#### CSSS508, Week 8

Strings

Breon Haskett

Nov 17, 2021

Updated: Nov 17, 2021



#### Data Today

We'll use data on food safety inspections in King County from <u>data.kingcounty.gov</u>.

Note these data are *fairly large*. You may want to save them and load them from a *local directory*.

I recommend specifying the column types so they read in correctly.

#### glimpse(restaurants)

```
## Rows: 258,630
## Columns: 23
## $ Name
                               <chr> "@ THE SHACK, LLC ", "10 MERCER R~
## $ Program Identifier
                               <chr> "SHACK COFFEE", "10 MERCER RESTAU~
## $ Inspection Date
                               <chr> NA, "01/24/2017", "01/24/2017", "~
## $ Description
                               <chr> "Seating 0-12 - Risk Category I".~
## $ Address
                               <chr> "2920 SW AVALON WAY", "10 MERCER \sim
## $ City
                               <chr> "Seattle", "Seattle", "Seattle", ~
                               <chr> "98126", "98109", "98109", "98109~
## $ Zip Code
## $ Phone
                               ## $ Longitude
                               <dbl> -122.3709, -122.3562, -122.3562, ~
## $ Latitude
                               <dbl> 47.57043, 47.62506, 47.62506, 47.~
## $ Inspection_Business_Name
                               <chr> NA, "10 MERCER RESTAURANT", "10 M~
## $ Inspection Type
                               <chr> NA, "Routine Inspection/Field Rev~
## $ Inspection Score
                               <int> NA, 10, 10, 10, 15, 15, 15, 0, 15~
## $ Inspection Result
                               <chr> NA, "Unsatisfactory", "Unsatisfac~
## $ Inspection_Closed_Business <chr> NA, "false", "false", "false", "f~
## $ Violation Type
                               <chr> NA, "blue", "blue", "red", "blue"~
## $ Violation Description
                               <chr> NA, "4300 - Non-food contact surf~
## $ Violation Points
                               <int> 0, 3, 2, 5, 5, 5, 5, 0, 5, 10, 25~
## $ Business ID
                               <chr> "PR0048053", "PR0049572", "PR0049~
## $ Inspection Serial Num
                               <chr> NA, "DAHSIBSJT", "DAHSIBSJT", "DA~
                               <chr> NA, "IV43WZVLN", "IVCQ1ZIV0", "IV~
## $ Violation Record ID
                               <int> NA, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, ~
## $ Grade
## $ Date
                               <date> NA, 2017-01-24, 2017-01-24, 2017~
```

#### Strings

A general programming term for a unit of character data is a **string**, which is defined as *a sequence of characters*. In R the terms "strings" and "character data" are mostly interchangeable.

In other languages, "string" often also refers to a *sequence* of numeric information, such as binary strings (e.g. "01110000 01101111 01101111 01110000"). We rarely use these in R.

Note that these are *sequences* of numbers rather than single numbers, and thus *strings*.

One thing that separates a string from a number is that the leading zeroes are meaningful: 01 != 1



-UW CS&SS

#### nchar()

We've seen the nchar() function to get the number of characters in a string. How many characters are in the ZIP codes?

```
restaurants %>%
  mutate(ZIP_length = nchar(Zip_Code)) %>%
  count(ZIP_length)
```

#### substr()

You should be familiar with substr() from the homeworks. We can use it to pull out just the first 5 digits of the ZIP code.

```
restaurants <- restaurants %>%
    mutate(ZIP_5 = substr(Zip_Code, 1, 5))
restaurants %>% distinct(ZIP_5) %>% head()

## # A tibble: 6 x 1

## ZIP_5

## <chr>
## 1 98126

## 2 98109

## 3 98101

## 4 98032
```

7 / 39

## 5 98102 ## 6 98004

#### paste()

We can combine parts of strings together using the paste() function, e.g. to make a whole mailing address:

```
## # A tibble: 6 x 1
## mailing_address
## <chr>
## 1 2920 SW AVALON WAY, Seattle, WA 98126
## 2 10 MERCER ST, Seattle, WA 98109
## 3 1001 FAIRVIEW AVE N Unit 1700A, SEATTLE, WA 98109
## 4 1225 1ST AVE, SEATTLE, WA 98101
## 5 18114 E VALLEY HWY, KENT, WA 98032
## 6 121 11TH AVE E, SEATTLE, WA 98102
```

