# Gut Reaction

"Hey! I have a bunch of time series!"

# We have a bunch of… things!

"How are these things related?"



https://flic.kr/p/7VUJkB

# Pearson Correlation

Covariance of two variables divided by the product of their standard deviations

$$r = \frac{\sum(X - \overline{X})(Y - \overline{Y})}{\sqrt{\sum(X - \overline{X})^2}\sqrt{\sum(Y - \overline{Y})^2}}$$
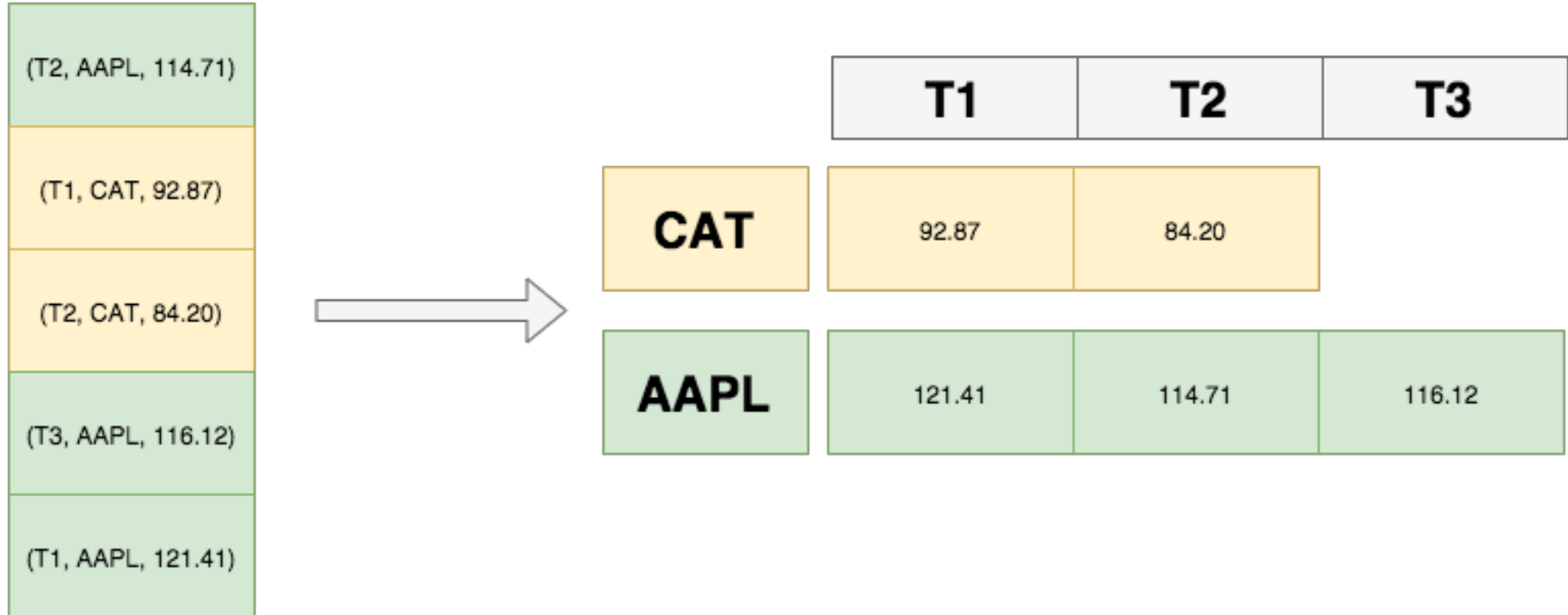
# Create a Spark RDD of observations

```
In [33]: wiki_df.head(5)

Out[33]: [Row(timestamp=datetime.datetime(2015, 8, 2, 17, 0), page=u'Morgan_Stanley', views=32.0),
          Row(timestamp=datetime.datetime(2015, 8, 2, 23, 0), page=u'Ford_Motor_Company', views=91.0),
          Row(timestamp=datetime.datetime(2015, 8, 3, 13, 0), page=u'Mattel', views=24.0),
          Row(timestamp=datetime.datetime(2015, 8, 3, 13, 0), page=u'Visa_Inc.', views=36.0),
          Row(timestamp=datetime.datetime(2015, 8, 3, 13, 0), page=u'Yum!_Brands', views=36.0)]
```

```
In [32]: ticker_df.head(5)

Out[32]: [Row(timestamp=datetime.datetime(2015, 9, 1, 6, 0), symbol=u'DOV', price=60.38130930615455),
          Row(timestamp=datetime.datetime(2015, 8, 7, 9, 0), symbol=u'BLL', price=68.8183106049063),
          Row(timestamp=datetime.datetime(2015, 9, 15, 6, 0), symbol=u'MDLZ', price=42.3693109177133
          1),
          Row(timestamp=datetime.datetime(2015, 8, 24, 12, 0), symbol=u'EMC', price=23.82834918849372
          5),
          Row(timestamp=datetime.datetime(2015, 8, 20, 11, 0), symbol=u'HRS', price=80.5737563743446
          1)]
```

