**Forecasting Disk Usage with Machine Learning – So easy, even a cave person can do it!**

Steve Koelpin and Alicia Dale
TransUnion

splunk> .conf19

# Steve Koelpin

Advisor - Splunk
TransUnion

splunk> .conf19

# Forward-Looking Statements

////////////////////////////////

During the course of this presentation, we may make forward-looking statements regarding future events or plans of the company. We caution you that such statements reflect our current expectations and estimates based on factors currently known to us and that actual events or results may differ materially. The forward-looking statements made in the this presentation are being made as of the time and date of its live presentation. If reviewed after its live presentation, it may not contain current or accurate information. We do not assume any obligation to update any forward-looking statements made herein.

In addition, any information about our roadmap outlines our general product direction and is subject to change at any time without notice. It is for informational purposes only, and shall not be incorporated into any contract or other commitment. Splunk undertakes no obligation either to develop the features or functionalities described or to include any such feature or functionality in a future release.

splunk> .conf19

# Agenda

1. **Objective**

2. **Common Challenges**

3. **Dark Data**

4. **Technical Deep Dive**

5. **Build a Score Mechanism**

6. **Scaling this Out**

7. **Apply this to your data with walk-through**

splunk> .conf19

# Objective

In need of a solution to get **ahead** of capacity constraints

splunk> .conf19

# Objective

## When will a server run out of disk?

Date Disk Reaches Capacity

# August 04 2030

Time Until Disk Reaches Capacity

# 11 Years, 1 Month, 2 Days

Forecasting Disk Usage - Macro View



Go from **Guessing** to **Knowing**,
allowing your organization to run more **Proactively**.

splunk> .conf19

# Common Challenges

Identifying Inefficiencies
Current VS Desired Process Flows
Total Lead Time

splunk> .conf19

# Common Challenges

- Reacting to space issues
- Inefficient work flow stream
- 30% of incidents due to low disk space

Description =        Instance G:
Object LogicalDisk
Counter % Free Space
Has a value 9.24244117736816
At time 2019-05-30T20:09:55.0000000-05:00

| 1<br>New Ticket | | 250<br>Open Tickets |
|---|---|---|
| | | ▲ Status |
| Logical Disk Free Space is low targeting G: | | Assigned |
| Logical Disk Free Space is low targeting G: | | Assigned |

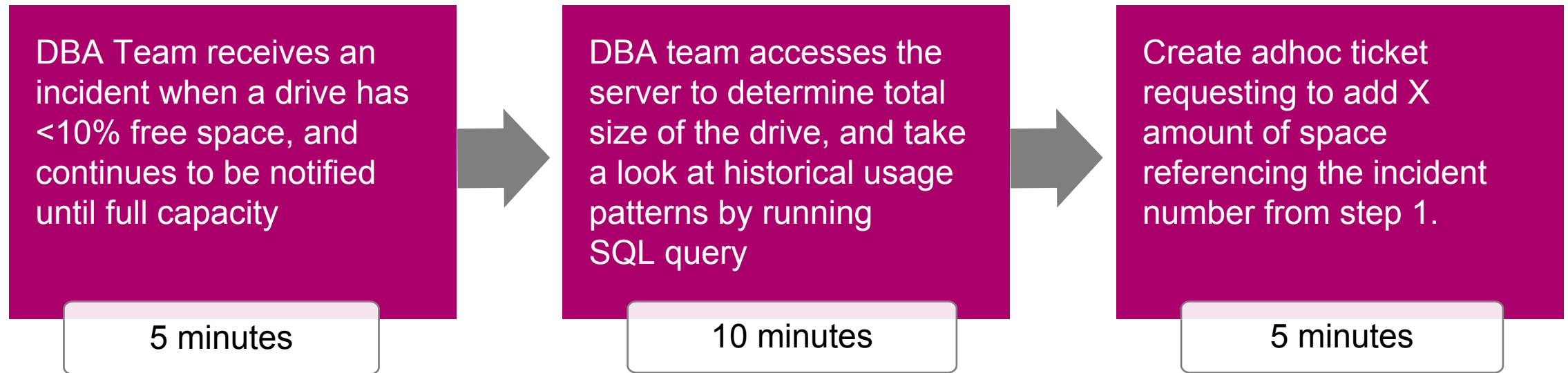splunk> .conf19

# Identifying Inefficiencies

## Inefficient Process

- DBA Team has, at times thousands of incident tickets
- **30%** of incidents are for low disk space
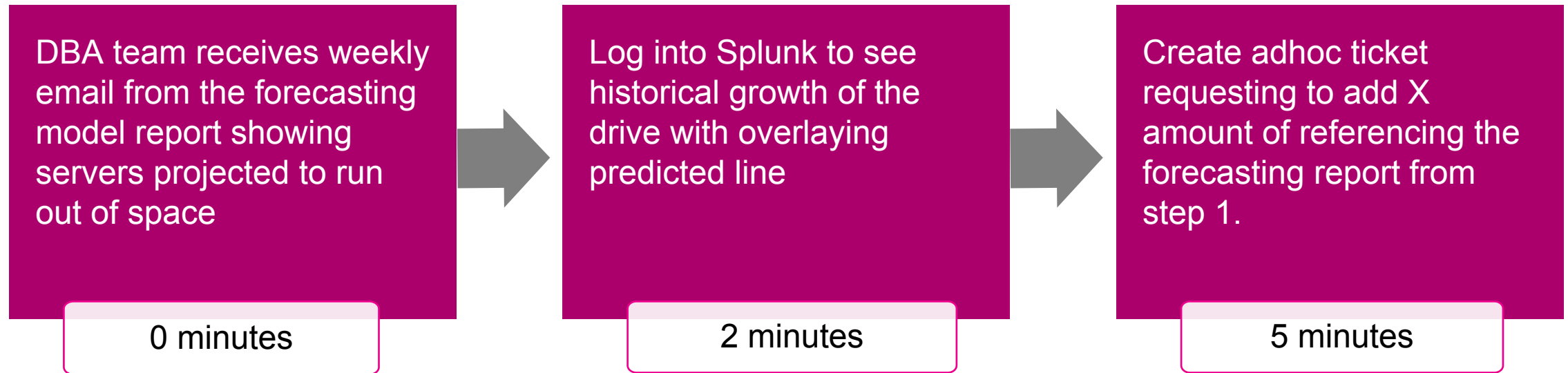
## What Leads to Inefficiency

- Difficult to distinguish the critical issues from the noise
- Using **%** for when a disk is at risk is convoluting

splunk> .conf19

# _Current_ Process Flow Chart

| | | |
|---|---|---|
| DBA Team receives an incident when a drive has <10% free space, and continues to be notified until full capacity | DBA team accesses the server to determine total size of the drive, and take a look at historical usage patterns by running SQL query | Create adhoc ticket requesting to add X amount of space referencing the incident number from step 1. |
| 5 minutes | 10 minutes | 5 minutes |

## Total Time **20 Minutes**

splunk> .conf19

# _Desired_ Process Flow Chart

| DBA team receives weekly email from the forecasting model report showing servers projected to run out of space | Log into Splunk to see historical growth of the drive with overlaying predicted line | Create adhoc ticket requesting to add X amount of referencing the forecasting report from step 1. |
|---|---|---|
| 0 minutes | 2 minutes | 5 minutes |

## Total Time 7 Minutes

splunk> .conf19

# Total Lead Time

Process Flow showcasing adding space to a serve

**Not Enough Space**

Storage opens a change order to add space

Final Change is scheduled to add space, and perform fail overs

**2 WEEKS**

**TOTAL 4 WEEKS**

DBA's request to Windows to add space to the drive

Windows checks the space available

**2 WEEKS**

**Enough Space**

Final Change is scheduled to add space, and perform server fail overs

**1 WEEK**

**TOTAL 3 WEEKS**

splunk> .conf19

# Dark Data

What you can do with your dark data

Sensor Sensei

splunk> .conf19

# Dark Data

## What you can do with your Dark Data

- How many of you have perfmon data available to you?
- How many have MLTK installed or the ability to get it installed?

```
index=perfmon                                                    // Performance data for Windows Servers
sourcetype="Perfmon:FreeDiskSpace"                               // Free Disk Space Values
counter="% Free Space" OR counter="Free Megabytes"               // Free Disk Space either % OR in MB
instance=G:                                                       // Drive Letter
| eval Free_GB=FreeMBytes/1024                                    // Convert MB to GB
| timechart span=1d max(Free_GB) AS Free_GB max(PercentFreeSpace) AS PercentFreeSpace // Bucket time with 1 day span and use MAX Values for "Free" fields
| eval Total_Capacity=round(100*'Free_GB'/'PercentFreeSpace',2)  // Calculate Total Capacity from "Free" Values available
| eval Used_GB=round('Total_Capacity'-'Free_GB',2)               // Calculate Used GB using total capacity from prior calculation and Free GB
| fields _time,Used_GB,Total_Capacity                            // Only showcase the fields Time, Used and Total in GB
```

splunk> .conf19

A server is forecasted to run out of disk in one month, if I add 100GB, **how much time** will that buy me?

