

Bayes Risk Consistency of Nonparametric Classification Rules for Spike Trains Data

Mirosław Pawlak, Member, IEEE, Mateusz Pabian, Student Member, IEEE,
and Dominik Rzepka, Member, IEEE

Abstract

Spike trains data find a growing list of applications in computational neuroscience, imaging, streaming data and finance. Machine learning strategies for spike trains are based on various neural network and probabilistic models. The probabilistic approach is relying on parametric or nonparametric specifications of the underlying spike generation model. In this paper we consider the two-class statistical classification problem for a class of spike train data characterized by nonparametrically specified intensity functions. We derive the optimal Bayes rule and next form the plug-in nonparametric kernel classifier. Asymptotical properties of the rules are established including the limit with respect to the increasing recording time interval and the size of a training set. In particular the convergence of the kernel classifier to the Bayes rule is proved. The obtained results are supported by a finite sample simulation studies.

Index Terms

Bayes risk consistency, kernel classifiers, spike trains data, stochastic integrals

I. Introduction

EVENT driven systems are often encountered in science and engineering. In such systems data are represented by point processes that define arrival times of events. In computational neuroscience and machine learning this type of data are called spike trains [?], [?], [?]. In optical communication systems one observes a train of impulses (representing a point process) emitted by photon-sensitive detectors. Signal detection and estimation methods for such the so-called Poisson regime channels have been extensively examined in communication and information theory [?], [?], [?]. On the other hand, the mathematical theory of point processes has been extensively studied in the statistical and stochastic processes literature [?], [?]. However, the research on event type processes from the statistical classification theory [?] perspective has been initiated very recently [?], [?]. Probabilistic spiking neural networks have been introduced for supervised and unsupervised learning problems [?]. Various simulation results have been reported supporting their usefulness without, however, any accuracy studies and fundamental limits.

In this paper, we develop the Bayes strategy [?] for the spiking data supervised classification problem. This strategy can be applied to research problems where event occurrence is the primary information carrier [?], [?]. We consider a class of temporal spiking processes that are characterized by random intensity functions. This intensity function plays the central role in our theory as it describes the local rate of occurrence of spikes. For Fuchs processes (See Sect. II) we derive the optimal Bayes rule in terms of classification functions. Next, Sect. III is the limit behavior of the Bayes rule with respect to the increasing length of the observation interval is examined in Section Sect. IV the plug-in parametric kernel classification rule from multiple replicates of spiking processes is proposed. This is followed by the asymptotical optimality result, i.e., the convergence of the kernel rule to the Bayes rule. This result can be considered as the counterpart of the result in [?] concerning the classical plug-in nonparametric classification rules defined in the finite-dimensional Euclidean space. The spike train data are characterized by the variable length discrete dictors of events and their number avgives an observation interval. The main mathematical tool in our asymptotic analysis is the theory of the martingale decomposition for counting processes [?].

It is also worth mentioning that the asymptotic optimality does not hold if one observes the long single realization of the underlying spiking process. In fact, the intensity estimation problem for spiking processes does not fall into the classic large sample analysis distance between sample points framework as the point process is causal in time [?]. Hence, for a fixed observation interval one must increase the number of events. This can be achieved by either scaling the intensity function or by using the replicates of the spiking process. The former approach can be based on the multiplicative intensity model due to Aalen [?], whereas the latter one (used in this paper) is the standard machine learning strategy, where the replicates form the training set. In this case the resulting kernel estimate will be obtained by aggregating kernel estimates from single realizations. Our asymptotic results are supported by simulation studies.

M. Pawlak is with the Department of Electrical and Computer Engineering, University of Manitoba, R3T 5V6 Winnipeg, Canada, and with the Department of Measurement and Electronics, AGH University of Krakow, 30-059 Krakow, Poland. E-mail: Miroslaw.Pawlak@umanitoba.ca.

M. Pabian and D. Rzepka are with the Department of Electrical and Computer Engineering, University of Manitoba, R3T 5V6 Winnipeg, Canada, and with the Department of Measurement and Electronics, AGH University of Krakow, 30-059 Krakow, Poland. E-mail: Miroslaw.Pawlak@umanitoba.ca.

A preliminary version of this paper was presented at IEEE ICASSP 2023.

M. Pabian and D. Rzepka are with the Department of Measurement and Electronics, AGH University of Krakow, 30-059 Krakow, Poland.

A preliminary version of this paper was presented at IEEE ICASSP 2023.

by aggregating kernel estimates from single realizations. Our asymptotic results share support (P) of simulation studies presented in Section 3.3. (The) preliminary version of the results developed in this paper has been reported in a preprint by the symbol $\mathbb{E}(A)$ denotes the indicator function of the set. An \mathbb{E} shall be the notation $a(R)$ for the convergence in probability, which leads (note) to implies the convergence with probability one. Also α is the Lipschitz constant of $f(t)$, i.e., $|f(t_1) - f(t_2)| \leq M_f |t_1 - t_2|$ for sufficiently large T . Furthermore, $\overline{\lim}$, $\underline{\lim}$ denote the limit superior and inferior, respectively. Also, by M_f we will denote the Lipschitz constant of a function $f(t)$, i.e., $f(t)$ meets the Lipschitz condition if $|f(t_1) - f(t_2)| \leq M_f |t_1 - t_2|$ for all t_1, t_2 .

II. Bayes Classification Rule

A temporal spiking process $\{N(t), t \geq 0\}$ consists of a sequence of random times $\{t_i\}$ of isolated events in time such that $N(0) = 0$. The process $N(t)$ can be defined by the counting function $N(t) = \sum_i \mathbf{1}_{(t_i \leq t)}$ which is the number of temporal spiking events $\{N(t) \geq 0\}$ consists of a sequence of random times $\{t_i\}$ of isolated events in time such that $N(0) = 0$. This process $N(t)$ can be defined by the counting function $N(t)$ is negative $\sum_i \mathbf{1}_{(t_i \leq t)}$ which is the number of events in $[0, t]$. We assume that the process $N(t)$ is observed over the time window $[0, T]$ and is characterized by the non-random intensity function $\lambda(u)$ that is defined for all $u \geq 0$. This is a non-negative function that describes the local arriving rate of events such that $\mathbb{E}[N(T)]$ is $\int_0^T \lambda(u) du$ is the average number of events in $[0, T]$. Hence, the observed process $N(t)$ can be represented by the variable length vector $\mathbf{x} = [t_1, t_2, \dots, t_N, T]$, where $0 < t_1 < t_2 < \dots < t_N < T$. After the event X and the $N(T)$. Writing (x, π_1, N) we emphasize the fact that the data vector consists of two parts the data of current process $\{t_i\}$ develop a and N being the number of events in $[0, T]$ and the continuous part of the event X is the last class is discrete. Without a loss of generality we consider a two-class classification problem of this paper is the development of a rigorous classification methodology for the aforementioned class of spiking processes based on the Bayes theory of classification [1]. Without loss of generality we consider a two-class classification problem [see Section 3 for the generalization to the multi-class case] where class labels are denoted as ω_1, ω_2 with the prior probabilities π_1, π_2 , respectively. In order to form the optimal Bayes rule we recall the following known result [?] on the joint occurrence density of $\mathbf{x} = \prod_{i=1}^N \lambda(t_i) \exp \left(- \int_0^T \lambda(u) du \right)$ (1)

for $N = N(T) \geq 1$, whereas if $N = 0$ then $f(\mathbf{x}) = \prod_{i=1}^N \lambda(t_i) \exp \left(- \int_0^T \lambda(u) du \right)$. It is worth noting that (??) is the continuous-discrete distribution and by virtue of (??) the marginal density of the occurrence times $\{t_1, \dots, t_N\}$ for $N \in \{0, 1, \dots\}$ is given by $\sum_{n=0}^{\infty} f(t_1, \dots, t_N; N = n) = \exp \left(- \int_0^T \lambda(u) du \right)$ (2)

$$\sum_{n=0}^{\infty} f(t_1, \dots, t_N; N = n) = \exp \left(- \int_0^T \lambda(u) du \right) \sum_{n=1}^{\infty} \prod_{j=1}^n \lambda(t_j), \quad (2)$$

which is defined over the simplex regions $\mathcal{C}_n = \{(t_1, \dots, t_n) : 0 \leq t_1 \leq \dots \leq t_n \leq T\}$, $n = 1, 2, \dots$. The formula in (??) defines the proper density over $\{\mathcal{C}_n\}$, i.e., we have

which is defined over the simplex regions $\mathcal{C}_n = \{(t_1, \dots, t_n) : 0 \leq t_1 \leq \dots \leq t_n \leq T\}$, $n = 1, 2, \dots$. The formula in (??) defines the proper density over $\{\mathcal{C}_n\}$, i.e., we have

$$\exp \left(- \int_0^T \lambda(u) du \right) + \exp \left(- \int_0^T \lambda(u) du \right) \sum_{n=1}^{\infty} \int_{\mathcal{C}_n} \prod_{j=1}^n \lambda(t_j) dt_1 \cdots dt_n = 1. \quad (3)$$

In the context of the classification problem the class occurrence densities in (??) will be denoted $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ depending whether \mathbf{x} comes from class ω_1 (denoted as $\mathbf{x} \in \omega_1$) or if $\mathbf{x} \in \omega_2$, respectively. The corresponding class intensities in the context of the classification problem the class occurrence densities in (??) will be denoted $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ depending whether \mathbf{x} comes from class ω_1 (denoted as $\mathbf{x} \in \omega_1$) or if $\mathbf{x} \in \omega_2$, respectively. The corresponding class intensities are $\lambda_1(t), \lambda_2(t)$ being the non-negative functions defined on $[0, \infty)$. Then using (??), one can form the optimal Bayes rule $\psi_T^*: \mathbf{x} \in \omega_1 \text{ if } \prod_{i=1}^N \frac{\lambda_1(t_i)}{\lambda_2(t_i)} \exp \left(\int_0^T [\lambda_2(u) - \lambda_1(u)] du \right) \geq \frac{\pi_2}{\pi_1}$. (4)

assuming that $N \geq 1$ and $\exp \left(\int_0^T [\lambda_2(t_i) \lambda_1(u)] du \right) \geq \frac{\pi_2}{\pi_1}$ if $N = 0$. Clearly if the reverse inequality in (??) holds, then we classify \mathbf{x} to ω_2 . The log transform of (??) gives the alternative convenient form of the rule ψ_T^* , i.e., $\mathbf{x} \in \omega_1$ if assuming that $N \geq 1$ and $\exp \left(\int_0^T [\lambda_2(u) - \lambda_1(u)] du \right) \geq \frac{\pi_2}{\pi_1}$ if $N = 0$. Clearly, if the reverse inequality in (??) holds, then we classify \mathbf{x} to ω_2 . The log transform of (??) gives the alternative convenient form of the rule ψ_T^* , i.e., $\mathbf{x} \in \omega_1$ if

$$\sum_{i=1}^N \log \left(\frac{\lambda_1(t_i)}{\lambda_2(t_i)} \right) \geq \gamma, \quad (5)$$

where $\gamma = \int_0^T [\lambda_1(u) - \lambda_2(u)] du + \log\left(\frac{\pi_2}{\pi_1}\right)$. The rule in (??) can be usefully written in terms of the stochastic integral of the log-ratio $\log\left(\frac{\lambda_1(t)}{\lambda_2(t)}\right)$ with respect to the increments of the counting process $N(t)$, i.e., $\mathbf{X} \in \omega_1$ if

$$\int_0^T \log\left(\frac{\lambda_1(t)}{\lambda_2(t)}\right) dN(t) \geq \gamma. \quad (6)$$

Here $N(t)$ is the aforementioned counting process with the intensity function $\lambda(t)$, where

$$\lambda(t) = \begin{cases} \lambda_1(t) & \text{if } \mathbf{X} \in \omega_1 \\ \lambda_2(t) & \text{if } \mathbf{X} \in \omega_2 \end{cases}. \quad (7)$$

For our further considerations it is useful to represent the class intensity functions on $[0, T]$ in terms of the so-called intensity factor and shape function [?]. Thus, let $\lambda_1(t) = \tau_1 p_1(t)$, $\lambda_2(t) = \tau_2 p_2(t)$, where

$$\tau_i = \int_0^T \lambda_i(u) du, \quad p_i(t) = \lambda_i(t)/\tau_i, \quad i = 1, 2. \quad (8)$$

Clearly $p_1(t)$, $p_2(t)$ are well-defined probability density functions on $[0, T]$. The representation in (??) allows to represent the classification problem in terms of the class intensity factors and shape densities and employ information-theoretic divergence measures. Using (??), we can rewrite the rule in (??) as follows: $\mathbf{X} \in \omega_1$ if

$$\sum_{i=1}^N \log\left(\frac{p_1(t_i)}{p_2(t_i)}\right) \geq \eta, \quad (9)$$

where $\eta = \tau_1 - \tau_2 + N \log\left(\frac{\tau_2}{\tau_1}\right) + \log\left(\frac{\pi_2}{\pi_1}\right)$. The Bayes rule ψ_T^* in (??) will be written as $W_T(\mathbf{X}) \geq \eta_T$ emphasizing the fact that the vector \mathbf{X} is observed within the time window $[0, T]$.

It is worth noting that if $\lambda_1(t) = \lambda_1$ and $\lambda_2(t) = \lambda_2$, i.e., if we have the homogeneous spike train data then the Bayes rule takes the following form $\psi_T^*: \mathbf{X} \in \omega_1$

$$N \log\left(\frac{\lambda_1}{\lambda_2}\right) + T(\lambda_2 - \lambda_1) \geq \log\left(\frac{\pi_2}{\pi_1}\right), \quad (10)$$

provided that $N \geq 1$. In the case $N = 0$ this reads as $\lambda_2 - \lambda_1 \geq \frac{1}{T} \log\left(\frac{\pi_2}{\pi_1}\right)$. The risk associated with the rule $\psi_T^*(\mathbf{x})$ in (??) (or (??)) is defined as $R_T^* = \mathbf{P}(W_T(\mathbf{X}) \neq W_T(\mathbf{Y}))$ and it is referred to as Bayes risk. Here $\mathbf{Y} \in \mathcal{Y}_{\omega_1, \omega_2}$ is this the true class label of \mathbf{X} . For future studies we express Bayes risk in terms of the decision function $W_T(\mathbf{X})$, i.e., we rewrite

$$R_T^* = \mathbf{P}(W_T(\mathbf{X}) \geq \eta_T | \mathbf{X} \in \omega_2) \pi_2 + \mathbf{P}(W_T(\mathbf{X}) < \eta_T | \mathbf{X} \in \omega_1) \pi_1. \quad (11)$$

It is an important question to evaluate the Bayes risk. This includes various bounds on R_T^* and the behavior of R_T^* as a function of T . In Section ?? and ?? we present results concerning such issues.

The presented results rely on the following local decomposition (see Appendix A) of the increment $dN(t)$ of the point process $N(t)$. Hence, we have

$$dN(t) = \lambda(t) dt + dM(t), \quad (12)$$

where $\lambda(t)$ is the intensity function of $N(t)$, and $dM(t)$ is a zero mean process with uncorrelated but non-stationary increments. The formula in (??) can be viewed as the signal plus noise decomposition, where the noise process $M(t)$ is real-valued martingale, see Appendix A. Appendix A gives the pertinent results concerning the martingale decomposition of the underlying spiking process.

The decomposition in (??) allows us to express the classification rule in (??) (or its version in (??)) in the convenient stochastic integral form. In fact, by virtue of (??) and (??) we write the left-hand side of (??) as

$$\begin{aligned} \int_0^T \log\left(\frac{\lambda_1(t)}{\lambda_2(t)}\right) dN(t) &= \int_0^T \log\left(\frac{\lambda_1(t)}{\lambda_2(t)}\right) \lambda(t) dt \\ &\quad + \int_0^T \log\left(\frac{\lambda_1(t)}{\lambda_2(t)}\right) dM(t), \end{aligned} \quad (13)$$

where $\lambda(t)$ is given in (??) and $M(t)$ is the corresponding noise process defined in (??). The first term in (??) is the bias term of the optimal decision function, whereas the second one is the zero mean random variable contributing to the statistical variability of the rule. In Section ?? we show that the normalized version of this term converges exponentially fast to zero as $T \rightarrow \infty$ with probability one.

III. The Bayes Rule and Risk: Bounds and Asymptotic Behavior

A. The Bayes Decision Function

In this section we examine the optimal decision function derived in (??) or its alternative form in (??). Owing to the decomposition in (??) and using (??) we can arrive to the following equivalent form of the rule ψ_T^* in (??), $\mathbf{X} \in \omega_1$ if

$$U_T(\mathbf{X}) \geq \alpha_T + \log\left(\frac{\pi_2}{\pi_1}\right), \quad (14)$$

where

$$U_T(\mathbf{X}) = \int_0^T g(t) dM(t). \quad (15)$$

Here $g(t) = \log\left(\frac{\lambda_1(t)}{\lambda_2(t)}\right) = \log\left(\frac{p_1(t)}{p_2(t)}\right) + \log\left(\frac{\tau_1}{\tau_2}\right)$ and

$$\begin{aligned} \alpha_T &= \tau_1 - \tau_2 + \log\left(\frac{\tau_2}{\tau_1}\right) \int_0^T \lambda(t) dt \\ &+ \int_0^T \log\left(\frac{p_2(t)}{p_1(t)}\right) \lambda(t) dt, \end{aligned} \quad (16)$$

where $\lambda(t)$ is specified in (??).

