

Manifold Identification for Ultimately Communication-Efficient Distributed Optimization

Yu-Sheng Li



國立臺灣大學
National Taiwan University

Joint work with Wei-Lin Chiang (NTU) and Ching-pei Lee (NUS)

Outline

Overview

Empirical Risk Minimization

The Proposed Algorithm

Experiments

Distributed Machine Learning

Read 1 MB sequentially from memory	3 μ s
Read 1 MB sequentially from network	22 μ s
Read 1 MB sequentially from disk (SSD)	49 μ s
Round trip in the same datacenter	500 μ s

(Latency Numbers Every Programmer Should Know.¹)

¹Originally by Jeff Dean in 2010, updated by Colin Scott at
https://colin-scott.github.io/personal_website/research/interactive_latency.html

Distributed Machine Learning

Read 1 MB sequentially from memory	3 μ s
Read 1 MB sequentially from network	22 μ s
Read 1 MB sequentially from disk (SSD)	49 μ s
Round trip in the same datacenter	500 μ s

(*Latency Numbers Every Programmer Should Know.*¹)

- ▶ **Inter-machine communication** may be more time-consuming than local computations within a machine

$$\text{Comm. cost} = (\# \text{ Comm. rounds}) \times (\text{Bytes communicated per round})$$

¹Originally by Jeff Dean in 2010, updated by Colin Scott at https://colin-scott.github.io/personal_website/research/interactive_latency.html

Sparsity-inducing Regularization

- ▶ To avoid overfitting and to force some desired structure of the solution, usually a sparsity-inducing regularizer is introduced

Sparsity-inducing Regularization

- ▶ To avoid overfitting and to force some desired structure of the solution, usually a sparsity-inducing regularizer is introduced
- ▶ Example: l_2 - vs. l_1 -regularized logistic regression on **news20**

	Relative reg. strength	Sparsity of solution	Test accuracy
l_2 -regularized			
	2^0	1,355,191 (100%)	99.7449%
	2^{10}	1,355,191 (100%)	97.0044%

Sparsity-inducing Regularization

- ▶ To avoid overfitting and to force some desired structure of the solution, usually a sparsity-inducing regularizer is introduced
- ▶ Example: l_2 - vs. l_1 -regularized logistic regression on **news20**

	Relative reg. strength	Sparsity of solution	Test accuracy
l_2 -regularized			
	2^0	1,355,191 (100%)	99.7449%
	2^{10}	1,355,191 (100%)	97.0044%
l_1 -regularized			
	2^0	67,071 (4.95%)	99.7499%
	2^2	42,020 (3.10%)	99.7499%
	2^4	14,524 (1.07%)	99.7449%
	2^6	5,432 (0.40%)	99.6749%
	2^8	1,472 (0.11%)	97.3495%
	2^{10}	546 (0.04%)	92.8936%

Our contributions

Recall:

$$\text{Comm. cost} = (\# \text{ Comm. rounds}) \times (\text{Bytes communicated per round})$$

Our contributions

Recall:

$$\text{Comm. cost} = (\# \text{ Comm. rounds}) \times (\text{Bytes communicated per round})$$

