
Bregman Divergence and Mirror Descent

1 Bregman Divergence

Motivation

- Generalize squared Euclidean distance to a class of distances that all share similar properties
- Lots of applications in machine learning, clustering, exponential family

Definition 1 (Bregman divergence) Let $\psi : \Omega \rightarrow \mathbb{R}$ be a function that is: a) strictly convex, b) continuously differentiable, c) defined on a closed convex set Ω . Then the Bregman divergence is defined as

$$\Delta_\psi(x, y) = \psi(x) - \psi(y) - \langle \nabla \psi(y), x - y \rangle, \quad \forall x, y \in \Omega. \quad (1)$$

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

Examples

- Euclidean distance. Let $\psi(x) = \frac{1}{2} \|x\|^2$. Then $\Delta_\psi(x, y) = \frac{1}{2} \|x - y\|^2$.
- $\Omega = \{x \in \mathbb{R}_+^n : \sum_i x_i = 1\}$, and $\psi(x) = \sum_i x_i \log x_i$. Then $\Delta_\psi(x, y) = \sum_i x_i \log \frac{x_i}{y_i}$ for $x, y \in \Omega$. This is called relative entropy, or Kullback–Leibler divergence between probability distributions x and y .
- ℓ_p norm. Let $p \geq 1$ and $\frac{1}{p} + \frac{1}{q} = 1$. $\psi(x) = \frac{1}{2} \|x\|_q^2$. Then $\Delta_\psi(x, y) = \frac{1}{2} \|x\|_q^2 + \frac{1}{2} \|y\|_q^2 - \langle x, \nabla \frac{1}{2} \|y\|_q^2 \rangle$. Note $\frac{1}{2} \|y\|_q^2$ is not necessarily continuously differentiable, which makes this case not precisely consistent with our definition.

Properties of Bregman divergence

- Strict convexity in the first argument x . Trivial by the strict convexity of ψ .
- Nonnegativity: $\Delta_\psi(x, y) \geq 0$ for all x, y . $\Delta_\psi(x, y) = 0$ if and only if $x = y$. Trivial by strict convexity.
- Asymmetry: in general, $\Delta_\psi(x, y) \neq \Delta_\psi(y, x)$. Eg, KL-divergence. Symmetrization not always useful.
- Non-convexity in the second argument. Let $\Omega = [1, \infty)$, $\psi(x) = -\log x$. Then $\Delta_\psi(x, y) = -\log x + \log y + \frac{x-y}{y}$. One can check its second order derivative in y is $\frac{1}{y^2}(\frac{2x}{y} - 1)$, which is negative when $2x < y$.
- Linearity in ψ . For any $a > 0$, $\Delta_{\psi+a\varphi}(x, y) = \Delta_\psi(x, y) + a\Delta_\varphi(x, y)$.
- Gradient in x : $\frac{\partial}{\partial x} \Delta_\psi(x, y) = \nabla \psi(x) - \nabla \psi(y)$. Gradient in y is trickier, and not commonly used.
- Generalized triangle inequality:

$$\Delta_\psi(x, y) + \Delta_\psi(y, z) = \psi(x) - \psi(y) - \langle \nabla \psi(y), x - y \rangle + \psi(y) - \psi(z) - \langle \nabla \psi(z), y - z \rangle \quad (2)$$

$$= \Delta_\psi(x, z) + \langle x - y, \nabla \psi(z) - \nabla \psi(y) \rangle. \quad (3)$$

- Special case: ψ is called strongly convex with respect to some norm with modulus σ if

$$\psi(x) \geq \psi(y) + \langle \nabla \psi(y), x - y \rangle + \frac{\sigma}{2} \|x - y\|^2. \quad (4)$$

Note the norm here is not necessarily Euclidean norm. When the norm is Euclidean, this condition is equivalent to $\psi(x) - \frac{\sigma}{2} \|x\|^2$ being convex. For example, the $\psi(x) = \sum_i x_i \log x_i$ used in KL-divergence is 1-strongly convex over the simplex $\Omega = \{x \in \mathbb{R}_+^n : \sum_i x_i = 1\}$, with respect to the ℓ_1 norm. When ψ is σ strong convex, we have

$$\Delta_\psi(x, y) \geq \frac{\sigma}{2} \|x - y\|^2. \quad (5)$$

Proof: By definition $\Delta_\psi(x, y) = \psi(x) - \psi(y) - \langle \nabla \psi(y), x - y \rangle \geq \frac{\sigma}{2} \|x - y\|^2$. ■

