



Lecture 9. Optimistic Online Mirror Descent

Advanced Optimization (Fall 2023)

Peng Zhao zhaop@lamda.nju.edu.cn Nanjing University

Beyond the Worst-Case Analysis

- All above regret guarantees hold against the worst case
 - Matching the *minimax optimality*
 - The environment is *fully adversarial*



- However, in practice:
 - We are not always interested in the *worst-case scenario*
 - Environments can exhibit *specific patterns*: gradual change, periodicity...



We are after some more *problem-dependent* guarantees.

Beyond the Worst-Case Analysis

- Beyond the worst-case analysis, achieving more adaptive results.
 - (1) *adaptivity*: achieving better guarantees in easy problem instances;
 - (2) *robustness*: maintaining the same worst-case guarantee.



Advanced Optimization (Fall 2023)

Small-Loss Bounds for PEA

Theorem 2. Suppose that $\forall t \in [T]$ and $i \in [N], 0 \leq \ell_{t,i} \leq 1$, then Hedge with *learning rate* $\eta \in (0, 1)$ *guarantees*

$$\operatorname{Regret}_{T} \leq \frac{\ln N}{\eta} + \eta \sum_{t=1}^{T} \langle \boldsymbol{p}_{t}, \boldsymbol{\ell}_{t} \rangle,$$

which can further imply
$$\operatorname{Regret}_{T} \leq \frac{1}{1-\eta} \left(\frac{\ln N}{\eta} + \eta L_{T,i^{\star}} \right) = \mathcal{O}\left(\sqrt{L_{T,i^{\star}} \log N} + \log N \right),$$

by setting $\eta = \min\left\{\frac{1}{2}, \sqrt{\frac{\ln N}{L_{T,i^{\star}}}}\right\}.$

- When $L_{T,i^*} = \mathcal{O}(T)$, it can recover the *minimax* $\mathcal{O}(\sqrt{T \log N})$ guarantee.
- When $L_{T,i^*} = \mathcal{O}(1)$, the regret bound is $\mathcal{O}(\log N)$, which is independent of T!

 \mathbf{T}

Small-Loss Bounds for PEA

• Addressing the unpleasant dependence on L_{T,i^*} via *self-confident tuning*.

Theorem 4. Suppose that $\forall t \in [T]$ and $i \in [N], 0 \leq \ell_{t,i} \leq 1$, then Hedge with adaptive learning rate $\eta_t = \sqrt{\frac{\ln N}{L_t + 1}}$ guarantees $\operatorname{Regret}_T \leq 8\sqrt{(L_{T,i^*} + 1)\ln N} + 3\ln N$ $= \mathcal{O}(\sqrt{L_{T,i^*}\log N} + \log N),$ where $L_t = \sum_{s=1}^t \langle \boldsymbol{p}_s, \boldsymbol{\ell}_s \rangle$ is cumulative loss the learner suffered at time t.

Key Analysis in Self-confident Tuning

Proof. From the potential-based proof, we already know that

$$\sum_{t=1}^{T} \langle \boldsymbol{p}_t, \boldsymbol{\ell}_t \rangle - L_{T,i^*} \leq \sqrt{(L_{T-1}+1) \ln N} + \sum_{t=1}^{T} \eta_{t-1} \langle \boldsymbol{p}_t, \boldsymbol{\ell}_t \rangle$$

$$\leq \sqrt{(L_{T-1}+1)\ln N} + \sum_{t=1}^{T} \frac{\langle \boldsymbol{p}_t, \boldsymbol{\ell}_t \rangle}{\sqrt{\sum_{s=1}^{t-1} \langle \boldsymbol{p}_s, \boldsymbol{\ell}_s \rangle + 1}} \begin{pmatrix} \eta_{t-1} = \sqrt{\frac{\ln N}{L_{t-1}+1}} \end{pmatrix} \\ (L_{t-1} = \sum_{s=1}^{t-1} \langle \boldsymbol{p}_s, \boldsymbol{\ell}_s \rangle) \end{pmatrix}$$

How to bound this term?

 \Box This is actually a common structure to handle.

Small-Loss Bound for PEA: Proof

Proof. From previous potential-based proof, we already known that

$$\sum_{t=1}^{T} \langle \boldsymbol{p}_t, \boldsymbol{\ell}_t \rangle - L_{T,i^{\star}} \leq \sqrt{(L_{T-1}+1) \ln N} + \sum_{t=1}^{T} \frac{\langle \boldsymbol{p}_t, \boldsymbol{\ell}_t \rangle}{\sqrt{\sum_{s=1}^{t-1} \langle \boldsymbol{p}_s, \boldsymbol{\ell}_s \rangle + 1}}$$

Lemma 2. Let
$$a_1, a_2, \ldots, a_T$$
 be non-negative real numbers. Then

$$\sum_{t=1}^{T} \frac{a_t}{\sqrt{1 + \sum_{s=1}^{t-1} a_s}} \leq 4\sqrt{1 + \sum_{t=1}^{T} a_t} + \max_{t \in [T]} a_t$$

$$\longrightarrow \qquad L_T - L_{T,i^*} \leq \sqrt{(L_{T-1} + 1) \ln N} + 4\sqrt{1 + L_T} + 1 \qquad (\ell_i \leq 1, \forall i \in [N])$$

$$\leq \sqrt{(L_T + 1) \ln N} + 4\sqrt{1 + L_T} + 1$$

Advanced Optimization (Fall 2023)

Small-Loss Bounds for OCO

Definition 4 (Small Loss). The small-loss quantity of the OCO problem (online function $f_t : \mathcal{X} \mapsto \mathbb{R}$) is defined as

$$F_T = \min_{\mathbf{x} \in \mathcal{X}} \sum_{t=1}^T f_t(\mathbf{x})$$

One essential property for small-loss bound for OCO: *self-bounding property*.

Corollary 1. For an L-smooth and non-negative function $f : \mathbb{R}^d \mapsto \mathbb{R}$, we have that $\|\nabla f(\mathbf{x})\|_2 \leq \sqrt{2Lf(\mathbf{x})}, \quad \forall \mathbf{x} \in \mathcal{X}.$

Advanced Optimization (Fall 2023)

Achieving Small-Loss Bound

• We show that under the *self-bounding condition*, OGD can yield the desired small-loss regret bound.

$$\mathbf{x}_{t+1} = \Pi_{\mathbf{x}\in\mathcal{X}}[\mathbf{x}_t - \eta_t \nabla f_t(\mathbf{x}_t)]$$

Theorem 6 (Small-loss Bound). Assume that f_t is L-smooth and non-negative for all $t \in [T]$, when setting $\eta_t = \frac{D}{\sqrt{1+G_t}}$, the regret of OGD to any comparator $\mathbf{u} \in \mathcal{X}$ is bounded as

$$\operatorname{Regret}_{T} = \sum_{t=1}^{T} f_{t}(\mathbf{x}_{t}) - \sum_{t=1}^{T} f_{t}(\mathbf{u}) \le \mathcal{O}\left(\sqrt{1+F_{T}}\right)$$

where $G_t = \sum_{s=1}^t \|\nabla f_s(\mathbf{x}_s)\|_2^2$ is the empirical cumulative gradient norm.

Proof.
$$\sum_{t=1}^{T} f_t(\mathbf{x}_t) - \sum_{t=1}^{T} f_t(\mathbf{u}) \le \sum_{t=1}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t)\|_2^2 + \sum_{t=1}^{T} \frac{1}{2\eta_t} \left(\|\mathbf{u} - \mathbf{x}_t\|_2^2 - \|\mathbf{u} - \mathbf{x}_{t+1}\|_2^2 \right)$$

$$\sum_{t=1}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t)\|_2^2 = D \sum_{t=2}^{T} \frac{\|\nabla f_t(\mathbf{x}_t)\|_2^2}{\sqrt{1+G_t}} + G^2 \le 2D \sqrt{1+\sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t)\|_2^2} + G^2$$
$$(\eta_1 \triangleq 1) \qquad (G_t = \sum_{s=1}^{t} \|\nabla f_s(\mathbf{x}_s)\|_2^2)$$

Lemma 1. Let a_1, a_2, \dots, a_T be non-negative real numbers. Then $\sum_{t=1}^{T} \frac{a_t}{\sqrt{1 + \sum_{s=1}^{t} a_s}} \leq 2\sqrt{1 + \sum_{t=1}^{T} a_t}$

Advanced Optimization (Fall 2023)

Several Remarks

- Remark 1: about the non-negative assumption When the online functions are non-negative, it is possible to redefine the small-loss quantity by incorporating each-round minimal function value.
- Remark 2: about the smoothness assumption

Smoothness is necessary to obtain small-loss regret bound by the first-order method (can be proved by the online-to-batch conversion and existing lower bounds for deterministic optimization).

• Remark 3: take care of the way dealing with variance term In OGD here we use Lemma 1, while in PEA Hedge for PEA we use Lemma 2.

Towards a Unified Framework

- Previous small-loss bounds seem to be ad-hoc designed.
- Is there a *unified framework* to get problem-dependent bounds?

• A reflection: Adaptive to the niceness of the environment. What does a "nice" environment actually mean?

☐ The environment is *"predictable"*.

Outline

- Optimistic Online Mirror Descent
 - A Unified Framework
 - Small-Loss bound
 - Gradient-Variance bound
 - Gradient-Variation bound

Optimistic Online Learning

• **Intuition**: what if the environment is *"predictable"*?

 \implies We can to some extent "*guess*" the next move.



If it is within the same season and no extreme weather

Guess: It still seems to rain on Friday?

