# Optimization for Machine Learning
# CSCI-599

## Lecture 3: Constrained and Non-Smooth Gradient Descent

**Sai Praneeth Karimireddy**

USC – `https://spkreddy.org/optmlspring2025.html`

January 29, 2025

**Q1**

f is non-smooth $\Rightarrow$ grad is undefined?

$\hookrightarrow$ subgradient : always exists everywhere
iff f is cvx

$\hookrightarrow$ if bounded subgrad $\|g(x)\|_2 \leq B \Rightarrow$
subgradient $\dfrac{RB}{\sqrt{T}}$
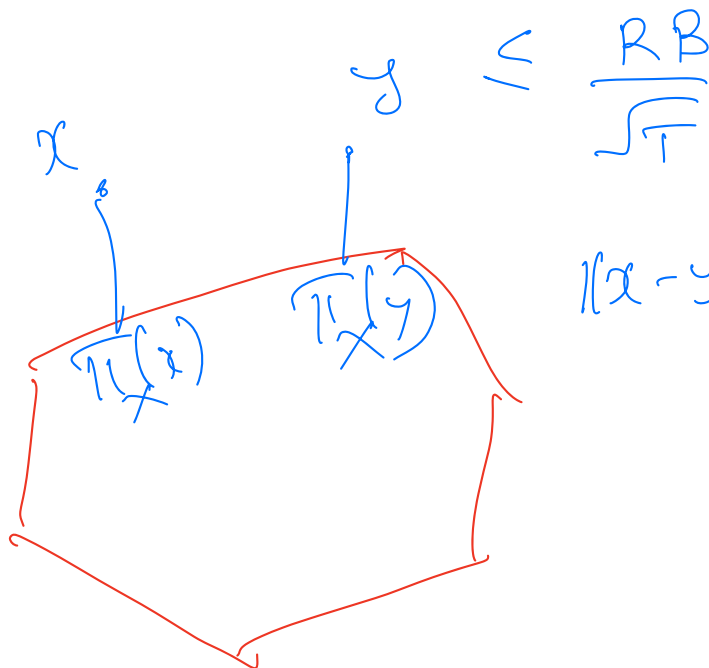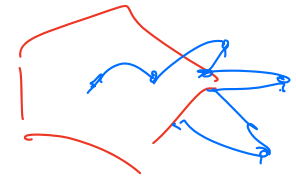
**Q2**

$\|x_0 - x^*\|_2 \leq R$ without knowing $x^*$?

constrained optimization

$\min\limits_{x \in X} f(x) \to$ cvx

$\to$ bounded $\Rightarrow \|x_0 - x^*\|_2 \leq \text{diam}(X)$

$\max\limits_{x \in X} \|g(x)\|_2 \leq B$

$x_{t+1} = \prod\limits_{X}\left(x_t - r\, g(x_t)\right)$



$y \leq \dfrac{RB}{\sqrt{T}}$

$x$

$\prod_X(y)$

$\prod_X(x)$

$\|x - y\|_2 \geq \|\prod_X(x) - \prod_X(y)\|_2$

$$x_{t+1} = \Pi_X \left( \underbrace{x_t - r D f(x_t)}_{y_{t+1}} \right) \qquad f \text{ is cvxt}$$

$$L\text{-smooth}$$

$$\Rightarrow \quad f(x_t) - f(x^*) \leq \frac{L \|x_0 - x^*\|_2^2}{T}$$

## Sufficient decrease

In Unconstrained case,



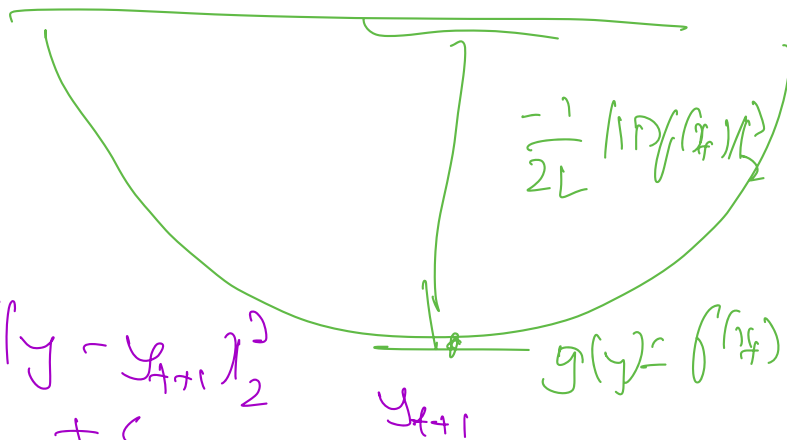$$\left( f(y_{t+1}) - f(x_t) \leq \frac{-1}{2L} \| Df(x_t) \|_2^2 \right)$$

In Constrained,

$$f(x_{t+1}) - f(x_t) \leq \frac{-1}{2L} \| Df(x_t) \|_2^2 + \frac{L}{2} \| y_{t+1} - x_{t+1} \|_2^2$$

Proof

Intuition



$$\frac{-1}{2L} \| Df(x_t) \|_2^2$$

$$g(y) = \frac{1}{2} \| y - y_{t+1} \|_2^2 + C$$

$$g(y) = f(x_t) + Df(x_t)^T (y - x_t) + \frac{L}{2} \| y - x_t \|_2^2$$

Formal

$$f(x_{t+1}) \leq f(x_t) + \frac{Df(x_t)^T(\boxed{x_{t+c}} - x_t)}{} + \frac{L}{2}\|x_{t+c} - x_t\|_2^2$$

$$= f(x_t) + \frac{L}{2}\left\| x_{t+c} - x_t + \frac{1}{L}Df(x_t)\right\|^2 - \frac{1}{2L}\| Df(x_t)\|_2^2$$

$$= f(x_t) + \frac{L}{2}\| x_{t+c} - y_{t+1}\|_2^2 - \frac{1}{2L}\|Df(x_t)\|_2^2$$

$$\Rightarrow \|Df(x_t)\|_2^2 \leq 2L\left(f(x_t) - f(x_{t+c})\right) + L^2\|x_{t+c} - y_{t+1}\|_2^2$$

$$\boxed{\|y_{t+c} - x^*\|_2^2 = \|x_t - x^*\|^2 - 2\gamma\langle Df(x_t), x_t - x^*\rangle + \gamma^2\|Df(x_t)\|_2^2}$$

$$\geq \|x_{t+1} - x^*\|_2^2 + \|x_{t+c} - y_{t+1}\|^2$$

$$\Rightarrow \|x_{t+c} - x^*\|^2 \leq \|x_t - x^*\|^2 - 2\gamma\langle Df(x_t), x_t - x^*\rangle + \gamma^2\|Df(x_t)\|_2^2$$

$$- \|x_{t+1} - y_{t+1}\|_2^2$$

Substituting our bound on grad,

$$\|x_{t+1} - x^*\|^2 \leq \|x_t - x^*\|^2 - 2\gamma \langle \nabla f(x_t), x_t - x^* \rangle$$

$$+ \gamma^2 \left( 2L \left( f(x_t) - f(x_{t+1}) \right) \right.$$
$$\left. + L^2 \|x_{t+1} - y_{t+1}\|_2^2 \right)$$

$$- \|x_{t+1} - y_{t+1}\|_2^2$$

$$\gamma \leq \frac{1}{L} \implies \frac{L}{2\gamma}$$

$$\sum_{t=0}^{T} \underbrace{\langle \nabla f(x_t), x_t - x^* \rangle}_{\geq f(x_t) - f(x^*)} \frac{L}{2\gamma} \sum_{t=0}^{T} \|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2$$

$$\geq f(x_t) - f(x^*) + \sum_{t=0}^{T} \frac{2}{L 2\gamma} \left( f(x_t) - f(x_{t+1}) \right)$$

$$\implies f(x_t) - f(x^*) \leq \frac{2\|x_0 - x^*\|^2}{T}$$

If strgly cvx,

$$\underbrace{\langle \nabla f(x_t), x_t - x^* \rangle} \leq \frac{L}{2}\left(\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2\right)$$

$$\geq f(x_t) - f(x^*) \qquad + 0\left(f(x_t) - f(x_{t+1})\right)$$

$$+ \frac{\mu}{2}\|x_t - x^*\|^2$$

$$\Rightarrow \quad \frac{\mu}{2}\|x_t - x^*\|^2 \leq \frac{L}{2}\|x_t - x^*\|^2 - \frac{L}{2}\|x_{t+1} - x^*\|^2$$

$$\Rightarrow \quad \frac{\not{A}}{\not{2}}\|x_{t+1} - x^*\|^2 \leq \frac{\not{L}}{\not{2}}\left(1 - \frac{\mu}{L}\right)\|x_t - x^*\|_2^2$$

$$\Rightarrow \quad \|x_{t+1} - x^*\|^2 \leq \left(1 - \frac{\mu}{L}\right)^t \|x_0 - x^*\|_2^2$$

$$Q1 = \frac{1}{2L}\|Df(x)\|^2$$

$$Q3 = \frac{1}{2\mu}\|Df\|^2$$

$$\frac{\mu}{L}$$

# Can we go even faster?

So far: Error decreases with $1/\sqrt{T}$, or $1/T$...

Could it decrease exponentially in $T$?

# Can we go even faster?

$\nabla f(x) = 2x$

$\nabla^2 f(x) = 2$

- On $f(x) := x^2$, Stepsize $\gamma := \frac{1}{2}$  ($f$ is $L=2$ - smooth)

$$x_{t+1} = x_t - \frac{1}{2} \nabla f(x_t) =$$

$x_0 := 1$

Theory

$x_{t+1} = x_t - \frac{1}{2}(2x_t) = 0$

$\frac{\textcircled{L} \|x_0 - x^*\|^2}{2T}$

$= \frac{2 \cdot 1}{2T} = \frac{1}{T}$

# Can we go even faster?

▶ On $f(x) := x^2$: Stepsize $\gamma := \frac{1}{2}$  ($f$ is $L{=}2$ - smooth)

$$x_{t+1} = x_t - \frac{1}{2}\nabla f(x_t) = x_t - x_t = 0,$$

  ▶ converged in one step!

# Can we go even faster?

▶ On $f(x) := x^2$: Stepsize $\gamma := \frac{1}{2}$ ($f$ is $L=2$ - smooth)

$$x_{t+1} = x_t - \frac{1}{2}\nabla f(x_t) = x_t - x_t = 0,$$

  ▶ converged in one step!

▶ Same $f(x) := x^2$: Stepsize $\gamma := \frac{1}{4}$ ($f$ is $L=4$ - smooth)

$$x_{t+1} = x_t - \frac{1}{4}\nabla f(x_t) =$$

$x_0 = 1$

$\|\nabla^2 f(x)\|_2 \leq 2$

$2 \leq 4$

$x_{t+1} = x_t - \frac{1}{4} \cdot 2 x_t$

$= \frac{1}{2} x_t$

$f(x_t) - f^* \stackrel{x}{=} x_t^2 = 2^{-t}$

# Can we go even faster?

▶ On $f(x) := x^2$: Stepsize $\gamma := \frac{1}{2}$ ($f$ is $L{=}2$ - smooth)

$$x_{t+1} = x_t - \frac{1}{2}\nabla f(x_t) = x_t - x_t = 0,$$

$$\left(1 - \frac{\mu}{L}\right)^t$$

$$\left(1 - \frac{2}{2}\right)^t \quad D^2 f(\lambda) = 2 \succeq 2$$

▶ ▶ converged in one step!

