



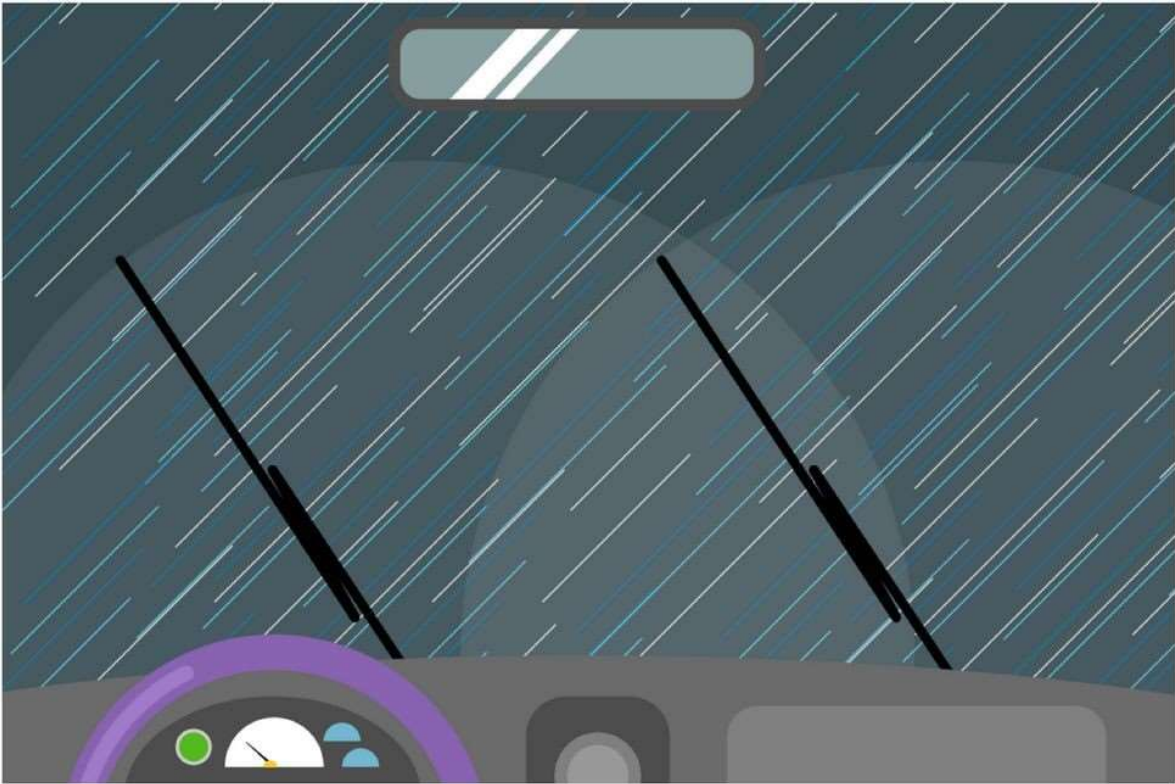
DATADOG

Automatically Detecting Anomalies and Outliers in Real-Time

Outline

- Monitoring
- Alerting
- Outlier vs. Anomaly Detection
- Outlier Detection Algorithms
- Anomaly Detection Algorithms

Monitor Everything



Monitor Everything

Datadog gathers performance data from all your application components.



AWS



Docker



CoreOS



Chef



Puppet



Github



Pagerduty



Nagios



Go



Postgres



Java



VMware



Redis



MySQL



Apache



Tomcat



MongoDB

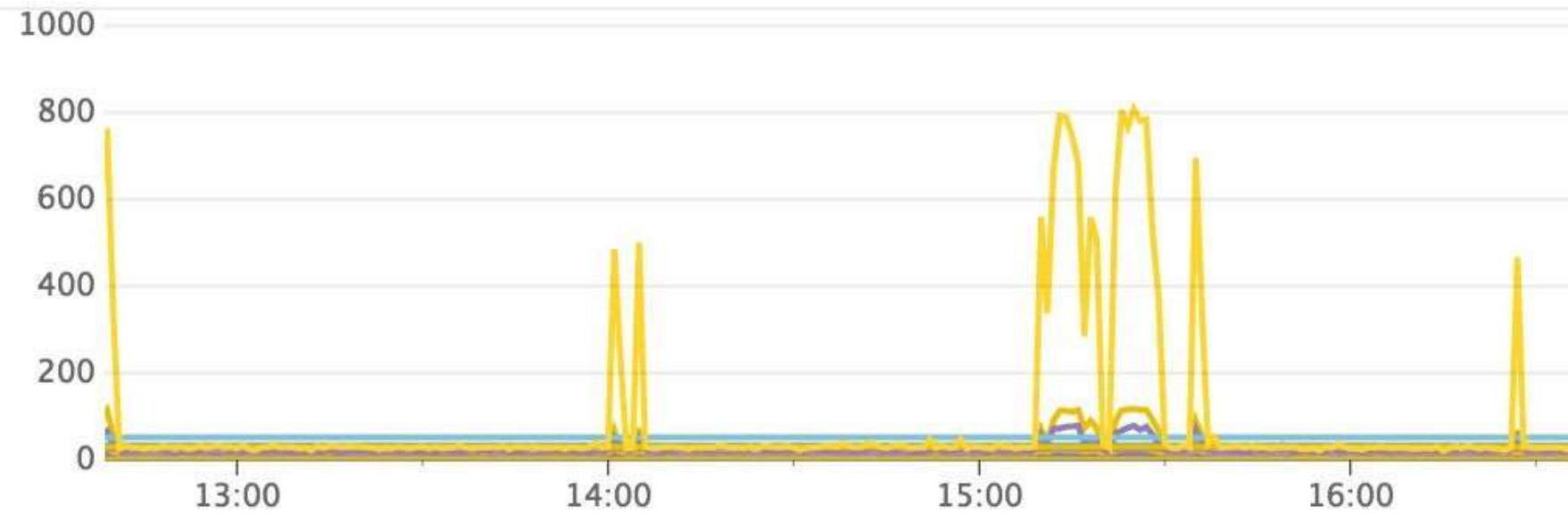


New Relic

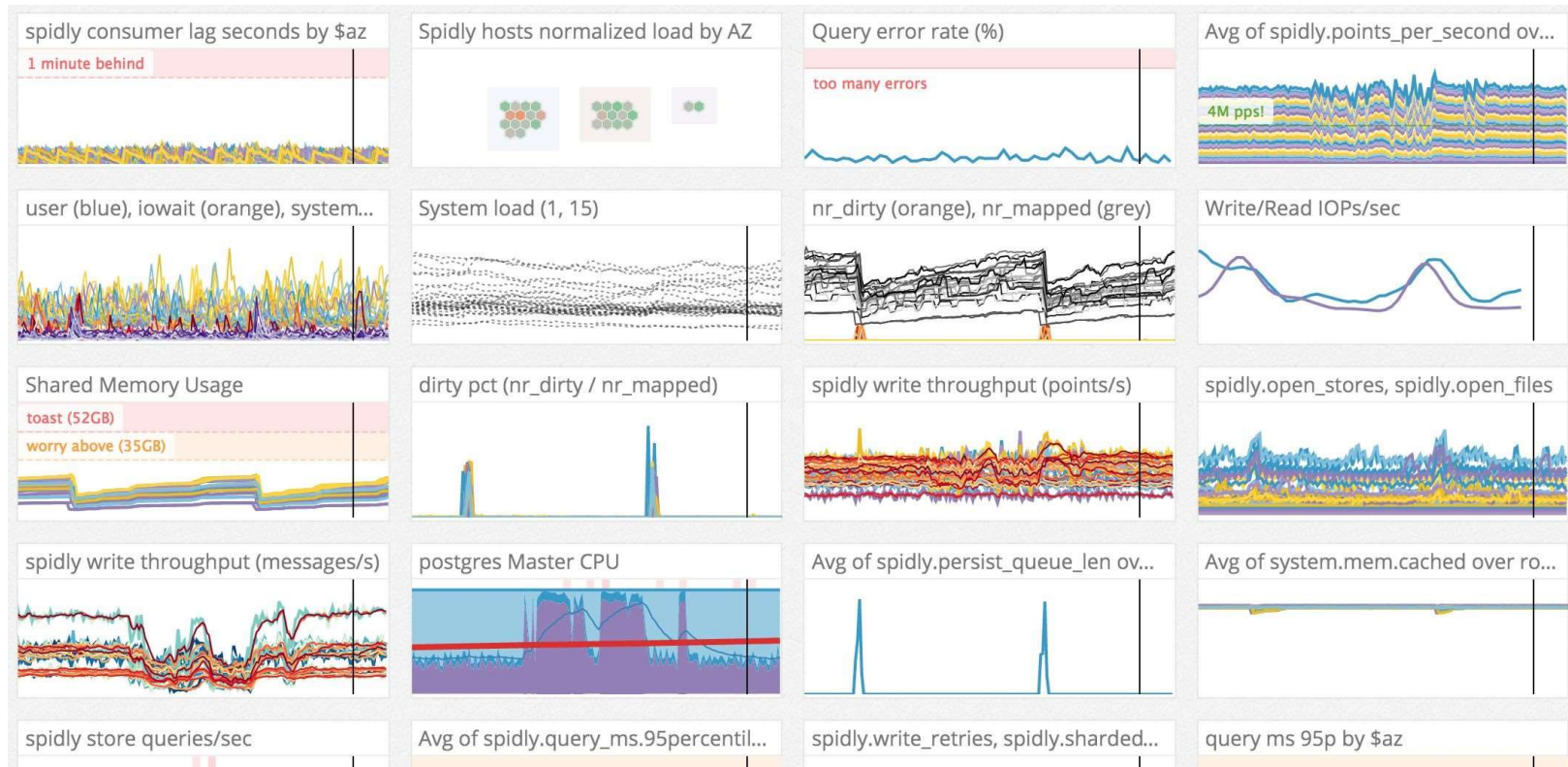
Monitor Everything



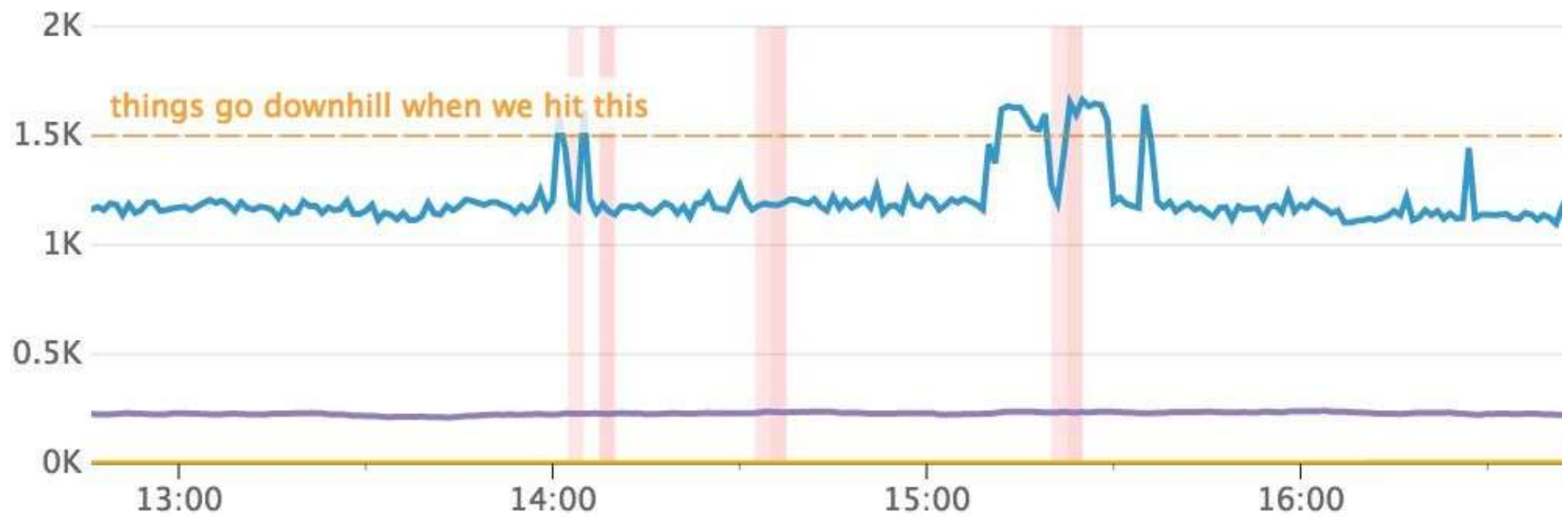
Monitor Everything



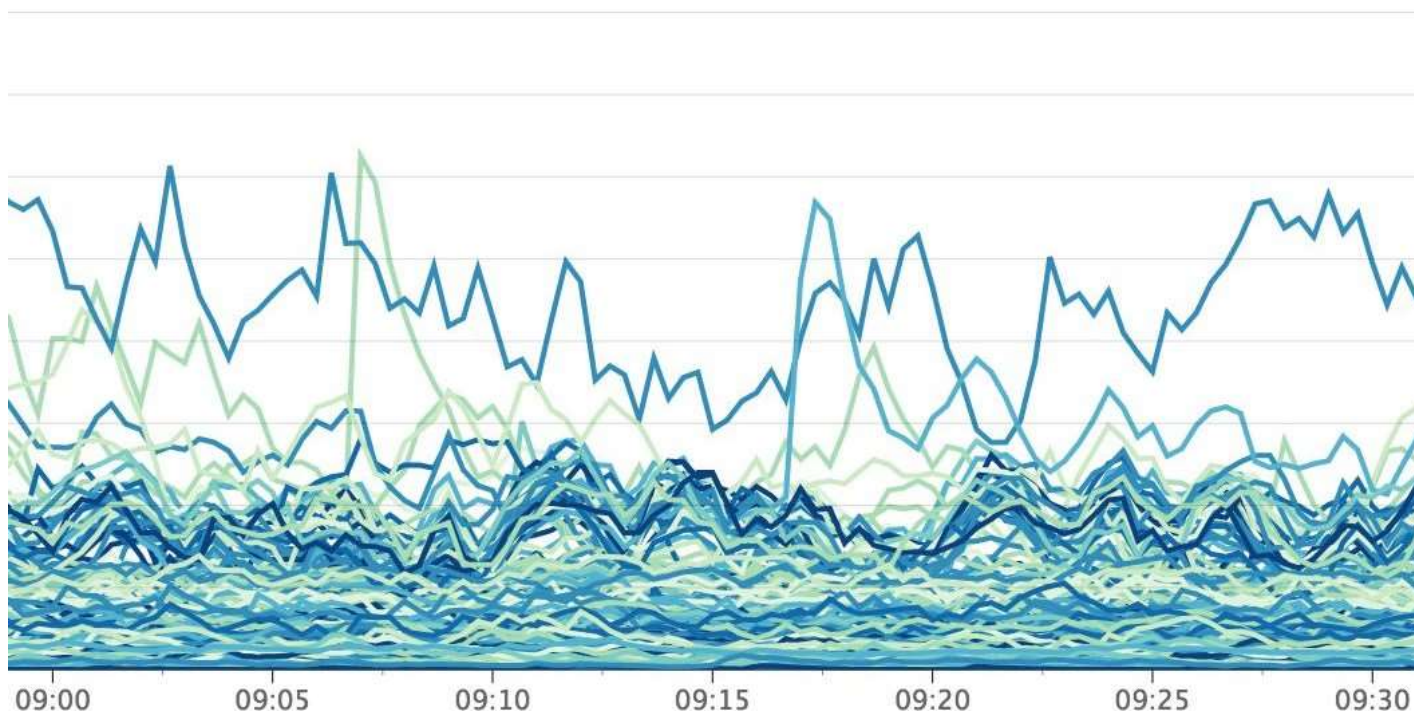
Monitor Everything?



