# Appendix: Classification by thresholding

In this appendix, we show how the bounds given in the first section of this monograph can be computed in practice on a simple example: the case when the classification is performed by comparing a series of measurements to threshold values. Let us mention that our description covers the case when the same measurement is compared to several thresholds, since it is enough to repeat a measurement in the list of measurements describing a pattern to cover this case.

## 5.1. Description of the model

Let us assume that the patterns we want to classify are described through $h$ real valued measurements normalized in the range $(0,1)$. In this setting the pattern space can thus be defined as $\mathfrak{X} = (0,1)^h$.

Consider the threshold set $\mathfrak{T} = (0,1)^h$ and the response set $\mathfrak{R} = \mathcal{Y}^{\{0,1\}^h}$. For any $t \in (0,1)^h$ and any $a : \{0,1\}^h \to \mathcal{Y}$, let

$$f_{(t,a)}(x) = a\Big\{ \big[\mathbb{1}(x^j \geq t_j)\big]_{j=1}^{h} \Big\}, \quad x \in \mathfrak{X},$$

where $x^j$ is the $j$th coordinate of $x \in \mathfrak{X}$. Thus our parameter set here is $\Theta = \mathfrak{T} \times \mathfrak{R}$. Let us consider the Lebesgue measure $L$ on $\mathfrak{T}$ and the uniform probability distribution $U$ on $\mathfrak{R}$. Let our prior distribution be $\pi = L \otimes U$. Let us define for any threshold sequence $t \in \mathfrak{T}$

$$\Delta_t = \Big\{ t' \in \mathfrak{T} : \overline{(t'_j, t_j)} \cap \{X_i^j; i = 1, \ldots, N\} = \varnothing, j = 1, \ldots, h \Big\},$$

where $X_i^j$ is the $j$th coordinate of the sample pattern $X_i$, and where the interval $\overline{(t'_j, t_j)}$ of the real line is defined as the convex hull of the two point set $\{t'_j, t_j\}$, whether $t'_j \leq t_j$ or not. We see that $\Delta_t$ is the set of thresholds giving the same response as $t$ on the training patterns. Let us consider for any $t \in \mathfrak{T}$ the middle

$$m(\Delta_t) = \frac{\int_{\Delta_t} t' L(dt')}{L(\Delta_t)}$$

of $\Delta_t$. The set $\Delta_t$ being a product of intervals, its middle is the point whose coordinates are the middle of these intervals. Let us introduce the finite set $T$ composed of the middles of the cells $\Delta_t$, which can be defined as

$$T = \{t \in \mathfrak{T} : t = m(\Delta_t)\}.$$

It is easy to see that $|T| \leq (N+1)^h$ and that $|\mathfrak{R}| = |\mathcal{Y}|^{2^h}$.

155

### 5.2. Computation of inductive bounds

For any parameter $(t, a) \in \mathcal{T} \times \mathcal{R} = \Theta$, let us consider the posterior distribution defined by its density

$$\frac{d\rho_{(t,a)}}{d\pi}(t', a') = \frac{\mathbb{1}(t' \in \Delta_t)\mathbb{1}(a' = a)}{\pi(\Delta_t \times \{a\})}.$$

In fact we are considering a finite number of posterior distributions, since $\rho_{(t,a)} = \rho_{(m(\Delta_t),a)}$, where $m(\Delta_t) \in T$. Moreover, for any exchangeable sample distribution $\mathbb{P} \in \mathcal{M}_+^1[(\mathcal{X} \times \mathcal{Y})^{N+1}]$ and any thresholds $t \in \mathcal{T}$,

$$\mathbb{P}\Big[ \overline{(X_{N+1}^j, t_j)} \cap \{X_i^j, i = 1, \dots, N\} = \varnothing \Big] \leq \frac{2}{N+1}.$$

Thus, for any $(t, a) \in \Theta$,

$$\mathbb{P}\Big\{\rho_{(t,a)}\big[f_.(X_{N+1})\big] \neq f_{(t,a)}(X_{N+1})\Big\} \leq \frac{2h}{N+1},$$

showing that the classification produced by $\rho_{(t,a)}$ on new examples is typically non-random; this result is only indicative, since it is concerned with a non-random choice of $(t, a)$.

Let us compute the various quantities needed to apply the results of the first section, focussing our attention on Theorem 2.1.3 (page 54).

