

DSCI 101 – Midterm Review

Data Wrangling – select()



Descriptions	Codes
Choose specific column name from the data frame	<pre data-bbox="1441 198 2342 309">new_data <- data %>% select(variable1, variable3)</pre>
<p>Select Columns by Name Patterns: use special helper functions like `starts_with()`, `ends_with()`, `contains()`, `matches()`, and `everything()` to select columns based on their names.</p>	
Exclude Columns: To exclude specific columns, you can use the "-" (minus) sign before the column name.	<pre data-bbox="1441 918 2342 1029">new_data <- data %>% select(-variable2,-variable4)</pre>
Select Columns by Index and Range:	<pre data-bbox="1441 1163 2048 1274">new_data <- data %>% select(c(1,3))</pre>

Data Wrangling – filter()

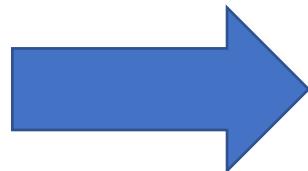
Variable 1	Variable 2	Variable 3	Variable 4
Yellow	Yellow	Yellow	Yellow
Red	Red	Red	Red
Blue	Blue	Blue	Blue
Yellow	Yellow	Yellow	Yellow
Blue	Blue	Blue	Blue
Yellow	Yellow	Yellow	Yellow
Yellow	Yellow	Yellow	Yellow
Red	Red	Red	Red
Red	Red	Red	Red
Yellow	Yellow	Yellow	Yellow



Variable 1	Variable 2	Variable 3	Variable 4
Yellow	Yellow	Yellow	Yellow
Yellow	Yellow	Yellow	Yellow
Yellow	Yellow	Yellow	Yellow
Yellow	Yellow	Yellow	Yellow
Red	Red	Red	Red
Red	Red	Red	Red
Yellow	Yellow	Yellow	Yellow

Descriptions	Codes
Simple Conditions - greater than (`>`), less than (`<`), or equal to (`==`)	<pre data-bbox="1479 195 2186 303">new_data <- data %>% filter(variable1 > 30)</pre> <pre data-bbox="1479 368 2314 476">new_data <- data %>% filter(variable2 == "yes")</pre>
Multiple Conditions - You can combine conditions using logical operators like `&` (AND) and ` ` (OR).	<pre data-bbox="1479 508 2237 663">new_data <- data %>% filter(variable1 > 30 & variable2 < 50000)</pre> <pre data-bbox="1479 728 2222 884">new_data <- data %>% filter(variable1 > 30 variable2 < 50000)</pre>
Exclusion - to exclude certain rows, you can use the `!=` operator (not equal to).	<pre data-bbox="1479 964 2222 1077">new_data <- data %>% filter(variable4 != 30)</pre>
Filter rows based on vector of conditions - The `%in%` operator is useful for filtering rows with values in a specified vector.	<pre data-bbox="1479 1170 2186 1326">new_data <- data %>% filter(variable3 %in% c("bananas", "grapes"))</pre>

Data Wrangling – mutate()

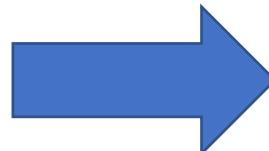


Descriptions	Codes
Create a new variable calculated from another variable	<pre>new_data <- data %>% mutate(variable5 = variable2 / 12)</pre>
We can use other functions within mutate to help us make the new variable like `ifelse()`.	<pre>new_data <- data %>% mutate(variable5 = ifelse(variable1 > 30, "old", "Young"))</pre>
If you have multiple condition you can use the `case_when()` function and list out your possible options.	<pre>new_data <- data %>% mutate(variable5 = case_when(variable1 < 30 ~ "Young", variable1 >= 30 & variable1 <= 50 ~ "Middle-aged", variable1 > 50 ~ "old", TRUE ~ "No Category"))</pre>

Data Wrangling – summarise() & group_by

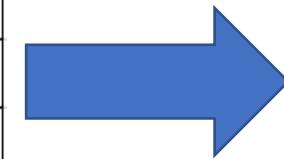
Variable 1	Variable 2	Variable 3	Variable 4
Yellow	Yellow		
Yellow	Yellow		
Blue	Blue		
Blue	Blue		
Red	Red		
Yellow	Yellow		
Red	Red		
Red	Red		
Yellow	Yellow		
Red	Red		

group_by()



