

# Announcement:

1. HW2 grades are out  
HW3 due on Apr. 15.

Piazza: errata page for HW3.

2. mid project progress report due Apr. 6.

Problem 1: SC.  
First-order optimality  
for general norms

23: guarantee of  
OGD.

Cor: ① map  $u, v \in \Omega$   $\|u - v\|_2 \leq B$ ; ②  $l(w, z)$  is  $\rho$ -Lip, or  
wrt  $w$ ; OGD, with  $f_t(w) = l(w, z_t)$  for  $z_1, \dots, z_T \stackrel{iid}{\sim}$   
 $\mathcal{D}$ , guarantees that:  $\bar{w}_T = \frac{1}{T} \sum_{t=1}^T w_t$ .

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

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

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

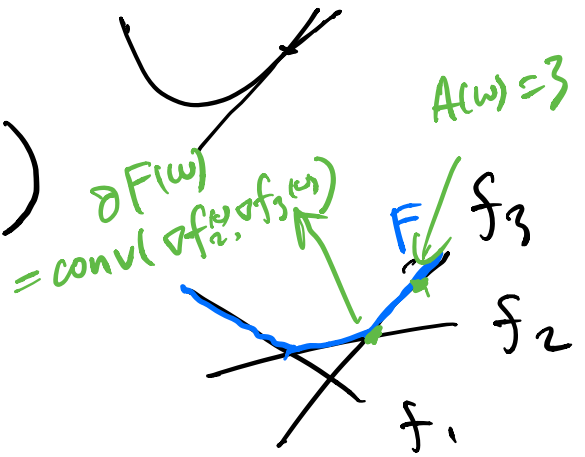
OGD + online to batch conversion competes with the  
guarantees of regularized loss minimization with only  
one pass over the data!

# Calculating subgradients:

Fact: For  $F(w) = \max_{i=1}^n (f_i(w))$

each  $f_i$  is convex

$$\partial F(w) = \text{conv} \left( \bigcup_{i \in A(w)} \partial f_i(w) \right)$$



$$A(w) = \left\{ j : j \in \underset{i}{\text{argmax}} f_i(w) \right\}$$

$$\text{conv}(\{x_1, \dots, x_n\}) = \left\{ \sum_{i=1}^n a_i x_i : a_i \in \Delta^{n-1} \right\}$$

Ex:  $F(w) = |w|$       $\partial F(0) = \text{conv}(\{1, -1\})$

$$F(w) = \max_{f_1, f_2} (0, 1 - \gamma \langle w, x \rangle)$$

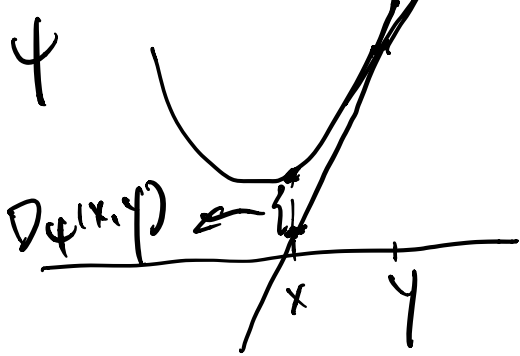
$$\partial F(w) = \begin{cases} \{0\} & 1 - \gamma \langle w, x \rangle < 0 \\ \{ -\gamma x \} & 1 - \gamma \langle w, x \rangle > 0 \\ \{ -\alpha \gamma x : \alpha \in [0, 1] \} & 1 - \gamma \langle w, x \rangle = 0 \end{cases}$$

strongly convex.

Def: if  $\psi$  is differentiable & strongly convex, then  $D_\psi(x, y) = \psi(x) - \psi(y) - \langle \nabla \psi(y), x - y \rangle$  is called the Bregman divergence induced by  $\psi$ .

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

①  $\psi$  is  $\lambda$ -SC wrt  $\|\cdot\|$ .  $x \neq y$ .  
 $D_\psi(x, y) \geq \frac{\lambda}{2} \|x - y\|^2 > 0$



②  $D_\psi$  can be asymmetric.

$D_\psi(x, y)$  not necessarily  $= D_\psi(y, x)$ .