It is worth noting that $U_T(\mathbf{X})$ in (??) represents the stochastic part of the Bayes rule. This takes the form of the stochastic integral with respect to the increments of the martingale process $M(t)$. It is known [?] that the martingale property is preserved under stochastic integration. Hence, since $\mathbb{E}[dM(t)] \equiv 0$ the process

$$\left\{ U_t(\mathbf{X}) \equiv \int_0^t g(u) dM(u), 0 \leq t \leq T \right\}$$

is a zero mean local martingale associated with the counting process $N(t)$; see Appendix A for further details. In addition, the integral in (??) is specified by the log-ratio $\log\left(\frac{\lambda_1(t)}{\lambda_2(t)}\right)$ and this is generally the unbounded function. To prevent this singularity it suffices to assume the class intensities $\lambda_1(t), \lambda_2(t)$ that are bounded away from zero. Moreover, intensity functions are commonly bounded. All these restrictions can be formalized by the following assumption that will be used in the paper. Hence, assume that there exist positive numbers δ and C such that

$$A1 : 0 \leq \delta \leq \lambda_i(t) \leq C, \quad i \equiv 1, 2, \text{ for all } t \geq 0. \quad (17)$$

We refer to [?], [?] for some weaker conditions for the existence of the aforementioned log-ratio.
We refer to [?], [?] for some weaker conditions for the existence of the aforementioned log-ratio.

In this section we present the preliminary results that characterize the Bayes rule specified by (??) and (??). This includes some bounds on the threshold α_T in (??) and the statistical properties of the stochastic term in (??). To do so, we recall that the Kullback-Leibler (KL) divergence [?] between densities $p(t)$ and $q(t)$ on $[0, T]$ is defined as follows

$$\mathbf{K}_T(p \parallel q) = \int_0^T \log\left(\frac{p(t)}{q(t)}\right) p(t) dt. \quad (18)$$

It is known that $\mathbf{K}_T(p \parallel q) \geq 0$ and $\mathbf{K}_T(p \parallel q) = 0$ if $p = q$.

The following lemma gives the upper and lower bounds for the threshold α_T in (??) in terms of the KL divergence between the class densities and the normalized square distance between the corresponding intensity factors. We will find these bounds useful in evaluating the Bayes risk.

Lemma 1. Let α_T be the threshold defined in (??). Then we have

(a) If $\mathbf{X} \in \omega_1$ then

(a) If $\mathbf{X} \in \omega_1$ then

$$\begin{aligned} &-\frac{(\tau_1 - \tau_2)^2}{(\tau_1 \tau_2)^2} - \tau_1 \mathbf{K}_T(p_1 \parallel p_2) \\ &\leq \alpha_T \leq -\tau_1 \mathbf{K}_T(p_1 \parallel p_2). \end{aligned} \quad (19)$$

(b) If $\mathbf{X} \in \omega_2$ then

(b) If $\mathbf{X} \in \omega_2$ then

$$\begin{aligned} \tau_2 \mathbf{K}_T(p_2 \parallel p_1) &\leq \alpha_T \\ \tau_2 \mathbf{K}_T(p_2 \parallel p_1) &\leq \frac{(\alpha_T - \tau_2)^2}{(\tau_1 \tau_2)^2} + \tau_2 \mathbf{K}_T(p_2 \parallel p_1), \\ &\leq \frac{(\alpha_T - \tau_2)^2}{(\tau_1 \tau_2)^2} + \tau_2 \mathbf{K}_T(p_2 \parallel p_1), \end{aligned} \quad (20)$$

where p_1, p_2, τ_1, τ_2 are defined in (??). The proof of Lemma ?? is given in Appendix B.

where p_1, p_2, τ_1, τ_2 are defined in (??). The proof of Lemma ?? is given in Appendix B.

As the KL divergence is non-negative, then Lemma ??(a) yields $\alpha_T \leq 0$ if $\mathbf{X} \in \omega_1$, whereas Lemma ??(b) gives $\alpha_T \geq 0$ for $\mathbf{X} \in \omega_2$. Also it is seen that α_T lies in the interval of the length $(\tau_1 - \tau_2)^2 / \tau_2$ and $(\tau_1 - \tau_2)^2 / \tau_1$ if $\mathbf{X} \in \omega_1$

As the KL divergence is non-negative, then Lemma ??(a) yields $\left\{ \int_0^T (\lambda_1(t) - \lambda_2(t)) dt \right\}^2$ whereas Lemma ??(b) gives $\sigma \geq 0$ for $\mathbf{X} \in \omega_2$. Also it is seen that σ lies in the interval of the length $(\tau_1 - \tau_2)^2 / \tau_2$ and $(\tau_1 - \tau_2)^2 / \tau_1$ if $\mathbf{X} \in \omega_1$ and $\mathbf{X} \in \omega_2$, respectively for \mathbf{X} is also worth noting that $(\tau_1 - \tau_2)^2 / \mathbf{X} \in \left\{ \int_0^T (\lambda_1(t) - \lambda_2(t)) dt \right\}^2$ represents the square of the difference of the average number of events on $[0, T]$ coming from classes ω_1 and ω_2 . As a result, if $\tau_1 = \tau_2$ then $\mathbf{X} \in \omega_1$ and $\mathbf{X} \in \omega_2$, respectively for \mathbf{X} is also worth noting that $(\tau_1 - \tau_2)^2 / \mathbf{X} \in \left\{ \int_0^T (\lambda_1(t) - \lambda_2(t)) dt \right\}^2$ represents the square of the difference of the average number of events on $[0, T]$ coming from classes ω_1 and ω_2 . As a result, if $\tau_1 = \tau_2$ then $\int_0^T K_T(p_1(t), p_2(t)) d\log \left(\frac{\lambda_1(t)}{\lambda_2(t)} \right)$ and the resp. $K_T(p_2(t), p_1(t))$ for $\mathbf{X} \in \omega_1 \cup \omega_2$ such that $\mathbb{E}[U_T(\mathbf{X})] = 0$. In the following lemma we evaluate the stochastic part $U_T(\mathbf{X})$ defined in (??). This is given in the form of the stochastic integral of the log-ratio $\log \left(\frac{\lambda_1(t)}{\lambda_2(t)} \right)$ with respect to the increments of $M(t)$ such that $\mathbb{E}[U_T(\mathbf{X})] = 0$. In the following Lemma 2. Let us consider the stochastic part $U_T(\mathbf{X})$ of the Bayes rule in (??). Then,

(a) If $\mathbf{X} \in \omega_1$ then

$$\bar{\text{Var}}[U_T(\mathbf{X})] = \tau_1 \int_0^T \left\{ \log \left(\frac{p_1(t)}{p_2(t)} \right) + \log \left(\frac{\tau_1}{\tau_2} \right) \right\}^2 p_1(t) dt. \quad (21)$$

(b) If $\mathbf{X} \in \omega_2$ then

$$\bar{\text{Var}}[U_T(\mathbf{X})] = \tau_2 \int_0^T \left\{ \log \left(\frac{p_1(t)}{p_2(t)} \right) + \log \left(\frac{\tau_1}{\tau_2} \right) \right\}^2 p_2(t) dt. \quad (21)$$

(b) If $\mathbf{X} \in \omega_2$ then

$$\bar{\text{Var}}[U_T(\mathbf{X})] = \tau_2 \int_0^T \left\{ \log \left(\frac{p_2(t)}{p_1(t)} \right) + \log \left(\frac{\tau_2}{\tau_1} \right) \right\}^2 p_2(t) dt. \quad (22)$$

The proof of Lemma ?? is given in Appendix B.

The formulas in Lemma ?? can be expressed in terms of the higher-order KL divergence between two class densities referred to as the KL variation [?]. Hence, let

The formulas in Lemma ?? can be expressed in terms of the higher-order KL divergence between two class densities referred to as the KL variation [?]. Hence, let $\mathbf{V}_T(p \parallel q) = \int_0^T \log^2 \left(\frac{p(t)}{q(t)} \right) p(t) dt$ be the KL variation between densities $p(t)$ and $q(t)$ on $[0, T]$. Note that $\mathbf{V}_T(p \parallel q) = 0$ if $p = q$. Moreover, the following result describes the relationship between $\mathbf{W}_T(p \parallel q)$ and the standard KL divergence in (??).

Lemma 3. For any pair of probability densities p, q on $[0, T]$ we have

be the KL variation between densities $p(t)$ and $q(t)$ on $[0, T]$. Note that $\mathbf{V}_T(p \parallel q) = 0$ if $p = q$. Moreover, the following result describes the relationship between $\mathbf{V}_T(p \parallel q)$ and the standard KL divergence in (??).

Lemma 3. For any pair of probability densities p, q on $[0, T]$ we have

the Cauchy-Schwarz inequality. Returning back to the formula in (??) we can obtain that

$$\mathbf{K}_T(p \parallel q) \leq \sqrt{\mathbf{V}_T(p \parallel q)}. \quad (24)$$

The bound in (??) results from the direct application of the Cauchy-Schwarz inequality. Returning back to the formula in (??) we can obtain that

$$\mathbf{Var}[U_T(\mathbf{X})] = \tau_1 \left\{ \mathbf{V}_T(p_1 \parallel p_2) + 2 \log \left(\frac{\tau_1}{\tau_2} \right) \mathbf{K}_T(p_1 \parallel p_2) + \log^2 \left(\frac{\tau_1}{\tau_2} \right) \right\}. \quad (25)$$

The analogous formula can be written for (??).

It is an interesting question to examine the behavior of $\mathbf{K}_T(p_1 \parallel p_2)$ and $\mathbf{V}_T(p_1 \parallel p_2)$ for an increasing value of the observation interval T . In particular, we wish to derive an analog of the law of large numbers, i.e., the limit behavior of $\mathbf{K}_T(p_1 \parallel p_2)$. To give an answer to such questions we need to put some condition on the growth of the assumed class of intensity functions. Hence, suppose that there exists positive number d such that

$\frac{1}{T} U_T(\mathbf{X}) = \frac{1}{T} \int_0^T g(t) dM(t)$

It is an interesting question to examine the behavior of the stochastic term $U_T(\mathbf{X})$ for an increasing value of the observation interval T . In particular, we wish to derive an analog of the law of large numbers, i.e., the limit behavior as $T \rightarrow \infty$, where $g(t)$ is the log-ratio $\log \left(\frac{\lambda_1(t)}{\lambda_2(t)} \right)$. To give an answer to such questions we need to put some condition on the growth of the assumed class of intensity functions. Hence, suppose that there exists positive number d such that

$A2 : \frac{1}{T} \int_0^T \lambda_i(u) du \rightarrow d, \quad i = 1, 2 \text{ as } T \rightarrow \infty$

as $T \rightarrow \infty$, where $g(t)$ is the log-ratio $\log \left(\frac{\lambda_1(t)}{\lambda_2(t)} \right)$. To give an answer to such questions we need to put some condition on the growth of the assumed class of intensity functions. Hence, suppose that there exists positive number d such that

Based on the assumption **A2** we wish to evaluate the limit behavior of $\mathbf{Var}[U_T(\mathbf{X})]$ as $T \rightarrow \infty$. It is clear that such

A2 it may not exist. Nevertheless, using the assumption **A1** we can find the upper and lower bounds for $\mathbf{Var}[U_T(\mathbf{X})]$.

In fact, recalling (??) and (??) we have that if $\mathbf{X} \in \omega_1$

The meaning of this condition is that the average number of events from the each class increases linearly with T . It is worth noting that for intensity functions that are integrable on $[0, \infty)$ the condition in (??) holds with $d = 0$.

Based on the assumption **A2** we wish to evaluate the limit behavior of $\mathbf{Var}[U_T(\mathbf{X})]$ as $T \rightarrow \infty$. It is clear that such

$$\leq \tau_1 \log^2 \left(\frac{C}{\delta} \right)$$

Smith may not exist. Nevertheless, using the assumption **A1** we can find the upper and lower bounds for $\text{Var}[U_T(\mathbf{X})]$. In fact, recalling (??) and (??) we have that if $\mathbf{X} \in \omega_1$

$$\begin{aligned}\text{Var}[U_T(\mathbf{X})] &\geq \tau_1 \left\{ K_T^2(p_1 \| p_2) \right. \\ \text{Var}[U_T(\mathbf{X})] &= \tau_1 \int \left\{ \log \left(\frac{\lambda_1(t)}{\lambda_2(t)} \right) \right\}^2 p_1(t) dt \\ &\quad + 2 \log \left(\frac{\tau_1}{\tau_2} \right) K_T(p_1 \| p_2) + \log^2 \left(\frac{\tau_1}{\tau_2} \right) \Big\} \\ &\leq \tau_1 \log^2 \left(\frac{C}{\delta} \right).\end{aligned}\quad (28)$$

The right-hand side of this inequality is equal to $\tau_1 \left(\int_0^T \log \left(\frac{\lambda_1(t)}{\lambda_2(t)} \right) p_1(t) dt \right)^2$ and by (??) this is not smaller than $\tau_1 \log^2 \left(\frac{C}{\delta} \right)$. Hence, if $\mathbf{X} \in \omega_1$ this gives the following bounds

$$\text{Var}[U_T(\mathbf{X})] \geq \tau_1 \left\{ K_T^2(p_1 \| p_2) \right. \\ \tau_1 \log^2 \left(\frac{C}{\delta} \right) \leq \text{Var}[U_T(\mathbf{X})] \leq \tau_1 \log^2 \left(\frac{C}{\delta} \right). \quad (29)$$

Analogously, we can show that if $\mathbf{X} \in \omega_2$, then

$$\begin{aligned}\text{Var}[U_T(\mathbf{X})] &\leq \tau_2 \left\{ K_T^2(p_1 \| p_2) \right. \\ &\quad + 2 \log \left(\frac{\tau_1}{\tau_2} \right) K_T(p_1 \| p_2) + \log^2 \left(\frac{\tau_1}{\tau_2} \right) \Big\}.\end{aligned}\quad (30)$$

The right-hand side of this inequality is equal to $\tau_2 \left(\int_0^T \log \left(\frac{\lambda_1(t)}{\lambda_2(t)} \right) p_2(t) dt \right)^2$ and by (??) this is not smaller than $\tau_2 \log^2 \left(\frac{C}{\delta} \right)$. Hence, if $\mathbf{X} \in \omega_2$ this gives the following bounds

The bounds in (??), (??) and the assumption in (??) lead to the following limit behavior of $\text{Var}[U_T(\mathbf{X})]$.

Lemma 4. Let the assumptions **A1**, **A2** hold. Then $\text{Var}[U_T(\mathbf{X})] \leq \mathbf{X} \log^2 \left(\frac{C}{\delta} \right)$ we have

$$\begin{aligned}\text{Analogously, we can show that } d \log^2 \left(\frac{\delta}{C} \right) \text{ then } \lim_{T \rightarrow \infty} \text{Var} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right] \\ \tau_2 \log^2 \left(\frac{\delta}{C} \right) \leq \lim_{T \rightarrow \infty} \text{Var} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right] \leq \tau_2 \log^2 \left(\frac{C}{\delta} \right).\end{aligned}\quad (31)$$

The bounds in (??), (??) and the assumption in (??) lead to the following limit behavior of $\text{Var}[U_T(\mathbf{X})]$. The question whether the inferior and superior limits in (??) are equal remains open. It should be noted that if (??) is in the form $\frac{1}{T} \int_0^T \lambda_i(u) du \rightarrow d_i$ $i = 1, 2$ then the result of Lemma ?? holds with d replaced by d_1 (if $\mathbf{X} \in \omega_1$) or d_2 (if $\mathbf{X} \in \omega_2$), respectively. To shed some light on the result in (??) let us consider the following simple example.

Example 1. Let us consider the classification problem with the intensity functions $\lambda_1(t)$ and $\lambda_2(t) = \mu \lambda_1(t)$ for some $\mu > 0$. Then we have $\tau_2 = \mu \tau_1$ and $p_2(t) = p_1(t)$. This implies that the condition in (??) reads as $\frac{1}{T} \int_0^T \lambda_1(u) du \rightarrow d$ and $\frac{1}{T} \int_0^T \lambda_2(u) du \rightarrow \mu d$. Then, a simple algebra gives the following analog of Lemma ??.

If $\mathbf{X} \in \omega_1$ then the question whether the inferior and superior limits in (??) are equal remains open. It should be noted that if (??) is in the form $\frac{1}{T} \int_0^T \lambda_i(u) du \rightarrow d_i$ $i = 1, 2$ then the result of Lemma ?? holds with d replaced by d_1 (if $\mathbf{X} \in \omega_1$) or d_2 (if $\mathbf{X} \in \omega_2$), respectively. To shed some light on the result in (??) let us consider the following simple example.

Example 1. Let us consider the classification problem with the intensity functions $\lambda_1(t)$ and $\lambda_2(t) = \mu \lambda_1(t)$ for some $\mu > 0$. Then we have $\tau_2 = \mu \tau_1$ and $p_2(t) = p_1(t)$. This implies that the condition in (??) reads as $\frac{1}{T} \int_0^T \lambda_1(u) du \rightarrow d$ and $\frac{1}{T} \int_0^T \lambda_2(u) du \rightarrow \mu d$. Then, a simple algebra gives the following analog of Lemma ??.

Note that the assumption **A1** is not required here. Also if $\mu = 1$ then the asymptotic constants are zero, i.e., this corresponds to the case $\lambda_1(t) = \lambda_2(t)$. Moreover, the asymptotic constants tend to infinity as $\mu \rightarrow \infty$.