First note that $\rho_{(t,a)}(r) = r[(t, a)]$. The entropy term is such that

$$\mathcal{K}(\rho_{t,a}, \pi) = -\log\big[\pi\big(\Delta_t \times \{r\}\big)\big] = -\log\big[L(\Delta_t)\big] + 2^h \log\big(|\mathcal{Y}|\big).$$

Let us notice accordingly that

$$\min_{(t,a)\in\Theta} \mathcal{K}(\rho_{(t,a)}, \pi) \leq h \log(N + 1) + 2^h \log\big(|\mathcal{Y}|\big).$$

Let us introduce the counters

$$b_y^t(c) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}\Big\{Y_i = y \text{ and } \big[\mathbb{1}(X_i^j \geq t_j)\big]_{j=1}^h = c\Big\},$$

$$t \in T, c \in \{0, 1\}^h, y \in \mathcal{Y},$$

$$b^t(c) = \sum_{y\in\mathcal{Y}} b_y^t(c) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}\Big\{\big[\mathbb{1}(X_i^j \geq t_j)\big]_{j=1}^h = c\Big\}, \qquad t \in T, c \in \{0, 1\}^h.$$

Since

$$r[(t, a)] = \sum_{c\in\{0,1\}^h} \big[b^t(c) - b_{a(c)}^t(c)\big],$$

the partition function of the Gibbs estimator can be computed as

$$\pi\big[\exp(-\lambda r)\big] = \sum_{t\in T} L(\Delta_t) \sum_{a\in\mathcal{R}} \frac{1}{|\mathcal{Y}|^{2^h}} \exp\Big[-\lambda \sum_{i=1}^N \mathbb{1}\big[Y_i \neq f_{(t,a)}(X_i)\big]\Big]$$

$$= \sum_{t\in T} L(\Delta_t) \sum_{a\in\mathcal{R}} \frac{1}{|\mathcal{Y}|^{2^h}} \exp\Big[-\lambda \sum_{c\in\{0,1\}^h} \big[b^t(c) - b_{a(c)}^t(c)\big]\Big]$$

$$= \sum_{t \in T} L(\Delta_t) \prod_{c \in \{0,1\}^h} \left[ \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} \exp\left( -\lambda \left[ b^t(c) - b^t_y(c) \right] \right) \right].$$

We see that the number of operations needed to compute $\pi \left[ \exp(-\lambda r) \right]$ is proportional to $|T| \times 2^h \times |\mathcal{Y}| \le (N+1)^h 2^h |\mathcal{Y}|$. An exact computation will therefore be feasible only for small values of $N$ and $h$. For higher values, a Monte Carlo approximation of this sum will have to be performed instead.

If we want to compute the bound provided by Theorem 2.1.3 (page 54) or by Theorem 2.2.2 (page 68), we need also to compute, for any fixed parameter $\theta \in \Theta$, quantities of the type

$$\pi_{\exp(-\lambda r)} \left\{ \exp \left[ \xi m'(\cdot, \theta) \right] \right\} = \pi_{\exp(-\lambda r)} \left\{ \exp \left[ \xi \rho_\theta(m') \right] \right\}, \quad \lambda, \xi \in \mathbb{R}_+.$$

We need to introduce

$$\overline{b}^t_y(\theta, c) = \frac{1}{N} \sum_{i=1}^N \left| \mathbb{1}\left[ f_\theta(X_i) \ne Y_i \right] - \mathbb{1}(y \ne Y_i) \right| \mathbb{1}\left\{ \left[ \mathbb{1}(X^j_i \ge t_j) \right]^h_{j=1} = c \right\}.$$

Similarly to what has been done previously, we obtain

$$\pi \left\{ \exp \left[ -\lambda r + \xi m'(\cdot, \theta) \right] \right\}$$
$$= \sum_{t \in T} L(\Delta_t) \prod_{c \in \{0,1\}^h} \left[ \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} \exp\left( -\lambda \left[ b^t(c) - b^t_y(c) \right] + \xi \overline{b}^t_y(\theta, c) \right) \right].$$

We can then compute

$$\pi_{\exp(-\lambda r)}(r) = -\frac{\partial}{\partial \lambda} \log \left\{ \pi \left[ \exp(-\lambda r) \right] \right\},$$

$$\pi_{\exp(-\lambda r)} \left\{ \exp \left[ \xi \rho_\theta(m') \right] \right\} = \frac{\pi \left\{ \exp \left[ -\lambda r + \xi m'(\cdot, \theta) \right] \right\}}{\pi \left[ \exp(-\lambda r) \right]},$$

$$\pi_{\exp(-\lambda r)} \left[ m'(\cdot, \theta) \right] = \frac{\partial}{\partial \xi} \bigg|_{\xi=0} \log \left[ \pi \left\{ \exp \left[ -\lambda r + \xi m'(\cdot, \theta) \right] \right\} \right].$$

