

Vegas

The Missing Matplotlib for
Scala/Spark

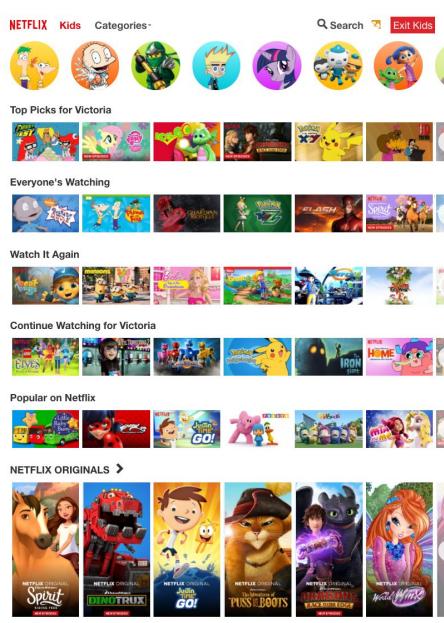
| DB Tsai
Roger Menezes

NETFLIX

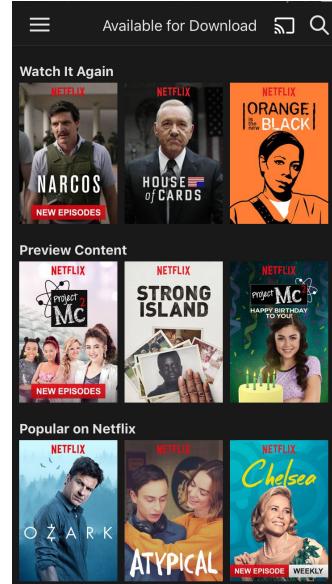
Netflix Recommendations



Homepage



Kids Page



Downloads Page

Every aspect
of the
Experience is
Machine
Learned



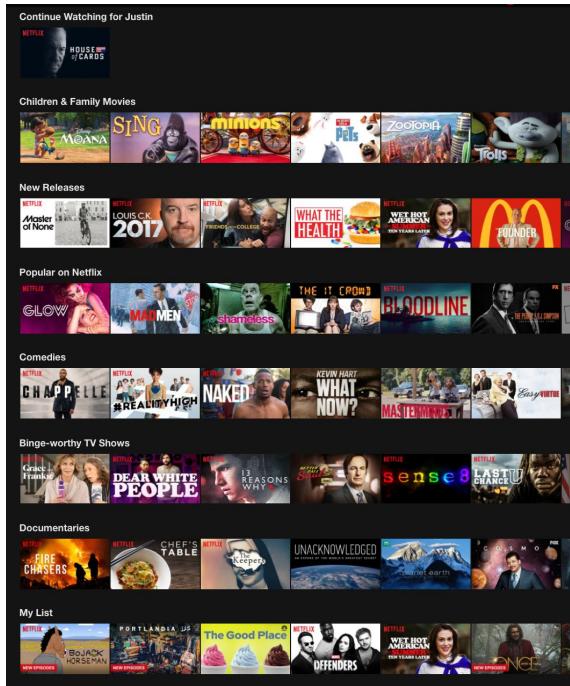
#netflixeverywhere

2017

- > 100M members
- > 190 countries



Multiple Devices



NETFLIX

NETFLIX ORIGINAL

**MARVEL
THE DEFENDERS**

Watch Now

They're not friends. But these four New Yorkers will fight as one to save the city they love.

PLAY

Continue Watching for Justin



My List



Popular on Netflix



Trending Now



Watch It Again



Search



NETFLIX ORIGINALS



Top Picks for Justin



Because you watched Maria Bamford: Old Baby ➔



Sims: 10 rows/page average

Children & Family Movies



Genres: 23 rows/page average

New Releases



Recently Added



NETFLIX

NETFLIX ORIGINAL

MARVEL

Search

**Billboard:**

Watch Now

They're not friends. But these four New Yorkers will fight as one to save the city they love.

▶ PLAY**✓ MY LIST**

Continue Watching for Justin

Continue Watching:

My List

My List:

PORTLANDIA

NEW EPISODES

The Good Place

NEW EPISODES

DEFENDERS

NEW EPISODES

WET HOT AMERICAN

TEN YEARS LATER

ONCE UPON A TIME

NEW EPISODES

Popular on Netflix

Popular on Netflix:

THE IT CROWD

BLOODLINE

Trending Now

Trending Now:

Watch It Again

Watch It Again:

NETFLIX ORIGINALS

**Originals Row**

NETFLIX ORIGINAL

NARCOS

NEW EPISODES

NETFLIX ORIGINAL

OZARK

NEW EPISODES

NETFLIX ORIGINAL

ATYPICAL

NEW EPISODES

NETFLIX ORIGINAL

ARRESTED DEVELOPMENT

NEW EPISODES

NETFLIX ORIGINAL

DEFENDERS

NEW EPISODES

NETFLIX ORIGINAL

Chelsea

Top Picks for Justin

Top Picks:

