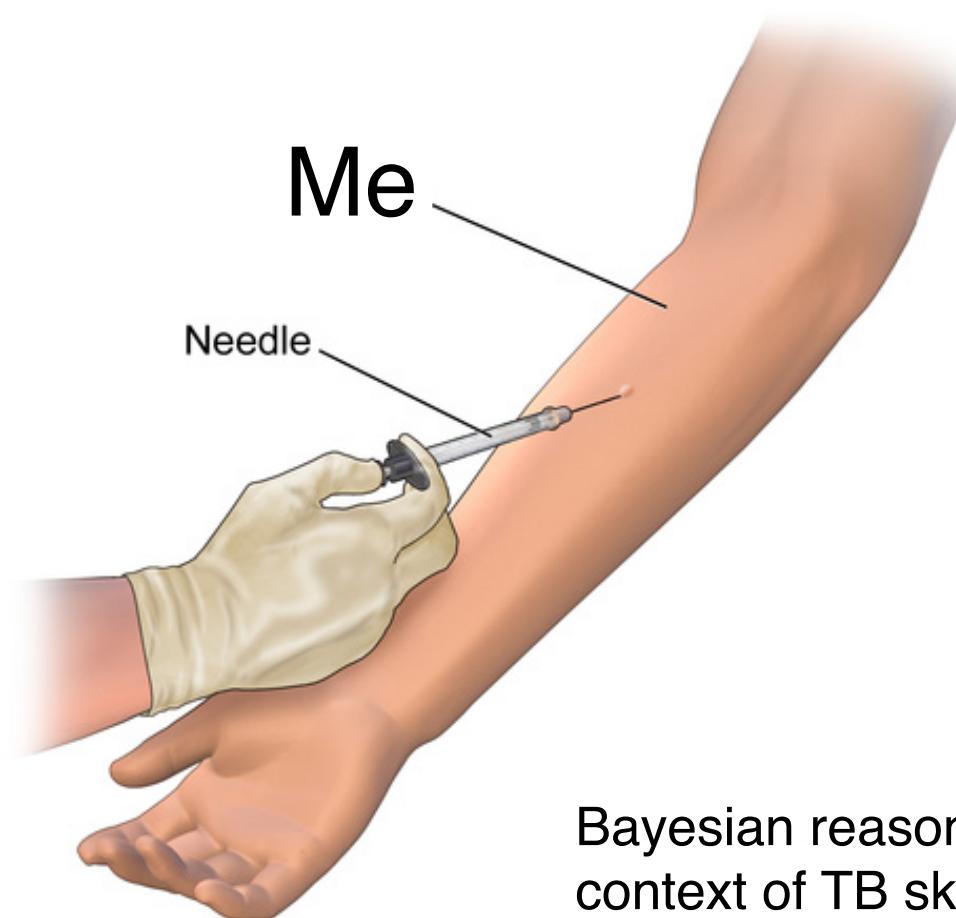


# A Probabilistic Grammar of Graphics

Xiaoying Pu  
Prelim presentation

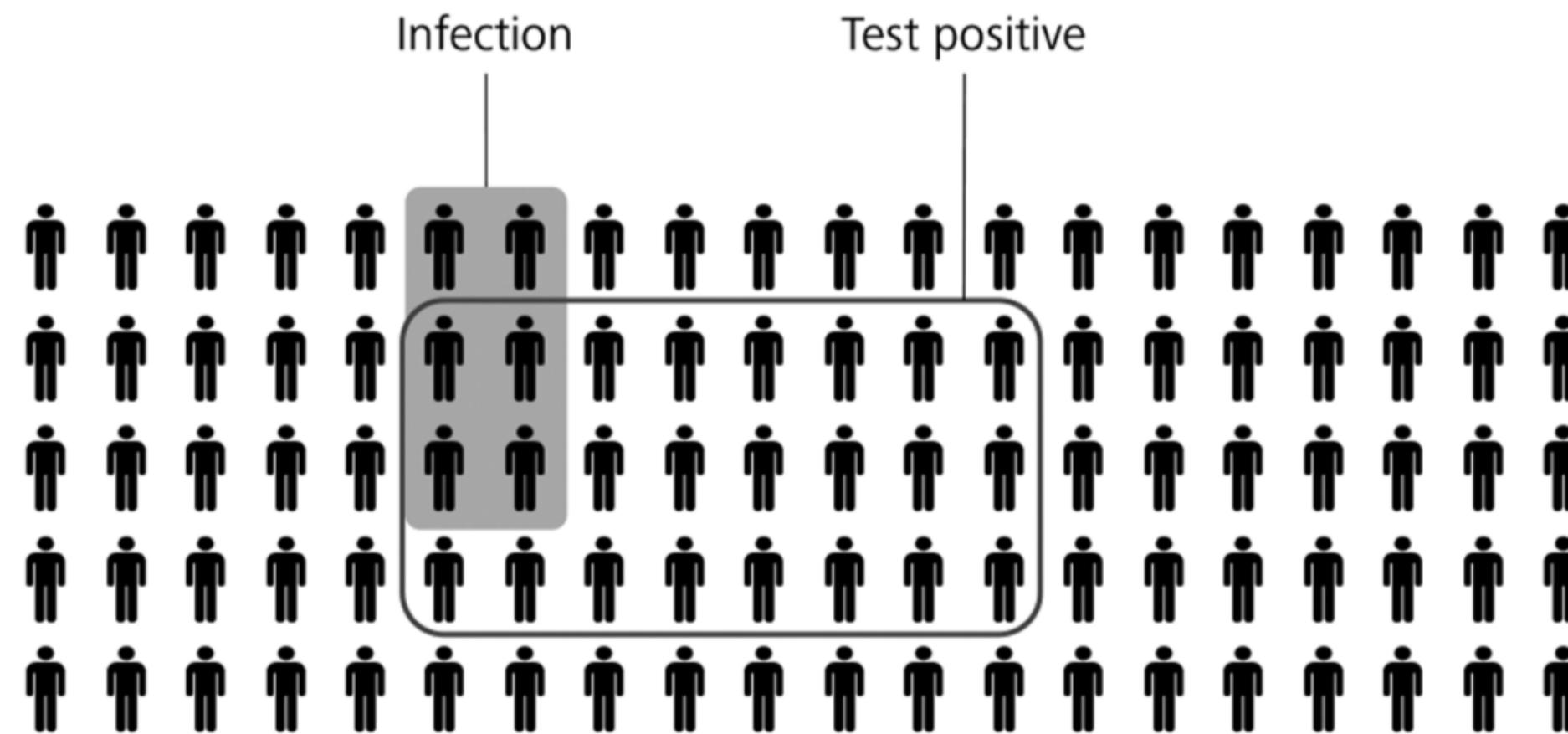
# So I got tested for tuberculosis and it came back positive



Bayesian reasoning in  
context of TB skin test

- $P(\text{test positive})?$
- $P(\text{infected})?$
- $P(\text{test positive AND infected})?$
- $P(\text{infected GIVEN test positive})?$
- arrrg,  $P(B|A) = P(B,A) / P(A)???$ !!!

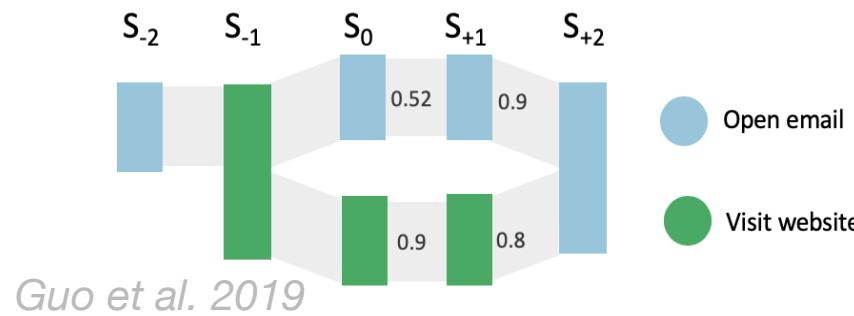
# This visualization can help a stress-ridden patient



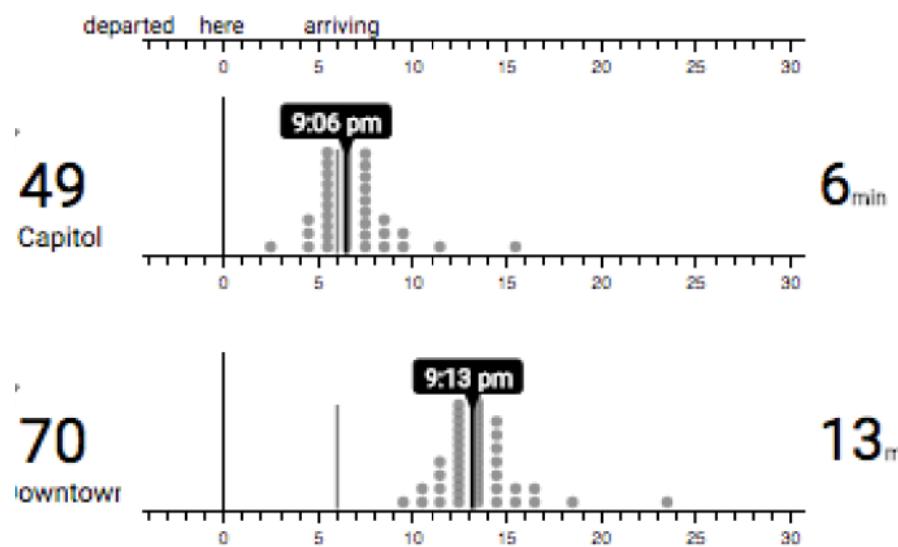
(Binder, Krauss, and Bruckmaier 2015)

- $P(\text{test positive}) = 24\%$
- $P(\text{infected}) = 6\%$
- $P(\text{test positive AND infected}) = 4\%$
- $P(\text{infected} \mid \text{test positive}) = 17\%$

## Machine learning

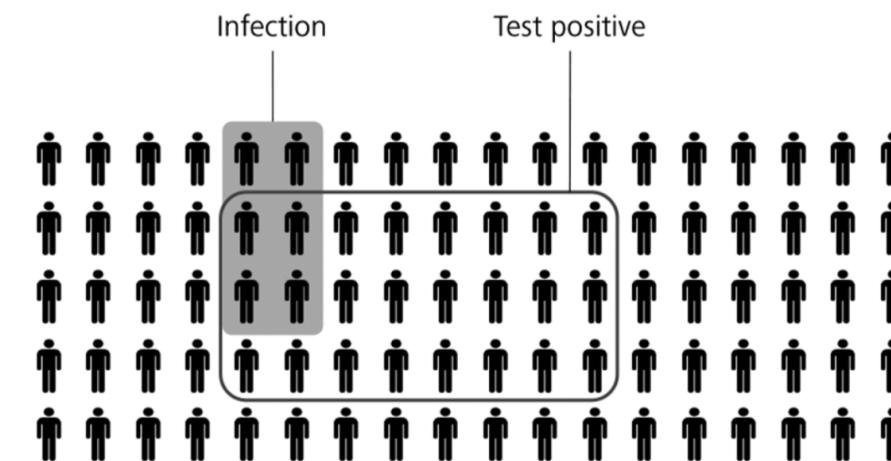


## Bus arrival time



(Fernandes et al. 2018)

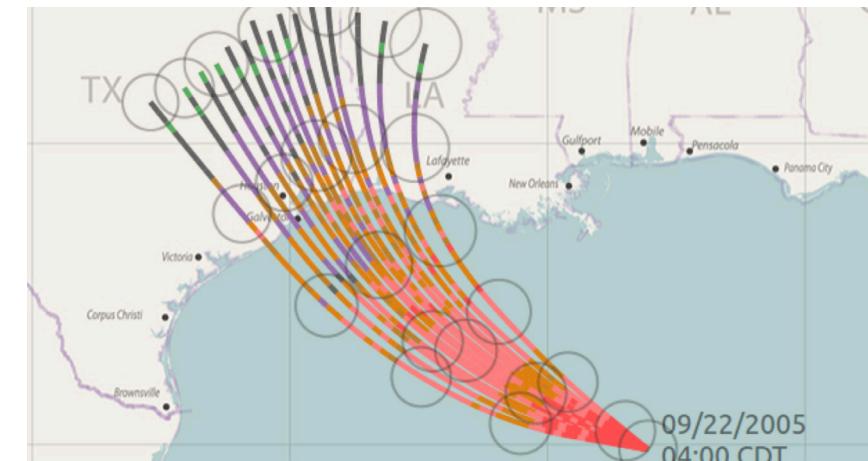
## Medical risk communication



(Binder, Krauss, and Bruckmaier 2015)

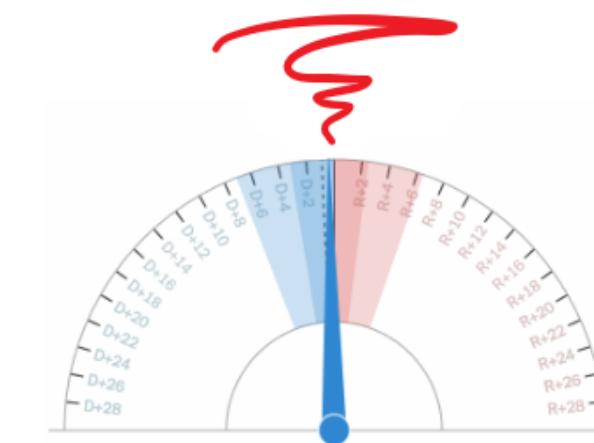
Probabilistic visualizations:  
same substrate, many benefits  
... but difficult to spec right

## Hurricane forecast



(Liu et al., 2018)

## Election forecast



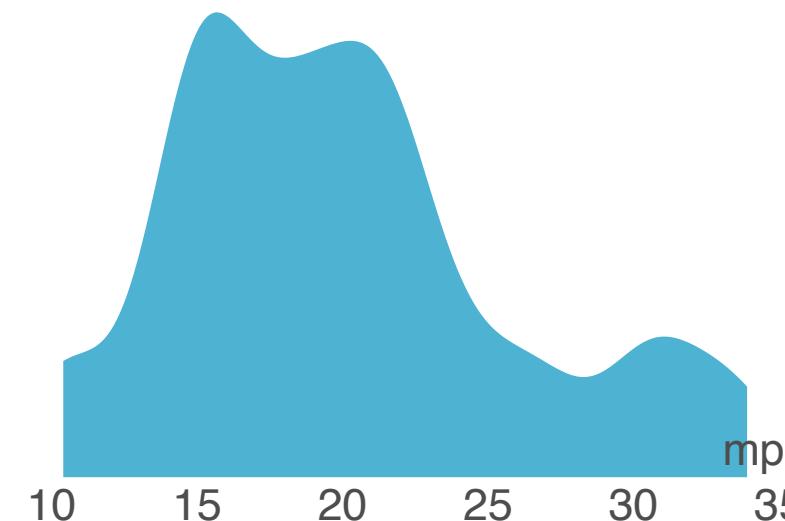
(New York Times, 2016)

# Motivating example: what could possibly go wrong?

	mpg	cyl
Mazda RX4	21.0	6
Mazda RX4 Wag	21.0	6
Datsun 710	22.8	4
Hornet 4 Drive	21.4	6

A user's mental process during data exploration

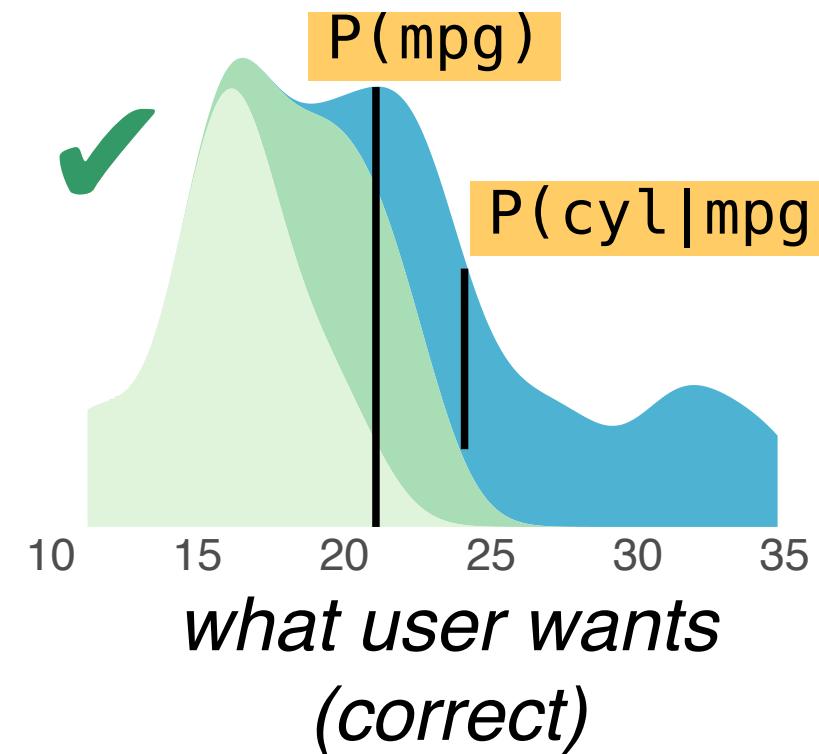
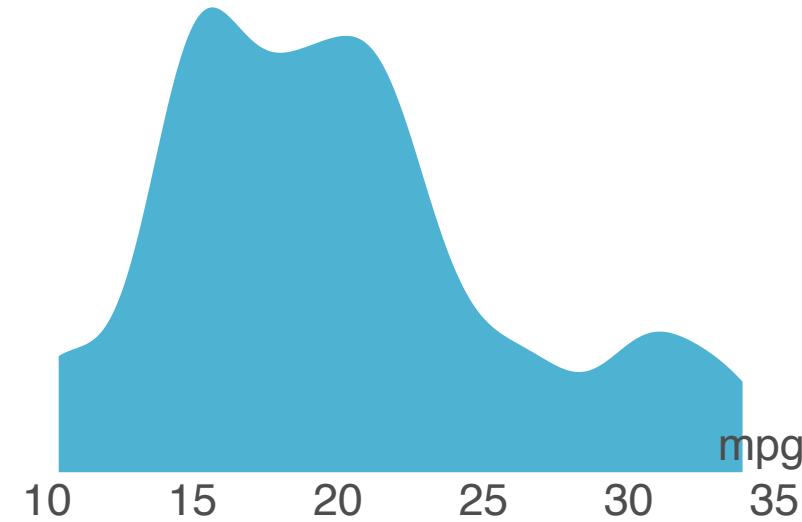
*I care about emissions.  
What's the mileage data like?*



