

A Survey of Algorithms and Analysis for Adaptive Online Learning

H. Brendan McMaha
(2017)

Presented by Jiin Seo
July 1, 2019

Outline

1. Introduction
2. FTRL Family
3. A General Analysis Technique
4. Additional Regularization Terms and Composite Objectives
5. Application to Specific Algorithms

Outline

1. Introduction

2. FTRL Family

3. A General Analysis Technique

4. Additional Regularization Terms and Composite Objectives

5. Application to Specific Algorithms

1. Introduction

Online Convex Optimization

- ▶ On each round $t \in \{1, 2, \dots\}$, select a point $x_t \in \mathbb{R}^n$.
For a convex loss function f_t ,

$$\text{Regret}(x^*) := \text{Regret}_T(x^*, f_t) \equiv \sum_{t=1}^T f_t(x_t) - \sum_{t=1}^T f_t(x^*)$$

- ▶ Main interest :

$$\text{Regret}(\mathcal{X}) \equiv \sup_{x^* \in \mathcal{X}} \text{Regret}(x^*)$$

- ▶ $\text{Regret}(x^*)$: frequently bounded by a function of $\|x^*\|$

Outline

1. Introduction

2. FTRL Family

3. A General Analysis Technique

4. Additional Regularization Terms and Composite Objectives

5. Application to Specific Algorithms

2. FTRL Family

Follow-the-Leader(FTL)

$$x_{t+1} = \arg \min_{x \in \mathbb{R}^n} f_{1:t}(x)$$

- ▶ f_t : strongly convex functions \Rightarrow sublinear regret

Follow-the-Regularized-Leader(FTRL)

- ▶ $r_t(x) \geq 0$: smoothing regularizer (adaptive member)

$$x_1 \in \arg \min_{x \in \mathbb{R}^n} r_0(x)$$
$$x_{t+1} = \arg \min_{x \in \mathbb{R}^n} (f_{1:t}(x) + r_{0:t}(x)) \quad , \text{ for } t = 1, 2, \dots$$

- ▶ Flexible decisions

- ① exact $f_{1:t}(x)$ or efficient lower bounds $\bar{f}_{1:t}(x)$
- ② r_t are minimized at fixed stationary point x_1 or current x_t

2. FTRL Family

Linearization of convex ftn

- ▶ For convex f_t and subgradient $g_t \in \partial f_t(x_t)$,

$$f_t(x_t) - f_t(x^*) \leq g_t \cdot (x_t - x^*), \forall x^*$$

- ▶ Let $\bar{f}_t(x) = g_t \cdot x$: (Linearization)

$$\text{Regret}(x^*, f_t) \leq \text{Regret}(x^*, \bar{f}_t)$$

- ▶ Update

$$x_{t+1} = \arg \min_x (g_{1:t} \cdot x + r_{0:t}(x))$$

2. FTRL Family

Regularization in FTRL

① FTRL-Centered

- ▶ Each r_t is minimized at a fixed point. ($x_1 = \arg \min_x r_0(x)$)
- ▶ $r_{0:t}$: prox-function

② FTRL-Proximal

- ▶ Each r_t is minimized at x_t .
- ▶ r_t : incremental proximal regularizers

2. FTRL Family

Notation and Definitions

- ▶ $\psi : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ an **extended-value convex function**, if

$$\psi(\theta x + (1 - \theta)y) \leq \theta\psi(x) + (1 - \theta)\psi(y), \quad \theta \in (0, 1)$$

$\text{dom } \psi \equiv \{x : \psi(x) < \infty\}$ is the convex set.

- ▶ ψ is **proper** if

$$\exists x \in \mathbb{R}^n \text{ s.t. } \psi(x) < +\infty \quad \text{and} \quad \forall x \in \mathbb{R}^n, \psi(x) > -\infty$$

- ▶ extended-value proper convex functions = "convex functions"
- ▶ $g (\in \partial\psi(x))$ is **subgradient** of ψ , if

$$\forall y \in \mathbb{R}^n, \psi(y) \geq \psi(x) + g \cdot (y - x)$$

3. A General Analysis Technique

Notation and Definitions

- ▶ **Convex conjugate** (or Fenchel conjugate) of $\psi : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$:

$$\psi^*(g) \equiv \sup_x [g \cdot x - \psi(x)]$$

- ▶ **Dual norm** :

$$\|x\|_* \equiv \sup_{y: \|y\| \leq 1} x \cdot y$$

- ▶ $\|\cdot\|_{(t)}$, $\|\cdot\|_{(t),*}$: functions of the round t

2. FTRL Family

Notation and Definitions

- ▶ $\psi : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ is σ -strongly convex w.r.t. a norm $\|\cdot\|$, if

$$\forall g \in \partial\psi(x), \quad \psi(y) \geq \psi(x) + g \cdot (y - x) + \frac{\sigma}{2} \|y - x\|^2, \forall x, y \in \mathbb{R}^n.$$

- ▶ A function ψ is σ -strongly-smooth w.r.t. a norm $\|\cdot\|$, if it is differentiable and for all x, y we have

$$\psi(y) \leq \psi(x) + \nabla\psi(x) \cdot (y - x) + \frac{\sigma}{2} \|y - x\|^2$$

Lemma 9 Let ψ be closed and convex. Then,

ψ is σ -strongly convex w.r.t. norm $\|\cdot\|$.

$\Leftrightarrow \psi^*$ is $\frac{1}{\sigma}$ -strongly smooth with respect to the dual norm $\|\cdot\|_*$.

