# Strings and Factors

# Strings

- "character" data type
- Values within either single or double quotes
- Useful when you want to subset, extract, or parse character strings
- Based on matching a pattern defined using "regular expression" syntax



stringr helps standardize string manipulations

# What is a string?

- Characters within single or double quotes (double preferred)
- Special characters require a leading backslash
  - Single or double quotes
  - Backslash
  - New Line
  - Tab
  - Unicode
- R will list all special characters with "?Quote"

# What is a string?

- Characters within single or double quotes (double preferred)
- Special characters require a leading backslash
    - Single or double quotes
    - Backslash
    - New Line
    - Tab
    - Unicode
- R will list all special characters with "?Quote"

Not a string?

   A)  "I am a string"         B) '23by51'         C) I_am_a_string

# Strings in Statistics and Modeling of ID

When might you encounter strings?

# Strings in Statistics and Modeling of ID

When might you encounter strings?

- Line level data
  - name, care notes, compound columns (27M, 35F)
- Compartment/class names output as part of a simulation
- Loading many data files with similar names
- …

# Regular Expressions, regex, regexp

- Developed in the 1950s, and for most computer languages
- Syntax used to describe a pattern that will be matched in a string
- Each character in an expression is either:
  - regular - literal meaning of the character
  - Metacharacter - special meaning dictated by language specific regex syntax
    - Metacharacters: `.^$\|*+?{}[]()`
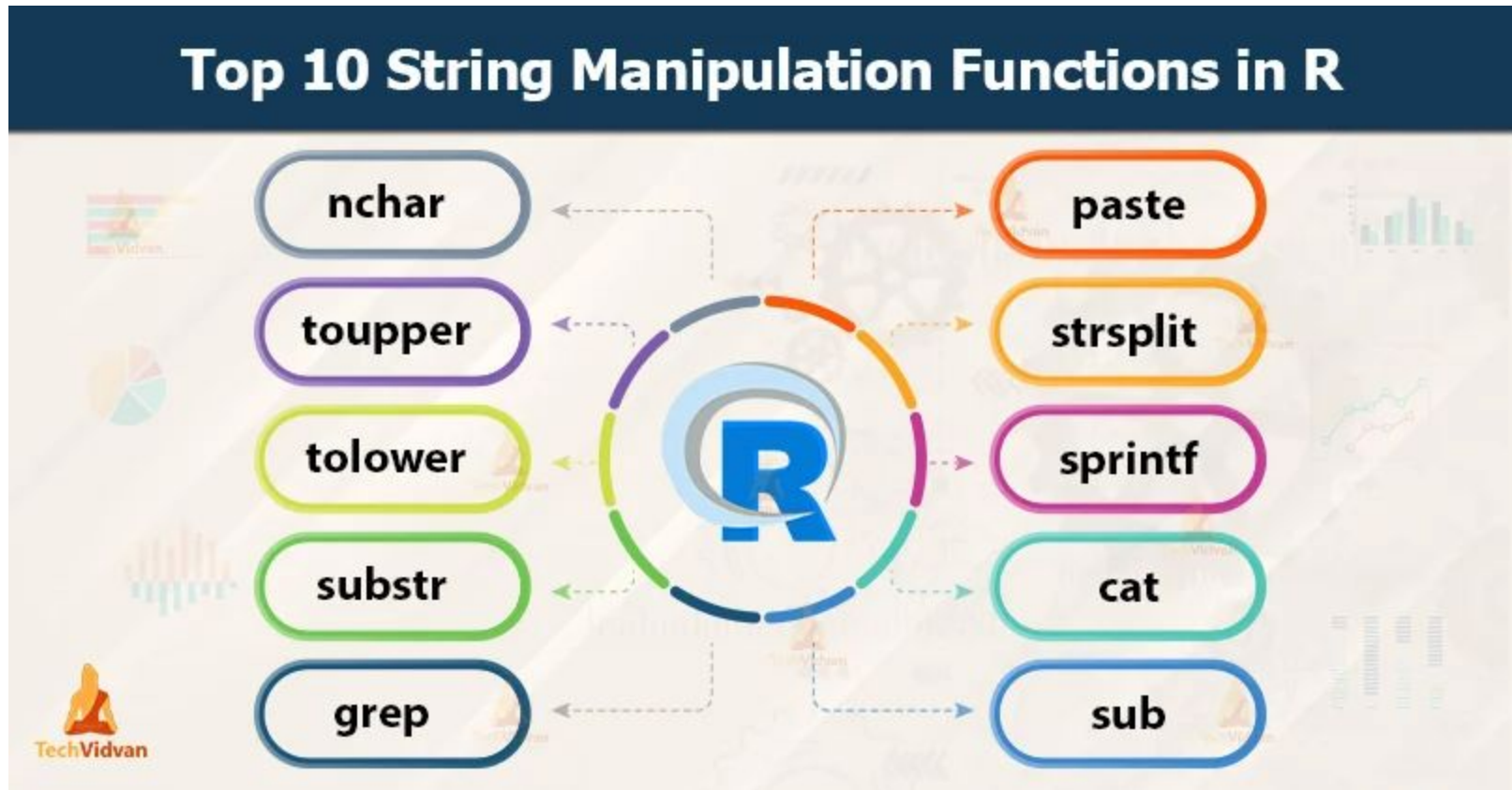    - Using a metacharacter as a regular character requires "escaping" the character

