

Lecture 6

Filter + Aggregate

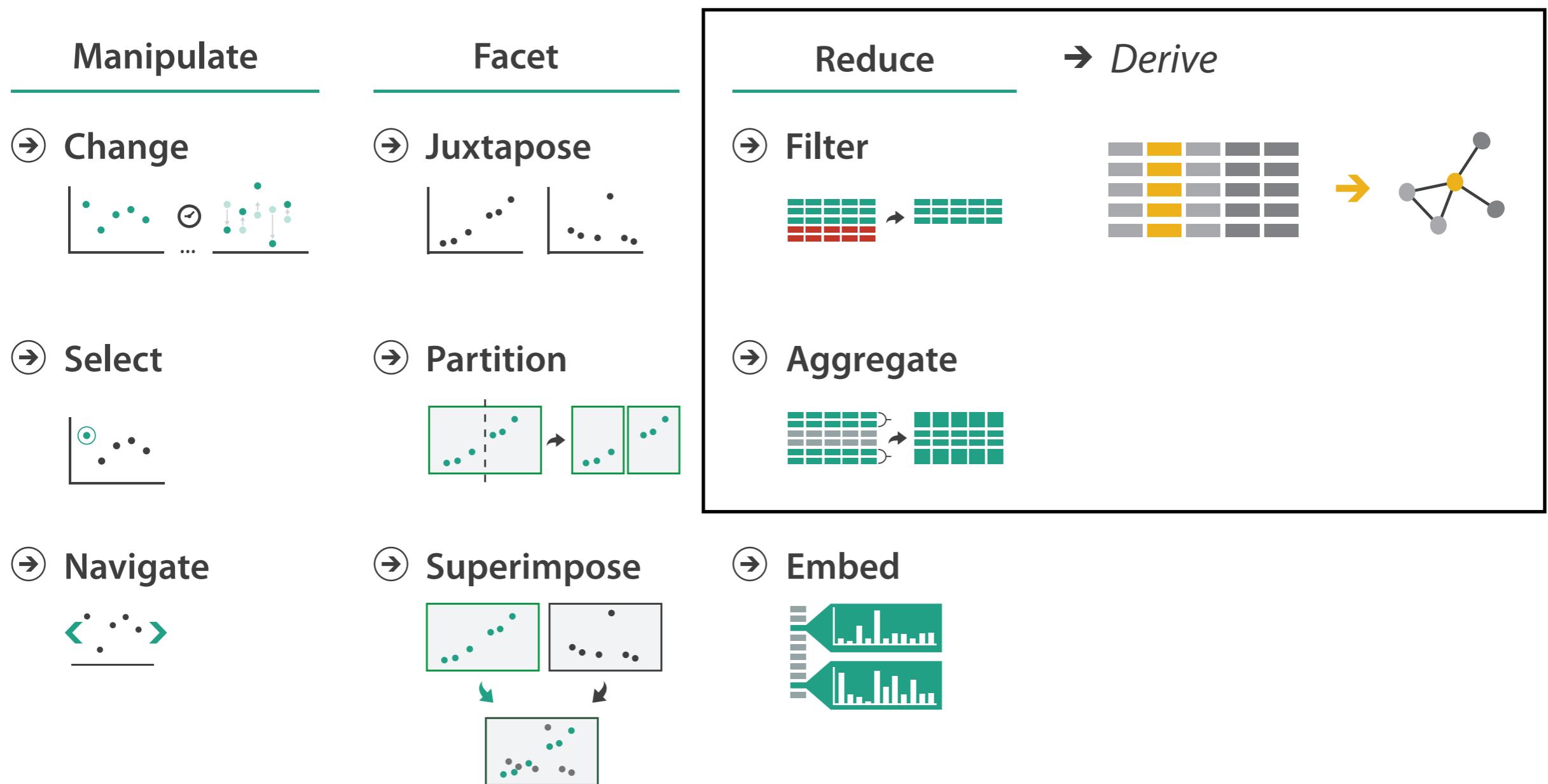
DSCI 532, Data Visualization II

January 21, 2019

Cydney Nielsen

Senior Designer, Microsoft
Adjunct Professor, UBC Department of Computer Science

Dealing with complexity



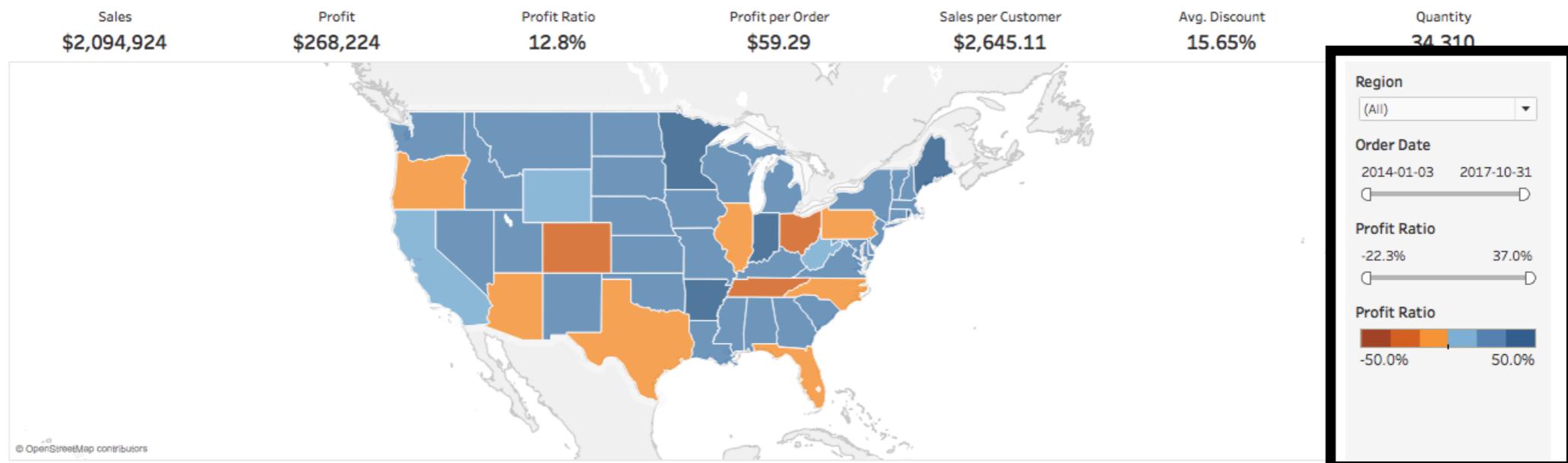
Filter

reducing complexity and creating focus

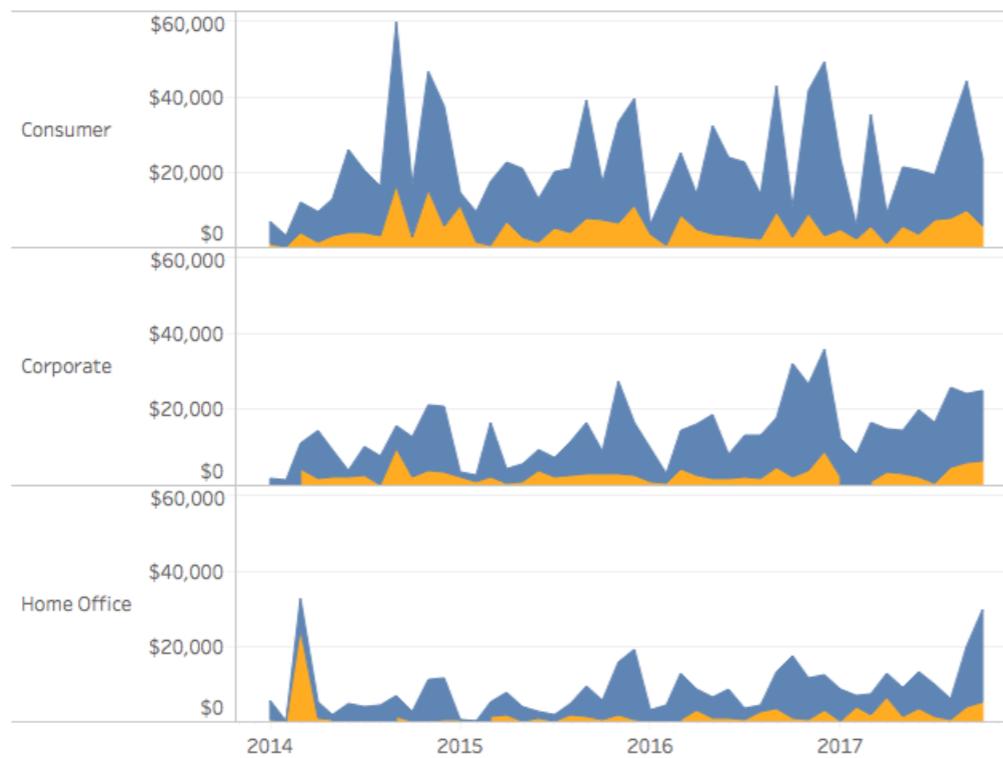
Dashboards

Filter widgets

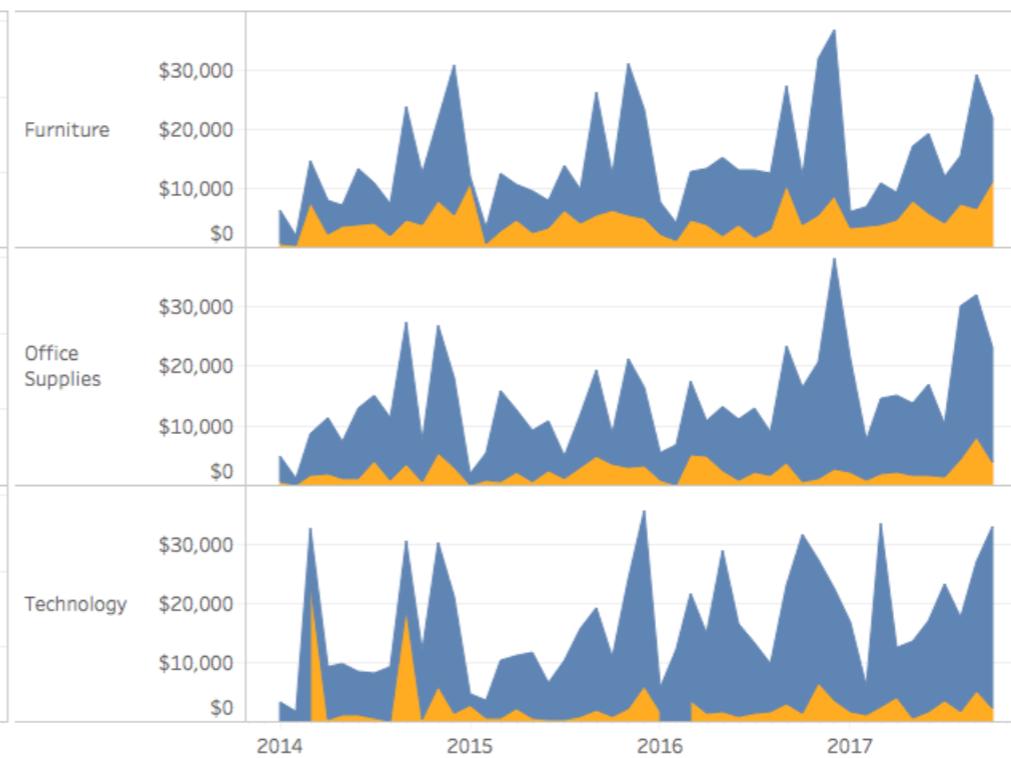
Executive Overview - Profitability (All)



Monthly Sales by Segment - States: All



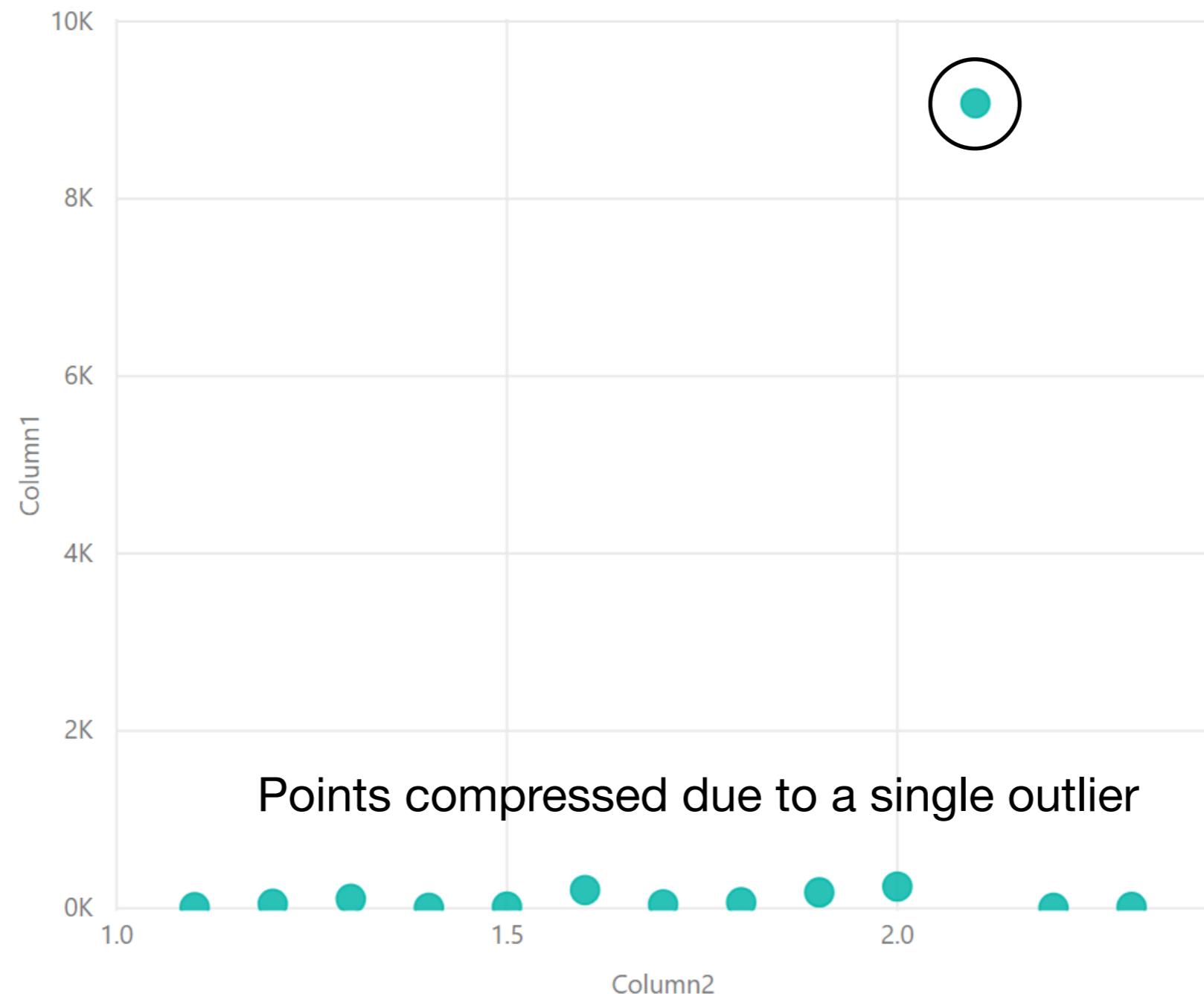
Monthly Sales by Product Category - States: All



Tableau

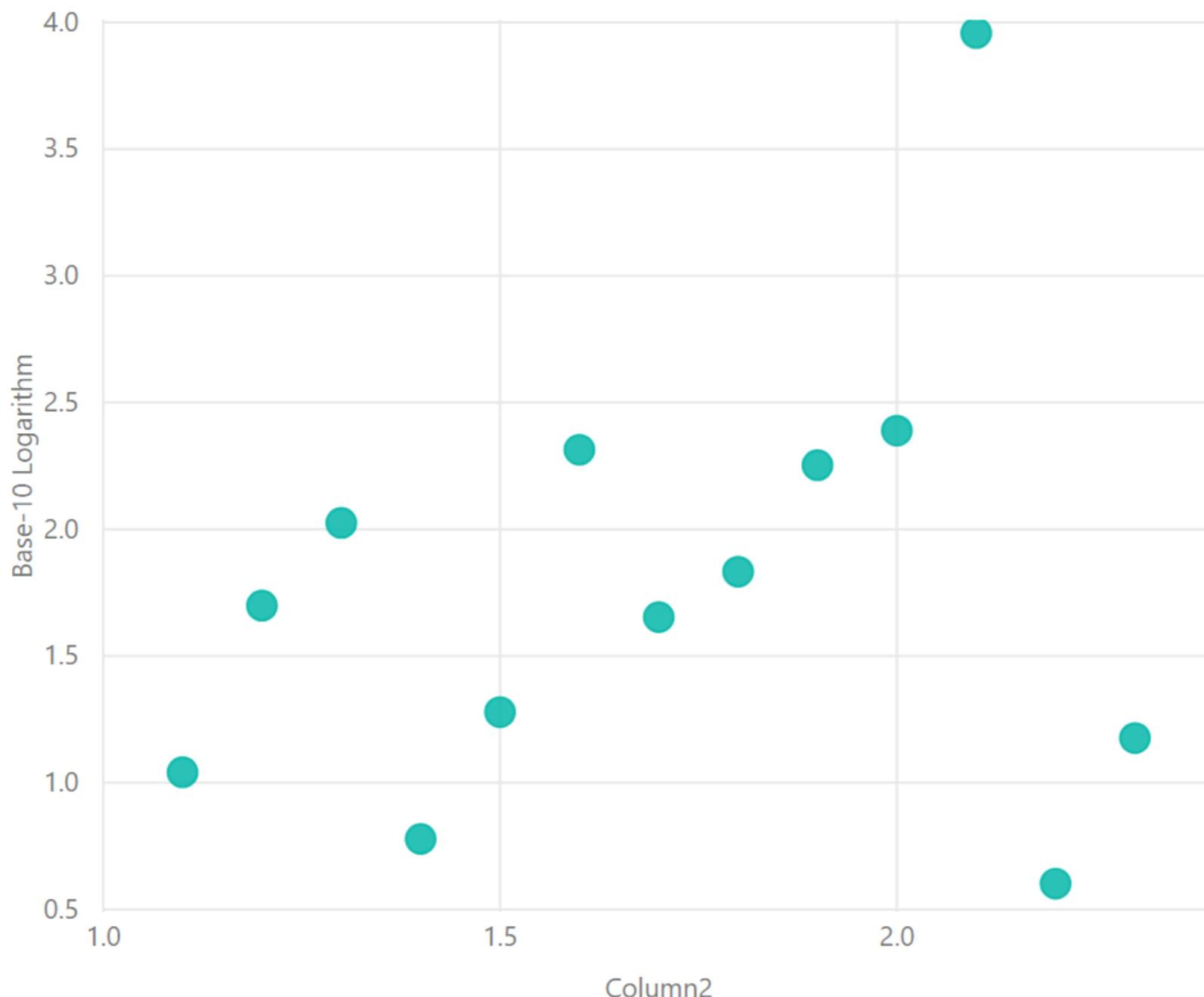
Dealing with outliers

Do you want to keep it or remove it?



