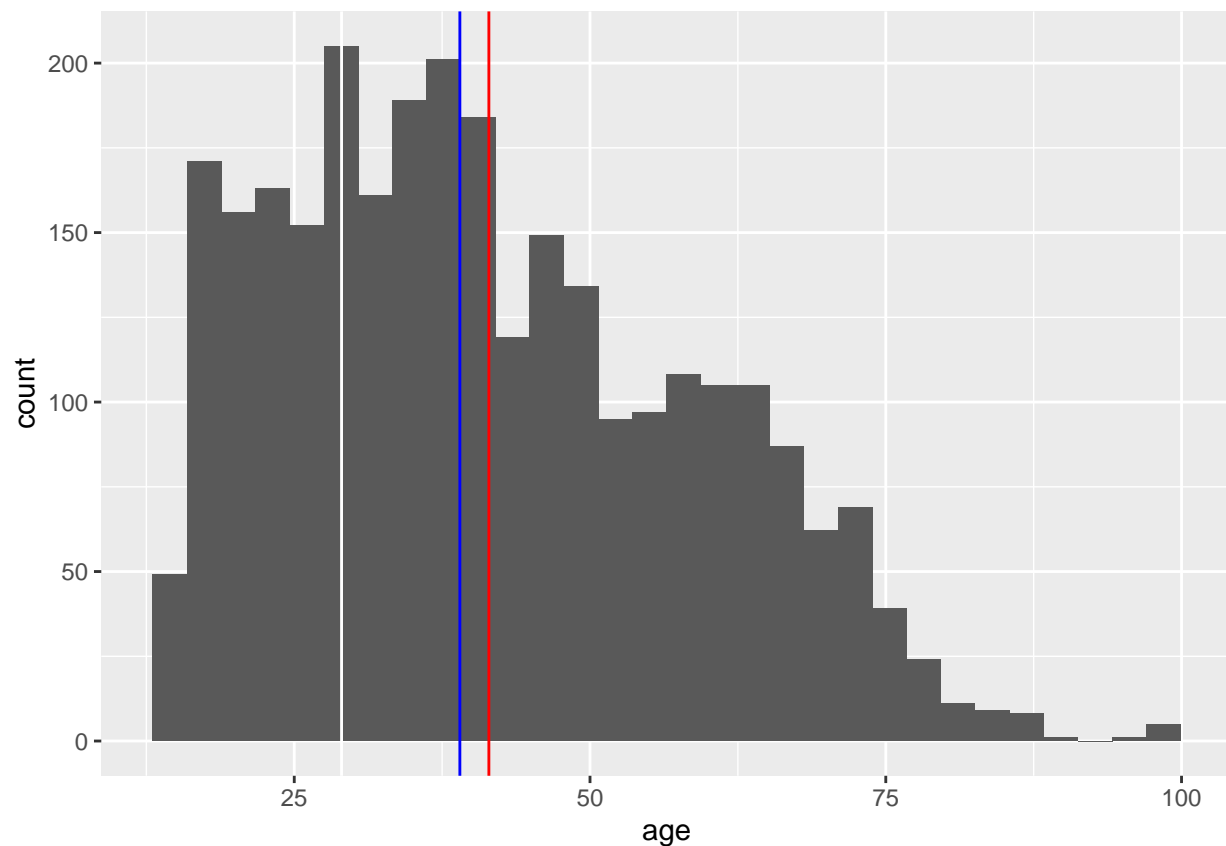
```
age_position
```

```
## # A tibble: 1 x 4
##   age_moyen age_median age_Q1 age_mode
##   <dbl>      <dbl> <dbl> <dbl>
## 1    41.5        39    54    29
```

Statistiques univariées : Histogramme

```
ggplot(crs96_small) +
  geom_histogram(aes(x = age)) +
  geom_vline(aes(xintercept = mean(age)), color = "red") +
  geom_vline(aes(xintercept = median(age)), color = "blue") +
  geom_vline(aes(xintercept = getmode(age)), color = "white")
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



On reviendra sur le cours prochain sur la visualisation, l'une des forces de tidyverse.

Statistique par groupe

Nous pouvons aussi regarder ces données selon le sexe des individus

```
age_position_sexe <-  
  crsc96_small %>%  
  group_by(sexq) %>%  
  summarise(age_moyen = mean(age),  
            age_median = median(age),  
            age_Q1 = quantile(age, prob = 0.25),  
            age_Q3 = quantile(age, prob = 0.75),  
            age_mode = getmode(age))  
  
age_position_sexe  
  
## # A tibble: 2 x 6  
##   sexq age_moyen age_median age_Q1 age_Q3 age_mode  
##   <dbl>   <dbl>     <dbl> <dbl> <dbl>   <dbl>  
## 1     1     40.8       39     27    52     44  
## 2     2     42.0       40     28    55     38
```

