

Pivoting

2023 年 6 月 18 日

0.1 Introduction

このヴィネットでは、新しい関数 `pivot_longer()` と `pivot_wider()` の使用法を説明します。その目的は `gather()` と `spread()` の使い勝手を向上させ、他のパッケージで見られる最先端の機能を取り入れることにあります。

以前から、`spread()` と `gather()` のデザインには何か根本的な問題があることは明らかでした。多くの人がこの名前を直感的に理解できず、どの方向が広がり、どの方向が集まるかを覚えるのが難しいのです。また、これらの関数の引数を覚えるのも意外と大変なようで、多くの人が(私も含めて!)毎回ドキュメントを参照しなければならないのです。

R でリシェイプを進めている他の R パッケージからヒントを得た重要な新機能が 2 つあります：

- `pivot_longer()` は、Matt Dowle と Arun Srinivasan による `data.table` パッケージが提供する拡張 `melt()` と `dcast()` 関数に触発されて、異なる型を持つ複数の値変数を扱うことができます。
- `pivot_longer()` と `pivot_wider()` は、John Mount と Nina Zumel による `cdata` パッケージに触発されて、列名に格納されたメタデータがデータ変数になる(あるいはその逆)方法を正確に指定するデータフレームを受け取ることができます。

このビネットでは、`pivot_longer()` と `pivot_wider()` を使用して、単純なものから複雑なものまで、様々なデータ再形成の課題を解決する様子をご覧くださいながら、そのキーアイデアを学びます。

まず始めに、必要なパッケージをいくつかロードします。実際の解析コードでは、`library(tidyverse)` で行うかと思いますが、このヴィネットはパッケージに埋め込まれているので、ここではそれができないのです。

```
1 library(tidyr)
2 library(dplyr)
3 library(readr)
```

0.2 Longer

`pivot_longer()` は、行数を増やし、列数を減らすことで、データセットを長くする。データセットが「長い形」であると表現することに意味があるとは思いません。長さは相対的な用語であり、(例えば)データセット A

はデータセット B より長いとしか言いようがありません。

自然界に存在するデータセットは、分析のしやすさよりもデータ入力のしやすさや比較のしやすさに最適化されていることが多いため、`pivot_longer()` は整頓によく用いられます。以下のセクションでは、現実的な幅広いデータセットに対して `pivot_longer()` を使用する方法を紹介します。

0.2.1 String data in column names

`relig_income` データセットには、宗教と年収を尋ねたアンケートに基づくカウントが格納されています：

```
1 relig_income
2 #> # A tibble: 18 x 11
3 #>   religion `<$10k` `<$10-20k` `<$20-30k` `<$30-40k` `<$40-50k` `<$50-75k` `<$75-100k`
4 #>   <chr>      <dbl>      <dbl>      <dbl>      <dbl>      <dbl>      <dbl>      <dbl>
5 #> 1 Agnostic      27        34        60        81        76       137       122
6 #> 2 Atheist       12        27        37        52        35        70        73
7 #> 3 Buddhist      27        21        30        34        33        58        62
8 #> 4 Catholic     418       617       732       670       638     1116     949
9 #> 5 Don't ~       15        14        15        11        10        35        21
10 #> 6 Evangel~     575      869     1064     982     881     1486     949
11 #> 7 Hindu         1         9         7         9        11        34        47
12 #> 8 Histori~     228      244      236      238      197      223      131
13 #> 9 Jehovah~     20        27        24        24        21        30        15
14 #> 10 Jewish      19        19        25        25        30        95        69
15 #> # i 8 more rows
16 #> # i 3 more variables: `<$100-150k` <dbl>, `>150k` <dbl>,
17 #> #   `Don't know/refused` <dbl>
```

このデータセットには 3 つの変数が含まれています：

- 行に格納されている `religion`、
- 列の名前に広がる `income` と、
- セルの値に格納される `count` です。

それを整理するために、`pivot_longer()` を使用します：

```
1 relig_income %>%
2   pivot_longer(
3     cols = !religion,
4     names_to = "income",
5     values_to = "count"
6   )
7 #> # A tibble: 180 x 3
```

```

8 #> religion income count
9 #> <chr> <chr> <dbl>
10 #> 1 Agnostic <$10k 27
11 #> 2 Agnostic $10-20k 34
12 #> 3 Agnostic $20-30k 60
13 #> 4 Agnostic $30-40k 81
14 #> 5 Agnostic $40-50k 76
15 #> 6 Agnostic $50-75k 137
16 #> 7 Agnostic $75-100k 122
17 #> 8 Agnostic $100-150k 109
18 #> 9 Agnostic >150k 84
19 #> 10 Agnostic Don't know/refused 96
20 #> # i 170 more rows

```

- 最初の引数は、整形するデータセット `relig_income` です。
- `cols` は、どのカラムの形を変える必要があるかを記述します。この場合この場合、`religion` を除いたすべてのカラムが対象となります。
- `names_to` は、カラム名に格納されたデータから作成される変数の名前です。`income` のように、カラム名に格納されたデータから作成される変数の名前を指定します。
- セルの値に格納されているデータから作成される変数の名前、つまり `values_to` を指定します。セル値に格納されたデータから生成される変数の名前、すなわち `count` を指定します。

`relig_income` には `names_to` と `values_to` のどちらのカラムも存在しないので、引用符で囲まれた文字列として提供します。

0.2.2 Numeric data in column names

`billboard` データセットには、2000 年における楽曲のビルボードランキングが記録されている。これは `relig_income` データと似た形式を持つが、カラム名にエンコードされているデータは文字列ではなく、実際には数値である。

```

1 billboard
2 #> # A tibble: 317 x 79
3 #> artist track date.entered wk1 wk2 wk3 wk4 wk5 wk6 wk7 wk8
4 #> <chr> <chr> <date> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
5 #> 1 2 Pac Baby~ 2000-02-26 87 82 72 77 87 94 99 NA
6 #> 2 2Ge+her The ~ 2000-09-02 91 87 92 NA NA NA NA NA
7 #> 3 3 Doors D~ Kryp~ 2000-04-08 81 70 68 67 66 57 54 53
8 #> 4 3 Doors D~ Loser 2000-10-21 76 76 72 69 67 65 55 59

```

```

9  #> 5 504 Boyz Wobb~ 2000-04-15 57 34 25 17 17 31 36 49
10 #> 6 98^0 Give~ 2000-08-19 51 39 34 26 26 19 2 2
11 #> 7 A*Teens Danc~ 2000-07-08 97 97 96 95 100 NA NA NA
12 #> 8 Aaliyah I Do~ 2000-01-29 84 62 51 41 38 35 35 38
13 #> 9 Aaliyah Try ~ 2000-03-18 59 53 38 28 21 18 16 14
14 #> 10 Adams, Yo~ Open~ 2000-08-26 76 76 74 69 68 67 61 58
15 #> # i 307 more rows
16 #> # i 68 more variables: wk9 <dbl>, wk10 <dbl>, wk11 <dbl>, wk12 <dbl>,
17 #> # wk13 <dbl>, wk14 <dbl>, wk15 <dbl>, wk16 <dbl>, wk17 <dbl>, wk18 <dbl>,
18 #> # wk19 <dbl>, wk20 <dbl>, wk21 <dbl>, wk22 <dbl>, wk23 <dbl>, wk24 <dbl>,
19 #> # wk25 <dbl>, wk26 <dbl>, wk27 <dbl>, wk28 <dbl>, wk29 <dbl>, wk30 <dbl>,
20 #> # wk31 <dbl>, wk32 <dbl>, wk33 <dbl>, wk34 <dbl>, wk35 <dbl>, wk36 <dbl>,
21 #> # wk37 <dbl>, wk38 <dbl>, wk39 <dbl>, wk40 <dbl>, wk41 <dbl>, wk42 <dbl>, ...

