

# Data Cleaning

## Recap on summarization

- `summary(x)`: quantile information
- `count(x)`: what unique values do you have?
  - `distinct()`: what are the distinct values?
  - `n_distinct()` with `pull()`: how many distinct values?
- `group_by()`: changes all subsequent functions
  - combine with `summarize()` to get statistics per group
  - combine with `mutate()` to add column
- `summarize()` with `n()` gives the count (NAs included)

▮ [Cheatsheet](#)

## Recap on data classes

- There are two types of number class objects: integer and double
- Logic class objects only have **TRUE** or **FALSE** (without quotes)
- `class()` can be used to test the class of an object `x`
- `as.CLASS_NAME(x)` can be used to change the class of an object `x`
- Factors are a special character class that has levels - more on that soon!
- tibbles show column classes!
- two dimensional object classes include: data frames, tibbles, matrices, and lists
- Dates can be handled with the `lubridate` package
- Make sure you choose the right function for the way the date is formatted!

▮ [Cheatsheet](#)

## Data Cleaning

In general, data cleaning is a process of investigating your data for inaccuracies, or recoding it in a way that makes it more manageable.

□ MOST IMPORTANT RULE - LOOK □ AT YOUR DATA! □

# Dealing with Missing Data

# Missing data types

One of the most important aspects of data cleaning is missing values.

Types of “missing” data:

- **NA** - general missing data
- **NaN** - stands for “**N**ot **a** **N**umber”, happens when you do  $0/0$ .
- **Inf** and **-Inf** - Infinity, happens when you divide a positive number (or negative number) by 0.

## Finding Missing data

- `is.na` - looks for NAN and NA
- `is.nan` - looks for NAN
- `is.infinite` - looks for Inf or -Inf

```
test <- c(0, NA, -1)  
test/0
```

```
[1] NaN  NA -Inf
```

```
test <- test/0  
is.na(test)
```

```
[1] TRUE  TRUE FALSE
```

```
is.nan(test)
```

```
[1] TRUE FALSE FALSE
```

```
is.infinite(test)
```

```
[1] FALSE FALSE  TRUE
```

## Useful checking functions

`any()` can help you check if there are any NA values in a vector

```
test
```

```
[1] NaN  NA -Inf
```

```
any(is.na(test))
```

```
[1] TRUE
```



## Finding NA values with `count()`

Check the values for your variables, are they what you expect?

`count()` is a great option because it helps you check if rare values make sense.

```
#library(dasehr)
yearly_co2 <- read_csv(file =
  "https://daseh.org/data/Yearly_CO2_Emissions_1000_tonnes.csv")
yearly_co2 %>% count(`1989`)
```

```
# A tibble: 164 × 2
```

```
  `1989`      n
  <dbl> <int>
1     22         1
2    47.7         1
3    51.3         1
4    58.7         1
5    62.3         2
6    66          1
7    69.7         1
8    77          1
9    80.7         1
10   103          2
# ▯ 154 more rows
```

# naniar

Sometimes you need to look at lots of data though... the [naniar package](#) is a good option.

```
#install.packages("naniar")  
library(naniar)
```



“Artwork by @allison\_horst”. <https://allisonhorst.com/>

## naniar: pct\_complete()

This can tell you if there are missing values in the dataset.

```
pct_complete(yearly_co2)
```

```
[1] 33.62421
```

Or for a particular variable:

```
yearly_co2 %>% select(`1989`) %>%  
  pct_complete()
```

```
[1] 89.58333
```

```
yearly_co2 %>% select(`1889`) %>%  
  pct_complete()
```

```
[1] 22.91667
```

## naniar:miss\_var\_summary()

To get the percent missing (and counts) for each variable as a table, use this function.

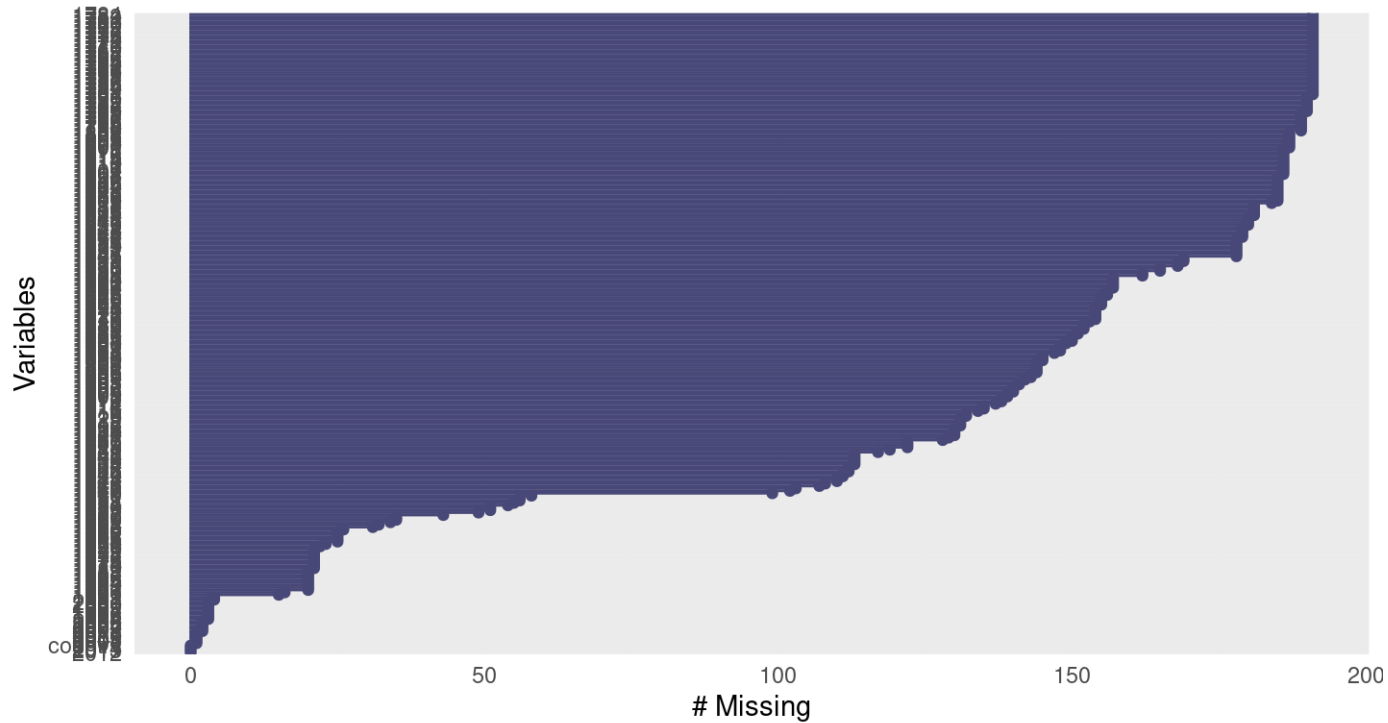
```
miss_var_summary(yearly_co2)
```

```
# A tibble: 265 × 3
  variable n_miss pct_miss
  <chr>    <int>    <num>
1 1751      191     99.5
2 1752      191     99.5
3 1753      191     99.5
4 1754      191     99.5
5 1755      191     99.5
6 1756      191     99.5
7 1757      191     99.5
8 1758      191     99.5
9 1759      191     99.5
10 1760      191     99.5
#   255 more rows
```