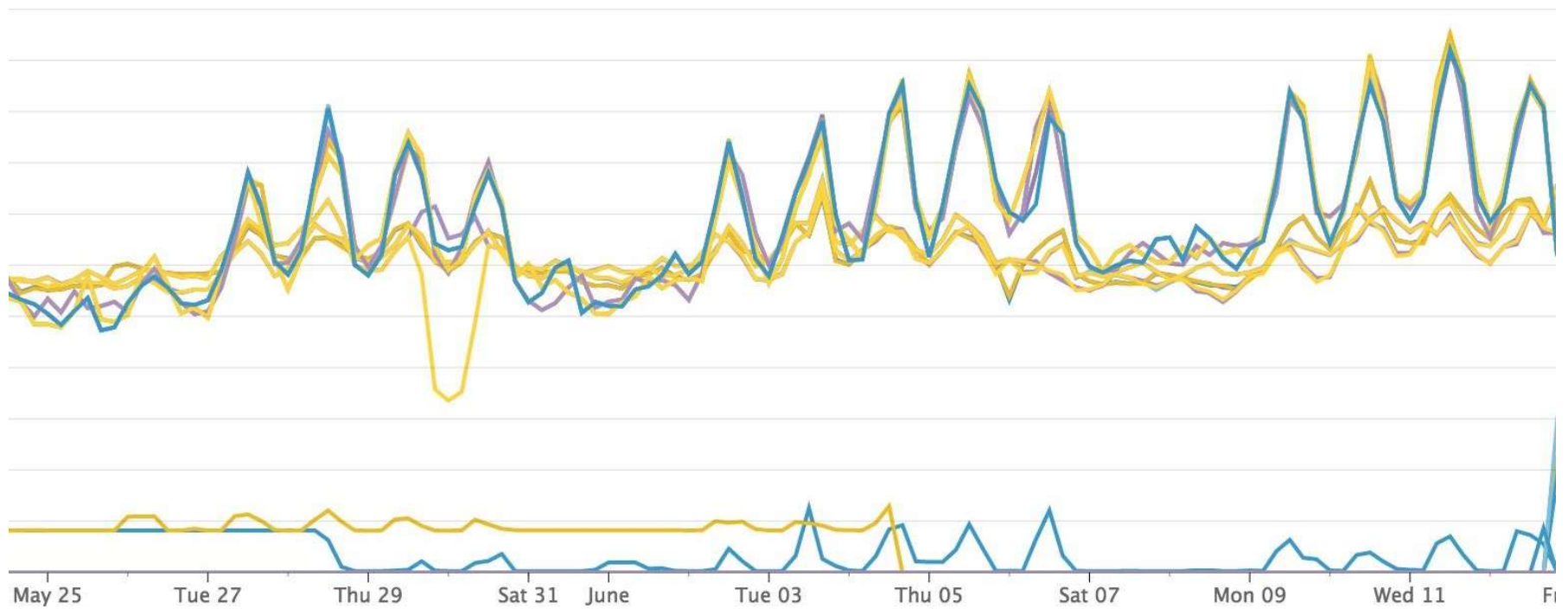
Alerting



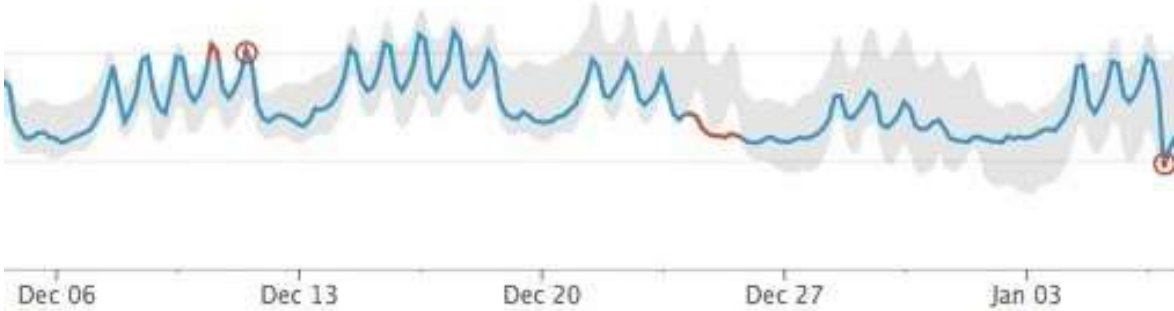
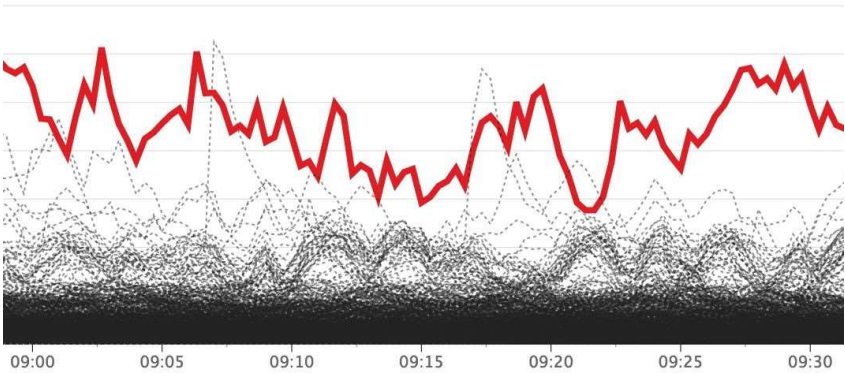
Alerting?



Alerting?