```

まず、`relig_income` データセットと同じ基本仕様で始めることができる。ここでは、名前を `week` という変数にし、値を `rank` という変数にしたい。また、`values_drop_na` を使用して、欠損値に対応する行を削除しています。すべての曲が 76 週すべてチャートにとどまるわけではないので、入力データの構造上、不必要な明示的な NA を作成せざるを得ない。

```

1 billboard %>%
2   pivot_longer(
3     cols = starts_with("wk"),
4     names_to = "week",
5     values_to = "rank",
6     values_drop_na = TRUE
7   )
8 #> # A tibble: 5,307 x 5
9 #>   artist track date.entered week rank
10 #>   <chr> <chr> <date> <chr> <dbl>
11 #> 1 2 Pac Baby Don't Cry (Keep... 2000-02-26 wk1 87
12 #> 2 2 Pac Baby Don't Cry (Keep... 2000-02-26 wk2 82
13 #> 3 2 Pac Baby Don't Cry (Keep... 2000-02-26 wk3 72
14 #> 4 2 Pac Baby Don't Cry (Keep... 2000-02-26 wk4 77
15 #> 5 2 Pac Baby Don't Cry (Keep... 2000-02-26 wk5 87
16 #> 6 2 Pac Baby Don't Cry (Keep... 2000-02-26 wk6 94
17 #> 7 2 Pac Baby Don't Cry (Keep... 2000-02-26 wk7 99
18 #> 8 2Ge+her The Hardest Part Of ... 2000-09-02 wk1 91
19 #> 9 2Ge+her The Hardest Part Of ... 2000-09-02 wk2 87
20 #> 10 2Ge+her The Hardest Part Of ... 2000-09-02 wk3 92

```

```
21 #> # i 5,297 more rows
```

各曲がどれくらいの期間チャートに入っていたかを簡単に調べることができればいいのですが、そのためには `week` 変数を整数に変換することが必要でしょう。`names_prefix` は `wk` という接頭辞を取り除き、`names_transform` は `week` を整数に変換します：

```
1 billboard %>%
2   pivot_longer(
3     cols = starts_with("wk"),
4     names_to = "week",
5     names_prefix = "wk",
6     names_transform = as.integer,
7     values_to = "rank",
8     values_drop_na = TRUE,
9   )
```

また、数値以外の成分を自動的に除去する `readr::parse_number()` を使えば、1 つの引数でこれを行うことができます：

```
1 billboard %>%
2   pivot_longer(
3     cols = starts_with("wk"),
4     names_to = "week",
5     names_transform = readr::parse_number,
6     values_to = "rank",
7     values_drop_na = TRUE,
8   )
```

0.2.3 Many variables in column names

より困難な状況は、カラム名に複数の変数が詰め込まれている場合に発生します。例えば、`who` データセットを見てみましょう：

```
1 who
2 #> # A tibble: 7,240 x 60
3 #>   country iso2 iso3  year new_sp_m014 new_sp_m1524 new_sp_m2534 new_sp_m3544
4 #>   <chr>   <chr> <chr> <dbl>      <dbl>      <dbl>      <dbl>      <dbl>
5 #> 1 Afghani~ AF   AFG   1980         NA         NA         NA         NA
6 #> 2 Afghani~ AF   AFG   1981         NA         NA         NA         NA
7 #> 3 Afghani~ AF   AFG   1982         NA         NA         NA         NA
```

```

8 #> 4 Afghani~ AF AFG 1983 NA NA NA NA
9 #> 5 Afghani~ AF AFG 1984 NA NA NA NA
10 #> 6 Afghani~ AF AFG 1985 NA NA NA NA
11 #> 7 Afghani~ AF AFG 1986 NA NA NA NA
12 #> 8 Afghani~ AF AFG 1987 NA NA NA NA
13 #> 9 Afghani~ AF AFG 1988 NA NA NA NA
14 #> 10 Afghani~ AF AFG 1989 NA NA NA NA
15 #> # i 7,230 more rows
16 #> # i 52 more variables: new_sp_m4554 <dbl>, new_sp_m5564 <dbl>,
17 #> # new_sp_m65 <dbl>, new_sp_f014 <dbl>, new_sp_f1524 <dbl>,
18 #> # new_sp_f2534 <dbl>, new_sp_f3544 <dbl>, new_sp_f4554 <dbl>,
19 #> # new_sp_f5564 <dbl>, new_sp_f65 <dbl>, new_sn_m014 <dbl>,
20 #> # new_sn_m1524 <dbl>, new_sn_m2534 <dbl>, new_sn_m3544 <dbl>,
21 #> # new_sn_m4554 <dbl>, new_sn_m5564 <dbl>, new_sn_m65 <dbl>, ...

```

country、iso2、iso3、year はすでに変数なので、そのままよい。しかし、new_sp_m014 から newrel_f65 までのカラムは、その名前に 4 つの変数をエンコードしています：

- 接頭辞の new_/new は、これらが新規症例のカウントであることを示す。このデータセットには新規のケースしか含まれていないので、ここでは一定なので無視することにする。
- sp/rel/ep は、症例がどのように診断されたかを説明する。
- m/f は性別を表します。
- 014/1524/2535/3544/4554/65 は、年齢層を表します。

names_to で複数のカラム名を指定し、names_sep または names_pattern を指定することで、これらの変数を分割することができます。ここでは names_pattern が最も自然にフィットします。これは extract と似たようなインターフェイスを持っています。グループを含む正規表現 () で定義を与えると、各グループをカラムに配置します。

```

1 who %>%
2   pivot_longer(
3     cols = new_sp_m014:newrel_f65,
4     names_to = c("diagnosis", "gender", "age"),
5     names_pattern = "new_?(.*)_(.)(.*)",
6     values_to = "count"
7   )
8 #> # A tibble: 405,440 x 8
9 #>   country    iso2 iso3   year diagnosis gender age   count
10 #>   <chr>      <chr> <chr> <dbl> <chr>      <chr> <chr> <dbl>
11 #> 1 Afghanistan AF    AFG   1980 sp        m    014    NA

```

```

12 #> 2 Afghanistan AF AFG 1980 sp m 1524 NA
13 #> 3 Afghanistan AF AFG 1980 sp m 2534 NA
14 #> 4 Afghanistan AF AFG 1980 sp m 3544 NA
15 #> 5 Afghanistan AF AFG 1980 sp m 4554 NA
16 #> 6 Afghanistan AF AFG 1980 sp m 5564 NA
17 #> 7 Afghanistan AF AFG 1980 sp m 65 NA
18 #> 8 Afghanistan AF AFG 1980 sp f 014 NA
19 #> 9 Afghanistan AF AFG 1980 sp f 1524 NA
20 #> 10 Afghanistan AF AFG 1980 sp f 2534 NA
21 #> # i 405,430 more rows

```

さらに一歩進んで、`readr` 関数を使用して性別と年齢を因数に変換することができます。これは、既知の値のセットを持つカテゴリ変数がある場合、良い練習になると思います。

```

1 who %>%
2   pivot_longer(
3     cols = new_sp_m014:newrel_f65,
4     names_to = c("diagnosis", "gender", "age"),
5     names_pattern = "new_?(.*)_(.)(.*)",
6     names_transform = list(
7       gender = ~ readr::parse_factor(.x, levels = c("f", "m")),
8       age = ~ readr::parse_factor(
9         .x,
10        levels = c("014", "1524", "2534", "3544", "4554", "5564", "65"),
11        ordered = TRUE
12      )
13    ),
14     values_to = "count",
15   )

```

この方法で実行すると、事後的に変更を行うよりも少し効率的です。`mutate()` では多くの繰り返しを変換する必要があるのに対し、`pivot_longer()` は各名前の 1 つの出現を変換するだけで済みます。

0.2.4 Multiple observations per row

これまで、1 行に 1 つの観測値を持つデータフレームを扱ってきましたが、重要なピボット問題の多くは、1 行に複数の観測値を持ちます。この場合、出力に表示したい列の名前が、入力列名の一部であることから、通常はこのケースを認識できます。このセクションでは、このようなデータをピボットする方法について学びます。