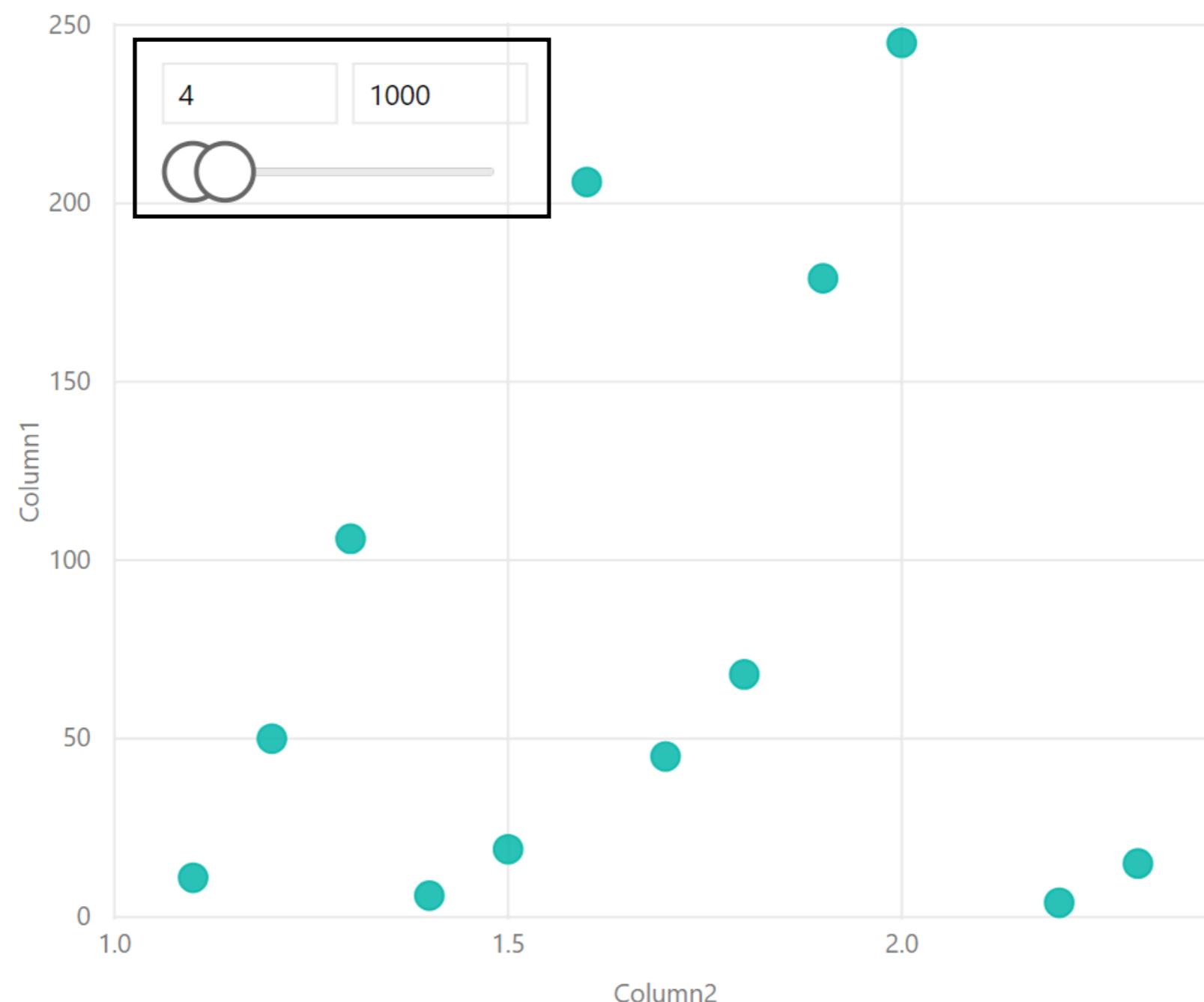
Dealing with outliers

log scale helps - appropriate if you want to keep all points



Dealing with outliers

interactive filtering is another solution - appropriate when want to remove outliers



Faceted search

[Home](#) / All / [Camping and hiking](#) / [Packs and bags](#)

CATEGORY

◀ All

Camping and hiking

92 items

Packs and bags (92)

Backpacking packs (56)

30 to 44 L

45 to 59 L

Daypacks (50)



SAVE 30%

★★★★★ (2)

Gregory Targhee 45 Backpack -
Unisex

\$181.00 | \$259.00

Compare



SAVE 20%

★★★★★ (3)

Gregory Zulu 40 Backpack -
Unisex

\$167.95 | \$209.00

Compare



SAVE 20%

★★★★★ (9)

Gregory Jade 28 Daypack -
Women's

\$127.95 | \$159.00

Compare

FEATURES

- Compression straps (83)
- Hydration compatible (73)
- Ice axe/pole attachment (70)
- Hipbelt pocket(s) (59)
- Ventilated back panel (44)

[MORE](#)



SAVE 29%

★★★★★ (13)

VOLUME

- 30 to 44 L (68)
- 45 to 59 L (31)
- 20 to 29 L (64)
- 10 to 19 L (46)
- 60 to 74 L (36)

[MORE](#)



SAVE 20%

★★★★★ (0)

Osprey Kestrel 48L Backpack -
Men's



SAVE 20%

★★★★★ (4)

Gregory Zulu 55 Backpack -
Unisex



ADVANCED SEARCH

[Donors](#) [Genes](#) [Mutations](#)

▼ [Donor](#)

e.g. DO45299, SA501608

[Upload Donor Set](#)

▼ [Primary Site](#)



	Blood	3,180
	Brain	2,506
	Breast	1,976
	Liver	1,840
	Kidney	1,605

[Select all](#)

▼ [17 more](#)

▼ [Project](#)



	ALL-US	1,267
	AML-US	358
	BLCA-CN	103
	BLCA-US	412
	BOCA-FR	100

[Select all](#)

▼ [79 more](#)

▼ [Study](#)

PCAWG 2,809

None 21,268

[Select all](#)

▼ [Gender](#)

Male 13,006

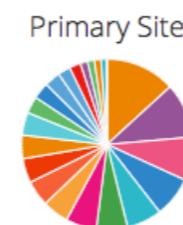
Donors
24,077

Genes
57,905

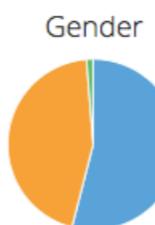
Mutations
77,462,290



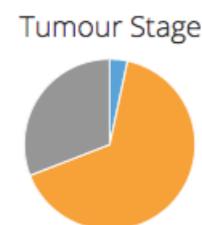
Vital Status



Disease Status



Relapse Type



Age

[Show More](#)

OncoGrid

Download Donor Data

View in Data Repositories

Save/Edit Donor Results

Donors

Showing 1 - 10 of 24,077 donors



	ID	Project	Site	Gender	Age	Stage	Survival (days)	Available Data Types:												# Mutations	# Genes
								SSM	CNSM	StSM	SGV	METH-A	METH-S	EXP-A	EXP-S	PEXP	miRNA-S	JCN			
	DO232761	PBCA-US	Brain	Male	3		2,992	✓	—	—	—	—	—	—	—	—	—	—	1,378,562	56,205	
	DO232224	PBCA-US	Brain	Female	7		547	✓	—	—	—	—	—	—	✓	—	—	—	—	1,015,534	55,338
	DO222843	MELA-AU	Skin	Male	76	IIC	907	✓	—	✓	—	—	—	—	—	—	—	—	—	964,360	51,565
	DO222837	MELA-AU	Skin	Male	82	IIB	1,110	✓	—	✓	—	—	—	—	—	—	—	—	—	786,166	49,858
	DO50970	LICA-FR	Liver	Female	72			✓	—	—	—	—	—	—	✓	—	—	—	—	610,827	48,874
	DO222363	MELA-AU	Skin	Male	81	IIC	154	✓	—	✓	—	—	—	—	—	—	—	—	—	775,848	47,839
	DO220886	MELA-AU	Skin	Male	56	IA/IB	7,730	✓	—	✓	—	—	—	—	—	—	—	—	—	530,806	46,681
	DO220877	MELA-AU	Skin	Male	25	IV	427	✓	—	✓	—	—	—	—	—	—	—	—	—	928,259	46,246
	DO222875	MELA-AU	Skin	Male	79	IIA	1,192	✓	—	✓	—	—	—	—	—	—	—	—	—	819,954	45,993
				Male	59	2	1,570	✓	✓	✓	—	—	—	—	—	—	—	—	—	665,096	45,747

<<< < 1 2 3 4 5 > >>

What works well?

What could be improved?

<https://dcc.icgc.org/search>

Dynamic filters

Before | No selections

Primary Site	
<input type="checkbox"/> Blood	3,180
<input type="checkbox"/> Brain	2,506
<input type="checkbox"/> Breast	1,976
<input type="checkbox"/> Liver	1,840
<input type="checkbox"/> Kidney	1,605
Select all	▼ 17 more
Project	
<input type="checkbox"/> ALL-US	1,267
<input type="checkbox"/> AML-US	358
<input type="checkbox"/> BLCA-CN	103
<input type="checkbox"/> BLCA-US	412
<input type="checkbox"/> BOCA-FR	100
Select all	▼ 79 more

After | Brain selected

Primary Site	
<input checked="" type="checkbox"/> Brain	2,506
<input type="checkbox"/> Blood	3,180
<input type="checkbox"/> Breast	1,976
<input type="checkbox"/> Liver	1,840
<input type="checkbox"/> Kidney	1,605
<input type="checkbox"/> Prostate	1,419
Select all / none	▼ 16 more
Project	
<input type="checkbox"/> GBM-CN	80
<input type="checkbox"/> GBM-US	596
<input type="checkbox"/> LGG-US	515
<input type="checkbox"/> PBCA-DE	554
<input type="checkbox"/> PBCA-US	649
Select all	▼ 1 more

Also selection moves to the top

Truncate list

Project list
changes based on
Primary Site

Query is visible

ADVANCED SEARCH

Donors Genes Mutations

Share Show PQL Primary Site IS Blood AND Project IN (AML-US , CLLE-ES)

Donor
e.g. DO45299, SA501608 Upload Donor Set

Primary Site
Blood 909
Select none

Project
CLLE-ES 551
AML-US 358
ALL-US 1,267

Donors 909 Genes 44,131 Mutations 423,980

Project Primary Site Gender

Vital Status Disease Status Relapse Type

Show More

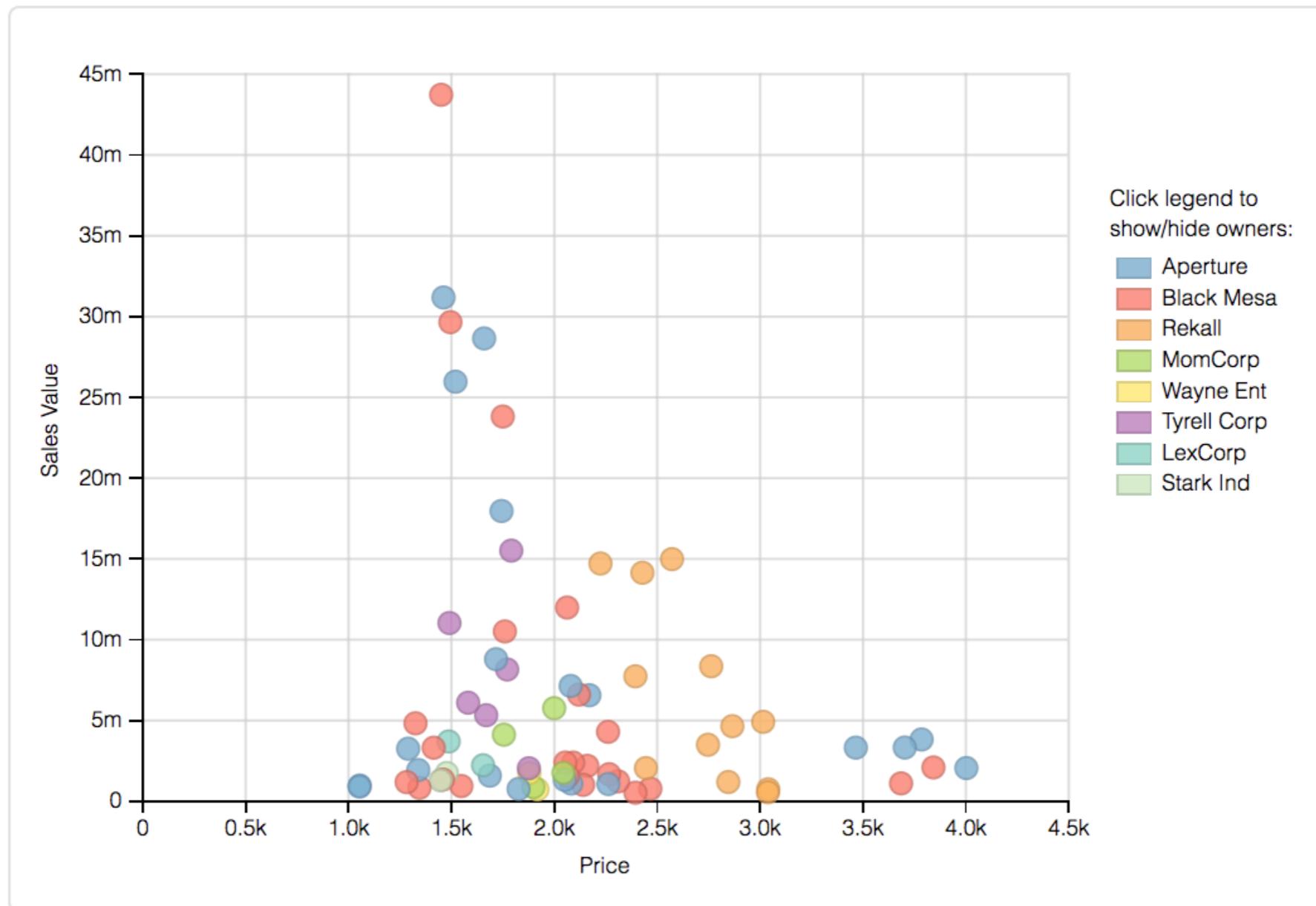
The screenshot shows the DCC ICGC search interface. At the top, there's a search bar with the query: "Primary Site IS Blood AND Project IN (AML-US , CLLE-ES)". Below the search bar, there are tabs for "Donors", "Genes", and "Mutations". On the left, there's a sidebar with a "Donor" section containing a text input "e.g. DO45299, SA501608" and a "Upload Donor Set" button. Below that is a "Project" section with three items: "CLLE-ES" (551), "AML-US" (358), and "ALL-US" (1,267). The main area has summary counts: "Donors 909", "Genes 44,131", and "Mutations 423,980". To the right, there are three pie charts: "Project" (orange and blue), "Primary Site" (orange), and "Gender" (orange and blue). Below the charts are sections for "Vital Status", "Disease Status", and "Relapse Type", each with a "Show More" button.

Query also captured in the page URL for easy sharing of the filter state

<https://dcc.icgc.org/search>

Interactive legend

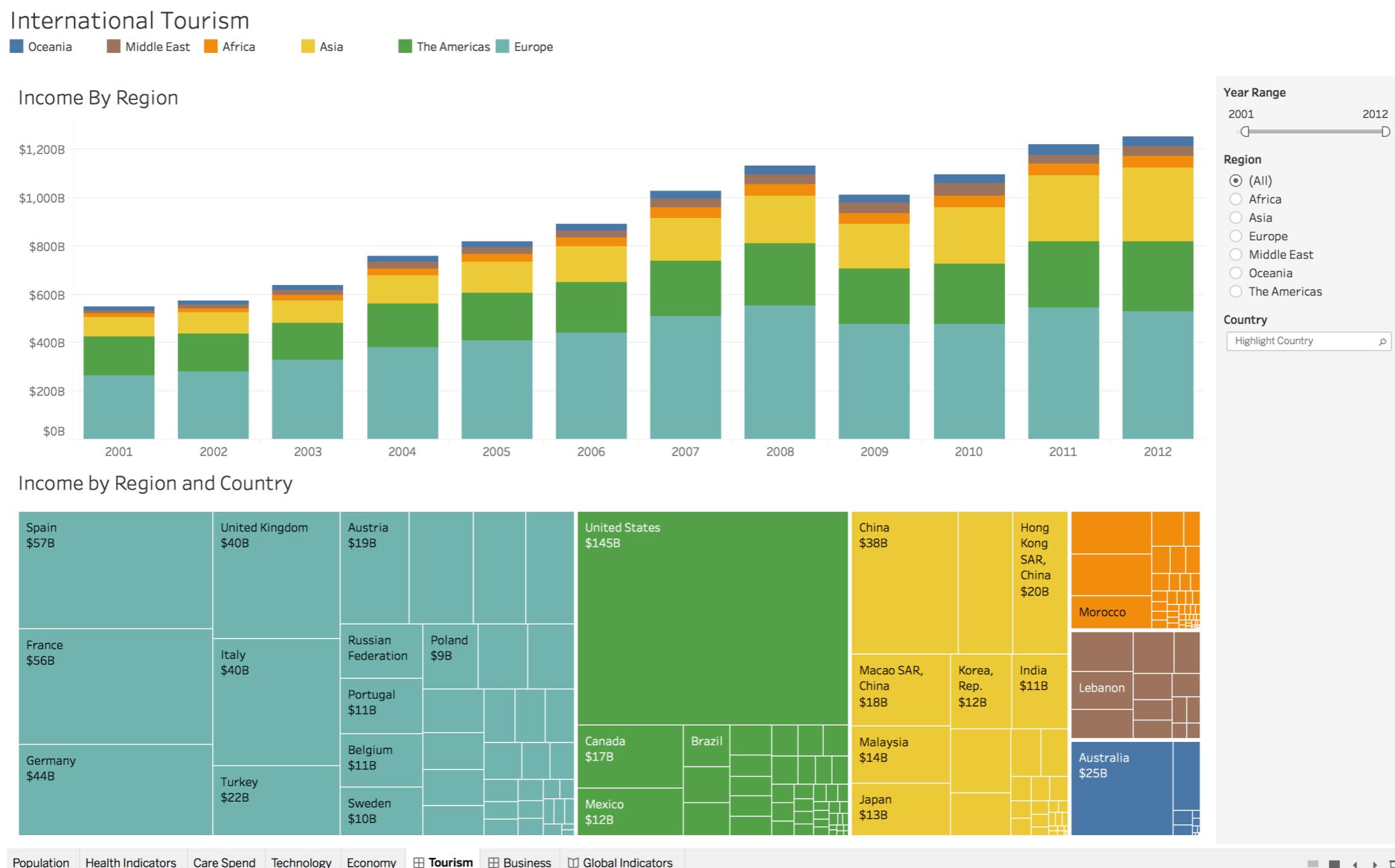
Legend also functions like a filter widget



http://dimplejs.org/advanced_examples_viewer.html?id=advanced_interactive_legends

Cross filtering

Selection on plots themselves serve as filters to linked views



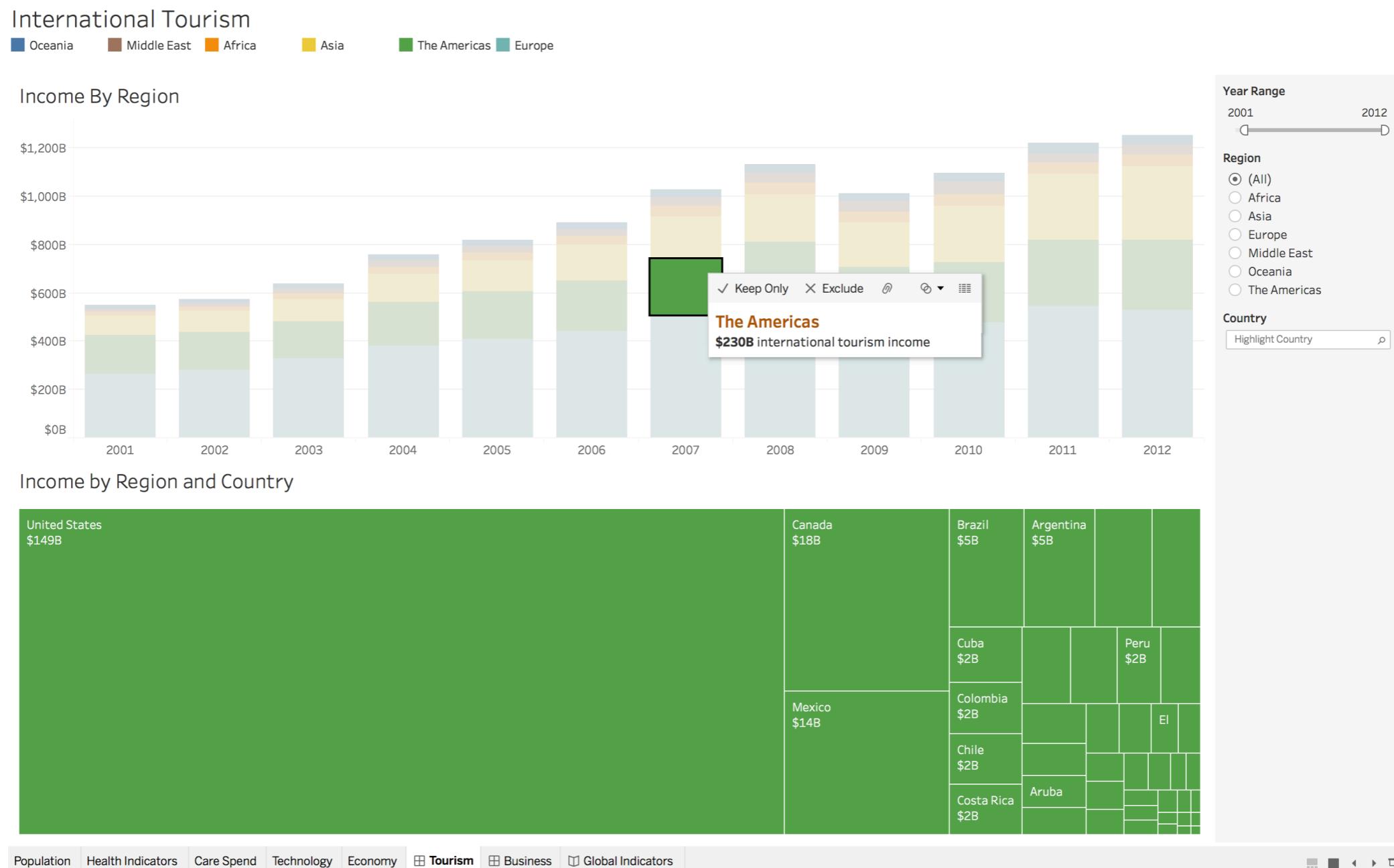
Tableau

World Indicators Sample Workbook

13

Cross filtering

Selection used to filter data in the linked views



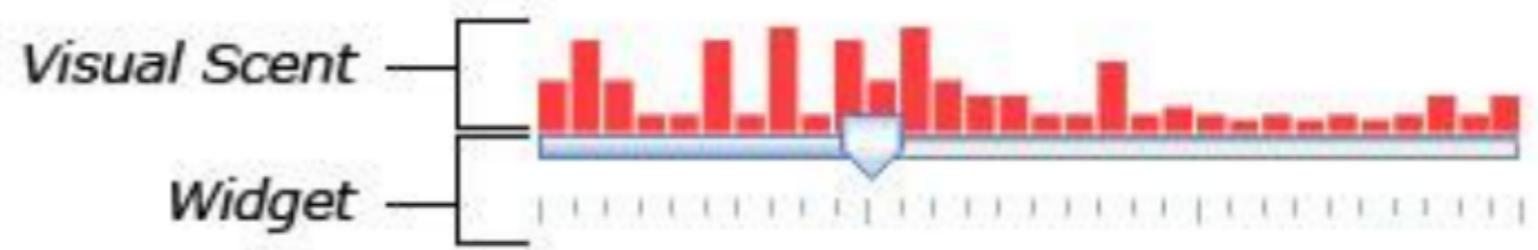
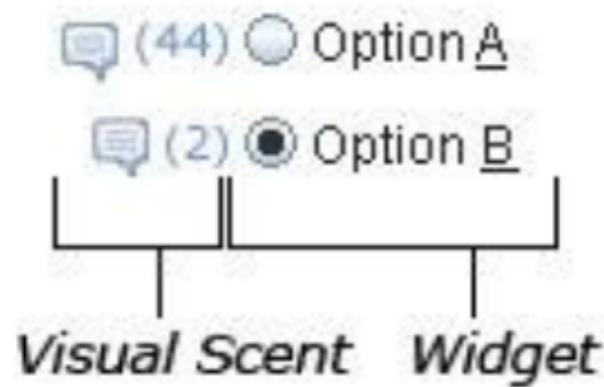
More space for selected values, but lose context