Important examples of strongly conv fns & Bregman Divergences

①  $\psi(w) = \frac{1}{2} \|w\|_A^2$  is 1-sc wrt  $\|\cdot\|_A$ .

$$D_\psi(x, y) \stackrel{?}{=} \frac{1}{2} \|x - y\|_A^2$$



②  $\Omega = \Delta^{d-1}$ ,  $\psi(w) = \sum_{i=1}^d w_i \ln w_i$  is 1-sc wrt  $\|\cdot\|_1$  (HW3).  
negative entropy

$$D_\psi(x, y) = \sum_{i=1}^d x_i \ln \frac{x_i}{y_i}$$

asymmetric.

(Relative entropy)  
K-L divergence  
Kullback Leibler

③  $\psi(w) = \frac{1}{2(p-1)} \|w\|_p^2$  is 1-sc wrt  $\|\cdot\|_p$ .

For  $p \in (1, 2]$ . (e.g. Shalev-Schwartz, 2007).

Online Mirror descent:

$$\left. \begin{array}{l} w_t \\ w_{t+1} \end{array} \right\} \quad \underline{g_t \in \partial f_t(w_t)}$$

$$w_{t+1} = \operatorname{argmin}_{w \in \Omega} \left( \langle w, \eta g_t \rangle + \frac{D\psi(w, w_t)}{\psi(w) - \langle \nabla\psi(w_t), w - w_t \rangle - \psi(w_t)} \right)$$

$$= \operatorname{argmin}_{w \in \Omega} \langle w, \eta g_t - \nabla\psi(w_t) \rangle + \psi(w)$$

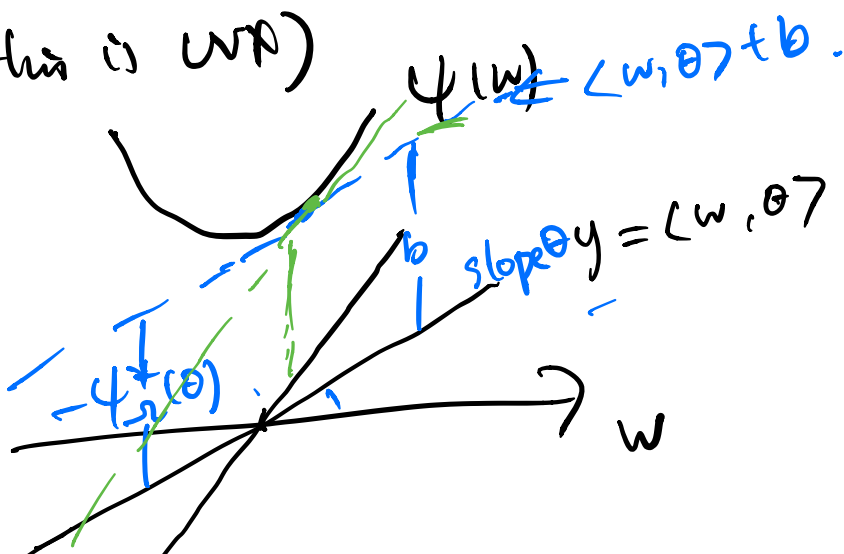
$$= \operatorname{argmax}_{w \in \Omega} \langle w, \underbrace{\nabla\psi(w_t) - \eta g_t}_{\theta} \rangle - \psi(w)$$

Fact  $\psi$  is strongly  $\mu$ ,  $\Omega$  convex set.

$$\psi_{\Omega}^*(\theta) = \max_{w \in \Omega} \langle w, \theta \rangle - \psi(w) \quad \text{is } \psi' \text{'s}$$

Fenchel conjugate wrt  $\Omega$

(this is WPA)

