Notes de cours | Pivot avancé

February 2024

Introduction
Objectifs d'apprentissage
Packages
Jeux de données
Du format large au format long
Comprendre names sep et ".value"
Type de valeur <i>avant</i> le séparateur
Un exemple qui n'est pas une série temporelle
Echapper le séparateur de point
Que faire quand vous n'avez pas un séparateur net ?
Du format long au format large
Bilan!
Références
Solutions des exercices

Introduction

Vous connaissez les opérations de pivot de base des jeux de données du format long au format large et vice versa. Cependant, comme c'est souvent le cas, les manipulations de base ne sont pas suffisantes pour le traitement des données que vous devez faire. Voyons maintenant le niveau suivant. Allons-y!

Objectifs d'apprentissage

- 1. Maîtriser le pivot complexe du format large au format long et du format long au format large
- 2. Savoir utiliser les séparateurs comme outil de pivot

Packages

```
r les packages
lire(pacman)) install.packages("pacman")
p_load(tidyverse, outbreaks, janitor, rio, here, knitr)
```

Jeux de données

Nous présenterons ces jeux de données au fur et à mesure, mais voici un aperçu :

- Données d'enquête d'une étude menée en Inde sur les dépenses des patients pour le traitement de la tuberculose
- Données d'une étude sur les biomarqueurs des entéropathogènes en Zambie
- Une enquête alimentaire au Vietnam

Du format large au format long

Parfois, vous avez plusieurs types de données au format large dans le même jeu de données. Considérez cet exemple factice de la taille et du poids des enfants sur deux ans :

Si vous pivotez toutes les colonnes des mesures, vous obtiendrez des données trop longues :

```
stats %>%
longer(2:5)
```

```
## # A tibble: 5 x 3
## enfant name value
## <chr> <chr> <chr> ## 1 A annee1_taille 80cm
## 2 A annee2_taille 85cm
## 3 A annee1_poids 5kg
```

```
## 4 A annee2_poids 10kg
## 5 B annee1 taille 85cm
```

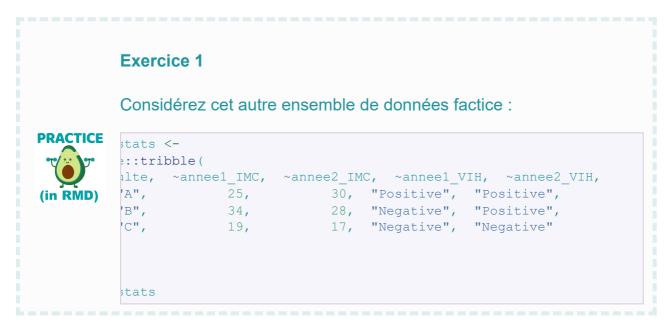
Ce n'est (généralement) pas ce que nous recherchons, car maintenant vous avez deux données différentes dans la même colonne- le poids et la taille.

Pour obtenir le bon format, vous devez utiliser l'argument names_sep et l'identifiant ".value" :

```
## # A tibble: 5 x 4
## enfant periode taille poids
## <chr> <chr> <chr> <chr> ## 1 A annee1 80cm 5kg
## 2 A annee2 85cm 10kg
## 3 B annee1 85cm 7kg
## 4 B annee2 90cm 12kg
## 5 C annee1 90cm 6kg
```

Maintenant, nous avons une ligne pour chaque combinaison enfant-période, un format long correct !

Ce que fait le code ci-dessus peut ne pas être clair, mais vous devriez déjà pouvoir répondre à l'exercice ci-dessous en reproduisant la syntaxe de l'exemple précédant. Après cet exercice, nous expliquerons l'argument names_sep et l'identifiant ".value" plus en détail.



```
## # A tibble: 3 × 5
   adulte anneel IMC annee2 IMC annee1 VIH
              <dbl> <dbl> <chr>
                            30 Positive
## 1 A
## 2 B
                   34
                             28 Negative
## 3 C
                  19
                             17 Negative
## annee2 VIH
##
   <chr>
## 1 Positive
## 2 Positive
## 3 Negative
```



Pivotez les données en un format long pour obtenir la structure suivante :

adulte annee IMC VIH

```
stats %>%
_longer(_____)
```

L'exemple ci-dessus <code>enfant_stats</code> a des nombres stockés en tant que caractères [...]

