

Scorpion

Explaining Away Outliers in Aggregate Queries

eugene wu and sam madden
MIT

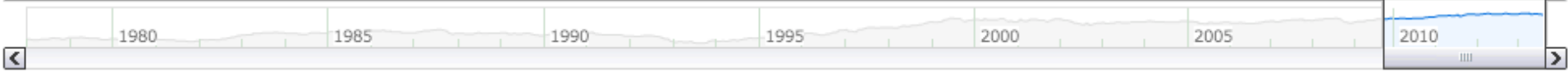


Compare: Dow Jones S&P 500 ORCL CSCO INTC SAP MSFT HPQ ACN

Zoom: [1d](#) [5d](#) [1m](#) [3m](#) [6m](#) [YTD](#) [1y](#) [5y](#) [10y](#) [All](#)

Sep 25, 2009 - Aug 22, 2013

● MSFT +29.91% ● ORCL +41.42% ● IBM +51.14%

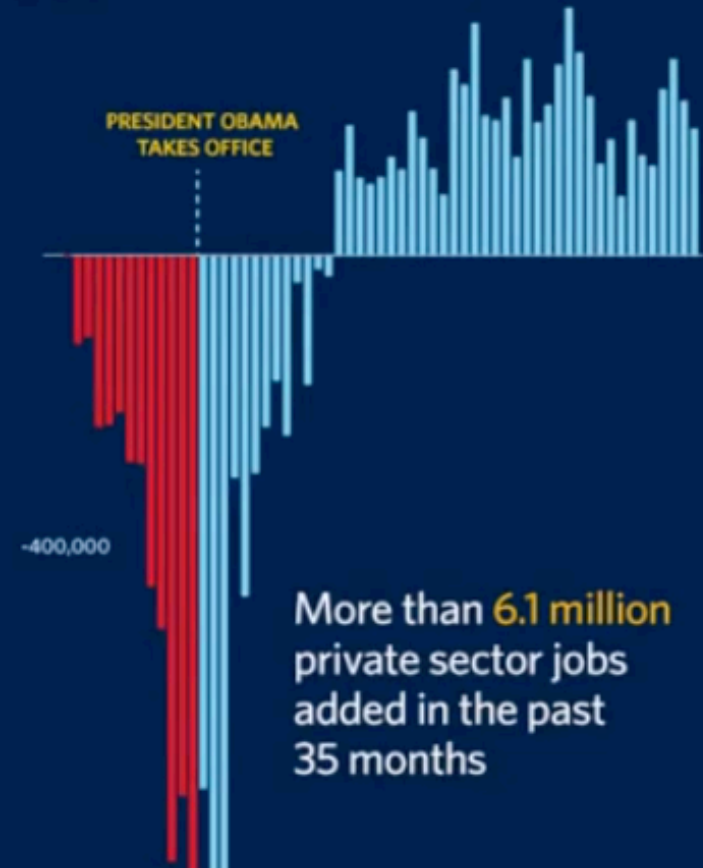




MONTHLY PRIVATE-SECTOR JOB LOSSES AND GAINS

350,000

PRESIDENT OBAMA
TAKES OFFICE



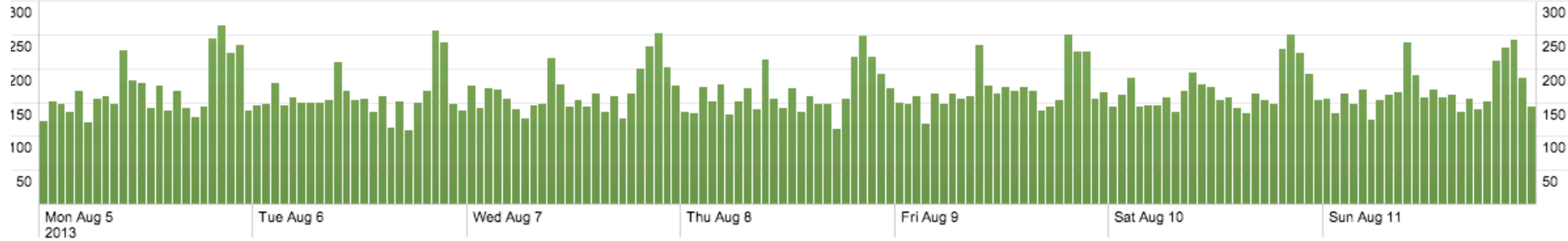
More than **6.1 million**
private sector jobs
added in the past
35 months

✓ 27,888 matching events

Navigation icons: Home, Stop, Refresh, Close, Info, Print, Save, Create

Hide Zoom out Zoom to selection Deselect

Linear scale 1 bar = 1 hour



Hide

27,888 events over all time

Export Options

« prev 1 2 3 4 5 6 7 8 9 10 next » 10 per page ▾

3 selected fields Edit

- host (1)
- source (1)
- sourcetype (1)

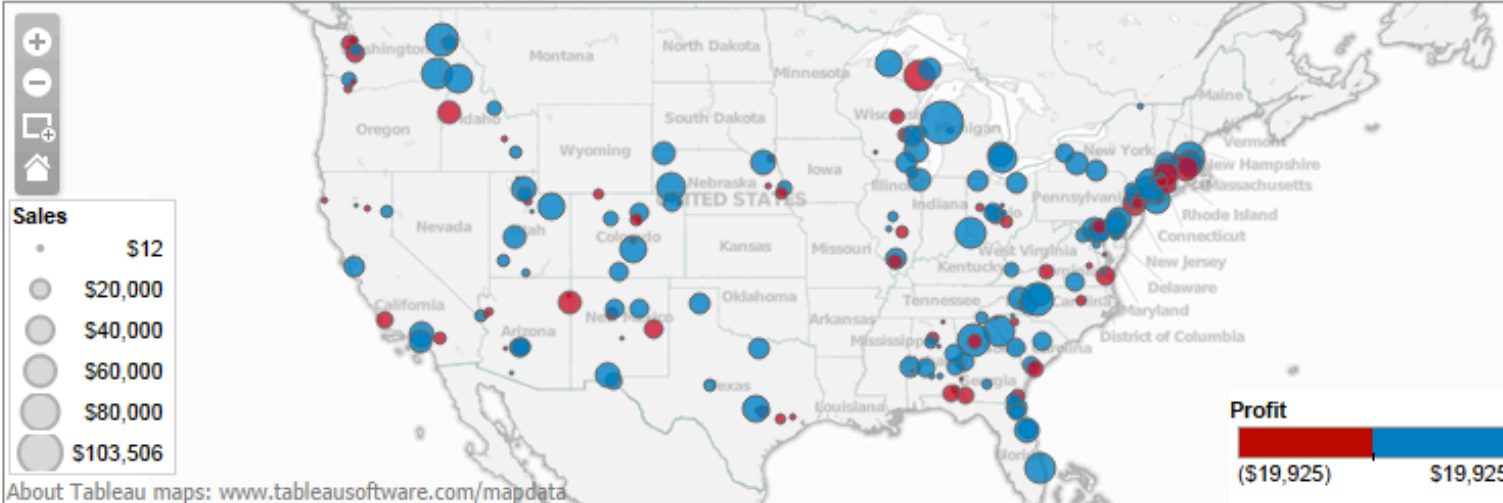
27 interesting fields

- action (2)
- # bytes (65)
- category_id (5)
- clientip (≥100)
- file (14)
- ident (1)

1	8/11/13 11:59:34.000 PM	178.19.3.39 - - [11/Aug/2013:23:59:34] "GET /flower_store/category.screen?category_id=CANDY HTTP/1.1" 200 10567 "http://mystore.splunk.com/flower_store/cart.do?action=purchase&itemId=EST-14&JSESSIONID=SD5SL10FF8ADFF3" "Mozilla/5.0 (X11; U; Linux i686; en-US; rv:1.8.0.10) Gecko/20070223 CentOS/1.5.0.10-0.1.e14.centos Firefox/1.5.0.10" 3187 1245 host=li-desktop sourcetype=access_combined_wcookie source=/Users/sirrice/Downloads/SampledData/apache3.splunk.com/access_combined.log ▾
2	8/11/13 11:59:15.000 PM	10.192.1.46 - - [11/Aug/2013:23:59:15] "GET /flower_store/category.screen?category_id=PLANTS HTTP/1.1" 200 10567 "http://mystore.splunk.com/flower_store/cart.do?action=purchase&itemId=EST-27&JSESSIONID=SD5SL10FF8ADFF3" "Mozilla/5.0 (X11; U; Linux i686; en-US; rv:1.8.0.10) Gecko/20070223 CentOS/1.5.0.10-0.1.e14.centos Firefox/1.5.0.10" 3352 1920 host=li-desktop sourcetype=access_combined_wcookie source=/Users/sirrice/Downloads/SampledData/apache3.splunk.com/access_combined.log ▾
3	8/11/13 11:59:15.000 PM	10.192.1.46 - - [11/Aug/2013:23:59:15] "GET /flower_store/category.screen?category_id=PLANTS HTTP/1.1" 200 10567 "http://mystore.splunk.com/flower_store/cart.do?action=purchase&itemId=EST-27&JSESSIONID=SD5SL10FF8ADFF3"

Executive Dashboard

Sales by Customer Location



Select Year:

(All)

Customer Region:

(All)

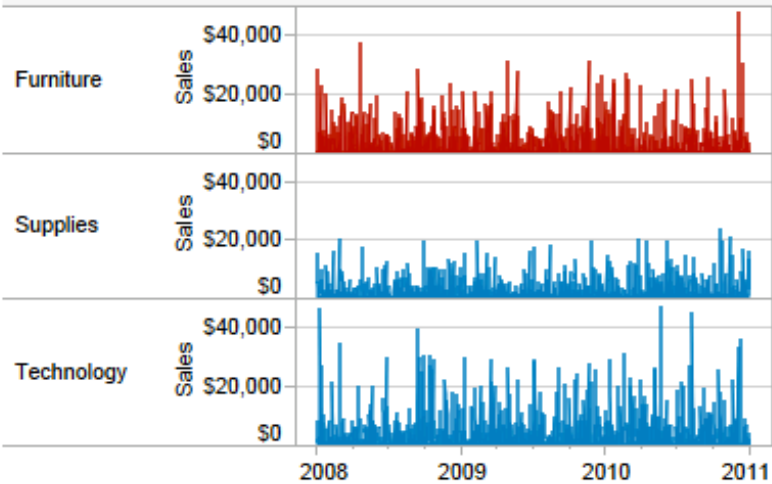
Product Category:

(All)

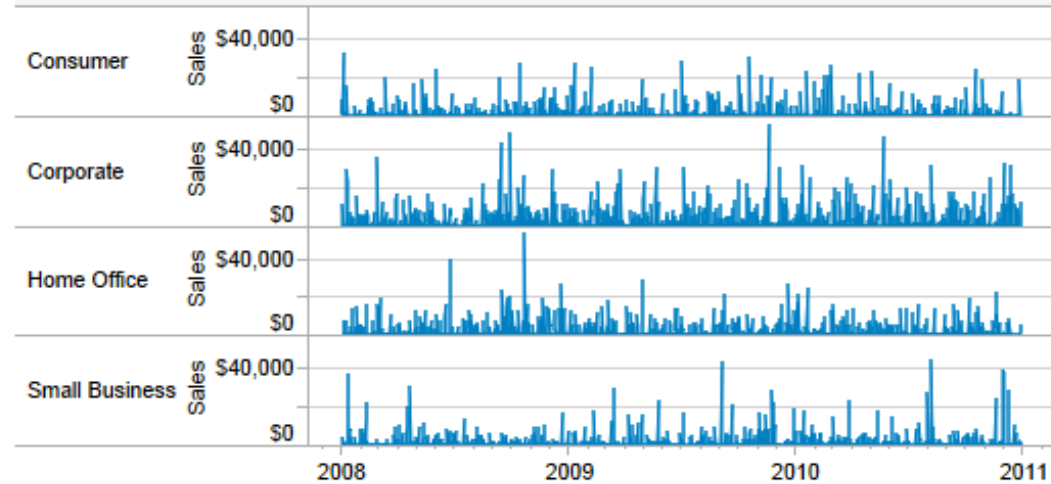
Customer Segment:

(All)

Sales by Product Category



Sales by Customer Segment



Find reports & more

Audience Overview

Jul 17, 2013 - Aug 16, 2013

- MY STUFF
- Dashboards
- Shortcuts
- Intelligence Events

STANDARD REPORTS

- Real-Time
- Audience
 - Overview
 - Demographics
 - Behavior
 - Technology
 - Mobile
 - Custom
 - Visitors Flow
- Traffic Sources
- Content
- Conversions

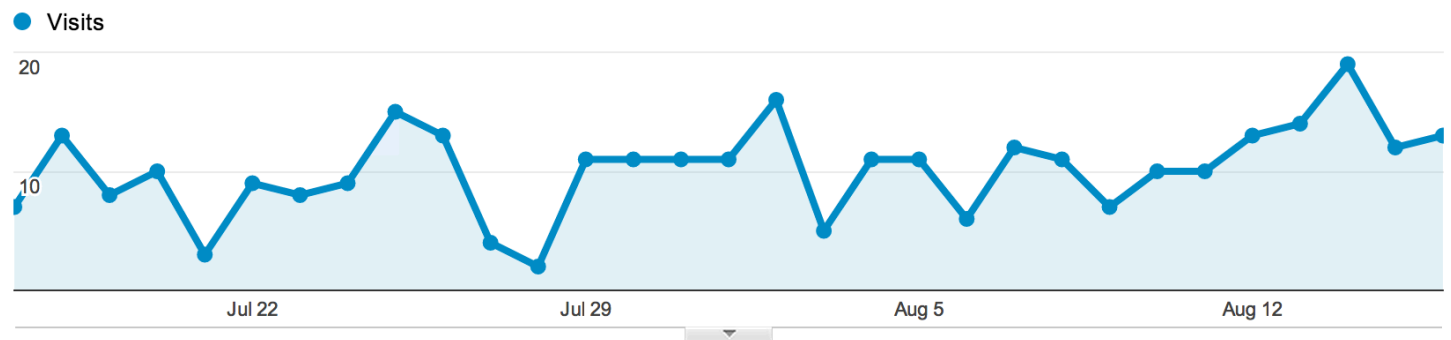
Advanced Segments | Email | Export | Add to Dashboard | Shortcut

% of visits: 100.00%

Overview

Visits vs. [Select a metric](#)

