

Bayes Risk Consistency of Nonparametric Classification Rules for Spike Trains Data

Mirosław Pawłak¹, M. Meshber¹, Mateusz Paliński¹, Małgorzata Pobłocka¹, Emanuele Domínguez², and Dominik Rzepka¹, 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 finite 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 specifications of the underlying spike generation model. Asymptotic 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

Index Terms

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

1. 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 [1]. In optical communication systems and engineering signal processing systems representing impulse emitted by photons that define arrival times of events. In communication and machine learning this type of data have been called spike trains [2] [3]. In optical communication systems one observes a train of impulses (representing point process) emitted by photodetectors. Signal detection and estimation methods for such processes have been extensively studied in the statistical and stochastic processes literature [4], [5]. However, the Poisson regime spike processes have been extensively examined in communication and information theory [2], [6], [7]. On the other hand, the probabilistic spiking neural networks have been extensively studied in learning problems [8]. Various spiking neuron literature have been reported supporting research interests, hypotheses from the accuracy studies and fundamental limits. In this paper we developed Bayes strategy [9] for the probabilistic spiking neural networks. This strategy has been introduced for applied research in supervised learning problems [9]. Various simulation results have been reported supporting their spiking process without however characterizing its statistical and fundamental limits. The intensity function plays the central role in our theory of point processes. In Section 2 the Bayes strategy [9] for the spiking process is proposed. This strategy can be applied to the Bayes consistency of spike trains. In Section 3 we develop optimal Bayes rule. This strategy can be applied to the Bayes consistency problem. In Section 4 the main idea of the Bayes rule is extended to the increasing length of the observation interval. In Section 5 the Bayes rule is characterized by nonparametric intensity functions. The intensity function plays the central role in our theory of point processes. This describes the local rate of potential optimal spikes. For such processes (Section 6) we derive the Bayes rule. The Bayes rule is considered as the continuous part of the result in the limit concerning the Bayes rule with respect to the increasing length of the observation interval is examined. In Section 7 the plug-in nonparametric kernel classification rule from multiple replicates of spiking process is proposed. This is given followed by the asymptotic analysis. In this paper we examine the Bayes rule to the Bayes decomposition. This result can be considered as the counterpart worth mentioning that the asymptotic plug-in nonparametric classifier long single realization finite-dimensional Euclidean space. The spike trains data are characterized by the spiking process does not consist of discrete classes of events in a manner similar to the case of a given observation interval. The main mathematical tool in this asymptotic analysis is the theory of the random variables in positions. This can be approached either scaling the intensity function or by using the probability measure. It is also worth mentioning that the Bayes rule approach does not make use of the multiplicative long single realization of the underlying spiking process. In this paper we propose the intensity standardization learning spiking process which depicts following the classical large sample setting. It will be obtained by averaging the point processes from single realizations. Besides, for a fixed observation interval it is supported by simulation studies presented in Section 8. The preliminary version of the resulting development in this paper has been reported in [10]. The preliminary version of the spiking process can be based on the multiplicative intensity model indicator function [11] where set \mathcal{A} is the set of points $\{P\}$ of \mathbb{R}^d which are probabilities, whereas \mathcal{A} denotes the the set of points with probability zero. In this case the resulting hypothesis bound will be $\underline{\lambda}$ if sufficiently large T . Furthermore, $\overline{\lim}_{T \rightarrow \infty}$, $\underline{\lim}_{T \rightarrow \infty}$ denote the limit superior and inferior, respectively. Also, by M_f we will denote the Lipschitz constant of a function $f(t)$, i.e., $f(t) \leq f(t_2) + M_f|t - t_2|$ for all t, t_2 .

¹M. Pawłak 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 Kraków, Poland. E-mail: Mirosław.Pawlak@umanitoba.ca.

²M. Pawłak 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 Kraków, Poland. E-mail: Mirosław.Pawlak@umanitoba.ca.

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

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

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

by aggregating kernel estimates from single realizations. Our asymptotic results are supported by simulation studies presented in Section ??.

