

Exercise 1: Class Prototypes (25 P)

Consider the linear model $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ 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

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

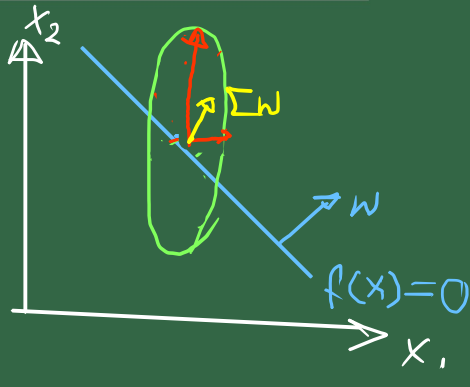
Find the prototype \mathbf{x}^* obtained by activation maximization subject to $\Omega(\mathbf{x}) = \log p(\mathbf{x})$ with $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ where $\boldsymbol{\mu}$ and Σ are the mean and covariance.

$$\frac{\partial}{\partial \mathbf{x}} \left(\mathbf{w}^T \mathbf{x} + b - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) + \text{const} \right)$$

$$= \mathbf{w} - \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \stackrel{!}{=} \mathbf{0}$$

$$\Leftrightarrow \mathbf{w} = \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \quad | \Sigma(\cdot)$$

$$\Leftrightarrow \Sigma \mathbf{w} = \mathbf{x} - \boldsymbol{\mu} \quad \Leftrightarrow \mathbf{x}^* = \Sigma \mathbf{w} + \boldsymbol{\mu}$$



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).

$\phi_i: S$	α_S	$f(x_{S \cup \{1\}}) - f(x_S)$	
$\{\}$	$1/3$	0	$= 0$
$\{2\}$	$1/6$	$1 - 0$	$= 1$
$\{3\}$	$1/6$	$1 - 0$	$= 1$
$\{2, 3\}$	$1/3$	$1 - 0$	$= 1$

$$\phi_1 = \frac{1}{3} \cdot 0 + \frac{1}{6} \cdot 1 + \frac{1}{6} \cdot 1 + \frac{1}{3} \cdot 1 = \frac{2}{3}$$

Conservation:

$$\phi_1 + \phi_2 + \phi_3 = \frac{1}{3}$$

Symmetry:

$$\phi_2 = \frac{1}{6} \quad \phi_3 = \frac{1}{6}$$

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

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

$$\tilde{\mathbf{x}} = \boldsymbol{\mu} + t(\mathbf{x} - \boldsymbol{\mu})$$

$$\mathbf{x} - \tilde{\mathbf{x}} = (1-t)(\mathbf{x} - \boldsymbol{\mu})$$

$$\phi_i = \frac{\tilde{x}_i - \mu_i}{\|\tilde{\mathbf{x}} - \boldsymbol{\mu}\|} \cdot (1-t)(x_i - \mu_i)$$

$$= \frac{\mu_i + t(x_i - \mu_i) - \mu_i}{\|\mu_i + t(x - \mu) - \mu\|} \cdot (1-t)(x_i - \mu_i) = \frac{t(x_i - \mu_i)}{\|t(\mathbf{x} - \boldsymbol{\mu})\|} \cdot (1-t)(x_i - \mu_i)$$

$$= \frac{(x_i - \mu_i)^2}{\|\mathbf{x} - \boldsymbol{\mu}\|^2} \cdot (1-t)$$

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

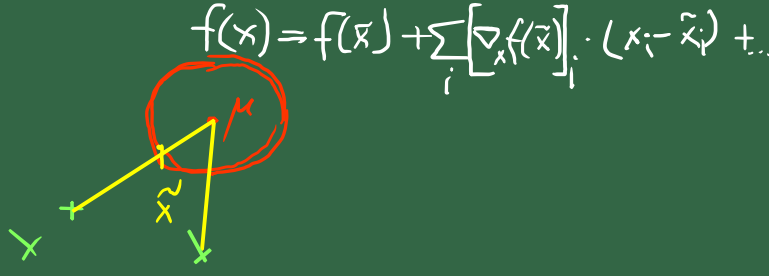
$$\|\mathbf{x} - \boldsymbol{\mu}\|^2 = \sum_i (x_i - \mu_i)^2$$

