# Lecture 11     Mirror descent

- Projected gradient descent

- Mirror descent

- MWU as a special case of Mirror descent

---

## Projected gradient descent

We want to solve the following optimization problem:

$$\min_{x \in C} f(x)$$

convex function

convex set

### GD

$$x_{t+1} \leftarrow x_t - \eta_t \nabla f(x_t)$$

Problem: We may end outside $C$ !!!

$\rightsquigarrow$

### Projected GD

$$x_{t+1/2} \leftarrow x_t - \eta_t \nabla f(x_t)$$

$$x_{t+1} = \Pi_C(x_{t+1/2})$$

projection in $C$ operator, using as a distance function $\| \cdot \|_2^2 \Rightarrow \Pi_C(x) = \arg\min_{y \in C} \|x-y\|_2^2$
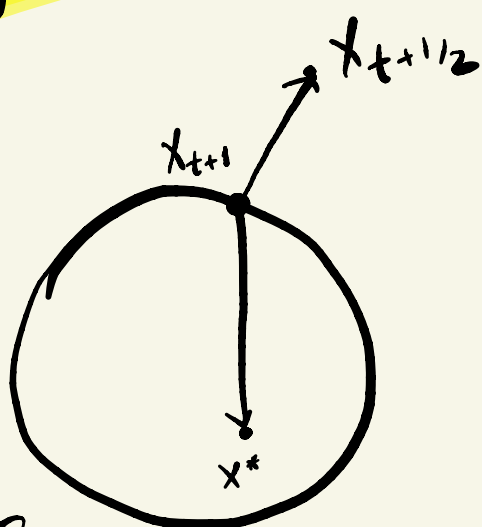
## Question:

Is it always good to project?

(in other words, do we get closer to the optimal point $x^*$ by projecting in $C$)

## Answer:

yes!!!

$$\left(x_{t+1/2} - x_{t+1}\right)\left(x^* - x_{t+1}\right) \le 0$$



$$\left(\begin{array}{c} \left(x - \Pi(x)\right)\left(y - \Pi(x)\right) \le 0 \\ \forall y \in C \quad, \quad \forall x \end{array}\right)$$

(remember that $C$ is convex)

$$\Rightarrow \quad \|x_{t+1/2} - x^*\|^2 \ge \|x_{t+1/2} - x_{t+1}\|^2 + \|x^* - x_{t+1}\|^2$$

$$\Rightarrow \quad \|x_{t+1/2} - x^*\|^2 \ge \|x_{t+1} - x^*\|^2$$

Important observation: (in order to derive later the Projected GD update on Mirror descent)

rule
$$x_{t+1/2} \leftarrow x_t - \eta_t \nabla f(x_t)$$
$$x_{t+1} \leftarrow \Pi_C\left(x_{t+1/2}\right) = \arg\min_{x \in C} \|x - x_{t+1/2}\|_2^2$$

is equivalent to a "regularized" updating rule

$$x_{t+1} = \arg\min_{x \in C} \left\{ f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{2\eta_t} \|x - x_t\|^2 \right\}$$

## Proof

$$\operatorname*{argmin}_{x \in C} \left\{ f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{2n_t} \|x - x_t\|_2^2 \right\} =$$

$$= \operatorname*{argmin}_{x \in C} \left\{ \langle \nabla f(x_t), x \rangle + \frac{1}{2n_t} \|x - x_t\|_2^2 \right\} =$$

$$= \operatorname*{argmin}_{x \in C} \left\{ n_t \langle \nabla f(x_t), x \rangle + \frac{1}{2} \|x - x_t\|_2^2 \right\}$$

and

$$\operatorname*{argmin}_{x \in C} \| x - x_{t+1/2} \|^2 = \operatorname*{argmin}_{x \in C} \| x - x_t + n_t \nabla f(x_t) \|^2 =$$

$$= \operatorname*{argmin}_{x \in C} \left\{ 2 \cdot \left( \frac{1}{2} \|x - x_t\|_2^2 + \frac{1}{2} \| n_t \nabla f(x_t) \|_2^2 + \langle x - x_t, n_t \nabla f(x_t) \rangle \right) \right\}$$

$$= \operatorname*{argmin}_{x \in C} \left\{ n_t \langle \nabla f(x_t), x \rangle + \frac{1}{2} \|x - x_t\|_2^2 \right\}$$

