

# Data Science Basics in R

Day 3: Exploratory data analysis

# Goals for today

- Define descriptive statistics & exploratory data analysis
- Create your first data visualization in R
- Identify options for visualization in R, including ggplot2
- Get creative and have fun exploring datasets

# Descriptive statistics

# Goals for today

- Define descriptive statistics & exploratory data analysis
- Create your first data visualization in R
- Identify options for visualization in R, including ggplot2
- Get creative and have fun exploring datasets

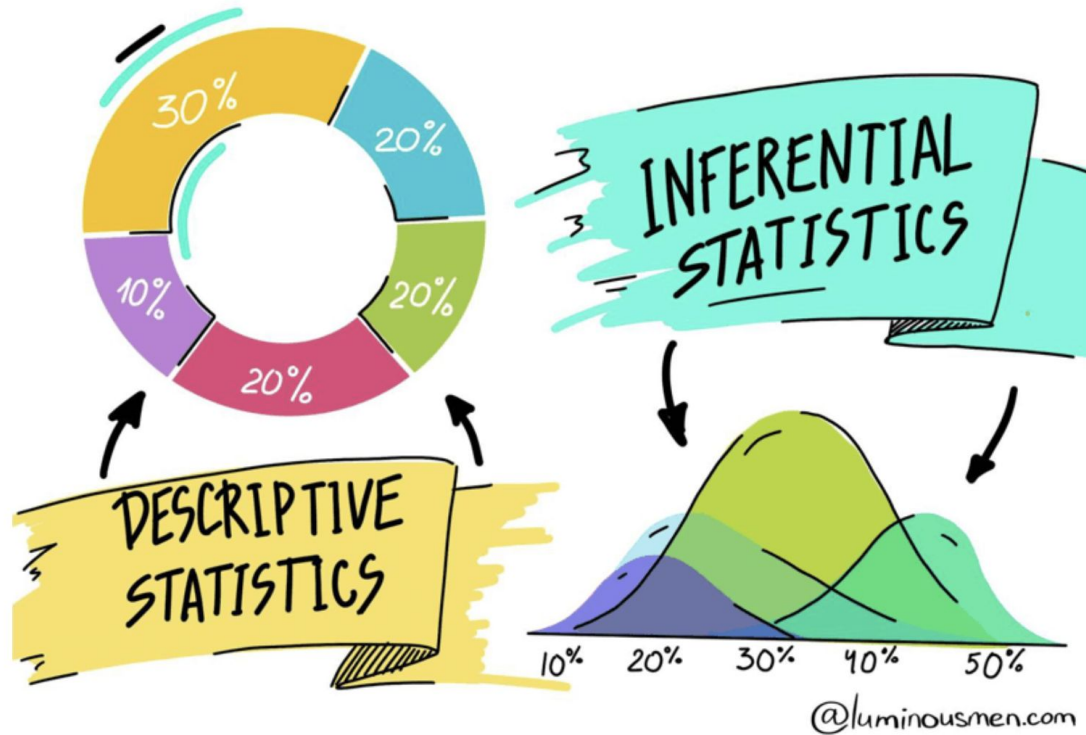
# Descriptive statistics

Descriptive statistics summarize data, and typically describe three types of things:

- center (e.g., mean, median)
- spread (e.g., standard deviation, interquartile range)
- counts & rates (e.g., summary tables)