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

$$\sum_i \phi_i = 1$$

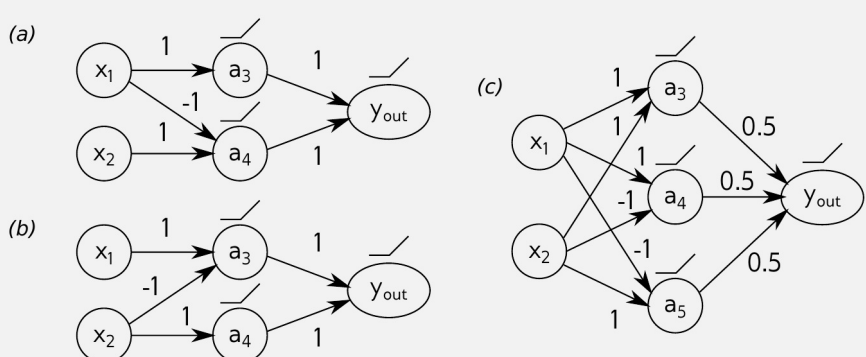
$$= \frac{(x_i - \mu_i)^2}{\|\mathbf{x} - \boldsymbol{\mu}\|^2} (\|\mathbf{x} - \boldsymbol{\mu}\| - \theta)$$

and: $f(\tilde{\mathbf{x}}) = 0$, solve for t .



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 $(\cdot)^+$ 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 $R_{\text{out}} = y$.

$$R_{\text{out}} = y$$

$$R_3 = \frac{a_3 \cdot \frac{1}{2}}{a_3 \cdot \frac{1}{2} + a_4 \cdot \frac{1}{2} + a_5 \cdot \frac{1}{2}} \cdot R_{\text{out}}$$

$$R_4 = \frac{a_4 \cdot \frac{1}{2}}{a_3 \cdot \frac{1}{2} + a_4 \cdot \frac{1}{2} + a_5 \cdot \frac{1}{2}} \cdot R_{\text{out}}$$

$$R_5 = \frac{a_5 \cdot \frac{1}{2}}{a_3 \cdot \frac{1}{2} + a_4 \cdot \frac{1}{2} + a_5 \cdot \frac{1}{2}} \cdot R_{\text{out}}$$

$$R_2 = \frac{x_2 \cdot 1}{x_2 \cdot 1 + x_1 \cdot 1} \cdot R_3 + \frac{x_2 \cdot 0}{x_2 \cdot 0 + x_1 \cdot 1} \cdot R_4$$

$$+ \frac{x_2 \cdot 1}{x_2 \cdot 1 + x_1 \cdot 0} \cdot R_5$$

$$R_1 = \frac{x_1 \cdot 1}{x_1 \cdot 1 + x_2 \cdot 1} \cdot R_3 + \frac{x_1 \cdot 1}{x_2 \cdot 0 + x_1 \cdot 1} \cdot R_4$$

$$+ \frac{x_1 \cdot 0}{x_2 \cdot 1 + x_1 \cdot 0} \cdot R_5$$

a)

$$R_{\text{out}} = y$$

$$R_3 = \frac{a_3 \cdot w_{3\text{out}}^+}{a_3 \cdot w_{3\text{out}}^+ + a_4 \cdot w_{4\text{out}}^+} \cdot R_{\text{out}}$$

$$= \frac{a_3 \cdot 1}{a_3 \cdot 1 + a_4 \cdot 1} \cdot R_{\text{out}}$$

$$R_4 = \frac{a_4 \cdot 1}{a_3 \cdot 1 + a_4 \cdot 1} \cdot R_{\text{out}}$$

$$R_2 = \frac{x_1 \cdot w_{24}^+}{x_2 \cdot w_{24}^+ + x_1 \cdot w_{14}^+} \cdot R_4$$

$$= R_4$$

$$R_1 = \frac{x_1 \cdot w_{13}^+}{x_1 \cdot w_{13}^+} \cdot R_3 + \frac{x_1 \cdot w_{14}^+}{x_1 \cdot w_{24}^+ + x_1 \cdot w_{14}^+} \cdot R_4$$

$$= R_3$$