▶ Same $f(x) := x^2$: Stepsize $\gamma := \frac{1}{4}$ ($f$ is $L{=}4$ - smooth)

$$\mu = 2$$

$$x_{t+1} = x_t - \frac{1}{4}\nabla f(x_t) = x_t - \frac{x_t}{2} = \frac{x_t}{2},$$

so $f(x_t) = \dfrac{1}{2^t}$

$$\left(1 - \frac{2}{4}\right)^t$$

$$\text{Error} \leq \left(1 - \frac{1}{4}\right)^t = \left(\frac{3}{4}\right)^t$$

$$\frac{1}{2^t}$$

# Can we go even faster?

- On $f(x) := x^2$: Stepsize $\gamma := \frac{1}{2}$ ($f$ is $L{=}2$ - smooth)

$$x_{t+1} = x_t - \frac{1}{2}\nabla f(x_t) = x_t - x_t = 0,$$

  - converged in one step!

- Same $f(x) := x^2$: Stepsize $\gamma := \frac{1}{4}$ ($f$ is $L{=}4$ - smooth)

$$x_{t+1} = x_t - \frac{1}{4}\nabla f(x_t) = x_t - \frac{x_t}{2} = \frac{x_t}{2},$$

so $f(x_t) = f\left(\frac{x_0}{2^t}\right) = \frac{1}{2^{2t}}x_0^2.$

  - Exponential in $t$ !

# Strongly convex functions

## "Not too flat"

## Definition

Let $f : \mathbf{dom}(f) \to \mathbb{R}$ be a differentiable function, $X \subseteq \mathbf{dom}(f)$ convex and $\mu \in \mathbb{R}_+, \mu > 0$. Function $f$ is called strongly convex (with parameter $\mu$) over $X$ if

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in X.$$

$$0^{th}$$

$$1^{st} \quad \left( \nabla f(y) - \nabla f(x) \right)^\top (y - x) \geq \mu \|y - x\|^2 \quad \forall y, x$$

## Lemma (Exercise 21)

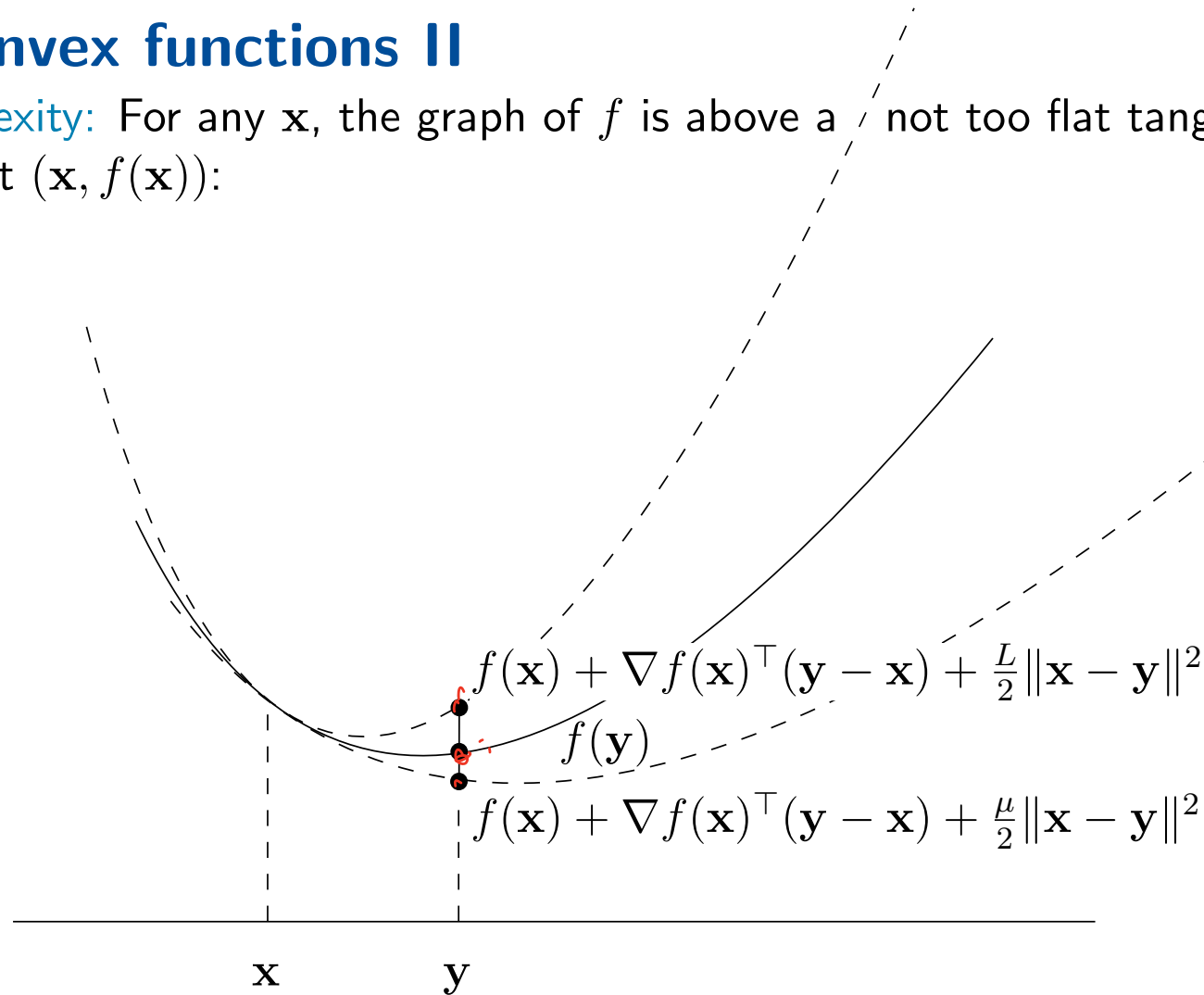*If $f$ is strongly convex with parameter $\mu > 0$, then $f$ is strictly convex and has a unique global minimum.*

$$2^{nd} \quad \nabla^2 f(x) \succeq \mu I \quad \forall x$$

# Strongly convex functions II

Strong convexity: For any $\mathbf{x}$, the graph of $f$ is above a / not too flat tangential paraboloid at $(\mathbf{x}, f(\mathbf{x}))$:

$$f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \tfrac{L}{2}\|\mathbf{x} - \mathbf{y}\|^2$$

$$f(\mathbf{y})$$

$$f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \tfrac{\mu}{2}\|\mathbf{x} - \mathbf{y}\|^2$$

$$\mathbf{x} \qquad \mathbf{y}$$

**Claim**

$$Df(x_t)^T(x_t - x^*) \geq f(x_t) - f(x^*) + \frac{\mu}{2}\|x_t - x^*\|^2$$

**Proof**

$$f(y) \geq f(x) + Df(x)^T(y-x) + \frac{\mu}{2}\|y-x\|^2$$

$$x = x_t, \quad y = x^*$$

$$\blacksquare$$

$$\langle Df(x_t), x_t - x^* \rangle \leq \frac{\gamma_t}{2}\|Df(x_t)\|^2 + \frac{1}{2\gamma}\left(\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2\right)$$

**Smoothness**

$$\|Df(x_t)\|^2 \leq 2L\left(f(x_t) - f^*\right)$$

**Strong Convexity**

$$\langle Df(x_t), x_t - x^* \rangle \geq \left(f(x_t) - f^*\right) + \frac{\mu}{2}\|x_t - x^*\|^2$$

$$\left(f(x_t) - f^*\right) + \frac{\mu}{2}\|x_t - x^*\|^2 \leq \left(f(x_t) - f^*\right) + \frac{L}{2}\left(\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2\right)$$

$$\Rightarrow \frac{L}{2}\|x_{t+1} - x^*\|^2 \leq \left(1 - \frac{\mu}{L}\right)\frac{L}{2}\|x_t - x^*\|^2 \leq \left(1 - \frac{\mu}{L}\right)^{t+1}\frac{L}{2}\|x_0 - x^*\|^2$$

# Smooth and strongly convex functions: $\mathcal{O}(\log(1/\varepsilon))$ steps

Want to show: $\lim_{t\to\infty} \mathbf{x}_t = \mathbf{x}^\star$

Vanilla Analysis:

$$\nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^\star) = \frac{\gamma}{2}\|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma}\left(\|\mathbf{x}_t - \mathbf{x}^\star\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2\right)$$

Now use stronger lower bound on left hand side, coming from strong convexity:

$$\nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^\star) \geq$$

# Smooth and strongly convex functions: $\mathcal{O}(\log(1/\varepsilon))$ steps

Want to show: $\lim_{t\to\infty} \mathbf{x}_t = \mathbf{x}^\star$

Vanilla Analysis:

$$\nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^\star) = \frac{\gamma}{2}\|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma}\left(\|\mathbf{x}_t - \mathbf{x}^\star\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2\right)$$

Now use stronger lower bound on left hand side, coming from strong convexity:

$$\nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^\star) \geq f(\mathbf{x}_t) - f(\mathbf{x}^\star) + \frac{\mu}{2}\|\mathbf{x}_t - \mathbf{x}^\star\|^2$$

Putting it together:

$$f(\mathbf{x}_t) - f(\mathbf{x}^\star) \leq$$

# Smooth and strongly convex functions: $\mathcal{O}(\log(1/\varepsilon))$ steps

Want to show: $\lim_{t\to\infty} \mathbf{x}_t = \mathbf{x}^\star$

Vanilla Analysis:

$$\nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^\star) = \frac{\gamma}{2}\|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma}\left(\|\mathbf{x}_t - \mathbf{x}^\star\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2\right)$$

Now use stronger lower bound on left hand side, coming from strong convexity:

$$\nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^\star) \geq f(\mathbf{x}_t) - f(\mathbf{x}^\star) + \frac{\mu}{2}\|\mathbf{x}_t - \mathbf{x}^\star\|^2$$

Putting it together:

$$f(\mathbf{x}_t) - f(\mathbf{x}^\star) \leq \frac{1}{2\gamma}\left(\gamma^2\|\nabla f(\mathbf{x}_t)\|^2 + \|\mathbf{x}_t - \mathbf{x}^\star\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2\right) - \frac{\mu}{2}\|\mathbf{x}_t - \mathbf{x}^\star\|^2.$$

