

Mirrorless Mirror Descent

Nati Srebro (TTIC)

Based on work with

Suriya Gunasekar (TTIC→MSR), Blake Woodworth (TTIC→ENS), Filip Hanzely (TTIC)

Plus: Bonus spotlight on acceleration, parallelization and interpolation learning!

A presentation in honor of Yurii Nesterov's 65th Birthday

$$\min_{x \in \mathbb{R}^d} f(x)$$

$$x_{k+1} \leftarrow \arg \min_x \langle \nabla f(x_k), x \rangle + \frac{1}{\eta_k} D_\Psi(x; x_k)$$

$$= \nabla \Psi^{-1}(\nabla \Psi(x_k) - \eta_k \nabla f(x_k))$$

$$D_\Psi(x; y) = \Psi(x) - (\Psi(y) + \langle \nabla \Psi(y), x - y \rangle)$$

Classical analysis:

$$D_\Psi(x; y) \geq \frac{1}{2} \|x - y\|^2$$

$$\|\nabla f(x) - \nabla f(y)\|_* \leq L \|x - y\|^2$$

Or, with relative smoothness:

$$\lambda \nabla^2 \Psi \preceq \nabla^2 f \preceq L \nabla^2 \Psi$$

[Bolte, Bauschke Teboulle][Lu Freund Nesterov]

→

$$k \leq \frac{LD_\Psi(x_0; x^*)}{\epsilon_k}$$

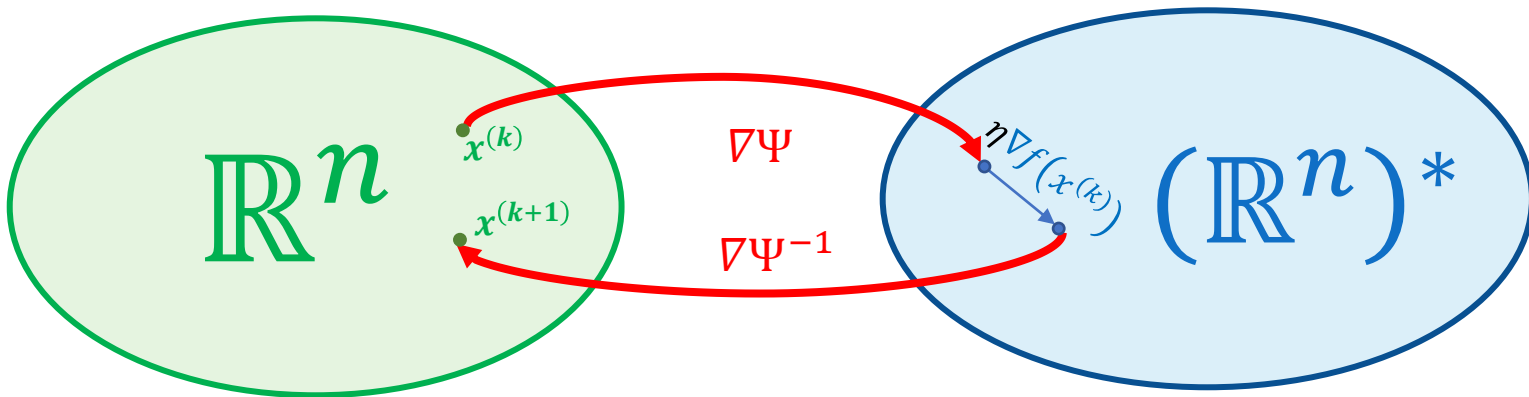
$\epsilon_k = f(x_k) - f(x^*)$

$$k \leq \frac{L}{\lambda} \log \left(\frac{LD_\Psi(x_0; x^*)}{\epsilon_k} \right)$$

With restarts (non-montone stepsize):

[as in, e.g., Roulet d'Aspremont '17, Woodworth S '21]

$$k \leq \frac{L}{\lambda} \log \left(\frac{\epsilon_0}{\epsilon_k} \right)$$



$$\dot{x}(t) = - \mathbf{H} \left(x(\lfloor t \rfloor_\eta) \right)^{-1} \nabla f \left(x(\lfloor t \rfloor_\eta) \right)$$

$$\lfloor t \rfloor_\eta = \eta \lfloor t/\eta \rfloor$$

$$\mathbf{H}(x) = \nabla^2 \Psi(x)$$

$\eta \rightarrow 0$

Forward Euler
Full Discretization
 $x_k = x(k\eta)$

$$x_{k+1} \leftarrow \nabla \Psi^{-1} \left(\nabla \Psi(x_k) - \eta \nabla f(x_k) \right)$$

$$x_{k+1} \leftarrow x_k - \eta \mathbf{H} (x_k)^{-1} \nabla f(x_k)$$

“Natural Gradient Descent” [Amari]

