# Data Science on Hadoop

Justin Erickson
Senior Director, Product Management

# Age of Machine Learning



NO Machine Learning

Machine Learning

Data volume

Cost of compute

Time

1950s  1960s  1970s  1980s  1990s  2000s  2010s

2

# The Enterprise Platform for
# Data Science and Machine Learning

**440x**

MORE DATA

**30B**

CONNECTED DEVICES

The data is now here



APACHE
**Spark**™

**cloudera**®

Modern Platform for Machine Learning and Advanced Analytics
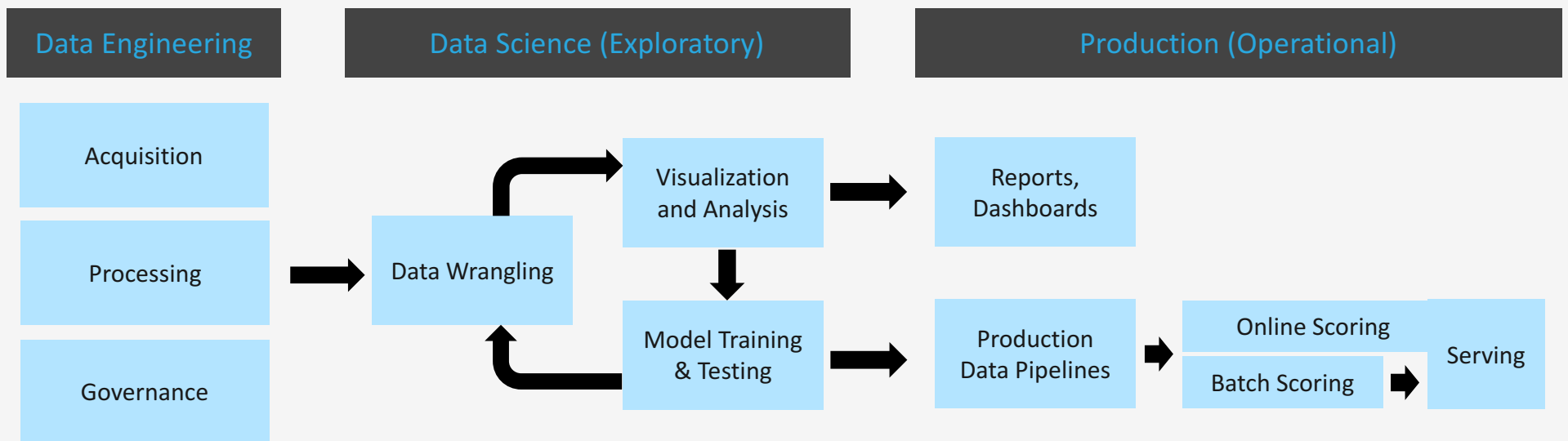
(intel)

Cloudera first to integrate Spark

**500**

Customers
Run Spark on

**cloudera**®

Leading adoption among enterprises

# Sample data science / machine learning workflow
From data to exploration to action

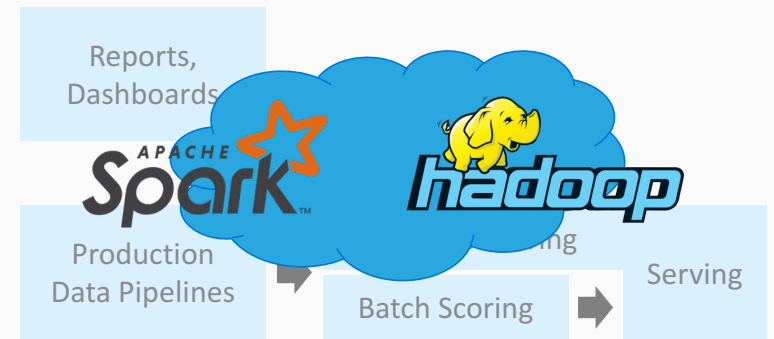| Data Engineering | Data Science (Exploratory) | Production (Operational) |
|---|---|---|

