

Differentially Quantized Gradient Descent

Chung-Yi Lin Victoria Kostina Babak Hassibi

Department of Electrical Engineering
California Institute of Technology

IEEE International Symposium on Information Theory
May 28, 2021

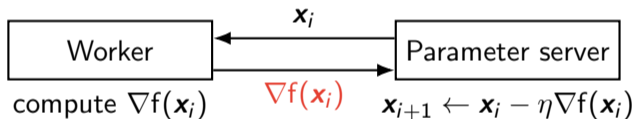
Unquantized gradient descent (GD)

$$\text{Solve } \mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

- with GD:

$$\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i - \eta \nabla f(\mathbf{x}_i)$$

- $\eta > 0$ is a constant stepsize.
- in distributed training



Exchanging **the gradients** carries high communication cost.

Previous works

Gradient quantization

- Stochastic gradient descent
 - Seide et al. 2014
 - Wen et al. 2017
 - Alistarh et al. 2017
 - Bernstein et al. 2018
 - Wu et al. 2018
 - Gandikota et al. 2019
 - Ramezani-Kebrya et al. 2019
 - Mayekar & Tyangi 2019, 2020
- GD
 - Luo & Tseng 1993
 - Friedlander & Schmidt 2012
 - Alistarh et al. 2016

Gradient sparsification

- Strom et al. 2015
- Aji & Heafield 2017
- Lin et al. 2018
- Wangni et al. 2018
- Wang et al. 2018
- Stich et al. 2018

Convergence lower bounds are **scarce**.

This work: the necessary and sufficient **bit rate** to achieve a target **convergence rate**.

Class of functions

$$\mathcal{F}_n(\mu, L, r) \triangleq \{f: \mathbb{R}^n \rightarrow \mathbb{R} \mid f \text{ satisfies the following.}\}$$

- f is L -smooth: $\|\nabla f(\mathbf{v}) - \nabla f(\mathbf{w})\| \leq L \|\mathbf{v} - \mathbf{w}\|$
- f is μ -strongly convex: $(\nabla f(\mathbf{v}) - \nabla f(\mathbf{w}))^\top (\mathbf{v} - \mathbf{w}) \geq \mu \|\mathbf{v} - \mathbf{w}\|^2$
- The minimizer $\|\mathbf{x}^*(f)\| \leq r$ for some $r > 0$

Unquantized gradient descent

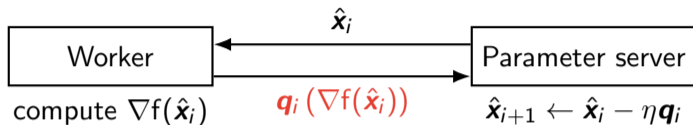
For any L -smooth and μ -strongly convex f , unquantized GD satisfies [Polyak, 1987]:

$$\|\mathbf{x}_T - \mathbf{x}^*(f)\| \leq \sigma^T \|\mathbf{x}_0 - \mathbf{x}^*(f)\|$$

- $\sigma \triangleq \frac{L-\mu}{L+\mu}$: contraction factor of GD.
- The bound is tight: $\exists f$ s.t. “=” holds.

Naively Quantized Gradient Descent (NQ-GD)

NQ-GD [Friedlander & Schmidt 2012, Alistarh et al. 2016] directly quantizes the gradient at $\hat{\mathbf{x}}_i$:



Theorem: NQ-GD

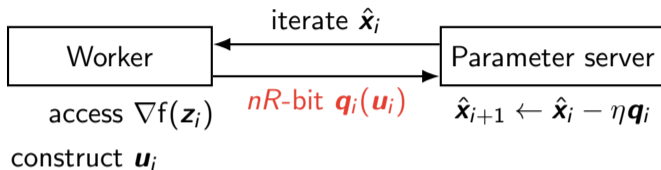
NQ-GD achieves the following contraction factor over \mathcal{F}_n

$$\sigma_{\text{NQ-GD}}(n, R) \leq \sigma + \rho_n 2^{-R}$$

ρ_n : **covering efficiency** of the quantizer

- uniform scalar quantizer: $\rho_n = \sqrt{n}$.
- $\rho_n \geq 1$.
- $\rho_n = 1 + o_n(1)$ is achievable with lattice quantizers [Rogers 1963].

Quantized gradient descent (QGD)



The worker, based on $\hat{\mathbf{x}}_i$ and $\mathbf{e}_0, \dots, \mathbf{e}_{i-1}$ ($\mathbf{e}_\ell \triangleq \mathbf{q}_\ell - \mathbf{u}_\ell$), decides:

- gradient query point \mathbf{z}_i
- quantizer's input \mathbf{u}_i .