# What could possibly go wrong?

A user's mental process during data exploration

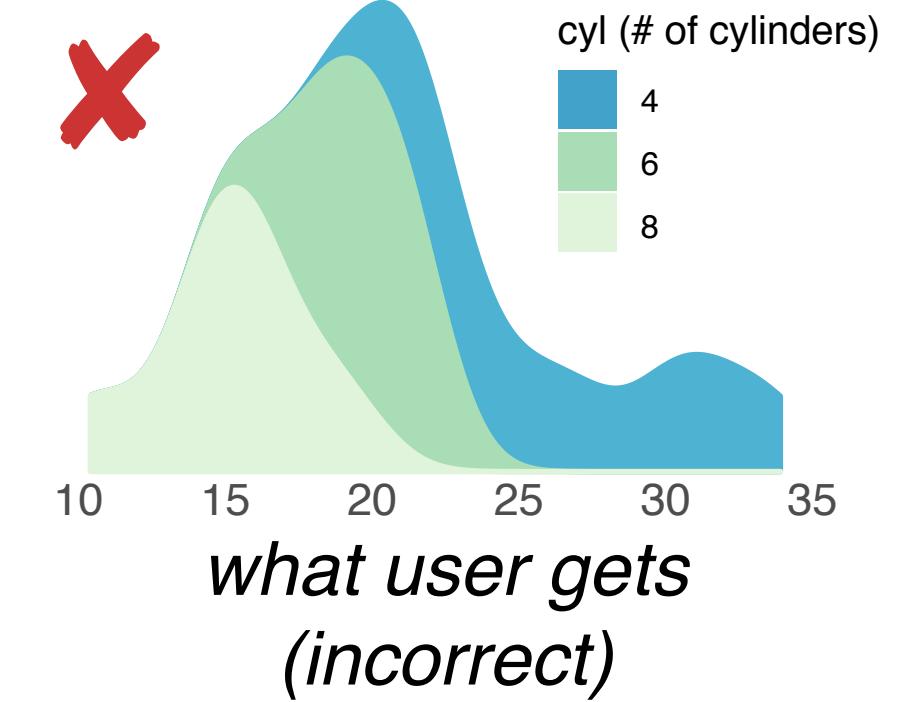
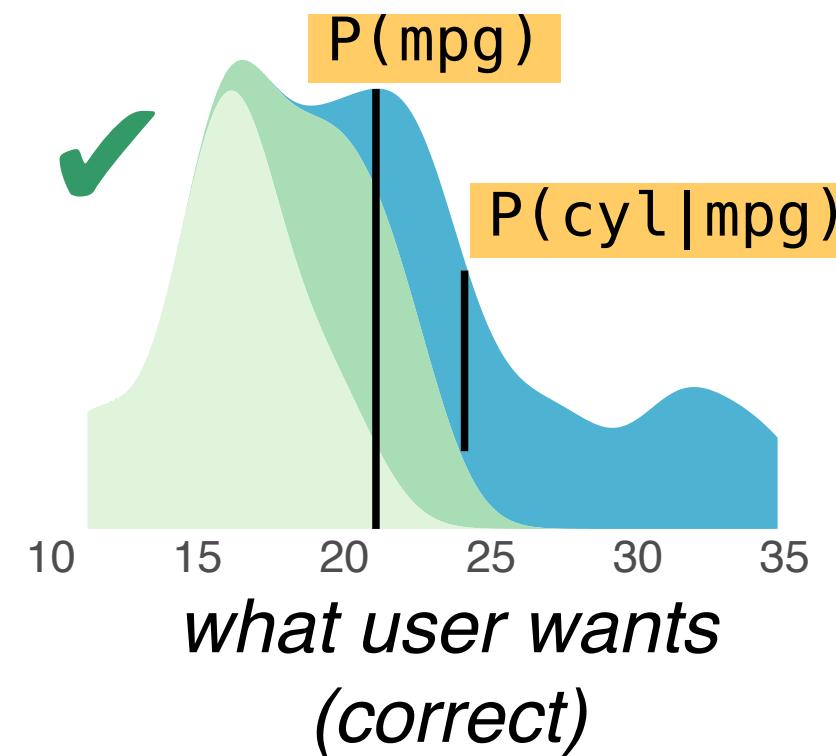
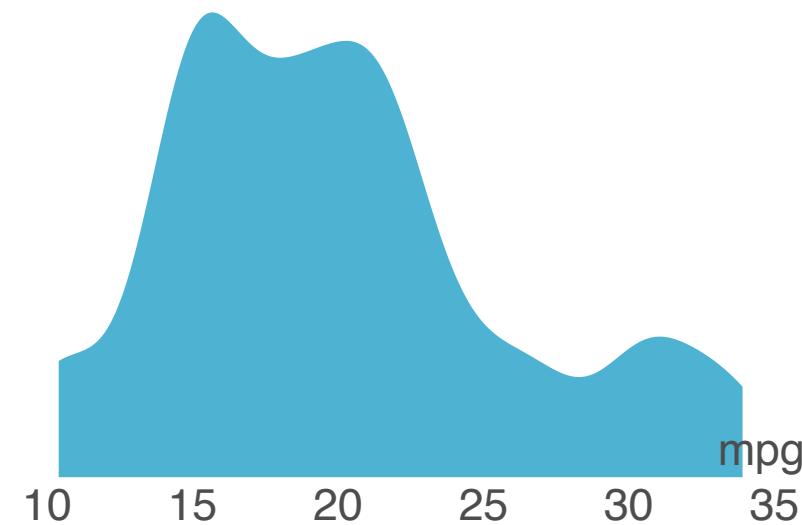
*What's the mileage data like? I want to see both mileage and cylinder counts... maybe there's a pattern*



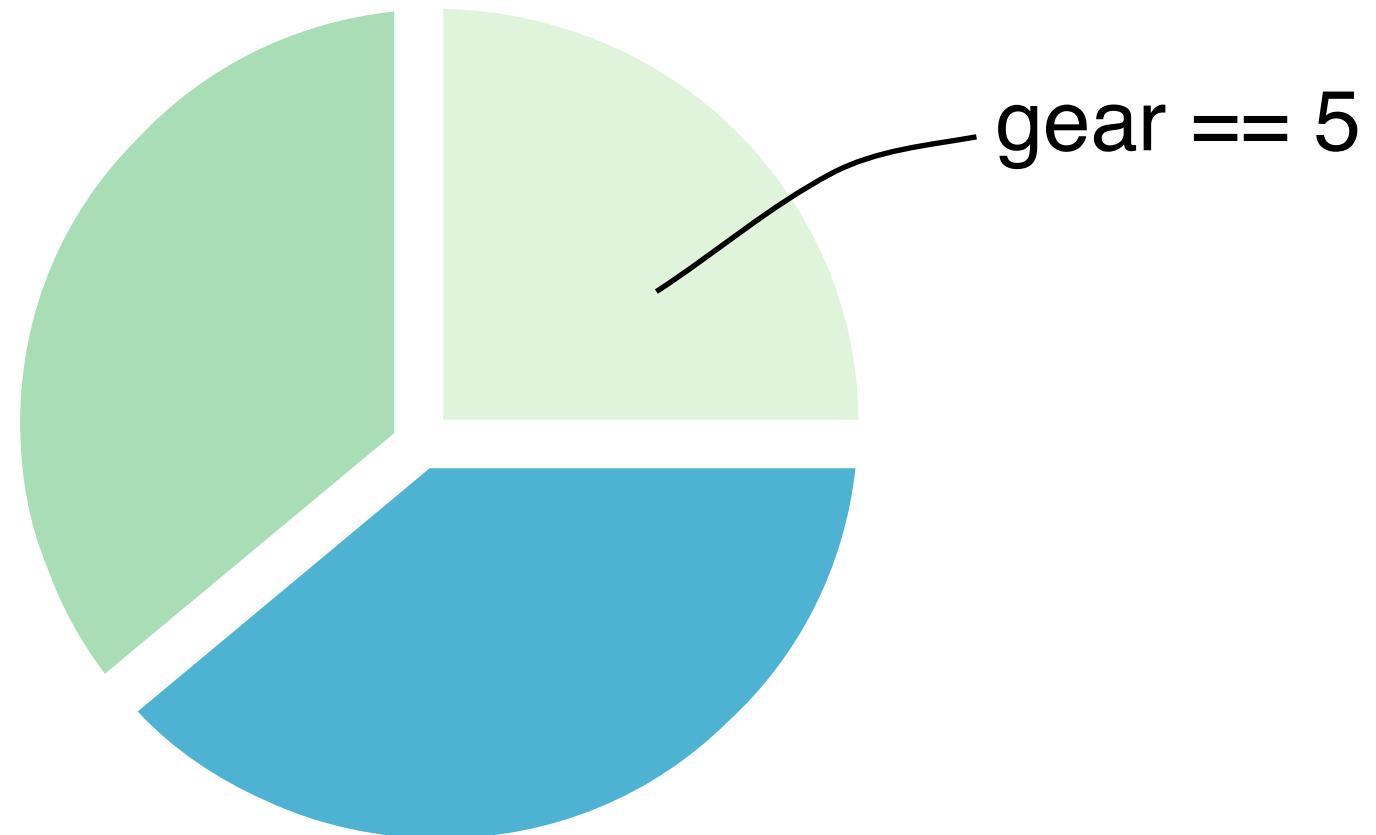
# What could possibly go wrong?

A user's mental process during data exploration

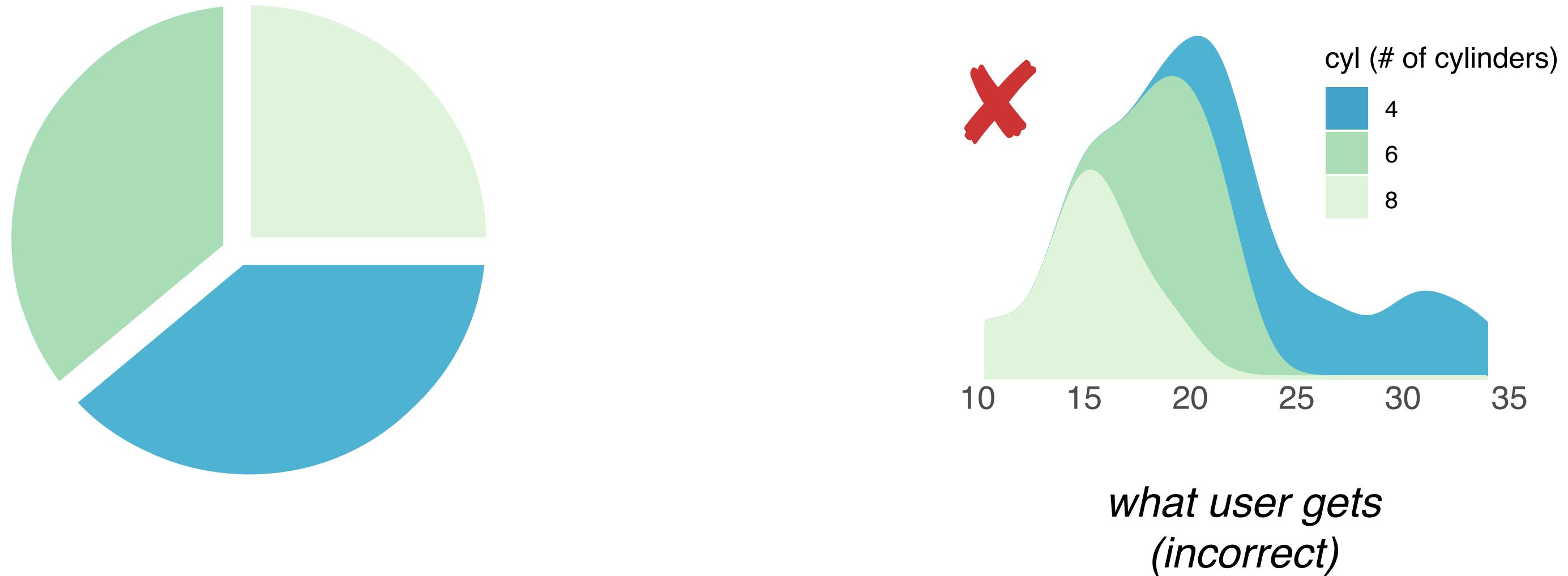
*What's the mileage data like? I want to see both mileage and cylinder counts... maybe there's a pattern*



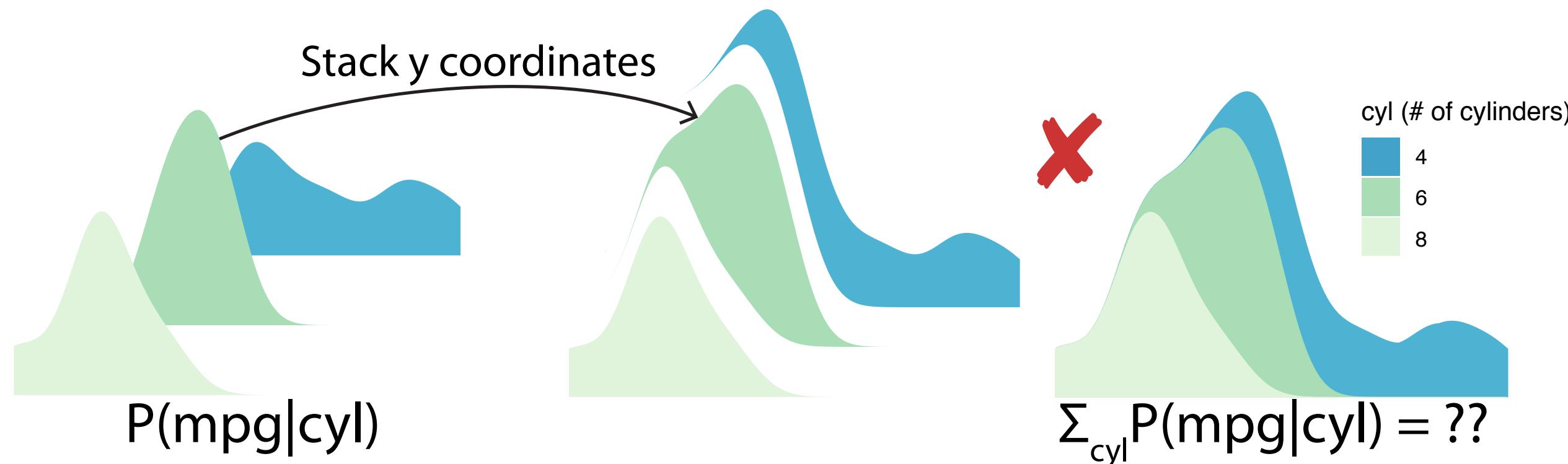
# Problem 1: vis shows incorrect probability distribution



# Problem 1: vis shows incorrect probability distribution



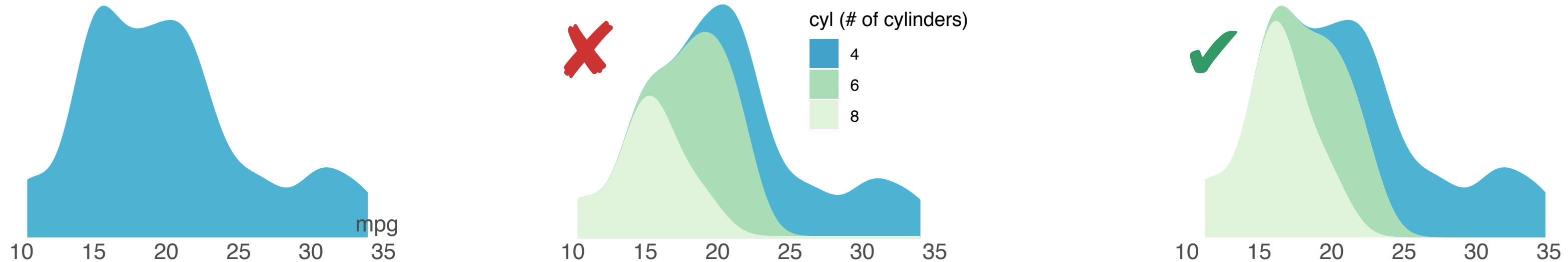
# Problem 1: vis shows incorrect probability distribution



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

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

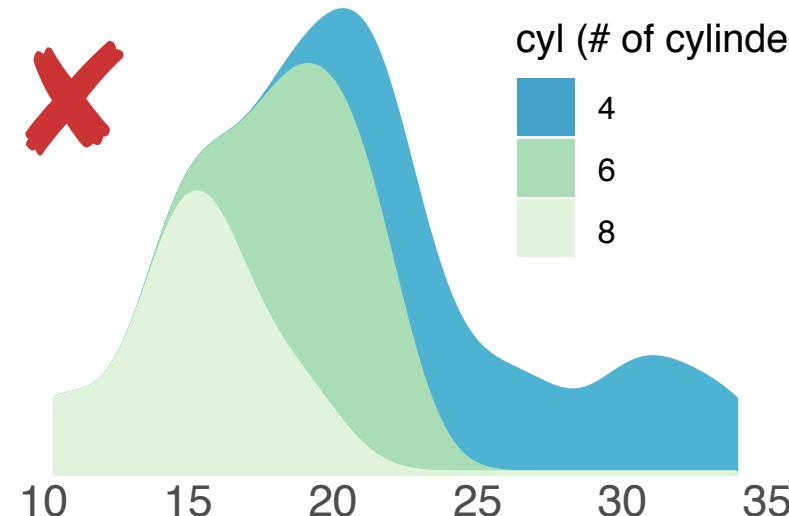
# Problem 1: vis shows incorrect probability distribution



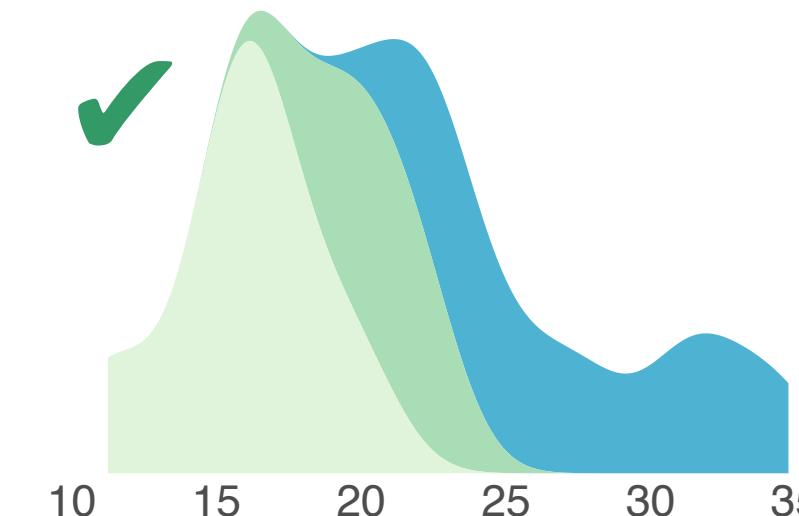
Two incorrect inferences

1. Wrong distribution of # of cylinders  $P(\text{cyl})$
2. Wrong overall distribution of mileage  $P(\text{mpg})$

# Wait we can fix this density plot

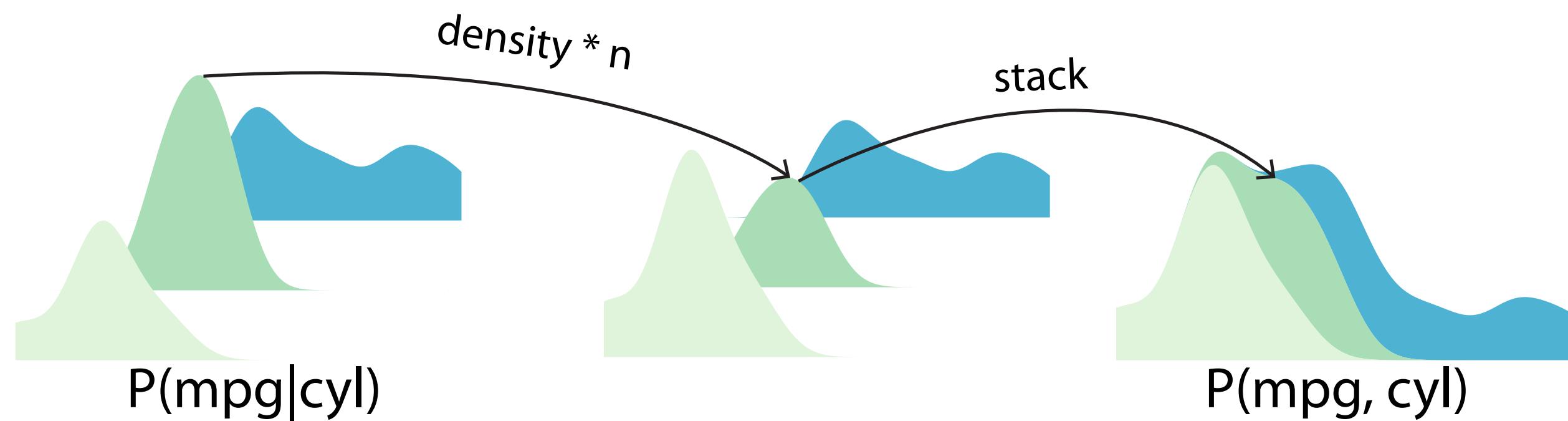


```
ggplot(mtcars) +  
  geom_density(  
    aes(x = mpg,  
        y = stat(density),  
        fill = cyl),  
    position = "stack")
```



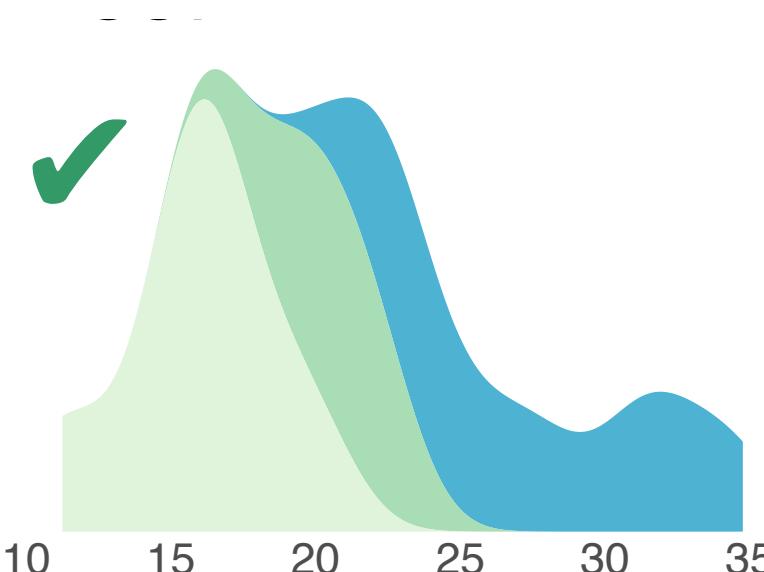
```
ggplot(mtcars) +  
  geom_density(  
    aes(x = mpg,  
        y = stat(density*n),  
        fill = cyl)),  
    position = "stack")
```

# Problem 2: specifying probability distributions is convoluted

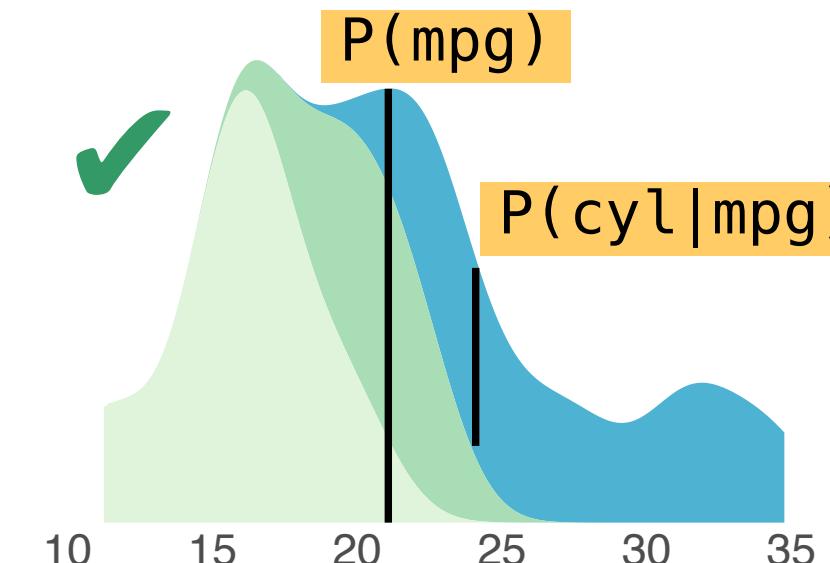


But how I am supposed to know `stat(density*n)` and `position`?