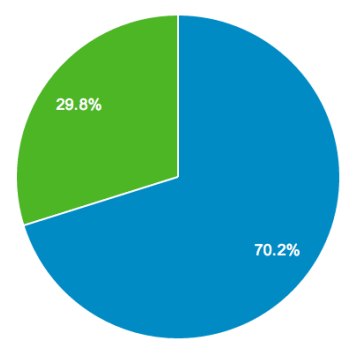
Hourly | Day | Week | Month



237 people visited this site



New Visitor | Returning Visitor





Table



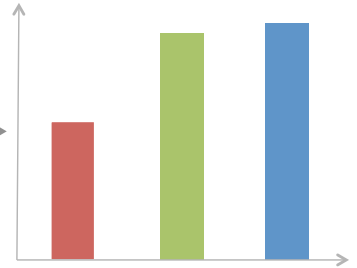
Split →

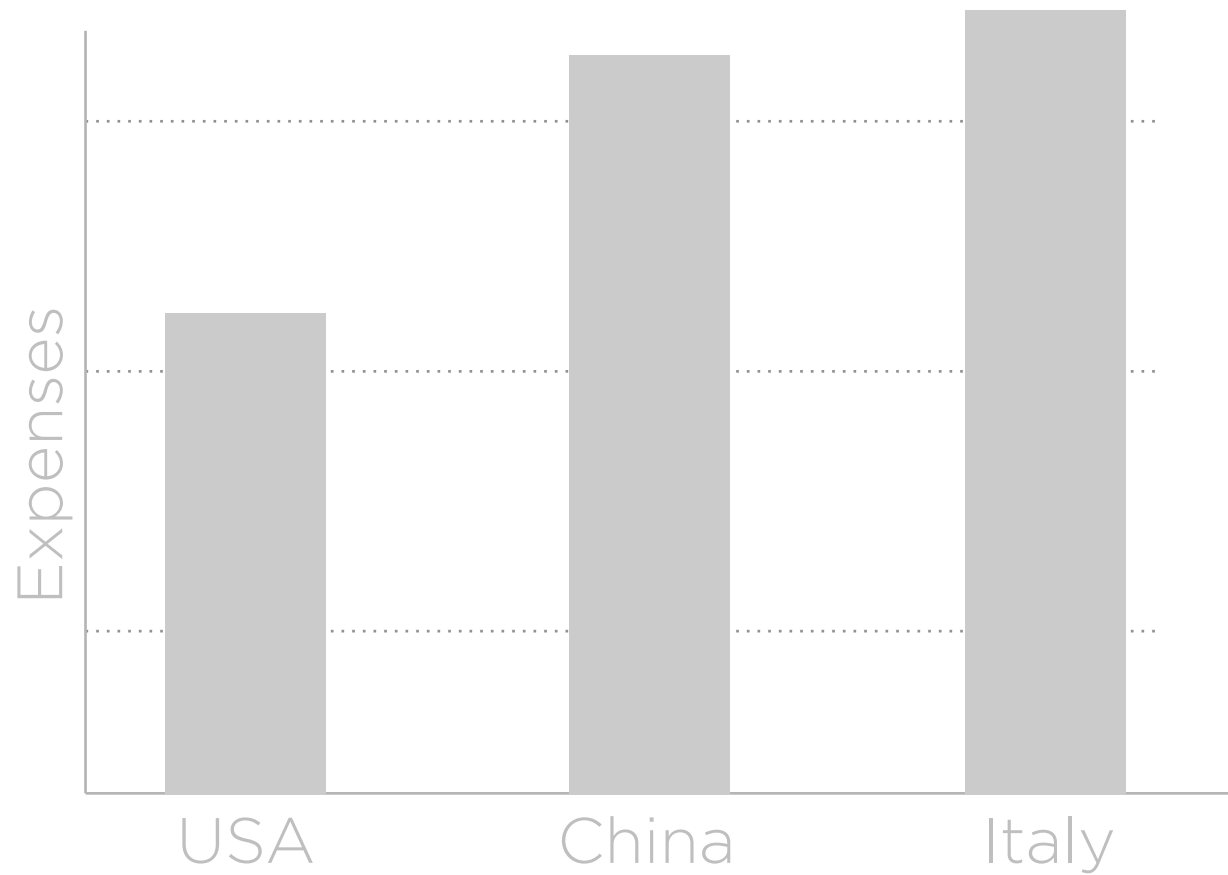


Aggregate →

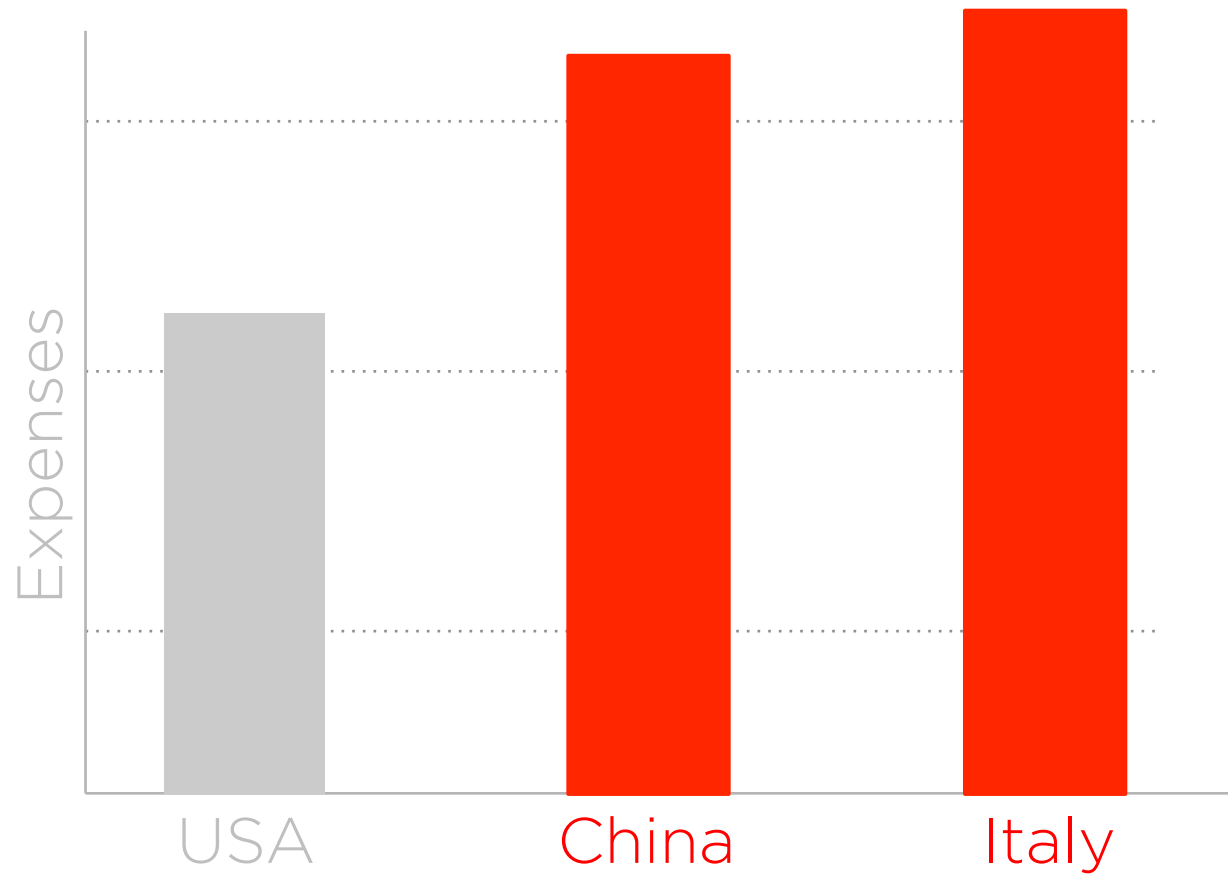


Visualize →

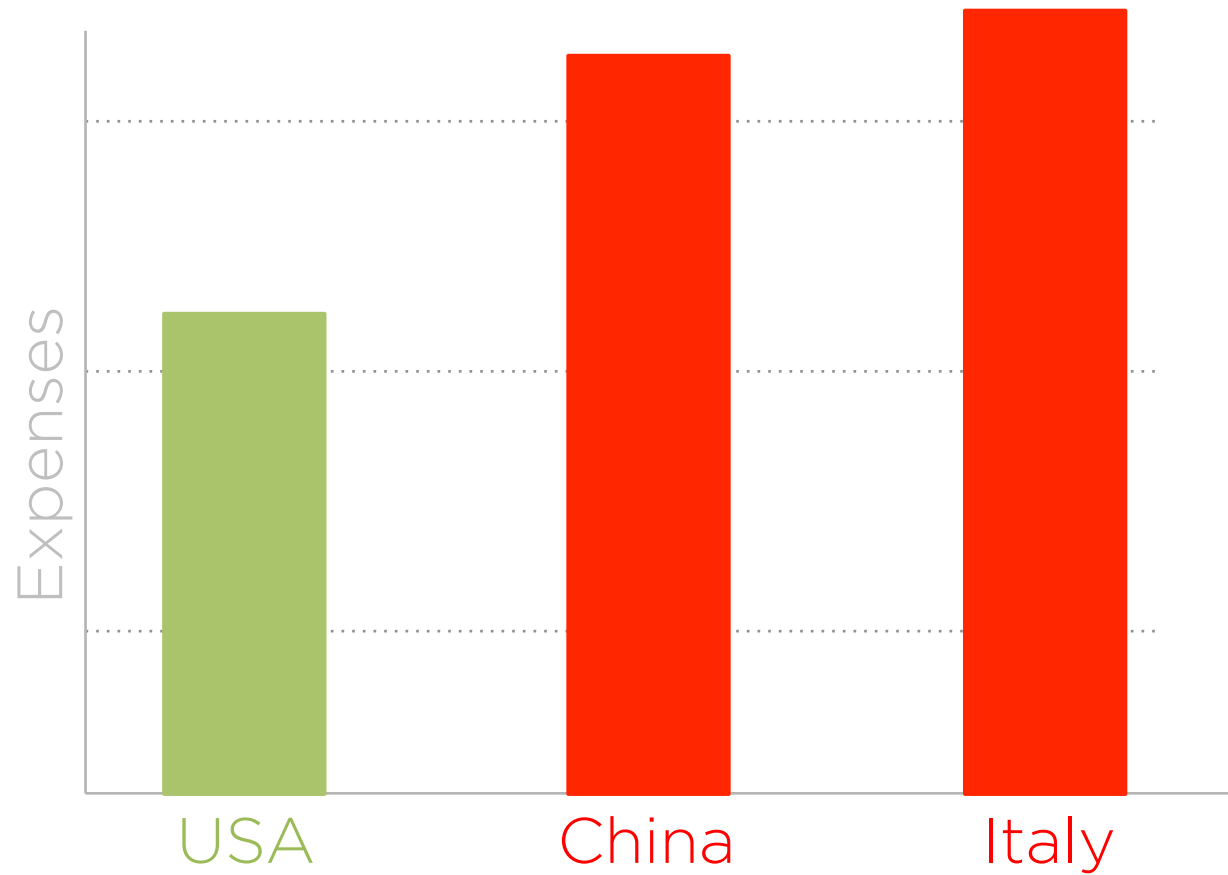




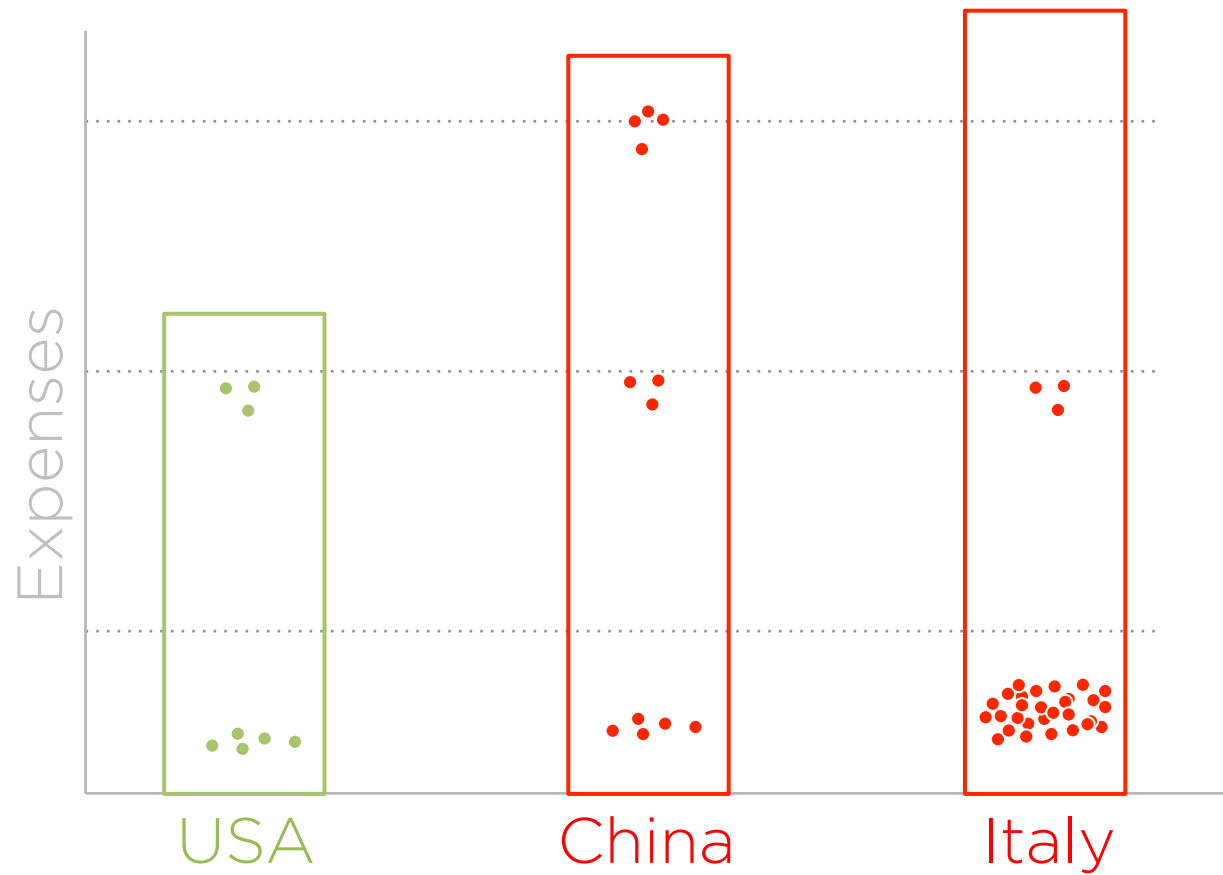
```
SELECT sum(cost)
FROM expenses
GROUPBY country
```



```
SELECT sum(cost)
FROM expenses
GROUPBY country
```



```
SELECT sum(cost)
FROM expenses
GROUPBY country
```

