

Mirror Descent part 2: Introduction to the algorithm

Andersen Ang

ECS, Uni. Southampton, UK
andersen.ang@soton.ac.uk

Homepage angms.science

Version: July 19, 2024

First draft: Sept 6, 2020

Content

- ▶ Bregman divergence

$$B_\phi(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x}) - \phi(\mathbf{y}) - \langle \nabla \phi(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle,$$

- ▶ Bregman projection of \mathbf{y} onto the same set \mathcal{X} , under the function ϕ

$$P_{\mathcal{X}, \phi}(\mathbf{y}) = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} B_\phi(\mathbf{x}, \mathbf{y}).$$

- ▶ Mirror descent

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{\alpha} B_\phi(\mathbf{x}, \mathbf{x}_k).$$

Read these first

- ▶ Gradient lives in the dual space

https://angms.science/doc/CVX/Grad_is_dual.pdf

- ▶ Mirror Descent part 1: Why mirror descent? Why gradient descent is not enough?

https://angms.science/doc/CVX/MD_motivation.pdf

Bregman divergence

- ▶ A key ingredient of mirror descent is the Bregman divergence.
- ▶ Given a norm $\|\cdot\|$ (any norm), the Bregman divergence defined on a function ϕ is

$$B_\phi(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x}) - \phi(\mathbf{y}) - \langle \nabla \phi(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle,$$

where ϕ

- ▶ is differentiable with respect to \mathbf{x} , i.e. $\nabla_{\mathbf{x}} \phi(\mathbf{x})|_{\mathbf{x}=\mathbf{x}_0}$ exists for all \mathbf{x}_0
 - ▶ is λ -strongly convex w.r.t. the norm $\|\cdot\|$, i.e., $\phi(\mathbf{x}) - \frac{\lambda}{2} \|\mathbf{x}\|^2$ is convex.
- ▶ Based on the definition of strong-convexity,

$$B_\phi(\mathbf{x}, \mathbf{y}) \geq \frac{\lambda}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

- ▶ See [here](#) for more on B_ϕ and [here](#) for more on strongly convex function.

Some key points about Bregman divergence

$$B_{\phi}(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x}) - \phi(\mathbf{y}) - \langle \nabla \phi(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle,$$

- ▶ Bregman divergence is not commutative: $B(\mathbf{x}, \mathbf{y}) \neq B(\mathbf{y}, \mathbf{x})$.
Easy to see it from the definition.
- ▶ Bregman divergence is nonnegative.
This is due to the assumption that $\phi(\mathbf{x})$ is strongly-convex
- ▶ Bregman divergence is like a generalization of the squared-Euclidean distance.
- ▶ See [here](#) for more on Bregman divergence.

Bregman projection vs Euclidean projection

- ▶ Given a set \mathcal{X} , the Euclidean projection of a point \mathbf{y} onto this set is defined as

$$P_{\mathcal{X}}(\mathbf{y}) = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$$

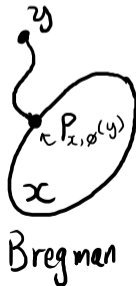
- ▶ The Bregman projection of \mathbf{y} onto the same set \mathcal{X} , under the function ϕ , is defined as

$$P_{\mathcal{X}, \phi}(\mathbf{y}) = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} B_{\phi}(\mathbf{x}, \mathbf{y}).$$

- ▶ Geometrically, Euclidean projection uses Euclidean geometry to give an orthogonal projection, while Bregman projection uses a curved geometry to give the projection.