Tableau

World Indicators Sample Workbook

Scented widgets

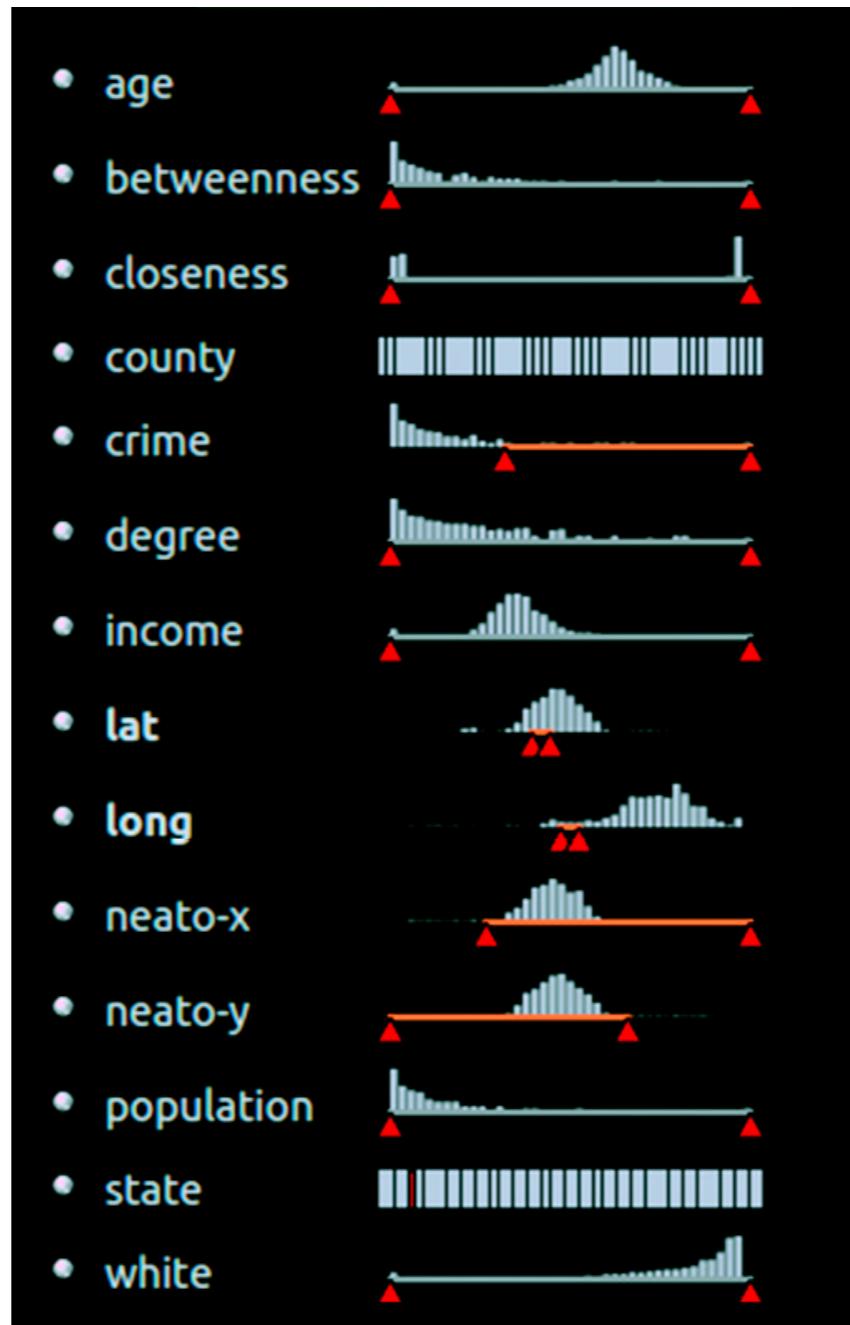
Embedded visualizations that help you decide what to filter



<http://vis.stanford.edu/files/2007-ScentedWidgets-InfoVis.pdf>

Scented widgets

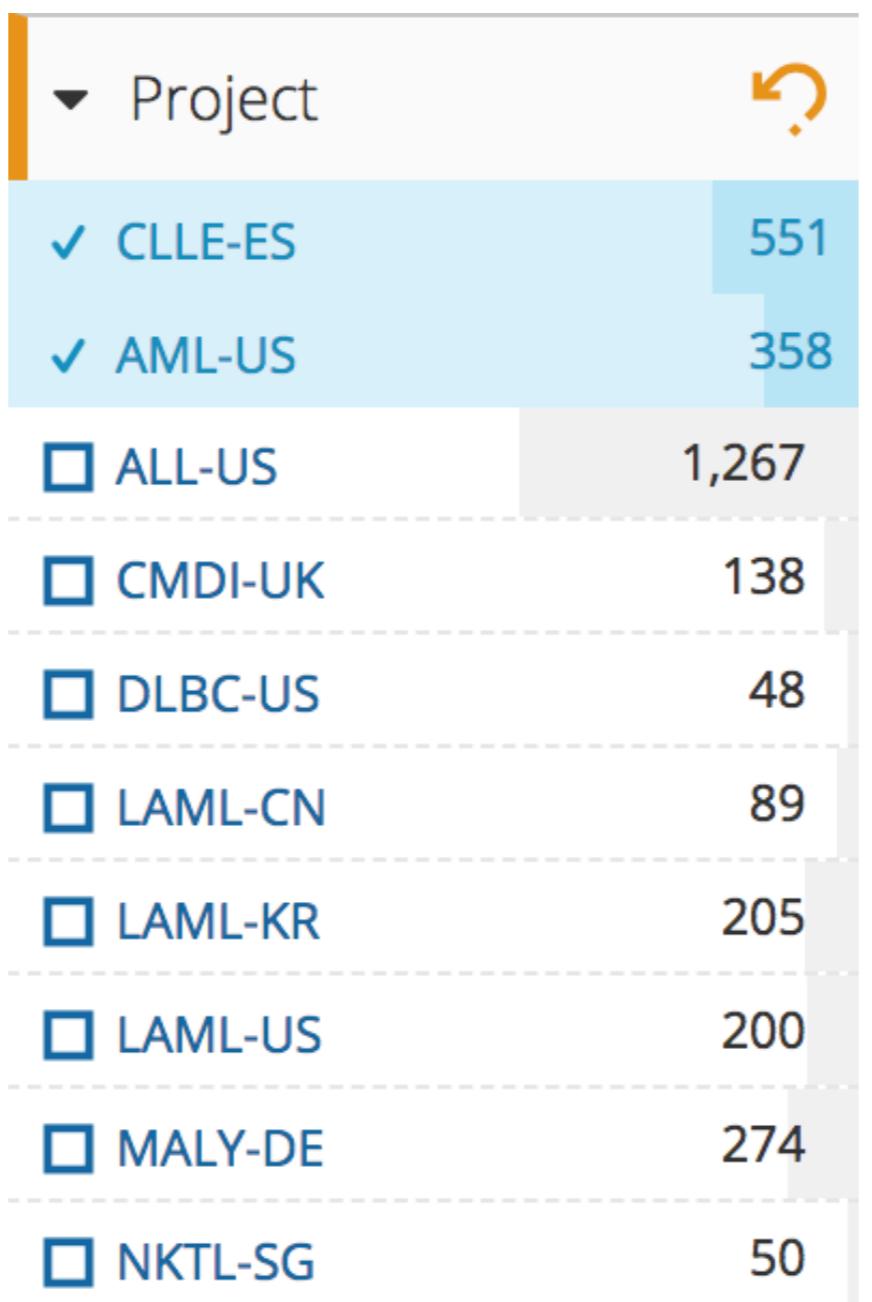
Embedded visualizations that help you decide what to filter



<https://www.win.tue.nl/~selzen/paper/InfoVis2014.pdf>

Scented widgets

Embedded visualizations that help you decide what to filter



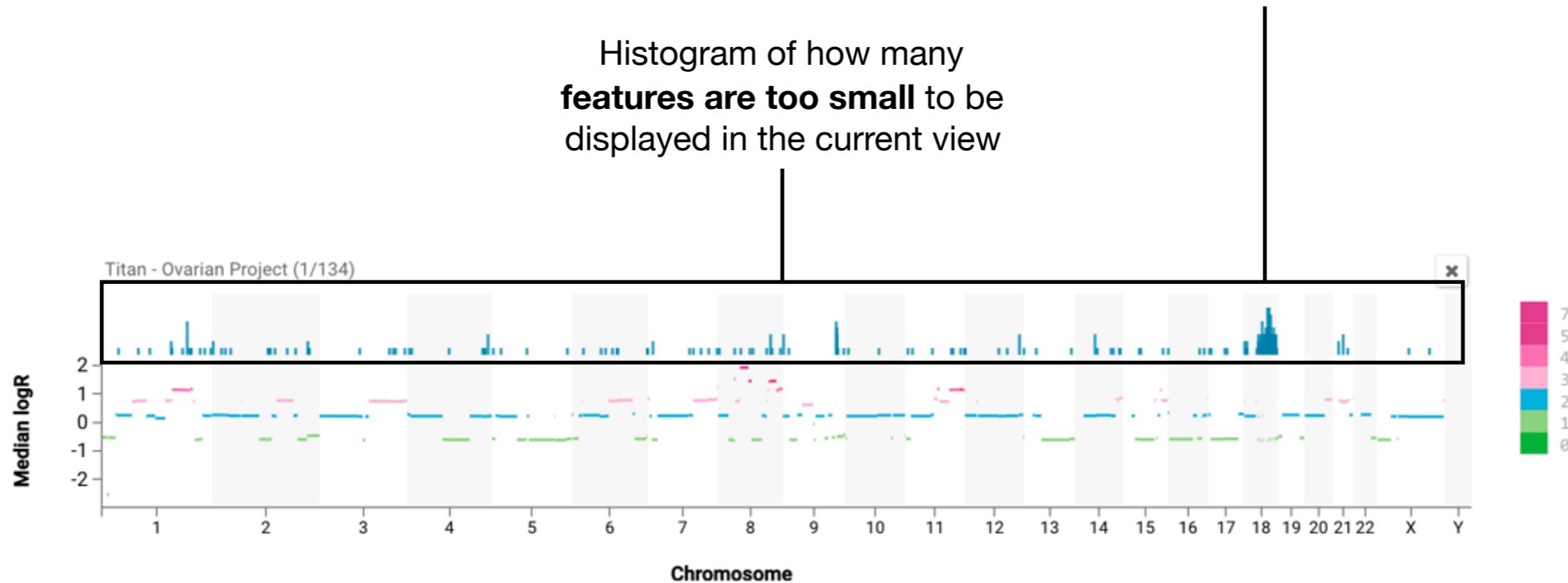
Bar chart shows
number of cases
in each category

Out of sight, out of mind

A visual reminder of what you removed

As you zoom into this peak,
small features become visible
(peak tells you where to look)

Histogram of how many
features are too small to be
displayed in the current view



← Human genome is >3 billion base pairs long →

Want to display in <1000 pixels

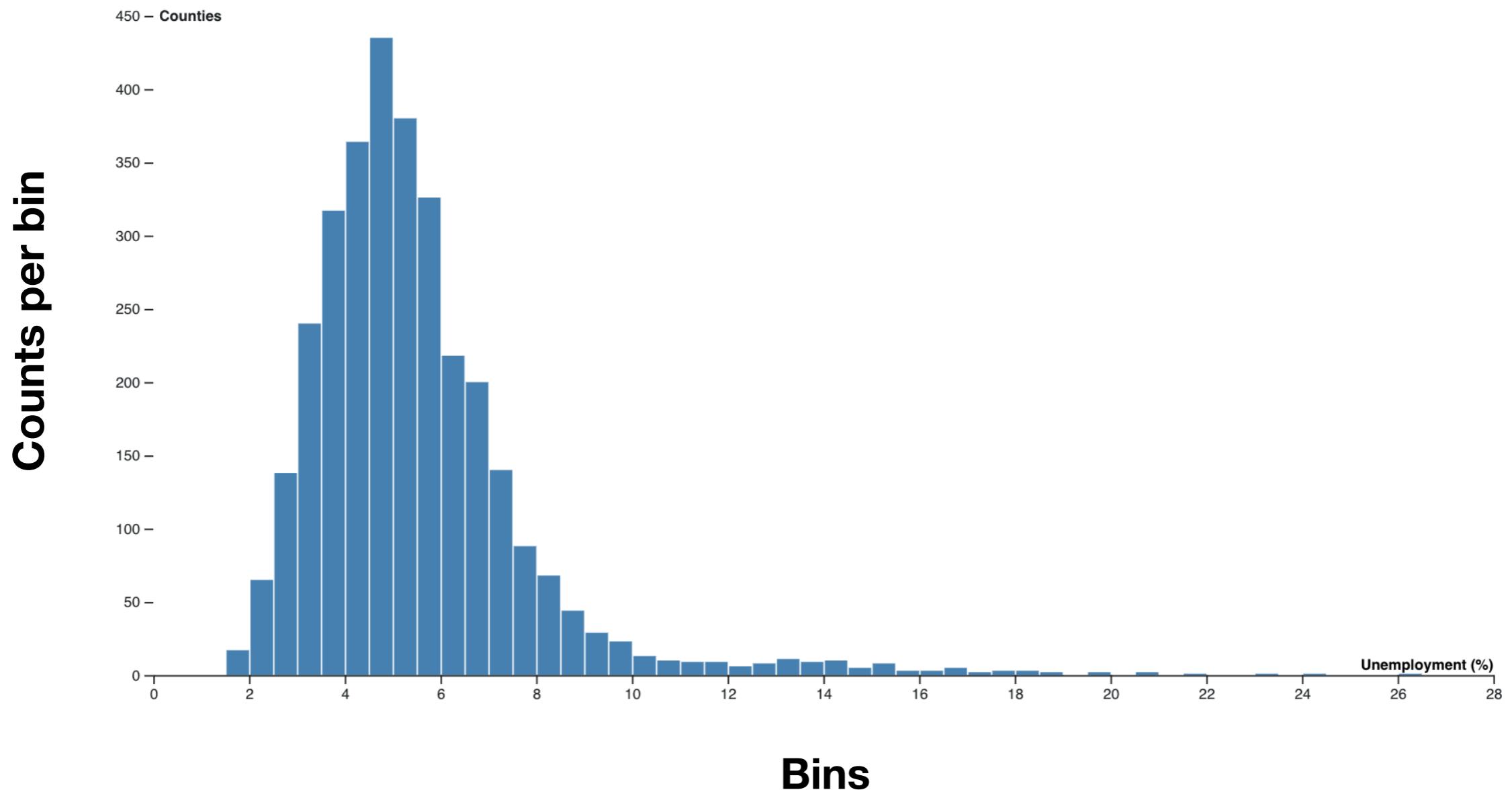
At that scale, small features won't be visible

Aggregation

Common aggregations

Histogram | Visualizing a distribution

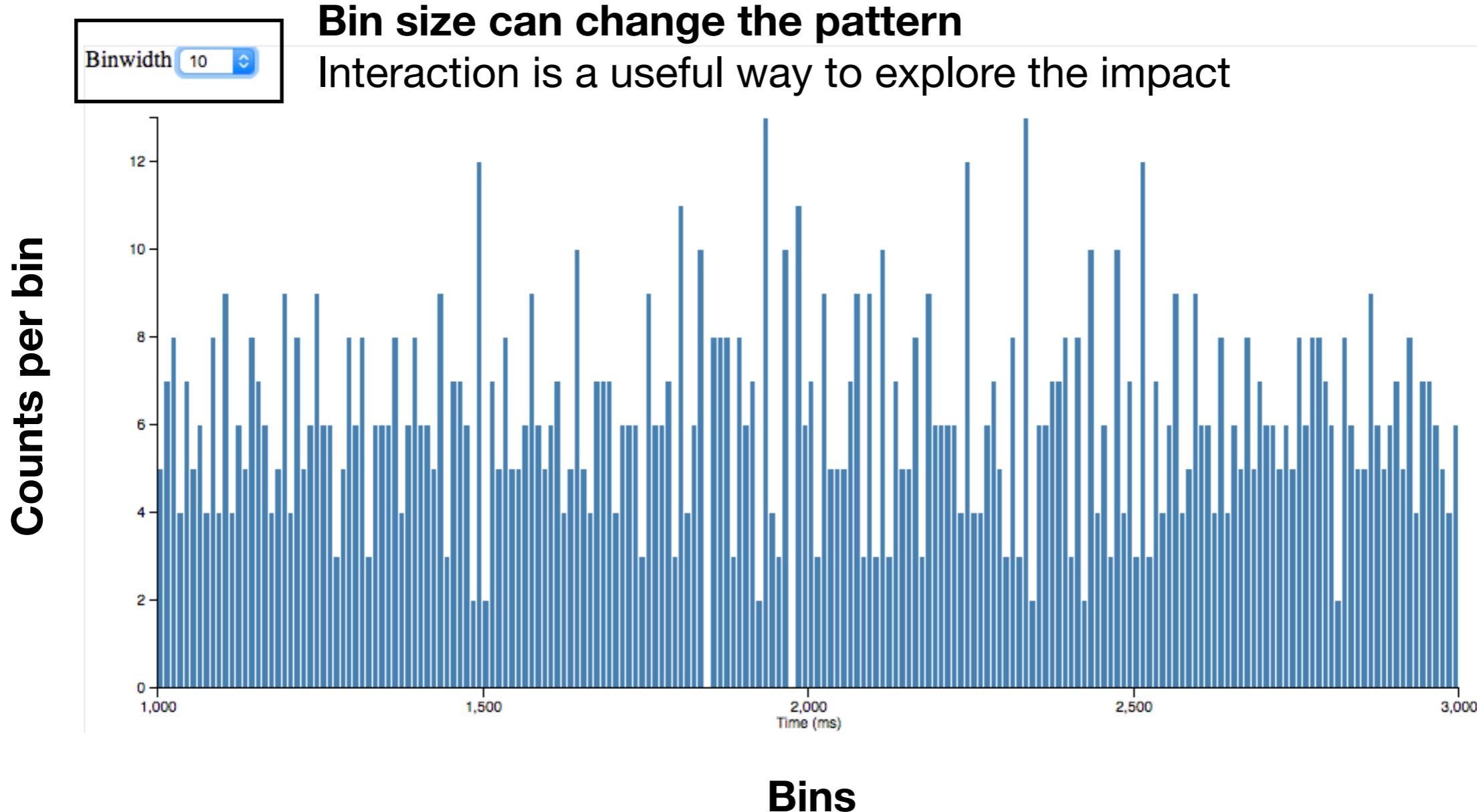
Unemployment rate by county, August 2016. Source: [Bureau of Labor Statistics](#).



