

# Przetwarzanie i wizualizacja danych

prof. UAM dr hab. Tomasz Górecki



tomasz.gorecki@amu.edu.pl

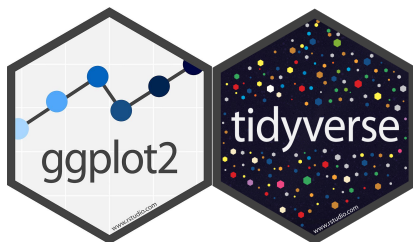
Zakład Statystyki Matematycznej i Analizy Danych  
Wydział Matematyki i Informatyki  
Uniwersytet im. Adama Mickiewicza w Poznaniu



## Biblioteka ggplot2

Pakiet **ggplot2** jest jednym z najbardziej zaawansowanych narzędzi do tworzenia wykresów statystycznych. Oznacza to, że konstrukcja pakietu jest na tyle elastyczna, że można z nim wykonać praktycznie każdą grafikę statystyczną.

-  Biecek, P. (2016). *Odkrywać! Ujawniać! Objaśniać!* Politechnika Warszawska.
-  Wickham, H. (2010). *ggplot2: Elegant graphics for data analysis*. Springer.



# ggplot2 – karta pomocy (1)

## Data Visualization with ggplot2 :: CHEAT SHEET



### 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.



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



Complete the template below to build a graph.

```
ggplot(data = <DATA>) +
  <GEOM_FUNCTION> (mapping = aes(<MAPPINGS>),
  stat = <STAT>, position = <POSITION>),
  <COORDINATE_FUNCTION> +
  <FACE_FUNCTION> +
  <SCALE_FUNCTION> +
  <THEME_FUNCTION>
```

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

ggplot(x = cty, y = hwy, data = mpg, geom = "point") Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

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.

### 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(eats, aes(x = long, y = lat))

a + geom_blank()
  (useful for expanding limits)

b + geom_curve(aes(yend = lat + 1,
  xend = long - 1, curvature = 2)) - x, vend, y, yend,
  alpha, angle, color, curvature, linetype, size

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

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

b + geom_rect(aes(xmin = long, ymin = lat, xmax =
  long + 1, ymax = lat + 1)) - x, xmin, ymin, 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 = 1.155, 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(fill))
  x, y, alpha, color, fill, linetype, size, weight
```

#### TWO VARIABLES

```
continuous x, continuous y
e <- ggplot(mpg, aes(cty, hwy))

e + geom_label(aes(label = cty, nudje_x = 1,
  nudje_y = 1, check_overlap = TRUE)) - x, y, label,
  alpha, angle, color, family, fontface, hjust,
  linewidth, size, vjust

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

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), nudje_x = 1,
  nudje_y = 1, check_overlap = TRUE)) - x, y, label,
  alpha, angle, color, family, fontface, hjust,
  linewidth, size, vjust
```

#### discrete x, continuous y

```
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(binwidth = "y", stackdir =
  "collor") - x, y, alpha, color, fill, group

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

#### discrete x, discrete y

```
g <- ggplot(diamonds, aes(carat, color))

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

#### THREE VARIABLES

```
seals52 ~ with(seals, qrt(delta_log2 * delta_lat2)) - ggplot(seals, aes(long, lat))

l + geom_contour(aes(z))
  x, y, z, alpha, colour, group, linetype, size, weight

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, weight
```

#### 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_density2d()
  x, y, alpha, colour, group, linetype, size

h + geom_hex()
  x, y, alpha, colour, 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, fill, linetype, size
```

#### visualizing error

```
df <- data.frame(murder = c("A", "B"), fit = 4.5, se = 1.2)
j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))

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

j + geom_errorbar()
  x, ymax, ymin, 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
```

#### maps

```
data = data.frame(murder = USAREsts$Murder,
  state = tolower(rownames(USAREsts)))
map = map_data("state")
k <- ggplot(data, aes(fill = murder))

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



# ggplot2 – karta pomocy (2)

## 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.



```

+ stat_bin(binwidth = 1, origin = 10)
x, y | count, ..ncount.., density, ..ndensity..
+ stat_count(width = 1) x, y | count, ..prop..
+ stat_density2d(ju = 1, kernel = "gaussian")
x, y | count, ..density.., scaled
+ stat_bin_2d(bins = 30, drop = T)
x, y, fill | count, ..density..
+ stat_bin_2d(bins = 30) x, y, fill | count, ..density..
+ stat_density_2d(contour = TRUE, n = 100)
x, y, color, size | level
+ stat_ellipse(level = 0.95, segments = 51, type = "r")
+ stat_contour(aes(z = z), x, y, z, order = 1, level = 1)
x, y, z | ..xval.., ..yval.., ..zval..
+ stat_summary_bin(aes(z = z), bins = 30, fun = max)
x, y, z | ..xval.., ..yval.., ..zval..
+ stat_summary_2d(aes(z = z), bins = 30, fun = mean)
x, y, z, fill | ..xval.., ..yval.., ..zval..
+ stat_boxplot(coef = 1.5) x, y | lower, ..midlow.., upper, ..width.., ..ymidmax..
+ stat_ydensity(kernel = "gaussian", scale = "area") x, y | density, ..scaled.., count, ..n.., violwidth, ..width..
+ stat_ecdf(n = 40) x, y | ..xval.., ..yval..
+ stat_quantile(quantiles = c(0.1, 0.9), formula = y ~ log(x), method = "k") x, y | quartile
+ stat_smooth(method = "lm", formula = y ~ x, se = T, level = 0.95) x, y | ..xval.., ..yval.., ..ymidmax..
ggplot() + stat_function(aes(x = -3.3), n = 99, fun = dnorm, arg = list(d = 0.5)) x | ..xval..
+ stat_identity(na.rm = TRUE)
ggplot() + stat_qq(aes(sample = 1:100), dist = qt, dparam = list(d = 5)) x, y | ..sample.., ..theoretical..
+ stat_sum(x, y, size) | ..n.., ..prop..
+ stat_summary(fun.data = "mean_cl_boot")
h = stat_summary_bin(fun.y = "mean", geom = "bar")
+ stat_unique()
    
