

## Lecture 10 — Sept 04

Lecturer: Aditya Gopalan

Scribe: Prakash Chandra

## 10.1 Recap

1. (Projected) Online Gradient Descent at time  $t$  :

$$y_t := w_{t-1} - \eta \nabla f_{t-1}(w_{t-1})$$

$$w_t := \prod_k y_k$$

2. Theorem:

$$\text{Regret}_T(\text{POGD}) \leq DG\sqrt{T}$$

where,  $D$  = diameter of convex space  $K$  and  $G$  = bound on  $\|\text{gradient}\|_2$

3.  $\sigma$ - strongly convex function:

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\sigma}{2} \|y - x\|^2$$

## 10.2 OGD with strongly convex losses

### 10.2.1 Logarithmic regret with time-varying learning rate

**Theorem 1:** Let  $\{f_t\}$  be  $\sigma$ - strongly convex OGD with time-varying step size:  $\eta_t = 1/\sigma_t$ , gives:

$$\text{Regret}_T(\text{POGD}) \leq \frac{G^2}{2\sigma} (1 + \log T) \quad (10.1)$$

**Notes:** 1. Extra curvature make regret  $O(\sqrt{T})$  to  $O(\log T)$ .

2. Strong convexity + bounded gradients is weaker than exp-concave.

**Proof:** start with the same approach as for “ $DG\sqrt{T}$ ”.

let,

$$w^* = \arg \min_{w \in K} \sum_{t=1}^T f_t(w)$$

$$f_t(w_t) - f_t(w^*) \leq \overbrace{\langle \nabla f_t(w_t), w_t - w^* \rangle}^{g(t)} - \frac{\sigma}{2} \|w_t - w^*\|^2 \quad (10.2)$$

hence,

$$\begin{aligned} \langle g(t), w_t - w^* \rangle &= \frac{1}{2\eta_t} [2\eta_t \langle g(t), w_t - w^* \rangle] \\ \Rightarrow \langle g(t), w_t - w^* \rangle &\leq \frac{\eta_t}{2} \|g(t)\|^2 \frac{\|w_t - w^*\|^2 - \|w_{t+1} - w^*\|^2}{2\eta_t} \end{aligned} \quad (10.3)$$

where,  $\eta_t = \frac{1}{\sigma t}$

use (10.2) together with (10.1) and sum over  $t = 1, 2, \dots, T$ .

$$\sum_{t=1}^T [f_t(w_t) - f_t(w^*)] \leq \frac{G^2}{2} \sum_{t=1}^T \eta_t + \frac{1}{2} \sum_{t=1}^T \left[ \frac{\beta_t - \beta_{t+1}}{\eta_t} - \sigma \beta_t \right]$$

where,  $\beta_t = \|w_t - w^*\|^2$

$$\begin{aligned} &= \frac{G^2}{2} \sum_{t=1}^T \left( \frac{1}{\sigma t} \right) + \frac{1}{2} \sum_{t=1}^T \left[ \left( \frac{1}{\eta_t} - \sigma \right) \beta_t - \frac{1}{\eta_t} \beta_{t+1} \right] \\ &\leq \frac{G^2}{2\sigma} (\log(T) + 1) + \frac{1}{2} \underbrace{\left( \frac{1}{\eta_1} - \sigma \right)}_{=0} \beta_1 + \sum_{t=0}^T \left[ \frac{1}{\eta_t} - \sigma - \frac{1}{\eta_{t-1}} \right] \beta_t \end{aligned}$$

as  $\frac{1}{\eta_t} - \sigma - \frac{1}{\eta_{t-1}} = \sigma t - \sigma - \sigma(t-1) = 0$

$$= \frac{G^2}{2\sigma} (1 + \log(T))$$

### 10.3 Impact of regularizer on FTRL performance

1. FTRL algorithm, regularizer =  $\|\cdot\|^2 \Rightarrow$  Gradient descent.
2. OGD  $N$ -expert gives regret =  $O\sqrt{NT}$ .
3. “What’s a good regularizer instead of” for my problem?

