



Predictive Analytics War Stories

Feb 13, 2015

Hobson Lane

slides.com/hobsonlane/data-analytics-war-stories/live

bit.ly/pawsvote
Choose Your Story

7707-2-TOTAL
(770) 728-6825

1. Only Nyquist Knows
2. The Meaning of Mean
3. Data Dearth
4. Question the Question
5. Deep Net Runs Aground
6. Escape the Maze

When your vehicle is out of control...

1. Only Nyquist Knows

Photo by
Eric Cutright
Public Domain



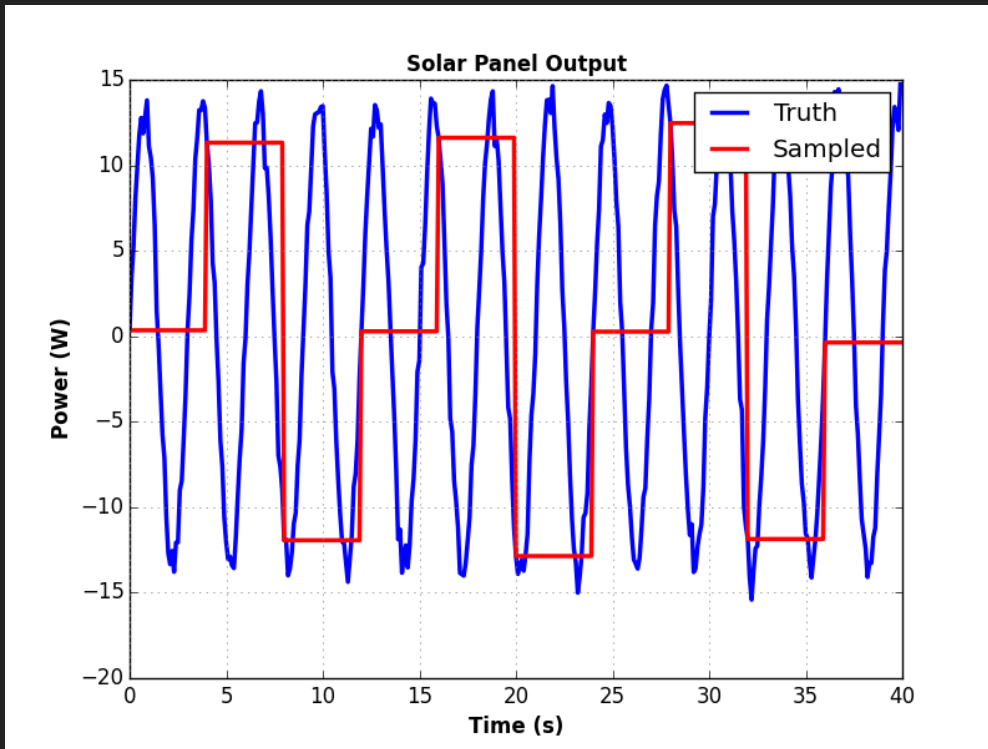
Photo by
Public Domain



Photo by
US Secret
Service

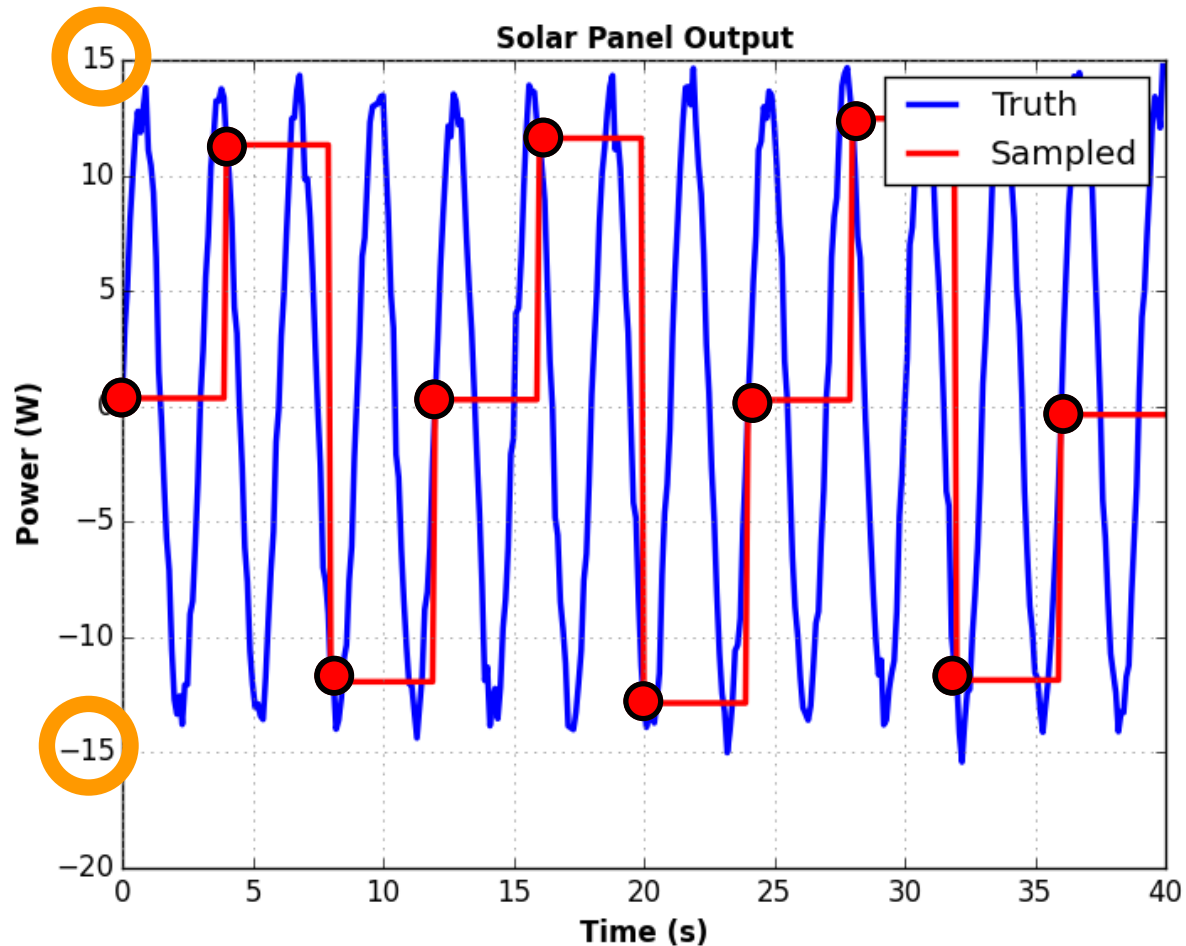
1. Only Nyquist Knows

- Nav sensors (gyro., accel) are "pegged"
- All you know is solar power:

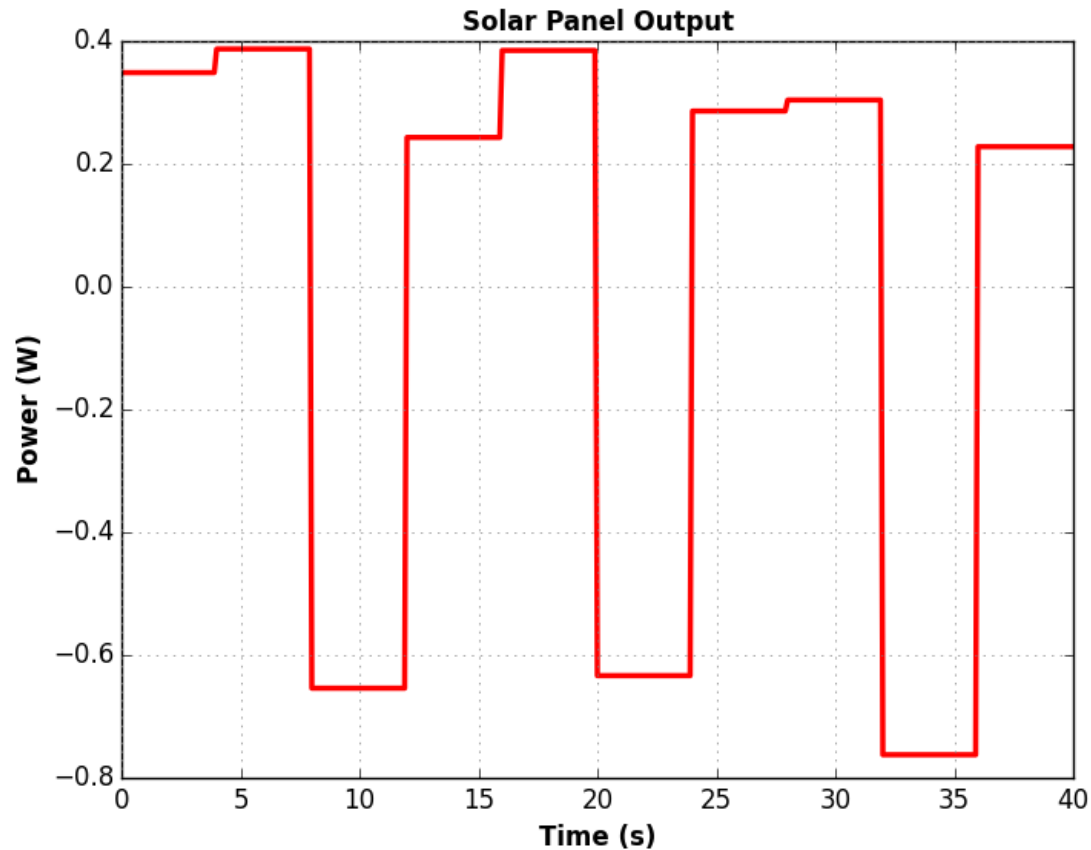


**How fast is
the tumble?**
4 sec !

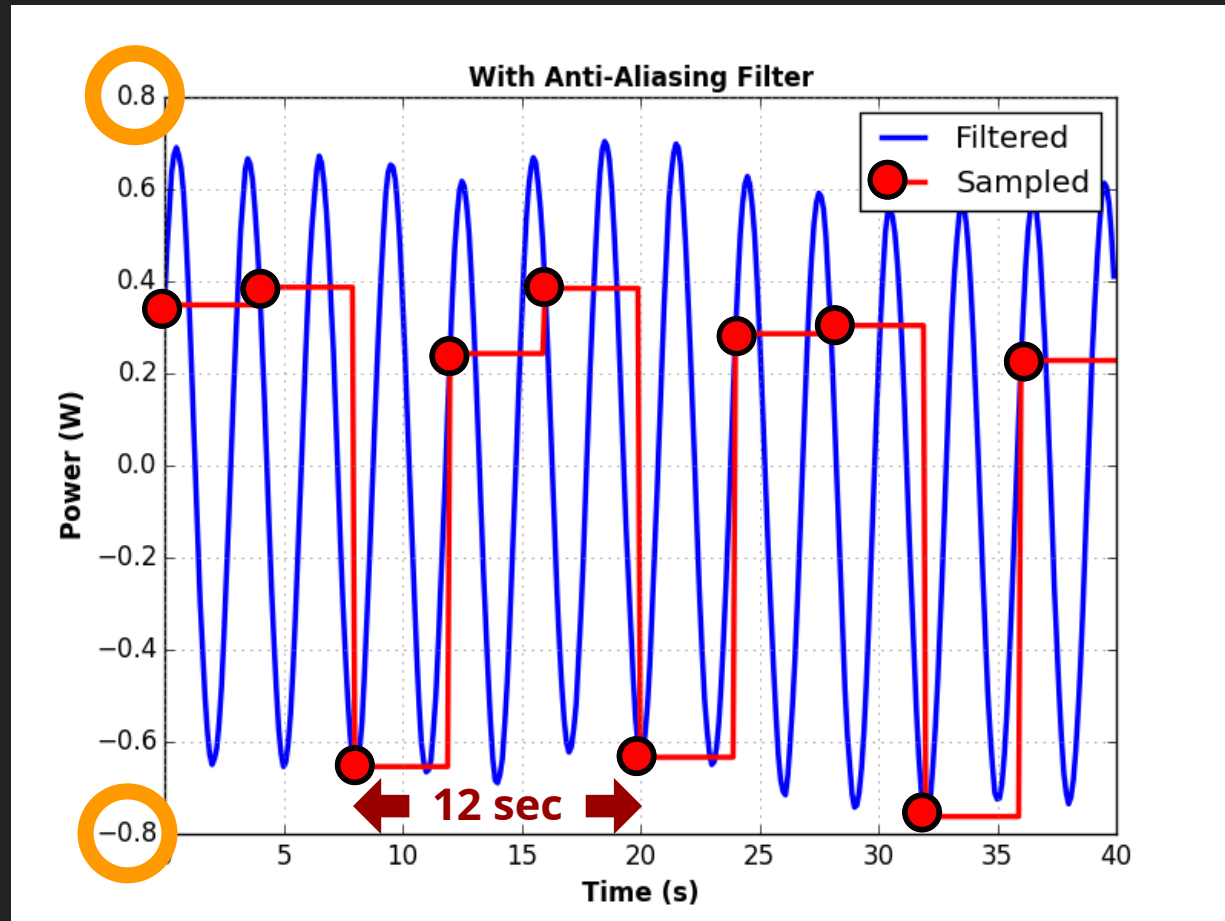
1. Only Nyquist Knows



Try an Anti-Aliasing Filter



Fail: Only Nyquist Knows



Workarounds

If Nyquist sampling (2x faster than truth) isn't possible....

- Use a different sensor
 - Postprocess existing signal (radio doppler)
- Sample **irregularly!**
 - Captures higher frequencies
 - Lomb-Scargle to post-process

```
spectrum = scipy.signal.lombscargle(sample_times, samples, frequencies)
```

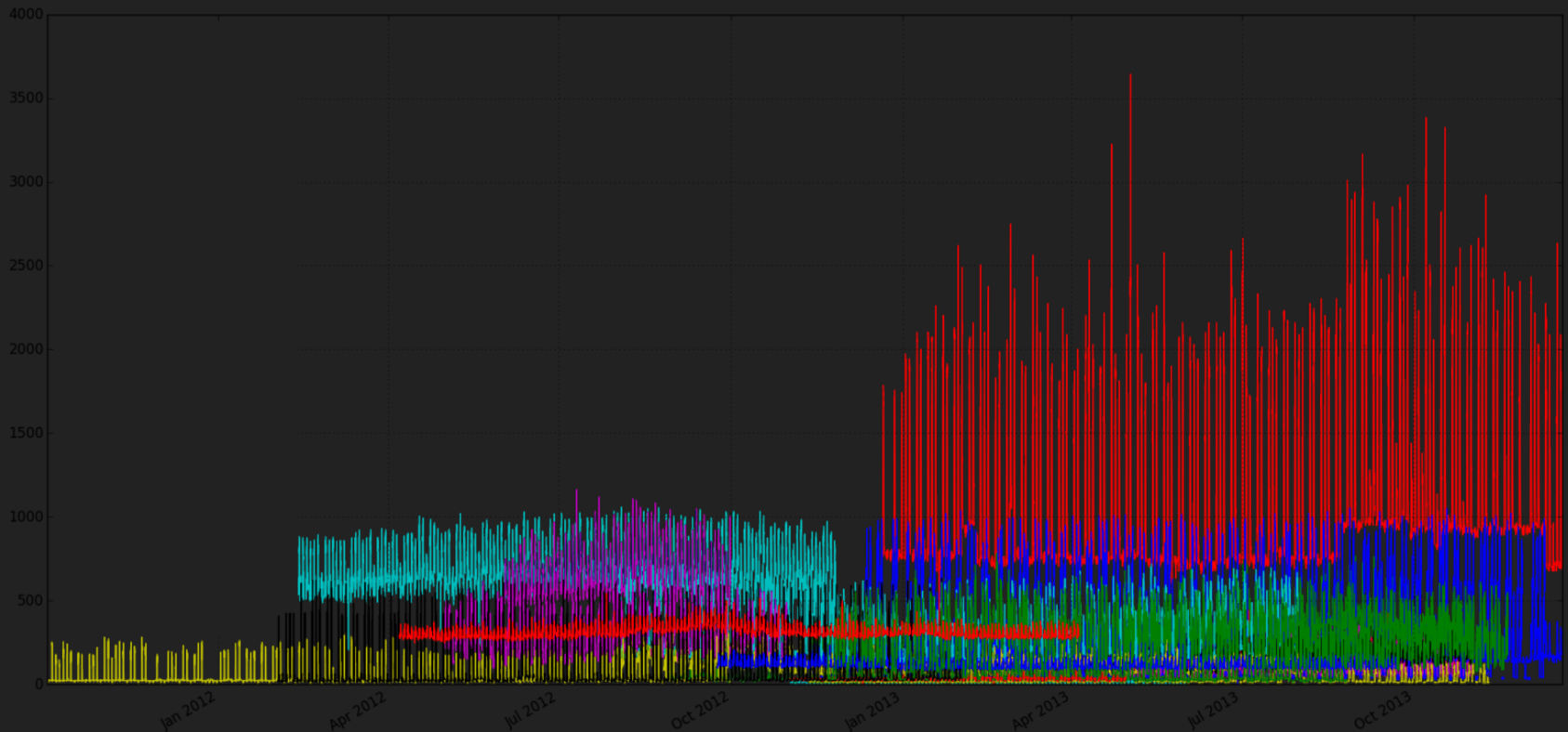
- Probabilistic modeling
 - Great for overwhelming data volume (**IoT**)

2. The Meaning of Mean

- Means don't tell the whole story
- Consider both μ and σ
- Meaning may be found in the means for each...
 - **group**, cluster, or class
- For us we started with grouping by **time of day**, but that wasn't enough...

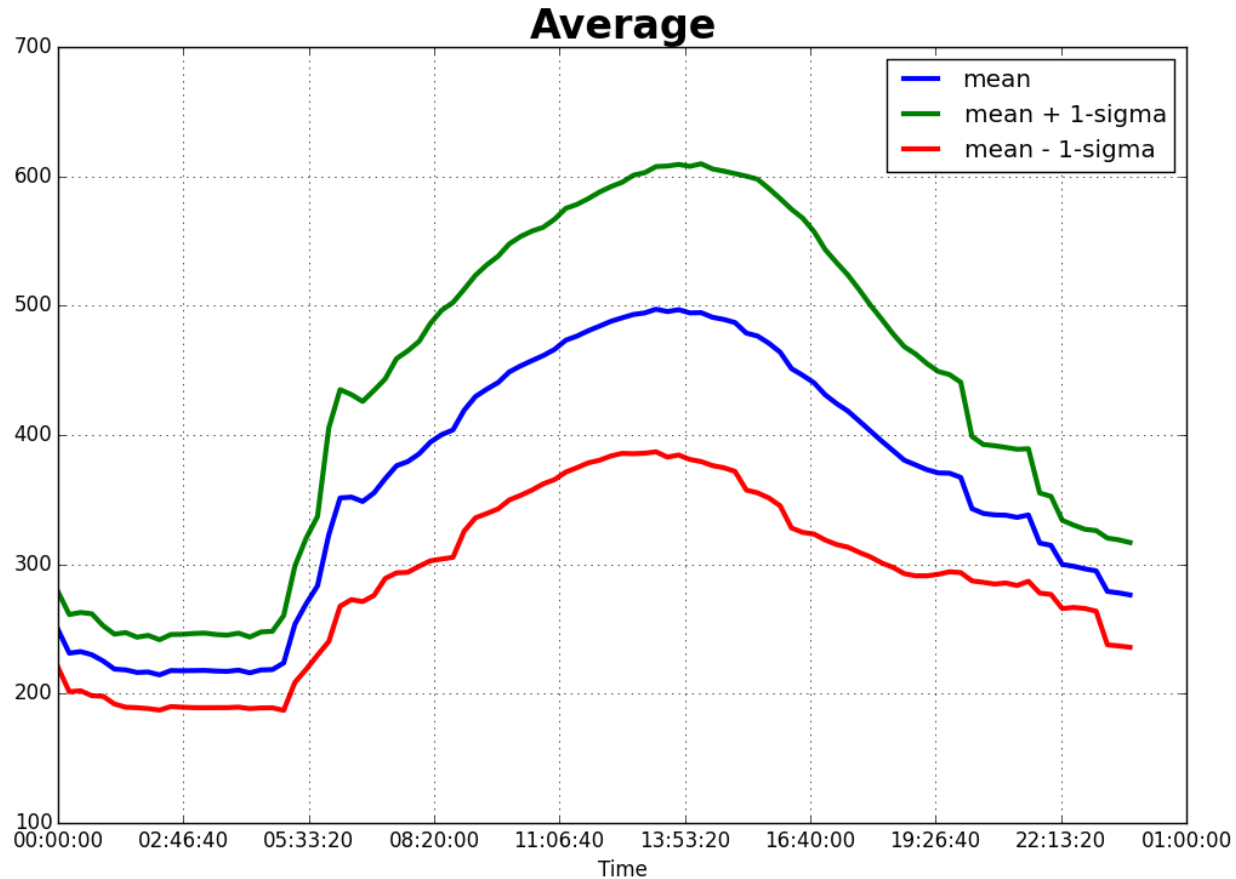
2. The Meaning of Mean

- Regression and classification required
- Many "fundamental frequencies"