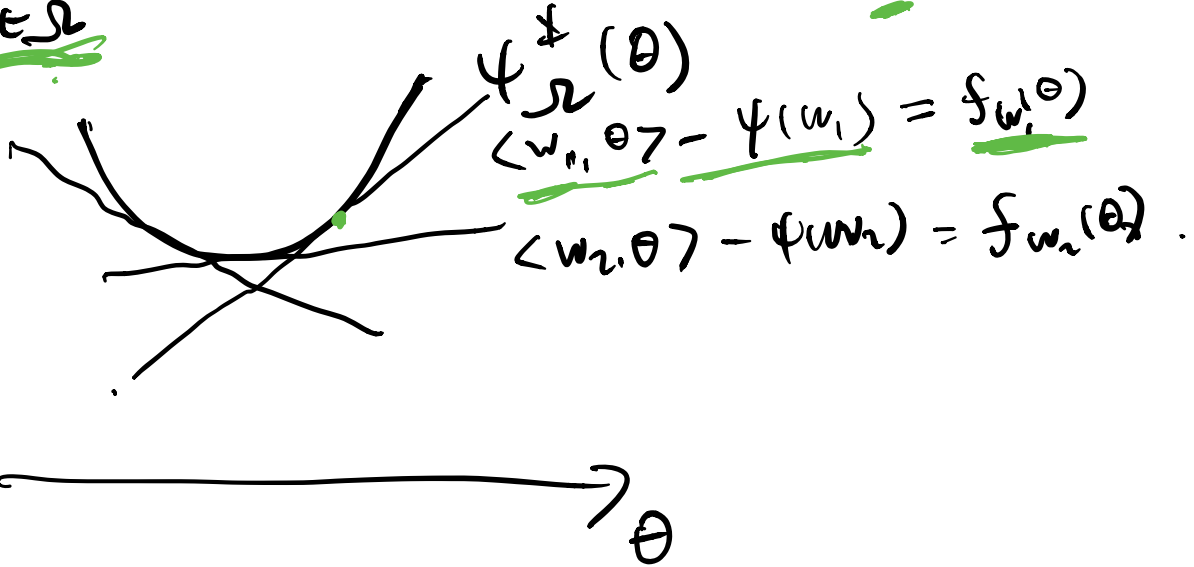


$$\min_w \psi(w) - \langle w, \theta \rangle + b = 0$$

$\nabla \psi_{\Omega}^*(\theta)$  exists  $\in \mathbb{R}^d$ .

$= \arg \max_{w \in \Omega} \langle w, \theta \rangle - \psi(w)$ .

$\nabla_{\theta} f_{w_1}(\theta) = w_1$

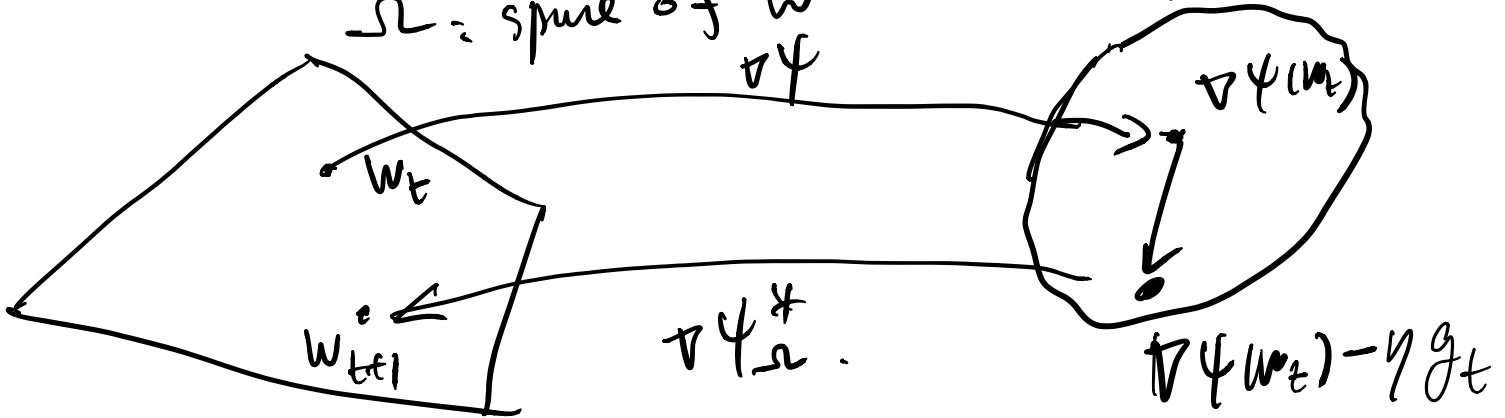


OMD:

$w_{t+1} = \nabla \psi_{\Omega}^* (\nabla \psi(w_t) - \eta g_t)$

$\Omega$ : space of  $w$

space of gradients



Two examples of OMD:

1:  $\Omega = \mathbb{R}^d$ .  $\psi(w) = \frac{1}{2} \|w\|_p^2$   $p \in (1, 2]$

$$\psi_{\Omega}^*(\theta) = \frac{1}{2} \|\theta\|_q^2 \quad \text{where } q \text{ is s.t.}$$

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

(exercise)

$$\nabla \psi_{\Omega}^* = (\nabla \psi)^{-1}$$

$$\text{OMD: } w_{t+1} = \nabla \psi_{\Omega}^* (\nabla \psi(w_t) - \eta g_t)$$

$$\Rightarrow \nabla \psi(w_{t+1}) = \nabla \psi(w_t) - \eta g_t.$$

$$\Rightarrow \nabla \psi(w_t) = -\eta \sum_{s=1}^{t-1} g_s.$$

$$\Rightarrow w_t = \nabla \psi_{\Omega}^* \left( -\eta \sum_{s=1}^{t-1} g_s \right).$$

$p$ -norm algorithm.

2. exponential weight algorithm.

$$\Omega = \Delta^{d-1}, \quad \psi(w) = \sum_i \frac{w_{(i)}}{w(i)}$$

distance generating fn

$$w_i = \left( \frac{1}{d}, \dots, \frac{1}{d} \right)$$

$$\psi_{\Omega}^*(\theta) = \max_{w \in \Delta^{d-1}} \langle w, \theta \rangle - \sum_{i=1}^d w(i) \ln w(i)$$

$$= \max_{\substack{w(1) - w(d-1) \geq 0 \\ \sum_{i=1}^{d-1} w(i) \leq 1}} \sum_{i=1}^{d-1} w(i) \theta(i) + \left( 1 - \sum_{i=1}^{d-1} w(i) \right) \theta(d)$$

$$- \sum_{i=1}^{d-1} w(i) \ln w(i) + \left( 1 - \sum_{i=1}^{d-1} w(i) \right) \ln \left( 1 - \sum_{i=1}^{d-1} w(i) \right)$$

$$= \dots \quad \left( \text{HW 3 } \psi(w) \in [-\ln d, 0] \right)$$

$$= \ln \left( \sum_{i=1}^d e^{\frac{\theta(i)}{d}} \right)$$

$$\nabla \psi_{\Omega}^*(\theta) = \left( \frac{e^{\theta(i)}}{\sum_{j=1}^d e^{\theta(j)}} \right)_{i=1}^d$$

$$\nabla \psi(w) = \left( \ln w(i) + 1 \right)_{i=1}^d$$

$$w_{t+1} = \frac{\nabla \psi_{\Omega}^* \left( \nabla \psi(w_t) - \eta g_t \right)}{1}$$

$$e^{x+y} = e^x \cdot e^y$$

$$= \left( \frac{w_t(i) \cdot e^{-\eta g_t(i)}}{e \cdot \sum_{j=1}^d w_t(j) e^{-\eta g_t(j)}} \right)_{i=1}^d$$

$$= \left( \frac{w_t(i) e^{-\eta g_t(i)}}{\sum_{j=1}^d w_t(j) e^{-\eta g_t(j)}} \right)_{i=1}^d$$

induction.

$$\Rightarrow w_t(i) \propto \frac{w_1(i) \cdot e^{-\eta \sum_{s=1}^{t-1} g_s(i)}}{\frac{1}{d}}$$

$$\Rightarrow w_t = \nabla \Psi_{\Omega}^* \left( -\eta \sum_{s=1}^{t-1} g_s \right)$$

This formula is not always true. just coincidence.

Guarantees of OMD:



Thm: If  $\psi$  is 1-SC wrt  $\|\cdot\|$ , then  
OMD w/  $\psi$  and learning rate  $\eta$  has  
regret guarantee:

$$\forall u \in \Omega: R_T(u) \leq \frac{D_\psi(u, w_1)}{\eta} + \frac{\eta}{2} \sum_{t=1}^T \|g_t\|_*^2$$

Specifically, if  $D_\psi(u, w_1) \leq H^2$  and

$\forall t, \|g_t\|_* \leq \rho$ , then

$$\eta = \frac{H}{\rho} \sqrt{\frac{1}{T}} \Rightarrow R_T(u) \leq H \cdot \rho \cdot \sqrt{T}.$$