Future scribing: Go to Piazza, edit link on overleaf repo.

HW3: errata page on Piazza.

Mid-project progress report due today on gradescope.

OCO for strongly cvx fns. / Kernel methods.

Motivation: fast algs for regularized loss minimization

$$\min_{w \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^{m} \left( \underbrace{\ell(w, (x_i, y_i))}_{\text{logistic / hinge}} + \frac{\lambda}{2} \|w\|_2^2 \right). \quad \|x_i\| \le B.$$

$$\underbrace{\phantom{\frac{1}{m} \sum_{i=1}^{m} \ell(w,(x_i,y_i))}}_{F_S(w)}$$

$$w. \qquad F_S(w) \le F_S(\hat{w}) + \varepsilon. \qquad \hat{w} = \operatorname{argmin}_w F_S(w).$$

— Ideal: gradient descent   (evaluating this would take $O(m)$.)

$$w_{t+1} \leftarrow w_t - \eta \underbrace{\nabla F_S(w_t)}$$   OGD's guarantee

$$f_t \equiv F_S. \qquad F_S(\bar{w}_T) - F_S(\hat{w}) \le \frac{1}{\underbrace{\sqrt{T}}_{\#\text{iters}}}.$$

$$T = O\left(\frac{1}{\varepsilon^2}\right).$$

running time $= O\left(\frac{m}{\varepsilon^2}\right).$

— Idea 2: stochastic gradient descent.

define $\hat{D} = \text{unif}\left((x_i, y_i)_{i=1}^{m}\right).$

use OGD on the regularized losses induced by random examples drawn from $\hat{D}$, and do online to batch conversion.

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

sample $i_t \sim \text{uniform}(\{1 \dots m\})$.

$$f_t(w) = \ell(w, (x_{i_t}, y_{i_t})) + \frac{\lambda}{2}\|w\|^2 \quad \lambda\text{-sc wrt } \|\cdot\|_2$$

$$w_{t+1} \leftarrow w_t - \eta \cdot g_t. \quad \text{where } g_t \in \partial f_t(w_t).$$

$$\{w_1 \dots w_T\} \longrightarrow \bar{w}_T.$$

$$\mathbb{E}[F_S(\bar{w}_T)] - F_S(\hat{w}) \overset{?}{\leq} \frac{1}{\sqrt{T}}.$$

$$\left(\mathbb{E}[L_D(\bar{w}_T)] - \min_w L_D(w) \leq \mathbb{E}[R_T(w^*)] \leq \frac{1}{\sqrt{T}}.\right.$$

$$T = O\left(\frac{1}{\varepsilon^2}\right). \quad \text{running time} = O\left(\frac{1}{\varepsilon^2}\right).$$

can we do even better?

**Thm** : $\Omega$ WX. $\{f_t\}_{t=1}^T$ are $\lambda$-sc wrt $\|\cdot\|_2$.

we run time varying step size OGD:

$$w_{t+1} \leftarrow w_t - \eta_t \cdot g_t. \qquad g_t \in \partial f_t(w_t)$$

$$\eta_t = \frac{1}{\lambda t}, \quad \text{and assume } \forall t. \ \|g_t\|_2 \leq L.$$

then,

$$R_T(\Omega) \leq \frac{(\ln T + 1) L^2}{2\lambda} \quad \left(\begin{array}{l}\text{much better than} \\ \text{the basic } O(\sqrt{T}) \\ \text{regret, by exploiting}\end{array}\right.$$

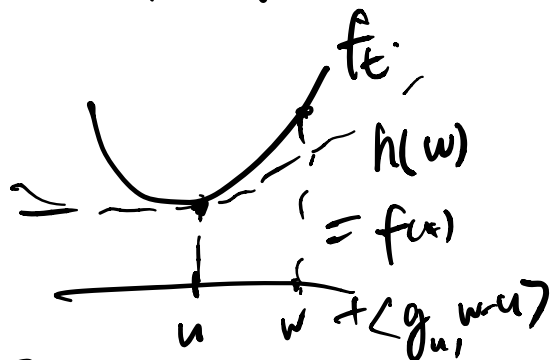This implies that running $\frac{1}{\lambda t}$-step-size SGD.

yields

$$F_S(\bar{w}_T) - F_S(\hat{w}) \leq \tilde{O}\left(\frac{1}{T}\right).$$

$$T = O\left(\frac{1}{\varepsilon}\right). \ll O\left(\frac{1}{\varepsilon^2}\right).$$

Pf: ① by quadratic lower bound property of $\lambda$-SC fns

$$u \in \Omega$$



$$\Rightarrow f_t(w_t) - f_t(u)$$

$$\leq \langle g_t, w_t - u \rangle - \frac{\lambda}{2}\|w_t - u\|^2.$$

( improving over the linearization step.

by utilizing $\lambda$-SC

② Recall that in OGD analysis:

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

Combining ① ② . summing over $t = 1, 2 \cdots T$.

