

- Nesterov's acceleration
 - Convergence for L -smooth convex
$$o\left(\frac{1}{\sqrt{\epsilon}}\right)$$
 - Interpretation via second-order ODE

Nestrov's method :

Yuri Nestrov 1983

$$y_{t+1} = x_t - \eta \nabla f(x_t)$$

$$x_{t+1} = y_{t+1} + \frac{t}{t+3} (y_{t+1} - y_t)$$

$$y_{t+1} = x_t - \eta_t \nabla f(x_t) \quad : \text{normal}$$

$$z_{t+1} = z_t - \eta_t \cdot \frac{t+1}{2} \nabla f(x_t) \quad : \text{aggressive}$$

$$x_{t+1} = \frac{t+1}{t+3} y_{t+1} + \frac{2}{t+3} z_{t+1} \quad : \text{average}$$

lower weight

Recall :

For $f : \mathbb{R}^d \rightarrow \mathbb{R}$ that is L smooth and convex,
gradient descent yields

$$f(\underline{x}_T) - f(\underline{x}^*) \leq \frac{L}{2T} \|\underline{x}_0 - \underline{x}^*\|^2, \quad T \geq 0$$

Iteration complexity : $O\left(\frac{1}{\epsilon}\right)$ $\epsilon \sim 10^{-6}$
 $T_{GD} \geq 10^6$

What does Nesterov's accelerated gradient yield?

Iteration complexity of $O\left(\frac{1}{\sqrt{\epsilon}}\right)$

$$T_{NAG} \geq 10^3$$

Let us define the energy function (potential or Lyapunov function) and assign to each time t :

$$\phi(t) = t(t+1) (f(\underline{y}_t) - f(\underline{x}^*)) + 2L \|\underline{z}_t - \underline{x}^*\|_2^2$$

aggressive step

If $\phi(t+1) \leq \phi(t)$, for each t , we have

$$\underbrace{T(T+1) (f(\underline{y}_T) - f(\underline{x}^*)) + 2L \|\underline{z}_T - \underline{x}^*\|_2^2}_{\phi(T)} \leq \underbrace{2L \|\underline{z}_0 - \underline{x}^*\|_2^2}_{\phi(0)}$$

$$\Rightarrow f(\underline{y}_T) - f(\underline{x}^*) \leq \frac{2L \|\underline{z}_0 - \underline{x}^*\|_2^2}{T(T+1)}$$

$$\Rightarrow T \approx O\left(\frac{1}{\sqrt{\epsilon}}\right)$$

Recall from the vanilla analysis and for L -Smooth functions:

$$\text{with } \eta_t = \eta = \frac{t+1}{2L} \quad \text{and} \quad \underline{g}_t = \nabla f(\underline{x}_t)$$

$$\underline{z}_{t+1} = \underline{z}_t - \eta_t \cdot \frac{t+1}{2} \nabla f(\underline{x}_t) : \text{aggressive}$$

$$\bullet \quad \underline{g}_t^T (\underline{z}_t - \underline{x}^*) = \frac{t+1}{4L} \|\underline{g}_t\|^2 + \frac{L}{t+1} (\|\underline{z}_t - \underline{x}^*\|^2 - \|\underline{z}_{t+1} - \underline{x}^*\|^2)$$

$$\underline{y}_{t+1} = \underline{x}_t - \eta_t \nabla f(\underline{x}_t) : \text{normal}$$

$$\bullet \quad f(\underline{y}_{t+1}) \leq f(\underline{x}_t) - \frac{1}{2L} \|\underline{g}_t\|_2^2 ; \quad \eta = \frac{1}{L}$$

$$\bullet \quad \underline{\text{Convexity}} : \quad f(\underline{x}_t) - f(\underline{w}) \leq \underline{g}_t^T (\underline{x}_t - \underline{w})$$

$$\Delta \leq t \left[\phi(\underline{x}_t) - \phi(\underline{y}_t) \right] + 2 \left[\phi(\underline{x}_t) - \phi(\underline{x}^*) \right] \\ - \frac{1}{2L} \|\underline{g}_t\|^2 - 2\underline{g}_t^\top (\underline{z}_t - \underline{x}^*)$$

$$\leq t \left[\phi(\underline{x}_t) - \phi(\underline{y}_t) \right] + 2 \left[\phi(\underline{x}_t) - \phi(\underline{x}^*) \right] \\ - 2\underline{g}_t^\top (\underline{z}_t - \underline{x}^*)$$