Forward Euler

$$\frac{d}{dt} \nabla \Psi(x(t)) = -\nabla f(x(t))$$

$$\dot{x}(t) = - \mathbf{H} \left(x(t) \right)^{-1} \nabla f \left(x(\lfloor t \rfloor_\eta) \right)$$

Partial Discretization
 $x_k = x(k\eta)$

$$\mathbf{H}(x) = \nabla^2 \Psi(x)$$

$\lfloor t \rfloor_\eta = \eta \lfloor t/\eta \rfloor$
 Forward Euler
 Full Discretization
 $x_k = x(k\eta)$

$$x_{k+1} \leftarrow \nabla \Psi^{-1} \left(\nabla \Psi(x_k) - \eta \nabla f(x_k) \right)$$

“Mirror Descent” / Bregman Gradient Scheme

$$x_{k+1} \leftarrow x_k - \eta \mathbf{H} \left(x_k \right)^{-1} \nabla f(x_k)$$

“Natural Gradient Descent” [Amari]

Proof: define $x(t) \doteq \nabla \Psi^{-1} \left(\nabla \Psi(x_k) - (t - k\eta) \nabla f(x_k) \right)$
 $k\eta \leq t \leq (k+1)\eta$

$$\begin{aligned} \dot{x}(t) &= \nabla^2 \Psi^{-1}(\dots) \left(-\nabla f(x_k) \right) \\ &= -\nabla^2 \Psi(x(t))^{-1} \nabla f(x(\lfloor t \rfloor_\eta)) \end{aligned}$$

Potential-Free (Truly Primal, Mirrorless) Gradient Scheme

$$\dot{x}(t) = -H(x(t))^{-1} \nabla f(x(\lfloor t \rfloor_\eta))$$

$$x_k = x(\eta k)$$

If $H(x)$ is a Hessian map: $H(x) = \nabla^2 \Psi(x)$:

$$\begin{aligned} \text{prox}_{\Psi/\eta}(g, y) &\stackrel{\text{def}}{=} \arg \min_x \eta \langle g, x \rangle + D\Psi(x; y) \\ &= \nabla \Psi^{-1}(\nabla \Psi(y) - \eta g) \\ &= \text{flow}_{(\nabla^2 \Psi)/\eta}(g, y) \end{aligned}$$

$$\begin{aligned} x_{k+1} &\leftarrow \text{prox}_{\Psi/\eta}(\nabla f(x_k), x_k) \\ &= \text{flow}_{(\nabla^2 \Psi)/\eta}(\nabla f(x_k), x_k) \end{aligned}$$

For any smooth $H(x) \succ 0$:

$$\begin{aligned} \text{flow}_{H/\eta}(g, y) &\stackrel{\text{def}}{=} x(\eta) \text{ where } x \text{ is the solution to} \\ &\dot{x}(\tau) = -H(x(\tau))^{-1} g \quad x(0) = y \end{aligned}$$

$$x_{k+1} \leftarrow \text{flow}_{H/\eta}(\nabla f(x_k), x_k)$$

Relative smoothness (and strong convexity):

$$\lambda H(x) \preceq \nabla^2 f(x) \preceq LH(x)$$

If H is a Hessian map: $k \leq \frac{L}{\lambda} \log \left(\frac{\epsilon_0}{\epsilon_k} \right)$

????? $k \leq \frac{L}{\lambda} \log \left(\frac{\epsilon_0}{\epsilon_k} \right)$??????