Bins

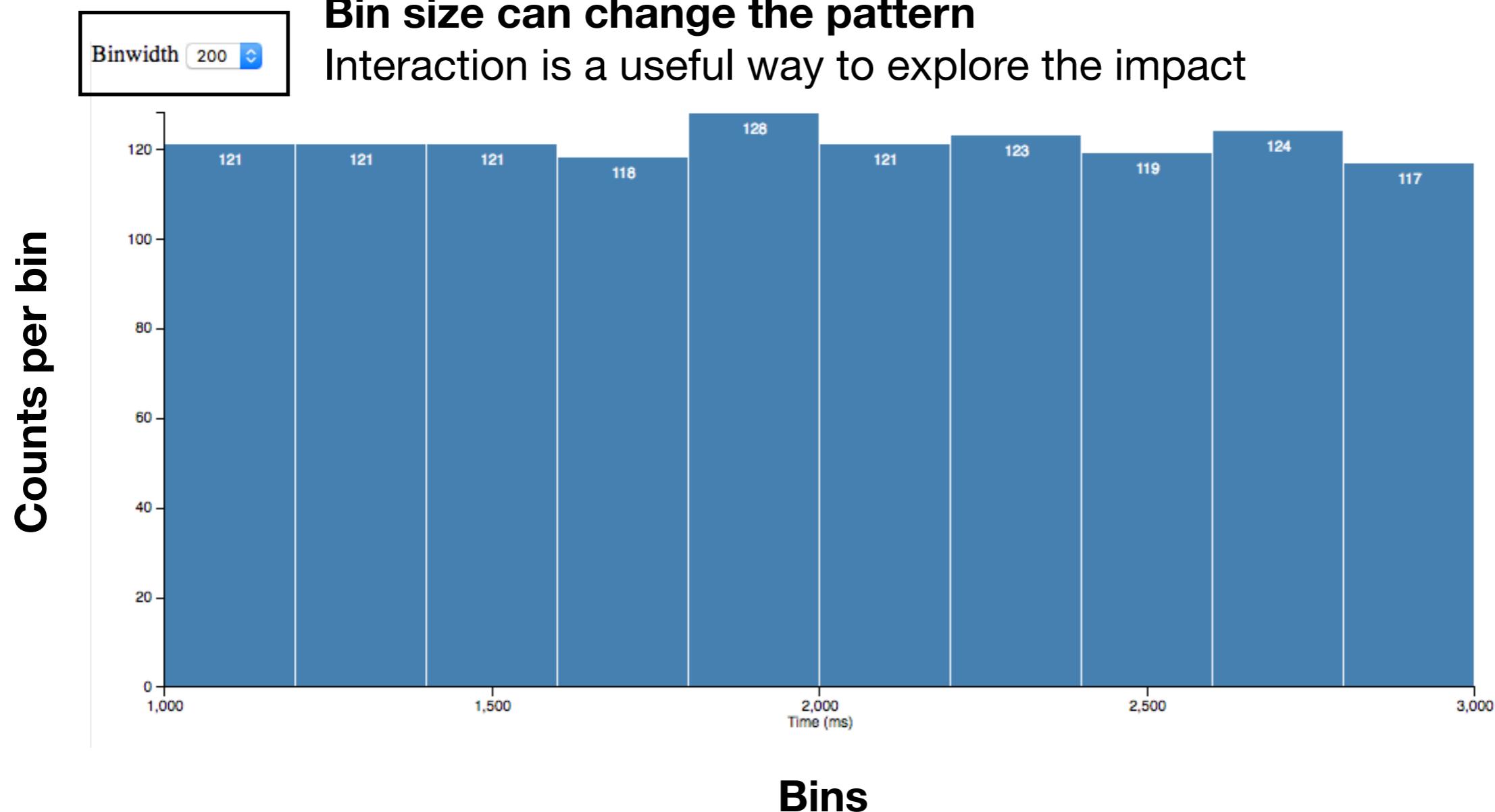
<https://beta.observablehq.com/@mbostock/d3-histogram>

Histogram | Visualizing a distribution



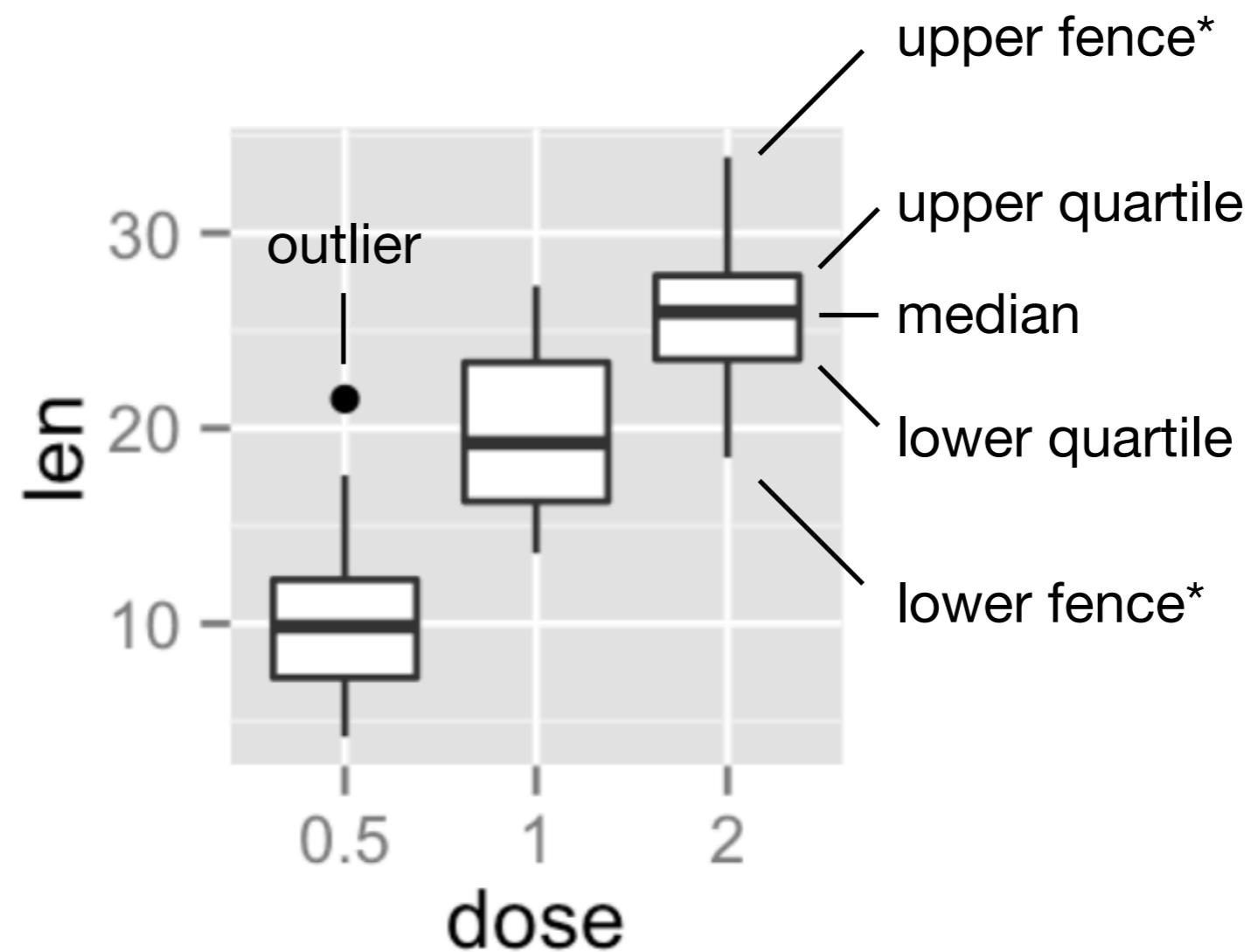
<http://bl.ocks.org/babsey/9b428f5f20e538c15b9e>

Histogram | Visualizing a distribution



<http://bl.ocks.org/babsey/9b428f5f20e538c15b9e>

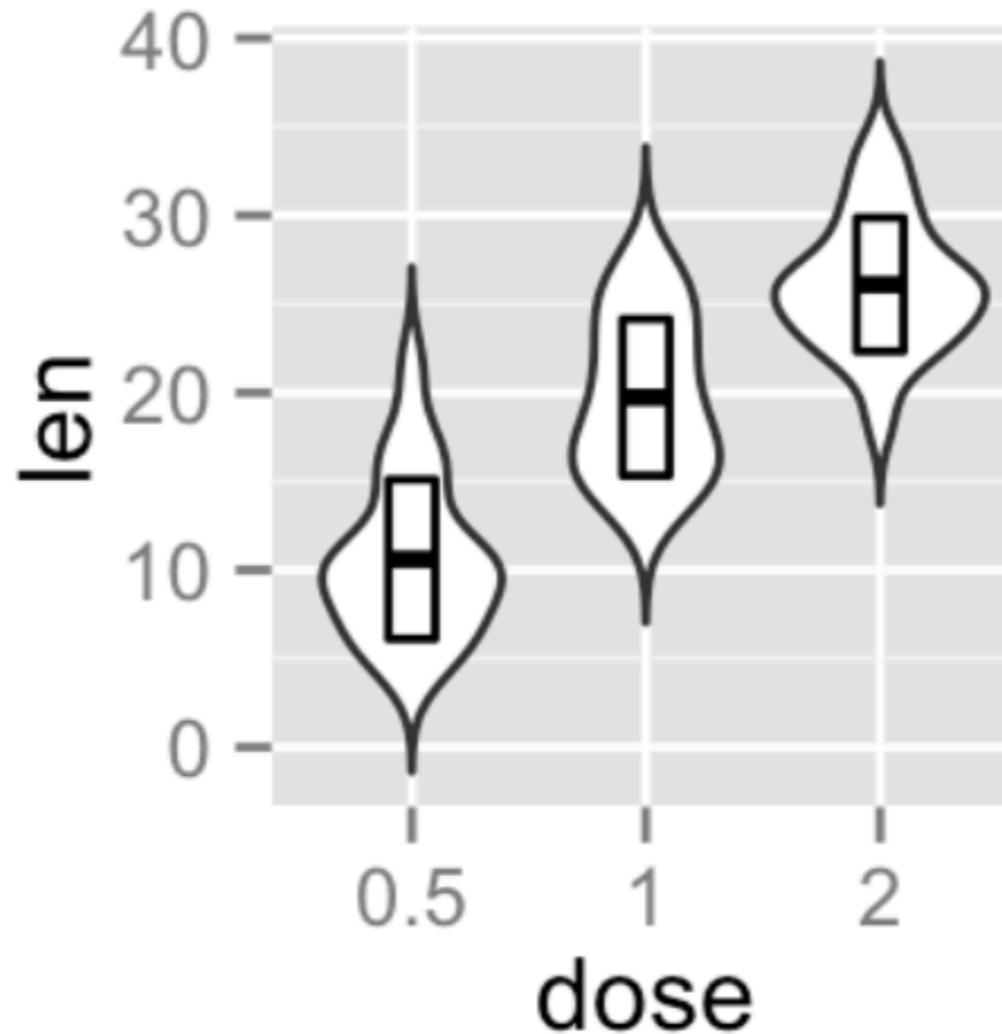
Box plots | Visualizing distributions



***fences** often extend from upper/lower quartile to the largest value no further than $+\/- 1.5 * \text{inter-quartile range}$

<http://www.sthda.com/english/wiki/ggplot2-violin-plot-quick-start-guide-...>

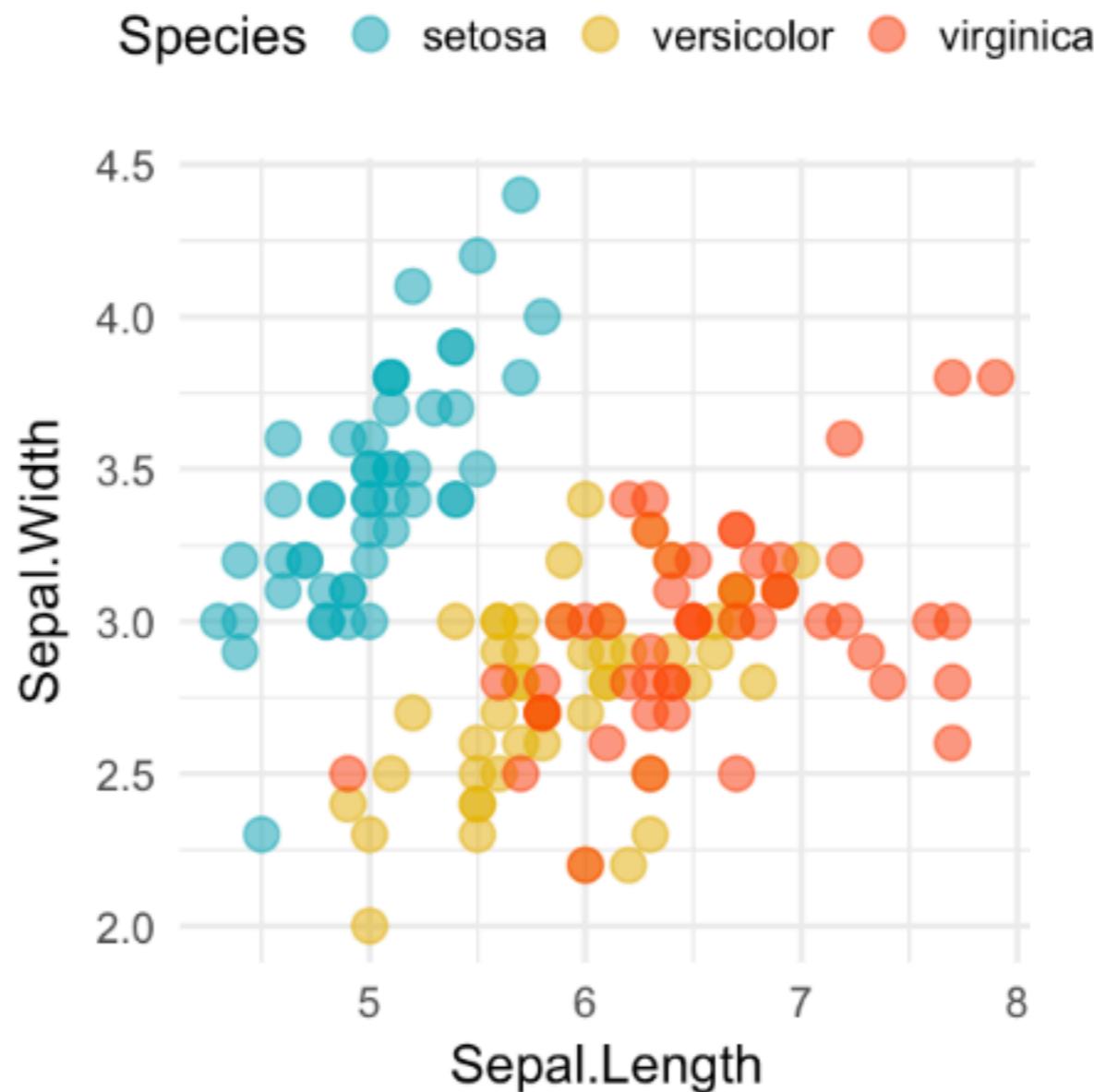
Violin plots | Visualizing distributions



Similar to a box plot but replaces the box with a mirrored density plot

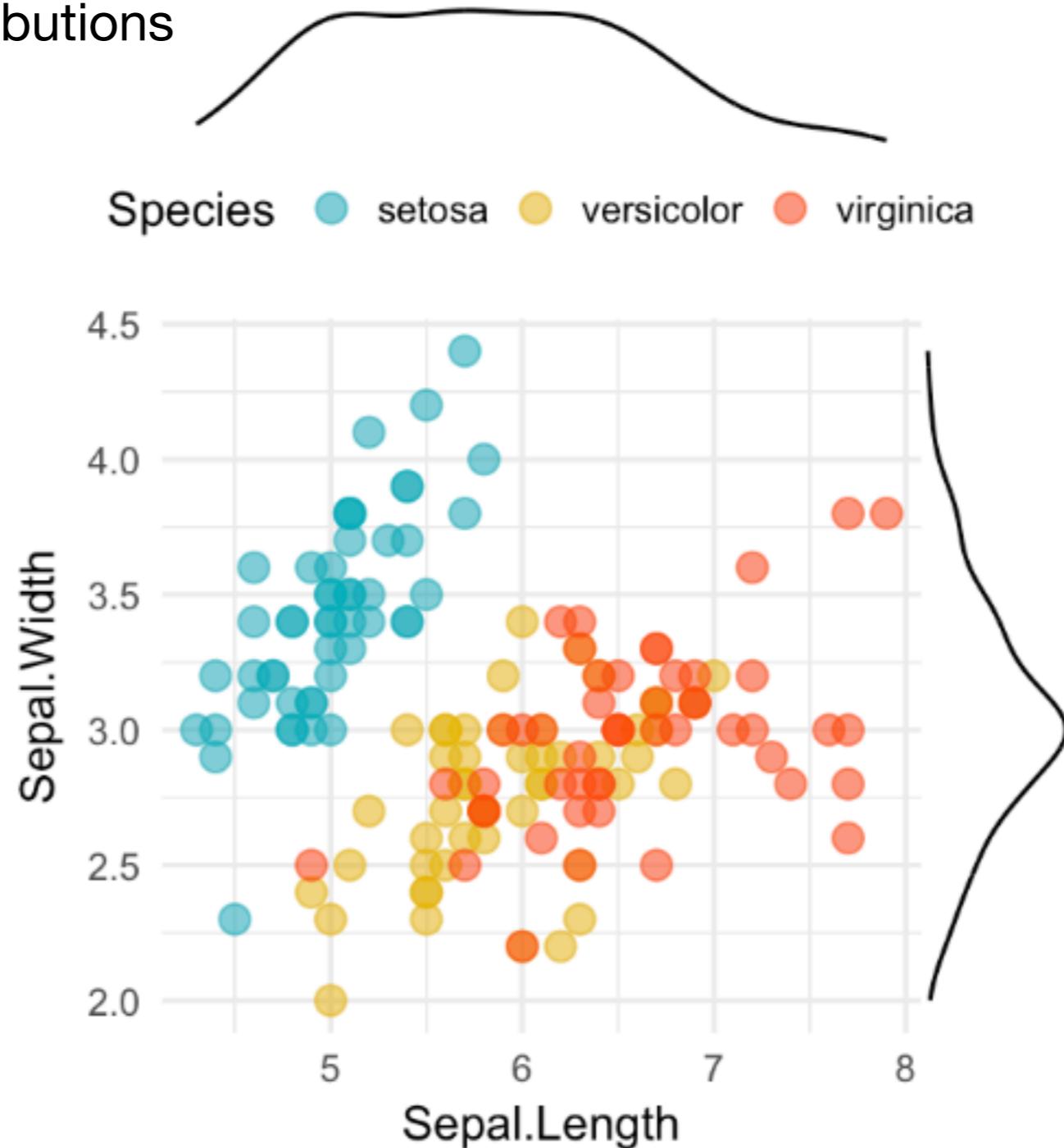
<http://www.sthda.com/english/wiki/ggplot2-violin-plot-quick-start-guide-...>

Scatter plots



Scatter plots

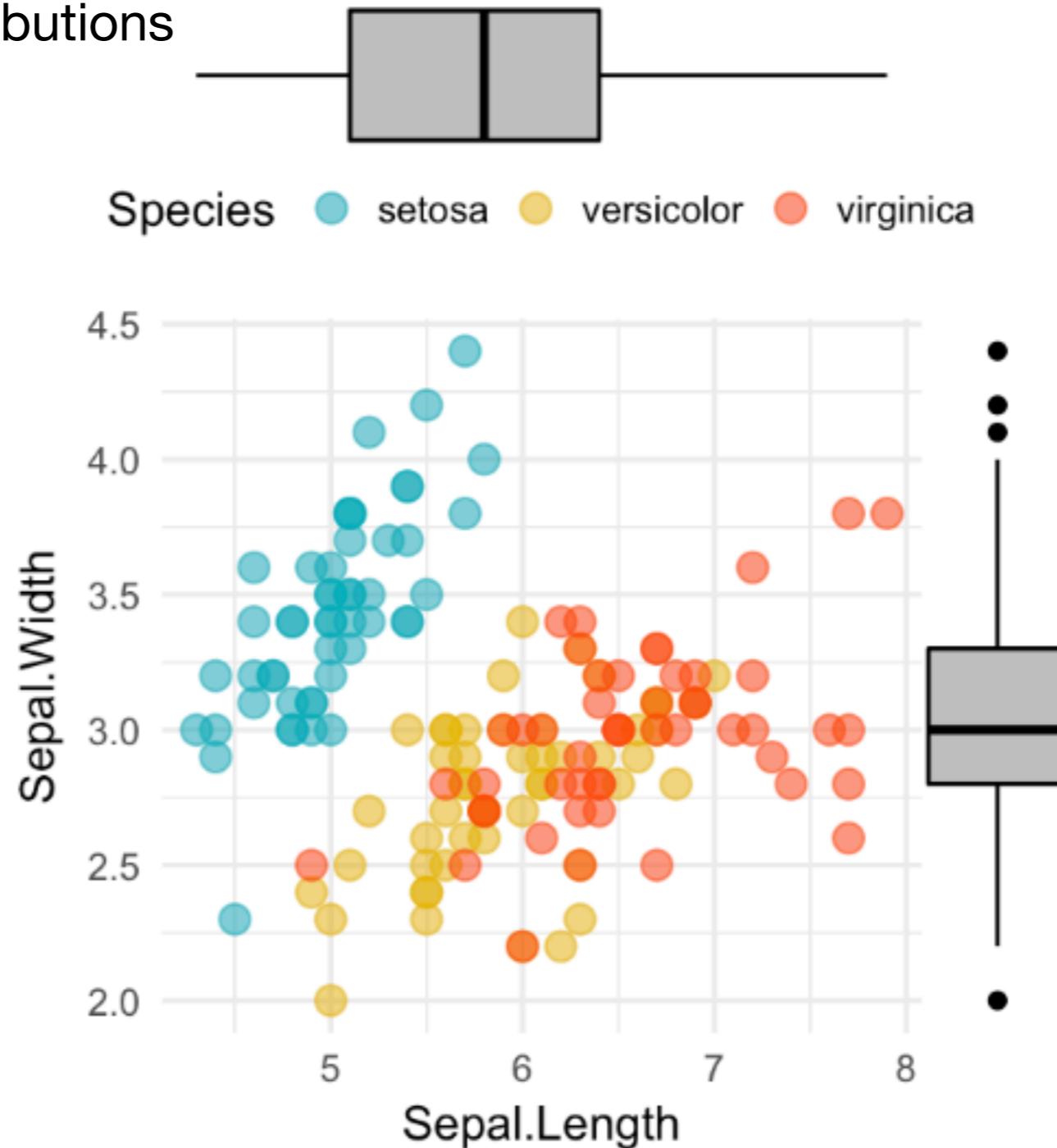
with marginal distributions



[http://www.sthda.com/english/articles/32-r-graphics-essentials/...](http://www.sthda.com/english/articles/32-r-graphics-essentials/)

