# R User Group

## Building a community of R users in Connecticut state government

Launched June 2024

Data & Policy Analytics

CONNECTICUT
Policy and Management

# **Table of contents**

# 2024 meeting schedule

**Monthly meetings**: second Tuesday of each month

- July 9
- August 13
- September 10
- October 8
- November 12
- December 10

# R resources

# R resources: Getting started

1. [Installing R and R Studio](#)
   1. [Download R](#)
   2. [Download R Studio](#)
2. [R Packages ](#)(what's an R package and how do I install one?)
3. [Short primers ](#)from the owners of R Studio

# R resources: Great references

1. [A Gentle Introduction to Tidy Statistics in R](#) – Introductory tutorial focusing on stats

2. [R for Data Science](#) – Free online textbook introducing R for data organization, analysis, manipulation, and visualization

3. [Cheat sheets](#) for popular R packages

# Data visualization with ggplot2 :: CHEATSHEET

## Basics

**ggplot2** is based on the **grammar of graphics**, the idea that you can build every graph from the same components: a **data** set, a **coordinate system**, and **geoms**—visual marks that represent data points.

data + geom (x = F y = A) | coordinate system = plot

To display values, map variables in the data to visual properties of the geom (**aesthetics**) like **size**, **color**, and **x** and **y** locations.

data + geom (x = F y = A color = F size = A) | coordinate system = plot

Complete the template below to build a graph.

```
ggplot (data = <DATA>) +                          required
  <GEOM_FUNCTION>(mapping = aes <MAPPINGS> ),
  stat = <STAT>, position = <POSITION>) +         Not
  <COORDINATE_FUNCTION> +                         required,
  <FACET_FUNCTION> +                              sensible
  <SCALE_FUNCTION> +                              defaults
  <THEME_FUNCTION>                                supplied
```

**ggplot**(data = mpg, **aes**(x = cty, y = hwy)) Begins a plot that you finish by adding layers to. Add one geom function per layer.

**last_plot()** Returns the last plot.

**ggsave**("plot.png", width = 5, height = 5) Saves last plot as 5' x 5' file named "plot.png" in working directory. Matches file type to file extension.

## Aes  Common aesthetic values.

**color** and **fill** - string ("red", "#RRGGBB")

**linetype** - integer or string (0 = "blank", 1 = "solid", 2 = "dashed", 3 = "dotted", 4 = "dotdash", 5 = "longdash", 6 = "twodash")

**size** - integer (in mm for size of points and text)

**linewidth** - integer (in mm for widths of lines)

**shape** - integer/shape name or a single character ("a")

## Geoms
Use a geom function to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

### GRAPHICAL PRIMITIVES
a <- ggplot(economics, aes(date, unemploy))
b <- ggplot(seals, aes(x = long, y = lat))

**a + geom_blank()** and **a + expand_limits()**
Ensure limits include values across all plots.

**b + geom_curve**(aes(yend = lat + 1, xend = long + 1), curvature = 1) - x, xend, y, yend, alpha, angle, color, curvature, linetype, size

**a + geom_path**(lineend = "butt", linejoin = "round", linemitre = 1)
x, y, alpha, color, group, linetype, size

**a + geom_polygon**(aes(alpha = 50)) - x, y, alpha, color, fill, group, subgroup, linetype, size

**b + geom_rect**(aes(xmin = long, ymin = lat, xmax = long + 1, ymax = lat + 1)) - xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size

**a + geom_ribbon**(aes(ymin = unemploy - 900, ymax = unemploy + 900)) - x, ymax, ymin, alpha, color, fill, group, linetype, size

### LINE SEGMENTS
common aesthetics: x, y, alpha, color, linetype, size

**b + geom_abline**(aes(intercept = 0, slope = 1))
**b + geom_hline**(aes(yintercept = lat))
**b + geom_vline**(aes(xintercept = long))