Add Storage (GB)

100

splunk> .conf19

# Solve Your Business Problems

**How much time until full capacity is reached?**

**Which day, month and year to reach capacity?**

**If we remove or add 500GB of space, how long will that last?**

11 Years, 24 Days

August 06 2030

Reclaim Storage (GB)

-500

splunk> .conf19

# Technical Deep Dive

Hands On

splunk> .conf19

# Capacity Planning - Forecast Disk Usage by Date - PROD

G: Drive

Enter Future Date

| Before 2020 ▼ | Hide Filters |

## Future Date Entered

# Jan 01 2020 - Wednesday

| Host ⇕ | | Forecasted Percent Full ⇕ |
|---|---|---|
| �_▔▔▔▔▔▔▔▔ | | 110.46 |
| ▔▔▔▔▔ | | 110.03 |
| ▔▔▔▔▔ | | 107.84 |
| ▔▔▔▔▔ | | 107.66 |
| ▔▔▔▔▔ | | 106.98 |
| ▔▔▔▔▔ | | 105.90 |
| ▔▔▔▔▔ | | 98.13 |
| ▔▔▔▔▔ | | 96.05 |
| ▔▔▔▔▔ | | 95.22 |
| ▔▔▔▔▔ | | 94.92 |
| ▔▔▔▔▔ | | 93.95 |
| ▔▔▔▔▔ | | 92.97 |
| ▔▔▔▔▔ | | 92.90 |

Edit · Export ▼ · ...

# Remember This Formula?

Slope Intercept Form

# Slope
Rise Over Run



$$y = 1x + 0$$

$$y = 2x + 0$$

$$y = 3x + 0$$

# Y - Intercept
Where it Crosses the Y-Axis



$$y = 1x - 2$$



$$y = 1x + 0$$
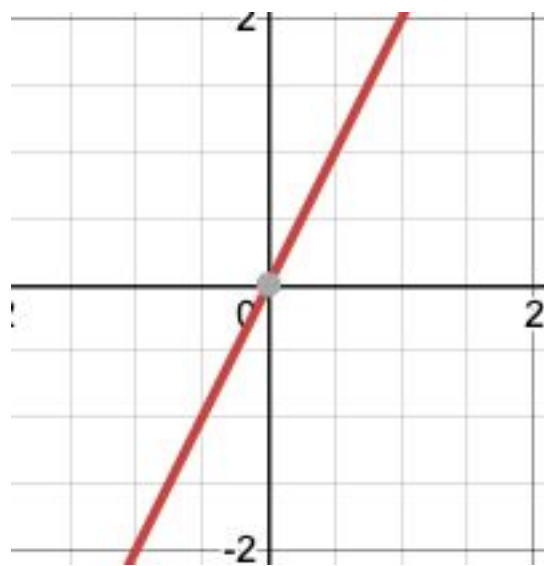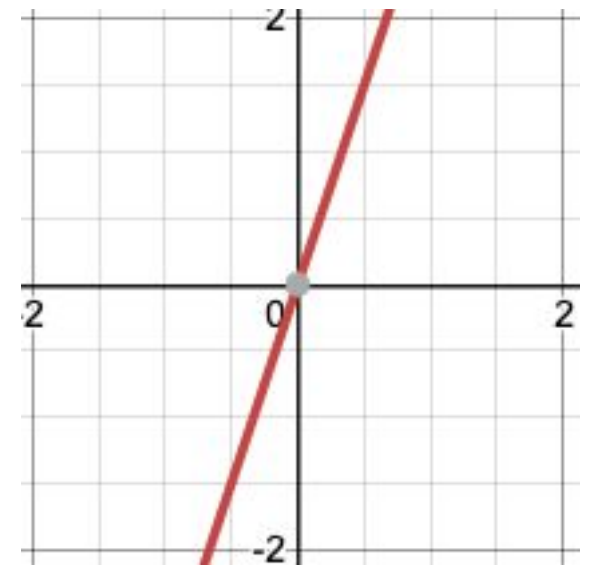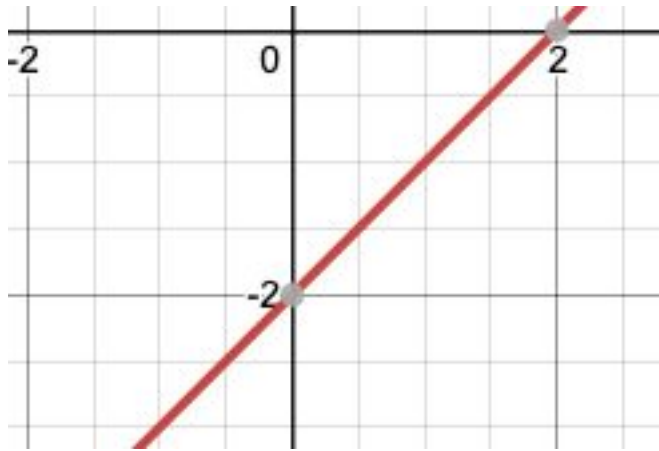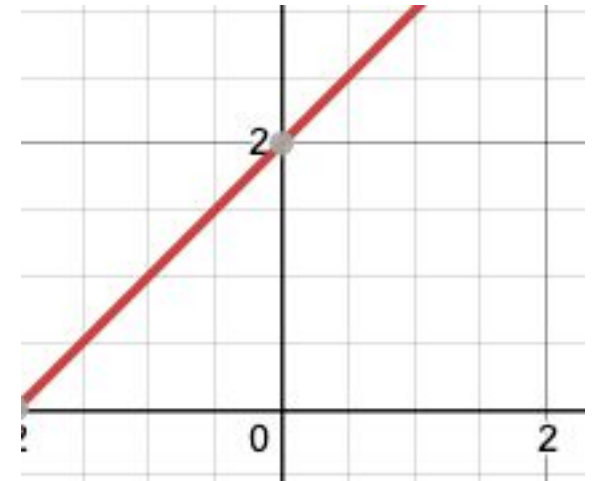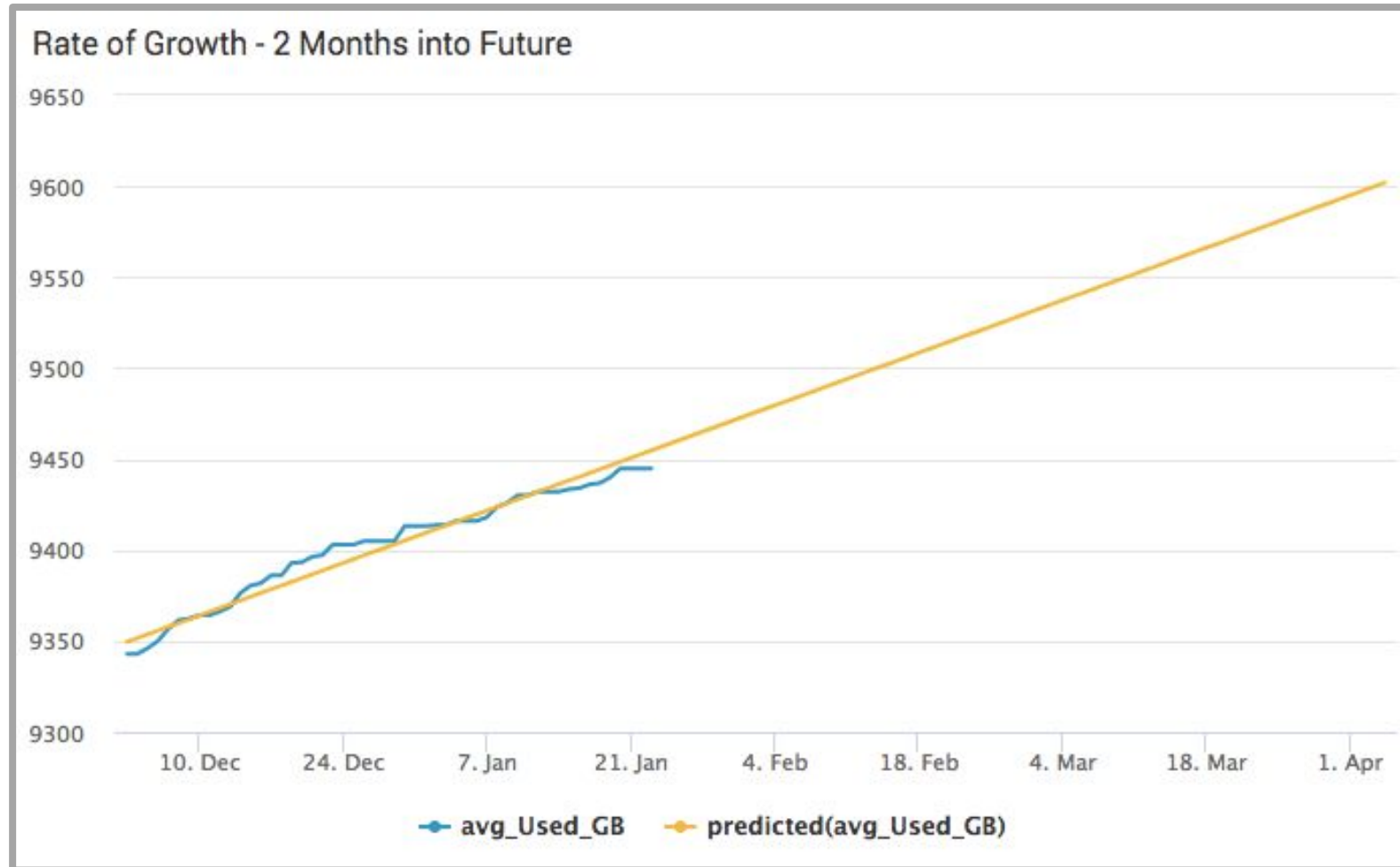


