# ORF522 – Linear and Nonlinear Optimization

## 17. Operator theory I

**Bartolomeo Stellato — Fall 2022**

# Today's lecture
## [Chapter 4, FMO][PA][PMO][LSMO]

**Operator theory I**

- Proximal gradient method

- Operators

- Monotone and cocoercive operators

- Fixed-point Iterations

# Proximal gradient method

# Remember: gradient descent interpretation

**Problem**

$$\text{minimize} \quad f(x)$$

**Iterations**

$$x^{k+1} = x^k - t\nabla f(x^k)$$

**Quadratic approximation**, replacing Hessian $\nabla^2 f(x^k)$ with $\dfrac{1}{t}I$

$$x^{k+1} = \operatorname*{argmin}_z f(x^k) + \nabla f(x^k)^T(z - x^k) + \frac{1}{2t}\|z - x^k\|_2^2$$

# Let's exploit the smooth part

minimize $\quad f(x) + g(x)$

$f(x)$ convex and smooth
$g(x)$ convex (may be not differentiable)

Quadratic approximation of $f$ while keeping $g$

$$x^{k+1} = \operatorname*{argmin}_{z} g(z) + \boxed{f(x^k) + \nabla f(x^k)^T(z - x^k) + \frac{1}{2t}\|z - x^k\|_2^2}$$

$\longleftarrow$ same as gradient descent

**Equivalent to**

$$x^{k+1} = \operatorname*{argmin}_{z} tg(z) + \frac{1}{2}\left\|z - (x^k - t\nabla f(x^k))\right\|_2^2 = \mathbf{prox}_{tg}\left(x^k - t\nabla f(x^k)\right)$$

**Proximal operator**

$\uparrow$
make $g$ small

$\uparrow$
stay close to gradient update

# Proximal gradient method

minimize $\quad f(x) + g(x)$

$f(x)$ convex and smooth
$g(x)$ convex (may be not differentiable)

**Iterations**

$$x^{k+1} = \mathbf{prox}_{tg}\left(x^k - t\nabla f(x^k)\right)$$

**Properties**

- Alternates between gradient updates of $f$ and proximal updates on $g$
- Useful if $\mathbf{prox}_{tg}$ is inespensive
- Can handle nonsmooth and constrained problems

# Special cases

## Generalized gradient descent

minimize $\quad f(x) + g(x)$

**Iterations**

$$x^{k+1} = \mathbf{prox}_{tg}\left(x^k - t\nabla f(x^k)\right)$$

**Smooth**

$$g(x) = 0 \implies \mathbf{prox}_{tg}(x) = x$$

**Gradient descent**

$$\implies \quad x^{k+1} = x^k - t\nabla f(x^k)$$

**Constraints**

$$g(x) = \mathcal{I}_C(x) \implies \mathbf{prox}_{tg}(x) = \Pi_C(x)$$

**Projected gradient descent**

$$\implies \quad x^{k+1} = \Pi_C(x^k - t\nabla f(x^k))$$

**Non smooth**

$$f(x) = 0$$

**Proximal minimization**

$$\implies \quad x^{k+1} = \mathbf{prox}_{tg}(x^k)$$

*Note:* useful if $\mathbf{prox}_{tg}$ is cheap

7

# What happens if we cannot evaluate the prox?

At every iteration, it can be very expensive to evaluate

$$\mathbf{prox}_g(x) = \operatorname*{argmin}_z \left( g(z) + \frac{1}{2}\|z - x\|_2^2 \right)$$

**Idea: solve it approximately!**

If you precisely control the $\mathbf{prox}_g(x)$ evaluation errors
you can obtain the same convergence guarantees (and rates)
as the exact evaluations.

[Schmidt et al. (2011), "Convergence rates of inexact proximal-gradient methods for convex optimization"]

# Example: Lasso

## Iterative Soft Thresholding Algorithm (ISTA)

$$\text{minimize} \quad \underbrace{(1/2)\|Ax - b\|_2^2}_{f(x)} + \underbrace{\lambda\|x\|_1}_{g(x)}$$

**Proximal gradient descent**

$$\nabla f(x) = A^T(Ax - b)$$

$$x^{k+1} = \mathbf{prox}_{tg}\left(x^k - t\nabla f(x^k)\right)$$

$$\mathbf{prox}_{tg}(x) = S_{\lambda t}(x) \quad \text{(component wise soft-thresholding)}$$

**Closed-form iterations**

$$x^{k+1} = S_{\lambda t}\left(x^k - tA^T(Ax^k - b)\right)$$
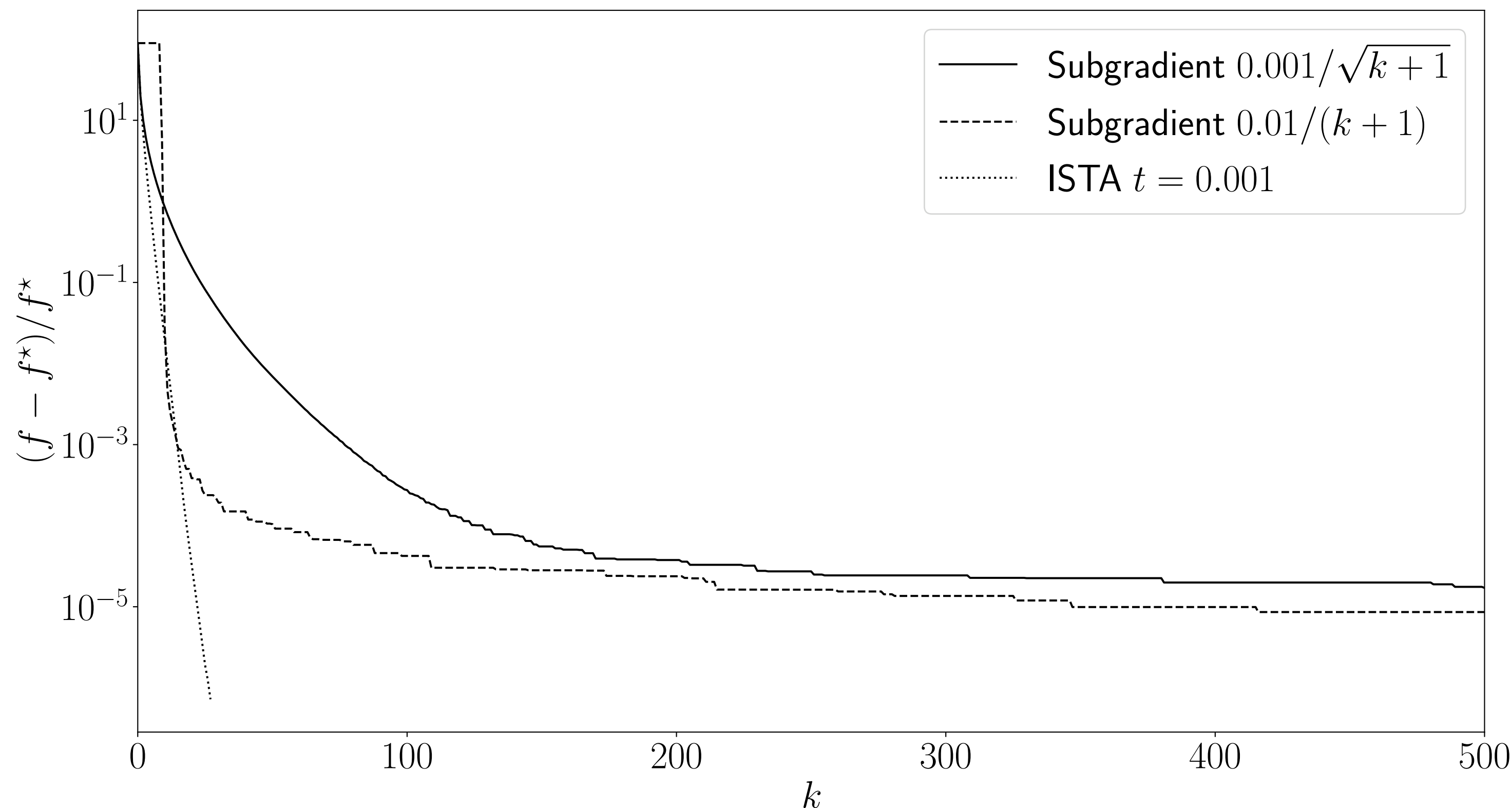
# Example: Lasso
## Iterative Soft Thresholding Algorithm (ISTA)

