

# 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**



**pipes**

$x \%>\% f(y)$  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))`



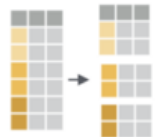
**count(x, ..., wt = NULL, sort = FALSE)**  
Count number of rows in each group defined by the variables in ... Also **tally()**.  
`count(iris, Species)`

### VARIATIONS

**summarise\_all()** - Apply funs to every column.  
**summarise\_at()** - Apply funs to specific columns.  
**summarise\_if()** - Apply funs to all cols of one type.

## Group Cases

Use **group\_by()** to create a "grouped" copy of a table. dplyr functions will manipulate each "group" separately and then combine the results.



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

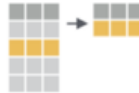
**group\_by(.data, ..., add = FALSE)**  
Returns copy of table grouped by ...  
`g_iris <- group_by(iris, Species)`

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

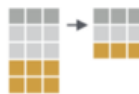
## Manipulate Cases

### EXTRACT CASES

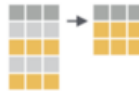
Row functions return a subset of rows as a new table.



**filter(.data, ...)** Extract rows that meet logical criteria. `filter(iris, Sepal.Length > 7)`



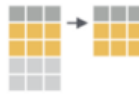
**distinct(.data, ..., .keep\_all = FALSE)** Remove rows with duplicate values.  
`distinct(iris, Species)`



**sample\_frac(tbl, size = 1, replace = FALSE, weight = NULL, .env = parent.frame())** Randomly select fraction of rows.  
`sample_frac(iris, 0.5, replace = TRUE)`



**sample\_n(tbl, size, replace = FALSE, weight = NULL, .env = parent.frame())** Randomly select size rows. `sample_n(iris, 10, replace = TRUE)`



**slice(.data, ...)** Select rows by position.  
`slice(iris, 10:15)`

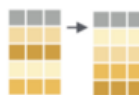
**top\_n(x, n, wt)** Select and order top n entries (by group if grouped data). `top_n(iris, 5, Sepal.Width)`

### Logical and boolean operators to use with filter()

<	<=	is.na()	%in%		xor()
>	>=	!is.na()	!	&	

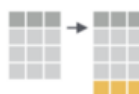
See `?base::logic` and `?Comparison` for help.

### ARRANGE CASES



**arrange(.data, ...)** Order rows by values of a column or columns (low to high), use with **desc()** to order from high to low.  
`arrange(mtcars, mpg)`  
`arrange(mtcars, desc(mpg))`

### ADD CASES



**add\_row(.data, ..., .before = NULL, .after = NULL)**  
Add one or more rows to a table.  
`add_row(faithful, eruptions = 1, waiting = 1)`

## 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(iris, Sepal.Length)`



**select(.data, ...)**  
Extract columns as a table. Also **select\_if()**.  
`select(iris, Sepal.Length, Species)`

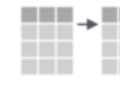
**Use these helpers with select (),**  
e.g. `select(iris, starts_with("Sepal"))`

<b>contains(match)</b>	<b>num_range(prefix, range)</b>	: e.g. mpg:cyl
<b>ends_with(match)</b>	<b>one_of(...)</b>	-, e.g. -Species
<b>matches(match)</b>	<b>starts_with(match)</b>	

### MAKE NEW VARIABLES

These apply **vectorized functions** to columns. Vectorized funs take vectors as input and return vectors of the same length as output (see back).

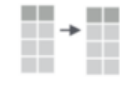
**vectorized function**



**mutate(.data, ...)**  
Compute new column(s).  
`mutate(mtcars, gpm = 1/mpg)`



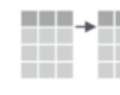
**transmute(.data, ...)**  
Compute new column(s), drop others.  
`transmute(mtcars, gpm = 1/mpg)`



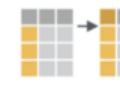
**mutate\_all(.tbl, .funs, ...)** Apply funs to every column. Use with **funs()**. Also **mutate\_if()**.  
`mutate_all(faithful, funs(log(.), log2(.)))`  
`mutate_if(iris, is.numeric, funs(log(.)))`



**mutate\_at(.tbl, .cols, .funs, ...)** Apply funs to specific columns. Use with **funs()**, **vars()** and the helper functions for **select()**.  
`mutate_at(iris, vars(-Species), funs(log(.)))`



**add\_column(.data, ..., .before = NULL, .after = NULL)** Add new column(s). Also **add\_count()**, **add\_tally()**. `add_column(mtcars, new = 1:32)`



**rename(.data, ...)** Rename columns.  
`rename(iris, Length = Sepal.Length)`