

Local convergence of adaptive tensor methods



Mathematical
Institute

KARL WELZEL

Mathematical Institute
University of Oxford

NA Internal Seminar, 21 November 2024



Oxford
Mathematics



Local convergence of adaptive tensor methods

Work in progress with



Raphael Hauser



Yang Liu



Coralia Cartis

1 Motivation

- Why higher-order methods?
- The AR p method

2 Numerical Illustrations

- Example with non-degenerate minimizer
- Example with degenerate minimizer

3 Theoretical Results

4 Conclusion

1 Motivation

- Why higher-order methods?
- The AR p method

2 Numerical Illustrations

- Example with non-degenerate minimizer
- Example with degenerate minimizer

3 Theoretical Results

4 Conclusion

Why higher-order methods?

- Goal: Find a local minimizer \mathbf{x}_* of a smooth function $f: \mathbb{R}^n \rightarrow \mathbb{R}$
- We compare rates of convergence when \mathbf{x}_k is close to \mathbf{x}_* (local convergence)
- Newton's method converges quadratically when $\nabla^2 f(\mathbf{x}_*) \succ \mathbf{0}$
- Newton's method converges only linearly for $f(x) = x^4$ for example
- Access to higher derivatives \rightarrow superlinear convergence for singular $\nabla^2 f(\mathbf{x}_*)$

The AR p method

Algorithm 1.1: Adaptive regularization algorithm using up to p th derivatives

Parameters: $\sigma_0 > 0$, $0 < \eta < 1$, $0 < \gamma_1 \leq 1 < \gamma_2$

```
1 for  $k = 0, 1, \dots$  do
2   Compute objective function and derivatives  $f(\mathbf{x}_k), \nabla f(\mathbf{x}_k), \dots, \nabla^p f(\mathbf{x}_k)$ 
3   Construct local Taylor expansion as  $t_k(\mathbf{y}) = \sum_{j=0}^p \frac{1}{j!} \nabla^j f(\mathbf{x}_k) [\mathbf{y} - \mathbf{x}_k]^j$ 
4   Construct local model as  $m_k(\mathbf{y}) = t_k(\mathbf{y}) + \sigma_k \|\mathbf{y} - \mathbf{x}_k\|^{p+1}$ 
5   Find a local minimizer  $\mathbf{y}_k \in \mathbb{R}^n$  of  $m_k$  that satisfies  $m_k(\mathbf{y}_k) < m_k(\mathbf{x}_k)$ 
6   if  $f(\mathbf{x}_k) - f(\mathbf{y}_k) \geq \eta(t_k(\mathbf{x}_k) - t_k(\mathbf{y}_k))$  then
7     Set  $\mathbf{x}_{k+1} = \mathbf{y}_k$  and  $\sigma_{k+1} = \gamma_1 \sigma_k$  // successful iteration
8   else
9     Set  $\mathbf{x}_{k+1} = \mathbf{x}_k$  and  $\sigma_{k+1} = \gamma_2 \sigma_k$  // unsuccessful iteration
10  end
11 end
```

The AR p method

- Assume $\nabla^p f$ is Lipschitz continuous: $\|\nabla^p f(\mathbf{x}) - \nabla^p f(\mathbf{y})\| \leq L_p \|\mathbf{x} - \mathbf{y}\|$
- “Ideal” regularization is $\sigma = \frac{L_p}{(p+1)!}$ but L_p is unknown
- Adaptive method is less dependent on user choice σ_0

1 Motivation

- Why higher-order methods?
- The AR p method

2 Numerical Illustrations

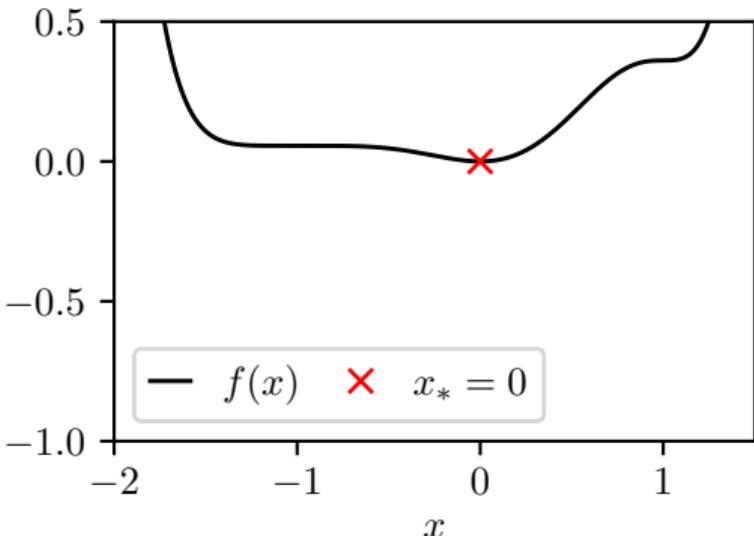
- Example with non-degenerate minimizer
- Example with degenerate minimizer

3 Theoretical Results

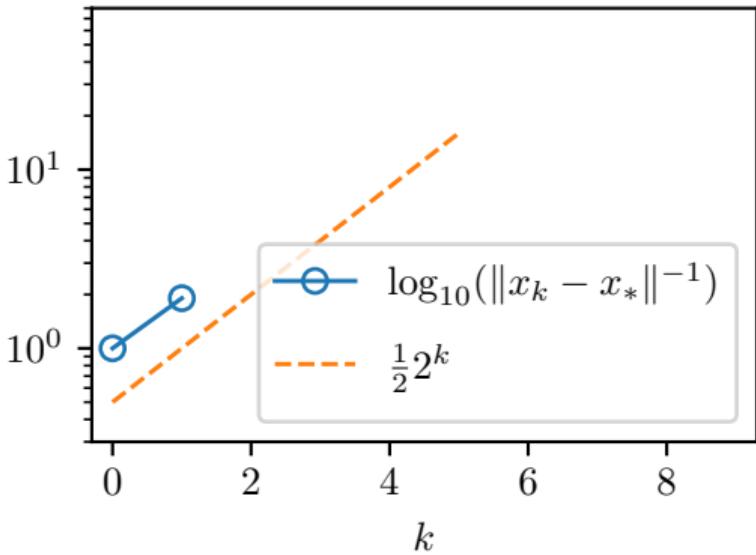
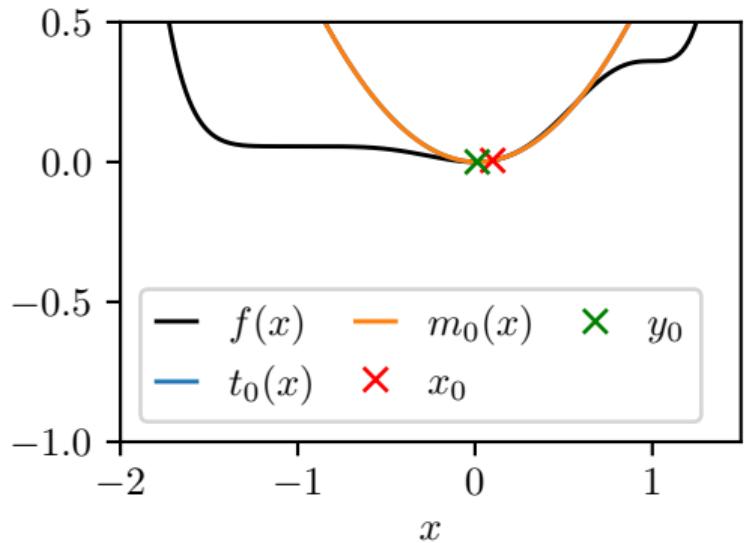
4 Conclusion

Example with non-degenerate minimizer

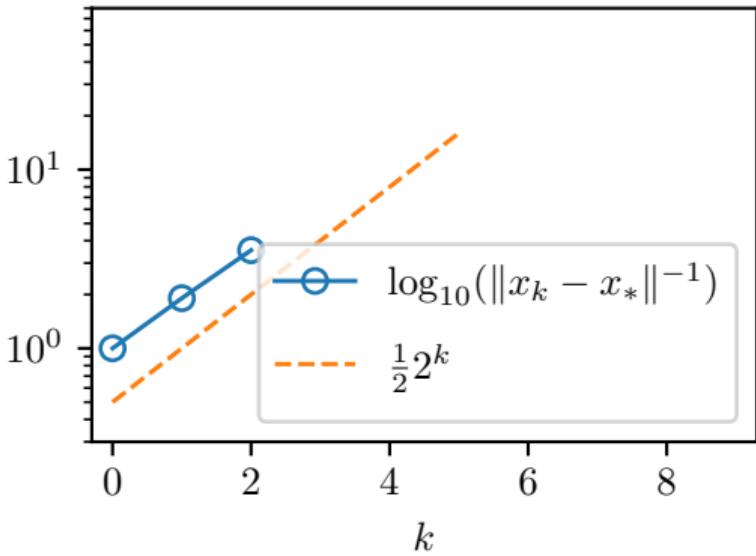
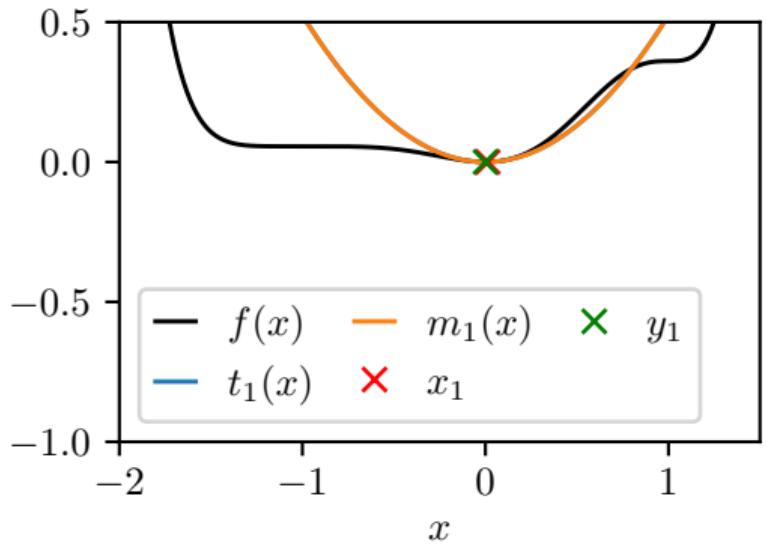
- $f(x) = \int_0^x (t+1)^4(t-1)^2 t \, dt$
- f has one local (and global) minimizer at $x_* = 0$
- f is not globally convex, but locally strongly convex: $f''(x_*) = 1$



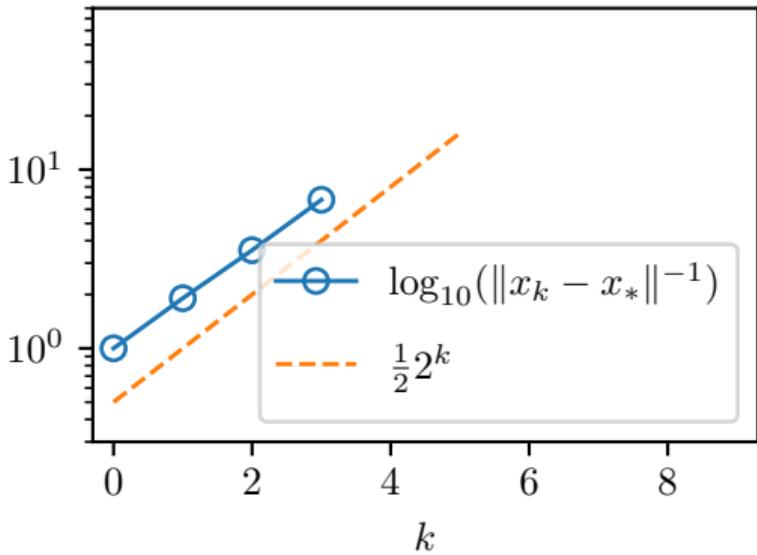
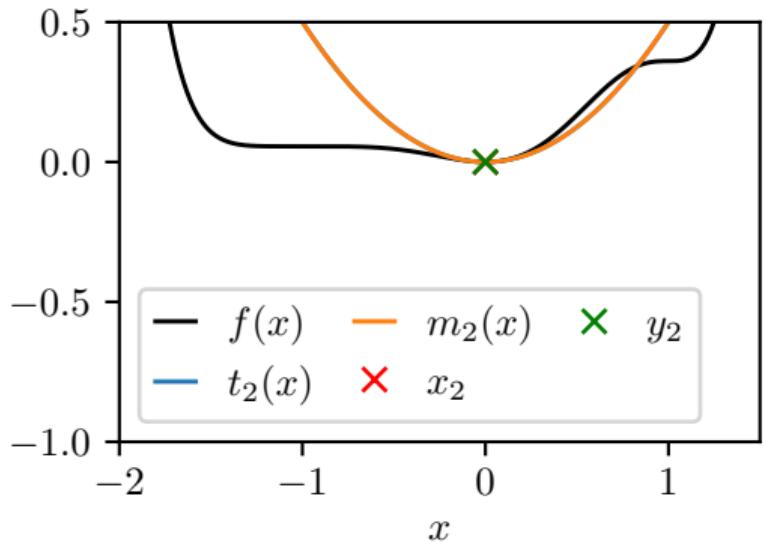
Newton's method ($p = 2$, $\sigma = 0$)



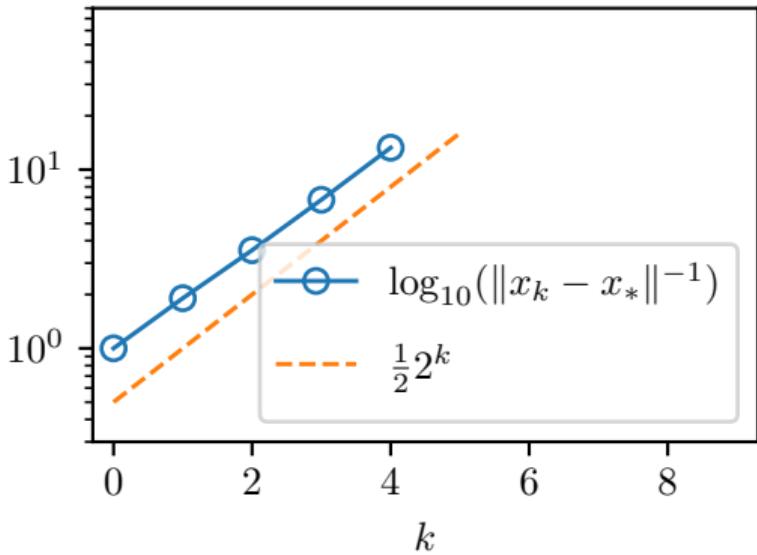
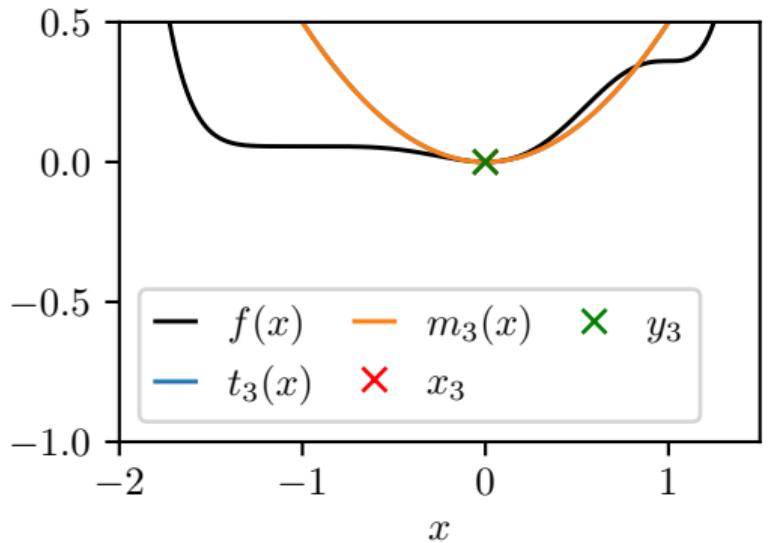
Newton's method ($p = 2$, $\sigma = 0$)