- Duality. Suppose ψ is strongly convex. Then

$$(\nabla \psi^*)(\nabla \psi(x)) = x, \quad \Delta_\psi(x, y) = \Delta_{\psi^*}(\nabla \psi(y), \nabla \psi(x)). \quad (6)$$

Proof: (for the first equality only) Recall

$$\psi^*(y) = \sup_{z \in \Omega} \{\langle z, y \rangle - \psi(z)\}. \quad (7)$$

sup must be attainable because ψ is strongly convex and Ω is closed. x is a maximizer if and only if $y = \nabla \psi(x)$. So

$$\psi^*(y) + \psi(x) = \langle x, y \rangle \Leftrightarrow y = \nabla \psi(x). \quad (8)$$

Since $\psi = \psi^{**}$, so $\psi^*(y) + \psi^{**}(x) = \langle x, y \rangle$, which means y is the maximizer in

$$\psi^{**}(x) = \sup_z \{\langle x, z \rangle - \psi^*(z)\}. \quad (9)$$

This means $x = \nabla \psi^*(y)$. To summarize, $(\nabla \psi^*)(\nabla \psi(x)) = x$. ■

- Mean of distribution. Suppose U is a random variable over an open set S with distribution μ . Then

$$\min_{x \in S} \mathbb{E}_{U \sim \mu} [\Delta_\psi(U, x)] \quad (10)$$

is optimized at $\bar{u} := \mathbb{E}_\mu[U] = \int_{u \in S} u \mu(u)$.

Proof: For any $x \in S$, we have

$$\mathbb{E}_{U \sim \mu} [\Delta_\psi(U, x)] - \mathbb{E}_{U \sim \mu} [\Delta_\psi(U, \bar{u})] \quad (11)$$

$$= \mathbb{E}_\mu [\psi(U) - \psi(x) - (U - x)' \nabla \psi(x) - \psi(U) + \psi(\bar{u}) + (U - \bar{u})' \nabla \psi(\bar{u})] \quad (12)$$

$$= \psi(\bar{u}) - \psi(x) + x' \nabla \psi(x) - \bar{u}' \nabla \psi(\bar{u}) + \mathbb{E}_\mu [-U' \nabla \psi(x) + U' \nabla \psi(\bar{u})] \quad (13)$$

$$= \psi(\bar{u}) - \psi(x) - (\bar{u} - x)' \nabla \psi(x) \quad (14)$$

$$= \Delta_\psi(\bar{u}, x). \quad (15)$$

This must be nonnegative, and is 0 if and only if $x = \bar{u}$. ■

- Pythagorean Theorem. If x^* is the projection of x_0 onto a convex set $C \in \Omega$:

$$x^* = \operatorname{argmin}_{x \in C} \Delta_\psi(x, x_0). \quad (16)$$

Then

$$\Delta_\psi(y, x_0) \geq \Delta_\psi(y, x^*) + \Delta_\psi(x^*, x_0). \quad (17)$$

In Euclidean case, it means the angle $\angle yx^*x_0$ is obtuse. More generally

Lemma 2 Suppose L is a proper convex function whose domain is an open set containing C . L is not necessarily differentiable. Let x^* be

$$x^* = \operatorname{argmin}_{x \in C} \{L(x) + \Delta_\psi(x^*, x_0)\}. \quad (18)$$

Then for any $y \in C$ we have

$$L(y) + \Delta_\psi(y, x_0) \geq L(x^*) + \Delta_\psi(x^*, x_0) + \Delta_\psi(y, x^*). \quad (19)$$

The projection in (16) is just a special case of $L = 0$. This property is the key to the analysis of many optimization algorithms using Bregman divergence.

Proof: Denote $J(x) = L(x) + \Delta_\psi(x, x_0)$. Since x^* minimizes J over C , there must exist a subgradient $d \in \partial J(x^*)$ such that

$$\langle d, x - x^* \rangle \geq 0, \quad \forall x \in C. \quad (20)$$

Since $\partial J(x^*) = \{g + \nabla_{x=x^*} \Delta_\psi(x, x_0) : g \in \partial L(x^*)\} = \{g + \nabla\psi(x^*) - \nabla\psi(x_0) : g \in \partial L(x^*)\}$. So there must be a subgradient $g \in \partial L(x^*)$ such that

$$\langle g + \nabla\psi(x^*) - \nabla\psi(x_0), x - x^* \rangle \geq 0, \quad \forall x \in C. \quad (21)$$

