

Online (sub) gradient descent algorithm:

Initializer $w_1 \in \Omega$. Parameter η .

For $t=1, 2, \dots, T$:

- choose w_t .

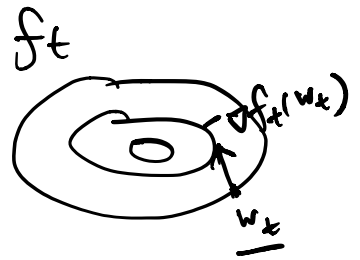
- receive f_t . suffer loss $f_t(w_t)$

- set $g_t \in \partial f_t(w_t)$

- update: $w_{t+1}' \leftarrow w_t - \eta g_t$. ($\eta > 0$)

Ω

$$w_{t+1} \leftarrow \Pi_{\Omega}(w_{t+1}')$$



$$= \operatorname{argmin}_{w \in \Omega} \|w - w_{t+1}'\|_2.$$

$$\|a + b\|_2^2 = \|a\|_2^2 + \|b\|_2^2 + 2\langle a, b \rangle$$

Remark:

$$w_{t+1} = \operatorname{argmin}_{w \in \Omega} \|w - w_t + \eta g_t\|_2^2$$

$$= \operatorname{argmin}_{w \in \Omega} \underbrace{\langle w, \eta g_t \rangle}_{\text{correctiveness}} + \underbrace{\frac{1}{2} \|w - w_t\|_2^2}_{\text{conservativeness}}$$

correctiveness conservativeness

(Kivinen & Warmuth '97).

OGD guarantees:

Then: OGD w/ initializer w_1 & step size $\eta > 0$

guarantees: $\forall u \in \Omega$

$$R_T(u) \leq \frac{\|u - w_1\|_2^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|g_t\|_2^2.$$

Moreover, if Ω has l_2 -diameter B ($\forall u, v \in \Omega$,

$\circledast B$.

$\|u - v\|_2 \leq B$), and $\|g_t\|_2 \leq \rho$. (which happens if all f_t 's are ρ -Lipschitz), then

$$R_T(\Omega) \leq \frac{B^2}{2\eta} + \frac{\eta}{2} \cdot T \cdot \rho^2$$

$\xrightarrow{\eta = \frac{B}{\rho} \cdot \sqrt{\frac{1}{T}}} \sqrt{T} \cdot B^2 \rho \rightarrow \sqrt{T} \cdot \rho \rightarrow \sqrt{T} \cdot \rho \cdot (B^2 + 1)$

$$\eta = \frac{B}{\rho} \cdot \sqrt{\frac{1}{T}} \quad B \cdot \rho \cdot \sqrt{T}$$

Cor: under the above setting, $l(w, z)$ is ρ -Lip w.r.t w ; $D \subset D$, with $f_t(w) = l(w, z_t)$ for $z_1, \dots, z_T \stackrel{iid}{\sim} D$, guarantees that: $\bar{w}_T = \frac{1}{T} \sum_{t=1}^T w_t$,

$$\textcircled{1} \quad \eta = \frac{B}{\rho} \cdot \sqrt{\frac{1}{T}} \Rightarrow \mathbb{E} L_D(\bar{w}_T) \leq \min_{w \in \Omega} L_D(w) + \frac{B\rho}{\sqrt{T}}.$$

$$\textcircled{2} \quad \eta = \frac{1}{\rho} \cdot \sqrt{\frac{1}{T}}, \quad \Omega = \mathbb{R}^d, \quad w_1 = 0, \quad B^2$$

$$\Rightarrow \mathbb{E} L_D(\bar{w}_T) \leq L_D(w^*) + \frac{(\|w^*\|_2^2 + 1)\rho}{\sqrt{T}} \quad \forall w^* \in \mathbb{R}^d.$$

Pf of corollary:

high prob. upper bound on

$$L_D(\bar{w}_T) \leq L_D(w^*) + \frac{R_T(w^*)}{T}$$

+ concentration

$$\frac{1}{T} \sum_{t=1}^T l_t(w_t, z_t) - \frac{1}{T} \sum_{t=1}^T l_t(w^*, z_t) \leq R_T(w^*)$$

