# Efficient Algorithms for Smooth Minimax Optimisation in Non-Euclidean Space Using Bregman Divergence Framework

Nitish V. Deshpande, Sandeep K. Routray
Department of Electrical Engineering, Indian Institute of Technology Kanpur
Emails: nitishvd@iitk.ac.in, sroutray@iitk.ac.in

EE698U Term Presentation

# Motivation

Minimax type of problems arise in several domains such as machine learning, optimization, statistics, communication, and game theory. However, a majority of results are established for the Euclidean norm due to its special self-dual nature.

**Motivation of analysis in Non-Euclidean Space**

(1) The quadratic proximity term $\frac{1}{2\eta_k} \|x - x_k\|_2^2$ is inappropriate for problems with highly inhomogeneous geometry

(2) For example, the quadratic minimization problem
$\min_{x \in \mathbb{R}^n} f(x) = \frac{1}{2} (x - x_0)^T Q (x - x_0)$ where $Q \succ 0$ is a diagonal matrix with high condition number. The inhomogeneous geometry leads to slow updates (since iteration complexity depends on condition number for the quadratic minimization problem) in the conventional GD algorithm.

(3) For the probability simplex problem, Euclidean distance is in general not recommended for measuring the distance between probability vectors

(4) To tackle these issues, the mirror gradient descent algorithm was introduced which adjusts the gradient updates to fit the problem geometry. The notion of mirror GD is to replace the quadratic proximity term $\frac{1}{2\eta_k} \|x - x_k\|_2^2$ by a class of general distance-like metric known as the Bregman divergence

## Prior work

(1) Work on Smooth Minimax Optimisation:[1] $\tilde{\mathcal{O}}(1/k^2)$ convergence rate for smooth, strongly-convex – concave problems, improving upon the previous best known rate of $\mathcal{O}(1/k)$

(2) Work on Non-Euclidean analysis:

   (1) General Norm[2] They analyse Nesterov's accelerated gradient descent using general norm for the unconstrained case. They develop a potential function based framework for proving convergence rate.

   (2) Relative strong convexity and smoothness:[3] They develop a notion of "relative smoothness" and relative strong convexity that is determined relative to a user-specified "reference function" $h(.)$. However, extension of this notion to accelerated gradient descent is still an open problem.

   (3) Riemannian space:[4] Proposed a Riemannian counterpart to Nesterov's AGD.

[1] Kiran K Thekumparampil et al. "Efficient algorithms for smooth minimax optimization". In: *Advances in Neural Information Processing Systems*. 2019, pp. 12680–12691.

[2] Nikhil Bansal and Anupam Gupta. "Potential-function proofs for first-order methods". In: *arXiv preprint arXiv:1712.04581* (2017).

[3] Haihao Lu, Robert M Freund, and Yurii Nesterov. "Relatively smooth convex optimization by first-order methods, and applications". In: *SIAM Journal on Optimization* 28.1 (2018), pp. 333–354.

[4] Kwangjun Ahn and Suvrit Sra. "From Nesterov's Estimate Sequence to Riemannian Acceleration". In: *arXiv preprint arXiv:2001.08876* (2020).

$$\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} g(x, y) \quad , \quad g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$$

**Assumptions**

**(A1)** $\mathcal{X}, \mathcal{Y} \subseteq \mathcal{V}$ where $\mathcal{V}$ is a normed vector space with an arbitrary norm $\|.\|$ on the underlying space

**(A2)** $g(x, .)$ is concave for every $x$ and $\sigma$-strongly convex for every $g(., y)$ for every $y$.

**(A3)** $\mathcal{Y}$ is a compact set , there exists a finite $D_{\mathcal{Y}} = \max_{y, y' \in \mathcal{Y}} \|y - y'\|$ also known as the diameter of $\mathcal{Y}$.