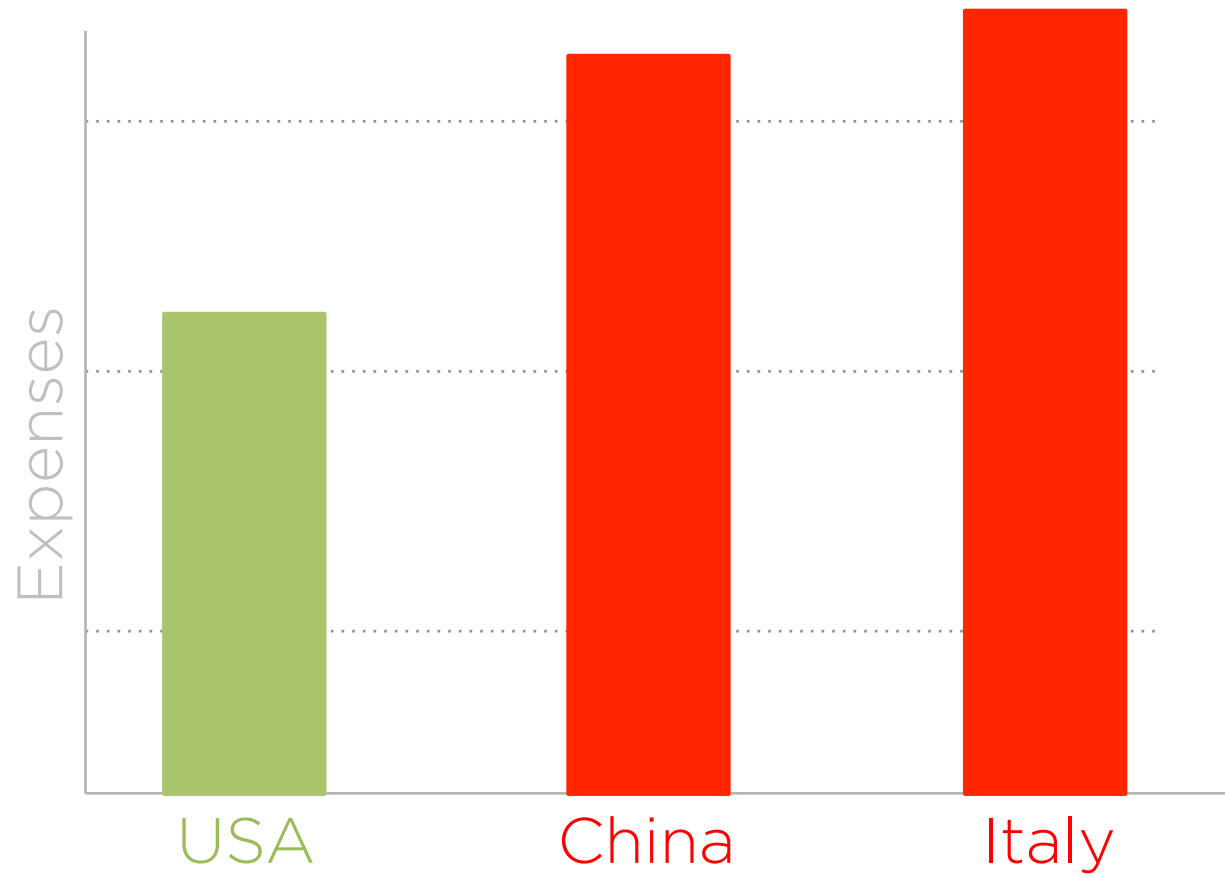


```
SELECT sum(cost)
FROM expenses
GROUPBY country
```

Given

Outlier and normal results

Understand Why



```
SELECT sum(cost)
FROM expenses
GROUPBY country
```

Given

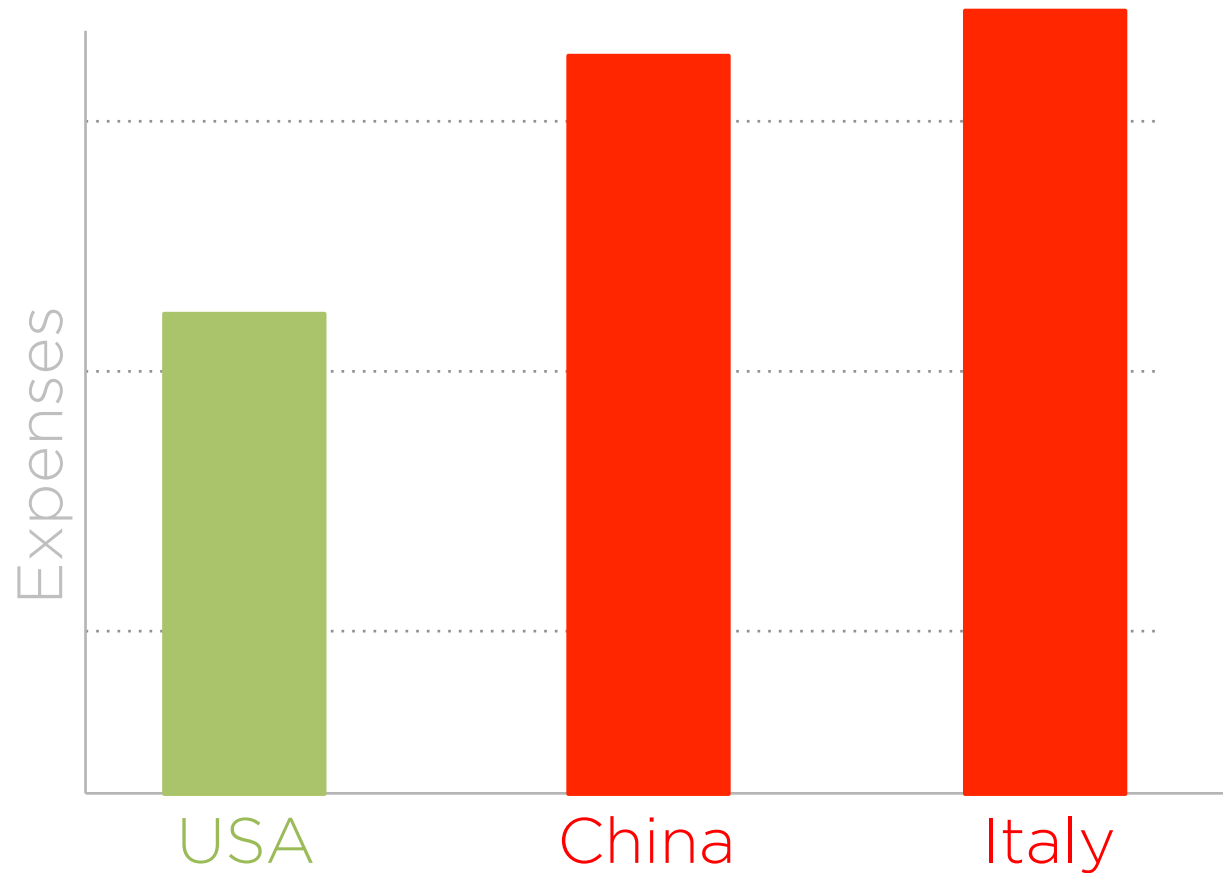
Outlier and normal results

What input properties

caused the outliers?

most caused the outliers?

caused outliers but didn't
affect normal outputs?



```
SELECT sum(cost)
FROM expenses
GROUPBY country
```

Can't Touch This



**I DON'T ALWAYS SOLVE
PROBLEMS**



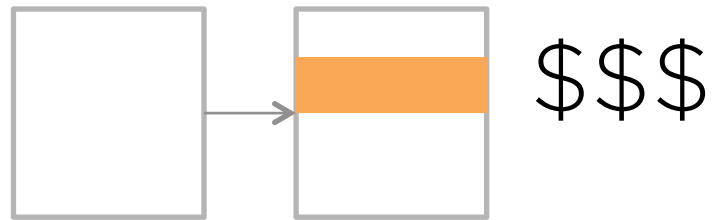
**BUT WHEN I DO.
I DATABASE**

Provenance



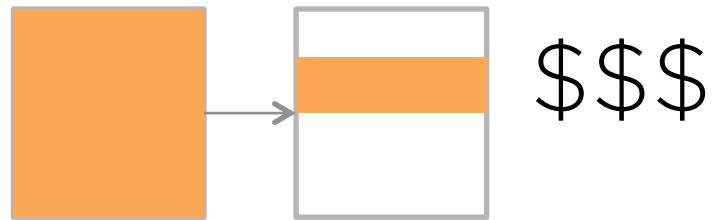
Provenance

```
SELECT SUM(cost)  
FROM sam's bank account
```



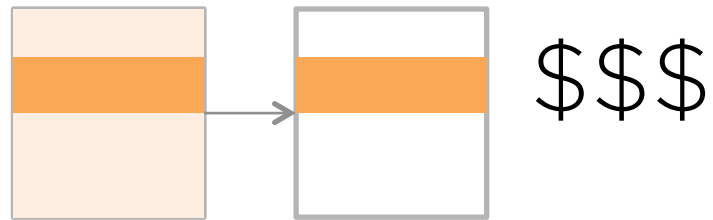
Provenance

```
SELECT SUM(cost)  
FROM sam's bank account
```



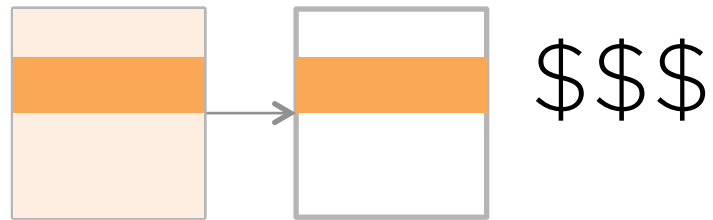
Provenance

```
SELECT SUM(cost)  
FROM sam's bank account
```



Provenance

```
SELECT SUM(cost)  
FROM sam's bank account
```



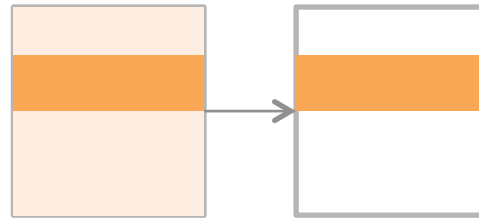
Darn!
Ya caught me

Proven



Provenance

```
SELECT SUM(cost)
FROM sam's bank account
```

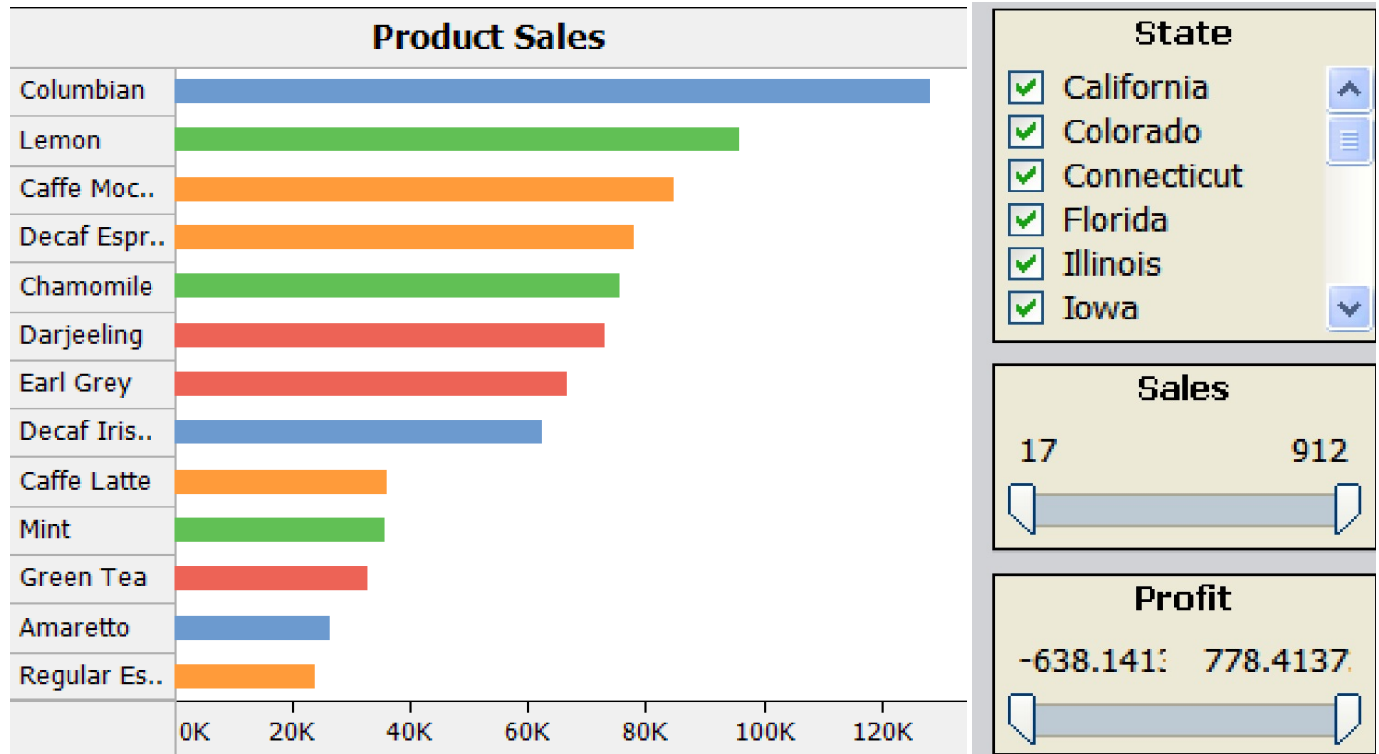


Filter for
“most influential”