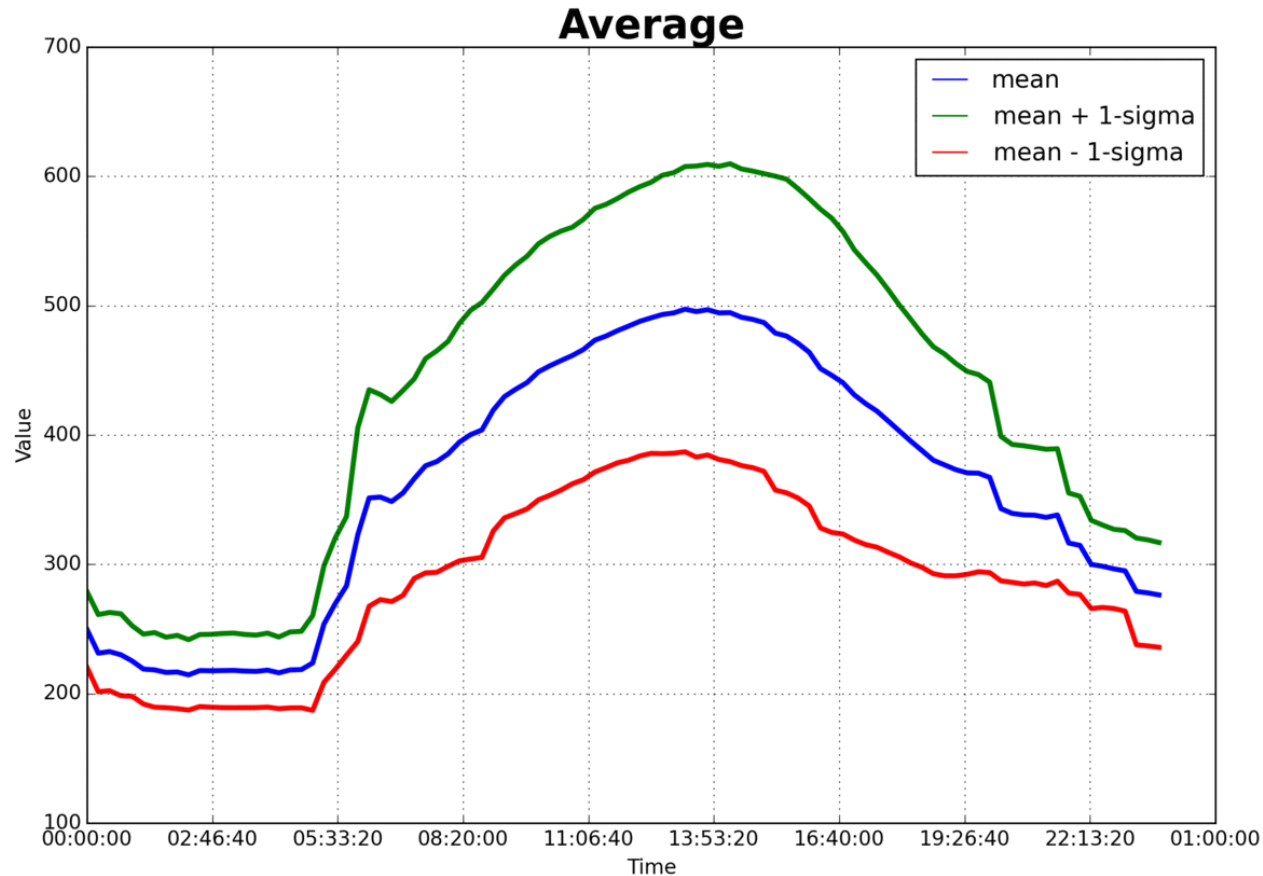


SMS: 7707-2-TOTAL or (770) 728-6825 MSGS: "1", "2", "3", "4", "5", or "6"

Mean for Each Time of Day

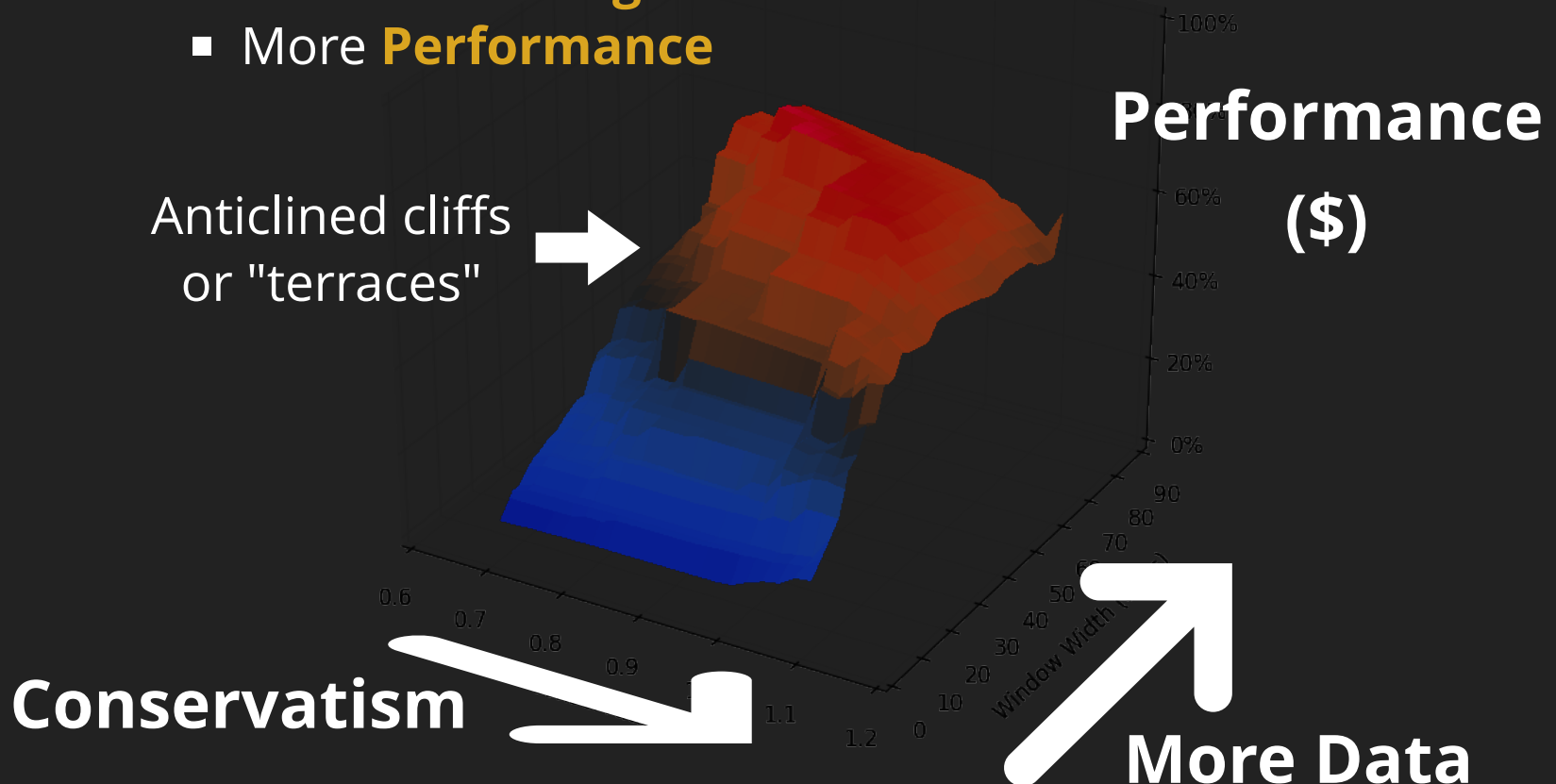


Classify Before Getting Mean



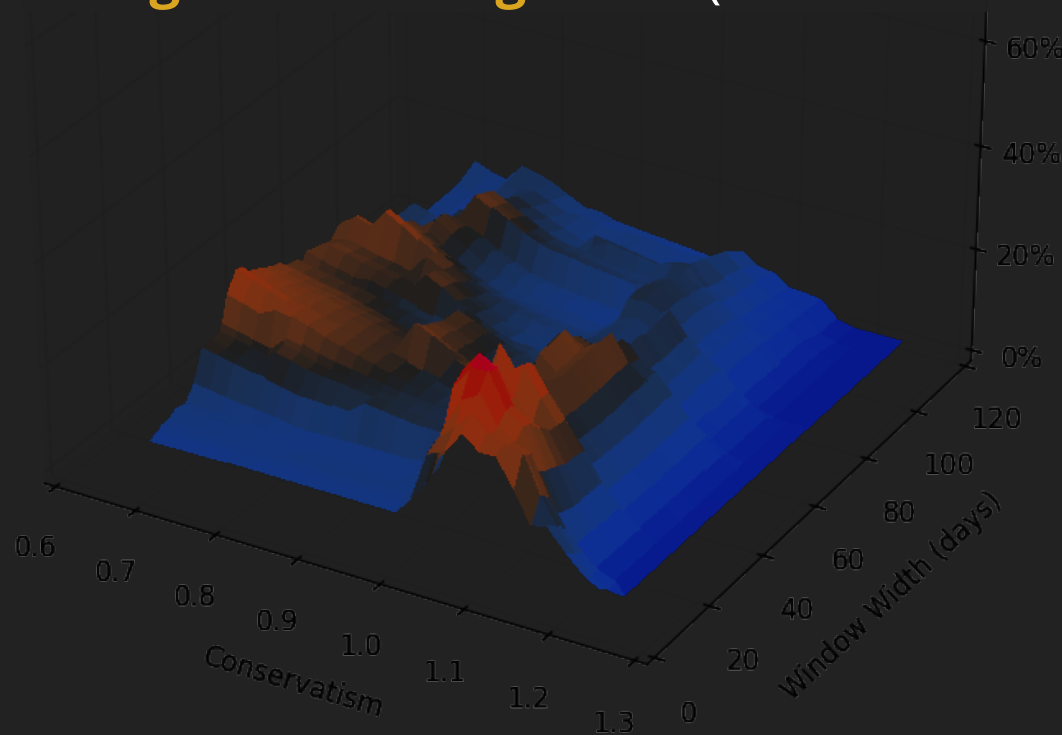
3. Data Dearth

- Tuning a 2-DOF predictive filter for performance
- More data gives algorithm more to work with
 - Less **Overfitting**
 - More **Performance**



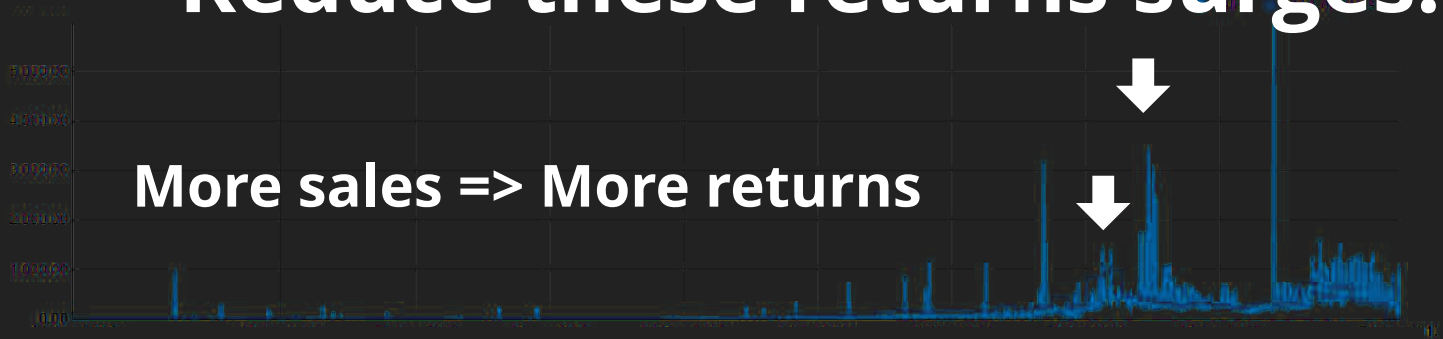
3. Data Dearth

- Sometimes **more of the same** doesn't help
 - Exogenous factors confound the smartest algorithm
- Make the **exogenous endogenous** (new data source)



4. Question the Question

Reduce these returns surges!