Acquisition

Processing

Governance

Data Wrangling

Visualization and Analysis

Model Training & Testing

Reports, Dashboards

Production Data Pipelines

Online Scoring

Batch Scoring

Serving

cloudera

# The good news



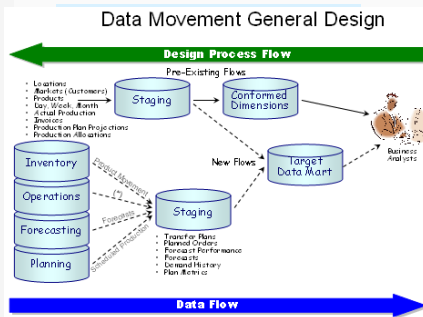| Data Engineering | Data Science (Exploratory) | Production (Operational) |
|---|---|---|
| Data has never been more plentiful | Open source data science and machine learning libraries are rapidly evolving | Commodity (and on-demand) compute makes scalable production machine learning affordable |

## cloudera

# The bad news

| Data Engineering | Data Science (Exploratory) | Production (Operational) |
|---|---|---|



Governance





Model Tr... & Testing

Data Pipelines



Batch Scoring    Serving

Data needs to move across multiple different systems

Teams have different, conflicting requests for languages & libraries

Most data science done at small scale, individually, and is difficult to replicate

Very few models reach production

# Additional challenges

## Access

For sensitive data, secure clusters are difficult to access. And IT typically doesn't want random packages installed on a secure cluster.

Popular open source tools don't easily connect to these environments, or always support Hadoop data formats.
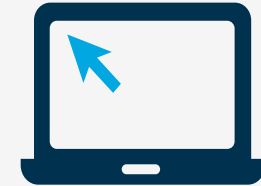
## Scale

Laptops rarely have capacity for medium, let alone big data. This leads to a lot of sampling.

Popular frameworks don't easily parallelize on a cluster. Typically code has to get rewritten for production.

## Developer Experience

Notebooks, while awesome, don't easily support virtual environment and dependency management, especially for teams. This makes sharing and reproducibility hard.
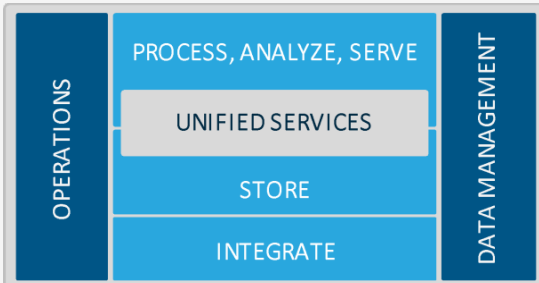
Notebooks are also challenging to "put into production."

**cloudera**

This year, our goal is to enable data science and machine learning at scale.

# Open data science in the enterprise

**Data Scientist**
explore, experiment, iterate

**IT**
drive adoption while
maintaining compliance

# Our goal: An open platform for data science at scale

**Help more data scientists use the power of Hadoop**

Use a powerful, familiar environment with direct access to Hadoop data and compute

Data Scientist
Data Engineer

**Make it easy and secure to add new users, use cases**

Offer secure self-service analytics and a faster path to production on common, affordable infrastructure

Enterprise Architect
Hadoop Admin

cloudera

# Introducing Cloudera Data Science Workbench
Self-service data science for the enterprise

Accelerates data science from development to production with:

- Secure self-service environments for data scientists to work against Cloudera clusters
- Support for Python, R, and Scala, plus project dependency isolation for multiple library versions
- Workflow automation, version control, collaboration and sharing

# Demo

cloudera

# With Cloudera Data Science Workbench…

## Data scientists can:

- Use R, Python, or Scala from a web browser, with no desktop footprint
- Install any library or framework within isolated project environments
- Directly access data in secure clusters with Spark and Impala
- Share insights with their team for reproducible, collaborative research
- Automate and monitor data pipelines using built-in job scheduling

## IT can:

- Give their data science team the freedom to work how they want, when they want
- Stay compliant with out-of-the-box support for full platform security, especially Kerberos
- Run on-premises or in the cloud, wherever data is managed

**cloudera**

# Solving Data Science is a Full-Stack Problem

- Support unlimited data
- Provide sufficient tools for Analysts
- Provide sufficient tools for Data Scientists + Data Engineers
- Enable real-time use cases
- Provide data governance
- Provide full-stack security
- Deploy in the cloud
- Integrate with partner tools
- Be easy for IT to deploy/maintain

- ✓ Hadoop
- ✓ Impala, Hive, Hue
- ✓ Spark, Data Science Workbench

- ✓ Kafka, Spark Streaming
- ✓ Navigator + Partners
- ✓ Kerberos, Sentry, Record Service, KMS/KTS
- ✓ Cloudera Director
- ✓ Rich Ecosystem
- ✓ Cloudera Manager + Director

