Extending the Enterprise Data Warehouse with
Hadoop
Robert Lancaster

Nov 7, 2012

*ORBITZ*®

# Who I Am

- Robert Lancaster
  - Solutions Architect, Hotel Supply Team
  - rlancaster@orbitz.com
  - @rob1lancaster
  - Organizer of Chicago Machine Learning Study Group
  - Co-organizer of Chicago Big Data.

# Launched in 2001

## Over 160 million bookings

# Some History…

# In 2009…

- The Machine Learning team is formed to improve site performance. For example, improving hotel search results.

- This required access to large volumes of behavioral data for analysis.

  - Fortunately, the required data was collected in session data stored in web analytics logs.

- The only archive of the required data went back about two weeks.

Non-transactional Data (e.g. searches)

Transactional data (e.g. bookings) and aggregated Non-transactional data

Data Warehouse

**Detailed non-transactional data (what every user sees, clicks, etc.)**

**Transactional data (e.g. bookings) and aggregated Non-transactional data**

Data Warehouse

Hadoop

# What is Hadoop?

- Distributed file system and parallel processing platform.

- Open source Apache project created by Doug Cutting.

- Modeled on papers published by Google on the Google File System and MapReduce.

- Intended to run on a cluster of relatively inexpensive machines (aka commodity hardware).

- Bring processing to the data.

# But Brought New Challenges…

- Most of these efforts are driven by development teams.

- The challenge now is unlocking the value of this data for non-technical users.

- Support for Hadoop via traditional BI/reporting tools still meager.

Both big (relatively)…

INFORMATICA
The Data Integration Company™

MicroStrategy®

QlikView

GREENPLUM®

ORBITZ®

# And small…

- Big Data team is formed under Business Intelligence team at Orbitz Worldwide.

- Allows the Big Data team to work more closely with the data warehouse and BI teams.

- Reflects the importance of big data to the future of the company.

- Our production cluster has grown 40-fold since it was launched.

# A View Shared Beyond Orbitz…

"We strongly believe that Hadoop is the nucleus of the next-generation cloud EDW…"

"…but that promise is still three to five years from fruition."*

*James Kobielus, Forrester Research,
"Hadoop, Is It Soup Yet?"

- Extraction and transformation of data for loading into the data warehouse – "ETL".

- Off-loading of analysis from the data warehouse.

## Proposed Processing

Raw logs → Hadoop → Dimensional model

Previous Processing in Data Warehouse



Several hours of processing

~20% original data size

- Moving to Hadoop:

  - Removed load from the data warehouse.

  - Facilitated adding additional attributes for processing.

  - Allowed processing to be run more frequently.



## Processing in Hadoop

- Facilitated analysis that allows for more personalized ad content.

- Allowed marketing team to analyze over a years worth of search data.

- Provided analysis that was difficult to perform in the data warehouse.

# Example Use Case: Selection Errors

- Multiple points of entry.

- Multiple paths through site.

- Goal: tie events together to form picture of customer behavior.

# Use Case – Selection Errors: Visualization

# Example Use Case: Beta Data

- Hotel Sort Optimization

- Compare A vs. B

- Web Analytics Data

  - What user saw.

  - How user behaved

- Server Log Data

  - Sorting behavior used.

# Example Use Case: RCDC

- Understand and improve cache behavior.

- Improve "coverage"
  - Traditionally search 1 page of hotels at a time.
  - Get "just enough" information to present to consumers.
  - Increase amount of availability information we have when consumer performs a search.

- Data needed to support needs beyond reporting.

# Use Case – RCDC: Visualization

# Conclusions

- Hadoop market is still immature, but growing quickly. Better tools are on the way.

  - Look beyond the usual (enterprise) suspects. Many of the most interesting companies in the big data space are small startups.

- Hadoop won't replace your EDW, but any organization with a large EDW should at least be exploring Hadoop as a complement to their BI infrastructure.

# Conclusions

- Work closely with your existing data management teams.

  - Your idea of what constitutes "big data" might quickly diverge from theirs.

- The flip-side to this is that Hadoop can be an excellent tool to off-load resource-consuming jobs from your data warehouse.

Thank you!

Questions?