**b + geom_segment**(aes(yend = lat + 1, xend = long + 1))
**b + geom_spoke**(aes(angle = 1:1155, radius = 1))

### ONE VARIABLE    continuous
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)

**c + geom_area**(stat = "bin")
x, y, alpha, color, fill, linetype, size

**c + geom_density**(kernel = "gaussian")
x, y, alpha, color, fill, group, linetype, size, weight

**c + geom_dotplot()**
x, y, alpha, color, fill

**c + geom_freqpoly()**
x, y, alpha, color, group, linetype, size

**c + geom_histogram**(binwidth = 5)
x, y, alpha, color, fill, linetype, size, weight

**c2 + geom_qq**(aes(sample = hwy))
x, y, alpha, color, fill, linetype, size, weight

### discrete
d <- ggplot(mpg, aes(fl))

**d + geom_bar()**
x, alpha, color, fill, linetype, size, weight

### TWO VARIABLES
**both continuous**
e <- ggplot(mpg, aes(cty, hwy))

**e + geom_label**(aes(label = cty), nudge_x = 1, nudge_y = 1) - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

**e + geom_point()**
x, y, alpha, color, fill, shape, size, stroke

**e + geom_quantile()**
x, y, alpha, color, group, linetype, size, weight

**e + geom_rug**(sides = "bl")
x, y, alpha, color, linetype, size

**e + geom_smooth**(method = lm)
x, y, alpha, color, fill, group, linetype, size, weight

**e + geom_text**(aes(label = cty), nudge_x = 1, nudge_y = 1) - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

**one discrete, one continuous**
f <- ggplot(mpg, aes(class, hwy))

**f + geom_col()**
x, y, alpha, color, fill, group, linetype, size

**f + geom_boxplot()**
x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight

**f + geom_dotplot**(binaxis = "y", stackdir = "center")
x, y, alpha, color, fill, group

**f + geom_violin**(scale = "area")
x, y, alpha, color, fill, group, linetype, size, weight

**both discrete**
g <- ggplot(diamonds, aes(cut, color))

**g + geom_count()**
x, y, alpha, color, fill, shape, size, stroke

**e + geom_jitter**(height = 2, width = 2)
x, y, alpha, color, fill, shape, size

### THREE VARIABLES
seals$z <- with(seals, sqrt(delta_long^2 + delta_lat^2)); l <- ggplot(seals, aes(long, lat))

**l + geom_contour**(aes(z = z))
x, y, z, alpha, color, group, linetype, size, weight

**l + geom_contour_filled**(aes(fill = z))
x, y, alpha, color, fill, group, linetype, size, subgroup

### continuous bivariate distribution
h <- ggplot(diamonds, aes(carat, price))

**h + geom_bin2d**(binwidth = c(0.25, 500))
x, y, alpha, color, fill, linetype, size, weight

**h + geom_density_2d()**
x, y, alpha, color, group, linetype, size

**h + geom_hex()**
x, y, alpha, color, fill, size

### continuous function
i <- ggplot(economics, aes(date, unemploy))

**i + geom_area()**
x, y, alpha, color, fill, linetype, size

**i + geom_line()**
x, y, alpha, color, group, linetype, size

**i + geom_step**(direction = "hv")
x, y, alpha, color, group, linetype, size

### visualizing error
df <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)
j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))

**j + geom_crossbar**(fatten = 2) - x, y, ymax, ymin, alpha, color, fill, group, linetype, size

**j + geom_errorbar()** - x, ymax, ymin, alpha, color, group, linetype, size, width
Also **geom_errorbarh()**.

**j + geom_linerange()**
x, ymin, ymax, alpha, color, group, linetype, size

**j + geom_pointrange()** - x, y, ymin, ymax, alpha, color, fill, group, linetype, shape, size

### maps
Draw the appropriate geometric object depending on the simple features present in the data. aes() arguments: map_id, alpha, color, fill, linetype, linewidth.

nc <- **sf::st_read**(system.file("shape/nc.shp", package = "sf"))

ggplot(nc) +
  **geom_sf**(aes(fill = AREA))

**l + geom_raster**(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE)
x, y, alpha, fill

**l + geom_tile**(aes(fill = z))
x, y, alpha, color, fill, linetype, size, width

# R resources: Fun stuff

1. R for cats – Fun tutorial about the basics of R programming… with cats!

2. Learning R – Episode of the PolicyViz podcast with Jonathan Schwabish from the Urban Institute where he talks about his approach to learning R

# Demos

# Census data

**Demo topic:** Working with census data in R

**Date:** July 9, 2024

**Presenters:** Coral Wonderly

**Packages used:** Tidycensus, Tidyverse, Insight

**Script link:** Markdown Version:
https://github.com/CTOpenData/r-user-group/blob/main/tidycensus.Rmd
PDF Version: https://github.com/CTOpenData/r-user-group/blob/main/tidycensus.pdf

