
Differential Geometric Regularization for Supervised Learning of Classifiers

Supplemental Materials

A. Proof of Theorem 1

Proof. For $\mathbf{f} : \mathbb{R}^N \rightarrow \Delta^{L-1} \subset \mathbb{R}^L$,

$$\{r_j = r_j(\mathbf{x}) = (0, \dots, \overset{j}{1}, \dots, 0, f_j^1, \dots, f_j^L) : j = 1, \dots, N\}$$

is a basis of the tangent space $T_{\mathbf{x}\text{gr}(\mathbf{f})}$ to $\text{gr}(\mathbf{f})$. Here $f_j^i = \partial_{x^j} f^i$. Let $\{e_i\}$ be an orthonormal frame of $T_{\mathbf{x}\text{gr}(\mathbf{f})}$. We have

$$e_i = B_i^j r_j$$

for some invertible matrix B_i^j .

Define the metric matrix g for the basis $\{r_j\}$ by

$$g = (g_{kj}) \text{ with } g_{kj} = r_k \cdot r_j = \delta_{kj} + f_k^i f_j^i.$$

Then

$$\begin{aligned} \delta_{ij} &= e_i \cdot e_j = B_i^k B_j^t r_k \cdot r_t = B_i^k B_j^t g_{kt} \\ &\Rightarrow I = (BB^T)g \Rightarrow BB^T = g^{-1}. \end{aligned}$$

Thus BB^T is computable in terms of derivatives of \mathbf{f} .

Let $D_u w$ be the \mathbb{R}^{N+L} directional derivative of w in the direction u . Then

$$\begin{aligned} \text{Tr II} &= P^\nu D_{e_i} e_i = P^\nu D_{B_i^j r_j} B_i^k r_k = B_i^j P^\nu D_{r_j} B_i^k r_k \\ &= B_i^j P^\nu [(D_{r_j} B_i^k) r_k] + B_i^j B_i^k D_{r_j} r_k \\ &= B_i^j B_i^k P^\nu D_{r_j} r_k \\ &= (g^{-1})^{jk} P^\nu D_{r_j} r_k, \end{aligned}$$

since $P^\nu r_k = 0$.

We have

$$\begin{aligned} r_k &= (0, \dots, 1, \dots, f_k^1(x^1, \dots, x^N), \dots, f_k^L(x^1, \dots, x^N)) \\ &= \partial_k^{\mathbb{R}^{N+L}} + \sum_{i=1}^L f_k^i \partial_{N+i}^{\mathbb{R}^{N+L}}, \end{aligned}$$

so in particular, $\partial_\ell^{\mathbb{R}^{N+L}} r_k = 0$ if $\ell > N$. Thus □

$$D_{r_j} r_k = (0, \dots, 0, f_{kj}^1, \dots, f_{kj}^L).$$

So far, we have

$$\text{Tr II} = (g^{-1})^{jk} P^\nu (0, \dots, 0, f_{kj}^1, \dots, f_{kj}^L).$$

Since g is given in terms of derivatives of \mathbf{f} , we need to write $P^\nu = I - P^T$ in terms of derivatives of \mathbf{f} . For any $u \in \mathbb{R}^{N+L}$, we have

$$\begin{aligned} P^T u &= (P^T u \cdot e_i) e_i = (u \cdot B_i^j r_j) B_i^k r_k \\ &= B_i^j B_i^k (u \cdot r_j) r_k \\ &= (g^{-1})^{jk} (u \cdot r_j) r_k. \end{aligned}$$

Thus

$$\text{Tr II} \tag{1}$$

$$= (g^{-1})^{jk} P^\nu (0, \dots, 0, f_{kj}^1, \dots, f_{kj}^L) \tag{2}$$

$$= (g^{-1})^{jk} (0, \dots, 0, f_{kj}^1, \dots, f_{kj}^L)$$

$$- P^T [(g^{-1})^{jk} (0, \dots, 0, f_{kj}^1, \dots, f_{kj}^L)]$$

$$= (g^{-1})^{jk} (0, \dots, 0, f_{kj}^1, \dots, f_{kj}^L)$$

$$- (g^{-1})^{jk} [(g^{-1})^{rs} (0, \dots, 0, f_{rs}^1, \dots, f_{rs}^L) \cdot r_j] r_k$$

$$= (g^{-1})^{jk} (0, \dots, 0, f_{kj}^1, \dots, f_{kj}^L)$$

$$- (g^{-1})^{jk} (g^{-1})^{rs} (f_{rs}^i f_j^i) r_k$$

$$= (g^{-1})^{ij} \left(0, \dots, - (g^{-1})^{rs} f_{rs}^a f_i^a, \dots, 0, \tag{3}$$

$$f_{ji}^1 - (g^{-1})^{rs} f_{rs}^a f_i^a f_j^1, \dots, f_{ji}^L - (g^{-1})^{rs} f_{rs}^a f_i^a f_j^L \right),$$

after a relabeling of indices. Therefore, the last L component of Tr II are given by

$$\begin{aligned} \text{Tr II}^L &= (g^{-1})^{ij} \left(f_{ji}^1 - (g^{-1})^{rs} f_{rs}^a f_i^a f_j^1, \dots, \right. \\ &\quad \left. f_{ji}^L - (g^{-1})^{rs} f_{rs}^a f_i^a f_j^L \right). \end{aligned}$$

B. An Easy Example with Bayes Consistency

We now give an example with a loss function that enables easy Bayes consistency proof under some mild initialization assumption. Related notation is summarized in §C.