Because you watched Maria Bamford: Old Baby ➤

Because You Watched:

Children & Family Movies

Genres:

New Releases

New Releases:

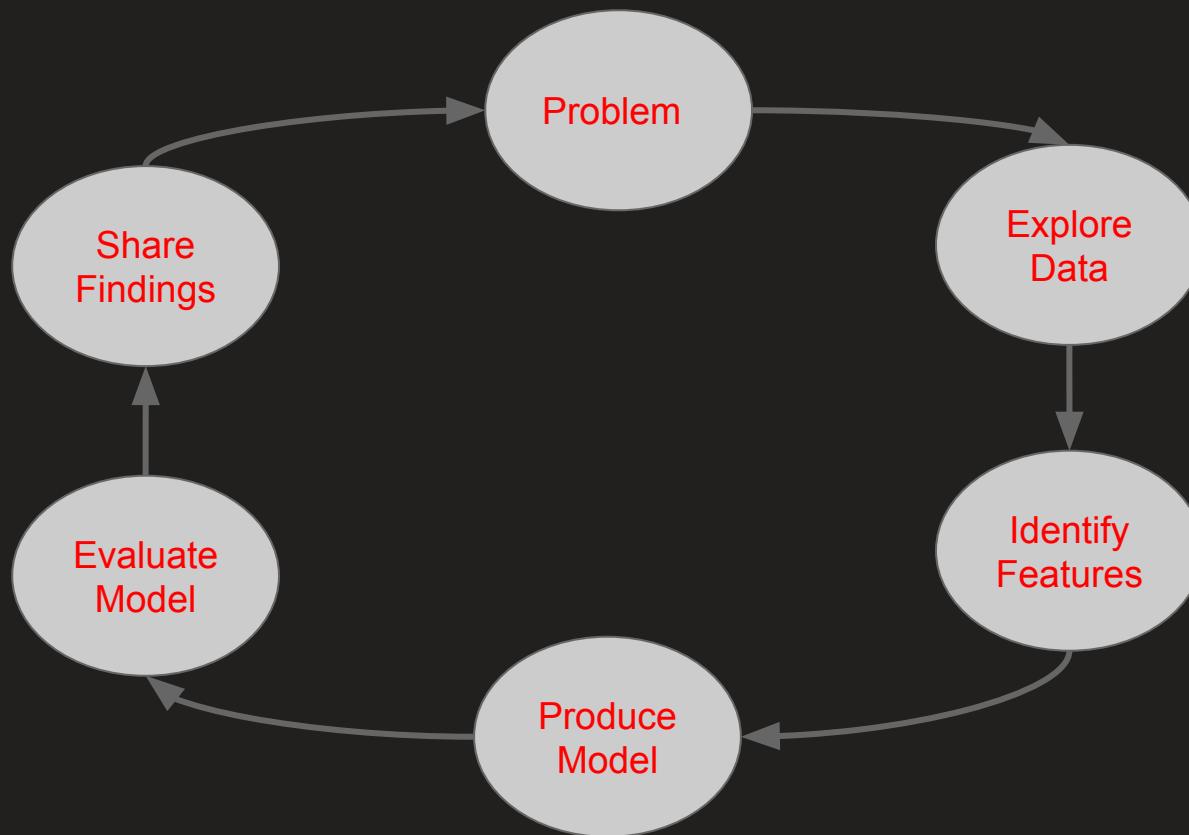
Recently Added

Recently Added:

Machine Learning at Netflix

- Optimize the Experimentation usecase vs Productionization
- Experimentation
 - Opportunity sizing, Data Exploration
 - Feature Identification and Selection
 - Tweaks to ML algos
 - Model Evaluation

Experimenter's loop



Notebooks

- Optimal for Experimentation
- Sharing reproducible research
 - Facilitates feedback loop with Product Managers
- End to end ML experiment.
 - Interactivity drives productivity

Python Notebooks

The image displays three distinct environments for Python notebooks:

- Top Left:** A screenshot of a web browser window titled "r_notebook_example" at "localhost:8889/notebooks/...". It shows a Jupyter notebook cell with code to generate a scatter plot of diamond prices based on carat weight. Below the cell is a scatter plot titled "Simple spectral analysis" illustrating the Discrete Fourier Transform.
- Top Middle:** A screenshot of the IPython Notebook interface titled "spectrogram". It shows a spectrogram of an audio signal with two subplots: "Raw audio signal" and "Spectrogram".
- Bottom Right:** A screenshot of the Jupyter nbviewer page for the same notebook. It includes a header with "JUPYTER" and "FAQ" links. The notebook content is displayed with syntax highlighting and output cells. Below the notebook are four maps of the United States showing spatial data with a legend for "HH", "LH", "LL", "HL", and "Non-significant" categories.

Python Notebooks

- Seamless Experience - ML experimentation
- Well known Scientific computing libraries
- Huge catalog of Visualization plotting libraries
 - Matplotlib, Seaborn, Bokeh, BQPlot, Lightning, etc.

