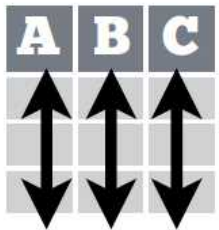


Tidy Data

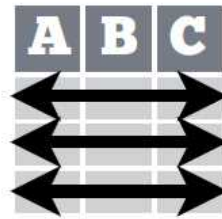


➤ Tidy data: a form of tabular data.

A table is tidy if:

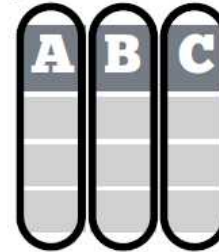


&

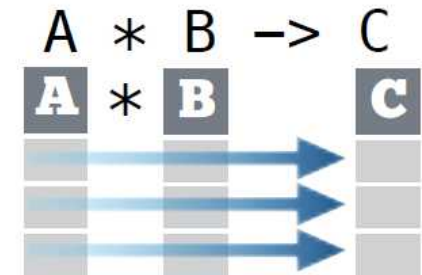


- ① Each **variable** is in its own **column**
- ② Each **observation**, or **case**, is in its own **row**
- ③ Each cell contains a single value

Tidy data:



Makes variables easy to access as vectors



Preserves cases during vectorized operations

Tidy Data



- tidyR 패키지: tidy data로 데이터를 정제하기 위한 다양한 함수 제공

기능	함수
Reshape data	gather(), spread()
Handle missing values	drop_na(), fill()
Expand tables	complete(), expand()
Split cells	Separate(), separate_rows(), unite()

- 실습 데이터셋

- TB(Tuberculosis) cases in Afghanistan, Brazil, and China between 1999 and 2000
- A subset of the data contained in the World Health Organization Global Tuberculosis Report
<https://www.who.int/tb/country/data/download/en/>
- Datasets: table1, table2, table3, table4a, table4b, table5
 - Four variables: country, year, cases and population
 - Each table organizes the values in a different layout



```
> table1
# 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
```

```
> table3
# A tibble: 6 x 3
  country    year rate
  <chr>    <int> <chr>
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
```

```
> 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
```

```
> table4a
# A tibble: 3 x 3
  country `1999` `2000`
  <chr>    <int> <int>
1 Afghanistan 745    2666
2 Brazil      37737 80488
3 China       212258 213766
```

```
> table4b
# A tibble: 3 x 3
  country `1999` `2000`
  <chr>    <int> <int>
1 Afghanistan 19987071 20595360
2 Brazil      172006362 174504898
3 China       1272915272 1280428583
```

```
> table5
# A tibble: 6 x 4
  country century year rate
  <chr>    <chr> <chr> <chr>
1 Afghanistan 19 99 745/19987071
2 Afghanistan 20 00 2666/20595360
3 Brazil      19 99 37737/172006362
4 Brazil      20 00 80488/174504898
5 China       19 99 212258/1272915272
6 China       20 00 213766/1280428583
```

Reshaping data in tidyr



- change the layout of values in a table

gather(data, key, value, ..., na.rm = FALSE,
convert = FALSE, factor_key = FALSE)

gather() moves column names into a **key** column, gathering the column values into a single **value** column.

table4a

country	1999	2000
A	0.7K	2K
B	37K	80K
C	212K	213K

→

country	year	cases
A	1999	0.7K
B	1999	37K
C	1999	212K
A	2000	2K
B	2000	80K
C	2000	213K

key value

```
gather(table4a, `1999`, `2000`,  
key = "year", value = "cases")
```

spread(data, key, value, fill = NA, convert = FALSE,
drop = TRUE, sep = NULL)

spread() moves the unique values of a **key** column into the column names, spreading the values of a **value** column across the new columns.

table2

country	year	type	count
A	1999	cases	0.7K
A	1999	pop	19M
A	2000	cases	2K
A	2000	pop	20M
B	1999	cases	37K
B	1999	pop	172M
B	2000	cases	80K
B	2000	pop	174M
C	1999	cases	212K
C	1999	pop	1T
C	2000	cases	213K
C	2000	pop	1T

key value

```
spread(table2, type, count)
```


Handling missing values in tidyr



drop_na(data, ...)

Drop rows containing NA's in ... columns.

X

x1	x2
A	1
B	NA
C	NA
D	3
E	NA

→

x1	x2
A	1
D	3

drop_na(x, x2)

fill(data, ..., .direction = c("down", "up"))

Fill in NA's in ... columns with most recent non-NA values.

X

x1	x2
A	1
B	NA
C	NA
D	3
E	NA

→

x1	x2
A	1
B	1
C	1
D	3
E	3

fill(x, x2)

replace_na(data, replace = list(), ...)

Replace NA's by column.

X

x1	x2
A	1
B	NA
C	NA
D	3
E	NA

→

x1	x2
A	1
B	2
C	2
D	3
E	2

replace_na(x, list(x2 = 2))

Expanding tables in tidyr



complete(data, ..., fill = list())

Adds to the data missing combinations of the values of the variables listed in ...

(ex) `complete(mtcars, cyl, gear, carb)`

Example:

```
df=tibble(  
  year=c(2010,2010,2010,2010,2012,2012,2012),  
  qtr=c(1,2,3,4,1,2,3),  
  revenue=c(10,20,30,40,NA,60,70)  
)  
df  
df %>% complete(year=full_seq(year,1), qtr)
```

expand(data, ...)

Create new tibble with all possible combinations of the values of the variables listed in ...

(ex) `expand(mtcars, cyl, gear, carb)`

Year	Qtr	Return
2010	1	10
2010	2	20
2010	3	30
2010	4	40
2011	1	NA
2011	2	NA
2011	3	NA
2011	4	NA
2012	1	NA
2012	2	60
2012	3	70
2012	4	NA

Splitting cells in tidyr



```
separate(data, col, into, sep = "[^[:alnum:]]+", remove = TRUE, convert = FALSE, extra = "warn", fill = "warn", ...)
```

Separate each cell in a column to make several columns.

table3

country	year	rate		country	year	cases	pop
A	1999	0.7K/19M	→	A	1999	0.7K	19M
A	2000	2K/20M		A	2000	2K	20M
B	1999	37K/172M		B	1999	37K	172
B	2000	80K/174M		B	2000	80K	174
C	1999	212K/1T		C	1999	212K	1T
C	2000	213K/1T		C	2000	213K	1T

```
separate(table3, rate, sep = "/",  
into = c("cases", "pop"))
```

```
separate_rows(data, ..., sep = "[^[:alnum:]]+", convert = FALSE)
```

Separate each cell in a column to make several rows.

table3

country	year	rate		country	year	rate
A	1999	0.7K/19M	→	A	1999	0.7K
A	2000	2K/20M		A	1999	19M
B	1999	37K/172M		A	2000	2K
B	2000	80K/174M		A	2000	20M
C	1999	212K/1T		B	1999	37K
C	2000	213K/1T		B	1999	172M
				B	2000	80K
				B	2000	174M
				C	1999	212K
				C	1999	1T
				C	2000	213K
				C	2000	1T

```
separate_rows(table3, rate, sep = "/")
```

Uniting cells in tidyr



```
unite(data, col, ..., sep = "_", remove = TRUE)
```

Collapse cells across several columns to make a single column.

table5

country	century	year		country	year
Afghan	19	99	→	Afghan	1999
Afghan	20	00		Afghan	2000
Brazil	19	99		Brazil	1999
Brazil	20	00		Brazil	2000
China	19	99		China	1999
China	20	00		China	2000

```
unite(table5, century, year,  
      col = "year", sep = "")
```