$$y = 1x + 2$$

splunk> .conf19

# Train The Model

Rate of Growth - 2 Months into Future

avg_Used_GB    predicted(avg_Used_GB)

# Create Future Empty Buckets

```
1  | makeresults count=100000         // Make 100,000 empty buckets into the future -- 273.97 years into the future
2  | streamstats count                // Create a new number for each emppty bucket
3  | eval earliest_time=now()         // Identify the earliest time as now
4  | eval time=case(count=100000,relative_time(earliest_time,"+100000d"),count=1,earliest_time) // New field called time with
5  | makecontinuous time span=1d      // Make the data contineous with one day buckets
6  | eval _time=time                  // Convert current time to the time field we made
7  | eval time_human=strftime(time, "%Y-%m-%d %H:%M:%S")  // Create a nicely formatted human readable time
8  | fields + time                    //  Only show me the field "time" with its future empty buckets and buckets in the past
```

| time ⇕ ✎ | time_human ⇕ |
|---|---|
| 1585544400 | 2020-03-30 00:00:00 |
| 1585630800 | 2020-03-31 00:00:00 |
| 1585717200 | 2020-04-01 00:00:00 |
| 1585803600 | 2020-04-02 00:00:00 |
| 1585890000 | 2020-04-03 00:00:00 |
| 1585976400 | 2020-04-04 00:00:00 |
| 1586062800 | 2020-04-05 00:00:00 |
| 1586149200 | 2020-04-06 00:00:00 |

splunk> .conf19

```
1   | makeresults count=100000        // Make 100,000 empty buckets into the future -- 273.97 years into the future
2   | streamstats count               // Create a new number for each emppty bucket
3   | eval earliest_time=now()        // Identify the earliest time as now
4   | eval time=case(count=100000,relative_time(earliest_time,"+100000d"),count=1,earliest_time) // New field called time with
5   | makecontinuous time span=1d     // Make the data contineous with one day buckets
6   | eval _time=time                 // Convert current time to the time field we made
7   | eval time_human=strftime(time, "%Y-%m-%d %H:%M:%S")  // Create a nicely formatted human readable time
8   | fields + time                   //  Only show me the field "time" with its future empty buckets and buckets in the past
9
10  | append                          // Append your query that identifes numeric value for GB used on a per day basis
11     [| search
12  index=xxx host=xxx sourcetype="Perfmon:FreeDiskSpace" counter="% Free Space" OR counter="Free Megabytes" instance=G:
13  | eval FreeGB=FreeMBytes/1024
14  | bucket span=1d _time
15  | stats
16        avg(FreeGB) AS FreeGB
17        avg(PercentFreeSpace) AS PercentFreeSpace
18        by _time, host
19  | eval host=lower(host)
20  | lookup Capacity_Planning_Forecasting_State.csv host OUTPUT Total_Capacity, y_intercept
21  | eval avg_Used_GB=round('Total_Capacity'-'FreeGB',2)
22  | timechart span=1d
23        max(FreeGB) AS FreeGB
24        max(PercentFreeSpace) AS PercentFreeSpace
25        max(Total_Capacity) AS Total_Capacity
26        max(avg_Used_GB) AS avg_Used_GB
27        max(y_intercept) AS y_intercept]
28  | sort + _time
29  | fields + _time  avg_Used_GB Total_Capacity y_intercept
30
31  | apply  Forecasting_$HOST$  // Apply your model
```

# Apply Model onto Future Empty Buckets

# Forecast Into The Future
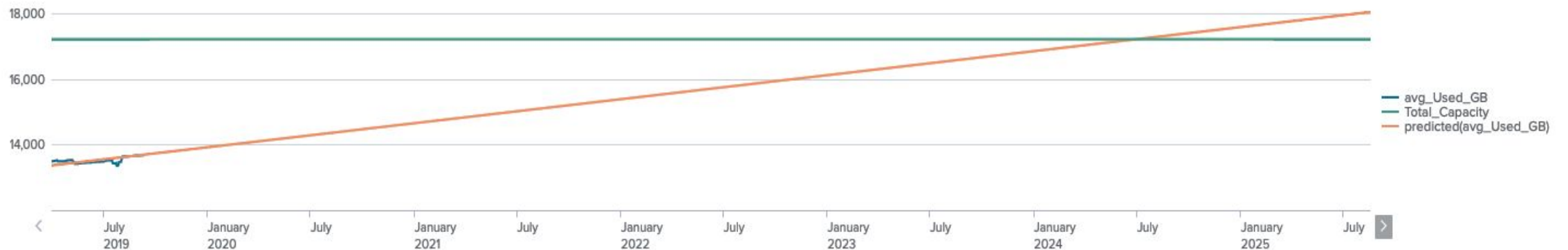
Fill those empty buckets using the apply command

| _time | avg_Used_GB | predicted(avg_Used_GB) |
|---|---|---|
| 2019-06-27 00:00:00 | 8476.35 | 8476.26 |
| 2019-06-28 00:00:00 | 8473.53 | 8478.41 |
| 2019-06-28 09:04:24 | | 8479.22 |
| 2019-06-29 00:00:00 | | 8480.55 |
| 2019-06-30 00:00:00 | | 8482.70 |
| 2019-07-01 00:00:00 | | 8484.85 |
| 2019-07-02 00:00:00 | | 8487.00 |
| 2019-07-03 00:00:00 | | 8489.15 |
| 2019-07-04 00:00:00 | | 8491.30 |
| 2019-07-05 00:00:00 | | 8493.45 |
| 2019-07-06 00:00:00 | | 8495.60 |
| 2019-07-07 00:00:00 | | 8497.75 |

splunk> .conf19

# Convert Disk Usage to Time

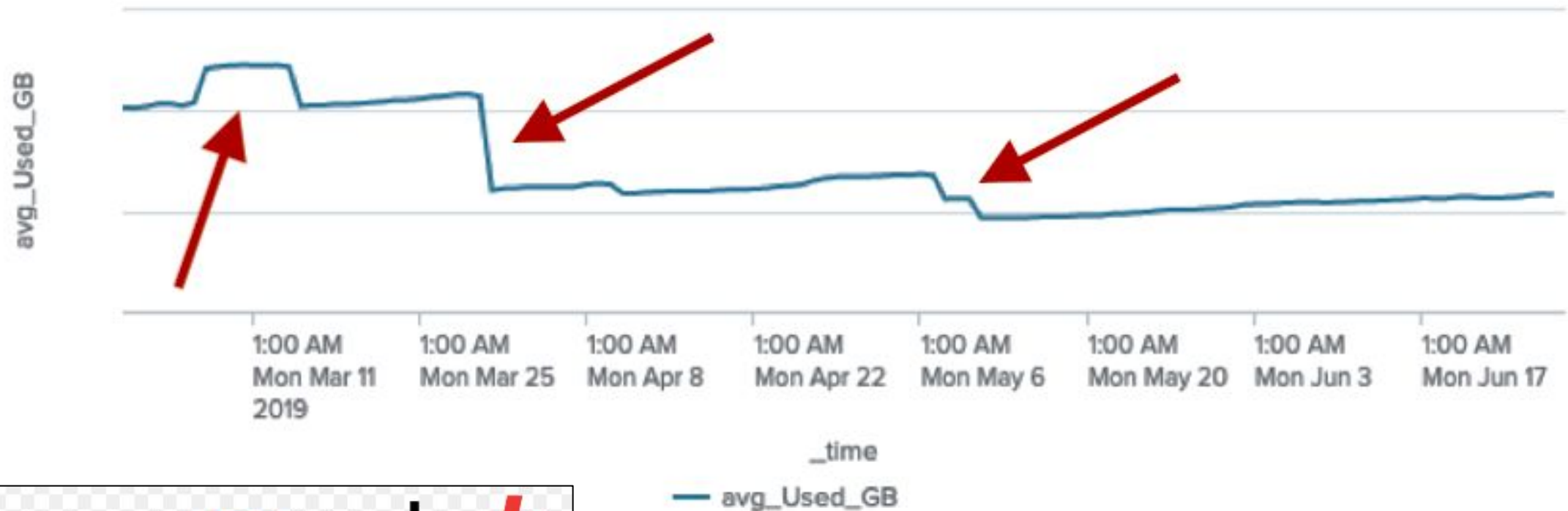Find intersection between forecasted disk usage and total capacity