## naniar plots

The `gg_miss_var()` function creates a nice plot about the number of missing values for each variable, (need a data frame).

```
gg_miss_var(yearly_co2)
```



# Missing Data Issues

Recall that mathematical operations with NA often result in NAs.

```
sum(c(1, 2, 3, NA))
```

```
[1] NA
```

```
mean(c(1, 2, 3, NA))
```

```
[1] NA
```

```
median(c(1, 2, 3, NA))
```

```
[1] NA
```

## Missing Data Issues

This is also true for logical data. Recall that **TRUE** is evaluated as 1 and **FALSE** is evaluated as 0.

```
x <- c(TRUE, TRUE, TRUE, TRUE, FALSE, NA)  
sum(x)
```

```
[1] NA
```

```
sum(x, na.rm = TRUE)
```

```
[1] 4
```

## **filter()** and missing data

Be **careful** with missing data using subsetting!

**filter()** removes missing values by default. Because R can't tell for sure if an NA value meets the condition. To keep them need to add `is.na()` conditional.

Think about if this is OK or not - it depends on your data!



## filter() and missing data

What if NA values represent values that are so low it is undetectable?

Filter will drop them from the data.

```
#CO_heat_ER <- read_csv(file = "https://daseh.org/data/Colorado_ER_heat_visits_AgeAdjusted_data.csv")
library(dasehr)
CO_heat_ER %>% filter(visits > 0)
```

```
# A tibble: 315 × 7
```

	county	rate	lower95cl	upper95cl	visits	year	gender
	<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<chr>
1	Statewide	5.64	4.70	6.59	140	2011	Female
2	Statewide	7.39	6.30	8.47	183	2011	Male
3	Statewide	6.51	5.80	7.23	323	2011	Both genders
4	Statewide	5.64	4.72	6.57	146	2012	Female
5	Statewide	7.56	6.48	8.65	193	2012	Male
6	Statewide	6.58	5.88	7.29	339	2012	Both genders
7	Statewide	4.94	4.06	5.82	124	2013	Female
8	Statewide	6.72	5.72	7.72	178	2013	Male
9	Statewide	5.82	5.16	6.49	302	2013	Both genders
10	Statewide	3.52	2.80	4.25	92	2014	Female

```
# [ ] 305 more rows
```

## filter() and missing data

`is.na()` can help us keep them.

```
CO_heat_ER %>% filter(visits > 0 | is.na(visits))
```

```
# A tibble: 1,147 × 7
```

	county <chr>	rate <dbl>	lower95cl <dbl>	upper95cl <dbl>	visits <dbl>	year <dbl>	gender <chr>
1	Statewide	5.64	4.70	6.59	140	2011	Female
2	Statewide	7.39	6.30	8.47	183	2011	Male
3	Statewide	6.51	5.80	7.23	323	2011	Both genders
4	Statewide	5.64	4.72	6.57	146	2012	Female
5	Statewide	7.56	6.48	8.65	193	2012	Male
6	Statewide	6.58	5.88	7.29	339	2012	Both genders
7	Statewide	4.94	4.06	5.82	124	2013	Female
8	Statewide	6.72	5.72	7.72	178	2013	Male
9	Statewide	5.82	5.16	6.49	302	2013	Both genders
10	Statewide	3.52	2.80	4.25	92	2014	Female

```
# 1,137 more rows
```

## To remove rows with NA values for a variable use `drop_na()`

A function from the `tidyr` package. (Need a data frame to start!)

Disclaimer: Don't do this unless you have thought about if dropping NA values makes sense based on knowing what these values mean in your data. **Also consider if you need those rows for values for other variables.**

```
dim(yearly_co2)
```

```
[1] 192 265
```

```
yearly_co2_drop <- yearly_co2 %>% drop_na(`1989`)  
dim(yearly_co2_drop)
```

```
[1] 172 265
```

# Let's take a look

Can still have NAs for other columns

yearly\_co2\_drop

```
# A tibble: 172 × 265
  country `1751` `1752` `1753` `1754` `1755` `1756` `1757` `1758` `1759` `1760`
  <chr>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 Afghan... NA     NA     NA     NA     NA     NA     NA     NA     NA     NA
2 Albania NA     NA     NA     NA     NA     NA     NA     NA     NA     NA
3 Algeria NA     NA     NA     NA     NA     NA     NA     NA     NA     NA
4 Angola  NA     NA     NA     NA     NA     NA     NA     NA     NA     NA
5 Antigu... NA     NA     NA     NA     NA     NA     NA     NA     NA     NA
6 Argent... NA     NA     NA     NA     NA     NA     NA     NA     NA     NA
7 Armenia NA     NA     NA     NA     NA     NA     NA     NA     NA     NA
8 Austra... NA     NA     NA     NA     NA     NA     NA     NA     NA     NA
9 Austria NA     NA     NA     NA     NA     NA     NA     NA     NA     NA
10 Azerba... NA     NA     NA     NA     NA     NA     NA     NA     NA     NA
#   162 more rows
#   254 more variables: `1761` <dbl>, `1762` <dbl>, `1763` <dbl>, `1764` <dbl>,
#   `1765` <dbl>, `1766` <dbl>, `1767` <dbl>, `1768` <dbl>, `1769` <dbl>,
#   `1770` <dbl>, `1771` <dbl>, `1772` <dbl>, `1773` <dbl>, `1774` <dbl>,
#   `1775` <dbl>, `1776` <dbl>, `1777` <dbl>, `1778` <dbl>, `1779` <dbl>,
#   `1780` <dbl>, `1781` <dbl>, `1782` <dbl>, `1783` <dbl>, `1784` <dbl>,
#   `1785` <dbl>, `1786` <dbl>, `1787` <dbl>, `1788` <dbl>, `1789` <dbl>, ...
```

