

Indexation et recodages avec R et questionr

Julien Barnier

28 novembre 2013

Au menu

PIZZA		BURGERS		FISH & CHIPS	
ALL PIZZAS HAVE TOMATOES, TOMATO SAUCE, SPICES, CHEESE, GARLIC & HERBS.		ALL BURGERS COME WITH MAYONNAISE (ONLY WHEN TOLERATED)		<p>Chips are available \$2.00</p> <p>Portions from \$1.50</p> <p>Fish: Deep Fry (along with bread)</p> <p>8oz. Fresh Fish of the Day (Seasonal)</p> <p>1.80</p> <p>1.90</p> <p>2.50</p> <p>3.00</p> <p>3.50</p> <p>4.00</p> <p>4.50</p> <p>5.00</p> <p>5.50</p> <p>6.00</p> <p>6.50</p> <p>7.00</p> <p>7.50</p> <p>8.00</p> <p>8.50</p> <p>9.00</p> <p>9.50</p> <p>10.00</p> <p>10.50</p> <p>11.00</p> <p>11.50</p> <p>12.00</p> <p>12.50</p> <p>13.00</p> <p>13.50</p> <p>14.00</p> <p>14.50</p> <p>15.00</p> <p>15.50</p> <p>16.00</p> <p>16.50</p> <p>17.00</p> <p>17.50</p> <p>18.00</p> <p>18.50</p> <p>19.00</p> <p>19.50</p> <p>20.00</p> <p>20.50</p> <p>21.00</p> <p>21.50</p> <p>22.00</p> <p>22.50</p> <p>23.00</p> <p>23.50</p> <p>24.00</p> <p>24.50</p> <p>25.00</p> <p>25.50</p> <p>26.00</p> <p>26.50</p> <p>27.00</p> <p>27.50</p> <p>28.00</p> <p>28.50</p> <p>29.00</p> <p>29.50</p> <p>30.00</p> <p>30.50</p> <p>31.00</p> <p>31.50</p> <p>32.00</p> <p>32.50</p> <p>33.00</p> <p>33.50</p> <p>34.00</p> <p>34.50</p> <p>35.00</p> <p>35.50</p> <p>36.00</p> <p>36.50</p> <p>37.00</p> <p>37.50</p> <p>38.00</p> <p>38.50</p> <p>39.00</p> <p>39.50</p> <p>40.00</p> <p>40.50</p> <p>41.00</p> <p>41.50</p> <p>42.00</p> <p>42.50</p> <p>43.00</p> <p>43.50</p> <p>44.00</p> <p>44.50</p> <p>45.00</p> <p>45.50</p> <p>46.00</p> <p>46.50</p> <p>47.00</p> <p>47.50</p> <p>48.00</p> <p>48.50</p> <p>49.00</p> <p>49.50</p> <p>50.00</p> <p>50.50</p> <p>51.00</p> <p>51.50</p> <p>52.00</p> <p>52.50</p> <p>53.00</p> <p>53.50</p> <p>54.00</p> <p>54.50</p> <p>55.00</p> <p>55.50</p> <p>56.00</p> <p>56.50</p> <p>57.00</p> <p>57.50</p> <p>58.00</p> <p>58.50</p> <p>59.00</p> <p>59.50</p> <p>60.00</p> <p>60.50</p> <p>61.00</p> <p>61.50</p> <p>62.00</p> <p>62.50</p> <p>63.00</p> <p>63.50</p> <p>64.00</p> <p>64.50</p> <p>65.00</p> <p>65.50</p> <p>66.00</p> <p>66.50</p> <p>67.00</p> <p>67.50</p> <p>68.00</p> <p>68.50</p> <p>69.00</p> <p>69.50</p> <p>70.00</p> <p>70.50</p> <p>71.00</p> <p>71.50</p> <p>72.00</p> <p>72.50</p> <p>73.00</p> <p>73.50</p> <p>74.00</p> <p>74.50</p> <p>75.00</p> <p>75.50</p> <p>76.00</p> <p>76.50</p> <p>77.00</p> <p>77.50</p> <p>78.00</p> <p>78.50</p> <p>79.00</p> <p>79.50</p> <p>80.00</p> <p>80.50</p> <p>81.00</p> <p>81.50</p> <p>82.00</p> <p>82.50</p> <p>83.00</p> <p>83.50</p> <p>84.00</p> <p>84.50</p> <p>85.00</p> <p>85.50</p> <p>86.00</p> <p>86.50</p> <p>87.00</p> <p>87.50</p> <p>88.00</p> <p>88.50</p> <p>89.00</p> <p>89.50</p> <p>90.00</p> <p>90.50</p> <p>91.00</p> <p>91.50</p> <p>92.00</p> <p>92.50</p> <p>93.00</p> <p>93.50</p> <p>94.00</p> <p>94.50</p> <p>95.00</p> <p>95.50</p> <p>96.00</p> <p>96.50</p> <p>97.00</p> <p>97.50</p> <p>98.00</p> <p>98.50</p> <p>99.00</p> <p>99.50</p> <p>100.00</p>	