以下の例は、この問題に対する `tidyr` の解決策のインスピレーションとして、[data.table vignette](#) から引用しています。

```

1 household
2 #> # A tibble: 5 x 5
3 #>   family dob_child1 dob_child2 name_child1 name_child2
4 #>   <int> <date>      <date>      <chr>      <chr>
5 #> 1     1 1998-11-26 2000-01-29 Susan      Jose
6 #> 2     2 1996-06-22 NA          Mark      <NA>
7 #> 3     3 2002-07-11 2004-04-05 Sam       Seth
8 #> 4     4 2004-10-10 2009-08-27 Craig     Khai
9 #> 5     5 2000-12-05 2005-02-28 Parker    Gracie

```

各子供について `name` と `dob` (date of birth) という 2 つの情報(または値)を持っていることに注意してください。これらは最終的に別のカラムに格納する必要があります。ここでも `names_to` に複数の変数を指定し、`names_sep` を使って各変数名を分割しています。これは `pivot_longer()` に、列名の一部が測定される「値」(出力では変数となる)を指定するものであることを伝えるものです。

```

1 household %>%
2   pivot_longer(
3     cols = !family,
4     names_to = c(".value", "child"),
5     names_sep = "_",
6     values_drop_na = TRUE
7   )
8 #> # A tibble: 9 x 4
9 #>   family child  dob      name
10 #>   <int> <chr>  <date>    <chr>
11 #> 1     1 child1 1998-11-26 Susan
12 #> 2     1 child2 2000-01-29 Jose
13 #> 3     2 child1 1996-06-22 Mark
14 #> 4     3 child1 2002-07-11 Sam
15 #> 5     3 child2 2004-04-05 Seth
16 #> 6     4 child1 2004-10-10 Craig
17 #> 7     4 child2 2009-08-27 Khai
18 #> 8     5 child1 2000-12-05 Parker
19 #> 9     5 child2 2005-02-28 Gracie

```

`values_drop_na = TRUE` の使用に注意してください：入力形状は、存在しないオブザベーションのための明示的な欠損変数を作成することを強制します。

同様の問題は、ベース R に組み込まれている `anscombe` データセットにも存在します：


```

1  anscombe
2  #>   x1 x2 x3 x4   y1   y2   y3   y4
3  #> 1  10 10 10  8  8.04 9.14  7.46  6.58
4  #> 2   8  8  8  8  6.95 8.14  6.77  5.76
5  #> 3  13 13 13  8  7.58 8.74 12.74  7.71
6  #> 4   9  9  9  8  8.81 8.77  7.11  8.84
7  #> 5  11 11 11  8  8.33 9.26  7.81  8.47
8  #> 6  14 14 14  8  9.96 8.10  8.84  7.04
9  #> 7   6  6  6  8  7.24 6.13  6.08  5.25
10 #> 8   4  4  4 19  4.26 3.10  5.39 12.50
11 #> 9  12 12 12  8 10.84 9.13  8.15  5.56
12 #> 10  7  7  7  8  4.82 7.26  6.42  7.91
13 #> 11  5  5  5  8  5.68 4.74  5.73  6.89

```

このデータセットには4組の変数(x1 と y1, x2 と y2 など)が含まれており, Anscombe のカルテット(平均、標準偏差など、要約統計量は同じだが、データが全く異なる4つのデータセットを集めたもの)の根底にある。私たちは、列 set、x、y を持つデータセットを作りたいと考えています。

```

1  anscombe %>%
2    pivot_longer(
3      cols = everything(),
4      cols_vary = "slowest",
5      names_to = c(".value", "set"),
6      names_pattern = "(.)(. )"
7    )
8  #> # A tibble: 44 x 3
9  #>   set      x      y
10 #>   <chr> <dbl> <dbl>
11 #> 1 1      10  8.04
12 #> 2 1       8  6.95
13 #> 3 1      13  7.58
14 #> 4 1       9  8.81
15 #> 5 1      11  8.33
16 #> 6 1      14  9.96
17 #> 7 1       6  7.24
18 #> 8 1       4  4.26
19 #> 9 1      12 10.8
20 #> 10 1       7  4.82
21 #> # i 34 more rows

```

`cols_vary` を "slowest" に設定すると、`x1` と `y1` の列の値をまとめて出力し、その後 `x2` と `y2` に移行します。この引数は、データセットのすべての列をピボット処理する場合に、より直感的に順序付けられた出力を生成することがよくあります。

パネルデータでも同様の事態が起こりうる。例えば、[Thomas Leeper](#) が提供したこのデータセットの例を見てみましょう。 `anscombe` と同じ手法で整頓することができます：

```
1  pnl <- tibble(  
2    x = 1:4,  
3    a = c(1, 1, 0, 0),  
4    b = c(0, 1, 1, 1),  
5    y1 = rnorm(4),  
6    y2 = rnorm(4),  
7    z1 = rep(3, 4),  
8    z2 = rep(-2, 4),  
9  )  
10  
11  pnl %>%  
12    pivot_longer(  
13      cols = !c(x, a, b),  
14      names_to = c(".value", "time"),  
15      names_pattern = "(.)(. )"   
16    )  
17  #> # A tibble: 8 x 6  
18  #>       x     a     b time      y      z  
19  #>   <int> <dbl> <dbl> <chr>   <dbl> <dbl>  
20  #> 1     1     1     0 1    -0.443     3  
21  #> 2     1     1     0 2     0.862    -2  
22  #> 3     2     1     1 1     0.280     3  
23  #> 4     2     1     1 2    -0.401    -2  
24  #> 5     3     0     1 1     1.08      3  
25  #> 6     3     0     1 2     0.501    -2  
26  #> 7     4     0     1 1    -0.782     3  
27  #> 8     4     0     1 2     0.575    -2
```

0.3 Wider

`pivot_wider()` は、`pivot_longer()` の逆で、列数を増やし、行数を減らすことでデータセットを**広く**するものです。整頓されたデータを作るために `pivot_wider()` が必要になることは比較的少ないですが、プレゼンテーション用の要約表や、他のツールで用いられる形式のデータを作成する際にはよく役に立ちます。

0.3.1 Capture-recapture data

[Myfanwy Johnston](#) が提供した `fish_encounters` データセットは、川を泳ぐ魚が自動監視局によって検出されたときの様子を記述しています：

```
1 fish_encounters
2 #> # A tibble: 114 x 3
3 #>   fish station seen
4 #>   <fct> <fct>   <int>
5 #> 1 4842 Release     1
6 #> 2 4842 I80_1      1
7 #> 3 4842 Lisbon     1
8 #> 4 4842 Rstr       1
9 #> 5 4842 Base_TD    1
10 #> 6 4842 BCE        1
11 #> 7 4842 BCW        1
12 #> 8 4842 BCE2       1
13 #> 9 4842 BCW2       1
14 #> 10 4842 MAE       1
15 #> # i 104 more rows
```

このデータを分析するために用いられる多くのツールは、各ステーションが列となる形でデータを必要とします：

```
1 fish_encounters %>%
2   pivot_wider(
3     names_from = station,
4     values_from = seen
5   )
6 #> # A tibble: 19 x 12
7 #>   fish Release I80_1 Lisbon Rstr Base_TD BCE BCW BCE2 BCW2 MAE MAW
8 #>   <fct>   <int> <int>   <int> <int>   <int> <int> <int> <int> <int> <int> <int>
9 #> 1 4842      1     1     1     1     1     1     1     1     1     1     1
10 #> 2 4843      1     1     1     1     1     1     1     1     1     1     1
11 #> 3 4844      1     1     1     1     1     1     1     1     1     1     1
12 #> 4 4845      1     1     1     1     1     NA     NA     NA     NA     NA     NA
13 #> 5 4847      1     1     1     NA     NA     NA     NA     NA     NA     NA     NA
14 #> 6 4848      1     1     1     1     NA     NA     NA     NA     NA     NA     NA
15 #> 7 4849      1     1     NA     NA     NA     NA     NA     NA     NA     NA     NA
16 #> 8 4850      1     1     NA     1     1     1     1     NA     NA     NA     NA
```

```

17 #> 9 4851      1      1      NA      NA      NA      NA      NA      NA      NA      NA      NA
18 #> 10 4854     1      1      NA      NA      NA      NA      NA      NA      NA      NA      NA
19 #> # i 9 more rows

