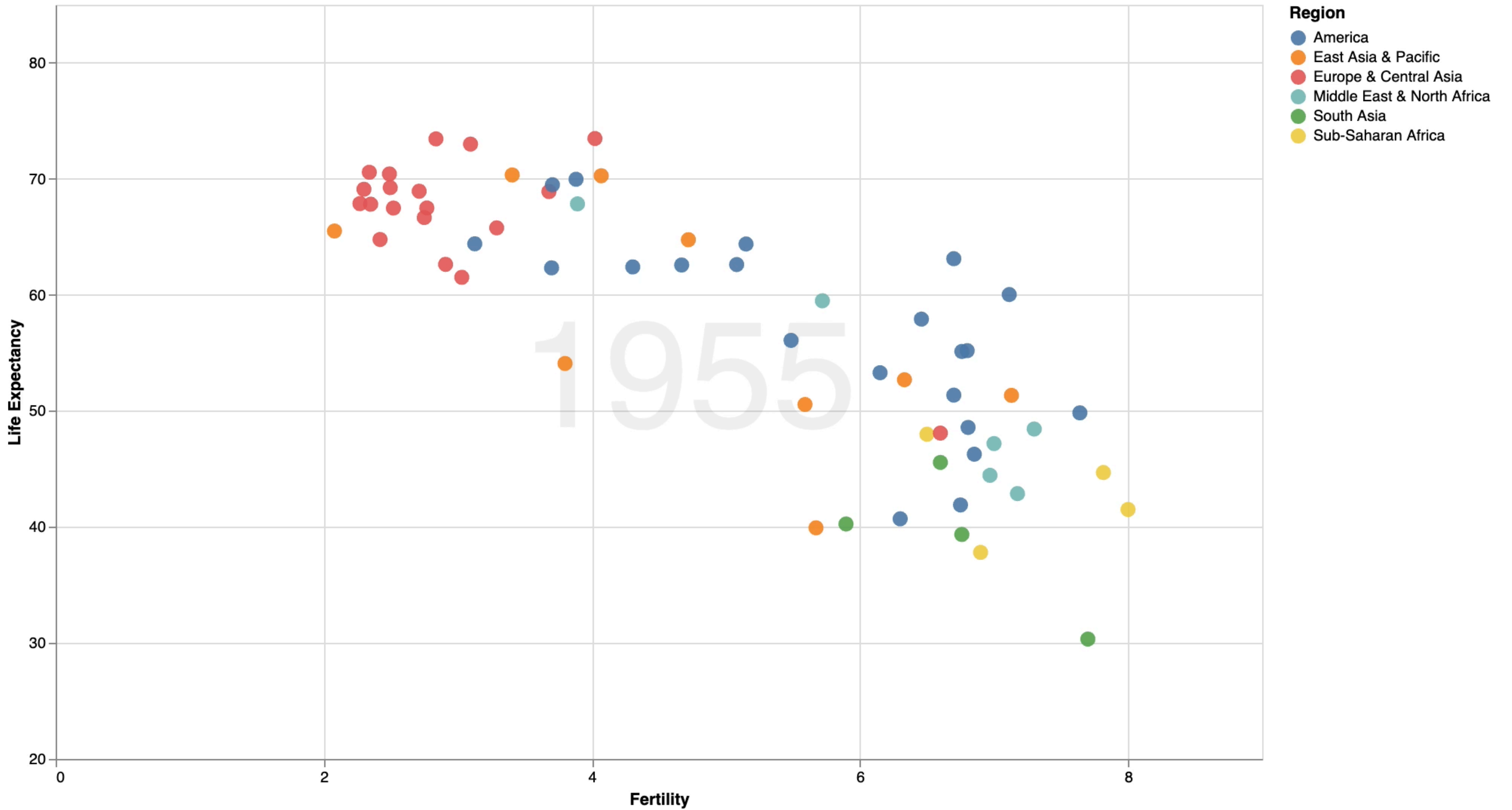# Enhancing Decision-Making through Interactive Data Visualization

Dominik Moritz @domoritz

Carnegie Mellon University

dig.cmu.edu

# 6 takeaways about data visualization

**Exposure**, the effective laying open of the data to display the unanticipated, is to us a major portion of data analysis. Formal statistics has given almost no guidance to exposure; indeed, it is not clear how the **informality** and **flexibility** appropriate to the **exploratory character of exposure** can be fitted into any of the structures of formal statistics so far proposed.

Data Analysis & Statistics. Tukey and Wilk. *1965.*

Effective Data Visualization. Heer. *2015.*

Nothing - not the careful logic of mathematics, not statistical models and theories, not the awesome arithmetic power of modern computers - nothing can substitute here for the **flexibility of the informed human** mind.

Accordingly, both approaches and techniques need to be structured so as to **facilitate human involvement and intervention**.

