

OmniSci and RAPIDS: An End-to-End GPU Data Science Workflow

ODSC Webinar - May 30, 2019





Randy Zwitch

Senior Director of Developer Advocacy

 @randyzwitch

 randy.zwitch@omnisci.com

 /in/randyzwitch/

 /randyzwitch



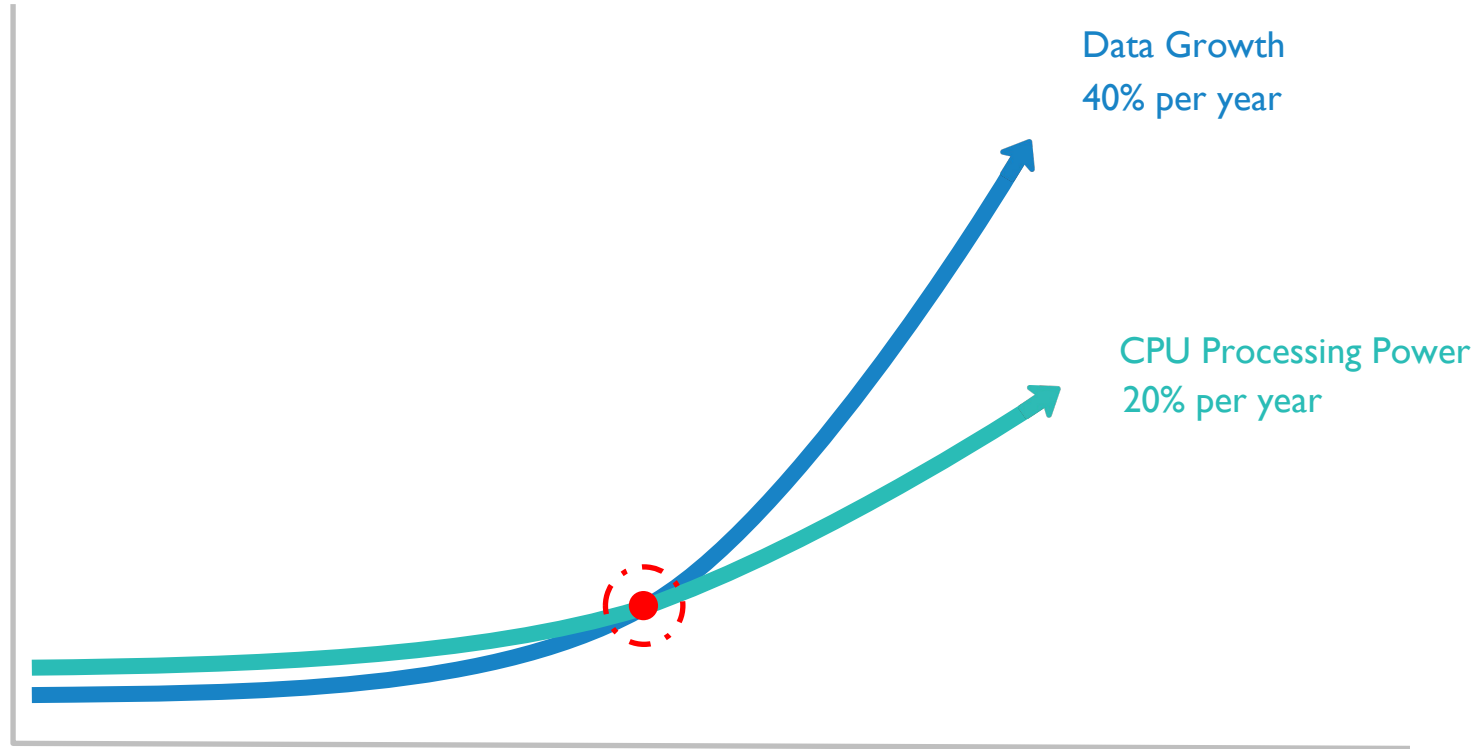
Volume

Agility

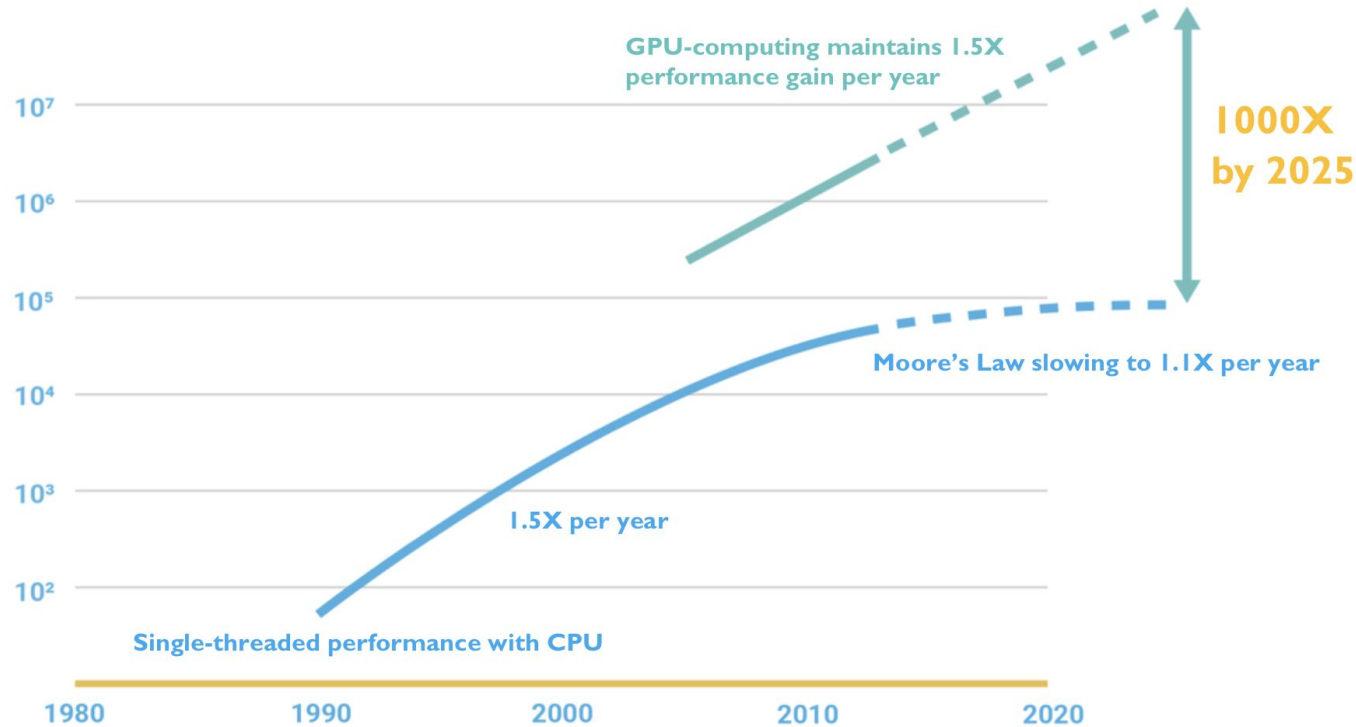
Spatio-

Temporal

Data Grows Faster Than CPU Processing

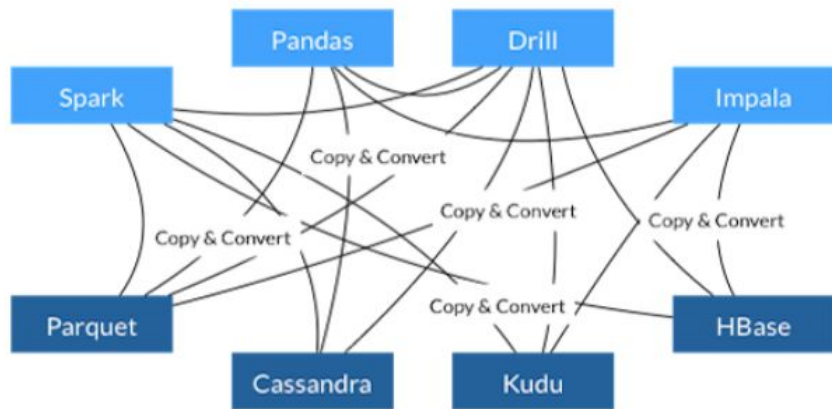


GPU Processing Keeps Moore's Law Alive

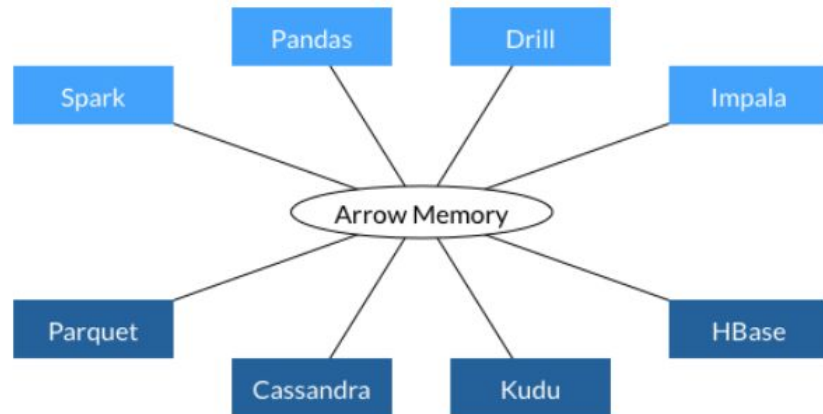


Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2015 by K. Rupp

Apache Arrow: Shared Memory Layout



- Each system has its own internal memory format
