

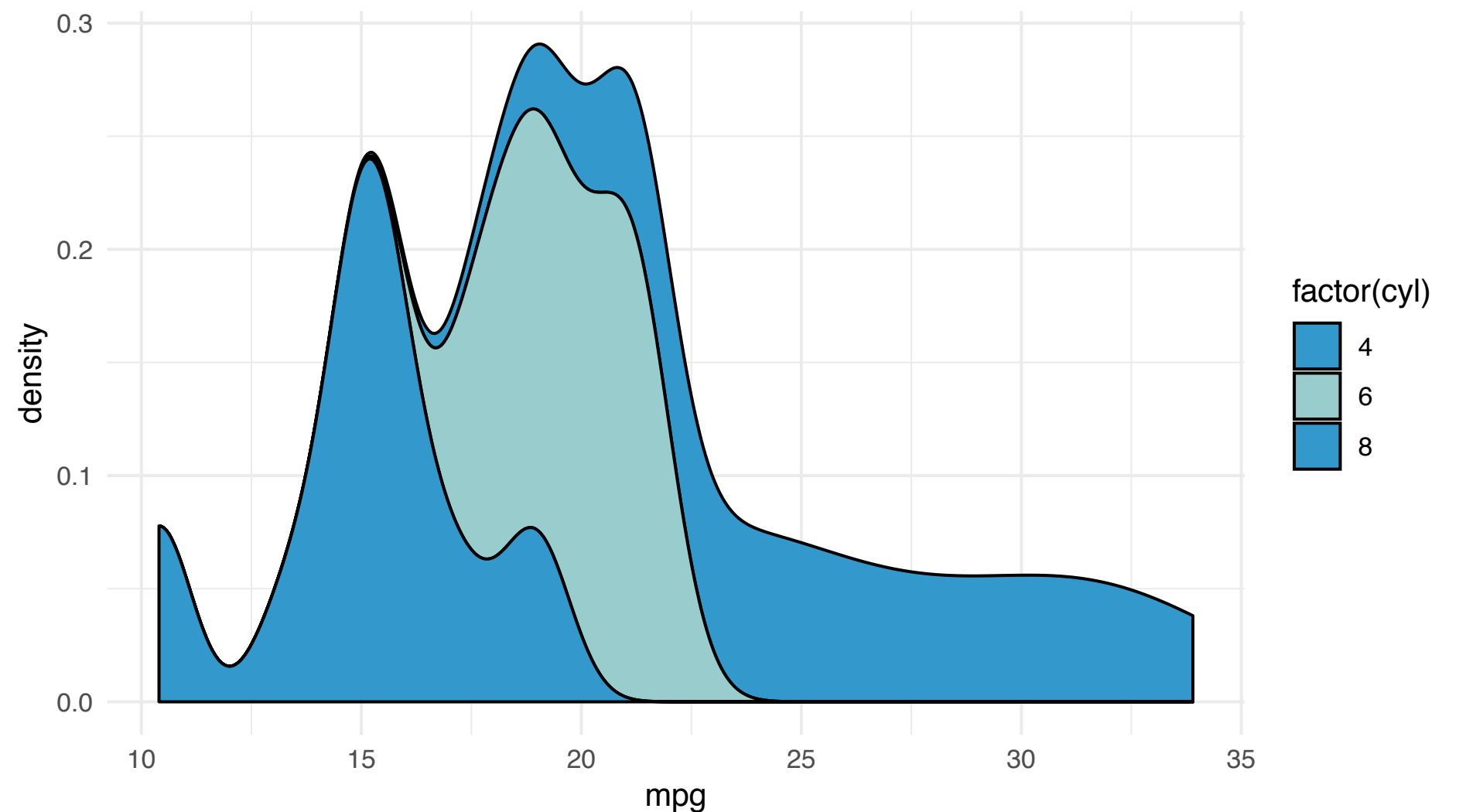
# A Probabilistic Grammar of Graphics

Xiaoying Pu  
Prelim presentation

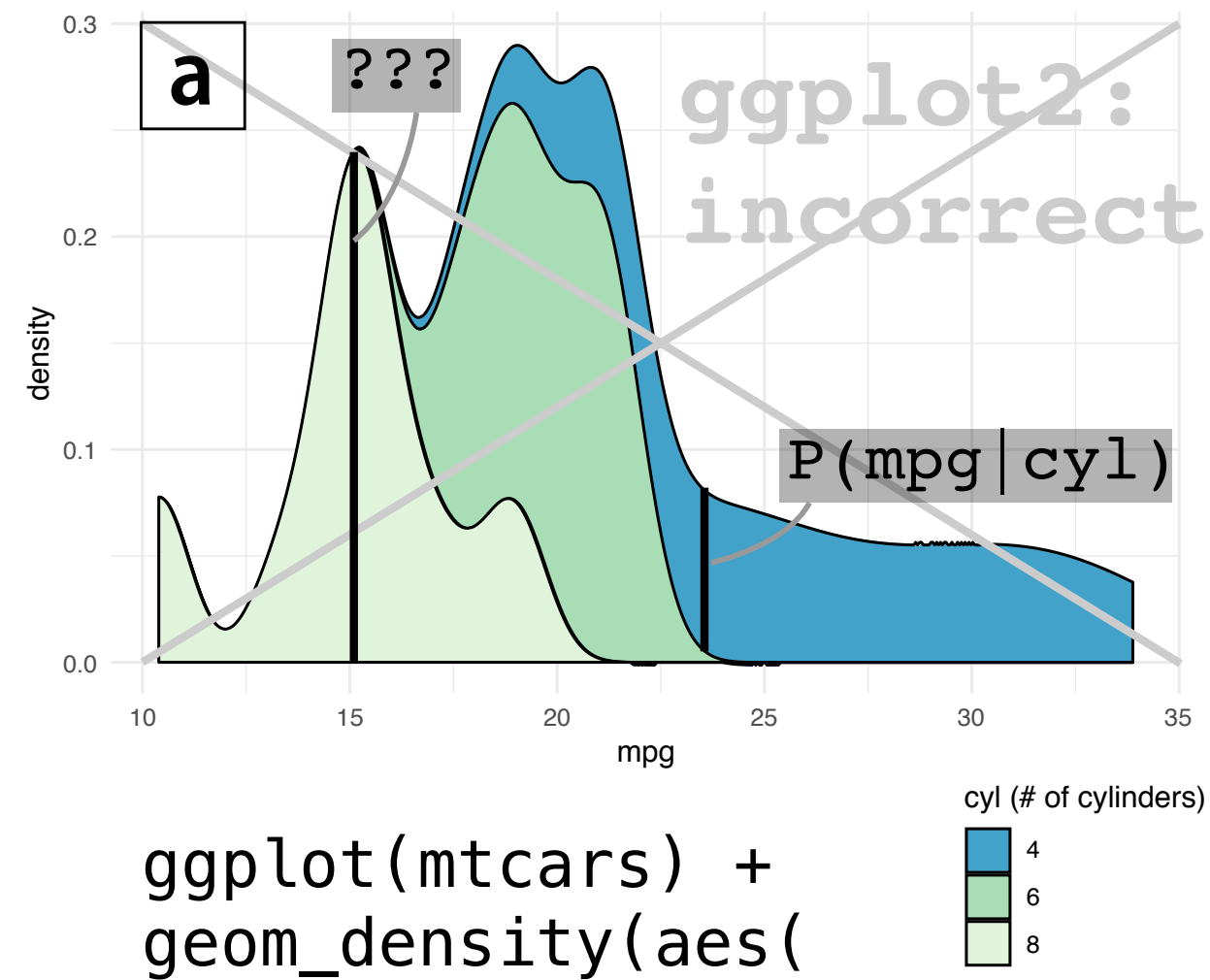
# What could possibly go wrong?

	mpg	cyl	am
Mazda RX4	21.0	6	1
Mazda RX4 Wag	21.0	6	1
Datsun 710	22.8	4	1
Hornet 4 Drive	21.4	6	0
Hornet Sportabout	18.7	8	0
Valiant	18.1	6	0

ggplot(mtcars) +  
geom\_density(aes(  
 x = mpg,  
 fill = cyl),  
 position = "stack")