**(A4)** $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is $L$-smooth

$$\max \left\{ \left\| \nabla_x g(x, y) - \nabla_x g(x', y') \right\|_*, \left\| \nabla_y g(x, y) - \nabla_y g(x', y') \right\|_* \right\} \leq L \left( \left\| x - x' \right\| + \left\| y - y' \right\| \right)$$

**Our goal** To find a $\epsilon$-primal-dual pair $(\hat{x}, \hat{y})$ defined as: $(\hat{x}, \hat{y})$ is an $\epsilon$-primal-dual pair of $g$ if the primal-dual gap is less than $\epsilon$: $\max_{y \in \mathcal{Y}} g(\hat{x}, y) - \min_{x \in \mathcal{X}} g(x, \hat{y}) \leq \epsilon$

# Contributions

(1) Using potential function based framework[5] and Bregman divergence framework, proved $\mathcal{O}(1/k^2)$ convergence rate for Nesterov's AGD for the *general norm* and *constrained case*.

(2) Proposed Generalized Conceptual Dual Implicit Accelerated Gradient Descent (GC-DIAG) which is adapted from the Conceptual Dual Implicit Accelerated Gradient (C-DIAG)[6] and proved $\mathcal{O}(1/k^2)$ convergence rate for the primal dual gap.

(3) Proved $\mathcal{O}(\frac{1}{k^4})$ convergence rate using Nesterov's AGD and restarting strategy which is an improvement over $\mathcal{O}(\frac{1}{k})$ for smooth and strongly convex functions with respect to an arbitrary norm.

| | Smooth and convex | Smooth and strongly convex |
|---|---|---|
| Mirror descent | $\mathcal{O}(\frac{1}{\sqrt{k}})$ | $\mathcal{O}(\frac{1}{k})$[7] |
| Nesterov's AGD | $\mathcal{O}(\frac{1}{k^2})$[8] | $\mathcal{O}(\frac{1}{k^4})$ |

Table: Comparison of oracle complexities with arbitrary norm

---

[5]Bansal and Gupta, "Potential-function proofs for first-order methods".

[6]Thekumparampil et al., "Efficient algorithms for smooth minimax optimization".

[7]Yhli. *Minimizing a Strongly Convex Function by Mirror Descent.* 2017. URL: http://yenhuanli.github.io/blog/2017/05/05/mirror-descent-str/.

[8]Y. Nesterov. "A method for solving the convex programming problem with convergence rate O(1/k²)".

# Nesterov's accelerated gradient ascent with general norm

---

**Algorithm 1** Nesterov's accelerated gradient ascent with general norm

---

**Input:** Smooth concave function $h(.)$, learning rate $\frac{1}{\beta}$, Bregman divergence $D_\psi(.\|.)$, initial point $y_0$ and $z_0$

**Output:** $y_K$

1: **for** $k = 0, 1, ..., K$ **do**

$$w_k \leftarrow (1 - \tau_k)y_k + \tau_k z_k \tag{1}$$

$$y_{k+1} \leftarrow \underset{y \in \mathcal{Y}}{\arg\min} \left\{ -\langle \nabla h(w_k), y - w_k \rangle + \frac{\beta}{2} \|y - w_k\|^2 \right\} \tag{2}$$

$$z_{k+1} \leftarrow \underset{z \in \mathcal{Y}}{\arg\min} \left\{ -\eta_k \langle \nabla h(w_k), z \rangle + D_\psi(z\|z_k) \right\} \tag{3}$$

2: **end for**

---

For Euclidean norm (2) becomes $y_{k+1} \leftarrow \mathcal{P}_\mathcal{Y}(\omega_k + \frac{1}{\beta}\nabla h(w_k))$ and (3) becomes
$z_{k+1} \leftarrow \mathcal{P}_\mathcal{Y}(z_k + \eta_k \nabla h(w_k))$
A major hurdle in general norm case: We cannot use properties of Projection operator like Non-expansiveness

# Analysis using Potential function

$$\Phi(k) = k(k+1)(h(y^*) - h(y_k)) + \frac{4\beta}{\mu_\psi} D_\psi(y^* \| z_k)$$

