- Gradient descent for unconstrained problems

| Lipschitz | $O\left(\dfrac{1}{\varepsilon^2}\right)$ $\varepsilon = 10^{-6}$ $\uparrow\uparrow 10^{12}$ !! |
|---|---|
| Smooth | $O\left(\dfrac{1}{\varepsilon}\right)$ $\uparrow\uparrow 10^{6}$ |
| Smooth & Strongly convex | $O\left(\log\left(\dfrac{1}{\varepsilon}\right)\right)$ $\uparrow\uparrow 6$ |

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

- Quadratic function

  – Exact and backtracking line search

- **Smooth** convex **functions** : $O\left(\frac{1}{\varepsilon}\right)$ : Sublinear convergence

$\rightarrow \quad f(y) \leq f(\underline{x}) + \nabla f^T(\underline{x})(y - \underline{x}) + \frac{L}{2} \| \underline{x} - \underline{y} \|^2$

$$\forall \underline{x}, \underline{y} \in dom \ f$$
$$= X$$

$\longrightarrow$ A bound on optimality gap : $f(\underline{x}) - f^*$

$$f^* = f(\underline{x}^*) \quad \text{where} \quad \underline{x}^* \text{ is a solution to } min. \ f(\underline{x})$$

$$\frac{1}{2L} \| \nabla f(\underline{x}) \|_2^2 \overset{(a)}{\leq} f(\underline{x}) - f^* \overset{(b)}{\leq} \frac{L}{2} \| \underline{x} - \underline{x}^* \|_2^2$$

**Gradient descent** : $\quad \underline{x}_{t+1} = \underline{x}_t - \frac{1}{L} \nabla f(\underline{x}_t)$

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

$$\implies \quad \underline{x}_{t+1} - \underline{x}_t = -\frac{1}{L} \nabla f(\underline{x}_t)$$

$$f(\underline{x}_{t+1}) \leq f(\underline{x}_t) - \frac{1}{L} \| \nabla f(\underline{x}_t) \|_2^2 + \frac{1}{2L} \| \nabla f(\underline{x}_t) \|_2^2$$

$$= f(\underline{x}_t) - \frac{1}{2L} \| \nabla f(\underline{x}_t) \|_2^2$$

$$\Rightarrow \quad \frac{1}{2L} \sum_{t=0}^{T-1} \| \nabla f(\underline{x}_t) \|_2^2 \leq \sum_{t=0}^{T-1} \left( f(\underline{x}_t) - f(\underline{x}_{t+1}) \right)$$

$$= f(\underline{x}_0) - f(\underline{x}_T)$$

(telescopic sum)

Recall:

$$\sum_{t=0}^{T-1} \left( f(\underline{x}_t) - f(x^*) \right) \leq \frac{\eta}{2} \sum_{t=0}^{T-1} \| \underline{g}_t \|_2^2 + \frac{1}{2\eta} \| \underline{x}_0 - \underline{x}^* \|^2$$

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

$$\sum_{t=0}^{T-1} \left( f(\underline{x}_t) - f(\underline{x}^*) \right) \leq \frac{1}{2L} \sum_{t=0}^{T-1} \| \nabla f(\underline{x}_t) \|_2^2$$

$$+ \frac{L}{2} \| \underline{x}_0 - x^* \|^2$$

$$\leq f(x_0) - f(\underline{x}_T) + \frac{L}{2} \| \underline{x}_0 - \underline{x}^* \|^2$$

$$\Rightarrow \quad \sum_{t=1}^{T} \left( f(\underline{x}_t) - f(\underline{x}^*) \right) \leq \frac{L}{2} \| \underline{x}_0 - \underline{x}^* \|_2^2$$

Since $\quad f(x_{t+1}) \leq f(x_t) \qquad \forall\, t \in [0, T]$

$$\frac{1}{T} \sum_{t=1}^{T} f(x_t) - f(x^*) = \left( \frac{1}{T} \sum_{t=1}^{T} f(x_t) \right) - f(x^*)$$

$$\geq f(x_T) - f(x^*)$$

$$\Rightarrow \quad f(x_T) - f(x^*) \leq \frac{1}{T} \sum_{t=1}^{T} f(x_t) - f(x^*)$$

$$\leq \frac{L}{2T} \| x_0 - x^* \|_2^2 \quad ; T > 0$$

$$\frac{LR^2}{2T} = \varepsilon$$

with $\qquad R^2 = \| x_0 - x^* \|_2^2$

to obtain $\qquad \min_{t=0 \,..\, T-1} f(x_t) - f(x^*) \leq \varepsilon$

we need $\qquad T \geq \dfrac{R^2 L}{2 \varepsilon}$

Previously:

$$T \geq \frac{R^2 B^2}{\varepsilon^2}$$
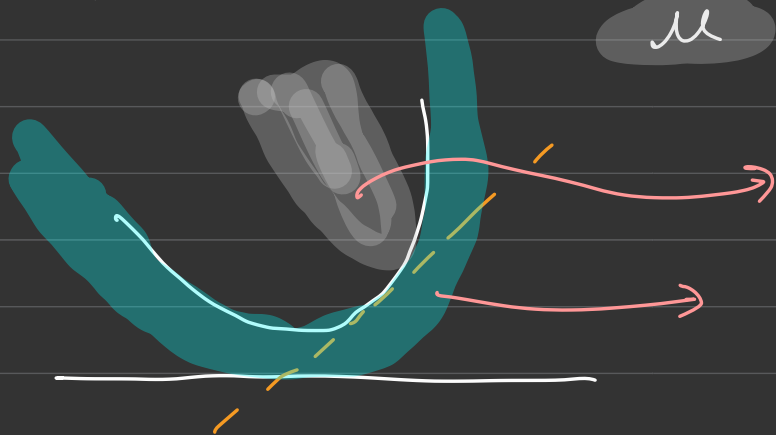
# L - Smooth and μ - Strongly convex functions :

A function $f : \mathbb{R}^n \to \mathbb{R}$ is

μ - Strongly convex and L - Smooth

if

$$\frac{\mu}{2} \| \underline{x} - \underline{y} \|_2^2 \leq f(\underline{y}) - f(\underline{x}) - \nabla f(\underline{x})^T (\underline{y} - \underline{x}) \leq \frac{L}{2} \| \underline{x} - \underline{y} \|_2^2$$

Define $\kappa = \dfrac{L}{\mu}$ is the condition number



L - Smooth

μ - Strong convex

$\underline{x}^*$ is the minimizer

Gradient descent with a fixed step size

$$\underline{x}_{t+1} = \underline{x}_t - \frac{1}{L} \nabla f(\underline{x}_t)$$

Start with arbitrary $\underline{x}_0 \in \mathbb{R}^n$

claim.

ⓐ Squared distances to $\underline{x}^*$ are geometrically decreasing

$$\| \underline{x}_{t+1} - \underline{x}^* \|^2 \leq \left( 1 - \frac{\mu}{L} \right) \| \underline{x}_t - \underline{x}^* \|^2, \quad t \geq 0$$

$$\leq \left( 1 - \frac{\mu}{L} \right)^t \| \underline{x}_0 - \underline{x}^* \|^2$$

ⓑ The error after $T$ iterations is exponentially small in $T$ :

$$f(\underline{x}_T) - f(\underline{x}^*) \leq \frac{L}{2} \left( 1 - \frac{\mu}{L} \right)^T \| \underline{x}_0 - \underline{x}^* \|^2 ; \quad T > 0$$

(a)

Recall $g_t = \nabla f(x_t)$

$$g_t^T(x_t - x^*) = \nabla f^T(x_t)(x_t - x^*)$$

(from $\mu$-strong convexity)

$$\geq f(x_t) - f(x^*) + \frac{\mu}{2}\|x_t - x^*\|_2^2$$

From vanilla analysis:

$$g_t^T(x_t - x^*) = \frac{\eta}{2}\|g_t\|^2 + \frac{1}{2\eta}\left[\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2\right]$$

$$\Rightarrow f(x_t) - f(x^*) \leq \frac{1}{2\eta}\left[\eta^2\|g_t\|^2 + \|x_t - x^*\|_2^2 - \|x_{t+1} - x^*\|_2^2\right]$$

$$- \frac{\mu}{2}\|x_t - x^*\|_2^2$$

We have a bound on $\|x_{t+1} - x^*\|_2^2$ :

$$\|x_{t+1} - x^*\|_2^2 \leq 2\eta\left[f(x_t) - f(x^*)\right] + \eta^2\|g_t\|^2 + (1 - \mu\eta)\|x_t - x^*\|_2^2$$

$\Downarrow$

this disappear as shown next

For L-smooth convex functions; for $\eta = \frac{1}{L}$ :

$f(\underline{x}^*) - f(\underline{x}_t) \leq f(\underline{x}_{t+1}) - f(\underline{x}_t) \leq \frac{-1}{2L} \|\nabla f(\underline{x}_t)\|_2^2$

$2\eta \left[ f(\underline{x}^*) - f(\underline{x}_t) \right] + \eta^2 \|\nabla f(\underline{x}_t)\|_2^2 \leq 0$

$\Rightarrow \| \underline{x}_{t+1} - \underline{x}^* \|_2^2 \leq \left( 1 - \mu\eta \right) \| \underline{x}_t - \underline{x}^* \|_2^2$

