

Fall 2017: Numerical Methods I Assignment 4 (due Nov 2, 2017)

1. **[Convexity, 1+1+1pt]** A set \mathcal{S} is convex, if for all $\mathbf{x}, \mathbf{y} \in \mathcal{S}$, also the convex combinations $\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}$ are in \mathcal{S} , for all $\lambda \in (0, 1)$. Convex functions are usually defined over convex sets (or over all of \mathbb{R}^n , which is, of course, a convex set).
 - (a) Show that for a convex function f , the set of global minimizers of f is a convex set.
 - (b) Let $\|\cdot\|$ be any norm on \mathbb{R}^n . Show that $\mathbf{x} \mapsto \|\mathbf{x}\|$ is a convex function.
 - (c) Let $Q \in \mathbb{R}^{n \times n}$ be symmetric and positive semidefinite, and $\mathbf{b} \in \mathbb{R}^n$. Prove that the function $\mathbf{x} \mapsto \frac{1}{2}\mathbf{x}^T Q \mathbf{x} - \mathbf{b}^T \mathbf{x}$ is convex.

2. **[Necessary and sufficient optimality conditions, 1+1+1pt]** Consider the function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ given by $f(x, y) = (x - y^2)(3x - y^2)$.
 - (a) Show that f satisfies the first and second-order necessary optimality conditions at $\mathbf{x}^* = (0, 0)$, but every neighborhood of \mathbf{x}^* contains a point \mathbf{z} with $f(\mathbf{z}) < f(\mathbf{x}^*)$.
 - (b) Nevertheless, show that for every direction $\mathbf{d} \in \mathbb{R}^2$, $\alpha^* = 0$ is a minimizer for the restriction to the one-dimensional function $\alpha \mapsto g(\alpha) := f(\mathbf{x}^* + \alpha \mathbf{d})$.
 - (c) Give an example of a twice continuously differentiable function that has a minimum at a point \mathbf{x}^* that satisfies first and second-order necessary conditions, but does not satisfy the second-order sufficient conditions.

3. **[Modified metric in steepest descent, 2pt]** Consider $f : \mathbb{R}^n \rightarrow \mathbb{R}$ continuously differentiable, and $\mathbf{x} \in \mathbb{R}^n$ with $\nabla f(\mathbf{x}) \neq 0$. For a symmetric positive definite matrix $A \in \mathbb{R}^{n \times n}$, we define the A -weighted norm $\|\mathbf{y}\|_A = \sqrt{\mathbf{y}^T A \mathbf{y}}$. Derive the unit norm steepest descent direction of f in \mathbf{x} with respect to the $\|\cdot\|_A$ -norm, i.e., find the solution to the problem¹

$$\min_{\|\mathbf{d}\|_A=1} \nabla f(\mathbf{x})^T \mathbf{d}.$$

Hint: Use the factorization $A = B^T B$ and the Cauchy-Schwarz inequality.

4. **[Gaussian mixtures interpolation/data fitting (continuation of problem 5 from the previous assignment), 1+1+2+1pt+2pt extra credit]** Let us try an approach for fitting data that avoids having to solve a *nonlinear* least squares problem. We use the same measurements as for the previous assignment, but choose a larger N and make several fixed choices for τ_j, σ_j , and only use a numerical method to find appropriate weights ω_j . We use two different values for σ_j , namely 1 and $1/2$, and 11 values for τ_j , namely $\tau_j \in \{-5, -4, \dots, 4, 5\}$. This gives the model function

$$\varphi_2(t; \mathbf{x}) = \sum_{j=1}^{11} \omega_j \exp\left(-\frac{1}{2}(t - (j - 6))^2\right) + \sum_{j=1}^{11} \omega_{j+11} \exp(-2(t - (j - 6))^2). \quad (1)$$

¹If f is twice differentiable with positive definite Hessian matrix, one can choose $A = \nabla^2 f(\mathbf{x})$. This shows that the Newton descent direction is the steepest descent direction in the metric where norms are weighted by the Hessian matrix.

