

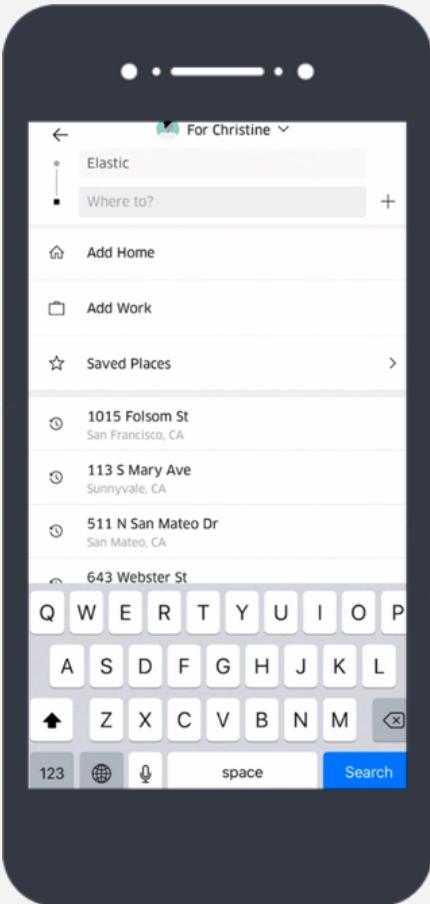


Machine Learning using Elasticsearch on Azure

Ravi Ramnani | Solutions Architect

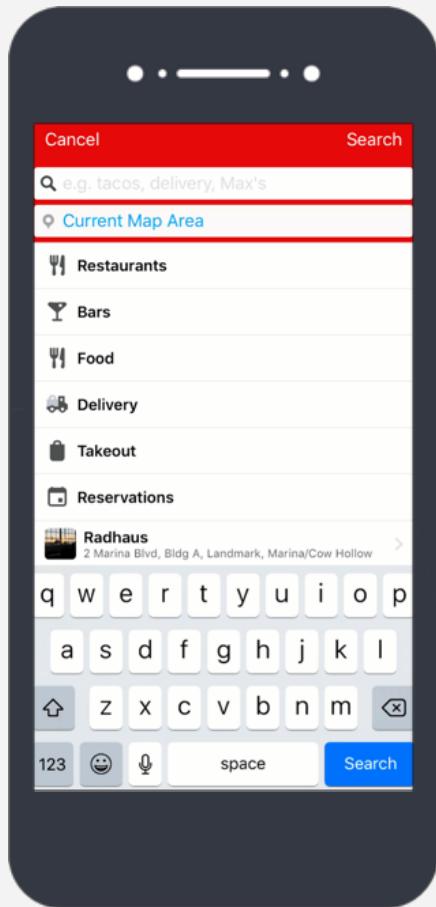
15th Dec, 2018

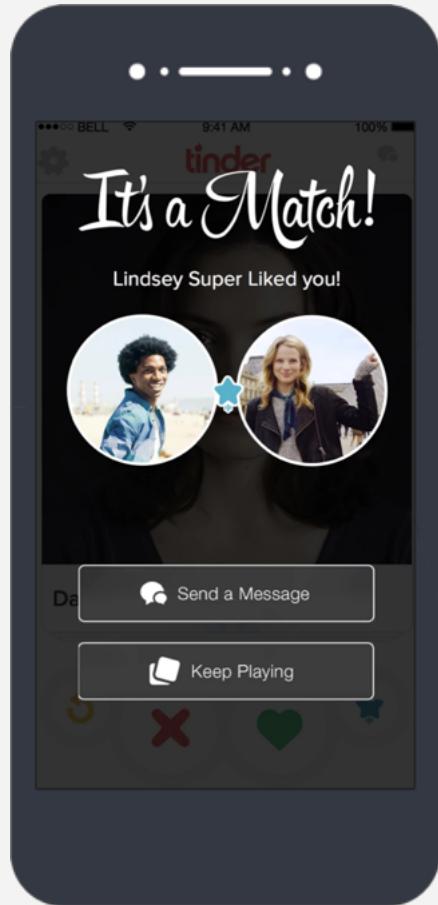
Elastic is a **search company**



Uber



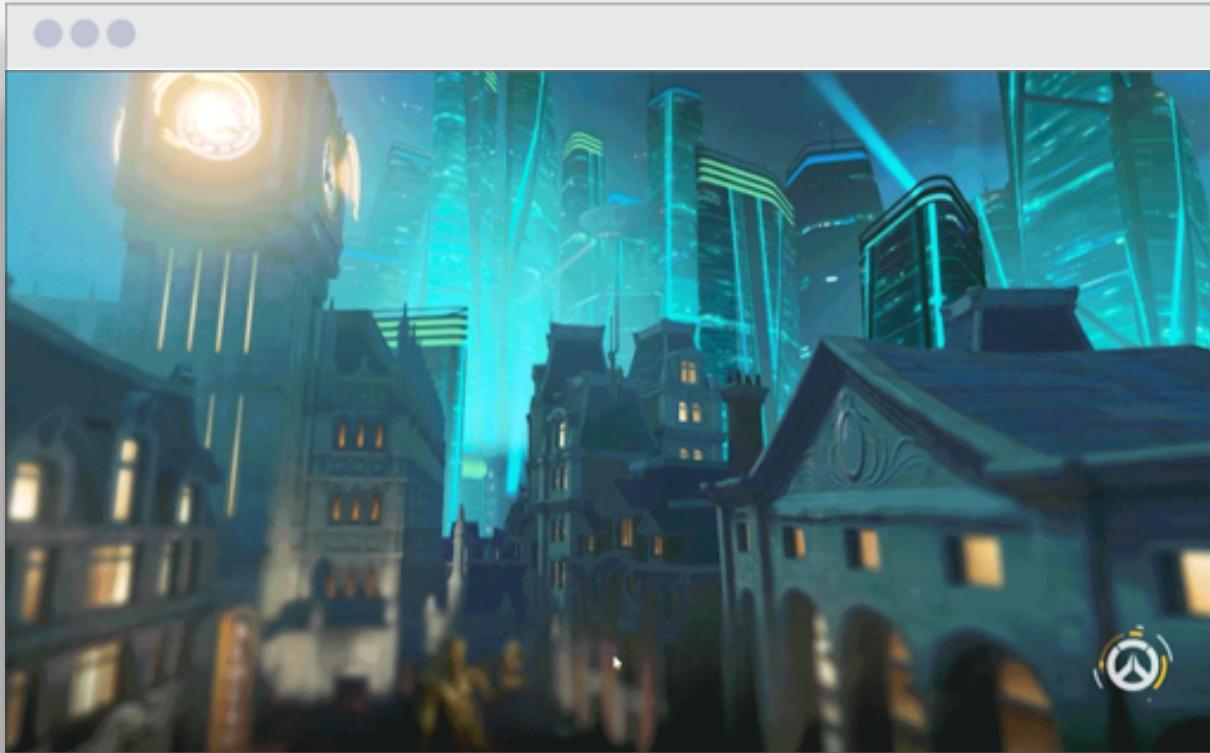




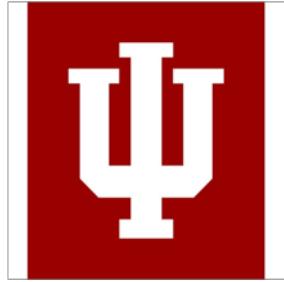
tinderTM



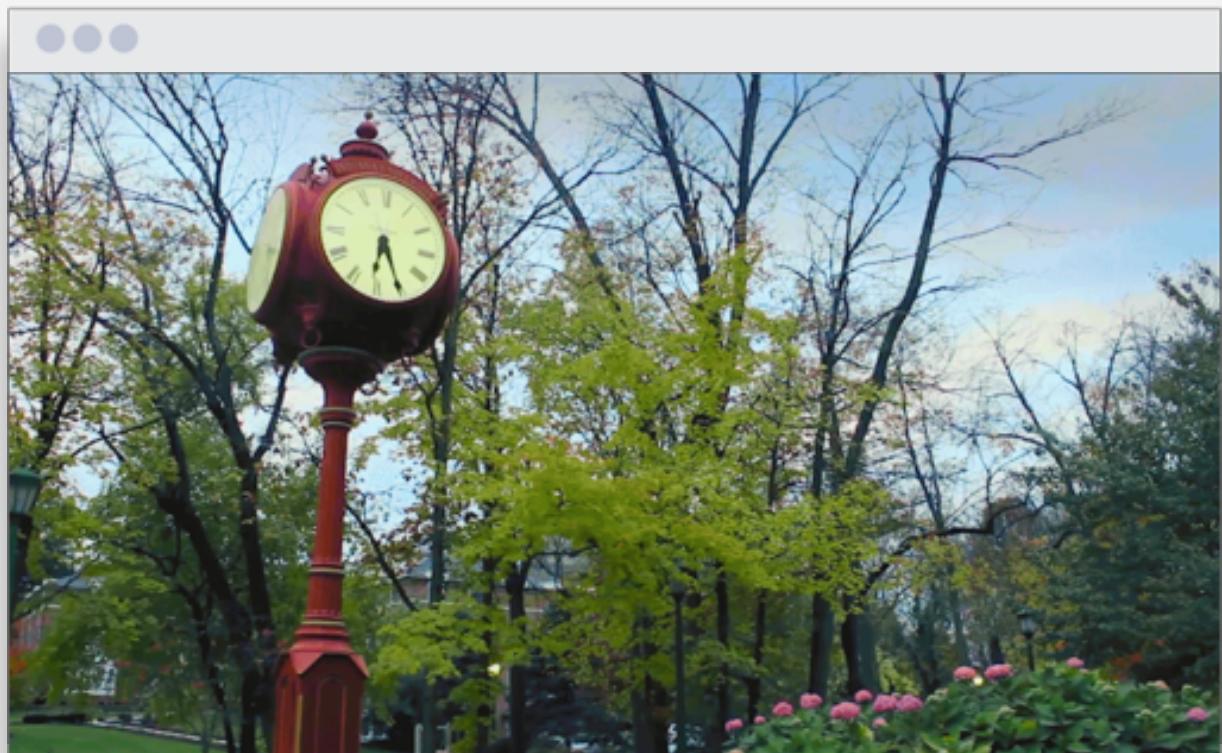
A screenshot of the Instacart website interface. At the top, there's a banner for Whole Foods Market showing a wooden cutting board with sliced bread. The Instacart logo is in the top left, and account and help links are in the top right. A shopping cart icon with a '4' is also present. Below the banner, there's a search bar and navigation links for Home, Departments, Coupons, Get \$50, and Your Items. To the right, delivery details show "Delivery to: 94123" and "Within 1 hour". There are three