For ease of reading, we change the notation for empirical penalty $\mathcal{P}_{\mathcal{T}_m}$ in this supplemental material to \mathcal{P}_D , i.e., $\mathcal{P} = \mathcal{P}_D + \lambda \mathcal{P}_G$. \mathcal{P}_D measures the deviation of $\text{gr}(\mathbf{f})$ from the mapped training points, a natural geometric distance penalty term is an L^2 distance in \mathbb{R}^L from $\mathbf{f}(\mathbf{x})$ to the averaged \mathbf{z} component of the k -nearest training points:

$$\mathcal{P}_D(\mathbf{f}) = R_{D, \mathcal{T}_m, k}(\mathbf{f}) = \int_{\mathcal{X}} d^2 \left(\mathbf{f}(\mathbf{x}), \frac{1}{k} \sum_{i=1}^k \tilde{\mathbf{z}}_i \right) d\mathbf{x}, \quad (4)$$

where d is the Euclidean distance in \mathbb{R}^L , $\tilde{\mathbf{z}}_i$ is the vector of the last L components of $(\tilde{\mathbf{x}}_i, \tilde{\mathbf{z}}_i) = (\tilde{\mathbf{x}}_i^1, \dots, \tilde{\mathbf{x}}_i^N, \tilde{\mathbf{z}}_i^1, \dots, \tilde{\mathbf{z}}_i^L)$, with $\tilde{\mathbf{x}}_i$ the i^{th} nearest neighbor of \mathbf{x} in \mathcal{T}_m , and $d\mathbf{x}$ is the Lebesgue measure. The gradient vector field is

$$\nabla(R_{D, \mathcal{T}_m, k})_{\mathbf{f}}(\mathbf{x}, \mathbf{f}(\mathbf{x})) = \frac{2}{k} \sum_{i=1}^k (\mathbf{f}(\mathbf{x}) - \tilde{\mathbf{z}}_i).$$

However, $\nabla(R_{D, \mathcal{T}_m, k})_{\mathbf{f}}$ is discontinuous on the set \mathcal{D} of points \mathbf{x} such that \mathbf{x} has equidistant training points among its k nearest neighbors. \mathcal{D} is the union of $(N-1)$ -dimensional hyperplanes in \mathcal{X} , so \mathcal{D} has measure zero. Such points will necessarily exist unless the last L components of the mapped training points are all 1 or all 0. To rectify this, we can smooth out $\nabla(R_{D, \mathcal{T}_m, k})_{\mathbf{f}}$ to a vector field

$$V_{D, \mathbf{f}, \phi} = \frac{2\phi(\mathbf{x})}{k} \sum_{i=1}^k (\mathbf{f}(\mathbf{x}) - \tilde{\mathbf{z}}_i). \quad (5)$$

Here $\phi(\mathbf{x})$ is a smooth damping function close to the singular function $\delta_{\mathcal{D}}$, which has $\delta_{\mathcal{D}}(\mathbf{x}) = 0$ for $\mathbf{x} \in \mathcal{D}$ and $\delta_{\mathcal{D}}(\mathbf{x}) = 1$ for $\mathbf{x} \notin \mathcal{D}$. Outside any open neighborhood of \mathcal{D} , $\nabla R_{D, \mathcal{T}_m, k} = V_{D, \mathbf{f}, \phi}$ for ϕ close enough to $\delta_{\mathcal{D}}$.

Recall the geometric penalty from the submission, i.e., $\mathcal{P}_G(\mathbf{f}) = \int_{\text{gr}(\mathbf{f})} \text{dvol}$, with the geometric gradient vector field being $V_{G, \mathbf{f}} = -\text{Tr } \Pi^L$.

Then the gradient vector field $V_{\text{tot}, \lambda, m, \mathbf{f}, \phi}$ of this example penalty \mathcal{P} is,

$$\begin{aligned} V_{\text{tot}, \lambda, m, \mathbf{f}, \phi} &= \nabla \mathcal{P}_{\mathbf{f}} = V_{D, \mathbf{f}, \phi} + \lambda V_{G, \mathbf{f}} \\ &= \frac{2\phi(\mathbf{x})}{k} \sum_{i=1}^k (\mathbf{f}(\mathbf{x}) - \tilde{\mathbf{z}}_i) - \lambda \text{Tr } \Pi^L \end{aligned} \quad (6)$$

