

Gradient Compression For Communication Limited Convex Optimization

Siddarth kumar

IIT Hyderabad

EE15BTECH11032

March 7, 2019

- Optimization using gradient descent.
 - Iteration complexity.
 - Communication cost of exchanging gradients.
- Compression technique for gradient is used to reduce compression.
- Discuss Gradient compression technique.

Problem formulation

Convex optimization problem,

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} f(x). \quad (1)$$

Function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ should satisfy below to conditions:

L-Lipschitz continuous gradient Condition,

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2. \quad (2)$$

Function is strongly convex,

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|y - x\|^2. \quad (3)$$

Proposed method,

$$x_{k+1} = x_k - \gamma_k Q(\nabla f(x_k)).$$

Q is quantization, γ is step size.

Instead of using traditional gradient descent,

$$x_{k+1} = x_k - \gamma \nabla f(x_k), \tag{4}$$

- Sparsification

- K-greedy Quantizer $Q: \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$[Q_G^K(g)]_i = \begin{cases} [g]_{\pi(i)} & \text{if } i \leq K \\ 0 & \text{otherwise} \end{cases}$$

where π is a permutation of $\{1, \dots, n\}$, and $|g_{\pi(k)}| \geq |g_{\pi(k+1)}|$
 g is gradient vector.

- Quantization

- Ternary Quantizer $Q_T: \mathbb{R}^n \rightarrow \mathbb{R}^n$ such that
 $[Q_T(g)]_i = \|g\| \text{sgn}(g_i)$

- Combination of Sparsification and Quantization (Dynamic gradient quantizer)
 - It is defined as $Q_D : \mathbb{R}^n \rightarrow \mathbb{R}^n$ as

$$[Q_D(g)]_i = \begin{cases} \|g\| \operatorname{sgn}(g_i) & \text{if } i \in I(g) \\ 0 & \text{otherwise} \end{cases}$$

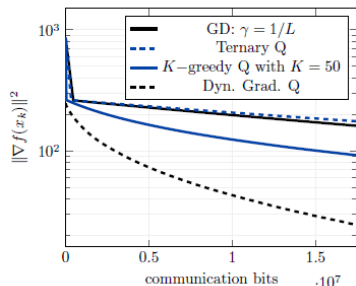
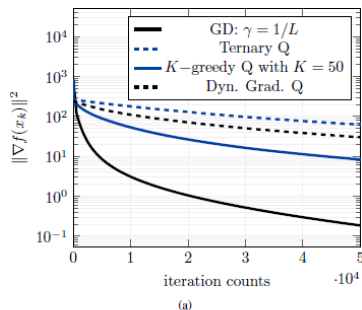
where $I(g)$ is the smallest subset of $\{1, \dots, n\}$ such that

$$\sum_{j \in I(g)} |g_j| \geq \|g\|.$$

- Number of bits required in different Quantizer
 - K-greedy Quantizer: $K(\log_2(n) + b)$ bits
 - Ternary Quantizer: $(2n + b)$ bits
 - Dynamic Quantizer: $|I(g)|(\log_2(n) + 1) + b$

Experimental results

- Let Function $f(x) = 0.5 \times ||Ax - b||$, where $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$, $m = 1000$ and $n = 800$, elements of matrix drawn from uniform distribution $[0,1]$,
Elements of b is sign of random variable drawn from $\mathcal{N}(0,1)$.
- Normalized each row of A by its Euclidean norm and computed the Lipschitz constant as $L = \lambda_{\max}(A^T A)$.



References

- <https://ieeexplore.ieee.org/document/8619625>
- https://en.wikipedia.org/wiki/Gradient_descent