Comme vous l'avez vu dans la leçon précédente, vous pouvez facilement extraire les nombres à partir du jeux de données de sortie au format long en utilisant la fonction parse_number() de readr:

```
SIDE NOTE
```

```
stats_long <-
stats %>%
longer(2:5,
          names_sep = "_",
          names_to = c("periode", ".value"))
stats_long
```

```
## # A tibble: 5 × 4
   enfant periode taille poids
   <chr> <chr> <chr> <chr>
           anneel 80cm
## 1 A
## 2 A
           annee2 85cm
                         10kg
## 3 B
           anneel 85cm
                         7kg
          annee2 90cm
## 4 B
                         12kg
## 5 C
          anneel 90cm
```

```
stats long %>%
                                                                      п
         (taille = parse number(taille),
                                                                      ī
           poids = parse number(poids))
SIDE NOTE
           ## # A tibble: 5 × 4
 ......
           ## enfant periode taille poids
           ## <chr> <chr> <dbl> <dbl>
           ## 1 A
                     annee1
                               80 5
           ## 2 A
                                      10
                      annee2
                                 85
                     annee1
           ## 3 B
           ## 4 B
                                90
                                     12
                     annee2
           ## 5 C
                     annee1
                                90
```

Comprendre names sep et ".value"

Maintenant, décomposons l'appel pivot longer () que nous avons vu ci-dessus :

```
tats
```

```
## # A tibble: 3 \times 5
## enfant annee1 taille annee2 taille
## <chr> <chr>
                       <chr>
## 1 A
          80cm
                        85cm
## 2 B
           85cm
                        90cm
## 3 C
          90cm
                        100cm
## anneel poids anneel poids
## <chr>
                <chr>
## 1 5kg
                 10kg
## 2 7kg
                12kg
## 3 6kg
                14kg
```

```
## # A tibble: 5 × 4
## enfant periode taille poids
## <chr> <chr> <chr> <chr>
          anneel 80cm
## 1 A
                      5kg
## 2 A
         annee2 85cm
                       10 kg
         anneel 85cm
## 3 B
                      7kg
## 4 B
         annee2 90cm 12kg
## 5 C
         anneel 90cm 6kg
```

Remarquez que les noms de colonnes dans le dataframe enfant_stats d'origine (anneel taille, anneel taille etc.) sont composés de trois parties :

- la période référencée : par exemple "annee1"
- un séparateur de soulignement, "_";
- et le type de valeur enregistrée "taille" ou "poids"

Nous pouvons faire un tableau avec ces parties :

```
nom_colonneperiodeseparateur".value"annee1_tailleannee1tailleannee2_tailleannee2tailleannee1_poidsannee1poidsannee2_poidsannee2poids
```

Sur la base de ce tableau, il devrait maintenant être plus facile de comprendre les arguments names sep et names to que nous avons fournis à pivot longer():

```
names sep = " ":
```

C'est le séparateur entre l'indicateur de période (année) et les valeurs (taille et poids) enregistrées.

Si nous utilisons un séparateur différent, l'argument va aussi changer. Par exemple, si le séparateur est un espace vide, " ", vous aurez <code>names_sep = " ", comme on le voit dans l'exemple ci-dessous :</code>

```
stats espace sep <-
:::tribble(
fant, ~`ann1 taille`, ~`ann2 taille`, ~`ann1 poids`, ~`ann2 poids`,
'A", "80cm", "85cm", "5kg", "10kg",
                      "90cm",
         "85cm", "90cm",
"90cm", "100cm",
                                     "7kg",
                                                  "12kg",
'B",
                                     "6kg",
'C",
                                                  "14kg"
stats espace sep %>%
longer(2:5,
      names sep = " ",
      names to = c("periode", ".value"))
```

```
## # A tibble: 5 × 4
## enfant periode taille poids
## <chr> <chr> <chr> <chr> defant chr> <chr> <chr> ## 1 A ann1 80cm 5kg
## 2 A ann2 85cm 10kg
## 3 B ann1 85cm 7kg
```

```
## 4 B ann2 90cm 12kg
## 5 C ann1 90cm 6kg

names to = c("periode", ".value")
```

Ensuite, l'argument names_to indique comment les données doivent être restructurées. Nous avons passé un vecteur de deux chaînes de caractères, "periode" et ".value" à cet argument. Voyons le rôle de chaque élément :

La chaîne "periode" indique que nous voulons placer les données de chaque année (ou période) dans une ligne séparée. Notez qu'il n'y a rien de spécial dans le mot "periode" utilisé ici ; nous pourrions changer cela par n'importe quelle autre chaîne. Donc, au lieu de "periode", vous auriez pu écrire "temps" ou "annee_de_mesure" ou autre chose :

```
## # A tibble: 5 × 4
## enfant annee_de_mesure taille poids
## <chr> <chr> <chr> ## 1 A annee1 80cm 5kg
## 2 A annee2 85cm 10kg
## 3 B annee1 85cm 7kg
## 4 B annee2 90cm 12kg
## 5 C annee1 90cm 6kg
```

Maintenant, **le placeholder ".value"** est un indicateur spécial, qui indique à pivot_longer() de créer une colonne séparée pour chaque valeur distincte qui apparaît après le séparateur. Dans notre exemple, ces valeurs sont "taille" et "poids".

