

Alexander Lex

Mastodon: @alexlex@vis.social

<http://alexander-lex.net>

Provenance as a Bridge Between Data Analysis Modalities



visualization
design lab



ACKNOWLEDGEMENTS



Kiran Gadhave



Zach Cutler

Also: Jochen Görtler, Carolina Nobre, Oliver Deussen, Miriah Meyer, Jeff Phillips, Marc Streit, Nils Gehlenborg



**CAREER: Enabling Reproducibility
of Interactive Visual Data Analysis**

WHAT IS PROVENANCE?

Provenance (from the French *provenir*, 'to come from/forth') is the **chronology** of the ownership, custody or location of a historical object. [Wikipedia]

In CS: a log, a record of everything that lead to a state

PURPOSES FOR PROVENANCE

Purposes for Provenance	
Recall	Maintaining or recovering memory and awareness of the current and previous states of the analysis
Replication	Reproducing the steps or workflow of a previous analysis
Action recovery	Maintaining the action history that allows undo/redo operations and branching actions during analysis
Collaborative communication	Communicating and sharing data, information, and ideas with others who are conducting the same analysis
Presentation	Communicating the insights or progression of the analysis with those who are not directly involved with the analysis themselves, such as general public, upper levels of management, or analysts focusing on other areas
Meta-analysis	Reviewing the analytic processes themselves in order to understand and improve aspects of the analysis (such as process efficiency, training efficiency, or analytic strategies)

[Ragan et al 2015]

Re-Application

Convert the user interactions into executable scripts.

[Xu et al 2020]

**PROVENANCE AS A
BRIDGE BETWEEN
DATA ANALYSIS
MODALITIES**

What are Modalities?

Interactive Visualization Systems

Code / Scripting

INTERACTIVE VISUALIZATION

Intuitive
Easy to use
Uses human perceptual capabilities



INTERACTIVE VISUALIZATION: DOWNSIDES

Limited Expressivity

Some operations are difficult

e.g., conditional queries..

Not reusable

need to redo analysis when data changes

Not reproducible



CODE / SCRIPTING

Flexible and powerful

you basically can do anything

Reusable

if your data changes, re-run

Reproducible

everything is documented

```
1 # Keep this cell
2 avy_df = pd.read_csv('./avalanches.csv')
3
4 # Remove NaN coordinates
5 avy_df = avy_df[avy_df['Coordinates']!=avy_df['Coordinates']]
6
7 # Split into latitude & longitude
8 avy_df[['lat', 'lon']] = avy_df['Coordinates'].str.split(',', expand=True)
9
10 # Remove values outside of Utah bounds
11 avy_df = avy_df[(36 < avy_df['lat'].astype('float')) & (avy_df['lat'].astype('float') < 42)]
12 avy_df = avy_df[(-114 < avy_df['lon'].astype('float')) & (avy_df['lon'].astype('float') < -108)]
13
14 # Keep columns we need
15 avy_df = avy_df[['Date', 'Region', 'Trigger', 'lat', 'lon']]
16 avy_df.head()
```


CODE / SCRIPTING: DOWNSIDES

It's hard

requires extensive training

reading documentation

not discoverable

Not everyone can do it

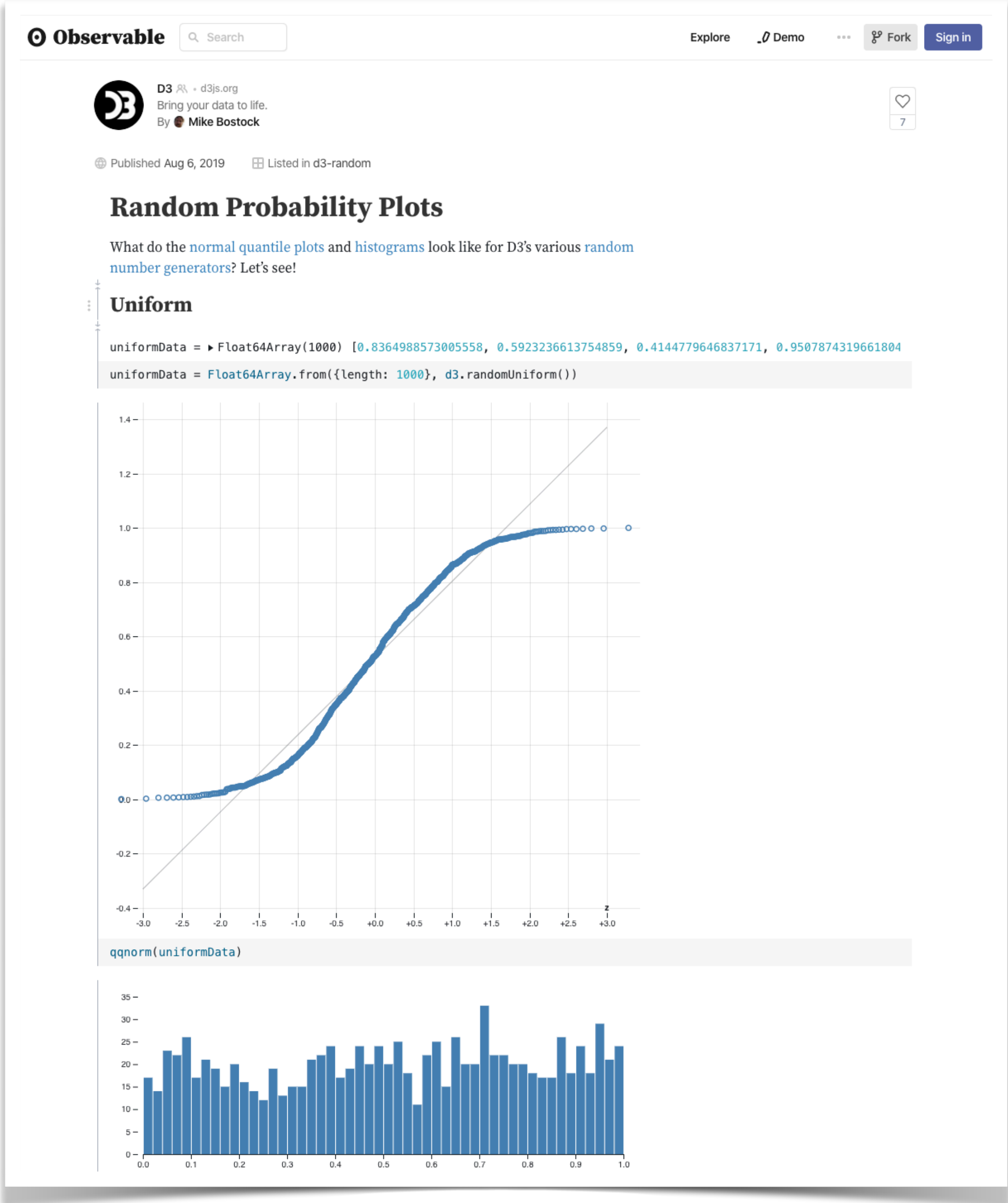
It's time consuming

Some operations are difficult

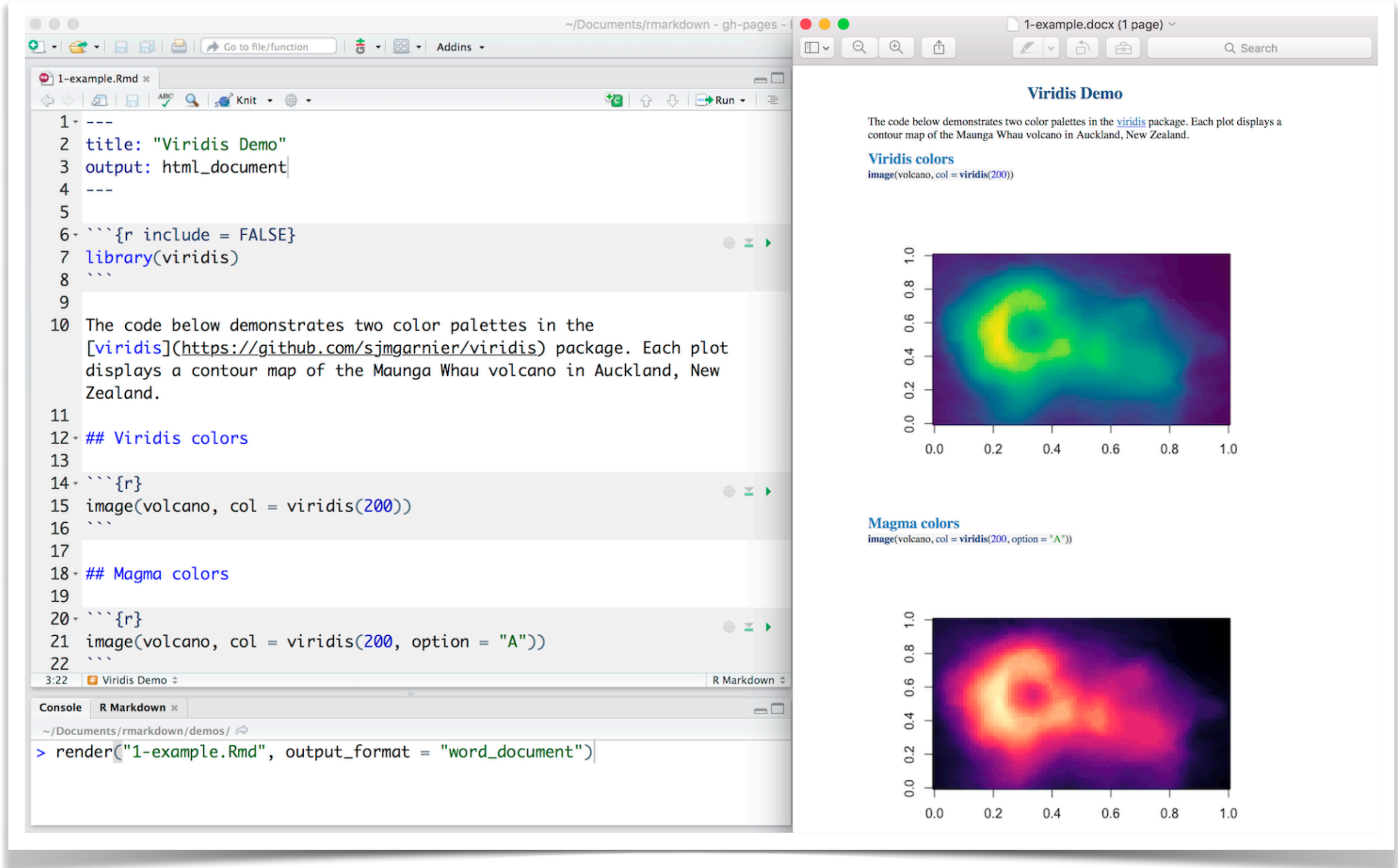
e.g., labeling data points

```
1 # Keep this cell
2 avy_df = pd.read_csv('./avalanches.csv')
3
4 # Remove NaN coordinates
5 avy_df = avy_df[avy_df['Coordinates']!=avy_df['Coordinates']]
6
7 # Split into latitude & longitude
8 avy_df[['lat', 'lon']] = avy_df['Coordinates'].str.split(',', expand=True)
9
10 # Remove values outside of Utah bounds
11 avy_df = avy_df[(36 < avy_df['lat'].astype('float')) & (avy_df['lat'].astype('float') < 42)]
12 avy_df = avy_df[(-114 < avy_df['lon'].astype('float')) & (avy_df['lon'].astype('float') < -108)]
13
14 # Keep columns we need
15 avy_df = avy_df[['Date', 'Region', 'Trigger', 'lat', 'lon']]
16 avy_df.head()
```

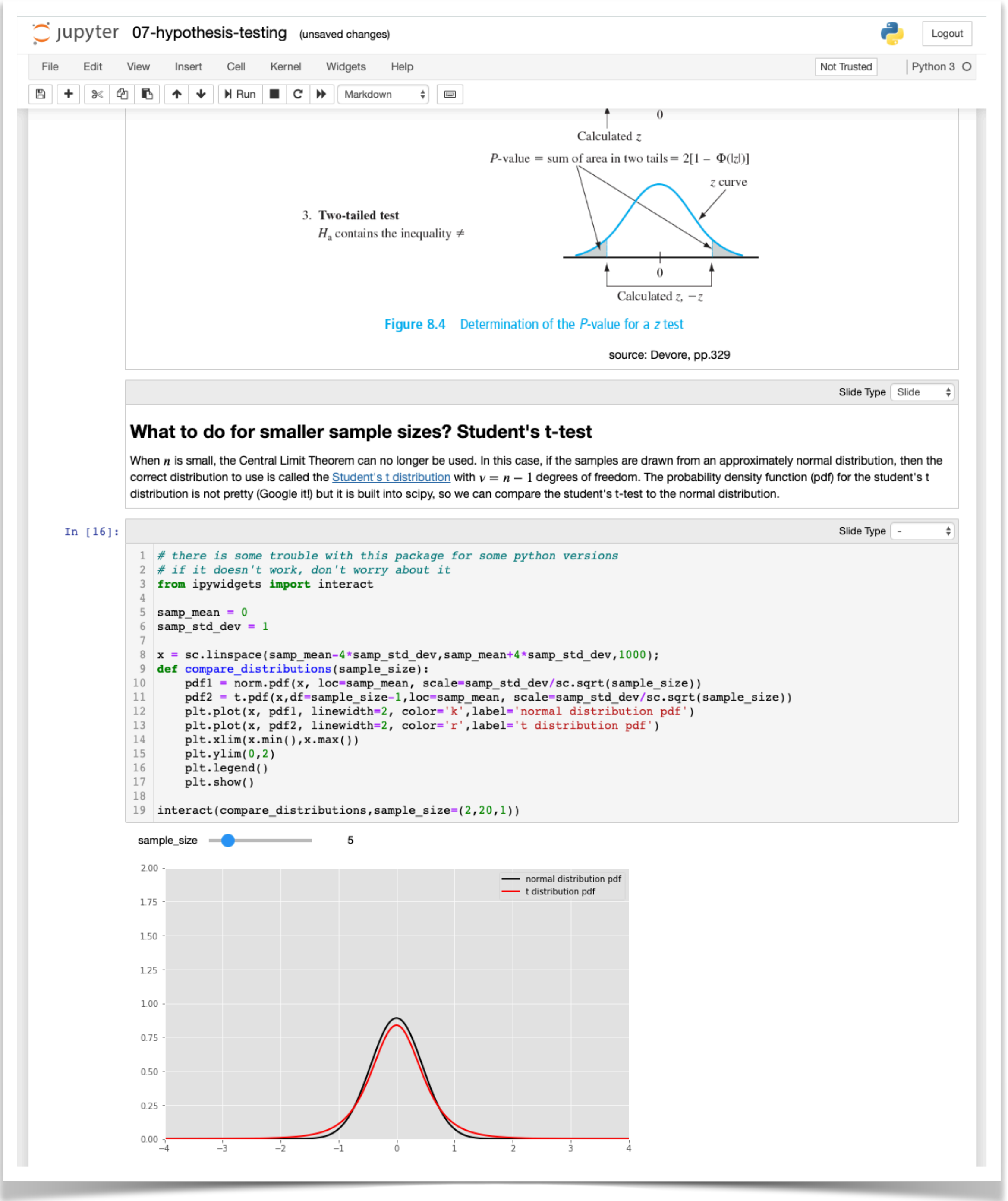
COMPUTATIONAL NOTEBOOKS: A MIDDLE GROUND?



Observable



R Markdown



Jupyter Notebooks

COMPUTATIONAL NOTEBOOKS: A MIDDLE GROUND?

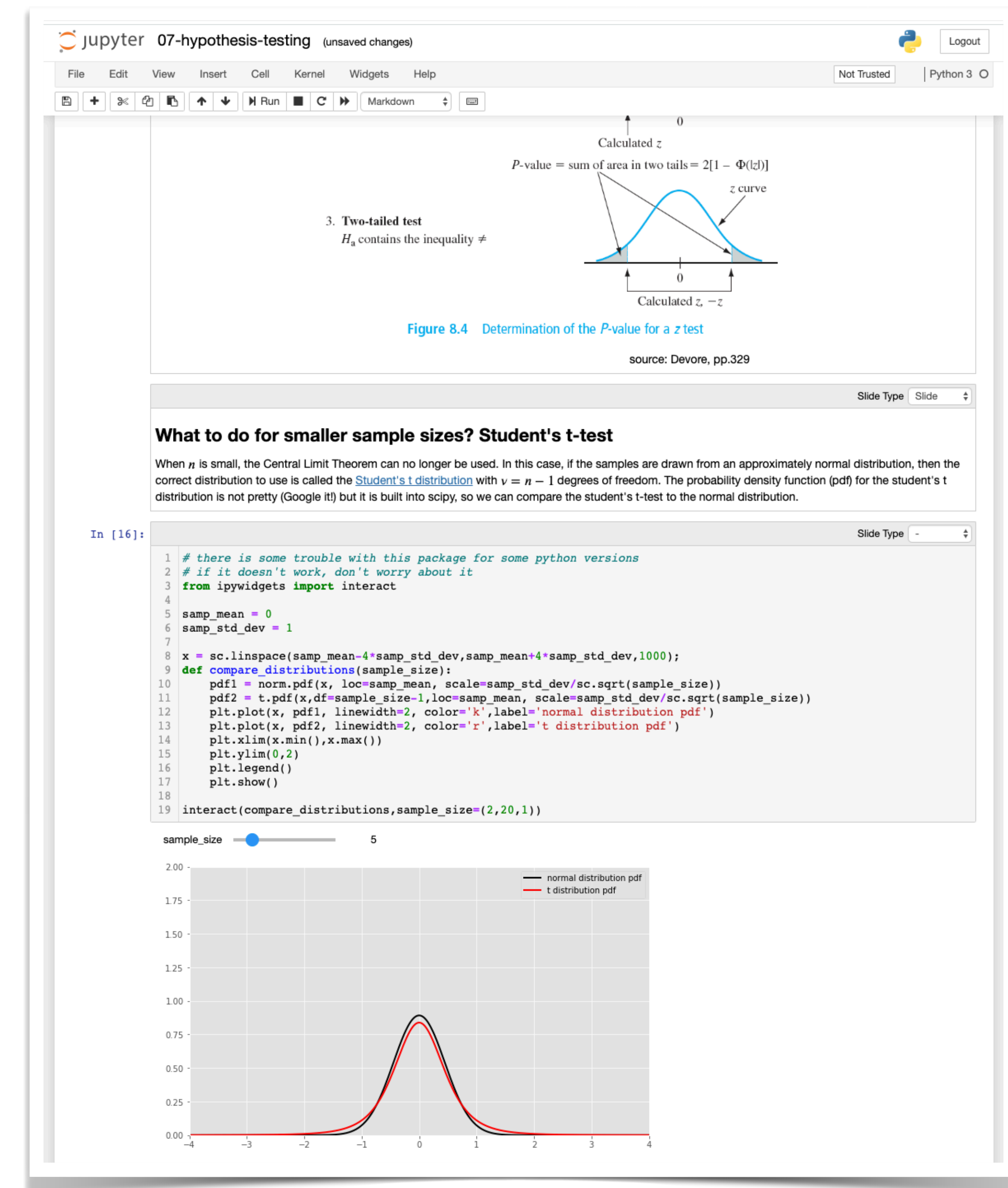
Yes!

Afford both scripting and interactive visualization

But visualizations are a dead end

can't "use" interaction in code

e.g., changing a label, or filtering a value

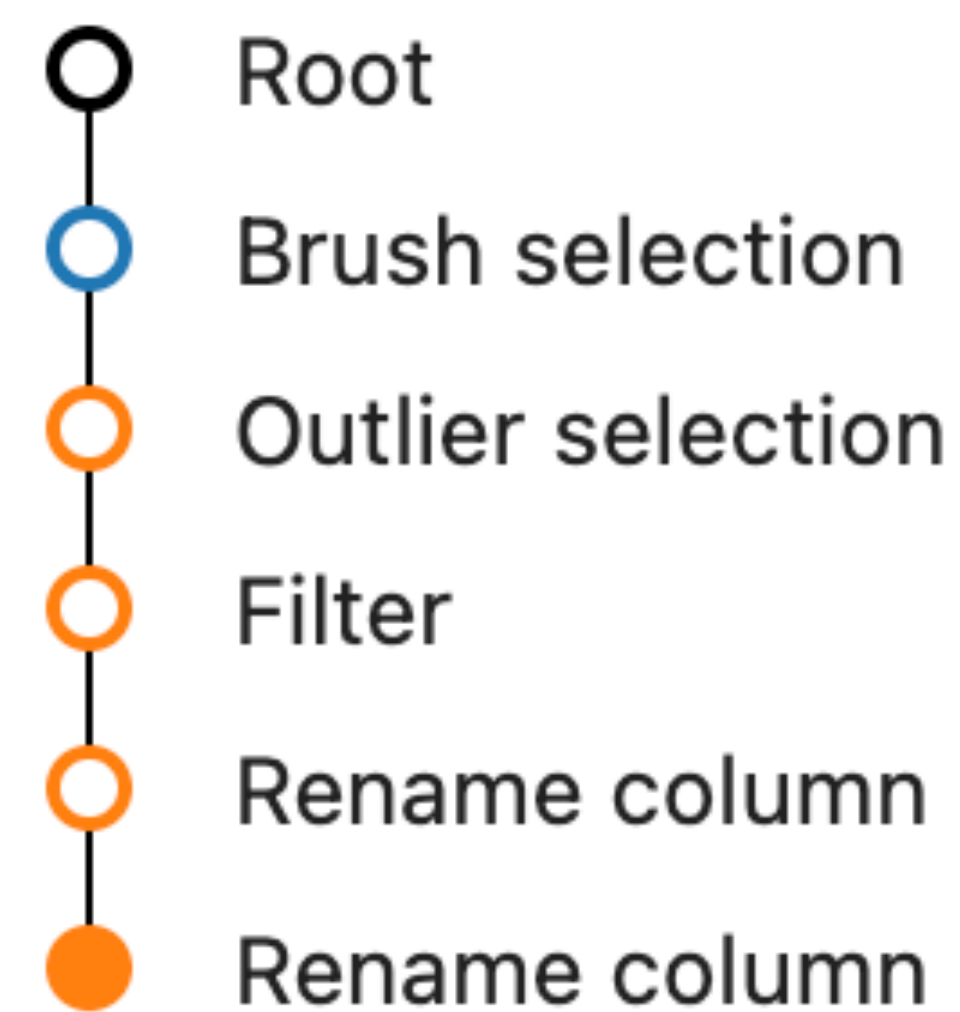


Jupyter Notebooks

THESIS: BRIDGING BETWEEN CODE AND VIS IS USEFUL

Use the best tool for each job

THE ROLE OF PROVENANCE IN BRIDGING BETWEEN CODE AND VISUALIZATIONS



Provenance: record of **actions that lead to a state**

> Design **actions** such that (some) **map to code operations**

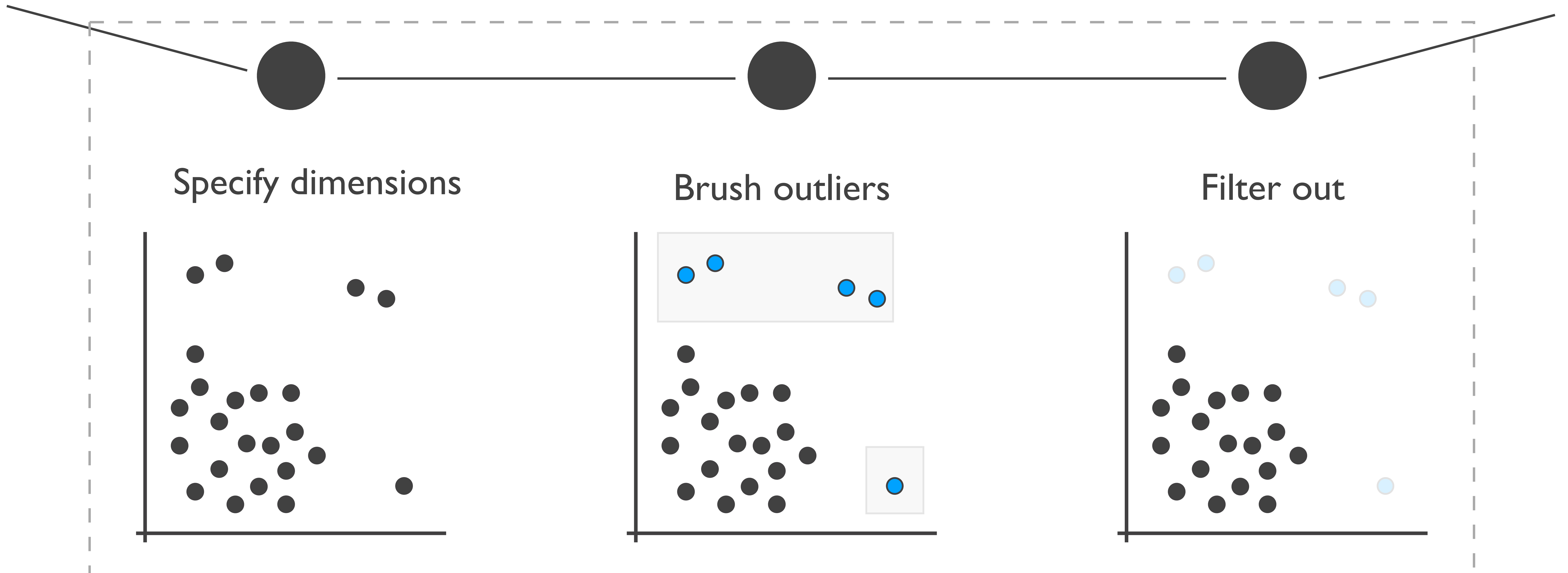
How?