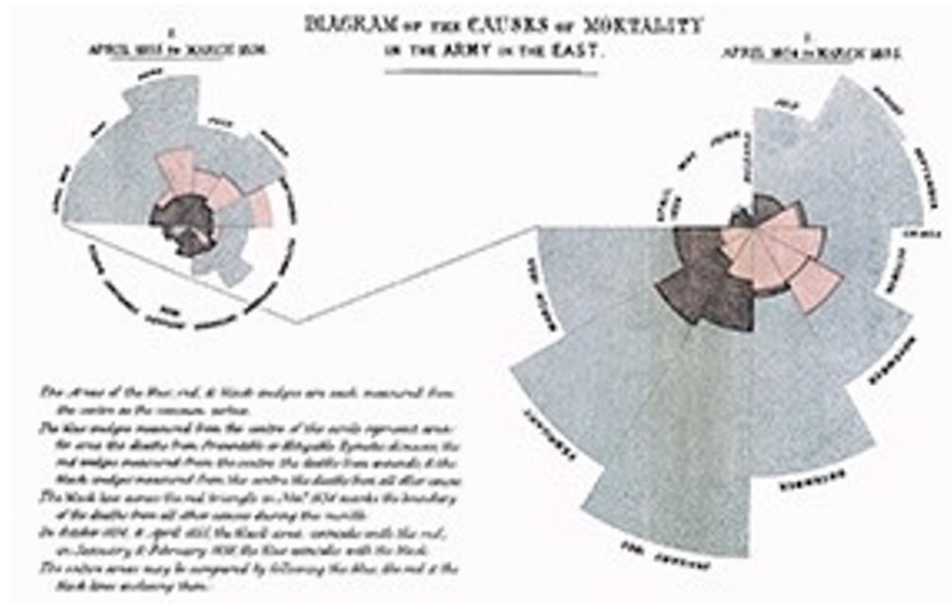
In a typical data analysis workflow, we explore these first! It's helpful to better understand your data, and to identify potential surprises.