### paste0()

paste0() is a shortcut for paste() without any separator.

```
paste(1:5, letters[1:5]) # sep is a space by default

## [1] "1 a" "2 b" "3 c" "4 d" "5 e"

paste(1:5, letters[1:5], sep ="")

## [1] "1a" "2b" "3c" "4d" "5e"

paste0(1:5, letters[1:5])

## [1] "1a" "2b" "3c" "4d" "5e"
```

#### paste() Practice

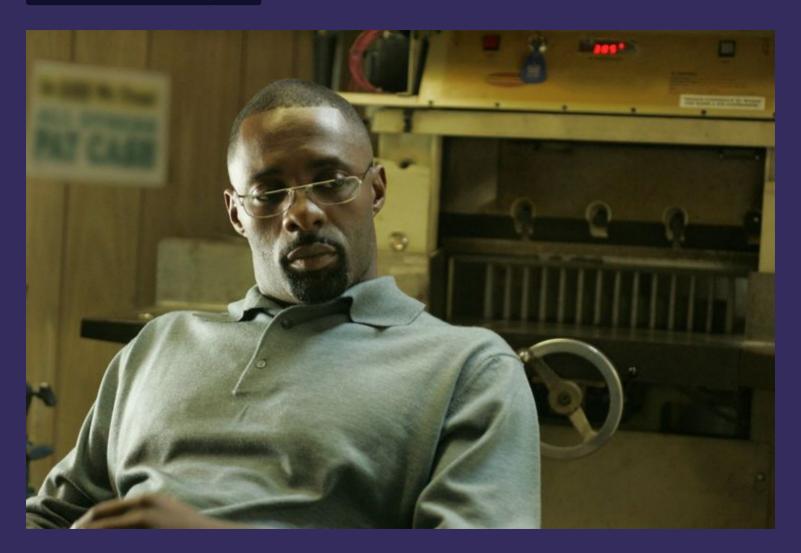
sep= controls what happens when doing entry-wise squishing of vectors you
give to paste(), while collapse= controls if/how they go from a vector to a
single string.

Here are some examples; make sure you understand how each set of arguments produces its results:

```
paste(letters[1:5], collapse = "!")
paste(1:5, letters[1:5], sep = "+")
paste0(1:5, letters[1:5], collapse = "???")
paste(1:5, "Z", sep = "*")
paste(1:5, "Z", sep = "*", collapse = " ~ ")
```

```
## [1] "a!b!c!d!e"
## [1] "1+a" "2+b" "3+c" "4+d" "5+e"
## [1] "1a???2b???3c???4d???5e"
## [1] "1*Z" "2*Z" "3*Z" "4*Z" "5*Z"
## [1] "1*Z ~ 2*Z ~ 3*Z ~ 4*Z ~ 5*Z"
```

# stringr



–UW CS&SS-

#### stringr

stringr is yet another R package from the Tidyverse (like ggplot2, dplyr, tidyr, lubridate, readr).

It provides functions that:

- Replace some basic string functions like paste() and nchar() in a way that's a bit less touchy with missing values or factors
- Remove whitespace or pad it out
- Perform tasks related to **pattern matching**: Detect, locate, extract, match, replace, split.
  - These functions use **regular expressions** to describe patterns
  - o Base R and stringi versions for these exist but are harder to use

Conveniently, *most* stringr functions begin with "str\_" to make RStudio auto-complete more useful.

library(stringr)

### stringr Equivalencies

• str\_sub() is like substr() but also lets you put in negative values to count backwards from the end (-1 is the end, -3 is third from end):

```
str_sub("Washington", 1, -3)
```

```
## [1] "Washingt"
```

• str\_c() ("string combine") is just like paste() but where the default is sep = "" (like paste0())

```
str_c(letters[1:5], 1:5)
```

```
## [1] "a1" "b2" "c3" "d4" "e5"
```

### stringr Equivalencies

• str\_length() is equivalent to nchar():

```
nchar("weasels")

## [1] 7

str_length("weasels")

## [1] 7
```

#### Changing Cases

str\_to\_upper(), str\_to\_lower(), str\_to\_title() convert cases, which
is often a good idea to do before searching for values:

```
head(unique(restaurants$City))

## [1] "Seattle" "SEATTLE" "KENT" "BELLEVUE" "KENMORE" "Issaquah"

restaurants <- restaurants %>%
    mutate(across(c(Name, Address, City), ~str_to_upper(.)))
head(unique(restaurants$City))

## [1] "SEATTLE" "KENT" "BELLEVUE" "KENMORE" "ISSAQUAH" "BURIEN"
```