An important consequence of Lemma ?? is the following weak law of large numbers for the average value of $U_T(\mathbf{X})$ defined in (??) whereas if $\mathbf{X} \in \omega_2$ then

Theorem 1. Let the conditions of Lemma ?? hold. Then $\lim_{T \rightarrow \infty} \text{Var} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right] = d \mu \log^2(\mu)$. Then for \mathbf{X} coming either from class ω_1 or class ω_2 we have

Note that the assumption **A1** is not required here. Also if $\mu = 1$ then the asymptotic constants are zero, i.e., this corresponds to the case $\lambda_1(t) = \lambda_2(t)$. Moreover, the asymptotic constants tend to infinity as $\mu \rightarrow \infty$.

An important consequence of Lemma ?? is the following weak law of large numbers for the average value of $U_T(\mathbf{X})$ defined in (??).

The proof of this fact is a direct application of Lemma ?? and the Chebyshev inequality. In fact, let us consider the case $\mathbf{X} \in \omega_1$. Then, for any $\epsilon > 0$ we have

Theorem 1. Let the conditions of Lemma ?? hold. Then for \mathbf{X} coming either from class ω_1 or class ω_2 we have

$$\frac{1}{T} U_T(\mathbf{X}) = \frac{1}{T} \int_0^T \log \left(\frac{\lambda_1(t)}{\lambda_2(t)} \right) dM(t)^2 \rightarrow 0. \quad (P) \quad (35)$$

The right-hand side of (??) is equal to $\text{Var} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right] / T \epsilon^2$, where due to (??) the limit superior of $\text{Var} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right]$ as $T \rightarrow \infty$.

The next goal is to strengthen the result of Theorem ?? by establishing the strong law of large numbers. This will be done directly from the exponential inequality for the average of $U_T(\mathbf{X})$ defined in (??). Our main tools here are exponential inequalities for martingales of counting processes established recently in [?], see also [?] for earlier results.

Hence, we employ the following adapted to our needs version of Theorem 5 in [?], see Appendix B for details. $\text{P} \left(\frac{|U_T(\mathbf{X})|}{T} \geq \epsilon \right) \leq \frac{\text{Var}[U_T(\mathbf{X})]}{T \epsilon^2} \quad (35)$

The right-hand side of (??) is equaling $\text{Var}_{\mathbb{P}} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right] / T \epsilon^2$, whereupon in (??) the limit superior of $\text{Var}_{\mathbb{P}} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right]$ is bounded by a finite constant. This confirms the claim of Theorem ??.

(a) Our next goal is to strengthen the result of Theorem ?? by establishing the strong law of large numbers. This will result directly from the exponential inequality for the average of $U_T(\mathbf{X})$ defined in (??). Our main tools here are exponential inequalities for martingales of counting processes established recently in [?], see also [?] for earlier results. Hence, we employ the following adapted to our needs version of Theorem 5 in [?], see Appendix B for details.

Lemma 5. Let $N(t)$ be the counting process allowing the decomposition in (??). Let $U_T = \int_0^T g(t)dM(t)$ be the stochastic integral of the real-valued function $g(t)$ with respect to the martingale $M(t)$ increments. Suppose that it is worth noting that this bound holds for any finite T .

(a) $|g(t)| \leq u_T$ for all $t \in [0, T]$.

(b) $\int_0^T g(t)\lambda(t)dt \leq u_T$, where by the assumption **A1** we have some finite constants. Then, for each $\frac{C}{\delta} > 0$ we have condition (a) in Lemma ?? is met with $u_T = \log(\frac{C}{\delta})$ for all $T > 0$. By virtue of the property (??) in Appendix A the integral in the condition (b) of Lemma ?? reads as

$$\mathbb{P}(|U_T| \geq \epsilon) \leq 2 \exp \left[-\frac{\epsilon^2}{1 + 2u_T + u_T \epsilon} \right]. \quad (36)$$

It is worth noting that this bound holds for any finite T .
 $\text{Var}[U_T(\mathbf{X})] = T \text{Var} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right] = T \theta_T,$ (37)

where due to (??) the limit superior of θ_T is bounded by a finite constant.

The preceding discussion gives the following exponential bound for the average value of $U_T(\mathbf{X})$ defined in (??). The bound is valid for any finite $T > 0$. By virtue of the property (??) in Appendix A the integral in the condition (b) of Lemma ?? reads as

$$\text{Var}[U_T(\mathbf{X})] = T \text{Var} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right] = T \theta_T, \quad (37)$$

$$\mathbb{P} \left(\frac{1}{T} |U_T(\mathbf{X})| \geq \epsilon \right) \leq 2 \exp \left[-T \frac{\epsilon^2}{2\theta_T + ue} \right], \quad (38)$$

where due to (??) the limit superior of θ_T is bounded by a finite constant. While the preceding discussion gives the following exponential bound for the average value of $U_T(\mathbf{X})$ defined in (??). The bound is valid for any finite $T > 0$.

The exponential bound in (??) and the Borel-Cantelli lemma yield the following strong version of Theorem ??.

Lemma 6. Suppose that the assumption **A1** holds. Then for \mathbf{X} coming either from class ω_1 or class ω_2 and every $\epsilon > 0$ we have $U_T(\mathbf{X})/T$ is a function of the continuous parameter T . Nevertheless, one can discretize ξ_T by finding a sequence of times T_n , such that $T_n \rightarrow \infty$ as $n \rightarrow \infty$ and then employ the standard Borel-Cantelli lemma. We refer to [?] for details for such discretization strategy.

where $u_T = 2 \log(\frac{C}{\delta})$ and the factor **A1** is defined in (??). Then for \mathbf{X} coming either from class ω_1 or class ω_2 we have

The exponential bound in (??) and the Borel-Cantelli lemma yield the following strong version of Theorem ??.

We should note, however, that the Borel-Cantelli lemma applies to a sequence of random variables, while the random variable $\xi_T = U_T(\mathbf{X})/T$ is a function of the continuous parameter T . Nevertheless, one can discretize ξ_T by finding a sequence of times T_n , such that $T_n \rightarrow \infty$ as $n \rightarrow \infty$ and then employ the standard Borel-Cantelli lemma. We refer to [?] for details for such discretization strategy.

B. The Bayes Risk

Theorem 2. Let the assumptions **A1** and **A2** hold. Then for \mathbf{X} coming either from class ω_1 or class ω_2 we have. In this section we wish to evaluate the Bayes risk defined in (??). Our analysis will employ the results obtained in Section ??.

Owing to (??) it suffices to consider the probability of misclassification $\mathbb{P}(W_T(\mathbf{X}) \geq \eta_T | \mathbf{X} \in \omega_2)$. The analysis of the probability $\mathbb{P}(W_T(\mathbf{X}) \leq \eta_T | \mathbf{X} \in \omega_1)$ is analogous. By virtue of (??) we can write

as $T \rightarrow \infty$.

$$\begin{aligned} & \mathbb{P}(W_T(\mathbf{X}) \geq \eta_T | \mathbf{X} \in \omega_2) \\ &= \mathbb{P} \left(U_T(\mathbf{X}) \geq \alpha_T + \log \left(\frac{\pi_2}{\pi_1} \right) | \mathbf{X} \in \omega_2 \right), \end{aligned} \quad (40)$$

B. The Bayes Risk

In this section we wish to evaluate the Bayes risk defined in (??). Our analysis will employ the results obtained in Section ??.

Owing to (??) it suffices to consider the probability of misclassification $\mathbb{P}(W_T(\mathbf{X}) \geq \eta_T | \mathbf{X} \in \omega_2)$. The analysis of the probability $\mathbb{P}(W_T(\mathbf{X}) \leq \eta_T | \mathbf{X} \in \omega_1)$ is analogous. By virtue of (??) we can write

The first result reveals that the Bayes risk tends to zero as $T \rightarrow \infty$ under the assumptions **A1** and **A2**. This is the direct consequence of the weak law of large numbers established in Theorem ??, see (??).

Hence, we have the following convergence result that also gives the upper bound for the Bayes risk.

where $U_T(\mathbf{X})$ is defined in (??) and α_T (under the fact that $\mathbf{X} \in \omega_2$) is given by

$$\alpha_T = \tau_1 - \tau_2 + \tau_2 \log \left(\frac{\pi_2}{\pi_1} \right) + \pi_2 K_T(p_2 \parallel p_1). \quad (41)$$

Furthermore,

$$R_T^* \leq (\pi_1 a_T + \pi_2 b_T) \frac{1}{T}, \quad (42)$$

The proof of Theorem ?? is deferred to Appendix B, where also the explicit expressions for a_T and b_T are given. Hence, we have the following convergence result that also gives the upper bound for the Bayes risk. The bound in (??) is obtained by utilizing only the second moment of the stochastic integral $U_T(\mathbf{X})$ in (??).

Theorem 3. Let the assumptions A_1 and A_2 hold. Then, we have

Remark 1. Hence under the assumptions **A1** and **A2** the Bayes risk tends to zero with the rate $1/T$. The proof of Theorem ?? reveals also the following form of $\text{tLR}^*_n \asymp n^{1/2} T^{-1/2}$.

Furthermore,

$$\mathbf{R}_T^* \leq \frac{1}{\theta} \left(\frac{\log(C/\delta)}{\log(\delta/C)} \right)^2 \frac{1}{T}, \quad (42)$$

Hence, for large T one can write $R_T^* \prec c_1 \frac{1}{T}$, for some finite constants a_T, b_T .

The proof of Theorem ?? is deferred to Appendix B, where also the explicit expressions for a_T and b_T are given. The bound in ?? is obtained by utilizing only the second moment of the stochastic integral $U_T(\mathbf{X})$ in ??.

Theorem 4. Let the assumptions A_1 and A_2 hold. Then, we have
Remark 1. Hence under the assumptions A_1 and A_2 the Bayes r

Remark 1. Hence under the assumptions **A1** and **A2** the Bayes risk tends to zero with the rate $1/T$. The proof of Theorem ?? reveals also the following form of the asymptotic constant $R_T \leq \pi_1 \exp\{-\lambda_1 T\} + \pi_2 \exp\{-\lambda_2 T\}$, (44)

for some finite constants A_T, B_T .

$$c_1 = \frac{1}{d} \left(\frac{\log(C/\delta)}{\log(\delta/C)} \right)^2. \quad (43)$$

The proof of Theorem ?? is deferred to Appendix B, where also the explicit expressions for A_T and B_T are presented. Hence, for large T one can write $R_T^* \prec c_1 \frac{1}{T}$.

Remark 2. The proof of Theorem ?? shows that using the exponential inequality for the martingale process the Bayes risk tends to zero with the exponential rate and the following asymptotic constant (??). Hence, we have the following result.

Theorem 4. Let the assumptions **A1** and **A2** hold. Then we have

where (δ, C) characterizes the assumption $\mathbf{A1}$, whereas T appears in the assumption $\mathbf{A1}$. Hence, for large T one (44) writes $R_T^* \prec \exp[-c_1 T]$. It is also worth noting that larger d in the assumption $\mathbf{A2}$ makes the bounds in (??) and (??) tighter. In fact, the constant c_1 in (??) decreases with d , whereas the constant c_2 in (??) increases with d .

The proof of Theorem ?? is deferred to Appendix B, where also the explicit expressions for A_T and B_T are presented. Example 2. Consider the classification problem discussed in Example ?? . Then, using the results in (??) and (??) and Remark 2. The proof of Theorem ?? shows that using the exponential inequality for the martingale process the Bayes risk tends to zero with the exponential rate and the following asymptotic constant

$$R_T \prec \pi_1 \exp[-c_1(\mu)dT] + \pi_2 \exp[-c_2(\mu)dT]. \quad (46)$$

The asymptotic constants $c_1(\mu)$, $c_2(\mu)$ can be written in the explicit form and they obey the following properties (45)

where (δ, C) characterizes the assumption **A1**, whereas d appears in the assumption **A1**. Hence, for large T one can write $R_T^* \prec \exp[-c_2 T]$. It is also worth noting that larger d in the assumption **A2** makes the bounds in (??) and (??) tighter. In fact, the constant c_1 in (??) decreases with d , whereas the constant c_2 in (??) increases with d .
 $\lim_{\mu \rightarrow \infty} c_1(\mu) = \lim_{\mu \rightarrow \infty} c_2(\mu) = \infty$.

Example 2. Consider the classification problem discussed in Example ???. Then, using the results in (??) and (??) and some algebra it can show the following counterpart of the result of Theorem (??). On the other hand, the latter limit exhibits that if $\mu \rightarrow \infty$ then $\mathbf{R}_T^* \rightarrow 0$. Again the assumption A1 is not needed here.

Remark 3. In [2] the following upper bound for the Bayes risk is given

The asymptotic constants $c_1(\mu)$, $c_2(\mu)$ can be written in the explicit form and they obey the following properties

$$\lim_{\mu \rightarrow 0} \frac{R_T^*}{c_1(\mu)} = \lim_{\mu \rightarrow 0} \frac{c_2(\mu)}{c_1(\mu)} = 0,$$

Remark 3. In [?] the following upper bound for the Bayes risk is given

where $\beta(T) = \int_0^T \left[\frac{1}{2}\lambda_1(u) + \frac{1}{2}\lambda_2(u) - \sqrt{\lambda_1(u)\lambda_2(u)} \right] du$ is a positive 1/2 factor. This is the classical Bhattacharyya bound (47) extended to the classification problem for point processes. The behavior of $\beta(T)$ under the condition A2 is an interesting open question. In the special case examined in Examples ??, ?? we can show that $\beta(T)$ behaves asymptotically as the results obtained in this paper under the multiplicative class intensity model in (??).

In the following example we give some numerical illustration of the aforementioned results.

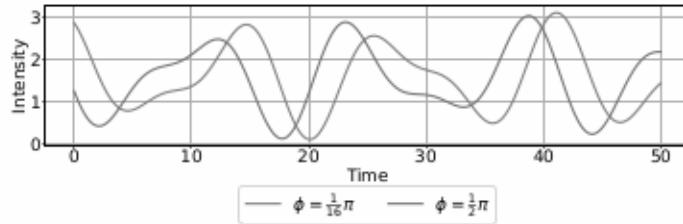


Figure 1. Intensity functions $\lambda_1(t) = \lambda\left(t; \frac{\pi}{16}\right)$ and $\lambda_2(t) = \lambda\left(t; \frac{\pi}{2}\right)$.

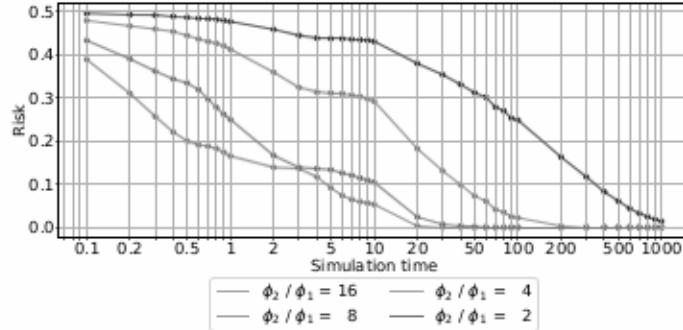


Figure 2. The Bayes risk R_T^* versus T for a two-class classification problem.

Remark 4.3. The convergence of the Bayes risk R_T^* to zero is determined by the condition in **A2**. This is due to the fact that the class intensity functions $\lambda_1(t), \lambda_2(t)$ grow with increasing T . If **A2** does not hold, e.g., if $\lambda_1(t), \lambda_2(t)$ are compactly supported then the convergence of R_T^* to zero is impossible. In this case in order to enforce the grow of $\lambda_1(t), \lambda_2(t)$ one could use the multiplicative model due to Aalen [?], i.e., we consider

$$\lambda_i(t) = d \gamma_i(t) \cos\left(\frac{\pi}{3\sqrt{2}}t + \frac{\pi}{4} + \phi\right) \quad (47)$$

whereas $\gamma_i(t)$ is a fixed function and d is a parameter that is allowed to grow. It is an interesting alternative to give the results obtained in this paper. Under the multiplicative class intensity model in (47) of class intensities defined in (??) is parametrized by ϕ , i.e., we set $\lambda_1(t) = \lambda_1(t; \phi_1)$ and $\lambda_2(t) = \lambda_1(t; \phi_2)$. The slowest decay of R_T^* is seen for very close intensities, i.e., when $\phi_2/\phi_1 = 2$ (in red), whereas the fast rate of convergence is observed for distant intensities, i.e., when $\phi_2/\phi_1 = 16$ (in blue). Nevertheless, since $\lambda(t; \phi)$ in (??) meets the assumptions **A1** and **A2** we can observe the exponential rate of convergence.

$$\lambda(t; \phi) = 1.6 + \cos\left(\frac{\pi}{4\sqrt{3}}t + \phi\right) \quad (48)$$

IV. Nonparametric Classification Rules

A. Plug-in Classifiers

In practice one does not know the true class intensities functions and must rely on some training data in order to form various choices of ϕ define $\lambda_1(t), \lambda_2(t)$. Figure ?? depicts $\lambda_1(t) = \lambda\left(t; \frac{\pi}{16}\right)$ and $\lambda_2(t) = \lambda\left(t; \frac{\pi}{2}\right)$: i.e., the classifier that a data-driven classification rule. In this paper we apply the plug-in strategy to design a classifier, i.e., the classifier that Figure ?? illustrates the fact that the Bayes risk tends to zero as T gets larger. The model of class intensities defined is the empirical counterpart of the optimal Bayes rule in (??) or equivalently in (??). We have already pointed out in (??) is parametrized by ϕ , i.e., we set $\lambda_1(t) = \lambda_1(t; \phi_1)$ and $\lambda_2(t) = \lambda_1(t; \phi_2)$. The slowest decay of R_T^* is seen that the single-sample based intensity function estimate cannot be consistent unless there is a certain mechanism that for very close intensities, i.e., when $\phi_2/\phi_1 = 2$ (in red), whereas the fast rate of convergence is observed for distant makes the intensity function increase, e.g., the multiplicative model in (??). In this paper we consider the intensity intensities, i.e., when $\phi_2/\phi_1 = 16$ (in blue). Nevertheless, since $\lambda(t; \phi)$ in (??) meets the assumptions **A1** and **A2** we model based on the increasing number of replicates of the class spiking processes. Hence, contrary to the results of can observe the exponential rate of convergence. Section ?? the observation interval $[0, T]$ is kept constant.

