Some things you learn running Apache Spark in production for three years William Benton (@willb) Red Hat, Inc.

Some things you learn running Apache Spark in production for three years William Benton (@willb) Red Hat, Inc.

About me

Since 2013: data engineering and data science focus

Since 2008: distributed systems focus

- Ancient history: compiler/VM design, static analysis, logic programming

Forecast

Introducing Apache Spark

How my team has used Spark

Lessons we've learned

From analytics as a workload to insightful applications

Meet Apache Spark

Parallel execution models



Parallel execution models





Parallel execution models





A user-friendly abstraction

Partitioned collections are spread across multiple processors or machines.

Immutable collections are copied, never modified in place.

Lazy operations only execute when necessary.

A user-friendly abstraction This means we'll have failures Partitioned collections are spread across multiple processors or machines.

Immutable collections are copied, never modified in place.

Lazy operations only execute when necessary.



A user-friendly abstraction This means we'll have failures Partitioned collections are spread across multiple processors or machines.

Immutable collections are copied, never modified in place.

Lazy operations only execute when necessary.

These mean we'll always have a recipe for how to recover



































A simple example program

- file = sc.textFile("file://...")
- counts = file.flatMap(lambda l: l.split(" ")) .map(lambda w: (w, 1)).reduceByKey(lambda x, y: x + y)
- # computation actually occurs here counts.saveAsTextFile("file://...")

Spark core (collections and scheduler)

Graph SQL

Spark core (collections and scheduler)



Graph SQL

Spark core (collections and scheduler)

ad hoc



Mesos

YARN























 $v("Madrid") - v("Spain") + v("France") \approx v("Paris")$

A more interesting example

read text file as a data frame df = spark.read.text("file://...")

fit and use a Word2Vec model

model = w2v(df)

model.findSynonyms("data", 5).show()

- .select(split("value", "\s+").alias("text"))
- w2v = Word2Vec(inputCol="text", outputCol="result")

How we've used Spark



| ••• | |
|-----|--|
| ••• | |
| ••• | |
| ••• | |
| ••• | |
| ••• | |

| ••• | |
|-----|-----|
| ••• | |
| ••• | |
| ••• | |
| ••• | |
| ••• | III |



| ••• | |
|-----|--|
| ••• | |
| ••• | |
| ••• | |
| ••• | |
| ••• | |



| ••• | |
|-----|--|
| ••• | |
| ••• | |
| ••• | |
| ••• | |
| ••• | |





Prototyping new techniques


Machine configuration analysis



Machine configuration analysis







000

Ш





| ••• | |
|-----|--|
| ••• | |
| ••• | |



| ••• | |
|-----|-----|
| ••• | III |
| ••• | III |
| ••• | III |

| ••• | III |
|-----|-----|
| ••• | III |
| ••• | III |

| ••• | III |
|-----|-----|
| ••• | III |
| ••• | III |



| ••• | III |
|-----|-----|
| ••• | III |
| ••• | |
| ••• | III |



Machine configuration analysis

| | - | |
|--|------|--|
| | | |
| | •••• | |
| | | |
| | r | |
| | _ | |
| | | |
| | ľ | |
| | | |
| | | |
| | r | |
| | | |
| | | |

• • •

- not libhugetlbfs && not fontpackages-filesystem &&
 - httpclient && transitional-eap6-jars
- b11-kit && not libtheora && not perl-Text-ParseWords &&
 jboss-metadata-appclient
- not python-suds && rt61pci-firmware &&
 - avahi && apache-commons-io-ea
- not sanlock-python && shrinkwrap-parent &&
 - not cli-tools-zend-server
- not perl-parent && not python-stevedore && jdom &&
 - not openshift-origin-cartridge-<mark>diy &&</mark>
 - not redh<mark>at-sso-log</mark>in-module-eap<mark>6 && iwl60</mark>50-firmware-41
- ot pytall<mark>oc && not</mark> libldb && xor<mark>g-x11-font</mark>s-Typ &&
 - not nodejs010-gyp && hibernate4-entitymanager









Do people work on Fedora for love or money?



Do people work on Fedora for love or money?

Is there anyone the community couldn't live without?



Do people work on Fedora for love or money?

Is there anyone the community couldn't live without?

How do we characterize breadth and depth of community engagement?



























en en fan de De fan

i fa ai

Out of 310 million log records, we identified 0.0012% as outliers.

en en fan de De fan

i fa ai



Thirty most extreme outliers

- Can not communicate with power supply 2. 10
- Power supply 2 failed. 9
- Power supply redundancy is lost. 8
- Drive A is removed. 1
- Can not communicate with power supply 1. 1
- Power supply 1 failed.

Modeling infrastructure costs



SYSTEM METRICS



CLOUD SPENDING

SELECT COUNT(value), MIN(value), MAX(value), AVG(value), items.key_, items.hostid **FROM** history, items WHERE history.itemid = items.itemid **GROUP BY** history.itemid

120gb of data on one node with 40 threads and 384gb of RAM



SELECT COUNT(value), MIN(value), MAX(value), AVG(value), items.key_, items.hostid **FROM** history, items WHERE history.itemid = items.itemid **GROUP BY** history.itemid

RDBMS: 15 hours

120gb of data on one node with 40 threads and 384gb of RAM





SELECT COUNT(value), MIN(value), MAX(value), AVG(value), items.key_, items.hostid **FROM** history, items WHERE history.itemid = items.itemid **GROUP BY** history.itemid

RDBMS: 15 hours Spark: 15 minutes

120gb of data on one node with 40 threads and 384gb of RAM





Lessons we learned

METALESSON: HOW TO MASTER DECLARATIVE PROGRAMMING

Three steps to mastery

Understand the programming model.