Data Analysis & Statistics. Tukey and Wilk. *1965.*

Effective Data Visualization. Heer. *2015.*

# Set 1    Set 2    Set 3    Set 4

| X | Y | X | Y | X | Y | X | Y |
|---|---|---|---|---|---|---|---|
| 10 | 8.04 | 10 | 9.14 | 10 | 7.46 | 8 | 6.58 |
| 8 | 6.95 | 8 | 8.14 | 8 | 6.77 | 8 | 5.76 |
| 13 | 7.58 | 13 | 8.74 | 13 | 12.74 | 8 | 7.71 |
| 9 | 8.81 | 9 | 8.77 | 9 | 7.11 | 8 | 8.84 |
| 11 | 8.33 | 11 | 9.26 | 11 | 7.81 | 8 | 8.47 |
| 14 | 9.96 | 14 | 8.1 | 14 | 8.84 | 8 | 7.04 |
| 6 | 7.24 | 6 | 6.13 | 6 | 6.08 | 8 | 5.25 |
| 4 | 4.26 | 4 | 3.1 | 4 | 5.39 | 19 | 12.5 |
| 12 | 10.84 | 12 | 9.11 | 12 | 8.15 | 8 | 5.56 |
| 7 | 4.82 | 7 | 7.26 | 7 | 6.42 | 8 | 7.91 |
| 5 | 5.68 | 5 | 4.74 | 5 | 5.73 | 8 | 6.89 |

## Summary Statistics          ## Linear Regression

$u_X = 9.0$     $\sigma_X = 3.317$     $Y^2 = 3 + 0.5\,X$

$u_Y = 7.5$     $\sigma_Y = 2.03$      $R^2 = 0.67$

[Anscombe 1973]
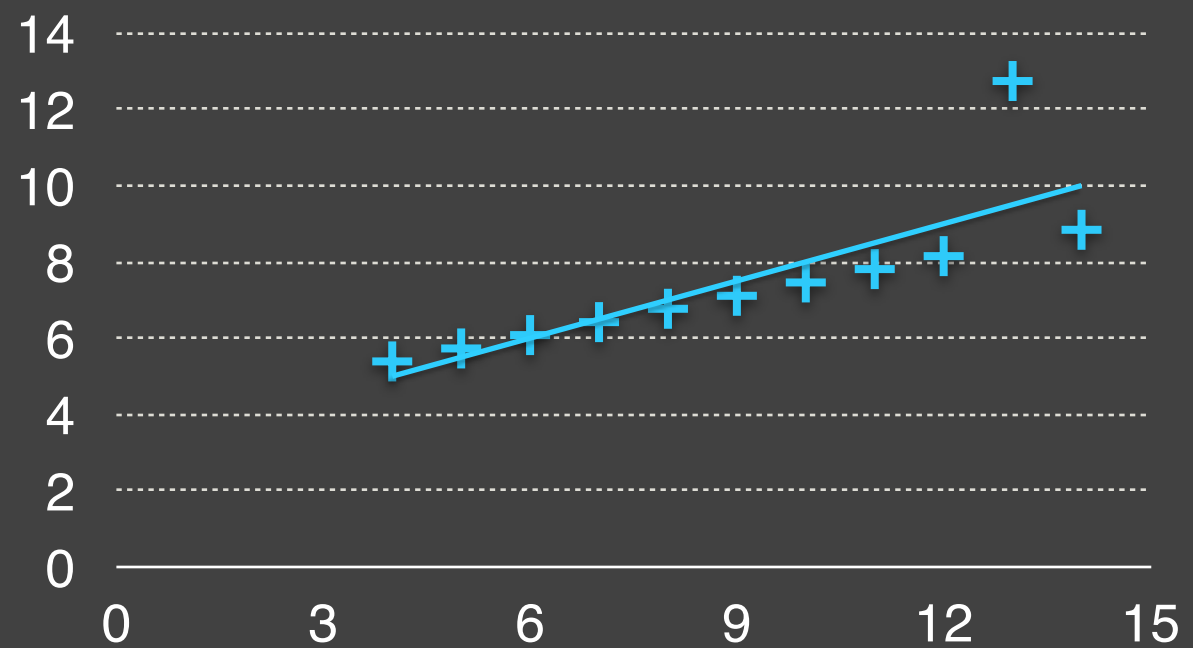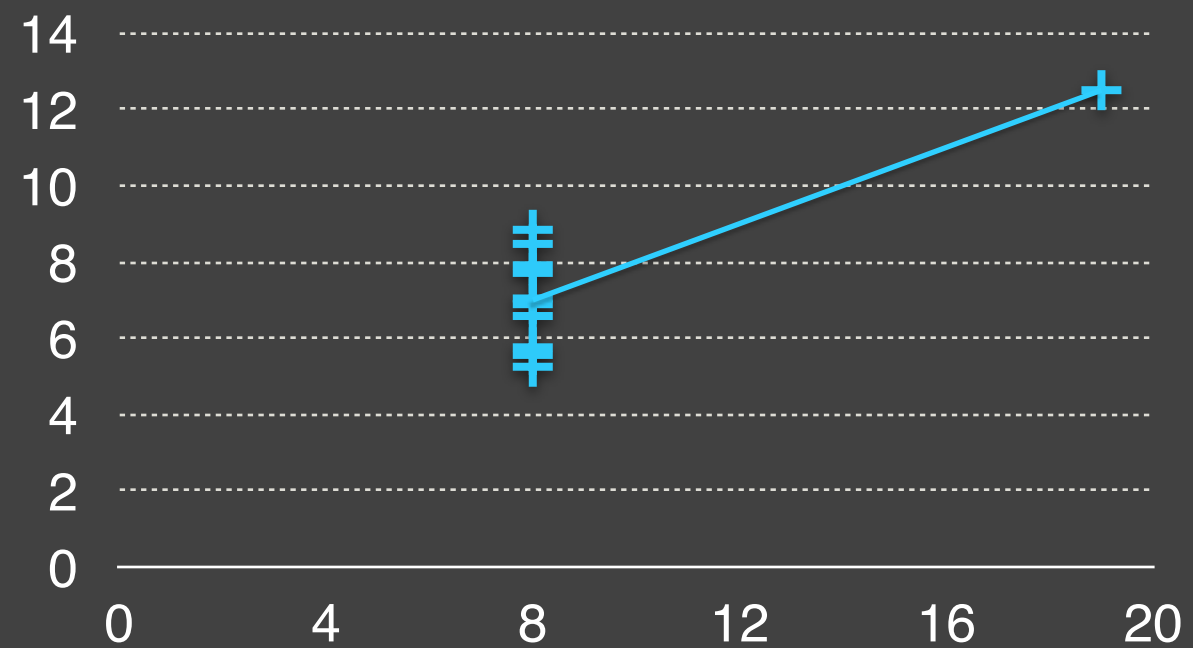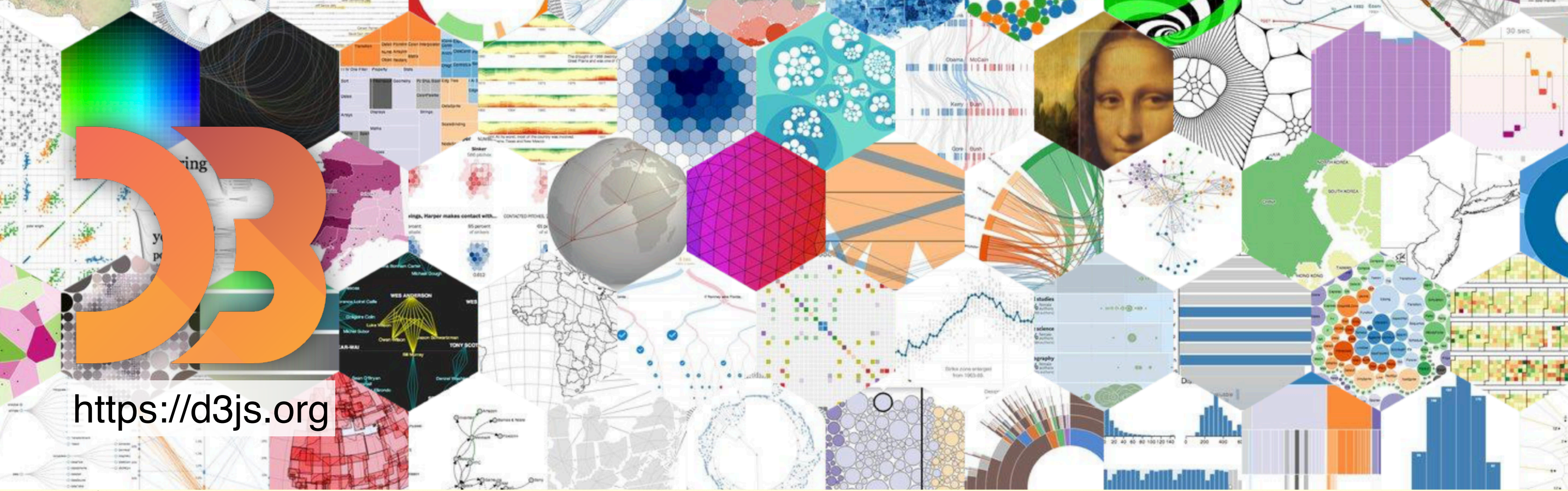
Takeaway:

Machine Learning, AI, and Statistics are people problems. For them to be effective, we need to design for human involvement. 🤖 👩‍💻

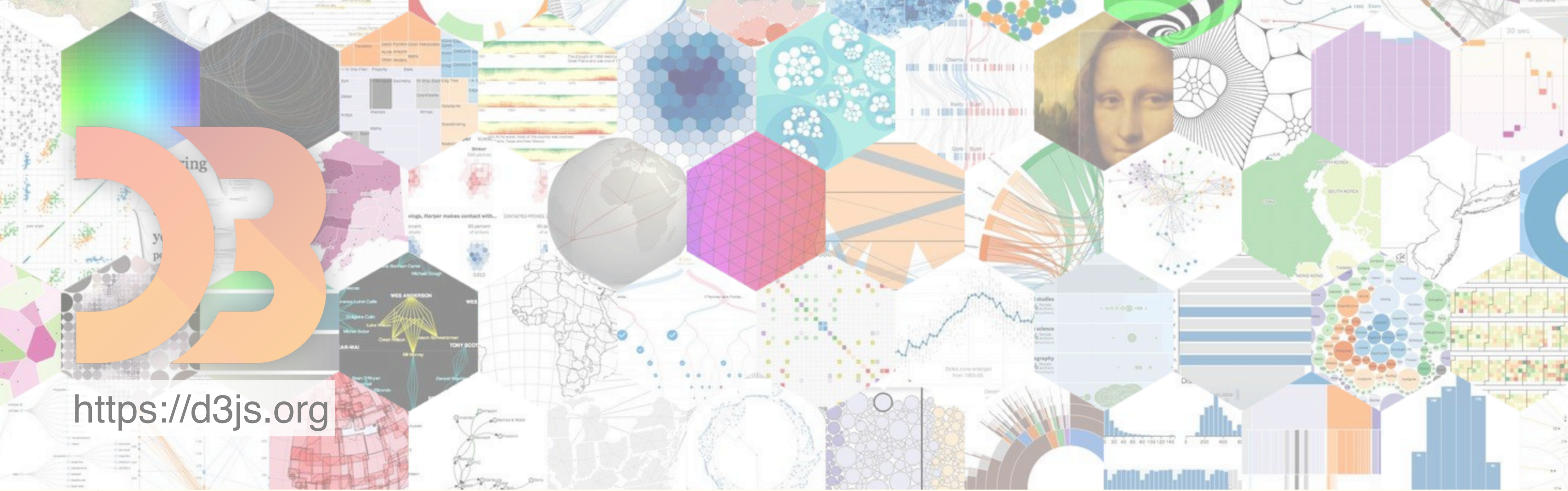| | H | I | J | K | L | M | N | | | X | Y | Z | AA | AB | AC |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 17 | | | | | | | | | Cr | | | | | | |
| 18 | | | | | | | | | | | | | | | |
| 19 | | | | | | | | | | | | | | | |
| 20 | | | | | | | | | | | | | | | |
| 21 | | | | | | | | | | | | | | | |
| 22 | | | | | | | | | | | | | | | |
| 23 | | | | | | | | | | | | | | | |
| 24 | | | | | | | | | | | | | | | |
| 25 | | | | | | | | | | | | | | | |
| 26 | | | | | | | | | | | | | | | |
| 27 | | | | | | | | | | | | | | | |
| 28 | | | | | | | | | | | | | | | |
| 29 | | | | | | | | | | | | | | | |
| 30 | | | | | | | | | | | | | | | |
| 31 | | | | | | | | | | | | | | | |
| 32 | | | | | | | | | | | | | | | |
| 33 | | | | | | | | | | | | | | | |
| 34 | | | | | | | | | | | | | | | |
| 35 | | | | | | | | | | | | | | | |
| 36 | | | | | | | | | | | | | | | |
| 37 | | | | | | | | | | | | | | | |

