

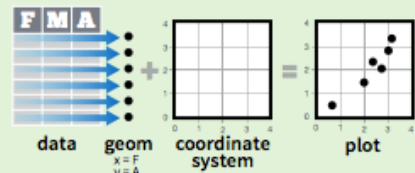
# Trực quan hóa số liệu với ggplot2

Cheat Sheet

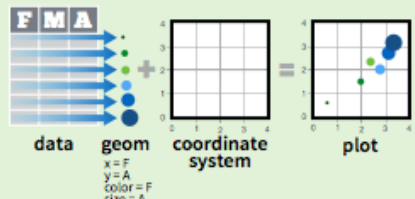


## Kiến thức cơ bản

**ggplot2** dựa trên khái niệm “**ngữ pháp của biểu đồ**”, trong đó tất cả các biểu đồ đều có thể được xây dựng từ những thành phần giống nhau: **data** - tập dữ liệu, **geoms** - mô tả cách thức thể hiện dữ liệu, và một hệ tọa độ (**coordinate**)



Để hiển thị các điểm dữ liệu, cần phải sắp xếp các biến trong dữ liệu với các thuộc tính hình học (geom) như kích cỡ, màu sắc, trục tọa độ x & y



Vẽ biểu đồ với ggplot() hoặc qplot()

Sắp xếp thuộc tính aes

Dữ liệu

geom

**qplot**(x = cty, y = hwy, color = cyl, data = mpg, geom = "point")  
Tạo một biểu đồ hoàn chỉnh với dữ liệu, geom & thuộc tính cho trước. Hỗ trợ nhiều chế độ mặc định

**ggplot**(data = mpg, aes(x = cty, y = hwy))  
Thêm các lớp (layer) vào biểu đồ đã tạo, hỗ trợ nhiều loại biểu đồ hơn qplot().

Dữ liệu

```
ggplot(mpg, aes(hwy, cty)) +  
  geom_point(aes(color = cyl)) +  
  geom_smooth(method = "lm") +  
  coord_cartesian() +  
  scale_color_gradient() +  
  theme_bw()
```

Thêm các lớp với dấu +

Lớp (layer) =  
geom + default  
stat + các thuộc  
tính khác

Các thành phần khác

Thêm lớp mới trong biểu đồ với hàm **geom\_\*()** hoặc **stat\_\*()**. Mỗi hàm sẽ xác định một "geom", là một nhóm các thuộc tính hình học, các tính toán mặc định và sự sắp xếp vị trí trong biểu đồ.

**last\_plot()**

Trả về biểu đồ đã tạo gần nhất

**ggsave**("plot.png", width = 5, height = 5)  
Lưu biểu đồ đã tạo gần nhất với kích thước 5' x 5', lưu với tên "plot.png" tại thư mục làm việc

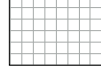
## Geoms

Sử dụng geom để biểu diễn các điểm dữ liệu, sử dụng các thuộc tính của aes để biểu diễn các biến. Mỗi hàm sẽ tạo ra một lớp

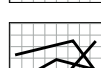
### Các thuộc tính hình học cơ bản

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

**a + geom\_blank()**



**a + geom\_curve**(aes(yend = lat + delta\_lat, xend = long + delta\_long, curvature = z))  
x, xend, y, yend, alpha, angle, color, curvature, linetype, size



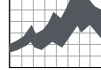
**b + geom\_path**(lineend="butt", linejoin="round", linemitre=1)  
x, y, alpha, color, group, linetype, size



**b + geom\_polygon**(aes(group = group))  
x, y, alpha, color, fill, group, linetype, size



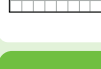
**a + geom\_rect**(aes(xmin = long, ymin = lat, xmax = long + delta\_long, ymax = lat + delta\_lat))  
xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size



**b + geom\_ribbon**(aes(ymin = unemploy - 900, ymax = unemploy + 900))  
x, ymax, ymin, alpha, color, fill, group, linetype, size



**a + geom\_segment**(aes(yend = lat + delta\_lat, xend = long + delta\_long))  
x, xend, y, yend, alpha, color, linetype, size



### Một biến

**Biến liên tục**

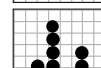
```
c <- ggplot(mpg, aes(hwy))
```



**c + geom\_area**(stat = "bin")  
x, y, alpha, color, fill, linetype, size, weight



**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

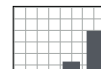


**c + geom\_histogram**(binwidth = 5)  
x, y, alpha, color, fill, linetype, size, weight  
**a + geom\_histogram**(aes(y = ..density..))