```
age_position_sexe <-  
  age_position_sexe %>%  
  mutate(écart = age_moyen - age_median)  
  
age_position_age4 <-  
  crsc96_small %>%  
  group_by(age4) %>%  
  summarise(age_moyen = mean(age),  
            age_median = median(age),  
            age_Q1 = quantile(age, prob = 0.25),  
            age_Q3 = quantile(age, prob = 0.75),  
            age_mode = getmode(age))  
  
age_position_age4  
  
## # A tibble: 5 x 5  
##   age4      age_moyen age_median age_Q1 age_mode  
##   <chr>         <dbl>     <dbl> <dbl>   <dbl>  
## 1 adolescent  17.1       17     18     18  
## 2 adulte     44.7       44     50     38  
## 3 ainé       68.5       67     72     60  
## 4 jeune      26.7       27     30     29  
## 5 <NA>       44.0       34     59     34
```

Statistiques univariés sur plusieurs variables, solution alternative

On regarde les statistiques pas sur une seule variable mais sur plusieurs variables. On peut combiner plusieurs tableaux avec les fonctions `binds_row` et `binds_col`.

```
position_q <-  
  crsc96_small %>%  
  select(num_range("q", c(1:4, 44, 95))) %>%  
  summarise_each(funs(mean, median))
```

```
## Warning: `summarise_each()` was deprecated in dplyr 0.7.0.
## Please use `across()` instead.

## Warning: `funs()` was deprecated in dplyr 0.8.0.
## Please use a list of either functions or lambdas:
##
##   # Simple named list:
##   list(mean = mean, median = median)
##
##   # Auto named with `tibble::lst()`:
##   tibble::lst(mean, median)
##
##   # Using lambdas
##   list(~ mean(., trim = .2), ~ median(., na.rm = TRUE))

position_q

## # A tibble: 1 x 12
##   q1_mean q2_mean q3_mean q4_mean q44_mean q95_mean q1_median q2_median
##   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1    1.91    3.28    3.68    4.52    3.94    4.43     2       4
## # ... with 4 more variables: q3_median <dbl>, q4_median <dbl>,
## #   q44_median <dbl>, q95_median <dbl>
```

Le problème, c'est la longueur du fichier.

EXERCICE

Calculer les paramètres de dispersion de la variable age et commenter.

Application 1: Données abérantes ou extrêmes

```
#install.packages("carData")
library(carData)
library(tidyverse)
data(package = "carData")

data("Davis", package = "carData")
save(Davis, file = "Davis.Rdata")

load(file = "Davis.Rdata")