Therefore using the property of subgradient, we have for all $y \in C$ that

$$L(y) \geq L(x^*) + \langle g, y - x^* \rangle \quad (22)$$

$$\geq L(x^*) + \langle \nabla\psi(x_0) - \nabla\psi(x^*), y - x^* \rangle \quad (23)$$

$$= L(x^*) - \langle \nabla\psi(x_0), x^* - x_0 \rangle + \psi(x^*) - \psi(x_0) \quad (24)$$

$$+ \langle \nabla\psi(x_0), y - x_0 \rangle - \psi(y) + \psi(x_0) \quad (25)$$

$$- \langle \nabla\psi(x^*), y - x^* \rangle + \psi(y) - \psi(x^*) \quad (26)$$

$$= L(x^*) + \Delta_\psi(x^*, x_0) - \Delta_\psi(y, x_0) + \Delta_\psi(y, x^*). \quad (27)$$

Rearranging completes the proof. ■

2 Mirror Descent

Why bother? Because the rate of convergence of subgradient descent often depends on the dimension of the problem.

Suppose we want to minimize a function f over a set C . Recall the subgradient descent rule

$$x_{k+\frac{1}{2}} = x_k - \alpha_k g_k, \quad \text{where } g_k \in \partial f(x_k) \quad (28)$$

$$x_{k+1} = \operatorname{argmin}_{x \in C} \frac{1}{2} \|x - x_{k+\frac{1}{2}}\|^2 = \operatorname{argmin}_{x \in C} \frac{1}{2} \|x - (x_k - \alpha_k g_k)\|^2. \quad (29)$$

This can be interpreted as follows. First approximate f around x_k by a first-order Taylor expansion

$$f(x) \approx f(x_k) + \langle g_k, x - x_k \rangle. \quad (30)$$

Then penalize the displacement by $\frac{1}{2\alpha_k} \|x - x_k\|^2$. So the update rule is to find a regularized minimizer of the model

$$x_{k+1} = \operatorname{argmin}_{x \in C} \left\{ f(x_k) + \langle g_k, x - x_k \rangle + \frac{1}{2\alpha_k} \|x - x_k\|^2 \right\}. \quad (31)$$

It is trivial to see this is exactly equivalent to (29). To generalize the method beyond Euclidean distance, it is straightforward to use the Bregman divergence as a measure of displacement:

$$x_{k+1} = \operatorname{argmin}_{x \in C} \left\{ f(x_k) + \langle g_k, x - x_k \rangle + \frac{1}{\alpha_k} \Delta_\psi(x, x_k) \right\} \quad (32)$$

$$= \operatorname{argmin}_{x \in C} \{ \alpha_k f(x_k) + \alpha_k \langle g_k, x - x_k \rangle + \Delta_\psi(x, x_k) \}. \quad (33)$$

Mirror descent interpretation. Suppose the constraint set C is the whole space (*i.e.* no constraint). Then we can take gradient with respect to x and find the optimality condition

$$g_k + \frac{1}{\alpha_k} (\nabla\psi(x_{k+1}) - \nabla\psi(x_k)) = 0 \quad (34)$$

$$\Leftrightarrow \nabla\psi(x_{k+1}) = \nabla\psi(x_k) - \alpha_k g_k \quad (35)$$

$$\Leftrightarrow x_{k+1} = (\nabla\psi)^{-1}(\nabla\psi(x_k) - \alpha_k g_k) = (\nabla\psi^*)(\nabla\psi(x_k) - \alpha_k g_k). \quad (36)$$

For example, in KL-divergence over simplex, the update rule becomes

$$x_{k+1}(i) = x_k(i) \exp(-\alpha_k g_k(i)). \quad (37)$$

Rate of convergence. Recall in unconstrained subgradient descent we followed 4 steps.

1. Bound on single update

$$\|x_{k+1} - x^*\|_2^2 = \|x_k - \alpha_k g_k - x^*\|_2^2 \quad (38)$$

$$= \|x_k - x^*\|_2^2 - 2\alpha_k \langle g_k, x_k - x^* \rangle + \alpha_k^2 \|g_k\|_2^2 \quad (39)$$

$$\leq \|x_k - x^*\|_2^2 - 2\alpha_k (f(x_k) - f(x^*)) + \alpha_k^2 \|g_k\|_2^2. \quad (40)$$

2. Telescope

$$\|x_{T+1} - x^*\|_2^2 \leq \|x_1 - x^*\|_2^2 - 2 \sum_{k=1}^T \alpha_k (f(x_k) - f(x^*)) + \sum_{k=1}^T \alpha_k^2 \|g_k\|_2^2. \quad (41)$$

