

5.7 Online mirror descent

Quadratic regularization is imposing a certain prior knowledge where there is a good w_t in a small L_2 ball.

Now we develop a general way of obtaining regret bounds for general regularizers.

Goal: analyse FTRL for \forall convex loss and regularizers.

[Algorithm 5] Online mirror descent (OMD)

Let $\psi: \mathbb{R}^d \rightarrow \mathbb{R}$ be the regularizer

Let f_1, \dots, f_T : sequence of losses.

On each iteration $t=1, \dots, T$, the learner chooses w_t s.t.

$$w_t \in \operatorname{argmin}_{w \in \mathbb{R}^d} \{ \psi(w) - w \cdot \theta_t \}. \quad (5.2.2)$$

where $z_t \in \partial f_t(w_t)$, $\theta_t = - \sum_{i=1}^{t-1} z_i$ (negative sum)

Technical Note: $S = \mathbb{R}^d$. And we can always fold constraints by setting $\psi(w) = \infty$ if w violates the constraints.

[Examples of regularizers]

- Quadratic : $\psi(w) = \frac{1}{2\eta} \|w\|_2^2$

- Non-spherical quadratic : $\psi(w) = \frac{1}{2\eta} w^T A w$

- Entropic : $\psi(w) = \frac{1}{\eta} \sum_{j=1}^d w_j \log w_j$ if $w \in \Delta_d$

Difference: the slope when $w \rightarrow$ boundary.

[Definition 27] Fenchel conjugate

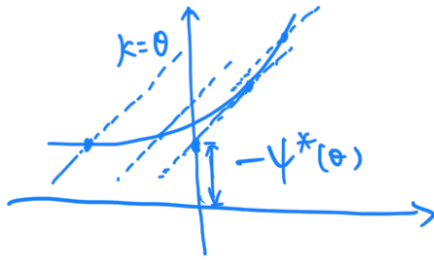
The Fenchel conjugate of a function $\psi: \mathbb{R}^d \rightarrow \mathbb{R}$ is

$$\psi^*(\theta) := \sup_{w \in \mathbb{R}^d} \{ w \cdot \theta - \psi(w) \}$$

Intuition: For scalars $w, \theta \in \mathbb{R}$, fixed θ (as slope), $-\psi^*(\theta)$

is the maximum value of $w \cdot \theta - \psi(w)$ over all $w \in \mathbb{R}$.

is the position where supporting hyperplane of ψ with slope θ hit the vertical axis.



Useful facts:

- ① ψ^* is always convex, since it's a supremum over a collection of linear function $\theta \mapsto w \cdot \theta - f(w)$
- ② $\psi^*(\theta) \geq w \cdot \theta - \psi(w)$ for all $w \in \mathbb{R}^d$ (French-Young inequality)
- ③ If $r(w) = a\psi(w)$ with $a > 0$, then $r^*(\theta) = a\psi^*(\theta/a)$
- ④ $\psi^{**} = \psi$ iff ψ is convex (and low semi-continuous)
- ⑤ If ψ is differentiable, then

$$\nabla \psi^*(\theta) = \arg \max_{w \in \mathbb{R}^d} \{w \cdot \theta - \psi(w)\} \quad \{w \text{ is the gradient}\}$$

Mirroring

(i) OMD update (522): $w_t \in \arg \min_{w \in S} \{\psi(w) - w \cdot \theta_t\}$

We have $w_t = \nabla \psi^*(\theta_t)$ and $-\psi^*(\theta_t)$ is the corresponding value of regularized loss

(ii) Since w_t attains $\sup \{w \cdot \theta_t - \psi(w)\}$, differentiate w.r.t. w and

we have $\theta_t = \nabla \psi(w_t)$

(iii) One-to-one mapping:

$$w_t = \nabla \psi^*(\theta_t), \quad \theta_t = \nabla \psi(w_t)$$

(iv) OMD updates: $\theta_{t+1} = \theta_t - z_t$, $\theta_1, \dots, \theta_T$

by mirrored: $w_t = \nabla \psi^*(\theta_t)$, w_1, \dots, w_T .

[Example 32] Quadratic

$$\text{Let } \psi(w) = \frac{1}{2\eta} \|w\|_2^2$$

$$\text{Then } \psi^*(\theta) = \sup_{w \in S} \{w \cdot \theta - \frac{1}{2\eta} \|w\|_2^2\} = \frac{\eta}{2} \|\theta\|_2^2$$

[Example 33] Entropic

$$\text{Let } \psi(w) = \frac{1}{\eta} \sum_{j=1}^d w_j \log w_j \text{ for } w \in \Delta_d$$

$$\max \{w \cdot \theta - \psi(w)\} \text{ s.t. } w \in \Delta_d.$$

$$\text{Let } \theta_i - \frac{1}{\eta} (\log w_i + 1) = c, \quad w_i = e^{\eta \theta_i - \eta c - 1}$$

$$\sum w_i = 1.$$

$$\Rightarrow w_i = \frac{e^{\eta \theta_i}}{\sum_j e^{\eta \theta_j}}$$

$$\Rightarrow \psi^*(w) = \sum_i \frac{\theta_i e^{\eta \theta_i}}{\sum_j e^{\eta \theta_j}} - \frac{1}{\eta} \sum_i \frac{e^{\eta \theta_i}}{\sum_j e^{\eta \theta_j}} [\eta \theta_i - \log \sum_j e^{\eta \theta_j}]$$

$$\Rightarrow \psi^*(w) = \frac{1}{\eta} \log \left(\sum_{j=1}^d e^{\eta \theta_j} \right)$$