

Gradient Descent

(51)

We start by looking at convex functions that are Lipschitz cont.

Def: Let S be an open and convex set. Let $f: S \rightarrow \mathbb{R}$.

We say that f is convex if for all $x, y \in S$ and $\lambda \in [0, 1]$

$$f(z) \leq \lambda f(x) + \bar{\lambda} f(y) \text{ where} \\ z = \lambda x + \bar{\lambda} y.$$

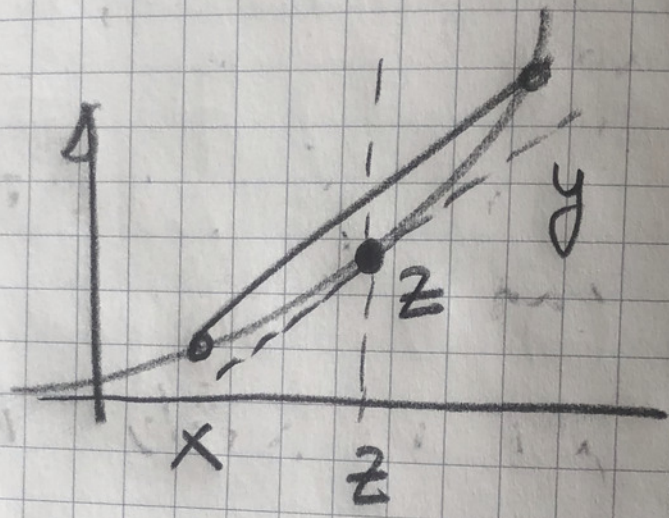
An alternative characterisation is the following.

Lemma: Let S be an open and convex set. Let $f: S \rightarrow \mathbb{R}$.

Then f is convex if $\forall z \in S$ there exists \underline{v} so that $\forall x \in S$

$$f(x) \geq f(z) + \langle \underline{v}, x - z \rangle. (*)$$

0.5.5.6.0



Proof:

(52)

Assume the statement of the Lemma.

Let $z = \lambda x + \bar{\lambda} y$ for some $x, y \in S$
and $\lambda \in [0, 1]$. By the lemma,

$$f(x) \geq f(z) + \langle v_z, x - z \rangle \quad | \quad \lambda$$

$$f(y) \geq f(z) + \langle v_z, y - z \rangle \quad | \quad \bar{\lambda}$$

$\lambda f(x) + \bar{\lambda} f(y) \geq f(z)$, hence convex
according to def.

Conversely, assume def.

If f has gradient at z then clear that
the promised v_z is the gradient of
 f at z and is unique.

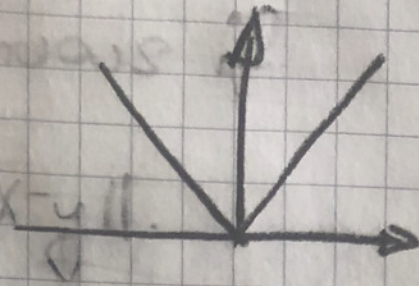
Example: $f: S \rightarrow \mathbb{R}$

(53)

$$f(x) = |x|$$

$$|f(x) - f(y)| \leq \rho \|x - y\|$$

for all $x, y \in S$.



$$\partial f(x) = \begin{cases} \{-1\}, & x < 0 \\ \{1\}, & x > 0 \\ [-1, 1], & x = 0 \end{cases}$$

Example: Let g_1, \dots, g_r be convex differentiable functions. Let

$$g(z) = \max_{i \in \{1, \dots, r\}} g_i(z)$$

Let $D = \text{conv} \{ \nabla g_i(z) \}$. Then

$$\nabla g(z) \subseteq D.$$

Subgradient:

Consider the previous lemma.
Any v_z that fulfills (*) is called a subgradient of f at z . If f is differentiable at z then v_z is unique.

We say $v_z \in \partial f(z)$ and call $\partial f(z)$ the differential set.

We have for all $x \in S$,

(54)

$$g(x) \geq g_j(x)$$

$$\geq g_j(z) + \langle \nabla g_j(z), x-z \rangle$$

$$= g(z) + \langle \nabla g_j(z), x-z \rangle$$



We say that f is

(55)

ρ -Lipschitz if for all $x, y \in S$

$$|f(x) - f(y)| \leq \rho \|x - y\|.$$

Lemma: Let S be an open convex set and let $f: S \rightarrow \mathbb{R}$

be a convex function. Then f

is ρ -Lipschitz iff for all

$z \in S$ and $v_z \in \partial f(z)$, $\|v_z\| \leq \rho$.

Proof: Assume that for all $z \in S$

and $v_z \in \partial f(z)$, $\|v_z\| \leq \rho$.