~~Provenance~~

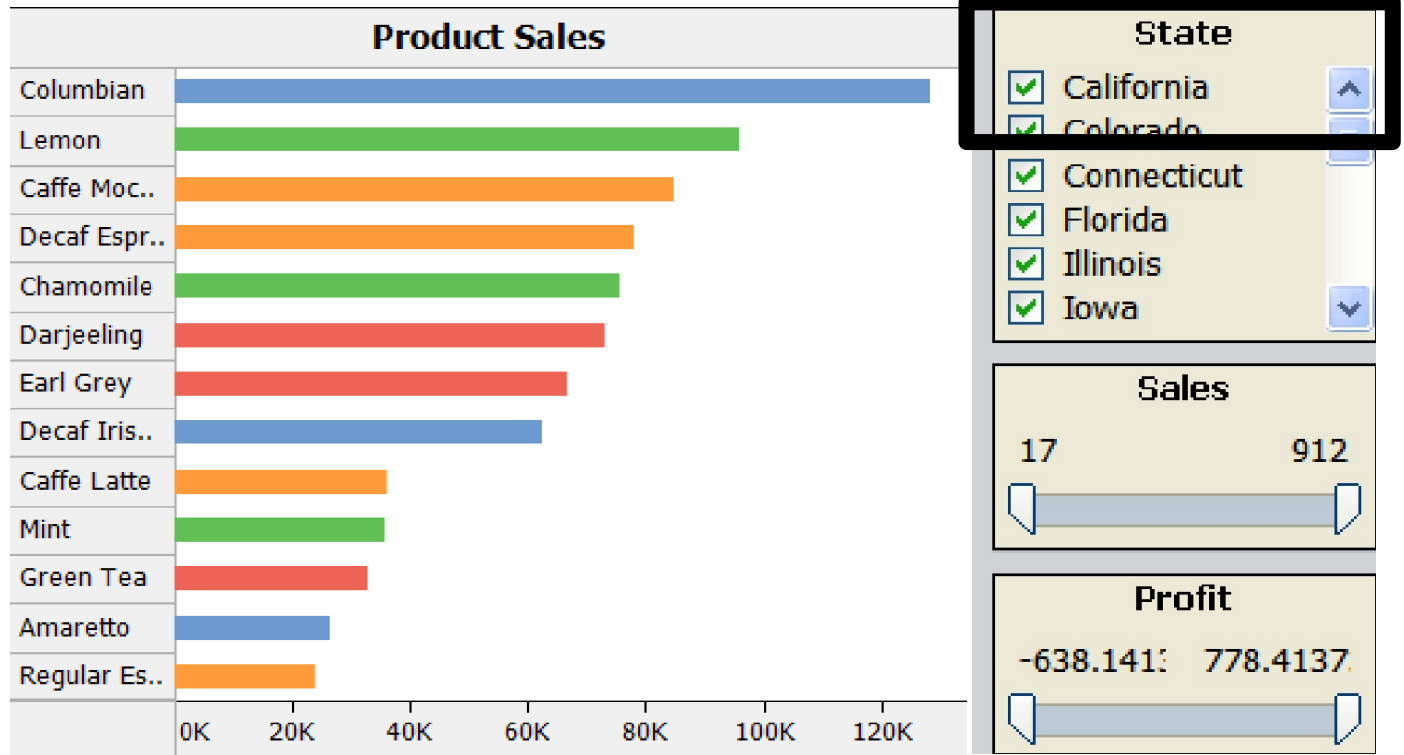


Provenance



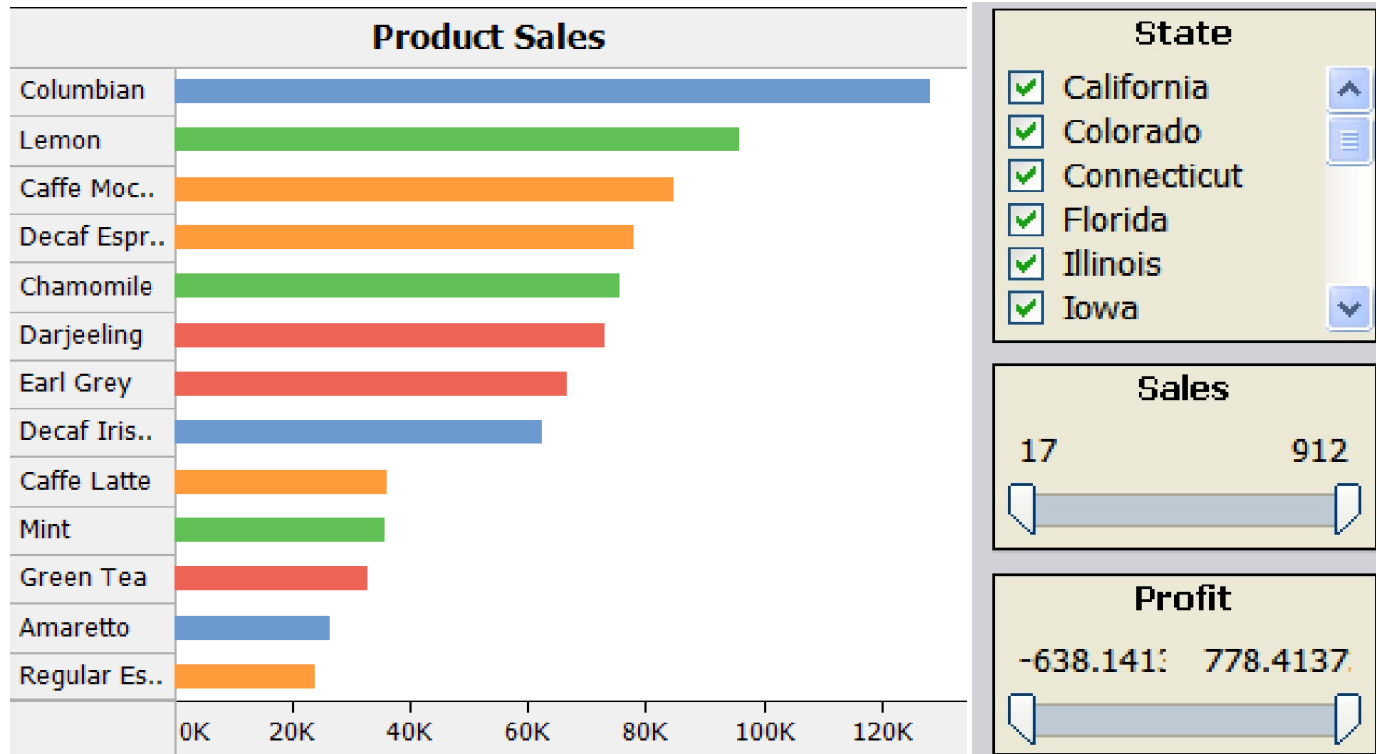
Faceting

Provenance



Faceting

Provenance



Faceting

Dimensionality :(

Dealing with multiple outliers?

~~Provenance~~

Faceting



~~Provenance~~

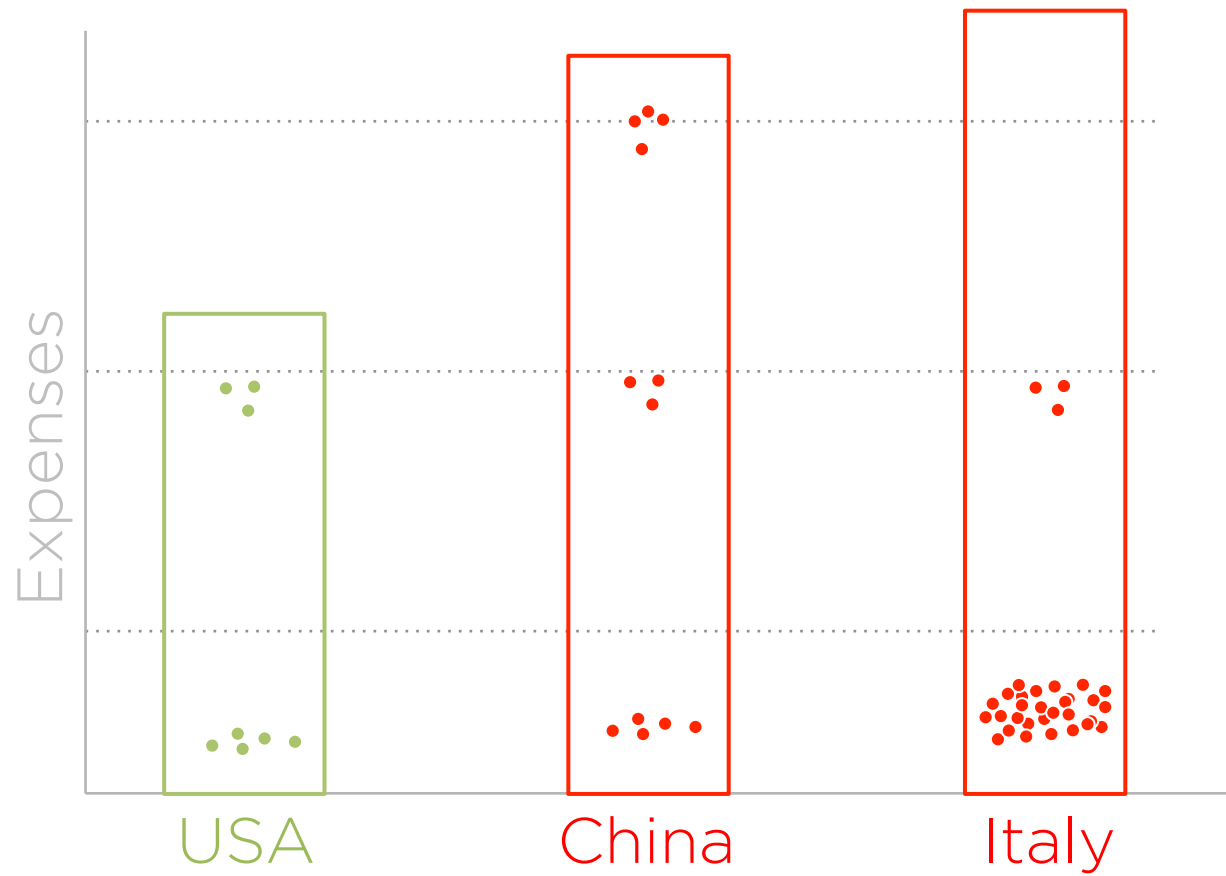
Faceting

Scorpion!

Given

Outlier and normal results

Understand Why

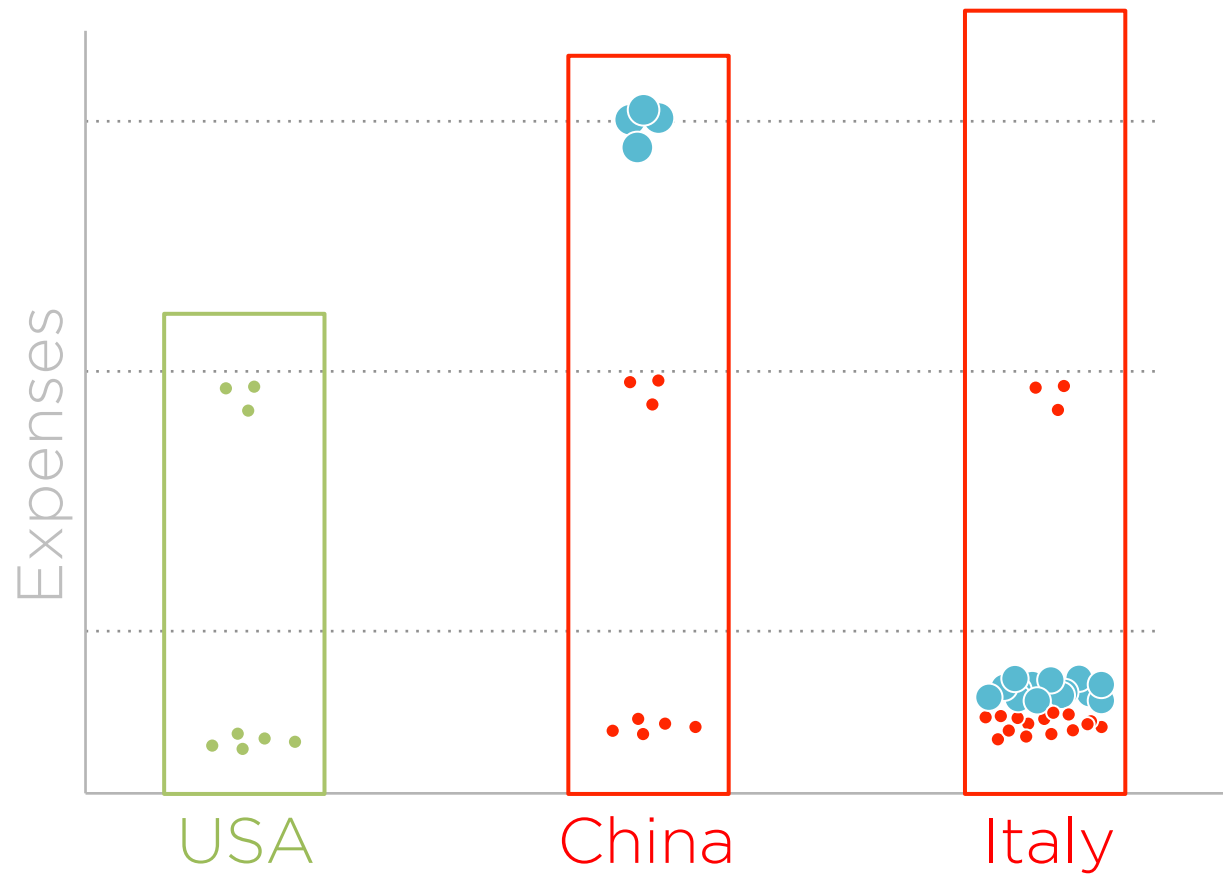


Given

Outlier and normal results

Find

Predicates correlated with outliers



Desc = "toilets"

Given

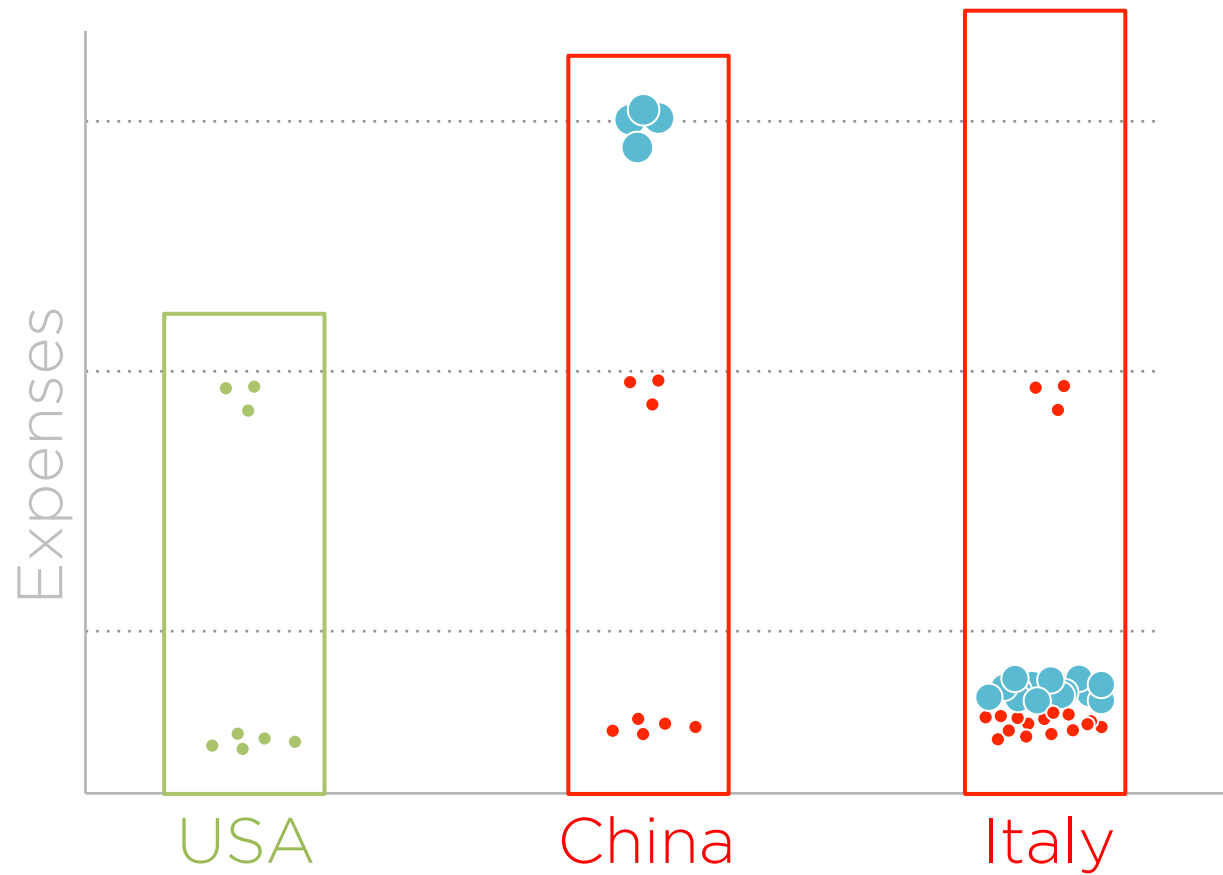
Outlier and normal results

Find

Predicates correlated with outliers

s.t.