- ▶ Focusing on the small subproblem
⇒ fewer bytes to communicate

Our contributions

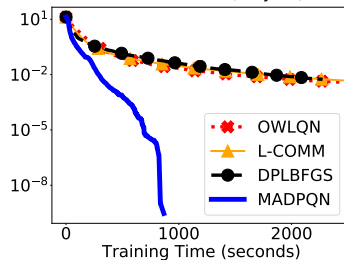
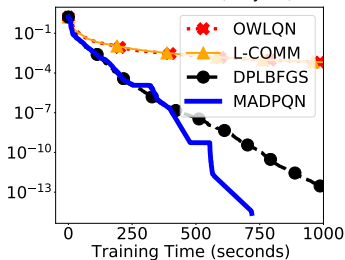
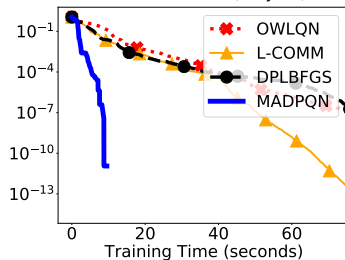
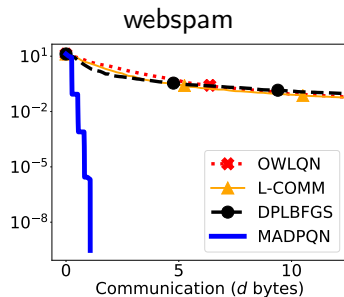
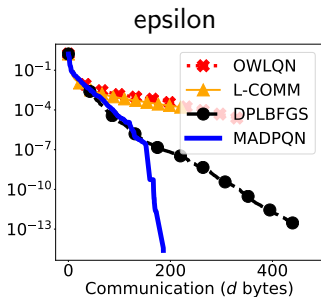
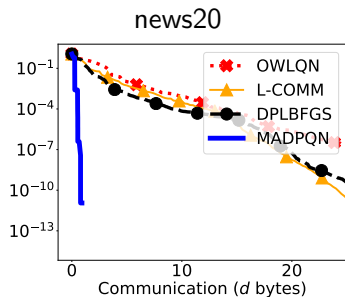
Recall:

$$\text{Comm. cost} = (\# \text{ Comm. rounds}) \times (\text{Bytes communicated per round})$$

- ▶ Focusing on the small subproblem
⇒ fewer bytes to communicate
- ▶ Acceleration by smooth optimization in the correct manifold
⇒ fewer rounds of communication

Results (ours: MADPQN)

y-axis: relative distance to the optimal value (log-scaled)
x-axis: communication costs (upper), training time (lower)



Outline

Overview

Empirical Risk Minimization

The Proposed Algorithm

Experiments

Outline

Overview

Empirical Risk Minimization

The Proposed Algorithm

Experiments

Distributed Empirical Risk Minimization (ERM)

- ▶ Train a model by minimizing a function that measures the performance on training data

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) := \sum_{k=1}^K f_k(\mathbf{w})$$

- ▶ There are K machines, and f_k is exclusively available on machine k

Distributed Empirical Risk Minimization (ERM)

- ▶ Train a model by minimizing a function that measures the performance on training data

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) := \sum_{k=1}^K f_k(\mathbf{w})$$

- ▶ There are K machines, and f_k is exclusively available on machine k
- ▶ Synchronize \mathbf{w} or $\nabla f(\mathbf{w})$ by communication: communication cost per iteration is $O(d)$
- ▶ How to reduce the $O(d)$ cost?

Sparsity-inducing Regularizer

- ▶ If \mathbf{w} is sparse throughout the training process, we only need to synchronize a shorter vector
- ▶ Regularized ERM:

$$\min_{\mathbf{w}} f(\mathbf{w}) + R(\mathbf{w})$$

Sparsity-inducing Regularizer

- ▶ If \mathbf{w} is sparse throughout the training process, we only need to synchronize a shorter vector
- ▶ Regularized ERM:

$$\min_{\mathbf{w}} f(\mathbf{w}) + R(\mathbf{w})$$

- ▶ An ideal regularization term for forcing sparsity is the ℓ_0 norm:

$$\|\mathbf{w}\|_0 = \text{number of nonzeros in } \mathbf{w}$$

Sparsity-inducing Regularizer

- ▶ If \mathbf{w} is sparse throughout the training process, we only need to synchronize a shorter vector
- ▶ Regularized ERM:

$$\min_{\mathbf{w}} f(\mathbf{w}) + R(\mathbf{w})$$

- ▶ An ideal regularization term for forcing sparsity is the ℓ_0 norm:

$$\|\mathbf{w}\|_0 = \text{number of nonzeros in } \mathbf{w}$$

- ▶ But this norm is not continuous and hence hard to optimize
- ▶ A good surrogate is the ℓ_1 norm $\|\mathbf{w}\|_1 = \sum_{i=1}^d |w_i|$
- ▶ Our algorithm works for other partly smooth R , e.g. group-LASSO