Correlation != Causation

(a. la. Tyler Vigen)

Multiple interacting causes

**Normalize return rate for sales
(lag-compensated)**

4. Question the **6σ** Question

"Cost of quality"

"Customer reject rate"

"Defect rate"

$$\text{Reject rate} = \frac{\text{Rejects (last quarter)}}{\text{Sales (last quarter)}}$$

Simple equation everyone can agree on

But it's **Wrong!**

And it's **Late!**

4. Better "Question"

$$\text{Reject rate} = \frac{\text{Rejects (last quarter)}}{\text{Sales (qtr before last)}}$$

Even Better

$$\text{Reject rate} = \frac{\text{Rejects (last quarter)}}{\text{Sales (estimate lagged quarter)}}$$

Correct

$$\text{Reject rate} = \frac{\text{Rejects (last week)}}{\text{Sales (integral of lagged sales)}}$$

$$r_r = \sum_k \alpha s_{n-k}$$

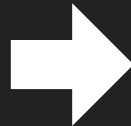
"Birth-Death Process"

Sale

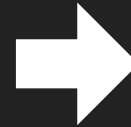
Lag

Reject

$S(t)$



$H(t, \tau)$



$R(t)$

Product enters
"pipeline" arbitrarily

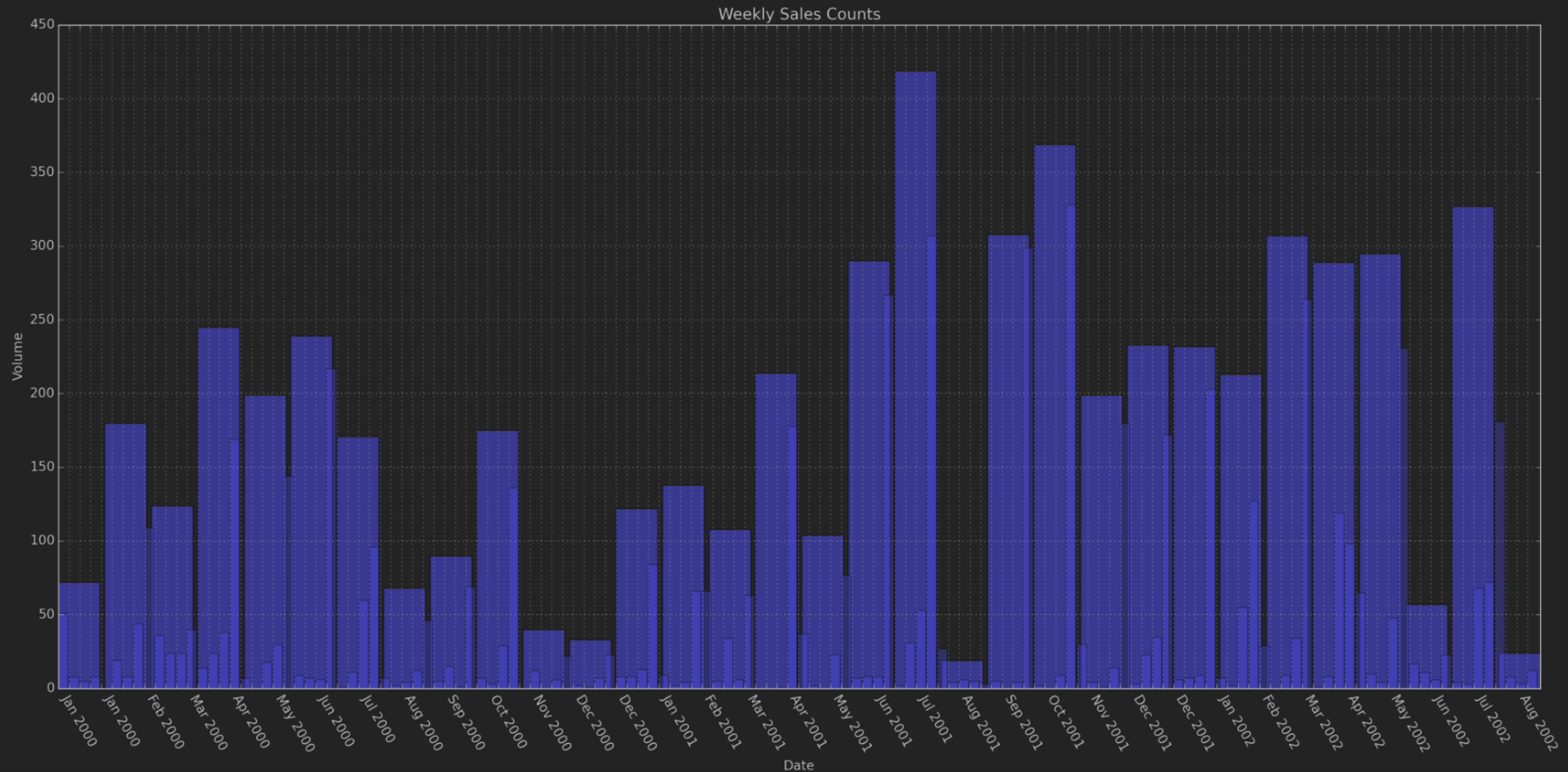
Flow rate
(Reject rate)