Our Previous Efforts

• Review OMD update:

OMD updates:
$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_t \langle \mathbf{x}, \nabla f_t(\mathbf{x}_t) \rangle + \mathcal{D}_{\psi}(\mathbf{x}, \mathbf{x}_t) \right\}$$

- This framework provides a unified framework for the algorithmic design and regret analysis for the worst-case scenarios.
- We aim to encode *"predictable"* information in the update such that the overall algorithm can adapt to the niceness of environments.

Optimistic Online Mirror Descent

• We introduce a sequence of *optimistic vector* $\{M_t\}_{t=1}^T$ serving as the available predictable information of *future gradients*.

Optimistic Online Mirror Descent

$$\mathbf{x}_{t} = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_{t} \left\langle \mathbf{M}_{t}, \mathbf{x} \right\rangle + \mathcal{D}_{\psi}(\mathbf{x}, \widehat{\mathbf{x}}_{t}) \right\}$$
$$\widehat{\mathbf{x}}_{t+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_{t} \left\langle \nabla f_{t}(\mathbf{x}_{t}), \mathbf{x} \right\rangle + \mathcal{D}_{\psi}(\mathbf{x}, \widehat{\mathbf{x}}_{t}) \right\}$$

where $M_t \in \mathbb{R}^d$ is the optimistic vector at each round.

essentially two-step mirror-descent updates

Understand Optimistic OMD

$$\mathbf{x}_{t} = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_{t} \left\langle M_{t}, \mathbf{x} \right\rangle + \mathcal{D}_{\psi}(\mathbf{x}, \widehat{\mathbf{x}}_{t}) \right\}$$
$$\widehat{\mathbf{x}}_{t+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_{t} \left\langle \nabla f_{t}(\mathbf{x}_{t}), \mathbf{x} \right\rangle + \mathcal{D}_{\psi}(\mathbf{x}, \widehat{\mathbf{x}}_{t}) \right\}$$



Optimistic OMD: Regret Analysis



The proof still relies on the *stability lemma* and the *Bregman proximal inequality*, but now it requires taking the two-step updates (with optimism) into account.

- The key is to have a proper regret decomposition.
- Due to the two-step updates, we need to incorporate optimism and intermediate decision in regret analysis.

$$\mathbf{x}_{t} = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_{t} \left\langle \mathbf{M}_{t}, \mathbf{x} \right\rangle + \mathcal{D}_{\psi}(\mathbf{x}, \widehat{\mathbf{x}}_{t}) \right\}$$
$$\widehat{\mathbf{x}}_{t+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_{t} \left\langle \nabla f_{t}(\mathbf{x}_{t}), \mathbf{x} \right\rangle + \mathcal{D}_{\psi}(\mathbf{x}, \widehat{\mathbf{x}}_{t}) \right\}$$

$$\begin{array}{c} & & \\ & & \\ \hline \end{array} \\ & = \underbrace{\langle \nabla f_t(\mathbf{x}_t) - f_t(\mathbf{u}) \leq \langle \nabla f_t(\mathbf{x}_t), \mathbf{x}_t - \mathbf{u} \rangle}_{\text{term (a)}} + \underbrace{\langle M_t, \mathbf{x}_t - \widehat{\mathbf{x}}_{t+1} \rangle}_{\text{term (b)}} + \underbrace{\langle \nabla f_t(\mathbf{x}_t), \widehat{\mathbf{x}}_{t+1} - \mathbf{u} \rangle}_{\text{term (c)}} \end{aligned}$$

Advanced Optimization (Fall 2023)

$$Proof. f_t(\mathbf{x}_t) - f_t(\mathbf{u}) \leq \underbrace{\langle \nabla f_t(\mathbf{x}_t) - M_t, \mathbf{x}_t - \widehat{\mathbf{x}}_{t+1} \rangle}_{\text{term (a)}} + \underbrace{\langle M_t, \mathbf{x}_t - \widehat{\mathbf{x}}_{t+1} \rangle}_{\text{term (b)}} + \underbrace{\langle \nabla f_t(\mathbf{x}_t), \widehat{\mathbf{x}}_{t+1} - \mathbf{u} \rangle}_{\text{term (c)}}$$

For term (a), we use the *stability lemma*.

Lemma 2 (Stability Lemma). *Consider the following updates:*

 $\begin{cases} \mathbf{x} = \arg\min_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{g}, \mathbf{x} \rangle + \mathcal{D}_{\psi}(\mathbf{x}, \mathbf{c}) \\ \mathbf{x}' = \arg\min_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{g}', \mathbf{x} \rangle + \mathcal{D}_{\psi}(\mathbf{x}, \mathbf{c}) \end{cases}$

When the regularizer $\psi : \mathcal{X} \mapsto \mathbb{R}$ *is a* λ *-strongly convex function with respect to norm* $\| \cdot \|$ *, we have*

$$\lambda \left\| \mathbf{x} - \mathbf{x}' \right\| \le \left\| \mathbf{g} - \mathbf{g}' \right\|_{\star}$$

$$\operatorname{term} (\mathbf{a}) = \langle \nabla f_t(\mathbf{x}_t) - M_t, \mathbf{x}_t - \widehat{\mathbf{x}}_{t+1} \rangle$$
$$\leq \| \nabla f_t(\mathbf{x}_t) - M_t \|_{\star} \| \mathbf{x}_t - \widehat{\mathbf{x}}_{t+1} \| \leq \eta_t \| \nabla f_t(\mathbf{x}_t) - M_t \|_{\star}^2$$

Advanced Optimization (Fall 2023)

$$Proof. f_t(\mathbf{x}_t) - f_t(\mathbf{u}) \leq \underbrace{\langle \nabla f_t(\mathbf{x}_t) - M_t, \mathbf{x}_t - \widehat{\mathbf{x}}_{t+1} \rangle}_{\text{term (a)}} + \underbrace{\langle M_t, \mathbf{x}_t - \widehat{\mathbf{x}}_{t+1} \rangle}_{\text{term (b)}} + \underbrace{\langle \nabla f_t(\mathbf{x}_t), \widehat{\mathbf{x}}_{t+1} - \mathbf{u} \rangle}_{\text{term (c)}}$$

For term (b), we adopt the *Bregman Proximal inequality*.

Lemma 3 (Bregman Proximal Inequality). Consider convex optimization problem with the following update form $\mathbf{x}_{t+1} = \underset{\mathbf{x} \in \mathcal{X}}{\arg\min} \left\{ \langle \mathbf{g}_t, \mathbf{x} \rangle + \mathcal{D}_{\psi}(\mathbf{x}, \mathbf{x}_t) \right\}.$

Then, it satisfies the following inequality for any $\mathbf{u} \in \mathcal{X}$ *:*

$$\langle \mathbf{g}_t, \mathbf{x}_{t+1} - \mathbf{u} \rangle \leq \mathcal{D}_{\psi}(\mathbf{u}, \mathbf{x}_t) - \mathcal{D}_{\psi}(\mathbf{u}, \mathbf{x}_{t+1}) - \mathcal{D}_{\psi}(\mathbf{x}_{t+1}, \mathbf{x}_t).$$

Thus, according to update rule: $\mathbf{x}_t = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_t \langle M_t, \mathbf{x} \rangle + \mathcal{D}_{\psi}(\mathbf{x}, \widehat{\mathbf{x}}_t) \right\}$

$$\texttt{term}\left(\texttt{b}\right) = \langle M_t, \mathbf{x}_t - \widehat{\mathbf{x}}_{t+1} \rangle \leq \frac{1}{\eta_t} \bigg(\mathcal{D}_{\psi}(\widehat{\mathbf{x}}_{t+1}, \widehat{\mathbf{x}}_t) - \mathcal{D}_{\psi}(\widehat{\mathbf{x}}_{t+1}, \mathbf{x}_t) - \mathcal{D}_{\psi}(\mathbf{x}_t, \widehat{\mathbf{x}}_t) \bigg)$$

Advanced Optimization (Fall 2023)

$$Proof. f_t(\mathbf{x}_t) - f_t(\mathbf{u}) \leq \underbrace{\langle \nabla f_t(\mathbf{x}_t) - M_t, \mathbf{x}_t - \widehat{\mathbf{x}}_{t+1} \rangle}_{\text{term (a)}} + \underbrace{\langle M_t, \mathbf{x}_t - \widehat{\mathbf{x}}_{t+1} \rangle}_{\text{term (b)}} + \underbrace{\langle \nabla f_t(\mathbf{x}_t), \widehat{\mathbf{x}}_{t+1} - \mathbf{u} \rangle}_{\text{term (c)}}$$

For term (c), we also adopt the *Bregman Proximal inequality*.

