

Supplementary Material for Optimal Bayes Classifiers for Functional Data and Density Ratios

By Xiongtao Dai, Hans-Georg Müller

Department of Statistics, University of California, Davis, California 95616, U.S.A.

dai@ucdavis.edu hgmuller@ucdavis.edu

and Fang Yao

*Department of Statistical Sciences, University of Toronto, 100 St. George Street,
 Toronto, Ontario M5S 3G3, Canada*

fyao@utstat.toronto.edu

1. An Alternative Nonparametric Regression Estimate

An alternative approach for estimating the density ratios is via nonparametric regression. This is motivated by Bayes' theorem, as follows,

$$\begin{aligned} \frac{f_{j1}(u)}{f_{j0}(u)} &= \frac{\text{pr}(Y = 1 \mid j = u)p_j(u)}{\text{pr}(Y = 0 \mid j = u)p_j(u)} \frac{\text{pr}(Y = 1)}{\text{pr}(Y = 0)} \\ &= \frac{\text{pr}(Y = 1 \mid j = u)}{\text{pr}(Y = 0 \mid j = u)} \frac{\pi_1}{\pi_0} = \frac{\pi_0 \text{pr}(Y = 1 \mid j = u)}{\pi_1 (1 - \text{pr}(Y = 1 \mid j = u))} \end{aligned} \quad (1)$$

where $p_j(\cdot)$ is the marginal density of the j th projection. This reduces the construction of nonparametric Bayes classifiers to a sequence of nonparametric regressions $E(Y \mid j = u) = \text{pr}(Y = 1 \mid j = u)$. These again can be implemented by a kernel method (Nadaraya, 1964; Watson, 1964), smoothing the scatter plots of the pooled estimated scores $\hat{\xi}_{ijk}$ of group k , which leads to the nonparametric estimators

$$\hat{E}(Y \mid \hat{j} = u) = \frac{\sum_{k=0}^1 \sum_{i=1}^{n_k} K\left(\frac{u - \hat{\xi}_{ijk}}{h_j}\right)}{\sum_{k=0}^1 \sum_{i=1}^{n_k} K\left(\frac{u - \hat{\xi}_{ijk}}{h_j}\right)}$$

where $h_j = \{(\hat{j}_0 + \hat{j}_1)/2\}^{1/2}$ is the bandwidth. This results in estimates $\hat{E}(Y \mid \hat{j} = u) = \hat{\text{pr}}(Y = 1 \mid \hat{j} = u)$ that we plug-in at the right hand side of (1), which then yields an alternative estimate of the density ratio, replacing the two kernel density estimates $\hat{f}_{j1}(u)$ $\hat{f}_{j0}(u)$ by just one nonparametric regression estimate $\hat{E}(Y \mid \hat{j} = u)$.

The estimated criterion function based on kernel regression is

$$\hat{R}_J(x) = \log \frac{\hat{\pi}_1}{\hat{\pi}_0} + \sum_{j \leq J} \log \frac{\hat{\pi}_0 \hat{E}(Y \mid \hat{j} = u)}{\hat{\pi}_1 \{1 - \hat{E}(Y \mid \hat{j} = u)\}}$$

2. Perfect Classification when the Mean and the Covariance Functions are the Same

Let the projection scores ξ_j be independent random variables with mean 0 and variance ν_j that follow normal distributions under Π_1 and Laplace distributions under Π_0 . Then

$$\begin{aligned}
 J(X) &= \sum_{j=1}^J \log \frac{\frac{1}{(2\pi\nu_j)^{1/2}} \exp(-\frac{\xi_j^2}{2\nu_j})}{\frac{1}{(2\nu_j)^{1/2}} \exp\{-\frac{|\xi_j|}{(\nu_j/2)^{1/2}}\}} \\
 &= \sum_{j=1}^J \left(-\frac{1}{2} \log \pi - \frac{\xi_j^2}{2\nu_j} + \sqrt{2} |\xi_j| \nu_j^{-1/2} \right)
 \end{aligned} \tag{2}$$

Since centred normal and Laplace distributions are in scale families, $\xi_j = \nu_j^{-1/2} \eta_j$ have a common standard distribution η_k under Π_k , irrespective of k . Denoting the summand of (2) by η_j , this implies $\eta_j = -(\log \pi + \frac{\xi_j^2}{2\nu_j}) + \sqrt{2} |\xi_j| \nu_j^{-1/2}$ are independent and identically distributed. Note that $E_{\Pi_0}(\eta_j) = -(\log \pi + 1) + 1 = 0$, $E_{\Pi_1}(\eta_j) = -(\log \pi + 1) + (\pi/2)^{-1/2}$, and η_j has finite variance under either population. So the misclassification error under Π_0 is

$$\begin{aligned}
 \Pr_{\Pi_0}(J(X) \leq 0) &= \Pr_{\Pi_0}\left(\sum_{j=1}^J \eta_j - E_{\Pi_0}\left(\sum_{j=1}^J \eta_j\right) \leq -E_{\Pi_0}\left(\sum_{j=1}^J \eta_j\right)\right) \\
 &\leq \frac{\text{var}_{\Pi_0}\left(\sum_{j=1}^J \eta_j\right)}{E_{\Pi_0}\left(\sum_{j=1}^J \eta_j\right)^2} \\
 &= \frac{J \text{var}_{\Pi_0}(\eta_1)}{J^2 E_{\Pi_0}(\eta_1)^2} \rightarrow 0
 \end{aligned}$$

as $J \rightarrow \infty$, where the inequality is due to Chebyshev's inequality and the last equality is due to η_j are independently and identically distributed. Similarly, the misclassification error under Π_1 also goes to zero as $J \rightarrow \infty$. Therefore perfect classification occurs under this non-Gaussian case where both the mean and the covariance functions are the same.