The Regularized Problem

- ▶ Now the problem becomes

$$\min_{\mathbf{w}} f(\mathbf{w}) + \|\mathbf{w}\|_1,$$

which is harder to minimize than $f(\mathbf{w})$ alone since $\|\mathbf{w}\|_1$ is **not differentiable**

- ▶ As the gradient may not even exist, gradient descent or Newton method cannot be directly applied

Proximal Quasi-Newton

- ▶ Proximal gradient is a simple algorithm that solves

$$\min_{\mathbf{w}'} \nabla f(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w}) + \frac{1}{2\alpha} \|\mathbf{w}' - \mathbf{w}\|_2^2 + \|\mathbf{w}'\|_1,$$

where α is the step size for the current iteration

- ▶ Each calculation of ∇f requires one round of communication

Proximal Quasi-Newton

- ▶ Proximal gradient is a simple algorithm that solves

$$\min_{\mathbf{w}'} \nabla f(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w}) + \frac{1}{2\alpha} \|\mathbf{w}' - \mathbf{w}\|_2^2 + \|\mathbf{w}'\|_1,$$

where α is the step size for the current iteration

- ▶ Each calculation of ∇f requires one round of communication
- ▶ To reduce the amount of communication, we include some second-order information:

reducing iterations \Rightarrow reducing rounds of communication

- ▶ Replace the term $\|\mathbf{w}' - \mathbf{w}\|_2^2/2\alpha$ with $(\mathbf{w}' - \mathbf{w})^\top H(\mathbf{w}' - \mathbf{w})/2$ for some $H \approx \nabla^2 f(\mathbf{w})$

Outline

Overview

Empirical Risk Minimization

The Proposed Algorithm

Experiments

Utilizing Sparsity

- ▶ Even if we only update the nonzero entries of \boldsymbol{w} , if we still compute the whole gradient $\nabla f(\boldsymbol{w})$, then the communication cost remains $O(d)$

Utilizing Sparsity

- ▶ Even if we only update the nonzero entries of \mathbf{w} , if we still compute the whole gradient $\nabla f(\mathbf{w})$, then the communication cost remains $O(d)$
- ▶ Guess: if $w_i = 0$ at some iteration and it is likely to stay 0 at the next iteration, it remains 0 at the final solution
- ▶ Then we only solve the subproblem with respect to the coordinates that are likely to be nonzero

Utilizing Sparsity

- ▶ Even if we only update the nonzero entries of \mathbf{w} , if we still compute the whole gradient $\nabla f(\mathbf{w})$, then the communication cost remains $O(d)$
- ▶ Guess: if $w_i = 0$ at some iteration and it is likely to stay 0 at the next iteration, it remains 0 at the final solution
- ▶ Then we only solve the subproblem with respect to the coordinates that are likely to be nonzero
- ▶ A progressive shrinking approach: once we guess $w_i = 0$, we remove those coordinates from our problem in future iterations
- ▶ So the number of nonzeros in \mathbf{w} (i.e. $\|\mathbf{w}\|_0$) gradually decreases

Convergence Issue

- ▶ What if our guess was wrong at some iteration?

Convergence Issue

- ▶ What if our guess was wrong at some iteration?
- ▶ Need to double-check: when some stopping criterion is met, we restart with all coordinates
- ▶ Training is terminated only when our model can hardly be improved using all coordinates

More Acceleration by Smooth Optimization

- ▶ $|w_i|$ becomes twice-differentiable when $w_i \neq 0$
- ▶ If the coordinates where $w_i \neq 0$ are fixed, the proximal approach is not needed anymore
- ▶ The problem can then be transformed into a smooth one for faster convergence

More Acceleration by Smooth Optimization

- ▶ $|w_i|$ becomes twice-differentiable when $w_i \neq 0$
- ▶ If the coordinates where $w_i \neq 0$ are fixed, the proximal approach is not needed anymore
- ▶ The problem can then be transformed into a smooth one for faster convergence
- ▶ When the nonzero pattern (manifold) does not change for some iterations, it is likely to be the final pattern

