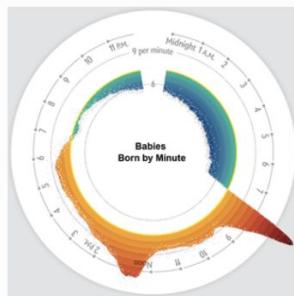


# Data Visualization for Scientific Discovery

[zan.armstrong@gmail.com](mailto:zan.armstrong@gmail.com)  
@zanstrong

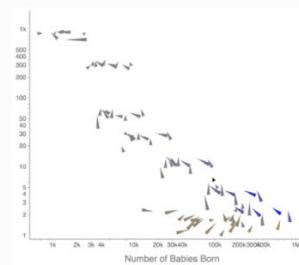
I'm a data visualization designer and engineer.



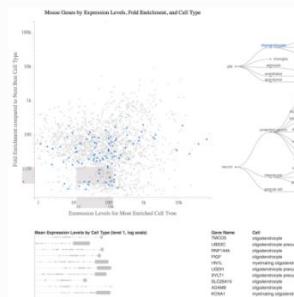
SCIENTIFIC AMERICAN: WHY ARE SO MANY BABIES BORN AROUND 8AM?



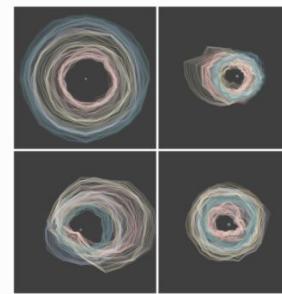
WITH THE TEAM AT STAMEN



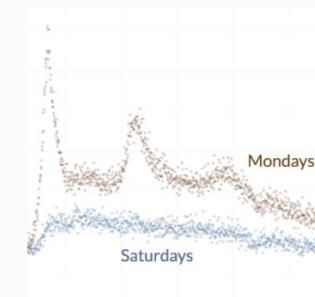
RESEARCH: VISUALIZING STATISTICAL MIX EFFECTS AND SIMPSON'S PARADOX



GENETIC EXPRESSION TOOL



WEATHER CIRCLES



OPENVISCONF: EVERYTHING IS SEASONAL



WHICH IS BIGGER?



SCHOOL FOR POETIC COMPUTATION

format specifier	resulting formatted num
<code>d3.format("")</code>	2509.31
<code>d3.format("s")</code>	2.50931k
<code>d3.format(",%")</code>	250,931%
<code>d3.format("+,%")</code>	+250,931%
<code>d3.format(".1%")</code>	250,931.0%
<code>d3.format("4r")</code>	2509
<code>d3.format(".4r")</code>	2509.3100
<code>d3.format(".4n")</code>	2,509
<code>d3.format(".3n")</code>	2,51e+3

TINY TOOLS

I'm a [data visualization designer](#) and engineer.

I am part of **Google's Accelerated Science** team.

My goal is to use data visualization to help a team of scientists, statisticians/data scientists, and machine learning engineers analyze their data and make discoveries.

# All you need to know about Data Viz

Pie Charts are bad.  
Non-zero bar charts are bad.  
Dual y-axis are bad.  
Rainbows are bad.  
Chart-junk is bad.  
A high ink-to-data ratio is bad.

The End.

Any questions?

## Oft-quoted "rules" of data viz.

- Pie Charts are bad.
- Non-zero bar charts are bad.
- Dual y-axis are bad.
- Rainbows are bad.
- Chart-junk is bad.
- A high ink-to-data ratio is bad.



is SOY  
**BAD**  
for YOU?



ARE  
PRESERVATIVES  
BAD FOR  
YOU?

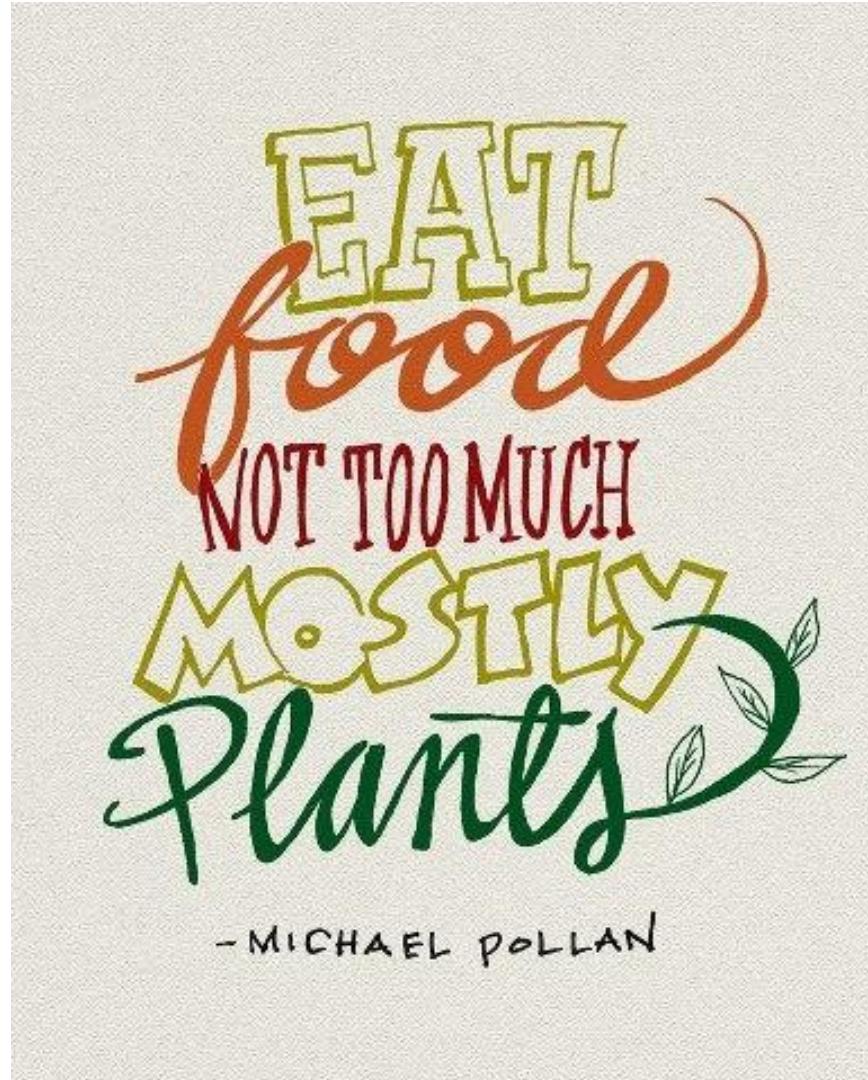


How Food Coloring Is Made  
& Is It Bad For You?



Better to have a few good  
guidelines than a bunch of  
rules telling you what not  
to do

Better to have a few  
good guidelines



Better to have a few  
good guidelines

Create small multiples (many, related  
charts)

Use color intentionally

Optimize your charts for your  
time/energy/attention

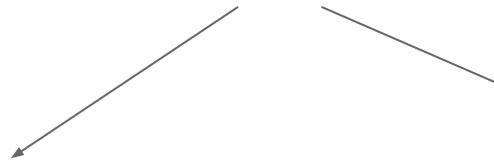
Is it good or bad depends  
on your goals, context,  
and constraints.

The answer to "is it good" depends on goals, context, and constraints.

Is a candy bar good for me?

The answer to "is it good" depends on goals, context, and constraints.

Is a candy bar good for me?



Am I at home drinking tea?



Am I on a week-long backpacking trip?



no :(

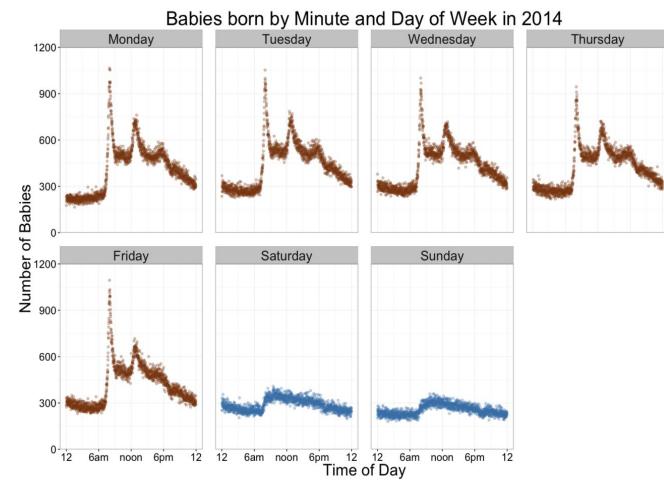
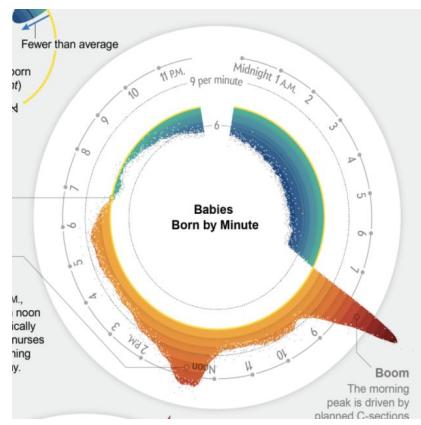


yes!



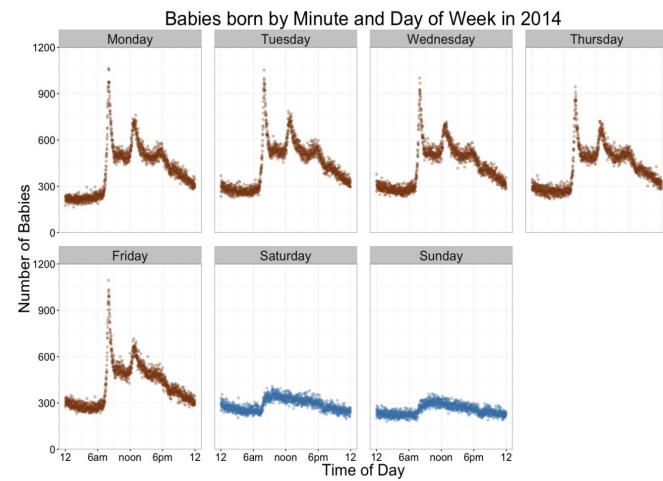
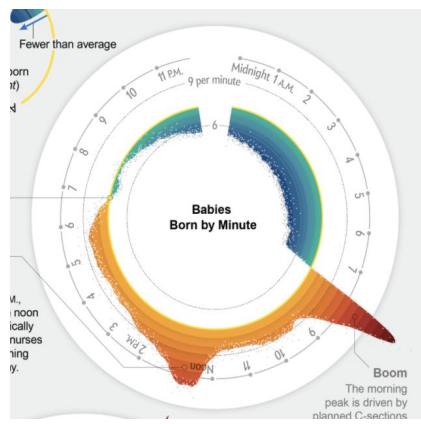
The answer to "is it good" depends on goals, context, and constraints.

Which visualization is better?



The answer to "is it good" depends on goals, context, and constraints.

Which visualization is better?



It depends!

What's the point?

How much time to create?

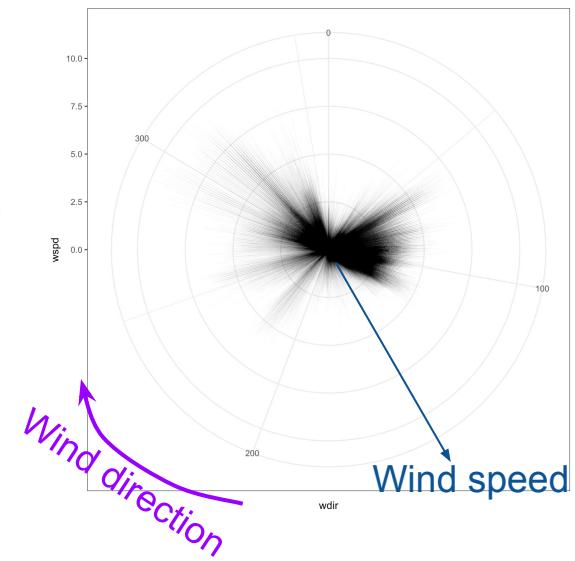
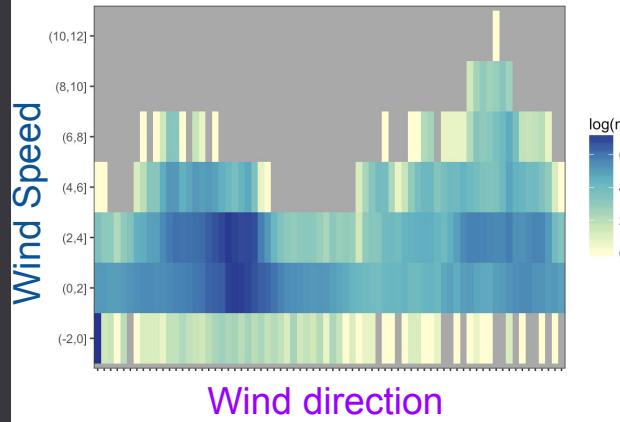
Audience? Who? Their goal?

Do I care about day of week?

Do I care about total values or difference to the mean?

The answer to "is it good" depends on goals, context, and constraints.

Which visualization is better?



For communication (publication, journalism, etc)

The answer to "is it good" depends on goals, context, and constraints.

For scientific analysis/discovery (exploratory & "pipelines")

The answer to "is it good" depends on goals, context, and constraints.

For communication (publication, journalism, etc)

Does it attract the audience's attention?

Get the message across?

It is true to the data? (not misleading)

How much space on the page or screen?

Is it immediately understandable?

Can everybody interpret it?

Most viz "rules"  
optimized for  
communication

For scientific analysis/discovery (exploratory & "pipelines")

The answer to "is it good" depends on goals, context, and constraints.

For communication (publication, journalism, etc)

**Does it attract the audience's attention?**

**Get the message across?**

**It is true to the data? (not misleading)**

How much space on the page or screen?

Is it immediately understandable?

Can everybody interpret it?

Most viz "rules"  
optimized for  
communication

For scientific analysis/discovery (exploratory & "pipelines")

**How much of your time/energy is needed to interpret?**

**If something in your data is important, will you see it in the viz?**

How long does it take to create?

How much mental effort does it take to create?

How hard is it to ask a new question?

Will I look at this type of chart/tool often? Aka - is a learning curve for interpretation ok?

Goals, context, and  
constraints for  
analysis and  
discovery

If something in  
your data is  
important, will  
you notice it?

*It's ok to risk something looking important that  
isn't, because you're going to do more analysis.*

Goals, context, and  
constraints for  
analysis and  
discovery

**Your time &  
energy is your  
most limited  
resource.**

*Number of charts is (generally) not a limiting  
resource.*

Is it good or bad depends  
on your goals, context,  
and constraints

Is it good or bad depends  
on your goals, context,  
and constraints\*

\*for food and viz... and engineering, ML, and stats

# Create small multiples

What are small  
multiples?

Lots of small charts.

What are small  
multiples?

What are small  
multiples?