# Descriptive statistics



# Goals of data visualization

# Florence Nightingale's polar area chart, "Diagram of the causes of mortality in the army in the East".



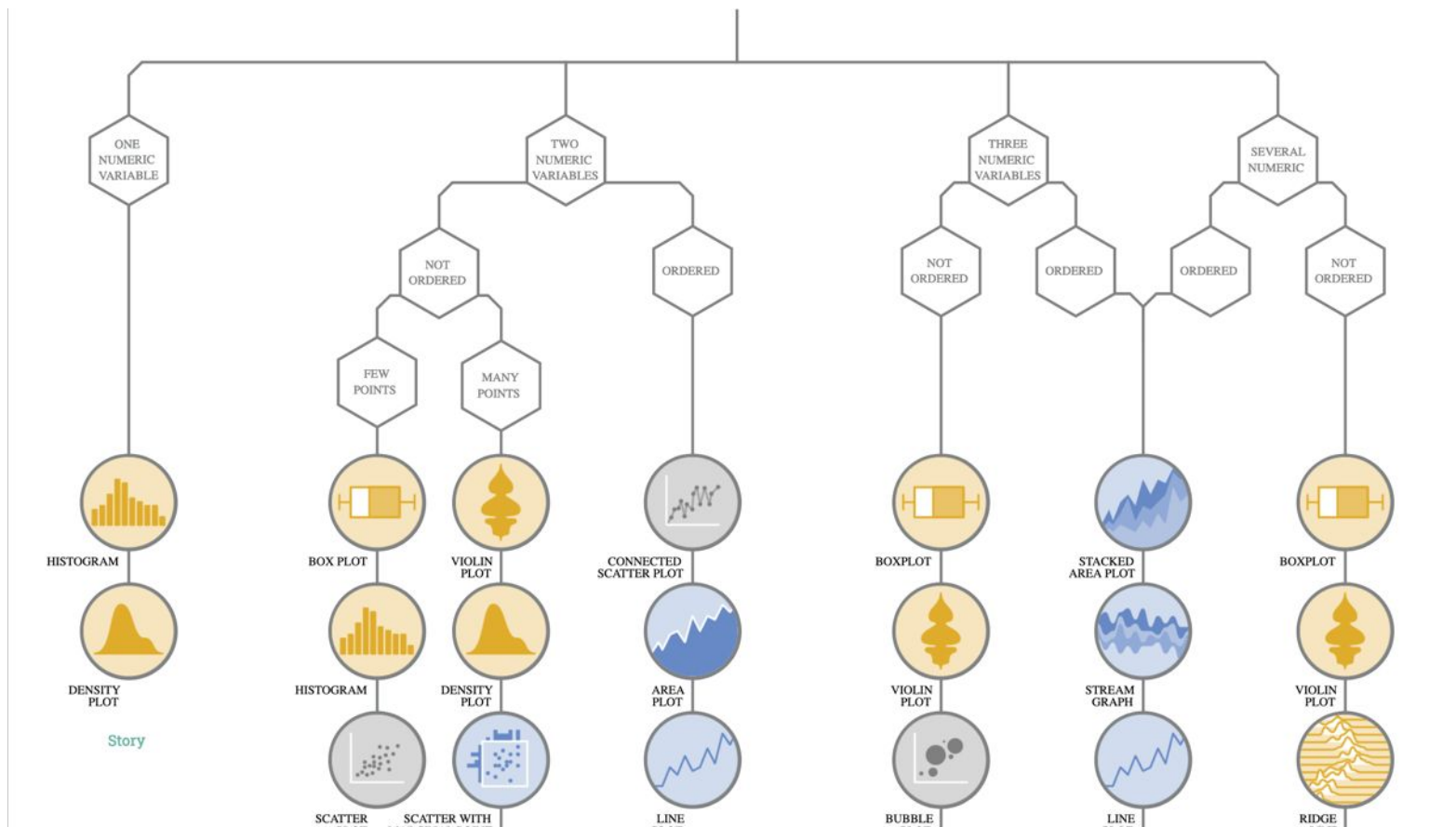


# Goals of data visualization

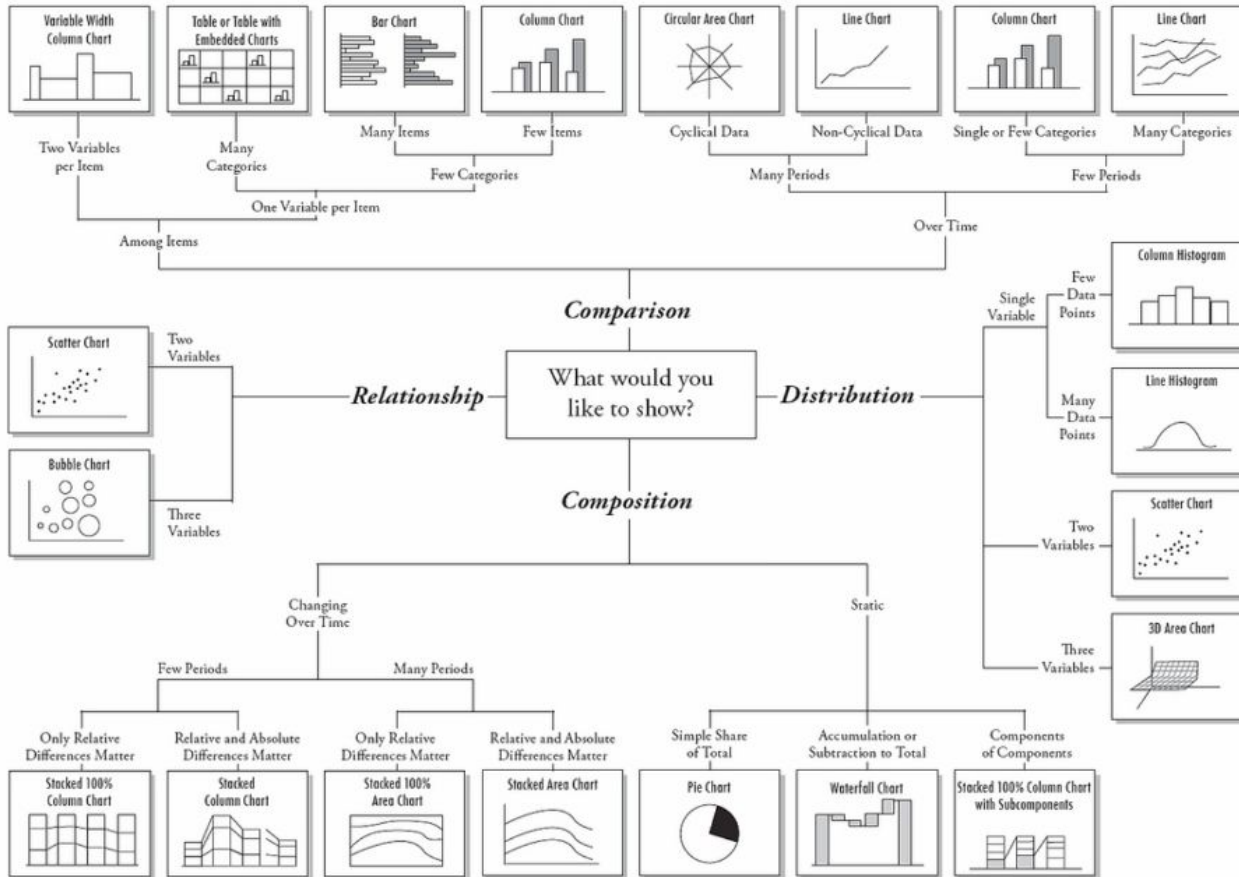
**What makes a data visualization *good*?**

