# Beyond "Regular" Regular Expressions

Cary Petterborg

Senior Engineer, LDS Church

# Boiler Plate

# Disclaimer

splunk> .conf2016

# My Disclaimer

During the course of this presentation, I may make references to my employer, the Church of Jesus Christ of Latter-day Saints. This should not be taken as an endorsement of Splunk or Splunk products by the LDS Church.

# About Me

# Who is Cary Petterborg?

- Splunk user and administrator for 4.5 years

- Monitoring Engineer for 9 years

- Web developer for 22 years

- Software engineer for 36 years

- Many languages from assembly to Ruby

- Application development including Flight Sim, DB systems, and Web

- Works for the LDS Church in Salt Lake City

- Favorite cologne while in Florida – Deep Woods OFF

# Why Regular Expressions?

- Using Regular Expressions since the mid 80's

- Started using regex with lex/yacc/sed/grep for software development

- Realized the power of regex quickly

- Taught classes on regex

- Love working with regex stuff in Splunk and other utilities

- Regex is an important skill, and I want to share my skill

- While everyone was playing Pokemon GO on their phone this last summer, I was having **fun with regular expressions**

# One day, you too…

# One day, you too…

# Where are regular expressions useful?

# What do Regular Expressions do to you?

# In Splunk

- The **rex** and **regex** search commands

- In *props.conf, transforms.conf* and other *.conf* files

- Field extractions

- Data feeds

- Splunk regular expressions are PCRE (Perl Compatible Regular Expressions) and use the PCRE C library.

# Elsewhere

Though mostly similar, there are differences in the various implementations, but many of the concepts carry across from one form to another:

- Shell commands

- Utilities like *sed, lex, grep, egrep*

- Programming languages (*perl, python, php, C*, etc.) through libraries or built-in functionality

- **Sed** mode available in Splunk commands like **rex** and in *transforms.conf* files

# The Splunk Field Extraction Tool

.conf2016

splunk>

# Field Extraction Tool

- GUI tool in the web UI of Splunk

- Simple to use

- Pretty good *for a start*, but not *always* good for a final result

- Not able to optimize or do anything complex

- Limited in order or number of named capture groups

- Makes mistakes if you don't have regular data

- **Demo**

Search     Pivot     Reports     Alerts     Dashboards

🔍 New Search     Save As ▾   Close

sourcetype="linux_secure"     All time ▾   🔍

✓ 600 events (before 9/8/16 3:35:29.000 PM)     No Event Sampling ▾     Job ▾   ‖ ◻ ↗ 🖨 ⬇   ● Smart Mode ▾

| Events (600) | Patterns | Statistics | Visualization |

Format Timeline ▾     — Zoom Out     + Zoom to Selection     ✕ Deselect     1 hour per column

List ▾     Format ▾     20 Per Page ▾     ‹ Prev   1   2   3   4   5   6   7   8   9   …   Next ›

| | i | Time | Event |
|---|---|---|---|

‹ Hide Fields     ≡ All Fields

**Selected Fields**

a host 1
a source 1
a sourcetype 1

**Interesting Fields**

# date_hour 18
# date_mday 3
# date_minute 60
a date_month 1
# date_second 60
a date_wday 3

| | 6/8/16 9:39:12.000 AM | Wed Jun 08 2016 09:39:12 www1 sshd[4092]: Failed password for invalid user local from 109.169.32.135 port 2172 ssh2 |
|---|---|---|
| | | host = Carys-MacBook-Pro.local   source = secure.log   sourcetype = linux_secure |
| | 6/8/16 9:39:12.000 AM | Wed Jun 08 2016 09:39:12 www1 sshd[38618]: Failed password for nsharpe from 10.2.10.163 port 8317 ssh2 |
| | | host = Carys-MacBook-Pro.local   source = secure.log   sourcetype = linux_secure |
| | 6/8/16 9:38:00.000 AM | Wed Jun 08 2016 09:38:00 www1 sshd[3860]: Failed password for invalid user dean from 109.169.32.135 port 3860 ssh2 |
| | | host = Carys-MacBook-Pro.local   source = secure.log   sourcetype = linux_secure |
| | 6/8/16 9:37:15.000 AM | Wed Jun 08 2016 09:37:15 www1 sshd[3799]: Failed password for invalid user operator from 109.169.32.135 port 1735 ssh2 |
| | | host = Carys-MacBook-Pro.local   source = secure.log   sourcetype = linux_secure |
| | 6/8/16 9:37:00.000 AM | Wed Jun 08 2016 09:37:00 www1 sshd[4454]: Failed password for invalid user itmadmin from 109.169.32.135 port 1490 ssh2 |

# Examples of Field Extractor Deficiencies

# Multi-format Security Data in...

# Simple FET extraction of Port



**Select Fields**

Highlight one or more values in the sample event to create fields. You can indicate one value is required, meaning it must exist in an event for the regular expres
part of an existing extraction, first turn off the existing extractions. Learn more ⬈

Wed Jun 08 2016 09:39:12 www1 sshd[4092]: Failed password for invalid user local from 109.169.32.135 port  2172  ssh2

Hide Regular Expression ⌄

^(?:[^\.\n]*\.){3}\d+\s+\w+\s+(?P<port>\d+)

splunk> .conf2016

# Intelligent extraction of Port

**Regular Expression**

```
port\s+(?P<port>\d+)
```

# Named Capture Groups

- (?P<name>…) === (?<name>…)
- The **P** is optional (came from **P**ython), but it is usually considered more correct
- Splunk FET will use (?P<name>…), so why not make things similar?