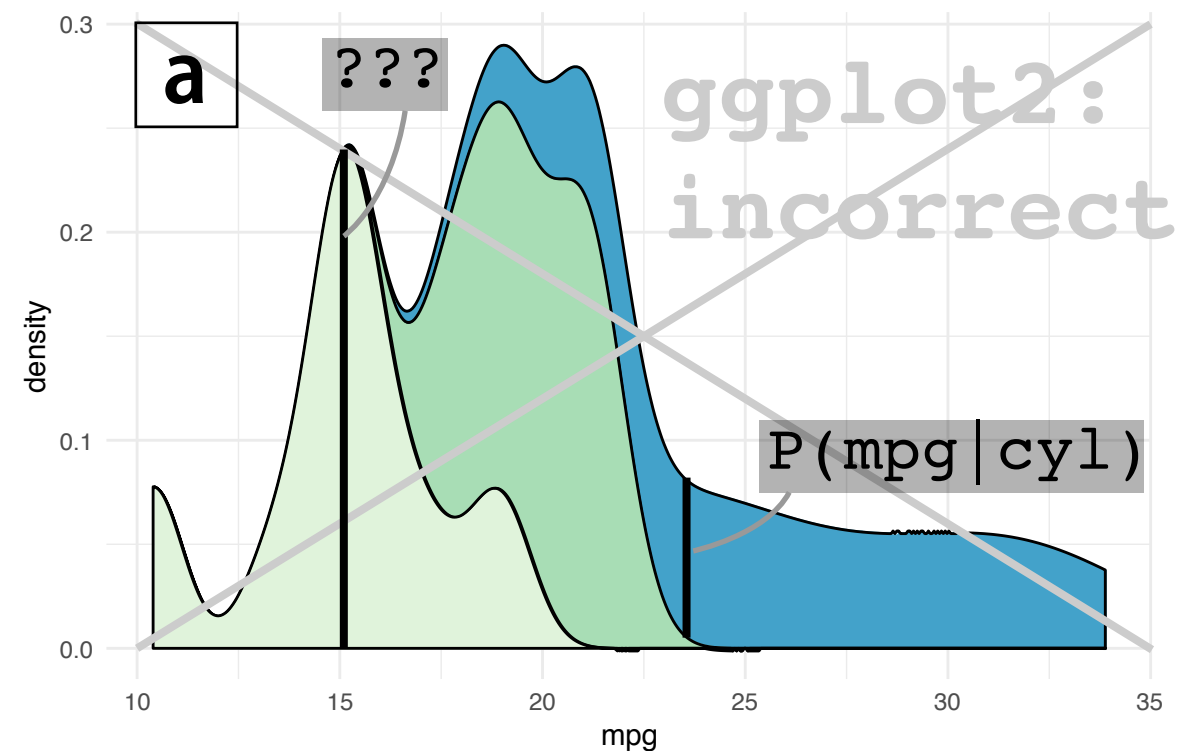


# Problem 1: visualization shows incorrect probability distrib

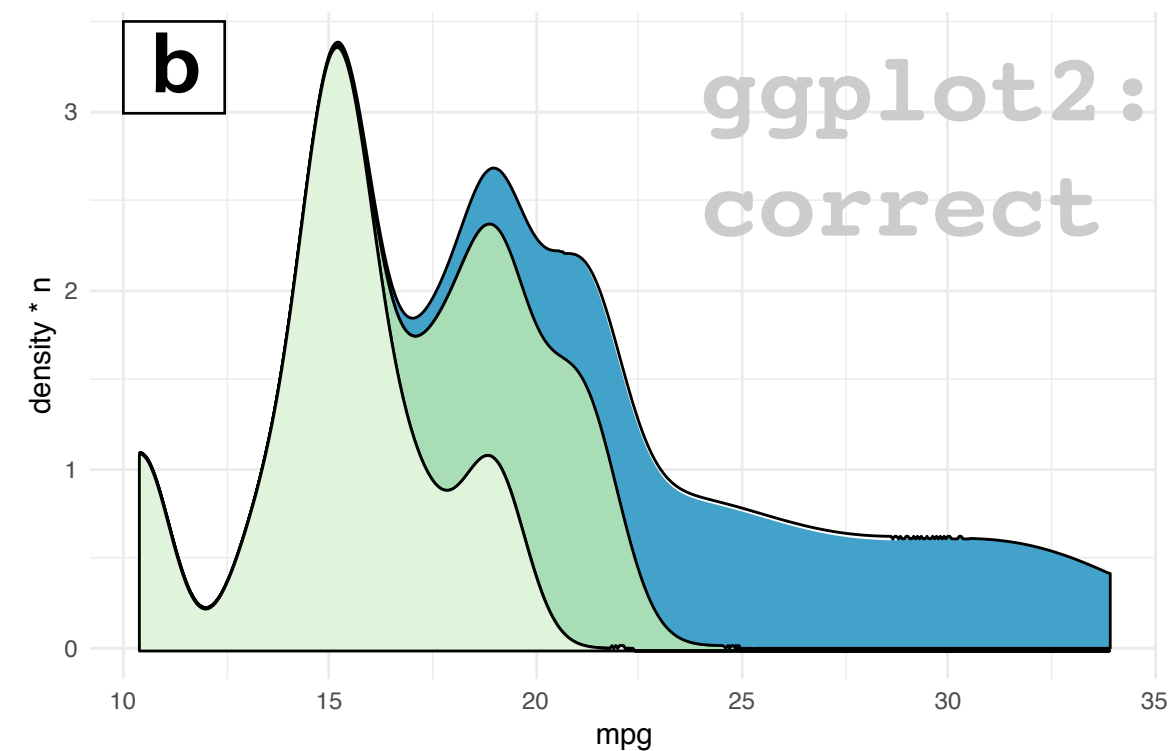
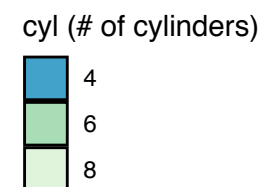


```
ggplot(mtcars) +  
  geom_density(aes(  
    x = mpg,  
    fill = cyl),  
  position = "stack")
```