3. Simulation Results without Pre-smoothing

The misclassification results when using predictor functions sampled with noise that are not presmoothed are reported in Table 1. When the covariances are the same but the means differ, the centroid method is the overall best if we use the noisy predictors while the Gaussian implementation of the proposed Bayes classifiers has comparable performance. This is expected because our method estimates more parameters than the centroid method while both assume the correct model for the simulated data. All methods gain performance from pre-smoothing due to the presence of noise in the predictor functions. The logistic method benefits the most from pre-smoothing and becomes the winner when only a mean difference is present.

Table 1. Misclassification rates (%), with standard errors in brackets for raw predictors

n			Centroid	Gaussian	NPD	NPR	Logistic
Scenario A (Gaussian case)							
50	same	diff	49.3 (0.12)	23.8 (0.18)	24.5 (0.21)	26.7 (0.22)	49.4 (0.12)
	diff	same	40.2 (0.16)	41.5 (0.16)	43.4 (0.17)	42.4 (0.18)	40.7 (0.16)
	diff	diff	37.9 (0.17)	20.8 (0.18)	21.2 (0.20)	23.3 (0.22)	38.8 (0.17)
100	same	diff	49.1 (0.13)	17.2 (0.11)	18.6 (0.12)	20.0 (0.13)	49.3 (0.13)
	diff	same	37.8 (0.13)	39.2 (0.13)	41.4 (0.15)	40.2 (0.16)	38.3 (0.13)
	diff	diff	35.3 (0.14)	14.6 (0.1)	15.8 (0.10)	17.1 (0.12)	35.8 (0.15)
Scenario B (exponential case)							
50	same	diff	49.0 (0.13)	30.2 (0.19)	31.2 (0.22)	33.5 (0.23)	49.2 (0.13)
	diff	same	38.3 (0.21)	40.6 (0.21)	39.5 (0.22)	38.6 (0.21)	38.7 (0.23)
	diff	diff	35.0 (0.20)	23.3 (0.18)	23.5 (0.21)	24.3 (0.22)	35.7 (0.22)
100	same	diff	48.8 (0.14)	26.0 (0.13)	25.4 (0.14)	26.7 (0.16)	48.9 (0.13)
	diff	same	35.8 (0.16)	38.6 (0.19)	36.3 (0.18)	35.7 (0.16)	35.9 (0.16)
	diff	diff	32.4 (0.14)	18.7 (0.13)	16.7 (0.13)	17.0 (0.14)	32.7 (0.15)
Scenario C (dependent case)							
50	same	diff	48.9 (0.14)	33.3 (0.19)	35.3 (0.22)	37.3 (0.22)	49.1 (0.14)
	diff	same	39.3 (0.22)	42.1 (0.21)	41.0 (0.22)	40.1 (0.22)	39.2 (0.23)
	diff	diff	36.0 (0.21)	27.3 (0.20)	28.6 (0.21)	29.3 (0.23)	36.7 (0.23)
100	same	diff	49.1 (0.13)	29.8 (0.14)	30.6 (0.14)	31.8 (0.15)	49.0 (0.13)
	diff	same	36.4 (0.17)	39.8 (0.20)	37.9 (0.18)	37.1 (0.17)	36.3 (0.16)
	diff	diff	33.3 (0.16)	24.1 (0.15)	22.6 (0.15)	22.9 (0.16)	33.5 (0.16)

Centroid method; Gaussian, NPD, and NPR correspond to the Gaussian, nonparametric density, and nonparametric regression implementations of the proposed Bayes classifiers, respectively; Logistic, functional logistic regression.

4. Proofs

55

4.1. Theorem A1 and Theorem A2

Let $\mathcal{S}(c) = \{x : \|x\| \leq c\}$ be a bounded set of all square integrable functions for $c > 0$, where $\|\cdot\|$ is the L^2 norm. We will use the following lemma:

Lemma 1. Under Conditions A1–A4, for any $\alpha \geq 1$, $\beta = 0, 1$,

$$\sup_{x \in \mathcal{S}(c)} |\hat{f}_{jk}(\hat{x}_j) - f_{jk}(x_j)| = O_p \left(\frac{n^{-1/2}}{\log n} \right) + \frac{n^{-1/2}}{\log n}$$

