

## Exercise Sheet 11

### Exercise 1: Activation Maximization (20 P)

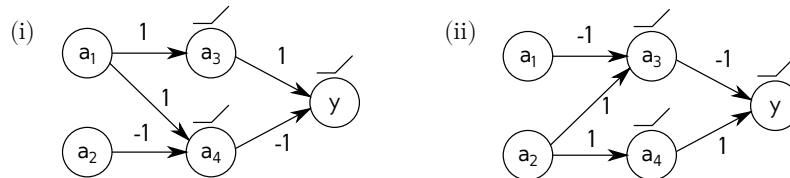
Consider the linear model  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$  mapping some input  $\mathbf{x}$  to an output  $f(\mathbf{x})$ . We would like to interpret the function  $f$  by building a prototype  $\mathbf{x}^*$  in the input domain which produces a large value  $f$ . Activation maximization produces such interpretation by optimizing

$$\max_{\mathbf{x}} [f(\mathbf{x}) - \Omega(\mathbf{x})].$$

- Find the prototype  $\mathbf{x}^*$  obtained by activation maximization subject to the penalty  $\Omega(\mathbf{x}) = \lambda \|\mathbf{x}\|^2$ .
- Find the prototype  $\mathbf{x}^*$  obtained by activation maximization subject to the penalty  $\Omega(\mathbf{x}) = -\log p(\mathbf{x})$  with  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$  where  $\boldsymbol{\mu}$  and  $\Sigma$  are the mean and covariance.
- Find the prototype  $\mathbf{x}^*$  obtained when the data is generated as (i)  $\mathbf{z} \sim \mathcal{N}(0, I)$  and (ii)  $\mathbf{x} = A\mathbf{z} + \mathbf{c}$ , with  $A$  and  $\mathbf{c}$  the parameters of the generator. Here, we optimize  $f$  w.r.t. the code  $\mathbf{z}$  subject to the penalty  $\Omega(\mathbf{z}) = \lambda \|\mathbf{z}\|^2$ .

### Exercise 2: 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  $y = \min(a_1, a_2)$ , where  $a_1, a_2 \in \mathbb{R}^+$  are the input activations. This function can be implemented as a ReLU network in multiple ways. Two examples are given below.



- Show that these two networks implement the ‘min’ function on the relevant domain.
- We consider the LRP- $\gamma$  propagation rule:

$$R_j = \sum_k \frac{a_j \cdot (w_{jk} + \gamma w_{jk}^+)}{\sum_j a_j \cdot (w_{jk} + \gamma w_{jk}^+)} R_k$$

where  $()^+$  denotes the positive part. For each network, give for the case  $a_1 = a_2$  an analytic solution for the scores  $R_1$  obtained by application this propagation rule at each layer. More specifically, express  $R_1$  as a function of the input activations.

### Exercise 3: Neuralization (20 P)

Consider the one-class SVM that predicts for every new data point  $\mathbf{x}$  the ‘inlierness’ score:

$$f(\mathbf{x}) = \sum_{i=1}^M \alpha_i k(\mathbf{x}, \mathbf{u}_i)$$

where  $(\mathbf{u}_i)_{i=1}^M$  is the collection of support vectors, and  $\alpha_i > 0$  are their weightings. We use the Gaussian kernel  $k(\mathbf{x}, \mathbf{x}') = \exp(-\gamma \|\mathbf{x} - \mathbf{x}'\|^2)$ .

Because we are typically interested in the degree of anomaly of a particular data point, we can also define the score  $o(\mathbf{x}) = -\frac{1}{\gamma} \log f(\mathbf{x})$  which grows with the degree of anomaly of the data point.

(a) Show that the outlier score  $o(\mathbf{x})$  can be rewritten as a two-layer neural network:

$$h_i = \|\mathbf{x} - \mathbf{u}_i\|^2 - \gamma^{-1} \log \alpha_i \quad (\text{layer 1})$$

$$o(\mathbf{x}) = -\frac{1}{\gamma} \log \sum_{i=1}^M \exp(-\gamma h_i) \quad (\text{layer 2})$$

(b) Show that the layer 2 converges to a min-pooling (i.e.  $o(\mathbf{x}) = \min_{i=1}^N \{h_i\}$ ) in the limit of  $\gamma \rightarrow \infty$ .

