

Distributed Power Management under Limited Communication

Na Li Harvard University

Rutgers, 08/22/2017

Acknowledgment :

Harvard Univ: Guannan Qu, Chinwendu Enyioha, Vahid Tarokh
KTH: Sindri Magnusson, Carlo Fishchione
Caltech: Steven Low
NREL: Changhong Zhao
Univ. of Colorado, Boulder: Lijun Chen



<u>A Vision of Future(IoT)?</u>

All devices are **connected and coordinated** to

.....

- Maximize social welfare
- Satisfy operation constraints

Distributed Optimization

Devices communicate, compute decisions, & communicate, ... until reach an efficient point (Iterative, two-way comm)



Sensing, Communication, Computation



Sources: Gigaom

Communication Challenges

- Lack reliability
- Unaccepted delays
- Vulnerable to malicious attacks
- Leak privacy
- Limited bandwidth
 (e.g. Power Line Comm.)
- High deployment cost







How about reducing communication needs?

Reduce communication in power management

Extract information from physical measurements (Feedback)

- Load frequency control
- Power allocation in buildings/data centers

Recover information from local computation

• Quantized dual gradient for power allocation

This talk: Limited communication in power systems

- Extra information from physical measurements (Feedback)
 - Load frequency control
 - Power allocation in buildings/data centers

- Recover information from local computation
 - Quantized dual gradient for power allocation

Power Systems



Source: Graphic courtesy of North American Electric Reliability Corportion (NERC)

Optimal Load Control

- Balance total generation and load
- Keep frequency deviation small
- Minimize aggregate load disutility



Distributed Optimization (e.g. ADMM) Applies. But...





- Hard to get the real-time disturbance information
- Heavily relies on iterative communication

Can loads response in real-time and closed-loop?

> Network physical dynamics help!

Physical dynamics: Swing Dynamics



Variables denote the deviations from their reference (steady state) values

Network dynamics

DC approximation of power flow

$$\dot{P}_{ij} = b_{ij} \left(\omega_i - \omega_j \right)$$



- Lossless (resistance=0)
- Fixed voltage magnitudes $|V_i|$
- Small deviation of angles

System Model Recap

$$M_{i}\dot{\omega}_{i} = B_{i}^{b} - \sum_{l \in i} d_{l} - D_{i}\omega_{i} - \sum_{j:i \to j} P_{ij} + \sum_{k:k \to i} P_{ki}$$

$$\dot{P}_{ij} = b_{ij} (\omega_{i} - \omega_{j}) \text{Load Control}$$

$$P_{i}^{m} \rightarrow \sum_{l \in I} d_{l} + D_{i}\omega_{l}$$

$$(\omega_{i} - \omega_{j}) \text{Load Control}$$

$$P_{i}^{m} \rightarrow \sum_{l \in I} d_{l} + D_{i}\omega_{l}$$

$$(\omega_{i} - \omega_{j}) \text{Load Control}$$

$$P_{i}^{m} \rightarrow \sum_{l \in I} d_{l} + D_{i}\omega_{l}$$

$$(\omega_{i} - \omega_{j}) \text{Load Control}$$

$$P_{i}^{m} \rightarrow \sum_{l \in I} d_{l} + D_{i}\omega_{l}$$

$$(\omega_{i} - \omega_{j}) \text{Load Control}$$

$$P_{i}^{m} \rightarrow \sum_{l \in I} d_{l} + D_{i}\omega_{l}$$

i

Load frequency control



Converge to the optimal solution (Primal-Dual Gradient Flow)



Primal-Dual Gradient Flow: Arrow etc 1958, Feijer and Paganini 2010, Zhao, Low etc 2013, You, Chen etc 2014, Cherukuri, Mallada, Cortes, 2015, etc



Frequency: a locally measurable signal ("price" of imbalance)

Completely decentralized; no explicit communication necessary



Dynamic simulation of IEEE 68-bus system (New England)



- Power System Toolbox (RPI)
- Detailed generation model
- Exciter model, power system stabilizer model
- Nonzero resistance lines