Scala Notebooks

- Zeppelin, Jupyter, Databricks, Spark-Notebooks, ...
- Computing library gap filling up
- Lack of Visualization Libraries
 - Main friction point in adoption
 - End to End ML use case not convincing

Introducing Vegas

- Visualization Library in Scala
- Mainly built for the notebook use case
- Scala wrapper around Vega-Lite
 - Missing Matplotlib for the Scala/Spark world.

NETFLIX

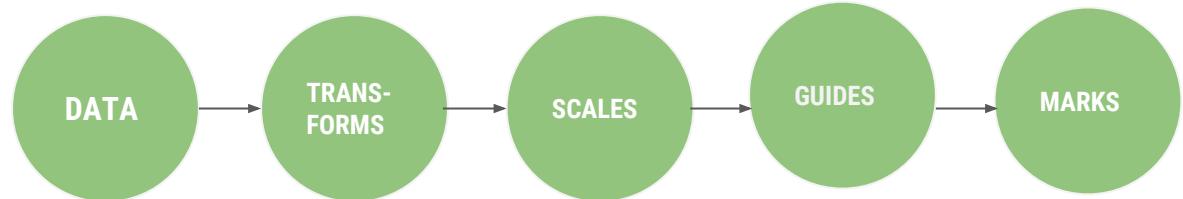
DECLARATIVE STATISTICAL VISUALIZATION GRAMMAR IN SCALA

VEGAS

You tell it **WHAT** should be done with the data, and it knows **HOW** to do it!

Operations such as *filtering, aggregation, faceting* are built into the visualization, rather than putting the burden on the user to massage the data into shape.

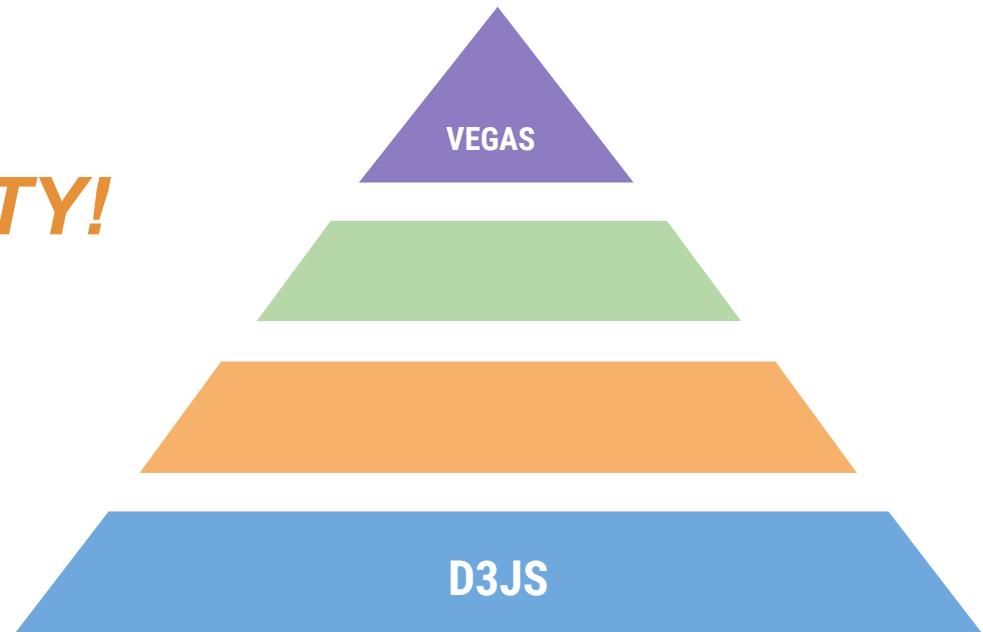
Complex visualizations can be built with a few high level abstractions:



cf : Altair Talk by Brian Granger in PyData 2016 <https://youtu.be/v5mrwq7yJc4>

Added Bonus of Declarative Visualizations:

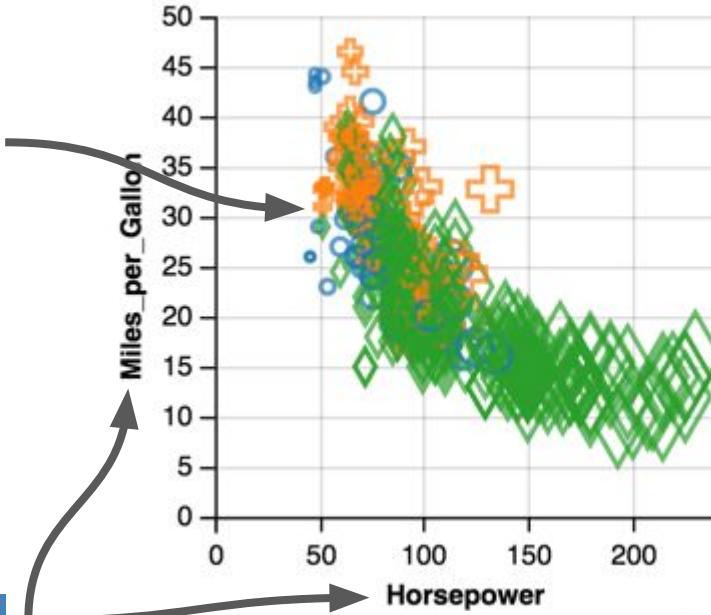
INTERACTIVITY!