We proceed by analysing the convergence speed of Projected Gradient Descent when applied to Lipschitz, convex functions

## Theorem

$$\min_{x \in C} f(x)$$

convex

convex set

$L$-Lipschitz $\Rightarrow |f(x) - f(y)| \leq L \cdot \|x-y\|_2$

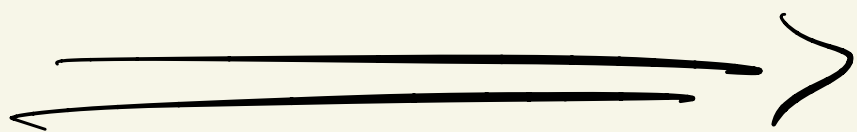$\Rightarrow \|\nabla f(x)\|_2 \leq L$

the Projected G.D. with stepsize

$$\eta = \frac{\|x_0 - x^*\|}{L \cdot \sqrt{T}} \quad \text{satisfies:}$$

$$f\left(\frac{1}{T} \sum_{t=1}^{T} x_t\right) - f(x^*) \leq \frac{L\|x_0 - x^*\|}{\sqrt{T}}$$

and

$$\min_{t \in [T]} f(x_t) - f(x^*) \leq \frac{L \cdot \|x_0 - x^*\|}{\sqrt{T}}$$

proof

# proof

$$\|x_{t+1} - x^*\|^2 \leq \|x_{t+1/2} - x^*\|^2 \quad \text{(Projection)}$$

$$= \|x_t - \eta \nabla f(x_t) - x^*\|^2 \quad \left(\begin{array}{c}\text{definition} \\ \text{of } x_{t+1/2}\end{array}\right)$$

$$= \|x_t - x^*\|^2 + \eta^2 \|\nabla f(x_t)\|^2 - 2\eta \langle \nabla f(x_t), x_t - x^* \rangle \leq$$

convexity of $f$

$f(x^*) \geq f(x_t) + \langle \nabla f(x_t), x^* - x_t \rangle$

$$\leq \|x_t - x^*\|^2 + \eta^2 \|\nabla f(x_t)\|^2 + 2\eta \left(f(x^*) - f(x_t)\right)$$

$\Big\downarrow$ adding for $t = 0$ to $T$

$$2\eta \sum_{t=0}^{T} \left(f(x_t) - f(x^*)\right) + \|x_{T+1} - x^*\|^2 \leq \|x_0 - x^*\|^2 + \eta^2 \sum_{t=0}^{T} \|\nabla f(x_t)\|^2$$

$$\implies \sum_{t=0}^{T} \left(f(x_t) - f(x^*)\right) \leq \frac{\|x_0 - x^*\|^2}{2\eta} + \eta \cdot \frac{\sum_{t=0}^{T} \|\nabla f(x_t)\|^2}{2}$$

$\|x_{T+1} - x^*\|^2 \geq 0$

$\|\nabla f(x_t)\|^2 \leq \alpha^2$

$$\implies \frac{1}{T} \sum f(x_t) - f(x^*) \leq \frac{\|x_0 - x^*\|^2}{2\eta T} + \frac{\eta \cdot \alpha^2}{2}$$

by noting that $\quad \min_{t \in [T]} f(x_t) \leq \frac{1}{T} \sum_{t=0}^{T} f(x_t)$

$$f\left(\frac{1}{T} \sum_{t=0}^{T} x_t\right) \leq \frac{1}{T} \sum_{t=0}^{T} f(x_t) \quad \text{(Jensen)}$$

and setting $\eta = \|x_0 - x^*\|_2 / \alpha \cdot \sqrt{T}$ we get the theorem 🔲

# Observations

(1) Thus after $T$ iterations the error is

$$\frac{\|x_0 - x^*\|_2 \cdot d}{\sqrt{T}}. \text{ To get an error}$$

of $\varepsilon$ we need $T = O\left(\frac{\|x_0 - x^*\|_2^2 \cdot d^2}{\varepsilon^2}\right)$

iterations.

