

## Machine Learning Exercise Sheet 06

### Optimization

## Homework

### 1 Convexity of functions

**Problem 1:** Given  $n$  convex functions  $g_i : \mathbb{R}^{d_i} \rightarrow \mathbb{R}$  for  $i \in \{1, \dots, n\}$ , prove or disprove that the function

- a)  $h(\mathbf{x}) = g_2(g_1(\mathbf{x}))$  is convex (here  $d_1 \in \mathbb{N}$ ,  $d_2 = 1$ ),
- b)  $h(\mathbf{x}) = g_2(g_1(\mathbf{x}))$  is convex if  $g_2$  is non-decreasing (here  $d_1 \in \mathbb{N}$ ,  $d_2 = 1$ ),
- c)  $h(\mathbf{x}) = \max(g_1(\mathbf{x}), \dots, g_n(\mathbf{x}))$  is convex (here all  $d_i \in \mathbb{N}$ ).

Notation: for  $x_0, x_1 \in \mathbb{R}^d$  and  $\lambda \in \mathbb{R}$  we introduce  $x_\lambda = \lambda x_1 + (1 - \lambda)x_0$ .

- a) Statement is false: consider  $g_1(\mathbf{x}) = \|\mathbf{x}\|_2^2$  and  $g_2(x) = -x$ . Both functions are convex, but for  $h(\mathbf{x}) = -\|\mathbf{x}\|_2^2$  it holds with  $\mathbf{x}_0 = \mathbf{0}$ ,  $\mathbf{x}_1 \neq \mathbf{0}$  and  $\lambda \in (0, 1)$  that

$$h(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_0) = -\lambda^2 \|\mathbf{x}_1\|_2^2 \text{ and } \lambda h(\mathbf{x}_1) + (1 - \lambda) h(\mathbf{x}_0) = -\lambda \|\mathbf{x}_1\|_2^2.$$

Since  $\lambda \in (0, 1)$  we get  $\lambda^2 < \lambda \Rightarrow -\lambda^2 \|\mathbf{x}_1\|_2^2 > -\lambda \|\mathbf{x}_1\|_2^2 \Rightarrow h(\mathbf{x}_\lambda) > \lambda h(\mathbf{x}_1) + (1 - \lambda) h(\mathbf{x}_0)$  contradicting the definition of convexity.

- b) Statement is true: consider  $\mathbf{x}_0, \mathbf{x}_1 \in \mathbb{R}^d$  and  $\lambda \in (0, 1)$ . For  $h$  we prove convexity using definitions of the given properties of  $g_1$  and  $g_2$ :

$$\begin{aligned} g_1 \text{ convex} &\Rightarrow g_1(\mathbf{x}_\lambda) \leq \lambda g_1(\mathbf{x}_1) + (1 - \lambda) g_1(\mathbf{x}_0) \\ g_2 \text{ non-decreasing} &\Rightarrow g_2(g_1(\mathbf{x}_\lambda)) \leq g_2(\lambda g_1(\mathbf{x}_1) + (1 - \lambda) g_1(\mathbf{x}_0)) \quad (1) \\ g_2 \text{ convex} &\Rightarrow g_2(\lambda g_1(\mathbf{x}_1) + (1 - \lambda) g_1(\mathbf{x}_0)) \leq \lambda g_2(g_1(\mathbf{x}_1)) + (1 - \lambda) g_2(g_1(\mathbf{x}_0)) \quad (2) \\ (1) \text{ and } (2) &\Rightarrow h(\mathbf{x}_\lambda) \leq \lambda h(\mathbf{x}_1) + (1 - \lambda) h(\mathbf{x}_0) \end{aligned}$$

