

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

$$\leq \frac{R^2 L}{2C\sqrt{t}} + \frac{CL}{2\sqrt{t}} = \frac{1}{\sqrt{t}} (\sim)$$

(Still $\sim \frac{1}{\sqrt{t}}$)

An interpretation of GD

Recall that by Taylor:

$$f(z) = f(x^{(t)}) + \nabla f(x^{(t)})^T (z - x^{(t)})$$

$$+ \frac{1}{2} (z - x^{(t)})^T \nabla^2 f(\xi) (z - x^{(t)})$$

goal
is to
min.
 $f(z)$

ξ is between $x^{(t)}$ & z

GD approximates $\nabla^2 f(\xi)$ by
 $\frac{1}{\mu} I$
Identity matrix

$$f(z) \approx f(x^{(t)}) + \nabla f(x^{(t)})^T (z - x^{(t)}) + \frac{1}{2\mu} \|z - x^{(t)}\|^2$$

$\underbrace{\hspace{15em}}_{g(z)}$

Notice that $g(z)$ is convex
 so to minimize it we want
 to solve

$$\nabla g(z^*) = 0$$

$$\Leftrightarrow \nabla f(x^{(t)}) + \frac{1}{\mu} (z^* - x^{(t)}) = 0$$

$$\Rightarrow z^* = x^{(t)} - \mu \nabla f(x^{(t)})$$

\uparrow
 set $x^{(t+1)}$ to this, and repeat
 \Rightarrow GD!

Another interpretation of GD
& the method of Steepest
descent:

One way we can think of GD,
is that we pick at each
step, $t+1$, $x^{(t+1)}$

so that

$$f(x^{(t+1)}) \approx f(x^{(t)}) - \mu \|\nabla f(x^{(t)})\|_2^2$$

$$\left(\begin{aligned} &\text{bec, in general, } f(z) \approx f(x^{(t)}) - \mu \nabla f(x^{(t)})^T (z - x^{(t)}) \\ &\& \text{ in GD, } x^{(t+1)} = x^{(t)} - \mu \nabla f(x^{(t)}) \end{aligned} \right)$$


Another idea: think of GD as
picking a direction that minimizes
the function $h(\vec{p}) = \nabla f(x^{(t)})^T \frac{\vec{p}}{\|\vec{p}\|}$

In other words, if I think of algorithms that update x via

$$x^{(t+1)} = x^{(t)} - \mu \vec{p}$$

GD is the one that picks \vec{p} so that $h(p)$ is minimized
i.e. we solve

$$\min_P h(p) \Leftrightarrow \min_P \frac{\nabla f(x^{(t)})^T \vec{p}}{\|\vec{p}\|_2}$$

$$\|\vec{p}\|_2 = \sqrt{\sum_{i=1}^N p_i^2}$$


This suggests that we can come up with other algorithms by replacing $\|\vec{p}\|_2$ in the denominator, by other norms.

Def'n Let $x \in \mathbb{R}^N$, then
we define

$$\|x\|_1 := \sum_{i=1}^N |x_i| \rightarrow \begin{matrix} 1\text{-norm} \\ \text{or } \ell_1\text{-norm} \end{matrix}$$

$$\|x\|_2 := \sqrt{\sum_{i=1}^N |x_i|^2} \rightarrow \begin{matrix} 2\text{-norm} \\ \text{or} \\ \ell_2\text{-norm} \end{matrix}$$

$$\|x\|_\infty = \max_i |x_i|$$

more generally for $p \in (1, \infty)$

$$\|x\|_p = \left(\sum_{i=1}^p |x_i|^p \right)^{1/p}$$

Now, consider algorithms of
the type

$$x^{(t+1)} = x^{(t)} - \mu P$$

where P minimizes $\frac{\nabla F(x^{(t)})^T P}{\|P\|_1}$

$$\propto \frac{\nabla F(x^{(t)})^T P}{\|P\|_\infty} \quad (\text{or some other norm})$$

For 1-norm :

$$\min_P \frac{\nabla F(x^{(t)})^T P}{\|P\|_1} \Rightarrow$$

$$P^* = - \text{sign} \left(\frac{\partial F}{\partial x_{j^*}} \right) \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ \delta \\ \vdots \end{pmatrix}$$


\hookrightarrow a vector

j^{th} entry \rightarrow

$$\text{where } \delta = \|\nabla F(x)\|_\infty \\ = \left| \frac{\partial F}{\partial x_{j^*}}(x^{(t)}) \right|$$

where j^* indexes the entry of $\nabla F(x^{(t)})$ with the largest magnitude.

Example : suppose

$$\nabla f(x^{(t)}) = \begin{pmatrix} 1 \\ -2 \\ -0.5 \end{pmatrix}$$


then $P^* = \begin{pmatrix} 0 \\ 2 \\ 0 \end{pmatrix}$

This algorithm is called

"coordinate descent" (move along
one coordinate at each
iteration)

For ∞ -norm :


$$\min_P \frac{\nabla f(x^{(t)})^T P}{\|P\|_\infty}$$

$$\Rightarrow P^* = -\delta \begin{pmatrix} \text{sign}(\nabla f(x^{(t)}))_1 \\ \text{sign}(\nabla f(x^{(t)}))_2 \\ \vdots \\ \text{sign}(\nabla f(x^{(t)}))_n \end{pmatrix}$$

where $\delta = \|\nabla f(x^{(t)})\|_1$

$$\text{So } P^* = -\|\nabla f(x^{(t)})\|_1 \begin{pmatrix} \pm 1 \\ \pm 1 \\ \pm 1 \\ \vdots \\ \pm 1 \end{pmatrix}$$

Example: suppose

$$\nabla f(x^{(t)}) = \begin{pmatrix} 1 \\ -2 \\ -0.5 \end{pmatrix}$$


then $P^* = -3.5 \begin{pmatrix} 1 \\ -1 \\ -1 \end{pmatrix}$.

(GD would have given

$$P^* = \begin{pmatrix} -1 \\ 2 \\ 0.5 \end{pmatrix})$$

Aside: the proof of the optimality of the P^* chosen in the examples above relies on Hölder's inequality:

$$\begin{cases} |x^T y| \leq \|x\|_\infty \|y\|_1 \\ |x^T y| \leq \|x\|_1 \|y\|_\infty \end{cases} \quad \left\{ \begin{array}{l} |x^T y| \leq \|x\|_2 \|y\|_2 \\ |x^T y| \leq \|x\|_2 \|y\|_2 \end{array} \right.$$