# Choosing a data visualization

(we'll talk more about this tomorrow)



# Chart Suggestions—A Thought-Starter



Graphic by Andrew Abela

[https://extremepresentation.typepad.com/blog/2006/09/choosing\\_a\\_good.html](https://extremepresentation.typepad.com/blog/2006/09/choosing_a_good.html)

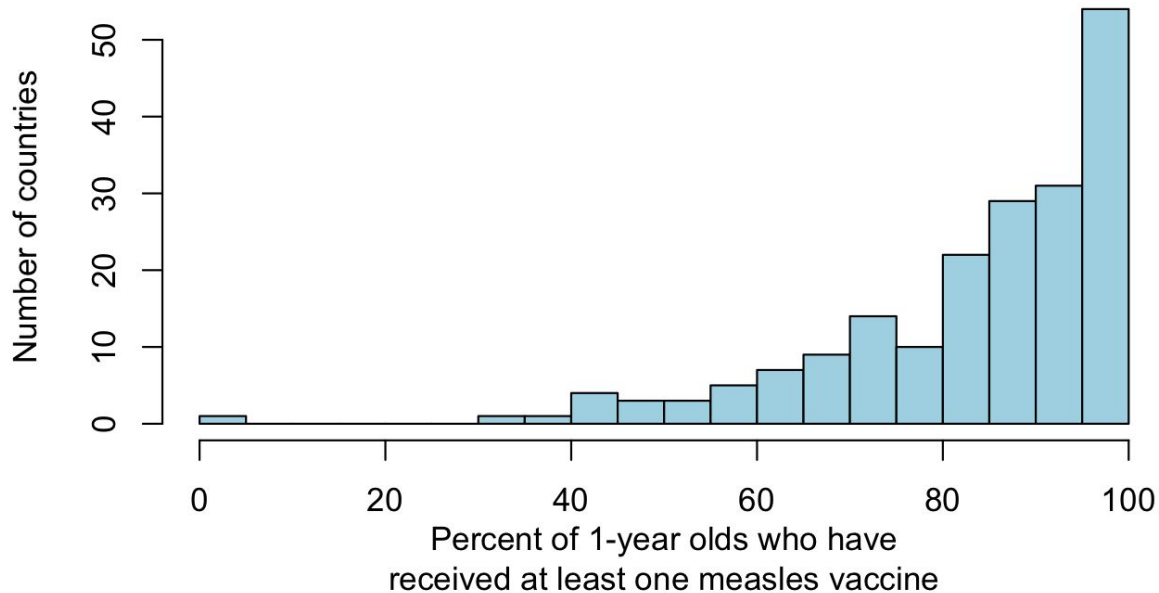
© 2006 A. Abela — a.vabela@gmail.com

# Plots in base R

*(code in github for live demo)*

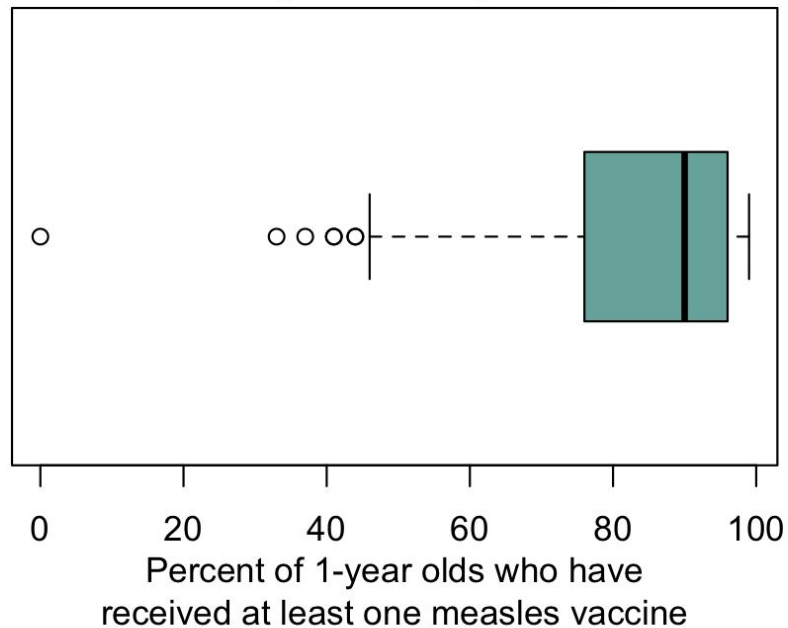
