

Local convergence of adaptively regularized tensor methods

KARL WELZEL
Mathematical Institute
University of Oxford

EUROPT 2025, 1 July 2025



Oxford
Mathematics



Mathematical
Institute

Local convergence of adaptive tensor methods

Collaboration with



Raphael Hauser



Yang Liu



Coralia Cartis

The AR p method

- Unconstrained minimization $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$, $f: \mathbb{R}^n \rightarrow \mathbb{R}$ smooth, non-convex
- Tensor methods use stronger oracles
 - Gradient descent uses $f(\mathbf{x}_k), \nabla f(\mathbf{x}_k)$
 - Newton's method uses $f(\mathbf{x}_k), \nabla f(\mathbf{x}_k), \nabla^2 f(\mathbf{x}_k)$
 - p th-order tensor methods use $f(\mathbf{x}_k), \nabla f(\mathbf{x}_k), \dots, \nabla^p f(\mathbf{x}_k)$
- AR p has optimal worst-case complexity, improves with p [Cartis, Gould, and Toint 2020]
- This talk: local convergence

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
```

Known results

Assumptions:

- $\nabla^p f$ is Lipschitz continuous: $\|\nabla^p f(\mathbf{x}) - \nabla^p f(\mathbf{y})\| \leq L_p \|\mathbf{x} - \mathbf{y}\|$
- f is uniformly convex of order q : $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{\mu_q}{q} \|\mathbf{y} - \mathbf{x}\|^q$