60 *Proof.* We prove the statement for $\alpha = 0$; the proof for $\alpha = 1$ is analogous. Let the sample mean of the j th estimated projection \hat{x}_j be \bar{x}_j . Observe

$$\begin{aligned} & \sup_{x \in \mathcal{S}(c)} \left| \hat{x}_j - \bar{x}_j \right| - \frac{x_j}{\sqrt{j}} \\ & \leq \sup_{x \in \mathcal{S}(c)} \left| \hat{x}_j - \bar{x}_j \right| - \frac{x_j}{\sqrt{j}} + \sup_{x \in \mathcal{S}(c)} \left| \bar{x}_j - \frac{x_j}{\sqrt{j}} \right| \\ & = o_p\{(n \log n)^{-1/2}\} + O_p\left(\frac{n}{\log n}\right)^{-1/2} = O_p\left(\frac{n}{\log n}\right)^{-1/2} \end{aligned} \quad (3)$$

65 where the first rate is due to Theorem 3.1 in Delaigle & Hall (2010), and the second to, for example, Theorem 2 in Stone (1983). Then

$$\begin{aligned} \sup_{x \in \mathcal{S}(c)} |\hat{f}_j(x) - f_j(x)| &= \sup_{x \in \mathcal{S}(c)} \left| \frac{1}{\sqrt{j}} \hat{x}_j - \frac{1}{\sqrt{j}} \bar{x}_j \right| \\ &\leq \sup_{x \in \mathcal{S}(c)} \left| \hat{x}_j - \bar{x}_j \right| + \left| \frac{1}{\sqrt{j}} \bar{x}_j - \frac{1}{\sqrt{j}} \right| \\ &= O_p \left(\sup_{x \in \mathcal{S}(c)} \left| \hat{x}_j - \bar{x}_j \right| \right) + O_p \left(\frac{1}{\sqrt{j}} \right) \\ &= O_p \left(\frac{n}{\log n} \right)^{-1/2} \end{aligned}$$

70

where the second equality follows from the consistency of \hat{x}_j and Condition A4, and the third equality follows from (3) and the fact that \bar{x}_j converges at a root- n rate. \square

Proof of Theorem A1. For simplicity we consider the case where the supports of f_{j0} and f_{j1} are in common. The case where the supports differ can be proven in two steps: 75 First consider classifying elements x whose projections x_j are in the intersection of the supports of f_{j0} and f_{j1} ; next consider classifying an element x for which a projection score x_j is not contained in the intersection of the supports, in which case $\hat{f}_j(x)$ will be $\pm\infty$, whence $\hat{f}_j(x)$ will also diverge to $\pm\infty$, respectively, and thus consistency is obtained.

80 Now fix $\alpha \in [0, 1]$. Set c be such that $\Pr(\|X\| \leq c) = \Pr\{X \in \mathcal{S}(c)\} \leq \frac{1}{2}$. First we prove there exists an event \mathcal{E}_n such that $\hat{f}_j(X) - f_j(X) \rightarrow 0$ on \mathcal{E}_n with $\Pr(\mathcal{E}_n) \geq 1 - \frac{1}{n}$. By Lemma 1 there exists $M_{jk} > 0$ such that the event

$$\mathcal{E}_n = \sup_{x \in \mathcal{S}(c)} |\hat{f}_{jk}(\hat{x}_j) - f_{jk}(x_j)| \leq M_{jk} + \frac{n}{\log n}^{-1/2} \quad (\alpha = 1/2; \alpha = 0/1)$$

has probability $\Pr(\mathcal{E}_n) \geq 1 - 2^{-(j+2)}$. Letting $\mathcal{E}_n = \bigcap_{j \geq 1, k=0,1} \mathcal{E}_{n,jk}$

85 $\bigcap_{j \geq 1, k=0,1} \{j \in \text{supp}(f_{jk})\} \cap \{\|X\| \leq c\}$, we have $\Pr(\mathcal{E}_n) \geq 1 - \frac{1}{n}$, where supp means the support of a density. Let a_n be some increasing sequence such that $a_n \rightarrow \infty$ and

$a_n \{ + (n / \log n)^{-1/2} \} = o(1)$. Define $\mathcal{U}_{jk} = \{x : x_j \in \text{supp}(f_{jk})\}$, $\mathcal{U} = \bigcap_{j \geq 1, k=0,1} \mathcal{U}_{jk}$,

$$d_{jk} = \min\{1, \inf_{x \in \mathcal{S}(c) \cap \mathcal{U}} f_{jk}(x)\} \quad J = \sup \{J' \geq 1 : \frac{M_{jk}}{d_{jk}} \leq a_n\}_{j \leq J', k=0,1}$$

The d_{jk} are bounded away from 0 by Condition A5, and J is nondecreasing and tends to infinity as $n \rightarrow \infty$. ~~On~~ we have

$$\begin{aligned} \frac{1}{d_{jk}} \sup_{x \in \mathcal{S}(c)} |\hat{f}_{jk}(\hat{x}_j) - f_{jk}(x_j)| &\leq \frac{M_{jk}}{d_{jk}} + \frac{n^{-1/2}}{\log n} \\ &\leq a_n + \frac{n^{-1/2}}{\log n} = o(1) \end{aligned} \quad (4) \quad 90$$

where the first and second inequalities are due to the property ~~of~~ and J , respectively, and the last equality is by the definition of a_n .