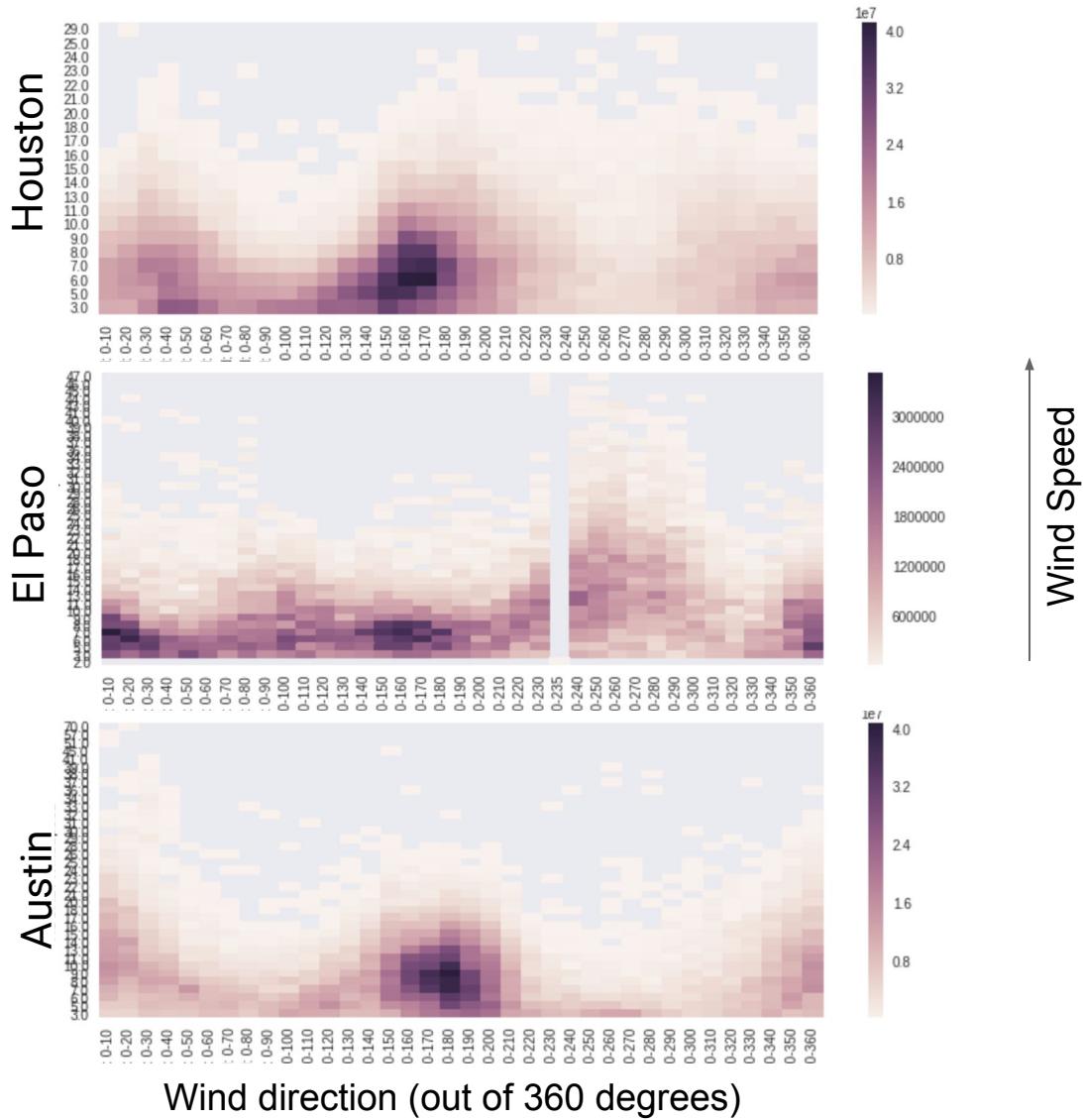
Lots of small charts.  
Similar to each other.

What are small  
multiples?

Lots of small charts.  
Similar to each other.  
That you look at at the same  
time.

# What are small multiples?

Lots of small charts.  
Similar to each other.  
That you look at at the same time.



# Why small multiples?

**You can't compare things you can't see.**

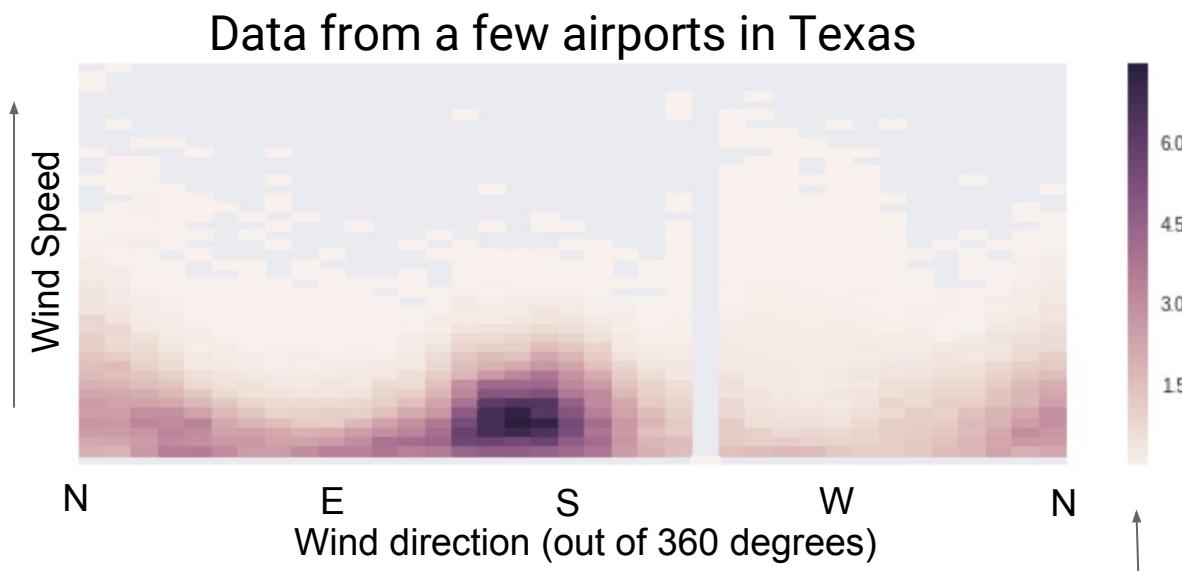
# Create small multiples

Same form, related data

# Small Multiples

Same form,  
Related Data

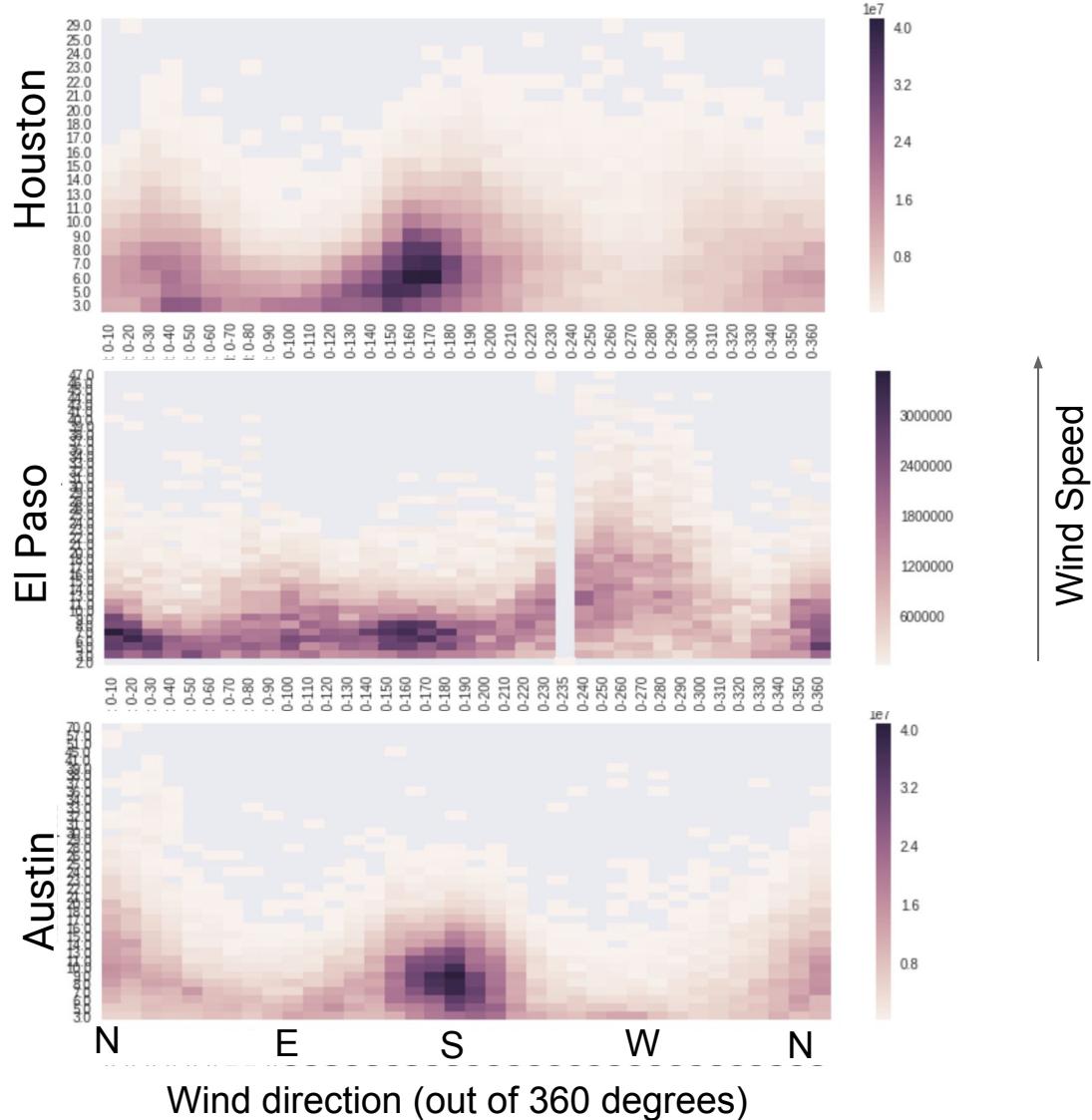
Let's say we're interested in the relationship between wind direction and wind speed in Texas.



# Small Multiples

Same form,  
Related Data

Small multiples provide **context** and **contrast**.



# Create small multiples

Same data, different form

## Small Multiples

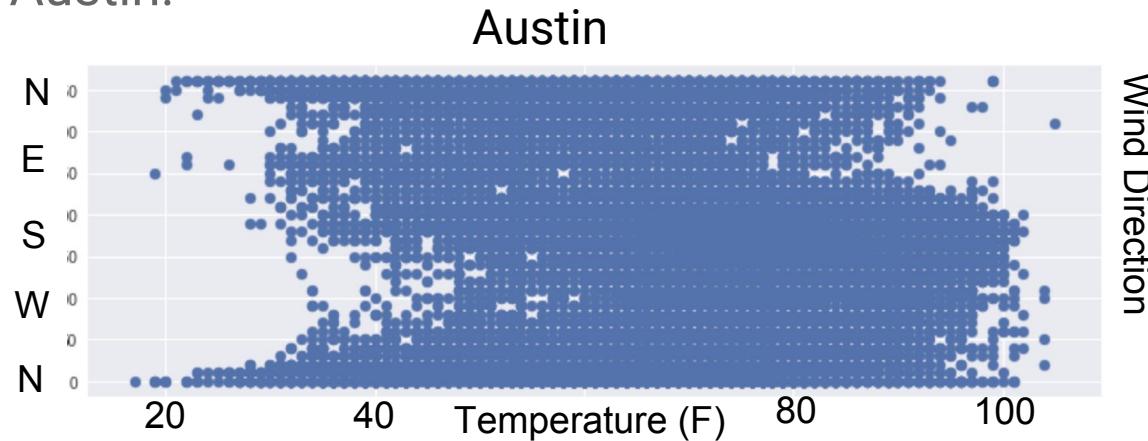
Same data,  
Different form

Let's say we're interested in the relationship between wind direction and temperature in Austin.

# Small Multiples

Same data,  
Different form

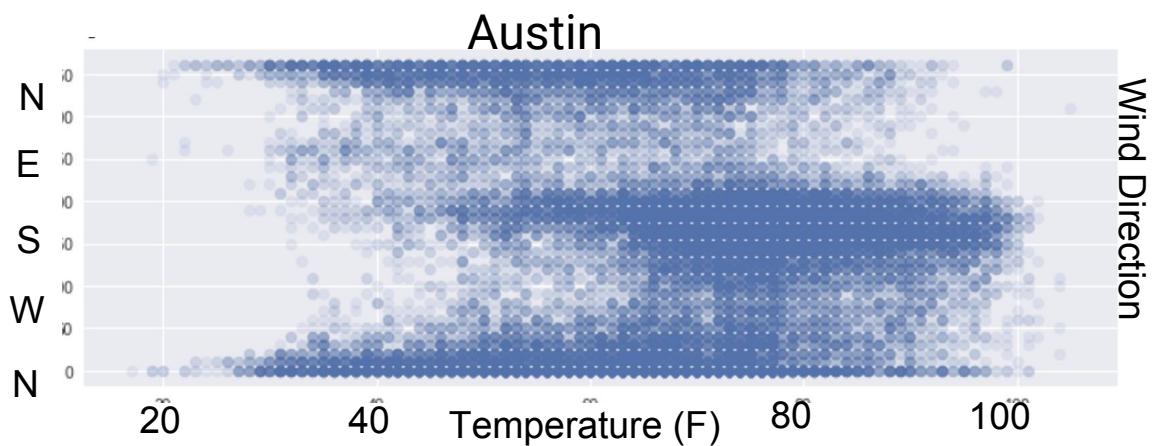
Let's say we're interested in the relationship between wind direction and temperature in Austin.