```

このデータセットでは、魚がステーションで検出された場合のみ記録され、検出されなかった場合は記録されません(これはこの種のデータではよくあることです)。つまり、出力データは `NA` で埋め尽くされます。しかし、この場合、記録がないということは魚が「目撃(`seen`)されなかった」ということなので、`pivot_wider()` に欠損値をゼロで埋めてもらうことができます：

```

1 fish_encounters %>%
2   pivot_wider(
3     names_from = station,
4     values_from = seen,
5     values_fill = 0
6   )
7 #> # A tibble: 19 x 12
8 #>   fish Release I80_1 Lisbon Rstr Base_TD BCE BCW BCE2 BCW2 MAE MAW
9 #>   <fct>   <int> <int>   <int> <int>   <int> <int> <int> <int> <int> <int> <int>
10 #> 1 4842      1      1      1      1      1      1      1      1      1      1      1      1
11 #> 2 4843      1      1      1      1      1      1      1      1      1      1      1      1
12 #> 3 4844      1      1      1      1      1      1      1      1      1      1      1      1
13 #> 4 4845      1      1      1      1      1      0      0      0      0      0      0      0
14 #> 5 4847      1      1      1      0      0      0      0      0      0      0      0      0
15 #> 6 4848      1      1      1      1      0      0      0      0      0      0      0      0
16 #> 7 4849      1      1      0      0      0      0      0      0      0      0      0      0
17 #> 8 4850      1      1      0      1      1      1      1      0      0      0      0      0
18 #> 9 4851      1      1      0      0      0      0      0      0      0      0      0      0
19 #> 10 4854     1      1      0      0      0      0      0      0      0      0      0      0
20 #> # i 9 more rows

```

0.3.2 Aggregation

また、`pivot_wider()` を使って単純な集計を行うこともできます。例えば、ベース R に組み込まれている `warpbreaks` データセットを見てみましょう(より良いプリント手法のために `tibble` に変換されています)：

```

1 warpbreaks <- warpbreaks %>%
2   as_tibble() %>%
3   select(wool, tension, breaks)
4 warpbreaks
5 #> # A tibble: 54 x 3

```

```

6 #>      wool  tension breaks
7 #>      <fct> <fct>      <dbl>
8 #>  1 A      L          26
9 #>  2 A      L          30
10 #>  3 A      L          54
11 #>  4 A      L          25
12 #>  5 A      L          70
13 #>  6 A      L          52
14 #>  7 A      L          51
15 #>  8 A      L          26
16 #>  9 A      L          67
17 #> 10 A      M          18
18 #> # i 44 more rows

```

これは、`wool` (A、B)と `tension` (L、M、H)の組み合わせごとに9つの複製を持つ設計実験である：

```

1 warpbreaks %>%
2   count(wool, tension)
3 #> # A tibble: 6 x 3
4 #>   wool  tension     n
5 #>   <fct> <fct>   <int>
6 #> 1 A      L         9
7 #> 2 A      M         9
8 #> 3 A      H         9
9 #> 4 B      L         9
10 #> 5 B      M         9
11 #> 6 B      H         9

```

`wool` の水準を列にピボットさせようとするとうなるでしょうか？

```

1 warpbreaks %>%
2   pivot_wider(
3     names_from = wool,
4     values_from = breaks
5   )
6 #> # A tibble: 3 x 3
7 #>   tension A      B
8 #>   <fct>   <list> <list>
9 #> 1 L      <dbl [9]> <dbl [9]>
10 #> 2 M      <dbl [9]> <dbl [9]>

```

```
11 #> 3 H          <dbl [9]> <dbl [9]>
```

出力の各セルが、入力複数のセルに対応しているという警告が表示されます。デフォルトの動作では、個々の値をすべて含むリストカラムが生成されます。より有用な出力は、例えば、ウールとテンションの組み合わせごとの平均値などの要約統計です：

```
1 warpbreaks %>%
2   pivot_wider(
3     names_from = wool,
4     values_from = breaks,
5     values_fn = mean
6   )
7 #> # A tibble: 3 x 3
8 #>   tension      A      B
9 #>   <fct>    <dbl> <dbl>
10 #> 1 L      44.6  28.2
11 #> 2 M      24    28.8
12 #> 3 H      24.6  18.8
```

より複雑な要約操作の場合は、再形成の前に要約することをお勧めしますが、単純なケースでは、`pivot_wider()` 内で要約するのが便利が多いです。

0.3.3 Generate column name from multiple variables

<https://stackoverflow.com/questions/24929954> のように、製品、国、年の組み合わせを含む情報があると想像してください。整頓された形では、次のようになります：

```
1 production <-
2   expand_grid(
3     product = c("A", "B"),
4     country = c("AI", "EI"),
5     year = 2000:2014
6   ) %>%
7   filter((product == "A" & country == "AI") | product == "B") %>%
8   mutate(production = rnorm(nrow(.)))
9 production
10 #> # A tibble: 45 x 4
11 #>   product country  year production
12 #>   <chr>    <chr>  <int>      <dbl>
13 #> 1 A      AI      2000    -0.922
14 #> 2 A      AI      2001     0.407
```

```

15 #> 3 A      AI      2002      2.23
16 #> 4 A      AI      2003      1.51
17 #> 5 A      AI      2004      2.77
18 #> 6 A      AI      2005      1.84
19 #> 7 A      AI      2006     -0.0340
20 #> 8 A      AI      2007     -0.203
21 #> 9 A      AI      2008     -1.46
22 #> 10 A     AI      2009     -0.146
23 #> # i 35 more rows

```

データを広げて、`product` と `country` の組み合わせごとに1つのカラムを持つようにしたいです。このとき、`names_from` に複数の変数を指定することがポイントです：

```

1 production %>%
2   pivot_wider(
3     names_from = c(product, country),
4     values_from = production
5   )
6 #> # A tibble: 15 x 4
7 #>   year    A_AI    B_AI    B_EI
8 #>   <int>  <dbl>  <dbl>  <dbl>
9 #> 1  2000 -0.922  0.425  0.324
10 #> 2  2001  0.407  1.47   0.854
11 #> 3  2002  2.23   1.01   1.95
12 #> 4  2003  1.51  -0.658 -2.16
13 #> 5  2004  2.77   0.613 -1.38
14 #> 6  2005  1.84   0.403 -0.169
15 #> 7  2006 -0.0340 -1.20   0.377
16 #> 8  2007 -0.203  0.0362 -0.0955
17 #> 9  2008 -1.46  -0.263 -0.451
18 #> 10 2009 -0.146  2.22   0.214
19 #> # i 5 more rows

```

`names_from` または `values_from` が複数の変数を選択する場合、`names_sep` や `names_prefix`、または `names_glue` を使って、出力する列名の構成方法を制御することができます：

```

1 production %>%
2   pivot_wider(
3     names_from = c(product, country),
4     values_from = production,