# Problem 2 can be solved with ...



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



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

Details later

# PGoG

Given

1. The need to visualize *probability distributions*
2. Specifying probability distributions is convoluted and error-prone

## A Probabilistic Grammar of Graphics

- A visualization grammar that makes probability distributions first-class citizens
- Unifies a meaningful set of probabilistic visualizations
- Cognitively ergonomic and guaranteed to be correct

# Outline

PGoG in context of

- visualization specification grammar/languages
- formats for communicating probability distributions

Design Requirements  
for PGoG

PGoG abstract grammar

PGoG implementation

Evaluation in terms of

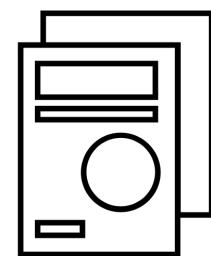
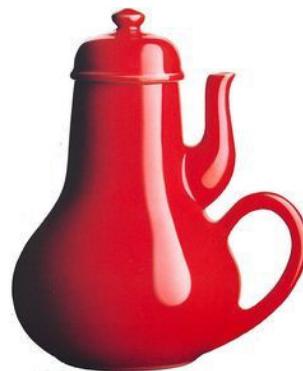
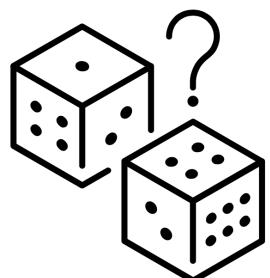
- Expressiveness
- Generativeness
- Cognitive ergonomics

Future: quantitative uncertainty communication!

# Related work: motivations and theories behind PGoG

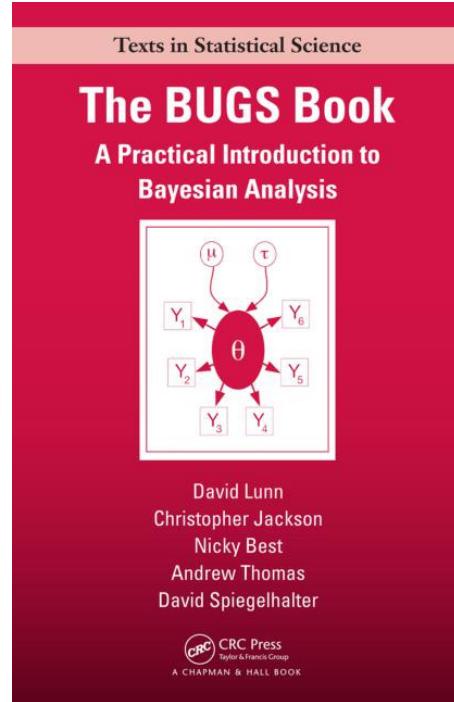
Probabilistic visualizations are important but difficult/error-prone to specify and explore...

What should a grammar for probabilistic visualizations look like?



- Probability first: benefits of having a “probabilistic” grammar
- Usability and correctness: current visualization specs are inadequate for probabilistic visualizations
- Theoretical grounding: how uncertainty communication can inform visualization design

# Related: probabilistic programming



BUGS



Stan

Widely cited and applied  
to many domains

$$\theta \sim \text{beta}(1, 1)$$
$$y \sim \text{bernoulli}(\theta)$$

```
model {  
    theta ~ beta(1, 1); //prior  
    y ~ bernoulli(theta); //likelihood  
}
```

Sticks with what users know (probabilities)  
Avoids implementation details (stat(density\*n))

# Related: how to specify a visualization (Grammar of Graphics)

Data +

example\_df

---

A	B	C
1	2	a
2	1	a
3	4	b
4	2	b

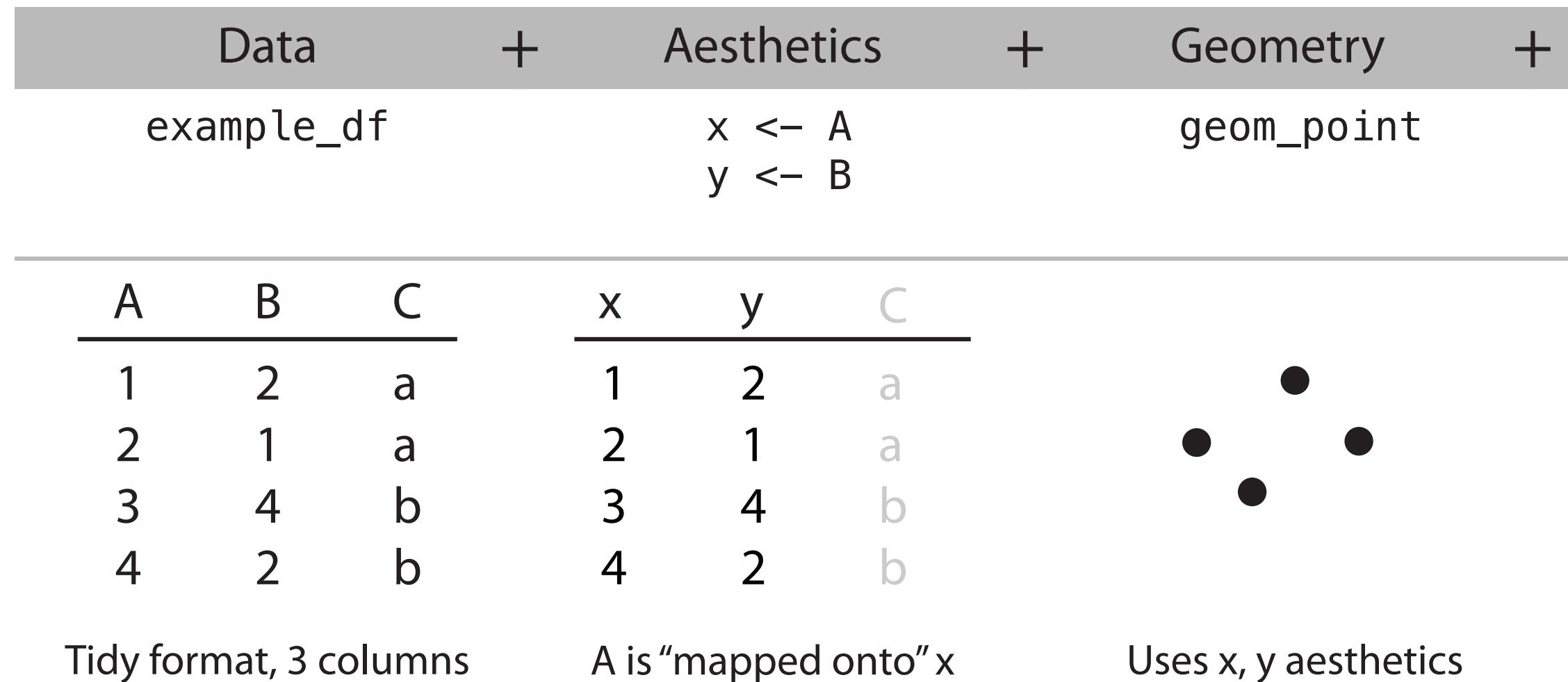
Tidy format, 3 columns

# Related: how to specify a visualization (Grammar of Graphics)

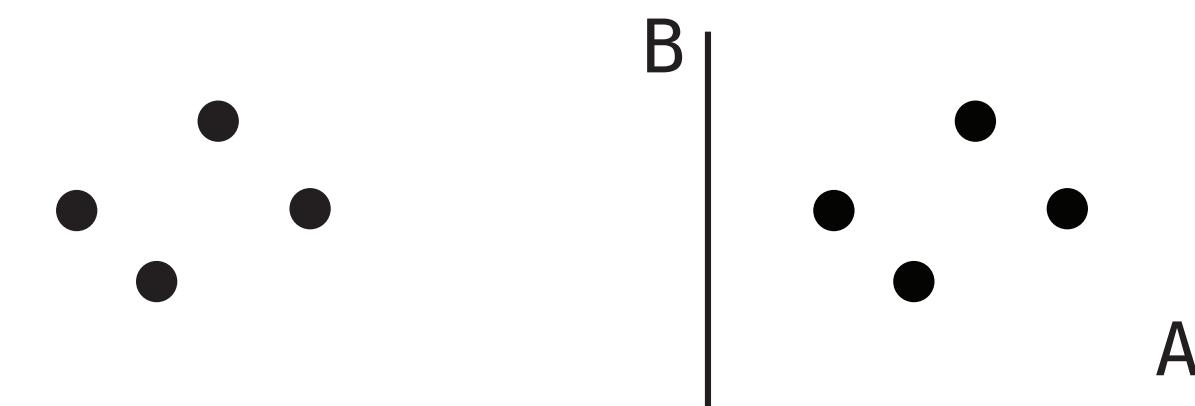
Data	+	Aesthetics	+
example_df		x <- A y <- B	
<hr/>			
A	B	C	x
1	2	a	1
2	1	a	2
3	4	b	3
4	2	b	4
			y
			2
			1
			4
			2
		C	
		a	
		a	
		b	
		b	

Tidy format, 3 columns      A is “mapped onto” x

# Related: how to specify a visualization (Grammar of Graphics)



# Related: how to specify a visualization (Grammar of Graphics)

Data	+	Aesthetics	+	Geometry	+ ... = A plot																														
example_df		x <- A y <- B		geom_point																															
<table><thead><tr><th>A</th><th>B</th><th>C</th></tr></thead><tbody><tr><td>1</td><td>2</td><td>a</td></tr><tr><td>2</td><td>1</td><td>a</td></tr><tr><td>3</td><td>4</td><td>b</td></tr><tr><td>4</td><td>2</td><td>b</td></tr></tbody></table>	A	B	C	1	2	a	2	1	a	3	4	b	4	2	b		<table><thead><tr><th>x</th><th>y</th><th>C</th></tr></thead><tbody><tr><td>1</td><td>2</td><td>a</td></tr><tr><td>2</td><td>1</td><td>a</td></tr><tr><td>3</td><td>4</td><td>b</td></tr><tr><td>4</td><td>2</td><td>b</td></tr></tbody></table>	x	y	C	1	2	a	2	1	a	3	4	b	4	2	b			
A	B	C																																	
1	2	a																																	
2	1	a																																	
3	4	b																																	
4	2	b																																	
x	y	C																																	
1	2	a																																	
2	1	a																																	
3	4	b																																	
4	2	b																																	

Tidy format, 3 columns      A is “mapped onto” x      Uses x, y aesthetics      A scatter plot

# Related: how to specify a visualization (Grammar of Graphics)

Data	+	Aesthetics	+	Geometry	+ ... =	A plot																														
example_df		x <- A y <- B color <- C		geom_point																																
<table><thead><tr><th>A</th><th>B</th><th>C</th></tr></thead><tbody><tr><td>1</td><td>2</td><td>a</td></tr><tr><td>2</td><td>1</td><td>a</td></tr><tr><td>3</td><td>4</td><td>b</td></tr><tr><td>4</td><td>2</td><td>b</td></tr></tbody></table>	A	B	C	1	2	a	2	1	a	3	4	b	4	2	b		<table><thead><tr><th>x</th><th>y</th><th>color</th></tr></thead><tbody><tr><td>1</td><td>2</td><td>a</td></tr><tr><td>2</td><td>1</td><td>a</td></tr><tr><td>3</td><td>4</td><td>b</td></tr><tr><td>4</td><td>2</td><td>b</td></tr></tbody></table>	x	y	color	1	2	a	2	1	a	3	4	b	4	2	b				
A	B	C																																		
1	2	a																																		
2	1	a																																		
3	4	b																																		
4	2	b																																		
x	y	color																																		
1	2	a																																		
2	1	a																																		
3	4	b																																		
4	2	b																																		
Tidy format, 3 columns		A is "mapped onto" x		Uses x, y aesthetics		A scatter plot																														

# Related: how to specify a visualization (layout-based is viscous)

(Blackwell et al. 2001)

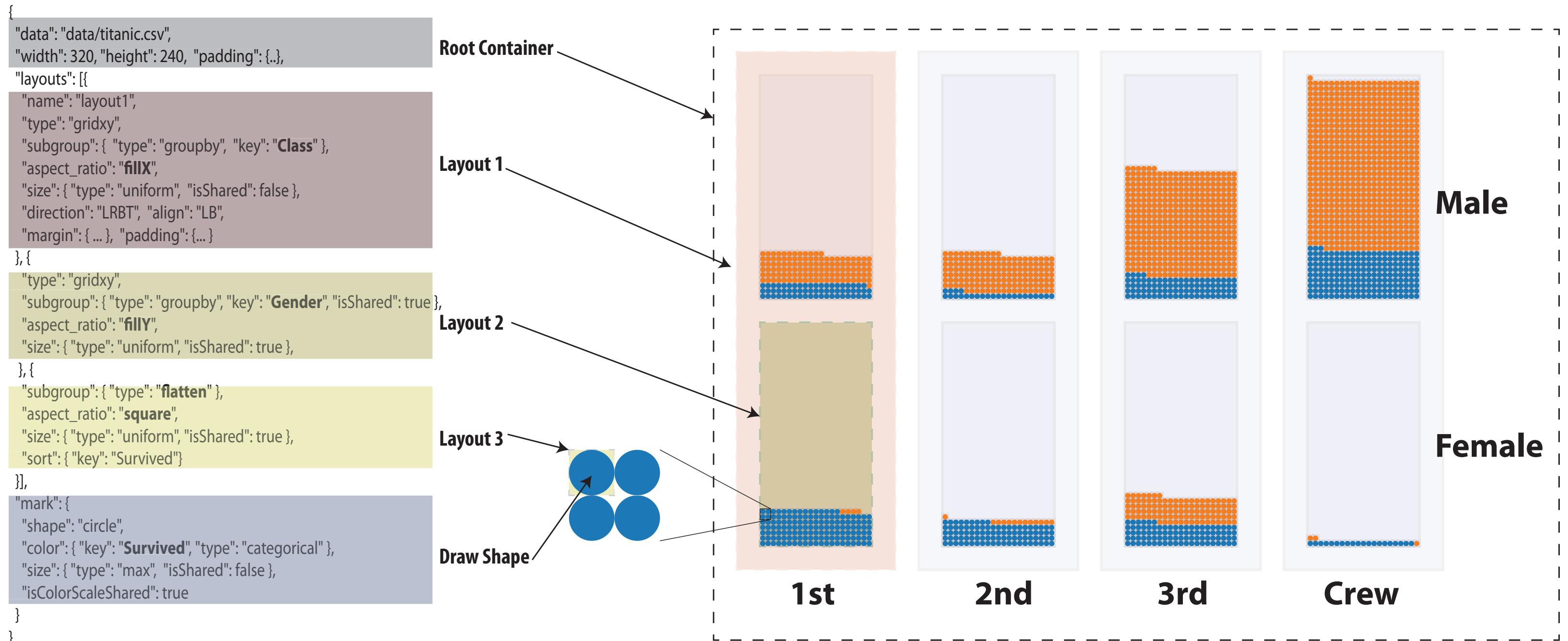
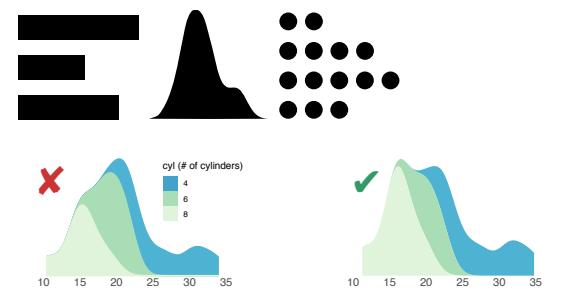


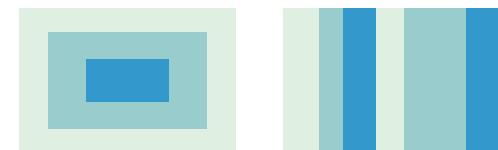
Fig. 6. Example grammar to generate a unit column chart for survivors of the Titanic by passenger class. (Park et al. 2017)

# Related: how to specify a visualization in general

## Grammar of Graphics



## Layout-based

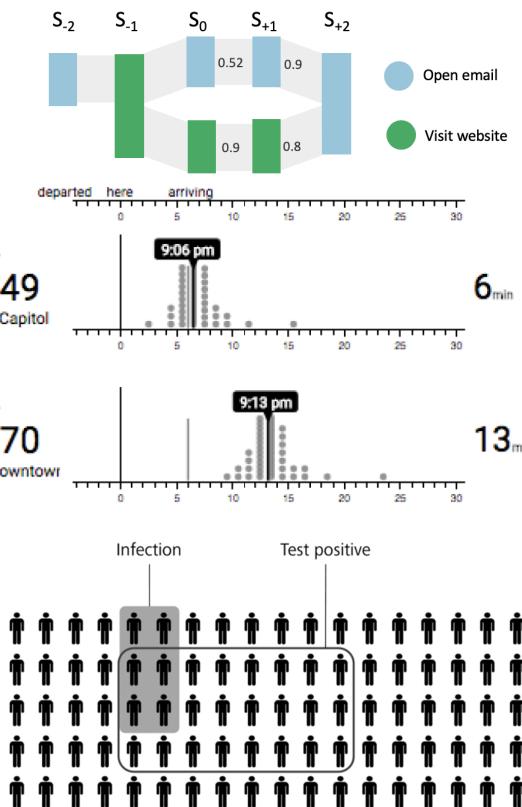


## Constraints-based

```
encoding(e1).  
:- not channel(e1,x).  
:- not field(e1,horsepower).  
:- not bin(e1,_).
```

(Moritz et al. 2019)

Correct?  
Easy to use/  
explore  
designs?



