

Strategies & Principles for Distributed Machine Learning

Eric Xing

epxing@cs.cmu.edu

School of Computer Science Carnegie Mellon University

Acknowledgement: Wei Dai, Qirong Ho, Jin Kyu Kim, Abhimanu Kumar, Seunghak Lee, Jinliang Wei, Pengtao Xie, Yaoliang Yu, Hao Zhang, Xun Zheng James Cipar, Henggang Cui, and, Phil Gibbons, Greg Ganger, Garth Gibson



Machine Learning: -- a view from outside





Inside ML ...









Massive Data



1B+ USERS 30+ PETABYTES



32 million pages

You Tube

100+ hours video uploaded every minute 645 million users 500 million tweets / day

twitter3



Growing Model Complexity



Google Brain Deep Learning for images: 1~10 Billion model parameters Multi-task Regression for simplest wholegenome analysis: 100 million ~ 1 Billion model parameters



Collaborative filteringfor Video recommendation:1~10 Billion1~10 BillionmodelNETFUNXparameters

6



centers won't fit into memory of single machine

7





Naïve MapReduce not best for ML

- Hadoop can execute iterative-convergent, data-parallel ML...
 - \circ map() to distribute data samples i, compute update $\Delta(D_i)$
 - reduce() to combine updates $\Delta(D_i)$
 - Iterative ML algo = repeat map()+reduce() again and again
- But reduce() writes to HDFS before starting next iteration's map() very slow iterations!





Spark: Faster MapR on Data-Parallel

- Spark's solution: Resilient Distributed Datasets (RDDs)
 - \circ Input data \rightarrow load as RDD \rightarrow apply transforms \rightarrow output result
 - RDD transforms strict superset of MapR
 - RDDs cached in memory, avoid disk I/O



• Spark ML library supports data-parallel ML algos, like Hadoop

- Spark and Hadoop: comparable first iter timings...
- o But Spark's later iters are much faster

10





Intrinsic Properties of ML Programs [Xing et al., 2015]





4 Principles of ML System Design

How to execute distributed-parallel ML programs? ML program equations tell us "What to Compute". But...

- 1. How to Distribute?
- 2. How to Bridge Computation and Communication?
- 3. How to Communicate?
- 4. What to Communicate?



Principles of ML system Design [Xing et al., to appear 2016]

1. How to Distribute:

Scheduling and Balancing workloads





Example: Model Distribution

Lasso via coordinate descent:





- How to correctly divide computational workload across workers?
- What is the best order to update parameters?



Model Dependencies

• Concurrent updates of β may induce errors





Avoid Dependency Errors via Structure-Aware Parallelization (SAP)

[Lee et al., 2014] [Kim et al, 2016]





A Structure-aware Dynamic Scheduler (Strads) [Lee et al., 2014] [Kim et al, 2016]





SAP Scheduling: Faster, Better Convergence across algorithms

• SAP on Strads achieves better speed and objective





$$\mathbb{E}\left[f(X^{(t)}) - f(X^*)\right]$$

 $\frac{\mathcal{O}(M)}{\mathcal{O}(P^2\rho)}\frac{1}{t} = \mathcal{O}\left(\frac{1}{Pt}\right)$

where t is # of iterations

Take-away: SAP minimizes ρ by searching for feature subsets $X^{(1)}$, $X^{(2)}$, ..., $X^{(B)}$ w/o cross-correlation => as close to P-fold speedup as possible



Correctly Measuring Parallel

Performance [blinded, to appear]

YahooLDA progress per iteration



80GB data, 2M words, 1K topics, <mark>100 machines</mark>

	YahooLDA data throughput
25 machines	39.7 M/s (1x)
50 machines	78 M/s (1.96x)
100 machines	151 M/s (3.8x)

- YahooLDA attains near-ideal throughput $(1 \rightarrow 3.8x)...$
- ... but progress per iteration gets worse with more machines
- YahooLDA only <2x speedup from 25 \rightarrow 100 machines
 - 6.7x slower compared to SAP-LDA



Correctly Measuring Parallel

Performance [blinded, to appear]

SAP-LDA progress per iteration



80GB data, 2M words, 1K topics, <mark>100 machines</mark>

	SAP-LDA data throughput
25 machines	58.3 M/s (1x)
50 machines	114 M/s (1.96x)
100 machines	204 M/s (3.5x)