# Wait we can fix this density plot

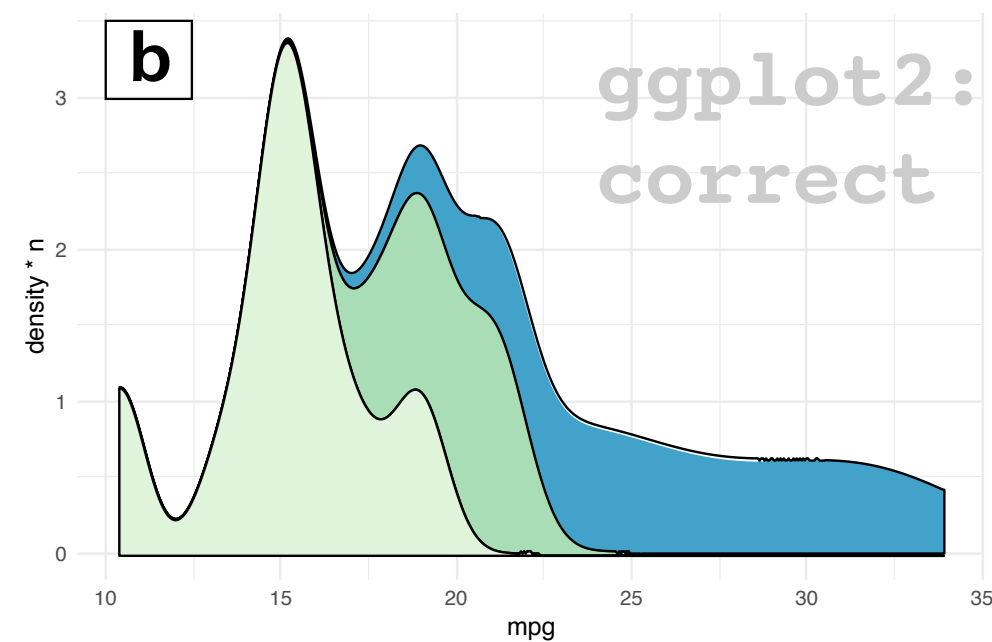


```
ggplot(mtcars) +  
  geom_density(aes(  
    x = mpg,  
    fill = cyl),  
  position = "stack")
```