Rewriting:

$$\|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2 \leq$$

# Smooth and strongly convex functions: $\mathcal{O}(\log(1/\varepsilon))$ steps

Want to show: $\lim_{t \to \infty} \mathbf{x}_t = \mathbf{x}^\star$

Vanilla Analysis:

$$\nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^\star) = \frac{\gamma}{2} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} \left( \|\mathbf{x}_t - \mathbf{x}^\star\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2 \right)$$

Now use stronger lower bound on left hand side, coming from strong convexity:

$$\nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^\star) \geq f(\mathbf{x}_t) - f(\mathbf{x}^\star) + \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}^\star\|^2$$

Putting it together:

$$f(\mathbf{x}_t) - f(\mathbf{x}^\star) \leq \frac{1}{2\gamma} \left( \gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 + \|\mathbf{x}_t - \mathbf{x}^\star\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2 \right) - \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}^\star\|^2.$$

Rewriting:

$$\|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2 \leq 2\gamma(f(\mathbf{x}^\star) - f(\mathbf{x}_t)) + \gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 + (1 - \mu\gamma)\|\mathbf{x}_t - \mathbf{x}^\star\|^2.$$

# Smooth and strongly convex functions: $\mathcal{O}(\log(1/\varepsilon))$ steps II

$$\underline{\|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2} \leq 2\gamma(f(\mathbf{x}^\star) - f(\mathbf{x}_t)) + \gamma^2\|\nabla f(\mathbf{x}_t)\|^2 + \underline{(1 - \mu\gamma)\|\mathbf{x}_t - \mathbf{x}^\star\|^2}.$$

**Squared distance to $\mathbf{x}^\star$ goes down by a constant factor, up to some "noise".**

## Theorem

Let $f : \mathbb{R}^d \to \mathbb{R}$ be differentiable with a global minimum $\mathbf{x}^\star$; suppose that $f$ is smooth with parameter $L$ and strongly convex with parameter $\mu > 0$. Choosing $\gamma := \frac{1}{L}$, gradient descent with arbitrary $\mathbf{x}_0$ satisfies the following two properties.

(i) Squared distances to $\mathbf{x}^\star$ are geometrically decreasing:

$$\|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2 \leq \left(1 - \frac{\mu}{L}\right)\|\mathbf{x}_t - \mathbf{x}^\star\|^2, \quad t \geq 0.$$

(ii) The absolute error after $T$ iterations is exponentially small in $T$:

$$f(\mathbf{x}_T) - f(\mathbf{x}^\star) \leq \frac{L}{2}\left(1 - \frac{\mu}{L}\right)^T \|\mathbf{x}_0 - \mathbf{x}^\star\|^2, \quad T > 0.$$

# Smooth and strongly convex functions: $\mathcal{O}(\log(1/\varepsilon))$ steps III

$$\underline{\|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2} \leq 2\gamma(f(\mathbf{x}^\star) - f(\mathbf{x}_t)) + \gamma^2\|\nabla f(\mathbf{x}_t)\|^2 + \underline{(1 - \mu\gamma)\|\mathbf{x}_t - \mathbf{x}^\star\|^2}.$$

Proof of (i).

Bounding the noise:

$$2\gamma(f(\mathbf{x}^\star) - f(\mathbf{x}_t)) + \gamma^2\|\nabla f(\mathbf{x}_t)\|^2 \quad =$$

# Smooth and strongly convex functions: $\mathcal{O}(\log(1/\varepsilon))$ steps III

$$\underline{\|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2} \leq 2\gamma(f(\mathbf{x}^\star) - f(\mathbf{x}_t)) + \gamma^2\|\nabla f(\mathbf{x}_t)\|^2 + \underline{(1 - \mu\gamma)\|\mathbf{x}_t - \mathbf{x}^\star\|^2}.$$

Proof of (i).

Bounding the noise: $\gamma = 1/L$ , sufficient decrease

$$
\begin{aligned}
2\gamma(f(\mathbf{x}^\star) - f(\mathbf{x}_t)) + \gamma^2\|\nabla f(\mathbf{x}_t)\|^2 &= \frac{2}{L}(f(\mathbf{x}^\star) - f(\mathbf{x}_t)) + \frac{1}{L^2}\|\nabla f(\mathbf{x}_t)\|^2 \\
&\leq \frac{2}{L}(f(\mathbf{x}_{t+1}) - f(\mathbf{x}_t)) + \frac{1}{L^2}\|\nabla f(\mathbf{x}_t)\|^2 \\
&\leq -\frac{1}{L^2}\|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{L^2}\|\nabla f(\mathbf{x}_t)\|^2 = 0.
\end{aligned}
$$

Hence, the noise is nonpositive, and we get (i):

$$\|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2 \leq (1 - \mu\gamma)\|\mathbf{x}_t - \mathbf{x}^\star\|^2 = \left(1 - \frac{\mu}{L}\right)\|\mathbf{x}_t - \mathbf{x}^\star\|^2.$$

# Smooth and strongly convex functions: $\mathcal{O}(\log(1/\varepsilon))$ steps III

Proof of (ii).

From (i):

$$\|\mathbf{x}_T - \mathbf{x}^\star\|^2 \leq \left(1 - \frac{\mu}{L}\right)^T \|\mathbf{x}_0 - \mathbf{x}^\star\|^2.$$

Smoothness together with $\nabla f(\mathbf{x}^\star) = \mathbf{0}$:

$$f(\mathbf{x}_T) - f(\mathbf{x}^\star) \leq \nabla f(\mathbf{x}^\star)^\top (\mathbf{x}_T - \mathbf{x}^\star) + \frac{L}{2}\|\mathbf{x}_T - \mathbf{x}^\star\|^2 = \frac{L}{2}\|\mathbf{x}_T - \mathbf{x}^\star\|^2.$$

Putting it together:

$$f(\mathbf{x}_T) - f(\mathbf{x}^\star) \leq \frac{L}{2}\|\mathbf{x}_T - \mathbf{x}^\star\|^2 \leq \frac{L}{2}\left(1 - \frac{\mu}{L}\right)^T \|\mathbf{x}_0 - \mathbf{x}^\star\|^2.$$

$\square$

# Smooth and strongly convex functions: $\mathcal{O}(\log(1/\varepsilon))$ steps IV

$R^2 := \|\mathbf{x}_0 - \mathbf{x}^\star\|^2.$

$$T \geq \frac{L}{\mu} \ln\left(\frac{R^2 L}{2\varepsilon}\right) \quad \Rightarrow \quad \text{error } \leq \frac{L}{2}\left(1 - \frac{\mu}{L}\right)^T R^2 \leq \varepsilon.$$

**Conclusion:** To reach absolute error at most $\varepsilon$, we only need $\mathcal{O}(\log\frac{1}{\varepsilon})$ iterations, e.g.

▶ $\frac{L}{\mu}\ln(50 \cdot R^2 L)$ iterations for error $0.01$ ...

▶ ... as opposed to $50 \cdot R^2 L$ in the smooth case

In Practice:

What if we don't know the smoothness parameter $L$?

$\rightarrow$ (similar to) **Exercise 15**

$$f(x) = |x| \qquad \lor \qquad \frac{\partial |x|}{\partial x} = \begin{cases} 1 & x > 0 \\ -1 & x < 0 \\ [-1,1] & x = 0 \end{cases}$$

$$f(x) = \|x\|_1$$
$$= \sum_i |x_i|$$

$$g(x) = \left[ \quad \begin{cases} 1 & x_i > 0 \\ -1 & x_i < 0 \\ [-1,1] & x_i = 0 \end{cases} \quad \right]$$

$$f(x) = \frac{1}{2} x^T A x - bx + \lambda \|x\|_1$$

regular

$$A x - b + \lambda g(x)$$

$$\min_x f(x) + \lambda \|x\|_1$$

$$\min_x f(x)$$
$$s.t.$$
$$\|x\|_1 \le \lambda$$

# Chapter 3

## Projected Gradient Descent

$$\tilde{f}(x) = f(x) + \mathbb{1}\{x \in C\}$$

# Constrained Optimization

Constrained Optimization Problem *Smooth*

$$\begin{aligned} \text{minimize} \quad & f(\mathbf{x}) \\ \text{subject to} \quad & \mathbf{x} \in X \end{aligned}$$



$f$

$f(\boldsymbol{x})$

$\boldsymbol{x}$

$X \subseteq \mathbb{R}^d$

Solving Constrained Optimization Problems

A  Projected Gradient Descent
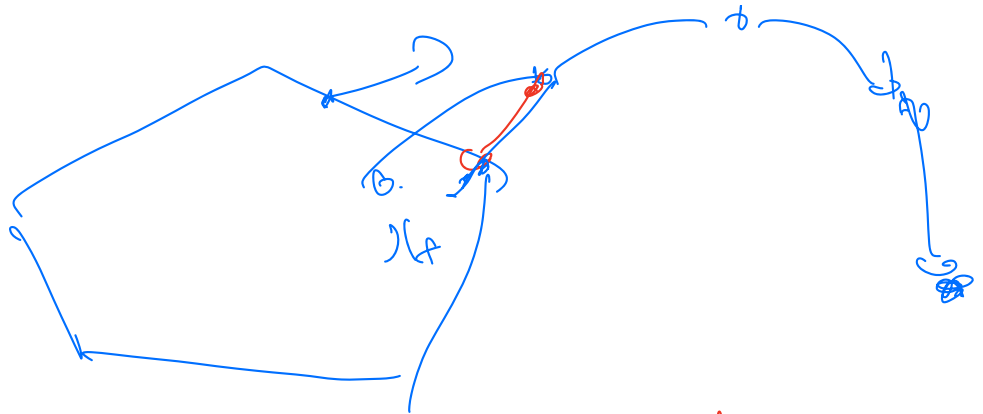
B  Transform it into an *unconstrained* problem

$\rightarrow$ $x^*$ is opt of $\min\limits_{x \in X} f(x)$

if

~~$Df(x^*) = 0$~~ ?



$$\langle Df(x^*), x - x^* \rangle \geq 0$$

$$\forall x \in X$$

$\rightarrow$

$$x_{t+1} = x_t - r Df(x_t)$$



$$\Pi_X(y) = \arg\min\limits_{x \in X} \|x - y\|_2$$

$\hookrightarrow C \times X$

$$x_{t+1} = \Pi_X(x_t - r Df(x_t))$$

$$x^* = \Pi_{\Pi}\left(y^* - \gamma \nabla f(y^x)\right)$$
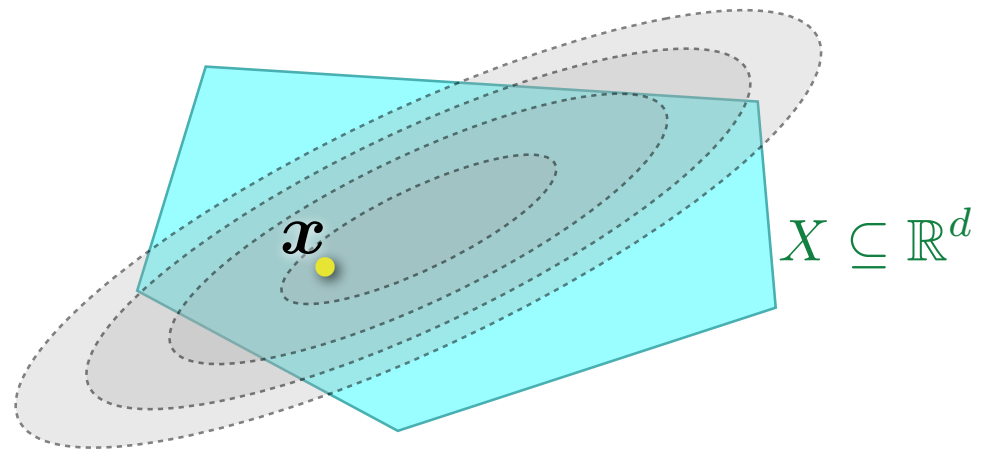
$\gamma \nabla f(x^x)$