Removing predicate from inputs “fixes” outliers & maintains normal results



Desc = “toilets”

Given

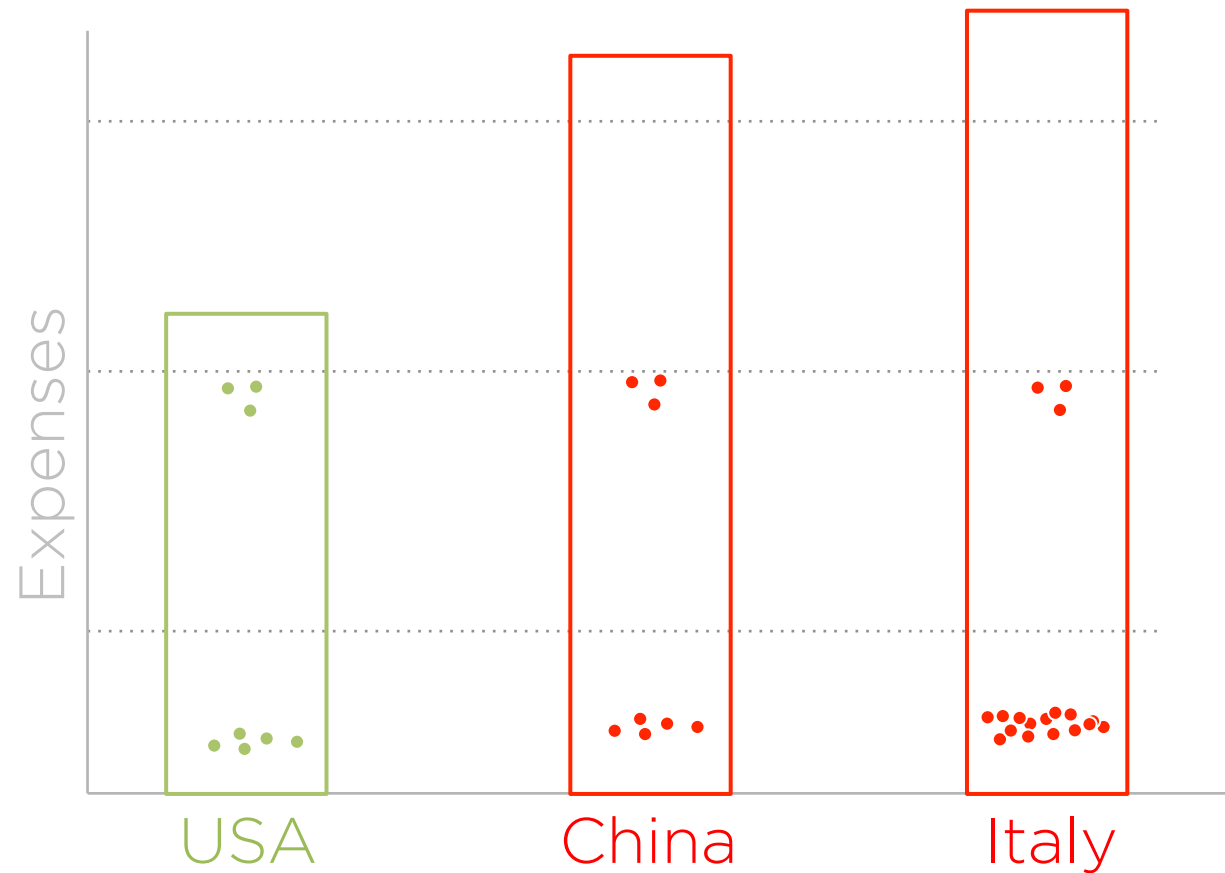
Outlier and normal results

Find

Predicates correlated with outliers

s.t.

Removing predicate from inputs “fixes” outliers & maintains normal results



~~Dese = "toilets"~~

Given

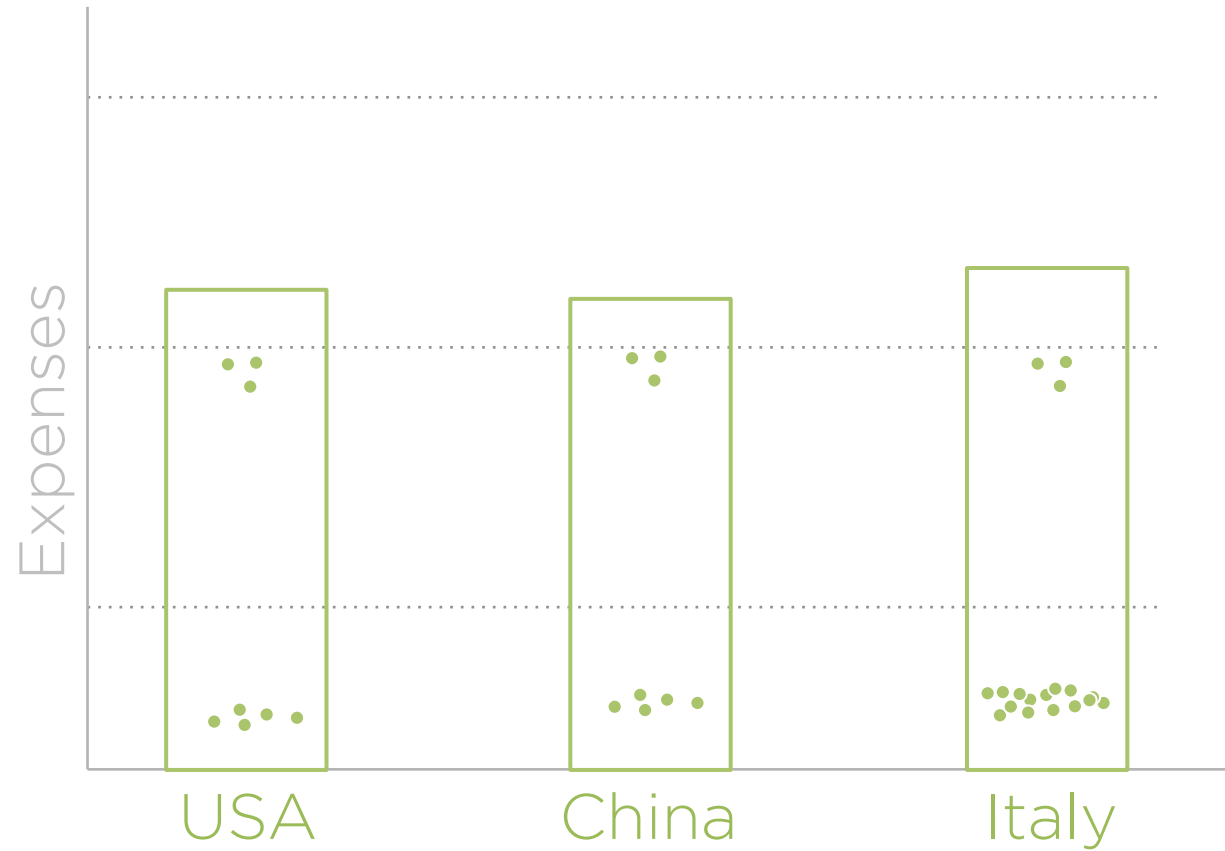
Outlier and normal results

Find

Predicates correlated with outliers

s.t.

Removing predicate from inputs “fixes” outliers & maintains normal results



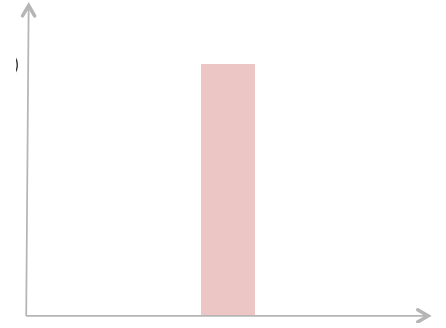
~~Dese = "toilets"~~

Formalize “influence” as metric

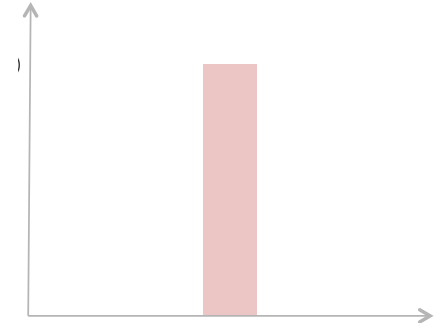
Predicate search heuristics

Some results

T



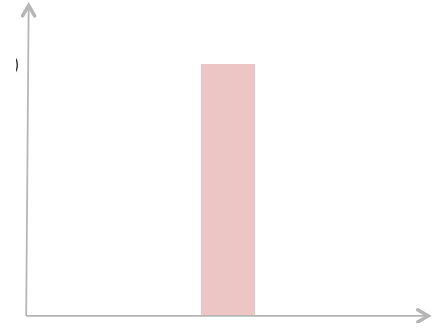
T



p

Desc = "toilet"

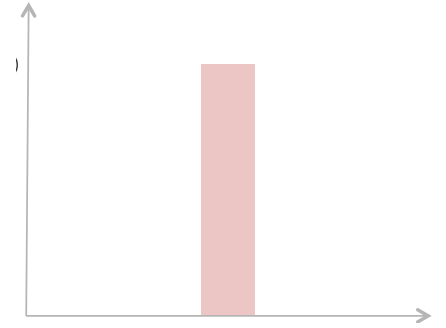
T



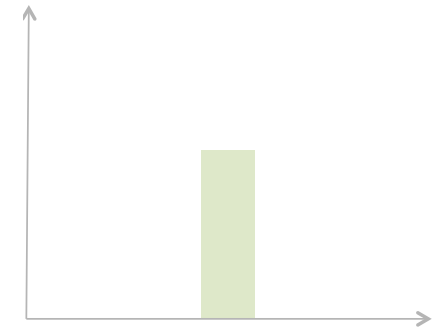
$p(T)$



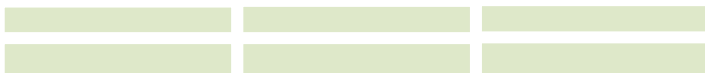
T

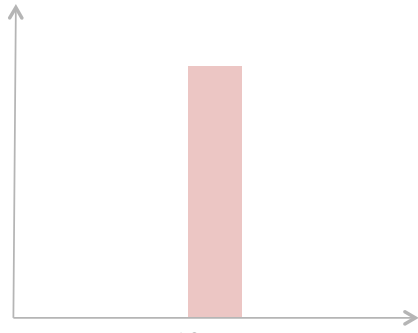


p(T)

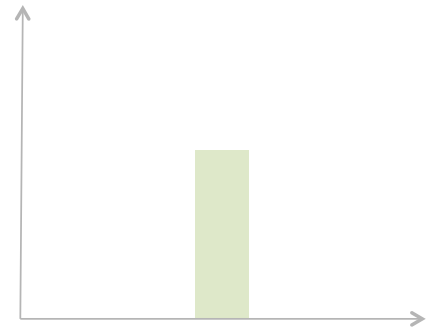


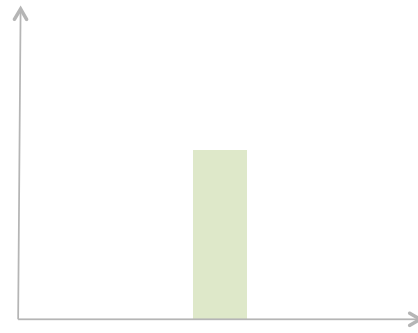
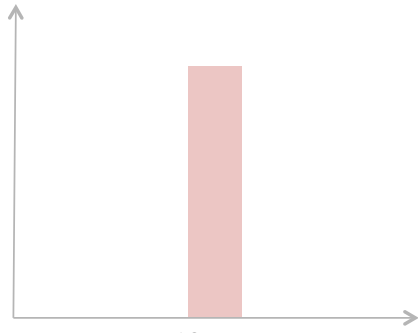
T - p(T)





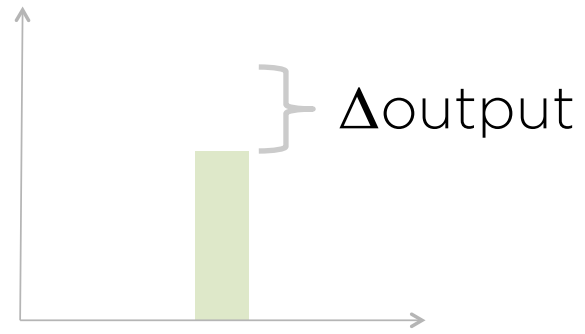
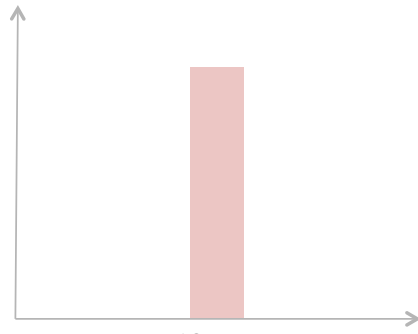
p(T)



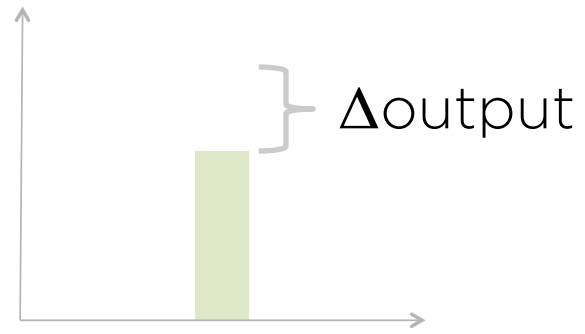
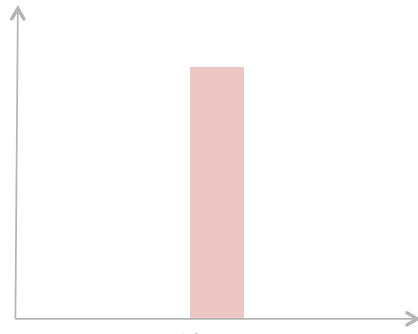