This is all we need to compute $B(\rho_\theta, \beta, \gamma)$ (and also $B(\pi_{\exp(-\lambda r)}, \beta, \gamma)$) in Theorem 2.1.3 (page 54), using the approximation

$$\log \left\{ \pi_{\exp(-\lambda_1 r)} \left[ \exp \left\{ \xi \pi_{\exp(-\lambda_2 r)}(m') \right\} \right] \right\}$$
$$\le \log \left\{ \pi_{\exp(-\lambda_1 r)} \left[ \exp \left\{ \xi m'(\cdot, \theta) \right\} \right] \right\} + \xi \pi_{\exp(-\lambda_2 r)} \left[ m'(\cdot, \theta) \right], \quad \xi \ge 0.$$

Let us also explain how to apply the posterior distribution $\rho_{(t,a)}$, in other words our randomized estimated classification rule, to a new pattern $X_{N+1}$:

$$\rho_{(t,a)} \left[ f.(X_{N+1}) = y \right] = L(\Delta_t)^{-1} \int_{\Delta_t} \mathbb{1}\left[ a\left\{ \left[ \mathbb{1}(X^j_{N+1} \ge t'_j) \right]^h_{j=1} \right\} = y \right] L(dt')$$
$$= L(\Delta_t)^{-1} \sum_{c \in \{0,1\}^h} L\left( \left\{ t' \in \Delta_t : \left[ \mathbb{1}(X^j_{N+1} \ge t'_j) \right]^h_{j=1} = c \right\} \right) \mathbb{1}\left[ a(c) = y \right].$$

Let us define for short

$$\Delta_t(c) = \left\{ t' \in \Delta_t : \left[ \mathbb{1}(X^j_{N+1} \ge t'_j) \right]^h_{j=1} = c \right\}, \qquad c \in \{0,1\}^h.$$

With this notation

$$\rho_{(t,a)}\big[f.(X_{N+1}) = y\big] = L\big(\Delta_t\big)^{-1} \sum_{c \in \{0,1\}^h} L\big[\Delta_t(c)\big] \mathbb{1}\big[a(c) = y\big].$$

We can compute in the same way the probabilities for the label of the new pattern under the Gibbs posterior distribution:

$$\pi_{\exp(-\lambda r)}\big[f.(X_{N+1}) = y'\big]$$

$$= \left\{ \sum_{t \in T} \prod_{c \in \{0,1\}^h} \left[ \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} \exp\left(-\lambda\big[b^t(c) - b^t_y(c)\big]\right) \right] \right.$$

$$\times \sum_{c \in \{0,1\}^h} L\big[\Delta_t(c)\big] \frac{\sum_{y \in \mathcal{Y}} \mathbb{1}(y = y') \exp\{-\lambda\big[b^t(c) - b^t_y(c)\big]\}}{\sum_{y \in \mathcal{Y}} \exp\{-\lambda\big[b^t(x) - b^t_y(c)\big]\}} \right\}$$

$$\times \left\{ \sum_{t \in T} L(\Delta_t) \prod_{c \in \{0,1\}^h} \left[ \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} \exp\left(-\lambda\big[b^t(c) - b^t_y(c)\big]\right) \right] \right\}^{-1}.$$

### 5.3. Transductive bounds

In the case when we observe the patterns of a shadow sample $(X_i)_{i=N+1}^{(k+1)N}$ on top of the training sample $(X_i, Y_i)_{i=1}^N$, we can introduce the set of thresholds responding as $t$ on the extended sample $(X_i)_{i=1}^{(k+1)N}$

$$\overline{\Delta}_t = \Big\{ t' \in \mathcal{T} : \overline{(t'_j, t_j)} \cap \big\{ X_i^j; i = 1, \ldots, (k+1)N \big\} = \varnothing, j = 1, \ldots, h \Big\},$$

consider the set

$$\overline{T} = \big\{ t \in \mathcal{T} : t = m(\overline{\Delta}_t) \big\},$$

of the middle points of the cells $\overline{\Delta}_t$, $t \in \mathcal{T}$, and replace the Lebesgue measure $L \in \mathcal{M}_+^1\big[(0,1)^h\big]$ of the previous section with the uniform probability measure $\overline{L}$ on $\overline{T}$. We can then consider $\pi = \overline{L} \otimes U$, where $U$ is as previously the uniform probability measure on $\mathcal{R}$. This gives obviously an exchangeable posterior distribution and therefore qualifies $\pi$ for transductive bounds. Let us notice that $|\overline{T}| \leq \big[(k+1)N + 1\big]^h$, and therefore that $\pi(t,a) \geq \big[(k+1)N + 1\big]^{-h} |\mathcal{Y}|^{-2^h}$, for any $(t,a) \in \overline{T} \times \mathcal{R}$.

