

Lecture 19: June 3

Lecturer: Yu-Xiang Wang

Scribes: Sanae Amani Geshnigani

Note: *LaTeX template courtesy of UC Berkeley EECS dept.*

Disclaimer: *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.*

This lecture's notes illustrate some uses of various \LaTeX macros. Take a look at this and imitate.

19.1 Remainder of Lecture 18

Examples of ADMM:

$$\begin{aligned} \min_{\theta} \|\theta - y\|_2^2 + \lambda \|D^{(k+1)}\theta\|_1 &\Leftrightarrow \min_{\theta} \|\theta - y\|_2^2 + \lambda \|D^{(1)}z\|_1 \\ &\text{s.t. } D^{(k)}\theta = z \end{aligned}$$

The Lagrangian L is $L = \|\theta - y\|^2 + \lambda \|D^{(1)}z\|_1 + u^T(D^{(k)}\theta - z) + \frac{\rho}{2} \|D^{(k)}\theta - z\|^2$.

Update rule is:

1. Find $\operatorname{argmin}_{\theta} L(\theta, z, u) = \operatorname{argmin}_{\theta} \theta \left(\frac{\rho}{2} D^{(k)T} D^{(k)} + I \right) \theta + \tilde{b}^T \theta$. — a linear system that is banded diagonal, which can be solved in $O(kn)$ time by Gaussian elimination.
2. $\operatorname{argmin}_z L(\theta, z, u) = \operatorname{prox}_{\|D^{(1)}\cdot\|_1}(\tilde{b}) \rightarrow$ Fused lasso/TV denoising \rightarrow can be solved by Dynamic Programming algorithm in $O(n)$.
3. $u^+ = u + (D^{(k)}\theta - z)$

Further generalization: D is the incidence matrix of a graph. The linear system can be solved by fast Laplacian solvers. The prox-operator can be solved by graph-cut, parametric max-flow and etc.

19.2 Recap of Online Convex Optimization

For $t = 1, 2, \dots, T$ **then**

Player chooses $x_t \in \mathcal{K}$

Adversary chooses $f_t \in \mathcal{F}$ (strongly convex functions)

Player incurs a loss $f_t(x_t)$

Player receives feedback:

$$\begin{cases} \nabla f_t(x_t) \in \text{full information setting} \\ f_t(x_t) \in \text{bandit setting} \\ \nabla f_t(x_t) + z_t \in \text{noisy gradient setting} \end{cases}$$

End For

Goal:

$$\text{Regret} = \sum_{t=1}^T f_t(x_t) - \min_x \sum_{t=1}^T f_t(x)$$

$$\text{no-regret algorithm} \iff \lim_{T \rightarrow \infty} \frac{\text{Regret}_T}{T} = 0.$$

Definition of "Static Regret" with respect to parameter u [1]:

$$\text{Regret}(u) = \sum_{t=1}^T f_t(x_t) - \sum_{t=1}^T f_t(u)$$

Example 1. Let $f_t(x) = (x - t/T)^2$, $\mathcal{K} = [0, 1]$. What is the best expert in the hindsight?

$$\min_x \sum_{t=1}^T f_t(x) = \sum_{t=1}^T (x - t/T)^2 \Rightarrow \nabla = \sum_{t=1}^T 2(x - t/T) = 0 \Rightarrow x^* = \frac{T+1}{2T} \asymp \frac{1}{2}.$$

$$\text{Then, the optimal value is: } \sum_{t=1}^T (\frac{T+1}{2T} - t/T)^2 = \frac{1}{T^2} \int_1^T (\frac{T+1}{2} - t)^2 dt = \frac{1}{T^2} [\frac{1}{3}(t - \frac{T+1}{2})^3]_1^T \asymp \mathcal{O}(T)!$$

The conclusion is that even though we can get $O(\log T)$ regret in this case, it doesn't mean much because we are comparing to a very weak baseline.

Question: Can we do any better?

19.3 Dynamic Regret [Zinkevich, 2003]

Dynamic regret competes against an arbitrary sequence of competitors in the *hindsight*.

$$\text{D.Regret} = \sum_{t=1}^T f_t(x_t) - \min_{(u_1, \dots, u_T)} \sum_{t=1}^T f_t(u_t)$$

Example 2.

$$f_t(x) = \begin{cases} (x - 1)^2 & , \text{w.p } 1/2 \\ (x + 1)^2 & , \text{w.p } 1/2 \end{cases}$$

So, the best dynamic competitor in the hindsight when f_1, \dots, f_T are known is $u_t = \text{argmin}_x f_t(x)$ which gives $\min_{(u_1, \dots, u_T)} \sum_{t=1}^T f_t(u_t) = 0$. This is not possible to achieve. Why?