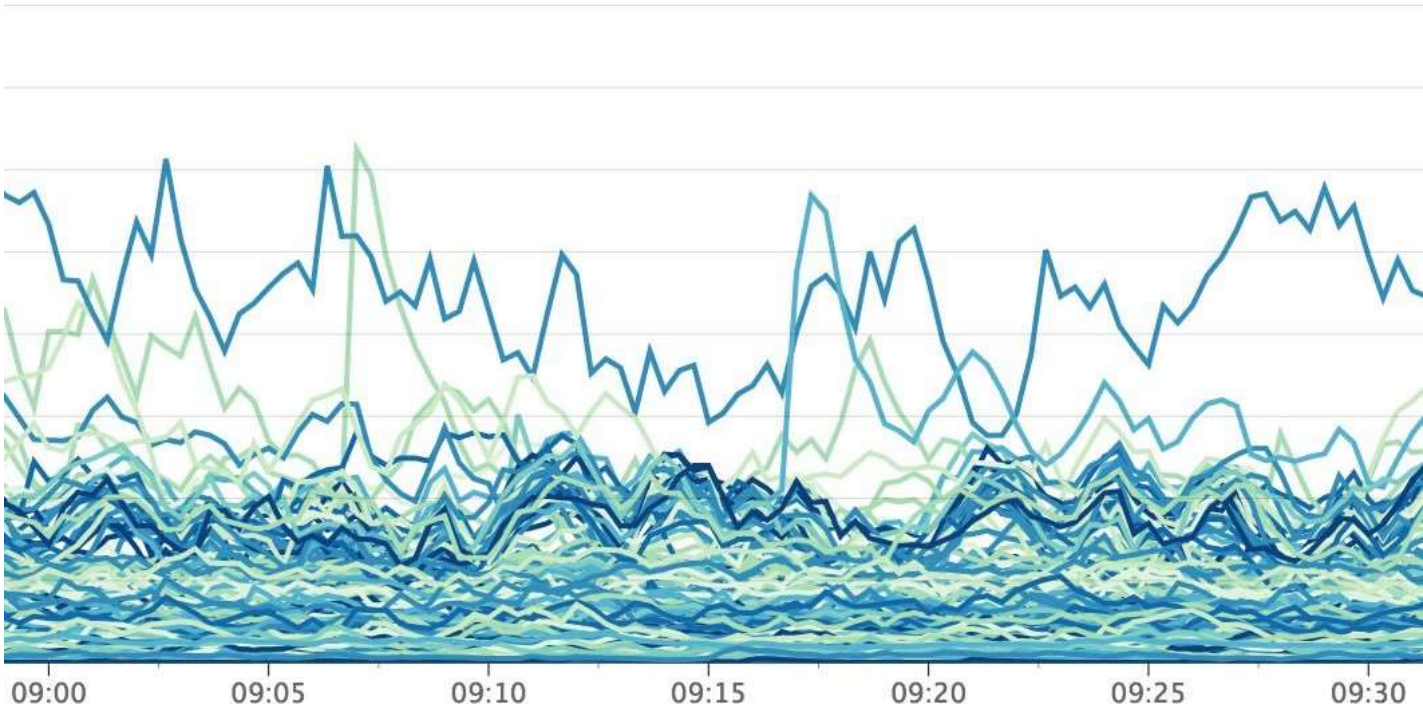


Outlier and Anomaly Detection

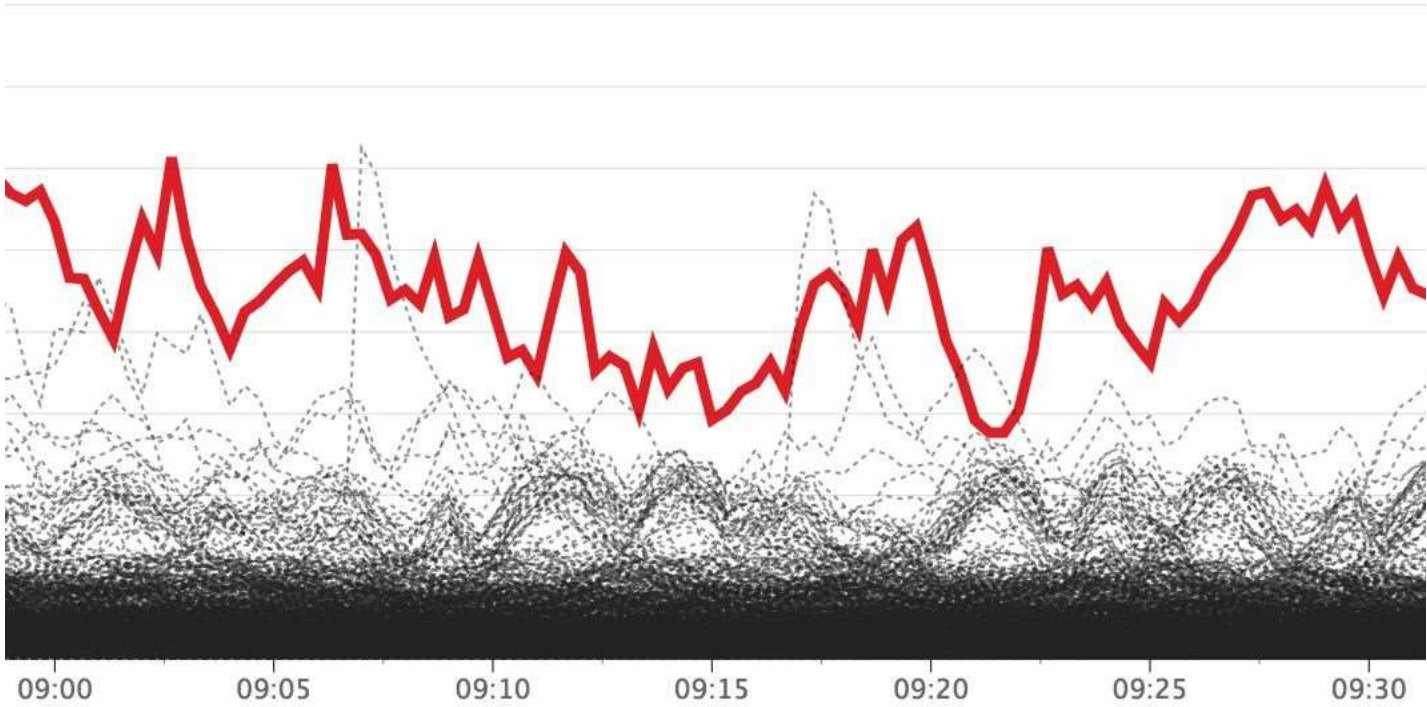


Outlier Detection

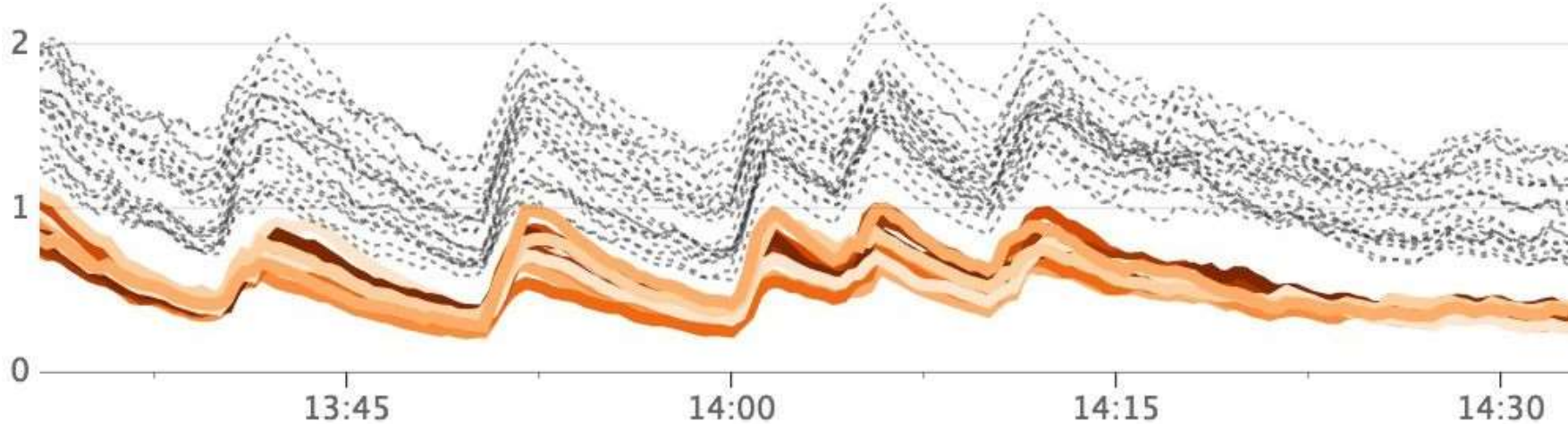
Outlier Detection



Outlier Detection



Outlier Detection



Outlier Detection Algorithms

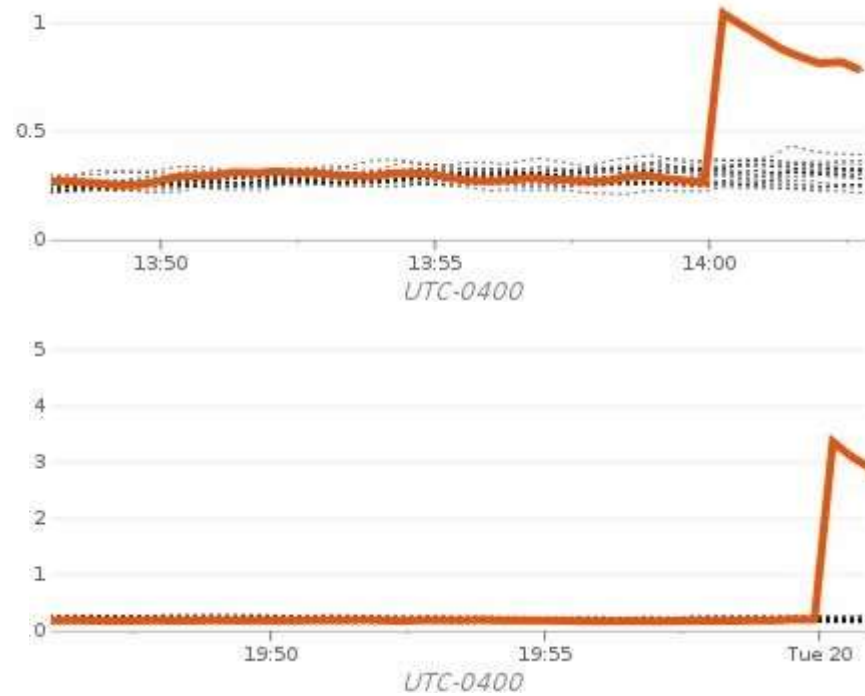
MAD

median absolute deviation

DBSCAN

density-based spatial clustering of applications with noise

Robust Outlier Detection Algorithms



Median Absolute Deviation

$$\text{MAD}(D) = \text{median}(\{|d_i - \text{median}(D)|\})$$

Median Absolute Deviation

$$\text{MAD}(D) = \text{median}(\{|d_i - \text{median}(D)|\})$$

$$D = \{1, 2, 3, 4, 5, 6, 100\}$$

Median Absolute Deviation

$$\text{MAD}(D) = \text{median}(\{|d_i - \text{median}(D)|\})$$

$$D = \{1, 2, 3, 4, 5, 6, 100\}$$

median = 4

Median Absolute Deviation

$$\text{MAD}(D) = \text{median}(\{|d_i - \text{median}(D)|\})$$

$$D = \{1, 2, 3, 4, 5, 6, 100\}$$

$$\text{median} = 4$$

$$\text{deviations} = \{-3, -2, -1, 0, 1, 2, 96\}$$

Median Absolute Deviation

$$\text{MAD}(D) = \text{median}(\{|d_i - \text{median}(D)|\})$$

$$D = \{1, 2, 3, 4, 5, 6, 100\}$$

$$\text{median} = 4$$

$$\text{deviations} = \{-3, -2, -1, 0, 1, 2, 96\}$$

$$\text{abs deviations} = \{0, 1, 1, 2, 2, 3, 96\}$$

Median Absolute Deviation

$$\text{MAD}(D) = \text{median}(\{|d_i - \text{median}(D)|\})$$

$$D = \{1, 2, 3, 4, 5, 6, 100\}$$

$$\text{median} = 4$$

$$\text{deviations} = \{-3, -2, -1, 0, 1, 2, 96\}$$

$$\text{abs deviations} = \{0, 1, 1, 2, 2, 3, 96\}$$

$$\text{MAD} = 2$$

Median Absolute Deviation

$$\text{MAD}(D) = \text{median}(\{|d_i - \text{median}(D)|\})$$

$$D = \{1, 2, 3, 4, 5, 6, 100\}$$

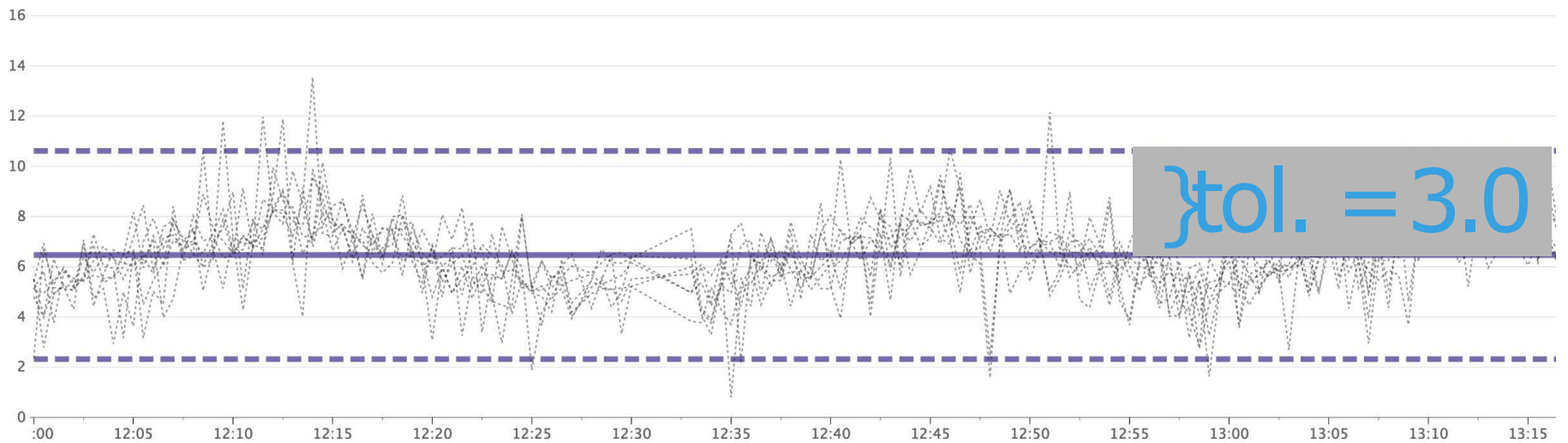
$$\text{median} = 4$$

$$\text{deviations} = \{-3, -2, -1, 0, 1, 2, 96\}$$

$$\text{abs deviations} = \{0, 1, 1, 2, 2, 3, 96\}$$

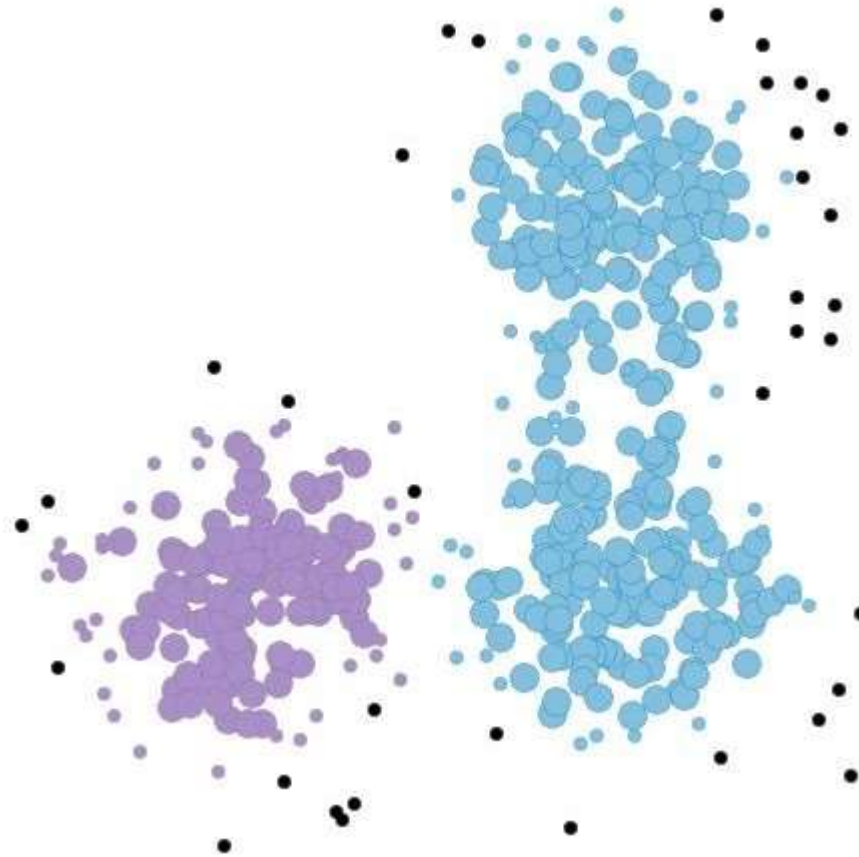
$$\text{MAD} = 2 \quad (\text{std dev} = 33.8)$$