For any $(t,a) \in \mathcal{T} \times \mathcal{R}$ we may similarly to the inductive case consider the posterior distribution $\rho_{(t,a)}$ defined by

$$\frac{d\rho_{(t,a)}}{d\pi}(t', a') = \frac{\mathbb{1}(t' \in \Delta_t)\mathbb{1}(a' = a)}{\pi(\Delta_t \times \{a\})},$$

but we may also consider $\delta_{(m(\overline{\Delta}_t),a)}$, which is such that $r_i\{[m(\overline{\Delta}_t), a]\} = r_i[(t,a)]$, $i = 1, 2$, whereas only $\rho_{(t,a)}(r_1) = r_1[(t,a)]$, while

$$\rho_{(t,a)}(r_2) = \frac{1}{|\overline{T} \cap \Delta_t|} \sum_{t' \in \overline{T} \cap \Delta_t} r_2[(t', a)].$$

We get

$$\mathcal{K}(\rho_{(t,a)}, \pi) = -\log\big[\overline{L}(\Delta_t)\big] + 2^h \log(|\mathcal{Y}|)$$
$$\leq \log\big(|\overline{T}|\big) + 2^h \log(|\mathcal{Y}|) = \mathcal{K}(\delta_{[m(\overline{\Delta}_t),a]}, \pi)$$
$$\leq h \log\big[(k+1)N + 1\big] + 2^h \log(|\mathcal{Y}|),$$

whereas we had no such uniform bound in the inductive case. Similarly to the inductive case

$$\pi\big[\exp(-\lambda r_1)\big] = \sum_{t \in T} \overline{L}(\Delta_t) \prod_{c \in \{0,1\}^h} \left[ \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} \exp\Big(-\lambda\big[b^t(c) - b_y^t(c)\big]\Big) \right].$$

Moreover, for any $\theta \in \Theta$,

$$\pi\big\{\exp\big[-\lambda r_1 + \xi \rho_\theta(m')\big]\big\} = \pi\big\{\exp\big[-\lambda r_1 + \xi m'(\cdot, \theta)\big]\big\}$$
$$= \sum_{t \in T} \overline{L}(\Delta_t) \prod_{c \in \{0,1\}^h} \left[ \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} \exp\Big(-\lambda\big[b^t(c) - b_y^t(c)\big] + \xi \overline{b}(\theta, c)\Big) \right].$$

The bound for the transductive counterpart to Theorems 2.1.3 (page 54) or 2.2.2 (page 68), obtained as explained page 115, can be computed as in the inductive case, from these two partition functions and the above entropy computation.

Let us mention finally that, using the same notation as in the inductive case,

$$\pi_{\exp(-\lambda r_1)}\big[f.(X_{N+1}) = y'\big]$$
$$= \left\{ \sum_{t \in T} \prod_{c \in \{0,1\}^h} \left[ \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} \exp\Big(-\lambda\big[b^t(c) - b_y^t(c)\big]\Big) \right] \right.$$
$$\times \sum_{c \in \{0,1\}^h} \overline{L}\big[\Delta_t(c)\big] \frac{\sum_{y \in \mathcal{Y}} \mathbb{1}(y = y') \exp\big\{-\lambda\big[b^t(c) - b_y^t(c)\big]\big\}}{\sum_{y \in \mathcal{Y}} \exp\big\{-\lambda\big[b^t(x) - b_y^t(c)\big]\big\}} \right\}$$
$$\times \left\{ \sum_{t \in T} \overline{L}(\Delta_t) \prod_{c \in \{0,1\}^h} \left[ \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} \exp\Big(-\lambda\big[b^t(c) - b_y^t(c)\big]\Big) \right] \right\}^{-1}.$$

To conclude this appendix on classification by thresholding, note that similar factorized computations are feasible in the important case of *classification trees*. This can be achieved using some variant of the *context tree weighting method* discovered by Willems et al. (1995) and successfully used in lossless compression theory. The interested reader can find a description of this algorithm applied to classification trees in Catoni (2004, page 62).