Hence, let $D_L = \{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_L, Y_L)\}$ be the learning sequence being a sample of L independent observations of the labeled spiking processes (\mathbf{X}, Y) . Here \mathbf{X}_j is the variable-length vector, i.e., $\mathbf{X}_j = [t_1^{[j]}, \dots, t_{N^{[j]}}^{[j]}, N^{[j]}]$ and

A. Plug-in Classifiers
In practice one does not know the true class intensities functions and must rely on some training data in order to form a data-driven classification rule. In this paper we apply the plug-in strategy to design a classifier, i.e., the classifier that

We wish to form the plug-in classification rule based on the optimal decision given in (??). This requires estimating is the empirical counterpart of the optimal Bayes rule in (??) or equivalently in (??). We have already pointed out the class intensity functions $\lambda_1(t), \lambda_2(t)$, or equivalently the shape densities $p_1(t), p_2(t)$ and the corresponding intensity that the single-sample based intensity function estimate cannot be consistent unless there is a certain mechanism that factors τ_1, τ_2 . It is known that the prior probabilities can be estimated by $\pi_1 = L_1/L$ and $\pi_2 = L_2/L$. In order to makes the intensity function increase, e.g., the multiplicative model in (??). In this paper we consider the intensity estimate $\{\tau_i, p_i(t), i = 1, 2\}$ one can begin with the use of the single sample \mathbf{X}_j . Note that $E[N^{[j]}|Y_j = \omega_i] = \tau_i$ and model based on the increasing number of replicates of the class spiking processes. Hence, contrary to the results of one can form the unbiased estimate of τ_i as $\tau_i^{[j]} = N^{[j]}$. However, $\text{Var}[N^{[j]}|Y_j = \omega_i] = \tau_i$ and this is an inconsistent Section ?? the observation interval $[0, T]$ is kept constant.

est. Hence, let $\mathbf{X} = \{\mathbf{x}_j\}_{j=1}^L$ result (\mathbf{X}, Y) from the L independent observations, the aggregated spiking processes $(\widehat{\tau}_i^{[j]})$, leads to (\mathbf{X}, Y) . Here \mathbf{x}_j is the variable length vector, i.e., $\mathbf{x}_j = [t_1, \dots, t_{N_j}; \mathbf{y}_j]$ and $Y_j \in \{\omega_1, \omega_2\}$, where $N^{[j]} = N^{[j]}(T)$. Hence, all data are measured in the fixed time window $[0, T]$. Let L_1, L_2 be the number of training data of classes ω_1 and ω_2 , respectively. $\sum_{j=1}^{N^{[j]}} \mathbf{1}(Y_j = \omega_i)$ (49)

We wish to form the plug-in classification rule based on the optimal decision given in (??). This requires estimating the class intensity functions $\lambda_1(t), \lambda_2(t)$, or equivalently the shape densities $p_1(t), p_2(t)$ and the corresponding intensity factors τ_1, τ_2 . It is known that the prior probabilities can be estimated by $\pi_1 = L_1/L$ and $\pi_2 = L_2/L$. In order to be a certain nonparametric estimate of $p_i(t)$ based on the single sample \mathbf{x}_j from the class ω_i . Then, the aggregated estimate $\{\widehat{\tau}_i, \widehat{p}_i(t)\}, i=1, 2$ one can begin with the use of the single sample \mathbf{x}_j . Note that $\mathbb{E}[N^{[j]}|Y_j = \omega_i] = \tau_i$ and estimate of $\widehat{p}_i(t)$ takes the following form

$$\text{one can form the unbiased estimate of } \tau_i \text{ as } \widehat{\tau}_i^{[j]} = N^{[j]}. \text{ However, } \text{Var}[N^{[j]}|Y_j = \omega_i] = \tau_i \text{ and this is an inconsistent estimate of } \tau_i. \text{ The latter fact results from the local Poisson behavior of the spiking process, see Appendix A. Nevertheless, the aggregation of } \{\widehat{\tau}_i^{[j]}\} \text{ leads to a consistent estimate of } \widehat{\tau}_i \text{ for the increased size of the training set. Hence, let}$$

Plugging (??) and (??) into (??) gives us the following empirical classification rule $\widehat{\psi}_{L,T}$: classify $\mathbf{X} = [t_1, \dots, t_N; N] \in \omega_1$ if

$$\widehat{\tau}_i = \frac{1}{L_i} \sum_{j=1}^{L_i} N^{[j]} \mathbf{1}(Y_j = \omega_i) \quad (49)$$

$$W_{L,T}(\mathbf{X}) \geq \widehat{\eta}_{L,T}, \quad (50)$$

be an estimate of τ_i , $i=1, 2$. In the analogous way we can deal with the problem of estimating $p_i(t)$. Let $\widehat{p}_i^{[j]}(t)$ where $W_{L,T}(\mathbf{X}) = \sum_{i=1}^N \log\left(\frac{\widehat{p}_i^{[j]}(t)}{p_2(t)}\right)$, $\widehat{\eta}_{L,T} = \widehat{\tau}_1 - \widehat{\tau}_2 + N \log\left(\frac{\widehat{p}_1^{[j]}(t)}{\widehat{p}_2^{[j]}(t)}\right) + \log\left(\frac{\tau_2}{\tau_1}\right)$. In Section ?? we propose a concrete kernel-type estimate of the shape densities.

In this section we present a general result on the convergence of the rule $\widehat{\psi}_{L,T}$ to the Bayes decision ψ_T^* . This result is in the spirit of the Bayes risk consistency theorem established in [?] in the context of the standard fixed dimension data sets. Let us first consider the pointwise behavior of the rule $\widehat{\psi}_{L,T}$ in (??). Hence, let $\mathbf{P}(\widehat{\psi}_{L,T}(\mathbf{x}) = \psi_T^*(\mathbf{x}))$ be the probability that the empirical rule makes the same decisions as the optimal Bayes rule for a fixed test vector \mathbf{x} . Our Plugging (??) and (??) into (??) gives us the following empirical classification rule $\widehat{\psi}_{L,T}$: classify $\mathbf{X} = [t_1, \dots, t_N; N] \in \omega_1$ if

Theorem 5. Suppose that for $i = 1, 2$ and $L \rightarrow \infty$ the following property holds

$$W_{L,T}(\mathbf{X}) \geq \widehat{\eta}_{L,T}, \quad (51)$$

where $\widehat{W}_{L,T}(\mathbf{X}) = \sum_{i=1}^N \log\left(\frac{\widehat{p}_i(t_i)}{p_2(t_i)}\right)$, $\widehat{\eta}_{L,T} = \widehat{\tau}_1 - \widehat{\tau}_2 + N \log\left(\frac{\widehat{p}_1(t_i)}{\widehat{p}_2(t_i)}\right) + \log\left(\frac{\tau_2}{\tau_1}\right)$. In Section ?? we propose a concrete kernel-type estimate of the shape densities.

In this section we present a general result of the convergence of the rule $\widehat{\psi}_{L,T}$ to the Bayes decision ψ_T^* . This result is in the spirit of the Bayes risk consistency theorem established in [?] in the context of the standard fixed dimension as $L \rightarrow \infty$. The proof of Theorem ?? is given in Appendix G. This result assures that $\widehat{\psi}_{L,T}$ converges to ψ_T^* as long as one can construct uniformly consistent estimates of $p_i(t)$, $i=1, 2$. Clearly, the uniform convergence of estimates of the class intensity functions $\lambda_i(t)$ also implies the local consistency result of Theorem ??.

The proof of Theorem ?? reveals also that the 0-1 distance between $\widehat{\psi}_{L,T}(\mathbf{x})$ and $\psi_T^*(\mathbf{x})$ tends to zero. Hence, we have

$$\widehat{p}_i(t) \xrightarrow{P} (p_i(t)(\mathbf{x})) \text{ uniformly on } P[0, T]. \quad (52)$$

Then, $\widehat{\psi}_{L,T}(\mathbf{x})$, where

$$\rho(\widehat{\psi}_{L,T}(\mathbf{x}), \psi_T^*(\mathbf{x})) = \mathbf{1}(\widehat{\psi}_{L,T}(\mathbf{x}) \neq \psi_T^*(\mathbf{x})).$$

The proof of Theorem ?? is given in Appendix G. This result assures that $\widehat{\psi}_{L,T}(\mathbf{x})$ tends to the optimal decision as long as one can construct uniformly consistent estimates of $p_i(t)$ in the proof of Theorem ??.

The proof of Theorem ?? reveals also that the 0-1 distance between $\widehat{\psi}_{L,T}(\mathbf{x})$ and $\psi_T^*(\mathbf{x})$ tends to zero. Hence, we have assumption **A1**. Let (??) hold. Then, we have

$$\begin{aligned} &\rho(\widehat{\psi}_{L,T}(\mathbf{x}), \psi_T^*(\mathbf{x})) \rightarrow 0 \quad (P) \\ &W_{L,T}(\mathbf{x}) \rightarrow W_T(\mathbf{x}) \quad (P), \end{aligned} \quad (53)$$

as $L \rightarrow \infty$, where
as $L \rightarrow \infty$.

$$\rho(\widehat{\psi}_{L,T}(\mathbf{x}), \psi_T^*(\mathbf{x})) = \mathbf{1}(\widehat{\psi}_{L,T}(\mathbf{x}) \neq \psi_T^*(\mathbf{x})).$$

The proof of Lemma ?? is postponed to Appendix C. The convergence in (??) is uniform with respect to \mathbf{x} .

The condition in (??) of Theorem ?? assures that the Bayesian function $\widehat{W}_{L,T}(\mathbf{x})$ in (??) tends to the optimal decision function $W_T(\mathbf{x})$ in (??). This is the convergence needed in the proof of Theorem ?? and is summarized in the following lemma. τ_i is weakly consistent. Hence, the preceding discussion gives the following consistency result

Lemma 7. Let the class intensities $\lambda_1(t), \lambda_2(t)$ be uniformly continuous on $[0, \infty)$ such that restricted to $[0, T]$ satisfy the assumption **A1**. Let (??) hold. Then, we have as $L \rightarrow \infty$.

The local consistency of $\widehat{\psi}_{L,T}$ leads to the $\widehat{W}_{L,T}(\mathbf{x}) \rightarrow W_T(\mathbf{x})$ characterized by the conditional risk. Hence, (54)

$\text{RS}(\widehat{\psi}_{L,T}(\mathbf{X}) \neq Y | \mathbf{D}_L) = \mathbb{E} \left[\mathbf{1} (\widehat{\psi}_{L,T}(\mathbf{X}) \neq Y) | \mathbf{D}_L \right]$ be the conditional risk associated with the rule $\widehat{\psi}_{L,T}$. Since $\mathbf{R}_T^* = \mathbb{E}[\mathbf{1}(\psi_T^*(\mathbf{X}) \neq Y)]$ then one can write.

The proof of Lemma ?? is postponed to Appendix C. The convergence in (??) is uniform with respect to \mathbf{x} . The classification rule $\widehat{\psi}_{L,T}$ in (??) is also characterized by the threshold value $\widehat{\eta}_{L,T}$. Note that $\pi_1 = L_1/L$, $\pi_2 = L_2/L$ are weakly consistent estimates of the prior probabilities π_1, π_2 . Also the aggregated estimate τ_i in (??) of the intensity factor τ_i is weakly consistent. Hence, the preceding discussion gives the following consistency result

Recalling the definition of the distance in (??) the above is bounded by (55)

$$\text{as } L \rightarrow \infty. \quad \mathbb{E} \left[\rho(\widehat{\psi}_{L,T}(\mathbf{X}), \psi_T^*(\mathbf{X})) | \mathbf{D}_L \right].$$

The local consistency of $\widehat{\psi}_{L,T}$ leads to the global convergence characterized by the conditional risk. Hence, let Owing to (??) and Lebesgue's dominated convergence theorem we obtain the main result of this section. $\widehat{\mathbf{P}}(\psi_{L,T}(\mathbf{X}) \neq Y | \mathbf{D}_L) = \mathbb{E}[\mathbf{1}(\psi_{L,T}(\mathbf{X}) \neq Y) | \mathbf{D}_L]$ be the conditional risk associated with the rule $\widehat{\psi}_{L,T}$. Since Theorem ?? Consider the class of plug-in classifiers defined in (??). Suppose that the conditions of Theorem ?? hold. Then, we have the following Bayes risk consistency result

$$0 \leq \mathbf{R}_{L,T} - \mathbf{R}_T^* \\ = \mathbb{E} \left[\mathbf{1} (\widehat{\psi}_{L,T}(\mathbf{X}) \neq Y) - \mathbf{1} (\psi_T^*(\mathbf{X}) \neq Y) | \mathbf{D}_L \right]. \quad (56)$$

as $L \rightarrow \infty$.

Recalling the definition of the distance in (??) the above is bounded by

B. Kernel Classifiers

$$\mathbb{E} \left[\rho(\widehat{\psi}_{L,T}(\mathbf{X}), \psi_T^*(\mathbf{X})) | \mathbf{D}_L \right].$$

It is known [?], [?] that the intensity function of a point process can be efficiently estimated by a class of kernel methods [?], [?]. In particular, the standard single sample kernel estimate of $\lambda_i(t)$ takes the form

Theorem 6. Consider the class of plug-in classifiers defined in (??). Suppose that the conditions of Theorem ?? hold. Then, we have the following Bayes risk consistency result

$$\lambda_i^{[j]}(t) = \sum K_h(t - t_l^{[j]}), \quad (57)$$

$$\mathbf{R}_{L,T} \rightarrow \mathbf{R}_T^*(P) \quad (56)$$

where the sample $\mathbf{X}_j = [t_1^{[j]}, \dots, t_{N^{[j]}}^{[j]}; N^{[j]}]$ comes from the class ω_i .

Here $K_h(t) = h^{-1}K(t/h)$, where the kernel $K(t)$ is assumed to be a compactly supported on $[-1, 1]$, symmetric probability classifier function. For instance, one can choose the so-called Epanechnikov kernel

It is known [?], [?] that the intensity function of a point process can be efficiently estimated by a class of kernel methods [?], [?]. In particular, the standard single sample kernel estimate of $\lambda_i(t)$ takes the form

The crucial tuning parameter h is called the bandwidth as it controls the level of smoothing via the scaled kernel $K_h(t)$.

The parameter τ_i can be estimated (from a single sample) by $\widehat{\tau}_i^{[j]} = N^{[j]}$. Therefore (??) yields the following estimate of the shape density

where the sample $\mathbf{X}_j = [t_1^{[j]}, \dots, t_{N^{[j]}}^{[j]}; N^{[j]}]$ comes from the class ω_i .

Here $K_h(t) = h^{-1}K(t/h)$, where the kernel $K(t)$ is assumed to be a compactly supported on $[-1, 1]$, symmetric probability density function. For instance, one can choose the so-called Epanechnikov kernel

As we have already pointed in Section ?? the estimates $K(t) = \frac{3}{4}(1-t^2)\mathbf{1}(|t| \leq 1)$ cannot be consistent by merely increasing T . To overcome this problem one can utilize the observed multiple training vectors and aggregate the single-sample estimates. The crucial tuning parameter h is called the bandwidth kernel as it controls the level of smoothing via the scaled kernel $K_h(t)$.

The parameter τ_i can be estimated (from a single sample) by $\widehat{\tau}_i^{[j]} = N^{[j]}$. Therefore (??) yields the following estimate of the shape density

$$p_i(t) = \frac{1}{L_i} \sum_{j=1}^L p_i^{[j]}(t) \mathbf{1}(Y_j = \omega_i). \quad (58)$$

Moreover, the aggregated estimate $\widehat{\tau}_i$ of τ_i is defined in (??) and plugging $\widehat{p}_i(t)$ and $\widehat{\tau}_i$, $i = 1, 2$ into (??) we obtain the kernel classification rule. The aggregated kernel estimate $\widehat{\lambda}_i(t)$ of $\lambda_i(t)$ is defined in the analogous way, see (??).

Theorem ?? and Theorem ?? reveal that the sufficient condition for the Bayes risk consistency is the convergence property in (??). Note that the statistical behavior of $p_i(t)$ and $\lambda_i(t)$ is the same and therefore we can verify the requirement in (??) for the kernel intensity estimate. Hence, with some abuse of the notation let $\{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_L\}$ be the data set from the fixed class (ω_1 or ω_2) of the counting process $N(t)$ characterized by the class intensity function $\lambda(t)$. Thus, one observes the L copies $\{N^{[j]}(t)\}$ of the counting process $N(t)$, where $N^{[j]}(t)$ is represented by the feature vector $\mathbf{X}_j = [t_1^{[j]}, \dots, t_{N^{[j]}}^{[j]}; N^{[j]}]$ with $N^{[j]}(t) = \sum_{i=1}^{N^{[j]}} \mathbf{1}(Y_i = \omega_i)$. The local martingale decomposition in (??) for $N^{[j]}(t)$ reads

Moreover, the aggregated estimate $\widehat{\tau}_i$ of τ_i is defined in (??). Plugging $\widehat{p}_i(t)$ and $\widehat{\tau}_i$, $i = 1, 2$ into (??) we obtain the kernel classification rule. The aggregated kernel estimate $\widehat{\lambda}_i(t)$ of $\lambda_i(t)$ is defined in the analogous way, see (??).

