

Deep Analytics Pipeline

A Benchmark Proposal

Milind Bhandarkar
Chief Scientist, Pivotal Software
(Twitter: @techmilind)

Quest for “Typical” Workloads

- Benchmarks most relevant if representative
- Tune systems for broadly applicable workloads
- Designing optimized systems: Make common tasks fast, other tasks possible

Encouraging Early Results

- Analyzed characteristics of 1M+ real Hadoop jobs on production clusters at Yahoo, 100+ features
- Identified 8 Job types
- Verified with GridMix 3
- Characterization of Hadoop Jobs Using Unsupervised Learning, Sonali Aggarwal, Shashank Phadke & Milind Bhandarkar, in 2010 IEEE Second International Conference on Cloud Computing Technology and Science, Indianapolis, Indiana, December 2010 (<http://doi.ieeecomputersociety.org/10.1109/CloudCom.2010.20>)

Big Data Benchmark Community

- Formed in January 2012 under the Center for Large Data Systems Umbrella (clds.sdsc.edu)
 - 200+ members on the mailing list
- Workshops on Big Data Benchmarking (San Jose, 2012, Pune, India 2012, Xian, China 2013, San Jose CA 2013, Germany 2014)

Benchmark Proposals

- BigBench : Rabl et al, Published in SIGMOD 2013, Full Specification published in WBDB proceedings (Springer-Verlag)
 - Data Model enhanced from TPC-DS
 - Added Text Mining, Relevance
- TPC-BD : Formed September 2013, Terasort + TPC-rigor

BigData Top 100

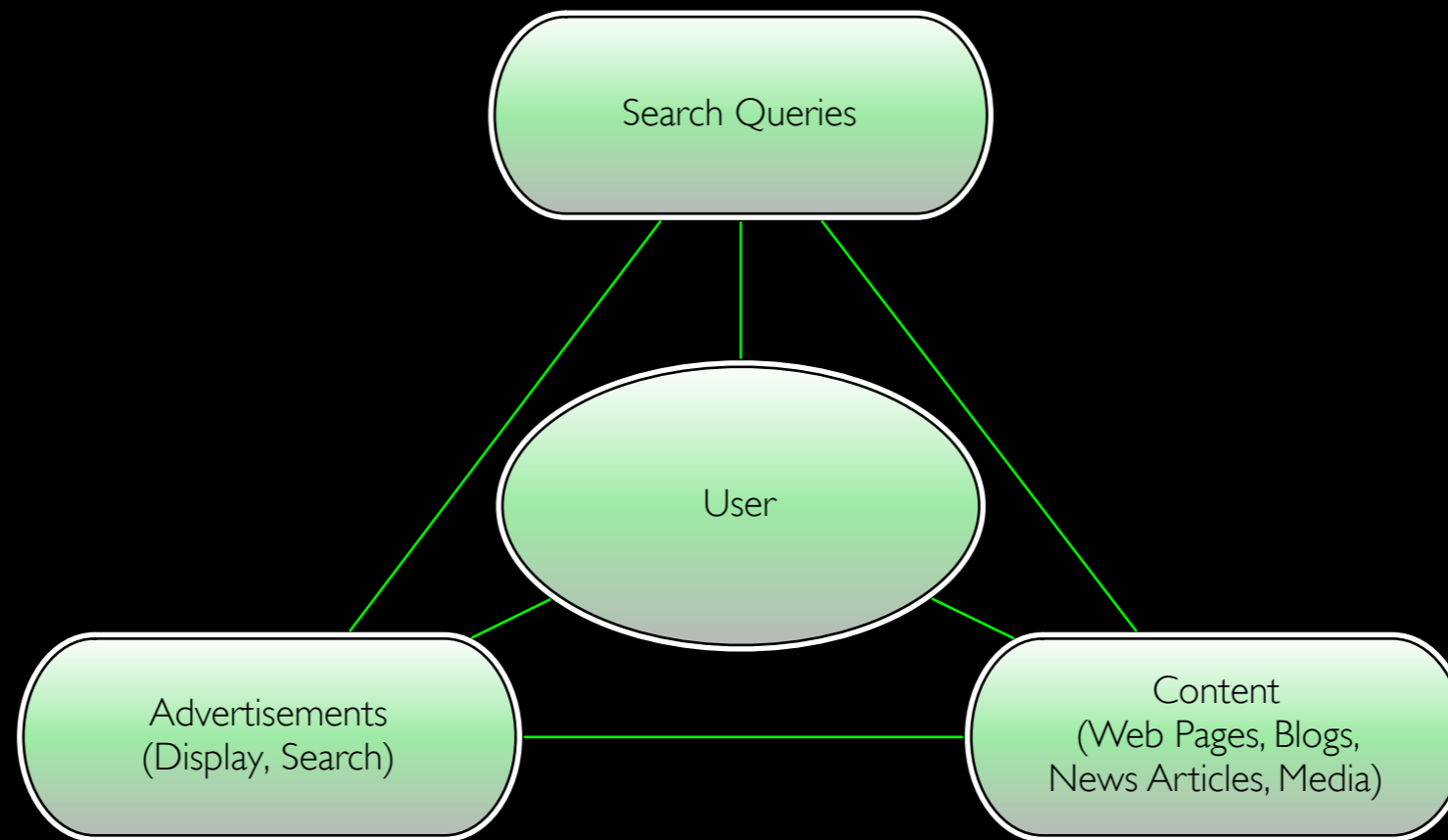
- Modeled after Top500 & Graph500 in HPC Community
- Proposal Presented at Strata Conference in February 2013
- Workload : BigBench (most likely)

