Data visualization with ggplot2

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



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



Complete the template below to build a graph.

<THEME_FUNCTION> <COORDINATE_FUNCTION> + stat = <STAT>, position = <POSITION>) + <SCALE_FUNCTION> + <FACET_FUNCTION> + <GEOM_FUNCTION>(mapping = aes (<MAPPINGS>) Not

gplot(data = mpg, aes(x = cty, y = hwy)) Begins a plot hat you finish by adding layers to. Add one geom

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

c + geom_area(stat = "bin")

x, y, alpha, color, fill, linetype, size

ast_plot() Returns the last plot

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

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",

linejoin - string ("round", "mitre", or "bevel") lineend - string ("round", "butt", or "square"

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

x, alpha, color, fill, linetype, size, weight

d + geom_bar()

X

Studio

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

ggplot(economics, aes(date, unemploy)) ggplot(seals, aes(x = long, y = lat))

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

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

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

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

b + geom_rect(aes(xmin = long, ymin = lat, xmax = long + 1, ymax = lat + 1)) - xmax, xm

ymax = unemploy + 900)) - x, ymax, ymin,

ymax, ymin, alpha, color, fill, linetype, size

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

e <- ggplot(mpg, aes(cty, hwy)) both continuous

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

e + geom_quantile()

e + geom_rug(sides = "bl")

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

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

i <- ggplot(economics, aes(date, unemploy))</pre> continuous function

i+geom_area()

i + geom_line()

x, y, alpha, color, fill, linetype, size

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

x, y, alpha, color, group, linetype, size

common aesthetics: x, y, alpha, color, linetype, size

LINE SEGMENTS

one discrete, one continuous
f <- ggplot(mpg, aes(class, hwy))</pre>

f + geom_boxplot()

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

b + geom_hline(aes(yintercept = lat))
b + geom_vline(aes(xintercept = long))

b + geom_abline(aes(intercept = 0, slope = 1))

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

color, fill, group, linetype, shape, size, weight x, y, lower, middle, upper, ymax, ymin, alpha

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

both discrete

x, y, alpha, color, fill, group, linetype, size, weight

f + geom_violin(scale = "area")

g <- ggplot(diamonds, aes(cut, color)) g + geom_count()

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

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

c + geom_freqpoly()

x, y, alpha, color, fill, shape, size, stroke

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

TWO VARIABLES

h <- ggplot(diamonds, aes(carat, price))

continuous bivariate distribution

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

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

h + geom_density_2d()

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

x, y, alpha, color, group, linetype, size, weight

x, y, alpha, color, linetype, size

visualizing error

 $df \leftarrow data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)$ $j \leftarrow ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))$

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

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

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

k <- ggplot(data, aes(fill = murder))

k + geom_map(aes(map_id = state), map = map) + expand_limits(x = map\$long, y = map\$lat)
map_id, alpha, color, fill, linetype, size

THREE VARIABLES

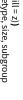
c2 + geom_qq(aes(sample = hwy))
x, y, alpha, color, fill, linetype, size, weight

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

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

l+ geom_contour(aes(z = z))





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



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

Stats An alternative way to build a layer.

A stat builds new variables to plot (e.g., count, prop)



Visualize a stat by changing the default stat of a geom function, geom_bar(stat="count") or by using a stat function, stat_count(geom="bar"), which calls a default geom to make a layer (equivalent to a geom function).
Use **..name..** syntax to map stat variables to aesthetics.

y=..count..



- **x, y** | ...count.., ..ncount.., ..density.., ..ndensity.. c + stat_bin(binwidth = 1, boundary = 10)
- c + stat_density(adjust = 1, kernel = "gaussian")
 x, y | ...count.., ..density.., ..scaled.. e + stat_bin_2d(bins = 30, drop = T) ..count.., ..density..

c + **stat_count(**width = **1) x, y** | ...count.., ..prop.

