

## EXERCISE 4 - SOLUTION

Date issued: 8th May 2023

Date due: 16th May 2023

### Homework Problem 4.1 (Efficient Step Sizes in Quadratic Steepest Descent)

4 Points

Show that both constant step sizes (as in § 4.3) as well as the Cauchy step size are efficient for the steepest descent method (cf. Algorithm 4.6) for solving quadratic optimization problems of the type

$$\text{minimize } f(x) := \frac{1}{2}x^\top Ax - b^\top x + c \quad \text{where } x \in \mathbb{R}^n$$

with s. p. d.  $A \in \mathbb{R}^{n \times n}$ ,  $b \in \mathbb{R}^n$ ,  $c \in \mathbb{R}$ .

### Solution.

For efficiency of step sizes  $\alpha$ , we need to show that there exists a  $\theta > 0$  such that, as long as  $d^{(k)} \neq 0$ , we have that

$$f(x^{(k)} + \alpha^{(k)} d^{(k)}) \leq f(x^{(k)}) - \theta \left( \frac{f'(x^{(k)})^\top d^{(k)}}{\|d^{(k)}\|_M} \right)^2 \quad (5.11)$$

for all  $k \geq 0$ .

In the steepest descent scheme for quadratic functionals, we use

$$d^{(k)} = -M^{-1}r^{(k)} = -M^{-1}(Ax^{(k)} - b) = -M^{-1}f'(x^{(k)})^\top$$

and therefore obtain

$$\begin{aligned} f(x^{(k+1)}) - f(x^{(k)}) &= f(x^{(k)} + \alpha^{(k)} d^{(k)}) - f(x^{(k)}) \\ &= \frac{1}{2} d^{(k)\top} A d^{(k)} \alpha^{(k)^2} + (Ax^{(k)} - b)^\top d^{(k)} \alpha^{(k)} \\ &= \frac{1}{2} \|d^{(k)}\|_A^2 \alpha^{(k)^2} - \|d^{(k)}\|_M^2 \alpha^{(k)}. \end{aligned}$$

(2 Points)

The Cauchy step size corresponding to  $d^{(k)}$  are

$$\alpha^{(k)} \left( d^{(k)} \right) = - \frac{f'(x^{(k)})d^{(k)}}{d^{(k)\top} A d^{(k)}} = \frac{\|d^{(k)}\|_M^2}{\|d^{(k)}\|_A^2}$$

(**Note:** This shows how the Cauchy step sizes measure the distortion of the  $M$  vs. the  $A$  isolines)  
Accordingly:

$$\begin{aligned} f(x^{(k+1)}) - f(x^{(k)}) &= \frac{1}{2} \|d^{(k)}\|_A^2 \alpha^{(k)2} - \|d^{(k)}\|_M^2 \alpha^{(k)} \\ &= -\frac{1}{2} \frac{\|d^{(k)}\|_M^4}{\|d^{(k)}\|_A^2} \\ &= -\frac{1}{2} \left( \frac{f'(x^{(k)})d^{(k)}}{\|d^{(k)}\|_A} \right)^2 \\ &\leq -\frac{1}{2} \frac{1}{\lambda_{\max}(A; M)} \left( \frac{f'(x^{(k)})d^{(k)}}{\|d^{(k)}\|_M} \right)^2 \end{aligned}$$

shows that the Cauchy step sizes are efficient for the  $M$ -steepest descent directions with  $\theta = \frac{1}{2\lambda_{\max}(A; M)^2}$ .  
(1 Point)

For  $\alpha^{(k)} = \alpha$  constant, we obtain

$$\begin{aligned} f(x^{(k+1)}) - f(x^{(k)}) &= \frac{1}{2} \|d^{(k)}\|_A^2 \alpha^2 - \|d^{(k)}\|_M^2 \alpha \\ &\leq \alpha \left( \frac{\alpha \lambda_{\max}(A; M)}{2} - 1 \right) \|d^{(k)}\|_M^2 \\ &= \alpha \left( \frac{\alpha \lambda_{\max}(A; M)}{2} - 1 \right) \left( \frac{f'(x^{(k)})d^{(k)}}{\|d^{(k)}\|_M} \right)^2 \end{aligned}$$

where

$$\theta := -\alpha^{(k)} \left( \frac{\alpha^{(k)} \lambda_{\max}(A; M)}{2} - 1 \right)$$

is positive whenever  $\alpha \in (0, \frac{2}{\lambda_{\max}(A; M)})$ , which is exactly the condition we derived for the convergence of the steepest descent scheme with constant step sizes in § 4.3. (1 Point)