B.1. Consistency analysis

For a training set \mathcal{T}_m , we let $\mathbf{f}_{\mathcal{T}_m} = (f_{\mathcal{T}_m}^1, \dots, f_{\mathcal{T}_m}^L)$ be the class probability estimator given by our approach. We denote the generalization risk of the corresponding plug-in classifier $h_{\mathbf{f}_{\mathcal{T}_m}}$ by $R_P(\mathbf{f}_{\mathcal{T}_m}) = \mathbb{E}_P[\mathbb{1}_{h_{\mathbf{f}_{\mathcal{T}_m}}(\mathbf{x}) \neq y}]$. The Bayes risk is defined by $R_P^* = \inf_{h: \mathcal{X} \rightarrow \mathcal{Y}} R_P(h) = \mathbb{E}_P[\mathbb{1}_{h_{\eta}(\mathbf{x}) \neq y}]$. Our algorithm is Bayes consistent if

$\lim_{m \rightarrow \infty} R_P(\mathbf{f}_{\mathcal{T}_m}) = R_P^*$ holds in probability for all distributions P on $\mathcal{X} \times \mathcal{Y}$. Usually, gradient flow methods are applied to a convex functional, so that a flow line approaches the unique global minimum. If the domain of the functional is an infinite dimensional manifold of (e.g. smooth) functions, we always assume that flow lines exist and that the actual minimum exists in this manifold.

Because our functionals are not convex, we can only hope to prove Bayes consistency for the set of initial estimators in the stable manifold of a stable fixed point (or sink) of the vector field (Guckenheimer & Worfolk, 1993). Recall that a stable fixed point \mathbf{f}_0 has a maximal open neighborhood, the stable manifold $\mathcal{S}_{\mathbf{f}_0}$, on which flow lines tend towards \mathbf{f}_0 . For the manifold \mathcal{M} , the stable manifold for a stable critical point of the vector field $V_{\text{tot}, \lambda, m, \mathbf{f}, \phi}$ is infinite dimensional.

The proof of Bayes consistency for multiclass (including binary) classification follows these steps:

Step 1: $\lim_{\lambda \rightarrow 0} R_{D, P, \lambda}^* = 0$.

Step 2: $\lim_{n \rightarrow \infty} R_{D, P}(\mathbf{f}_n) = 0 \Rightarrow \lim_{n \rightarrow \infty} R_P(\mathbf{f}_n) = R_P^*$.

Step 3: For all $\mathbf{f} \in \mathcal{M} = \text{Maps}(\mathcal{X}, \Delta^{L-1})$, $|R_{D, \mathcal{T}_m}(\mathbf{f}) - R_{D, P}(\mathbf{f})| \xrightarrow{m \rightarrow \infty} 0$ in probability.

Proofs of these steps are provided in following subsections. For the notation see §C. $R_{D, P, \lambda}^*$ is the minimum of the regularized D risk $\hat{R}_{D, P, \lambda}(\mathbf{f})$ for \mathbf{f} : $R_{D, P, \lambda}(\mathbf{f}) = R_{D, P}(\mathbf{f}) + \lambda \mathcal{P}_G(\mathbf{f})$, with $R_{D, P}(\mathbf{f}) = \int_{\mathcal{X}} d^2(\mathbf{f}(\mathbf{x}), \boldsymbol{\eta}(\mathbf{x})) d\mathbf{x}$ the D -risk. Also, $R_{D, \mathcal{T}_m, \lambda}(\mathbf{f}) = R_{D, \mathcal{T}_m}(\mathbf{f}) + \lambda \mathcal{P}_G(\mathbf{f})$, with $R_{D, \mathcal{T}_m}(\mathbf{f}) = \int_{\mathcal{X}} d^2 \left(\mathbf{f}(\mathbf{x}), \frac{1}{k} \sum_{i=1}^k \tilde{\mathbf{z}}_i \right) d\mathbf{x}$ the empirical D -risk.

Theorem 2 (Bayes Consistency). *Let m be the size of the training data set. Let $\mathbf{f}_{1, \lambda, m} \in \mathcal{S}_{\mathbf{f}_{D, \mathcal{T}_m, \lambda}}$, the stable manifold for the global minimum $\mathbf{f}_{D, \mathcal{T}_m, \lambda}$ of $R_{D, \mathcal{T}_m, \lambda}$, and let $\mathbf{f}_{n, \lambda, m, \phi}$ be a sequence of functions on the flow line of $V_{\text{tot}, \lambda, m, \mathbf{f}, \phi}$ starting with $\mathbf{f}_{1, \lambda, m}$ with the flow time $t_n \rightarrow \infty$ as $n \rightarrow \infty$. Then $R_P(\mathbf{f}_{n, \lambda, m, \phi}) \xrightarrow[\lambda \rightarrow 0, \phi \rightarrow \delta_{\mathcal{D}}]{m, n \rightarrow \infty} R_P^*$ in probability for all distributions P on $\mathcal{X} \times \mathcal{Y}$, if $k/m \rightarrow 0$ as $m \rightarrow \infty$.*