$= \left( 1 - \frac{\mu}{L} \right) \| \underline{x}_t - \underline{x}^* \|_2^2$

$\| \underline{x}_T - \underline{x}^* \|_2^2 \leq \left( 1 - \frac{\mu}{L} \right)^T \| \underline{x}_0 - \underline{x}^* \|_2^2$

## (b)

from smoothness:

$$f(\underline{x}_T) - f(\underline{x}^*) \leq \nabla f^T(\underline{x}^*) \underbrace{(\underline{x}_T - \underline{x}^*)}_{\nabla f \cdot (\underline{x}^*) = 0} + \frac{L}{2} \|\underline{x}_T - \underline{x}^*\|_2^2$$

$$= \frac{L}{2} \| \underline{x}_T - \underline{x}^* \|_2^2$$

$$\underset{(b)}{\leq} \frac{L}{2} \left( 1 - \frac{\mu}{L} \right)^T \underbrace{\| \underline{x}_0 - \underline{x}^* \|_2^2}_{R^2}$$

To find the number of iterations:

$$\frac{L}{2} \left( 1 - \frac{\mu}{L} \right)^T R^2 = \varepsilon \implies \left( 1 - \frac{\mu}{L} \right)^T = \frac{2\varepsilon}{R^2 L}$$

$$\implies T \ln \left( 1 - \frac{\mu}{L} \right) = \ln \left( \frac{2\varepsilon}{R^2 L} \right)$$

since $\ln(1+x) \leq x$

$$T \left( -\frac{\mu}{L} \right) \leq \ln \left( \frac{2\varepsilon}{R^2 L} \right) \implies T \geq \frac{L}{\mu} \ln \left( \frac{R^2 L}{2\varepsilon} \right)$$

# Summary:

Gradient descent with fixed step size

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

$\varepsilon = 10^{-6}$

| | | |
|---|---|---|
| Lipschitz | $O\left(\dfrac{1}{\varepsilon^2}\right)$ | $10^{12}$ |
| Smooth | $O\left(\dfrac{1}{\varepsilon}\right)$ | $10^6$ |
| Smooth & Strongly convex | $O\left(\log\left(\dfrac{1}{\varepsilon}\right)\right)$ | 6 |

$$\log(10^6)$$

# A similar result:

Suppose $f$ is $\mu$-strongly convex and $L$-smooth.
Then gradient descent with $\eta_t = \eta = \dfrac{2}{\mu + L}$
satisfies

a. $\quad \| x_T - x^* \|_2^2 \leq \left( \dfrac{K-1}{K+1} \right)^{2T} \| x_0 - x^* \|_2^2$

b. $\quad f(x_T) - f(x^*) \leq \dfrac{L}{2} \left( \dfrac{K-1}{K+1} \right)^{2T} \| x_0 - x^* \|_2^2$

Homework 2.

# Example: Quadratic minimization:

$$\text{minimize}_{\underline{x}} \quad f(\underline{x}) = \frac{1}{2}(\underline{x} - \underline{x}^*)^T Q (\underline{x} - \underline{x}^*)$$

$$Q > 0 \qquad \nabla f(\underline{x}) = Q(\underline{x} - \underline{x}^*)$$

$$\underline{x}_{t+1} - \underline{x}^* = \underline{x}_t - \underline{x}^* - \eta_t \nabla f(\underline{x}_t)$$

$$= (I - \eta_t Q)(\underline{x}_t - \underline{x}^*)$$

We have

$$\|\underline{x}_{t+1} - \underline{x}^*\|_2 \leq \|I - \eta_t Q\| \, \|\underline{x}_t - \underline{x}^*\|$$

$$\|I - \eta_t Q\| = \max\left\{ |1 - \eta_t \lambda_1(Q)|, \; |1 - \eta_t \lambda_n(Q)| \right\}$$

$\eta$ that yields $\quad |1 - \eta_t \lambda_1(Q)| = |1 - \eta_t \lambda_n(Q)|$

$$\Rightarrow \quad \eta = \frac{2}{\lambda_1(Q) + \lambda_n(Q)}$$

So

$$\|I - \eta \cdot Q\| = 1 - \frac{2\lambda_n(Q)}{\lambda_1(Q) + \lambda_n(Q)} = \frac{\lambda_1(Q) - \lambda_n(Q)}{\lambda_1(Q) + \lambda_n(Q)}$$

$$\Rightarrow \quad \|x_t - x^*\|_2 \leq \left(\frac{\lambda_1(Q) - \lambda_n(Q)}{\lambda_1(Q) + \lambda_n(Q)}\right) \|x_t - x^*\|_2$$

$$= \left(\frac{\lambda_1(Q) - \lambda_n(Q)}{\lambda_1(Q) + \lambda_n(Q)}\right)^t \|x_0 - x^*\|_2$$