promotional boxes: one for "Coupon savings" (up to 40% off everyday essentials) with a "Shop Coupons" button; one for "Free Delivery" with select Holiday Baking & Meals items, featuring a "Save Now" button and a Mrs. Fields cookie jar icon; and one for "Free Delivery" with select Tropicana, IZZE, Naked, and more items, also with a "Save Now" button and a Tropicana logo. Below these is a "Trending Near You" section with images of a red bell pepper, green beans, raw meat, a kiwi, and a carton of Met & Fit yogurt, each with a plus sign icon.



ACTIVISION
BLIZZARD



INDIANA UNIVERSITY



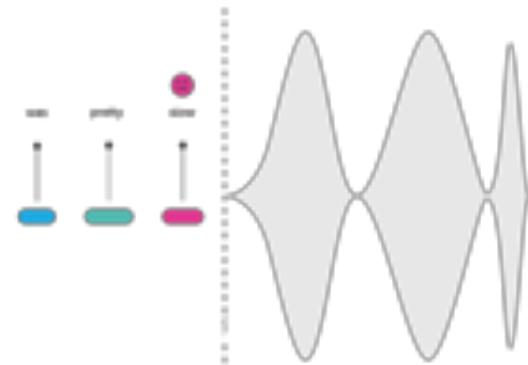
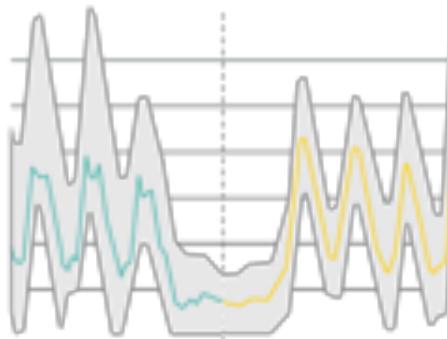
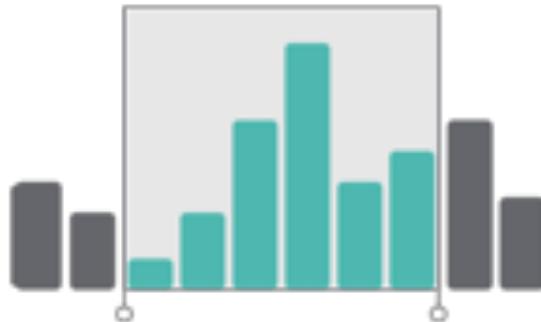
Customers across various industries, segments, and geographies

TECHNOLOGY	FINANCE	TELCO	CONSUMER	HEALTHCARE	PUBLIC SECTOR	AUTOMOTIVE / TRANSPORTATION	RETAIL
							
							
							
							
							

