

Exercise Sheet 14

Exercise 1: Class Prototypes (25 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 $\Omega(\mathbf{x}) = \log p(\mathbf{x})$ with $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ where $\boldsymbol{\mu}$ and Σ are the mean and covariance.

Exercise 2: Shapley Values (25 P)

Consider the function $f(\mathbf{x}) = \min(x_1, \max(x_2, x_3))$. Compute the Shapley values ϕ_1, ϕ_2, ϕ_3 for the prediction $f(\mathbf{x})$ with $\mathbf{x} = (1, 1, 1)$. (We assume a reference point $\tilde{\mathbf{x}} = \mathbf{0}$, i.e. we set features to zero when removing them from the coalition).

Exercise 3: Taylor Expansions (25 P)

Consider the simple radial basis function

$$f(\mathbf{x}) = \|\mathbf{x} - \boldsymbol{\mu}\| - \theta$$

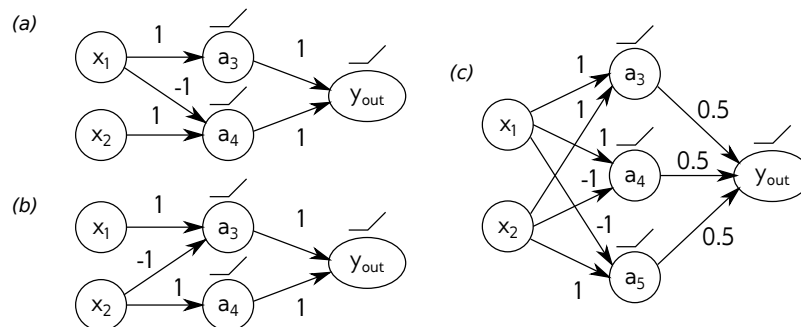
with $\theta > 0$. For the purpose of extracting an explanation, we would like to build a first-order Taylor expansion of the function at some root point $\tilde{\mathbf{x}}$. We choose this root point to be taken on the segment connecting $\boldsymbol{\mu}$ and \mathbf{x} (we assume that $f(\mathbf{x}) > 0$ so that there is always a root point on this segment).

Show that the first-order terms of the Taylor expansion are given by

$$\phi_i = \frac{(x_i - \mu_i)^2}{\|\mathbf{x} - \boldsymbol{\mu}\|^2} \cdot (\|\mathbf{x} - \boldsymbol{\mu}\| - \theta)$$

Exercise 4: Layer-Wise Relevance Propagation (25 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 $y = \max(x_1, x_2)$, where $x_1, x_2 \in \mathbb{R}^+$ are the input activations. This function can be implemented as a ReLU network in multiple ways. Three examples are given below.



We consider the propagation rule:

$$R_j = \sum_k \frac{a_j w_{jk}^+}{\sum_j a_j w_{jk}^+} R_k$$

where j and k are indices for two consecutive layers and where $()^+$ denotes the positive part. This propagation rule is applied to both layers.

Give for each network the computational steps that lead to the scores R_1 and R_2 , and the obtained relevance values. More specifically, express R_1 and R_2 as a function of R_3 and R_4 (and R_5), and express the latter relevances as a function of $R_{\text{out}} = y$.

Exercise 1: Class Prototypes (25 P)

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

$$\max_x [f(x) + \Omega(x)].$$

Find the prototype x^* obtained by activation maximization subject to $\Omega(x) = \log p(x)$ with $x \sim \mathcal{N}(\mu, \Sigma)$ where μ and Σ are the mean and covariance.