- c) Statement is true: first, note that for the function  $\max : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\mathbf{y}, \mathbf{z} \in \mathbb{R}^n$  and  $\lambda \geq 0$  we have

$$\max(\lambda \mathbf{y}_1, \dots, \lambda \mathbf{y}_n) = \lambda \max(\mathbf{y}_1, \dots, \mathbf{y}_n) \quad (3)$$

$$\max(\mathbf{y}_1 + \mathbf{z}_1, \dots, \mathbf{y}_n + \mathbf{z}_n) = \bar{\mathbf{y}}_i + \bar{\mathbf{z}}_i \leq \max(\mathbf{y}_1, \dots, \mathbf{y}_n) + \max(\mathbf{z}_1, \dots, \mathbf{z}_n) \quad (4)$$

$$\text{if } \mathbf{y}_i \geq \mathbf{z}_i \text{ for all } i = 1, \dots, N \text{ then } \max(\mathbf{y}_1, \dots, \mathbf{y}_n) \geq \max(\mathbf{z}_1, \dots, \mathbf{z}_n) \quad (5)$$

Now it follows ( $\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_\lambda$  and  $\lambda$  as in b) arbitrary) that

$$\begin{aligned}
 g_i \text{ convex} &\Rightarrow g_i(\mathbf{x}_\lambda) \leq \lambda g_i(\mathbf{x}_1) + (1 - \lambda)g_i(\mathbf{x}_0) \\
 (5) &\Rightarrow \overbrace{\max(g_1(\mathbf{x}_\lambda), \dots, g_n(\mathbf{x}_\lambda))}^{h(\mathbf{x}_\lambda)} \leq \max(\lambda g_1(\mathbf{x}_1) + (1 - \lambda)g_1(\mathbf{x}_0), \dots, \lambda g_n(\mathbf{x}_1) + (1 - \lambda)g_n(\mathbf{x}_0)) \quad (6) \\
 (3) + (4) &\Rightarrow \max(\lambda g_1(\mathbf{x}_1) + (1 - \lambda)g_1(\mathbf{x}_0), \dots, \lambda g_n(\mathbf{x}_1) + (1 - \lambda)g_n(\mathbf{x}_0)) \leq \lambda \underbrace{\max(g_1(\mathbf{x}_1), \dots, g_n(\mathbf{x}_1))}_{h(\mathbf{x}_1)} + (1 - \lambda) \underbrace{\max(g_1(\mathbf{x}_0), \dots, g_n(\mathbf{x}_0))}_{h(\mathbf{x}_0)} \quad (7) \\
 (6) + (7) &\Rightarrow h(\mathbf{x}_\lambda) \leq \lambda h(\mathbf{x}_1) + (1 - \lambda)h(\mathbf{x}_0)
 \end{aligned}$$

## 2 Optimization / Gradient descent

**Problem 2:** You are given the following objective function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$

$$f(x_1, x_2) = 0.5x_1^2 + x_2^2 + 2x_1 + x_2 + \cos(\sin(\sqrt{\pi})).$$

- a) Compute the minimizer  $\mathbf{x}^*$  of  $f$  analytically.

As  $f$  is a sum of convex functions, it is convex. To find the global minimizer, we compute the gradient and set it to zero

$$\nabla f(x_1, x_2) = \begin{pmatrix} x_1 + 2 \\ 2x_2 + 1 \end{pmatrix} \stackrel{!}{=} \begin{pmatrix} 0 \\ 0 \end{pmatrix} \Rightarrow \begin{pmatrix} x_1^* \\ x_2^* \end{pmatrix} = \begin{pmatrix} -2 \\ -\frac{1}{2} \end{pmatrix}.$$

- b) Perform 2 steps of gradient descent on  $f$  starting from the point  $\mathbf{x}^{(0)} = (0, 0)$  with a constant learning rate  $\tau = 1$ .