La chaîne ".value" ne peut pas être remplacée arbitrairement. Par exemple, ceci ne fonctionnera pas :

```
## # A tibble: 5 × 4
## enfant periode valeurs value
## <chr> <chr> <chr> <chr> ## 1 A annee1 taille 80cm
## 2 A annee2 taille 85cm
## 3 A annee1 poids 5kg
```

```
## 4 A annee2 poids 10kg
## 5 B annee1 taille 85cm
```

Autrement dit, le placeholder ".value" indique à pivot_longer() que nous voulons séparer les valeurs "taille" et "poids" dans deux colonnes séparées, car nous avons deux types de valeurs après le séparateur "_" dans les noms de colonnes.

Cela signifie que si vous avez un jeu de données au format large avec trois types de valeurs, vous obtiendrez trois colonnes séparées, une pour chaque type de valeur. Par exemple, considérez le jeu de données fictif ci-dessous qui montre les enregistrements d'enfants, à deux moments, pour les variables suivantes :

- âge en mois,
- pourcentage de graisse corporelle
- IMC

```
stats_trois_valeurs <-
stats_trois_valeurs <-
stats_trois_valeurs <-
stats_trois_valeurs <-
stats_trois_valeurs <-
stats_trois_valeurs <--
stats_
```

Ici, dans les noms de colonnes, il y a trois types de valeurs qui apparaissent après le séparateur "_" : age, graisse et imc; la chaîne ".value" indique à pivot_longer() de créer une nouvelle colonne pour chaque type de valeur :

```
## # A tibble: 4 × 5
## enfant temps age graisse imc
## <chr> <chr> <chr> <chr> <chr> <chr> 13%
14
```

```
## 2 a t2 8 mois 15% 15
## 3 b t1 7 mois 15% 16
## 4 b t2 9 mois 17% 18
```

Exercice 2

Un pédiatre enregistre les informations suivantes pour un ensemble d'enfants sur deux ans :

- périmètre cranien ;
- circonférence du cou ; et
- tour de hanches

le tout en centimètres.

Voici le tableau de sortie :

```
ice stats <-
:::tribble(
ant, ~ann1 tete, ~ann2 tete, ~ann1 cou, ~ann2 cou, ~ann1 hanche, ~ann2 l
    45, 48, 23, 24, 51,
52,
'b",
         48,
                 50,
                         24,
                                 26,
                                           52,
52,
'C",
         50, 52,
                         24,
                               27,
                                            53,
54
ice stats
```

```
## # A tibble: 3 × 7
## enfant ann1_tete ann2_tete ann1_cou ann2_cou
## <chr> <dbl> <dbl> <dbl> <dbl>
                     48
## 1 a
            45
                            23
## 2 b 48 50
## 3 c 50 52
                              24
                                      26
                              24
                                      27
## ann1 hanche ann2 hanche
## <dbl> <dbl>
        51
                  52
## 1
## 2
          52
                     52
## 3
          53
                     54
```

Pivotez les données en un format long pour obtenir la structure suivante :

PRACTICE

(in RMD)



Type de valeur avant le séparateur

Dans tous les exemples que nous avons utilisés jusqu'à présent, les noms de colonnes étaient construits de telle sorte que le type de valeur venait après le séparateur. Rappelez-vous notre tableau :

nom_colonne	periode	separateur	".value"
anneel_taille	annee1	_	taille
annee2_taille	annee2	_	taille
anneel_poids	annee1	_	poids
annee2_poids	annee2	_	poids

Mais bien sûr, les noms de colonnes pourraient être construits différemment, avec les types de valeurs venant avant le séparateur, comme dans cet exemple :

```
## # A tibble: 3 \times 5
## enfant taille anneel taille annee2
##
    <chr> <chr> <chr>
## 1 A 80cm
## 2 B 85cm
## 3 C 90cm
                         85cm
          85cm
                         90cm
          90cm
                         100cm
## poids anneel poids annee2
## <chr>
                 <chr>
## 1 5kg
                 10kg
## 2 7kg
                12kg
## 3 6kg
                 14kg
```

Ici, les types de valeurs (taille et poids) viennent avant le "_" séparateur.