More Acceleration by Smooth Optimization

- ▶ $|w_i|$ becomes twice-differentiable when $w_i \neq 0$
- ▶ If the coordinates where $w_i \neq 0$ are fixed, the proximal approach is not needed anymore
- ▶ The problem can then be transformed into a smooth one for faster convergence
- ▶ When the nonzero pattern (manifold) does not change for some iterations, it is likely to be the final pattern
- ▶ Example with $d = 5$:

$$\{1, 2, 3, 4, 5\}$$

More Acceleration by Smooth Optimization

- ▶ $|w_i|$ becomes twice-differentiable when $w_i \neq 0$
- ▶ If the coordinates where $w_i \neq 0$ are fixed, the proximal approach is not needed anymore
- ▶ The problem can then be transformed into a smooth one for faster convergence
- ▶ When the nonzero pattern (manifold) does not change for some iterations, it is likely to be the final pattern
- ▶ Example with $d = 5$:

$$\{1, 2, 3, 4, 5\} \rightarrow \{2, 3, 5\}$$

More Acceleration by Smooth Optimization

- ▶ $|w_i|$ becomes twice-differentiable when $w_i \neq 0$
- ▶ If the coordinates where $w_i \neq 0$ are fixed, the proximal approach is not needed anymore
- ▶ The problem can then be transformed into a smooth one for faster convergence
- ▶ When the nonzero pattern (manifold) does not change for some iterations, it is likely to be the final pattern
- ▶ Example with $d = 5$:

$$\{1, 2, 3, 4, 5\} \rightarrow \{2, 3, 5\} \rightarrow \{2, 5\}$$

More Acceleration by Smooth Optimization

- ▶ $|w_i|$ becomes twice-differentiable when $w_i \neq 0$
- ▶ If the coordinates where $w_i \neq 0$ are fixed, the proximal approach is not needed anymore
- ▶ The problem can then be transformed into a smooth one for faster convergence
- ▶ When the nonzero pattern (manifold) does not change for some iterations, it is likely to be the final pattern
- ▶ Example with $d = 5$:

$$\{1, 2, 3, 4, 5\} \rightarrow \{2, 3, 5\} \rightarrow \{2, 5\} \rightarrow \{2, 5\}$$

More Acceleration by Smooth Optimization

- ▶ $|w_i|$ becomes twice-differentiable when $w_i \neq 0$
- ▶ If the coordinates where $w_i \neq 0$ are fixed, the proximal approach is not needed anymore
- ▶ The problem can then be transformed into a smooth one for faster convergence
- ▶ When the nonzero pattern (manifold) does not change for some iterations, it is likely to be the final pattern
- ▶ Example with $d = 5$:

$$\{1, 2, 3, 4, 5\} \rightarrow \{2, 3, 5\} \rightarrow \{2, 5\} \rightarrow \{2, 5\} \rightarrow \{2, 5\}$$

More Acceleration by Smooth Optimization

- ▶ $|w_i|$ becomes twice-differentiable when $w_i \neq 0$
- ▶ If the coordinates where $w_i \neq 0$ are fixed, the proximal approach is not needed anymore
- ▶ The problem can then be transformed into a smooth one for faster convergence
- ▶ When the nonzero pattern (manifold) does not change for some iterations, it is likely to be the final pattern
- ▶ Example with $d = 5$:

$$\{1, 2, 3, 4, 5\} \rightarrow \{2, 3, 5\} \rightarrow \{2, 5\} \rightarrow \{2, 5\} \rightarrow \{2, 5\} \xrightarrow{\text{accelerate}} \dots$$