#### Homework Problem 4.2 (Efficiency of Wolfe-Powell Step Sizes for $C^{1,1}$ Functions) 5 Points

Let  $f \in C^1$  and let  $x^{(0)} \in \mathbb{R}^n$  be an initial iterate of the generic descent scheme (Algorithm 5.2). Further assume that  $f'$  is Lipschitz continuous on the sublevel set  $\mathcal{M}_f(x^{(0)}) := \{x \in \mathbb{R}^n \mid f(x) \leq f(x^{(0)})\}$ .

Show that step sizes  $\alpha^{(k)}$  that satisfy the Wolfe-Powell-conditions at  $x^{(k)}$  for the descent direction  $d^{(k)}$  for all  $k$  are efficient and that there is a  $c > 0$  such that

$$f(x^{(k)} + \alpha^{(k)} d^{(k)}) - f(x^{(k)}) \leq -c \left( \cos \angle(-\nabla_M f(x^{(k)}), d^{(k)}) \|f'(x^{(k)})^\top\|_{M^{-1}} \right)^2$$

for all  $k \geq 0$ .

**Solution.**

For efficiency of step sizes  $\alpha$ , we need to show that there exists a  $\theta > 0$  such that, as long as  $d^{(k)} \neq 0$ , we have that

$$f(x^{(k)} + \alpha^{(k)} d^{(k)}) \leq f(x^{(k)}) - \theta \left( \frac{f'(x^{(k)}) d^{(k)}}{\|d^{(k)}\|_M} \right)^2 \quad (5.11)$$

for all  $k \geq 0$ .

Because the Armijo conditions are satisfied, we actually have a descent scheme.

The curvature condition of the Wolfe-Powell step-lengths states that

$$f'(x^{(k)} + \alpha^{(k)} d^{(k)}) d^{(k)} \geq \tau f'(x^{(k)}) d^{(k)} \quad \text{or} \quad \varphi'(\alpha^{(k)}) \geq \tau \varphi'(0) \quad (5.17)$$

for a  $\tau \in (\sigma, 1)$  and all  $k \geq 0$ . Subtracting  $f'(x^{(k)}) d^{(k)}$  from both sides, applying Cauchy-Schwarz's inequality and using the Lipschitz continuity of  $f'$  (measured in the  $M^{-1}$  and the  $M$  norm, respectively), we obtain that

$$\begin{aligned} (\tau - 1) f'(x^{(k)}) d^{(k)} &\leq \left( f'(x^{(k)} + \alpha^{(k)} d^{(k)}) - f'(x^{(k)}) \right) d^{(k)} \\ &= \left( f'(x^{(k)} + \alpha^{(k)} d^{(k)}) - f'(x^{(k)}) \right) M^{-1} M d^{(k)} \\ &\leq \left\| M^{-1} \left( f'(x^{(k)} + \alpha^{(k)} d^{(k)}) - f'(x^{(k)}) \right)^\top \right\|_M \left\| d^{(k)} \right\|_M \\ &= \left\| \left( f'(x^{(k)} + \alpha^{(k)} d^{(k)}) - f'(x^{(k)}) \right)^\top \right\|_{M^{-1}} \left\| d^{(k)} \right\|_M \\ &\leq L_{M^{-1}, M} \alpha^{(k)} \left\| d^{(k)} \right\|_M^2 \end{aligned}$$

(2 Points)