Comment notre commande pivot_longer() peut-elle s'adapter à cela ? C'est simple! Il suffit d'inverser l'ordre du vecteur donné à l'argument names to:

```
Donc, au lieu de names_to = c("temps", ".value"), vous aurez names_to =
c(".value", "temps"):
```

```
## # A tibble: 5 × 4

## cenfant temps taille poids

## < <chr> <chr> <chr> <chr> <chr> ## 1 A annee1 80cm 5kg

## 2 A annee2 85cm 10kg

## 3 B annee1 85cm 7kg

## 4 B annee2 90cm 12kg

## 5 C annee1 90cm 6kg
```

Et voilà!

Exercice 3

Considérez le jeu de données suivant de la Zambie sur les entéropathogènes et leurs biomarqueurs.

```
thogenes_zambie_large<-
read_csv(here("data/fr_enteropathogenes_zambie_large.csv"))
thogenes_zambie_large</pre>
```

```
## # A tibble: 5 × 7
## ID LPS 1 LPS 2 LBP 1 LBP 2 IFABP 1
  <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 1002 222. 390. 38414. 6840. 1294.
## 2 1003 181. NA 26888. NA
## 3 1004 257. 221. 49183. 5426.
## 4 1005 NA 369. NA 1938.
## 5 1006 275. NA 61758. NA
## IFABP 2
##
## 1
     610.
## 2
      NA
## 3
       0
## 4 1010.
## 5 NA
```



Ce jeu de données se compose des colonnes suivantes :

- LPS_1 et LPS_2 : niveau des lipopolysaccharides, mesuré par Pyrochrome LAL, en EU/mL
- LBP_1 et LBP_2 : niveau des protéines de liaison au LPS, en pg/mL
- IFABP_1 et IFAPB_2 : niveau des protéines de liaison aux acides gras de type intestinal, en pg/mL

Pivotez le jeu de données pour qu'il ressemble à la structure suivante .

ID numero echantillon LPS LBP IFABP

```
PRACTICE | thogenes_zambie_large %>% | longer(_____)
```

Un exemple qui n'est pas une série temporelle

Jusqu'à présent, nous avons utilisé des ensembles de données personne-période (séries temporelles) pour illustrer l'idée de pivots complexes avec plusieurs types de valeurs.

Mais comme nous l'avons mentionné, tous les jeux de données nécessitant une restructuration ne sont pas forcément des données de séries temporelles. Voyons un exemple rapide qui n'est pas une série temporelle.

Vous pourriez mesurer la taille (cm) et le poids (kg) d'une série de couples parentaux dans un tableau comme celui-ci :

```
## # A tibble: 3 × 5
## couple pere_taille pere_poids mere_taille
80
                   90
                          150
           185
## 2 b
## 3 c
           182
                  93
                         143
## mere_poids
## <dbl>
## 1
## 2
       76
## 3
       78
```

lci, nous avons deux types de valeurs différents (poids et taille) pour chaque personne du couple.

Pour pivoter à une ligne par personne, nous aurons encore besoin des arguments names sep et names to:

```
## # A tibble: 5 × 4
## couple personne taille poids
## <chr> <chr> <chr> <chr> <chr> <dbl> <dbl> <80
## 1 a pere 180 80
## 2 a mere 160 70
## 3 b pere 185 90
## 4 b mere 150 76
## 5 c pere 182 93</pre>
```

Le séparateur est un trait de soulignement, "_", donc nous avons utilisé names_sep = "_" et comme les types de valeurs viennent après le séparateur, l'identifiant ".value" a été placé en deuxième dans l'argument names to.

Echapper le séparateur de point

Un cas spécial que vous pourriez rencontrer est un ensemble de données où le séparateur est un point.

```
## # A tibble: 5 x 4
## enfant periode taille poids
## <chr> <chr> <chr> <chr> ## 1 A annee1 80cm 5kg
## 2 A annee2 85cm 10kg
## 3 B annee1 85cm 7kg
## 4 B annee2 90cm 12kg
## 5 C annee1 90cm 6kg
```

Ici, nous avons utilisé la chaîne "\." pour indiquer un point "." parce que le "." est un caractère spécial dans R qui dans certains cas doit être échappé.