**cloudera**

# The importance of an open ecosystem



Open Ecosystem



Black Box

FLIGHT RECORDER DO NOT OPEN

**cloudera**
# Thank You

Justin Erickson

Project quick find

**Projects**
**Jobs**
**Sessions**
**Settings**
**Admin**

| 0 | 2 | 3 | 6 |
|---|---|---|---|
| sessions running | jobs running | 0    vCPU    80.00 | 0    GB    263.86 |

## Projects

Creator ▾    **+ New Project**

### 👥 Product Overview
Introduction to Cloudera Data Science Workbench
● By **Matt Brandwein**. Last worked on just now.  ⑂ forked from **Product Demo**

**1**
running

### 👥 Data Analysis in Python
● By **Matt Brandwein**. Last worked on just now.

**1**
running

### 🔒 Health Data Demo
By **Matt Brandwein**. Last worked on 14 minutes ago.  ⑂ forked from **Health Data Demo**

0
running

### 🔒 Intel BigDL Experiments
By **Matt Brandwein**. Last worked on 2 weeks ago.

0
running

### 🔒 Deep Learning with TensorFlow
By **Matt Brandwein**. Last worked on 3 weeks ago.

0
running

### 👥 RImpala
By **Tristan Zajonc**. Last worked on 5 weeks ago.

0
running

### 👥 R Analysis
By **Tristan Zajonc**. Last worked on January 24.

0
running

### 👥 HDFS IO
Reading and writing data from HDFS.
By **Tristan Zajonc**. Last worked on January 23.

0
running

Q Project quick find        +    🐞 DST ▾    ▦

**Overview**

**Jobs**

**Sessions**

**Files**

**Team**

**Settings**

## Product Overview 👥

Introduction to Cloudera Data Science Workbench

ੴ forked from **Product Demo**

| 0 | Fork | ▭ Open Workbench |

## Jobs

Creator ▾

| Name | Runs / Failures | Duration | Status | Latest Run | Actions |
|------|-----------------|----------|--------|------------|---------|
| **Nightly Report** | 1 / 0 | 00:08 | Success | 7 minutes ago | ▶ |

## Files

⬇ Download  ✚ New  ⬆ Upload

| ☐ | Name ⌃ | Size | Last Modified |
|---|--------|------|---------------|
| ☐ | 📁 data | - | 17 hours ago |
| ☐ | 📁 img | - | 17 hours ago |
| ☐ | 📁 R | - | 17 hours ago |
| ☐ | 📄 1_python.py | 2.95 kB | 17 hours ago |
| ☐ | 📄 2_tensorflow.py | 3.38 kB | 17 hours ago |
| ☐ | 📄 3_sparklyr.R | 2.23 kB | 17 hours ago |
| ☐ | 📄 3a.R | 356 B | 17 hours ago |
| ☐ | 📄 README.md | 134 B | 17 hours ago |
| ☐ | 📄 utils.py | 1.04 kB | 17 hours ago |
| ☐ | 📄 utils.pyc | 1.43 kB | 17 hours ago |

2_tensorflow.py

```
1   import tensorflow as tf
2   import numpy as np
3   import matplotlib
4   from matplotlib import pyplot as plt
5   import utils
6
7   ### Import MNIST data
8   from tensorflow.examples.tutorials.mnist import input_data
9   mnist = input_data.read_data_sets('data/MNIST', one_hot=True)
10
11  ### View Data
12  for i in xrange(0, 3):
13    tmp = mnist.train.images[i]
14    tmp = tmp.reshape((28,28))
15    plt.imshow(tmp, cmap = cm.Greys)
16    plt.show()
17
18  ### Parameters
19  learning_rate = 0.01
20  training_epochs = 5
21  batch_size = 100
22  display_step = 1
23  logs_path = '/tmp/tensorboard'
24
25  ### Cleanup old logs
26  if tf.gfile.Exists(logs_path):
27    tf.gfile.DeleteRecursively(logs_path)
28  tf.gfile.MakeDirs(logs_path)
29
30  ### Model
31  # Use a single-layer perceptron as example $pred = softmax(W
32  x = tf.placeholder('float', [None, 784], name='data')
33  y = tf.placeholder('float', [None, 10], name='label')
34
35  # Model bias and weight variables: W, b
36  W = tf.Variable(tf.zeros([784,10]), name='weights')
37  b = tf.Variable(tf.zeros([10]), name='bias')
38
39  # Put the model ops into scopes for tensorboard
40  with tf.name_scope('Model'):
41      logits = tf.matmul(x,W)+b
42      pred = tf.nn.softmax( logits )
43  with tf.name_scope('Loss'):
44      cost = tf.reduce_mean( tf.nn.softmax_cross_entropy_with_lc
```