  **BUT**

- Do it the way you feel most comfortable

# Goal: *user* from all entries using one regex

```
1 sshd[3697]: Failed password for invalid user whois from 10.1.10.172 port 124
1 sudo:  djohnson ;   TTY=pts/0 ; PWD=/home/djohnson ; USER=root ; COMMAND=/bin/
1 sshd[4333]: Failed password for root from 10.1.10.172 port 2772 ssh2
1 sshd[1632]: Failed password for invalid user elena_andubasquet from 10.1.10.
1 sshd[1155]: Failed password for invalid user yp from 10.1.10.172 port 1822 s
1 sshd[2021]: Failed password for myuan from 10.1.10.172 port 4468 ssh2
1 sshd[81145]: pam_unix(sshd:session): session opened for user djohnson by (ui
1 sshd[4986]: Failed password for invalid user mysql from 128.241.220.82 port
1 sshd[4761]: Failed password for invalid user local from 128.241.220.82 port
```

# There is no user automatically extracted

| | 6/8/16 9:39:12.000 AM | Wed Jun 08 2016 09:39:12 www1 sshd[4092]: Failed password for invalid user local from 109.169.32.1 |
|---|---|---|

**Event Actions ⌄**

| Type | ☑ Field | Value | Actions |
|---|---|---|---|
| Selected | ☑ host ⌄ | Carys-MacBook-Pro.local | ⌄ |
| | ☑ source ⌄ | secure.log | ⌄ |
| | ☑ sourcetype ⌄ | security | ⌄ |
| Event | ☐ index ⌄ | main | ⌄ |
| | ☐ linecount ⌄ | 1 | ⌄ |
| | ☐ splunk_server ⌄ | Carys-MacBook-Pro.local | ⌄ |
| Time ⊕ | _time ⌄ | 2016-06-08T09:39:12.000-06:00 | |
| Default | ☐ punct ⌄ | ___.:__[]:_____..__ | ⌄ |

splunk> .conf2016

# Field Extractor Failed



or invalid user local from 109.169.32.135 port 2172

for nsharpe from 10.2.10.163 port 8317 ssh2

or invalid user dean from 109.169.32.135 port 3860

or invalid user operator from 109.169.32.135 port '

or invalid user itmadmin from 109.169.32.135 port '

or mail from 84.34.159.23 port 1190 ssh2

or sync from 84.34.159.23 port 2530 ssh2

or invalid user inet from 10.3.10.46 port 1516 ssh:

ssion): session closed for user myuan by (uid=0)

or invalid user administrator from 10.3.10.46 port

# Field Extractor Failed After More Lines Added

⚠️ **The extraction failed. If you are extracting multiple fields, try removing one or more fields. Start with extractions tha**

## Select Fields

Highlight one or more values in the sample event to create fields. You can indicate one value is required, meaning it must exist in an

```
Wed Jun 08 2016 09:39:12 www1 sshd[4092]: Failed password for invalid user local from 109.169.32
```

```
⊗  Wed Jun 08 2016 02:13:48 www3 sshd[68976]: Accepted password for djohnson from 10.3.10.46 por
```

splunk> .conf2016

# Let's look at the generated REGEX

Hide Regular Expression ∨

`^\w+\s+\w+\s+\d+\s+\d+\s+\d+:\d+:\d+\s+\w+\d+\s+\w+\[\d+\]:\s+\w+\s+\w+\s+\w+\s+\w+\s+\w+\s+(?P<user>[^ ]+)`

splunk> .conf2016

# My tool of choice: regex101.com

# Add the data and a regex:



REGULAR EXPRESSION

/ for (?P<user>\S+) from

TEST STRING

```
Mon Jun 06 2016 02:29:19 www1 sshd[3849]: Failed password for root from 128.241.22
Mon Jun 06 2016 02:29:24 www1 sshd[4267]: Failed password for invalid user adminis
Mon Jun 06 2016 02:29:44 www1 sshd[5001]: Failed password for invalid user george
Mon Jun 06 2016 02:30:13 www1 sshd[1638]: pam_unix(sshd:session): session opened f
Mon Jun 06 2016 02:30:13 www1 sshd[1480]: Failed password for invalid user yp from
Mon Jun 06 2016 02:30:22 www1 sshd[2291]: Failed password for invalid user email f
Mon Jun 06 2016 02:30:26 www1 sshd[4761]: Failed password for invalid user local f
Mon Jun 06 2016 02:30:50 www1 sshd[4986]: Failed password for invalid user mysql f
Mon Jun 06 2016 02:31:19 www1 sshd[81145]: pam_unix(sshd:session): session opened
Mon Jun 06 2016 02:31:19 www1 sshd[2021]: Failed password for myuan from 10.1.10.1
Mon Jun 06 2016 02:31:28 www1 sshd[1155]: Failed password for invalid user yp from
Mon Jun 06 2016 02:31:32 www1 sshd[1632]: Failed password for invalid user elena_c
Mon Jun 06 2016 02:31:54 www1 sshd[4333]: Failed password for root from 10.1.10.17
Mon Jun 06 2016 02:32:00 www1 sudo: djohnson ;  TTY=pts/0 ; PWD=/home/djohnson ; L
Mon Jun 06 2016 02:32:00 www1 sshd[3697]: Failed password for invalid user whois f
Mon Jun 06 2016 02:32:19 www1 sshd[5985]: Failed password for invalid user testuse
```

