

Intro to DS

- data import

```
# using readr in tidyverse
library(readr)
data <- read_csv("file.csv") # Reads a CSV file from the current directory
data <- read_csv("../file.csv") # Reads a CSV file from the parent directory
data <- read_csv("~/path/file.csv") # Reads a CSV file using an absolute path

# Reads an Excel sheet
library(readxl)
data <- read_excel("data.xlsx", sheet = "Sheet1")

# Reads an RDS file
data <- readRDS("data.rds")
saveRDS(data, "file.rds") # Saves a dataframe as an RDS file

# Reads a tab-delimited text file(.txt or .tsv)
# using readr in tidyverse
library(readr)
data <- read_delim("file.txt", delim = "\t")
# read_delim() reads files with a specified delimiter(分隔符)
# delim = "\t" tells R that the columns are separated by tabs

# Reads a stata file
library(haven)
data <- read_dta("file.dta")

# using basic R
read.csv("file.csv") # Base R version
```

- components of graphics
 - a frame/coordinate system
 - x-axis and y-axis specifications and etc.
 - a layer
 - geometric elements such as lines and points
 - each type is a separate layer
 - scales
 - aesthetics added such as colors, sizes and shapes
 - faceting
 - splitting up of the data into multiple subplots
 - a theme

- additional controls of the plot aesthetics such as font types, backgrounds and color scheme
- examine data properties using visualization
 - typical outcome
 - variability & range
 - shape
 - outliers
- univariate - data visualization

```
# only categorical x
# bar graph
ggplot(data, aes(x = x)) +
  geom_bar(color = "orange", fill = "blue") +
  labs(x = "Rating", y = "Number of hikes") +
  theme_minimal()

# quantitative x
# histogram
# divide up the observed range of the variable into bins of equal width
# count up the number of cases that fall into each bin
ggplot(data, aes(x = x)) +
  geom_histogram(color = "white", binwidth = 5) +
  labs(x = "Elevation (feet)", y = "Number of hikes")

# quantitative x
# density plot
ggplot(data, aes(x = x)) +
  geom_density()
```

- bivariate - data visualization

```
# scatter plot
# numerical x
# numerical y
ggplot(data, aes(x=x, y=y)) +
  geom_point()

# Scatter + trend line
ggplot(data, aes(x=x, y=y))+
# size sets the size of the points
# alpha sets transparency (1 = fully opaque, 0 = fully transparent)
geom_point(size = 2, alpha=0.7) +
# geom_smooth() adds a smoothed line to represent a trend using LOESS
# method = "lm" specifies that a linear model
# se = TRUE includes the shaded confidence interval around the regression line
```

```
# span = 0.5 controls 50% of the nearest points is used to compute the local
smooth fit
geom_smooth(method="lm",se=TRUE, span = 0.5)

# Line charts
# numerical x or data x
# numerical y
ggplot(data, aes(x=x, y=y)) +
geom_line()

# jitter plot
# alternative for scatter plots when categorical variables are involved
# categorical x
# numerical y
ggplot(data, aes(x = x, y = y)) +
# width sets jitter along x-axis (0 ~ infy)
# you can adjust both width and height
geom_jitter(width = 0.2)

# violin plot
# categorical x
# numerical y
ggplot(elections, aes(y = y, x = x)) +
geom_violin()

# box plot
# categorical x
# numerical y
ggplot(data, aes(x = x, y = y)) +
geom_boxplot()

# jitter plot
# alternative for scatter plots when categorical variables are involved
# numerical x
# categorical y
ggplot(data, aes(x = x, y = y)) +
# height sets jitter along y-axis (0 ~ infy)
geom_jitter(height = 0.2) +
coord_flip()

# density plot
# The density plots are on top of each other
# numerical x
# categorical y
ggplot(data, aes(x = x, fill = y)) +
# alpha = 0.5 adds transparency
```

```
geom_density(alpha = 0.5) +  
# scale_fill_manual defines what colors to use for the fill categories  
scale_fill_manual(values = c("blue", "purple", "red")) +  
# facet_wrap separates the density plots into facets for each category  
facet_wrap(~ y)
```

```
# histogram plot  
# The histogram plots are on top of each other  
# numerical x  
# categorical y  
ggplot(data, aes(x = x, fill = y)) +  
# color sets the color for frame  
geom_histogram(color = "white") +  
scale_fill_manual(values = c("blue", "purple", "red"))
```

```
# ridge plot  
# numerical x  
# categorical y  
# used for y with many categories  
# Install ggridges package  
library(ggridges)  
ggplot(data, aes(x = x, y = y)) +  
geom_density_ridges()
```

```
# count plot  
# categorical x  
# categorical y  
ggplot(data, aes(x = x, y = y)) +  
geom_count()
```

```
# A side-by-side bar plot  
# the same as count plot  
# y-axis = counts  
# categorical x  
# categorical y  
ggplot(data, aes(x = x, fill = y)) +  
  geom_bar(position = "dodge")
```