**cloudera**

‹#›
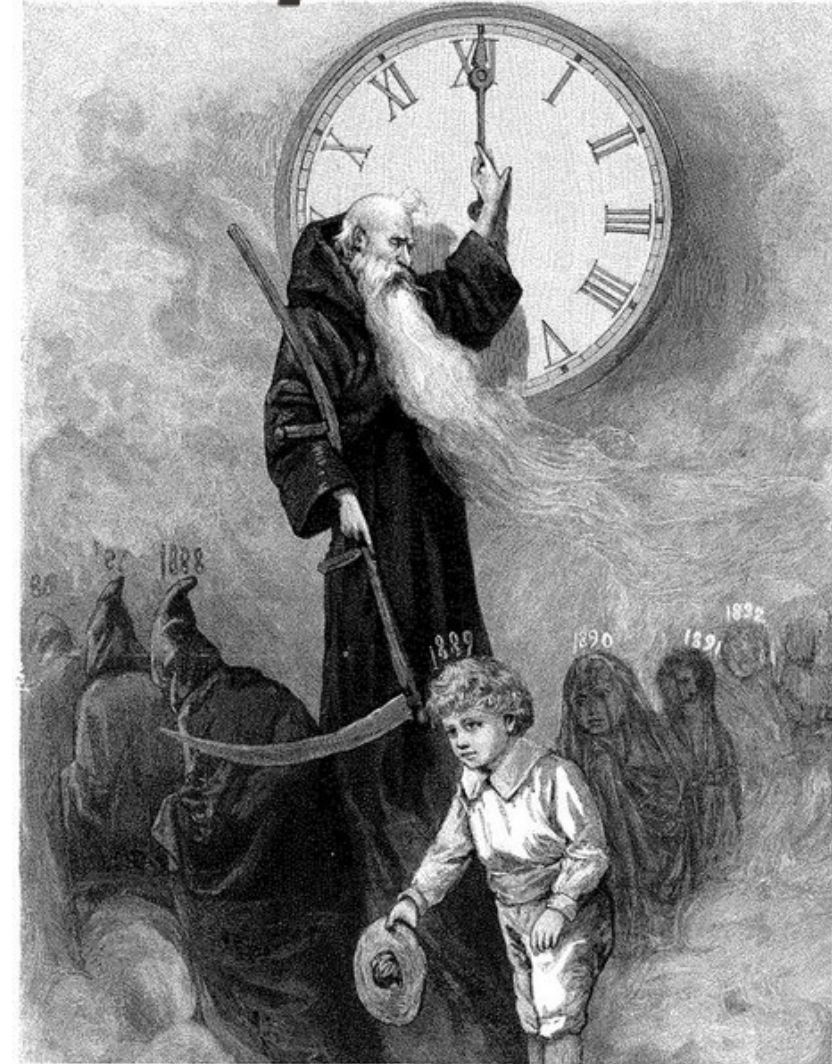
# A time series is more than a pile of observations.