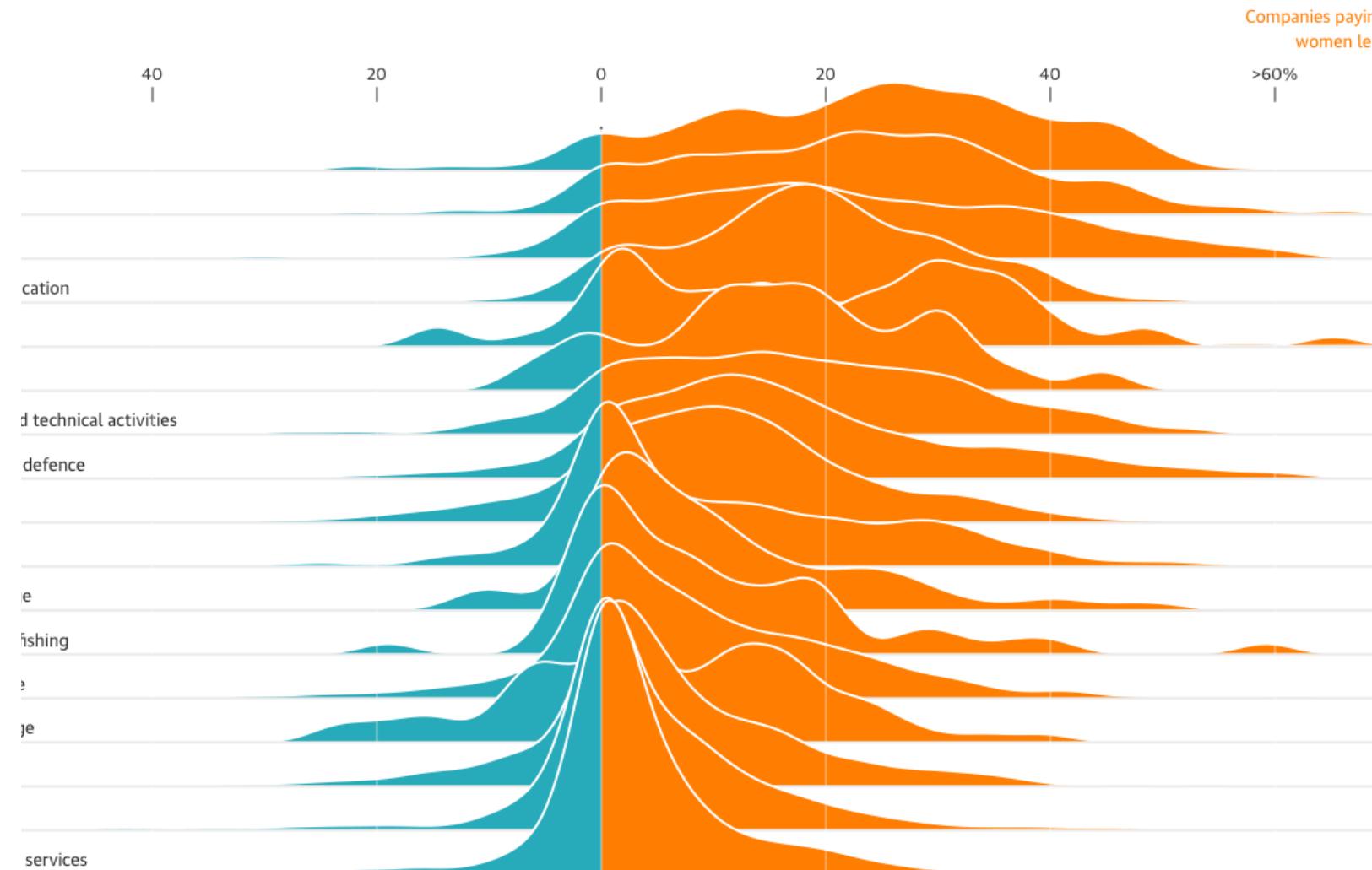
Need a closer integration between statistics and visualization

(Heer and Shneiderman 2012)

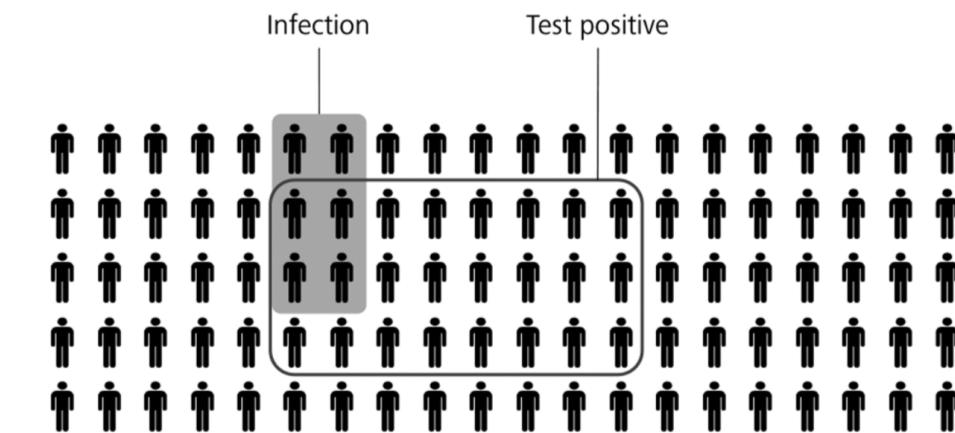
# Related: communicating/visualizing uncertainty

Probabilistic visualizations  
are often used to communicate uncertainty data

**Women are likely to be underpaid in certain sectors**



<https://www.theguardian.com/news/ng-interactive/2018/apr/05/women-are-paid-less-than-men-heres-how-to-fix-it>

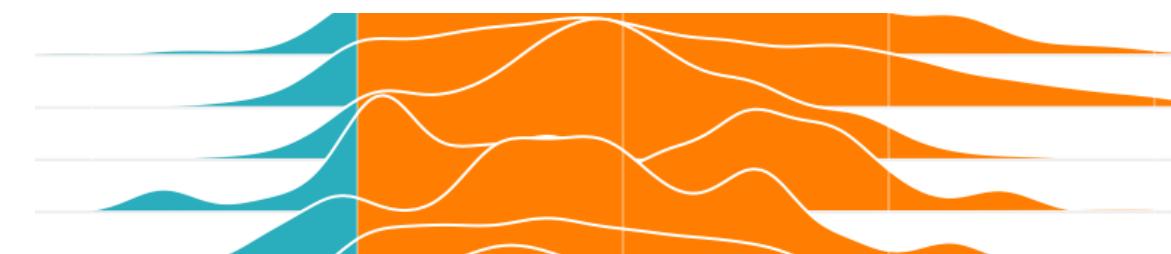
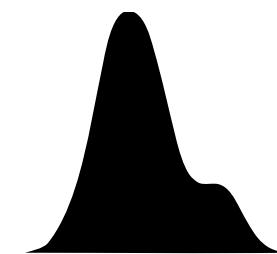


# Related: communicating/visualizing uncertainty

## Probabilistic visualizations

are often used to communicate uncertainty data

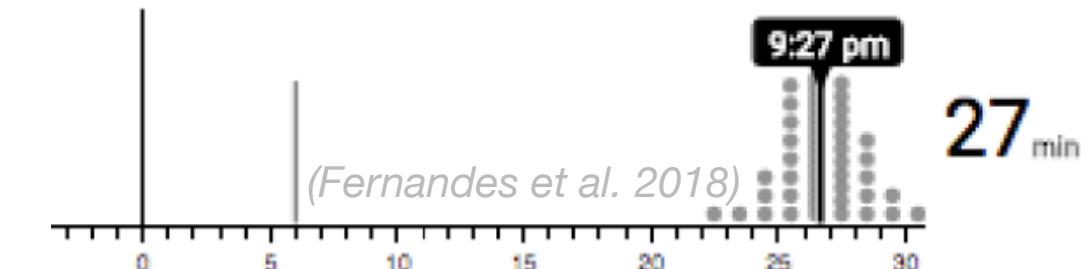
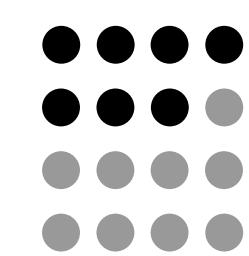
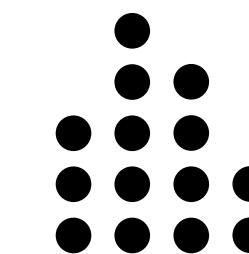
*Probability format    X%*



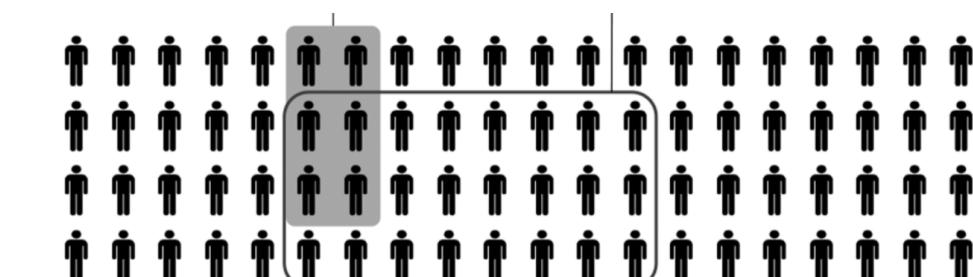
(Gigerenzer and Hoffrage 1995)

(Guo et al. 2019)

*Frequency format    X-in-100*



(Fernandes et al. 2018)

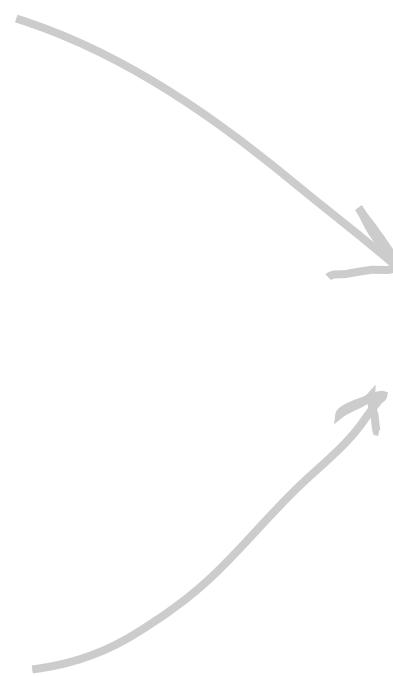


(Binder, Krauss, and Bruckmaier 2015)

*Which one to choose?*

A closer integration  
between statistics and  
visualization

A need to support  
probability and  
frequency formats



## Design Requirements for a probabilistic Grammar of Graphics

A closer integration  
between statistics and  
visualization

- 
- 1 Guaranteeing correctness  
of distributions expressed in  
visualization
  - 2 Enabling specification close to  
probability expressions, such  
as  $P(A|B)$ , which target users  
know

## Design Requirements for a probabilistic Grammar of Graphics

A need to support  
probability and  
frequency formats

- 
- 3 Facilitating exploration *with coherent and reusable grammar components*
  - 4 (and automation in the future)

## Design Requirements for a probabilistic Grammar of Graphics



# The design process

Defaults

Data  $\dashrightarrow A$

Aesthetics  $\rightarrow x \leftarrow A$

Layer

Data

Aesthetics

Geom  $\dashrightarrow \text{geom\_bar}$

Stat

Position

Scale



geom\_density



geom\_points



geom\_rect



geom\_...

(Wickham 2010)

# What is the Probabilistic Grammar of Graphics?

Grammar	ggplot2	PGoG
Defaults		
Data	$\text{Data} \dashrightarrow A$	$P(A B, \dots)$
Aesthetics	$\text{Aesthetics} \dashrightarrow x \leftarrow A$	$\text{height} \leftarrow P(A B, \dots)$
Layer		
Data		
Aesthetics		
Geom	$\text{Geom} \dashrightarrow \text{geom\_bar}$	$\text{geom\_bloc}$
Stat		
Position		$\text{geom\_icon}$
Scale		
Coord	$\text{geom\_density}$	
Facet		
	$\text{geom\_points}$	
	$\text{geom\_rect}$	
	$\text{geom\_...}$	

(Wickham 2010)

# What is the Probabilistic Grammar of Graphics?

1. The PGoG **grammar** is an extension to *Grammar of Graphics*
2. Probability distributions are first class citizens (data) and other grammar components (aesthetics and geometries) are theoretically informed.

# PGoG Grammar/*data*

	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

Simple variable

mpg

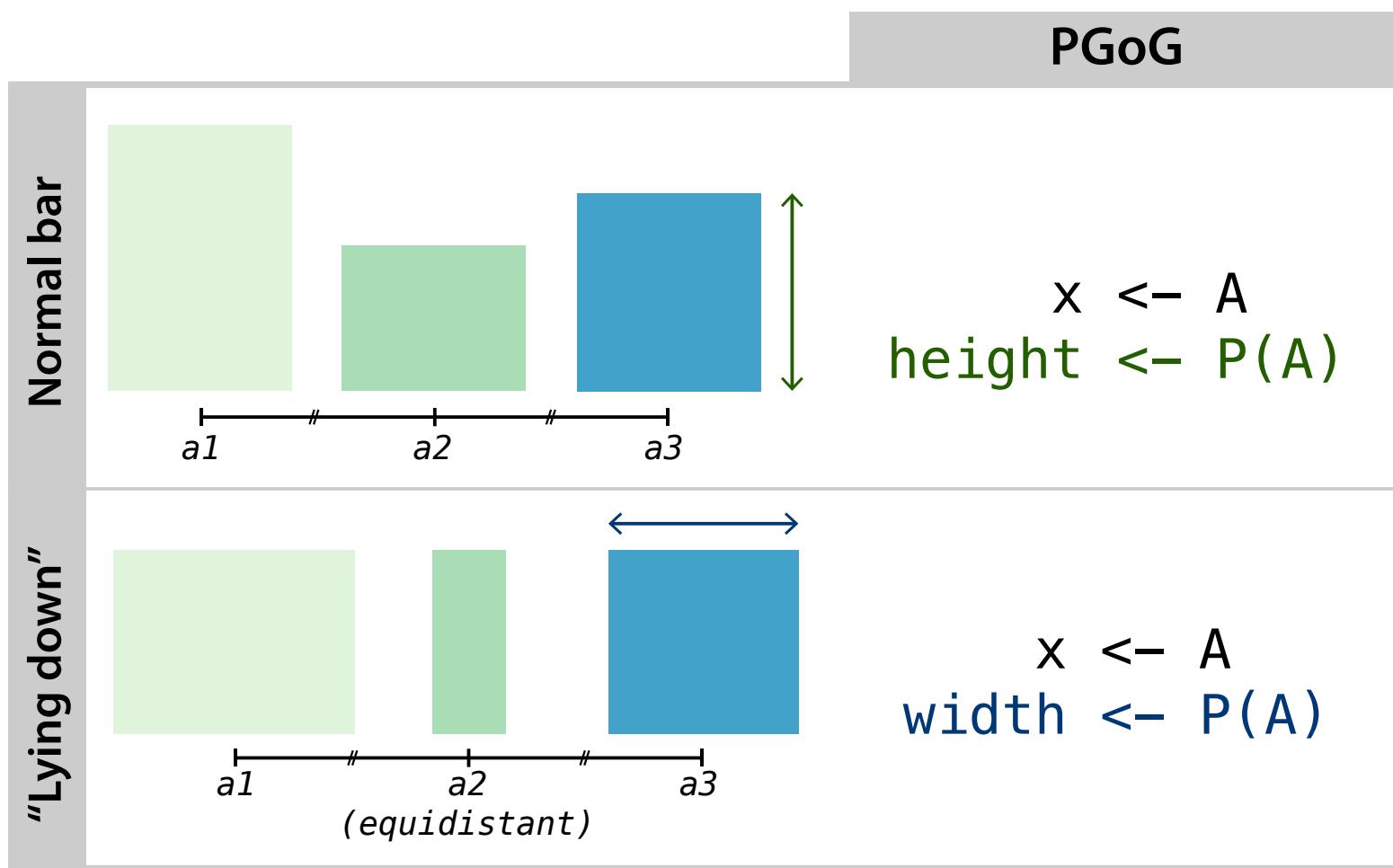
A column in tidy dataset

Probabilistic variable

$P(\text{mpg} | \text{cyl})$

In the form of  $P(A...|B...)$ , where A, B and ... are variables in columns

# PGoG Grammar/aesthetics 1/3

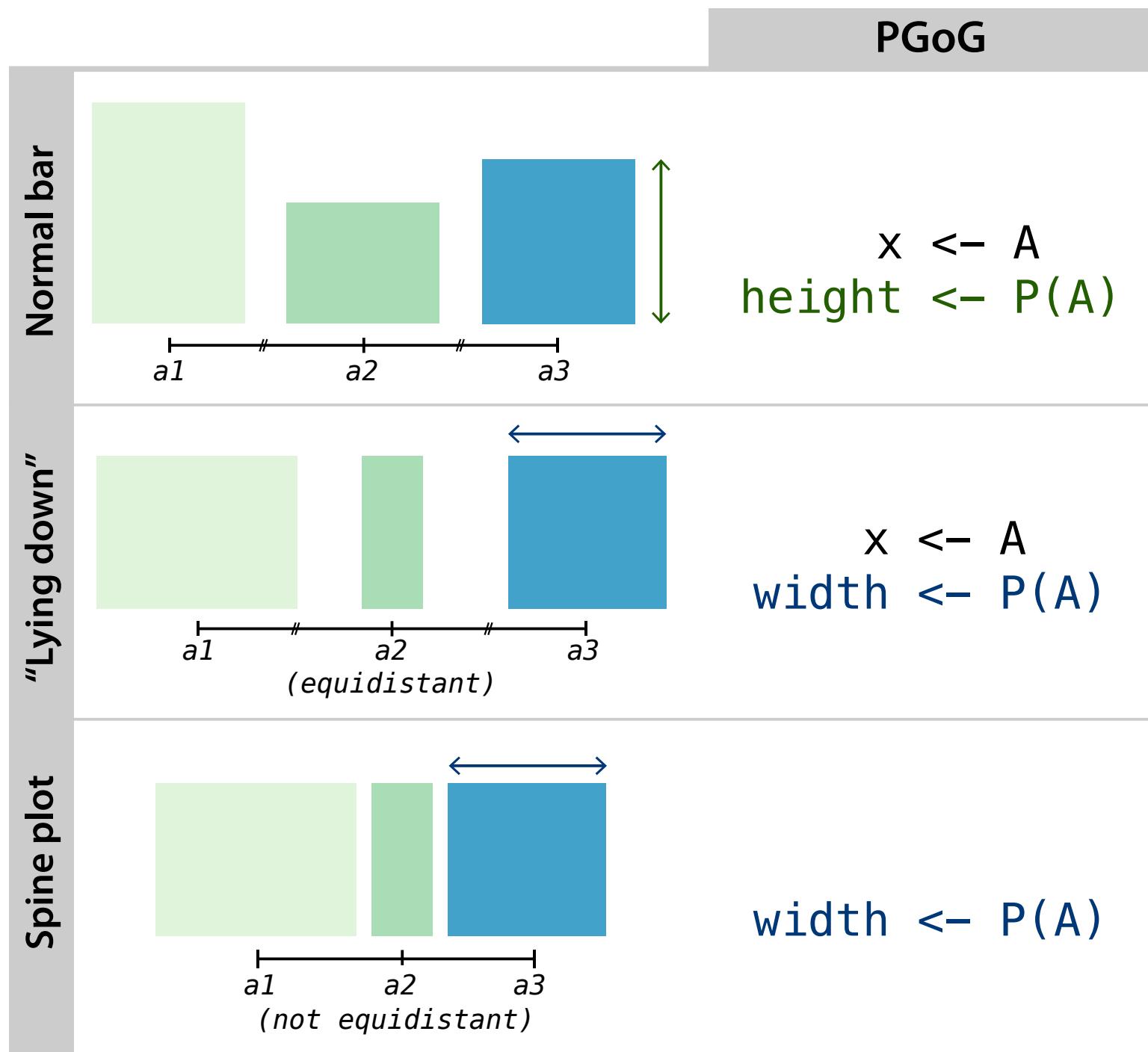


Probabilistic aesthetics

**width, height**

- Works with probabilistic variables only
- Expresses the probability value by length

# PGoG Grammar/aesthetics 2/3



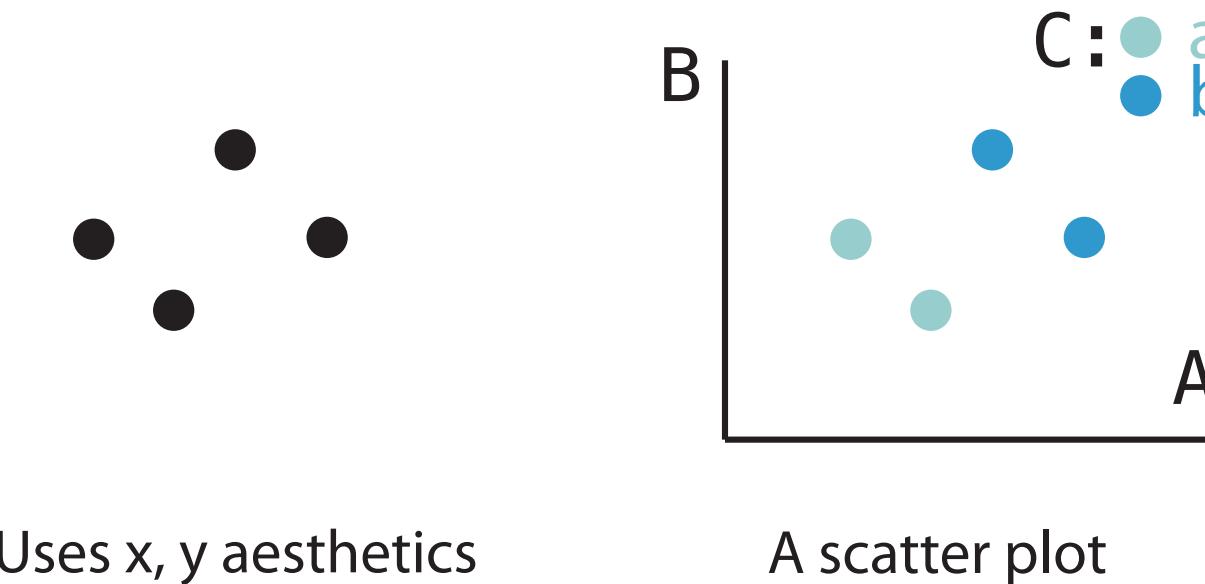
Probabilistic aesthetics

Coordinate aesthetics

**X, Y**

- For discrete vars: equidistant partitions
- For continuous vars: as one would expect

# PGoG Grammar/aesthetics 3/3



Uses x, y aesthetics

A scatter plot

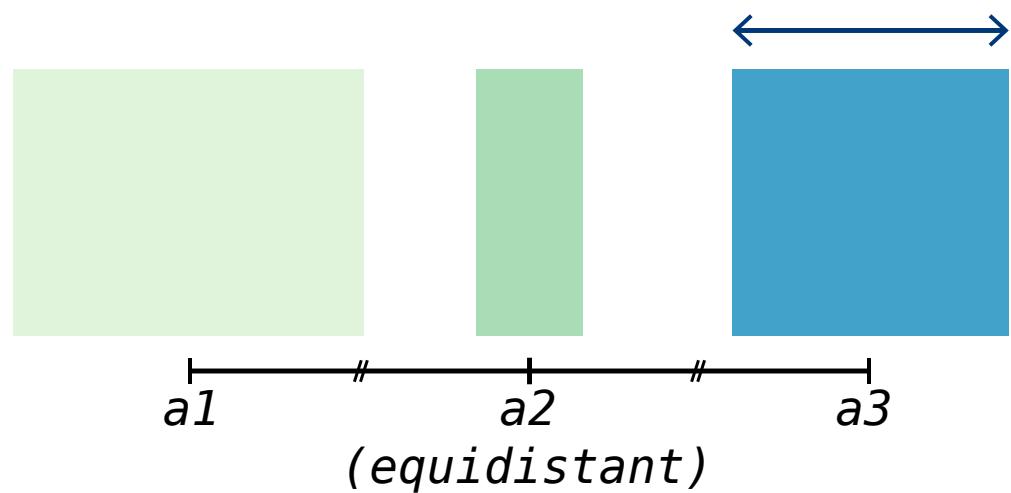
Probabilistic aesthetics

Coordinate aesthetics

Visual aesthetics

fill, color, alpha, ...