This intersection will represent when forecasted disk reaches capacity
Take the _time of that intersection

Convert number of days from "today" to calculate the exact Year, Month, and Day disk will run out

# External Forces That Could Skew the Results
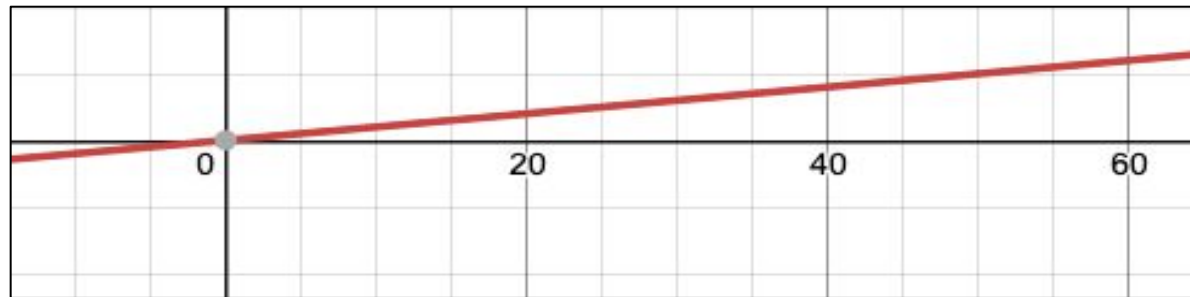
Historical Disk Usage



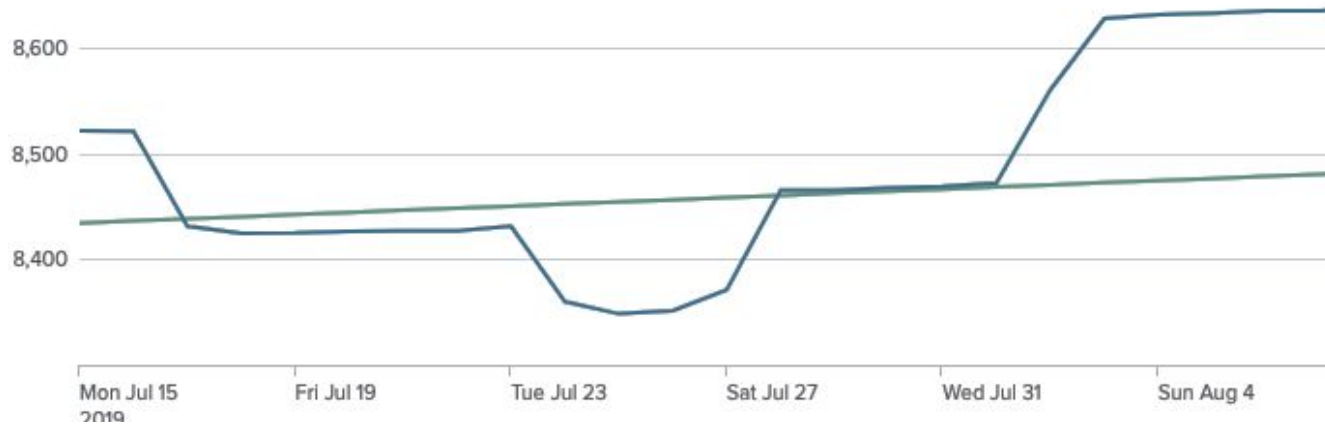$$y = mx + b$$

splunk>  .conf19

# Y – Intercept

What is This and Why Do We Care?

# Scheduled Report to Adjust the Y-Intercept

**Before**



**After**



*With auto-adjusting the Y-Intercept, the predicted values are now using the latest Y-Intercept to forecast.*

splunk> .conf19

# What If Machine

**Current Situation**
without changing total capacity

If we **Reclaim 1TB of Capacity**
what will that look like?

Add Storage (GB)

0

Add Storage (GB)

-1000

## January 03 2032

**12 Years, 0 Months, 25 Days**

## August 21, 2030

**10 Years, 11 Months, 15 Days**

splunk> .conf19

# Build a Scoring Mechanism

## Shorten the Feedback Loops

splunk> .conf19

# Build a Scoring Mechanism to Test Accuracy

# Scaling This Out

Monitor Accuracy with Scheduled Reports

Efficiency Techniques Used

Common Challenges

Overfitting the Data

splunk> .conf19

# Scale your ML Project

**Scaling Touch-Points**

## Schedule **Daily** Residual Reports to ensure accuracy

The scheduled report 'Capacity_Planning_Daily_Residual_Report' has run.

Report:  Capacity_Planning_Daily_Residual_Report

| _time | orig_host | Application Name | abs_residual | abs_days | avg_Used_GB | predicted(avg_Used_GB) | Total_Capacity |
|---|---|---|---|---|---|---|---|
| Sun Aug 25 00:00:00 2019 | ◆≋𝖒&⃛ | ☊⃛✶■𝖒◆ | 0.38 | 578 | 321.83 | 321.45 | 749.87 |
| Sun Aug 25 00:00:00 2019 | ◆≋𝖒&⃛ | ☊⃛✶■𝖒◆ | 115.56 | 230 | 10310.04 | 10425.6 | 14335.87 |
| Sun Aug 25 00:00:00 2019 | ◆≋𝖒&⃛ | ☊⃛✶■𝖒◆ | 6.56 | 57 | 12350.11 | 12343.55 | 15359.87 |
| Sun Aug 25 00:00:00 2019 | ◆≋𝖒&⃛ | ☊⃛✶■𝖒◆ | 0.06 | 55 | 5336.35 | 5336.29 | 7841.46 |
| Sun Aug 25 00:00:00 2019 | ◆≋𝖒&⃛ | ☊⃛✶■𝖒◆ | 12.28 | 15 | 8633.09 | 8620.81 | 9380 |
| Sun Aug 25 00:00:00 2019 | ◆≋𝖒&⃛ | ☊⃛✶■𝖒◆ | 21.28 | 10.1 | 8521.92 | 8543.2 | 9419.87 |
| Sun Aug 25 00:00:00 2019 | ◆≋𝖒&⃛ | ☊⃛✶■𝖒◆ | 21.99 | 6.28 | 7206.63 | 7184.64 | 8191.87109 |

splunk> .conf19

# Scale your ML Project
Scaling Touch-Points

Get familiar with the map command

```
| map maxsearches=1000 search="     //map is looping through the host list, and will not exceed 1000 hosts as input
| makeresults count=100000          //creates 100000 empty buckets that will be used to hold future forecasted values
| streamstats count as count        // streamstats is used to produce a cumulative count of the buckets
| eval earliest_time=now()          // the empty buckets will be created from "NOW" until 1000000 buckets have been reached
```

© 2019 SPLUNK INC.

