

# Processing Large Amounts of Data in R (AKA "The Tidyverse Tutorial")

By Travis Oishi

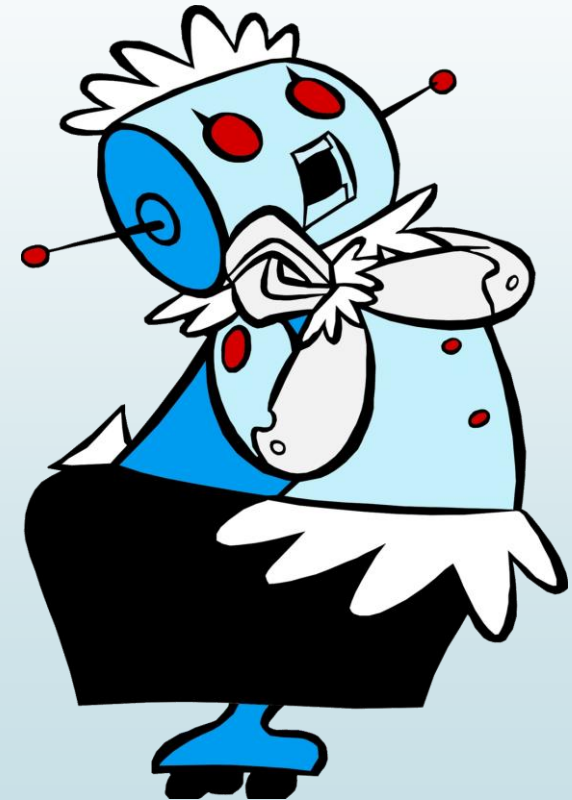
# About Me

- Central Shenandoah Health District Epidemiologist
- Jr Epi for Central Shenandoah from Oct 2019 - March 2022
- Limited experience in R when I first started in 2019
- Started coding around Spring of 2020 to improve epi processes
- [travis.oishi@vdh.virginia.gov](mailto:travis.oishi@vdh.virginia.gov)



# Benefits of R

- Completing tasks that you need to do more than once
- Processing complex equations
- Navigating large datasets
- Large scale data manipulation
- Less prone to human error



# What is Tidyverse?

- Collection of packages designed for data science
- Quality of life improvements
- Packages complement one another




# dplyr

- Data manipulation
- List of verbs that help you solve the most common data manipulation challenges
- Simplifies and streamlines codes
- Uses piping, `%>%` to create a sequence of actions that flow together



