

# *Analiza danych*




prof. UAM dr hab. Tomasz Górecki

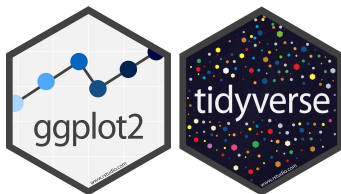
tomasz.gorecki@amu.edu.pl  
<http://drizzt.home.amu.edu.pl>

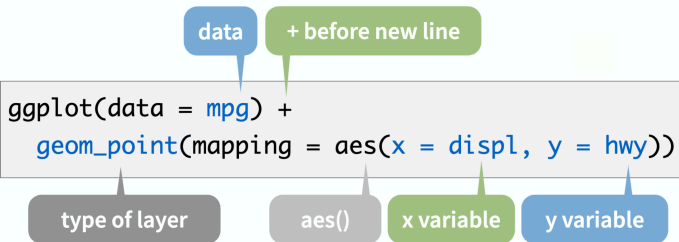
Zakład Statystyki Matematycznej i Analizy Danych  
Wydział Matematyki i Informatyki  
Uniwersytet im. Adama Mickiewicza w Poznaniu



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

-  Biecek, P. (2016). *Odkrywać! Ujawniać! Objaśniać!* Politechnika Warszawska.
-  Healy, K. (2019). *Data visualization. A practical introduction.* Princeton University Press.
-  Wickham, H. (2010). *ggplot2: Elegant graphics for data analysis.* Springer.





<i>Data</i>	<i>{variables of interest}</i>				
<i>Aesthetics</i>	<i>x-axis</i> <i>y-axis</i>	<i>colour</i> <i>fill</i>	<i>size</i> <i>labels</i>	<i>alpha</i> <i>shape</i>	<i>line width</i> <i>line type</i>
<i>Geometries</i>	<i>point</i>	<i>line</i>	<i>histogram</i>	<i>bar</i>	<i>boxplot</i>
<i>Facets</i>	<i>columns</i>	<i>rows</i>			
<i>Statistics</i>	<i>binning</i>	<i>smoothing</i>	<i>descriptive</i>	<i>inferential</i>	
<i>Coordinates</i>	<i>cartesian</i>	<i>fixed</i>	<i>polar</i>	<i>limits</i>	
<i>Themes</i>	<i>non-data ink</i>				

## 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 **x** and **y** locations.



Complete the template below to build a graph.

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

- aesthetic mappings** `data` `geom`
- plot()** `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(mpg, aes(cty, y = hwy))**
- a + geom\_blank()**  
(Useful for expanding limits)
- b + geom\_curve(aes(yend = lat + 1, xend = long - 1, curvature = 2))** `x`, `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))** `small_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\_abline(aes(intercept = lat))**
- b + geom\_vline(aes(xintercept = long))**
- b + geom\_segment(aes(yend = lat - 1, xend = long + 1))**
- b + geom\_spline(aes(angle = 1.1515, 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))**
- d + geom\_bar()**  
`x`, `alpha`, `color`, `fill`, `linetype`, `size`, `weight`

#### TWO VARIABLES

- continuous x, continuous y**
- e = ggplot(mpg, aes(cty, hwy))**
- e + geom\_label(aes(label = cty), nudges\_x = 1, nudges\_y = 1, check\_overlap = TRUE)** `x`, `y`, `label`, `alpha`, `angle`, `color`, `family`, `fontface`, `hjust`, `linetype`, `size`
- 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(styles = "tbl")** `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), nudges\_x = 1, nudges\_y = 1, check\_overlap = TRUE)** `x`, `y`, `label`, `alpha`, `angle`, `color`, `family`, `fontface`, `hjust`, `linetype`, `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 = "center")** `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

- seals5 = with(seals, sqrt(delta\_long^2 + delta\_lat^2))** `ggplot(seals, aes(long, lat))`
- h = geom\_contour(aes(z))**  
`x`, `y`, `z`, `alpha`, `colour`, `group`, `linetype`, `size`, `weight`
- h + geom\_raster(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE)**  
`x`, `y`, `alpha`, `color`, `fill`
- h + geom\_tile(aes(fill = z))** `x`, `y`, `alpha`, `color`, `fill`, `linetype`, `size`, `width`

#### 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`, `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(g = c("A", "B"), fit = 4.5, se = 1.2)**
- df = ggplot(df, aes(g, 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

- data = data.frame(murder = USArrests\$Murder, state = tolower(row.names(USArrests)))**
- map = map\_data("state")**
- 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)**