# PGoG Grammar/Example for conditional

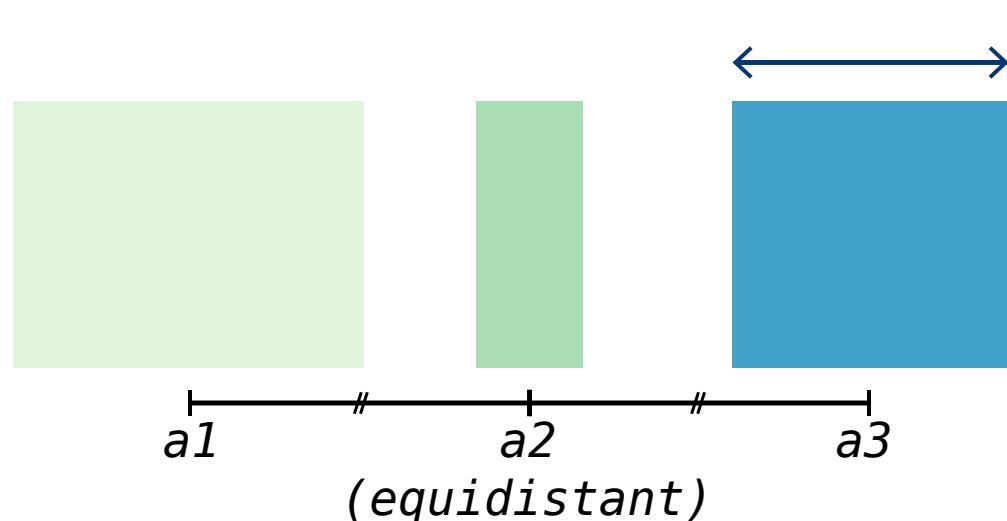


$x \leftarrow A$   
 $\text{width} \leftarrow P(A)$

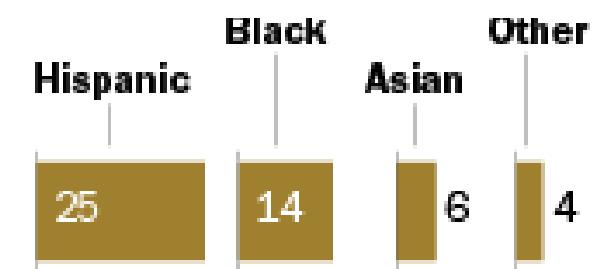
$x \leftarrow \text{race}$

$\text{width} \leftarrow P(\text{race} | \text{generation})$

# PGoG Grammar/Example for conditional



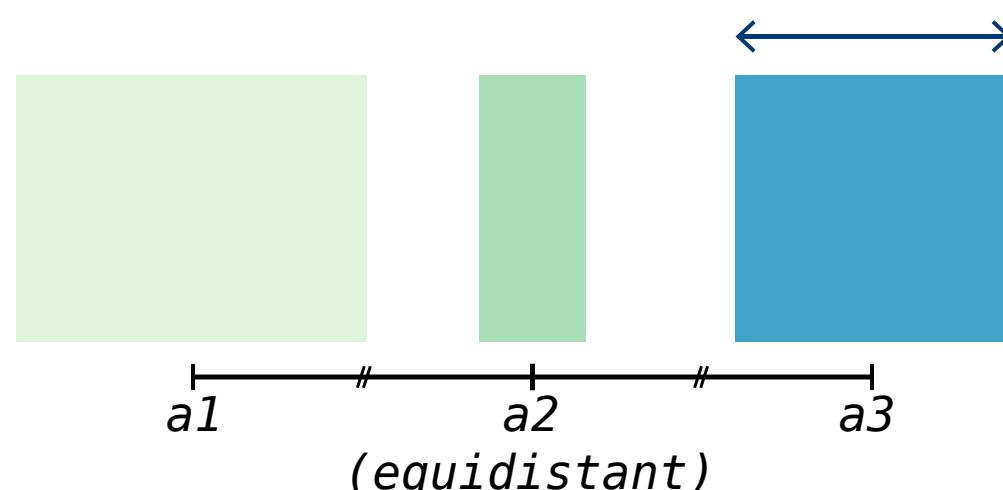
```
x <- A  
width <- P(A)
```



$x <- \text{race}$

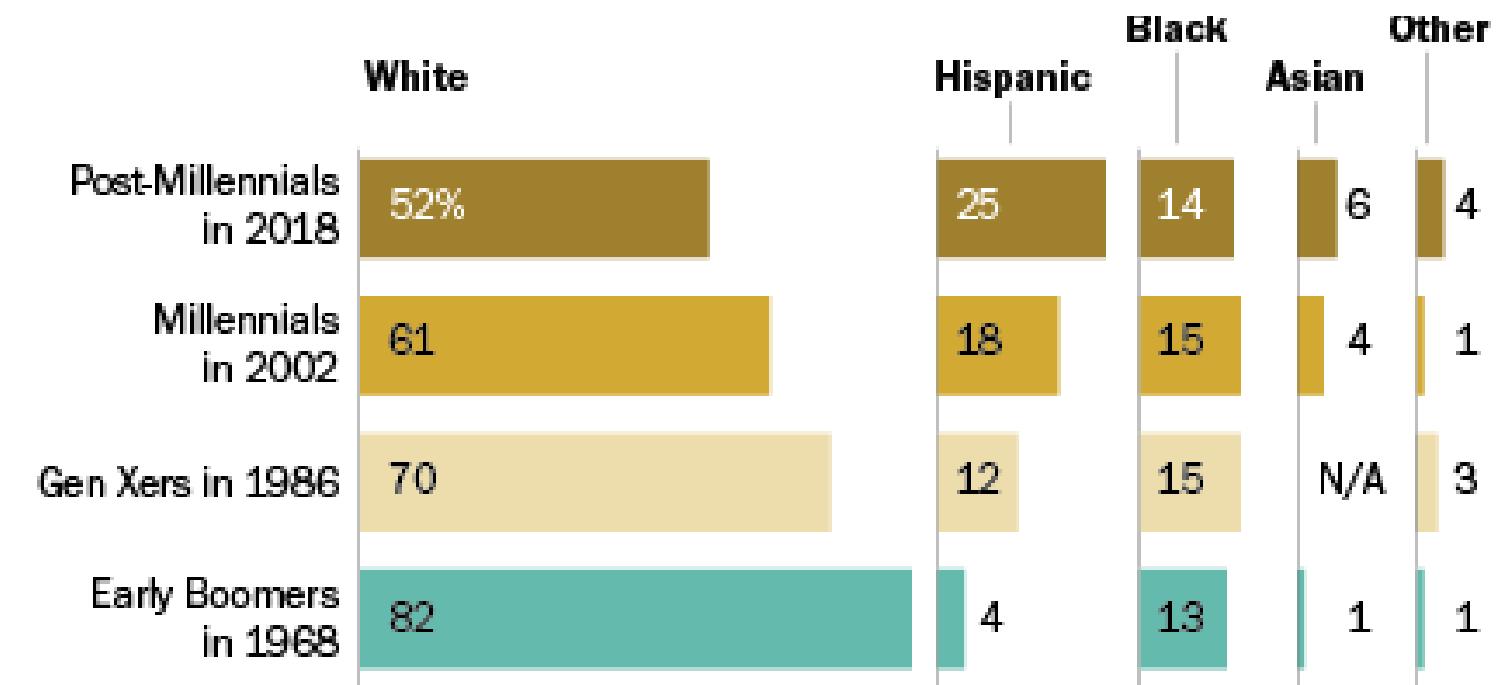
$\text{width} <- \text{P}(\text{race}|\text{generation})$

# PGoG Grammar/Example for conditional



$x \leftarrow A$   
 $width \leftarrow P(A)$

<http://www.pewresearch.org/fact-tank/2018/12/13/18-striking-findings-from-2018/>



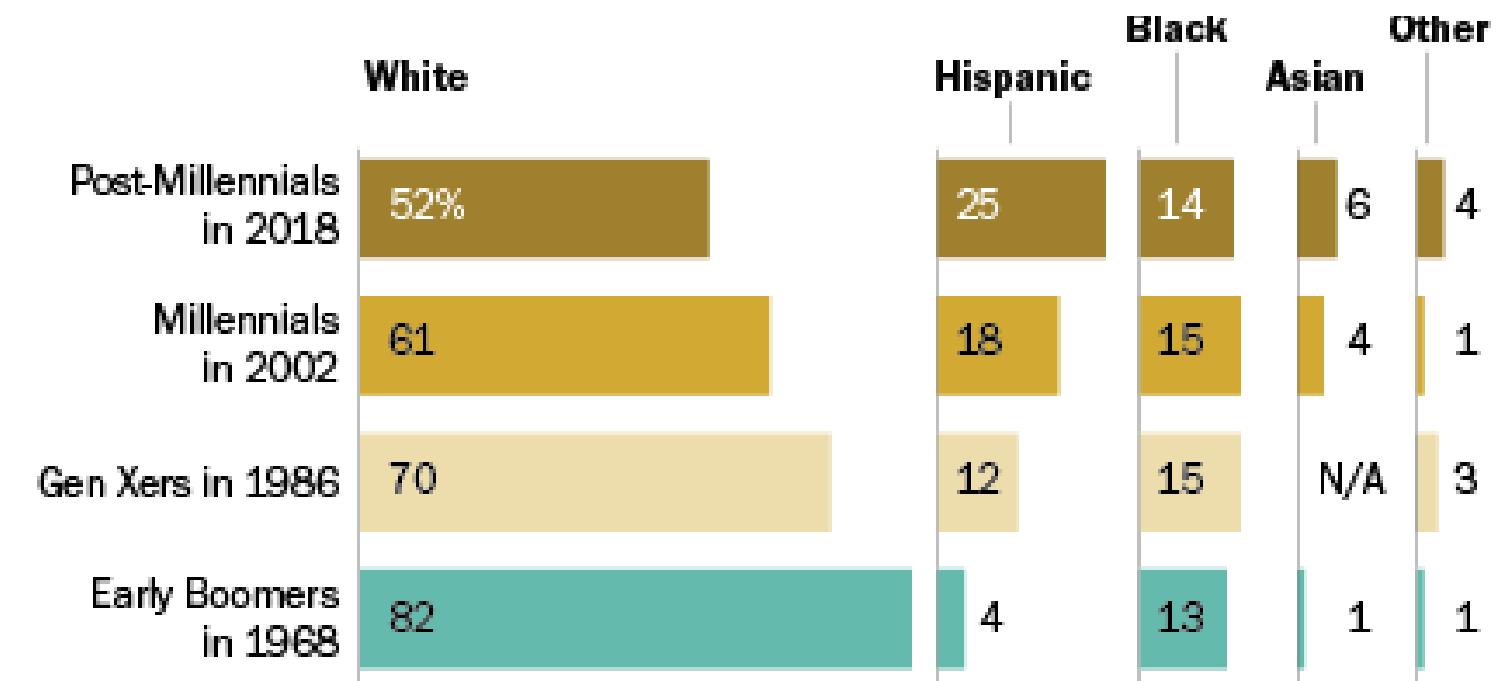
$x \leftarrow \text{race}$   
 $y \leftarrow \text{generation}$   
 $width \leftarrow P(\text{race}|\text{generation})$

# PGoG Grammar/Example for conditional

```
population %>%
  group_by(generation, race) %>%
  summarize(count.gen = n()) %>%
  mutate( prop.cyl = count.gen/
    sum(count.gen)) %>%
  ggplot(aes(x = factor(generation), y =
prop.gen)) +
  geom_col() +
  coord_flip() +
  facet_grid(.~race)
```

**Base ggplot**

<http://www.pewresearch.org/fact-tank/2018/12/13/18-striking-findings-from-2018/>



`x <- race`  
`y <- generation`  
`width <- P(race|generation)`

**PGoG**

# PGoG Grammar/*Example for joint*

Math

$$P(\text{cyl} \mid \text{mpg}) \ P(\text{mpg}) = P(\text{mpg}, \text{cyl})$$

# PGoG Grammar/*Example for joint*

Math

$$P(\text{cyl} \mid \text{mpg}) \ P(\text{mpg}) = P(\text{mpg}, \text{cyl})$$

Coord aes

```
x <- mpg
```

# PGoG Grammar/*Example for joint*

Math

$$P(\text{cyl} \mid \text{mpg}) \ P(\text{mpg}) = P(\text{mpg}, \text{cyl})$$

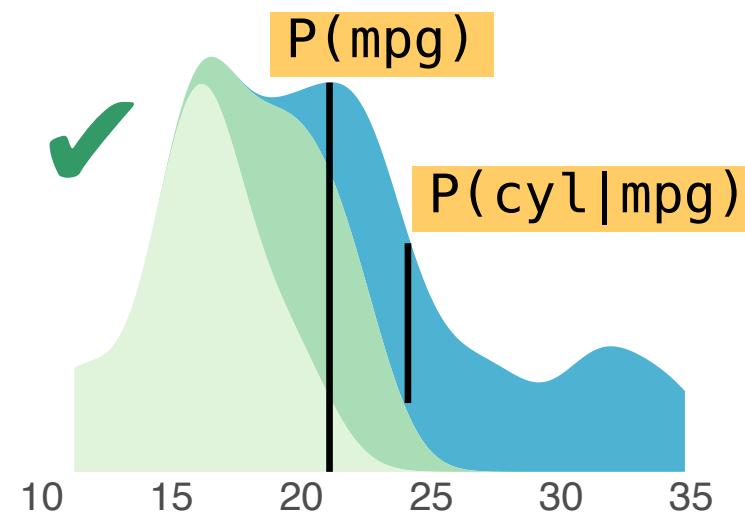
Coord aes

```
x <- mpg
```

Prob aes

```
height <- P(mpg) P(cyl | mpg)
```

# PGoG Grammar/Example for joint



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

Math

Coord aes  
Prob aes  
Visual aes

$$P(\text{cyl}|\text{mpg}) P(\text{mpg}) = P(\text{mpg}, \text{cyl})$$

`x <- mpg`

`height <- P(mpg) P(cyl|mpg)`  
`fill <- cyl`

# PGoG Grammar/*checking correctness* 1/2

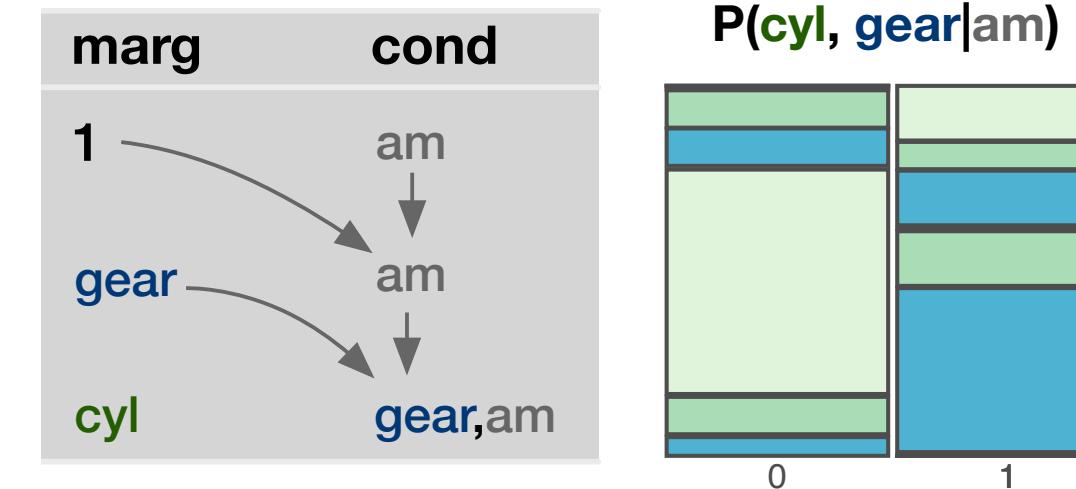
One of the **rules**: the  
probabilistic variables  
need to be valid factors of  
a **probability function**

```
x <- gear  
height <- P(gear|am)  
          P(cyl|gear, am)
```

# PGoG Grammar/*checking correctness*

One of the rules: the probabilistic variables need to be valid factors of a **probability function**

$x \leftarrow \text{gear}$   
height  $\leftarrow \begin{cases} P(\text{gear} | \text{am}) \\ P(\text{cyl} | \text{gear}, \text{am}) \end{cases}$



A “chain” data structure used for checking probabilistic variables

Grammar

ggplot2

PGoG

Defaults

Data  $\dashrightarrow A$

Aesthetics  $\dashrightarrow x \leftarrow A$

Layer

Data

Aesthetics

Geom  $\dashrightarrow \text{geom\_bar}$



Stat  
Position

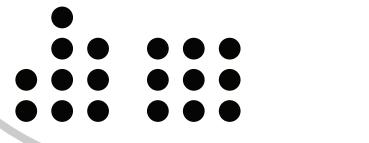
Scale



geom\_bloc



geom\_density



geom\_icon



geom\_points



Coord

geom\_rect



Facet

geom\_...

