

Ceci n'est pas une pipe.

Next generation programming in R

Florian Uhlitz
uhlitz@hu-berlin.de
uhlitz.github.io

data wrangling



data wrangling

readr



tidyr



dplyr



load

$\%>%$

format

$\%>%$

transform

data analysis



magrittr $\%>%$

report
rmarkdown

visualise
ggplot2

model
base
broom

adapted from H. Wickham





magrittr

In a pipe, the result of the left hand statement is handed over to the function on the right hand side:

$$\begin{aligned} & f(x, y) \\ \Leftrightarrow & x \%>% f(y) \end{aligned}$$

$$\begin{aligned} & f(x, y, z) \\ \Leftrightarrow & x \%>% f(y, z) \end{aligned}$$

$$\begin{aligned} & f2(f1(x), y) \\ \Leftrightarrow & f1(x) \%>% f2(y) \end{aligned}$$

stolen from
~~...idea similar to Unix pipe operator |~~



magrittr

```
summarise(  
  merge(  
    filter(  
      transform(  
        raw_interesting_data, somehow  
      ),  
      the_good_parts  
    ),  
    lookup_table  
  ),  
  result = mean(values)  
)
```

nested
functions

```
raw_interesting_data %>%  
  transform(somehow) %>%  
  filter(the_good_parts) %>%  
  merge(lookup_table) %>%  
  summarise(result = mean(values))
```

chain of
functions

data wrangling

readr



load

`%>%`

tidyr



format

`%>%`

dplyr



`%>%`

transform

data analysis



adapted from H. Wickham



readr, readxl, haven

readr::read_csv()
readr::read_tsv()
readr::read_log()
readr::read_delim()
readr::read_fwf()
readr::read_table()

~10x faster than
base R functions

readxl::read_excel()

+
progress bar

haven::read_sas()
haven::read_spss()
haven::read_stata()

```
> read_tsv("gff3v19.tsv")
```

```
|=====
```

| 68% 730 MB



tidyverse

Formatting

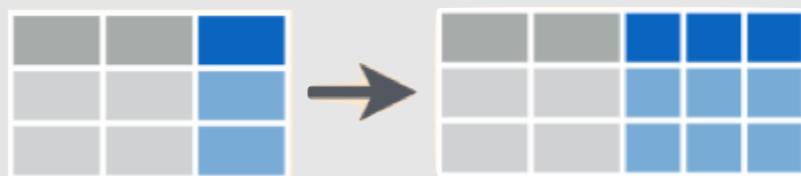
pivot_longer()



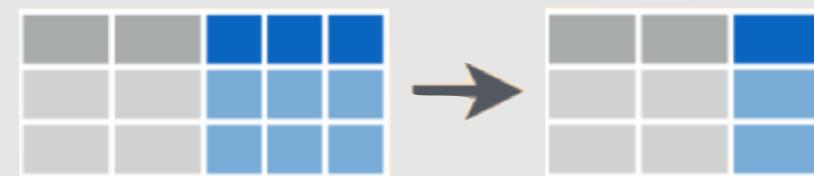
pivot_wider()



separate()



unite()





dplyr

Subsetting

`filter(x > 1)`

x				
1				
2				
3				
1				



x				
2				
3				

`select(B, C, E)`

A	B	C	D	E



B	C	E



dplyr

Transforming



`mutate(z = x + y)`

x	y		
1	4		
2	5		
3	6		

→

x	y		z
1	4		5
2	5		7
3	6		9

Summarising



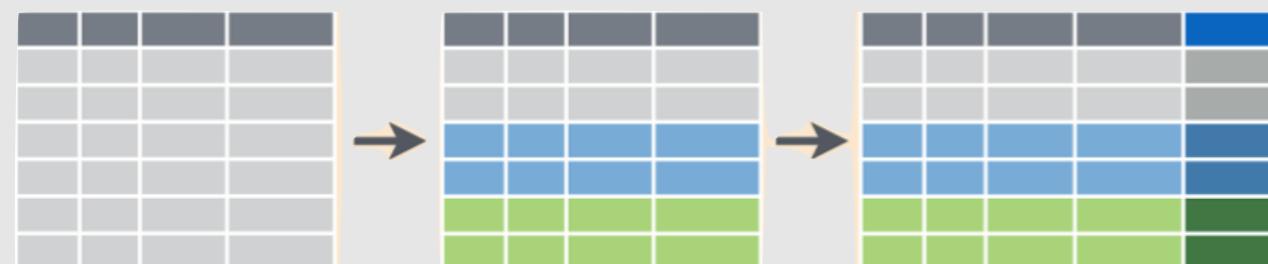
`summarise(A = sum(x), B = sum(y))`

x	y		
1	4		
2	5		
3	6		

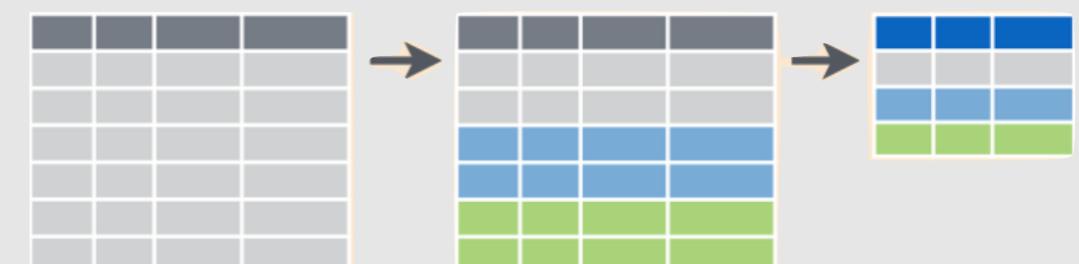
→

A	B
6	15

`group_by() %>% mutate()`



`group_by() %>% summarise()`



What's tidy data?



KEEP
CALM
AND
TIDY
UP

Tidy data principle

»Tidy data sets are all alike; every messy data set is messy in its own way.«

Hadley Wickham

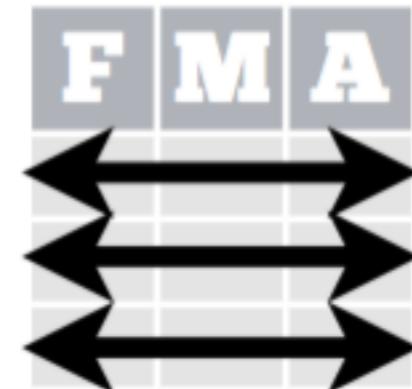
Tidy data definition

In a tidy
data set:



Each **variable** is saved
in its own **column**

&



Each **observation** is
saved in its own **row**

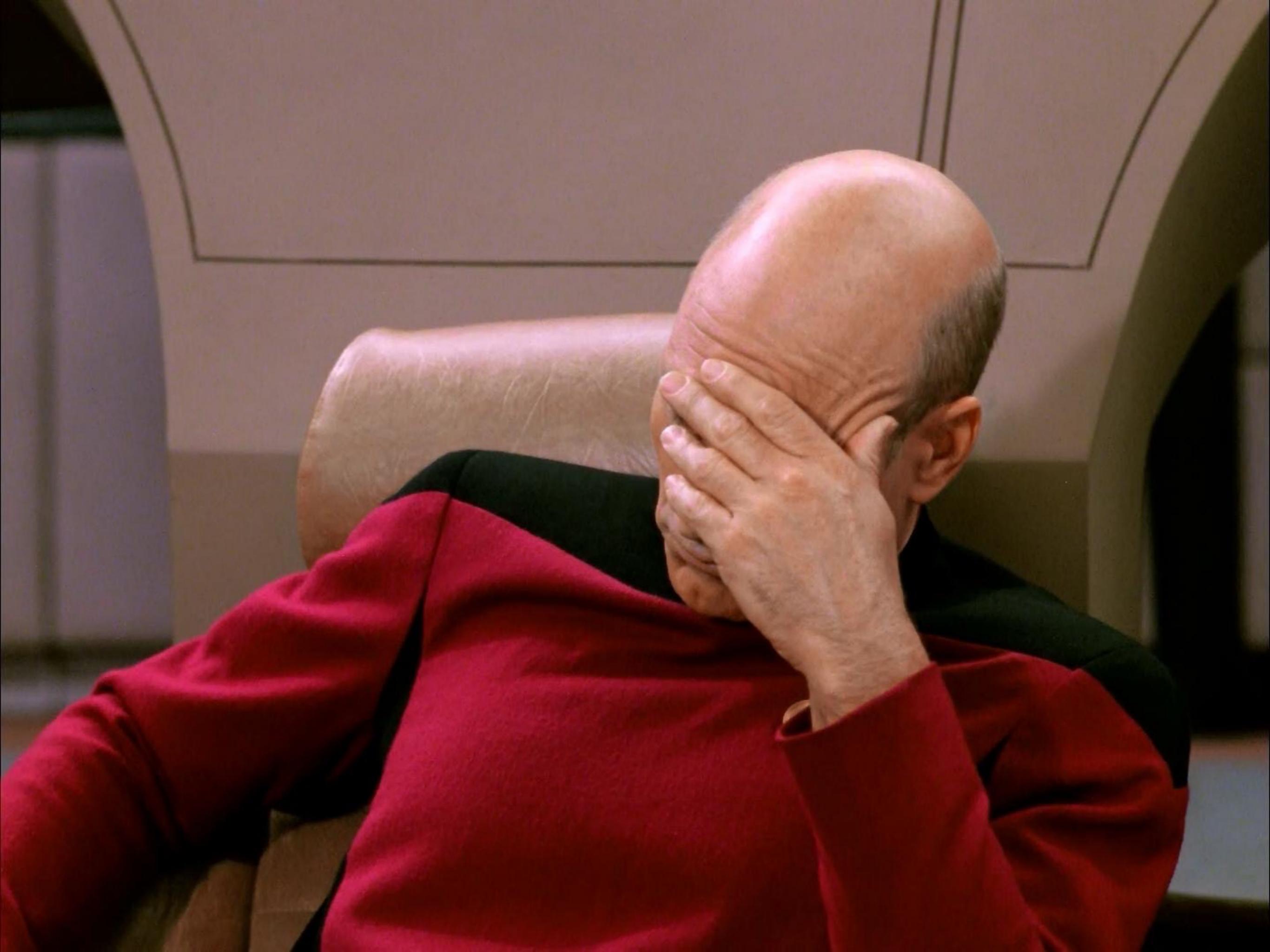
Wickham, H. (2014). Tidy Data. Journal of Statistical Software

Excel screenshot showing a data analysis of protein phosphorylation levels across various conditions and treatments.