Common trick: Lower opacity to see density

## Small Multiples

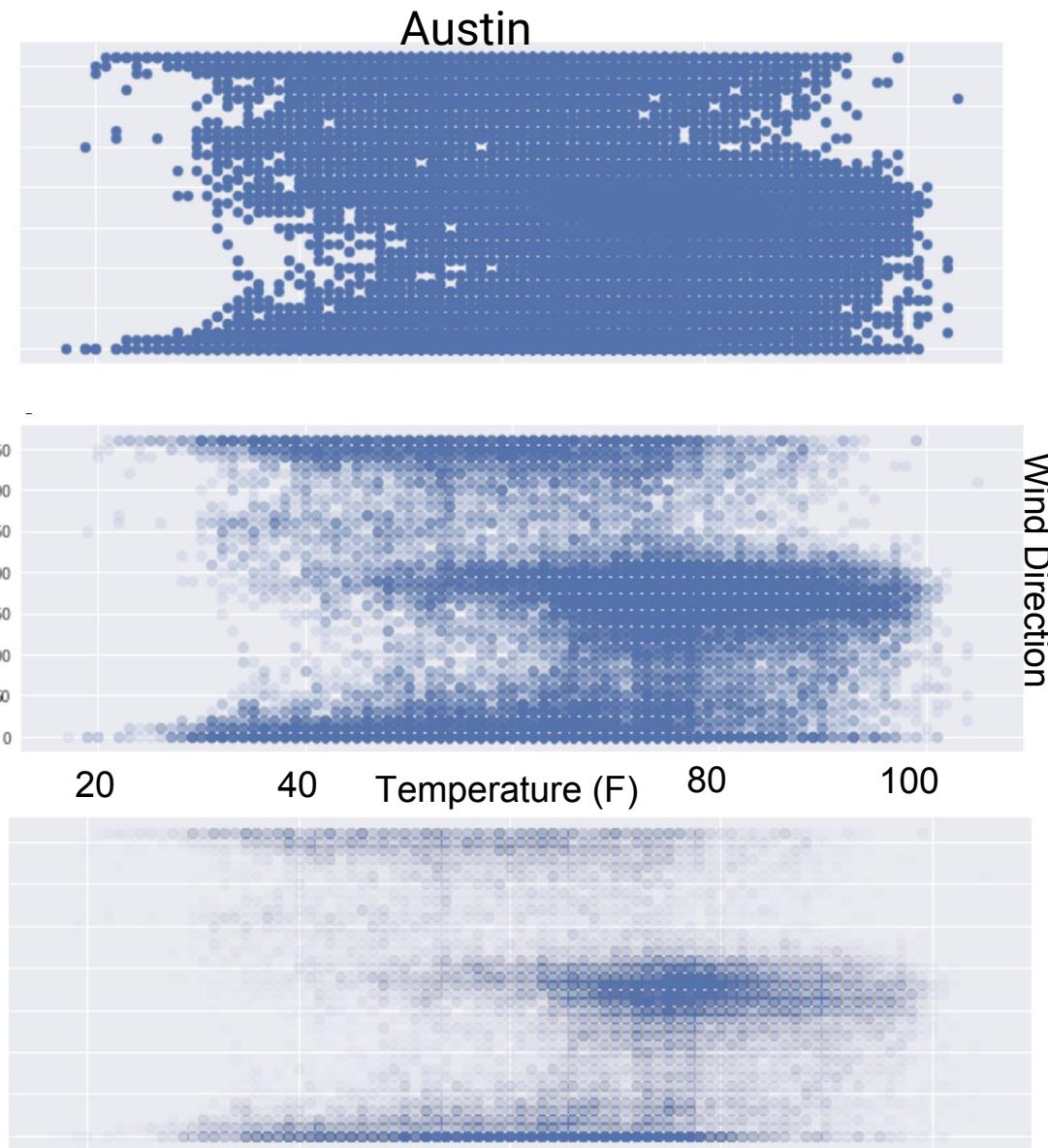
Same data,  
Different form



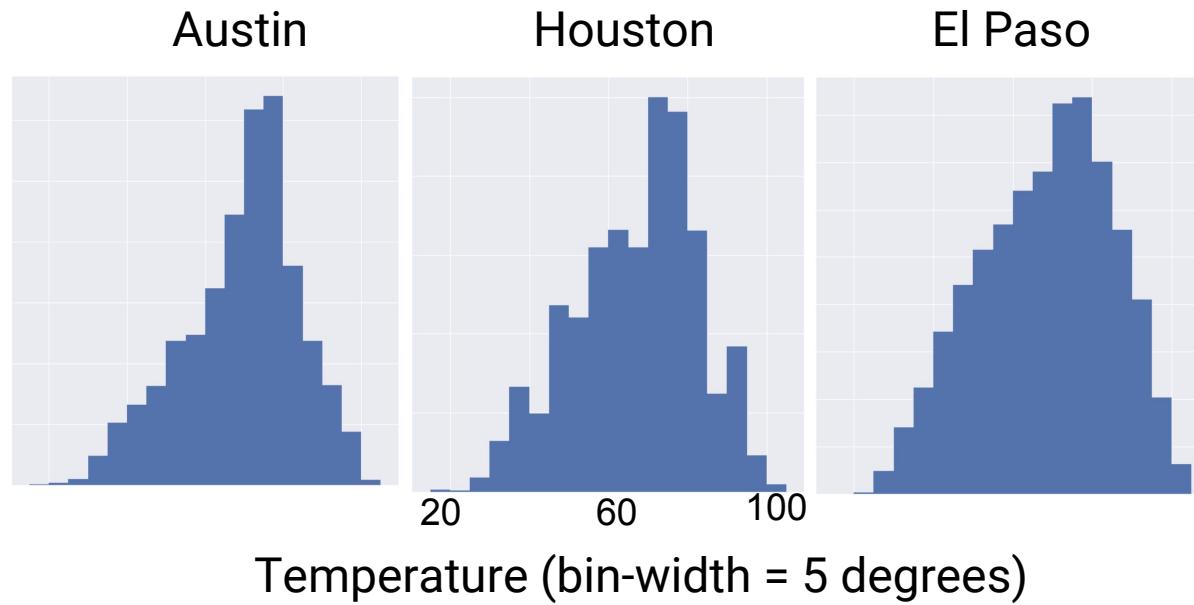
# Small Multiples

Same data,  
Different form

Let's look at 3 versions at once: see extremes,  
and highest density of data.

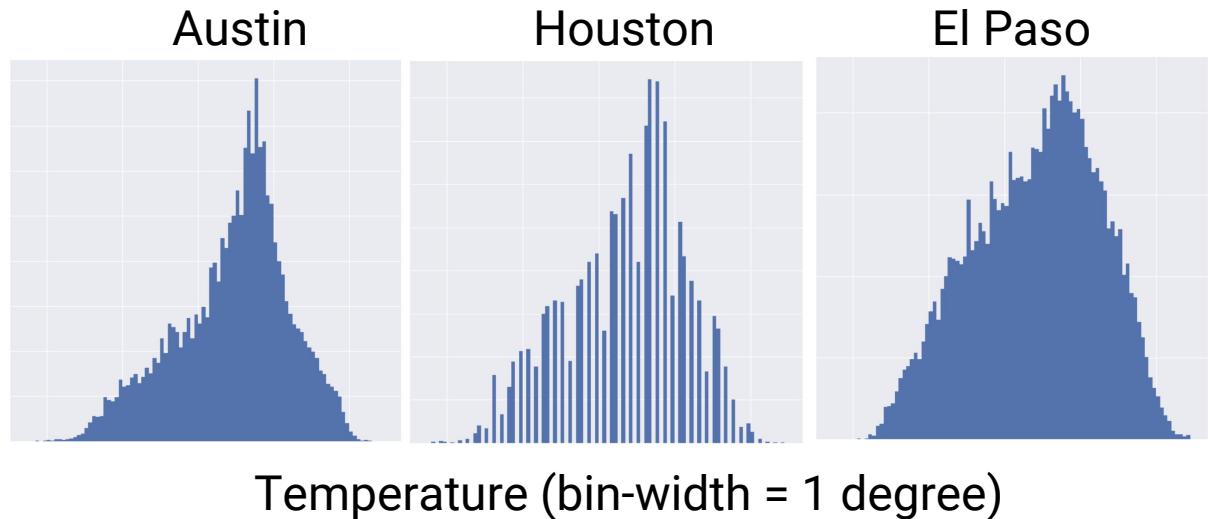
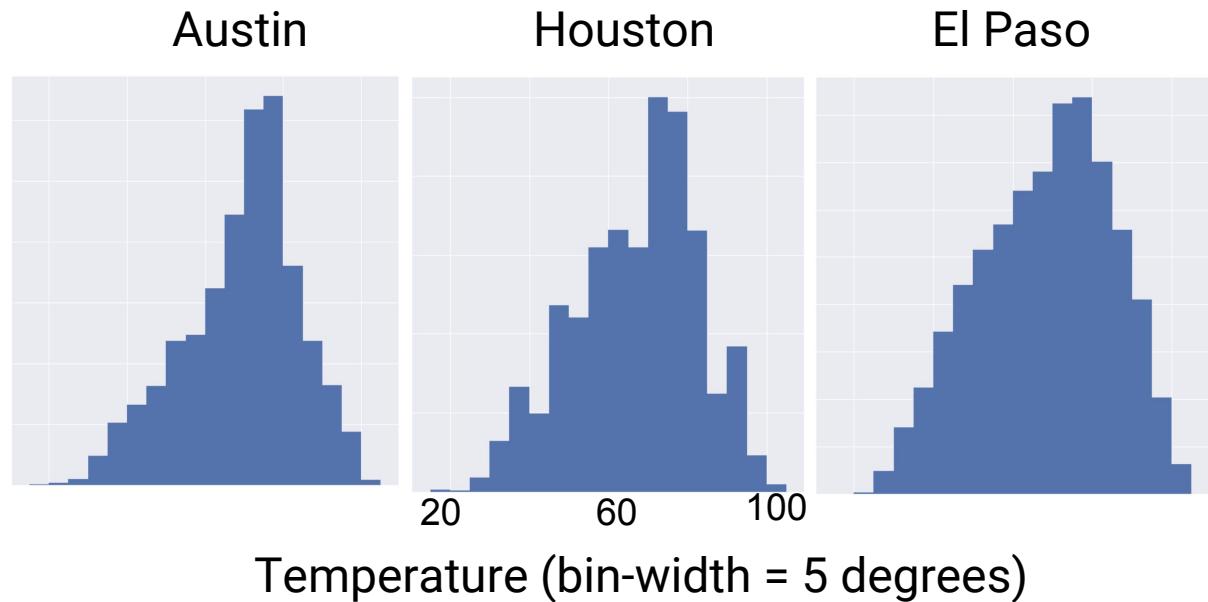


## Small Multiples



For more, see Amelia and Arun's exploration of [binwidth and histograms](#)

## Small Multiples



For more, see Amelia and Arun's exploration of [binwidth and histograms](#)

# Small Multiples Resources

Matplotlib: [Subplots Example](#) (most flexible)

Seaborn: [FacetGrid](#)

Altair: [Trellis Example](#)

# Use color intentionally

*The basic visual attributes ... are perceived without any conscious effort.*

*[This] is called "preattentive" processing...*

*Preattentive perception is done in parallel, but attentive processing is done serially and is, therefore, much slower.*

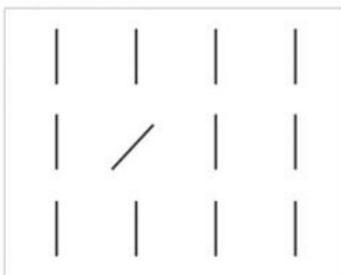
*The basic visual attributes ... are perceived without any conscious **effort**.*

*[This] is called "preattentive" processing...*

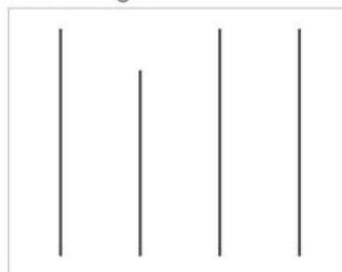
*Preattentive perception is done in parallel, but attentive processing is done serially and is, therefore, much slower.*

## Form

### Orientation



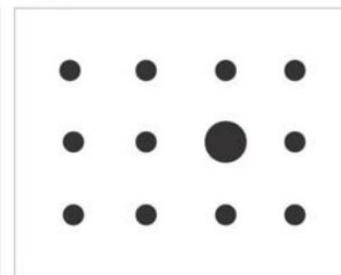
### Line Length



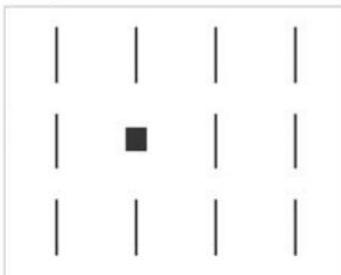
### Line Width



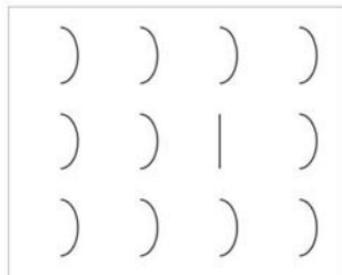
### Size



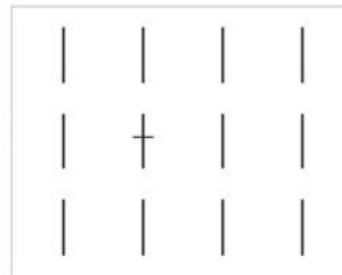
### Shape



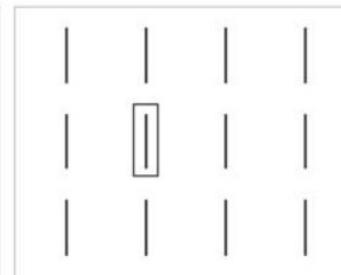
### Curvature



### Added Marks

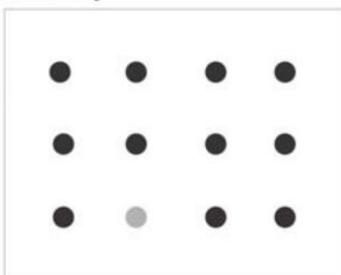


### Enclosure

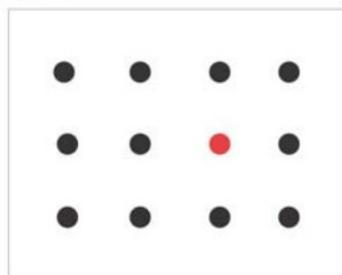


## Color

### Intensity



### Hue



## Spatial Position

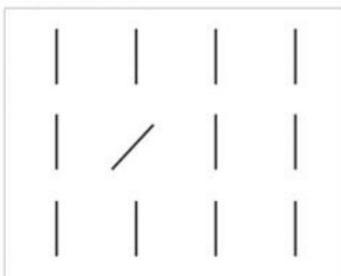
### 2-D Position



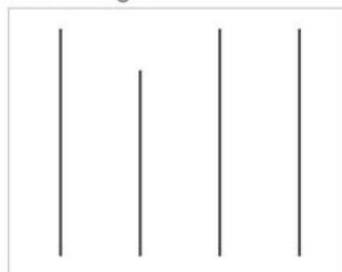
FIGURE 5: Preattentive attributes of visual perception most applicable to data presentation.

## Form

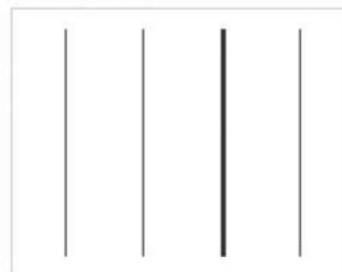
### Orientation



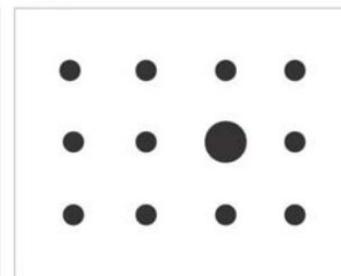
### Line Length



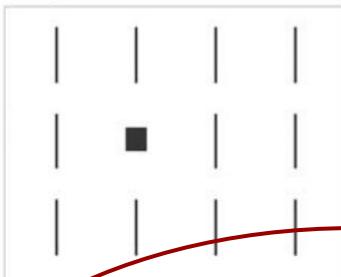
### Line Width



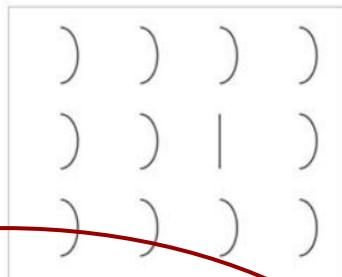
### Size



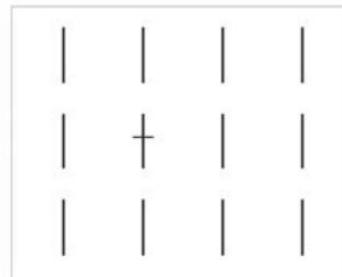
### Shape



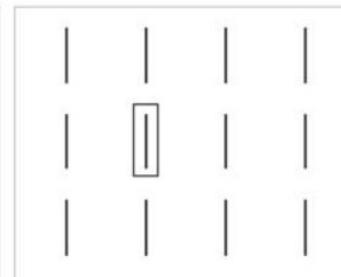
### Curvature



### Added Marks

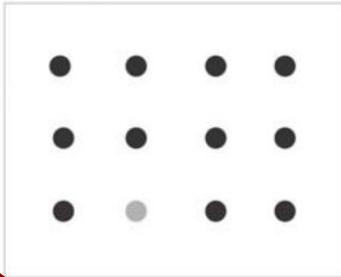


### Enclosure

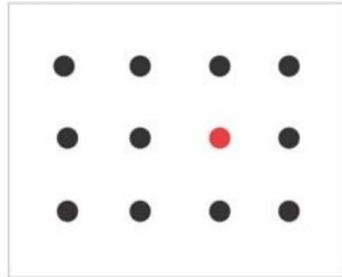


## Color

### Intensity



### Hue



## Spatial Position

### 2-D Position

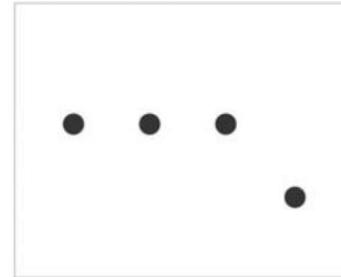


FIGURE 5: Preattentive attributes of visual perception most applicable to data presentation.

Any interesting patterns in this set of numbers\*?

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

\* Example inspired by Cole's [Storytelling with Data](#) book/blog and Colin Ware's books.

How many 7's are there?

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

Let's try the typical default: add categorical colors

Let's try the typical default: add categorical colors

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

0

0->9 is ascending. Let's look at a sequential scale.



A 4x8 grid of numbers ranging from 0 to 9. The numbers are arranged in four rows and eight columns. The colors of the numbers transition from light blue for lower values to dark blue for higher values, illustrating a sequential color scale. The sequence starts at 0 and ends at 9, with each row showing a different segment of the sequence.

6	8	4	7	2	6	6	9
2	6	6	9	7	4	3	8
5	6	9	9	7	6	3	3
6	8	3	4	6	7	2	6
3	9	5	6	2	5	9	6
2	9	6	2	9	3	5	8
4	7	5	2	3	6	5	3
3	4	6	7	5	2	9	3
4	5	3	4	5	3	4	2
5	3	4	5	3	4	2	6
4	4	3	5	2	6	3	5
3	5	3	4	3	5	3	4
4	4	5	4	4	5	4	5
5	2	6	3	5	3	4	4
3	5	3	4	3	5	2	6
4	2	7	6	9	8	2	5
5	9	2	6	7	4	9	2
6	1	5	7	6	6	9	5
7	2	1	4	7	2	6	1
8	3	2	5	8	3	4	2
9	4	6	9	7	4	3	8

9

0

9

Viridis is also a popular colorscheme

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261



Interested in extremes? Diverging from 4.5?