Median Absolute Deviation



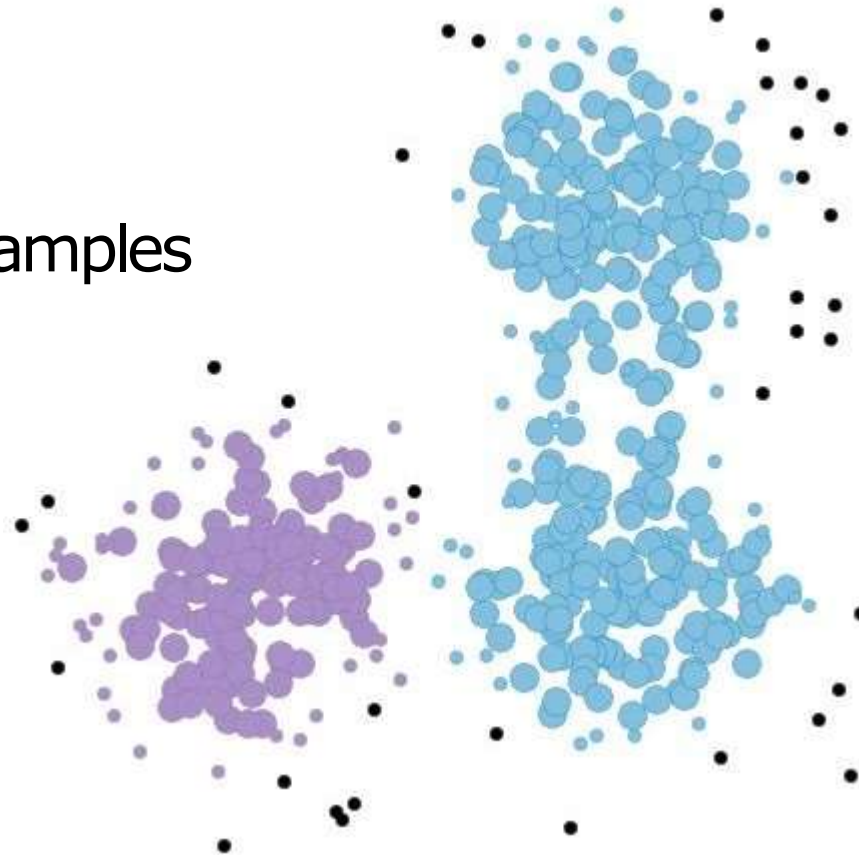
Parameters: Tolerance, Pct

DBSCAN

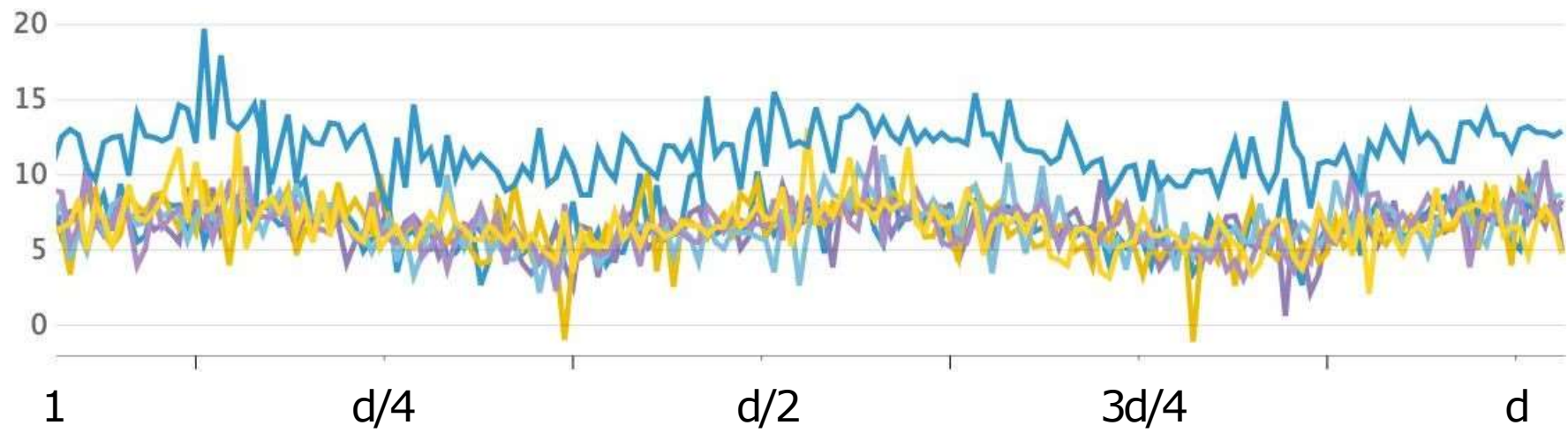


DBSCAN

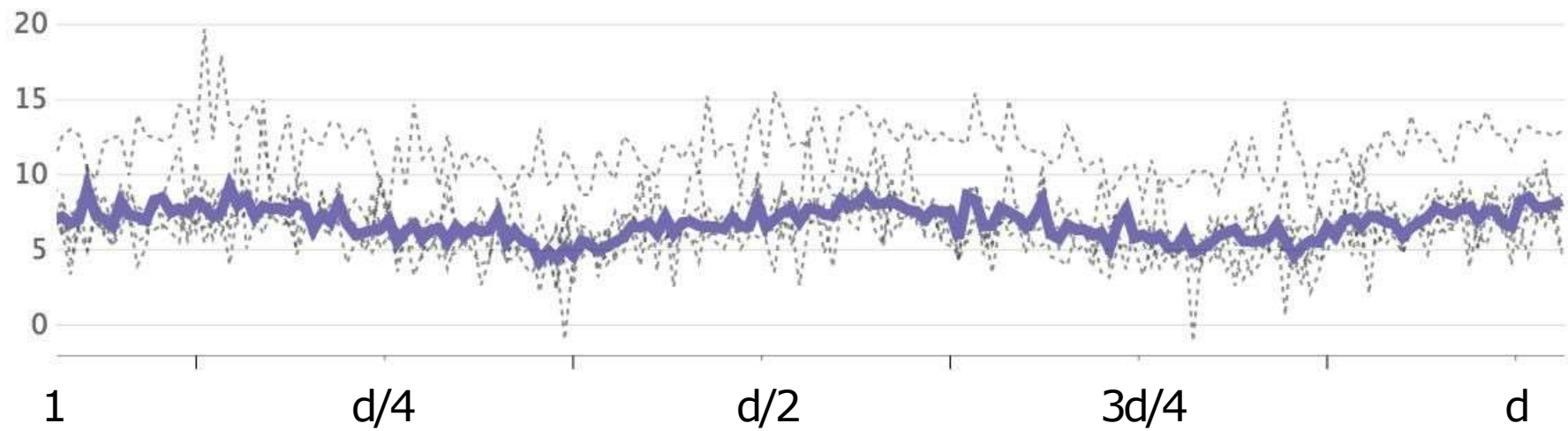
Parameters:
epsilon, min_samples



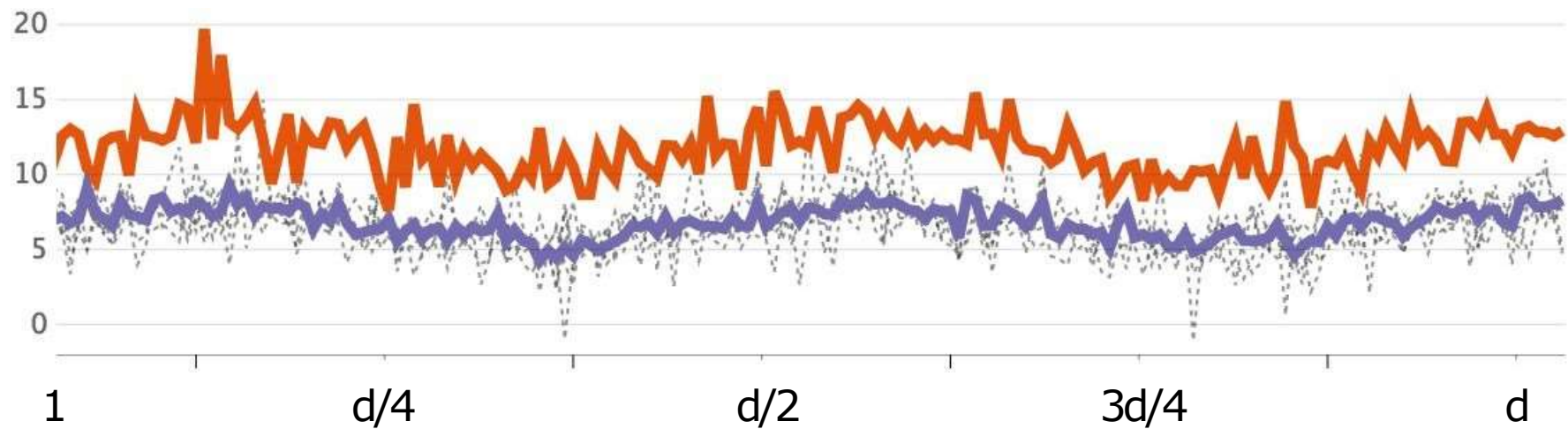
DBSCAN



DBSCAN

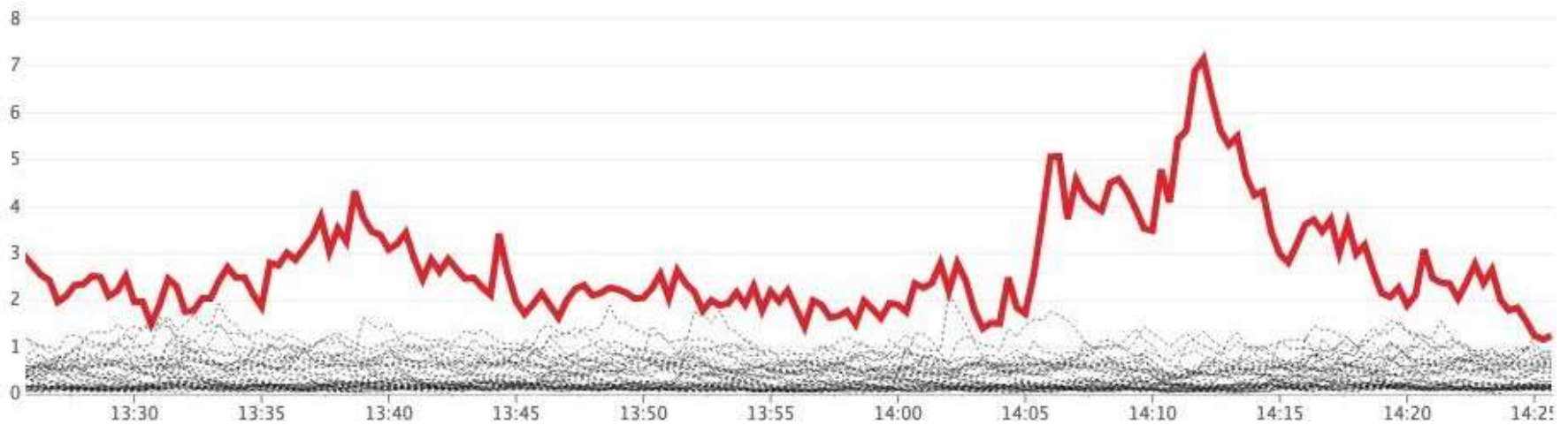


DBSCAN

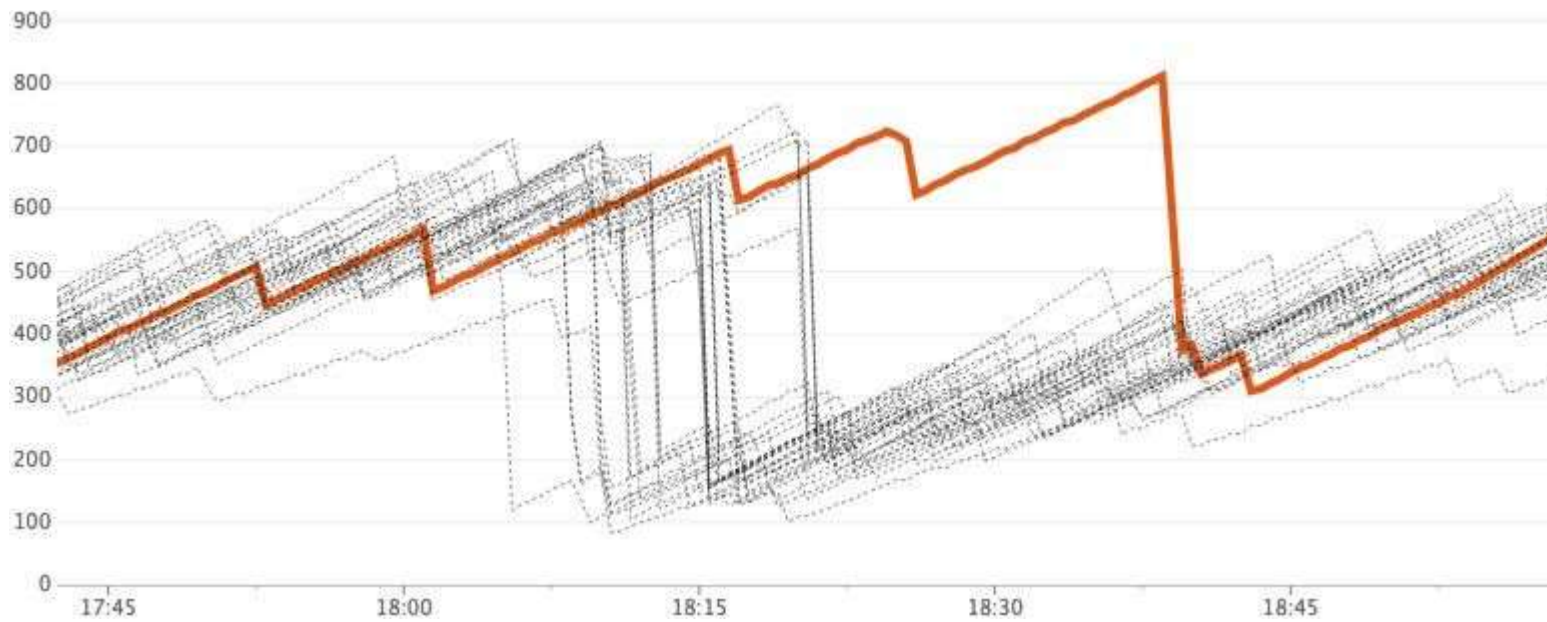


$\epsilon \sim \text{median}(\text{dist from median series}) \times \text{tolerance}$

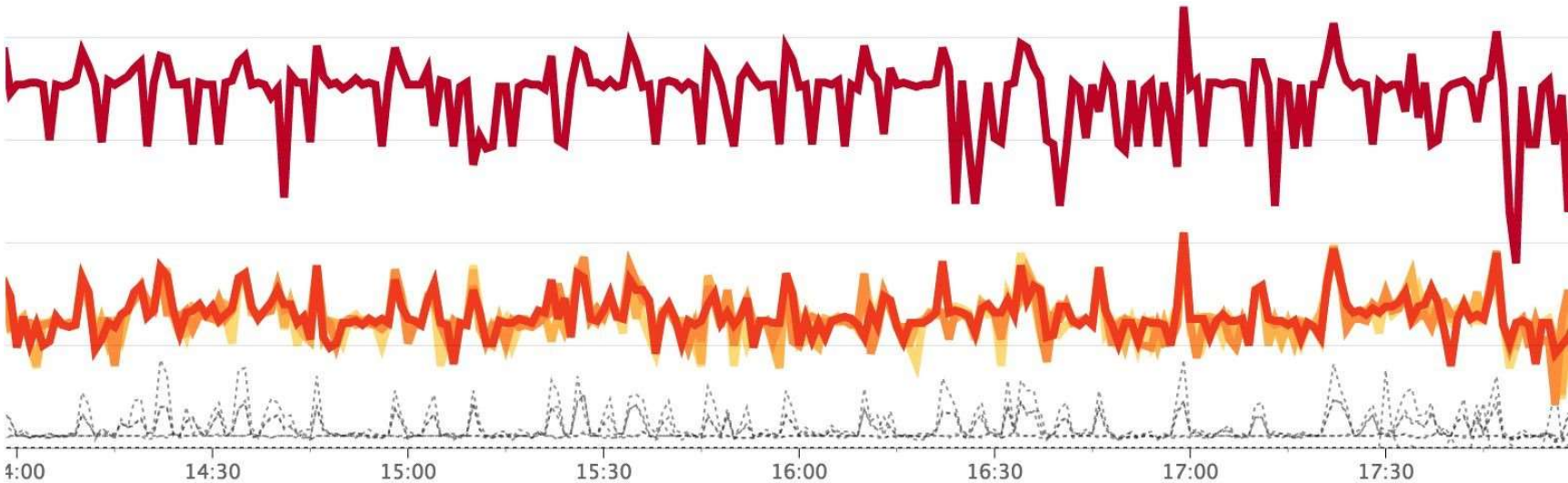
MAD or DBSCAN?



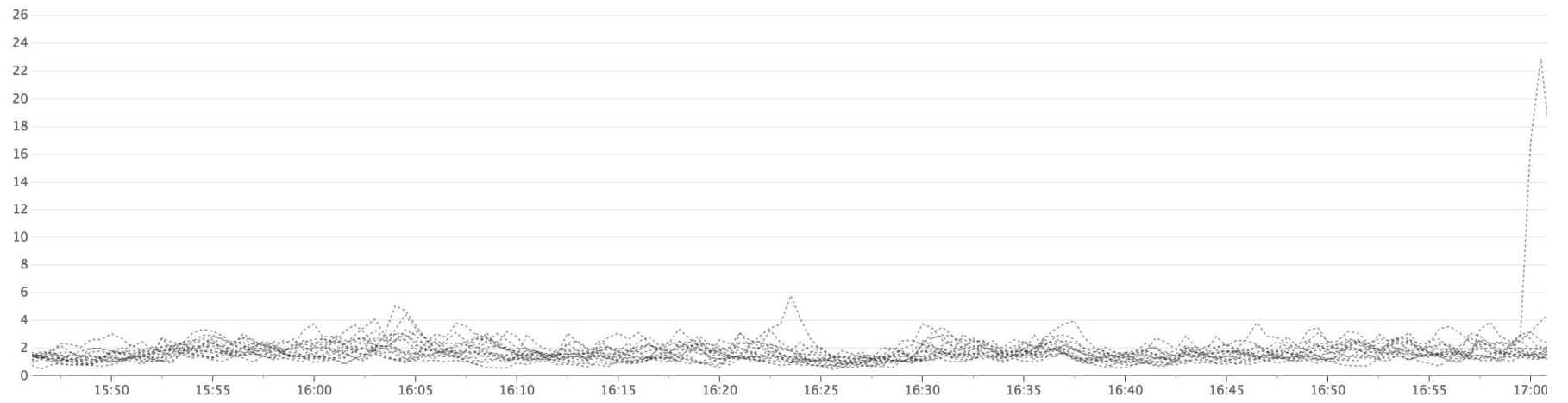
MAD or DBSCAN?