All products "die",
Question is **when**
And the portion that
happens too soon

$$r_r = \sum_k \alpha s_{n-k}$$

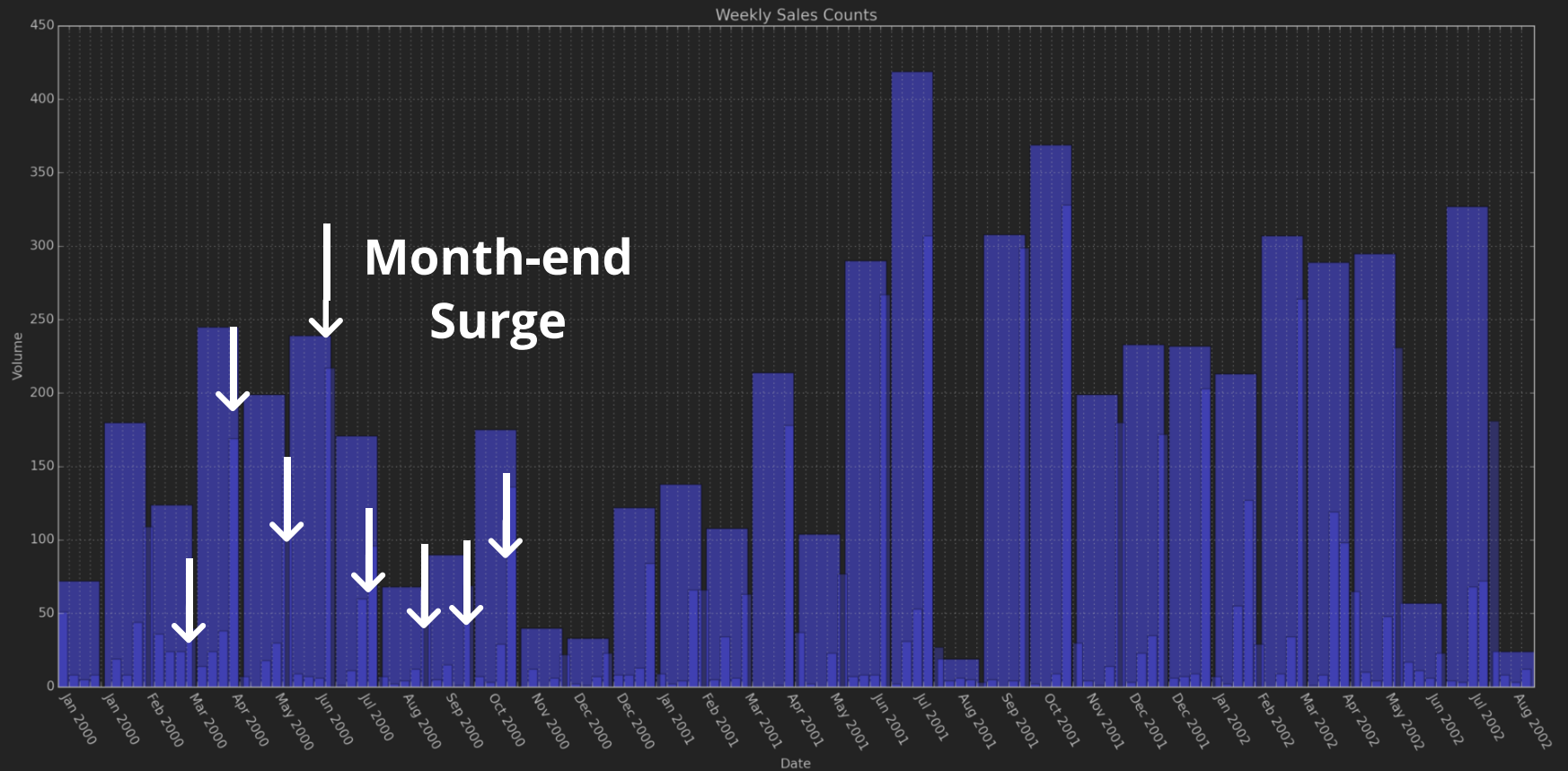
4. Question the Question

Histogram reveals trend and seasonality



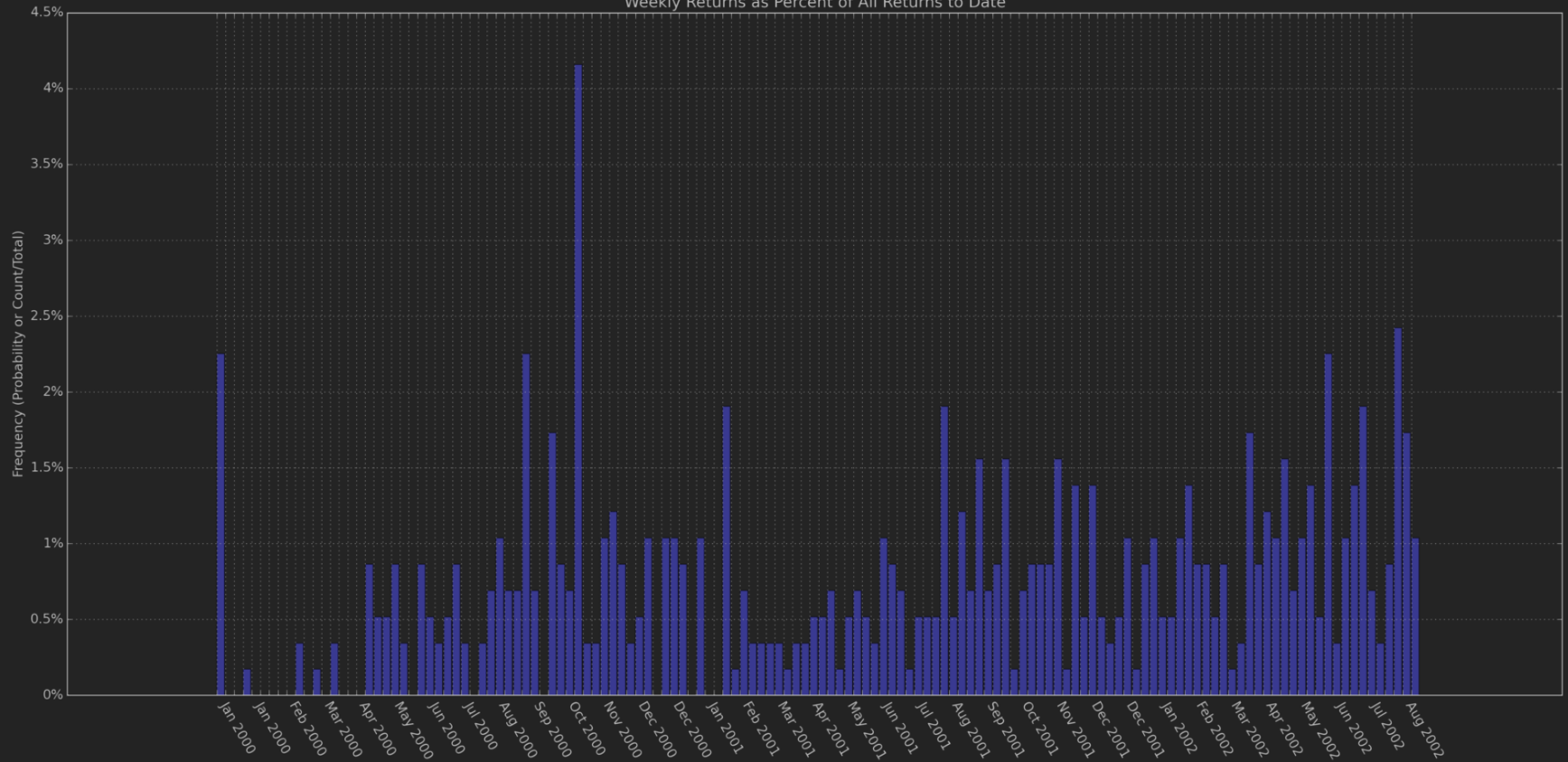
SMS: 7707-2-TOTAL or (770) 728-6825 MSGS: "1", "2", "3", "4", "5", or "6"