- ▶ fundamental duality between strongly convex and strongly smooth functions

2. FTRL Family

Setting1

- ▶ Algorithm :

$$x_1 \in \arg \min_{x \in \mathbb{R}^n} r_0(x)$$
$$x_{t+1} = \arg \min_{x \in \mathbb{R}^n} \underbrace{(f_{1:t}(x) + r_{0:t}(x))}_{h_{0:t}(x)}, \text{ for } t = 1, 2, \dots$$

- ▶ $r_t(x) \geq 0$, $f_t : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$, f_t, r_t : convex.
- ▶ $h_{0:t} = r_{0:t} + f_{1:t}$, $\text{dom}(h_{0:t}) \neq \emptyset$, $\partial f_t(x_t) \neq \emptyset$

2. FTRL Family

Theorem 1. General FTRL Bound

- ▶ Setting1
- ▶ r_t are chosen s.t. $h_{0:t} + f_{t+1} = r_{0:t} + f_{1:t+1}$ is 1-strongly convex w.r.t. $\|\cdot\|_{(t)}$.
- ▶ Then, $\forall x^* \in \mathbb{R}^n$ and $\forall T > 0$,

$$\text{Regret}_T(x^*) \leq r_{0:T-1}(x^*) + \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t-1),*}^2$$

Theorem 2. FTRL-Proximal Bound

- ▶ Setting1
- ▶ r_t are chosen s.t. $h_{0:t} = r_{0:t} + f_{1:t}$ is 1-strongly convex w.r.t. $\|\cdot\|_{(t)}$.
- ▶ r_t are proximal. (i.e. x_t is a minimizer of r_t)
- ▶ Then, $\forall x^* \in \mathbb{R}^n$ and $\forall T > 0$,

$$\text{Regret}_T(x^*) \leq r_{0:T}(x^*) + \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t),*}^2, \quad g_t \in \partial f_t(x_t)$$

2. FTRL Family

Non-Adaptive Algorithms

- ▶ Setting 1
- ▶ r_0 : fixed non-adaptive regularizer, 1-strongly convex w.r.t. $\|\cdot\|$
- ▶ $r_t(x) = 0$, $t \geq 1$
- ▶ Then, $\forall t$, $\|x\|_{(t),*} = \|x\|_*$ and so both thms provide the identical statement

$$\text{Regret}(x^*) \leq r_0(x^*) + \frac{1}{2} \sum_{t=1}^T \|g_t\|_*^2$$

2. FTRL Family

Incorporating a Feasible Set

- ▶ $x \in \mathbb{R}^n \rightarrow x_t \in \mathcal{X}$: convex feasible set

$$I_{\mathcal{X}}(x) = \begin{cases} 0 & , x \in \mathcal{X} \\ \infty & , \text{otherwise} \end{cases}$$

- ▶ $r_0 \equiv r_0 + I_{\mathcal{X}}$

$$\arg \min_{x \in \mathbb{R}^n} f_{1:t}(x) + r_{0:t}(x) + I_{\mathcal{X}}(x) = \arg \min_{x \in \mathcal{X}} f_{1:t}(x) + r_{0:t}(x)$$

- ▶ $I_{\mathcal{X}}$ appears in regret bound.

$$(\text{Thm1}) : \text{Regret}_T(x^*) \leq I_{\mathcal{X}}(x^*) + r_{0:T-1}(x^*) + \frac{1}{2} \sum_{t=1}^T \|\mathbf{g}_t\|_{(t-1),*}^2$$

Outline

1. Introduction
2. FTRL Family
- 3. A General Analysis Technique**
4. Additional Regularization Terms and Composite Objectives
5. Application to Specific Algorithms

3. A General Analysis Technique

Inductive Lemmas

Lemma 5 (Strong FTRL Lemma).

- ▶ f_t : seq. of arbitrary (possibly non-convex) loss functions
- ▶ r_t : arbitrary non-nega regularization functions, s.t.
- ▶ $h_{0:t}(x) \equiv f_{1:t}(x) + r_{0:t}(x)$
- ▶ $x_{t+1} = \arg \min_x h_{0:t}(x)$: well defined

Then, the algorithm that selects these x_t achieves

$$\text{Regret}(x^*) \leq r_{0:T}(x^*) + \sum_{t=1}^T [h_{0:t}(x_t) - h_{0:t}(x_{t+1}) - r_t(x_t)]$$

(Standard FTRL Lemma).

$$\text{Regret}(x^*) \leq r_{0:T}(x^*) + \sum_{t=1}^T [f_t(x_t) - f_t(x_{t+1})].$$

3. A General Analysis Technique

Proof of Lemma 5 (Strong FTRL Lemma).

► FTL algorithm against h_1, \dots, h_T

$$\begin{aligned} & \sum_{t=1}^T (f_t(x_t) + r_t(x_t)) - f_{1:T}(x^*) - r_{0:T}(x^*) \\ &= \sum_{t=1}^T h_t(x_t) - h_{0:T}(x^*) = \sum_{t=1}^T (h_{0:t}(x_t) - h_{0:t-1}(x_t)) - h_{0:T}(x^*) \\ &\leq \sum_{t=1}^T (h_{0:t}(x_t) - h_{0:t-1}(x_t)) - h_{0:T}(x_{T+1}) \quad (x_{T+1} \text{ minimizes } h_{0:T}) \\ &\leq \sum_{t=1}^T (h_{0:t}(x_t) - h_{0:t}(x_{t+1})) \quad (, -h_0(x_1) = -r_0(x_1) \leq 0) \end{aligned}$$

■ Proof of Lemma5

3. A General Analysis Technique

Inductive Lemmas

Corollary 6

- ▶ The same conditions as Lemma 5
- ▶ $f_t(x) = g_t \cdot x_t$: Linear loss

Then,

$$h_{0:t}(x_t) - h_{0:t}(x_{t+1}) - r_t(x_t) = r_{0:t}^*(-g_{1:t}) - r_{0:t-1}^*(-g_{1:t-1}) + g_t \cdot x_t$$

which implies

$$\text{Regret}(x^*) \leq r_{0:T}(x^*) + \sum_{t=1}^T [r_{0:t}^*(-g_{1:t}) - r_{0:t-1}^*(-g_{1:t-1}) + g_t \cdot x_t]$$

3. A General Analysis Technique

Proof of Corollary 6.