- 70-80% computation wasted on serialization and deserialization
- Similar functionality implemented in multiple projects



- All systems utilize the same memory format
- No overhead for cross-system communication
- Projects can share functionality (eg, Parquet-to-Arrow reader)

Source: <https://arrow.apache.org/>

RAPIDS and the GPU DataFrame

Born from the GPU Open Analytics Initiative – fusing Machine Learning and GPU analytics

CONTRIBUTORS



ADOPTERS



By adopting a common memory layout in Apache Arrow, these tools can work seamlessly with one another with zero-copy memory transfer!

OPEN SOURCE



Demo: Real-Time Bikeshare Data

GitHub: https://github.com/omnisci/gbfs_kafka

Goal: Monitor Bike Availability in Real-Time

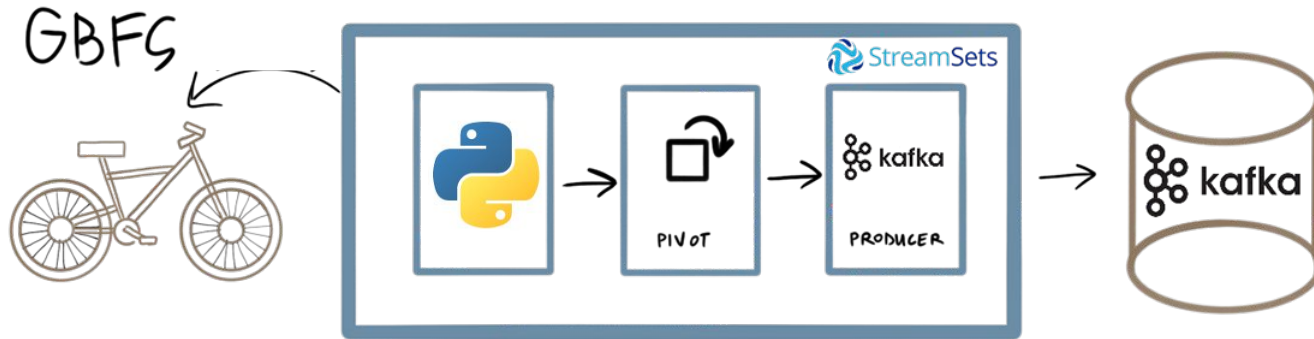
Architecture Considerations:

- Hundreds of API feeds conforming to GBFS specification:
<https://github.com/NABSA/gbfs>
- Each feed provides relatively small amount of info as JSON; need to pre-process before loading to OmniSci
- Feeds have different TTL values; want to be respectful when pinging API endpoints

Data Pre-Processing: Python and StreamSets

Using Azure HDInsight, we set up a managed Apache Kafka cluster

- Kafka serves several purposes: aggregating feeds into a single stream, buffers for a more consistent throughput
- [StreamSets](#) is a data pipeline tool



