

Chapter 7 Follow The Regularized Leader

Algorithm 7.1 Follow-the-Regularized-Leader

```
Require: Closed and non-empty set  $V \subseteq \mathbb{R}^d$ , a sequence of regularizers  $\psi_1, \dots, \psi_T : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ 
1: for  $t = 1$  to  $T$  do
2:   Output  $x_t \in \operatorname{argmin}_{x \in V} \psi_t(x) + \sum_{i=1}^{t-1} \ell_i(x)$ 
3:   Receive  $\ell_t : \mathbb{R}^d \rightarrow (-\infty, +\infty]$  and pay  $\ell_t(x_t)$ 
4: end for
```

7.2 FTRL Regret Bound using Strong Convexity

We consider a easy case to upper bound the regret for FTRL, where the losses plus regularizer are strongly convex

7.2.1 Properties of Strongly Convex Functions

We first give several properties of strongly convex functions.

μ -strongly convex:

$$f(x_0) + \langle Df(x_0), x - x_0 \rangle + \frac{\mu}{2} \|x - x_0\|^2 \leq f(x), \forall x$$

Lemma 7.5. Let $f : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ μ -strongly convex with respect to a norm $\|\cdot\|$. Then, for all $x, y \in \operatorname{dom} f$, $g \in \partial f(y)$, and $g' \in \partial f(x)$, we have

$$f(x) - f(y) \leq \langle g, x - y \rangle + \frac{1}{2\mu} \|g - g'\|_*^2.$$

$$f(x) \leq f(y) + \frac{1}{2\mu} \langle g, x - y \rangle + \frac{1}{2\mu} \|g - g'\|_*^2$$

Proof: Define $\phi(z) = f(z) - \langle g, z \rangle$

$$\partial\phi(z) = \partial f(z) - g$$

Thus $\partial\phi(y) = \partial f(y) - g$. Thus. $0 \in \partial\phi(y)$
 $(g \in \partial f(y))$

Based on Subgradient optimality condition

$$f(x^*) = \min_x f(x) \Leftrightarrow 0 \in \partial f(x^*)$$

Therefore, y is the minimizer of $\phi(z)$ ①

What's more, $\partial\phi(x) = \partial f(x) - g$

And $g' \in \partial f(x)$

Thus, $g' - g \in \partial\phi(x)$ ②

Hence, we can write

$$\phi(y) = \min_{z \in \text{dom}\phi} \phi(z) \quad (\text{based on ①})$$

$$\geq \min_{z \in \text{dom}\phi} \left(\phi(x) + \langle g' - g, z - x \rangle + \frac{\mu}{2} \|z - x\|^2 \right) \quad (\text{strongly-convex})$$

Subgradient of x (based on ②)

$$\geq \inf_{z \in \mathbb{R}^d} \left(\phi(x) + \langle g' - g, z - x \rangle + \frac{\mu}{2} \|z - x\|^2 \right)$$

$$= \phi(x) - \frac{1}{2\mu} \|g' - g\|_*^2 \quad (\text{discussed below})$$

The last equality is because:

$$f(x) = \frac{1}{2} \|x\|^2 \text{ 's conjugate is } f^*(\theta) = \frac{1}{2} \|\theta\|_*^2$$

Recall:

$$\text{the dual norm } \|z\|_* = \sup_{\|x\| \leq 1} z^T x$$

$$\text{the Fenchel conjugate } f^*(\theta) = \sup_{x \in \mathbb{R}^d} \langle \theta, x \rangle - f(x)$$

Thus,

$$\frac{1}{2} \|g' - g\|_*^2 = \sup_{z \in \mathbb{R}^d} \langle g' - g, z - x \rangle - \frac{1}{2} \|z - x\|^2$$

Multiply $-\frac{1}{u}$:

$$\begin{aligned} -\frac{1}{u} \|g' - g\|_*^2 &= \inf_{z \in \mathbb{R}^d} \frac{1}{u} \|z - x\|^2 + \frac{1}{u} \langle g' - g, z - x \rangle \\ &= \inf_{z \in \mathbb{R}^d} \langle g' - g, z - x \rangle + \frac{u}{2} \|z - x\|^2 \end{aligned}$$

This is where the last equality comes from.