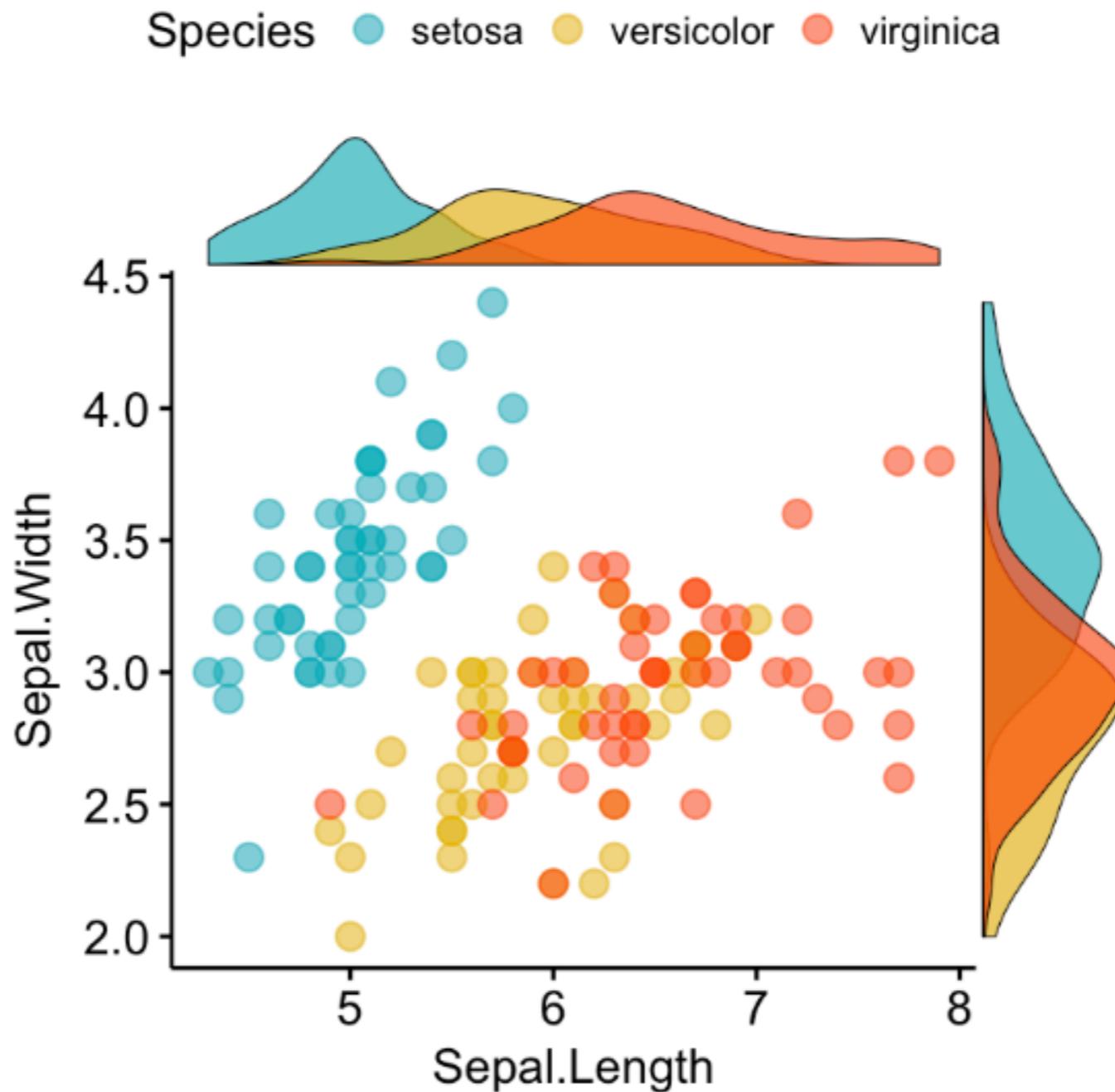
Scatter plots

with marginal distributions



Scatter plots

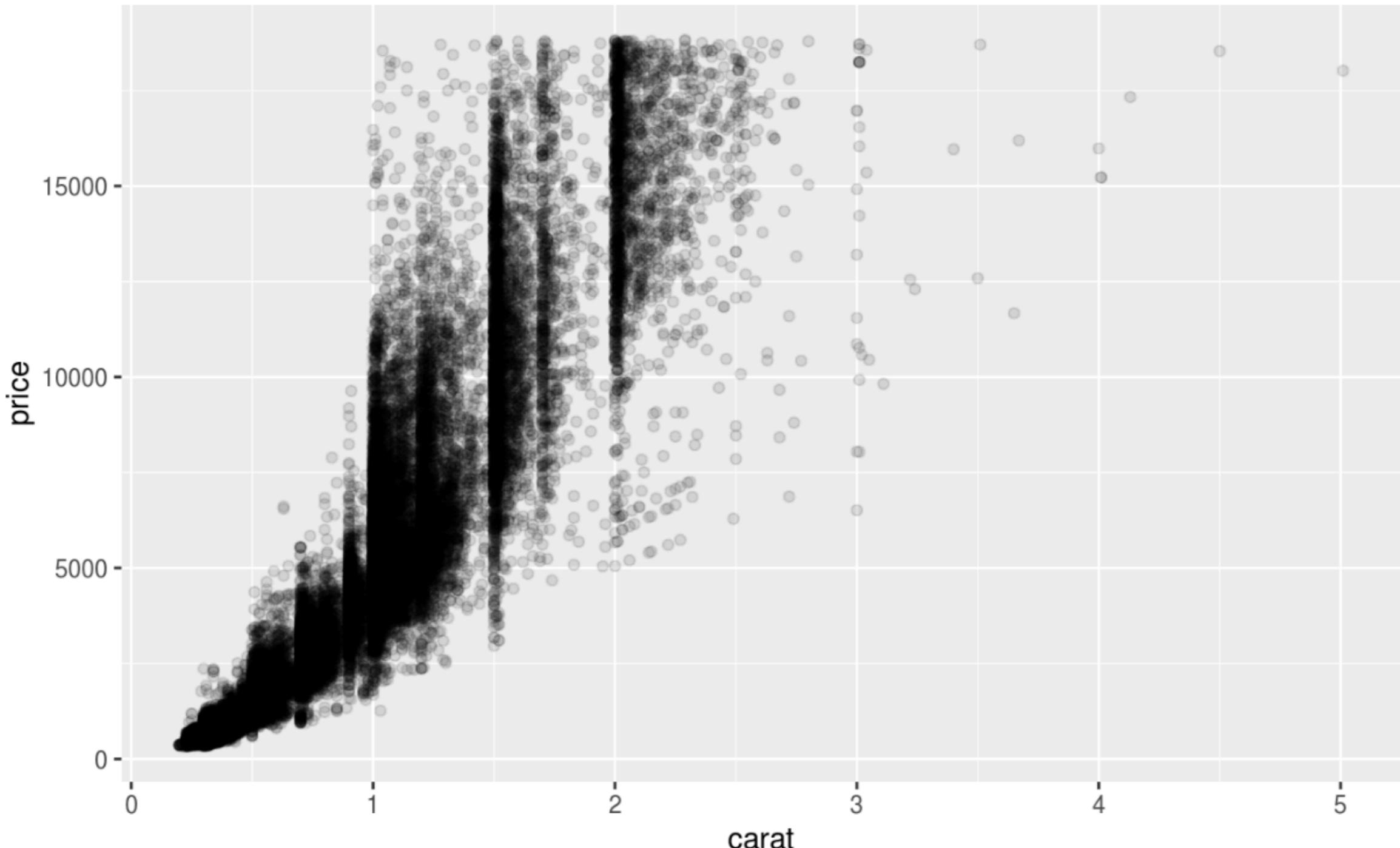
with marginal distributions



[http://www.sthda.com/english/articles/32-r-graphics-essentials/...](http://www.sthda.com/english/articles/32-r-graphics-essentials/)

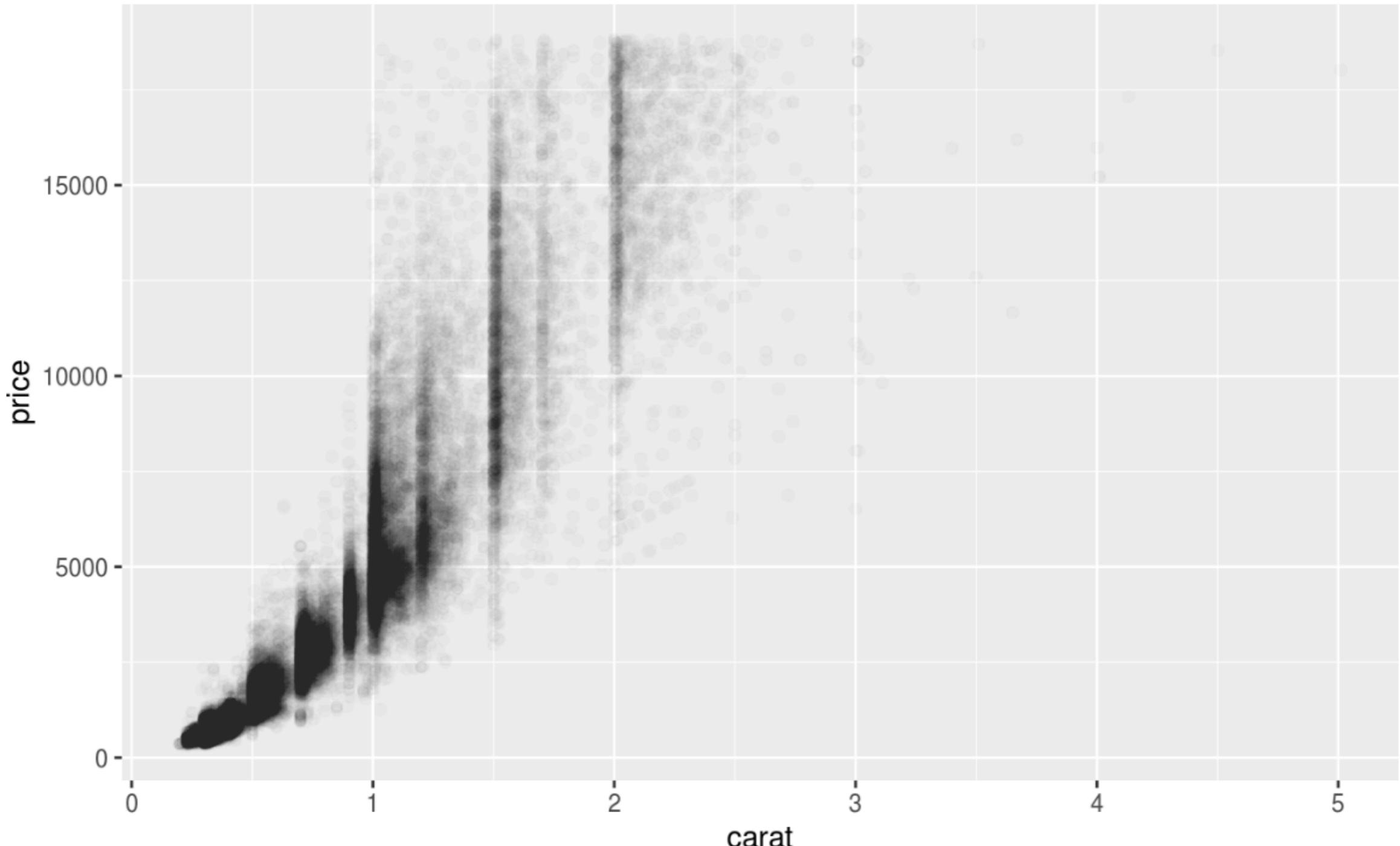
What if there are too many points to see them all?

Overplotting



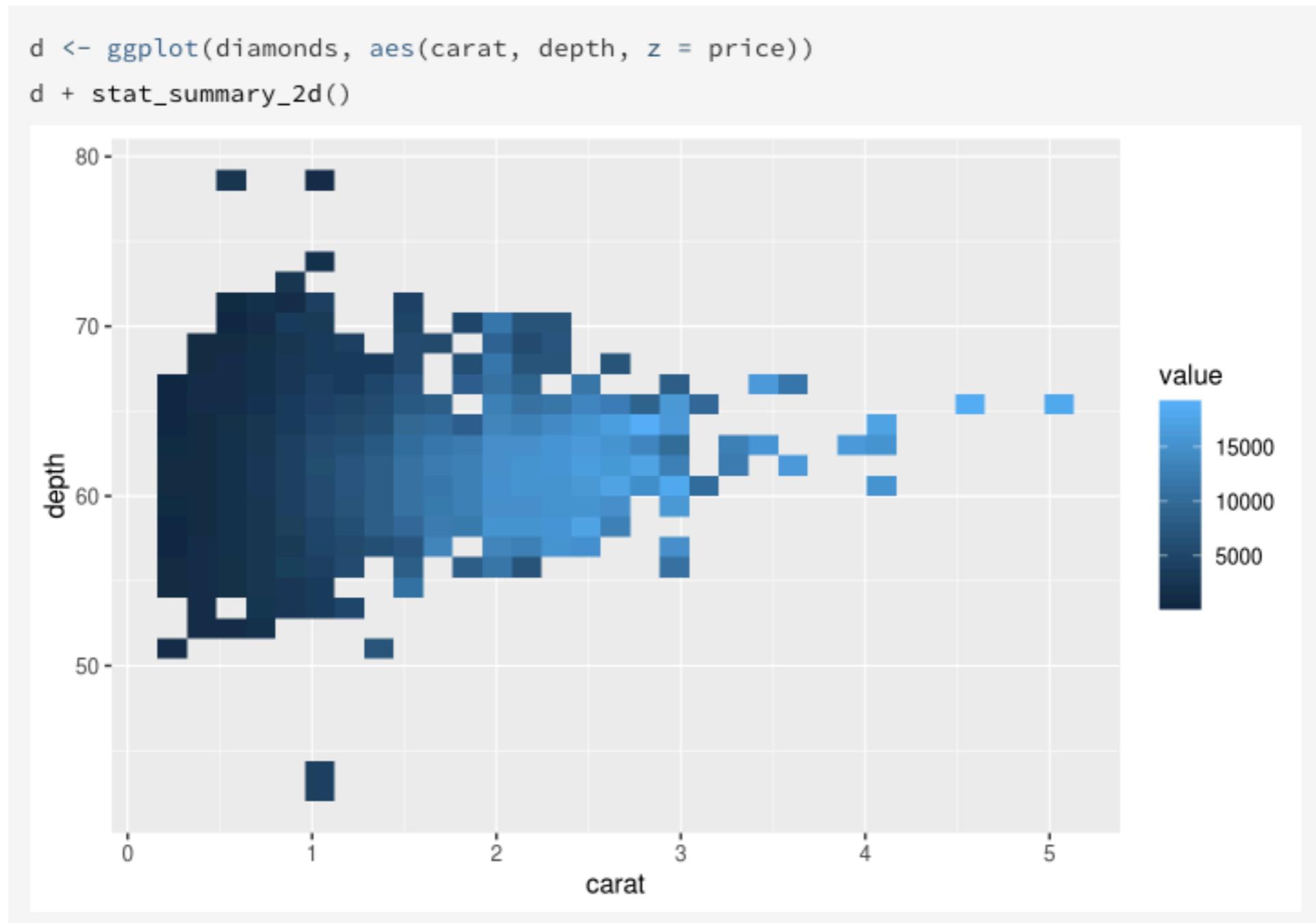
Overplotting

Changing the alpha value can help



Binning | Rectangular binning

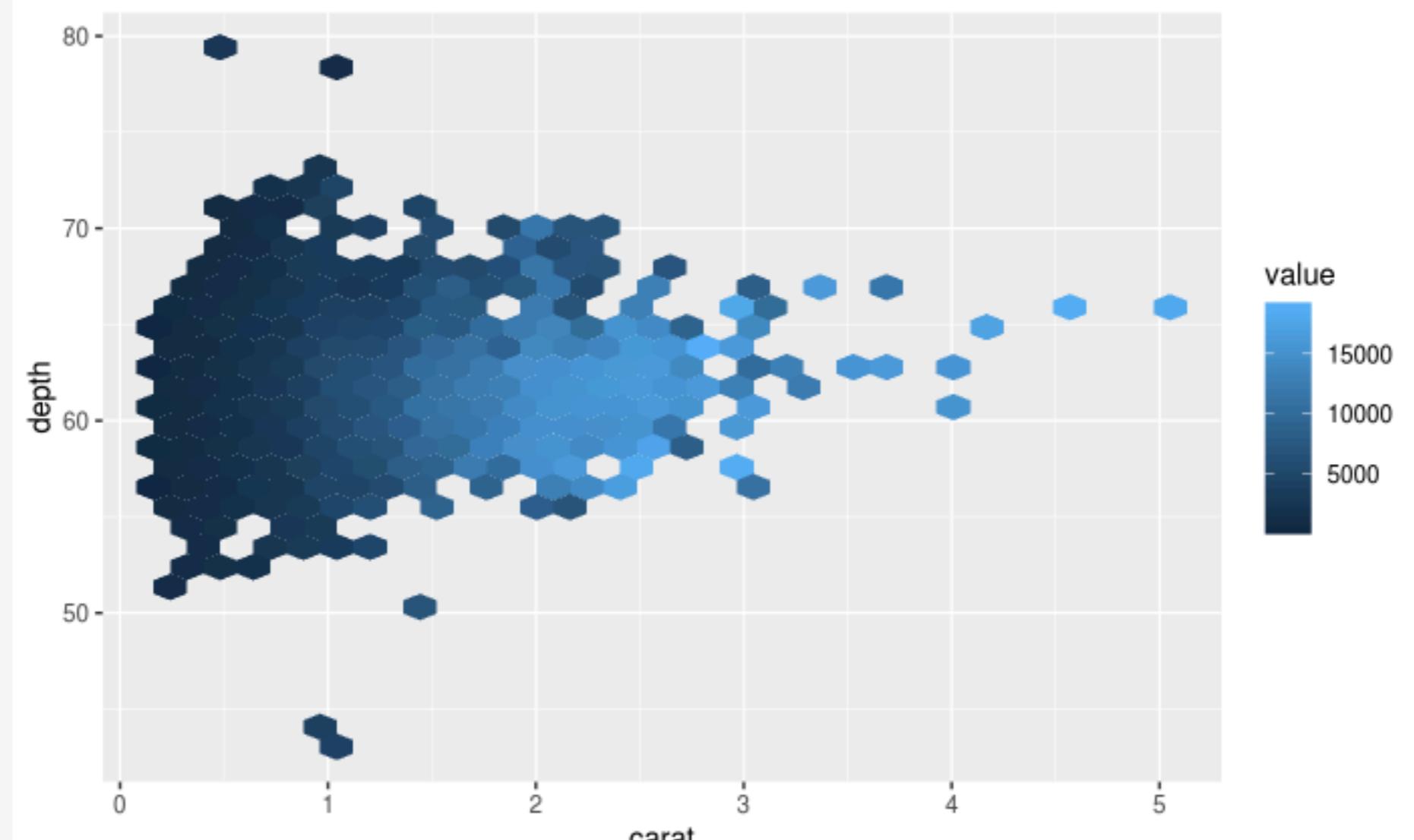
Great solution for large data sets



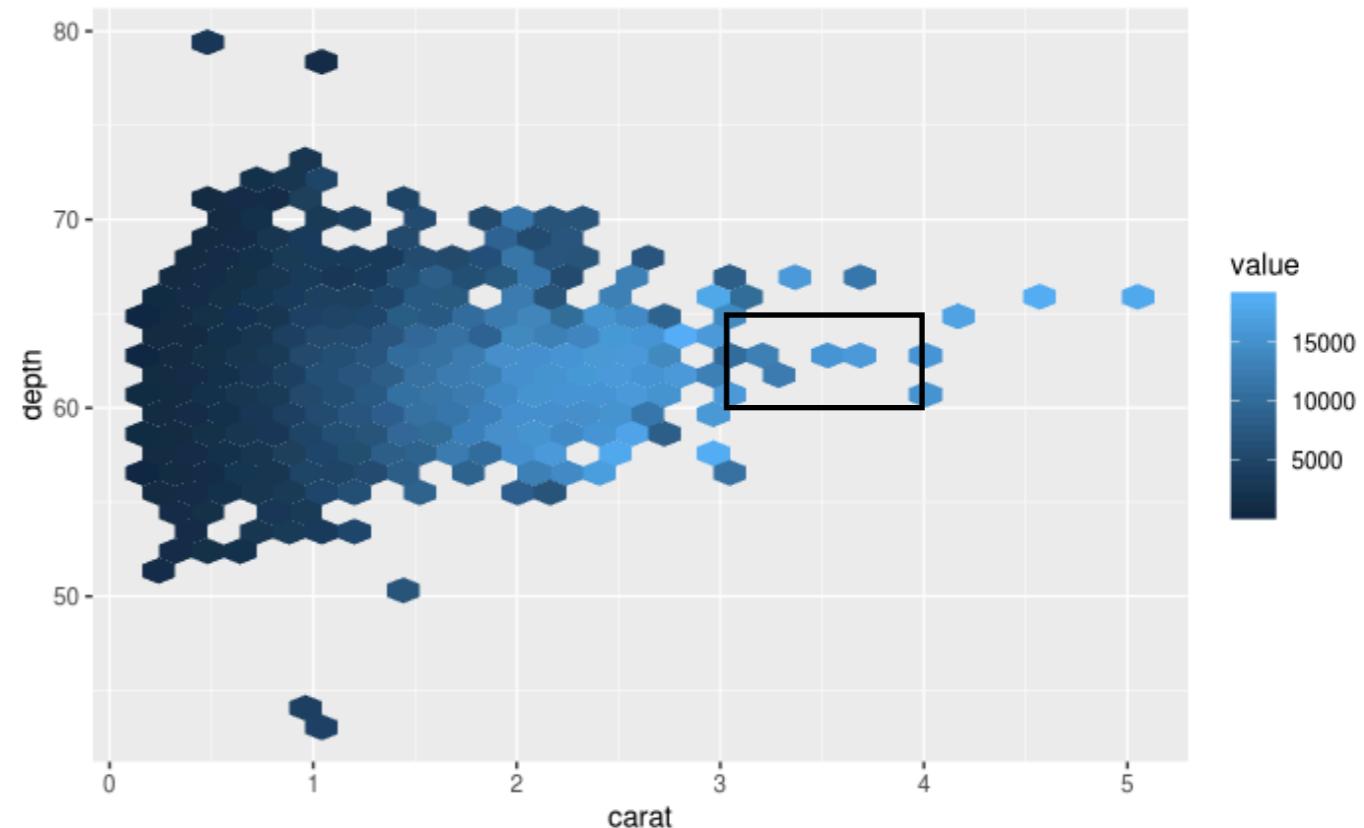
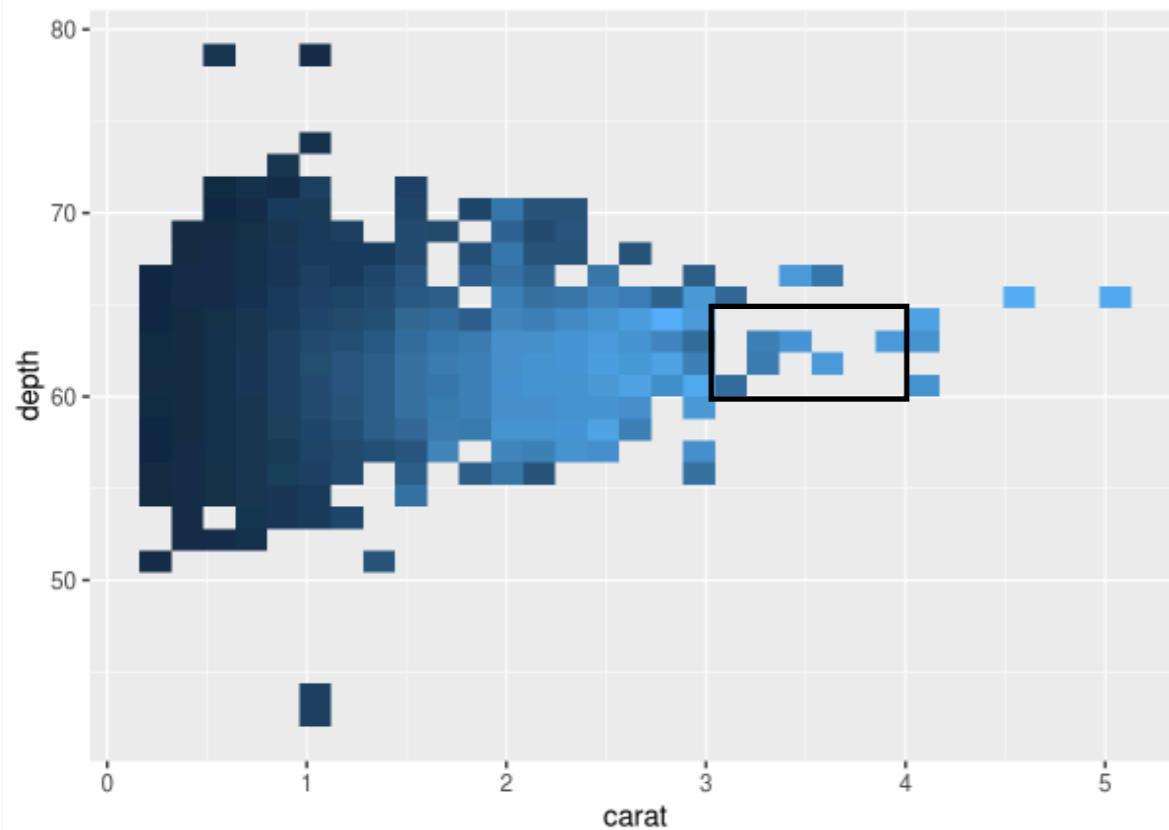
Binning | Hexagonal binning