# Efficiency Techniques

**Train during a change freeze**

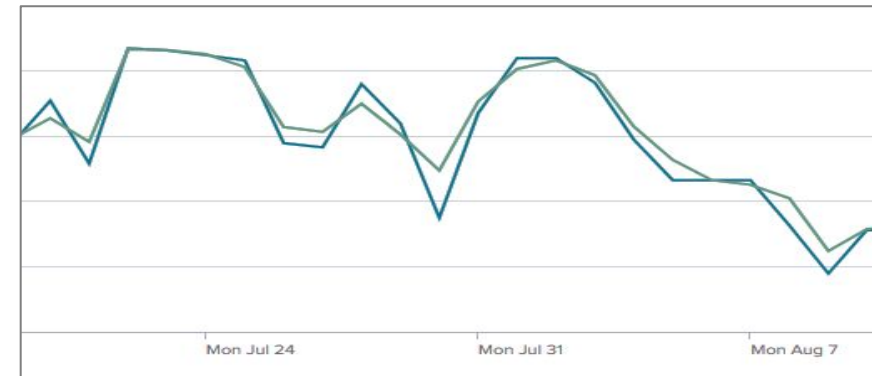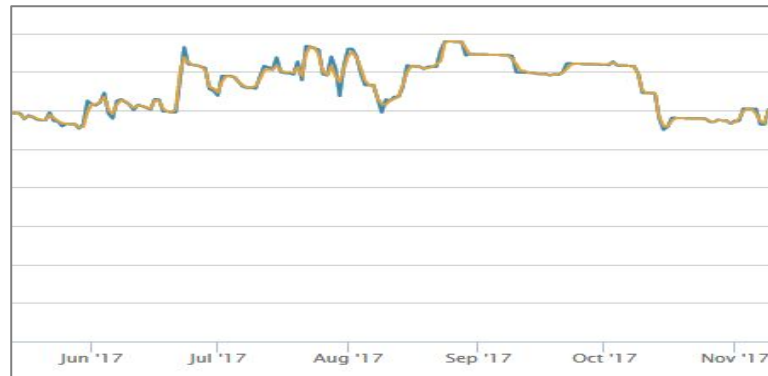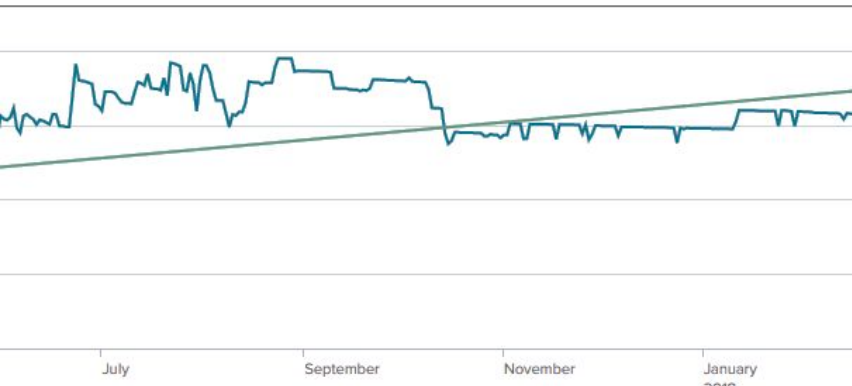**Use | loadjob for post-processing**

**Summary Indexing to spread out the load**



```
1 | loadjob 1568161600.78_CCED4_98DB6899KD    //Loads events or results of a previously completed search job
2 | eval residual = 'predicted(avg_Used_GB)' - 'avg_Used_GB' // calculating the residual values
3 | eval days='residual'/'slope'   // calculates how many days off is the prediction
```
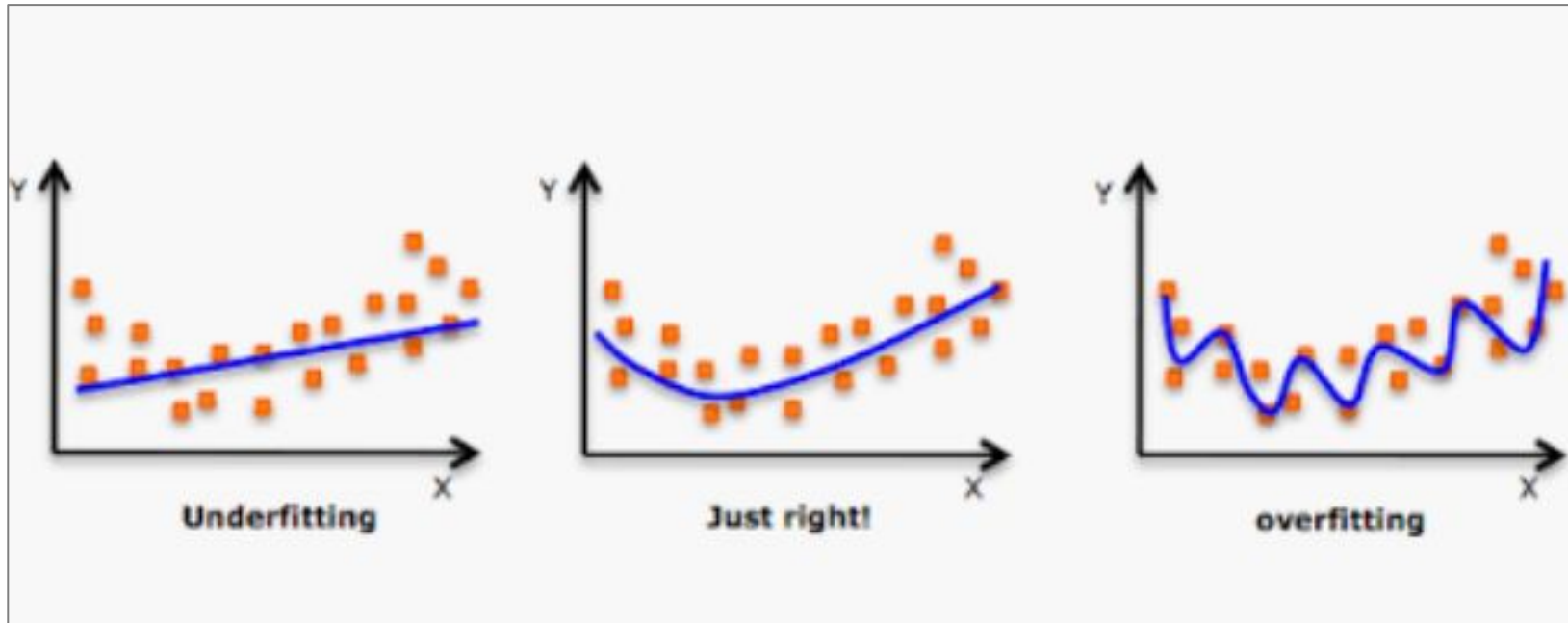
splunk> .conf19

# Challenges

- **25 failed prototype dashboards** before coming to the conclusion of using the linear regression
- Data didn't express **complete** linearity..
- Searches took longer than expected to run

# Overfitting
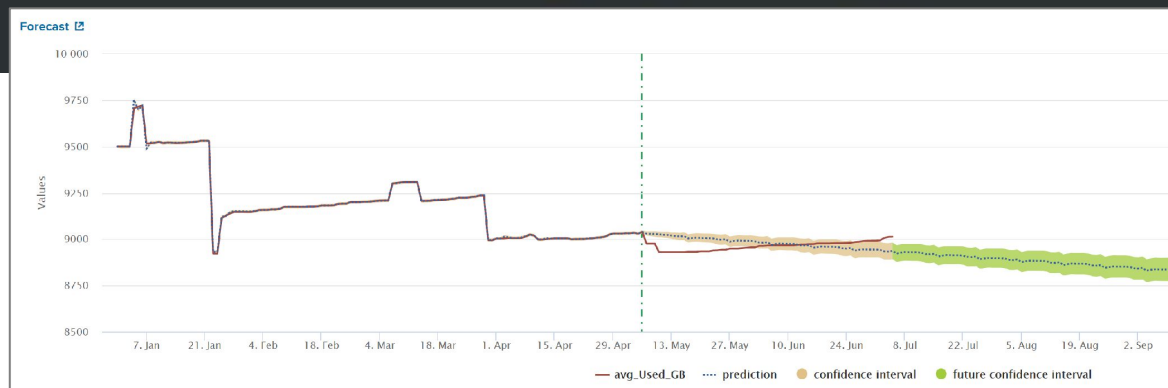
1. Using the predict command resulted in a really good prediction…almost too good to be true.
2. The best solution to an overfitting problem is avoidance.
3. Identify overfitting, then find ways to tackle overfitting and learn from the mistake.