$$\leq t \underline{g}_t^\top (\underline{x}_t - \underline{y}_t) + 2 \underline{g}_t^\top (\underline{x}_t - \underline{x}^*) \\ - 2 \underline{g}_t^\top (\underline{z}_t - \underline{x}^*)$$

$$= \underline{g}_t^\top \left[(t+2) \underline{x}_t - t \underline{y}_t - 2 \underline{z}_t \right]$$

$$= 0 \iff \underline{x}_{t+1} = \frac{t+1}{t+3} \underline{y}_{t+1} + \frac{2}{t+3} \underline{z}_{t+1} \quad : \text{average}$$

$$\Rightarrow \Delta \leq 0 \Rightarrow \phi(t+1) \leq \phi(t) \quad \text{for each } t. \quad \square$$

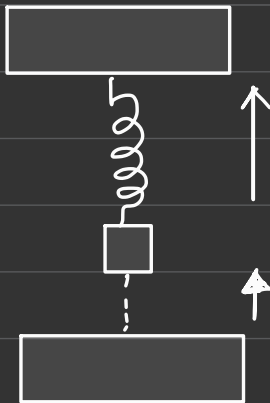
Interpretation using differential equations:

Second-order ODE:

$$\ddot{y}(\tau) + \frac{2}{\tau} \dot{y}(\tau) + \nabla \phi(y(\tau)) = 0$$

damping
co-efficient

Lyapunov



$F = -\nabla \phi(x)$ (Spring force)

$-\frac{\tau}{3} \dot{x}$ (time-varying damping)

$$\frac{t}{t+3} = 1 - \frac{3}{t+3}$$

$$y_{t+1} = x_t - \eta \nabla \phi(x_t)$$

$$x_{t+1} = y_{t+1} + \frac{t}{t+3} (y_{t+1} - y_t)$$

$$y_{t+1} = y_t + \frac{t-1}{t+2} (y_t - y_{t-1}) - \eta \nabla f(x_t)$$

$$\Rightarrow \frac{y_{t+1} - y_t}{\sqrt{\eta}} = \frac{t-1}{t+2} \frac{y_t - y_{t-1}}{\sqrt{\eta}} - \sqrt{\eta} \nabla f(x_t)$$

Let $t = \tau / \sqrt{\eta}$, $Y(\tau) \approx y_{\tau / \sqrt{\eta}} = \underline{y}_t$

and $Y(\tau + \sqrt{\eta}) \approx \underline{y}_{t+1}$

Then, using Taylor expansion:

$$\frac{y_{t+1} - y_t}{\sqrt{\eta}} = (Y(\tau + \sqrt{\eta}) - Y(\tau)) \cdot \frac{1}{\sqrt{\eta}} \approx \dot{Y}(\tau) + \frac{1}{2} \ddot{Y}(\tau) \sqrt{\eta}$$

Similarly, $\frac{y_t - y_{t-1}}{\sqrt{\eta}} \approx \dot{Y}(\tau) - \frac{1}{2} \ddot{Y}(\tau) \sqrt{\eta}$

So Newton's acceleration

$$\dot{\gamma}(\tau) + \frac{\sqrt{\eta}}{2} \ddot{\gamma}(\tau) \approx \left(1 - \frac{3\sqrt{\eta}}{2}\right) \left[\dot{\gamma}(\tau) - \frac{\sqrt{\eta}}{2} \ddot{\gamma}(\tau) \right]$$

$$- \sqrt{\eta} \nabla f(\gamma(\tau))$$

$$\Rightarrow \ddot{\gamma}(\tau) + \frac{3}{2} \dot{\gamma}(\tau) + \nabla f(\gamma(\tau)) \approx 0$$

For this second-order ODE

$$f(\gamma(\tau)) - f_{\text{opt}} \leq O\left(\frac{1}{\tau^2}\right)$$

Actually, 3 is the smallest constant that guarantees $O\left(\frac{1}{\tau^2}\right)$

Let f is L -smooth and μ -strongly convex, then
Nesterov's accelerated gradient satisfies

$$f(\underline{y}_T) - f(\underline{x}^*) \leq \frac{L + \mu}{2} \exp\left(-\frac{T}{\sqrt{\kappa}}\right) \|\underline{x}_0 - \underline{x}^*\|_2^2$$

with

$$t = \frac{3}{2} (\sqrt{\kappa} - 1)$$

$$\kappa = \frac{L}{\mu}$$

$$\eta = \frac{1}{L}$$