Cheatsheet: https://github.com/rstudio/cheatsheets/blob/main/strings.pdf

# Regular Expressions, regex, regexp

- Developed in the 1950s, and for most computer languages
- Syntax used to describe a pattern that will be matched in a string
- Each character in an expression is either:
  - regular - literal meaning of the character
  - Metacharacter - special meaning dictated by language specific regex syntax
    - Metacharacters: `.^$\|*+?{}[]()`
    - Using a metacharacter as a regular character requires "escaping" the character

Examples:

- "Summari[zs]e" would match both "Summarize" and "Summarise"

- ".a." would match any string with an "a" contained <u>within</u> the string, it would match "banana" but not "apple"

Cheatsheet: https://github.com/rstudio/cheatsheets/blob/main/strings.pdf

# Common Base R functions



**Top 10 String Manipulation Functions in R**

nchar

toupper

tolower

substr

grep

paste

strsplit

sprintf

cat

sub

TechVidvan

From: https://techvidvan.com/tutorials/r-string-manipulation/

# Systematic Names for String Manipulation Functions

All stringr functions start with "str_"

```
> ○ str_c          {stringr}    str_c(..., sep = "", collapse = NULL)
>
> ◆ str_conv       {stringr}    To understand how str_c works, you need to imagine that you are
>                               building up a matrix of strings. Each input argument forms a
> ◆ str_count      {stringr}    column, and is expanded to the length of the longest argument,
>                               using the usual recyling rules. The sep string is inserted between
> ◆ str_detect     {stringr}    each column. If collapse is NULL each row is collapsed into a single
>                               string. If non-NULL that string is inserted at the end of each row,
> ◆ str_dup        {stringr}    and the entire matrix collapsed to a single string.
>
> ◆ str_extract    {stringr}
>                               Press F1 for additional help
> ◆ str_extract_all {stringr}
> str_|
```

# Give it a try

Regex for strings with:

1. containing "A"
2. containing any number
3. ending with any uppercase letter
4. contains a question mark
5. any two repeated letters

Open Exercise.Rmd to check responses using str_view()

# Give it a try

| Pattern | Regex |
|---|---|
| containing "A" | A |
| containing any number | [:digit:] |
| ending with any uppercase letter | [:upper:]$ |
| contains a question mark | [?] |
| any two repeated letters | ([:alpha:][:alpha:])\\1 |