[Doikov and Nesterov 2022]: If $\sigma \geq \frac{pL_p}{(p+1)!}$ constant and $p > q - 1$ then

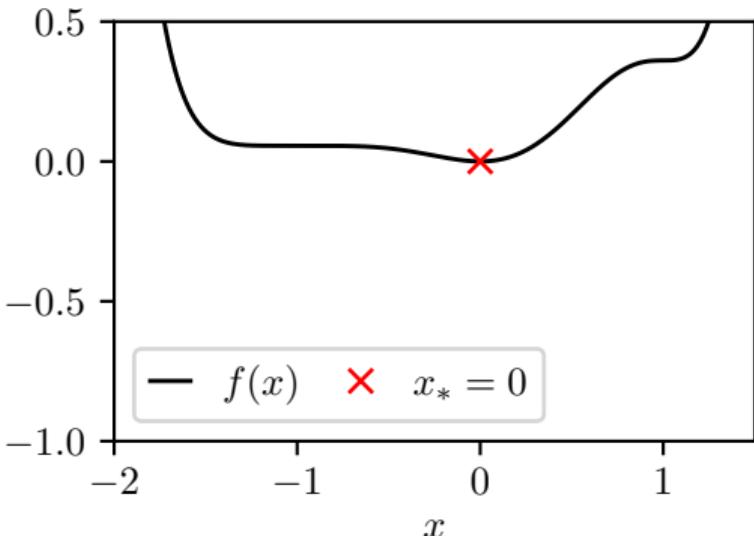
- $f(\mathbf{x}_k) \rightarrow f^*$ at a $\frac{p}{q-1}$ th-order rate and
- $\nabla f(\mathbf{x}_k) \rightarrow \mathbf{0}$ at a $\frac{p}{q-1}$ th-order rate.

[Cartis, Gould, and Toint 2022]: Using an adaptive σ_k , for $\nu = \frac{p+1}{p}$, $\mu = \frac{q}{q-1}$ the function values converge as follows (until $f(\mathbf{x}) - f^* < \varepsilon$)

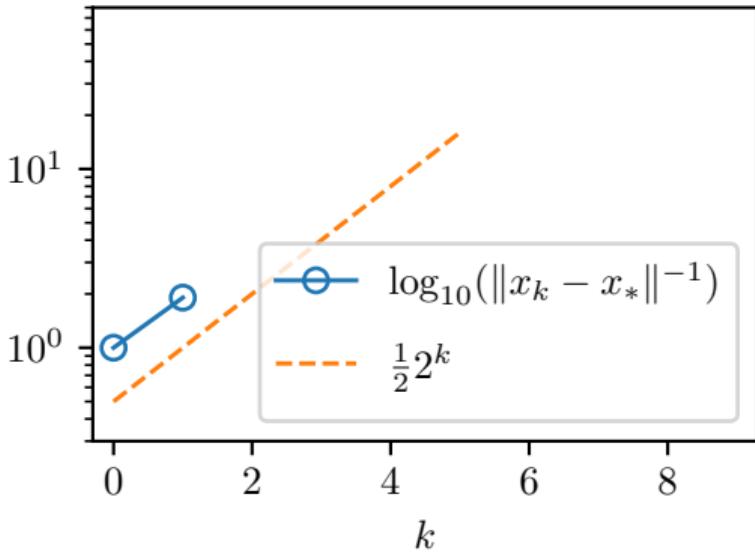
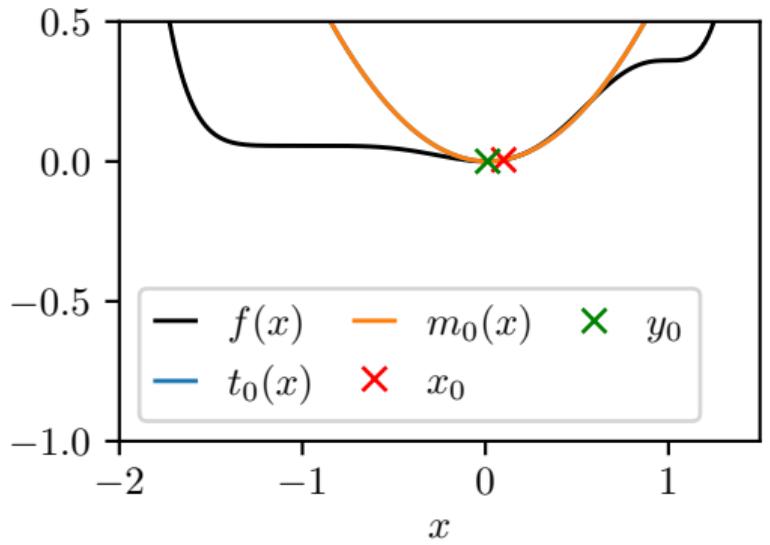
- $p < q - 1$ ($\nu > \mu$): sublinear, $O(\varepsilon^{-\frac{\nu-\mu}{\mu}})$ iterations
- $p = q - 1$ ($\nu = \mu$): linear, $O(\log(\varepsilon^{-1}))$ iterations
- $p > q - 1$ ($\nu < \mu$): superlinear, $O(\log(\log(\varepsilon^{-1})))$ iterations

Example with non-degenerate minimizer

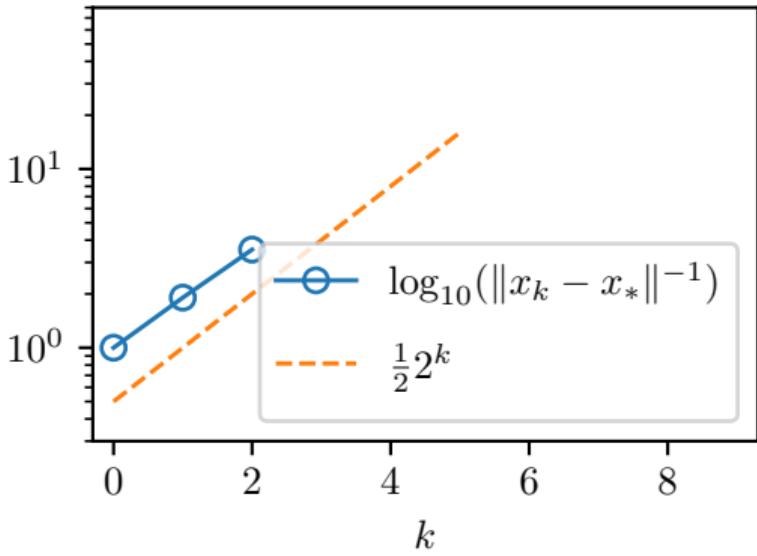
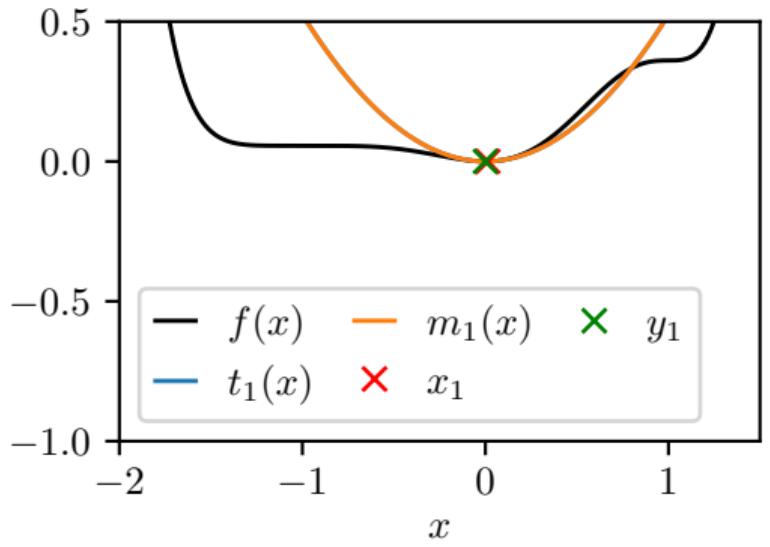
- $f(x) = \int_0^x (t+1)^4(t-1)^2 t \, dt$
- f has one local (and global) minimizer at $x_* = 0$
- f is not globally convex, but locally strongly convex: $f''(x_*) = 1$



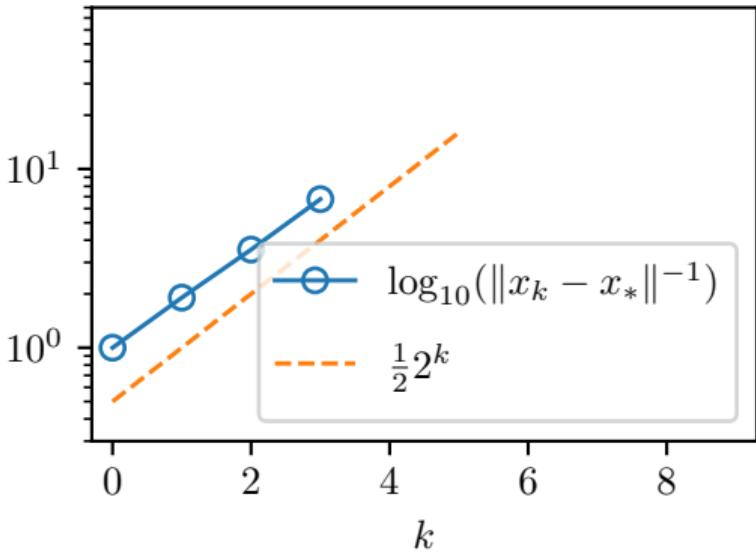
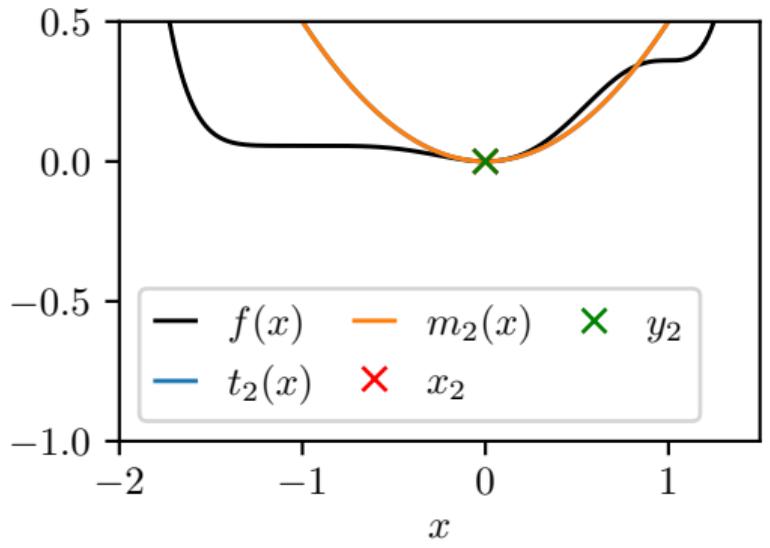
Newton's method ($p = 2$, $\sigma = 0$)



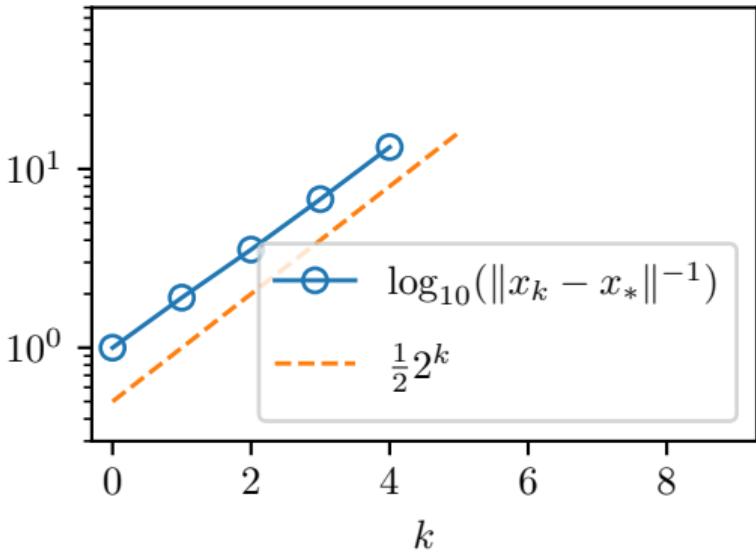
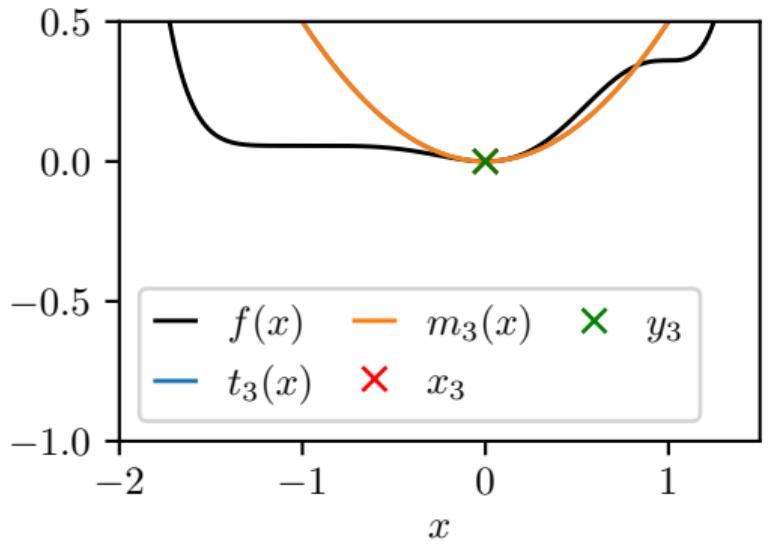