Variable 1	Variable 2	Variable 3	Variable 4
Yellow	Yellow		
Yellow	Yellow		
Blue	Blue		
Blue	Blue		
Red	Red		
Yellow	Yellow		
Red	Red		
Red	Red		
Yellow	Yellow		
Red	Red		

summarise()

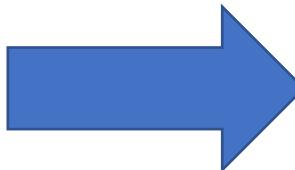


Variable 1	Variable 2	NEW VARIABLE
Yellow	Yellow	
Blue	Blue	
Red	Red	

Descriptions	Codes
Compute different statistics (mean, median..etc) grouped by different variables.	<pre data-bbox="1141 321 2396 548">new_data <- data %>% group_by(variable1, variable2) %>% summarise(NEW_VARIABLE1 = mean(variable1), NEW_VARIABLE2 = median(variable2))</pre>
If you do not group_by, it will make a computation on the entire column:	<pre data-bbox="1141 781 2396 951">new_data <- data %>% summarise(NEW_VARIABLE1 = mean(variable1), NEW_VARIABLE2 = median(variable2))</pre>

Data Wrangling – arrange()

Variable 1	Variable 2	Variable 3	Variable 4
Light Blue	Light Blue	Light Blue	Light Blue
Dark Blue	Dark Blue	Dark Blue	Dark Blue
Light Gray	Light Gray	Light Gray	Light Gray
Light Blue	Light Blue	Light Blue	Light Blue
Dark Blue	Dark Blue	Dark Blue	Dark Blue

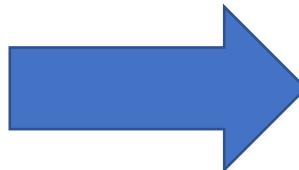


Variable 1	Variable 2	Variable 3	Variable 4
Light Blue	Light Blue	Light Blue	Light Blue
Light Gray	Light Gray	Light Gray	Light Gray
Light Blue	Light Blue	Light Blue	Light Blue
Dark Blue	Dark Blue	Dark Blue	Dark Blue

Descriptions	Codes
Arrange in ascending order (smallest to largest)	<code>new_data <- data %>% arrange(variable1)</code>
Arrange in descending order (largest to smallest)	<code>new_data <- data %>% arrange(desc(variable1))</code>

Data Wrangling – count()

Variable 1	Variable 2	Variable 3	Variable 4
Red			
Yellow			
Red			
Yellow			
Yellow			

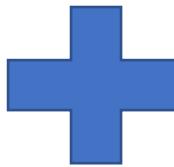


Variable 1	n
Red	2
Yellow	3

Descriptions	Codes
Count the occurrences of each "category"	<pre data-bbox="1121 198 1710 303">new_data <- data %>% count(variable1)</pre>
Can achieve the same thing as count using summarize. With summarize you can name the new variable whatever you want. With count, it automatically calls it “n”.	<pre data-bbox="1121 659 1966 822">new_data <- data %>% group_by(variable1) %>% summarise(NEW_VARIABLE = n())</pre>

Joining Datasets

Variable 1	Variable 2	Variable 3	Variable 4
Red	Light Blue	Light Green	Light Orange
Red	Light Blue	Light Green	Light Orange
Red	Light Blue	Light Green	Light Orange
Red	Light Blue	Light Green	Light Orange



Variable 1	Variable 5	Variable 6	Variable 7
Red	Cyan	Yellow	Light Green
Red	Cyan	Yellow	Light Green
Red	Cyan	Yellow	Light Green
Red	Cyan	Yellow	Light Green



Variable 1	Variable 2	Variable 3	Variable 4	Variable 5	Variable 6	Variable 7
Red	Light Blue	Light Green	Light Orange	Cyan	Yellow	Light Green
Red	Light Blue	Light Green	Light Orange	Cyan	Yellow	Light Green
Red	Light Blue	Light Green	Light Orange	Cyan	Yellow	Light Green
Red	Light Blue	Light Green	Light Orange	Cyan	Yellow	Light Green

Joining Datasets – inner_join()

Variable 1	Variable 2	Variable 3	Variable 4
#			
*			
&			
@			
%			



Variable 1	Variable 5	Variable 6	Variable 7
#			
*			
&			
^			
!			