**Our goal:**

$$\Phi(k+1) \leq \Phi(k)$$

## Lemma

Suppose $h(.)$ is an $L$-smooth function and the parameters of Algorithm 1 are chosen so that $\beta > L$, $\eta_k = \frac{(k+1)}{2\beta}\mu_\psi$ and $\tau_k = \frac{2}{k+2}$. Then, we have

$$\Phi(k+1) \leq \Phi(k)$$

$$\Phi(k+1) - \Phi(k) = (k+1)(k+2) \underbrace{(h(\mathsf{w}_k) - h(\mathsf{y}_{k+1}))}_{(a)}$$

$$\underbrace{-k(k+1)(h(\mathsf{w}_k) - h(\mathsf{y}_k)) + 2(k+1)(h(\mathsf{y}) - h(\mathsf{w}_k))}_{(b)} + \frac{4\beta}{\mu_\psi} \underbrace{(D_\psi(\mathsf{y}\|\mathsf{z}_{k+1}) - D_\psi(\mathsf{y}\|\mathsf{z}_k))}_{(c)}$$

$$(4)$$

**For bounding (a) and (b), we used:**

**(1)** The fact that $(-h(\mathsf{x}))$ is $L$-smooth and the choice of $\beta > L$.

**(2)** Concavity of $h(\mathsf{x})$ and the choice of $\tau_k = \frac{2}{k+2}$.

## Analysis using Potential function

$D_\psi(y\|z_{k+1}) - D_\psi(y\|z_k)$

$= (\psi(y) - \psi(z_{k+1}) - \langle \nabla\psi(z_{k+1}), y - z_{k+1}\rangle) - (\psi(y) - \psi(z_k) - \langle \nabla\psi(z_k), y - z_k\rangle)$

$= \psi(z_k) - \psi(z_{k+1}) + \langle \nabla\psi(z_k), z_{k+1} - z_k\rangle + \langle \nabla\psi(z_{k+1}) - \nabla\psi(z_k), z_{k+1} - y\rangle$

$$\leq -\frac{\mu_\psi}{2}\|z_{k+1} - z_k\|^2 + \underbrace{\langle \nabla\psi(z_{k+1}) - \nabla\psi(z_k), z_{k+1} - y\rangle}_{(d)} \tag{5}$$

Inequality is due to $\mu_\psi$-strongly convex function $\psi(.)$
From the update in (3) in Algorithm 1, we write the optimality condition as

$$\left\langle (-\eta_k \nabla h(w_k) + \nabla_z D_\psi(z\|z_k))\big|_{z=z_{k+1}}, y - z_{k+1}\right\rangle \geq 0, \quad \forall y \in \mathcal{Y}. \tag{6}$$

From definition of Bregman divergence, $\nabla_z D_\psi(z\|z_k)\big|_{z=z_{k+1}} = \nabla\psi(z_{k+1}) - \nabla\psi(z_k)$.
Hence, the term $(d)$ in Equation (5) can be bounded as

$$\langle \nabla\psi(z_{k+1}) - \nabla\psi(z_k), z_{k+1} - y\rangle \leq \langle \eta_k \nabla h(w_k), z_{k+1} - y\rangle. \tag{7}$$

Hence,

$$D_\psi(y\|z_{k+1}) - D_\psi(y\|z_k) \leq -\frac{\mu_\psi}{2}\|z_{k+1} - z_k\|^2 + \langle \eta_k \nabla h(w_k), z_{k+1} - y\rangle. \tag{8}$$

$$\Phi(k+1) - \Phi(k) = (k+1)(k+2)(h(w_k) - h(y_{k+1}))$$

$$- k(k+1)(h(w_k) - h(y_k)) + 2(k+1)(h(y) - h(w_k)) + \frac{4\beta}{\mu_\psi}(D_\psi(y\|z_{k+1}) - D_\psi(y\|z_k))$$

$$h(y_{k+1}) - h(w_k) \geq \tau_k \langle \nabla h(w_k), z_{k+1} - z_k \rangle - \frac{\beta}{2}\tau_k^2 \|z_{k+1} - z_k\|^2 \qquad (9)$$

$$- k(k+1)(h(w_k) - h(y_k)) + 2(k+1)(h(y) - h(w_k)) \leq 2(k+1) \langle \nabla h(w_k), y - z_k \rangle, \quad (10)$$

$$D_\psi(y\|z_{k+1}) - D_\psi(y\|z_k) \leq -\frac{\mu_\psi}{2} \|z_{k+1} - z_k\|^2 + \langle \eta_k \nabla h(w_k), z_{k+1} - y \rangle. \qquad (11)$$

$$\Phi(k+1) - \Phi(k) \leq (k+1)(k+2)\left(-\tau_k \langle \nabla h(w_k), z_{k+1} - z_k \rangle + \frac{\beta}{2}\tau_k^2 \|z_{k+1} - z_k\|^2\right)$$