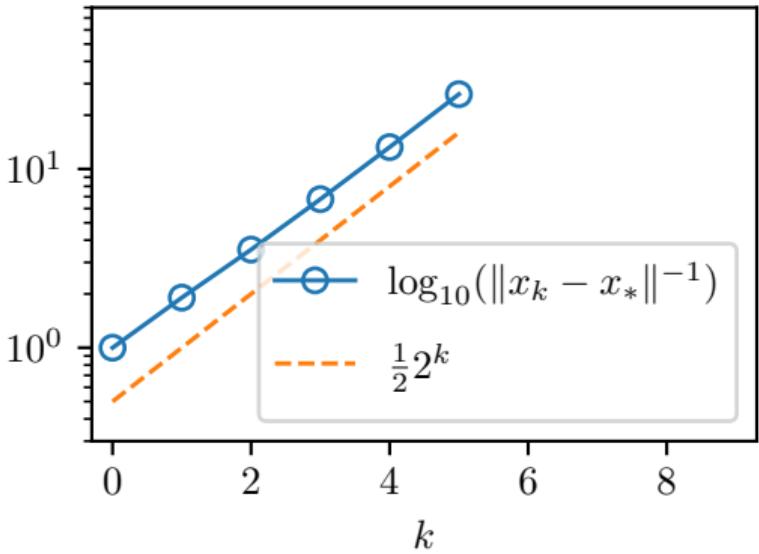
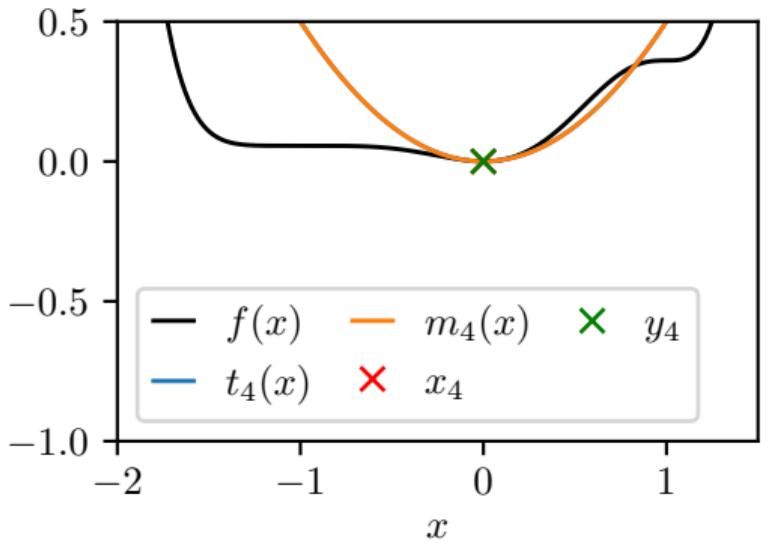
Newton's method ($p = 2$, $\sigma = 0$)



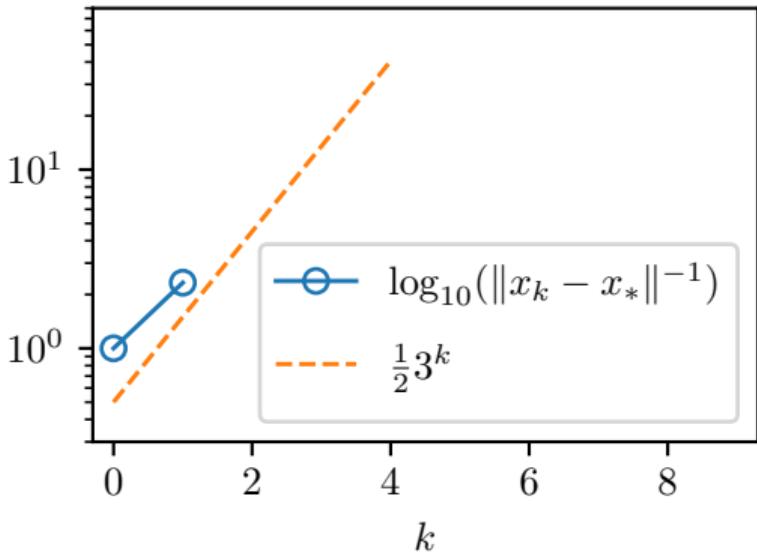
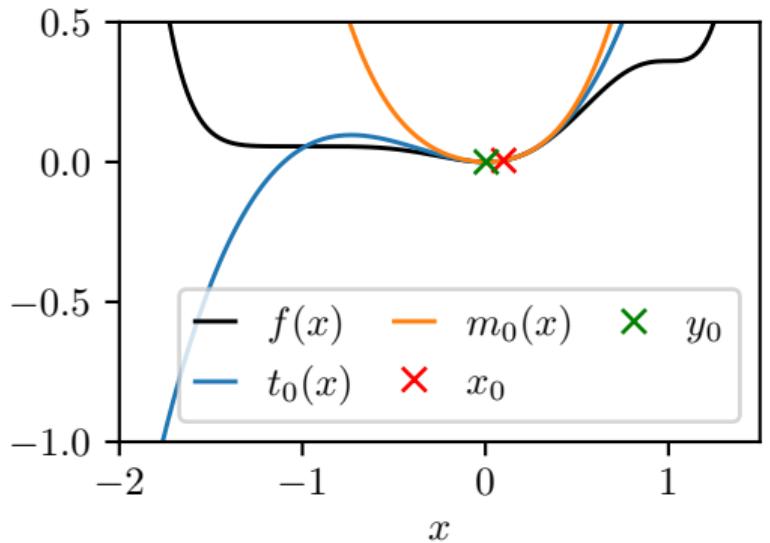
Newton's method ($p = 2$, $\sigma = 0$)

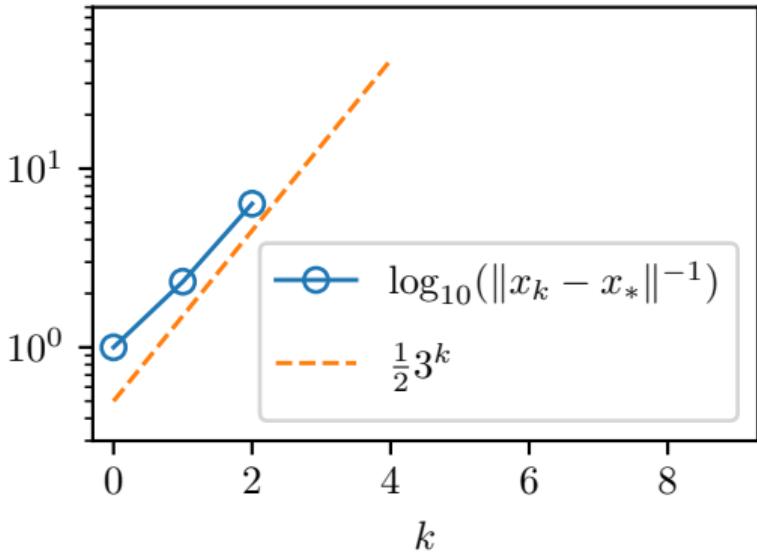
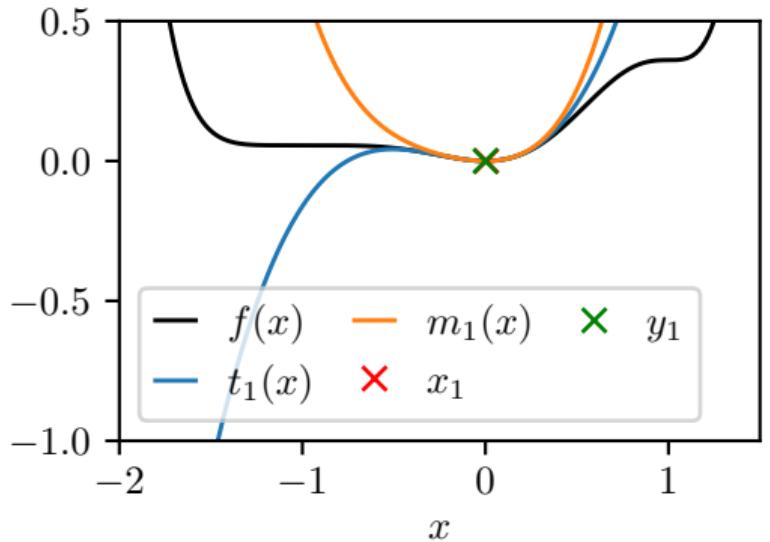


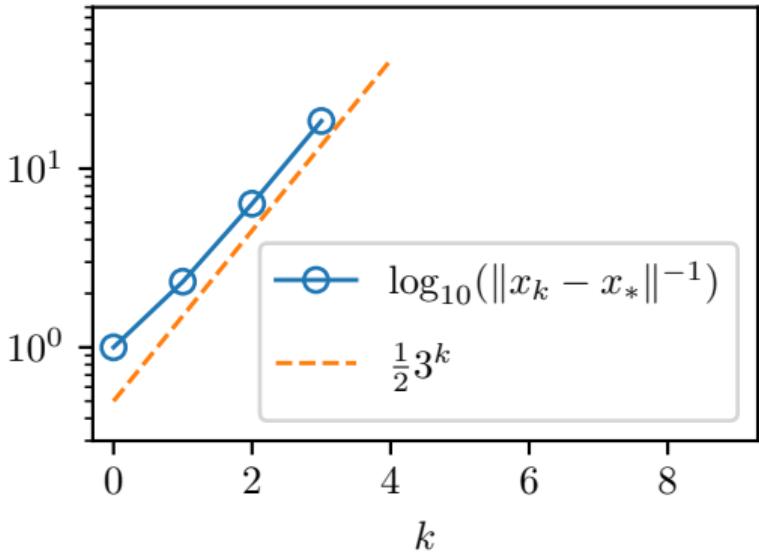
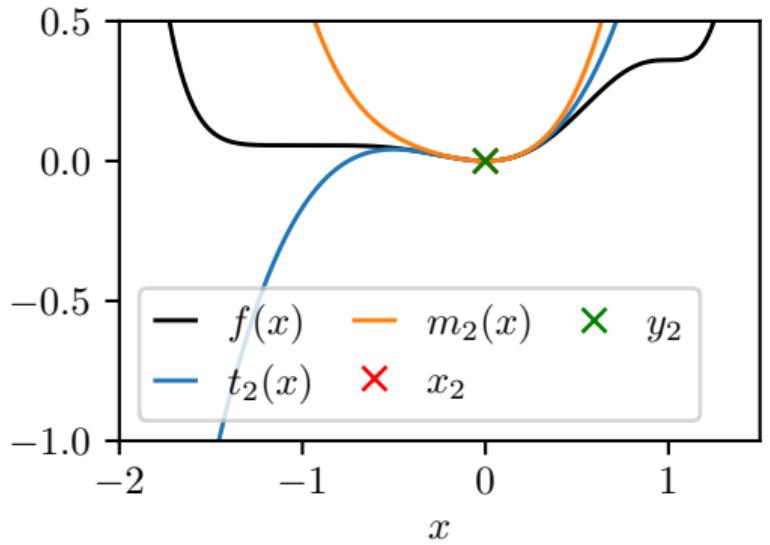
Newton's method ($p = 2$, $\sigma = 0$)

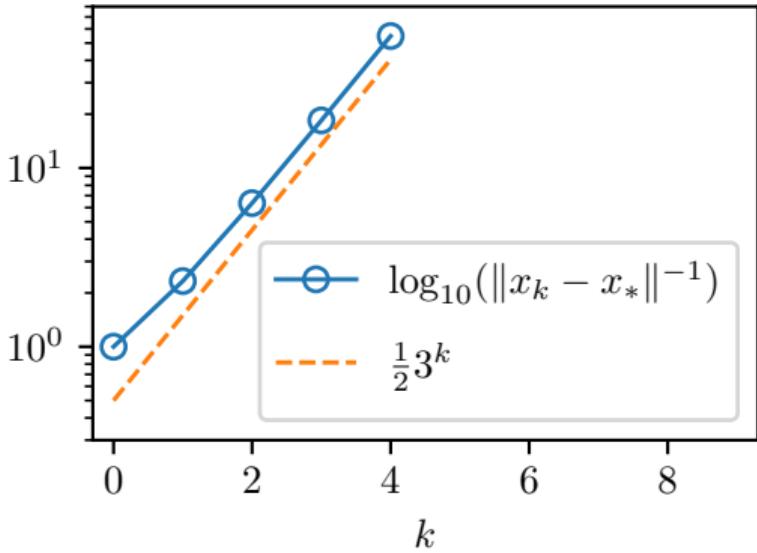
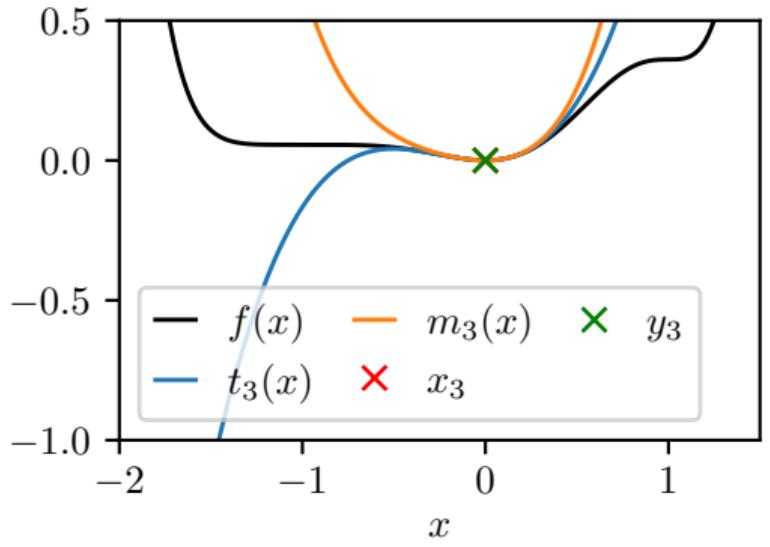


AR3, semi-adaptive σ_k

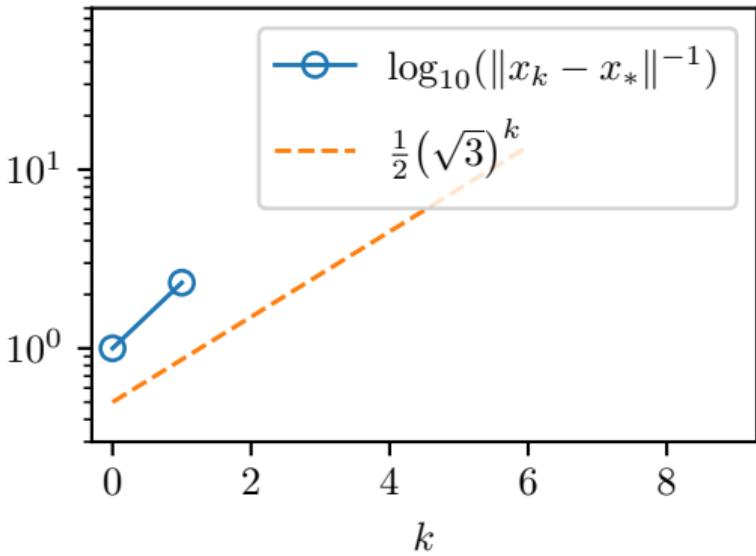
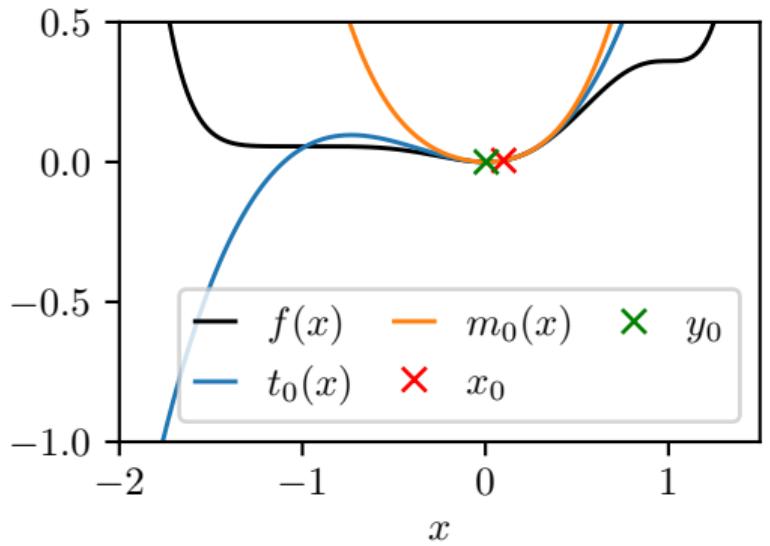