<p>CRISP, HOT, VANILLA BEAN ICE CREAM OR VANILLA SAUCE (1.50) (2.00) (2.50) (3.00)</p> <p>WASTE SPECIAL (1.50) (2.00) (2.50) (3.00)</p> <p>SEASON SUPREME (1.50) (2.00) (2.50) (3.00)</p> <p>VEGETARIAN (1.50) (2.00) (2.50) (3.00)</p> <p>CHICKEN (1.50) (2.00) (2.50) (3.00)</p> <p>TRADITIONAL (1.50) (2.00) (2.50) (3.00)</p> <p>PIZZA LOVERS (1.50) (2.00) (2.50) (3.00)</p> <p>CLASSIC (1.50) (2.00) (2.50) (3.00)</p> <p>SUPREMO (1.50) (2.00) (2.50) (3.00)</p> <p>GARLIC PIZZA VARIANO (1.50) (2.00) (2.50) (3.00)</p>	<p>HAMBURGER \$ 1.50</p> <p>CHEESEBURGER \$ 1.50</p> <p>BEAN & EGG BURGER \$ 1.50</p> <p>FISH BURGER \$ 1.50</p> <p>HAWAIIAN BURGER \$ 1.50</p> <p>TRIPLE BURGER \$ 1.50</p> <p>CHICKEN SCHNITZEL BURGER \$ 1.50</p> <p>VEGETARIAN BURGER \$ 1.50</p>	<p>Chips are available \$2.00</p> <p>Portions from \$1.50</p> <p>Fish: Deep Fry (along with bread)</p> <p>8oz. Fresh Fish of the Day (Seasonal)</p> <p>1.80</p> <p>1.90</p> <p>2.50</p> <p>3.00</p> <p>3.50</p> <p>4.00</p> <p>4.50</p> <p>5.00</p> <p>5.50</p> <p>6.00</p> <p>6.50</p> <p>7.00</p> <p>7.50</p> <p>8.00</p> <p>8.50</p> <p>9.00</p> <p>9.50</p> <p>10.00</p> <p>10.50</p> <p>11.00</p> <p>11.50</p> <p>12.00</p> <p>12.50</p> <p>13.00</p> <p>13.50</p> <p>14.00</p> <p>14.50</p> <p>15.00</p> <p>15.50</p> <p>16.00</p> <p>16.50</p> <p>17.00</p> <p>17.50</p> <p>18.00</p> <p>18.50</p> <p>19.00</p> <p>19.50</p> <p>20.00</p> <p>20.50</p> <p>21.00</p> <p>21.50</p> <p>22.00</p> <p>22.50</p> <p>23.00</p> <p>23.50</p> <p>24.00</p> <p>24.50</p> <p>25.00</p> <p>25.50</p> <p>26.00</p> <p>26.50</p> <p>27.00</p> <p>27.50</p> <p>28.00</p> <p>28.50</p> <p>29.00</p> <p>29.50</p> <p>30.00</p> <p>30.50</p> <p>31.00</p> <p>31.50</p> <p>32.00</p> <p>32.50</p> <p>33.00</p> <p>33.50</p> <p>34.00</p> <p>34.50</p> <p>35.00</p> <p>35.50</p> <p>36.00</p> <p>36.50</p> <p>37.00</p> <p>37.50</p> <p>38.00</p> <p>38.50</p> <p>39.00</p> <p>39.50</p> <p>40.00</p> <p>40.50</p> <p>41.00</p> <p>41.50</p> <p>42.00</p> <p>42.50</p> <p>43.00</p> <p>43.50</p> <p>44.00</p> <p>44.50</p> <p>45.00</p> <p>45.50</p> <p>46.00</p> <p>46.50</p> <p>47.00</p> <p>47.50</p> <p>48.00</p> <p>48.50</p> <p>49.00</p> <p>49.50</p> <p>50.00</p> <p>50.50</p> <p>51.00</p> <p>51.50</p> <p>52.00</p> <p>52.50</p> <p>53.00</p> <p>53.50</p> <p>54.00</p> <p>54.50</p> <p>55.00</p> <p>55.50</p> <p>56.00</p> <p>56.50</p> <p>57.00</p> <p>57.50</p> <p>58.00</p> <p>58.50</p> <p>59.00</p> <p>59.50</p> <p>60.00</p> <p>60.50</p> <p>61.00</p> <p>61.50</p> <p>62.00</p> <p>62.50</p> <p>63.00</p> <p>63.50</p> <p>64.00</p> <p>64.50</p> <p>65.00</p> <p>65.50</p> <p>66.00</p> <p>66.50</p> <p>67.00</p> <p>67.50</p> <p>68.00</p> <p>68.50</p> <p>69.00</p> <p>69.50</p> <p>70.00</p> <p>70.50</p> <p>71.00</p> <p>71.50</p> <p>72.00</p> <p>72.50</p> <p>73.00</p> <p>73.50</p> <p>74.00</p> <p>74.50</p> <p>75.00</p> <p>75.50</p> <p>76.00</p> <p>76.50</p> <p>77.00</p> <p>77.50</p> <p>78.00</p> <p>78.50</p> <p>79.00</p> <p>79.50</p> <p>80.00</p> <p>80.50</p> <p>81.00</p> <p>81.50</p> <p>82.00</p> <p>82.50</p> <p>83.00</p> <p>83.50</p> <p>84.00</p> <p>84.50</p> <p>85.00</p> <p>85.50</p> <p>86.00</p> <p>86.50</p> <p>87.00</p> <p>87.50</p> <p>88.00</p> <p>88.50</p> <p>89.00</p> <p>89.50</p> <p>90.00</p> <p>90.50</p> <p>91.00</p> <p>91.50</p> <p>92.00</p> <p>92.50</p> <p>93.00</p> <p>93.50</p> <p>94.00</p> <p>94.50</p> <p>95.00</p> <p>95.50</p> <p>96.00</p> <p>96.50</p> <p>97.00</p> <p>97.50</p> <p>98.00</p> <p>98.50</p> <p>99.00</p> <p>99.50</p> <p>100.00</p>			