Sample rate 250ms Step increase of loads on bus 1, 7, 27

Simulations









$$\begin{split} M_i \dot{\omega}_i &= P_i^m - \sum_{l \in i} d_l - D_i \omega_i - \sum_{j:i \to j} P_{ij} + \sum_{k:k \to i} P_{ki} \\ \dot{P}_{ij} &= b_{ij} \left(\omega_i - \omega_j \right) \\ \end{split}$$
Network Dynamics



 $P_i^m \xrightarrow{i}_{l \in i} d_l + D_i \omega_i$

This Idea Extends to General Systems

Network Dynamics:
$$\dot{x} = \sum_{j:i\sim j} A_{ij}x_j + B_iu_i + C_i\omega_i$$

How to design distributed,
closed-loop controller u?
Optimization: $\min_{x_i,u_i} \sum_i f_i(x_i) + \sum_i g_i(u_i)$
s. t. $\sum_{j:i\sim j} A_{ij}x_j + B_iu_i + C_iw_i = 0$
 $h_i(x_i, u_i) \le 0$

- [Li, Chen, Zhao, 2015]: Economic Automatic Generation Control
- [Zhang, Antonois, Li, 2016]: Sufficient and Necessary Conditions
- [Zhang, Malkawi, Li, 2016]: Thermal Control for HVAC

This talk: limited communication

- Load frequency control
- Decentralized voltage control (distribution network) (Qu, Li, Dahleh, 2014)
- Power allocation in buildings/data center

- Recover information from local computation
 - Quantized dual gradient for power allocation

This talk: limited communication

- Load frequency control
- Decentralized voltage control (distribution network) (Qu, Li, Dahleh, 2014)
- Power allocation in buildings/data center

- Recover information from local computation
 - Quantized dual gradient for power allocation

This talk: limited communication

- Load frequency control
- Decentralized voltage control (distribution network) (Qu, Li, Dahleh, 2014)
- Power allocation in buildings/data center

- Recover information from local computation
 - Quantized dual gradient for power allocation

Power management within buildings

Control center coordinates power consumption of appliances

- Maximize utility, minimize cost
- Satisfy operation constraints, e.g. power capacity constraints



Distributed Coordination under Two-way Comm.

Iterate Step 1: Appliances to center: Power request <u>Step 2</u>: Center to appliances: Coordination signal

Assume *perfect*, *reliable*, and *ubiquitous* communication resources



Reduce Communication Needs

<u>Q 1:</u> Is it possible to use only one-way comm.? <u>Q 2:</u> How many bits are needed?



Not just for the buildings/grids



Data Center

Communication cost is much higher than computation [Bolsens I., 2002]



Power allocation problem





A distributed algorithm: Dual gradient descent



A distributed algorithm: One-way comm.



Replace this with *true measurement* of total power consump. Q(t).

What's the problem here?



It might violates hard physical constraint $\sum q_i \leq Q$

Theorem: If the *step size* and *initial setting* are chosen properly, the constraint will hold all the time.

"Distributed resource allocation using one-way communication", Magnusson, Enyioha, Li, Fischione, Tarokh, 2016



- Load frequency control
- Power allocation in buildings/data center

- Recover information from local computation
 - Quantized dual gradient for power allocation



Extract information from physical measurements (Feedback)

- Load frequency control
- Power allocation in buildings/data center

Recover information from *local computation*

• Quantized dual gradient for power allocation

Recall: Dual Gradient with One-way Comm.



Further reduce comm.


Dual Gradient with One-bit One-way Comm.



This is quantized (normalized) gradient descent of dual function

Normalized Gradient
Descent [Shor 1985]:
$$\begin{array}{l} \min f(p) \\ p(t+1) = p(t) - \gamma(t) \frac{f'(p(t))}{||f'(p(t))||} \end{array}$$

Quantized (Normalized) Gradient Descent (QGD)

Problem:
$$\min_{p} f(p)$$

QGD: $p(t+1) = p(t) - \gamma(t)d(t), d(t) \in \mathcal{D} \in \mathcal{S}^{K-1}$

Question: When \mathcal{D} is a good quantization?

Definition: A quantization is proper (good) if and only if the algorithm is able to converge to the optimal points for any wellbehaviored *f*, *e.g.* convex smooth function.