3. Bounding by $\|x_1 - x^*\|_2^2 \leq R^2$ and $\|g_k\|_2^2 \leq G^2$:

$$2 \sum_{k=1}^T \alpha_k (f(x_k) - f(x^*)) \leq R^2 + G^2 \sum_{k=1}^T \alpha_k^2. \quad (42)$$

4. Denote $\epsilon_k = f(x_k) - f(x^*)$ and rearrange

$$\min_{k \in \{1, \dots, T\}} \epsilon_k \leq \frac{R^2 + G^2 \sum_{k=1}^T \alpha_k^2}{2 \sum_{k=1}^T \alpha_k}. \quad (43)$$

By setting the step size α_k judiciously, we can achieve

$$\min_{k \in \{1, \dots, T\}} \epsilon_k \leq \frac{RG}{\sqrt{T}}. \quad (44)$$

Suppose C is the simplex. Then $R \leq \sqrt{2}$. If each coordinate of each gradient g_i is upper bounded by M , then G can be at most $M\sqrt{n}$, *i.e.* depends on the dimension.

Clearly the step 2 to 4 can be easily extended by replacing $\|x_{k+1} - x^*\|_2^2$ with $\Delta_\psi(x^*, x_{k+1})$. So the only challenge left is to extend step 1. This is actually possible via Lemma 2.

We further assume ψ is σ strongly convex. In (33), consider $\alpha_k(f(x_k) + \langle g_k, x - x_k \rangle)$ as the L in Lemma 2. Then

$$\alpha_k (f(x_k) + \langle g_k, x^* - x_k \rangle) + \Delta_\psi(x^*, x_k) \geq \alpha_k (f(x_k) + \langle g_k, x_{k+1} - x_k \rangle) + \Delta_\psi(x_{k+1}, x_k) + \Delta_\psi(x^*, x_{k+1}) \quad (45)$$

Canceling some terms can rearranging, we obtain

$$\Delta_\psi(x^*, x_{k+1}) \leq \Delta_\psi(x^*, x_k) + \alpha_k \langle g_k, x^* - x_{k+1} \rangle - \Delta_\psi(x_{k+1}, x_k) \quad (46)$$

$$= \Delta_\psi(x^*, x_k) + \alpha_k \langle g_k, x^* - x_k \rangle + \alpha_k \langle g_k, x_k - x_{k+1} \rangle - \Delta_\psi(x_{k+1}, x_k) \quad (47)$$

$$\leq \Delta_\psi(x^*, x_k) - \alpha_k (f(x_k) - f(x^*)) + \alpha_k \langle g_k, x_k - x_{k+1} \rangle - \frac{\sigma}{2} \|x_k - x_{k+1}\|^2 \quad (48)$$

$$\leq \Delta_\psi(x^*, x_k) - \alpha_k (f(x_k) - f(x^*)) + \alpha_k \|g_k\|_* \|x_k - x_{k+1}\| - \frac{\sigma}{2} \|x_k - x_{k+1}\|^2 \quad (49)$$

$$\leq \Delta_\psi(x^*, x_k) - \alpha_k (f(x_k) - f(x^*)) + \frac{\alpha_k^2}{2\sigma} \|g_k\|_*^2 \quad (50)$$

Now compare with (40), we have successfully replaced $\|x_{k+1} - x^*\|_2^2$ with $\Delta_\psi(x^*, x_i)$. Again upper bound $\Delta_\psi(x^*, x_1)$ by R^2 and $\|g_k\|_*$ by G . Note the norm on g_k is the dual norm. To see the advantage of mirror descent, suppose C is the n dimensional simplex, and we use KL-divergence for which ψ is 1 strongly convex with respect to the ℓ_1 norm. The dual norm of the ℓ_1 norm is the ℓ_∞ norm. Then we can bound $\Delta_\psi(x^*, x_1)$ by using KL-divergence, and it is at most $\log n$. G can be upper bounded by M . So as for the value of RG , mirror descent is smaller than subgradient descent by an order of $O(\sqrt{\frac{n}{\log n}})$.