$$\text{Regret}_T^{\text{FTRL}}(u) \leq R(u) - R(w_1) + \sum_{t=1}^T [f_t(w_t) - f_t(w_{t+1})]$$

$\therefore$  controlling  $[f_t(w_t) - f_t(w_{t+1})]$   
 $\Rightarrow$  control Regret (general philosophy)

**Definition 2:** [Lipschitz continuity]

$f : \rightarrow \mathbb{R}$  is Lipschitz continuous w.r.t a norm  $\|\cdot\|_{\square}$  if

$$|f(x) - f(y)| \leq L\|x - y\|_{\square}$$

where,  $x, y \in K$

### 10.3.1 FTRL regret bound with Lipschitz losses + strongly convex regularizer

**Theorem 3:** [FTRL regret w/Lipschitz losses (w.r.t  $\|\cdot\|_{\square}$ ) + strongly-convex regularizer]

Suppose  $f_1, f_2, \dots$  is such that  $f_t$  is  $L_t$ -Lipschitz continuous (w.r.t  $\|\cdot\|_{\square}$ ). Let the regularizer  $R$  be  $\sigma$ -strongly convex w.r.t the same norm  $\|\cdot\|_{\square}$ .

Then,  $\forall u \in K$

$$\text{Regret}_T^{\text{FTRL}}(u) \leq R(u) - \min_{v \in K} R(v) + \frac{TL^2}{\sigma} \quad (10.4)$$

Let's apply this(theorem 3) for expert advice problem  $W/N$  experts. i.e  $K = \Delta_N$ ,  $f_t(\pi) := \langle \pi, Z_t = \text{loss vector at time } t \in \mathbb{R}^N \text{ or } \in [0, 1]^N \rangle$   
 use **entropic regularizer:**

$$R_{\eta}(w) := -\frac{1}{\eta} H(w) = \frac{1}{\eta} \sum_i w_i \log(w_i)$$

so

$$\text{FTRL} \Rightarrow \text{Exp-weight}$$

Let's first go through the following claim:

**Claim 4:**  $R_1(w) = -H(w)$  is  $\frac{1}{B}$ -stronger convex over  $K_B := w \in \mathbb{R}_+^N : \|w\|_1 \leq B$  w.r.t  $\|\cdot\|_1$

**Proof:**

Lipschitz continuity of  $\{f_t\}$

$$\begin{aligned} f_t(\pi) - f_t(\Psi) &= |\langle Z_t, \pi - \Psi \rangle| \\ &\leq \underbrace{\|\pi - \Psi\|_1 \|Z_t\|_{\infty}}_{\text{Holder's inequality}} \leq 1 \|\pi - \Psi\|_1 \end{aligned}$$

$\Rightarrow \{f_t\}$  are 1-Lipschitz continuous w.r.t  $\|\cdot\|_1$   
 $\Rightarrow R_1$  is  $\frac{1}{B}$ -strongly convex.

---

Let's use this claim (claim 4) for entropic regularizer:

hence,  $R_\eta = \frac{1}{\eta} R_1$

so  $R_\eta$  is  $\frac{1}{\eta}$ -strong convex over  $\Delta_N$ .

Apply theorem 3 with  $L = 1$ ,  $\sigma = \frac{1}{\eta}$  in above expert advice with N experts problem setup:

$$\begin{aligned} \text{Regret}_T^{\text{Exp-wts}(\eta)}(u) &\leq - \min_{v \in \Delta_N} R(v) + T\eta \\ &= \max_{v \in \Delta_N} (v) + T\eta \\ \Rightarrow \text{Regret}_T^{\text{Exp-wts}(\eta)}(u) &\leq \frac{\log(N)}{\eta} + T\eta \end{aligned}$$

note that if  $\eta = \sqrt{\frac{\log(N)}{T}}$  then  $\text{Regret}_T^{\text{Exp-wts}(\eta)}(u) \leq \sqrt{T \log(N)}$ .