Search is a **constant/foundation**



.54 seconds | 1,000,000,000 records



Technology **differentiation**



SCALE

Distributed by design

SPEED

Find matches in milliseconds

RELEVANCE

Get highly relevant results

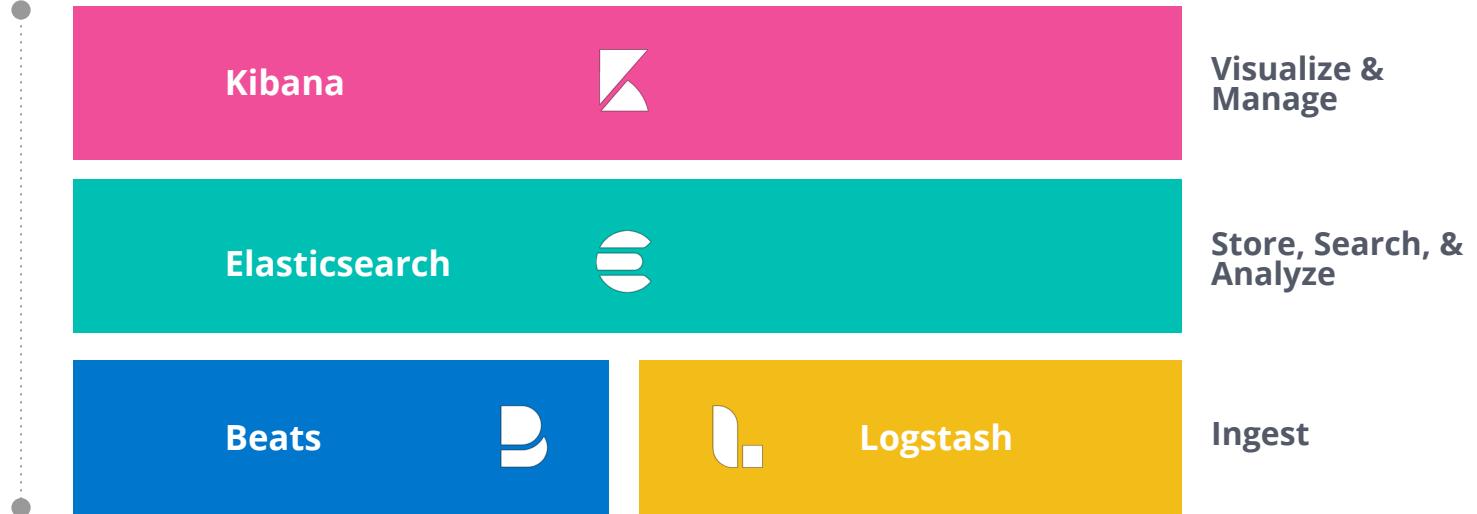
Elastic Products

Elastic Stack

SOLUTIONS



Elastic Stack



SaaS



Elastic cloud

SELF-MANAGED



Elastic cloud
Enterprise

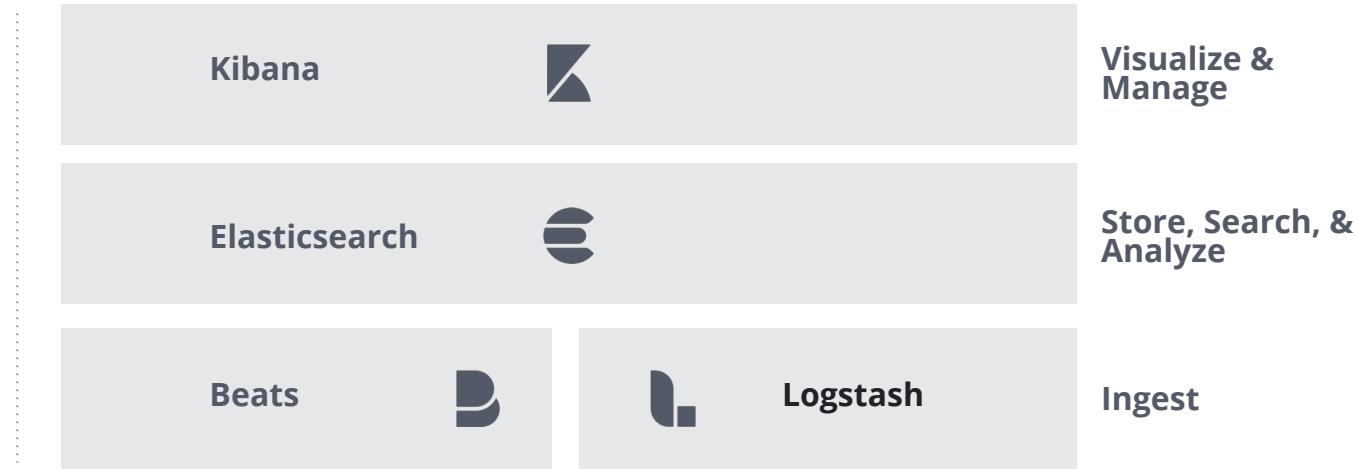


Standalone

Solutions



Elastic Stack



SaaS

SELF-MANAGED

AI and ML

AI and ML



Algorithmic Approaches

Supervised

- Labelled Data
- Driven by objective
- NNs, SVMs, Decision Trees

Unsupervised

- No Label
- No objective, only Data
- K-means, Apriori

Reinforcement

- No Label
- Driven by objective
- Q-Learning, SARSA, DQL

Anomaly Detection Concepts

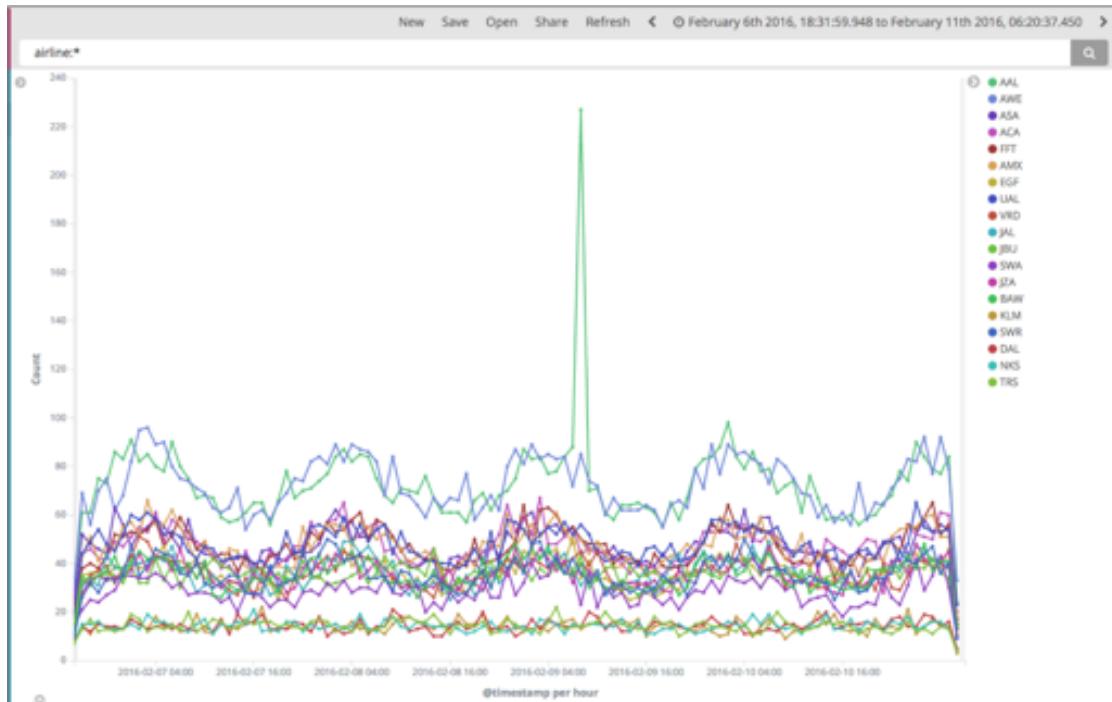
Terminology

- **Machine Learning**
 - Broad term, but X-Pack Machine Learning is automated anomaly detection for time-series data (for now).
- **Anomaly Detection**
 - Discovery of what's "weird" or "different", not what's "bad"
- **Unsupervised Learning**
 - Learning without human-labeled examples (without being "taught"). Rely only on the data
- **Bayesian**
 - An approach based on probability in which prior results are used to calculate probabilities of certain present or future events

What is “Abnormal”?

What's abnormal here?

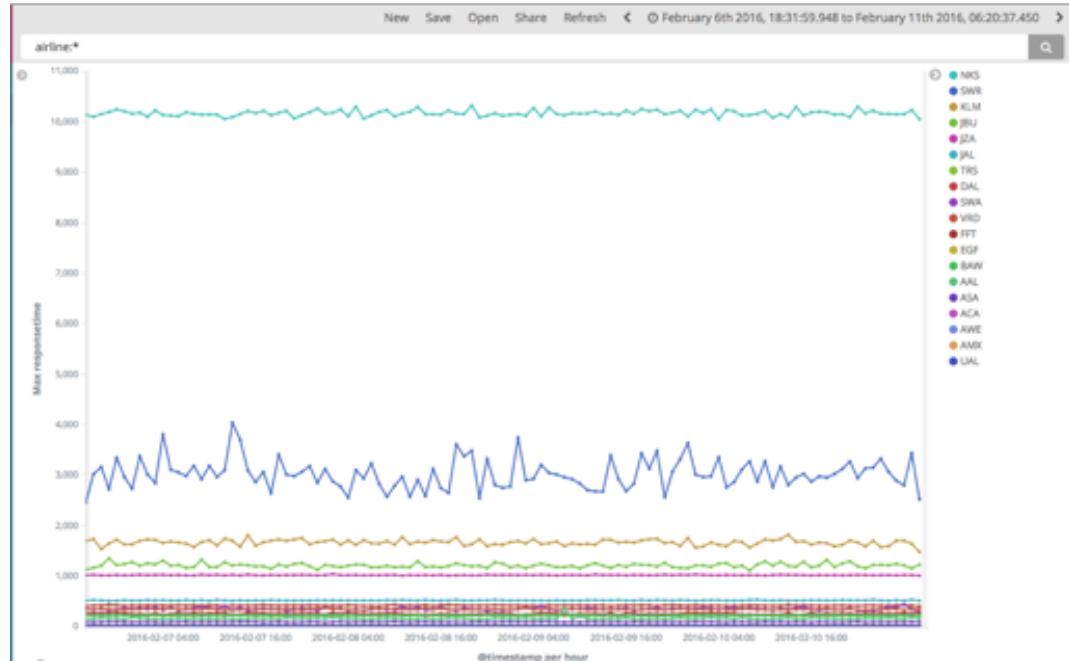
Why?



What is “Abnormal”?

What's abnormal here?

Why?



What is “Abnormal”?

What's abnormal here?

Why?

