## Lecture 12: More Online Convex Optimization

*Lecturer: Jacob Abernethy*                        *Scribes: Jihui Jin & Sihan Zeng*

**Disclaimer**: *These notes have not been subjected to the usual scrutiny reserved for formal publications.*

## 12.1   Reminder of Online Convex Optimization Framework

The online convex optimization framework involves the following protocol. Given a convex set $\mathcal{K} \subseteq \mathbb{R}^d$
For $t = 1, ..., T$:

1. Alg selects $x_t \in \mathcal{K}$

2. Nature reveals convex function $f_t : \mathcal{K} \to \mathbb{R}$

**Goal:** Minimize $\text{Regret}_T = \sum_{t=1}^T f_t(x_t) - \min_{x \in \mathcal{K}} \sum_{t=1}^T f_t(x)$

**Definition 12.1 (projection)** *For any $y \in \mathbb{R}^d$, the **projection** onto a set $\mathcal{K}$ is defined as*

$$\Pi_{\mathcal{K}}(y) = \operatorname*{argmin}_{x \in \mathcal{K}} \|x - y\|_2$$

**Lemma 12.2 (Pythagorean Theorem for Bregman Divergences)** *For any $x \in \mathcal{K}$,*

$$\|x - \Pi_{\mathcal{K}}(y)\| \le \|x - y\|$$

*, and the equation holds iff $y \in K$.*

**Proof: (Exercise)**                                                                                       ∎

## 12.2   Online Gradient Descent

**Assumptions:**

- $f_t$ is $G$−Lipschitz: $\|\nabla f_t\| \le G$

- $\mathcal{K}$ has diameter $D$: $\|x - y\| \le D \quad \forall x, y \in \mathcal{K}$

---

**Algorithm 1:** Online Gradient Descent

---
1 Init $x_1 \in \mathcal{K}$;
2 **for** $i \leftarrow 1$ **to** $T$ **do**
3     $\eta_t = \frac{D}{G\sqrt{t}}$;
4     $\tilde{x}_{t+1} = x_t + \eta_t \nabla f_t(x_t)$;
5     $x_{t+1} = \Pi_{\mathcal{K}}(\tilde{x}_{t+1})$;
6 **end**

---

**Theorem 12.3** *Assume that $\|\nabla f_t\| \le G$, $\|x - y\| \le D$ for all $x, y$ in $\mathcal{K}$. Then*

$$\text{Regret}_T(OGD) \le \frac{3}{2} G D \sqrt{T}$$

**Proof:** Let $x^* = \text{argmin}_{x \in \mathcal{K}} \sum_{t=1}^{T} f_t(x)$. Define $\nabla_t := \nabla f_t(x_t)$

By convexity of $f_t$ at $x_t$:

$$f(x_t) - f_t(x^*) \leq \nabla_t^\top (x_t - x^*)$$

By Lemma 12.2, we have that

$$\begin{aligned} \|x_{t+1} - x^*\|^2 &= \|\Pi_\mathcal{K}(x_t - \eta_t \nabla_t) - x^*\|^2 \\ &\leq \|x_t - \eta_t \nabla_t - x^*\|^2 \\ &= \|x_t - x^*\|^2 + \eta_t^2 \|\nabla_t\|^2 - 2\eta_t \nabla_t^\top (x_t - x^*) \end{aligned}$$

Combining these two gives that

$$f_t(x_t) - f_t(x^*) \leq \frac{\|x_t - x^*\|^2}{2\eta_t} + \frac{\eta_t \|\nabla_t\|^2}{2} - \frac{\|x_{t+1} - x^*\|^2}{2\eta_t}$$

$$\begin{aligned} \text{Regret}_T &= \sum_{t=1}^{T} f_t(x_t) - f_t(x^*) \\ &\leq \frac{1}{2} \sum_{t=1}^{T} \left( \frac{\|t_x - x^*\|^2 - \|x_{t+1} - x^*\|^2}{\eta_t} \right) + \frac{G^2}{2} \sum_{t=1}^{T} \eta_t \\ &\leq \frac{1}{2} \sum_{t=1}^{T} \|x_t - x^*\|^2 \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) + \frac{G^2}{2} \sum_{t=1}^{T} \eta_t \\ &\leq \frac{D^2}{2} \sum_{t=1}^{T} \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) + \frac{G^2}{2} \sum_{t=1}^{T} \eta_t \\ &= \frac{D^2}{2} \left( \frac{1}{\eta_T} - \frac{1}{\eta_0} \right) + \frac{G^2}{2} \sum_{t=1}^{T} \frac{D}{G\sqrt{t}} \\ &\leq \frac{1}{2} DG\sqrt{T} + \frac{G^2}{2} \left( 2\frac{D}{G}\sqrt{T} \right) \\ &= \frac{3}{2} DG\sqrt{T} \end{aligned}$$

Where in the second to last step we make use of the inequality that

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

∎

**Corollary 12.4** *Gradient descent algorithms for minimizing one convex function $f$ with "averaging" converges at a rate of $O(\frac{DG}{\sqrt{T}})$*

## 12.3   Stochastic Gradient Descent (SGD)

Let $h(\cdot; \cdot)$ be some convex loss function. Consider the SGD algorithm to minimize $f(x) = \mathbb{E}_\zeta[h(x; \zeta)]$, where $\zeta$ denotes samples from the distribution.

The steps of the SGD algorithm at each iteration $t$ is as follows:

$$(1) \text{ sample } \zeta_t$$
$$(2) \ x_{t+1} \leftarrow x_t - \eta_t \nabla h(x_t; \zeta_t)$$
$$(3) \text{ output:} \bar{x}_T = \frac{1}{T} \sum_{t=1}^{T} x_t$$

**Claim 12.5** *SGD converges at* $O(\frac{DG}{\sqrt{T}})$.

**Proof:**

$$\underset{\zeta_{1:T}}{\mathbb{E}} [f(\bar{x}_T)] \leq \underset{\zeta}{\mathbb{E}}[\frac{1}{T} \sum_{t=1}^{T} f(x_t)] \text{ by the convexity of the loss function}$$

$$= \underset{\zeta}{\mathbb{E}}[\frac{1}{T} \sum_{t=1}^{T} \underset{\zeta_{1:t}}{\mathbb{E}} [f(x_t)|\zeta_{1:t-1}]] \text{ since } x_t \text{ only depends on samples before } t$$

$$= \underset{\zeta}{\mathbb{E}}[\frac{1}{T} \sum_{t=1}^{T} \underset{\zeta_{1:t}}{\mathbb{E}} [h(x_t; \zeta_t)|\zeta_{1:t-1}]] \text{ using the definition of } f$$

$$= \underset{\zeta}{\mathbb{E}}[\frac{1}{T} \sum_{t=1}^{T} f_t(x_t)] \text{ using the fact that } \mathbb{E}[\mathbb{E}[A|B]] = \mathbb{E}[A]$$

$$\leq \underset{\zeta}{\mathbb{E}}[\frac{1}{T} \sum_{t=1}^{T} f_t(x^*)] + \frac{Regret_T(OGD)}{T} \text{ since we use the OGD protocol to update x at each iteration}$$

$$= f(x^*) + \frac{3DG}{2\sqrt{T}}$$

Therefore, the convergence rate $\mathbb{E}[f(\bar{x}_T)] - f(x^*) \leq \frac{3DG}{2\sqrt{T}}$.

∎