Corollary 7.6. Let $f : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ μ -strongly convex with respect to a norm $\|\cdot\|$. Let $x^* = \operatorname{argmin}_x f(x)$. Then, for all $x \in \operatorname{dom} f$, and $g \in \partial f(x)$, we have

$$f(x) - f(x^*) \leq \frac{1}{2\mu} \|g\|_*^2.$$

In words, the above lemma says that an upper bound to the suboptimality gap is proportional to the squared norm of the subgradient.

7.2.2 An Explicit Regret bound

Now, we state a Lemma quantifying the intuition on the "stability" of the predictions.

Lemma 7.7. With the notation and assumptions of Lemma 7.1, assume that F_t is proper and λ_t -strongly convex w.r.t. $\|\cdot\|$, and ℓ_t proper and convex. Also, assume that $\partial\ell_t(\mathbf{x}_t)$ is non-empty. Then, we have

$$F_t(\mathbf{x}_t) - F_{t+1}(\mathbf{x}_{t+1}) + \ell_t(\mathbf{x}_t) \leq \frac{\|\mathbf{g}_t\|_*^2}{2\lambda_t} + \psi_t(\mathbf{x}_{t+1}) - \psi_{t+1}(\mathbf{x}_{t+1}),$$

for all $\mathbf{g}_t \in \partial\ell_t(\mathbf{x}_t)$.

Show the stability

Proof: First we review Lemma 7.1.

Lemma 7.1. Let $V \subseteq \mathbb{R}^d$ be closed and non-empty. Denote by $F_t(\mathbf{x}) = \psi_t(\mathbf{x}) + \sum_{i=1}^{t-1} \ell_i(\mathbf{x})$. Assume that $\operatorname{argmin}_{\mathbf{x} \in V} F_t(\mathbf{x})$ is not empty and set $\mathbf{x}_t \in \operatorname{argmin}_{\mathbf{x} \in V} F_t(\mathbf{x})$. Then, for any \mathbf{u} , we have

$$\sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) = \psi_{T+1}(\mathbf{u}) - \min_{\mathbf{x} \in V} \psi_1(\mathbf{x}) + \sum_{t=1}^T [F_t(\mathbf{x}_t) - F_{t+1}(\mathbf{x}_{t+1}) + \ell_t(\mathbf{x}_t)] + F_{T+1}(\mathbf{x}_{T+1}) - F_{T+1}(\mathbf{u}).$$

This providing an equality for the regret, which is surprising.

Then we prove Lemma 7.7.

$$\begin{aligned} & F_t(\mathbf{x}_t) - F_{t+1}(\mathbf{x}_{t+1}) + \ell_t(\mathbf{x}_t) \\ &= (F_t(\mathbf{x}_t) + \ell_t(\mathbf{x}_t)) - (F_t(\mathbf{x}_{t+1}) + \ell_t(\mathbf{x}_{t+1})) + \psi_t(\mathbf{x}_{t+1}) - \psi_{t+1}(\mathbf{x}_{t+1}) \end{aligned}$$

because of:

$$F_{t+1}(\mathbf{x}) = \psi_{t+1}(\mathbf{x}) + \sum_{i=1}^t \ell_i(\mathbf{x}) = F_t(\mathbf{x}) - \psi_t(\mathbf{x}) + \psi_{t+1}(\mathbf{x}) + \ell_t(\mathbf{x})$$

$$\begin{aligned}
& (F_t(x_t) + l_t(x_t)) - (F_t(x_{t+1}) + l_t(x_{t+1})) + \psi_t(x_{t+1}) - \psi_{t+1}(x_{t+1}) \\
& \leq (F_t(x_t) + l_t(x_t)) - (F_t(x_t^*) + l_t(x_t^*)) + \psi_t(x_{t+1}) - \psi_{t+1}(x_{t+1}) \\
& \leq \frac{1}{2\lambda_t} \|g_t'\|_*^2 + \psi_t(x_{t+1}) - \psi_{t+1}(x_{t+1}) \quad (\text{By Corollary 7.6})
\end{aligned}$$

Corollary 7.6. Let $f : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ μ -strongly convex with respect to a norm $\|\cdot\|$. Let $x^* = \operatorname{argmin}_x f(x)$. Then, for all $x \in \operatorname{dom} f$, and $g \in \partial f(x)$, we have

$$f(x) - f(x^*) \leq \frac{1}{2\mu} \|g\|_*^2.$$

where

$$g_t' \in \partial(F_t(x_t) + l_t(x_t))$$

but we want find
 $g_t \in \partial l_t(x_t)$

What's more, $x_t = \operatorname{argmin}_{x \in V} F_t(x)$ (in our setting)

By Theorem 6.10:

Theorem 6.10. Let $f : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ proper. Then $x^* \in \operatorname{argmin}_{x \in \mathbb{R}^d} f(x)$ iff $0 \in \partial f(x^*)$.