- $O(\sqrt{T \log(N)})$  is the “right scaling” for Exp-wts:

-BOTTOMLINE: choice of regularizer is important:

-Can/should depend on Lipschitz continuity of losses:

-Also depend on structure of  $K$ , i.e. how large (w.r.t  $\|\cdot\|_\square$ ) does  $K$  look.

e.g for best experts:

$$\text{diam}_{\|\cdot\|}(K) \simeq \text{diam}_{\|\cdot\|_1}(K) = O(1)$$

-But, cost functions are 1-Lip. w.r.t  $\|\cdot\|_1$ , but  $\sqrt{N}$ -Lip. w.r.t  $\|\cdot\|_2$ .

---

### Proof of theorem 3:

$\forall t$ ,

$$\Phi_t(w) = \sum_{s=1}^{t-1} f_s(w) + R(w) \tag{10.5}$$

FTRL picks  $w_t$  such that:

$$w_t = \arg \min_{w \in K} \Phi_t(w) \tag{10.6}$$

hence,  $R$  is  $\sigma$ -strongly convex w.r.t  $\|\cdot\|_\square$

$\Rightarrow \Phi_t$  is  $\sigma$ -strongly convex w.r.t  $\|\cdot\|_\square$

because adding  $\sigma$ -str-cvx + cvx is  $\sigma$ -str-cvx

Let's use the following lemma:

**Lemma 5:** if  $f$  is  $\sigma$ -str-convex and  $x^* = \min_{x \in K}(x)$   
 then  $f(x) - f(x^*) \geq \frac{\sigma}{2} \|x - x^*\|_{\square}^2$

using lemma 5 and (10.6):

$$\Phi_t(w_{t+1}) - \Phi_t(w_t) \geq \frac{\sigma}{2} \|w_{t+1} - w_t\|_{\square}^2 \quad (10.7)$$

$$\Phi_{t+1}(w_t) - \Phi_{t+1}(w_{t+1}) \geq \frac{\sigma}{2} \|w_{t+1} - w_t\|_{\square}^2 \quad (10.8)$$

adding inequalities (10.7) and (10.8) and using (10.5):

$$\begin{aligned} f_t(w_t) - f_t(w_{t+1}) &\geq \sigma \|w_t - w_{t+1}\|_{\square}^2 \\ \Rightarrow \sigma \|w_t - w_{t+1}\|_{\square}^2 &\leq \underbrace{f_t(w_t) - f_t(w_{t+1})}_{\text{Lipschitz continuity}} \leq L_t \|w_t - w_{t+1}\|_{\square} \end{aligned}$$

$$\Rightarrow \|w_t - w_{t+1}\|_{\square} \leq \frac{L_t}{\sigma} \quad (10.9)$$

and,

$$\begin{aligned} \text{Regret}_T^{\text{FTRL}}(u) &\leq R(u) - R(w) + \sum_{t=1}^T [f_t(w_t) - f_t(w_{t+1})] \\ &\leq R(u) - \min_{v \in K} R(v) + \sum_{t=1}^T L_t \|w_t - w_{t+1}\|_{\square} \end{aligned}$$

from (10.9)

$$\leq R(u) - \min_{v \in K} R(v) + \sum_{t=1}^T L_t \frac{L_t}{\sigma}$$

hence,  $L^2 \geq \frac{1}{T} \sum_{t=1}^T L_t^2$

$$\Rightarrow \text{Regret}_T^{\text{FTRL}}(u) \leq R(u) - \min_{v \in K} R(v) + \frac{TL^2}{\sigma}$$

### 10.3.2 Different View of FTRL: “Online Mirror Descent”

Recall FTRL applied to linear costs  $\{f_t\}$ .

