

Week 2 Tutorial 1

(1.1) Determine the critical points of the function

$$f(x, y) = x^3 + y^3 - 3\alpha xy$$

with respect to $\alpha \in \mathbb{R} \setminus \{0\}$ and decide whether it is a minimum, maximum or saddle point.

The gradient of the function is given by

$$\nabla f(x, y) = \begin{pmatrix} 3x^2 + 3\alpha y \\ -3y^2 + 3\alpha x \end{pmatrix}$$

Setting it equal to zero, gives the following system

$$\begin{aligned} 3x^2 + 3\alpha y &= 0 \\ -3y^2 + 3\alpha x &= 0 \end{aligned}$$

which implies that $\alpha y = -x^2$ and $\alpha x = y^2$. Therefore $x = \frac{y^2}{\alpha}$, and $\alpha y + \frac{y^4}{\alpha} = y(y^3 + \alpha^3) = 0$. The solutions are therefore

$$y = 0 \text{ and } x = 0 \quad \text{or} \quad y = -\alpha \text{ and } x = \alpha.$$

The Hessian is

$$\nabla^2 f(x, y) = \begin{pmatrix} 6x & 3\alpha \\ 3\alpha & -6y \end{pmatrix}$$

Evaluating the Hessian at the two critical points gives

$$\begin{aligned} (0, 0) : & \begin{pmatrix} 0 & 3\alpha \\ 3\alpha & 0 \end{pmatrix} \\ (\alpha, -\alpha) : & \begin{pmatrix} 6\alpha & 3\alpha \\ 3\alpha & 6\alpha \end{pmatrix} \end{aligned}$$

In the first case the determinant is negative for all values of α considered, and therefore we have a saddle point.

In the second case the determinant is given by $25\alpha^2$. For $\alpha > 0$ we have that $f_{xx} > 0$ and therefore it's a local minimum, for $\alpha < 0$ we have that $f_{xx} < 0$ and therefore it's a local maximum.

For TAs Some stuff you can discuss before or while going through the example:

- Partial derivatives, gradients and Hessian
- 2nd partial derivative test
https://en.wikipedia.org/wiki/Second_partial_derivative_test

(1.2) Consider the **Rosenbrock function** in \mathbb{R}^2 ,

$$f(\mathbf{x}) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2.$$

Compute the gradient ∇f and the Hessian $\nabla^2 f$. Show that $\mathbf{x}^* = (1, 1)^\top$ is the only local minimizer of this function, and that the Hessian at this point is positive definite.

Using Python or another computing system, draw a contour plot of the Rosenbrock function.

The gradient of the Rosenbrock function is given by

$$\nabla f = \begin{pmatrix} -400x_1(x_2 - x_1^2) + 2(1 - x_1) \\ 200(x_2 - x_1^2) \end{pmatrix}$$

The Hessian by

$$\nabla^2 f = \begin{pmatrix} 1200x_1^2 - 400x_2 + 2 & -400x_1 \\ -400x_1 & 200 \end{pmatrix}$$

We see that for $(1, 1)$ the gradient vanishes, the eigenvalues of the Hessian are $\lambda_1 \approx 1001.6$ and $\lambda_2 = 0.399361$ hence it's positive definite.

For TAs: The Rosenbrock function is one of the benchmark problems in global optimisation. You can show them how it looks like (contour lines,) and maybe have a look at other benchmark problems

https://en.wikipedia.org/wiki/Rosenbrock_function

<https://www.sfu.ca/~ssurjano/optimization.html>

(1.3) Problems in optimization are often solved by iteration: if we want to minimise a function $f(\mathbf{x})$, we start with a point \mathbf{x}_0 and generate a sequence of points $\mathbf{x}_1, \mathbf{x}_2, \dots$ such that the \mathbf{x}_i approach a minimizer of f . One method of generating such a sequence for a differentiable function is by **gradient descent**:

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \eta \nabla f(\mathbf{x}_i), \quad i = 1, 2, \dots$$

The parameter ∇_i can be constant or change in time, and is called the **step length** or **learning rate** in machine learning.