$$+ 2(k+1) \langle \nabla h(w_k), y - z_k \rangle + \frac{4\beta}{\mu_\psi}\left(-\frac{\mu_\psi}{2} \|z_{k+1} - z_k\|^2 + \langle \eta_k \nabla h(w_k), z_{k+1} - y \rangle\right)$$

$$\leq 2\beta \|z_{k+1} - z_k\|^2 \left(\frac{k+1}{k+2} - 1\right) + \left(-2(k+1) + \frac{4\beta}{\mu_\psi}\eta_k\right) \langle \nabla h(w_k), z_{k+1} - y \rangle \overset{(7)}{\leq} 0, \quad (12)$$

inequality (7) follows from the choice of $\eta_k = \frac{(k+1)}{2\beta}\mu_\psi$.

# Generalized Conceptual Dual Implicit Accelerated Gradient Descent (GC-DIAG)

---

**Algorithm 2** Generalized Conceptual Dual Implicit Accelerated Gradient Descent (GC-DIAG) for strongly-convex-concave programming

---

    **Input:** $g$, $D_\psi, \mu_\psi$, $L$, $\sigma$, $\mathsf{x}_0$, $\mathsf{y}_0$, $K$,

    **Output:** $\bar{\mathsf{x}}_K$ $\mathsf{y}_K$

1: Set $\beta > L$, $\mathsf{z}_0 \leftarrow \mathsf{y}_0$

2: **for** $k = 0, 1, ..., K$ **do**

3:      $\tau_k \leftarrow \frac{2}{k+2}$, $\eta_k \leftarrow \frac{k+1}{2\beta} \mu_\psi$, $\mathsf{w}_k \leftarrow (1 - \tau_k)\, \mathsf{y}_k + \tau_k \mathsf{z}_k$

4:      Choose $\mathsf{x}_{k+1}$, $\mathsf{y}_{k+1}$, ensuring:

$$g(\mathsf{x}_{k+1}, \mathsf{y}_{k+1}) = \min_{\mathsf{x}} g(\mathsf{x}, \mathsf{y}_{k+1}),$$

$$\mathsf{y}_{k+1} \leftarrow \arg\min_{\mathsf{y} \in \mathcal{Y}} \left\{ -\langle \nabla_{\mathsf{y}} g(\mathsf{x}_{k+1}, \mathsf{w}_k), \mathsf{y} - \mathsf{w}_k \rangle + \frac{\beta}{2} \|\mathsf{y} - \mathsf{w}_k\|^2 \right\}$$

5:      $\mathsf{z}_{k+1} \leftarrow \arg\min_{\mathsf{z} \in \mathcal{Y}} \left\{ -\eta_k \langle \nabla_{\mathsf{y}} g(\mathsf{x}_{k+1}, \mathsf{w}_k), \mathsf{z} \rangle + D_\psi(\mathsf{z} \| \mathsf{z}_k) \right\}$,

6:      $\bar{\mathsf{x}}_{k+1} \leftarrow \frac{2}{(k+1)(k+2)} \sum_{i=1}^{k+1} i \mathsf{x}_i$

7: **end for**

8: **return** $\bar{\mathsf{x}}_K, \mathsf{y}_K$

$$\underbrace{(k+1)(k+2)\left(g(\mathsf{x}_{k+1},\mathsf{y}) - g(\mathsf{x}_{k+1},\mathsf{y}_{k+1})\right) + \frac{4\beta}{\mu_\psi}D_\psi(\mathsf{y}\|\mathsf{z}_{k+1})}_{\Phi_{k+1}}$$

$$\stackrel{(1)}{\leq} (k)(k+1)\left(g(\mathsf{x}_{k+1},\mathsf{y}) - g(\mathsf{x}_{k+1},\mathsf{y}_k)\right) + \frac{4\beta}{\mu_\psi}D_\psi(\mathsf{y}\|\mathsf{z}_k)$$

$$\stackrel{(2)}{=} (k)(k+1)\left(g(\mathsf{x}_{k+1},\mathsf{y}) - g(\mathsf{x}_k,\mathsf{y})\right) + (k)(k+1)\left(g(\mathsf{x}_k,\mathsf{y}) - g(\mathsf{x}_{k+1},\mathsf{y}_k)\right) + \frac{4\beta}{\mu_\psi}D_\psi(\mathsf{y}\|\mathsf{z}_k)$$

