

Mathematical Foundation of Machine Learning Spring 2024: Assignment II

1. (40%) [Besov smoothness of indicator functions] Let $f(x) = 1_{\tilde{\Omega}}(x)$, where $\tilde{\Omega} \subset [0,1]^2$ is a domain with a smooth boundary.
- Prove that for an isotropic tree \mathcal{T}_I , that applies non-adaptive subdivisions at dyadic partitions of the two variables, we get that $f \in B_\tau^\alpha(\mathcal{T}_I)$, for $\alpha < 1/(2\tau)$.
 - Assume that one can construct an anisotropic tree \mathcal{T}_A , such that from some level $k_0 \geq 1$, there are only $\leq c_1 2^m$ domains at the level $k_0 + 2m$ that intersect the boundary of $\tilde{\Omega}$, each of area $\leq c_2 2^{-3m}$. Prove that $f \in B_\tau^\alpha(\mathcal{T}_A)$, where $\alpha < 2/(3\tau)$. You can assume that there exists $0 < \rho < 1$, such that for any child Ω' of Ω , $|\Omega'| \leq \rho|\Omega|$.

Hints:

- For both cases you need to follow the proof of the case of the classic Besov smoothness with dyadic cubes. The case (a) is very similar to the classic case, because $B_\tau^{2\alpha} \sim B_\tau^\alpha(\mathcal{T}_I)$.
 - In both cases, you could first bound the contribution of the odd levels of the trees by the even levels. Then, it will be sufficient to only estimate the sum of contributions of even levels.
 - Use the ρ volume condition for (b) and the properties of the modulus. It is useful to observe that the condition $|\Omega'| \leq \rho|\Omega|$, also gives $|\Omega'| \geq (1-\rho)|\Omega|$.
2. (30%) [Convolutions] The convolution of $f, g \in L_1(\mathbb{R}^n)$ is defined by $f * g(x) := \int_{\mathbb{R}^n} f(x-y)g(y)dy$.
- Prove that $f * g \in L_1(\mathbb{R}^n)$,
 - Prove that $f * g = g * f$,
 - The Fourier Transform is defined by $\hat{f}(w) = \int_{\mathbb{R}^n} f(x)e^{-i\langle w, x \rangle} dx$, for $w \in \mathbb{R}^n$. Show that

$$(f * g)^\wedge(w) = \hat{f}(w)\hat{g}(w), \quad \forall w \in \mathbb{R}^n.$$
 - Let $f \in L_1(\mathbb{R}^2)$ be a piecewise constant function. Design a ‘filter’ $g \in L_\infty(\mathbb{R}^2)$, with support in $[-\varepsilon/2, \varepsilon/2]^2$, for some $\varepsilon > 0$, such that $f * g$ is ‘significant’ only in ε neighborhoods of points where f has ‘almost’ vertical edges.

3. (30%) [Back propagation] Let $f(x) = h_\theta(f_2(f_1(x, w_1), w_2))$, $w := (w_1, w_2) \in \mathbb{R}^2$, $\theta = (\beta, \beta_0) \in \mathbb{R}^2$, be a binary classification model consisting of 2 neural network layers and a logistic regression function

$$h_\theta(t) = \frac{1}{1 + e^{-(\beta t + \beta_0)}}.$$

Compute the gradient with respect to the negative log-likelihood loss, for the 4 weights at some point $(w^*, \theta^*) \in \mathbb{R}^4$, using the training dataset $\{x_i, y_i\}_{i \in I}$, $x_i \in \mathbb{R}$, $y_i \in \{0, 1\}$.