Line 1, Column 1                    ★    105 Lines    Python    Spaces 2

## Running Sessions

Python session, 1 vCPU (burstable), 2 GiB memory, on 3/9 at 10:09

## Start New Session

Select Engine Image

○ R (3.3.0-3)

⦿ Python (2.7.11-2)

○ Scala (2.11-5)

Select Engine Profile

2 vCPU (burstable), 4 GiB memory    ⬍

Launch Session

0.7.3 (b3c2553)

```python
24  display_charts(data, chart_type="stock", title="DJIA vs. Debt Query Volume", secon
25  seaborn.lmplot("debt", "djia", data=data, size=7)
26
27  # Detect if search volume is increasing or decreasing in
28  # any given week by forming a moving average and testing if the current value
29  # crosses the moving average of the past 3 weeks.
30  #
31  # Let's first compute the moving average.
32
33  data['debt_mavg'] = data.debt.rolling(window=3, center=False).mean()
34  data.head()
35
36  # Since we want to see if the current value is above the moving average of the
37  # *preceeding* weeks, we have to shift the moving average timeseries forward by one
38
39  data['debt_mavg'] = data.debt_mavg.shift(1)
40  data.head()
41
42  # Generate Orders
43  # ===============
44  #
45  # We use Google Trends to determine how many searches have been
46  # carried out for a specific search term such as debt in week,
47  # where Google defines weeks as ending on a Sunday, relative to the total
48  # number of searches carried out on Google during that time.
49  #
50  # We implement the strategy of selling when debt searchess exceed
51  # the moving average and buying when debt searchers fall below the moving
52  # average.
53
54  data['order'] = 0
55  data.loc[data.debt > data.debt_mavg, 'order'] = -1
56  data.loc[data.debt < data.debt_mavg, 'order'] = -1
57  data.head()
58
59  # Compute Returns
60  # ===============
61
62  data['ret_djia'] = data.djia.pct_change()
63  data.head()
64
65  # Returns at week `t` are relative to week `t-1`. However, we are buying at
66  # week `t` and selling at week `t+1`, so we have to adjust by shifting
67  # the returns upward.
68
69  data['ret_djia'] = data['ret_djia'].shift(-1)
70
71  # The algorithm that is used by the authors makes a decision every Monday of
```

> data.head()

|        | djia     | debt     |
|--------|----------|----------|
| **Date** |          |          |
| **2004-01-14** | 10485.18 | 0.210000 |
| **2004-01-22** | 10528.66 | 0.210000 |
| **2004-01-28** | 10702.51 | 0.210000 |
| **2004-02-04** | 10499.18 | 0.213333 |
| **2004-02-11** | 10579.03 | 0.200000 |

Show DJIA vs. debt related query volume.

> display_charts(data, chart_type="stock", title="DJIA vs. Debt Query Volume", secondary_y="debt")



DJIA vs. Debt Query Volume

Zoom  1m  3m  6m  YTD  1y  **All**        From  Jan 12, 2004   To  Mar 2, 2011

```
1   ## Connecting to Spark
2
3   library(sparklyr)
4   library(dplyr)
5
6   # The returned Spark connection (sc) provides a remote dplyr data source to the
7   sc
8
9   ## Using dplyr
10  # We can now use all of the available dplyr verbs against the tables within the
11
12  # # filter by departure delay
13  flights_tbl %>% filter(dep_delay == 2)
14
15  # Introduction to dplyr provides additional dplyr examples you can try. For exam
16  delay <- flights_tbl %>%
17    group_by(tailnum) %>%
18    summarise(count = n(), dist = mean(distance), delay = mean(arr_delay)) %>%
19    filter(count > 20, dist < 2000, !is.na(delay)) %>%
20    collect()
21
22  # # Plot delays
23  library(ggplot2)
24  ggplot(delay, aes(dist, delay)) +
25    geom_point(aes(size = count), alpha = 1/2) +
26    geom_smooth() +
27    scale_size_area(max_size = 2)
28
29  ## Machine Learning
30  # You can orchestrate machine learning algorithms in a Spark cluster via the mac
31
32  # In this example we'll use ml_linear_regression to fit a linear regression mode
33
34  # copy mtcars into spark
35  mtcars_tbl <- copy_to(sc, mtcars)
36
37  # transform our data set, and then partition into 'training', 'test'
38  partitions <- mtcars_tbl %>%
39    filter(hp >= 100) %>%
40    mutate(cyl8 = cyl == 8) %>%
41    sdf_partition(training = 0.5, test = 0.5, seed = 1099)
42
43  # fit a linear model to the training dataset
44  fit <- partitions$training %>%
45    ml_linear_regression(response = "mpg", features = c("wt", "cyl"))
46
47  # For linear regression models produced by Spark, we can use summary() to learn
```