$y^*$

normal cone

$f$

# Constrained Optimization

Solving Constrained Optimization Problems

$$
\begin{aligned}
\text{minimize} \qquad & f(\mathbf{x}) \\
\text{subject to} \qquad & \mathbf{x} \in X
\end{aligned}
$$

▶ Here: Projected Gradient Descent

$$x$$

$$X \subseteq \mathbb{R}^d$$

# Projected Gradient Descent
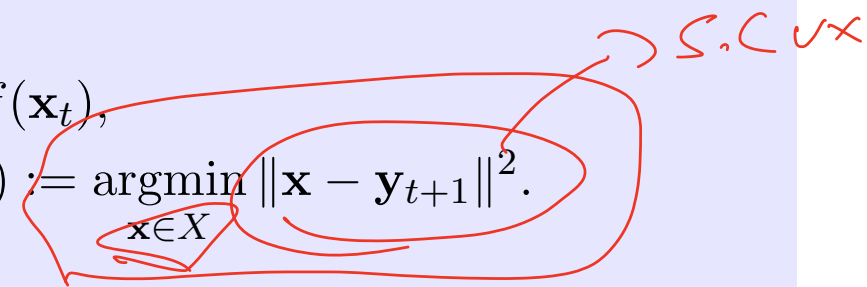
Idea: project onto $X$ after every step: $\Pi_X(\mathbf{y}) := \arg\min_{\mathbf{x} \in X} \|\mathbf{x} - \mathbf{y}\|$



$$x$$

$$X \subseteq \mathbb{R}^d$$

$$\Pi_X(\boldsymbol{y})$$

$$-\nabla f(\boldsymbol{x})$$

$$\boldsymbol{y}$$

Projected gradient descent: $\mathbf{x}_{t+1} := \Pi_X\big[\mathbf{x}_t - \gamma \nabla f(\mathbf{x}_t)\big]$

# The Algorithm

**Projected gradient descent:**

$$
\begin{aligned}
\mathbf{y}_{t+1} &:= \mathbf{x}_t - \gamma \nabla f(\mathbf{x}_t), \\
\mathbf{x}_{t+1} &:= \Pi_X(\mathbf{y}_{t+1}) := \operatorname*{argmin}_{\mathbf{x} \in X} \|\mathbf{x} - \mathbf{y}_{t+1}\|^2.
\end{aligned}
$$

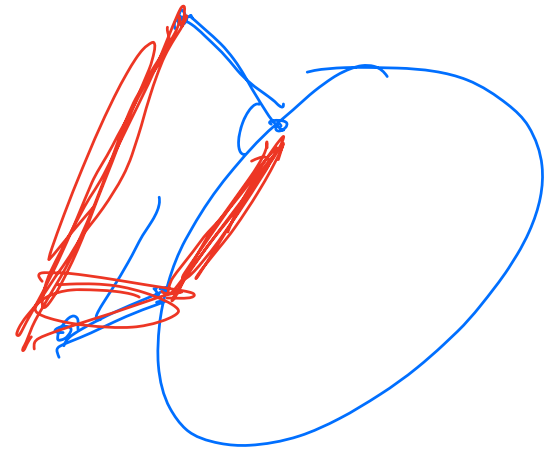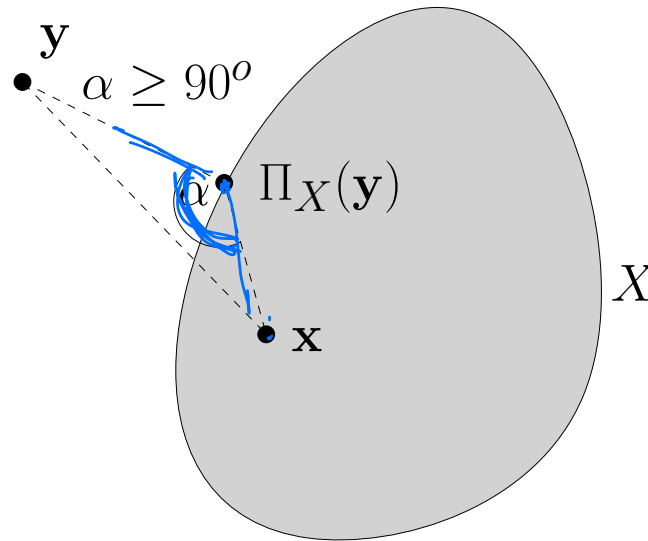for **timesteps** $t = 0, 1, \ldots,$ and **stepsize** $\gamma \geq 0$.

# Properties of Projection

Fact

Let $X \subseteq \mathbb{R}^d$ be closed and convex, $\mathbf{x} \in X, \mathbf{y} \in \mathbb{R}^d$. Then

(i) $(\mathbf{x} - \Pi_X(\mathbf{y}))^\top (\mathbf{y} - \Pi_X(\mathbf{y})) \leq 0.$

(ii) $\|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2.$

$\mathbf{y}$

$\alpha \geq 90^o$

$\alpha$

$\Pi_X(\mathbf{y})$

$\mathbf{x}$

$X$

# Properties of Projection II

**Fact**

*Let $X \subseteq \mathbb{R}^d$ be closed and convex, $\mathbf{x} \in X, \mathbf{y} \in \mathbb{R}^d$. Then*

(i) $(\mathbf{x} - \Pi_X(\mathbf{y}))^\top (\mathbf{y} - \Pi_X(\mathbf{y})) \leq 0.$

(ii) $\|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2.$

**Proof.**

(i) $\Pi_X(\mathbf{y})$ is minimizer of (differentiable) convex function $d_\mathbf{y}(\mathbf{x}) = \|\mathbf{x} - \mathbf{y}\|^2$ over $X$.
By first-order characterization of optimality (**Lemma 1.28**),

$$0 \quad \leq \quad \nabla d_\mathbf{y}(\Pi_X(\mathbf{y}))^\top (\mathbf{x} - \Pi_X(\mathbf{y}))$$

$$(a-b)^2 + (b-c)^2 \leq (a-c)^2$$

$$2\left(\|a-b\|^2 + \|b-c\|^2\right) \geq \left(\|a-b\| + \|b-c\|\right)^2 \geq \|a-c\|^2$$

# Properties of Projection II

### Fact

*Let $X \subseteq \mathbb{R}^d$ be closed and convex, $\mathbf{x} \in X, \mathbf{y} \in \mathbb{R}^d$. Then*

(i) $(\mathbf{x} - \Pi_X(\mathbf{y}))^\top (\mathbf{y} - \Pi_X(\mathbf{y})) \leq 0.$

(ii) $\|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2.$

*Handwritten annotations:*

$\Pi_X(y) = \arg\min_{z \in X} \|z - y\|_2^2$

$\frac{1}{2}\|x \cdot y\|^2$

$y = y_{t+1}$

$\Pi_X(y) = x_{t+1}$

$x = x^*$

### Proof.

(i) $\Pi_X(\mathbf{y})$ is minimizer of (differentiable) convex function $d_{\mathbf{y}}(\mathbf{x}) = \|\mathbf{x} - \mathbf{y}\|^2$ over $X$.

By first-order characterization of optimality (**Lemma 1.28**),

$$
\begin{aligned}
0 &\leq \nabla d_{\mathbf{y}}(\Pi_X(\mathbf{y}))^\top (\mathbf{x} - \Pi_X(\mathbf{y})) \\
&= 2(\Pi_X(\mathbf{y}) - \mathbf{y})^\top (\mathbf{x} - \Pi_X(\mathbf{y})) \\
\Leftrightarrow \quad 0 &\geq 2(\mathbf{y} - \Pi_X(\mathbf{y}))^\top (\mathbf{x} - \Pi_X(\mathbf{y})) \\
\Leftrightarrow \quad 0 &\geq (\mathbf{x} - \Pi_X(\mathbf{y}))^\top (\mathbf{y} - \Pi_X(\mathbf{y}))
\end{aligned}
$$

*Handwritten annotations:*

$\|x_{t+1} - x^*\|^2 + \|y_{t+1} - x_{t+1}\|^2 \leq \|y_{t+1} - x^*\|^2 \rightarrow$ Piece 2

# Properties of Projection III

**Fact**

*Let $X \subseteq \mathbb{R}^d$ be closed and convex, $\mathbf{x} \in X, \mathbf{y} \in \mathbb{R}^d$. Then*

  (i) $(\mathbf{x} - \Pi_X(\mathbf{y}))^\top (\mathbf{y} - \Pi_X(\mathbf{y})) \leq 0$.

  (ii) $\|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2$.

**Proof.**

(ii)

$$\mathbf{v} := (\mathbf{x} - \Pi_X(\mathbf{y})), \quad \mathbf{w} := (\mathbf{y} - \Pi_X(\mathbf{y})).$$

By (i),

$$0 \geq 2\mathbf{v}^\top \mathbf{w} \quad =$$

# Properties of Projection III

**Fact**

*Let $X \subseteq \mathbb{R}^d$ be closed and convex, $\mathbf{x} \in X, \mathbf{y} \in \mathbb{R}^d$. Then*

(i) $(\mathbf{x} - \Pi_X(\mathbf{y}))^\top (\mathbf{y} - \Pi_X(\mathbf{y})) \leq 0$.

(ii) $\|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2$.

**Proof.**

(ii)

$$\mathbf{v} := (\mathbf{x} - \Pi_X(\mathbf{y})), \quad \mathbf{w} := (\mathbf{y} - \Pi_X(\mathbf{y})).$$

By (i),

$$
\begin{aligned}
0 \geq 2\mathbf{v}^\top \mathbf{w} &= \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 - \|\mathbf{v} - \mathbf{w}\|^2 \\
&= \|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 - \|\mathbf{x} - \mathbf{y}\|^2.
\end{aligned}
$$

$\square$

# Results for **projected** gradient descent over closed and convex $X$

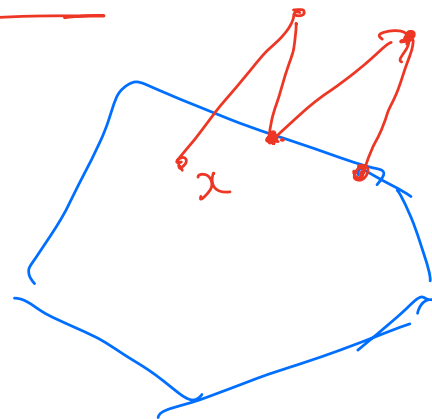The same number of steps as gradient over $\mathbb{R}^d$!