From (4) we infer that ~~on~~,

$$\sup_{x \in \mathcal{S}(c)} |\hat{f}_{jk}(\hat{x}_j) - f_{jk}(x_j)| \leq d_{jk} \cdot 2 \quad (5)$$

eventually and uniformly for all $j \leq J$. ~~Then on~~ it holds that

$$\begin{aligned} |\hat{\pi}_J(X) - \pi_J(X)| &\leq \sup_{x \in \mathcal{S}(c) \cap \mathcal{U}} |\hat{\pi}_J(x) - \pi_J(x)| \\ &\leq \sup_{j \leq J, k=0,1} \sup_{x \in \mathcal{S}(c) \cap \mathcal{U}} |\log \hat{f}_{jk}(\hat{x}_j) - \log f_{jk}(x_j)| \\ &\leq \sup_{j \leq J, k=0,1} \sup_{x \in \mathcal{S}(c)} |\hat{f}_{jk}(\hat{x}_j) - f_{jk}(x_j)| \frac{1}{\inf_{x \in \mathcal{S}(c) \cap \mathcal{U}} \pi_{3jk}} \\ &\leq \sup_{j \leq J, k=0,1} \sup_{x \in \mathcal{S}(c)} |\hat{f}_{jk}(\hat{x}_j) - f_{jk}(x_j)| \frac{2}{d_{jk}} \\ &= o(1) \end{aligned} \quad 95$$

where the third inequality is by Taylor's theorem, π_{3jk} is between $f_{jk}(x_j)$ and $\hat{f}_{jk}(\hat{x}_j)$, the last inequality is due to (5) which holds for large enough n , and the equality is due to (4). We conclude that ~~pr~~ $\{ \hat{\pi}_J(X) \geq 0 \} \neq \{ \pi_J(X) \geq 0 \} \} \rightarrow 0$ as $n \rightarrow \infty$ by noting that $\hat{\pi}_J(X)$ converges to $\pi_J(X)$ and thus has the same sign as $\pi_J(X)$ as $n \rightarrow \infty$. Notice that $\pi_J(X)$ has a continuous density and thus $\text{pr}\{ \pi_J(X) = 0 \} = 0$ by Condition A4. \square 100 105

Proof of Theorem A2. Recall we assume $\hat{\pi}_1 = \hat{\pi}_0$. Then

$$\begin{aligned} \hat{E}(Y \mid \hat{X}_j = u) &= \frac{\sum_{k=0}^1 \sum_{i=1}^{n_k} K\left(\frac{u - \hat{\xi}_{ijk}}{h_j}\right)}{\sum_{k=0}^1 \sum_{i=1}^{n_k} K\left(\frac{u - \hat{\xi}_{ijk}}{h_j}\right)} \\ &= \frac{\sum_{i=1}^{n_1} K\left(\frac{u - \hat{\xi}_{ij1}}{h_j}\right)}{\sum_{i=1}^{n_1} K\left(\frac{u - \hat{\xi}_{ij1}}{h_j}\right) + \sum_{i=1}^{n_0} K\left(\frac{u - \hat{\xi}_{ij0}}{h_j}\right)} \\ &= \frac{\hat{f}_{j1}(u)}{\hat{f}_{j1}(u) + \hat{f}_{j0}(u)} \end{aligned}$$

110 where \hat{f}_{jk} are the kernel density estimators with bandwidth h_j , implying

$$\begin{aligned} \hat{R}_J^R(x) &= \sum_{j=1}^J \log \frac{\hat{E}(Y \mid \hat{X}_j = \hat{x}_j)}{\{1 - \hat{E}(Y \mid \hat{X}_j = \hat{x}_j)\}} \\ &= \sum_{j=1}^J \log \frac{\hat{f}_{j1}(\hat{x}_j)}{\hat{f}_{j0}(\hat{x}_j)} \end{aligned}$$

Observe that \hat{R}_J^R has the same form as \hat{R}_J , so this result follows from Theorem A1. \square

4.2. Theorem 1

115 The proof of Theorem 1 requires the following key lemma, which is an extension of Lemma 1, changing the rate from $(n \log n)^{-1/2}$ to $(n \log n)^{-1/2} + (n^{2/5} \log n)^{-1}$. The remainder of the proof is omitted, since it is analogous to that of Theorem A1.