# Histogram

**Distribution of country-level  
measles vaccination rates  
for 1-year olds (MCV1)**



# Boxplot

**Distribution of country-level  
measles vaccination rates  
for 1-year olds (MCV1)**

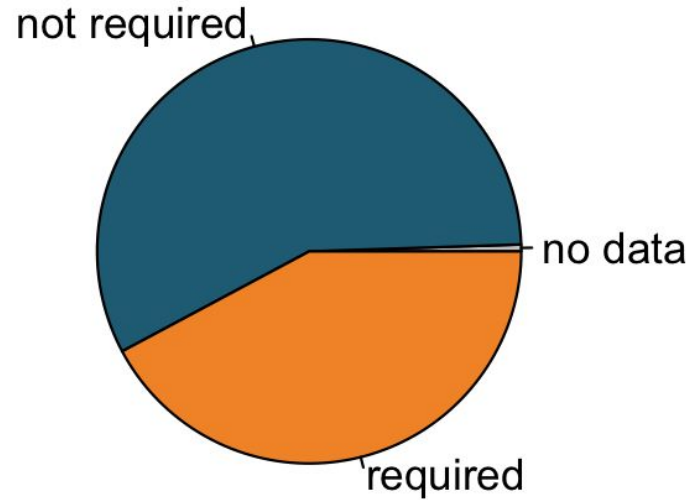




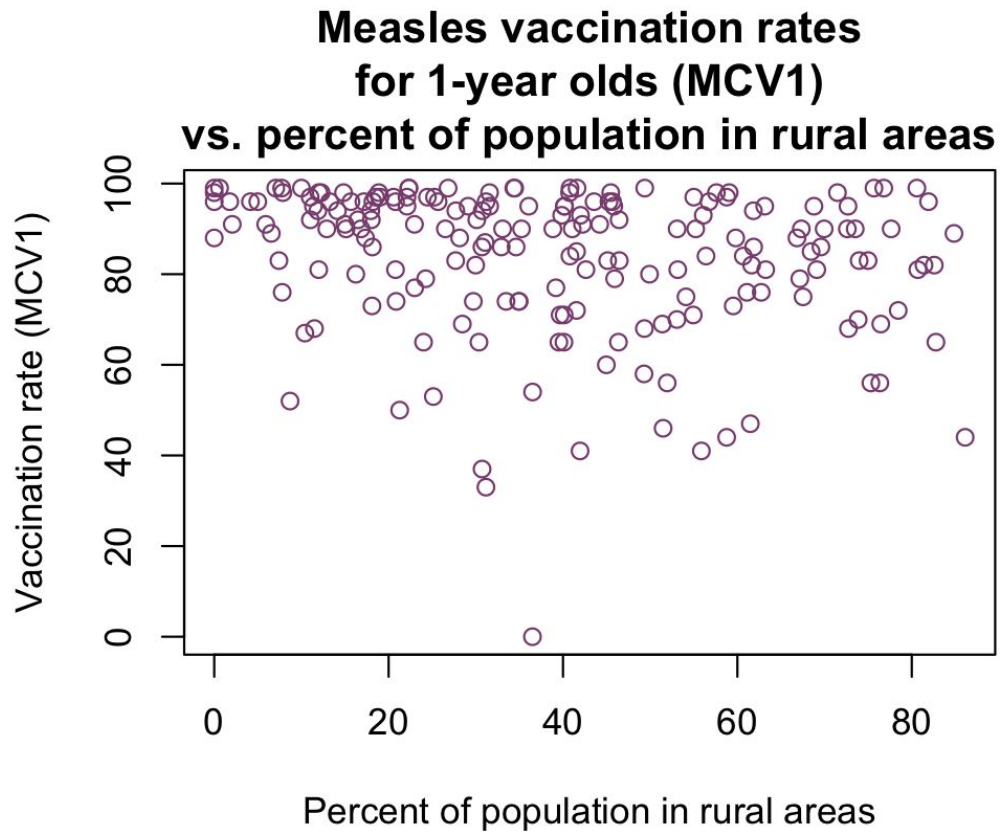
# Barchart

# Pie chart

## Policy requirement for measles vaccination



# Scatterplot



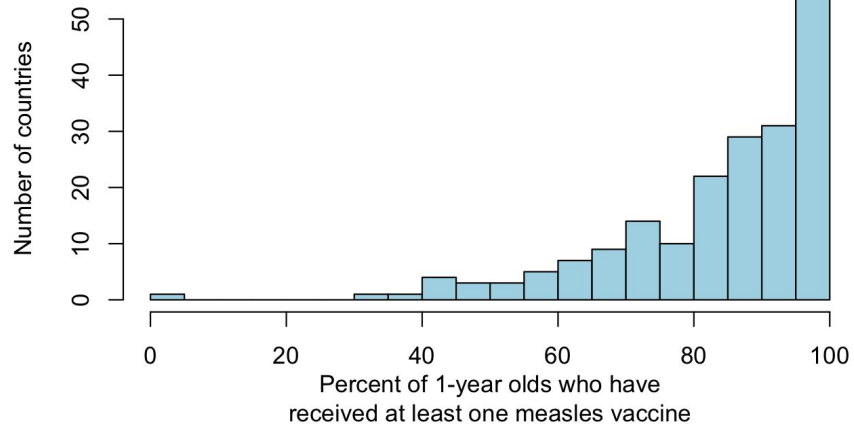
**10 minute break**

# Plots in ggplot2

*(code in github for live demo)*