Considérez à nouveau les données adulte_stats que vous avez vues ci-dessus. Maintenant, les noms des colonnes ont été légèrement modifiés.



```
## # A tibble: 3 × 5
## adulte IMC.annee1 IMC.annee2 VIH.annee1
## <chr> <dbl> <dbl> <chr>
               25 30 Positive
34 28 M
              25
## 1 A
## 2 B
## 3 C
                        17 Negative
                19
## VIH.annee2
  <chr>
##
## 1 Positive
## 2 Positive
## 3 Negative
```

Encore une fois, pivotez les données en un format long pour obtenir la structure suivante :



```
e_stats_point_sep %>%
longer(_____)
```

Que faire quand vous n'avez pas un séparateur net ?

Parfois, vous n'avez pas de séparateur net.

Considérez ces données d'une enquête menée en Inde qui examine les dépenses des patients pour le traitement de la tuberculose :

```
ces <- read_csv(here("data/fr_india_tb_pathways_and_costs_data.csv")) %>%
names() %>%
(id, premiere_visite_emplacement, premiere_visite_cout,
deuxieme_visite_emplacement, deuxieme_visite_cout,
troisieme_visite_emplacement, troisieme_visite_cout)

ces
```

```
## # A tibble: 5 \times 7
## id premiere visite em...¹ premiere visite...²
##
    <dbl> <chr>
                                             <dbl>
## 1 100202 GH
## 2 100396 Pvt. docto
                                              1500
## 3 100590 Pvt. docto
                                              2000
## 4 100687 Pvt. hospi
                                              20000
## 5 100784 Pvt. docto
                                              1000
## deuxieme visite emplacement deuxieme visite...3
## <chr>
## 1 <NA>
## 2 Pvt. clini
                                              1000
## 3 Pvt. docto
                                              3000
## 4 Pvt. hospi
                                              1500
## 5 GH
## # i abbreviated names: ¹premiere visite emplacement,
## # 2premiere visite cout, 3deuxieme visite cout
```

Il n'y a pas de séparateur net entre les indicateurs de temps (premier, deuxième, troisième) et le type de valeur (cout, emplacement). C'est-à-dire, au lieu de "premierevisite_emplacement", nous avons plutôt "premiere_visite_emplacement", donc le trait de soulignement est utilisé pour deux buts. Pour cette raison, si vous essayez notre stratégie de pivot habituelle, vous obtiendrez un message d'erreur :

La façon la plus directe de restructurer ce jeu de données avec succès serait d'utiliser un "regex" spécial (manipulation de chaînes de caractères), mais il est probable que vous n'ayez pas encore appris cela!

Alors pour l'instant, la solution que nous recommandons est de renommer manuellement vos colonnes pour insérer un séparateur clair, "___" :

```
:es_renomme <-
sites %>%
(premiere_visite_emplacement = premiere_visite_emplacement,
   premiere_visite_cout = premiere_visite_cout,
   deuxieme_visite_emplacement = deuxieme_visite_emplacement,
   deuxieme_visite_cout = deuxieme_visite_cout,
   troisieme_visite_emplacement = troisieme_visite_emplacement,
   troisieme_visite_cout = troisieme_visite_cout)

:es_renomme
```

```
## # A tibble: 5 × 7
     id premiere visite e...¹ premiere visit...²
## <dbl> <chr>
                                             <db1>
## 1 100202 GH
## 2 100396 Pvt. docto
                                              1500
## 3 100590 Pvt. docto
                                              2000
## 4 100687 Pvt. hospi
                                             20000
## 5 100784 Pvt. docto
                                              1000
   deuxieme visite emplacem...3 deuxieme visit...4
## <chr>
## 1 <NA>
                                              1000
## 2 Pvt. clini
## 3 Pvt. docto
                                              3000
## 4 Pvt. hospi
                                              1500
## 5 GH
## # i abbreviated names: ¹premiere visite emplacement,
## # 2premiere visite cout, ...
```

Maintenant, nous pouvons essayer le pivot :

```
## # A tibble: 5 \times 4
##
       id numero_visite visite_emplacement
##
    <dbl> <chr>
                 <chr>
## 1 100202 premiere
                       GH
## 2 100202 deuxieme
                       <NA>
## 3 100202 troisieme
                        <NA>
## 4 100396 premiere
                      Pvt. docto
## 5 100396 deuxieme
                      Pvt. clini
## visite cout
##
          <dbl>
## 1
## 2
             0
## 3
             0
```

```
## 4 1500
## 5 1000
```