**Biến rời rạc**

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



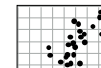
**d + geom\_bar**()  
x, alpha, color, fill, linetype, size, weight

### Biểu đồ hai biến

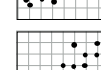
**Biến X liên tục, biến Y liên tục**  
**e <- ggplot(mpg, aes(cty, hwy))**



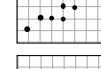
**e + geom\_label**(aes(label = cty), nudge\_x = 1, nudge\_y = 1, check\_overlap = TRUE)  
x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust



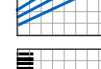
**e + geom\_jitter**(height = 2, width = 2)  
x, y, alpha, color, fill, shape, size



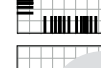
**e + geom\_point**()  
x, y, alpha, color, fill, shape, size, stroke



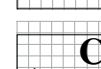
**e + geom\_quantile**()  
x, y, alpha, color, group, linetype, size, weight



**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



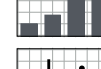
**e + geom\_text**(aes(label = cty), nudge\_x = 1, nudge\_y = 1, check\_overlap = TRUE)  
x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust



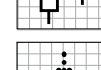
**Biến X rời rạc, biến Y liên tục**  
**f <- ggplot(mpg, aes(class, hwy))**



**f + geom\_bar**(stat = "identity")  
x, y, alpha, color, fill, linetype, size, weight



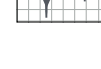
**f + geom\_boxplot**()  
x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight



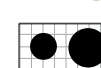
**f + geom\_dotplot**(binaxis = "y", stackdir = "center")  
x, y, alpha, color, fill, group



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



**Biến X rời rạc, biến Y rời rạc**  
**g <- ggplot(diamonds, aes(cut, color))**



**g + geom\_count**()  
x, y, alpha, color, fill, shape, size, stroke



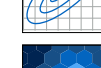
**Hai biến phân phối liên tục**  
**h <- ggplot(diamonds, aes(carat, price))**



**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

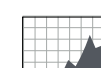


**h + geom\_hex**()  
x, y, alpha, colour, fill, size



**Hàm liên tục**

```
i <- ggplot(economics, aes(date, unemploy))
```



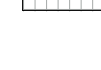
**i + geom\_area**()  
x, y, alpha, color, fill, linetype, size



**i + geom\_line**()  
x, y, alpha, color, group, linetype, size

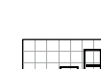


**i + geom\_step**(direction = "hv")  
x, y, alpha, color, group, linetype, size

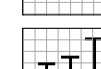


**Trực quan hóa sai số**

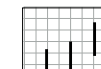
```
df <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)  
j <- ggplot(df, aes(grp, fit, ymin = fit-se, ymax = fit+se))
```



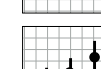
**j + geom\_crossbar**(fatten = 2)  
x, y, ymax, ymin, alpha, color, fill, group, linetype, size



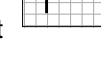
**j + geom\_errorbar**()  
x, ymax, ymin, alpha, color, group, linetype, size, width (also geom\_errorbarh())



**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



**Bản đồ**

```
data <- data.frame(murder = USArrests$Murder,  
  state = tolower(rownames(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

### Ba biến

```
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, 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, width

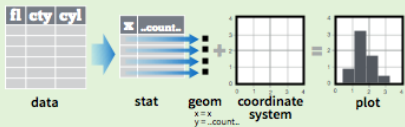




## Stats – cách thức khác để tạo biểu đồ

Một số biểu đồ hiển thị dữ liệu đã được biến đổi. Sử dụng **stat** để lựa chọn hình thức biến đổi dữ liệu, VD.