0

9



But, we know we want to count the number 7.

So, let's just show the 7s.

604 **7** 26692669 **7** 438  
5699 **7** 630683406 **7** 2  
3956259629358 **7** 24  
06 **7** 5236503962934  
5345345342645342  
5344352635344526  
35343502 **7** 6982562  
35 **7** 6695926 **7** 49261

So, let's just show the 7s.

604 726692669 7438  
5699 7630683406 72  
3956259629358 724  
0675236503962934  
5345345342645342  
5344352635344526  
35343502 76982562  
3576695926 749261

How many 3's are there?

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

There are lots :)

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

What about 1's?

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

It is the loneliest number.

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

# We can also use highlighting with "small multiples"

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

Or, for more general analysis, look at different color schemes that emphasize different aspects of the data.

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

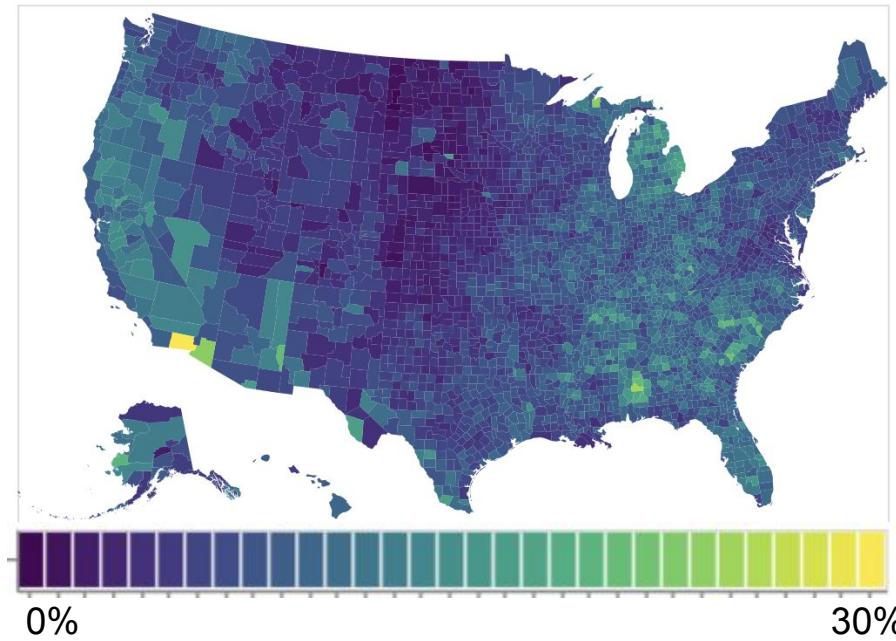
6047266926697438  
5699763068340672  
3956259629358724  
0675236503962934  
5345345342645342  
5344352635344526  
3534350276982562  
3576695926749261

# Use color intentionally

Make the mapping from numeric domain to color range be meaningful

## Unemployment Data by County in the US, ~2009

Use color intentionally:  
Mapping from numbers  
to colors



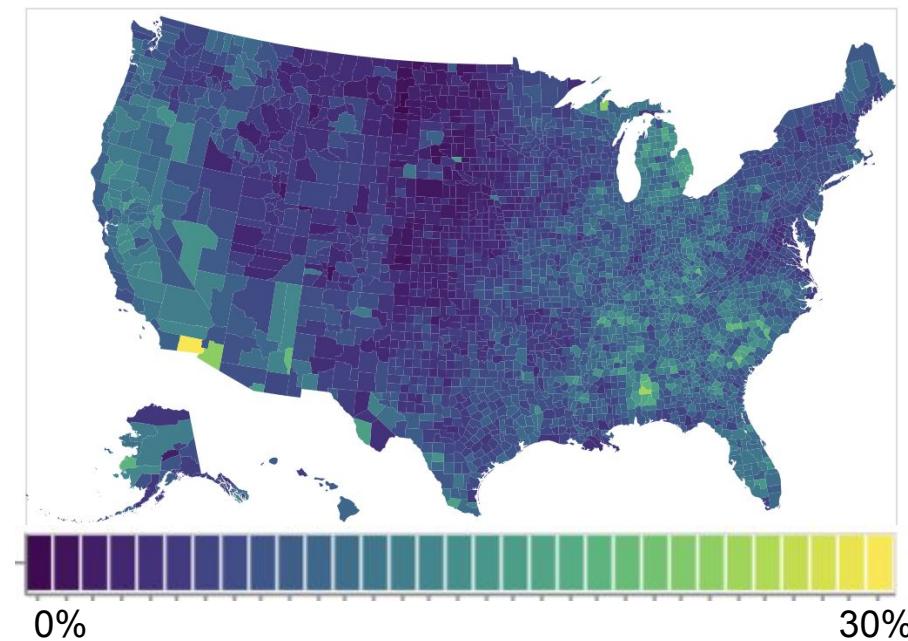
[Altair Examples](#)  
[Colab w/ this set of maps](#)

Use color intentionally:  
Mapping from numbers  
to colors

Start with the default color mapping (Altair)

**Typical Default Color Mapping**

Color range spans from 0 to max (30.1%).



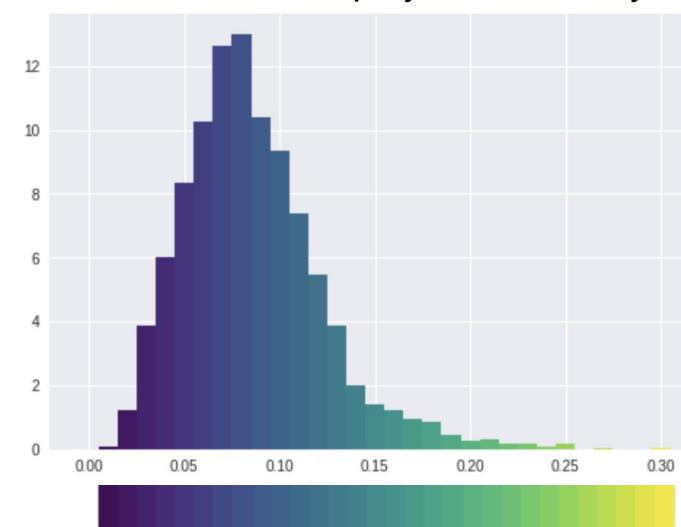
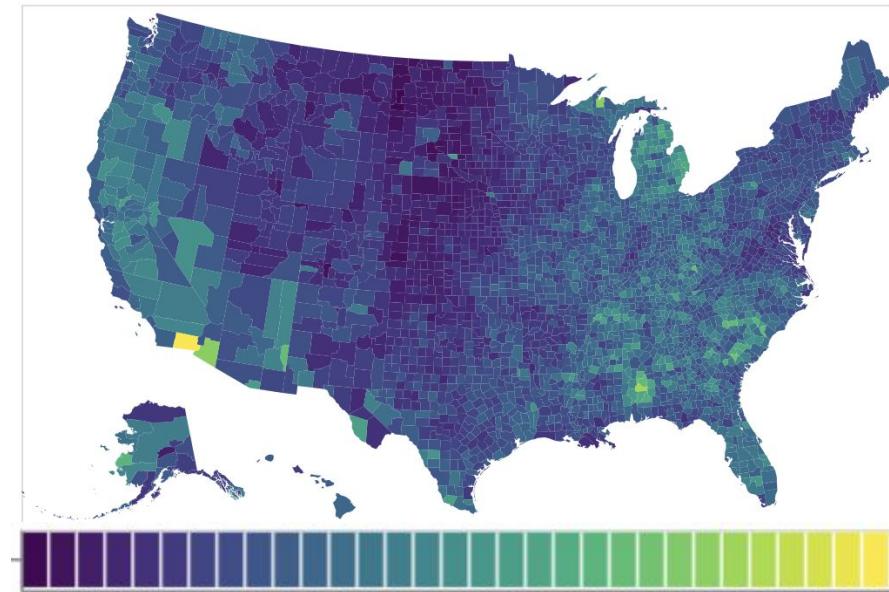
[Altair Examples](#)  
[Colab w/ this set of maps](#)

Use color intentionally:  
Mapping from numbers  
to colors

But, an extreme outlier overwhelms the map

**Typical Default  
Color Mapping**

Color range  
spans from  
from 0 to max  
(30.1%).

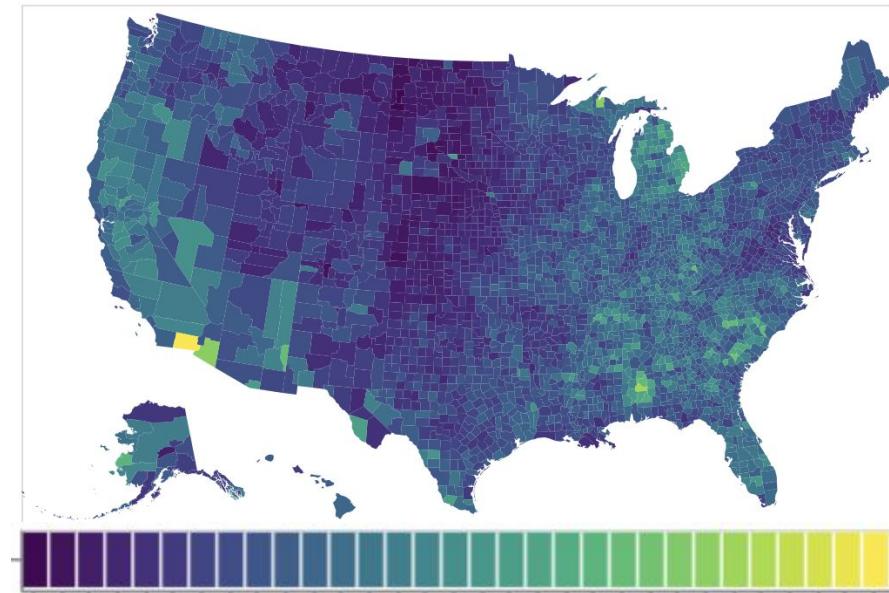


Use color intentionally:  
Mapping from numbers  
to colors

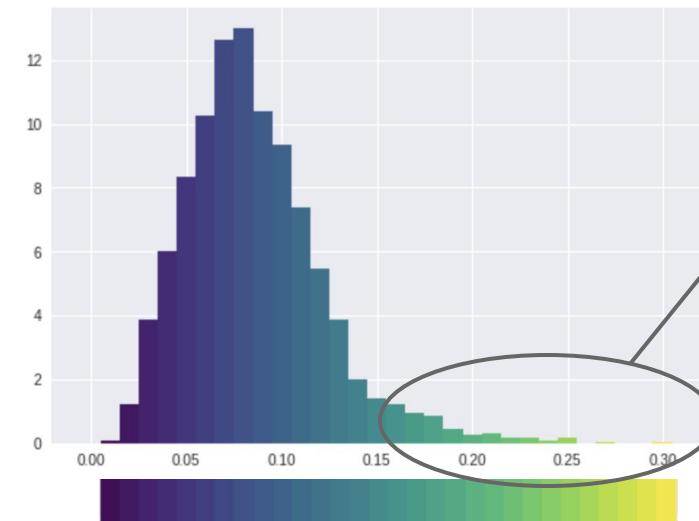
Harder to differentiate differences for bulk of the data

**Typical Default Color Mapping**

Color range spans from 0 to max (30.1%).



Distribution of unemployment rates by county in the dataset



Only 4.5% of the counties have  $>15\%$  unemployment, but we're using half our color scale on 15%-30%.

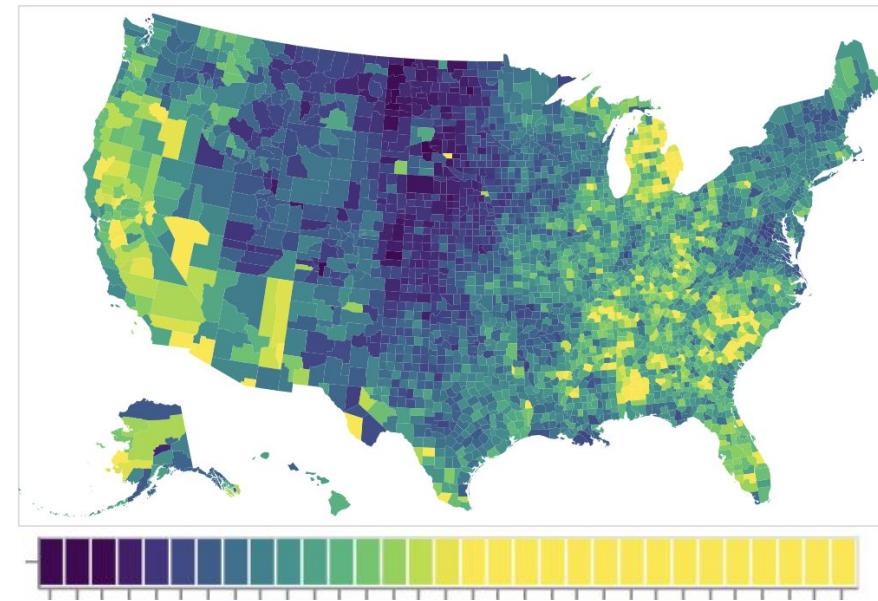
[Altair Examples](#)  
[Colab w/ this set of maps](#)

# Use color intentionally: Mapping from numbers to colors

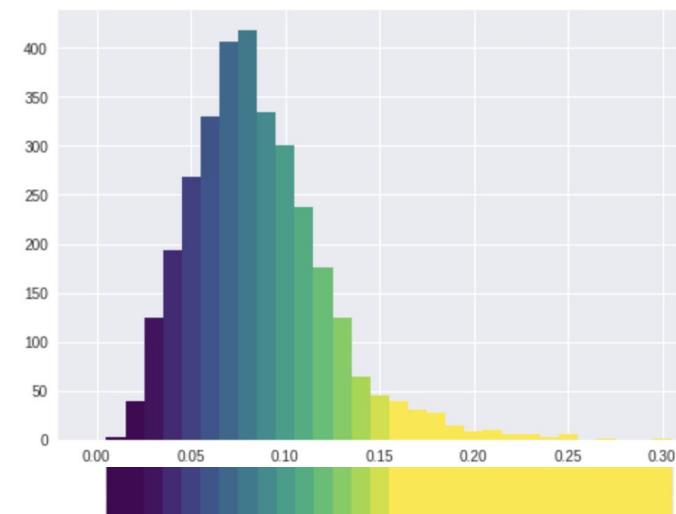
Use distribution of the data, rather than max value, to set color min/max

## Intentional Range, based on the data distribution

Color range is set to min of 2% and max of 15.5%.



Distribution of unemployment rates by county in the dataset



More than doubling the perceptual differentiation, with little loss on the extremes.

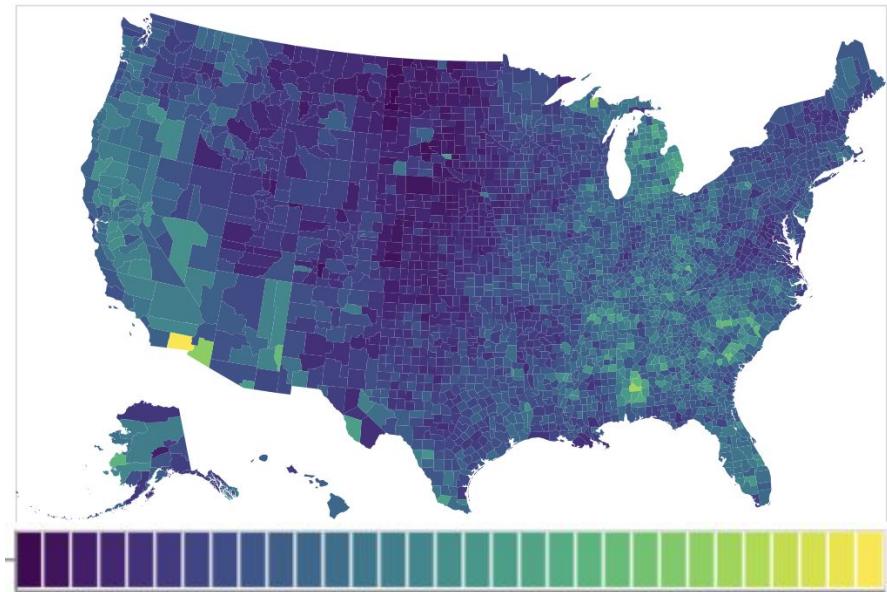
[Altair Examples](#)  
[Colab w/ this set of maps](#)

Use color intentionally:  
Mapping from numbers  
to colors

Let's compare.