Proof. In the notation of §C, if $\mathbf{f}_{D, \mathcal{T}_m, \lambda}$ is a global minimum for $R_{D, \mathcal{T}_m, \lambda}$, then outside of \mathcal{D} , $\mathbf{f}_{D, \mathcal{T}_m, \lambda}$ is the limit of critical points for the negative flow of $V_{\text{tot}, \lambda, m, \mathbf{f}, \phi}$ as $\phi \rightarrow \delta_{\mathcal{D}}$. To see this, fix an ϵ_i neighborhood \mathcal{D}_{ϵ_i} of \mathcal{D} . For a sequence $\phi_j \rightarrow \delta_{\mathcal{D}}$, $V_{\text{tot}, \lambda, m, \mathbf{f}, \phi_j}$ is independent of $j \geq j(\epsilon_i)$ on $\mathcal{X} \setminus \mathcal{D}_{\epsilon_i}$, so we find a function \mathbf{f}_i , a critical point of $V_{\text{tot}, \lambda, m, \mathbf{f}, \phi_j(\epsilon_i)}$, equal to $\mathbf{f}_{D, \mathcal{T}_m, \lambda}$ on $\mathcal{X} \setminus \mathcal{D}_{\epsilon_i}$. Since any $\mathbf{x} \notin \mathcal{D}$ lies outside some \mathcal{D}_{ϵ_i} , the sequence \mathbf{f}_i

converges at \mathbf{x} if we let $\epsilon_i \rightarrow 0$. Thus we can ignore the choice of ϕ in our proof, and drop ϕ from the notation.

For our algorithm, for fixed λ, m , we have as above

$$\lim_{n \rightarrow \infty} \mathbf{f}_{n,\lambda,m} = \mathbf{f}_{D,\mathcal{T}_m,\lambda}, \text{ so}$$

$$\lim_{n \rightarrow \infty} R_{D,\mathcal{T}_m,\lambda}(\mathbf{f}_{n,\lambda,m}) = R_{D,\mathcal{T}_m,\lambda}(\mathbf{f}_{D,\mathcal{T}_m,\lambda}),$$

for $\mathbf{f}_1 \in \mathcal{S}_{\mathbf{f}_{D,\mathcal{T}_m,\lambda}}$. By Step 2, it suffices to show $R_{D,P}(\mathbf{f}_{D,\mathcal{T}_m,\lambda}) \xrightarrow[\lambda \rightarrow 0]{m \rightarrow \infty} 0$. In probability, we have $\forall \delta > 0, \exists m > 0$ such that

$$\begin{aligned} 0 &\leq R_{D,P}(\mathbf{f}_{D,\mathcal{T}_m,\lambda}) \\ &\leq R_{D,P}(\mathbf{f}_{D,\mathcal{T}_m,\lambda}) + \lambda \mathcal{P}_G(\mathbf{f}_{D,\mathcal{T}_m,\lambda}) \\ &\leq R_{D,\mathcal{T}_m}(\mathbf{f}_{D,\mathcal{T}_m,\lambda}) + \lambda \mathcal{P}_G(\mathbf{f}_{D,\mathcal{T}_m,\lambda}) + \frac{\delta}{3} \quad (\text{Step 3}) \\ &= R_{D,\mathcal{T}_m,\lambda}(\mathbf{f}_{D,\mathcal{T}_m,\lambda}) + \frac{\delta}{3} \\ &\leq R_{D,\mathcal{T}_m,\lambda}(\mathbf{f}_{D,P,\lambda}) + \frac{\delta}{3} \quad (\text{minimality of } \mathbf{f}_{D,\mathcal{T}_m,\lambda}) \\ &= R_{D,\mathcal{T}_m}(\mathbf{f}_{D,P,\lambda}) + \lambda \mathcal{P}_G(\mathbf{f}_{D,P,\lambda}) + \frac{\delta}{3} \\ &\leq R_{D,P}(\mathbf{f}_{D,P,\lambda}) + \lambda \mathcal{P}_G(\mathbf{f}_{D,P,\lambda}) + \frac{2\delta}{3} \quad (\text{Step 3}) \\ &= R_{D,P,\lambda}(\mathbf{f}_{D,P,\lambda}) + \frac{2\delta}{3} = R_{D,P,\lambda}^* + \frac{2\delta}{3} \\ &\leq \delta, \quad (\text{Step 1}) \end{aligned}$$

for λ close to zero. \square

B.2. Step 1

Lemma 1. (Step 1) $\lim_{\lambda \rightarrow 0} R_{D,P,\lambda}^* = 0$.

Proof. After the smoothing procedure in §3.1 for the distance penalty term, the function $R_{D,P,\lambda} : \mathcal{M} \rightarrow \mathbb{R}$ is continuous in the Fréchet topology on \mathcal{M} . We check that the functions $R_{D,P,\lambda} : \mathcal{M} \rightarrow \mathbb{R}$ are equicontinuous in λ : for fixed $\mathbf{f}_0 \in \mathcal{M}$ and $\epsilon > 0$, there exists $\delta = \delta(\mathbf{f}_0, \epsilon)$ such that $|\lambda - \lambda'| < \delta \Rightarrow |R_{D,P,\lambda}(\mathbf{f}_0) - R_{D,P,\lambda'}(\mathbf{f}_0)| < \epsilon$. This is immediate:

$$|R_{D,P,\lambda}(\mathbf{f}_0) - R_{D,P,\lambda'}(\mathbf{f}_0)| = |(\lambda - \lambda') \mathcal{P}_G(\mathbf{f}_0)| < \epsilon,$$

if $\delta < \epsilon / |\mathcal{P}_G(\mathbf{f}_0)|$. It is standard that the infimum $\inf R_\lambda$ of an equicontinuous family of functions is continuous in λ , so $\lim_{\lambda \rightarrow 0} R_{D,P,\lambda}^* = R_{D,P,\lambda=0}^* = R_{D,P}(\boldsymbol{\eta}) = 0$. \square

B.3. Step 2

We assume that the class probability function $\boldsymbol{\eta}(\mathbf{x}) : \mathbb{R}^N \rightarrow \mathbb{R}^L$ is smooth, and that the marginal distribution $\mu(\mathbf{x})$ is continuous. We also let μ denote the corresponding measure on \mathcal{X} .