Newton's method ($p = 2$, $\sigma = 0$)



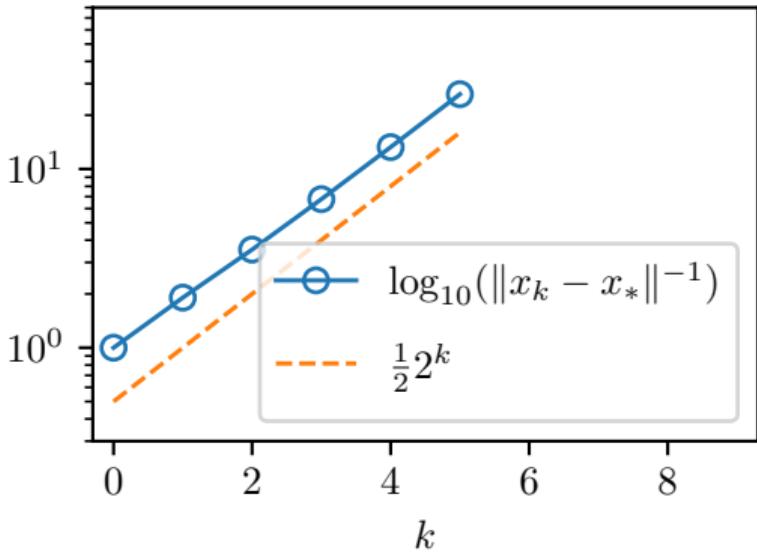
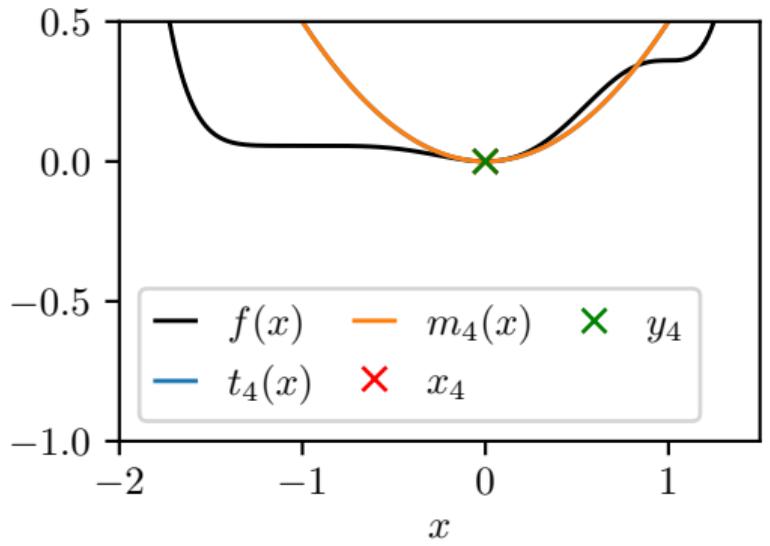
Newton's method ($p = 2$, $\sigma = 0$)



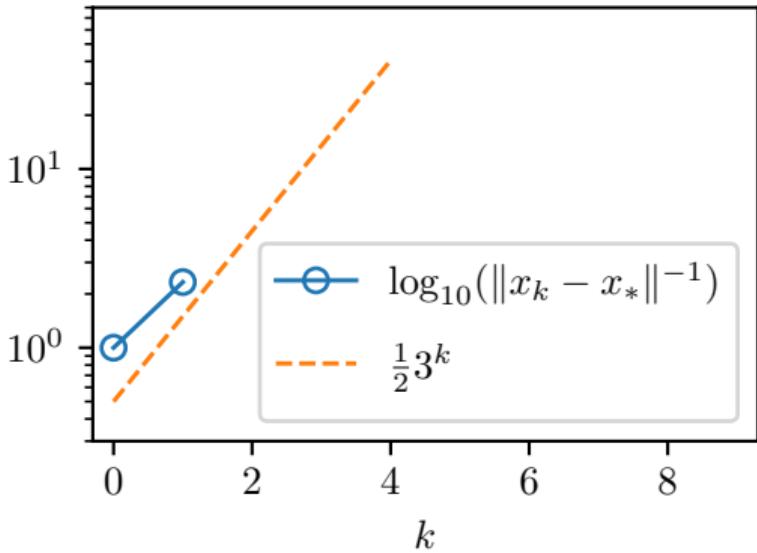
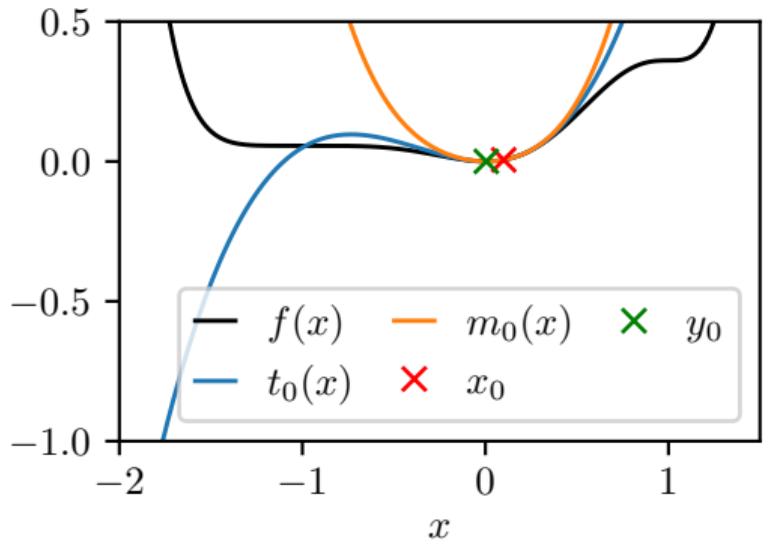
Newton's method ($p = 2$, $\sigma = 0$)

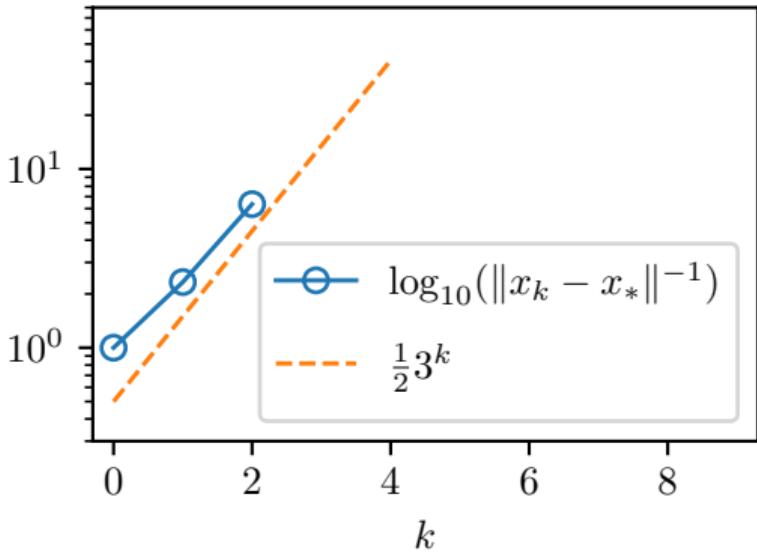
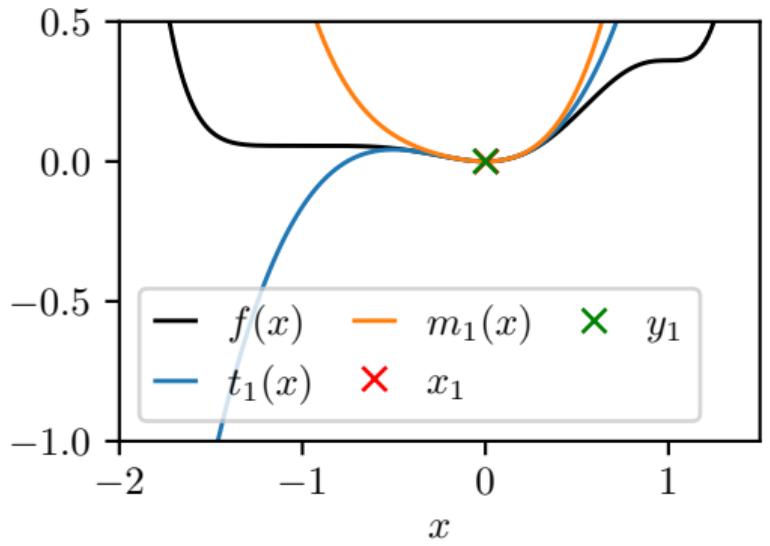


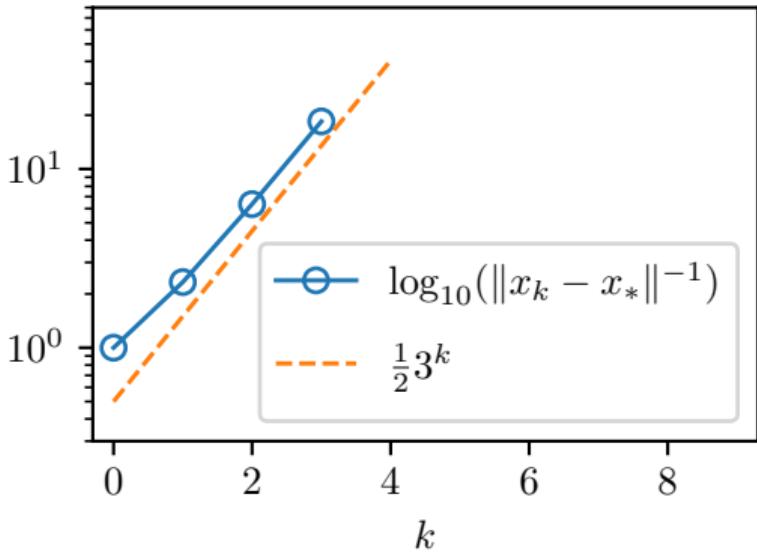
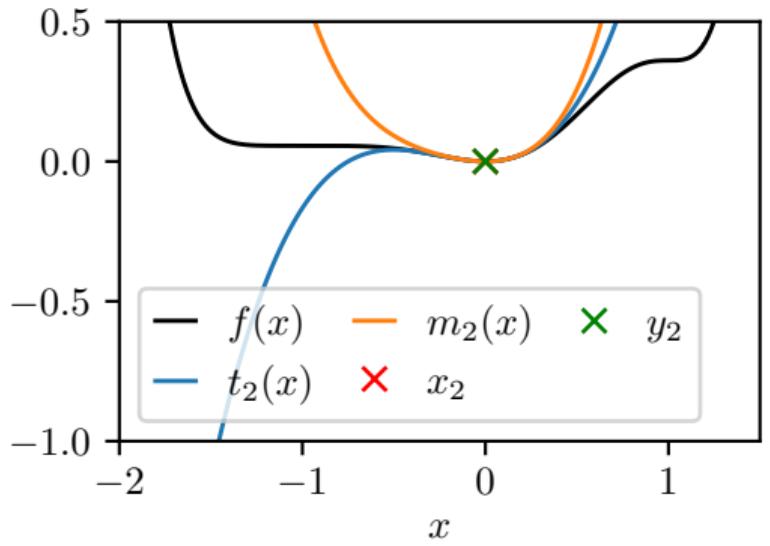
Newton's method ($p = 2, \sigma = 0$)

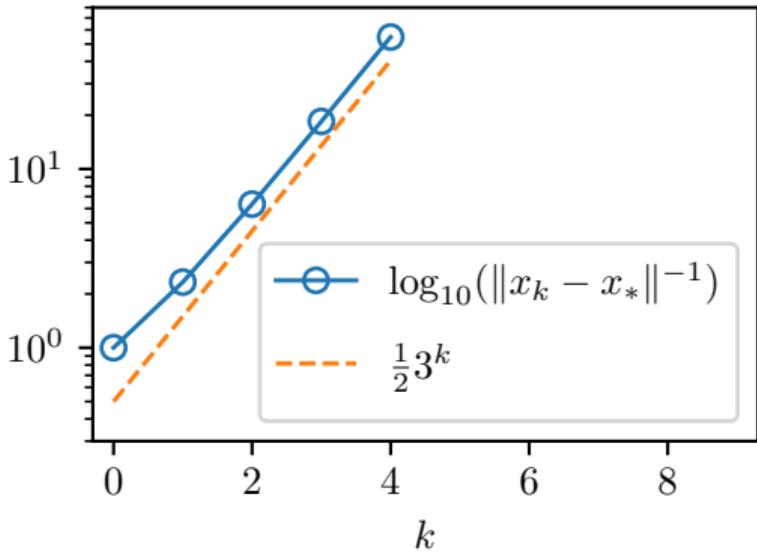
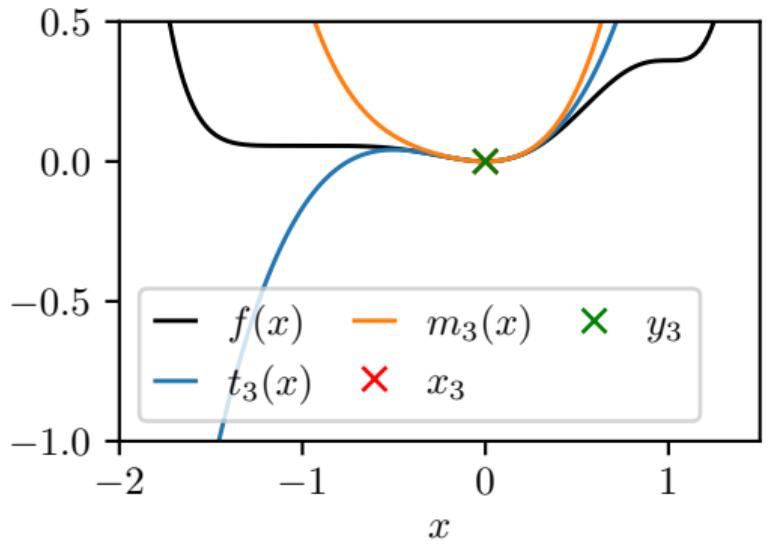


AR3, constant σ_k

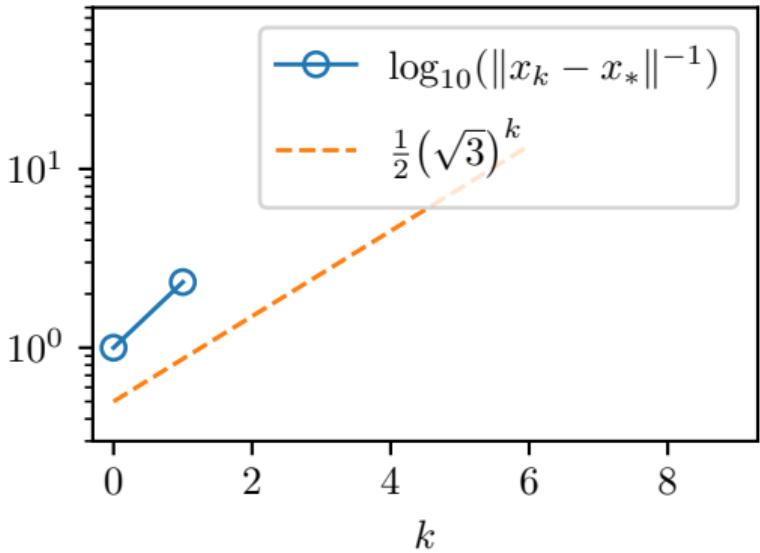
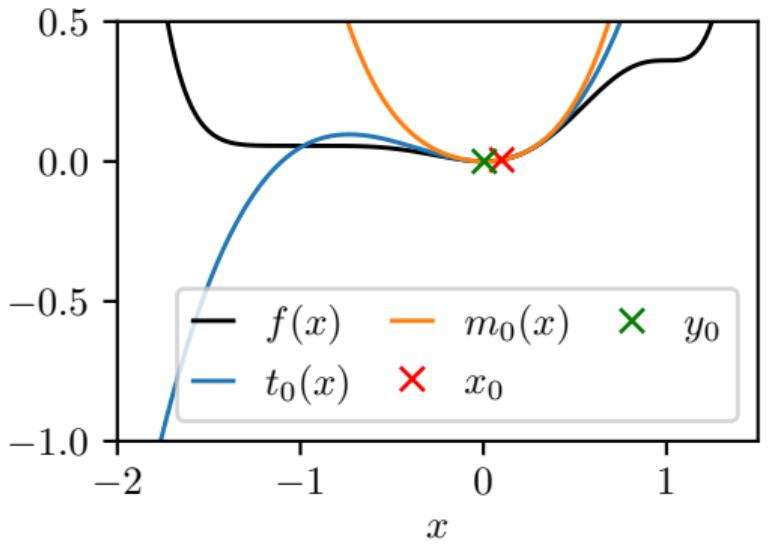




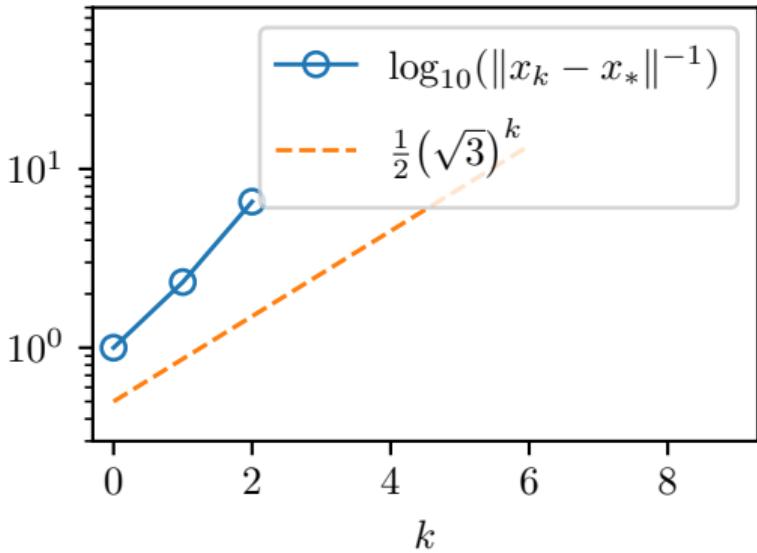
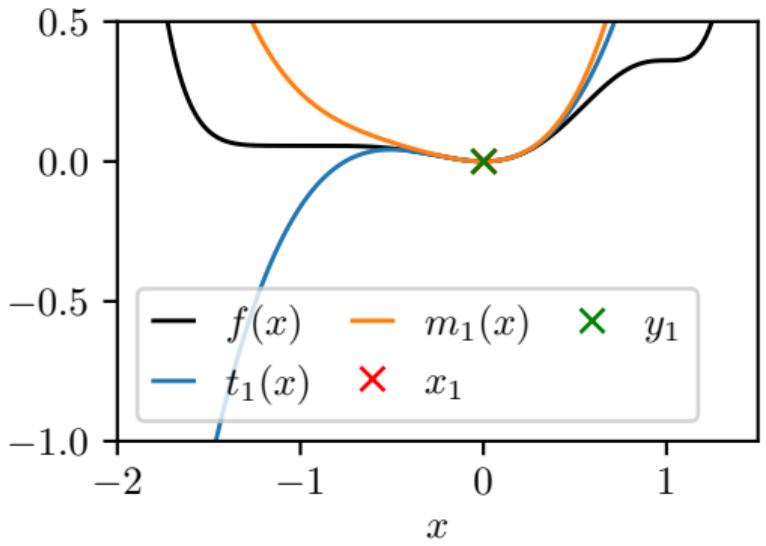




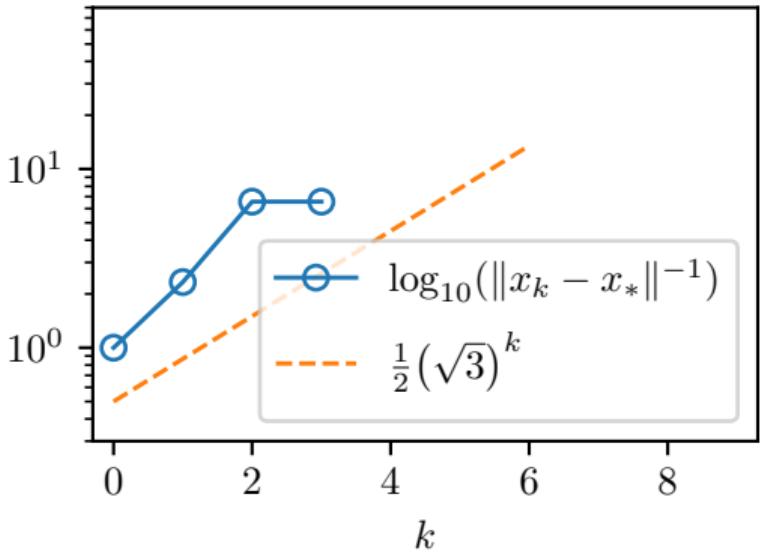
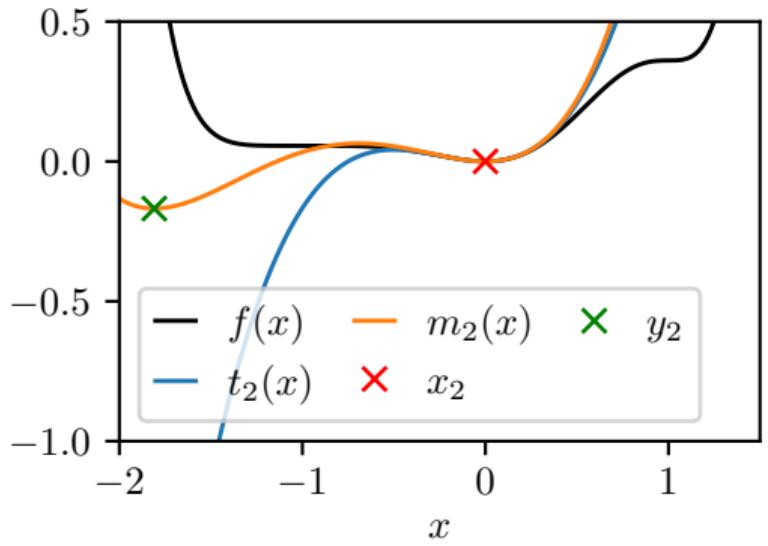
AR3, adaptive σ_k and global model minimizer



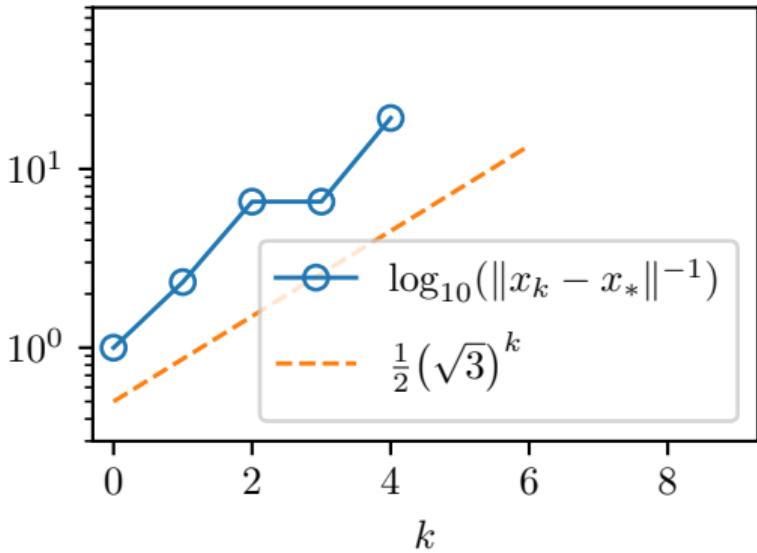
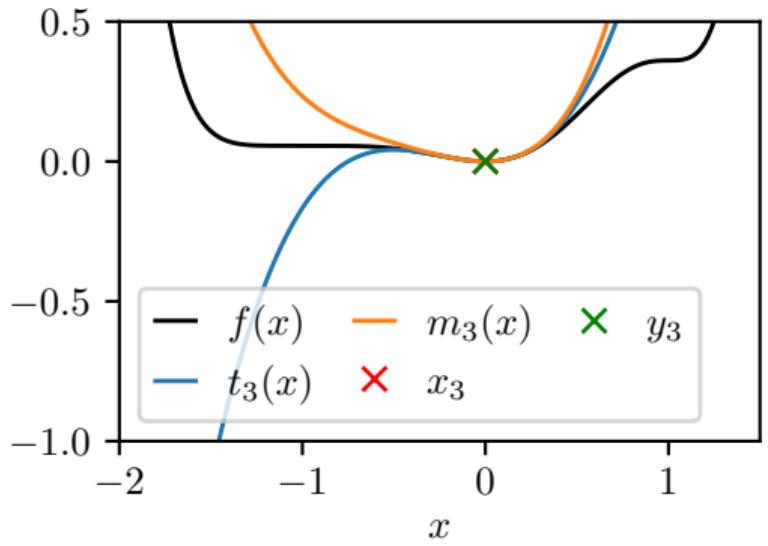
AR3, adaptive σ_k and global model minimizer



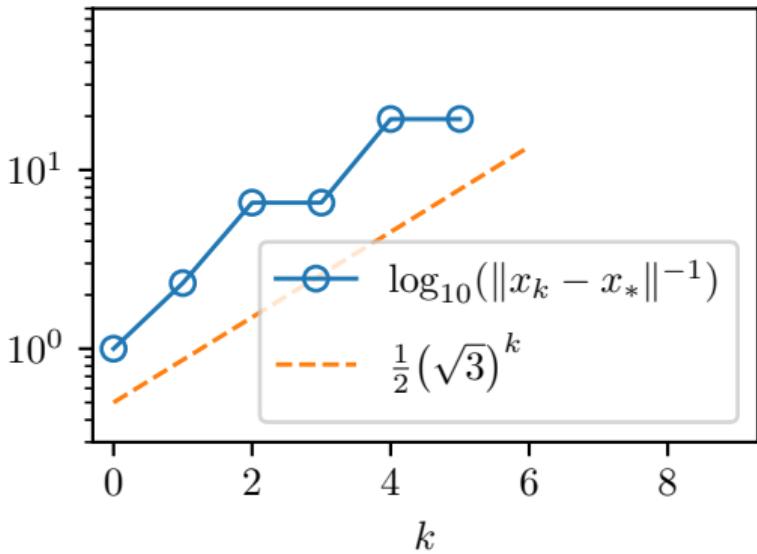
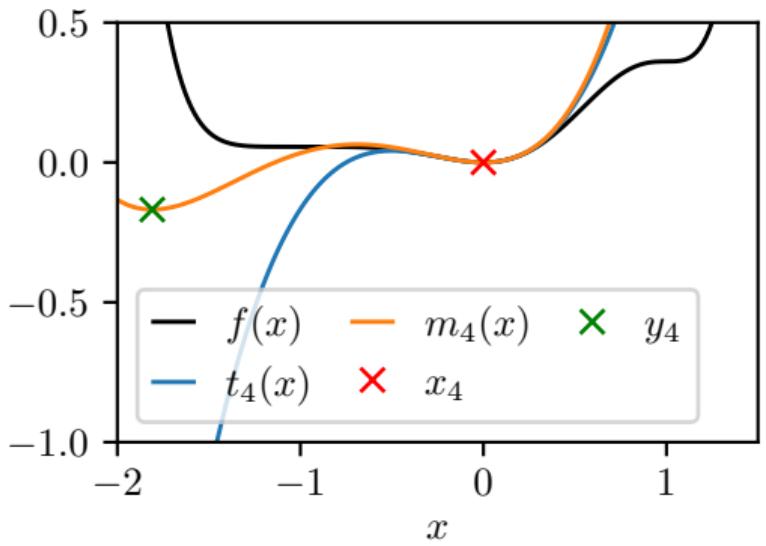
AR3, adaptive σ_k and global model minimizer



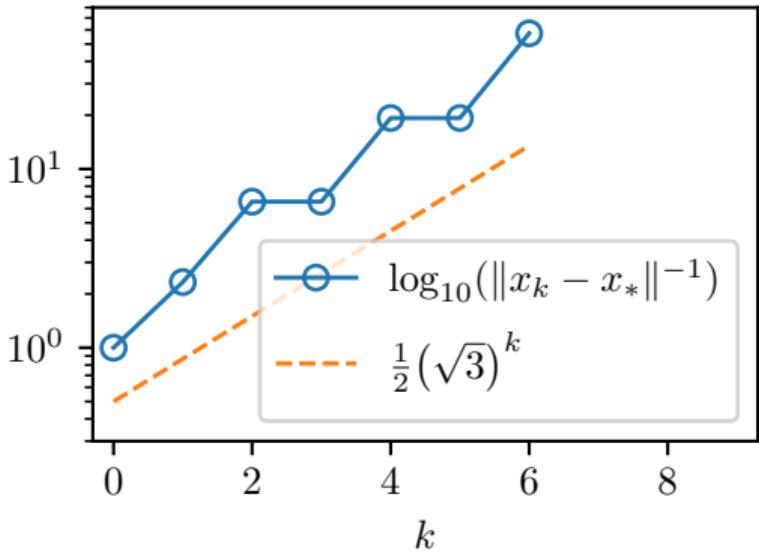
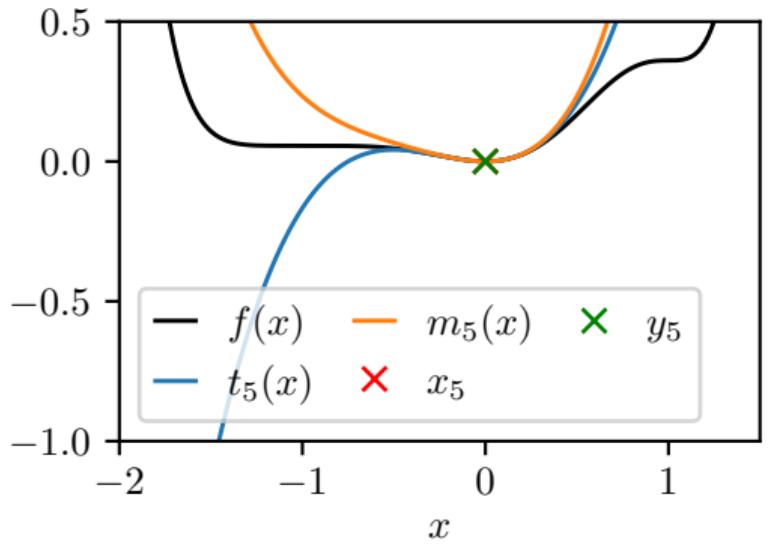
AR3, adaptive σ_k and global model minimizer