# Questions and switch to Rmd

# Factors - a type of string

- are categorical variables, variables that have a fixed and known set of possible values
- most often interact with factors when:
  - re-ordering
  - data visualizations
  - creating sub-groups within factors

Examples:

c("Jan","Feb", "Mar", "April")                    c("1","2","3")

c("low", "medium", "high")

# Ways to interact with factors

| Manipulation | Change order | Change values | Add/Drop levels | Combine factors |
|---|---|---|---|---|
| Original | 1. 'Feb'<br>2. 'Jan'<br>3. 'Mar' | 1. 'Jan'<br>2. 'Feb'<br>3. 'Mar' | 1. 'Jan'<br>2. 'Feb' | 1. 'Jan'<br>2. 'Feb'<br>3. 'Mar' |
| Modified | 1. 'Jan'<br>2. 'Feb'<br>3. 'Mar' | 1. 'January'<br>2. 'February'<br>3. 'March' | 1. 'Jan'<br>2. 'Feb'<br>3. 'Mar' | 1. Winter ('Jan', 'Feb')<br>2. Spring ('Mar') |

# Ways to interact with factors

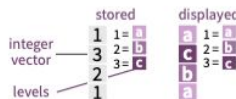| Manipulation | Change order | Change values | Add/Drop levels | Combine factors |
|---|---|---|---|---|
| Original | 1. 'Feb'<br>2. 'Jan'<br>3. 'Mar' | 1. 'Jan'<br>2. 'Feb'<br>3. 'Mar' | 1. 'Jan'<br>2. 'Feb' | 1. 'Jan'<br>2. 'Feb'<br>3. 'Mar' |
| Modified | 1. 'Jan'<br>2. 'Feb'<br>3. 'Mar' | 1. 'January'<br>2. 'February'<br>3. 'March' | 1. 'Jan'<br>2. 'Feb'<br>3. 'Mar' | 1. Winter ('Jan', 'Feb')<br>2. Spring ('Mar') |
| Forcats function | Use the cheat sheet to find the appropriate function | | | |

# Factors with forcats : : **CHEATSHEET**

The **forcats** package provides tools for working with factors, which are R's data structure for categorical data.
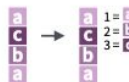
## Factors

R represents categorical data with factors. A **factor** is an integer vector with a **levels** attribute that stores a set of mappings between integers and categorical values. When you view a factor, R displays not the integers, but the levels associated with them.
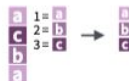
*Create a factor with factor()*

**factor**(x = character(), levels, labels = levels, exclude = NA, ordered = is.ordered(x), nmax = NA) Convert a vector to a factor. Also **as_factor()**.
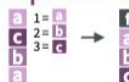f <- factor(c("a", "c", "b", "a"), levels = c("a", "b", "c"))

*Return its levels with levels()*

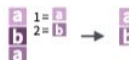**levels**(x) Return/set the levels of a factor. levels(f); levels(f) <- c("x","y","z")
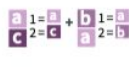
*Use unclass() to see its structure*

## Inspect Factors

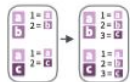**fct_count**(f, sort = FALSE, prop = FALSE) Count the number of values with each level. fct_count(f)

**fct_match**(f, lvls) Check for lvls in f. fct_match(f, "a")

**fct_unique**(f) Return the unique values, removing duplicates. fct_unique(f)

## Combine Factors

**fct_c**(...) Combine factors with different levels. Also **fct_cross()**.
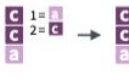f1 <- factor(c("a", "c"))
f2 <- factor(c("b", "a"))
fct_c(f1, f2)

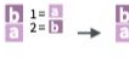**fct_unify**(fs, levels = lvls_union(fs)) Standardize levels across a list of factors. fct_unify(list(f2, f1))
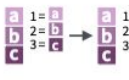
## Change the order of levels