Translate your actions to code

OR

Make code understand actions

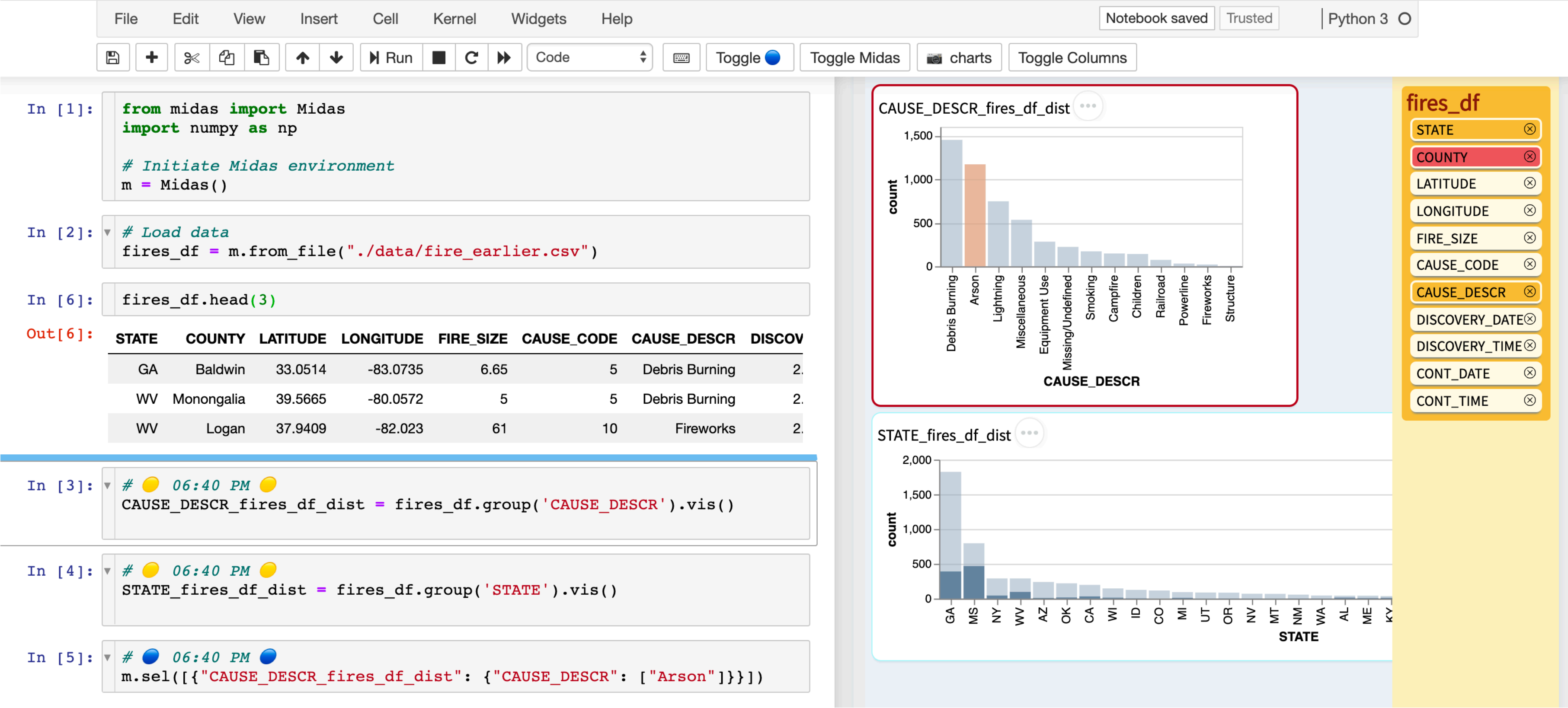
INTERACTIVE WORKFLOWS



“Filter Outliers” Workflow

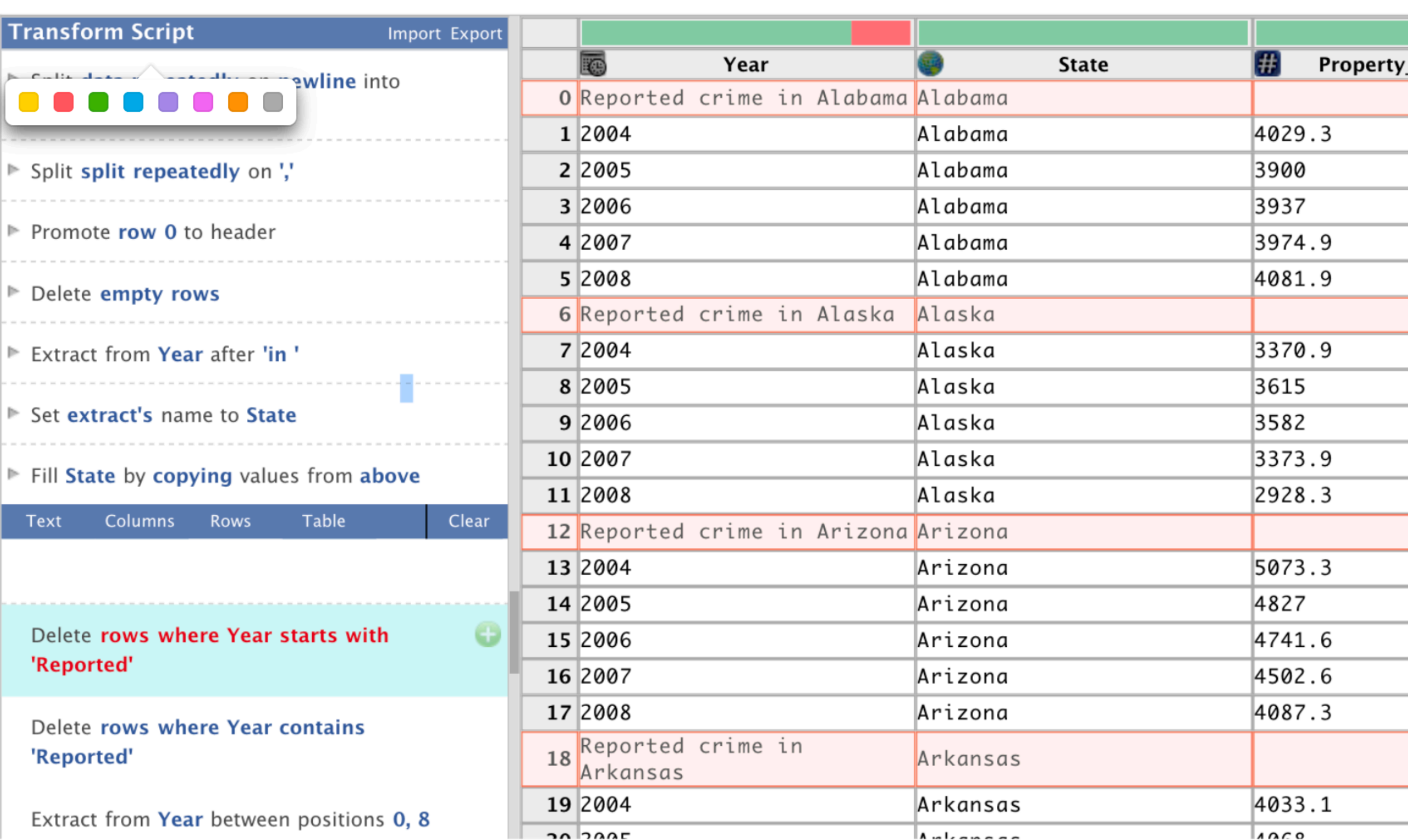
NOTABLE EXAMPLES FOR CODE SYNTHESIS

B2



[Wu et al. 2020]

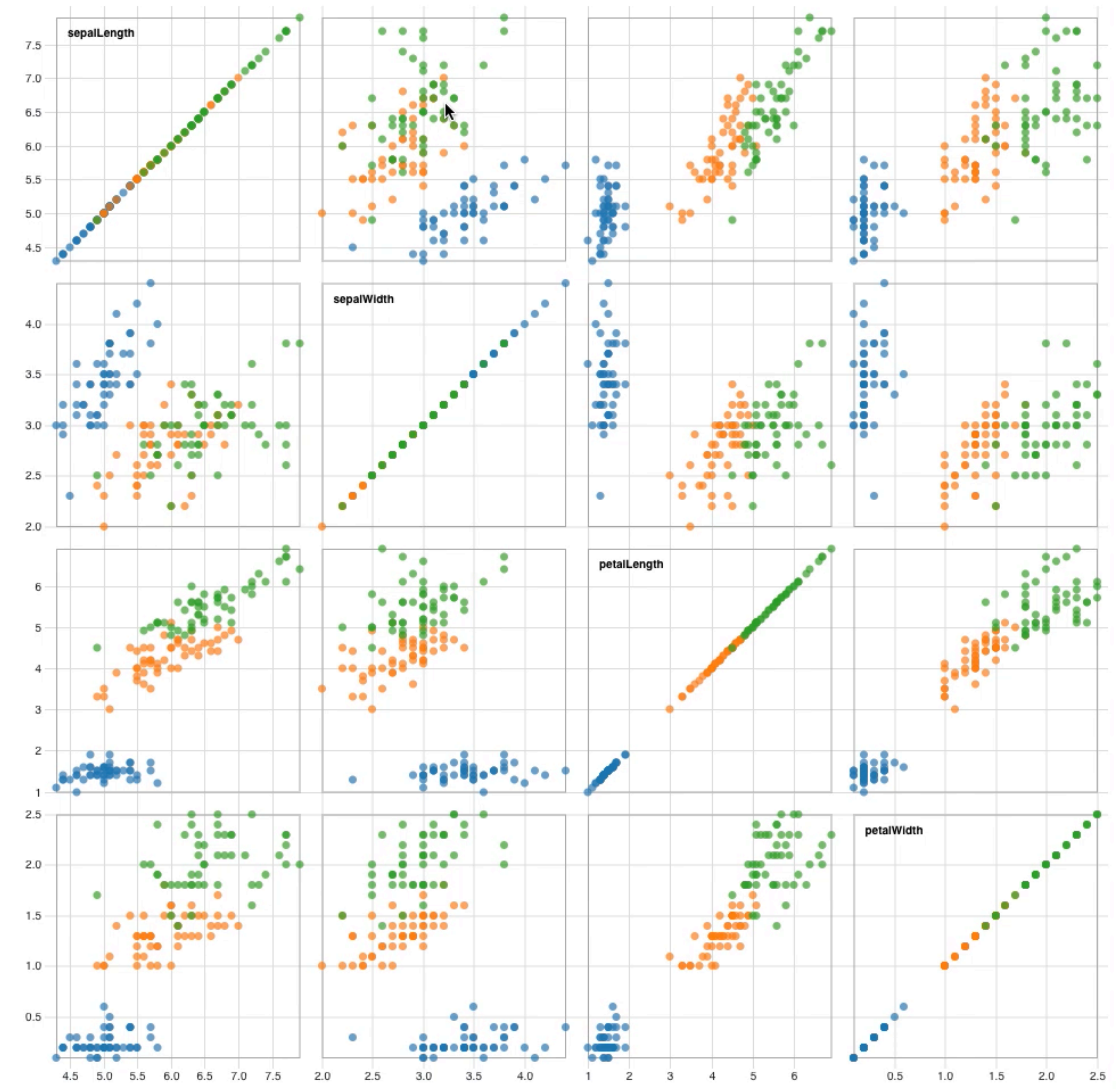
Wrangler



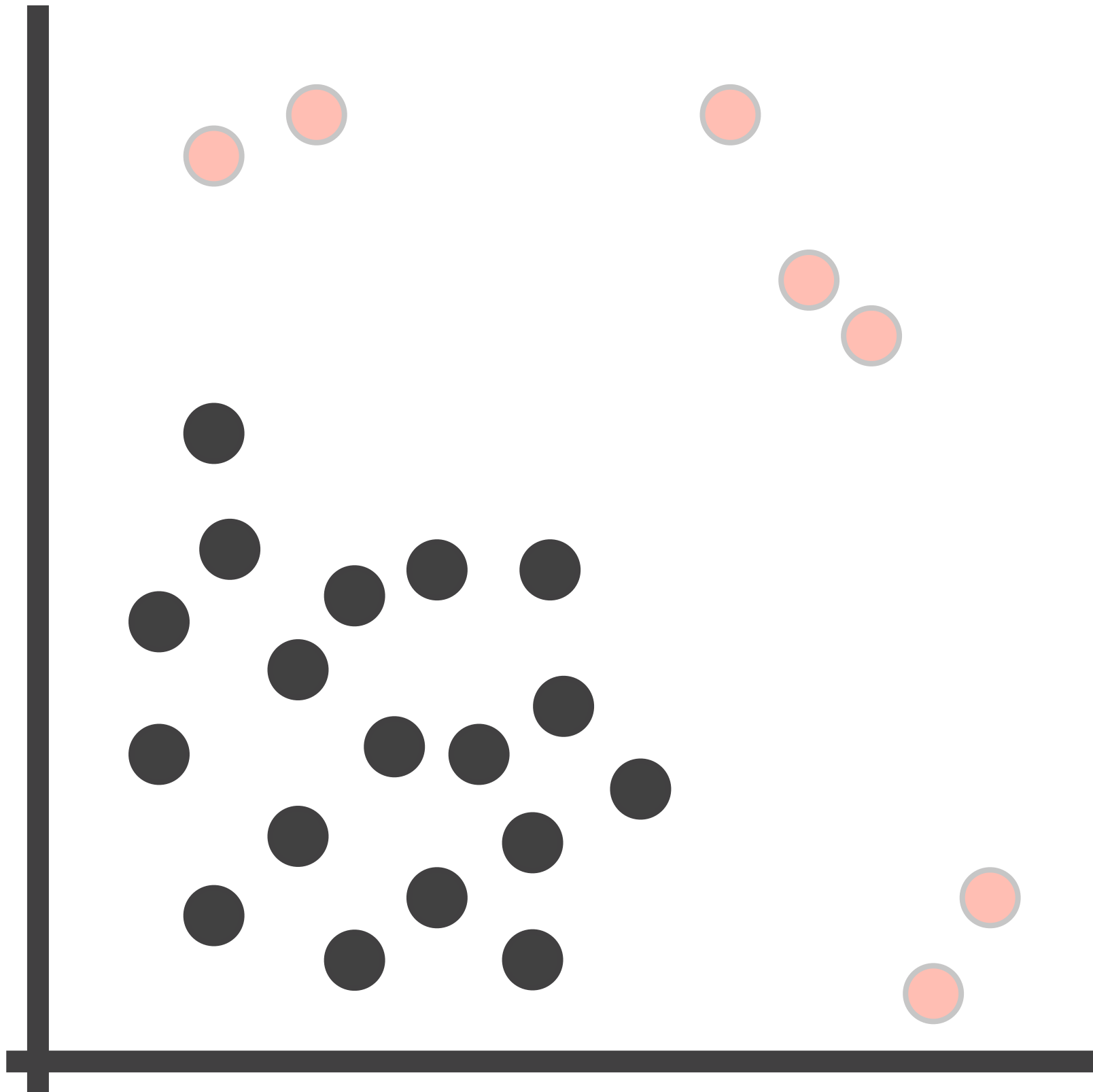
[Kandel et al. 2011]

**CHALLENGE: MAKING
WORKFLOWS ROBUST**

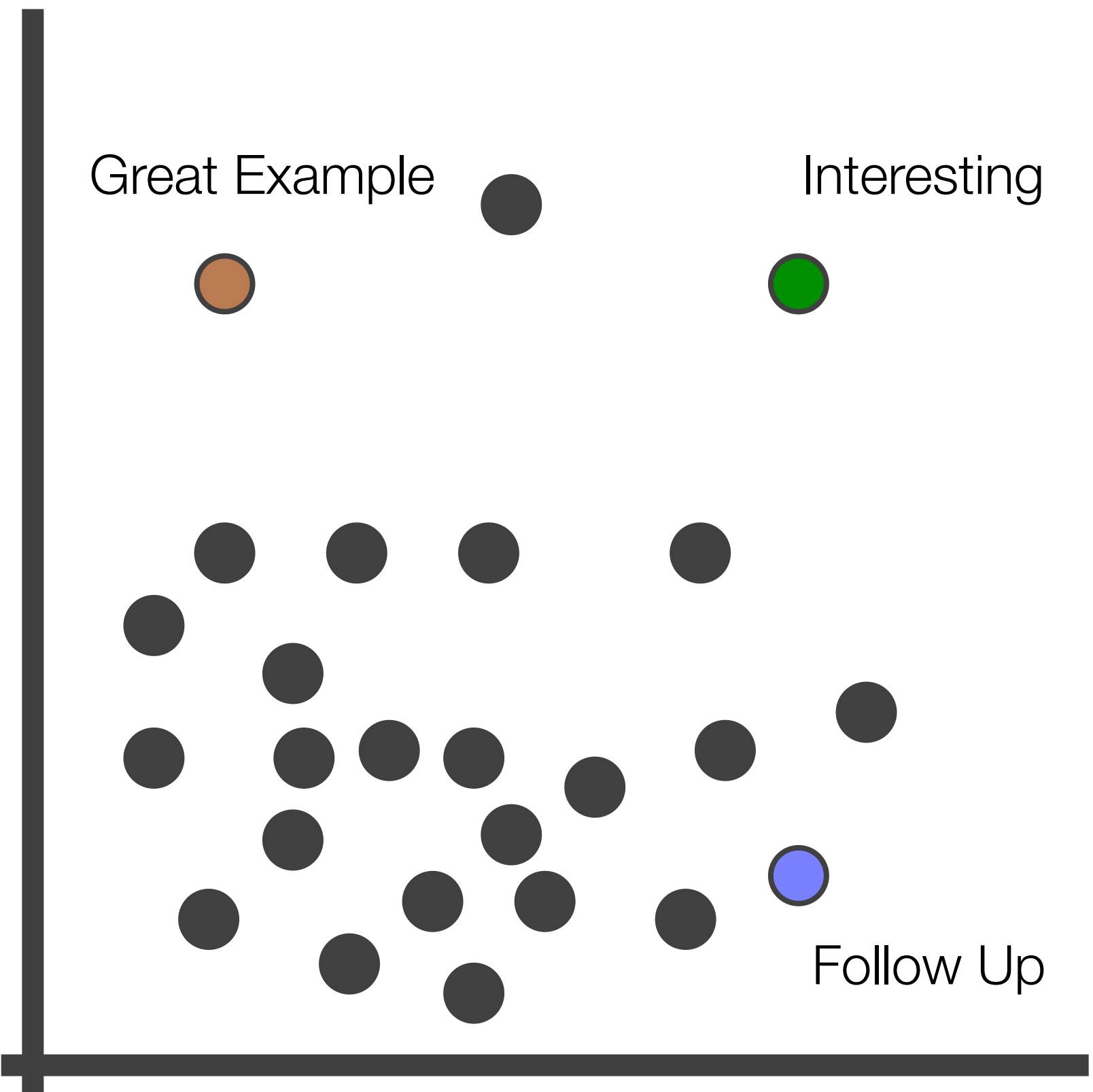
WHAT ARE SELECTIONS?



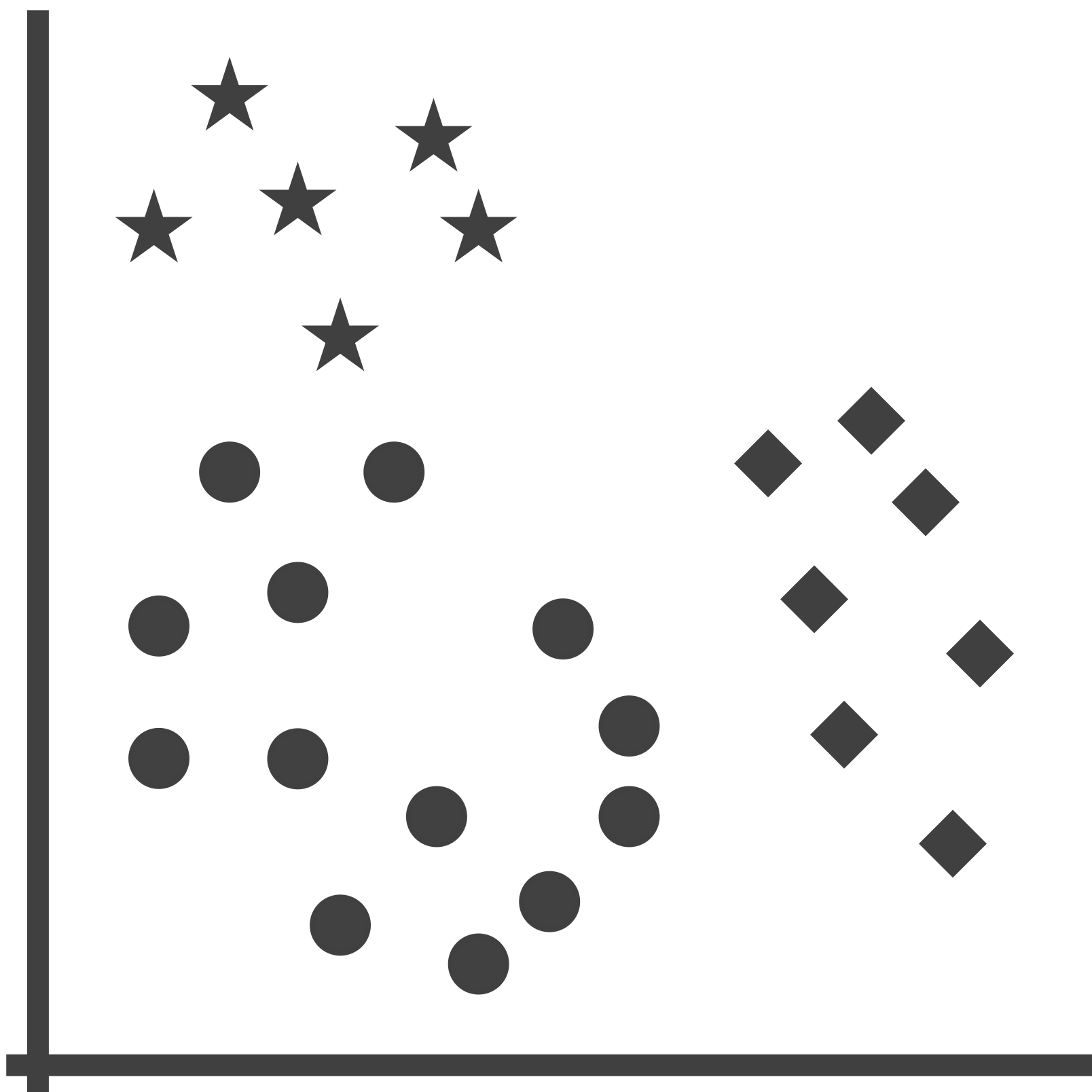
FROM SELECTIONS TO ADVANCED OPERATIONS



Filter

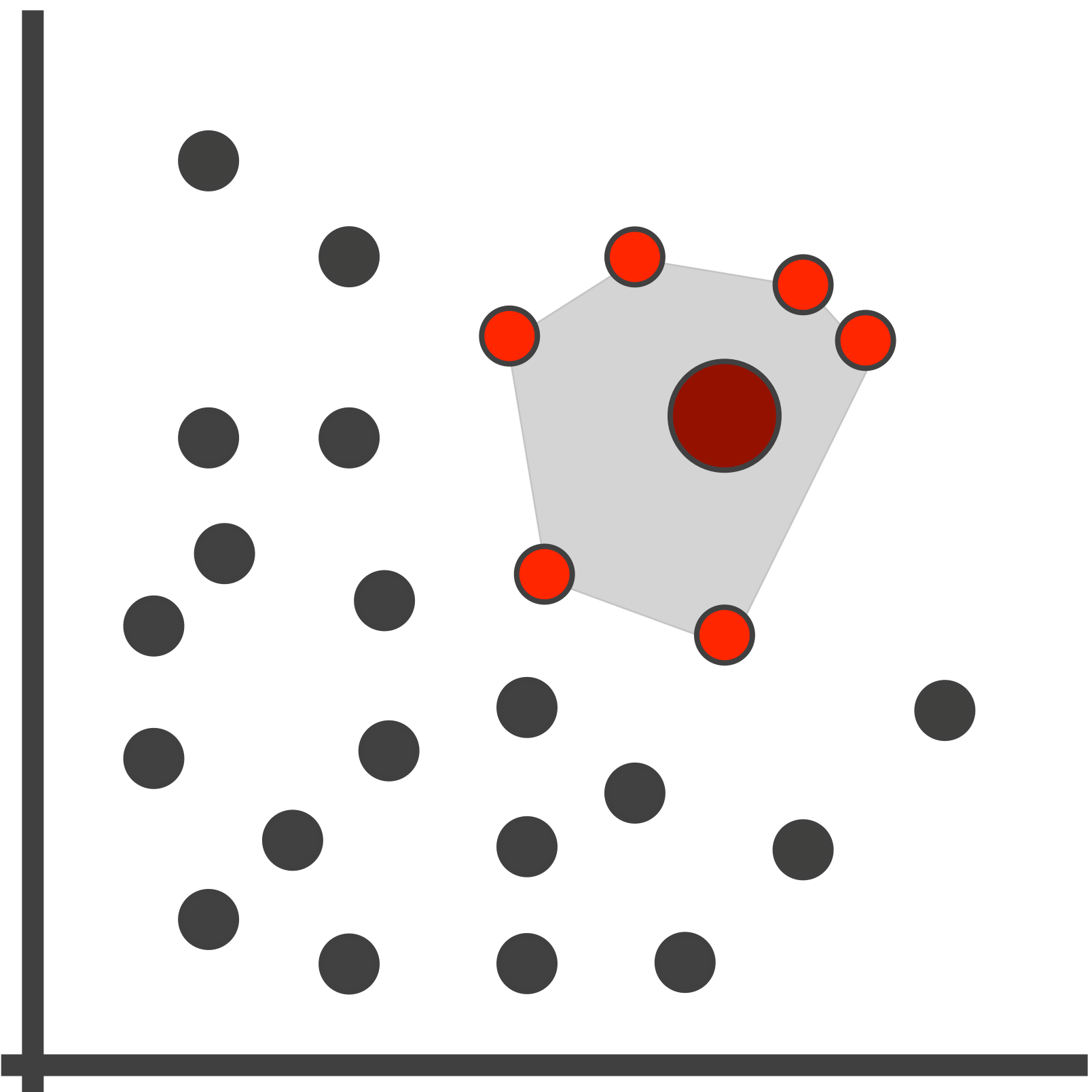


Label



Categorize

- Unassigned
- ◆ Category A
- ★ Category B



Aggregate

ROBUST SELECTIONS UNDERPIN MOST INTERACTIVE OPERATIONS

ROBUST SELECTIONS

Less Robust



More Robust

ID Based Selection:

