

Lecture 16: Mirror Descent

Lecturer: M. K. Hanawal

Scribes: Shubham Uttam

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.*

16.1 Recap

In the last class we saw the Online gradient descent algorithm

Algorithm 1 Online Gradient Descent

```

1: Start with  $w_1 \in \mathcal{K} \subset \mathbb{R}^d$ 
2: for rounds  $t=1, \dots, T$  do
3:   player plays  $w_t$ 
4:   environment plays  $w_t$ 
5:    $w_{t+1} \leftarrow \text{Proj}_k[w_t - \eta_t \nabla c_t(w_t)]$ 

```

Note that here the set \mathcal{K} is closed and bounded and the cost function c_t is strongly convex and continuously differentiable. When c_t is convex but not continuously differentiable then we get online sub-gradient descent algorithm by making some changes in the update statement.

Algorithm 2 Online Sub-Gradient Descent

```

1: Start with  $w_1 \in \mathcal{K} \subset \mathbb{R}^d$ 
2: for rounds  $t=1, \dots, T$  do
3:   player plays  $w_t$ 
4:   environment plays  $w_t$ 
5:    $w_{t+1} \leftarrow \text{Proj}_k[w_t - \eta_t \nabla c_t(w_t)]$ 

```

16.2 Some basics of Linear Algebra

Definition 16.1 Field: A field F is a set with two operations addition and multiplication,

$$+ : F \times F \rightarrow F \text{ and } \cdot : F \times F \rightarrow F$$

which obey the following axioms:

- $(F, +)$ is an **abelian group** under addition:

(1) Addition is associative. That is for every x, y and $z \in F$,

$$(x + y) + z = x + (y + z).$$

- (2) There is an identity element under addition. This element is often denoted $0 \in F$ and for every element $x \in F$,

$$0 + x = x + 0 = x.$$

- (3) Every element has an additive inverse. That is given $x \in F$, there is an element $-x \in F$ and

$$x + (-x) = x - x = 0$$

- (4) Addition is commutative. That is given x and $y \in F$,

$$x + y = y + x.$$

- Let $F^* = F - \{0\}$. Then (F^*, \cdot) is an abelian group under multiplication:

- (1) Multiplication is associative. That is for every x, y and $z \in F$,

$$(x \cdot y) \cdot z = x \cdot (y \cdot z)$$

- (2) There is an identity element under multiplication. This element is often denoted $1 \in F$ and for every element $x \in F$,

$$1 \cdot x = x \cdot 1 = x$$

- (3) Every element has an multiplicative inverse. That is given $x \in F$, there is an element $x^{-1} \in F$ and

$$x \cdot x^{-1} = x^{-1} \cdot x = 1$$

- (4) multiplication is commutative. That is given x and $y \in F$,

$$x \cdot y = y \cdot x$$

- addition and multiplication are compatible, i.e. F satisfies the distributive law. that is given x, y and $z \in F$,

$$x \cdot (y + z) = x \cdot y + x \cdot z$$

Definition 16.2 Vector Space: A vector space over a field F is a set V together with two operations, vector addition and vector multiplication that satisfy the eight axioms:

Let $u, v, w \in V$ and $a, b \in F$

- Associativity of addition:

$$u + (v + w) = (u + v) + w$$

- Commutativity of addition:

$$u + v = v + u$$

- Identity element of addition:

there exists an element $0 \in V$, called zero vector, such that $v + 0 = v$ for all $v \in V$

- Inverse elements of addition:

For every $v \in V$, there exists an element $-v \in V$, called the additive inverse of v , such that $v + (-v) = 0$.

- Compatibility of scalar multiplication with field multiplication:

$$a(bv) = (ab)v$$

- *Identity element of scalar multiplication:*
 $1v=v$, where 1 denotes the multiplicative identity in F .
- *Distributivity of scalar multiplication with respect to vector addition:*

$$a(u+v) = au + av$$

- *Distributivity of scalar multiplication with respect to field addition:*

$$(a+b)v = av + bv$$

Definition 16.3 Inner Product Space: An inner product space is a vector space V over the field F together with an inner product, i.e., with a map

$$\langle \cdot, \cdot \rangle: V \times V \rightarrow F$$

that satisfies the following three axioms for all vectors $x, y, z \in V$ and all scalars $a \in F$:

- *Conjugate symmetry:*

$$\langle x, y \rangle = \overline{\langle y, x \rangle}$$

- *Linearity in the first argument:*

$$\langle ax, y \rangle = a \langle x, y \rangle$$

$$\langle x + y, z \rangle = \langle x, z \rangle + \langle y, z \rangle$$

- *Positive-definiteness:*

$$\langle x, x \rangle \geq 0$$

$$\langle x, x \rangle = 0 \Leftrightarrow x = 0$$

Definition 16.4 Complete Space: A metric space M is called complete (or a Cauchy space) if every Cauchy sequence of points in M has a limit that is also in M or, alternatively, if every Cauchy sequence in M converges in M .