# **Predict Command**

**KALMAN FILTERING** – when using the predict command we were actually overfitting the data

```
index=tu_perfmon host=XXXXXXX sourcetype="Perfmon:FreeDiskSpace"
counter="% Free Space" OR counter="Free Megabytes" instance=G:
| eval Free_GB=FreeMBytes/1024
    | timechart span=1d max(Free_GB) AS Free_GB max(PercentFreeSpace) AS PercentFreeSpace
    | eval avg_Total_Capacity=round(100*'Free_GB'/'PercentFreeSpace',2)
    | eval avg_Used_GB=round('avg_Total_Capacity'-'Free_GB',2)
| fields  _time avg_Used_GB
| predict "avg_Used_GB" as prediction algorithm=LLP holdback=0 future_timespan=180 upper10=upper10 lower10=lower10
| `forecastviz(180, 0, "avg_Used_GB", 10)`
```

# ARIMA

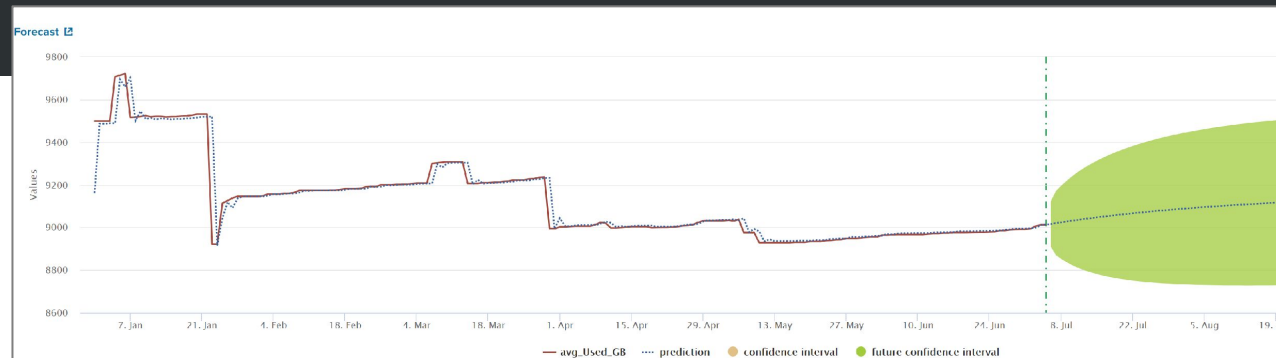ARIMA – gave large confidence interval (we want a smaller interval in order to trust the accuracy of the forecast)

```
index=tu_perfmon host=xxxxx sourcetype="Perfmon:FreeDiskSpace" counter="% Free Space" OR counter="Free Megabytes" instance=G:
| eval Free_GB=FreeMBytes/1024
    | timechart span=1d max(Free_GB) AS Free_GB max(PercentFreeSpace) AS PercentFreeSpace
    | eval avg_Total_Capacity=round(100*'Free_GB'/'PercentFreeSpace',2)
    | eval avg_Used_GB=round('avg_Total_Capacity'-'Free_GB',2)
| fields  _time avg_Used_GB avg_Total_Capacity
| fit ARIMA _time avg_Used_GB holdback=0 conf_interval=95 order=3-0-0 forecast_k=90 as prediction
| `forecastviz(90, 0, "avg_Used_GB", 95)`
```



splunk> .conf19

# Linear Regression

**Linear Regression**- showed incredible accuracy and didn't over fit the data like the predict command
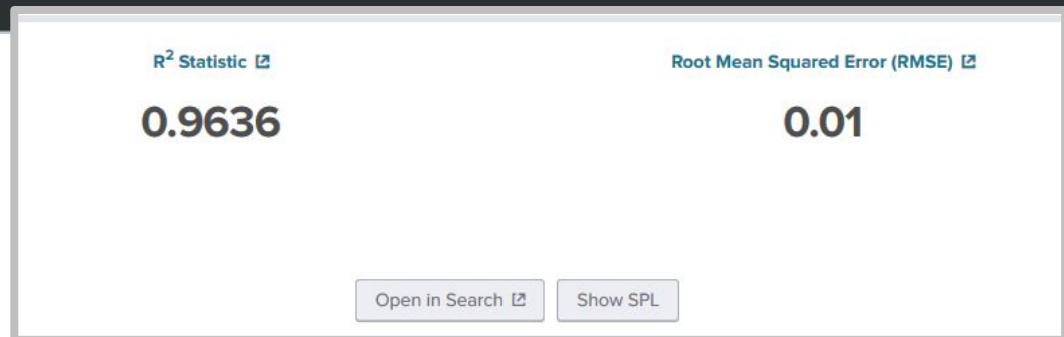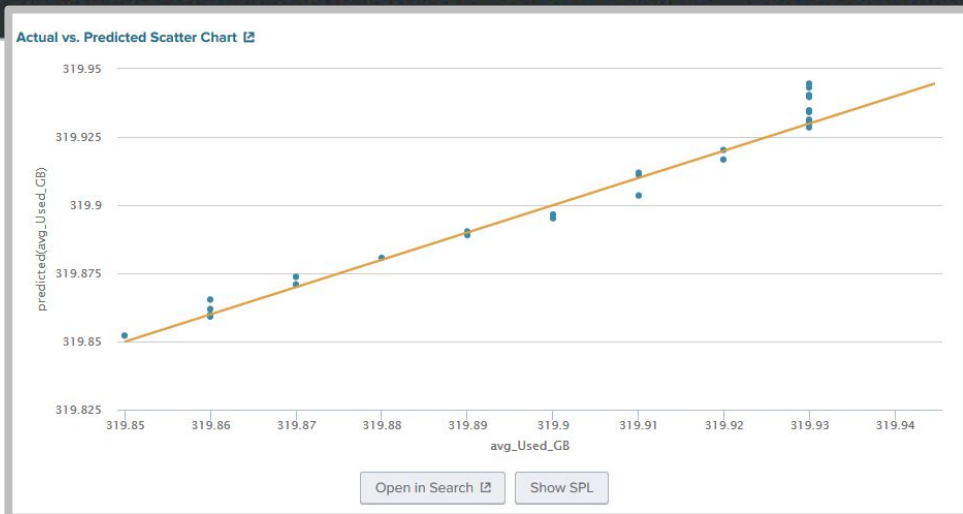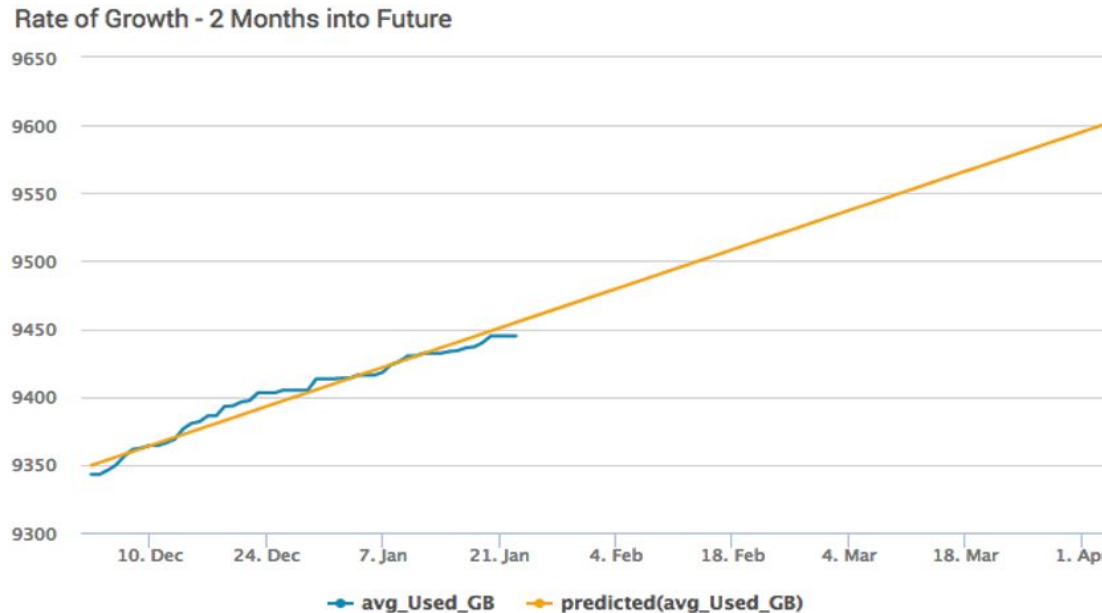
```
index=tu_perfmon host=xxxxxxx sourcetype="Perfmon:FreeDiskSpace" counter="% Free Space" OR counter="Free Megabytes" instance=G:
| eval Free_GB=FreeMBytes/1024
| timechart span=1d max(Free_GB) AS Free_GB max(PercentFreeSpace) AS PercentFreeSpace
| eval avg_Total_Capacity=round(100*'Free_GB'/'PercentFreeSpace',2)
| eval avg_Used_GB=round('avg_Total_Capacity'-'Free_GB',2)
| fields _time, avg_Used_GB
| fit LinearRegression "avg_Used_GB" from "_time" fit_intercept=true into "_exp_draft_589637e1779f4c7790dbd5a2a325cf82"
```



Actual vs. Predicted Scatter Chart



R² Statistic: 0.9636

Root Mean Squared Error (RMSE): 0.01

# What This *Will* Do

1. Assist with Prevention of Outages, and Budgeting
2. Reduce time when addressing capacity concerns
3. Allow for Re-usability within the organization