Au menu

- Manipulation de données : indexation et recodages

Au menu

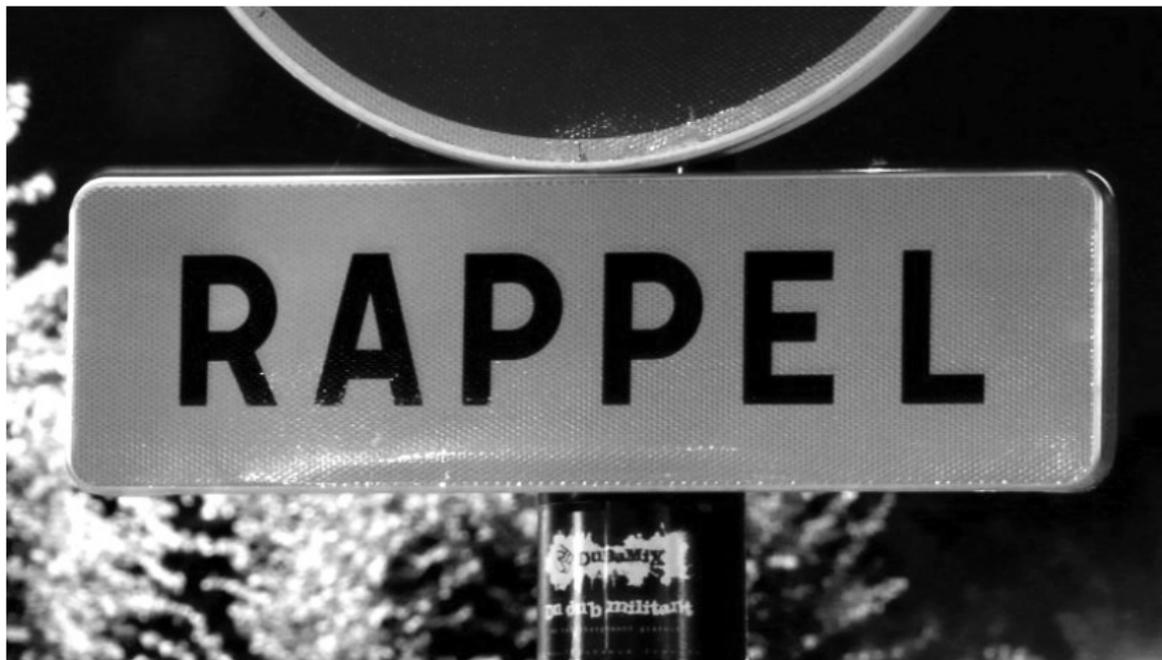
- Manipulation de données : indexation et recodages
- Le document *Introduction à R* :
<http://alea.fr.eu.org/pages/intro-R>