| | T1 | T2 | T3 |
|------|-------|--------|--------|
| CAT | 92.87 | 84.20 | |
| AAPL | 121.41 | 114.71 | 116.12 |

(T2, AAPL, 114.71)

(T1, CAT, 92.87)

(T2, CAT, 84.20)

(T3, AAPL, 116.12)

(T1, AAPL, 121.41)

cloudera

# Spark-Timeseries

# Spark-Timeseries

- Library for manipulating and analyzing large-scale time series data

- Github: https://github.com/cloudera/spark-timeseries

- Doc: http://cloudera.github.io/spark-timeseries/

# Observations

| Timestamp | Symbol | Price |
|-----------|--------|-------|
| 2015-04-10 | AAPL | 2.0 |
| 2015-04-11 | AAPL | 3.0 |
| 2015-04-10 | MSFT | 4.5 |
| 2015-04-11 | MSFT | 1.5 |
| 2015-04-10 | GOOG | 6.0 |

```
schema = StructType([
    StructField('timestamp', TimestampType()), \
    StructField('symbol', StringType()), \
    StructField('price', DoubleType()) \
])
ticker_obs = sqlCtx.createDataFrame(row_rdd,
schema)
```

# Time samples

| Timestamp | AAPL | MSFT | GOOG |
|-----------|------|------|------|
| 2015-04-10 | 2.0 | 4.5 | 6.0 |
| 2015-04-11 | 3.0 | 1.5 | NaN |

# Time series

| DateTimeIndex: [2015-04-10, 2015-04-11] | |
|:---:|:---:|
| **Symbol** | **Series** |
| AAPL | [2.0, 3.0] |
| MSFT | [4.5, 1.5] |
| GOOG | [6.0, NaN] |

# Transform observations into time series

# Time to code



https://flic.kr/p/8RcFUD

# Apple



https://flic.kr/p/8d4XB3

# What's the correlation between AAPL and [en.wikipedia.org/wiki/Apple_Inc.](en.wikipedia.org/wiki/Apple_Inc.)?

Wiki Page:

# Data Quality and NaNs



https://flic.kr/p/bnnHaP

cloudera

# Link symbols and pages

# Ask focused questions

cloudera

# Volatility

Measure of variety or variation of a time series

**cloudera**

‹#›

# Magic

- Magnitude
- Articulation
- Generality
- Interestingness
- Credibility



https://flic.kr/p/a4h8WE

# Magic

- Magnitude
- Articulation
- Generality
- Interestingness
- Credibility

What is our question? How are we quantifying and measuring?

https://flic.kr/p/a4h8WE

# Statistical Plotting: Seaborn

# Distribution of Volatility

cloudera

# Regression as a measure of effect

# Linear Regression using sklearn



http://scikit-learn.org/stable/documentation.html

# Linear Regression using sklearn

```python
from sklearn import linear_model

def regress(X, y):
    model = linear_model.LinearRegression()
    model.fit(X, y)
    score = model.score(X, y)
    return (score, model)


lag = 2
lead = 2


joined = regressions = wiki_daily_views.flatMap(get_page_symbol) \
    .join(ticker_daily_vol)

models = joined.mapValues(lambda x: regress(lead_and_lag(lead, lag, x[0]), x[1][lag:-lead]))
models.cache()
models.count()
```
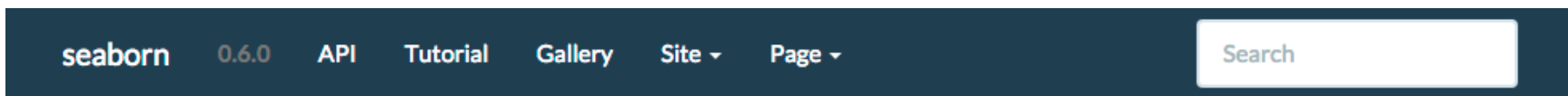
# Linear Regression using seaborn



‹#›

# Why regression when we have correlation?



https://flic.kr/p/8mNZ6r

# If you need to fit a linear model over all of your data

## pyspark.mllib.classification module

*class* `pyspark.mllib.classification.`**`LogisticRegressionModel`**(*weights, intercept, numFeatures, numClasses*)  [source]

Classification model trained using Multinomial/Binary Logistic Regression.