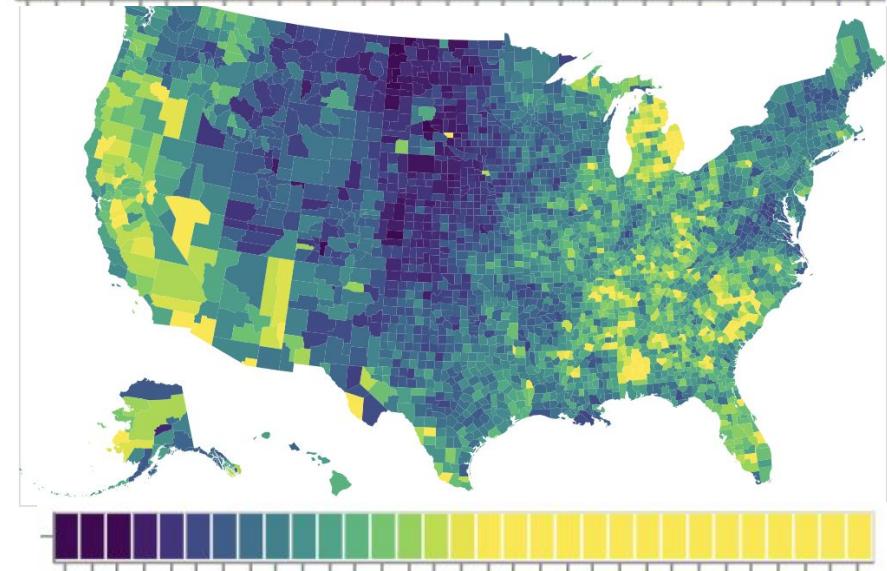
**Typical Default  
Color Mapping**

Color range  
spans from  
min to  
max (27%).



**Intentional  
Range, based  
on data  
distribution**

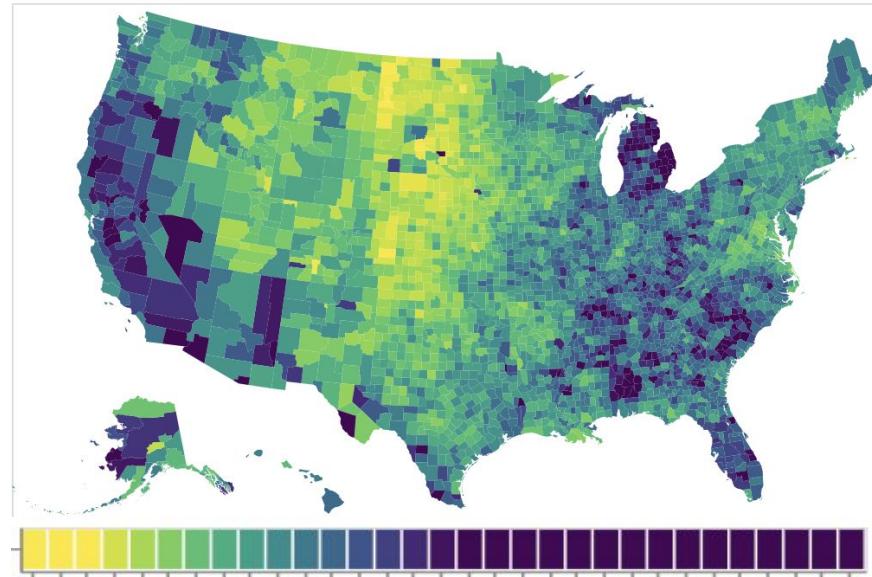
Color range is  
set to min of  
2% and max of  
15.5%.



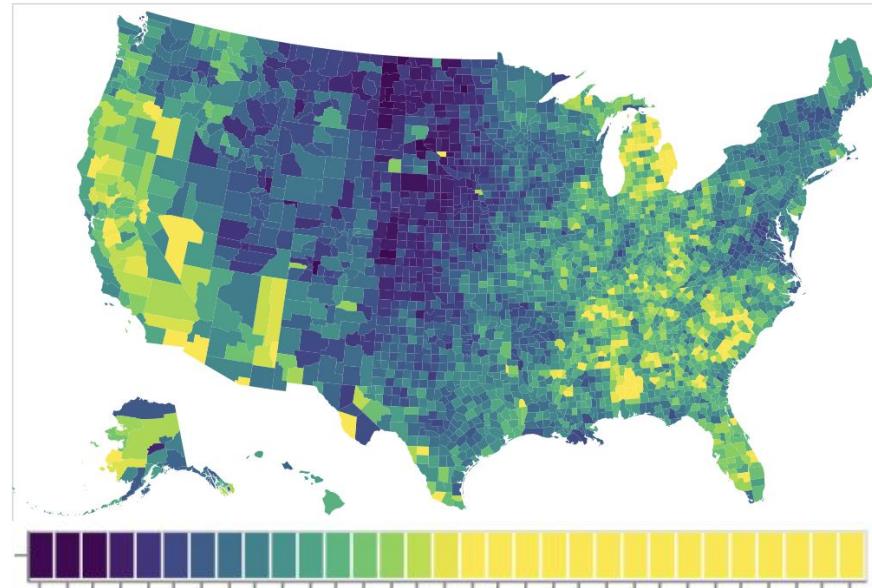
Use color intentionally:  
Mapping from numbers  
to colors

What if we flip the color scheme?

**Intentional Range, flip the color map**  
Lowest unemployment rates are yellow, highest rates are purple.



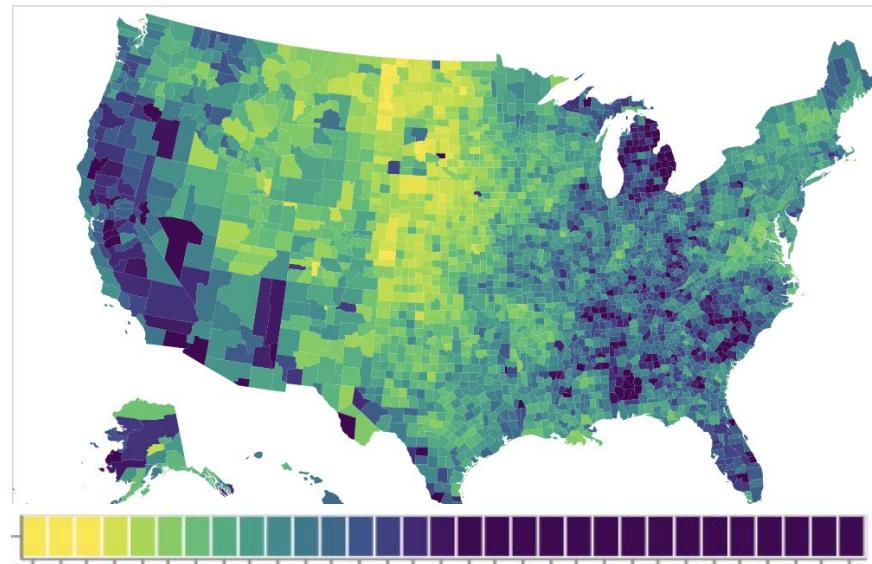
**Intentional Range**  
Highest unemployment rates are yellow, lowest are purple.



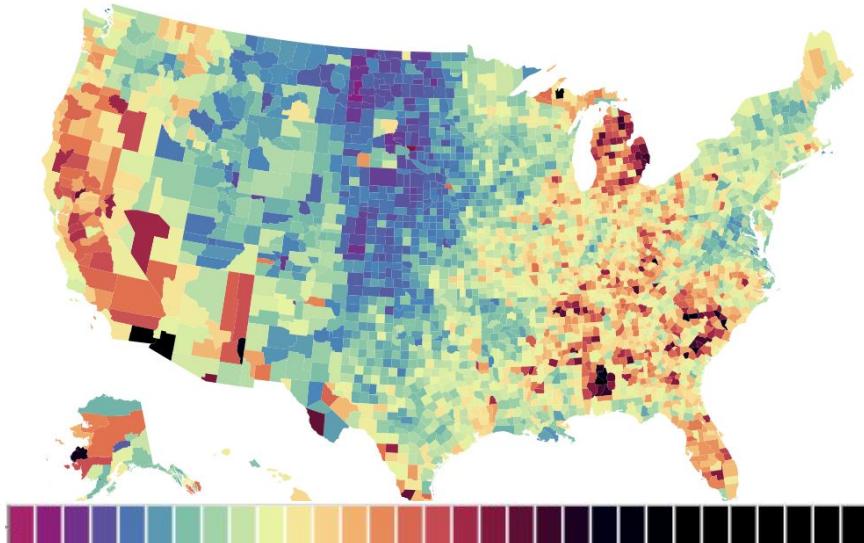
Use color intentionally:  
Mapping from numbers  
to colors

Or try a different color scheme?

**Intentional Range, flip the color map**  
Lowest unemployment rates are yellow, highest rates are purple.



**Another colorscheme?**  
Do we notice the same things?  
Different things?

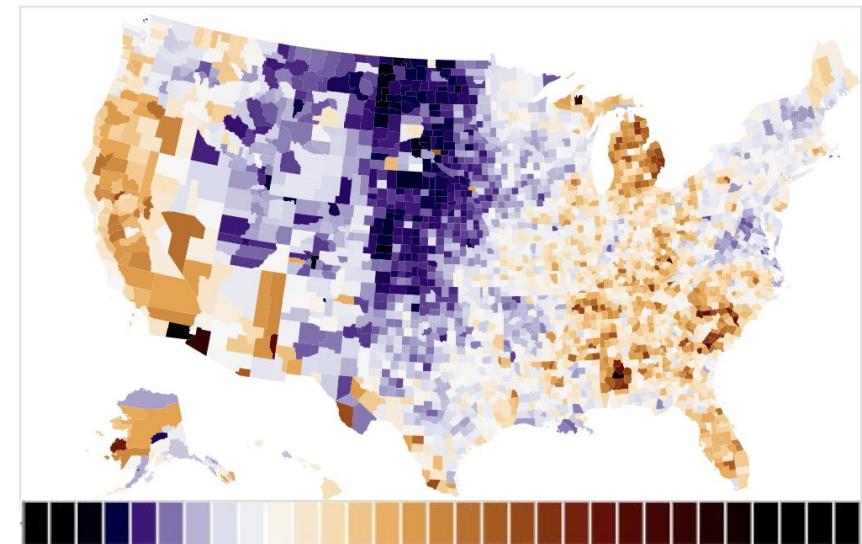


Use color intentionally:  
Mapping from numbers  
to colors

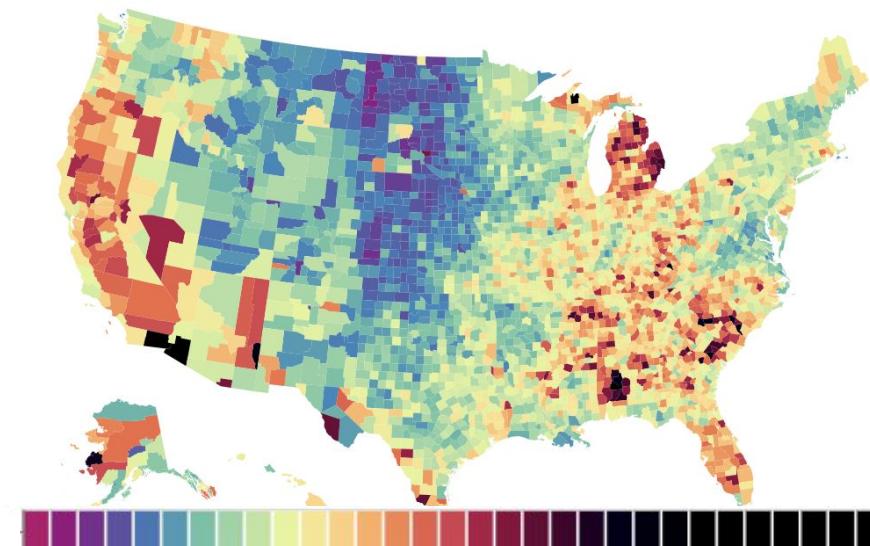
Perhaps diverging is more appropriate?  
Compare to the median.

**Diverging**

What if we intentionally compared to the median value? Purple is below 8.5%, median is shown as white, above the median as orange.



**Another colorscheme?**  
Do we notice the same things?  
Different things?

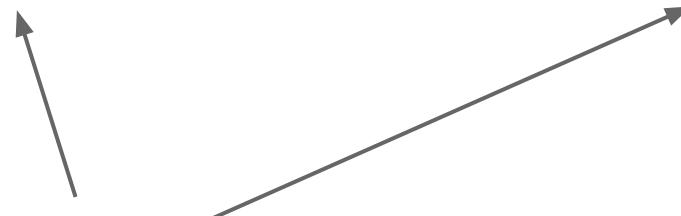
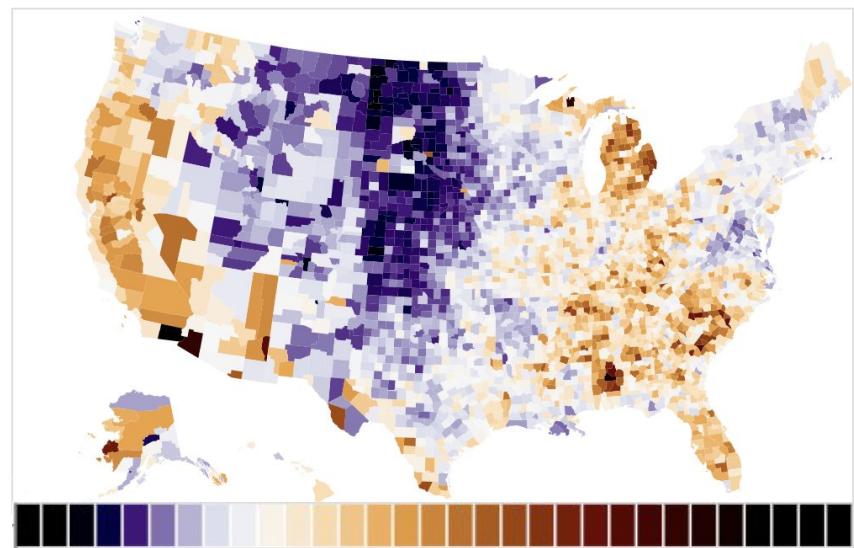


Use color intentionally:  
Mapping from numbers  
to colors

But, is the 8.5% median unemployment rate the most relevant "middle"? What about the "ideal" rate of 5.5% for a healthy economy?

**Diverging**

What if we intentionally compared to the median value (8.5%)?



*Very dark extreme ends start to look similar.*

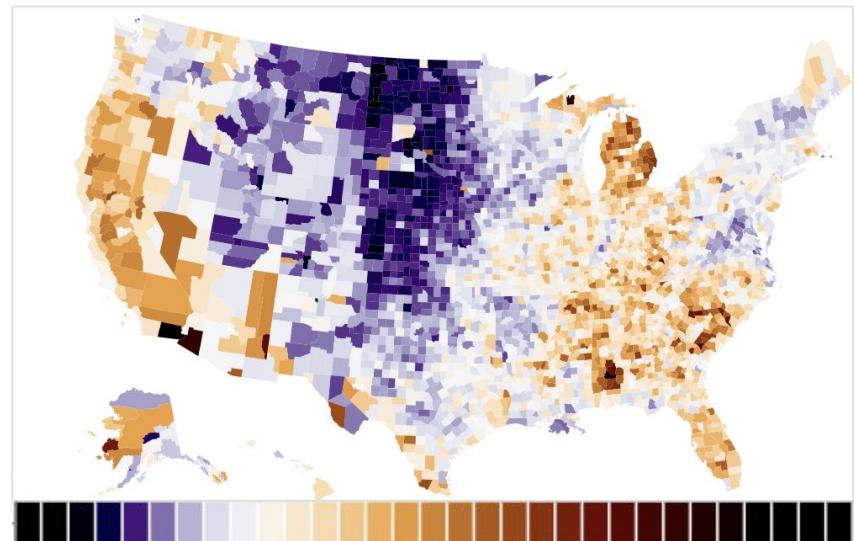
*We could "clamp" to a mid-dark max value to prevent this.*

Use color intentionally:  
Mapping from numbers  
to colors

But, is the 8.5% median unemployment rate the most relevant "middle"? What about the "ideal" rate of 5.5% for a healthy economy?