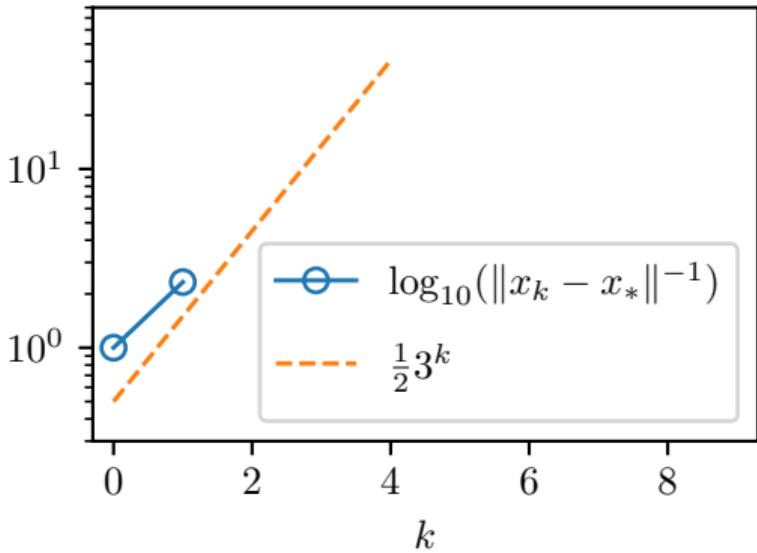
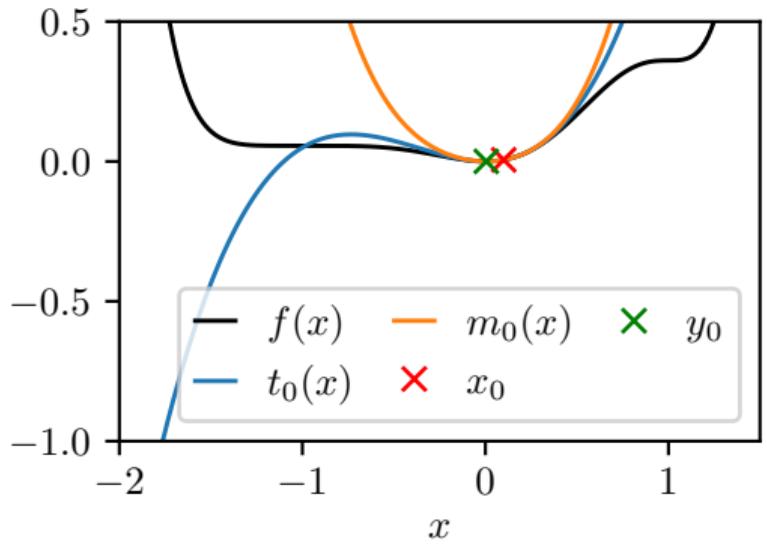
AR3, adaptive σ_k and global model minimizer



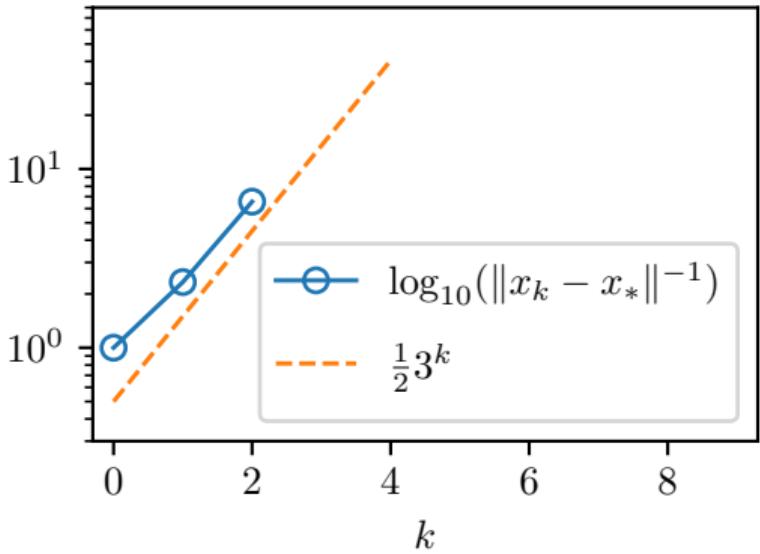
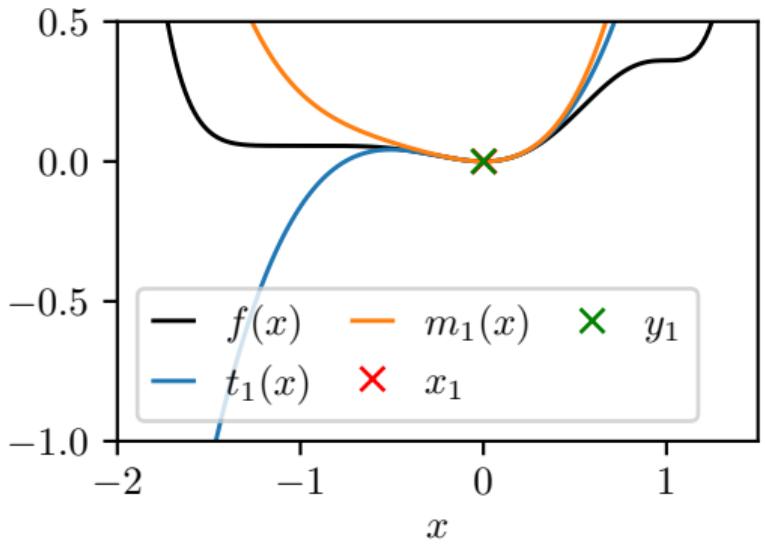
AR3, adaptive σ_k and global model minimizer



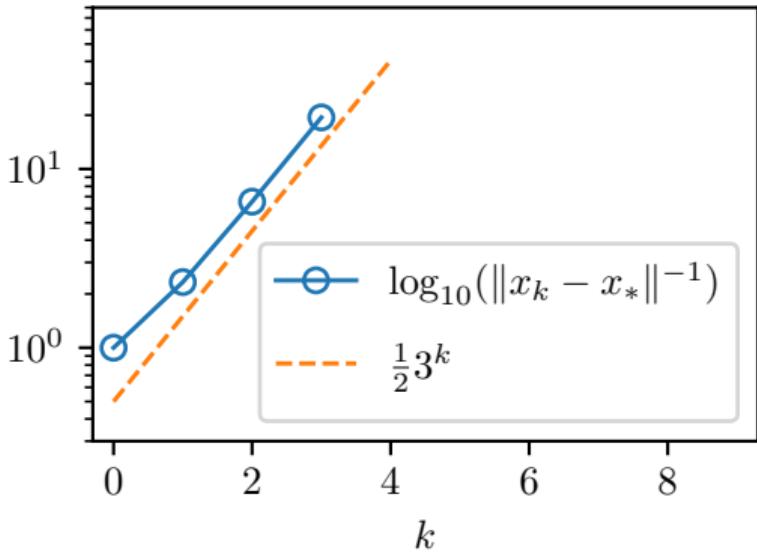
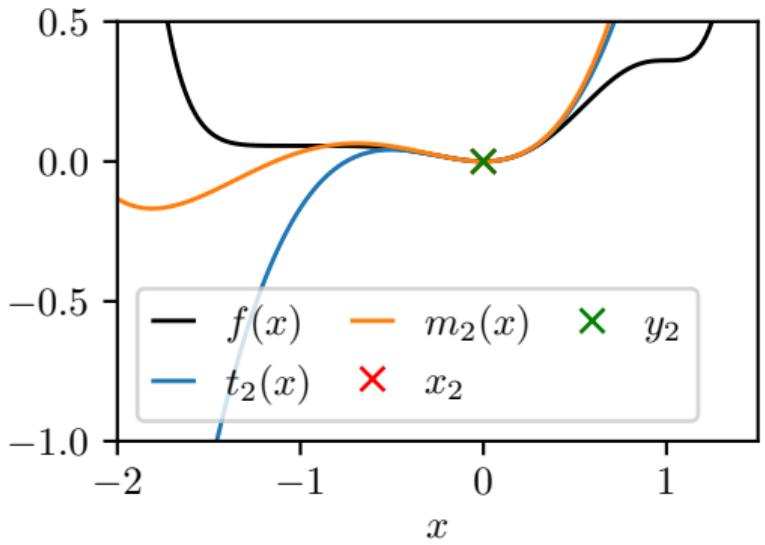
AR3, adaptive σ_k and right model minimizer



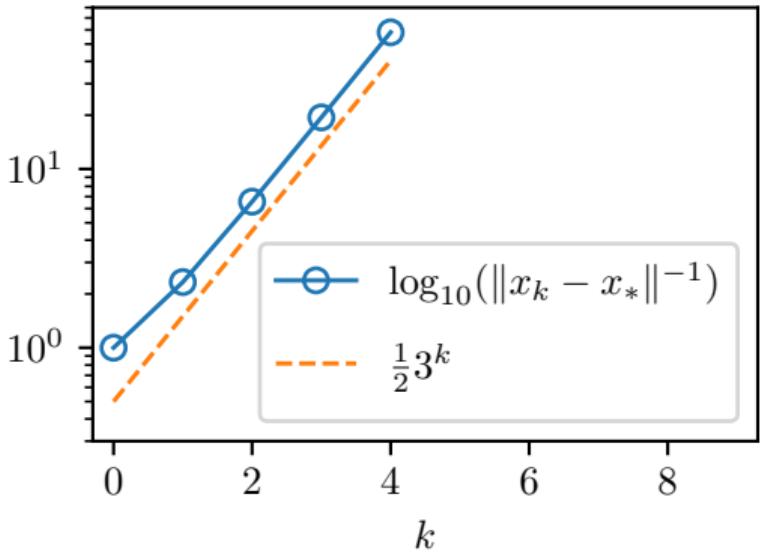
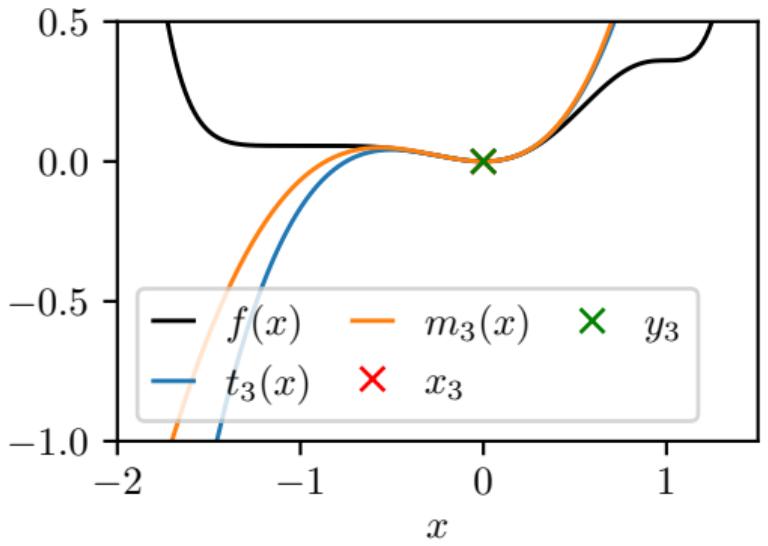
AR3, adaptive σ_k and right model minimizer



AR3, adaptive σ_k and right model minimizer

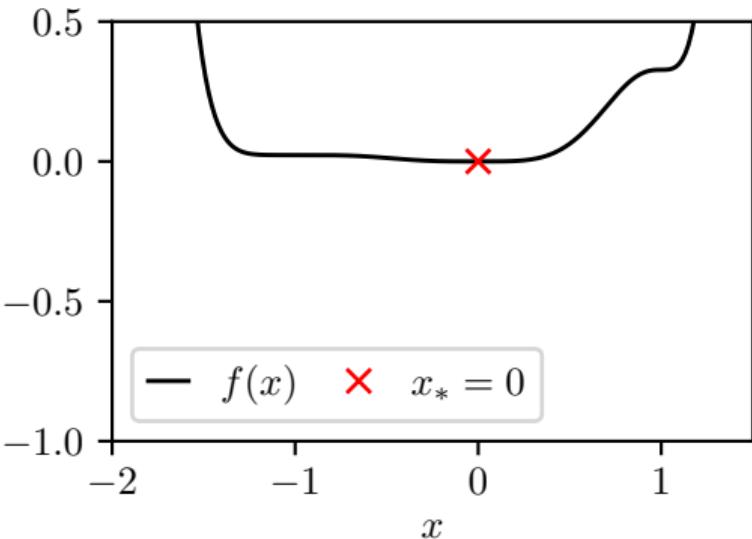


AR3, adaptive σ_k and right model minimizer

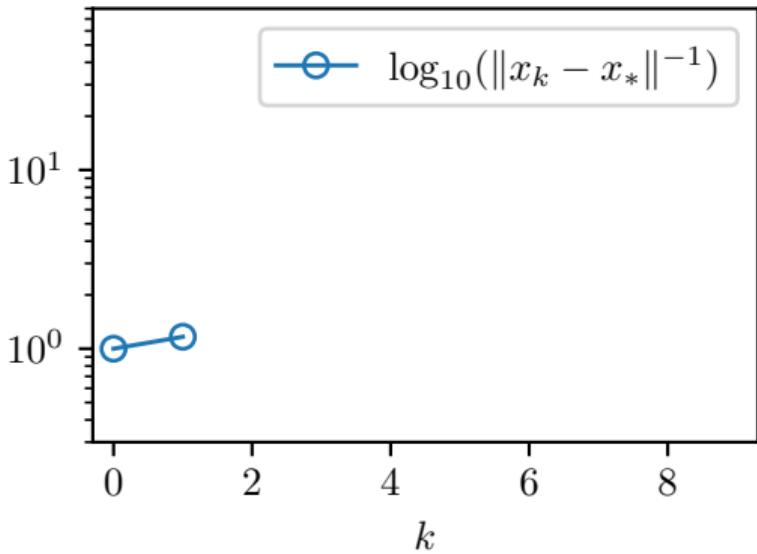
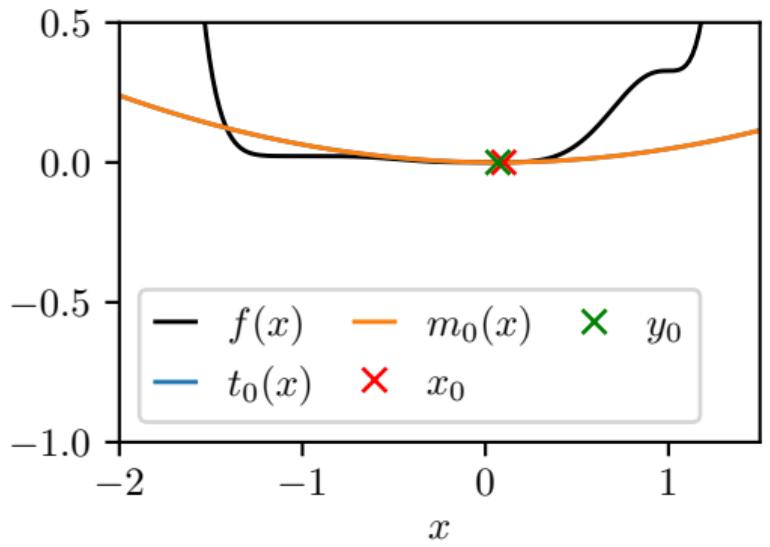


Example with degenerate minimizer

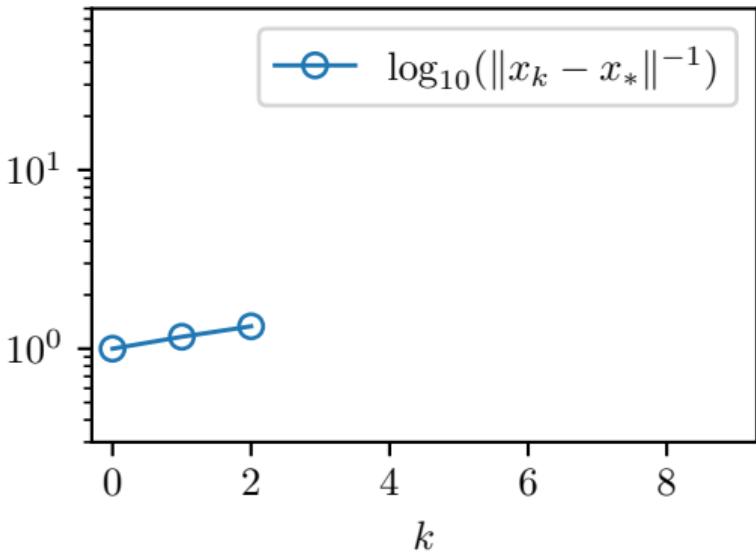
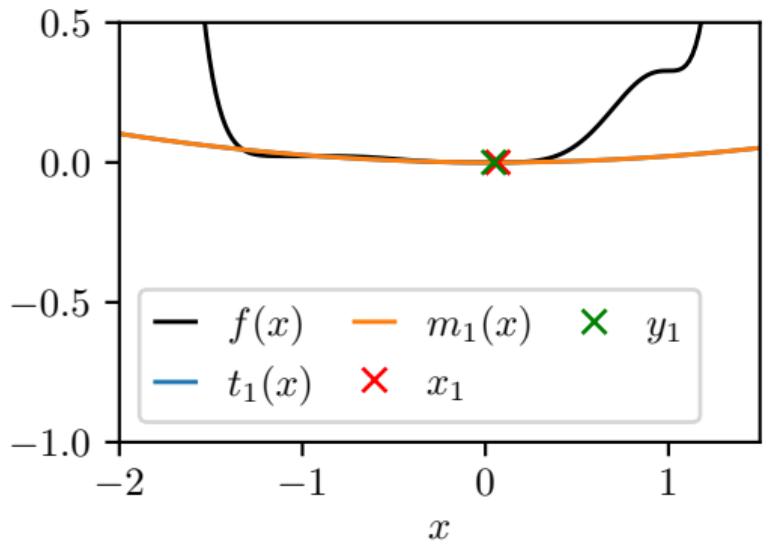
- $f(x) = \int_0^x 3(t+1)^4(t-1)^2t^3 dt$
- f has one local (and global) minimizer at $x_* = 0$
- f is locally uniformly convex around x_* with $q = 4$:
 - $f''(x_*) = 0$
 - $f'''(x_*) = 0$
 - $f''''(x_*) = 18$



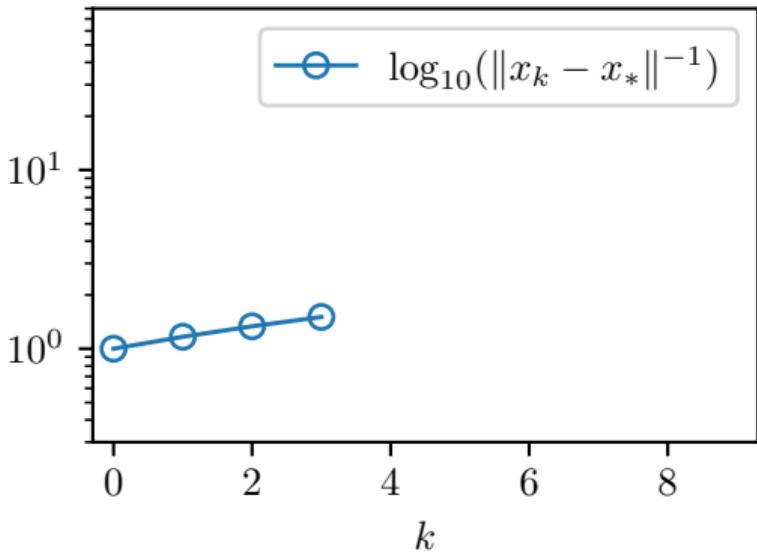
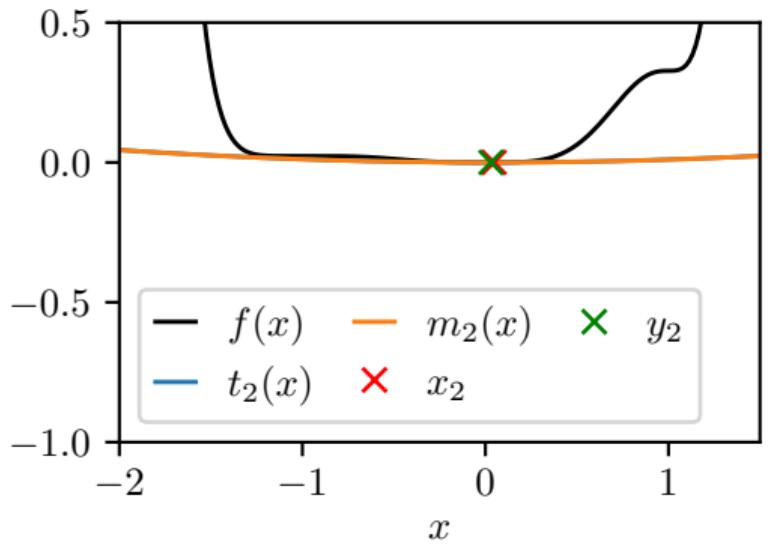
Newton's method ($p = 2$, $\sigma = 0$)



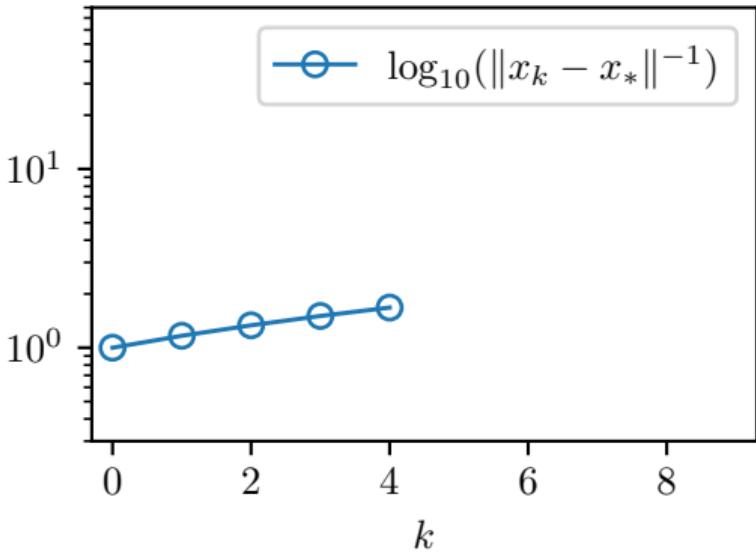
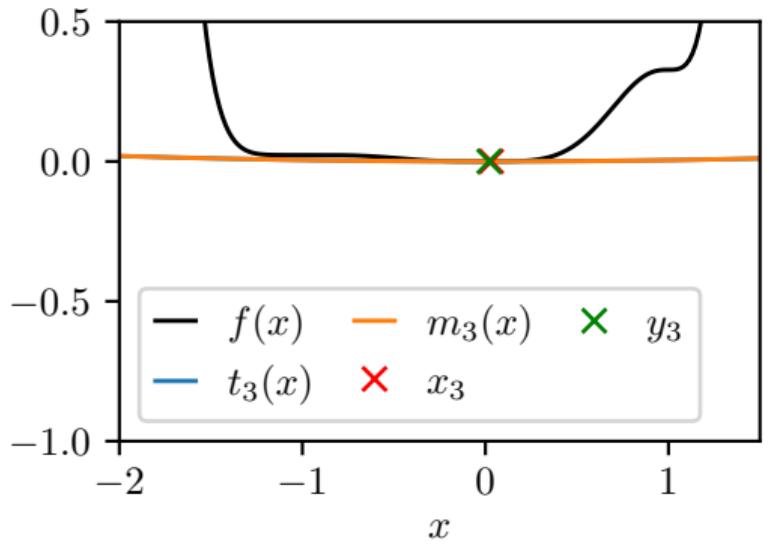
