# Optimization for Data Science

## Lec 07: Metric Gradient

Yaoliang Yu

UNIVERSITY OF WATERLOO | FACULTY OF MATHEMATICS
DAVID R. CHERITON SCHOOL
OF COMPUTER SCIENCE

Unconstrained minimization:

$$f_\star = \inf_{\mathbf{w} \in \mathbb{R}^d} f(w_1, \ldots, w_d)$$

- $f$: smooth w.r.t. a general norm $\|\cdot\|$ and possibly nonconvex

- For simplicity, no constraints on $\mathbf{w}$

# Gradient Compression

Typical problem in ML:

$$f(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^{m} f_i(\mathbf{w}; \mathcal{D}_i)$$

- Each $f_i$ represent a different user/study/processor

$$f'(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^{m} f_i'(\mathbf{w}; \mathcal{D}_i)$$

- For large $d$, communicating and aggregating the individual gradients are expensive
- Compress the gradients by simply taking its sign?

J. Bernstein, Y.-X. Wang, K. Azizzadenesheli, and A. Anandkumar. "signSGD: Compressed Optimisation for Non-Convex Problems". In: *Proceedings of the 35th International Conference on Machine Learning*. 2018, pp. 560–569, J. Bernstein, J. Zhao, K. Azizzadenesheli, and A. Anandkumar. "signSGD with Majority Vote is Communication Efficient and Fault Tolerant". In: *International Conference on Learning Representations*. 2019.

## Definition: Norm

Recall a norm $\| \cdot \|$ satisfies:

- definiteness: $\|\mathbf{x}\| \geq 0$ with 0 attained iff $\mathbf{x} = \mathbf{0}$
- positive homogeneity: $\|\lambda \mathbf{x}\| = |\lambda| \cdot \|\mathbf{x}\|$ for any $\lambda \in \mathbb{R}$
- triangle inequality: $\|\mathbf{x} + \mathbf{z}\| \leq \|\mathbf{x}\| + \|\mathbf{z}\|$

## Definition: Dual

The dual norm of a norm $\| \cdot \|$ is

$$\|\mathbf{w}^*\|_\circ := \max_{\|\mathbf{w}\| \leq 1} \langle \mathbf{w}; \mathbf{w}^* \rangle$$

## Example:

The dual of the $\ell_p$ norm $\|\mathbf{w}\|_p := (\sum_j |w_j|^p)^{1/p}$ is $\ell_q$ norm, where $1/p + 1/q = 1$.

## Definition: duality mapping

Let $q := \frac{1}{2}\|\cdot\|^2$ be "quadratic." We define the duality mapping

$$J = \partial q : V \to V^*, \quad j : V \to V^*, \quad \mathbf{w} \mapsto j(\mathbf{w}) \in J(\mathbf{w}),$$

where $j$ is an arbitrary single-valued selection of $J$.

$$\langle \mathbf{w}; j(\mathbf{w}) \rangle = \|\mathbf{w}\|^2 = \|j(\mathbf{w})\|_\circ^2$$

## Definition: metric gradient w.r.t. a norm

We define the metric gradient w.r.t. a norm $\|\cdot\|$ as

$$\blacktriangledown f = J^{-1}(f'), \quad \triangledown f = j^{-1}(f') : V \to V.$$

M. Golomb and R. A. Tapia. "The metric gradient in normed linear spaces". *Numerische Mathematik*, vol. 20 (1972), pp. 115–124.

# Steepest Descent

Another way to recognize the metric gradient is through Kantorovich's steepest descent. Fixing the current iterate $\mathbf{w}_t$, we look for a direction $\mathbf{d}$ such that the univariate function

$$\eta \mapsto h(\eta) := f(\mathbf{w}_t - \eta \mathbf{d})$$

decreases steepest.

Kantorovich proposed to find the direction $\mathbf{d}$ through the subproblem:

$$\operatorname*{argmin}_{\mathbf{d} \neq \mathbf{0}} \quad \frac{h'(\eta)|_{\eta=0}}{\|\mathbf{d}\|} = \frac{-\langle \mathbf{d}; f'(\mathbf{w}_t) \rangle}{\|\mathbf{d}\|} \implies \mathbf{d} = \frac{\overline{\nabla} f(\mathbf{w}_t)}{\|\overline{\nabla} f(\mathbf{w}_t)\|} = \frac{\overline{\nabla} f(\mathbf{w}_t)}{\|\nabla f(\mathbf{w}_t)\|_{\circ}},$$

which is exactly the normalized metric gradient!