Lemma 3 (Bregman Proximal Inequality). Consider convex optimization problem with the following update form $\mathbf{x}_{t+1} = \underset{\mathbf{x} \in \mathcal{X}}{\arg\min} \left\{ \langle \mathbf{g}_t, \mathbf{x} \rangle + \mathcal{D}_{\psi}(\mathbf{x}, \mathbf{x}_t) \right\}.$

Then, it satisfies the following inequality for any $\mathbf{u} \in \mathcal{X}$ *:*

$$\langle \mathbf{g}_t, \mathbf{x}_{t+1} - \mathbf{u} \rangle \leq \mathcal{D}_{\psi}(\mathbf{u}, \mathbf{x}_t) - \mathcal{D}_{\psi}(\mathbf{u}, \mathbf{x}_{t+1}) - \mathcal{D}_{\psi}(\mathbf{x}_{t+1}, \mathbf{x}_t).$$

Thus, according to update rule: $\widehat{\mathbf{x}}_{t+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_t \langle \nabla f_t(\mathbf{x}_t), \mathbf{x} \rangle + \mathcal{D}_{\psi}(\mathbf{x}, \widehat{\mathbf{x}}_t) \right\}$

$$\texttt{term}\left(\texttt{c}\right) = \langle \nabla f_t(\mathbf{x}_t), \widehat{\mathbf{x}}_{t+1} - \mathbf{u} \rangle \leq \frac{1}{\eta_t} \bigg(\mathcal{D}_{\psi}(\mathbf{u}, \widehat{\mathbf{x}}_t) - \mathcal{D}_{\psi}(\mathbf{u}, \widehat{\mathbf{x}}_{t+1}) - \mathcal{D}_{\psi}(\widehat{\mathbf{x}}_{t+1}, \widehat{\mathbf{x}}_t) \bigg)$$

$$Proof. f_t(\mathbf{x}_t) - f_t(\mathbf{u}) \leq \underbrace{\langle \nabla f_t(\mathbf{x}_t) - M_t, \mathbf{x}_t - \widehat{\mathbf{x}}_{t+1} \rangle}_{\text{term (a)}} + \underbrace{\langle M_t, \mathbf{x}_t - \widehat{\mathbf{x}}_{t+1} \rangle}_{\text{term (b)}} + \underbrace{\langle \nabla f_t(\mathbf{x}_t), \widehat{\mathbf{x}}_{t+1} - \mathbf{u} \rangle}_{\text{term (c)}}$$

Put the three terms together, we can finish the proof.

$$\begin{split} \operatorname{term}\left(\mathbf{a}\right) &\leq \eta_{t} \left\|\nabla f_{t}(\mathbf{x}_{t}) - M_{t}\right\|_{\star}^{2} \\ \operatorname{term}\left(\mathbf{b}\right) &\leq \frac{1}{\eta_{t}} \left(\mathcal{D}_{\psi}(\widehat{\mathbf{x}}_{t+1}, \widehat{\mathbf{x}}_{t}) - \mathcal{D}_{\psi}(\widehat{\mathbf{x}}_{t+1}, \mathbf{x}_{t}) - \mathcal{D}_{\psi}(\mathbf{x}_{t}, \widehat{\mathbf{x}}_{t}) \right) \\ \operatorname{term}\left(\mathbf{c}\right) &\leq \frac{1}{\eta_{t}} \left(\mathcal{D}_{\psi}(\mathbf{u}, \widehat{\mathbf{x}}_{t}) - \mathcal{D}_{\psi}(\mathbf{u}, \widehat{\mathbf{x}}_{t+1}) - \mathcal{D}_{\psi}(\widehat{\mathbf{x}}_{t+1}, \widehat{\mathbf{x}}_{t}) \right) \\ \hline \end{pmatrix} \\ \hline & = \sum f_{t}(\mathbf{x}_{t}) - f_{t}(\mathbf{u}) \leq \eta_{t} \left\|\nabla f_{t}(\mathbf{x}_{t}) - M_{t}\right\|_{\star}^{2} + \frac{1}{\eta_{t}} \left(\mathcal{D}_{\psi}(\mathbf{u}, \widehat{\mathbf{x}}_{t}) - \mathcal{D}_{\psi}(\mathbf{u}, \widehat{\mathbf{x}}_{t+1}) \right) \\ & - \frac{1}{\eta_{t}} \left(\mathcal{D}_{\psi}(\widehat{\mathbf{x}}_{t+1}, \mathbf{x}_{t}) + \mathcal{D}_{\psi}(\mathbf{x}_{t}, \widehat{\mathbf{x}}_{t}) \right) \end{split}$$

Advanced Optimization (Fall 2023)

Example: Optimistic OGD

• Consider the Euclidean regularizer $\mathcal{D}_{\psi}(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$, i.e.,

$$\mathbf{x}_{t} = \arg\min_{\mathbf{x}\in\mathcal{X}} \eta \langle \mathbf{M}_{t}, \mathbf{x} \rangle + \frac{1}{2} \|\mathbf{x} - \hat{\mathbf{x}}_{t}\|_{2}^{2}$$

$$\widehat{\mathbf{x}}_{t+1} = \arg\min_{\mathbf{x}\in\mathcal{X}} \eta \langle \nabla f_{t}(\mathbf{x}_{t}), \mathbf{x} \rangle + \frac{1}{2} \|\mathbf{x} - \hat{\mathbf{x}}_{t}\|_{2}^{2}$$

$$(quality of guess)$$

$$+ \frac{1}{2\eta} \sum_{t=1}^{T} f_{t}(\mathbf{u}) \leq \frac{1}{2\eta} \sum_{t=1}^{T} \left(\|\mathbf{u} - \hat{\mathbf{x}}_{t}\|_{2}^{2} - \|\mathbf{u} - \hat{\mathbf{x}}_{t+1}\|_{2}^{2} \right)$$

$$(telescoping term)$$

$$- \frac{1}{2\eta} \sum_{t=1}^{T} \left(\|\widehat{\mathbf{x}}_{t+1} - \mathbf{x}_{t}\|_{2}^{2} + \|\mathbf{x}_{t} - \widehat{\mathbf{x}}_{t}\|_{2}^{2} \right)$$

$$(negative term)$$

Advanced Optimization (Fall 2023)

Example: Optimistic OGD

• Consider the Euclidean regularizer $\mathcal{D}_{\psi}(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$, i.e.,:

Advanced Optimization (Fall 2023)

Optimistic OMD: Regret Analysis



- For problem-independent bounds, negative terms of OMD is usually dropped;
- For problem-dependent bounds, the *negative term* of optimistic OMD can be sometimes extremely crucial.

Optimistic OMD: Applications

• Small-Loss Bound

• Gradient-Variance Bound

Gradient-Variation Bound

Optimistic OMD: Applications

• Small-Loss Bound

• Gradient-Variance Bound

Gradient-Variation Bound

• Recall the guarantee of optimistic OGD:

$$\mathbf{x}_{t} = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_{t} \left\langle \mathbf{M}_{t}, \mathbf{x} \right\rangle + \mathcal{D}_{\psi}(\mathbf{x}, \widehat{\mathbf{x}}_{t}) \right\}$$
$$\widehat{\mathbf{x}}_{t+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_{t} \left\langle \nabla f_{t}(\mathbf{x}_{t}), \mathbf{x} \right\rangle + \mathcal{D}_{\psi}(\mathbf{x}, \widehat{\mathbf{x}}_{t}) \right\}$$

• Consider the Euclidean regularizer $\mathcal{D}_{\psi}(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$, i.e.,:

$$\sum_{t=1}^{T} f_t(\mathbf{x}_t) - \sum_{t=1}^{T} f_t(\mathbf{u}) \le \mathcal{O}\left(\sqrt{1 + \sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t) - M_t\|_2^2}\right)$$
Setting $M_t = 0 \implies \sum_{t=1}^{T} f_t(\mathbf{x}_t) - \sum_{t=1}^{T} f_t(\mathbf{u}) \le \mathcal{O}\left(\sqrt{1 + \sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t)\|_2^2}\right)$

Advanced Optimization (Fall 2023)

• Employing the *self-bounding property* of smooth and non-negative functions.

Corollary 1. For an L-smooth and non-negative function $f : \mathbb{R}^d \mapsto \mathbb{R}$, we have that $\|\nabla f(\mathbf{x})\|_2 \leq \sqrt{2Lf(\mathbf{x})}, \quad \forall \mathbf{x} \in \mathcal{X}.$

Setting $M_t = 0$ in Optimistic OMD (with Euclidean regularizer):

$$\sum_{t=1}^{T} f_t(\mathbf{x}_t) - \sum_{t=1}^{T} f_t(\mathbf{u}) \le \mathcal{O}\left(\sqrt{1 + \sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t)\|_2^2}\right) \le \mathcal{O}\left(\sqrt{1 + L\sum_{t=1}^{T} f_t(\mathbf{x}_t)}\right) \text{ (self-bounding property)}$$

$$\sum_{t=1}^{T} \operatorname{Regret}_T = \sum_{t=1}^{T} f_t(\mathbf{x}_t) - \sum_{t=1}^{T} f_t(\mathbf{u}) = \mathcal{O}\left(D\sqrt{L\sum_{t=1}^{T} f_t(\mathbf{u}) + 1} + G^2\right).$$

$$(converting trick) \qquad \Box$$

Advanced Optimization (Fall 2023)

- Since we are using optimistic OMD with a fixed step size, the algorithm requires $G_T \triangleq \sum_{t=1}^T \|\nabla f_t(\mathbf{x}_t)\|_2^2$ when achieving small-loss bound.
- This is can be rectified by the *self-confident tuning*. We can use the optimistic OMD with time-varying step sizes.