Maintenant, nettoyons le jeu de données :

```
## # A tibble: 5 × 4
## id numero visite visite emplacement
## <dbl> <dbl> <chr>
## 1 100202
                   1 GH
                   1 Pvt. docto
## 2 100396
                   2 Pvt. clini
## 3 100396
## 4 100396
                   3 Pvt. hospi
## 5 100590
                   1 Pvt. docto
## visite cout
     <dbl>
##
## 1
## 2
        1500
## 3
        1000
## 4
        2500
## 5
        2000
```

lci, nous avons d'abord supprimé les observations où nous n'avons pas d'information sur l'emplacement de la visite (c'est-à-dire que nous filtrons les lignes où la variable d'emplacement de la visite est définie à ""). Nous convertissons ensuite en valeurs numériques la variable du numero de la visite, où les chaînes "premiere" à "troisieme" sont converties en valeurs numériques 1 à 3. Enfin, nous nous assurons que la variable du coût de la visite est numérique en utilisant mutate() et la fonction d'aide as.numeric().

Exercice 5



Nous allons utiliser les données d'une enquête alimentaire au Vietnam. Des femmes de Hanoi ont été interrogées sur leurs achats alimentaires, et les données collectées ont servi à créer un profil nutritionnel de chaque femme. Ici, nous utiliserons un sous-ensemble de ces données de 61 ménages qui sont venus pour 2 visites, enregistrant :

- enerc_kcal_s_1 : l'apport énergétique de l'ingrédient/nourriture (Kcal) lors de la première visite (_2 pour la deuxième visite)
- sec_s_1 : l'apport sec de l'ingrédient/nourriture (g) lors de la première visite (2 pour la deuxième visite)
- eau_s_1 : l'apport en eau de l'ingrédient/nourriture (g) lors de la première visite (2 pour la deuxième visite)
- graisse_s_1 : l'apport en lipides de l'ingrédient/nourriture (g) lors de la première visite (2 pour la deuxième visite)

```
e_alimentaire_vietnam_large <-
read_csv(here("data/fr_diet_diversity_vietnam_wide.csv"))

e_alimentaire_vietnam_large</pre>
```



```
## # A tibble: 5 × 10
## ...1 menage id enerc kcal s 1 enerc kcal s 2
## <dbl> <dbl> <dbl> <dbl>
## 1 1
            348
                       2268.
      2
                        2775.
## 2
             354
                                    1240.
## 2 2 354
## 3 3 53
## 4 4 18
## 5 5 211
                        3104.
                                     2075.
                        2802.
                                     2146.
                       1298.
                                    1191.
## sec s 1 sec s 2 eau s 1 eau s 2 graisse s 1
## <dbl> <dbl> <dbl> <dbl> <dbl>
     548. 281. 4219. 1997.
## 1
     600.
            284. 2376. 3145.
                                   115.
## 3 646.
            451. 2808. 2305.
                                  127.
## 4 620.
            807. 3457. 1903.
                                   87.4
                        2269.
## 5 269.
            288.
                  2584.
                                    47.8
## # i 1 more variable: graisse s 2 <dbl>
```

Vous devrez d'abord vérifier si vous avez un opérateur net et renommer vos colonnes si nécessaire. Ensuite, rassemblez les données enregistrées sur les deux visite dans une colonne par type d'apport (énergétique, lipides, eau et poids sec). En d'autres termes, pivotez le jeu de données en un format long de cette forme :

menage_id visite enerc_kcal_s sec_s eau_s graisse_s

```
e_alimentaire_vietnam_large %>%
longer(_____)
```

Du format long au format large

Nous venons de voir comment effectuer certaines opérations complexes du format large au format long, qui, comme nous l'avons vu dans la leçon précédente, sont essentielles pour tracer et manipuler les données. Passons maintenant à la transformation inverse.

Il peut être utile de passer du format long au format large pour transformer et filtrer les données ou encore pour traiter des valeurs manquantes (NA). Dans ce format, vos mesures / données collectées deviennent les colonnes du jeu de données.

