

Online Learning

Zhe Li

November 18, 2014

1 Follow The Leader(FTL)

Follow The Leader algorithm is the most natural learning rule, which predicts the weight vector by attempting to minimize loss on all the previous trials. That is,

$$\forall t, w_t = \underset{w \in S}{\operatorname{argmin}} \sum_{i=1}^{t-1} f_i(w) \quad (1)$$

One introduces the regret to analyze the designed online learning algorithm, which measures how much difference the learner is compared with the best learner, which is defined by

$$\operatorname{Regret}_T(u) = \sum_{t=1}^T f_t(w_t) - f_t(u) \quad (2)$$

To analyze FTL, one can show the regret of FTL is upper bounded by the cumulative difference between the loss of w_t and w_{t+1} . Assume w_1, w_2, \dots be the sequence vectors produced by FTL, then, for all $u \in S$ we have

$$\operatorname{Regret}_T(u) = \sum_{t=1}^T (f_t(w_t) - f_t(u)) \leq \sum_{t=1}^T (f_t(w_t) - f_t(w_{t+1})) \quad (3)$$

This inequality is very important inequality used through all the following analysis for online learning algorithm. It can be understood as the loss suffered from any fixed predictive weight is greater or equal than loss from FTL since one can proof $\sum_{t=1}^T f_t(u) \geq \sum_{t=1}^T f_t(w_{t+1})$. Since it is valid for all u , one can choose u resulting in minimum loss, that is, u is the best learner, to analyze the designed online learning algorithm. FTL works well for Online Quadratic Optimization, but it fails to achieve sublinear regret in Online Linear Optimization, which can be shown by the devised example. FTL focuses mainly on minimizing loss incurred from past rounds without considering the generalization aspect. From the traditional machine learning view, FTL just focus on “training data”, even though one can get minimum error on “training data”, it does not necessary mean it will achieve the better performance. That’s the reason why one always add the penalty term with loss term in classical machine learning.

2 Follow The Regularized Leader

Follow The Regularized Leader (FTRL) might be inspired by this idea. FLRT predicts the weight vector at t trail by jointly minimizing the previous loss add regularizer term, formally,

$$\forall t, w_t = \underset{w \in S}{\operatorname{argmin}} \sum_{i=1}^{t-1} f_i(w) + R(w) \quad (4)$$

In the following, we will focus on the Euclidean regularization $R(w) = \frac{1}{2\eta} \|w\|^2$ and Entropic regularization $R(w) = \sum_{i=1}^d w_i \log w_i$. Specifically, for Online Linear Optimization, it turns out that FTRL with Euclidean regularization is Online Gradient Descent, since one can derive update rule $w_{t+1} = w_t - \eta \nabla f_t(w_t)$. It can be easily shown the regret from FTRL from the Inequality (3),

$$\operatorname{Regret}_T(u) = \sum_{t=1}^T (f_t(w_t) - f_t(u)) \leq R(u) - R(w_1) + \sum_{t=1}^T (f_t(w_t) - f_t(w_{t+1})) \quad (5)$$

Using this general form for Online Linear Optimization with Euclidean Norm $R(w) = \frac{1}{2\eta} \|w\|^2$

$$\begin{aligned} \operatorname{Regret}_T(u) &= \sum_{t=1}^T (f_t(w_t) - f_t(u)) \leq R(u) - R(w_1) + \sum_{t=1}^T (f_t(w_t) - f_t(w_{t+1})) \\ &= \frac{1}{2\eta} \|u\|^2 + \sum_{t=1}^T (w_t z_t - w_{t+1} z_t) \\ &= \frac{1}{2\eta} \|u\|^2 + \sum_{t=1}^T \langle w_t - w_{t+1}, z_t \rangle \quad (\text{update rule } w_{t+1} = w_t - \eta z_t) \\ &= \frac{1}{2\eta} \|u\|^2 + \eta \sum_{t=1}^T \|z_t\|_2^2 \end{aligned}$$

Let $U = \{u : \|u\| \leq B\}$ and let L be such that $\frac{1}{T} \sum_{t=1}^T \|z_t\|_2^2 \leq L^2$, and setting $\eta = \frac{B}{L\sqrt{2T}}$, we obtained the concise upper bound

$$\operatorname{Regret}_T(u) \leq BL\sqrt{2T} \quad (6)$$

We will generalize the above results to two aspects: any sequence of Lipschitz function instead of linear function with bounded norm and other regularization function. Firstly, we show an useful lemma

Lemma 1 *let $f : S \rightarrow \mathbb{R}$ be a convex function. Then, f is L -Lipschitz over S with respect to a norm $\|\cdot\|$ if and only if for all $w \in S$ and $z \in \partial f(w)$ we have that $\|z\|_\star \leq L$, where $\|\cdot\|_\star$ is the dual norm.*

Proof. Since f is convex function and $z \in \partial f(w)$ and $\|z\|_\star \leq L$, we have

$$f(w) - f(u) \leq (z, w - u) \leq \|z\|_\star \|w - u\| \leq L \|w - u\|$$

From definition of L -Lipschitz function, f is the L -Lipschitz function over S . For the other direction, we have f is L -Lipschitz function over S and $z \in \partial f(w)$, let u be such that $u - w = \underset{\|v\|=1}{\operatorname{argmin}}(v, z)$, therefore, $(u - w, z) = \|z\|_\star$, from the definition of the sub-gradient,

$$f(u) - f(w) \geq (z, u - w) = \|z\|_\star$$

Since Lipschitzness of f , we have

$$L \|u - w\| \geq f(u) - f(w)$$

Combine the above two inequalities, we conclude that $\|z\|_\star \leq L$. **Need to spend more time here**