pixtastock.com – 44869386

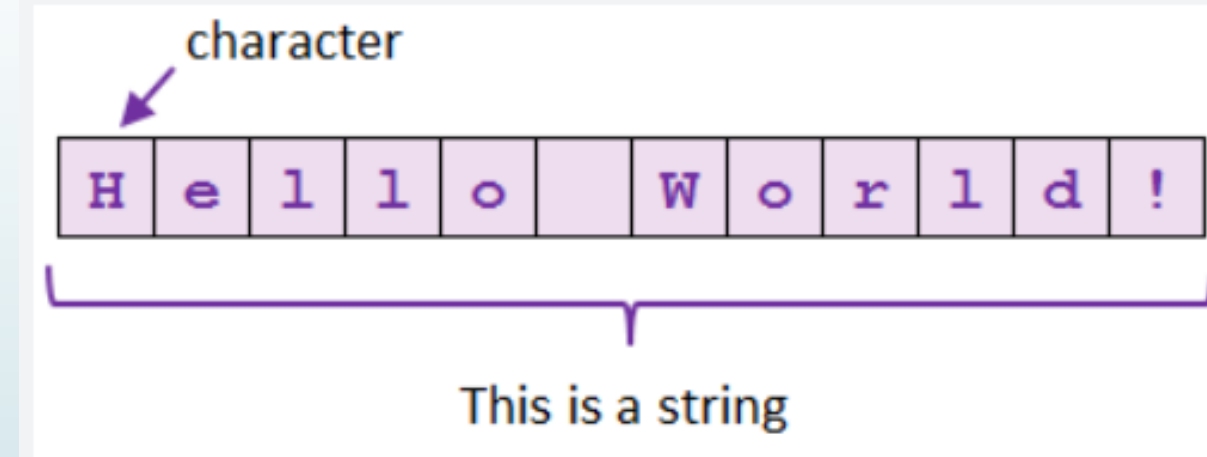




Operation	base R example	dplyr function	dplyr example
select some rows	<code>my_data[c(2,3,10),]</code>	<code>slice()</code>	<code>slice(my_data, c(2,3,10))</code>
select some columns	<code>my_data[, 1:2]</code> OR <code>my_data[, c("Var_1", "Var_2")]</code>	<code>select()</code>	<code>select(my_data, Var_1, Var_2)</code>
subset	<code>my_data[my_data\$Var_2&gt;80,]</code> OR <code>subset(my_data, Var_2&gt;80)</code>	<code>filter()</code>	<code>filter(my_data, Var_2&gt;80)</code>
order the rows	<code>my_data[order(my_data\$Var_2),]</code>	<code>arrange()</code>	<code>arrange(my_data, Var_2)</code>
add a column	<code>my_data\$logVar_2 &lt;- log(my_data\$Var_2)</code> OR <code>transform(my_data, logVar_2=log(Var_2))</code>	<code>mutate()</code>	<code>mutate(my_data, logVar_2 = log(Var_2))</code>
define groups of data	Done within other functions.	<code>group_by()</code>	<code>my_data %&gt;% group_by(Var_3)</code>
summarise the data	<code>aggregate(Var_2 ~ Var_3, data = my_data, FUN = mean)</code> OR <code>tapply(my_data\$Var_2, list(my_data\$Var_3), mean)</code>	<code>summarise()</code> AND <code>group_by()</code>	<code>my_data %&gt;% group_by(Var_3) %&gt;% meanVar_2 = mean(Var_2)</code>

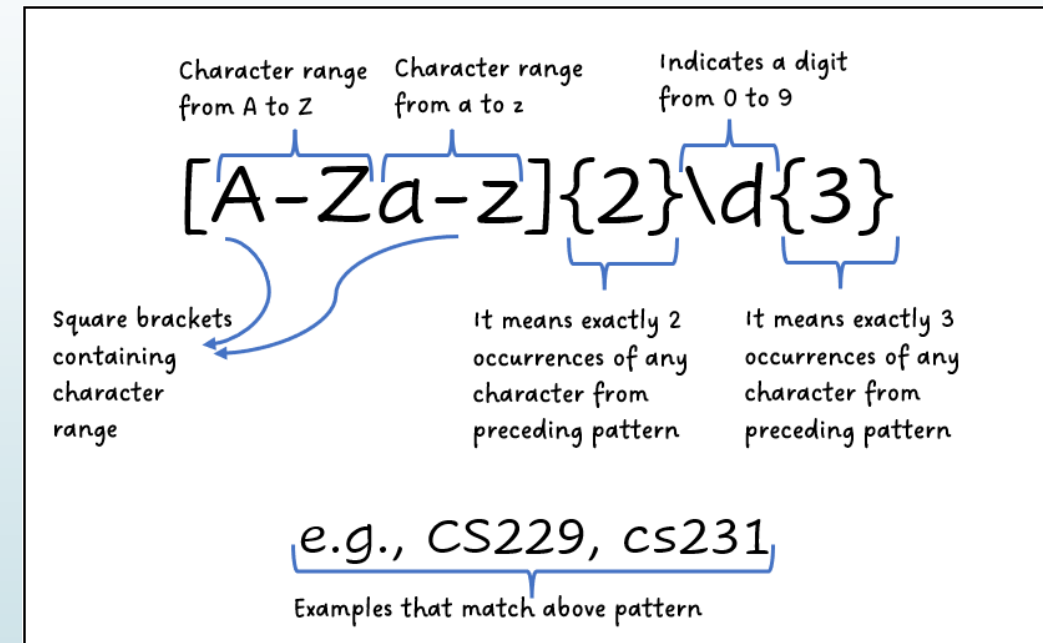
# stringr

- String identification and manipulation
- Pattern recognition
- Data cleaning
- Use of regular expressions