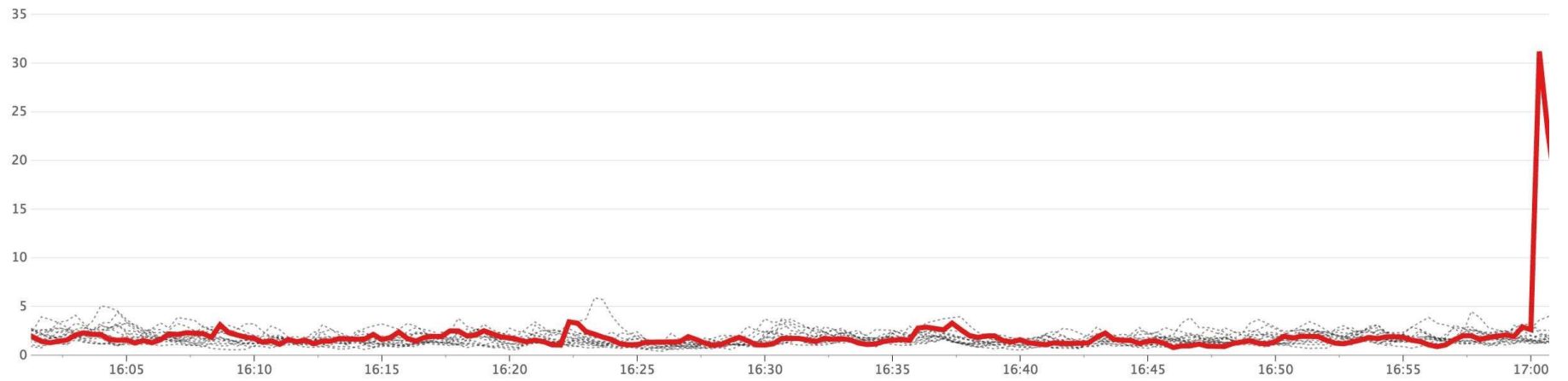
Some subtleties



Some subtleties

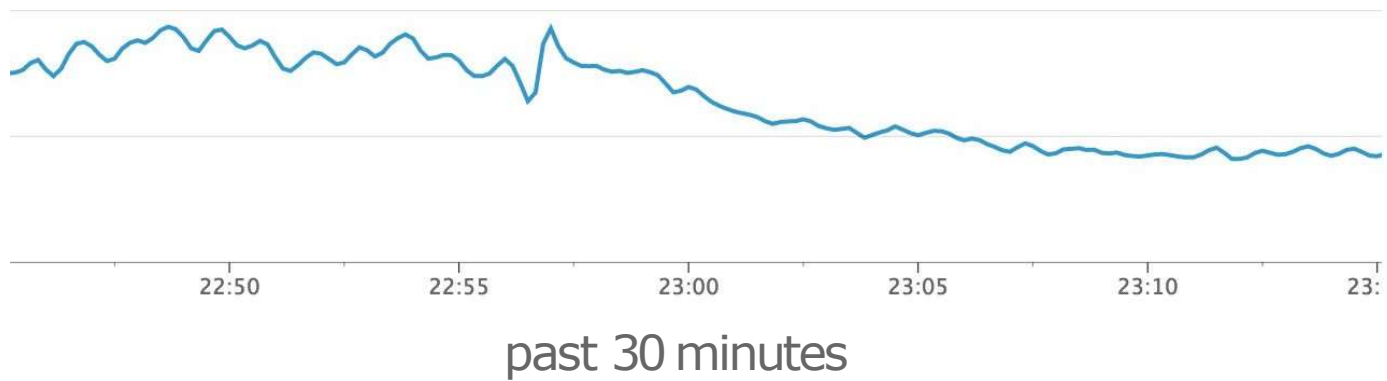


Some subtleties

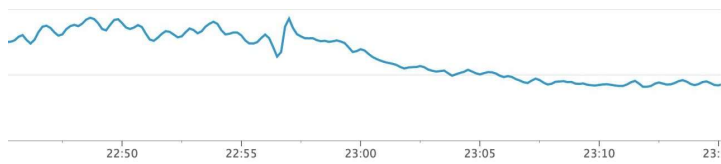


Anomaly Detection

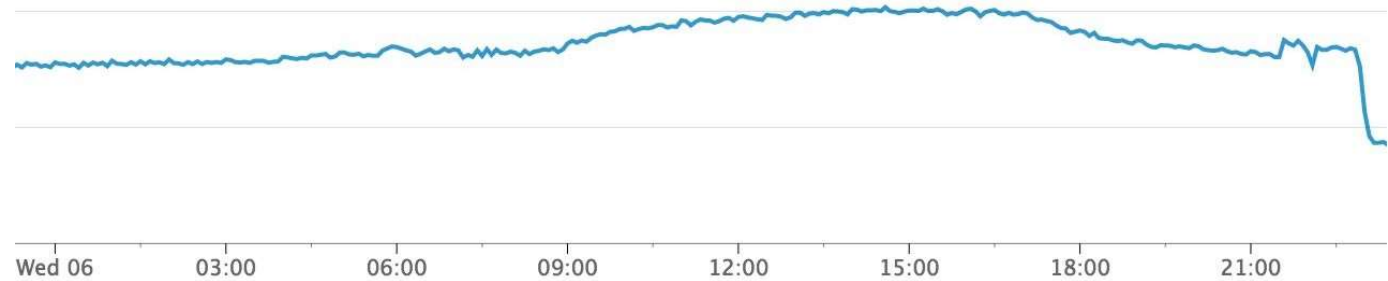
An Investigation



An Investigation

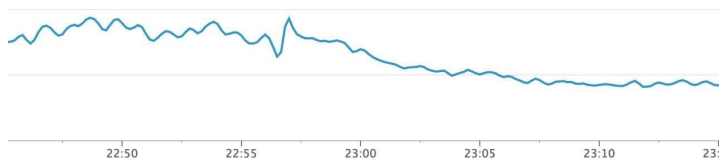


past 30 minutes

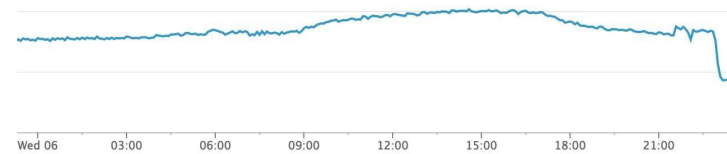


past day

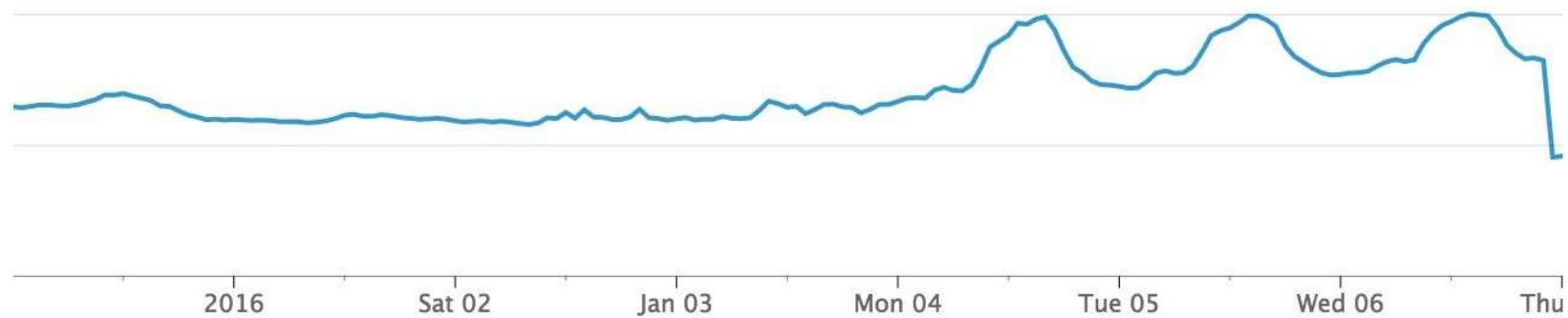
An Investigation



past 30 minutes

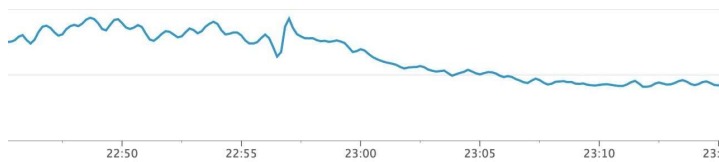


past day

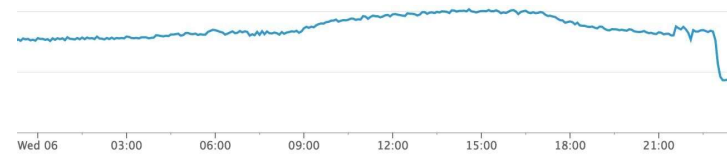


past week

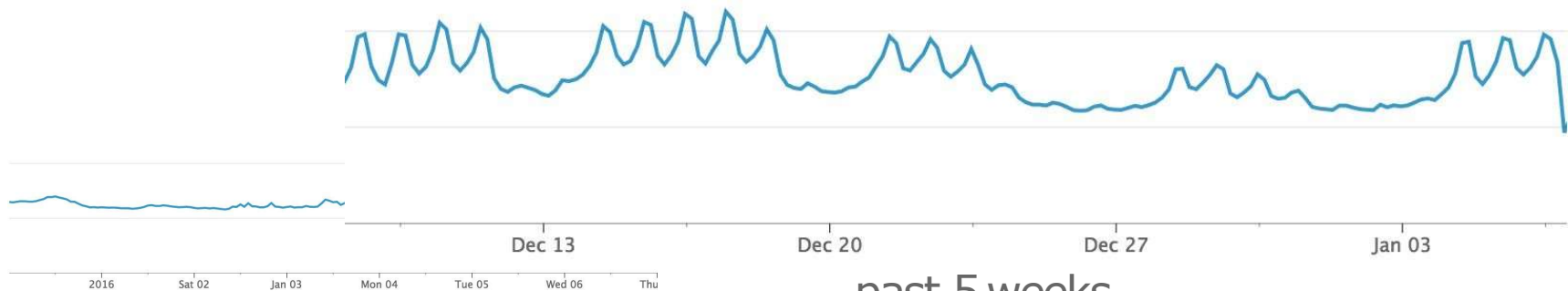
An Investigation



past 30 minutes



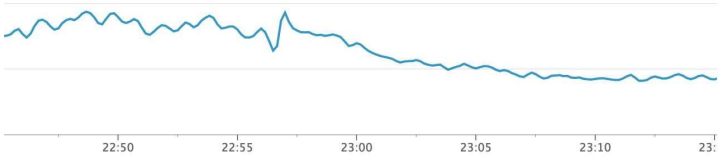
past day



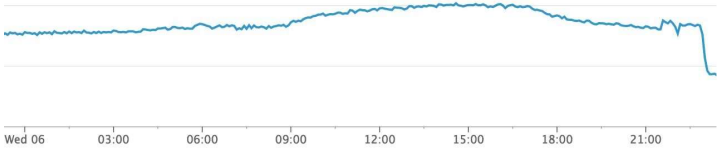
past week

past 5 weeks

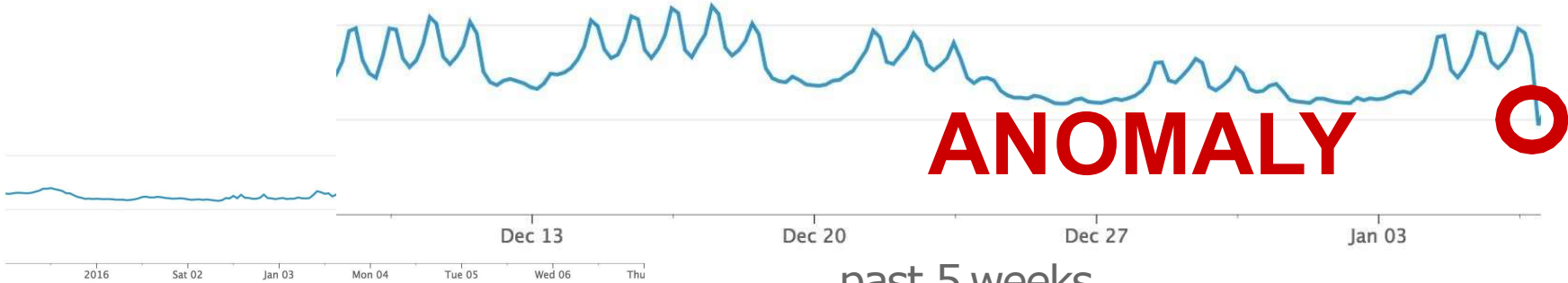
An Investigation



past 30 minutes



past day



ANOMALY

past week

past 5 weeks

Anomalies

A time series point is an anomaly if:

- Given the past points in the series ($\textcircled{R}\textcircled{R}\textcircled{R}\textcircled{R}\textcircled{R}$), the point in question (\textcircled{R}) is unlikely given your model of the past;

Anomalies

A time series point is an anomaly if:

- Given the past points in the series ($\textcircled{R}\textcircled{R}\textcircled{R}\textcircled{R}\textcircled{R}$), the point in question (\textcircled{R}) is unlikely given your model of the past;

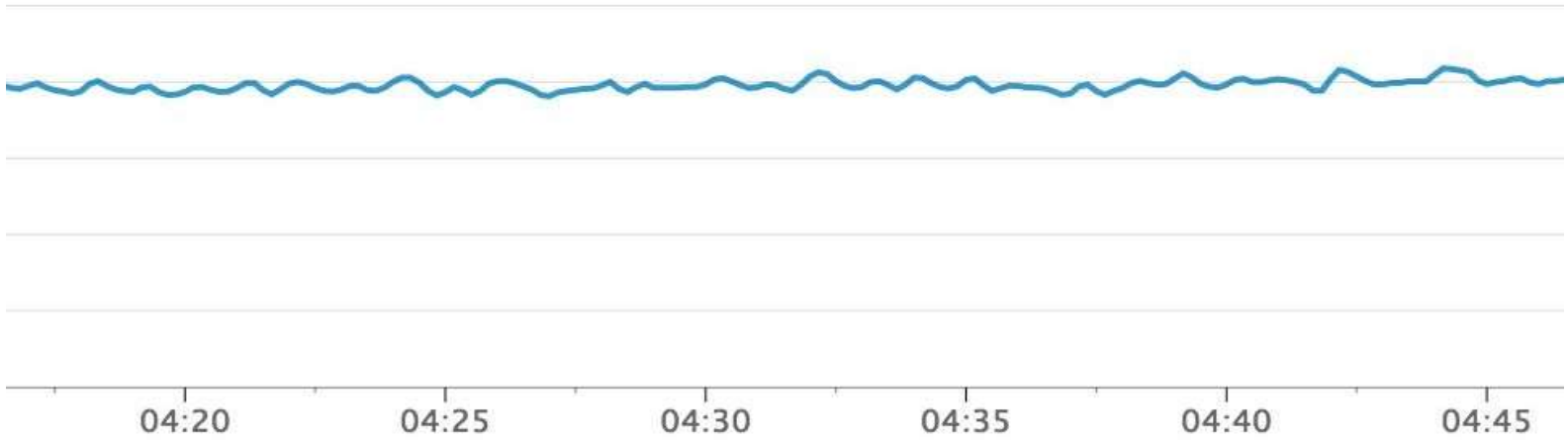
and you should alert on a set of anomalies if:

- they are a symptom of an issue you care about (\hat{I}).

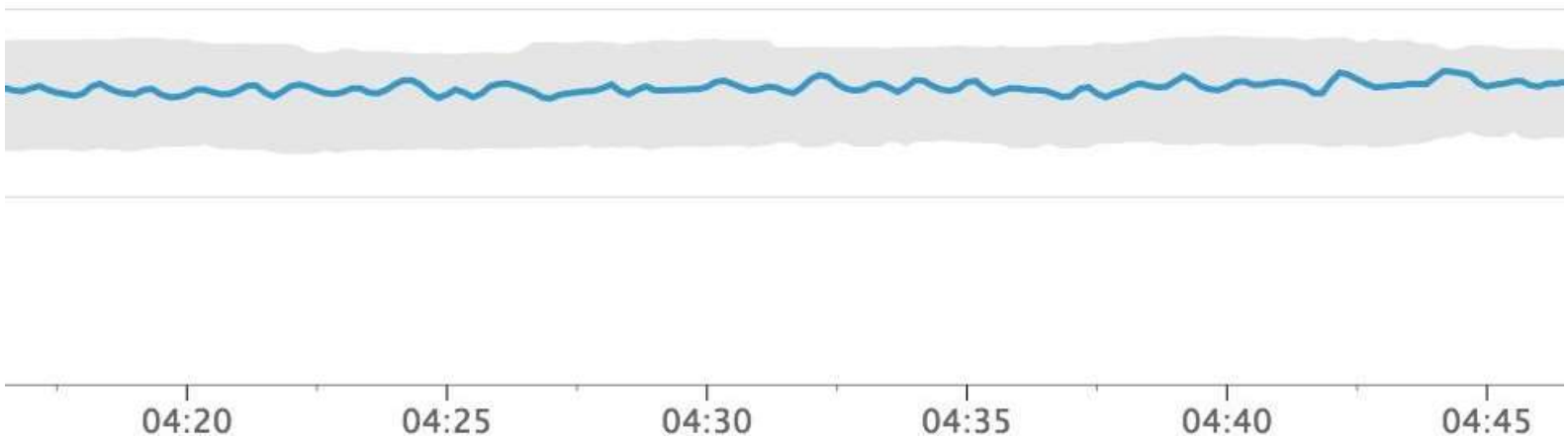
Our Approach

1. Extract as much signal as we can from the timeseries.
2. Use robust statistical measures when creating the model.
3. Give the user control over when they get alerted.

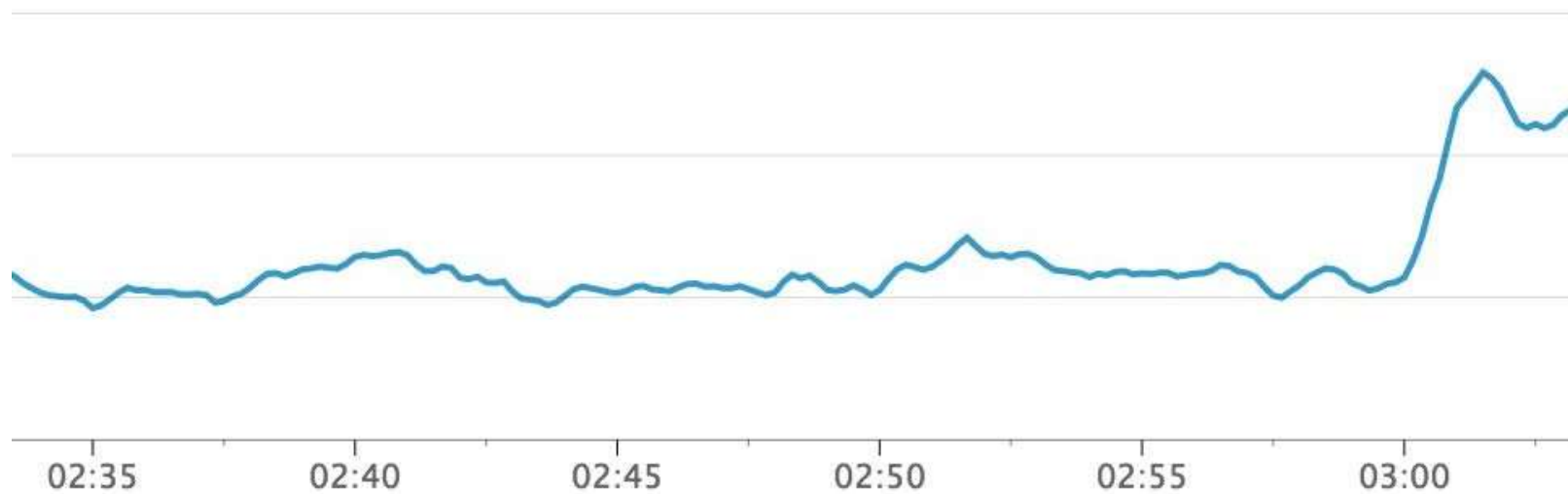
What's Normal?



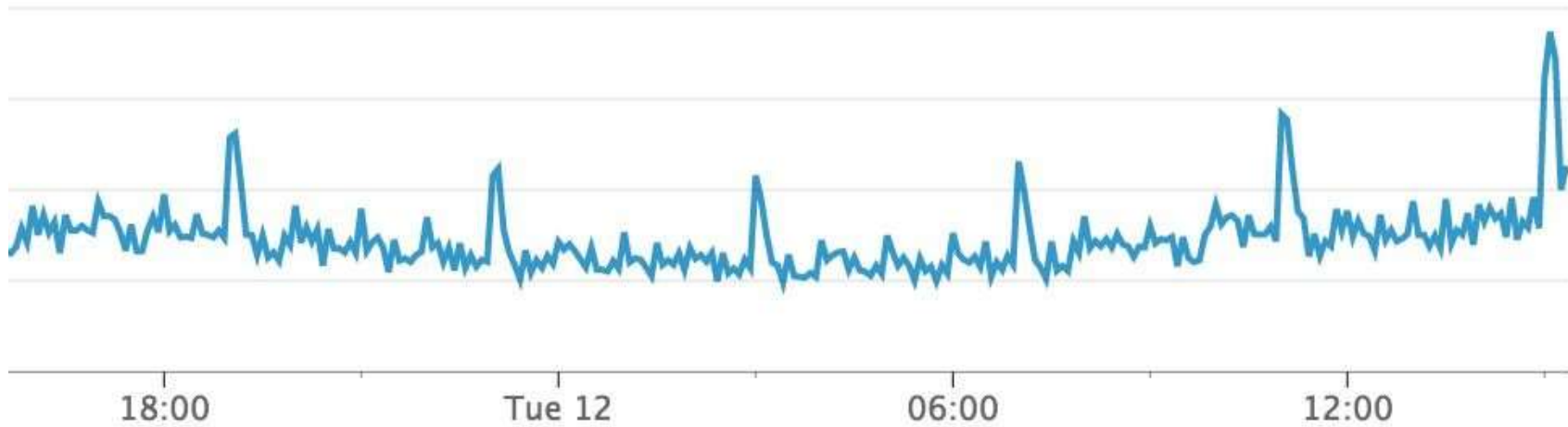
What's Normal?



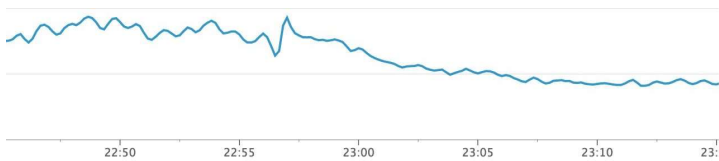
What's Normal?



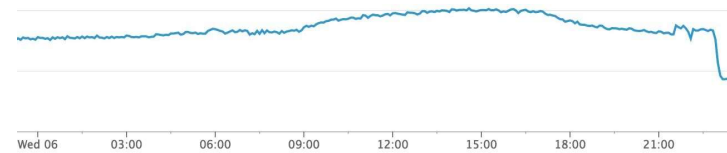
What's Normal?



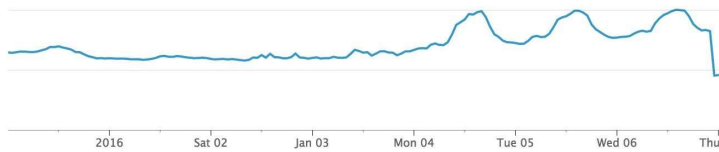
Past Performance...