summary(Davis)
```

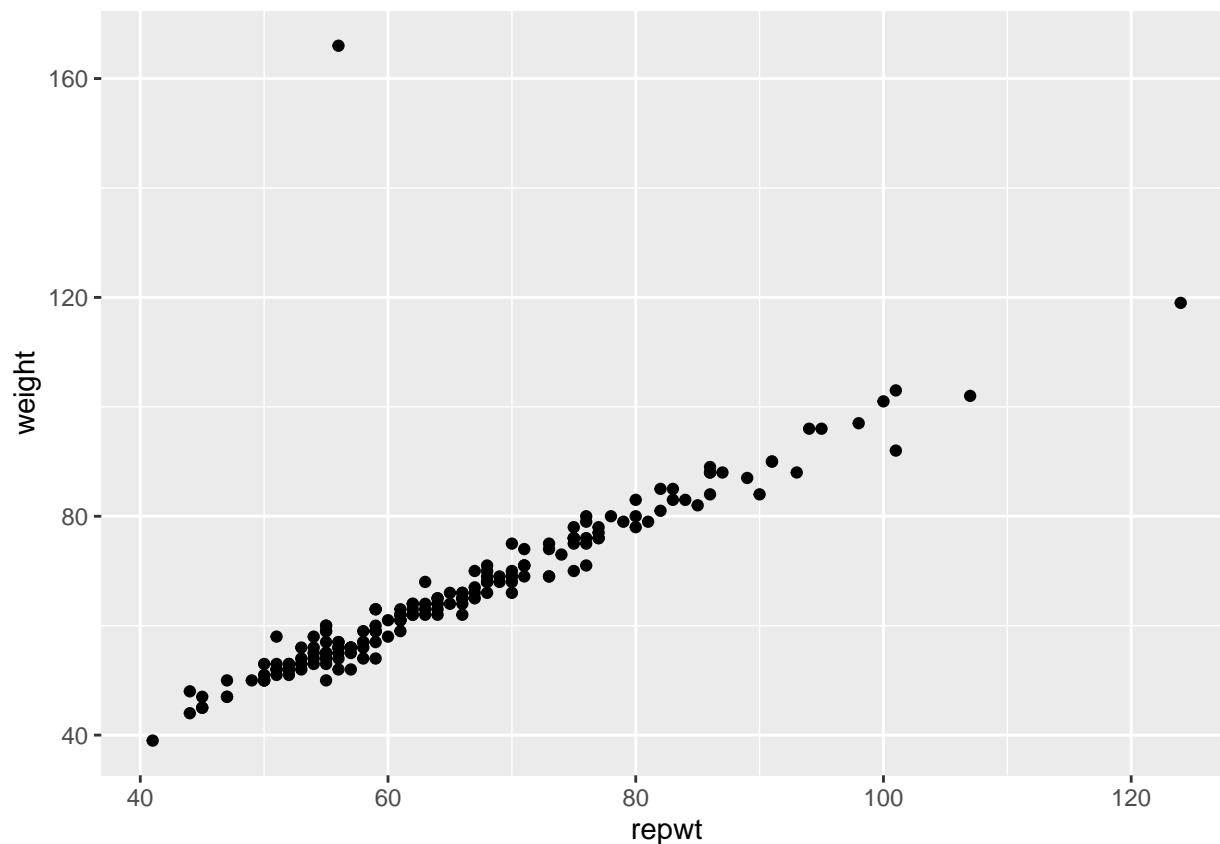
```
## sex      weight      height      repwt      repht
## F:112   Min.       : 39.0   Min.       : 57.0   Min.       : 41.00  Min.       :148.0
## M: 88   1st Qu.: 55.0   1st Qu.:164.0   1st Qu.: 55.00  1st Qu.:160.5
##        Median : 63.0   Median :169.5   Median : 63.00  Median :168.0
##        Mean   : 65.8   Mean    :170.0   Mean    : 65.62  Mean    :168.5
##        3rd Qu.: 74.0   3rd Qu.:177.2   3rd Qu.: 73.50  3rd Qu.:175.0
##        Max.    :166.0   Max.     :197.0   Max.     :124.00  Max.     :200.0
##                                     NA's      :17      NA's      :17
```

Ce fichier comprend les informations sur le poids et la taille de 200 individus ainsi que leur poids et taille autodéclaré. On veut voir de quelle manière les poids auto-déclarés sont fiables.

On ne l'a pas encore vu, mais on peut rapidement voir les variables deux par deux dans un graphique

```
ggplot(Davis) +  
  geom_point(aes(x = repwt, y = weight))
```

```
## Warning: Removed 17 rows containing missing values (geom_point).
```



Application 2 : Données manquantes

```
poids_moyen <-  
  Davis %>%  
  summarise(poids_moyen = mean(weight))  
poids_moyen  
  
##   poids_moyen  
## 1         65.8  
  
poids_moyen_reporte <-  
  Davis %>%  
  summarise(poids_moyen_reporte = mean(repwt))  
poids_moyen_reporte  
  
##   poids_moyen_reporte  
## 1                  NA
```

Qu'est-ce qui s'est passé. En fait, le poids reporté comporte des valeurs manquantes. Il faut indiquer dans le calcul de la moyenne qu'il y a des valeurs manquantes, et qu'il faut les enlever avant de calculer la moyenne, ou toute autre statistique.

```
poids_moyen_reporte <-
  Davis %>%
  summarise(poids_moyen_reporte = mean(repwt, na.rm = TRUE))
poids_moyen_reporte
```

```
##   poids_moyen_reporte
## 1                65.62295
```

Quel est le problème qui se pose quand des informations sont manquantes. Peut-on faire confiance aux résultats?

Remarques

1. Tous les objets que vous créez, vous pouvez les manipuler à votre guise
2. Les variables que vous créez, vous pouvez les réutiliser juste après
3. Interprétations des résultats

Statistiques bivariées : Association entre variables

Statistiques bivariées : Association entre variables

Existe-il une relation entre l'âge et l'opinion des gens? Existe-il une relation entre le sexe et l'opinion des gens?