$A \in \mathbf{R}^{500 \times 100}$

$$\text{minimize} \quad (1/2)\|Ax - b\|_2^2 + \lambda\|x\|_1$$

**Closed-form iterations**

$$x^{k+1} = S_{\lambda t}\left(x^k - tA^T(Ax^k - b)\right)$$



**Better convergence**

**Can we prove convergence generally?**

**Can we combine different operators?**
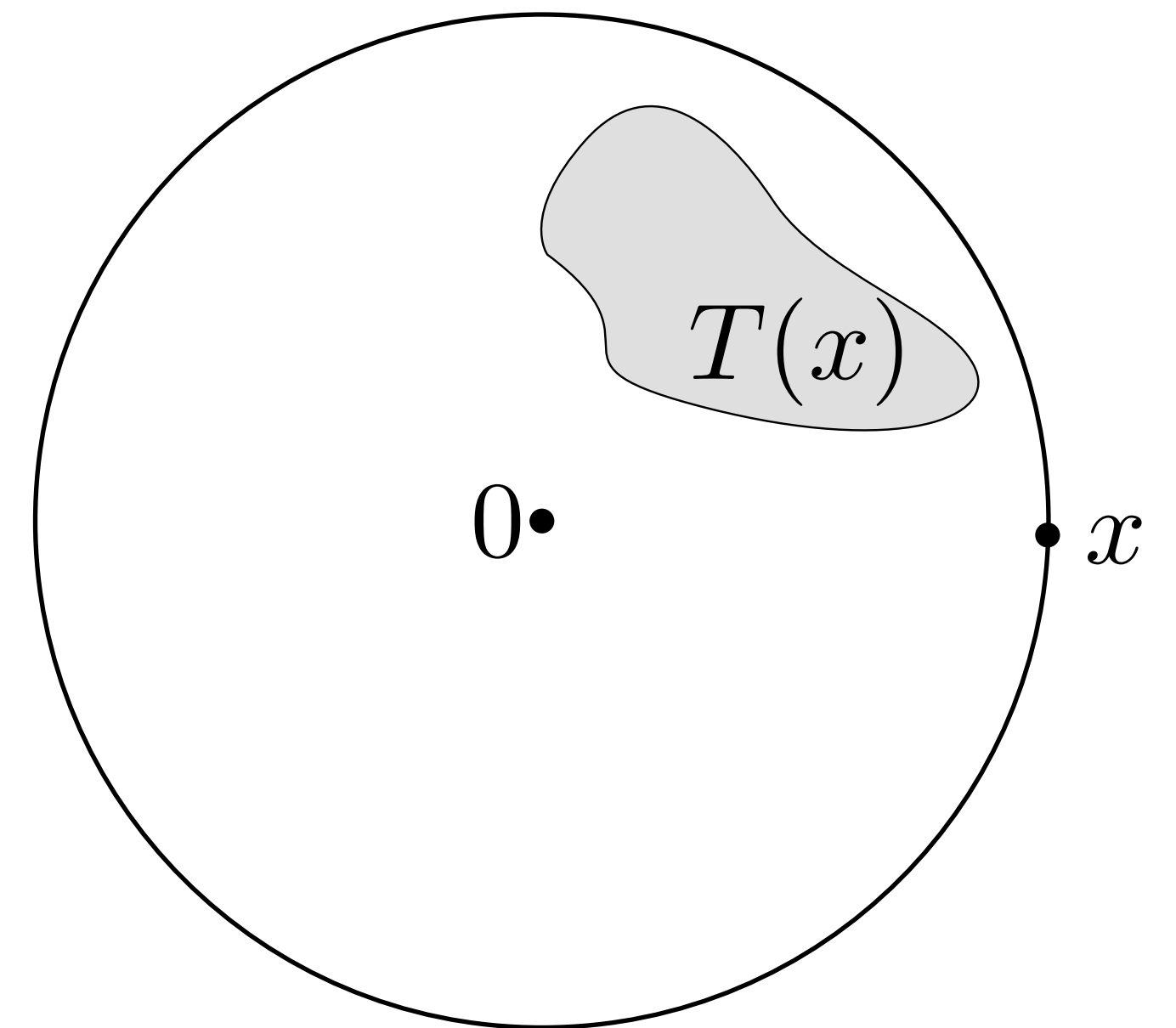
# Operators

# Operators

An operator $T$ maps each point in $\mathbf{R}^n$ to a subset of $\mathbf{R}^n$

- **set valued** $T(x)$ returns a set
- **single-valued** $T(x)$ (function) returns a singleton

The **domain** of $T$ is the set $\mathbf{dom}\, T = \{x \mid T(x) \neq \emptyset\}$

**Example**

- The subdifferential $\partial f$ is a set-valued operator
- The gradient $\nabla f$ is a single-valued operator

# Graph and inverse operators

**Graph**

The graph of an operator $T$ is defined as

$$\mathbf{gph}T = \{(x, y) \mid y \in T(x)\}$$

In other words, all the pairs of points $(x, y)$ such that $y \in T(x)$.

**Inverse**

The graph of the inverse operator $T^{-1}$ is defined as

$$\mathbf{gph}T^{-1} = \{(y, x) \mid (x, y) \in \mathbf{gph}T\}$$

Therefore, $y \in T(x)$ if and only if $x \in T^{-1}(y)$.

# Zeros

## Zero

$x$ is a **zero** of $T$ if $\qquad 0 \in T(x)$

## Zero set

The set of all the zeros $\qquad T^{-1}(0) = \{x \mid 0 \in T(x)\}$

**Example**
If $T = \partial f$ and $f : \mathbf{R}^n \to \mathbf{R}$, then
$0 \in T(x)$ means that $x$ minimizes $f$

Many problems
can be posed as finding zeros
of an operator

# Fixed points

$\bar{x}$ is a **fixed-point** of a single-valued operator $T$ if

$$\bar{x} = T(\bar{x})$$

**Set of fixed points**   $\mathbf{fix}\, T = \{x \in \mathbf{dom}\, T \mid x = T(x)\} = (I - T)^{-1}(0)$

**Examples**
- **Identity** $T(x) = x$. Any point is a fixed point
- **Zero operator** $T(x) = 0$. Only $0$ is a fixed point

# Lipschitz operators

An operator $T$ is $L$-Lipschitz if

$$\|T(x) - T(y)\| \leq L\|x - y\|, \quad \forall x, y \in \mathbf{dom}\, T$$

**Fact** If $T$ is Lipschitz, then it is single-valued

**Proof** If $y = T(x), z = T(x)$, then $\|y - z\| \leq L\|x - x\| = 0 \implies y = z$ ∎

For $L = 1$ we say $T$ is **nonexpansive**
For $L < 1$ we say $T$ is **contractive** (with contraction factor $L$)