**Note:** Note that we were able to use Lipschitz continuity because we are working with a descent scheme.

Rearranging the estimate yields the following bound on the step size

$$\alpha^{(k)} \geq \frac{(\tau - 1) f'(x^{(k)}) d^{(k)}}{L_{M^{-1},M} \|d^{(k)}\|_M^2}.$$

Inserting that into the Armijo-condition

$$f(x^{(k)} + \alpha^{(k)} d^{(k)}) \leq f(x^{(k)}) + \sigma \alpha^{(k)} f'(x^{(k)}) d^{(k)} \quad \text{or} \quad \varphi(\alpha^{(k)}) \leq \varphi(0) + \sigma \alpha^{(k)} \varphi'(0), \quad (5.12)$$

we immediately obtain that

$$f(x^{(k)} + \alpha^{(k)} d^{(k)}) \leq f(x^{(k)}) + \underbrace{\sigma \alpha^{(k)} f'(x^{(k)}) d^{(k)}}_{<0} \leq f(x^{(k)}) + \underbrace{\frac{\sigma(\tau - 1)}{L_{M^{-1},M}} \left( \frac{f'(x^{(k)}) d^{(k)}}{\|d^{(k)}\|_M} \right)^2}_{=: -c, c > 0}$$

for all  $k \geq 0$  which shows efficiency. (2 Points)

The additional statement simply follows from the fact that

$$\frac{f'(x^{(k)}) d^{(k)}}{\|d^{(k)}\|_M} = \frac{(\nabla_M f(x^{(k)}), d^{(k)})_M}{\|d^{(k)}\|_M} \frac{\|\nabla_M f(x^{(k)})\|_M}{\|\nabla_M f(x^{(k)})\|_M} = -\cos \angle(-\nabla_M f(x^{(k)}), d^{(k)}) \|f'(x^{(k)})^\top\|_{M^{-1}}$$

and the square. (1 Point)

#### Homework Problem 4.3 (Scaling Invariance of Armijo- and Curvature Conditions) 5 Points

Show the statement of [remark 5.21](#), i. e. that when a step length  $\alpha$  satisfies any of the Armijo- or curvature conditions (5.12), (5.17) and (5.18) for  $g(x) := \gamma f(Ax + b) + \delta$  at  $x \in \mathbb{R}^n$  with search direction  $d \in \mathbb{R}^n$ , where  $A \in \mathbb{R}^{n \times n}$  is non-singular,  $b \in \mathbb{R}^n$ ,  $\gamma > 0$  and  $\delta \in \mathbb{R}$ , then it satisfies the respective conditions for  $f$  at  $Ax + b$  with the search direction  $Ad$ .

#### Solution.

If  $\alpha$  satisfies the [Armijo condition](#) for  $g$  at  $x$  and direction  $d$  with parameter  $\sigma$ , then

$$\begin{aligned} g(x + \alpha d) &\leq g(x) + \sigma g'(x) d \\ \Rightarrow \gamma f(A(x + \alpha d) + b) + \delta &\leq \gamma f(Ax + b) + \delta + \sigma \gamma f'(Ax + b) Ad \end{aligned}$$

and deviding by  $\gamma > 0$  and subtracting  $\delta$  from both sides yields

$$\Rightarrow f(Ax + b + \alpha Ad) \leq f(Ax + b) + \sigma f'(Ax + b) Ad$$

meaning that  $\alpha$  satisfies the Armijo condition for  $f$  at  $Ax + b$  and direction  $Ad$  with parameter  $\sigma$ .  
(2 Points)

If the **curvature condition** is satisfied by  $\alpha$  for  $g$  at  $x$  with direction  $d$  and parameter  $\tau$ , then

$$\begin{aligned} g'(x + \alpha d)d &\geq \tau g'(x)d \\ \Rightarrow \gamma f'(A(x + \alpha d) + b)Ad &\geq \tau \gamma f'(Ax + b)Ad \end{aligned}$$

and deviding by  $\gamma > 0$  yields

$$\Rightarrow f'(Ax + b + \alpha Ad)Ad \geq \tau f'(Ax + b)Ad,$$

which means that the curvature condition is satisfied by  $\alpha$  for  $f$  at  $Ax + b$  with direction  $Ad$  and parameter  $\tau$ .  
(2 Points)