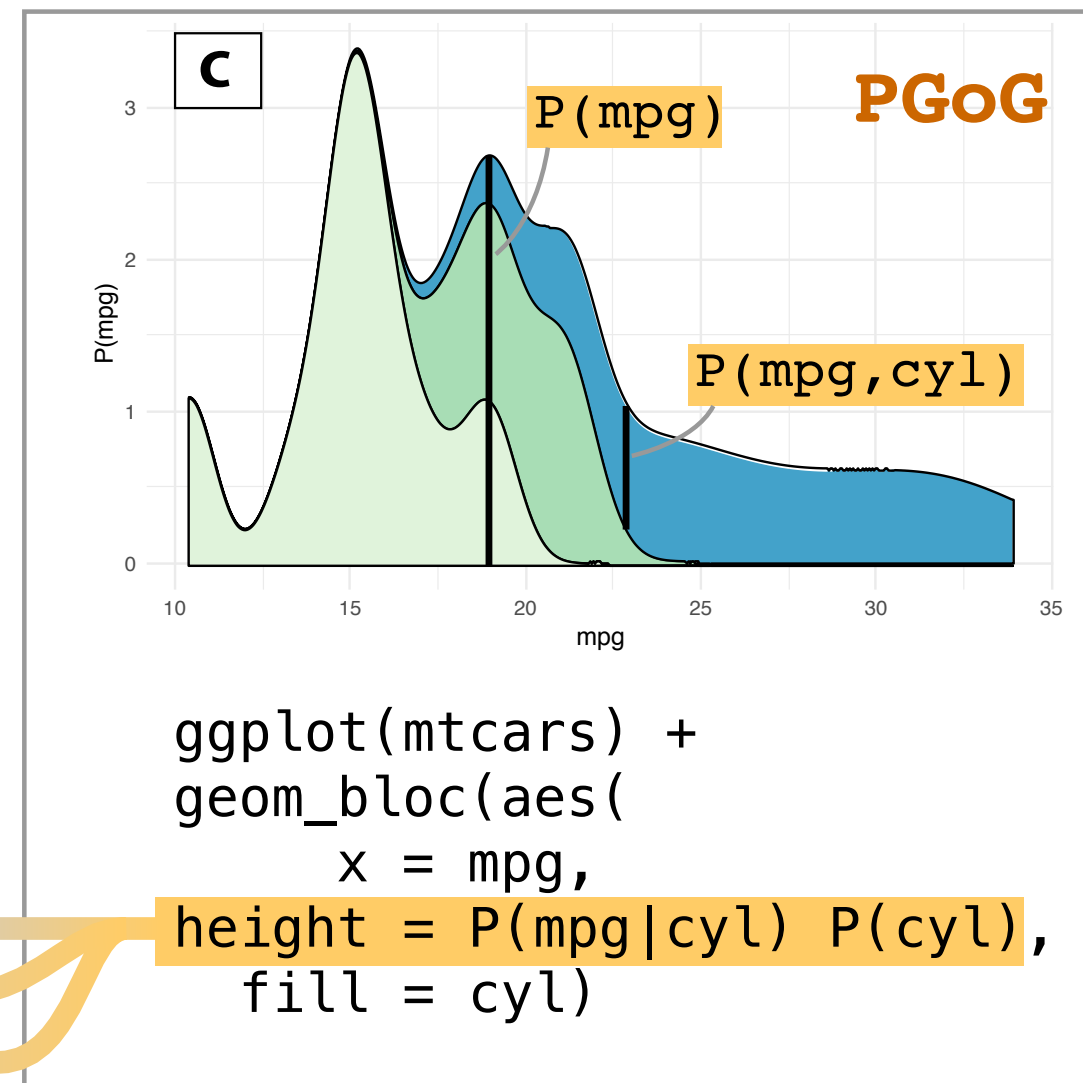


```
ggplot(mtcars)+  
  geom_density(aes(  
    x = mpg,  
    y = stat(density*n),  
    fill = cyl)) +  
  position = "stack")
```

# Problem 2: specifying probability distributions is convoluted



```
ggplot(mtcars)+  
  geom_density(aes(  
    x = mpg,  
    y = stat(density*n),  
    fill = cyl)) +  
  position = "stack")
```



```
ggplot(mtcars) +  
  geom_bloc(aes(  
    x = mpg,  
    height = P(mpg|cyl) P(cyl),  
    fill = cyl)
```

But what are `stat(density*n)` and `position`?