# Lipschitz operators examples

**Lipschitz affine functions**

$$T(x) = Ax + b$$

$\longleftrightarrow$

maximum singular value

$$L = \|A\|_2 = \sqrt{\lambda_{\max}(A^T A)}$$

**Lipschitz differentiable functions**

$T$ such that there exists derivative $DT$

$\longleftrightarrow$

derivative is bounded

$$\|DT\|_2 \leq L$$

# Lipschitz operators and fixed points

Given a $L$-Lipschitz operator $T$ and a fixed point $\bar{x} = T\bar{x}$,

$$\|Tx - \bar{x}\| = \|Tx - T\bar{x}\| \leq L\|x - \bar{x}\|$$

A contractive operator ($L < 1$) can have at most
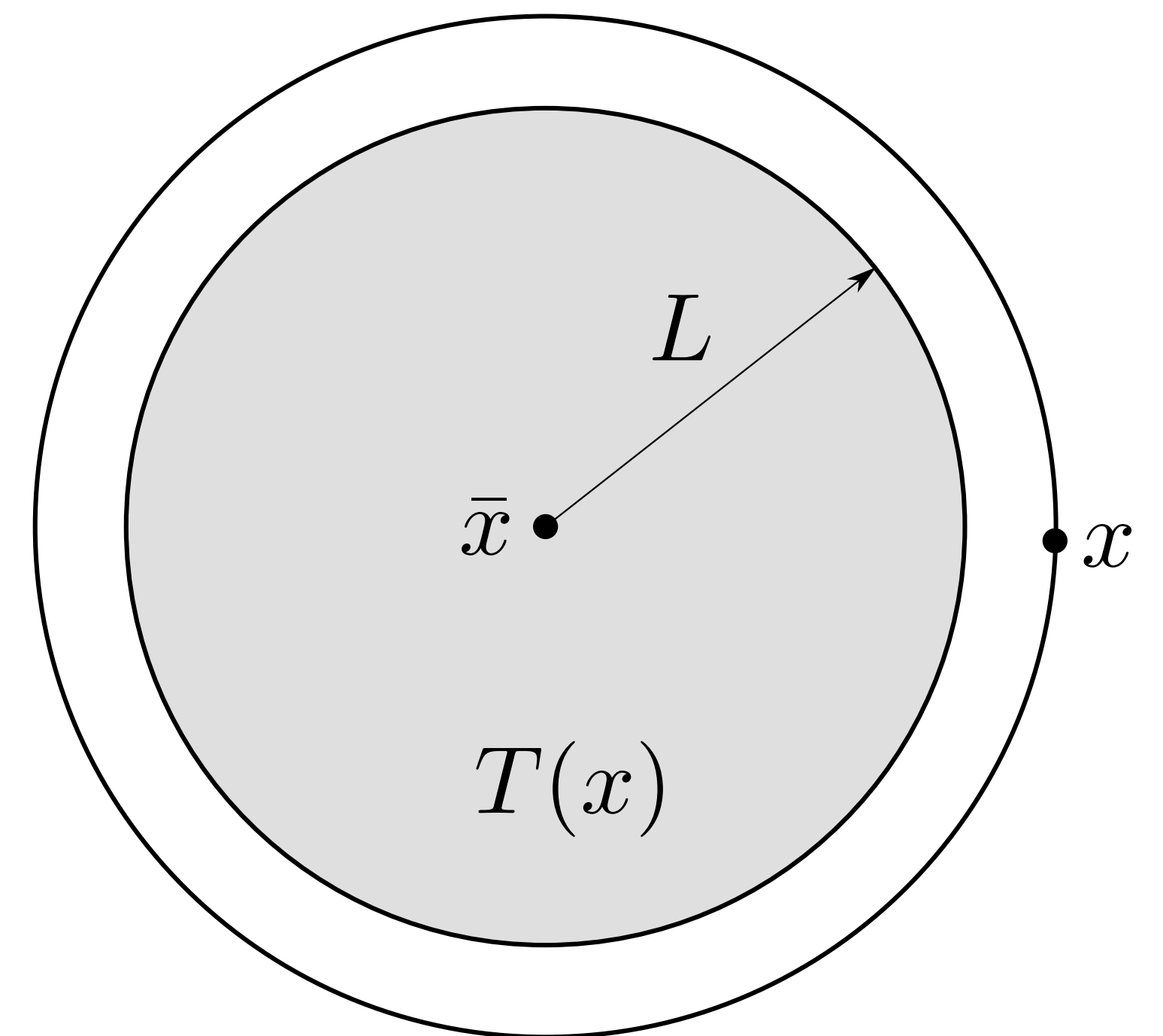one fixed point, i.e., $\mathbf{fix}\, T = \{\bar{x}\}$

**Proof**
If $\bar{x}, \bar{y} \in \mathbf{fix}\, T$ and $\bar{x} \neq \bar{y}$ then
$\|\bar{x} - \bar{y}\| = \|T(\bar{x}) - T(\bar{y})\| < \|\bar{x} - \bar{y}\|$ (contradiction) ■

A nonexpansive operator ($L = 1$) need not
have a fixed point

**Example** $T(x) = x + 2$

# Combining Lipschitz operators

$$T_1 \text{ is } L_1\text{-Lipschitz and } T_2 \text{ is } L_2\text{-Lipschitz}$$

The **composition** $T_1 T_2$ is $L_1 L_2$-Lipschitz

**Proof** $\|T_1 T_2 x - T_1 T_2 y\|_2 \leq L_1 \|T_2 x - T_2 y\|_2 \leq L_1 L_2 \|x - y\|_2$ ■

- Composition of *nonexpansive* is nonexpansive
- Composition of *nonexpansive* and *contractive* is contractive

The **weighted average** $\theta T_1 + (1 - \theta) T_2, \ \theta \in (0, 1)$ is $(\theta L_1 + (1 - \theta) L_2)$-Lipschitz

**Proof** (exercise)

- Weighted average of *nonexpansive* is nonexpansive
- Weighted average of *nonexpansive* and *contractive* is contractive