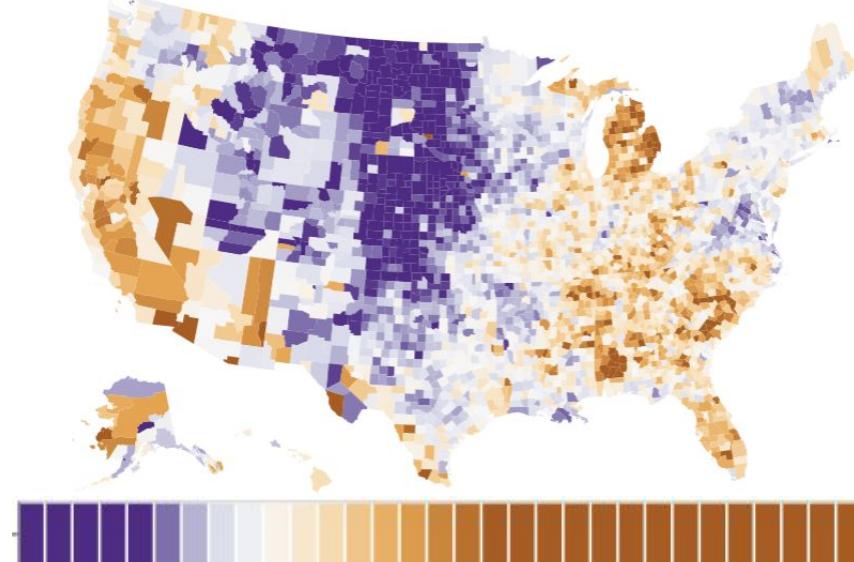
**Diverging**

What if we  
intentionally  
compared to  
the median  
value (8.5%)?



**Diverging**

From median,  
Clamped  
min/max  
values.

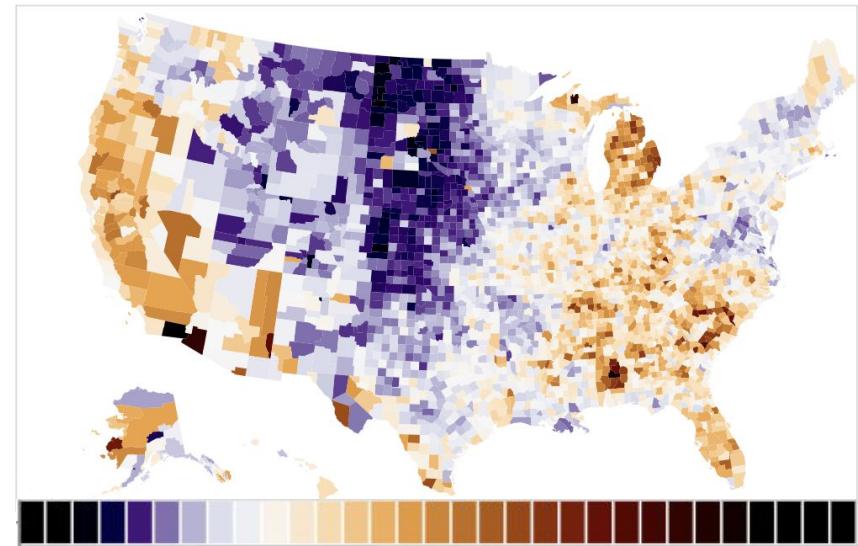


Use color intentionally:  
Mapping from numbers  
to colors

What should we set as the middle value?

**Diverging**

What if we  
intentionally  
compared to  
the median  
value (8.5%)?



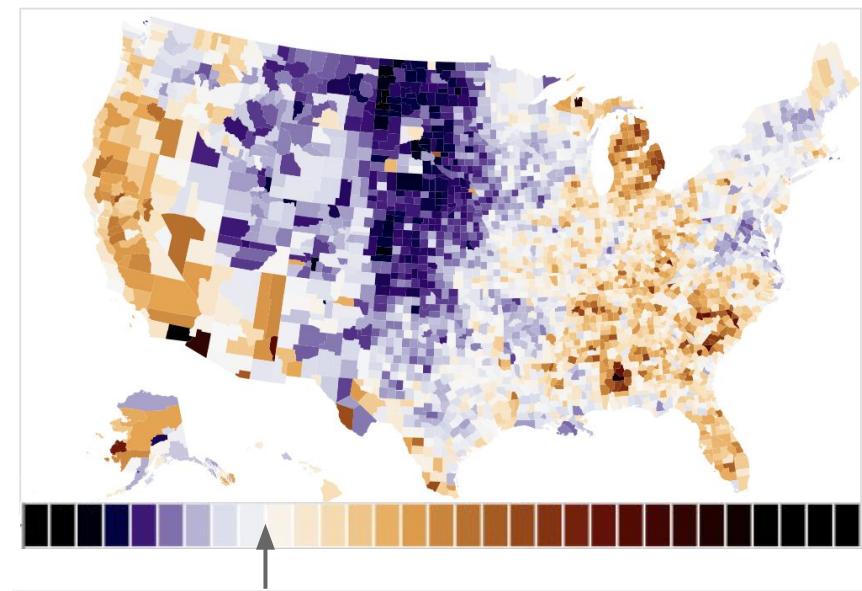
*Why should the "middle" be 8.5%?*

Use color intentionally:  
Mapping from numbers  
to colors

But, is the 8.5% median unemployment rate the most relevant "middle"? What about the "ideal" rate of 5.5% for a healthy economy?

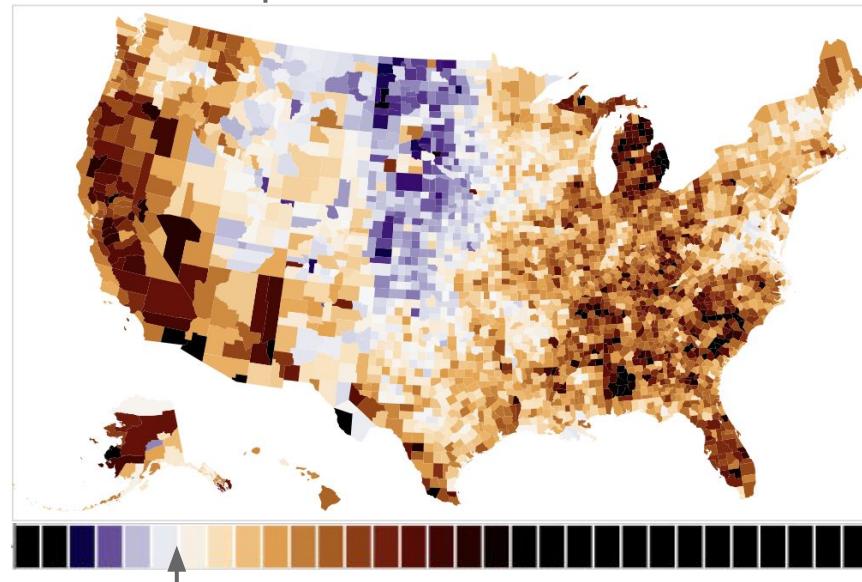
**Diverging**

What if we intentionally compared to the median value (8.5%)?



**Diverging:**

compare to ideal value according to economists (~5.5%) rather than median

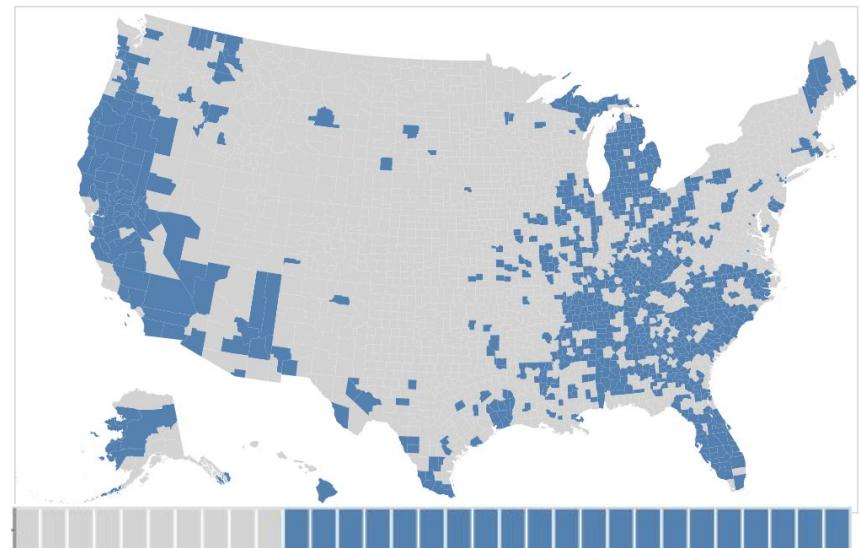


Use color intentionally:  
Mapping from numbers  
to colors

We could put in a hard cut-off, say to highlight areas with greater than 10% unemployment.

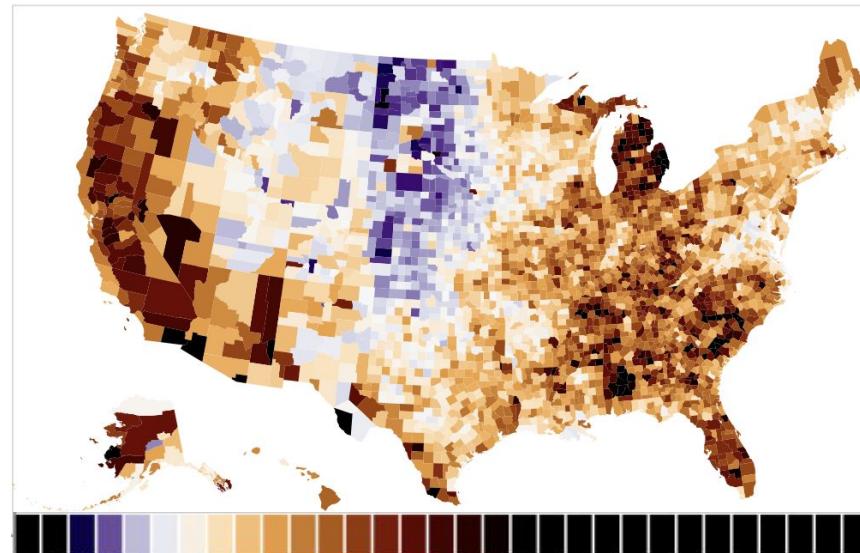
**Highlight**

Just focus on areas with >10% unemployment (in blue).



**Diverging:**

compare to ideal value according to economists (~5.5%) rather than median

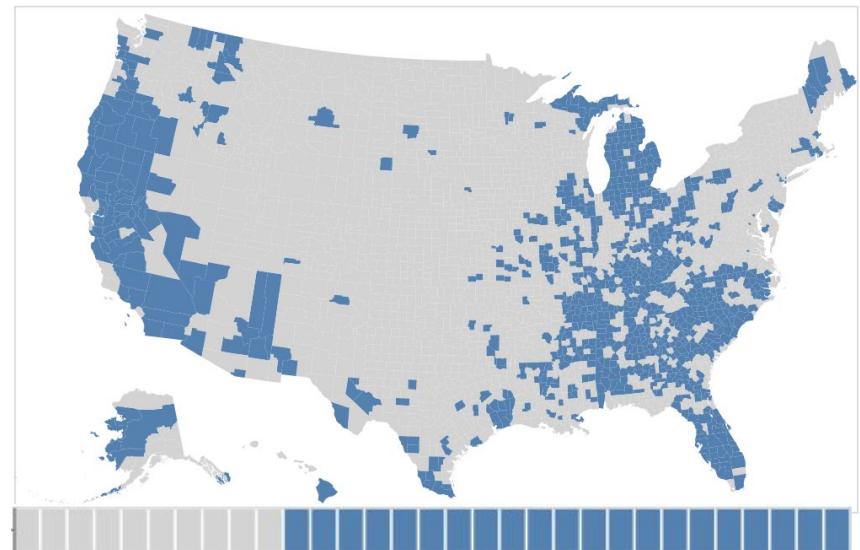


Use color intentionally:  
Mapping from numbers  
to colors

Or, a hard cut-off of <6%.

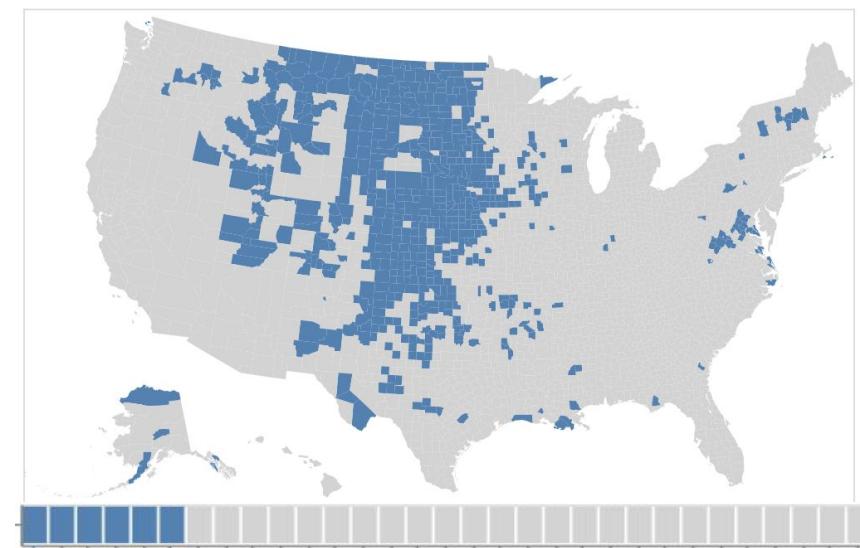
**Highlight**

Just focus on  
areas with  
>10%  
unemployment  
(in blue).