Great solution for large data sets

```
if (requireNamespace("hexbin")) {  
  d + stat_summary_hex()  
}
```



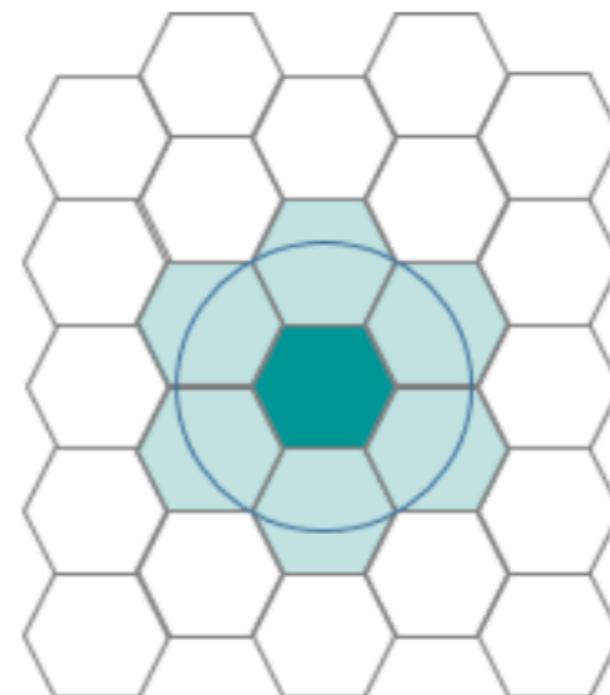
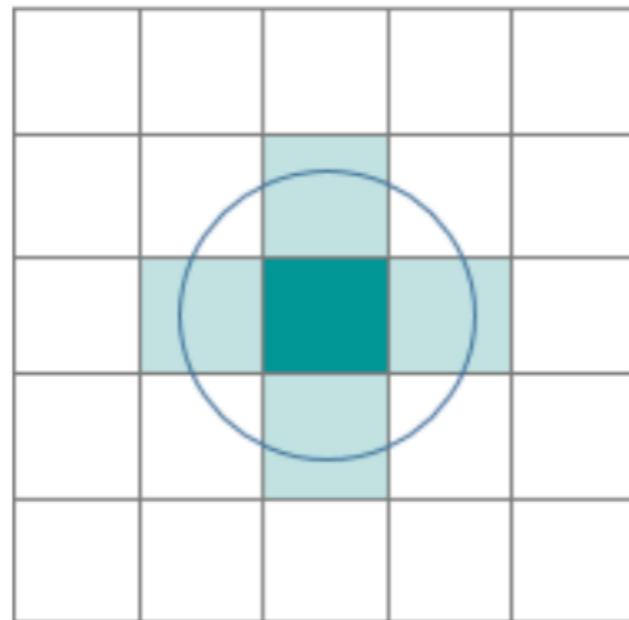
Rectangles versus Hexagons



- Rectangular binning can lead to artificially straight, parallel line patterns
- Circularity of a hexagon grid allows it to represent curved patterns in your data more naturally than square grids

<http://pro.arcgis.com/en/pro-app/tool-reference/spatial-statistics/h-whyhexagons.htm>

Rectangles versus Hexagons

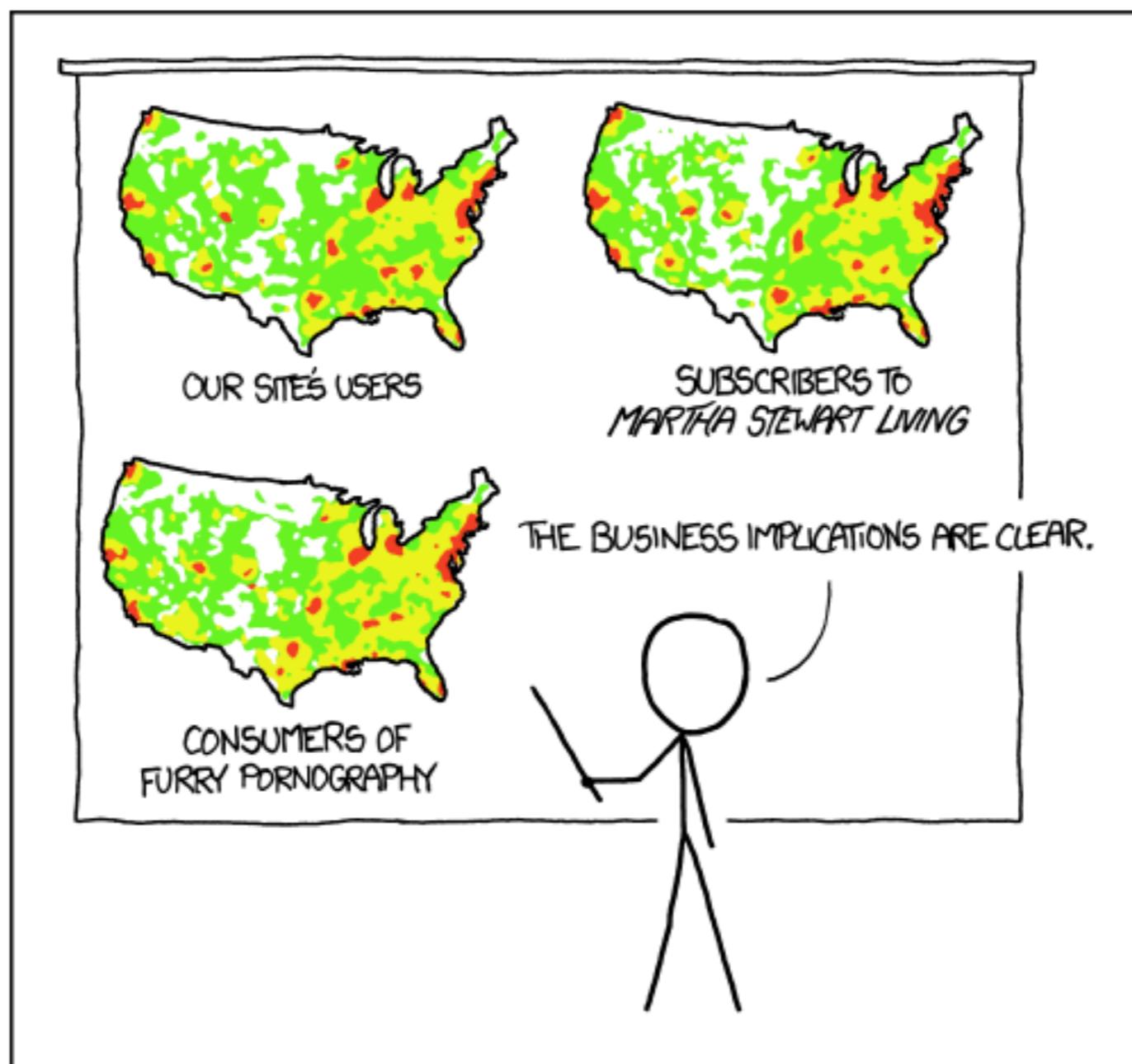


- Rectangular binning can lead to artificially straight, parallel line patterns
- Circularity of a hexagon grid allows it to represent curved patterns in your data more naturally than square grids
- Nice property that the distance between centroids is the same in all six directions with hexagons; this will impact calculations of number of neighbours within a given distance

<http://pro.arcgis.com/en/pro-app/tool-reference/spatial-statistics/h-whyhexagons.htm>

Aggregation on a geographic map

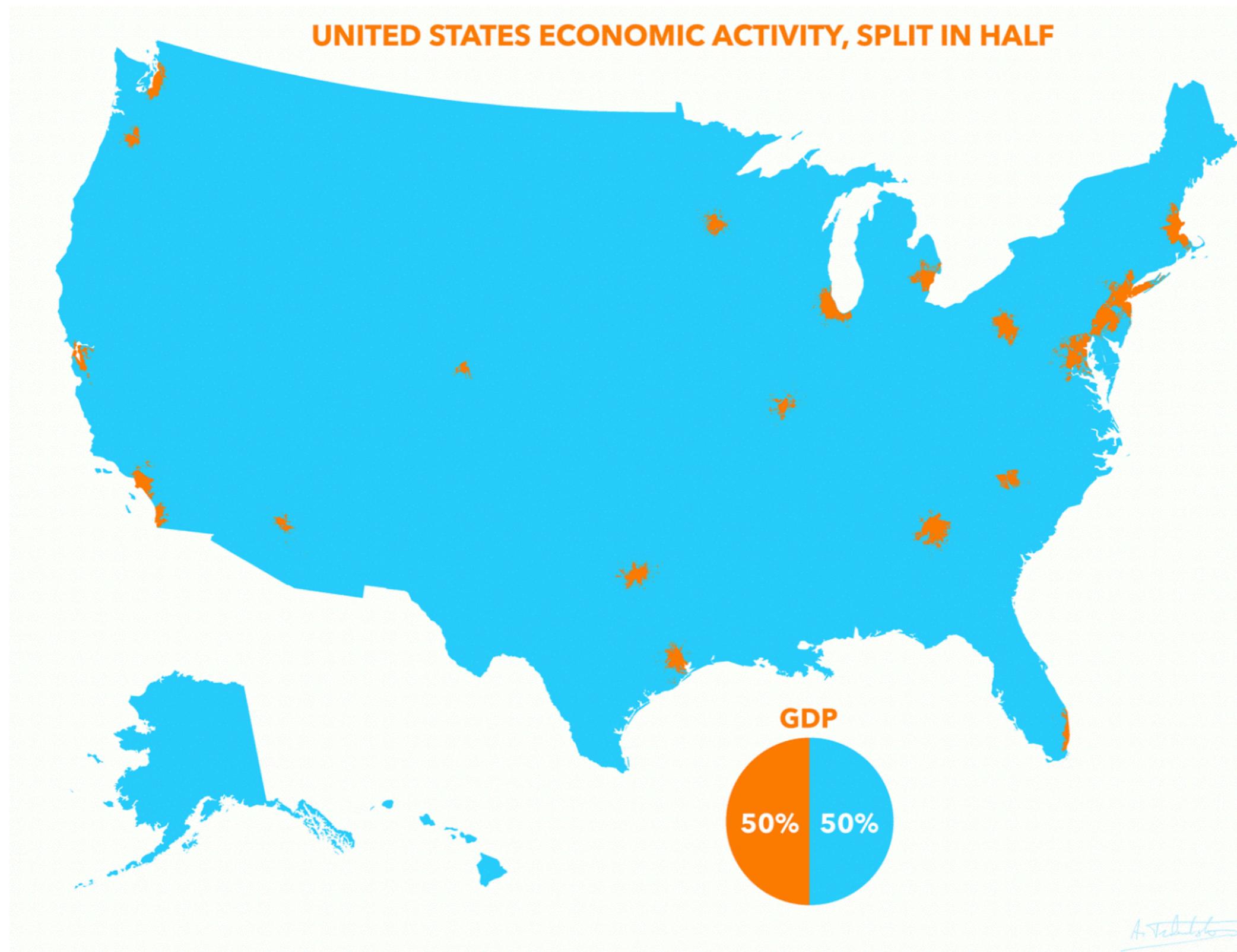
Choropleth maps



PET PEEVE #208:
GEOGRAPHIC PROFILE MAPS WHICH ARE
BASICALLY JUST POPULATION MAPS

<https://xkcd.com/1138/>

The incredible map that shows that half of the U.S. population produces half of the GDP



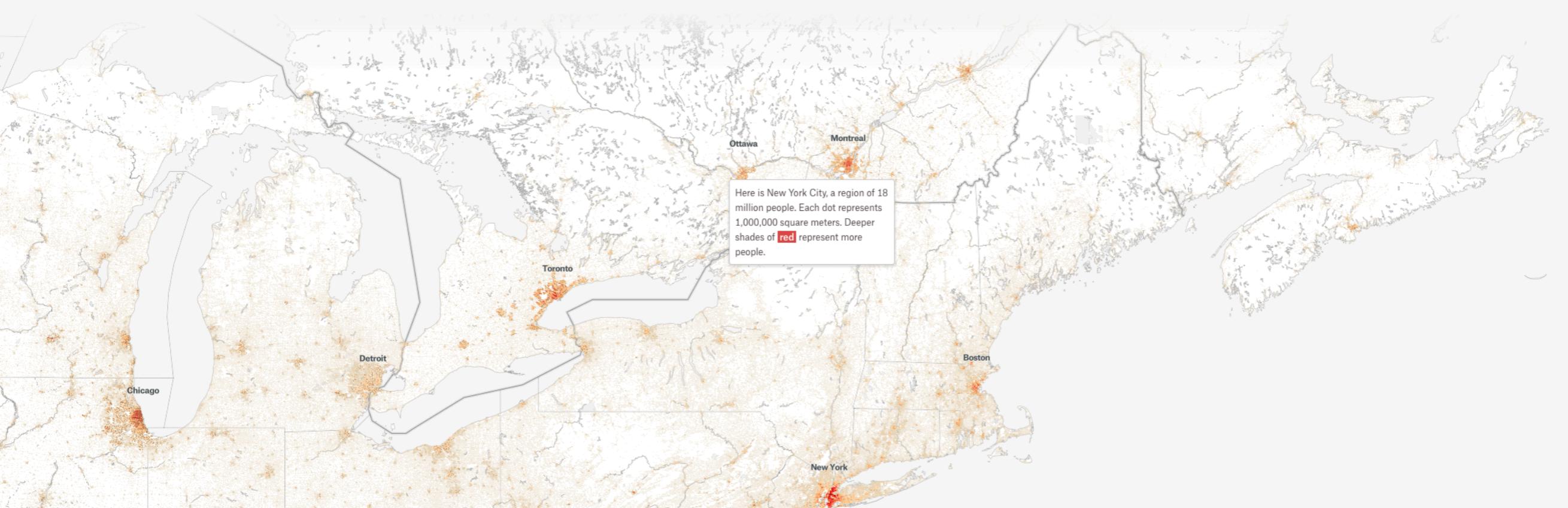
See discussion and links on Alberto Cairo's blog | the Functional Art
<http://www.thefunctionalart.com/2014/02/the-incredible-gdp-map-that-shows-that.html>

What does the world's population distribution look like?

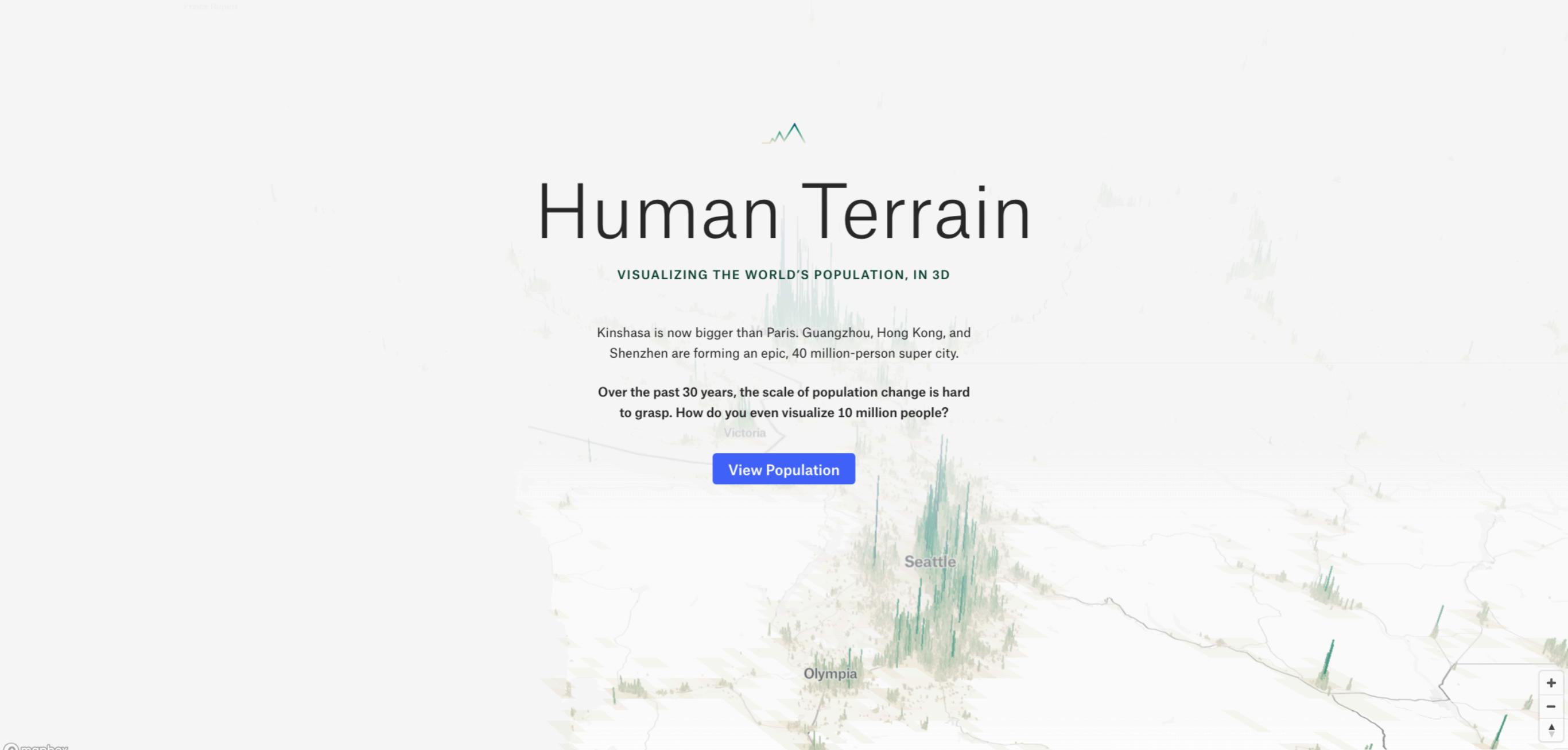
Population Mountains

By [Matt Daniels](#)

This is a story about how to perceive the population of cities.