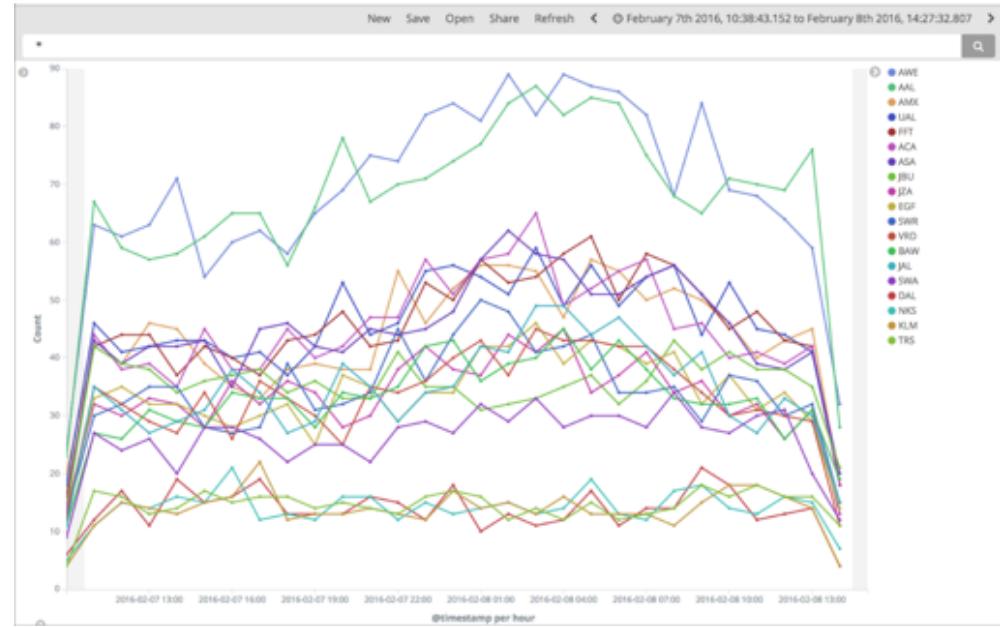


What is “Normal”?

In general, this question can be answered in two ways:

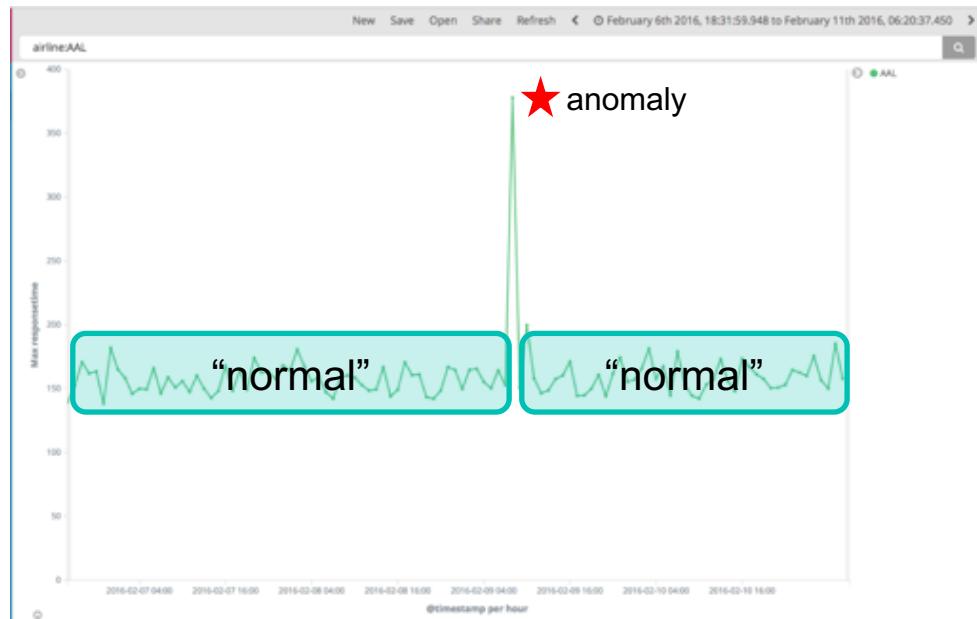
1) Something behaves in a consistent way with respect to itself, over time

2) Something behaves in a consistent way compared against similar entities



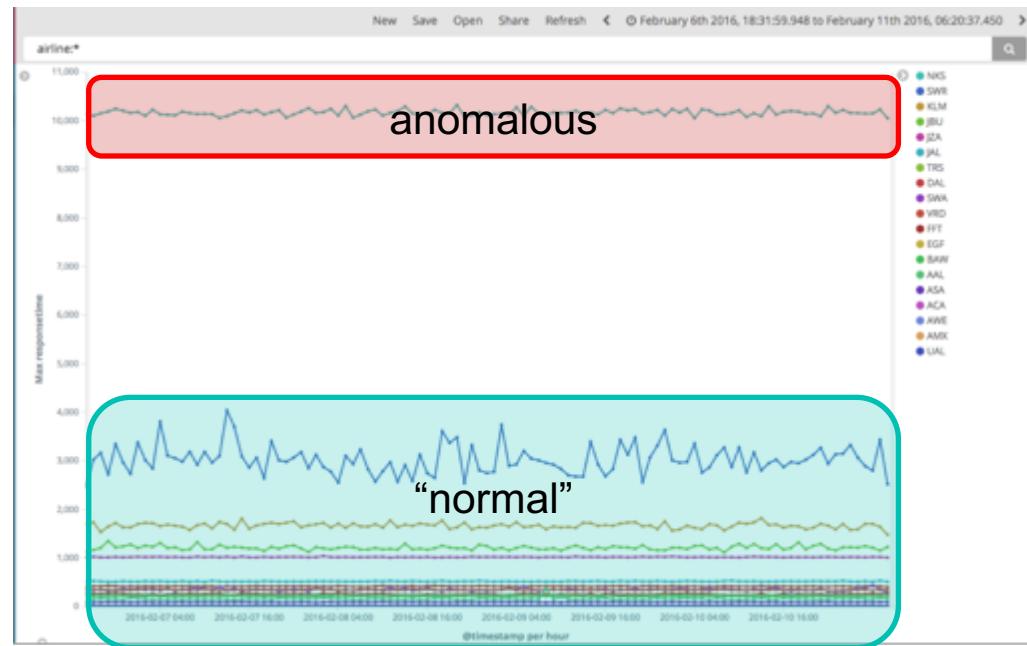
What is “Abnormal”?

1) If something changes its behavior, compared to its own history – that change is *anomalous*.



What is “Abnormal”?

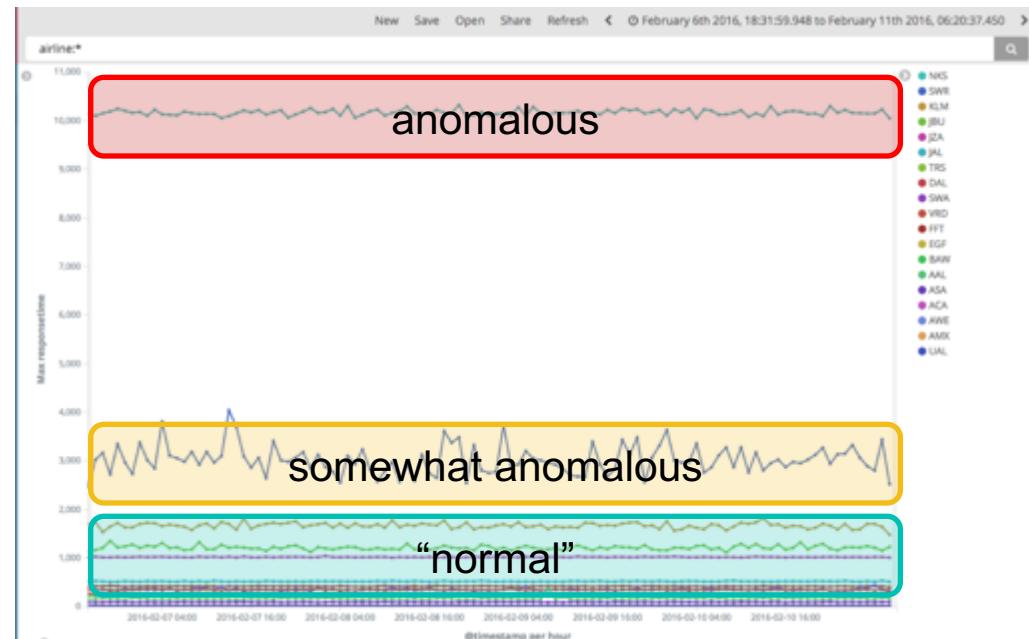
2) If something is drastically different than others within a population, then that entity is *anomalous*.



What is “Abnormal”?

