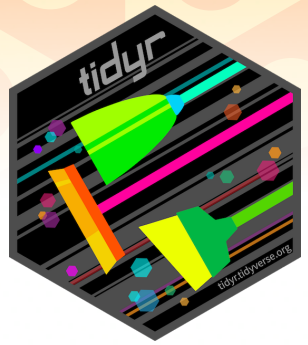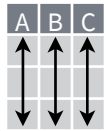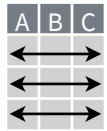# Data tidying with tidyr : : **CHEAT SHEET**

**Tidy data** is a way to organize tabular data in a consistent data structure across packages. A table is tidy if:
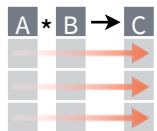


Each **variable** is in its own **column** & Each **observation**, or **case**, is in its own row

Access **variables** as **vectors** — Preserve **cases** in vectorized operations

## Tibbles

**AN ENHANCED DATA FRAME**

Tibbles are a table format provided by the **tibble** package. They inherit the data frame class, but have improved behaviors:

- **Subset** a new tibble with ], a vector with [[ and $.
- **No partial matching** when subsetting columns.
- **Display** concise views of the data on one screen.

**options**(tibble.print_max = n, tibble.print_min = m, tibble.width = Inf**)** Control default display settings.

**View()** or **glimpse()** View the entire data set.

**CONSTRUCT A TIBBLE**

**tibble(**…**)** Construct by columns.
tibble(x = 1:3, y = c("a", "b", "c"))

**tribble(**…**)** Construct by rows.
tribble(~x,  ~y,
         1, "a",
         2, "b",
         3, "c")

**Both make this tibble**

```
A tibble: 3 × 2
      x      y
  <int>  <chr>
1     1      a
2     2      b
3     3      c
```

**as_tibble(**x, …**)** Convert a data frame to a tibble.

**enframe(**x, name = "name", value = "value"**)** Convert a named vector to a tibble. Also **deframe()**.

**is_tibble(**x**)** Test whether x is a tibble.

## Reshape Data  - Pivot data to reorganize values into a new layout.

table4a →

**pivot_longer(**data, cols, names_to = "name", values_to = "value", values_drop_na = FALSE**)**

"Lengthen" data by collapsing several columns into two. Column names move to a new names_to column and values to a new values_to column.

pivot_longer(table4a, cols = 2:3, names_to ="year", values_to = "cases")

table2 →

**pivot_wider(**data, names_from = "name", values_from = "value"**)**

The inverse of pivot_longer(). "Widen" data by expanding two columns into several. One column provides the new column names, the other the values.

pivot_wider(table2, names_from = type, values_from = count)

## Split Cells  - Use these functions to split or combine cells into individual, isolated values.

table5 →

**unite(**data, col, …, sep = "_", remove = TRUE, na.rm = FALSE**)** Collapse cells across several columns into a single column.

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

table3 →

**separate(**data, col, into, sep = "[^[:alnum:]]+", remove = TRUE, convert = FALSE, extra = "warn", fill = "warn", …**)** Separate each cell in a column into several columns. Also **extract()**.

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

table3 →

**separate_rows(**data, …, sep = "[^[:alnum:].]+", convert = FALSE**)** Separate each cell in a column into several rows.

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

## Expand Tables

Create new combinations of variables or identify implicit missing values (combinations of variables not present in the data).

x →

**expand(**data, …**)** Create a new tibble with all possible combinations of the values of the variables listed in … Drop other variables.
expand(mtcars, cyl, gear, carb)

x →

**complete(**data, …, fill = list()**)** Add missing possible combinations of values of variables listed in … Fill remaining variables with NA.
complete(mtcars, cyl, gear, carb)

## Handle Missing Values

Drop or replace explicit missing values (NA).