Denote: $h_{\mathbf{f}}(\mathbf{x}) = \operatorname{argmax}\{f^\ell(\mathbf{x}), \ell \in \mathcal{Y}\}$, and,

$$\mathbb{1}_{h_{\mathbf{f}}(\mathbf{x}) \neq y} = \begin{cases} 1, & h_{\mathbf{f}}(\mathbf{x}) \neq y, \\ 0, & h_{\mathbf{f}}(\mathbf{x}) = y. \end{cases}$$

Lemma 2. (Step 2 for a subsequence)

$$\lim_{n \rightarrow \infty} R_{D,P}(\mathbf{f}_n) = 0 \Rightarrow \lim_{i \rightarrow \infty} R_P(\mathbf{f}_{n_i}) = R_P^*$$

for some subsequence $\{\mathbf{f}_{n_i}\}_{i=1}^\infty$ of $\{\mathbf{f}_n\}$.

Proof. The left hand side of the Lemma is

$$\int_{\mathcal{X}} d^2(\mathbf{f}_n(\mathbf{x}), \boldsymbol{\eta}(\mathbf{x})) d\mathbf{x} \rightarrow 0,$$

which is equivalent to

$$\int_{\mathcal{X}} d^2(\mathbf{f}_n(\mathbf{x}), \boldsymbol{\eta}(\mathbf{x})) \mu(\mathbf{x}) d\mathbf{x} \rightarrow 0, \quad (7)$$

since \mathcal{X} is compact and μ is continuous. Therefore, it suffices to show

$$\begin{aligned} \int_{\mathcal{X}} d^2(\mathbf{f}_n(\mathbf{x}), \boldsymbol{\eta}(\mathbf{x})) \mu(\mathbf{x}) d\mathbf{x} &\rightarrow 0 \\ \implies \mathbb{E}_P[\mathbb{1}_{h_{\mathbf{f}_n}(\mathbf{x}) \neq y}] &\rightarrow \mathbb{E}_P[\mathbb{1}_{h_{\boldsymbol{\eta}}(\mathbf{x}) \neq y}]. \end{aligned} \quad (8)$$

We recall that L^2 convergence implies pointwise convergence a.e, so (7) implies that a subsequence of \mathbf{f}_n , also denoted \mathbf{f}_n , has $\mathbf{f}_n \rightarrow \boldsymbol{\eta}(\mathbf{x})$ pointwise a.e. on \mathcal{X} . (By our assumption on $\mu(\mathbf{x})$, these statements hold for either μ or Lebesgue measure.) By Egorov's theorem, for any $\epsilon > 0$, there exists a set $B_\epsilon \subset \mathcal{X}$ with $\mu(B_\epsilon) < \epsilon$ such that $\mathbf{f}_n \rightarrow \boldsymbol{\eta}(\mathbf{x})$ uniformly on $\mathcal{X} \setminus B_\epsilon$.

Fix $\delta > 0$ and set

$$Z_\delta = \{\mathbf{x} \in \mathcal{X} : \#\{\operatorname{argmax}_{\ell \in \mathcal{Y}} \eta^\ell(\mathbf{x})\} = 1,$$

$$|\max_{\ell \in \mathcal{Y}} \eta^\ell(\mathbf{x}) - \operatorname{submax}_{\ell \in \mathcal{Y}} \eta^\ell(\mathbf{x})| < \delta\},$$

where submax denotes the second largest element in $\{\eta^1(\mathbf{x}), \dots, \eta^L(\mathbf{x})\}$. For the moment, assume that $Z_0 = \{\mathbf{x} \in \mathcal{X} : \#\{\operatorname{argmax}_{\ell \in \mathcal{Y}} \eta^\ell(\mathbf{x})\} > 1\}$ has $\mu(Z_0) = 0$.

It follows easily¹ that $\mu(Z_\delta) \rightarrow 0$ as $\delta \rightarrow 0$. On $\mathcal{X} \setminus (Z_\delta \cup B_\epsilon)$, we have $\mathbb{1}_{h_{\mathbf{f}_n}(\mathbf{x}) \neq y} = \mathbb{1}_{h_{\boldsymbol{\eta}}(\mathbf{x}) \neq y}$ for $n > N_\delta$. Thus

$$\mathbb{E}_P[\mathbb{1}_{\mathcal{X} \setminus (Z_\delta \cup B_\epsilon)} \mathbb{1}_{h_{\mathbf{f}_n}(\mathbf{x}) \neq y}] = \mathbb{E}_P[\mathbb{1}_{\mathcal{X} \setminus (Z_\delta \cup B_\epsilon)} \mathbb{1}_{h_{\boldsymbol{\eta}}(\mathbf{x}) \neq y}].$$