Theorem 6 (Small-loss Bound). Assume that $\psi(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|_2^2$ and f_t is L-smooth and non-negative for all $t \in [T]$, when setting $\eta_t = \frac{D}{\sqrt{1+G_t}}$ and $M_t = \mathbf{0}$, the regret of Optimistic OMD to any comparator $\mathbf{u} \in \mathcal{X}$ is bounded as

$$\operatorname{Regret}_{T} = \sum_{t=1}^{T} f_{t}(\mathbf{x}_{t}) - \sum_{t=1}^{T} f_{t}(\mathbf{u}) \le \mathcal{O}\left(\sqrt{1+F_{T}}\right),$$

where $G_t = \sum_{s=1}^t \|\nabla f_s(\mathbf{x}_s)\|_2^2$ is the empirical cumulative gradient norm.

$$\begin{aligned} \textbf{Proof.} \quad \sum_{t=1}^{T} f_t(\mathbf{x}_t) - \sum_{t=1}^{T} f_t(\mathbf{u}) &\leq \sum_{t=1}^{T} \eta_t \left\| \nabla f_t(\mathbf{x}_t) - M_t \right\|_2^2 \qquad (\text{quality of guess, term}(\mathbf{a})) \\ &+ \sum_{t=1}^{T} \frac{1}{2\eta_t} \left(\left\| \mathbf{u} - \widehat{\mathbf{x}}_t \right\|_2^2 - \left\| \mathbf{u} - \widehat{\mathbf{x}}_{t+1} \right\|_2^2 \right) \quad (\text{telescoping term, term}(\mathbf{b})) \\ &- \sum_{t=1}^{T} \frac{1}{2\eta_t} \left(\left\| \widehat{\mathbf{x}}_{t+1} - \mathbf{x}_t \right\|_2^2 + \left\| \mathbf{x}_t - \widehat{\mathbf{x}}_t \right\|_2^2 \right) \quad (\text{negative term, term}(\mathbf{c})) \end{aligned}$$

For term (a),

$$\sum_{t=1}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t) - M_t\|_2^2 = D \sum_{t=2}^{T} \frac{\|\nabla f_t(\mathbf{x}_t)\|_2^2}{\sqrt{1+G_t}} + G^2 \le 2D \sqrt{1 + \sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t)\|_2^2} + G^2$$
(self-confident tuning lemma)
$$\le D \sqrt{1 + 2L \sum_{t=1}^{T} f_t(\mathbf{x}_t)} + G^2$$
(self-bounding property)

Advanced Optimization (Fall 2023)

$$\begin{array}{l} \textbf{Proof.} \quad \texttt{term} \left(\texttt{b} \right) = \sum_{t=1}^{T} \frac{1}{2\eta_t} \left(\| \mathbf{u} - \widehat{\mathbf{x}}_t \|_2^2 - \| \mathbf{u} - \widehat{\mathbf{x}}_{t+1} \|_2^2 \right) \\ & \leq \frac{1}{2\eta_T} \sum_{t=1}^{T} \left(\| \mathbf{u} - \widehat{\mathbf{x}}_t \|_2^2 - \| \mathbf{u} - \widehat{\mathbf{x}}_{t+1} \|_2^2 \right) \left(\{\eta_1, \dots, \eta_T\} \text{ decreasing step size} \right) \\ & \leq \frac{1}{2\eta_T} \| \mathbf{u} - \widehat{\mathbf{x}}_1 \|_2^2 \qquad (\texttt{telescoping}) \\ & \leq \frac{D}{2} \sqrt{1 + 2L \sum_{t=1}^{T} f_t(\mathbf{x}_t)} + \frac{D}{2} \quad (\texttt{by def of } \eta_T = \frac{D}{\sqrt{1+G_T}} \text{ and domain boundedness}) \\ \end{array}$$

Advanced Optimization (Fall 2023)

Optimistic OMD: Applications

• Small-Loss Bound

• Gradient-Variance Bound

Gradient-Variation Bound

Definition 3 (Gradient Variance). Let *T* be the time horizon and $\mathcal{X} \subseteq \mathbb{R}^d$ be the feasible domain. For the function sequence f_1, \ldots, f_T with $f_t : \mathcal{X} \mapsto \mathbb{R}$ for $t \in [T]$, its gradient variance is defined as $\operatorname{Var}_T = \sup_{\{\mathbf{x}_1, \ldots, \mathbf{x}_T\} \in \mathcal{X}} \sum_{t=1}^T \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_T\|_2^2$ where $\boldsymbol{\mu}_T \triangleq \arg \min_{\boldsymbol{\mu}} \sum_{t=1}^T \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}\|_2^2 = \frac{1}{T} \sum_{t=1}^T \nabla f_t(\mathbf{x}_t)$.



Implicit prior on the enviornment: there exists a *latent mean gradient* $\mathbb{E}[\nabla f_t(\mathbf{x}_t)]$.

e.g. SGD (sampled from a set of data) e.g. Classification (sampled from training set)

Definition 3 (Gradient Variance). Let *T* be the time horizon and $\mathcal{X} \subseteq \mathbb{R}^d$ be the feasible domain. For the function sequence f_1, \ldots, f_T with $f_t : \mathcal{X} \mapsto \mathbb{R}$ for $t \in [T]$, its gradient variance is defined as $\operatorname{Var}_T = \sup_{\{\mathbf{x}_1, \ldots, \mathbf{x}_T\} \in \mathcal{X}} \sum_{t=1}^T \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_T\|_2^2$ where $\boldsymbol{\mu}_T \triangleq \arg \min_{\boldsymbol{\mu}} \sum_{t=1}^T \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}\|_2^2 = \frac{1}{T} \sum_{t=1}^T \nabla f_t(\mathbf{x}_t)$.

Optimistic Online Mirror Descent

$$\mathbf{x}_{t} = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_{t} \left\langle M_{t}, \mathbf{x} \right\rangle + \frac{1}{2} \|\mathbf{x} - \widehat{\mathbf{x}}_{t}\|_{2}^{2} \right\} \text{ How to choose Mt?}$$
$$\widehat{\mathbf{x}}_{t+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_{t} \left\langle \nabla f_{t}(\mathbf{x}_{t}), \mathbf{x} \right\rangle + \frac{1}{2} \|\mathbf{x} - \widehat{\mathbf{x}}_{t}\|_{2}^{2} \right\}$$

Advanced Optimization (Fall 2023)

Definition 3 (Gradient Variance). Let T be the time horizon and $\mathcal{X} \subseteq \mathbb{R}^d$ be the feasible domain. For the function sequence f_1, \ldots, f_T with $f_t : \mathcal{X} \mapsto \mathbb{R}$ for $t \in [T]$, its gradient variance is defined as $\operatorname{Var}_{T} = \sup_{\{\mathbf{x}_{1},...,\mathbf{x}_{T}\}\in\mathcal{X}} \sum_{t=1} \left\|\nabla f_{t}(\mathbf{x}_{t}) - \boldsymbol{\mu}_{T}\right\|_{2}^{2}$ where $\mu_T \triangleq \frac{1}{T} \sum_{t=1}^{T} \nabla f_t(\mathbf{x}_t)$ is the gradient mean. **Optimistic Online Mirror Descent** self-confident estimate of gradient mean: $\boldsymbol{\mu}_t = \frac{1}{t} \sum_{s=1}^t \nabla f_s(\mathbf{x}_s)$ $\mathbf{x}_{t} = \arg\min_{\mathbf{x}\in\mathcal{X}} \left\{ \eta_{t} \left\langle M_{t}, \mathbf{x} \right\rangle + \frac{1}{2} \|\mathbf{x} - \widehat{\mathbf{x}}_{t}\|_{2}^{2} \right\}$ $\widehat{\mathbf{x}}_{t+1} = \arg\min_{\mathbf{x}\in\mathcal{X}} \left\{ \eta_t \left\langle \nabla f_t(\mathbf{x}_t), \mathbf{x} \right\rangle + \frac{1}{2} \|\mathbf{x} - \widehat{\mathbf{x}}_t\|_2^2 \right\}$

Advanced Optimization (Fall 2023)

Theorem 5 (gradient-variance bound). Assume that $\psi(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|_2^2$, when setting $\eta_t =$ $\sqrt{1+\widetilde{\operatorname{Var}}_{t-1}}$ and $M_t = \widetilde{\mu}_{t-1}$, the regret of Optimistic OMD to any comparator $\mathbf{u} \in \mathcal{X}$ is bounded as $\operatorname{Regret}_{T} = \sum_{t=1}^{T} f_{t}(\mathbf{x}_{t}) - \sum_{t=1}^{T} f_{t}(\mathbf{u}) \leq \widetilde{\mathcal{O}}\left(\sqrt{1 + \operatorname{Var}_{T}}\right)$ where $\widetilde{\operatorname{Var}}_{t-1} = \sum_{s=1}^{t-1} \|\nabla f_s(\mathbf{x}_s) - \boldsymbol{\mu}_s\|_2^2$ is the self-confident estimate of variance Var_T , and $\boldsymbol{\mu}_t = \frac{1}{t} \sum_{s=1}^t \nabla f_s(\mathbf{x}_s)$ is the empirical gradient mean. **Proof.** $\sum_{t=1}^{T} f_t(\mathbf{x}_t) - \sum_{t=1}^{T} f_t(\mathbf{u}) \le \sum_{t=1}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t) - M_t\|_2^2 + \sum_{t=1}^{T} \frac{1}{2\eta_t} \left(\|\mathbf{u} - \widehat{\mathbf{x}}_t\|_2^2 - \|\mathbf{u} - \widehat{\mathbf{x}}_{t+1}\|_2^2 \right)$ $-\sum_{t=1}^{T} \frac{1}{2\eta_t} \left(\|\widehat{\mathbf{x}}_{t+1} - \mathbf{x}_t\|_2^2 + \|\mathbf{x}_t - \widehat{\mathbf{x}}_t\|_2^2 \right)$ (negative term)

Advanced Optimization (Fall 2023)