**Chart** ▸
Sparklines...
Table

Add-ins ▸

Page Break
Reset All Page Breaks

Function...

Name ▸

New Comment

Column
Bar
Line
**Area**
Pie
Treemap
Sunburst
Histogram
Pareto
Box and Whisker

https://d3js.org

| Chart | > | | Column |
| Sparklines... | | | Bar |
| Table | | | Line |
| | | | **Area** |
| Add-ins | > | | Pie |
| | | | Treemap |
| Page Break | | | Sunburst |
| Reset All Page Breaks | | | Histogram |
| Function... | | | Pareto |
| Name | > | | Box and Whisker |
| New Comment | | | |

https://d3js.org

Excel

| Chart | > | Column |
| Sparklines... | | Bar |
| Table | | Line |
| | | **Area** |
| Add-ins | > | Pie |
| | | Treemap |
| Page Break | | Sunburst |
| Reset All Page Breaks | | Histogram |
| Function... | | Pareto |
| Name | > | |
| New Comment | | Box and Whisker |

How do we make visualizations in the midst of an analysis?

Vega-Lite

Run
Auto

Commands · Export · Share · Gist · Examples

Help · Settings · Sign in with

**VEGA-LITE** · CONFIG

```
1   {
2     "$schema": "https://vega.github.io/schema/vega-lite/v4.json",
3     "description": "Drag the sliders to highlight points.",
4     "data": {"url": "data/gapminder.json"},
5     "width": 600,
6     "height": 400,
7     "layer": [
8       {
9         "transform": [
10          {"filter": {"field": "country", "equal": "Afghanistan"}},
11          {"filter": {"selection": "year"}}
12        ],
13        "mark": {
14          "type": "text",
15          "fontSize": 100,
16          "x": 420,
17          "y": 250,
18          "opacity": 0.06
19        },
20        "encoding": {"text": {"field": "year"}}
21      },
22      {
23        "transform": [
24          {
25            "lookup": "cluster",
26            "from": {
27              "key": "id",
28              "fields": ["name"],
29              "data": {
30                "values": [
31                  {"id": 0, "name": "South Asia"},
32                  {"id": 1, "name": "Europe & Central Asia"},
33                  {"id": 2, "name": "Sub-Saharan Africa"},
34                  {"id": 3, "name": "America"},
35                  {"id": 4, "name": "East Asia & Pacific"},
36                  {"id": 5, "name": "Middle East & North Africa"}
```

Compiled Vega · Extended Vega-Lite Spec

**Region**
- America
- East Asia & Pacific
- Europe & Central Asia
- Middle East & North Africa
- South Asia
- Sub-Saharan Africa

Life Expectancy

1955

Fertility

Year ————————— 1955

Vega 5.17.1, Vega-Lite 4.17.0, Vega-Tooltip 0.24.2, Editor 0.92.2

LOGS · **DATA VIEWER** · SIGNAL VIEWER

year_store

| unit | fields | values |
|---|---|---|
| "layer_1_layer_1" | [{"type":"E","field":"year"}] | [1955] |

# Vega-Lite

A high-level declarative grammar for interactive multi-view charts.