Goals:

- characterize the tradeoff between how fast any QGD algorithm converges and R .
- propose an algorithm that achieves it.

Quantized Gradient Descent: worst-case contraction factor

For a QGD algorithm A operating at R bits per problem dimension, **worst-case** (over $f \in \mathcal{F}_n(\mu, L, r)$) contraction factor:

$$\sigma_A(n, R) \triangleq \sup_{f \in \mathcal{F}_n} \limsup_{T \rightarrow \infty} \|\hat{\mathbf{x}}_T(R) - \mathbf{x}^*(f)\|^{\frac{1}{T}}$$

- unquantized GD: $\sigma_{\text{GD}}(n, \infty) = \sigma = \frac{L-\mu}{L+\mu}$.

Theorem: converse

The contraction factor of any QGD algorithm A operating at R bits per problem dimension satisfies

$$\sigma_A(n, R) \geq \max \left\{ \sigma, 2^{-R} \right\}$$

Proof combines two converses:

- Reduction to unquantized GD: $\sigma_A(n, R) \geq \sigma$;
- Volume division argument: $\sigma_A(n, R) \geq 2^{-R}$.

Differentially Quantized Gradient Descent (DQ-GD)

Algorithm 1: DQ-GD

Initialize $\mathbf{e}_{-1} \leftarrow \mathbf{0}$

for $i = 0$ **to** $T - 1$ **do**

Worker:

$$\mathbf{z}_i \leftarrow \hat{\mathbf{x}}_i + \eta \mathbf{e}_{i-1}$$

$$\mathbf{u}_i \leftarrow \nabla f(\mathbf{z}_i) - \mathbf{e}_{i-1}$$

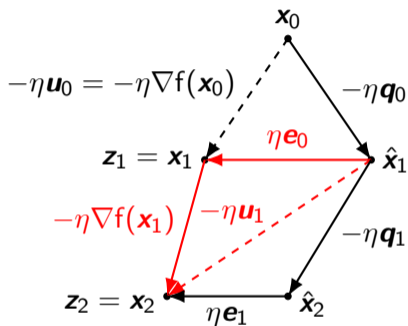
$$\mathbf{q}_i = \mathbf{q}_i(\mathbf{u}_i)$$

$$\mathbf{e}_i \leftarrow \mathbf{q}_i - \mathbf{u}_i$$

Parameter server:

$$\hat{\mathbf{x}}_{i+1} \leftarrow \hat{\mathbf{x}}_i - \eta \mathbf{q}_i$$

end

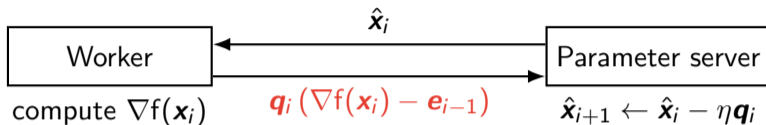


Differential quantization¹ directs the quantized trajectory to the unquantized trajectory.

¹The idea of error compensation dates back to $\Sigma\Delta$ modulation [Gray, 1989].

Differentially Quantized Gradient Descent (DQ-GD)

DQ-GD computes the gradient at \mathbf{x}_i and compensates the previous quantization error:



Theorem: DQ-GD

DQ-GD achieves the following contraction factor over \mathcal{F}_n

$$\sigma_{\text{DQ-GD}}(n, R) \leq \max \left\{ \sigma, \rho_n 2^{-R} \right\}$$

- Since $\rho_n \rightarrow 1$ is achievable [Rogers, 1963], DQ-GD attains the converse as $n \rightarrow \infty$
- $R \geq \log_2 \rho_n / \sigma$: achieves the contraction factor σ of unquantized GD
- $R < \log_2 \rho_n / \sigma$: achieved contraction factor is only $\rho_n 2^{-R}$

Proof sketch

$$\text{Induction: } \hat{\mathbf{x}}_i = \mathbf{x}_i - \eta \mathbf{e}_{i-1}$$

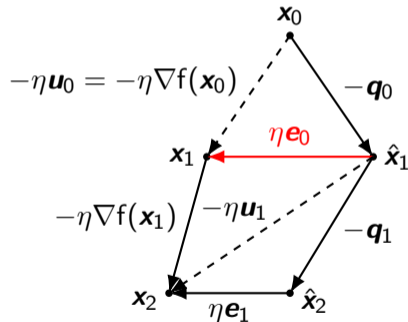
- $\|\hat{\mathbf{x}}_i - \mathbf{x}^*(f)\| \leq \|\mathbf{x}_i - \mathbf{x}^*(f)\| + \eta \|\mathbf{e}_{i-1}\|$
- 1st term: convergence of GD

$$\|\mathbf{x}_T - \mathbf{x}^*(f)\| \leq \sigma^T \|\mathbf{x}_0 - \mathbf{x}^*(f)\|.$$

