

## Scal(a)ing up Machine Learning and Big Graph Analytics



Tyson Condie, Markus Weimer and Raghu Ramakrishnan Yahoo! Research

> Vinayak Borkar, Yingyi Bu and Mike Carey U.C. Irvine

Joshua Rosen and Neoklis Polyzotis





### Develop a platform for Machine Learning and Graph analytics on a cluster of shared-nothing machines

- Programming environment
  - > Specific to ML and Graph domain
- Data-parallel runtime
  - › Supports ML and Graph analytic workflows





# **Characteristic Workflow for ML Analytics**

- Feature Extraction
  - > ETL workflow
- Modeling
  - > Iterative algorithm that fits a model to the data
- Evaluation
  - › Checks model fidelity



## What is a Model?

- Global Model
  - › Aggregate summary of all data points
  - > Small in size compared to data set
  - > e.g., regression, classification

- Local Model
  - > Interdependent parameters for each data point
  - > Size proportional to the data set
  - › e.g., topic (graphical) models, clustering, PageRank



## Outline

- Case Study
- Modeling in the Cloud
- Our work
- Project roadmap





#### YAHOO! MAIL

### **Outbound Spam Filter for Yahoo! Mail**



# **Ingredient: Client Classification**

#### Goal:

> Train a spam classifier that can classify **outbound** messages

#### Problem:

› Other mail services won't tell us if our users send spam

#### Solution:

 Take the mail from Yahoo! to Yahoo! and observe how it was classified by our users



### **EMail Flow**





### **EMail Flow: Data Sets**





**Step I Extract Features** 

• NET\_DB
, [IP, ...]

[OUT\_ID, USER\_FEATURE NET\_FEATURES]

11/1/11



NIOC





**Step II: Finding a Label** 

- OUTBOUND\_LOG
  - > [OUT\_ID, FROM, TO, TIME]

- WEB\_LOG
  > [IN\_ID,VOTE]

This needs a fuzzy join



## **Step III: Creating the training data**

- JOIN labels and features
- Subsample the training data
- Copy to a single machine



## **Step IV: Training the model**

- Until a satisfactory model is found:
  - > Train the model using sequential code
  - > Evaluate the found model in Pig again
  - Apply insights to feature extraction, label generation and learning algorithm



### Recap



## Take-Away

### Usability is bad (took an intern 3 months)

- Many tools and technologies needed
- › Fractured workflow, captured in e.g. Oozie

### Subpar solution

- Subsampling hurts classifier fidelity
- › Copying data to a single machine is slow



### **Now: Parallelize the middle**

Extraction

Feature



## What do people do?

#### Gradient Boosted Decision Trees

- > Popular non-linear modeling algorithm
- > Used in Yahoo! search engine to learn the ranking function

#### Gradient Descent

- > A method for minimizing objective functions
- > Written as a sum of differentiable functions

#### Latent Dirichlet Allocation (LDA)

- > Example of a topic model
- > Usually expressed via Graphical Model



# What do people do?

#### Gradient Boosted Decision Trees

- > Popular non-linear modeling algorithm
- > Used in Yahoo! search engine to learn the ranking function
- Solution: Fake mappers that implement MPI
- Gradient Descent
  - > A method for minimizing objective functions
  - > Written as a sum of differentiable functions
  - > Solution: Fake mappers that run aggregation trees
- Latent Dirichlet Allocation (LDA)
  - > Example of a topic model
  - > Usually expressed via Graphical Model
  - > Solution: Fake mappers that implement this graphical model



# Framework I (Global Models): Spark

- Data Model: Resilient Distributed Datasets
  - > Support for hash and range partitioning
  - > Transformations or Actions can be added to them
  - The results can be materialized in memory (RDD.cache()) or to disk (RDD.save(...))
- Caching allows fast iteration



## Framework II (Local Models): Pregel

21

- Graph-oriented API
  - > UDF that implements single vertex update
  - Input messages from neighboring vertices
  - > Output: outgoing messages and new vertex state
- Bulk synchronous parallel computing
  - Runtime executes all vertices with messages
  - Model updated and cached in memory
  - Checkpoint for fault-tolerance
  - > Continues until everyone votes to halt