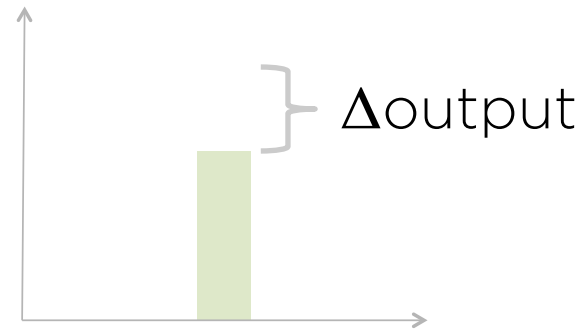
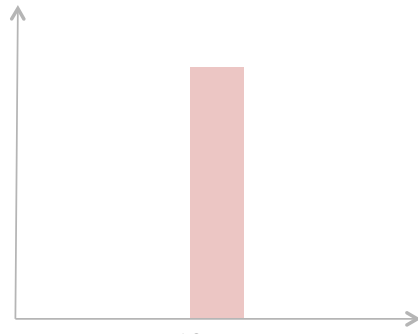
p(T)





p(T) 





$$\frac{\Delta\text{output}}{|p(T)|}$$

**Influence
Metric**

$$\frac{\Delta f(x)}{\Delta x}$$

**Sensitivity
Analysis**

$$\frac{\Delta \text{output}}{|\rho(T)|}$$

**Influence
Metric**

Δ Output

$$\frac{\Delta \text{output}}{|p(T)|}$$

“High vs Low”

$|p(T)|$

Δ Normal

Multiple Outputs

Δ Output

$$\frac{\Delta \text{output}}{|p(T)|}$$

“High vs Low”

$$\frac{\Delta \text{output} \cdot V}{|p(T)|}$$

$|p(T)|$

Δ Normal

Multiple Outputs

Δ Output

$$\frac{\Delta \text{output}}{|p(T)|}$$

“High vs Low”

$$\frac{\Delta \text{output} \cdot V}{|p(T)|}$$

$|p(T)|$

$$\frac{\Delta \text{output} \cdot V}{|p(T)|^c}$$

Δ Normal

Multiple Outputs

Δ Output

$$\frac{\Delta_{\text{output}}}{|p(T)|}$$

“High vs Low”

$$\frac{\Delta_{\text{output}} \cdot V}{|p(T)|}$$

$|p(T)|$

$$\frac{\Delta_{\text{output}} \cdot V}{|p(T)|^c}$$

Δ Normal

$$\frac{\Delta_{\text{outlier}} \cdot V}{|p(T)|^c} - |\Delta_{\text{Normal}}|$$

Multiple Outputs

Δ Output

$$\frac{\Delta_{\text{output}}}{|p(T)|}$$

“High vs Low”

$$\frac{\Delta_{\text{output}} \cdot V}{|p(T)|}$$

$|p(T)|$

$$\frac{\Delta_{\text{output}} \cdot V}{|p(T)|^c}$$

Δ Normal

$$\frac{\Delta_{\text{outlier}} \cdot V}{|p(T)|^c} - |\Delta_{\text{Normal}}|$$

Multiple Outputs

$$\text{mean}_{\text{outlier}} \frac{\Delta_{\text{outlier}} \cdot V}{|p(T)|^c} - \text{max}_{\text{normal}} |\Delta_{\text{Normal}}|$$

Δ output

$$\frac{\Delta_{\text{outlier}}}{|P(T)|}$$

“High vs Low”

$$\frac{\Delta_{\text{outlier}} \cdot V}{|P(T)|}$$

influence(p)

Δ Normal

$$\frac{\Delta_{\text{outlier}} \cdot V}{|P(T)|^c} - |\Delta_{\text{Hold-out}}|$$

Multiple Outputs

$$\text{mean}_{\text{outlier}} \frac{\Delta_{\text{outlier}} \cdot V}{|P(T)|^c} - \text{max}_{\text{normal}} |\Delta_{\text{Hold-out}}|$$

Formalize “influence” as metric

Predicate search heuristics

Some results

$$p^* = \underset{p \in \text{predicates}}{\operatorname{argmax}} \text{influence}(p)$$

$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \underbrace{\text{influence}(p)}_{O(\text{agg}(T-p(T)))}$$

$$\text{SUM}(\{1,2,3,4,5\}) = 15$$

$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$


 $O(\text{agg}(T-p(T)))$

$$\text{SUM}(\{1,2,3,\overbrace{4,5}^p\}) = 15$$

$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$

$O(\text{agg}(T-p(T)))$

$$\text{SUM}(\{1,2,3,\overbrace{4,5}^p\}) = 15$$

$\underline{\quad}$
 $\{4,5\}$

$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$

$\underbrace{\hspace{10em}}$
 $O(\text{agg}(T-p(T)))$

$$\text{SUM}(\{1,2,3,\overbrace{4,5}^p\}) = 15$$

$$\text{SUM}(\{1,2,3\}) = 6$$

$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$

$$O(\text{agg}(T - p(T)))$$

$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$

$O(\text{exponential})$ $O(\text{agg}(T-p(T)))$

Operator Properties

$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$

$\underbrace{\hspace{10em}}_{O(\text{exponential})} \quad \underbrace{\hspace{10em}}_{O(\text{agg}(T-p(T)))}$

Operator Properties

$$p^* = \underset{p \in \text{predicates}}{\operatorname{argmax}} \text{influence}(p)$$

$O(\text{exponential})$

$O(\text{agg}(p(T)))$

**Incrementally
removable**

$$\text{SUM}(\{1,2,3,\overbrace{4,5}^p\}) = 15$$

$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$

$O(\text{exponential})$ $O(\text{agg}(p(T)))$

**Incrementally
removable**

$$\text{SUM}(\{1,2,3,\overbrace{4,5}^p\}) = 15$$

$$15 - \text{SUM}(\{4,5\}) = 6$$

$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$

$O(\text{exponential})$ $O(\text{agg}(p(T)))$

**Incrementally
removable**

SUM
COUNT
AVG
STDDEV

$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$

$O(\text{exponential})$ $O(\text{agg}(p(T)))$

**Incrementally
removable**

~~MEDIAN~~

~~MODE~~

SUM
COUNT
AVG
STDDEV

$$p^* = \underset{p \in \text{predicates}}{\operatorname{argmax}} \text{influence}(p)$$

$O(\text{exponential})$

$O(\text{agg}(p(T)))$

**Incrementally
removable**

Least
influence



Most
influence

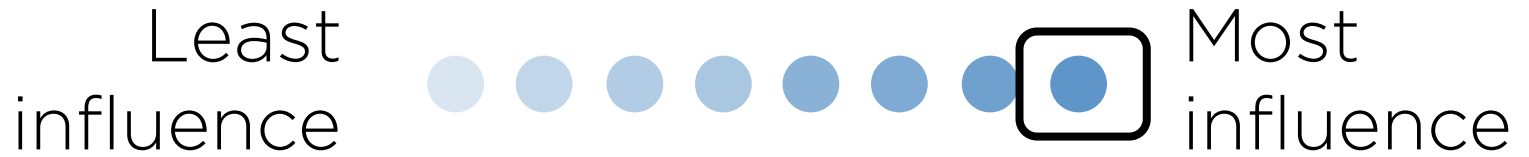
$$p^* = \underset{p \in \text{predicates}}{\operatorname{argmax}} \text{influence}(p)$$

$O(\text{exponential})$

Independent

$O(\text{agg}(p(T)))$

**Incrementally
removable**



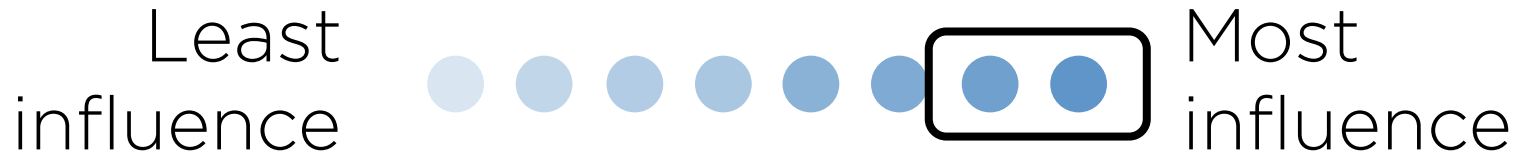
$$p^* = \underset{p \in \text{predicates}}{\operatorname{argmax}} \text{influence}(p)$$

$O(\text{exponential})$

Independent

$O(\text{agg}(p(T)))$

**Incrementally
removable**



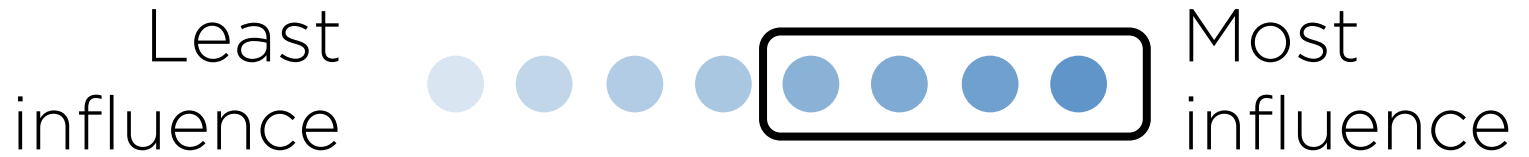
$$p^* = \underset{p \in \text{predicates}}{\operatorname{argmax}} \text{influence}(p)$$

$O(\text{exponential})$

Independent

$O(\text{agg}(p(T)))$

**Incrementally
removable**



$$p^* = \underset{p \in \text{predicates}}{\operatorname{argmax}} \text{influence}(p)$$

$O(\text{exponential})$

Independent

$O(\text{agg}(p(T)))$

**Incrementally
removable**



$$p^* = \underset{p \in \text{predicates}}{\operatorname{argmax}} \text{influence}(p)$$

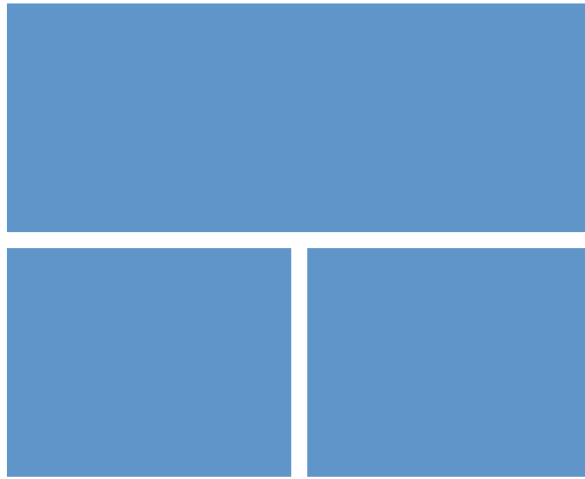
$O(\text{exponential})$

$O(\text{agg}(p(T)))$

Top Down

Independent

**Incrementally
removable**



$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$

$O(\text{exponential})$ $O(\text{agg}(p(T)))$

Top Down

Independent

**Incrementally
removable**



$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$

$O(\text{exponential})$

$O(\text{agg}(p(T)))$

Top Down

Independent

**Incrementally
removable**



$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$

$O(\text{exponential})$

$O(\text{agg}(p(T)))$

Top Down

Independent

**Incrementally
removable**

$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$

$O(\text{exponential})$

$O(\text{agg}(p(T)))$

Top Down

Independent

**Incrementally
removable**

Anti-monotonic

$$p' \subset p$$

$$p^* = \underset{p \in \text{predicates}}{\operatorname{argmax}} \text{influence}(p)$$