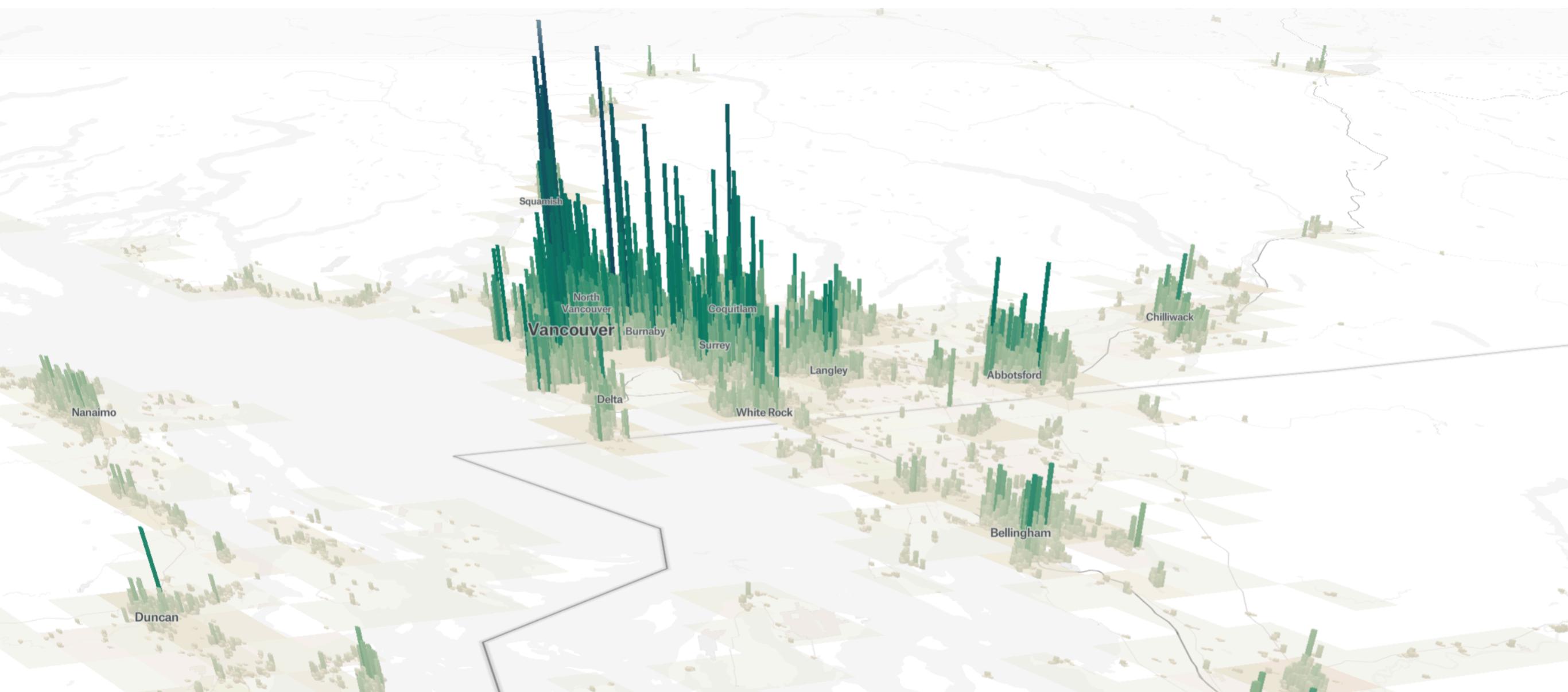
<https://pudding.cool/2018/12/3d-cities-story/>

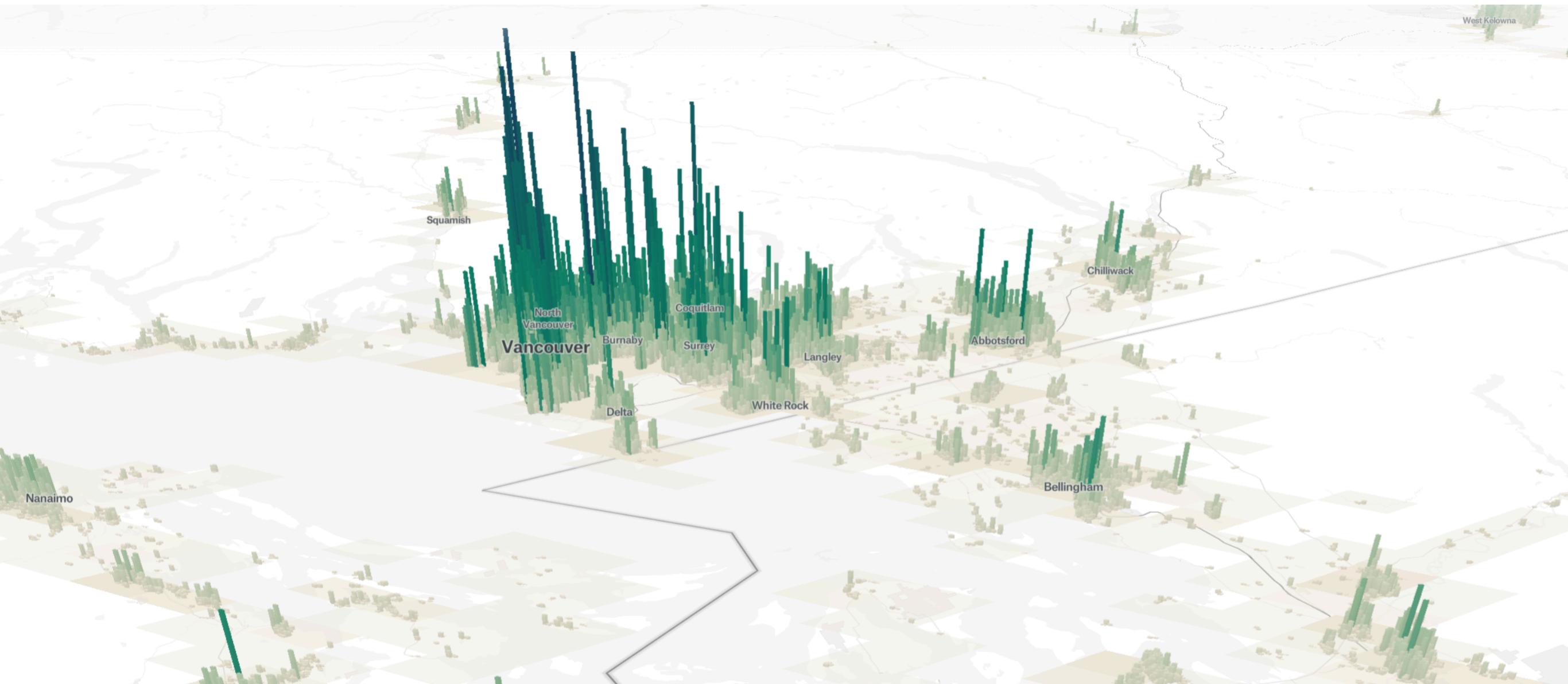


What works well?

What could be improved?

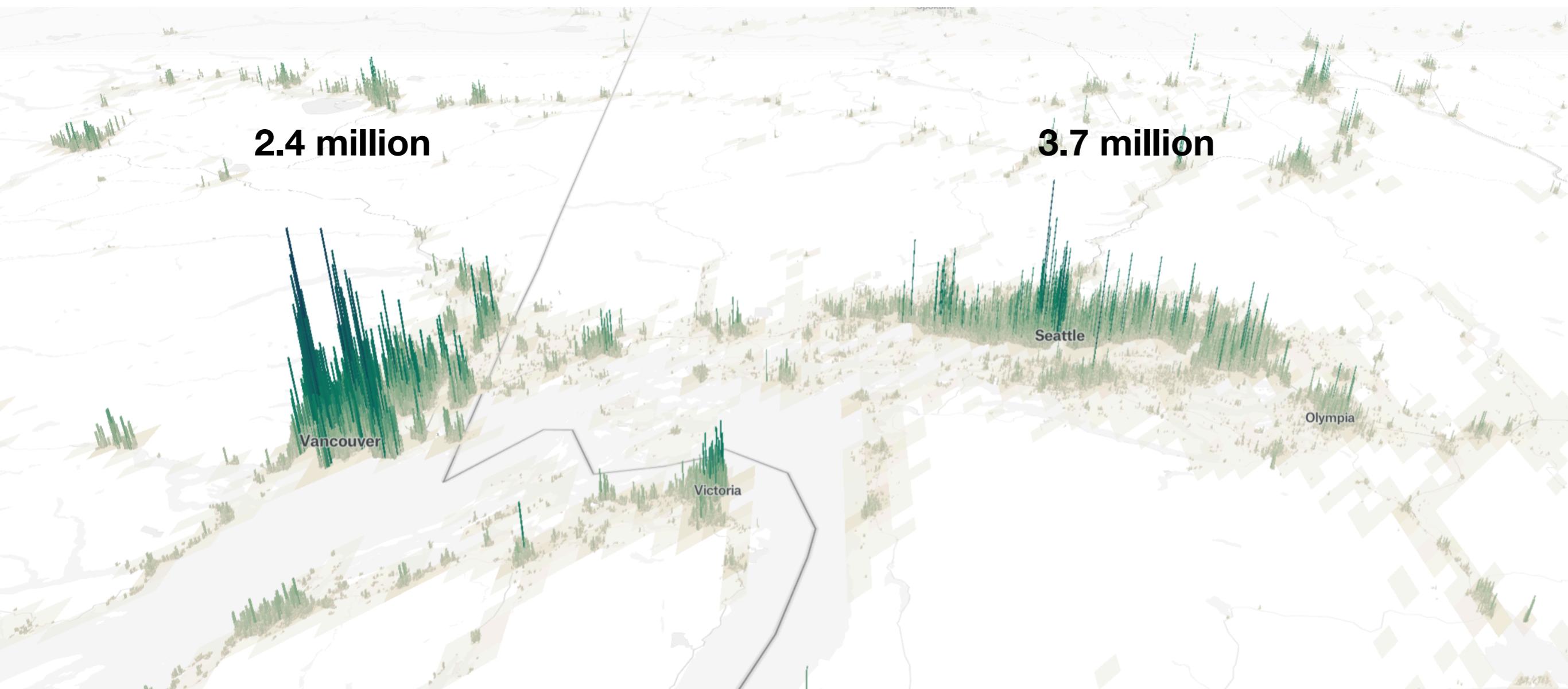
https://pudding.cool/2018/10/city_3d/





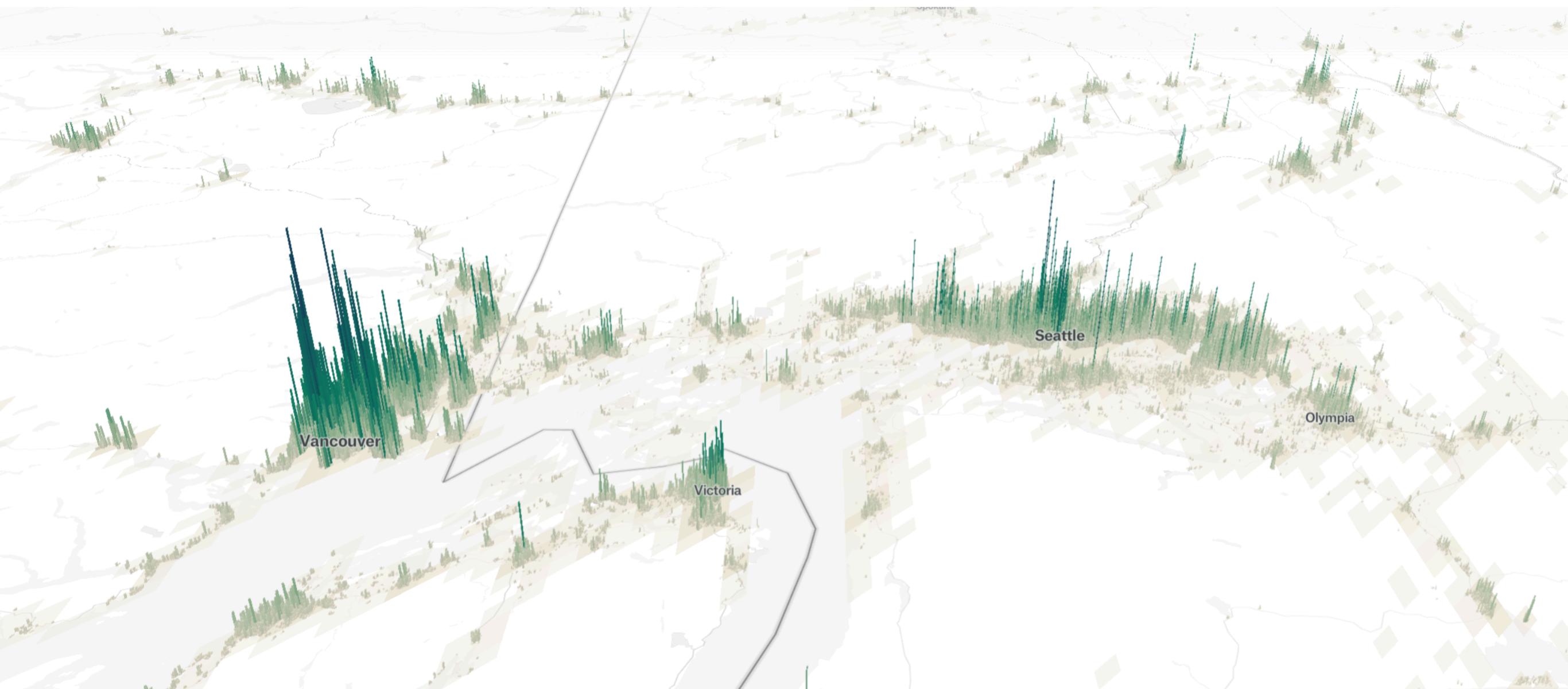
Occlusion problem

- tall bars block shorter ones and it's hard to tell where the labels map



Which city is bigger?

- Colour + bar height emphasize tall bars (areas with dense population)



Other discussion points

- A scale would be helpful!

Zone effects | Where you draw the line matters

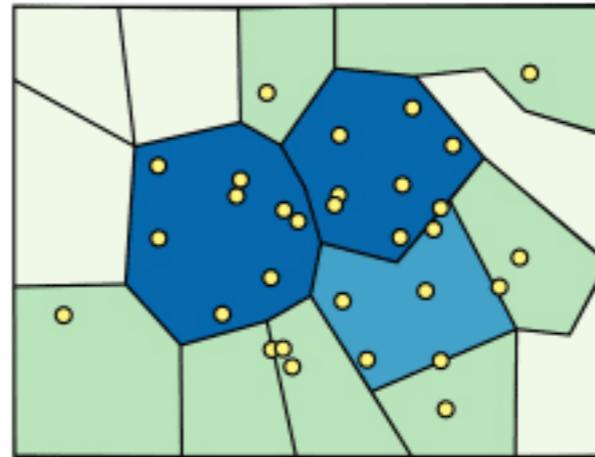
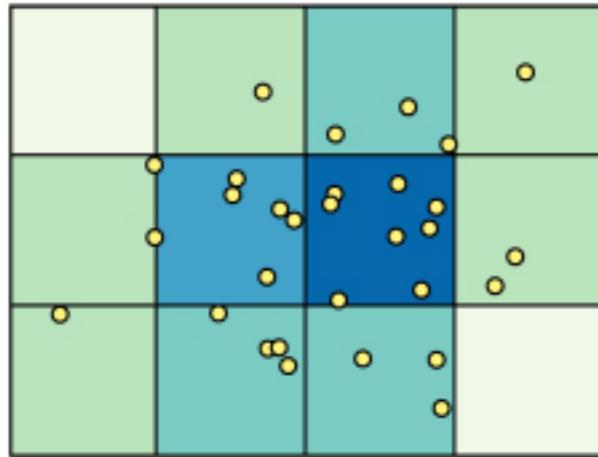
“...the areal units (zonal objects) used in many geographical studies are arbitrary, modifiable, and subject to the whims and fancies of whoever is doing, or did, the aggregating.”

- Openshaw (1984)

Modifiable Areal Unit Problem

<https://blog.cartographica.com/blog/2011/5/19/the-modifiable-areal-unit-problem-in-gis.html>

Zone effects | Where you draw the line matters

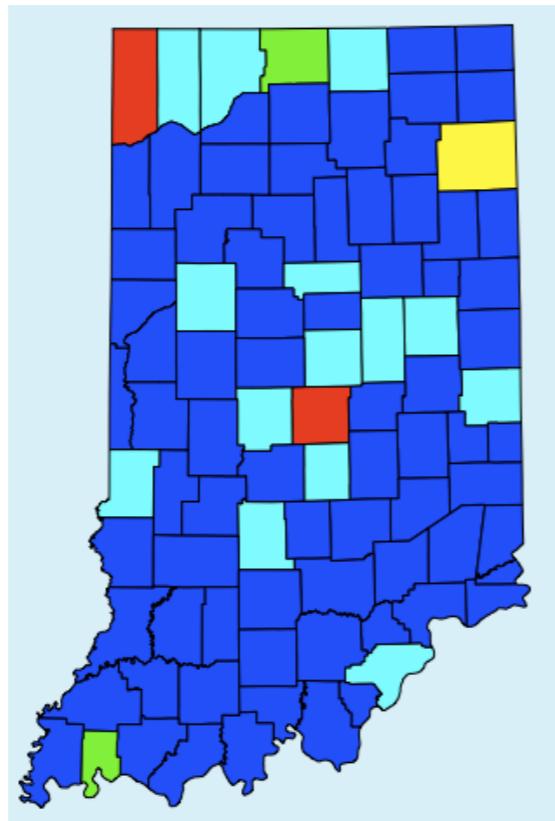


Different boundaries give different results

- gerrymandering (manipulating voting district boundaries) is one example

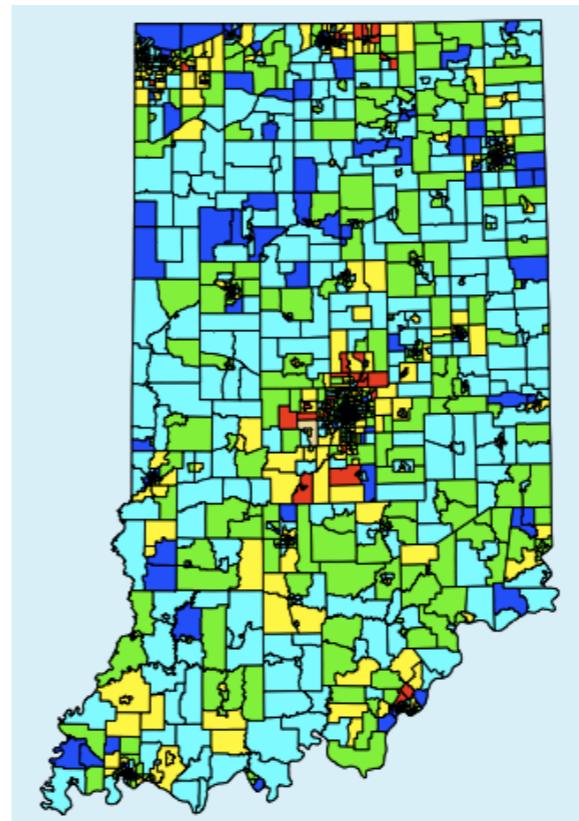
Scale effects | Bin size matters

Household counts across the state of Indiana



Counties

Larger bins, less variation



Census Tracts

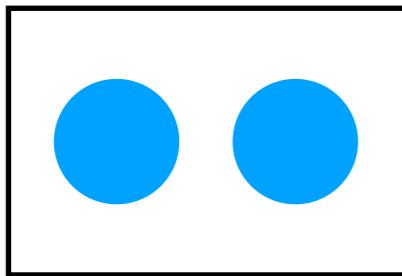
Smaller bins, more variation

Modifiable Areal Unit Problem

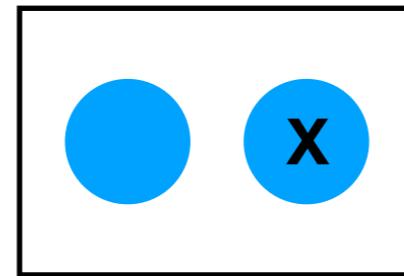
<https://blog.cartographica.com/blog/2011/5/19/the-modifiable-areal-unit-problem-in-gis.html>

When populations are low, variation tends to be high

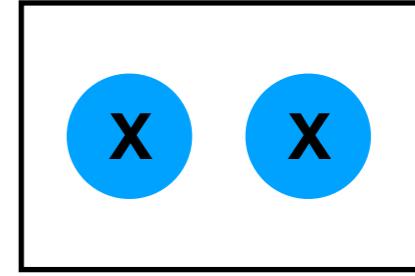
Imagine a population of 2 people



0% crime rate
Lowest in the country!