More Acceleration by Smooth Optimization

- ▶ $|w_i|$ becomes twice-differentiable when $w_i \neq 0$
- ▶ If the coordinates where $w_i \neq 0$ are fixed, the proximal approach is not needed anymore
- ▶ The problem can then be transformed into a smooth one for faster convergence
- ▶ When the nonzero pattern (manifold) does not change for some iterations, it is likely to be the final pattern
- ▶ Example with $d = 5$:

$$\begin{array}{c} \{1, 2, 3, 4, 5\} \rightarrow \{2, 3, 5\} \rightarrow \{2, 5\} \rightarrow \{2, 5\} \rightarrow \{2, 5\} \xrightarrow{\text{accelerate}} \dots \\ \xrightarrow{\text{restart}} \{1, 2, 3, 4, 5\} \rightarrow \dots \end{array}$$

More Acceleration by Smooth Optimization

- ▶ $|w_i|$ becomes twice-differentiable when $w_i \neq 0$
- ▶ If the coordinates where $w_i \neq 0$ are fixed, the proximal approach is not needed anymore
- ▶ The problem can then be transformed into a smooth one for faster convergence
- ▶ When the nonzero pattern (manifold) does not change for some iterations, it is likely to be the final pattern
- ▶ Example with $d = 5$:

$$\begin{aligned} & \{1, 2, 3, 4, 5\} \rightarrow \{2, 3, 5\} \rightarrow \{2, 5\} \rightarrow \{2, 5\} \rightarrow \{2, 5\} \xrightarrow{\text{accelerate}} \dots \\ \xrightarrow{\text{restart}} & \{1, 2, 3, 4, 5\} \rightarrow \dots \\ \xrightarrow{\text{restart}} & \{1, 2, 3, 4, 5\} \rightarrow \dots \end{aligned}$$

More Acceleration by Smooth Optimization

- ▶ $|w_i|$ becomes twice-differentiable when $w_i \neq 0$
- ▶ If the coordinates where $w_i \neq 0$ are fixed, the proximal approach is not needed anymore
- ▶ The problem can then be transformed into a smooth one for faster convergence
- ▶ When the nonzero pattern (manifold) does not change for some iterations, it is likely to be the final pattern
- ▶ Example with $d = 5$:

$$\begin{aligned} & \{1, 2, 3, 4, 5\} \rightarrow \{2, 3, 5\} \rightarrow \{2, 5\} \rightarrow \{2, 5\} \rightarrow \{2, 5\} \xrightarrow{\text{accelerate}} \dots \\ \xrightarrow{\text{restart}} & \{1, 2, 3, 4, 5\} \rightarrow \dots \\ \xrightarrow{\text{restart}} & \{1, 2, 3, 4, 5\} \rightarrow \dots \\ \xrightarrow{\text{restart}} & \{1, 2, 3, 4, 5\} \rightarrow \text{terminated} \end{aligned}$$

Theoretical Guarantees

Theorem

If a cluster point \mathbf{w}^ of $\{\mathbf{w}$ after each restart $\}$ satisfies*

$$\mathbf{0} \in \text{relint}(\nabla f(\mathbf{w}^*) + \partial R(\mathbf{w}^*)),$$

then the manifold of \mathbf{w}^ will be identified within finite restarts.*

Outline

Overview

Empirical Risk Minimization

The Proposed Algorithm

Experiments

Settings

- ▶ We show the effectiveness of the proposed approach by ℓ_1 -regularized logistic regression

$$\min_{\mathbf{w}} \sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{x}_i^\top \mathbf{w})) + \|\mathbf{w}\|_1,$$

where there are n instances with features $\mathbf{x}_i \in \mathbb{R}^d$ and labels $y_i \in \{-1, 1\}$

Settings

- ▶ We show the effectiveness of the proposed approach by ℓ_1 -regularized logistic regression

$$\min_{\mathbf{w}} \sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{x}_i^\top \mathbf{w})) + \|\mathbf{w}\|_1,$$