```
crsc96_small <-
  crsc96_small %>%
  mutate(sexe = factor(sexq, labels = c("Homme", "Femme")))

freq(crsc96_small$sexe)
```

```
## Frequencies
## crsc96_small$sexe
## Type: Factor
##
##           Freq  % Valid  % Valid Cum.  % Total  % Total Cum.
## -----
##      Homme  1361    47.60      47.60    47.60    47.60
##      Femme  1498    52.40     100.00    52.40   100.00
##      <NA>     0     0.00     100.00    0.00   100.00
##      Total  2859   100.00     100.00   100.00   100.00
```

Association

- Sexe avec la variable q2

q2: "An unmarried girl of 18 should not have sexual relations" Une jeune fille non mariée de 18 ans ne devrait pas avoir de relations sexuelles - Se fait à base de tableaux croisés (contingency table)

```
qlabel <- c("totally agree", "agree somewhat", "DK/NA", "disagree somewhat", "totally disagree")

crsc96_small <-
  crsc96_small %>%
    mutate(q2_new = factor(q2, labels = qlabel),
           q3_new = factor(q3, labels = qlabel))

ctable(crsc96_small$sexe, crsc96_small$q2_new)

## Cross-Tabulation, Row Proportions
## sexe * q2_new
## Data Frame: crsc96_small
##
## -----
##      q2_new    totally agree    agree somewhat    DK/NA    disagree somewhat    totally disagree
##  sexe
##  Homme          208 (15.3%)      304 (22.3%)      12 (0.9%)          418 (30.7%)          419 (30.8%)
##  Femme          308 (20.6%)      332 (22.2%)      14 (0.9%)          476 (31.8%)          368 (24.6%)
##  Total          516 (18.0%)      636 (22.2%)      26 (0.9%)          894 (31.3%)          787 (27.5%)
## -----
```

Association

Les colonnes et les lignes d'un tableau croisés, ne sont pas identiques.

```
ctable(crsc96_small$q2_new, crsc96_small$sexe)
```

```
## Cross-Tabulation, Row Proportions
## q2_new * sexe
## Data Frame: crsc96_small
##
## -----
##      sexe      Homme      Femme      Total
##  q2_new
##  totally agree      208 (40.3%)      308 (59.7%)      516 (100.0%)
##  agree somewhat      304 (47.8%)      332 (52.2%)      636 (100.0%)
##  DK/NA              12 (46.2%)      14 (53.8%)      26 (100.0%)
##  disagree somewhat      418 (46.8%)      476 (53.2%)      894 (100.0%)
##  totally disagree      419 (53.2%)      368 (46.8%)      787 (100.0%)
##  Total              1361 (47.6%)      1498 (52.4%)      2859 (100.0%)
## -----
```

Lequel des deux tableaux donne une indication sur l'association entre les deux variables?

Association

Aussi, est-il important de préciser si vous calculez des proportions lignes ou des proportions colonnes.

```
ctable(crsc96_small$sexe, crsc96_small$q2_new, "r")
```

```
## Cross-Tabulation, Row Proportions
## sexe * q2_new
## Data Frame: crsc96_small
##
```

```
## -----
##      q2_new  totally agree  agree somewhat      DK/NA  disagree somewhat  totally disagree
##  sexe
##  Homme      208 (15.3%)    304 (22.3%)    12 (0.9%)      418 (30.7%)    419 (30.8%)
##  Femme      308 (20.6%)    332 (22.2%)    14 (0.9%)      476 (31.8%)    368 (24.6%)
##  Total      516 (18.0%)    636 (22.2%)    26 (0.9%)      894 (31.3%)    787 (27.5%)
## -----
```

```
ctable(crsc96_small$sexe, crsc96_small$q2_new, "c")
```

```
## Cross-Tabulation, Column Proportions
```

```
## sexe * q2_new
```

```
## Data Frame: crsc96_small
```

```
## -----
##      q2_new  totally agree  agree somewhat      DK/NA  disagree somewhat  totally disagree
##  sexe
##  Homme      208 ( 40.3%)    304 ( 47.8%)    12 ( 46.2%)      418 ( 46.8%)    419 ( 53.2%)
##  Femme      308 ( 59.7%)    332 ( 52.2%)    14 ( 53.8%)      476 ( 53.2%)    368 ( 46.8%)
##  Total      516 (100.0%)    636 (100.0%)    26 (100.0%)      894 (100.0%)    787 (100.0%)
## -----
```

Questions

Que faites vous si une des variables est quantitative, par exemple l'âge et q2_new