# Monotone cocoercive operators

# Monotone operators

An operator $T$ on $\mathbf{R}^n$ is **monotone** if

$$(u - v)^T (x - y) \geq 0, \quad \forall (x, u), (y, v) \in \mathbf{gph}\, T$$
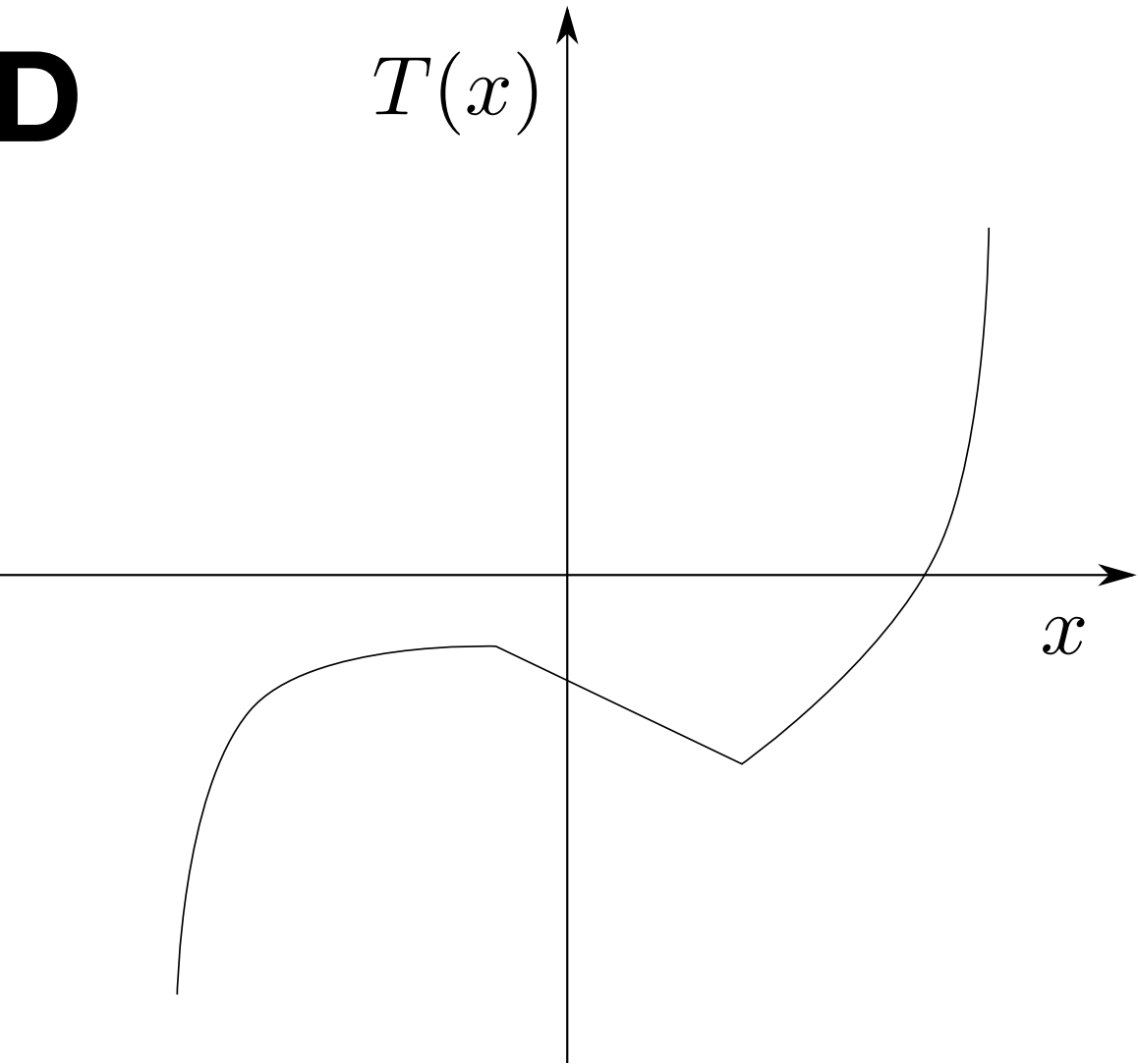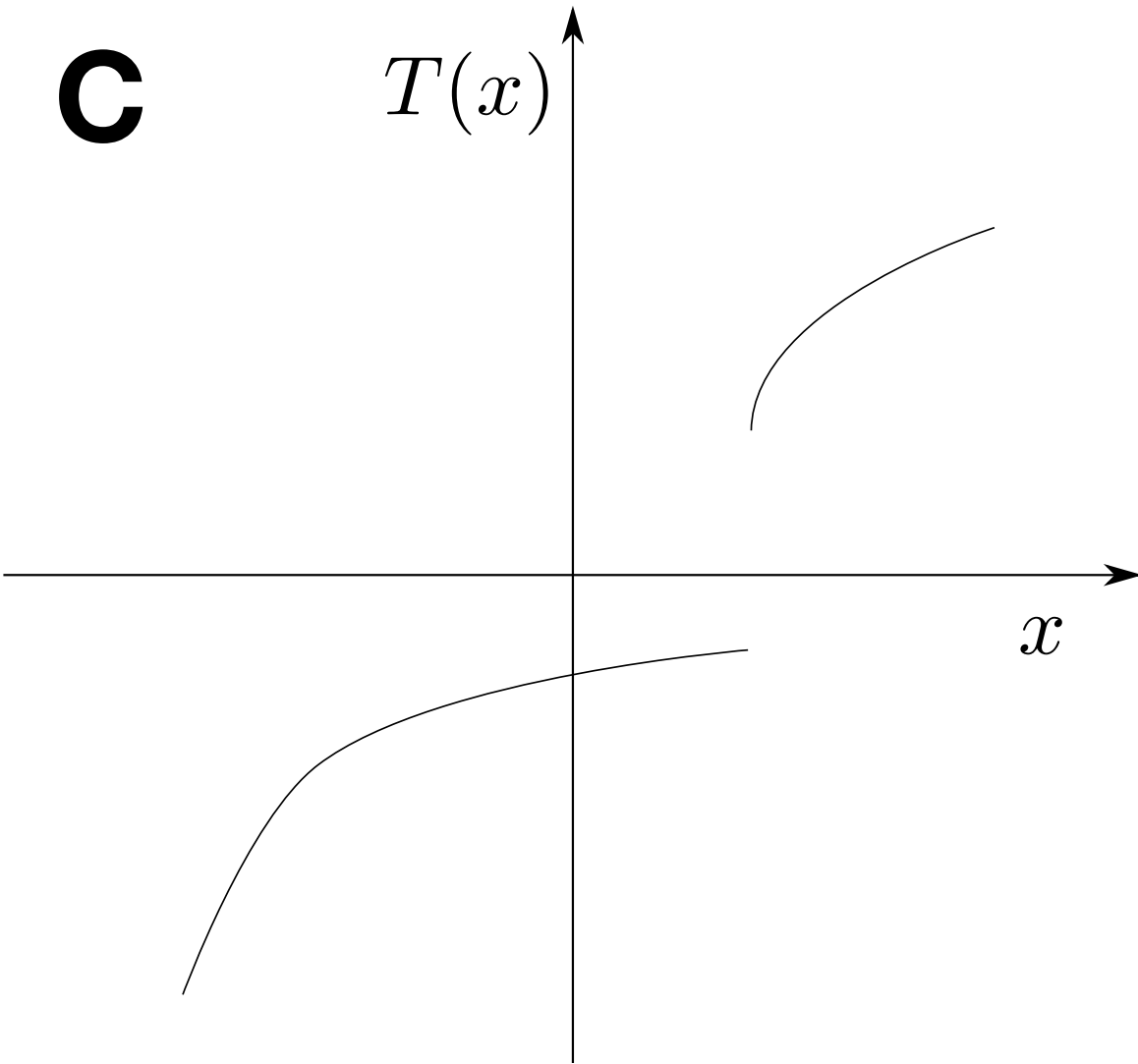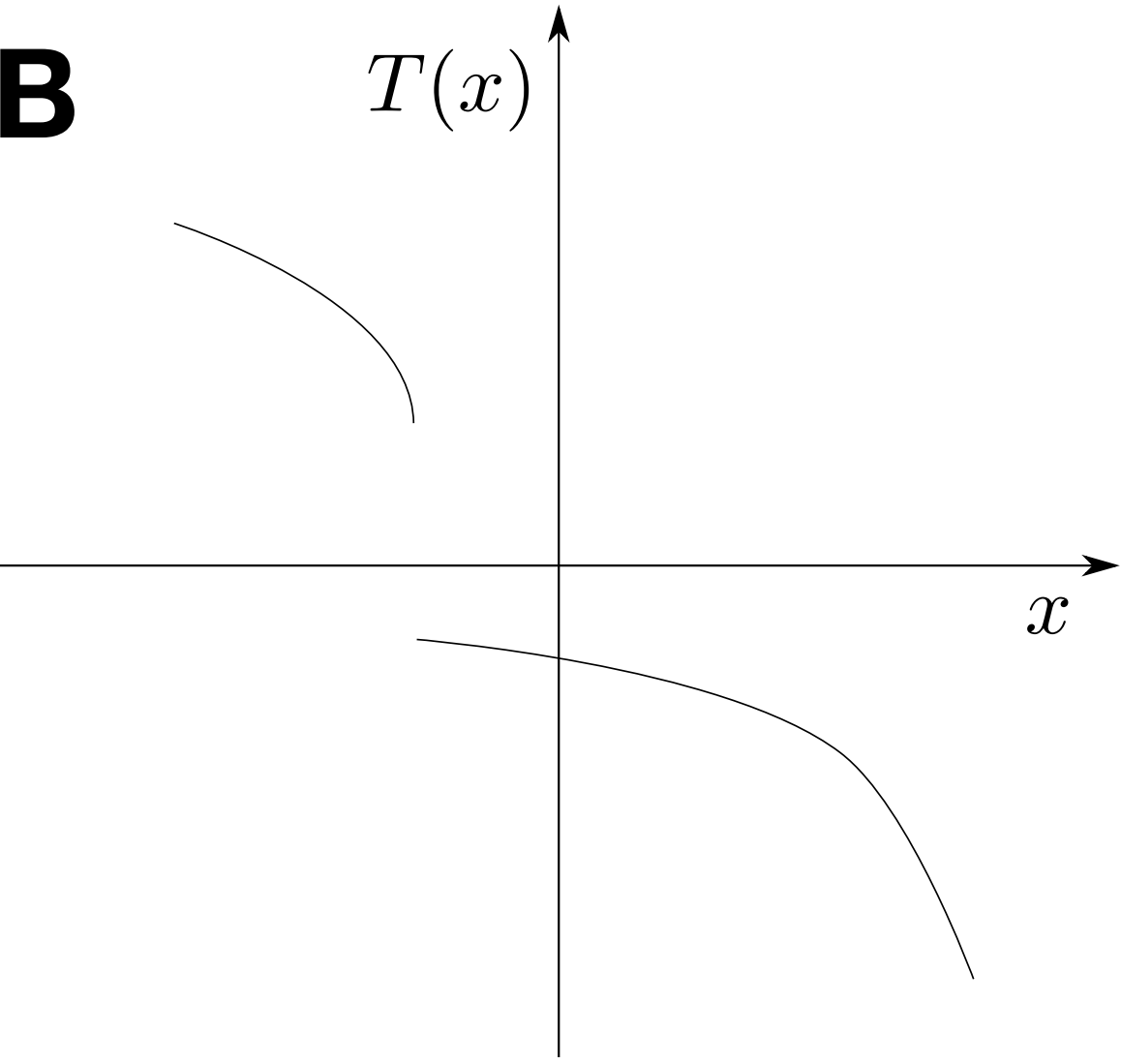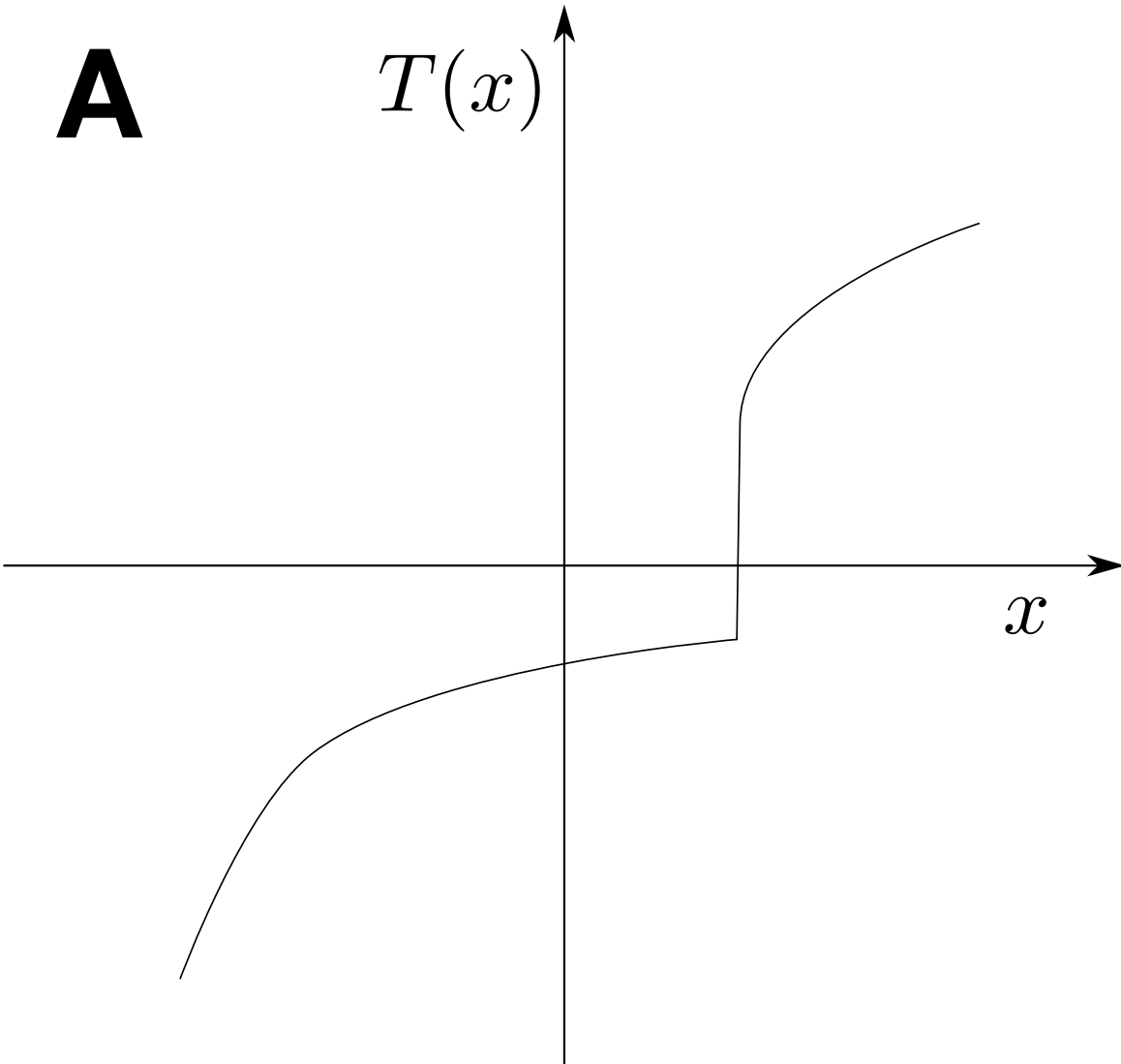
$T$ is **maximal monotone** if
$\nexists (\bar{x}, \bar{u}) \notin \mathbf{gph}\, T$ such that