Definition 16.5 Hilbert Space: A complete space with an inner product is called a Hilbert space.

Definition 16.6 Linear Functional: A map $f: V \rightarrow F$ over a vector space V over F is called a linear functional if $f(\alpha u + \beta v) = \alpha f(u) + \beta f(v)$, $\forall u, v \in V$, $\forall \alpha, \beta \in F$.

Definition 16.7 Dual Space: Let H be a Hilbert space. Then the set of all linear functional over H forms the dual space H^* .

Riesz Representation Theorem:

Let H be a Hilbert space, and let H^* be its dual space. If x is an element of H , then the function φ_x , for all y in H defined by:

$$\varphi_x(y) = \langle y, x \rangle$$

Theorem: The mapping $\phi: H \rightarrow H^*$ defined by $\phi(x) = \varphi_x$ is an isometric isomorphism, meaning that:

- ϕ is bijective.
- The norm of x and ϕ_x agree: $\|x\| = \|\phi(x)\|$.
- ϕ is additive: $\phi(x_1 + x_2) = \phi(x_1) + \phi(x_2)$.
- If the base field is \mathbb{R} , then $\phi(\lambda x) = \lambda \phi(x)$ for all real numbers λ .
- If the base field is \mathbb{C} , then $\phi(\lambda x) = \bar{\lambda} \phi(x)$ for all complex numbers λ , where $\bar{\lambda}$ denotes the complex conjugation of λ .

Definition 16.8 Banach space: A Banach space is a vector space X over the field \mathbb{R} , or over the field \mathbb{C} , which is equipped with a norm and which is complete with respect to that norm, that is to say, for every Cauchy sequence $\{x_n\}$ in X , there exists an element x in X such that

$$\lim_{n \rightarrow \infty} x_n = x,$$

or equivalently:

$$\lim_{n \rightarrow \infty} \|x_n - x\|_x = 0$$

If B is Banach space, then Riesz Representation theorem may not hold.

If g is a linear functional on \mathbb{R}^d , then:

$$g(y) = y_1 g(e_1) + y_2 g(e_2) + \dots + y_d g(e_d)$$

where $y = y_1 e_1 + y_2 e_2 + \dots + y_d e_d$ where $\{e_1, e_2, \dots, e_d\}$ is the set of basis of \mathbb{R}^d .

16.3 Mirror Descent

Definition 16.9 Bregman Divergence: Let $f : \omega \rightarrow \mathbb{R}$ be a function that is strictly convex, continuously differentiable and defined on a closed convex set ω . Then the Bregman divergence is defined as:

$$B_f(x, y) := f(x) - f(y) - \langle \nabla f(y), x - y \rangle, \forall x, y \in \omega$$

That is, the difference between the value of f at x and the first order Taylor expansion of f around y evaluated at point x .

Examples of Bregman divergence

- Euclidean distance. Let $f(x) = \frac{1}{2} \|x\|^2$. Then Bregman divergence $B_f(x, y) = \frac{1}{2} \|x - y\|^2$
- Non Euclidean distance, Mahalanobis distance. We have $f(x) = \frac{1}{2} x^T A x$, for some positive semidefinite matrix A . Then Bregman divergence is

$$B_f(x, y) = \frac{1}{2} (x - y)^T A (x - y)$$

- Kullback-Leibler divergence. We have $f(p) = \sum p_i \log p_i$. Then Bregman divergence is

$$B_f(p, q) = \sum p_i \log \frac{p_i}{q_i} - \sum p_i + \sum q_i$$

Two components which we will use in mirror descent procedure:

1. Mirror function, ϕ

It is a convex function defined as, $\phi : V \rightarrow V^*$ where V^* is the dual space of V .