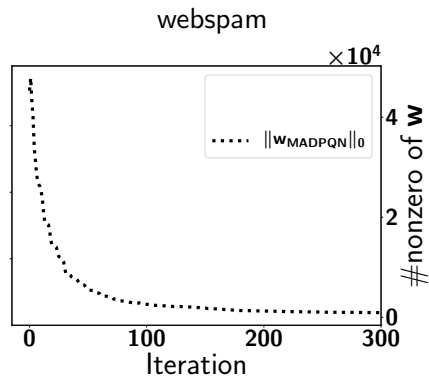
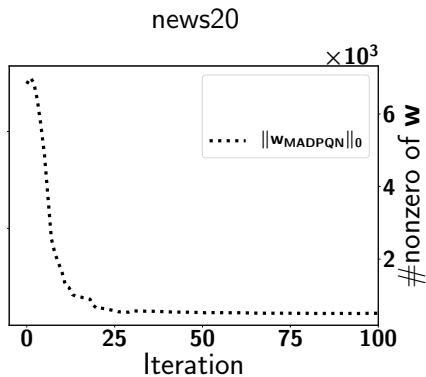
where there are n instances with features $\mathbf{x}_i \in \mathbb{R}^d$ and labels $y_i \in \{-1, 1\}$

- ▶ The instances are evenly split across $K = 10$ machines, connected by Intel MPI in a 1Gbps network environment

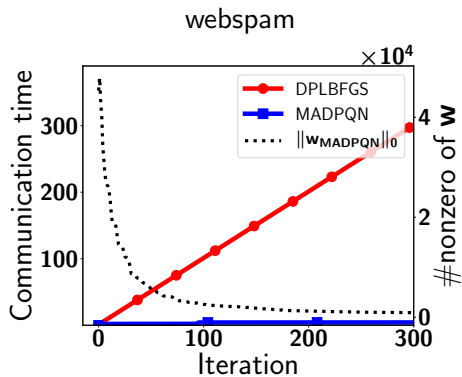
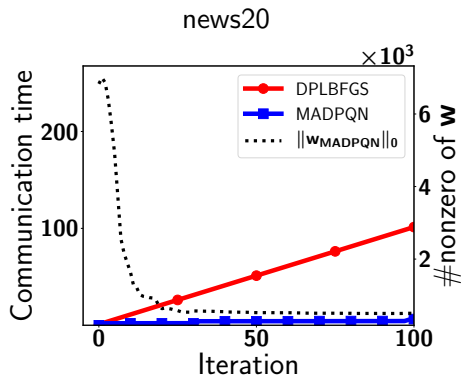
Data Statistics

Data set	Instances (n)	Features (d)	Nonzeros in optimal w^*
news20	19,996	1,355,191	506
epsilon	400,000	2,000	1,463
webspam	350,000	16,609,143	793
url	2,396,130	3,231,961	25,399
avazu-site	25,832,830	999,962	11,858
KDD2010-b	19,264,097	29,890,096	2,005,632

Results



Results



- ▶ DPLBFGS:
a distributed proximal quasi-Newton method (Lee et al. 2019)
- ▶ Manifold-Aware Distributed Proximal Quasi-Newton (MADPQN):
DPLBFGS + manifold selection + further acceleration

Comparison with state of the art

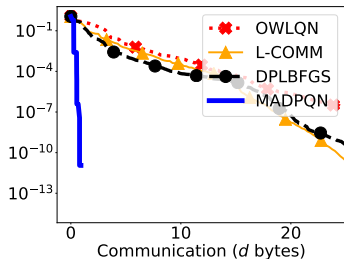
- ▶ OWLQN ([Andrew and Gao 2007](#)):
an extension of a quasi-Newton method, LBFGS, which is the most commonly used distributed method
- ▶ L-COMM ([Chiang et al. 2018](#)):
an extension of the common directions method ([Wang et al. 2016](#))
- ▶ DPLBFGS ([Lee et al. 2019](#)):
a distributed proximal LBFGS method
- ▶ MADPQN:
Our proposed Manifold-Aware Distributed Proximal Quasi-Newton method