$$P(A|B)$$

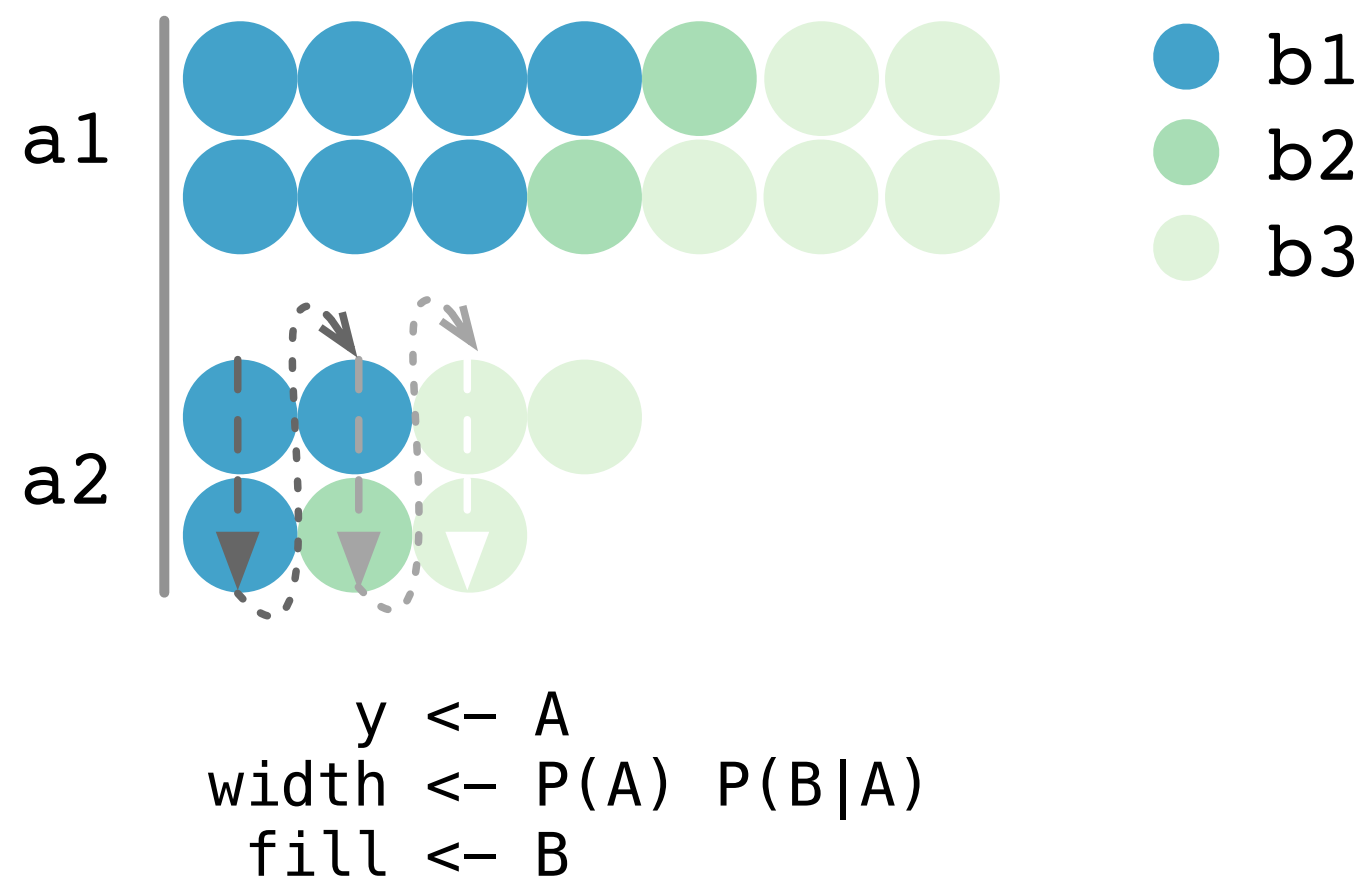
Given







1. visualization shows incorrect *probability distribution*
2. specifying *probability distribution* is convoluted

## A Probabilistic Grammar of Graphics

- A high-level visualization grammar
- Makes probability distributions first-class citizens (thus solving the two problems)
- Covers a meaningful set of probabilistic visualizations

aes	$P(\text{cyl}, \text{gear} \text{am})$	marg	cond
x.cond	$\leftarrow P(1 \text{am})$	1	am
x.height	$\leftarrow P(\text{gear} \text{am})$	gear	am
height	$\leftarrow P(\text{cyl} \text{gear}, \text{am})$	cyl	gear, am



Grammar	ggplot2	PGoG
Defaults		
Data	-----> A	P(A)
Mapping	----> x ← A	height ← P(A)
Layer		
Data		
Mapping		
Geom	----> geom_bar	geom_bloc
Stat		
Position	geom_density	geom_icon
Scale		
	geom_points	
Coord		
Facet	geom_rect	
		
	geom_...	