## To remove rows with **NA** values for a data frame use **drop\_na()**

This function of the `tidyr` package drops rows with **any** missing data in **any** column when used on a df.

```
yearly_co2_drop <- yearly_co2 %>% drop_na()  
yearly_co2_drop
```

```
# A tibble: 1 × 265
```

```
  country `1751` `1752` `1753` `1754` `1755` `1756` `1757` `1758` `1759` `1760`  
  <chr>    <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>  
1 United ... 9360 9360 9360 9370 9370 10000 10000 10000 10000 10000  
#   254 more variables: `1761` <dbl>, `1762` <dbl>, `1763` <dbl>, `1764` <dbl>,  
#   `1765` <dbl>, `1766` <dbl>, `1767` <dbl>, `1768` <dbl>, `1769` <dbl>,  
#   `1770` <dbl>, `1771` <dbl>, `1772` <dbl>, `1773` <dbl>, `1774` <dbl>,  
#   `1775` <dbl>, `1776` <dbl>, `1777` <dbl>, `1778` <dbl>, `1779` <dbl>,  
#   `1780` <dbl>, `1781` <dbl>, `1782` <dbl>, `1783` <dbl>, `1784` <dbl>,  
#   `1785` <dbl>, `1786` <dbl>, `1787` <dbl>, `1788` <dbl>, `1789` <dbl>,  
#   `1790` <dbl>, `1791` <dbl>, `1792` <dbl>, `1793` <dbl>, `1794` <dbl>, ...
```

# Drop columns with any missing values

Use the `miss_var_which()` function from `naniar`

```
miss_var_which(yearly_co2)# which columns have missing values
```

```
[1] "1751" "1752" "1753" "1754" "1755" "1756" "1757" "1758" "1759" "1760"  
[11] "1761" "1762" "1763" "1764" "1765" "1766" "1767" "1768" "1769" "1770"  
[21] "1771" "1772" "1773" "1774" "1775" "1776" "1777" "1778" "1779" "1780"  
[31] "1781" "1782" "1783" "1784" "1785" "1786" "1787" "1788" "1789" "1790"  
[41] "1791" "1792" "1793" "1794" "1795" "1796" "1797" "1798" "1799" "1800"  
[51] "1801" "1802" "1803" "1804" "1805" "1806" "1807" "1808" "1809" "1810"  
[61] "1811" "1812" "1813" "1814" "1815" "1816" "1817" "1818" "1819" "1820"  
[71] "1821" "1822" "1823" "1824" "1825" "1826" "1827" "1828" "1829" "1830"  
[81] "1831" "1832" "1833" "1834" "1835" "1836" "1837" "1838" "1839" "1840"  
[91] "1841" "1842" "1843" "1844" "1845" "1846" "1847" "1848" "1849" "1850"  
[101] "1851" "1852" "1853" "1854" "1855" "1856" "1857" "1858" "1859" "1860"  
[111] "1861" "1862" "1863" "1864" "1865" "1866" "1867" "1868" "1869" "1870"  
[121] "1871" "1872" "1873" "1874" "1875" "1876" "1877" "1878" "1879" "1880"  
[131] "1881" "1882" "1883" "1884" "1885" "1886" "1887" "1888" "1889" "1890"  
[141] "1891" "1892" "1893" "1894" "1895" "1896" "1897" "1898" "1899" "1900"  
[151] "1901" "1902" "1903" "1904" "1905" "1906" "1907" "1908" "1909" "1910"  
[161] "1911" "1912" "1913" "1914" "1915" "1916" "1917" "1918" "1919" "1920"  
[171] "1921" "1922" "1923" "1924" "1925" "1926" "1927" "1928" "1929" "1930"  
[181] "1931" "1932" "1933" "1934" "1935" "1936" "1937" "1938" "1939" "1940"  
[191] "1941" "1942" "1943" "1944" "1945" "1946" "1947" "1948" "1949" "1950"  
[201] "1951" "1952" "1953" "1954" "1955" "1956" "1957" "1958" "1959" "1960"  
[211] "1961" "1962" "1963" "1964" "1965" "1966" "1967" "1968" "1969" "1970"  
[221] "1971" "1972" "1973" "1974" "1975" "1976" "1977" "1978" "1979" "1980"  
[231] "1981" "1982" "1983" "1984" "1985" "1986" "1987" "1988" "1989" "1990"  
[241] "1991" "1992" "1993" "1994" "1995" "1996" "1997" "1998" "1999" "2000"
```

# Drop columns with any missing values

`miss_var_which` and function from `naniar` (need a data frame)

```
yearly_co2_drop <- yearly_co2 %>% select(!miss_var_which(yearly_co2))
yearly_co2_drop
```

```
# A tibble: 192 × 4
```

```
  country      `2012` `2013` `2014`
  <chr>      <dbl> <dbl> <dbl>
1 Afghanistan  10800  10000   9810
2 Albania       4910   5060   5720
3 Algeria     130000 134000 145000
4 Andorra        488    477    462
5 Angola       33400  32600  34800
6 Antigua and Barbuda  524    524    532
7 Argentina    192000 190000 204000
8 Armenia       5690   5500   5530
9 Australia    388000 372000 361000
10 Austria      62300  62500  58700
# ▯ 182 more rows
```

## Change a value to be **NA**

Let's say we think that all 0 values should be **NA**.

```
covid_ww <- covid_wastewater
```

```
#covid_ww <- read_csv(file = "https://daseh.org/data/SARS-CoV-2_Wastewater_Data.
```

```
covid_ww %>% count(is.na(rna_pct_change_15d))
```

```
# A tibble: 2 × 2
```

```
  `is.na(rna_pct_change_15d)`      n  
  <lgl>                          <int>  
1 FALSE                          706442  
2 TRUE                             69617
```



## Change a value to be **NA**

The `na_if()` function of `dplyr` can be helpful for changing all 0 values to **NA**.

```
covid_ww <- covid_ww %>%  
  mutate(rna_pct_change_15d = na_if(rna_pct_change_15d, 0))
```

```
covid_ww %>% count(is.na(rna_pct_change_15d))
```

```
# A tibble: 2 × 2  
  `is.na(rna_pct_change_15d)`      n  
  <lgl>                          <int>  
1 FALSE                          694836  
2 TRUE                             81223
```

## Change **NA** to be a value

The `replace_na()` function (part of the `tidyr` package), can do the opposite of `na_if()`. (note that you must use numeric values as replacement - we will show how to replace with character strings soon)

```
covid_ww %>%  
  mutate(rna_pct_change_15d = replace_na(rna_pct_change_15d, 0)) %>%  
  count(is.na(rna_pct_change_15d))
```

```
# A tibble: 1 × 2  
  `is.na(rna_pct_change_15d)`      n  
  <lgl>                          <int>  
1 FALSE                          776059
```

## Think about **NA**

### THINK ABOUT YOUR DATA FIRST!

- Sometimes removing **NA** values leads to distorted math - be careful!
- Think about what your **NA** means for your data (are you sure ?).
  - Is an **NA** for values so low they could not be reported?
  - Or is it if it was too low and also if there was a different issue (like no one reported)?