# ggplot2 - karta pomocy (2)

## Stats

An alternative way to build a layer

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



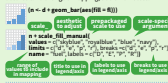
Visualize a stat by changing the default stat of a geom function, `geom_bar(stat="count")` or by using a stat function, `stat_count(geom="bar")`, which calls a default geom to make a layer (equivalent to a geom function). Use `name_` syntax to map stat variables to aesthetics.

```
geom to use | stat function | geom mappings
+ stat_density2d(aes(fill = level),
  geom = "polygon") | variable created by stat
```

```
c + stat_bin(nwidth = 1, origin = 10)
x, y | count, ...ncount, ...density, ...ndensity.
+ stat_count(width = 1) x, y | count, ...prop.
+ stat_density(x = x, y = y, kernel = "gaussian")
x, y | count, ...density, ...scaled.
+ stat_bin_2d(bins = 30, drop = T)
x, y, fill | count, ...density.
+ stat_bin_hex(bins = 30, x, y, fill | count, ...density.
+ stat_density_2d(contour = TRUE, n = 100)
x, y, color, size | level.
+ stat_ellipse(level = 0.55, segments = 51, type = "r")
+ stat_contour(aes(z = z), x, y, order | level.
+ stat_summary_hex(aes(z = z), bins = 30, fun = max)
x, y, fill | value.
+ stat_summary_2d(aes(z = z), bins = 30, fun = mean)
x, y, z, fill | value.
+ stat_boxplot(coeff = 1.5, x, y | lower,
  middle, upper, width, ymin, ymax.
+ stat_ydensity(kernel = "gaussian", scaled = TRUE, width = 10)
density, x, count, ...count, ...width, width.
+ stat_ecdf(n = 40) x, y | x, y.
+ stat_quantile(quantiles = c(0.1, 0.5), formula = y ~
  log(x), method = "r") x, y | quantile.
+ stat_smooth(method = "lm", formula = y ~ x, se = T,
  level = 0.5) x, y | se, x, y, ymin, ymax.
+ stat_identity(n.m = TRUE)
ggplot() + stat_function(aes(x = -3.3), n = 99, fun =
  dnorm, arg = list(0.5)) | x, y.
+ stat_summary(fun = TRUE)
ggplot() + stat_qq(aes(sample = 1:100), dist = qt,
  p.param = list(d = 5)) sample, x, y | sample, theoretical.
+ stat_samel(x, y, size | count, ...prop.
+ stat_summary(fun.data = "mean_cdot_box")
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()` - treat data x values as date times. Use same arguments as `scale_x_date()`. See "time" for format lists.  
`scale_*_width()` - treat data values as dates.`scale_*_datetime()` - treat data x values as date times. Use same arguments as `scale_x_datetime()`. See "time" for format lists.

### 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)  
`scale_*_discrete(values = c())`  
`scale_*_fill_brewer(palette = "Blues")`  
`scale_*_fill_brewer(palette = "Blues")`  
`scale_*_fill_brewer(palette = "Blues")`  
`scale_*_fill_grey(start = 0.2, end = 0.8, na.value = "red")`

### COLOR AND FILL SCALES (CONTINUOUS)

`scale_*_color_discrete(values = c())`  
`scale_*_fill_discrete(values = c())`  
`scale_*_color_gradient(low = "red", high = "yellow")`  
`scale_*_color_gradient(low = "red", high = "blue", mid = "white", midpoint = 25)`  
`scale_*_color_gradient(colors = topo.colors(6))`  
Also: `rainbow()`, `heat.colors()`, `terrain.colors()`, `cm.colors()`, `RColorBrewer::brewer.pal()`

### SHAPE AND SIZE SCALES

`scale_*_geom_point(aes(shape = fl, size = cyl))`  
`scale_*_shape()`  
`scale_*_size()`  
`scale_*_shape_manual(values = c())`  
`scale_*_size_manual(values = c())`  
`scale_*_radius(range = c(1, 6))`  
`scale_*_size_area(max_size = 6)`

## Coordinate Systems

`r <- d + geom_bar()`  
`r + coord_cartesian(dim = c(0, 5))`  
`lim, ylim`  
The default cartesian coordinate system  
`r + coord_fixed(ratio = 1/2)`  