```

## Scales

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



**GENERAL PURPOSE SCALES**  
Use with most aesthetics  
`scale_*_continuous()` - map `cont` values to visual ones  
`scale_*_discrete()` - map discrete values to visual ones  
`scale_*_identity()` - use data values as visual ones  
`scale_*_manual(values = c())` - map discrete values to manually chosen visual ones  
`scale_*_date(date_labels = "mm/yy/d")`, `date_breaks = "2 weeks"` - treat data values as dates.  
`scale_*_datetime()` - treat data x values as date times. Use same arguments as `scale_x_date()`. See `?sprintf` for label formats.

**X & Y LOCATION SCALES**  
Use with `x` or `y` aesthetics (x shown here)  
`scale_x_log10()` - Plot on log10 scale  
`scale_x_reverse()` - Reverse direction of x axis  
`scale_x_sqrt()` - Plot x on square root scale

**COLOR AND FILL SCALES (DISCRETE)**  
`n = d - geom_bar(aes(fill = R))`  
`scale_fill_brewer(palette = "blues")`  
For palette choices: `RColorBrewer::brewer.pal()`  
`scale_fill_grey(stat = 0.2, end = 0.8, na.value = "red")`

**COLOR AND FILL SCALES (CONTINUOUS)**  
`o = c - geom_dotplot(aes(fill = ..x..))`  
`scale_fill_distiller(palette = "blues")`  
`scale_fill_gradient(low = "red", high = "yellow")`  
`scale_fill_gradient2(low = "red", high = "blue", mid = "white", midpoint = 35)`  
`scale_fill_gradientn(colors = topo.colors(8))`  
Also: `rainbow()`, `heat.colors()`, `terrain.colors()`, `cm.colors()`, `RColorBrewer::brewer.pal()`

**SHAPE AND SIZE SCALES**  
`p = c - geom_point(aes(shape = fl, size = cyl))`  
`scale_shape_manual(values = c(1, 7))`  
`scale_shape_size(values = c(1, 8))`  
`scale_size_area(max_size = 4)`

## Coordinate Systems

```

r = d - geom_bar()
+ coord_cartesian(lim = c(0, 5))
xlim, ylim
The default cartesian coordinate system
+ coord_fixed(ratio = 1/2)
ratio, xlim, ylim
Cartesian coordinates with fixed aspect ratio between x and y units
+ coord_flip()
xlim, ylim
Flipped Cartesian coordinates
+ coord_polar(theta = "x", direction = 1)
theta, start, direction
Polar coordinates
+ coord_trans(proj = "sqrt")
xlim, ylim, size, alpha
Coordinates. Set xlims and ylims to the name of window function.
+ coord_quickmap()
+ coord_map(proj = "ortho",
direction = "c", 14, projection, orientation,
xlim, ylim)
Map projections from the mapping package
(Indicator: false, aboquata, bgrange, etc.)
    
```

## Position Adjustments

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

```

s = ggplot(mpg, aes(fl, drv))
+ geom_bar(position = "dodge")
Arrange elements side by side
+ geom_bar(position = "stack")
Stack elements on top of one another
+ geom_point(position = "jitter")
Add random noise to x and y position of each element to avoid overlapping
+ geom_label(position = "nudge")
Nudge labels away from points
+ geom_bar(position = "stack")
Stack elements on top of one another
    
```

Each position adjustment can be recast as a function with **manual width and height** arguments

```

+ geom_bar(position = position_dodge(width = 1))
    
```

## Themes

```

+ theme_bw()
White background with gray borders
+ theme_gray()
Gray background (default theme)
+ theme_dark()
Dark for contrast
+ theme_classic()
+ theme_light()
+ theme_minimal()
Minimal theme
+ theme_void()
Empty theme
    