Now, the only model parameters are $\mathbf{x} = (\omega_1 \dots, \omega_{22})^T$, and since these enter linearly in φ_2 , this is now a *linear* least squares problem, namely

$$\min_{\mathbf{x} \in \mathbb{R}^{22}} \frac{1}{2} \sum_{i=1}^{10} (\varphi_2(t_i; \mathbf{x}) - f_i)^2. \quad (2)$$

- (a) Specify the entries of the matrix $A \in \mathbb{R}^{10 \times 22}$ and the vector $\mathbf{f} \in \mathbb{R}^{10}$ such that (2) can be written as

$$\min_{\mathbf{x} \in \mathbb{R}^{22}} \frac{1}{2} \|A\mathbf{x} - \mathbf{f}\|^2.$$

- (b) Comparing the number of unknowns and linear equations, we see that this is an *underdetermined* system, which does not have a unique solution \mathbf{x} . Find a solution to this underdetermined system using the built-in function (backslash in MATLAB). Amongst all of the possible solutions to this system, the built-in function should return a vector \mathbf{x} that has some sparsity, i.e., a vector where several entries are zero, which means that many of the coefficients ω_j in (1) are zero. Because the terms with zero weight can, of course, be neglected in (1), we have found a model with small(ish) N . Let's aim at improving the result by looking at alternative approaches.

- (c) To make the solution \mathbf{x} unique, we can also consider the following modified least squares problem

$$\min_{\mathbf{x} \in \mathbb{R}^{22}} \frac{1}{2} \|A\mathbf{x} - \mathbf{f}\|^2 + \frac{\gamma}{2} \|\mathbf{x}\|^2, \quad (3)$$

where γ is a (small) positive number. Show that the minimizer of (3) satisfies the modified normal equation

$$(A^T A + \gamma I)\mathbf{x} = A^T \mathbf{f}, \quad (4)$$

where $I \in \mathbb{R}^{22 \times 22}$ denotes the identity.² Prove that solutions of (4) are unique for $\gamma > 0$.

- (d) Solve (4) for different $\gamma > 0$ and plot the resulting model curves on top of the data points (t_i, f_i) . We are told that N should be small in the model function, which would mean that many of the components in \mathbf{x} are zero or close to zero. Is this the case for the solutions of (4) that you found?³
- (e) Let us try to obtain a vector \mathbf{x} that is sparser, i.e., has more zeros (and thus the model $\varphi(t)$ has a small N as desired). We add a different term to the minimization problem compared to (3), namely:

$$\min_{\mathbf{x} \in \mathbb{R}^{22}} \frac{1}{2} \|A\mathbf{x} - \mathbf{f}\|^2 + \gamma \|\mathbf{x}\|_1, \quad (5)$$

where $\|\mathbf{x}\|_1 = \sum_{j=1}^{22} |x_j|$ is the 1-norm. This choice tends to "favor" sparse vectors, i.e., by enforcing that $\|\mathbf{x}\|_1$ be small, many entries of the computed \mathbf{x} are zero. We

²Hint: At the minimizer of (3), the gradient with respect to \mathbf{x} of $1/2\|A\mathbf{x} - \mathbf{f}\|^2 + \gamma/2\|\mathbf{x}\|^2$ must be zero. Alternatively, you can consider the an extended least squares problem with the augmented matrix $\tilde{A} = [A, \sqrt{\gamma}I]^T \in \mathbb{R}^{32 \times 22}$ and use the results from class for overdetermined systems.

³It shouldn't be. By adding the additional term in (3), we say that we prefer vectors \mathbf{x} with small (squared) 2-norm, but not that we are looking for solutions \mathbf{x} that contain a large number of zeros.

would like to take derivatives of the function in (5) but, as we have seen before, this is problematic since the 1-norm is not differentiable. We thus replace (5) by

$$\min_{\mathbf{x} \in \mathbb{R}^{22}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{f}\|^2 + \gamma \sum_{j=1}^{22} \sqrt{x_j^2 + \varepsilon}, \quad (6)$$

where ε is a (small) positive number. Show⁴ that the solution of (6) satisfies the nonlinear equation

$$\mathbf{A}^T \mathbf{A}\mathbf{x} + \gamma D(\mathbf{x})\mathbf{x} = \mathbf{A}^T \mathbf{f}, \quad (7)$$

where $D(\mathbf{x}) = \text{diag}(1/\sqrt{x_1^2 + \varepsilon}, \dots, 1/\sqrt{x_{22}^2 + \varepsilon}) \in \mathbb{R}^{22 \times 22}$. To solve this nonlinear equation, choose an initial guess \mathbf{x}^0 and use the fixed point-like iteration

$$(\mathbf{A}^T \mathbf{A} + \gamma D(\mathbf{x}^{k-1}))\mathbf{x}^k = \mathbf{A}^T \mathbf{f} \quad \text{for } k = 1, 2, \dots$$

Try different values for γ and ε ⁵, and plot the model functions which result from the computed parameters. Are the vectors \mathbf{x} obtained with this method sparser? Since we had to introduce the parameter ε to make the problem differentiable, most entries in \mathbf{x} will not exactly be zero, but they should be rather small. Based on your results, what do you think are the parameters of the Gaussian mixture used to generate the original data?