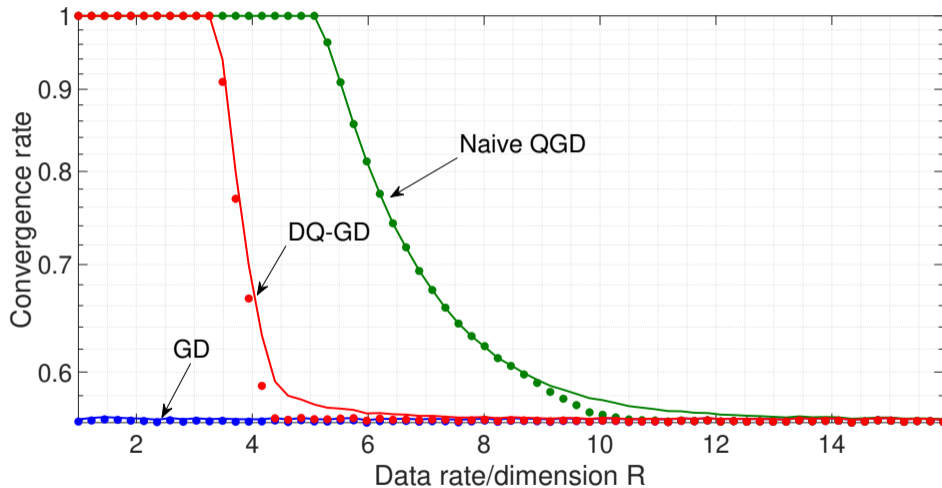
- 2nd term: choose the *dynamic range* r_i of the quantizer carefully so that

$$\sup_{\|\mathbf{u}_i\| \leq r_i} \|\mathbf{e}_i\| \leq \frac{\rho_n}{2R} r_i$$

for each iteration i .

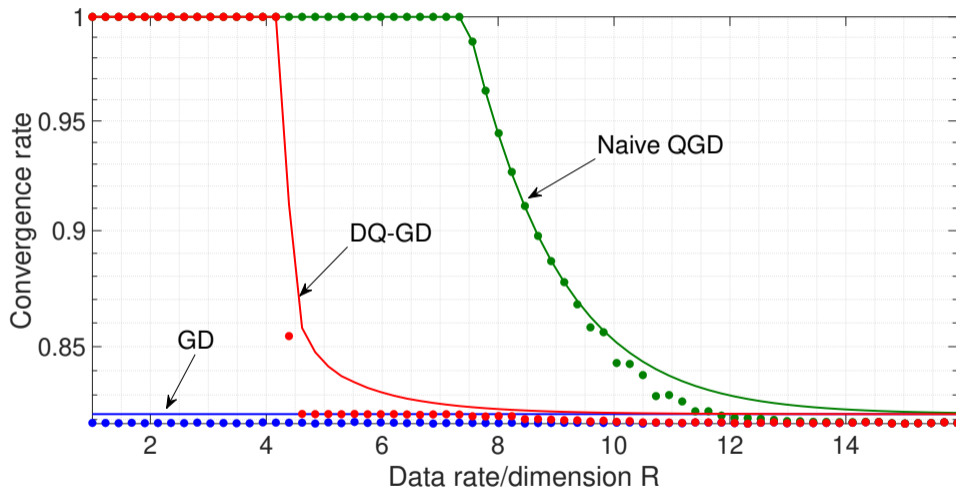


Least-squares problems: Gaussian ensemble



$f(x) = \|y - Ax\|^2 / 2$. $A \in \mathbb{R}^{1000 \times 100}$ and $y \in \mathbb{R}^{1000}$ iid standard normal entries. $\kappa(A) \approx 1.8862$ on average.

Least-squares problems: Real-world matrix



$\mathbf{A} \in \mathbb{R}^{958 \times 292}$ with $\kappa(\mathbf{A}) \approx 3.2014$ (SuiteSparse matrix collection) and $\mathbf{y} \in \mathbb{R}^{958}$ iid standard normal entries.