### str\_trim() Whitespace

Extra leading or trailing whitespace is common in text data:

head(unique(restaurants\$Name), 4)

```
## [1] "@ THE SHACK, LLC " "10 MERCER RESTAURANT" ## [3] "100 LB CLAM" "1000 SPIRITS"
```

Any character column is potentially affected. We can use the str\_trim() function in stringr to clean them up all at once:

```
restaurants <- restaurants %>%
  mutate(across(where(is.character), ~str_trim(.)))
head(unique(restaurants$Name), 4)
```

```
## [1] "@ THE SHACK, LLC" "10 MERCER RESTAURANT" ## [3] "100 LB CLAM" "1000 SPIRITS"
```

across(where(x),  $\sim$  y) applies function y to every column for which function x returns TRUE.

Regular Expressions and Pattern Matching

·UW CS&SS**-**

# What are Regular Expressions?

**Regular expressions** or **regex**es are how we describe patterns we are looking for in text in a way that a computer can understand. We write an **expression**, apply it to a string input, and then compute with **matches** we find.

- Literal characters are defined snippets to search for like SEA or 206
- Metacharacters let us be flexible in describing patterns:
  - backslash \, caret ^, dollar sign \$, period ., pipe |, question mark
     ?, asterisk \*, plus sign +, parentheses ( and ), square brackets [ and ], curly braces { and }
  - To treat a metacharacter as a literal character, you must escape it with two preceding backslashs \\, e.g. to match (206) including the parentheses, you'd use \\(206\\) in your regex

#### str\_detect()

I want to get inspections for coffee shops. I'll say a coffee shop is anything that has "COFFEE", "ESPRESSO", or "ROASTER" in the name. The regex for this is COFFEE | ESPRESSO | ROASTER because | is a metacharacter that means "OR". Use the str\_detect() function, which returns TRUE if it finds what you're looking for and FALSE if it doesn't (similar to grepl()):

```
coffee <- restaurants %>%
  filter(str_detect(Name, "COFFEE|ESPRESSO|ROASTER"))
coffee %>% distinct(Name) %>% head()
```

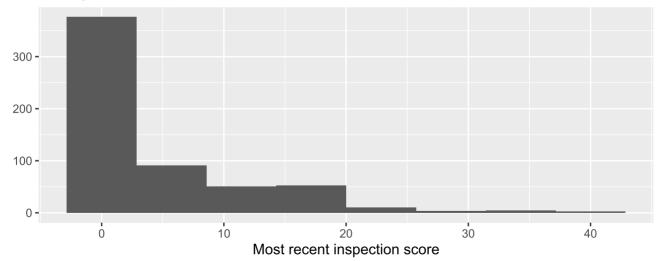
```
## # A tibble: 6 x 1
## Name
## <chr>
## 1 2 SISTERS ESPRESSO
## 2 701 COFFEE
## 3 909 COFFEE AND WINE
## 4 AJ'S ESPRESSO
## 5 ALKI HOMEFRONT SMOOTHIES & ESPRESSO
## 6 ALL CITY COFFEE
```

#### Will My Coffee Kill Me?

Let's take each unique business identifier, keep the most recent inspection score, and look at a histogram of scores:

```
coffee %>% select(Business_ID, Name, Inspection_Score, Date) %>%
    group_by(Business_ID) %>% filter(Date == max(Date)) %>%
    distinct(.keep_all=TRUE) %>% ggplot(aes(Inspection_Score)) +
    geom_histogram(bins=8) + xlab("Most recent inspection score") + ylab("") +
    ggtitle("Histogram of inspection scores for Seattle coffee shops")
```

#### Histogram of inspection scores for Seattle coffee shops



#### str\_detect(): Patterns

Let's look for phone numbers whose first three digits are "206" using str\_detect().

We will want it to work whether they have parentheses around the beginning or not, but NOT to match "206" occurring elsewhere:

## [1] TRUE TRUE TRUE FALSE

- ^ is a metacharacter meaning "look only at the *beginning* of the string"
- \\(? means look for a left parenthesis (\\(), but it's optional (?)
- 206 is the literal string to look for after the optional parenthesis

#### str\_view()

stringr also has a function called str\_view() that allows you to see in the viewer pane *exactly* what text is being selected with a regular expression.

```
str_view(phone_test_examples, area_code_206_pattern)
```

This will generate a small web page in the viewer pane (but not in Markdown docs).

```
2061234567
(206)1234567
(206)123-4567
555-206-1234
```

Just be careful to not load an entire long vector / variable or it may crash RStudio as it tries to render a massive page!