VEGAS CODE EXPANDS OUT TO D3JS CODE!

Anatomy of a plot: Channels

SHAPE CHANNEL



X/Y CHANNEL

COLOR CHANNEL

Origin
Europe
Japan
USA

Horsepower
50
100
150
200

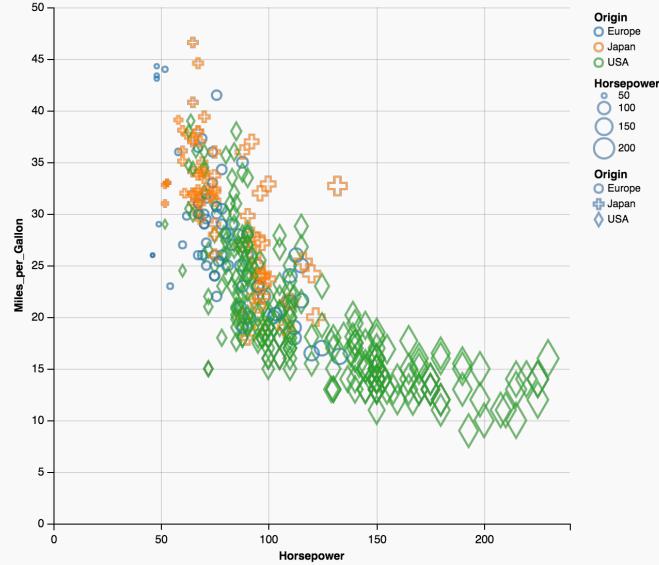
Origin
Europe
Japan
USA

SIZE CHANNEL

```
1 | carsDf.printSchema
```

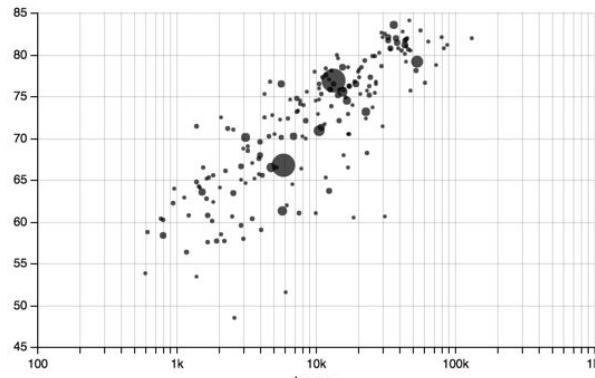
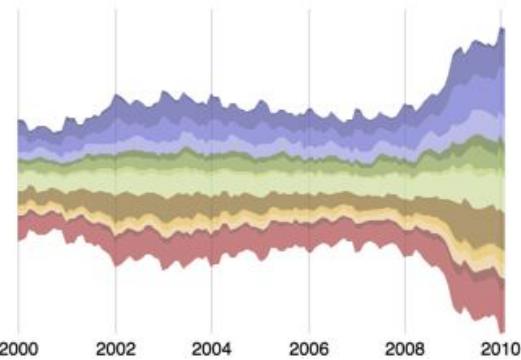
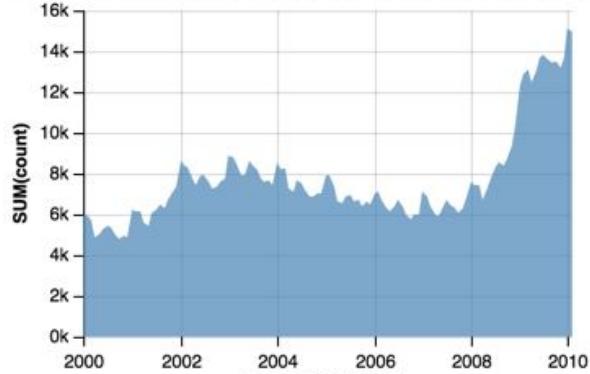
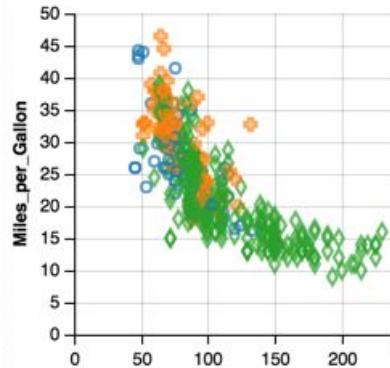
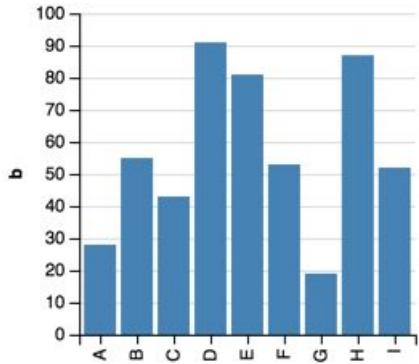
```
root
| -- Acceleration: double (nullable = true)
| -- Cylinders: long (nullable = true)
| -- Displacement: double (nullable = true)
| -- Horsepower: long (nullable = true)
| -- Miles_per_Gallon: double (nullable = true)
| -- Name: string (nullable = true)
| -- Origin: string (nullable = true)
| -- Weight_in_lbs: long (nullable = true)
| -- Year: string (nullable = true)
```

```
Vegas("Miles_per_Gallon vs horsepower").
  withDataFrame(carsDf).
  mark(Point).
  encodeX("Horsepower", Quantitative).
  encodeY("Miles_per_Gallon", Quantitative).
  encodeColor("Origin", Nominal).
  encodeSize("Horsepower", Quantitative).
  encodeShape("Origin", Nominal).
  show
```