Understand the execution model.

Understand when to let the environment work for you.

Three steps to mastery "What does this mean?"

Understand the execution model.

Understand when to let the environment work for you.

Three steps to mastery "What does this mean?" "What does this do?"

Understand when to let the environment work for you.

Three steps to mastery "What does this mean?" "What does this do?" "How can I get out of its way?"

Avoid shuffles when possible

Cache only when necessary



Avoid shuffles when possible

Cache only when necessary





LESSON: LEARN THE API; USE THE RIGHT METHODS

Spark driver (application)



Spark driver (application)

Spark workers





Spark workers





Spark workers


Take advantage of the model





Take advantage of the model



Spark driver (application)





Take advantage of the model



Spark driver (application)



Spark driver (application)





Spark driver (application)





Spark driver (application)







Spark driver (application)









Spark driver (application)



Spark driver (application)

Aggregate at the workers instead of in the driver!



LESSON: LET SPARK WORK FOR YOU WHENEVER POSSIBLE

Two favorite features

Query planning makes dumb code run faster. Typed APIs prevent really dumb code from running at all.

Query planning

- - A.ID = B.ID AND
 - uncommon(A.X) AND
 - extremelyRare(B.Y)

SELECT * FROM A, B WHERE













FILTER



FILTER









FILTER

FILTER





FILTER FILTER





LESSON: STORAGE FORMATS MATTER MORE THAN LOCALITY















"Disks are too slow"

"Memories are too small"

"Locality is king"

"For the workloads from Facebook and Bing, we see that 96% and 89% of the active jobs respectively can have their data entirely fit in memory, given an allowance of 32GB memory per server for caching"

- "PACMan: Coordinated Memory Caching for Parallel Jobs." G. Ananthanarayanan et al., in Proceedings of NSDI '12.





from local disks is only about 8% faster than [and] this 8% number is decreasing."

-Tom Phelan, "The Elephant in the Big Data Room: Data Locality is Irrelevant for Hadoop" (goo.gl/MnCKuM)

"Recent studies have shown that reading data

- reading it from remote disks over the network ...



"Three out of ten hours of job runtime were spent moving files from the staging directory to the final directory in HDFS... We were essentially compressing, serializing, and replicating three copies for a single read."

- "Apache Spark @Scale: a 60+ TB production use case" Facebook Engineering Blog Post



Do you need node locality?

perform worse than iterative processing in memory.

compute and storage.

- Working set sizes typically fit in cluster memory even if raw data don't.
- **I/O-heavy frameworks** designed for colocated compute and storage
- Colocating compute and storage prevents independent scale-out of



An optimization that matters

An optimization that matters

FORMAT ORIENTED ROW

An optimization that matters

FORMAT ORIENT ROW


FORMAT ORIENT ROW



FORMAT μ ORIENT ROW





FORMAT μ ORIENT ROW





FORMAT μ ORIENT ROW-





FORMAT ORIENT ROW



FORMAT UMNAR CO

10% of the space 1-10% of the time

LESSON: THINGS TO CONSIDER WHEN PREDICTING THE FUTURE

Feature engineering



| nountain bike | 0 | 1 | 0.35 | 1 | 1 | 1 |
|---------------|---|---|------|---|---|---|
| | | | | | | |

| LABEL | DROP | FLAT | TIRE | TIRE | FRONT | REA |
|-------|----------------|------|------|-------|-------------|-----|
| | HANDLEBAR TYPE | | SIZE | KNOBS | SUSPENSION? | |

| clocross bike | 1 | 0 | 0.13 | 1 | 0 | 0 |
|---------------|---|---|------|---|---|---|
| | | | | | | |









> bird, animal, outdoors







food, macro, fruit







outdoors, racing, bike, wout van aert



outdoors, racing, bike, wout van aert





outdoors, racing, bike, wout van aert

Accurate predictions are only part of a model's value!



From analytics as a workload to insightful applications



Mesos

Spark executor

Spark executor

Spark executor

Spark executor

Spark executor

Spark executor

Mesos

Spark executor

Spark executor

Spark executor

Spark executor

Spark executor

Spark executor

Mesos

1 Spark executor

1 Spark executor

2 Spark executor

3 Spark executor

3 Spark executor

4 Spark executor

2 3 4

Mesos

1 Spark executor

1 Spark executor

2 Spark executor

3 Spark executor

3 Spark executor

4 Spark executor













































Kubernetes













app 6



Object stores

Databases

Kubernetes











app 6



Object stores

Databases

Kubernetes









http://radanalytics.io



app 6



Databases

Object stores

Takeaways

Let Spark work for you.

Use efficient storage formats but don't worry about data locality (yet).

ETL input data to Parquet early in your process.

Feature engineering effort often trumps fancy models.

Prefer interpretable models and easy-to-implement algorithms.



willb@redhat.com • @willb https://chapeau.freevariable.com http://radanalytics.io