$$(\bar{u} - u)^T (\bar{x} - x) \geq 0$$

*Equivalently:* $\nexists$ monotone $R$
such that $\mathbf{gph}\, T \subset \mathbf{gph}\, R$

# Monotone operators in 1D

**A** $T(x)$

**B** $T(x)$

$x$

$x$

**C** $T(x)$

**D** $T(x)$

$x$

$x$

|   | Monotone | Max Monotone |
|---|----------|--------------|
| **A** |  |  |
| **B** |  |  |
| **C** |  |  |
| **D** |  |  |

**Monotonicity**

$$y > x \quad \Rightarrow \quad T(y) \geq T(x)$$

**Continuity**

If $T$ single-valued, continuous and monotone, then it's maximal monotone [22]

# Monotone operator properties

- **sum** $T + R$ is monotone

- **nonnegative scaling** $\alpha T$ with $\alpha \geq 0$ is monotone

- **inverse** $T^{-1}$ is monotone

- **congruence** for $M \in \mathbf{R}^{n \times m}$, then $M^T T(Mz)$ is monotone on $\mathbf{R}^m$

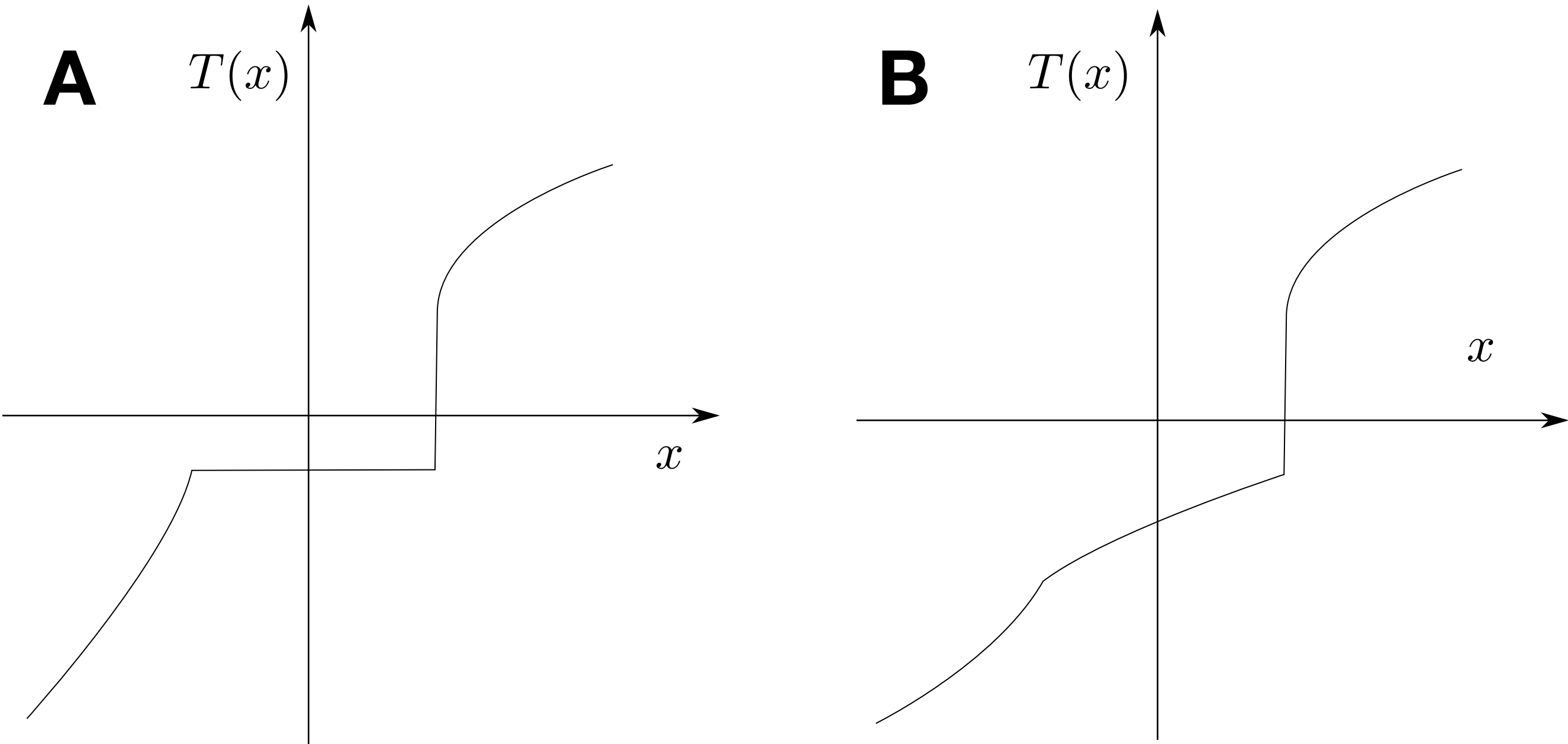**Affine function** $T(x) = Ax + b$ is maximal monotone
$$\iff \quad A + A^T \succeq 0$$

# Strongly monotone operators

An operator $T$ on $\mathbf{R}^n$ is $\mu$-**strongly monotone** if

$$(u - v)^T (x - y) \geq \mu \|x - y\|^2, \quad \mu > 0 \qquad \text{(also called } \mu\text{-\textbf{coercive})}$$

$$\forall (x, u), (y, v) \in \mathbf{gph}T$$

**Let's fill the table**

A

B

| | Monotone | Strongly Monotone |
|---|---|---|
| A | | |
| B | | |

The slope is at least $\mu$

# Cocoercive operators

An operator $T$ is $\beta$-**cocoercive**, $\beta > 0$, if

$$(T(x) - T(y))^T (x - y) \geq \beta \|T(x) - T(y)\|^2$$

If $T$ is $\beta$-**cocoercive**, then $T$ is $(1/\beta)$-**Lipschitz**

**Proof** $\beta\|T(x) - T(y)\|^2 \leq (T(x) - T(y))^T (x - y) \leq \|T(x) - T(y)\|\|x - y\|$

$$\implies \|T(x) - T(y)\| \leq (1/\beta)\|x - y\| \qquad \blacksquare$$

If $T$ is $\mu$-**strongly monotone** if and only if $T^{-1}$ is $\mu$-**cocoercive**