**Parameters:**
- **weights** – Weights computed for every feature.
- **intercept** – Intercept computed for this model. (Only used in Binary Logistic Regression. In Multinomial Logistic Regression, the intercepts will not be a single value, so the intercepts will be part of the weights.)
- **numFeatures** – the dimension of the features.
- **numClasses** – the number of possible outcomes for k classes classification problem in Multinomial Logistic Regression. By default, it is binary logistic regression so numClasses will be set to 2.

```
>>> data = [
...     LabeledPoint(0.0, [0.0, 1.0]),
...     LabeledPoint(1.0, [1.0, 0.0]),
... ]
>>> lrm = LogisticRegressionWithSGD.train(sc.parallelize(data), iterations=10)
>>> lrm.predict([1.0, 0.0])
1
>>> lrm.predict([0.0, 1.0])
0
>>> lrm.predict(sc.parallelize([[1.0, 0.0], [0.0, 1.0]])).collect()
[1, 0]
>>> lrm.clearThreshold()
>>> lrm.predict([0.0, 1.0])
0.279...
```

https://flic.kr/p/8mNZ6r

# Identifying Outliers

**cloudera**

# When are "unusual" events occurring?



https://flic.kr/p/bn9tsC

**cloudera**

‹#›

# Tukey's criterion

- Tukey's boxplot criterion for outlier identification
  - Non-parametric: doesn't need to assume particular probability distribution over daily purchases
  - *Robust to outliers*: focusing on interquartile range means we sidestep effects of anomalies



Upper Whisker

Q3: Upper Quartile

Median

Q1: Lower Quartile

Lower Whisker

cloudera

# Black Monday
August 24th, 2015

stock tickers

Investing Guide

# Trading was halted 1,200 times Monday

Recommend 3k

By Matt Egan   @mattmegan5

🏠 / Black Monday vs. the Flash Crash and the SEC Response

**NEWS**

## Black Monday vs. the Flash Crash and the SEC Response

Calder Lamb   Sep 23, 2015

Share

**TRADER TALK**

# What happened during the Aug 24 'flash crash'

Bob Pisani | @BobPisani
Friday, 25 Sep 2015 | 3:59 PM ET

**CNBC**

# What happened?

Which stocks had the largest overnight drop
from end of August 23rd to open August 24th?

Which stocks saw the most volatility that day?

Hue:  cdh4.srowen.com:8888  srowen/srowen
ipython:  cdh1.srowen.com:8880  user11.ipynb
Spark: cdh1.srowen.com:4040 (redirects)
HDFS: /user/srowen/user11

Hue:  cdh4.srowen.com:8888  srowen/srowen
ipython:  cdh1.srowen.com:8880  user12.ipynb
Spark: cdh1.srowen.com:4041 (redirects)
HDFS: /user/srowen/user12

Hue:  cdh4.srowen.com:8888  srowen/srowen
ipython:  cdh1.srowen.com:8880  user13.ipynb
Spark: cdh1.srowen.com:4042 (redirects)
HDFS: /user/srowen/user13

Hue:  cdh4.srowen.com:8888  srowen/srowen
ipython:  cdh1.srowen.com:8880  user14.ipynb
Spark: cdh1.srowen.com:4043 (redirects)
HDFS: /user/srowen/user14

Hue:  cdh4.srowen.com:8888  srowen/srowen
ipython:  cdh1.srowen.com:8880  user15.ipynb
Spark: cdh1.srowen.com:4044 (redirects)
HDFS: /user/srowen/user15

Hue:  cdh4.srowen.com:8888  srowen/srowen
ipython:  cdh1.srowen.com:8880  user16.ipynb
Spark: cdh1.srowen.com:4045 (redirects)
HDFS: /user/srowen/user16

Hue:  cdh4.srowen.com:8888  srowen/srowen
ipython:  cdh1.srowen.com:8880  user17.ipynb
Spark: cdh1.srowen.com:4046 (redirects)
HDFS: /user/srowen/user17