Sales



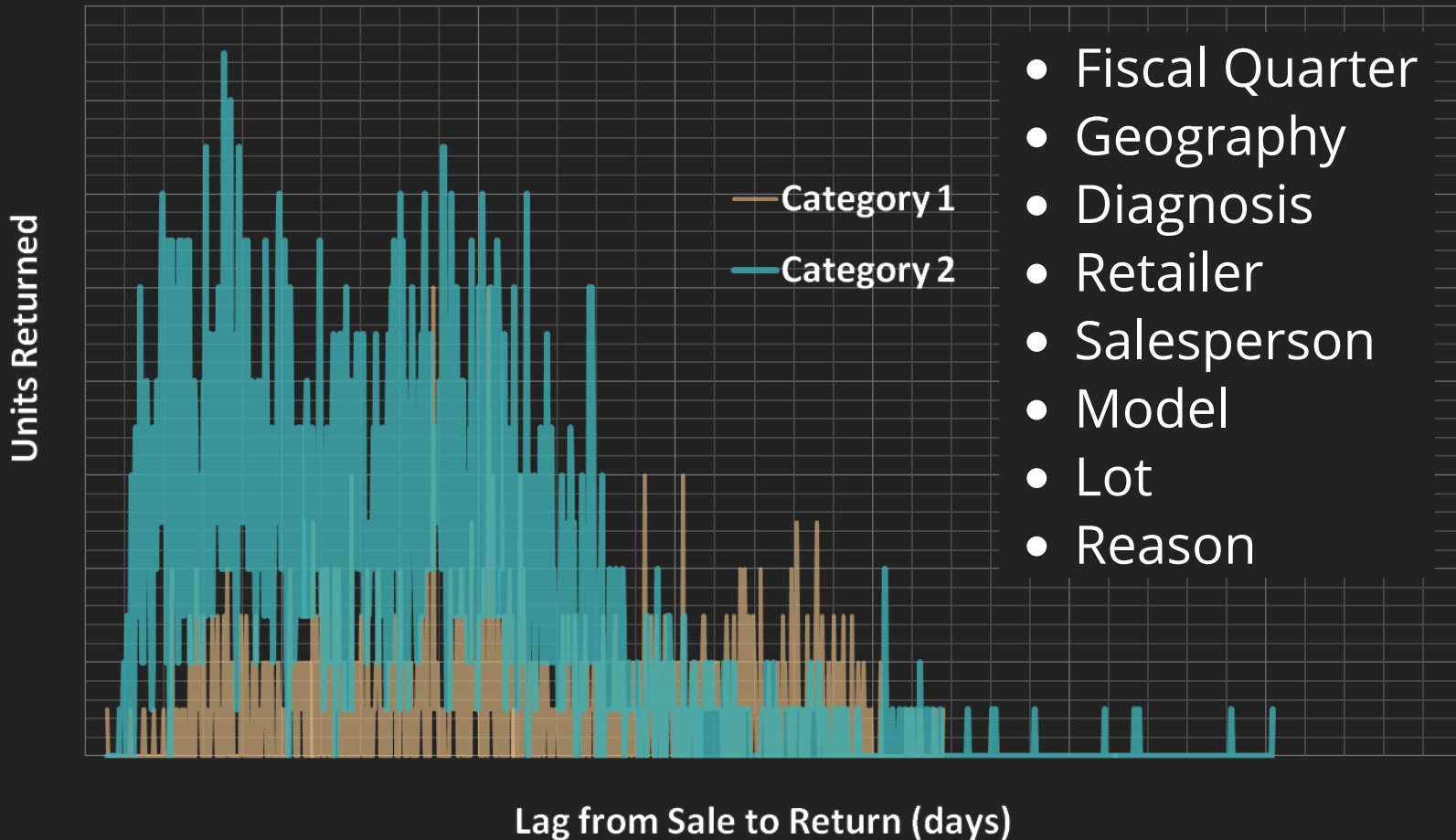
Rejects

Weekly Returns as Percent of All Returns to Date



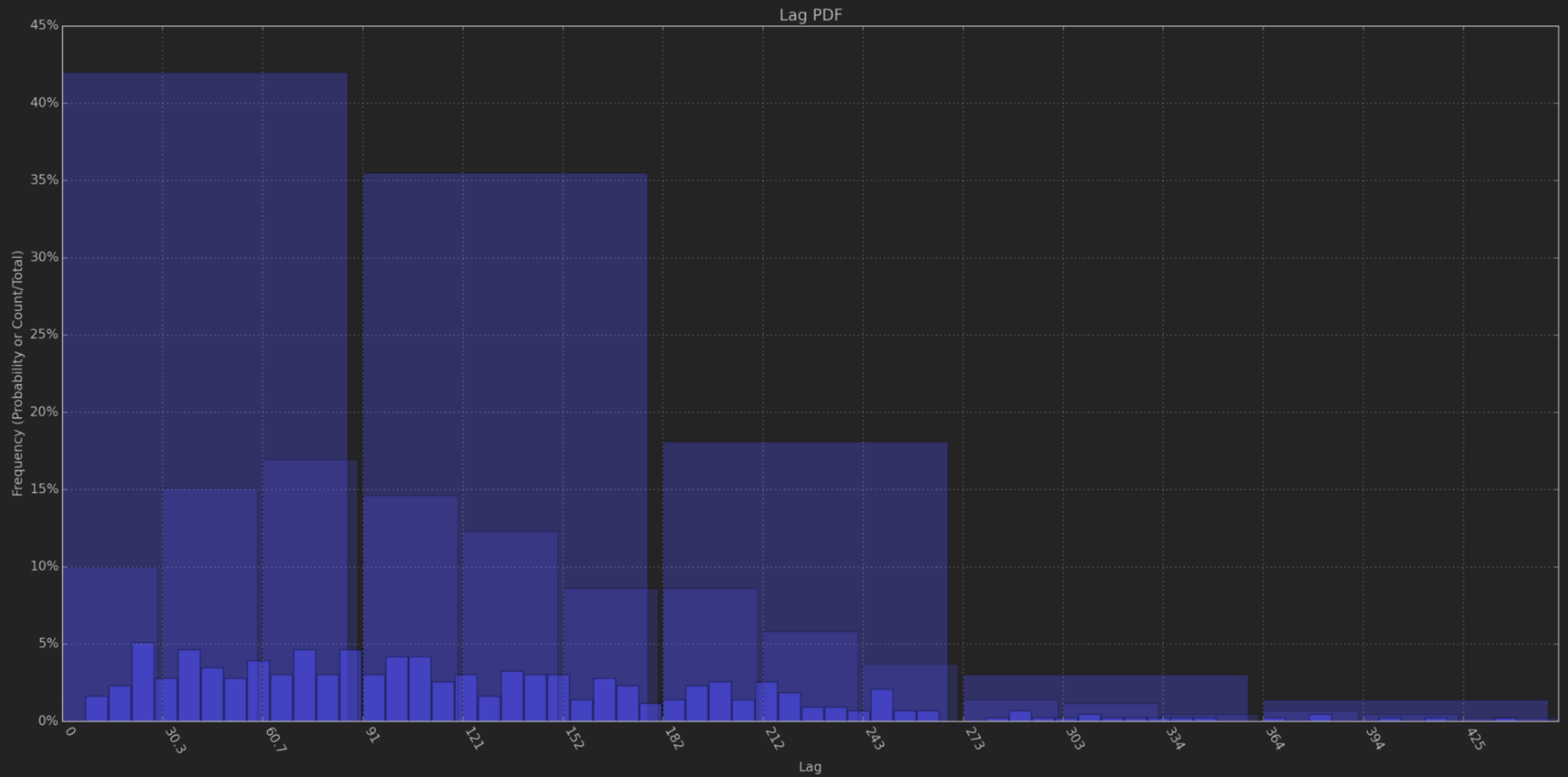
4. Question the Question

Returns Lag Histogram

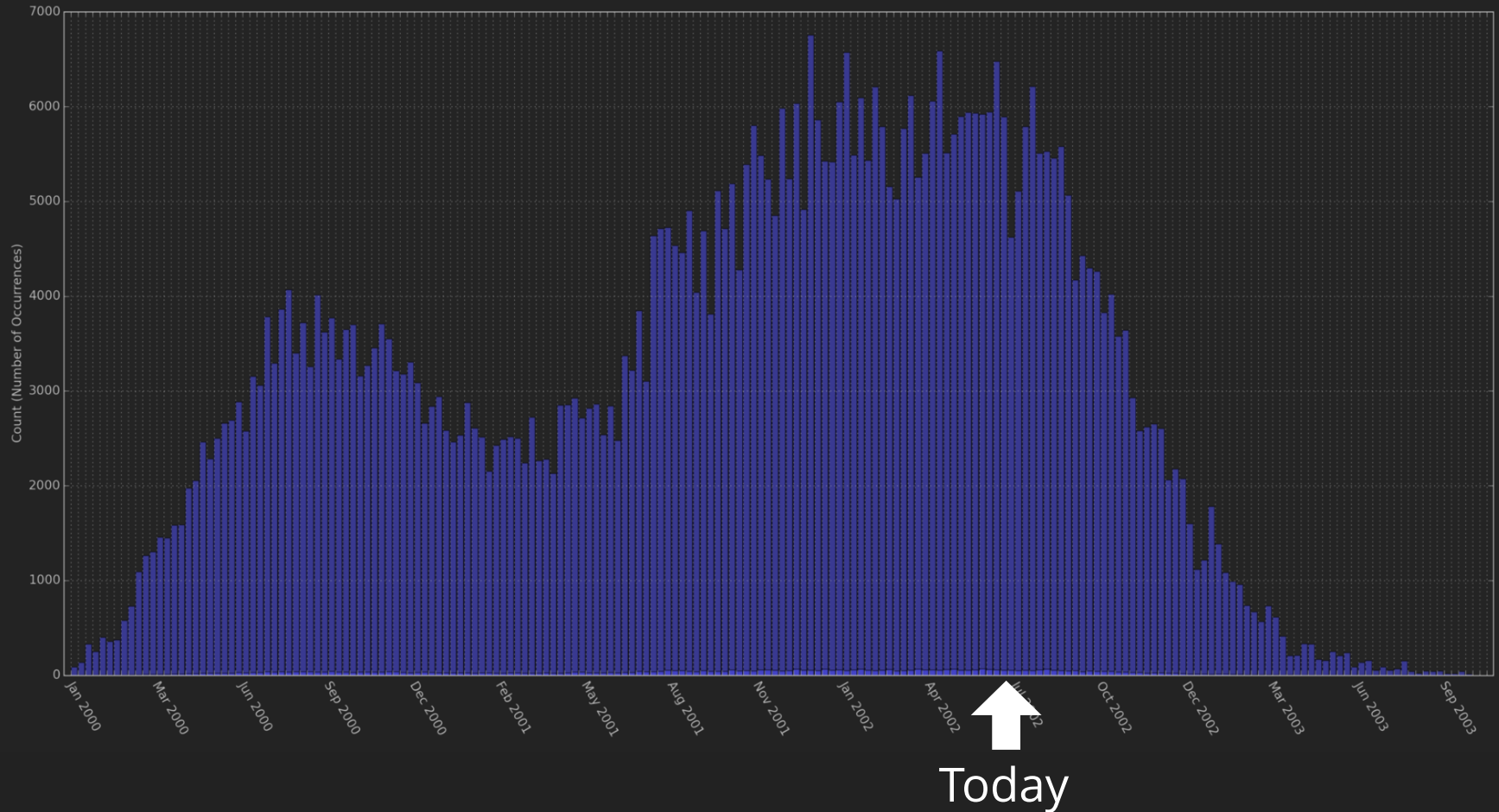


SMS: 7707-2-TOTAL or (770) 728-6825 MSGS: "1", "2", "3", "**4**", "5", or "6"