**fct_relevel**(.f, ..., after = 0L) Manually reorder factor levels. fct_relevel(f, c("b", "c", "a"))

**fct_infreq**(f, ordered = NA) Reorder levels by the frequency in which they appear in the data (highest frequency first). Also **fct_inseq()**.
f3 <- factor(c("c", "c", "a"))
fct_infreq(f3)

**fct_inorder**(f, ordered = NA) Reorder levels by order in which they appear in the data. fct_inorder(f2)

**fct_rev**(f) Reverse level order.
f4 <- factor(c("a","b","c"))
fct_rev(f4)

**fct_shift**(f) Shift levels to left or right, wrapping around end. fct_shift(f4)

**fct_shuffle**(f, n = 1L) Randomly permute order of factor levels. fct_shuffle(f4)

**fct_reorder**(.f, .x, .fun = median, ..., .desc = FALSE) Reorder levels by their relationship with another variable.
boxplot(
  PlantGrowth,
  weight ~ fct_reorder(group, weight)
)

**fct_reorder2**(.f, .x, .y, .fun = last2, ..., .desc = TRUE) Reorder levels by their final values when plotted with two other variables.
ggplot(
  diamonds,
  aes(
    carat, price,
    color = fct_reorder2(color, carat, price)
  )) +
  geom_smooth()

## Change the value of levels

**fct_recode**(.f, ...) Manually change levels. Also **fct_relabel()** which obeys purrr::map syntax to apply a function or expression to each level.
fct_recode(f, v = "a", x = "b", z = "c")
fct_relabel(f, ~ paste0("x", .x))

**fct_anon**(f, prefix = "") Anonymize levels with random integers. fct_anon(f)

**fct_collapse**(.f, ..., other_level = NULL) Collapse levels into manually defined groups. fct_collapse(f, x = c("a", "b"))

**fct_lump_min**(f, min, w = NULL, other_level = "Other") Lumps together factors that appear fewer than min times. Also **fct_lump_n()**, **fct_lump_prop()**, and **fct_lump_lowfreq()**. fct_lump_min(f, min = 2)

**fct_other**(f, keep, drop, other_level = "Other") Replace levels with "other." fct_other(f, keep = c("a", "b"))

## Add or drop levels

**fct_drop**(f, only) Drop unused levels.
f5 <- factor(c("a","b"),c("a","b","x"))
f6 <- fct_drop(f5)

**fct_expand**(f, ...) Add levels to a factor. fct_expand(f6, "x")

**fct_na_value_to_level**(f, level = "(Missing)") Assigns a level to NAs to ensure they appear in plots, etc.
f7 <- factor(c("a", "b", NA))
fct_na_value_to_level(f7, level = "(Missing)")

# Ways to interact with factors

| Manipulation | Change order | Change values | Add/Drop levels | Combine factors |
|---|---|---|---|---|
| Original | 1. 'Feb'<br>2. 'Jan'<br>3. 'Mar' | 1. 'Jan'<br>2. 'Feb'<br>3. 'Mar' | 1. 'Jan'<br>2. 'Feb' | 1. 'Jan'<br>2. 'Feb'<br>3. 'Mar' |
| Modified | 1. 'Jan'<br>2. 'Feb'<br>3. 'Mar' | 1. 'January'<br>2. 'February'<br>3. 'March' | 1. 'Jan'<br>2. 'Feb'<br>3. 'Mar' | 1. Winter ('Jan', 'Feb')<br>2. Spring ('Mar') |
| Forcats function | fct_relevel() | fct_recode() | fct_drop() | fct_collapse() |

# Questions and switch back to Rmd

# Functions: A way to automate a process

Writing a function can be a time saver or sink

Knowing what **could** be a function vs what **should** be a function comes with experience

You're always going to be wrong at some point

What are some potentially good candidate tasks?