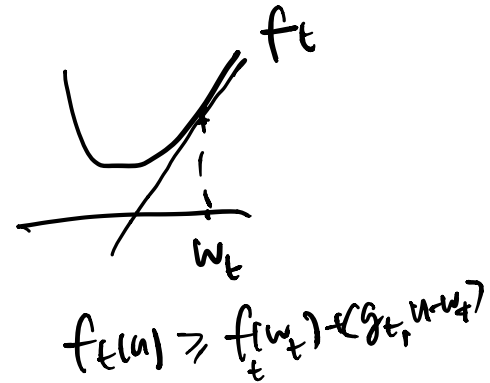
take expectation on LHS:

$$\mathbb{E}[\text{LHS}] = \mathbb{E} \left[\sum_{t=1}^T L_D(w_t) \right] - T \cdot L_D(w^*)$$

pf of OGD guarantees:

step 1: "linearization"

$$R_T(u) = \sum_{t=1}^T (f_t(w_t) - f_t(u))$$



$$\leq \sum_{t=1}^T \langle g_t, w_t - u \rangle$$

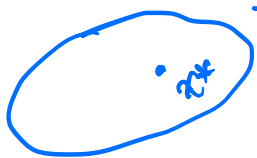
step 2: use optimality condition on w_t :

First order optimality condition:

$f \mapsto \text{conv}$ in conv domain Ω , f differentiable.

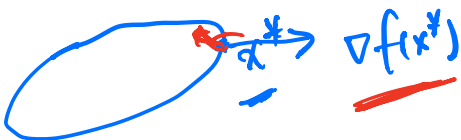
$$x^* = \underset{x \in \Omega}{\text{argmin}} f(x) \iff \forall y \in \Omega, \langle \nabla f(x^*), y - x^* \rangle \geq 0$$

①



x^* in interior of Ω
 $\nabla f(x^*) = 0$

②



x^* we need $\forall y \in \Omega$
 $\langle \nabla f(x^*), y - x^* \rangle \geq 0$

Df idea: (\Rightarrow) if $\exists y, \langle \nabla f(x^*), y - x^* \rangle < 0$.

$$f(x^* + \alpha(y - x^*)) = f(x^*) + \alpha \langle \nabla f(x^*), y - x^* \rangle + o(\alpha)$$

$< f(x^*)$ for small $\alpha > 0$.

(\Leftarrow) $\forall y, f(y) \geq f(x^*) + \underbrace{\langle \nabla f(x^*), y - x^* \rangle}_{\geq 0}$.

$$w_{t+1} = \operatorname{argmin}_{w \in \Omega} \underbrace{\langle \eta g_t, w \rangle + \frac{1}{2} \|w - w_t\|^2}_{f(w)}$$

First order optimality:

$$\langle \eta g_t + w_{t+1} - w_t, \underbrace{u - w_{t+1}}_a \rangle \geq 0, \quad \forall u \in \Omega.$$

$$\Rightarrow \langle g_t, w_{t+1} - u \rangle \leq \frac{1}{\eta} \langle w_{t+1} - w_t, u - w_{t+1} \rangle.$$

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

$$= \frac{1}{2\eta} (\|u - w_t\|^2 - \|u - w_{t+1}\|^2 - \|w_t - w_{t+1}\|^2)$$

step 3: bounding $\langle g_t, w_t - u \rangle$.

$$\langle g_t, w_t - u \rangle = \langle g_t, w_{t+1} - u \rangle + \underbrace{\langle g_t, w_t - w_{t+1} \rangle}_{\|g_t\| \|w_t - w_{t+1}\|}$$

Cauchy-Schwarz & AM-GM.

$$\leq \langle g_t, w_{t+1} - u \rangle + \frac{\eta}{2} \|g_t\|_2^2 + \frac{1}{2\eta} \|w_t - w_{t+1}\|^2$$

$$\leq \frac{\eta}{2} \|g_t\|_2^2 + \frac{1}{2\eta} (\underbrace{\|u - w_t\|^2}_{D_t} - \underbrace{\|u - w_{t+1}\|^2}_{D_{t+1}})$$

step 4: summing over t

$$\sum_{t=1}^T \langle g_t, w_t - u \rangle \quad \underbrace{D_1 - D_2 + D_2 - D_3 + \dots - D_{T+1}}_{\text{telescoping}}$$

$$\leq \frac{\eta}{2} \sum_{t=1}^T \|g_t\|_2^2 + \frac{1}{2\eta} \|u - w_1\|^2 \quad \cancel{\text{X}}$$

Online Mirror descent: exploiting different geometry in data.

— can we develop algs w/ regrets that scale w/ other geometric measures of data. (L₁, L₂, etc).

Back ground on norms:

Def: $f_n = \|\cdot\|: \mathbb{R}^d \rightarrow \mathbb{R}$ said to be a norm, if

1. $\forall a \in \mathbb{R}, x \in \mathbb{R}^d, \|ax\| = |a| \cdot \|x\|$ homogeneity
2. $\forall x, y \in \mathbb{R}^d, \|x+y\| \leq \|x\| + \|y\|$.

① $\|\cdot\|_*$ is also a norm.