Take time t , suppose the player knows the adversary is doing the above random sampling:

$$\min_x \mathbb{E}[f_t(x)] = \min_x \frac{1}{2}(x - 1)^2 + \frac{1}{2}(x + 1)^2 \stackrel{x=0}{=} 1 \Rightarrow \text{D.Regret of any player} = T!$$

What do we do?

1. Restrict the family of competitor class u_1, \dots, u_T :

$$\left\{ (u_1, \dots, u_T) \left| \sum_{t=2}^T \|u_t - u_{t-1}\|_2 \leq P_T \right. \right\} \quad \text{path constraint [Zinkevich, 2003]} \quad (19.1)$$

2. Make assumptions on sequence of f_1, \dots, f_T such that they change slowly, e.g., what Besbes et al. [2015] assume

$$\sum_{t=2}^T \|f_t - f_{t-1}\|_\infty \leq V_T \Leftrightarrow \sum_{t=2}^T \sup_x |f_t(x) - f_{t-1}(x)| \leq V_T \quad (19.2)$$

A generalization of the above is:

$$\left(\sum_{t=2}^T \|x_t - x_{t-1}\|_p^q \right)^{1/q} \leq V_T(p, q) \quad (19.3)$$

which is the topic of [Chen et al., 2019].

An alternative assumption about the function-variation is the following. Let $x_t^* = \operatorname{argmin}_x f_t(x)$. Then an assumption can be made on x_1^*, \dots, x_T^* changing slowly, i.e.,

$$\sum_{t=2}^T \|x_t^* - x_{t-1}^*\|_2 \leq U_T$$

which is studied in [Yang et al., 2016].

Note that the original path-length-regret is very general because it parameterizes the regret with the path-length instead of making any assumptions about the functions (which we typically can only assume, but not verify.)

19.4 Path-length constraint, full information feedback

In the first class of assumption:

$$\operatorname{Regret}(T, u_1, \dots, u_T) = \sum_{t=1}^T f_t(x_t) - f_t(u_t) \leq \operatorname{function} \left(T, P_T(u, 1, \dots, u_T) = \sum_{t=2}^T \|u_t - u_{t-1}\|_2 \right)$$

Online Gradient Descent (OGD):

$x_{t+1} = \operatorname{proj}_{\mathcal{K}}(x_t - y_t \nabla_t(x_t)) \Rightarrow \operatorname{S.Regret} \leq \mathcal{O}(GD\sqrt{T})$,
where $\nabla \|f_t(x)\|_2 \leq G$, i.e., f_t is G -lipschitz, and $\|u_1 - u_2\|_2 \leq D \forall u_1, u_2 \in \mathcal{K}$.

Theorem 19.1 (Thm. 2 [Zinkevich, 2003]). *If $\eta_t = \eta$, then:*

$$G \cdot \operatorname{Regret}(T, P_T) \leq \frac{D^2}{2\eta} + \frac{P_T D}{\eta} + \frac{T\eta G^2}{2} \asymp \sqrt{T(D^2 + P_T D)G^2}. \quad (19.4)$$

Proof. We rely on:

1. Convexity of $f_t, \forall u \in \mathcal{K} : f_t(x_t) - f_t(u) \leq \langle g_t, x_t - u \rangle$.
2. By OGD algorithm:

$$\|x_t - u\|_2^2 = \|\Pi_{\mathcal{K}}(x_t - \eta g_t) - u\|_2^2 \stackrel{u \in \mathcal{K}}{\leq} \|x_t - \eta g_t - u\|_2^2 = \|x_t - u\|_2^2 + \eta^2 \|g_t\|_2^2 - 2\eta \langle g_t, x_t - u \rangle.$$

3. $\|x_{t+1} - u_t\|_2^2 = \|x_{t+1} + u_{t+1} - u_{t+1} - u_t\|_2^2 = \|x_{t+1} - u_{t+1}\|_2^2 + \|u_{t+1} - u_t\|_2^2 + 2\langle x_{t+1} - u_{t+1}, u_{t+1} - u_t \rangle$.