$$\begin{aligned} f(x) + \Omega(x) &= w^T x + b + \log p(x) \\ &= w^T x + b + \log \frac{1}{(2\pi)^R \det(\Sigma)} \cdot \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right) \\ &= w^T x + b - \frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu) + \underbrace{\log \frac{1}{(2\pi)^R \det(\Sigma)}}_{= \text{const}} \end{aligned}$$

$$\Rightarrow \max_x \{f(x) + \Omega(x)\} = \max_x \left\{ w^T x + b - \frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu) \right\}$$

$$\nabla_x \left(w^T x + b - \frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu) \right) = w - \Sigma^{-1}(x-\mu) = \vec{0}$$

$$\Sigma \cdot \Sigma^{-1}(x-\mu) = \Sigma w \quad (\text{To simplify the calculation, we assume } \Sigma \text{ is PD})$$

$$x - \mu = \Sigma w$$

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

Exercise 2: Shapley Values (25 P)

Consider the function $f(x) = \min(x_1, \max(x_2, x_3))$. Compute the Shapley values ϕ_1, ϕ_2, ϕ_3 for the prediction $f(x)$ with $x = (1, 1, 1)$. (We assume a reference point $\tilde{x} = \mathbf{0}$, i.e. we set features to zero when removing them from the coalition).

$$\phi_i = \sum_{S: i \notin S} \underbrace{\frac{|S|!(d-|S|-1)!}{d!}}_{\alpha_S} \underbrace{[f(x_{S \cup \{i\}}) - f(x_S)]}_{\Delta_S}$$

$\phi_1: S$	α_S	Δ_S	
$\{\}$	$\frac{0!(3-0-1)!}{3!} = \frac{1}{3}$	$f(1,0,0) - f(0,0,0) = 0 - 0 = 0$	$\phi_1 = \frac{1}{3} \cdot 0 + \frac{1}{6} \cdot 1 + \frac{1}{6} \cdot 1 + \frac{1}{3} \cdot 1$
$\{2\}$	$\frac{1!(3-2-1)!}{2!} = \frac{1}{6}$	$f(1,1,0) - f(0,1,0) = 1 - 0 = 1$	
$\{3\}$	$\frac{1}{6}$	$f(1,0,1) - f(0,0,1) = 1 - 0 = 1$	
$\{2,3\}$	$\frac{2!(3-2-1)!}{2!} = \frac{1}{3}$	$f(1,1,1) - f(0,1,1) = 1 - 0 = 1$	

$$= \frac{2}{3}$$

$\phi_2: S$	α_S	Δ_S	
$\{\}$	$\frac{1}{3}$	$f(0,1,0) - f(0,0,0) = 0 - 0 = 0$	$\phi_2 = \frac{1}{3} \cdot 0 + \frac{1}{6} \cdot 1 + \frac{1}{6} \cdot 0 + \frac{1}{3} \cdot 0$
$\{1\}$	$\frac{1}{6}$	$f(1,1,0) - f(1,0,0) = 1 - 0 = 1$	
$\{3\}$	$\frac{1}{6}$	$f(0,1,1) - f(0,0,1) = 0 - 0 = 0$	
$\{1,3\}$	$\frac{1}{3}$	$f(1,1,1) - f(1,0,1) = 1 - 1 = 0$	

$$= \frac{1}{6}$$

$$\phi_1 = \frac{2}{3} \rightarrow \phi_2 + \phi_3 = \frac{1}{3} \xrightarrow{\text{Symmetric}} \phi_2 = \phi_3 = \frac{1}{6}$$

$\phi_3: S$	α_S	Δ_S	
$\{\}$	$\frac{1}{3}$	$f(0,0,1) - f(0,0,0) = 0 - 0 = 0$	
$\{1\}$	$\frac{1}{6}$	$f(1,0,1) - f(1,0,0) = 1 - 0 = 1$	$\phi_3 = \frac{1}{6}$
$\{2\}$	$\frac{1}{6}$	$f(0,1,1) - f(0,1,0) = 0$	
$\{1,2\}$	$\frac{1}{3}$	$f(1,1,1) - f(1,1,0) = 1 - 1 = 0$	

Exercise 3: Taylor Expansions (25 P)

Consider the simple radial basis function

$$f(x) = \|x - \mu\| - \theta$$

with $\theta > 0$. For the purpose of extracting an explanation, we would like to build a first-order Taylor expansion of the function at some root point \tilde{x} . We choose this root point to be taken on the segment connecting μ and x (we assume that $f(x) > 0$ so that there is always a root point on this segment).