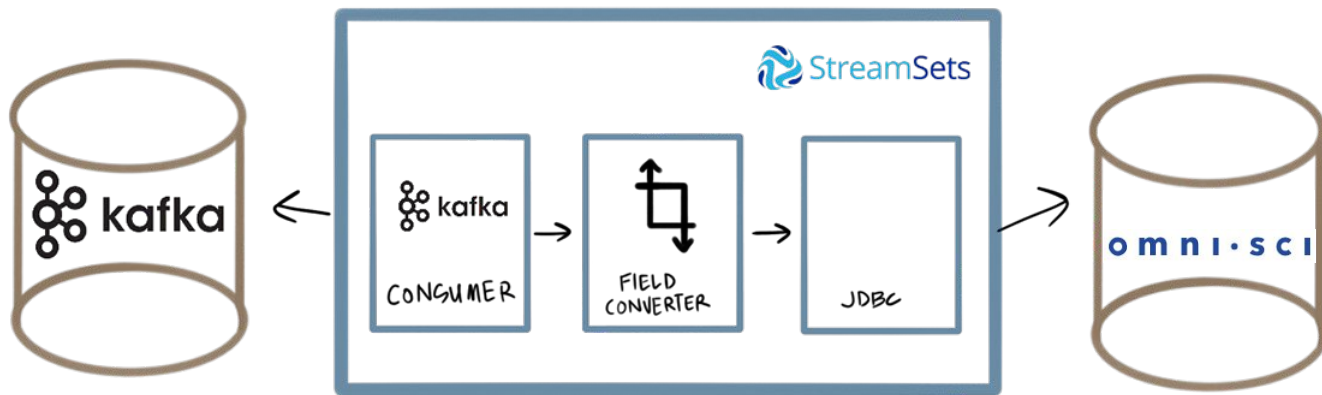
API Access Code Using Python

```
72 lines (54 sloc) | 2.57 KB
Raw Blame History
1 import pandas as pd
2 import requests
3 from multiprocessing import Pool
4 import datetime
5 import simplejson as json
6
7 from credentials import hosts
8
9 #round to minute, set the same for all records
10 #this is so that each load represents same time period
11 start_time = int(datetime.datetime.now().replace(second=0, microsecond=0).timestamp())
12
13 # read in endpoints
14 endpoints_df = pd.read_csv("data/gbfs_endpoints.csv")
15
16 # simple function to be p.map[pd]
17 def get_url(url):
18
19     r = None
20     try:
21         r = (requests.get(url, timeout=2).json(), url)
22     except:
23         pass
24
25     return r
26
27 def feedlist(df, feedtype):
28     return df[df["feedtype"] == feedtype].feedurl
29
30 # no reason for 12 specifically, just "some" parallelism
31 p = Pool(12)
32
33 ##### free_bike_status
34 bike_status = p.map(get_url, feedlist(endpoints_df, "free_bike_status"))
```