Newton's method ($p = 2$, $\sigma = 0$)



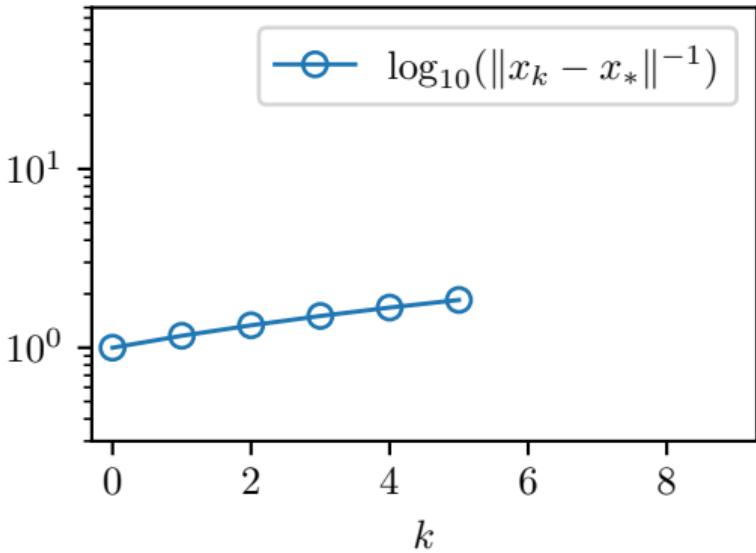
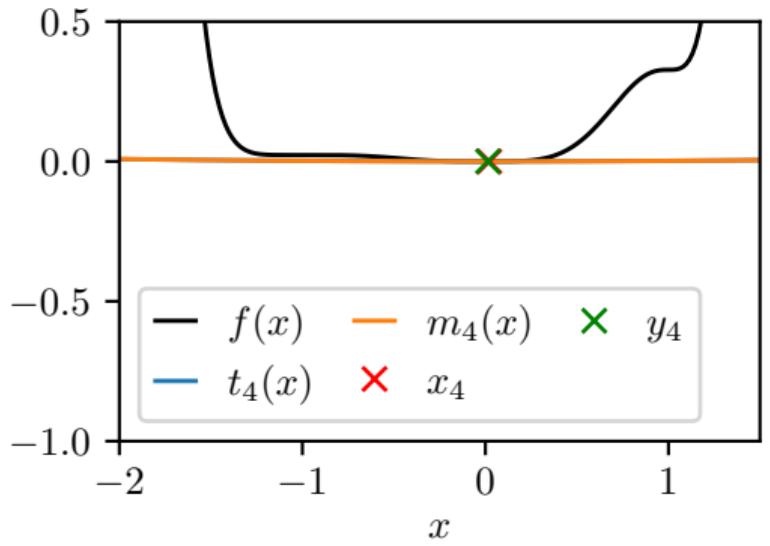
Newton's method ($p = 2$, $\sigma = 0$)



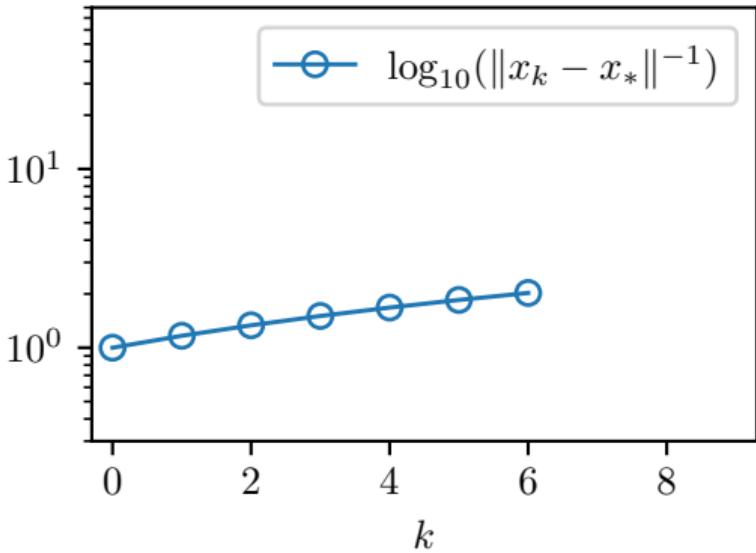
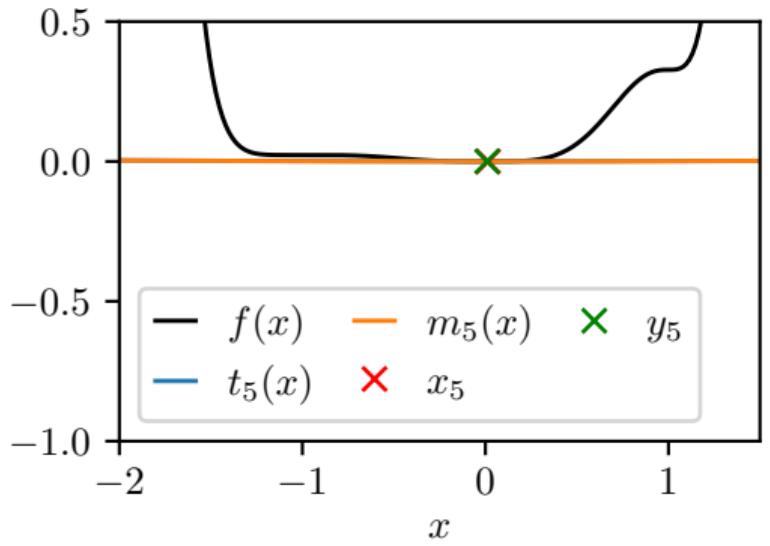
Newton's method ($p = 2$, $\sigma = 0$)



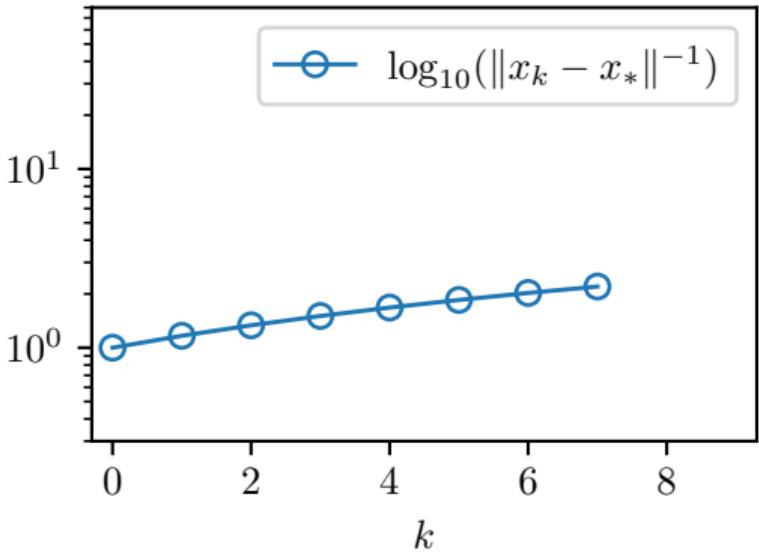
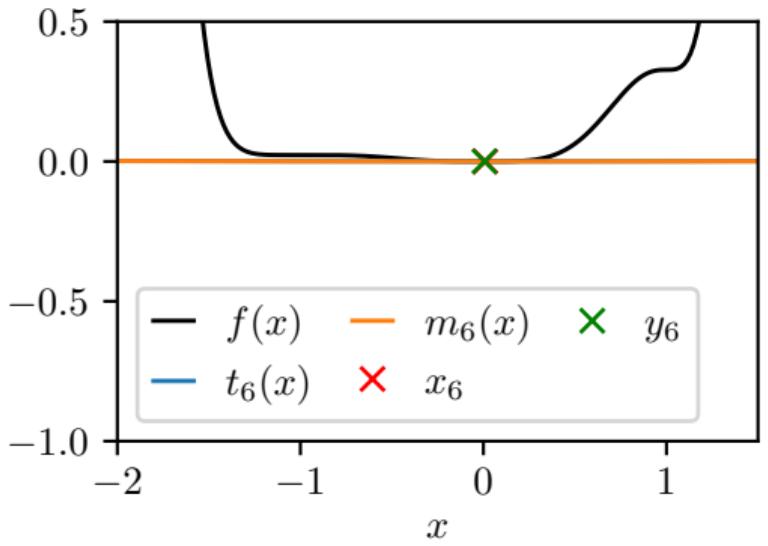
Newton's method ($p = 2$, $\sigma = 0$)



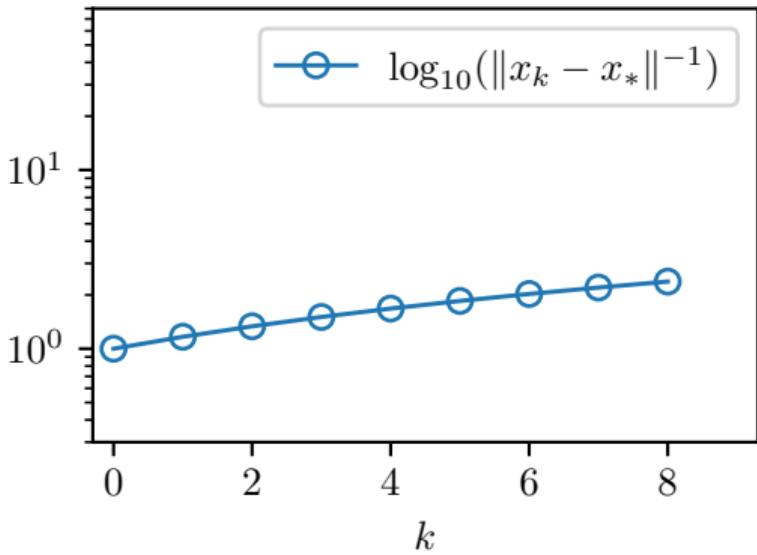
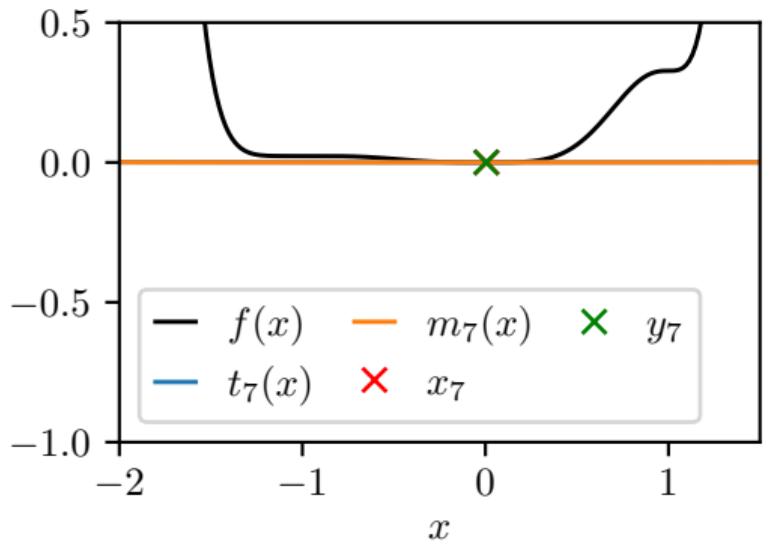
Newton's method ($p = 2, \sigma = 0$)



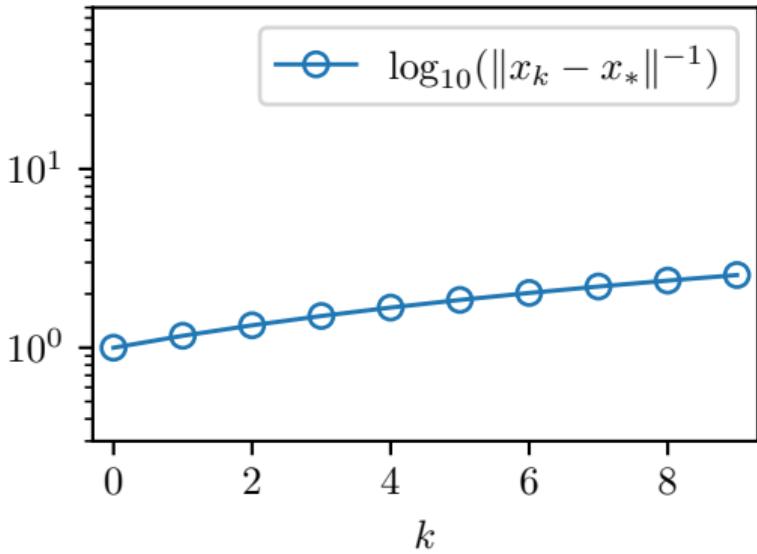
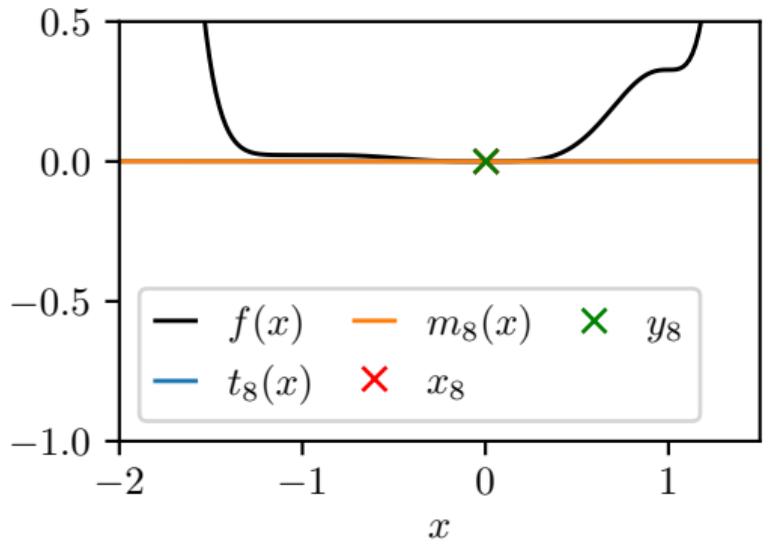
Newton's method ($p = 2$, $\sigma = 0$)



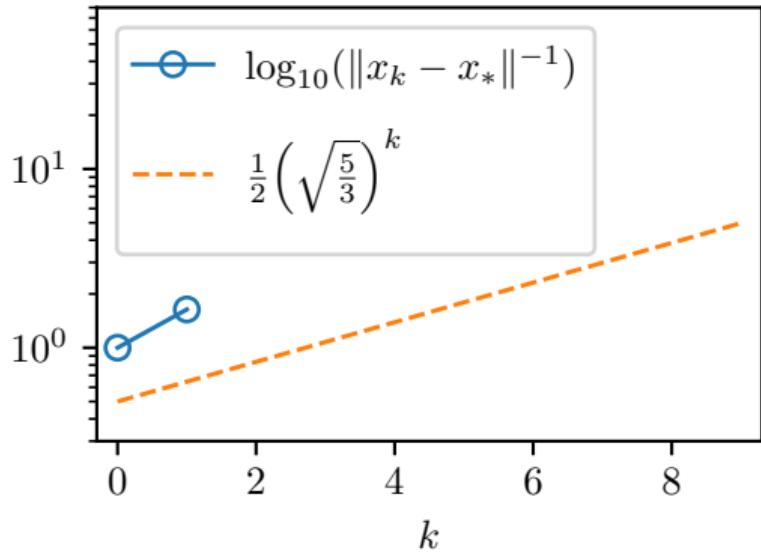
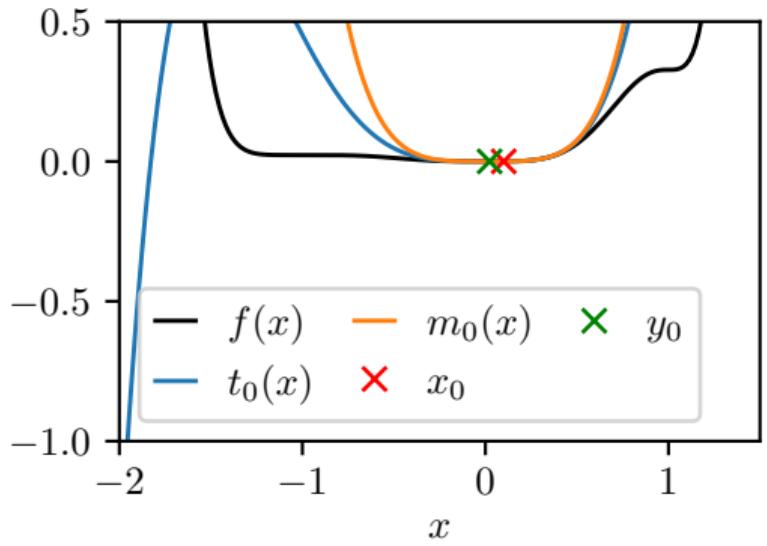
Newton's method ($p = 2$, $\sigma = 0$)



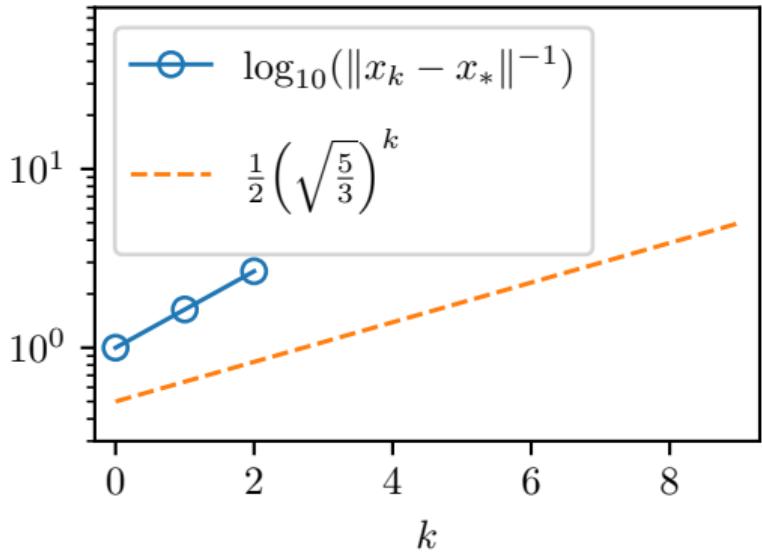
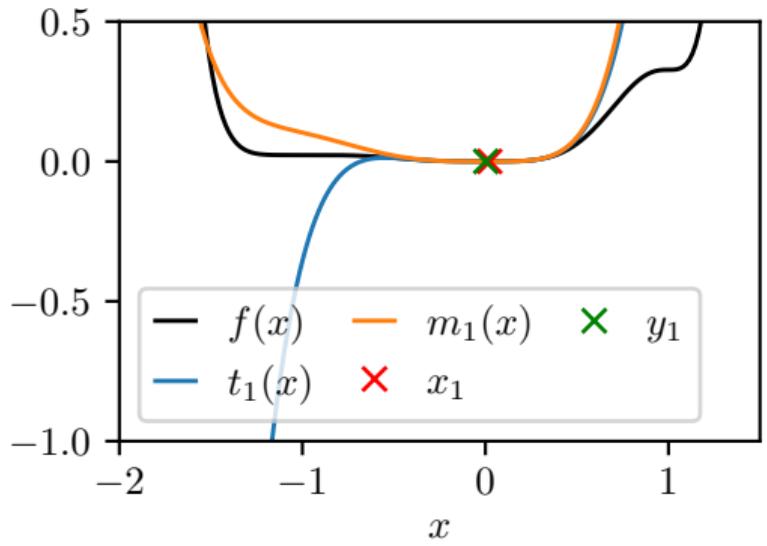
Newton's method ($p = 2, \sigma = 0$)



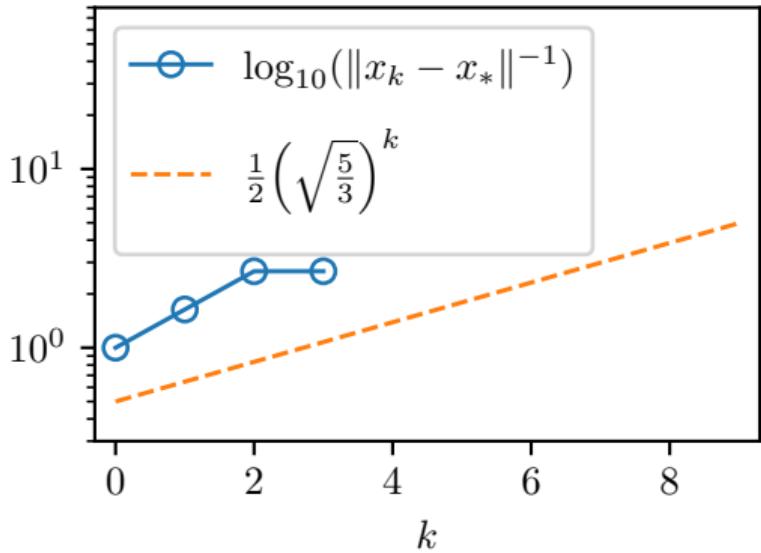
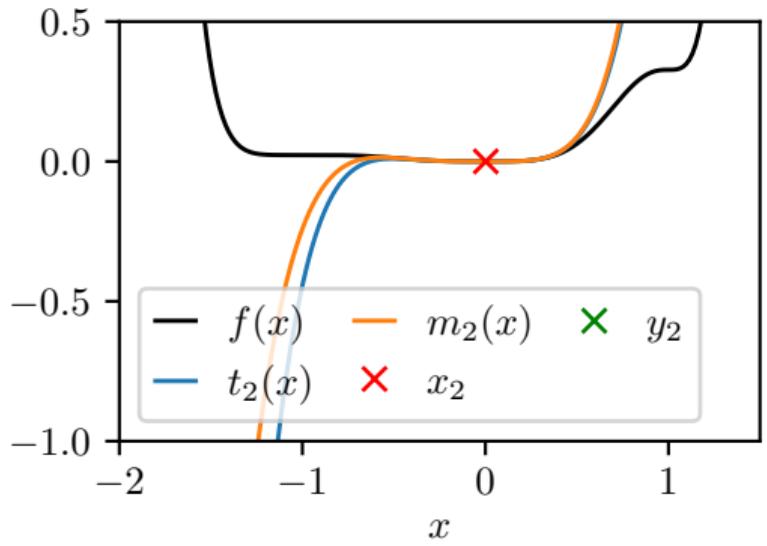
AR5, adaptive σ_k and global model minimizer



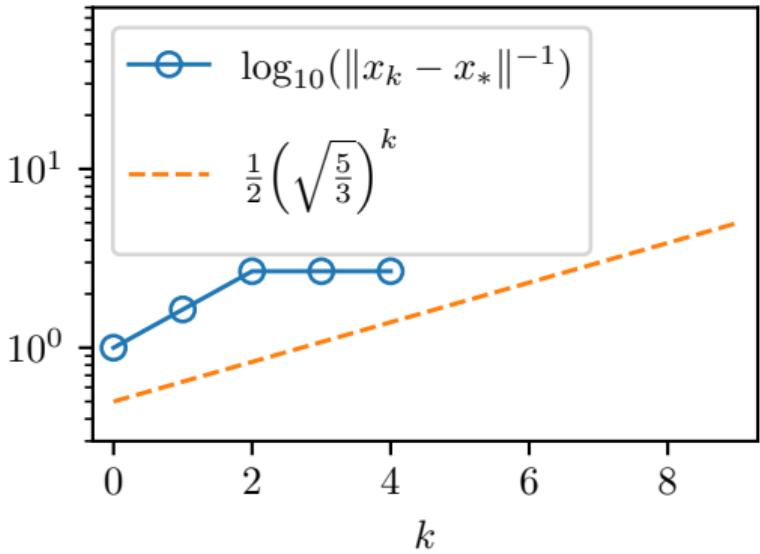
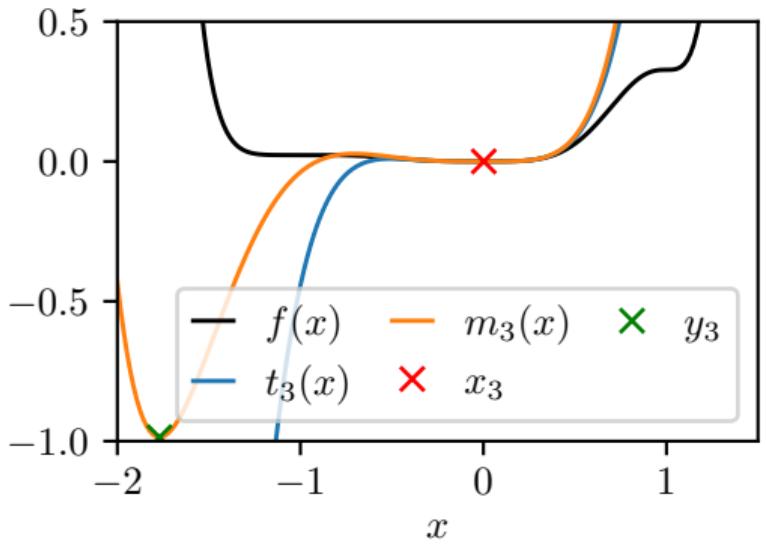
AR5, adaptive σ_k and global model minimizer



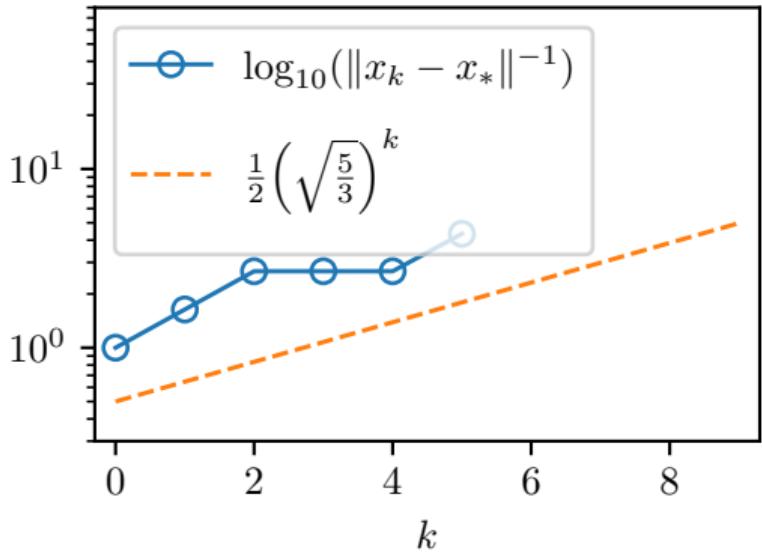
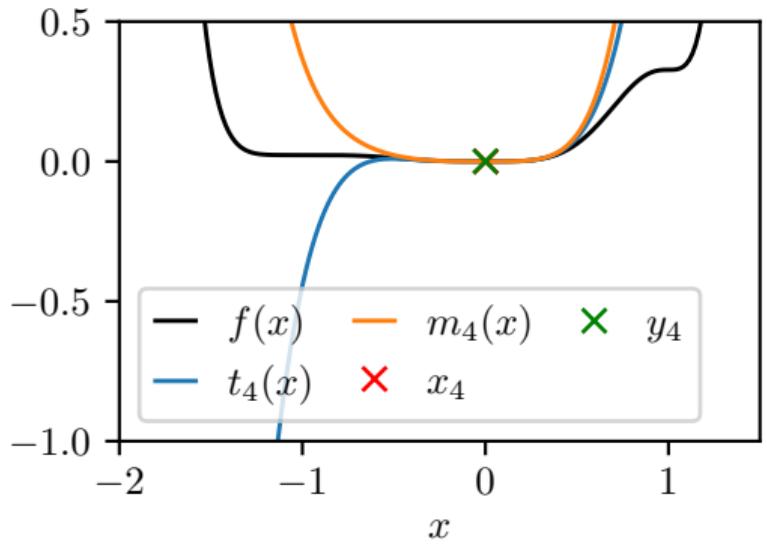
AR5, adaptive σ_k and global model minimizer



