

Time-Reversed Dissipation Induces Duality Between Minimizing Gradient Norm and Function Value

Jaeyeon Kim[†], Asuman Ozdaglar[‡], Chanwoo Park[‡], Ernest K.Ryu[†]

[†]: SNU Mathematics , [‡]: MIT EECS

2023. 07. 29

Duality between methods

$f(x)$ reducing method $\xleftrightarrow{\text{Duality}}$ $\|\nabla f(x)\|^2$ reducing method

- We observed a **symmetric phenomenon** between pairs of **first-order methods**.
 - They have similar coefficients.
 - They share the convergence rate.
- Can we find a **duality framework** that generalizes this phenomenon?

First Order Methods and H-dual

Definition

f is L -smooth convex. An N -step First Order Methods with $H = \{h_{k,i}\}_{0 \leq i < k \leq N}$ is:

$$x_{k+1} = x_k - \frac{1}{L} \sum_{k=0}^i h_{k+1,i} \nabla f(x_i), \quad k = 0, 1, \dots, N-1.$$

Its H-dual is defined as

$$x_{k+1} = x_k - \frac{1}{L} \sum_{k=0}^i h_{N-i, N-1-k} \nabla f(x_i), \quad k = 0, 1, \dots, N-1,$$

$$\begin{bmatrix} h_{10} & 0 & \cdots & 0 \\ h_{20} & h_{21} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ h_{n0} & h_{n1} & \cdots & h_{n(n-1)} \end{bmatrix} \xrightarrow{\text{Anti-Transpose}} \begin{bmatrix} h_{n(n-1)} & 0 & \cdots & 0 \\ h_{n(n-2)} & h_{(n-1)(n-2)} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ h_{n0} & h_{(n-1)0} & \cdots & h_{10} \end{bmatrix}$$

Example 1 : OGM and OGM-G

(OGM)¹ and (OGM-G)^{2,3}

$$x_{k+1} = x_k - \frac{1}{L} \sum_{i=0}^k \left(\frac{\theta_i^2 (2\theta_i - 1)}{\theta_k^2 \theta_{k+1}} + \delta_{k,i} \right) \nabla f(x_i), \quad (\text{OGM})$$

$$y_{k+1} = y_k - \frac{1}{L} \sum_{i=0}^k \left(\frac{\theta_{N-k-1}^2 (2\theta_{N-k-1} - 1)}{\theta_{N-i-1}^2 \theta_{N-i}} + \delta_{N-k-1, N-i-1} \right) \nabla f(x_i) \quad (\text{OGM-G})$$

- They are **H-dual** of each other.
- Convergence rates:

$$f(x_N) - f_\star \leq \frac{1}{\theta_N^2} \frac{L}{2} \|x_0 - x_\star\|^2 \quad (\text{OGM}),$$

$$\frac{1}{2L} \|\nabla f(y_N)\|^2 \leq \frac{1}{\theta_N^2} (f(y_0) - f_\star) \quad (\text{OGM-G}).$$

¹[Kim and Fessler, 2016]

²[Kim and Fessler, 2021]

³ $\theta_{-1} = \theta_0 = 1, \theta_{i+1}^2 - \theta_i^2 = \theta_i^2$ for $0 \leq i \leq N-2, \theta_N^2 - \theta_{N-1}^2 = 2\theta_{N-1}^2$

These three pairs of methods share the **same factor of convergence rate** and their H -matrices are in the **Anti-Transpose** relationship.

$$\left[(OGM), (OGM-G) \right], \quad \left[(GD), (GD) \right], \quad \left[(OBL-F_b), (OBL-G_b) \right]^4$$

Main theorem:

A method reduces *function value* \leftrightarrow Its **H-dual** reduces *gradient norm*

⁴[Park and Ryu, 2021]

Energy function Structure and H-duality

$$x_{\star} := \operatorname{argmin} f(x).$$

$$f_{\star} := f(x_{\star}).$$

$$\llbracket x, y \rrbracket := f(y) - f(x) + \langle \nabla f(y), x - y \rangle + \frac{1}{2L} \|\nabla f(x) - \nabla f(y)\|^2 \leq 0.$$

$$\mathcal{U}_k := \frac{L}{2} \|x_0 - x_{\star}\|^2 + \sum_{i=0}^{k-1} u_i \llbracket x_i, x_{i+1} \rrbracket + \sum_{i=0}^k (u_i - u_{i-1}) \llbracket x_{\star}, x_i \rrbracket,$$

$$\mathcal{V}_k := v_0 (\llbracket y_N, y_{\star} \rrbracket + f(y_0) - f_{\star}) + \sum_{i=0}^{k-1} v_{i+1} \llbracket y_i, y_{i+1} \rrbracket + \sum_{i=0}^{k-1} (v_{i+1} - v_i) \llbracket y_N, y_i \rrbracket.$$

$u_{-1} = 0$. $\{u_i\}_{i=0}^N, \{v_i\}_{i=0}^N$ are **free variables**. If they're positive and *monotonically increasing*, $\{\mathcal{U}_k\}$ and $\{\mathcal{V}_k\}$ are *monotonically decreasing*.