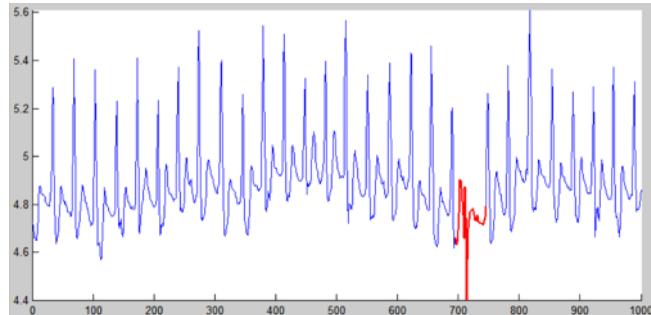
2) If something is drastically different than others within a population, then that entity is *anomalous*.

There's also the concept of being “somewhat anomalous”



In Summary, Anomalousness is:

1) When an entities' ***behavior changes*** significantly and suddenly



2) When an entity is drastically ***different than others*** within a population



How to Learn “Normal”

An Analogy

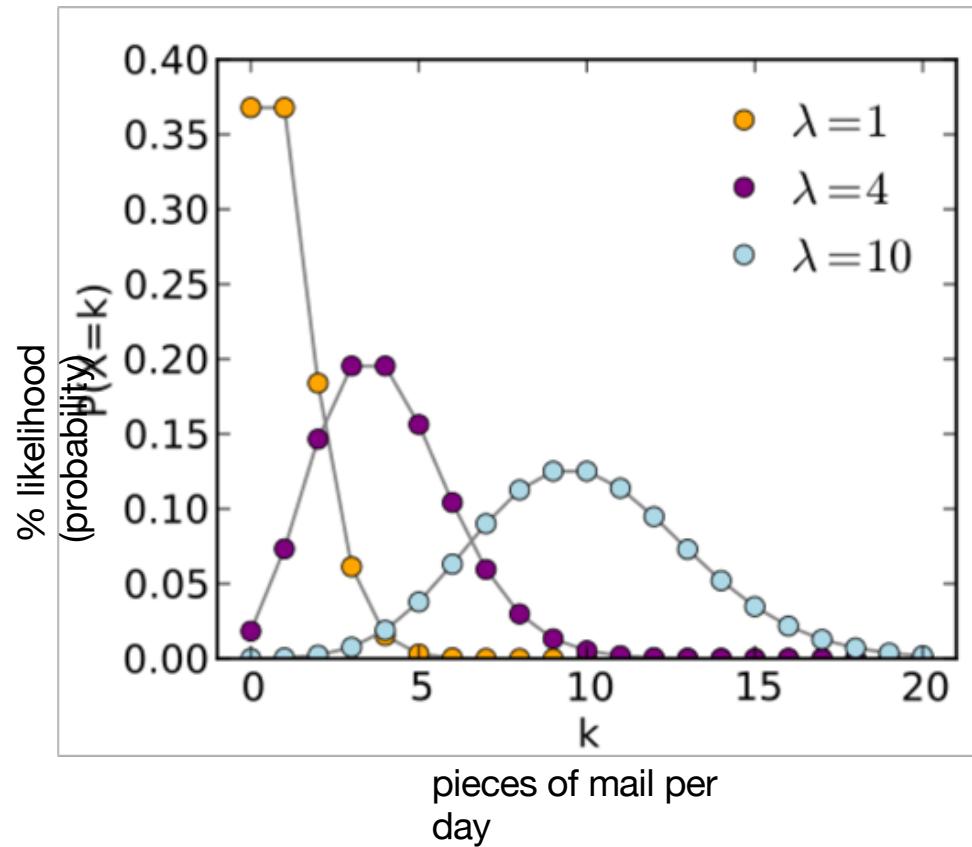
- How could I learn how much Postal-mail you get daily and how could I use that information to predict how much you might get tomorrow?



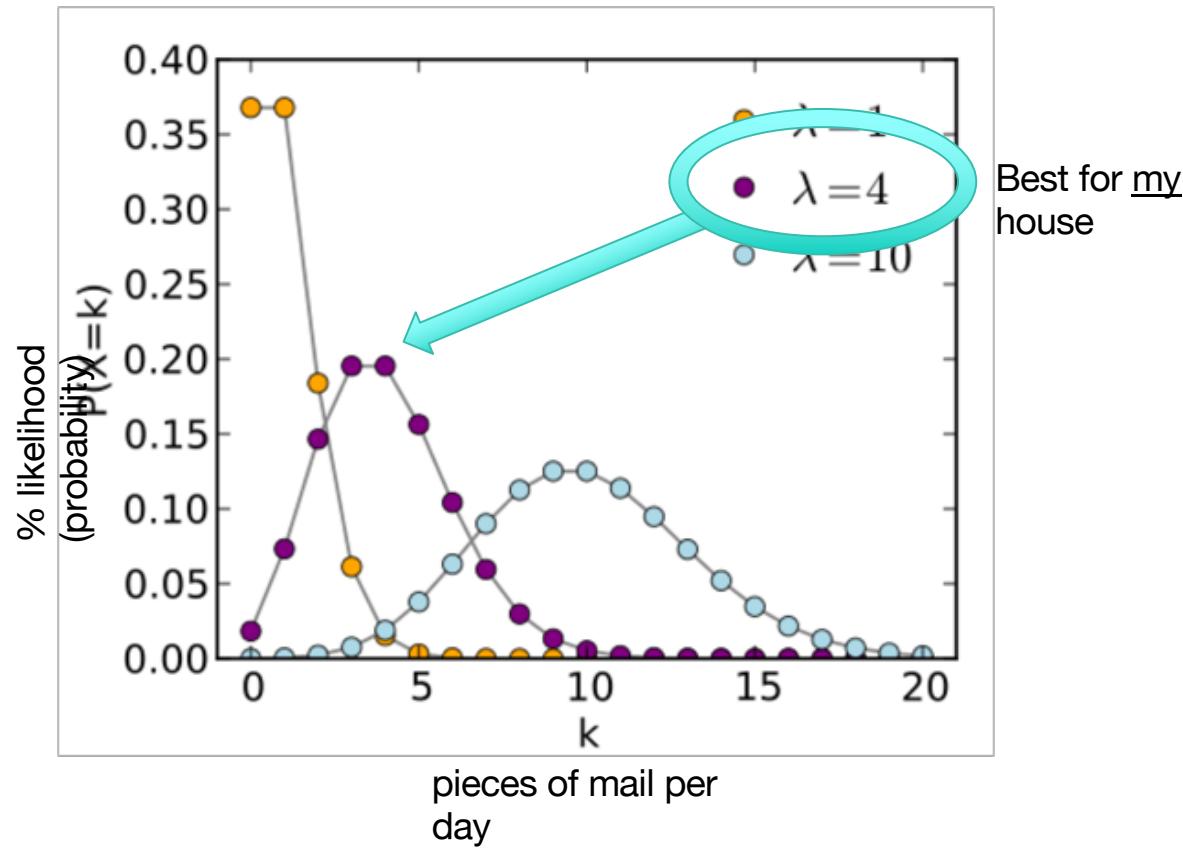
I could stand on your front porch and...

- Watch your mail delivery volume for a while, but how long?
 - 1 day?
 - 1 week?
 - 1 month?
- Notice, that you intuitively feel like you'll gain accuracy in your predictions with more data that you see.
- Use those observations to create a model...