# Vega-Lite

```
{
  data: {
    url: "seattle-weather.csv"
  },
  mark: "bar",
  encoding: {
    x: {
      timeUnit: "month",
      field: "date"
    },
    y: {
      aggregate: "count"
    },
    color: {
      field: "weather"
    }
  }
}
```
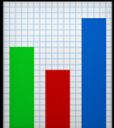
# Vega-Lite

A high-level grammar for creating interactive multi-view charts.
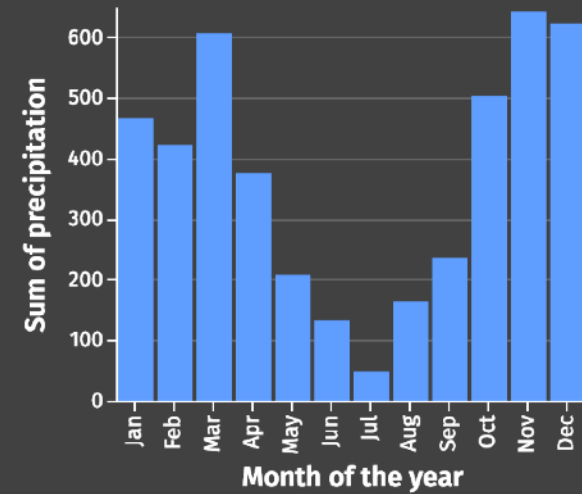
2.1M monthly downloads from CDN.
Used at , Microsoft, Google, Netflix, etc.

vega.github.io/vega-lite

Takeaway:

Grammar-based visualization tools (such as Vega-Lite) support flexible interactive visualization and exploration. 🗣️ 📊

# How has Barley Yield Changed at Different Sites?

| Year | Site | Median Yield |
|------|------|-------------|
| 1931 | "University Farm" | 36.58 |
| 1931 | "Waseca" | 52.71 |
| 1931 | "Morris" | 28.73 |
| 1931 | "Crookston" | 42.85 |
| 1931 | "Grand Rapids" | 29.71 |
| 1931 | "Duluth" | 30.63 |
| 1932 | "University Farm" | 27.75 |
| 1932 | "Waseca" | 39.88 |
| 1932 | "Morris" | 43.36 |
| 1932 | "Crookston" | 32.09 |
| 1932 | "Grand Rapids" | 20.26 |
| 1932 | "Duluth" | 24.28 |

# How has Barley Yield Changed at Different Sites?

| Year | Site | Median Yield |
|------|------|------|
| 1931 | "University Farm" | 36.58 |
| 1931 | "Waseca" | 52.71 |
| 1931 | "Morris" | 28.73 |
| 1931 | "Crookston" | 42.85 |
| 1931 | "Grand Rapids" | 29.71 |
| 1931 | "Duluth" | 30.63 |
| 1932 | "University Farm" | 27.75 |
| 1932 | "Waseca" | 39.88 |
| 1932 | "Morris" | 43.36 |
| 1932 | "Crookston" | 32.09 |
| 1932 | "Grand Rapids" | 20.26 |
| 1932 | "Duluth" | 24.28 |

# How has Barley Yield Changed at Different Sites?

# Draco

A formal model of visualization design with learned constraints. Can be used to automatically create "good" visualizations.

New version is in progress at github.com/cmudig/draco2



"Visualize temperature in the weather dataset."

Takeaway:

When we design for perception, otherwise hidden patterns emerge. 👁 🧠

Design Visualizations → Explore Data
Make Decisions

# Movie Data

| | |
|---|---|
| **Title** | String (N) |
| **IMDB Rating** | Number (Q) |
| **Rotten Tomatoes Rating** | Number (Q) |
| **MPAA Rating** | String (O) |
| **Release Date** | Date (T) |

Collection Date

Movies from the Future?