## Framework: Pros/Cons

### Pros

- > Adds the notion of a working set for fast iterations
- > Exports a clean API specific to the problem domain
- Abstracts (many) low-level implementation details

### Cons

- › Point solutions that mainly focus on modeling
- Limited optimization support
- › No external memory operators



# ScalOps + HyracksML

- ScalOps
  - Scala DSL for the entire analytic pipeline
  - Supports Pig Latin
  - Looping construct that captures iteration
- HyracksML + Algebricks
  - Data-parallel runtime that runs on a cluster of shared-nothing machines
  - > Deductive database extensions to directly support iteration
  - Relational algebra and query optimizer that captures the entire workflow

Programming Model ScalOps

Compiler/Optimizer Algebricks

Runtime Engine *HyracksML* 

Scalable File System e.g. HyracksES\_HDES



### **Spam Filter: ScalOps**

object OutboundSpamFilter extends Query {

```
// Join Targets
  case class Features(outID: String, features: Vector)
  case class Label(outID: String, label: Double)
  def getLabels(inboundLogs: Queryable[YMail.InboundRecord], outboundLogs: Queryable[YMail.OutboundRecord], webLogs: Queryable[YMail.WebRecord]) =
    val inoutmap = join(inboundLogs by (YMail.keyFor()), outboundLogs by (YMail.keyFor()))
    join(inoutmap by ( . 1.inID), webLogs by ( .inID)).map(e => Label(e. 1. 2.outID, e. 2.vote))
  }
  def getFeatures (outboundLogs: Oueryable[YMail.OutboundRecord], userDB: Oueryable[YMail.User], netDB: Oueryable[YMail.IP]) = {
    case class IPInformation (outID: String, ipInfo: YMail.IP)
    case class UserInformation(outID: String, userInfo: YMail.User)
    val users = join(outboundLogs by ( .from), userDB by ( .id)).map(e => UserInformation(e. 1.outID, e. 2))
    val ips = join(outboundLogs by ( .ip), netDB by ( .ip)).map(e => IPInformation(e. 1.outID, e. 2))
    join(outboundLogs by ( .outID), users by ( .outID), ips by ( .outID)).map(e => Features(e. 1.outID, YMail.extractFeatures(e. 1, e. 2.userInfo,
e. 3.ipInfo)))
  def run {
    val inboundLogs = load[YMail.InboundRecord]("hdfs://ymail/inbound.dat").filter(YMail.isYMail)
   val outboundLogs = load[YMail.OutboundRecord]("hdfs://ymail/outbound.dat").filter(YMail.isYMail)
    val webLogs = load[YMail.WebRecord]("hdfs://ymail/web.dat")
    val userDB = load[YMail.User]("hdfs://y/users.dat")
    val netDB = load[YMail.IP]("hdfs://y/network.dat")
    // Extract labels
    val labels = getLabels(inboundLogs, outboundLogs, webLogs)
    // Extract Features
    val features = getFeatures(outboundLogs, userDB, netDB)
    // Assemble a dataset
    val dataset = join(labels by ( .outID), features by ( .outID)).map(e => new Example(e. 1.outID, e. 1.label, e. 2.features))
    // Train a model
    val model = ML.train(dataset)
    // Evaluate
    val score = dataset.map(ML.evaluate(model, )).reduce( + ) / dataset.count
```



### **Spam Filter: Hyracks RQ**













# **Conclusion: Target Algorithms**

### Linear Models

- › Batch Gradient Descent
- › Stochastic Gradient Descent
- › Logistic Regression

### Clustering

- › K-Means
- Matrix Factorization

### Graph Analysis

- › PageRank
- › LDA (Graphical Models)



# **Coming Attractions**

- Recursive data-parallel runtime (HyracksML)
  - Deductive Databases, BOOM, Flying Fixed-Point



Yingyi Bu and Joshua Rosen

- Extended Relational Algebra (Algebricks)
  - **Captures the entire workflow in an optimizer friendly representation**



Vinayak Borkar

- High-level language (ScalOps)
  - > Abstractions for Machine learning and Graph-based algorithms



**Markus Weimer** 

Executive Producers



Mike Carey, Tyson Condie, Neoklis Polyzotis and Raghu Ramakrishnan