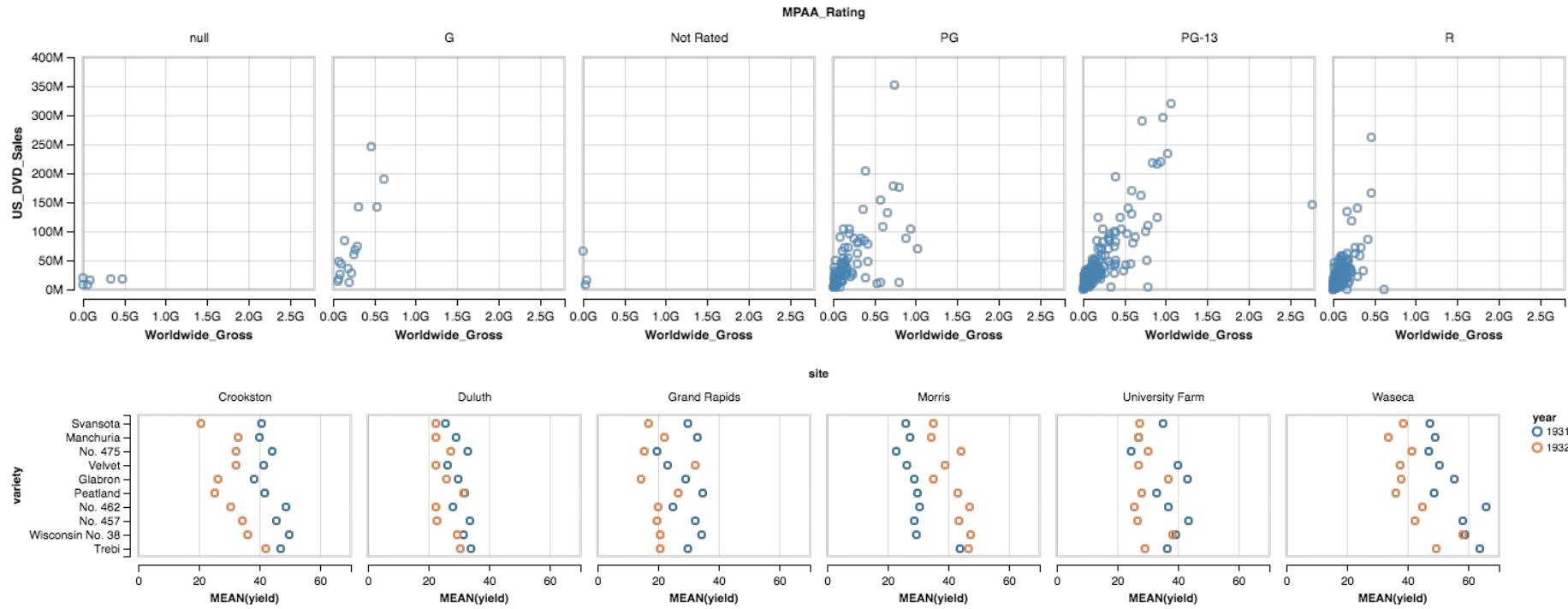


Features...