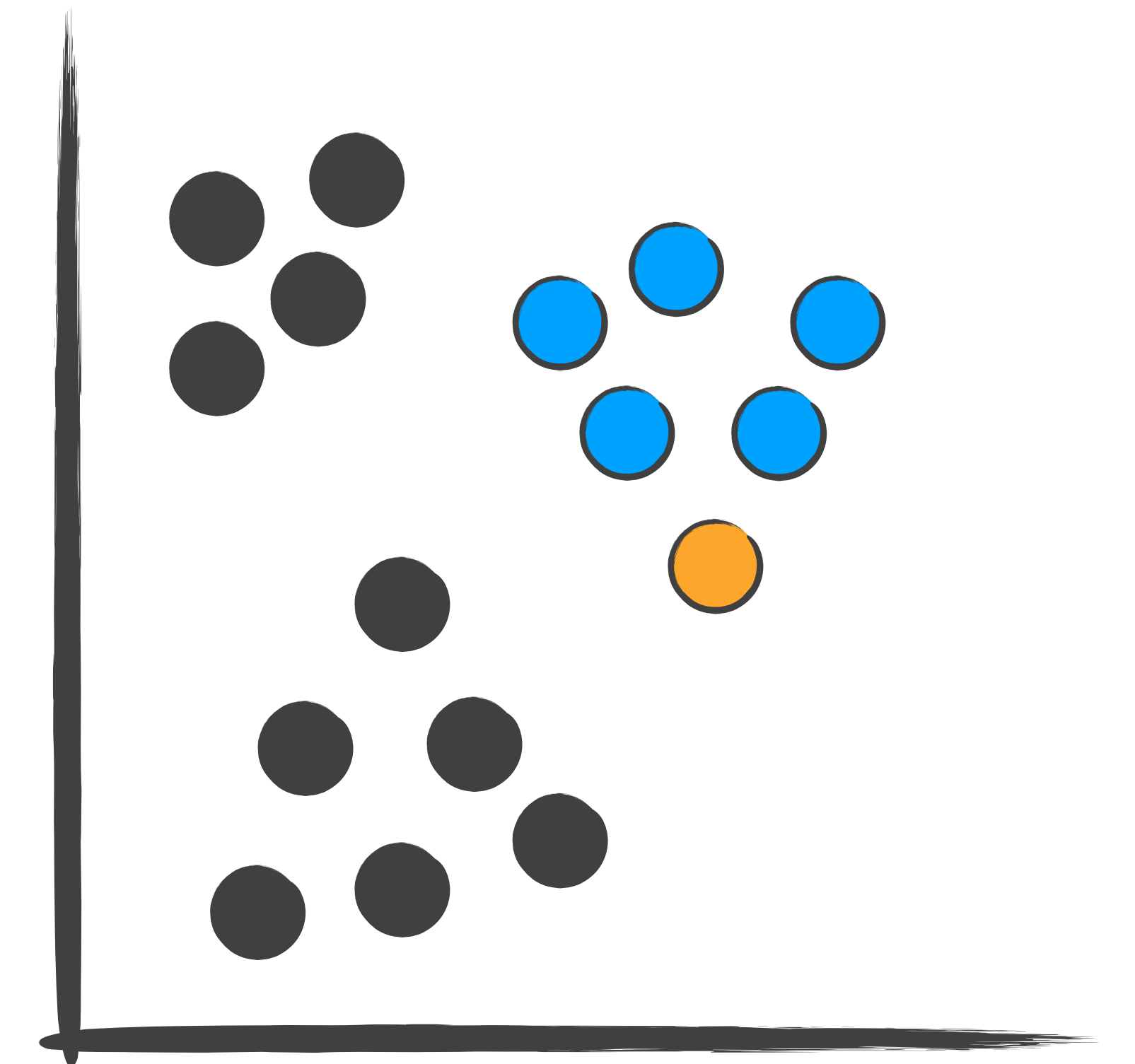
Elements 7, 9, 13, 18, 22

Region-Based Selection:

Elements that are > 1.5 in x
and > 2 in y

Semantic Selection:

Elements in K-Means cluster centered at $[2, 3]$



Meaningful, higher level concept:

improves reproducibility

Robust to changes and updates in dataset:

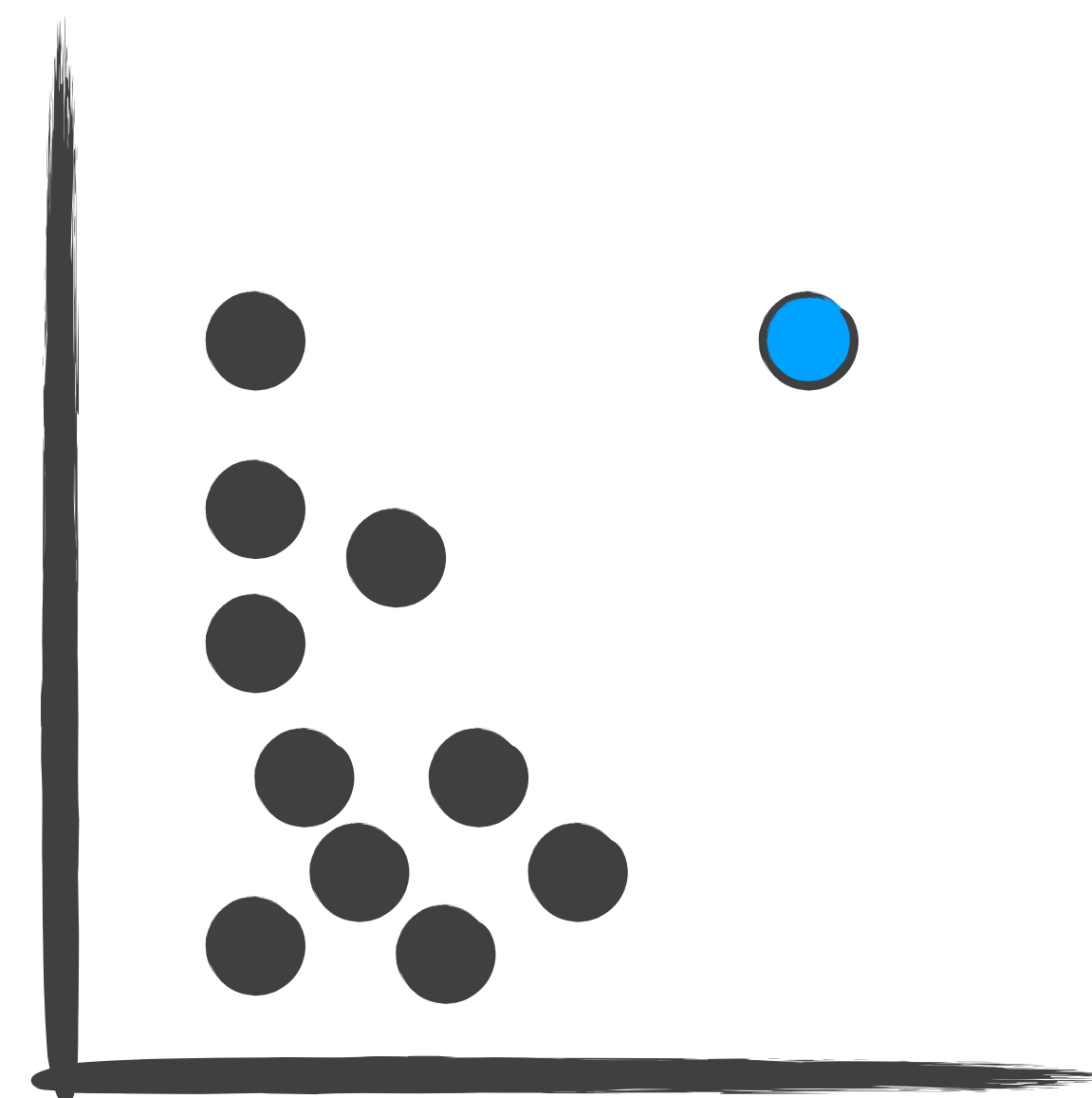
enables re-usability

LEVERAGING INTENT FOR ROBUST SELECTIONS

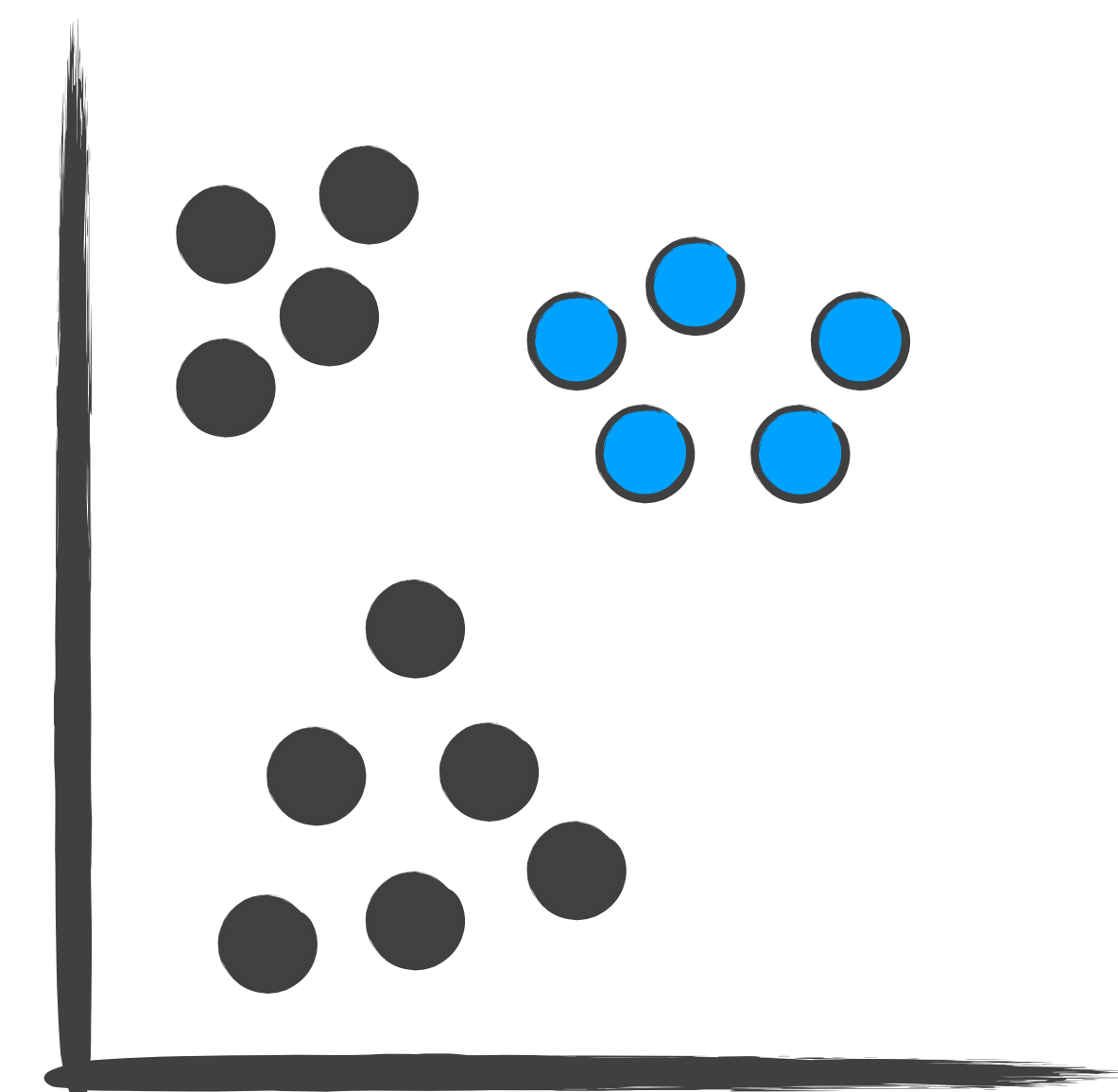
Intent is the user's reason for selecting in a visualization.

Domain Specific Intent: Capture through Annotation

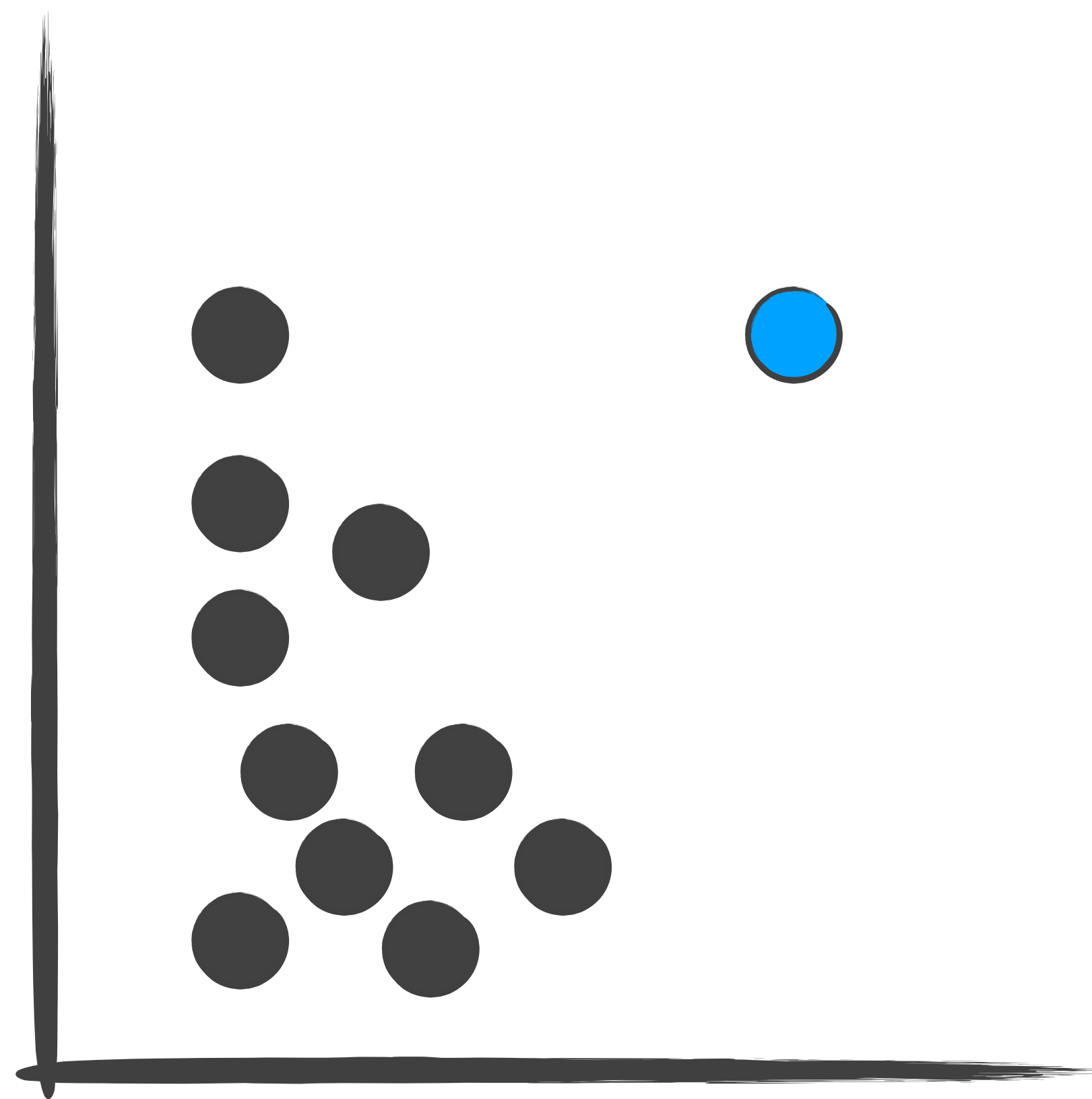
Pattern-Based Intent: Capture Automatically