Recap: main result

Optimal contraction factor over $f \in \mathcal{F}_n$ and $A \in \text{QGD}$

$$\lim_{n \rightarrow \infty} \inf_{\text{QGD } A} \sigma_A(n, R) = \max \left\{ \sigma, 2^{-R} \right\}.$$

Phase transition

- $R \geq \log_2 1/\sigma$: contraction factor σ of unquantized GD is achievable
- $R < \log_2 1/\sigma$: only 2^{-R} is achievable

Extension: gradient methods with momentum

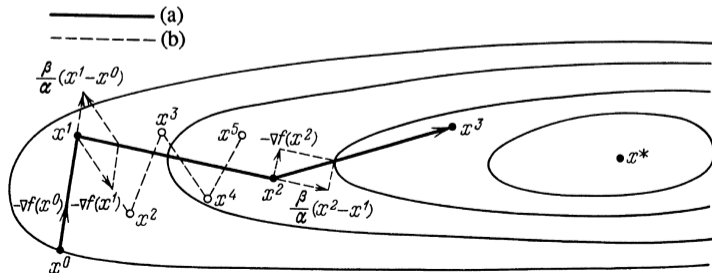
- Accelerated Gradient Descent [Nesterov, 1982]:

$$\mathbf{y}_{i+1} \leftarrow \mathbf{x}_i - \eta \nabla f(\mathbf{x}_i)$$

$$\mathbf{x}_{i+1} \leftarrow \mathbf{y}_{i+1} + \gamma (\mathbf{y}_{i+1} - \mathbf{y}_i)$$

- Heavy Ball Method [Polyak, 1987]:

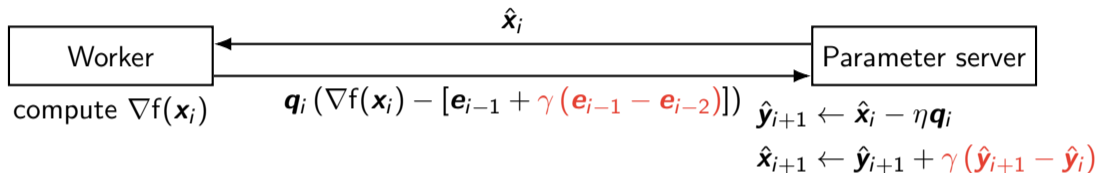
$$\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i - \eta \nabla f(\mathbf{x}_i) + \gamma (\mathbf{x}_i - \mathbf{x}_{i-1})$$



(a) Heavy Ball Method
(b) Gradient Descent

Differentially Quantized Accelerated Gradient Descent (DQ-AGD)

DQ-AGD computes the gradient at \mathbf{x}_i and compensates two past quantization errors:



Theorem: DQ-AGD

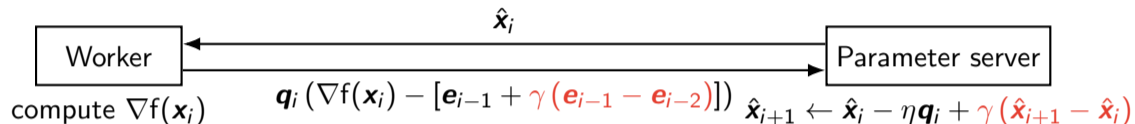
DQ-AGD achieves the following contraction factor over \mathcal{F}_n

$$\sigma_{\text{DQ-AGD}}(n, R) \leq \max \{ \sigma_{\text{AGD}}, \phi_{\text{DQ-AGD}}(n, R) \}$$

- σ_{AGD} : contraction factor of unquantized AGD
- $\phi_{\text{DQ-AGD}}(n, R)$: exponentially decreasing function of R

Differentially Quantized Heavy Ball Method (DQ-HB)

DQ-AGD computes the gradient at \mathbf{x}_i and compensates two past quantization errors:



Theorem: DQ-HB

DQ-HB achieves the following contraction factor over $f \in \mathcal{F}_n$ that are twice continuously differentiable:

$$\sigma_{\text{DQ-HB}}(n, R) \leq \max \{ \sigma_{\text{HB}}, \phi_{\text{DQ-HB}}(n, R) \}$$

- σ_{HB} : contraction factor of unquantized HB
- $\phi_{\text{DQ-HB}}(n, R)$: exponentially decreasing function of R

Conclusion

- Introduced **Differential Quantization**.
- Differential Quantization is **substantially better** than naive quantization.
- If $R \geq R_A$, differentially quantized $A \in \{\text{GD}, \text{AGD}, \text{HB}\}$ attains the contraction factor of the unquantized A .
- In the limit of $n \rightarrow \infty$, DQ-GD attains the **optimal contraction factor** within the class of QGD algorithms.
- Multiworker?