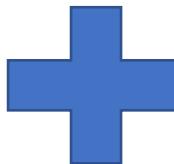


Variable 1	Variable 2	Variable 3	Variable 4	Variable 5	Variable 6	Variable 7
#						
*						
&						

```
new_data <- left_data %>%
  inner_join(right_data, by = c("variable1" = "variable1"))
```

Joining Datasets – left_join()

Variable 1	Variable 2	Variable 3	Variable 4
#			
*			
&			
@			
%			



Variable 1	Variable 5	Variable 6	Variable 7
#			
*			
&			
^			
!			



Variable 1	Variable 2	Variable 3	Variable 4	Variable 5	Variable 6	Variable 7
#						
*						
&						
@				NA	NA	NA
%				NA	NA	NA

```
new_data <- left_data %>%
  left_join(right_data, by = c("variable1" = "variable1"))
```

Joining Datasets – right_join()

Variable 1	Variable 2	Variable 3	Variable 4
#			
*			
&			
@			
%			



Variable 1	Variable 5	Variable 6	Variable 7
#			
*			
&			
^			
!			

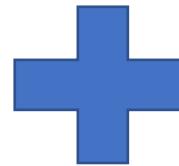


Variable 1	Variable 2	Variable 3	Variable 4	Variable 5	Variable 6	Variable 7
#						
*						
&						
^	NA	NA	NA			
!	NA	NA	NA			

```
new_data <- left_data %>%
  right_join(right_data, by = c("variable1" = "variable1"))
```

Joining Datasets – full_join()

Variable 1	Variable 2	Variable 3	Variable 4
#			
*			
&			
@			
%			



Variable 1	Variable 5	Variable 6	Variable 7
#			
*			
&			
^			
!			



Variable 1	Variable 2	Variable 3	Variable 4	Variable 5	Variable 6	Variable 7
#						
*						
&						
@				NA	NA	NA
%				NA	NA	NA
^	NA	NA	NA			
!	NA	NA	NA			

```
new_data <- left_data %>%
  full_join(right_data, by = c("variable1" = "variable1"))
```

Pivot – pivot_longer() & pivot_wider

Pivot_Variable	Variable1	Variable2	Variable3

pivot_longer()

pivot_wider()

Pivot_Variable	NEW_VARIABLE	NEW_VARIABLE2
	Variable1	
	Variable1	
	Variable1	
	Variable2	
	Variable2	
	Variable2	
	Variable3	
	Variable3	
	Variable3	

```
new_data <- data %>%  
pivot_longer(cols = -Pivot_Variable,  
names_to = "NEW_VARIABLE",  
values_to = "NEW_VARIABLE2")
```

```
new_data <- data %>%  
pivot_wider(names_from = NEW_VARIABLE,  
values_from = NEW_VARIABLE2)
```

Data Visualization – ggplot()

- **Data:** The dataset you're working with.
- **Aesthetics Mapping (aes):** How data variables map to plot aesthetics like position, color, shape, etc.

x = variable1

y = variable2

color = variable3

fill = variable4

shape = variable5

Data Visualization – ggplot()

- **Geometric Objects (geom):** The visual elements to represent the data (points, lines, bars, etc.).

`geom_histogram()`

`geom_density()`

`geom_boxplot()`

`geom_line()`

`geom_bar()`

`geom_point()`

Data Visualization – ggplot()

- **Facets (facet_wrap or facet_grid):** Splitting data into subplots based on a variable.

`facet_wrap(~ Variable1)`

`facet_grid(Variable1 ~ Variable2)`

Data Visualization – ggplot()

- **Theme:** Controlling the overall appearance of the plot.

`theme_gray()`

`theme_bw()`

`theme_minimal()`

`theme_void()`

Data Visualization – ggplot()

- **Add ons**

coord_flip()

geom_smooth(method = "lm", se = FALSE)

geom_smooth(method = "loess", se = TRUE)

scale_fill_manual(values = c("red", "green", "blue"))

scale_color_manual(values = c("red", "green", "blue"))

labs()

For loops & If/Else Statements

```
x<-NA  
y<-NA  
  
for (i in 1:5) {  
  y[i] <- 1+i  
  
  if (y[i]>3) {  
    x[i] <- "Yes"  
  }else if (y[i]<3){  
    x[i] <- "No"  
  }else{  
    x[i]<- "Maybe"  
  }  
}  
y  
x
```

Functions

```
circumference <- function(r){  
  c = 2*pi*r  
  return(c)  
}
```

```
circumference(5)  
circumference(r = 5)
```