# Refine the regex – better matches, but not all

# Refine the regex again – almost there

REGULAR EXPRESSION

`/ for ((invalid user )|(user ))?(?P<user>\S+) (from|by)|`

TEST STRING

```
Mon Jun 06 2016 02:29:19 www1 sshd[3849]: Failed password for root from 128.241.220.82 port 2253 ssh2
Mon Jun 06 2016 02:29:24 www1 sshd[4267]: Failed password for invalid user administrator from 128.241.22
Mon Jun 06 2016 02:29:44 www1 sshd[5001]: Failed password for invalid user george from 128.241.220.82 po
Mon Jun 06 2016 02:30:13 www1 sshd[1638]: pam_unix(sshd:session): session opened for user djohnson by (u
Mon Jun 06 2016 02:30:13 www1 sshd[1480]: Failed password for invalid user yp from 128.241.220.82 port 2
Mon Jun 06 2016 02:30:22 www1 sshd[2291]: Failed password for invalid user email from 128.241.220.82 por
Mon Jun 06 2016 02:30:26 www1 sshd[4761]: Failed password for invalid user local from 128.241.220.82 por
Mon Jun 06 2016 02:30:50 www1 sshd[4986]: Failed password for invalid user mysql from 128.241.220.82 por
Mon Jun 06 2016 02:31:19 www1 sshd[81145]: pam_unix(sshd:session): session opened for user djohnson by (
Mon Jun 06 2016 02:31:19 www1 sshd[2021]: Failed password for myuan from 10.1.10.172 port 4468 ssh2
Mon Jun 06 2016 02:31:28 www1 sshd[1155]: Failed password for invalid user yp from 10.1.10.172 port 1822
Mon Jun 06 2016 02:31:32 www1 sshd[1632]: Failed password for invalid user elena_andubasquet from 10.1.1
Mon Jun 06 2016 02:31:54 www1 sshd[4333]: Failed password for root from 10.1.10.172 port 2772 ssh2
Mon Jun 06 2016 02:32:00 www1 sudo: djohnson ;   TTY=pts/0 ; PWD=/home/djohnson ; USER=root ; COMMAND=/bi
Mon Jun 06 2016 02:32:00 www1 sshd[3697]: Failed password for invalid user whois from 10.1.10.172 port 2
```

# And FINALLY – we got them all

REGULAR EXPRESSION

`/ ((for ((invalid user )|(user ))?)|(sudo: ))(?P<user>\S+) (from|by)?`

TEST STRING

```
Mon Jun 06 2016 02:29:19 www1 sshd[3849]: Failed password for root from 128.241.220.82 port 2253 ssh2
Mon Jun 06 2016 02:29:24 www1 sshd[4267]: Failed password for invalid user administrator from 128.241.2:
Mon Jun 06 2016 02:29:44 www1 sshd[5001]: Failed password for invalid user george from 128.241.220.82 p
Mon Jun 06 2016 02:30:13 www1 sshd[1638]: pam_unix(sshd:session): session opened for user djohnson by (
Mon Jun 06 2016 02:30:13 www1 sshd[1480]: Failed password for invalid user yp from 128.241.220.82 port :
Mon Jun 06 2016 02:30:22 www1 sshd[2291]: Failed password for invalid user email from 128.241.220.82 po
Mon Jun 06 2016 02:30:26 www1 sshd[4761]: Failed password for invalid user local from 128.241.220.82 po
Mon Jun 06 2016 02:30:50 www1 sshd[4986]: Failed password for invalid user mysql from 128.241.220.82 po
Mon Jun 06 2016 02:31:19 www1 sshd[81145]: pam_unix(sshd:session): session opened for user djohnson by (
Mon Jun 06 2016 02:31:19 www1 sshd[2021]: Failed password for myuan from 10.1.10.172 port 4468 ssh2
Mon Jun 06 2016 02:31:28 www1 sshd[1155]: Failed password for invalid user yp from 10.1.10.172 port 182:
Mon Jun 06 2016 02:31:32 www1 sshd[1632]: Failed password for invalid user elena_andubasquet from 10.1.:
Mon Jun 06 2016 02:31:54 www1 sshd[4333]: Failed password for root from 10.1.10.172 port 2772 ssh2
Mon Jun 06 2016 02:32:00 www1 sudo: djohnson ;    TTY=pts/0 ; PWD=/home/djohnson ; USER=root ; COMMAND=/b:
Mon Jun 06 2016 02:32:00 www1 sshd[3697]: Failed password for invalid user whois from 10.1.10.172 port :
```

splunk> .conf2016

# Back to the Field Extractor – use our regex

Use the event listing below to validate the field extractions produced by your regular expression.

## Regular Expression

```
((for ((invalid user )|(user ))?)|(sudo: ))(?P<user>\S+) (from|by)?
```

# And the results are MUCH better

```
1 sshd[3697]: Failed password for invalid user whois from 10.1.10.172 port 124
1 sudo: djohnson ;   TTY=pts/0 ; PWD=/home/djohnson ; USER=root ; COMMAND=/bin/
1 sshd[4333]: Failed password for root from 10.1.10.172 port 2772 ssh2
1 sshd[1632]: Failed password for invalid user elena_andubasquet from 10.1.10.
1 sshd[1155]: Failed password for invalid user yp from 10.1.10.172 port 1822 s
1 sshd[2021]: Failed password for myuan from 10.1.10.172 port 4468 ssh2
1 sshd[81145]: pam_unix(sshd:session): session opened for user djohnson by (ui
1 sshd[4986]: Failed password for invalid user mysql from 128.241.220.82 port
1 sshd[4761]: Failed password for invalid user local from 128.241.220.82 port
```

# Dissecting the Path to Success

# Note the progression in the regex

➢ for (?P<user>\S+) from

➢ for (invalid user )?(?P<user>\S+) from

➢ for ((invalid user )|(user ))?(?P<user>\S+) (from|by)

➢ ((for ((invalid user )|(user ))?)|(sudo: ))(?P<user>\S+) (from|by)?