$O(\text{exponential})$

$O(\text{agg}(p(T)))$

Top Down

Independent

**Incrementally
removable**

Anti-monotonic

$$p' \subset p$$



$$\text{influence}(p') \leq \text{influence}(p)$$

$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$

$O(\text{exponential})$

$O(\text{agg}(p(T)))$

Top Down

Independent

**Incrementally
removable**

Anti-monotonic



$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$

$O(\text{exponential})$

$O(\text{agg}(p(T)))$

Top Down

Independent

**Incrementally
removable**

Bottom Up

Anti-monotonic



$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$

$O(\text{exponential})$

$O(\text{agg}(p(T)))$

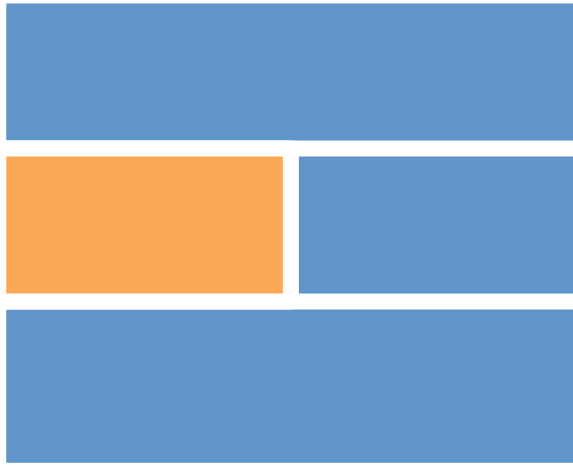
Top Down

Independent

**Incrementally
removable**

Bottom Up

Anti-monotonic



$$p^* = \underset{p \in \text{predicates}}{\text{argmax}} \text{influence}(p)$$

$O(\text{exponential})$

$O(\text{agg}(p(T)))$

Top Down

Independent

**Incrementally
removable**

Bottom Up

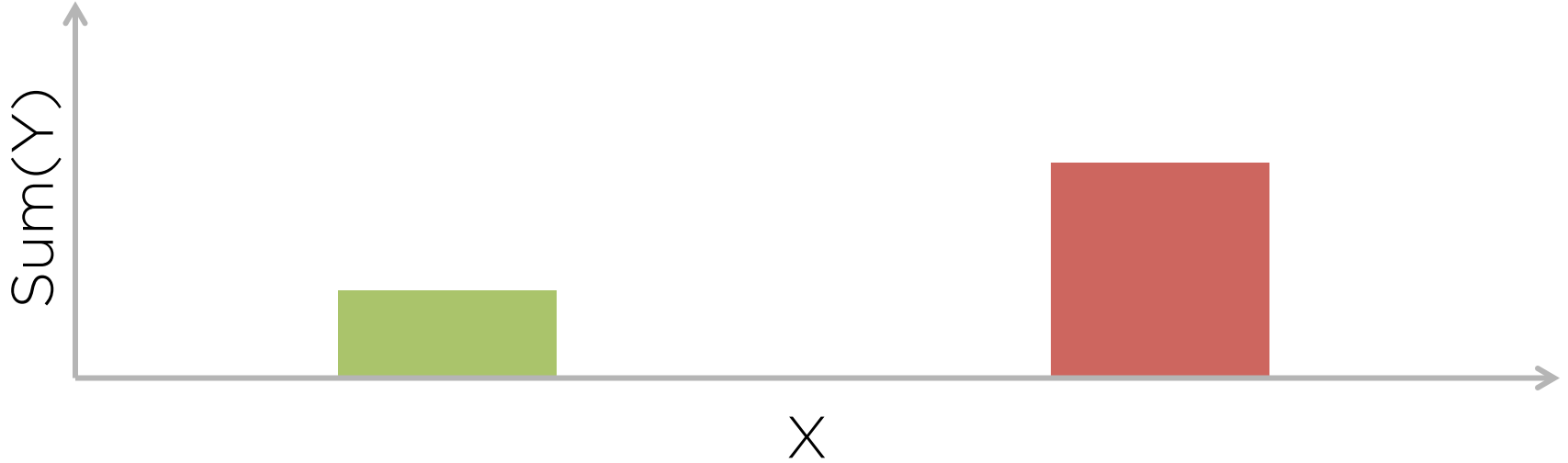
Anti-monotonic

Formalize “influence” as metric

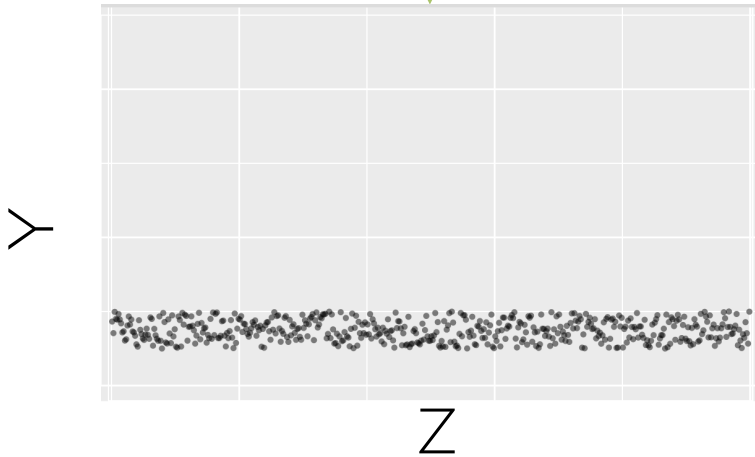
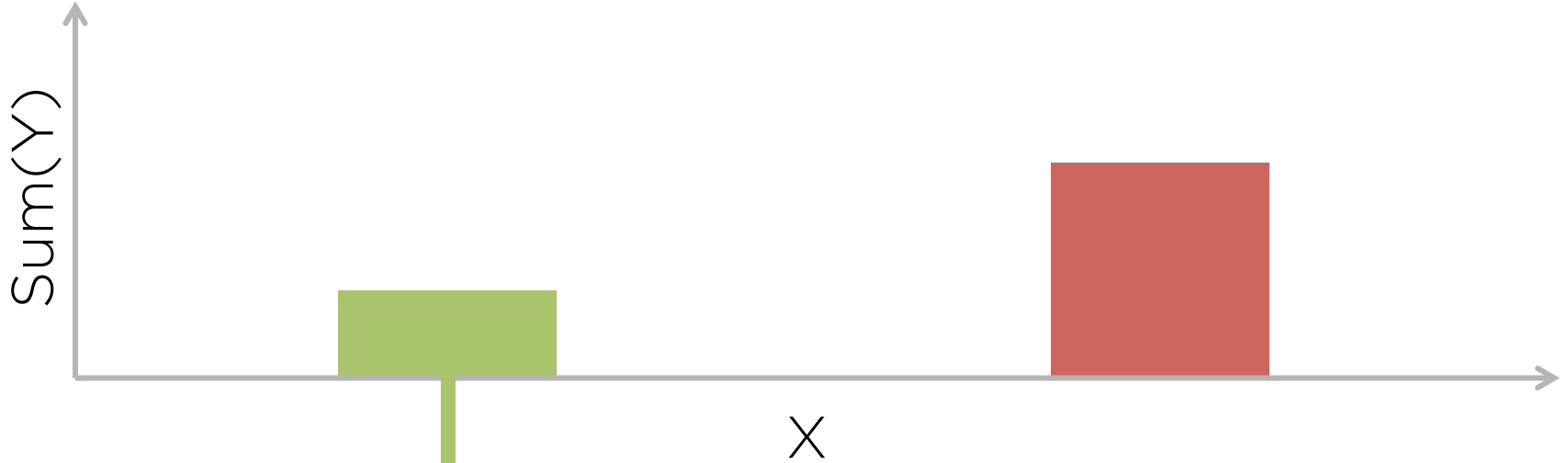
Predicate search heuristics