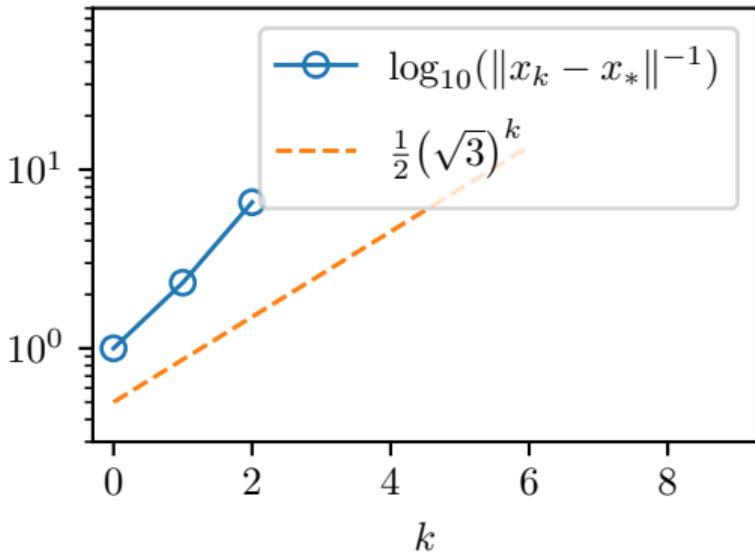
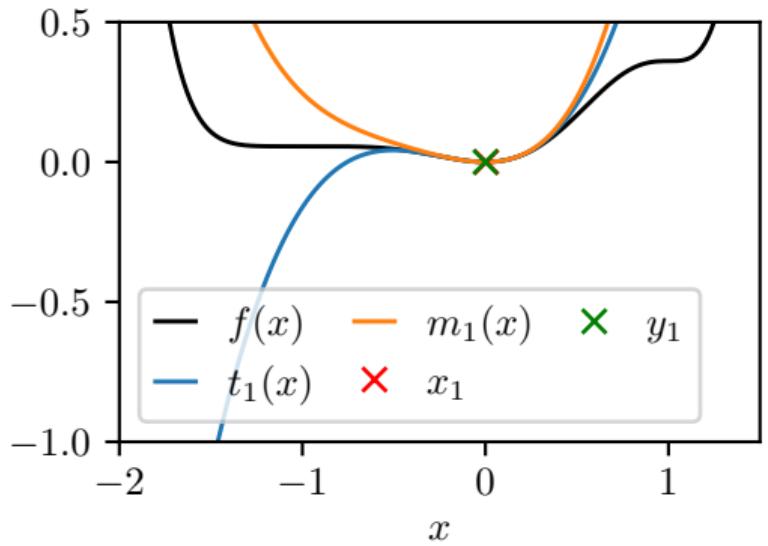




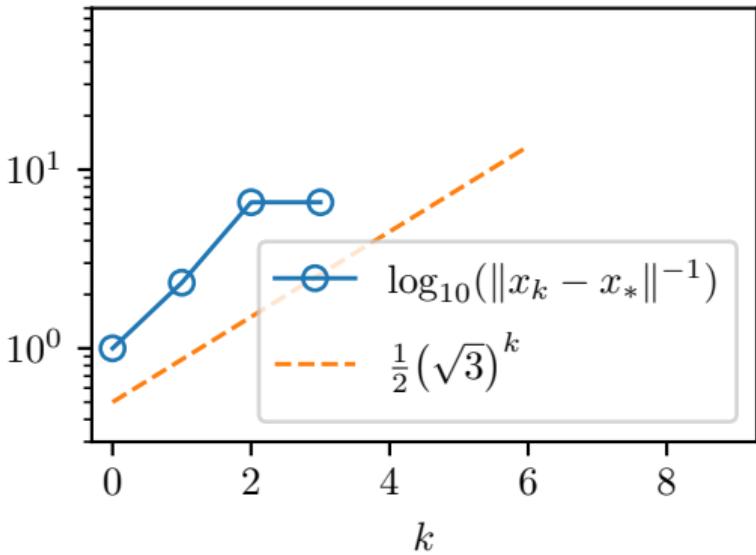
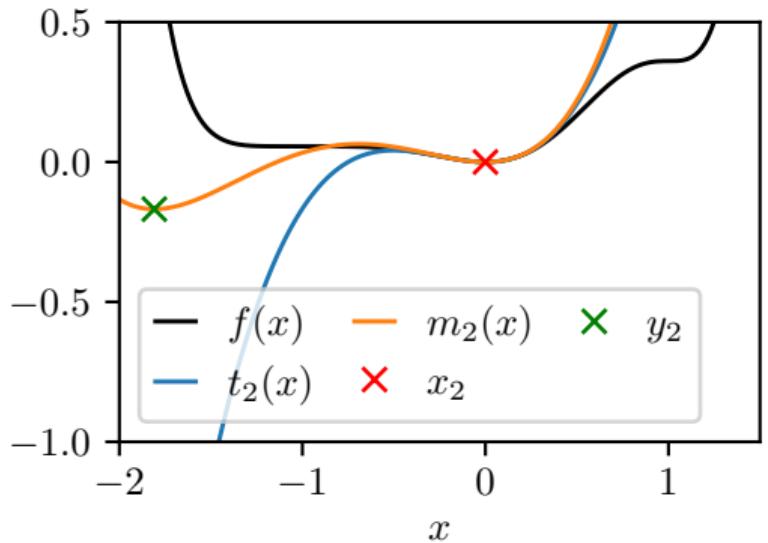
AR3, adaptive σ_k and global model minimizer



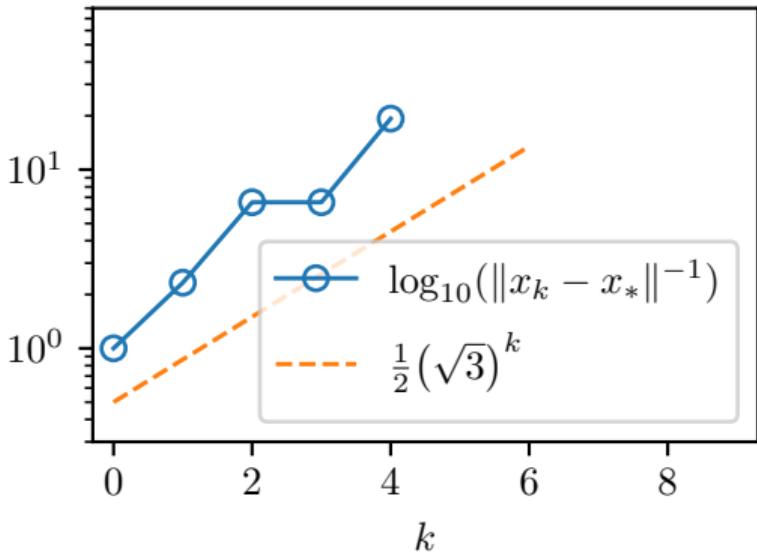
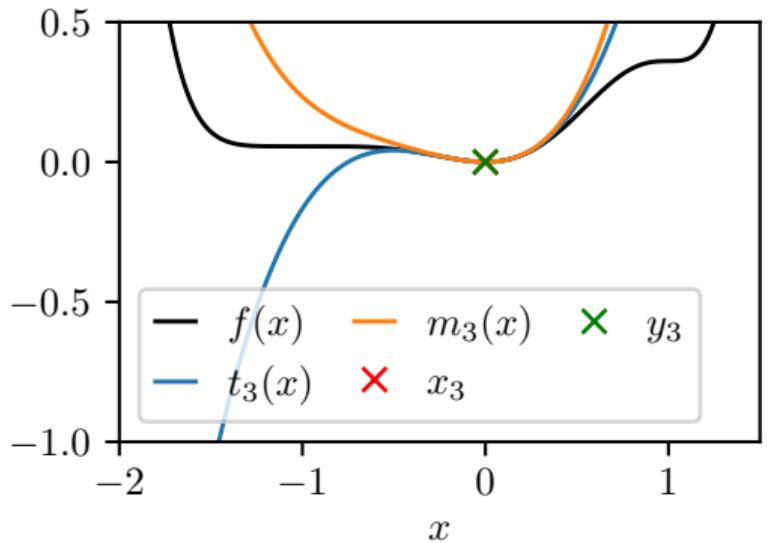
AR3, adaptive σ_k and global model minimizer



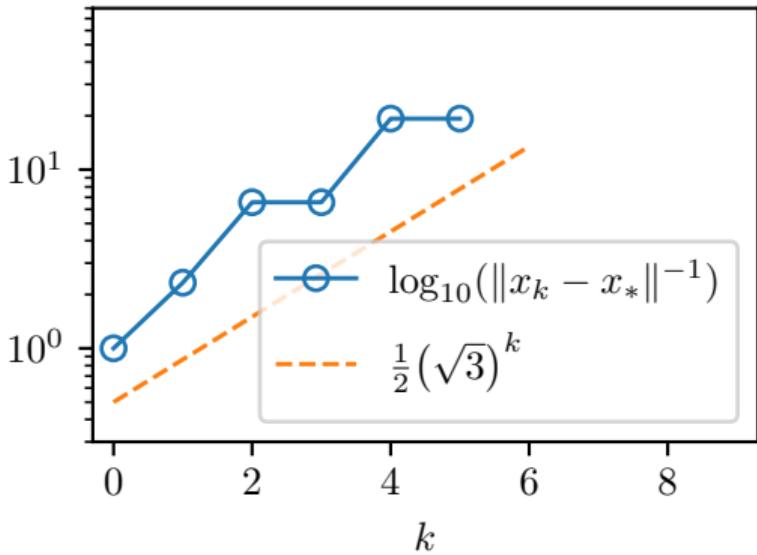
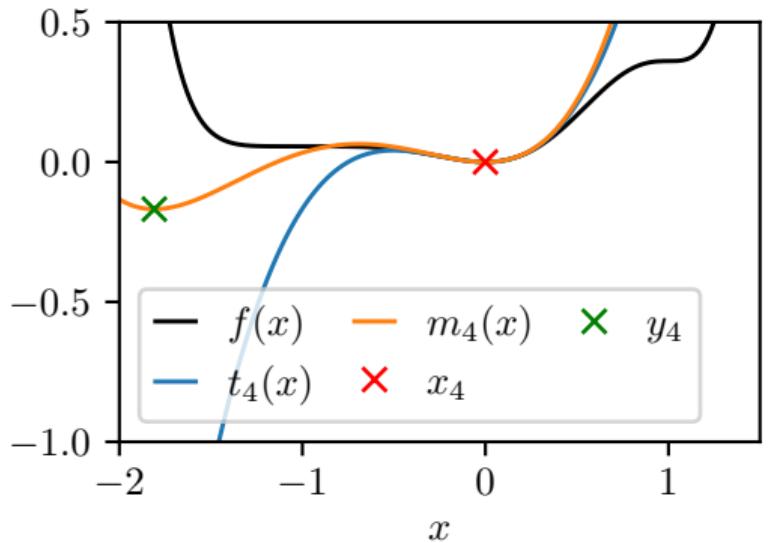
AR3, adaptive σ_k and global model minimizer



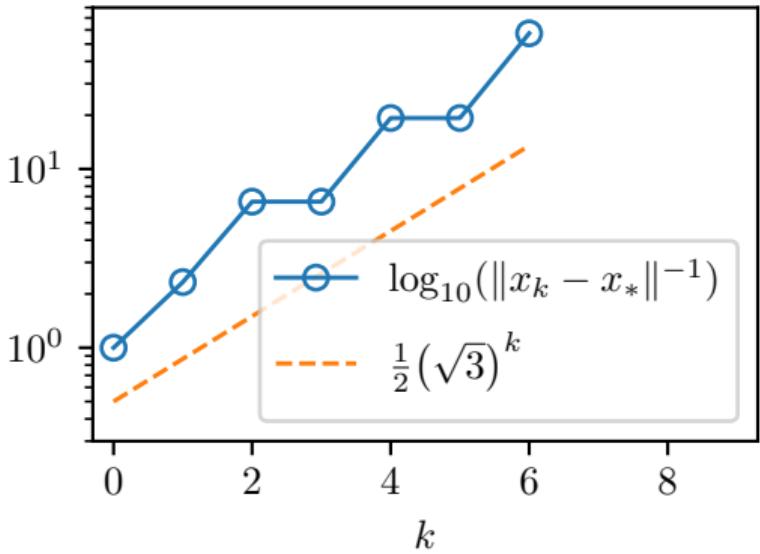
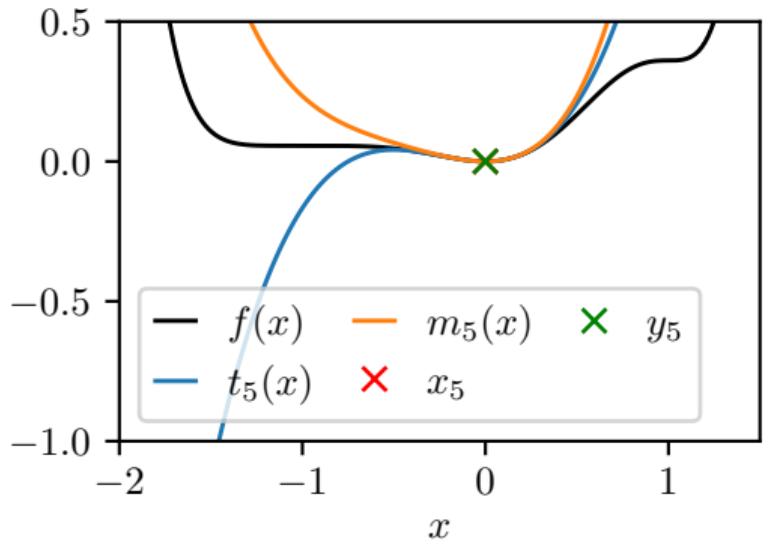
AR3, adaptive σ_k and global model minimizer



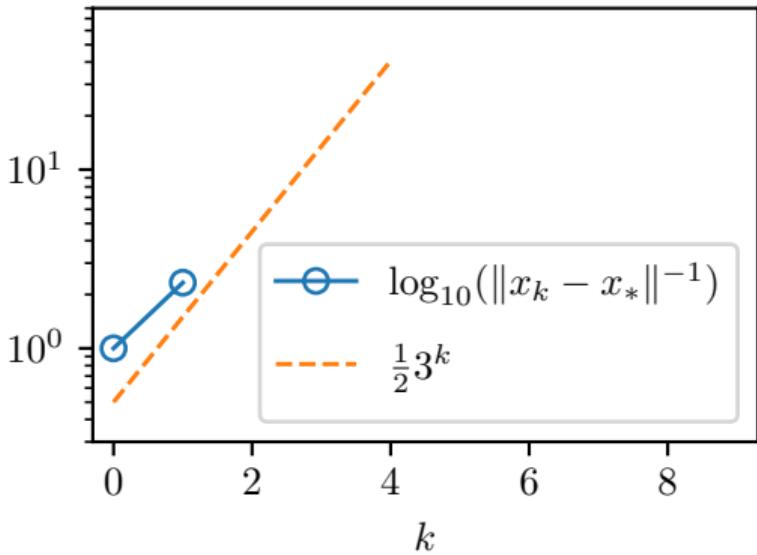
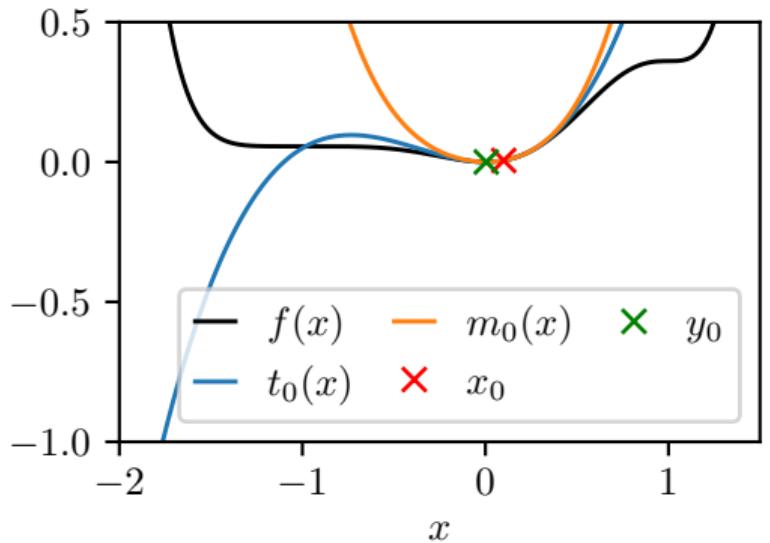
AR3, adaptive σ_k and global model minimizer