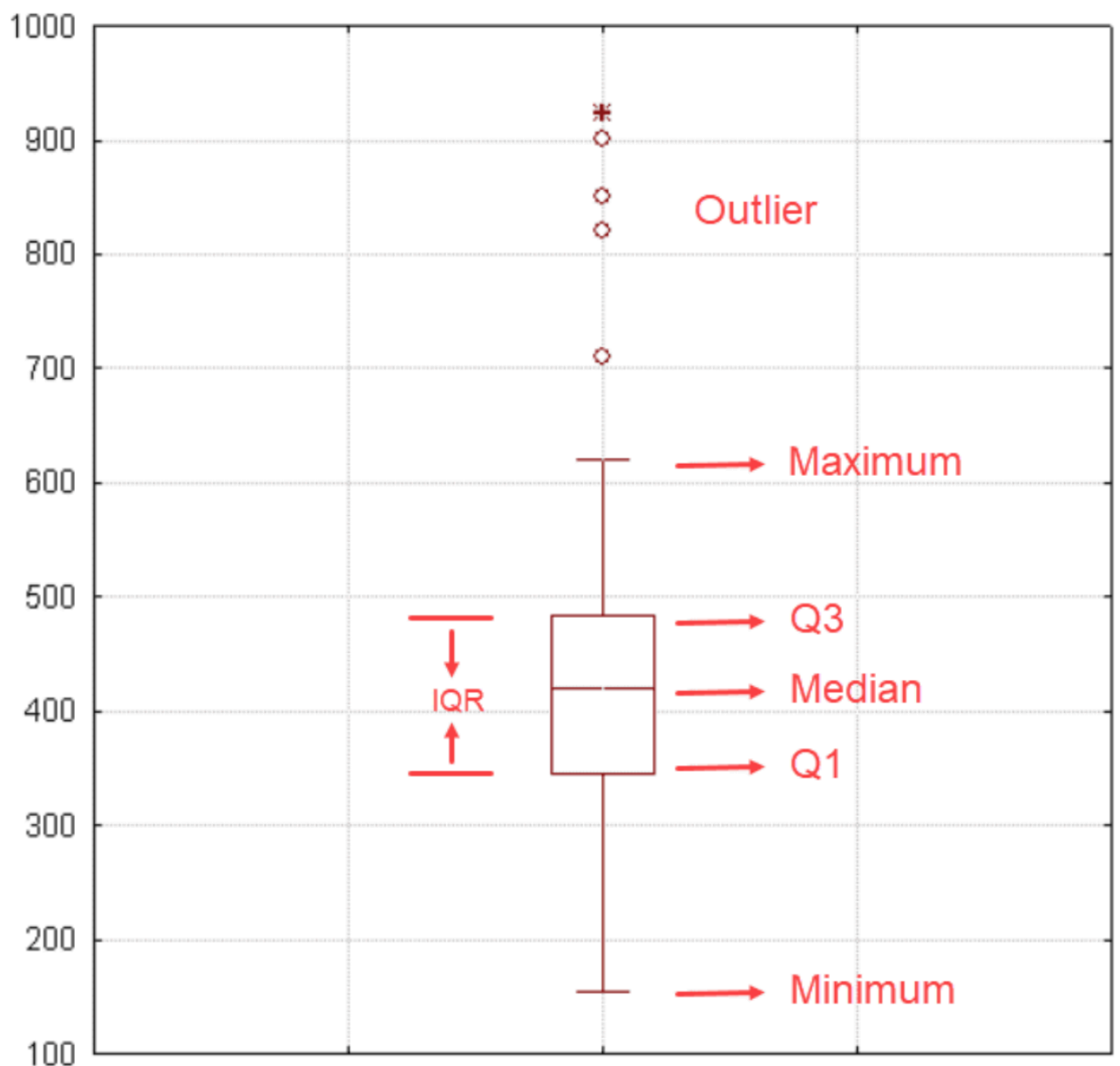
```
# A stacked bar plot  
# y-axis = counts  
# categorical x  
# categorical y  
ggplot(data, aes(x = x, fill = y)) +  
  geom_bar()
```

```
# A faceted bar plot
```

```
# categorical x
# categorical y
ggplot(data, aes(x = x)) +
  geom_bar() +
  facet_wrap(~ y)

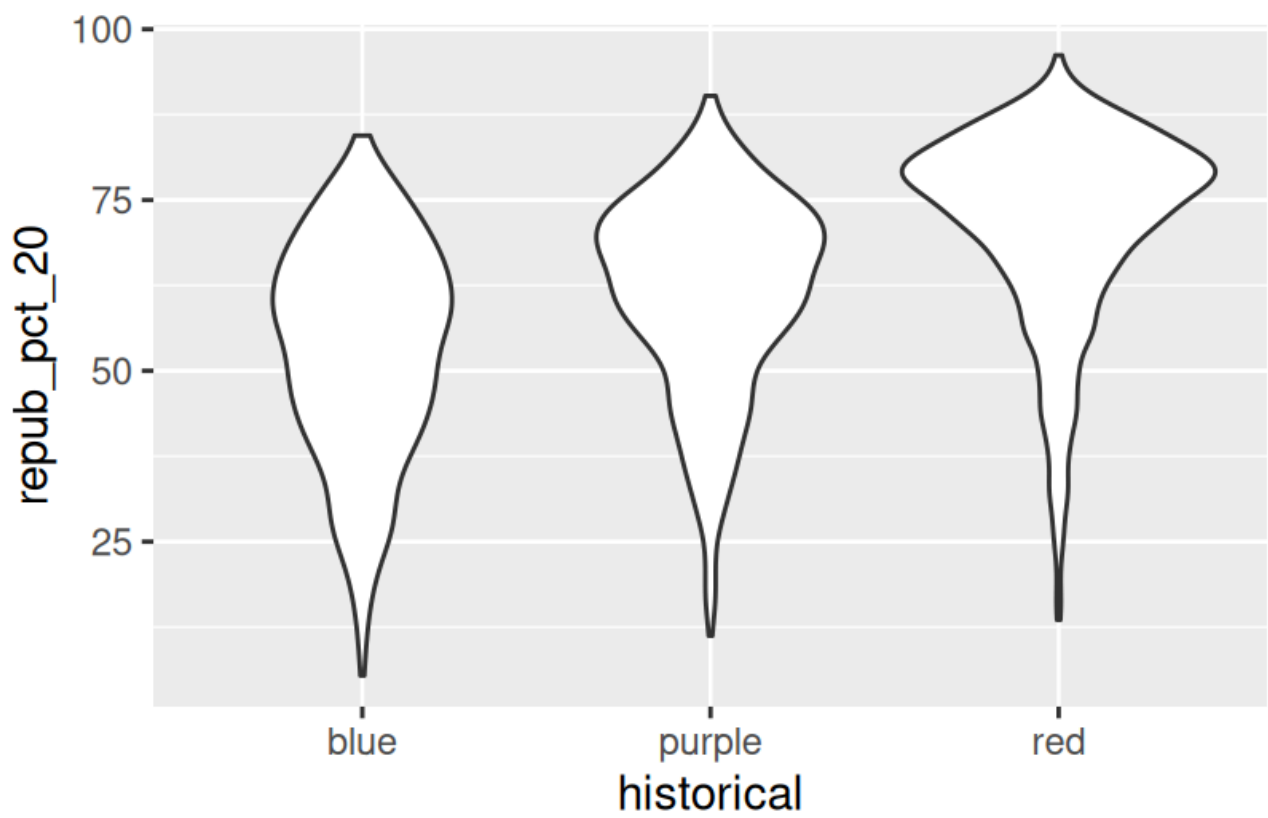
# A proportional bar plot
# y-axis = percentage of each y category in x's category groups
# categorical x
# categorical y
ggplot(elections, aes(x = x, fill = y)) +
  geom_bar(position = "fill")
```