(2) we can also use a variable step

size $n_t = \frac{\|x_0 - x^*\|_2}{d \sqrt{t+1}}$ and get that

$$\min_{t \in [T]} f(x_t) - f(x^*) \leq \frac{\|x_0 - x^*\|_2^2 + d^2 \cdot \sum n_t^2}{\sum n_t} =$$

$$= O\left(\frac{\|x_0 - x^*\|^2 \cdot d \cdot \log T}{\sqrt{T}}\right)$$

(3) We still get that

$$\frac{1}{T} \sum_{t=0}^{T} \left(f_t(x_t) - f_t(x^*)\right) \leq \frac{\|x_0 - x^*\|_2 \cdot d}{\sqrt{T}}$$

where the function $f$ changes every

time (online setting)

What happens if we do not have "very" useful information w.r.t the $\ell_2$-norm of the gradient?

For example we know that $\|\nabla f(x)\|_\infty \leq 1$ but only that $\|\nabla f(x)\|_2 \leq \sqrt{n}$.

# Mirror descent

Overview of what we will do:

(1) We will define Bregman divergence which "should be viewed" as a distance function

(2) We will redefine the projected G.D update rule to use the Bregman divergence as the projection function (instead of $\|\cdot\|_2$)

(3) choosing appropriate Bregman divergen as we get bounds w.r.t other norms than the $\ell_2$-norm

# Bregman Divergence

Let $\psi: C \to \mathbb{R}$ be a strictly convex function, continuously differentiable, over a closed convex set $C$. Then we define the Bregman divergence as

$$D_\varphi(x,y) := \psi(x) - \psi(y) - \langle \nabla \psi(y), x-y \rangle$$

## observations and examples

① $D_\varphi(x,y)$ measures how good is the first order approximation of $\psi(x)$ at $y$.

② $D_\varphi(x,y) \geq 0 \quad \forall \, x,y \in C$ ($\varphi$ is strictly convex)

③ $D_\varphi(x,y) = 0$ iff $x=y$ ($\varphi$ is strictly convex)

④ $\nabla_x D_\varphi(x,y) = \nabla \varphi(x) - \nabla \psi(y)$

⑤ Generalized law of cosines
  (usually this property is called
   Generalized pythagorean theorem)

$$D_\psi(x,y) + D_\varphi(y,z) =$$

$$= \psi(x) - \varphi(y) - \langle \nabla\psi(y), x-y \rangle$$

$$+ \psi(y) - \psi(z) - \langle \nabla\varphi(z), y-z \rangle =$$

$$= \left(\psi(x) - \psi(z) - \langle \nabla\varphi(z), x-z \rangle\right) + D_\varphi(x,z)$$

$$+ \langle \nabla\varphi(z), x-z \rangle - \langle \nabla\psi(y), x-y \rangle - \langle \nabla\varphi(z), y-z \rangle =$$

$$= D_\varphi(x,z) + \langle \nabla\psi(z) - \nabla\varphi(y), x-y \rangle$$

"angle" between

z-y and x-y (try $\psi(x) = \frac{1}{2}\|x\|^2$ to re-derive the law of cosines)

⑥ Let $\psi(x) = \frac{1}{2}\|x\|^2 \Rightarrow D_\psi(x,y) =$

$$= \frac{1}{2}\|x\|^2 - \frac{1}{2}\|y\|^2 - \langle y, x-y \rangle = \frac{1}{2}\|x\|^2 - \frac{1}{2}\|y\|^2 + \|y\|^2$$
$$- \langle x, y \rangle =$$

$$= \frac{1}{2}\|x\|^2 + \frac{1}{2}\|y\|^2 - \langle x, y \rangle = \frac{1}{2}\|x-y\|^2$$