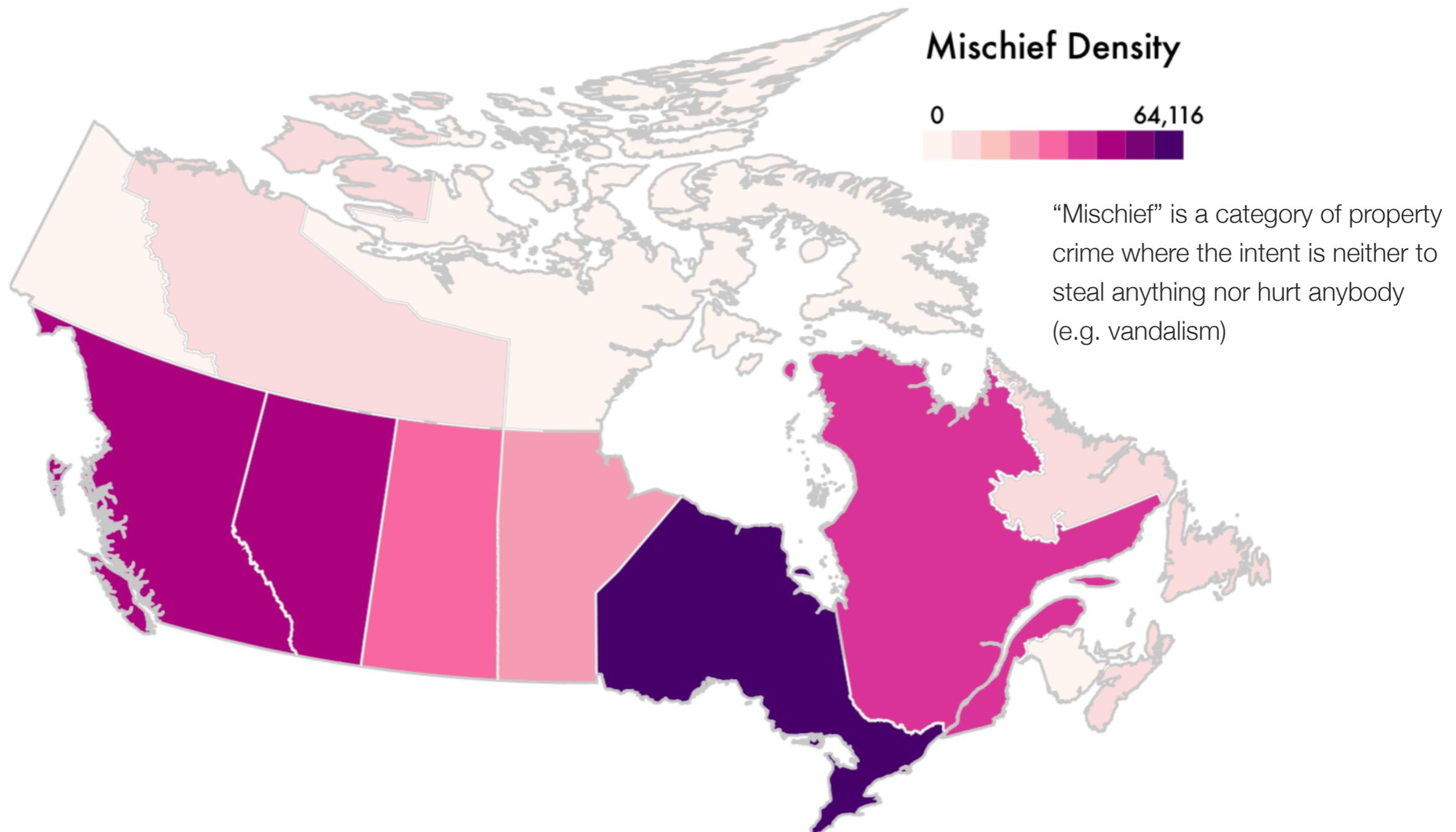
50% crime rate
Highest in the country!



100% crime rate
What kind of a place is this!

<https://medium.com/@uwdata/surprise-maps-showing-the-unexpected-e92b67398865>

Choropleth map | Event counts

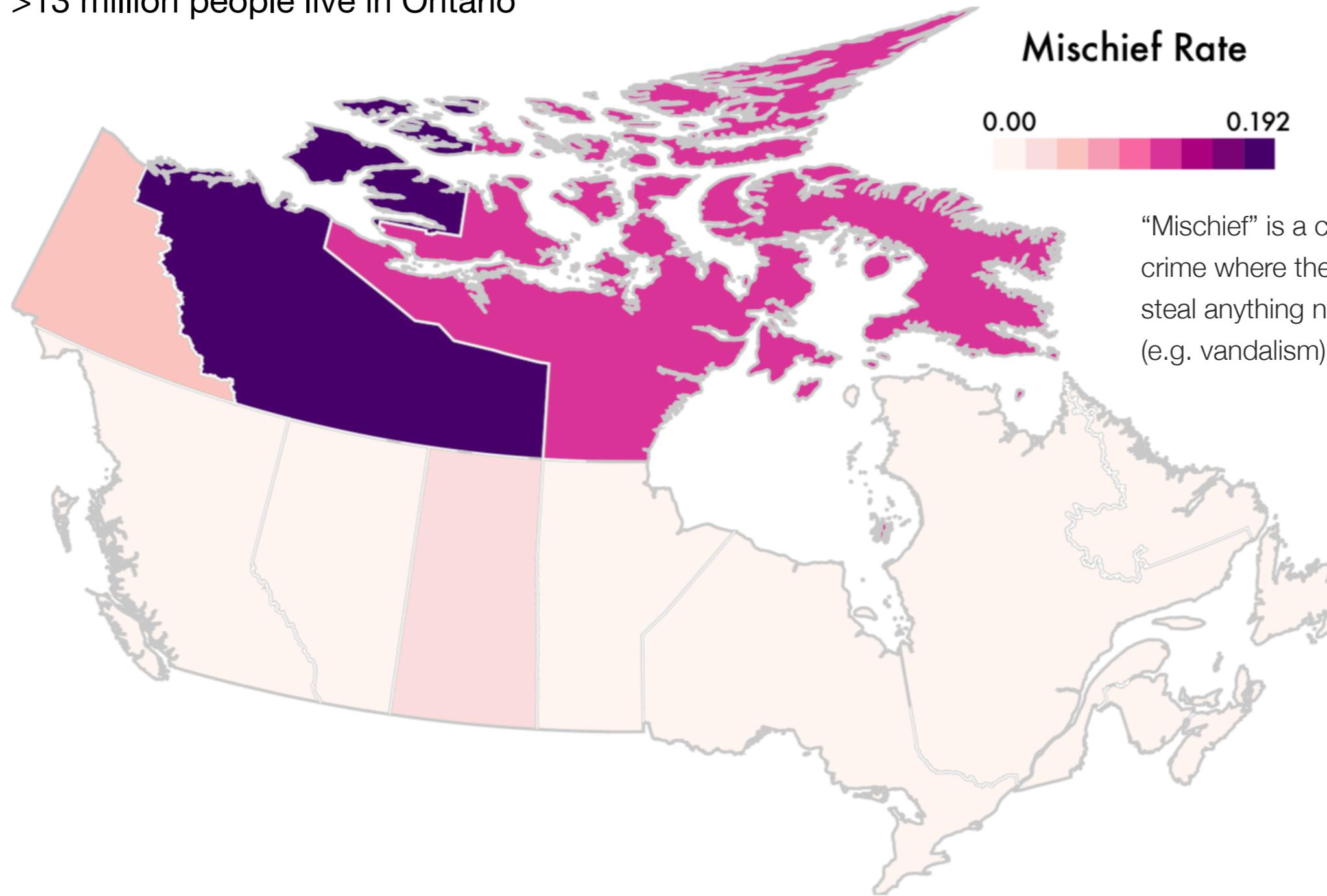


<https://medium.com/@uwdata/surprise-maps-showing-the-unexpected-e92b67398865>

Choropleth map | Per capita rate of crime

< 44,000 people live in the Northwest Territories

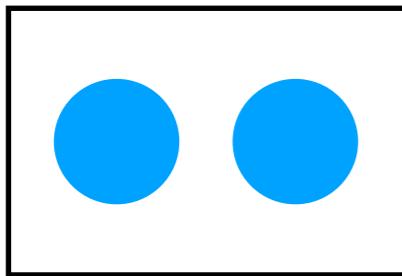
>13 million people live in Ontario



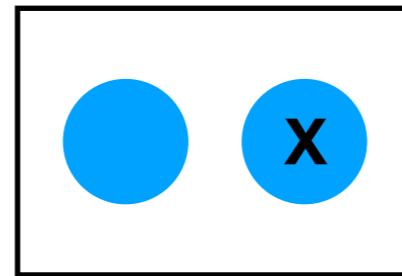
<https://medium.com/@uwdata/surprise-maps-showing-the-unexpected-e92b67398865>

When populations are low, variation tends to be high

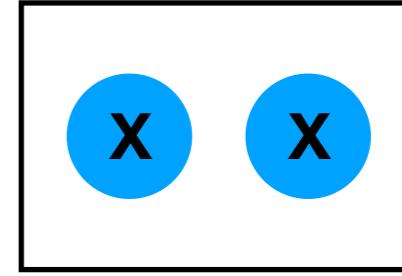
Imagine a population of 2 people



0% crime rate
Lowest in the country!



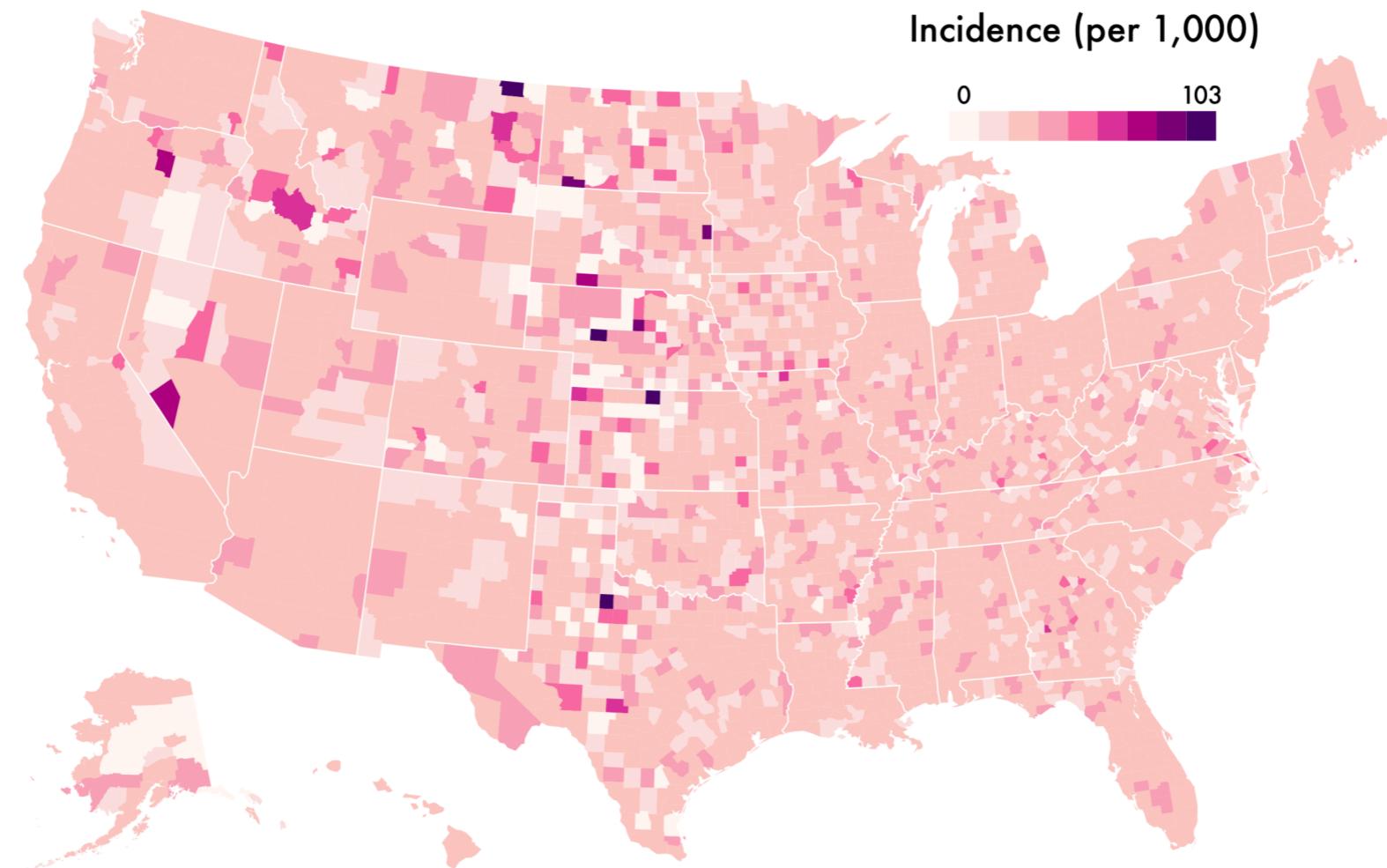
50% crime rate
Highest in the country!



100% crime rate
What kind of a place is this!

None of these cases offer much evidence that this two-person province is really the safest or most dangerous place to live

When populations are low, variation tends to be high



Random incidence

- Flip a coin for every person in the US
(0.1% chance of success for each citizen, regardless of location)
- The apparent geographic patterns are an artifact of the high variance in counties with low population

<https://medium.com/@uwdata/surprise-maps-showing-the-unexpected-e92b67398865>

Surprise map

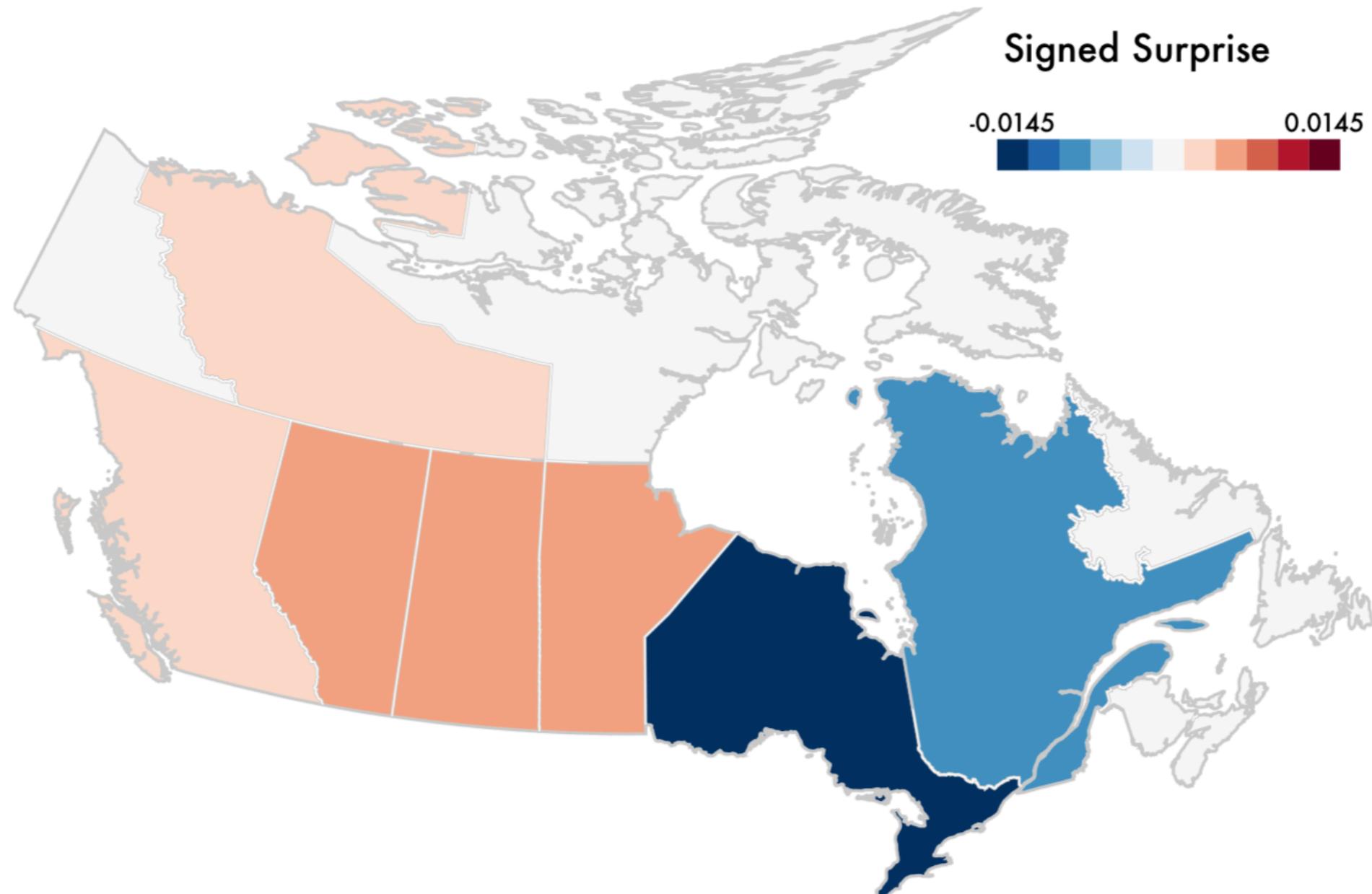
- **Let's assume two things**
 - If there is no big geographical differences in mischief, we'd expect each province to have the same per capita rate
 - If there is no big geographical difference in mischief, we'd expect variability to increase for smaller populations
- **Given these two models of how we expect the data to appear, we can measure deviations from these models**

<https://medium.com/@uwdata/surprise-maps-showing-the-unexpected-e92b67398865>

Surprise map

Bluish regions - where we have *less* mischief than we'd expect, given our models

Reddish regions - where we have *more* mischief than we'd expect, given our models

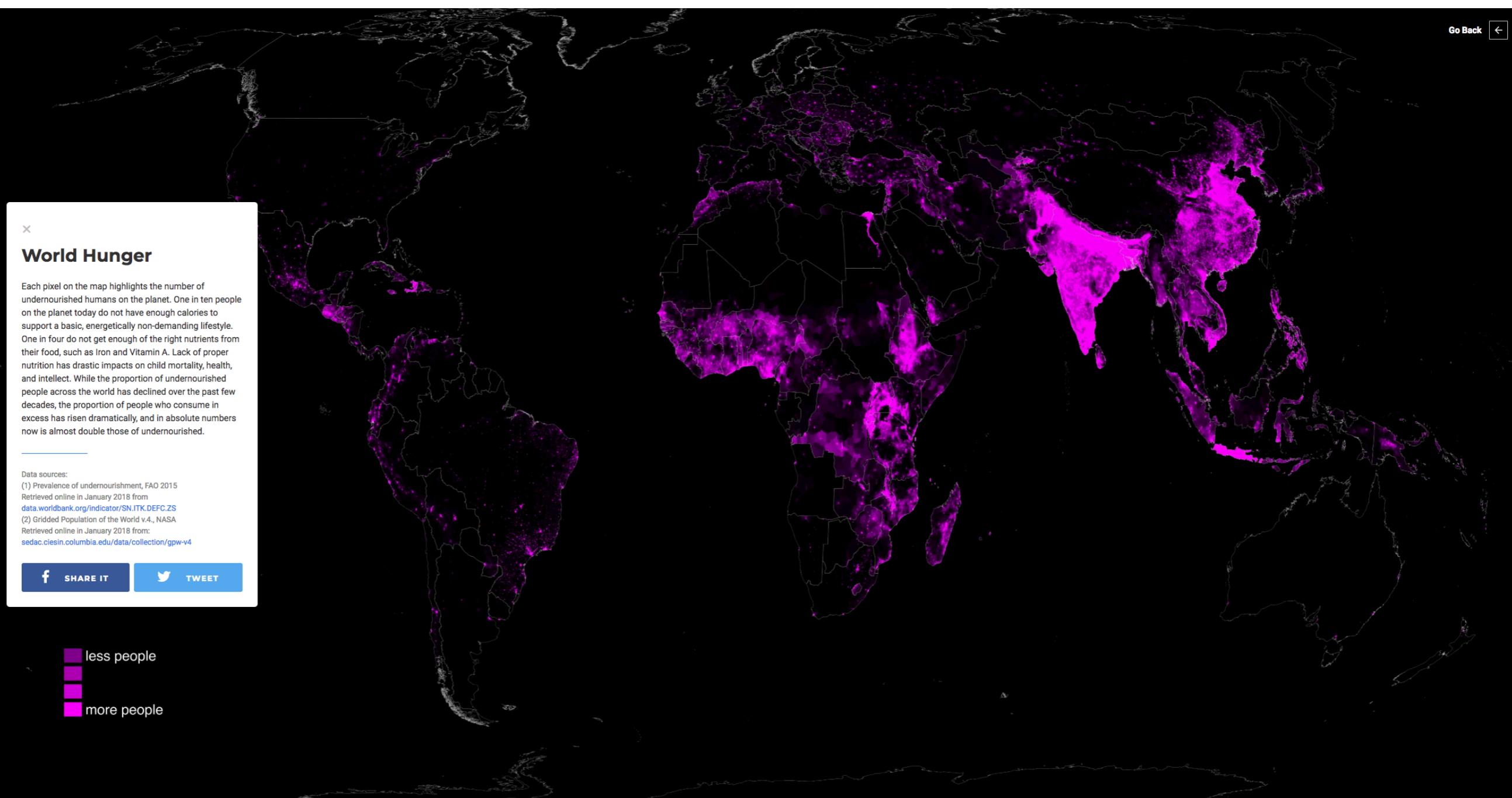


<https://medium.com/@uwdata/surprise-maps-showing-the-unexpected-e92b67398865>

Interested to read more...

Surprise maps

- Paper | <http://idl.cs.washington.edu/papers/surprise-maps/>
- Code | <https://github.com/uwdata/bayesian-surprise>



What works well?

What could be improved?

<https://www.colours-of-food-security.com/exhibition/181/>

Discussion points

- Beautiful and engaging visualization
 - Dramatic black background and bright colours
- Several maps just show population; could benefit from per capita normalization
- Difficult to see values on very small areas
- Do we need to show these data on a map?
 - Map helps spot differences between continents quickly (probably faster than reading labels on a bar chart)
 - Could complement with a small bar chart of countries with highest values (link to map); bars give equal visual weight to differently sized regions
- Could add zooming/panning interaction
- Consider visually distinguishing areas where no measure was taken (no data) and where measured data has a value of zero

Keypoints

Key points

- **Filtering**
 - Consider helping your user to make informed filtering decisions
 - Embed visuals in your widgets (scented widgets)
 - If appropriate, give visual cues as to how much data has been removed
- **Aggregation**
 - Binning is a great way to deal with large data sets (e.g. histograms, binned scatterplots)
 - Be aware that how you define your bins (shape and size) can change the resulting pattern
 - Avoid making population maps
 - Normalize per capita
 - Consider a surprise map approach

What I've been reading...



FLOWINGDATA

<https://flowingdata.com>