The data is organized into several sections:

- Top Row:** AA29, Start, Layout, Tabellen, Diagramme, SmartArt, Formeln, Daten, Überprüfen.
- Search Bar:** Auf dem Blatt suchen.
- Header Row 1:** Contains numerical values (e.g., 63, 7, 16, 27, 12, 5, 8, 4, 2, 3, 3) and labels (e.g., max/min).
- Header Row 2:** Contains numerical values (e.g., 63, 7, 16, 27, 12, 5, 8, 4, 2, 3, 3) and labels (e.g., max/min).
- Section 1 (Rows 3-16):** Compares protein phosphorylation levels (e.g., p-Akt, p-c-Jun, p-ERK2, p-GSK3a/B, p-JNK, p-MEK1, p-p70 S6 Kinase, p-Stat3 (Ser727), p-NF-kB p65, p-p38 MAPK, p-HSP27) under various conditions (Type: X1-X12; Well: A9-D10). Red arrows point to specific rows (X1-X4, X7-X10, X17-X20).
- Section 2 (Rows 17-28):** Compares protein phosphorylation levels under different conditions (Type: X13-X24; Well: E10-H11). Red arrows point to specific rows (X17-X20).
- Section 3 (Rows 29-36):** Normalization of data. It includes a note: "mit wnt und bsa normiert, da höchste std 18%" and provides normalized values for the same set of proteins across different conditions (Wnt oder BSA, Wnt und BSA) and their respective MW and Std%.
- Section 4 (Rows 37-46):** Provides a summary of normalized values for each protein across different conditions (e.g., EGF, IGF, FGF, IFN, BSA) and their respective MW and Std%.

Bottom navigation bar: XII, I, III, IV, V, VI, VII, VIII, IX, X, XI, XII, XIII, XIV, XV, Vgl IX X XI XII.



Microsoft Excel - Sheet1

Auf dem Blatt suchen

Start Layout Tabellen Diagramme SmartArt Formeln Daten Überprüfen

H37 fx

	A	B	C	D	E	F	G	H
1	folder	sample_id			group	condition		
2	Sample_CH_DD_001	1	emp.20 -dox 24h	CH_DD_001	empty empty_cl20	(-)dox	24h	
3	Sample_CH_DD_002	2	emp.20 +dox 24h	CH_DD_002		(+)dox		
4	Sample_CH_DD_003	3	emp.20 -dox 48h	CH_DD_003		(-)dox	48h	
5	Sample_CH_DD_004	4	emp.20 +dox 48h	CH_DD_004		(+)dox		
6	Sample_CH_DD_005	5	emp.20 -dox 72h	CH_DD_005		(-)dox	72h	
7	Sample_CH_DD_006	6	emp.20 +dox 72h	CH_DD_006		(+)dox		
8	Sample_CH_DD_007	7	wt_15 -dox 24h	CH_DD_007	NRas wt_cl15	(-)dox	24h	
9	Sample_CH_DD_008	8	wt_15 +dox 24h	CH_DD_008		(+)dox		
10	Sample_CH_DD_009	9	wt_15 -dox 48h	CH_DD_009		(-)dox	48h	
11	Sample_CH_DD_010	10	wt_15 +dox 48h	CH_DD_010		(+)dox		
12	Sample_CH_DD_011	11	wt_15 -dox 72h	CH_DD_011		(-)dox	72h	
13	Sample_CH_DD_012	12	wt_15 +dox 72h	CH_DD_012		(+)dox		
14	Sample_CH_DD_013	13	12D_6 -dox 24h	CH_DD_013	NRas V12D_cl6	(-)dox	24h	
15	Sample_CH_DD_014	14	12D_6 +dox 24h	CH_DD_014		(+)dox		
16	Sample_CH_DD_015	15	12D_6 -dox 48h	CH_DD_015		(-)dox	48h	
17	Sample_CH_DD_016	16	12D_6 +dox 48h	CH_DD_016		(+)dox		
18	Sample_CH_DD_017	17	12D_6 -dox 72h	CH_DD_017		(-)dox	72h	
19	Sample_CH_DD_018	18	12D_6 +dox 72h	CH_DD_018		(+)dox		
20	Sample_CH_DD_019	19	61K_2 -dox 24h	CH_DD_019	NRas Q61K_cl2	(-)dox	24h	
21	Sample_CH_DD_020	20	61K_2 +dox 24h	CH_DD_020		(+)dox		
22	Sample_CH_DD_021	21	61K_2 -dox 48h	CH_DD_021		(-)dox	48h	
23	Sample_CH_DD_022	22	61K_2 +dox 48h	CH_DD_022		(+)dox		
24	Sample_CH_DD_023	23	61K_2 -dox 72h	CH_DD_023		(-)dox	72h	
25	Sample_CH_DD_024	24	61K_2 +dox 72h	CH_DD_024		(+)dox		
26	Sample_CH_DD_025	25	KRas I -dox 24h	CH_DD_025	KRas G12V_subcl2	(-)dox	24h	
27	Sample_CH_DD_026	26	KRas I +dox 24h	CH_DD_026		(+)dox		
28	Sample_CH_DD_027	27	KRas I -dox 48h	CH_DD_027		(-)dox	48h	
29	Sample_CH_DD_028	28	KRas I +dox 48h	CH_DD_028		(+)dox		
30	Sample_CH_DD_029	29	KRas I -dox 72h	CH_DD_029		(-)dox	72h	
31	Sample_CH_DD_030	30	KRas I +dox 72h	CH_DD_030		(+)dox		
32	Sample_CH_DD_031	31	KRas II -dox 24h	CH_DD_031	KRas G12V_subcl2	(-)dox	24h	
33	Sample_CH_DD_032	32	KRas II +dox 24h	CH_DD_032		(+)dox		
34	Sample_CH_DD_033	33	KRas II -dox 48h	CH_DD_033		(-)dox	48h	
35	Sample_CH_DD_034	34	KRas II +dox 48h	CH_DD_034		(+)dox		
36	Sample_CH_DD_035	35	KRas II -dox 72h	CH_DD_035		(-)dox	72h	
37	Sample_CH_DD_036	36	KRas II +dox 72h	CH_DD_036		(+)dox		
38								
39	(+)-dox--> oncogene is induced							
40								
41								

Microsoft Excel - Sheet1

Auf dem Blatt suchen

Start Layout Tabellen Diagramme SmartArt Formeln Daten Überprüfen

H13 fx

	A	B	C	D	E	F	G	H
1	folder	sample_id			group	condition		
2	Sample_CH_DD_001	1	emp.20 -dox 24h	CH_DD_001	empty empty_cl20	(-)dox	24h	
3	Sample_CH_DD_002	2	emp.20 +dox 24h	CH_DD_002		(+)dox		
4	Sample_CH_DD_003	3	emp.20 -dox 48h	CH_DD_003		(-)dox	48h	
5	Sample_CH_DD_004	4	emp.20 +dox 48h	CH_DD_004		(+)dox		
6	Sample_CH_DD_005	5	emp.20 -dox 72h	CH_DD_005		(-)dox	72h	
7	Sample_CH_DD_006	6	emp.20 +dox 72h	CH_DD_006		(+)dox		
8	Sample_CH_DD_007	7	wt_15 -dox 24h	CH_DD_007	NRas wt_cl15	(-)dox	24h	
9	Sample_CH_DD_008	8	wt_15 +dox 24h	CH_DD_008		(+)dox		
10	Sample_CH_DD_009	9	wt_15 -dox 48h	CH_DD_009		(-)dox	48h	
11	Sample_CH_DD_010	10	wt_15 +dox 48h	CH_DD_010		(+)dox		
12	Sample_CH_DD_011	11	wt_15 -dox 72h	CH_DD_011		(-)dox	72h	
13	Sample_CH_DD_012	12	wt_15 +dox 72h	CH_DD_012		(+)dox		
14	Sample_CH_DD_013	13	12D_6 -dox 24h	CH_DD_013	NRas V12D_cl6	(-)dox	24h	
15	Sample_CH_DD_014	14	12D_6 +dox 24h	CH_DD_014		(+)dox		
16	Sample_CH_DD_015	15	12D_6 -dox 48h	CH_DD_015		(-)dox	48h	
17	Sample_CH_DD_016	16	12D_6 +dox 48h	CH_DD_016		(+)dox		
18	Sample_CH_DD_017	17	12D_6 -dox 72h	CH_DD_017		(-)dox	72h	
19	Sample_CH_DD_018	18	12D_6 +dox 72h	CH_DD_018		(+)dox		
20	Sample_CH_DD_019	19	61K_2 -dox 24h	CH_DD_019	NRas Q61K_cl2	(-)dox	24h	
21	Sample_CH_DD_020	20	61K_2 +dox 24h	CH_DD_020		(+)dox		
22	Sample_CH_DD_021	21	61K_2 -dox 48h	CH_DD_021		(-)dox	48h	
23	Sample_CH_DD_022	22	61K_2 +dox 48h	CH_DD_022		(+)dox		
24	Sample_CH_DD_023	23	61K_2 -dox 72h	CH_DD_023		(-)dox	72h	
25	Sample_CH_DD_024	24	61K_2 +dox 72h	CH_DD_024		(+)dox		
26	Sample_CH_DD_025	25	KRas I -dox 24h	CH_DD_025	KRas G12V_subcl2	(-)dox	24h	
27	Sample_CH_DD_026	26	KRas I +dox 24h	CH_DD_026		(+)dox		
28	Sample_CH_DD_027	27	KRas I -dox 48h	CH_DD_027		(-)dox	48h	
29	Sample_CH_DD_028	28	KRas I +dox 48h	CH_DD_028		(+)dox		
30	Sample_CH_DD_029	29	KRas I -dox 72h	CH_DD_029		(-)dox	72h	
31	Sample_CH_DD_030	30	KRas I +dox 72h	CH_DD_030		(+)dox		
32	Sample_CH_DD_031	31	KRas II -dox 24h	CH_DD_031	KRas G12V_subcl2	(-)dox	24h	
33	Sample_CH_DD_032	32	KRas II +dox 24h	CH_DD_032		(+)dox		
34	Sample_CH_DD_033	33	KRas II -dox 48h	CH_DD_033		(-)dox	48h	
35	Sample_CH_DD_034	34	KRas II +dox 48h	CH_DD_034		(+)dox		
36	Sample_CH_DD_035	35	KRas II -dox 72h	CH_DD_035		(-)dox	72h	
37	Sample_CH_DD_036	36	KRas II +dox 72h	CH_DD_036		(+)dox		
38								
39	(+)-dox--> oncogene is induced							
40								
41								