(7) Let $\varphi(x) = \sum_{i=1}^{n} x_i \log x_i$, $x \in \Delta_1^n$

simplex

$D_\varphi(x,y) = \sum x_i \log x_i - \sum y_i \log y_i$

$$- \langle 1 + \log y, x - y \rangle =$$

$$= \sum x_i \log x_i - \sum y_i \log y_i - \sum (x_i - y_i)^{\nearrow 0 \text{(prob. distri)}}$$

$$- \sum (x_i - y_i) \log y_i =$$

$$= \sum x_i \log x_i - \sum x_i \log y_i =$$

$$= \sum x_i \log x_i/y_i = \text{Kd}(x \| y)$$ 🙂

## Projections with Bregman divergence

Let $\Pi(x) = \underset{y \in C}{\arg\min} \; D_\varphi(y, x)$ be the

projection operation using Bregman

divergence of $\varphi$.

- $\Pi(x)$ is uniquely determined because

$$g(y) = D_\psi(y,x) = \psi(y) - \psi(x) - \langle \nabla \psi(x), y-x \rangle$$

is strictly convex, and $C$ is a closed convex set $\Rightarrow$ unique minimizer

- $\langle \nabla \psi(x) - \nabla \psi(\Pi(x)), y - \Pi(x) \rangle \leq 0 \quad \forall x, \forall y \in C$

proof

$\Pi(x)$ is a minimizer of $D_\psi(\cdot, x)$ in $C \Rightarrow$

$\Rightarrow \langle \nabla D_\psi(\Pi(x), x), y - \Pi(x) \rangle \geq 0 \quad \forall y \in C$

$\Rightarrow \langle \nabla \psi(\Pi(x)) - \nabla \psi(x), y - \Pi(x) \rangle \geq 0 \quad \forall y \in C$

- $D_\psi(y, \Pi(x)) + D_\psi(\Pi(x), x) \leq D_\psi(y, x)$

  $\forall x, \forall y \in C$

$\Rightarrow D_\psi(x^*, x) \geq D_\psi(x^*, \Pi(x)), \quad x^* \in C$

In order to construct Mirror descent we will slightly change the update rule of projected G.D.

P.G.D.

$$X_{t+1} = \underset{x \in C}{\text{argmin}} \left\{ f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{2\eta_t} \|x - x_t\|_2^2 \right\}$$

MD

$$X_{b+1} = \underset{x \in C}{\text{argmin}} \left\{ f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{\eta_t} D_\varphi(x, x_t) \right\}$$
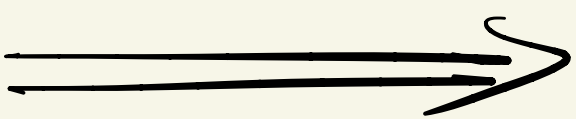
which is equivalent to

$$X_{t+\frac{1}{2}} = \text{argmin} \left\{ f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{\eta_t} D_\varphi(x, x_t) \right\}$$

ⓧ → uncoustrained

$$X_{t+1} = \underset{x \in C}{\text{argmin}} \, D_\varphi(x, x_{t+1/2})$$

→ projection step

proof

⟹

- $X_{t+1} = \underset{x \in C}{\arg\min} \left\{ f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{n_t} D_\psi(x, x_t) \right\} =$

$= \underset{x \in C}{\arg\min} \left\{ n_t \langle \nabla f(x_t), x \rangle + \psi(x) - \langle \nabla \psi(x_t), x \rangle \right\}$

$= \underset{x \in C}{\arg\min} \left\{ \psi(x) - \langle \nabla \psi(x_t) - n_t \nabla f(x_t), x \rangle \right\}$

- $X_{t+\frac{1}{2}} = \underset{x}{\arg\min} \left\{ f(x_t) + \langle \nabla f(x_t), x - x_t \rangle + \frac{1}{n_t} D_\psi(x, x_t) \right\}$

$X_{t+1} = \underset{x \in C}{\arg\min} D_\psi(x, x_{t+1/2})$

the first step is unconstrained, therefore
the gradient should be equal to $0 \Rightarrow$

$\Rightarrow \nabla f(x_t) + \frac{1}{n_t} \left( \nabla \psi(x) - \nabla \psi(x_t) \right) = 0 \Rightarrow$