Show that the first-order terms of the Taylor expansion are given by

$$\phi_i = \frac{(x_i - \mu_i)^2}{\|x - \mu\|^2} \cdot (\|x - \mu\| - \theta)$$

$$f(x) = \|x - \mu\| - \theta = \left((x - \mu)^T (x - \mu) \right)^{\frac{1}{2}} - \theta$$

$$\nabla f(x) = \frac{1}{2} \cdot (x - \mu)^T \cdot (x - \mu)^{-\frac{1}{2}} \cdot 2(x - \mu)$$

$$= \frac{x - \mu}{\|x - \mu\|}$$

$$f(x) = \underbrace{f(\tilde{x})}_{=0} + \sum_{i=1}^d \underbrace{[\nabla f(\tilde{x})]_i \cdot (x_i - \tilde{x}_i)}_{\phi_i} + \dots$$

$$\text{In Taylor expansion we know } f(\tilde{x}) = \|\tilde{x} - \mu\| - \theta = 0$$

$$\Rightarrow \|\tilde{x} - \mu\| = \theta \quad \dots \quad (1)$$

we know that \tilde{x} lies on the sphere centered at μ with radius θ .

$$\Rightarrow \tilde{x} - \mu = \theta \cdot \frac{x - \mu}{\|x - \mu\|} \quad \dots \quad (2)$$

$$\phi_i = [\nabla f(\tilde{x})]_i \cdot (x_i - \tilde{x}_i) = \frac{\tilde{x}_i - \mu_i}{\|\tilde{x} - \mu\|} \cdot (x_i - \tilde{x}_i)$$

$$= \frac{\tilde{x}_i - \mu_i}{\theta} (x_i - \tilde{x}_i) \quad | \quad (1)$$

$$= \frac{\theta \cdot \frac{x_i - \mu_i}{\|x - \mu\|}}{\theta} (x_i - \tilde{x}_i) \quad | \quad (2)$$

$$= \frac{x_i - \mu_i}{\|x - \mu\|} (x_i - \tilde{x}_i)$$

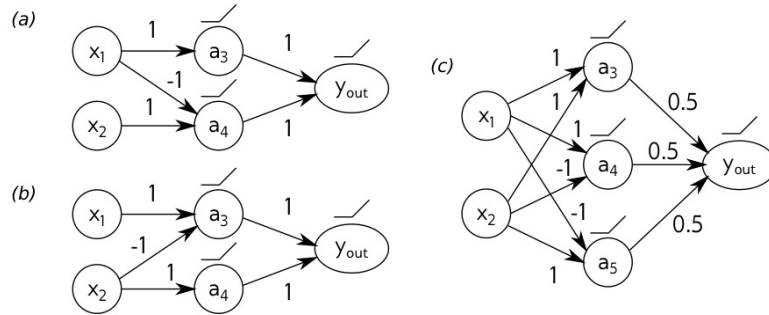
$$= \frac{x_i - \mu_i}{\|x - \mu\|} \cdot \left(x_i - \frac{x_i - \mu_i}{\|x - \mu\|} \cdot \theta - \mu_i \right) \quad | \quad (2)$$

$$= \frac{(x_i - \mu_i)^2}{\|x - \mu\|^2} \cdot \left(1 - \frac{\theta}{\|x - \mu\|} \right)$$

$$= \frac{(x_i - \mu_i)^2}{\|x - \mu\|^2} (\|x - \mu\| - \theta)$$

Exercise 4: Layer-Wise Relevance Propagation (25 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 $y = \max(x_1, x_2)$, where $x_1, x_2 \in \mathbb{R}^+$ are the input activations. This function can be implemented as a ReLU network in multiple ways. Three examples are given below.