Acceleration 1: f is strongly convex. We say f is strongly convex with respect to another convex function ψ with modulus λ if

$$f(x) \geq f(y) + \langle g, x - y \rangle + \lambda \Delta_\psi(x, y) \quad \forall g \in \partial f(y). \quad (51)$$

Note we do not assume f is differentiable. Now in the step from (47) to (48), we can plug in the definition of strong convexity:

$$\Delta_\psi(x^*, x_{k+1}) = \dots + \alpha_k \langle g_k, x^* - x_k \rangle + \dots \quad (\text{copy of (47)}) \quad (52)$$

$$\leq \dots - \alpha_k (f(x_k) - f(x^*) + \lambda \Delta_\psi(x^*, x_k)) + \dots \quad (53)$$

$$\leq \dots \quad (54)$$

$$\leq (1 - \lambda \alpha_k) \Delta_\psi(x^*, x_k) - \alpha_k (f(x_k) - f(x^*)) + \frac{\alpha_k^2}{2\sigma} \|g_k\|_*^2 \quad (55)$$

Denote $\delta_k = \Delta_\psi(x^*, x_k)$. Set $\alpha_k = \frac{1}{\lambda k}$. Then

$$\delta_{k+1} \leq \frac{k-1}{k} \delta_k - \frac{1}{\lambda k} \epsilon_k + \frac{G^2}{2\sigma \lambda^2 k^2} \Rightarrow k \delta_{k+1} \leq (k-1) \delta_k - \frac{1}{\lambda} \epsilon_k + \frac{G^2}{2\sigma \lambda^2 k} \quad (56)$$

Now telescope (sum up both sides from $k = 1$ to T)

$$T \delta_{T+1} \leq -\frac{1}{\lambda} \sum_{k=1}^T \epsilon_k + \frac{G^2}{2\sigma \lambda^2} \sum_{k=1}^T \frac{1}{k} \Rightarrow \min_{i \in \{1, \dots, T\}} \epsilon_k \leq \frac{G^2}{2\sigma \lambda} \frac{1}{T} \sum_{k=1}^T \frac{1}{k} \leq \frac{G^2}{2\sigma \lambda} \frac{O(\log T)}{T}. \quad (57)$$

Acceleration 2: f has Lipschitz continuous gradient. If the gradient of f is Lipschitz continuous, there exists $L > 0$ such that

$$\|\nabla f(x) - \nabla f(y)\|_* \leq L \|x - y\|, \quad \forall x, y. \quad (58)$$

Sometimes we just directly say f is smooth. It is also known that this is equivalent to

$$f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{L}{2} \|x - y\|^2. \quad (59)$$

We bound the $\langle g_k, x^* - x_{k+1} \rangle$ term in (46) as follows

$$\langle g_k, x^* - x_{k+1} \rangle = \langle g_k, x^* - x_k \rangle + \langle g_k, x_k - x_{k+1} \rangle \quad (60)$$

$$\leq f(x^*) - f(x_k) + f(x_k) - f(x_{k+1}) + \frac{L}{2} \|x_k - x_{k+1}\|^2 \quad (61)$$

$$= f(x^*) - f(x_{k+1}) + \frac{L}{2} \|x_k - x_{k+1}\|^2. \quad (62)$$

Plug into (46), we get

$$\Delta_\psi(x^*, x_{k+1}) \leq \Delta_\psi(x^*, x_k) + \alpha_k \left(f(x^*) - f(x_{k+1}) + \frac{L}{2} \|x_k - x_{k+1}\|^2 \right) - \frac{\sigma}{2} \|x_k - x_{k+1}\|^2. \quad (63)$$

Set $\alpha_k = \frac{\sigma}{L}$, we get

$$\Delta_\psi(x^*, x_{k+1}) \leq \Delta_\psi(x^*, x_k) - \frac{\sigma}{L} (f(x_{k+1}) - f(x^*)). \quad (64)$$

Telescope we get

$$\min_{k \in \{2, \dots, T+1\}} f(x_k) - f(x^*) \leq \frac{L \Delta(x^*, x_1)}{\sigma T} \leq \frac{LR^2}{\sigma T}. \quad (65)$$

This gives $O(\frac{1}{T})$ convergence rate. But if we are smarter, like Nesterov, the rate can be improved to $O(\frac{1}{T^2})$. We will not go into the details but the algorithm and proof are again based on Lemma 2. This is often called accelerated proximal gradient method.

2.1 Composite Objective

Suppose the objective function is $h(x) = f(x) + r(x)$, where f is smooth and $r(x)$ is simple, like $\|x\|_1$. If we directly apply the above rates for optimizing h , we get $O(\frac{1}{\sqrt{T}})$ rate of convergence because h is not smooth.