x →

**drop_na(**data, …**)** Drop rows containing NA's in … columns.
drop_na(x, x2)

x →

**fill(**data, …, .direction = "down"**)** Fill in NA's in … columns using the next or previous value.
fill(x, x2)

x →

**replace_na(**data, replace**)** Specify a value to replace NA in selected columns.
replace_na(x, list(x2 = 2))

**posit**

# Nested Data

A **nested data frame** stores individual tables as a list-column of data frames within a larger organizing data frame. List-columns can also be lists of vectors or lists of varying data types.
Use a nested data frame to:
- Preserve relationships between observations and subsets of data. Preserve the type of the variables being nested (factors and datetimes aren't coerced to character).
- Manipulate many sub-tables at once with **purrr** functions like map(), map2(), or pmap() or with **dplyr** rowwise() grouping.
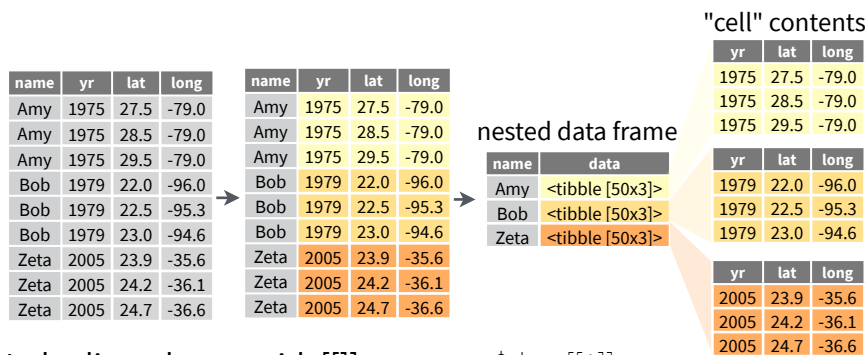
## CREATE NESTED DATA

**nest(data, …)** Moves groups of cells into a list-column of a data frame. Use alone or with dplyr::**group_by()**:

1. Group the data frame with **group_by()** and use **nest()** to move the groups into a list-column.
   ```
   n_storms <- storms %>%
       group_by(name) %>%
       nest()
   ```

2. Use **nest(new_col = c(x, y))** to specify the columns to group using dplyr::**select()** syntax.
   ```
   n_storms <- storms %>%
       nest(data = c(year:long))
   ```

| name | yr | lat | long |
|------|------|------|------|
| Amy | 1975 | 27.5 | -79.0 |
| Amy | 1975 | 28.5 | -79.0 |
| Amy | 1975 | 29.5 | -79.0 |
| Bob | 1979 | 22.0 | -96.0 |
| Bob | 1979 | 22.5 | -95.3 |
| Bob | 1979 | 23.0 | -94.6 |
| Zeta | 2005 | 23.9 | -35.6 |
| Zeta | 2005 | 24.2 | -36.1 |
| Zeta | 2005 | 24.7 | -36.6 |

nested data frame

| name | data |
|------|------|
| Amy | <tibble [50x3]> |
| Bob | <tibble [50x3]> |
| Zeta | <tibble [50x3]> |

"cell" contents

| yr | lat | long |
|------|------|------|
| 1975 | 27.5 | -79.0 |
| 1975 | 28.5 | -79.0 |
| 1975 | 29.5 | -79.0 |
| 1979 | 22.0 | -96.0 |
| 1979 | 22.5 | -95.3 |
| 1979 | 23.0 | -94.6 |
| 2005 | 23.9 | -35.6 |
| 2005 | 24.2 | -36.1 |
| 2005 | 24.7 | -36.6 |

Index list-columns with [[]]. n_storms$data[[1]]

## CREATE TIBBLES WITH LIST-COLUMNS

tibble::**tribble(**…**)** Makes list-columns when needed.
```
tribble( ~max, ~seq,
            3,   1:3,
            4,   1:4,
            5,   1:5)
```

| max | seq |
|-----|-----|
| 3 | <int [3]> |
| 4 | <int [4]> |
| 5 | <int [5]> |

tibble::**tibble(**…**)** Saves list input as list-columns.
```
tibble(max = c(3, 4, 5), seq = list(1:3, 1:4, 1:5))
```

tibble::**enframe(**x, name="name", value="value"**)**
Converts multi-level list to a tibble with list-cols.
```
enframe(list('3'=1:3, '4'=1:4, '5'=1:5), 'max', 'seq')
```

## OUTPUT LIST-COLUMNS FROM OTHER FUNCTIONS

dplyr::**mutate()**, **transmute()**, and **summarise()** will output list-columns if they return a list.
```
mtcars %>%
    group_by(cyl) %>%
    summarise(q = list(quantile(mpg)))
```
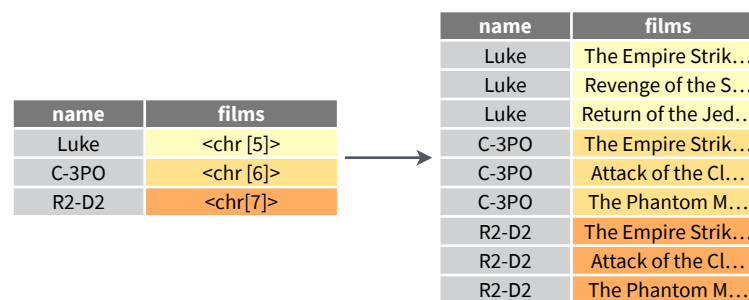
## RESHAPE NESTED DATA