If the **strong curvature condition** is satisfied by  $\alpha$  for  $g$  at  $x$  with direction  $d$  and parameter  $\tau$ , then

$$\begin{aligned} |g'(x + \alpha d)d| &\leq -\tau g'(x)d \\ \Rightarrow |\gamma f'(A(x + \alpha d) + b)Ad| &\leq -\tau \gamma f'(Ax + b)Ad \end{aligned}$$

and deviding by  $\gamma > 0$  yields

$$\Rightarrow |f'(Ax + b + \alpha Ad)Ad| \leq -\tau f'(Ax + b)Ad$$

which means that the strong curvature condition is satisfied by  $\alpha$  for  $f$  at  $Ax + b$  with direction  $Ad$  and parameter  $\tau$ .  
(1 Point)

**Note:** We did not require  $A$  to be nonsingular anywhere in the proof. This requirement is merely needed to show the inverse by applying the result we just proved to the inverted form that generates  $f$  from  $g$ .

**Homework Problem 4.4** (Implementation of Nonlinear Steepest Descent and Armijo Backtracking)  
8 Points

Implement the  $M$ -steepest descent method as outlined in [Algorithm 5.22](#) with the original and the modified (interpolating) Armijo backtracking as outlined in [Algorithms 5.11](#) and [5.15](#).

Visualize and examine the effect of the parameters of the step size strategy on the behavior of the algorithm when applied to quadratic, strongly convex functions, **Rosenbrock's** and/or **Himmelblau's** functions.

### Solution.

Figures 0.1 and 0.3 show the behavior of the iterates and Figures 0.2 and 0.4 show the (approximate)  $f''(x^*)$ -error in a semilogarithmic plot for the steepest descent method with armijo backtracking applied to quadratic optimization problem and the minimization of the Rosenbrock function (with minimizer  $x^* = (1, 1)$ ), respectively, for backtracking parameters  $\beta \in \{0.01, 0.5, 0.99\}$  and  $\sigma \in \{10^{-2}, 0.3, 0.7\}$ .

We can observe that when  $\beta$  is chosen very small (top rows), the choice of  $\sigma$  does not influence the behavior greatly. This is due to one backtracking step reducing the trial step size so much that we essentially end up with gradient flow like behavior. Only when  $\sigma$  is really large (very relaxed acceptance of trial steps) we can get lucky with the first trials in the rosenbrock case.

For rather strict acceptance of trial step sizes (large  $\sigma$ , right columns) we also observe gradient flow type behavior (as this forces the step sizes to be small).

When  $\beta$  is chosen large (bottom rows) we can expect fine adjustments of the step sizes. When  $\sigma$  is small an we are lenient with accepting step sizes, we end up with pretty extreme zig-zagging. We observe the best behavior for moderate parameter choices. Small  $\beta$  and large  $\sigma$  dominate the behavior in any case.

Near minimizers with s.p.d.  $f''(x^*)$  (this is the sufficient second order condition) we can expect nonlinear problems to behave almost quadratically, so the difference in function values  $f(x^{(k)}) - f(x^*) \approx \|x^{(k)} - x^*\|_{f''(x^*)}^2$ , which coincides with the A-error in the quadratic case. Looking at those plots over the iterations in semilogarithmic axes shows the expected behavior once close to the minimizer. Zig-zagging and sudden drops are experienced well outside close neighbourhoods of the minimizer.

Note that the steepest descent method is struggling strongly with the rosenbrock function due to zig zagging along the barely-sloping banana-shaped valley leading to the minimizer, as no local curvature information is used in the model hessians (identity preconditioner for all plots). Depending on the safe-guarding conditions on might actually end up with the step size computation failing because step sizes become to small and the function almost appears locally constant.

Additionally, using interpolation over simple backtracking can improve the number of iterations needed until the termination criteria is used, but it actually may increase them as well.