from the tutorial which plots data on flight delays:

```
> delay <- flights_tbl %>%
    group_by(tailnum) %>%
    summarise(count = n(), dist = mean(distance), delay = mean(arr_delay)) %>%
    filter(count > 20, dist < 2000, !is.na(delay)) %>%
    collect()
```
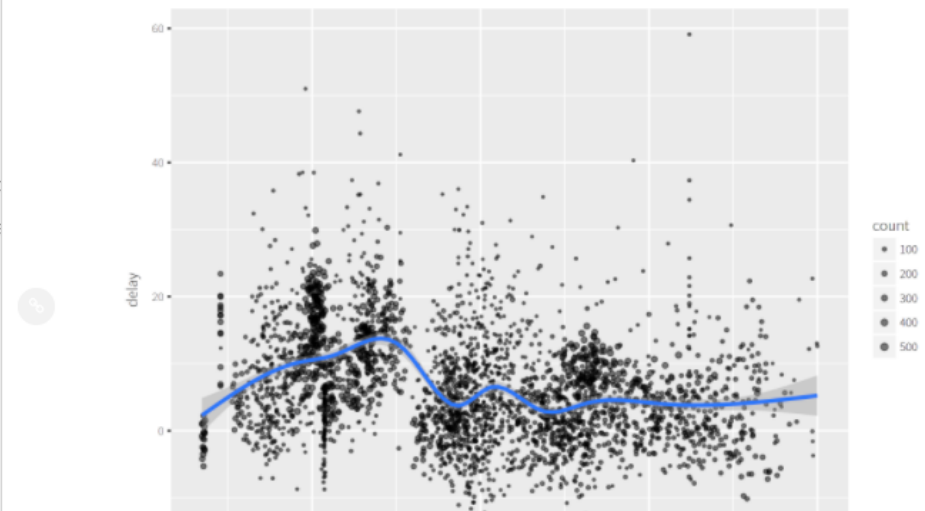
# Plot delays

```
> library(ggplot2)
> ggplot(delay, aes(dist, delay)) +
    geom_point(aes(size = count), alpha = 1/2) +
    geom_smooth() +
    scale_size_area(max_size = 2)
```

`geom_smooth()` using method = 'gam'

```
1   ## Connecting to Spark
2
3   library(sparklyr)
4   library(dplyr)
5
6   # The returned Spark connecti                                              n(arr_delay)) %>%
7   sc
8
9   ## Using dplyr
10  # We can now use all of the a
11
12  # # filter by departure dela
13  flights_tbl %>% filter(dep_de
14
15  # Introduction to dplyr provi
16  delay <- flights_tbl %>%
17    group_by(tailnum) %>%
18    summarise(count = n(), dist
19    filter(count > 20, dist < 2
20    collect()
21
22  # # Plot delays
23  library(ggplot2)
24  ggplot(delay, aes(dist, delay
25    geom_point(aes(size = count
26    geom_smooth() +
27    scale_size_area(max_size =
28
29  ## Machine Learning
30  # You can orchestrate machine
31
32  # In this example we'll use m
33
34  # copy mtcars into spark
35  mtcars_tbl <- copy_to(sc, mt
36
37  # transform our data set, and
38  partitions <- mtcars_tbl %>%
39    filter(hp >= 100) %>%
40    mutate(cyl8 = cyl == 8) %>%
41    sdf_partition(training = 0.
42
43  # fit a linear model to the t
44  fit <- partitions$training %>%
45    ml_linear_regression(respon
46
47  # ression model
```