① $h_{0:t}(x_{t+1}) \rightarrow$ Fenchel conjugate:

$$\begin{aligned} r_{0:t}^*(-g_{1:t}) &= \max_x [-g_{1:t} \cdot x - r_{0:t}(x)] && : \text{Def. of Fenchel conjugate} \\ &= -\min_x [g_{1:t} \cdot x + r_{0:t}(x)] \\ &= -h_{0:t}(x_{t+1}) \end{aligned}$$

②

$$\begin{aligned} h_{0:t}(x_t) - r_t(x_t) &= g_{1:t} \cdot x_t + r_{0:t}(x_t) - r_t(x_t) \\ &= g_{1:t-1} \cdot x_t + r_{0:t-1}(x_t) + g_t \cdot x_t \\ &= h_{0:t-1}(x_t) + g_t \cdot x_t \\ &= -r_{0:t-1}^*(-g_{1:t-1}) + g_t \cdot x_t \quad (\text{using 1 with } t \rightarrow t-1) \end{aligned}$$

■ Proof of Coro6

3. A General Analysis Technique

Tools from Convex Analysis

Lemma 7.

- ▶ $\phi_1 : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$: convex function s.t. $\exists x_1 = \arg \min_x \phi_1(x)$
- ▶ ψ : convex function s.t. $\phi_2(x) = \phi_1(x) + \psi(x)$ is strongly convex w.r.t. norm $\|\cdot\|$
- ▶ $x_2 = \arg \min_x \phi_2(x)$.

Then,

$$\forall b \in \partial\psi(x_1), \quad \|x_1 - x_2\| \leq \|b\|_*$$

and

$$\forall x', \quad \phi_2(x_1) - \phi_2(x') \leq \frac{1}{2} \|b\|_*^2$$

3. A General Analysis Technique

Proof of Lemma 7.

- 1 $\phi_2(x) = \phi_1(x) + \psi(x)$: 1-strongly convex w.r.t. $\|\cdot\|$
- 2 Construct ϕ'_1 : 1-strongly convex w.r.t. $\|\cdot\|$, but minimizer is x_1 .

$$\begin{aligned}\phi'_1(x) &\equiv \phi_1(x) + \psi(x) - b \cdot x, && \text{ : 1-strongly convex w.r.t. } \|\cdot\| \\ 0 &\in \partial(\psi(x) - b \cdot x)|_{x=x_1} && \text{ : } x_1 = \arg \min \phi'_1(x)\end{aligned}$$

$$\Rightarrow \phi_2(x) = \phi'_1(x) + b \cdot x$$

- 3 Applying Lemma 16 to ϕ'_1 and ϕ_2

$$\forall x', \quad \phi_2(x_1) - \phi_2(x') \leq \phi_2(x_1) - \phi_2(x_2) \leq \frac{1}{2} \|b\|_*^2$$

$$\text{and } \|x_1 - x_2\| \leq \|b\|_*$$

■ Proof of Lemma7

3. A General Analysis Technique

Tools from Convex Analysis

Lemma 16

- ▶ $\phi_1 : \mathbb{R}^n \rightarrow \mathbb{R}$: strongly convex w.r.t. norm $\|\cdot\|$
- ▶ $\phi_2(x) := \phi_1(x) + b \cdot x$ for $b \in \mathbb{R}^n$.
- ▶ $x_1 = \arg \min_x \phi_1(x)$, $x_2 = \arg \min_x \phi_2(x)$

Then,

$$\phi_2(x_1) - \phi_2(x_2) \leq \frac{1}{2} \|b\|_{\star}^2, \quad \text{and} \quad \|x_1 - x_2\| \leq \|b\|_{\star}$$

3. A General Analysis Technique

Proof of Lemma 16.

① Fenchel conjugate:

$$\begin{aligned} -\phi_1^*(0) &= -\max_x [0 \cdot x - \phi_1(x)] = \min_x [\phi_1(x)] = \phi_1(x_1) \\ -\phi_1^*(-b) &= -\max_x [-b \cdot x - \phi_1(x)] = \min_x [b \cdot x + \phi_1(x)] \\ &= b \cdot x_2 + \phi_1(x_2) \end{aligned}$$

② ϕ_1^* : strongly-smooth (by Lemma 9)

$$\begin{aligned} \phi_1^*(-b) &\leq \phi_1^*(0) + \nabla \phi_1^*(0) \cdot (-b - 0) + \frac{1}{2} \| -b - 0 \|^2 \\ &= \phi_1^*(0) + x_1 \cdot (-b - 0) + \frac{1}{2} \|b\|_*^2 \quad (, x_1 = \nabla \phi_1^*(0)) \end{aligned}$$

3. A General Analysis Technique

Proof of Lemma 16.

③ 1st bound : Combining 1 and 2

$$\begin{aligned}\phi_2(x_1) - \phi_2(x_2) &= \phi_1(x_1) + b \cdot x_1 - \phi_1(x_2) - b \cdot x_2 \\ &= -\phi_1^*(0) + b \cdot x_1 + \phi_1^*(-b) \\ &\leq -\phi_1^*(0) + b \cdot x_1 + \phi_1^*(0) + x_1 \cdot (-b) + \frac{1}{2} \|b\|_*^2 \\ &= \frac{1}{2} \|b\|_*^2\end{aligned}$$

④ 2nd bound : using Lemma 15

$$\nabla \phi_1^*(0) = x_1, \text{ and } \nabla \phi_1^*(-b) = x_2$$

$$\|x_1 - x_2\| = \|\nabla \phi_1^*(0) - \nabla \phi_1^*(-b)\| \leq \|b\|_*$$

■ Proof of Lemma16

3. A General Analysis Technique

Tools from Convex Analysis

Lemma 15.

- ▶ $\psi : 1$ -strongly convex w.r.t. $\|\cdot\|$
($\Leftrightarrow \psi^* : 1$ -strongly smooth w.r.t. $\|\cdot\|_*$)

Then,

$$\|\nabla\psi^*(z) - \nabla\psi^*(z')\| \leq \|z - z'\|_*$$

and

$$\arg \min_x [g \cdot x + \psi(x)] = \nabla\psi^*(-g)$$

3. A General Analysis Technique

Tools from Convex Analysis

Corollary 8.