¹Let A_k be sets with $A_{k+1} \subset A_k$ and with $\mu(\cap_{k=1}^\infty A_k) = 0$. If $\mu(A_k) \not\rightarrow 0$, then there exists a subsequence, also called A_k , with $\mu(A_k) > K > 0$ for some K . We claim $\mu(\cap A_k) \geq K$, a contradiction. For the claim, let $Z = \cap A_k$. If $\mu(Z) \geq \mu(A_k)$ for all k , we are done. If not, since the A_k are nested, we can replace A_k by a set, also called A_k , of measure K and such that the new A_k are still nested. For the relabeled $Z = \cap A_k$, $Z \subset A_k$ for all k , and we may assume $\mu(Z) < K$. Thus there exists $Z' \subset A_1$ with $Z' \cap Z = \emptyset$ and $\mu(Z') > 0$. Since $\mu(A_i) = K$, we must have $A_i \cap Z' \neq \emptyset$ for all i . Thus $\cap A_i$ is strictly larger than Z , a contradiction. In summary, the claim must hold, so we get a contradiction to assuming $\mu(A_k) \not\rightarrow 0$.

(Here $\mathbb{1}_A$ is the characteristic function of a set A .)

As $\delta \rightarrow 0$,

$$\mathbb{E}_P[\mathbb{1}_{\mathcal{X} \setminus (Z_\delta \cup B_\epsilon)} \mathbb{1}_{h_{f_n}(\mathbf{x}) \neq y}] \rightarrow \mathbb{E}_P[\mathbb{1}_{\mathcal{X} \setminus B_\epsilon} \mathbb{1}_{h_{f_n}(\mathbf{x}) \neq y}].$$

and similarly for f_n replaced by $\eta(\mathbf{x})$. During this process, N_δ presumably goes to ∞ , but that precisely means

$$\lim_{n \rightarrow \infty} \mathbb{E}_P[\mathbb{1}_{\mathcal{X} \setminus B_\epsilon} \mathbb{1}_{h_{f_n}(\mathbf{x}) \neq y}] = \mathbb{E}_P[\mathbb{1}_{\mathcal{X} \setminus B_\epsilon} \mathbb{1}_{h_\eta(\mathbf{x}) \neq y}].$$

Since

$$\left| \mathbb{E}_P[\mathbb{1}_{\mathcal{X} \setminus B_\epsilon} \mathbb{1}_{h_{f_n}(\mathbf{x}) \neq y}] - \mathbb{E}_P[\mathbb{1}_{h_{f_n}(\mathbf{x}) \neq y}] \right| < \epsilon,$$

and similarly for $\eta(\mathbf{x})$, we get

$$\begin{aligned} & \left| \lim_{n \rightarrow \infty} \mathbb{E}_P[\mathbb{1}_{h_{f_n}(\mathbf{x}) \neq y}] - \mathbb{E}_P[\mathbb{1}_{h_\eta(\mathbf{x}) \neq y}] \right| \\ & \leq \left| \lim_{n \rightarrow \infty} \mathbb{E}_P[\mathbb{1}_{h_{f_n}(\mathbf{x}) \neq y}] - \lim_{n \rightarrow \infty} \mathbb{E}_P[\mathbb{1}_{\mathcal{X} \setminus B_\epsilon} \mathbb{1}_{h_{f_n}(\mathbf{x}) \neq y}] \right| \\ & \quad + \left| \lim_{n \rightarrow \infty} \mathbb{E}_P[\mathbb{1}_{\mathcal{X} \setminus B_\epsilon} \mathbb{1}_{h_{f_n}(\mathbf{x}) \neq y}] - \mathbb{E}_P[\mathbb{1}_{\mathcal{X} \setminus B_\epsilon} \mathbb{1}_{h_\eta(\mathbf{x}) \neq y}] \right| \\ & \quad + \left| \lim_{n \rightarrow \infty} \mathbb{E}_P[\mathbb{1}_{\mathcal{X} \setminus B_\epsilon} \mathbb{1}_{h_\eta(\mathbf{x}) \neq y}] - \mathbb{E}_P[\mathbb{1}_{h_\eta(\mathbf{x}) \neq y}] \right| \\ & \leq 3\epsilon. \end{aligned}$$

(Strictly speaking, $\lim_{n \rightarrow \infty} \mathbb{E}_P[\mathbb{1}_{h_{f_n}(\mathbf{x}) \neq y}]$ is first lim sup and then lim inf to show that the limit exists.) Since ϵ is arbitrary, the proof is complete if $\mu(Z_0) = 0$.