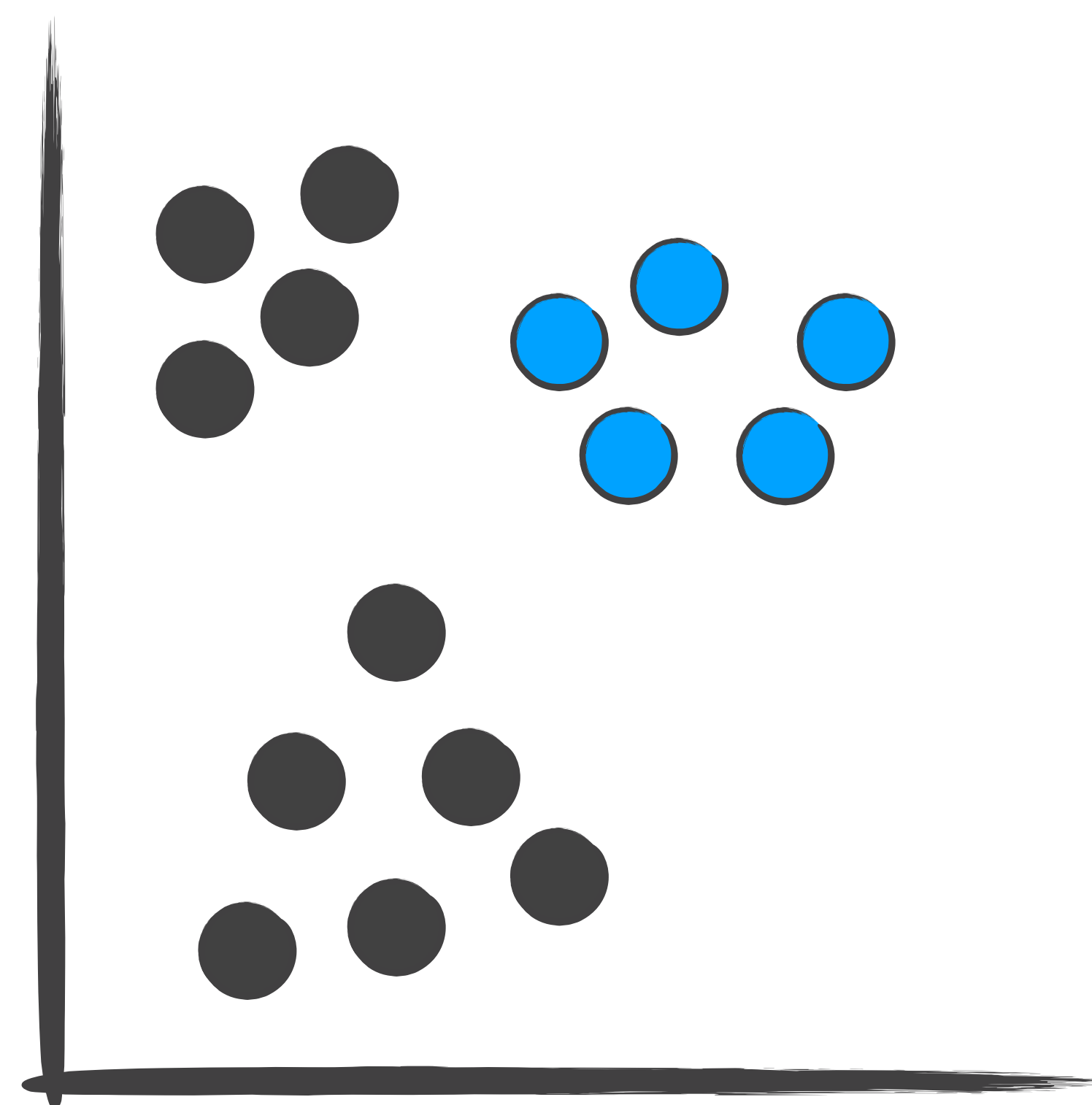
Outlier



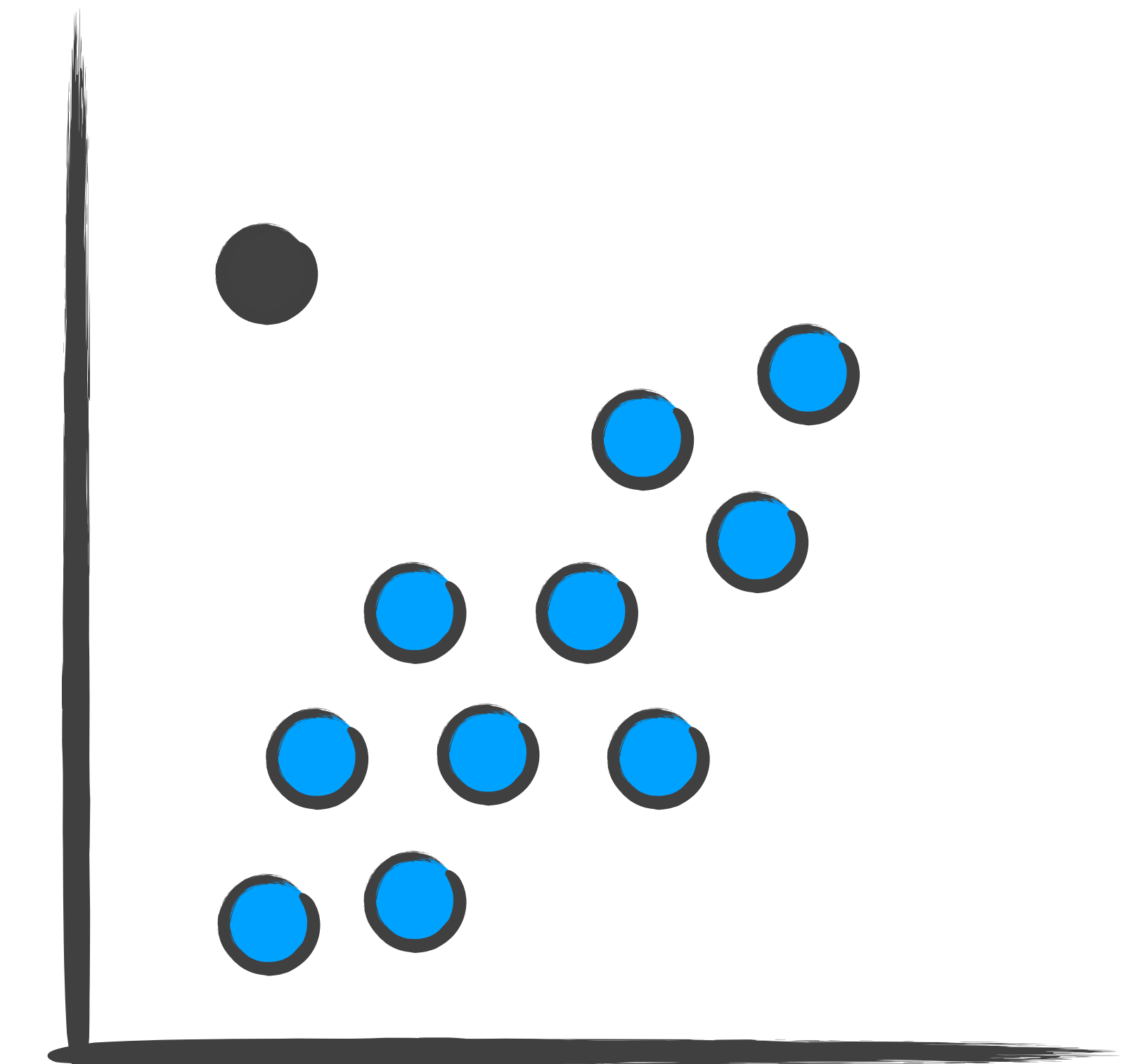
Clusters



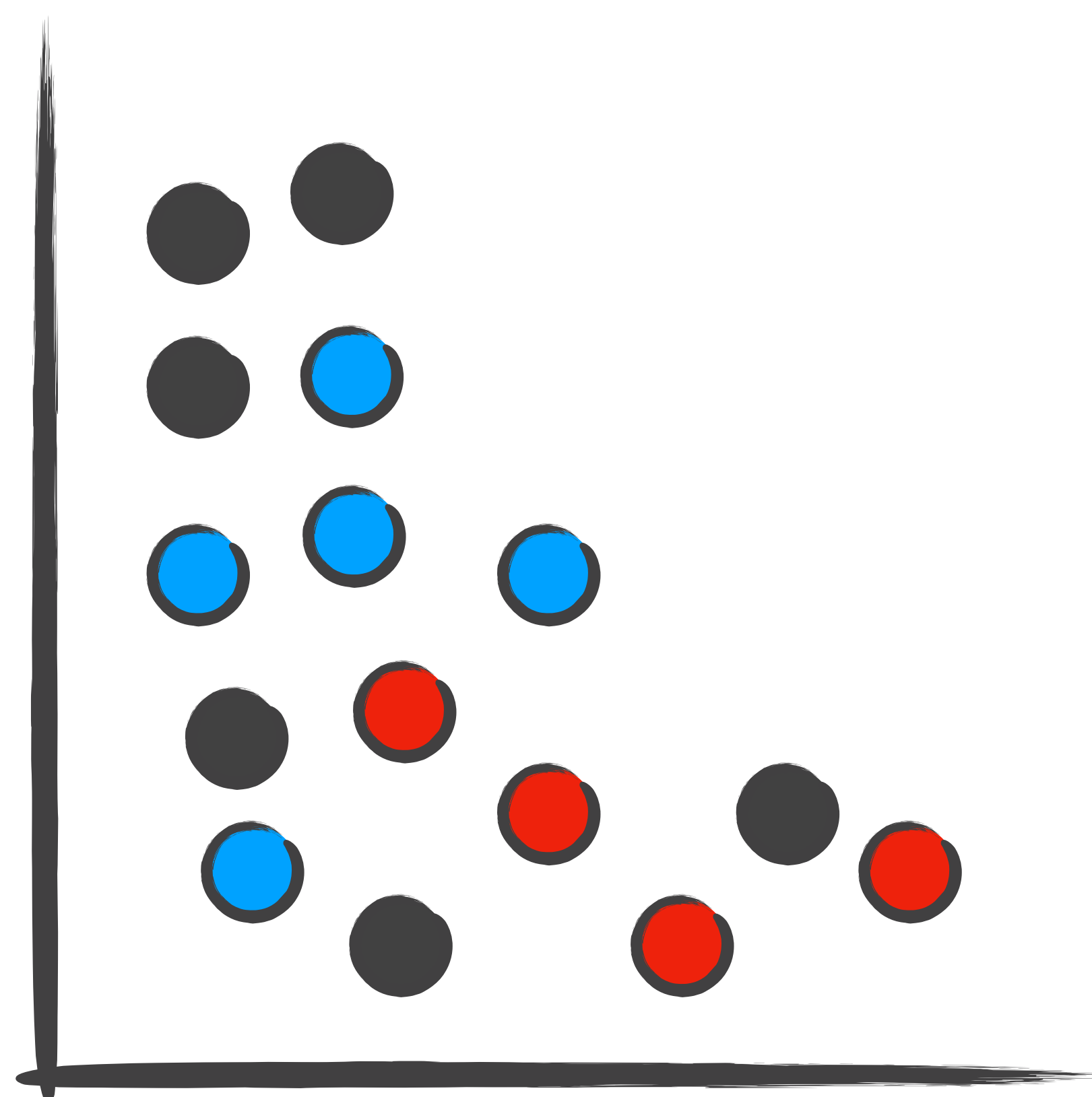
Outlier



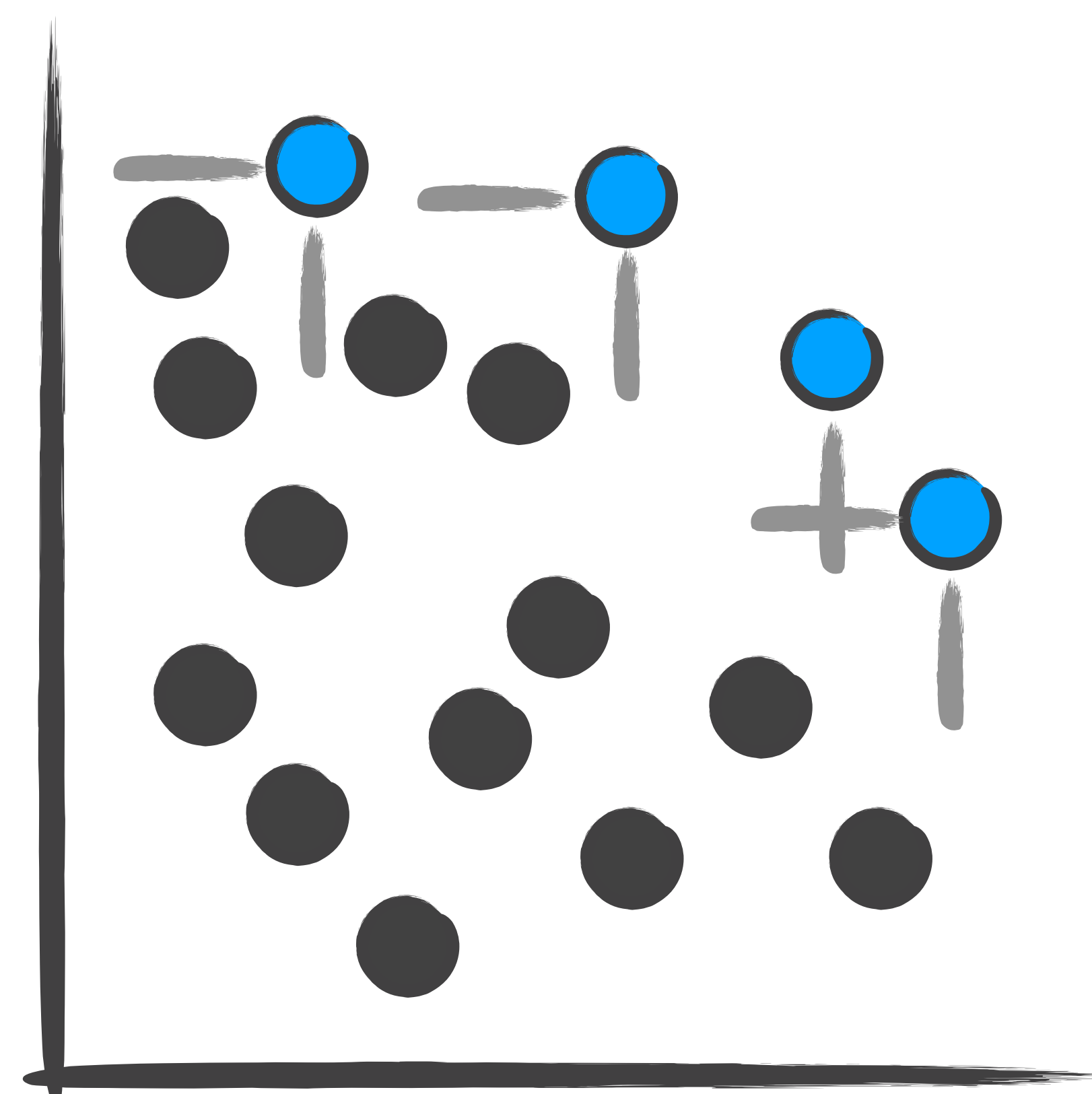
Clusters



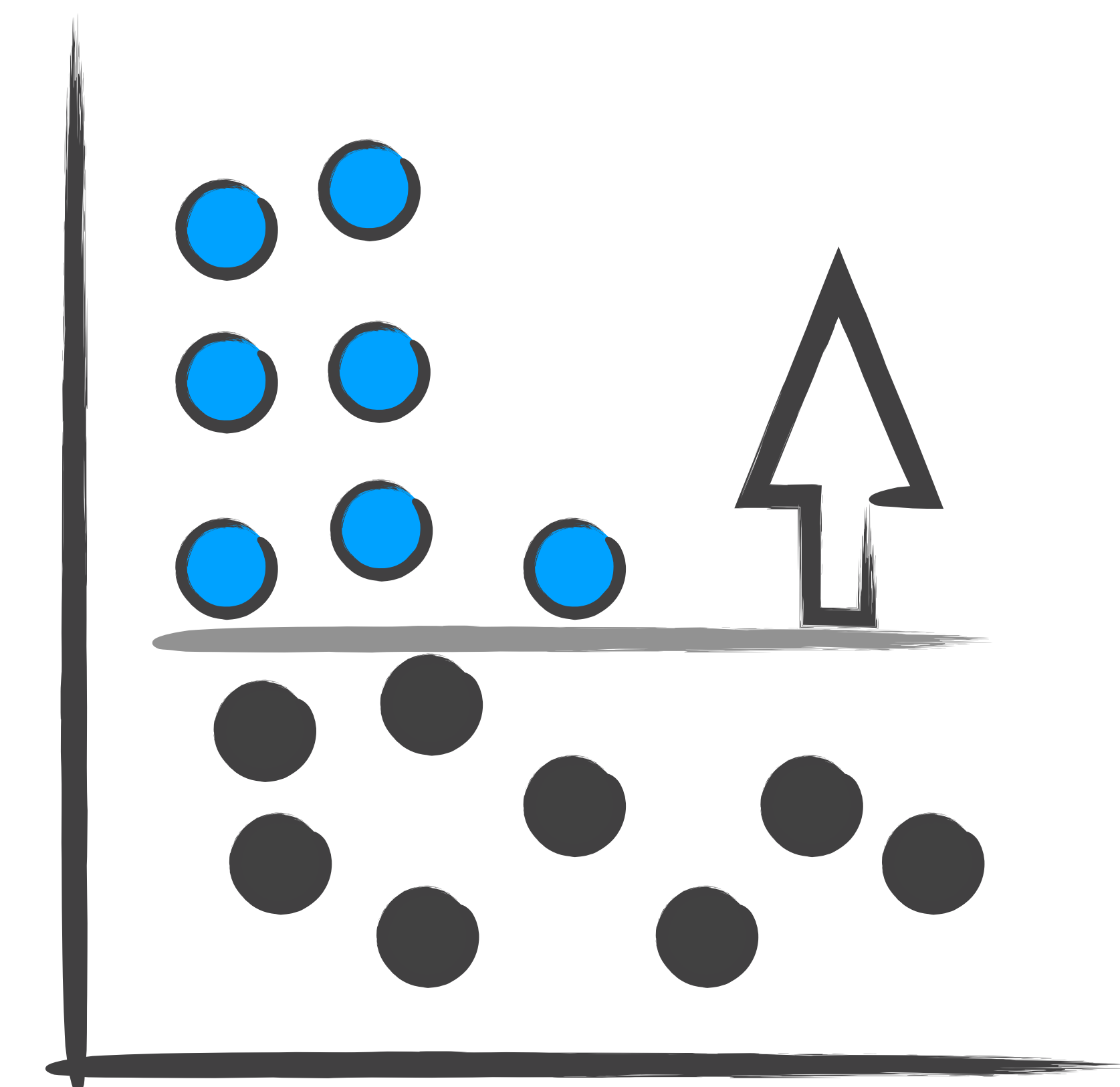
Correlation



Categories

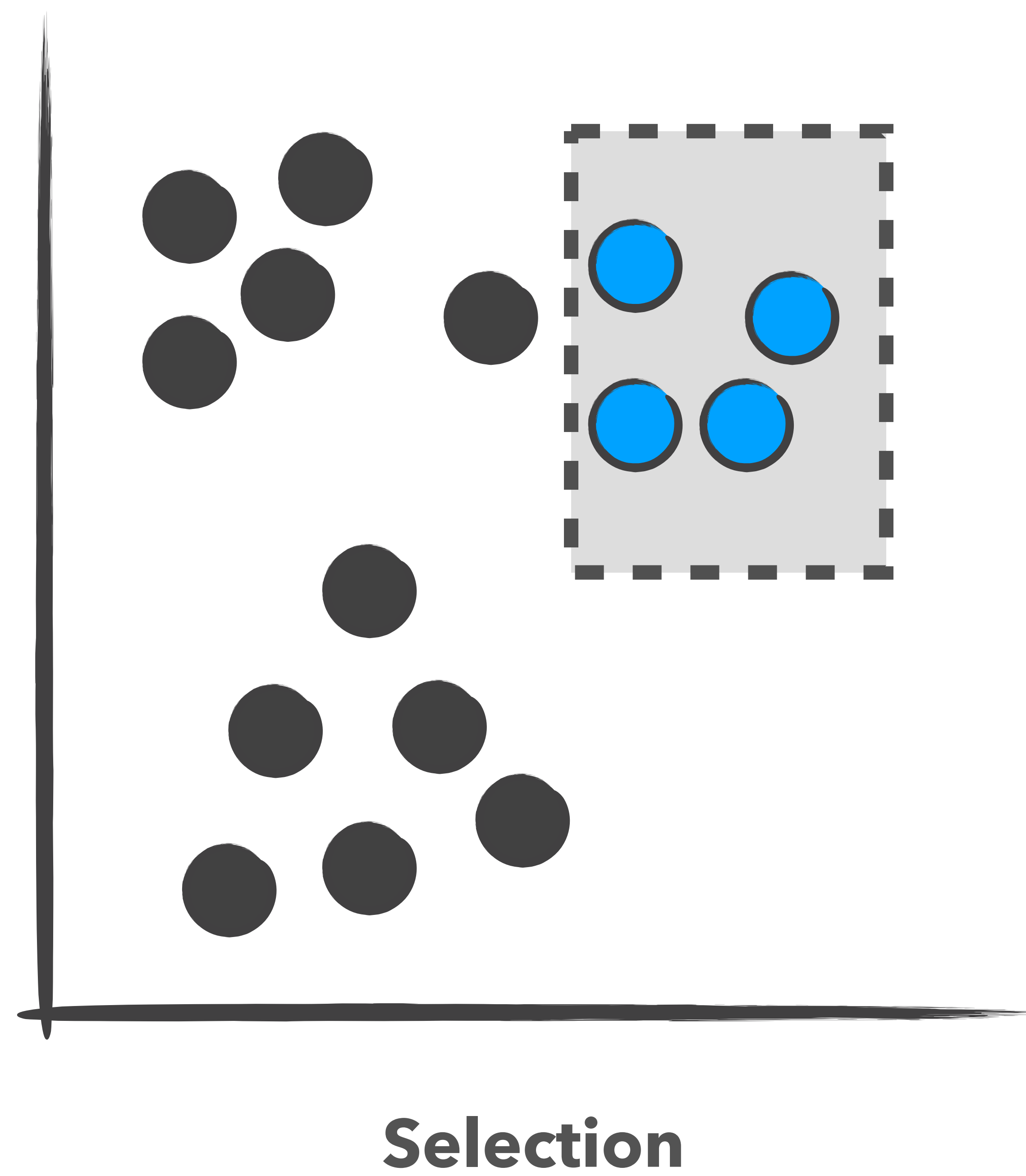


Multivariate Optimization

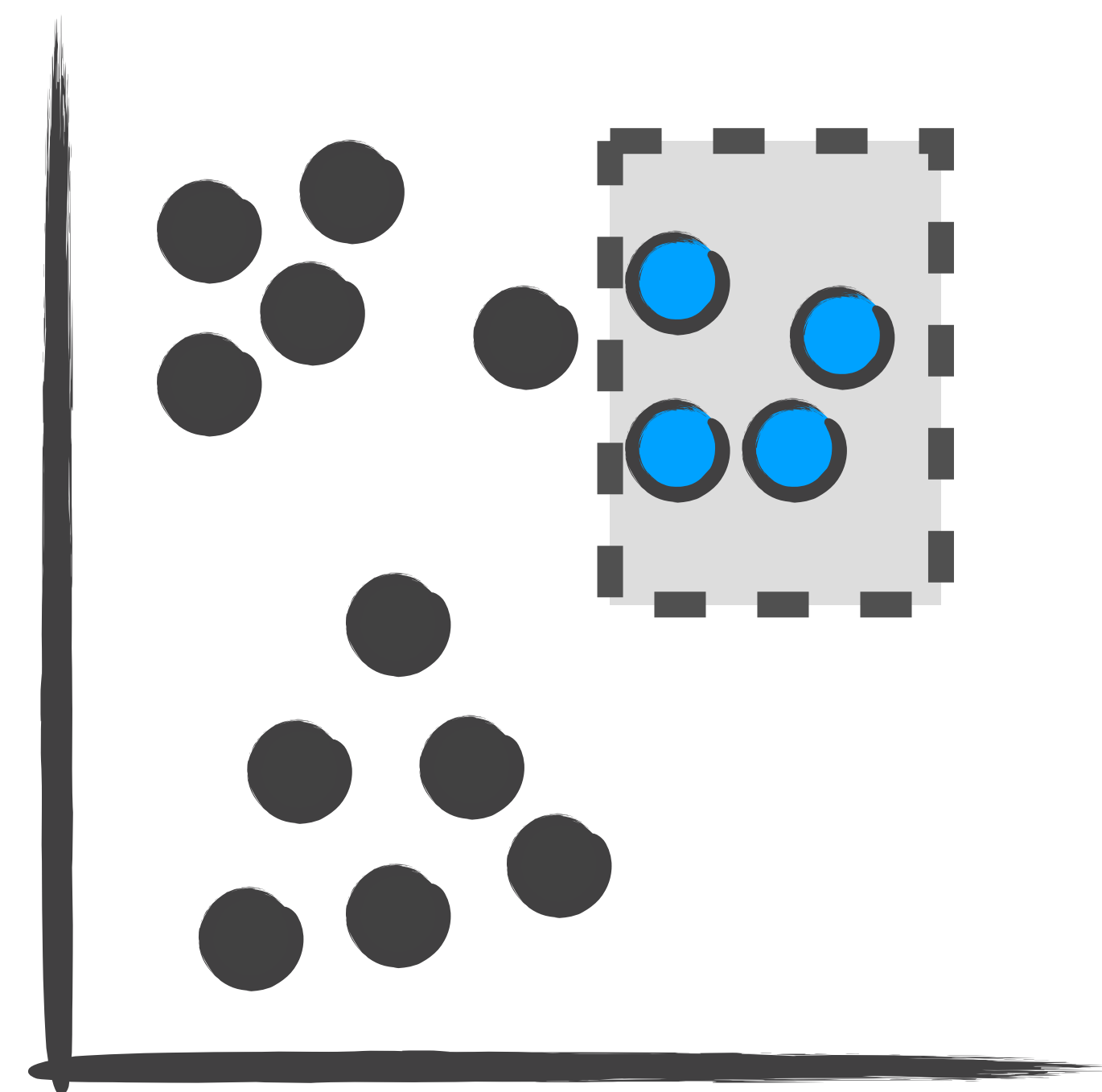


Ranges

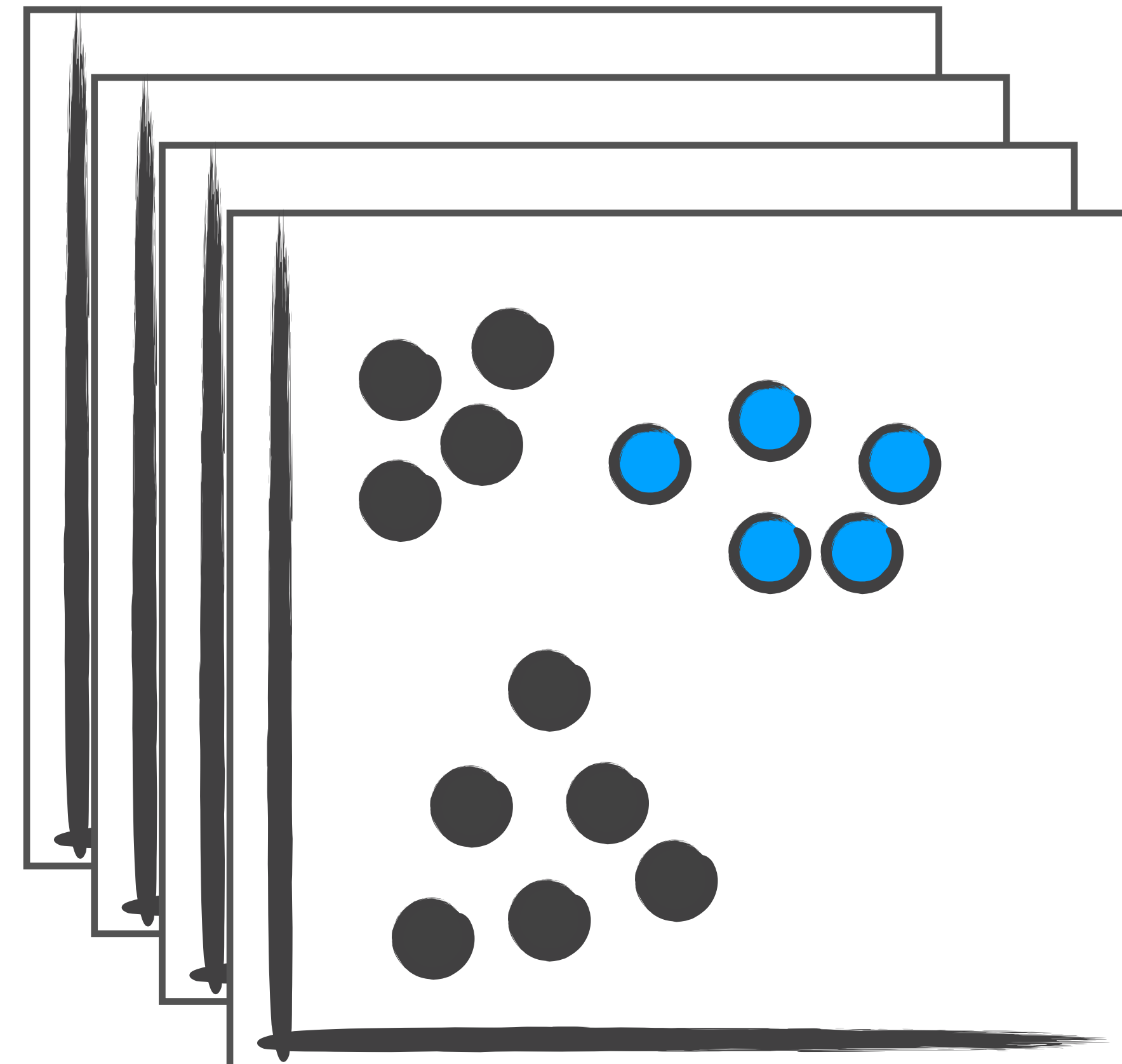
HOW DO WE INFER INTENT?



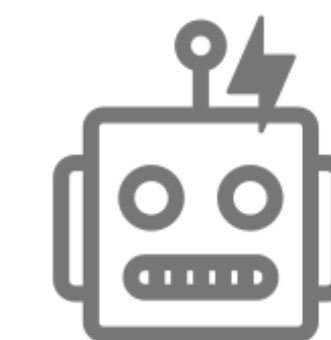
HOW DO WE INFER INTENT?