This gives the analogous decomposition for the aggregated counting process i.e., we have Theorem ?? and Theorem ?? reveal that the sufficient condition for the Bayes risk consistency is the convergence property in (??). Note that the statistical behavior of $\widehat{p}_i(t)$ and $\widehat{\lambda}_i(t)$ is the same and therefore we can verify (??)

requirement in (??) for the kernel intensity estimate. Hence, with some abuse of the notation let $\{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_L\}$ be the data set from the fixed class (ω_1 or ω_2) of the counting process $N(t)$ characterized by the class intensity function $\lambda(t)$. Thus, one observes the L copies $\{N^{[j]}(t)\}_{j=1}^L$ of the counting process $N(t)$, where $N^{[j]}(t)$ is represented by the feature vector $\mathbf{X}_j = [t_1^{[j]}, \dots, t_{N^{[j]}}^{[j]}; N^{[j]}]$ with $N^{[j]} = N^{[j]}(T)$. The local martingale decomposition in (??) for $N^{[j]}(t)$ reads

$$dN^{[j]}(t) = d\bar{M}_L(t) + \frac{1}{L} \sum_{j=1}^L dM^{[j]}(t), \quad j = 1, \dots, L.$$

This gives the analogous decomposition for the aggregated counting process, i.e., we have

It is important to note that the aggregated residual process $d\bar{M}_L(t)$ meets all the properties listed in Appendix A. Hence, $\mathbb{E}[d\bar{N}_L(t)] = 0$ and the properties in (??) and (??) follows

where

$$\begin{aligned} \text{Var}[d\bar{M}_L(t)] &= \frac{1}{L} \lambda(t) dt, \\ \text{Var}\left[\int_0^T g(u) d\bar{N}_L(u)\right] &= \frac{1}{L} \sum_{j=1}^L \frac{1}{L} \int_0^T g^2(u) \lambda^2(u) du. \end{aligned} \quad (60)$$

The single-sample kernel estimate of $\lambda(t)$ is $\hat{\lambda}(t) = \frac{1}{L} \sum_{j=1}^L \hat{M}^{[j]}(t)$, whereas its aggregated version takes the form

It is important to note that the aggregated residual process $d\bar{M}_L(t)$ meets all the properties listed in Appendix A. Hence, $\mathbb{E}[d\bar{M}_L(t)] = 0$ and the properties in (??) and (??) are as follows

This due to (??) can be written in the convenient stochastic integral form

$$\begin{aligned} \text{Var}[d\bar{M}_L(t)] &= \frac{1}{L} \lambda(t) dt, \\ \text{Var}\left[\int_0^T \frac{\hat{\lambda}(t)}{g(u)} d\bar{N}_L(u)\right] &= \frac{1}{L} \int_0^T g^2(u) \lambda(u) du. \end{aligned} \quad (60)$$

Employing this identity along with (??) and the aforementioned properties of $d\bar{M}_L(t)$ (see (??)) yield the following The single-sample kernel estimate of $\lambda(t)$ as in (??), whereas its aggregated version takes the form

$$\mathbb{E}[\hat{\lambda}(t)] \hat{\lambda}(t) = \int_0^T \frac{1}{L} \sum_{j=1}^L \frac{K\left(\frac{t-s}{h}\right)}{h} \lambda(s) ds, \quad (63)$$

This due to (??) can be written in the convenient stochastic integral form $\frac{1}{Lh} \int_0^T \frac{1}{h} K\left(\frac{t-s}{h}\right) \lambda(s) ds$.

These formulas and the standard analysis developed in the context of kernel estimates [?], [?] reveal that if

Employing this identity along with (??) and the aforementioned properties of $d\bar{M}_L(t)$ (see (??)) yield the following identities for the bias and the variance of $\hat{\lambda}(t)$

$$\mathbb{E}[\hat{\lambda}(t)] = \int_0^T \frac{1}{Lh} \int_0^T \frac{1}{h} K\left(\frac{t-s}{h}\right) \lambda(s) ds, \quad (63)$$

at $t \in (0, T)$ where $\lambda(t)$ is continuous. This is the pointwise convergence that holds at interior points of $[0, T]$. It is known [?], [?] that the convergence fails at the boundary points near $t = 0$, $t = T$. This enforces us to confine the required uniform convergence to the interval $[\epsilon, T - \epsilon]$ for arbitrarily small $\epsilon > 0$. Yet another option is to introduce the boundary modified kernels [?], [?] that are able to restore the convergence property at the boundary points. These formulas and the standard analysis developed in the context of kernel estimates [?], [?] reveal that if

Lemma 8. Let $\lambda(t)$ be Lipschitz continuous on $(L, \infty) \setminus 0$ and the kernel function $K(t)$ be Lipschitz continuous on $[-1, 1]$. Suppose that

$$h(L) \rightarrow 0 \text{ and } Lh^3(L) \rightarrow \infty \text{ as } L \rightarrow \infty. \quad (65)$$

Then for arbitrarily small $\epsilon > 0$ at $t \in (0, T)$ where $\lambda(t)$ is continuous. This is the pointwise convergence that holds at interior points of $[0, T]$. It is known [?], [?] that the convergence fails at the boundary points near $t = 0, \infty, t = T$. This enforces us to confine the required uniform convergence to the interval $[\epsilon, T - \epsilon]$ for arbitrarily small $\epsilon > 0$. Yet another option is to introduce the boundary modified kernels [?], [?] that are able to restore the convergence property at the boundary points. The following lemma gives the sufficient conditions for the uniform convergence where one needs that $Lh(L) \rightarrow 0$ in order to prove that (??) can be replaced by the weaker condition $Lh(L)/\log(L) \rightarrow \infty$. This is the case for the uniform convergence Lemma 8. Let $\lambda(t)$ be Lipschitz continuous on $[0, \infty)$. Let the kernel function $K(t)$ be Lipschitz continuous on $[-1, 1]$. Suppose that Our proof is based on more elementary techniques. The proof of Lemma ?? is given in Appendix D. The result of Lemma ?? applies directly to the shape densities and by using Theorem ?? and Theorem ?? we can formulate the The following arbitrarily small step result for the kernel classifier.

$$\sup_{t \in [\epsilon, T - \epsilon]} |\hat{\lambda}(t) - \lambda(t)| \rightarrow 0 \text{ (P) as } L \rightarrow \infty. \quad (67)$$

Theorem 7 Let the class intensities $\lambda(t)$, $t \in (0, T)$ satisfy the conditions of Lemma ?? and the condition $Lh^2(L) \rightarrow \infty$. Then the kernel density estimate $\hat{R}_{L,T} \rightarrow R_T^*(P)$ holds under the stronger condition $Lh^2(L) \log(L) \rightarrow \infty$. This is the case for the uniform convergence of the kernel density estimate where advanced tools from the empirical processes theory have been utilized [?], [?]. Our proof is based on more elementary techniques. The proof of Lemma ?? is given in Appendix D. The result of Lemma ?? applies directly to the shape densities and by using Theorem ?? and Theorem ?? we can formulate the following Bayes risk consistency result for the kernel classifier. Nevertheless, the issue of the rate of convergence would also be essential. This question is left for further research.

Theorem 7 Let the class intensities $\lambda(t)$, $t \in (0, T)$ satisfy the conditions of Lemma ?? and the condition $Lh^2(L) \rightarrow \infty$. If (??) holds, then the kernel classification rule is Bayes Risk consistent analysis we have to the expression in (??) and (??) shows that if $\lambda(t)$ has two continuous derivatives for $t \in (0, T)$ then

$$\mathbb{E} [\hat{\lambda}(t)] = \lambda(t) + \mathcal{O}(h^2)$$

as $L \rightarrow \infty$.

and

The convergence in Theorem ?? is an important property of the kernel classifier. Nevertheless, the issue of the rate of convergence would also be essential. This question is left for further research.

The selection of the bandwidth h is the most important issue in determining the finite sample accuracy of the kernel classification rule. The standard analysis applied to the expression in (??) and (??) shows that if $\lambda(t)$ has two continuous derivatives for $t \in (0, T)$ then

$$\begin{aligned} \mathbb{E} [\hat{\lambda}(t) - \lambda(t)]^2 &= \mathcal{O}\left(\frac{1}{Lh}\right) + \mathcal{O}(h^4), \quad t \in (0, T). \\ \mathbb{E} [\hat{\lambda}(t)] &= \lambda(t) + \mathcal{O}(h^2) \end{aligned}$$

and minimum of the error yields the asymptotically optimal choice of the bandwidth, i.e., $h^* = cL^{-1/5}$ for some positive constant c . This is the asymptotically optimal choice of h that optimizes the kernel intensity estimate. An optimal bandwidth for the kernel classifier may be quite different it is seen from the restriction in (??). See also [?] for the general theory of plug-in nonparametric classifiers.

This leads to the following asymptotical formula for the mean squared error. In practical applications one can specify the bandwidth using some resampling techniques like cross-validation [?], [?]. In our experimental studies we choose separate bandwidth for each class. This is done by finding the maximum of the cross-validated log-likelihood of the kernel estimate of the shape densities. Hence, let $\hat{p}_i(t; h)$ be the kernel estimate in (??) specified by the bandwidth h . Then, the likelihood function of $\hat{p}_i(t; h)$ specified by test data is given by the minimum of the error yields the asymptotically optimal choice of the bandwidth, i.e., $h^* = cL^{-1/5}$ for some positive constant c . This is the asymptotically optimal choice of h that optimizes the kernel intensity estimate. An optimal bandwidth for the kernel classifier may be quite different it is seen from the restriction in (??). See also [?] for the general theory of plug-in nonparametric classifiers.

In practical applications one can specify the bandwidth using some resampling techniques like cross-validation [?], [?]. where $t_r^{[l]}$ represents the r -th observation of the l -th test sample. We use the test sample of size q (per class). Also In our experimental studies we choose separate bandwidth for each class. This is done by finding the maximum of the $\tilde{p}_i(t; h)$ is the version of $\hat{p}_i(t; h)$ in (??) determined from the $L - q$ size training set. Then, the bandwidth is selected cross-validated log-likelihood of the kernel estimate of the shape densities. Hence, let $p_i(t; h)$ be the kernel estimate as the one that maximizes $\text{CV}(h)$ in (??). This is equivalent to the following choice in (??) specified by the bandwidth h . Then, the likelihood function of $p_i(t; h)$ specified by test data is given by

$$\hat{h}_i = \underset{h}{\text{CV}} \max \sum_{l=1}^p \prod_{r=1}^{N^{[l]}} \log(p_i(t_r^{[l]}; h)). \quad (68)$$

where $t_r^{[l]}$ represents the r -th observation of the l -th test sample. We use the test sample of size q (per class). Also $\tilde{p}_i(t; h)$ is the version of $\hat{p}_i(t; h)$ in (??) determined from the $L - q$ size training set. Then, the bandwidth is selected as the one that maximizes $\text{CV}(h)$ in (??). This is equivalent to the following choice in order to gain insight into the behavior of $R_{L,T}$ with respect to the training set size L and the observation window size T .

In all experiments the kernel classifier \hat{h}_i is given by $\sum_i \sum_j \log(\hat{p}_i(t_r^{[l]}; h))$ with the estimated $\hat{p}_i(t; h)$ specified by (??) and (??), respectively. The Gaussian kernel is employed, whereas the bandwidth is selected by the log-likelihood method in (??). When selecting the bandwidth, we consider a grid of ten evenly logarithmically spaced points $h \in \{10^{-1}, 10^1\}$. Additionally, we employ a 5-fold cross validation in order to avoid biasing the selected bandwidth \hat{h}_i with the test data. In order to assess the proposed methodology, we conduct a simulation study. We limit the scope of our experiments to time-dependent intensity functions defined in (??), and use these in simulations in order to gain insight into the behavior of $R_{L,T}$ with respect to the training set size L and the observations window size T . Unless noted otherwise, we perform experiments the kernel classifier is given by (??) with the estimated $\hat{p}_i(t; h)$ specified by (??) and (??), respectively. The Gaussian kernel is employed, whereas the bandwidth is selected by the log-likelihood method in (??). When plotting the bandwidth, we consider a grid of ten evenly logarithmically spaced points for $10^{-1}, 10^1$. Additionally, we employ a 5-fold cross-validation in order to avoid biasing the selected bandwidth \hat{h}_i with the test data. Finally, we denote by $R_{L,T}^*$ the empirical risk evaluated using a average over the simulation runs with a testing set size of 10^4 simulation space slice in subsequent analysis, i.e., with the value of T fixed.

We shall focus on the simulation results obtained for the intensity function specified by (??). Unless noted otherwise, we refer to the intensity function pair parameterized by $\theta_1 = \pi/4$ and $\theta_2 = \pi/4$. We observe

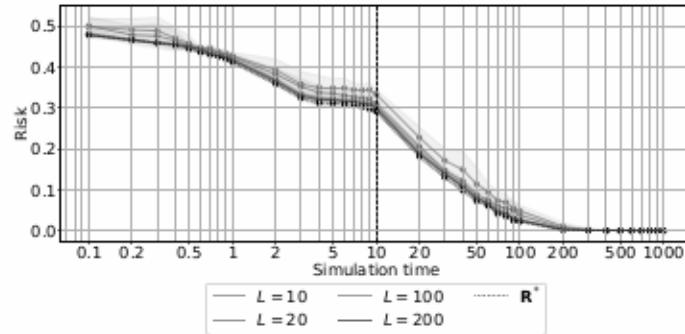


Figure 3. The average risk $E[R_{L,T}]$ versus T for different values of L . The vertical dashed line at $T = 10$ denotes the simulation space slice presented in Figure ??-b-??.

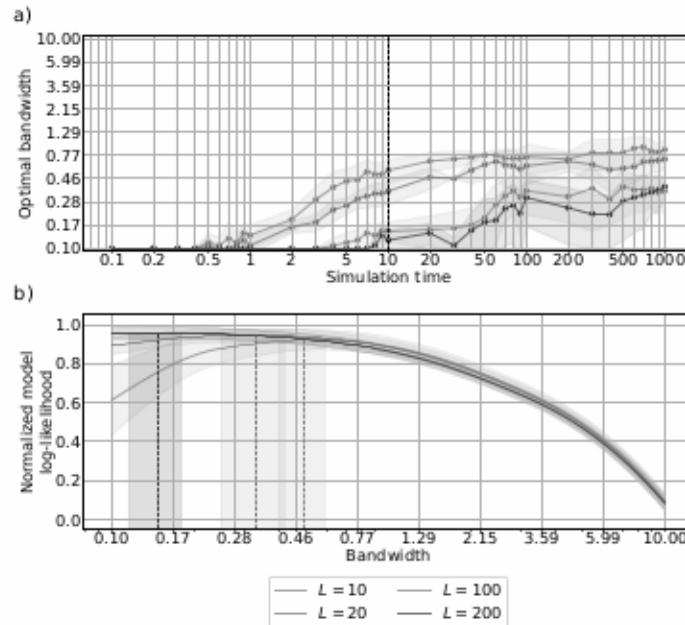


Figure 4. a) The average empirical bandwidth $E[\hat{h}_i]$ versus T for different values of L . The vertical dashed line at $T = 10$ denotes the simulation space slice presented in the lower subfigure. b) The average normalized model log-likelihood on test data versus h for different values of L . The vertical dashed lines denote function maxima. Note that the curves for $L = 100$ and $L = 200$ overlap.

an Figure ?? depicts the Bayes risk versus T for the size of training data varying from Line 10 to $L = 200$. The Bayes risk R^* is also plotted for comparison. The convergence of $E[R_{L,T}]$ to the zero analog is as it was observed for the Bayes risk (see Figure ??). Also the small variability of the difference $E[R_{L,T}] - R^*$ for all L is observed for fixed T (Figure ??). The small variability of the risk with respect to the training data size L . The vertical dashed line at $T = 10$ denotes the simulation space slice presented in Figure ??-b. In subsequent analysis, we will use the value of the fixed risk for different values of the intensity function at the value of the optimal bandwidth fixed by $T = 10$, according to the log-likelihood threshold. For brevity, in Figure ?? we show only the results for λ_1 , noting that the ones obtained for λ_2 are analogous. We observe that in the Bayes risk with λ_1 which aligns with the notion that as the observation window increases, the distribution of events in time becomes sparser, yielding much larger bandwidths. On the other hand, the obtained results also show that Gaussian type L increases. Another way to view this property is to analyze the model log-likelihood versus h for fixed T (Figure ??-b).