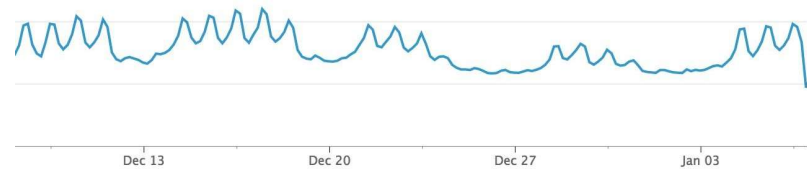
past 30 minutes



past day

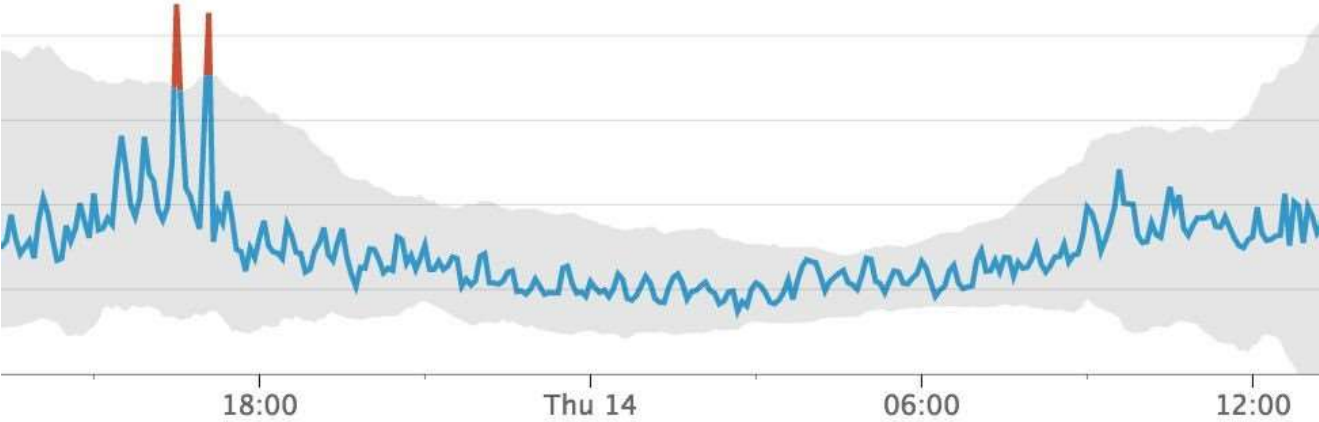


past week

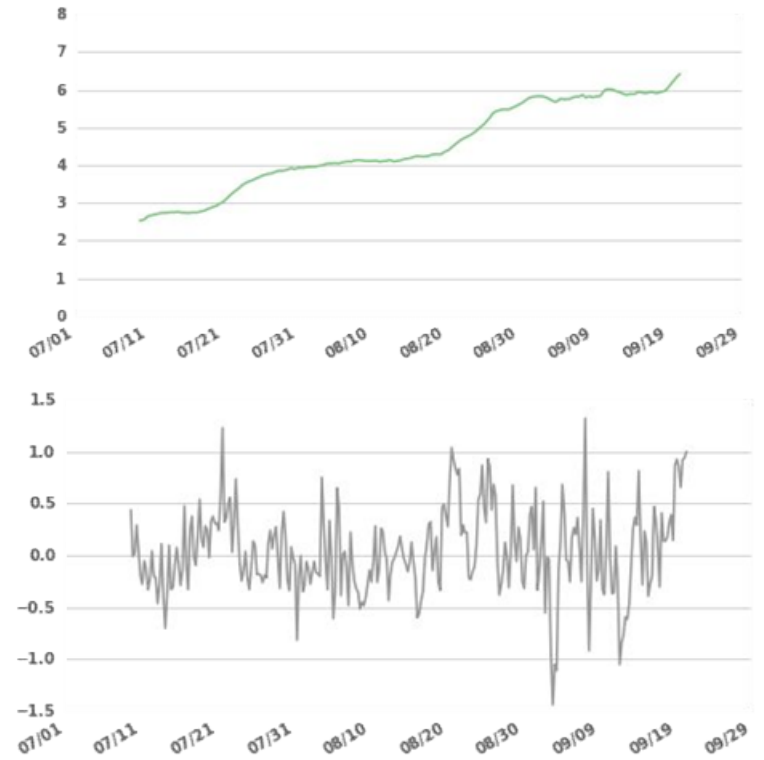
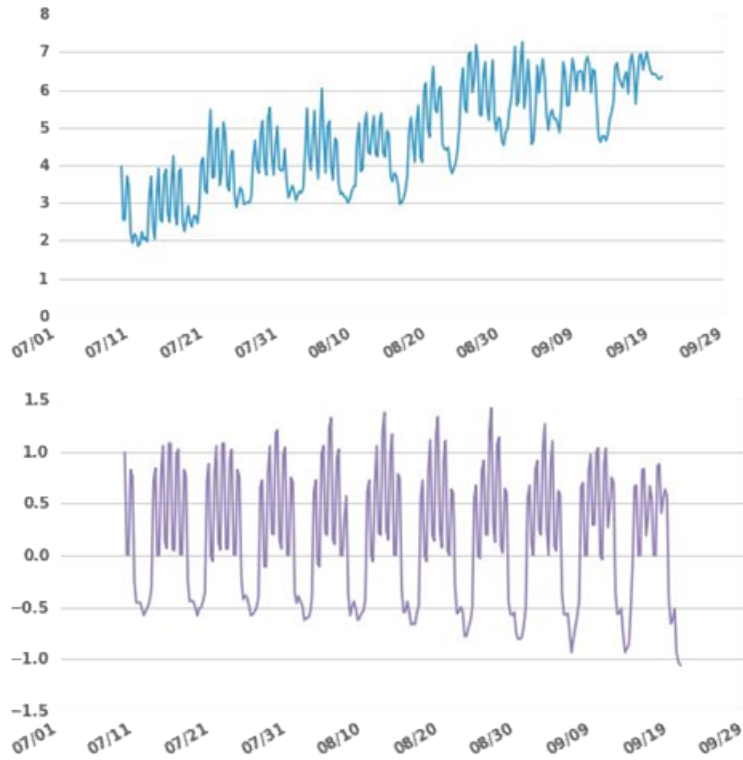


past 5 weeks

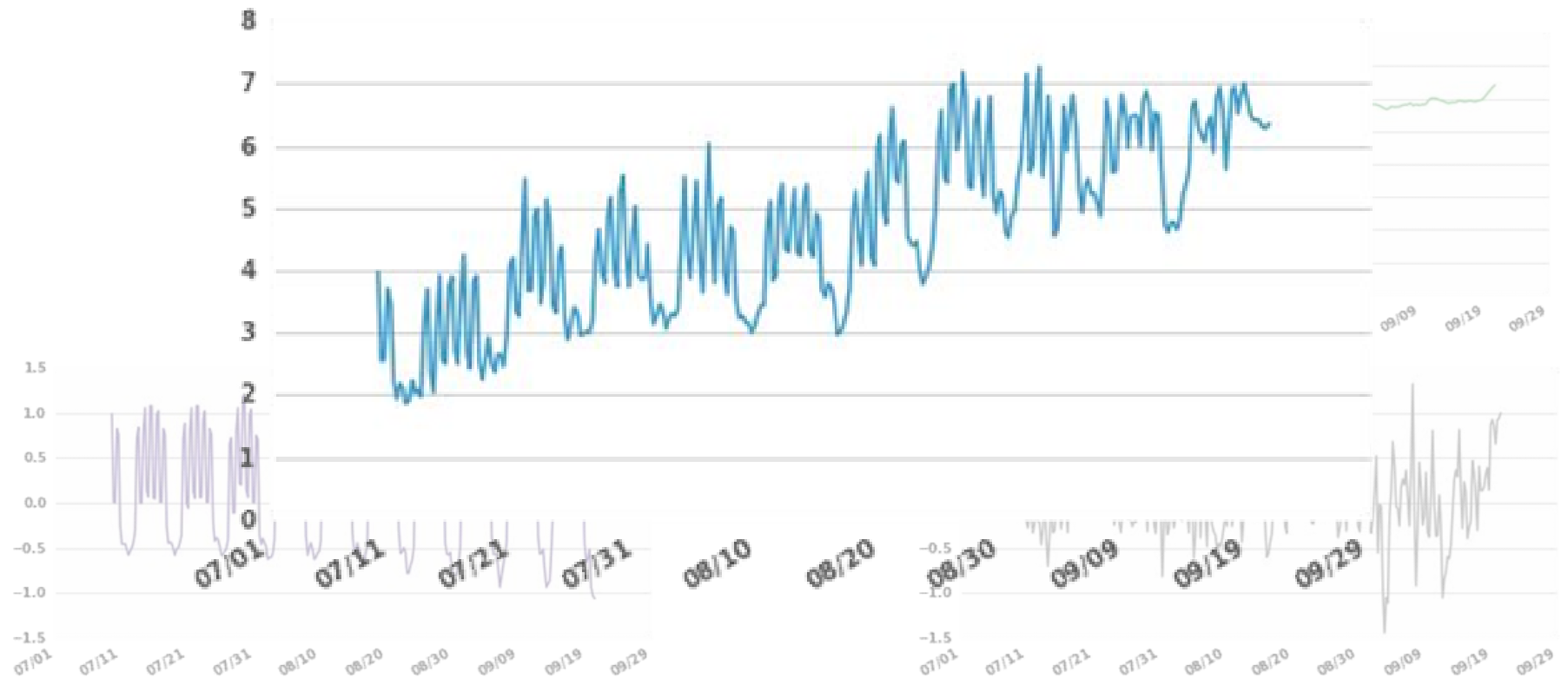
Past Performance...