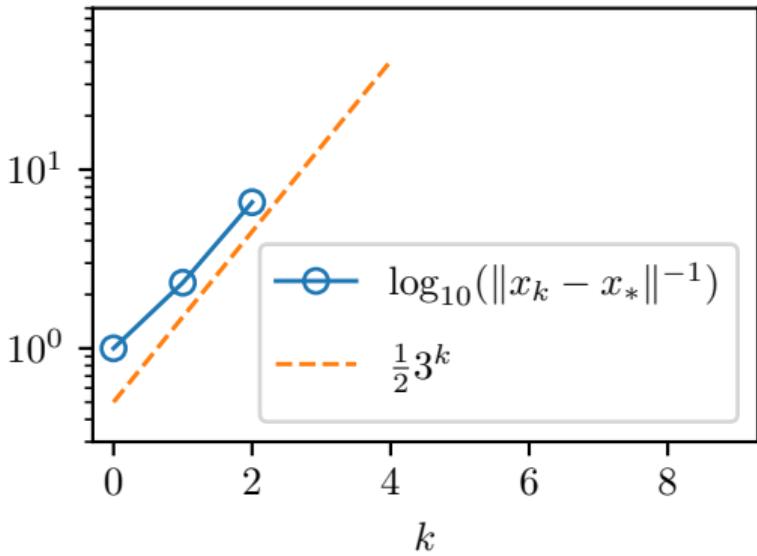
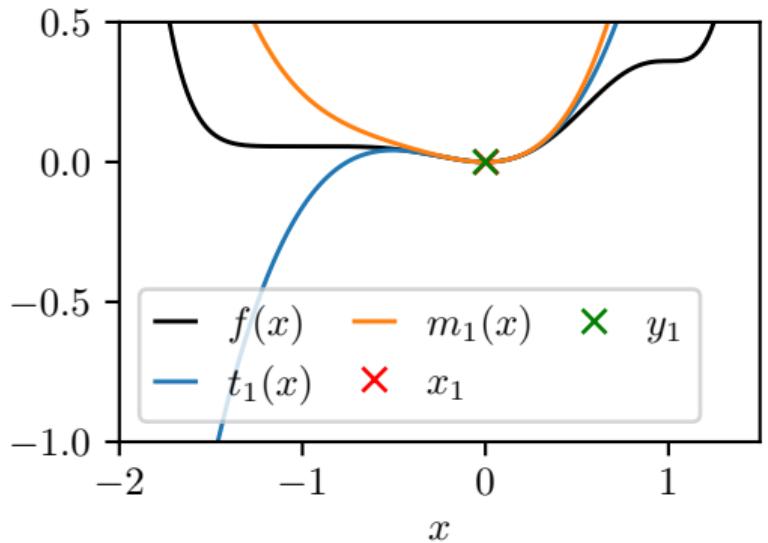
AR3, adaptive σ_k and global model minimizer



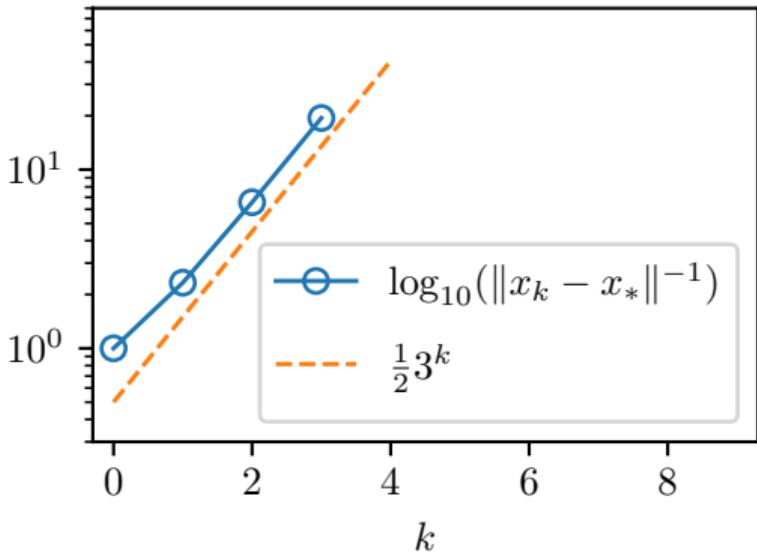
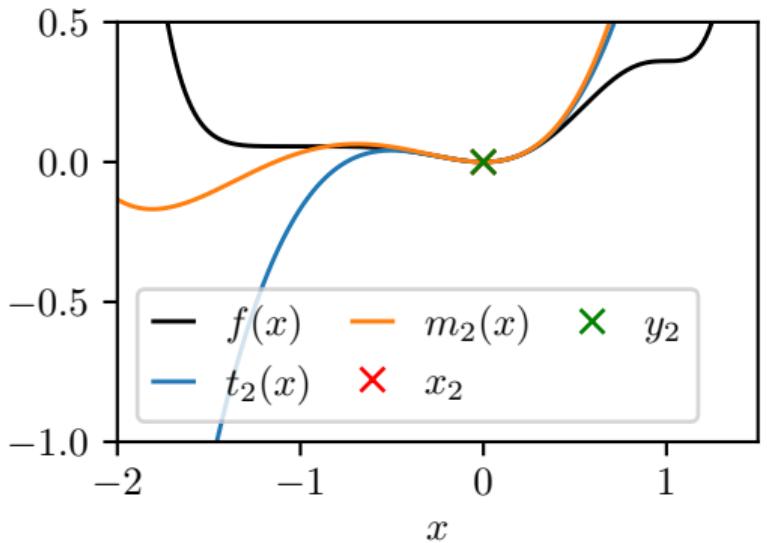
AR3, adaptive σ_k and right model minimizer



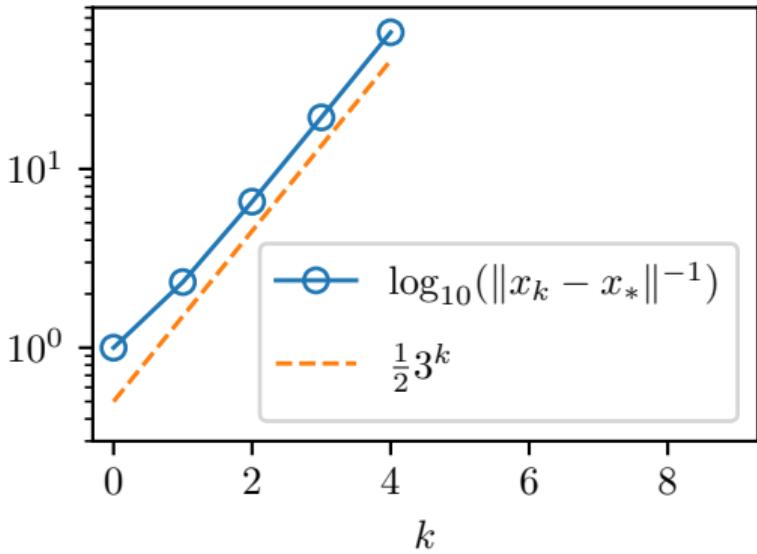
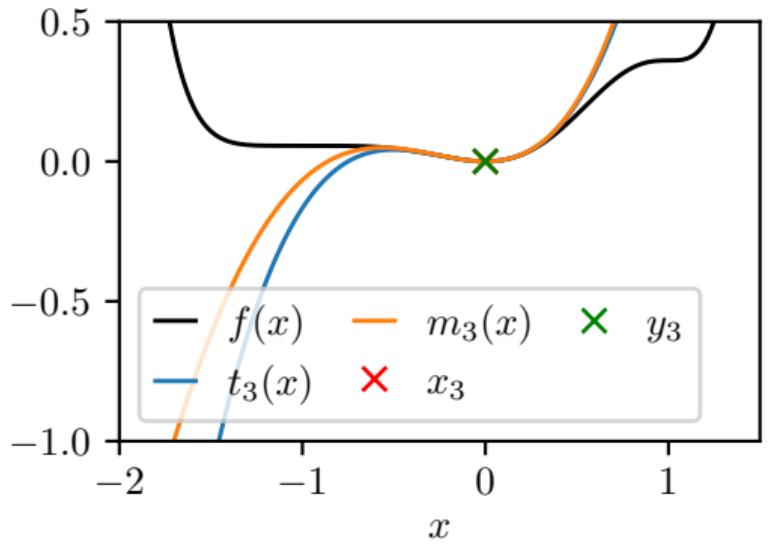
AR3, adaptive σ_k and right model minimizer



AR3, adaptive σ_k and right model minimizer



AR3, adaptive σ_k and right model minimizer



Definition

A function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is uniformly convex of order $q \geq 2$ inside $\Omega \subset \mathbb{R}^n$ if

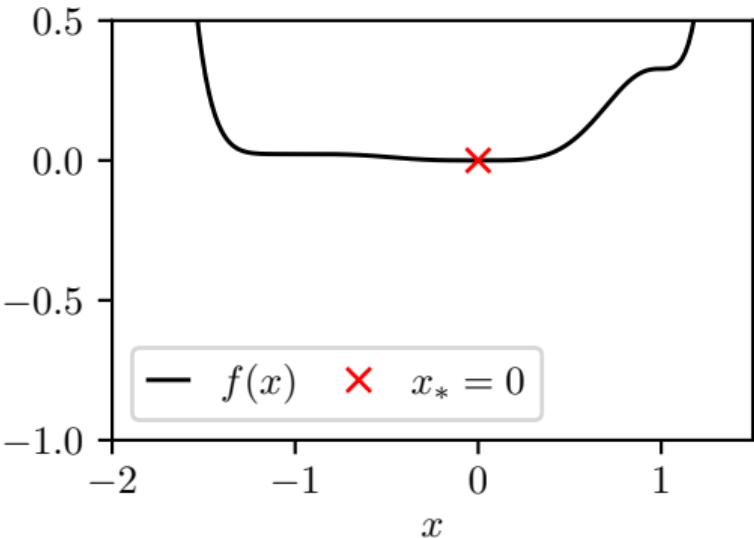
$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{\mu_q}{q} \|\mathbf{y} - \mathbf{x}\|^q \quad \forall \mathbf{x}, \mathbf{y} \in \Omega$$

for some $\mu_q > 0$.

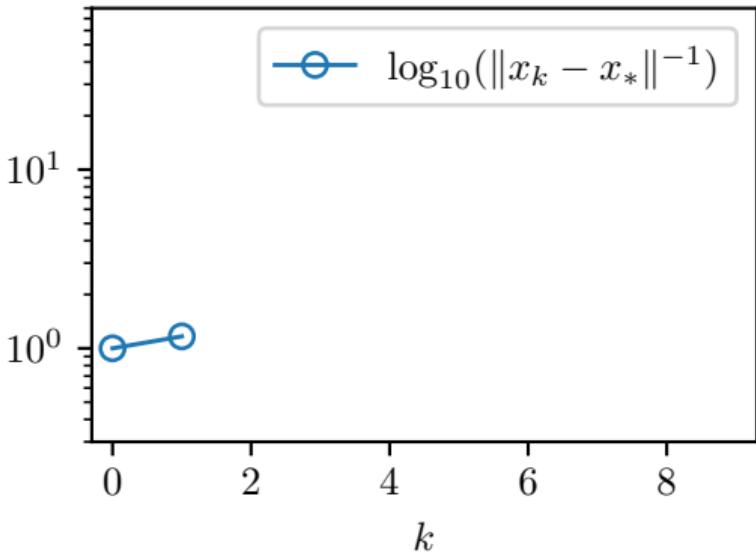
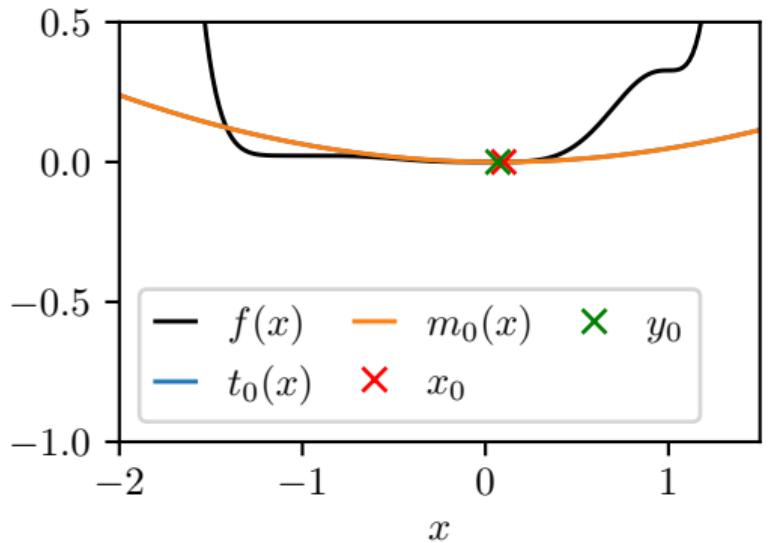
- Uniformly convex of order q around \mathbf{x}_* $\implies f(\mathbf{x}) - f(\mathbf{x}_*) \geq \frac{\mu_q}{q} \|\mathbf{x} - \mathbf{x}_*\|^q$
- Uniformly convex of order q around \mathbf{x}_* $\implies \mathbf{x}_*$ is an isolated minimizer
- Uniformly convex of order 2 \iff strongly convex

Example with degenerate minimizer

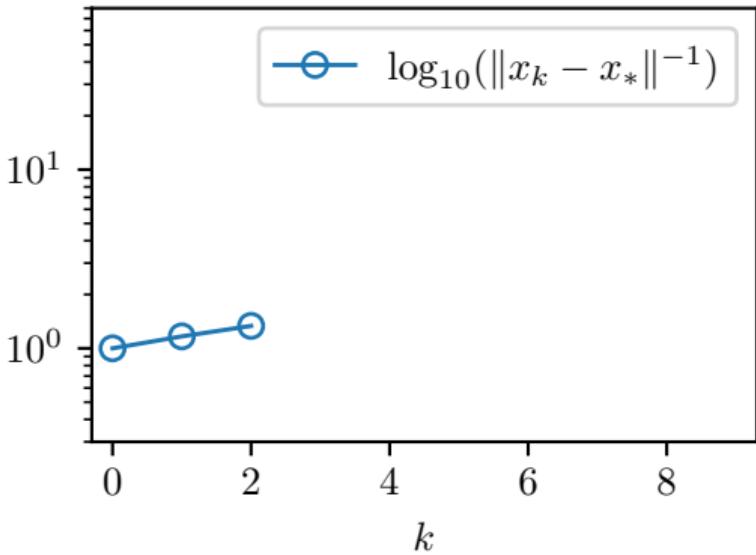
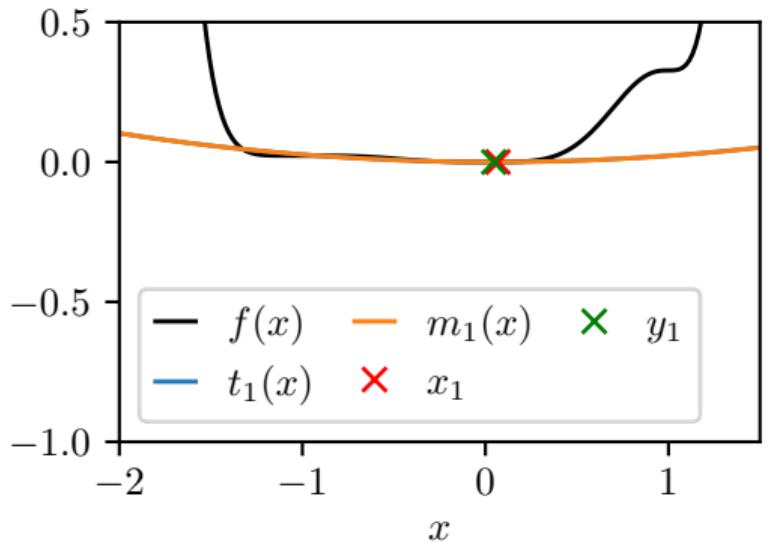
- $f(x) = \int_0^x 3(t+1)^4(t-1)^2t^3 dt$
- f has one local (and global) minimizer at $x_* = 0$
- f is locally uniformly convex around x_* with $q = 4$:
 - $f''(x_*) = 0$
 - $f'''(x_*) = 0$
 - $f''''(x_*) = 18$