- box plot
 - interpretation
 - box: 50% of the data
 - median: central line
 - example



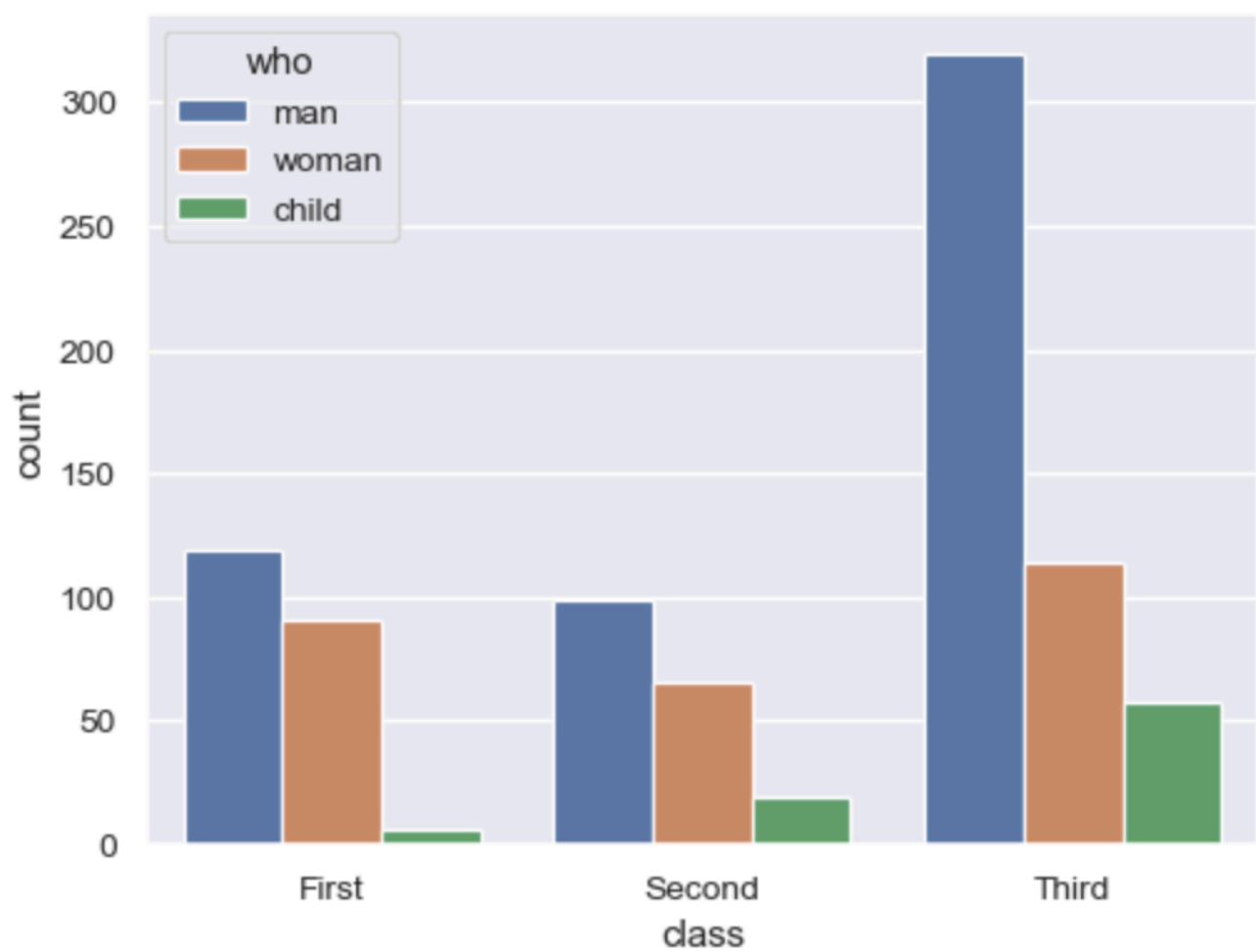
- violin plot

- example



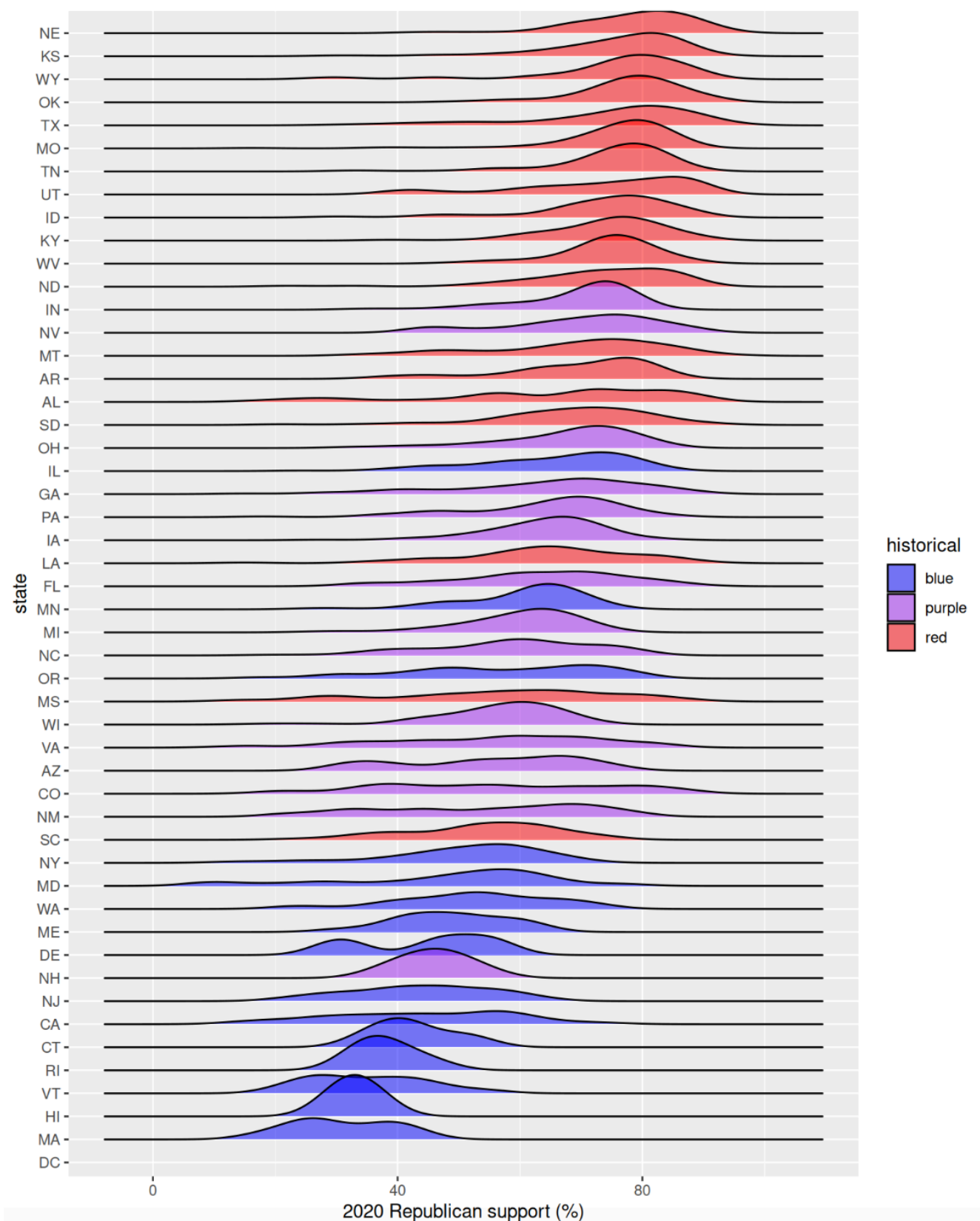
- count plot

- example



- ridge plot

- example



- multivariate - data visualization

```
# ridge plot
# numerical x
# categorical y
# categorical z
# reorder y based on x
# desc. = TRUE sets the reorder from largest to smallest from top to bottom
# desc. = FALSE sets the reorder from smallest to largest from top to bottom
```

```

ggplot(data, aes(x = x, y = fct_reorder(y, x, desc. = TRUE), fill = z)) +
  geom_density_ridges(alpha = 0.5) +
  labs(y = "state", x = "2020 Republican support (%)") +
  # note z have only 3 categories / colors
  scale_fill_manual(values = c("blue", "purple", "red"))

# scatter plot
# numerical x
# numerical y
# categorical z
# applicable when x and y has the same scale and unit such as election info in
different years
ggplot(data, aes(x = x, y = fct_reorder(z, x))) +
  geom_point(color = "red") +
  geom_point(aes(x = y, y = z))

# scatter plot
# numerical x
# numerical y
# categorical g
# size is not a variable
# categorical sh
ggplot(data, aes(x=x, y=y, color=g, size=s, shape=sh)) +
  geom_point()

# scatter plot + trend
# numerical x
# numerical y
# numerical z
# cut(z, 2) divides z variable into 2 groups and assign different colors for each
# note that there are 2 color groups thus 2 fitted lines
ggplot(data, aes(y = y, x = x, color = cut(z, 2))) +
  geom_point() +
  geom_smooth(se = FALSE, method = "lm")

# heatmap
# categorical x
# categorical y
# numerical z
ggplot(data, aes(x = x, y = y)) +
  geom_tile(aes(fill = z)) +
  scale_fill_gradient(low = "white", high = "red")

# heatmap
# scale the data before viz
# Load the gplots package needed for heatmaps