Proof. For term (a), $\sum_{t=1}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t) - M_t\|_2^2 = \sum_{t=2}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_{t-1}\|_2^2 + G^2$ $(\eta_1 \triangleq 1)$ $\leq 2\sum_{t=2}^{I} \eta_t \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_t\|_2^2 + 2\sum_{t=2}^{I} \eta_t \|\boldsymbol{\mu}_t - \boldsymbol{\mu}_{t-1}\|_2^2 + G^2$ $\leq 2D \sum_{t=2}^{T} \frac{\|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_t\|_2^2}{\sqrt{1 + \sum_{s=1}^{t-1} \|\nabla f_s(\mathbf{x}_s) - \boldsymbol{\mu}_s\|_2^2}} + 2D \sum_{t=2}^{T} \frac{9G^2}{t^2} + G^2 \quad \begin{aligned} (\boldsymbol{\mu}_t = \frac{(t-1)\boldsymbol{\mu}_{t-1} + \nabla f_t(\mathbf{x}_t)}{t}) \\ (\|\boldsymbol{\mu}_t\|_2 \leq G, \forall t \in [T]) \\ (\eta_t \leq 1, \forall t \in [T]) \end{aligned}$ $\leq 2D \sum_{t=2}^{T} \frac{\|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_t\|_2^2}{\sqrt{1 + \sum_{s=1}^{t-1} \|\nabla f_s(\mathbf{x}_s) - \boldsymbol{\mu}_s\|_2^2}} + 18DG^2 \cdot \frac{\pi^2}{6} + G^2 \quad (\sum_{x=1}^{\infty} \frac{1}{x^2} = \frac{\pi^2}{6})$

Advanced Optimization (Fall 2023)

Proof.
$$\sum_{t=1}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t) - M_t\|_2^2 \le 2D \sum_{t=2}^{T} \frac{\|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_t\|_2^2}{\sqrt{1 + \sum_{s=1}^{t-1} \|\nabla f_s(\mathbf{x}_s) - \boldsymbol{\mu}_s\|_2^2}} + 18DG^2 \cdot \frac{\pi^2}{6} + G^2$$

Lemma 2. Let
$$a_1, a_2, \dots, a_T$$
 be non-negative real numbers. Then

$$\sum_{t=1}^T \frac{a_t}{\sqrt{1 + \sum_{s=1}^{t-1} a_s}} \leq 4\sqrt{1 + \sum_{t=1}^T a_t} + \max_{t \in [T]} a_t$$

$$\sum_{t=1}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t) - M_t\|_2^2 \le 8D \sqrt{1 + \sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_t\|_2^2 + 8DG^2 + 18DG^2 \cdot \frac{\pi^2}{6} + G^2 }$$
Recall that our goal is to obtain $\mathcal{O}\left(\sqrt{\sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_T\|_2^2}\right)$

Proof.
$$\sum_{t=1}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t) - M_t\|_2^2 \le 8D\sqrt{1 + \sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_t\|_2^2 + 8DG^2 + 18DG^2 \cdot \frac{\pi^2}{6} + G^2}$$

We need to measure the gap between $\sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_t\|_2^2$ and $\sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_T\|_2^2$

Let us consider *another online learning process*: the online function is $h_t : \mathbb{R}^d \mapsto \mathbb{R}$,

$$h_t(\mathbf{a}) = \frac{1}{2} \left\| \nabla f_t(\mathbf{x}_t) - \mathbf{a} \right\|_2^2,$$

which is evidently a 1-strongly convex function with respect to $\|\cdot\|_2$.

Consider OGD over $\{h_t\}_{t=1}^T$ with step size $\{\eta_t\}_{t=1}^T$, which updates by

$$\mathbf{a}_{t+1} = \mathbf{a}_t - \eta_t \nabla h_t(\mathbf{a}_t) = \mathbf{a}_t - \eta_t \left(\mathbf{a}_t - \nabla f_t(\mathbf{x}_t)\right) = (1 - \eta_t)\mathbf{a}_t + \eta_t \nabla f_t(\mathbf{x}_t) \qquad (\bigstar)$$

Advanced Optimization (Fall 2023)

Proof.
$$\sum_{t=1}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t) - M_t\|_2^2 \le 8D \sqrt{1 + \sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_t\|_2^2 + 8DG^2 + 18DG^2 \cdot \frac{\pi^2}{6} + G^2}$$

We need to measure the gap between $\sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_t\|_2^2$ and $\sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_T\|_2^2$

Consider OGD over $\{h_t\}_{t=1}^T$ with step size $\{\eta_t\}_{t=1}^T$, which updates by

$$\mathbf{a}_{t+1} = (1 - \eta_t)\mathbf{a}_t + \eta_t \nabla f_t(\mathbf{x}_t) \qquad (\bigstar)$$

On the other hand, by definition of gradient mean, we have

$$\boldsymbol{\mu}_t = \frac{t-1}{t} \boldsymbol{\mu}_{t-1} + \frac{1}{t} \nabla f_t(\mathbf{x}_t) \qquad (\boldsymbol{\mu}_t = \frac{1}{t} \sum_{s=1}^t \nabla f_s(\mathbf{x}_s))$$

Thus, set $\mathbf{a}_1 = \mathbf{0}$, $\eta_t = \frac{1}{t}$, then $\{\mathbf{a}_{t+1}\}_{t=1}^{T-1}$ sequence is *equivalent* to $\{\boldsymbol{\mu}_t\}_{t=1}^{T-1}$ sequence. More specifically, we have $\mathbf{a}_{t+1} = \boldsymbol{\mu}_t$ for $t = 1, \dots, T-1$.

Proof.
$$h_t(\mathbf{a}) = \frac{1}{2} \|\nabla f_t(\mathbf{x}_t) - \mathbf{a}\|_2^2$$
, $\mathbf{a}_{t+1} \stackrel{(\bigstar)}{=} (1 - \eta_t)\mathbf{a}_t + \eta_t \nabla f_t(\mathbf{x}_t)$, $\boldsymbol{\mu}_t = \frac{t - 1}{t} \boldsymbol{\mu}_{t-1} + \frac{1}{t} \nabla f_t(\mathbf{x}_t)$

Thus, set $\mathbf{a}_1 = \mathbf{0}$, $\eta_t = \frac{1}{t}$, then $\{\mathbf{a}_{t+1}\}_{t=1}^{T-1}$ sequence is *equivalent* to $\{\boldsymbol{\mu}_t\}_{t=1}^{T-1}$ sequence. Since (\bigstar) is essentially OGD for 1-strongly convex, whose guarantee is:

$$\operatorname{Regret}(\{h_t\}_{t=1}^{T-1}) = \sum_{t=1}^{T-1} h_t(\boldsymbol{\mu}_t) - \sum_{t=1}^{T-1} h_t(\boldsymbol{\mu}) \quad (\text{holds for any point } \boldsymbol{\mu} \text{ in } \mathbb{R}^d)$$
$$= \sum_{t=1}^{T-1} \frac{1}{2} \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_t\|_2^2 - \sum_{t=1}^{T-1} \frac{1}{2} \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_T\|_2^2 \quad (\text{taking } \boldsymbol{\mu}_T \text{ as the comparator})$$
$$\leq \frac{(2G)^2}{2\alpha} (1 + \ln(T - 1)) \quad (\text{regret bound of } \alpha \text{-strongly convex function does not rely} \text{ on domain diameter})$$

Proof.
$$h_t(\mathbf{a}) = \frac{1}{2} \|\nabla f_t(\mathbf{x}_t) - \mathbf{a}\|_2^2$$
, $\mathbf{a}_{t+1} \stackrel{(\bigstar)}{=} (1 - \eta_t)\mathbf{a}_t + \eta_t \nabla f_t(\mathbf{x}_t)$, $\boldsymbol{\mu}_t = \frac{t - 1}{t} \boldsymbol{\mu}_{t-1} + \frac{1}{t} \nabla f_t(\mathbf{x}_t)$

Thus, set $\mathbf{a}_1 = \mathbf{0}$, $\eta_t = \frac{1}{t}$, then $\{\mathbf{a}_{t+1}\}_{t=1}^{T-1}$ sequence is *equivalent* to $\{\boldsymbol{\mu}_t\}_{t=1}^{T-1}$ sequence. Since (\bigstar) is essentially OGD for 1-strongly convex, whose guarantee is:

$$\operatorname{Regret}(\{h_t\}_{t=1}^{T-1}) = \sum_{t=1}^{T-1} h_t(\boldsymbol{\mu}_t) - \sum_{t=1}^{T-1} h_t(\boldsymbol{\mu}_T) = \sum_{t=1}^{T-1} \frac{1}{2} \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_t\|_2^2 - \sum_{t=1}^{T-1} \frac{1}{2} \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_T\|_2^2 \le 2G^2(1+\ln T)$$

$$\Longrightarrow \sum_{t=1}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t) - M_t\|_2^2 \le 8D \sqrt{1 + \sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_t\|_2^2} + 8DG^2 + 18DG^2 \cdot \frac{\pi^2}{6} + G^2$$

$$\le 8D \sqrt{1 + \sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_T\|_2^2} + 4G^2(1+\ln T) + 8DG^2 + 18DG^2 \cdot \frac{\pi^2}{6} + G^2$$

Advanced Optimization (Fall 2023)

Proof. We then analyze term (b) in the same way as before:

$$\begin{aligned} \texttt{term} \left(\texttt{b} \right) &= \sum_{t=1}^{T} \frac{1}{2\eta_{t}} \left(\| \mathbf{u} - \widehat{\mathbf{x}}_{t} \|_{2}^{2} - \| \mathbf{u} - \widehat{\mathbf{x}}_{t+1} \|_{2}^{2} \right) \\ &= \sum_{t=2}^{T} \left(\frac{1}{2\eta_{t}} - \frac{1}{2\eta_{t-1}} \right) \| \mathbf{u} - \widehat{\mathbf{x}}_{t} \|_{2}^{2} + \frac{1}{2\eta_{1}} \| \mathbf{u} - \widehat{\mathbf{x}}_{1} \|_{2}^{2} \\ &\leq \sum_{t=2}^{T} \left(\frac{1}{2\eta_{t}} - \frac{1}{2\eta_{t-1}} \right) D^{2} + \frac{1}{2\eta_{1}} D^{2} \quad (\eta_{t} \le \eta_{t-1} \text{ and } \| \mathbf{x} - \mathbf{y} \|_{2} \le D, \forall \mathbf{x}, \mathbf{y} \in \mathcal{X}) \\ &\leq \frac{D^{2}}{2\eta_{T}} + \frac{1}{2\eta_{1}} D^{2} \le \frac{D}{2} \sqrt{1 + \operatorname{Var}_{T}} + \frac{D}{2} \qquad (\frac{1}{\eta_{T}} = \frac{\sqrt{1 + \operatorname{Var}_{T-1}}}{D} \le \frac{\sqrt{1 + \operatorname{Var}_{T}}}{D}) \end{aligned}$$