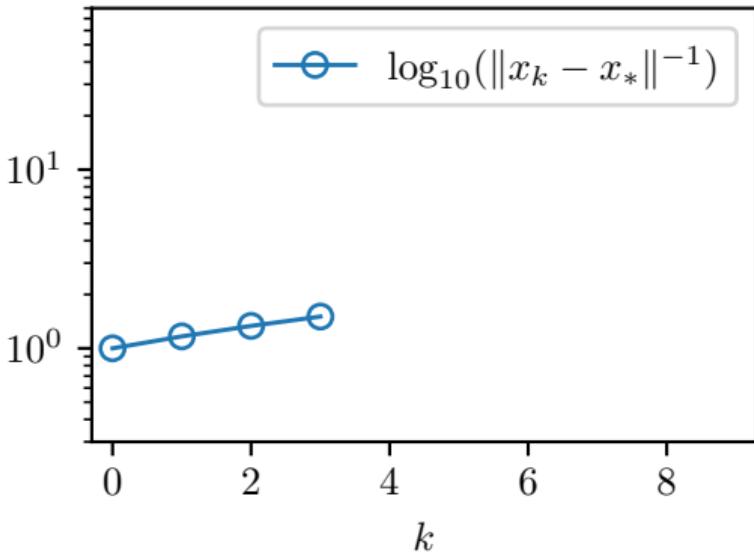
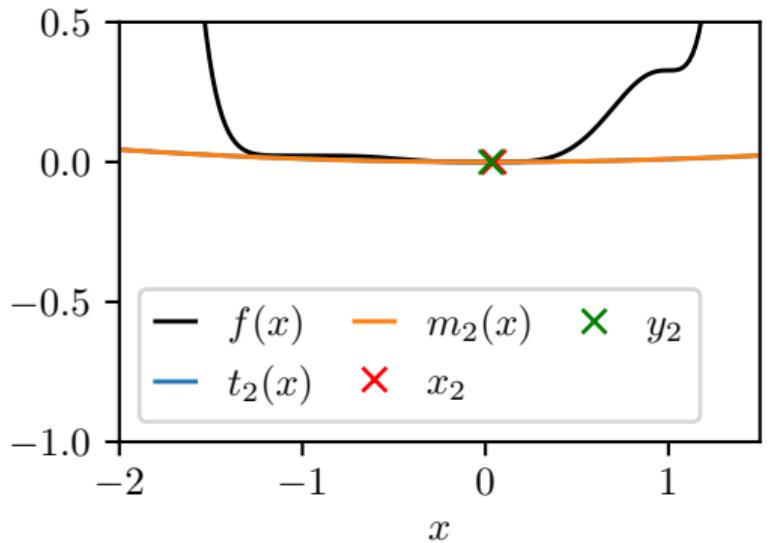
Newton's method ($p = 2$, $\sigma = 0$)



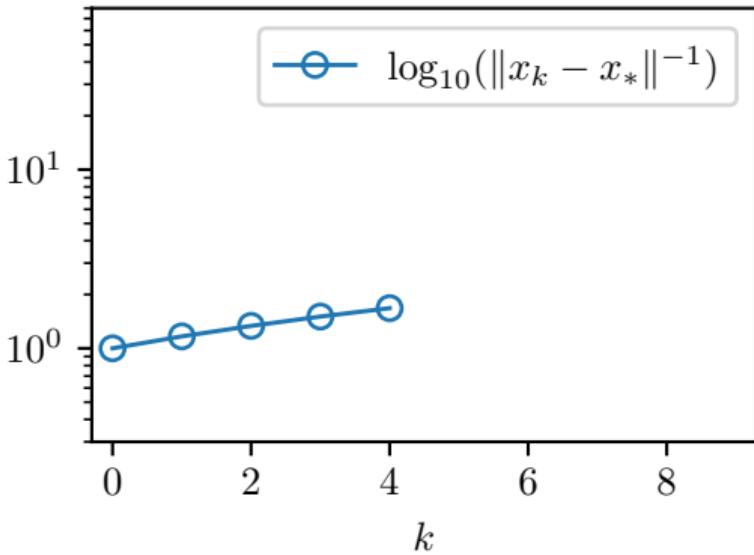
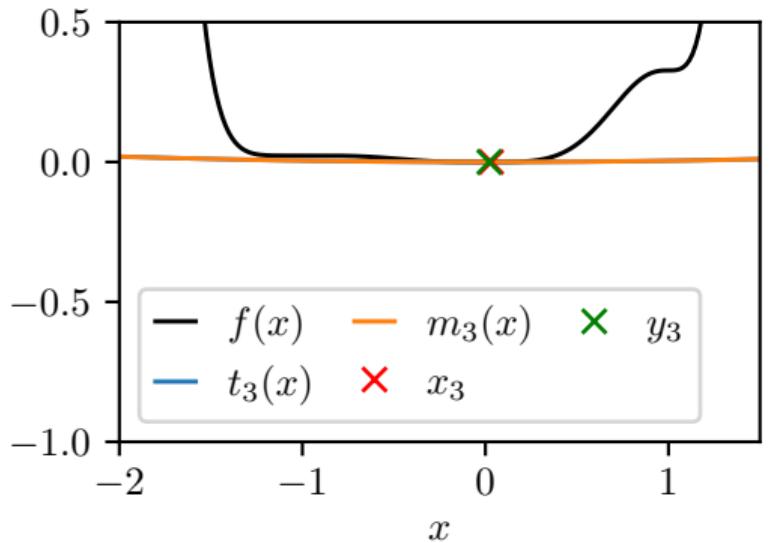
Newton's method ($p = 2$, $\sigma = 0$)



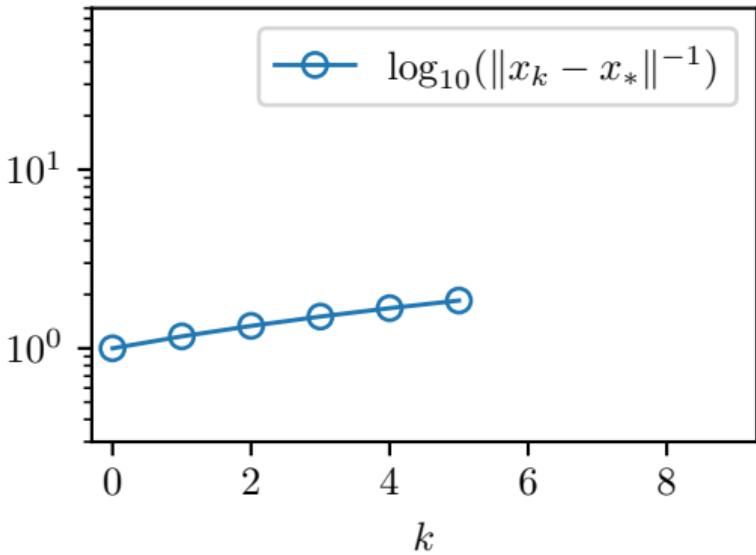
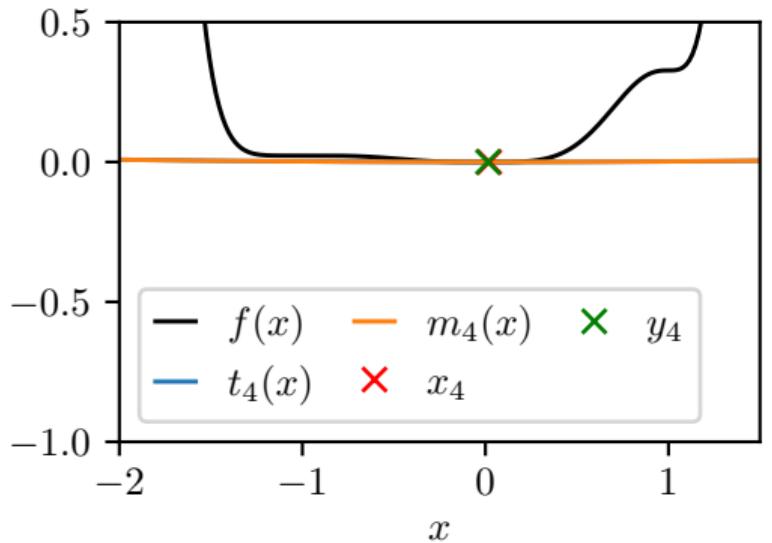
Newton's method ($p = 2$, $\sigma = 0$)



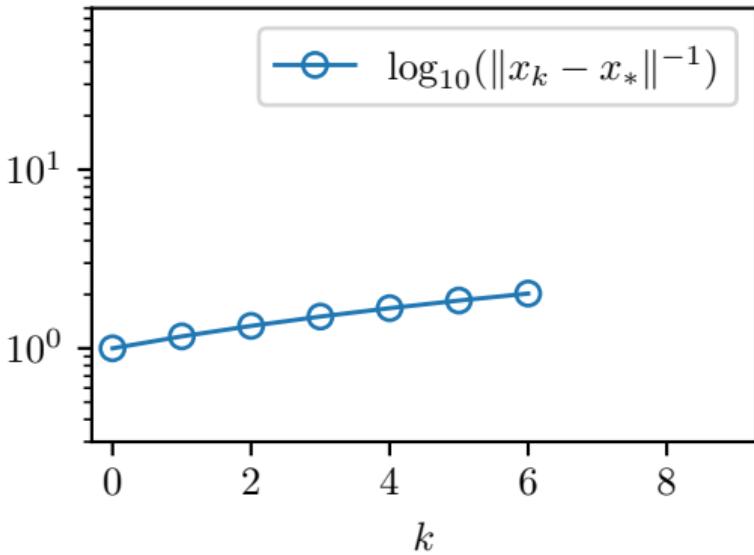
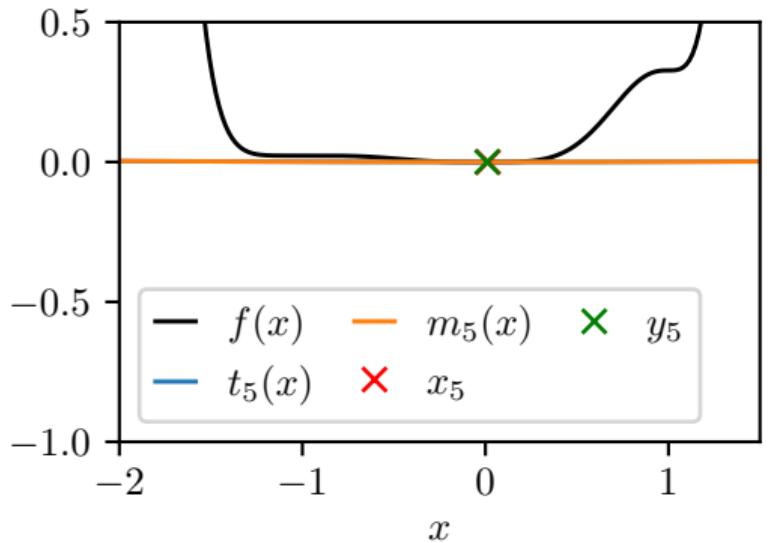
Newton's method ($p = 2$, $\sigma = 0$)



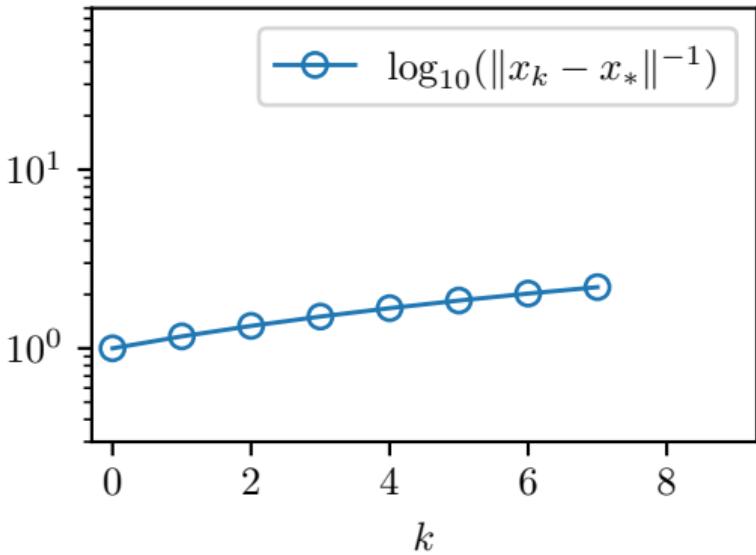
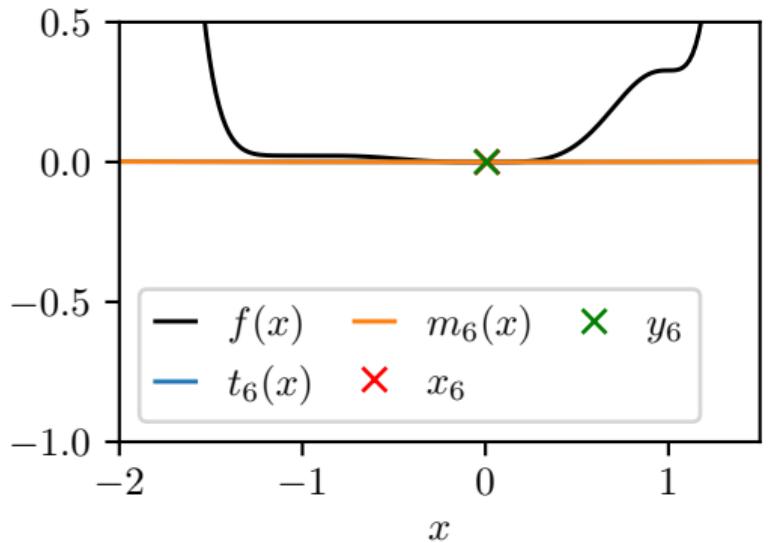
Newton's method ($p = 2$, $\sigma = 0$)



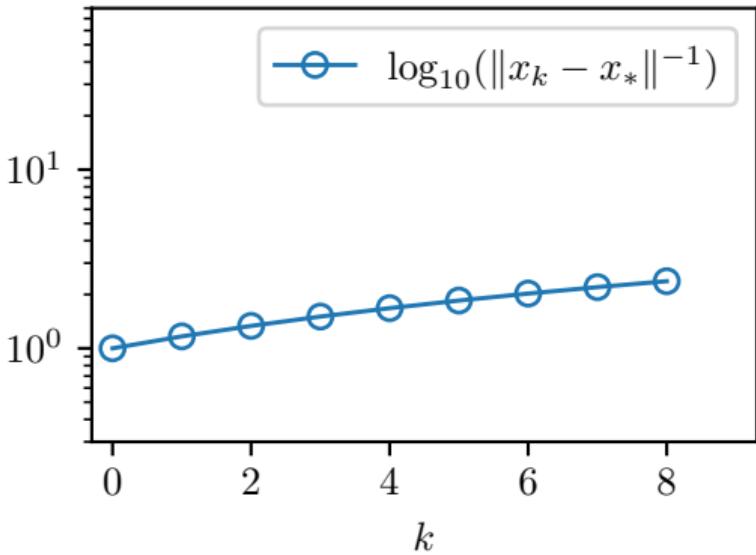
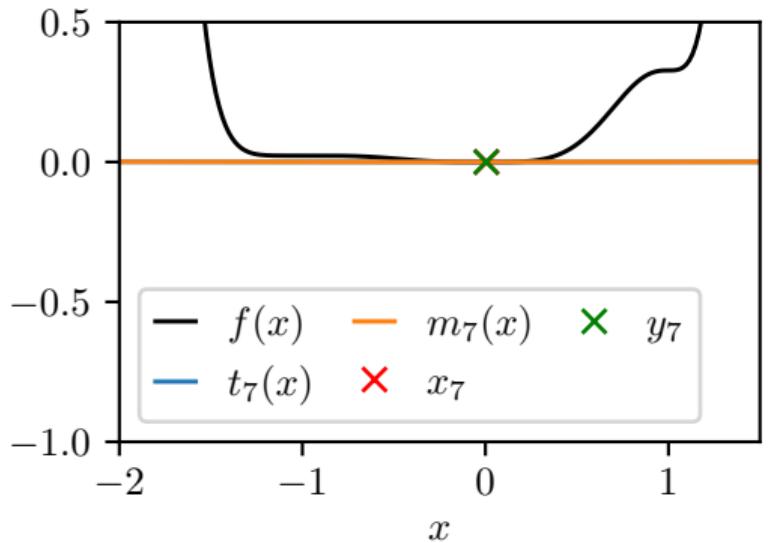
Newton's method ($p = 2$, $\sigma = 0$)