- ▶ $\phi_1 : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$: a convex function s.t. $\exists x_1 = \arg \min_x \phi_1(x)$
- ▶ ψ, Ψ : convex functions s.t. $\phi_2(x) = \phi_1(x) + \psi(x) + \Psi(x)$: strongly convex w.r.t. norm $\|\cdot\|$.
- ▶ $x_2 = \arg \min_x \phi_2(x)$.

Then, $\forall b \in \partial\psi(x_1)$ and $\forall x'$,

$$\phi_2(x_1) - \phi_2(x') \leq \frac{1}{2} \|b\|_*^2 + \Psi(x_1) - \Psi(x_2)$$

3. A General Analysis Technique

Proof of Corollary 8.

$$\textcircled{1} \quad x'_2 := \arg \min_x [\phi_1(x) + \psi(x)]$$

$$\phi_1(x_1) + \psi(x_1) - \phi_1(x'_2) - \psi(x'_2) \leq \frac{1}{2} \|b\|_*^2 \quad (\text{by Lemma 7})$$

$$\phi_1(x'_2) + \psi(x'_2) \leq \phi_1(x_2) + \psi(x_2) \quad (\text{by Def})$$

$\textcircled{2}$

$$\begin{aligned} \forall x', \quad & \phi_2(x_1) - \phi_2(x') \\ &= \phi_1(x_1) + \psi(x_1) + \Psi(x_1) - \phi_1(x_2) - \psi(x_2) - \Psi(x_2) \\ &\leq \phi_1(x_1) + \psi(x_1) + \Psi(x_1) - \phi_1(x'_2) - \psi(x'_2) - \Psi(x_2) \\ &\leq \frac{1}{2} \|b\|_*^2 + \Psi(x_1) - \Psi(x_2) \end{aligned}$$

■ Proof of Corollary 8

3. A General Analysis Technique

Theorem 1. General FTRL Bound

- ▶ Setting1
- ▶ r_t are chosen s.t. $h_{0:t} + f_{t+1} = r_{0:t} + f_{1:t+1}$ is 1-strongly convex w.r.t. $\|\cdot\|_{(t)}$.
- ▶ Then, $\forall x^* \in \mathbb{R}^n$ and $\forall T > 0$,

$$\text{Regret}_T(x^*) \leq r_{0:T-1}(x^*) + \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t-1),*}^2$$

3. A General Analysis Technique

General FTRL including FTRL-Centered (Proof of Theorem 1)

① Bound the stability terms : $h_{0:t}(x_t) - h_{0:t}(x_{t+1})$

▶ For fixed t ,

$$\phi_1(x) := h_{0:t-1}(x) \quad , x_t = \arg \min_x \phi_1(x)$$

$$\phi_2(x) := h_{0:t-1}(x) + f_t(x) \quad : \text{1-strongly-convex w.r.t. } \|\cdot\|_{(t-1)}$$

▶ applying Lemma 7 :

$$\phi_2(x_t) - \phi_2(x_{t+1}) \leq \frac{1}{2} \|g_t\|_{(t-1),*}^2 \quad \text{for } g_t \in \partial f_t(x_t)$$

$$\begin{aligned} h_{0:t}(x_t) - h_{0:t}(x_{t+1}) - r_t(x_t) &= \phi_2(x_t) + r_t(x_t) - \phi_2(x_{t+1}) - r_t(x_{t+1}) - r_t(x_t) \\ &\leq \frac{1}{2} \|g_t\|_{(t-1),*}^2 \quad , r_t(x) \geq 0 \end{aligned}$$

3. A General Analysis Technique

General FTRL including FTRL-Centered (Proof of Theorem 1)

- ② Sum these up for all rounds. & Use Lemma 5.

$$\begin{aligned} \text{Regret}(x^*) &\leq r_{0:T}(x^*) + \sum_{t=1}^T [h_{0:t}(x_t) - h_{0:t}(x_{t+1}) - r_t(x_t)] \\ &\leq r_{0:T}(x^*) + \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t-1),*}^2 \\ &= r_{0:T-1}(x^*) + \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t-1),*}^2, \quad r_T(x) := 0 \text{ (WLOG)} \end{aligned}$$

■ Proof of Thm1

3. A General Analysis Technique

Theorem 2. FTRL-Proximal Bound

- ▶ Setting 1
- ▶ r_t are chosen s.t. $h_{0:t} = r_{0:t} + f_{1:t}$ is 1-strongly convex w.r.t. $\|\cdot\|_{(t)}$.
- ▶ r_t are proximal. (i.e. x_t is a minimizer of r_t)
- ▶ Then, $\forall x^* \in \mathbb{R}^n$ and $\forall T > 0$,

$$\text{Regret}_T(x^*) \leq r_{0:T}(x^*) + \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t),*}^2, \quad g_t \in \partial f_t(x_t)$$

3. A General Analysis Technique

FTRL-Proximal (Proof of Theorem 2)

① Bound the stability terms : $h_{0:t}(x_t) - h_{0:t}(x_{t+1})$

▶ For fixed t ,

$$\phi_1(x) := f_{1:t-1}(x) + r_{0:t}(x) = h_{0:t}(x) - f_t(x)$$

$$x_t = \arg \min_x \phi_1(x) , \quad x_{t+1} = \arg \min_x [\phi_1(x) + f_t(x)]$$

▶ applying Lemma 7 and $\phi_2(x) = \phi_1(x) + f_t(x) = h_{0:t}(x)$

$$\begin{aligned} h_{0:t}(x_t) - h_{0:t}(x_{t+1}) - r_t(x_t) & \\ & \leq h_{0:t}(x_t) - h_{0:t}(x_{t+1}) \quad , (r_t(x) \geq 0) \\ & = \phi_1(x_t) + f_t(x_t) - \phi_1(x_{t+1}) - f_t(x_{t+1}) \\ & \leq \frac{1}{2} \|g_t\|_{(t),*}^2 \quad , (\text{Lemma 7}) \end{aligned}$$