Advanced Optimization (Fall 2023)

Proof. Finally, putting three terms together achieves

 $\begin{aligned} &\texttt{term} \ (\texttt{a}) \leq 8D\sqrt{1 + \text{Var}_T + 4G^2(1 + \ln T)} + (39D + 1)G^2 \\ &\texttt{term} \ (\texttt{b}) \leq \frac{D^2}{2\eta_T} + \frac{1}{2\eta_1}D^2 \leq \frac{D}{2}\sqrt{1 + \text{Var}_T} + \frac{D}{2} \\ &\texttt{term} \ (\texttt{c}) \geq 0 \end{aligned}$

$$\square \qquad \text{Regret}_T = \texttt{term} (\texttt{a}) + \texttt{term} (\texttt{b}) - \texttt{term} (\texttt{c}) \\ \leq 9D\sqrt{1 + \text{Var}_T + 4G^2(1 + \ln T)} + 39DG^2 + G^2 = \widetilde{\mathcal{O}}\Big(\sqrt{1 + \text{Var}_T}\Big). \ \square$$

Advanced Optimization (Fall 2023)

Optimistic OMD: Applications

• Small-Loss Bound

• Gradient-Variance Bound

Gradient-Variation Bound

Gradient-Variation Bound

Definition 3 (Gradient Variation). Let *T* be the time horizon and $\mathcal{X} \subseteq \mathbb{R}^d$ be the feasible domain. For the function sequence f_1, \ldots, f_T with $f_t : \mathcal{X} \mapsto \mathbb{R}$ for $t \in [T]$, its gradient variation is defined as

$$V_T = \sum_{t=2}^{T} \sup_{\mathbf{x} \in \mathcal{X}} \left\| \nabla f_t(\mathbf{x}) - \nabla f_{t-1}(\mathbf{x}) \right\|_2^2$$

Gradient variation characterizes online functions' *shifting intensity*.

- *Adaptivity*: it can be small in slowly changing environments.
- *Robustness*: $V_T \leq 4G^2T$ in the worst case. $(\|\nabla f_t(\mathbf{x})\| \leq G, \forall \mathbf{x} \in \mathcal{X} \text{ and } t \in [T])$

Gradient-Variation Bound

Definition 3 (Gradient Variation). Let *T* be the time horizon and $\mathcal{X} \subseteq \mathbb{R}^d$ be the feasible domain. For the function sequence f_1, \ldots, f_T with $f_t : \mathcal{X} \mapsto \mathbb{R}$ for $t \in [T]$, its gradient variation is defined as

$$V_T = \sum_{t=2}^{T} \sup_{\mathbf{x} \in \mathcal{X}} \|\nabla f_t(\mathbf{x}) - \nabla f_{t-1}(\mathbf{x})\|_2^2$$



Implicit assumption:

Gradient (online function) shifts slowly

e.g., age forecasting by portraits

Optimistic OMD for Gradient-Variation Bound

Optimistic Online Mirror Descent

$$\mathbf{x}_{t} = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_{t} \left\langle \mathbf{M}_{t}, \mathbf{x} \right\rangle + \frac{1}{2} \|\mathbf{x} - \widehat{\mathbf{x}}_{t}\|_{2}^{2} \right\}$$
$$\widehat{\mathbf{x}}_{t+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_{t} \left\langle \nabla f_{t}(\mathbf{x}_{t}), \mathbf{x} \right\rangle + \frac{1}{2} \|\mathbf{x} - \widehat{\mathbf{x}}_{t}\|_{2}^{2} \right\}$$

Question: How to choose M_t ?



Imposing a prior on the change of the online functions

setting M_t as the **last-round gradient** $M_t = \nabla f_{t-1}(\mathbf{x}_{t-1})$

Optimistic OMD for Gradient-Variation Bound

7

Optimistic Online Mirror Descent

$$\mathbf{x}_{t} = \underset{\mathbf{x} \in \mathcal{X}}{\arg\min \eta_{t} \langle \mathbf{M}_{t}, \mathbf{x} \rangle} + \frac{1}{2} \|\mathbf{x} - \widehat{\mathbf{x}}_{t}\|_{2}^{2}$$
$$\widehat{\mathbf{x}}_{t+1} = \underset{\mathbf{x} \in \mathcal{X}}{\arg\min \eta_{t} \langle \nabla f_{t}(\mathbf{x}_{t}), \mathbf{x} \rangle} + \frac{1}{2} \|\mathbf{x} - \widehat{\mathbf{x}}_{t}\|_{2}^{2}$$

Optimistic OMD for Gradient-Variation Bound

$$\mathbf{x}_{t} = \underset{\mathbf{x}\in\mathcal{X}}{\arg\min \eta_{t}} \left\langle \nabla f_{t-1}(\mathbf{x}_{t-1}), \mathbf{x} \right\rangle + \frac{1}{2} \|\mathbf{x} - \widehat{\mathbf{x}}_{t}\|_{2}^{2}$$
$$\widehat{\mathbf{x}}_{t+1} = \underset{\mathbf{x}\in\mathcal{X}}{\arg\min \eta_{t}} \left\langle \nabla f_{t}(\mathbf{x}_{t}), \mathbf{x} \right\rangle + \frac{1}{2} \|\mathbf{x} - \widehat{\mathbf{x}}_{t}\|_{2}^{2}$$

Advanced Optimization (Fall 2023)

Gradient-Variation Bound

Theorem 4 (Gradient Variation Regret Bound). Assume that $\psi(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|_2^2$ and f_t is *L*-smooth for all $t \in [T]$, when setting $\eta_t = \min\{\frac{1}{4L}, \frac{D}{\sqrt{1+\tilde{V}_{t-1}}}\}$ and $M_t = \nabla f_{t-1}(\mathbf{x}_{t-1})$, the regret of Optimistic OMD to any comparator $\mathbf{u} \in \mathcal{X}$ is $\operatorname{Regret}_{T} = \sum_{t=1}^{T} f_{t}(\mathbf{x}_{t}) - \sum_{t=1}^{T} f_{t}(\mathbf{u}) \leq \mathcal{O}\left(\sqrt{1+V_{T}}\right)$ where $\widetilde{V}_{t-1} = \sum_{s=2}^{t-1} \|\nabla f_s(\mathbf{x}_{s-1}) - \nabla f_{s-1}(\mathbf{x}_{s-1})\|_2^2$ is the empirical estimates of V_t . **Proof.** $\sum_{t=1}^{T} f_t(\mathbf{x}_t) - \sum_{t=1}^{T} f_t(\mathbf{u}) \le \sum_{t=1}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t) - M_t\|_2^2 + \sum_{t=1}^{T} \frac{1}{2\eta_t} \left(\|\mathbf{u} - \widehat{\mathbf{x}}_t\|_2^2 - \|\mathbf{u} - \widehat{\mathbf{x}}_{t+1}\|_2^2 \right)$ $-\sum_{t=1}^{I} \frac{1}{2\eta_{t}} \left(\|\widehat{\mathbf{x}}_{t+1} - \mathbf{x}_{t}\|_{2}^{2} + \|\mathbf{x}_{t} - \widehat{\mathbf{x}}_{t}\|_{2}^{2} \right)$ (negative term)

Advanced Optimization (Fall 2023)

Proof. For term 1,

$$\sum_{t=1}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t) - M_t\|_2^2 \leq \sum_{t=2}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2^2 + G^2 \qquad (\eta_1 \triangleq 1)$$

$$\leq 2\sum_{t=2}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t) - \nabla f_t(\mathbf{x}_{t-1})\|_2^2 + 2\sum_{t=2}^{T} \eta_t \|\nabla f_t(\mathbf{x}_{t-1}) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2^2 + G^2$$

$$\leq 2\sum_{t=2}^{T} \eta_t L^2 \|\mathbf{x}_t - \mathbf{x}_{t-1}\|_2^2 + 2D\sum_{t=2}^{T} \frac{\|\nabla f_t(\mathbf{x}_{t-1}) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2^2}{\sqrt{1 + \sum_{s=2}^{t-1} \|\nabla f_s(\mathbf{x}_{s-1}) - \nabla f_{s-1}(\mathbf{x}_{s-1})\|_2^2}} + G^2$$