Office ribbon tabs: Start, Layout, Tabellen, Diagramme, SmartArt, Formeln, Daten, Überprüfen.

Cell D39 contains the formula: =fx

Cell A39 contains the text: (+)dox--> oncogene is induced

	A	B	C	D	E	F	G	H
1	folder	sample_id			group	condition		
2	Sample_CH_DD_001	1	emp.20 -dox 24h	CH_DD_001	empty empty_cl20	(-)dox	24h	
3	Sample_CH_DD_002	2	emp.20 +dox 24h	CH_DD_002		(+)dox		
4	Sample_CH_DD_003	3	emp.20 -dox 48h	CH_DD_003		(-)dox	48h	
5	Sample_CH_DD_004	4	emp.20 +dox 48h	CH_DD_004		(+)dox		
6	Sample_CH_DD_005	5	emp.20 -dox 72h	CH_DD_005		(-)dox	72h	
7	Sample_CH_DD_006	6	emp.20 +dox 72h	CH_DD_006		(+)dox		
8	Sample_CH_DD_007	7	wt_15 -dox 24h	CH_DD_007	NRas wt_cl15	(-)dox	24h	
9	Sample_CH_DD_008	8	wt_15 +dox 24h	CH_DD_008		(+)dox		
10	Sample_CH_DD_009	9	wt_15 -dox 48h	CH_DD_009		(-)dox	48h	
11	Sample_CH_DD_010	10	wt_15 +dox 48h	CH_DD_010		(+)dox		
12	Sample_CH_DD_011	11	wt_15 -dox 72h	CH_DD_011		(-)dox	72h	
13	Sample_CH_DD_012	12	wt_15 +dox 72h	CH_DD_012		(+)dox		
14	Sample_CH_DD_013	13	12D_6 -dox 24h	CH_DD_013	NRas V12D_cl6	(-)dox	24h	
15	Sample_CH_DD_014	14	12D_6 +dox 24h	CH_DD_014		(+)dox		
16	Sample_CH_DD_015	15	12D_6 -dox 48h	CH_DD_015		(-)dox	48h	
17	Sample_CH_DD_016	16	12D_6 +dox 48h	CH_DD_016		(+)dox		
18	Sample_CH_DD_017	17	12D_6 -dox 72h	CH_DD_017		(-)dox	72h	
19	Sample_CH_DD_018	18	12D_6 +dox 72h	CH_DD_018		(+)dox		
20	Sample_CH_DD_019	19	61K_2 -dox 24h	CH_DD_019	NRas Q61K_cl2	(-)dox	24h	
21	Sample_CH_DD_020	20	61K_2 +dox 24h	CH_DD_020		(+)dox		
22	Sample_CH_DD_021	21	61K_2 -dox 48h	CH_DD_021		(-)dox	48h	
23	Sample_CH_DD_022	22	61K_2 +dox 48h	CH_DD_022		(+)dox		
24	Sample_CH_DD_023	23	61K_2 -dox 72h	CH_DD_023		(-)dox	72h	
25	Sample_CH_DD_024	24	61K_2 +dox 72h	CH_DD_024		(+)dox		
26	Sample_CH_DD_025	25	KRas I -dox 24h	CH_DD_025	KRas G12V_subcl2	(-)dox	24h	
27	Sample_CH_DD_026	26	KRas I +dox 24h	CH_DD_026		(+)dox		
28	Sample_CH_DD_027	27	KRas I -dox 48h	CH_DD_027		(-)dox	48h	
29	Sample_CH_DD_028	28	KRas I +dox 48h	CH_DD_028		(+)dox		
30	Sample_CH_DD_029	29	KRas I -dox 72h	CH_DD_029		(-)dox	72h	
31	Sample_CH_DD_030	30	KRas I +dox 72h	CH_DD_030		(+)dox		
32	Sample_CH_DD_031	31	KRas II -dox 24h	CH_DD_031	KRas G12V_subcl2	(-)dox	24h	
33	Sample_CH_DD_032	32	KRas II +dox 24h	CH_DD_032		(+)dox		
34	Sample_CH_DD_033	33	KRas II -dox 48h	CH_DD_033		(-)dox	48h	
35	Sample_CH_DD_034	34	KRas II +dox 48h	CH_DD_034		(+)dox		
36	Sample_CH_DD_035	35	KRas II -dox 72h	CH_DD_035		(-)dox	72h	
37	Sample_CH_DD_036	36	KRas II +dox 72h	CH_DD_036		(+)dox		
38								
39	(+)-dox--> oncogene is induced							
40								
41								

Microsoft Excel - Sheet1

Auf dem Blatt suchen

Start Layout Tabellen Diagramme SmartArt Formeln Daten Überprüfen

H17 fx

	A	B	C	D	E	F	G	H
1	folder	sample_id			group	condition		
2	Sample_CH_DD_001	1	emp.20 -dox 24h	CH_DD_001	empty empty_cl20	(-)dox	24h	
3	Sample_CH_DD_002	2	emp.20 +dox 24h	CH_DD_002		(+)dox		
4	Sample_CH_DD_003	3	emp.20 -dox 48h	CH_DD_003		(-)dox	48h	
5	Sample_CH_DD_004	4	emp.20 +dox 48h	CH_DD_004		(+)dox		
6	Sample_CH_DD_005	5	emp.20 -dox 72h	CH_DD_005		(-)dox	72h	
7	Sample_CH_DD_006	6	emp.20 +dox 72h	CH_DD_006		(+)dox		
8	Sample_CH_DD_007	7	wt_15 -dox 24h	CH_DD_007	NRas wt_cl15	(-)dox	24h	
9	Sample_CH_DD_008	8	wt_15 +dox 24h	CH_DD_008		(+)dox		
10	Sample_CH_DD_009	9	wt_15 -dox 48h	CH_DD_009		(-)dox	48h	
11	Sample_CH_DD_010	10	wt_15 +dox 48h	CH_DD_010		(+)dox		
12	Sample_CH_DD_011	11	wt_15 -dox 72h	CH_DD_011		(-)dox	72h	
13	Sample_CH_DD_012	12	wt_15 +dox 72h	CH_DD_012		(+)dox		
14	Sample_CH_DD_013	13	12D_6 -dox 24h	CH_DD_013	NRas V12D_cl6	(-)dox	24h	
15	Sample_CH_DD_014	14	12D_6 +dox 24h	CH_DD_014		(+)dox		
16	Sample_CH_DD_015	15	12D_6 -dox 48h	CH_DD_015		(-)dox	48h	
17	Sample_CH_DD_016	16	12D_6 +dox 48h	CH_DD_016		(+)dox		
18	Sample_CH_DD_017	17	12D_6 -dox 72h	CH_DD_017		(-)dox	72h	
19	Sample_CH_DD_018	18	12D_6 +dox 72h	CH_DD_018		(+)dox		
20	Sample_CH_DD_019	19	61K_2 -dox 24h	CH_DD_019	NRas Q61K_cl2	(-)dox	24h	
21	Sample_CH_DD_020	20	61K_2 +dox 24h	CH_DD_020		(+)dox		
22	Sample_CH_DD_021	21	61K_2 -dox 48h	CH_DD_021		(-)dox	48h	
23	Sample_CH_DD_022	22	61K_2 +dox 48h	CH_DD_022		(+)dox		
24	Sample_CH_DD_023	23	61K_2 -dox 72h	CH_DD_023		(-)dox	72h	
25	Sample_CH_DD_024	24	61K_2 +dox 72h	CH_DD_024		(+)dox		
26	Sample_CH_DD_025	25	KRas I -dox 24h	CH_DD_025	KRas G12V_subcl2	(-)dox	24h	
27	Sample_CH_DD_026	26	KRas I +dox 24h	CH_DD_026		(+)dox		
28	Sample_CH_DD_027	27	KRas I -dox 48h	CH_DD_027		(-)dox	48h	
29	Sample_CH_DD_028	28	KRas I +dox 48h	CH_DD_028		(+)dox		
30	Sample_CH_DD_029	29	KRas I -dox 72h	CH_DD_029		(-)dox	72h	
31	Sample_CH_DD_030	30	KRas I +dox 72h	CH_DD_030		(+)dox		
32	Sample_CH_DD_031	31	KRas II -dox 24h	CH_DD_031	KRas G12V_subcl2	(-)dox	24h	
33	Sample_CH_DD_032	32	KRas II +dox 24h	CH_DD_032		(+)dox		
34	Sample_CH_DD_033	33	KRas II -dox 48h	CH_DD_033		(-)dox	48h	
35	Sample_CH_DD_034	34	KRas II +dox 48h	CH_DD_034		(+)dox		
36	Sample_CH_DD_035	35	KRas II -dox 72h	CH_DD_035		(-)dox	72h	
37	Sample_CH_DD_036	36	KRas II +dox 72h	CH_DD_036		(+)dox		
38								
39	(+)-dox--> oncogene is induced							
40								
41								