Au menu

- Manipulation de données : indexation et recodages
- Le document *Introduction à R* :
<http://alea.fr.eu.org/pages/intro-R>
- L'extension `questionr`

```
R> install.packages("questionr")
```

Rappels



Rappels

Écrire un script R, c'est appliquer des *fonctions* à des *objets*.

```
R> x <- c(150, 182, 197, 171)
```

```
R> mean(x)
```

```
[1] 175
```

Les objets

Les objets en R servent à stocker des données. Ils peuvent être de différents types :

Les objets

Les objets en R servent à stocker des données. Ils peuvent être de différents types :

- Nombre

```
R> x <- 3
```

Les objets

Les objets en R servent à stocker des données. Ils peuvent être de différents types :

- Nombre

```
R> x <- 3
```

- Chaîne de caractère

```
R> x <- "Chihuahua"
```

Les objets

Les objets en R servent à stocker des données. Ils peuvent être de différents types :

- Nombre

```
R> x <- 3
```

- Chaîne de caractère

```
R> x <- "Chihuahua"
```

- Vecteur

```
R> x <- c(150, 182, 197, 171)
```

```
R> y <- c("rouge", "vert", "rouge")
```

Les objets

Ils peuvent aussi stocker des données beaucoup plus complexes :

- Un tableau de données (`data.frame`)

```
R> df <- read.csv("fichier.csv", stringsAsFactors = FALSE)
```

- Le résultat d'un tri croisé ou d'un test

```
R> tab <- table(df$sexe, df$revenu)
```

```
R> test <- chisq.test(tab)
```

- Un graphique
- etc. etc.

Les fonctions

Elles prennent une suite d'arguments entre parenthèses, et renvoient un résultat :

```
R> c(150, 182, 197, 171)
```

```
[1] 150 182 197 171
```

```
R> mean(x)
```

```
[1] 175
```

Les fonctions

Elles prennent une suite d'arguments entre parenthèses, et renvoient un résultat :

```
R> c(150, 182, 197, 171)
```

```
[1] 150 182 197 171
```

```
R> mean(x)
```

```
[1] 175
```

Certains arguments sont nommés :

```
R> mean(x, na.rm = TRUE)
```

Travail avec des données

Indexation



Indexation

- Indexer, c'est extraire une partie des valeurs d'un objet.
- L'opérateur d'indexation est le crochet [].
- Il existe plusieurs types d'indexation, selon ce qu'on met entre les crochets.

Indexation par position

Si on utilise des nombres entiers, on indique les positions des éléments à sélectionner.

Indexation par position

Si on utilise des nombres entiers, on indique les positions des éléments à sélectionner.

```
R> x <- c(2, 5, 8, 11)
```

Indexation par position

Si on utilise des nombres entiers, on indique les positions des éléments à sélectionner.

```
R> x <- c(2, 5, 8, 11)
```

```
R> x[1]
```

```
[1] 2
```

Indexation par position

Si on utilise des nombres entiers, on indique les positions des éléments à sélectionner.

```
R> x <- c(2, 5, 8, 11)
```

```
R> x[1]
```

```
[1] 2
```

```
R> x[3]
```

```
[1] 8
```

Indexation par position

On peut sélectionner plusieurs éléments en indexant avec un vecteur.