## Think about **NA**

If it is something more like a zero then you might want it included in your data like a zero instead of an **NA**.

Example: - survey reports **NA** if student has never tried cigarettes - survey reports 0 if student has tried cigarettes but did not smoke that week

□ You might want to keep the **NA** values so that you know the original sample size.

## Word of caution

□ Calculating percentages will give you a different result depending on your choice to include NA values.!

This is because the denominator changes.

## Word of caution - Percentages with NA

```
count(yearly_co2, `1800`) %>% mutate(percent = (n/(sum(n)) *100))
```

```
# A tibble: 6 × 3
```

```
  `1800`      n percent
  <dbl> <int>   <dbl>
1    3.67     1  0.521
2   253     1  0.521
3   407     1  0.521
4   796     1  0.521
5 26700     1  0.521
6    NA    187  97.4
```

## Word of caution - Percentages with NA

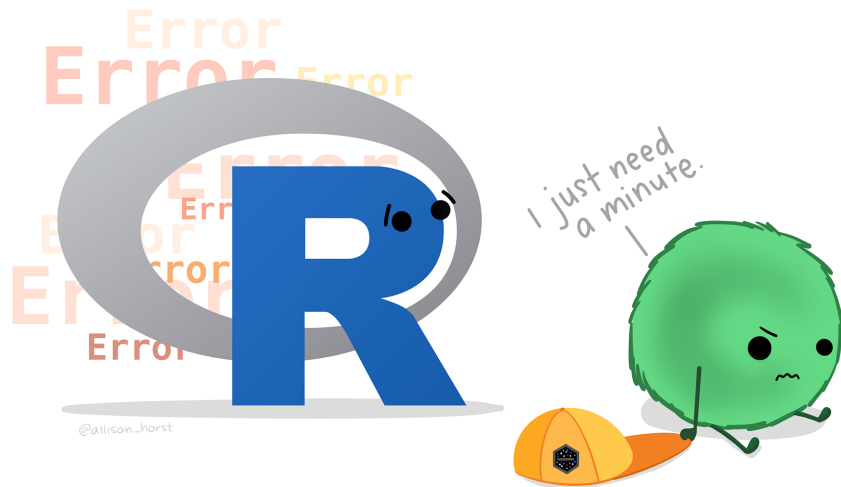
```
yearly_co2 %>% drop_na(`1800`) %>%  
  count(`1800`) %>% mutate(percent = (n/(sum(n)) *100))
```

```
# A tibble: 5 × 3  
  `1800`      n percent  
  <dbl> <int>   <dbl>  
1     3.67     1     20  
2    253     1     20  
3    407     1     20  
4    796     1     20  
5 26700     1     20
```

Should you be dividing by the total count with NA values included?  
It depends on your data and what NA might mean.  
Pay attention to your data and your NA values!

# Don't forget about the common issues

- Extra or Missing commas
- Extra or Missing parentheses
- Case sensitivity
- Spelling





## Summary

- `is.na()`, `any(is.na())`, `all(is.na())`, `count()`, and functions from `nanjar` like `gg_miss_var()` and `miss_var_summary` can help determine if we have NA values
- `miss_var_which()` can help you drop columns that have any missing values.
- `filter()` automatically removes NA values - can't confirm or deny if condition is met (need `| is.na()` to keep them)
- `drop_na()` can help you remove NA values from a variable or an entire data frame
- NA values can change your calculation results
- think about what NA values represent - don't drop them if you shouldn't
- `na_if()` will make NA values for a particular value
- `replace_na()` will replace `NA values with a particular value

# Lab Part 1

- [Class Website](#)
- [Lab](#)

# Recoding Variables

## Example of Recoding

Let's make some data about microplastics:

<https://www.medicalnewstoday.com/articles/what-do-we-know-about-microplastics-in-food#Common-microplastics-in-food>

<https://www.sciencedirect.com/science/article/pii/S0045653523028072>

<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5708064/>

```
set.seed(124)
data_diet <- tibble(Diet = rep(c("A", "B", "B"),
                             times = 4),
                   microplastic = c("Dioxin",
                                    "Dioxin",
                                    "Other",
                                    "Bisphenol A",
                                    "BPA",
                                    "dioxin",
                                    "bpa",
                                    "0",
                                    "dioxin",
                                    "bpa",
                                    "BPA",
                                    "0"),
                   blood_level_start_nM = sample(1:8, size = 12, replace = TRUE),
                   blood_level_change_nM = sample(-.5:4, size = 12, replace = TRUE))
```

Say we have some data about samples in a diet study:

```
data_diet
```

```
# A tibble: 12 × 4
```

```
  Diet microplastic blood_level_start_nM blood_level_change_nM
  <chr> <chr>          <int>          <dbl>
1 A     Dioxin           1            2.5
2 B     Dioxin           7            2.5
3 B     Other            2            0.5
4 A     Bisphenol A      3           -0.5
5 B     BPA              5            0.5
6 B     dioxin           8            0.5
7 A     bpa              6            2.5
8 B     0                5            3.5
9 B     dioxin           2            0.5
10 A    bpa              1            1.5
11 B    BPA              7            1.5
12 B    0                3            2.5
```

# Oh dear...

This needs lots of recoding.

```
data_diet %>%  
  count(microplastic)
```

```
# A tibble: 7 × 2  
  microplastic      n  
  <chr>          <int>  
1 BPA              2  
2 Bisphenol A      1  
3 Dioxin            2  
4 0                 2  
5 Other             1  
6 bpa               2  
7 dioxin            2
```

## **dp1yr can help!**

Using Excel to find all of the different ways `microplastic` has been coded, could be hectic! In `dp1yr` you can use the `case_when` function.

## Or you can use `case_when()`

The `case_when()` function of `dplyr` can help us to do this as well.

It is more flexible and powerful.

(need `mutate` here too!)



## Or you can use `case_when()`

Need quotes for conditions and new values!

```
data_diet %>%  
  mutate(microplastic_recoded = case_when(  
    microplastic == "0" ~ "Other",  
    microplastic == "Bisphenol A" ~ "BPA",  
    microplastic == "bpa" ~ "BPA",  
    microplastic == "dioxin" ~ "Dioxin")) %>%  
  count(microplastic, microplastic_recoded)
```

```
# A tibble: 7 × 3  
  microplastic microplastic_recoded     n  
  <chr>         <chr>         <int>  
1 BPA          <NA>           2  
2 Bisphenol A BPA            1  
3 Dioxin       <NA>           2  
4 0            Other          2  
5 Other        <NA>           1  
6 bpa          BPA            2  
7 dioxin       Dioxin         2
```

# What happened?

We seem to have NA values!