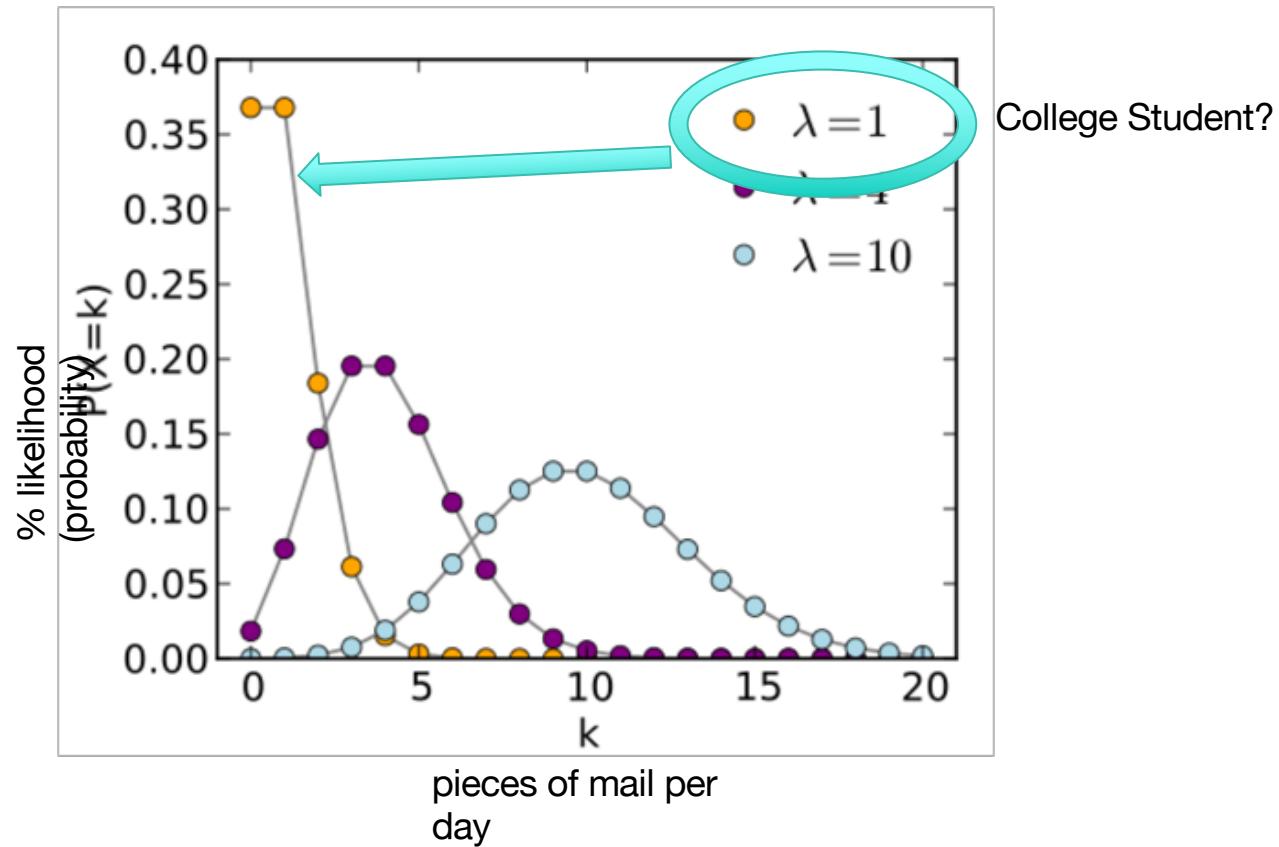
Probability Distribution Function



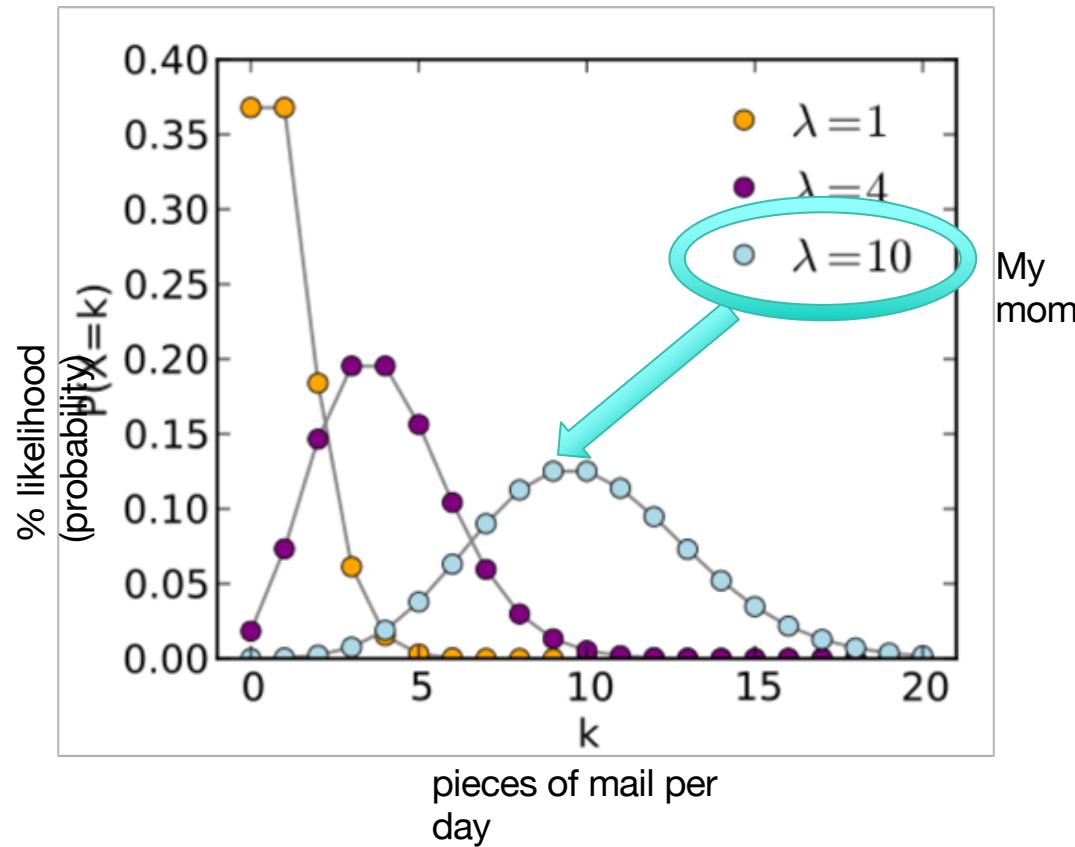
Probability Distribution Function



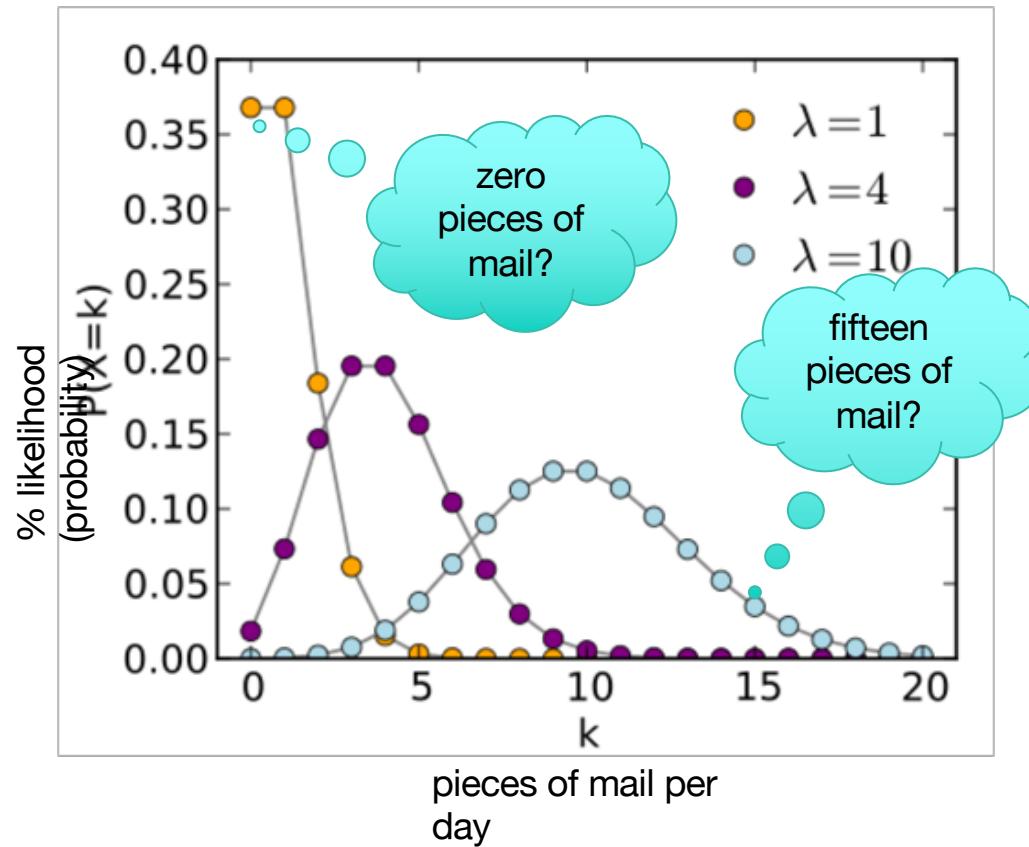
Probability Distribution Function



Probability Distribution Function



Using the Model to Find What is Unexpected



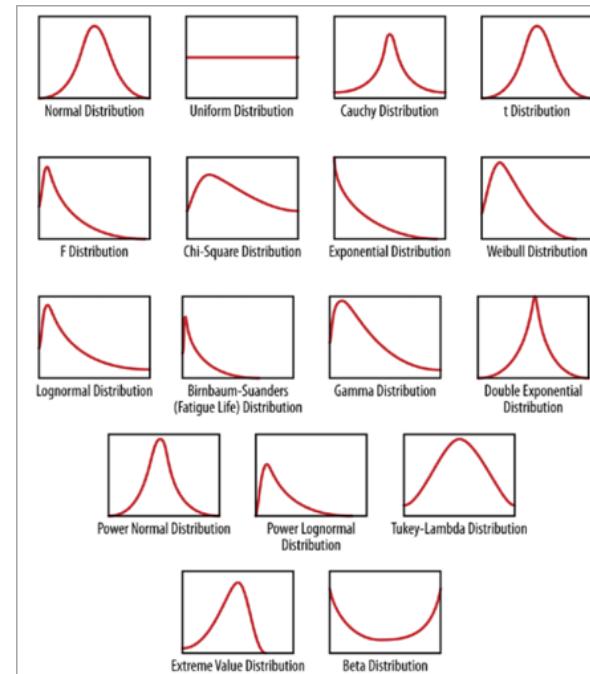
Relate back to IT/Security data

- # Pieces of mail = # events of a certain type
 - number of failed logins
 - number of errors in a log
 - number of events with certain status codes
- Or, metrics
 - response time
 - utilization
 - orders per minute

=> Every kind of data will need its own unique “model” (probability distribution function)

How does one pick a model?

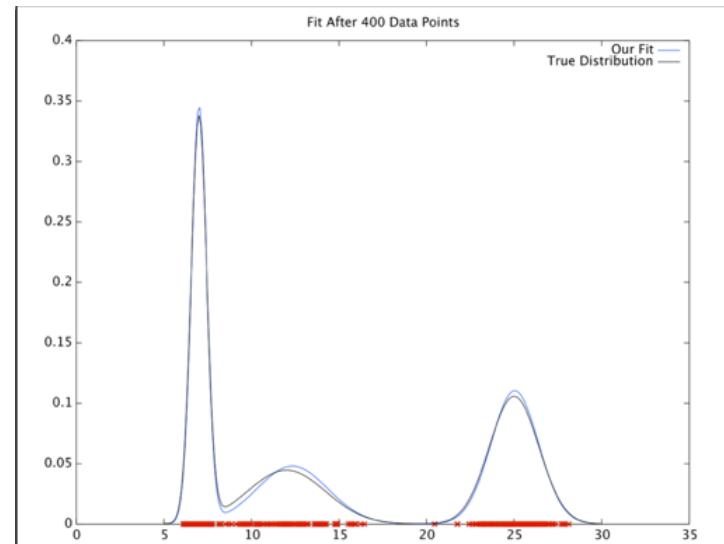
Which one best fits
your data?



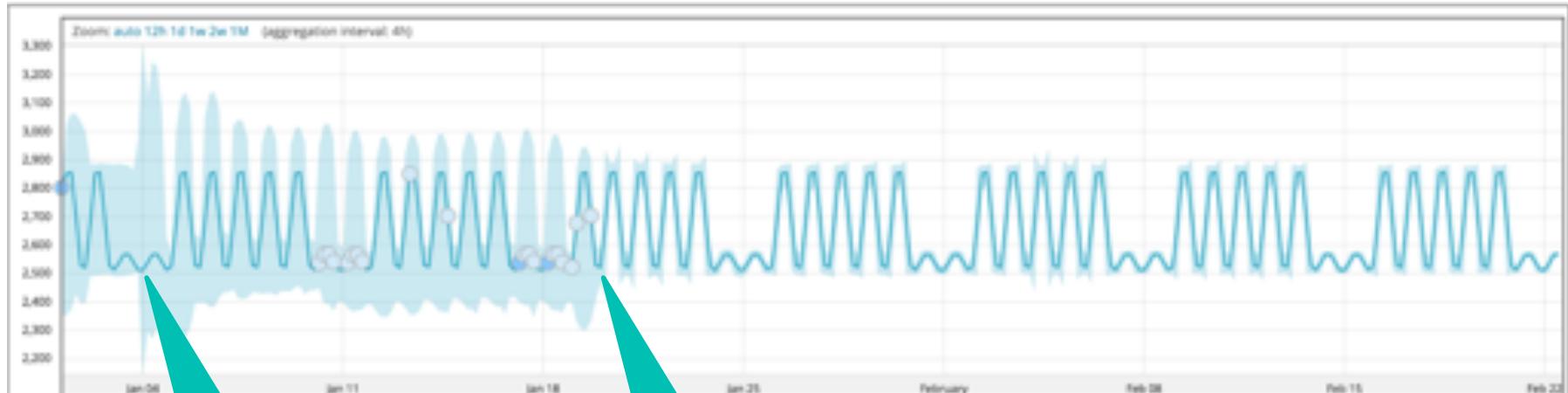
source: "Doing Data Science"