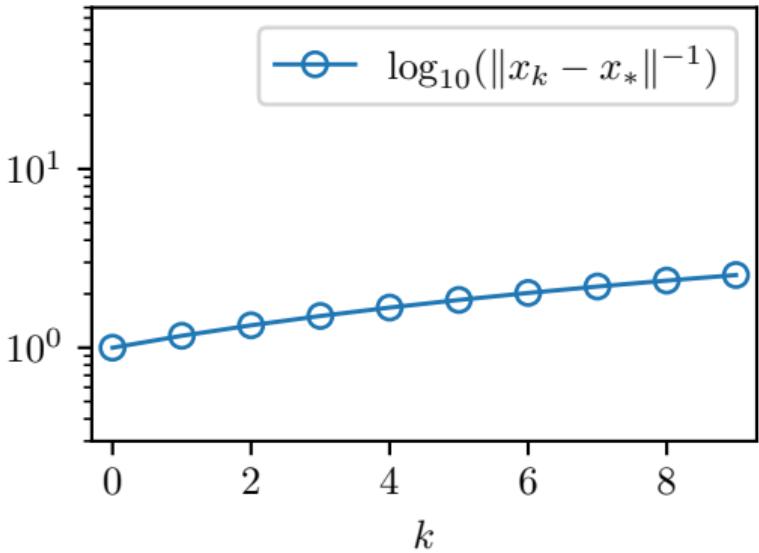
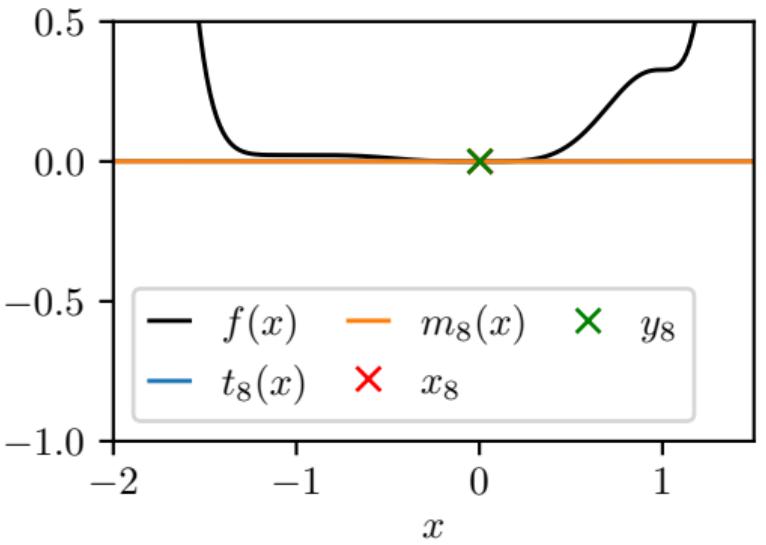
Newton's method ($p = 2$, $\sigma = 0$)



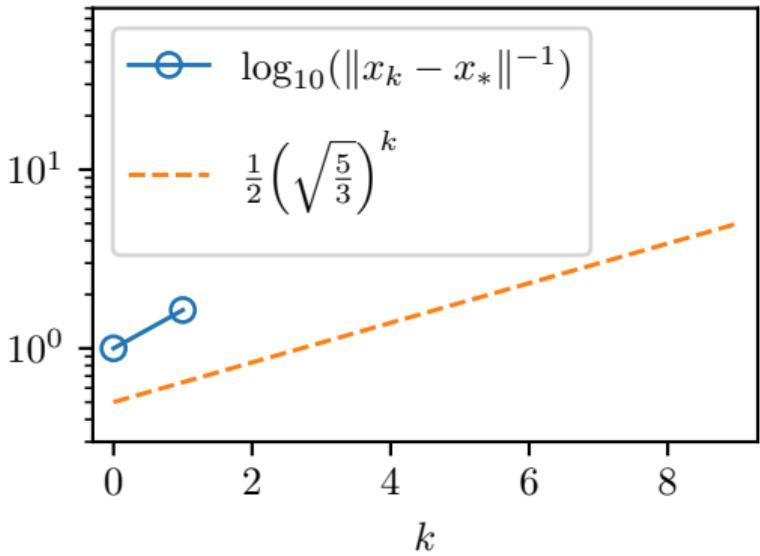
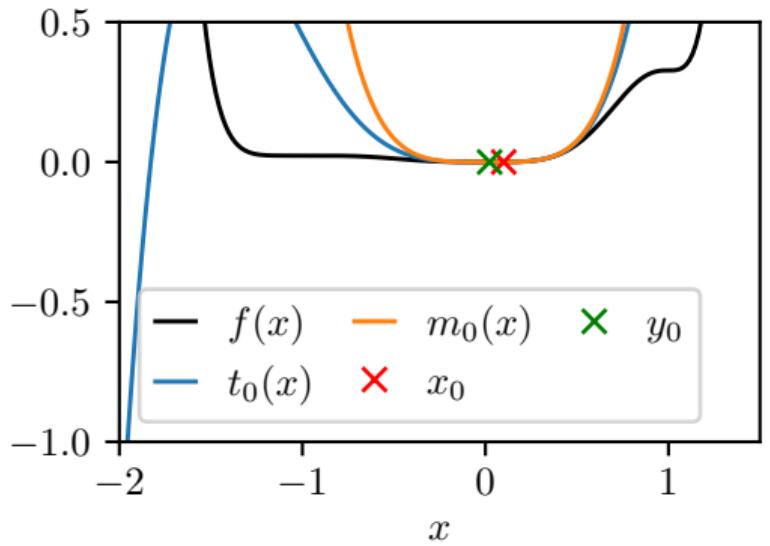
Newton's method ($p = 2$, $\sigma = 0$)



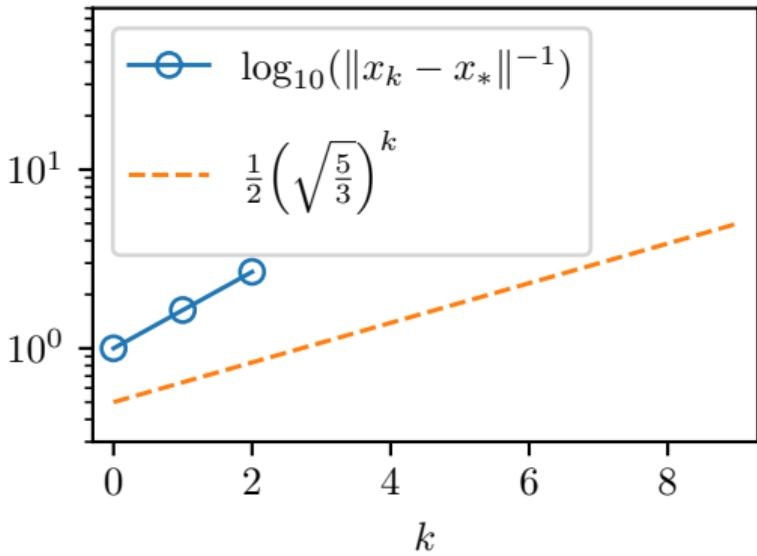
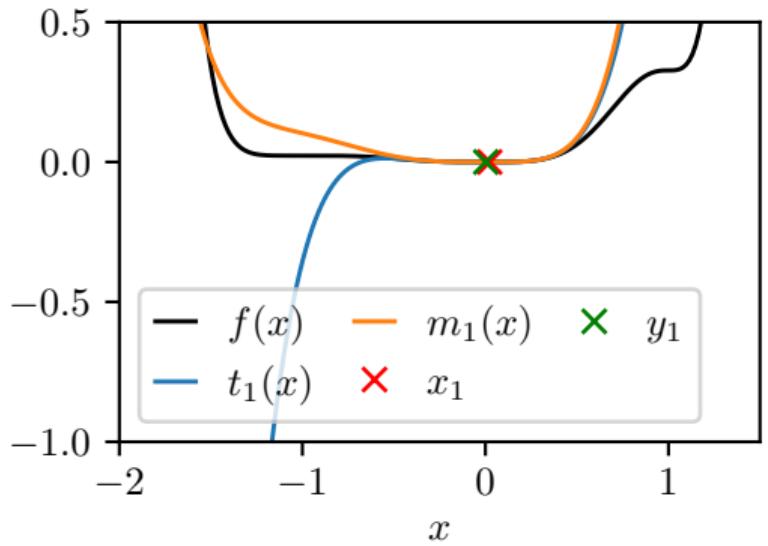
Newton's method ($p = 2$, $\sigma = 0$)



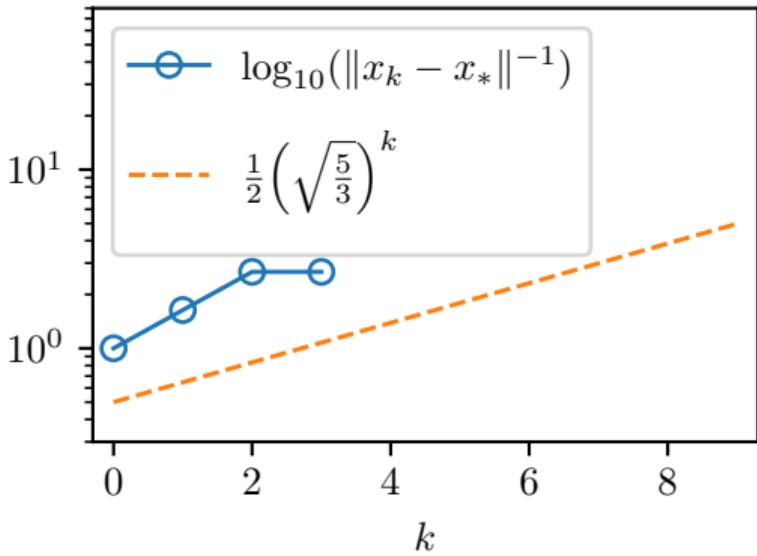
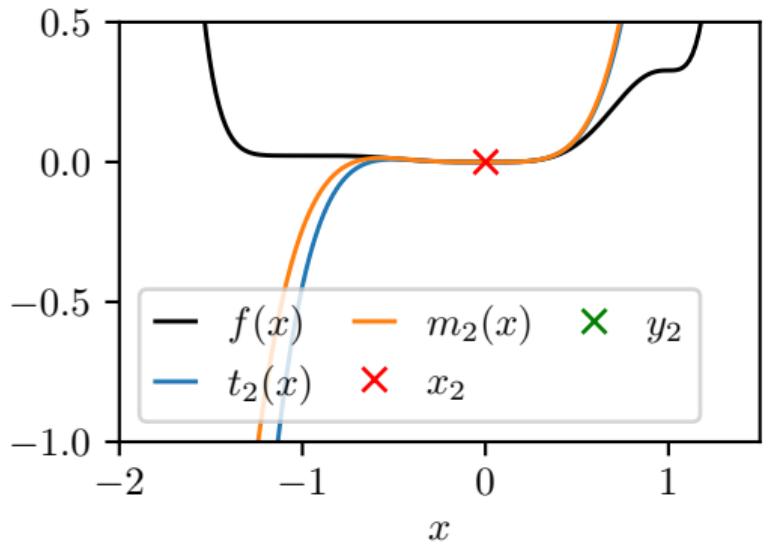
AR5, adaptive σ_k and global model minimizer



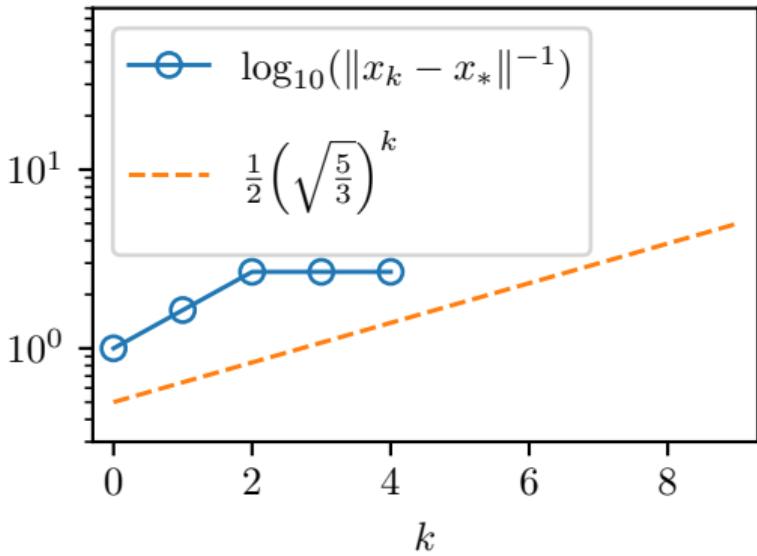
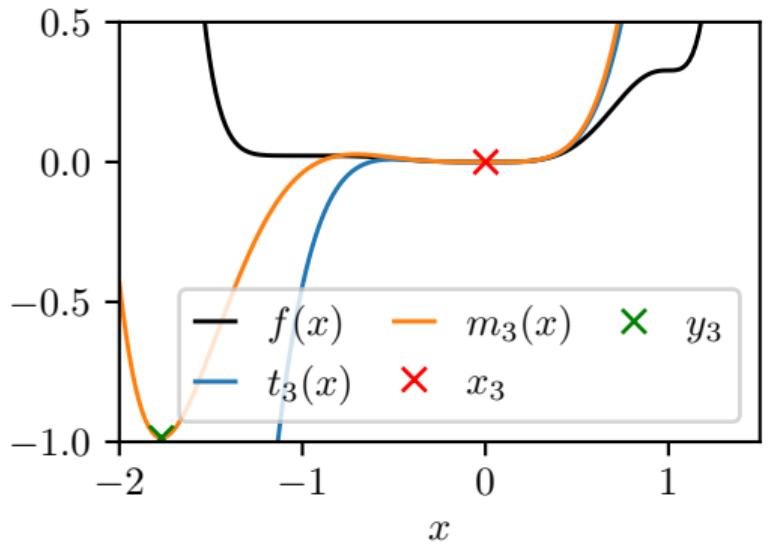
AR5, adaptive σ_k and global model minimizer



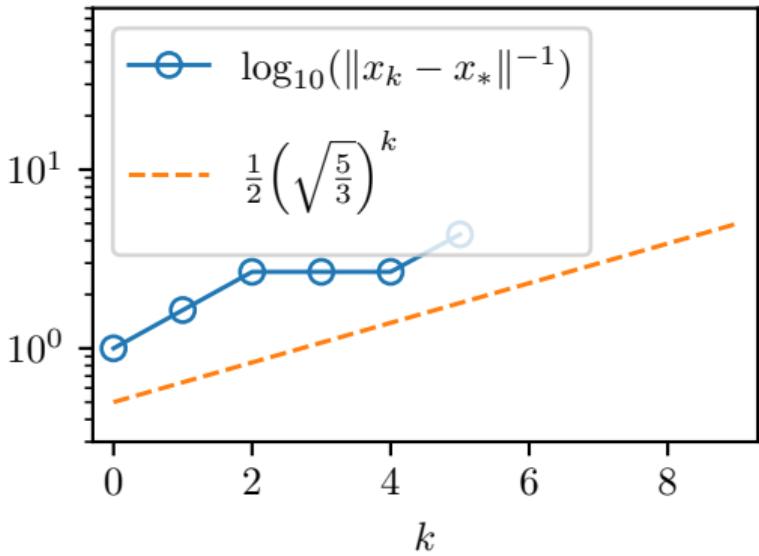
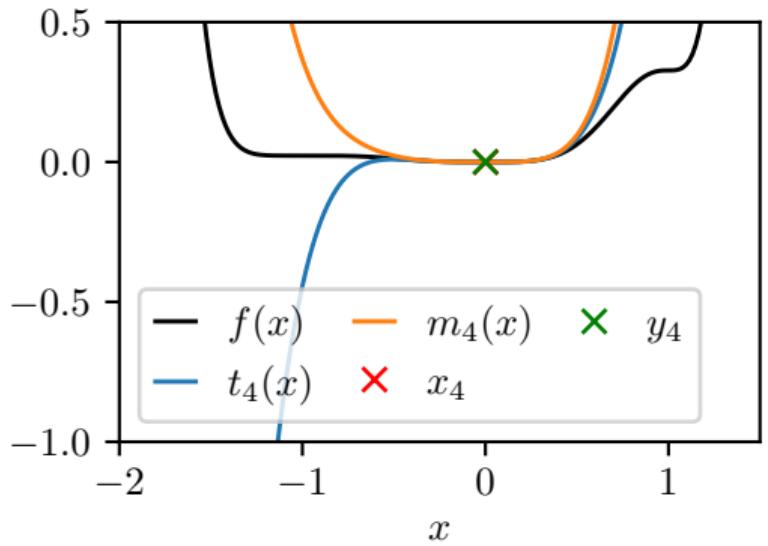
AR5, adaptive σ_k and global model minimizer