Once we have the above lemma, if we run OGD on a sequence convex and L -Lipschitz function f_1, f_2, \dots then for all u , we have

$$\operatorname{Regret}_T(u) \leq \frac{1}{2\eta} \|u\|^2 + \eta \sum_{t=1}^T \|z_t\|_2^2 \leq \frac{1}{2\eta} \|u\|^2 + \eta T L^2$$

and in particular, if $U = u : \|u\|_2 \leq B$ and setting $\eta = \frac{B}{L\sqrt{2T}}$, then

$$\operatorname{Regret}_T(u) \leq B L \sqrt{2T}$$

3 Analyzing FTRL with Strong Convex Regularizers

From the above, we know

$$\operatorname{Regret}_T(u) = \sum_{t=1}^T (f_t(w_t) - f_t(u)) \leq R(u) - R(w_1) + \sum_{t=1}^T (f_t(w_t) - f_t(w_{t+1})) \quad (7)$$

If we add a strongly convex regularizers $R(w)$ to the function f_t , $F_t(w) = \sum_{i=1}^{t-1} f_i(w) + R(w)$ will be strongly convex function. From strongly convexity of $F_t(w)$ and Lipschitzness of f_t , we can obtain a general lemma:

Lemma 2 *Let f_1, \dots, f_T be a sequence of convex function such that f_t is L_t Lipschitz with respect to some norm $\|\cdot\|$, Let L be such that $\frac{1}{T} \sum_{t=1}^T L_t^2 \leq L^2$, Assume that FTLR is run on the sequence with δ -strongly convex regularization function with respect to same norm, then*

$$\operatorname{Regret}_T(u) \leq R(u) - \min_{v \in S} R(v) + \frac{T L^2}{\delta} \quad (8)$$

4 Online Mirror Descent

One possible disadvantage of the FTRL is that is more involved to solve an optimization problem at each online round. Using the linear function in which $f_t(w) = (w, z_t)$ with some regularization function $R(w)$. From FTRL, we know

$$\begin{aligned} w_{t+1} &= \underset{w}{\operatorname{argmin}} R(w) + \sum_{i=1}^t (w, z_i) \\ &= \underset{w}{\operatorname{argmin}} R(w) + (w, z_{1:t}) \quad (z_{1:t} = \sum_{i=1}^t z_i) \\ &= \underset{w}{\operatorname{argmax}} (w, -z_{1:t}) - R(w) \end{aligned}$$

Letting

$$g(\theta) = \underset{w}{\operatorname{argmax}} (w, \theta) - R(w), \quad (9)$$

So the FTRL prediction can be based on the following recursive update rule:

1. $\theta_{t+1} = \theta_t - z_t$
2. $w_{t+1} = g(\theta_{t+1})$

From this general online mirror descent framework, we can derive the Normalized Exponentiated Gradient which has the strong relationship with Weighted Majority algorithm in the expert learning scenario need to check this again

4.1 Normized Exponentiated Gradient

Let i^{th} component of $g(\theta)$ be

$$g_i(\theta) = \frac{e^{\eta\theta[i]}}{\sum_j e^{\eta\theta[j]}} \quad (10)$$

therefore,

$$\begin{aligned} w_{t+1}[i] &= \frac{e^{\eta\theta_{t+1}[i]}}{\sum_j e^{\eta\theta_{t+1}[j]}} \\ &= \frac{e^{\eta\theta_t[i]} e^{-\eta z_t[i]}}{\sum_j e^{\eta\theta_t[j]} e^{-\eta z_t[j]}} \\ &= \frac{\frac{e^{\eta\theta_t[i]}}{\sum_j e^{\eta\theta_t[j]}} e^{-\eta z_t[i]}}{\sum_j \frac{e^{\eta\theta_t[j]}}{\sum_j e^{\eta\theta_t[j]}} e^{-\eta z_t[j]}} \\ &= \frac{w_t[i] e^{-\eta z_t[i]}}{\sum_j w_t[j] e^{-\eta z_t[j]}} \end{aligned}$$

5 Appendix

For entropic regularization function $R(w) = \sum_{i=1}^d w_i \log w_i$, where $w \in S = \{w : w \geq 0, \sum_{i=1}^d w_i \leq B\}$, we attempt to obtain the maximum and minimum value of $R(w)$, Firstly, consider

$$\sum_{w \in S_{i=1}}^d w_i \log w_i \quad (11)$$

Construct the Lagrange function

$$L(w, \lambda) = \sum_{i=1}^d w_i \log w_i - \lambda \left(\sum_{i=1}^d w_i - B \right) \quad (12)$$

Computing the gradient of function L w.r.t w and set it to zero, we get

$$w_i = e^{\lambda-1} \quad (13)$$

that means every component of w is equal, that is, $w_i = \frac{B}{d}$, Plugging this into the objective function, we get minimum value of $R(w)$,

$$R(w)_{min} = B \log \frac{B}{d} \quad (14)$$

For maximum value, this convex function can achieve it at boundary, So the maximum value of $R(w)$ is

$$R(w)_{min} = B \log B \quad (15)$$

Here, we define $0 \log 0 = 0$

6 reference

<http://courses.cs.washington.edu/courses/cse599s/12sp/scribes.html>
<http://www.cs.huji.ac.il/~shais/papers/0Lsurvey.pdf>
<http://www.cs.tufts.edu/~roni/Teaching/CLT/LN/>