- Ideal rate: progress per iter preserved from $25 \rightarrow 100$ machines
 - Thanks to dependency checking
- Near-ideal throughput: data rate $1x \rightarrow 3.5x$ from $25 \rightarrow 100$ machines
 - Thanks to load balancing
- Convergence Speed = rate x throughput
 - Therefore near-ideal 3.5x speedup from $25 \rightarrow 100$ machines

Principles of ML system Design [Xing et al., to appear 2016]

2. How to Bridge Computation and Communication: Bridging Models and Bounded Asynchrony





The Bulk Synchronous Parallel Bridging Model [Valiant & McColl]



- Perform barrier in order to communicate parameters
- Mimics sequential computation "serializable" property
- Enjoys same theoretical guarantees as sequential execution



The Bulk Synchronous Parallel Bridging Model [Valiant & McColl]



The success of the von Neumann model of sequential computation is attributable to the fact it is an efficient bridge between software and hardware... an analogous bridge is required for parallel computation if that is to become as widely used – Leslie G. Valiant

- Numerous implementations since 90s (list by Bill McColl):
 - Oxford BSP Toolset ('98), Paderborn University BSP Library ('01), Bulk Synchronous Parallel ML ('03), BSPonMPI ('06), ScientificPython ('07), Apache Hama ('08), Apache Pregel ('09), MulticoreBSP ('11), BSPedupack ('11), Apache Giraph ('11), GoldenOrb ('11), Stanford GPS Project ('11) ...

But There Is No Ideal Distributed System!

• Two distributed challenges:

Networks are slow

SAN UNG

• "Identical" machines rarely perform equally

Result: BSP barriers can be slow



Is there a better way to interleave computation and communication?

Safe/slow (BSP) vs. Fast/risky (Async)?

S12/1412(G

- Challenge 1: Need "Partial" synchronicity
 - Spread network comms evenly (don't sync unless needed)
 - Threads usually shouldn't wait but mustn't drift too far apart!

Challenge 2: Need straggler tolerance

Slow threads must somehow catch up







Stale Synchronous Parallel (SSP)

• Fastest/slowest workers not allowed to drift >s iterations apart

Consequence

- Fast like async, yet correct like BSP
- Why? Workers' local view of model parameters "not too stale" ($\leq s$ iterations old)





SSP Data-Parallel Async Speed, BSP Guarantee



Lasso

Matrix Fact.



- Massive Data Parallelism
- Effective across different algorithms

T $\sqrt{T} \left(\frac{\eta L}{\tau} + \eta + \frac{2\eta L}{\eta} \frac{\mu_{\gamma}}{\tau} \right) \geq \tau \left[-\frac{1}{\tau} \exp \left(2\bar{\eta}_T \sigma_{\gamma} + \frac{2}{2} \eta L^2 (2s+1) P \tau \right) \right]$

SSP Data Parallel Convergence Theorem [Ho et al., 2013, Dai et al., 2015]

ମ

Let observed staleness be γt Let staleness mean, variance be $\mu_{\gamma} = \mathbb{E}[\gamma_t]$, $\sigma_{\gamma} = var(\gamma_t)$

Theorem: Given L-Lipschitz objective f_t and step size h_t ,

$$P\left[\frac{R[X]}{T} - \frac{\mathcal{O}(F^2 + \mu_{\gamma}L^2)}{\sqrt{T}} \ge \tau\right] \le \exp\left\{\frac{-T\tau^2}{\mathcal{O}(\bar{\eta}_T\sigma_{\gamma} + L^2sP\tau)}\right\}$$

where
$$R[X] := \sum_{t=1}^T f_t(\tilde{x}_t) - f(x^*) \qquad \bar{\eta}_T = \frac{\eta^2 L^4(\ln T + 1)}{T} = o(T)$$

Explanation: the distance between true optima and current estimate decreases exponentially with more SSP iterations. *Lower staleness mean, variance* μ_{γ} , σ_{γ} *improve the convergence rate.*





Instical Artificial Intelligence & Megrative Genomics

Proximal Gradient under SSP

• Model (e.g. SVM, Lasso ...):

 $\min_{\mathbf{a} \in \mathbb{R}^d} \mathcal{L}(\mathbf{a}, \underline{D}), \text{ where } \mathcal{L}(\mathbf{a}, D) = f(\mathbf{a}, D) + g(\mathbf{a})$ data *D*, model *a*

- Model parallel
 - Model dimension d too large to fit in a single worker
 - Divide model among *P* workers $\mathbf{a} = (a_1, a_2, \dots, a_P)$