- e + stat_density_2d(contour = TRUE, n = 100) x, y, color, size | ..level.. e + stat_bin_hex(bins = 30) x, y, fill | ..count.., ..density..
- e + stat_ellipse(level = 0.95, segments = 51, type = "t")
- l + stat_summary_hex(aes(z = z), bins = 30, fun = max)
 x, y, z, fill | ..value.. $l + stat_contour(aes(z = z)) x, y, z, order | ..level..$
- l + **stat_summary_2d(**aes(z = z), bins = 30, fun = mean)
- f + stat_boxplot(coef = 1.5)
- **x,y** | ..lower.., ..middle.., ..upper.., ..width.. , ..ymin.., ..ymax.
- f + stat_ydensity(kernel = "gaussian", scale = "area") x, y | ..density.., ..scaled.., ..count.., ..n.., ..violinwidth.., ..width..
- **e** + **stat_quantile**(quantiles = c(0.1, 0.9), formula = $y \sim log(x)$, method = "rq") **x, y** | ..quantile.. e + stat_smooth(method = "lm", formula = $y \sim x$, se = T, evel = 0.95) x, y | ..se.., ..x.., ..y.., ..ymin.., ..ymax..

 $e + stat_ecdf(n = 40) x, y \mid ..x., ..y.$

n = 20, geom = "point") $\mathbf{x} \mid ..x.., ..y..$ ggplot() + xlim(-5, 5) + stat_function(fun = dnorm.

- **x, y, sample** | ..sample.., ..theoretical.. **ggplot() + stat_qq(**aes(sample = 1:100)**)**
- **e + stat_sum() x, y, size** | ..n.., ..prop.
- e + stat_summary(fun.data = "mean_cl_boot")
- h + stat_summary_bin(fun = "mean", geom = "bar")
- e + stat_identity()
- e + stat_unique()

Override defaults with scales package

aesthetic. To change a mapping, add a new scale Scales map data values to the visual values of an









GENERAL PURPOSE SCALES

Use with most aesthetics

scale_*_date(date_labels = "%m/%d"),
date_breaks = "2 weeks") - Treat data values as dates. scale_*_datetime() - Treat data values as date times. Same as scale_*_date(). See ?strptime for label formats. **scale_*_manual**(values = c()) - Map discrete values to manually chosen visual ones. scale_*_identity() - Use data values as visual ones. scale_*_binned() - Map continuous values to discrete bins scale_*_discrete() - Map discrete values to visual ones. scale_*_continuous() - Map cont' values to visual ones.

X & Y LOCATION SCALES

scale_x_reverse() - Reverse the direction of the x axis, scale_x_sqrt() - Plot x on square root scale. scale_x_log10() - Plot x on log10 scale. Use with x or y aesthetics (x shown here)

COLOR AND FILL SCALES (DISCRETE)

RColorBrewer::display.brewer.all() n + scale_fill_brewer(palette = "Blues") end = 0.8, na.value = "red") n + scale_fill_grey(start = 0.2, For palette choices:

COLOR AND FILL SCALES (CONTINUOUS)

- o <- c + geom_dotplot(aes(fill = ..x..))
- o + scale_fill_distiller(palette = "Blues")

o + scale_fill_gradient(low="red", high="yellow")

- o + scale_fill_gradient2(low = "red", high = "blue", mid = "white", midpoint = 25)
- o + scale_fill_gradientn(colors = topo.colors(6)) Also: rainbow(), heat.colors(), terrain.colors(), cm.colors(), RColorBrewer::brewer.pal()

SHAPE AND SIZE SCALES

p <- e + geom_point(aes(shape = fl, size = cyl))</pre>

- **\rightarrow** p + scale_shape_manual(values = c(3:7)) p + scale_shape() + scale_size()
- p + scale_ra dius(range = c(1,6))
 p + scale_size_area(max_size = 6) $\Delta \nabla \Diamond \Box \Diamond \Diamond \Diamond$