Existing ggplot2 packages

Changes

Syntax



geom:mosaic→bar  
+fill +y  
-divider

```
geom_mosaic
  x = cyl,
  mpg*
divider = hspine,
hspine
```

Probabilistic Grammar of Graphics (PGoG)

Changes

Syntax

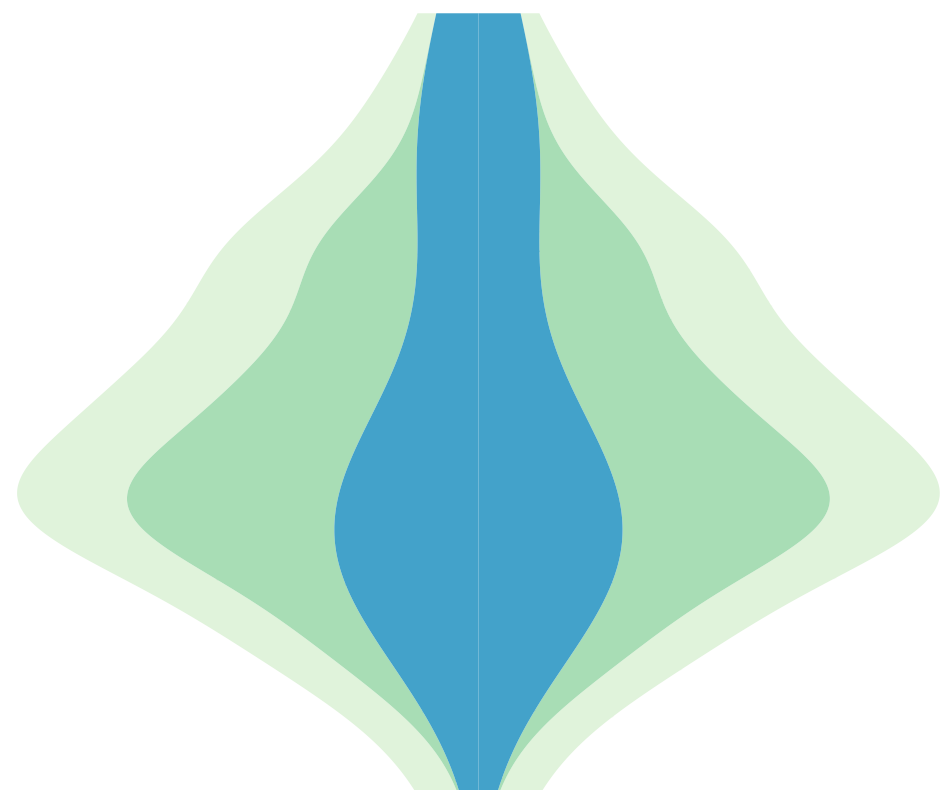


+x

```
geom_bloc
  h <- P(mpg*)
      P(cyl | mpg*)
fill<- cyl
```

mpg\*: discretized miles per gallon





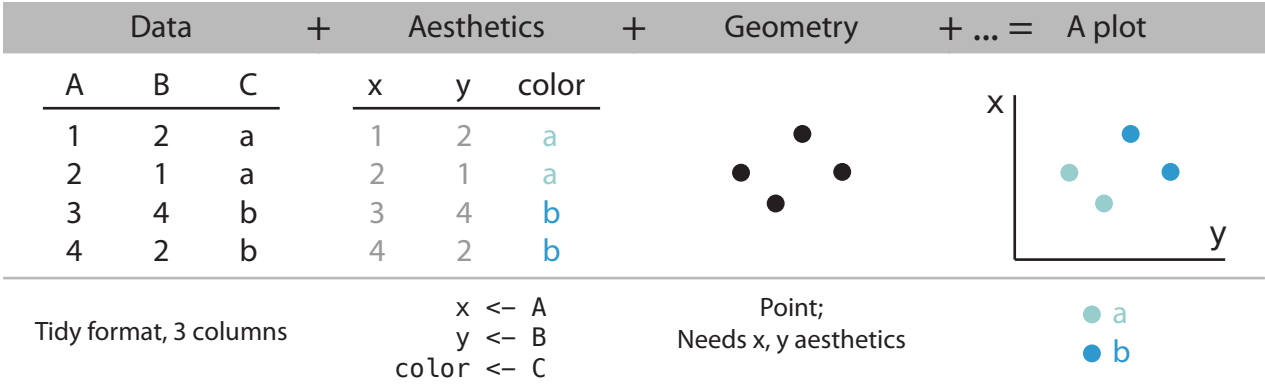
## Onion plot

`geom_bloc:`

`y`  $\leftarrow$  mpg

`width`  $\leftarrow$   $P(\text{mpg})$   $P(\text{cyl})$

`direction`  $\leftarrow$  both



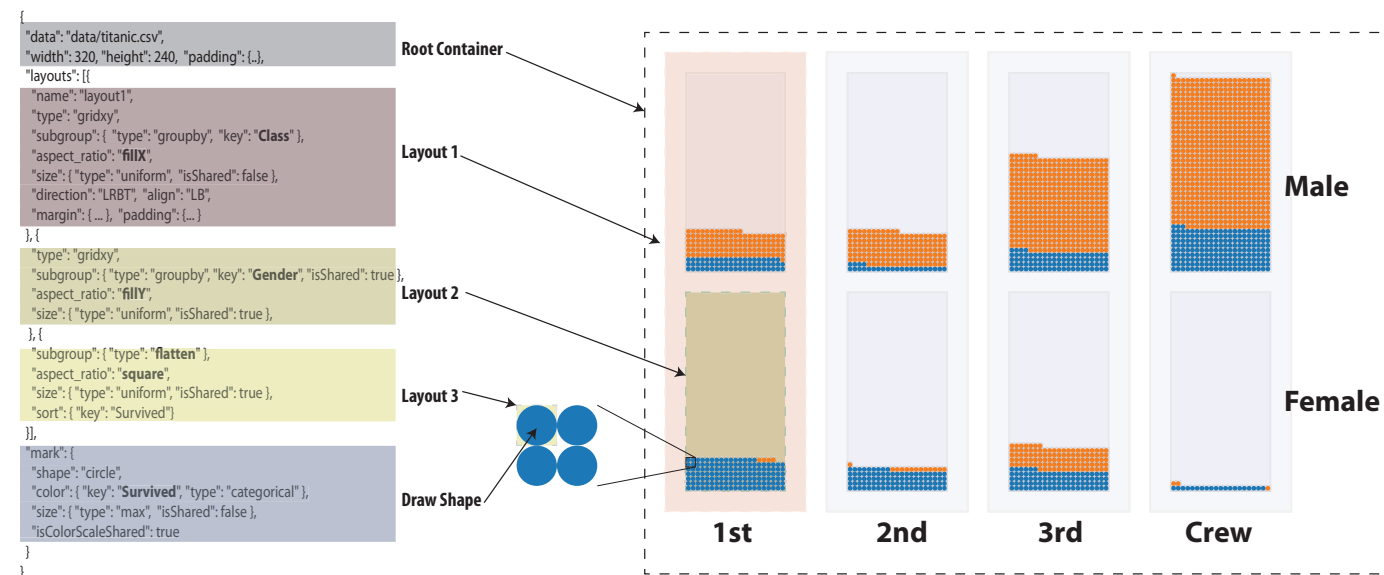
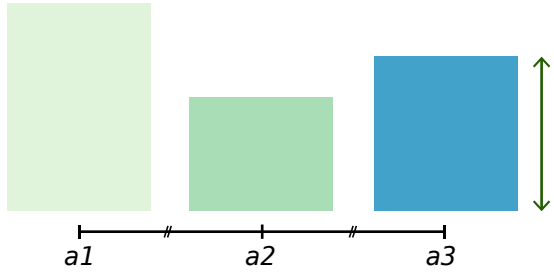
