# 2: Data Wrangling with `dplyr`

Suppose we have the following toy data set, named `df`. The first two columns are numeric while the third column is categorical.

| x1 | x2 | cat1 |
|----|----|------|
| 1  | 3  | Yes  |
| 7  | NA | Yes  |
| 4  | 2  | No   |

$y_1$    New var

0.333    small x1

NA    large x1

2    large x1

**Create New Variables with `mutate()` (Perhaps with `case_when()` or `if_else()`)** — name of new variable

$$df \; \%>\% \; mutate(y1 = x1 / x2)$$ — operations on old variables

$$df \; \%>\% \; mutate(newvar = if\_else(x_1 > 3,$$
$$true = "large \; x1",$$
$$false = "small \; x1"))$$

**Choose Rows to Keep with `filter()`** — condition

$$df \; \%>\% \; filter(x1 > 6)$$

$$\Rightarrow \quad \frac{x1}{7} \quad \frac{x2}{NA} \quad \frac{cat1}{Yes}$$

$$df \; \%>\% \; filter(cat1 == "No")$$

$$df \; \%>\% \; filter(x_1 == max(x1))$$

$$\Rightarrow \quad \frac{x1}{7} \quad \frac{x2}{NA} \quad \frac{cat1}{Yes}$$

**Choose Columns to Keep with `select()`**

$$df \; \%>\% \; select(x2, cat1)$$

$$\Rightarrow \quad 
\begin{array}{cc}
x2 & cat1 \\
\hline
3 & Yes \\
NA & Yes \\
2 & No \\
\end{array}$$

| x1 | x2 | cat1 |
|----|----|------|
| 1  | 3  | Yes  |
| 7  | NA | Yes  |
| 4  | 2  | No   |

**Order/Sort Your Data Set with arrange()**

$$df \ \%>\% \ arrange(x1)$$

$$=7 \qquad \begin{array}{ccc} x1 & x2 & cat1 \\ 1 & 3 & Yes \\ 4 & 2 & No \\ 7 & NA & Yes \end{array}$$

**Obtain Numerical Summaries with summarise()**

$$df \ \%>\% \ summarise(x1m = median(x1))$$

$$=7 \qquad \frac{x1m}{4}$$

**Obtain Numerical Summaries by Group with group_by() and summarise()**

$$df \ \%>\% \ group\_by(cat1) \ \%>\%$$

$$summarise(Numns = min(x1))$$

$$=> \qquad \begin{array}{cc} cat1 & Numns \\ Yes & 1 \\ No & 4 \end{array}$$

mean(), median(),
max(), min(),
sd(), var(),
quantile(), etc.