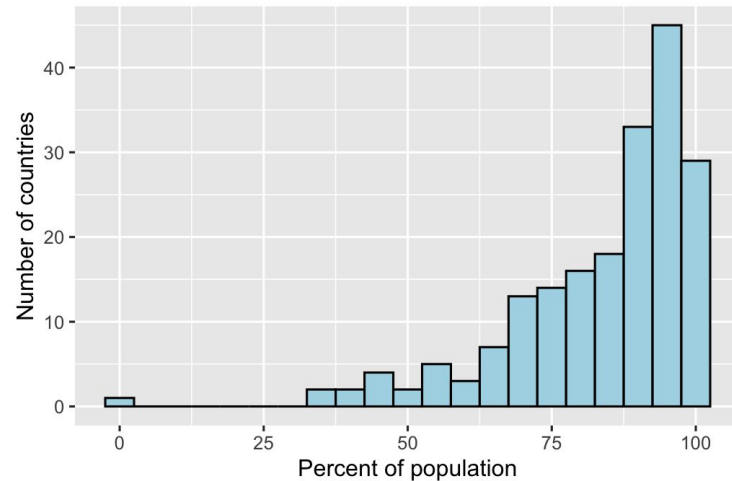
# Histogram

**Distribution of country-level  
measles vaccination rates  
for 1-year olds (MCV1)**



Base R

**Distribution of country-level measles vaccination rates  
for 1-year olds (MCV1)**

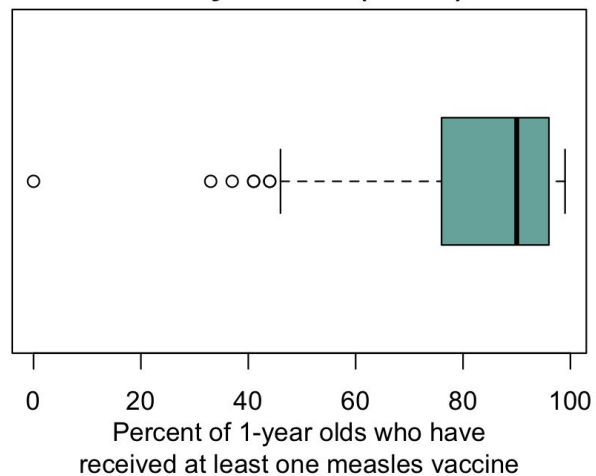


MCV1 is defined as the percentage of children under one year of age who have received at least one dose of measles-containing vaccine in a given year.