Proof. We have that

$$x^* \in \operatorname{argmin}_{x \in \mathbb{R}^d} f(x) \Leftrightarrow \forall y \in \mathbb{R}^d, f(y) \geq f(x^*) = f(x^*) + \langle 0, y - x^* \rangle \Leftrightarrow 0 \in \partial f(x^*).$$

□

Thus, $0 \in \partial(F_t(x_t) + i_V(x_t))$

By Theorem 2.19, we can select $g_t' \in \partial l_t(x_t)$

Theorem 2.19 ([Bauschke and Combettes, 2011, Theorem 18.5]). Let $(f_i)_{i \in I}$ be a finite set of convex functions from \mathbb{R}^d to $(-\infty, +\infty]$ and suppose $x \in \bigcap_{i \in I} \operatorname{dom} f_i$ and f_i continuous at x . Set $F = \max_{i \in I} f_i$ and let $A(x) = \{i \in I \mid f_i(x) = F(x)\}$ the set of the active functions. Then

$$\partial F(x) = \operatorname{conv} \bigcup_{i \in A(x)} \partial f_i(x).$$

function composition's subgradient

Let's see some immediate applications of FTRL.

Corollary 7.8. Let ℓ_t a sequence of convex loss functions. Let $\psi : V \rightarrow \mathbb{R}$ a μ -strongly convex function w.r.t. $\|\cdot\|$. Set the sequence of regularizers as $\psi_t(\mathbf{x}) = \frac{1}{\eta_{t-1}}(\psi(\mathbf{x}) - \min_{\mathbf{z}} \psi(\mathbf{z}))$, where $\eta_{t+1} \leq \eta_t$, $t = 1, \dots, T$. Then, FTRL guarantees

$$\sum_{t=1}^T \ell(\mathbf{x}_t) - \sum_{t=1}^T \ell_t(\mathbf{u}) \leq \frac{\psi(\mathbf{u}) - \min_{\mathbf{x} \in V} \psi(\mathbf{x})}{\eta_{T-1}} + \frac{1}{2\mu} \sum_{t=1}^T \eta_{t-1} \|\mathbf{g}_t\|_*^2,$$

for all $\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t)$. Moreover, if the functions ℓ_t are L -Lipschitz, setting $\eta_{t-1} = \frac{\alpha \sqrt{\mu}}{L \sqrt{t}}$ we get

$$\sum_{t=1}^T \ell(\mathbf{x}_t) - \sum_{t=1}^T \ell_t(\mathbf{u}) \leq \left(\frac{\psi(\mathbf{u}) - \min_{\mathbf{x}} \psi(\mathbf{x})}{\alpha} + \alpha \right) \frac{L \sqrt{T}}{\sqrt{\mu}}.$$

Proof: This corollary is immediate from Lemma 7.1, 7.7

Lemma 7.1. Let $V \subseteq \mathbb{R}^d$ be closed and non-empty. Denote by $F_t(\mathbf{x}) = \psi_t(\mathbf{x}) + \sum_{i=1}^{t-1} \ell_i(\mathbf{x})$. Assume that $\arg\min_{\mathbf{x} \in V} F_t(\mathbf{x})$ is not empty and set $\mathbf{x}_t \in \arg\min_{\mathbf{x} \in V} F_t(\mathbf{x})$. Then, for any \mathbf{u} , we have

$$\sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) = \psi_{T+1}(\mathbf{u}) - \min_{\mathbf{x} \in V} \psi_1(\mathbf{x}) + \sum_{t=1}^T [F_t(\mathbf{x}_t) - F_{t+1}(\mathbf{x}_{t+1}) + \ell_t(\mathbf{x}_t)] + F_{T+1}(\mathbf{x}_{T+1}) - F_{T+1}(\mathbf{u}).$$

$$\leq \underbrace{\psi_{T+1}(\mathbf{u}) - \min_{\mathbf{x} \in V} \psi_1(\mathbf{x})}_{\text{by Lemma 7.1}} + \underbrace{F_{T+1}(\mathbf{x}_{T+1}) - F_{T+1}(\mathbf{u})}_{\leq 0 \text{ by assumption}}$$

set $\Psi_T = \psi_{T+1}$

$$+ \sum_{t=1}^T \left[\frac{\|\mathbf{g}_t\|_*^2}{\sum \mathbf{u}} + \underbrace{\psi_t(\mathbf{x}_{t+1}) - \psi_{t+1}(\mathbf{x}_{t+1})}_{\text{by Lemma 7.7}} \right]$$