Cloudera Data Science Workbench Terminal

ⓘ cdsw.edh.cloudera.com/terminal/9ixh0n8och0jtsum?p=21jtgimbk5mn7cxj

```
Cloudera Data Science Workbench Terminal
Your project files are located in /home/sense
sense@9ixh0n8och0jtsum:~$ ls
R  README.md  analysis.py  analysis.r  data  models
sense@9ixh0n8och0jtsum:~$ hdfs dfs -ls /
Found 15 items
drwxr-xr-x   - cops      edh-ingestion        0 2015-12-22 15:47 /backups
drwxrwxr-x   - hdfs      edh-ingestion        0 2015-11-05 15:51 /data
drwxr-xr-x   - impala    supergroup           0 2015-07-17 00:38 /datestamp=
drwxr-xr-x   - hdfs      supergroup           0 2014-12-19 15:37 /etc
drwx------+  - hbase     hbase                0 2017-01-23 17:10 /hbase
drwxr-xr-x   - cops      edh-ingestion        0 2016-08-04 15:50 /jobs
drwxr-xr-x   - hdfs      supergroup           0 2014-08-30 18:15 /jobtracker
drwxr-xr-x   - hdfs      supergroup           0 2017-01-31 21:02 /projects
drwxrwxr-x   - cops      edh-ingestion        0 2014-10-10 19:08 /schemas
drwxrwxr-x   - solr      edh-ingestion        0 2017-02-10 00:12 /solr
drwxr-xr-x   - hdfs      supergroup           0 2017-03-08 20:07 /system
-rw-r--r--   3 someuser  somegroup            0 2014-08-02 04:21 /test.txt
drwxrwxrwt   - hdfs      supergroup           0 2017-03-13 17:36 /tmp
drwxr-xr-x   - hdfs      supergroup           0 2017-03-13 15:00 /user
drwxr-xr-x   - impala    supergroup           0 2015-09-24 17:23 /users
sense@9ixh0n8och0jtsum:~$ pip list
abstract-rendering (0.5.1)
alabaster (0.7.3)
argcomplete (0.8.9)
astropy (1.0.3)
Babel (1.3)
backports.ssl-match-hostname (3.4.0.2)
bcolz (0.9.0)
beautifulsoup4 (4.3.2)
binstar (0.11.0)
bitarray (0.8.1)
blaze (0.8.0)
blz (0.6.2)
bokeh (0.9.0)
```
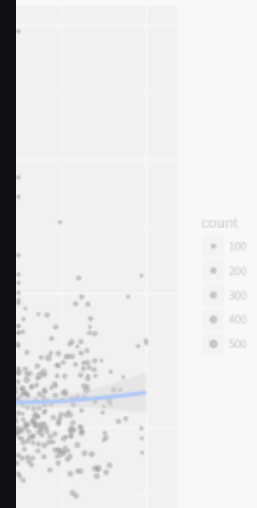
count
• 100
• 200
• 300
• 400
• 500

cloudera

Line 52, Column 2                              ★   52 Lines   R   Spaces 2

Project quick find

mbrandwein ▾

**Overview**

**Jobs**

**Sessions**

**Files**

**Team**

**Settings**

# Jobs

**+ New Job**

Job Dependencies for Mos

---

### Cluster Metadata

Success   Run

**Overview**    History    Dependencies    Settings

**Script:** bin/transformation.py
**Schedule:** after Most Recent Collection
**Engine Profile:** 1 vCPU (burstable), 2 GiB memory
**Created By:** Ricky Saltzer

**Latest Run:** 27 minutes ago
**Duration:** 00:59
**Runs:** 487
**Failures:** 2

**Job History**

Mar 6 18:00   Mar 7 00:00   Mar 7 06:00   Mar 7 12:00   Mar 7 18:00   Mar 8 00:00   Mar 8 06:00   Mar 8 12:00

---

Creator ▾

| Name | Runs / Failures | Duration | Status | Latest Run | Actions |
|------|-----------------|----------|--------|------------|---------|
| Cluster Metadata | 487 / 2 | 00:59 | Success | 27 minutes ago | ▶ |
| Most Recent Collection | 466 / 1 | 21:47 | Success | 49 minutes ago | ▶ |
| Unification | 342 / 0 | 00:11 | Success | 49 minutes ago | ▶ |
| Case Issue Clarification | 30 / 2 | 00:19 | Success | 10 hours ago | ▶ |
| SFDC Tasks | 16 / 1 | 00:51 | Success | 18 hours ago | ▶ |
| Cluster Configs | 15 / 0 | 00:34 | Success | 19 hours ago | ▶ |

0.7.3 (b3c2553)