$$\stackrel{(3)}{\leq} (k)(k+1)\left(g(\mathsf{x}_{k+1},\mathsf{y}) - g(\mathsf{x}_k,\mathsf{y})\right) + \underbrace{(k)(k+1)\left(g(\mathsf{x}_k,\mathsf{y}) - g(\mathsf{x}_k,\mathsf{y}_k)\right) + \frac{4\beta}{\mu_\psi}D_\psi(\mathsf{y}\|\mathsf{z}_k)}_{\Phi_k}$$

$$\stackrel{(4)}{\leq} (k)(k+1)\,g(\mathsf{x}_{k+1},\mathsf{y}) - \sum_{i=1}^{k}(2i)(g(\mathsf{x}_i,\mathsf{y})) + \frac{4\beta}{\mu_\psi}D_\psi(\mathsf{y}\|\mathsf{z}_0), \tag{13}$$

Equality (2) follows by adding and subtracting $k(k+1)g(\mathsf{x}_k,\mathsf{y})$. Inequality (3) follows from the update rule $g(\mathsf{x}_k,\mathsf{y}_k) = \min_\mathsf{x} g(\mathsf{x},\mathsf{y}_k)$ in Step 4 of Algorithm 2 and hence $g(\mathsf{x}_k,\mathsf{y}_k) \leq g(\mathsf{x}_{k+1},\mathsf{y}_k)$. Inequality (4) follows from the recurrence relation established in inequality (3).

Now, we write

$$(k+1)(k+2)\left(g(x_{k+1}, y) - g(x_{k+1}, y_{k+1})\right) + \frac{4\beta}{\mu_\psi} D_\psi(y\|z_{k+1}) \tag{14}$$

$$\leq (k)(k+1)\, g(x_{k+1}, y) - \sum_{i=1}^{k}(2i)(g(x_i, y)) + \frac{4\beta}{\mu_\psi} D_\psi(y\|z_0). \tag{15}$$

Rearranging the terms, we get

$$\sum_{i=1}^{k+1}(2i)(g(x_i, y)) - (k+1)(k+2)g(x_{k+1}, y_{k+1}) \leq \frac{4\beta}{\mu_\psi} D_\psi(y\|z_0) - \frac{4\beta}{\mu_\psi} D_\psi(y\|z_{k+1}) \tag{16}$$

$$\sum_{i=1}^{k+1}(2i)(g(x_i, y)) - (k+1)(k+2)g(x, y_{k+1}) \overset{(5)}{\leq} \frac{4\beta}{\mu_\psi} D_\psi(y\|z_0) \tag{17}$$

$$g(\bar{x}_{k+1}, y)) - g(x, y_{k+1}) \overset{(6)}{\leq} \frac{4\beta}{\mu_\psi(k+1)(k+2)} D_\psi(y\|z_0) \tag{18}$$

$$\max_{y \in \mathcal{Y}} g(\bar{x}_{k+1}, y)) - \min_{x \in \mathcal{X}} g(x, y_{k+1}) \overset{(7)}{\leq} \frac{4\beta}{\mu_\psi(k+1)(k+2)} D_\psi(y\|z_0) \tag{19}$$

where inequality (5) follows from the fact that $g(x, y_{k+1}) \geq g(x_{k+1}, y_{k+1}) \forall x \in \mathcal{X}$.
Inequality (6) follows by defining a convex combination of $\bar{x}_{k+1} = \frac{1}{(k+1)(k+2)} \sum_{i=1}^{k+1}(2i)x_i$
and from the fact that $g(., y)$ is convex for every y.