Some results

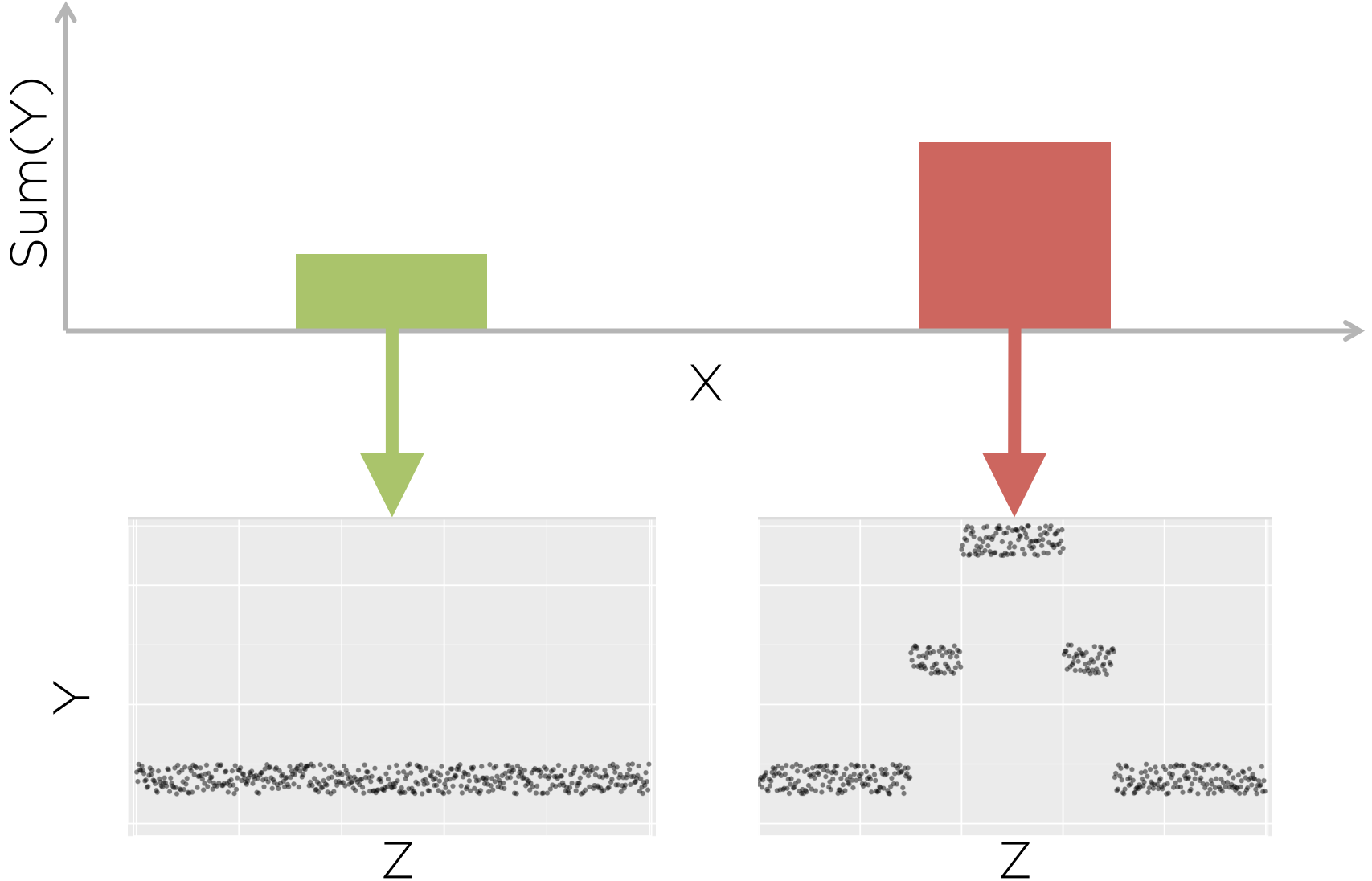

```
SELECT sum(Y) GROUPBY X
```



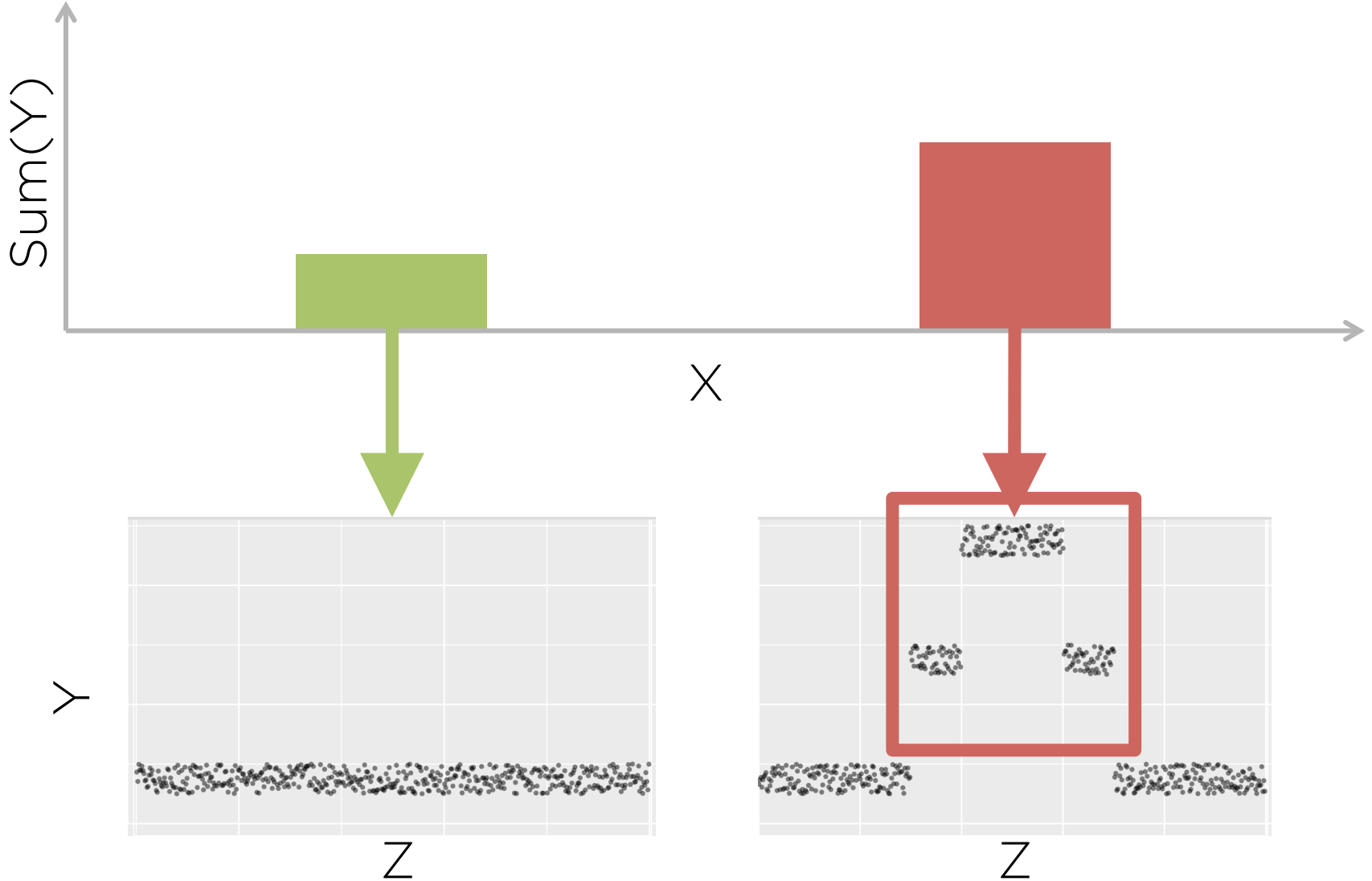
SELECT sum(Y) GROUPBY X



SELECT sum(Y) GROUPBY X



SELECT sum(Y) GROUPBY X

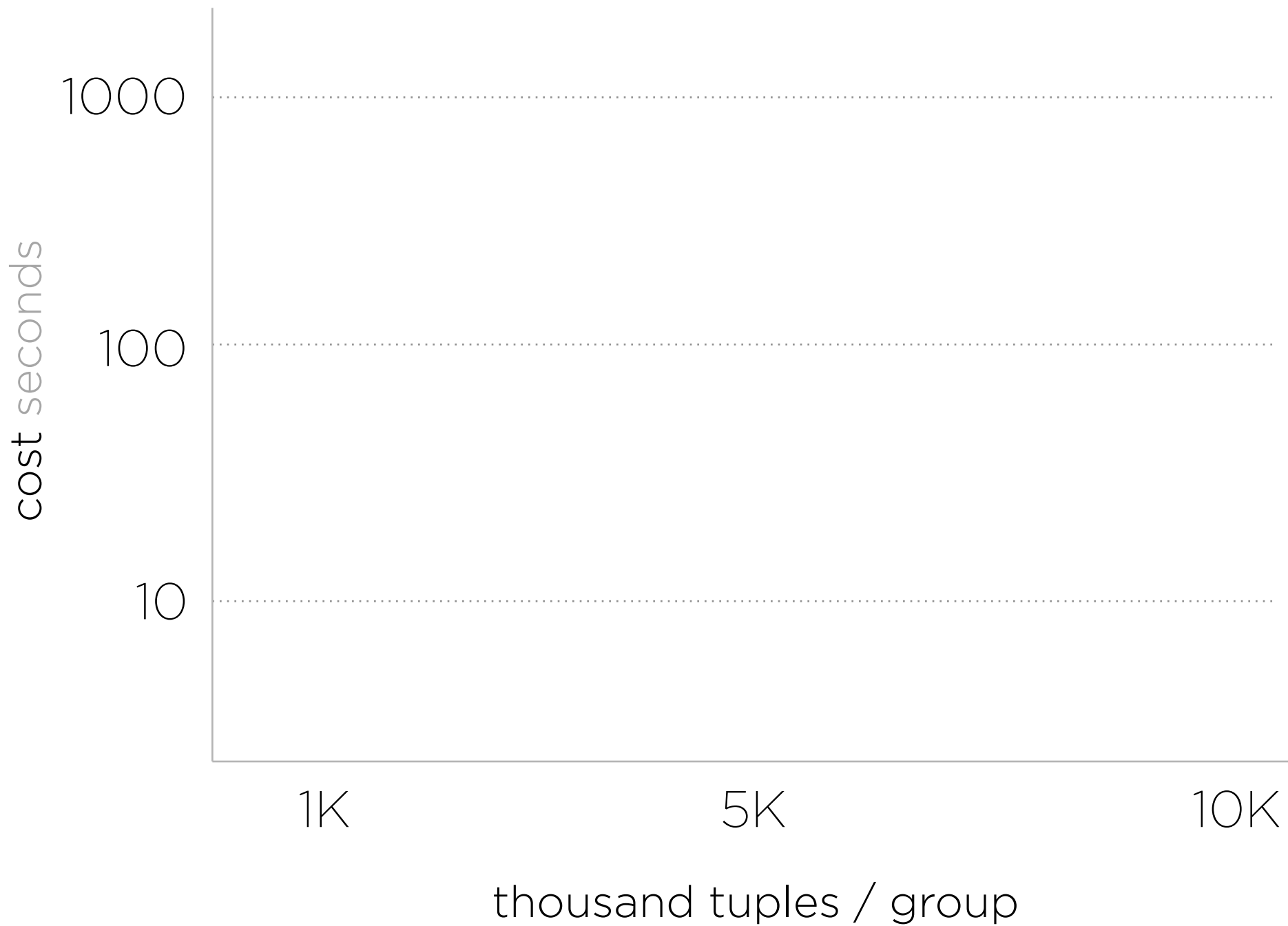


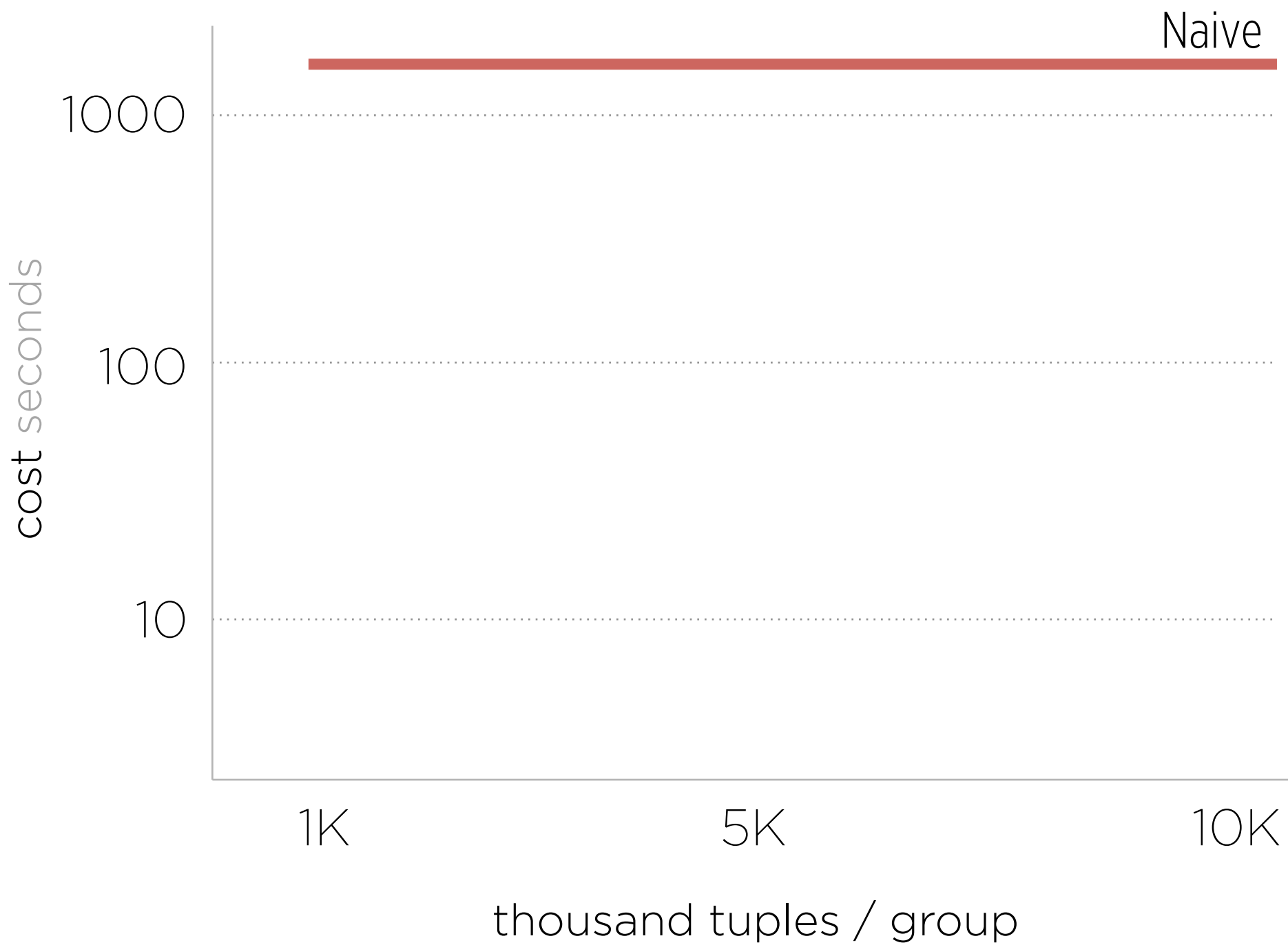
1K

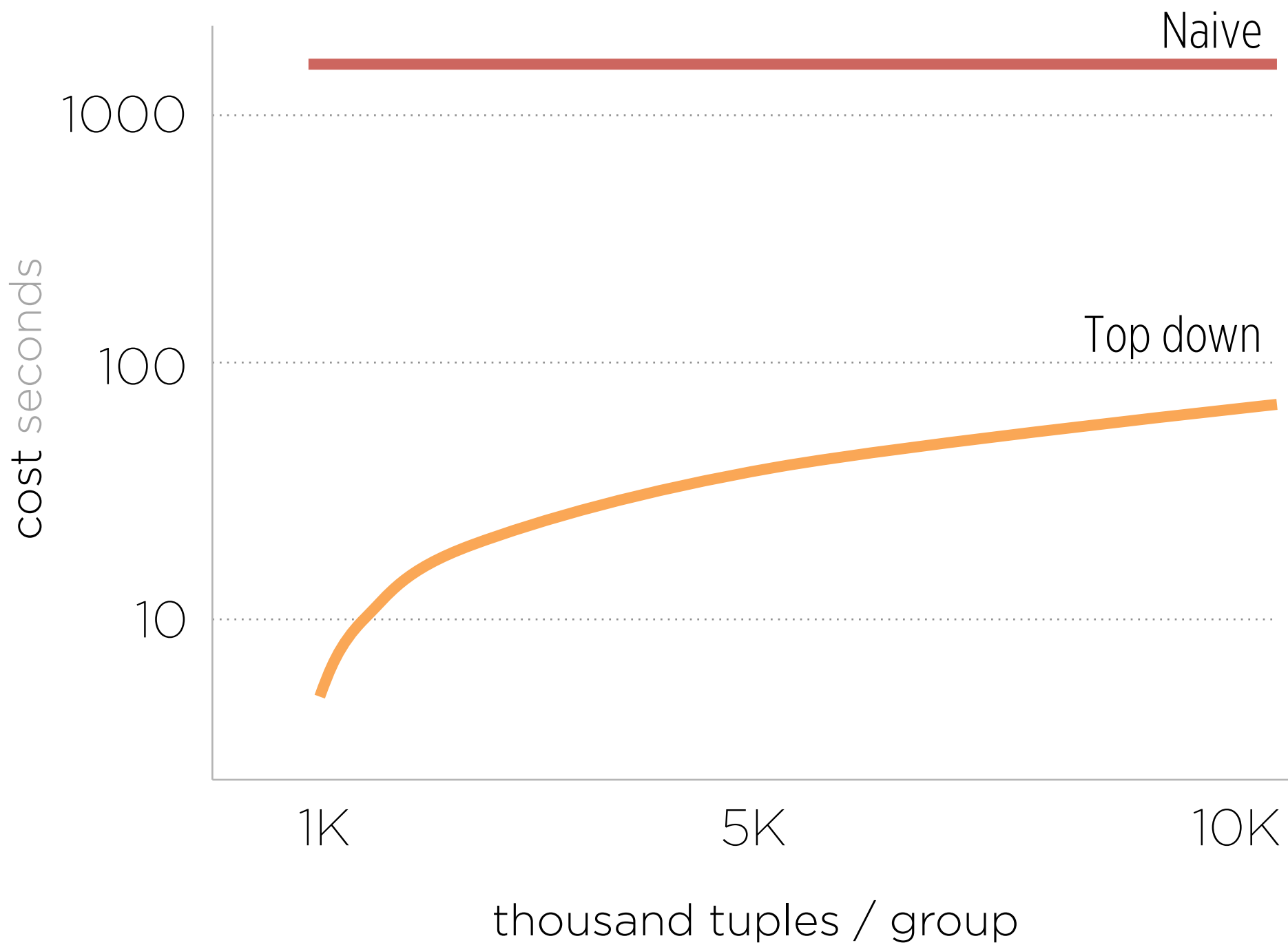
5K

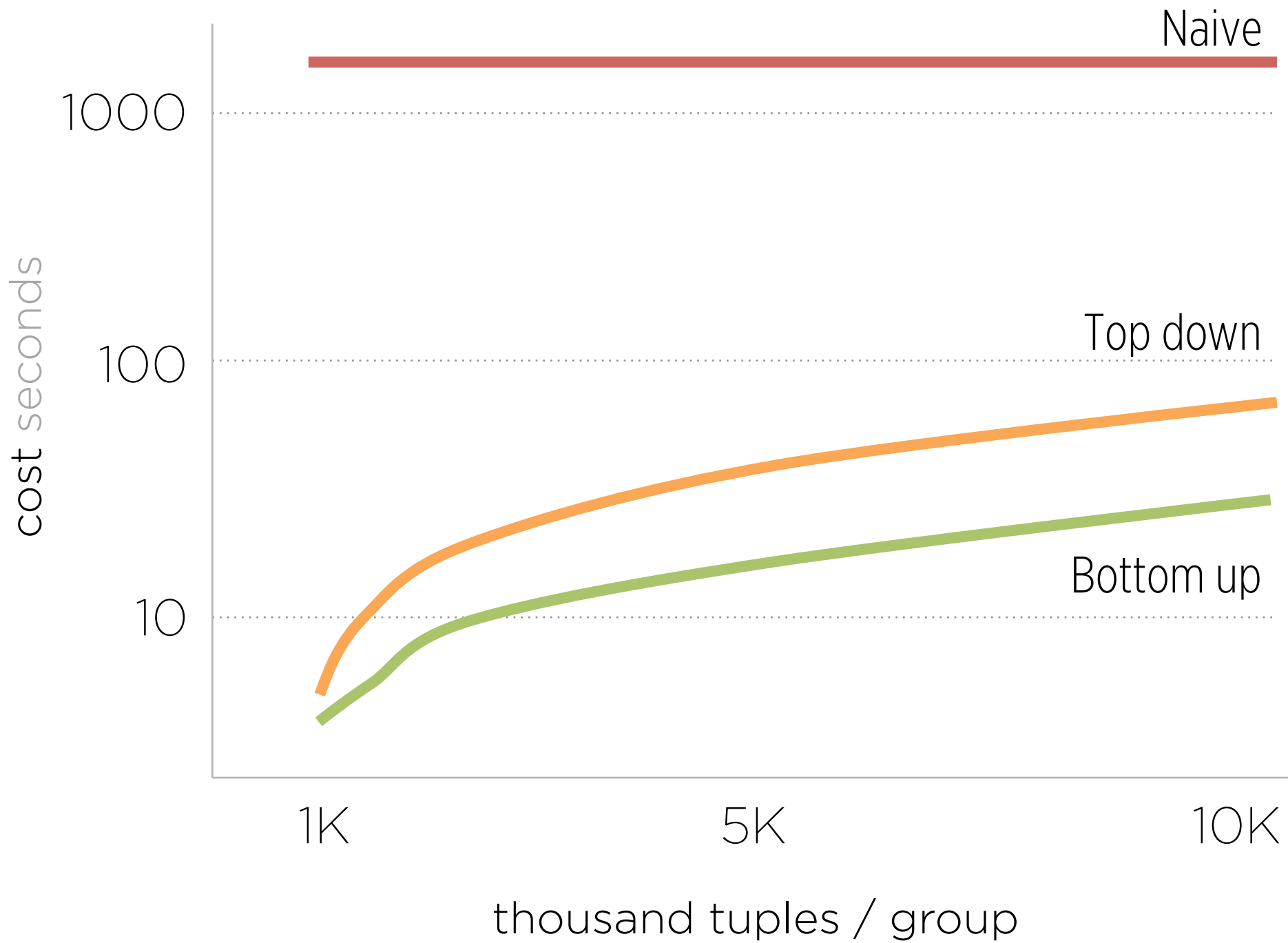
10K

thousand tuples / group









influence metric

that is

accessible to end-users

for

Data cleaning

Data exploration

Provenance reduction

scorpion

eugenewu@mit.edu



scorpion

eugnewu@mit.edu

GET OVER HERE!



scorpion

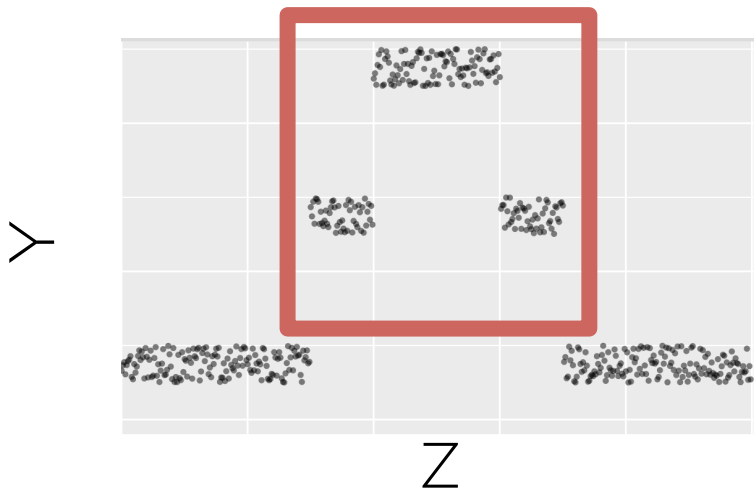
eugenewu@mit.edu



C-parameter

$$\frac{\Delta \text{output} \cdot V}{|p(T)|^c}$$

Low C



High C