# Explanation

➢ for (?P<user>\S+) from



EXPLANATION                                                      ⊖

▲ / for (?P<user>\S+) from /

  for    matches the characters  for   literally (case sensitive)

▲ (?P<user>\S+) **Named capturing group** user

  ▲ \S+ **match any non-white space character** [^\r\n\t\f ]

      Quantifier: + Between one and unlimited times, as

      many times as possible, giving back as needed [greedy]

  from matches the characters  from literally (case sensitive)

# Explanation (cont)

➢ for (invalid user )?(?P<user>\S+) from

EXPLANATION                                                    ⊖

◢ / for (invalid user )?(?P<user>\S+) from /

  for  matches the characters  for  literally (case sensitive)

◢ 1st Capturing group (invalid user )?

  Quantifier: ? Between zero and one time, as many times as possible, giving back as needed [greedy]

  *Note:* A repeated capturing group will only capture the last iteration. Put a capturing group around the repeated group to capture all iterations or use a non-capturing group instead if you're not interested in the data

  invalid user  matches the characters invalid user  literally (case sensitive)

◢ (?P<user>\S+) Named capturing group user

  ◢ \S+ match any non-white space character [^\r\n\t\f ]

    Quantifier: + Between one and unlimited times, as many times as possible, giving back as needed [greedy]

  from matches the characters  from literally (case sensitive)

# Explanation (cont)

➤ for ((invalid user )|(user ))?(?P<user>\S+) (from|by)



1st **Alternative:** (invalid user )
  2nd **Capturing group** (invalid user )
    invalid user matches the characters invalid user literally (case sensitive)
2nd **Alternative:** (user )
  3rd **Capturing group** (user )
    user matches the characters user literally (case sensitive)

5th **Capturing group** (from|by)
  1st **Alternative:** from
    from matches the characters from literally (case sensitive)
  2nd **Alternative:** by
    by matches the characters by literally (case sensitive)

# Explanation (cont)

> ((for ((invalid user )|(user ))?)|(sudo: ))(?P<user>\S+) (from|by)?

◢ 1st **Capturing group** ((for ((invalid user )|(user ))?)|(sudo: ))
  ◢ 1st **Alternative:** (for ((invalid user )|(user ))?)
    ◢ 2nd **Capturing group** (for ((invalid user )|(user ))?)
      for  matches the characters for  literally (case sensitive)
    ◢ 3rd **Capturing group** ((invalid user )|(user ))?
      Quantifier:  ? Between zero and one time, as many times as possible, giving back as needed [greedy]
      *Note: A repeated capturing group will only capture the last iteration. Put a capturing group around the repeated group to capture all iterations or use a non-capturing group instead if you're not interested in the data*
      ◢ 1st **Alternative:** (invalid user )
        ◢ 4th **Capturing group** (invalid user )
          invalid user  matches the characters invalid user  literally (case sensitive)
      ◢ 2nd **Alternative:** (user )
        ◢ 5th **Capturing group** (user )
          user  matches the characters user  literally (case sensitive)
  ◢ 2nd **Alternative:** (sudo: )
    ◢ 6th **Capturing group** (sudo: )
      sudo:  matches the characters sudo:  literally (case sensitive)

# Explanation (cont)

➤ ((for ((invalid user )|(user ))?)|(sudo: ))(?P<user>\S+) <span style="color:red">(from|by)?</span>



8th **Capturing group** `(from|by)?`
  `Quantifier:` `?` Between `zero` and `one` time, as many times as possible, giving back as needed `[greedy]`
  *Note:* *A repeated capturing group will only capture the last iteration. Put a capturing group around the*
  *repeated group to capture all iterations or use a non-capturing group instead if you're not interested in the data*
1st **Alternative:** `from`
    `from` matches the characters `from` literally (case sensitive)
2nd **Alternative:** `by`
    `by` matches the characters `by` literally (case sensitive)

# It's not hard, it just takes understanding

- Start by giving defining a simple part of the regex that will work on one of the line **types**.

- You might want to define another regex for a different type of line and try to combine the two regex's after

  Or

- You can just modify the one that you had to find an additional instance

- Remember – You can only have ONE and only one named capture group *for a given name* **IN THE SAME REGULAR EXPRESSION**, but you can have multiple named capture groups with different names, or different regular expression with the same capture group name

splunk> .conf2016

# One named capture group with a single name

- This will fail:

- (for invalid user (?P<user>\S+))|(for (?P<user>\S+))



EXPLANATION                                                            ⊖

/(for invalid user (?P<user>\S+))|(for (?P<user>\S+))/g
Errors are explained from left to right. Move the mouse cursor over them to see the error highlighted in your pattern
(?P<user> Subpattern name declared more than once
) Unmatched parenthesis

splunk> .conf2016

# Use parenthesis!

- Use parenthesis to define capture groups to define groups of text that may be optional or when you want to pick from **One Of**

REGULAR EXPRESSION

/ `for ((invalid user )|(user ))?(?P<user>\S+) (from|by)`

TEST STRING

- Not only does it make it more readable later, but this one wouldn't work without them
  - **(from|by) user** is not the same as **from|by user**

# Use the best character class

- Sometime you need to make a field be able to match data that you haven't seen in the data yet, so in this case be general

- The **\S** is best because **\s** will be the delimiter (a space in this case) because you want to catch any potential case that you don't see in the data – yet.

$$\S+$$

- If you have a delimiter that you can count on, use something like this to match the field value (int this case be specific about what it is NOT):

$$[^,]+$$

# Now you can wake up – seriously!

.conf2016

splunk>

# Transforms.conf Gotchas

# Problem

You can't index Social Security Numbers in your data, but phone numbers are okay.