Quantized (Normalized) Gradient Descent (QGD)

Problem:
$$\min_{p} f(p)$$

QGD: $p(t+1) = p(t) - \gamma(t)d(t), d(t) \in \mathcal{D} \in \mathcal{S}^{K-1}$

Questions:

- A) How to determine a quantization is proper?
- **B)** What is the minimal size of the quantization to be proper?
- C) How to choose d(t) and $\epsilon(t)$, given a good quantization?
- *D) What* are the *connections* between the *fineness* of the quantization to the *convergence* of the algorithm?

Descent direction



Proposition. A set \mathcal{D} is a proper quantization if and only if there exists $\theta \in [0, \pi/2)$, for every $g \in \mathcal{S}^{K-1}$ there exists $d \in \mathcal{D}$ such that $\arg(q, d) \leq \theta$.

Such quantization \mathcal{D} is also called as θ -cover.

- Minimal proper quantization, $|\mathcal{D}| = K + 1$.
- Only $\log_2(K+1)$ bits at each iteration.



Red: Quantization direction Blue: Gradient direction

"Convergence of limited communications gradient methods", Magnusson, Enyioha, Li, Fischione, Tarokh, Transactions on Automatic Control, 2017

Convergence rate

One Stopping Criterion: $\mathcal{P}(\epsilon) := \{ p \in \mathbb{R}^K : ||\nabla f(p)|| \le \epsilon \}$

Theorem: For any $\epsilon > 0$, if $\gamma \in (0, 2\cos(\theta)\epsilon/L)$, then there exists $T \in \mathbb{N}$ such that $p(T) \in \mathcal{P}(\epsilon)$, with T bounded by

$$T \le \frac{2(f(p(0)) - f^*)}{\gamma(2\cos(\theta)\epsilon - L\gamma)}.$$

- Finer quantization, larger stepsize is allowed
- Finer quantization, faster convergence

Quantization size



Simulation

f is the dual function of a two-dimension resource allocation problem.



(3) Infinite bandwidth: normalized gradient(2) Infinite bandwidth: gradient

Summary of QGD

Problem:
$$\min_{p} f(p)$$

QGD: $p(t+1) = p(t) - \gamma(t)d(t), d(t) \in \mathcal{D} \in \mathcal{S}^{K}$

- A) Proper quantization = θ -cover
- B) Minimal size of proper quantization is K+1
- C) Pick the quantized direction closest to gradient direction
- D) θ plays an important role in the convergence

"Convergence of limited communications gradient methods", Magnusson, Enyioha, Li, Fischione, Tarokh, Transactions on Automatic Control, 2017

Extension to Constrained Case

Problem:
$$\min_{p \in \mathcal{P}} f(p)$$

QGD: $p(t+1) = [p(t) - \gamma(t)d(t)]_{\mathcal{P}}, d(t) \in \mathcal{D} \in \mathcal{S}^{K}$

Can the results of unconstrained case extend?

Θ-cover does not work for constrained case



Get stuck at non-optimal points

Not necessarily a descent direction

Extension to Constrained Case

Problem:
$$\min_{p \in \mathbb{R}^n_+} f(p)$$

QGD: $p(t+1) = [p(t) - \gamma(t) \operatorname{sign}(\nabla f(x))]_+$

- This is a θ -cover with $\cos(\theta) = \frac{1}{\sqrt{K}}$.
- 1 bit for one dimension (constraint). N bits in total.
- For this special quantization, similar convergence results hold.
- Applications: TCP flow control, Optimal network flow, Power management, Voltage control in distribution networks, etc

Communication Complexity





Communication Complexity



Question:

What *minimal* bits (in total) are needed to achieve *E*-optimal solution? *E*-complexity (a min max definition)

What *optimal* accuracy is able to be achieved using *b*-bits (in total)? *b*-complexity (a min max definition)

Is there a simple coding scheme that reaches the complexity? Yes.

"Communication Complexity of Distributed Resource Allocation Optimization", Magnusson, Enyioha, Li, Fischione, Tarokh, submitted, 2017

Summary: Limited Communication

Extract information from physical measurements (Feedback)

- Load frequency control
- Power allocation in buildings/data center

Recover information from local computation

• Quantized dual (normalized) gradient for power allocation

Question:

How to choose the right algorithms and integrate them together?