**Common Analysis Pitfalls:**

Overlook data quality issues

Fixate on specific relationships

Other cognitive biases

datavoyager

**Add Dataset**                                          close

Change Dataset | Paste or Upload Data | From URL

Barley                          Cars

Crimea                          Driving

Iris                            Jobs

Population                      Movies
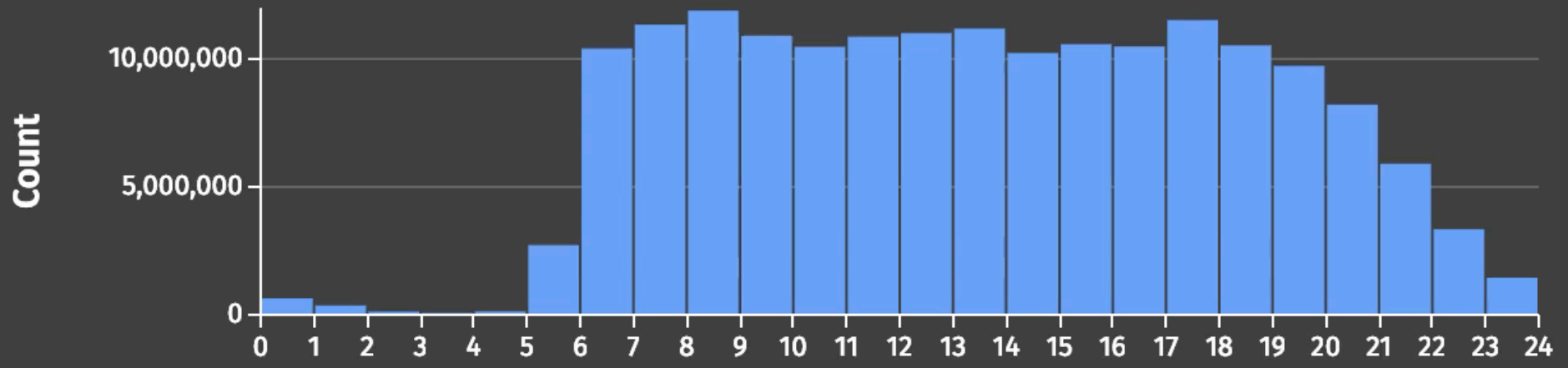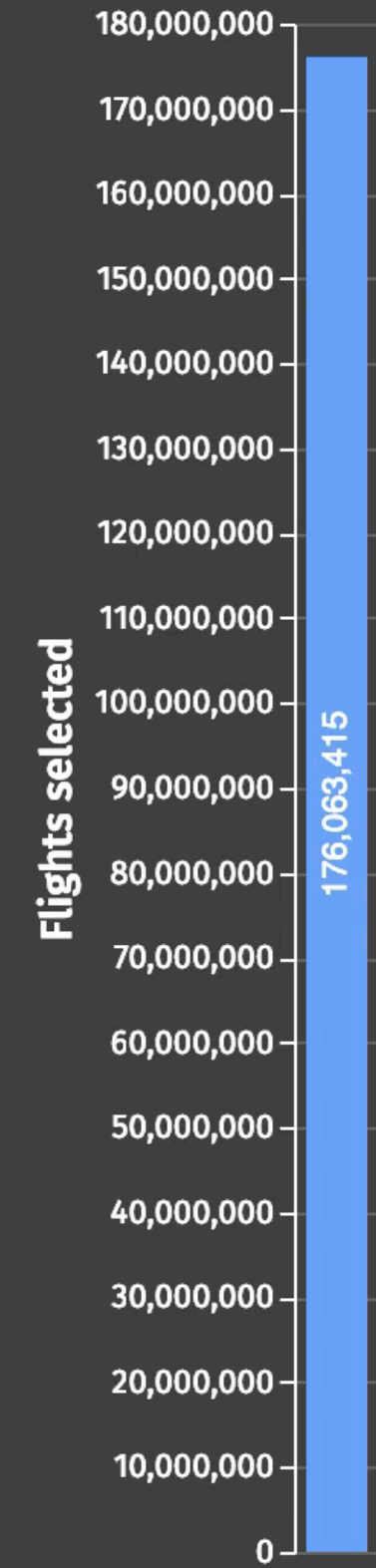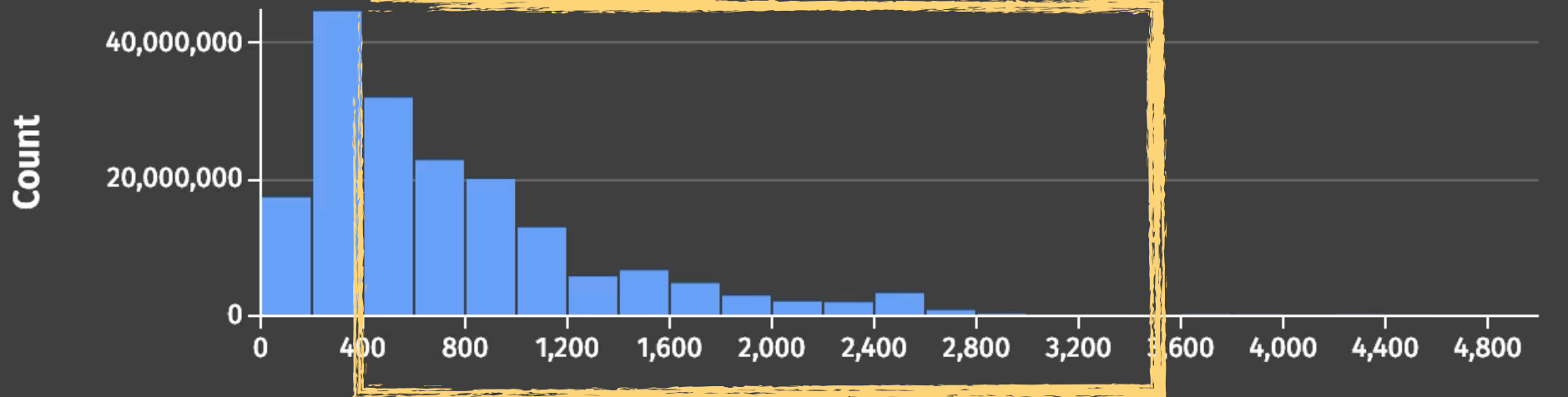
Birdstrikes                     Burtin

Campaigns

Takeaway:

UI tools can encourage best-practices. 🖥️ ✅

Arrival Delay in Minutes

Departure Time

Distance in Miles

Flights selected — 176,063,415

Arrival Delay in Minutes

Departure Time
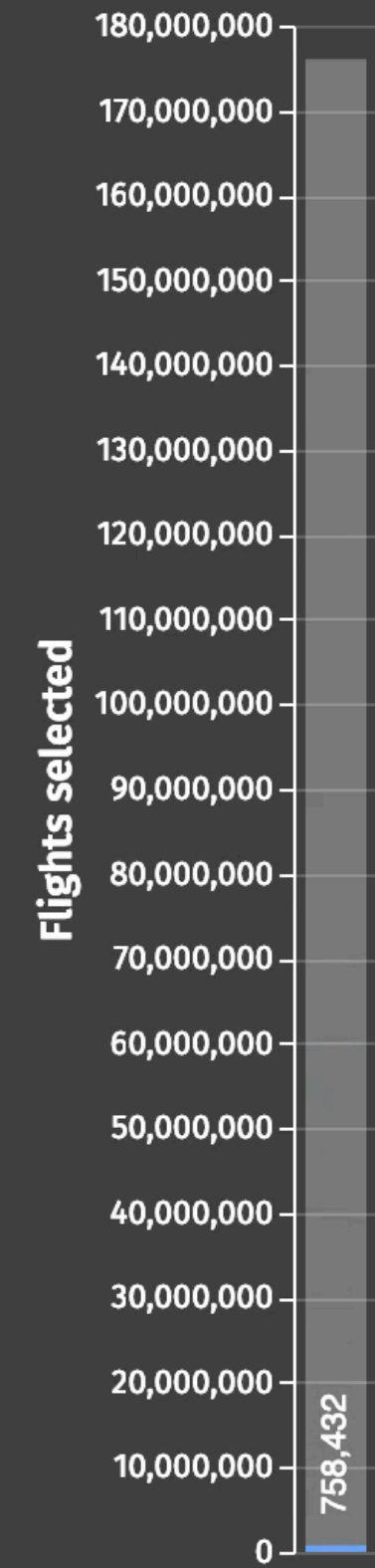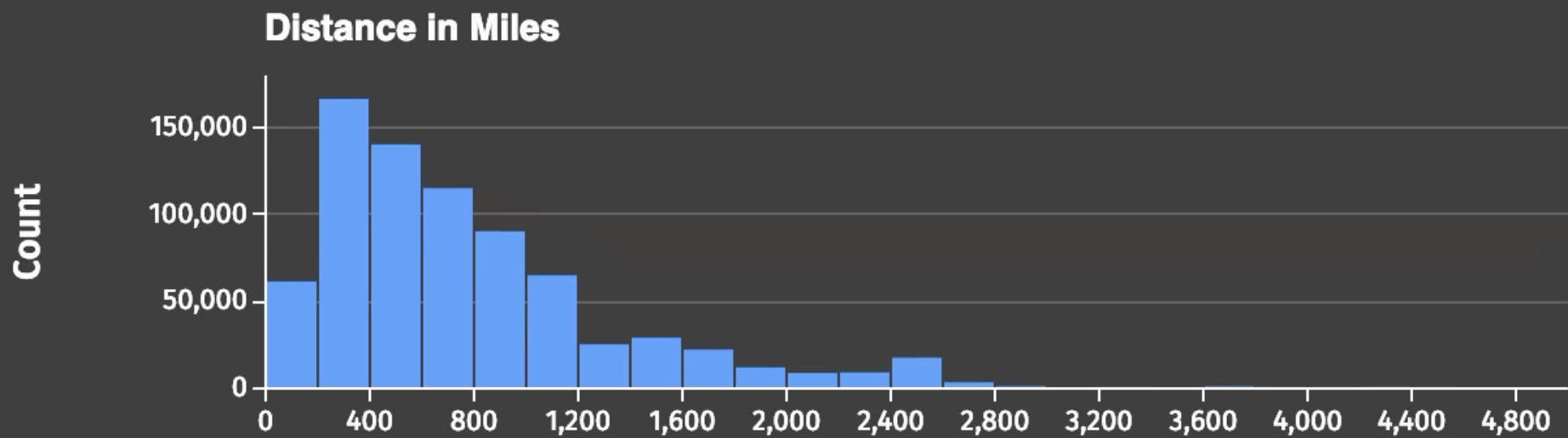
Distance in Miles

Flights selected: 176,063,415

Takeaway:

Interactivity enables us to find patterns that exist across multiple dimensions. 👆 🔢

# How do we interact with billion+record datasets in real-time?

# How do we interact with billion+record datasets in real-time?

Delays reduce engagement
and lead to fewer observations.