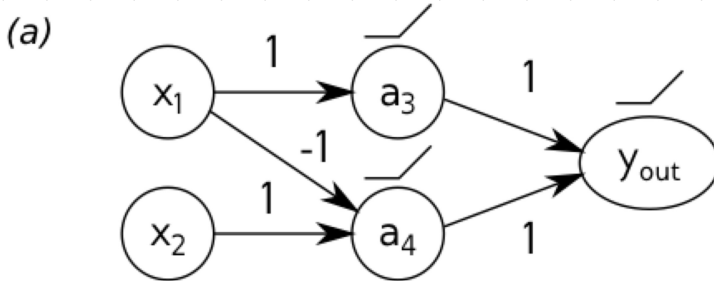
We consider the propagation rule:

$$R_j = \sum_k \frac{a_j w_{jk}^+}{\sum_j a_j w_{jk}^+} R_k$$

j is all nodes in left side

where j and k are indices for two consecutive layers and where $()^+$ denotes the positive part. This propagation rule is applied to both layers.

Give for each network the computational steps that lead to the scores R_1 and R_2 , and the obtained relevance values. More specifically, express R_1 and R_2 as a function of R_3 and R_4 (and R_5), and express the latter relevances as a function of $R_{out} = y$.



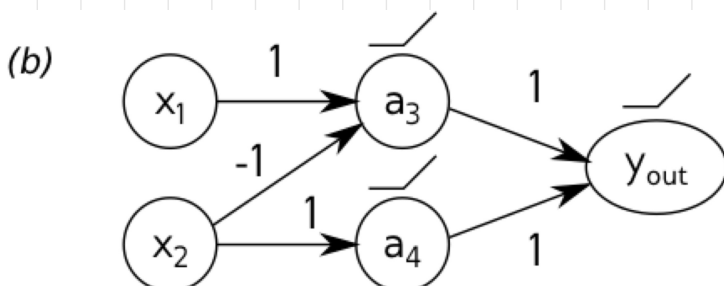
$$R_5 = R_{out} = y$$

$$R_4 = \frac{a_4 w_{45}}{\sum_{k=3}^4 a_k w_{k5}} \cdot R_5 = \frac{a_4}{a_3 + a_4} y$$

$$R_3 = \frac{a_3 w_{35}}{\sum_{k=3}^4 a_k w_{k5}} \cdot R_5 = \frac{a_3}{a_3 + a_4} y$$

$$R_2 = \sum_{k=3}^4 \frac{a_2 w_{2k}}{\sum_{j=1}^2 a_j w_{jk}} R_k = R_3 \cdot \frac{a_2 w_{23}}{a_1 w_{13} + a_2 w_{23}} + R_4 \cdot \frac{a_2 w_{24}}{a_1 w_{14} + a_2 w_{24}} = R_4$$

$$R_1 = \sum_{k=3}^4 \frac{a_1 w_{1k}}{\sum_{j=1}^2 a_j w_{jk}} R_k = R_3 \cdot \frac{a_1 w_{13}}{a_1 w_{13} + a_2 w_{23}} + R_4 \cdot \frac{a_1 w_{14}}{a_1 w_{14} + a_2 w_{24}} = R_3$$