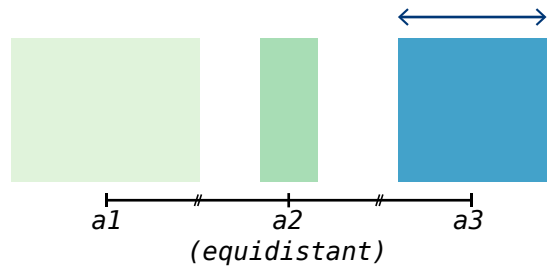
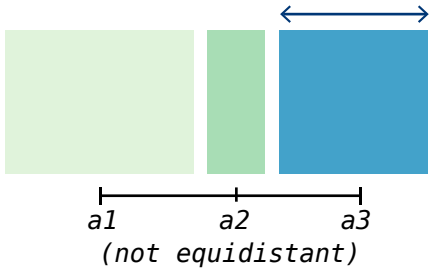


Fig. 6. Example grammar to generate a unit column chart for survivors of the Titanic by passenger class.

		PGoG	Product Plot
Normal bar	 <p><math>a1</math> <math>a2</math> <math>a3</math></p>	$x \leftarrow A$ $height \leftarrow P(A)$	$\sim A$ $hbar$
"Lying down"	 <p><math>a1</math> <math>a2</math> <math>a3</math> (equidistant)</p>	$x \leftarrow A$ $width \leftarrow P(A)$	Does not exist
Spine plot	 <p><math>a1</math> <math>a2</math> <math>a3</math> (not equidistant)</p>	$width \leftarrow P(A)$	$\sim A$ $hspine$