PGoG

$P(A|B, \dots)$

height  $\leftarrow P(A|B, \dots)$

PGoG Grammar/  
*geometries 1/2*

---

**ggplot2**

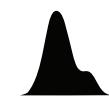
geom\_bar



geom\_mosaic\*



geom\_density



geom\_violin

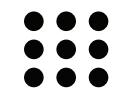


geom\_density\_ridges\*



---

geom\_waffle\*



**PGoG**

geom\_bloc

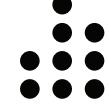
Look at all those geometries  
in **ggplot2** we have  
replaced

Also, probability and  
frequency formats

---

**geom\_icon**

geom\_dotplot

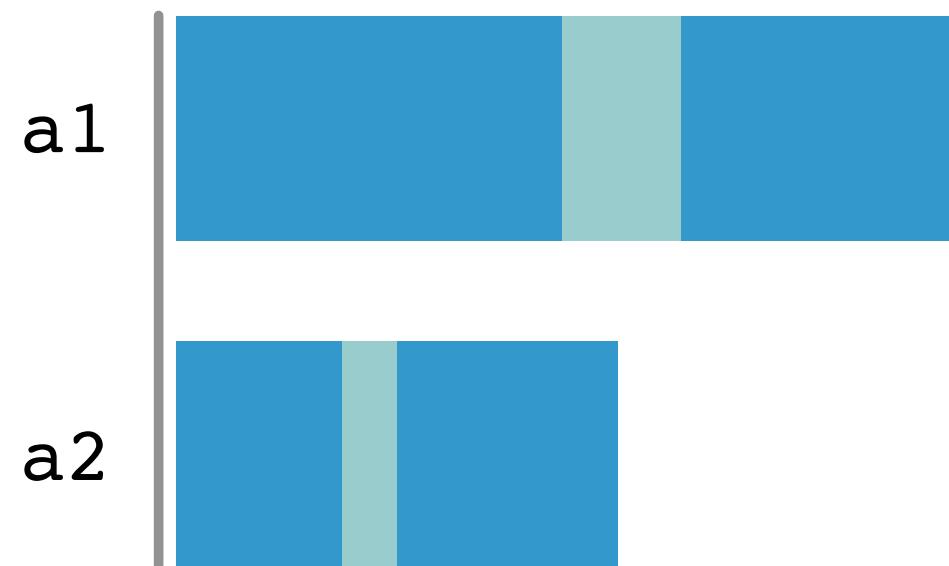


\* **ggplot2** extensions

# PGoG Grammar/*geometries*

`geom_bloc`: recursive sub-partition to support many probabilistic variables

```
y <- A  
width <- P(A) P(B|A)  
fill <- B
```



`geom_icon` needs a new way to pack icons

```
y <- A  
width <- P(A) P(B|A)  
fill <- B
```

# PGoG Grammar/*geometries*

`geom_bloc`: recursive sub-partition to support many probabilistic variables

```
y <- A  
width <- P(A) P(B|A)  
fill <- B
```



`geom_icon` needs a new way to pack icons

```
y <- A  
width <- P(A) P(B|A)  
fill <- B
```

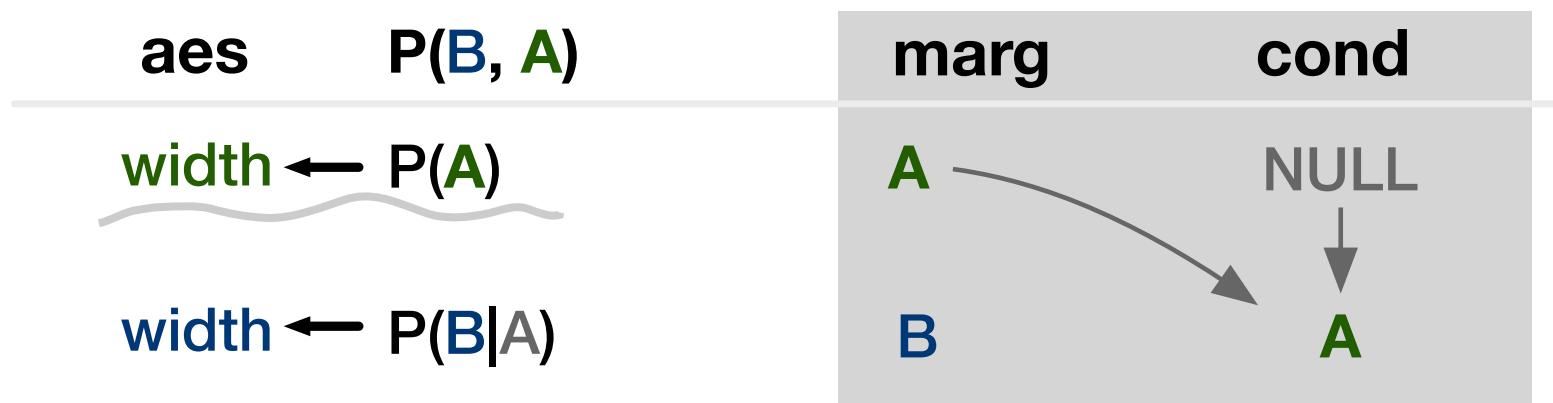
aes	P(B, A)	marg	cond
	width <- P(A)	A	NULL
	width <- P(B A)	B	A

Probability structure (the “chain”) determines the visualization structure

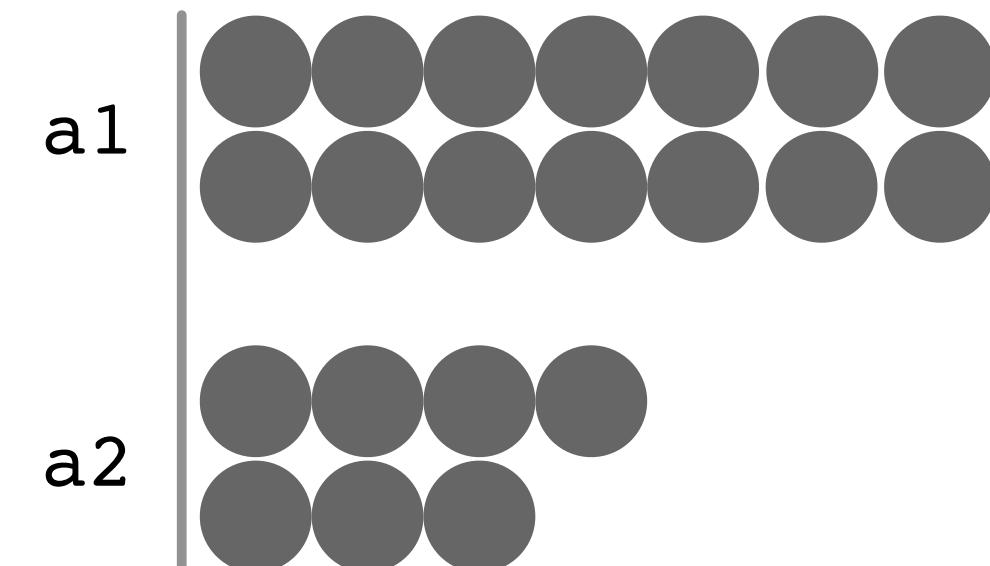
# PGoG Grammar/*geometries*

`geom_icon` needs a new way to pack icons

```
y <- A  
width <- P(A) P(B|A)  
fill <- B
```

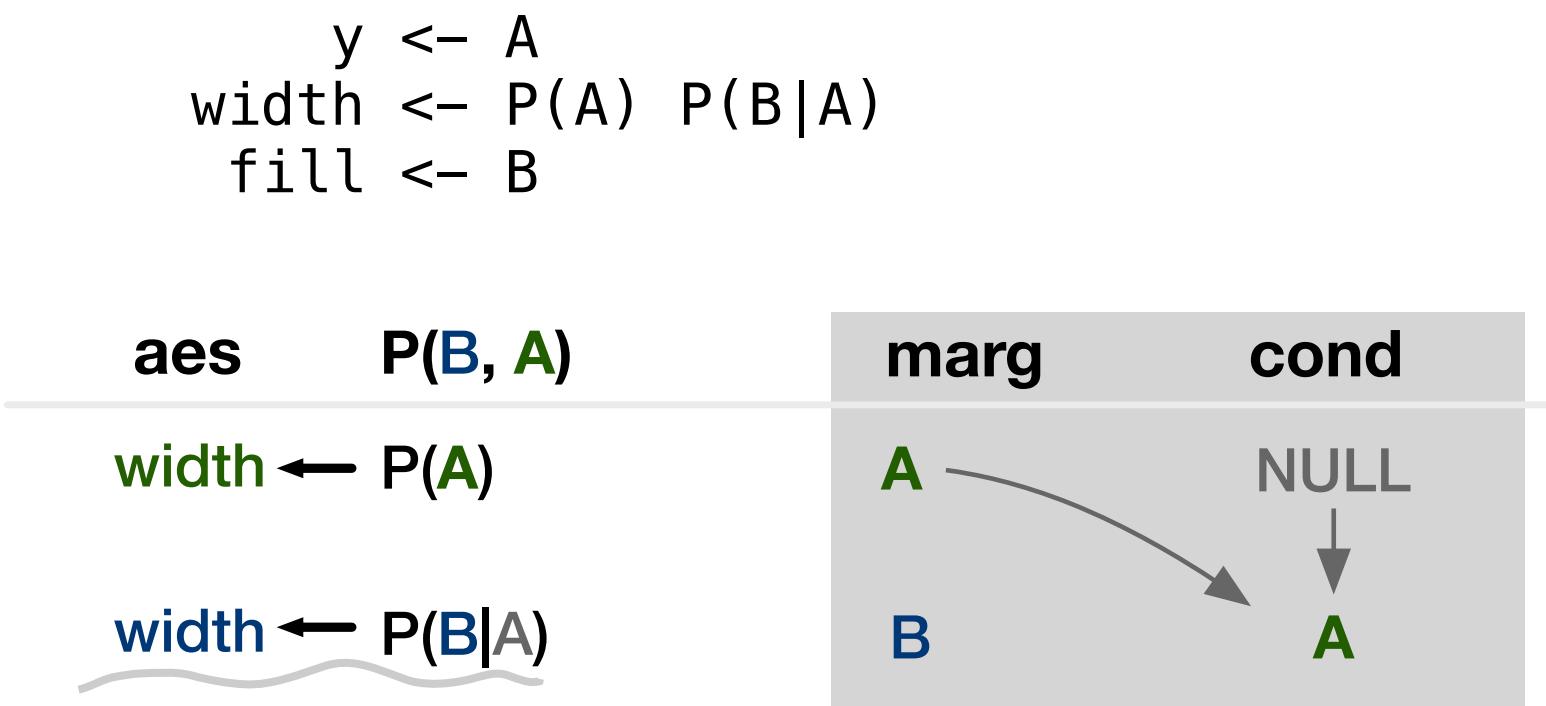


Probability structure (the “chain”) determines the visualization structure

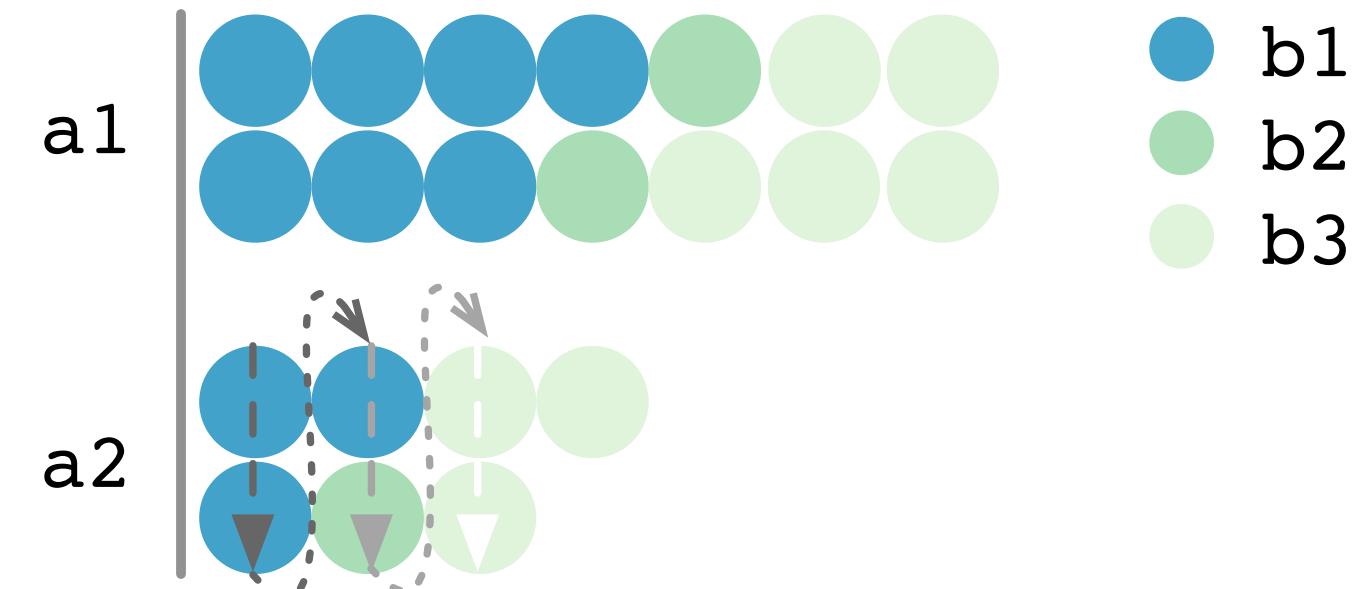


# PGoG Grammar/*geometries*

geom\_icon needs a new way to pack icons



Probability structure (the “chain”) determines the visualization structure



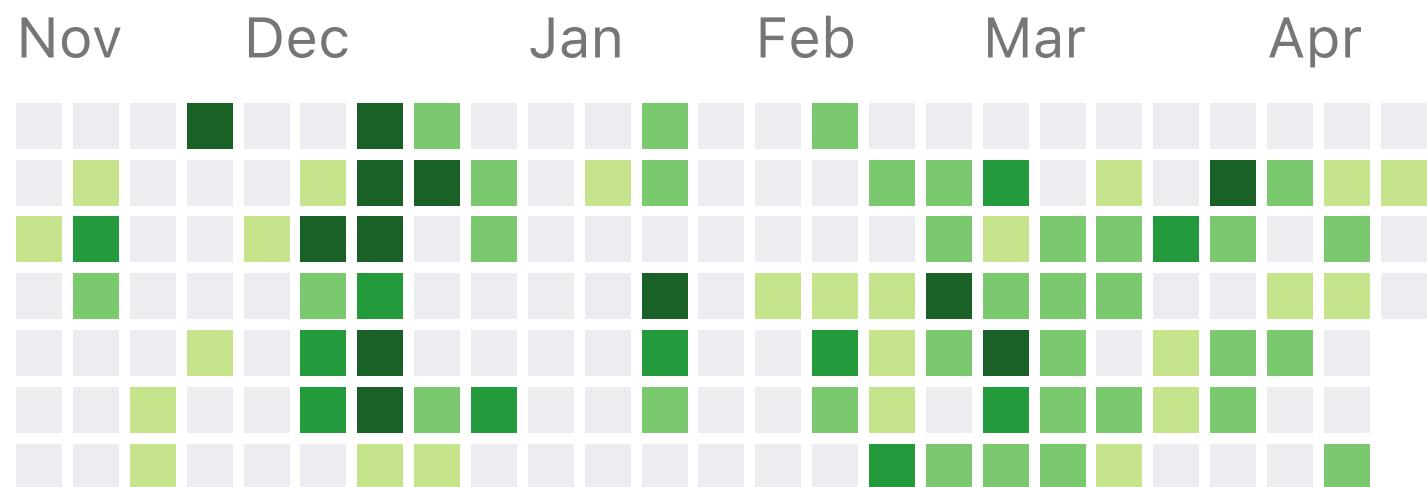
# Implementing the Grammar

## *Why in R*

- Grammar of graphics is implemented in ggplot2
- Metaprogramming features in R helps parsing
- PGoG grammar is transferable

## *Current progress*

- geom\_bloc with discrete variables
- geom\_bloc with up to two continuous variables
- geom\_icon with up to two variables
- TODO: coloring, extra parameters



# Evaluation of the Grammar

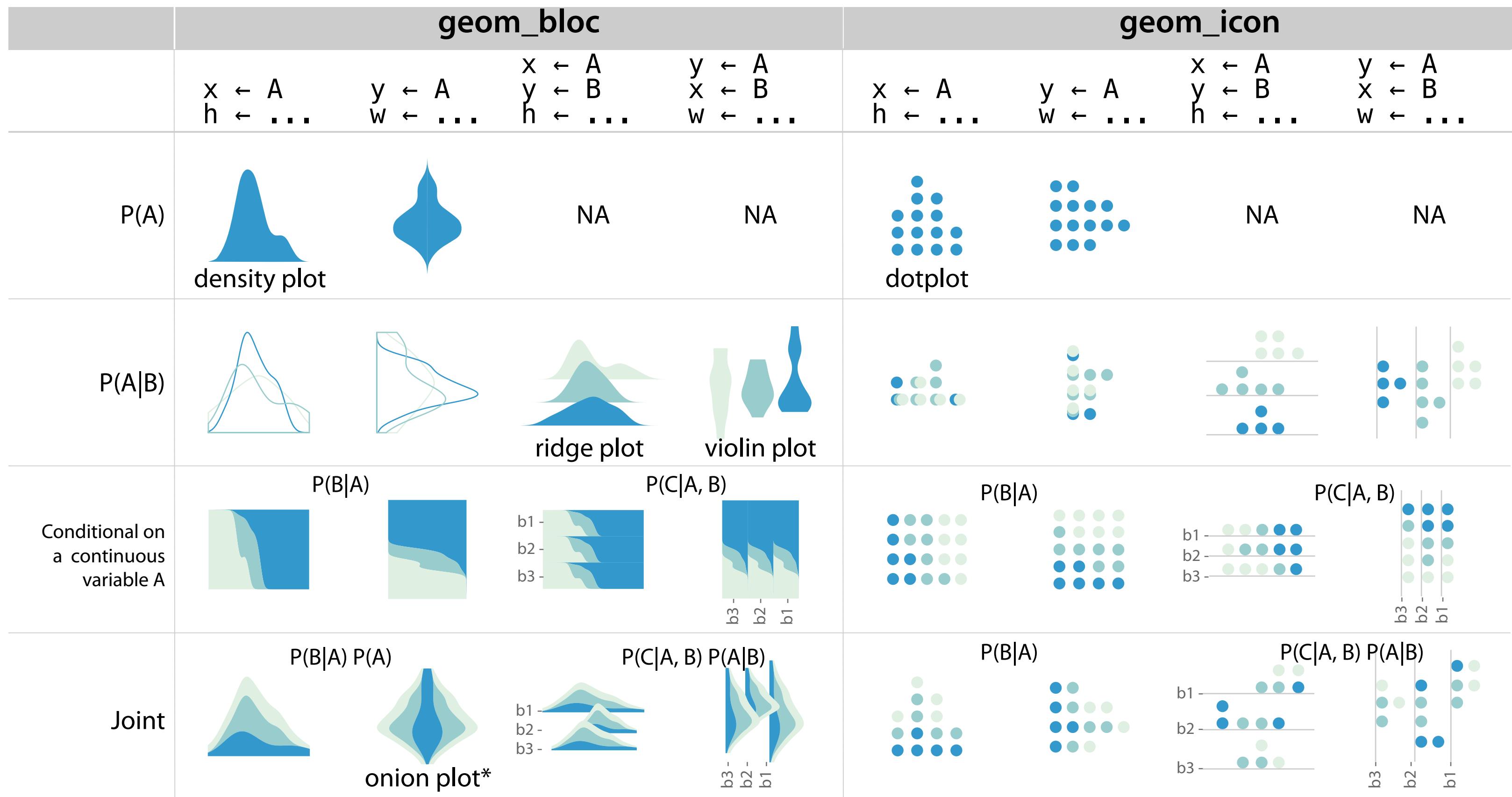
*Usefulness:* is it a good visualization grammar

- Expressive?
- Generative?

*Usability:* is it easy to express probabilities as visualizations

- Cognitively ergonomic?