We already know how to compute the gradient from a).

$$\begin{aligned}
 \text{first step} \quad \begin{pmatrix} x_1^{(1)} \\ x_2^{(1)} \end{pmatrix} &= \begin{pmatrix} x_1^{(0)} \\ x_2^{(0)} \end{pmatrix} - \tau \begin{pmatrix} x_1^{(0)} + 2 \\ 2x_2^{(0)} + 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} - 1 \begin{pmatrix} 0 + 2 \\ 0 + 1 \end{pmatrix} = \begin{pmatrix} -2 \\ -1 \end{pmatrix} \\
 \text{second step} \quad \begin{pmatrix} x_1^{(2)} \\ x_2^{(2)} \end{pmatrix} &= \begin{pmatrix} x_1^{(1)} \\ x_2^{(1)} \end{pmatrix} - \tau \begin{pmatrix} x_1^{(1)} + 2 \\ 2x_2^{(1)} + 1 \end{pmatrix} = \begin{pmatrix} -2 \\ -1 \end{pmatrix} - 1 \begin{pmatrix} -2 + 2 \\ -2 + 1 \end{pmatrix} = \begin{pmatrix} -2 \\ 0 \end{pmatrix}
 \end{aligned}$$

- c) Will the gradient descent procedure from Problem b) ever converge to the true minimizer  $\mathbf{x}^*$ ? Why or why not? If the answer is no, how can we fix it?

By performing one more iteration of gradient descent we observe that

$$\begin{pmatrix} x_1^{(3)} \\ x_2^{(3)} \end{pmatrix} = \begin{pmatrix} x_1^{(2)} \\ x_2^{(2)} \end{pmatrix} - \tau \begin{pmatrix} x_1^{(2)} + 2 \\ 2x_2^{(2)} + 1 \end{pmatrix} = \begin{pmatrix} -2 \\ 0 \end{pmatrix} - 1 \begin{pmatrix} -2 + 2 \\ 0 + 1 \end{pmatrix} = \begin{pmatrix} -2 \\ -1 \end{pmatrix} = \begin{pmatrix} x_1^{(1)} \\ x_2^{(1)} \end{pmatrix}.$$

That is, we are stuck iterating between  $\mathbf{x}^{(1)}$  and  $\mathbf{x}^{(2)}$  forever. We can fix this by decreasing the learning rate (adaptive stepsize, etc.).

**Problem 3:** Load the notebook `06_homework_optimization.ipynb` from Piazza. Fill in the missing code and run the notebook. Convert the evaluated notebook to pdf and add it to the printout of your homework.

*Note: We suggest that you use Anaconda for installing Python and Jupyter, as well as for managing packages. We recommend that you use Python 3.*

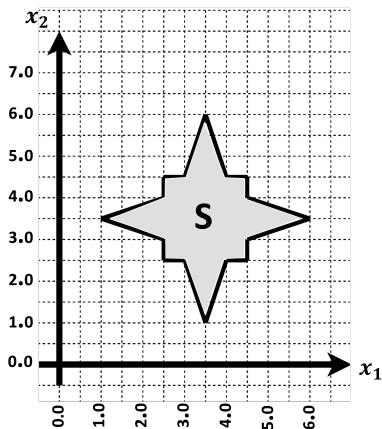
*For more information on Jupyter notebooks and how to convert them to other formats, consult the Jupyter documentation and nbconvert documentation.*

The solution notebook is uploaded on Piazza.

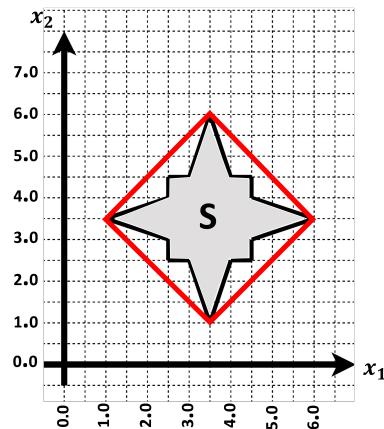
**Problem 4:** Let  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  be the following convex function:

$$f(x_1, x_2) = e^{x_1+x_2} - 5 \cdot \log(x_2)$$

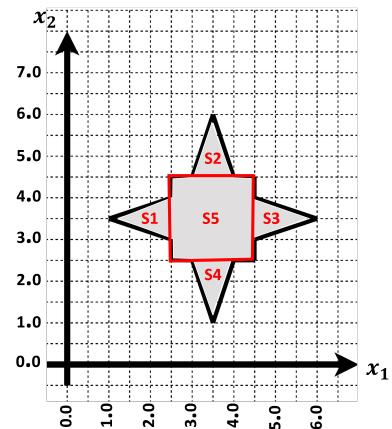
- a) Consider the following shaded region  $S$  in  $\mathbb{R}^2$ . Is this region convex? Why?
- b) Find the maximizer  $\mathbf{x}^*$  of  $f$  over the shaded region  $S$ . Justify your answer.
- c) Assume that we are given an algorithm `ConvOpt`( $f, D$ ) that takes as input a convex function  $f$  and convex region  $D$ , and returns the minimum of  $f$  over  $D$ . Using the `ConvOpt` algorithm, how would you find the global minimum of  $f$  over the shaded region  $S$ ?



(a) initial non-convex set



(b) convex hull



(c) union of convex sets

- a) It is not because we can choose two points in  $\mathbf{S}$  such that the line connecting the points does not completely resides in  $\mathbf{S}$ , for example  $(1.0, 3.5)^T$  and  $(3.5, 6.0)^T$  (see Figure 1b).
- b) According to the lecture, maximizer of the convex function over a non-convex set can be computed by solving the task over its convex hull (see Figure 1b). Since convex hull is convex we need to compute the value of  $f$  only on the four corners  $(1.0, 3.5)^T, (3.5, 1.0)^T, (6.0, 3.5)^T$  and  $(3.5, 6.0)^T$ . Afterwards, we pick the corner at which the value of  $f$  is the largest.

$$\begin{aligned} f(1.0, 3.5) &= e^{4.5} - 5 \cdot \log(3.5) = 83.7533, & f(3.5, 1.0) &= e^{4.5} - 5 \cdot \log(1.0) = 90.01 \\ f(6.0, 3.5) &= e^{9.5} - 5 \cdot \log(3.5) = 13353.46, & f(3.5, 6.0) &= e^{9.5} - 5 \cdot \log(6.0) = 13350.76. \end{aligned}$$

Therefore, the maximizer of  $f$  over  $\mathbf{S}$  is  $(6.0, 3.5)$ .

- c) We can partition the shaded region  $\mathbf{S}$  to the following five convex regions  $\mathbf{S}_1, \dots, \mathbf{S}_5$  (see Figure 1c). Afterwards, we run the `ConvOpt` algorithm separately for the 5 regions and obtain

$$m_i = \min_{\mathbf{x} \in \mathbf{S}_i} f(\mathbf{x}) = \text{ConvOpt}(f, \mathbf{S}_i).$$

Finally, the minimum over the whole  $\mathbf{S}$  can be computed as the smallest of these values, that is  $\min_{\mathbf{x} \in \mathbf{S}} f(\mathbf{x}) = \min(m_1, \dots, m_5)$ .

## In-class Exercises

**Problem 5:** Prove or disprove whether the following functions  $f : D \rightarrow \mathbb{R}$  are convex

- a)  $D = (1, \infty)$  and  $f(x) = \log(x) - x^3$ ,
- b)  $D = \mathbb{R}^+$  and  $f(x) = -\min(\log(3x+1), -x^4 - 3x^2 + 8x - 42)$ ,
- c)  $D = (-10, 10) \times (-10, 10)$  and  $f(x, y) = y \cdot x^3 - y \cdot x^2 + y^2 + y + 4$ .

a) The second derivative of  $f$  is  $\frac{d^2 f(x)}{dx^2} = \frac{d}{dx} \left( \frac{1}{x} - 3x^2 \right) = -\frac{1}{x^2} - 6x$ , which is negative on the given set  $D$  and therefore  $f$  is not convex.

b) Transform min to max

$$-\min\{\log(3x+1), -x^4 - 3x^2 + 8x - 42\} = \max\{-\log(3x+1), x^4 + 3x^2 - 8x + 42\}.$$