Tradeoff: Efficiency, Robustness, Communication, Sensing,

Computation, Convergence speed

Thank you!

Accelerated Distributed Nesterov Gradient Descent

Guannan Qu, Na Li, John A. Paulson School of Engineering and Applied Sciences, Harvard University





<i>u</i> -strongly convex and <i>L</i> -smooth cost functions		
$\begin{array}{c} \hline \textbf{Proposed Algorithm} \\ \hline \textbf{Initialize}_{i}^{T_{i}}(0) = y_{i}(0) = v_{i}(0) \text{ arbitrary } s_{i}(0) = \nabla f_{i}(y_{i}(0)) \\ \hline \textbf{Step size } \eta > 0, \text{ and set } \alpha = \sqrt{\mu \eta} \end{array}$		
$x_{i}(t+1) = \sum_{j} w_{ij}y_{j}(t) - \eta s_{i}(t)$ $v_{i}(t+1) = (1-\alpha)\sum_{j} w_{ij}v_{j}(t) + \alpha \sum_{j} w_{ij}y_{j}(t) - \frac{\eta}{\alpha}s_{i}(t)$		
$y_i(t+1) = \frac{x_i(t+1)^j + \alpha v_i(t+1)}{1+\alpha} \int_{i}^{j} v_{ij}(t) dt \\ s_i(t+1) = \sum_i w_{ij} s_j(t) + \nabla f_i(y_i(t+1)) - \nabla f_i(y_i(t))$		
Summary of Results Theorem. Assume f_i is <i>L</i> -smooth and μ -strongly convex. When		
$0 < \eta < \frac{\sigma^3 (1 - \sigma)^3}{61909L} (\frac{\mu}{L})^{3/7}$		
where σ is the second largest singular value of the W , we have • $f(\bar{x}(t)) - f^* = O((1 - \sqrt{\mu\eta})^t)$		
• $ y(t) - 1x^* = O((1 - \sqrt{\mu\eta})^{t/2})$ Convergence Rate: $\left[1 - \frac{1}{\sqrt{61000}} (\frac{\mu}{L})^{5/7} \sigma^{3/2} (1 - \sigma)^{3/2} \right]^t$		
Dependence on the condition number $(\frac{\mu}{L})^{5/7}$ strictly better than that of GD $(\frac{\mu}{L})$! Nesterov does bring acceleration!		
Acc-DNGD-SC 10 ⁻⁰ 10 ⁻⁰ 1		
U DUU IUUU IDUU 10 ⁰ 10 ¹ 10 ² 10 ³ 10 ⁴		

For more detailed results, see Guannan Quand Na Li, "Accelerated Distributed Nesterov Gradient Descent," *arXiv preprint arXiv:1705.07176*(2017).

$$\label{eq:started_st$$

Summary of Results

(Throughout this section, f_i is convex and *L*-smooth.)

Theorem. When using vanishing step size $\eta_t = \frac{\eta}{(t+t_0)^{0.6+\epsilon}}$ with small η and large t_0 , and if the set of minimizers are compact, then $f(\bar{x}(t)) - f^* = O(\frac{1}{t^{1.4-\epsilon}})$. **Remark.** In simulation, $\eta = \frac{1}{2L}$ and $t_0 = 1$ suffice. **Theorem.** When using small enough fixed step size

 $\eta_t = \eta$, and assuming $f_i(x) = h_i(A_ix)$, for some matrix A_i and some h_i that is strongly convex and smooth, we have $f(\bar{x}(t)) - f^* = O(\frac{1}{t^2})$.

Conjecture. When using small enough fixed step size $\eta_t = \eta$, we have $f(\bar{x}(t)) - f^* = O(\frac{1}{t^2})$.

Nesterov does bring acceleration, from $O(\frac{1}{t})$ to $O(\frac{1}{t^{1.4-\epsilon}})$, or even $O(\frac{1}{t^2})$!

