# Ex 11

Exercises for the course

## Machine Learning 2

Summer semester 2017

Abteilung Maschinelles Lernen
Institut für Softwaretechnik und theoretische Informatik
Fakultät IV, Technische Universität Berlin
Prof. Dr. Klaus-Robert Müller
Email: klaus-robert.mueller@tu-berlin.de

## Exercise Sheet 11

In this exercise sheet, we refer to the sections of the paper "*Methods for interpreting and understanding deep neural networks*" linked via ISIS.

### Exercise 1: Experts and Prototypes (40 P)

Consider the linear model $y = \boldsymbol{w}^\top \boldsymbol{x} + b$ mapping some input $\boldsymbol{x}$ to an output $y$. We would like to interpret the output $y$ by building a prototype $\boldsymbol{x}^*$ in the input domain following the activation maximization techniques outlined in Section 3.

(a) *Find* the prototype $\boldsymbol{x}^*$ obtained by activation maximization as formulated in Section 3.1.

(b) *Find* the prototype $\boldsymbol{x}^*$ obtained by activation maximization as formulated in Section 3.2. We assume that the data is represented by the Gaussian expert $p(\boldsymbol{x}) = \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ where $\boldsymbol{\mu}$ and $\Sigma$ are the mean and covariance.

(c) *Find* the prototype $\boldsymbol{x}^*$ obtained by activation maximization as formulated in Section 3.3. The data is generated as (i) $\boldsymbol{z} \sim \mathcal{N}(0, I)$ and (ii) $\boldsymbol{x} = A\boldsymbol{z} + \boldsymbol{c}$, where $A$ and $\boldsymbol{c}$ are the parameters of the generator.

(d) *Relate* the prototypes obtained for the three approaches above, in particular under which regularizers, experts and generator, the found prototypes are mutually equivalent.
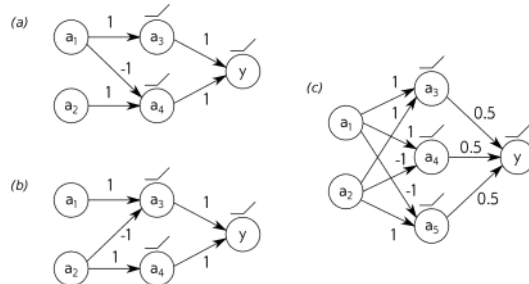
### Exercise 2: Sensitivity Analysis and Taylor Decomposition (30 P)

Let us consider a data point $\boldsymbol{x}$ and its prediction by a homogeneous linear model $f(\boldsymbol{x}) = \boldsymbol{w}^\top \boldsymbol{x}$. We would like to explain the prediction using the methods described in Section 4.

(a) *Compute* the explanation for the prediction $f(\boldsymbol{x})$ using sensitivity analysis as described in Section 4.1.

(b) *Compute* the explanation for the prediction $f(\boldsymbol{x})$ using Taylor decomposition (Section 4.2) at root point $\tilde{\boldsymbol{x}} = \boldsymbol{0}$.

(c) *Compute* the explanation for the prediction $f(\boldsymbol{x})$ using Taylor decomposition at root point $\tilde{\boldsymbol{x}}$ chosen to be the nearest (in the Euclidean sense) from $\boldsymbol{x}$. (Hint: You can use the Lagrange multipliers to find this root point.)

### Exercise 3: Layer-Wise Relevance Propagation (30 P)

We would like to test the dependence of layer-wise relevance propagation (LRP) on the structure of the neural network. For this, we consider the function $y = \max(a_1, a_2)$, where $a_1, a_2 \in \mathbb{R}^+$ are the input activations. This function can be implemented as a ReLU network in multiple ways. Three examples are given below.



Because of the positive activations, an appropriate rule for both layers is LRP-$\alpha_1\beta_0$ defined in Section 5.1.

(a) *Give* for each network an analytic solution for the obtained scores $R_1$ and $R_2$ obtained by application this propagation rule at each layer.

(b) *Discuss* which implementation of the "max" function (a, b, or c) gives the most intuitive explanations.

---

**1**

linear model    $y = \vec{w} \cdot \vec{x} + b$    $\longrightarrow$ find a prototype

**a)**

Consider a DNN classifier mapping data points $\boldsymbol{x}$ to a set of classes $(\omega_c)_c$. The output neurons encode the modeled class probabilities $p(\omega_c | \boldsymbol{x})$. A prototype $\boldsymbol{x}^*$ representative of the class $\omega_c$ can be found by optimizing:

$$\max_{\boldsymbol{x}} \ \log p(\omega_c | \boldsymbol{x}) - \lambda \|\boldsymbol{x}\|^2.$$

wie hängt  $y$  mit  $p(\omega_c | x)$  zusammen?     $y = \log p(\omega_+ | x)$ ?

$\longrightarrow$ prototype could be    $\underset{x}{\arg\max} \ y(\vec{x}) - \lambda \|\vec{x}\|^2$    // $\lambda$ ist hier kein Lagrange-Multiplikator