"I SPEND A LOT OF TIME ON THIS TASK. I SHOULD WRITE A PROGRAM AUTOMATING IT!"

THEORY:

WORK — WRITING CODE, WORK ON ORIGINAL TASK, AUTOMATION TAKES OVER, FREE TIME

TIME

REALITY:

WORK — WRITING CODE, DEBUGGING, RETHINKING, ONGOING DEVELOPMENT, NO TIME FOR ORIGINAL TASK ANYMORE

TIME

# Functions: A way to automate a process

Writing a function can be a time saver or sink

Knowing what **could** be a function vs what **should** be a function comes with experience

You're always going to be wrong at some point

What are some potentially good candidate tasks?

- Processing standardized data
- ODEs/Dynamic Models
- Common data visualizations



"I SPEND A LOT OF TIME ON THIS TASK.
I SHOULD WRITE A PROGRAM AUTOMATING IT!"

THEORY:

WORK

WRITING CODE

WORK ON ORIGINAL TASK

AUTOMATION TAKES OVER

FREE TIME

TIME

REALITY:

WORK

WRITING CODE

DEBUGGING

RETHINKING

ONGOING DEVELOPMENT

NO TIME FOR ORIGINAL TASK ANYMORE

TIME

# Anatomy of a function

```
FunctionName <- function(arguement1,
arguement2, ...){
  #some analysis
  return(outputOfFunction)
}
```

# Document with Comments!

```r
CalculateSampleCovariance <- function(x, y, verbose = TRUE) {
  # Computes the sample covariance between two vectors.
  #
  # Args:
  #   x: One of two vectors whose sample covariance is to be calculated.
  #   y: The other vector. x and y must have the same length, greater than one,
  #      with no missing values.
  #   verbose: If TRUE, prints sample covariance; if not, not. Default is TRUE.
  #
  # Returns:
  #   The sample covariance between x and y.
  n <- length(x)
  # Error handling
  if (n <= 1 || n != length(y)) {
    stop("Arguments x and y have invalid lengths: ",
         length(x), " and ", length(y), ".")
  }
  if (TRUE %in% is.na(x) || TRUE %in% is.na(y)) {
    stop(" Arguments x and y must not have missing values.")
  }
  covariance <- var(x, y)
  if (verbose)
    cat("Covariance = ", round(covariance, 4), ".\n", sep = "")
  return(covariance)
}
```

```r
#' Add together two numbers.
#'
#' @param x A number.
#' @param y A number.
#'
#' @return The sum of x and y.
#'
#' @examples
#' add(1, 1)
#' add(10, 1)
#'
add <- function(x, y) {
    # general comment
    x + y  # inline comment
}


#' Next function
```

# Document with Comments!

```
CalculateSampleCovariance <- function(x, y, verbose = TRUE) {
  # Computes the sample covariance between two vectors.
  #
  # Args:
  #   x: One of two vectors whose sample covariance is to be calculated.
  #   y: The other vector. x and y must have the same length, greater than one,
  #       with no missing values.
  #   verbose: If TRUE, prints sample covariance; if not, not. Default is TRUE.
  #
  # Returns:
  #   The sample covariance between x and y.
  n <- length(x)
  # Error handling
  if (n <= 1 || n != length(y)) {
    stop("Arguments x and y have invalid lengths: ",
          length(x), " and ", length(y), ".")
  }
  if (TRUE %in% is.na(x) || TRUE %in% is.na(y)) {
    stop(" Arguments x and y must not have missing values.")
  }
  covariance <- var(x, y)
  if (verbose)
    cat("Covariance = ", round(covariance, 4), ".\n", sep = "")
  return(covariance)
}
```

```
#' Add together two numbers.
#'
#' @param x A number.
#' @param y A number.
#'
#' @return The sum of x and y.
#'
#' @examples
#' add(1, 1)
#' add(10, 1)
#'
add <- function(x, y) {
    # general comment
    x + y  # inline comment
}


#' Next function
```