```



```

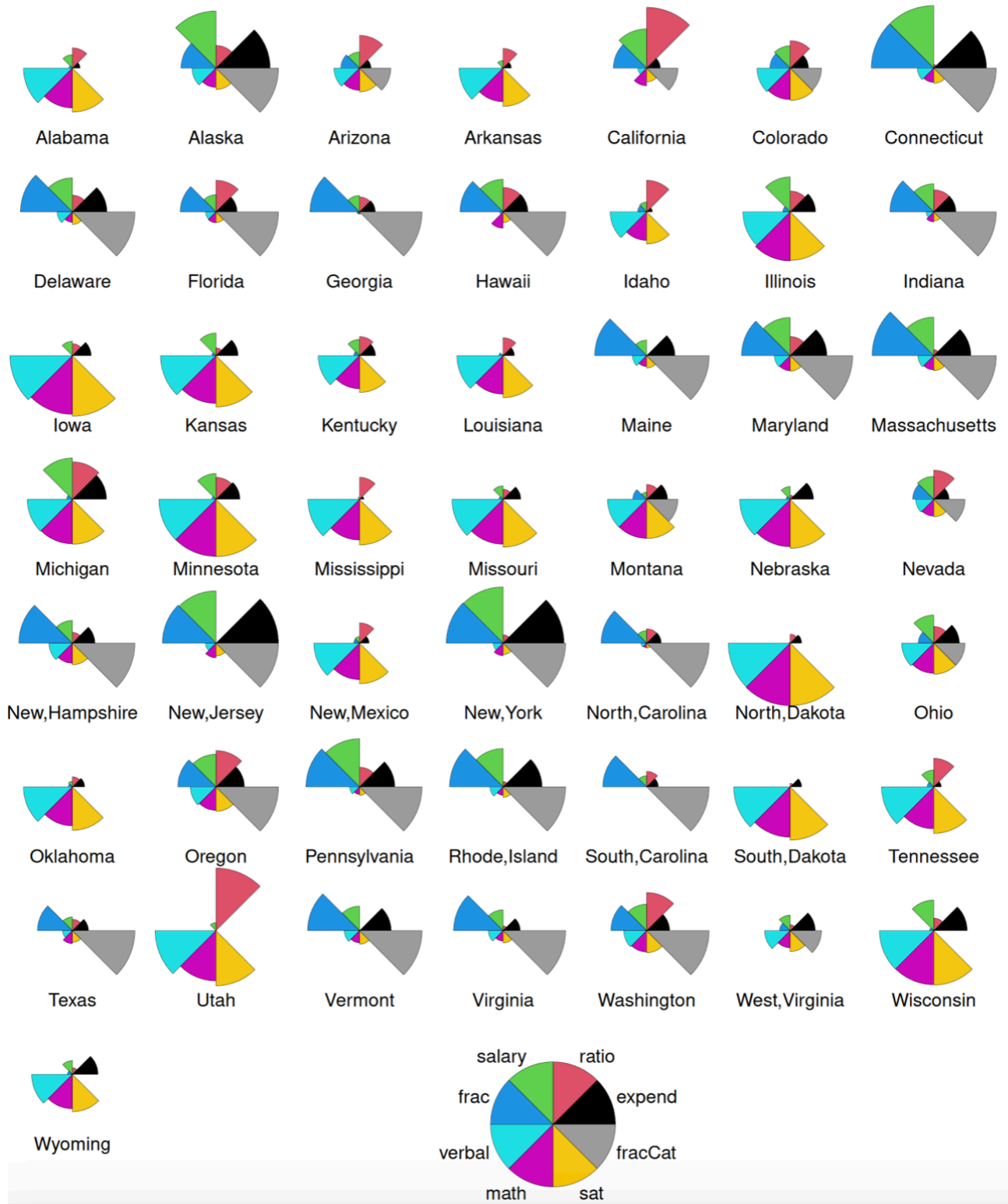
library(gplots)
# heatmap.2() = heatmap function
# data: a numerical matrix where rows and columns are heatmap axes and value
# represented by color intensity
heatmap.2(data,
# dendrogram : tree-like diagram that shows hierachical clustering
# if = "row" cluster rows only
# if = "column" cluster columns only
# if = "both" cluster columns and rows
dendrogram = "none",
# Rowv = NA prevent reordering of rows based on hierachical clustering
# if = TRUE rows will be reordered
Rowv = NA,
# standardize each column
# if = "row" standardize each row
# if = "none" no standardization
scale = "column",
# adjust the size of the color key
keysize = 0.7,
# if = "density" it would shouw value density as a histogram
# if = "histogram" it would display a frequency histogram
density.info = "none",
# define the color scale for the heatmap
# hcl.colors(256) generates 256 gradient colors using the Hue-Chroma-Luminance
# color model
col = hcl.colors(256),
# adjust margin around the heatmap
# 10 units of margin for row label
# 20 units of margin for column labels
margins = c(10, 20),
# colsep = c(1:7) adds vertical seperators after the first 7 columns
# rowsep = (1:50) adds horizontal seperators after the first 50 rows
# sepwidth = c(0.05, 0.05) defines the thinkness of column and row seperators
colsep = c(1:7), rowsep = (1:50), sepwidth = c(0.05, 0.05),
# sepcolor defines the color of the seperator
# trace = "none" removes trace lines
# if = "row" adds trace lines along rows
# if = "column" adds trace lines along columns
# if = "both" adds trace lines in both directions
sepcolor = "white", trace = "none"
)

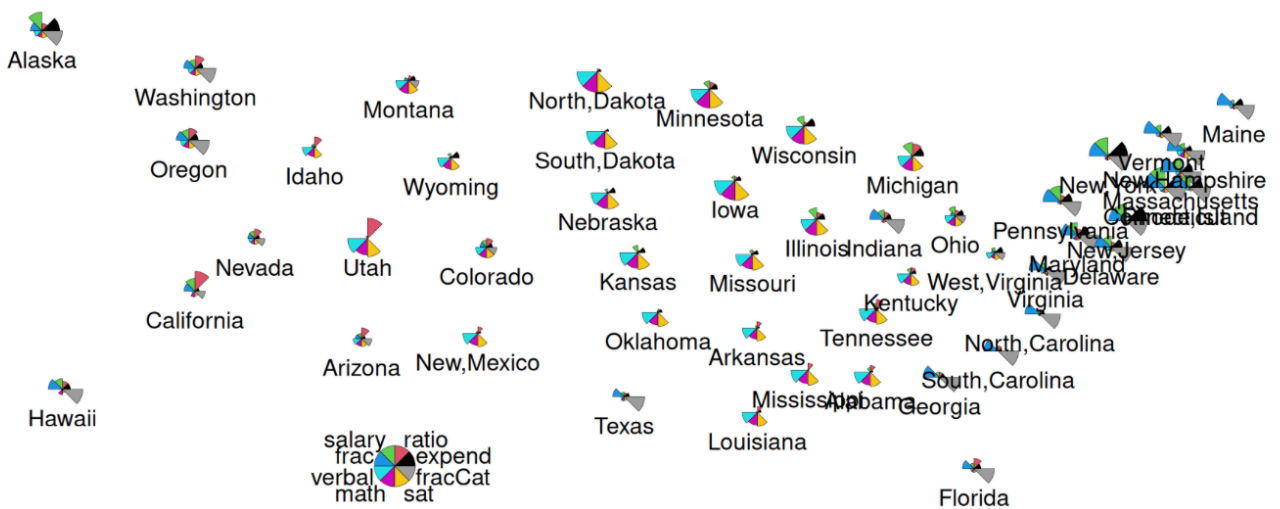
# star plots
# data must only contain numerical scaled values
stars(data,
  flip.labels = FALSE,