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

$$\leq \sum_{t=1}^{T} \frac{\|u - w_t\|_2^2 - \|u - w_{t+1}\|_2^2}{2\eta_t} - \sum_{t=1}^{T} \frac{\lambda}{2} \|u - w_t\|_2^2$$

$$+ \sum_{t=1}^{T} \eta_t \cdot \|g_t\|_2^2 \; .$$

$\eta_t = \frac{1}{\lambda t}$

$$\underline{\underline{\frac{\lambda \|u - w_1\|_2^2}{2\eta_1}}} - \frac{\lambda \|u - w_2\|_2^2}{2\eta_1} + \frac{\|u - w_2\|_2^2}{2\eta_2} - \frac{\|u - w_3\|_2^2}{2\eta_2} +$$

$$\cdots + \frac{\|u - w_T\|^2}{2\eta_T} - \underbrace{\frac{\|u - w_{T+1}\|^2}{2\eta_T}}_{\leq 0}$$

coefficient of $\|u - w_1\|_2^2$ : $\frac{1}{2\eta_1} - \frac{\lambda}{2} = 0$ .

$$\|u - w_2\|_2^2 \quad : \quad -\frac{1}{2\eta_1} + \frac{1}{2\eta_2} - \frac{\lambda}{2} = 0$$

$$\|u - w_T\|_2^2 \quad -\frac{1}{2\eta_{T-1}} + \frac{1}{2\eta_T} - \frac{\lambda}{2} = 0$$

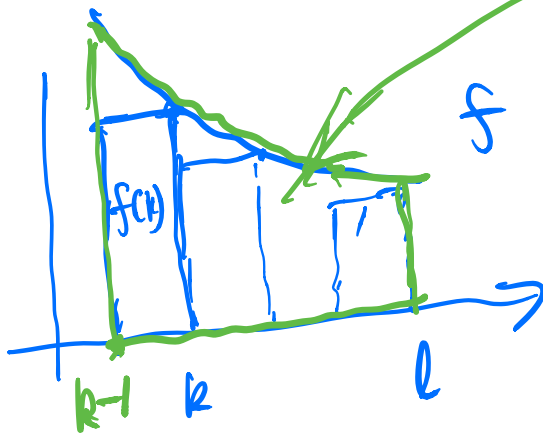$$R_T(u) \leq \sum_{t=1}^{T} \frac{\eta_t}{2} \|g_t\|_2^2 \; . \qquad \leq \int_1^T \frac{1}{x} dx = \ln T$$

$$= \frac{L^2}{2\lambda} \left( \sum_{t=1}^{T} \frac{1}{t} \right) \longrightarrow = 1 + \sum_{t=2}^{T} \frac{1}{t}$$

$$\leq 1 + \ln T. \quad k=2, \; \ell=T.$$

$$f(x) = \frac{1}{x}$$

**Fact**: $f$ decreasing, then

$$\sum_{i=k}^{\ell} f(i) \leq \int_{k-1}^{\ell} f(x) \, dx.$$



$f(x)$

$f$

$k-1 \quad k \quad \quad \ell$

regularized.

Instantiating this result to unconstrained loss
minimization:

$\Omega = \mathbb{R}^d.$

For $t = 1, 2, \cdots T$.

$\beta$-Lip

sample $i_t \sim$ uniform $(\{1 \cdots n\})$.

$\ell_t(w)$

$$f_t(w) = \boxed{\ell(w, (x_{i_t}, y_{i_t}))} + \frac{\lambda}{2} \|w\|^2. \quad \lambda\text{-sc wrt} \quad \|\cdot\|_2$$

$\hookleftarrow \; w_{t+1} \leftarrow w_t - \frac{1}{\lambda t} g_t.$ where $g_t \in \partial f_t(w_t).$

(Fact: if $f, g$ cvx. $h = f + g$.

calculating $g_t$:

— $V_t \in \partial l_t(w_t)$

— $g_t = V_t + \lambda w_t$.