**Algorithm 3** Generalized Dual Implicit Accelerated Gradient Descent (G-DIAG) for strongly-convex-concave programming

---

**Input:** $g$, $D_\psi$, $\mu_\psi$, $L$, $\sigma$, $x_0$, $y_0$, $K$, $\left\{\epsilon_{\text{step}}^{(k)}\right\}_{k=1}^{K}$

**Output:** $\bar{x}_K$ $y_K$

1: Set $\beta \leftarrow ?$, $z_0 \leftarrow y_0$
2: **for** $k = 0, 1, ..., K$ **do**
3:     $\tau_k \leftarrow ?$, $\eta_k \leftarrow ?$, $w_k \leftarrow (1 - \tau_k)\, y_k + \tau_k z_k$
4:     $x_{k+1}$, $y_{k+1} \leftarrow \texttt{Imp} - \texttt{STEP}(g, L, \sigma, x_0, w_k, \beta, \epsilon_{\text{step}}^{(k+1)})$, ensuring:

$$g(x_{k+1}, y_{k+1}) \leq \min_x g(x, y_{k+1}) + \epsilon_{\text{step}}^{(k+1)}$$

$$y_{k+1} \leftarrow \underset{y \in \mathcal{Y}}{\arg \min} \left\{ -\langle \nabla_y g(x_{k+1}, w_k), y - w_k \rangle + \frac{\beta}{2} \|y - w_k\|^2 \right\}$$

5:     $z_{k+1} \leftarrow \arg \min_{z \in \mathcal{Y}} \left\{ -\eta_k \langle \nabla_y g(x_{k+1}, w_k), z \rangle + D_\psi(z \| z_k) \right\}$,
6:     $\bar{x}_{k+1} \leftarrow \frac{2}{(k+1)(k+2)} \sum_{i=1}^{k+1} i x_i$
7: **end for**
8: **return** $\bar{x}_K, y_K$

---

**Algorithm 4** Imp Step subroutine in G-DIAG

1: $\text{Imp} - \text{STEP}(g, L, \sigma, x_0, w, \beta, \epsilon_{\text{step}}^{(k+1)})$:
2:      Set $\epsilon_{\text{mp}} \leftarrow ?$, $R \leftarrow ?$, $\epsilon_{\text{agd}} \leftarrow ?$, $y_0 \leftarrow w$
3:      **for** $r = 0, 1, ..., R$ **do**
4:          Starting at $x_0$ use generalized AGD (Algorithm 1 with $-g(., y_r)$ to compute $x_r$ such that:

$$g(\hat{x}_r, y_r) \leq \min_x g(x, y_r) + \epsilon_{\text{agd}}, \tag{20}$$

5:          $y_{k+1} \leftarrow \arg\min_{y \in \mathcal{Y}} \left\{ -\langle \nabla_y g(\hat{x}_r, w), y - w \rangle + \frac{\beta}{2} \|y - w\|^2 \right\}$
6:      **end for**
7:      **return** $\hat{x}_R, y_{R+1}$

In the $\text{Imp} - \text{Step}$ of the original DIAG algorithm,[9] AGD (with $\ell^2$ norm) is used to efficiently calculate (in logarithmic number of steps) an estimate for $x_r = \arg\min_{x \in \mathcal{X}} g(x, y_r) + \epsilon_{\text{agd}}$. The proof uses the guarantee on AGD for strongly convex case.[10]

---

[9] Thekumparampil et al., "Efficient algorithms for smooth minimax optimization".
[10] Bansal and Gupta, "Potential-function proofs for first-order methods", Equation (5.68).

We faced the the following issues

**(A) Analyzing AGD for strongly convex case in arbitrary norm**

**(1)** In our setting, strong convexity is defined with respect to some arbitrary norm. We have not found any literature that tackles this problem using AGD.

**(2)** It may happen that, for arbitrary norm, finding an $\epsilon_{\mathrm{agd}}$ estimate $x_r$ would not be possible in linear time.

**(3)** As an example, we found a blog[11] that discusses the case of mirror descent for the strongly convex case with respect to a general norm.

**(4)** They note that the improved oracle complexity (due to strong convexity) is $\mathcal{O}(1/k)$ only. So, it may not be possible to achieve the linear convergence in an arbitrary norm case.

**(B) Proving that the** $\mathtt{Imp - Step}$ **convergence:** Let $x^*(y) = \arg\min_{x \in \mathcal{X}} g(x, y)$, then for proving that there exists a fixed point of the iterations of the $\mathtt{Imp - Step}$, we require to show $y^+ = \arg\min_{y \in \mathcal{Y}} \left\{ -\langle \nabla_y g(x^*(y), w), y - w \rangle + \frac{\beta}{2} \|y - w\|^2 \right\}$ is a contraction, that is $\|y_1^+ - y_2^+\| \le \alpha \|y_1 - y_2\|, \ \alpha \le 1$.