Indexation par position

On peut sélectionner plusieurs éléments en indexant avec un vecteur.

```
R> x <- c(2, 5, 8, 11)
```

Indexation par position

On peut sélectionner plusieurs éléments en indexant avec un vecteur.

```
R> x <- c(2, 5, 8, 11)
```

```
R> x[c(1, 4)]
```

```
[1] 2 11
```

Indexation par position

On peut sélectionner plusieurs éléments en indexant avec un vecteur.

```
R> x <- c(2, 5, 8, 11)
```

```
R> x[c(1, 4)]
```

```
[1] 2 11
```

```
R> x[c(2, 2, 3, 1)]
```

```
[1] 5 5 8 2
```

Indexation par position

Pour l'instant on a indexé seulement un vecteur

```
R> x[positions]
```

Cette indexation peut également se faire sur un *data frame*.

```
R> df[lignes, colonnes]
```

Indexation par position

```
    sexe age
1 homme 25
2 femme 37
3 femme 21
```

Indexation par position

```
      sexe age  
1 homme  25  
2 femme  37  
3 femme  21
```

```
R> df[2, 1]
```

```
[1] "femme"
```

Indexation par position

```
      sexe age
1 homme  25
2 femme  37
3 femme  21
```

```
R> df[2, 1]
```

```
[1] "femme"
```

```
R> df[c(1, 3), 2]
```

```
[1] 25 21
```

Indexation par position

Cas particulier : si on n'indique rien, on sélectionne tout.

Indexation par position

Cas particulier : si on n'indique rien, on sélectionne tout.

C'est vrai pour un vecteur.

```
R> x <- c(2, 5, 8, 11)
```

```
R> x[]
```

```
[1] 2 5 8 11
```

Indexation par position

C'est vrai aussi, et c'est plus intéressant, pour un *data frame*.

```
   sexe age
1 homme 25
2 femme 37
3 femme 21
```

Indexation par position

C'est vrai aussi, et c'est plus intéressant, pour un *data frame*.

```
  sexe age
1 homme 25
2 femme 37
3 femme 21
```

```
R> df[1, ]
```

```
  sexe age
1 homme 25
```

Indexation par position

C'est vrai aussi, et c'est plus intéressant, pour un *data frame*.

```
  sexe age
1 homme 25
2 femme 37
3 femme 21
```

```
R> df[1, ]
```

```
  sexe age
1 homme 25
```

```
R> df[, 2]
```

```
[1] 25 37 21
```

Exemple sur des données

Indexation par nom

Si les éléments du vecteur sont nommés, on peut les sélectionner en utilisant leurs noms.

Indexation par nom

Si les éléments du vecteur sont nommés, on peut les sélectionner en utilisant leurs noms.

```
R> xn <- c(a = 2, b = 57, c = 34)
```

Indexation par nom

Si les éléments du vecteur sont nommés, on peut les sélectionner en utilisant leurs noms.

```
R> xn <- c(a = 2, b = 57, c = 34)
```

```
R> xn["a"]
```

```
a
```

```
2
```

```
R> xn[c("a", "a", "c")]
```

```
a  a  c
```

```
2  2 34
```

Indexation par nom

On peut également indexer par nom pour les *data frame*, notamment sur les colonnes.

```
   sexe age
1 homme 25
2 femme 37
3 femme 21
```

Indexation par nom

On peut également indexer par nom pour les *data frame*, notamment sur les colonnes.

```
  sexe age
1 homme 25
2 femme 37
3 femme 21
```

```
R> df[, "sexe"]
```

```
[1] "homme" "femme" "femme"
```

ou l'équivalent :

```
R> df$sexe
```

```
[1] "homme" "femme" "femme"
```

Exemple sur des données

Indexation logique

Troisième type d'indexation : l'indexation par un vecteur de valeurs logiques.