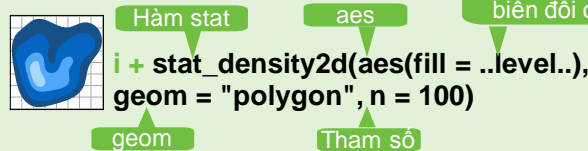
**a + geom\_bar(stat = "count")**



Mỗi **stat** sẽ tạo thêm các biến mới ứng với các thuộc tính hình học. Các biến này sử dụng cấu trúc thông thường **..name..**.

Hàm **stat** và **geom** đều kết hợp một stat với một geom để tạo một lớp (layer) mới, VD.

**stat\_count(geom="bar")** cho ra kết quả tương tự như **geom\_bar(stat="count")**



**c + stat\_bin(binwidth = 1, origin = 10)**  
x, y | ..count.., ..ncount.., ..density.., ..ndensity..  
**c + stat\_count(width = 1)**  
x, y, | ..count.., ..prop..  
**c + stat\_density(adjust = 1, kernel = "gaussian")**  
x, y, | ..count.., ..density.., ..scaled..

**e + stat\_bin\_2d(bins = 30, drop = TRUE)**  
x, y, fill | ..count.., ..density..  
**e + stat\_bin\_hex(bins = 30)**  
x, y, fill | ..count.., ..density..  
**e + stat\_density\_2d(contour = TRUE, n = 100)**  
x, y, color, size | ..level..  
**e + stat\_ellipse(level = 0.95, segments = 51, type = "t")**

**l + stat\_contour(aes(z = z))**  
x, y, z, order | ..level..  
**l + stat\_summary\_hex(aes(z = z), bins = 30, fun = mean)**  
x, y, z, fill | ..value..  
**l + stat\_summary\_2d(aes(z = z), bins = 30, fun = mean)**  
x, y, z, fill | ..value..

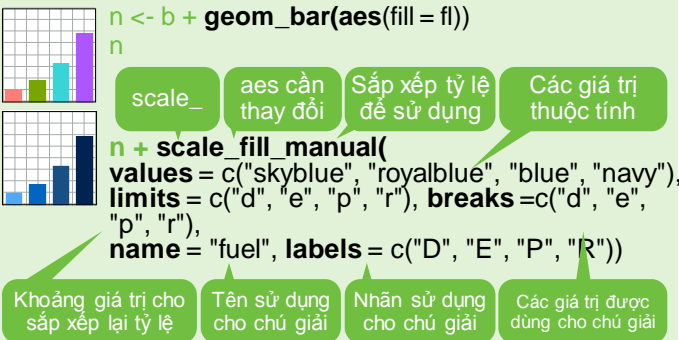
**f + stat\_boxplot(coef = 1.5)**  
x, y | ..lower.., ..middle.., ..upper.., ..width.., ..ymin.., ..ymax..  
**f + stat\_ydensity(adjust = 1, kernel = "gaussian", scale = "area")**  
x, y | ..density.., ..scaled.., ..count.., ..n.., ..violinwidth.., ..width..

**e + stat\_ecdf(n = 40)**  
x, y | ..x.., ..y..  
**e + stat\_quantile(quantiles = c(0.25, 0.5, 0.75), formula = y ~ log(x), method = "rq")**  
x, y | ..quantile..  
**e + stat\_smooth(method = "auto", formula = y ~ x, se = TRUE, n = 80, fullrange = FALSE, level = 0.95)**  
x, y | ..se.., ..x.., ..y.., ..ymin.., ..ymax..

**ggplot() + stat\_function(aes(x = -3:3), fun = dnorm, n = 101, args = list(sd=0.5))**  
x | ..x.., ..y..  
**e + stat\_identity(na.rm = TRUE)**  
**ggplot() + stat\_qq(aes(sample=1:100), distribution = qt, dparams = list(df=5))**  
sample, x, y | ..sample.., ..theoretical..  
**e + stat\_sum()**  
x, y, size | ..n.., ..prop..  
**e + stat\_summary(fun.data = "mean\_cl\_boot")**  
**h + stat\_summary\_bin(fun.y = "mean", geom = "bar")**  
**e + stat\_unique()**

## Scales – Tỷ lệ

**Scales** – Tỷ lệ quy định cách thức biểu đồ sắp xếp dữ liệu với các thuộc tính hình học trên biểu đồ. Để thay đổi cách sắp xếp này, cần thay đổi tỷ lệ.