We didn't specify what happens to values that were already **Other** or **Dioxin** or **BPA**.

```
data_diet %>%
  mutate(microplastic_recoded = case_when(
    microplastic == "0" ~ "Other",
    microplastic == "Bisphenol A" ~ "BPA",
    microplastic == "bpa" ~ "BPA",
    microplastic == "dioxin" ~ "Dioxin")) %>%
  count(microplastic, microplastic_recoded)
```

## case\_when() drops unspecified values

Note that automatically values not reassigned explicitly by case\_when() will be NA unless otherwise specified.

*# General Format - this is not code!*

```
{data_input} %>%  
  mutate({variable_to_fix} = case_when({Variable_fixing}  
    /some condition/ ~ {value_for_con},  
    TRUE ~ {value_for_not_meeting_condition})
```

{value\_for\_not\_meeting\_condition} could be something new or it can be the original values of the column

## case\_when with TRUE ~ original variable name

```
data_diet %>%
  mutate(microplastic_recoded = case_when(
    microplastic == "0" ~ "Other",
    microplastic == "Bisphenol A" ~ "BPA",
    microplastic == "bpa" ~ "BPA",
    microplastic == "dioxin" ~ "Dioxin",
    TRUE ~ microplastic)) %>%

  count(microplastic_recoded)

# A tibble: 3 × 2
  microplastic_recoded     n
  <chr>                 <int>
1 BPA                    5
2 Dioxin                  4
3 Other                   3
```

## Typically it is good practice to include the TRUE statement

You never know if you might be missing something - and if a value already was an NA it will stay that way.

```
data_diet %>%
  mutate(microplastic_recoded = case_when(
    microplastic == "0" ~ "Other",
    microplastic == "Bisphenol A" ~ "BPA",
    microplastic == "bpa" ~ "BPA",
    microplastic == "dioxin" ~ "Dioxin",
    TRUE ~ microplastic)) %>%
  count(microplastic, microplastic_recoded)
```

## But maybe we want NA?

Perhaps we want values that are 0 or Other to actually be NA, then `case_when` can be helpful for this. We simply specify everything else.

```
data_diet %>%
  mutate(microplastic_recoded = case_when(
    microplastic == "Bisphenol A" ~ "BPA",
    microplastic == "bpa" ~ "BPA",
    microplastic == "BPA" ~ "BPA",
    microplastic == "dioxin" ~ "Dioxin",
    microplastic == "Dioxin" ~ "Dioxin")) %>%
  count(microplastic, microplastic_recoded)
```

```
# A tibble: 7 × 3
  microplastic microplastic_recoded     n
  <chr>         <chr>         <int>
1 BPA          BPA             2
2 Bisphenol A  BPA             1
3 Dioxin       Dioxin          2
4 0            <NA>            2
5 Other        <NA>            1
6 bpa          BPA             2
7 dioxin       Dioxin          2
```

## case\_when() can also overwrite/update a variable

You need to specify what we want in the first part of mutate.

```
data_diet %>%
  mutate(microplastic = case_when(
    microplastic == "Bisphenol A" ~ "BPA",
    microplastic == "bpa" ~ "BPA",
    microplastic == "BPA" ~ "BPA",
    microplastic == "dioxin" ~ "Dioxin",
    microplastic == "Dioxin" ~ "Dioxin")) %>%
  count(microplastic)
```

# A tibble: 3 × 2

	microplastic	n
	<chr>	<int>
1	BPA	5
2	Dioxin	4
3	<NA>	3

## More complicated case\_when()

case\_when can do complicated statements and can match many patterns at a time.

```
data_diet %>%
  mutate(microplastic_recoded = case_when(
    microplastic == "Dioxin" ~ "Dioxin", # keep it the same!
    microplastic == "dioxin" ~ "Dioxin",
    microplastic %in% c("BPA", "bpa", "Bisphenol A") ~ "BPA",
    microplastic %in% c("0", "Other") ~ "Other")) %>%

  count(microplastic, microplastic_recoded)
```

```
# A tibble: 7 × 3
  microplastic microplastic_recoded     n
  <chr>         <chr>         <int>
1 BPA          BPA             2
2 Bisphenol A BPA             1
3 Dioxin       Dioxin          2
4 0            Other           2
5 Other        Other           1
6 bpa          BPA             2
7 dioxin       Dioxin          2
```



## Another reason for `case_when()`

`case_when` can do very sophisticated comparisons!

Here we create a new variable called `Effect`.

```
data_diet <- data_diet %>%  
  mutate(Effect = case_when(blood_level_change_nM > 0 ~ "Increase",  
                             blood_level_change_nM == 0 ~ "Same",  
                             blood_level_change_nM < 0 ~ "Decrease"))  
  
head(data_diet)
```

```
# A tibble: 6 × 5  
  Diet      microplastic blood_level_start_nM blood_level_change_nM Effect  
  <chr>    <chr>                <int>                <dbl> <chr>  
1 A      Dioxin                 1                 2.5 Increase  
2 B      Dioxin                 7                 2.5 Increase  
3 B      Other                  2                 0.5 Increase  
4 A      Bisphenol A           3                -0.5 Decrease  
5 B      BPA                    5                 0.5 Increase  
6 B      dioxin                 8                 0.5 Increase
```

## Now it is easier to see what is happening

```
data_diet %>%  
  count(Diet, Effect)  
  
# A tibble: 3 × 3  
  Diet Effect      n  
  <chr> <chr>   <int>  
1 A     Decrease 1  
2 A     Increase 3  
3 B     Increase 8
```

# Working with strings

# Strings in R

- R can do much more than find exact matches for a whole string!



# The `stringr` package

The `stringr` package:

- Modifying or finding **part** or all of a character string
- We will not cover `grep` or `gsub` - base R functions
  - are used on forums for answers
- Almost all functions start with `str_*`

## **stringr**

`str_detect`, and `str_replace` search for matches to argument pattern within each element of a **character vector** (not data frame or tibble!).

