

3.3.2 Spread and Gather

Note: this section discusses the soon-to-be-deprecated functions `spread()` and `gather()`. These functions will soon be replaced by `pivot_wider()` and `pivot_longer()`. See section [3.3.1 \(page 60\)](#) for code when this happens.

Unfortunately, most of the data you will find in the “wild” is not tidy. So, we need tools to help us tidy unruly data.

The main tools in `tidyr` are the ideas of `spread()` and `gather()`. `gather()` “lengthens” our data, increasing the number of rows and decreasing the number of columns. `spread()` does the opposite, increasing the number of columns and decreasing the number of rows.

These two functions resolve one of two common problems:

1. One variable might be spread across multiple columns. (`gather()`)
2. One observation might be scattered across multiple rows. (`spread()`)

A common issue with data is when values are used as column names.

table4a

```
## # A tibble: 3 x 3
##   country    `1999` `2000`
## * <chr>      <int> <int>
## 1 Afghanistan    745   2666
## 2 Brazil        37737  80488
## 3 China         212258 213766
```

We can fix this using `gather()`.

```
table4a %>%
  gather(-country, key = "year", value = "cases")
```

```
## # A tibble: 6 x 3
##   country    year  cases
##   <chr>      <chr> <int>
## 1 Afghanistan 1999    745
## 2 Brazil      1999  37737
```

```
## 5 Brazil      2000    80488
## 6 China       2000   213766
```

Notice we specified with columns we wanted to consolidate by telling the function the column we *didn't* want to change (`-country`). We can use the `dplyr::select()` syntax here for specifying the columns to pivot.

We can do the same thing with `table4b` and then **join** the databases together by specifying unique identifying attributes.

```
table4a %>%
  gather(-country, key = "year", value = "cases") %>%
  left_join(table4b %>% gather(-country, key = "year", value =
    "population"))
```

```
## Joining, by = c("country", "year")
```

```
## # A tibble: 6 x 4
##   country      year  cases population
##   <chr>      <chr>  <int>      <int>
## 1 Afghanistan 1999      745    19987071
## 2 Brazil      1999    37737    172006362
## 3 China       1999   212258   1272915272
## 4 Afghanistan 2000     2666    20595360
## 5 Brazil      2000    80488   174504898
## 6 China       2000   213766   1280428583
```

If, instead, variables don't have their own column, we can `spread()`.

```
table2
```

```
## # A tibble: 12 x 4
##   country      year type      count
##   <chr>      <int> <chr>      <int>
## 1 Afghanistan 1999 cases         745
## 2 Afghanistan 1999 population 19987071
## 3 Afghanistan 2000 cases         2666
## 4 Afghanistan 2000 population 20595360
## 5 Brazil      1999 cases         37737
## 6 Brazil      1999 population 172006362
## 7 Brazil      2000 cases         80488
## 8 Brazil      2000 population 174504898
```

```
## 9 China      1999 cases      212258
## 10 China     1999 population 1272915272
## 11 China     2000 cases      213766
## 12 China     2000 population 1280428583
```

```
table2 %>%
  spread(key = type, value = count)
```

```
## # A tibble: 6 x 4
##   country      year  cases population
##   <chr>      <int> <int>      <int>
## 1 Afghanistan 1999     745    19987071
## 2 Afghanistan 2000    2666    20595360
## 3 Brazil      1999   37737   172006362
## 4 Brazil      2000   80488   174504898
## 5 China       1999  212258  1272915272
## 6 China       2000  213766  1280428583
```