• Algorithm:
$$\forall p, a_p(t+1) = a_p(t) + \underbrace{\gamma_p(t)}_{t} \cdot F_p(\mathbf{a}^p(t))$$
 workers can skip updates

$$= a_p(0) + \sum_{k=0}^{t} \gamma_p(k) \cdot F_p(\mathbf{a}^p(t))$$
staleness
(local) $\mathbf{a}^p(t) = (a_1(\tau_1^p(t)), \dots, a_P(\underline{\tau_P^p(t)}))$
(global) $\mathbf{a}(t) = (a_1(t), \dots, a_P(t)).$

$$\mathbf{a}^p(t+1) := F_p(\mathbf{a}^p(t)) = \operatorname{prox}_{g_p}^{\eta}(a_p(t) - \eta \nabla_p f(\mathbf{a}^p(t))) - a_p(t)$$
proximal step wrt g

worker p keeps local copy of the full model (can be avoided for linear models)



SSP Model-Parallel Async Speed, BSP Guarantee

Lasso: 1M samples, 100M features, 100 machines



- Massive Model Parallelism
- Effective across different algorithms



SSP Model Parallel Convergence Theorem [Zhou et al., 2016]

Theorem: Given that the SSP delay is bounded, with appropriate step size and under mild technical conditions, then

$$\sum_{t=0}^{\infty} \|\mathbf{a}(t+1) - \mathbf{a}(t)\| < \infty \qquad \sum_{t=0}^{\infty} \|\mathbf{a}^{p}(t+1) - \mathbf{a}^{p}(t)\| < \infty$$

In particular, the global and local sequences converge to the same critical point, with rate $O(t^{-1})$:

$$\mathcal{L}\left(\frac{1}{t}\sum_{k=1}^{t}\mathbf{a}(k)\right) - \inf\mathcal{L} \leq O\left(t^{-1}\right)$$

Explanation: Finite length guarantees that the algorithm stops (the updates must eventually go to zero). Furthermore, the algorithm converges at rate $O(t^{-1})$ to the optimal value; same as BSP model parallel.



Principles of ML system Design [Xing et al., to appear 2016]

3. How to Communicate:

Managed Communication and Topologies





Managed Communication [Wei et al., 2015]

- SSP only
 - Communicates only at iteration boundary
 - Ensures bounded staleness consistency







- Matrix Factorization, Netflix data, rank = 400
- 8 machines * 16 cores, 1GbE ethernet



- Latent Dirichlet Allocation, NYTimes, # topics = 1000,
- 16 machines * 16 cores, 1GbE ethernet



Topology: Master-Slave



- Used with **centralized storage** paradigm
- Topology = **bipartite graph**: Servers (masters) to Workers (slaves)
- **Disadvantage:** need to code/manage clients and servers separately
- Advantage: bipartite topology far smaller than full N² P2P connections



Topology: Peer-to-Peer (P2P)

worker 1 worker 2 ML App Model copy ML App Model copy ML App Model copy ML App Model copy ML App Model copy

worker 3

worker 4

- Used with **decentralized storage** paradigm
- Workers update local parameter view by broadcasting/receiving
- Disadvantage: expensive unless updates ΔW are lightweight; expensive for large # of workers
- Advantage: only need worker code (no central server code); if ∆W is low rank, comms reduction possible



Halton Sequence Topology [Li et al., 2015]



- Used with decentralized storage paradigm
- Like P2P topology, but route messages through many workers
 - e.g. to send message from 1 to 6, use 1->2->3->6
- **Disadvantage:** incur higher SSP staleness due to routing, e.g. 1->2->3->6 = staleness 3
- Advantage: support bigger messages; support more machines than P2P topology



Random Partial Broadcasting and Diverse Mini-Batch Selection

- Random Partial Broadcasting
 - Each machine randomly selects Q<<P machines to send messages (instead of full broadcast)
 - Message cost reduced: from O(P²) to O(PQ), scales linearly with machine count P!



- Diverse Mini-batch Selection
 - Choose training data samples that maximize diversity score

$$s = \frac{1}{K(K-1)} \sum_{k=1}^{K} \sum_{j \neq i}^{K} (-\mathbf{x}_{k}^{\top} \mathbf{x}_{j})$$

- Pick few diverse samples instead of many random samples
- Using few but diverse samples further reduces comms costs without hurting output quality!

