

## Practical Session 5

### Optimization

---

## 1 Convexity

**Problem 1:** Show that affine functions of the form  $\mathbf{w}^T \mathbf{x} + b$  are both convex and concave.

A function is convex iff  $\lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) \geq f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y})$  and concave iff  $\lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) \leq f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y})$ . Hence, for a function to be both convex and concave it must hold that

$$\lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) = f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}).$$

We have

$$\begin{aligned}\lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) &= \lambda[\mathbf{w}^T \mathbf{x} + b] + (1 - \lambda)[\mathbf{w}^T \mathbf{y} + b] \\ &= \lambda \mathbf{w}^T \mathbf{x} + (1 - \lambda)\mathbf{w}^T \mathbf{y} + \lambda b + (1 - \lambda)b \\ &= \lambda \mathbf{w}^T \mathbf{x} + (1 - \lambda)\mathbf{w}^T \mathbf{y} + b \\ &= \mathbf{w}^T [\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}] + b \\ &= f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y})\end{aligned}\quad \square$$

**Problem 2:** Show that a twice differentiable function  $f(\mathbf{x})$  with a convex domain is convex if and only if its Hessian or second derivative is positive semidefinite:  $\nabla^2 f(\mathbf{x}) \geq 0$  for all  $\mathbf{x} \in \text{dom}(f)$ .

We first assume we have a single dimension ( $n = 1$ ). Suppose  $f: \mathbb{R} \rightarrow \mathbb{R}$  is convex. Let  $x, y \in \text{dom}(f)$ .

By the first-order condition,

$$\begin{aligned}f'(x)(y - x) &\leq f(y) - f(x), \\ f'(y)(x - y) &\leq f(x) - f(y) \quad \Leftrightarrow \quad f(y) - f(x) \leq f'(y)(y - x).\end{aligned}$$

Hence,

$$f'(x)(y - x) \leq f'(y)(y - x).$$

Subtracting the lefthand side from the righthand side and dividing by  $(y - x)^2$  gives:

$$\frac{f'(y) - f'(x)}{y - x} \geq 0$$

Taking the limit for  $y \rightarrow x$  yields  $f''(x) \geq 0$ , for any  $x \in \text{dom}(f)$ .

---

Conversely, suppose  $f''(z) \geq 0$  for all  $z \in \text{dom}(f)$ . Consider two arbitrary points  $x, y \in \text{dom}(f)$ . Without loss of generality we assume that  $x < y$ . We have

$$\begin{aligned} 0 &\leq \int_x^y f''(z)(y-z) \, dz \\ &= (f'(z)(y-z)) \Big|_{z=x}^{z=y} + \int_x^y f'(z) \, dz \\ &= -f'(x)(y-x) + f(y) - f(x) \end{aligned}$$

i.e.  $f(y) \geq f(x) + f'(x)(y-x)$ , which is the first order convexity condition and shows that  $f$  is convex.

To generalize to  $n > 1$ , we note that a function is convex if and only if it is convex on all lines, i.e. iff the one-dimensional function  $g(t) = f(\mathbf{x}_0 + t\mathbf{v})$  is convex in  $t$  for all  $\mathbf{x}_0 \in \text{dom}(f)$  and all  $\mathbf{v} \in \mathbb{R}^n$ , for values satisfying  $\mathbf{x}_0 + t\mathbf{v} \in \text{dom}(f)$ . Therefore,  $f$  is convex if and only if

$$g''(t) = \mathbf{v}^T \nabla^2 f(\mathbf{x}_0 + t\mathbf{v}) \mathbf{v} \geq 0.$$

In other words, it is necessary and sufficient that  $\mathbf{v}^T \nabla^2 f(\mathbf{x}) \mathbf{v} \geq 0$  for all  $\mathbf{x} = \mathbf{x}_0 + t\mathbf{v} \in \text{dom}(f)$ , which is exactly the definition of the Hessian  $\nabla^2 f(\mathbf{x})$  being positive semi-definite.  $\square$

*Note:* For the strictly convex case one can show that if  $\nabla^2 f(\mathbf{x})$  is positive definite then  $f$  is strictly convex. The converse, however, is not true: For example, the function  $f: \mathbb{R} \rightarrow \mathbb{R}$  given by  $f(x) = x^4$  is strictly convex but has zero second derivative at  $x = 0$ .

## 2 Logistic Regression

**Problem 3:** Prove that the objective function of logistic regression

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

is convex. What is the benefit of having a convex function for optimization?

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

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

are convex, our objective function as given in Eq.1 must also be convex since any linear combination (with positive constants) of two or more convex combinations is also convex. Since  $y_i$  and  $1 - y_i$  are positive this holds.

To prove that the first function is convex we will use the second-order condition of convexity.

*Reminder:* 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.

To compute the Hessian matrix we first calculate the derivative of the sigmoid function:

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

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.  $\square$

The benefit of having a convex function for optimization is that, subject to relatively mild assumptions, stochastic gradient descent converges almost surely to a global minimum.

### 3 Optimization methods

**Problem 4:** Discuss the following topics:

- Condition number
- Consistency, convergence, stability
- Stiffness

- Condition number: How much does output vary depending on the input, more precisely the maximum ratio of the relative error in the output  $\hat{y}$  due to the relative error in the input  $x$ . Ill-conditioned: High condition number. This is a property of the problem itself, not of the algorithm.
- Consistency: Local discretization error (error due to a single step)  $l(\delta t) \rightarrow 0$  for  $\delta t \rightarrow 0$
- Convergence: Global discretization error (overall error)  $e(\delta t) \rightarrow 0$  for  $\delta t \rightarrow 0$
- Stability: Algorithms that do not magnify approximation errors. Instabilities can be caused e.g. by nearby singularities (e.g. very small eigenvalues), truncation errors, or loss of significance. Stability + consistency = convergence
- Stiffness: Local property of the algorithm's solution imposes extremely high resolution.