If $\mu(Z_0) > 0$, we rerun the proof with \mathcal{X} replaced by Z_0 . As above, $f_n|_{Z_0}$ converges uniformly to $\eta(\mathbf{x})$ off a set of measure ϵ . The argument above, without the set Z_δ , gives

$$\int_{Z_0} \mathbb{1}_{h_{f_n}(\mathbf{x}) \neq y} \mu(\mathbf{x}) d\mathbf{x} \rightarrow \int_{Z_0} \mathbb{1}_{h_\eta(\mathbf{x}) \neq y} \mu(\mathbf{x}) d\mathbf{x}.$$

We then proceed with the proof above on $\mathcal{X} \setminus Z_0$. \square

Corollary 3. (Step 2 in general) For our algorithm, $\lim_{n \rightarrow \infty} R_{D,P}(\mathbf{f}_{n,\lambda,m}) = 0 \Rightarrow \lim_{i \rightarrow \infty} R_P(\mathbf{f}_{n,\lambda,m}) = R_P^*$.

Proof. Choose $\mathbf{f}_{1,\lambda,m}$ as in Theorem 2. Since $V_{tot,\lambda,m,\mathbf{f}_{n,\lambda,m}}$ has pointwise length going to zero as $n \rightarrow \infty$, $\{\mathbf{f}_{n,\lambda,m}(\mathbf{x})\}$ is a Cauchy sequence for all \mathbf{x} . This implies that $\mathbf{f}_{n,\lambda,m}$, and not just a subsequence, converges pointwise to η . \square

B.4. Step 3

Lemma 4. (Step 3) If $k \rightarrow \infty$ and $k/m \rightarrow 0$ as $m \rightarrow \infty$, then for $\mathbf{f} \in \text{Maps}(\mathcal{X}, \Delta^{L-1})$,

$$|R_{D,\mathcal{T}_m}(\mathbf{f}) - R_{D,P}(\mathbf{f})| \xrightarrow{m \rightarrow \infty} 0 \text{ in probability,}$$

for all distributions P that generate \mathcal{T}_m .

Proof. Since $R_{D,P}(\mathbf{f})$ is a constant for fixed \mathbf{f} and P , convergence in probability will follow from weak convergence, i.e.,

$$\mathbb{E}_{\mathcal{T}_m}[|R_{D,\mathcal{T}_m}(\mathbf{f}) - R_{D,P}(\mathbf{f})|] \xrightarrow{m \rightarrow \infty} 0.$$

We have

$$\begin{aligned} & |R_{D,\mathcal{T}_m}(\mathbf{f}) - R_{D,P}(\mathbf{f})| \\ & = \left| \int_{\mathcal{X}} \left[d^2 \left(\mathbf{f}(\mathbf{x}), \frac{1}{k} \sum_{i=1}^k \tilde{z}_i \right) - d^2(\mathbf{f}(\mathbf{x}), \eta(\mathbf{x})) \right] d\mathbf{x} \right| \\ & \leq \int_{\mathcal{X}} \left| d^2 \left(\mathbf{f}(\mathbf{x}), \frac{1}{k} \sum_{i=1}^k \tilde{z}_i \right) - d^2(\mathbf{f}(\mathbf{x}), \eta(\mathbf{x})) \right| d\mathbf{x}. \end{aligned}$$

Set $\mathbf{a} = \mathbf{f}(\mathbf{x}) - \frac{1}{k} \sum_{i=1}^k \tilde{z}_i$, $\mathbf{b} = \mathbf{f}(\mathbf{x}) - \eta(\mathbf{x})$. Then

$$\begin{aligned} & \left| \|\mathbf{a}\|_2^2 - \|\mathbf{b}\|_2^2 \right| \\ & = \left| \sum_{\ell=1}^L a_\ell^2 - \sum_{\ell=1}^L b_\ell^2 \right| = \left| \sum_{\ell=1}^L (a_\ell^2 - b_\ell^2) \right| \\ & \leq \sum_{\ell=1}^L |a_\ell^2 - b_\ell^2| \leq 2 \sum_{\ell=1}^L |a_\ell - b_\ell| \max\{|a_\ell|, |b_\ell|\} \\ & \leq 2 \sum_{\ell=1}^L |a_\ell - b_\ell|, \end{aligned}$$

since $f^\ell(\mathbf{x}), \frac{1}{k} \sum_{i=1}^k \tilde{z}_i^\ell, \eta^\ell(\mathbf{x}) \in [0, 1]$. Therefore, it suffices to show that

$$\sum_{\ell=1}^L \mathbb{E}_{\mathcal{T}_m} \left[\int_{\mathcal{X}} \left| \left((f^\ell(\mathbf{x}) - \frac{1}{k} \sum_{i=1}^k \tilde{z}_i^\ell) - (f^\ell(\mathbf{x}) - \eta^\ell(\mathbf{x})) \right) \right| d\mathbf{x} \right] \xrightarrow{m \rightarrow \infty} 0,$$

so the result follows if

$$\lim_{m \rightarrow \infty} \mathbb{E}_{\mathcal{T}_m, \mathbf{x}} \left[\left| \eta^\ell(\mathbf{x}) - \frac{1}{k} \sum_{i=1}^k \tilde{z}_i^\ell \right| \right] = 0 \text{ for all } \ell. \quad (9)$$

By Jensen's inequality ($\mathbb{E}[f]^2 \leq \mathbb{E}[f^2]$), (9) follows if

$$\lim_{m \rightarrow \infty} \mathbb{E}_{\mathcal{T}_m, \mathbf{x}} \left[\left(\eta^\ell(\mathbf{x}) - \frac{1}{k} \sum_{i=1}^k \tilde{z}_i^\ell \right)^2 \right] = 0 \text{ for all } \ell. \quad (10)$$