Lemma 2. *Under Conditions A1–A4 and A6–A9, for any $\alpha \geq 1$, $\beta = 0$, 1,*

$$\sup_{x \in \mathcal{S}(c)} |\tilde{f}_{jk}(\tilde{x}_j) - f_{jk}(x_j)| = O_p \left(\frac{n^{-1/2}}{\log n} + (n^{2/5} \log n)^{-1} \right)$$

Proof. Given $x \in \mathcal{S}(c)$, by triangle inequality

$$|\tilde{f}_{jk}(\tilde{x}_j) - f_{jk}(x_j)| \leq |\tilde{f}_{jk}(\tilde{x}_j) - \hat{f}_{jk}(\hat{x}_j)| + |\hat{f}_{jk}(\hat{x}_j) - f_{jk}(x_j)|$$

120 The rate for the second term can be derived from Lemma 1, so we focus only on the first term. For fixed j, k and $\tilde{x}_j = \frac{1}{2}(\tilde{x}_j + x_j)$,

$$\begin{aligned} |\tilde{f}_{jk}(\tilde{x}_j) - \hat{f}_{jk}(\hat{x}_j)| &= \frac{1}{n_k} \sum_{i=1}^{n_k} K \frac{\int_{\mathcal{T}} \{\tilde{X}_i^{(k)}(t) - x(t)\} \tilde{w}_j(t) dt}{\int_{\mathcal{T}} \tilde{w}_j(t) dt} - K \frac{\int_{\mathcal{T}} \{X_i^{(k)}(t) - x(t)\} \hat{w}_j(t) dt}{\int_{\mathcal{T}} \hat{w}_j(t) dt} \\ &\leq \frac{1}{n_k} \sum_{i=1}^{n_k} \int_{\mathcal{T}} \{ \tilde{X}_i^{(k)}(t) - X_i^{(k)}(t) \} \tilde{w}_j(t) - \hat{w}_j(t) dt \cdot |K'(\cdot_{4jk})| \\ &\leq \frac{c_3}{n_k} \sum_{i=1}^{n_k} \int_{\mathcal{T}} \{ \tilde{X}_i^{(k)}(t) - X_i^{(k)}(t) \} \tilde{w}_j(t) - \hat{w}_j(t) dt \end{aligned} \tag{6}$$

for a constant $c_3 > 0$, where the first inequality is by Taylor's theorem, μ_{4jk} is a mean value, and the last inequality is by Condition A4. The summand in (6) is

$$\begin{aligned}
 & \int_{\mathcal{T}} \{\tilde{X}_i^{(k)}(t) - x(t)\} \tilde{y}_j(t) - \{X_i^{(k)}(t) - x(t)\} \hat{y}_j(t) dt \\
 = & \int_{\mathcal{T}} \{\tilde{X}_i^{(k)}(t) - X_i^{(k)}(t)\} \tilde{y}_j(t) + \{X_i^{(k)}(t) - x(t)\} \{\tilde{y}_j(t) - \hat{y}_j(t)\} dt \\
 \leq & \int_{\mathcal{T}} \{\tilde{X}_i^{(k)}(t) - X_i^{(k)}(t)\} \tilde{y}_j(t) dt + \int_{\mathcal{T}} \{X_i^{(k)}(t) - x(t)\} \{\tilde{y}_j(t) - \hat{y}_j(t)\} dt \\
 \leq & \|\tilde{X}_i^{(k)} - X_i^{(k)}\| \cdot \|\tilde{y}_j\| + \|X_i^{(k)} - x\| \cdot \|\tilde{y}_j - \hat{y}_j\| \\
 \leq & \|\tilde{X}_i^{(k)} - X_i^{(k)}\| + (\|X_i^{(k)}\| + c) \|\tilde{y}_j - \hat{y}_j\|
 \end{aligned} \tag{130}$$

where the second and third inequalities follow from Cauchy-Schwarz inequality and from $\|x\| \leq c$, respectively. Plugging the previous result into (6),

$$|\tilde{f}_{jk}(\tilde{x}_j) - \hat{f}_{jk}(\hat{x}_j)| \leq \frac{c_3}{2} \frac{1}{n_k} \sum_{i=1}^{n_k} \|\tilde{X}_i^{(k)} - X_i^{(k)}\| + \|\tilde{y}_j - \hat{y}_j\| \left(\frac{1}{n_k} \sum_{i=1}^{n_k} \|X_i^{(k)}\| + c \right) \tag{7}$$

Since $(\tilde{X}_i^{(k)} | X_i^{(k)})$ are identically distributed ($\epsilon = 1$ n_k), and by Condition A6 the first term in the brackets has expected value equal to

$$E(\|\tilde{X}_i^{(k)} - X_i^{(k)}\|) = E_{X_i^{(k)}}\{E_{\epsilon_i}(\|\tilde{X}_i^{(k)} - X_i^{(k)}\| | X_i^{(k)})\} = O\{(\sigma)^{-1/2} + \sigma^2\} = O(\sigma^{-2/5})$$

where more details about the second equality can be found in the Supplementary Material of Kong et al. (2016). Also $E(\frac{1}{n_k} \sum_{i=1}^{n_k} \|X_i^{(k)}\| + c) = O(1)$ by Condition A1. So

$$\frac{1}{n_k} \sum_{i=1}^{n_k} \|\tilde{X}_i^{(k)} - X_i^{(k)}\| = O_p(\sigma^{-2/5}) \quad \frac{1}{n_k} \sum_{i=1}^{n_k} \|X_i^{(k)}\| + c = O_p(1) \tag{8}$$