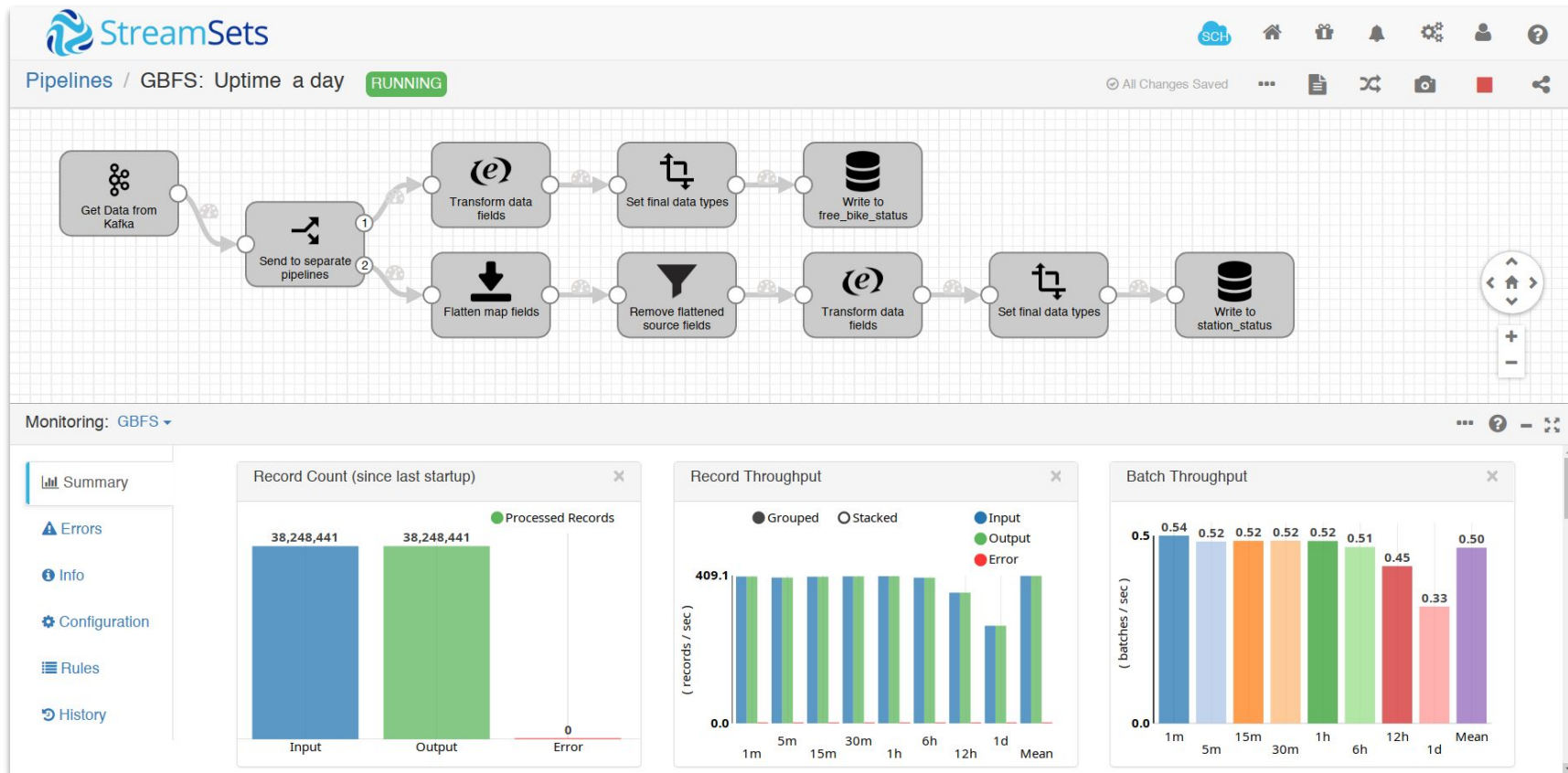
Data Ingestion: StreamSets to OmniSci

With feeds aggregated to single Kafka Producer (topic), ingest to OmniSci via JDBC

- OmniSci supports data streaming directly from Kafka, but using StreamSets allows for additional transformation
- Using JDBC along with StreamSets also allows StreamSets to manage retries and Kafka offsets