Cette fois-ci, nous allons utiliser le jeux de données originel sur les entéropathogènes en Zambie. En effet, ce que vous manipuliez jusqu'à présent était un jeu de données **préparé pour vous**, en format large. **Le jeu de données originel est au format long** et nous allons maintenant voir la préparation des données que j'ai faite au préalable, en coulisses. Vous êtes presque en train de devenir l'enseignant de cette leçon ;)

```
thogenes_zambie_long <-
read_csv(here("data/fr_enteropathogenes_zambie_long.csv"))
thogenes_zambie_long</pre>
```

```
## # A tibble: 5 × 5

## ID group LPS LBP IFABP

## 2 <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> 

## 1 1002 1 222. 38414. 1294.

## 2 1002 2 390. 6840. 610.

## 3 1003 1 181. 26888. 22.5

## 4 1004 2 221. 5426. 0

## 5 1004 1 257. 49183. 0
```

Voici comment nous le convertissons du format long au format large :

```
athogenes_zambie_large <-
pathogenes_zambie_long %>%
wider(
ss_from = group,
les_from = c(LPS, LBP, IFABP)
athogenes_zambie_large
```

```
## # A tibble: 5 × 7
## ID LPS_1 LPS_2 LBP_1 LBP_2 IFABP_1
## <dbl> = 400. 1294.
## 2 1003 181. NA 26888. NA 22.5
## 3 1004 257. 221. 49183. 5426.
```

```
## 4 1005 NA 369. NA 1938. 0
## 5 1006 275. NA 61758. NA 0
## IFABP_2
## <dbl>
## 1 610.
## 2 NA
## 3 0
## 4 1010.
## 5 NA
```

Vous pouvez voir que les valeurs de la variable group (1 ou 2) sont ajoutées aux noms des valeurs (LPS, LBP, IFABP) pour créer les nouvelles colonnes représentant différents groupes de données : par exemple, LPS 1 et LPS 2.

Nous considérons que c'est une option "avancée" du pivot car nous pivotons plusieurs variables en même temps, mais comme vous pouvez le voir, la syntaxe est assez simple. Nous utilisons les mêmes arguments <code>names_from</code> et <code>values_from</code> qu'avec les pivots plus simples que nous avons vus dans la leçon précédente.

Voyons un autre exemple, en utilisant les données de l'enquête alimentaire du Vietnam que vous avez manipulées précédemment :

```
ce_alimentaire_vietnam_long <-
read_csv(here("data/fr_diet_diversity_vietnam_long.csv"))
ce_alimentaire_vietnam_long</pre>
```

```
## # A tibble: 5 × 6
## numero visite menage id enerc kcal s sec s
1 348
1 354
## 1
                               2268. 548.
                               2775. 600.
## 2
              1
                      53
## 3
                               3104. 646.
                      18
              1
## 4
                               2802. 620.
## 5
                     211
                               1298. 269.
              1
## eau_s graisse_s
## <dbl> <dbl>
## 1 4219.
## 1 4219. 78.4

## 2 2376. 115.

## 3 2808. 127.

## 4 3457. 87.4

## 5 2584. 47.8
             78.4
```

lci, nous allons utiliser la variable numero_visite pour créer une nouvelle variable pour les différents apports enregistrés lors des deux visites :

```
te_alimentaire_vietnam_large <-
site_alimentaire_vietnam_long %>%
wider(
es_from = numero_visite,
ses_from = c(enerc_kcal_s, sec_s, eau_s, graisse_s)

te_alimentaire_vietnam_large
```

```
## # A tibble: 5 × 9
## menage id enerc kcal s 1 enerc kcal s 2
## <dbl> <dbl> <dbl>
## 1
        348
                   2268.
## 2
        354
                    2775.
                                 1240.
         53
                    3104.
## 3
                                  2075.
        18
211
                    2802.
## 4
                                  2146.
                    1298.
## 5
                                  1191.
## sec s 1 sec s 2 eau s 1 eau s 2 graisse s 1
    <dbl> <dbl> <dbl> <dbl> <dbl> <
##
## 1 548. 281. 4219. 1997.
## 2 600. 284. 2376. 3145.
                                      78.4
                                     115.
## 3 646. 451. 2808. 2305.
                                     127.
## 4 620. 807. 3457. 1903.
## 5 269. 288. 2584. 2269.
                                      87.4
                                      47.8
## # i 1 more variable: graisse s 2 <dbl>
```

Vous pouvez voir que les valeurs de la variable <code>numero_visite</code> (1 ou 2) sont ajoutées aux noms des valeurs (<code>enerc_kcal_s</code>, <code>sec_s</code>, <code>graisse_s</code>, <code>eau_s</code>) pour créer les nouvelles colonnes représentant différents groupes de données : par exemple, <code>eau_s_1</code> et <code>eau_s_2</code>. Nous avons pivoté en format large toutes ces variables en même temps. Maintenant, chaque mesure de l'apport par visite est représentée comme une seule variable (c'est-à-dire une colonne) dans le jeu de données.