O'Neil & Schutt

Machine Learning picks it for you

- ML uses sophisticated machine-learning techniques to best-fit the right statistical model for your data.
- Better models = better outlier detection = less false alarms
- Anomalies occur when observation is in low probability area



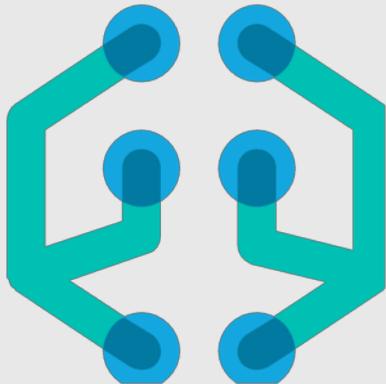
The Model's Evolution in Time



After 2 full days,
daily periodicity
has been learned.

After 2 full weeks,
weekly periodicity
has been learned.

Elastic Machine Learning



- Online Unsupervised Learning
- Index Visualizer
- Anomaly Detection
- Forecasting

Anomaly Detection

- Single or Multi-metric time series
- Outliers in population
- Rare events categorization

What you can do?

Operational Analytics

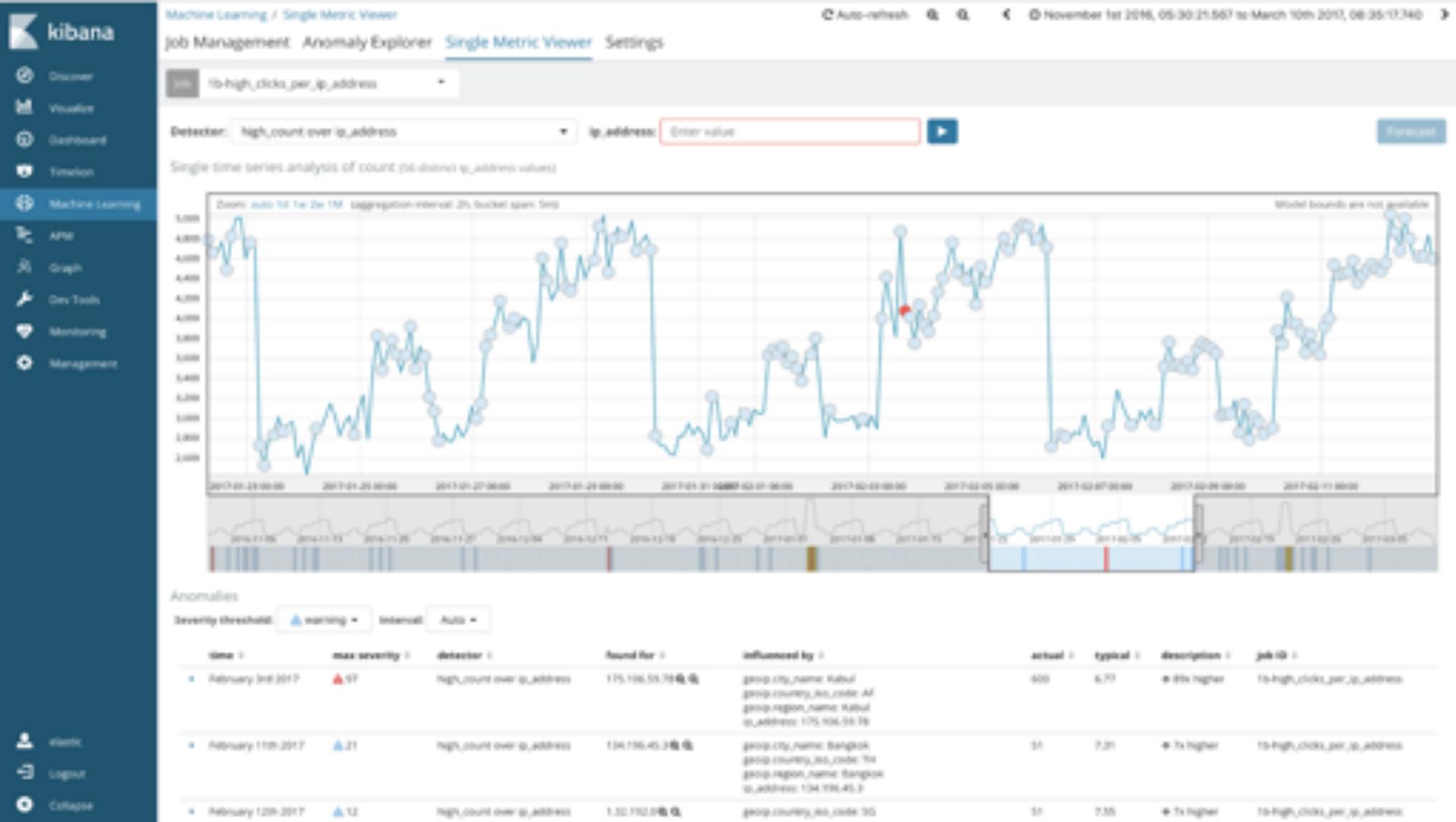
- Drop in Orders?
- Unusual Traffic
- Early hints before things go wrong