Take $u = u_t$. Then following (2), we have: $\|x_t - u_t\|_2^2 \leq \|x_t - u_t\|_2^2 + \eta^2 \|g_t\|_2^2 - 2\eta \langle g_t, x_t - u_t \rangle$. Furthermore, we have:

$$\begin{aligned} f_t(x_t) - f_t(u_t) &\stackrel{(1)}{\leq} \langle g_t, x_t - u_t \rangle \\ &\stackrel{(2)}{\leq} \frac{1}{2\eta} (\|x_t - u_t\|_2^2 + \eta^2 \|g_t\|_2^2 - \|x_{t+1} - u_t\|_2^2) \\ &\stackrel{(3)}{\leq} \frac{1}{2\eta} (\|x_t - u_t\|_2^2 + \eta^2 \|g_t\|_2^2 - \|x_{t+1} - u_{t+1}\|_2^2 - \|u_{t+1} - u_t\|_2^2 - 2\langle x_{t+1} - u_{t+1}, u_{t+1} - u_t \rangle) \\ &\leq \frac{1}{2\eta} (\|x_t - u_t\|_2^2 - \|x_{t+1} - u_{t+1}\|_2^2 + 2D \|u_{t+1} - u_t\|_2^2 + \eta^2 G^2). \end{aligned}$$

Now, sum up $t = 1, \dots, T$, Telescope:

$$\begin{aligned} \sum_{t=1}^T f_t(x_t) - f_t(u_t) &\leq \frac{1}{2\eta} \left(\|x_1 - u_1\|_2^2 - \|x_{T+1} - u_{T+1}\|_2^2 + 2D \sum_{t=1}^T \|u_{t+1} - u_t\|_2^2 + T\eta^2 G^2 \right) \\ &\leq \frac{D^2}{2\eta} + \frac{\eta G^2 T}{2} + \frac{2DP_T}{2\eta} \stackrel{\eta = \frac{\sqrt{D^2 + DP_T}}{\sqrt{T}G}}{\leq} \sqrt{TG^2(D^2 + DP_T)}. \end{aligned}$$

□

- This bound is optimal in T, P_T but the OGD cannot compete against all u_1, \dots, u_T , simultaneously. So, we want to design an *Adaptive Algorithm* when P_T is *not* an input which is done by [Zhang et al., 2018]
- if $P_T = 1$, we get static regret bound. If $P_T = T$ (Example 2), then we get D.Regret = T

19.5 Function variation constraint, noisy gradient feedback

Now, suppose we make assumption on function variation and have noisy gradient feedback:

$$\sum_{t=2}^T \sup_x |f_t(x) - f_{t-1}(x)| \leq V_T.$$

The feedback model is $g_t = \nabla f_t(x_t) + z_t$, where z_t is independent sub-Gaussian noise. The dynamic regret is defined as follows:

$$\text{D.Regret} = \mathbb{E} \left[\sum_{t=1}^T f_t(x_t) \right] - \sum_{t=1}^T f_t(x_t^*)$$

, where $x_t^* = \operatorname{argmin}_x f_t(x)$. The final regret bound would be a function of T and V_T .

The algorithm takes any sub-routine A that is OCO with static regret bound. We partition the time horizon into batches of size Δ_T . Then, the algorithm runs A in each sequence of Δ_T rounds and restarts at the end of that time slot.

Theorem 19.2 (Prop. 2 [Besbes et al., 2015]). $\text{Regret}(T, V_T) \leq \lceil \frac{T}{\Delta_T} \rceil \text{Regret}_A(\Delta_T) + 2\Delta_T V_T$.

Proof. Let $x_j^* = \operatorname{argmin}_x \sum_{t_j=1}^{\Delta_T} f_{t_j}(x)$.

$$\begin{aligned} \sum_{t=1}^T f_t(x_t) - f_t(x_t^*) &= \sum_{j=1}^{\frac{T}{\Delta_T}} \sum_{t_j=1}^{\Delta_T} f_{t_j}(x_{t_j}) - f_{t_j}(x_{t_j}^*) + f_{t_j}(x_{t_j}^*) - f_{t_j}(x_{t_j}^*) \\ &\leq \frac{T}{\Delta_T} \operatorname{Regret}_A(\Delta_T) + \underbrace{\sum_{j=1}^{\frac{T}{\Delta_T}} \sum_{t_j=1}^{\Delta_T} f_{t_j}(x_{t_j}^*) - f_{t_j}(x_{t_j}^*)}_{(\star)} \end{aligned}$$

Now, we bound (\star) .

$$\begin{aligned} (\star) &= \sum_j \Delta_T \left(\frac{1}{\Delta_T} \sum_{t_j=1}^{\Delta_T} f_{t_j} \right) (x_j^*) - \sum_{t_j=1}^{\Delta_T} f_{t_j}(x_{t_j}^*) \\ &\stackrel{x_j^* \text{ is optimal}}{\leq} \sum_j \Delta_T \left(\frac{1}{\Delta_T} \sum_{t_j=1}^{\Delta_T} f_{t_j} \right) \left(\frac{1}{\Delta_T} \sum_j x_{t_j}^* \right) - \sum_{t_j=1}^{\Delta_T} f_{t_j}(x_{t_j}^*) \\ &\stackrel{\frac{1}{\Delta_T} \sum_{t_j=1}^{\Delta_T} f_{t_j} \text{ is convex}}{\leq} \sum_j \left[\left(\frac{1}{\Delta_T} \sum_{t_j=1}^{\Delta_T} f_{t_j} \right) (x_{t_j}^*) - f_{t_j}(x_{t_j}^*) \right] \end{aligned}$$