	A	B	C	D	E	F	G	H
1	folder	sample_id			group	condition		
2	Sample_CH_DD_001	1	emp.20 -dox 24h	CH_DD_001	empty empty_cl20	(-)dox	24h	
3	Sample_CH_DD_002	2	emp.20 +dox 24h	CH_DD_002		(+)dox		
4	Sample_CH_DD_003	3	emp.20 -dox 48h	CH_DD_003		(-)dox	48h	
5	Sample_CH_DD_004	4	emp.20 +dox 48h	CH_DD_004		(+)dox		
6	Sample_CH_DD_005	5	emp.20 -dox 72h	CH_DD_005		(-)dox	72h	
7	Sample_CH_DD_006	6	emp.20 +dox 72h	CH_DD_006		(+)dox		
8	Sample_CH_DD_007	7	wt_15 -dox 24h	CH_DD_007	NRas wt_cl15	(-)dox	24h	
9	Sample_CH_DD_008	8	wt_15 +dox 24h	CH_DD_008		(+)dox		
10	Sample_CH_DD_009	9	wt_15 -dox 48h	CH_DD_009		(-)dox	48h	
11	Sample_CH_DD_010	10	wt_15 +dox 48h	CH_DD_010		(+)dox		
12	Sample_CH_DD_011	11	wt_15 -dox 72h	CH_DD_011		(-)dox	72h	
13	Sample_CH_DD_012	12	wt_15 +dox 72h	CH_DD_012		(+)dox		
14	Sample_CH_DD_013	13	12D_6 -dox 24h	CH_DD_013	NRas V12D_cl6	(-)dox	24h	
15	Sample_CH_DD_014	14	12D_6 +dox 24h	CH_DD_014		(+)dox		
16	Sample_CH_DD_015	15	12D_6 -dox 48h	CH_DD_015		(-)dox	48h	
17	Sample_CH_DD_016	16	12D_6 +dox 48h	CH_DD_016		(+)dox		
18	Sample_CH_DD_017	17	12D_6 -dox 72h	CH_DD_017		(-)dox	72h	
19	Sample_CH_DD_018	18	12D_6 +dox 72h	CH_DD_018		(+)dox		
20	Sample_CH_DD_019	19	61K_2 -dox 24h	CH_DD_019	NRas Q61K_cl2	(-)dox	24h	
21	Sample_CH_DD_020	20	61K_2 +dox 24h	CH_DD_020		(+)dox		
22	Sample_CH_DD_021	21	61K_2 -dox 48h	CH_DD_021		(-)dox	48h	
23	Sample_CH_DD_022	22	61K_2 +dox 48h	CH_DD_022		(+)dox		
24	Sample_CH_DD_023	23	61K_2 -dox 72h	CH_DD_023		(-)dox	72h	
25	Sample_CH_DD_024	24	61K_2 +dox 72h	CH_DD_024		(+)dox		
26	Sample_CH_DD_025	25	KRas I -dox 24h	CH_DD_025	KRas G12V_subcl2	(-)dox	24h	
27	Sample_CH_DD_026	26	KRas I +dox 24h	CH_DD_026		(+)dox		
28	Sample_CH_DD_027	27	KRas I -dox 48h	CH_DD_027		(-)dox	48h	
29	Sample_CH_DD_028	28	KRas I +dox 48h	CH_DD_028		(+)dox		
30	Sample_CH_DD_029	29	KRas I -dox 72h	CH_DD_029		(-)dox	72h	
31	Sample_CH_DD_030	30	KRas I +dox 72h	CH_DD_030		(+)dox		
32	Sample_CH_DD_031	31	KRas II -dox 24h	CH_DD_031	KRas G12V_subcl2	(-)dox	24h	
33	Sample_CH_DD_032	32	KRas II +dox 24h	CH_DD_032		(+)dox		
34	Sample_CH_DD_033	33	KRas II -dox 48h	CH_DD_033		(-)dox	48h	
35	Sample_CH_DD_034	34	KRas II +dox 48h	CH_DD_034		(+)dox		
36	Sample_CH_DD_035	35	KRas II -dox 72h	CH_DD_035		(-)dox	72h	
37	Sample_CH_DD_036	36	KRas II +dox 72h	CH_DD_036		(+)dox		
38								
39			(+)dox--> oncogene is induced					
40								
41								

```

read_excel("untidy_data.xlsx") %>%
  set_colnames(mynames) %>%
  slice(1:36) %>%
  fill(group, condition) %>%
  separate(group, into = c("Gene", "Mutation", "clone"), sep = "_") %>%
  write_tsv("tidy_data.tsv")
  
```

The screenshot shows a software application window with a toolbar at the top containing icons for file operations like Open, Save, and Print, along with a 'Go to file/function' search bar and an 'Addins' dropdown. The main area is a data grid titled 'Filter' with a magnifying glass icon. The grid has columns labeled: Sample_ID, Gene, Mutation, Clone, dox, time, replicate, Sample_Name, and dox_logical. The data consists of 19 rows numbered 1 to 18, with row 19 being a dashed separator. The data entries are as follows:

	Sample_ID	Gene	Mutation	Clone	dox	time	replicate	Sample_Name	dox_logical
1	CH_DD_001	empty	empty	cl20	(-)	24h	1	emp_20 -dox 24h	FALSE
2	CH_DD_002	empty	empty	cl20	(+)	24h	1	emp_20 +dox 24h	TRUE
3	CH_DD_003	empty	empty	cl20	(-)	48h	1	emp_20 -dox 48h	FALSE
4	CH_DD_004	empty	empty	cl20	(+)	48h	1	emp_20 +dox 48h	TRUE
5	CH_DD_005	empty	empty	cl20	(-)	72h	1	emp_20 -dox 72h	FALSE
6	CH_DD_006	empty	empty	cl20	(+)	72h	1	emp_20 +dox 72h	TRUE
7	CH_DD_007	NRas	wt	cl15	(-)	24h	1	wt_15 -dox 24h	FALSE
8	CH_DD_008	NRas	wt	cl15	(+)	24h	1	wt_15 +dox 24h	TRUE
9	CH_DD_009	NRas	wt	cl15	(-)	48h	1	wt_15 -dox 48h	FALSE
10	CH_DD_010	NRas	wt	cl15	(+)	48h	1	wt_15 +dox 48h	TRUE
11	CH_DD_011	NRas	wt	cl15	(-)	72h	1	wt_15 -dox 72h	FALSE
12	CH_DD_012	NRas	wt	cl15	(+)	72h	1	wt_15 +dox 72h	TRUE
13	CH_DD_013	NRas	G12D	cl6	(-)	24h	1	12D_6 -dox 24h	FALSE
14	CH_DD_014	NRas	G12D	cl6	(+)	24h	1	12D_6 +dox 24h	TRUE
15	CH_DD_015	NRas	G12D	cl6	(-)	48h	1	12D_6 -dox 48h	FALSE
16	CH_DD_016	NRas	G12D	cl6	(+)	48h	1	12D_6 +dox 48h	TRUE
17	CH_DD_017	NRas	G12D	cl6	(-)	72h	1	12D_6 -dox 72h	FALSE
18	CH_DD_018	NRas	G12D	cl6	(+)	72h	1	12D_6 +dox 72h	TRUE

Showing 1 to 19 of 36 entries

```
read_excel("untidy_data.xlsx") %>%
  set_colnames(mynames) %>%
  slice(1:36) %>%
  fill(group, condition) %>%
  separate(group, into = c("Gene", "Mutation", "clone"), sep = "_") %>%
  write_tsv("tidy_data.tsv")
```

read_excel

The screenshot shows the RStudio interface with the following details:

- Toolbar:** Includes icons for New Project, Open Project, Save Project, Print, Go to file/function, and a search bar.
- Console Tab:** Labeled "Console ~/Projects/R/Current/dplyr_tut/".
- Code Input:** Shows the command `> read_excel("annotation_rnaseq_collapsed.xlsx")` and its output: "Source: local data frame [38 x 7]".
- Data Preview:** Displays the first 10 rows of the data frame. The columns are labeled as follows:

	folder	sample_id	NA	NA	group	condition	NA					
	(chr)	(dbl)	(chr)	(chr)	(chr)	(chr)	(chr)					
1	Sample_CH_DD_001		1	emp.20	-dox	24h	CH_DD_001	empty	empty_c120	(-)dox	24h	
2	Sample_CH_DD_002		2	emp.20	+dox	24h	CH_DD_002			NA	(+)dox	NA
3	Sample_CH_DD_003		3	emp.20	-dox	48h	CH_DD_003			NA	(-)dox	48h
4	Sample_CH_DD_004		4	emp.20	+dox	48h	CH_DD_004			NA	(+)dox	NA
5	Sample_CH_DD_005		5	emp.20	-dox	72h	CH_DD_005			NA	(-)dox	72h
6	Sample_CH_DD_006		6	emp.20	+dox	72h	CH_DD_006			NA	(+)dox	NA
7	Sample_CH_DD_007		7	wt_15	-dox	24h	CH_DD_007	NRas	wt_c115	(-)dox	24h	
8	Sample_CH_DD_008		8	wt_15	+dox	24h	CH_DD_008			NA	(+)dox	NA
9	Sample_CH_DD_009		9	wt_15	-dox	48h	CH_DD_009			NA	(-)dox	48h
10	Sample_CH_DD_010		10	wt_15	+dox	48h	CH_DD_010			NA	(+)dox	NA
- Ellipsis:** Shows three ellipsis characters ("...") at the bottom of the data preview.
- Bottom Navigation:** Includes tabs for Environment, Files, Plots, Packages, Help, Git, and Viewer.



```
read_excel %>% set_colnames
```

The screenshot shows the RStudio interface with the following details:

- Toolbar:** Includes icons for file operations (New, Open, Save, Print), Go to file/function, and a search bar.
- Console Tab:** Labeled "dplyr_tut".
- Console Output:** Displays R code and its execution results. The code reads an Excel file and sets column names. The resulting data frame has 38 rows and 7 columns, with sample data for rows 1 through 10.
- Data Frame Preview:** Shows the first 10 rows of the data frame with columns X1 through X7.

```
>
>
>
>
> read_excel("annotation_rnaseq_collapsed.xlsx") %>%
+   set_colnames(paste0("X", 1:7))
Source: local data frame [38 x 7]

      X1     X2     X3     X4     X5     X6     X7
      (chr) (dbl)  (chr)  (chr)  (chr)  (chr)  (chr)
1 Sample_CH_DD_001    1 emp.20 -dox 24h CH_DD_001 empty empty_cl20 (-)dox 24h
2 Sample_CH_DD_002    2 emp.20 +dox 24h CH_DD_002 NA (+)dox NA
3 Sample_CH_DD_003    3 emp.20 -dox 48h CH_DD_003 NA (-)dox 48h
4 Sample_CH_DD_004    4 emp.20 +dox 48h CH_DD_004 NA (+)dox NA
5 Sample_CH_DD_005    5 emp.20 -dox 72h CH_DD_005 NA (-)dox 72h
6 Sample_CH_DD_006    6 emp.20 +dox 72h CH_DD_006 NA (+)dox NA
7 Sample_CH_DD_007    7 wt_15 -dox 24h CH_DD_007 NRas wt_cl15 (-)dox 24h
8 Sample_CH_DD_008    8 wt_15 +dox 24h CH_DD_008 NA (+)dox NA
9 Sample_CH_DD_009    9 wt_15 -dox 48h CH_DD_009 NA (-)dox 48h
10 Sample_CH_DD_010   10 wt_15 +dox 48h CH_DD_010 NA (+)dox NA
...
> |
```



```
read_excel %>% set_colnames %>% tail
```

The screenshot shows the RStudio interface with the following details:

- Toolbar:** Includes icons for file operations (New, Open, Save, Print), Go to file/function, and a search bar.
- Console Tab:** Labeled "dplyr_tut".
- Console Area:** Shows the R code and its output.

```
>  
>  
>  
> read_excel("annotation_rnaseq_collapsed.xlsx") %>%  
+   set_colnames(paste0("X",1:7)) %>%  
+   tail(10)  
Source: local data frame [10 x 7]  
  
          X1      X2          X3      X4          X5      X6      X7  
          (chr)  (dbl)      (chr)  (chr)      (chr)  (chr)  (chr)  
1 Sample_CH_DD_029     29 KRas I -dox 72h CH_DD_029      NA (-)dox 72h  
2 Sample_CH_DD_030     30 KRas I +dox 72h CH_DD_030      NA (+)dox NA  
3 Sample_CH_DD_031     31 KRas II -dox 24h CH_DD_031 KRas G12V_subcl2 (-)dox 24h  
4 Sample_CH_DD_032     32 KRas II +dox 24h CH_DD_032      NA (+)dox NA  
5 Sample_CH_DD_033     33 KRas II -dox 484h CH_DD_033      NA (-)dox 48h  
6 Sample_CH_DD_034     34 KRas II +dox 48h CH_DD_034      NA (+)dox NA  
7 Sample_CH_DD_035     35 KRas II -dox 72h CH_DD_035      NA (-)dox 72h  
8 Sample_CH_DD_036     36 KRas II +dox 72h CH_DD_036      NA (+)dox NA  
9 (+)dox--> oncogene is induced      NA      NA      NA      NA      NA      NA  
10                      NA      NA      NA      NA      NA      NA      NA
```
- Environment, Files, Plots, Packages, Help, Git, Viewer:** Standard RStudio navigation tabs.



```
read_excel %>% set_colnames
```

The screenshot shows the RStudio interface with the following details:

- Toolbar:** Includes icons for file operations (New, Open, Save, Print), Go to file/function, and a search bar.
- Console Tab:** Labeled "dplyr_tut".
- Console Output:** Displays R code and its execution results. The code reads an Excel file and sets column names. The resulting data frame has 38 rows and 7 columns, with sample data for rows 1 through 10.
- Data Frame Preview:** Shows the first 10 rows of the data frame with columns X1 through X7.

```
>
>
>
>
> read_excel("annotation_rnaseq_collapsed.xlsx") %>%
+   set_colnames(paste0("X", 1:7))
Source: local data frame [38 x 7]

      X1     X2     X3     X4     X5     X6     X7
      (chr) (dbl)  (chr)  (chr)  (chr)  (chr)  (chr)
1 Sample_CH_DD_001    1 emp.20 -dox 24h CH_DD_001 empty empty_cl20 (-)dox 24h
2 Sample_CH_DD_002    2 emp.20 +dox 24h CH_DD_002 NA (+)dox NA
3 Sample_CH_DD_003    3 emp.20 -dox 48h CH_DD_003 NA (-)dox 48h
4 Sample_CH_DD_004    4 emp.20 +dox 48h CH_DD_004 NA (+)dox NA
5 Sample_CH_DD_005    5 emp.20 -dox 72h CH_DD_005 NA (-)dox 72h
6 Sample_CH_DD_006    6 emp.20 +dox 72h CH_DD_006 NA (+)dox NA
7 Sample_CH_DD_007    7 wt_15 -dox 24h CH_DD_007 NRas wt_cl15 (-)dox 24h
8 Sample_CH_DD_008    8 wt_15 +dox 24h CH_DD_008 NA (+)dox NA
9 Sample_CH_DD_009    9 wt_15 -dox 48h CH_DD_009 NA (-)dox 48h
10 Sample_CH_DD_010   10 wt_15 +dox 48h CH_DD_010 NA (+)dox NA
...
> |
```



```
read_excel %>% set_colnames %>% slice
```

The screenshot shows the RStudio IDE interface with the following details:

- Toolbar:** Includes icons for file operations (New, Open, Save, Print), Go to file/function, and a search bar.
- Console Tab:** Labeled "dplyr_tut".
- Console Output:**

```
>
>
>
> read_excel("annotation_rnaseq_collapsed.xlsx") %>%
+   set_colnames(paste0("X",1:7)) %>%
+   slice(1:36)
Source: local data frame [36 x 7]

      X1      X2      X3      X4      X5      X6      X7
      (chr)  (dbl)  (chr)  (chr)  (chr)  (chr)  (chr)
1 Sample_CH_DD_001    1 emp.20 -dox 24h CH_DD_001 empty empty_cl20 (-)dox 24h
2 Sample_CH_DD_002    2 emp.20 +dox 24h CH_DD_002 NA (+)dox NA
3 Sample_CH_DD_003    3 emp.20 -dox 48h CH_DD_003 NA (-)dox 48h
4 Sample_CH_DD_004    4 emp.20 +dox 48h CH_DD_004 NA (+)dox NA
5 Sample_CH_DD_005    5 emp.20 -dox 72h CH_DD_005 NA (-)dox 72h
6 Sample_CH_DD_006    6 emp.20 +dox 72h CH_DD_006 NA (+)dox NA
7 Sample_CH_DD_007    7 wt_15 -dox 24h CH_DD_007 NRas wt_cl15 (-)dox 24h
8 Sample_CH_DD_008    8 wt_15 +dox 24h CH_DD_008 NA (+)dox NA
9 Sample_CH_DD_009    9 wt_15 -dox 48h CH_DD_009 NA (-)dox 48h
10 Sample_CH_DD_010   10 wt_15 +dox 48h CH_DD_010 NA (+)dox NA
...
> |
```
- Data View:** A preview of the data frame structure is shown below the console output.