```

```
key.loc = c(10, 1.5),
cex = 1,
draw.segments = TRUE
)
```

- starplot
 - example





- map - data visualization
 - types
 - point maps
 - plotting locations of individual observations
 - contour maps
 - plotting the *density* or distribution of observations
 - choropleth maps
 - methods
 - ggplot() + geom_map()
 - Plot data points on top of a map
 - create choropleth maps
 - leaflet
 - create interactive maps
 - general steps
 - create a map widget
 - add a base map
 - add layers
 - print the map widgets to display it

```
# leaflet

# Load the leaflet package
library(leaflet)

# create a plotting frame (no map)
leaflet(data = fave_places)

# add the map background
leaflet(data = fave_places) |>
  addTiles()
```

```

# add map background and markers
# longitude and latitude refer to the variables in our data
leaflet(data = fave_places) |>
  addTiles() |>
  addMarkers(lng = ~longitude, lat = ~latitude)

# if we named them "longitude" and "latitude", the function automatically
recognizes these variables
leaflet(data = fave_places) |>
  addTiles() |>
  addMarkers()

# Load package needed to change color
library(gplots)

# We can add colored circles instead of markers at each location
leaflet(data = fave_places) |>
  addTiles() |>
# colors need to be in hex form
# col2hex is used to convert the color into the hex form
  addCircles(color = col2hex("red"))

leaflet(data = fave_places) |>
# can change the background
  addProviderTiles("USGS") |>
# weight sets the thickness to make the lines and circles
  addCircles(weight = 10, opacity = 1, color = col2hex("yellow")) |>
# connect the dots, in their order in the dataset, with green lines
  addPolylines(
    lng = ~longitude,
    lat = ~latitude,
    color = col2hex("green")
  )

```

```

# ggplot

# Load the package
library(rnaturalearth)

# Get info about country boundaries across the world
# in a "sf" or simple feature format
# world maps
# the continent of Africa: ne_countries(continent = 'Africa', returnclass = 'sf')
# a set of countries: ne_countries(country = c('france', 'united kingdom',
'germany'), returnclass = 'sf')
# boundaries within a country: ne_states(country = 'united states of america',

```

```

returnclass = 'sf')
world_boundaries <- ne_countries(returnclass = "sf")

# produce a basic map
# initialize a ggplot object
ggplot(world_boundaries) +
# adds a simple features layers for maps
  geom_sf()

# Load package needed to change map theme
library(mosaic)

# Add a point for each Starbucks
# NOTE: The Starbucks info is in our starbucks data, not world_boundaries
ggplot(world_boundaries) +
  geom_sf() +
  geom_point(
    data = starbucks,
    aes(x = Longitude, y = Latitude),
    alpha = 0.3, size = 0.2, color = "darkgreen"
  ) +
  theme_map()

ggplot(cma_boundaries) +
  geom_sf() +
  geom_point(
    data = starbucks_cma,
    aes(x = Longitude, y = Latitude),
    alpha = 0.3,
    size = 0.2,
    color = "darkgreen"
  ) +
    # restrict the x-axis to the range between -179.14 and -50
  coord_sf(xlim = c(-179.14, -50)) +
  theme_map()

# create background maps of state and county level
# Load packages
library(sf)
library(maps)

# Get the boundaries
# extract county-level boundary data
# st_as_sf() converts the output into a simple feature (sf) object
# maps::map() uses the map function from the maps package
# fill = TRUE ensures that county boundaries are closed polygons

```

```

# plot = FALSE prevents the function from automatically plotting the map
midwest_boundaries <- st_as_sf(maps::map("county", region = c("minnesota",
"wisconsin", "north dakota", "south dakota"), fill = TRUE, plot = FALSE))

# Get the latitude and longitude coordinates of state boundaries
states_map <- map_data("state")

# create choropleth maps
# map_id specifies which variable in our dataset indicates the region
# map_id and the region variable in our mapping background must have the same
possible outcomes in order to be matched up
# expand_limits assures that the map covers the entire area by pulling
longitudes and latitudes
ggplot(elections_by_state, aes(map_id = state_name, fill = repub_pct_20)) +
  geom_map(map = states_map) +
  expand_limits(x = states_map$long, y = states_map$lat) +
  theme_map()

ggplot(elections_by_state, aes(map_id = state_name, fill = repub_pct_20)) +
  geom_map(map = states_map) +
  expand_limits(x = states_map$long, y = states_map$lat) +
  theme_map() +
  scale_fill_gradientn(name = "% Republican", colors = c("blue", "purple", "red"),
values = scales::rescale(seq(0, 100, by = 5)))

ggplot(elections_by_state, aes(map_id = state_name, fill = repub_20_categories)) +
  geom_map(map = states_map) +
  expand_limits(x = states_map$long, y = states_map$lat) +
  theme_map()

# Load package needed for refining color palette
library(RColorBrewer)

# use the refined colors
ggplot(elections_by_state, aes(map_id = state_name, fill = repub_20_categories)) +
  geom_map(map = states_map) +
  expand_limits(x = states_map$long, y = states_map$lat) +
  theme_map() +
  scale_fill_manual(values = rev(brewer.pal(8, "RdBu")), name = "% Republican")

# add points
ggplot(elections_by_state, aes(map_id = state_name, fill = repub_20_categories)) +
  geom_map(map = states_map) +
  geom_point(
    data = starbucks_us,