▶ Lipschitz convex functions over $X$: $\mathcal{O}(1/\varepsilon^2)$ steps

▶ Smooth convex functions over $X$: $\mathcal{O}(1/\varepsilon)$ steps

▶ Smooth and strongly convex functions over $X$: $\mathcal{O}(\log(1/\varepsilon))$ steps

We will adapt the previous proofs for gradient descent.

BUT:

▶ Each step involves a projection onto $X$

▶ may or may not be efficient (in relevant cases, it is)...

$$x_{t+1} = \Pi_X(x_t - \gamma g_t)$$

# Lipschitz convex functions over $X$: $\mathcal{O}(1/\varepsilon^2)$ steps

Assume that all gradients of $f$ are bounded in norm over closed and convex $X$.

▶ Equivalent to $f$ being Lipschitz over $X$ (Theorem 1.10; Exercise 12).
▶ Many interesting functions are Lipschitz over bounded sets $X$.

## Theorem (same as the unconstrained one, but more useful)

*Let $f : \mathbb{R}^d \to \mathbb{R}$ be convex and differentiable, $X \subseteq \mathbb{R}^d$ closed and convex, $\mathbf{x}^\star$ a minimizer of $f$ over $X$; furthermore, suppose that $\|\mathbf{x}_0 - \mathbf{x}^\star\| \leq R$ with $\mathbf{x}_0 \in X$, and that $\|\nabla f(\mathbf{x})\| \leq B$ for all $\mathbf{x} \in X$. Choosing the constant stepsize*

$$\gamma := \frac{R}{B\sqrt{T}},$$

*projected gradient descent yields*

$$\frac{1}{T} \sum_{t=0}^{T-1} f(\mathbf{x}_t) - f(\mathbf{x}^\star) \leq \frac{RB}{\sqrt{T}}.$$

$$x_{t+1} = \prod_X \underbrace{\left( x_t - \gamma g_t \right)}_{y_{t+1}}$$

Proof.

▶ Replace $\mathbf{x}_{t+1}$ in the vanilla analysis with $\mathbf{y}_{t+1}$ (the unprojected gradient step):

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star) = \frac{1}{2\gamma} \left( \gamma^2 \|\mathbf{g}_t\|^2 + \|\mathbf{x}_t - \mathbf{x}^\star\|^2 - \|\underline{\mathbf{y}_{t+1}} - \mathbf{x}^\star\|^2 \right).$$
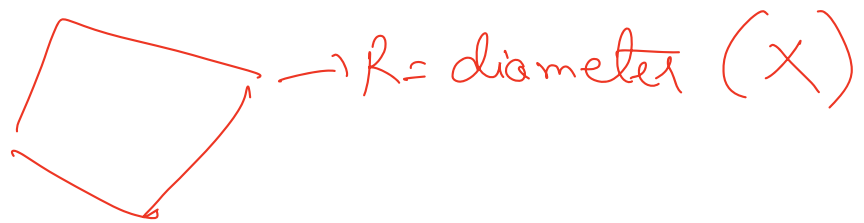
▶ Use Fact (ii):  $\|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2.$

▶ With $\mathbf{x} = \mathbf{x}^\star, \mathbf{y} = \mathbf{y}_{t+1}$, we have $\Pi_X(\mathbf{y}) = \mathbf{x}_{t+1}$ , and hence

$$\|x_{t+1} - x^*\|^2 = \|\Pi_X(x_t - \gamma g_t) - \Pi_X(x^*)\|^2$$

$$\Pi_X(x^*) = x^*$$

$$\leq \|x_t - \gamma g_t - x^*\|^2$$

$$= \|x_t - x^*\|_2^2 - 2\gamma g_t^\top (x_t - x^*) + \gamma^2 \|g\|^2$$

$$\Rightarrow \quad \sum_t \underbrace{g_t^T(x_t - \ell^*)}_{\geq f(x_t) - f(x^*)} \leq \frac{1}{2\gamma}\Big(\|x_t - x^*\|_2^2 \underbrace{- \|x_{t+c} - x^*\|^2\Big)}_{\|x_0 - x^*\| \leq \boxed{R}} + \frac{\gamma}{2}\|g_t\|^2$$

$$\|g_t\| \leq B$$

$$\Rightarrow \quad f(\overline{x_T}) - f(x^*) \leq \frac{RB}{\sqrt{T}}$$



$\longrightarrow R = $ diameter $(X)$

# Lipschitz convex functions: $\mathcal{O}(1/\varepsilon^2)$ steps II

▶ Replace $\mathbf{x}_{t+1}$ in the vanilla analysis with $\mathbf{y}_{t+1}$ (the unprojected gradient step):

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star) = \frac{1}{2\gamma}\left(\gamma^2 \|\mathbf{g}_t\|^2 + \|\mathbf{x}_t - \mathbf{x}^\star\|^2 - \|\underline{\mathbf{y}_{t+1}} - \mathbf{x}^\star\|^2\right).$$

▶ Use Fact (ii):    $\|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2.$

▶ With $\mathbf{x} = \mathbf{x}^\star, \mathbf{y} = \mathbf{y}_{t+1}$, we have $\Pi_X(\mathbf{y}) = \mathbf{x}_{t+1}$ , and hence

$$\|\mathbf{x}^\star - \mathbf{x}_{t+1}\|^2 \qquad\qquad\qquad\qquad \leq \|\mathbf{x}^\star - \mathbf{y}_{t+1}\|^2$$

▶ We go back to the original vanilla analyis and continue from there as before:

# Lipschitz convex functions: $\mathcal{O}(1/\varepsilon^2)$ steps II

▶ Replace $\mathbf{x}_{t+1}$ in the vanilla analysis with $\mathbf{y}_{t+1}$ (the unprojected gradient step):

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star) = \frac{1}{2\gamma} \left( \gamma^2 \|\mathbf{g}_t\|^2 + \|\mathbf{x}_t - \mathbf{x}^\star\|^2 - \|\underline{\mathbf{y}_{t+1}} - \mathbf{x}^\star\|^2 \right).$$

▶ Use Fact (ii):  $\|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2.$

▶ With $\mathbf{x} = \mathbf{x}^\star, \mathbf{y} = \mathbf{y}_{t+1}$, we have $\Pi_X(\mathbf{y}) = \mathbf{x}_{t+1}$ , and hence

$$\|\mathbf{x}^\star - \mathbf{x}_{t+1}\|^2 \qquad\qquad\qquad \leq \|\mathbf{x}^\star - \mathbf{y}_{t+1}\|^2$$

▶ We go back to the original vanilla analyis and continue from there as before:

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star) \leq \frac{1}{2\gamma} \left( \gamma^2 \|\mathbf{g}_t\|^2 + \|\mathbf{x}_t - \mathbf{x}^\star\|^2 - \|\underline{\mathbf{x}_{t+1}} - \mathbf{x}^\star\|^2 \right).$$
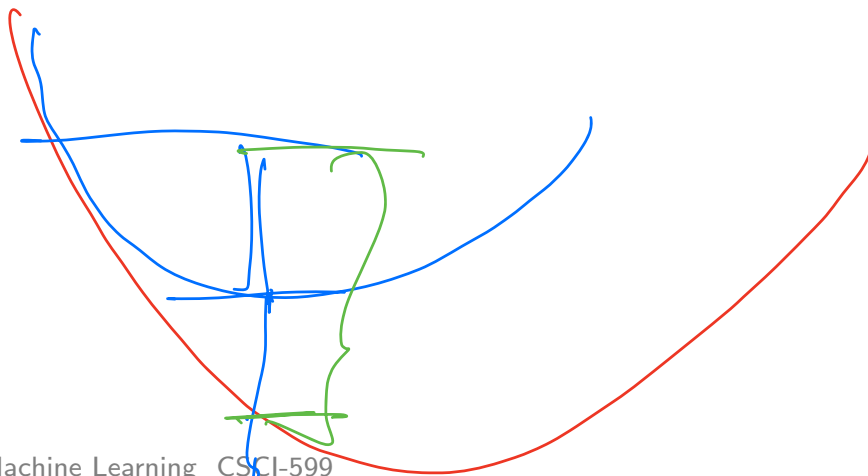
# Smooth functions over $X$

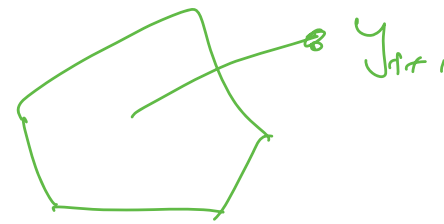$$f(x_{t+1}) \leq f(x_t) - \frac{1}{2C} \|\nabla f(x_t)\|^2$$

Recall:

$f$ is called smooth (with parameter $L$) over $X$ if

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2}\|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in X.$$

$$y_{t+1} = x_t - \gamma g_t$$

$$x_{t+1} = \Pi_X(y_{t+1})$$

# Sufficient decrease

### Lemma
Let $f : \mathbb{R}^d \to \mathbb{R}$ be differentiable and smooth with parameter $L$ over $X$. Choosing stepsize

$$\gamma := \frac{1}{L},$$

*projected* gradient descent with arbitrary $\mathbf{x}_0 \in X$ satisfies

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L}\|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2}\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2, \quad t \geq 0.$$

### Remark
More specifically, this already holds if $f$ is smooth with parameter $L$ over the line segment connecting $\mathbf{x}_t$ and $\mathbf{x}_{t+1}$.

$$f(x_{t+1}) \leq f(x_t) - \nabla f(x_t)^{\top}(x_{t+1} - x_t) + \frac{L}{2}\|x_{t+1} - x_t\|^2$$

# Sufficient decrease II

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L}\|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2}\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2.$$

Proof.
Use smoothness                                                                :

# Sufficient decrease II

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L}\|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2}\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2.$$

Proof.

Use smoothness                                                                                          :

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2}\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2$$

# Sufficient decrease II

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L}\|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2}\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2.$$

Proof.
Use smoothness

$$\mathbf{y}_{t+1} = \mathbf{y}_t - \frac{1}{L}\nabla f(\mathbf{x}_{t-1})$$

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top(\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2}\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2$$

$$= f(\mathbf{x}_t) - L(\mathbf{y}_{t+1} - \mathbf{x}_t)^\top(\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2}\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2$$

# Sufficient decrease II

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L}\|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2}\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2.$$

Proof.

Use smoothness, $\mathbf{y}_{t+1} - \mathbf{x}_t = -\nabla f(\mathbf{x}_t)/L$                         :

$$
\begin{aligned}
f(\mathbf{x}_{t+1}) &\leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2}\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\
&= f(\mathbf{x}_t) - L(\mathbf{y}_{t+1} - \mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2}\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\
&= f(\mathbf{x}_t) - \frac{L}{2}\left(\|\mathbf{y}_{t+1} - \mathbf{x}_t\|^2 + \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 - \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2\right) + \frac{L}{2}\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2
\end{aligned}
$$

# Sufficient decrease II

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L}\|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2}\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2.$$

Proof.

