

Recall :

GD: $x^{(t+1)} = x^{(t)} - \mu^{(t)} \nabla f(x^{(t)})$

We developed necessary & sufficient conditions for optimality;

- Today, we'd like to build towards making the choice of $\mu^{(t)}$ rigorous.
- Along the way, we will develop an understanding of the rate of convergence of GD.

(we'd like to know how $f(x^{(t)}) - f(x^*)$ behaves as a function of t .)

Def'n: A function $f: \Omega \rightarrow \mathbb{R}$
 $\hookrightarrow \mathbb{C} \mathbb{R}^n$
is L -Lipschitz if $\forall x, y \in \Omega$
 $|f(x) - f(y)| \leq L \|x - y\|$
"distance"
between $f(x)$ & $f(y)$ distance between
 x & y

So how much f changes
between x & y is determined
by how far x & y are from
each other.

Example: $f: \mathbb{R} \rightarrow \mathbb{R}$

$f(x) = |x|$ is Lipschitz

because $|f(x) - f(y)| = ||x| - |y||$

$$\leq |x - y|.$$

Moreover the Lipschitz constant is 1.

Lemma: If $f: \Omega \rightarrow \mathbb{R}$ is L -Lipschitz, convex, and differentiable then $\|\nabla f(x)\| \leq L \quad \forall x \in \Omega$

Proof: For any $x, y \in \Omega$

$$f(x) - f(y) \geq \nabla f(y)^T (x - y)$$

\uparrow convexity

$$L\|x - y\| \geq |f(x) - f(y)| \geq |\nabla f(y)^T (x - y)|$$

\uparrow Lipschitz.

Pick $x = y + \nabla f(y)$

So :

$$L\|\nabla f(y)\| \geq \|\nabla f(y)\|^2$$

$$\Rightarrow \|\nabla f(y)\| \leq L.$$



Theorem : $\|\nabla f(x)\| \leq L \forall x \in \Omega$
then $|f(x) - f(y)| \leq L \|x - y\|$

Proof : By Taylor's theorem :

$$f(x) - f(y) = \nabla f(tx + (1-t)y)^T (x - y)$$

for some $t \in (0, 1)$

$$\begin{aligned} \text{So } |f(x) - f(y)| &= |\nabla f(tx + (1-t)y)^T (x - y)| \\ &\leq \|\nabla f(\quad)\| \|x - y\| \\ &\leq L \|x - y\| \end{aligned}$$

$$\begin{aligned} \text{bec } w^T z &= \|w\| \|z\| \cos \theta \\ &\& \quad |w^T z| \leq \|w\| \|z\| \end{aligned}$$



We are now ready to present our result on the choice of μ & the convergence of GD.

Theorem: Let f be convex, differentiable, L -Lipschitz

& let $\left\{ \begin{array}{l} \|x^{(0)} - x^*\| \leq R \\ \|\nabla f(x)\| \leq L \end{array} \right.$

Choose $\mu = \frac{R}{L\sqrt{t}}$

Conv.
rate
theorem

Then

$$f\left(\frac{1}{t} \sum_{s=0}^{t-1} x^{(s)}\right) - \underbrace{f(x^*)}_{\text{optimal value}} \leq \frac{RL}{\sqrt{t}}$$

gap between $\frac{1}{t} \sum_{s=0}^{t-1} x^{(s)}$ & opt. value

Example: Consider

$$f(x_1, x_2) = \sin(x_1) + x_2$$

and notice that $\nabla f(x_1, x_2) = \begin{pmatrix} \cos x_1 \\ 1 \end{pmatrix}$

$$\Rightarrow \|\nabla f\| \leq \sqrt{\cos^2(x_1) + 1} \leq \sqrt{2}$$

$\Rightarrow f$ is Lipschitz with $L = \sqrt{2}$

So if we run GD for t -steps

with $\mu = \frac{R}{\sqrt{2}\sqrt{t}}$ and we'll get

within $\frac{R\sqrt{2}}{\sqrt{t}}$ of a local min.