```

```

5     names_sep = ".",
6     names_prefix = "prod."
7   )
8 #> # A tibble: 15 x 4
9 #>   year prod.A.AI prod.B.AI prod.B.EI
10 #>   <int>    <dbl>    <dbl>    <dbl>
11 #> 1  2000   -0.922    0.425    0.324
12 #> 2  2001    0.407    1.47     0.854
13 #> 3  2002    2.23     1.01     1.95
14 #> 4  2003    1.51   -0.658   -2.16
15 #> 5  2004    2.77    0.613   -1.38
16 #> 6  2005    1.84    0.403   -0.169
17 #> 7  2006   -0.0340  -1.20    0.377
18 #> 8  2007   -0.203    0.0362  -0.0955
19 #> 9  2008   -1.46   -0.263   -0.451
20 #> 10 2009   -0.146    2.22    0.214
21 #> # i 5 more rows
22
23 production %>%
24   pivot_wider(
25     names_from = c(product, country),
26     values_from = production,
27     names_glue = "prod_{product}_{country}"
28   )
29 #> # A tibble: 15 x 4
30 #>   year prod_A_AI prod_B_AI prod_B_EI
31 #>   <int>    <dbl>    <dbl>    <dbl>
32 #> 1  2000   -0.922    0.425    0.324
33 #> 2  2001    0.407    1.47     0.854
34 #> 3  2002    2.23     1.01     1.95
35 #> 4  2003    1.51   -0.658   -2.16
36 #> 5  2004    2.77    0.613   -1.38
37 #> 6  2005    1.84    0.403   -0.169
38 #> 7  2006   -0.0340  -1.20    0.377
39 #> 8  2007   -0.203    0.0362  -0.0955
40 #> 9  2008   -1.46   -0.263   -0.451
41 #> 10 2009   -0.146    2.22    0.214
42 #> # i 5 more rows

```


0.3.4 Tidy census

`us_rent_income` データセットには、2017 年のアメリカの各州の所得と家賃の中央値に関する情報が含まれています(American Community Survey より、[tidycensus](#) パッケージで取得)。

```
1 us_rent_income
2 #> # A tibble: 104 x 5
3 #>   GEOID NAME      variable estimate   moe
4 #>   <chr> <chr>      <chr>      <dbl> <dbl>
5 #> 1 01    Alabama    income    24476  136
6 #> 2 01    Alabama    rent       747    3
7 #> 3 02    Alaska     income    32940  508
8 #> 4 02    Alaska     rent     1200   13
9 #> 5 04    Arizona    income    27517  148
10 #> 6 04    Arizona    rent       972    4
11 #> 7 05    Arkansas   income    23789  165
12 #> 8 05    Arkansas   rent       709    5
13 #> 9 06    California income    29454  109
14 #> 10 06   California rent      1358    3
15 #> # i 94 more rows
```

ここで `estimate` と `moe` はどちらも値のカラムなので、`values_from` に渡すことができます：

```
1 us_rent_income %>%
2   pivot_wider(
3     names_from = variable,
4     values_from = c(estimate, moe)
5   )
6 #> # A tibble: 52 x 6
7 #>   GEOID NAME      estimate_income estimate_rent moe_income moe_rent
8 #>   <chr> <chr>      <dbl>      <dbl>      <dbl>   <dbl>
9 #> 1 01    Alabama    24476      747      136     3
10 #> 2 02    Alaska    32940     1200     508    13
11 #> 3 04    Arizona    27517      972     148    4
12 #> 4 05    Arkansas   23789      709     165    5
13 #> 5 06    California  29454     1358     109    3
14 #> 6 08    Colorado    32401     1125     109    5
15 #> 7 09    Connecticut  35326     1123     195    5
16 #> 8 10    Delaware    31560     1076     247   10
17 #> 9 11    District of Columbia  43198     1424     681   17
```

```

18 #> 10 12 Florida 25952 1077 70 3
19 #> # i 42 more rows

```

なお、出力列には、変数名が自動的に付加されます。

0.3.5 Implicit missing values

時折、名前変数が因子として符号化されているデータに出会うことがあります、すべてのデータが表現されるわけではありません。

```

1 weekdays <- c("Mon", "Tue", "Wed", "Thu", "Fri", "Sat", "Sun")
2
3 daily <- tibble(
4   day = factor(c("Tue", "Thu", "Fri", "Mon"), levels = weekdays),
5   value = c(2, 3, 1, 5)
6 )
7
8 daily
9 #> # A tibble: 4 x 2
10 #>   day    value
11 #>   <fct> <dbl>
12 #> 1 Tue      2
13 #> 2 Thu      3
14 #> 3 Fri      1
15 #> 4 Mon      5

```

`pivot_wider()` のデフォルトでは、実際にデータで表現されている値から列を生成しますが、将来的にデータが変更された場合に備えて、可能性のある各レベルの列を含めるとよいかもしれません。

```

1 daily %>%
2   pivot_wider(
3     names_from = day,
4     values_from = value
5   )
6 #> # A tibble: 1 x 4
7 #>   Tue  Thu  Fri  Mon
8 #>   <dbl> <dbl> <dbl> <dbl>
9 #> 1     2     3     1     5

```

`names_expand` 引数は、暗黙の因子レベルを明示的なものに変え、結果で表現することを強制します。また、レベルの順序を使用してカラム名をソートするので、この場合、より直感的な結果が得られます。

```

1 daily %>%
2   pivot_wider(
3     names_from = day,
4     values_from = value,
5     names_expand = TRUE
6   )
7 #> # A tibble: 1 x 7
8 #>   Mon  Tue  Wed  Thu  Fri  Sat  Sun
9 #>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
10 #> 1     5     2   NA     3     1   NA   NA

```

If multiple `names_from` columns are provided, `names_expand` will generate a Cartesian product of all possible combinations of the `names_from` values. Notice that the following data has omitted some rows where the percentage value would be 0. `names_expand` allows us to make those explicit during the pivot.

複数の `names_from` カラムを指定すると、`names_expand` は `names_from` 値のすべての可能な組み合わせのデカルト積を生成します。以下のデータでは、パーセント値が 0 となる行がいくつか省略されていることに注意してください。 `names_expand` を使用すると、ピボットの際にこれらを明示することができます。

```

1 percentages <- tibble(
2   year = c(2018, 2019, 2020, 2020),
3   type = factor(c("A", "B", "A", "B"), levels = c("A", "B")),
4   percentage = c(100, 100, 40, 60)
5 )
6
7 percentages
8 #> # A tibble: 4 x 3
9 #>   year type  percentage
10 #>   <dbl> <fct>      <dbl>
11 #> 1  2018 A          100
12 #> 2  2019 B          100
13 #> 3  2020 A           40
14 #> 4  2020 B           60
15
16 percentages %>%
17   pivot_wider(
18     names_from = c(year, type),
19     values_from = percentage,
20     names_expand = TRUE,
21     values_fill = 0

```

```

22   )
23   #> # A tibble: 1 x 6
24   #>   `2018_A` `2018_B` `2019_A` `2019_B` `2020_A` `2020_B`
25   #>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
26   #> 1     100       0       0     100     40     60

```

A related problem can occur when there are implicit missing factor levels or combinations in the `id_cols`. In this case, there are missing rows (rather than columns) that you'd like to explicitly represent. For this example, we'll modify our daily data with a `type` column, and pivot on that instead, keeping `day` as an `id` column.

```

1  daily <- mutate(daily, type = factor(c("A", "B", "B", "A")))
2  daily
3  #> # A tibble: 4 x 3
4  #>   day   value type
5  #>   <fct> <dbl> <fct>
6  #> 1 Tue       2 A
7  #> 2 Thu       3 B
8  #> 3 Fri       1 B
9  #> 4 Mon       5 A

```

All of our `type` levels are represented in the columns, but we are missing some rows related to the unrepresented `day` factor levels.

```

1  daily %>%
2    pivot_wider(
3      names_from = type,
4      values_from = value,
5      values_fill = 0
6    )
7  #> # A tibble: 4 x 3
8  #>   day     A     B
9  #>   <fct> <dbl> <dbl>
10 #> 1 Tue       2     0
11 #> 2 Thu       0     3
12 #> 3 Fri       0     1
13 #> 4 Mon       5     0