```

```

    aes(x = Longitude, y = Latitude),
    size = 0.05,
    alpha = 0.2,
    inherit.aes = FALSE
  ) +
  expand_limits(x = states_map$long, y = states_map$lat) +
  theme_map() +
  scale_fill_manual(values = rev(brewer.pal(8, "RdBu")), name = "% Republican")

ggplot(elections_by_counties, aes(map_id = county_fips, fill =
repub_20_categories)) +
  geom_map(map = county_map) +
  scale_fill_manual(values = rev(brewer.pal(10, "RdBu")), name = "% Republican") +
  expand_limits(x = county_map$long, y = county_map$lat) +
  theme_map() +
# moves the legend to the right
  theme(legend.position = "right") +
# ensures equal aspect ratio between xy axes
  coord_equal()

ggplot(elections_by_counties, aes(map_id = county_fips, fill = median_rent)) +
  geom_map(map = county_map) +
  expand_limits(x = county_map$long, y = county_map$lat) +
  theme_map() +
  theme(legend.position = "right") +
  coord_equal() +
  scale_fill_gradientn(name = "median rent", colors = c("white", "lightgreen",
"darkgreen"))

ggplot(elections_by_counties, aes(map_id = county_fips, fill = median_age)) +
  geom_map(map = county_map) +
  expand_limits(x = county_map$long, y = county_map$lat) +
  theme_map() +
  theme(legend.position = "right") +
  coord_equal() +
# scale_fill_gradientn() specifies a custom color gradient for a continuous
variable mapped to fill
# name = "median age" sets the legend title for the color scale
# colors = terrain.colors(10) uses 10 colors from the "terrain" color palette to
create a smooth gradient
  scale_fill_gradientn(name = "median age", colors = terrain.colors(10))

```

- addition - data visualization

```

# Faceting
# create separate plots for each level of category

```

```
# scales="free y" allows the y-axis scale to vary independently in each facet
facet_wrap(~ category, scales="free y")

# Faceting
# creates a grid layout
# scales="free" allows both axes to adjust independently
facet_grid(row ~ column, scales="free")

# Customization
# title is at the top and caption is at the low end
labs(title="Title", x="X-axis", y="Y-axis", caption="Source")

# themes
# no background color
# light gray gridlines
# no box around the plot
# modern and clean appearance
theme_minimal()

# themes
# white background with no gridlines
# black axes
# no additional distractions
theme_classic()

# themes
# white background with black gridlines
# box around plot
theme_bw()

# swaps x and y axes
coord_flip()

# converts Cartesian coordinates to polar
# used to create pie charts, radial bar plots, and circular graphs
coord_polar()

# transform the x-axis to a logarithmic scale
scale_x_log10()

# restrict y-axis range between 0 and 100
scale_y_continuous(limits=c(0,100))

# applies a color palette from RColorBrewer package for discrete/categorical
variables
scale_fill_brewer(palette="Set1")
```



```
# applies the viridis color scale for continous data
# viridis is colorblind-friendly
scale_color_viridis_c()

# applies the viridis color scale for discrete data
# viridis is colorblind-friendly
scale_color_viridis_d()

# add text annotate
annotate("text", x=5, y=20, label="Point")

# remove legend
theme(legend.position="none")
```

- data exploration

```
# What type of variables do we have?
str(data)

# number of data points
nrow(data)

# number of variables
ncol(data)

# return the number of rows and the number of columns
dim(data)

# Check out the first 6 rows
head(data)

# Check out the first 5 rows
head(data, n=5)

# Check out the last 5 rows
tail(data, n=5)

# shows column data types & first values
glimpse(data)

# Statistical summary of each variable
summary(data)

# types for each variable and first few values for each
str(data)
```

```

colnames(data)

rownames(data)

# shows the column's values as a dataframe (no repetitions)
distinct(data, var)

# Extract the column's values as a vector
pull(data, var)

# Select rows by position
data %>% slice(1:10)

```

- data wrangling

```

# arrange the rows according to some column var
# from lowest to highest
arrange(var)
# from highest to lowest
arrange(desc(var))
# arrange var1 first and then var2
arrange(var1,var2)

# filter out or obtain a subset of the rows
filter(var = ??)
# %in% c(???, ???) means a list of multiple values
filter(state_name %in% c("Hawaii", "Delaware"))

# select a subset of columns
select (var1, var2)

# mutate or create a column
mutate(var1 = ??)
# Define repub_win_20, whether the Repub won in 2020 (TRUE or FALSE)
mutate(elections_small, repub_win_20 = repub_pct_20 > dem_pct_20)

# calculate a numerical summary of a column
summarize()

# group the rows by a specified column
group_by()

```

- pipe functions

```
# Without a pipe
filter(elections, state_name == "Minnesota")

# With a pipe
elections |>
  filter(state_name == "Minnesota")
```

- String Function

```
str length("abc") # Length: 3
str sub("abcdef", 2, 4) # Substring:

"bcd"

str detect("abc", "a") # Detect: TRUE
str count("ababc", "ab") # Count: 2 str replace("abc", "a", "x") # Re- place first
str replace all("aab", "a", "x") # Replace all
str trim(" abc ") # Remove whites- pace
str to lower("ABC");
str to upper("abc")
str to title("hello world") # "Hello World"
str c("a", "b", sep="-") # Concate- nate
str glue("Name: {name}") # String interpolation
str extract("abc123", "\d+") # Ex- tract: "123"
str extract all("a1b2c3", "\d") # c("1","2","3")
str match("name: Alice", "(\w+): (\w+)")
str split("a,b,c", ",") # Split to list
```