Avec ce format, il est facile de faire la somme de l'apport énergétique par ménage par exemple :

```
e_alimentaire_vietnam_large %>%
(menage_id, enerc_kcal_s_1, enerc_kcal_s_2) %>%
(energie_totale_kcal = enerc_kcal_s_1 + enerc_kcal_s_2) %>%
ge (menage_id)
```

```
## # A tibble: 5 \times 4
## menage id enerc kcal s 1 enerc kcal s 2
   ##
## 1
       14
                1040.
                           1663.
       17
## 2
                2100.
                           1286.
## 3
        18
                2802.
                            2146.
## 4
        22
                 3187.
                            1582.
## 5 24
                2359.
                            2026.
## energie_totale_kcal
##
             <dbl>
```

```
## 1 2704.

## 2 3386.

## 3 4948.

## 4 4769.

## 5 4385.
```

Cependant, vous pourriez obtenir un résultat similaire avec le format long :

```
e_alimentaire_vietnam_long %>%
by(menage_id) %>%
ize(energie_totale = sum(enerc_kcal_s))
```

```
## # A tibble: 5 \times 2
## menage id energie totale
##
    ## 2
## 2
                 2704.
                 3386.
## 3
       18
                 4948.
## 4
        22
                 4769.
## 5
        24
                 4385.
```

Exercice 6



Prenez le jeu de données tb_visites_long que nous avons manipulé plus haut et pivotez-le à nouveau au format large.

```
es_long %>%
wider(_____)
```

Bilan!

Vos compétences en manipulation de données viennent d'être renforcées avec le pivot avancé. Cette compétence s'avérera souvent essentielle lors de la manipulation des données du monde réel. Je ne doute pas que vous la mettrez bientôt en pratique. Elle est également essentielle, comme nous l'avons vu, pour la conception des graphiques. J'espère donc que le pivot vous sera utile non seulement pour votre manipulation de données, mais aussi pour pour la conception des graphiques.

Contributeurs

Les membres suivants de l'équipe ont contribué à cette leçon :



KENE DAVID NWOSU

Data analyst, the GRAPH Network Passionate about world improvement



LAURE VANCAUWENBERGHE

Data analyst, the GRAPH Network A firm believer in science for good, striving to ally programming, health and education



CAMILLE BEATRICE VALERA

Project Manager and Scientific Collaborator, The GRAPH Network



IMANE BENSOUDA KORACHI

R Developer and Instructor, the GRAPH Network

Références

Solutions des exercices

Exercice 1

```
## 4 B annee2 28 Positive
## 5 C annee1 19 Negative
```

Exercice 2

Exercice 3

Exercice 4

```
## 4 B annee2 28 Positive
## 5 C annee1 19 Negative
```

Exercice 5

```
te_alimentaire_vietnam_large%>%
ename(
enerc_kcal_s__1 = enerc_kcal_s_1,
enerc_kcal_s__2 = enerc_kcal_s_2,
sec_s__1 = sec_s_1,
sec_s__2 = sec_s_2,
eau_s__1 = eau_s_1,
eau_s__2 = eau_s_2,
graisse_s__1 = graisse_s_1,
graisse_s__2 = graisse_s_2
%>%
.vot_longer(2:9, names_sep = "__", names_to = c(".value", "visite"))
```

```
## # A tibble: 5 × 6
## menage id visite enerc kcal s sec s eau s
##
    <dbl> <dbl> <dbl> <dbl> <dbl>
                        2268. 548. 4219.
## 1
       348 1
## 2
        348 2
                       1386. 281. 1997.
## 3
        354 1
                       2775. 600. 2376.
                       1240. 284. 3145.
## 4
        354 2
## 5
                        3104. 646. 2808.
         53 1
## graisse_s
    <dbl>
##
## 1
       78.4
## 2
       67.7
## 3
      115.
## 4
       45.3
## 5
      127.
```

Exercice 6

```
## # A tibble: 5 × 7
## id visite emplacement...¹ visite emplacem...²
## 1 100202 GH
                            <NA>
## 2 100396 Pvt. docto
                           Pvt. clini
                           Pvt. docto
## 3 100590 Pvt. docto
## 4 100687 Pvt. hospi
                           Pvt. hospi
## 5 100784 Pvt. docto
                           GH
  visite emplacement troisi...3 visite cout pre...4
## <chr>
                                       <dbl>
## 1 <NA>
```

```
## 2 Pvt. hospi 1500
## 3 PHC 2000
## 4 PHC 20000
## 5 <NA> 1000
## # i abbreviated names: 'visite_emplacement_premiere,
## # 'visite_emplacement_emplacement_emplacement, ...
```