Use smoothness, $\mathbf{y}_{t+1} - \mathbf{x}_t = -\nabla f(\mathbf{x}_t)/L$ , $2\mathbf{v}^\top\mathbf{w} = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 - \|\mathbf{v}-\mathbf{w}\|^2$:

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2}\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2$$

$$= f(\mathbf{x}_t) - L(\mathbf{y}_{t+1} - \mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2}\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2$$

$$= f(\mathbf{x}_t) - \frac{L}{2}\left(\|\mathbf{y}_{t+1} - \mathbf{x}_t\|^2 + \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 - \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2\right) + \frac{L}{2}\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2$$

$$= f(\mathbf{x}_t) - \frac{L}{2}\|\mathbf{y}_{t+1} - \mathbf{x}_t\|^2 + \frac{L}{2}\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2$$

# Sufficient decrease II

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L}\|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2}\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2.$$

Proof.
Use smoothness, $\mathbf{y}_{t+1} - \mathbf{x}_t = -\nabla f(\mathbf{x}_t)/L$ , $2\mathbf{v}^\top \mathbf{w} = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 - \|\mathbf{v} - \mathbf{w}\|^2$:

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2}\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2$$

$$= f(\mathbf{x}_t) - L(\mathbf{y}_{t+1} - \mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2}\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2$$

$$= f(\mathbf{x}_t) - \frac{L}{2}\left(\|\mathbf{y}_{t+1} - \mathbf{x}_t\|^2 + \underline{\|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2} - \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2\right) + \frac{L}{2}\underline{\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2}$$

$$= f(\mathbf{x}_t) - \frac{L}{2}\|\mathbf{y}_{t+1} - \mathbf{x}_t\|^2 + \frac{L}{2}\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2$$

$$= f(\mathbf{x}_t) - \frac{1}{2L}\|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2}\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2.$$

# Smooth convex functions over $X$: $\mathcal{O}(1/\varepsilon)$ steps

## Theorem

*Let $f : \mathbb{R}^d \to \mathbb{R}$ be convex and differentiable. Let $X \subseteq \mathbb{R}^d$ be a closed convex set, and assume that there is a minimizer $\mathbf{x}^\star$ of $f$ over $X$; furthermore, suppose that $f$ is smooth over $X$ with parameter $L$. Choosing stepsize*

$$\gamma := \frac{1}{L},$$

*projected* gradient descent yields

$$f(\mathbf{x}_T) - f(\mathbf{x}^\star) \leq \frac{L}{2T}\|\mathbf{x}_0 - \mathbf{x}^\star\|^2, \quad T > 0.$$

# Smooth convex functions over $X$: $\mathcal{O}(1/\varepsilon)$ steps II

$$f(\mathbf{x}_T) - f(\mathbf{x}^\star) \leq \frac{L}{2T}\|\mathbf{x}_0 - \mathbf{x}^\star\|^2, \quad T > 0.$$

Proof.
As before, use sufficient decrease to bound sum of squared gradients in vanilla analysis:

$$\frac{1}{2L}\|\nabla f(\mathbf{x}_t)\|^2 \leq f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}) + \frac{L}{2}\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2$$

But now: extra term $\frac{L}{2}\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2$.

Compensate in the vanilla analysis itself! □

# Recall: Constrained vanilla analysis

Proof.

▶ Replace $\mathbf{x}_{t+1}$ in the vanilla analysis with $\mathbf{y}_{t+1}$ (the unprojected gradient step):

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star) = \frac{1}{2\gamma} \left( \gamma^2 \|\mathbf{g}_t\|^2 + \|\mathbf{x}_t - \mathbf{x}^\star\|^2 - \|\mathbf{y}_{t+1} - \mathbf{x}^\star\|^2 \right).$$

▶ Use Fact (ii):    $\|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2.$
With $\mathbf{x} = \mathbf{x}^\star, \mathbf{y} = \mathbf{y}_{t+1}$, we have $\Pi_X(\mathbf{y}) = \mathbf{x}_{t+1}$, and hence

# Recall: Constrained vanilla analysis

Proof.

▶ Replace $\mathbf{x}_{t+1}$ in the vanilla analysis with $\mathbf{y}_{t+1}$ (the unprojected gradient step):

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star) = \frac{1}{2\gamma} \left( \gamma^2 \|\mathbf{g}_t\|^2 + \|\mathbf{x}_t - \mathbf{x}^\star\|^2 - \|\mathbf{y}_{t+1} - \mathbf{x}^\star\|^2 \right).$$

▶ Use Fact (ii): $\quad \|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 \le \|\mathbf{x} - \mathbf{y}\|^2$.
With $\mathbf{x} = \mathbf{x}^\star, \mathbf{y} = \mathbf{y}_{t+1}$, we have $\Pi_X(\mathbf{y}) = \mathbf{x}_{t+1}$, and hence

$$\|\mathbf{x}^\star - \mathbf{x}_{t+1}\|^2 + \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \quad \le \|\mathbf{x}^\star - \mathbf{y}_{t+1}\|^2$$

▶ We get back to the vanilla analysis. . . but with a saving!

# Recall: Constrained vanilla analysis

Proof.

▶ Replace $\mathbf{x}_{t+1}$ in the vanilla analysis with $\mathbf{y}_{t+1}$ (the unprojected gradient step):

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star) = \frac{1}{2\gamma} \left( \gamma^2 \|\mathbf{g}_t\|^2 + \|\mathbf{x}_t - \mathbf{x}^\star\|^2 - \|\mathbf{y}_{t+1} - \mathbf{x}^\star\|^2 \right).$$

▶ Use Fact (ii):   $\|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 \le \|\mathbf{x} - \mathbf{y}\|^2$.
With $\mathbf{x} = \mathbf{x}^\star, \mathbf{y} = \mathbf{y}_{t+1}$, we have $\Pi_X(\mathbf{y}) = \mathbf{x}_{t+1}$, and hence

$$\|\mathbf{x}^\star - \mathbf{x}_{t+1}\|^2 + \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \quad \le \|\mathbf{x}^\star - \mathbf{y}_{t+1}\|^2$$

▶ We get back to the vanilla analysis... but with a saving!

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star) \le \frac{1}{2\gamma} \left( \gamma^2 \|\mathbf{g}_t\|^2 + \|\mathbf{x}_t - \mathbf{x}^\star\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2 - \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \right)$$

$\square$

# Smooth convex functions over $X$: $\mathcal{O}(1/\varepsilon)$ steps III

$$f(\mathbf{x}_T) - f(\mathbf{x}^\star) \leq \frac{L}{2T}\|\mathbf{x}_0 - \mathbf{x}^\star\|^2, \quad T > 0.$$

Proof.

Use $f(\mathbf{x}_t) - f(\mathbf{x}^\star) \leq \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star)$ (convexity)

$$\sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^\star)) \leq \sum_{t=0}^{T-1} \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star)$$

# Smooth convex functions over $X$: $\mathcal{O}(1/\varepsilon)$ steps III

$$f(\mathbf{x}_T) - f(\mathbf{x}^\star) \leq \frac{L}{2T}\|\mathbf{x}_0 - \mathbf{x}^\star\|^2, \quad T > 0.$$

Proof.
Use $f(\mathbf{x}_t) - f(\mathbf{x}^\star) \leq \mathbf{g}_t^\top(\mathbf{x}_t - \mathbf{x}^\star)$ (convexity), vanilla analysis with saving, $\gamma = 1/L$:

$$\sum_{t=0}^{T-1}(f(\mathbf{x}_t) - f(\mathbf{x}^\star)) \leq \sum_{t=0}^{T-1}\mathbf{g}_t^\top(\mathbf{x}_t - \mathbf{x}^\star)$$

$$\leq \frac{1}{2L}\sum_{t=0}^{T-1}\|\mathbf{g}_t\|^2 + \frac{L}{2}\|\mathbf{x}_0 - \mathbf{x}^\star\|^2 - \frac{L}{2}\sum_{t=0}^{T-1}\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2.$$

# Smooth convex functions over $X$: $\mathcal{O}(1/\varepsilon)$ steps III

$$f(\mathbf{x}_T) - f(\mathbf{x}^\star) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^\star\|^2, \quad T > 0.$$

Proof.
Use $f(\mathbf{x}_t) - f(\mathbf{x}^\star) \leq \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star)$ (convexity), vanilla analysis with saving, $\gamma = 1/L$:

$$
\begin{aligned}
\sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^\star)) \;&\leq\; \sum_{t=0}^{T-1} \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star) \\
&\leq\; \frac{1}{2L} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2 + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^\star\|^2 - \frac{L}{2} \sum_{t=0}^{T-1} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2.
\end{aligned}
$$

Use sufficient decrease to bound $\frac{1}{2L} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2$ by

# Smooth convex functions over $X$: $\mathcal{O}(1/\varepsilon)$ steps III

$$f(\mathbf{x}_T) - f(\mathbf{x}^\star) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^\star\|^2, \quad T > 0.$$

Proof.

Use $f(\mathbf{x}_t) - f(\mathbf{x}^\star) \leq \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star)$ (convexity), vanilla analysis with saving, $\gamma = 1/L$:

$$
\begin{aligned}
\sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^\star)) \;&\leq\; \sum_{t=0}^{T-1} \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star) \\
&\leq\; \frac{1}{2L} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2 + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^\star\|^2 - \frac{L}{2} \sum_{t=0}^{T-1} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 .
\end{aligned}
$$

Use sufficient decrease to bound $\frac{1}{2L} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2$ by

$$\sum_{t=0}^{T-1} \left( f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}) + \frac{L}{2} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \right)$$

# Smooth convex functions over $X$: $\mathcal{O}(1/\varepsilon)$ steps III

$$f(\mathbf{x}_T) - f(\mathbf{x}^\star) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^\star\|^2, \quad T > 0.$$

Proof.