Preliminaries		
• f_i is L-smooth:	$\langle \nabla f_i(x) - \nabla f_i(y), x - y \rangle \le L x - y ^2$	
• f_i is μ -strongly convex:	$\langle \nabla f_i(x) - \nabla f_i(y), x - y \rangle \ge \mu x - y ^2$	
• $W = [w_{ij}]$ is a consensus averaging matrix:		
 Graph is connected 	• $w_{ii} > 0$, $w_{ij} > 0$ if i , j adjacent.	
• W is doubly stochast	ic • $w_{ij} = 0$ if i, j NOT adjacent.	

Back up: Reverse and Forward Engineering

Frequency control



Frequency control

- ➢ Goal: Balance the grid in an optimal (cost-effective) way
- Traditionally this is done at the generation side.



Frequency response



Advantages of load-side control

- faster (no/low inertia!)
- no waste or emission
- more resources (large #)
- Iocalize disturbance

Idea dates back to 1970s (Schweppe et al (1979, 1980))

Hierarchical Control at Different Time-scales



Challenges



Distributed Economically-Efficient Control



- Rebalance power
- Stabilize frequency
- Restore nominal frequency

Re-dispatch power optimally (min cost/disutility)

Distributed Economically-Efficient Control



Advantages:

For the *control*: Stable and more economically-efficient For the *optimization*: Save sensing/communication/computation

Problem setup

System Dynamics & Existing Control

$$\dot{x}_{i} = \sum_{j \in N(i)} A_{ij} x_{j} + B_{i} u_{i} + C_{i} w_{i}, \ i = 1, \cdots, N$$
$$\dot{u}_{i} = \sum_{j \in N(i)} D_{ij} x_{j} + \sum_{j \in N(i)} E_{ij} u_{j} + F_{i} w_{i}, \ i = 1, \cdots, N$$

Example:

Frequency dynamics, Voltage dynamics Primary/Secondary frequency/Voltage control Inverter dynamics/control (Model limitation: linear approximation)

Problem setup

Optimization Problem

$$\min_{x \in \mathbb{R}^{n}, u \in \mathbb{R}^{m}} \sum_{i=1}^{N} f_{i}(x_{i}) + \sum_{i=1}^{N} g_{i}(u_{i})$$

s.t.
$$\sum_{j \in N(i)} A_{ij}x_{j} + B_{i}u_{i} + C_{i}w_{i} = 0$$
$$\sum_{j \in N(i)} D_{ij}x_{j} + \sum_{j \in N(i)} E_{ij}u_{j} + F_{i}w_{i} = 0$$
$$h_{i}(x_{i}, u_{i}) \leq 0, \ i = 1, \dots, N$$

Example: Economic Dispatch, Optimal Load Response.

Problem setup



How to (re)design the control *u* to reach the optimal solution?

- Distributed
- Closed-loop (state-feedback)

Tool: reverse/forward engineering

Reverse



Forward

Equivalent

Economically Efficient State System Dynamics & Existing Control $\min_{x \in \mathbb{R}^n, u \in \mathbb{R}^m} \sum_{i=1}^N f_i(x_i) + \sum_{i=1}^N g_i(u_i)$ $\dot{x}_{i} = \sum A_{ij}x_{j} + B_{i}u_{i} + C_{i}w_{i}, i = 1, \dots, N$ $i \in N(i)$ s.t. $\sum A_{ii}x_{i} + B_{i}u_{i} + C_{i}w_{i} = 0$ $\dot{u}_{i} = \sum D_{ij}x_{j} + \sum E_{ij}u_{j} + F_{i}w_{i}, i = 1, \dots, N$ $i \in N(i)$ $\sum D_{ij}x_j + \sum E_{ij}u_j + F_iw_i = 0$ $h_i(x_i, u_i) \leq 0, i = 1, \cdots, N$ **Modified**