Principles of ML system Design [Xing et al., to appear 2016]

4. What to Communicate: Exploiting Structure in ML Updates





Matrix-Parameterized Models (MPMs)



Distance Metric Learning, Sparse Coding, Distance Metric Learning, Group Lasso, Neural Network, etc.



Big MPMs

Multiclass Logistic Regression on Wikipedia



#classes=325K

Billions of params = 10-100 GBs, costly network synchronization

What do we actually need to communicate?

Distance Metric Learning on ImageNet





.3B

Latent dim. = 50K

Sparse Coding on ImageNet

Feature dim. = 172K



Dic. Size=50K



#neurons in layer 1 = 33K



Full Updates

- Let matrix parameters be *W*. Need to send parallel worker updates ∠*W* to other machines...
 - Primal stochastic gradient descent (SGD)

$$\min_{W} \frac{1}{N} \sum_{i=1}^{N} f_i(Wa_i; b_i) + h(W)$$
$$\Delta W = \frac{\partial f(Wa_i, b_i)}{\partial W}$$

• Stochastic dual coordinate ascent (SDCA)

$$\min_{Z} \frac{1}{N} \sum_{i=1}^{N} f_i^* (-z_i) + h^* (\frac{1}{N} Z A^T)$$
$$\Delta W = (\Delta z_i) a_i$$



Sufficient Factor (SF) Updates [Xie et al., 2015]

- Full parameter matrix update ∠W can be computed as outer product of two vectors uv^T (called sufficient factors)
 - Primal stochastic gradient descent (SGD)

$$\min_{W} \frac{1}{N} \sum_{i=1}^{N} f_i(Wa_i; b_i) + h(W)$$

$$\Delta W = uv^{\mathrm{T}} \quad u = \frac{\partial f(Wa_i, b_i)}{\partial (Wa_i)} \quad v = a_i$$

• Stochastic dual coordinate ascent (SDCA)

$$\min_{Z} \frac{1}{N} \sum_{i=1}^{N} f_i^*(-z_i) + h^*(\frac{1}{N} Z A^T)$$
$$\Delta W = uv^T \quad u = \Delta z_i \quad v = a_i$$

Send the lightweight SF updates (*u*,*v*), instead of the expensive full-matrix ∠W updates!





SFB Convergence Theorem

[Xie et al., 2015]

Theorem 1. Let $\{\mathbf{W}_p^c\}$, p = 1, ..., P, and $\{\mathbf{W}^c\}$ be the local sequences and the auxiliary sequence generated by SFB for problem (P) (with $h \equiv 0$), respectively. Under Assumption 1 and set the learning rate $\eta_c^{-1} = \frac{L_F}{2} + 2sL + \sqrt{c}$, then we have

- $\liminf_{c \to \infty} \mathbb{E} \|\nabla F(\mathbf{W}^c)\| = 0$, hence there exists a subsequence of $\nabla F(\mathbf{W}^c)$ that almost surely vanishes;
- $\lim_{c \to \infty} \max_p \|\mathbf{W}^c \mathbf{W}_p^c\| = 0$, i.e. the maximal disagreement between all local sequences and the auxiliary sequence converges to 0 (almost surely);
- There exists a common subsequence of $\{\mathbf{W}_p^c\}$ and $\{\mathbf{W}_p^c\}$ that converges almost surely to a stationary point of F, with the rate $\min_{c \leq C} \mathbb{E} \|\sum_{p=1}^P \nabla F_p(\mathbf{W}_p^c)\|_2^2 \leq O\left(\frac{(L+L_F)\sigma^2 P_s \log C}{\sqrt{C}}\right).$

Explanation: Parameter copies W_p on different workers p converge to the same optima, *i.e. all workers reach the same (correct) answer.*

✓ Does not need central parameter server or key-value store

 \checkmark Works with SSP bridging model (staleness = s)



Why is SFB faster?





Near-linear scalability





SFB communication up to 100x smaller than PS and Spark





Summary

1. How to Distribute?

- Structure-Aware Parallelization
- Work Prioritization



2. How to Bridge Computation and Communication?

- BSP Bridging Model
- SSP Bridging Model for Data and Model Parallel

3. How to Communicate?

- Managed comms interleave comms/compute, prioritized comms
- Parameter Storage: Centralized vs Decentralized
- Communication Topologies: Master-Slave, P2P, Halton Sequence

4. What to Communicate?

- Full Matrix updates
- Sufficient Factor updates
- Hybrid FM+SF updates (as in a DL model)