Use $f(\mathbf{x}_t) - f(\mathbf{x}^\star) \leq \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star)$ (convexity), vanilla analysis with saving, $\gamma = 1/L$:

$$
\begin{aligned}
\sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^\star)) &\leq \sum_{t=0}^{T-1} \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star) \\
&\leq \frac{1}{2L} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2 + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^\star\|^2 - \frac{L}{2} \sum_{t=0}^{T-1} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2.
\end{aligned}
$$

Use sufficient decrease to bound $\frac{1}{2L} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2$ by

$$\sum_{t=0}^{T-1} \left( f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}) + \frac{L}{2} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \right) = f(\mathbf{x}_0) - f(\mathbf{x}_T) + \frac{L}{2} \sum_{t=0}^{T-1} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2.$$

# Smooth convex functions over $X$: $\mathcal{O}(1/\varepsilon)$ steps III

$$f(\mathbf{x}_T) - f(\mathbf{x}^\star) \le \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^\star\|^2, \quad T > 0.$$

**Proof.**

Use $f(\mathbf{x}_t) - f(\mathbf{x}^\star) \le \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star)$ (convexity), vanilla analysis with saving, $\gamma = 1/L$:

$$
\begin{aligned}
\sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^\star)) \;&\le\; \sum_{t=0}^{T-1} \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star) \\
&\le\; \frac{1}{2L} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2 + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^\star\|^2 - \frac{L}{2} \sum_{t=0}^{T-1} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2.
\end{aligned}
$$

Use sufficient decrease to bound $\frac{1}{2L} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2$ by

$$\sum_{t=0}^{T-1} \left( f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}) + \frac{L}{2} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \right) = f(\mathbf{x}_0) - f(\mathbf{x}_T) + \frac{L}{2} \sum_{t=0}^{T-1} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2.$$

# Smooth convex functions over $X$: $\mathcal{O}(1/\varepsilon)$ steps IV

$$f(\mathbf{x}_T) - f(\mathbf{x}^\star) \leq \frac{L}{2T}\|\mathbf{x}_0 - \mathbf{x}^\star\|^2, \quad T > 0.$$

**Proof.**

Putting it together: extra terms cancel, and as in unconstrained case, we get

$$\sum_{t=1}^{T} (f(\mathbf{x}_t) - f(\mathbf{x}^\star)) \leq \frac{L}{2}\|\mathbf{x}_0 - \mathbf{x}^\star\|^2.$$

**Exercise ??**: again, we make progress in every step (not immediate from sufficient decrease here). Hence,

$$f(\mathbf{x}_T) - f(\mathbf{x}^\star) \leq \frac{1}{T}\sum_{t=1}^{T} (f(\mathbf{x}_t) - f(\mathbf{x}^\star)) \leq \frac{L}{2T}\|\mathbf{x}_0 - \mathbf{x}^\star\|^2.$$

$\square$

# Smooth and strongly convex functions over $X$

Recall:

$f$ is strongly convex (with parameter $\mu$) over $X$ if

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in X.$$

# Smooth and strongly convex functions over $X$

**Exercise ??**: a strongly convex function has a unique minimizer $\mathbf{x}^\star$ of $f$ over $X$.

We prove that projected gradient descent converges to $\mathbf{x}^\star$.

# Smooth and strongly convex functions over $X$: $\mathcal{O}(\log(1/\varepsilon))$ steps

## Theorem

*Let $f : \mathbb{R}^d \to \mathbb{R}$ be convex and differentiable. Let $X \subseteq \mathbb{R}^d$ be a nonempty closed and convex set and suppose that $f$ is smooth over $X$ with parameter $L$ and strongly convex over $X$ with parameter $\mu > 0$. Choosing $\gamma := \frac{1}{L}$, projected gradient descent with arbitrary $\mathbf{x}_0$ satisfies the following two properties.*

(i) *Squared distances to $\mathbf{x}^\star$ are geometrically decreasing:*

$$\|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2 \leq \left(1 - \frac{\mu}{L}\right) \|\mathbf{x}_t - \mathbf{x}^\star\|^2, \quad t \geq 0.$$

(ii) *The absolute error after $T$ iterations is exponentially small in $T$:*

$$
\begin{aligned}
f(\mathbf{x}_T) - f(\mathbf{x}^\star) \;\leq\; & \|\nabla f(\mathbf{x}^\star)\| \left(1 - \frac{\mu}{L}\right)^{T/2} \|\mathbf{x}_0 - \mathbf{x}^\star\| \\
& + \; \frac{L}{2} \left(1 - \frac{\mu}{L}\right)^{T} \|\mathbf{x}_0 - \mathbf{x}^\star\|^2, \quad T > 0.
\end{aligned}
$$

## Theorem

*Let $f : \mathbb{R}^d \to \mathbb{R}$ be convex and differentiable. Let $X \subseteq \mathbb{R}^d$ be a nonempty closed and convex set and suppose that $f$ is smooth over $X$ with parameter $L$ and strongly convex over $X$ with parameter $\mu > 0$. Choosing $\gamma := \frac{1}{L}$, projected gradient descent with arbitrary $\mathbf{x}_0$ satisfies the following two properties.*

(i) *Squared distances to $\mathbf{x}^\star$ are geometrically decreasing:*

$$\|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2 \leq \left(1 - \frac{\mu}{L}\right) \|\mathbf{x}_t - \mathbf{x}^\star\|^2, \quad t \geq 0.$$

(ii) *The absolute error after $T$ iterations is exponentially small in $T$:*

$$\begin{aligned} f(\mathbf{x}_T) - f(\mathbf{x}^\star) \quad &\leq \quad \|\nabla f(\mathbf{x}^\star)\| \left(1 - \frac{\mu}{L}\right)^{T/2} \|\mathbf{x}_0 - \mathbf{x}^\star\| \quad \leftarrow \textit{in general, } \nabla f(\mathbf{x}^\star) \neq \mathbf{0}! \\ &\quad + \quad \frac{L}{2}\left(1 - \frac{\mu}{L}\right)^T \|\mathbf{x}_0 - \mathbf{x}^\star\|^2, \quad T > 0. \leftarrow \textit{as in unconstrained case} \end{aligned}$$

Proof.

(i) Geometric decrease plus noise: $\underline{\|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2} \leq \cdots$

▶ unconstrained case:

$$2\gamma(f(\mathbf{x}^\star) - f(\mathbf{x}_t)) + \gamma^2\|\nabla f(\mathbf{x}_t)\|^2 \qquad\qquad + \underline{(1 - \mu\gamma)\|\mathbf{x}_t - \mathbf{x}^\star\|^2}.$$

▶ constrained case (vanilla analysis with a saving):

$$2\gamma(f(\mathbf{x}^\star) - f(\mathbf{x}_t)) + \gamma^2\|\nabla f(\mathbf{x}_t)\|^2 - \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 + \underline{(1 - \mu\gamma)\|\mathbf{x}_t - \mathbf{x}^\star\|^2}.$$

## Proof.

To bound the noise, we use sufficient decrease.

▶ unconstrained case:

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L}\|\nabla f(\mathbf{x}_t)\|^2 \qquad , \quad t \geq 0.$$

▶ constrained case:

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L}\|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2}\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2, \quad t \geq 0.$$

Putting it together, the terms $\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2$ cancel, and we get

# Smooth and strongly convex functions over $X$: $\mathcal{O}(\log(1/\varepsilon))$ steps II

Proof.

To bound the noise, we use sufficient decrease.

▶ unconstrained case:

$$f(\mathbf{x}_{t+1}) \le f(\mathbf{x}_t) - \frac{1}{2L}\|\nabla f(\mathbf{x}_t)\|^2 \qquad\qquad , \quad t \ge 0.$$

▶ constrained case:

$$f(\mathbf{x}_{t+1}) \le f(\mathbf{x}_t) - \frac{1}{2L}\|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2}\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2, \quad t \ge 0.$$

Putting it together, the terms $\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2$ cancel, and we get

$$\|\mathbf{x}_{t+1} - \mathbf{x}^\star\|^2 \le (1 - \mu\gamma)\|\mathbf{x}_t - \mathbf{x}^\star\|^2 = \left(1 - \frac{\mu}{L}\right)\|\mathbf{x}_t - \mathbf{x}^\star\|^2.$$

in both cases. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ □

Proof.

(ii) Error bound from smoothness:

$$f(\mathbf{x}_T) - f(\mathbf{x}^\star) \;\; \leq \;\; \nabla f(\mathbf{x}^\star)^\top (\mathbf{x}_T - \mathbf{x}^\star) + \frac{L}{2} \|\mathbf{x}^\star - \mathbf{x}_T\|^2$$

Proof.

(ii) Error bound from smoothness:

$$
\begin{aligned}
f(\mathbf{x}_T) - f(\mathbf{x}^\star) &\leq \nabla f(\mathbf{x}^\star)^\top (\mathbf{x}_T - \mathbf{x}^\star) + \frac{L}{2}\|\mathbf{x}^\star - \mathbf{x}_T\|^2 \\
&\leq \|\nabla f(\mathbf{x}^\star)\| \, \|\mathbf{x}_T - \mathbf{x}^\star\| + \frac{L}{2}\|\mathbf{x}^\star - \mathbf{x}_T\|^2 \text{ (Cauchy-Schwarz)}
\end{aligned}
$$

Proof.

(ii) Error bound from smoothness:

$$
\begin{aligned}
f(\mathbf{x}_T) - f(\mathbf{x}^\star) \;\;&\leq\;\; \nabla f(\mathbf{x}^\star)^\top (\mathbf{x}_T - \mathbf{x}^\star) + \frac{L}{2} \|\mathbf{x}^\star - \mathbf{x}_T\|^2 \\[2mm]
&\leq\;\; \|\nabla f(\mathbf{x}^\star)\| \, \|\mathbf{x}_T - \mathbf{x}^\star\| + \frac{L}{2} \|\mathbf{x}^\star - \mathbf{x}_T\|^2 \;\; \text{(Cauchy-Schwarz)} \\[2mm]
&\leq\;\; \|\nabla f(\mathbf{x}^\star)\| \left(1 - \frac{\mu}{L}\right)^{T/2} \|\mathbf{x}_0 - \mathbf{x}^\star\| + \frac{L}{2} \left(1 - \frac{\mu}{L}\right)^{T} \|\mathbf{x}_0 - \mathbf{x}^\star\|^2 . \;\; \text{(i)}
\end{aligned}
$$

$\square$

constrained error bound $\approx \sqrt{\text{unconstrained error bound}}$

required number of steps roughly doubles.
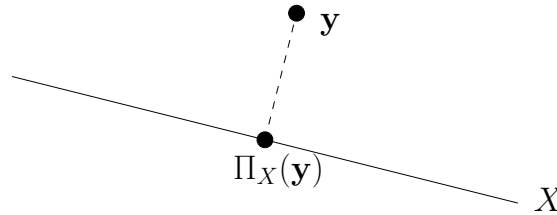
# The Projection Step: $\Pi_X(\mathbf{y}) := \operatorname{argmin}_{\mathbf{x} \in X} \|\mathbf{x} - \mathbf{y}\|$

Computing $\Pi_X(\mathbf{y})$ is an optimization problem itself.

Q1   When do we want to constrain?

$\quad\quad\quad\hookrightarrow$ bound $\quad \| x_0 - x^* \|_2 \leq R$

$\quad\quad\quad\hookrightarrow$ Regularizer in ML

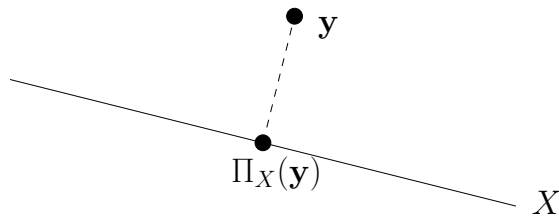Q2   Can we project onto constraint?