1. *Supports most plot types*



2. Trellis plots

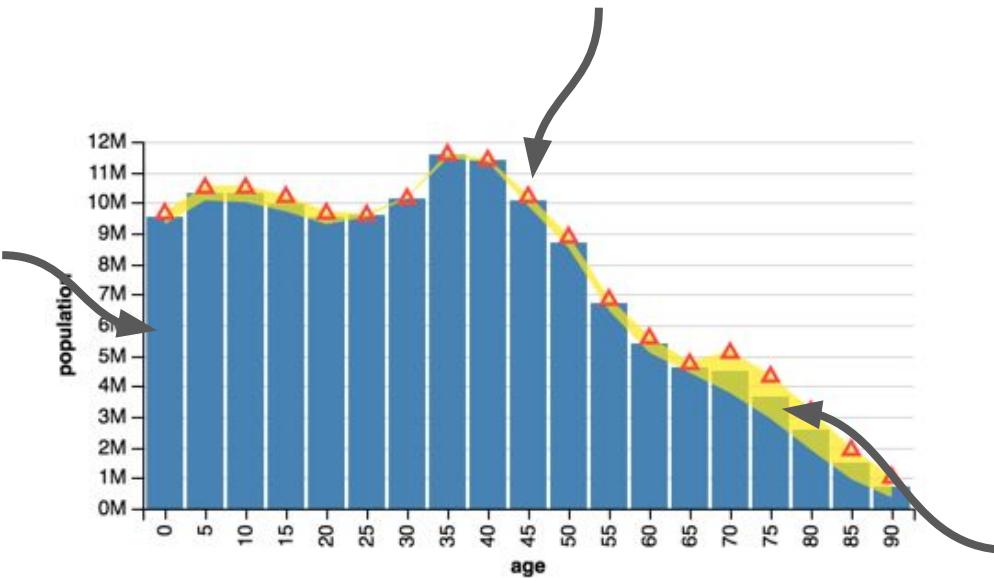


3. Layers

LAYER 1.

LAYER 2.

LAYER 3.



4. Notebook and Consoles

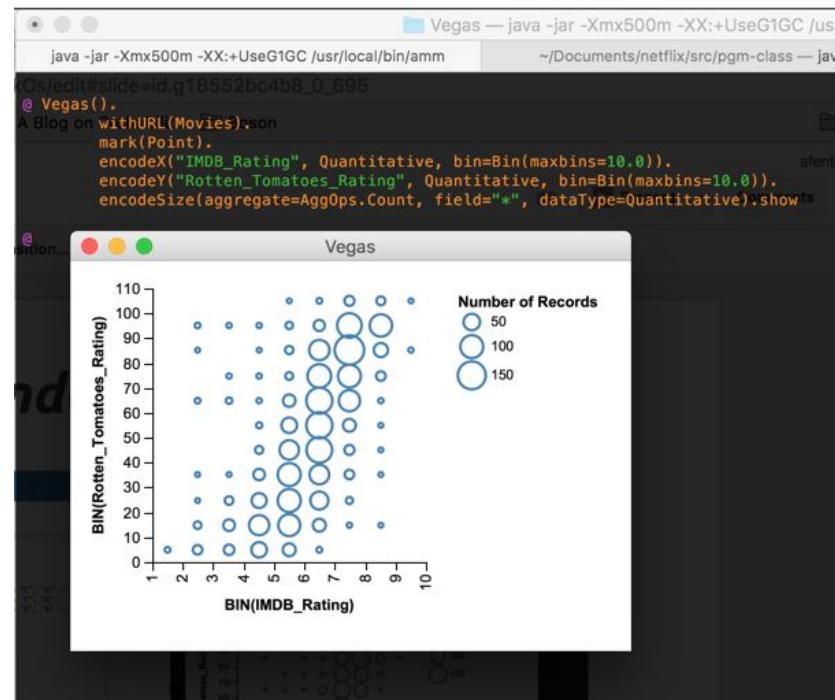
Zepplin Notebook ▾

Vegas Examples

A simple bar chart with embedded data.

```
Vegas("A simple bar chart with embedded data.").  
withData(SeqC  
  Map("a" -> "A", "b" -> 28), Map("a" -> "B", "b" -> 55), Map("a" -> "C", "b" -> 43),  
  Map("a" -> "D", "b" -> 91), Map("a" -> "E", "b" -> 81), Map("a" -> "F", "b" -> 53),  
  Map("a" -> "G", "b" -> 19), Map("a" -> "H", "b" -> 87), Map("a" -> "I", "b" -> 52)  
)).  
encodeX("a", Ordinal).  
encodeY("b", Quantitative).  
mark(Bar).  
show
```

Category	Value
A	28
B	55
C	43
D	91
E	81
F	53
G	19
H	87
I	52



5. *Built-in spark support*

Vegas

```
.withDataFrame(myDataFrame)  
.encodeX("population")  
.encodeY("age")
```

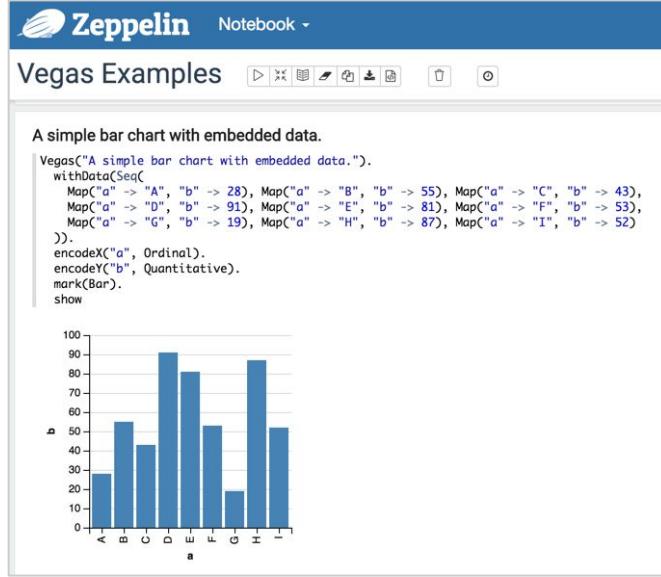
PASS IN DF.