**Proof** $(T(x) - T(x))^T (x - y) \geq \mu\|x - y\|^2$

Inverse: $u = T(x)$ and $v = T(y)$ if and only if $x \in T^{-1}(u)$ and $y \in T^{-1}(v)$

$$(u - v)^T (T^{-1}(u) - T^{-1}(v)) \geq \mu\|T^{-1}(u) - T^{-1}(v)\|^2 \qquad \blacksquare$$

# Cocoercive and nonexpansive operators

If $T$ is $\beta$-**cocoercive**   if and only if   $I - 2\beta T$ is **nonexpansive**

**Proof**   $\|(I-2\beta T)(y) - (I - 2\beta T)(x)\|^2 =$

$$= \|y - 2\beta T(y) - x - 2\beta T(x)\|^2$$

$$= \|y - x\|^2 - 4\beta(T(y) - T(x))^T(y - x) + 4\beta^2\|T(y) - T(x)\|^2$$

$$= |y - x\|^2 - 4\beta\left((T(y) - T(x))^T(y - x) - \beta\|T(y) - T(x)\|^2\right)$$

$$\leq \|y - x\|^2 \qquad \blacksquare \qquad\qquad \text{(cocoercive)}$$

# Summary of monotone and cocoercive operators

**Monotone**

$$(T(x) - T(y))^T (x - y) \geq 0$$

**Lipschitz**

$$\|F(x) - F(y)\| \leq L\|x - y\|$$

$$\mu = 0 \uparrow$$

$$L = 1/\mu \uparrow$$

**Strongly monotone**

$$(T(x) - T(y))^T (x - y) \geq \mu\|x - y\|^2 \quad \longleftrightarrow \quad (F(x) - F(y))^T (x - y) \geq \mu\|F(x) - F(y)\|^2$$

$$F = T^{-1}$$

**Cocoercive**

$$G = I - 2\mu F \updownarrow$$

**Nonexpansive**

$$\|G(x) - G(y)\| \leq \|x - y\|$$

# Fixed point iterations

# Fixed point iteration

**Apply operator**

$$x^{k+1} = T(x^k)$$

until you reach $\bar{x} \in \mathbf{fix}\, T$

**Main approach**

1. Find a suitable $T$ such that $\bar{x} \in \mathbf{fix}\, T$ solve your problem
2. Show that the fixed point iteration converges

**Fixed point residual to terminate**
$$r^k = T(x^k) - x^k$$

# Contractive fixed point iterations

**Contraction mapping theorem**
If $T$ is $L$-Lipschitz with $L < 1$ (contraction), the iteration

$$x^{k+1} = T(x^k)$$

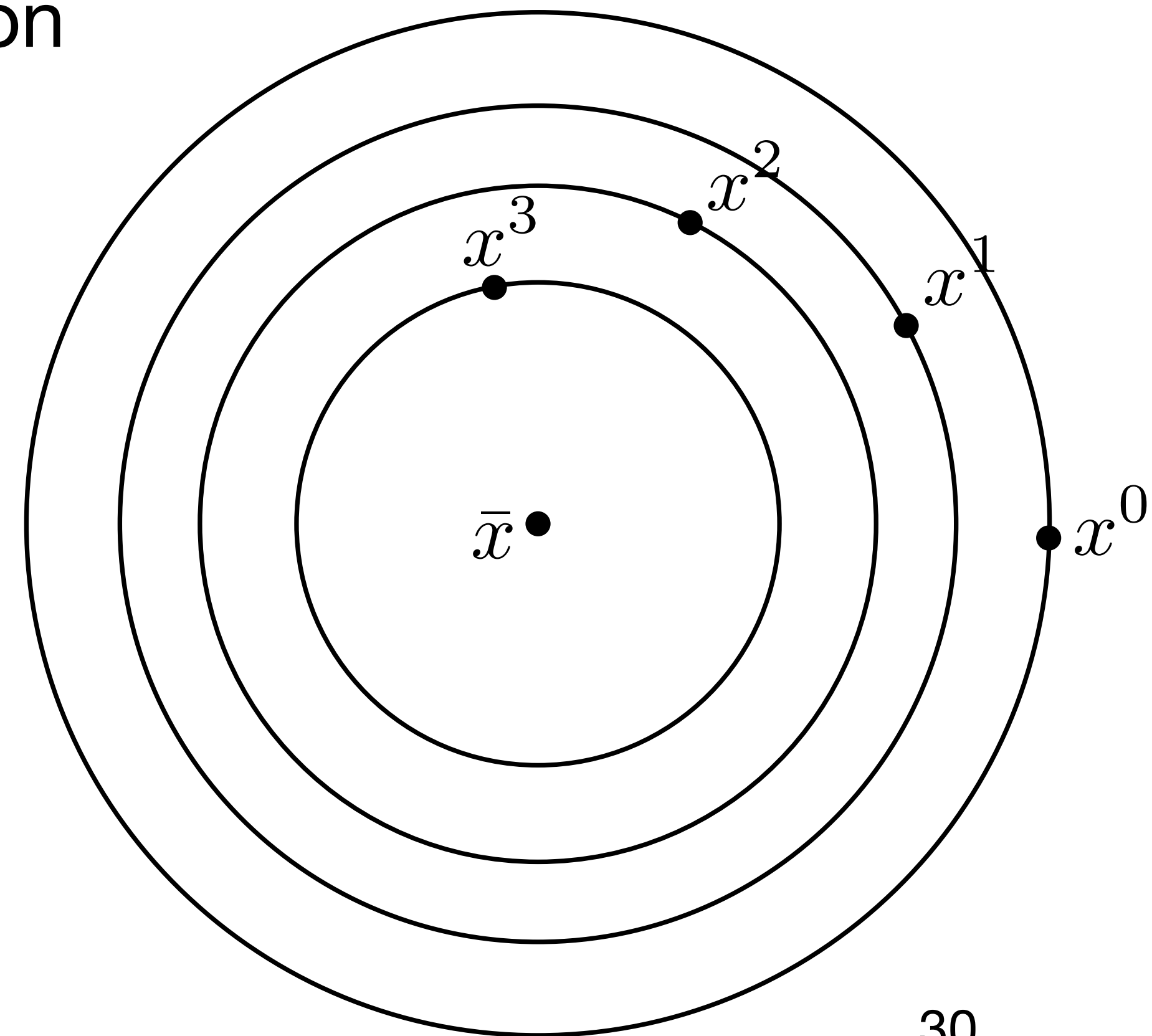converges to $\bar{x}$, the unique fixed point of $T$

**Properties**

- Distance to $\bar{x}$ decreases at each step

$$\|x^{k+1} - \bar{x}\| \leq L\|x^k - \bar{x}\|$$

(iteration is **Fejer monotone**)

- Linear convergence rate $L$

# Contraction mapping theorem

**Proof**

The sequence $x^k$ is Cauchy

$$\|x^{k+\ell} - x^k\| \leq \|x^{k+\ell} - x^{k+\ell-1}\| + \cdots + \|x^{k+1} - x^k\|$$

$$\leq (L^{\ell-1} + \cdots + 1)\|x^{k+1} - x^k\| \quad \text{(Lipschitz constant)}$$

$$\leq \frac{1}{1-L}\|x^{k+1} - x^k\| \quad \text{(geometric series)}$$

$$\leq \frac{L^k}{1-L}\|x^1 - x^0\| \quad \text{(Lipschitz constant)}$$

Therefore it converges to a point $\bar{x}$ which must be the (unique) fixed point of $T$