(8 Points)

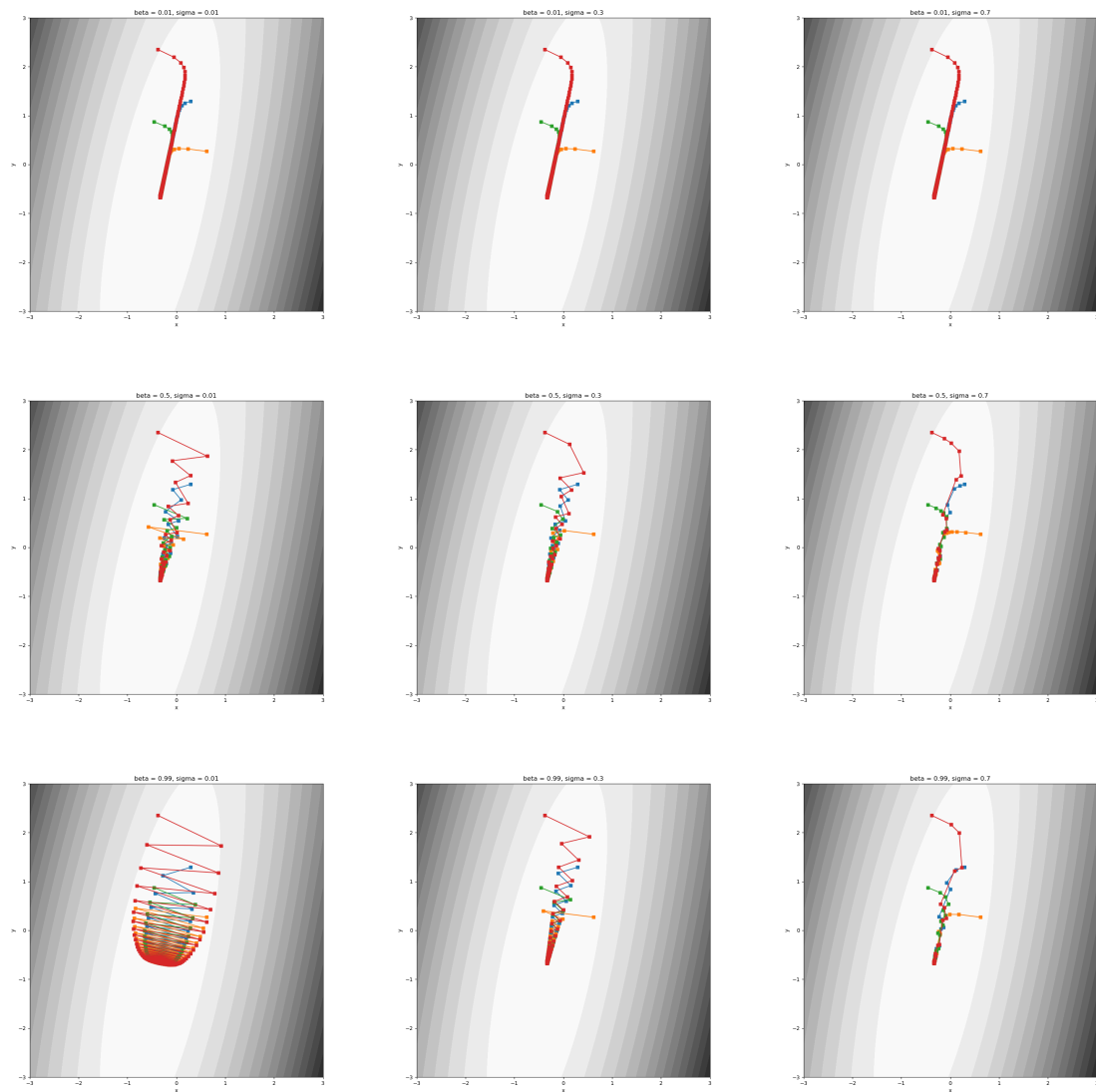


Figure 0.1: Iterates for steepest descent with armijo step length rule applied to quadratic optimization.