3. A General Analysis Technique

FTRL-Proximal (Proof of Theorem 2)

- ② Sum up these for all rounds. & Use Lemma 5.

$$\begin{aligned}\text{Regret}(x^*) &\leq r_{0:T}(x^*) + \sum_{t=1}^T [h_{0:t}(x_t) - h_{0:t}(x_{t+1}) - r_t(x_t)] \\ &\leq r_{0:T}(x^*) + \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t),*}^2\end{aligned}$$

■ Proof of Thm2

Outline

1. Introduction
2. FTRL Family
3. A General Analysis Technique
4. Additional Regularization Terms and Composite Objectives
5. Application to Specific Algorithms

4. Additional Regularization Terms and Composite Objectives

generalized FTRL

- ▶ Ψ : non-negative convex function
- ▶ $\alpha_t \geq 0$ for $t \geq 1$: non-increasing in t
- ▶ $x_1 = \arg \min_x \Psi(x) = \arg \min_x r_0(x)$, w.l.o.g. $\Psi(x_1) = 0$
- ▶ Objective function :

$$h_t(x) = g_t \cdot x + \underbrace{\alpha_t \Psi(x)}_{\text{: additional regularization}} + r_t(x)$$

- ▶ FTRL update

$$x_{t+1} = \arg \min_x h_{0:t}(x) = \arg \min_x g_{1:t} \cdot x + \alpha_{1:t} \Psi(x) + r_{0:t}(x)$$

4. Additional Regularization Terms and Composite Objectives

generalized FTRL

non-smooth regularizer Ψ

- ▶ L_2 penalty : $\Psi(x) = \|x\|_2^2$
- ▶ L_1 penalty : $\Psi(x) = \|x\|_1 \Rightarrow$ sparse solutions
- ▶ indicator function : $\Psi(x) = I_{\mathcal{X}}(x) \Rightarrow$ feasible solutions

4. Additional Regularization Terms and Composite Objectives

generalized FTRL

① Ψ terms in the objective g_t

- ▶ $f_t^\Psi(x) \equiv g_t \cdot x + \alpha_t \Psi(x)$
- ▶ Regret

$$\text{Regret}(x^*, f^\Psi) \equiv \sum_{t=1}^T f_t^\Psi(x_t) - f_t^\Psi(x^*)$$

② Ψ terms in the regularizer r_t

- ▶ $f_t(x) = g_t \cdot x$
- ▶ Regret

$$\text{Regret}(x^*, g_t) \equiv \sum_{t=1}^T g_t \cdot x_t - g_t \cdot x^*$$

4. Additional Regularization Terms and Composite Objectives

Theorem 10. FTRL-Proximal Bounds for Composite Objectives

- ▶ Ψ : a non-negative convex function,
 $x_1 = \arg \min_x \Psi(x), \quad \Psi(x_1) = 0$
 - ▶ $\alpha_t \geq 0$ be a non-increasing seq. of constants
 - ▶ Setting1
 - ▶ $h_t(x) := g_t \cdot x + \alpha_t \Psi(x) + r_t(x)$
 - ▶ r_t are chosen s.t. $h_{0:t} : 1$ -strongly-convex w.r. to $\|\cdot\|_{(t),*}$
 - ▶ r_t : proximal (x_t is a global minimizer of r_t)
- ① Case1 : $f_t^\Psi(x) = g_t \cdot x + \alpha_t \Psi(x)$

$$\text{Regret}(x^*, f^\Psi) \leq r_{0:T}(x^*) + \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t),*}^2$$

- ② Case2 : $f_t(x) = g_t \cdot x$

$$\text{Regret}(x^*, g_t) \leq r_{0:T}(x^*) + \alpha_{1:T} \Psi(x^*) + \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t),*}^2$$

4. Additional Regularization Terms and Composite Objectives

Proof of Theorem 10.

① Bound the stability terms : $h_{0:t}(x_t) - h_{0:t}(x_{t+1})$

► For fixed t ,

$$\begin{aligned}\phi_1(x) &:= h_{0:t-1}(x) + r_t(x) & , x_t &= \arg \min_x \phi_1(x) \\ \phi_2(x) &= \phi_1(x) + g_t \cdot x + \alpha_t \Psi(x) = h_{0:t}(x) & , x_{t+1} &= \arg \min_x \phi_2(x)\end{aligned}$$

► applying Corollary 8

$$\begin{aligned}h_{0:t}(x_t) - h_{0:t}(x_{t+1}) - r_t(x_t) & \\ & \leq h_{0:t}(x_t) - h_{0:t}(x_{t+1}) \quad , (r_t(x) \geq 0) \\ & = \phi_1(x_t) + f_t(x_t) + \alpha_t \Psi(x_t) - \phi_1(x_{t+1}) - f_t(x_{t+1}) - \alpha_t \Psi(x_{t+1}) \\ & \leq \frac{1}{2} \|g_t\|_{(t),*}^2 + \alpha_t \Psi(x_t) - \alpha_t \Psi(x_{t+1}) \quad , (\text{Corollary 8})\end{aligned}$$

4. Additional Regularization Terms and Composite Objectives

Proof of Theorem 10.

- ② Sum up these for all rounds. & Use Lemma 5.

$$\begin{aligned} & \sum_{t=1}^T \alpha_t \Psi(x_t) - \alpha_t \Psi(x_{t+1}) \\ &= \alpha_1 \Psi(x_1) - \alpha_T \Psi(x_{T+1}) + \sum_{t=2}^T [\alpha_t \Psi(x_t) - \alpha_{t-1} \Psi(x_t)] \\ &\leq 0 \quad (\Psi(x) \geq 0, \alpha_t \leq \alpha_{t-1}, \text{ and } \Psi(x_1) = 0) \end{aligned}$$

$$\begin{aligned} & \sum_{t=1}^T [h_{0:t}(x_t) - h_{0:t}(x_{t+1}) - r_t(x_t)] \\ &\leq \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t),*}^2 + \sum_{t=1}^T [\alpha_t \Psi(x_t) - \alpha_t \Psi(x_{t+1})] \\ &\leq \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t),*}^2 \end{aligned}$$

4. Additional Regularization Terms and Composite Objectives

Proof of Theorem 10.