ggplot

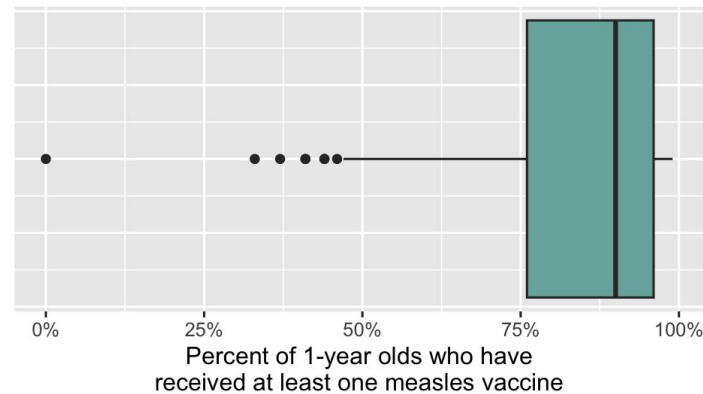
# Boxplot

**Distribution of country-level  
measles vaccination rates  
for 1-year olds (MCV1)**



Base R

Distribution of country-level measles vaccination rates  
for 1-year olds (MCV1)



ggplot

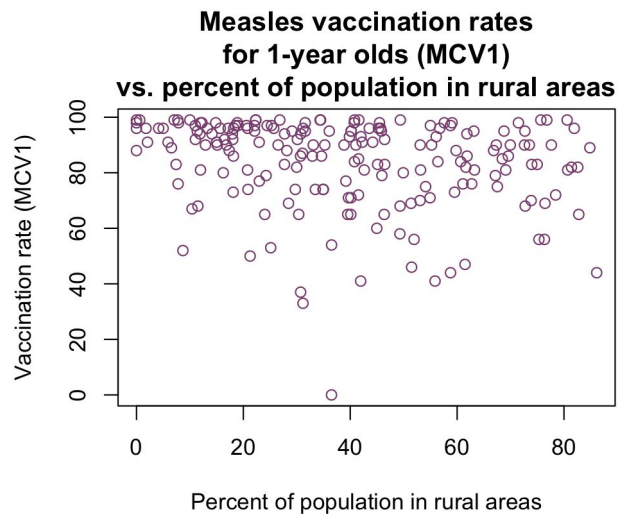
# Bar chart

Base R

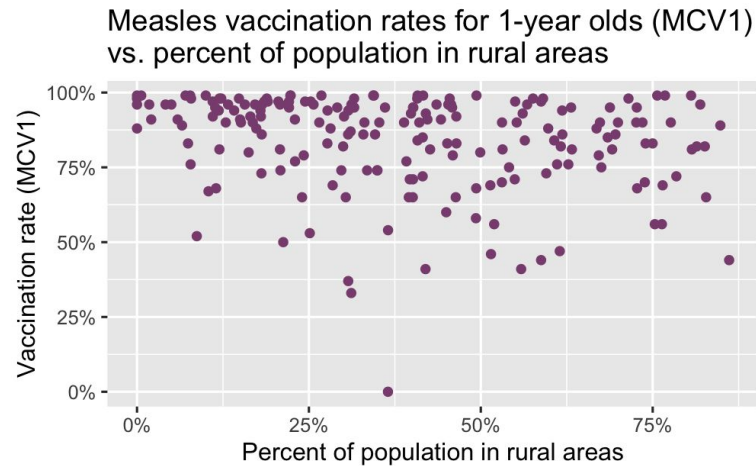
ggplot



# Scatterplot



Base R



ggplot

**Thank you!**

**See you tomorrow.**

***Please come with a fully charged laptop.***