---

[11] Yhli, *Minimizing a Strongly Convex Function by Mirror Descent*.

# Convergence analysis for generalized AGD with strong concavity

(1) In case of mirror descent, the restarting strategy discussed in[12] improves the convergence rate from $\mathcal{O}(1/\sqrt{k})$ to $\mathcal{O}(1/k)$ through introduction of strong convexity (or concavity in our case).

(2) We try to come up with the convergence rate for generalized AGD with strong concavity following a similar procedure.

(3) Using Lemma 7, the oracle complexity for generalized AGD can be found by carrying out a telescoping sum.

## Theorem

Suppose $h(.)$ is a $L$-smooth function and the parameters of Algorithm 1 are chosen as per Lemma 7, then the following holds

$$h(y^*) - h(y_T) \leq \frac{4\beta}{\mu_\psi} \cdot \frac{D_\psi(y^*||y_0)}{T(T+1)} \tag{21}$$

If we also assume $\Omega = \max_{y \in \mathcal{Y}} D_\psi(y||y_0)$, then the following bound would hold

$$h(y^*) - h(y_T) \leq \frac{4\beta}{\mu_\psi} \cdot \frac{\Omega}{T(T+1)} \tag{22}$$

---

[12]Yhli, *Minimizing a Strongly Convex Function by Mirror Descent.*

# Convergence analysis for generalized AGD with strong concavity

Now, assume that $h(.)$ is also $\sigma$-strongly concave with respect to some norm $\|.\|$. Further, $\psi(.)$ is chosen such that it is $\mu_\psi$-strongly convex on the whole $\mathbb{R}^d$, *instead of only on $\mathcal{Y}$*. For any $R > 0$ and u, define $\psi_{R,u}(y) := \psi(R^{-1}(y - u))$. Let $D_{\psi,R,z}(.\|.)$ denote the corresponding Bregman divergence.

## Corollary

Suppose
$$\Omega = \max\left\{ D_\psi(v\|0)|\, \|v\| \leq 1, v \in \mathbb{R}^d \right\}, \quad R_0 = \|y^* - y_0\|$$

If we apply Algorithm 1 with $\eta_k = \frac{(k+1)}{2R_0\beta}$, $\tau_k = \frac{2}{k+2}$, learning rate $\frac{1}{R_0\beta}$ for some $\beta > L$ and Bregman divergence $D_{\psi,R,y_0}(.\|.)$ for $T$ iterations. Then, the following bounds hold

$$h(y^*) - h(y_T(R_0, y_0)) \leq \frac{4R_0\beta}{\mu_\psi} \cdot \frac{\Omega}{T(T+1)} \tag{23}$$

$$\|y^* - y_T(R_0, y_0)\|^2 \leq \frac{8R_0\beta}{\mu_\psi\mu} \cdot \frac{\Omega}{T(T+1)} \tag{24}$$

where $y^*$ is the unique maximizer of $h(.)$ on $\mathcal{Y}$.

## Proof.

Consider the norm $\|.\|_{R_0} := R_0^{-1} \|.\|$. Note that $h(.)$ is $R_0 L$-smooth and $\psi_{R,u}(.)$ is $\mu_\psi$-strongly convex with respect to $\|.\|_{R_0}$. In this case, $\Omega$ would become

$$D_{\psi, R_0, y_0}(y^*||y_0) = D_\psi(R_0^{-1}(y^* - y_0)||0) = \max\left\{ D_\psi(v||0) \mid \|v\| \le 1, v \in \mathbb{R}^d \right\} := \Omega$$

The bound in (23) follows directly from Theorem 17. The bound in Equation (24) is obtained from (23) using strong-concavity of $h(.)$ and the optimality of $y^*$.

$$h(y^*) - h(y_T) \ge \langle \nabla h(y^*), y^* - y_T \rangle + \frac{\mu}{2} \|y^* - y_T\|^2, \quad \langle \nabla h(y^*), y^* - y_T \rangle \ge 0$$

$\square$

## Convergence analysis for generalized AGD with strong concavity

The error bounds given by Equations (23) (24) depend on $R_0$ with smaller $R_0$ giving smaller error bounds. Also, the bound of the distance between $y^*$ and $y_T$ (24) is strictly decreasing with iterations $T$. These observations can be used to design the following restarting strategy.