$\Rightarrow \nabla \psi(x) = \nabla \psi(x_t) - n_t \nabla f(x_t)$

$\Rightarrow \nabla \psi(x_{t+1/2}) = \nabla \psi(x_t) - n_t \nabla f(x_t)$

$\Rightarrow X_{t+1/2} = (\nabla \psi)^{-1} \left( \nabla \psi(x_t) - n_t \nabla f(x_t) \right)$

$X_{t+1} = \underset{x \in C}{\arg\min} D_\psi(x, x_{t+1/2}) = \quad \Longrightarrow$

$$= \underset{x \in C}{\arg\min} \left\{ \psi(x) - \psi(x_{t+1/2}) - \langle \nabla \psi(x_{t+\frac{1}{2}}), x - x_{t+\frac{1}{2}} \rangle \right\}$$

$$= \underset{x \in C}{\arg\min} \left\{ \psi(x) - \langle \nabla \psi(x_{t+1/2}), x \rangle \right\} =$$

$$= \underset{x \in C}{\arg\min} \left\{ \psi(x) - \langle \nabla \psi(x_t) - n_t \nabla f(x_t), x \rangle \right\}$$

## Mirror descent

$$x_{t+1/2} = (\nabla \psi)^{-1} \left( \nabla \psi(x_t) - n_t \nabla f(x_t) \right)$$

$$x_{t+1} = \underset{x \in C}{\arg\min} D_\psi(x, x_{t+1/2})$$

## Theorem (convergence speed)

strictly convex, etc...

$$D_\psi(x, y) \geq a \|x - y\|^2 \quad \longleftarrow \text{arbitrary norm}$$

f is $d$-Lipschitz w.r.t $\|\cdot\|_\#$  then MD with

stepsize $n = \frac{\sqrt{a \cdot D_\psi(x^*, x_0)}}{d\sqrt{T}}$  satisfies

$$f\left(\frac{1}{T} \Sigma x_t\right) - f(x^*) \leq 2 \frac{\sqrt{D_\psi(x^*, x_0) \cdot d}}{\sqrt{aT}}$$

$$\underset{t \in [T]}{\min} f(x_t) - f(x^*) \leq \quad ''$$

## proof

As in the P.GD algorithm where we
proved that

$$\|x_{t+1}-x^*\|^2 \leq \|x_t-x^*\|^2 + 2n\left(f(x^*)-f(x_t)\right) + n^2\|\nabla f(x_t)\|_2^2$$

here it suffice to prove that

$$D_\varphi(x^*, x_{t+1}) \leq D_\varphi(x^*, x_t) + n\left(f(x^*)-f(x_t)\right) + \frac{n^2}{a}\|\nabla f(x_t)\|_*^2$$

( then by summing for $t=0,...,T$ we get a
telescopic sum etc etc )

we start by taking the generalized
law of cosines:

$$D_\varphi(x^*, x_{t+1}) + D_\varphi(x_{t+1}, x_t) = D_\varphi(x^*, x_t) + \langle \nabla\varphi(x_t) - \nabla\varphi(x_{t+1}), x^* - x_{t+1}\rangle$$

moreover since $x_{t+1} = \Pi_C(x_{t+1/2})$ we have

$$\langle \nabla_{x_{t+1}} D_\varphi(x_{t+1}, x_{t+1/2}), y - x_{t+1}\rangle \geq 0 \quad (\text{because } x_{t+1} \text{ is a minimizer})$$

$$\langle \nabla\varphi(x_{t+1}) - \nabla\varphi(x_{t+1/2}), x_{t+1} - y\rangle \leq 0 \quad \forall y \in C$$

$$\Rightarrow \langle \nabla\varphi(x_{t+1}), x_{t+1} - y\rangle \leq \langle \nabla\varphi(x_{t+1/2}), x_{t+1} - y\rangle \quad \forall y \in C$$

by setting $y = x^*$ and using the law of cosines
we get $\Rightarrow$

$$D_\varphi(x^*, x_{t+1}) + D_\varphi(x_{t+1}, x_t) = D_\varphi(x^*, x_t) + \langle \nabla \psi(x_t) - \nabla \psi(x_{t+1}), x^* - x_{t+1} \rangle$$