### Cách sử dụng thường dùng

Sử dụng với các giá trị aes:

alpha, color, fill, linetype, shape, size

**scale\_\*\_continuous()** – Sử dụng cho các biến liên tục  
**scale\_\*\_discrete()** – Sử dụng cho các biến rời rạc  
**scale\_\*\_identity()** – Sử dụng giá trị của tập dữ liệu  
**scale\_\*\_manual(values = c())** – Sắp xếp các biến rời rạc với các giá trị tùy biến

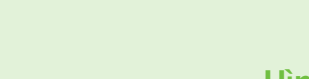
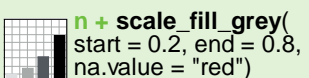
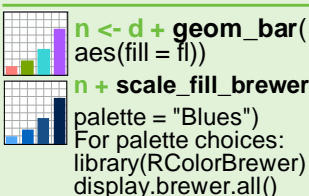
### X and Y location scales

Sử dụng với các thuộc tính của trục x hoặc y (phần dưới đây chỉ mô tả trục hoành x)

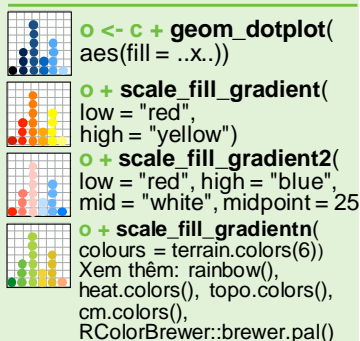
**scale\_x\_date(date\_labels = "%m/%d"), date\_breaks = "2 weeks")** - Coi x như biến ngày tháng. Xem thêm **?strptime** về nhãn (label)  
**scale\_x\_datetime()** - Coi x như biến ngày tháng, sử dụng các tham số như **scale\_x\_date()**  
**scale\_x\_log10()** – Thể hiện x với tỷ lệ log10  
**scale\_x\_reverse()** – Giữ nguyên hướng của trục x  
**scale\_x\_sqrt()** – Thể hiện x với tỷ lệ căn bậc hai

### Màu sắc

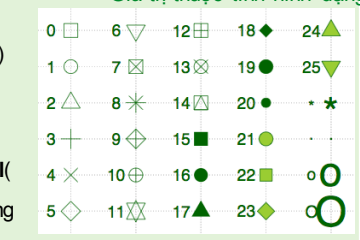
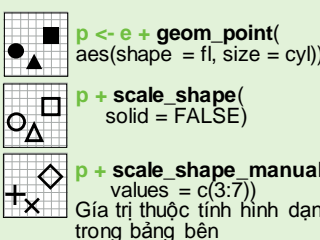
Biến rời rạc



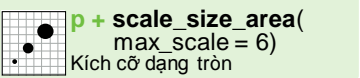
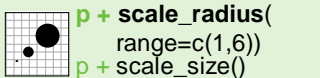
Biến liên tục



### Hình dạng

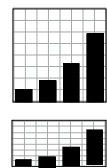


### Kích cỡ



## Coordinate – Hệ tọa độ

**r <- d + geom\_bar()**



**r + coord\_cartesian(xlim = c(0, 5))**

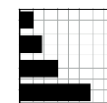
xlim, ylim

Hệ tọa độ Đề-các mặc định

**r + coord\_fixed(ratio = 1/2)**

ratio, xlim, ylim

Hệ tọa độ Đề-các, tỷ lệ x và y cố định



**r + coord\_flip()**

xlim, ylim

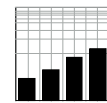
Đổi trục tọa độ



**r + coord\_polar(theta = "x", direction=1)**

theta, start, direction

Hệ tọa độ cực

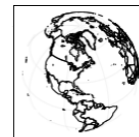


**r + coord\_trans(ytrans = "sqrt")**

xtrans, ytrans, limx, limy

Biến đổi hệ tọa độ Đề-các,

**π + coord\_map(projection = "ortho", orientation=c(41, -74, 0))**  
projection, orientation, xlim, ylim



Sử dụng packages **mapproj** (mercator (mặc định), azequalarea, lagrange,...)