#### **Exercise 4: Programming (30 P)**

Download the programming files on ISIS and follow the instructions.

## Exercise Sheet 11

### 1 Activation Maximization

$$a) \quad \frac{\partial}{\partial x} \max_x w^T x + b - \lambda \|x\|^2 \stackrel{!}{=} 0$$

$$\Leftrightarrow w^T - 2\lambda x \stackrel{!}{=} 0 \quad \Leftrightarrow x^* = \frac{w^T}{2\lambda}$$

$$b) \quad \max_x w^T x + b + \log(p(x))$$

$$= \max_x w^T x + b + \log\left(\frac{\exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu))}{(2\pi)^{\frac{d}{2}} |\Sigma|}\right)$$

$$= \max_x w^T x + b + (\log(\exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu))) - \log((2\pi)^{\frac{d}{2}} |\Sigma|))$$

$$\Rightarrow \frac{\partial}{\partial x} w^T x - \frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu) \stackrel{!}{=} 0$$

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

$$\Leftrightarrow w^T - \Sigma^{-1}x + \Sigma^{-1}\mu = 0$$

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

$$\Leftrightarrow x^* = \mu + \Sigma w^T$$

$$c) \quad \frac{\partial}{\partial z} \max_z w^T (Az + c) + b - \lambda \|z\|^2$$

$$\Rightarrow \frac{\partial}{\partial z} w^T (Az + c) + b - \lambda \|z\|^2 \stackrel{!}{=} 0$$

$$\Leftrightarrow w^T A - 2\lambda z = 0$$

$$\Leftrightarrow z = \frac{w^T A}{2\lambda} \quad | \quad x = Az + c \Leftrightarrow z = \frac{x-c}{A}$$

$$\Leftrightarrow \frac{x-c}{A} = \frac{w^T A}{2\lambda}$$

$$\Leftrightarrow x = \frac{w^T A A}{2\lambda} + c$$



## 2 Layer - Info Relevance Propagation

a)  $y = a_3 - a_4$   
 $\quad \quad \quad \underbrace{\quad}_{a_1} \quad \underbrace{\quad}_{\max(0, a_1 - a_2)}$

ReLU harm. effect  
 i.e. as in positive

if  $a_1 \geq a_2$  :  $y = a_1 - (a_1 - a_2) = a_2$

if  $a_1 \leq a_2$  :  $y = a_1 - 0 = a_1$

b) i)  $R_y = a$

$R_3 = a$

$R_1 = a$

ii)  $R_y = a$

$R_3 = 0$

$R_1 = 0$

## 3 Minimization

a) 
$$O(x) = -\frac{1}{p} \log \left( \sum_{i=1}^M \alpha_i e^{-\mu \|x - \mu_i\|^2} \right)$$
  

$$= -\frac{1}{p} \log \left( \sum_{i=1}^M e^{-\mu \|x - \mu_i\|^2} \right) - \frac{1}{p} \log(\alpha_i)$$

b)  $\min_{i=1}^M \{h_i\} = m$

$$\begin{aligned} O(x) - m &= -\frac{1}{p} \log \left( \sum_{i=1}^M e^{-\mu h_i} \right) + \frac{1}{p} \log e^{-\mu m} \\ &= -\frac{1}{p} \log \left( \sum_{i=1}^M e^{-\mu h_i} e^{\mu m} \right) \\ &= -\frac{1}{p} \log \left( \sum_{i=1}^M e^{-\mu (h_i - m)} \right) \\ &= -\frac{1}{p} \log e^{-\mu \cdot 0} + \sum_{i=1}^M e^{-\mu (h_i - m)} \\ &= -\frac{1}{p} \log (1 + \sum_{i=1}^M e^{-\mu (h_i - m)}) \end{aligned}$$

$O(x) - m = 0 \quad \text{for } \mu \rightarrow \infty$

$\Leftrightarrow O(x) = m \rightarrow \text{min - pooling}$