② by defn of dual norm,

$$\langle x, z \rangle = \|x\| \left\langle \frac{x}{\|x\|}, z \right\rangle$$

$$\leq \|x\| \cdot \|z\|_* \quad (\text{generalizes Cauchy-Schwarz})$$

Ex: $\|\cdot\|$ $\|\cdot\|_*$

$$\|\cdot\|_2 \longrightarrow \|\cdot\|_2$$

$$\|\cdot\|_1$$

$$\|\cdot\|_\infty$$

$$\|\cdot\|_p, p \in [1, \infty]$$

$$\|\cdot\|_q$$

$$\langle x, y \rangle \leq \|x\|_p \|y\|_q$$

(Hölder's inequality)

q is p 's conjugate exponent,

$$\frac{1}{p} + \frac{1}{q} = 1.$$

$$\|\cdot\|_A$$

$$\|\cdot\|_{A^{-1}} \quad (\text{exercise})$$

General: Lipschitz property (wrt arbitrary norms)

Recall: $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is L -lip wrt $\|\cdot\|$,

$$\forall x, y, |f(x) - f(y)| \leq L \|x - y\|$$

Relating Lipschitzness to gradient norm:

Lemma: ① $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is differentiable, then

$$f \text{ L-Lip wrt } \|\cdot\| \Leftrightarrow \forall x, \|\nabla f(x)\|_* \leq L.$$

② $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex, then

$$f \text{ L-Lip wrt } \|\cdot\| \Leftrightarrow \forall x, \forall g \in \partial f(x), \|g\|_* \leq L.$$

Pf: ① left as exercise

② (\Rightarrow) $\forall x, \forall g \in \partial f(x)$. we have

$$\forall y, f(y) \geq f(x) + \langle g, y-x \rangle.$$

pick $\underline{g^*}$ be $\|g^*\| \leq 1, \langle g, g^* \rangle = \|g\|_*$

$$\underline{y} = x + g^*.$$

$$\|g\|_* = \langle g, y-x \rangle \leq \underline{f(y) - f(x)} \leq L \|y-x\| \leq L.$$

(\Leftarrow) $\forall x, y$

$\partial f(x)$

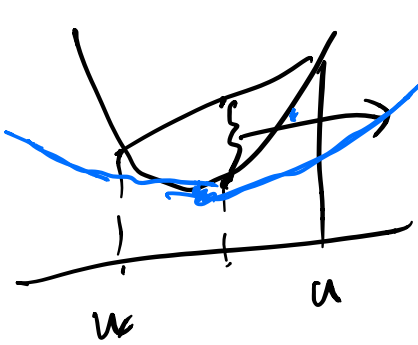
$$\begin{aligned} \langle g_y, x-y \rangle &\leq \underbrace{f(x) - f(y)}_{\leq} \leq \underbrace{\langle g_x, x-y \rangle}_{\leq} \\ &\leq \|g_x\|_* \cdot \|x-y\| \\ &\leq L \|x-y\|. \end{aligned}$$

Strong convexity (for general norms):

Recall f is λ -SC w.r.t $\|\cdot\|$ (general) if, $\forall u, w$,

$\alpha \in (0,1)$,

$$f(\alpha w + (1-\alpha)u) \leq \alpha f(w) + (1-\alpha)f(u) - \frac{\lambda}{2} \alpha(1-\alpha) \|w-u\|^2$$



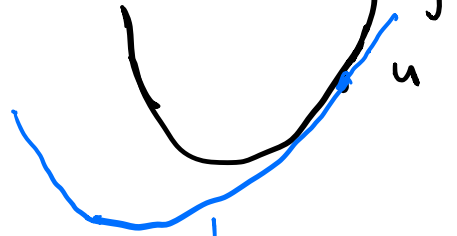
general norm

Lemma: If f is λ -SC, then $\forall u, w$,

$\forall g \in \partial f(u)$,

$$f(w) - f(u) \geq \langle g, w-u \rangle + \frac{\lambda}{2} \|w-u\|^2$$

f



quadratic approximation at $u \Rightarrow \langle g, d(w-u) \rangle$

Pf:
$$\frac{|f(dw + (1-d)u) - f(u)|}{d} \leq \frac{f(w) - f(u) - \frac{\lambda(\epsilon)}{2} \|w-u\|^2}{d}$$

$$\leq \langle g, w-u \rangle$$

pick $d = 0$.

\Rightarrow Lemma statement \star

Def: if ψ is differentiable & strongly convex.

then
$$D\psi(x, y) = \psi(x) - \psi(y) - \langle \nabla \psi(y), x-y \rangle.$$

\Rightarrow is called the Bregman divergence induced by

ψ .

$$\psi(w) = \frac{1}{2} \|w\|_2^2 \Rightarrow D\psi(x, y) = \frac{1}{2} \|x-y\|_2^2.$$

Next lecture:

Online Mirror descent:

$$w_{t+1} = \underset{w \in \Omega}{\operatorname{argmin}} \langle \eta g_t, w \rangle + D_\psi(w, w_t)$$