L. V. Kantorovich. "On an effective method of solving extremal problems for quadratic functionals". *Soviet Mathematics Doklady*, vol. 48, no. 7 (1945), pp. 595–600.

**Algorithm 1:** Metric gradient descent for unconstrained smooth minimization

**Input:** $\mathbf{w}_0$, norm $\|\cdot\|$

1 **for** $t = 0, 1, \ldots$ **do**
2      $\mathbf{g}_t \leftarrow \overline{\nabla} f(\mathbf{w}_t)$          // compute any metric gradient
3      **if** $\|\mathbf{g}_t\| = 0$ **then**
4          **break**
5      choose step size $\eta_t > 0$
6      $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t \mathbf{g}_t$          // update

Key insight (note the similarity as before):

$$f(\mathbf{w}) \leq f(\mathbf{w}_t) + \langle \mathbf{w} - \mathbf{w}_t; f'(\mathbf{w}_t) \rangle + \tfrac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|^2,$$

i.e. L-smoothness w.r.t. a general norm $\|\cdot\|$.

Apply polar decomposition on the RHS:

$$\min_{\lambda \geq 0} \min_{\|\mathbf{w} - \mathbf{w}_t\| = \lambda} f(\mathbf{w}_t) + \langle \mathbf{w} - \mathbf{w}_t ; f'(\mathbf{w}_t) \rangle + \frac{1}{2\eta_t}\lambda^2 \quad \equiv \quad \min_{\lambda \geq 0} -\lambda \|f'(\mathbf{w}_t)\|_\circ + \frac{1}{2\eta_t}\lambda^2.$$

Thus, $\lambda = \eta_t \|f'(\mathbf{w}_t)\|_\circ$ and

$$\mathbf{w} - \mathbf{w}_t = \lambda \frac{-\nabla f(\mathbf{w}_t)}{\|f'(\mathbf{w}_t)\|_\circ}, \qquad i.e. \qquad \mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t)$$

## Theorem: convergence of metric gradient descent for L-smooth functions

Let $f : \mathbb{R}^d \to \mathbb{R}$ be L-smooth w.r.t. a general norm $\|\cdot\|$ and bounded from below (i.e. $f_\star > -\infty$). If the step size $\eta_t \in [\alpha, \frac{2}{L} - \beta]$ for some $\alpha, \beta > 0$, then the sequence $\{\mathbf{w}_t\}$ generated satisfies $\overline{\nabla} f(\mathbf{w}_t) \to \mathbf{0}$. Moreover,

$$\min_{0 \leq t \leq T-1} \|\overline{\nabla} f(\mathbf{w}_t)\| \leq \sqrt{\frac{f(\mathbf{w}_0) - f_\star}{\alpha \beta L T / 2}}.$$

- The proof is literally the same as that of gradient descent
- Choosing $\alpha = \beta = \frac{1}{L}$, the bound reduces to

$$\min_{0 \leq t \leq T-1} \|\overline{\nabla} f(\mathbf{w}_t)\| \leq \sqrt{\frac{2L[f(\mathbf{w}_0) - f_\star]}{T}}.$$

- Obviously, LHS depends on the norm and so does RHS (through $L = L_{\|\cdot\|}$)

# $\ell_p$ norm metric gradient

Let $\mathsf{V} = \mathbb{R}^d$ be equipped with the $\ell_p$ norm, whose dual is $\ell_q$ norm with $1/p + 1/q = 1$.

$$\blacktriangledown f(\mathbf{w}) := \left[ \underset{\|\mathbf{z}\|_p \le \|f'(\mathbf{w})\|_q}{\operatorname{argmax}} \langle \mathbf{z}; f'(\mathbf{w}) \rangle \right] = \|f'(\mathbf{w})\|_q^{1-q/p} \cdot \operatorname{sign}(f'(\mathbf{w})) \cdot |f'(\mathbf{w})|^{q/p}$$