Lag



Lagged Sales



Predicted Returns

Lag

Sales

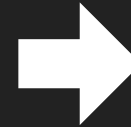
Process

Rejects

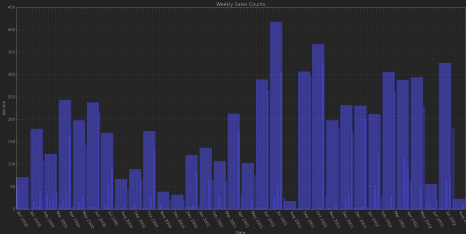
$S(t)$



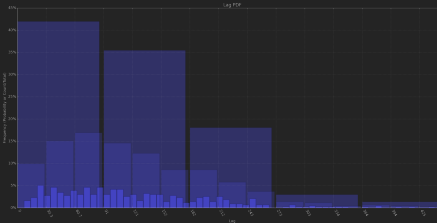
$H(t, \tau)$



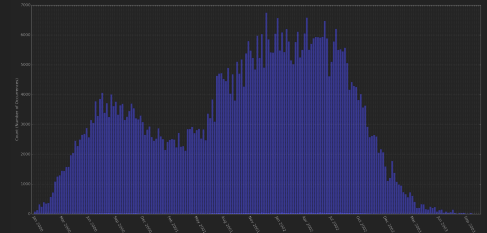
$R(t)$



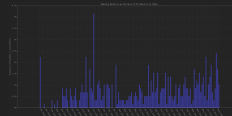
*



=

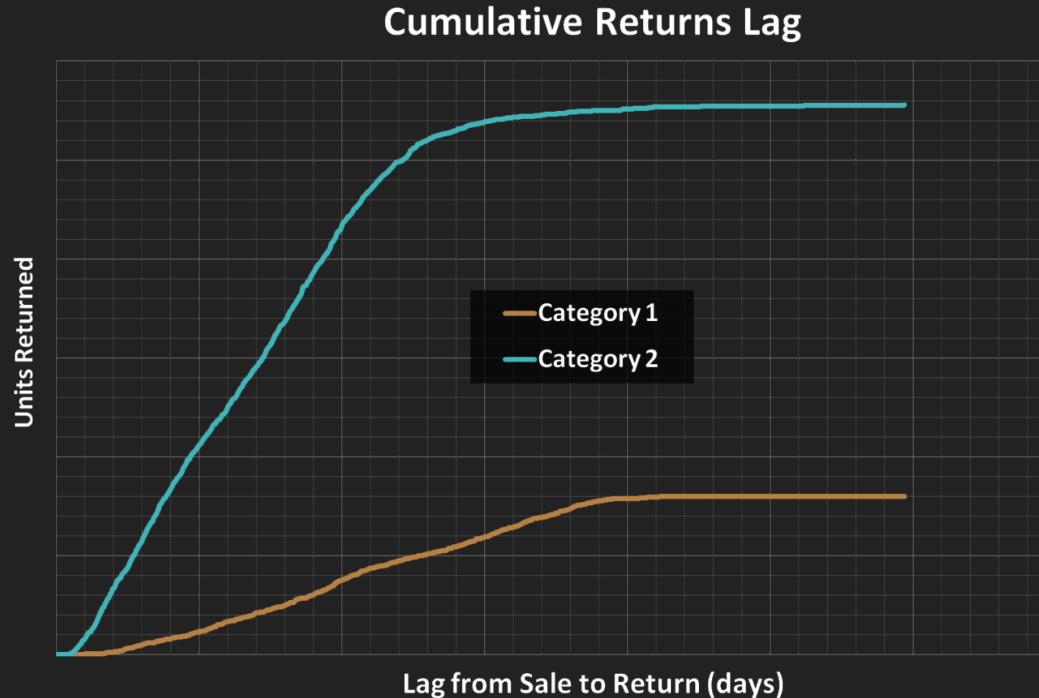


÷



4. Analyze the Question

Cumulative histograms focus attention on final total

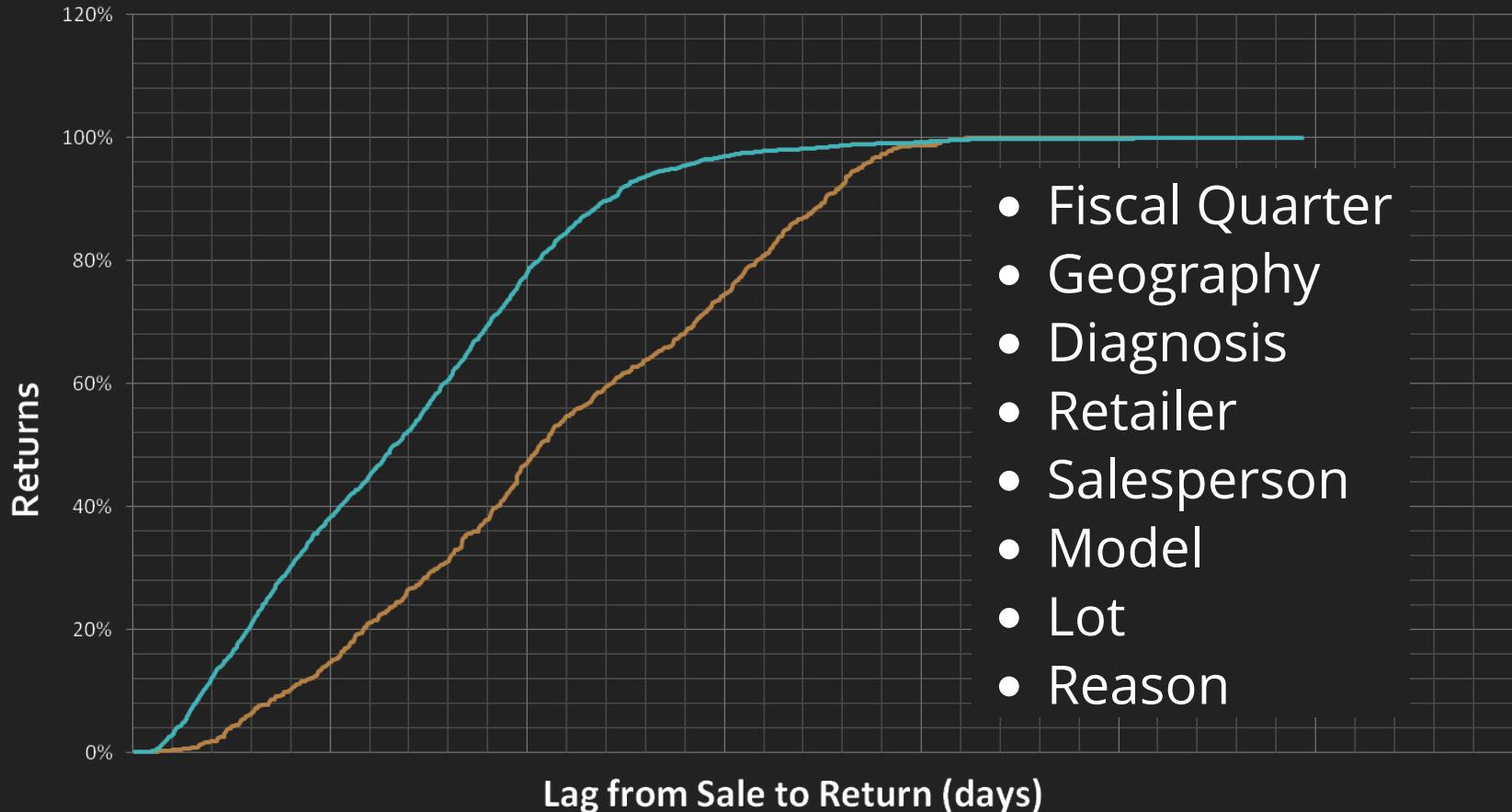


Product returns stop when...

- You stop accepting returns
- You stop counting
- You stop selling

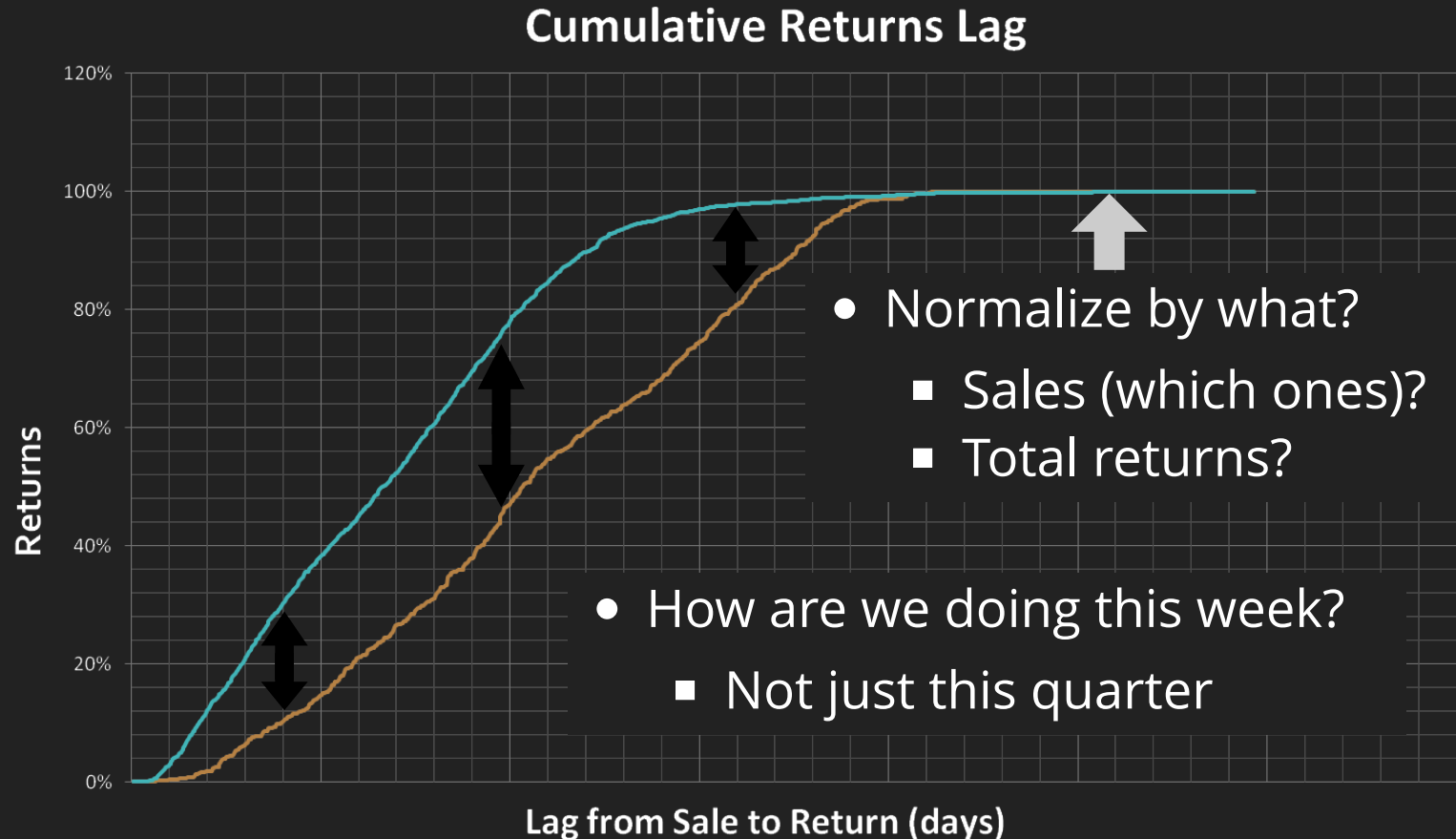
4. Normalize & Compare

Cumulative Returns Lag



4. Analyze the Question

Normalize histograms to compare categories



4. Question the Question

Unsupervised natural language processing?

President inaugural speeches

Target category = political party

4. Question the Question

What are the US Presidents' political parties based on speeches?

4. Question the Question

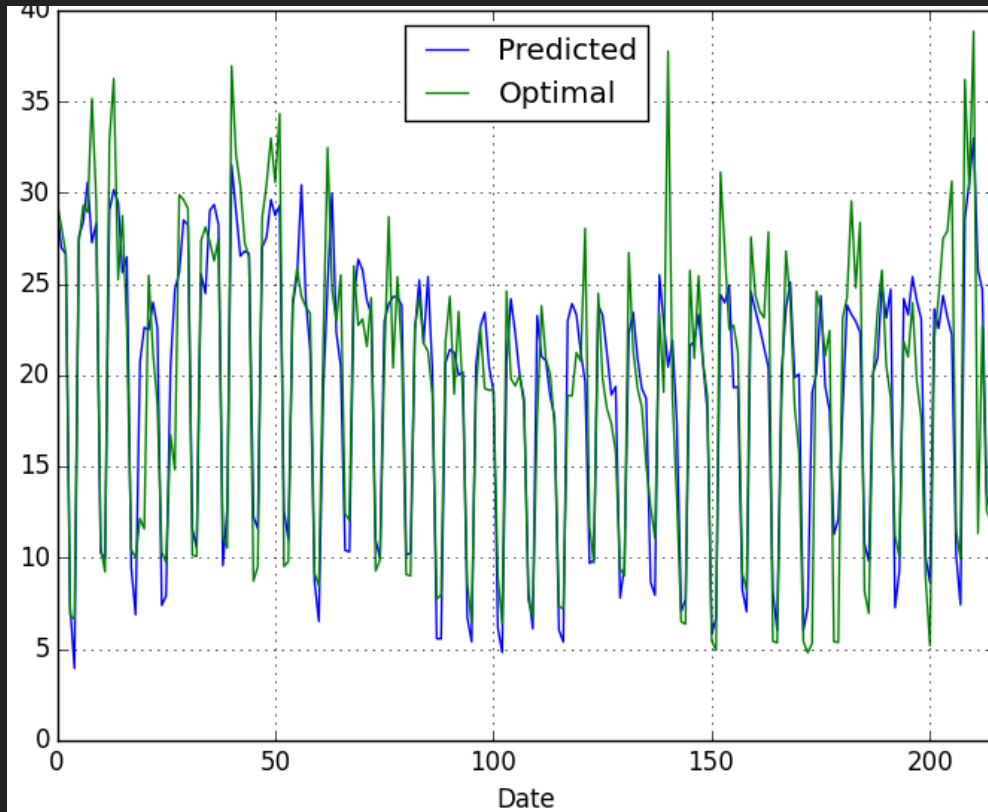
What are the US Presidents' political parties based on speeches?

4. Question the Question

- The category you're interested in will not likely be the most important "factor" in the NLP statistics
- Dimension reduction (SVD, PCA) can identify factors
 - Word-sets that are most significant
- These represent the "themes"
 - Interpretation of these "themes" is up to you
 - Statistics \neq Meaning

5. Deep Nets Run Aground

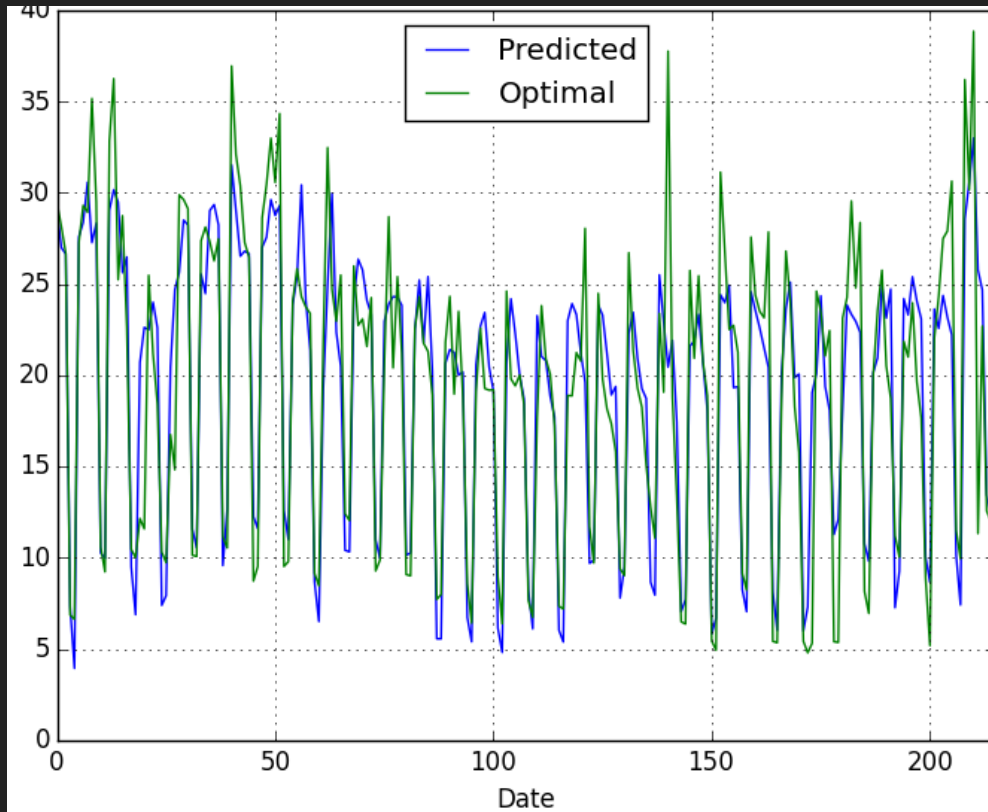
Deep net performs well!