- Use transforms.conf to clean up the data

- Leave the phone numbers alone, but take out the SSN's

# SSN vs Phone #

Distinguishing in regex

| **SSN** | **Phone #** |
|---|---|
| 123-45-6789 | 800-123-4567 |
| • \d{3}-\d{2}-\d{4} | • \d{3}-\d{3}-\d{4} |

# Be as specific in your matches as possible

- You could use:

$$\text{\textbackslash d+-\textbackslash d+-\textbackslash d+}$$

- But it will mistake a phone number for a SSN:

REGULAR EXPRESSION

/ \d+-\d+-\d+

TEST STRING

My phone number is 800-555-1212.
My social security number is 123-45-6789.

splunk> .conf2016

# This is a better match

- You could use:

$$\backslash d+-\backslash d\backslash d-\backslash d+$$

- Simple change, but it will NOT mistake a phone number for a SSN, but…:



```
REGULAR EXPRESSION

/  \d+-\d\d-\d+

TEST STRING

My phone number is 800-555-1212.
My social security number is 123-45-6789.
My birthday is 6-11-1958.
I want to order and item with a part number of 71-34-912.
```

# We're so close

- This is exactly what every SSN looks like (in regex):

## \d{3}-\d{2}-\d{4}

- This matches better, but it needs more work:

REGULAR EXPRESSION

/ \d{3}-\d{2}-\d{4}

TEST STRING

```
My phone number is 800-555-1212.
My social security number is 123-45-6789.
My birthday is 6-11-1958.
I want to order item with part number 8871-34-91268.
```

splunk> .conf2016

# So let's make it specific

- This is exactly what every SSN looks like (in regex):

$$\D(?P<ssn>\d\{3\}-\d\{2\}-\d\{4\})\D$$

- This matches only the SSN:

REGULAR EXPRESSION

/ `\D(?P<ssn>\d{3}-\d{2}-\d{4})\D`

TEST STRING

```
My phone number is 800-555-1212.
My social security number is 123-45-6789.
My birthday is 6-11-1958.
I want to order item with part number 8871-34-91268.
```

# And finally:

- Don't let it match when there are digits at the front or the end of the group:

$$(.*)(?<!\backslash d)(\backslash d\{3\}-\backslash d\{2\}-\backslash d\{4,4\})(\$|\backslash D)(.*)$$

This matches only the SSN:

REGULAR EXPRESSION

```
/ (.*)(?<!\d)(\d{3}-\d{2}-\d{4,4})($|\D)(.*)
```

TEST STRING

```
My phone number is 800-555-1212.
My social security number is 123-45-6789.
My birthday is 6-11-1958.
I want to order item with part number 871-34-91269.
I want to order item with part number 8871-34-9126.
```

# Regex101.com explains



EXPLANATION                                                                    ⊖

◢ / (.*)(?<!\d)(\d{3}-\d{2}-\d{4,4})($|\D)(.*) / g
  ◢ **1st** **Capturing group** (.*)
    ◢ **.*** **matches any character (except newline)**
        **Quantifier:** **\*** Between **zero** and **unlimited** times, as many times as possible, giving back
        as needed **[greedy]**
  ◢ (?<!\d) **Negative Lookbehind - Assert that it is impossible to match the regex below**
        \d match a digit **[0-9]**
  ◢ **2nd** **Capturing group** (\d{3}-\d{2}-\d{4,4})
    ◢ \d{3} **match a digit** **[0-9]**
        **Quantifier:** **{3}** Exactly **3** times
        **-** matches the character **-** literally
    ◢ \d{2} **match a digit** **[0-9]**
        **Quantifier:** **{2}** Exactly **2** times
        **-** matches the character **-** literally
    ◢ \d{4,4} **match a digit** **[0-9]**
        **Quantifier:** **{4,4}** Exactly **4** times
  ◢ **3rd** **Capturing group** ($|\D)
    ◢ **1st** **Alternative:** **$**
        **$** assert position at end of the string
    ◢ **2nd** **Alternative:** **\D**
        \D match any character that's not a digit **[^0-9]**
  ◢ **4th** **Capturing group** (.*)
    ◢ **.*** **matches any character (except newline)**
        **Quantifier:** **\*** Between **zero** and **unlimited** times, as many times as possible, giving back
        as needed **[greedy]**
    **g modifier:** global. All matches (don't return on first match)

splunk> .conf2016

# The entry

[noSSN]
REGEX = (?m)^(.*)(?<!\d)(\d{3}-\d{2}-\d{4,4})($|\D)(.*)$
FORMAT = $1###-##-####$3$4
DEST_KEY = _raw

# Why this REGEX

- Multi-line match, otherwise this will only work on single line events:

  (?m)

- Beginning of line and everything up to the last non-digit before our SSN in a capture group.

  ^(.*)

- Here is the SSN (only match if not preceded by a digit)

  (?<!\d)(\d{3}-\d{2}-\d{4,4})

- The first non-digit if not the end of the line in one capture group, then everything from the last non-digit after the SSN to the end of the line in the last capture group

  ($|\D)(.*)$

# Why this FORMAT

- $1 is to be replaced with the first capture group

- $2 would be the SSN, but it isn't included in the FORMAT

- $3 is to be replaced with the single non-digit after the SSN (if it exists)

- $4 is to be replaced with the end capture group

- The SSN is substituted with hashes

$1###-##-####$3$4

# REX and REGEX commands

.conf2016

splunk>

# When to use REX

- You **don't always** want to extract the data

- You want to extract data from a field that is already extracted

- You don't have access to field extractions (permissions, etc.)