Change of parameterization:

$$\tilde{x} = \phi(x)$$

$$\tilde{f}(\tilde{x}) = f(\phi^{-1}(\tilde{x}))$$

$$\tilde{H}(\tilde{x}) = \nabla\phi^{-1}(\tilde{x})H(\phi^{-1}(\tilde{x}))\nabla\phi^{-1}(\tilde{x})$$

$$\min_x f(x)$$

$$\min_{\tilde{x}} \tilde{f}(\tilde{x})$$

$$(MD) \quad \dot{x}(t) = -H(x(t))^{-1} \nabla f(x(\lfloor t \rfloor_\eta))$$

$$(\widetilde{MD}) \quad \dot{\tilde{x}}(t) = -\tilde{H}(\tilde{x}(t))^{-1} \nabla \tilde{f}(\tilde{x}(\lfloor t \rfloor_\eta))$$

Is the solution invariant? Do the solutions of (MD) and (\widetilde{MD}) satisfy $\tilde{x}(t) = \phi(x(t))$?

(It is invariant when $\eta = 0$)

Let $x(t)$ be the solution of (MD) , set $\tilde{x}(t) \doteq \phi(x(t))$ and check if (\widetilde{MD}) is satisfied:

$$\dot{\tilde{x}}(t) = -\tilde{H}(\tilde{x}(t))^{-1} \left(\nabla\phi^{-1}(\tilde{x}(\lfloor t \rfloor_\eta)) \nabla\phi(\tilde{x}(t)) \right) \nabla\tilde{f}(\tilde{x}(\lfloor t \rfloor_\eta))$$

Mirrorless Mirror Descent

- Primal-Only, Potential-Free, derivation of Primal Gradient Scheme, as partial discretization of Riemannian Gradient Flow:

$$\begin{aligned}\dot{x}(t) &= -H(x(t))^{-1} \nabla f(x(\lfloor t \rfloor_\eta)) \\ x^{(k)} &= x(\eta k)\end{aligned}$$

- Valid first-order method for any metric tensor $H(x)$
- Does relative smoothness guarantee hold even if H is not a Hessian map?

$$\lambda H(x) \preceq \nabla^2 f(x) \preceq LH(x) \quad \rightarrow \quad k \leq \frac{L}{\lambda} \log \left(\frac{\epsilon_0}{\epsilon_k} \right) \quad \text{????}$$

<https://arxiv.org/abs/2004.01025>

An Even More Optimal Stochastic Optimization Algorithm: Interpolation and Parallellization

Blake Woodworth (TTIC→ENS), Nati Srebro (TTIC)



$$\min F(x)$$

$$F(x) = \mathbb{E}_{z \sim \mathcal{D}}[f(x, z)]$$

- Accelerated SGD: $f(x_k) - f(x^*) \leq O\left(\frac{L\|x^*\|^2}{k^2} + \frac{\sigma\|x^*\|}{\sqrt{k}}\right) \quad \forall_x \mathbb{E}[\|\nabla f(x, z_i) - \nabla F(x)\|^2] \leq \sigma^2 \quad \nabla^2 f \preceq L$
- SGD (Without acceleration): $f(x_k) - f(x^*) \leq O\left(\frac{L\|x^*\|^2}{k} + \frac{\sigma_*\|x^*\|}{\sqrt{k}}\right) \quad \mathbb{E}[\|\nabla f(x^*, z_i) - \nabla F(x^*)\|^2] \leq \sigma_*^2$

- Is it possible to get: $\frac{L\|x^*\|^2}{k^2} + \frac{\sigma_*\|x^*\|}{\sqrt{k}} ?$

Interpolation (realizable) learning:
 $f(x, z) \geq 0, F(x^*) = 0 \rightarrow \sigma_* = 0$

- Answer: **No!**

- For any stochastic optimization method (even if not first order) that uses k samples $z_1, \dots, z_k \sim iid \mathcal{D}$, there exists a problem with $f(x_k) - f(x^*) \geq \Omega\left(\frac{L\|x^*\|^2}{k} + \frac{\sigma_*\|x^*\|}{\sqrt{k}}\right)$

- But using k mini-batch grad estimates $g_i = \frac{1}{b} \sum_{j=1}^b \nabla f(x_i, z_{i,j})$,

accelerated MB-SGD ensures: $f(x_k) - f(x^*) \leq O\left(\frac{L\|x^*\|^2}{k^2} + \frac{L\|x^*\|^2}{bk} + \frac{\sigma_*\|x^*\|}{\sqrt{bk}}\right)$

...and this is optimal (in terms of L, σ_*, b , for methods using mini-batch grad estimates)

Linear parallel speedup in broad regime (even when $\sigma_* = 0$)

When $\sigma_* = 0$, parallelizing on $b = k$ machines yields optimal parallel runtime