Proof. For term (a), $\sum_{t=1}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t) - M_t\|_2^2 \le \sum_{t=0}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2^2 + G^2$ $(\eta_1 \triangleq 1)$ $\leq 2\sum_{t=2}^{T} \eta_t \|\nabla f_t(\mathbf{x}_t) - \nabla f_t(\mathbf{x}_{t-1})\|_2^2 + 2\sum_{t=2}^{T} \eta_t \|\nabla f_t(\mathbf{x}_{t-1}) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2^2 + G^2$ $\leq 2\sum_{t=2}^{T} \eta_t L^2 \|\mathbf{x}_t - \mathbf{x}_{t-1}\|_2^2 + 2D \sum_{t=2}^{T} \frac{\|\nabla f_t(\mathbf{x}_{t-1}) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2^2}{\sqrt{1 + \sum_{s=2}^{t-1} \|\nabla f_s(\mathbf{x}_{s-1}) - \nabla f_{s-1}(\mathbf{x}_{s-1})\|_2^2}} + G^2$ **Lemma 2.** Let a_1, a_2, \ldots, a_T be non-negative real numbers. Then $\sum_{t=1}^{T} \frac{a_t}{\sqrt{1+\sum_{t=1}^{t-1} a_t}} \le 4\sqrt{1+\sum_{t=1}^{T} a_t + \max_{t\in[T]} a_t}$

Advanced Optimization (Fall 2023)

Proof. term (a)
$$\leq 2 \sum_{t=2}^{T} \eta_t L^2 \|\mathbf{x}_t - \mathbf{x}_{t-1}\|_2^2 + 4D \sqrt{1 + \sum_{t=2}^{T} \|\nabla f_t(\mathbf{x}_{t-1}) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2^2 + (4D+1)G^2}$$

$$\leq 2 \sum_{t=2}^{T} \eta_t L^2 \|\mathbf{x}_t - \mathbf{x}_{t-1}\|_2^2 + 4D\sqrt{1 + V_T} + (4D+1)G^2 \\ (V_T = \sum_{t=2}^{T} \sup_{\mathbf{x} \in \mathcal{X}} \|\nabla f_t(\mathbf{x}) - \nabla f_{t-1}(\mathbf{x})\|_2^2)$$
This term depends on our algorithm, how to deal with it?

Proof. For the term (c), we have

$$\begin{aligned} \texttt{term} \left(\mathbf{c} \right) &= \sum_{t=1}^{T} \frac{1}{2\eta_{t}} \left(\| \widehat{\mathbf{x}}_{t+1} - \mathbf{x}_{t} \|_{2}^{2} + \| \mathbf{x}_{t} - \widehat{\mathbf{x}}_{t} \|_{2}^{2} \right) \\ &\geq \sum_{t=2}^{T} \frac{1}{2\eta_{t}} \left(\| \widehat{\mathbf{x}}_{t} - \mathbf{x}_{t-1} \|_{2}^{2} + \| \widehat{\mathbf{x}}_{t} - \mathbf{x}_{t} \|_{2}^{2} \right) \quad \left(\frac{1}{\eta_{t}} \geq \frac{1}{\eta_{t-1}} \right) \\ &\geq \sum_{t=2}^{T} \frac{1}{4\eta_{t}} \left\| \mathbf{x}_{t} - \mathbf{x}_{t-1} \right\|_{2}^{2} \qquad (a^{2} + b^{2} \geq (a+b)^{2}/2) \end{aligned}$$

Does this term look familiar?

Proof. We then analysis term (b),

$$\begin{aligned} \mathsf{term} \left(\mathsf{b} \right) &= \sum_{t=1}^{T} \frac{1}{2\eta_{t}} \left(\| \mathbf{u} - \widehat{\mathbf{x}}_{t} \|_{2}^{2} - \| \mathbf{u} - \widehat{\mathbf{x}}_{t+1} \|_{2}^{2} \right) \\ &\leq \sum_{t=2}^{T} \left(\frac{1}{2\eta_{t}} - \frac{1}{2\eta_{t-1}} \right) \| \mathbf{u} - \widehat{\mathbf{x}}_{t} \|_{2}^{2} + \frac{1}{2\eta_{1}} \| \mathbf{u} - \widehat{\mathbf{x}}_{1} \|_{2}^{2} \\ &\leq \sum_{t=2}^{T} \left(\frac{1}{2\eta_{t}} - \frac{1}{2\eta_{t-1}} \right) D^{2} + \frac{1}{2\eta_{1}} D^{2} \quad (\eta_{t} \le \eta_{t-1} \text{ and } \| \mathbf{x} - \mathbf{y} \|_{2} \le D, \forall \mathbf{x}, \mathbf{y} \in \mathcal{X}) \\ &\leq \frac{D^{2}}{2\eta_{T}} \qquad \text{noting that } \eta_{T} = \min \left\{ \frac{1}{4L}, \frac{D}{\sqrt{1 + \widetilde{V}_{T-1}}} \right\} \ge \min \left\{ \frac{1}{4L}, \frac{D}{\sqrt{1 + V_{T}}} \right\} \\ &\leq \frac{1}{2} \max\{4LD, D\sqrt{1 + V_{T}}\} \end{aligned}$$

Advanced Optimization (Fall 2023)

Proof. Finally, putting three terms together yields

$$\begin{aligned} \operatorname{term} \left(\mathbf{a} \right) &\leq 2 \sum_{t=2}^{T} \eta_{t} L^{2} \| \mathbf{x}_{t} - \mathbf{x}_{t-1} \|_{2}^{2} + 4D\sqrt{1 + V_{T}} + (4D+1)G^{2} \\ \operatorname{term} \left(\mathbf{b} \right) &\leq \frac{1}{2} \max\{4LD, D\sqrt{1 + V_{T}}\} \\ \operatorname{term} \left(\mathbf{c} \right) &\geq \sum_{t=2}^{T} \frac{1}{4\eta_{t}} \| \mathbf{x}_{t} - \mathbf{x}_{t-1} \|_{2}^{2} \quad \left(\eta_{t} = \min\{\frac{1}{4L}, \frac{D}{\sqrt{1 + \tilde{V}_{t-1}}}\} \right) \\ \hline & \Longrightarrow \operatorname{Regret}_{T} = \operatorname{term} \left(\mathbf{a} \right) + \operatorname{term} \left(\mathbf{b} \right) - \operatorname{term} \left(\mathbf{c} \right) \\ &\leq 5D\sqrt{1 + V_{T}} + (4D+1)G^{2} + 2LD = \mathcal{O}\left(\sqrt{1 + V_{T}}\right). \end{aligned}$$

A Summary of Problem-dependent Bounds

$$\mathbf{x}_{t} = \underset{\mathbf{x}\in\mathcal{X}}{\arg\min \eta} \langle \mathbf{M}_{t}, \mathbf{x} \rangle + \frac{1}{2} \|\mathbf{x} - \widehat{\mathbf{x}}_{t}\|_{2}^{2}$$
$$\widehat{\mathbf{x}}_{t+1} = \underset{\mathbf{x}\in\mathcal{X}}{\arg\min \eta} \langle \nabla f_{t}(\mathbf{x}_{t}), \mathbf{x} \rangle + \frac{1}{2} \|\mathbf{x} - \widehat{\mathbf{x}}_{t}\|_{2}^{2}$$

Different priors are imposed by designing suitable M_t for specific environments.

	Assumption(s)	Setting of Optimism	Setting of η_t	Problem-dependent Regret Bound
Small-loss Bound	<i>L-</i> Smooth + Non-negative	$M_t = 0$	$pprox rac{D}{\sqrt{1+G_t}}$	$\mathcal{O}\left(\sqrt{1+F_T}\right)$
Variance Bound	_	$M_t = \widetilde{\mu}_{t-1}$	$\approx \frac{D}{\sqrt{1 + \operatorname{Var}_{t-1}}}$	$\widetilde{\mathcal{O}}\left(\sqrt{1 + \operatorname{Var}_T}\right)$
Variation Bound	L-Smooth	$M_t = \nabla f_{t-1}(\mathbf{x}_{t-1})$	$\approx \frac{D}{\sqrt{1+\widetilde{V}_{t-1}}}$	$\mathcal{O}\left(\sqrt{1+V_T}\right)$

Gradient-Variation Algorithm: Implications

By using algorithm for gradient-variation Bound (OMD with $M_t = \nabla f_{t-1}(\mathbf{x}_{t-1})$):

$$\begin{aligned} \left| \sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t) - \nabla f_{t-1}(\mathbf{x}_{t-1})\|_2^2 &\leq 3\sum_{t=1}^{T} \|\nabla f_t(\mathbf{x}_t) - \boldsymbol{\mu}_t\|_2^2 \qquad (\leq 3 \operatorname{Var}_T) \\ &+ \left| 3\sum_{t=1}^{T} \|\nabla f_{t-1}(\mathbf{x}_{t-1}) - \boldsymbol{\mu}_{t-1}\|_2^2 \\ &+ \left| 3\sum_{t=1}^{T} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}_{t-1}\|_2^2 \\ &+ \left| 3\sum_{t=1}^{T} \|\boldsymbol{\mu}_t - \boldsymbol{\mu}_{t-1}\|_2^2 \\ &\qquad (\|\boldsymbol{\mu}_t\|_2 \leq G, \forall t \in [T]) \\ &\qquad (\leq 3 \cdot \frac{\pi^2}{6}) \end{aligned}$$

 \implies Optimistic OMD with last-round gradient as optimism (enjoying V_T -bound) can also attain gradient-variance bound (scaling with Var_T)

Advanced Optimization (Fall 2023)

Gradient-Variation Algorithm: Implications

By using algorithm for gradient-variation Bound (OMD with $M_t = \nabla f_{t-1}(\mathbf{x}_{t-1})$):

$$\leq 8LF_T^X \qquad (F_T^X \triangleq \sum_{t=1}^T f_t(\mathbf{x}_t))$$

further use converting trick to attain F_T bound

 \Box Optimistic OMD with last-round gradient as optimism (enjoying V_T -bound) can also attain small-loss bound (scaling with F_T)

Advanced Optimization (Fall 2023)

Gradient-Variation Bound Reflection

Definition 3 (Gradient Variation). Let *T* be the time horizon and $\mathcal{X} \subseteq \mathbb{R}^d$ be the feasible domain. For the function sequence f_1, \ldots, f_T with $f_t : \mathcal{X} \mapsto \mathbb{R}$ for $t \in [T]$, its gradient variation is defined as

$$V_T = \sum_{t=2}^T \sup_{\mathbf{x} \in \mathcal{X}} \|\nabla f_t(\mathbf{x}) - \nabla f_{t-1}(\mathbf{x})\|_2^2$$

- This gradient-variation notion tightly connects the *offline optimization* and *online optimization*.
- The gradient variation reveals the importance of *smoothness* for the firstorder methods, as well as the crucial role of the *negative term* in analysis.