It will be nice if we can enjoy the $O(\frac{1}{T})$ rate as in smooth optimization. Fortunately this is possible thanks to the simplicity of $r(x)$, and we only need to extend the proximal operator (33) as follows:

$$x_{k+1} = \operatorname{argmin}_{x \in C} \left\{ f(x_k) + \langle g_k, x - x_k \rangle + r(x) + \frac{1}{\alpha_k} \Delta_\psi(x, x_k) \right\} \quad (66)$$

$$= \operatorname{argmin}_{x \in C} \left\{ \alpha_k f(x_k) + \alpha_k \langle g_k, x - x_k \rangle + \alpha_k r(x) + \Delta_\psi(x, x_k) \right\}. \quad (67)$$

Algorithm 1: Protocol of online learning

- 1 The player initializes a model x_1 .
 - 2 **for** $k = 1, 2, \dots$ **do**
 - 3 The player proposes a model x_k .
 - 4 The rival picks a function f_k .
 - 5 The player suffers a loss $f_k(x_k)$.
 - 6 The player gets access to f_k and use it to update its model to x_{k+1} .
-

Here we use a first-order Taylor approximation of f around x_k , but keep $r(x)$ exact. Assuming this proximal operator can be computed efficiently, then we can show all the above rates carry over. We here only show the case of general f (not necessarily smooth or strongly convex), and leave the rest two cases as an exercise. In fact we can again achieve $O(\frac{1}{T^2})$ rate when f has Lipschitz continuous gradient.

Consider $\alpha_k(f(x_k) + \langle g_k, x - x_k \rangle + r(x))$ as the L in Lemma 2. Then

$$\alpha_k(f(x_k) + \langle g_k, x^* - x_k \rangle + r(x^*)) + \Delta_\psi(x^*, x_k) \quad (68)$$

$$\geq \alpha_k(f(x_k) + \langle g_k, x_{k+1} - x_k \rangle + r(x_{k+1})) + \Delta_\psi(x_{k+1}, x_k) + \Delta_\psi(x^*, x_{k+1}). \quad (69)$$

Following exactly the derivations from (46) to (50), we obtain

$$\Delta_\psi(x^*, x_{k+1}) \leq \Delta_\psi(x^*, x_k) + \alpha_k \langle g_k, x^* - x_{k+1} \rangle + \alpha_k(r(x^*) - r(x_{k+1})) - \Delta_\psi(x_{k+1}, x_k) \quad (70)$$

$$\leq \dots \quad (71)$$

$$\leq \Delta_\psi(x^*, x_k) - \alpha_k(f(x_k) + r(x_{k+1}) - f(x^*) - r(x^*)) + \frac{\alpha_k^2}{2\sigma} \|g_k\|_*^2. \quad (72)$$

This is almost the same as (50), except that we want to have $r(x_k)$ here, not $r(x_{k+1})$. Fortunately this is not a problem as long as we use a slightly different way of telescoping. Denote $\delta_k = \Delta_\psi(x^*, x_k)$ and then

$$f(x_k) + r(x_{k+1}) - f(x^*) - r(x^*) \leq \frac{1}{\alpha_k}(\delta_k - \delta_{k+1}) + \frac{\alpha_k}{2\sigma} \|g_k\|_*^2. \quad (73)$$

Summing up from $k = 1$ to T we obtain

$$r(x_{T+1}) - r(x_1) + \sum_{k=1}^T (h(x_k) - h(x^*)) \leq \frac{\delta_1}{\alpha_1} + \sum_{k=2}^T \delta_k \left(\frac{1}{\alpha_k} - \frac{1}{\alpha_{k-1}} \right) - \frac{\delta_{T+1}}{\alpha_T} + \frac{G^2}{2\sigma} \sum_{k=1}^T \alpha_k \quad (74)$$

$$\leq R^2 \left(\frac{1}{\alpha_1} + \sum_{k=2}^T \left(\frac{1}{\alpha_k} - \frac{1}{\alpha_{k-1}} \right) \right) + \frac{G^2}{2\sigma} \sum_{k=1}^T \alpha_k \quad (75)$$

$$= \frac{R^2}{\alpha_T} + \frac{G^2}{2\sigma} \sum_{k=1}^T \alpha_k. \quad (76)$$

Suppose we choose $x_1 = \operatorname{argmin}_x r(x)$, which ensures $r(x_{T+1}) - r(x_1) \geq 0$. Setting $\alpha_k = \frac{R}{G} \sqrt{\frac{\sigma}{k}}$, we get

$$\sum_{k=1}^T (h(x_k) - h(x^*)) \leq \frac{RG}{\sqrt{\sigma}} \left(\sqrt{T} + \frac{1}{2} \sum_{k=1}^T \frac{1}{\sqrt{k}} \right) = \frac{RG}{\sqrt{\sigma}} O(\sqrt{T}). \quad (77)$$

Therefore $\min_{k=1, \dots, T} \{h(x_k) - h(x^*)\}$ decays at the rate of $O(\frac{RG}{\sqrt{\sigma T}})$.