# Cleaning address data

**Demo topic:** Cleaning a column with unstandardized town name data

**Date:** June 4, 2024

**Presenters:** Sarah Hurley and Pauline Zaldonis

**Packages used:** dplyr, readr, stringr

**Script link:** https://github.com/CTOpenData/r-user-group/blob/main/address_data_cleaning.R

# Appendix

# Who we are

# What is R?

Data & Policy Analytics

# What is R?

- Programming language for statistical computing and data visualization

- Open-source and free to use
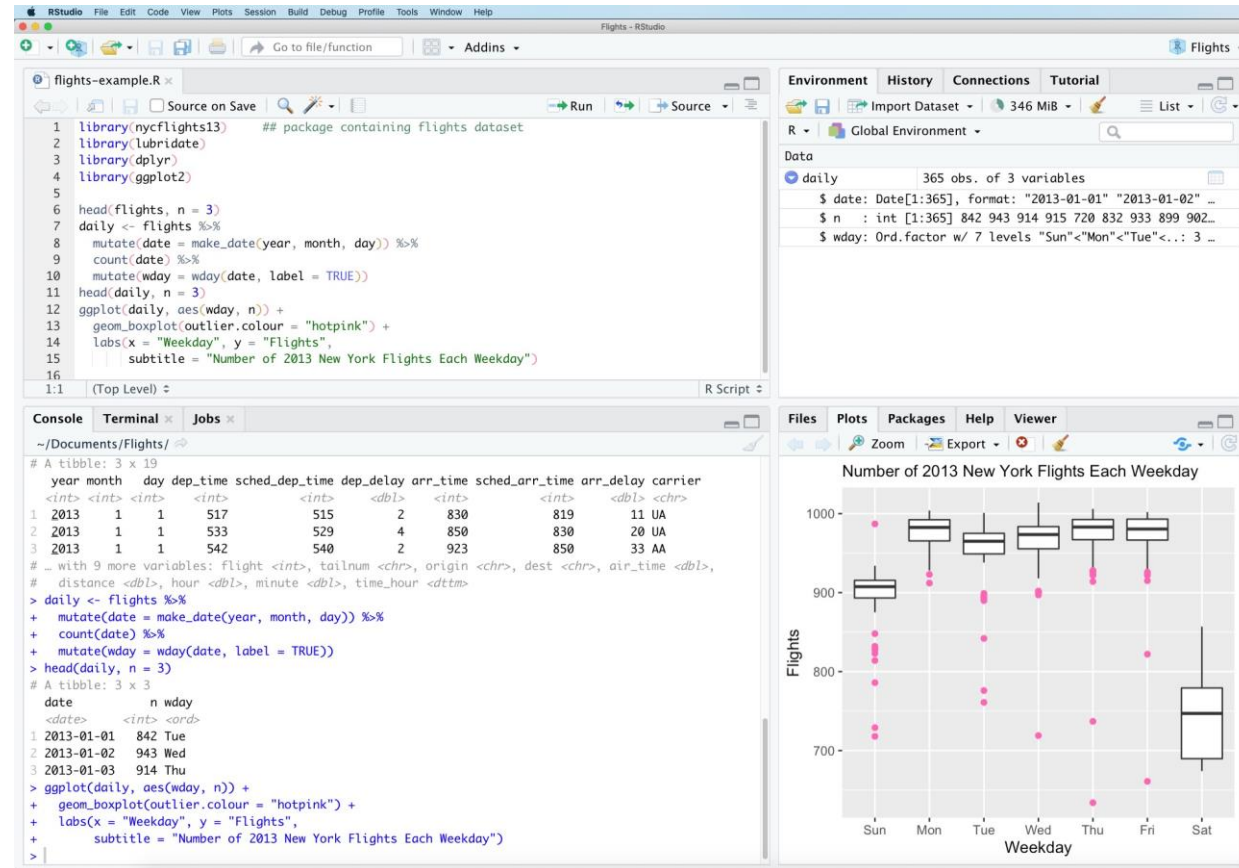
- Created by statisticians for statisticians

# Key features of R

- Statistical analysis: linear and nonlinear modeling, time-series analysis, classification, clustering

- Data visualization: high-quality graphs and plots

- Extensible with thousands of packages available

- Comprehensive R Archive Network (CRAN) hosts user-contributed packages

- Large and active user community

- Extensive documentation and support available

# R Studio

- R Studio is the integrated development environment (IDE) for working with R

- User-friendly interface for writing and debugging R code

- Enhances productivity and ease of use

# Why use R?

# Why use R?

- **Interoperability**
  - Integrates with other programming languages and tools
  - Compatibility with Python, SQL, Hadoop, etc.
  - Facilitates seamless workflow in data science projects
- **Reproducible processes**
  - Create processes that you can quickly repeat & reproduce results
  - Tools like R Markdown for creating dynamic documents
    - Combines code, output, and narrative text
    - Ensures reproducible research and reports

# Discussion

# Discussion questions

1. What R projects have you been working on?

2. What questions do you have?

3. What do you want to learn about R?

4. What would you be willing to demo at a future R meeting?