Rate of Growth - 2 Months into Future

splunk> .conf19

# What This Will *NOT* Do

1. Business logic needs to be accounted for when looking at the visualization.
2. This will not account for external changes in growth that do not have an established pattern

Walk-through

splunk> .conf19

# Where to Start?

| Find your perfmon data | → | Create a search to identify Total Capacity and Used GB over time | → | Install the MLTK app |
|---|---|---|---|---|

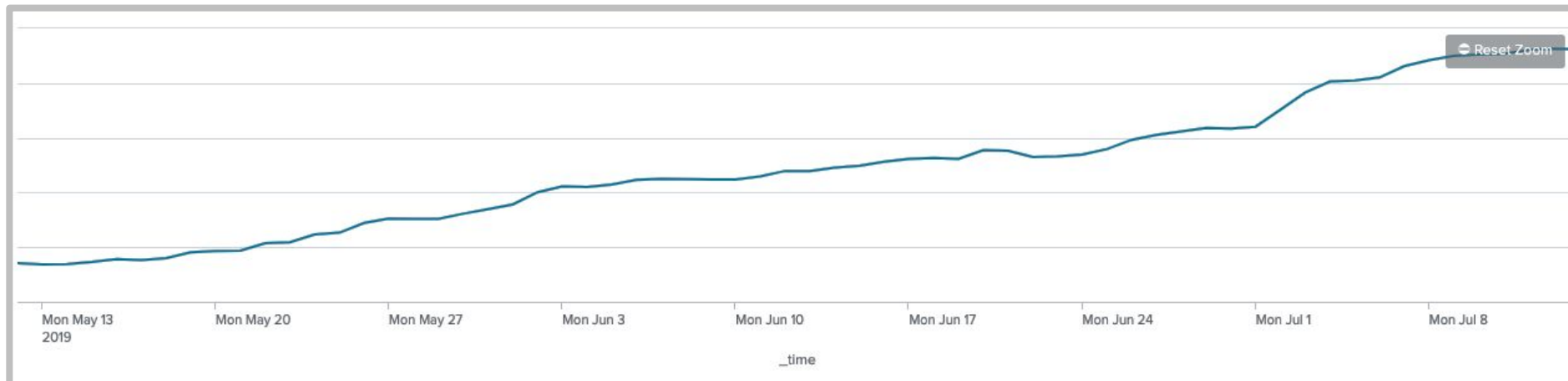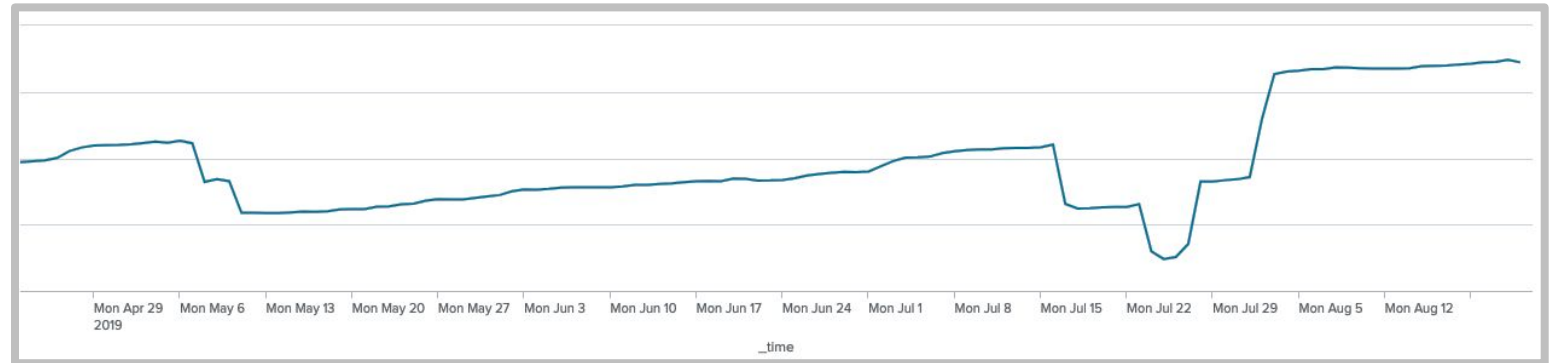Splunk Machine Learning Toolkit

```
1    index=xxx host=xxx sourcetype="Perfmon:FreeDiskSpace" counter="% Free Space" OR counter="Free Megabytes" instance=G:
2    | eval FreeMBytes=if(counter=="Free Megabytes",Value,null())
3    | eval FreeGB=FreeMBytes/1024
4    | eval PercentFreeSpace=if(counter=="% Free Space",Value,null())
5    | bucket span=1d _time
6    | stats
7            avg(FreeGB) AS FreeGB
8            avg(PercentFreeSpace) AS PercentFreeSpace
9            by _time
10   | eval avg_Total_Capacity=(100*FreeGB)/PercentFreeSpace
11   | eval avg_Used_GB=round('avg_Total_Capacity'-'FreeGB',2)
12   | timechart span=1d
13           max(FreeGB) AS FreeGB
14           max(PercentFreeSpace) AS PercentFreeSpace
15           max(avg_Total_Capacity) AS avg_Total_Capacity
16           max(avg_Used_GB) AS avg_Used_GB
17   | sort + _time
```

splunk> .conf19

# Find Target Host
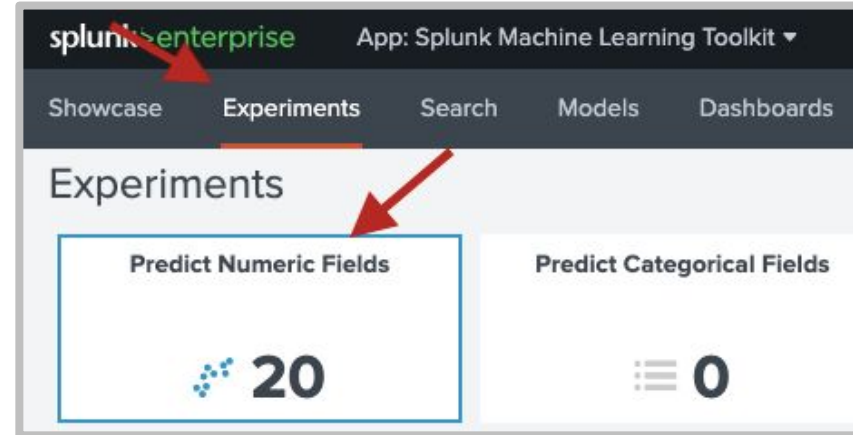
1. Identify a single server and instance to test against

2. Validate that your growth has a good constant linear slope. Ex.(find a time during a change freeze)

# Using the MLTK

1.  Navigate to the MLTK app > Experiments > Predict Numerical Fields

2.  Enter your SPL into the search

3.  Select "Linear Regression" as the algorithm

4.  Select field to predict and single _time feature

# Interpret Results

**R^2 >** how much the estimates deviate from the actual values in the data set on average.

- * Better Models have R^2 closer to 1.



**R² Statistic** ↗  **Root Mean Squared Error (RMSE)** ↗

0.9506   6.75

Open in Search ↗   Show SPL

**RMSE >** indicates how close the observed data points are to the model's predicted values.

- * Better models have a smaller RMSE.



Actual vs. Predicted Line Chart ↗

Sort by: Default Sort ▼

avg_Used_GB   predicted(avg_Used_GB)

Actual vs. Predicted Scatter Chart ↗

predicted(avg_Used_GB)

avg_Used_GB

© 2019 SPLUNK INC.