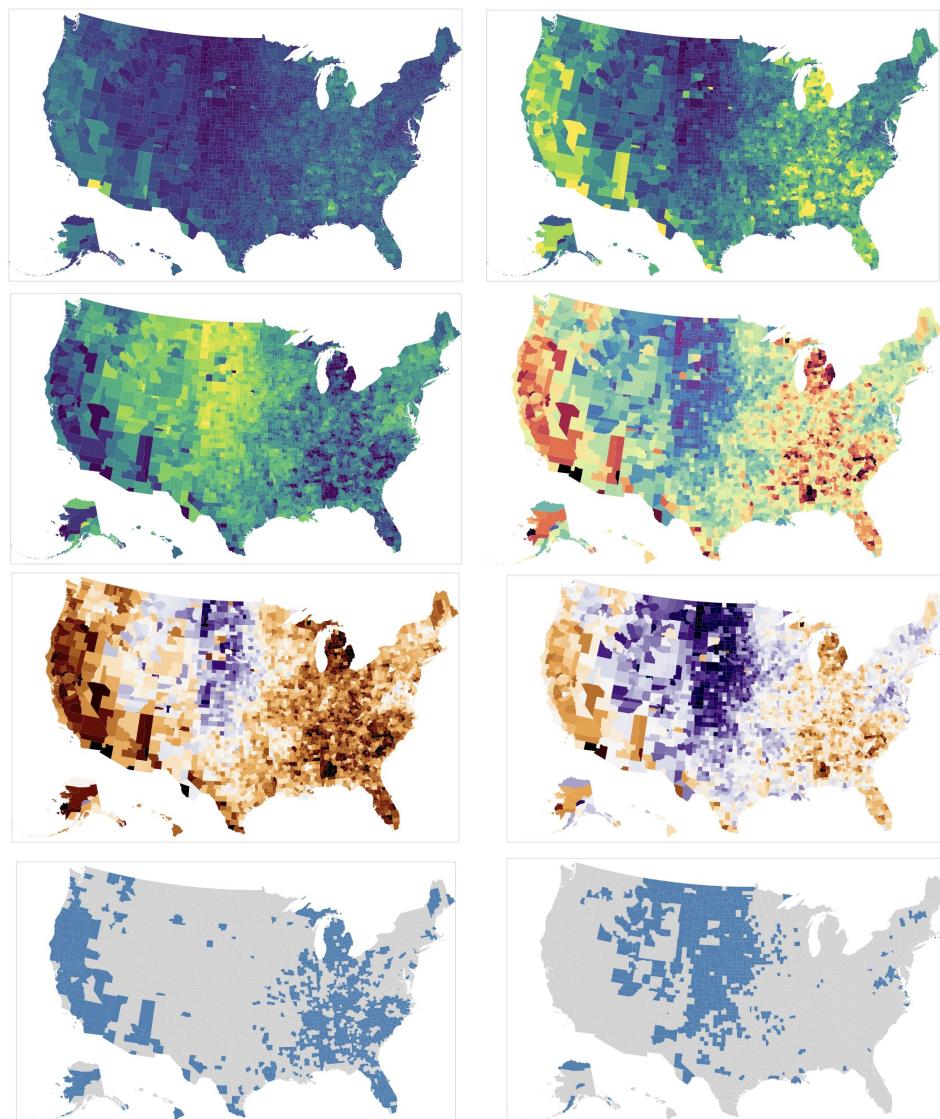
**Highlight:**

Focus on areas  
with <6%  
unemployment



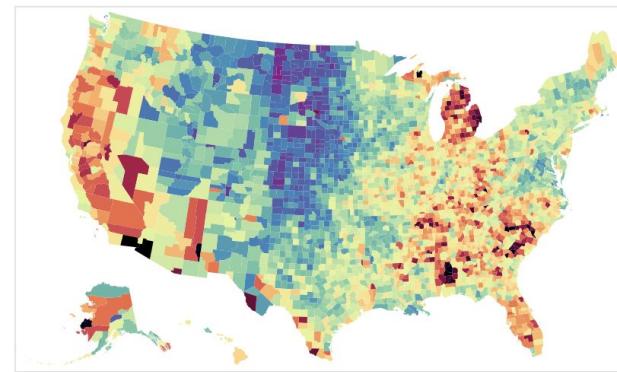
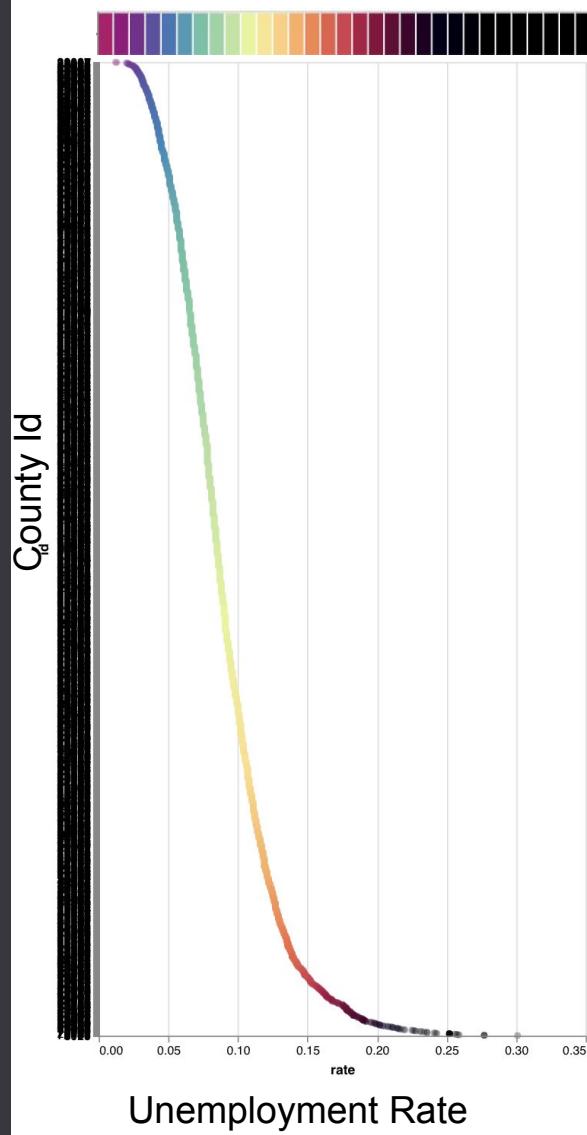
Use color intentionally:  
Mapping from numbers  
to colors

**Different mappings ask different questions** of our data. The "right" color mapping depends on what question you want to ask your data (like parameters/algorithms in stats/ML).



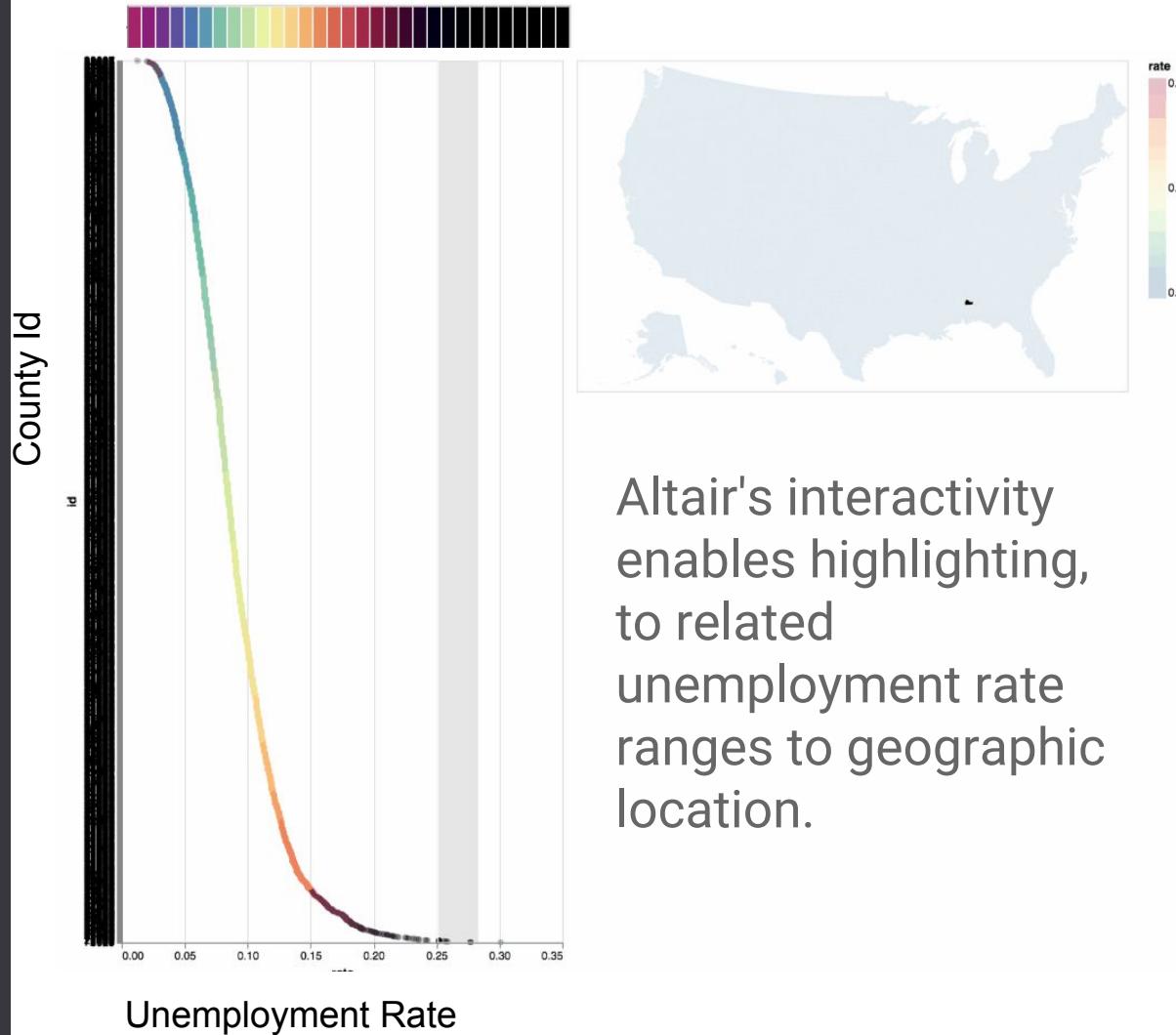
Using colormaps  
intentionally:

If you want to optimize for **quantitative** rather than qualitative comparisons; try a chart that relies on position rather than color.



Using colormaps  
intentionally:

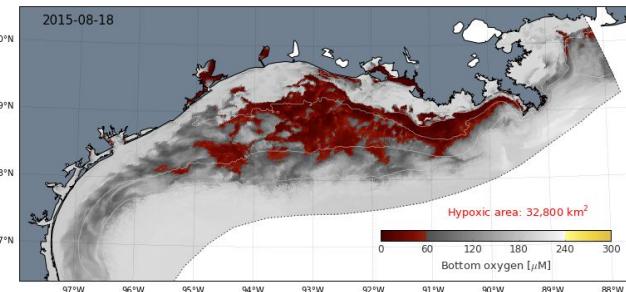
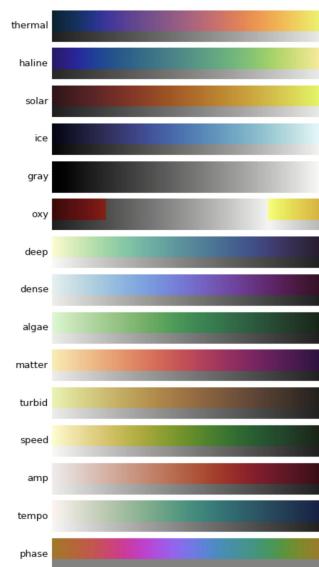
If you want to optimize for quantitative rather than qualitative comparisons; try a chart that relies on position rather than color.



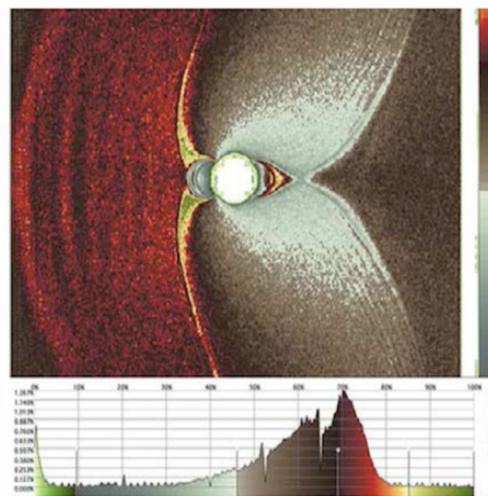
# Use color intentionally

## Resources

### cmocean color scales



### sciviscolor



# Use color intentionally

With a **highlight** color to focus attention

# Use color intentionally

6047266926697438	6047266926697438	6047266926697438
5699763068340672	5699763068340672	5699763068340672
3956259629358724	3956259629358724	3956259629358724
0675236503962934	0675236503962934	0675236503962934
5345345342645342	5345345342645342	5345345342645342
5344352635344526	5344352635344526	5344352635344526
3534350276982562	3534350276982562	3534350276982562
3576695926749261	3576695926749261	3576695926749261

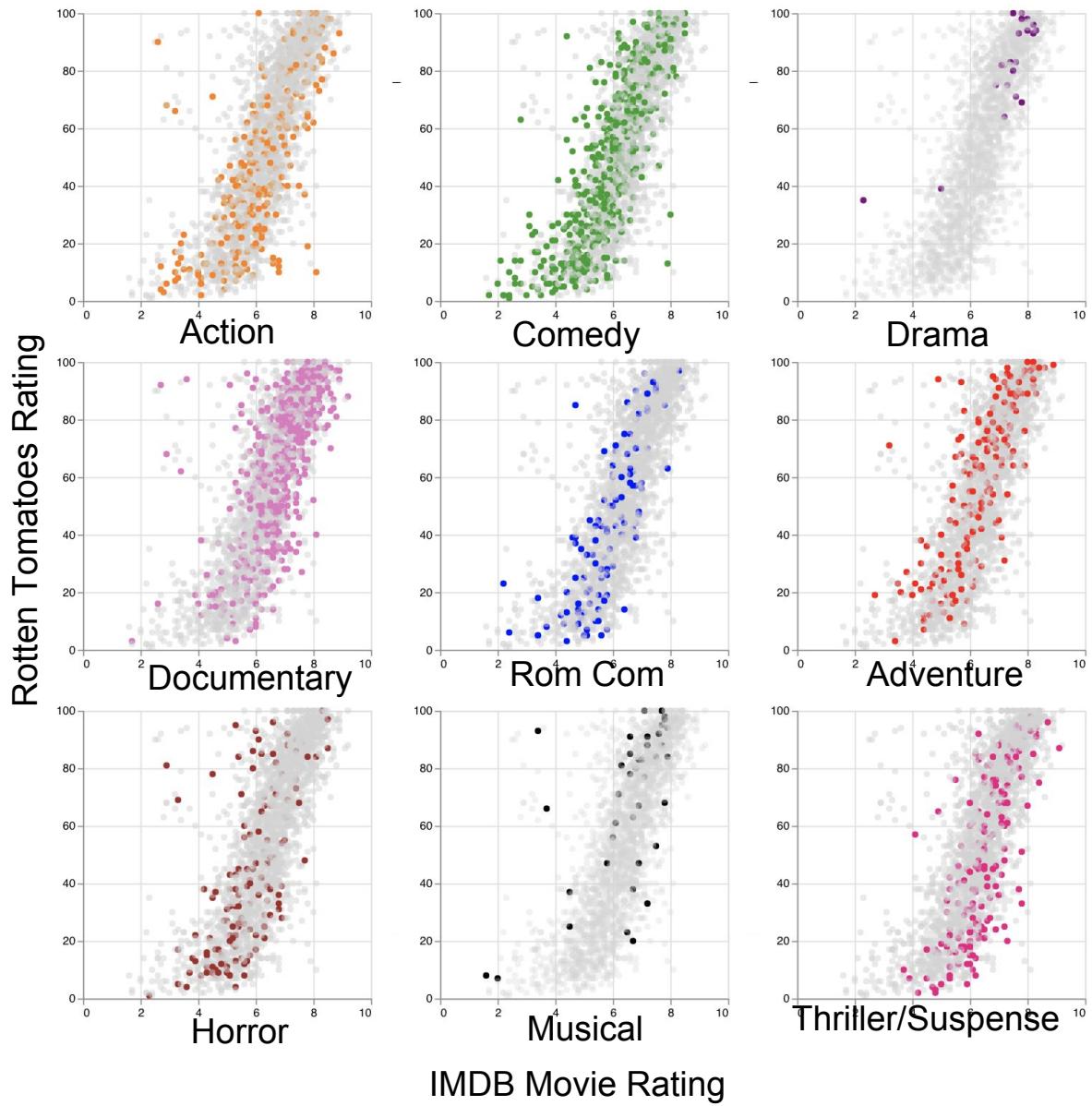
6047266926697438	6047266926697438	6047266926697438
5699763068340672	5699763068340672	5699763068340672
3956259629358724	3956259629358724	3956259629358724
0675236503962934	0675236503962934	0675236503962934
5345345342645342	5345345342645342	5345345342645342
5344352635344526	5344352635344526	5344352635344526
3534350276982562	3534350276982562	3534350276982562
3576695926749261	3576695926749261	3576695926749261

6047266926697438	6047266926697438	6047266926697438
5699763068340672	5699763068340672	5699763068340672
3956259629358724	3956259629358724	3956259629358724
0675236503962934	0675236503962934	0675236503962934
5345345342645342	5345345342645342	5345345342645342
5344352635344526	5344352635344526	5344352635344526
3534350276982562	3534350276982562	3534350276982562
3576695926749261	3576695926749261	3576695926749261

6047266926697438	6047266926697438	6047266926697438
5699763068340672	5699763068340672	5699763068340672
3956259629358724	3956259629358724	3956259629358724
0675236503962934	0675236503962934	0675236503962934
5345345342645342	5345345342645342	5345345342645342
5344352635344526	5344352635344526	5344352635344526
3534350276982562	3534350276982562	3534350276982562
3576695926749261	3576695926749261	3576695926749261