Hue:  cdh4.srowen.com:8888  srowen/srowen
ipython:  cdh1.srowen.com:8880  user18.ipynb
Spark: cdh1.srowen.com:4047 (redirects)
HDFS: /user/srowen/user18

Hue:   cdh4.srowen.com:8888   srowen/srowen
ipython:   cdh2.srowen.com:8880   user21.ipynb
Spark: cdh2.srowen.com:4040 (redirects)
HDFS: /user/srowen/user21

Hue:   cdh4.srowen.com:8888   srowen/srowen
ipython:   cdh2.srowen.com:8880   user22.ipynb
Spark: cdh2.srowen.com:4041 (redirects)
HDFS: /user/srowen/user22

Hue:   cdh4.srowen.com:8888   srowen/srowen
ipython:   cdh2.srowen.com:8880   user23.ipynb
Spark: cdh2.srowen.com:4042 (redirects)
HDFS: /user/srowen/user23

Hue:   cdh4.srowen.com:8888   srowen/srowen
ipython:   cdh2.srowen.com:8880   user24.ipynb
Spark: cdh2.srowen.com:4043 (redirects)
HDFS: /user/srowen/user24

Hue:   cdh4.srowen.com:8888   srowen/srowen
ipython:   cdh2.srowen.com:8880   user25.ipynb
Spark: cdh2.srowen.com:4044 (redirects)
HDFS: /user/srowen/user25

Hue:   cdh4.srowen.com:8888   srowen/srowen
ipython:   cdh2.srowen.com:8880   user26.ipynb
Spark: cdh2.srowen.com:4045 (redirects)
HDFS: /user/srowen/user26

Hue:   cdh4.srowen.com:8888   srowen/srowen
ipython:   cdh2.srowen.com:8880   user27.ipynb
Spark: cdh2.srowen.com:4046 (redirects)
HDFS: /user/srowen/user27

Hue:   cdh4.srowen.com:8888   srowen/srowen
ipython:   cdh2.srowen.com:8880   user28.ipynb
Spark: cdh2.srowen.com:4047 (redirects)
HDFS: /user/srowen/user28

*Hue:  cdh4.srowen.com:8888  srowen/srowen*
*ipython:  cdh3.srowen.com:8880  user31.ipynb*
*Spark: cdh3.srowen.com:4040 (redirects)*
*HDFS: /user/srowen/user31*

*Hue:  cdh4.srowen.com:8888  srowen/srowen*
*ipython:  cdh3.srowen.com:8880  user32.ipynb*
*Spark: cdh3.srowen.com:4041 (redirects)*
*HDFS: /user/srowen/user32*

*Hue:  cdh4.srowen.com:8888  srowen/srowen*
*ipython:  cdh3.srowen.com:8880  user33.ipynb*
*Spark: cdh3.srowen.com:4042 (redirects)*
*HDFS: /user/srowen/user33*

*Hue:  cdh4.srowen.com:8888  srowen/srowen*
*ipython:  cdh3.srowen.com:8880  user34.ipynb*
*Spark: cdh3.srowen.com:4043 (redirects)*
*HDFS: /user/srowen/user34*

*Hue:  cdh4.srowen.com:8888  srowen/srowen*
*ipython:  cdh3.srowen.com:8880  user35.ipynb*
*Spark: cdh3.srowen.com:4044 (redirects)*
*HDFS: /user/srowen/user35*

*Hue:  cdh4.srowen.com:8888  srowen/srowen*
*ipython:  cdh3.srowen.com:8880  user36.ipynb*
*Spark: cdh3.srowen.com:4045 (redirects)*
*HDFS: /user/srowen/user36*

*Hue:  cdh4.srowen.com:8888  srowen/srowen*
*ipython:  cdh3.srowen.com:8880  user37.ipynb*
*Spark: cdh3.srowen.com:4046 (redirects)*
*HDFS: /user/srowen/user37*