It remains to be shown $\|\tilde{y}_j - \hat{y}_j\| = O_p(\sigma^{-2/5})$. Let $\tilde{\Delta}_k = \tilde{G}_k - \hat{G}_k$ and for a square-integrable function $A(s, t)$ denote $\|A\|_F = \left\{ \int_{\mathcal{T}} \int_{\mathcal{T}} A(s, t)^2 ds dt \right\}^{1/2}$ be the Frobenius norm. In their Supplementary Material, Kong et al. (2016) show that $\|\tilde{\Delta}_k\|_F = O_p(\sigma^{-2/5})$, so $\|\tilde{\Delta}\|_F = \|\tilde{\Delta}_0 + \tilde{\Delta}_1\|_F \leq 2 = O_p(\sigma^{-2/5})$. By standard perturbation theory for operators (Bosq, 2000), for a fixed

$$\|\tilde{y}_j - \hat{y}_j\| = O(\|\tilde{\Delta}\|_F / \sup_{k \neq j} |\hat{y}_j - \hat{y}_k|) = O_p(\sigma^{-2/5}) \tag{9}$$

Inserting (8) and (9) into (7) we arrive at the conclusion. \square

4.3. Theorem 2

Assuming X is Gaussian under $\theta = 0, 1$, whence the criterion function $J(x)$ defined in (3) in the main text becomes

$$G_J(x) = \frac{1}{2} \sum_{j=1}^J (\log y_{j0} - \log y_{j1}) - \frac{1}{j_1} (x_j - y_j)^2 - \frac{1}{j_0} x_j^2 \quad 0 \tag{10}$$

Letting $r_j = r_{j0}^{1/2}$, then

$$r_j \sim N(0, 1) \text{ under } \Pi_0 \quad r_j \sim N(r_j^{-1}, r_j^{-1}) \text{ under } \Pi_1$$

$$G_J(X) = \frac{1}{2} \sum_{j=1}^J \{\log r_j - r_j(r_j - r_j)^2 + \frac{1}{r_j}\}$$

Under Gaussian assumptions, our Bayes classifier is a special case of the quadratic discriminant, a non-Bayes classifier because it uses two different sets of projections. The perfect classification properties for the functional quadratic discriminant were discussed in Delaigle & Hall (2013) in the context of truncated functional observations or fragments. We use the following auxiliary result.

Lemma 3. *Assume the predictors come from Gaussian processes. If $\sum_{j=1}^{\infty} \frac{1}{r_j} < \infty$ and $\sum_{j=1}^{\infty} (r_j - 1)^2 < \infty$, then $G_J(X)$ converges almost surely to a random variable as $J \rightarrow \infty$, in which case perfect classification does not occur. Otherwise perfect classification occurs.*

This lemma is similar to Theorem 1 of Delaigle & Hall (2013), but uses more transparent conditions and a proof technique based on the optimality property of Bayes classifiers which will be reused in the proof of Theorem 2. Lemma 3 states perfect classification occurs for Gaussian processes if and only if there are sufficient differences between the two groups in the mean or the covariance functions. This perfect classification phenomenon arises for the non-degenerate infinite dimensional case because we have infinitely many independent projection scores r_j for classification.

Proof of Lemma 3. Case 1: Assume $\sum_{j=1}^{\infty} (r_j - 1)^2 = \infty$ and that there exists a subsequence r_{j_l} of r_j that goes to ∞ or 0 as $l \rightarrow \infty$. Correspondingly take a subsequence $r_{j_l} \rightarrow \infty$ or $r_{j_l} \rightarrow 0$ for all $l = 1, 2, \dots$. Denoting the summand $(\log r_{j_0} - \log r_{j_1}) - \frac{1}{2} \left\{ \frac{(r_{j_0} - r_{j_1})^2}{r_{j_0} + r_{j_1}} \right\}$ of (10) as G_j , for any $l \leq J$ the misclassification rate $\Pr\{I\{G_J(X) \geq 0\} \neq Y\}$ is smaller than or equal to $\Pr\{I\{G_j \geq 0\} \neq Y\}$, since the former is the Bayes classifier using the first J projections, which minimizes the misclassification error among the class. Thus the misclassification rate of $G_J(X)$ is bounded above by that of the classifier $I\{\log r_j - r_j(r_j - r_j)^2 + \frac{1}{r_j} \geq 0\}$ for any $l \leq J$. Let \Pr_{Π_k} denote the conditional probability measure under group k . If there exists $r_{j_l} \rightarrow 0$,

$$\Pr_{\Pi_0} \{\log r_{j_l} - r_{j_l}(r_{j_l} - r_{j_l})^2 + \frac{1}{r_{j_l}} \geq 0\} \leq \Pr_{\Pi_0}(\log r_{j_l} + \frac{1}{r_{j_l}} \geq 0) \rightarrow 0$$

observing $\frac{1}{r_{j_l}} \sim \frac{1}{r_{j_l}}$ under Π_0 and $\log r_{j_l} \rightarrow -\infty$.