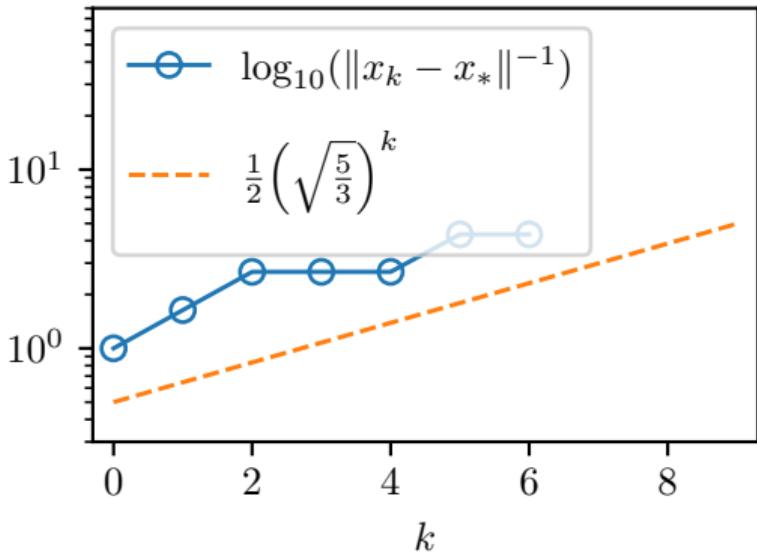
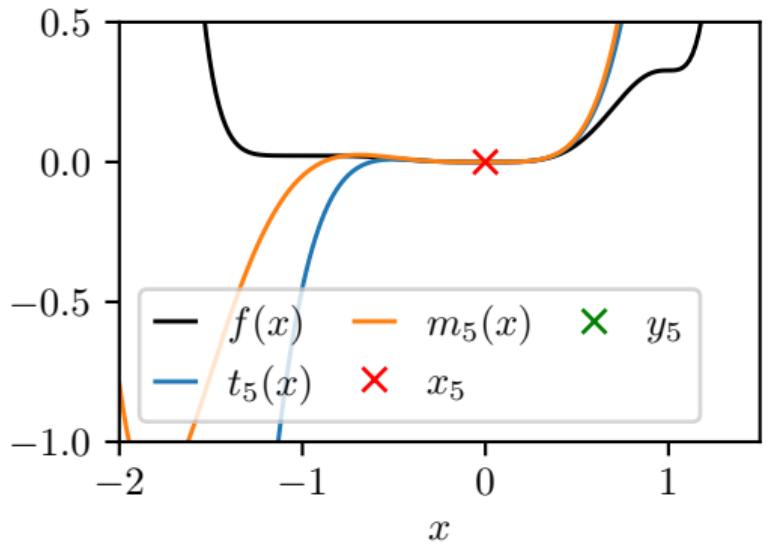
AR5, adaptive σ_k and global model minimizer



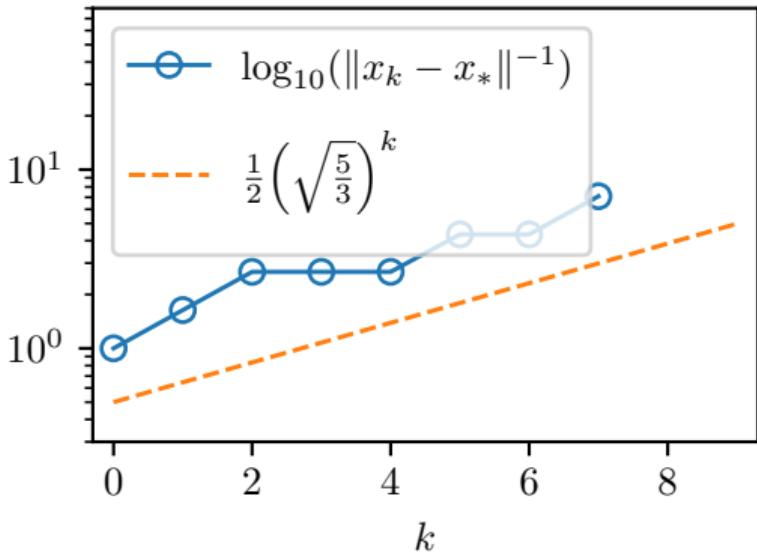
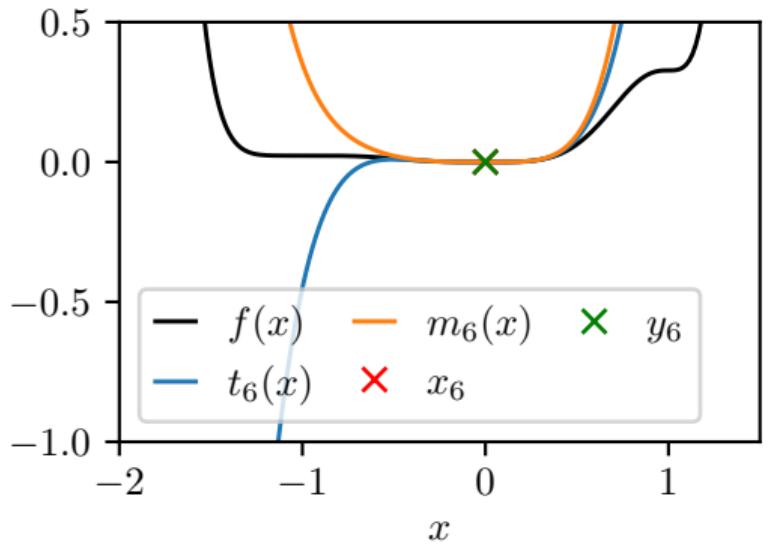
AR5, adaptive σ_k and global model minimizer



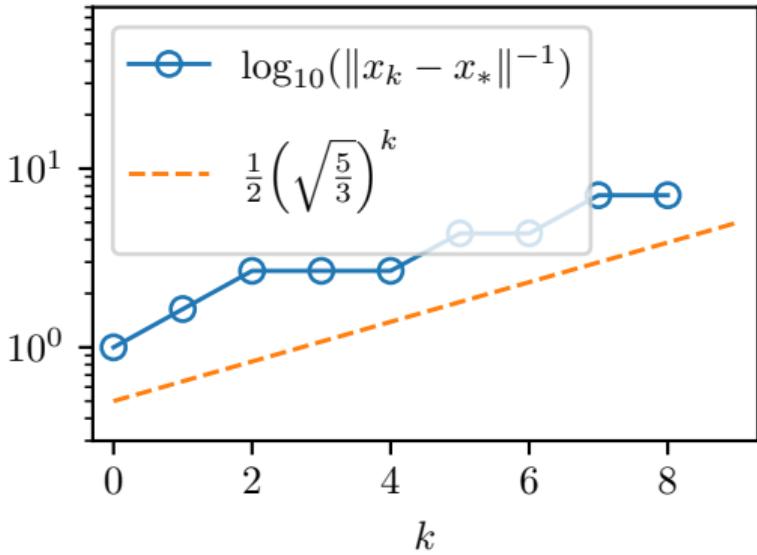
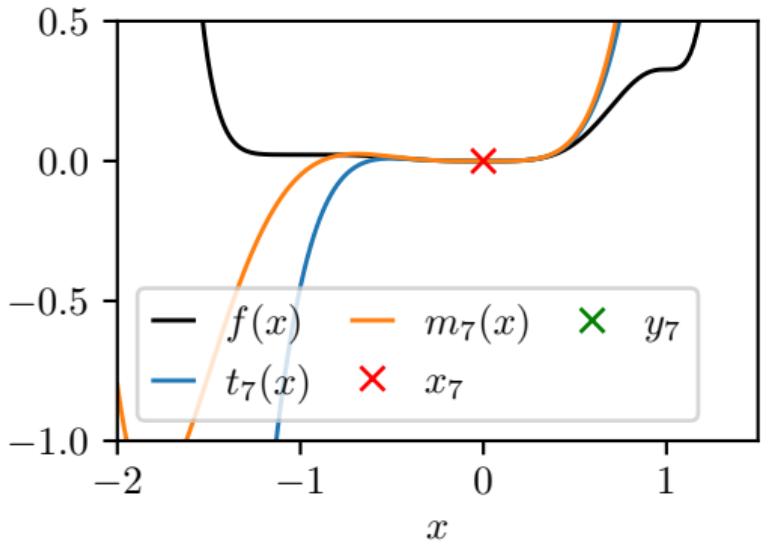
AR5, adaptive σ_k and global model minimizer



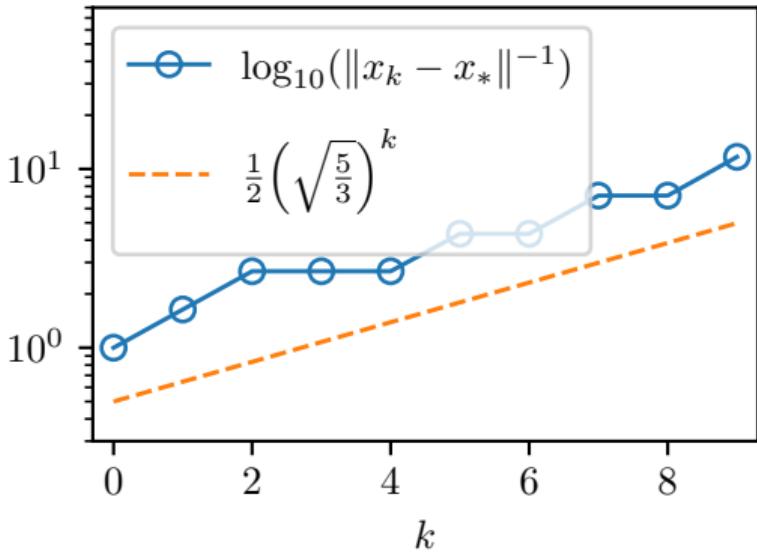
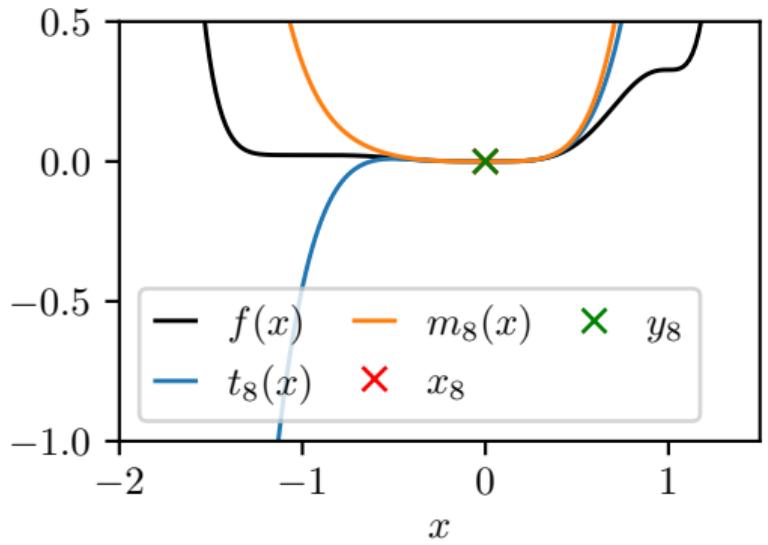
AR5, adaptive σ_k and global model minimizer



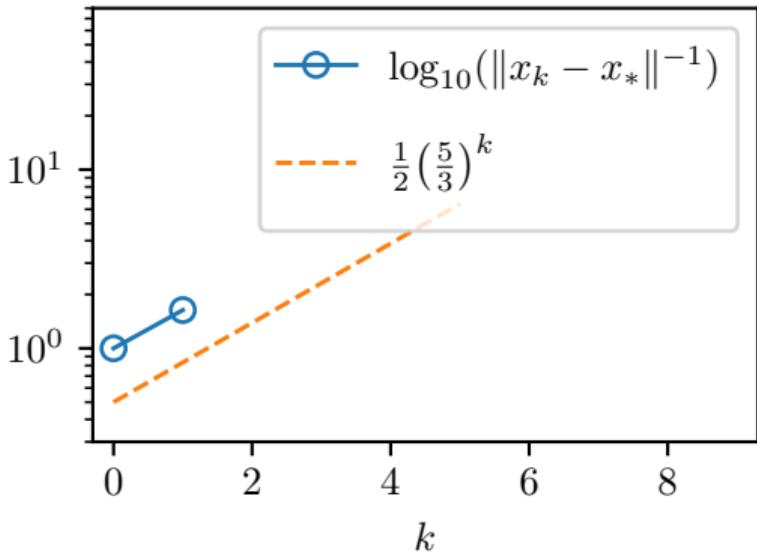
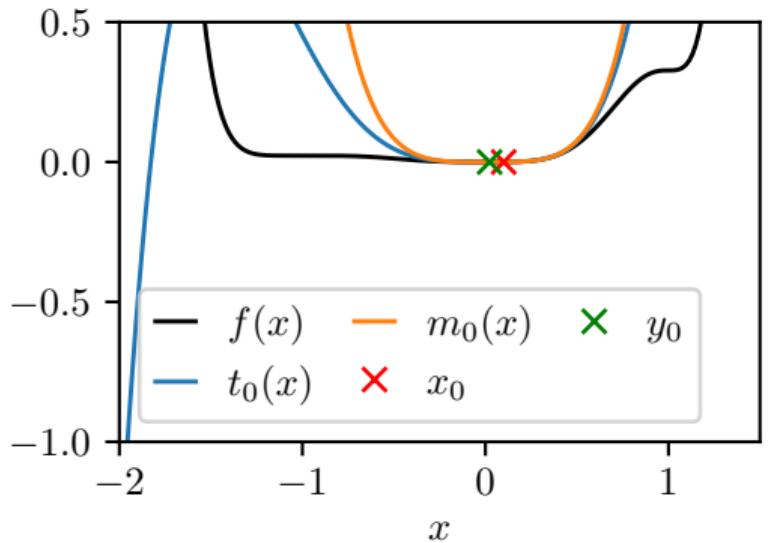
AR5, adaptive σ_k and global model minimizer