Indexation logique

Troisième type d'indexation : l'indexation par un vecteur de valeurs logiques.

```
R> x <- c(2, 5, 8, 11)
```

```
R> x[c(TRUE, TRUE, FALSE, FALSE)]
```

```
[1] 2 5
```

Indexation logique

Troisième type d'indexation : l'indexation par un vecteur de valeurs logiques.

```
R> x <- c(2, 5, 8, 11)
R> x[c(TRUE, TRUE, FALSE, FALSE)]
```

```
[1] 2 5
```

Le vecteur de valeur logique peut être recyclé.

```
R> y <- 1:20
R> y[c(TRUE, FALSE)]
```

```
[1] 1 3 5 7 9 11 13 15 17 19
```

Indexation logique

Là aussi, ça fonctionne sur les *data frame*.

```
   sexe age
1 homme  25
2 femme  37
3 femme  21
```

Indexation logique

Là aussi, ça fonctionne sur les *data frame*.

```
  sexe age
1 homme 25
2 femme 37
3 femme 21
```

```
R> df[c(TRUE, TRUE, FALSE), c(FALSE, TRUE)]
```

```
[1] 25 37
```



Tests

Un test est une opération logique qui renvoie vrai (TRUE) ou faux (FALSE).

```
R> 1 == 2
```

```
[1] FALSE
```

Tests

Un test est une opération logique qui renvoie vrai (TRUE) ou faux (FALSE).

```
R> 1 == 2
```

```
[1] FALSE
```

```
R> "Yorkshire" == "Chihuahua"
```

```
[1] FALSE
```

Tests

Un test est une opération logique qui renvoie vrai (TRUE) ou faux (FALSE).

```
R> 1 == 2
```

```
[1] FALSE
```

```
R> "Yorkshire" == "Chihuahua"
```

```
[1] FALSE
```

```
R> 3 < 5
```

```
[1] TRUE
```

Tests

Un test est une opération logique qui renvoie vrai (TRUE) ou faux (FALSE).

```
R> 1 == 2
```

```
[1] FALSE
```

```
R> "Yorkshire" == "Chihuahua"
```

```
[1] FALSE
```

```
R> 3 < 5
```

```
[1] TRUE
```

```
R> 1.5 >= mean(c(3, 4, 5))
```

```
[1] FALSE
```

Tests

On peut appliquer un test directement à un vecteur.

Tests

On peut appliquer un test directement à un vecteur.

```
R> x <- c(2, 5, 8, 11)
```

```
R> y <- c("Rouge", "Vert", "Rouge", "Bleu")
```

Tests

On peut appliquer un test directement à un vecteur.

```
R> x <- c(2, 5, 8, 11)
```

```
R> y <- c("Rouge", "Vert", "Rouge", "Bleu")
```

```
R> x == 5
```

```
[1] FALSE TRUE FALSE FALSE
```

Tests

On peut appliquer un test directement à un vecteur.

```
R> x <- c(2, 5, 8, 11)
```

```
R> y <- c("Rouge", "Vert", "Rouge", "Bleu")
```

```
R> x == 5
```

```
[1] FALSE TRUE FALSE FALSE
```

```
R> x > 6
```

```
[1] FALSE FALSE TRUE TRUE
```

Tests

On peut appliquer un test directement à un vecteur.

```
R> x <- c(2, 5, 8, 11)
```

```
R> y <- c("Rouge", "Vert", "Rouge", "Bleu")
```

```
R> x == 5
```

```
[1] FALSE TRUE FALSE FALSE
```

```
R> x > 6
```

```
[1] FALSE FALSE TRUE TRUE
```

```
R> y == "Rouge"
```

```
[1] TRUE FALSE TRUE FALSE
```

Indexation par test

On peut indexer avec un vecteur logique

et

Un test renvoie un vecteur logique

donc

On peut indexer avec un test

Indexation par test

```
R> x <- c(2, 5, 8, 11)
```

Indexation par test

```
R> x <- c(2, 5, 8, 11)
```

```
R> x == 5
```

```
[1] FALSE TRUE FALSE FALSE
```

Indexation par test

```
R> x <- c(2, 5, 8, 11)
```

```
R> x == 5
```

```
[1] FALSE TRUE FALSE FALSE
```

```
R> x[x == 5]
```

```
[1] 5
```

Indexation par test

```
R> x <- c(2, 5, 8, 11)
```

```
R> x == 5
```

```
[1] FALSE TRUE FALSE FALSE
```

```
R> x[x == 5]
```

```
[1] 5
```

```
R> x[x > 7]
```

```
[1] 8 11
```

Indexation par test

Le test peut tout à fait porter sur un autre vecteur.