③ Case1 : $f_t(x) \leftarrow g_t \cdot x + \alpha_t \Psi(x)$ and $r_t \leftarrow r_t$

$$\text{Regret}(x^*, f^\Psi) \leq r_{0:T}(x^*) + \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t),*}^2 \quad (\text{Lemma 5})$$

④ Case2 : $f_t(x) \leftarrow g_t \cdot x$ and $r_t(x) \leftarrow \alpha_t \Psi(x) + r_t(x)$

$$\text{Regret}(x^*, g_t) \leq r_{0:T}(x^*) + \alpha_{1:T} \Psi(x^*) + \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t),*}^2 \quad (\text{Lemma 5})$$

■ Proof of Thm10

Outline

1. Introduction
2. FTRL Family
3. A General Analysis Technique
4. Additional Regularization Terms and Composite Objectives
5. Application to Specific Algorithms

5. Application to Specific Algorithms

Lemma 3.

- ▶ Setting 1.
- ▶ $r_{0:t}(x) \geq r_{0:t-1}(x)$,
- ▶ $r_{0:t}^*(x) \leq r_{0:t-1}^*(x)$, where $r_{0:t}^*$ is the convex-conjugate of $r_{0:t}$.
- ▶ If r_t is σ_t -strongly convex w.r.t. a norm $\|\cdot\|$ for $\sigma_t \geq 0$,
- ▶ Then,

$$r_{0:t} : \begin{cases} \sigma_{0:t}\text{-strongly convex w.r.t. } \|\cdot\| \\ 1\text{-strongly-convex w.r.t. } \|x\|_{(t)} = \sqrt{\sigma_{0:t}}\|x\| \quad , \|x\|_{(t),*} = \frac{1}{\sqrt{\sigma_{0:t}}}\|x\| \end{cases}$$

Learning rate schedule (η_t)

$$\eta_t = \frac{1}{\sigma_{0:t}} \Leftrightarrow \sigma_t = \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \quad , \text{ where } \sigma_0 = \frac{1}{\eta_0}$$

(Assumption) : $f_t(x) = g_t \cdot x$

5. Application to Specific Algorithms

Non-adaptive algorithm

- ▶ OGD with constant η

$$(\text{OGD}) : x_{t+1} = x_t - \eta g_t \quad , \eta > 0 \Leftrightarrow x_{t+1} = -\eta g_{1:t}$$

- ▶ FTRL with $r_0(x) = \frac{1}{2\eta} \|x\|_2^2$ and $r_t(x) = 0$ for $t \geq 1$

$$x_{t+1} = \arg \min_x g_{1:t} \cdot x + \frac{1}{2\eta} \|x\|_2^2 \Leftrightarrow x_{t+1} = -\eta g_{1:t}$$

- ▶ Regret

$$\begin{aligned} \text{Regret}_T(x^*) &\leq r_{0:T-1}(x^*) + \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t-1),*}^2 \\ &= \frac{1}{2\eta} \|x^*\|_2^2 + \frac{1}{2} \sum_{t=1}^T \eta \|g_t\|_2^2 \quad , \text{using } \|x\|_{(t),*} = \sqrt{\eta} \|x\|_2 \end{aligned}$$

5. Application to Specific Algorithms

Non-adaptive algorithm

- ▶ (Assumption) $\|x^*\|_2 \leq R, \|g_t\|_2 \leq G$
we want to minimize regret after T' rounds.
- ▶ Choosing $\eta = \frac{R}{G\sqrt{T'}}$ minimize bound of $\text{Regret}_T(x^*)$

$$\begin{aligned}\text{Regret}_T(x^*) &\leq \frac{RG}{2}\sqrt{T'} + \frac{RG}{2}\frac{T}{\sqrt{T}} \\ &= RG\sqrt{T} \quad (\text{, when } T = T')\end{aligned}$$

- ▶ bound : $\mathcal{O}(\sqrt{T})$ when $T = \mathcal{O}(T')$

5. Application to Specific Algorithms

Non-adaptive algorithm

- ▶ OGD with Constant η and Feasible Set \mathcal{X}
- ▶ (Assumption) $\|x^*\|_2 \leq R$
- ▶ $r_0(x) = \frac{1}{2\eta}\|x\|_2^2 + l_{\mathcal{X}}(x)$

$$x_{t+1} = \arg \min_{x \in \mathbb{R}^n} g_{1:t} \cdot x + \frac{1}{2\eta}\|x\|_2^2 + l_{\mathcal{X}}(x) = \arg \min_{x \in \mathcal{X}} g_{1:t} \cdot x + \frac{1}{2\eta}\|x\|_2^2$$

- ▶ equiv. to 2-step update (Lazy projection)

$$\Leftrightarrow \begin{cases} u_{t+1} = \arg \min_{x \in \mathbb{R}^n} g_{1:t} \cdot x + \frac{1}{2\eta}\|x\|^2 \\ x_{t+1} = \Pi_{\mathcal{X}}(u_{t+1}) \end{cases}, \Pi_{\mathcal{X}}(u) \equiv \arg \min_{x \in \mathcal{X}} \|x - u\|_2$$

Appendix C.2.