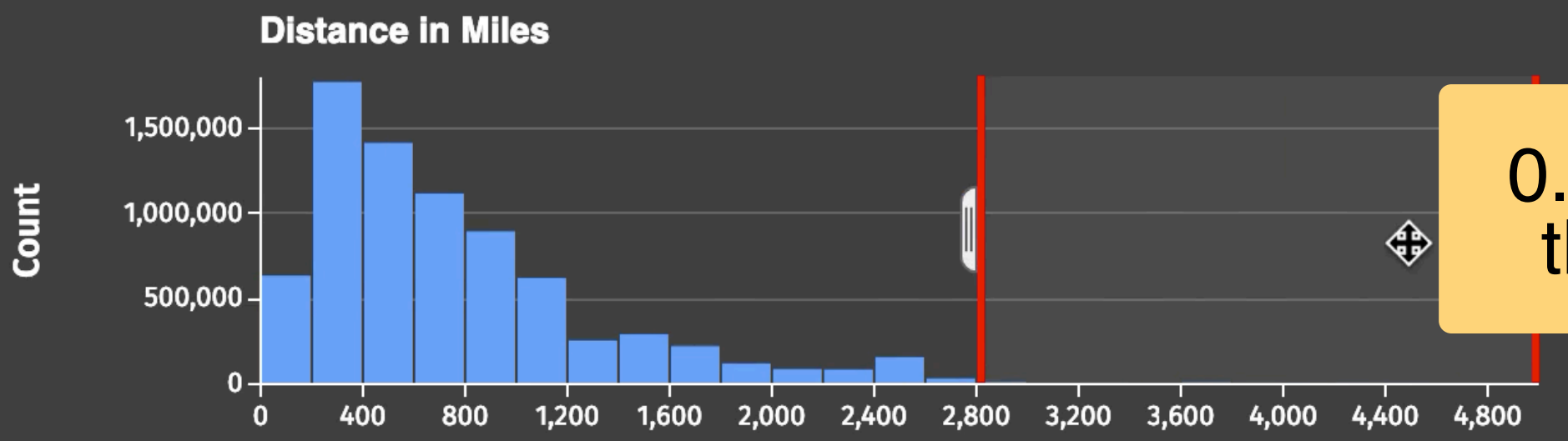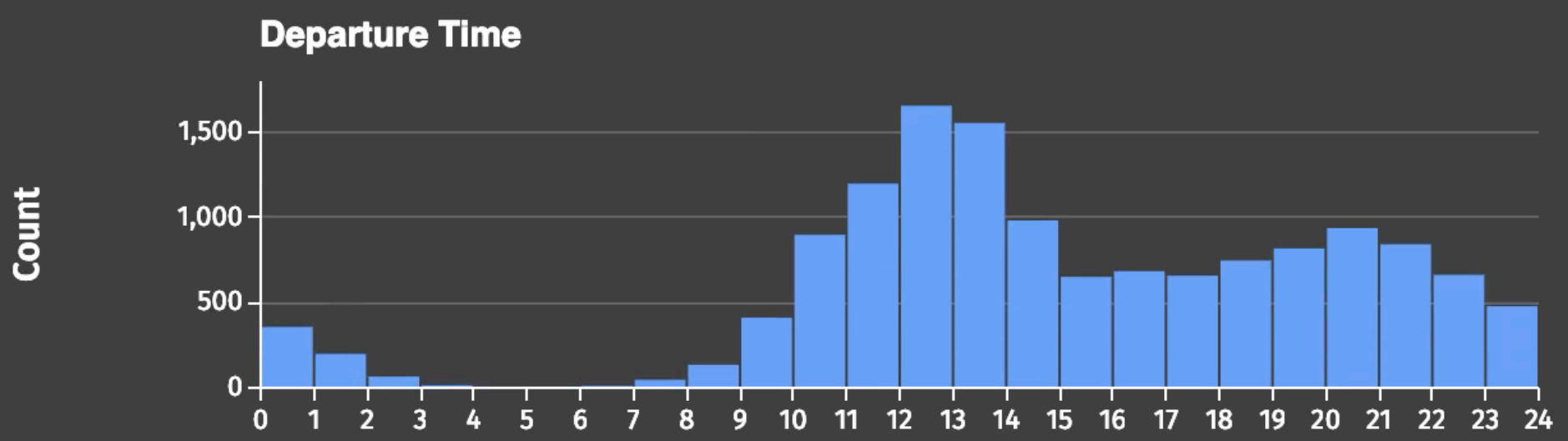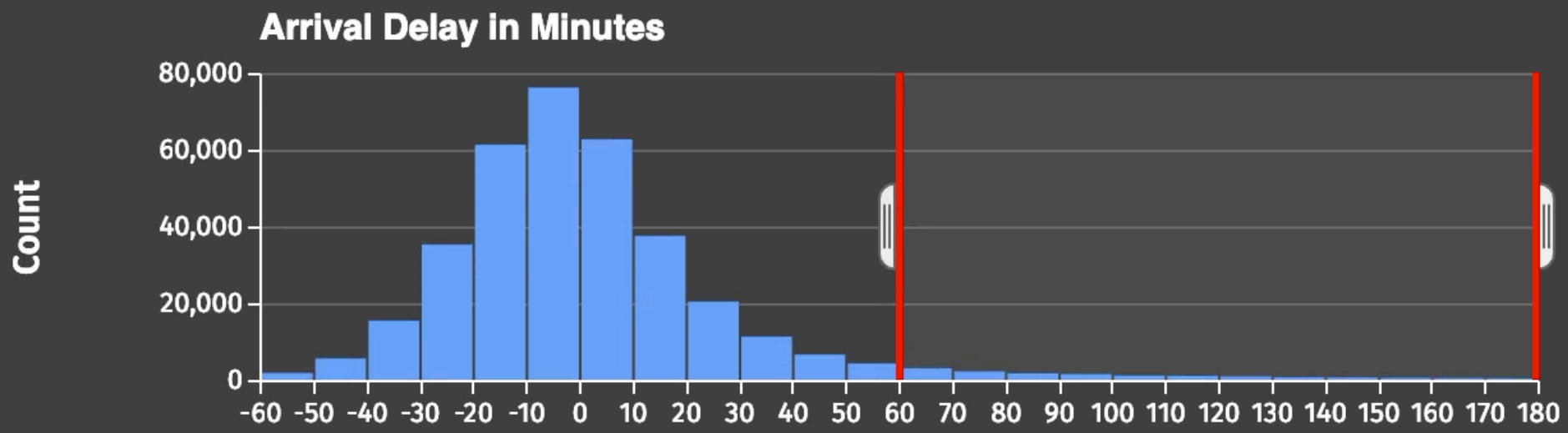The Effect of Interactive Latency. Liu, Heer. *IEEE Infovis 2014.*

Delays may bias analysts
towards convenient data.

# Falcon

Zero-latency crossfiltering for massive datasets.

Leverages smart prefetching and precomputation designed around human perception.

Available as open source at
github.com/vega/falcon

Takeaway:

Interactivity should be real-time regardless of the scale of the data. ⏱️ 🌍

Machine Learning, AI, and Statistics are people problems. For them to be effective, we need to **design for human involvement**. 🤖 🧑‍💻

**Grammar-based visualization tools** (such as Vega-Lite) support flexible interactive visualization and exploration. 🗣️ 📊

When we **design for perception**, otherwise hidden patterns emerge. 👁️ 🧠

UI tools can encourage best-practices. 🖥️ ✅

**Interactivity** enables us to find patterns that exist **across multiple dimensions**. 👆 🔢

Interactivity should be **real-time** regardless of the scale of the data. ⏱️ 🌎

Research Mission

Empower everyone to effectively
analyze and communicate data,
by *designing interactive systems that
richly integrate the strengths
of both people and machines*.

dig.cmu.edu