Indexation par test

Le test peut tout à fait porter sur un autre vecteur.

```
R> taille <- c(151, 187, 178, 169)
```

```
R> sexe <- c("Homme", "Femme", "Homme", "Femme")
```

Indexation par test

Le test peut tout à fait porter sur un autre vecteur.

```
R> taille <- c(151, 187, 178, 169)
```

```
R> sexe <- c("Homme", "Femme", "Homme", "Femme")
```

```
R> taille[sexe == "Homme"]
```

```
[1] 151 178
```

Indexation par test

Le test peut tout à fait porter sur un autre vecteur.

```
R> taille <- c(151, 187, 178, 169)
```

```
R> sexe <- c("Homme", "Femme", "Homme", "Femme")
```

```
R> taille[sexe == "Homme"]
```

```
[1] 151 178
```

```
R> sexe[taille <= 180]
```

```
[1] "Homme" "Homme" "Femme"
```

Indexation par test

Ça marche évidemment très bien sur les *data frame*.

```
   sexe age
1 homme  25
2 femme  37
3 femme  21
```

Indexation par test

Ça marche évidemment très bien sur les *data frame*.

```
  sexe age
1 homme 25
2 femme 37
3 femme 21
```

```
R> df[df$age < 30, ]
```

```
  sexe age
1 homme 25
3 femme 21
```

Indexation par test

Ça marche évidemment très bien sur les *data frame*.

```
  sexe age
1 homme 25
2 femme 37
3 femme 21
```

```
R> df[df$age < 30, ]
```

```
  sexe age
1 homme 25
3 femme 21
```

```
R> df[df$sexe == "homme", "age"]
```

```
[1] 25
```

Exemples sur les données

Recodage



Recodage

On a utilisé l'indexation pour extraire des données. Mais on peut l'utiliser pour remplacer les données sélectionnées.

Recodage

On a utilisé l'indexation pour extraire des données. Mais on peut l'utiliser pour remplacer les données sélectionnées.

```
R> x <- c(2, 5, 8, 11)
```

```
R> x[1]
```

```
[1] 2
```

Recodage

On a utilisé l'indexation pour extraire des données. Mais on peut l'utiliser pour remplacer les données sélectionnées.

```
R> x <- c(2, 5, 8, 11)
```

```
R> x[1]
```

```
[1] 2
```

```
R> x[1] <- 4
```

```
R> x
```

```
[1] 4 5 8 11
```

Recodage

On peut remplacer plusieurs valeurs d'un coup.

```
R> x <- c(2, 5, 8, 11)
```

Recodage

On peut remplacer plusieurs valeurs d'un coup.

```
R> x <- c(2, 5, 8, 11)
```

```
R> x[c(2, 3)] <- c(100, 200)
```

```
R> x
```

```
[1]  2 100 200 11
```

Recodage

On peut remplacer plusieurs valeurs d'un coup.

```
R> x <- c(2, 5, 8, 11)
```

```
R> x[c(2, 3)] <- c(100, 200)
```

```
R> x
```

```
[1] 2 100 200 11
```

```
R> x[1:3] <- 0
```

```
R> x
```

```
[1] 0 0 0 11
```

Recodage

C'est encore plus intéressant avec l'indexation par test.

```
R> x <- c(2, 5, 8, 11)
```

Recodage

C'est encore plus intéressant avec l'indexation par test.

```
R> x <- c(2, 5, 8, 11)
```

```
R> x[x == 5] <- 66
```

```
R> x
```

```
[1] 2 66 8 11
```

Recodage

C'est encore plus intéressant avec l'indexation par test.

```
R> x <- c(2, 5, 8, 11)
```

```
R> x[x == 5] <- 66
```

```
R> x
```

```
[1] 2 66 8 11
```

```
R> x[x < 10] <- 0
```

```
R> x
```

```
[1] 0 66 0 11
```

Recodage

Ça s'applique aux *data frame*.

```
  sexe age
1 homme 25
2 femme 37
3 femme 21
```

Recodage

Ça s'applique aux *data frame*.

```
  sexe age
1 homme 25
2 femme 37
3 femme 21
```

```
R> df$sexe[df$sexe == "femme"] <- "f"
```

Recodage

Ça s'applique aux *data frame*.

```
  sexe age
1 homme 25
2 femme 37
3 femme 21
```

```
R> df$sexe[df$sexe == "femme"] <- "f"
```

```
R> df
```

```
  sexe age
1 homme 25
2     f  37
3     f  21
```

Exemples sur les données

On sait désormais recoder une variable qualitative.