# Exact line Search:

$$\eta_t = \underset{\eta \geq 0}{\arg\min} \; f\left(\underline{x}_t - \eta \nabla f(\underline{x}_t)\right)$$
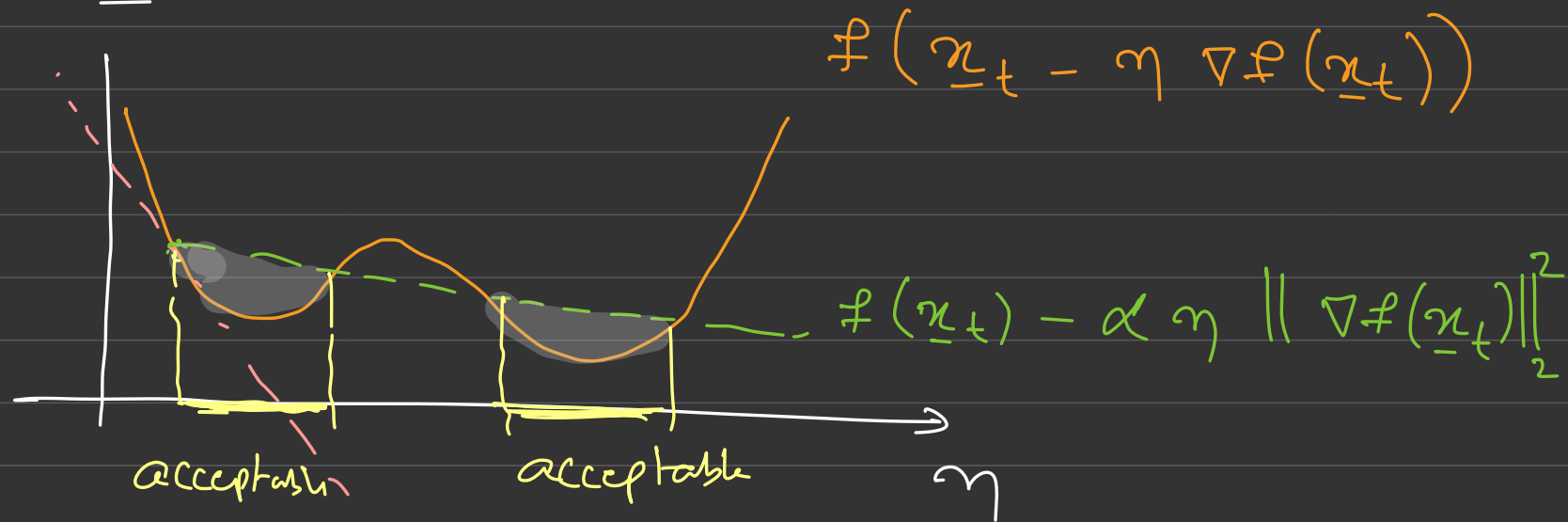
$$\eta_t = \frac{\underline{g}_t^\top \underline{g}_t}{\underline{g}_t^\top \underline{8} \, \underline{g}_t}$$

$$f(\underline{x}_t) - f(\underline{x}^*) \leq \left(\frac{\lambda_1(\underline{8}) - \lambda_n(\underline{8})}{\lambda_1(\underline{8}) - \lambda_n(\underline{8})}\right)^{2t} \left(f(x_0) - f(x^*)\right)$$

(Homework 2)

- Convergence rate is not faster than fixed step size

# Backtracking line Search:

$$f(\underline{x}_t - \eta \nabla f(\underline{x}_t))$$



$$f(\underline{x}_t) - \alpha \eta \| \nabla f(\underline{x}_t) \|_2^2$$

acceptable     acceptable

$$f(\underline{x}_t) - \eta \| \nabla f(\underline{x}_t) \|^2$$

**Armijo Condition:** Ensures sufficient decrease in the objective value

$$0 < \alpha < 1$$

$$f(\underline{x}_t - \eta \nabla f(\underline{x}_t)) < f(\underline{x}_t) - \alpha \eta \| \nabla f(\underline{x}_t) \|_2^2$$

- $\eta = 1$, $0 < \alpha \leq \frac{1}{2}$, $0 \leq \beta < 1$

  while $f(\underline{x}_t - \eta \nabla f(\underline{x}_t)) > f(\underline{x}_t) - \alpha \eta \| \nabla f(\underline{x}_t) \|_2^2$

  $$\eta \leftarrow \beta \eta$$

$f$ is $\mu$ - strongly convex and $L$ - smooth:

$$f(\underline{x}_t) - f(\underline{x}^*) \leq \left(1 - \min\left\{2\mu\alpha, \frac{2\beta\alpha\mu}{L}\right\}\right)^t \left(f(\underline{x}_0) - f(\underline{x}^*)\right)$$