Energy function structure and H-duality

If

$$u_N(f(x_N) - f_\star) \leq \mathcal{U}_N, \quad (\text{C1})$$

$$f(x_N) - f_\star \leq \frac{\mathcal{U}_N}{u_N} \leq \dots \leq \frac{\mathcal{U}_{-1}}{u_N} = \frac{1}{u_N} \frac{L}{2} \|x_0 - x_\star\|^2.$$

If

$$\frac{1}{2L} \|\nabla f(y_N)\|^2 \leq \mathcal{V}_N, \quad (\text{C2})$$

$$\frac{1}{2L} \|\nabla f(y_N)\|^2 \leq \mathcal{V}_N \leq \dots \leq \mathcal{V}_0 \leq v_0 (f(y_0) - f_\star).$$

Theorem

Assume $v_i = \frac{1}{u_{N-i}} > 0$ for $i = 0, \dots, N$.

$$[(\text{C1})] \quad \Leftrightarrow \quad [(\text{C2}) \text{ for the } H\text{-dual}].$$

Time-Reversed Dissipation

- Continuous-time limit $N \rightarrow \infty$ of first-order method becomes ODE.
- The notion of **H-dual** becomes:

$$\ddot{X}(t) + \gamma'(t)\dot{X}(t) + \nabla f(X(t)) = 0, \quad \ddot{Y}(t) + \gamma'(T-t)\dot{Y}(t) + \nabla f(Y(t)) = 0.$$

- Friction terms are **time-reversed**.

Applications of H-duality


- Constructing *gradient-norm* reducing method is quite hard rather than constructing the *function-value* reducing method.
- Utilizing **H-duality** makes it easier.
- In the composite minimization setting, we obtain (*Super FISTA-G*) (*SFG*), which reduces $\min_{v \in \partial F(x)} \|v\|^2$.
 - $F := f + h$, h is convex.
 - $h = 0$: $\min_{v \in \partial F(x)} \|v\|^2 = \|\nabla f(x)\|^2$.
- (*SFG*) is state-of-the-art : 5.26 times faster than (*FISTA-G*)⁵.


⁵[Lee et al., 2021]


Conclusion


- In this work, we defined the notion of **H-duality**.
- We established that the **H-dual** of a method reducing *function value* is another method reducing *gradient norm*.
- We are currently working on further generalizing **H-duality** in the *Mirror Descent* setting.


References


 Drori, Y. (2017).
The exact information-based complexity of smooth convex minimization.
Journal of Complexity, 39:1–16.

 Drori, Y. and Teboulle, M. (2014).
Performance of first-order methods for smooth convex minimization: a novel approach.
Mathematical Programming, 145(1):451–482.

 Kim, D. and Fessler, J. A. (2016).
Optimized first-order methods for smooth convex minimization.
Mathematical Programming, 159(1):81–107.

 Kim, D. and Fessler, J. A. (2021).
Optimizing the efficiency of first-order methods for decreasing the gradient of smooth convex functions.
Journal of Optimization Theory and Applications, 188(1):192–219.

 Lee, J., Park, C., and Ryu, E. K. (2021).
A geometric structure of acceleration and its role in making gradients small fast.
NeurIPS.

 Park, C. and Ryu, E. K. (2021).
Optimal first-order algorithms as a function of inequalities.
arXiv preprint arXiv:2110.11035.

Concluding Remark : (OGM) and (OGM-G)

$$x_{k+1} = x_k^+ + \frac{\theta_k - 1}{\theta_{k+1}}(x_k^+ - x_{k-1}^+) + \frac{\theta_k}{\theta_{k+1}}(x_k^+ - x_k) \quad (\text{OGM})$$

$$y_{k+1} = y_k^+ + \frac{(\theta_{N-k} - 1)(2\theta_{N-k-1} - 1)}{2\theta_{N-k}(2\theta_{N-k} - 1)}(y_k^+ - y_{k-1}^+) + \frac{2\theta_{N-k-1} - 1}{2\theta_{N-k} - 1}(y_k^+ - y_k) \quad (\text{OGM-G})$$

$$x_{k+1} = x_k^+ + \frac{\theta_k - 1}{\theta_{k+1}}(x_k^+ - x_{k-1}^+) \quad (\text{Nesterov's FGM})$$

(OGM) is exactly optimal ⁶.

⁶[Drori, 2017]

Concluding Remark : (SFG)

$$y_{k+1} = y_k^{\oplus,4} + \frac{(N-k+1)(2N-2k-1)}{(N-k+3)(2N-2k+1)} (y_k^{\oplus,4} - y_{k-1}^{\oplus,4}) + \frac{(4N-4k-1)(2N-2k-1)}{6(N-k+3)(2N-2k+1)} (y_k^{\oplus,4} - y_k)$$
$$y_N = y_{N-1}^{\oplus,4} + \frac{3}{10} (y_{N-1}^{\oplus,4} - y_{N-2}^{\oplus,4}) + \frac{3}{40} (y_{N-1}^{\oplus,4} - y_{N-1}) \quad (\text{SFG})$$

for $k = 0, \dots, N-2$, where

Theorem

(SFG) exhibits

$$\min_{v \in \partial F(y_N^{\oplus,4})} \|v\|^2 \leq \frac{50L}{(N+2)(N+3)} (F(y_0) - F_\star).$$

(SFG) is the state-of-the-art : 5 times faster than *FISTA-G*⁷.

⁷[Lee et al., 2021]

Example 2 : Gradient Descent

$$x_{k+1} = x_k - \frac{1}{L} \nabla f(x_k), \quad k = 0, \dots, N-1. \quad (GD)$$

- (GD) is **H-dual** of itself.
- Convergence rate⁸:

$$f(x_N) - f_\star \leq \frac{1}{2N+1} \frac{L}{2} \|x_0 - x_\star\|^2,$$
$$\frac{1}{2L} \|\nabla f(x_N)\|^2 \leq \frac{1}{2N+1} (f(x_0) - f_\star).$$

⁸[Drori and Teboulle, 2014, Kim and Fessler, 2021]