Let $\frac{1}{\Delta_T} \sum_{t_j=1}^{\Delta_T} f_{t_j} = \tilde{f}_j$ and $V_T = \sum_j V_T^{(j)}$. For $\forall i_1, i_2 \in j$ -th Bin, we have $\|f_{i_1} - f_{i_2}\|_\infty \leq V_T^{(j)}$. Hence:

$$(\star) \leq \sum_j V_T^{(j)} \leq \Delta_T V_T \Rightarrow \text{D.Regret} \leq \frac{T}{\Delta_T} \operatorname{Regret}_A(\Delta_T) + \Delta_T V_T$$

If A is *Convex* OGD \Rightarrow D.Regret $\asymp T^{2/3} V_T^{1/3}$.

If A is *Strongly Convex* OGD \Rightarrow D.Regret $\asymp \sqrt{\frac{T \log TV_T G}{m}} \asymp \sqrt{TV_T}$. \square

19.6 A natural family of definitions for the variational functionals?

Now, we get back to Example 1: $f_t(x) = (x - t/T)^2 = (x - \Theta_t)^2$. Suppose, the player has access to $g_t = 2(x_t - \Theta_t) + z_t$. So, $\hat{\Theta}_t = -\frac{g_t}{2} + x_t = \Theta_t + \text{noise}$. This is called a non-parametric regression problem where we assume Θ_t changes slowly and want to design an algorithm A to minimize $MSE = \mathbb{E} \left[\sum_t (A_t(g_1, \dots, g_t) - \Theta_t)^2 \right]$.

- If $\sum_t^T |\Theta_t - \Theta_{t-1}| \leq P_T \Rightarrow \Theta \in \text{TV}(P_T)$.
- If $\sqrt{\sum_t^T |\Theta_t - \Theta_{t-1}|^2} \leq P_T \Rightarrow \Theta \in \text{Sobolev}(\frac{P_T}{\sqrt{T}})$.
- If $\left(\sum_t^T (\Theta_t - \Theta_{t-1})^p \right)^{1/p} \leq P_T \Rightarrow \Theta \stackrel{p \rightarrow \infty}{\in} \text{Holder Class}(\frac{P_T}{T})$.

The scaling is chosen to match their definitions in the corresponding continuous function class.

For the sequences in the first setting, the optimal rate of offline problem is $\mathcal{O}(T^{1/3}P_T^{2/3})$, a well-known information-theoretic lower bound due to [Donoho et al., 1990]. If we think of this feedback model as $\nabla(y_t - x_t)^2 \Leftrightarrow y_t = \Theta_t + \text{noise}$, OGD and Restarting OGD can achieve regret $\mathcal{O}(T^{1/2}P_T^{1/2})$, which is also a lower bound for all *linear smoothers* — a class of algorithms that subsume OGD and Restarting OGD. On the other hand, with a more careful design of the algorithm one can have an optimal algorithm for this problem [Baby and Wang, 2019].

References

- Dheeraj Baby and Yu-Xiang Wang. Online forecasting of total-variation-bounded sequences. In *Advances in Neural Information Processing Systems*, pages 11069–11079, 2019.
- Omar Besbes, Yonatan Gur, and Assaf Zeevi. Non-stationary stochastic optimization. *Operations research*, 63(5):1227–1244, 2015.
- Xi Chen, Yining Wang, and Yu-Xiang Wang. Nonstationary stochastic optimization under l_p, q -variation measures. *Operations Research*, 67(6):1752–1765, 2019.
- David L Donoho, Richard C Liu, and Brenda MacGibbon. Minimax risk over hyperrectangles, and implications. *The Annals of Statistics*, pages 1416–1437, 1990.
- Tianbao Yang, Lijun Zhang, Rong Jin, and Jinfeng Yi. Tracking slowly moving clairvoyant: Optimal dynamic regret of online learning with true and noisy gradient. In *International Conference on Machine Learning*, pages 449–457, 2016.
- Lijun Zhang, Shiyin Lu, and Zhi-Hua Zhou. Adaptive online learning in dynamic environments. In *Advances in neural information processing systems*, pages 1323–1333, 2018.
- Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the 20th international conference on machine learning (icml-03)*, pages 928–936, 2003.