Optimization Problem

Forward

System Dynamics & Modified Control

$$x_{i}^{k} = \sum_{j \in N(i)} A_{ij}x_{j} + B_{i}u_{i} + C_{i}w_{i},$$

$$u_{i}^{k} = \sum_{j \in N(i)} D_{ij}x_{j} + \sum_{j \in N(i)} E_{ij}u_{j} + F_{i}w_{i} + g(z_{i}),$$

$$x_{i}^{k} = f\left(x_{j}, u_{j}, z_{j} : j \in N_{i}\right)$$
Solve
$$\sum_{j \in N(i)} D_{ij}x_{j} + \sum_{j \in N(i)} E_{ij}u_{j} + F_{i}w_{i} = 0$$

$$\sum_{j \in N(i)} D_{ij}x_{j} + \sum_{j \in N(i)} E_{ij}u_{j} + F_{i}w_{i} = 0$$

$$h_{i}(x_{i}, u_{i}) \leq 0, i = 1, \cdots, N$$
Equivalent
Optimization Problem

Distributed Economically-Efficient Control

System Dynamics & Existing Control

$$\dot{x}_{i} = \sum_{j \in N(i)} A_{ij}x_{j} + B_{i}u_{i} + C_{i}w_{i}, i = 1, \dots, N$$

$$\dot{u}_{i} = \sum_{j \in N(i)} D_{ij}x_{j} + \sum_{j \in N(i)} E_{ij}u_{j} + F_{i}w_{i}, i = 1, \dots, N$$
Economical Efficient State

$$\min_{x \in \mathbb{R}^{n}, u \in \mathbb{R}^{m}} \sum_{i=1}^{N} f_{i}(x_{i}) + \sum_{i=1}^{N} g_{i}(u_{i})$$
s.t.
$$\sum_{j \in N(i)} A_{ij}x_{j} + B_{i}u_{i} + C_{i}w_{i} = 0$$

$$\sum_{j \in N(i)} D_{ij}x_{j} + \sum_{j \in N(i)} E_{ij}u_{j} + F_{i}w_{i} = 0$$

$$h_{i}(x_{i}, u_{i}) \leq 0, i = 1, \dots, N$$

Sufficient and necessary conditions are available at [Zhang, Antonois, Li, 2015, 2016]

Application to optimal load control for primary freq. control

Simulation of IEEE a 68-bus system



Application to automatic generation control

Simulation of a 4-bus system



Distributed Economically-Efficient Control



Advantages:

For the *control*: Stable and more economically-efficient For the *optimization*:

A large amount of sensing, comm. and comp. is saved

Thank you!

Load frequency control

System Dynamics

$$\dot{\omega}_{i} = -\frac{1}{M_{i}} \left(\sum_{l \in L(i)} d_{l}(t) + D_{i}\omega_{i}(t) - P_{i}^{m} + \sum_{j:i \to j} P_{ij}(t) - \sum_{k:k \to i} P_{ji}(t) \right)$$
$$\dot{P}_{ij} = b_{ij} \left(\omega_{i}(t) - \omega_{j}(t) \right)$$

 Λ

$$d_l(t) = \left[C_l^{-1}(\omega_i(t))\right]_{\underline{d}_l}^{\overline{d}_l} \text{ for } l \in L(i)$$

Frequency: a locally measurable signal ("price" of imbalance)

Completely decentralized; no explicit communication necessary
<u>Step 1</u>: Each appliance *i* updates the power request $q_i(t)$ & sends to the control center $q_i(t) = \arg \max_{q_i} U_i(q_i) - p(t)q_i$

<u>Step 2</u>: Control center updates the signal p(t) & sends to each appliance

$$p(t+1) = [p(t) + \gamma (\sum_{i} q_i(t) - Q)]^+$$

Replace this with *true measurement* of total power consump. Q(t).

Normalized Gradient

Problem: $\min_{p} f(p)$ Gradient Descent: $p(t+1) = p(t) - \gamma(t)f'(p(t))$ Normalized Gradient Descent: $p(t+1) = p(t) - \gamma(t)\frac{f'(p(t))}{||f'(p(t))||}$

Questions: Reduce the communication?

If $p \in \mathbb{R}$, normalized gradient reduces to binary signal, ± 1 .

If $p \in \mathbb{R}^{K}$, K > 1, normalized gradient is in \mathcal{S}^{K-1} (unit spher Primal Problem $\max_{q_i} \sum_{i} U_i(q_i)$ s.t. $Aq \leq B$ e.g. Network constraints; Multi-resource allocation

Quantized (Normalized) Gradient Descent (QGD)

Problem:
$$\min_{p} f(p)$$

QGD: $p(t+1) = p(t) - \gamma(t)d(t), d(t) \in \mathcal{D} \in \mathcal{S}^{K-1}$

Question: When \mathcal{D} is a good quantization?