5. **[Behavior of descent methods, 6pt]** Consider the unconstrained optimization problem

$$\min f(x, y) \equiv -\cos x \cos(y/5).$$

- Find and classify all stationary points in the region $-\pi/2 \leq x \leq \pi/2, -5\pi/2 \leq y \leq 5\pi/2$.
- There is a portion of the problem region within which the Hessian matrix of $f(x, y)$ is positive definite. Give expressions for this portion. You should be able to do this analytically.
- Derive expressions for the search directions associated with the steepest descent and Newton methods.
- Write a program that performs both iterations, both without a line search and with an exact line search. Note that you will not be able to find the value of the optimal step length analytically; instead, determine it numerically.⁶
- Run your program for various initial guesses within the region. Verify the following:⁷

⁴Again, you can show that by computing the gradient of the function in (6) and using that it has to be zero at the minimizer.

⁵To get good results, these should be small, e.g., $\gamma \sim 10^{-3}$ and $\varepsilon \sim 10^{-9}$

⁶You can use a built-in one-dimensional minimization function (`fzero` in MATLAB). While we use exact line search here, this is usually too costly as it requires a large number of function evaluations. As we've discussed, one thus uses an inexact step size that satisfies, for instance, the Wolfe conditions to guarantee convergence to a stationary point.

⁷One way to illustrate this (non)convergence is to randomly choose initializations and draw them as dots with different colors depending on whether the method, started from that initialization, converged.

- Steepest descent converges to the minimum x^* for any starting point within the region.
- Newton's method with line search converges to the minimum only for initial points for which the Hessian matrix is positive definite.
- Newton's method without line search has an even more restricted radius of convergence.

(f) What do you observe about the convergence rate in these cases?

6. **[Globalization of Newton descent, 3+1pt]** As we have seen, the Newton direction for solving a minimization problem is only a descent direction if the Hessian matrix is positive definite. This is not always the case, in particular far from the minimizer. To guarantee a descent direction in Newton's method, a simple idea is as follows (where we choose $0 < \alpha_1 < 1$ and $\alpha_2 > 0$):

- Compute a direction \mathbf{d}^k by solving the Newton equation

$$\nabla^2 f(\mathbf{x}^k) \mathbf{d}^k = -\nabla f(\mathbf{x}^k).$$

If that is possible and \mathbf{d}^k satisfies⁸

$$\frac{-\nabla f(\mathbf{x}^k)^T \mathbf{d}^k}{\|\nabla f(\mathbf{x}^k)\| \|\mathbf{d}^k\|} \geq \min(\alpha_1, \alpha_2 \|\nabla f(\mathbf{x}^k)\|), \quad (8)$$

then use \mathbf{d}^k as descent direction.

- Otherwise, use the steepest descent direction $\mathbf{d}^k = -\nabla f(\mathbf{x}^k)$.

To illustrate this globalization, let $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ be defined by

$$f(\mathbf{x}) = \frac{1}{2}(x_1^2 + x_2^2) \exp(x_1^2 - x_2^2).$$

- Using the initial iterate $\mathbf{x}^0 = (1, 1)^T$, find a local minimum of f using the modified Newton method described above, combined with Armijo line search with backtracking.⁹ [Hand in a listing of your implementation.](#)¹⁰
- Carry out the computation also with the modified Newton matrix $\nabla^2 f(\mathbf{x}) + 3\mathbf{I}$, which has been made to be positive definite by adding a multiple of the identity. Discuss your findings.

⁸This is a condition on the angle between the negative gradient and the Newton directions, which must less than 90° . However, using this condition, the angle may approach 90° at the same speed as $\|\nabla f(\mathbf{x}^k)\|$ approaches zero. Recall from class that what is required to guarantee convergence for a descent method with Wolfe line search is that an infinite sum that involves the square of the right hand side in (8) and the norm of the gradient is finite.

⁹For the Armijo linesearch, one usually uses c to be rather small, e.g., $c = 10^{-4}$. Use $\alpha = 0.1$ and $\alpha_2 = 0.5$ for the constants in (8).

¹⁰It is sufficient to hand in a listing of the important parts of your implementation, i.e., Armijo line search and computation of the descent direction.