2. Bregman projection

$$Proj_K(u) = \arg \min_{v \in K} B_f(v, u)$$

Definition 16.10 Fenchel-Legendre dual: A function $f : \mathcal{H} \rightarrow \mathbb{R}$, define $f^* : \mathcal{H}^* \rightarrow \mathbb{R}$ such that

$$f^*(u) = \sup_{x \in \mathcal{H}} \{f(x) - \langle u, x \rangle\}$$

If function f is convex and differentiable, then $f^*(\nabla f(x)) = f(x) - \langle \nabla f(x), x \rangle$

Motivation of mirror descent procedure: If there is a case where $\nabla c_t(w_t) \notin K$, ($K \neq \mathbb{R}^d$) then nothing can be said about $w_t - \eta_t \nabla c_t(w_t)$, so we conduct update in dual space instead of real space.

16.3.1 Online Mirror Descent Procedure

- Update in the dual space
- Take a mirror function, ϕ .
- Mirror Update: $\nabla \phi(\tilde{w}_{t+1}) \leftarrow \nabla \phi(w_t) - \eta_t \nabla c_t(w_t)$
- $w_{t+1} \leftarrow \text{BregmanProj}_K^\phi[\tilde{w}_t]$

If we take $\phi(x) = \frac{1}{2}||x||^2$, then it becomes Online Descent Algorithm and if we take ϕ to be a negative entropy term in mirror update then algorithm we get is known as Online Exponential Gradient algorithm.

16.3.2 Regret of OMD

Let environment reveals a convex regret loss function as c_t in round t , then the regret for online mirror descent is given by:

$$\text{Regret}_{\text{OMD}}(u, T) = \sum_{t=1}^T c_t(w_t) - \sum_{t=1}^T c_t(u)$$

16.3.3 Regret Bound for OMD

$$\begin{aligned}
\text{Consider } c_t(w_t) - c_t(u) &\leq \langle \nabla c_t(w_t), w_t - u \rangle \quad (\text{By convexity of } c_t) \\
&= \frac{1}{n} \langle \nabla \phi(w_t) - \nabla \phi(\tilde{w}_{t+1}), w_t - u \rangle \quad (\text{since } \nabla \phi(\tilde{w}_{t+1}) = \nabla \phi(w_t) - \eta \nabla \phi(w_t)) \\
&= \frac{1}{n} [B_\phi(u, w_t) - B_\phi(u, \tilde{w}_{t+1}) + B_\phi(w_t, \tilde{w}_{t+1})] \quad (\text{using Bregman divergence}) \\
&= \frac{1}{n} [B_\phi(u, w_t) - B_\phi(w_{t+1}, \tilde{w}_{t+1}) - B_\phi(u, w_{t+1}) + B_\phi(w_t, \tilde{w}_{t+1})] \\
&\quad (\text{as } B_\phi(u, \tilde{w}_{t+1}) = B_\phi(w_{t+1}, \tilde{w}_{t+1}) + B_\phi(u, w_{t+1}))
\end{aligned}$$

Summing over T:

$$\begin{aligned}
\sum_{t=1}^T (c_t(w_t) - c_t(u)) &\leq \frac{1}{n} \sum_{t=1}^T [B_\phi(u, w_t) - B_\phi(u, w_{t+1})] + \frac{1}{n} \sum_{t=1}^T [B_\phi(w_t, \tilde{w}_{t+1}) - B_\phi(w_{t+1}, \tilde{w}_{t+1})] \\
&= \frac{1}{n} [B_\phi(u, w_1) - B_\phi(u, w_{T+1})] + \frac{1}{n} \sum_{t=1}^T [B_\phi(w_t, \tilde{w}_{t+1}) - B_\phi(w_{t+1}, \tilde{w}_{t+1})] \\
&\leq \frac{1}{n} [B_\phi(u, w_1)] + \frac{1}{n} \sum_{t=1}^T [B_\phi(w_t, \tilde{w}_{t+1}) - B_\phi(w_{t+1}, \tilde{w}_{t+1})] \\
&\quad (\text{as for some convex } \phi, B_\phi(u, v) \geq 0)
\end{aligned}$$

For $w_1 : \phi(u) - \phi(w_1) \leq D^2, \quad \forall u \in K$

$B_\phi(u, w) = \phi(u) - \phi(w) - \langle \nabla \phi(w), u - w \rangle \leq D^2 \quad (\text{as } \langle \nabla \phi(w), u - w \rangle \geq 0)$