If there exist $r_{j_l} \rightarrow \infty$, then there exists a sequence $M \rightarrow \infty$ such that $(\log r_{j_l} + M)/r_{j_l} \rightarrow 0$. For any $l = 1, 2, \dots$,

$$\begin{aligned} \Pr_{\Pi_0} \{\log r_j - r_j(r_j - r_j)^2 + \frac{1}{r_j} \geq 0\} &\leq \Pr_{\Pi_0} \{\log r_j - r_j(r_j - r_j)^2 + M \geq 0\} + \Pr_{\Pi_0}(\frac{1}{r_j} \geq M) \\ &= \Pr_{\Pi_0} \left\{ \frac{(r_j - r_j)^2}{r_j} \leq \frac{\log r_j + M}{r_j} \right\} + o(1) \\ &= \Pr_{\Pi_0} \left\{ |r_j - r_j| \leq \frac{\log r_j + M}{r_j} \right\}^{1/2} + o(1) \end{aligned} \tag{11}$$

Plugging the sequence r_{j_i} for r_j into (11) we have $\{(\log r_j + M) / r_j\}^{1/2} \rightarrow 0$ as $j \rightarrow \infty$ and $M \rightarrow \infty$. Since r_{j_i} are standard normal and thus have uniformly bounded densities, (11) goes to zero and we have $\text{pr}_{\Pi_0} \{\log r_{j_i} - r_{j_i} (r_{j_i} - r_{j_i})^2 + \frac{2}{j_i} \geq 0\} \rightarrow 0$ as $j \rightarrow \infty$ and $M \rightarrow \infty$. Using similar arguments we can also prove $\text{pr}_{\Pi_1} \{\log r_{j_i} - r_{j_i} (r_{j_i} - r_{j_i})^2 + \frac{2}{j_i} \leq 0\} \rightarrow 0$ as $j \rightarrow \infty$. By Bayes' theorem $\text{pr}\{I\{ \sum_{j=1}^J G_j \geq 0\} \neq Y\} = \text{pr}(Y = 0) \text{pr}\{ \sum_{j=1}^J G_j \geq 0 \mid Y = 0\} + \text{pr}(Y = 1) \text{pr}\{ \sum_{j=1}^J G_j \leq 0 \mid Y = 1\} \rightarrow 0$ as $J \rightarrow \infty$. Therefore

$$\text{pr}\{I\{ \sum_{j=1}^J G_j(X) \geq 0\} \neq Y\} \leq \text{pr}\{I\{ \sum_{j=1}^J G_j \geq 0\} \neq Y\} \rightarrow 0 \quad J \rightarrow \infty$$

which means perfect classification occurs.

Case 2: Assume $\sum_{j=1}^{\infty} (r_j - 1)^2 = \infty$, and there exists M_1 and M_2 such that $0 < M_1 \leq r_j \leq M_2 < \infty$ for all $j \geq 1$. Letting E_{Π_k} and var_{Π_k} be the conditional expectation and variance under group k , respectively, we have

$$\begin{aligned} E_{\Pi_0} \{\log r_j - r_j (r_j - r_j)^2 + \frac{2}{j}\} &= \log r_j - (r_j - 1) - \frac{2}{j} r_j \\ E_{\Pi_1} \{\log r_j - r_j (r_j - r_j)^2 + \frac{2}{j}\} &= -\log r_j^{-1} + (r_j^{-1} - 1) + \frac{2}{j} \\ \text{var}_{\Pi_0} \{\log r_j - r_j (r_j - r_j)^2 + \frac{2}{j}\} &= 2(1 - r_j)^2 + 4 \frac{2}{j} r_j^2 \\ \text{var}_{\Pi_1} \{\log r_j - r_j (r_j - r_j)^2 + \frac{2}{j}\} &= 2(r_j^{-1} - 1)^2 + 4 \frac{2}{j} r_j^{-1} \end{aligned}$$

Then

$$\begin{aligned} \text{pr}_{\Pi_0} \left\{ \sum_{j=1}^J \{\log r_j - r_j (r_j - r_j)^2 + \frac{2}{j}\} \geq 0 \right\} &\leq \frac{\sum_{j=1}^J \{2(1 - r_j)^2 + 4 \frac{2}{j} r_j^2\}}{\left\{ -\sum_{j=1}^J (r_j - 1 - \log r_j + \frac{2}{j} r_j) \right\}^2} \\ &\leq \frac{\sum_{j=1}^J \{2(1 - r_j)^2 + 4M_2^2 \frac{2}{j}\}}{\left[\sum_{j=1}^J \left\{ \frac{1}{M_2} (r_j - 1)^2 + M_1 \frac{2}{j} \right\} \right]^2} \\ &= \frac{4M_2^2 M_1}{\sum_{j=1}^J \left\{ \frac{1}{M_2} (r_j - 1)^2 + M_1 \frac{2}{j} \right\}} \times \frac{\sum_{j=1}^J \{2(1 - r_j)^2 + 4M_2^2 \frac{2}{j}\}}{\sum_{j=1}^J \{4 \frac{M_2}{M_1} (r_j - 1)^2 + 4M_2^2 \frac{2}{j}\}} \\ &\leq \frac{4M_2^2 M_1}{\sum_{j=1}^J \left\{ \frac{1}{M_2} (r_j - 1)^2 + M_1 \frac{2}{j} \right\}} \rightarrow 0 \quad J \rightarrow \infty \end{aligned}$$

where Chebyshev's inequality is used for the first inequality, and Taylor expansion in the second inequality. Analogously the misclassification rate under Π_1 also can be proven to go to zero.