5. Application to Specific Algorithms

Dual Averaging

- ▶ adaptive FTRL-Centered algorithm with linearized loss function
- ▶ bound : $\mathcal{O}(\sqrt{T})$, $\forall T$
- ▶ $r_t(x) = \frac{\sigma_t}{2} \|x\|_2^2$, $\sigma_t \geq 0$
- ▶ $r_{0:t} : 1$ -strongly-convex w.r.t. the norm $\|x\|_{(t)} = \sqrt{\sigma_{0:t}} \|x\|_2$
- ▶ dual norm : $\|x\|_{(t),*} = \frac{1}{\sqrt{\sigma_{0:t}}} \|x\|_2 = \sqrt{\eta_t} \|x\|_2$

$$\begin{aligned} \forall T, \text{Regret}_T(x^*) &\leq r_{0:T-1}(x^*) + \frac{1}{2} \sum_{t=1}^T \|\mathbf{g}_t\|_{(t-1),*}^2 \\ &= \frac{1}{2\eta_{T-1}} \|x^*\|_2^2 + \frac{1}{2} \sum_{t=1}^T \eta_{t-1} \|\mathbf{g}_t\|_2^2 \end{aligned}$$

5. Application to Specific Algorithms

Dual Averaging

- ▶ (Assumption) $\|x^*\|_2 \leq R, \|g_t\|_2 \leq G$
- ▶ Choosing $\eta_t = \frac{R}{\sqrt{2}G\sqrt{t+1}}$ and using $\sum_{t=1}^T \frac{1}{\sqrt{t}} \leq 2\sqrt{T}$

$$\forall T, \text{Regret}_T(x^*) \leq \frac{\sqrt{2}}{2} \left(R + \frac{\|x^*\|_2^2}{R} \right) G\sqrt{T}$$

- ▶ When $\|x^*\| \leq R$, we have $\text{Regret} \leq \sqrt{2}RG\sqrt{T}$
- ▶ RDA(Regularized Dual Averaging algorithm)
linearized f_t + non-smooth regularization (adding $\lambda_1\|x\|_1 + \lambda_2\|x\|_2^2$ to r_0)

5. Application to Specific Algorithms

FTRL-Proximal

- ▶ $\mathcal{X} \subseteq \{x \mid \|x\|_2 \leq R\}$
- ▶ $r_0(x) = I_{\mathcal{X}}(x)$, $r_t(x) = \frac{\sigma_t}{2} \|x - x_t\|_2^2, t > 1$
- ▶ $r_{0:t}(x^*) \leq \frac{\sigma_{1:t}}{2} (2R)^2, \forall x^* \in \mathcal{X}$, since each $x_t \in \mathcal{X}$
- ▶ $r_{0:t}$ is 1-strongly-convex w.r.t. the norm $\|x\|_{(t)} = \sqrt{\sigma_{1:t}} \|x\|_2$
- ▶ dual norm : $\|x\|_{(t),*} = \frac{1}{\sqrt{\sigma_{1:t}}} \|x\|_2$

$$\begin{aligned} \forall T, \text{Regret}_T(x^*) &\leq r_{0:T}(x^*) + \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t),*}^2 \\ &= \frac{1}{2\eta_T} (2R)^2 + \frac{1}{2} \sum_{t=1}^T \eta_t \|g_t\|_2^2, \text{ where } \eta_t = \frac{1}{\sigma_{1:t}} \end{aligned}$$

- ▶ (Assumption) $\|x^*\| \leq R, \quad \|g_t\|_2 \leq G$
- ▶ Choosing $\eta_t = \frac{\sqrt{2}R}{G\sqrt{t}}$.

$$(x^*) \leq 2\sqrt{2}RG\sqrt{T}$$

- ▶ One would hope x_t is closer to x^* than 0

5. Application to Specific Algorithms

FTRL-Proximal with Diagonal Matrix Learning Rates

- ▶ 1-dimensional + AdaGrad FTRL-Proximal algorithm : adaptive to seq. of g_t
- ▶ Let $r_0 = I_{\mathcal{X}}$ with $\mathcal{X} = [-R, R]$ and $\eta_t = \frac{\sqrt{2}R}{\sqrt{\sum_{s=1}^t g_s^2}}$.

$$\begin{aligned}\forall x^* \in \mathcal{X}, \quad \text{Regret}(x^*) &\leq \frac{1}{2\eta_T} (2R)^2 + \frac{1}{2} \sum_{t=1}^T \eta_t \|g_t\|^2 \\ &\leq 2\sqrt{2}R \sqrt{\sum_{t=1}^T g_t^2}\end{aligned}$$

- ▶ Fully adaptive bound : do not need to know T and G
- ▶ R is chosen in advance.

5. Application to Specific Algorithms

FTRL-Proximal with Diagonal Matrix Learning Rates

- ▶ n-dimensional AdaGrad-style algorithm
- ▶ $\mathcal{X} \subseteq [-R_\infty, R_\infty]^n$
- ▶ per-coordinate basis ; For coordinate i , $\eta_{t,i} = \frac{\sqrt{2}R_\infty}{\sqrt{\sum_{s=1}^t g_{s,i}^2}}$
- ▶ $r_0 = I_{\mathcal{X}}$, $r_t(x) = \frac{1}{2} \left\| Q_t^{\frac{1}{2}} (x - x_t) \right\|_2^2, t \geq 1$
 $Q_t = \text{diag}(\sigma_{t,i}) : \text{P.S.D.}, \quad \sigma_{t,i} = \eta_{t,i}^{-1} - \eta_{t-1,i}^{-1}$
- ▶ $r_{0:t}$ is 1-strongly-convex w.r.t. the norm $\|x\|_{(t)} = \left\| (Q_{1:t})^{\frac{1}{2}} x \right\|_2$
- ▶ dual norm : $\|g\|_{(t),*} = \left\| (Q_{1:t})^{-\frac{1}{2}} g \right\|_2$

$$\begin{aligned} \forall x^* \in \mathcal{X}, \quad \text{Regret}_T(x^*) &\leq r_{0:T}(x^*) + \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t),*}^2 \\ &= r_{0:T}(x^*) + \frac{1}{2} \sum_{t=1}^T \left\| (Q_{1:t})^{-\frac{1}{2}} g_t \right\|_2^2 \end{aligned}$$