`ratio, xlim, ylim`  
Cartesian coordinates with fixed aspect ratio between axes  
`r + coord_flip()`  
`lim, ylim`  
Flipped Cartesian coordinates  
`r + coord_polar(data = "x", direction = 1)`  
`break, coord`  
Polar coordinates  
`r + coord_trans(trans = "sqrt")`  
`trans, ylim, lim, ylim`  
Map projections from the `maptools` package (Mercator, Yoffm, Albers, etc.)  
`r + coord_quickmap()`  
`projection = "ortho"`  
Orthographic projection, orientation, `lim, ylim`  
Map projections from the `maptools` package (Mercator, Yoffm, Albers, etc.)

## Position Adjustments

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

`s + ggplot(mpg, aes(fl, drv))`  
`s + geom_bar(position = "dodge")`  
Arrange elements side by side  
`s + geom_bar(position = "fill")`  
Stack elements on top of one another, normalize height  
`s + geom_point(position = "jitter")`  
Add random noise to x and y position of each element to avoid overlapping  
`s + geom_label(position = "nudge")`  
Nudge labels away from geoms  
`s + 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  
`s + geom_bar(position = position_dodge(width = 1))`

## Themes

`r + theme_bw()`  
White background with grid lines  
`r + theme_gray()`  
Gray background (default theme)  
`r + theme_dark()`  
Dark for contrast  
`r + theme_classic()`  
White background with grid lines  
`r + theme_light()`  
White background with grid lines  
`r + theme_minimal()`  
Minimal themes  
`r + 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()`  
`t + facet_grid(~ fl)`  
Facet into columns based on fl  
`t + facet_grid(~ drv)`  
Facet into rows based on drv  
`t + facet_grid(~ cyl)`  
Facet into both rows and columns  
`t + facet_wrap(~ fl)`  
Wrap facets into a rectangular layout

Scale to let axis limits vary across facets  
`t + facet_grid(drv ~ fl, scales = "free")`  
x and y axis limits adjust to individual facets  
`free_x` - x axis limits adjust  
`free_y` - y axis limits adjust  
Use `vars()` to specify multiple variables  
Set `labeler` to adjust facet labels  
`t + facet_grid(~ fl, labeler = label_both)`  
`fl` `drv` `fl` `drv`  
`t + facet_grid(~ fl, labeler = label_bquote(alpha = .))`  
`fl` `drv` `fl` `drv`  
`t + facet_grid(~ fl, labeler = label_parsed)`  
`fl` `drv` `fl` `drv`

## Labels

`label` - x = "New x axis label", y = "New y axis label", title = "Add a title above the plot", subtitle = "Add a subtitle below title", caption = "Add a caption below plot", text = "New text label", `annotate` - `geom = "text", x = 8, y = 9, label = "A"`  
`geom to place` `manual values for geom's aesthetics`

## Legends

`n + theme(legend.position = "bottom")`  
Place legend at "bottom", "top", "left", or "right"  
`n + guides(fill = "none")`  
Set legend type for each aesthetic: colorbar, legend, or none for legend  
`n + scale_fill_discrete(aes("Title", labels = c("A", "B", "C", "D")))`  
Set legend title and labels with a scale function.

## Zooming

Without clipping (preferred)  
`t + coord_cartesian(lim = c(0, 100), ylim = c(10, 20))`  
With clipping (removes unseen data points)  
`t + xlim(0, 100) + ylim(10, 20)`  
`t + scale_x_continuous(limits = c(0, 100)) + scale_y_continuous(limits = c(0, 100))`



## ggplot2 aesthetics cheat sheet

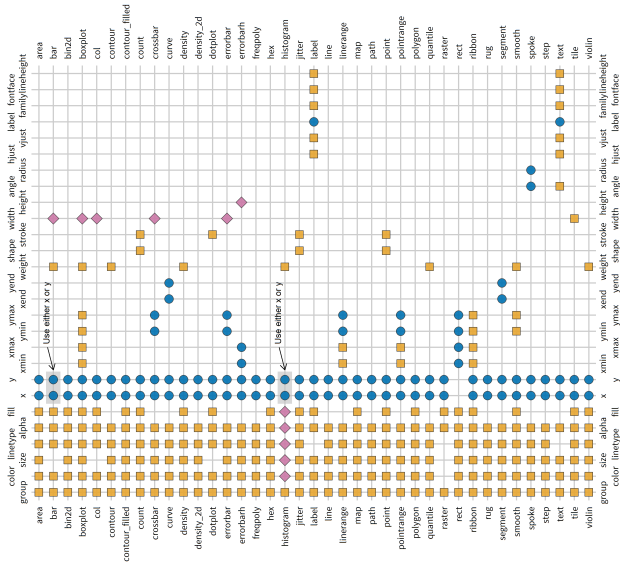
Use this table to find the right aesthetics for your geoms:

**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(aes(fill = class), width = .9)`



- usually must be inside aes()
- can be inside aes()
- ◆ must be outside aes()