The convergence is linear (geometric) with rate $L$

$$\|x^k - \bar{x}\| = \|T(x^{k-1}) - T(\bar{x})\| \leq L\|x^{k-1} - \bar{x}\| \leq L^k\|x^0 - x^\star\| \quad \blacksquare \quad 31$$
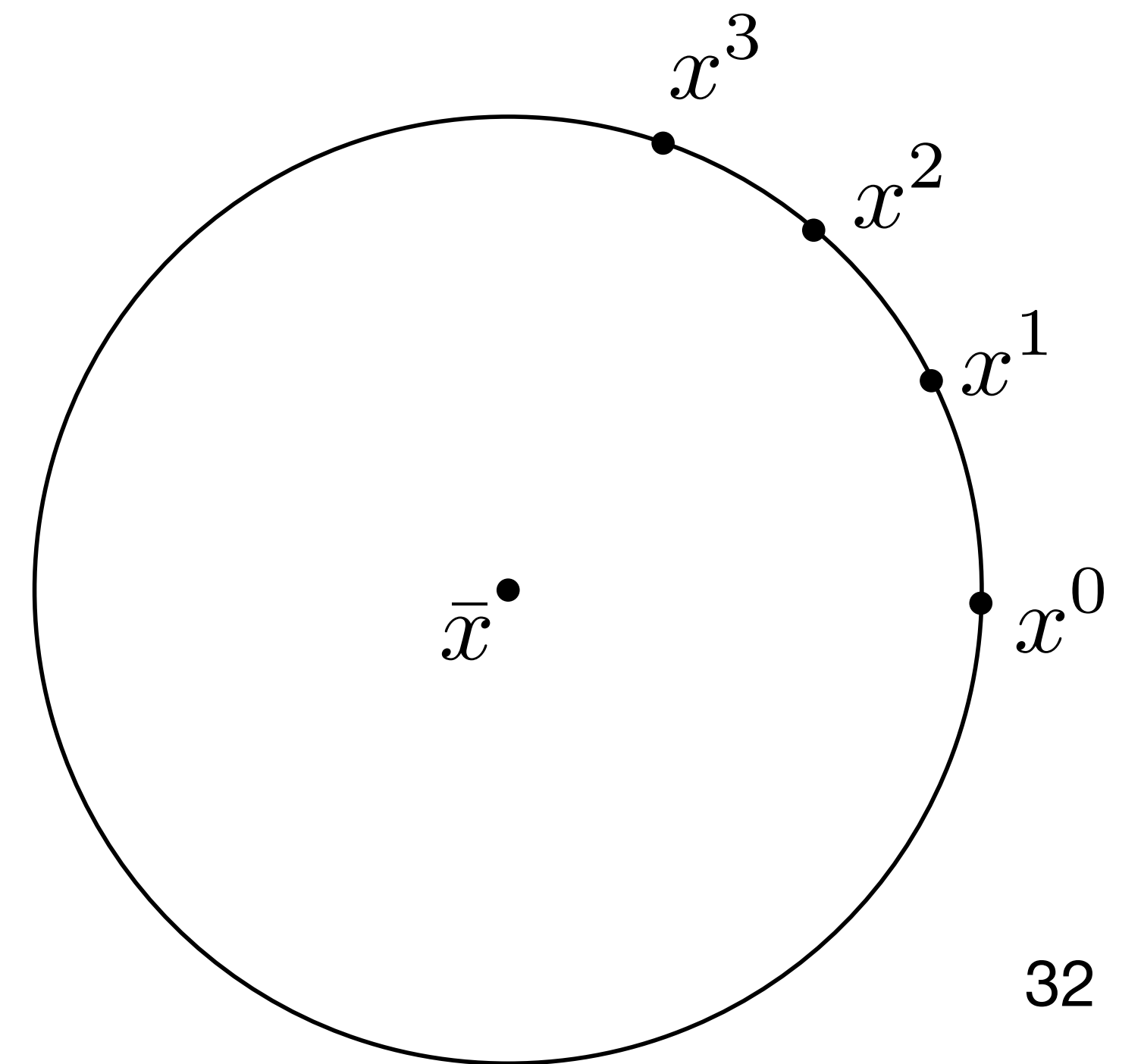
# Nonexpansive fixed point iterations

If $T$ is $L$-Lipschitz with $L = 1$ (nonexpansive), the iteration

$$x^{k+1} = T(x^k)$$

need not converge to a fixed point, even if one exists.

**Example**
- Let $T$ be a rotation around the origin
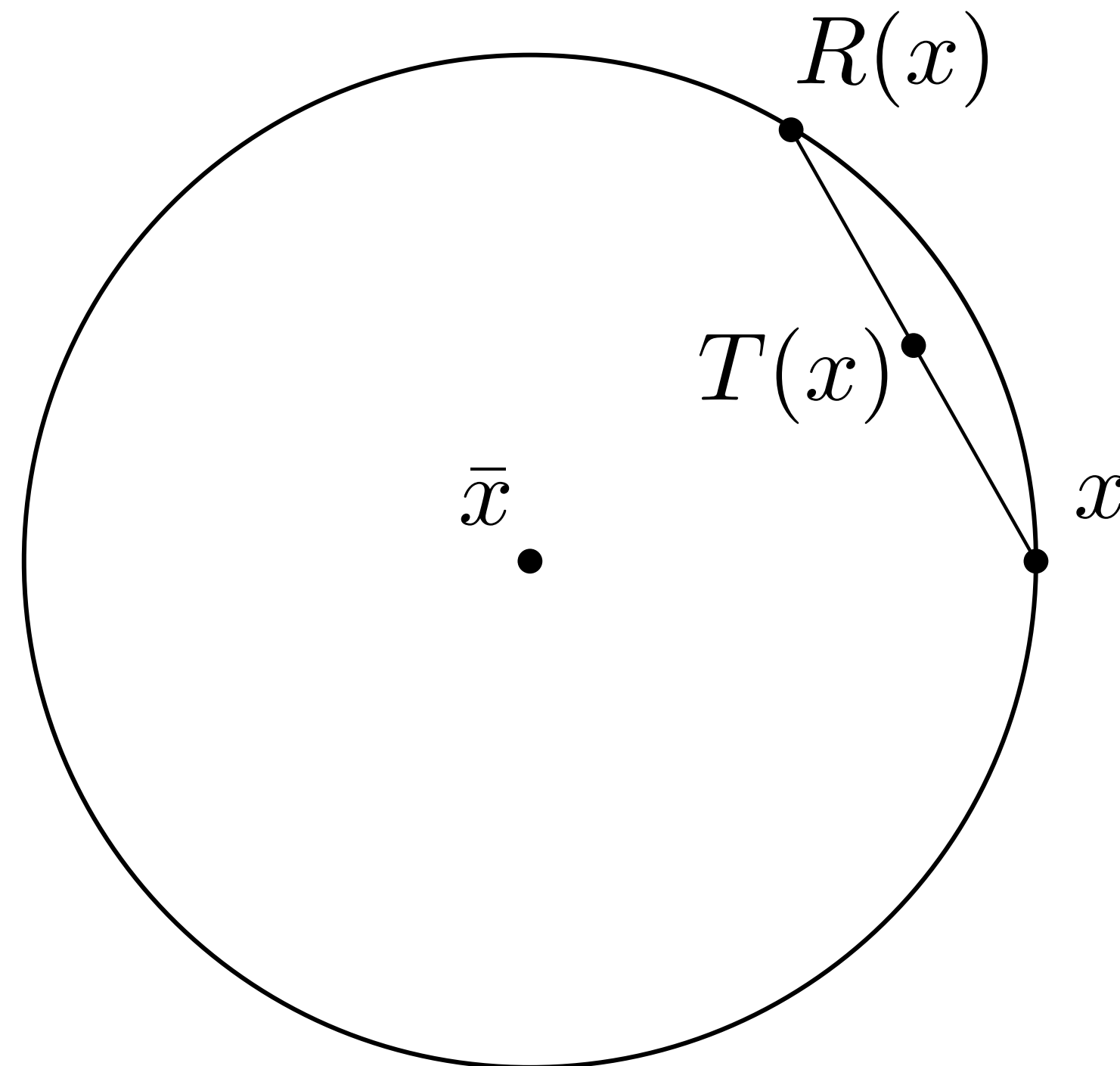- $T$ is nonexpansive and has a fixed point $\bar{x} = 0$
- $\|x^k\|$ never decreases

# Averaged operators

We say that an operator $T$ is $\alpha-$**averaged** with $\alpha \in (0,1)$ if

$$T = (1 - \alpha)I + \alpha R$$

and $R$ is nonexpansive.