Security Analytics

- MITM detection
- Malware? Insider Threat?
- Maliciously running processes

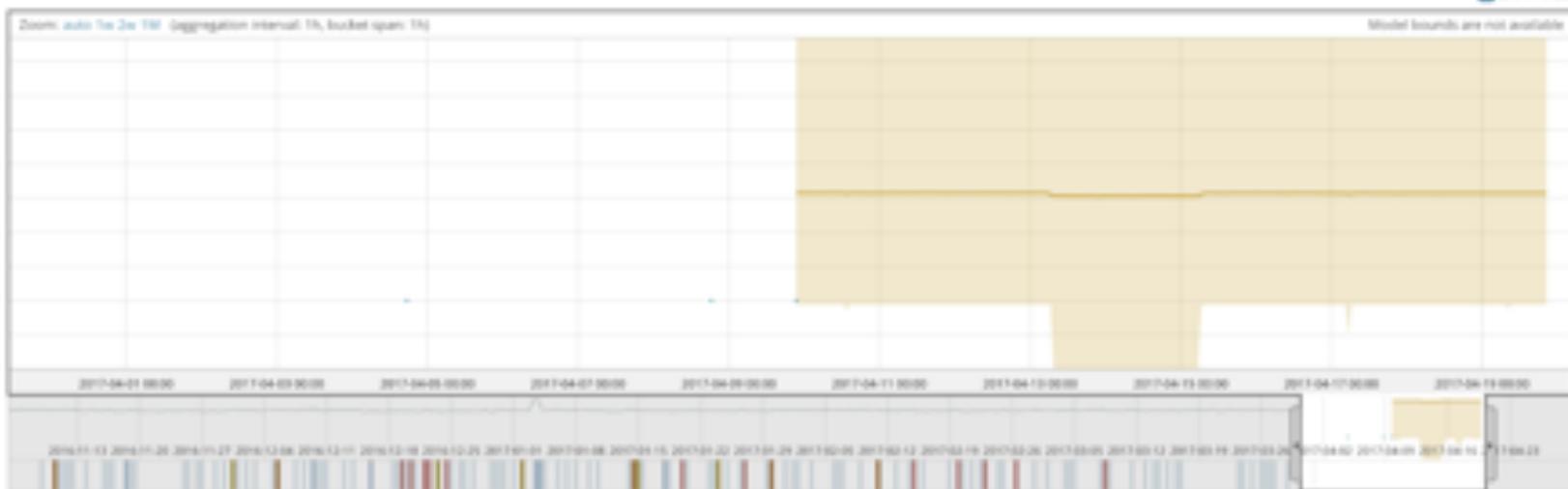
Business Analytics

- Latency in response times
- Low click through rate on ads



Job **Za-high_purchase_quantity**Detector: **high_meantotal_quantity/high_purchase_quantity** (or) user:

Single time series analysis of avg total_quantity (52 distinct user values)

 show forecast

Anomalies

Severity threshold:

No matching results found

security-analytics-winlogbeat-*

Use full security-analytics-winlogbeat-* data

Search... (e.g. status:200 AND extension:json)

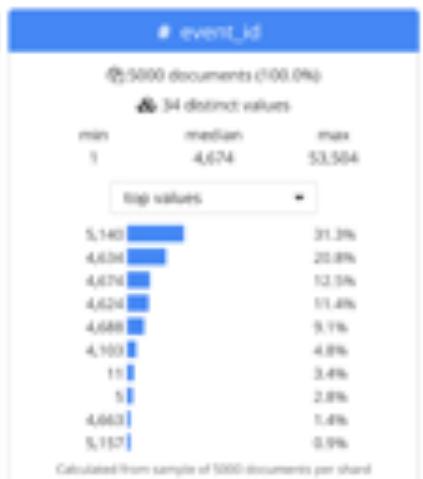


Sample 5000 → documents per shard from a total of 46293 documents

Metrics

9 fields exist in documents (29 in total) show empty fields

Filter



Create job

Use the Advanced job wizard to create a job to find anomalies in this data:



Advanced

Use the full range of options to create a job for more advanced use cases.



Microsoft preferred solution. [Learn more](#)

Elastic's addition to the Microsoft Azure Marketplace means developers can utilize preconfigured templates built by Elastic to more easily and quickly deploy an Elasticsearch cluster on Azure.

Elastic is the company behind the Elastic Stack, a suite of products that include Elasticsearch, Logstash, Kibana and Beats, which are focused on scalability and ease-of-use to help you make sense of your data.

This is a Bring-Your-Own-License (BYOL) solution template. The BYOL model gives users the option to add critical commercial enhancements such as

- cluster and data security
- user authentication
- stack monitoring
- alerting and notifications
- graph
- machine learning capabilities

through an Elastic subscription purchased directly from Elastic. This subscription also provides global, Enterprise class support for Elasticsearch on Azure. You can learn more about Elastic Stack features at www.elastic.co.

The solution template installs with a 30 day trial license for Elastic platinum features. For information about an Elastic subscription, contact azuremarketplace@elastic.co. Be sure to also check out our [documentation](#) to learn more about this offering.

Elasticsearch cluster with up to 50 data nodes (assuming you have a sufficiently large VM quota), with optional Kibana and Logstash. Elastic Stack platinum features like monitoring, security, alerting, graph and machine learning are automatically activated as part of a 30 day trial license.

Save for later

Create



kibana

- Discover
- Visualize
- Dashboard
- Timeline
- Canvas
- Machine Learning
- Infrastructure
- Logs
- Graph
- Dev Tools
- Monitoring
- Management

- profile user
- Logout
- Default
- Collaborate

Add Data to Kibana

Use these solutions to quickly turn your data into pre-built dashboards and monitoring systems.



Logging

Ingest logs from popular data sources and easily visualize in preconfigured dashboards.

[Add log data](#)

Metrics

Collect metrics from the operating system and services running on your servers.

[Add metric data](#)

Security analytics

Centralize security events for interactive investigation in ready-to-go visualizations.

[Add security events](#)

[Add sample data](#)
Load a data set and a Kibana dashboard

[Upload data from log file](#)
Import a CSV, NDJSON, or log file

[Use Elasticsearch data](#)
Connect to your Elasticsearch index

Visualize and Explore Data



Canvas

Showcase your data in a pixel-perfect way.



Dashboard

Display and share a collection of visualizations and saved searches.



Discover

Interactively explore your data by querying and filtering.



Graph

Surface and analyze relevant relationships in your

Manage and Administer the Elastic Stack



Console

Skip UIs, and use this JSON interface to work with your data directly.



Index Patterns

Manage the index patterns that help retrieve your data from Elasticsearch.



Logstash Pipelines

Create, delete, update, and clone data ingestion



Monitoring

Track the real-time health and performance of your