

Exercise Sheet 12 (theory part)

Exercise 1: Activation Maximization (10 + 10 + 10 P)

Consider the linear model $f: \mathbb{R}^d \rightarrow \mathbb{R}$ with $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$. We would like to interpret the function f by building a prototype \mathbf{x}^* in the input domain which produces a large value of f . Activation maximization produces such prototype by computing

$$\mathbf{x}^* = \arg \max_{\mathbf{x}} [f(\mathbf{x}) - \Omega(\mathbf{x})].$$

where Ω is a regularization function.

(a) Find the prototype \mathbf{x}^* obtained by activation maximization subject to the penalty $\Omega(\mathbf{x}) = \lambda \|\mathbf{x}\|^2$ where λ is a hyperparameter.

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{w}^\top \mathbf{x} + b - \lambda \|\mathbf{x}\|^2) = \mathbf{w} - 2\lambda \mathbf{x} = 0 \quad \Rightarrow \quad \mathbf{x}^* = \frac{1}{2\lambda} \mathbf{w}$$

(b) 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 Σ are the mean and covariance.

$$\begin{aligned} \frac{\partial}{\partial \mathbf{x}} (\mathbf{w}^\top \mathbf{x} + b - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}) + \text{const.}) \\ = \mathbf{w} - \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}) = 0 \quad \Rightarrow \quad \mathbf{x}^* = \Sigma \mathbf{w} + \boldsymbol{\mu} \end{aligned}$$

(c) 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 $\Omega(\mathbf{z}) = \lambda \|\mathbf{z}\|^2$.

$$\begin{aligned} \frac{\partial}{\partial \mathbf{z}} (\mathbf{w}^\top (A\mathbf{z} + \mathbf{c}) + b - \lambda \|\mathbf{z}\|^2) = A^\top \mathbf{w} - 2\lambda \mathbf{z} = 0 \quad \Rightarrow \quad \mathbf{z}^* = \frac{A^\top}{2\lambda} \mathbf{w} \\ \Rightarrow \quad \mathbf{x}^* = \frac{AA^\top}{2\lambda} \mathbf{w} + \mathbf{c} \end{aligned}$$

Exercise 2: Attribution (10 + 10 P)

Consider the function $f: \mathbb{R}_+^3 \rightarrow \mathbb{R}_+$ with

$$f(\mathbf{x}) = \min(x_1, \max(x_2, x_3))$$

implementing some min-max pooling between three input features. For the data point $\mathbf{x} = (3, 2, 1)$ and its prediction $f(\mathbf{x}) = 2$, we would like to perform an attribution of the prediction to the input features. We investigate the Shapley value and Gradient \times Input methods for attribution.

(a) Recall that the Shapley value method identifies the contributions R_1, \dots, R_d of features x_1, \dots, x_d as:

$$R_i = \sum_{S: i \notin S} \frac{|S|!(d-|S|-1)!}{d!} [f(\mathbf{x}_{S \cup \{i\}}) - f(\mathbf{x}_S)]$$

where $(\mathbf{x}_S)_S$ are all possible subsets of features contained in the input \mathbf{x} . Compute the Shapley values associated to the prediction above. (We assume a reference point $\tilde{\mathbf{x}} = \mathbf{0}$, i.e. we set features to zero when removing them).

We first evaluate the function, and we get $f(3, 2, 1) = 2$.

Consider R_1 , and analyze all coalitions that don't contain that feature:

\mathcal{S}	$\alpha_{\mathcal{S}}$	$f(\mathbf{x}_{\mathcal{S} \cup \{1\}}) - f(\mathbf{x}_{\mathcal{S}})$
\emptyset	$1/3$	0
$\{2\}$	$1/6$	2
$\{3\}$	$1/6$	1
$\{2, 3\}$	$1/3$	2

Therefore $R_1 = 1/3 \cdot 0 + 1/6 \cdot 2 + 1/6 \cdot 1 + 1/3 \cdot 2 = 7/6$.

Consider now R_2 , and analyze all coalitions that don't contain that feature:

\mathcal{S}	$\alpha_{\mathcal{S}}$	$f(\mathbf{x}_{\mathcal{S} \cup \{2\}}) - f(\mathbf{x}_{\mathcal{S}})$
\emptyset	$1/3$	0
$\{1\}$	$1/6$	2
$\{3\}$	$1/6$	0
$\{1, 3\}$	$1/3$	1

Therefore $R_2 = 1/3 \cdot 0 + 1/6 \cdot 2 + 1/6 \cdot 1 + 1/3 \cdot 1 = 2/3$.

Consider now R_3 , and analyze all coalitions that don't contain that feature:

\mathcal{S}	$\alpha_{\mathcal{S}}$	$f(\mathbf{x}_{\mathcal{S} \cup \{3\}}) - f(\mathbf{x}_{\mathcal{S}})$
\emptyset	$1/3$	0
$\{1\}$	$1/6$	1
$\{2\}$	$1/6$	0
$\{1, 2\}$	$1/3$	0

Therefore $R_3 = 1/3 \cdot 0 + 1/6 \cdot 1 + 1/6 \cdot 1 + 1/3 \cdot 0 = 1/6$.

We can test our result by verifying the conservation property:

$$R_1 + R_2 + R_3 = 7/6 + 2/3 + 1/6 = 12/6 = f(\mathbf{x})$$

(b) The Gradient \times Input method attributes to the input features according to the formula:

$$R_i = x_i \cdot [\nabla f(\mathbf{x})]_i$$

Compute the Gradient \times Input attribution associated to the prediction above.

$$R_1 = 3 \cdot 0 = 0$$

$$R_2 = 2 \cdot 1 = 2$$

$$R_3 = 1 \cdot 0 = 0$$