# Identifying Patterns in Strings: Regular Expressions

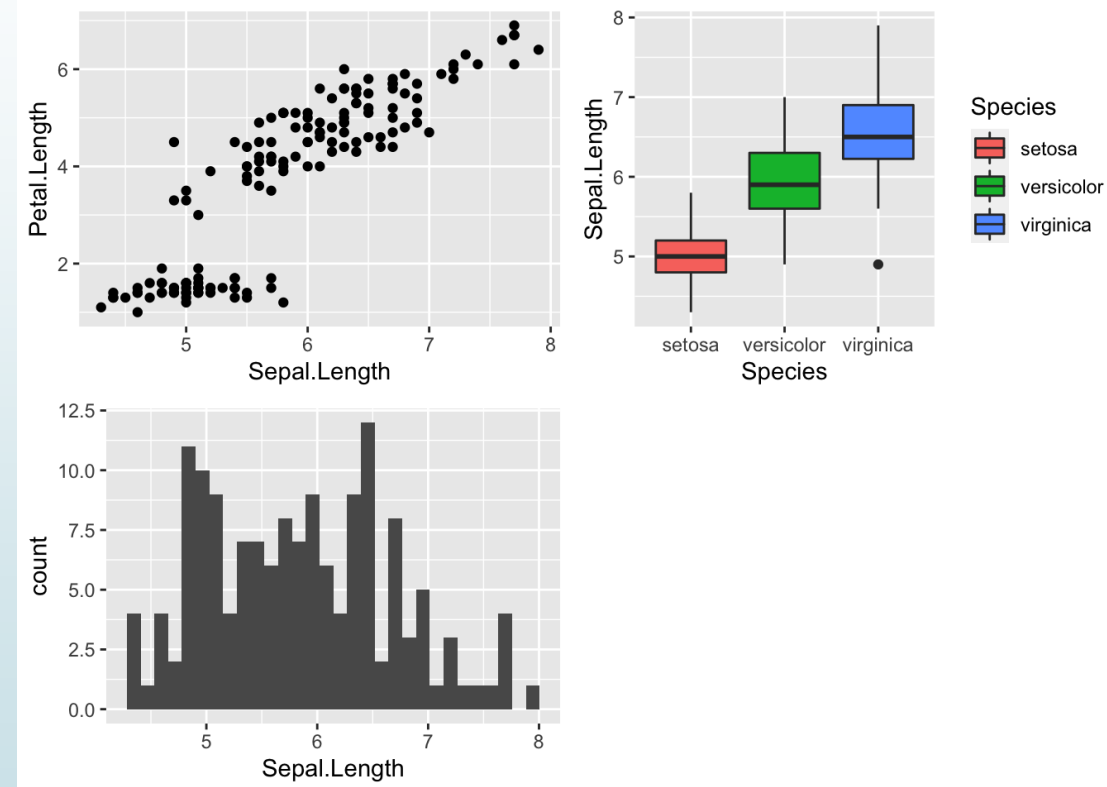
- Method for the computer to process text
- Identifies specific patterns
- Practical uses:
  - Address processing
  - Identifying emails