# Hard vs soft coding

Hard coding:  values or functions are directly embedded into the code

Soft coding:  values or functions are set in parameters that is reference by the code

# Hard vs soft coding

Hard coding:  values or functions are directly embedded into the code

Soft coding:  values or functions are set in parameters that is reference by the code

```
travel.data %>%
filter((!AddressState %in% c("ALBERTA", "BC", "NA","CANADA", "IN TRANSIT TO
CANADA")))
```

# Hard vs soft coding

Hard coding:  values or functions are directly embedded into the code

Soft coding:  values or functions are set in parameters that is reference by the code

```
travel.data %>%
filter((!AddressState %in% c("ALBERTA", "BC", "NA","CANADA", "IN TRANSIT TO
CANADA")))
```

# Hard vs soft coding

Hard coding:  values or functions are directly embedded into the code

Soft coding:  values or functions are set in parameters that is reference by the code

This is an example of hard coding

```
travel.data %>%
filter((!AddressState %in% c("ALBERTA", "BC", "NA","CANADA", "IN TRANSIT TO CANADA")))
```

# Hard vs soft coding

Hard coding:  values or functions are directly embedded into the code

Soft coding:  values or functions are set in parameters that is reference by the code

How would you convert to soft coding?

```
travel.data %>%
filter((!AddressState %in% c("ALBERTA", "BC", "NA","CANADA", "IN TRANSIT TO
CANADA")))
```

# Hard vs soft coding

Hard coding:  values or functions are <u>directly embedded</u> into the code

Soft coding:  values or functions are set in <u>parameters</u> that is reference by the code

Converted to soft coding

```
drop.names <- c("ALBERTA", "BC", "NA","CANADA", "IN TRANSIT TO CANADA")

travel.data %>%
 filter((!AddressState %in% drop.names))
```

# Hard vs soft coding

Hard coding:  values or functions are <u>directly embedded</u> into the code

```
travel.data %>%
filter((!AddressState %in% c("ALBERTA", "BC",
"NA","CANADA", "IN TRANSIT TO CANADA")))
```

Use when:

Soft coding:  values or functions are set in <u>parameters</u> that is reference by the code

```
drop.names <- c("ALBERTA", "BC",
"NA","CANADA", "IN TRANSIT TO CANADA")

travel.data %>%
 filter((!AddressState %in% drop.names))
```

Use when:

# Hard vs soft coding

Hard coding:  values or functions are <u>directly embedded</u> into the code

```
travel.data %>%
filter((!AddressState %in% c("ALBERTA", "BC",
"NA","CANADA", "IN TRANSIT TO CANADA")))
```

Use when:

Value:
- never changes
- is referenced only once
- Is short, and direct embedding makes the code easier to understand

Soft coding:  values or functions are set in <u>parameters</u> that is reference by the code

```
drop.names <- c("ALBERTA", "BC",
"NA","CANADA", "IN TRANSIT TO CANADA")

travel.data %>%
 filter((!AddressState %in% drop.names))
```

Use when:

Value:
- changes often
- is referenced multiple times
- is ugly, making it hard to read code

# Hard vs soft coding

Hard coding:  values or functions are <u>directly embedded</u> into the code

Soft coding:  values or functions are set in <u>parameters</u> that is reference by the code

```
travel.data %>%
filter((!AddressState %
"NA","CANADA", "IN T
```

```
ERTA", "BC",
RANSIT TO CANADA")

%in% drop.names))
```

**Take home:**
1.  Soft coding is essential for flexible functions
2. Soft coding is strategy to make code easier to read and/or modify,  but comes at the cost of increased nesting structure

Use when:

Value:
- never changes
- is referenced only once
- Is short, and direct embedding makes the code easier to understand

Value:
- changes often
- is referenced multiple times
- is ugly, making it hard to read code

# Questions and switch back to Rmd