Case 3: Assume $\sum_{j=1}^{\infty} (r_j - 1)^2 < \infty$ and $\sum_{j=1}^{\infty} \frac{2}{j} = \infty$. The proof is essentially the same as in Case 2.

Case 4: Assume $\sum_{j=1}^{\infty} (r_j - 1)^2 < \infty$ and $\sum_{j=1}^{\infty} \frac{2}{j} < \infty$. Then the mean and variance of $\sum_{j=1}^J G_j(X)$ converges, so $\sum_{j=1}^J G_j(X)$ converges to a random variable under either population by Billingsley (1995). Therefore misclassification does not occur. \square

We can then proceed to prove Theorem 2, which does not assume Gaussianity.

Proof of Theorem 2. Case 1: Assume $\sum_{j=1}^{\infty} (r_j - 1)^2 = \infty$ and there exists a subsequence r_{j_i} of r_j that goes to 0 or ∞ as $j \rightarrow \infty$. By the optimality of Bayes classifiers, the Bayes classifier $I\{ \sum_{j=1}^J G_j(X) \geq 0\}$ using the first J components has smaller misclassification error than that of $I\{ \sum_{j=1}^J G_j \geq 0\}$, where G_j is the j th component in the summand of (3) in the main text, for all $J \leq J$. Since $I\{ \sum_{j=1}^J G_j \geq 0\}$ is the Bayes classifier using only

the j th projection, it has a smaller misclassification error than the non-Bayes classifier ($\mathbb{P}(G_j \geq 0)$), where $G_j = \log r_j - r_j(j - j)^2 + \frac{2}{j}$ is the j th summand in (10). Under Conditions A10–A11, we prove that the misclassification error converges to zeros by adopting the same argument as in Lemma 3 Case 1.

220 Case 2: Assume $\sum_{j=1}^{\infty} (r_j - 1)^2 = \infty$, and there exists M_1 and M_2 such that $0 < M_1 \leq r_j \leq M_2 < \infty$ for all $j \geq 1$. By some algebra,

$$E_{\Pi_0} \{ \log r_j - r_j(j - j)^2 + \frac{2}{j} \} = \log r_j - (r_j - 1) - \frac{2}{j} r_j$$

$$E_{\Pi_1} \{ \log r_j - r_j(j - j)^2 + \frac{2}{j} \} = -\log r_j^{-1} + (r_j^{-1} - 1) + \frac{2}{j}$$

$$\text{var}_{\Pi_0} \{ \log r_j - r_j(j - j)^2 + \frac{2}{j} \} \leq (2C_M - 1)(1 - r_j)^2 + 4(C_M + 1) \frac{2}{j} r_j^2$$

225
$$\text{var}_{\Pi_1} \{ \log r_j - r_j(j - j)^2 + \frac{2}{j} \} \leq (2C_M - 1)(r_j^{-1} - 1)^2 + 4(C_M + 1) \frac{2}{j} r_j^{-1}$$

The expectations are the same as in the Gaussian case because the first two moments of j do not depend on distributional assumptions. The inequalities in the variance calculation are due to $2ab \leq a^2 + b^2$ for all $a, b \in \mathbb{R}$. The same Chebyshev's inequality argument can be applied as for Theorem A1.

230 Case 3: Assume $\sum_{j=1}^{\infty} (r_j - 1)^2 < \infty$ and $\sum_{j=1}^{\infty} \frac{2}{j} = \infty$. The proof is essentially the same as that for Case 2. \square

References

- Billingsley, P.** (1995). *Probability and Measure*. New York: John Wiley & Sons Inc., 3rd ed.
- 235 **Bosq, D.** (2000). *Nonparametric Functional Data Analysis: Theory and Applications*. New York: Springer-Verlag.
- Delaigle, A. & Hall, P.** (2010). Defining probability density for a distribution of random functions. *Annals of Statistics* **38**, 1171–1193.
- Delaigle, A. & Hall, P.** (2013). Classification using censored functional data. *Journal of American Statistical Association* **108**, 1269–1283.
- 240 **Kong, D., Xue, K., Yao, F. & Zhang, H. H.** (2016). Partially functional linear regression in high dimensions. *Biometrika* **103**, 147–159.
- Nadaraya, E.** (1964). On estimating regression. *Journal of the Royal Statistical Society B* **9**, 141–142.
- Stone, C. J.** (1983). Optimal uniform rate of convergence for nonparametric estimators of a density function or its derivatives. In *Nonparametric Statistics*, M. H. Rizvi, J. S. Rustagi & D. Siegmund, eds. pp. 393–406.
- 245 **Watson, G. S.** (1964). Smooth regression analysis. *Annals of Statistics* **26**, 359–372.