$$\Downarrow$$

$$D_\varphi(x^*, x_{t+1}) + D_\varphi(x_{t+1}, x_t) \leq D_\varphi(x^*, x_t) + \langle \nabla \psi(x_t) - \nabla \psi(x_{t+1/2}), x^* - x_{t+1} \rangle$$

$$\Downarrow \quad \boxed{\nabla \varphi(x_{t+1/2}) = \nabla \psi(x_t) - \eta_t \nabla f(x_t)}$$

$$D_\varphi(x^*, x_{t+1}) + D_\varphi(x_{t+1}, x_t) \leq D_\varphi(x^*, x_t) + \eta_t \langle \nabla f(x_t), x^* - x_{t+1} \rangle$$

$$\Downarrow \quad \text{rearranging}$$

$$D_\varphi(x^*, x_{t+1}) \leq D_\varphi(x^*, x_t) + \boxed{\eta_t \langle \nabla f(x_t), x^* - x_t \rangle} - \eta_t \langle \nabla f(x_t), x_{t+1} - x_t \rangle$$

$$\Downarrow \quad \text{first order condition} \quad - D_\psi(x_{t+1}, x_t)$$

$$D_\varphi(x^*, x_{t+1}) \leq D_\varphi(x^*, x_t) + \eta_t \big( f(x^*) - f(x_t) \big) - \eta_t \langle \nabla f(x_t), x_{t+1} - x_t \rangle$$
$$- D_\varphi(x_{t+1}, x_t)$$

$$\boxed{D_\varphi(x, y) \geq \alpha \cdot \|x - y\|^2}$$

$$\Downarrow \quad \boxed{\eta_t \langle \nabla f(x_t), x_{t+1} - x_t \rangle \leq}$$
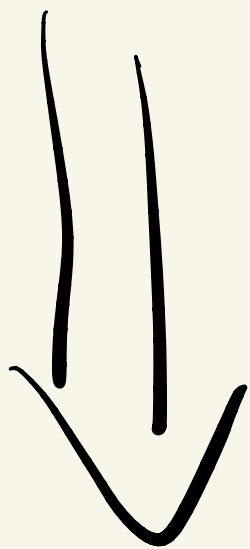$$\boxed{\leq \eta_t \cdot \|\nabla f(x_t)\|_* \|x_{t+1} - x_t\|}$$

$$D_\varphi(x^*, x_{t+1}) \leq D_\varphi(x^*, x_t) + \eta_t \big( f(x^*) - f(x_t) \big) + \eta_t \|\nabla f(x_t)\|_* \|x_{t+1} - x_t\|$$
$$- \alpha \|x_{t+1} - x_t\|^2$$

$$D_\varphi(x^*, x_{t+1}) \leq D_\varphi(x^*|x_t) + n_t(f(x^*) - f(x_t)) + n_t\|\nabla f(x_t)\|_* \|x_{t+1} - x_t\|$$
$$- a\|x_{t+1} - x_t\|^2$$

$$n_t\|\nabla f(x_t)\|_* \|x_{t+1} - x_t\|$$
$$-a\|x_{t+1} - x_t\|^2$$
$$\leq \frac{n_t^2\|\nabla f(x_t)\|_*^2}{a}$$

$$D_\varphi(x^*, x_{t+1}) \leq D_\varphi(x^*|x_t) + n_t(f(x^*) - f(x_t)) + \frac{n_t^2\|\nabla f(x_t)\|_*^2}{a}$$

to make everything work
remember we need to select
$\varphi$ s.t $D_\varphi(x,y) \geq a\|x-y\|^2$

Thus, to achieve $\varepsilon$-optimality

we need $2\dfrac{\sqrt{D_\varphi(x^*|x_0) \cdot d}}{\sqrt{aT}} \leq \varepsilon \implies$

$$\implies T = O\left(\frac{D_\varphi(x^*, x_0) d^2}{a \varepsilon^2}\right)$$

$f$ is $d$- Lipschitz w.r.t
norm $\|\cdot\|$ which may not be
the $\ell_2$-norm !!!