*Hue:  cdh4.srowen.com:8888  srowen/srowen*
*ipython:  cdh3.srowen.com:8880  user38.ipynb*
*Spark: cdh3.srowen.com:4047 (redirects)*
*HDFS: /user/srowen/user38*

| | |
|---|---|
| Hue:   cdh4.srowen.com:8888   srowen/srowen<br>ipython:   cdh4.srowen.com:8880   user41.ipynb<br>Spark: cdh4.srowen.com:4040 (redirects)<br>HDFS: /user/srowen/user41 | Hue:   cdh4.srowen.com:8888   srowen/srowen<br>ipython:   cdh4.srowen.com:8880   user42.ipynb<br>Spark: cdh4.srowen.com:4041 (redirects)<br>HDFS: /user/srowen/user42 |
| Hue:   cdh4.srowen.com:8888   srowen/srowen<br>ipython:   cdh4.srowen.com:8880   user43.ipynb<br>Spark: cdh4.srowen.com:4042 (redirects)<br>HDFS: /user/srowen/user43 | Hue:   cdh4.srowen.com:8888   srowen/srowen<br>ipython:   cdh4.srowen.com:8880   user44.ipynb<br>Spark: cdh4.srowen.com:4043 (redirects)<br>HDFS: /user/srowen/user44 |
| Hue:   cdh4.srowen.com:8888   srowen/srowen<br>ipython:   cdh4.srowen.com:8880   user45.ipynb<br>Spark: cdh4.srowen.com:4044 (redirects)<br>HDFS: /user/srowen/user45 | Hue:   cdh4.srowen.com:8888   srowen/srowen<br>ipython:   cdh4.srowen.com:8880   user46.ipynb<br>Spark: cdh4.srowen.com:4045 (redirects)<br>HDFS: /user/srowen/user46 |
| Hue:   cdh4.srowen.com:8888   srowen/srowen<br>ipython:   cdh4.srowen.com:8880   user47.ipynb<br>Spark: cdh4.srowen.com:4046 (redirects)<br>HDFS: /user/srowen/user47 | Hue:   cdh4.srowen.com:8888   srowen/srowen<br>ipython:   cdh4.srowen.com:8880   user48.ipynb<br>Spark: cdh4.srowen.com:4047 (redirects)<br>HDFS: /user/srowen/user48 |

Hue:   cdh4.srowen.com:8888   srowen/srowen
ipython:   cdh5.srowen.com:8880   user51.ipynb
Spark: cdh5.srowen.com:4040 (redirects)
HDFS: /user/srowen/user51

Hue:   cdh4.srowen.com:8888   srowen/srowen
ipython:   cdh5.srowen.com:8880   user52.ipynb
Spark: cdh5.srowen.com:4041 (redirects)
HDFS: /user/srowen/user52

Hue:   cdh4.srowen.com:8888   srowen/srowen
ipython:   cdh5.srowen.com:8880   user53.ipynb
Spark: cdh5.srowen.com:4042 (redirects)
HDFS: /user/srowen/user53

Hue:   cdh4.srowen.com:8888   srowen/srowen
ipython:   cdh5.srowen.com:8880   user54.ipynb
Spark: cdh5.srowen.com:4043 (redirects)
HDFS: /user/srowen/user54

Hue:   cdh4.srowen.com:8888   srowen/srowen
ipython:   cdh5.srowen.com:8880   user55.ipynb
Spark: cdh5.srowen.com:4044 (redirects)
HDFS: /user/srowen/user55

Hue:   cdh4.srowen.com:8888   srowen/srowen
ipython:   cdh5.srowen.com:8880   user56.ipynb
Spark: cdh5.srowen.com:4045 (redirects)
HDFS: /user/srowen/user56

Hue:   cdh4.srowen.com:8888   srowen/srowen
ipython:   cdh5.srowen.com:8880   user57.ipynb
Spark: cdh5.srowen.com:4046 (redirects)
HDFS: /user/srowen/user57

Hue:   cdh4.srowen.com:8888   srowen/srowen
ipython:   cdh5.srowen.com:8880   user58.ipynb
Spark: cdh5.srowen.com:4047 (redirects)
HDFS: /user/srowen/user58