Finally, Figure ?? shows the convergence of the empirical kernel rule risk to the Bayes risk for different values of L which does not satisfy the assumptions **A1** and **A2**. While the intensity function has an infinite support in practice, it is extremely unlikely for events to occur outside of some narrow time interval. In Figure ?? we consider the classification problem with $\lambda_1(t) = \lambda_1(t; 300, 20)$ and $\lambda_2(t) = \lambda_2(t; 600, 40)$. For such specified intensity functions we can evaluate that $\int_0^T \lambda_1(t) dt = 1.9 < 1.69$ for all T . Hence, the average algorithm fails to converge. Consider the following Gaussian-type intensity function λ_2 does not hold. Note that the empirical risk does not converge to the Bayes risk that takes very small values for $T > 0.5$. Note also that the risk $E[R_{L,T}]$ is the smallest around the maximum of (??) at $t = 0.5$. Afterwards $R_{L,T}$ slightly increases and reaches a plateau because no new events can be observed.

$$\lambda(t; a, b) = a \exp[-b(t - 0.5)^2], \quad (69)$$

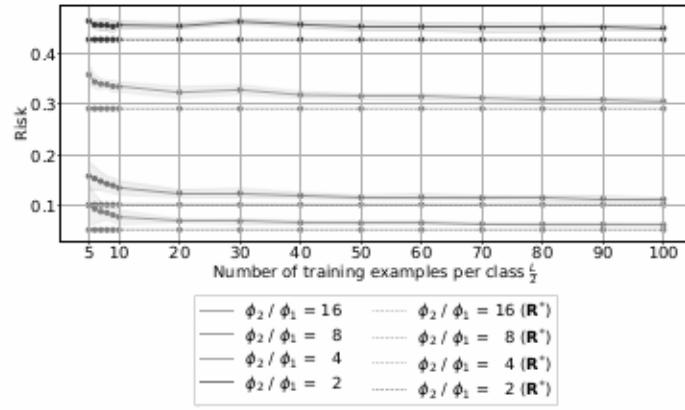


Figure 5. The average risk $E[\mathbf{R}_{L,T}]$ versus L at given T for different values intensity function pairs parametrized by ϕ_1 , ϕ_2 . The horizontal dashed lines denote estimated Bayes risk \mathbf{R}_T^* for the associated intensity function pair.

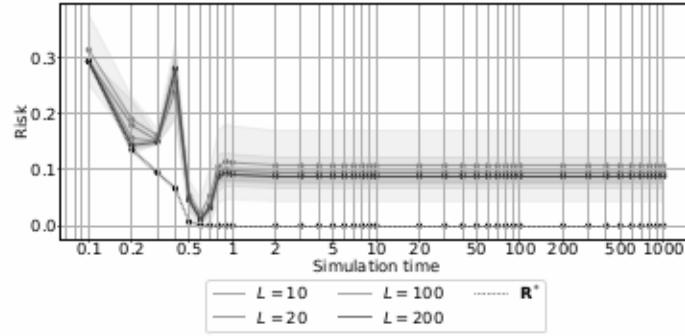


Figure 6. The average risk $E[\mathbf{R}_{L,T}]$ versus T for different values of L for the Gaussian type intensity function.

which does not satisfy the assumptions **A1** and **A2** while the Intensity function has an infinite support, in practice it is extremely unlikely for events to occur outside of some narrow time interval. In Figure ?? we consider the classification problem with $\lambda_1(t) = \lambda_1(t; 300, 20)$ and $\lambda_2(t) = \lambda_2(t; 600, 40)$. For such specified intensity functions we can evaluate that $\int_0^T \lambda_1(t) dt \leq 119$ and $\int_0^T \lambda_2(t) dt \leq 169$ for all T . Hence, the average number of events from each class is finite and consequently the condition **A2** does not hold. Note that the empirical risk does not converge to the Bayes risk that takes very small values for $T > 0.5$. Note also that the risk $E[\mathbf{R}_{\infty}]$ is the smallest around the sufficient conditions for the 0.5. Afterwards \mathbf{R}_{∞} slightly increases and reaches a plateau because no new events can be observed.

There are various ways to extend and generalize the results obtained in this paper. First of all, the log transformed version of the Bayes rule in (??) holds for a general class of point processes such as the Hawkes self-excited process [?] and multivariate or marked point processes [?]. Hence, the extension of our results to this type of point processes is a natural topic for future research. The two-class classification problem studied in this paper has straightforward generalization to the multi-class situation with the class labels denoted as $\{\omega_i\}$. In fact, the Bayes rule in (??) for the two-class classification problem generalizes to the multi-class case. We then introduced a general class of plug-in empirical classification rules and formulated the sufficient conditions for their convergence (as the amount of data grows) to the Bayes risk. This optimality property is confirmed and verified for the plug-in kernel classifier derived from the aggregated data.

There are various ways to extend and generalize the results obtained in this paper. First of all, the log transformed version of the Bayes rule in (??) holds for a general class of point processes such as the Hawkes self-excited process [?] and multivariate or marked point processes [?]. Hence, the extension of our results to this type of point processes where $\gamma_{ik} = \int_0^T (\lambda_i(u) - \lambda_j(u)) du + \log(\tau_k / \tau_i)$. Here $\{\lambda_i(t)\}$ are class intensity functions and $\{\tau_i\}$ are prior probabilities. Utilizing the martingale decomposition (see (??)) for point processes would allow us to generalize our asymptotic generalization to the multi-class situation with the class labels denoted as $\{\omega_i, \dots, \omega_c\}$. In fact, the Bayes rule in (??) results to the multi-class case. Also designing nonparametric plug-in classification rules with the desirable asymptotic optimality property would be of a great practical topic for further research.

$$\mathbf{X} \in \omega_i \text{ if } \sum_{s=1}^N \log \left(\frac{\lambda_i(t_s)}{\lambda_k(t_s)} \right) \geq \gamma_{ik},$$

The asymptotic theory of the classification problem examined in this paper is based on martingale methods. This appendix gives brief summary of the essential facts concerning the counting processes theory and their martingale

where equation $\int_0^T \lambda_t(u) du$ for $\lambda_t(u)$ is the edge of $\lambda_t(u)$. Here $\lambda_t(u)$ is a function of $N(s)$ intensity function and $\lambda_t(u)$ which are probability densities. Utilizing the martingale decomposition (see (??)) for point processes $N(t)$ could allow $dN(t)$ to generalize to the asymptotic version (t) to the threshold classification. Also designing nonparametric $p(N(t))$ in classification implies with the desirable asymptotic optimality property (would be of great practical topic for further research).

$$\mathbb{E}[dN(t)|\mathbf{F}_t] = \lambda(t)dt, \quad (70)$$

where \mathbf{F}_t denotes the history of $N(t)$ in the interval $[0, t]$. Note that $\lambda(t)$ is generally random due to the dependence on the values of $N(t)$ prior to the time t . The formula in (70) implies that the residual process $M(t) = N(t) - \lambda(t)dt$ satisfies the property $\mathbb{E}[dM(t)|\mathbf{F}_t] = 0$. Hence, let $N(t)$ be a spike train process which can be considered as a counting process of the occurrences in the interval $[0, t]$ such that $N(0) = 0$. By $dN(t)$ we denote the increment of $N(t)$ over the small interval $[t, t+dt]$. The evolution of $N(t)$ in time is completely characterized by the local intensity function $\lambda(t)$. This is defined as

This confirms the fact that the process

$$\mathbb{E}[dN(t)|\mathbf{F}_t] = \lambda(t)dt, \quad (70)$$

where \mathbf{F}_t denotes the history of $N(t)$ in the interval $[0, t]$. Note that $\lambda(t)$ is generally random due to the dependence on the values of $N(t)$ prior to the time t . The formula in (??) implies that the residual process is a zero mean local martingale.

The formula in (??) can be written as

$$dM(t) = dN(t) - \lambda(t)dt \quad (71)$$

satisfies the property

$$dN(t) = \lambda(t)dt + dM(t). \quad (74)$$

This can be viewed as the local signal plus noise decomposition of $N(t)$. Moreover, the noise process $dM(t)$ in (??) is a zero mean martingale that has uncorrelated but nonstationary increments [?]. Based on these facts it can be shown that $dM(t)$ has the following second order property

$$\mathbb{E}[dM(t)|\mathbf{F}_t] = 0 \quad \text{Var}[dM(t)|\mathbf{F}_t] = \int_0^t \lambda(u)du. \quad (73)$$

Also $\text{Var}[dM(t)|\mathbf{F}_t] = \text{Var}[dN(t)|\mathbf{F}_t]$.

The fact that $dM(t)$ has uncorrelated increments and that it reveals a piecewise constant sample paths allow us to define the stochastic Stieltjes type integral with respect to $dM(t)$. Hence, let

This can be viewed as the local signal plus noise decomposition of $N(t)$. Moreover, the noise process $dM(t)$ in (??) is a zero mean martingale that has uncorrelated but nonstationary increments [?]. Based on these facts it can be shown that $dM(t)$ has the following second order property

define the stochastic integral of the measurable function $g(t)$ with respect to the increments of the martingale $M(t)$. It is known [?] that the martingale property is preserved under stochastic integration. Since $\mathbb{E}[I(t)] = 0$ the integral $I(t) = \int_0^t g(u)dM(u)$ is a zero mean martingale with respect to the history of the counting process $N(t)$. The variance of $I(t)$ is given by

Also $\text{Var}[dM(t)|\mathbf{F}_t] = \text{Var}[dN(t)|\mathbf{F}_t]$.

The fact that $dM(t)$ has uncorrelated increments and that it reveals a piecewise constant sample paths allow us to define the stochastic Stieltjes type integral with respect to $dM(t)$. Hence, let

The uncorrelated increments property of the martingale process allows us to establish the following generalized version of (??)

define the stochastic integral of the measurable function $g(t)$ with respect to the increments of the martingale $M(t)$. It is known [?] that the martingale property is preserved under stochastic integration. Since $\mathbb{E}[I(t)] = 0$ the integral $I(t) = \int_0^t g_1(u)g_2(u)\lambda(u)du$ is a zero mean martingale with respect to the history of the counting process $N(t)$. The variance of $I(t)$ is given by

$$= \int_0^t g_1(u)g_2(u)\lambda(u)du \quad \text{Var}[I(t)|\mathbf{F}_t] = \int_0^t g^2(u)\lambda(u)du. \quad (76)$$

where $g_1(t), g_2(t)$ are measurable functions.

The uncorrelated increments property of the martingale process allows us to establish the following generalized version of (??)

Appendix B

To prove the results of this section we need the following elementary inequalities

$$\text{Cov}\left[\int_0^t g_1(u)dM(u), \int_0^t g_2(u)dM(u), |\mathbf{F}_t\right] \leq \log(x) \leq x - 1, \quad x > 0. \quad (77)$$

The tighter version of this inequality for $x \geq 1$ reads as follows

$$\text{where } g_1(t), g_2(t) \text{ are measurable functions} \quad \frac{x-1}{x+1} \leq \log(x) \leq \frac{x^2-1}{2x}, \quad x \geq 1. \quad (79)$$

Proof of Lemma ???. Let $\mathbf{X} \in \omega_1$. Then, the formula for the threshold value α_T in (??) becomes

$$\alpha_T = \tau_1 - \tau_2 + \tau_1 \log\left(\frac{\tau_2}{\tau_1}\right) - \tau_1 \int_0^T \log\left(\frac{p_1(t)}{p_2(t)}\right) p_1(t)dt. \quad (80)$$

Here $\mathbf{K}_T(p_1 \parallel p_2) = \int_0^T \log\left(\frac{p_1(t)}{p_2(t)}\right) p_1(t) dt$ is the Kullback-Leibler divergence between the densities $p_1(t)$ and $p_2(t)$. To prove the results of this section we need the following elementary inequalities

$$\alpha_T \leq \tau_1 \frac{x-1}{x+1} \leq \log\left(\frac{\tau_2}{\tau_1}\right) \leq x - 1, \quad \mathbf{K}_T(p_1 \parallel p_2). \quad (78)$$

The tighter version of this inequality for $x \geq 1$ reads as follows

As $\mathbf{K}_T(p_1 \parallel p_2) \geq 0$ we conclude that $\alpha_T \leq \frac{x-1}{x+1} \leq \log(x) \leq \frac{x^2}{2x}$, $x \geq 1$. Concerning the lower bound for α_T in (??) we again use (??). Hence

Proof of Lemma ???. Let $\mathbf{X} \in \omega_1$. Then, the formula for the threshold value α_T in (??) becomes

$$\alpha_T = \tau_1 - \tau_2 + \frac{(\tau_1 - \tau_2)^2}{\tau_2} \log\left(\frac{\tau_2}{\tau_1}\right) \mathbf{K}_T(p_1 \parallel p_2) \int_0^T \log\left(\frac{p_1(t)}{p_2(t)}\right) p_1(t) dt. \quad (80)$$

This confirms the inequalities in Lemma ???. The case when $\mathbf{X} \in \omega_2$ can be proved in the analogous way by noting that α_T is now equal to

$$\begin{aligned} \alpha_T &= \tau_1 - \tau_2 + \tau_2 \log\left(\frac{\tau_2}{\tau_1}\right) + \tau_2 \mathbf{K}_T(p_2 \parallel p_1), \\ \alpha_T &\leq \tau_1 - \tau_2 + \tau_1 \left\{ \frac{\tau_2}{\tau_1} \log\left(\frac{\tau_2}{\tau_1}\right) - \tau_1 \mathbf{K}_T(p_1 \parallel p_2) \right\}. \end{aligned}$$

Then, the application of (??) gives the result in Lemma ???. \square

Proof of Lemma ???. The result of Lemma ?? is implied by the straightforward application of the identity in (??) to As $\mathbf{K}_T(p_1 \parallel p_2) \geq 0$ we conclude that $\alpha_T \leq 0$. Concerning the lower bound for α_T in (??) we again use (??). Hence,

$$\alpha_T \geq \tau_1 - \tau_2 + \log\left(\frac{\lambda_1(t)}{\lambda_2(t)}\right) dM(t) \mathbf{K}_T(p_1 \parallel p_2)$$

Here $M(t)$ is the local martingale corresponding to the intensity $\lambda_1(t)$ or $\lambda_2(t)$ depending whether $\mathbf{X} \in \omega_1$ or $\mathbf{X} \in \omega_2$, respectively. \square

This confirms the inequalities in Lemma ???. The case when $\mathbf{X} \in \omega_2$ can be proved in the analogous way by noting that α_T is now equal to

$$\alpha_T = \tau_1 - \tau_2 + \tau_2 \log\left(\frac{\tau_2}{\tau_1}\right) + \tau_2 \mathbf{K}_T(p_2 \parallel p_1). \quad (81)$$

Then, the application of (??) gives the result in Lemma ???. \square

Proof of Lemma ???. The result of Lemma ?? is implied by the straightforward application of the identity in (??) to the stochastic integral $\frac{x^2}{2+x} \leq J(x) \leq \frac{x^2}{2}$.

The application of the above lower bound in (??) leads to the version of (??) given in (??). \square

Here $M(t)$ is the local martingale corresponding to the intensity $\lambda_1(t)$ or $\lambda_2(t)$ depending whether $\mathbf{R}_T^* \mathbf{X} \in \omega_1$ or $\mathbf{R}_T^* \mathbf{X} \in \omega_2$. Respectively, (??) it suffices to consider the probability of misclassification $\mathbf{P}(W_T(\mathbf{X}) \geq \eta_T | \mathbf{X} \in \omega_2)$. As it has been observed in (??) this probability is equivalent to the following probability.

Proof of Lemma ???. The result in (??) of Lemma ?? is the version of Theorem 5 in [?] that says that under the conditions (a) and (b) of Lemma ?? we have $\mathbf{P}\left(\frac{1}{T} U_T(\mathbf{X}) \geq \frac{1}{T} \alpha_T + \frac{1}{T} \kappa | \mathbf{X} \in \omega_2\right)$,

$$\text{where } \kappa = \log(\pi_2/\pi_1). \text{ Since } \mathbf{X} \in \omega_2 \text{ then } \mathbf{P}(|U_T| \geq \epsilon) \leq 2 \exp\left[-\frac{v_T}{u_T^2} J\left(\epsilon \frac{u_T}{v_T}\right)\right], \quad (81)$$

where $J(x) = (1+x) \log(1+x) - x$. Using the inequalities in (??) we can easily obtain that

$$\text{where } J(x) = (1+x) \log(1+x) - x \geq 0. \quad (83)$$

Then, by the Chebyshev inequality the probability in (??) is bounded by

The application of the above lower bound in (??) leads to the version of (??) given in (??). \square

Proof of Theorem ???. We will prove the result in (??). This clearly implies the convergence $\mathbf{R}_T^* \rightarrow 0$ as $T \rightarrow \infty$. By virtue of (??) it suffices to consider the probability of misclassification $\mathbf{P}(W_T(\mathbf{X}) \geq \eta_T | \mathbf{X} \in \omega_2)$. As it has been observed in (??) this probability is equivalent to the following probability

$$\mathbf{P}\left(\frac{1}{T} U_T(\mathbf{X}) \geq \frac{1}{T} \alpha_T + \frac{1}{T} \kappa | \mathbf{X} \in \omega_2\right), \quad (82)$$

By the assumptions **A1** and **A2** we have $\mathbf{P}\left(\frac{1}{T} U_T(\mathbf{X}) \geq \frac{1}{T} \alpha_T + \frac{1}{T} \kappa | \mathbf{X} \in \omega_2\right)$,

$$\text{where } \kappa = \log(\pi_2/\pi_1). \text{ Since } \mathbf{X} \in \omega_2 \text{ then } \lim_{T \rightarrow \infty} \alpha_T/T \leq d \log\left(\frac{C}{\delta}\right) \quad (86)$$

$$\alpha_T = \tau_1 - \tau_2 + \tau_2 \log\left(\frac{\tau_2}{\tau_1}\right) + \tau_2 \mathbf{K}_T(p_2 \parallel p_1). \quad (83)$$

Then, by the Chebyshev inequality the probability in (??) is bounded by

$$\lim_{T \rightarrow \infty} \alpha_T / T \geq \frac{d \log \left(\frac{\delta}{C} \right)}{b_T T}. \quad (84)$$

Then by the result of Lemma ?? and (??), we get

$$\lim_{T \rightarrow \infty} b_T b_T^{-1} \frac{d \log \left(\frac{C}{\sqrt{T}} U_T(\mathbf{X}) \right) \log(C/\delta)}{d^2 \log^2 \left(\frac{\delta}{C} \right) + T d \log \left(\delta/C \right)}^2. \quad (85)$$

By the assumption **A1** and **A2** we have $\mathbf{P}(W_T(\mathbf{X}) \geq \eta_T | \mathbf{X} \in \omega_2)$ is bounded by b_T/T , where the superior limit of b_T is given in (??). In the analogous way one can show that the probability of misclassification $\mathbf{P}(W_T(\mathbf{X}) < \eta_T | \mathbf{X} \in \omega_1)$ is bounded by a_T/T , where the superior limit of a_T is also given by (??). This concludes the proof of Theorem ??.