# Averaged operators fixed points

We say that an operator $T$ is $\alpha-$**averaged** with $\alpha \in (0,1)$ if

$$T = (1-\alpha)I + \alpha R$$

**Fact** If $T$ is $\alpha$-averaged, then $\mathbf{fix}\, T = \mathbf{fix}\, R$

**Proof** $\quad \bar{x} = T(\bar{x}) = (1-\alpha)I(\bar{x}) + \alpha R(\bar{x})$

$$= (1-\alpha)\bar{x} + \alpha R(\bar{x})$$

$$\Longleftrightarrow \quad \alpha\bar{x} = \alpha R(\bar{x})$$

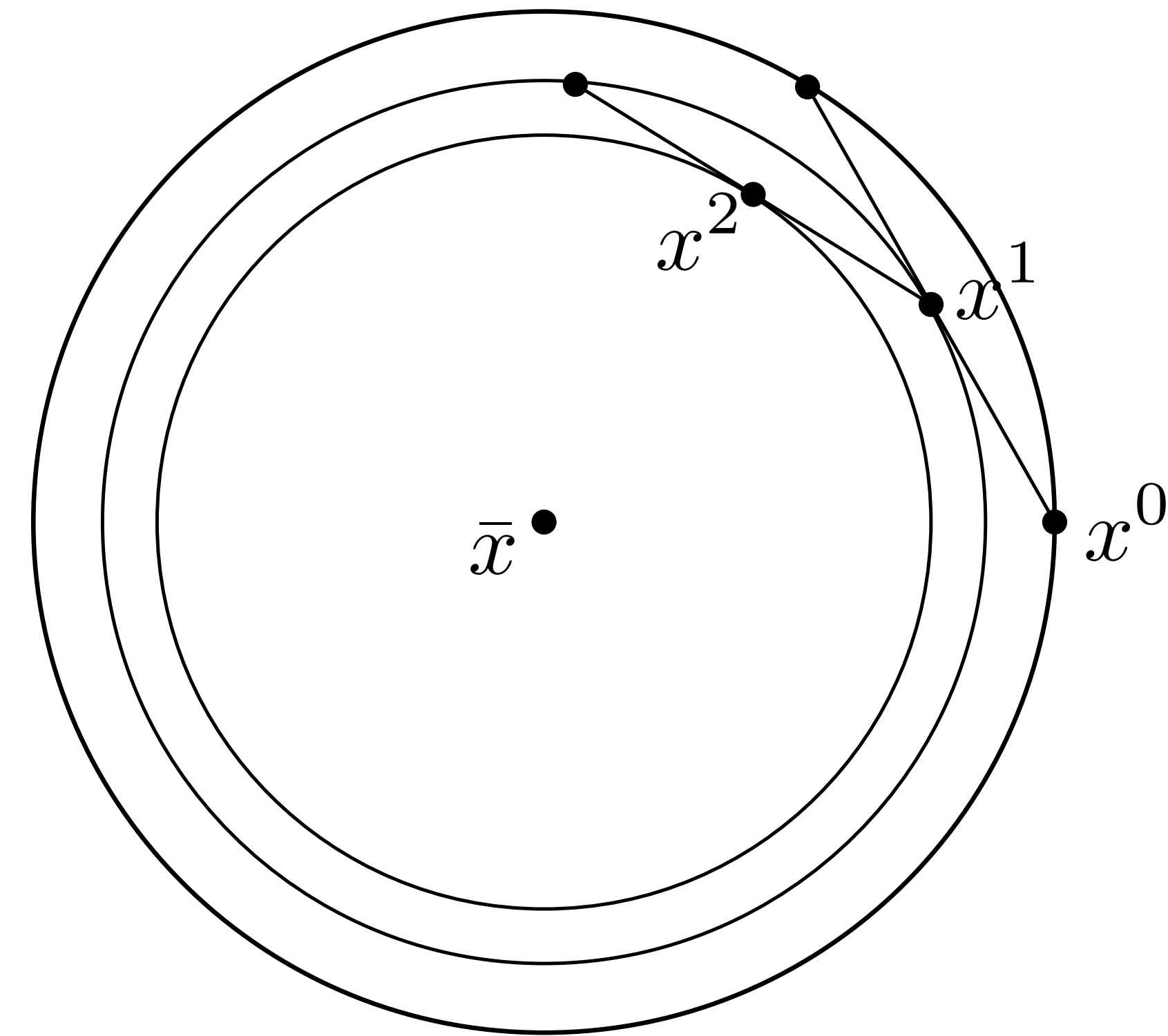$$\Longleftrightarrow \quad \bar{x} = R(\bar{x}) \qquad \blacksquare$$

# Averaged fixed point iterations

If $T = (1 - \alpha)I + \alpha R$ is $\alpha$-averaged
($\alpha \in (0,1)$ and $R$ nonexpansive), the iteration

$$x^{k+1} = T(x^k)$$

converges to $\bar{x} \in \mathbf{fix}\,T$

       (also called damped, averaged
       or Mann-Krasnosel'skii iteration)

**Properties**
- Distance to $\bar{x}$ decreases at each step (**Fejer monotone**)
- Sublinear convergence to fixed-point residual

$$\|R(x^k) - x^k\| \leq \frac{1}{\sqrt{(k+1)\alpha(1-\alpha)}} \|x^0 - \bar{x}\|$$

# Averaged fixed point iterations

**Proof**

Use the identity (proof by expanding)

$$\|(1-\alpha)a + \alpha b\|^2 = (1-\alpha)\|a\|^2 + \alpha\|b\|^2 - \alpha(1-\alpha)\|a-b\|^2$$

and apply it to

$$x^{k+1} - \bar{x} = (1-\alpha)\underbrace{(x^k - \bar{x})}_{a} + \alpha\underbrace{(R(x^k) - \bar{x})}_{b}$$

obtaining

$$\|x^{k+1} - \bar{x}\|^2 = (1-\alpha)\|x^k - \bar{x}\|^2 + \alpha\|R(x^k) - \bar{x}\|^2 - \alpha(1-\alpha)\|x^k - R(x^k)\|^2$$

$$\leq (1-\alpha)\|x^k - \bar{x}\|^2 + \alpha\|x^k - \bar{x}\|^2 - \alpha(1-\alpha)\|x^k - R(x^k)\|^2 \text{ (nonexpansive)}$$

$$= \|x^k - \bar{x}\|^2 \underbrace{- \alpha(1-\alpha)\|x^k - R(x^k)\|^2}_{\leq 0}$$

Iterations are Fejer monotone

36

# Averaged fixed point iterations

**Proof (continued)**

iterate righthand side over $k$ steps

$$\|x^{k+1} - \bar{x}\|^2 \le \|x^0 - \bar{x}\|^2 - \alpha(1-\alpha) \sum_{i=0}^{k} \|x^i - R(x^i)\|^2$$

Since $\|x^{k+1} - \bar{x}\|^2 \ge 0$, we have $\quad \sum_{i=0}^{k} \|x^i - R(x^i)\|^2 \le \dfrac{1}{\alpha(1-\alpha)} \|x^0 - \bar{x}\|^2$

Using $\displaystyle\sum_{i=0}^{k} \|x^i - R(x^i)\|^2 \ge (k+1) \min_{i=0,\ldots,k} \|x^i - R(x^i)\|^2$, we obtain

$$\min_{i=0,\ldots,k} \|x^i - R(x^i)\|^2 \le \frac{1}{(k+1)\alpha(1-\alpha)} \|x^0 - \bar{x}\|^2$$

$(R$ is nonexpansive $\to \min$ at $k)$ $\quad \|x^k - R(x^k)\|^2 \le \dfrac{1}{(k+1)\alpha(1-\alpha)} \|x^0 - \bar{x}\|^2$ ∎ 37

# Average fixed point iteration convergence rates

$$\|R(x^k) - x^k\| \leq \frac{1}{\sqrt{(k+1)\alpha(1-\alpha)}} \|x^0 - \bar{x}\|$$

Righthand side minimized when $\alpha = 1/2$

**Iterations**

$$\|R(x^k) - x^k\| \leq \frac{2}{\sqrt{k+1}} \|x^0 - \bar{x}\|$$

$$x^{k+1} = (1/2)x^k + (1/2)R(x^k)$$

**Remarks**
- Sublinear convergence (same as subgrad method), in general not the actual rate
- $\alpha = 1/2$ is very common for averaged operators

# How to design an algorithm

## Problem

$$\text{minimize} \quad f(x)$$

## Algorithm (operator) construction

1. Find a suitable $T$ such that $\bar{x} \in \mathbf{fix}\, T$ solve your problem
2. Show that the fixed point iteration converges

If $T$ is contractive $\implies$ **linear convergence**

If $T$ is averaged $\implies$ **sublinear convergence**

Most first order algorithms can be constructed in this way

# Operator theory

Today, we learned to:

- **Define** and **evaluate** proximal operators for various common functions

- **Apply** proximal operators to generalize gradient descent (vanilla, projected, proximal)

- **Define** monotone and cocoercive operators and their relations

- **Use operator theory** to construct general fixed-point iterations and prove their convergence

# Next lecture

- Operators in optimization algorithms