Then, since f is convex and v_z

is a subgradient of f at z , $\forall x \in S$

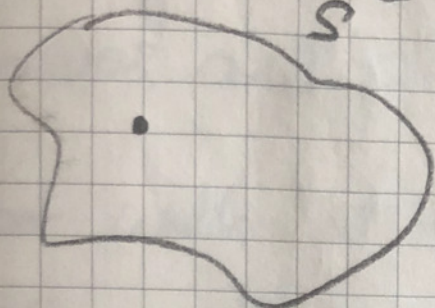
$$f(x) \geq f(z) + \langle v_z, x - z \rangle$$

$$f(z) - f(x) \leq \langle v_z, z - x \rangle$$

$$\leq \rho \|z - x\|$$

We get the second inequality by switching roles of x and z .

For the converse assume that (56)
 $z \in S$ and $v_z \in \partial f(z)$.



$$\text{Let } x = z + \epsilon \frac{v_z}{\|v_z\|}.$$

Then $\|x - z\| = \epsilon$. By Lipschitz

$$|f(x) - f(z)| \leq \rho \|x - z\| = \epsilon \rho.$$

By convexity,

$$f(x) \geq f(z) + \langle v_z, x - z \rangle$$

$$\epsilon \rho \geq f(x) - f(z) \geq \|v_z\| \epsilon$$

$$\text{Hence } \rho \geq \|v_z\|.$$

Let us now look at the 57
GD algorithm applied to a
convex and ρ -Lipschitz function f .

Start with $w^{(0)} = 0$. Then at each
step $w^{(t)} = w^{(t-1)} - \eta \nabla f(w^{(t-1)})$
for some step size η . If the
function f does not have a
gradient at $w^{(t)}$ then pick any
subgradient from $\partial f(w^{(t)})$.

After the step T , output the
estimate

$$\bar{w} = \frac{1}{T} \sum_{t=1}^T w^{(t)},$$

How close will we be to the
optimal value?

Lemma: Apply the

(58)

above algorithm. Let w^* be

$w^* = \arg \min_{w: \|w\| \leq B} f(w)$. Then

$$f(\bar{w}) - f(w^*) \leq \frac{B\rho}{\sqrt{T}}.$$

In other words, we need at most $\frac{B^2 \rho^2}{\epsilon^2}$ steps to get ϵ close.

Here the stepsize is

$$\eta = \sqrt{\frac{B^2}{\rho^2 T}}.$$

Proof:

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

convexity

$$\leq \frac{1}{T} \sum_{t=1}^T (f(w^{(t)}) - f(w^*))$$

convex

$$\leq \frac{1}{T} \sum_{t=1}^T \langle f(w^{(t)} - w^*), \nabla f(w^{(t)}) \rangle$$

see the
next page

$$\leq \frac{1}{T} \left[\frac{\|w^*\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\nabla f(w^{(t)})\|^2 \right]$$

$$\leq \frac{1}{T} \left[\frac{B^2}{2\eta} + \frac{\eta}{2} T \rho^2 \right]$$

$$= \frac{B^2 \rho \sqrt{T}}{2T} + \frac{B T}{2\rho \sqrt{T} T} \rho$$

$$= \frac{\rho B}{2\sqrt{T}} + \frac{\rho B}{2\sqrt{T}} = \frac{\rho B}{\sqrt{T}}$$

Note: We want both error terms
to be of equal size

$$\Rightarrow \eta = \sqrt{\frac{B^2}{\rho^2 T}}$$

$$\frac{1}{2\gamma} \sum_{t=1}^T \langle w^{(t)} - w^*, \nabla f(w^{(t)}) \rangle$$

(54)

$$= \frac{1}{2\gamma} \sum_{t=1}^T -\|w^{(t)} - w^* - \gamma \nabla f(w^{(t)})\|^2 + \|w^{(t)} - w^*\|^2 + \gamma^2 \|\nabla f(w^{(t)})\|^2$$

$$= \frac{1}{2\gamma} \sum_{t=1}^T -\|w^{(t+1)} - w^*\|^2 + \|w^{(t)} - w^*\|^2 + \gamma^2 \|\nabla f(w^{(t)})\|^2$$

GD step

$$= \frac{1}{2\gamma} (\|w^{(1)} - w^*\|^2 - \|w^{(T+1)} - w^*\|^2) + \frac{\gamma}{2} \sum_{t=1}^T \|\nabla f(w^{(t)})\|^2$$

$w^{(0)} = 0$

$$\leq \frac{1}{2\gamma} \|w^*\|^2 + \frac{\gamma}{2} \sum_{t=1}^T \|\nabla f(w^{(t)})\|^2$$

Note: $\langle a, b \rangle = -\|\frac{a-b}{2}\|^2 + \|\frac{a}{2}\|^2 + \|\frac{b}{2}\|^2$

parallelogram law