

Last time:

- Grad Descent
- Convergence in a simple $c \times c$ setting

Today: More structure

- Convergence Rates of GD for:
 - strongly convex fun
 - smooth
 - nonconvex smooth
- What do these mean for practical setups.

Reminder:

str. convexity: $f(x)$ is a str-cvx
if $f(x) - \frac{\gamma}{2}\|x\|^2$ is cvx

Lm: If $f(\cdot)$ is γ -str convex and β -smooth. Then:

$$(\nabla f(x) - \nabla f(y))^T(x-y) \geq \frac{\gamma\beta}{\gamma+\beta} \|x-y\|^2 + \frac{1}{\beta+\gamma} \|\nabla f(x) - \nabla f(y)\|^2$$

Cor: $\langle \nabla f(x), x-x^* \rangle \geq c\|x-x^*\|^2 + c'\|\nabla f(x)\|^2$
 (Strong correlation towards opt)
 "attraction"

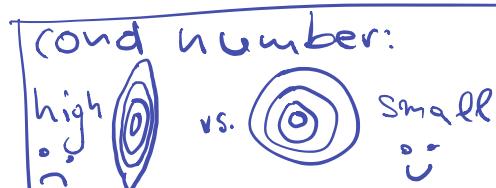
Theorem:

Let f be β -smooth and γ -strongly convex. The Grad. Descent

with $\gamma = \frac{2}{\gamma+\beta}$ obtains

$$\|x_t - x^*\|^2 \leq e^{-2t/\kappa} \|x_0 - x^*\|^2$$

where $\kappa = \beta/\gamma$



"stretch of function"

Proof:

$$\begin{aligned}
 \|x_{k+1} - x^*\|^2 &= \|x_k - \gamma \nabla f(x_k) - x^*\|^2 \\
 \Delta_{k+1} &= \Delta_k - 2\gamma \langle \nabla f(x_k), x_k - x^* \rangle + \gamma^2 \|\nabla f(x_k)\|^2 \\
 &\quad \text{apply co-coercivity}
 \end{aligned}$$

$$\begin{aligned}
 \gamma &= \frac{2}{\alpha + \beta} \\
 &= \|x_k - x^*\| - \frac{4\beta}{(\alpha + \beta)^2} \|x_k - x^*\|^2 \\
 &\quad + \left[\frac{4}{(\alpha + \beta)^2} - \frac{4}{(\alpha + \beta)^2} \right] \|\nabla f(x_k)\|^2
 \end{aligned}$$

$$\begin{aligned}
 &= \left(1 - \frac{4\beta}{(\alpha + \beta)^2} \right) \|x_k - x^*\|^2 \\
 &= \left(\frac{\alpha^2 + 2\beta + \beta^2 - 4\beta}{(\alpha + \beta)^2} \right) \|x_k - x^*\|^2
 \end{aligned}$$

$$\begin{aligned}
 &= \left(\frac{\gamma - \beta}{\gamma + \beta} \right)^2 \|x_t - x^*\|^2 = \left(\frac{1 - \beta/\gamma}{1 + \beta/\gamma} \right)^2 \|x_t - x^*\|^2 \\
 &= \left(\frac{k-1}{k+1} \right)^2 \|x_t - x^*\|^2 = \left(\frac{k-1}{k+1} \right)^{2t} \|x_0 - x^*\|^2 \\
 &= e^{2t \log(1 - \frac{1}{k})} \|x_0 - x^*\|^2 \leq e^{-2t/k} \|x_0 - x^*\|^2
 \end{aligned}$$

□

Comparison of Conv. Rates:

<u>L-Lip:</u> γ -Str. Conv. + B -Smooth	$\frac{R \cdot L}{T}$ $B R e^{-2 \frac{T}{B}}$	}	Structure helps us analyze convergence!
<u>B-Smooth:</u> Str-Conv + L-Lip	$\frac{B \cdot R^2}{T}$ $\frac{L^2}{\gamma \cdot T}$		

Remark:
Not always
tight.

Q: What do these "constants" look like for some real problems?

Computational Complexity of GD:

- Let one ∇f eval be our unit of comp. cost.
- Goal: Measure the total cost wrt $\nabla f(x)$ evals.
- Total cost: $O([\# \nabla f \text{ evals}] \cdot [\text{#iter}(\epsilon)])$

Let's see an example:

Log. reg. + regularization

$$f(w) = \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i \langle x_i, w \rangle}) + \frac{\lambda}{2} \|w\|^2$$

- Reminder:
- $\log(1 + e^x)$ is 1-Lip $\frac{1}{4}$ -Smooth
 - $g_1(g_2(x))$ is $L g_1 L g_2$ -Lip
 - $x^T w + b$ is $\|x\| -$ Lip
 - $g(x^T w + b)$ is $\beta \|x\|^2$ -smooth

$$\frac{1}{n} \sum_i \log(1 + e^{-g_i(x_i, w)})$$

is • $\frac{1}{n} \sum_i \|x_i\| - \text{Lip}$ (L)

• $\frac{1}{n} \sum_i \frac{1}{4} \|x_i\|^2 - \text{smooth}$ (B.)

$\frac{1}{2}\|w\|^2$ • is ↗ str. CVX

• not Lip!

• 2-smooth (B₂)

Assm Let $\|x_i\|_2 = O(d)$, $\|\omega^*\| = O(d)$

Then, $R = \|\omega_0 - \omega^*\| = O(\sqrt{d})$, $L = O(\sqrt{d})$

$B = O(d)$, $\lambda = O(1)$

Rates: • If $\lambda = 0$. $\frac{\beta R^2}{T} = \frac{d^2}{T}$

• If $\lambda \neq 0$ $\beta R^2 e^{-2+\frac{\lambda}{B}} = d^2 e^{-c \frac{T}{d}}$

$$\text{For } \varepsilon - \text{error} \Rightarrow T = O\left(\frac{d^2}{\varepsilon}\right) \quad [\lambda=0]$$

or if $\lambda \neq 0$ $T = O\left(d \log\left(\frac{d}{\varepsilon}\right)\right)$

Q: Overall complexity?

Time to compute $\nabla f(x)$ is proportional
to computing $\langle x_i, w \rangle + b$, i.e.,
 $O(n \cdot d)$.

Cor. For $\lambda = O(1)$ and $B = O(d)$

GD takes $O(nd^2 \log d/\varepsilon)$ on
logistic regression.

Cost / iteration: $O(nd)$ it requires
1 pass over the data!

Too expensive! Can we make this $O(d)$?

Answer: S (d). Next time!