	X1	X2	X3	X4	X5	X6	X7			
	(chr)	(dbl)	(chr)	(chr)	(chr)	(chr)	(chr)			
1	Sample_CH_DD_001	1	emp.20	-dox	24h	CH_DD_001	empty	empty_cl20	(-)dox	24h
2	Sample_CH_DD_002	2	emp.20	+dox	24h	CH_DD_002	NA	(+)dox	NA	
3	Sample_CH_DD_003	3	emp.20	-dox	48h	CH_DD_003	NA	(-)dox	48h	
4	Sample_CH_DD_004	4	emp.20	+dox	48h	CH_DD_004	NA	(+)dox	NA	
5	Sample_CH_DD_005	5	emp.20	-dox	72h	CH_DD_005	NA	(-)dox	72h	
6	Sample_CH_DD_006	6	emp.20	+dox	72h	CH_DD_006	NA	(+)dox	NA	
7	Sample_CH_DD_007	7	wt_15	-dox	24h	CH_DD_007	NRas	wt_cl15	(-)dox	24h
8	Sample_CH_DD_008	8	wt_15	+dox	24h	CH_DD_008	NA	(+)dox	NA	
9	Sample_CH_DD_009	9	wt_15	-dox	48h	CH_DD_009	NA	(-)dox	48h	
10	Sample_CH_DD_010	10	wt_15	+dox	48h	CH_DD_010	NA	(+)dox	NA	
...



```
read_excel %>% set_colnames %>% slice %>% fill
```

The screenshot shows the RStudio interface with the following details:

- Toolbar:** Includes icons for file operations (New, Open, Save, Print), Go to file/function, and a search bar.
- Console Tab:** Labeled "dplyr_tut".
- Console Output:** Displays the R code used to read and manipulate the data, followed by the resulting data frame structure and its contents.
- Data Frame Preview:** Shows the first 10 rows of the data frame with columns X1 through X7.

```
>  
>  
> read_excel("annotation_rnaseq_collapsed.xlsx") %>%  
+   set_colnames(paste0("X",1:7)) %>%  
+   slice(1:36) %>%  
+   fill(X5:X7)  
Source: local data frame [36 x 7]  
  
          X1      X2          X3      X4          X5      X6      X7  
          (chr)  (dbl)      (chr)  (chr)      (chr)  (chr)  (chr)  
1 Sample_CH_DD_001    1 emp.20 -dox 24h CH_DD_001 empty empty_cl20 (-)dox 24h  
2 Sample_CH_DD_002    2 emp.20 +dox 24h CH_DD_002 empty empty_cl20 (+)dox 24h  
3 Sample_CH_DD_003    3 emp.20 -dox 48h CH_DD_003 empty empty_cl20 (-)dox 48h  
4 Sample_CH_DD_004    4 emp.20 +dox 48h CH_DD_004 empty empty_cl20 (+)dox 48h  
5 Sample_CH_DD_005    5 emp.20 -dox 72h CH_DD_005 empty empty_cl20 (-)dox 72h  
6 Sample_CH_DD_006    6 emp.20 +dox 72h CH_DD_006 empty empty_cl20 (+)dox 72h  
7 Sample_CH_DD_007    7 wt_15 -dox 24h CH_DD_007 NRas wt_cl15 (-)dox 24h  
8 Sample_CH_DD_008    8 wt_15 +dox 24h CH_DD_008 NRas wt_cl15 (+)dox 24h  
9 Sample_CH_DD_009    9 wt_15 -dox 48h CH_DD_009 NRas wt_cl15 (-)dox 48h  
10 Sample_CH_DD_010   10 wt_15 +dox 48h CH_DD_010 NRas wt_cl15 (+)dox 48h  
..     ...  ...      ...  ...      ...  ...  ...  
> |
```



```
read_excel %>% set_colnames %>% slice %>% fill %>% select
```

The screenshot shows the RStudio interface with the following details:

- Toolbar:** Includes icons for New Project, Open Project, Save, Print, Go to file/function, and a Git icon.
- Console Tab:** Labeled "Console ~/Projects/R/Current/dplyr_tut/".
- Code Input:** Shows the R code used to process an Excel file:

```
>  
> read_excel("annotation_rnaseq_collapsed.xlsx") %>%  
+   set_colnames(paste0("X",1:7)) %>%  
+   slice(1:36) %>%  
+   fill(X5:X7) %>%  
+   select(X1, X5)  
Source: local data frame [36 x 2]
```

- Data Output:** Displays the resulting data frame:

	X1	X5
	(chr)	(chr)
1	Sample_CH_DD_001	empty empty_cl20
2	Sample_CH_DD_002	empty empty_cl20
3	Sample_CH_DD_003	empty empty_cl20
4	Sample_CH_DD_004	empty empty_cl20
5	Sample_CH_DD_005	empty empty_cl20
6	Sample_CH_DD_006	empty empty_cl20
7	Sample_CH_DD_007	NRas wt_cl15
8	Sample_CH_DD_008	NRas wt_cl15
9	Sample_CH_DD_009	NRas wt_cl15
10	Sample_CH_DD_010	NRas wt_cl15
..



```
read_excel %>% set_colnames %>% slice %>% fill %>% select %>% distinct
```

The screenshot shows the RStudio interface with the following details:

- Toolbar:** Includes icons for New Project, Open Project, Save, Print, Go to file/function, and a search bar.
- Console Tab:** Labeled "Console ~/Projects/R/Current/dplyr_tut/".
- Code Input:** Shows a sequence of R code using the pipe operator (%>%). The code reads an Excel file, sets column names, slices the first 36 rows, fills missing values in columns X5:X7, selects columns X1 through X5, and finds distinct values in column X5.
- Output:** The output indicates the data is a local data frame with 5 rows and 2 columns. It then displays the data frame structure and its contents:

	X1	X5
	(chr)	(chr)
1	Sample_CH_DD_001	empty empty_cl20
2	Sample_CH_DD_007	NRas wt_cl15
3	Sample_CH_DD_013	NRas V12D_cl6
4	Sample_CH_DD_019	NRas Q61K_cl2
5	Sample_CH_DD_025	KRas G12V_subcl2



```
read_excel %>% set_colnames %>% slice %>% fill %>% select %>% distinct %>%  
separate
```

Caution!

readr, tidy & dplyr do “clever” stuff.
(heuristics like predicting a column class by
looking at the first 1000 entries)

```
>  
>  
>  
>  
>  
> read_excel("annotation_rnaseq_collapsed.xlsx") %>%  
+   set_colnames(paste0("X",1:7)) %>%  
+   slice(1:36) %>%  
+   fill(X5:X7) %>%  
+   select(X1, X5) %>%  
+   distinct(X5) %>%  
+   separate(X5, into = c("Gene", "Mutation", "Clone"))  
Source: local data frame [5 x 4]
```

	X1	Gene	Mutation	Clone
	(chr)	(chr)	(chr)	(chr)
1	Sample_CH_DD_001	empty	empty	c120
2	Sample_CH_DD_007	NRas	wt	c115
3	Sample_CH_DD_013	NRas	V12D	c16
4	Sample_CH_DD_019	NRas	Q61K	c12
5	Sample_CH_DD_025	KRas	G12V	subc12

```
read_excel %>% set_colnames %>% slice %>% fill %>% select %>% distinct  
separate
```

The screenshot shows the RStudio interface with the following details:

- Toolbar:** Includes icons for file operations (New, Open, Save, Print), Go to file/function, and a search bar.
- Session Tab:** Labeled "dplyr_tut".
- Console Area:** Displays the R code and its execution results.
 - Code input:

```
>  
>  
>  
>  
>  
> read_excel("annotation_rnaseq_collapsed.xlsx") %>%  
+   set_colnames(paste0("X",1:7)) %>%  
+   slice(1:36) %>%  
+   fill(X5:X7) %>%  
+   select(X1, X5) %>%  
+   distinct(X5) %>%  
+   separate(X5, into = c("Gene", "Mutation", "Clone"), sep = " |_")
```
 - Output:

```
Source: local data frame [5 x 4]
```

	X1	Gene	Mutation	Clone
	(chr)	(chr)	(chr)	(chr)
1	Sample_CH_DD_001	empty	empty	c120
2	Sample_CH_DD_007	NRas	wt	c115
3	Sample_CH_DD_013	NRas	V12D	c16

```
4 Sample_CH_DD_019 NRas Q61K c12  
5 Sample_CH_DD_025 KRas G12V subc12
```



```
read_excel %>% set_colnames %>% slice %>% fill %>% select %>% distinct  
separate %>% unite
```

The screenshot shows the RStudio interface with the following details:

- Toolbar:** Includes icons for file operations (New, Open, Save, Print), Go to file/function, and a search bar.
- Session Tab:** Shows the current session name: `dplyr_tut`.
- Console Area:** Displays the R code and its execution results.
 - Code input:

```
>  
>  
>  
>  
> read_excel("annotation_rnaseq_collapsed.xlsx") %>%  
+   set_colnames(paste0("X",1:7)) %>%  
+   slice(1:36) %>%  
+   fill(X5:X7) %>%  
+   select(X1, X5) %>%  
+   distinct(X5) %>%  
+   separate(X5, into = c("Gene", "Mutation", "Clone"), sep = " |_") %>%  
+   unite(G_M_C, Gene, Mutation, Clone)
```
 - Output:

```
Source: local data frame [5 x 2]
```

	X1	G_M_C
	(chr)	(chr)
1	Sample_CH_DD_001	empty_empty_cl20
2	Sample_CH_DD_007	NRas_wt_cl15
3	Sample_CH_DD_013	NRas_V12D_cl6

```
4 Sample_CH_DD_019      NRas_Q61K_cl2  
5 Sample_CH_DD_025 KRas_G12V_subcl2
```



```
read_excel %>% set_colnames %>% slice %>% fill %>% select %>% distinct  
separate %>% unite
```

The screenshot shows the RStudio interface with the following details:

- Toolbar:** Includes icons for file operations (New, Open, Save, Print), Go to file/function, and a search bar.
- Console Tab:** Labeled "dplyr_tut".
- Console Output:** Displays the R code used to process the data, followed by the source information: "Source: local data frame [5 x 2]".
- Data Frame Preview:** Shows the resulting data frame with two columns: "X1" (chr) and "G_M_C" (chr). The data consists of five rows:

	X1	G_M_C
	(chr)	(chr)
1	Sample_CH_DD_001	empty (-_M) empty (-_M) cl20
2	Sample_CH_DD_007	NRas (-_M) wt (-_M) cl15
3	Sample_CH_DD_013	NRas (-_M) V12D (-_M) cl6
4	Sample_CH_DD_019	NRas (-_M) Q61K (-_M) cl2
5	Sample_CH_DD_025	KRas (-_M) G12V (-_M) subcl2



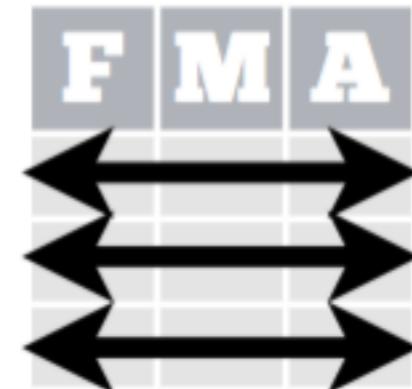
Tidy data definition

In a tidy
data set:



Each **variable** is saved
in its own **column**

&



Each **observation** is
saved in its own **row**

Wickham, H. (2014). Tidy Data. Journal of Statistical Software

read_tsv

```
read_tsv %>% gather(key, value, -variable)
```

Console ~/Projects/R/Current/dplyr_tut/ ↵

```
>
>
>
>
>
> read_tsv("mutation_data.tsv",
+           col_types = paste(rep("c", 22), collapse = "")) %>%
+           gather(Gene, Status, -STUDY_ID, -CASE_ID)
Source: local data frame [288,320 x 4]

  STUDY_ID      CASE_ID     Gene Status
  (chr)        (chr)   (fctr)  (chr)
1 acc_tcga TCGA.OR.A5J1.01 AKAP9    NA
2 acc_tcga TCGA.OR.A5J2.01 AKAP9    NA
3 acc_tcga TCGA.OR.A5J3.01 AKAP9    NA
4 acc_tcga TCGA.OR.A5J4.01 AKAP9    NA
5 acc_tcga TCGA.OR.A5J5.01 AKAP9    NA
6 acc_tcga TCGA.OR.A5J6.01 AKAP9    NA
7 acc_tcga TCGA.OR.A5J7.01 AKAP9    NA
8 acc_tcga TCGA.OR.A5J8.01 AKAP9    NA
9 acc_tcga TCGA.OR.A5J9.01 AKAP9    NA
10 acc_tcga TCGA.OR.A5JA.01 AKAP9   NA
..   ...     ...   ...
> |
```

```
read_tsv %>% gather %>% spread(key, value)
```

The screenshot shows the RStudio interface with the following details:

- Toolbar:** Includes icons for file operations (New, Open, Save, Print), Go to file/function, and a search bar.
- Project Bar:** Shows the current project is "dplyr_tut".
- Console:** Displays the R code and its output. The code reads a TSV file, gathers data by Gene and Status, and then spreads the data by Gene. The output shows the resulting data frame structure with columns for STUDY_ID, CASE_ID, and various genes (AKAP9, APC, ARID1A, ATM, ATRX, BRAF, BRCA2, CDKN2A, EGFR, KMT2C, KMT2D). It also highlights specific mutations: R894Q, P754A, and P2354Lfs*30.
- Data Preview:** Below the console, a preview of the first 10 rows of the data frame is shown. Row 10 contains mutations E1365* and E1828*.
- Message:** A message at the bottom states: "Variables not shown: KRAS (chr), NF1 (chr), PBRM1 (chr), PIK3CA (chr), PTEN (chr), SETD2 (chr), SPEN (chr), TP53 (chr), VHL (chr)".
- Bottom:** Shows the R prompt (>).

```
read_tsv %>% gather
```

```
>
>
>
> read_tsv("mutation_data.tsv",
+           col_types = paste(rep("c", 22), collapse = "")) %>%
+   gather(Gene, Status, -STUDY_ID, -CASE_ID) %>%
+   spread(Gene, Status) %>%
+   gather(Gene, Status, -STUDY_ID, -CASE_ID, -BRAF)
Source: local data frame [273,904 x 5]

  STUDY_ID      CASE_ID BRAF    Gene Status
  (chr)        (chr) (chr) (fctr) (chr)
1 acc_tcga TCGA.OR.A5J1.01    NA AKAP9     NA
2 acc_tcga TCGA.OR.A5J2.01    NA AKAP9     NA
3 acc_tcga TCGA.OR.A5J3.01    NA AKAP9     NA
4 acc_tcga TCGA.OR.A5J4.01    NA AKAP9     NA
5 acc_tcga TCGA.OR.A5J5.01    NA AKAP9     NA
6 acc_tcga TCGA.OR.A5J6.01    NA AKAP9     NA
7 acc_tcga TCGA.OR.A5J7.01    NA AKAP9     NA
8 acc_tcga TCGA.OR.A5J8.01    NA AKAP9     NA
9 acc_tcga TCGA.OR.A5J9.01    NA AKAP9     NA
10 acc_tcga TCGA.OR.A5JA.01   NA AKAP9     NA
...
...
```

```
read_tsv %>% gather %>% filter
```

The screenshot shows the RStudio interface with the following details:

- Toolbar:** Includes icons for file operations (New, Open, Save, Print), Go to file/function, and a search bar.
- Console Tab:** Labeled "dplyr_tut".
- Console Area:** Displays the R code and its output.
 - Code:

```
>  
>  
> read_tsv("mutation_data.tsv",  
+           col_types = paste(rep("c", 22), collapse = "")) %>%  
+           gather(Gene, Status, -STUDY_ID, -CASE_ID) %>%  
+           spread(Gene, Status) %>%  
+           gather(Gene, Status, -STUDY_ID, -CASE_ID, -BRAF) %>%  
+           filter(!is.na(BRAF))
```
 - Output:

```
Source: local data frame [15,998 x 5]
```
 - Data Preview:

	STUDY_ID (chr)	CASE_ID (chr)	BRAF (chr)	Gene (fctr)	Status (chr)
1	acc_tcga	TCGA.OR.A5JB.01	D594G	AKAP9	NA
2	blca_tcga	TCGA.GD.A30P.01	S364L	AKAP9	NA
3	blca_tcga_pub	TCGA.GD.A30P.01	S364L	AKAP9	NA
4	brca_tcga	TCGA.AC.A2FF.01	D449N	AKAP9	D3873N,R1947C
5	brca_tcga	TCGA.AN.A0FZ.01	K698R	AKAP9	NA
6	brca_tcga	TCGA.C8.A1HJ.01	E309*	AKAP9	NA
7	brca_tcga	TCGA.E2.A14U.01	E433K	AKAP9	NA
8	brca_tcga	TCGA.E2.A15K.06	L537S	AKAP9	NA
9	brca_tcga_pub	TCGA.AN.A0FN.01	K698R	AKAP9	NA
10	brca_tcga_pub	TCGA.C8.A12P.01	E309*	AKAP9	NA
..

```
read_tsv %>% gather %>% filter %>% group_by
```

The screenshot shows the RStudio interface with the following details:

- Toolbar:** Includes icons for file operations (New, Open, Save, Print), Go to file/function, and a search bar.
- Console Tab:** Labeled "dplyr_tut".
- Console Output:**

```
> read_tsv("mutation_data.tsv",
+           col_types = paste(rep("c", 22), collapse = "")) %>%
+           gather(Gene, Status, -STUDY_ID, -CASE_ID) %>%
+           spread(Gene, Status) %>%
+           gather(Gene, Status, -STUDY_ID, -CASE_ID, -BRAF) %>%
+           filter(!is.na(BRAF)) %>%
+           group_by(BRAF)
```

Source: local data frame [15,998 x 5]
Groups: BRAF [93]
- Data View:** Displays the first 10 rows of the resulting data frame.

	STUDY_ID (chr)	CASE_ID (chr)	BRAF (chr)	Gene (fctr)	Status (chr)
1	acc_tcga	TCGA.OR.A5JB.01	D594G	AKAP9	NA
2	blca_tcga	TCGA.GD.A30P.01	S364L	AKAP9	NA
3	blca_tcga_pub	TCGA.GD.A30P.01	S364L	AKAP9	NA
4	brca_tcga	TCGA.AC.A2FF.01	D449N	AKAP9	D3873N,R1947C
5	brca_tcga	TCGA.AN.A0FZ.01	K698R	AKAP9	NA
6	brca_tcga	TCGA.C8.A1HJ.01	E309*	AKAP9	NA
7	brca_tcga	TCGA.E2.A14U.01	E433K	AKAP9	NA
8	brca_tcga	TCGA.E2.A15K.06	L537S	AKAP9	NA
9	brca_tcga_pub	TCGA.AN.A0FN.01	K698R	AKAP9	NA
10	brca_tcga_pub	TCGA.C8.A12P.01	E309*	AKAP9	NA
..

```
read_tsv %>% gather %>% filter %>% group_by %>% summarise %>% arrange
```

The screenshot shows the RStudio interface with the following details:

- Toolbar:** Includes icons for file operations (New, Open, Save, Print), Go to file/function, and a search bar.
- Session Tab:** Shows "dplyr_tut".
- Console Area:**
 - Text input: `+ col_types = paste(rep("c", 22), collapse = "")) %>%`
 - Text output:

```
+   gather(Gene, Status, -STUDY_ID, -CASE_ID) %>%
+   spread(Gene, Status) %>%
+   gather(Gene, Status, -STUDY_ID, -CASE_ID, -BRAF) %>%
+   filter(!is.na(BRAF)) %>%
+   group_by(BRAF) %>%
+   summarise(n = n()) %>%
+   arrange(-n)
```
 - Text output: `Source: local data frame [93 x 2]`
- Data View:** Displays a data frame with columns "BRAF" and "n".