$\max(g_1(x), g_2(x))$  is convex if both  $g_1$  and  $g_2$  are convex on  $D = \mathbb{R}^+$  (see Exercise Sheet 6, Problem 1c).  $g_1(x) = -\log(3x+1)$  is convex since the second derivative is positive on  $\mathbb{R}^+$ :

$$\frac{d^2}{dx^2}(-\log(3x+1)) = \frac{d}{dx} \left( -\frac{3}{3x+1} \right) = \frac{9}{(3x+1)^2} > 0$$

$g_2(x) = x^4 + 3x^2 - 8x + 42$  is also convex:

$$\frac{d^2}{dx^2}(x^4 + 3x^2 - 8x + 42) = \frac{d}{dx}(4x^3 + 6x - 8) = 12x^2 + 6 > 0$$

Thus  $f$  is convex.

- c) For the function  $f(x, y)$  to be convex (on  $D$ ) it has to hold for all  $x_1, x_2, y \in D$  and  $\lambda \in (0, 1)$  that

$$\lambda f(x_1, y) + (1 - \lambda)f(x_2, y) \geq f(\lambda x_1 + (1 - \lambda)x_2, y).$$

It does not hold in our case, consider  $y = 1, x_1 = -4, x_2 = 0$  and  $\lambda = 0.5$ :

$$\begin{aligned} 0.5f(-4, 1) + 0.5f(0, 1) &= 0.5 \cdot (-74) + 0.5 \cdot 6 = -34 \\ f(0.5 \cdot (-4) + 0.5 \cdot 0, 0.5 \cdot 1 + 0.5 \cdot 1) &= f(-2, 1) = -6 \color{red} > -34 \end{aligned}$$

Thus  $f(x, y)$  is not convex.

**Problem 6:** Prove that the following functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  are convex:

- a)  $D \subset \mathbb{R}^d$  bounded and closed set and  $f$  is defined by  $f(\mathbf{x}) = \max_{\mathbf{w} \in D} \mathbf{x}^T \mathbf{w}$ .
- b)  $f$  is the objective function of logistic regression, that is

$$f(\mathbf{w}) = -\ln p(\mathbf{y} | \mathbf{w}, \mathbf{X}) = -\sum_{i=1}^N (y_i \ln \sigma(\mathbf{w}^T \mathbf{x}_i) + (1 - y_i) \ln(1 - \sigma(\mathbf{w}^T \mathbf{x}_i))) .$$

- a) We will need some facts similar to (3) and (4) that we used in Problem 1c. Let  $g_1$  and  $g_2$  be two functions that have finite maxima on  $D$  (e.g. later  $g_1(\mathbf{w}) = \lambda \mathbf{x}_1^T \mathbf{w}$  and  $g_2(\mathbf{w}) = (1 - \lambda) \mathbf{x}_0^T \mathbf{w}$ ), then

$$\max_{\mathbf{w} \in D} g_1(\mathbf{w}) + g_2(\mathbf{w}) = g_1(\bar{\mathbf{w}}) + g_2(\bar{\mathbf{w}}) \leq \max_{\mathbf{w} \in D} g_1(\mathbf{w}) + \max_{\mathbf{w} \in D} g_2(\mathbf{w}).$$

With  $\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_\lambda$  and  $\lambda$  as before we get

$$\begin{aligned} f(\mathbf{x}_\lambda) &= \max_{\mathbf{w} \in D} \mathbf{x}_\lambda^T \mathbf{w} = \max_{\mathbf{w} \in D} \overbrace{\lambda \mathbf{x}_1^T \mathbf{w}}^{=g_1(\mathbf{w})} + \overbrace{(1 - \lambda) \mathbf{x}_0^T \mathbf{w}}^{=g_2(\mathbf{w})} \leq \max_{\mathbf{w} \in D} \lambda \mathbf{x}_1^T \mathbf{w} + \max_{\mathbf{w} \in D} (1 - \lambda) \mathbf{x}_0^T \mathbf{w} \\ &= \lambda \max_{\mathbf{w} \in D} \mathbf{x}_1^T \mathbf{w} + (1 - \lambda) \max_{\mathbf{w} \in D} \mathbf{x}_0^T \mathbf{w} \\ &= \lambda f(\mathbf{x}_1) + (1 - \lambda) f(\mathbf{x}_0). \end{aligned}$$