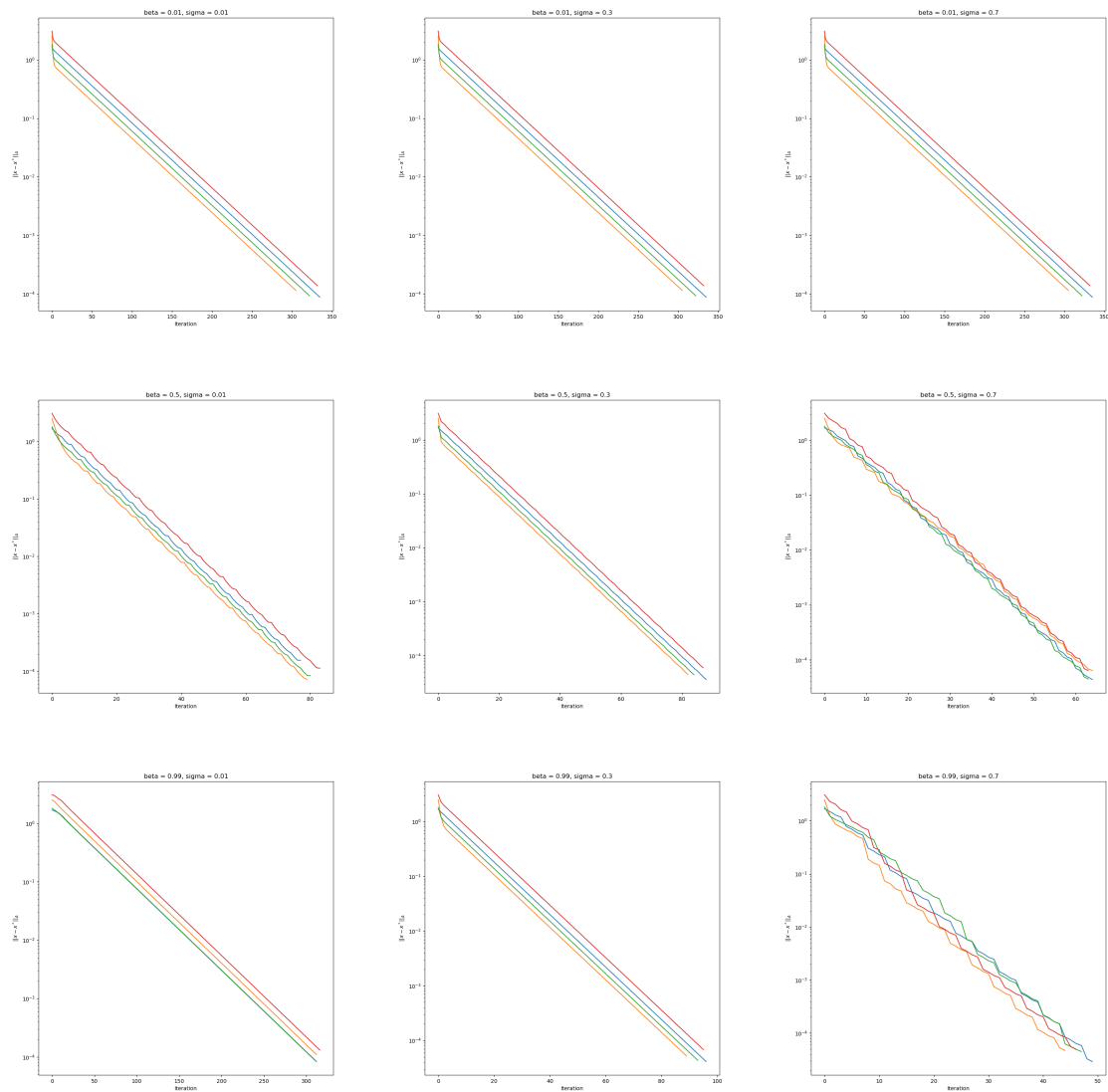


Figure 0.2: A- norm of errors for steepest descent with armijo step length rule applied to quadratic optimization.