≤ 0 by assumption

Lemma 7.7. With the notation and assumptions of Lemma 7.1 assume that F_t is proper and λ_t -strongly convex w.r.t. $\|\cdot\|$, and ℓ_t proper and convex. Also, assume that $\partial \ell_t(\mathbf{x}_t)$ is non-empty. Then, we have

$$F_t(\mathbf{x}_t) - F_{t+1}(\mathbf{x}_{t+1}) + \ell_t(\mathbf{x}_t) \leq \frac{\|\mathbf{g}_t\|_*^2}{2\lambda_t} + \psi_t(\mathbf{x}_{t+1}) - \psi_{t+1}(\mathbf{x}_{t+1}),$$

for all $\mathbf{g}_t \in \partial \ell_t(\mathbf{x}_t)$.

\leq the right-hand side.

7.3 FTRL with Linearized Losses

In linear case, the same bound with OMD
We consider the case in which the losses are linear,

i.e., $l_t(x) = \langle g_t, x \rangle$, $t=1, \dots, T$, we have

$$\begin{aligned} x_{t+1} &\in \underset{x \in V}{\operatorname{argmin}} \psi_{t+1}(x) + \sum_{i=1}^t \langle g_i, x \rangle \\ &= \underset{x \in V}{\operatorname{argmax}} \left\langle -\sum_{i=1}^t g_i, x \right\rangle - \psi_{t+1}(x) \end{aligned}$$

Denote by $\psi_{v,t}(x) = \psi_t(x) + i_v(x)$

Now, we assume $\psi_{v,t}$ to be proper, convex and closed. using the Theorem 5.5, we have that $x_{t+1} \in \partial \psi_{v,t+1}^* \left(-\sum_{i=1}^t g_i \right)$

Theorem 5.5 ([Rockafellar, 1970] Corollary 23.5.1 and Theorem 23.5). Let $f : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ be convex, proper, and closed. Then

1. $x \in \partial f^*(\theta)$ iff $\theta \in \partial f(x)$.
2. $\langle \theta, x \rangle - f(x)$ achieves its supremum in x at $x = x^*$ iff $x^* \in \partial f^*(\theta)$.

Moreover, if $\psi_{v,t+1}$ is strongly convex, we know that

$\psi_{v,t+1}^*$ is differentiable and we get

$$x_{t+1} = \nabla \psi_{v,t+1}^* \left(-\sum_{i=1}^t g_i \right)$$

$$f(\theta) = \sup_x \langle \theta, x \rangle - f(x)$$

In turn, the update can be written in the following way

$$\theta_{t+1} = \theta_t - g_t$$

$$x_{t+1} = D\psi_{V,t+1}^*(\theta_{t+1})$$

As shown in Figure below.

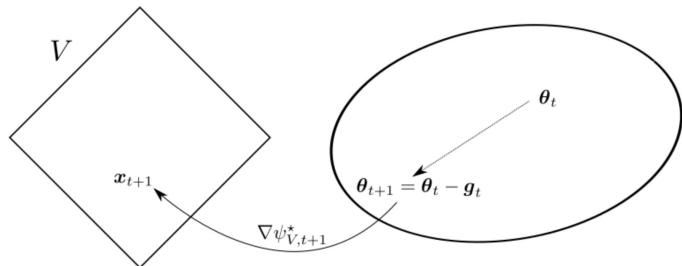


Figure 7.1: Dual mapping for FTRL with linear losses.

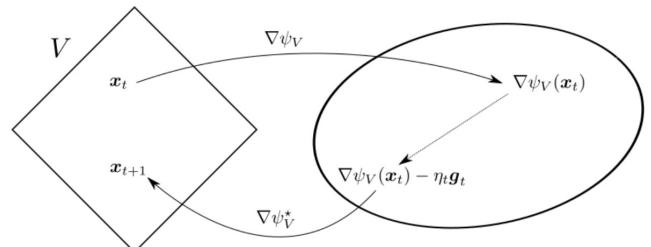


Figure 6.4: OMD update in terms of duality mappings.