Decomposition



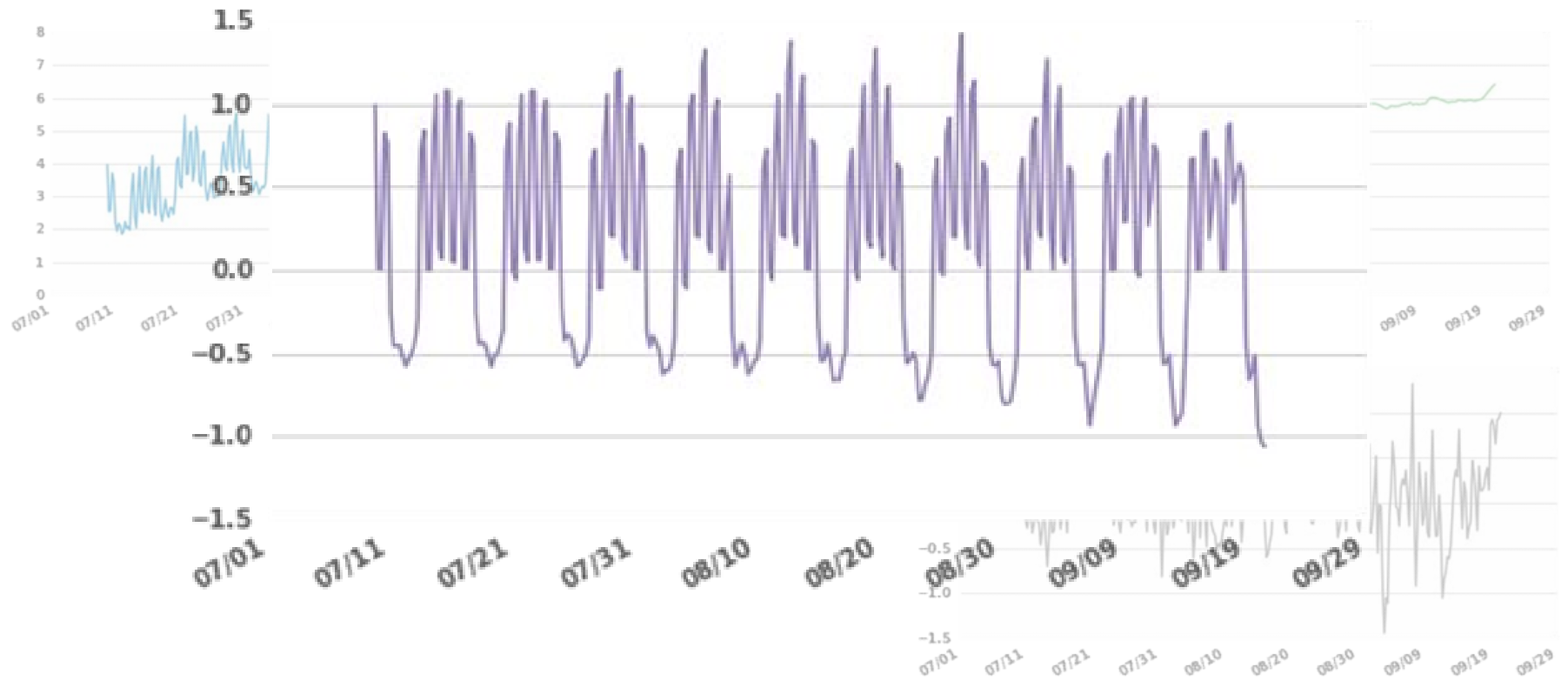
Decomposition



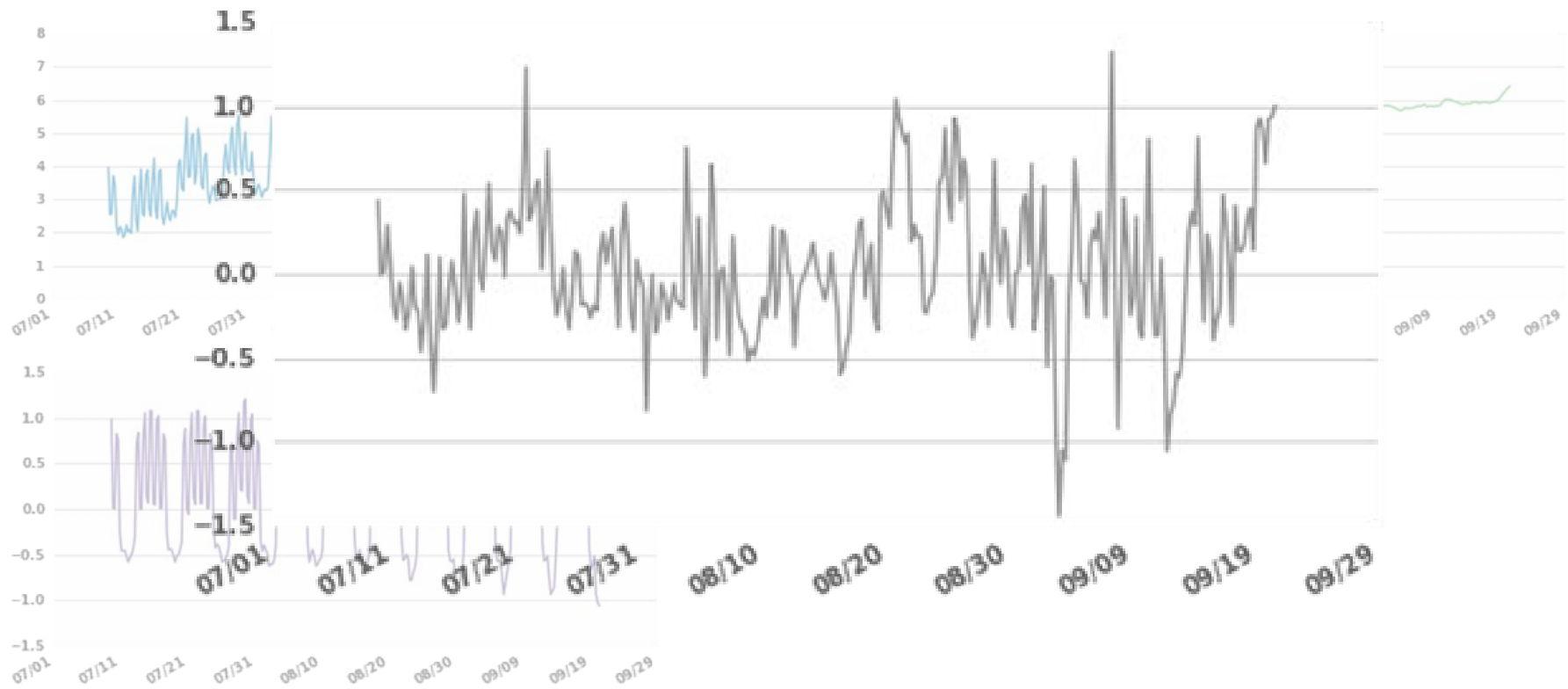
Decomposition



Decomposition



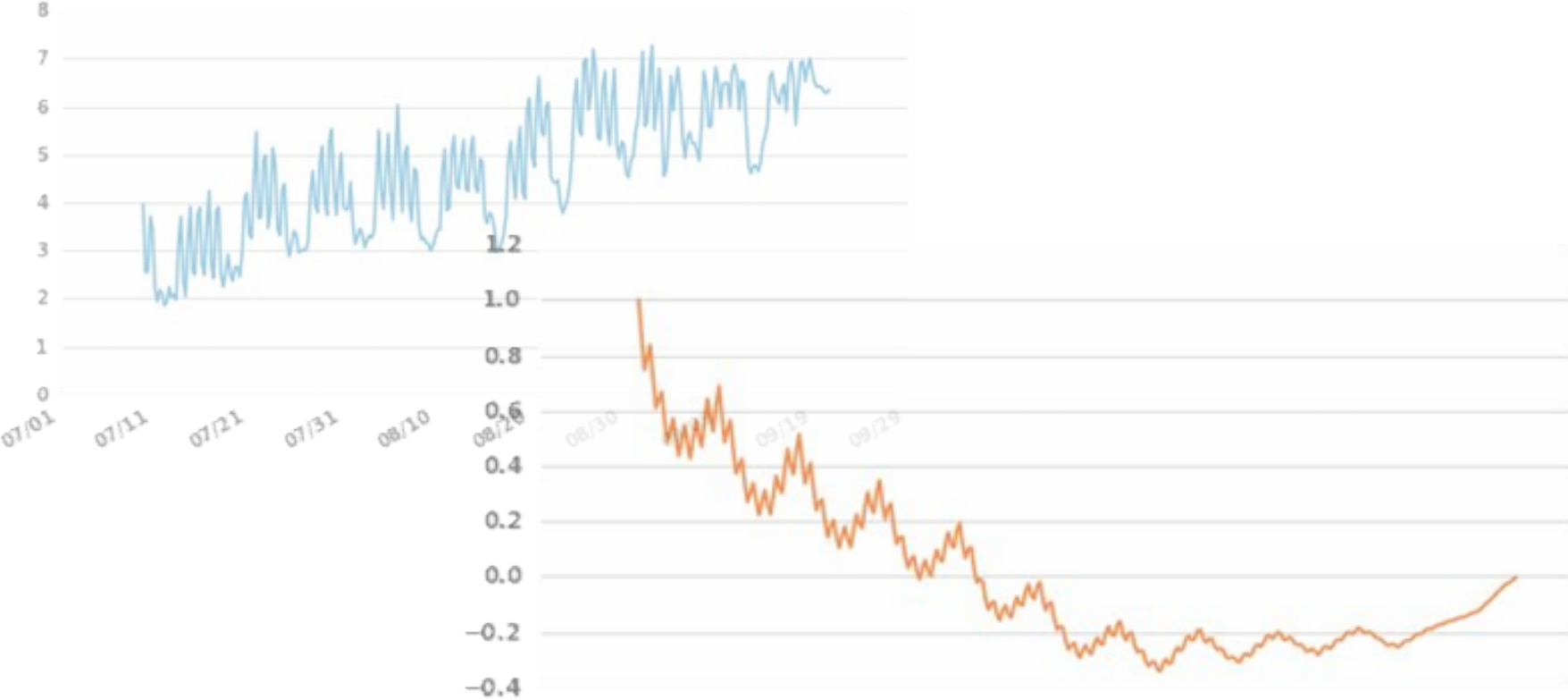
Decomposition



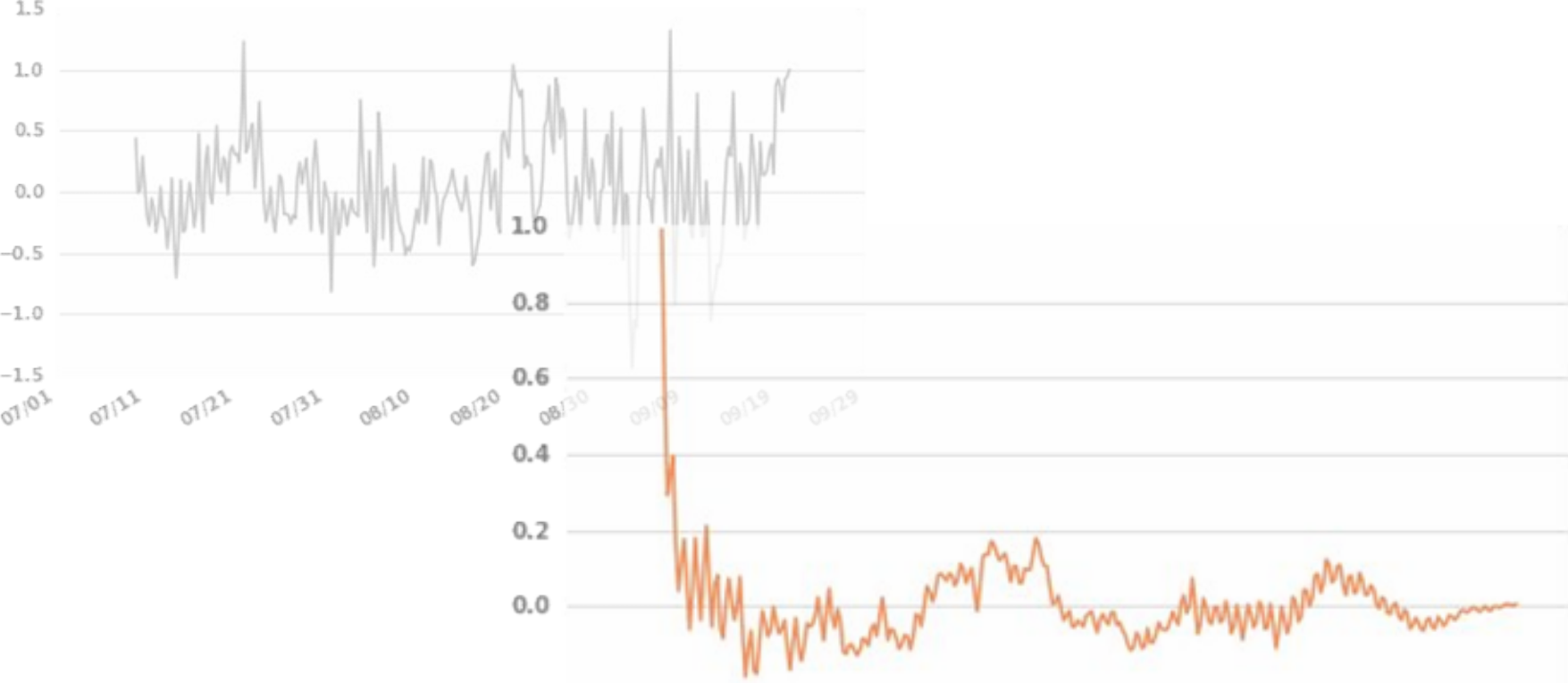
Autocorrelation



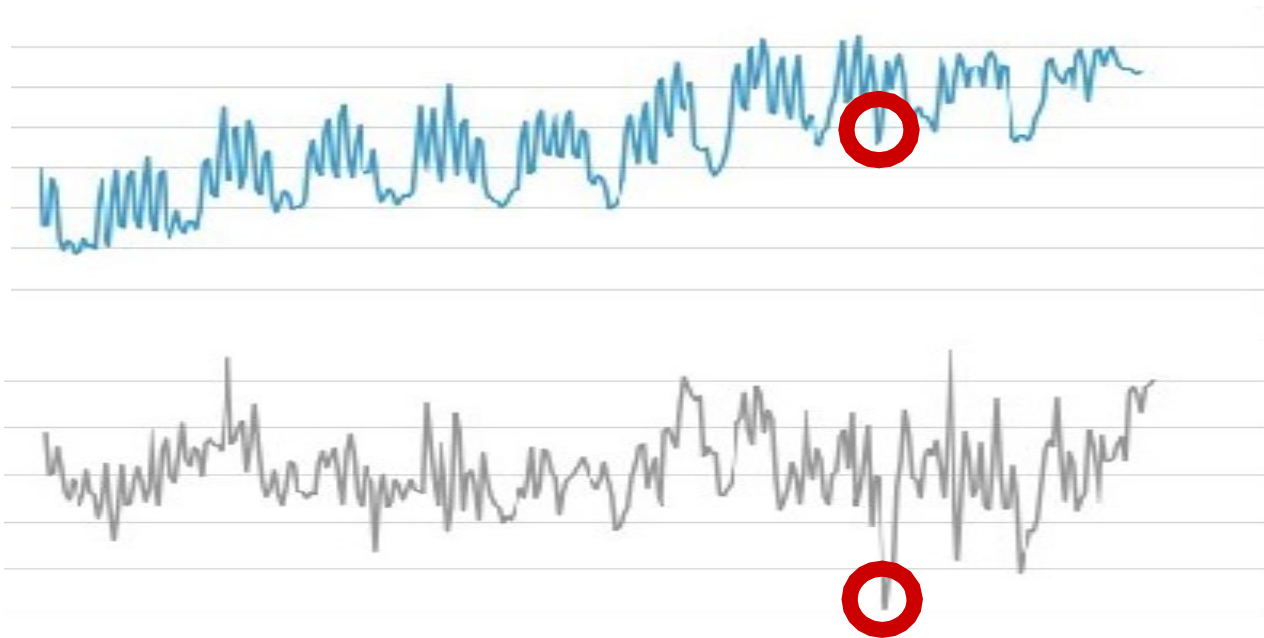
Signal vs. Noise



Signal vs. Noise



Signal vs. Noise vs. Signal



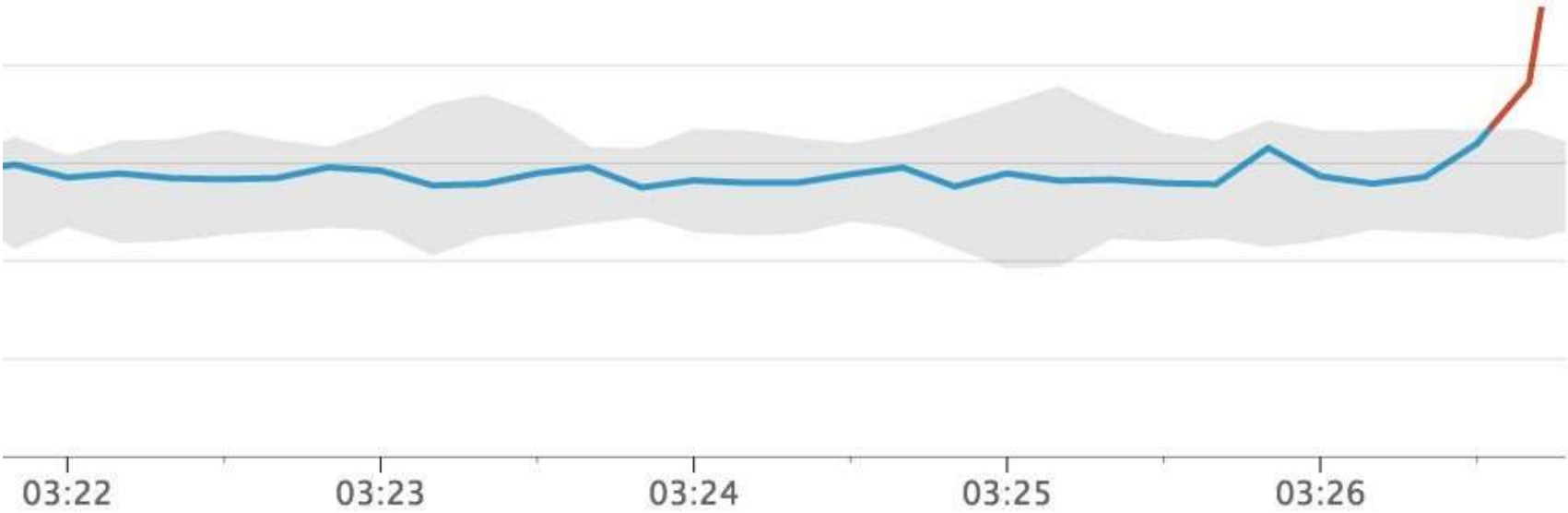
Real-time Anomaly Detection



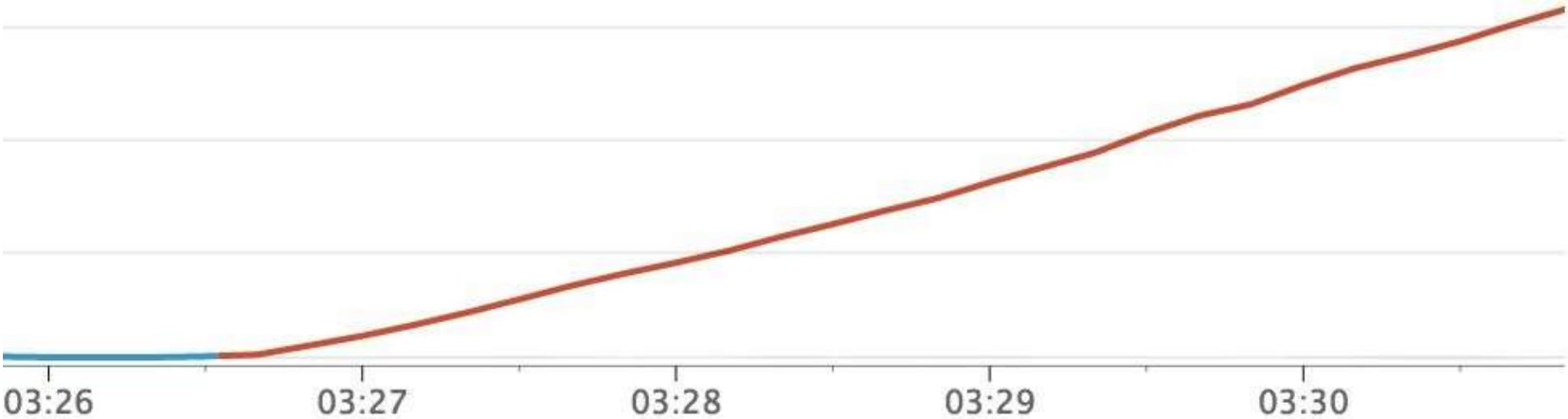
Anomaly Detection



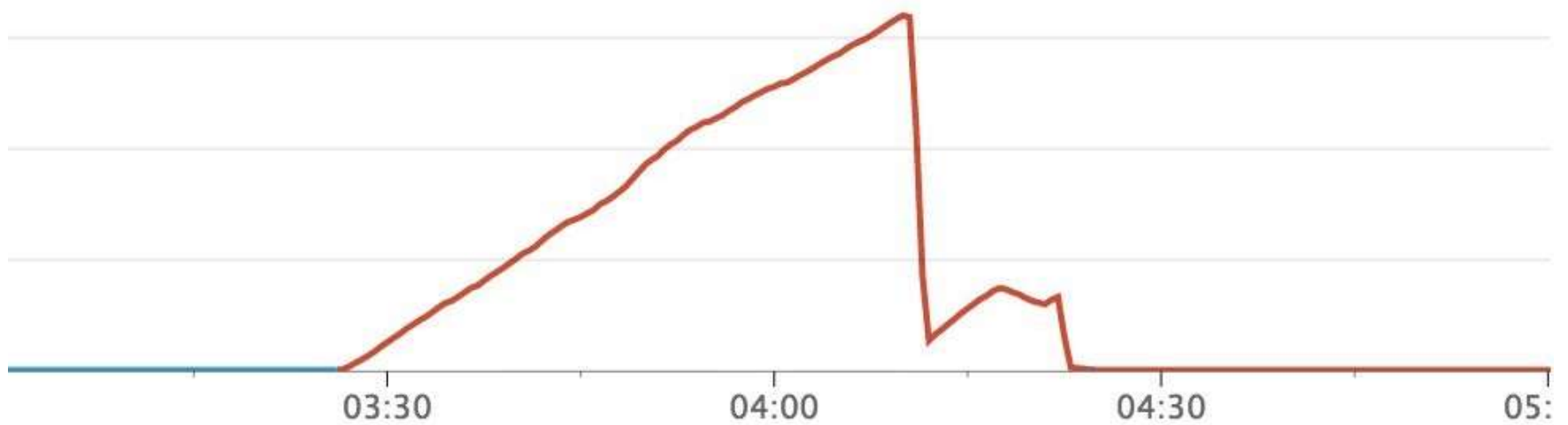
Robust Anomaly Detection



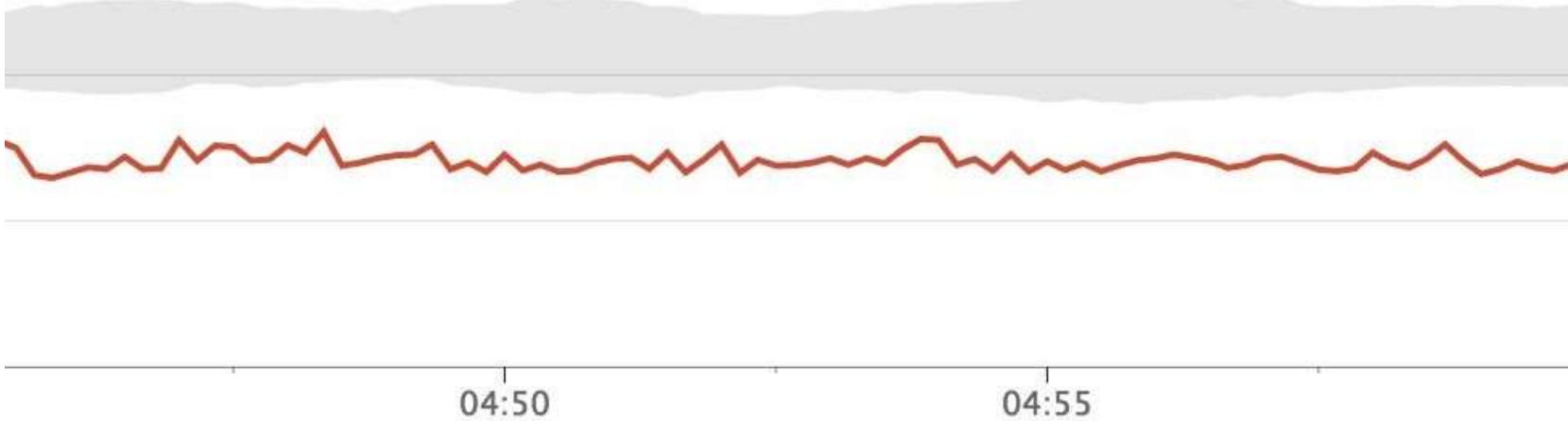
Robust Anomaly Detection



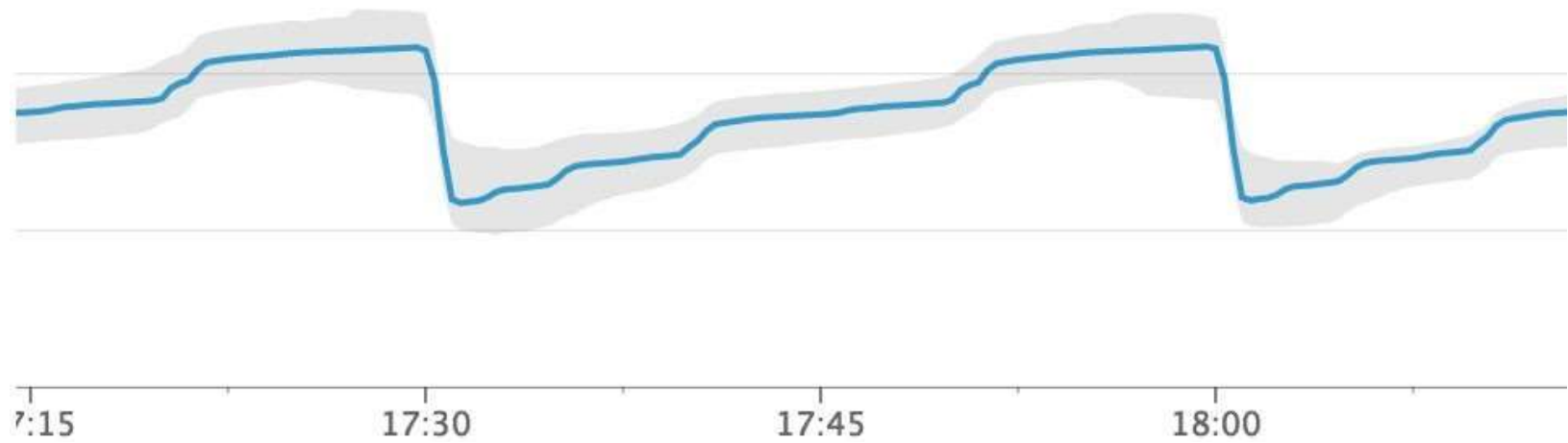
Robust Anomaly Detection



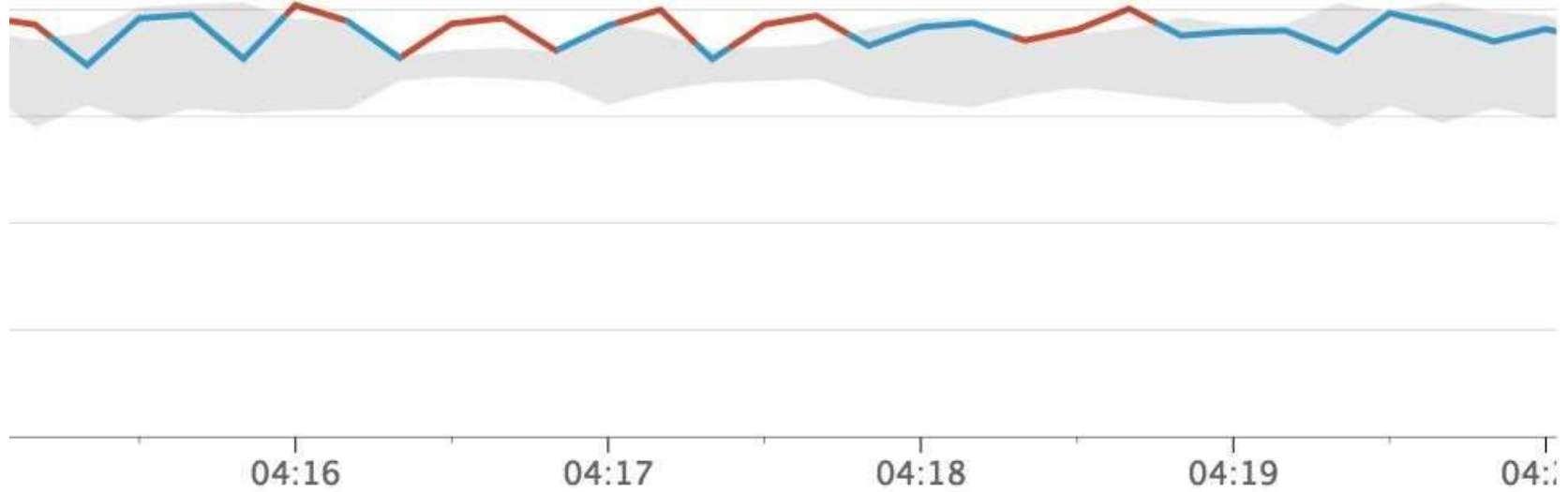
Robust Anomaly Detection



Alerting



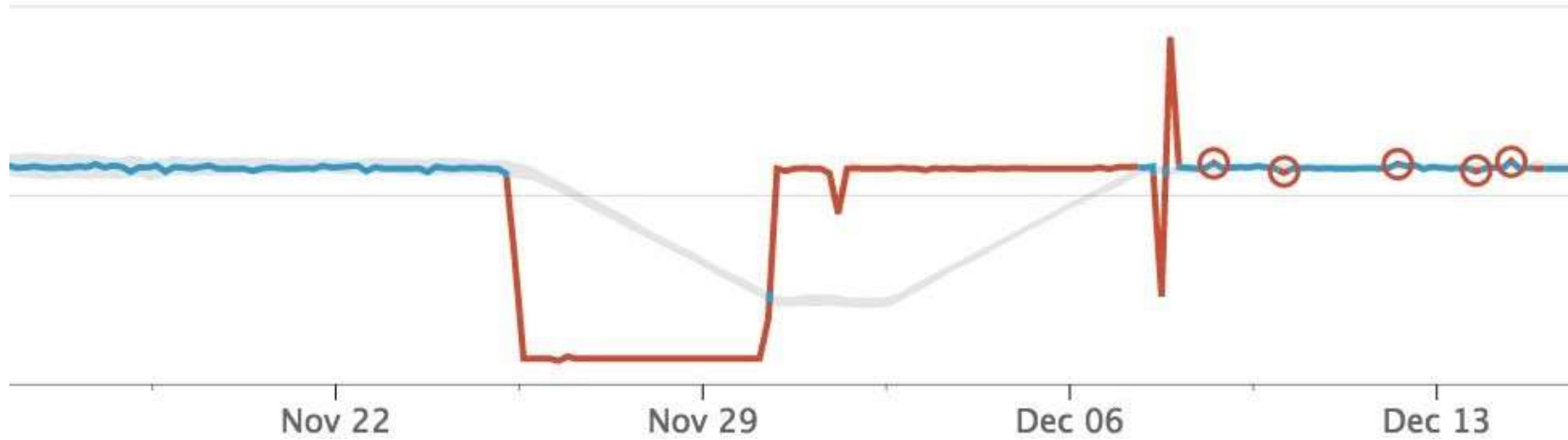
Alerting



Recap

- Extract as much signal as you can.
- Use robust statistical measures.
- Alert judiciously.
- Don't over-optimize.

Anomalies or Noise?



Thanks!



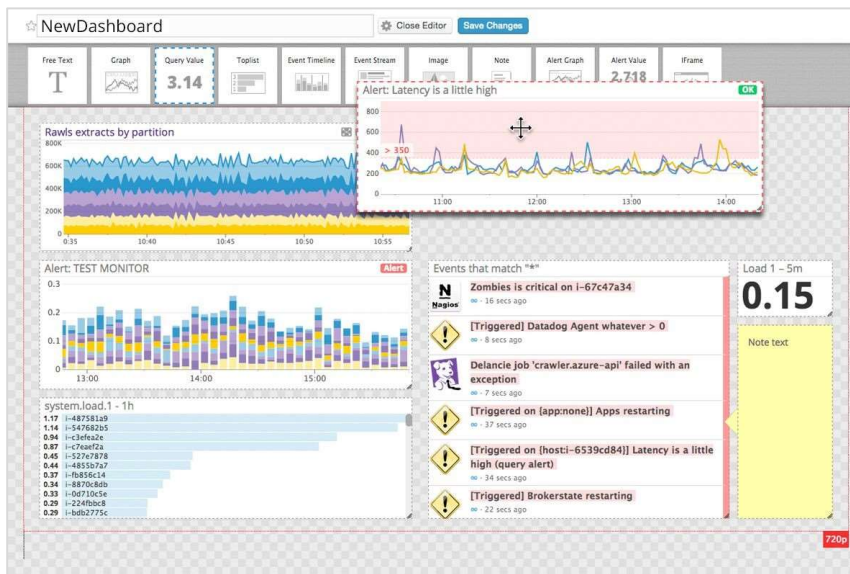
Appendix

See It All In One Place

Your Servers, Your Clouds, Your Metrics, Your Apps, Your team. Together.

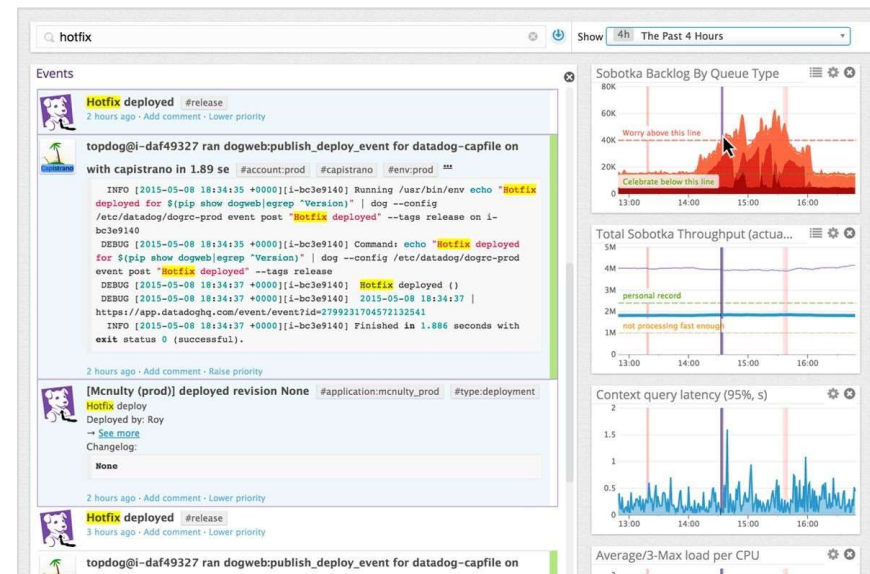
DASHBOARDS

Build Real-Time Interactive Dashboards



CORRELATION

Search And Correlate Metrics And Events

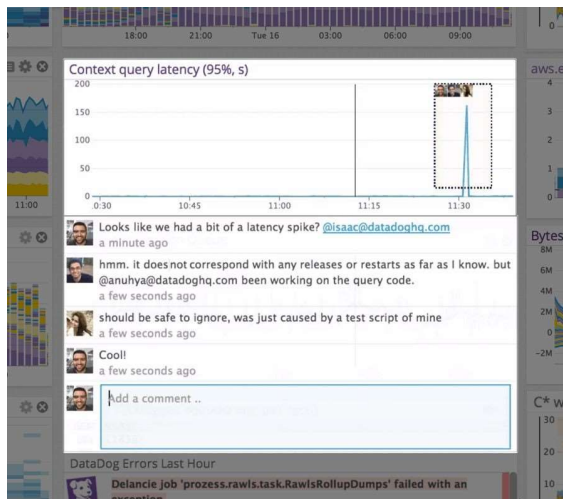


See It All In One Place