# Save and Publish

# Apply

Apply your saved model onto your data

```
1   | makeresults count=100000
2   | streamstats count as count
3   | eval earliest_time=now()
4   | eval time=case(count=100000,relative_time(earliest_time,"+100000d"),count=1,earliest_time)
5   | makecontinuous time span=1d
6   | eval timeAsANumber=time
7   | eval _time=time
8   | eval time_human=strftime(time, "%Y-%m-%d %H:%M:%S")
9   | fields + time
10
11
12  | append
13      [| search
14  index=xxx host=xxx sourcetype="Perfmon:FreeDiskSpace" counter="% Free Space" OR counter="Free Megabytes"
15  | eval FreeMBytes=if(counter=="Free Megabytes",Value,null())
16  | eval FreeGB=FreeMBytes/1024
17  | eval PercentFreeSpace=if(counter=="% Free Space",Value,null())
18  | bucket span=1d _time
19  | stats
20      avg(FreeGB) AS FreeGB
21      avg(PercentFreeSpace) AS PercentFreeSpace
22      by _time, host
23  | eval host=lower(host)
24  | lookup Capacity_Planning_Forecasting_State.csv host OUTPUT Total_Capacity, y_intercept
25  | eval avg_Used_GB=round('Total_Capacity'-'FreeGB',2)
26  | timechart span=1d
27      max(FreeGB) AS FreeGB
28      max(PercentFreeSpace) AS PercentFreeSpace
29      max(Total_Capacity) AS Total_Capacity
30      max(avg_Used_GB) AS avg_Used_GB
31      max(y_intercept) AS y_intercept]
32  | sort + _time
33  | fields + _time  avg_Used_GB Total_Capacity y_intercept
34
35
36  | apply Forecasting_Demo
```

# Auto-Adjust the Y-Intercept

Start with gathering the data from a lookup file

```
1   | inputlookup Capacity_Planning_Forecasting_State.csv // lookup contains drive letter, host name, total capacity, and y-intercept value
2   | map maxsearches=100  // map is looping through the host list, and will not exceed 100 hosts as input
```

Use the same search string that we have been using to gather Total Capacity as well as Used GB

```
3   search="search
4   index=tu_perfmon host=$host$ sourcetype=\"Perfmon:FreeDiskSpace\" counter=\"% Free Space\" OR counter=\"Free Megabytes\" instance=G:
5   | eval FreeGB=FreeMBytes/1024
6   | bucket span=1d _time
7   | stats
8        avg(FreeGB) AS FreeGB
9        avg(PercentFreeSpace) AS PercentFreeSpace
10       by _time, host
11  | eval host=lower(host)
12  | lookup Capacity_Planning_Forecasting_State.csv host OUTPUT Total_Capacity, y_intercept
13  | eval avg_Used_GB=round('Total_Capacity'-'FreeGB',2)
14  | bin _time span=1d
15  | stats
16       max(FreeGB) AS FreeGB
17       max(PercentFreeSpace) AS PercentFreeSpace
18       max(Total_Capacity) AS Total_Capacity
19       max(avg_Used_GB) AS avg_Used_GB
20       max(y_intercept) AS y_intercept by _time,host
21  | sort + _time
22  | apply  Forecasting_$host$"
```

splunk> .conf19

# Auto-Adjust the Y-Intercept cont…

- Filling the future buckets with the current total capacity, and y-intercept value

- Calculating % Capacity

- Outputlookup conducts the y-intercept adjustment in the lookup file.

```
23  | filldown Total_Capacity // fill the empty buckets with the current Total Capacity value
24  | filldown y_intercept     // fill the empty buckets with the current y-intercept value
25  | eval one_mon=now()+2592000
26  | eval Total_Capacity=if(Total_Capacity="",'Total_Capacity','Total_Capacity')
27  | eval Percent_Capacity=Total_Capacity*1.00
28  | eval predicted(avg_Used_GB)='predicted(avg_Used_GB)'-'y_intercept'-0
29  | eval avg_Used_GB='avg_Used_GB'-0
30  | fields -   Percent_Capacity one_mon
31  | eval residual = 'predicted(avg_Used_GB)' - 'avg_Used_GB'
32  | eval y_intercept_new=residual
33  | eval predicted(avg_Used_GB)='predicted(avg_Used_GB)'-'y_intercept_new'
34  | eval residual_new='predicted(avg_Used_GB)'-avg_Used_GB
35  | eval y_intercept=y_intercept+y_intercept_new
36  | eval drive="G"
37  | fields + drive Total_Capacity host y_intercept
38  | fields - _time
39  | outputlookup override_if_empty=false Capacity_Planning_Forecasting_State.csv
```

splunk> .conf19

# Build a Forecasting Dashboard

## Base Search

```
| makeresults count=100000
| streamstats count as count
| eval earliest_time=now()
| eval time=case(count=100000,relative_time(earliest_time,"+100000d"),count=1,earliest_time)
| makecontinuous time span=1d
| eval _time=time
| eval time_human=strftime(time, "%Y-%m-%d %H:%M:%S")
| fields + time
| append
[| search
index=tu_perfmon host=$HOST$ sourcetype="Perfmon:FreeDiskSpace" counter="% Free Space" OR counter="Free Megabytes"
instance=G:
| eval FreeGB=FreeMBytes/1024
| bucket span=1d _time
| stats avg(FreeGB) AS FreeGB  avg(PercentFreeSpace) AS PercentFreeSpace by _time, host
| eval host=lower(host)
| lookup Capacity_Planning_Forecasting_State.csv host OUTPUT Total_Capacity, y_intercept|
eval avg_Used_GB=round('Total_Capacity'-'FreeGB',2)
| timechart span=1d max(FreeGB) AS FreeGB  max(PercentFreeSpace) AS PercentFreeSpace  max(Total_Capacity) AS
Total_Capacity       max(avg_Used_GB) AS avg_Used_GB  max(y_intercept) AS y_intercept]
| sort + _time
| fields + _time  avg_Used_GB Total_Capacity y_intercept
| apply  Forecasting_$HOST$
| filldown Total_Capacity| filldown y_intercept
```

splunk> .conf19

# Build a Forecasting Dashboard

```
1   | fields *
2   | eval Total_Capacity=if(Total_Capacity="",'Total_Capacity','Total_Capacity')
3   | eval Percent_Capacity=Total_Capacity*$TOTAL_CAP$
4
5   | eval predicted(avg_Used_GB)='predicted(avg_Used_GB)'-'y_intercept'-$OFFSET$
6
7   | eval intersect=if('predicted(avg_Used_GB)'>=Percent_Capacity,_time,"")
8   | eval intersect=strftime(intersect,"%B %d %Y")
9   | where isnotnull(intersect)
10  | head 1
11  | fields + intersect
```

Date Disk Reaches Capacity

# July 21 2030

splunk> .conf19

# Build a Forecasting Dashboard

```
33  | fields *
34  | eval Total_Capacity=if(Total_Capacity="",'Total_Capacity','Total_Capacity')
35  | eval Percent_Capacity=Total_Capacity*1.00
36  | eval predicted(avg_Used_GB)='predicted(avg_Used_GB)'-'y_intercept'-0
37  | eval intersect=if('predicted(avg_Used_GB)'>=Percent_Capacity,_time,"")
38  | eval TIME=now()
39  | eval how_many_years=intersect-TIME
40  | eval intersect=strftime(intersect,"%B %d %Y")
41  | where isnotnull(intersect)
42  | head 1
43  | fields + intersect _time TIME how_many_years
44  | eval TIME=strftime(TIME, "%Y-%m-%d")
45  | eval how_many_years=floor(how_many_years/60/60/24)*60*60*24
46  | `duration(how_many_years )`
47  | fields - duration_*
48  | rename duration AS "Time Until Capacity is Reached"
49  | fields + "Time Until Capacity is Reached"
```

Time Until Disk Reaches Capacity

# 10 Years, 10 Months, 13 Days

splunk> .conf19

# Build a Forecasting Dashboard



```
36 | eval residual = 'predicted(avg_Used_GB)' - 'avg_Used_GB'
37 | eval zero=0
38 | fields + _time residual zero
```

Residuals to Test Accuracy of Forecast

# Build a Forecasting Dashboard

```
37  | eval residual = 'predicted(avg_Used_GB)' - 'avg_Used_GB'
38  | eval zero=0
39  | fields + _time residual zero predicted(avg_Used_GB) slope
40  | eval days='residual'/'slope'
41
42  | sort - _time
43  | eval epoch_time=strptime(_time, "%s")
44  | eval epoch_secs='days'*86400
45  | eval tolerance=epoch_time-epoch_secs
46  | eval Actual_Date=strftime(tolerance, "%Y-%m-%d")
47  | fields _time residual days Actual_Date update
48  | eval days_abs=abs(days)
49  | eval SLA=if(days_abs>7.1,"ERROR","GOOD")
50  | eval date_day=strftime(_time, "%d")
51  | eval today=now()
52  | eval today=strftime(today, "%d")
53  | where date_day!=today
54  | fields _time residual  days SLA
```

### Accuracy with Number of Days Tolerance

Day column shows the range of days the prediction is off - SLA goes in ERROR status if prediction is off by more than a week

| _time ‡ | residual ‡ | days ‡ | SLA ‡ |
|---|---|---|---|
| 2019-09-12 | -0.00 | -0.00148 | GOOD |
| 2019-09-11 | -0.08 | -0.04 | GOOD |
| 2019-09-10 | -1.27 | -0.592 | GOOD |
| 2019-09-09 | 2.32 | 1.08 | GOOD |
| 2019-09-08 | 2.54 | 1.18 | GOOD |
| 2019-09-07 | -0.23 | -0.11 | GOOD |

# Q&A

Steve Koelpin | Splunk Advisor

splunk> .conf19