	BRAF	n
	(chr)	(int)
1	V600E	12160
2	V600K	380
3	Fusion	171
4	K601E	171
5	N581S	152
6	P403Lfs*8	133
7	G469V	114
8	G466V	95
9	A762V	76
10	G469A	76
..

```
read_tsv %>% gather %>% filter %>% group_by %>% summarise %>% arrange
```

The screenshot shows the RStudio interface with the following details:

- Toolbar:** Includes icons for file operations (New, Open, Save, Print), Go to file/function, and a search bar.
- Session Tab:** Labeled "dplyr_tut".
- Console Area:**
 - Text output:

```
+   gather(Gene, Status, -STUDY_ID, -CASE_ID) %>%
+   spread(Gene, Status) %>%
+   gather(Gene, Status, -STUDY_ID, -CASE_ID, -BRAF) %>%
+   filter(!is.na(BRAF)) %>%
+   group_by(BRAF, STUDY_ID) %>%
+   summarise(n = n()) %>%
+   ungroup %>%
+   arrange(-n)
```

Source: local data frame [163 x 3]

	BRAF	STUDY_ID	n
	(chr)	(chr)	(int)
1	V600E	thca_tcga	4503
2	V600E	thca_tcga_pub	4465
3	V600E	skcm_tcga	2451
4	V600K	skcm_tcga	380
5	V600E	coadread_tcga_pub	361
6	Fusion	thca_tcga_pub	171
7	V600E	gbm_tcga	95
8	V600E	gbm_tcga_pub2013	95
9	V600E	luad_tcga	76
10	V600E	luad_tcga_pub	76
..
 - Environment Tab:** Shows the current environment variables.
 - Help Tab:** Shows available help topics.

```
read_tsv %>% gather %>% filter %>% group_by %>% summarise %>% arrange
```

The screenshot shows the RStudio interface with the following details:

- Toolbar:** Includes icons for file operations (New, Open, Save, Print), Go to file/function, and a search bar.
- Console Tab:** Labeled "dplyr_tut".
- Console Output:** Displays the R code used to generate the data frame, followed by the source information: "Source: local data frame [239 x 3]".
- Data Frame Preview:** A table showing the first 10 rows of the data frame. The columns are "Gene", "STUDY_ID", and "n". The "Gene" column is a factor with levels KMT2D, KMT2C, TP53, CDKN2A, SPEN, PTEN, AKAP9, APC, ARID1A, and ARID1A. The "STUDY_ID" column has two distinct values: "skcm_tcga" and "stad_tcga". The "n" column contains integer values ranging from 13 to 28.
- Bottom Line:** Shows the R prompt ">".

	Gene	STUDY_ID	n
	(fctr)	(chr)	(int)
1	KMT2D	skcm_tcga	28
2	KMT2C	skcm_tcga	27
3	TP53	skcm_tcga	27
4	CDKN2A	skcm_tcga	25
5	SPEN	skcm_tcga	23
6	PTEN	skcm_tcga	21
7	AKAP9	skcm_tcga	16
8	APC	skcm_tcga	14
9	ARID1A	stad_tcga	13
10	ARID1A	stad_tcga_pub	13
..

Data Transformation with dplyr :: CHEAT SHEET



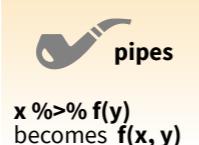
dplyr functions work with pipes and expect **tidy data**. In tidy data:



Each **variable** is in its own **column**



Each **observation**, or **case**, is in its own **row**



`x %>% f(y)` becomes `f(x, y)`

Summarise Cases

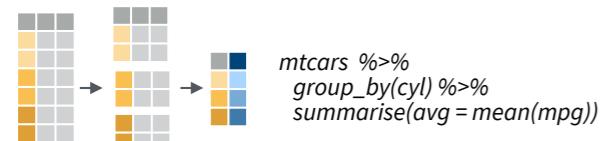
These apply **summary functions** to columns to create a new table of summary statistics. Summary functions take vectors as input and return one value (see back).

summary function

- `summarise(.data, ...)`
Compute table of summaries.
`summarise(mtcars, avg = mean(mpg))`
- `count(x, ..., wt = NULL, sort = FALSE)`
Count number of rows in each group defined by the variables in ... Also `tally()`.
`count(mtcars, cyl)`

Group Cases

Use `group_by(.data, ..., .add = FALSE)` to create a "grouped" copy of a table grouped by columns in ... dplyr functions will manipulate each "group" separately and combine the results.



`mtcars %>%
group_by(cyl) %>%
summarise(avg = mean(mpg))`

Use `rowwise(.data, ...)` to group data into individual rows. dplyr functions will compute results for each row. Also used to apply functions to list-columns without purrr functions.



`ungroup(x, ...)` Returns ungrouped copy of table.
`ungroup(g_mtcars)`



RStudio® is a trademark of RStudio, Inc

Manipulate Cases

EXTRACT CASES

Row functions return a subset of rows as a new table.



`filter(.data, ...)` Extract rows that meet logical criteria.
`filter(mtcars, mpg > 20)`

Data Import :: CHEAT SHEET

R's **tidyverse** is built around **tidy data** stored in **tibbles**, which are enhanced data frames.



The front side of this sheet shows how to read text files into R with `readr`.



The reverse side shows how to create tibbles with `tibble` and to layout tidy data with `tidyR`.

OTHER TYPES OF DATA

Try one of the following packages to import other types of files

- `haven` - SPSS, Stata, and SAS files
- `readxl` - excel files (.xls and .xlsx)
- `DBI` - databases
- `jsonlite` - json
- `xml2` - XML
- `httr` - Web APIs
- `rvest` - HTML (Web Scraping)

Save Data

Save `x`, an R object, to `path`, a file path, as:

Comma delimited file

`write_csv(x, path, na = "NA", append = FALSE, col_names = !append)`

File with arbitrary delimiter

`write_delim(x, path, delim = " ", na = "NA", append = FALSE, col_names = !append)`

CSV for excel

`write_excel_csv(x, path, na = "NA", append = FALSE, col_names = !append)`

String to file

`write_file(x, path, append = FALSE)`

String vector to file, one element per line

`write_lines(x, path, na = "NA", append = FALSE)`

Object to RDS file

`write_rds(x, path, compress = c("none", "gz", "bz2", "xz", ...))`

Tab delimited files

`write_tsv(x, path, na = "NA", append = FALSE, col_names = !append)`



Manipulate Variables

EXTRACT VARIABLES

Column functions return a set of columns as a new vector or table.



`pull(.data, var = -1)` Extract column values as a vector. Choose by name or index.
`pull(mtcars, wt)`

Data types

`readr` functions guess the types of each column and convert types when appropriate (but will NOT convert strings to factors automatically).

A message shows the type of each column in the result.

Parsed with column specification:
cols(
age = col_integer(),
sex = col_character(),
earn = col_double()
)

age is an integer
sex is a character
earn is a double (numeric)

1. Use `problems()` to diagnose problems.
`x <- read_csv("file.csv"); problems(x)`

2. Use a `col_` function to guide parsing.

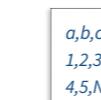
- `col_guess()` - the default
- `col_character()`
- `col_double()`, `col_euro_double()`
- `col_datetime(format = "")` Also `col_date(format = "")`, `col_time(format = "")`
- `col_factor(levels, ordered = FALSE)`
- `col_integer()`
- `col_logical()`
- `col_number()`, `col_numeric()`
- `col_skip()`

`x <- read_csv("file.csv", col_types = cols(A = col_double(), B = col_logical(), C = col_factor()))`

3. Else, read in as character vectors then parse with a `parse_` function.

- `parse_guess()`
- `parse_character()`
- `parse_datetime()` Also `parse_date()` and `parse_time()`
- `parse_double()`
- `parse_factor()`
- `parse_integer()`
- `parse_logical()`
- `parse_number()`
- `x$A <- parse_number(x$A)`

USEFUL ARGUMENTS



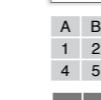
Example file

`write_file("a,b,c\n1,2,3\n4,5,NA","file.csv")
f <- "file.csv"`



Skip lines

`read_csv(f, skip = 1)`



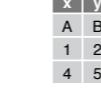
No header

`read_csv(f, col_names = FALSE)`



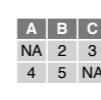
Read in a subset

`read_csv(f, n_max = 1)`



Provide header

`read_csv(f, col_names = c("x", "y", "z"))`



Missing Values

`read_csv(f, na = c("1", "."))`

Read Non-Tabular Data

Read a file into a single string

`read_file(locale = default_locale())`

Read each line into its own string

`read_lines(file, skip = 0, n_max = -1L, na = character(), locale = default_locale(), progress = interactive())`

Read Apache style log files

`read_log(file, col_names = FALSE, col_types = NULL, skip = 0, n_max = -1, progress = interactive())`

Read a file into a raw vector

`read_file_raw(file)`

Read each line into a raw vector

`read_lines_raw(file, skip = 0, n_max = -1L, progress = interactive())`



Ceci n'est pas une pipe.