

Exercise Sheet 9 - Bonus

Exercise 1: Analysis of a similarity models (0 P)

We consider here similarity models of type $y(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$ with the dot product on a feature map $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^h$ and satisfying first-order positive homogeneity i.e. $\forall \mathbf{x}, \forall t > 0 : \phi(t\mathbf{x}) = t\phi(\mathbf{x})$. In the following we focus on Linear/ReLU layers:

$$a_k = \left(\sum_j a_j w_{jk} \right)^+ \\ a_{k'} = \left(\sum_{j'} a_{j'} w_{j'k'} \right)^+,$$

with activations a_j and weights w_{jk} and $(\cdot)^+$ indicating the ReLU function. Further assume root points $(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}') = (\varepsilon \tilde{\mathbf{x}}, \varepsilon \tilde{\mathbf{x}}')$ with ε almost zero.

(a) Write down the Taylor expansion of function $y(\mathbf{x}, \mathbf{x}')$ up to second-order terms.

$$\begin{aligned} y(\mathbf{x}, \mathbf{x}') &= y(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}') \\ &+ \sum_i [\nabla y(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}')]_i (x_i - \tilde{x}_i) \\ &+ \sum_{i'} [\nabla y(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}')]_{i'} (x'_{i'} - \tilde{x}'_{i'}) \\ &+ \sum_{ii'} [\nabla^2 y(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}')]_{ii'} (x_i - \tilde{x}_i) (x'_{i'} - \tilde{x}'_{i'}) \end{aligned}$$

(b) Analyse zero-order terms. Why do they vanish?

With $\phi(t\mathbf{x}) = t\phi(\mathbf{x})$, and $(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}') = (\varepsilon \tilde{\mathbf{x}}, \varepsilon \tilde{\mathbf{x}}')$ with ε almost zero we see that zero-order terms vanish and $y(\tilde{\mathbf{x}}, \tilde{\mathbf{x}}') = 0$.

Now, assume the following propagation rule for the Linear/ReLU layer to identify relevant interaction between a pair of neurons j and j' :

$$\begin{aligned} R_{jj'} &= \sum_{kk'} R_{jj' \leftarrow kk'} \\ &= \sum_{kk'} \frac{a_j a_{j'} \rho(w_{jk}) \rho(w_{j'k'})}{\sum_{jj'} a_j a_{j'} \rho(w_{jk}) \rho(w_{j'k'})} R_{kk'} \end{aligned}$$

- (i) If neurons (j, j') jointly activate, i.e. $a_j a_{j'}$ is non-zero.
- (ii) If pairs of neurons in the layer above jointly react, i.e. $R_{kk'}$ is non-zero (or relevant).
- (iii) If these reacting pairs are themselves relevant, i.e. the term $a_j a_{j'} \rho(w_{jk}) \rho(w_{j'k'})$ is non-zero.

(c) Show that $R_{jj'}$ factorizes as $R_{jj'} = \sum_{m=1}^h R_{jm} R_{j'm}$. Use the factorization of the subsequent layer $R_{kk'} = \sum_{m=1}^h R_{km} \cdot R_{k'm}$.

$$\begin{aligned} R_{jj'} &= \sum_{kk'} \frac{a_j a_{j'} \rho(w_{jk}) \rho(w_{j'k'})}{\sum_{jj'} a_j a_{j'} \rho(w_{jk}) \rho(w_{j'k'})} \sum_{m=1}^h R_{km} R_{k'm} \\ &= \sum_{m=1}^h \sum_{kk'} \frac{a_j \rho(w_{jk}) a_{j'} \rho(w_{j'k'})}{\sum_j a_j \rho(w_{jk}) \sum_{j'} a_{j'} \rho(w_{j'k'})} R_{km} R_{k'm} \\ &= \sum_{m=1}^h \underbrace{\left(\sum_k \frac{a_j \rho(w_{jk})}{\sum_j a_j \rho(w_{jk})} R_{km} \right)}_{R_{jm}} \cdot \underbrace{\left(\sum_{k'} \frac{a_{j'} \rho(w_{j'k'})}{\sum_{j'} a_{j'} \rho(w_{j'k'})} R_{k'm} \right)}_{R_{j'm}} \end{aligned}$$