W CS&SS-----

#### str\_detect()

Perhaps we want to know how many phone numbers aren't in the 206 area code?

str\_detect() returns NA for rows with missing (NA) phone numbers--you
can't search for text in a missing value.

#### str\_extract()

str\_extract() extracts substrings that match a regex.

Let's extract the <u>directional part of Seattle</u> of addresses: N, NW, SE, none, etc.

```
## [1] " W" " NW " NA NA
```

- The first space will match a space character, then
- (N|NW|NE|S|SW|SE|W|E) matches one of the directions in the group
- ( |\$) is a group saying either there is a space after, or it's the end of the address string (\$ means the end of the string)

#### Where are the Addresses?

```
## # A tibble: 9 x 2
##
     city region
     <chr>
##
                  <int>
## 1 NF
                   2086
## 2 S
                   1764
## 3 <NA>
                   1745
## 4 N
                    879
## 5 SE
                    868
## 6 SW
                    705
## 7 F
                    538
## 8 NW
                    438
## 9 W
                    235
```

A common operation is to str\_extract() something with extra spaces and then str\_trim() them out.

## str\_replace(): Replacing

Maybe we want to do a street-level analysis of inspections (e.g. compare The Ave to Pike Street). How can we remove building numbers?

```
number_pattern <- "^[0-9]*-?[A-Z]? (1/2 )?"
number_examples <-
    c("2812 THORNDYKE AVE W", "1ST AVE", "10A 1ST AVE",
        "10-A 1ST AVE", "5201-B UNIVERSITY WAY NE",
        "7040 1/2 15TH AVE NW")
str_replace(number_examples, number_pattern, replacement = "")</pre>
```

```
## [1] "THORNDYKE AVE W" "1ST AVE" "1ST AVE" "1ST AVE" "UNIVERSITY WAY NE" "15TH AVE NW"
```

We can also use the shortcut str\_remove():

```
str_remove(number_examples, number_pattern)

## [1] "THORNDYKE AVE W" "1ST AVE" "1ST AVE"

## [4] "1ST AVE" "UNIVERSITY WAY NE" "15TH AVE NW"
```

# How Does the Building Number regex Work?

Let's break down  $"^[0-9]*-?[A-Z]? (1/2)?"$ :

- ^[0-9] means look for a digit between 0 and 9 ([0-9]) at the beginning (^)
- \* means potentially match more digits after that
- -? means optionally (?) match a hyphen (-)
- [A-Z]? means optionally match (?) a letter ([A-Z])
- Then we match a space ( )
- (1/2 )? optionally matches a 1/2 followed by a space since this is apparently a thing with some address numbers

#### Removing Street Numbers

```
restaurants <- restaurants %>%
  mutate(street_only = str_remove(Address, number_pattern))
restaurants %>% distinct(street_only) %>% head(10)
## # A tibble: 10 x 1
```

```
##
     street only
     <chr>
##
##
   1 SW AVAION WAY
   2 MFRCFR ST
##
##
   3 FAIRVIEW AVE N UNIT 1700A
##
   4 1ST AVF
##
   5 E VALLEY HWY
## 6 11TH AVE E
## 7 112TH AVE NE #125
## 8 NE BOTHELL WAY
##
   9 NW GILMAN BL C-08
## 10 NE 20TH ST STE 300
```

#### How About Units/Suites Too?

Getting rid of unit/suite references is tricky, but a decent attempt would be to drop anything including and after "#", "STE", "SUITE", "SHOP", "UNIT":

```
unit_pattern <- " (#|STE|SUITE|SHOP|UNIT).*$"
unit_examples <-
   c("1ST AVE", "RAINIER AVE S #A", "FAUNTLEROY WAY SW STE 108",
     "4TH AVE #100C", "NW 54TH ST")
str_remove(unit_examples, unit_pattern)</pre>
```

```
## [1] "1ST AVE" "RAINIER AVE S" "FAUNTLEROY WAY SW" ## [4] "4TH AVE" "NW 54TH ST"
```

#### How'd the Unit regex Work?

Breaking down " (#|STE|SUITE|SHOP|UNIT).\*\$":

- First we match a space
- (#|STE|SUITE|SHOP|UNIT) matches one of those words
- .\*\$ matches any character (.) after those words, zero or more times (\*), until the end of the string (\$)

#### Removing Units/Suites

```
## # A tibble: 11 x 1
##
     street only
##
     <chr>
##
   1 SW AVAION WAY
## 2 MERCER ST
   3 FATRVIEW AVE N
##
##
   4 1ST AVE
##
   5 E VALLEY HWY
##
   6 11TH AVE E
##
   7 112TH AVE NE
##
   8 NE BOTHELL WAY
   9 NW GILMAN BL C-08
##
## 10 NE 20TH ST
## 11 S ORCAS ST
```

For serious work, we would want to also look into special cases like "C-08" here.