Your Servers, Your Clouds, Your Metrics, Your Apps, Your team. Together.

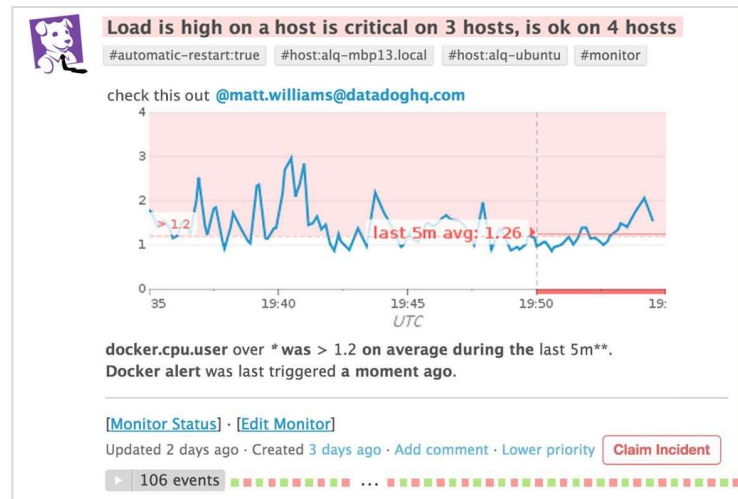
COLLABORATION

Share What You Saw, Write What You Did



METRIC ALERTS

Get Alerted On Critical Issues



DEVELOPER API

Instrument Your Apps, Write New Integrations

```
worker.py
18 def init_worker():
19     options = {
20         'api_key': DD_API_KEY,
21         'app_key': DD_APP_KEY,
22     }
23     initialize(**options)
24
25
26 def process_job(job):
27     # measure job pending time
28     q_time = job['meta']['queued_time']
29     api.Metric.send(
30         'job.in_flight_time',
31         time.time() - q_time,
32         tags=job['meta']['tags']
33     )
34
35
36 try:
37     s_time = time.time()
38     job_to_klass(job).execute()
39     api.Metric.send(
40         'job.execution_time',
41         time.time() - s_time,
42         tags=job['meta']['tags']
43     )
44     log.exception("Job %s succeeded", job.t
```

Flexible Pricing

To Match Your Dynamic Infrastructure.

Free

Up to 5 Hosts

1 Day retention

Custom metrics and events

Discussion group supported

Pro

Up to 500 Hosts

\$15 Per Host / Month

13 Month retention

Custom metrics and events

Metric alerts*

Email supported

Enterprise

500+ Hosts

Contact us for pricing:
+1 866 329 4466
sales@datadoghq.com

Customized retention

Custom metrics and events

Metric alerts*

Email and phone supported