Coordinate Systems

 $r \leftarrow d + geom_bar()$

The default cartesian coordinate system $\mathbf{r} + \mathbf{coord}_{\mathbf{cartesian}}(\mathbf{xlim} = c(0, 5)) - \mathbf{xlim}, \mathbf{ylim}$

ratio, xlim, ylim - Cartesian coordinates with r + coord_fixed(ratio = 1/2)

ggplot(mpg, aes(y = fl)) + **geom_bar()**Flip cartesian coordinates by switching fixed aspect ratio between x and y units

theta, start, direction - Polar coordinates. r + coord_polar(theta = "x", direction=1)

x and y aesthetic mappings.

Transformed cartesian coordinates. Set xtrans and ytrans to the name of a window function. r + coord_trans(y = "sqrt") - x, y, xlim, ylim

π + coord_quickmap()

 π + coord_map(projection = "ortho", orientation = c(41, -74, 0)) - projection, xlim, ylim Map projections from the mapproj package (mercator (default), azequalarea, lagrange, etc.).

Position Adjustments

Position adjustments determine how to arrange geoms that would otherwise occupy the same space.

s <- ggplot(mpg, aes(fl, fill = drv))



another, normalize height. s + geom_bar(position = "fill") Stack elements on top of one

e + geom_label(position = "nudge") e + geom_point(position = "jitter")
Add random noise to X and Y position of each element to avoid overplotting.

s + geom_bar(position = "stack") Nudge labels away from points. Stack elements on top of one another.

s + geom_bar(position = position_dodge(width = 1)) with manual width and height arguments: Each position adjustment can be recast as a function



r + theme_gray() (default theme). Grey background

Dark for contrast. r + theme_dark() r + theme_minimal(r + theme_void() **Empty theme** Minimal theme.

as axis, legend, panel, and facet properties. r + theme() Customize aspects of the theme such · + theme(panel.background = element_rect(fill = "blue"))

Faceting

discrete variables. values of one or more subplots based on the Facets divide a plot into



t <- ggplot(mpg, aes(cty, hwy)) + geom_point(

t + facet_grid(cols = vars(fl)) Facet into columns based on fl.

Facet into rows based on year. t + facet_grid(rows = vars(year))

Facet into both rows and columns t + facet_grid(rows = vars(year), cols = vars(fl))

Wrap facets into a rectangular layout. t + facet_wrap(vars(fl))

Set scales to let axis limits vary across facets

x and y axis limits adjust to individual facets:
"free_x" - x axis limits adjust "free_y" - y axis limits adjust

Set **labeller** to adjust facet label:

t + facet_grid(rows = vars(fl), t + facet_grid(cols = vars(fl), labeller = label_both) α^c α^d labeller = label_bquote(alpha ^ .(fl))) α^e fl:p **≓**:

<u>abel</u> ls and Legends

Use **labs()** to label the elements of your plot.

t + labs(x = "New x axis label", y = "New y axis label", subtitle = "Add a subtitle below title" **title** ="Add a title above the plot"

Places a geom with manually selected aesthetics. **t + annotate**(geom = "text", x = 8, y = 9, label = "A")

n + guides(fill = "none") Set legend type for each aesthetic: colorbar, legend, or none (no legend). p + guides(x = guide_axis(n.dodge = 2)) Avoid crowded
or overlapping labels with guide_axis(n.dodge or angle).

n + scale_fill_discrete(name = "Title", labels = c("A", "B", "C", "D", "E")) Set legend title and labels with a scale function. n + theme(legend.position = "bottom")
Place legend at "bottom", "top", "left", or "right".

Zooming

Without clipping (preferred):

 $t + coord_cartesian(xlim = c(0, 100), ylim = c(10, 20))$ With clipping (removes unseen data points):

 $t + scale_x_continuous(limits = c(0, 100)) + scale_y_continuous(limits = c(0, 100))$ t + xlim(0, 100) + ylim(10, 20)