updating $w_t$:

$$w_{t+1} \leftarrow w_t - \frac{1}{\lambda t}(\lambda w_t + V_t)$$

$$= (1-\tfrac{1}{t}) w_t - \frac{1}{\lambda t} V_t.$$

observation:

$$\underbrace{t \cdot w_{t+1}}_{A_{t+1}} = \underbrace{(t-1) w_t}_{A_t} - \frac{1}{\lambda} V_t.$$

$\Rightarrow \qquad A_{t+1} = \sum_{s=1}^{t} -\frac{1}{\lambda} V_s.$

$\qquad\qquad\qquad\qquad\qquad \textcolor{blue}{\in \partial_t \phi(x_{it})}.$

$\Rightarrow \qquad w_{t+1} = -\frac{1}{\lambda t} \sum_{s=1}^{t} V_s.$

$$R_T(\underline{\underline{\Omega}}) \le O\left(\frac{\overset{\mathbb{R}^d}{\| }\ln T}{T}\right)$$

suppose additionally, $\ell_t$'s are $B$-Lip. wrt $\|\cdot\|_2$ $\forall t$  $f_t$ is are B-Lip.

then. $\|V_t\|_2 \le B$

$\Rightarrow \|W_t\|_2 \le \dfrac{B}{\lambda}$.

$\Rightarrow g_t = \lambda \underset{\underset{B}{\overline{\wedge}}}{W_t} + \underset{\underset{B}{\wedge}}{V_t}$ satisfies $\|g_t\|_2 \le 2B$.

Thm

$\Rightarrow \forall u \in \underset{=\mathbb{R}^d}{\underline{\Omega}},\quad R_T(u) \le \dfrac{4B^2 \cdot (\ln T + 1)}{2\lambda}$. $L$

online to batch conversion

$\Rightarrow \mathbb{E}\, F_s(\bar{w}_T) - F_s(\hat{w}) \le \dfrac{2B^2(\ln T + 1)}{\lambda T}$.

$\Rightarrow$ setting $T = \tilde{O}\left(\dfrac{B^2}{\lambda \varepsilon}\right)$ ensures

$\mathbb{E}\, F_s(\bar{w}_T) - F_s(\hat{w}) \le \varepsilon$.

kernel methods : a brief introduction

Suppose all examples are transformed by a nonlinear map: $\phi: \mathbb{R}^d \to \mathbb{R}^N$. N is very large.

Goal: find $w$ that approx minimizes

$$\frac{1}{m} \sum_{i=1}^{m} f( \underbrace{y_i \langle \underbrace{w}_{\in \mathbb{R}^N}, \phi(x_i) \rangle}_{\ell(w, \phi(x_i), y_i)} ) + \frac{\lambda}{2} \|w\|_2^2 .$$

Silver lining: assume that $\langle \phi(x), \phi(z) \rangle = K(x, z)$ (kernel fn by $\phi$). can be evaluated efficiently.

$$K(x, z) = (1 + \langle x, z \rangle)^{\ell} . \quad \text{(polynomial kernel)}$$

the $\phi$ induced by $K$ will have $N = O(d^{\ell})$.

— can we develop efficient trainig & test algs w/ runnig time independent from N?

— key idea:
   keep track of coefficient of $w_t$. $\alpha_t \in \mathbb{R}^m$
   maintain invariant that $w_t = \sum_{i=1}^{m} \alpha_t(i) \phi(x_i)$.

$(m << N)$.

$$W = \sum_{i=1}^{m} \alpha(i) \phi(x_i).$$

prediction:

$$\langle w, \phi(x) \rangle = \sum_{i=1}^{m} \alpha_i \langle \phi(x_i), \phi(x) \rangle$$

$$= \sum_{i=1}^{m} \alpha_i k(x_i, x) \quad \text{. can be}$$

done efficiently

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

sample $i_t \sim$ uniform $(\{1 \dots m\})$.

$$f_t(w) = \overbrace{\ell(w, (x_{i_t}, y_{i_t})) + \frac{\lambda}{2} \|w\|^2}^{\ell_t(w)} \quad \begin{array}{l} \lambda - \text{SC wrt} \\ \|\cdot\|_2 \end{array}$$

$\Omega = \mathbb{R}^d$.

calculating $g_t$:

— $V_t \in \partial \ell_t(w_t)$

— $g_t = V_t + \lambda w_t$.

updating $\underline{w_t}$:

$$w_{t+1} \leftarrow w_t - \frac{1}{\lambda t}(\lambda w_t + V_t)$$

$$= (1 - \frac{1}{t}) w_t - \frac{1}{\lambda t} V_t.$$

Q: can we modify the alg so that instead
of keep $w_t$'s, we keep $\alpha_t$'s.