Data Ingestion: StreamSets to OmniSci



Demo: Python and OmniSci for Data Science

GitHub: <https://github.com/omnisci/odscwebinar>

Exploring the Data Using Python

ODSC Webinar: OmniSci and RAPIDS

An End-to-End GPU Data Science Workflow

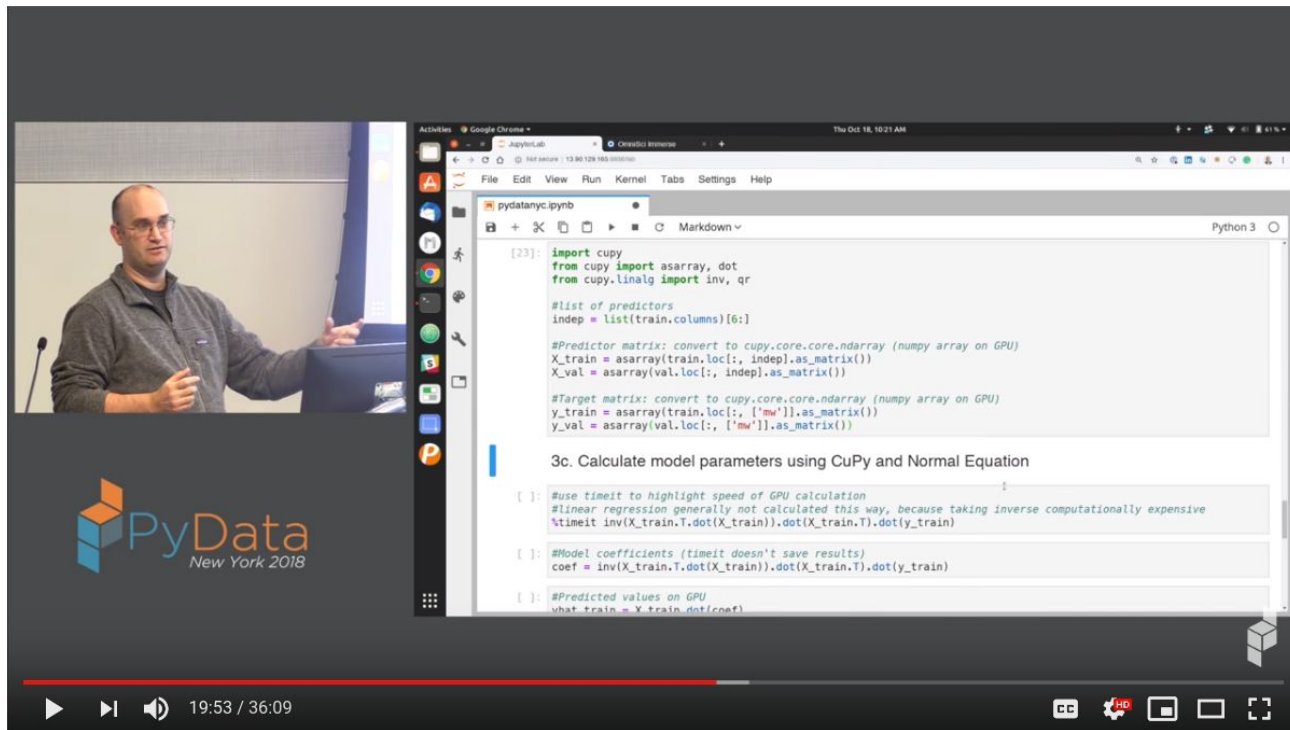
May 30, 2019

```
In [ ]: 1 #import pymapd to connect to OmniSci
        2 #importing pandas is for convenience with the pd.read_sql method
        3 #ibis is alternate method for querying using pandas-like API
        4 import pymapd
        5 import pandas as pd
        6 import ibis
        7 from credentials import credentials
```

```
In [ ]: 1 # Connect to OmniSciDB, get list of tables in database
        2 conn = pymapd.connect(host="localhost",
        3                       dbname=credentials["dbname"],
        4                       user=credentials["user"],
        5                       password=credentials["password"])
        6 conn.get_tables()
```

```
In [ ]: 1 query = """SELECT
        2 date trunc(minute, accessed_on) accessed_on,
        3 COUNT(*) AS records
        4 FROM free_bike_status
        5 WHERE accessed_on BETWEEN TIMESTAMP(0) '2019-05-30 00:33:21' AND TIMESTAMP(0) '2019-05-30 00:43:21'
        6 GROUP BY 1
        7 ORDER BY 1 DESC
        8 LIMIT 60"""
        9
       10 df = pd.read_sql(query, conn)
```


Machine Learning With OmniSci and cudf



The video player displays a presentation slide on the left and a Jupyter Notebook on the right. The slide features a man speaking and the PyData New York 2018 logo. The Jupyter Notebook shows the following code:

```
[23]: import cupy
      from cupy import asarray, dot
      from cupy.linalg import inv, qr

      #list of predictors
      indep = list(train.columns)[6:]

      #Predictor matrix: convert to cupy.core.ndarray (numpy array on GPU)
      X_train = asarray(train.loc[:, indep].as_matrix())
      X_val = asarray(val.loc[:, indep].as_matrix())

      #Target matrix: convert to cupy.core.ndarray (numpy array on GPU)
      y_train = asarray(train.loc[:, ['mv']].as_matrix())
      y_val = asarray(val.loc[:, ['mv']].as_matrix())
```

3c. Calculate model parameters using CuPy and Normal Equation

```
[ ]: #use timeit to highlight speed of GPU calculation
      #linear regression generally not calculated this way, because taking inverse computationally expensive
      %timeit inv(X_train.T.dot(X_train)).dot(X_train.T).dot(y_train)

[ ]: #Model coefficients (timeit doesn't save results)
      coef = inv(X_train.T.dot(X_train)).dot(X_train.T).dot(y_train)

[ ]: #Predicted values on GPU
      what_train = X_train.dot(coef)
```

<https://www.youtube.com/embed/gQszQcFHcZc>

Questions



Randy Zwitch

Senior Director of Developer Advocacy

 @randyzwitch

 randy.zwitch@omnisci.com

 /in/randyzwitch/

 /randyzwitch