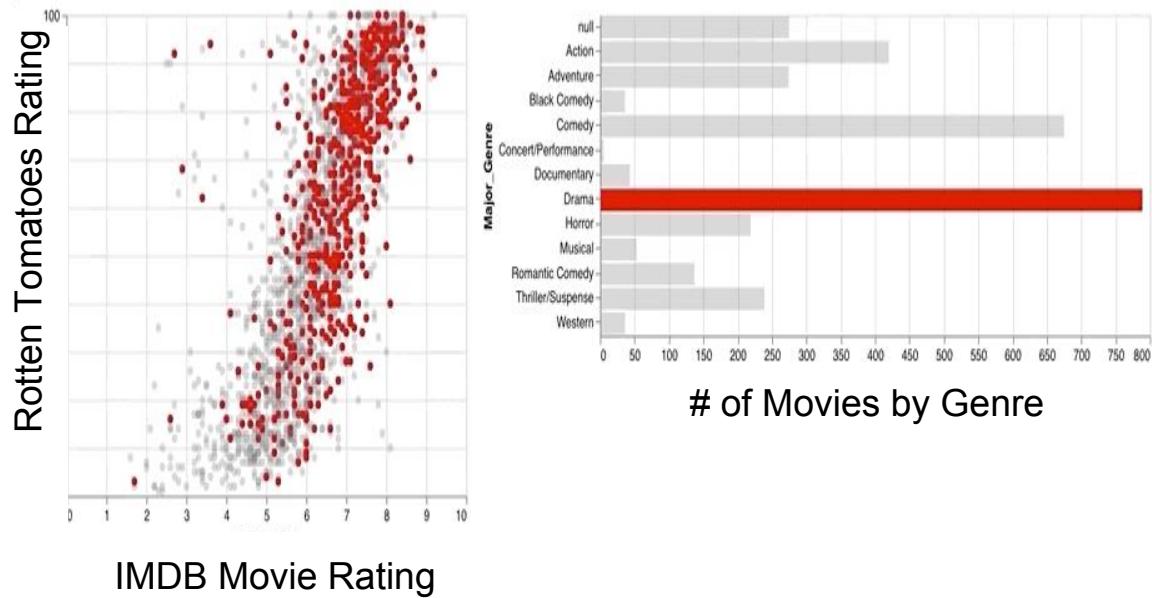
# Use color intentionally

## Interactivity example



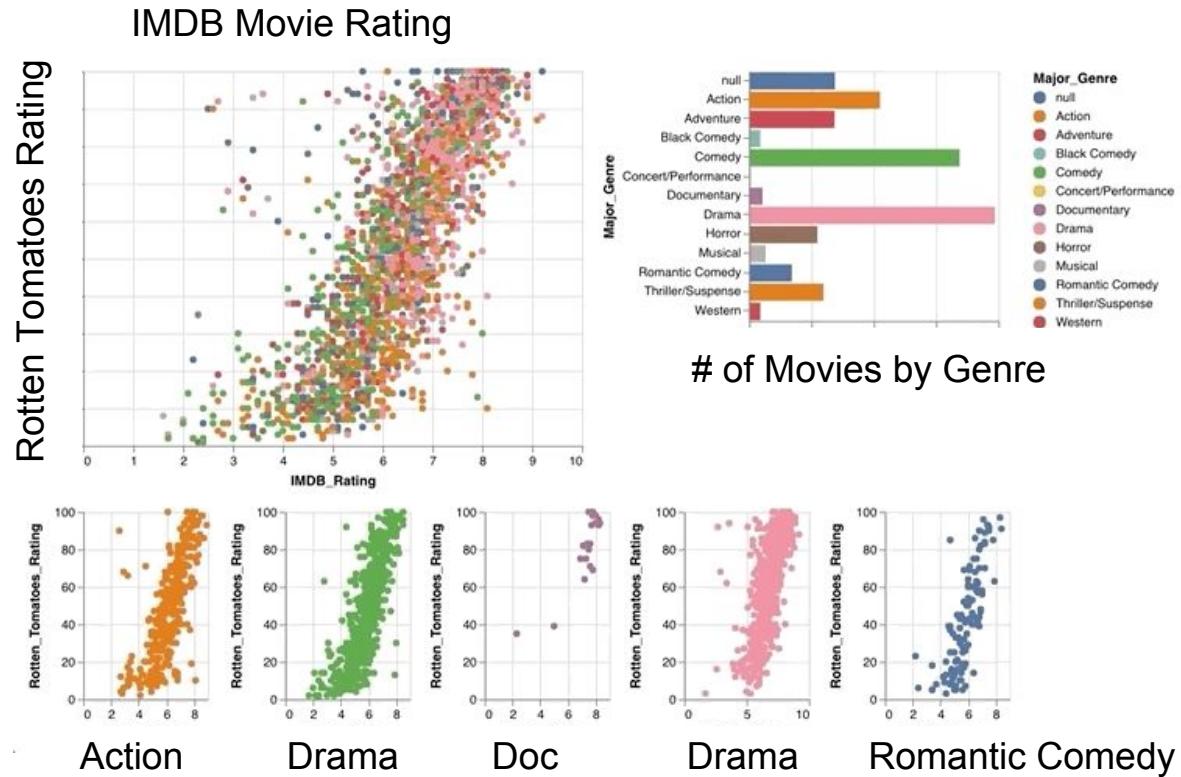
# Use color intentionally

## Interactivity example



# Use color intentionally

## Interactivity example



Optimize for your time,  
energy, and attention

Optimize for time,  
energy, attention

In science, we often have "data processing pipelines" (processes you run regularly on similar types of data).

Optimize for time,  
energy, attention

In science, we often have "data processing pipelines" (processes you run regularly on similar types of data).

Want to know:

Is the data ok? Broken?

Is there anything "interesting" here?

Optimize for time,  
energy, attention

In science, we often have "data processing pipelines" (processes you run regularly on similar types of data).

Want to know:

Is the data ok? Broken?

Is there anything "interesting" here?

Create:

Many charts

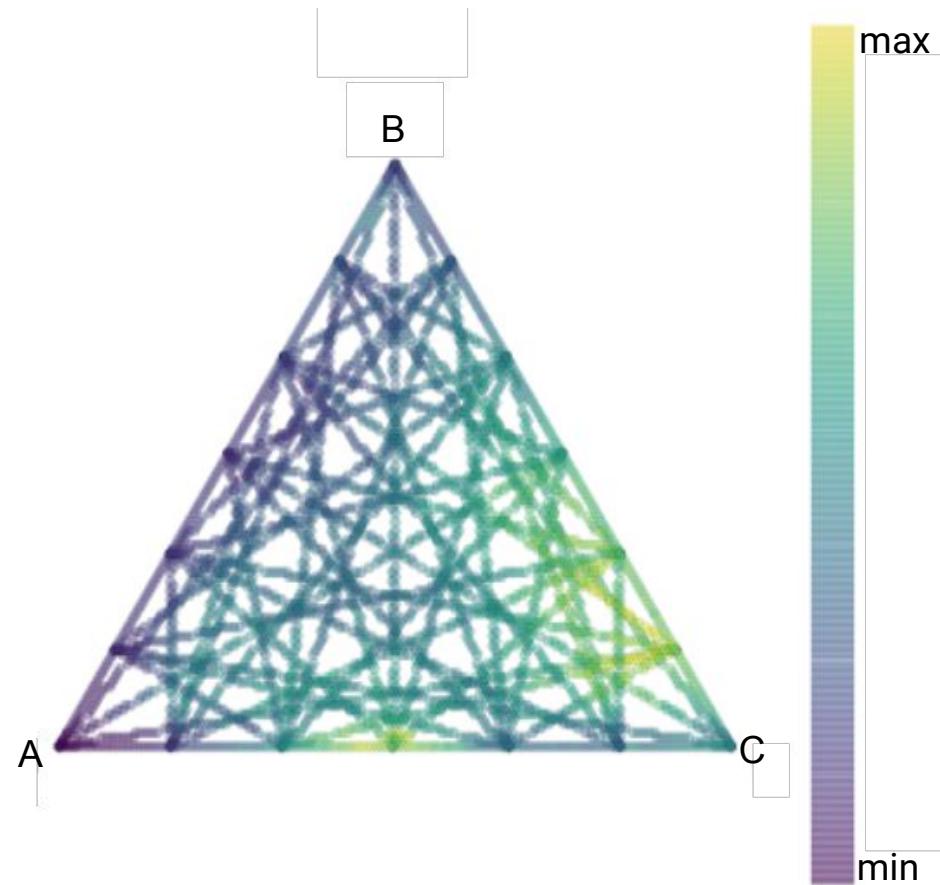
Mostly boring

Important stuff is obvious

Optimize for time,  
energy, attention

## Case Study - Collaboration w/ Google Accelerated Science & CalTech

Looking for areas of interest (local min, max,  
slope changes) in 3-part mix

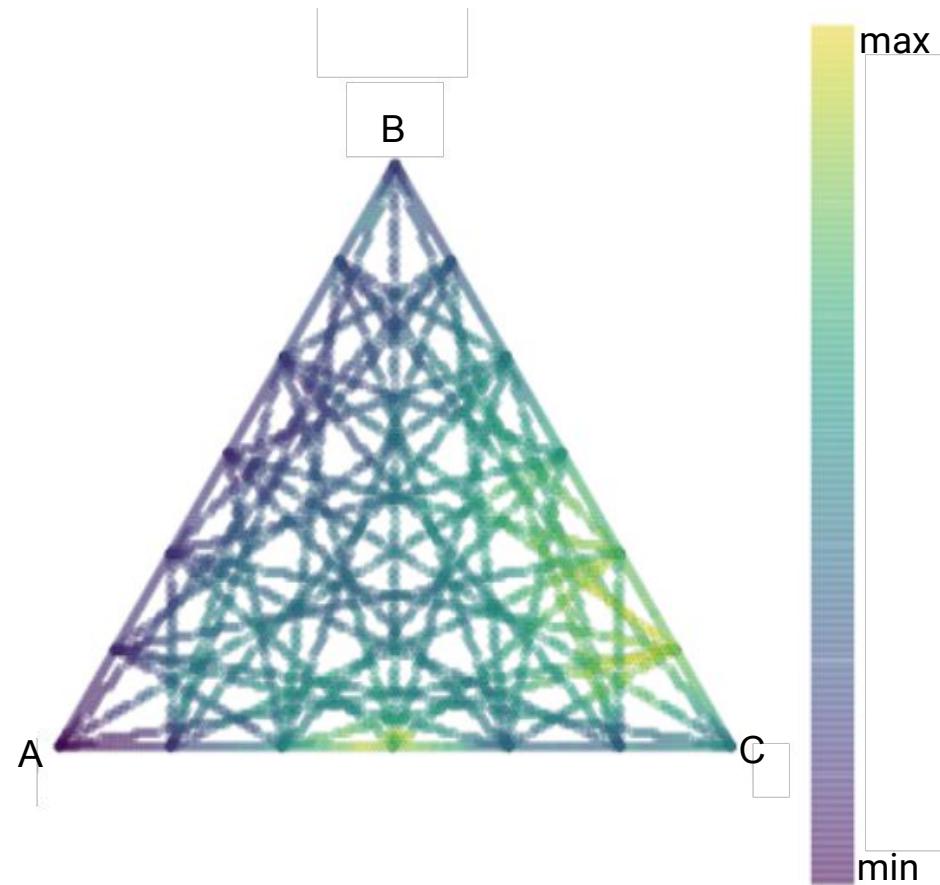


Optimize for time,  
energy, attention

Issues w/ original:

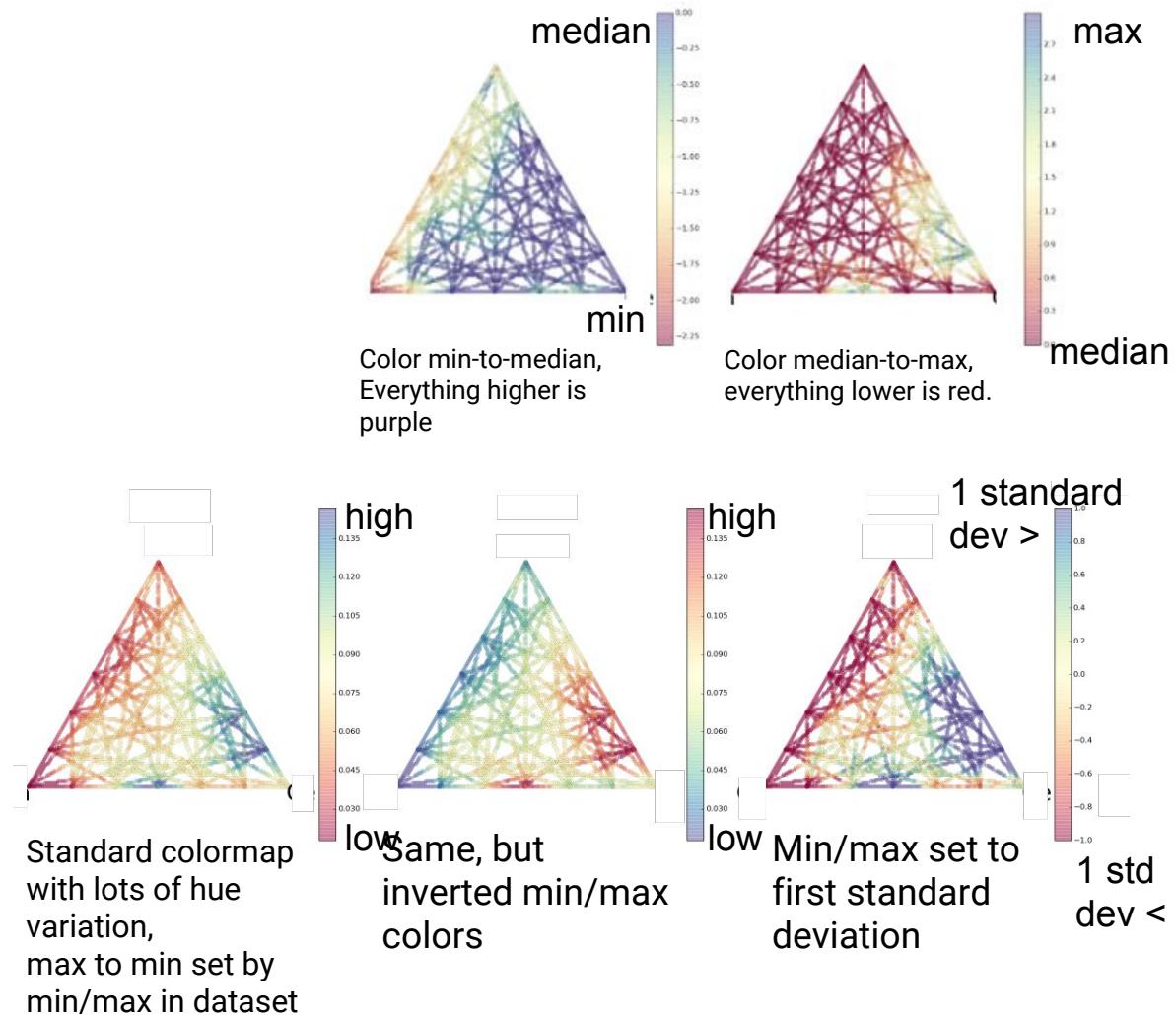
1 - Max values more eye-catching than min, middle. Hard to see "interesting" in mid-colorscale.

2 - Extreme values wash out the color scale



Optimize for time,  
energy, attention

# Small multiples; Intentional color Optimize for seeing if something is important, quickly



In summary:

Is it good or bad depends  
on your goals, context, and  
constraints.

In summary:

Create small multiples (many, related charts)

Use color intentionally

Optimize your charts for your  
time/energy/attention

# Resources

[cmocean color scales](#)

[sciviscolor](#)

[Varying binwidths and histograms](#)

[Colin Ware's Information Visualization book](#)

[Cole Nussbaumer - Storytelling with Data](#)

## Examples from the talk

[colab with geographic maps/colors](#)

[Colab with small multiples](#)