Scalable Machine Learning to Exploit Big Data for Knowledge Discovery

Una-May O'Reilly MIT

MIT ILP-EPOCH Taiwan Symposium Big Data: Technologies and Applications





Lots of Data Everywhere





Knowledge Mining Opportunities



The GigaBeats Project







GigaBeats Project



Machine Learning Primer













Agenda **Distributed Computation** +-Scalable Machine Learning SCALE FlexGP **EC-Star**





SCALE









DCAP Protocol







SCALE Demonstration: WD-1

Forecasting ABP

- 10 levels
- 95K exemplars
- 67K/28K training/test split
- 7 dimensions
 - Stats, trends on MAP
- 1225 learner tasks (10 fold cross validation)
- 80 nodes, ~2 days, ~4000 node hours

Classifier	Number of Instances	Errors
Neural Networks	720	0
Discrete Bayes	18	0
Naive Bayes	64	0
Support Vector Machines (SVM)	66	34
Decision Trees	324	0





SCALE: WD-1 Running time



SCALE: WD-1 Accuracy results

F1 Score over all Classifiers, All Classes



SCALE: WD-1 Comparison SVM v DT







Scaling up: from SCALE to FlexGP

SCALE

- Modest 10's of features
- Assumes all training data fits into RAM

FlexGP

- 100's of features
- Big Data
- Big Data requires multidimensional factoring, filtering then fusion













FlexGP Learner





Genetic Programming Symbolic Regression



FlexGP Filter and Fusion



Factoring is Better







FlexGP: Data Factoring Size Study



Resource and System Management Layer

Completely decentralized launch

- Start up
 - Cascading launch is decentralized and allows elastic retraction and expansion

Completely decentralized network support

- IP discovery and launch of gossip protocol added to start up cascade
 - Launch Ip discovery and ongoing gossip protocol enables communication network among nodes at algorithm level
 - Support for Monitoring/reporting/harvesting current best

Remarks

- Essential design for resilience to node failure
 - Launch or running node loss will not halt computation,
 - lost launch branch or node can be integrated seamlessly
 - Node loss in a communication network won't break the network

Statistical oversight of learning algorithm's execution parms and data

 distribution for parms and data added to start up cascade, negotiated between parent and child at launch of child



FlexGP



System Layer: SCALE vs FlexGP

SCALE:

- LAMPs pre-defines the tasks
- Every learner has to know IP of task handler
- Task handler is a bottleneck and central point of failure

FlexGP

- Autonomous task specification
- Learners gossip to learn each others' IP
- No central task handler or point of failure





FlexGP System Layer









ECStar

- Goal: compute very cost effectively on *VAST* number of nodes...with a lot of training data
 - Runs on thousand to 10'Ks 100K's million nodes
 - Vast requires cost effective -> volunteer
- Domain: learn from time series
 - Finance, medical signals domain
- Solution is strategy or classifier expressed as rule sets





EC-Star Paradigm Digital Directed Evolution of Models



EC-Star - dedicated replicable Graduates Migrants Evolutionary Coordinator Condition and President and Volunteer Evolutionary Evolutionary Engine Engine compute Population Population Breeding Breeding Pool Pool Training Case Server dedicated _ replicable randomly sampled

EC-Star Divide and Conquer

genet

C U.M. O'Reilly

BigData factoring strategy

- Use fewer training samples and cull poor models
- Increase training samples on better models
- Local and distributed randomization
- Oversampling on superior models



System Layer Comparison

	Scale	FlexGP	EC-Star
ML domain	Classification	Regression Classification	Rule Learning
Resource Scale	10's to 100	100's to 1000	10^3 to 10^6
Resource Type	Cloud	Cloud	Volunteer and Dedicated
Fusion	External	External	Integrated
Local Algorithm	Different	Same	Same
Server:Client ratio	1: many	Decentralized	Few: many





Scalable Machine Learning Comparison

	SCALE	FlexGP	EC-Star
Factor	Algorithm	Algorithm and Data	Data: Under to oversampling
Filter		Correlation Accuracy	Layered competition
Fuse		Non- parametric output space approaches	Migration and ancestral properties





Automation

• "In the end, the biggest bottleneck is not data or CPU cycles, but human cycles."





Looking Forward



ML requires a lot of Human Effort





Looking Forward



Compressing the ML Endeavor

- From desktop
 - Multi-core
- To GPU
- To Cloud
 Seamlessly!
 Rapidly!
 Flexibly!
 Scalably!

Wrap Up

DCAP is open source: https://github.com/byterial/dcap FlexGP is documented in publications and thesis

Owen Derby, MEng, 2013

Thanks to...

- ALFA group members
 - Large team of students
 - Postdoc: Dr. Erik Hemberg
 - Research Scientist: Dr. Kalyan Veeramachaneni
- CSAIL Quanta cloud
- Our collaborators and sponsors