Selection

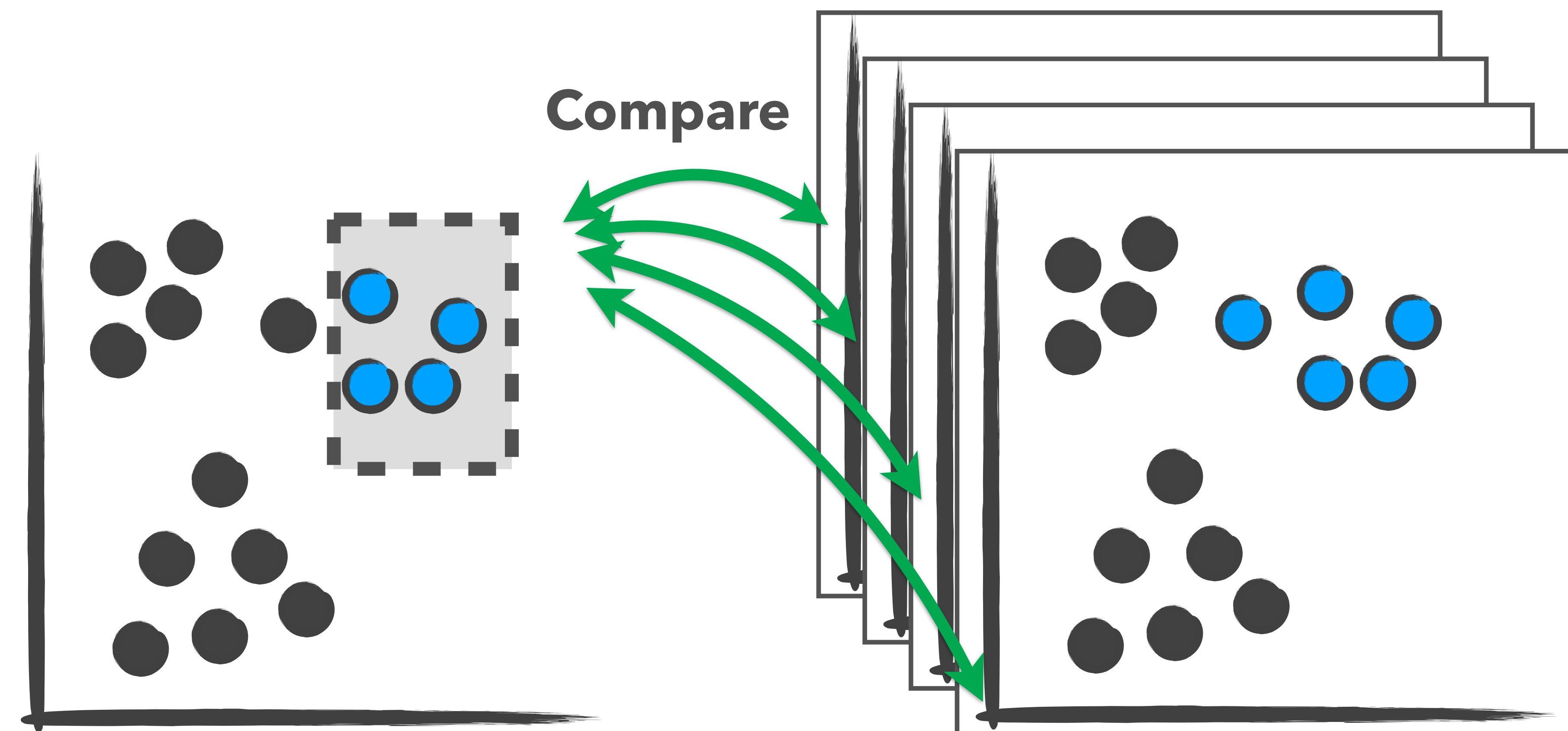


Predictions



K-Means
DBScan
Regression
Outlier Detection
Skyline
Decision Trees / Ranges
Categories

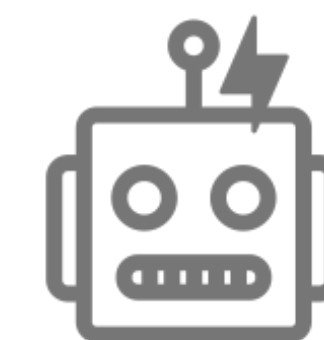
HOW DO WE INFER INTENT?



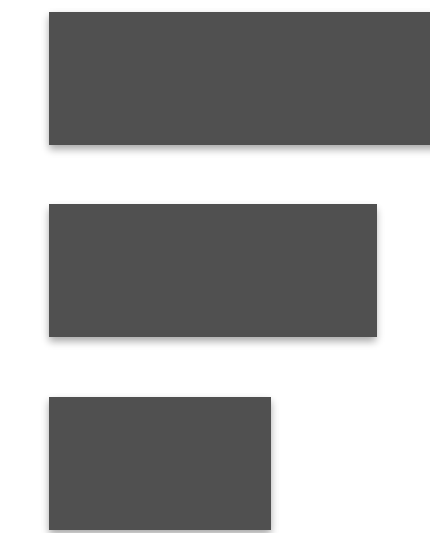
Selection

Predictions

K-Means
DBScan
Regression
Outlier Detection
Skyline
Decision Trees / Ranges
Categories

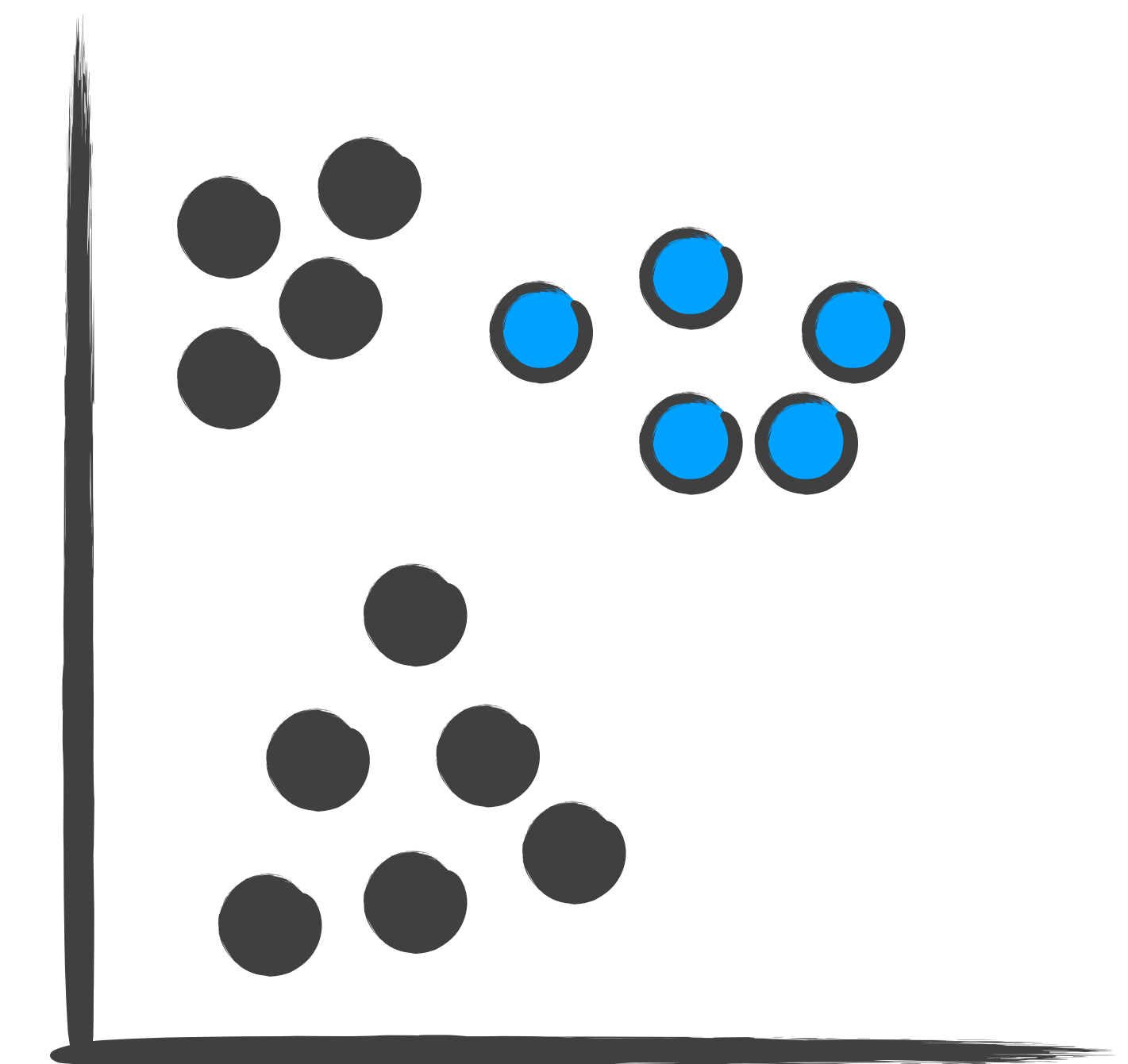
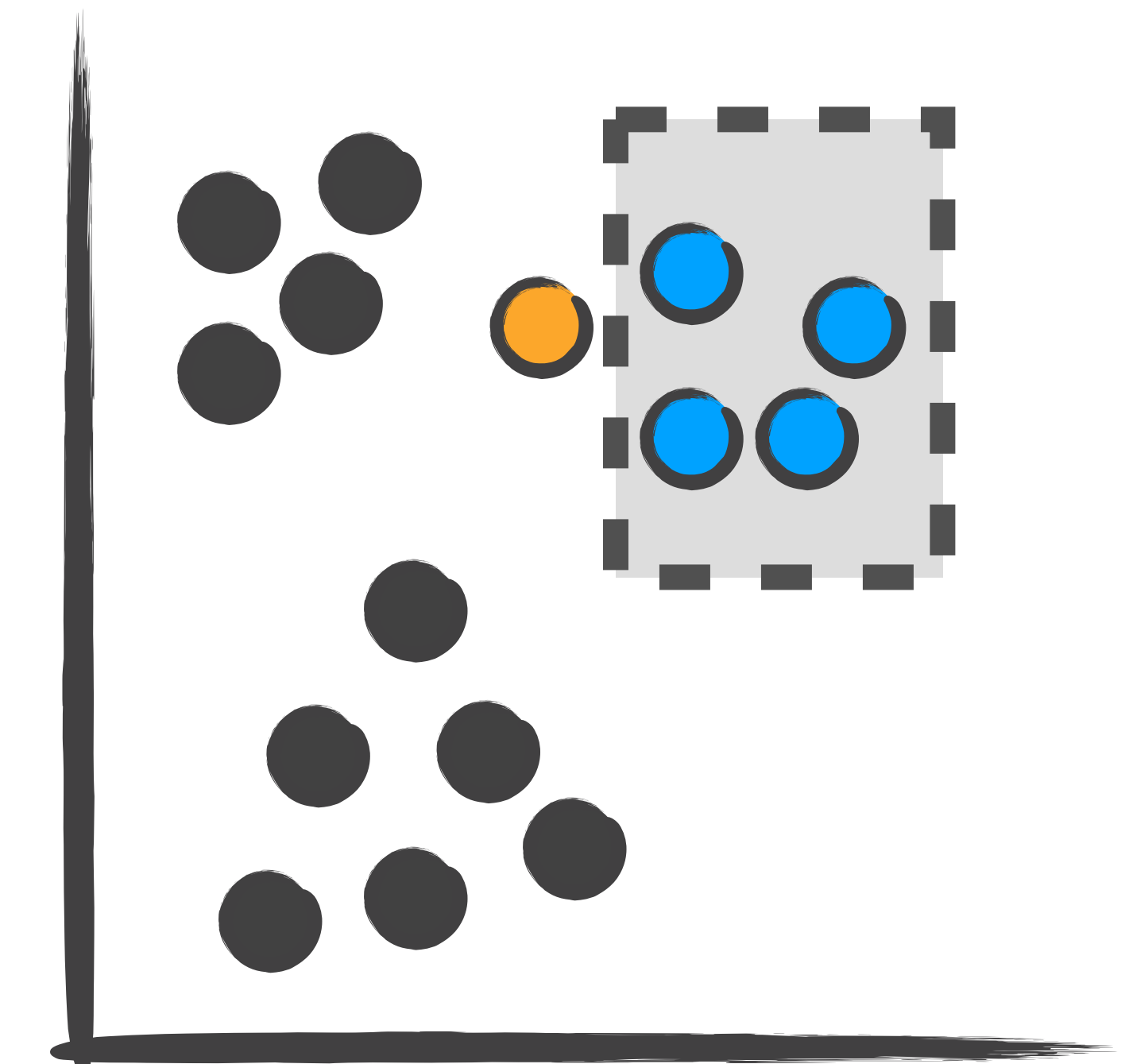
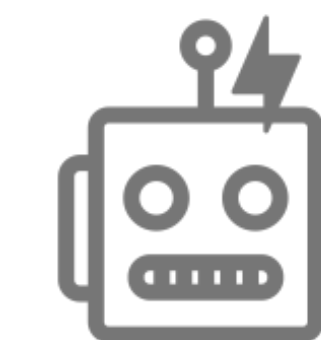


1. Range
2. Cluster
3. Outlier



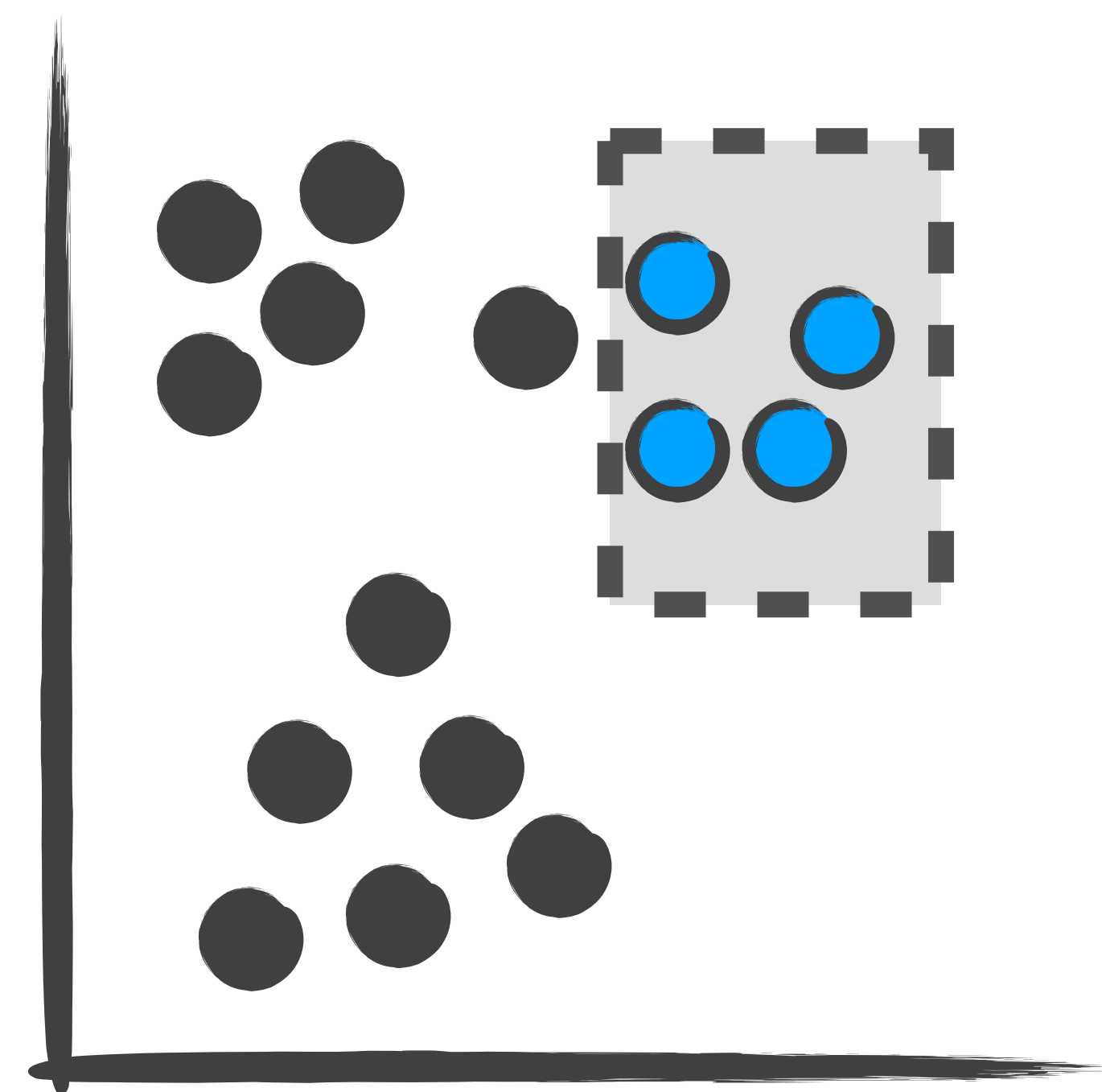
Ranking

Jaccard Distance
Naive Bayes
Classifier
Heuristic
Measures

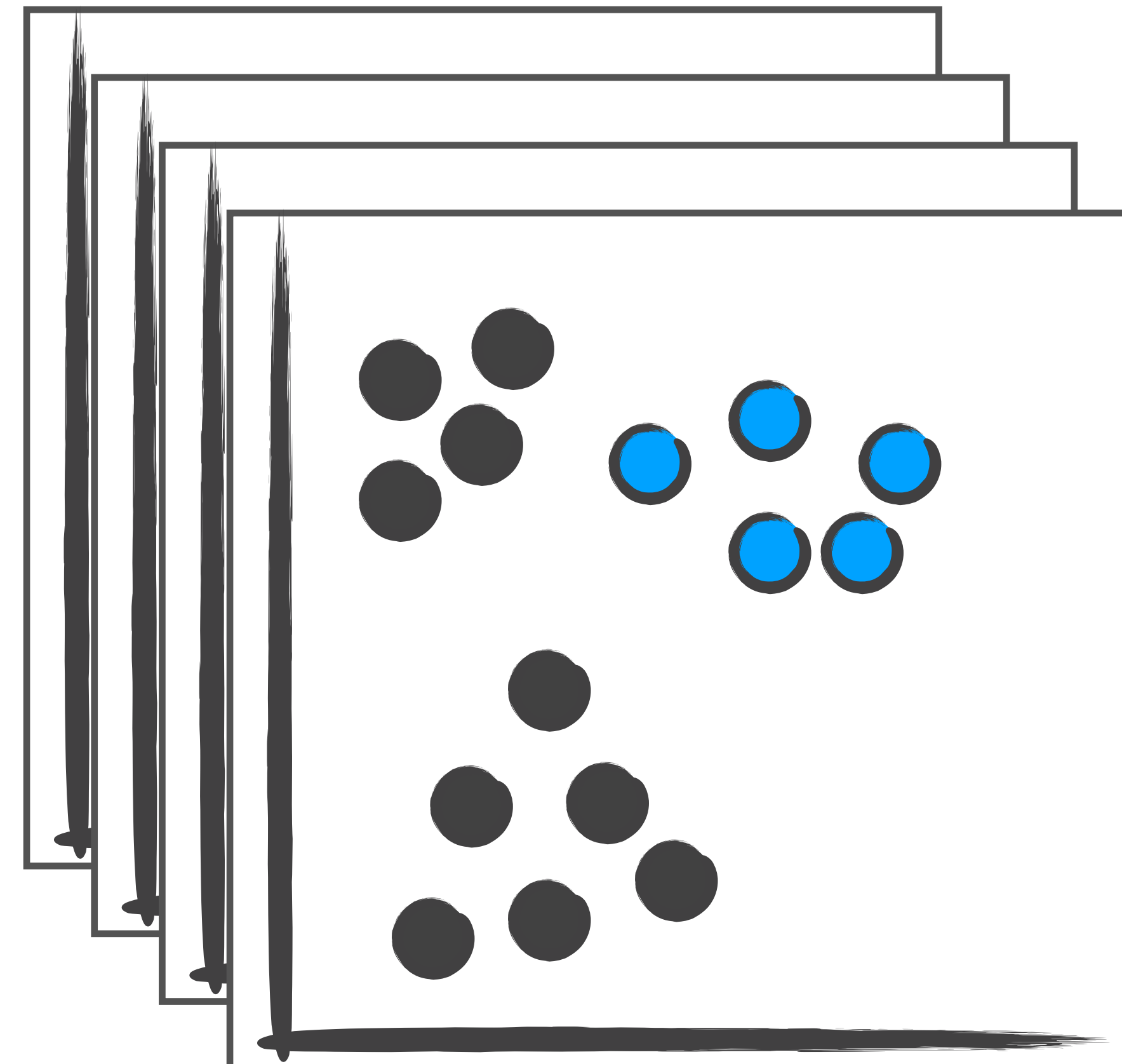


$$J(S, C) = \frac{|S \cap C|}{|S \cup C|}$$

HOW DO WE INFER INTENT?

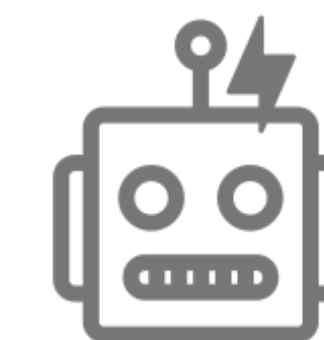


Selection



Predictions

K-Means
DBScan
Regression
Outlier Detection
Skyline
Decision Trees / Ranges
Categories

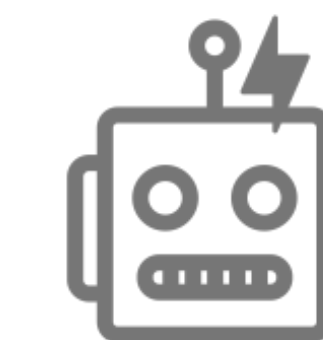


1. Range 
2. Cluster 
3. Outlier 

I think this cluster...