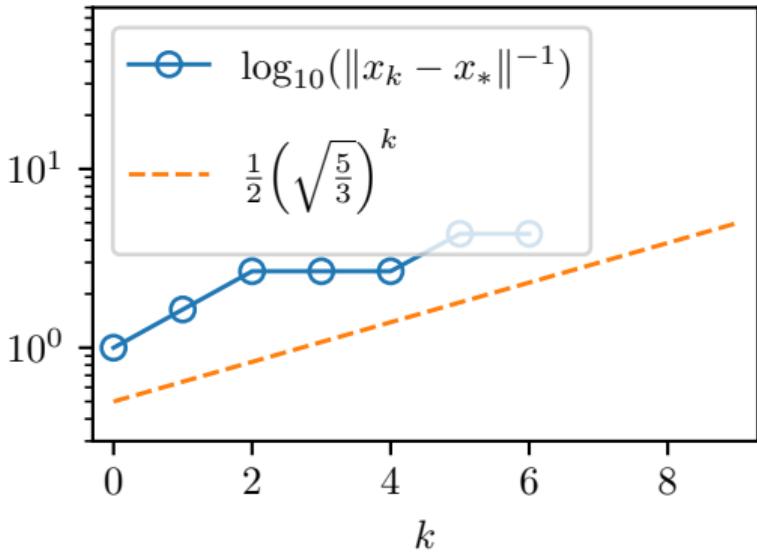
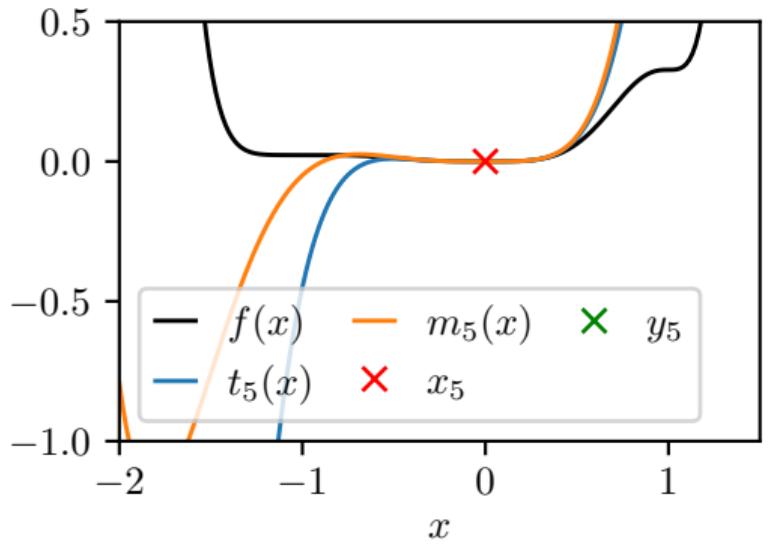
AR5, adaptive σ_k and global model minimizer



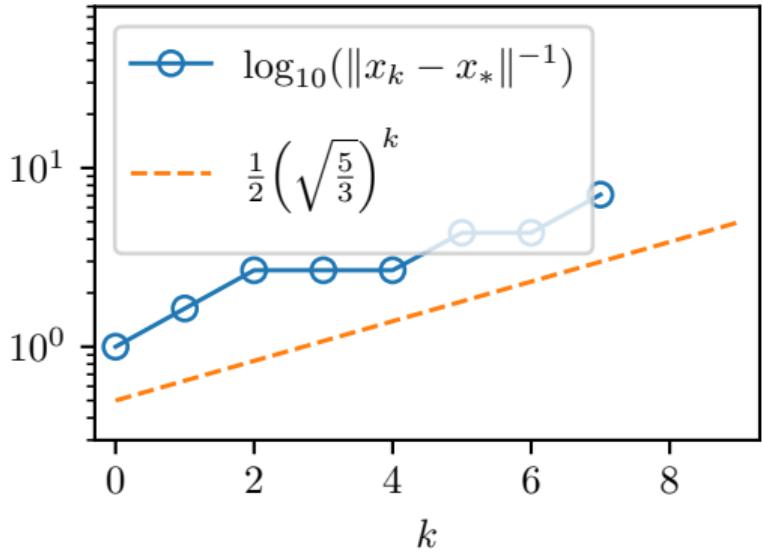
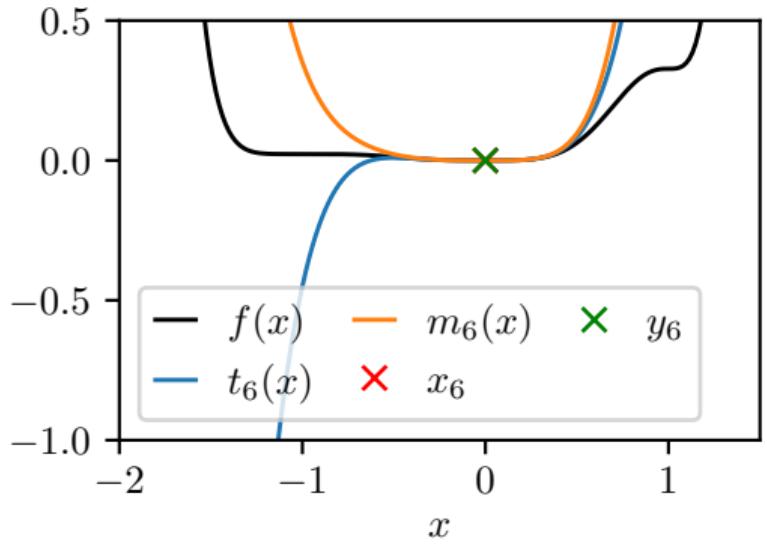
AR5, adaptive σ_k and global model minimizer



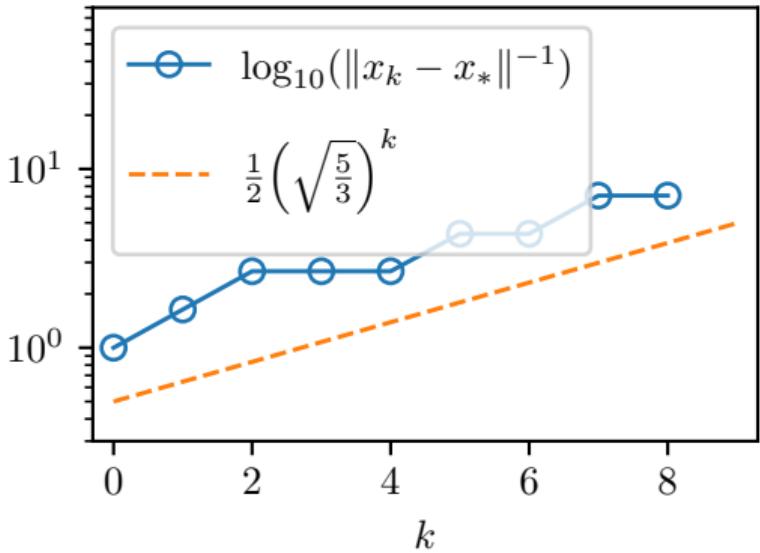
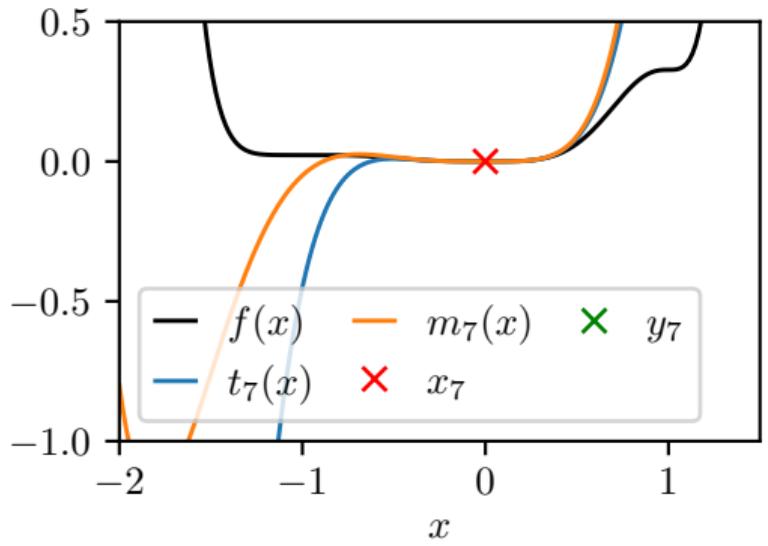
AR5, adaptive σ_k and global model minimizer



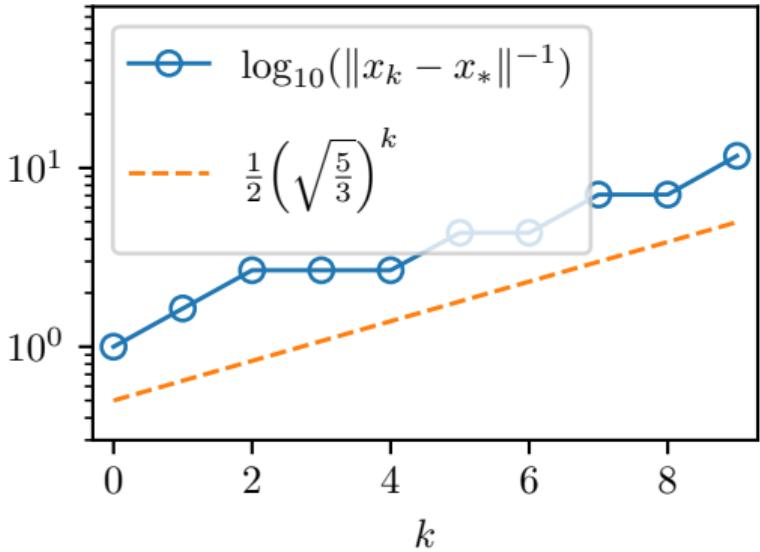
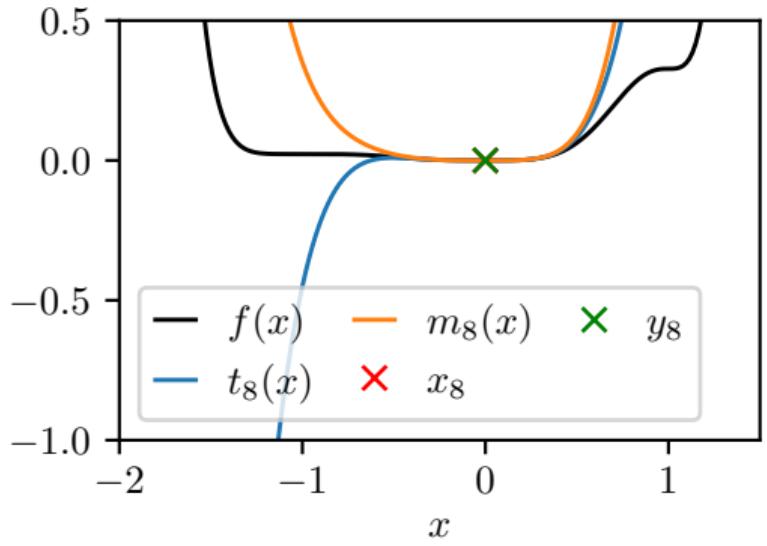
AR5, adaptive σ_k and global model minimizer



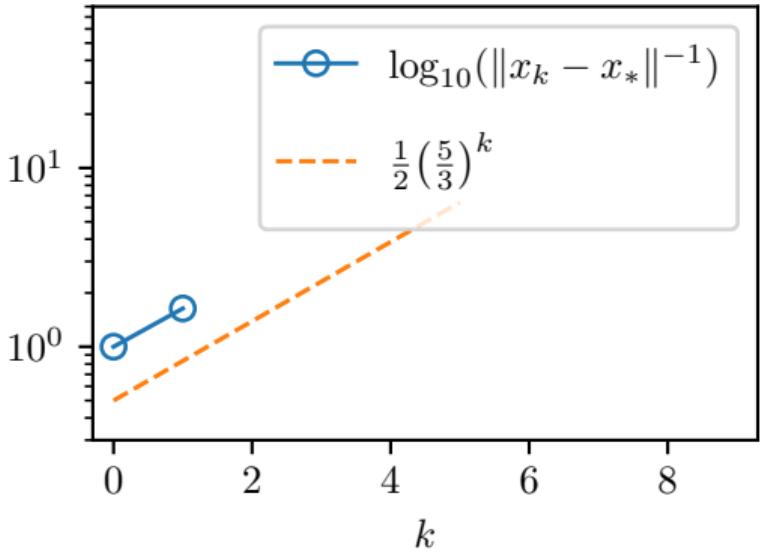
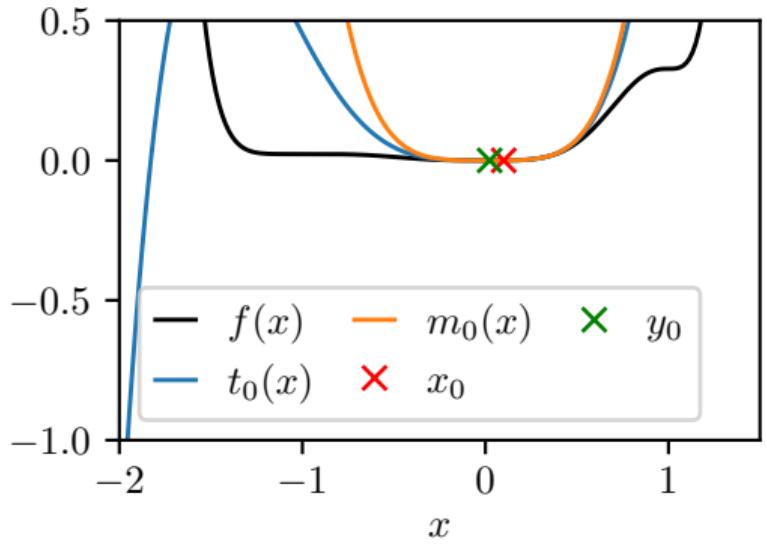
AR5, adaptive σ_k and global model minimizer



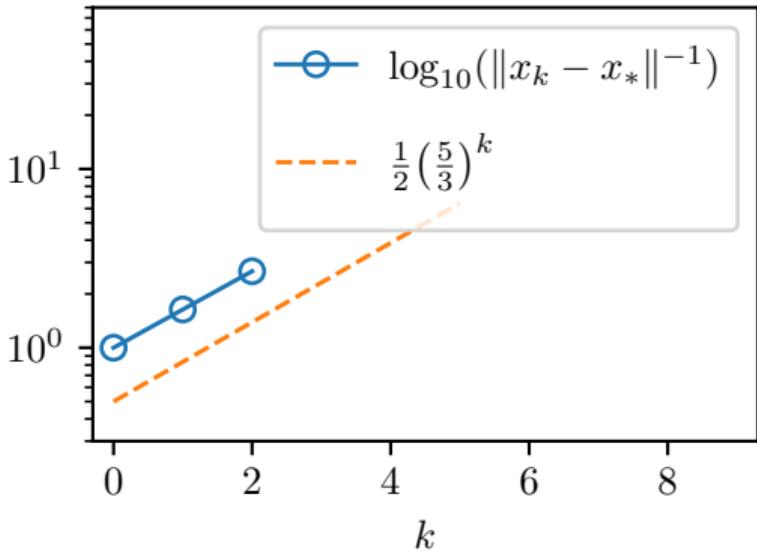
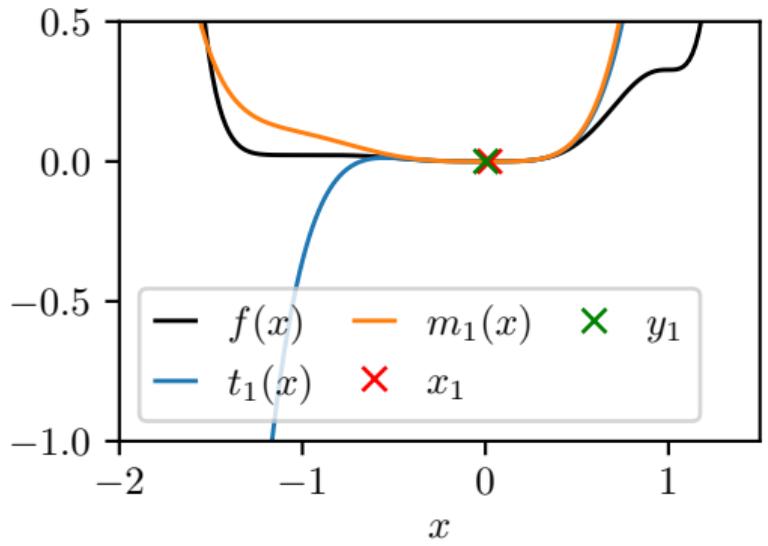
AR5, adaptive σ_k and global model minimizer



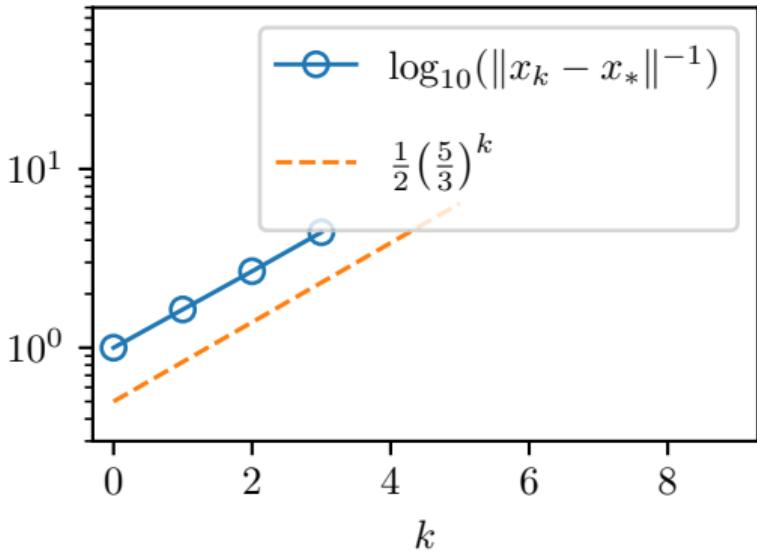
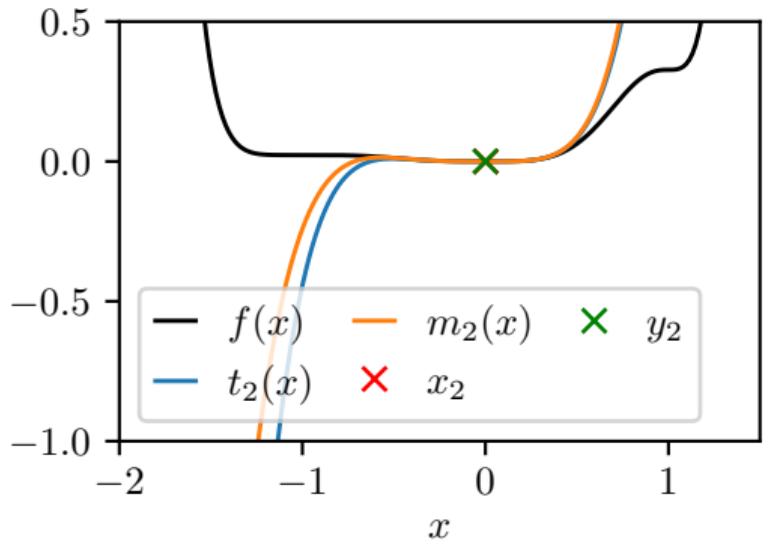
AR5, adaptive σ_k and right model minimizer



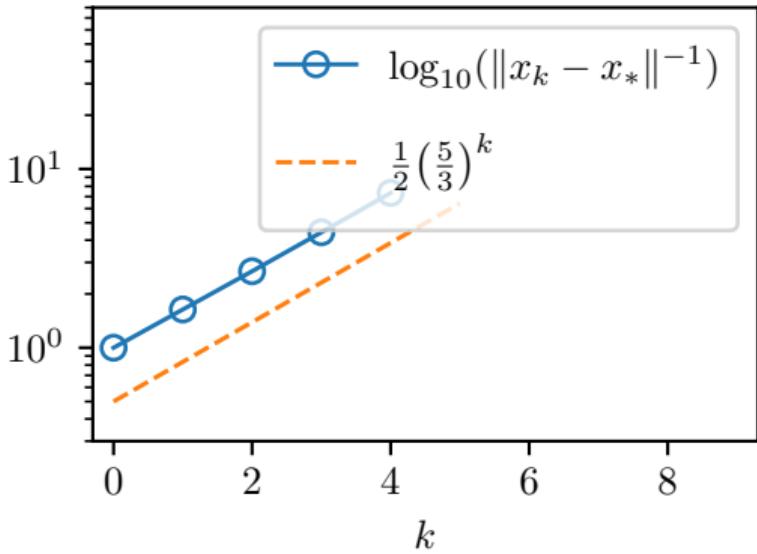
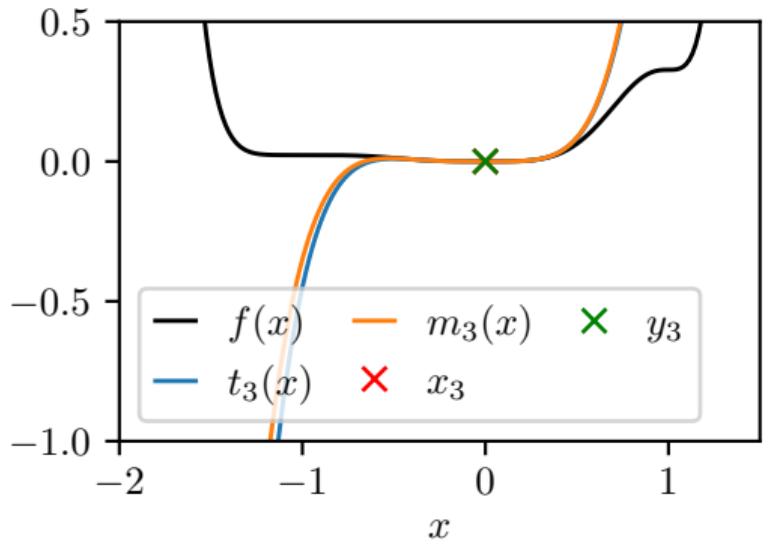
AR5, adaptive σ_k and right model minimizer



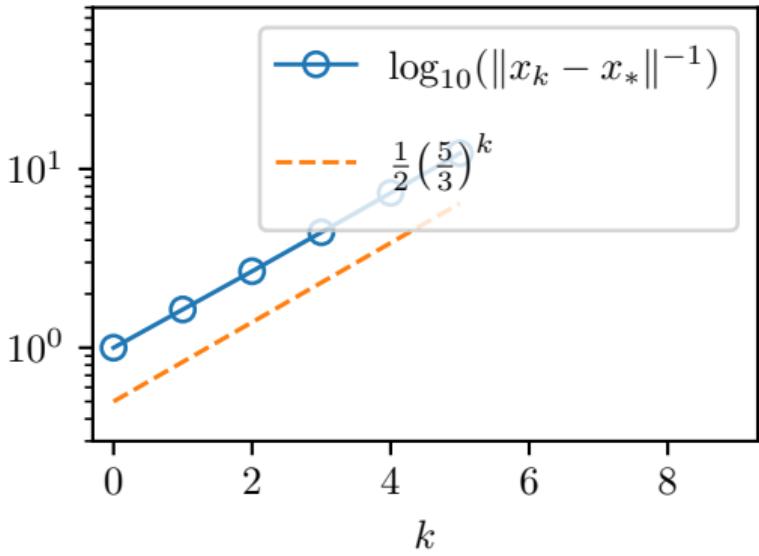
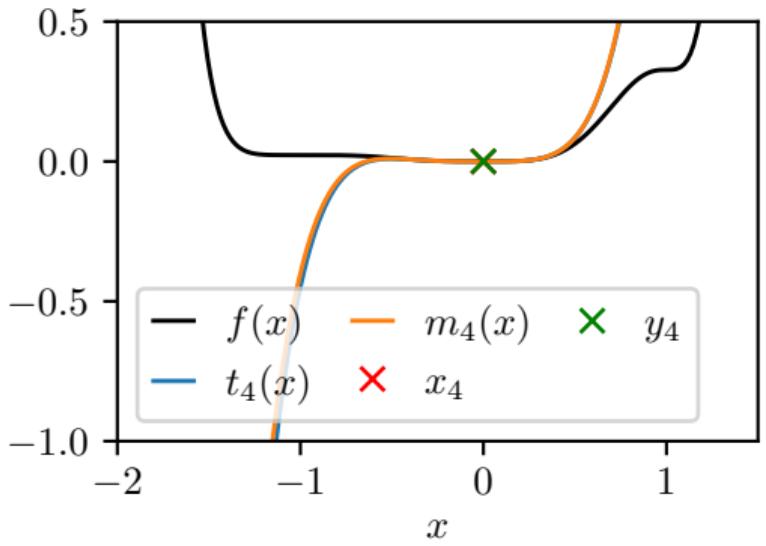
AR5, adaptive σ_k and right model minimizer



AR5, adaptive σ_k and right model minimizer



AR5, adaptive σ_k and right model minimizer



Assumptions