5. Application to Specific Algorithms

FTRL-Proximal with Diagonal Matrix Learning Rates

- ▶ Bound of diagonal FTRL-Proximal is at least as good than bound of FTRL proximal.
- ▶ Less useful because of $\mathcal{O}(n^2)$ requirements
- ▶ Update simplifies to

$$x_{t+1} = \arg \min_{x \in \mathcal{X}} (g_{1:t} - a_{1:t}) \cdot x + \sum_{i=1}^n \frac{1}{2\eta_{t,i}} x_i^2 \quad \text{where} \quad a_t = \sigma_t x_t$$

5. Application to Specific Algorithms

$$\sum_{t=1}^T 1/\sqrt{t} \leq 2\sqrt{T}$$

Lemma 4. For any non-negative real numbers a_1, a_2, \dots, a_n

$$\sum_{i=1}^n \frac{a_i}{\sqrt{\sum_{j=1}^i a_j}} \leq 2\sqrt{\sum_{i=1}^n a_i}$$

Proof)

- ▶ $n = 1$: Clear!
- ▶ Assume the lemma holds for some $n - 1$,

$$\begin{aligned} \sum_{i=1}^n \frac{x_i}{\sqrt{\sum_{j=1}^i x_j}} &\leq 2\sqrt{\sum_{i=1}^{n-1} x_i} + \frac{x_n}{\sqrt{\sum_{i=1}^n x_i}}, \quad (Z \equiv \sum_{i=1}^n x_i \text{ and } x \equiv x_n) \\ &= 2\sqrt{Z - x} + \frac{x}{\sqrt{Z}} \quad (\text{RHS is maximized at } x = 0) \\ &\leq 2\sqrt{Z} \end{aligned}$$

- ▶ using derivative of RHS w.r.t. $x = \frac{-1}{\sqrt{Z-x}} + \frac{1}{\sqrt{Z}} < 0$ for $x > 0$

5. Application to Specific Algorithms

AdaGrad Dual Averaging

- ▶ off-by-one difference
- ▶ per-coordinate learning rates $\eta_t \approx \frac{R}{\sqrt{G^2 + \sum_{s=1}^t g_s^2}}$, $|g_t| \leq G$
- ▶ almost fully adaptive

5. Application to Specific Algorithms

AdaGrad Dual Averaging with the Entropic Regularizer ???

- ▶ probability distribution selection problems
 $x_t \in \Delta_n$ with $\Delta_n = \{x \mid \sum_{i=1}^n x_i = 1 \text{ and } x_i \geq 0\}$
- ▶ (assumption) $\|g_t\|_\infty \leq G_\infty$
- ▶ Using entropic regularizer

$$h(x) = I_\Delta(x) + \log n + \sum_{i=1}^n x_i \log x_i$$

- ▶ adaptive regularization schedule:

$$r_{0:t}(x) = \frac{1}{\eta_t} h(x) \quad \text{where} \quad \eta_t = \frac{\sqrt{\log n}}{\sqrt{G_\infty^2 + \sum_{s=1}^t \|g_s\|_\infty^2}}, \quad t \geq 0$$

- ▶ $x_1 = (1/n, \dots, 1/n)$, where $h(x) = 0$
- ▶ h is maximized at the corners of Δ_n

5. Application to Specific Algorithms

AdaGrad Dual Averaging with the Entropic Regularizer

- ▶ h is 1-strongly-convex w.r.t. L_1
- ▶ $r_{0:t}$ is 1-strongly-convex w.r.t. $\|x\|_{(t)} = \frac{1}{\sqrt{\eta t}} \|x\|_1$, and $\|g\|_{(t),*}^2 = \eta t \|g\|_\infty^2$

$$\begin{aligned} \text{Regret}(x^*) &\leq r_{0:T-1}(x^*) + \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t-1),*}^2 \\ &\leq \frac{\log n}{\eta_{T-1}} + \frac{1}{2} \sum_{t=1}^T \eta_{t-1} \|g_t\|_\infty^2 \quad \text{since } \forall t, \|g_t\|_\infty \leq G_\infty \\ &\leq \frac{\log n}{\eta_{T-1}} + \frac{\sqrt{\log n}}{2} \sum_{t=1}^T \frac{\|g_t\|_\infty^2}{\sqrt{\sum_{s=1}^t \|g_s\|_\infty^2}} \quad \text{Lemma 4 and } \|g_T\|_\infty \\ &\leq 2G_\infty \sqrt{T \log n} \end{aligned}$$

5. Application to Specific Algorithms

Strongly Convex Functions

- ▶ Suppose f_t is 1 -strongly-convex w.r.t. a norm $\|\cdot\|$, and $r_t(x) \equiv 0 \forall t$.
:Exact FTL
- ▶ (Assumption) $\|g_t\| \leq G$
- ▶ $\|x\|_{(t)} = \sqrt{t}\|x\|$
- ▶ $h_{0:t}(x)$ is 1 -strongly-convex w.r.t. $\|\cdot\|_{(t)}$ (by Lemma 3)
- ▶ using $\sum_{t=1}^T \frac{1}{t} \leq 1 + \log T$

$$\text{Regret}(x^*) \leq \frac{1}{2} \sum_{t=1}^T \|g_t\|_{(t),*}^2 = \frac{1}{2} \sum_{t=1}^T \frac{1}{t} \|g_t\|^2 \leq \frac{G^2}{2} (1 + \log T)$$

- ▶ \bar{f}_t, f_t : 1 -strongly-convex w.r.t. $\|\cdot\|$
$$f_t(x) \geq \underbrace{f_t(x_t) + g_t \cdot (x - x_t) + \frac{1}{2} \|x - x_t\|^2}_{\equiv \bar{f}_t}$$

5. Application to Specific Algorithms

Strongly Convex Functions

- ▶ FTL on $\bar{f}_t \Leftrightarrow$ OGD
- ▶ $x_{t+1} = \arg \min_x \bar{f}_{1:t}(x) \Leftrightarrow x_{t+1} = x_t - \frac{1}{t} g_t$