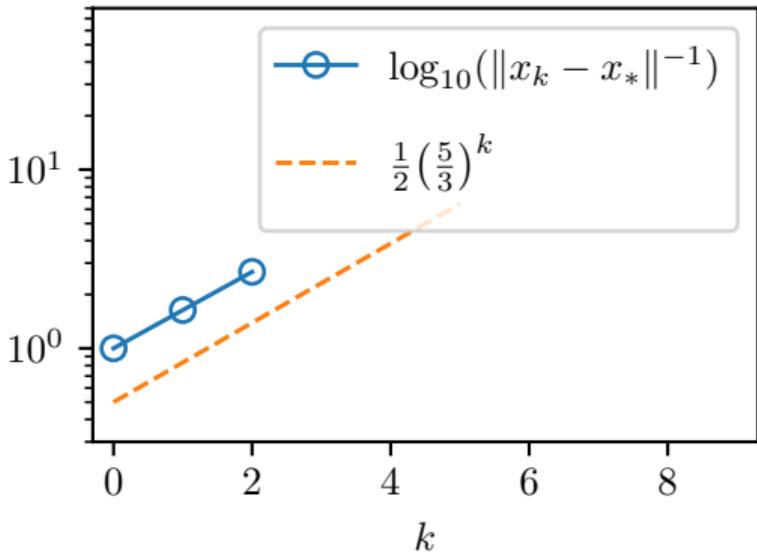
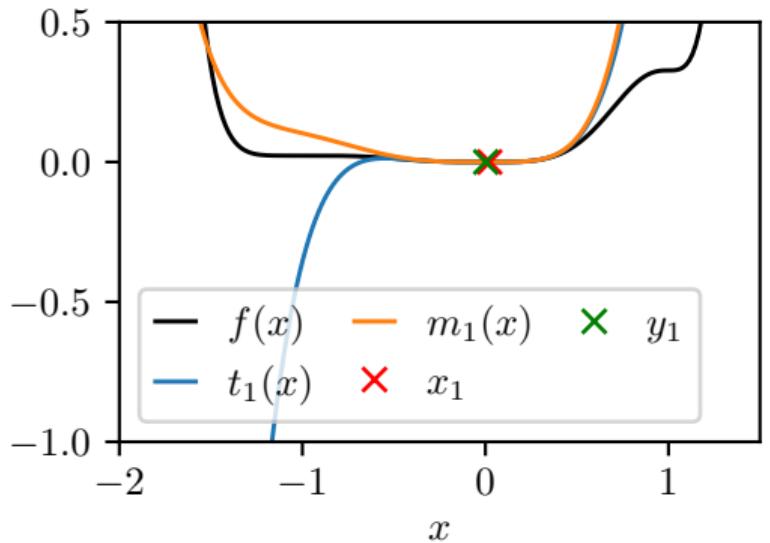
AR5, adaptive σ_k and global model minimizer



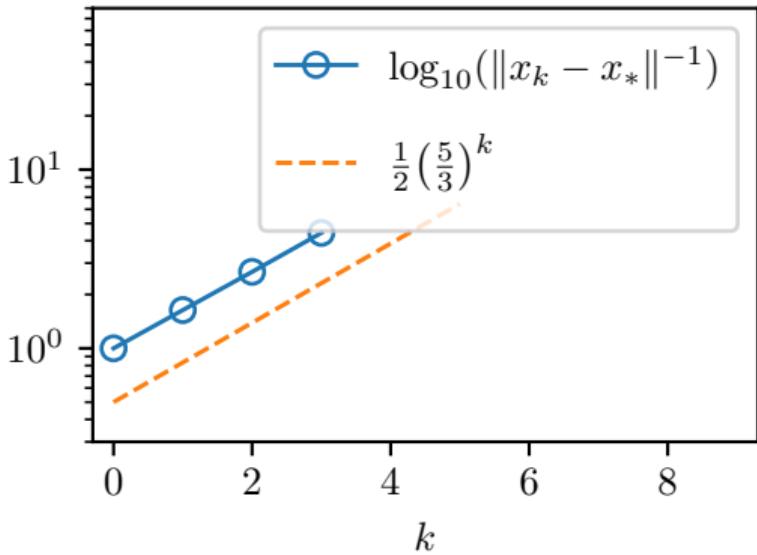
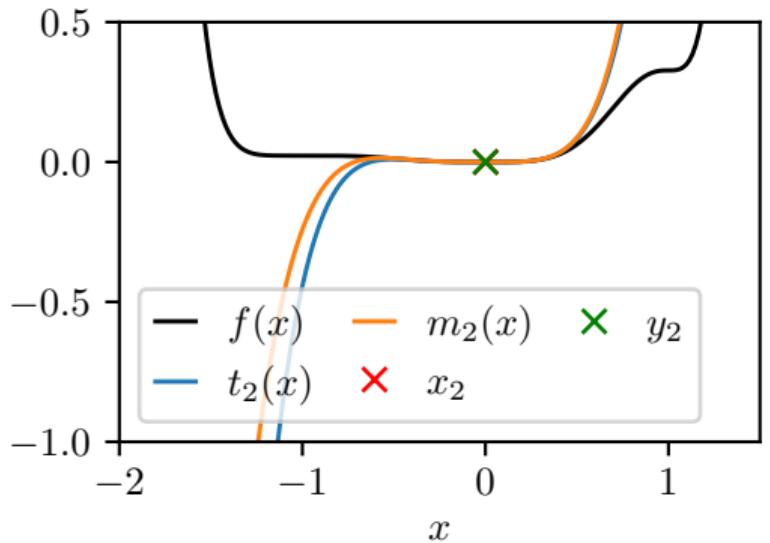
AR5, adaptive σ_k and right model minimizer



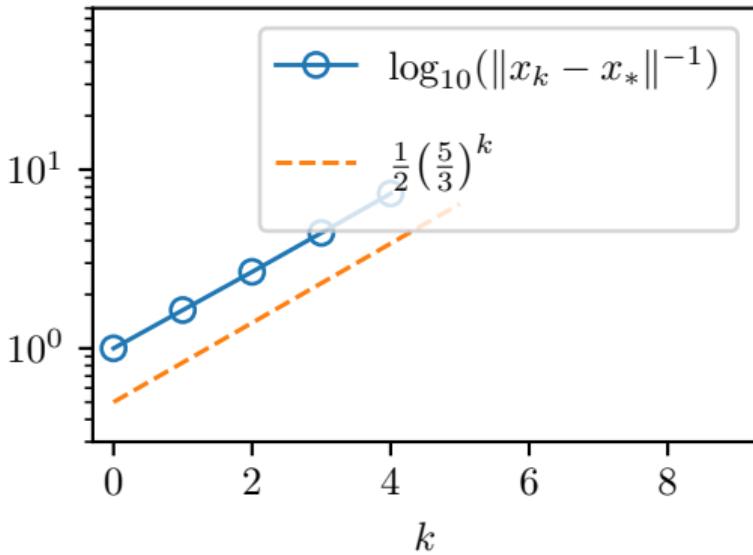
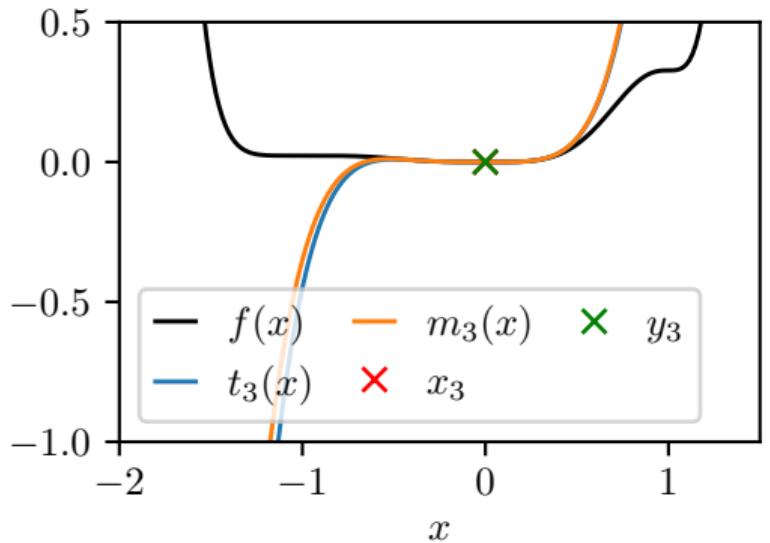
AR5, adaptive σ_k and right model minimizer



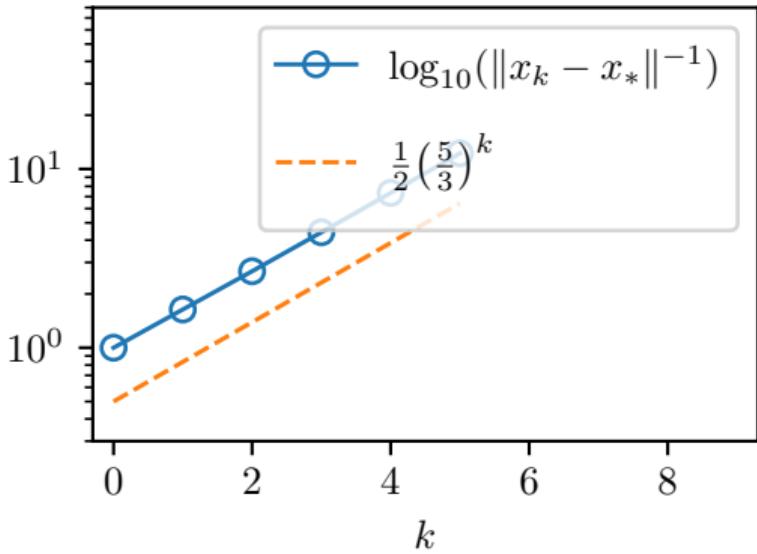
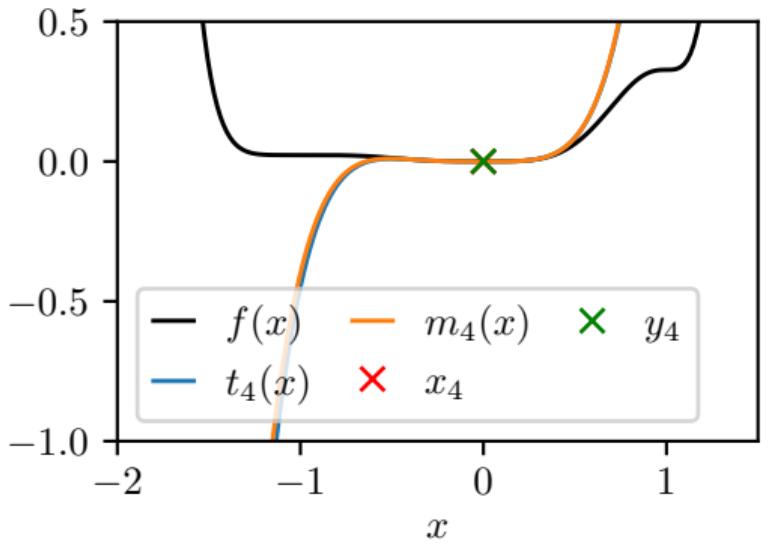
AR5, adaptive σ_k and right model minimizer



AR5, adaptive σ_k and right model minimizer



AR5, adaptive σ_k and right model minimizer



Outline

1 Motivation

- Why higher-order methods?
- The AR p method

2 Numerical Illustrations

- Example with non-degenerate minimizer
- Example with degenerate minimizer

3 Theoretical Results

4 Conclusion

Lemma

Let \mathbf{x}_* be a local minimizer of f , let $\nabla^p f$ be Lipschitz continuous and let f be uniformly convex of order q inside $\mathcal{B}(\mathbf{x}_*, r_\mu)$. Assume that $\mathbf{x}_k, \mathbf{y}_k \in \mathcal{B}(\mathbf{x}_*, r_\mu)$ and that iteration k is successful, then the gradients satisfy

$$\|\nabla f(\mathbf{y}_k)\| \leq (L_p/p! + (p+1)\sigma_k) \left(\frac{q}{\mu_q} \right)^{\frac{p}{q-1}} \|\nabla f(\mathbf{x}_k)\|^{\frac{p}{q-1}}.$$

Convergence for successful iterations

Theorem

Let \mathbf{x}_* be a local minimizer of f , let $\nabla^p f$ be Lipschitz continuous and let f be uniformly convex of order q inside $\mathcal{B}(\mathbf{x}_*, r_\mu)$. Assume that all iterates stay inside $\mathcal{B}(\mathbf{x}_*, r_\mu)$, that \mathbf{x}_0 is close enough to \mathbf{x}_* and that $p > q - 1$, then

$$\begin{aligned}\|\nabla f(\mathbf{x}_{k_i})\| \rightarrow 0 && \text{at a } Q\text{-}\frac{p}{q-1}\text{th-order rate} \\ f(\mathbf{x}_{k_i}) \rightarrow f(\mathbf{x}_*) && \text{at an } R\text{-}\frac{p}{q-1}\text{th-order rate} \\ \mathbf{x}_{k_i} \rightarrow \mathbf{x}_* && \text{at an } R\text{-}\frac{p}{q-1}\text{th-order rate}\end{aligned}$$

where $\{k_1, k_2, \dots\}$ are the successful iterations.

Theorem

Let \mathbf{x}_* be a local minimizer of f , let $\nabla^p f$ be Lipschitz continuous and let f be uniformly convex of order q inside $\mathcal{B}(\mathbf{x}_*, r_\mu)$. Assume that all iterates stay inside $\mathcal{B}(\mathbf{x}_*, r_\mu)$, that \mathbf{x}_0 is close enough to \mathbf{x}_* and that $p > q - 1$, then

$$\|\nabla f(\mathbf{x}_k)\| \rightarrow 0 \quad \text{at an } R^{-\sqrt{\frac{p}{q-1}} \text{ th-order rate}}$$

$$f(\mathbf{x}_k) \rightarrow f(\mathbf{x}_*) \quad \text{at an } R^{-\sqrt{\frac{p}{q-1}} \text{ th-order rate}}$$

$$\mathbf{x}_k \rightarrow \mathbf{x}_* \quad \text{at an } R^{-\sqrt{\frac{p}{q-1}} \text{ th-order rate}}$$

where $\gamma_1 = \gamma_2^{-\alpha}$. ($\alpha = 0$ and $\alpha = 1$ in experiments)

Lemma

Let \mathbf{x}_* be a local minimizer of f , let $\nabla^p f$ be Lipschitz continuous and let f be uniformly convex of order q inside $\mathcal{B}(\mathbf{x}_*, r_\mu)$. If $p > q - 1$, there exists radii $r_x, r_y > 0$ with $r_x \leq r_y \leq r_\mu$ such that for any $\mathbf{x}_k \in \mathcal{B}(\mathbf{x}_*, r_x)$ there exists a local minimizer of m_k inside $\mathcal{B}(\mathbf{x}_*, r_y)$ and any such minimizer will give a successful iteration.

Theorem

Let \mathbf{x}_* be a local minimizer of f , let $\nabla^p f$ be Lipschitz continuous, let f be uniformly convex of order q inside $\mathcal{B}(\mathbf{x}_*, r_\mu)$ and let r_x and r_y be chosen according to the previous lemma. Assume that $\mathbf{y}_k \in \mathcal{B}(\mathbf{x}_*, r_y)$ in every iteration, that \mathbf{x}_0 is close enough to \mathbf{x}_* and that $p > q - 1$, then all iterations are successful and

$$\begin{aligned}\|\nabla f(\mathbf{x}_k)\| &\rightarrow 0 && \text{at a } Q\text{-}\frac{p}{q-1}\text{th-order rate} \\ f(\mathbf{x}_k) &\rightarrow f(\mathbf{x}_*) && \text{at an } R\text{-}\frac{p}{q-1}\text{th-order rate} \\ \mathbf{x}_k &\rightarrow \mathbf{x}_* && \text{at an } R\text{-}\frac{p}{q-1}\text{th-order rate.}\end{aligned}$$

Comparison with Doikov and Nesterov, 2022

Doikov, N., Nesterov, Y. Local convergence of tensor methods. *Math. Program.* 193, 315–336 (2022). <https://doi.org/10.1007/s10107-020-01606-x>

■ Assumptions:

- $\nabla^p f$ Lipschitz continuous
- f uniformly convex of order q
- $p > q - 1$
- σ_k constant and $\geq \frac{pL_p}{(p+1)!}$
- $\|\nabla f(\mathbf{x}_k)\|$ and $f(\mathbf{x}_k)$ converge with $\frac{p}{q-1}$ th-order rate
- Also covers composite optimization case with convex non-differentiable component

■ Assumptions

- $\nabla^p f$ Lipschitz continuous
- f locally uniformly convex of order q
- $p > q - 1$
- σ_k adaptively chosen
- $\|\nabla f(\mathbf{x}_k)\|$, $f(\mathbf{x}_k)$ and \mathbf{x}_k converge with $\sqrt[1+\alpha]{\frac{p}{q-1}}$ th-order rate
- $\|\nabla f(\mathbf{x}_k)\|$, $f(\mathbf{x}_k)$ and \mathbf{x}_k converge with $\frac{p}{q-1}$ th-order rate when choosing the “right” minimizer

Outline

1 Motivation

- Why higher-order methods?
- The AR p method

2 Numerical Illustrations

- Example with non-degenerate minimizer
- Example with degenerate minimizer

3 Theoretical Results

4 Conclusion

Conclusion

- For tensor methods ($p > 3$) choosing the right model minimizer is crucial
- Tensor methods achieve $\frac{p}{q-1}$ th-order rates in theory and experiments
- Tensor methods achieve superlinear convergence even for degenerate minimizers