- `str_detect` - returns TRUE if pattern is found
- `str_replace` - replaces pattern with replacement

## `str_detect()`

The `string` argument specifies what to check

The `pattern` argument specifies what to check for (case sensitive)

```
Effect <- pull(data_diet) %>% head(n = 6)
```

```
Effect
```

```
[1] "Increase" "Increase" "Increase" "Decrease" "Increase" "Increase"
```

```
str_detect(string = Effect, pattern = "d")
```

```
[1] FALSE FALSE FALSE FALSE FALSE FALSE
```

```
str_detect(string = Effect, pattern = "D")
```

```
[1] FALSE FALSE FALSE TRUE FALSE FALSE
```

## **str\_replace()**

The `string` argument specifies what to check

The `pattern` argument specifies what to check for

The `replacement` argument specifies what to replace the pattern with

```
str_replace(string = Effect, pattern = "D", replacement = "d")
```

```
[1] "Increase" "Increase" "Increase" "decrease" "Increase" "Increase"
```



**str\_replace()** only replaces the first instance of the pattern in each value

str\_replace\_all() can be used to replace all instances within each value

```
str_replace(string = Effect, pattern = "e", replacement = "E")
```

```
[1] "IncrEase" "IncrEase" "IncrEase" "DEcrease" "IncrEase" "IncrEase"
```

```
str_replace_all(string = Effect, pattern = "e", replacement = "E")
```

```
[1] "IncrEasE" "IncrEasE" "IncrEasE" "DEcrEasE" "IncrEasE" "IncrEasE"
```

## Subsetting part of a string

`str_sub()` allows you to subset part of a string

The `string` argument specifies what strings to work with

The `start` argument specifies position of where to start

The `end` argument specifies position of where to end

```
str_sub(string = Effect, start = 1, end = 3)
```

```
[1] "Inc" "Inc" "Inc" "Dec" "Inc" "Inc"
```

# filter and stringr functions

```
head(data_diet, n = 4)
```

```
# A tibble: 4 × 5
```

```
Diet microplastic blood_level_start_nM blood_level_change_nM Effect
<chr> <chr> <int> <dbl> <chr>
1 A Dioxin 1 2.5 Increase
2 B Dioxin 7 2.5 Increase
3 B Other 2 0.5 Increase
4 A Bisphenol A 3 -0.5 Decrease
```

```
data_diet %>%
```

```
  filter(str_detect(string = microplastic,
                    pattern = "B"))
```

```
# A tibble: 3 × 5
```

```
Diet microplastic blood_level_start_nM blood_level_change_nM Effect
<chr> <chr> <int> <dbl> <chr>
1 A Bisphenol A 3 -0.5 Decrease
2 B BPA 5 0.5 Increase
3 B BPA 7 1.5 Increase
```

## OK back to our original problem

```
count(data_diet, microplastic)
```

```
# A tibble: 7 × 2
  microplastic      n
  <chr>          <int>
1 BPA              2
2 Bisphenol A     1
3 Dioxin           2
4 0                2
5 Other            1
6 bpa              2
7 dioxin           2
```

## case\_when() made an improvement

But we still might miss a strange value - like a misspelling

```
data_diet %>%  
  mutate(microplastic_recoded = case_when(  
    microplastic %in% c("Dioxin", "dioxin") ~ "Dioxin",  
    microplastic %in% c("BPA", "bpa", "Bisphenol A") ~ "BPA",  
    microplastic %in% c("0", "Other") ~ "Other",  
    TRUE ~ microplastic))
```

## case\_when() improved with stringr

^ indicates the beginning of a character string \$ indicates the end

```
data_diet %>%  
  mutate(microplastic_recoded = case_when(  
    str_detect(string = microplastic, pattern = "^b|B") ~ "BPA",  
    str_detect(string = microplastic, pattern = "^o|^O") ~ "Other",  
    str_detect(string = microplastic, pattern = "^d|^D") ~ "Dioxin",  
    TRUE ~ microplastic)) %>%  
  count(microplastic, microplastic_recoded)
```

```
# A tibble: 7 × 3  
  microplastic microplastic_recoded     n  
  <chr>         <chr>         <int>  
1 BPA           BPA             2  
2 Bisphenol A  BPA             1  
3 Dioxin       Dioxin          2  
4 0            Other           2  
5 Other        Other           1  
6 bpa          BPA             2  
7 dioxin       Dioxin          2
```

This is a more robust solution! It will catch typos as long as first letter is correct.

That's better!



Separating and uniting data



## Uniting columns

The `unite()` function can help combine columns

The `col` argument specifies new column name

The `sep` argument specifies what separator to use when combining -default is "\_"

The `remove` argument specifies if you want to drop the old columns

```
diet_comb <- data_diet %>%  
  unite(Diet, Effect, col = "change", remove = TRUE)
```

```
diet_comb
```

```
# A tibble: 12 × 4
```

```
  change      microplastic blood_level_start_nM blood_level_change_nM  
  <chr>      <chr>                <int>                <dbl>  
1 A_Increase Dioxin                  1                    2.5  
2 B_Increase Dioxin                  7                    2.5  
3 B_Increase Other                 2                    0.5  
4 A_Decrease Bisphenol A          3                   -0.5  
5 B_Increase BPA                 5                    0.5  
6 B_Increase dioxin              8                    0.5  
7 A_Increase bpa                 6                    2.5  
8 B_Increase 0                   5                    3.5  
9 B_Increase dioxin              2                    0.5  
10 A_Increase bpa                 1                    1.5  
11 B_Increase BPA                 7                    1.5  
12 B_Increase 0                   3                    2.5
```

## Separating columns based on a separator

The `separate()` function from `tidyr` can split a column into multiple columns.

The `col` argument specifies what column to work with

The `into` argument specifies names of new columns

The `sep` argument specifies what to separate by

```
diet_comb <- diet_comb %>%  
  separate(col = change, into = c("Diet", "Change"), sep = "_" )  
diet_comb
```

```
# A tibble: 12 × 5
```

	Diet <chr>	Change <chr>	microplastic <chr>	blood_level_start_nM <int>	blood_level_change_nM <dbl>
1	A	Increase	Dioxin	1	2.5
2	B	Increase	Dioxin	7	2.5
3	B	Increase	Other	2	0.5
4	A	Decrease	Bisphenol A	3	-0.5
5	B	Increase	BPA	5	0.5
6	B	Increase	dioxin	8	0.5
7	A	Increase	bpa	6	2.5
8	B	Increase	0	5	3.5
9	B	Increase	dioxin	2	0.5
10	A	Increase	bpa	1	1.5
11	B	Increase	BPA	7	1.5
12	B	Increase	0	3	2.5

## Summary

- `case_when()` requires `mutate()` when working with dataframes/tibbles
- `case_when()` can recode **entire values** based on **conditions** (need quotes for conditions and new values)
  - remember `case_when()` needs `TRUE ~ variable` to keep values that aren't specified by conditions, otherwise will be `NA`

**Note:** you might see the `recode()` function, it only does some of what `case_when()` can do, so we skipped it, but it is in the extra slides at the end.