Offline Scenario

• Online algorithm with *gradient-variation* regret bound:

$$\operatorname{Regret}_{T} \triangleq \sum_{t=1}^{T} f_{t}(\mathbf{x}_{t}) - \min_{\mathbf{x} \in \mathcal{X}} \sum_{t=1}^{T} f_{t}(\mathbf{x}) \leq \mathcal{O}\left(\sqrt{1+V_{T}}\right).$$

• For an offline optimization problem $\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$

When the function is convex and *smooth*, we can use this *gradient-variation* algorithm to obtain an averaged model with error bound as

$$\varepsilon_T \triangleq f\left(\frac{1}{T}\sum_{t=1}^T \mathbf{x}_t\right) - \min_{\mathbf{x}\in\mathcal{X}} f(\mathbf{x}) \le \mathcal{O}\left(\frac{\sqrt{1+V_T(f,\ldots,f)}}{T}\right) = \mathcal{O}\left(\frac{1}{T}\right).$$

Advanced Optimization (Fall 2023)

Offline Scenario

• Online algorithm with *problem-independent* bound:

$$\operatorname{Regret}_{T} \triangleq \sum_{t=1}^{T} f_{t}(\mathbf{x}_{t}) - \min_{\mathbf{x} \in \mathcal{X}} \sum_{t=1}^{T} f_{t}(\mathbf{x}) \leq \mathcal{O}(\sqrt{T}).$$

• For an offline optimization problem $\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$

When the function is convex and *Lipschitz*, we can use this *problem-independent* algorithm to obtain an averaged model with error bound as

$$\varepsilon_T \triangleq f\left(\frac{1}{T}\sum_{t=1}^T \mathbf{x}_t\right) - \min_{\mathbf{x}\in\mathcal{X}} f(\mathbf{x}) \le \mathcal{O}\left(\frac{\sqrt{T}}{T}\right) = \mathcal{O}\left(\sqrt{\frac{1}{T}}\right).$$

Advanced Optimization (Fall 2023)

Variation-type Bounds: History Bits

Extracting Certainty from Uncertainty: Regret Bounded by Variation in Costs

Elad Hazan IBM Almaden 650 Harry Rd, San Jose, CA 95120 hazan@us.ibm.com

Abstract

Prediction from expert advice is a fundamental prob lem in machine learning. A major pillar of the field is the existence of learning algorithms whose average loss approaches that of the best expert in hindsight (in other words, whose average regret approaches zero). Traditionally the regret of online algorithms was bounded in terms of the number of prediction rounds.

Cesa-Bianchi, Mansour and Stoltz [4] posed the question whether it is be possible to bound the regret of an online algorithm by the *variation* of the observed costs. In this paper we resolve this question, and prove such bounds in the fully adversaial setting, in two important online learning scenarios: prediction from expert advice, and online linear optimization.

1 Introduction

A cornerstone of modern machine learning are algorithms for prediction from expert advice. The seminal work of Littlestone and Warnuth [12], Vovk [13] and Freund and Schapier [6] gave algorithms which, under fully adversarial cost sequences, attain average cost approaching that of the best expert in hindsight.

To be more precise, consider a prediction setting in which an online learner has access to ne experts. Inertirely, the learner may chose the advice of any expert deterministically prendomly. After chosoing a course of action, an adversary reveals the cost of following the advice of the different experts, from which the expected cost of the online learner is derived. The classic results mentioned abver give algorithms which sequentially produce randomized decisions, such that the difference between the (expected) cost of the algorithm and the best expert in hindsight grows like $O(\sqrt{T \log n})$, where T is the number of prediction iterations. This extra additive cost is known as the regret of the online learning algorithm.

noweet, a prior it is not clear why omme rearing a gorithms should have high regret (growing with the number of iterations) in an unchanging environment. As an extreme example, consider a setting in which there are only two experts. Suppose that the first expert always incurs cost 1, whereas Satyen Kale Microsoft Research I Microsoft Way, Redmond, WA 98052 sakale@microsoft.com

the second expert always incurs cost $\frac{1}{2}$. One would expect to "figure out" this pattern quickly, and focus on the second expert, thus incurring a total cost that is at most $\frac{1}{2}$ plus at most a constant extra cost (irrespective of the number of rounds T), thus having only constant regret. However, any straightforward application of previously known analyses of expert learning algorithms only gives a regret bound of $\Theta(\sqrt{T})$ in this simple case (or very simple variations of it.)

More generally, the natural bound on the regret of a "good" learning algorithm should depend on variation in the sequence of costs, rather than purely on the number of iterations. If the cost sequence has low variation, we expect our algorithm to be able to perform better.

This intuition has a direct analog in the stochastic setting: here, the sequence of experts' costs are independently sampled from a distribution. In this situation, a natural bound on the rate of convergence to the optimal expert is controlled by the variance of the distribution (low variance should imply faster convergence). This was formalized by Cesa-Bianchi, Mansson and Stoltz [4], who assert that "proving such a rate in the full distortarratia setting would be a fundamental re-

In this paper we prove the first such regret bounds on online learning algorithms in two important scenarios: prediction from expert advice, and the more general framework of online linear optimization. Our algorithms have regret bounded by the variation of the cost sequence, in a manner that is made precise in the following sections. Thus, our bounds are tighter than *all* previous bounds, and hence yield better bounds on the applications of previous bounds (see, for example, the applications of 14).

1.1 Online linear optimization

Online linear optimization [10] is a general framework for online learning which has received much attention recently. In this framework the decision set is an arbitrary bounded, closed, convex set in Euclidean space $K \subseteq \mathbb{R}^n$ rather than a fixed set of experts, and the costs are determined by adversarially constructed vectors, $f_1, f_2, \dots, \in \mathbb{R}^n$, such that the cost of point $x \in K$ is given by $f_1 \cdot x$. The online learner iteratively chooses a point in the enverse set $x_1 \in K$, and then the cost vector f_1 is revealed and the cost $f_1 \cdot x_1$ is occurred. The performance of online learner it-is inseasured by the regret, which is defined as the difference in the total cost of the sequence of points chosen by the algeorithm, viz. JMLR: Workshop and Conference Proceedings vol 23 (2012) 6.1-6.20 25th Annual Conference on Learning Theory

Online Optimization with Gradual Variations

Chao-Kai Chiang ^{1,2}	CHAOKAI@IIS.SINICA.EDU.TW
Tianbao Yang ³	YANGTIA1@MSU.EDU
Chia-Jung Lee ¹	LEECJ@HS.SINICA.EDU.TW
Iehrdad Mahdavi ³	MAHDAVIM@CSE.MSU.EDU
Chi-Jen Lu ¹	CJLU@IIS.SINICA.EDU.TW
Rong Jin ³	RONGJIN@CSE.MSU.EDU
benghuo Zhu ⁴	ZSH@SV.NEC-LABS.COM
Institute of Information Science,	
cademia Sinica, Taipei, Taiwan.	
Department of Computer Science and Information Engineering,	
lational Taiwan University, Taipei, Taiwan.	
Department of Computer Science and Engineering	
Iichigan State University, East Lansing, MI, 48824, USA	
NEC Laboratories America	
Cupertino, CA, 95014, USA	

Editor: Shie Mannor, Nathan Srebro, Robert C. Williamson

Abstract

We study the online convex optimization problem, in which an online algorithm has to make repeated decisions with convex loss functions and hopes to achieve a small regret. We consider a natural restriction of this problem in which the loss functions have a small deviation, measured by the sum of the distances between every two consecutive loss functions, according to some distance metrics. We show that for the linear and general smooth convex loss functions, an online algorithm modified from the gradient descend algorithm can achieve a regret which only scales as the square root of the deviation. For the closely related problem of prediction with expert advice, we show that an online algorithm modified from the multiplicative update algorithm can also achieve a similar regret bound for a different measure of deviation. Finally, for loss functions which are strictly convex, we show that an online algorithm modified from the online Newton step algorithm can also aregret which is only logarithmic in terms of the deviation, and as an application, we can also have such a logarithmic regret for the portfolio management problem. Keywords: Online Learning, Regret, Convex Optimization, Deviation.

1. Introduction

We study the online convex optimization problem in which a player has to make decisions iteratively for a number of rounds in the following way. In round t, the player has to choose a point x_t from some convex feasible set $X \subseteq \mathbb{R}^N$, and after that the player receives a convex loss function f_t and suffers the corresponding loss $f_t(x_t) \in [0, 1]$. The player would like to have an online algorithm that can minimize its regret, which is the difference between the total loss it suffers and that of the best fixed point in hindsight. It is known

© 2012 C.-K. Chiang, T. Yang, C.-J. Lee, M. Mahdavi, C.-J. Lu, R. Jin & S. Zhu.



Extracting Certainty from Uncertainty: Regret Bounded by Variation in Costs. COLT 2008. Online Optimization with Gradual Variations. COLT 2012.

Advanced Optimization (Fall 2023)

Summary



Advanced Optimization (Fall 2023)