### Restarting strategy

(1) Set $y_0 \in \mathcal{Y}$, $l = 0$
(2) Set $T_l$ such that $\|y^* - y_{T_l}(R_l, y_l)\|^2 \leq 2^{-1} R_l^2$.
(3) Compute $y_{l+1} = y_{T_l}(R_l, y_l)$ using generalized AGD as per Corollary 18.
(4) Set $R_{l+1}^2 = 2^{-1} R_l^2$, $l = l + 1$. Go to step 2.

By Corollary 18, it suffices to choose $T_l$ such that

$$\frac{8R_l\beta}{\mu_\psi\mu} \cdot \frac{\Omega}{T_l(T_l+1)} \leq \frac{8R_l\beta}{\mu_\psi\mu} \cdot \frac{\Omega}{T_l^2} \leq 2^{-1} R_l^2$$

$$\implies T_l = \left\lceil \sqrt{\frac{16\beta\Omega}{R_l\mu_\psi\mu}} \right\rceil \tag{25}$$

The total number of AGD iterations required to get $y_L$ is defined as $M_L = \sum_{l=0}^{L-1} T_l$.

# Convergence analysis for generalized AGD with strong concavity

## Proposition

Let $L^*$ be the largest L such that

$$L \geq \sqrt{\frac{8\beta\Omega}{R_0\mu_\psi\mu}}2^{(L+1)/4}$$

Then, the proposed restarting strategy guarantees the following bound

$$h(y^*) - h(y_L) \leq 2^{-(0.5M_L+1)}\mu R_0^2, \quad \text{for } L \leq L^* \tag{26}$$

$$\|y^* - y_L\|^2 \leq 2^{-0.5M_L}R_0^2, \quad \text{for } L \leq L^* \tag{27}$$

and

$$h(y^*) - h(y_L) \leq \frac{1024\beta^2\Omega^2}{\mu_\psi^2\mu M_L^4}, \quad \text{for } L > L^* \tag{28}$$

$$\|y^* - y_L\|^2 \leq \frac{2048\beta^2\Omega^2}{\mu_\psi^2\mu^2 M_L^4}, \quad \text{for } L > L^* \tag{29}$$

# Convergence analysis for generalized AGD with strong concavity

## Proof.

Using the proposed restarting strategy, it follows from Corollary 18 that

$$h(y^*) - h(y_L) \leq 2^{-(L+1)} \mu R_0^2 \tag{30}$$

$$\|y^* - y_L\|^2 \leq 2^{-L} R_0^2 \tag{31}$$

By the choice of $T_l$ given in Equation (25), it holds that

$$M_L \leq L + \sum_{l=0}^{L-1} \sqrt{\frac{16\beta\Omega}{R_l \mu_\psi \mu}} = L + \sum_{l=0}^{L-1} \sqrt{\frac{16\beta\Omega}{R_0 \mu_\psi \mu}} 2^{l/4} \leq L + \sqrt{\frac{8\beta\Omega}{R_0 \mu_\psi \mu}} 2^{(L+1)/4}$$

Therefore, depending on value of $L$, the following holds

$$M_L \leq 2L, \quad \text{for } L \leq L^* \tag{32}$$

$$M_L \leq \sqrt{\frac{32\beta\Omega}{R_0 \mu_\psi \mu}} 2^{(L+1)/4}, \quad \text{for } L > L^* \tag{33}$$

The bounds given in Equations (26)-(29) follow by eliminating $L$ from Equations (30), (31) using the relations in Equations (32), (33). $\qquad \square$

**(1)** We extended the framework of Conceptual Dual Implicit Accelerated Gradient Descent to arbitrary norm case and proved $\mathcal{O}(1/k^2)$ convergence rate using Bregman divergence framework

**(2)** A key intermediate step involved proving convergence guarantee for Nesterov's AGD for a strongly convex and smooth function with respect to an arbitrary norm. We improved the convergence to $\mathcal{O}(\frac{1}{k^4})$ using the restarting strategy

**(3)** We plan to use the notion of relative smoothness and strong convexity to prove the contraction bound required for the inexact version of Dual Implicit Accelerated Gradient Descent.

Thank You