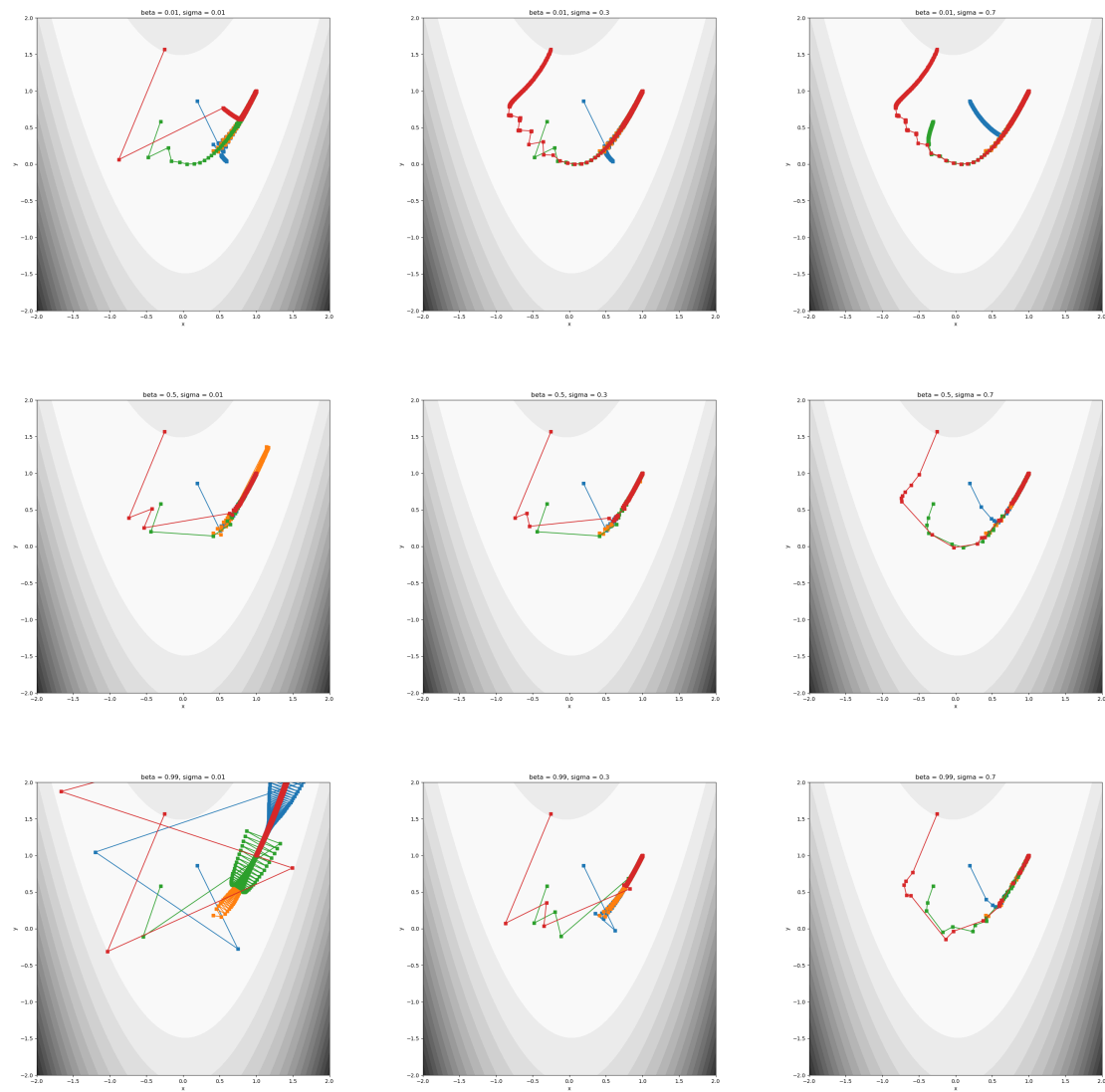


Figure 0.3: Iterates for steepest descent with armijo step length rule applied to rosenbrock optimization.