Proof of Theorem ??. Consider again $\mathbf{P}(W_T(\mathbf{X}) \geq \eta_T | \mathbf{X} \in \omega_2)$ or equivalently the probability in (??). We wish to use the exponential inequality in (??) of Lemma ???. Then, the probability in (??) is bounded by

$$\lim_{T \rightarrow \infty} \alpha_T / T \geq d \log \left(\frac{e_T^2}{C} \right). \quad (87)$$

Then by the result of Lemma ?? and (??), we get $\exp \left[-T \frac{e_T^2}{2\theta_T + u\epsilon_T} \right]$,

where $u = \log \left(\frac{C}{\delta} \right)$ characterizes the assumption $\lim_{T \rightarrow \infty} b_T \leq \frac{d \log^2 \left(\frac{C}{\delta} \right) \frac{1}{T} \alpha_T}{d^2 \log^2 \left(\frac{C}{\delta} \right)} + \frac{1}{d} \left(\frac{\log(C/\delta)}{\log(\delta/C)} \right)^2 = \text{Var} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right]$. This defines the exponential factor

Hence, the probability of misclassification $\mathbf{P}(W_T(\mathbf{X}) \geq \eta_T | \mathbf{X} \in \omega_2)$ is bounded by b_T/T , where the superior limit of b_T is given in (??). In the analogous way one can show that the probability of misclassification $\mathbf{P}(W_T(\mathbf{X}) < \eta_T | \mathbf{X} \in \omega_1)$ is bounded by a_T/T , where the superior limit of a_T is also given by (??). This concludes the proof of Theorem ??.

Proof of Theorem ??. Consider again $\mathbf{P}(W_T(\mathbf{X}) \geq \eta_T | \mathbf{X} \in \omega_2)$ or equivalently the probability in (??). We wish to use the exponential inequality in (??) of Lemma ???. Then, the probability in (??) is bounded by

$$\exp \left[\frac{1}{3} T \frac{\log(e_T^2/C)}{2\theta_T + u\epsilon_T} \right], \quad (89)$$

This combined with (??) gives the required bound. Since the analogous analysis can be carried out for the probability of misclassification $\mathbf{P}(W_T(\mathbf{X}) < \eta_T | \mathbf{X} \in \omega_1)$ therefore the proof of Theorem ?? has been completed.

$$B_T = \frac{e_T^2}{2\theta_T + u\epsilon_T}.$$

Proof of Theorem ??. The proof of Theorem ?? is in the spirit of the proof of Theorem 1 in [?]. Hence, the consistency results established in (??) and (??) imply that for the selected $\delta > 0$ there exists l_0 such that for $L > l_0$ and $\epsilon > 0$ we have

$$\begin{aligned} \lim_{T \rightarrow \infty} B_T &\geq \frac{d^2 \log^2(\delta/C)}{\mathbf{P} \left(\left| \widehat{W}_{L,T}(\mathbf{x}) - W_T(\mathbf{x}) \right| < \epsilon \right)^2} \geq \frac{ud \log(C/\delta)}{\delta/2}, \\ \mathbf{P} \left(\left| \widehat{\eta}_{L,T} - \eta_T \right| < \frac{1}{3} \log(C/\delta) \right) &> 1 - \delta/2. \end{aligned} \quad (90)$$

This combined with (??) gives the required bound. Since the analogous analysis can be carried out for the probability of misclassification $\mathbf{P}(W_T(\mathbf{X}) < \eta_T | \mathbf{X} \in \omega_1)$ therefore the proof of Theorem ?? has been completed.

$$\mathbf{P} \left(\psi_{L,T}^*(\mathbf{x}) = \psi_T^*(\mathbf{x}) \right) = \mathbf{P} \left(\widehat{W}_{L,T}(\mathbf{x}) > \widehat{\eta}_{L,T} \right).$$

The right-hand side of this equality is not smaller than

Proof of Theorem ??. The proof of Theorem ?? is in the spirit of the proof of Theorem 1 in [?]. Hence, the consistency results established in (??) and (??) imply that for the selected $\delta > 0$ there exists l_0 such that for $L > l_0$ and $\epsilon > 0$ we have

for $0 < \epsilon < \frac{1}{2} (W_T(\mathbf{x}) - \eta_T)$. Moreover, (??) is bounded from below by

$$\begin{aligned} \mathbf{P} \left(\left| \widehat{W}_{L,T}(\mathbf{x}) - W_T(\mathbf{x}) \right| < \epsilon \right) &> 1 - \delta/2, \\ \mathbf{P} \left(\left| \widehat{\eta}_{L,T} - \eta_T \right| < \epsilon \right) &> 1 - \delta/2. \end{aligned} \quad (91)$$

In turn by the elementary inequality $\mathbf{P}(A \cap B) \geq \mathbf{P}(A) + \mathbf{P}(B) - 1$, the lower bound for (??) is

Let $\psi_T^*(\mathbf{x}) = \omega_1$, i.e., we have $W_T(\mathbf{x}) > \eta_T$. Then,

$$\begin{aligned} \mathbf{P} \left(\left| \widehat{W}_{L,T}(\mathbf{x}) - W_T(\mathbf{x}) \right| < \epsilon \right) + \mathbf{P} \left(\left| \widehat{\eta}_{L,T} - \eta_T \right| < \epsilon \right) - 1, \\ \mathbf{P} \left(\psi_{L,T}^*(\mathbf{x}) = \psi_T^*(\mathbf{x}) \right) = \mathbf{P} \left(\widehat{W}_{L,T}(\mathbf{x}) > \widehat{\eta}_{L,T} \right). \end{aligned}$$

Recalling (??) we have shown that for $L > l_0$

The right-hand side of this equality is not smaller than

$$\mathbf{P} \left(\left| \left(\widehat{W}_{L,T}(\mathbf{x}) - \widehat{\eta}_{L,T} \right) - (W_T(\mathbf{x}) - \eta_T) \right| < 2\epsilon \right) > 1 - \delta. \quad (91)$$

Since we can choose an arbitrary small δ , this confirms the claimed convergence. \square

Proof of Lemma ??. The proof will be based on the following version of Helly's theorem [?] for the Stieltjes integral.

for $0 < \epsilon < \frac{1}{2}(W_T(\mathbf{x}) - \eta_T)$. Moreover, (??) is bounded from below by

$$\mathbf{P}\left(\left|\widehat{W}_{L,T}(\mathbf{x}) - W_T(\mathbf{x})\right| < \epsilon, |\widehat{\eta}_{L,T} - \eta_T| < \epsilon\right). \quad (92)$$

If $g(x)$ is a function of bounded variation on $[a, b]$ then

In turn by the elementary inequality $\mathbf{P}(A \cap B) \geq \mathbf{P}(A) + \mathbf{P}(B) - 1$, the lower bound for (??) is

$$\mathbf{P}\left(\left|\widehat{W}_{L,T}(\mathbf{x}) - W_T(\mathbf{x})\right| > \epsilon\right) + \mathbf{P}(|\widehat{\eta}_{L,T} - \eta_T| > \epsilon) - 1. \quad (93)$$

Considering (??) we have shown that for $L \rightarrow \infty$ in (??) and its empirical counterpart $\widehat{W}_{L,T}(\mathbf{x})$ in (??). Then, we can write (see (??))

$$\widehat{W}_{L,T}(\mathbf{x}) - W_{L,T}(\mathbf{x}) = \mathbf{P}\left(\widehat{\psi}_{L,T}(\mathbf{x}) = \psi_T^*(\mathbf{x})\right) > 1 - \delta.$$

Since we can choose an arbitrary small δ , this confirms the claimed convergence. \square

Proof of Lemma ???. The proof will be based on the following version of Helly's theorem [?] for the Stieltjes integral. Let

We wish to prove that $\left|\widehat{W}_{L,T}(\mathbf{x}) - W_{L,T}(\mathbf{x})\right| \rightarrow 0$ as $L \rightarrow \infty$. Owing to Helly's theorem it suffices to show that

If $g(x)$ is a function of bounded variation on $[a, b]$ then $\left|\log\left(\frac{p_1(t)}{p_2(t)}\right)\right| \rightarrow 0$ (P)

$$\int_a^b g(x) dg(x) \rightarrow \int_a^b f(x) dg(x) \text{ as } L \rightarrow \infty. \quad (93)$$

Observe that the left-hand side of (??) is equal to $\left|\log\left(\frac{\widehat{p}_1(t)p_2(t)}{\widehat{p}_2(t)p_1(t)}\right)\right|$. Then, using (??) this is bounded by Consider the optimal decision function $W_T(\mathbf{x})$ in (??) and its empirical counterpart $\widehat{W}_{L,T}(\mathbf{x})$ in (??). Then, we can write (see (??))

$$\begin{aligned} & \left|\frac{\widehat{p}_1(t)p_2(t) - \widehat{p}_2(t)p_1(t)}{\widehat{W}_{L,T}(\mathbf{x})(t) p_{L,T}(\mathbf{x})}\right| \\ &= \left|\frac{(\widehat{p}_1(t) - p_1(t))p_2(t) + (p_2(t) - \widehat{p}_2(t))p_1(t)}{\left(\frac{\log\left(\frac{\widehat{p}_1(t)p_2(t)}{\widehat{p}_2(t)p_1(t)}\right)}{\widehat{p}_2(t)}\right) p_1(t) + \left(\frac{\log\left(\frac{\widehat{p}_1(t)p_2(t)}{\widehat{p}_2(t)p_1(t)}\right)}{\widehat{p}_2(t)}\right) p_2(t) N(t)}\right|. \end{aligned} \quad (94)$$

This is not greater than We wish to prove that $\left|\widehat{W}_{L,T}(\mathbf{x}) - W_{L,T}(\mathbf{x})\right| \rightarrow 0$ as $L \rightarrow \infty$. Owing to Helly's theorem it suffices to show that

By the assumption **A1** limited to the interval $[0, T]$ and the fact that $p_i(t) = \lambda_i(t)/\tau_i$ the above expression does (95) exceed uniformly on $[0, T]$ as $L \rightarrow \infty$

$$\left(\frac{C}{\delta}\right) T \left\{ |\widehat{p}_1(t) - p_1(t)| + |\widehat{p}_2(t) - p_2(t)| \right\}.$$

Observe that the left-hand side of (??) is equal to $\left|\log\left(\frac{\widehat{p}_1(t)p_2(t)}{\widehat{p}_2(t)p_1(t)}\right)\right|$. Then, using (??) this is bounded by

This by recalling the assumption in (??) proves (??). The proof of Lemma ?? has been completed. \square

$$\left|\frac{\widehat{p}_1(t)p_2(t) - \widehat{p}_2(t)p_1(t)}{\widehat{p}_2(t)p_1(t)}\right| \xrightarrow{\text{Appendix D}}$$

Proof of Lemma ???. We wish to show $\left|\frac{(\widehat{p}_1(t) - p_1(t))p_2(t) + (p_2(t) - \widehat{p}_2(t))p_1(t)}{(\widehat{p}_2(t) - p_2(t))p_1(t) + p_1(t)p_2(t)}\right|$

This is not greater than $\sup_{t \in T_\epsilon} |\lambda(t) - \widehat{\lambda}(t)| \rightarrow 0$ (P) as $L \rightarrow \infty$. \square

where $T_\epsilon = [\epsilon, T - \epsilon]$ for small $\epsilon > 0$. We begin with the standard bounding into the variance and bias terms

By the assumption **A1** limited to the interval $[0, T]$ and the fact that $p_i(t) = \lambda_i(t)/\tau_i$ the above expression does (97) exceed

Owing to (??) the bias term is equal to $\left(\frac{C}{\delta}\right)^3 T \left\{ |\widehat{p}_1(t) - p_1(t)| + |\widehat{p}_2(t) - p_2(t)| \right\}$.

$$\mathbb{E}[\widehat{\lambda}(t)] - \lambda(t) = \int_{(T-h)h}^T K(s)\lambda(t+hs) ds - \lambda(t)$$

This by recalling the assumption in (??) proves (??). The proof of Lemma ?? has been completed. \square

for $t \in T_\epsilon$. Since $K(t)$ and $\lambda(t)$ are positive and $K(t)$ is a density function supported on $[-1, 1]$ then we have $\lambda(t)$ Appendix D

Proof of Lemma ???. We wish to show that $\left|\mathbb{E}[\widehat{\lambda}(t)] - \lambda(t)\right| = \int_{-1}^1 K(s) |\lambda(t+hs) - \lambda(t)| ds$

$$\sup_{t \in T_\epsilon} \left|\widehat{\lambda}(t) - \lambda(t)\right| \rightarrow 0 \text{ (P) as } L \rightarrow \infty. \quad (96)$$

where $T_\epsilon = [\epsilon, T - \epsilon]$ for small $\epsilon > 0$. We begin with the standard bounding into the variance and bias terms uniformly in $t \in T_\epsilon$, where M_λ is the Lipschitz constant of $\lambda(t)$.

Let us consider the stochastic part $\lambda(t) - \mathbb{E}[\widehat{\lambda}(t)]$ as the interval $[0, T]$ is compact one $\lambda(t)$ define a finite partition of T_ϵ (97) disjoint equal size intervals, i.e., $T_\epsilon = \bigcup_{j=1}^{q(L)} \mathbf{U}_j$, where the size of \mathbf{U}_j is denoted as $\Delta(L)$. Clearly the number of intervals

Owing to (T) and the fact that $\lambda(t)$ is bounded between the middle point of \mathbf{U}_j . Then, the uniform norm of the stochastic term in (??) can be bounded as follows

$$\mathbb{E} \left[\left| \hat{\lambda}(t) - \lambda(t) \right| \right] = \mathbb{E} \left[\int_{-t/h}^{(T-t)/h} K(s) \lambda(t+hs) ds - \lambda(t) \right]$$

for $t \in T_\epsilon$. Since $K(t)$ and $\lambda(t)$ are positive and $K(t)$ is a density function supported on $[-1, 1]$ then we have

$$\begin{aligned} & \leq \max_{1 \leq j \leq q(L)} \sup_{t \in T_\epsilon \cap \mathbf{U}_j} \left| \hat{\lambda}(t) - \lambda(u_j) \right| \\ & + \mathbb{E} \left[\left| \hat{\lambda}(t) - \lambda(t) \right| \right] = \sup_{1 \leq j \leq q(L)} \mathbb{E} \left[\left| \hat{\lambda}(t) - \mathbb{E} \left[\hat{\lambda}(u_j) \right] \right| \right] \\ & + \max_{1 \leq j \leq q(L)} \left| \hat{\lambda}(u_j) M_\lambda h \mathbb{E} \int_1^1 |\hat{\lambda}(K_j(s))| s ds \right| \end{aligned} \quad (98)$$

uniformly in $t \in T_\epsilon$, where M_λ is the Lipschitz constant of $\lambda(t)$.

Let us consider the stochastic part in (??). As the interval T_ϵ is compact, one can define a finite partition of T_ϵ into disjoint equal size intervals, i.e., $T_\epsilon = \bigcup_{j=1}^{q(L)} \mathbf{U}_j$, where the size of \mathbf{U}_j is denoted as $\Delta(L)$. Clearly the number of intervals is of order $T_\epsilon/\Delta(L)$. Let $u_j \in \mathbf{U}_j$ be the middle point of \mathbf{U}_j . Then, the uniform norm of the stochastic term in (??) can be bounded as follows

for $t, u_j \in \mathbf{U}_j$. Noting that $|t - u_j| \leq \Delta(L)$ and using the fact that $K(t)$ is Lipschitz we get

$$\sup_{t \in \mathbf{U}_j} \left| \hat{\lambda}(t) - \mathbb{E} \left[\hat{\lambda}(t) \right] \right| \leq M_K \frac{\Delta(L)}{h^2} \int_0^T d\bar{N}_L(s). \quad (99)$$

Note that $\int_0^T d\bar{N}_L(s) = \bar{N}_L(T)$ and we know, see (??), that $\mathbb{E} [\bar{N}_L(T)] = \int_0^T \lambda(t) dt$ and $\text{Var} [\bar{N}_L(T)] = \frac{1}{L} \int_0^T \lambda(t) dt$. This proves that

$$+ \max_{1 \leq j \leq q(L)} \sup_{t \in \mathbf{U}_j} \left| \mathbb{E} \left[\hat{\lambda}(t) \right] - \mathbb{E} \left[\hat{\lambda}(u_j) \right] \right| \cdot \quad (98)$$

$$+ \max_{1 \leq j \leq q(L)} \left| \hat{\lambda}(u_j) - \mathbb{E} \left[\hat{\lambda}(u_j) \right] \right| \quad (100)$$

uniformly in $t \in T_\epsilon$. Concerning the term A_1 in (??) we can use (??). Then, we obtain

$$= A_1 + A_2 + A_3.$$

Consider first the term A_1 . By virtue of (??) we have

$$\hat{\lambda}(t) - \hat{\lambda}(u_j) = \int_0^T \frac{[K_h(t-s) - K_h(u_j-s)] \lambda(s) ds}{[K_h(t-s) - K_h(u_j-s)] d\bar{N}_L(s)}.$$

This gives

for $t, u_j \in \mathbf{U}_j$. Noting that $|t - u_j| \leq \Delta(L)$ and using the fact that $K(t)$ is Lipschitz we get

$$A_2 = \mathcal{O} \left(\frac{\Delta(L)}{h^2} \right). \quad (101)$$

Hence, we have shown that the terms A_1 and A_2 are of order $\mathcal{O} \left(\frac{\Delta(L)}{h^2} \right)$, where $\Delta(L)$ is to be selected.