# Summary continued

**dplyr::case\_when()** **IF ELSE...**  
(but you love it?)

df %>% **ADD COLUMN 'danger'**  
mutate(danger = case\_when(**IF type is kraken THEN danger is extreme!**  
type == "kraken" ~ "extreme!",  
**TRUE ~ "high"**))  
**OTHERWISE, danger is high.**



@allison\_horst

"Artwork by @allison\_horst". <https://allisonhorst.com/>

## Summary Continued

- `stringr` package has great functions for looking for specific **parts of values** especially `filter()` and `str_detect()` combined
- `stringr` also has other useful string functions like `str_detect()` (finding patterns in a column or vector), `str_subset()` (parsing text), `str_replace()` (replacing the first instance in values), `str_replace_all()` (replacing all instances in each value) and **more!**
- `separate()` can split columns into additional columns
- `unite()` can combine columns
- `:` can indicate when you want to start and end with columns next to one another

## Lab Part 2

- ▢ [Class Website](#)
- ▢ [Lab](#)



Image by [Gerd Altmann](#) from [Pixabay](#)

**Extra Slides**

## recode() function

This is similar to `case_when()` but it can't do as much.

(need mutate for data frames/tibbles!)

*# General Format - this is not code!*

```
{data_input} %>%  
  mutate({variable_to_fix_or_new} = recode({Variable_fixing}, {old_value} = {new  
                                           {another_old_value} = {new_value}))
```



## recode() function

Need quotes for new values! Tolerates quotes for old values.

```
data_diet %>%  
  mutate(microplastic_recoded = recode(microplastic,  
    "Bisphenol A" = "BPA",  
    "bpa" = "BPA",  
    "dioxin" = "Dioxin",  
    "0" = "Other")) %>%  
  count(microplastic, microplastic_recoded)
```

## recode()

```
data_diet %>%  
  mutate(microplastic_recoded = recode(microplastic,  
    "Bisphenol A" = "BPA",  
    "bpa" = "BPA",  
    "dioxin" = "Dioxin",  
    "0" = "Other")) %>%  
  count(microplastic, microplastic_recoded)
```

```
# A tibble: 7 × 3  
  microplastic microplastic_recoded     n  
  <chr>         <chr>         <int>  
1 BPA           BPA             2  
2 Bisphenol A  BPA             1  
3 Dioxin       Dioxin          2  
4 0            Other           2  
5 Other        Other           1  
6 bpa          BPA             2  
7 dioxin       Dioxin          2
```

## Can update or overwrite variables with recode too!

Just use the same variable name to change the variable within mutate.

```
data_diet %>%  
  mutate(microplastic = recode(microplastic,  
                                "Bisphenol A" = "BPA",  
                                "bpa" = "BPA",  
                                "dioxin" = "Dioxin",  
                                "0" = "Other")) %>%  
  count(microplastic)
```

```
# A tibble: 3 × 2  
  microplastic    n  
  <chr>         <int>  
1 BPA           5  
2 Dioxin        4  
3 Other         3
```

## String Splitting

- `str_split(string, pattern)` - splits strings up - returns list!

```
library(stringr)
x <- c("I really like writing R code")
df <- tibble(x = c("I really", "like writing", "R code programs"))
y <- unlist(str_split(x, " "))
y
```

```
[1] "I"      "really" "like"    "writing" "R"      "code"
```

```
length(y)
```

```
[1] 6
```

## A bit on Regular Expressions

- <http://www.regular-expressions.info/reference.html>
- They can use to match a large number of strings in one statement
- `.` matches any single character
- `*` means repeat as many (even if 0) more times the last character
- `?` makes the last thing optional
- `^` matches start of vector `^a` - starts with "a"
- `$` matches end of vector `b$` - ends with "b"

## Let's look at modifiers for `stringr`

?modifiers

- `fixed` - match everything exactly
- `ignore_case` is an option to not have to use `tolower`

## Using a fixed expression

One example case is when you want to split on a period ".". In regular expressions . means **ANY** character, so we need to specify that we want R to interpret "." as simply a period.

```
str_split("I.like.strings", ".")
```

```
[[1]]  
[1] "" "" "" "" "" "" "" "" "" "" "" "" "" "" ""
```

```
str_split("I.like.strings", fixed("."))
```

```
[[1]]  
[1] "I"      "like"   "strings"
```

```
str_split("I.like.strings", "\\.")
```

```
[[1]]  
[1] "I"      "like"   "strings"
```

## Pasting strings with `paste` and `paste0`

Paste can be very useful for joining vectors together:

```
paste("Visit", 1:5, sep = "_")
```

```
[1] "Visit_1" "Visit_2" "Visit_3" "Visit_4" "Visit_5"
```

```
paste("Visit", 1:5, sep = "_", collapse = "_")
```

```
[1] "Visit_1_Visit_2_Visit_3_Visit_4_Visit_5"
```

*# and `paste0` can be even simpler see `?paste0`*

```
paste0("Visit",1:5) # no space!
```

```
[1] "Visit1" "Visit2" "Visit3" "Visit4" "Visit5"
```



Comparison of **stringr** to base R -  
not covered

# Splitting Strings

# Substringing

stringr

- `str_split(string, pattern)` - splits strings up - returns list!

## Splitting String:

In `stringr`, `str_split` splits a vector on a string into a `list`

```
library(stringr)
x <- c("I really", "like writing", "R code programs")
y <- str_split(x, pattern = " ") # returns a list
y
```

```
[[1]]
[1] "I"      "really"
```

```
[[2]]
[1] "like"   "writing"
```

```
[[3]]
[1] "R"      "code"   "programs"
```

## 'Find' functions: stringr compared to base R

Base R does not use these functions. Here is a "translator" of the `stringr` function to base R functions

- `str_detect` - similar to `grep1` (return logical)
- `grep(value = FALSE)` is similar to `which(str_detect())`
- `str_subset` - similar to `grep(value = TRUE)` - return value of matched
- `str_replace` - similar to `sub` - replace one time
- `str_replace_all` - similar to `gsub` - replace many times

# Important Comparisons

Base R:

- Argument order is (pattern, x)
- Uses option (fixed = TRUE)

stringr

- Argument order is (string, pattern) aka (x, pattern)
- Uses function fixed(pattern)

## some data to work with

```
Sal <- read_csv(file =  
  "https://daseh.org/data/Baltimore_City_Employee_Salaries_FY2015.csv")
```

## Showing difference in `str_extract`

`str_extract` extracts just the matched string

```
ss <- str_extract(Sal$Name, "Rawling")
```

```
Warning: Unknown or uninitialised column: `Name`.
```

```
head(ss)
```

```
character(0)
```

```
ss[ !is.na(ss)]
```

```
character(0)
```