~~Proof of theorem:~~

First, note that

$$f(x^*) \geq f(x^{(s)}) + \nabla f(x^{(s)})^T (x^* - x^{(s)})$$

$$\Leftrightarrow f(x^*) - f(x^{(s)}) \geq \nabla f(x^{(s)})^T (x^* - x^{(s)})$$

①

Next, note that

$$x^{(s+1)} = x^{(s)} - \mu \nabla f(x^{(s)})$$

$$\Rightarrow \nabla f(x^{(s)}) = \frac{x^{(s)} - x^{(s+1)}}{\mu}$$

②

① & ② \Rightarrow

$$f(x^{(s)}) - f(x^*) \leq \frac{1}{\mu} \langle x^{(s)} - x^{(s+1)}, x^{(s)} - x^* \rangle$$

③

Fact:

$$a^T b = \frac{\|a\|^2 + \|b\|^2 - \|a-b\|^2}{2}$$

or $\langle a, b \rangle$

Using this fact in (3) :

$$f(x^{(s)}) - f(x^*) \leq \frac{1}{2\mu} \left(\underbrace{\|x^{(s)} - x^*\|^2}_a + \underbrace{\|x^{(s)} - x^{(s+1)}\|^2}_b - \underbrace{\|x^{(s+1)} - x^*\|^2}_{a-b} \right)$$

$$\text{But } x^{(s)} - x^{(s+1)} = \mu \nabla f(x^{(s)})$$

So

$$f(x^{(s)}) - f(x^*) \leq \frac{1}{2\mu} \left(\|x^{(s)} - x^*\|^2 - \|x^{(s+1)} - x^*\|^2 \right) + \frac{\mu}{2} \|\nabla f(x^{(s)})\|^2$$

(4)

Summing from $s=0$ to $t-1$:

$$\begin{aligned}
& \sum_{s=0}^{t-1} (f(x^{(s)}) - f(x^*)) \\
& \leq \frac{1}{2\mu} \left(\underbrace{\|x^{(0)} - x^*\|^2}_{\leq R^2} - \underbrace{\|x^{(t)} - x^*\|^2}_{\geq 0} \right) \\
& \quad + \frac{\mu}{2} \sum_{s=0}^{t-1} \underbrace{\|\nabla f(x^{(s)})\|^2}_{\leq L^2}
\end{aligned}$$

$$\leq \frac{1}{2\mu} \left(R^2 + \frac{\mu}{2} L^2 t \right)$$

dividing by t :

$$\frac{1}{t} \sum_{s=0}^{t-1} f(x^{(s)}) - f(x^*)$$

$$\leq \frac{1}{2\mu t} R^2 + \frac{\mu}{2} L^2$$

But convexity of f then gives

$$f\left(\frac{1}{t} \sum_{s=0}^{t-1} x^{(s)}\right) - f(x^*) \leq \frac{1}{2\mu t} R^2 + \frac{\mu}{2} L^2$$

by plugging in
 $\mu = \frac{R}{L\sqrt{t}}$

$$\leq \frac{RL}{\sqrt{t}} \quad \square$$

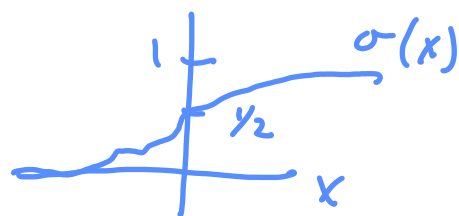
Note about HW2 (don't need to solve)

Problem $\{(x_i, y_i)\}_{i=1}^n$, $y_i \in \{-1, 1\}$

Want to build classifier

Model probability of label w/ sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



$w^T x$ selects + weights features

$$P(y_i | w^T x_i) = \sigma(y_i w^T x_i)$$

Cross entropy loss to learn w looks to

minimize $-E_{(x,y)} \log(p(x,y))$

$$L = -\frac{1}{n} \sum_{i=1}^n \log(\sigma(y_i w^T x_i)) = \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i w^T x_i})$$