2.2 Online learning

The protocol of online learning is shown in Algorithm 1. The player's goal of online learning is to minimize the regret, the minimal possible loss $\sum_k f_k(x)$ over all possible x :

$$\text{Regret} = \sum_{k=1}^T f_k(x_k) - \min_x \sum_{k=1}^T f_k(x). \quad (78)$$

Note there is no assumption made on how the rival picks f_k , and it can adversarial. After obtaining f_k at iteration k , let us update the model into x_{k+1} by using the mirror descent rule on function f_k only:

$$x_{k+1} = \operatorname{argmin}_{x \in C} \left\{ f_k(x_k) + \langle g_k, x - x_k \rangle + \frac{1}{\alpha_k} \Delta_\psi(x, x_k) \right\}, \quad \text{where } g_k \in \partial f_k(x_k). \quad (79)$$

Then it is easy to derive the regret bound. Using f_k in step (50), we have

$$f_k(x_k) - f_k(x^*) \leq \frac{1}{\alpha_k} (\Delta_\psi(x^*, x_k) - \Delta_\psi(x^*, x_{k+1})) + \frac{\alpha_k}{2\sigma} \|g_k\|_*^2. \quad (80)$$

Summing up from $k = 1$ to n and using the same process as in (74) to (77), we get

$$\sum_{k=1}^T (f_k(x_k) - f_k(x^*)) \leq \frac{RG}{\sqrt{\sigma}} O(\sqrt{T}). \quad (81)$$

So the regret grows in the order of $O(\sqrt{T})$.

f is strongly convex. Exactly use (55) with f_k in place of f , and we can derive the $O(\log T)$ regret bound immediately.

f has Lipschitz continuous gradient. The result in (64) can NOT be extended to the online setting because if we replace f by f_k we will get $f_k(x_{k+1}) - f_k(x^*)$ on the right-hand side. Telescoping will not give a regret bound. In fact, it is known that in the online setting, having a Lipschitz continuous gradient itself cannot reduce the regret bound from $O(\sqrt{T})$ (as in nonsmooth objective) to $O(\log T)$.

Composite objective. In the online setting, both the player and the rival know $r(x)$, and the rival changes $f_k(x)$ at each iteration. The loss incurred at each iteration is $h_k(x_k) = f_k(x_k) + r(x_k)$. The update rule is

$$x_{k+1} = \operatorname{argmin}_{x \in C} \left\{ f_k(x_k) + \langle g_k, x - x_k \rangle + r(x) + \frac{1}{\alpha_k} \Delta_\psi(x, x_k) \right\}, \quad \text{where } g_k \in \partial f_k(x_k). \quad (82)$$

Note in this setting, (73) becomes

$$f_k(x_k) + r(x_{k+1}) - f_k(x^*) - r(x^*) \leq \frac{1}{\alpha_k} (\delta_k - \delta_{k+1}) + \frac{\alpha_k}{2\sigma} \|g_k\|_*^2. \quad (83)$$

Although we have $r(x_{k+1})$ here rather than $r(x_k)$, it is fine because r does not change through iterations. Choosing $x_1 = \operatorname{argmin}_x r(x)$ and telescoping in the same way as from (74) to (77), we immediately obtain

$$\sum_{k=1}^T (h_k(x_k) - h_k(x^*)) \leq \frac{G}{\sqrt{\sigma}} O(\sqrt{T}). \quad (84)$$

So the regret grows at $O(\sqrt{T})$.

When f_k are strongly convex, we can get $O(\log T)$ regret for the composite case. But as expected, having Lipschitz continuity of ∇f_k alone cannot reduce the regret from $O(\sqrt{T})$ to $O(\log T)$.

2.3 Stochastic optimization

Let us consider optimizing a function which takes a form of expectation

$$\min_x F(x) := \mathbb{E}_{\omega \sim p} [f(x; \omega)], \quad (85)$$

where p is a distribution of ω . This subsumes a lot of machine learning models. For example, the SVM objective is

$$F(x) = \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - c_i \langle a_i, x \rangle\} + \frac{\lambda}{2} \|x\|^2. \quad (86)$$

It can be interpreted as (85) where ω is uniformly distributed in $\{1, 2, \dots, m\}$ (i.e. $p(\omega = i) = \frac{1}{m}$), and

$$f(x; i) = \max\{0, 1 - c_i \langle a_i, x \rangle\} + \frac{\lambda}{2} \|x\|^2. \quad (87)$$