Ranking

Jaccard Distance
Naive Bayes
Classifier
Heuristic
Measures



Confirming Intent & Annotation 

Clusters ▾

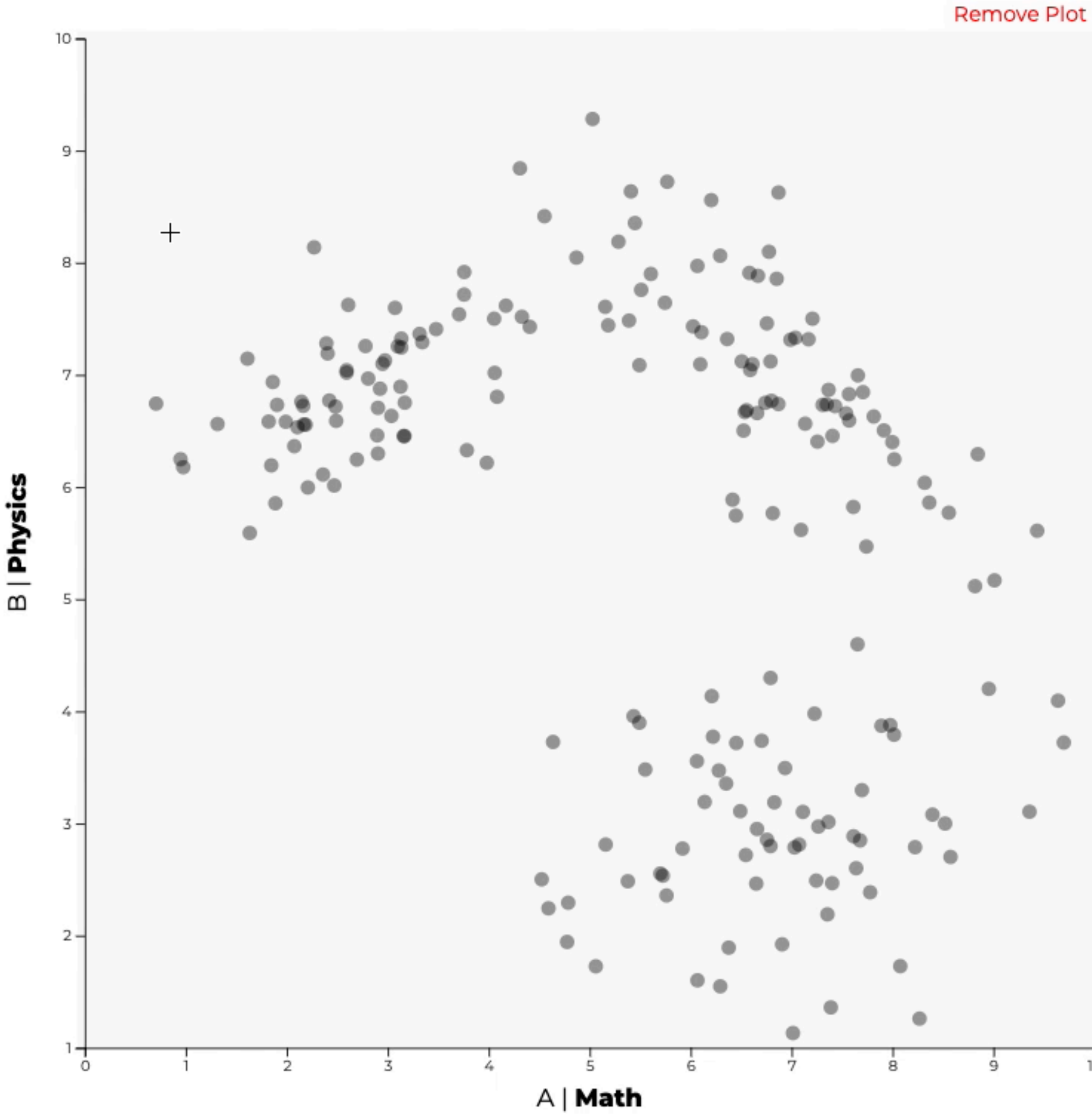
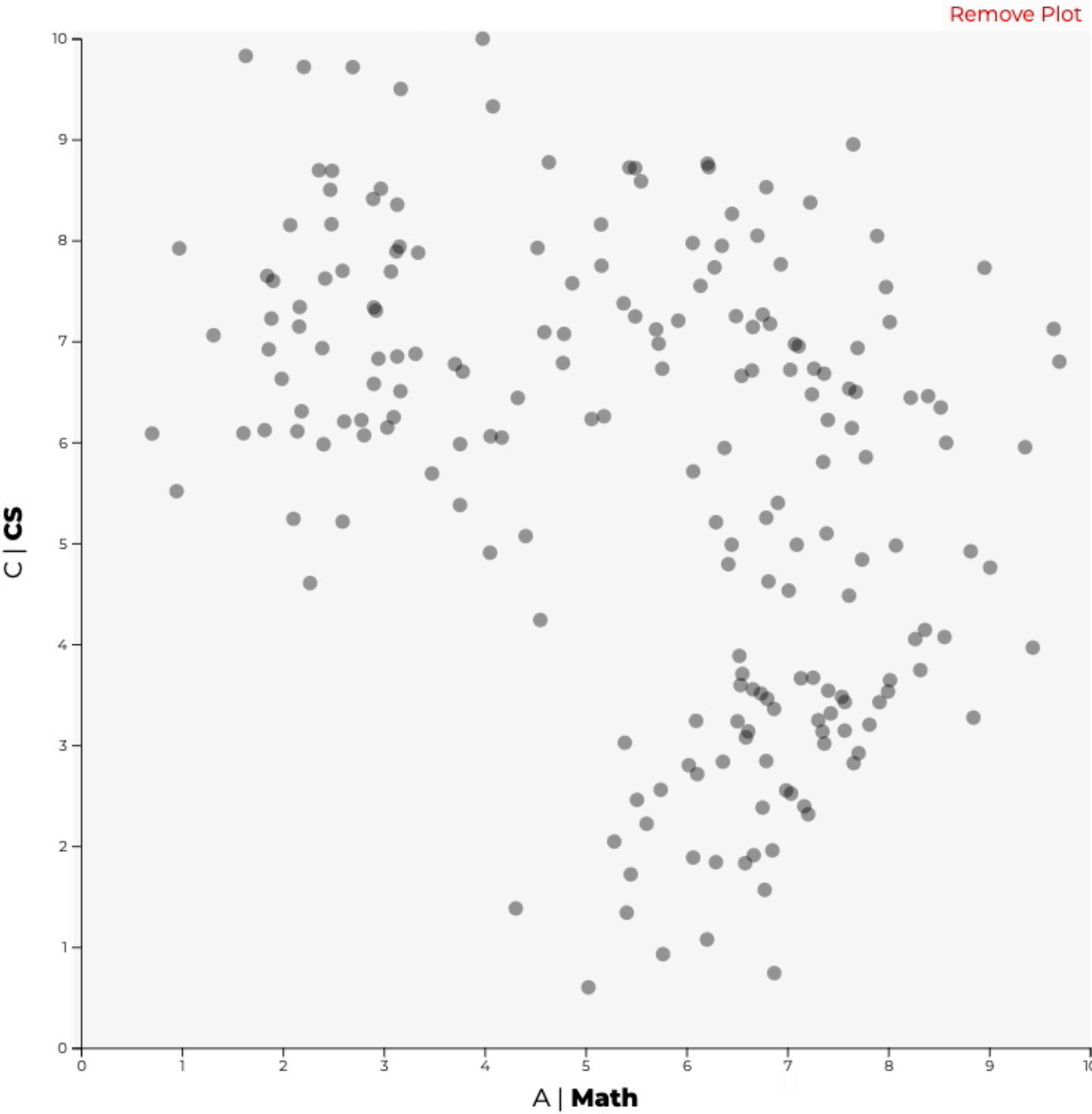
+ Add plot

☐ Show Categories

☒ Union

Invert Selections

Clear Selections



Visualization and Selection

Intent

Please interact

Annotate

Predictions

Time required: 0.01 seconds

Selections

0	0	0	0
UNION	INTERSECTION	INDIVIDUAL	TOTAL

Annotation of Intent and Predictions

REFLECTION

Robustness is easier in code:

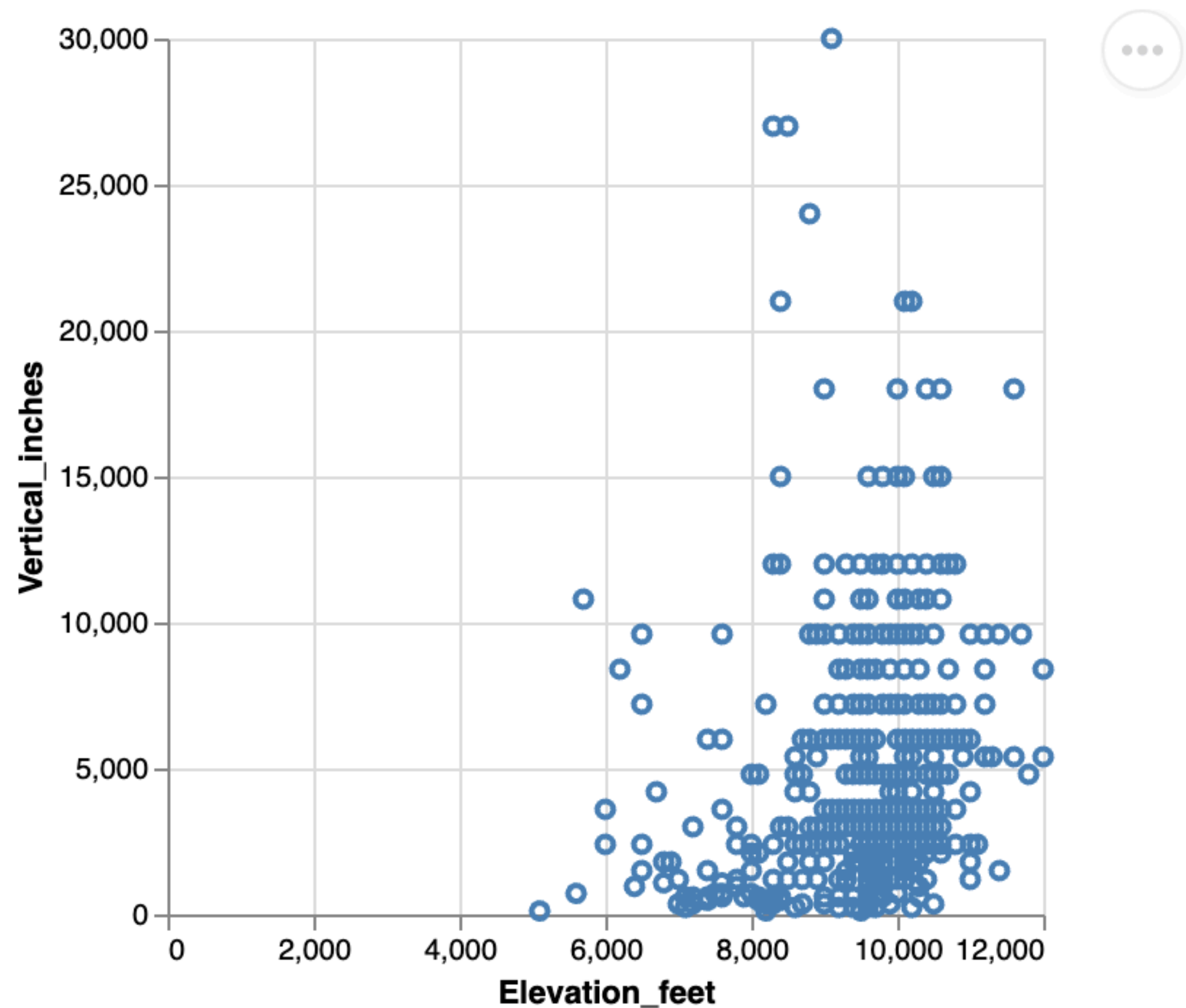
developers write rules, rather than lists of items.

Have to do extra work to **get to the "rules" when trying to create robust workflows from interactions**