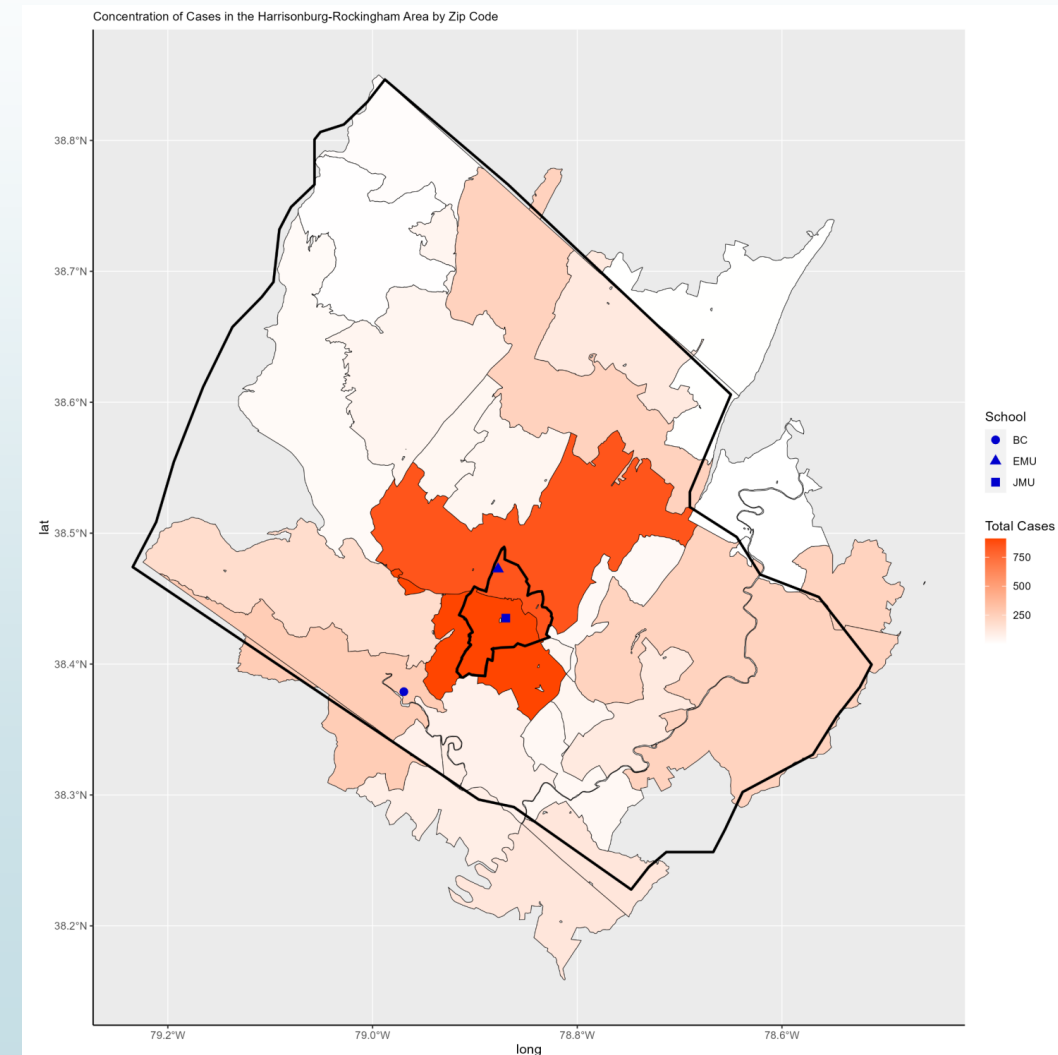
# ggplot2

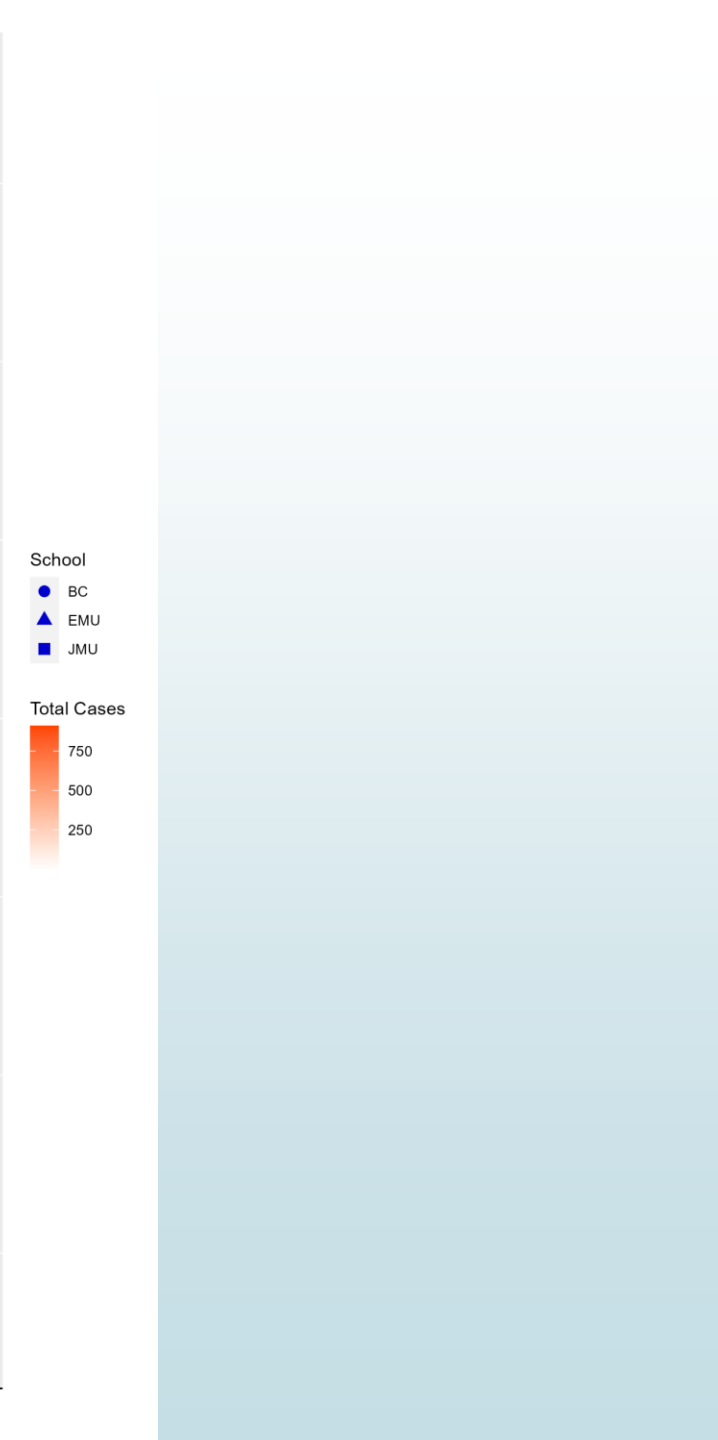
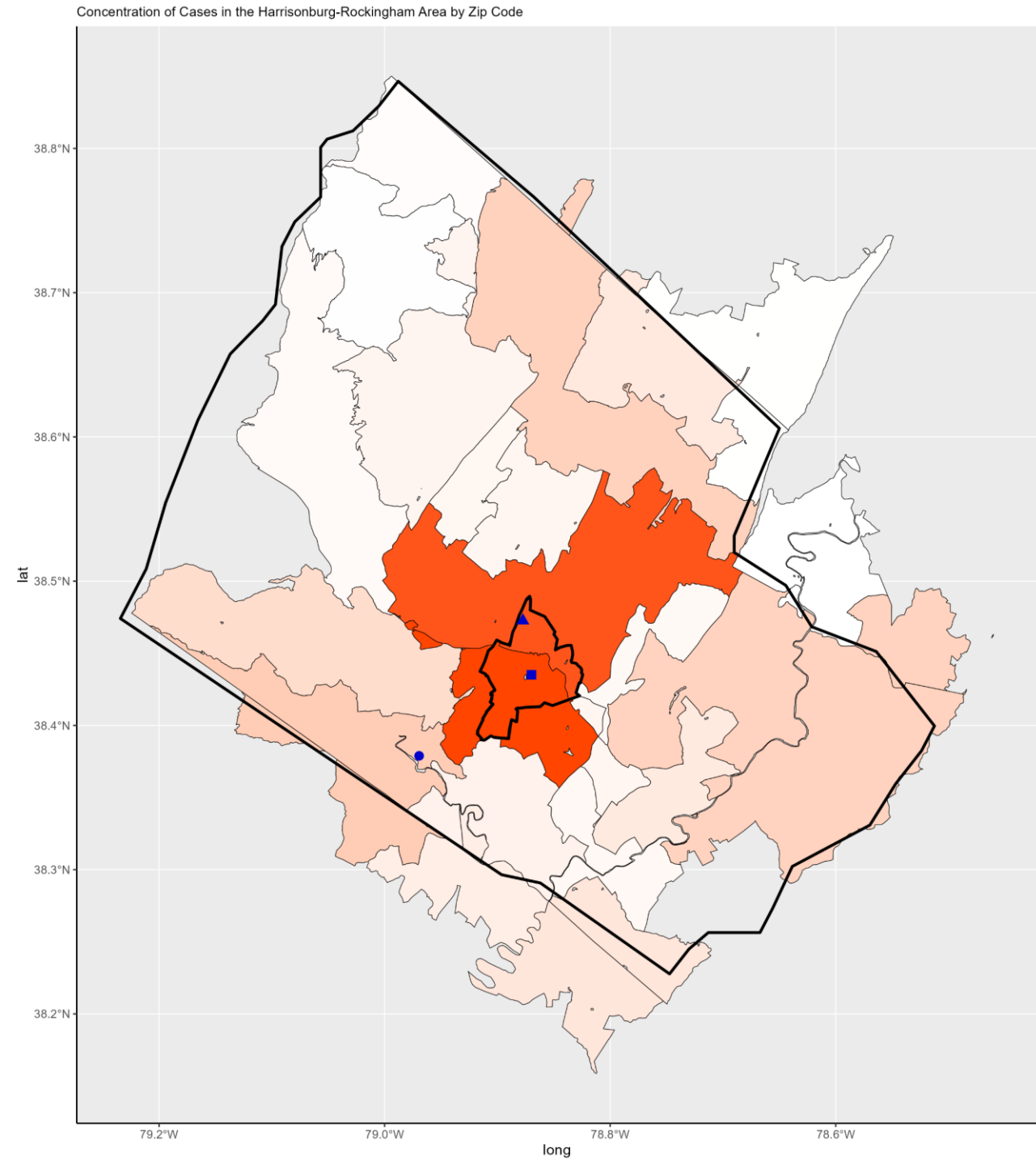
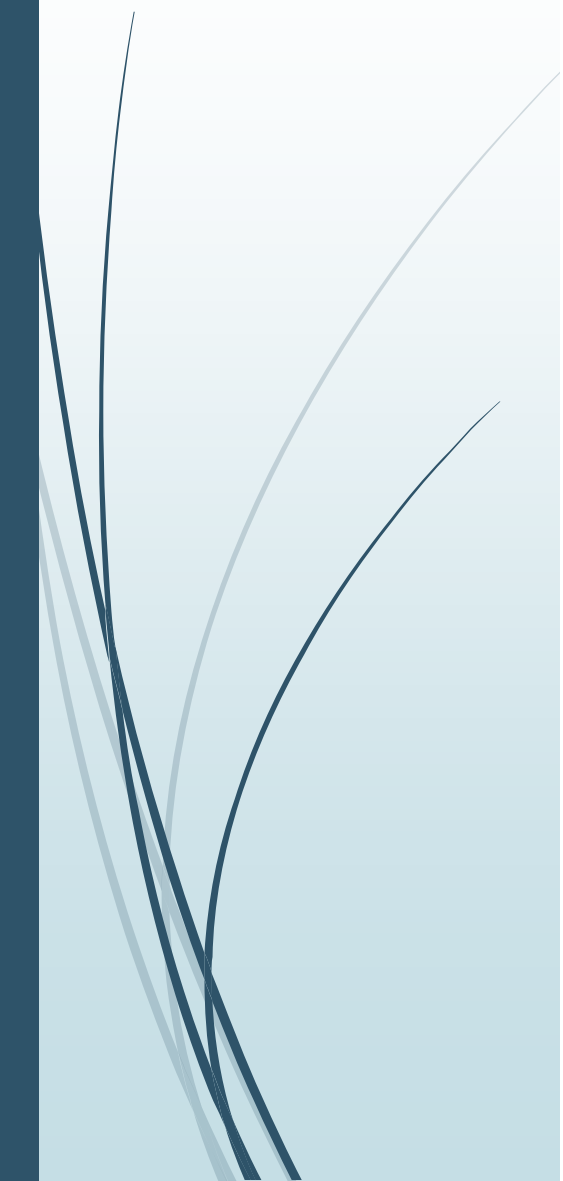
- Data visualization
- Control the aesthetics of plots
- Create a variety of plots based on one or more data sources



# Tidyverse in Action: Mapping Concentration of Disease Cases

- Identify complete zip codes reported in Rockingham County case reports (**stringr**)
- Summarize the total cases in each zip code (**dplyr**)
- Match zip code data to a zip code map shape file (**dplyr**)
- Plot the map (**ggplot2**)





# Type in the Chat

- 1: You are familiar with story of *The Odyssey*
- 0: You are not familiar with the story of *The Odyssey*





# Learn More

- <https://study.com/learn/lesson/tidyverse-packages-examples-r-programming.html#:~:text=Tidyverse%20is%20an%20R%20programming,continually%20being%20modified%20and%20improved.>
- <https://www.tidyverse.org/>