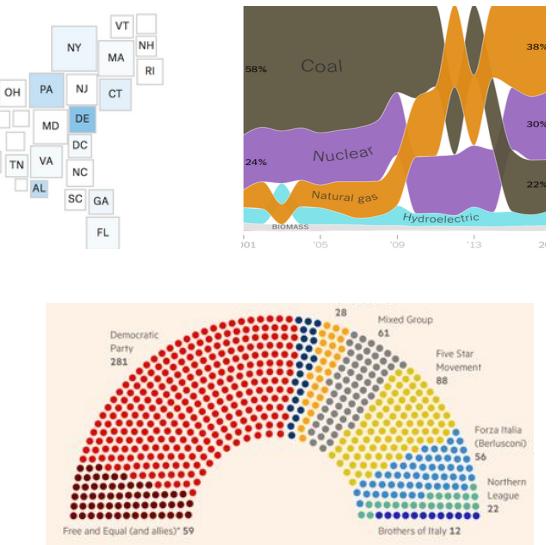
# Expressiveness of the grammar



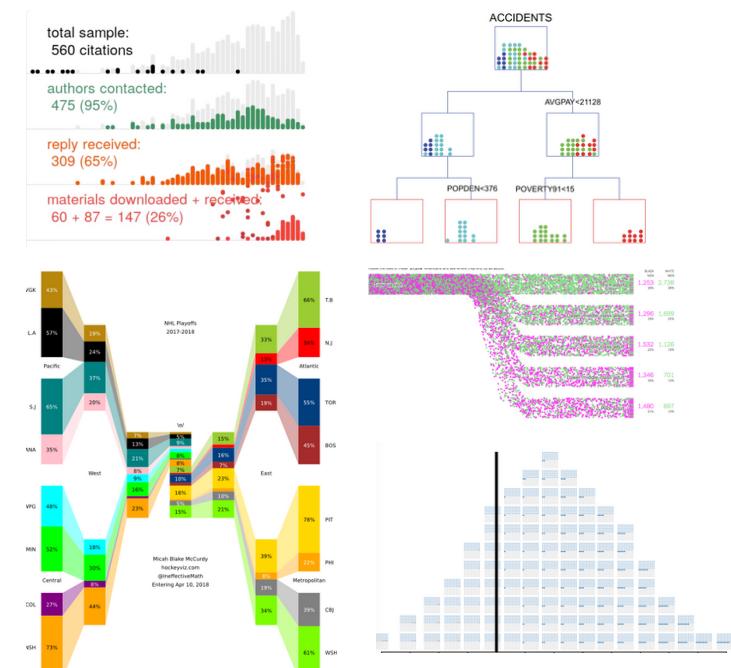
# Expressiveness of the grammar, in the wild

Can reasonably reproduce

Time and space

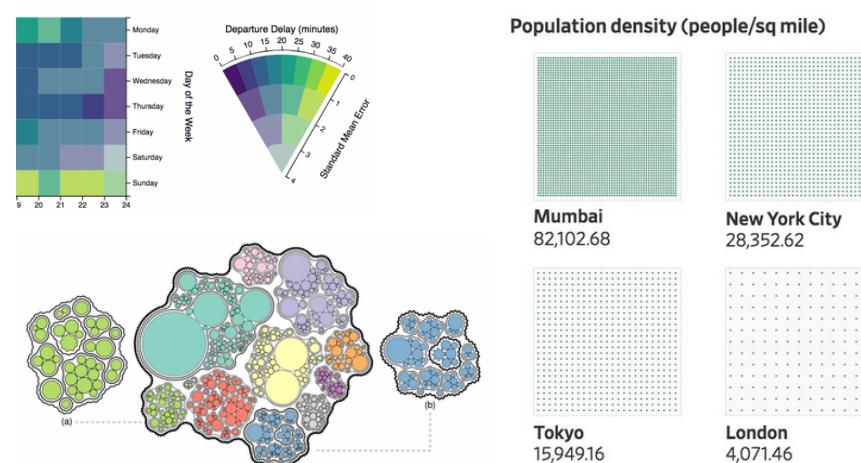


Hierarchy/structure

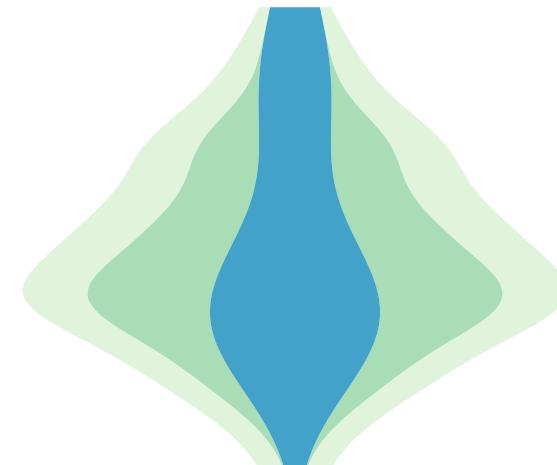


A Probabilistic Grammar of Graphics

Additional channels



# Generativeness from the combination of aesthetics



## Onion plot

`geom_bloc:`

`y ← mpg`

`width ← P(mpg) P(cyl|mpg)`

`direction ← both`

# Cognitive ergonomics

(Blackwell et al. 2001)

- ... concerns the usability of notational systems
- Evaluated with *Cognitive Dimensions of Notations*

*Pro:*

Short edit distances (Kim et al. 2017)

- *Low viscosity*
- *No premature commitment*
- Close to probability expressions

*Con:*

Specifying probability expressions can be difficult

- *Hidden dependencies*
- *Error prone-ness*

# Cognitive ergonomics

*Pro:* Short edit distances, close to probability expressions

- *Low viscosity*
- *No premature commitment*

Existing ggplot2 packages

Changes

Syntax



---

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

# Cognitive ergonomics

Pro: Short edit distances, close to probability expressions

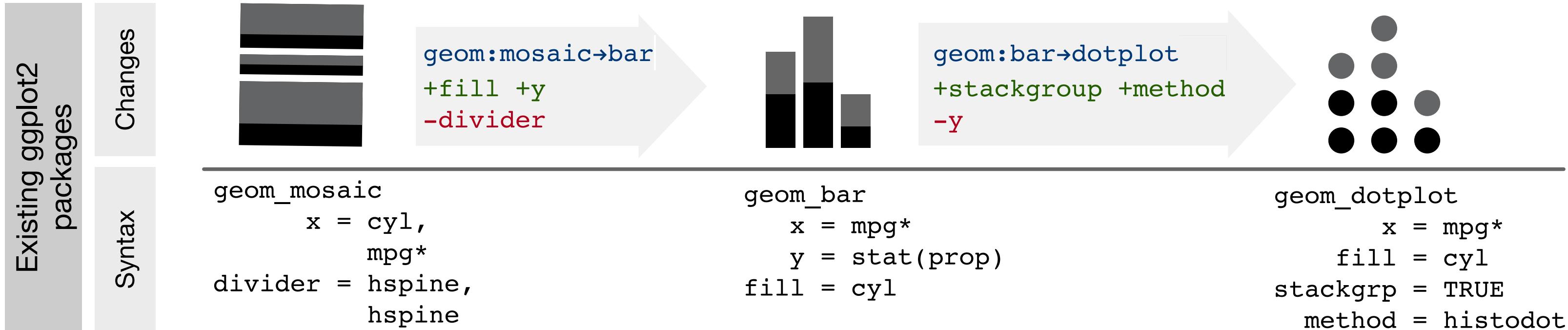
- *Low viscosity*
- *No premature commitment*



# Cognitive ergonomics

Pro: Short edit distances, close to probability expressions

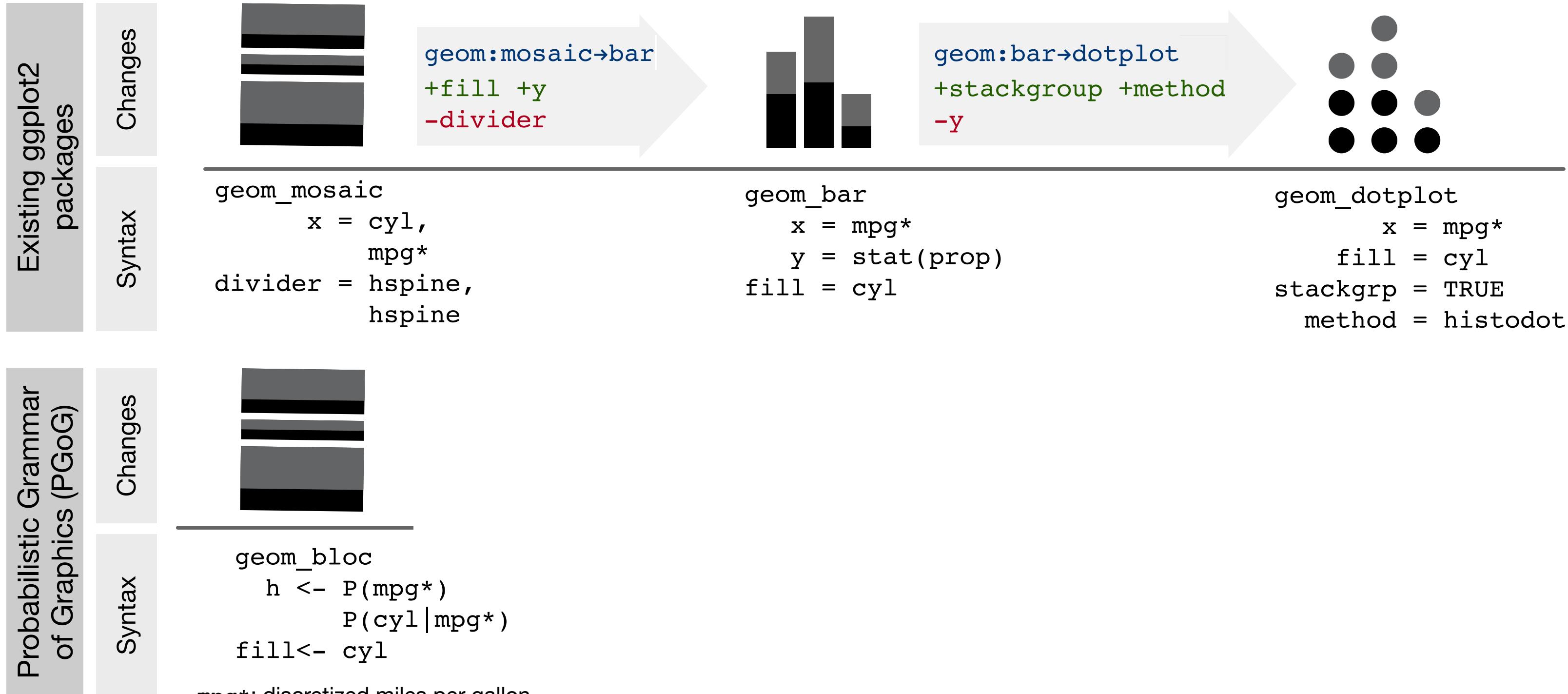
- *Low viscosity*
- *No premature commitment*



# Cognitive ergonomics

Pro: Short edit distances, close to probability expressions

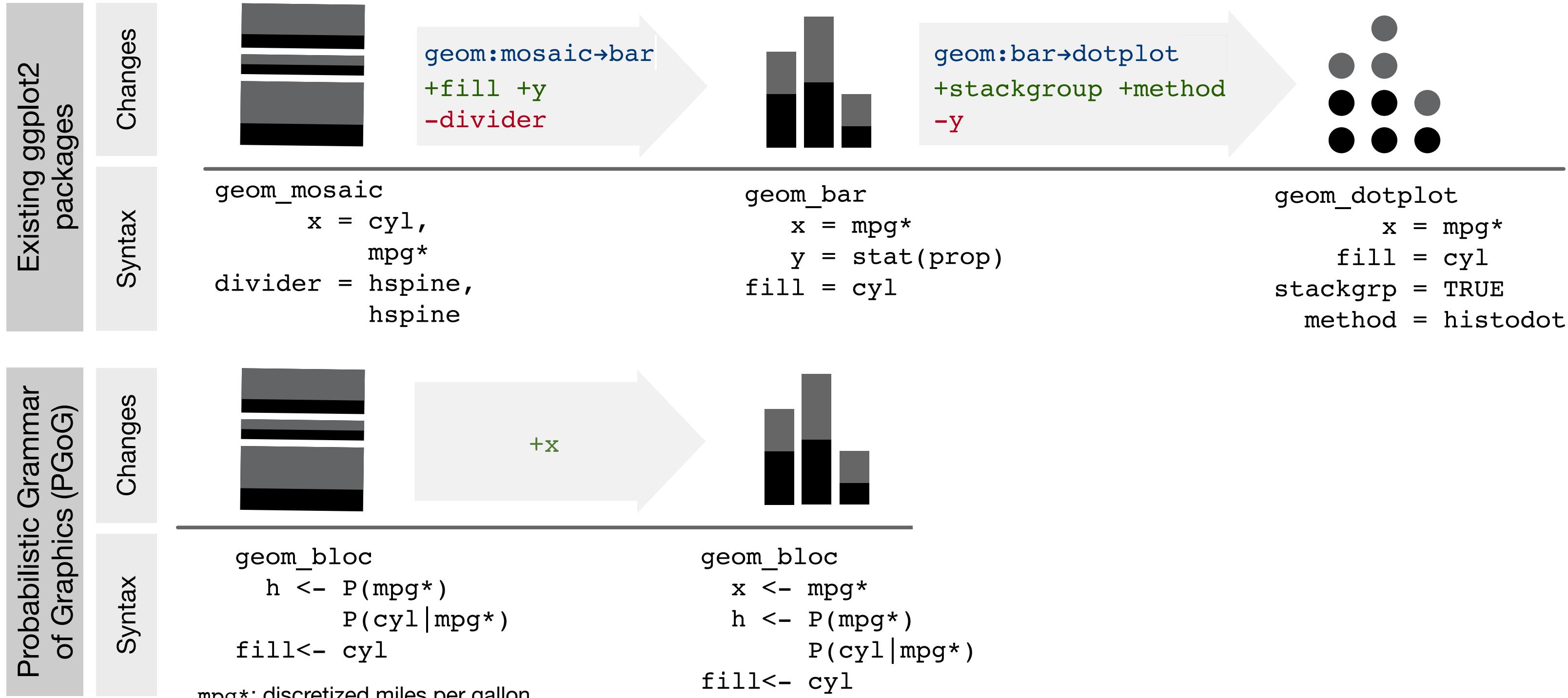
- *Low viscosity*
- *No premature commitment*



# Cognitive ergonomics

Pro: Short edit distances, close to probability expressions

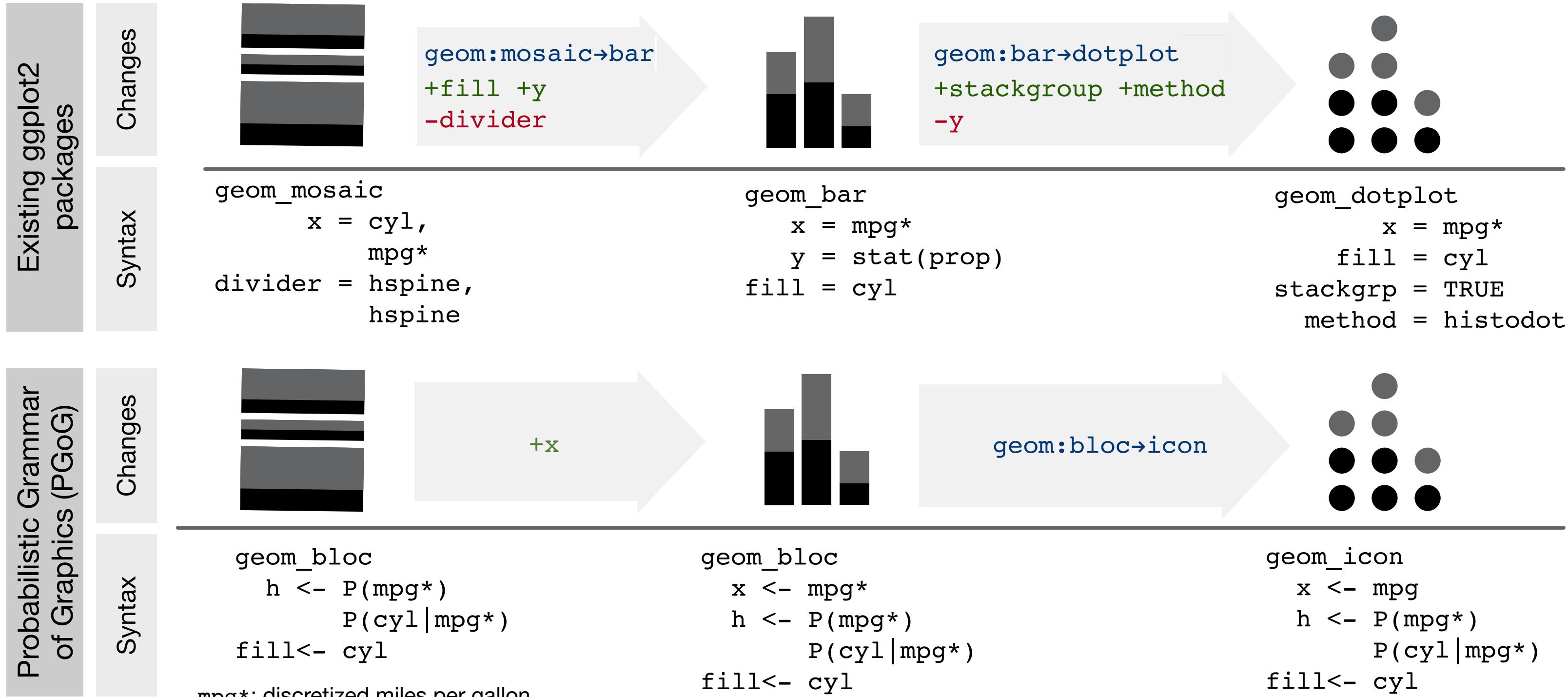
- *Low viscosity*
- *No premature commitment*



# Cognitive ergonomics

Pro: Short edit distances, close to probability expressions

- *Low viscosity*
- *No premature commitment*



# Cognitive ergonomics

*Con:* specifying probability expressions can be difficult

- *Error prone-ness \**
- *Hidden dependencies*

Math

$$P(\text{cyl}|\text{mpg}) \ P(\text{mpg}) = P(\text{mpg}, \text{cyl})$$



$$P(\text{mpg}|\text{cyl}) \ P(\text{mpg}) ?$$

# Cognitive ergonomics

Con: specifying probability expressions can be difficult

- *Error prone-ness*
- *Hidden dependencies*

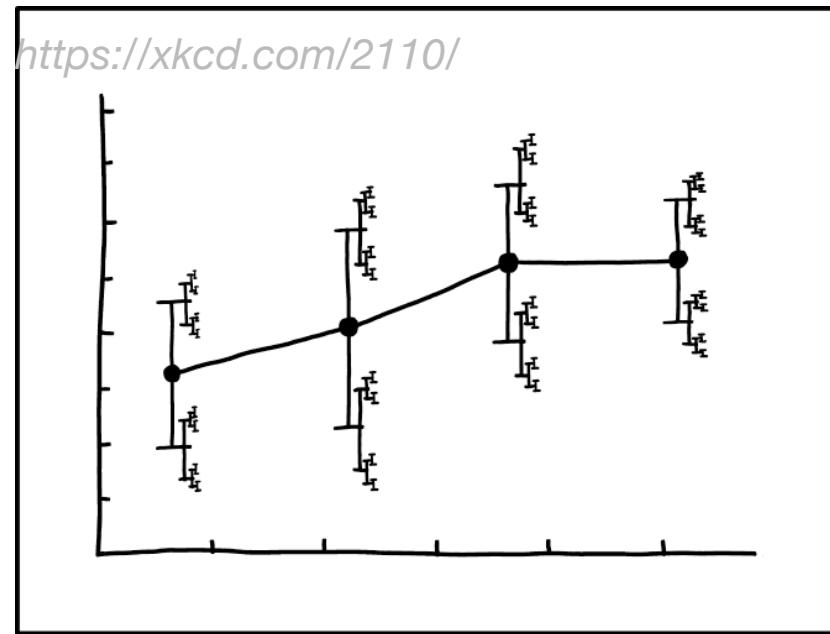
	$P(\text{mpg}   \text{cyl}) \ P(\text{mpg}) ?$
Math	$P(\text{cyl}   \text{mpg}) \ P(\text{mpg}) = P(\text{mpg}, \text{cyl})$
Coord aes	$x \leftarrow \text{mpg}$
Prob aes	$\text{height} \leftarrow P(\text{cyl}   \text{mpg}) \ P(\text{mpg})$
Visual aes	$\text{fill} \leftarrow \text{cyl}$

# Further questions

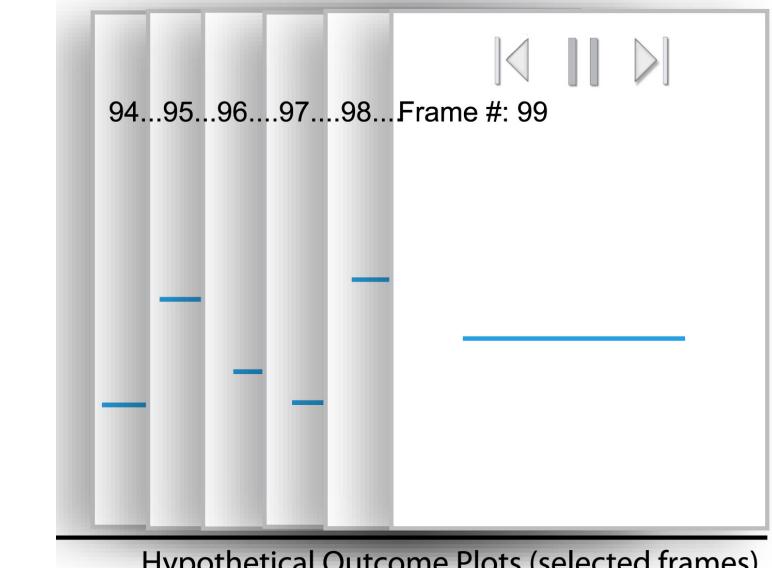
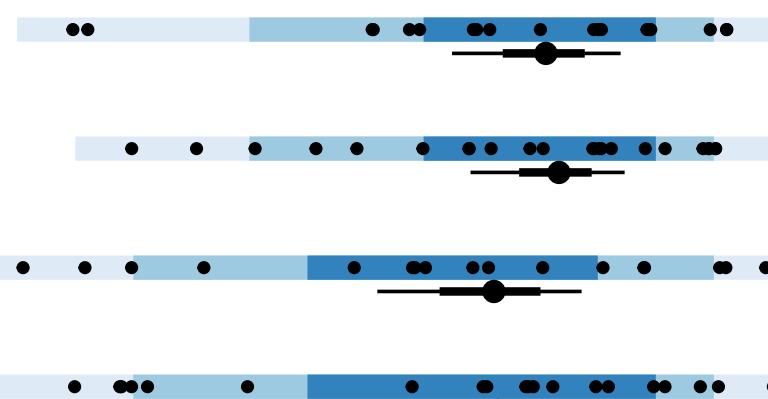
Based on *the Cognitive Dimensions of Notations*:

<i>error-proneness/ closeness of mapping:</i>	“Are users able to factorize and supply probability distributions?”	Rate correctness
<i>premature commitment:</i>	“Do users explore more visualization designs with PGoG than a baseline system”	Count visualizations explored
<i>hidden dependencies:</i>	“Can users learn to use the PGoG components to replicate existing probabilistic visualizations”	Record task completion time, NASA-TLX, rate visualization completion

# Future work: more uncertainty vizes -> more concepts

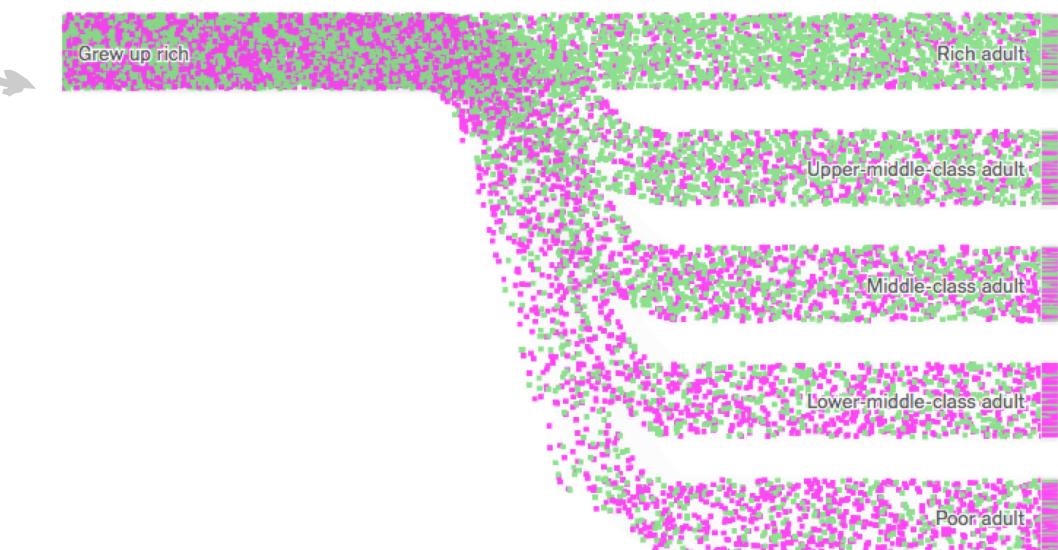


- Uncertainty sources: **aleatory** or epistemic
- Data structure: **hierarchical**, sequential, etc.
- **Summary statistics**, confidence intervals, etc.
- Visualization techniques such as **linking**

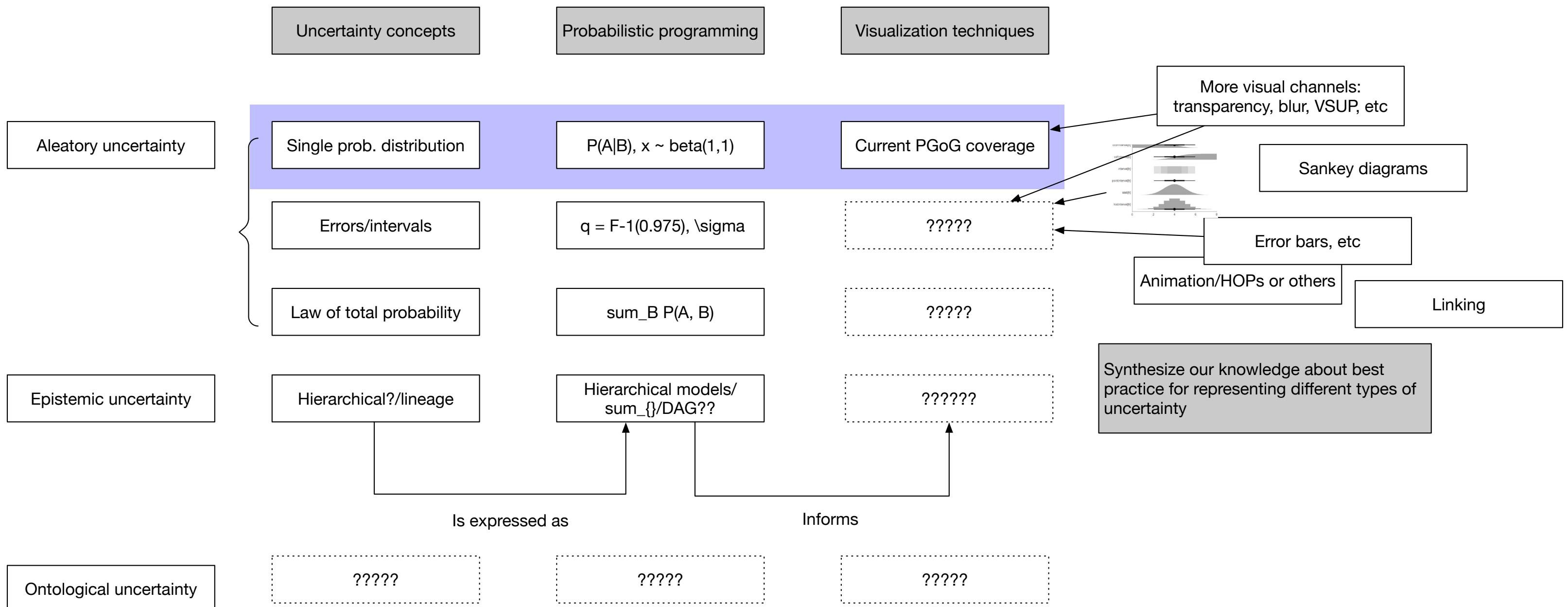


(Hullman, Resnick, and Adar 2015)

[https://www.nytimes.com/interactive/2018/03/27/up-  
shot/make-your-own-mobility-animation.html](https://www.nytimes.com/interactive/2018/03/27/up-shot/make-your-own-mobility-animation.html)

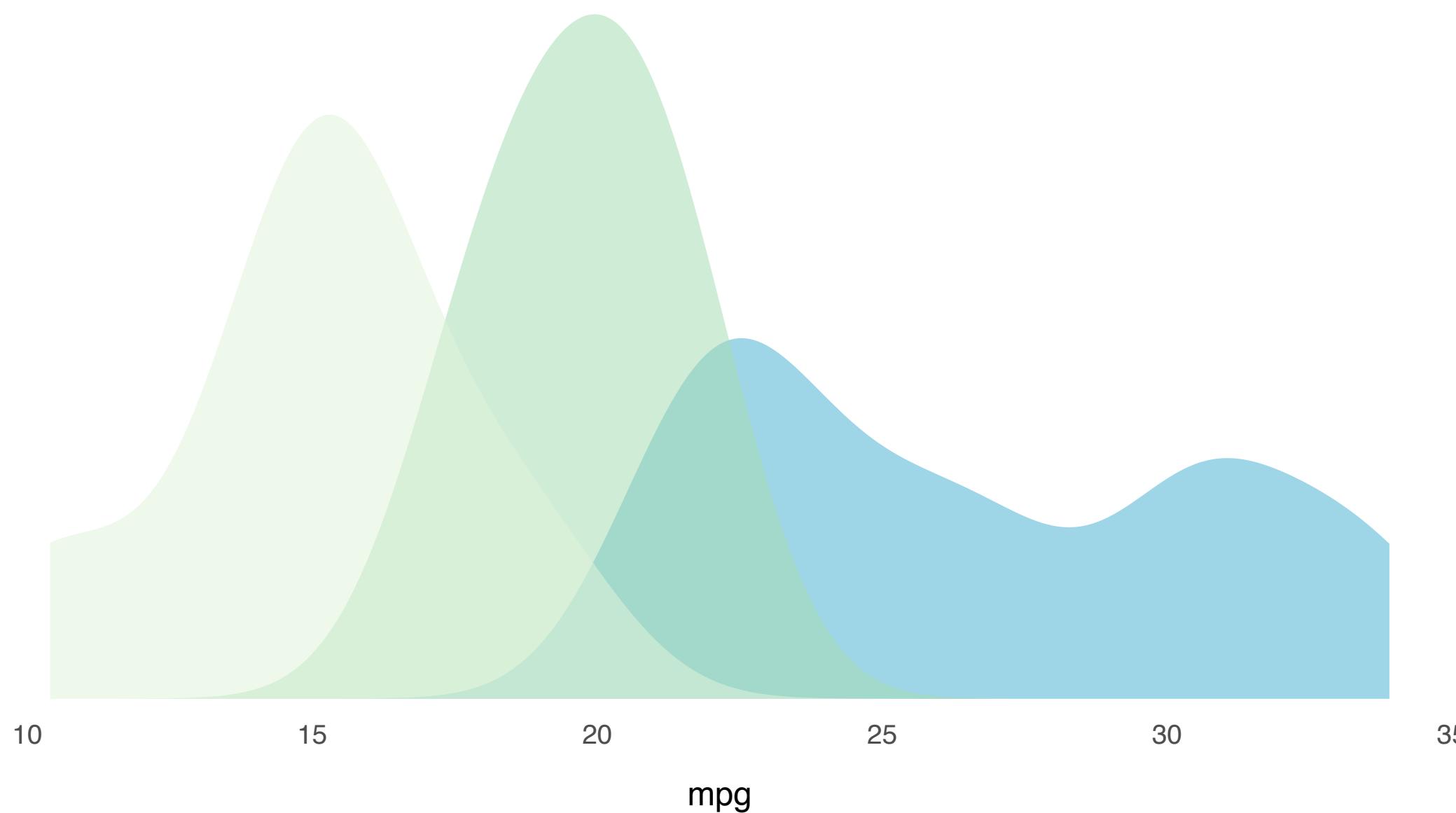


# Future work: an overview inspired by uncertainty taxonomies



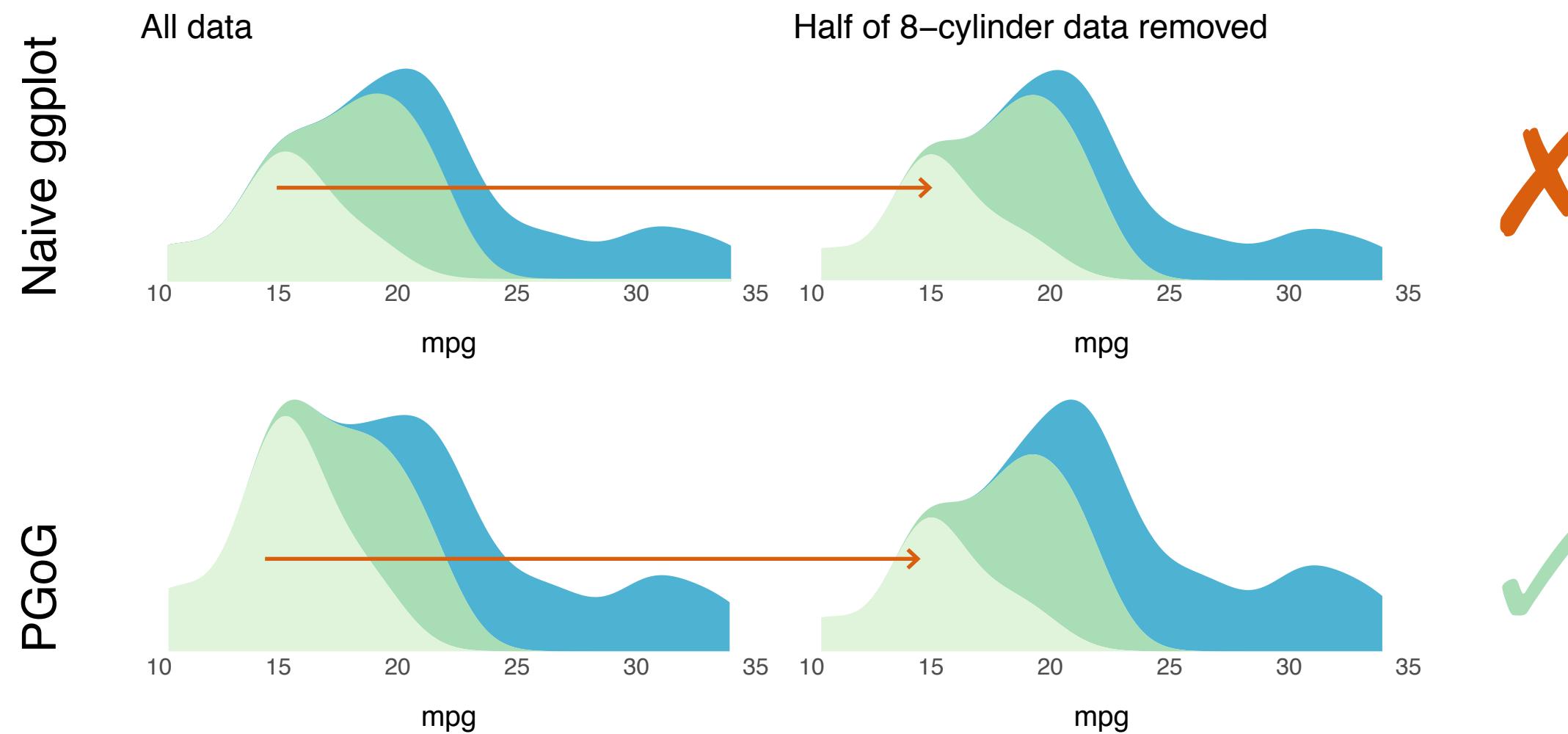
# Conclusions

Probabilistic Grammar of Graphics is a visualization grammar that treats probability distributions as first-class citizens. It **shifts our thinking** about specifications for **probabilistic visualizations** and could facilitate **uncertainty** communication in the future.



The non-ambiguous way of presenting  $P(\text{mpg}|\text{cyl})$

# Using an algebraic process for visualization design [kindlmann\_algebraic\_2014]



$\alpha$ : data symmetry  $\sim$  reduced 7/32 data points

$\omega$ : visualization symmetry  $\sim$  does the visualization change accordingly?

$\alpha \neq \omega$ : violation of the principle of data-visualization correspondence

# 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,
```

# What could possibly go wrong?

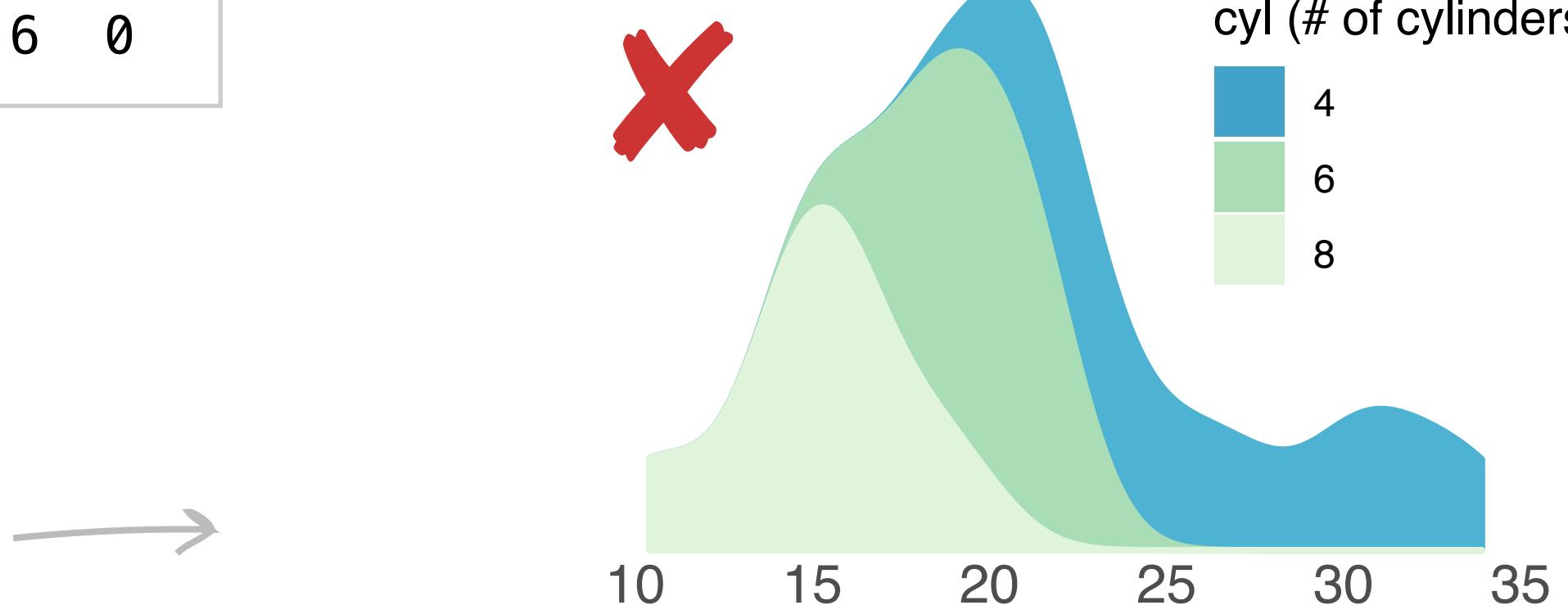
	mpg	cyl	am
Mazda RX4	21.0	6	1
Mazda RX4 Wag	21.0	6	1
Datsun 710	22.8	4	1
Hornet 4 Drive	21.4	6	0
Hornet Sportabout	18.7	8	0
Valiant	18.1	6	0

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

# 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")
```



## Experimental evaluation challenges

- Recruiting participants with appropriate level of expertise: if they are experts in an existing system and learning the new PGoG, we might underestimate the expressiveness or usability of PGoG. Some measures such as correctness of visualizations produced might be less sensitive to this problem.
- Recruiting enough participants: a power analysis will be conducted based on pilot studies. If for example  $N = 200$  is needed, recruitment might be difficult, assuming in-person experiments.



## Expert evaluation of the Grammar

- PGoG is new
  - PGoG is a formalism/abstraction
  - Using the PGoG prototype needs knowledge of ggplot
- > Expert evaluation is the appropriate first step

# How we collected probabilistic visualizations in the wild

## Methods of collection (informal)

- Google keyword searches: e.g., “icon array, faceted”
- Reputable news outlets
- “Best of” visualization awards/collections

## Criteria for inclusion

- Does it describe a probability distribution (counts also OK)
- Is it systematically generated (as opposed to hand-drawn or with inconsistent data -> visual mapping)

## How much does PGoG cover in the end?

- A lot, if judging by how much the PGoG aes/geoms apply
- A few, if “coverage” means strict replication of the original

- Viscosity: Resistance to Change
  - In layout-based grammar, need to rewrite layers, but for GoG, aesthetics mapping is easier to change
- Visibility: Ability to View Components Easily
  - Probabilities as first-class citizens!
- Premature Commitment: constraints on the Order of Doing Things
  - In PGoG, data, geoms, and aesthetics are quite independent
- Hidden dependencies
  - Probability rules that the user should know, like factoring a joint
- Role-expressiveness: The Purpose of an Entity is Readily Inferred
  - In addition to the probabilistic variables, PGoG has intuitively named aesthetics such as width and height, assuming that users understand the basic concepts of grammar of graphics

- Error-proneness: The Notation Invites Mistakes and the System Gives Little Protection
  - Foreseeable errors in probabilistic variables and implied mappings
- Abstraction: Types and Availability of Abstraction Mechanisms
  - None? PGoG implementation does not have a template-type functionality
- Secondary Notation: Extra Information in Means Other Than Formal Syntax
  - Yes, comments are supported in R implementation
- Closeness of Mapping: Closeness of Representation to Domain
  - $P(A | B)$  yay
- Consistency: Similar Semantics Are Expressed in Similar Syntactic Forms
  - Yes? Every visualization is specified with similar syntax
- Diffuseness: Verbosity of Language
  - Yes? Only template-based visualization specs are more succinct, but they are not flexible
- Hard Mental Operations: High Demand on Cognitive Resources
  - For users familiar with probabilities, PGoG removes mental efforts in data wrangling and lower level vis specs.
- Provisionality: Degree of Commitment to Actions or Marks
  - Similar to premature commitment
- Progressive Evaluation: Work-to-Date Can be Checked at Any Time
  - Difficult, because a visualization must have data, geoms, and aesthetics specified. R implementation does provide parsed probability “chain”, however.