$f_t(\cdot) = \langle Z_t, \cdot \rangle$  over  $K \subseteq \mathbb{R}^d$ , regularizer  $R$

$= \arg \min_{w \in K} [\langle w_t, Z_{1:t} \rangle + R(w)]$  where  $Z_{1:t} = \sum_{i=1}^t Z_i$

$= \arg \max_{w \in K} [\langle w_t, -Z_{1:t} \rangle - R(w)]$

let  $h : \mathbb{R}^d \rightarrow K$

$h(\theta) := \arg \max_{w \in K} [\langle w, \theta \rangle - R(w)]$

Called “link function” or “prox function”

$\therefore$  FTRL can equivalently written as:

$\theta_1 = 0$  and  $\forall t = 1, 2, 3, \dots$

(1) *predict* :  $w_t = h(\theta_t)$

(2) *update* :  $\theta_{t+1} = \theta_t - Z_t$

- If dealing with general convex  $\{f_t\}$

We can feed the gradients  $\nabla f_t(w_t)$

i.e. we use linear functions:  $\tilde{f}_t \equiv \langle \nabla f_t(w_t), \cdot \rangle$

to get regret:

$$\begin{aligned} \text{Regret}_T(u) &= \sum_{t=1}^T [f_t(x_t) - f_t(u)] \\ &\leq \sum_{t=1}^T [\langle \nabla f_t(x_t), x_t \rangle - \langle \nabla f_t(x_t), u \rangle] \\ &= \sum_{t=1}^T [\tilde{f}_t(x_t) - \tilde{f}_t(u)] \end{aligned}$$

This generic reduction gives a general algorithm [mirror descent]:

**Algorithm:**

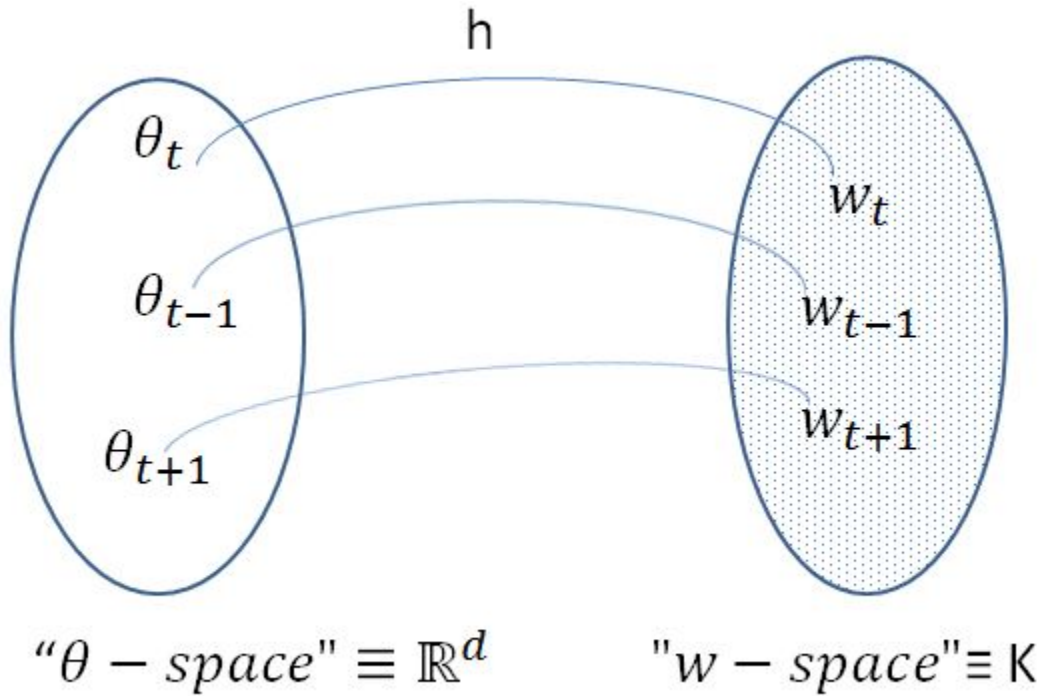
$\theta_1 = 0$

and  $\forall t = 1, 2, 3, \dots$

(1) *predict* :  $w_t = h(\theta_t)$

(2) *update* :  $\theta_{t+1} = \theta_t - \nabla f_t(w_t)$

**Geometric Picture:** [“Dual” spaces figure:]



**Interpretation:**  $\theta$  is updated in dual space and prediction is linked/mirrored to the primal space ( $K$ ).

---