Finally let us consider the term A_3 in (??). First we note that for $\delta > 0$ we know, see (??), that $\mathbb{E} [\bar{N}_L(T)] = \int_0^T \lambda(t) dt$ and $\text{Var} [\bar{N}_L(T)] = \frac{1}{L} \int_0^T \lambda(t) dt$. This proves that

$$\begin{aligned} & \mathbb{P} \left(\max_{1 \leq j \leq q(L)} \left| \hat{\lambda}(u_j) - \mathbb{E} \left[\hat{\lambda}(u_j) \right] \right| \geq \delta \right) \\ & \leq q(L) \sup_{t \in T_\epsilon} \mathbb{P} \left(\left| \hat{\lambda}(t) - \mathbb{E} \left[\hat{\lambda}(t) \right] \right| \geq \delta \right). \end{aligned} \quad (102)$$

uniformly in $t \in T_\epsilon$. Concerning the term A_3 in (??) we can use (??). Then, we obtain

By virtue of (??) and (??) we have

$$\begin{aligned} & \mathbb{E} \left[\hat{\lambda}(t) \right] - \mathbb{E} \left[\hat{\lambda}(u_j) \right] \\ & \hat{\lambda}(t) \int_0^T \mathbb{E} \left[\hat{\lambda}(t-s) \right] K_h(u_j-s) d\bar{M}_L(s). \end{aligned}$$

This, (??) and Chebyshev inequality yield

$$\mathbb{P} \left(\left| \hat{\lambda}(t) - \mathbb{E} \left[\hat{\lambda}(t) \right] \right| \geq \delta \right) = \mathcal{O} \left(\frac{\Delta(L)^2}{L h^3} \right) K_h^2(t-s) \lambda(s) ds / \delta^2. \quad (101)$$

Note that the right-hand side of this inequality is of order $\mathcal{O}(1/Lh^3)$ uniformly in L . This proves (??) and the fact that $q(L) = \mathcal{O}(1/\Delta(L))$ lead to the following uniform bound

Finally, let us consider the term A_3 in (??). First we note that for $\delta > 0$

$$\mathbb{P} \left(\max_{t \in T_\epsilon} \left| \hat{\lambda}(t) - \mathbb{E} \left[\hat{\lambda}(t) \right] \right| \geq \delta \right)$$

or equivalently $A_3 = \mathcal{O}_P \left(1/\sqrt{\Delta(L)Lh} \right)$. Hence, balancing $A_3 = \mathcal{O}_P \left(1/\sqrt{\Delta(L)Lh} \right)$ versus $A_1, A_2 = \mathcal{O}(\Delta(L)/h^2)$ gives the choice $\Delta(L) = h/L^{1/3}$. This yields the convergence in (??) if $h(L) \rightarrow 0$ and $Lh^3(L) \rightarrow \infty$ as $L \rightarrow \infty$.

The proof of Lemma ?? has been completed. \square

By virtue of (??) and (??) we have

This work was supported by the Polish National Center of Science under Grant DEC-2017/27/B/ST7/03082 and NSERC Grant 319732.

This, (??) and Chebyshev inequality yield

Acknowledgment

$$\lambda(t) - \mathbb{E}[\lambda(t)] = \int_0^T K_h(t-s)dM_L(s).$$

References

- [1] W. Gerstner, W. M. Kistler, R. Nauj, and L. Paninski, *Neuronal Dynamics: From Single Neurons to Networks and Models of Cognition*. Cambridge University Press, 2014.
- [2] H. Jang, O. Simeone, B. Gardner, and A. Grunig, "An introduction to probabilistic spiking neural networks: Probabilistic models, learning rules, and applications," *IEEE Signal Processing Magazine*, vol. 36, no. 6, pp. 64–77, 2019.
- Note that the right hand side of this inequality is of order $O(1/\Delta(L))$, uniformly in $T \in T_\delta$. This, (??) and the fact that $q(L) = O(\Delta(L))$ lead to the following uniform bound
- [3] O. Shchur, A. G. Türkmen, T. Jamuszkowski, and S. Günnemann, "Neural temporal point processes: A review," arXiv preprint arXiv:2104.03528, 2021.
- [4] I. Bar-David, "Communication under the Poisson regime," *IEEE Transactions on Information Theory*, vol. 15, pp. 31–37, 1969.
- [5] D. Guo, S. Shamai, and S. Verdú, "Mutual information and conditional mean estimation in Poisson channels," *IEEE Transactions on Information Theory*, vol. 54, pp. 1827–1840, 2008.
- [6] or equivalently $\mathcal{O}(L) = \mathcal{O}_P\left(\frac{1}{\sqrt{\Delta(L)L}}\right)$. Hence, balancing $A_3 = \mathcal{O}_P\left(\frac{1}{\sqrt{\Delta(L)L}}\right)$ versus $A_4, A_2 = \mathcal{O}(\Delta(L)/L^2)$, gives the choice $\Delta(L) = h/L^{1/3}$. This yields the convergence in (??) if
- [7] D. J. Daley and D. Vere-Jones, *An Introduction to the Theory of Point Processes*. Springer, 2003.
- [8] P. K. Andersen, O. Borgan, R. D. Gill, and N. Keiding, *Statistical Models Based on Counting Processes*. Springer, 2012.
- [9] L. Devroye, L. Györfi, and G. Lugosi, *A Probabilistic Theory of Pattern Recognition*. Springer, 1996.
- [10] A. Chalapudis, L. Forzani, P. Llop, and L. Moreno, "On the classification problem for Poisson point processes," *Journal of Multivariate Analysis*, vol. 163, pp. 1–15, 2017.
- [11] X. Rong and V. Solo, "On the error rate for classifying point processes," in *60th IEEE Conference on Decision and Control (CDC)*. IEEE, 2021, pp. 120–125.
- [12] M. Mazza, S. Cresci, M. Avvenuti, W. Quattrociocchi, and M. Tesconi, "RTburst: Exploiting temporal patterns for botnet detection on This work was supported by the Polish National Center of Science under Grant DEC-2017/27/B/ST7/03082 and NSERC Grant 319732."
- [13] H. Jang, O. Simeone, B. Gardner, and A. Grunig, "Event labeling combining ensemble detectors and background knowledge," *Progress in Artificial Intelligence*, vol. 2, no. 2, pp. 113–127, 2014.
- [14] W. Greblicki, "Asymptotically optimal pattern recognition procedures with density estimates," *IEEE Transactions on Information Theory*, vol. 24, no. 2, pp. 250–251, 1978.
- [15] P. Diggle and J. S. Marron, "Equivalence of smoothing parameter selectors in density and intensity estimation," *Journal of the American Statistical Association*, vol. 83, no. 403, pp. 793–800, 1988.
- [16] W. Greblicki, W. M. Kistler, R. Nauj, and L. Paninski, *Neuronal Dynamics: From Single Neurons to Networks and Models of Cognition*. Cambridge University Press, 2014.
- [17] M. Pawłak, M. Pabian, and D. Rzepka, "Asymptotically optimal nonparametric classification rules for spike train data," in *ICASSP 2023-IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 64–77, 2023.
- [18] O. Shchur, A. G. Türkmen, T. Jamuszkowski, and S. Günnemann, "Neural temporal point processes: A review," arXiv preprint arXiv:2104.03528, 2021.
- [19] I. Bar-David, "Communication under the Poisson regime," *IEEE Transactions on Information Theory*, vol. 15, pp. 31–37, 1969.
- [20] D. Guo, S. Shamai, and S. Verdú, "Mutual information and conditional mean estimation in Poisson channels," *IEEE Transactions on Information Theory*, vol. 54, pp. 1827–1840, 2008.
- [21] N. Mørkøv Geel, "Exponential inequalities for martingales, with application to maximum likelihood estimation for counting processes," *The Annals of Statistics*, pp. 1779–1801, 1995.
- [22] P. Diggle and A. Vered-Jones, "Optimal correlogram for detection and estimation in optical receivers," *IEEE Transactions on Information Theory*, vol. 67, pp. 5200–5210, 2021.
- [23] P. J. Green and A. van der Vaart, *Fundamentals of Nonparametric Bayesian Inference*. Cambridge University Press, 2017.
- [24] D. J. Daley and D. Vere-Jones, *An Introduction to the Theory of Point Processes*. Springer, 2003.
- [25] R. Le Guével, "Exponential inequalities for the supremum of some counting processes and their square martingales," *Comptes Rendus. Mathématique*, vol. 309, no. 8, pp. 969–974, 2004.
- [26] P. K. Andersen, O. Borgan, R. D. Gill, and N. Keiding, *Statistical Models Based on Counting Processes*. Springer, 2012.
- [27] D. Vere-Jones, "On the estimation of frequency of point processes," *Journal of Applied Probability*, vol. 19, pp. 383–394, 1982.
- [28] A. Chalapudis, L. Forzani, P. Llop, and L. Moreno, "On the classification problem for Poisson point processes," *Journal of Multivariate Analysis*, vol. 153, pp. 1–15, 2017.
- [29] M. P. Wand and M. C. Jones, *Kernel Smoothing*. CRC Press, 1994.
- [30] W. Greblicki and M. Pawłak, *Nonparametric System Identification*. Cambridge University Press, 2008.
- [31] X. Rong and V. Solo, "On the error rate for classifying point processes," in *60th IEEE Conference on Decision and Control (CDC)*. IEEE, 2021, pp. 1449–1455, 2020.
- [32] M. Mazza, S. Cresci, M. Avvenuti, W. Quattrociocchi, and M. Tesconi, "RTburst: Exploiting temporal patterns for botnet detection on twitter," in *Proceedings of the 10th ACM conference on web science*, 2019, pp. 183–192.
- [33] H. Jang, O. Simeone, B. Gardner, and A. Grunig, "Event labeling combining ensemble detectors and background knowledge," *Progress in Artificial Intelligence*, vol. 2, no. 2, pp. 113–127, 2014.
- [34] W.-Y. Audibert and A. B. Tsybakov, "Fast learning rates for plug-in classifiers," *The Annals of Statistics*, vol. 35, pp. 608–633, 2007.
- [35] D. Guo, S. Shamai, and S. Verdú, "Martingale procedures, model selection, inference, and control: an overview," *SIAM Review*, vol. 65, pp. 331–374, 2023.
- [36] T. Apostol, *Mathematical Analysis*. Addison-Wesley, 1974.
- [37] L. Birgé and P. Massart, "Rates of convergence for minimum contrast estimators," *Probability Theory and Related Fields*, vol. 97, pp. 113–150, 1993. Miroslaw Pawlak received the Ph.D. (under the supervision of Prof. Greblicki) and D.Sc. degrees in computer engineering from Wroclaw University of Technology, Wroclaw, Poland, in 1984 and 2006, respectively. He is currently a Professor at the Department of Electrical and Computer Engineering, University of Manitoba, Winnipeg, MB, The Annals of Statistics, pp. 1779–1801, 1995.
- [38] S. Ghosal and A. Van der Vaart, *Fundamentals of Nonparametric Bayesian Inference*. Cambridge University Press, 2017.
- [39] R. Le Guével, "Exponential inequalities for the supremum of some counting processes and their square martingales," *Comptes Rendus. Mathématique*, vol. 339, no. 8, pp. 969–974, 2004.
- [40] D. Vere-Jones, "On the estimation of frequency of point processes (Capbridge Univ. Press, 2008)." Journal of Applied Probability, Vol. 19, pp. 383–394, 1982. His research interests include statistical signal processing, machine learning, and nonparametric modeling. Dr. Pawlak has been a Associate Editor for *Signal Processing*, *Machine Learning*, *Pattern Recognition*, *International Journal of Approximate Reasoning*, *Image Analysis by Moments* (Wroclaw Univ. Technol. Press, 2006), and *Journal of Nonparametric System Identification*.
- [41] M. P. Wand and M. C. Jones, *Kernel Smoothing*. CRC Press, 1994.
- [42] W. Greblicki and M. Pawłak, *Nonparametric System Identification*. Cambridge University Press, 2008.
- [43] M. D. Cattaneo, M. P. Wand, and M. C. Jones, "Simple local polynomial density estimators," *Journal of the American Statistical Association*, vol. 115, no. 531, pp. 1449–1455, 2020.
- [44] E. Giné and A. Guillou, "Rates of strong uniform consistency for multivariate kernel density estimators," *Annales de l'Institut Henri Poincaré (B) Probability and Statistics*, vol. 38, no. 6, pp. 907–921, 2002.

- [28] E. Masry, "Multivariete regression estimation with local polynomial fitting for time series," *Statistica et Proximitate in Proximity of Science in Applications*, vol. 65, no. 1, pp. Kielce, Poland, in 2019, where he is currently pursuing the Ph.D. degree at the Department of Measurement and
- [29] J.-Y. Audibert and A. Tsybakov, 2017. *Fast learning rates for plug-in classifiers*. *The Annals of Statistics*, vols. 35, pp. 608–633.
- [30] R. Lima, "Hawkes processes in modeling inference, prediction and control which interview," *SIAM Review*, vol. 65, pp. 331–374, 2023 and analysis of
- [31] T. Apostol, *Mathematical Analysis*. Addison-Wesley, 1974. From 1974 to 2022 he was a Machine Learning Researcher at Comarch Healthcare, where he has been involved in the development of computer-aided medical diagnostics systems. Since 2022 he has been working as Model Developer at UBS in Kraków, Poland. His research interests include machine learning and event-based systems, particularly spiking neural networks.



Miroslaw Pawlak received the Ph.D. (under the supervision of Prof. Greblicki) and D.Sc. degrees in computer engineering from Wroclaw University of Technology, Wroclaw, Poland, in 1984 and 2006, respectively. He is currently a Professor with the Department of Electrical and Computer Engineering, University of Manitoba, Winnipeg, MB, Canada. He has held a number of visiting positions in North American, Australian, and European Universities. He was with the University of Ulm and University of Goettingen as an Alexander von Humboldt Foundation Fellow. Among his publications in these areas are the books *Image Analysis by Moments* (Wroclaw Univ. Technol. Press, 2006), and *Nonparametric System Identification* (Cambridge Univ. Press, 2008), coauthored with Prof. Włodzimierz Greblicki. His research interests include statistical signal processing, machine learning, and nonparametric modeling. Dr. Pawlak has been an Associate Editor for the *Journal of Pattern Recognition and Applications*, *Pattern Recognition*, *International Journal on Sampling Theory in Signal and Image Processing*, *Opuscula Mathematica* and *Statistics in Transition-New Series*.



Mateusz Pabian received the M.Sc. degree in biomedical engineering from the AGH University of Science and Technology, Kraków, Poland, in 2019, where he is currently pursuing the Ph.D. degree at the Department of Measurement and Electronics. In 2014, he completed an internship at NTT from the Institute of Sciences and Technology of the University of Vienna, Austria. In 2014–2015, he participated in the program which focused on experimental design in signal acquisition and analysis of metabolites evoked potentials. From 2017 to 2022, he was a Machine Learning Researcher at Comarch Healthcare, where he has been involved in the development of computer-aided medical diagnostics systems. Since 2022, he has been working as a Model Developer at UBS in Kraków, Poland. His research interests include machine learning and event-based systems, particularly spiking neural networks and neuromorphic machine learning. He was a Visiting Student and Postdoc Researcher in the University of Manitoba, Winnipeg, Canada, from 2014 to 2023, and in The City College of New York, USA, in 2015. Since 2015, he is working as Signal Processing and Machine Learning Researcher in Comarch Healthcare and Fitech, developing algorithms for diagnostics and quality assurance systems. His research interests include signal processing and machine learning in biomedicine, wireless communication and industrial inspection, and event-based systems.



Dominik Rzepka received his M.Sc. and Ph.D. degree in electrical engineering from the AGH University of Science and Technology, Kraków, Poland, in 2009 and 2018 respectively. From 2007 to 2011, he was with the Wireless Sensor and Control Networks Group, AGH University of Science and Technology, where he was involved in the design of the low-power algorithms for the processing of radio signals and software-defined radio. In 2011, he joined the Event-Based Control and Signal Processing Group at AGH University of Science and Technology, where he currently works on methods of signals reconstruction from event-triggered samples and on neuromorphic machine learning. He was a Visiting Student and Postdoc Researcher in the University of Manitoba, Winnipeg, Canada, from 2014 to 2023, and in The City College of New York, USA, in 2015. Since 2015, he is working as Signal Processing and Machine Learning Researcher in Comarch Healthcare and Fitech, developing algorithms for diagnostics and quality assurance systems. His research interests include signal processing and machine learning in biomedicine, wireless communication and industrial inspection, and event-based systems.