```

We can use `id_expand` in the same way that we used `names_expand`, which will expand out (and sort) the implicit missing rows in the `id_cols`.

```

1 daily %>%
2   pivot_wider(
3     names_from = type,
4     values_from = value,
5     values_fill = 0,
6     id_expand = TRUE
7   )
8 #> # A tibble: 7 x 3
9 #>   day      A      B
10 #>   <fct> <dbl> <dbl>
11 #> 1 Mon      5      0
12 #> 2 Tue      2      0
13 #> 3 Wed      0      0
14 #> 4 Thu      0      3
15 #> 5 Fri      0      1
16 #> 6 Sat      0      0
17 #> 7 Sun      0      0

```

0.3.6 Unused columns

Imagine you've found yourself in a situation where you have columns in your data that are completely unrelated to the pivoting process, but you'd still like to retain their information somehow. For example, in updates we'd like to pivot on the system column to create one row summaries of each county's system updates.

```

1 updates <- tibble(
2   county = c("Wake", "Wake", "Wake", "Guilford", "Guilford"),
3   date = c(as.Date("2020-01-01") + 0:2, as.Date("2020-01-03") + 0:1),
4   system = c("A", "B", "C", "A", "C"),
5   value = c(3.2, 4, 5.5, 2, 1.2)
6 )
7
8 updates
9 #> # A tibble: 5 x 4
10 #>   county  date      system value
11 #>   <chr>   <date>   <chr>   <dbl>
12 #> 1 Wake   2020-01-01 A        3.2
13 #> 2 Wake   2020-01-02 B         4
14 #> 3 Wake   2020-01-03 C        5.5

```

```

15 #> 4 Guilford 2020-01-03 A      2
16 #> 5 Guilford 2020-01-04 C      1.2

```

We could do that with a typical `pivot_wider()` call, but we completely lose all information about the date column.

```

1 updates %>%
2   pivot_wider(
3     id_cols = county,
4     names_from = system,
5     values_from = value
6   )
7 #> # A tibble: 2 x 4
8 #>   county      A      B      C
9 #>   <chr>    <dbl> <dbl> <dbl>
10 #> 1 Wake      3.2      4    5.5
11 #> 2 Guilford   2      NA    1.2

```

For this example, we'd like to retain the most recent update date across all systems in a particular county. To accomplish that we can use the `unused_fn` argument, which allows us to summarize values from the columns not utilized in the pivoting process.

```

1 updates %>%
2   pivot_wider(
3     id_cols = county,
4     names_from = system,
5     values_from = value,
6     unused_fn = list(date = max)
7   )
8 #> # A tibble: 2 x 5
9 #>   county      A      B      C date
10 #>   <chr>    <dbl> <dbl> <dbl> <date>
11 #> 1 Wake      3.2      4    5.5 2020-01-03
12 #> 2 Guilford   2      NA    1.2 2020-01-04

```

You can also retain the data but delay the aggregation entirely by using `list()` as the summary function.

```

1 updates %>%
2   pivot_wider(
3     id_cols = county,

```

```

4     names_from = system,
5     values_from = value,
6     unused_fn = list(date = list)
7   )
8 #> # A tibble: 2 x 5
9 #>   county      A      B      C date
10 #>   <chr>    <dbl> <dbl> <dbl> <list>
11 #> 1 Wake      3.2      4    5.5 <date [3]>
12 #> 2 Guilford  2        NA    1.2 <date [2]>

```

0.3.7 Contact list

A final challenge is inspired by [Jiena Gu](#). Imagine you have a contact list that you’ve copied and pasted from a website:

```

1 contacts <- tribble(
2   ~field, ~value,
3   "name", "Jiena McLellan",
4   "company", "Toyota",
5   "name", "John Smith",
6   "company", "google",
7   "email", "john@google.com",
8   "name", "Huxley Ratcliffe"
9 )

```

This is challenging because there’s no variable that identifies which observations belong together. We can fix this by noting that every contact starts with a name, so we can create a unique id by counting every time we see “name” as the field:

```

1 contacts <- contacts %>%
2   mutate(
3     person_id = cumsum(field == "name")
4   )
5 contacts
6 #> # A tibble: 6 x 3
7 #>   field  value      person_id
8 #>   <chr>  <chr>         <int>
9 #> 1 name   Jiena McLellan      1
10 #> 2 company Toyota          1
11 #> 3 name   John Smith         2

```

```

12 #> 4 company google                2
13 #> 5 email   john@google.com        2
14 #> 6 name    Huxley Ratcliffe        3

```

Now that we have a unique identifier for each person, we can pivot field and value into the columns:

```

1 contacts %>%
2   pivot_wider(
3     names_from = field,
4     values_from = value
5   )
6 #> # A tibble: 3 x 4
7 #>   person_id name          company email
8 #>   <int> <chr>          <chr> <chr>
9 #> 1     1 Jiena McLellan Toyota <NA>
10 #> 2     2 John Smith      google john@google.com
11 #> 3     3 Huxley Ratcliffe <NA> <NA>

```

0.4 Longer, then wider

Some problems can't be solved by pivoting in a single direction. The examples in this section show how you might combine `pivot_longer()` and `pivot_wider()` to solve more complex problems.

0.4.1 World bank

`world_bank_pop` contains data from the World Bank about population per country from 2000 to 2018.

```

1 world_bank_pop
2 #> # A tibble: 1,064 x 20
3 #>   country indicator `2000` `2001` `2002` `2003` `2004` `2005` `2006`
4 #>   <chr> <chr>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
5 #> 1 ABW SP.URB.TOTL 4.16e4 4.20e+4 4.22e+4 4.23e+4 4.23e+4 4.24e+4 4.26e+4
6 #> 2 ABW SP.URB.GROW 1.66e0 9.56e-1 4.01e-1 1.97e-1 9.46e-2 1.94e-1 3.67e-1
7 #> 3 ABW SP.POP.TOTL 8.91e4 9.07e+4 9.18e+4 9.27e+4 9.35e+4 9.45e+4 9.56e+4
8 #> 4 ABW SP.POP.GROW 2.54e0 1.77e+0 1.19e+0 9.97e-1 9.01e-1 1.00e+0 1.18e+0
9 #> 5 AFE SP.URB.TOTL 1.16e8 1.20e+8 1.24e+8 1.29e+8 1.34e+8 1.39e+8 1.44e+8
10 #> 6 AFE SP.URB.GROW 3.60e0 3.66e+0 3.72e+0 3.71e+0 3.74e+0 3.81e+0 3.81e+0
11 #> 7 AFE SP.POP.TOTL 4.02e8 4.12e+8 4.23e+8 4.34e+8 4.45e+8 4.57e+8 4.70e+8
12 #> 8 AFE SP.POP.GROW 2.58e0 2.59e+0 2.61e+0 2.62e+0 2.64e+0 2.67e+0 2.70e+0
13 #> 9 AFG SP.URB.TOTL 4.31e6 4.36e+6 4.67e+6 5.06e+6 5.30e+6 5.54e+6 5.83e+6

```



```

14 #> 10 AFG      SP.URB.GROW      1.86e0 1.15e+0 6.86e+0 7.95e+0 4.59e+0 4.47e+0 5.03e+0
15 #> # i 1,054 more rows
16 #> # i 11 more variables: `2007` <dbl>, `2008` <dbl>, `2009` <dbl>, `2010` <dbl>,
17 #> #   `2011` <dbl>, `2012` <dbl>, `2013` <dbl>, `2014` <dbl>, `2015` <dbl>,
18 #> #   `2016` <dbl>, `2017` <dbl>