Let $\eta_{k,m}^\ell(\mathbf{x}) = \frac{1}{k} \sum_{i=1}^k \tilde{z}_i^\ell$. Then $\eta_{k,m}^\ell$ is actually an estimate of the class probability $\eta^\ell(\mathbf{x})$ by the k -Nearest Neighbor rule. Following the proof of Stone's Theorem (Stone, 1977; Devroye et al., 1996), if $k \xrightarrow{m \rightarrow \infty} \infty$ and $k/m \xrightarrow{m \rightarrow \infty} 0$, (10) holds for all distributions P . \square

C. Notation

$$\begin{aligned}
 h_{\mathbf{f}}(\mathbf{x}) = \operatorname{argmax}\{f^\ell(\mathbf{x}), \ell \in \mathcal{Y}\} & : \text{plug-in classifier of estimator } \mathbf{f} : \mathcal{X} \rightarrow \Delta^{L-1} \\
 \mathbb{1}_{h_{\mathbf{f}}(\mathbf{x}) \neq y} & = \begin{cases} 1, & h_{\mathbf{f}}(\mathbf{x}) \neq y, \\ 0, & h_{\mathbf{f}}(\mathbf{x}) = y. \end{cases} \\
 R_P(\mathbf{f}) = \mathbb{E}_P[\mathbb{1}_{h_{\mathbf{f}}(\mathbf{x}) \neq y}] & : \text{generalization risk for the estimator } \mathbf{f} \\
 \boldsymbol{\eta}(\mathbf{x}) = (\eta^1(\mathbf{x}), \dots, \eta^L(\mathbf{x})) & : \text{class probability function: } \eta^\ell(\mathbf{x}) = P(y = \ell | \mathbf{x}) \\
 R_P^* = R_P(\boldsymbol{\eta}) & : \text{Bayes risk} \\
 (D\text{-risk for our } \mathcal{P}_D) R_{D,P}(\mathbf{f}) & = \int_{\mathcal{X}} d^2(\mathbf{f}(\mathbf{x}), \boldsymbol{\eta}(\mathbf{x})) d\mathbf{x} \\
 (\text{empirical } D\text{-risk}) R_{D,\mathcal{T}_m}(\mathbf{f}) = R_{D,\mathcal{T}_m,k}(\mathbf{f}) & = \int_{\mathcal{X}} d^2\left(\mathbf{f}(\mathbf{x}), \frac{1}{k} \sum_{i=1}^k \tilde{\mathbf{z}}_i\right) d\mathbf{x} \\
 & \text{where } \tilde{\mathbf{z}}_i \text{ is the vector of the last } L \text{ components of} \\
 & (\tilde{\mathbf{x}}_i, \tilde{\mathbf{z}}_i), \text{ with } \tilde{\mathbf{x}}_i \text{ the } i^{\text{th}} \text{ nearest neighbor of } \mathbf{x} \text{ in } \mathcal{T}_m \\
 (\text{volume penalty term}) \mathcal{P}_G(\mathbf{f}) & = \int_{\operatorname{gr}(\mathbf{f})} \operatorname{dvol} \\
 R_{D,P,\lambda}(\mathbf{f}) = R_{D,P}(\mathbf{f}) + \lambda \mathcal{P}_G(\mathbf{f}) & : \text{regularized } D\text{-risk for estimator } \mathbf{f} \\
 R_{D,\mathcal{T}_m,\lambda}(\mathbf{f}) = R_{D,\mathcal{T}_m}(\mathbf{f}) + \lambda \mathcal{P}_G(\mathbf{f}) & : \text{regularized empirical } D\text{-risk for estimator } \mathbf{f} \\
 \mathbf{f}_{D,P,\lambda} & = \text{function attaining the global minimum for } R_{D,P,\lambda} \\
 R_{D,P,\lambda}^* = R_{D,P,\lambda}(\mathbf{f}_{D,P,\lambda}) & : \text{minimum value for } R_{D,P,\lambda} \\
 \mathbf{f}_{D,\mathcal{T}_m,\lambda} = \mathbf{f}_{D,\mathcal{T}_m,k,\lambda} & : \text{function attaining the global minimum for } R_{D,\mathcal{T}_m,\lambda}(\mathbf{f})
 \end{aligned}$$

Note that we assume $\mathbf{f}_{D,P,\lambda}$ and $\mathbf{f}_{D,\mathcal{T}_m,\lambda}$ exist.

References

- Devroye, Luc, Györfi, László, and Lugosi, Gábor. *A probabilistic theory of pattern recognition*. Springer, 1996.
- Guckenheimer, John and Worfolk, Patrick. Dynamical systems: some computational problems. In *Bifurcations and periodic orbits of vector fields (Montreal, PQ, 1992)*, volume 408 of *NATO Adv. Sci. Inst. Ser. C Math. Phys. Sci.*, pp. 241–277. Kluwer Acad. Publ., Dordrecht, 1993.
- Stone, Charles. Consistent nonparametric regression. *Annals of Statistics*, pp. 595–620, 1977.