Using Python or another computing system, compute and plot the sequence of points \mathbf{x}_k , starting with $\mathbf{x}_0 = (0, 0)^\top$, for the gradient descent algorithm for the problem

$$\text{minimize} \quad \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$$

with data

$$\mathbf{A} = \begin{pmatrix} 1 & 2 \\ 2 & 1 \\ -1 & 0 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 10 \\ -1 \\ 0 \end{pmatrix}.$$

Experiment with different step lengths and try to find an optimal one.

For TAs: This is a bit of an outlook (we will cover iterative algorithms in week 2). It could also be useful to write out the optimality conditions (time permitting) and calculate the minimiser explicitly. For information, here is a snippet from the lecture notes on iterative algorithms describing the derivation of the optimal step length and the resulting algorithm:

We add the factor $1/2$ because it makes the expressions a bit simpler down the line:

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2.$$

The Hessian is symmetric and positive semidefinite, with the gradient given by

$$\nabla f(\mathbf{x}) = \mathbf{A}^\top (\mathbf{Ax} - \mathbf{b}).$$

The method of gradient descent proceeds as

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{A}^\top (\mathbf{Ax}_k - \mathbf{b}).$$

To find the best α_k , we compute the minimum of the function

$$\alpha \mapsto \varphi(\alpha) = f(\mathbf{x}_k - \alpha \mathbf{A}^\top (\mathbf{Ax}_k - \mathbf{b})). \quad (0.1)$$

If we set $\mathbf{r}_k := \mathbf{A}^\top (\mathbf{b} - \mathbf{Ax}_k) = -\nabla f(\mathbf{x}_k)$ and compute the minimum of (0.1) by setting the derivative to zero,

$$\begin{aligned} \varphi'(\alpha) &= \frac{d}{d\alpha} f(\mathbf{x}_k + \alpha \mathbf{r}_k) = \langle \nabla f(\mathbf{x}_k + \alpha \mathbf{r}_k), \mathbf{r}_k \rangle \\ &= \langle \mathbf{A}^\top (\mathbf{A}(\mathbf{x}_k + \alpha \mathbf{r}_k) - \mathbf{b}), \mathbf{r}_k \rangle \\ &= \langle \mathbf{A}^\top (\mathbf{Ax}_k - \mathbf{b}), \mathbf{r}_k \rangle + \alpha^2 \langle \mathbf{A}^\top \mathbf{Ar}_k, \mathbf{r}_k \rangle \\ &= -\mathbf{r}_k^\top \mathbf{r}_k + \alpha \mathbf{r}_k^\top \mathbf{A}^\top \mathbf{Ar}_k = 0, \end{aligned}$$

we get the step length

$$\alpha_k = \frac{\mathbf{r}_k^\top \mathbf{r}_k}{\mathbf{r}_k^\top \mathbf{A}^\top \mathbf{Ar}_k} = \frac{\|\mathbf{r}_k\|_2^2}{\|\mathbf{Ar}_k\|_2^2}.$$

Note also that when we have \mathbf{r}_k and α_k , we can compute the next \mathbf{r}_k as

$$\begin{aligned} \mathbf{r}_{k+1} &= \mathbf{A}^\top (\mathbf{b} - \mathbf{Ax}_{k+1}) \\ &= \mathbf{A}^\top (\mathbf{b} - \mathbf{A}(\mathbf{x}_k + \alpha_k \mathbf{r}_k)) \\ &= \mathbf{A}^\top (\mathbf{b} - \mathbf{Ax}_k - \alpha_k \mathbf{Ar}_k) = \mathbf{r}_k - \alpha_k \mathbf{A}^\top \mathbf{Ar}_k. \end{aligned}$$

The gradient descent algorithm for the linear least squares problem proceeds by first computing $\mathbf{r}_0 = \mathbf{A}^\top (\mathbf{b} - \mathbf{Ax}_0)$, and then at each step

$$\begin{aligned} \alpha_k &= \frac{\mathbf{r}_k^\top \mathbf{r}_k}{\mathbf{r}_k^\top \mathbf{A}^\top \mathbf{Ar}_k} \\ \mathbf{x}_{k+1} &= \mathbf{x}_k + \alpha_k \mathbf{r}_k \\ \mathbf{r}_{k+1} &= \mathbf{r}_k - \alpha_k \mathbf{A}^\top \mathbf{Ar}_k. \end{aligned}$$