**HOW CAN WE INTEGRATE
MODALITIES IN PRACTICE?**

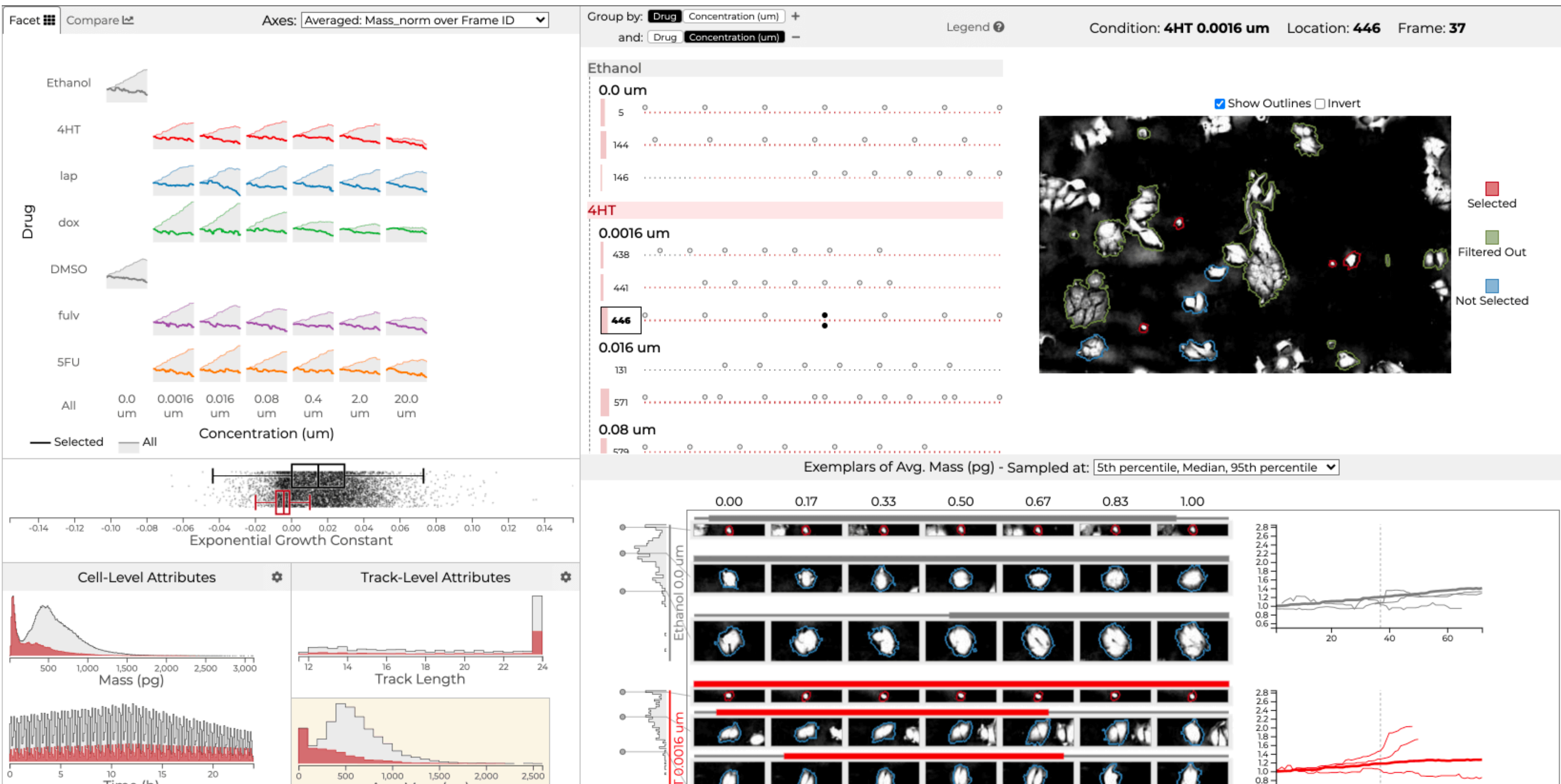
SPECTRUM OF VISUALIZATION SYSTEMS

Generic Charting



Easy to integrate
in **notebooks**

Specialized Systems

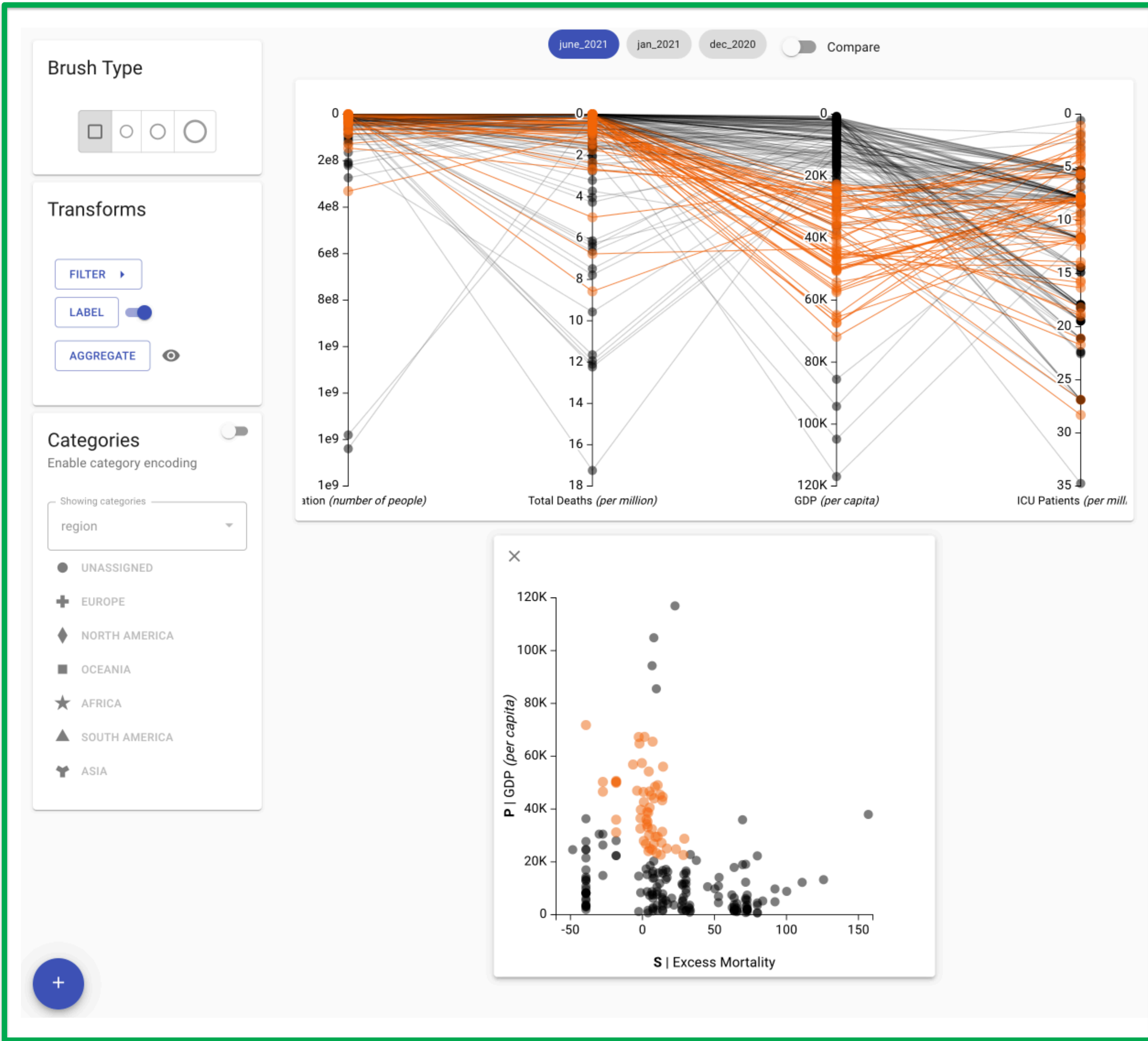


Require
standalone application

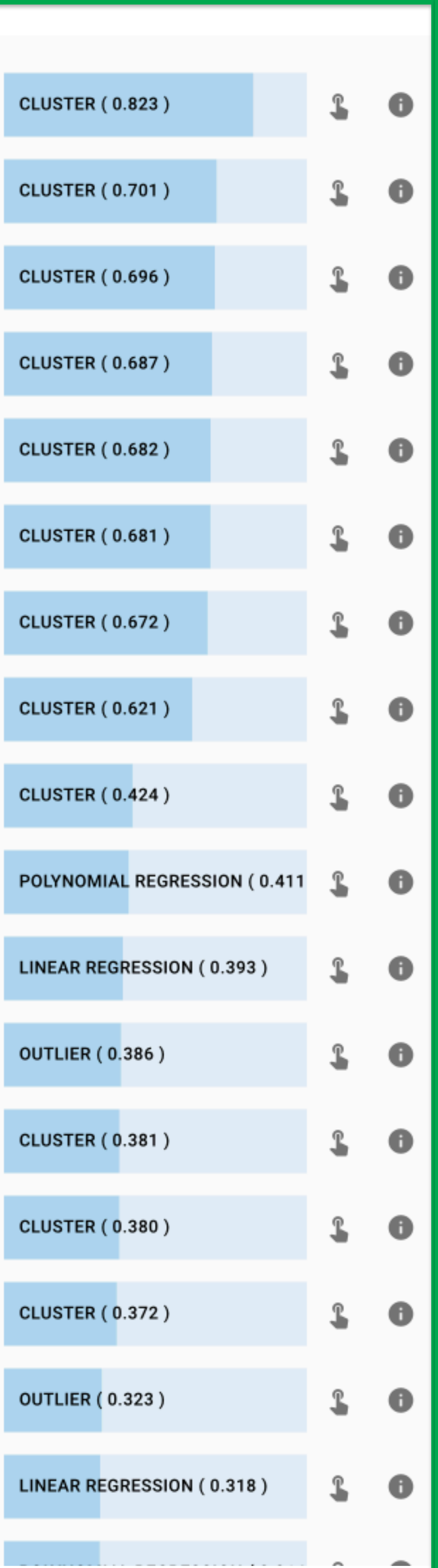
STANDALONE SYSTEMS USING WORKFLOWS

VISUALIZATION SYSTEM

Visualization Interface



Predictions



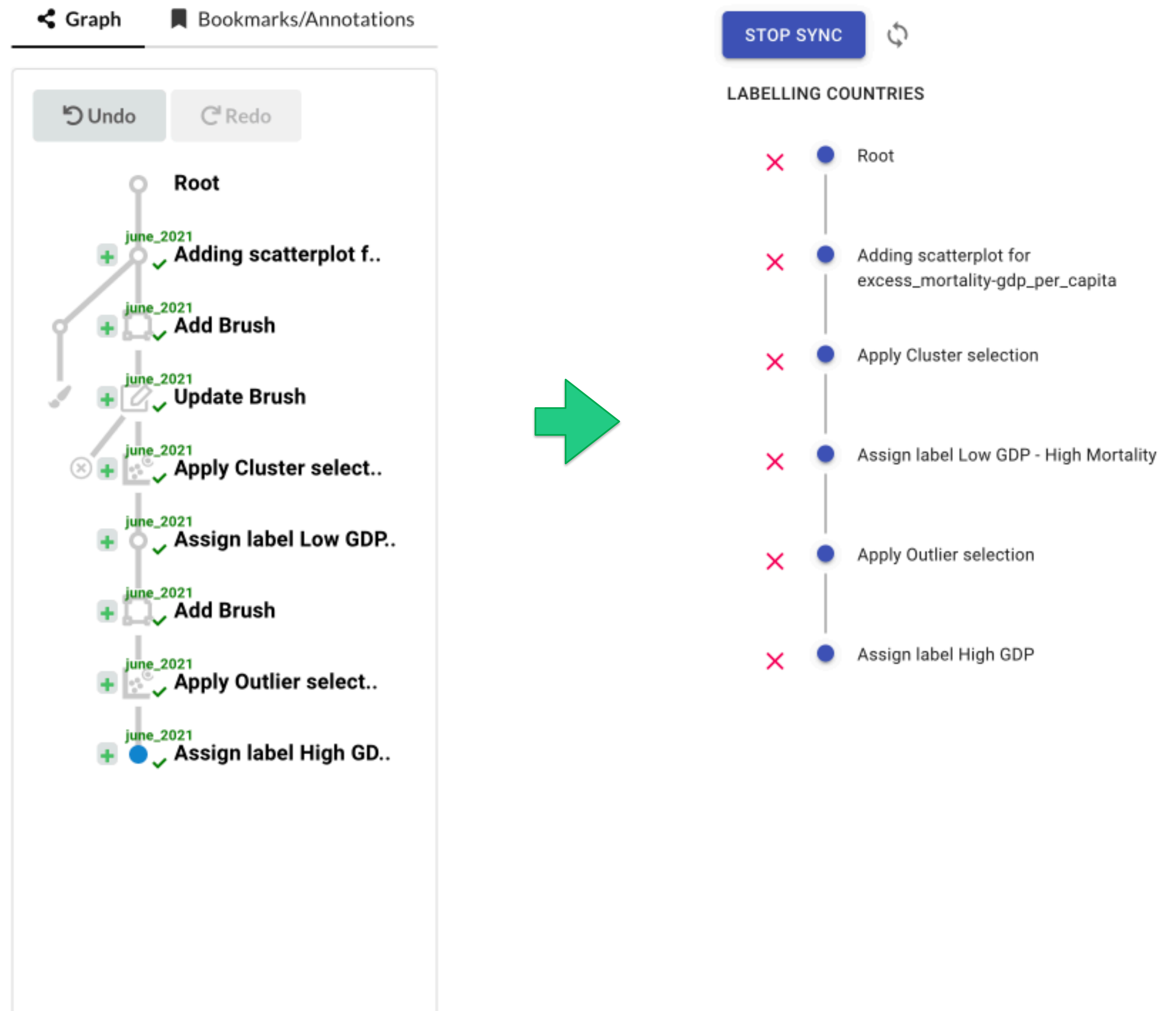
Provenance



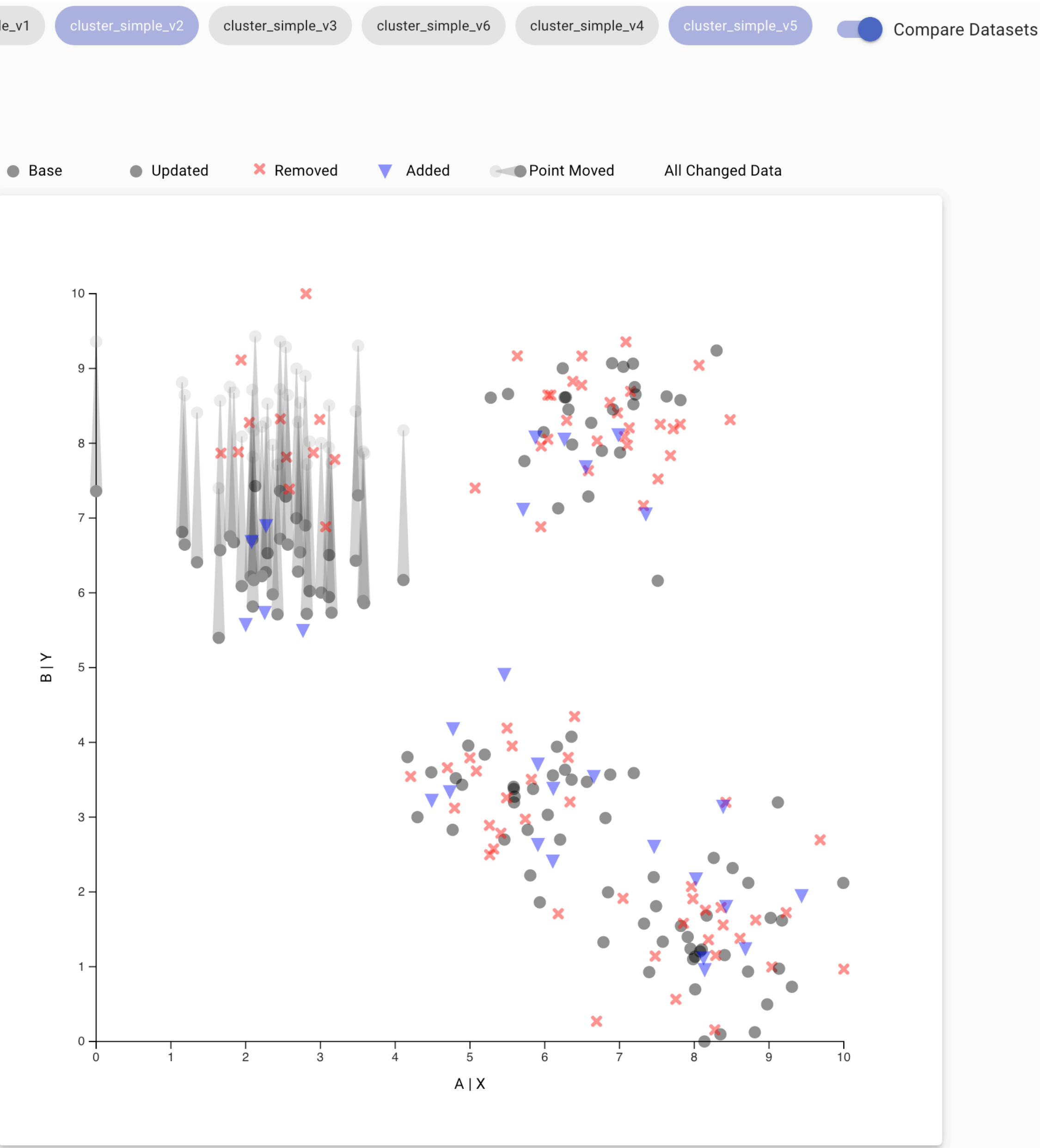
WORKFLOW EDITOR

Interaction logs
aren't clean

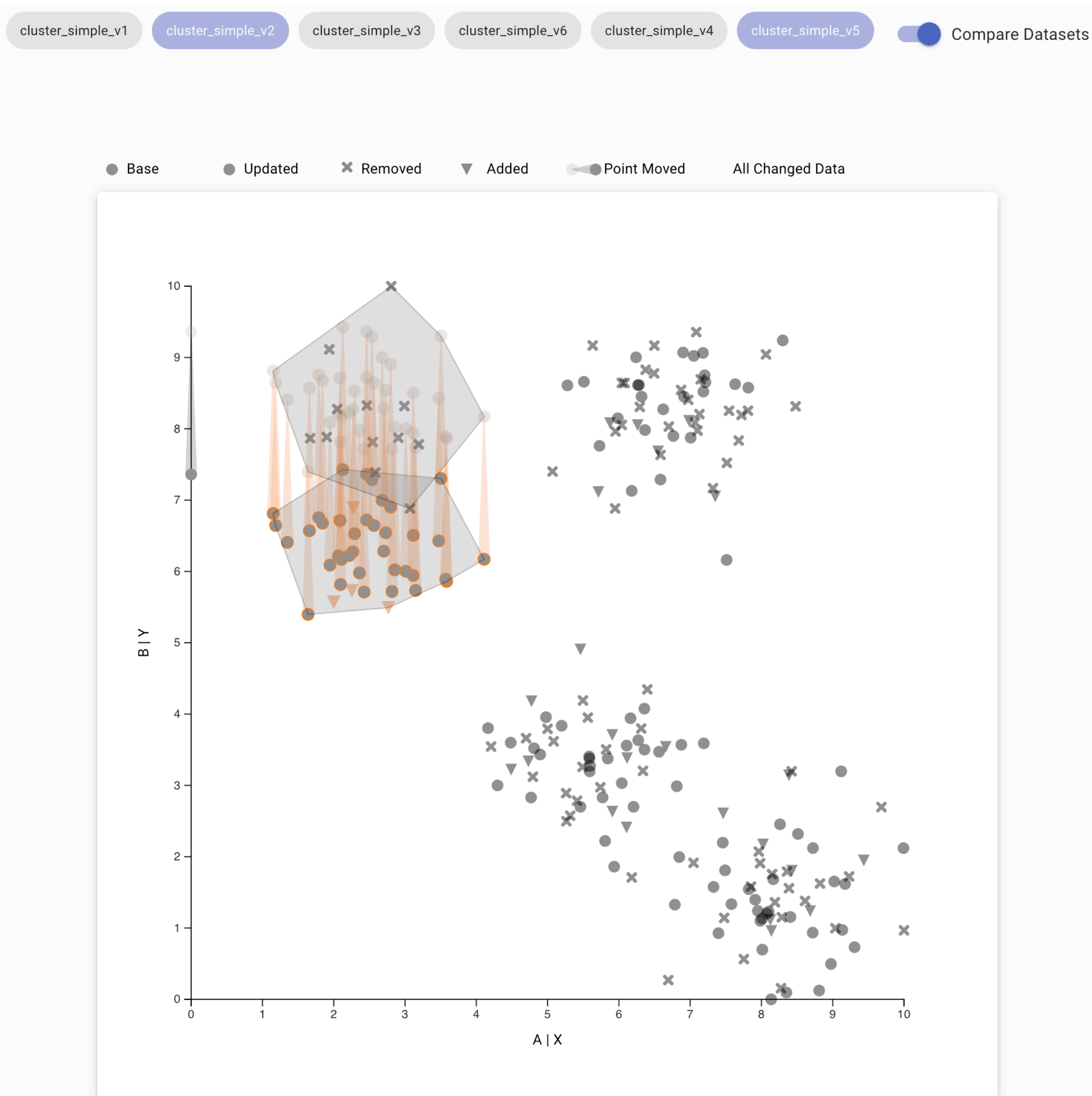
Need to tidy them up



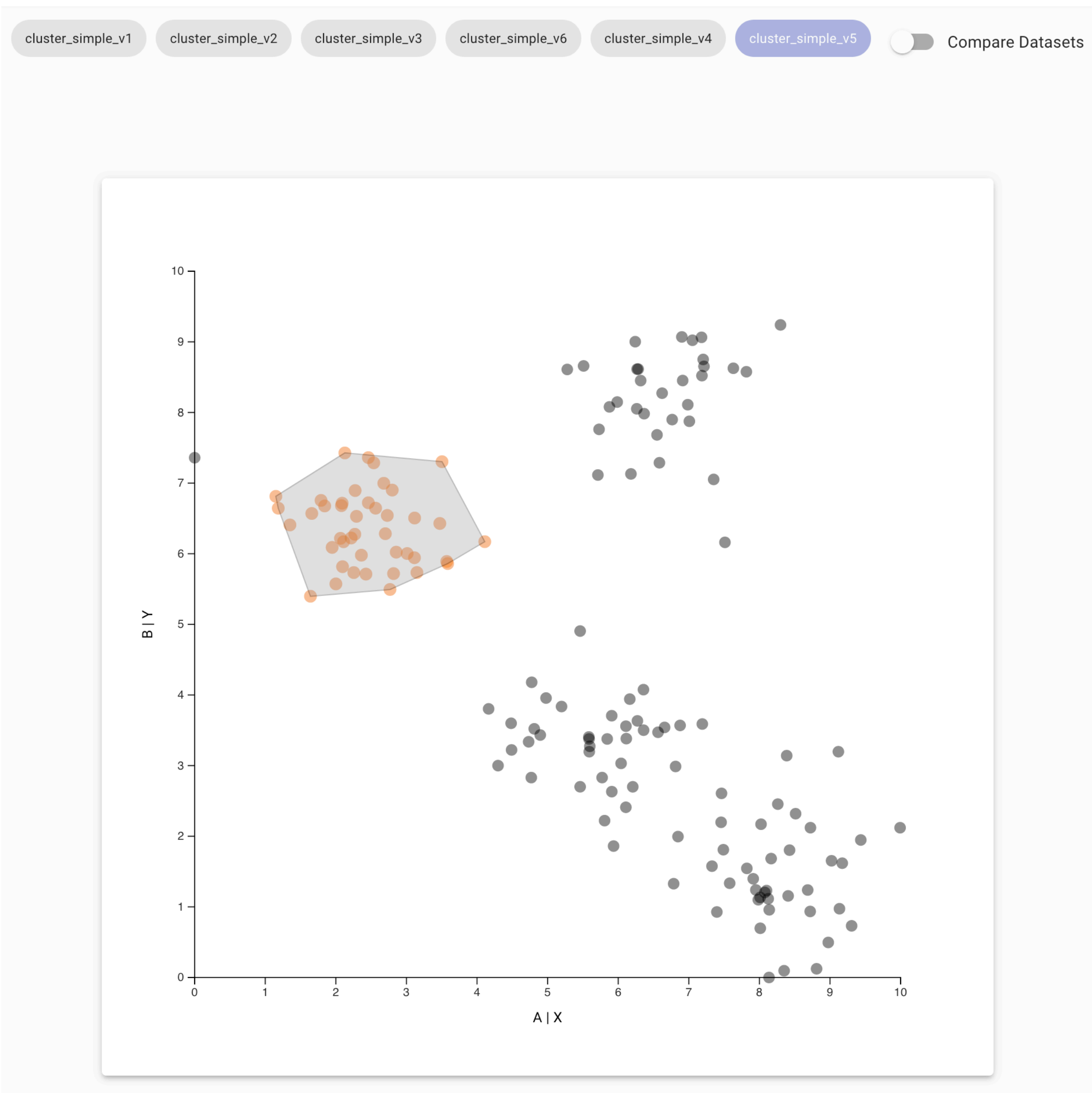
REUSING SELECTIONS ON UPDATED DATASETS



Changed Dataset

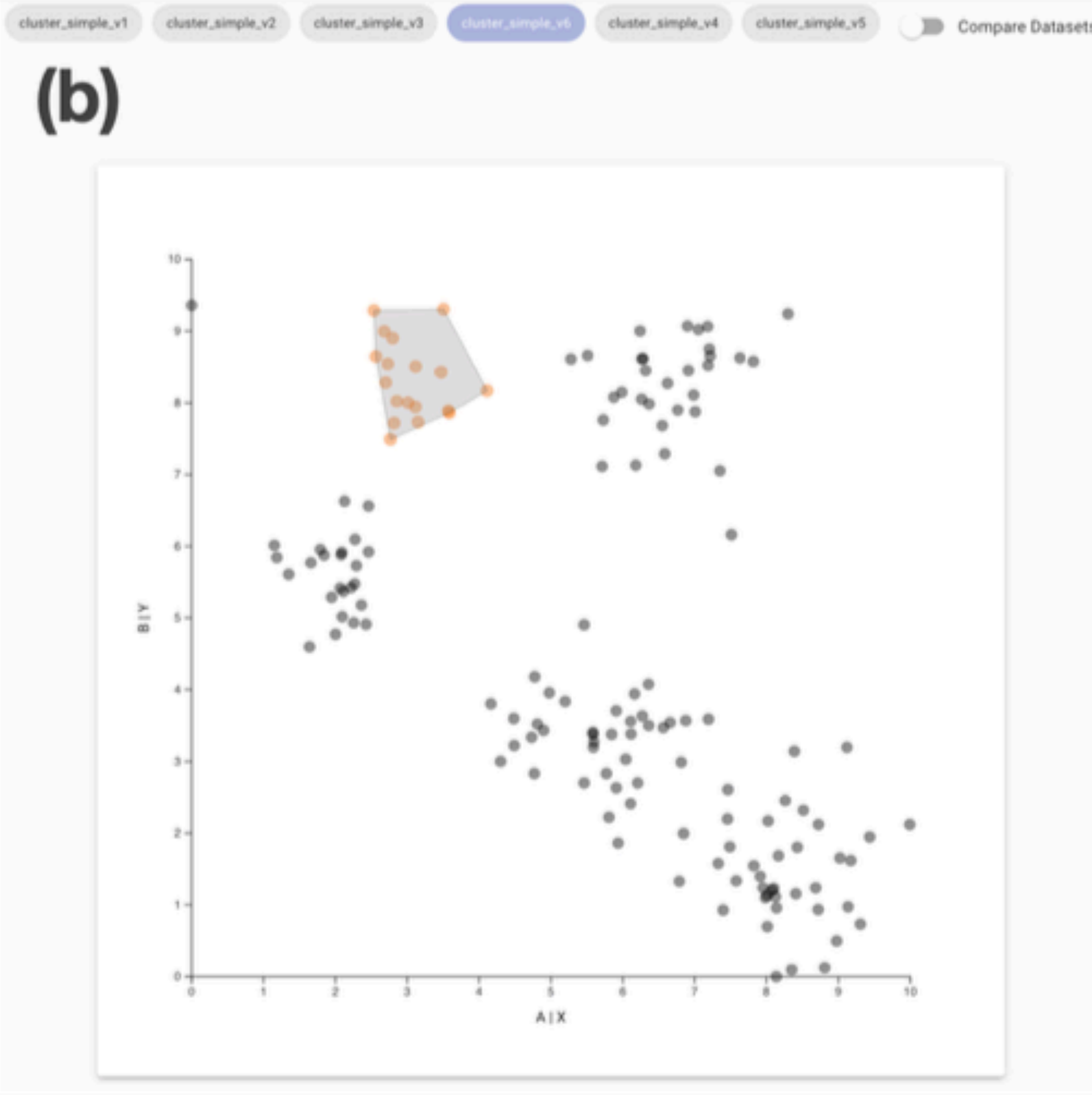
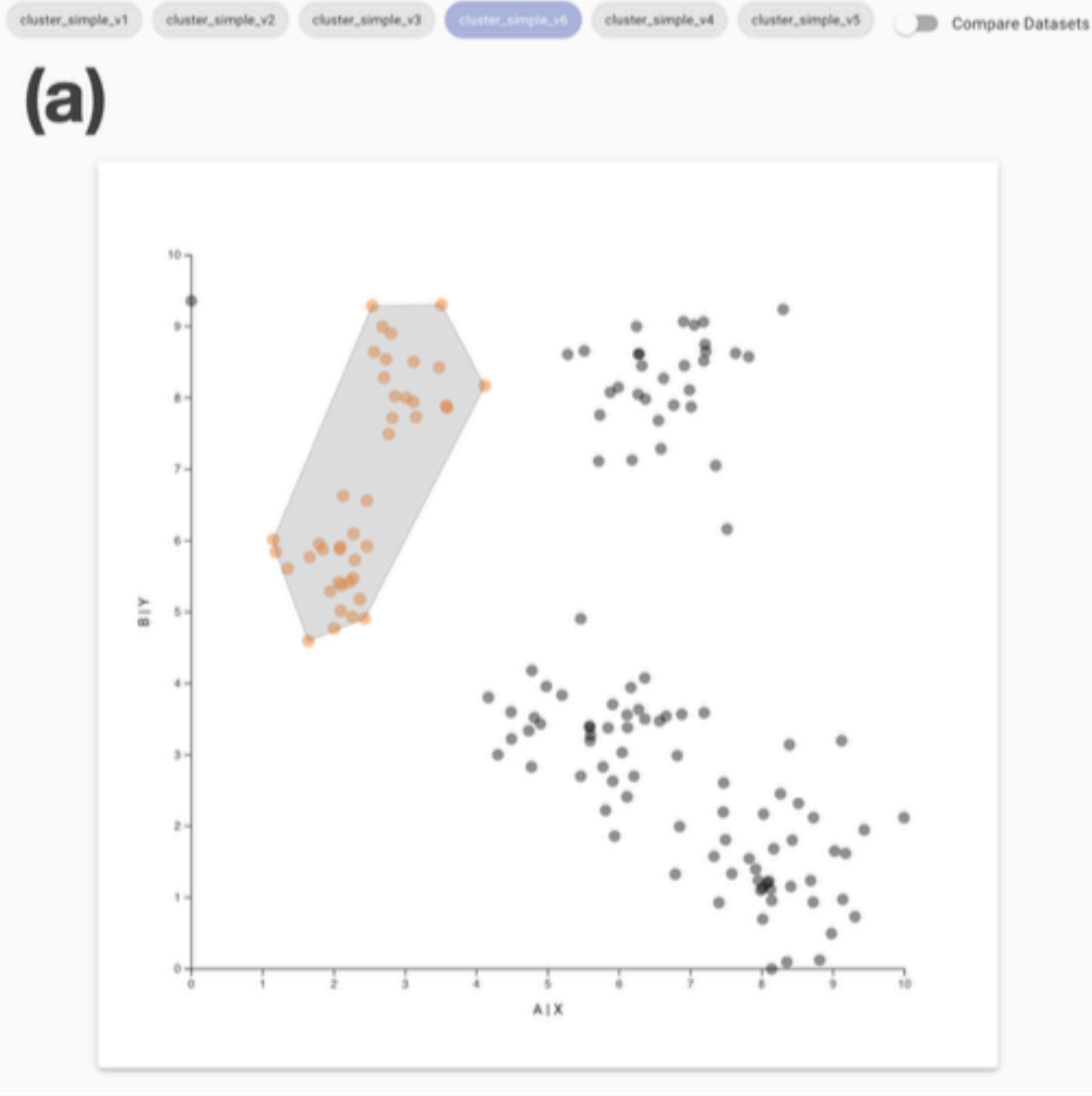


Tracking A Selected Cluster



Selected Cluster on Changed Dataset

HUMAN REVIEWS



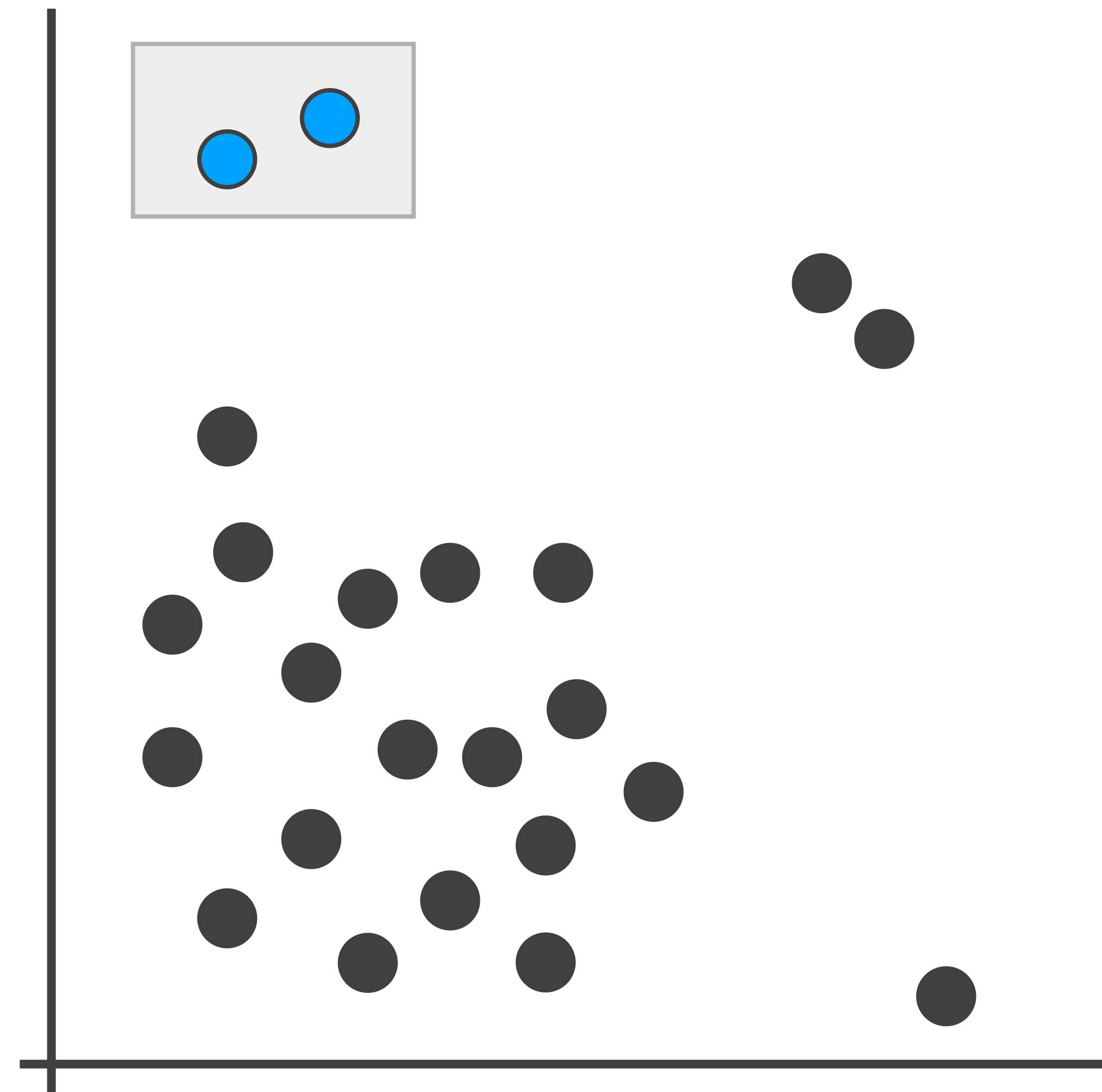
WORKFLOWS: OPPORTUNITY FOR BRIDGING THE GAP

A workflow is a **series of executable steps**

Can be abstracted and **re-used in different environments**

> **Bridging between our tool and Python**

Interactive Visual Analysis



Workflow
Database

Computational Analysis



Library that tracks and re-executes actions

USING WORKFLOW IN A COMPUTATIONAL NOTEBOOK

```
# Installing the reapply-workflows adds a module called backend
# This module exposes the Reapply class which initializes the library

from backend import Reapply
```

```
# Here we load the reapply_workflows library.
r = Reapply()

# We add a workflow from our workflow database.
workflow = r.load_workflow("workflow1617808681620")
```

```
# Print the workflow name
print("Workflow: ", workflow.name, "\n")

# Description of the workflow and the operations in it
workflow.describe
```

Workflow: Deleting Cluster

```
| Root
+--| Add Plot
    +--| Added brush to: X-Y
        +--| Cluster Selection
            +--| Filter: Out
```

```
# Prints the reapply results for all interactions, along with review status.
```

```
# Apply the workflow to target dataset.
# apply function requires the target dataset
# and the label column as arguments.
res = workflow.apply(target, "Label")

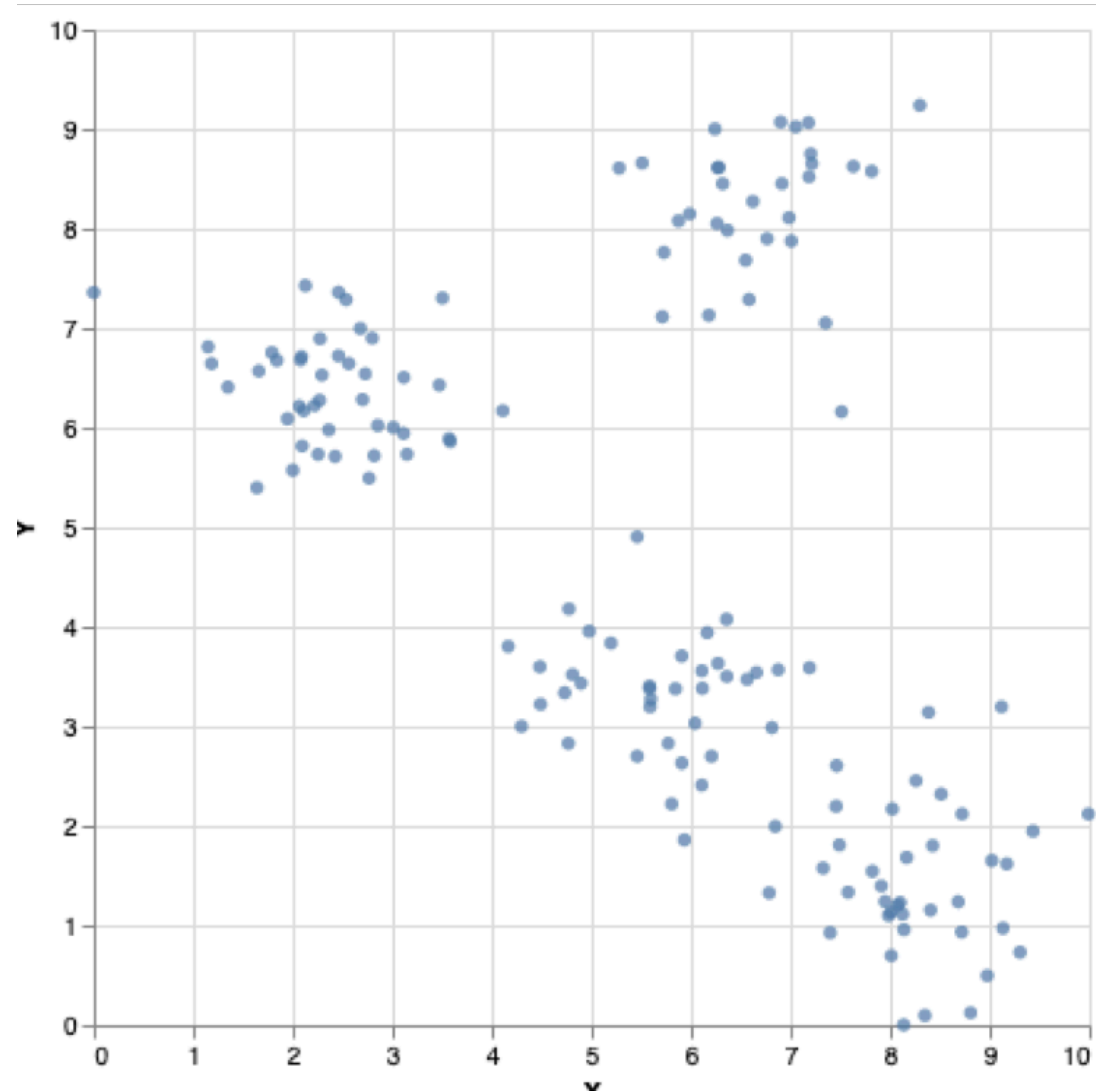
# Results is an array of datasets for each interaction
# we grab the final one.
result_dataset = res.results[-1]['data']
result_dataset
```