MAPPED COLUMNS

6. *Visual statistics*

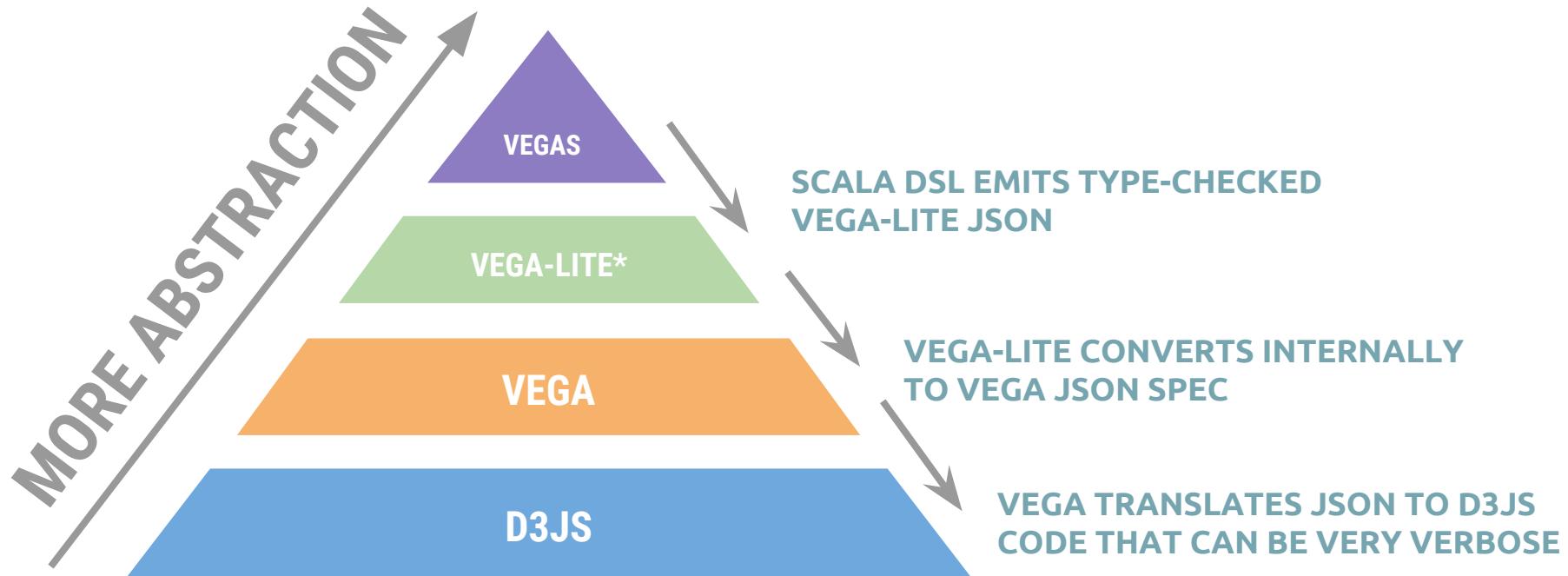
- Advanced Binning
- Sorting
- Scaling
- Custom Transforms
- Time Series
- Aggregation
- Filtering
- Math functions (log, etc)
- Descriptive Statistics

How It Works !



1. Specify in Scala
2. Embed HTML
(iFrame)
3. Render within
iFrame using JS

A SCALA DSL FOR VEGA-LITE



* Vega-Lite

Example 3

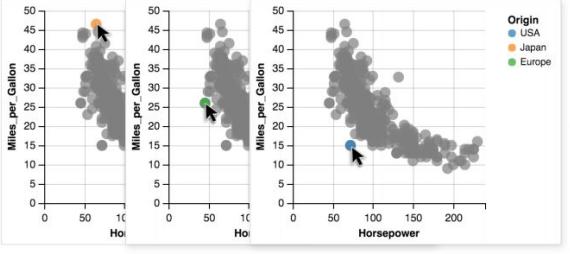
Other Channels + Transforms

What's coming

1. Interactive selections

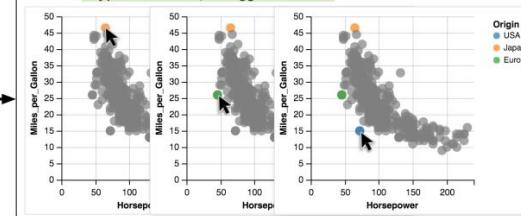
(a) Highlight a single point on click

```
{  
  "data": {"url": "data/cars.json"},  
  "mark": "circle",  
  "select": {  
    "id": {"type": "point"}  
  },  
  "encoding": {  
    "x": {"field": "Horsepower", "type": "Q"},  
    "y": {"field": "MPG", "type": "Q"},  
    "color": [  
      {"if": {"id": "id", "field": "Origin", "type": "N"},  
      {"value": "grey"}  
    ],  
    "size": {"value": 100}  
  }  
}
```