```

My goal is to produce a tidy dataset where each variable is in a column. It's not obvious exactly what steps are needed yet, but I'll start with the most obvious problem: year is spread across multiple columns.

```

1 pop2 <- world_bank_pop %>%
2   pivot_longer(
3     cols = `2000`:`2017`,
4     names_to = "year",
5     values_to = "value"
6   )
7 pop2
8 #> # A tibble: 19,152 x 4
9 #>   country indicator   year  value
10 #>   <chr>    <chr>      <chr> <dbl>
11 #> 1 ABW      SP.URB.TOTL 2000  41625
12 #> 2 ABW      SP.URB.TOTL 2001  42025
13 #> 3 ABW      SP.URB.TOTL 2002  42194
14 #> 4 ABW      SP.URB.TOTL 2003  42277
15 #> 5 ABW      SP.URB.TOTL 2004  42317
16 #> 6 ABW      SP.URB.TOTL 2005  42399
17 #> 7 ABW      SP.URB.TOTL 2006  42555
18 #> 8 ABW      SP.URB.TOTL 2007  42729
19 #> 9 ABW      SP.URB.TOTL 2008  42906
20 #> 10 ABW     SP.URB.TOTL 2009  43079
21 #> # i 19,142 more rows

```

Next we need to consider the indicator variable:

```

1 pop2 %>%
2   count(indicator)
3 #> # A tibble: 4 x 2
4 #>   indicator      n
5 #>   <chr>      <int>
6 #> 1 SP.POP.GROW  4788
7 #> 2 SP.POP.TOTL  4788

```

```

8 #> 3 SP.URB.GROW 4788
9 #> 4 SP.URB.TOTL 4788

```

Here SP.POP.GROW is population growth, SP.POP.TOTL is total population, and SP.URB.* are the same but only for urban areas. Let's split this up into two variables: area (total or urban) and the actual variable (population or growth):

```

1 pop3 <- pop2 %>%
2   separate(indicator, c(NA, "area", "variable"))
3 pop3
4 #> # A tibble: 19,152 x 5
5 #>   country area variable year value
6 #>   <chr>   <chr> <chr>   <chr> <dbl>
7 #> 1 ABW    URB    TOTL    2000 41625
8 #> 2 ABW    URB    TOTL    2001 42025
9 #> 3 ABW    URB    TOTL    2002 42194
10 #> 4 ABW    URB    TOTL    2003 42277
11 #> 5 ABW    URB    TOTL    2004 42317
12 #> 6 ABW    URB    TOTL    2005 42399
13 #> 7 ABW    URB    TOTL    2006 42555
14 #> 8 ABW    URB    TOTL    2007 42729
15 #> 9 ABW    URB    TOTL    2008 42906
16 #> 10 ABW    URB    TOTL    2009 43079
17 #> # i 19,142 more rows

```

Now we can complete the tidying by pivoting variable and value to make TOTL and GROW columns:

```

1 pop3 %>%
2   pivot_wider(
3     names_from = variable,
4     values_from = value
5   )
6 #> # A tibble: 9,576 x 5
7 #>   country area year  TOTL  GROW
8 #>   <chr>   <chr> <chr> <dbl> <dbl>
9 #> 1 ABW    URB    2000 41625 1.66
10 #> 2 ABW    URB    2001 42025 0.956
11 #> 3 ABW    URB    2002 42194 0.401
12 #> 4 ABW    URB    2003 42277 0.197
13 #> 5 ABW    URB    2004 42317 0.0946

```

```

14 #> 6 ABW      URB      2005  42399 0.194
15 #> 7 ABW      URB      2006  42555 0.367
16 #> 8 ABW      URB      2007  42729 0.408
17 #> 9 ABW      URB      2008  42906 0.413
18 #> 10 ABW     URB      2009  43079 0.402
19 #> # i 9,566 more rows

```

0.4.2 Multi-choice

Based on a suggestion by [Maxime Wack](https://github.com/tidyverse/tidyr/issues/384), <https://github.com/tidyverse/tidyr/issues/384>, the final example shows how to deal with a common way of recording multiple choice data. Often you will get such data as follows:

```

1 multi <- tribble(
2   ~id, ~choice1, ~choice2, ~choice3,
3   1, "A", "B", "C",
4   2, "C", "B", NA,
5   3, "D", NA, NA,
6   4, "B", "D", NA
7 )

```

But the actual order isn't important, and you'd prefer to have the individual questions in the columns. You can achieve the desired transformation in two steps. First, you make the data longer, eliminating the explicit NAs, and adding a column to indicate that this choice was chosen:

```

1 multi2 <- multi %>%
2   pivot_longer(
3     cols = !id,
4     values_drop_na = TRUE
5   ) %>%
6   mutate(checked = TRUE)
7 multi2
8 #> # A tibble: 8 x 4
9 #>       id name      value checked
10 #>   <dbl> <chr>    <chr> <lgl>
11 #> 1     1 choice1 A      TRUE
12 #> 2     1 choice2 B      TRUE
13 #> 3     1 choice3 C      TRUE
14 #> 4     2 choice1 C      TRUE
15 #> 5     2 choice2 B      TRUE

```

```

16 #> 6      3 choice1 D      TRUE
17 #> 7      4 choice1 B      TRUE
18 #> 8      4 choice2 D      TRUE

```

Then you make the data wider, filling in the missing observations with FALSE:

```

1 multi2 %>%
2   pivot_wider(
3     id_cols = id,
4     names_from = value,
5     values_from = checked,
6     values_fill = FALSE
7   )
8 #> # A tibble: 4 x 5
9 #>       id A      B      C      D
10 #>   <dbl> <lgl> <lgl> <lgl> <lgl>
11 #> 1     1 TRUE  TRUE  TRUE  FALSE
12 #> 2     2 FALSE TRUE  TRUE  FALSE
13 #> 3     3 FALSE FALSE FALSE TRUE
14 #> 4     4 FALSE TRUE  FALSE TRUE

```

0.5 Manual specs

The arguments to `pivot_longer()` and `pivot_wider()` allow you to pivot a wide range of datasets. But the creativity that people apply to their data structures is seemingly endless, so it's quite possible that you will encounter a dataset that you can't immediately see how to reshape with `pivot_longer()` and `pivot_wider()`. To gain more control over pivoting, you can instead create a "spec" data frame that describes exactly how data stored in the column names becomes variables (and vice versa). This section introduces you to the spec data structure, and show you how to use it when `pivot_longer()` and `pivot_wider()` are insufficient.

0.5.1 Longer

To see how this works, let's return to the simplest case of pivoting applied to the `relig_income` dataset. Now pivoting happens in two steps: we first create a spec object (using `build_longer_spec()`) then use that to describe the pivoting operation:

```

1 spec <- relig_income %>%
2   build_longer_spec(
3     cols = !religion,

```

```

4     names_to = "income",
5     values_to = "count"
6   )
7   pivot_longer_spec(relig_income, spec)
8   #> # A tibble: 180 x 3
9   #>   religion income          count
10  #>   <chr>   <chr>         <dbl>
11  #> 1 Agnostic <$10k          27
12  #> 2 Agnostic $10-20k         34
13  #> 3 Agnostic $20-30k         60
14  #> 4 Agnostic $30-40k         81
15  #> 5 Agnostic $40-50k         76
16  #> 6 Agnostic $50-75k        137
17  #> 7 Agnostic $75-100k       122
18  #> 8 Agnostic $100-150k      109
19  #> 9 Agnostic >150k          84
20  #> 10 Agnostic Don't know/refused 96
21  #> # i 170 more rows

```

(This gives the same result as before, just with more code. There's no need to use it here, it is presented as a simple example for using `spec`.)

What does `spec` look like? It's a data frame with one row for each column in the wide format version of the data that is not present in the long format, and two special columns that start with `.`:

- `.name` gives the name of the column.
- `.value` gives the name of the column that the values in the cells will go into.