- statistical functions

```
mean(x, na.rm=TRUE)
```

```
median(x, na.rm=TRUE)
```

```
min(x)
```

```
max(x)
```

```
range(x)
```

```
sd(x)
```

```
var(x)
```

```
IQR(x)
```

```
quantile(x, probs=c(0.25,0.75))
```

```
cor(x, y)
```

```
cov(x, y)
```

```
t.test(x, y); wilcox.test(x, y)
```

```
lm(y~ x, data=df) # Linear regression
```

```
glm(y~ x, family=binomial, data=df) # Logistic regression  
summary(model) # Model summary  
predict(model, newdata) # Make predictions
```

```
scale(x) # Standardize data
```

```
cut(x, breaks=c(0,5,10)) # Categorize
```

```
numeric
```

```
cut number(x, n=4) # Equal sized groups
```

```
is.na(df$col); !is.na(df$col); sum(is.na(df$col))
```

```
filter(!is.na(col)) # Keep non-NA rows  
drop na(df) # Remove rows with any NA
```

```
drop na(df, col1, col2) # Specific columns
```

```
replace na(df, list(col1=0, col2="Unknown"))
```

```
fill(df, col1, .direction="down") # Fill NA with prev val
```

```
complete(df, x, y) # Create all combinations
```

```
na if(df$col, -99) # Replace values with NA
```

- color palette

```
# Load package needed for refining color palette
```

```
library(RColorBrewer)
```

```
# Now fix the colors
```

```
ggplot(elections_by_state, aes(map_id = state_name, fill = repub_20_categories)) +  
  geom_map(map = states_map) +  
  expand_limits(x = states_map$long, y = states_map$lat) +  
  theme_map() +  
  scale_fill_manual(values = rev(brewer.pal(8, "RdBu")), name = "% Republican")
```

- Date & Time (lubridate)

```
# Convert string in year-month-day format to date object
```

```
ymd("2023-02-24")
```

```
# Convert string in month/day/year format to date object
mdy("02/24/2023")

# Convert string in day-month-year format to date object
dmy("24-02-2023")

# Convert string with date and time to datetime object
ymd_hms("2023-02-24 13:45:30")

# Get current date and time
now()

# Get current date
today()

# Convert string to date object
as.date("2023-02-24")

# Extract year from date object
year(date)

# Extract month number from date object
month(date)

# Extract day number from date object
day(date)

# Get month name from date object
month(date, label=TRUE)

# Get day of week as number
wday(date)

# Get day of week as name
wday(date, label=TRUE)

# Extract hour from datetime object
hour(dt)

# Extract minute from datetime object
minute(dt)

# Extract second from datetime object
second(dt)

# Add 7 days to a date
```

```
date + days(7)

# Add 1 month to a date
date + months(1)

# Add 1 year to a date
date + years(1)

# Calculate time interval between two dates
interval(date1, date2)

# Convert interval to duration in seconds
as.duration(interval)

# Round date down to start of month
floor_date(date, "month")

# Round date up to end of month
ceiling_date(date, "month")

# Round date to nearest week
round_date(date, "week")

# Check if year is a leap year
leap_year(year)

# Set timezone for a time object
with_tz(time, "America/New_York")
```

missing values

```
# Filter non-missing
filter(!is.na(column))

# Replace missing
mutate(col_clean = replace_na(col, "refused"))

# Drop missing rows
drop_na()
```

reshape

```
# Convert data from wide to long format
# Take columns q1, q2, q3, q4 and stack them into two columns: "quarter" and
"revenue"
pivot_longer(df, cols = c(q1, q2, q3, q4), names_to = "quarter", values_to =
```

```
"revenue")
```

```
# Convert all columns to long format except id & year  
pivot_longer(df, cols = -c(id, year), names_to = c("quarter", "type"), names_sep =  
" ", values_to = "value")
```

```
# Convert data from long to wide format  
# Spread values in the "revenue" column across new columns named by values in the  
"quarter" column  
pivot_wider(df, names_from = quarter, values_from = revenue)
```

```
# Create complex wide format with multiple components in column names  
pivot_wider(df, id_cols = id, names_from = c(year, quarter), names_sep = " ",  
values_from = value)
```

```
# Split a column into multiple columns based on a separator  
separate(df, col, into = c("year", "month"), sep = "-")
```

```
# Combine multiple columns into a single column  
unite(df, "date", c("year", "month", "day"), sep = "-")
```

join

```
# Keep matching rows  
inner_join(table1 , table2)  
  
# Keep all rows from left table  
left_join(table1 , table2)  
  
# Keep all rows from both tables  
full_join(table1 , table2)  
  
# Rows from left that match right  
semi_join(table1 , table2)  
  
# Rows from left that don't match right  
anti_join(table1 , table2)  
  
# Join by specific column  
inner_join(table1 , table2 , by = "id")
```

factors

```
# Create a factor with three levels  
f <- factor(c("low", "high", "medium"))
```

```
# Get all levels of factor f
levels(f)

# Manually reorder factor levels
fct_relevel(f, "medium", "low", "high")

# Reorder factor levels based on values of a numeric variable x
fct_reorder(f, x)

# Reorder factor for line plots, using both x and y values
fct_reorder2(f, x, y)

# Order factor levels by frequency (most frequent first)
fct_infreq(f)

# Reverse the order of factor levels
fct_rev(fct_infreq(f))

# Keep top 2 most frequent levels, combine rest as "Other"
fct_lump_n(f, n=2)

# Lump together levels with <10% frequency as "Other"
fct_lump_prop(f, prop=0.1)

# Rename specific factor levels
fct_recode(f, "Low"="low", "High"="high")

# Collapse multiple factor levels into one
fct_collapse(f, Other=c("rare1", "rare2"))

# Keep only specified levels, recode rest as "Other"
fct_other(f, keep=c("a", "b"))

# Count occurrences of each factor level
fct_count(f)

# Convert character vector to factor
as.factor(c("a", "b"))
```