Results

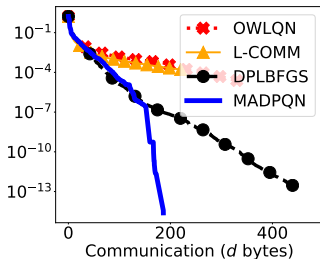
y-axis: relative distance to the optimal value (log-scaled)

x-axis: communication costs (upper), training time (lower)

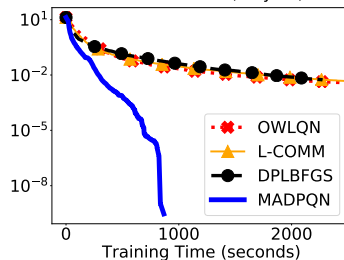
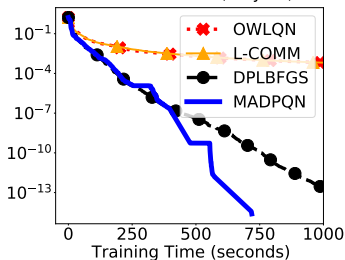
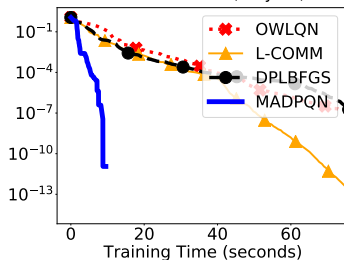
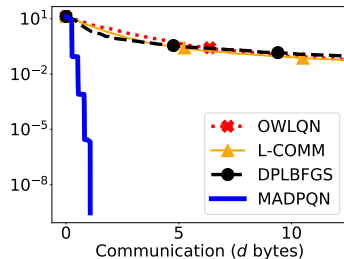
news20



epsilon



webspam

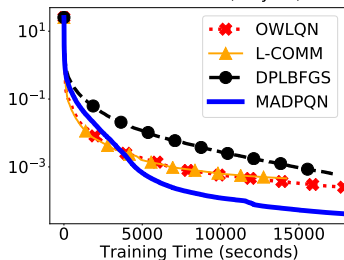
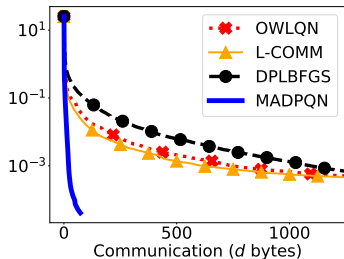


Results

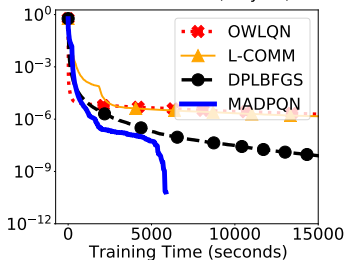
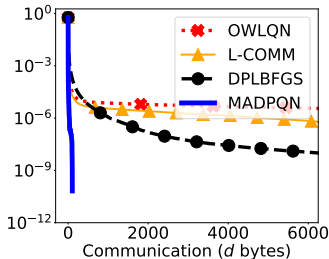
y-axis: relative distance to the optimal value (log-scaled)

x-axis: communication costs (upper), training time (lower)

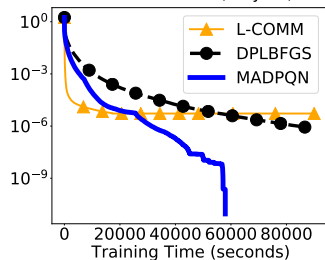
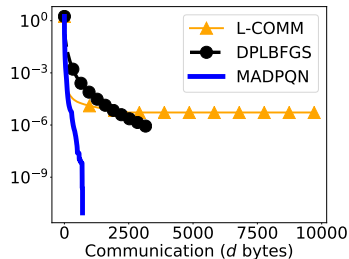
url



avazu-site



KDD2010-b



Conclusions

- ▶ Communication may be the bottleneck in distributed machine learning
- ▶ Communication cost can be reduced by utilizing the sparsity pattern throughout training
- ▶ Second-order information further improves convergence in later stage
- ▶ Theoretical support on manifold identification and superlinear convergence
- ▶ Source code to be released soon