## Showing difference in `str_extract` and `str_extract_all`

`str_extract_all` extracts all the matched strings

```
head(str_extract(Sal$AgencyID, "\\d"))
```

```
[1] "0" "2" "6" "9" "4" "9"
```

```
head(str_extract_all(Sal$AgencyID, "\\d"), 2)
```

```
[[1]]  
[1] "0" "3" "0" "3" "1"
```

```
[[2]]  
[1] "2" "9" "0" "4" "5"
```

## Using Regular Expressions

- Look for any name that starts with:
  - Payne at the beginning,
  - Leonard and then an S
  - Spence then capital C

```
head(grep("^Payne.*", x = Sal$name, value = TRUE), 3)
```

```
[1] "Payne El,Boaz L"      "Payne El,Jackie"  
[3] "Payne Johnson,Nickole A"
```

```
head(grep("Leonard.?S", x = Sal$name, value = TRUE))
```

```
[1] "Payne,Leonard S"      "Szumlanski,Leonard S"
```

```
head(grep("Spence.*C.*", x = Sal$name, value = TRUE))
```

```
[1] "Spencer,Charles A"    "Spencer,Clarence W"  "Spencer,Michael C"
```

## Using Regular Expressions: `stringr`

```
head(str_subset( Sal$name, "^Payne.*"), 3)
```

```
[1] "Payne El,Boaz L"      "Payne El,Jackie"  
[3] "Payne Johnson,Nickole A"
```

```
head(str_subset( Sal$name, "Leonard.?S"))
```

```
[1] "Payne,Leonard S"      "Szumlanski,Leonard S"
```

```
head(str_subset( Sal$name, "Spence.*C.*"))
```

```
[1] "Spencer,Charles A"    "Spencer,Clarence W"  "Spencer,Michael C"
```

# Replace

Let's say we wanted to sort the data set by Annual Salary:

```
class(Sal$AnnualSalary)
```

```
[1] "character"
```

```
sort(c("1", "2", "10")) # not sort correctly (order simply ranks the data)
```

```
[1] "1" "10" "2"
```

```
order(c("1", "2", "10"))
```

```
[1] 1 3 2
```

## Replace

So we must change the annual pay into a numeric:

```
head(Sal$AnnualSalary, 4)
```

```
[1] "$55314.00" "$74000.00" "$64500.00" "$46309.00"
```

```
head(as.numeric(Sal$AnnualSalary), 4)
```

```
Warning in head(as.numeric(Sal$AnnualSalary), 4): NAs introduced by coercion
```

```
[1] NA NA NA NA
```

R didn't like the \$ so it thought turned them all to NA.

`sub()` and `gsub()` can do the replacing part in base R.

## Replacing and subbing

Now we can replace the \$ with nothing (used `fixed=TRUE` because \$ means ending):

```
Sal$AnnualSalary <- as.numeric(gsub(pattern = "$", replacement="",
                                   Sal$AnnualSalary, fixed=TRUE))
Sal <- Sal[order(Sal$AnnualSalary, decreasing=TRUE), ]
Sal[1:5, c("name", "AnnualSalary", "JobTitle")]
```

```
# A tibble: 5 × 3
  name           AnnualSalary JobTitle
<chr>           <dbl> <chr>
1 Mosby, Marilyn J 238772 STATE'S ATTORNEY
2 Batts, Anthony W 211785 Police Commissioner
3 Wen, Leana       200000 Executive Director III
4 Raymond, Henry J 192500 Executive Director III
5 Swift, Michael   187200 CONTRACT SERV SPEC II
```

## Replacing and subbing: `stringr`

We can do the same thing (with 2 piping operations!) in `dplyr`

```
dplyr_sal <- Sal
dplyr_sal <- dplyr_sal %>% mutate(
  AnnualSalary = AnnualSalary %>%
    str_replace(
      fixed("$"),
      "" ) %>%
    as.numeric) %>%
  arrange(desc(AnnualSalary))
check_sal = Sal
rownames(check_sal) = NULL
all.equal(check_sal, dplyr_sal)
```

```
[1] TRUE
```

## Creating Two-way Tables

A two-way table. If you pass in 2 vectors, `table` creates a 2-dimensional table.

```
tab <- table(c(0, 1, 2, 3, 2, 3, 3, 2, 2, 3),  
            c(0, 1, 2, 3, 2, 3, 3, 4, 4, 3),  
            useNA = "always")
```

```
tab
```

	0	1	2	3	4	<NA>
0	1	0	0	0	0	0
1	0	1	0	0	0	0
2	0	0	2	0	2	0
3	0	0	0	4	0	0
<NA>	0	0	0	0	0	0



## Creating Two-way Tables

```
tab_df <- tibble(x = c(0, 1, 2, 3, 2, 3, 3, 2, 2, 3),  
                y = c(0, 1, 2, 3, 2, 3, 3, 4, 4, 3))  
tab_df %>% count(x, y)
```

```
# A tibble: 5 × 3  
  x     y     n  
<dbl> <dbl> <int>  
1     0     0     1  
2     1     1     1  
3     2     2     2  
4     2     4     2  
5     3     3     4
```

# Creating Two-way Tables

```
tab_df %>%  
  count(x, y) %>%  
  group_by(x) %>% mutate(pct_x = n / sum(n))
```

```
# A tibble: 5 × 4  
# Groups:   x [4]  
   x     y     n pct_x  
<dbl> <dbl> <int> <dbl>  
1     0     0     1     1  
2     1     1     1     1  
3     2     2     2     0.5  
4     2     4     2     0.5  
5     3     3     4     1
```

# Creating Two-way Tables

```
library(scales)
tab_df %>%
  count(x, y) %>%
  group_by(x) %>% mutate(pct_x = percent(n / sum(n)))
```

```
# A tibble: 5 × 4
# Groups:   x [4]
   x     y     n pct_x
<dbl> <dbl> <int> <chr>
1     0     0     1 100%
2     1     1     1 100%
3     2     2     2  50%
4     2     4     2  50%
5     3     3     4 100%
```

# Removing columns with threshold of percent missing values

```
is.na(df) %>% head(n = 3)
```

```
      X  
[1,] FALSE  
[2,] FALSE  
[3,] FALSE
```

```
colMeans(is.na(df))#TRUE and FALSE treated like 0 and 1
```

```
X  
0
```

```
which(colMeans(is.na(df)) < 0.2) #the location of the columns <.2
```

```
X  
1
```

```
df %>% select(which(colMeans(is.na(df)) < 0.2))# remove if over 20% missing
```

```
# A tibble: 3 × 1
```

```
      X  
  <chr>  
1 I really  
2 like writing  
3 R code programs
```