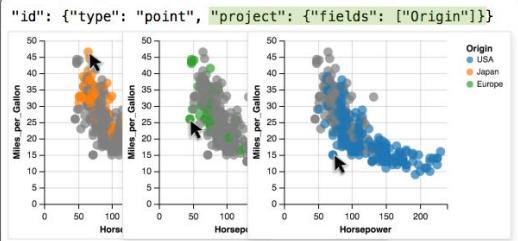


(b) Highlight a list of individual points

```
"id": {"type": "list", "toggle": true}
```

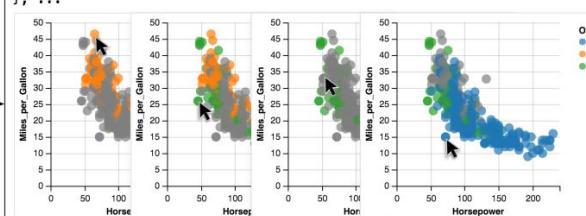


(d) Highlight a single Origin

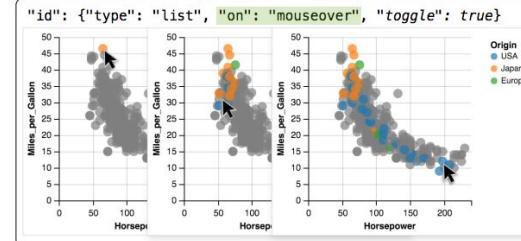


(e) Highlight a list of Origins

```
"select": {  
  "id": {"type": "list", "toggle": true, "project": {"fields": ["Origin"]}}  
}, ...
```



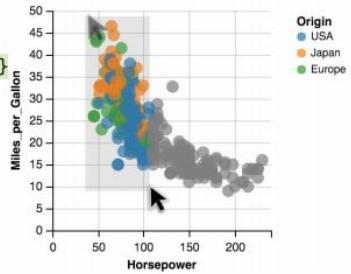
(c) "Paintbrush": highlight multiple points on hover



2. Selections transforms

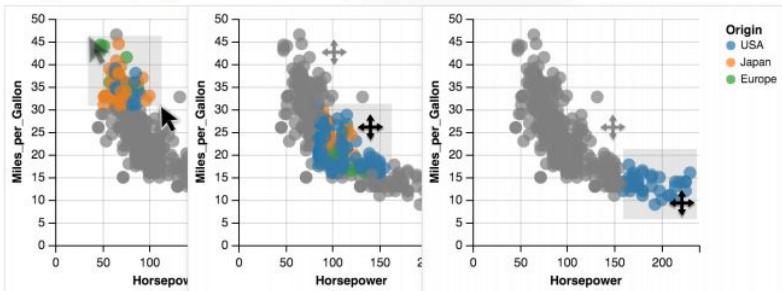
(a) Rectangular brush

```
"select": {  
  "region": {"type": "interval"}  
},  
...  
  "color": [  
    {"if": "region", ...}  
  ]  
...
```



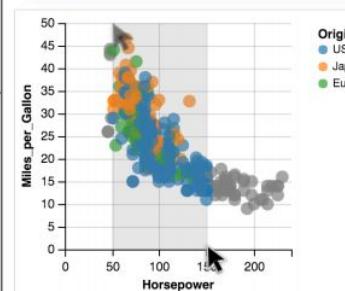
(b) Moving the brush

```
"region": {"type": "interval", "translate": true}
```

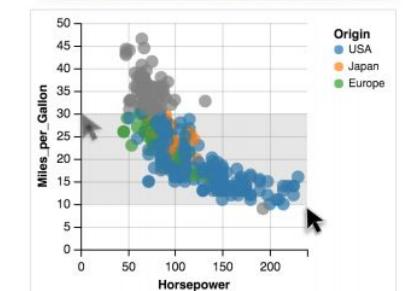


(c) Single-dimension brush

```
"region": {"type": "interval",  
  "project": {"channels": ["x"]}}
```



```
"region": {"type": "interval",  
  "project": {"channels": ["y"]}}
```



Contributors



Aish



DB



Roger



Sudeep



Jeremy

NETFLIX

Thank you.

NETFLIX



The missing Matplotlib for Scala/Spark

<http://vegas-viz.org>

@NetflixResearch
@rogermenezes @dbtsai

NETFLIX