Therefore we get:

$$\begin{aligned}
\sum_{t=1}^T (c_t(w_t) - c_t(u)) &\leq \frac{1}{n} D^2 + \frac{1}{n} \sum_{t=1}^T [B_\phi(w_t, \tilde{w}_{t+1}) - B_\phi(w_{t+1}, \tilde{w}_{t+1})] \\
&= \frac{1}{n} D^2 + \frac{1}{n} \sum_{t=1}^T [\phi(w_t) - \phi(w_{t+1}) - \langle \nabla \phi(\tilde{w}_{t+1}), w_t - w_{t+1} \rangle]
\end{aligned}$$

Assumption: ϕ is strongly convex with modulus α

i.e. $\phi(y) \geq \phi(x) + \langle \nabla \phi(x), y - x \rangle + \frac{\alpha}{2} \|y - x\|^2$

$$\text{Therefore, } \sum_{t=1}^T (c_t(w_t) - c_t(u)) \leq \frac{1}{\eta} D^2 + \sum_{t=1}^T [\langle \nabla \phi(w_t), w_t - w_{t+1} \rangle - \frac{\alpha}{2} \|w_t - w_{t+1}\|^2 - \langle \nabla \phi(\tilde{w}_{t+1}), w_t - w_{t+1} \rangle]$$

Recall update step:

$$\nabla \phi(\tilde{w}_{t+1}) = \nabla \phi(w_t) - \eta \nabla c_t(w_t)$$

$$\nabla \phi(w_t) - \nabla \phi(\tilde{w}_{t+1}) = \eta \nabla c_t(w_t)$$

$$\begin{aligned}
\text{Therefore, } \sum_{t=1}^T (c_t(w_t) - c_t(u)) &\leq \frac{1}{\eta} D^2 + \sum_{t=1}^T [\langle \nabla \phi(w_t) - \nabla \phi(\tilde{w}_{t+1}), w_t - w_{t+1} \rangle - \frac{\alpha}{2} \|w_t - w_{t+1}\|^2] \\
&= \frac{1}{\eta} D^2 + \sum_{t=1}^T [\eta \langle \nabla c_t(w_t), w_t - w_{t+1} \rangle - \frac{\alpha}{2} \|w_t - w_{t+1}\|^2] \\
&= \frac{1}{\eta} D^2 + \sum_{t=1}^T [\langle \eta \nabla c_t(w_t), w_t - w_{t+1} \rangle - \frac{\alpha}{2} \|w_t - w_{t+1}\|^2]
\end{aligned}$$

(16.1)

Now by Cauchy Schwartz Inequality, we have:

$$| \langle x, y \rangle |^2 \leq \|x\| \|y\|$$

So then we have:

$$\begin{aligned} \langle \eta \nabla c_t(w_t), w_t - w_{t+1} \rangle &= \eta \langle \nabla c_t(w_t), w_t - w_{t+1} \rangle \\ \langle \eta \nabla c_t(w_t), w_t - w_{t+1} \rangle &\leq \eta \| \nabla c_t(w_t) \| \|w_t - w_{t+1}\| \end{aligned}$$

Now as $\| \nabla c_t(w_t) \|$ is bounded by G , we have

$$\langle \eta \nabla c_t(w_t), w_t - w_{t+1} \rangle \leq \eta G \|w_t - w_{t+1}\|$$

So we have:

$$\begin{aligned} \langle \eta \nabla c_t(w_t), w_t - w_{t+1} \rangle &\leq -\frac{\alpha}{2} \|w_t - w_{t+1}\|^2 \leq \frac{\eta^2 G^2}{2\alpha} + \frac{\alpha}{2} \|w_t - w_{t+1}\|^2 - \eta G \|w_t - w_{t+1}\|^2 \\ &= \left(\frac{\eta G}{\sqrt{2\alpha}} - \sqrt{\frac{\alpha}{2}} \|w_t - w_{t+1}\| \right)^2 \end{aligned}$$

Therefore the regret bound we have is:

$$Reg = \sum_{t=1}^T (c_t(w_t) - c_t(u)) \leq \frac{1}{\eta} D^2 + \sum_{t=1}^T \frac{\eta^2 G^2}{2\alpha} = \frac{1}{\eta} \left[D^2 + \frac{\eta^2 G^2 T}{2\alpha} \right]$$

16.4 Follow the Regularized Leader

Follow the regularized leader minimizes the loss of all past rounds plus a regularization term. The goal of the regularization term is to stabilize the solution. Formally, for a regularization function, $R : S \rightarrow \mathbb{R}$ we define the weight for any round t as:

$$w_t = \arg \min_{w \in S} \sum_{i=1}^{t-1} f_i(w) + R(w)$$