- When $p = q = 2$, we have $\blacktriangledown f = \nabla f$

- When $p = 1, q = \infty$, we have $\blacktriangledown f = \operatorname{conv}\{\nabla_j f \cdot \mathbf{e}_j : |\nabla_j f| = \|\nabla f\|_\infty\}$

- When $p = \infty, q = 1$, we have $\blacktriangledown f = \operatorname{conv}\{\|\nabla f\|_1 \cdot \operatorname{sign}(\nabla f)\}$, $\operatorname{sign}(0) \in [-1, 1]$

$$\boxed{\text{metric gradient indeed depends on the norm}}$$

# Sign gradient descent

Let us equip the input space $V$ (where $\mathbf{w}$ lives) with the $\ell_\infty$ norm , and the gradient space $V^*$ (where $f'(\mathbf{w})$ lives) with the corresponding dual $\ell_1$ norm.

We obtain the so-called sign gradient descent algorithm, where in each iteration we only update with the sign of the gradient:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \|\nabla f(\mathbf{w}_t)\|_1 \cdot \operatorname{sign}(\nabla f(\mathbf{w}_t)),$$

which is particularly appealing in distributed and low-resource devices.

# Coordinate gradient descent

Let us equip the input space $V$ (where $\mathbf{w}$ lives) with the $\ell_1$ norm , and the gradient space $V^*$ (where $f'(\mathbf{w})$ lives) with the corresponding dual $\ell_\infty$ norm.

We obtain the so-called greedy coordinate gradient descent algorithm, where in each iteration we only take a gradient step along one (block of) coordinate(s):

$$w_{j,t+1} = w_{j,t} - \eta_t \nabla_j f(\mathbf{w}_t), \quad \text{where} \quad |\nabla_j f(\mathbf{w}_t)| = \|\nabla f(\mathbf{w}_t)\|_\infty.$$

- Compute all derivatives to figure out which one is largest

- Most of the computational effort is wasted...

R. V. Southwell. "Stress-Calculation in Frameworks by the Method of "Systematic Relaxation of Constraints". I and II". *Proceedings of the Royal Society of London. Series A, Mathematical and Physical Sciences*, vol. 151, no. 872 (1935), pp. 56–95.

# Alternatives

An obvious alternative is to update the coordinates cyclically:

$$\textbf{for} \quad j = 1, \ldots, d$$
$$w_j \leftarrow w_j - \eta \nabla_j f(\mathbf{w})$$

- computing the gradient $\nabla f$ vs. computing a single component $\nabla_j f$?

- L-smoothness is w.r.t. different norms!

- Can randomize our choice of the coordinates [0]

- Might as well go to the extreme:

$$w_j \leftarrow \underset{w}{\operatorname{argmin}} \ f(w_1, \ldots, w_{j-1}, w, w_{j+1}, \ldots, w_d)$$

Y. Nesterov. "Efficiency of Coordinate Descent Methods on Huge-Scale Optimization Problems". *SIAM Journal on Optimization*, vol. 22, no. 2 (2012), pp. 341–362.

## Definition: metric projection

We define the metric projection w.r.t. an arbitrary norm and a closed set $C$:

$$\mathrm{P}_C(\mathbf{w}) = \operatorname*{argmin}_{\mathbf{z} \in C} \|\mathbf{w} - \mathbf{z}\|.$$

However, the metric projection may no longer be nonexpansive even when $C$ is convex.

---

**Algorithm 2:** Metric projected gradient descent

**Input:** $\mathbf{w}_0 \in C$, norm $\|\cdot\|$

1 **for** $t = 0, 1, \ldots$ **do**

2     $\mathbf{g}_t \leftarrow \overline{\nabla} f(\mathbf{w}_t)$        // compute any metric gradient

3     $\eta_t \leftarrow \operatorname{argmin}_{\eta \geq 0} f(\mathrm{P}_C(\mathbf{w}_t - \eta \mathbf{g}_t))$        // Cauchy's rule

4     $\mathbf{w}_{t+1} \leftarrow \mathrm{P}_C(\mathbf{w}_t - \eta_t \mathbf{g}_t)$        // update

---

G. P. McCormick. "Anti-Zig-Zagging by Bending". *Management Science*, vol. 15, no. 5 (1969), pp. 315–320.