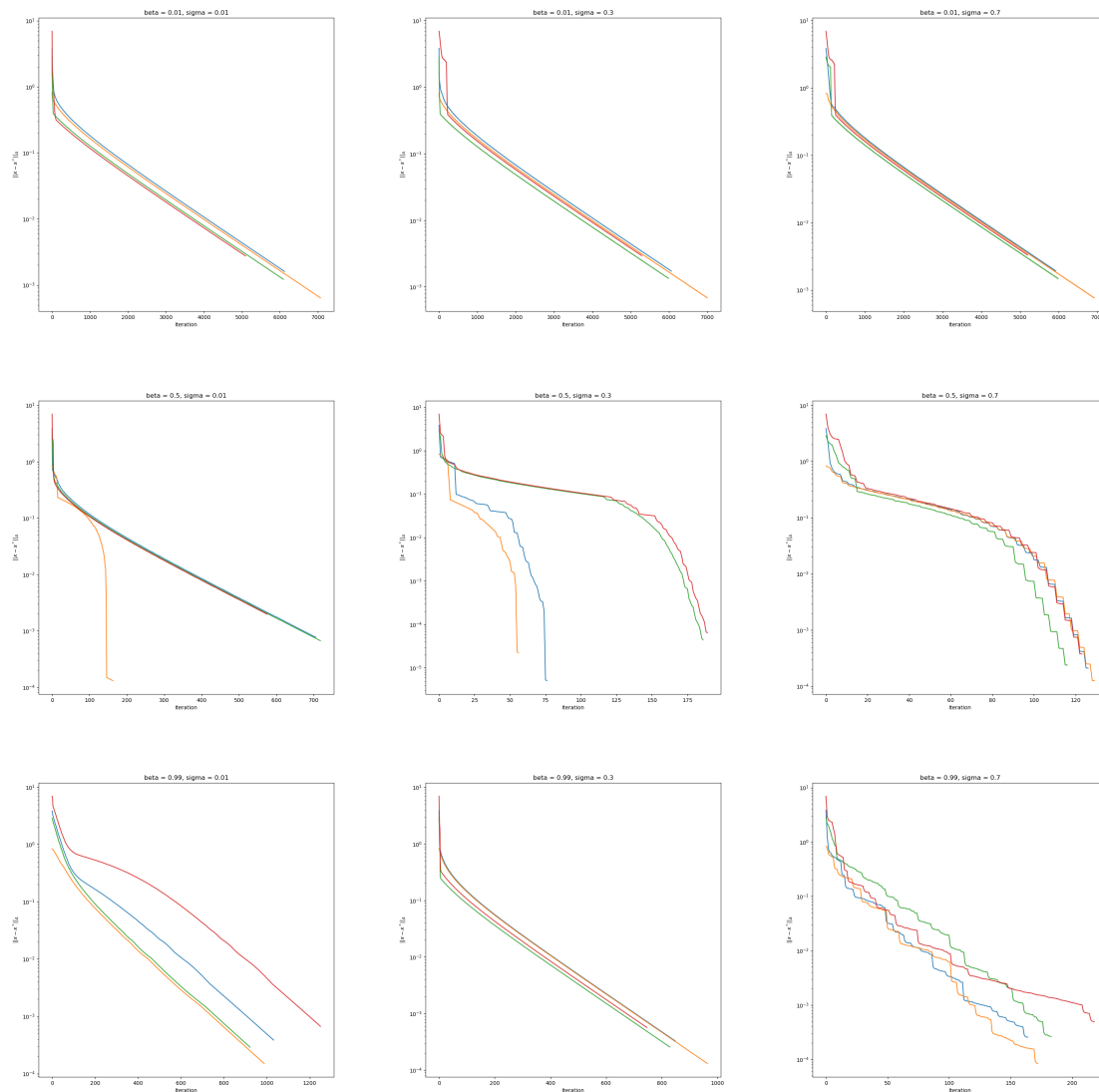


Figure 0.4: Approximate  $f''(x^*)$ -error for steepest descent with armijo step length rule applied to rosenbrock optimization.

Please submit your solutions as a single pdf and an archive of programs via [moodle](#).