$\rightarrow$ prototype could be $\underset{x}{\arg\max}\ y(\vec{x}) - \lambda\underbrace{\|\vec{x}\|^2}$     // $\lambda$ ist hier kein Lagrange-Multiplikator

$$\text{mit}\quad L = y(\vec{x}) - \lambda\|\vec{x}\|^2 \quad\rightarrow\quad \frac{\partial L}{\partial \vec{x}} = \vec{w} - 2\lambda\vec{x} \overset{!}{=} 0 \qquad \rightarrow \vec{x}^* = \frac{\vec{w}}{2\lambda}$$

b)   $p(x) \sim \mathcal{N}(\vec{\mu}, \Sigma)$    $\underset{x}{\max}\ \log p(\omega_c|x) + \log p(x).$

$$\arg\underset{\vec{x}}{\max}\ \vec{w}\cdot\vec{x} + b + \log\left(\frac{1}{\sqrt{|2\pi\Sigma|}}\exp\left(-\frac{1}{2}(\vec{x}-\vec{\mu})^T\Sigma^{-1}(\vec{x}-\vec{\mu})\right)\right)$$

$$\frac{\partial L}{\partial \vec{x}} = \vec{w}^T + \frac{1}{\mathcal{N}(\mu,\Sigma)}\cdot\mathcal{N}(\mu,\Sigma)\cdot\frac{\partial}{\partial\vec{x}}\left(-\frac{1}{2}(\vec{x}-\vec{\mu})^T\Sigma^{-1}(\vec{x}-\vec{\mu})\right)$$

$$= \vec{w}^T - \frac{1}{2}\frac{\partial}{\partial x}\left[(\vec{x}-\vec{\mu})^T\Sigma^{-1}(\vec{x}-\vec{\mu})\right]$$

$$= \vec{w}^T - (\vec{x}-\vec{\mu})^T\Sigma^{-1} \overset{!}{=} 0$$

$$\vec{w}^T = (\vec{x}-\vec{\mu})^T\Sigma^{-1} = \vec{x}^T\Sigma^{-1} - \vec{\mu}^T\Sigma^{-1}$$

$$\vec{w}^T + \vec{\mu}^T\Sigma^{-1} = \vec{x}^T\Sigma^{-1} \qquad |\cdot\Sigma$$

$$\vec{w}^T\Sigma + \vec{\mu}^T = \vec{x}^T \qquad\qquad |^T$$

$$\Sigma^T\vec{w} + \vec{\mu} = \vec{x}^*$$

use: $\dfrac{\partial(\vec{x}^T A\vec{x})}{\partial\vec{x}} = 2\vec{x}^T A$

if $A$ is symmetric

c)

$\vec{x} = A\vec{z} + \vec{c}$    code

$\underset{z\in Z}{\max}\ \log p(\omega_c\,|\,g(z)) - \lambda\|z\|^2,$

optimize in code space

$$\frac{\partial}{\partial\vec{z}}\left(\vec{w}\cdot(A\vec{z}+\vec{c}) + b) - \lambda\|z\|^2\right)$$

$$= \frac{\partial}{\partial\vec{z}}\left(\underbrace{\vec{w}^T\cdot A}_{\vec{a}^T}\vec{z}\right) - 2\lambda\vec{z}$$

$$= A\cdot\vec{w} - 2\lambda\vec{z} \quad\rightarrow\quad \vec{z}^* = \frac{A^T\cdot\vec{w}}{2\lambda} \quad \text{(like part a if } A = \mathbb{1})$$

$\dfrac{\partial(\vec{a}^T\vec{x})}{\partial\vec{x}} = \vec{a}$

2    $f(x) = \vec{w}^T\vec{x}$

a)   $R_i(x) = \left(\frac{\partial f}{\partial x_i}\right)^2.$    $\rightarrow \frac{\partial f}{\partial x_i} = \frac{\partial}{\partial x_i}\sum_j w_j x_j = w_i$

$\rightarrow R_i(\vec{x}) = w_i^2$

b)   $R_i(x) = \frac{\partial f}{\partial x_i}\cdot x_i.$    $\rightarrow R_i(\vec{x}) = \frac{\partial f}{\partial x_i}\Big|_{\vec{x}=\vec{x}_0}\cdot(x_i - x_{0i}) = w_i x_i$

**b)** $R_i(\boldsymbol{x}) = \frac{\partial f}{\partial x_i} \cdot x_i.$ $\quad\longrightarrow\quad$ $R_i(\vec{x}) = \left.\frac{\partial f}{\partial x_i}\right|_{\vec{x}=\vec{x}_0} \cdot (x_i - x_{0i}) = w_i x_i$

**c)**

$\quad$ min distance while $\vec{\omega}^T \vec{x} = 0$

$\quad\to\quad L = \|\vec{x} - \vec{x}_0\|^2 + \lambda \vec{\omega}^T \vec{x}_0$

$\frac{\partial L}{\partial \vec{x}_0} = -2(\vec{x} - \vec{x}_0) + \lambda \vec{\omega} \stackrel{!}{=} 0 \quad\to\quad \vec{x} - \frac{\lambda}{2}\vec{\omega} = \vec{x}_0 \quad \Big|\ \vec{\omega}^T \cdot$

$\quad\quad \vec{\omega}^T \vec{x} = \frac{\lambda}{2}\|\vec{\omega}\|^2$

$\Longrightarrow \quad \vec{x} - \frac{\lambda}{2}\vec{\omega} = \vec{x}_0 \ , \quad\quad \lambda = \frac{2 \cdot \vec{\omega}^T \vec{x}}{\omega^2}$

$\quad\to\quad \vec{x}_0 = \vec{x} - \frac{(\vec{\omega}^T \cdot \vec{x})\,\vec{\omega}}{\|\vec{\omega}\|^2} = \vec{x} - (\hat{\omega}^T \vec{x})\,\hat{\omega}$

$R_i(\vec{x}) = \left.\frac{\partial f}{\partial x_i}\right|_{\vec{x}=\vec{x}_0} \cdot (x_i - x_{0i}) = w_i \cdot \left( x_i - x_i + \frac{1}{\omega^2}(\vec{\omega}^T \cdot \vec{x})\,w_i \right)$

$\quad\quad\quad\quad\quad\quad = \frac{w_i^2}{\|\vec{\omega}\|^2}(\vec{\omega}^T \cdot \vec{x})$