- b) First, notice that if we can prove that the following two functions

$$g_1(\mathbf{w}) = -\ln \sigma(\mathbf{w}^T \mathbf{x}_i) \quad \text{and} \quad g_2(\mathbf{w}) = -\ln(1 - \sigma(\mathbf{w}^T \mathbf{x}_i))$$

are convex for all  $\mathbf{x}_i \in \mathbb{R}^d$ , the whole function  $f$  must also be convex since any linear combination (with non-negative weights, which is true since  $y_i$  and  $1 - y_i$  are non-negative) of two or more convex functions is also convex.

To prove that the first function is convex we will use the second-order condition of convexity. A function  $f(x)$  which is twice-differentiable is convex if and only if its Hessian matrix (matrix of second-order partial derivatives) is positive semi-definite (see Problem 8). To compute the Hessian matrix we first calculate the derivative of the sigmoid function:

$$\begin{aligned} \frac{d}{dx} \sigma(x) &= \frac{e^{-x}}{(1 + e^{-x})^2} = \sigma(x) \frac{e^{-x}}{1 + e^{-x}} = \sigma(x) \left( 1 + \frac{e^{-x} - 1 - e^{-x}}{1 + e^{-x}} \right) = \sigma(x) \left( 1 - \frac{1}{1 + e^{-x}} \right) \\ &= \sigma(x)(1 - \sigma(x)). \end{aligned}$$

Using this, we can derive the Hessian:

$$\begin{aligned} \nabla_{\mathbf{w}}^2 [-\ln \sigma(\mathbf{w}^T \mathbf{x}_i)] &= \nabla_{\mathbf{w}} [\nabla_{\mathbf{w}} (-\ln \sigma(\mathbf{w}^T \mathbf{x}_i))] \\ &= \nabla_{\mathbf{w}} \left[ -\mathbf{x}_i \frac{\sigma(\mathbf{w}^T \mathbf{x}_i)(1 - \sigma(\mathbf{w}^T \mathbf{x}_i))}{\sigma(\mathbf{w}^T \mathbf{x}_i)} \right] \\ &= \nabla_{\mathbf{w}} [\mathbf{x}_i (\sigma(\mathbf{w}^T \mathbf{x}_i) - 1)] \\ &= \sigma(\mathbf{w}^T \mathbf{x}_i)(1 - \sigma(\mathbf{w}^T \mathbf{x}_i)) \mathbf{x}_i \mathbf{x}_i^T. \end{aligned}$$

Next, we show that this Hessian matrix is positive semi-definite:

$$\begin{aligned} \forall \mathbf{z} : \quad &\mathbf{z}^T \nabla_{\mathbf{w}}^2 [-\ln \sigma(\mathbf{w}^T \mathbf{x}_i)] \mathbf{z} \\ &= \mathbf{z}^T [\sigma(\mathbf{w}^T \mathbf{x}_i)(1 - \sigma(\mathbf{w}^T \mathbf{x}_i)) \mathbf{x}_i \mathbf{x}_i^T] \mathbf{z} \\ &= \sigma(\mathbf{w}^T \mathbf{x}_i)(1 - \sigma(\mathbf{w}^T \mathbf{x}_i)) (\mathbf{x}_i^T \mathbf{z})^2 \geq 0 \end{aligned}$$

To prove that the second function is convex, we first notice:

$$\begin{aligned}-\ln(1 - \sigma(\mathbf{w}^T \mathbf{x}_i)) &= -\ln\left(1 - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}}\right) = -\ln\left(\frac{e^{-\mathbf{w}^T \mathbf{x}_i}}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}}\right) \\ &= \mathbf{w}^T \mathbf{x}_i - \ln\left(\frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}}\right) = \mathbf{w}^T \mathbf{x}_i - \ln \sigma(\mathbf{w}^T \mathbf{x}_i)\end{aligned}$$

This is a sum of two convex functions, since the affine function  $\mathbf{w}^T \mathbf{x}_i$  is convex and we just showed that  $-\ln \sigma(\mathbf{w}^T \mathbf{x}_i)$  is convex. Hence,  $-\ln(1 - \sigma(\mathbf{w}^T \mathbf{x}_i))$  is convex as well.