When m is large, it can be costly to calculate F and its subgradient. So a simple idea is to base the updates on a single randomly chosen data point. It can be considered as a special case of online learning in Algorithm 1, where the rival in step 4 now randomly picks f_k as $f(x; \omega_k)$ with ω_k being drawn independently from p . Ideally we hope that by using the mirror descent updates, x_k will gradually approach the minimizer

Algorithm 2: Protocol of online learning

- 1 The player initializes a model x_1 .
 - 2 **for** $k = 1, 2, \dots$ **do**
 - 3 The player proposes a model x_k .
 - 4 The rival randomly draws a ω_k from p , which defines a function $f_k(x) := f(x; \omega_k)$.
 - 5 The player suffers a loss $f_k(x_k)$.
 - 6 The player gets access to f_k and use it to update its model to x_{k+1} by, *e.g.*, mirror descent (79).
-

of $F(x)$. Intuitively this is quite reasonable, and by using f_k we can compute an unbiased estimate of $F(x_k)$ and a subgradient of $F(x_k)$ (because ω_k are sampled iid from p). This is a particular case of stochastic optimization, and we recap it in Algorithm 2.

In fact, the method is valid in a more general setting. For simplicity, let us just say the rival plays ω_k at iteration k . Then an online learning algorithm \mathcal{A} is simply a deterministic mapping from an ordered set $\{\omega_1, \dots, \omega_k\}$ to x_{k+1} . Denote as $\mathcal{A}(\omega_0)$ the initial model x_1 . Then the following theorem is the key for online to batch conversion.

Theorem 3 *Suppose an online learning algorithm \mathcal{A} has regret bound R_k after running Algorithm 1 for k iterations. Suppose $\omega_1, \dots, \omega_{T+1}$ are drawn iid from p . Define $\hat{x} = \mathcal{A}(\omega_{j+1}, \dots, \omega_T)$ where j is drawn uniformly random from $\{0, \dots, T\}$. Then*

$$\mathbb{E}[F(\hat{x})] - \min_x F(x) \leq \frac{R_{T+1}}{T+1}, \quad (88)$$

where the expectation is with respect to the randomness of $\omega_1, \dots, \omega_T$, and j .

Similarly we can have high probability bounds, which can be stated in the form like (not exactly true)

$$F(\hat{x}) - \min_x F(x) \leq \frac{R_{T+1}}{T+1} \log \frac{1}{\delta} \quad (89)$$

with probability $1 - \delta$, where the probability is with respect to the randomness of $\omega_1, \dots, \omega_T$, and j .

Proof of Theorem 3.

$$\mathbb{E}[F(\hat{x})] = \mathbb{E}_{j, \omega_1, \dots, \omega_{T+1}} [f(\hat{x}; \omega_{T+1})] = \mathbb{E}_{j, \omega_1, \dots, \omega_{T+1}} [f(\mathcal{A}(\omega_{j+1}, \dots, \omega_T); \omega_{T+1})] \quad (90)$$

$$= \mathbb{E}_{\omega_1, \dots, \omega_{T+1}} \left[\frac{1}{T+1} \sum_{j=0}^T f(\mathcal{A}(\omega_{j+1}, \dots, \omega_T); \omega_{T+1}) \right] \quad (\text{as } j \text{ is drawn uniformly random}) \quad (91)$$

$$= \frac{1}{T+1} \mathbb{E}_{\omega_1, \dots, \omega_{T+1}} \left[\sum_{j=0}^T f(\mathcal{A}(\omega_1, \dots, \omega_{T-j}); \omega_{T+1-j}) \right] \quad (\text{shift iteration index by iid of } \omega_i) \quad (92)$$

$$= \frac{1}{T+1} \mathbb{E}_{\omega_1, \dots, \omega_{T+1}} \left[\sum_{s=1}^{T+1} f(\mathcal{A}(\omega_1, \dots, \omega_{s-1}); \omega_s) \right] \quad (\text{change of variable } s = T - j + 1) \quad (93)$$

$$\leq \frac{1}{T+1} \mathbb{E}_{\omega_1, \dots, \omega_{T+1}} \left[\min_x \sum_{s=1}^{T+1} f(x; \omega_s) + R_{T+1} \right] \quad (\text{apply regret bound}) \quad (94)$$

$$\leq \min_x \mathbb{E}_{\omega} [f(x; \omega)] + \frac{R_{T+1}}{T+1} \quad (\text{expectation of min is smaller than min of expectation}) \quad (95)$$

$$= \min_x F(x) + \frac{R_{T+1}}{T+1}. \quad (96)$$