Drivers for Big Data

- Ubiquitous Connectivity
- Sensors Everywhere
- Democratization of Content

Big Data Sources

- Events
 - Direct - Human Initiated
 - Indirect - Machine Initiated
- Software Sensors (Clickstreams, Locations)
- Public Content (blogs, tweets, Status updates, images, videos)



Online: Major Data Sources

“User” Modeling

- Objective: Determine User-Interests by mining user-activities
- Large dimensionality of possible user activities
- Typical user has sparse activity vector
- Event attributes change over time

User-Modeling Pipeline

- Data Acquisition, Normalization, Sessionization
- Feature and Target Generation
- Model Training
- Offline Scoring & Evaluation
- Batch Scoring & Upload to serving

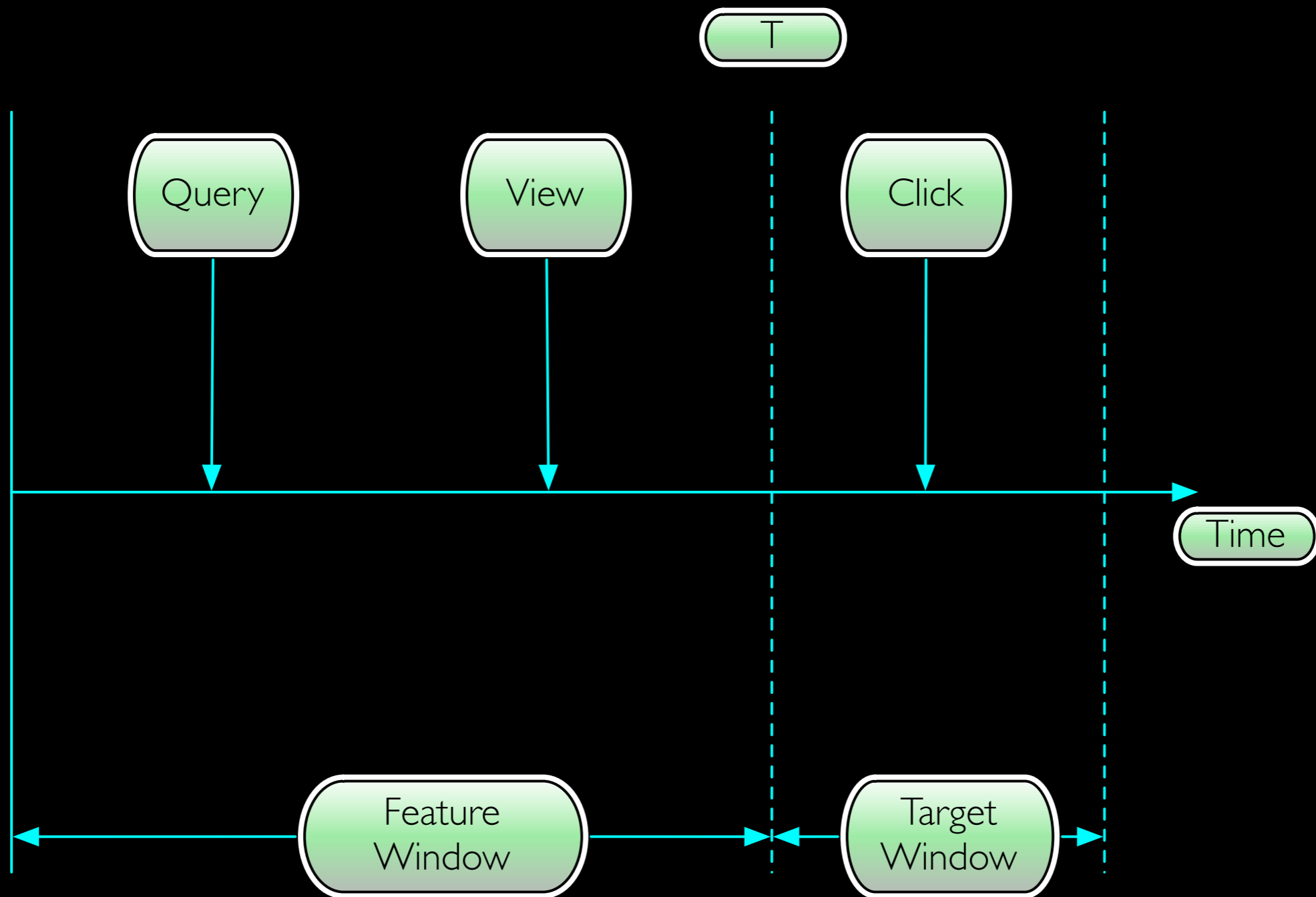
Data Acquisition

User	Time	Event	Source
U0	T0	Visited Auto website	Web Server logs
U0	T1	Searched for "Car Insurance"	Search Logs
U0	T2	Browsed stock quotes	Web Server Logs
U0	T3	Saw ad for "discount brokerage", did not click	Ad Logs
U0	T4	Checked Mail	Web Server Logs
U0	T5	Clicked Ad for "Auto Insurance"	Ad Logs, Click Logs

Normalization

User	Time	Event	Tag
U0	T0	View	Category: Autos, Tag: Mercedes Benz
U0	T1	Query	Category: Insurance, Tag: Auto
U0	T2	View	Category: Finance, Tag: EMC
U0	T3	View-Click	Category: Finance, Tag: Brokerage
U0	T4	Browse	Irrelevant Event, Dropped
U0	T5	View+Click	Category: Insurance, Tag: Auto

Features & Targets



Targets

- User-Actions of Interest