## Điều chỉnh vị trí

Cách thức sắp xếp các thuộc tính hình học (geom) trên biểu đồ

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

**s + geom\_bar(position = "dodge")**

Đặt các giá trị cạnh nhau

**s + geom\_bar(position = "fill")**

Đặt các giá trị chồng lên nhau, thay đổi tỷ lệ theo phần trăm

**e + geom\_point(position = "jitter")**

Thêm các yếu tố ngẫu nhiên (random noise) để tránh chồng lấn các điểm trên biểu đồ

**e + geom\_label(position = "nudge")**

Đặt các nhãn bên cạnh các điểm

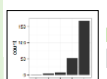
**s + geom\_bar(position = "stack")**

Đặt các giá trị chồng lên nhau

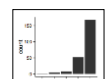
Vị trí trong biểu đồ có thể được thay đổi lại thành một hàm với các tham số của chiều dài và chiều rộng

**s + geom\_bar(position = position\_dodge(width = 1))**

## Themes – Hình nền trong biểu đồ



**r + theme\_bw()**  
Nền trắng



**r + theme\_classic()**  
Nền classic



**r + theme\_gray()**  
Nền xám (theme mặc định)



**r + theme\_minimal()**  
Nền minimal



**r + theme\_dark()**  
Nền tối



**r + theme\_void()**  
Để trống hình nền

## Faceting – Chia nhỏ biểu đồ

Chia nhỏ biểu đồ dựa trên giá trị của một hoặc nhiều biến rời rạc

**t <- ggplot(mpg, aes(cty, hwy)) + geom\_point()**



**t + facet\_grid(. ~ fl)**  
Cột chứa biến fl



**t + facet\_grid(year ~ .)**  
Hàng chứa biến year



**t + facet\_grid(year ~ fl)**  
Chia nhỏ biểu đồ theo cả hàng và cột



**t + facet\_wrap(~ fl)**  
Tự động sắp xếp biểu đồ

Quy định tỷ lệ để giới hạn các trục của biểu đồ khi sử dụng facet

**t + facet\_grid(drv ~ fl, scales = "free")**

Giới hạn trục x & y theo từng biểu đồ

- "free\_x" – Tự động điều chỉnh giới hạn trục x
- "free\_y" – Tự động điều chỉnh giới hạn trục y

Đặt nhãn, tiêu đề cho các biểu đồ khi dùng facet

**t + facet\_grid(. ~ fl, labeller = label\_both)**

**t + facet\_grid(fl ~ ., labeller = label\_bquote(alpha ^ .(fl)))**

**t + facet\_grid(. ~ fl, labeller = label\_parsed)**

## Labels – Tiêu đề & nhãn

**t + ggtitle("New Plot Title")**

Thêm tên biểu đồ

**t + xlab("New X label")**

Thay đổi tên trục x

**t + ylab("New Y label")**

Thay đổi tên trục y

**t + labs(title = "New title", x = "New x", y = "New y")**

Thay đổi tên biểu đồ và các trục x, y

## Chú giải

**n + theme(legend.position = "bottom")**

Thay đổi vị trí chú giải: "up", "bottom", "right", "left"

**n + guides(fill = "none")**

Quy định chú giải cho mỗi thuộc tính: colorbar, legend, hoặc "none" (không để chú giải)

**n + scale\_fill\_discrete(name = "Title", labels = c("A", "B", "C", "D", "E"))**

Sử dụng hàm tỷ lệ (scale) cho tiêu đề & nhãn trong chú giải

## Zooming – Phóng to biểu đồ



Không thay đổi dữ liệu (Nên dùng)

**t + coord\_cartesian(xlim = c(0, 100), ylim = c(10, 20))**

Thay đổi dữ liệu

(Loại bỏ các dữ liệu ngoài vùng phân tích)



**t + xlim(0, 100) + ylim(10, 20)**  
**t + scale\_x\_continuous(limits = c(0, 100))**  
**+ scale\_y\_continuous(limits = c(0, 100))**