 \mathcal{F} : set of convex and differential functions whose gradients are L-Lipschitz continuous and \mathcal{P}^* is nonempty and bounded.

 \mathcal{P}^* : set of optimizers

Definition [Proper Quantization]

A set \mathcal{D} is a proper quantization if for every $f \in \mathcal{F}$ and every initiation $p(0) \in \mathbb{R}^K$, there are $d(t) \in \mathcal{D}$ and $\epsilon(t) \in \mathbb{R}_+$, where $t \in \mathbb{N}$ such that

 $\lim_{t \to \infty} \mathsf{dist}(p(t), \mathcal{P}^*) = 0.$

A quantization is **proper if and only if** the algorithm is able to **converge** to the optimal points for **any** well-behaviored f.

Quantized (Normalized) Gradient Descent (QGD)

Problem:
$$\min_{p} f(p)$$

QGD: $p(t+1) = p(t) - \gamma(t)d(t), d(t) \in \mathcal{D} \in \mathcal{S}^{K-1}$

Questions:

- A) How to determine a quantization is proper?
- B) What is the minimal size of the quantization to be proper?
- C) How to choose d(t) and $\varepsilon(t)$, Given a proper quantization?
- *D) What* are the *connections* between the *fineness* of the quantization to the *convergence* of the algorithm?

Descent direction



Proper quantization

Proposition. A set \mathcal{D} is a proper quantization if and only if there exists $\theta > 0$, for every $g \in \mathcal{S}^{K-1}$ there exists $d \in \mathcal{D}$ such that $\cos(\operatorname{ang}(g, d)) \ge \theta$.

Such quantization \mathcal{D} is also called as θ -cover.

Minimal proper quantization, $|\mathcal{D}| = N + 1$.

Only $\log_2(N+1)$ bits are needed at each iteration.

Convergence rate

Stopping Criterion: $\mathcal{P}(\epsilon) := \left\{ p \in \mathbb{R}^K : ||\nabla f(p)|| \le \epsilon \right\}$

Theorem: For any $\epsilon > 0$, if $\gamma \in (0, 2\theta \epsilon/L)$, then there exists $T \in \mathbb{N}$ such that $p(T) \in \mathcal{P}(\epsilon)$, with T bounded by

$$T \le \frac{2(f(p(0)) - f^*)}{\gamma(2\theta\epsilon - L\gamma)}.$$

- Finer quantization, larger stepsize is allowed
- Finer quantization, faster convergence

Convergence rate

Stopping Criterion: $F(\epsilon) := \sup \{f(p) : p \in \mathcal{P}(\epsilon)\}$

Theorem:

1) For any $\epsilon > 0$, if $\gamma \in (0, 2\theta \epsilon/L)$, and $T \in \mathbb{N}$ such that $p(T) \in \mathcal{P}(\epsilon)$, then for all $t \ge T$ $f(p(t)) \le F(\epsilon) + \left(\epsilon + \frac{L}{2}\gamma\right)\gamma.$

2) If the stepsize $\gamma(t)$ is non-summable and square summable, then $\lim_{t\to\infty} \operatorname{dist}(p(t), \mathcal{P}^*) = 0$.

Quantization size



Quantization size



Simulation

f is the dual function of a two-dimension resource allocation problem.



(3) Infinite bandwidth: normalized gradient(2) Infinite bandwidth: gradient

Summary of QGD

Problem:
$$\min_{p} f(p)$$

QGD: $p(t+1) = p(t) - \gamma(t)d(t), d(t) \in \mathcal{D} \in \mathcal{S}^{K}$

- A) Proper quantization = θ -cover
- B) Minimal size of proper quantization is N+1
- C) Pick the quantized direction closest to gradient direction
- D) θ plays an important role in determines the convergence

Proper quantization

Minimal proper quantization, $|\mathcal{D}| = N + 1$. Only $\log_2(N+1)$ bits are needed at each iteration.

Finer quantization, larger stepsize is allowed Finer quantization, faster convergence

Communication Complexity