```

## Faceting

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



```

t = ggplot(mpg, aes(cty, hwy)) + geom_point()
+ facet_grid(~ fl)
Facet into columns based on fl
+ facet_grid(~ cyl)
Facet into both rows and columns
+ facet_wrap(~ fl)
Wrap facets into a rectangular layout
+ scale_x_continuous(limits = c(10, 100))
Scale x-axis limits
+ scale_y_continuous(limits = c(10, 100))
Scale y-axis limits
+ free_x()
Free x-axis limits
+ free_y()
Free y-axis limits
+ label_both()
Label both axes
+ label_x()
Label x-axis
+ label_y()
Label y-axis
+ label_parsed()
Parse labels
    
```

```

+ facet_grid(~ fl, labeller = label_both)
+ facet_grid(~ fl, labeller = label_both,
direction = "c", 14, projection, orientation,
xlim, ylim)
+ facet_grid(~ fl, labeller = label_parsed)
    
```

## Labels

```

+ label(x = "New x axis label", y = "New y axis label",
title = "Add a title above the plot",
subtitle = "Add a subtitle below the plot",
caption = "Add a caption below plot",
aes())
+ annotate(geom = "text", x = 8, y = 9, label = "X")
+ geom_text()
Manual values for geom's aesthetics
    
```

## Legends

```

+ theme(legend.position = "bottom")
Place legend at "bottom", "top", "left", or "right"
+ guides(fill = "none")
Set legend type for each aesthetic: color, legend, or none (no legend)
+ scale_fill_discrete(name = "Title")
Set legend title and labels with a scale function.
    
```

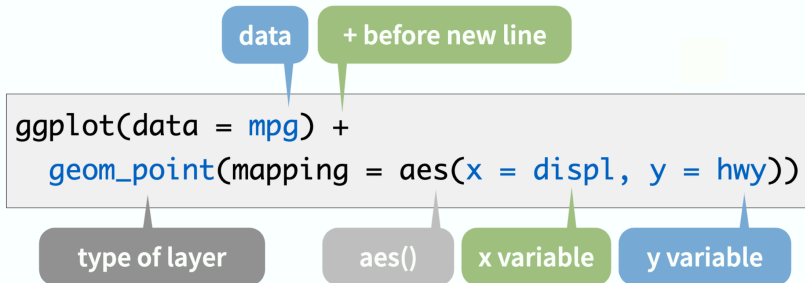
## Zooming

```

+ coord_cartesian(lim = c(10, 20))
xlim = c(10, 100), ylim = c(10, 200)
Without clipping (preferred)
+ coord_cartesian(lim = c(10, 20))
xlim = c(10, 100), ylim = c(10, 200)
With clipping (removes unseen data points)
+ xlim(10, 100) + ylim(10, 20)
+ scale_x_continuous(limits = c(10, 100))
+ scale_y_continuous(limits = c(10, 100))
    
```



## Biblioteka ggplot2 – podstawy



- `ggplot(data, aes(x, y, color, shape, size, label))`,
- `geom_point`, `geom_text(hjust, vjust)`, `geom_line`, `geom_ribbon`,
- `xlim`, `ylim`,
- `geom_smooth`, `geom_boxplot`, `geom_hist`, `geom_bar`

## Biblioteka ggplot2 – szczegóły

- ggtitle, xlab, ylab,
- theme,
- scale\_XXX\_yyy,
- grid.layout, grid.arrange,
- theme\_bw.

# ggplot2 – karta pomocy (estetyki)

## ggplot2 aesthetics cheat sheet

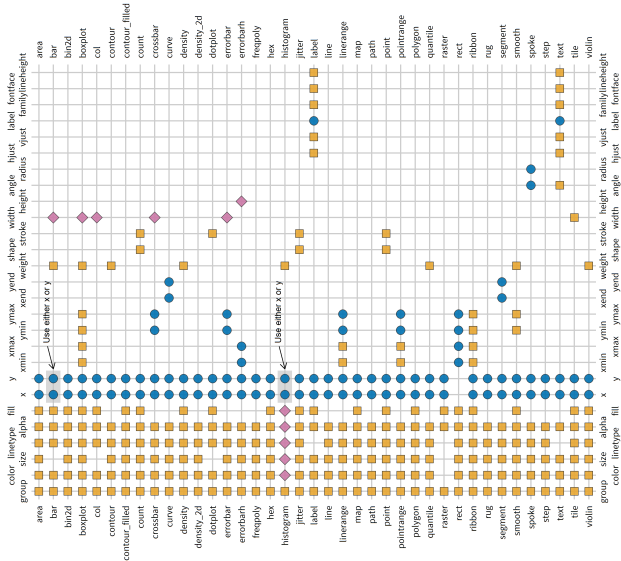
Use this table to find the right aesthetics for your geoms:

**Aesthetics that usually must be mapped to the data: use inside aes()**

**Aesthetics that can be mapped to the data: use in or outside aes()**

**Aesthetics that cannot be mapped to the data: use outside aes()**

e.g., `ggplot(mpg, aes(x = class, y = displ)) + geom_col(fill = class, width = .9)`



● usually must be inside aes()    ■ can be inside aes()    ◆ must be outside aes()

idea and design: Christian Burkhardt  
 design advice: Iida Aarnio

# Mapy

- Mapy wykonane za pomocą pakietów ggplot2 i sf.
- Mapy Google.