- You are in a hurry or you are doing a proof-of-concept

# REX example

- **index=ics.eup sourcetype=voice* Description | rex "Description=(?P<description>[^\^]+)" | rex field=description "From (?P<start>.+) to (?P<end>.+?):\s"**

- Aug  2 20:32:28 l13772 CPCMl13772: %local7-2-ALARM: 16$Description= Number of AuthenticationFailed events exceeds configured threshold during configured interval of time 1 within 3 minutes  on cluster StandAloneCluster.  There are 2 AuthenticationFailed events (up to 30) received during the monitoring interval From Wed Aug 03 10:25:00 PHT 2016 to Wed Aug 03 10:28:00 PHT 2016:   TimeStamp : 8/3/16 10:26 AM LoginFrom : 10.127.34.40 Interface : VMREST UserID : JacobMD AppID : Cisco Tomcat ClusterID :  NodeID : APPHMANAOVM001  TimeStamp : Wed Aug 03 10:26:13 PHT 2016 TimeStam::Status=2,cleared^Severity=minor^Acknowledged=no^CUSTOMER=Cisco Prime Collaboration^Private IP Address=10.160.17.24^Default Alarm Name=AuthenticationFailed^Managed Object=10.160.17.24^Managed Object Type=Unity Connection^MODE=2;Alarm ID=343815480^Component=10.160.17.24\x00000

# REX commands

- index=ics.eup sourcetype=voice* Description
  | rex "Description=(?P<description>[^\^]+)"
  | rex field=description "From (?P<start>.+) to (?P<end>.+?):\s"

- SYNTAX:
  | rex [field=*fieldname*] "regex"

- Also available:
  | rex mode=sed

# First rex – get the description

# Second rex – get the start and end

REGULAR EXPRESSION
1 MATCH - 1016 STEPS

```
From (?P<start>.+) to (?P<end>.+?):\s
```
g

TEST STRING

```
 Number of AuthenticationFailed events exceeds configured threshold during
configured interval of time 1 within 3 minutes  on cluster StandAloneClust
er.  There are 2 AuthenticationFailed events (up to 30) received during th
e monitoring interval From Wed Aug 03 10:25:00 PHT 2016 to Wed Aug 03 10:2
8:00 PHT 2016:    TimeStamp : 8/3/16 10:26 AM LoginFrom : 10.127.34.40 Inte
rface : VMREST UserID : JacobMD AppID : Cisco Tomcat ClusterID :  NodeID :
APPHMANAOVM001  TimeStamp : Wed Aug 03 10:26:13 PHT 2016  TimeStam::Status
=2,cleared
```

splunk> .conf2016

# When to use REGEX

- To filter out events/data that you don't want included in the pipeline

- This is like **search** on steroids, but doesn't replace **search**

- Only used as a filter

# Regex example

# Breakdown

- Search:
  **sourcetype=security | regex "10\.\d+\.\d+\.\d+"**

- Only internal (10.*) IP addresses make it though the **regex** filter

- Search produces events, regex then limits those results passed on through the pipeline by a fancy regular expression

- *Yes, there are other way to do this, but this is a regex example*

# Rex vs Regex

- Use **rex** to extract fields
- Use **regex** to limit results
- Yes, you can use them in the same search:

```
sourcetype=security | rex "from (?P<src_ip>\d+\.\d+\.\d+\.\d+)"
| regex src_ip="(?<!10)\.\d+\.\d+\.\d+"
```

Greedy vs. Lazy matches

# GREEDY vs LAZY

# Where is the difference?

- Greedy – Grab as much as you can

- Lazy – Grab as little as you can

- The lazy match will continue only as far as it needs to:

    <.+?> will match <12345>, while

    <.+> will match both <12345> and <12345><67890>

- SYNTAX: place a **?** after a **\*** or **+**

The lazy match only goes to the first instance of a match following the multiple match

# GREEDY



REGULAR EXPRESSION

/ \((?P<cmd>.*)session

# LAZY



REGULAR EXPRESSION

`/ \((?P<cmd>.*?)session`

```
led password for invalid user ge
_unix(sshd:session): session ope
led password for invalid user yp
led password for invalid user en
led password for invalid user lo
led password for invalid user my
m_unix(sshd:session): session op
led password for myuan from 10.1
led password for invalid user yp
```

# Second Look – Greedy

# Second Look - Lazy

# Choose your path wisely

- Greedy may cross long segments

- Lazy may stop prematurely

- Try it on various data sets to make sure it will do what you want

# Other Notes

# Performance considerations

- Many field extractions can be costly

- Complex regular expressions can be costly

- Use the Job Inspector to see if there is a difference in doing one complex field extraction vs. many simple field extractions

# Maintenance considerations

- The complex field extractions may be easier to maintain than multiple simple extractions
- Your own field extractions will probably be easier to maintain than those created using the Field Extraction Tool

# Tools

# Regex Web Page

http://regex101.com/

# Other relevant presentations

- **Become a regular expressions ninja and unlock your Splunk potential**
  - by *Gabriel Vasseur*

# Splunk Answers and Docs

Learn from others – ask questions – get answers

- http://answers.splunk.com/

Splunk Documentation

- https://docs.splunk.com/Documentation/Splunk/6.4.3/Knowledge/AboutSplunkregularexpressions

# Questions?

# Contact Me

- cary.petterborg@ldschurch.org

- carypetterborg@gmail.com

-  Channel: **The Splunk Hacker**

THANK YOU

.conf2016

splunk>

# Optional Content: Multiple dates with the same name

.conf2016

splunk>

# Social Media Feed Problem

Props.conf file must get the right timestamp, which is very difficult

- Twitter feeds in JSON format have multiple **postedTime** fields
- The first **postedTime** isn't necessarily the most recent
- You can't use the **postedTime** JSON key to determine the right timestamp
- You can't rely on the first timestamp to be the right one