Euclidean



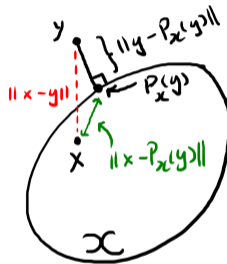
Bregman

Bregman projection vs Euclidean projection

- ▶ For Euclidean projection, for all $\mathbf{x} \in \mathcal{X}$, we have

$$\|\mathbf{x} - \mathbf{y}\|_2^2 \geq \|\mathbf{x} - P_{\mathcal{X}}(\mathbf{y})\|_2^2 + \|P_{\mathcal{X}}(\mathbf{y}) - \mathbf{y}\|_2^2.$$

Triangle inequality of Euclidean projection



- ▶ Similarly, for Bregman projection,

$$B_{\phi}(\mathbf{x}, \mathbf{y}) \geq B_{\phi}(\mathbf{x} - P_{\mathcal{X}}(\mathbf{y})) + B_{\phi}(P_{\mathcal{X}}(\mathbf{y}) - \mathbf{y}).$$

Re-look at the gradient descent algorithm

- The gradient descent update (with $\alpha > 0$ being the stepsize)

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k), \quad (1)$$

can be treated as the minimization of a local quadratic model of f at the point \mathbf{x}_k :

$$\begin{aligned} \mathbf{x}_{k+1} &= \operatorname{argmin}_{\mathbf{x}} Q(\mathbf{x}, \mathbf{x}_k) \\ &= \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\alpha} \|\mathbf{x} - \mathbf{x}_k\|_2^2. \end{aligned}$$

That is, take $\nabla Q(\mathbf{x}, \mathbf{x}_k) = \mathbf{0}$ gives (1).

- See [here](#) for more.

The mirror descent algorithm

- Mirror descent = gradient descent with the ℓ_2 -norm term $\frac{1}{2} \|\cdot\|_2^2$ replaced by the Bregman distance

$$\begin{aligned}\mathbf{x}_{k+1} &= \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{\alpha} B_\phi(\mathbf{x}, \mathbf{x}_k) \\ &= \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{\phi(\mathbf{x}) - \phi(\mathbf{x}_k) - \langle \nabla \phi(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle}{\alpha} \\ &= \operatorname{argmin}_{\mathbf{x}} \langle \nabla f(\mathbf{x}_k), \mathbf{x} \rangle + \frac{\phi(\mathbf{x}) - \langle \nabla \phi(\mathbf{x}_k), \mathbf{x} \rangle}{\alpha} \\ &= \operatorname{argmin}_{\mathbf{x}} \left\langle \nabla f(\mathbf{x}_k) - \frac{1}{\alpha} \nabla \phi(\mathbf{x}_k), \mathbf{x} \right\rangle + \frac{1}{\alpha} \phi(\mathbf{x}).\end{aligned}$$

Note: there is no an “explicit” formula of MD, because it all depends on the choice of ϕ .

- The whole picture:

1. You are given a function f to minimize
2. You pick a norm $\|\cdot\|$ and a function ϕ that is λ -strongly convex w.r.t. the chosen $\|\cdot\|$.
3. You minimize f using mirror descent.

You are free to choose any $\|\cdot\|$ and ϕ as you like, and when $\|\cdot\|$ and ϕ are “well-chosen”, mirror descent converges faster than gradient descent!

Operator theory of the mirror descent algorithm

- ▶ From the compact mirror descent update

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \underbrace{\left\langle \nabla f(\mathbf{x}_k) - \frac{1}{\alpha} \nabla \phi(\mathbf{x}_k), \mathbf{x} \right\rangle + \frac{1}{\alpha} \phi(\mathbf{x})}_{Q(\mathbf{x}, \mathbf{x}_k)}, \quad (*)$$

- ▶ Similar to what we did in slide 7, to find the explicit solution of (*), we take $\nabla Q(\mathbf{x}, \mathbf{x}_k)|_{\mathbf{x}=\mathbf{x}_{k+1}} = \mathbf{0}$, which gives

$$\begin{aligned} \nabla f(\mathbf{x}_k) - \frac{1}{\alpha} \nabla \phi(\mathbf{x}_k) + \frac{1}{\alpha} \nabla \phi(\mathbf{x}_{k+1}) &= \mathbf{0} \\ \iff \nabla \phi(\mathbf{x}_{k+1}) &= \nabla \phi(\mathbf{x}_k) - \alpha \nabla f(\mathbf{x}_k). \end{aligned}$$

- ▶ View $\nabla \phi$ as an operator, we have

$$\mathbf{x}_{k+1} = (\nabla \phi)^{-1} (\nabla \phi(\mathbf{x}_k) - \alpha \nabla f(\mathbf{x}_k)).$$

As a side note, $(\nabla \phi)^{-1} = \nabla(\phi^*)$, where ϕ^* is the conjugate of ϕ .