# The Projection Step: $\Pi_X(\mathbf{y}) := \operatorname{argmin}_{\mathbf{x} \in X} \|\mathbf{x} - \mathbf{y}\|$

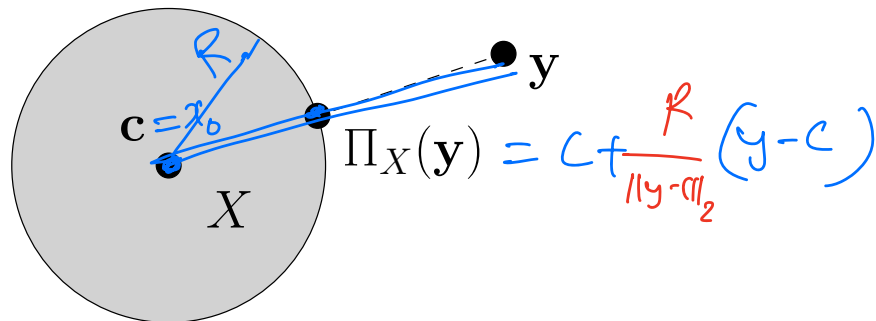Computing $\Pi_X(\mathbf{y})$ is an optimization problem itself.

It can efficiently be solved in relevant cases:

— $f$ is $\mu$-s.cvx, but non-smooth



$$\frac{1}{\mu} \|g(x)\|$$

— Weight decay,



$$\|x - x_0\|_2^2 \leq R$$

— $x$ is sparse $\Rightarrow$

$$\|x\|_1 \leq R$$

— $X$ is a matrix and low rank

$$\Rightarrow \quad Tr(X) \leq R$$

$X$ is symm psd

$$Tr(X) = \sum_{i \in [n]} \lambda_i (X)$$



$$Tr(X^T X)$$

# The Projection Step: $\Pi_X(\mathbf{y}) := \arg\min_{\mathbf{x} \in X} \|\mathbf{x} - \mathbf{y}\|$

Computing $\Pi_X(\mathbf{y})$ is an optimization problem itself.

It can efficiently be solved in relevant cases:

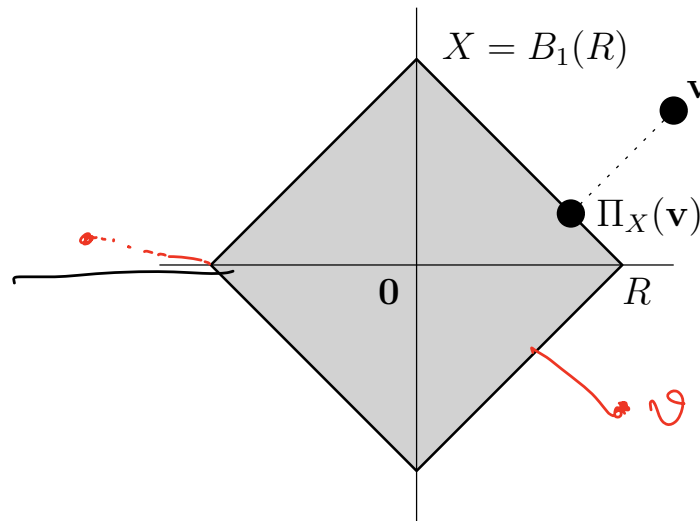▶ Projecting onto an affine subspace (leads to system of linear equations, similar to least squares)

# The Projection Step: $\Pi_X(\mathbf{y}) := \arg\min_{\mathbf{x} \in X} \|\mathbf{x} - \mathbf{y}\|$

Computing $\Pi_X(\mathbf{y})$ is an optimization problem itself.

It can efficiently be solved in relevant cases:

▶ Projecting onto an affine subspace (leads to system of linear equations, similar to least squares)



▶ Projecting onto a Euclidean ball with center $\mathbf{c}$ (simply scale the vector $\mathbf{y} - \mathbf{c}$)



$$\Pi_X(\mathbf{y}) = C + \frac{R}{\|y - c\|_2}(y - c)$$

# Projecting onto $\ell_1$-balls (needed in Lasso)

W.l.o.g. restrict to center at $\mathbf{0}$: $B_1(R) = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_1 = \sum_{i=1}^{d} |x_i| \leq R\}$.
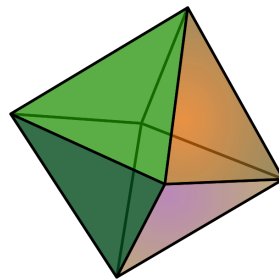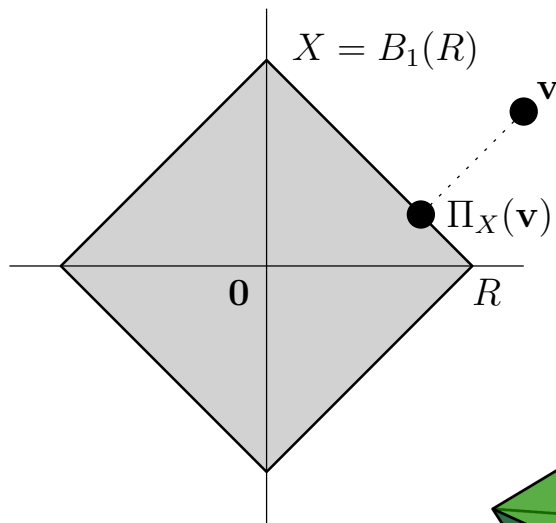


$\|x\|_1 \leq R$

$(\pm 1)^d$

$\underset{\lambda}{\text{argmin}} \, \|\mathcal{V} - \lambda \mathcal{I}\|_1 \leq R$

# Projecting onto $\ell_1$-balls (needed in Lasso)

W.l.o.g. restrict to center at $\mathbf{0}$: $B_1(R) = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_1 = \sum_{i=1}^{d} |x_i| \leq R\}$.
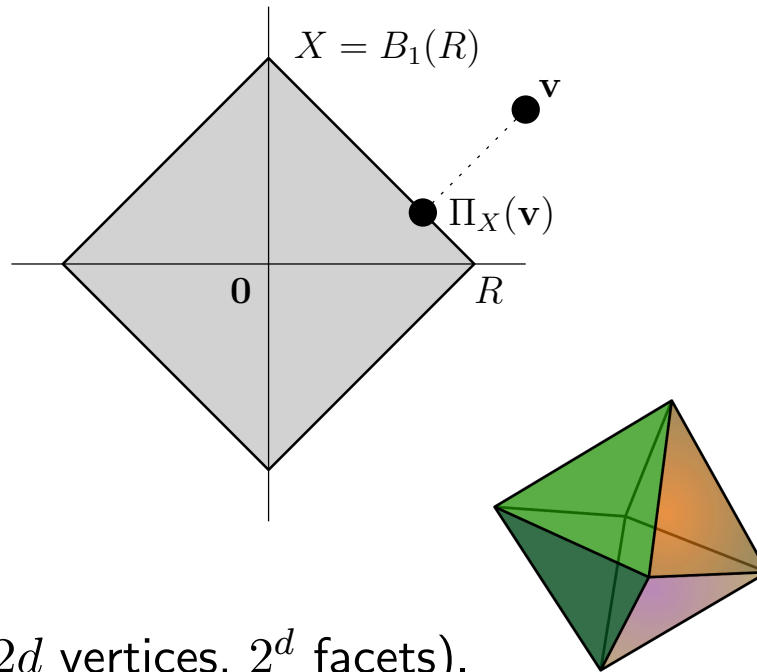


$B_1(R)$ is the cross polytope ($2d$ vertices, $2^d$ facets). (octahedron, $d = 3$)

# Projecting onto $\ell_1$-balls (needed in Lasso)

W.l.o.g. restrict to center at $\mathbf{0}$: $B_1(R) = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_1 = \sum_{i=1}^{d} |x_i| \leq R\}$.



$B_1(R)$ is the cross polytope ($2d$ vertices, $2^d$ facets).            (octahedron, $d = 3$)

Section **??**: projection can be computed in $\mathcal{O}(d \log d)$ time (can be improved to $\mathcal{O}(d)$)

## Ex 1

$f$ : $L$ smooth + cvx, we don't know $L$.

$$x_{t+1} = x_t - \frac{1}{L} \nabla f(x_t)$$

### Case 1

guess $\hat{L} \ll L$

$$\hat{x}_{t+1} = x_t - \frac{1}{\hat{L}} \nabla f(x_t)$$

$\nabla f$

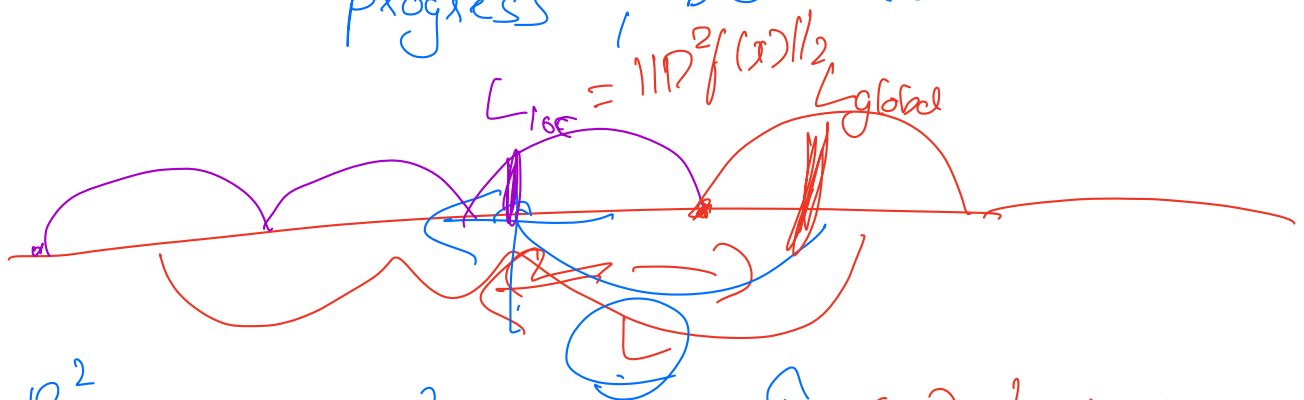$$f(\hat{x}_{t+1}) \leq f(x_t) - \frac{1}{2\hat{L}} \|\nabla f(x_t)\|_2^2$$

$$> \implies \hat{L} \text{ is too large}$$

we need to increase $\hat{L}$

### Case 2

guess $\hat{L} \gg L$

progress, but could do better



$L_{loc} = \|D^2 f(x)\|_2$   $L_{global}$

$$\frac{L R^2}{T}, \quad \frac{2 L R^2}{T} \qquad \hat{L} < 2 L_{global}$$

$$L_{local} \leq L \qquad \log \frac{L_y a_\varepsilon}{10^{-6}}$$

$$L \geq L_{global} \Rightarrow$$

$$f(x_{f+1}) \leq f(x_f) - \frac{1}{2L} \|\mathrm{D}f(x_f)\|_2^2$$



adaptive
algorithms

"Adam"