**unnest(data, cols, …, keep_empty = FALSE)** Flatten nested columns back to regular columns. The inverse of nest().
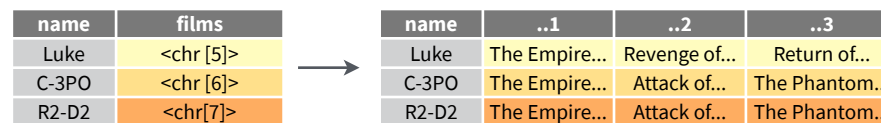n_storms %>% unnest(data)

**unnest_longer(data, col, values_to = NULL, indices_to = NULL)** Turn each element of a list-column into a row.

```
starwars %>%
    select(name, films) %>%
    unnest_longer(films)
```
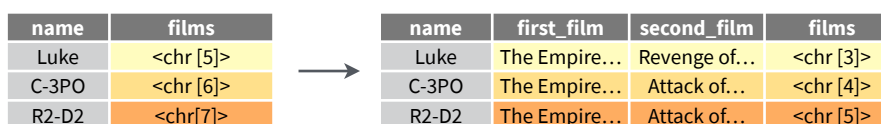
| name | films |
|------|-------|
| Luke | <chr [5]> |
| C-3PO | <chr [6]> |
| R2-D2 | <chr[7]> |

| name | films |
|------|-------|
| Luke | The Empire Strik… |
| Luke | Revenge of the S… |
| Luke | Return of the Jed… |
| C-3PO | The Empire Strik… |
| C-3PO | Attack of the Cl… |
| C-3PO | The Phantom M… |
| R2-D2 | The Empire Strik… |
| R2-D2 | Attack of the Cl… |
| R2-D2 | The Phantom M… |

**unnest_wider(data, col)** Turn each element of a list-column into a regular column.

```
starwars %>%
    select(name, films) %>%
    unnest_wider(films)
```

| name | films |
|------|-------|
| Luke | <chr [5]> |
| C-3PO | <chr [6]> |
| R2-D2 | <chr[7]> |

| name | ..1 | ..2 | ..3 |
|------|------|------|------|
| Luke | The Empire… | Revenge of… | Return of… |
| C-3PO | The Empire… | Attack of… | The Phantom… |
| R2-D2 | The Empire… | Attack of… | The Phantom… |

**hoist(**.data, .col, …, .remove = TRUE) Selectively pull list components out into their own top-level columns. Uses purrr::pluck() syntax for selecting from lists.
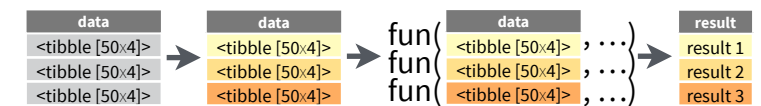
```
starwars %>%
    select(name, films) %>%
    hoist(films, first_film = 1, second_film = 2)
```

| name | films |
|------|-------|
| Luke | <chr [5]> |
| C-3PO | <chr [6]> |
| R2-D2 | <chr[7]> |

| name | first_film | second_film | films |
|------|------------|-------------|-------|
| Luke | The Empire… | Revenge of… | <chr [3]> |
| C-3PO | The Empire… | Attack of… | <chr [4]> |
| R2-D2 | The Empire… | Attack of… | <chr [5]> |

## TRANSFORM NESTED DATA

A vectorized function takes a vector, transforms each element in parallel, and returns a vector of the same length. By themselves vectorized functions cannot work with lists, such as list-columns.

dplyr::**rowwise(**.data, …**)** Group data so that each row is one group, and within the groups, elements of list-columns appear directly (accessed with [[ ), not as lists of length one. **When you use rowwise(), dplyr functions will seem to apply functions to list-columns in a vectorized fashion.**

| data |
|------|
| <tibble [50×4]> |
| <tibble [50×4]> |
| <tibble [50×4]> |

| data |
|------|
| <tibble [50×4]> |
| <tibble [50×4]> |
| <tibble [50×4]> |

fun( <tibble [50×4]> , …)
fun( <tibble [50×4]> , …)
fun( <tibble [50×4]> , …)

| result |
|--------|
| result 1 |
| result 2 |
| result 3 |

Apply a function to a list-column and **create a new list-column.**

```
n_storms %>%
    rowwise() %>%
    mutate(n = list(dim(data)))
```
> dim() returns two values per row
> wrap with list to tell mutate to create a list-column

Apply a function to a list-column and **create a regular column.**

```
n_storms %>%
    rowwise() %>%
    mutate(n = nrow(data))
```
> nrow() returns one integer per row

Collapse **multiple list-columns** into a single list-column.

```
starwars %>%
    rowwise() %>%
    mutate(transport = list(append(vehicles, starships)))
```
> append() returns a list for each row, so col type must be list

Apply a function to **multiple list-columns.**

```
starwars %>%
    rowwise() %>%
    mutate(n_transports = length(c(vehicles, starships)))
```
> length() returns one integer per row

See **purrr** package for more list functions.

posit™