# What do you do?

# Original twitter feed:

{"id":"tag:search.twitter.com,2005:754700573774151685","objectType":"activity","actor":{"objectType":"person","id":"id:twitter.com:4889755594","link":"http://www.twitter.com/searcherdonate","displayName":"Вероника Гусева","**postedTime**":"2016-02-08T20:33:04.000Z","image":"https://abs.twimg.com/sticky/default_profile_images/default_profile_3_normal.png","summary":null,"links":[{"href":null,"rel":"me"}],"friendsCount":3,"followersCount":81,"listedCount":27,"statusesCount":26831,"twitterTimeZone":"Pacific Time (US & Canada)","verified":false,"utcOffset":"-25200","preferredUsername":"searcherdonate","languages":["ru"],"favoritesCount":0},"verb":"post","**postedTime**":"2016-07-17T15:33:33.000Z","generator":{"displayName":"searcherdonate","link":"http://searcherdonate.com"},"provider":{"objectType":"service","displayName":"Twitter","link":"http://www.twitter.com"},"link":"http://twitter.com/searcherdonate/statuses/754700573774151685","body":"Donate Car for Tax Credit\n\nhttps://t.co/0XtalQOyZa https://t.co/YRGGzqgIRg","object":{"objectType":"note","id":"object:search.twitter.com,2005:754700573774151685","summary":"Donate Car for Tax Credit\n\nhttps://t.co/0XtalQOyZa https://t.co/YRGGzqgIRg","link":"http://twitter.com/searcherdonate/statuses/754700573774151685","**postedTime**":"2016-07-17T15:33:33.000Z"},"favoritesCount":0,"twitter_entities":{"hashtags":[],"urls":[{"url":"https://t.co/0XtalQOyZa","expanded_url":"http://searcherdonate.com/2016/07/17/donate-car-for-tax-credit-475/","display_url":"searcherdonate.com/2016/07/17/don…","indices":[27,50]}],"user_mentions":[],"symbols":[],"media":[{"id":754700571177783296,"id_str":"754700571177783296","indices":[51,74],"media_url":"http://pbs.twimg.com/media/Cnk8K8DWEAAU9Zw.jpg","media_url_https":"https://pbs.twimg.com/media/Cnk8K8DWEAAU9Zw.jpg","url":"https://t.co/YRGGzqgIRg","display_url":"pic.twitter.com/YRGGzqgIRg","expanded_url":"http://twitter.com/searcherdonate/status/754700573774151685/photo/1","type":"photo","sizes":{"medium":{"w":480,"h":360,"resize":"fit"},"small":{"w":480,"h":360,"resize":"fit"},"thumb":{"w":150,"h":150,"resize":"crop"},"large":{"w":480,"h":360,"resize":"fit"}}}]},"twitter_extended_entities":{"media":[{"id":754700571177783296,"id_str":"754700571177783296","indices":[51,74],"media_url":"http://pbs.twimg.com/media/Cnk8K8DWEAAU9Zw.jpg","media_url_https":"https://pbs.twimg.com/media/Cnk8K8DWEAAU9Zw.jpg","url":"https://t.co/YRGGzqgIRg","display_url":"pic.twitter.com/YRGGzqgIRg","expanded_url":"http://twitter.com/searcherdonate/status/754700573774151685/photo/1","type":"photo","sizes":{"medium":{"w":480,"h":360,"resize":"fit"},"small":{"w":480,"h":360,"resize":"fit"},"thumb":{"w":150,"h":150,"resize":"crop"},"large":{"w":480,"h":360,"resize":"fit"}}}]},"twitter_filter_level":"low","twitter_lang":"en","retweetCount":0,"gnip":{"matching_rules":[{"value":"Donate","tag":null}],"urls":[{"url":"https://t.co/0XtalQOyZa","expanded_url":"http://searcherdonate.com/2016/07/17/donate-car-for-tax-credit-475/","expanded_status":200},{"url":"https://t.co/YRGGzqgIRg","expanded_url":"http://twitter.com/searcherdonate/status/754700573774151685/photo/1","expanded_status":200}],"klout_score":42,"language":{"value":"en"}}}}

splunk> .conf2016

# Examination of dates:

Note newest date for postedTime:

- "displayName":"Вероника Гусева","**postedTime**":"2016-02-08T20:33:04.000Z"

- "verb":"post","**postedTime**":"**2016-07-17T15:33:33.000Z**"

- "link":"http://twitter.com/searcherdonate/statuses/754700573774151685","**postedTime**":"2016-07-17T15:33:33.000Z"