#### Where Does Danger Lurk?

Let's get the number of 45+ point inspections occurring on every street.

```
restaurants %>%
  filter(Inspection_Score > 45) %>%
  distinct(Business_ID, Date, Inspection_Score, street_only) %>%
  count(street_only) %>%
  arrange(desc(n)) %>%
  head(n=5)
```

```
## # A tibble: 5 x 2
     street_only
##
                           n
     <chr>
                       <int>
##
## 1 UNIVERSITY WAY NE
                         108
## 2 S JACKSON ST
                         105
## 3 PACIFIC HWY S
                          90
## 4 NE 24TH ST
                          76
## 5 RAINIER AVE S
                          70
```

#### Splitting up Strings

You can split up strings using tidyr::separate(), seen in Week 5. Another option is str\_split(), which will split strings based on a pattern separating parts and put these components in a list. str\_split\_fixed() will do that but with a matrix instead (thus can't have varying numbers of separators):

```
head(str_split_fixed(restaurants$Violation_Description, " - ", n = 2))
```

```
## [1,]
## [1,] ""
## [2,] "4300"
## [3,] "4800"
## [5,] "4100"
## [6,] "2120"
## [7,2]
## [1,] ""
## [2,] "Non-food contact surfaces maintained and clean"
## [3,] "Physical facilities properly installed,..."
## [4,] "Proper shellstock ID; wild mushroom ID; parasite destruction procedures for fish"
## [5,] "Warewashing facilities properly installed,..."
## [6,] "Proper cold holding temperatures ( 42 degrees F to 45 degrees F)"
```

#### Making Sentences

Maybe we have a report or website where we need text dynamically generated from data.

Lets prep some recent scores first.

#### With paste()

We can give *many* arguments to string a sentence together.

```
## # A tibble: 3 x 1
## text_desc
## <chr>
## 1 Supreme Bean Again is located at 14424 Ambaum Bl Sw in Burien and r~
## 2 Mandarin Garden is located at 40 E Sunset Way in Issaquah and recei~
## 3 Flapjacks Waffle House is located at 13806 1st Ave S in Burien and ~
```

#### With glue

Or we can use str\_glue, paste()'s more sophisticated sibling which uses
the glue package. Variables and functions just go inside { } and you can
create temporary variables for convenience.

```
## # A tibble: 3 x 1
## text_desc
## <glue>
## 1 Supreme Bean Again is located at 14424 Ambaum Bl Sw in Burien and r~
## 2 Mandarin Garden is located at 40 E Sunset Way in Issaquah and recei~
## 3 Flapjacks Waffle House is located at 13806 1st Ave S in Burien and ~
```

## str\_wrap() and \n

The previous output will work fine for in-line Markdown, but it runs off the edge of the console. It also won't wrap in many tables and images.

We can add regular linebreaks using str\_wrap() or manually with "\n".

```
score_text %>%
  pull(text_desc) %>%
  str_wrap(width = 70) %>%
  paste0("\n\n") %>% # add two linebreaks as a paragraph break
  cat() # cat combines text and prints it
```

```
## Supreme Bean Again is located at 14424 Ambaum Bl Sw in Burien and
## received a score of 10 on January 24th, 2017.
##
## Mandarin Garden is located at 40 E Sunset Way in Issaquah and received
## a score of 72 on March 9th, 2007.
##
## Flapjacks Waffle House is located at 13806 1st Ave S in Burien and
## received a score of 45 on October 3rd, 2008.
```

3 / / 39

# Other Useful **stringr**Functions

str\_pad(string, width, side, pad): Adds "padding" to any string to make it a given minimum width.

str\_subset(string, pattern): Returns all elements that contain matches of the pattern.

str\_which(string, pattern): Returns numeric indices of elements that match the pattern.

str\_replace\_all(string, pattern, replacement): Performs multiple
replacements simultaneously

str\_squish(string): Trims spaces around a string but also removes
duplicate spaces inside it.

### Coming Up

Homework 6, Part 2 is due next week, and peer reviews due the week after.

UW CS&SS.