Problèmes :

- la syntaxe n'est pas forcément intuitive au début
- difficile de s'en souvenir si on ne l'utilise que ponctuellement
- sensible aux fautes de frappe

On sait désormais recoder une variable qualitative.

Problèmes :

- la syntaxe n'est pas forcément intuitive au début
- difficile de s'en souvenir si on ne l'utilise que ponctuellement
- sensible aux fautes de frappe

Fonction `irec`

irec

Facteurs



Facteurs

Un facteur (`factor`) est une variable comportant une liste définie de valeurs possibles.

```
R> x <- c("Rouge", "Vert", "Rouge", "Bleu")
```

Facteurs

Un facteur (`factor`) est une variable comportant une liste définie de valeurs possibles.

```
R> x <- c("Rouge", "Vert", "Rouge", "Bleu")
```

```
R> f <- factor(x)
```

```
R> f
```

```
[1] Rouge Vert  Rouge Bleu  
Levels: Bleu Rouge Vert
```

Facteurs

Un facteur (`factor`) est une variable comportant une liste définie de valeurs possibles.

```
R> x <- c("Rouge", "Vert", "Rouge", "Bleu")
```

```
R> f <- factor(x)
```

```
R> f
```

```
[1] Rouge Vert  Rouge Bleu
```

```
Levels: Bleu Rouge Vert
```

```
R> f[1] <- "Chihuahua"
```

```
Warning: invalid factor level, NA generated
```

Facteurs

Par défaut les variables qualitatives importées dans R sont souvent converties en facteurs.

C'est notamment le cas des fichiers importés avec `read.table`, `read.csv`...

Facteurs

Par défaut les variables qualitatives importées dans R sont souvent converties en facteurs.

C'est notamment le cas des fichiers importés avec `read.table`, `read.csv`...

Solutions :

- utiliser l'option `StringsAsFactors=FALSE` de `read.csv`
- convertir le facteur avec `as.character`
- `irec` gère ce cas de figure

Exemple avec les données

Facteurs

Un intérêt : les facteurs permettent d'ordonner les niveaux dans les tableaux (tri à plat, tri croisé).

Facteurs

Un intérêt : les facteurs permettent d'ordonner les niveaux dans les tableaux (tri à plat, tri croisé).

Par défaut les niveaux sont ordonnés par ordre alphabétique.

```
R> dipl <- c("Collège", "Bac", "Sup", "Bac", "Collège")  
R> freq(dipl)
```

	n	%
Bac	2	40
Collège	2	40
Sup	1	20
NA	0	0

Facteurs

factor permet d'ordonner les niveaux.

```
R> diplf <- factor(dipl, levels = c("Collège", "Bac",  
+   "Sup"))  
R> freq(diplf)
```

	n	%
Collège	2	40
Bac	2	40
Sup	1	20
NA	0	0

Facteurs

Même problèmes que pour les recodages :

- la syntaxe n'est pas forcément intuitive au début
- difficile de s'en souvenir si on ne l'utilise que ponctuellement
- sensible aux fautes de frappe

Facteurs

Même problèmes que pour les recodages :

- la syntaxe n'est pas forcément intuitive au début
- difficile de s'en souvenir si on ne l'utilise que ponctuellement
- sensible aux fautes de frappe

Fonction `iorder`

Découpage d'une variable numérique en classes



Découpage d'une variable numérique

On peut convertir une variable quantitative en variable qualitative avec la fonction `cut`.

```
R> cut(var, breaks, include.lowest = FALSE, right = TRUE)
```

Découpage d'une variable numérique

Même problèmes que pour les recodages :

- la syntaxe n'est pas forcément intuitive au début
- difficile de s'en souvenir si on ne l'utilise que ponctuellement
- sensible aux fautes de frappe

Découpage d'une variable numérique

Même problèmes que pour les recodages :

- la syntaxe n'est pas forcément intuitive au début
- difficile de s'en souvenir si on ne l'utilise que ponctuellement
- sensible aux fautes de frappe

Fonction `icut`

C'est fini !