- \mathbf{x}_k and \mathbf{y}_k are generated by ARp (Algorithm 1)
- \mathbf{x}_* is a local minimizer of f
- $\nabla^p f$ is Lipschitz continuous with constant L_p
- f is uniformly convex of order q inside $\mathcal{B}(\mathbf{x}_*, r_\mu)$
- $p > q - 1$

Lemma

If $\mathbf{x}_k, \mathbf{y}_k \in \mathcal{B}(\mathbf{x}_*, r_\mu)$ and iteration k is successful, then the gradients satisfy

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

Theorem

If all iterates stay inside $\mathcal{B}(\mathbf{x}_*, r_\mu)$ and \mathbf{x}_0 is close enough to \mathbf{x}_* , then

$$\|\nabla f(\mathbf{x}_{k_i})\| \rightarrow 0 \quad \text{at a } Q-\frac{p}{q-1}\text{th-order rate}$$

$$f(\mathbf{x}_{k_i}) \rightarrow f(\mathbf{x}_*) \quad \text{at a } Q-\frac{p}{q-1}\text{th-order rate}$$

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

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

Theorem

If all iterates stay inside $\mathcal{B}(\mathbf{x}_*, r_\mu)$ and \mathbf{x}_0 is close enough to \mathbf{x}_* , then

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

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

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

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

Lemma