**Problem 7:** Prove that for differentiable convex functions each local minimum is a global minimum. More specifically, given a differentiable convex function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , prove that if there is a local minimum at  $\mathbf{x}^*$  then  $\nabla f(\mathbf{x}^*) = 0$  and if  $\nabla f(\mathbf{x}^*) = 0$  then  $\mathbf{x}^*$  is a global minimum.

At a (local) minimum point  $\boldsymbol{\theta}^*$  the gradient must be zero. Otherwise we could follow the gradient to get an even lower value.

$$f(\boldsymbol{\theta}^* - \epsilon \nabla f(\boldsymbol{\theta}^*)) = f(\boldsymbol{\theta}^*) - \epsilon \|\nabla f(\boldsymbol{\theta}^*)\|_2^2 + O(\epsilon^2 \|\nabla f(\boldsymbol{\theta}^*)\|_2^2) < f(\boldsymbol{\theta}^*)$$

for sufficiently small  $\epsilon$ .

Now, using the first order criterion we have

$$f(\mathbf{y}) \geq f(\mathbf{x}) + (\mathbf{y} - \mathbf{x}) \nabla f(\mathbf{x})$$

If we replace  $\mathbf{x}$  with  $\boldsymbol{\theta}^*$  and plug in  $\nabla f(\boldsymbol{\theta}^*) = 0$  we get:  $f(\mathbf{y}) \geq f(\boldsymbol{\theta}^*)$  for all  $\mathbf{y}$ , meaning  $\boldsymbol{\theta}^*$  is a global minimum.

**Problem 8:** Show that a twice differentiable function  $f : D \rightarrow \mathbb{R}$  with a convex domain  $D \subset \mathbb{R}^d$  is convex if and only if its Hessian is positive semi-definite on the whole domain  $D$ .

For all  $\mathbf{x}_0, \mathbf{x}_1 \in D$  and some  $\mathbf{x}_\lambda$  as before for a  $\lambda \in [0, 1]$  (depending on  $\mathbf{x}_0, \mathbf{x}_1$ ) it holds that

$$f(\mathbf{x}_1) = f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T (\mathbf{x}_1 - \mathbf{x}_0) + \frac{1}{2} (\mathbf{x}_1 - \mathbf{x}_0)^T \nabla^2 f(\mathbf{x}_\lambda) (\mathbf{x}_1 - \mathbf{x}_0). \quad (8)$$

If  $\nabla^2 f(\mathbf{x})$  is positive semi-definite for all  $\mathbf{x} \in D$  we get that the quadratic term in (8) is non-negative and therefore  $f(\mathbf{x}_1) \geq f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T (\mathbf{x}_1 - \mathbf{x}_0)$  which is the first-order convexity condition.

Now assume  $\nabla^2 f(\mathbf{x})$  is not positive semi-definite for some  $\mathbf{x} \in D$ , so there exists  $d$  such that  $d^T \nabla^2 f(\mathbf{x}) d < 0$ . Consider (8) for  $\mathbf{x}_0 = \mathbf{x}$  and  $\mathbf{x}_1 = \mathbf{x} + \epsilon \mathbf{d}$  for  $\epsilon > 0$ . For sufficiently small  $\epsilon$  we get  $\mathbf{d}^T \nabla^2 f(\mathbf{x}_\lambda) \mathbf{d} < 0$  since  $\nabla^2 f$  is continuous. Therefore, the quadratic term in (8) for our choice of  $\mathbf{x}_0$  and  $\mathbf{x}_1$  is now negative and we get  $f(\mathbf{x}_1) < f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T (\mathbf{x}_1 - \mathbf{x}_0)$  contradicting the first-order convexity condition.