Compare to Online Minor Descent in a similar way:

$$\theta_{t+1} = D\psi(x_t) - \eta_t g_t$$

$$x_{t+1} = D\psi_V^*(\theta_{t+1})$$

Differences:

- In OMD, the state is kept in x_t , so we need to transform it into a dual variable before making the update and then back to the primal variable.
- In FTRL with linear losses, the state is kept directly in the dual space, updated and then transformed in the primal variable. The primal variable is only used to predict, but not directly in the update.
- In OMD, the samples are weighted by the learning rates that is typically decreasing $\downarrow \eta_t$
- In FTRL with linear losses, all the subgradients have the same weight, but the regularizer is typically increasing over time. \uparrow regularizer

Algorithm 6.1 Online Mirror Descent

Require: Non-empty closed convex $V \subseteq X \subseteq \mathbb{R}^d$, $\psi : X \rightarrow \mathbb{R}$ strictly convex and continuously differentiable on $\text{int } X$, $x_1 \in V$ such that ψ is differentiable in $x_1, \eta_1, \dots, \eta_T > 0$

- 1: **for** $t = 1$ **to** T **do**
 - 2: Output x_t
 - 3: Receive $\ell_t : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ and pay $\ell_t(x_t)$
 - 4: Set $g_t \in \partial \ell_t(x_t)$
 - 5: $x_{t+1} = \operatorname{argmin}_{x \in V} \langle g_t, x \rangle + \frac{1}{\eta_t} B_\psi(x; x_t)$
 - 6: **end for**
-

FRTL

Algorithm 7.2 Follow-the-Regularized-Leader on Linearized Losses

Require: Convex, closed, and non-empty set $V \subseteq \mathbb{R}^d$, a sequence of regularizers $\psi_1, \dots, \psi_T : \mathbb{R}^d \rightarrow (-\infty, +\infty]$

- 1: **for** $t = 1$ **to** T **do**
 - 2: Output $x_t \in \operatorname{argmin}_{x \in V} \psi_t(x) + \sum_{i=1}^{t-1} \langle g_i, x \rangle$
 - 3: Receive $\ell_t : \mathbb{R}^d \rightarrow (-\infty, +\infty]$ and pay $\ell_t(x_t)$
 - 4: Set $g_t \in \partial \ell_t(x_t)$
 - 5: **end for**
-

7.3.1 FTRL with Linearized Losses Can Be Equivalent to OMD

In certain cases, OMD and FTRL can be equivalent.

For example, consider that $V = X = \text{dom } \psi$

The output of OMD is:

$$x_{t+1} = \underset{x}{\operatorname{argmin}} \langle \eta g_t, x \rangle + B_\psi(x; x_t)$$

Recall Bregman divergence:

$$B_\psi(x; y) = \psi(x) - \psi(y) - \langle \nabla \psi(y), x - y \rangle$$

We assume that $x_{t+1} \in \text{int dom } \psi$ for $t=1, \dots, T$

This implies that $\eta g_t + \nabla \psi(x_{t+1}) - \nabla \psi(x_t) = 0$

That is: $\nabla \psi(x_{t+1}) = \nabla \psi(x_t) - \eta g_t$

Assume $x_* = \min_{x \in V} \psi(x)$, we have

$$\nabla \psi(x_{t+1}) = -\eta \sum_{i=1}^t g_i$$

On the other hand, consider FTRL, with linearized losses with regularizers $\psi_t = \frac{1}{\eta} \psi$, then

$$x_{t+1} = \arg \min_x \frac{1}{\eta} \psi(x) + \sum_{i=1}^t \langle g_i, x \rangle$$

$$= \arg \min_x \psi(x) + \eta \sum_{i=1}^t \langle g_i, x \rangle$$

Assume $x_{t+1} \in \text{int dom } \psi$.

this implies that $\nabla \psi(x_{t+1}) = -\eta \sum_{i=1}^t g_i$

Thus the predictions of FTRL and OMD are the same.

This equivalence immediately gives us some intuition on the role of ψ in both algorithm: The same function is inducing the Bregman divergence, that is our similarity measure, and is the regularizer in FTRL. Moreover, the inverse of the growth rate of the regularizers in FTRL takes the role of the learning rate in OMD.

Example 7.10. Consider $\psi(x) = \frac{1}{2} \|x\|_2^2$ and $V = \mathbb{R}^d$, then it satisfies the conditions above to have the predictions of OMD equal to the ones of FTRL.