There exists radii $r_x, r_y > 0$ with $r_x \leq r_y \leq r_\mu$ such that for any $x_k \in \mathcal{B}(x_, r_x)$ there exists a local minimizer of m_k inside $\mathcal{B}(x_*, r_y)$ and any such minimizer will give a successful iteration.*

Theorem

Assume that $\mathbf{y}_k \in \mathcal{B}(\mathbf{x}_*, r_y)$ in every iteration. If \mathbf{x}_0 is close enough to \mathbf{x}_* , then all iterations are successful and

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

Conclusion

- For tensor methods ($p \geq 3$) 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

-  Cartis, Coralia, Nicholas I. M. Gould, and Philippe L. Toint (Jan. 2020). “Sharp Worst-Case Evaluation Complexity Bounds for Arbitrary-Order Nonconvex Optimization with Inexpensive Constraints”. In: *SIAM Journal on Optimization* 30.1, pp. 513–541. ISSN: 1052-6234, 1095-7189. DOI: 10.1137/17M1144854.
-  — (Jan. 2022). *Evaluation Complexity of Algorithms for Nonconvex Optimization: Theory, Computation and Perspectives*. Philadelphia, PA: Society for Industrial and Applied Mathematics. ISBN: 978-1-61197-698-4 978-1-61197-699-1. DOI: 10.1137/1.9781611976991.
-  Doikov, Nikita and Yurii Nesterov (May 2022). “Local Convergence of Tensor Methods”. In: *Mathematical Programming* 193.1, pp. 315–336. ISSN: 1436-4646. DOI: 10.1007/s10107-020-01606-x.