SMS: 7707-2-TOTAL or (770) 728-6825 MSGS: "1", "2", "3", "4", "5", "6"

5. Deep Nets Run Aground

Not so fast... it's **overfitting**



SMS: 7707-2-TOTAL or (770) 728-6825 MSGS: "1", "2", "3", "4", "5", "6"

5. Deep Nets Run Aground

$$\mathbf{p} \rightarrow \mathbf{W}_{S^k, S^{(k+1)}}^k \rightarrow \mathbf{a}$$

$$\mathbf{a} = \mathbf{W}_{S^k, S^{(k+1)}}^k \mathbf{p}$$

$$\mathbf{W}_{n,m}^k = \begin{bmatrix} w_{1,1} & w_{1,2} & \dots & w_{1,m} \\ w_{2,1} & w_{2,2} & \dots & w_{2,m} \\ \vdots & \vdots & & \vdots \\ w_{n,1} & w_{n,2} & \dots & w_{n,m} \end{bmatrix}$$

- Conventional Hebb rule

$$\mathbf{W}^{new} = \mathbf{W}^{old} + \mathbf{t}_q \mathbf{p}_q^T$$

- Hebb "delta" rule

$$\mathbf{W}^{new} = \mathbf{W}^{old} + \alpha (\mathbf{t}_q - \mathbf{a}_q) \mathbf{p}_q^T$$

5. Shallow Data

$$\mathbf{p} \rightarrow \mathbf{W}_{S^k, S^{(k+1)}}^k \rightarrow \mathbf{a}$$

$$\mathbf{a} = \mathbf{W}_{S^k, S^{(k+1)}}^k \mathbf{p}$$

$$\mathbf{W}_{n,m}^k = \begin{bmatrix} w_{1,1} & w_{1,2} & \dots & w_{1,m} \\ w_{2,1} & w_{2,2} & \dots & w_{2,m} \\ \vdots & \vdots & & \vdots \\ w_{n,1} & w_{n,2} & \dots & w_{n,m} \end{bmatrix}$$

- Model degree:

$$\sum_k S^k S^{(k+1)}$$

- Training data DOF:

$$S^1 S^3 N_{samples}$$

(independent samples)

5. Shallow Data

$$\mathbf{p} \rightarrow \mathbf{W}_{S^k, S^{(k+1)}}^k \rightarrow \mathbf{a}$$

$$\mathbf{a} = \mathbf{W}_{S^k, S^{(k+1)}}^k \mathbf{p}$$

$$\mathbf{W}_{n,m}^k = \begin{bmatrix} w_{1,1} & w_{1,2} & \dots & w_{1,m} \\ w_{2,1} & w_{2,2} & \dots & w_{2,m} \\ \vdots & \vdots & & \vdots \\ w_{n,1} & w_{n,2} & \dots & w_{n,m} \end{bmatrix}$$

- Model degree:

$$S^1 S^2 + S^2 S^3$$

(1 hidden layer)

- Training data DOF:

$$(S^1 + S^3) N_{samples} \text{ (independent samples)}$$

5. Bottom Line

$$N_{hidden} \ll N_{training}$$

bit.ly/nntune

6. Escape the Maze

Find Connections

(Actionable Insight)



18 databases

> 10k tables

> 100k fields

> 10M records/table

6. Escape from the Maze

- Tight heuristics vital for efficient graph search
- "Always turn right" is not good enough



6. Escape from the Maze

- Don't bother with "exhaustive" correlation search

$$\text{complexity} \approx O(\underset{10^7}{M}^2 \underset{10^5}{N}^2) \approx 10^{24}$$

- Find db relationships using meta-data
 - min, max, median
 - #records
 - #distinct
 - for reals: mean, std

$$\text{complexity} \approx O(MN \log(N)) \approx 10^{13}$$

Human Heuristics

- **Business knowledge narrows search:**
 - Repair technicians
 - Product designers
 - Factory managers
 - Suppliers
 - Sales channels
 - Call center

Accidental "Experiments"

- **Look for differences in**
 - Model
 - Lot
 - Product
 - Sales Channel
 - Customer Demographic
 - Region/Culture
- **Look for ...**
 - New/deleted features
 - Documentation updates
 - Cost-saving parts changes
 - Production facilities (outsourced vs insourced)

Kruskal's Algorithm

Minimum Spanning Tree

1. Add lowest cost edge with new node
2. Repeat until all nodes accounted for

Produces one graph for each connected subgraph

Built into python graph library (`networkx`):

```
def minimum_spanning_zipcodes():  
    zipcode_query_sequence = []  
    G = build_graph(api.db, limit=1000000)  
    for CG in nx.connected_component_subgraphs(G):  
        for edge in nx.minimum_spanning_edges(CG):  
            zipcode_query_sequence += [edge[2]['zipcode']]  
    return zipcode_query_sequence
```

A* Algorithm

Minimum Path to Goal

Provably optimal and optimally efficient

But typical data relationship graph has large branching factor

```
from networkx.algorithms.shortest_paths import astar_path  
astar_path(G, source, target, heuristic=None)
```

Built into python graph library (`networkx`)

A* Algorithm

Minimum Path to Goal

Provably optimal and optimally efficient

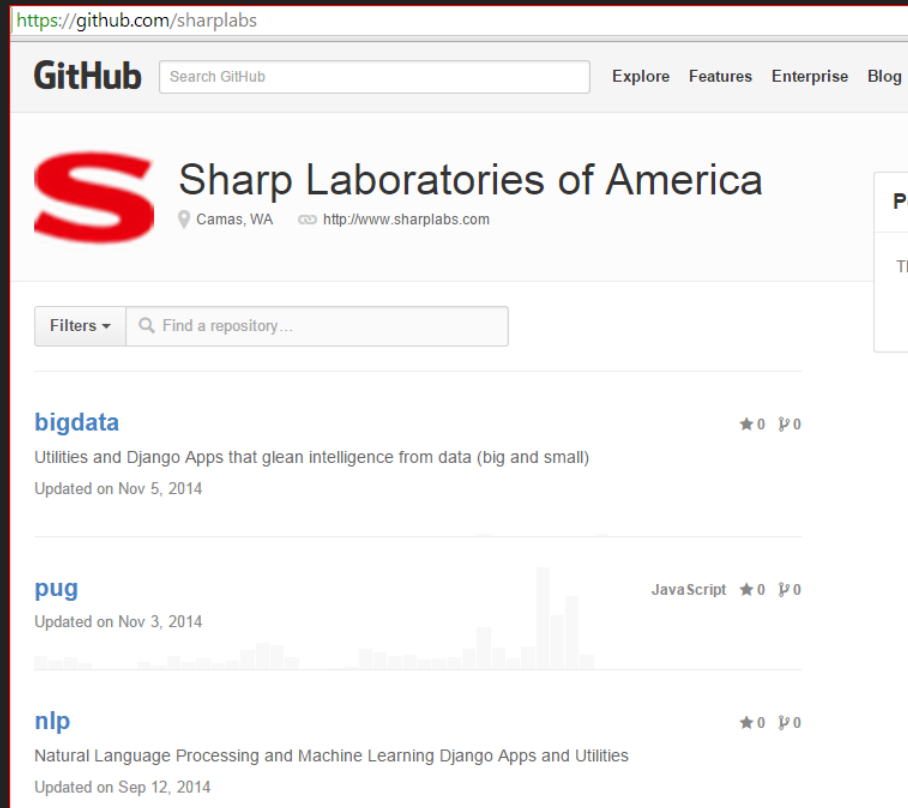
Built into python graph library (`networkx`)

```
from networkx.algorithms.shortest_paths import astar_path  
astar_path(G, source, target, heuristic=None)
```

You better have a good heuristic!

It's Open Source!

github.com/sharplabs



The screenshot shows the GitHub profile page for Sharp Laboratories of America. The page header includes the GitHub logo, a search bar, and navigation links for Explore, Features, Enterprise, and Blog. The profile section displays the organization's name, a red 'S' logo, location (Camas, WA), and website (http://www.sharplabs.com). Below this is a search bar with a 'Filters' dropdown and a 'Find a repository...' button. The repository list shows three items: 'bigdata' (Utilities and Django Apps that glean intelligence from data (big and small), updated Nov 5, 2014), 'pug' (JavaScript, updated Nov 3, 2014, with a commit history bar chart), and 'nlp' (Natural Language Processing and Machine Learning Django Apps and Utilities, updated Sep 12, 2014). Each repository has a star icon and a fork icon, both showing 0.

<https://github.com/sharplabs>

GitHub Search GitHub Explore Features Enterprise Blog

S Sharp Laboratories of America
Camas, WA <http://www.sharplabs.com>

Filters Find a repository...

bigdata ★ 0 0
Utilities and Django Apps that glean intelligence from data (big and small)
Updated on Nov 5, 2014

pug JavaScript ★ 0 0
Updated on Nov 3, 2014

nlp ★ 0 0
Natural Language Processing and Machine Learning Django Apps and Utilities
Updated on Sep 12, 2014

bit.ly/pawsvote
Choose Your Story

7707-2-TOTAL
(770) 728-6825

1. Only Nyquist Knows
 2. The Meaning of Mean
 3. Data Dearth
 4. Question the Question
 5. Deep Net Runs Aground
 6. Escape the Maze
- Consider sample rate
 - Classify before mean
 - Explore data sources
 - Reject rate metric
 - $\text{data} > \text{nodes} \times \text{inputs}$
 - Lazy correlation

References

- "Data Driven Documents"
 - 2011, Mike Bostock
- "Data Science with `pug`"
 - 2014, Lane, Zen, Kowalski, PDX Python U.G.
- "Neural Network Design"
 - 2014, Hagan, Demuth, et. al., OKSU
- "Forecasting Product Returns"
 - 2001, Toktay, INSEAD
- `scipy.ransac`
 - 2014, Andrew D. Straw
- "Choose Your Own Adventure Presentation"
 - 2014, Matt Makai