There is also one column in `spec` for each column present in the long format of the data that is not present in the wide format of the data. This corresponds to the `names_to` argument in `pivot_longer()` and `build_longer_spec()` and the `names_from` argument in `pivot_wider()` and `build_wider_spec()`. In this example, the `income` column is a character vector of the names of columns being pivoted.

```

1   spec
2   #> # A tibble: 10 x 3
3   #>   .name          .value income
4   #>   <chr>         <chr>  <chr>
5   #> 1 <$10k         count  <$10k
6   #> 2 $10-20k       count  $10-20k
7   #> 3 $20-30k       count  $20-30k
8   #> 4 $30-40k       count  $30-40k

```

```

9 #> 5 $40-50k          count $40-50k
10 #> 6 $50-75k          count $50-75k
11 #> 7 $75-100k         count $75-100k
12 #> 8 $100-150k        count $100-150k
13 #> 9 >150k            count >150k
14 #> 10 Don't know/refused count Don't know/refused

```

0.5.2 Wider

Below we widen `us_rent_income` with `pivot_wider()`. The result is ok, but I think it could be improved:

```

1 us_rent_income %>%
2   pivot_wider(
3     names_from = variable,
4     values_from = c(estimate, moe)
5   )
6 #> # A tibble: 52 x 6
7 #>   GEOID NAME          estimate_income estimate_rent moe_income moe_rent
8 #>   <chr> <chr>          <dbl>          <dbl>      <dbl>  <dbl>
9 #> 1 01 Alabama          24476           747        136     3
10 #> 2 02 Alaska           32940          1200        508    13
11 #> 3 04 Arizona           27517           972        148     4
12 #> 4 05 Arkansas          23789           709        165     5
13 #> 5 06 California         29454          1358        109     3
14 #> 6 08 Colorado           32401          1125        109     5
15 #> 7 09 Connecticut        35326          1123        195     5
16 #> 8 10 Delaware           31560          1076        247    10
17 #> 9 11 District of Columbia 43198          1424        681    17
18 #> 10 12 Florida           25952          1077         70     3
19 #> # i 42 more rows

```

I think it would be better to have columns `income`, `rent`, `income_moe`, and `rent_moe`, which we can achieve with a manual spec. The current spec looks like this:

```

1 spec1 <- us_rent_income %>%
2   build_wider_spec(
3     names_from = variable,
4     values_from = c(estimate, moe)
5   )
6 spec1

```

```

7 #> # A tibble: 4 x 3
8 #>   .name      .value variable
9 #>   <chr>      <chr>   <chr>
10 #> 1 estimate_income estimate income
11 #> 2 estimate_rent   estimate rent
12 #> 3 moe_income      moe      income
13 #> 4 moe_rent        moe      rent

```

For this case, we mutate spec to carefully construct the column names:

```

1 spec2 <- spec1 %>%
2   mutate(
3     .name = paste0(variable, ifelse(.value == "moe", "_moe", ""))
4   )
5 spec2
6 #> # A tibble: 4 x 3
7 #>   .name      .value variable
8 #>   <chr>      <chr>   <chr>
9 #> 1 income      estimate income
10 #> 2 rent         estimate rent
11 #> 3 income_moe moe      income
12 #> 4 rent_moe    moe      rent

```

Supplying this spec to `pivot_wider()` gives us the result we're looking for:

```

1 us_rent_income %>%
2   pivot_wider_spec(spec2)
3 #> # A tibble: 52 x 6
4 #>   GEOID NAME      income rent income_moe rent_moe
5 #>   <chr> <chr>      <dbl> <dbl>      <dbl>      <dbl>
6 #> 1 01    Alabama    24476   747      136         3
7 #> 2 02    Alaska     32940  1200      508        13
8 #> 3 04    Arizona     27517   972      148         4
9 #> 4 05    Arkansas     23789   709      165         5
10 #> 5 06    California    29454  1358      109         3
11 #> 6 08    Colorado     32401  1125      109         5
12 #> 7 09    Connecticut    35326  1123      195         5
13 #> 8 10    Delaware     31560  1076      247        10
14 #> 9 11    District of Columbia 43198  1424      681        17
15 #> 10 12    Florida     25952  1077       70         3

```

```
16 #> # i 42 more rows
```

0.5.3 By hand

Sometimes it's not possible (or not convenient) to compute the spec, and instead it's more convenient to construct the spec "by hand". For example, take this construction data, which is lightly modified from Table 5 "completions" found at <https://www.census.gov/construction/nrc/index.html>:

```
1 construction
2 #> # A tibble: 9 x 9
3 #>   Year Month `1 unit` `2 to 4 units` `5 units or more` Northeast Midwest South
4 #>   <dbl> <chr>   <dbl> <lgl>           <dbl>    <dbl>    <dbl> <dbl>
5 #> 1  2018 Janua~      859 NA              348      114      169    596
6 #> 2  2018 Febru~      882 NA              400      138      160    655
7 #> 3  2018 March      862 NA              356      150      154    595
8 #> 4  2018 April      797 NA              447      144      196    613
9 #> 5  2018 May        875 NA              364       90      169    673
10 #> 6  2018 June       867 NA              342       76      170    610
11 #> 7  2018 July       829 NA              360      108      183    594
12 #> 8  2018 August     939 NA              286       90      205    649
13 #> 9  2018 Septe~     835 NA              304      117      175    560
14 #> # i 1 more variable: West <dbl>
```

This sort of data is not uncommon from government agencies: the column names actually belong to different variables, and here we have summaries for number of units (1, 2-4, 5+) and regions of the country (NE, NW, midwest, S, W). We can most easily describe that with a tibble:

```
1 spec <- tribble(
2   ~.name,           ~.value, ~units, ~region,
3   "1 unit",         "n",    "1",    NA,
4   "2 to 4 units",   "n",    "2-4",   NA,
5   "5 units or more", "n",    "5+",    NA,
6   "Northeast",      "n",    NA,      "Northeast",
7   "Midwest",        "n",    NA,      "Midwest",
8   "South",          "n",    NA,      "South",
9   "West",           "n",    NA,      "West",
10  )
```

Which yields the following longer form:


```

1 construction %>% pivot_longer(spec)
2 #> # A tibble: 63 x 5
3 #>   Year Month   units region      n
4 #>   <dbl> <chr>   <chr> <chr>   <dbl>
5 #> 1  2018 January 1      <NA>      859
6 #> 2  2018 January 2-4    <NA>      NA
7 #> 3  2018 January 5+     <NA>     348
8 #> 4  2018 January <NA> Northeast 114
9 #> 5  2018 January <NA> Midwest   169
10 #> 6  2018 January <NA> South     596
11 #> 7  2018 January <NA> West      339
12 #> 8  2018 February 1      <NA>     882
13 #> 9  2018 February 2-4    <NA>      NA
14 #> 10 2018 February 5+     <NA>     400
15 #> # i 53 more rows

```

Note that there is no overlap between the units and region variables; here the data would really be most naturally described in two independent tables.

0.5.4 Theory

One neat property of the spec is that you need the same spec for `pivot_longer()` and `pivot_wider()`. This makes it very clear that the two operations are symmetric:

```

1 construction %>%
2   pivot_longer(spec) %>%
3   pivot_wider(spec)
4 #> # A tibble: 9 x 9
5 #>   Year Month `1 unit` `2 to 4 units` `5 units or more` Northeast Midwest South
6 #>   <dbl> <chr>   <dbl>         <dbl>         <dbl>         <dbl> <dbl> <dbl>
7 #> 1  2018 Janua~      859           NA           348         114      169    596
8 #> 2  2018 Febru~      882           NA           400         138      160    655
9 #> 3  2018 March      862           NA           356         150      154    595
10 #> 4  2018 April      797           NA           447         144      196    613
11 #> 5  2018 May       875           NA           364          90      169    673
12 #> 6  2018 June      867           NA           342          76      170    610
13 #> 7  2018 July      829           NA           360         108      183    594
14 #> 8  2018 August     939           NA           286          90      205    649
15 #> 9  2018 Septe~      835           NA           304         117      175    560
16 #> # i 1 more variable: West <dbl>

```

The pivoting spec allows us to be more precise about exactly how `pivot_longer(df, spec = spec)` changes the shape of `df`: it will have `nrow(df) * nrow(spec)` rows, and `ncol(df) - nrow(spec) + ncol(spec) - 2` columns.