Why mirror descent is called mirror descent

$$\mathbf{x}_{k+1} = (\nabla\phi)^{-1}(\nabla\phi(\mathbf{x}_k) - \alpha\nabla f(\mathbf{x}_k)).$$

- ▶ $\nabla\phi$ and $(\nabla\phi)^{-1}$ are the mirror maps:

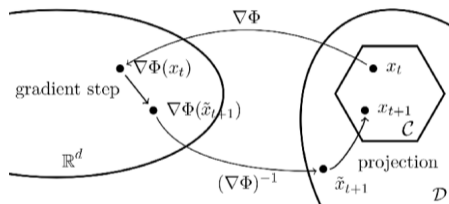


Figure: The “geometry” of mirror descent¹.

¹Figure from S. Bubeck, Convex optimization: Algorithms and complexity. Foundations and Trends in Machine Learning, 2015.

What about constrained problem

- ▶ The mirror descent update $\mathbf{x}_{k+1} = (\nabla\phi)^{-1}(\nabla\phi(\mathbf{x}_k) - \alpha\nabla f(\mathbf{x}_k))$ we discussed is for unconstrained problem: this update comes from replacing a term in the ordinary gradient descent step.
- ▶ For constrained problem $\min_{\mathbf{x} \in \mathcal{C}} f(\mathbf{x})$, the projected gradient descent adds a Euclidean projection on the gradient step

$$\mathbf{x}_{k+1} = P_{\mathcal{C}}(\mathbf{x}_k - \alpha\nabla f(\mathbf{x}_k)).$$

- ▶ For mirror descent, we can do the same: by adding a projection, we arrive at the projected mirror descent

$$\mathbf{x}_{k+1} = P_{\mathcal{C},\phi}\left(\left(\nabla\phi\right)^{-1}\left(\nabla\phi(\mathbf{x}_k) - \alpha\nabla f(\mathbf{x}_k)\right)\right).$$

This corresponds to the $\tilde{\mathbf{x}}$ in the figure in the previous slide.

- ▶ Note: the projection in projected mirror descent is the Bregman projection, not the Euclidean one.

The deep reasons why we consider mirror descent

- ▶ Motivation 1. Primal and dual vector spaces
 - ▶ In fact, for a function f , the point \mathbf{x} and the gradient $\nabla f(\mathbf{x})$ live in different space.
 - ▶ \mathbf{x} live in the primal space
 - ▶ $\nabla f(\mathbf{x})$ live in the dual space
 - ▶ For optimization in \mathbb{R}^n , because \mathbb{R}^n is self-dual, so we take for granted that thinking \mathbf{x} and $\nabla f(\mathbf{x})$ live in the same space, which in fact technically they live in different space.
 - ▶ Gradient descent does not take into account of the primal and dual vector spaces, but mirror descent does \implies why mirror descent is important.

- ▶ Motivation 2. Mirror descent using Bregman divergence can handle the case that f is L -smooth with respect to non-Euclidean norm.
 - ▶ There are some functions f that the gradient ∇f is L -Lipschitz in, say ℓ_∞ norm.

Last page - summary

- ▶ Bregman divergence
- ▶ Derivation of the mirror descent: replacing the ℓ_2 norm in gradient descent by Bregman divergence.
- ▶ Different expressions of the mirror descent update

Next

- ▶ Example of Mirror Descent algorithm
- ▶ The convergence theory of mirror descent on convex optimization problem

End of document