 - Clicks on Ads & Content
 - Site & Page visits
 - Conversion Events
 - Purchases, Quote requests
 - Sign-Up for membership etc

Features

- Summary of user activities over a time-window
- Aggregates, moving averages, rates over various time-windows
- Incrementally updated

Joining Targets & Features

- Target rates very low: 0.01% ~ 1%
- First, construct targets
- Filter user activity without targets
- Join feature vector with targets

Model Training

- Regressions
- Boosted Decision Trees
- Naive Bayes
- Support Vector Machines
- Maximum Entropy modeling
- Constrained Random Fields

Offline Scoring & Evaluation

- Apply model weights to features
- Pleasantly parallel
- Sort by scores and compute metrics
- Evaluate metrics

Batch Scoring

- Apply models to features from all user activity
- Upload scores to serving systems

Issues

- Different modeling techniques for different kinds of data
- Different notions of a “session”
- Widely varying number of events per entity

Proposal: 5 Classes

- Tiny (100K entities, 10 events per entity)
- Small (1M entities, 10 events per entity)
- Medium (10M entities, 100 events per entity)
- Large (100M entities, 1000 events per entity)
- Huge (1B entities, 1000 events per entity)

Proposal: Publish results for every stage

- Data pipelines constructed by mix-and-match of various stages
- Different modeling techniques per class
- Need to publish performance numbers for every stage

Questions ?