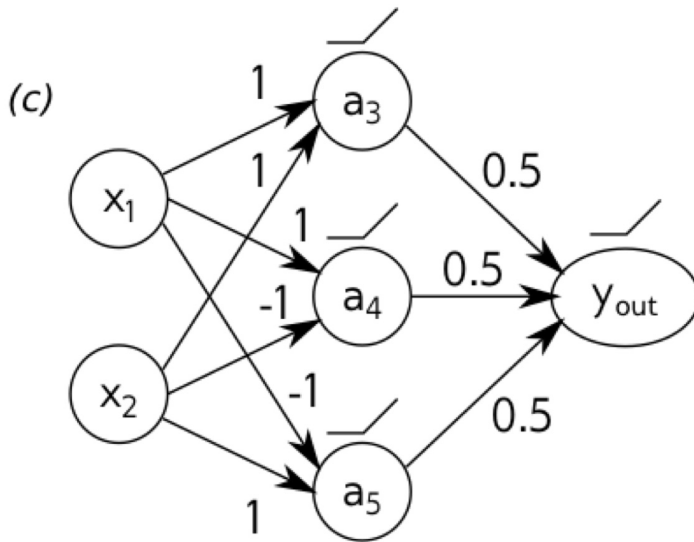
$$R_5 = R_{out} = Y$$

$$R_4 = \frac{a_4 W_{45}}{\sum_{k=3}^4 a_k W_{k5}} \cdot R_5 = \frac{a_4}{a_3 + a_4} Y$$

$$R_3 = \frac{a_3 V_3}{\sum_{k=3}^4 a_k V_k} \cdot R_5 = \frac{a_3}{a_3 + a_4} Y$$

$$R_2 = \sum_{k=3}^4 \frac{a_2 W_{2k}}{\sum_{j=1}^2 a_j W_{jk}} R_k = R_3 \cdot \frac{a_2 W_{23}^+}{a_1 W_{13}^+ + a_2 W_{23}^+} + R_4 \cdot \frac{a_2 W_{24}^+}{a_1 W_{14}^+ + a_2 W_{24}^+} = R_4$$

$$R_1 = \sum_{k=3}^4 \frac{a_1 W_{1k}}{\sum_{j=1}^2 a_j W_{jk}} R_k = R_3 \cdot \frac{a_1 W_{13}^+}{a_1 W_{13}^+ + a_2 W_{23}^+} + R_4 \cdot \frac{a_1 W_{14}^+}{a_1 W_{14}^+ + a_2 W_{24}^+} = R_3$$



$$R_6 = R_{out} = Y$$

$$R_5 = \frac{a_5 W_{56}}{\sum_{k=3}^5 a_k W_{k6}} \cdot R_6 = \frac{0.5 \cdot a_5}{0.5 \cdot (a_3 + a_4 + a_5)} Y = \frac{a_5}{a_3 + a_4 + a_5} Y$$

$$R_4 = \frac{a_4 W_{46}}{\sum_{k=3}^5 a_k W_{k6}} \cdot R_6 = \frac{0.5 \cdot a_4}{0.5 \cdot (a_3 + a_4 + a_5)} Y = \frac{a_4}{a_3 + a_4 + a_5} Y$$

$$R_3 = \frac{a_3 W_{36}}{\sum_{k=3}^5 a_k W_{k6}} \cdot R_6 = \frac{0.5 \cdot a_3}{0.5 \cdot (a_3 + a_4 + a_5)} Y = \frac{a_3}{a_3 + a_4 + a_5} Y$$

$$R_2 = \sum_{k=3}^5 \frac{a_2 W_{2k}}{\sum_{j=1}^2 a_j W_{jk}} R_k = R_3 \cdot \frac{a_2 W_{23}^+}{a_1 W_{13}^+ + a_2 W_{23}^+} + R_4 \cdot \frac{a_2 W_{24}^+}{a_1 W_{14}^+ + a_2 W_{24}^+} + R_5 \cdot \frac{a_2 W_{25}^+}{a_1 W_{15}^+ + a_2 W_{25}^+}$$

$$= R_3 \cdot \frac{x_1}{x_1 + x_2} + R_5$$

$$R_1 = \sum_{k=3}^5 \frac{a_1 W_{1k}}{\sum_{j=1}^2 a_j W_{jk}} R_k = R_3 \cdot \frac{a_1 W_{13}^+}{a_1 W_{13}^+ + a_2 W_{23}^+} + R_4 \cdot \frac{a_1 W_{14}^+}{a_1 W_{14}^+ + a_2 W_{24}^+} + R_5 \cdot \frac{a_1 W_{15}^+}{a_1 W_{15}^+ + a_2 W_{25}^+}$$

$$= R_3 \cdot \frac{x_1}{x_1 + x_2} + R_4$$