This workflow has not been reviewed for all interactions.
Please go to following url: <https://reapply-workflows.git/>

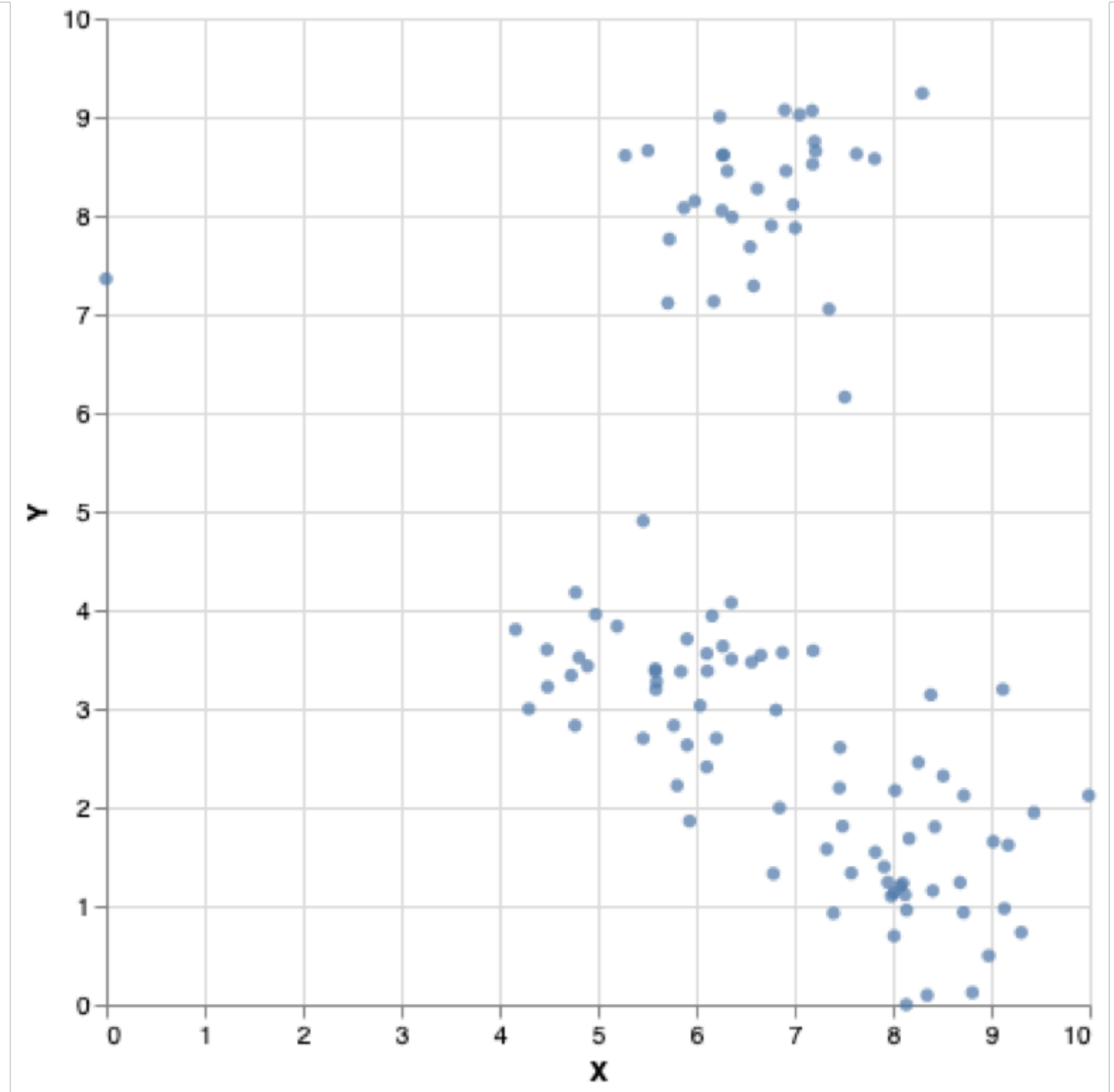
	Label	X	Y
3	P52	6.58351	7.28796
5	P171	4.77421	4.17980
8	P199	8.34966	0.09550
9	P183	8.42670	1.80299
10	P61	4.29760	2.99981
...
141	P138	7.35179	7.05215
142	P46	6.62171	8.27311

BEFORE AND AFTER

150 rows × 3 columns



108 rows × 3 columns



Brush Type



Transforms

FILTER ▸

LABEL ☒

AGGREGATE ☐

Categories

Enable category encoding

Showing categories ▾

Case Category ▾

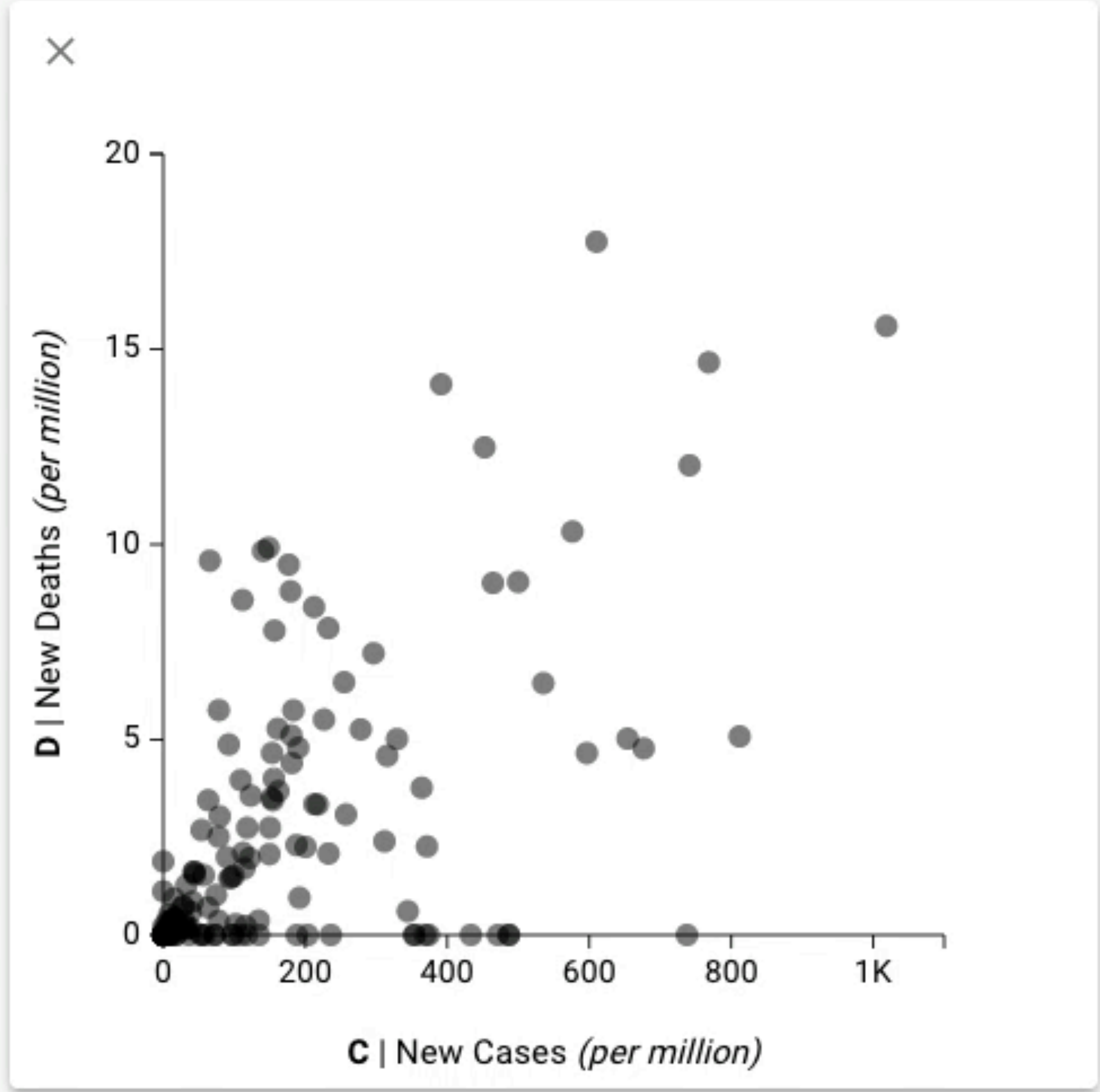
- UNASSIGNED
- + HIGH DEATHS - HIGH CASES
- ◆ LOW DEATH - HIGH CASES
- HIGH DEATHS - LOW CASES

jan_2021

june_2021

dec_2020

☐ Compare



NO VALID SELECTIONS

Undo

Redo

Root

jan_2021 Adding scatterplot f..

jan_2021 Adding scatterplot f..

jan_2021 Remove Scatterplot



REFLECTIONS

Useful for

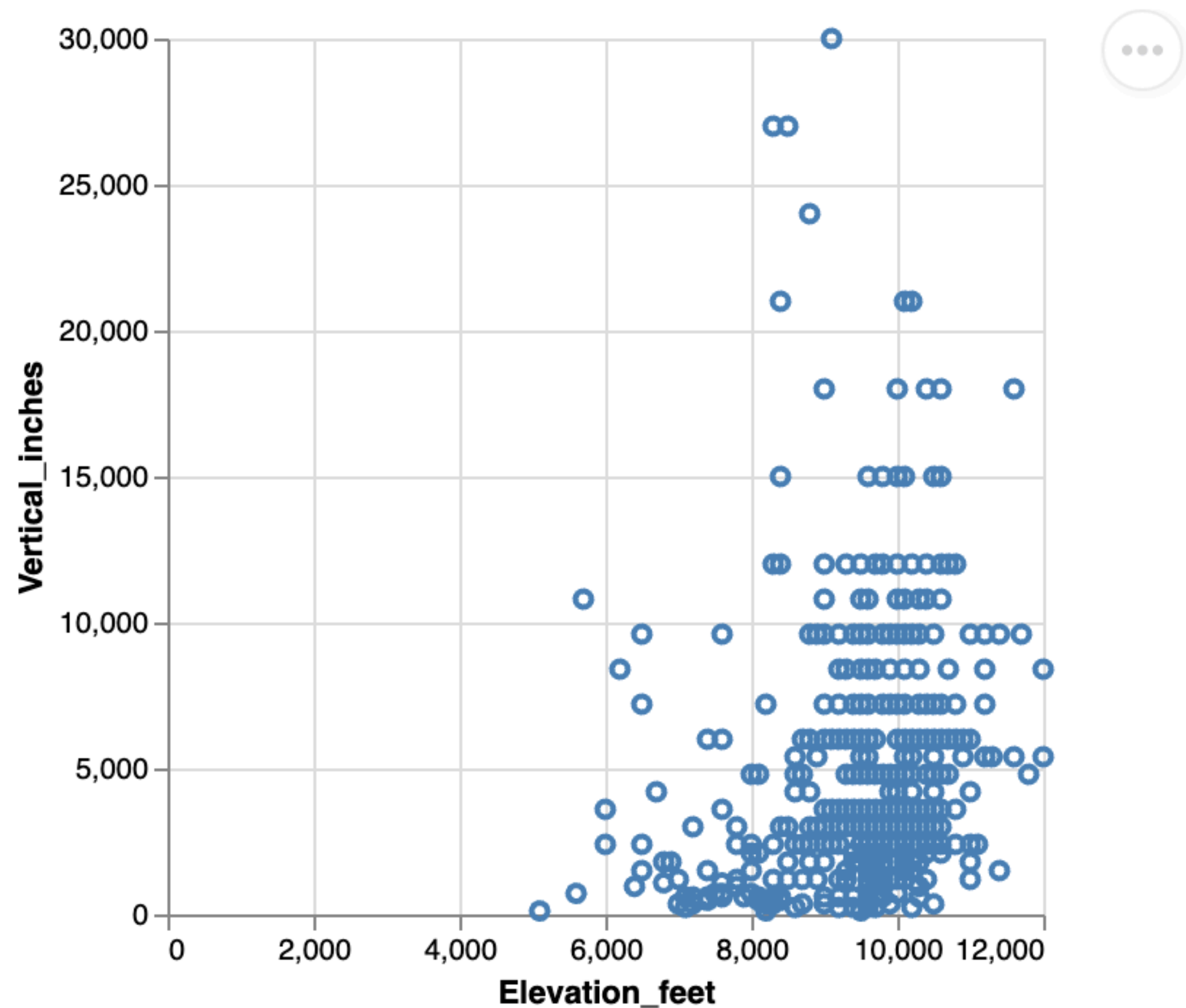
**using a standalone system with
updating data**

**using scripting after using a
standalone system**

**creating template workflows
(teaching)**

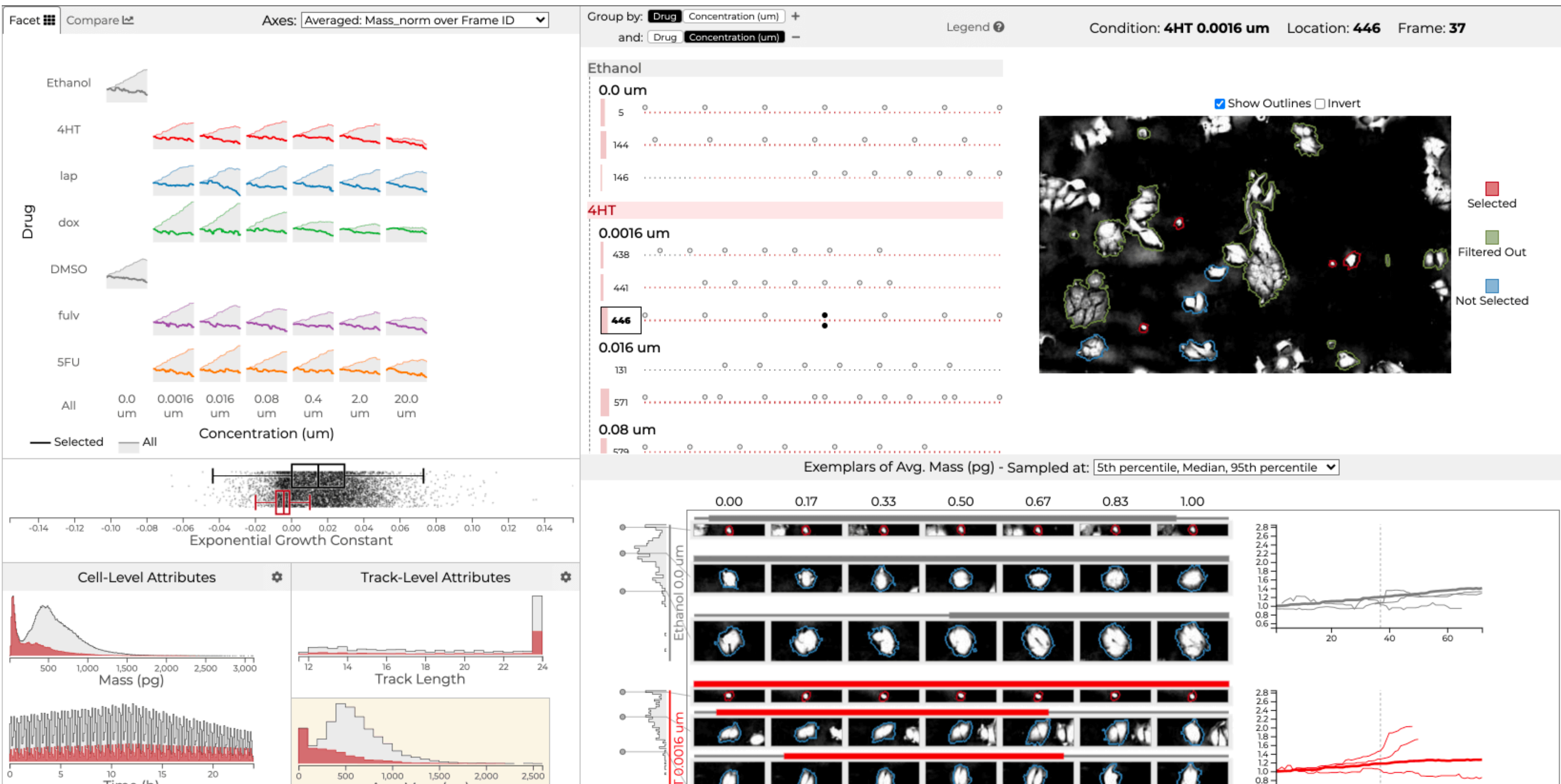
SPECTRUM OF VISUALIZATION SYSTEMS

Generic Charting



Easy to integrate
in **notebooks**

Specialized Systems



Require
standalone application

TRACKING PROVENANCE IN INTERACTIVE PLOTS IN NOTEBOOKS

WHY?

Interactive plots are now common in notebooks (e.g., Vega-Altair)

Operations such as **filtering outliers or **changing column labels** are more efficient to do in interactive plots**

PRINCIPLE

Track events in interactive visualizations

Map them to data frame operations

Operations then applied to data frame

SCOPE

“Listen” to all events by Vega-Altair

works for all charts that can be created
with Vega-Altair

no extra effort on developer!

**Custom table visualization for data
frames**

instead of `df.head()` -> use interactive
table

VEGA-ALTAIR EXAMPLE

Standard Python and Vega-Altair Code

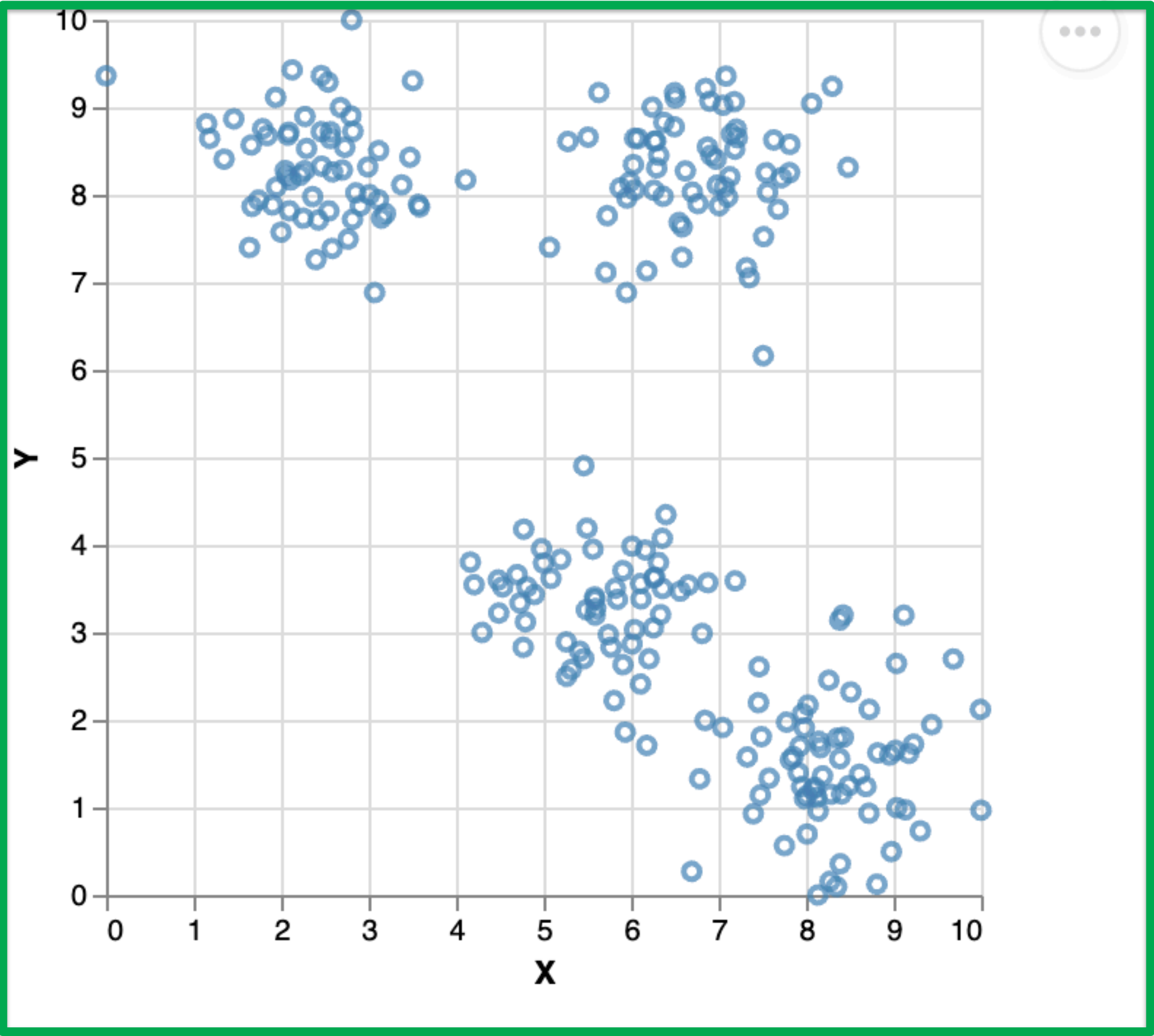
```
[3]: df = pd.read_csv("./cluster_simple_v1.csv")

brush = alt.selection_interval(name="brs")

alt.Chart("./cluster_simple_v1.csv", name="sp").mark_point().encode(
    x='X:Q',
    y='Y:Q',
    color=alt.condition(brush, alt.value('steelblue'), alt.value('grey')),
).add_params(brush)
```

[3]:

Toolbar Injected



Plot Generated by Vega Altair

Track	Description	Intent	Selections
○	Root		
○	Brush selection		
○	Aggregate by: mean		
○	Brush selection		
○	Aggregate by: mean		

Provenance Tracking Injected

TABLE EXAMPLE

```
[27]: table = PR.vis.interactive_table(df)
table
```



[27]:



<input type="checkbox"/>	index int64	Name object	Miles_per... float64	Cylinders int64	Displacem... float64	Horsepower float64	Weight_in... int64	Accelerati... float64	Year datetime6...	Origin object
<input type="checkbox"/>	1	chevrolet chev...	18	8	307	130	3504	12	0	USA
<input type="checkbox"/>	2	buick skylark ...	15	8	350	165	3693	11.5	0	USA
<input type="checkbox"/>	3	plymouth satel...	18	8	318	150	3436	11	0	USA
<input type="checkbox"/>	4	amc rebel sst	16	8	304	150	3433	12	0	USA
<input type="checkbox"/>	5	ford torino	17	8	302	140	3449	10.5	0	USA
<input type="checkbox"/>	6	ford galaxie 500	15	8	429	198	4341	10	0	USA
<input type="checkbox"/>	7	chevrolet impala	14	8	454	220	4354	9	0	USA
<input type="checkbox"/>	8	plymouth fury iii	14	8	440	215	4312	8.5	0	USA
<input type="checkbox"/>	9	pontiac catalina	14	8	455	225	4425	10	0	USA
<input type="checkbox"/>	10	amc ambassa...	15	8	390	190	3850	8.5	0	USA

Ttrack

Predictions

●

Root

OPERATIONS

Editing column names, cells

Sorting rows/columns

Dropping rows/columns

Filtering (in/out) items (with intent predictions)

Labeling items

Categorizing items

Grouping / Aggregating items

REFLECTIONS

Provenance enables **undo/redo** in notebook

Branching enables **alternative explorations**

Minimal invasiveness makes it **easy to adopt**

Fully **reusable**, just like code

**User interaction when useful,
use coding when more efficient!**

TAKE-AWAYS

Provenance can be used to **bridge analysis modalities!**

Translation to meaningful operation isn't always easy - more work needed!

Low-level charts are ripe target for **integrating interactivity in notebooks**

Alexander Lex

Mastodon: [@alexlex@vis.social](https://mastodon.social/@alexlex)

<http://alexander-lex.net>

Thanks to: **Kiran Gadhave, Zach Cutler**, Carolina Nobre, Marc Streit, Jochen Görtler, Oliver Deussen, Miriah Meyer, Jeff Phillips, Samuel Gratzl, Holger Stitz, Nils Gehlenborg, Hendrik Strobelt, Romain Vuillemot, Hanspeter Pfister, and many others!



visualization
design lab