Before latest **postedTime**  comes **"verb":"post"**

# Original tweet feed:

{"id":"tag:search.twitter.com,2005:754693294882295810","objectType":"activity","actor":{"objectType":"person","id":"id:twitter.com:2987589649","link":"http://www.twitter.com/lauramessner14","displayName":"laura.the.escaper","**postedTime**":"2015-01-17T22:22:12.000Z","image":"https://pbs.twimg.com/profile_images/713992821766758400/3HDVKsXR_normal.jpg","summary":"forever escaper • MDE is my Everything • Chris Ryan is the cutest • YouTube is my life • Amanda is my potato •  [follow me on insta laura_messner14]","links":[{"href":null,"rel":"me"}],"friendsCount":1534,"followersCount":1281,"listedCount":24,"statusesCount":55756,"twitterTimeZone":null,"verified":false,"utcOffset":null,"preferredUsername":"lauramessner14","languages":["en"],"favoritesCount":3244},"verb":"share","**postedTime**":"2016-07-17T15:04:37.000Z","generator":{"displayName":"Twitter for iPhone","link":"http://twitter.com/download/iphone"},"provider":{"objectType":"service","displayName":"Twitter","link":"http://www.twitter.com"},"link":"http://twitter.com/lauramessner14/statuses/754693294882295810","body":"RT @VansWarpedTour: Holmdel, donate 3 canned goods, $5, or a used cell phone at the @feedourchildren tent and get express entry! https://t.…","object":{"id":"tag:search.twitter.com,2005:754646851442253824","objectType":"activity","actor":{"objectType":"person","id":"id:twitter.com:17503591","link":"http://www.twitter.com/VansWarpedTour","displayName":"Vans Warped Tour","**postedTime**":"2008-11-20T03:36:43.000Z","image":"https://pbs.twimg.com/profile_images/638426457211965442/-WWswoJQ_normal.jpg","summary":"On the road through August 13th. Tickets available now and complete lineup details, Warped 101 and more can be found on the Warped site","links":[{"href":"http://www.vanswarpedtour.com","rel":"me"}],"friendsCount":118,"followersCount":939364,"listedCount":3417,"statusesCount":24198,"twitterTimeZone":"Pacific Time (US & Canada)","verified":true,"utcOffset":"-25200","preferredUsername":"VansWarpedTour","languages":["en"],"favoritesCount":1912},"verb":"post","**postedTime**":"2016-07-17T12:00:04.000Z","generator":{"displayName":"Twitter Ads","link":"https://ads.twitter.com"},"provider":{"objectType":"service","displayName":"Twitter","link":"http://www.twitter.com"},"link":"http://twitter.com/VansWarpedTour/statuses/754646851442253824","body":"Holmdel, donate 3 canned goods, $5, or a used cell phone at the @feedourchildren tent and get express entry! https://t.co/n6UrMrM6Oq","object":{"objectType":"note","id":"object:search.twitter.com,2005:754646851442253824","summary":"Holmdel, donate 3 canned goods, $5, or a used cell phone at the @feedourchildren tent and get express entry! https://t.co/n6UrMrM6Oq","link":"http://twitter.com/VansWarpedTour/statuses/754646851442253824","**postedTime**":"2016-07-17T12:00:04.000Z"},"favoritesCount":221,"twitter_entities":{"hashtags":[],"urls":[{"url":"https://t.co/n6UrMrM6Oq","expanded_url":"http://smarturl.it/FeedOurChildrenNow","display_url":"smarturl.it/FeedOurChildre…","indices":[109,132]}],"user_mentions":[{"screen_name":"FeedOurChildren","name":"FeedOurChildrenNOW!","id":116930972,"id_str":"116930972","indices":[64,80]}],"symbols":[]},"twitter_filter_level":"low","twitter_lang":"en","favoritesCount":0,"twitter_entities":{"hashtags":[],"urls":[{"url":"https://t.co/n6UrMrM6Oq","expanded_url":"http://smarturl.it/FeedOurChildrenNow","display_url":"smarturl.it/FeedOurChildre…","indices":[139,140]}],"user_mentions":[{"screen_name":"VansWarpedTour","name":"Vans Warped Tour","id":17503591,"id_str":"17503591","indices":[3,18]},{"screen_name":"FeedOurChildren","name":"FeedOurChildrenNOW!","id":116930972,"id_str":"116930972","indices":[84,100]}],"symbols":[]},"twitter_filter_level":"low","twitter_lang":"en","retweetCount":61,"gnip":{"matching_rules":[{"value":"Donate","tag":null}],"urls":[{"url":"https://t.co/n6UrMrM6Oq","expanded_url":"https://www.facebook.com/warpedtour/photos/p.10153673105303435/10153673105303435/?type=3&theater","expanded_status":403}],"klout_score":43,"language":{"value":"en"}}}}

# Examination of dates:

Note newest date for postedTime:

- "displayName":"laura.the.escaper","**postedTime**":"2015-01-17T22:22:12.000Z"
- "verb":"share","**postedTime**":"**2016-07-17T15:04:37.000Z**"
- ,"displayName":"Vans Warped Tour","**postedTime**":"2008-11-20T03:36:43.000Z"
- "verb":"post","**postedTime**":"2016-07-17T12:00:04.000Z"
- "link":"http://twitter.com/VansWarpedTour/statuses/754646851442253824","**postedTime**":"2016-07-17T12:00:04.000Z"

All retweeted items have **"verb":"share"** before **postedTime**, which is the latest time for retweets

# Sed – the stream editor

```
sed \
-e 's/"verb":"share","postedTime":"/"verb":"share","truePostedTime":"/' \
-e '/truePostedTime/!s/"verb":"post","postedTime":"/"verb":"post","truePostedTime":"/' \
-e '/^$/d'
```

- You need to change the right postedTime to another name
  - **postedTime** -> **truePostedTime**
- Modify for the retweeted **postedTime** first
- Then if there is no **truePostedTime**, modify for original post
- Then get rid of blank lines (not discussed, but helpful)

# Use Unbuffered Output

- Buffered output will break lines at something like 4096 bytes.

- This will make Splunk break events up improperly

- Definitely breaks JSON formatting

- Use unbuffered output to make the entire event be written/output

# Curl & Sed Unbuffered

Shortened to make it easier to read:

```
curl -N -s --compressed –u$USERINFO $GNIPSTREAM | \
sed –u \
 -e \ 's/"share","postedTime"/"share","truePostedTime"/' -e \
'/truePostedTime/!s/"post","postedTime"/"post","truePostedTime"/' -e '/^$/d' \
>>/app/gnip/log/inputstream.log 2>>/app/gnip/log/errors
```