The preliminary version of the results developed in this paper has been reported in [?]. A temporal spiking process $\{N(t), t \geq 0\}$ consists of a sequence of isolated events in time such that $N(0) = 0$. The process $N(t)$ can be defined by the counting function $N(t) = \sum_{i=1}^N \mathbf{1}_{\{t_i \leq t\}}$ which is the number of events in $[0, t]$. We assume that the process is observed on the time window $[0, T]$ and is characterized by the non-random intensity function $\lambda(t)$ that is defined for all $t \geq 0$. This is a non-negative function that describes the local arriving rate of events such that $\mathbb{E}[N(T)] = \int_0^T \lambda(u)du$ is the average number of events in $[0, T]$. Hence, the observed on $[0, T]$ process $N(t)$ can be represented by the variable-length vector $\mathbf{X} = [t_1, \dots, t_N; N]$, where $0 < t_1 < \dots < t_N < T$ are the event times and $N = N(T)$. Writing $[t_1, \dots, t_N; N]$ we emphasize the fact that the data vector consists of two parts: the occurrence times $\{t_1, \dots, t_N\}$ and N being the number of events in $[0, T]$. The former is the continuous part of the vector \mathbf{X} , whereas the latter is its discrete part.

A **temporal splitting process** $\{N(t), t \geq 0\}$ consists of a sequence of random times for the arrival of isolated events in spiking processes based on the Bayes theorem of classification. Defining the counting function consider a two-class classification problem (see Section ??) for the generalization that the process is observed over a short time window $[0, T]$ and with the prior probabilities non-randomly respectively. If function $\lambda(t)$ is defined by as well we refer to this following known function [?] that provides the local density of \mathbf{X} of events such that $\mathbb{E}[N(T)] = \int_0^T \lambda(u)du$ is the average number of events in $[0, T]$. Hence, the observed on $[0, T]$ process $N(t)$ can be represented by the variable-length vector $\mathbf{X} = [t_1, \dots, t_N; N]$, where $0 < t_1 < \dots < t_N < T$ are the event times and $N = N(T)$. Writing $[t_1, \dots, t_N; N]$ we emphasize the fact that the data vector consists of two parts: the occurrence times $\{t_1, \dots, t_N\}$ and N being the number of events in $[0, T]$. The former is the continuous part of the vector \mathbf{X} , whereas the latter is its discrete part.

The goal of this paper is to develop a rigorous classification methodology for the aforementioned class of spiking 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 processes based on the Bayes theory of classification [?]. Without a loss of generality we consider a two-class classification problem (see Section ?? for the generalization to the multi-class case) where class labels are denoted as ω_1, ω_2 with the priori 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 $\mathbf{X} = (T_1, T_2, \dots, T_N)$,

$$f(\mathbf{x}) = \prod_{i=1}^n \lambda(t_i) \exp\left(\sum_{n=1}^{\infty} \prod_{j=1}^n \lambda(t_j) \int_0^T \lambda(u) du\right) \quad (1)$$

which is defined over the simplex regions $\mathbb{C}_n = \{(t_1, \dots, t_n) : 0 \leq t_1 \leq \dots \leq t_n \leq T\}$, $n = 1, 2, \dots$. The formula in (??) defines the $N(T)$ per \mathbb{C}_n probability if $\{N_m\}_{m=1}^N$ are i.i.d. with $N_m \sim \text{Pois}(T)$. The function $f(\mathbf{x}) = \exp\left(-\int_0^T \lambda(u)du\right)$. 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

$$\begin{aligned} & \exp\left(-\int_0^T \lambda(u)du\right) \\ & + \sum_{n=0}^{\infty} f\left(\left(t_1, \dots, t_N; N\right)\right) = \sum_{n=1}^{\infty} \int_{\mathbb{C}_n} \prod_{j=1}^n \frac{\left(-\int_{t_j}^T \lambda(u)du\right)}{\lambda(t_j)} dt_1 \cdots dt_n = 1 \end{aligned} \quad (2)$$

In the context of the classification problem (see [Ex. 1](#)) the class occurrence densities in [\(??\)](#) will be denoted $f_1(x)$ and $f_2(x)$ depending whether X comes from class ω_1 (denoted as $X \in \omega_1$) or if $X \in \omega_2$, respectively. The corresponding class intensities are $\lambda_1(t)$ and $\lambda_2(t)$ being the then negative functions C_n defined on $(0, \infty)$. Then using [\(??\)](#) one can form the Optimal Bayes rule which is defined over the simple regions $C_n = \{t_1, t_2, \dots, t_n\}$, $n \geq 1$, $t_i < t_{i+1}$, $i = 1, 2, \dots, n-1$. The Optimal Bayes rule defines the proper density over $\{C_n\}$, i.e., we have

$$\exp \left(\prod_{i=1}^N \frac{\lambda_1(t_i)}{\int_0^{t_i} \lambda_2(u) du} \exp \left(\int_0^{t_i} [\lambda_2(u) - \lambda_1(u)] du \right) \right) \geq \frac{\pi_2}{\pi_1}. \quad (4)$$

assuming that $N \geq 1$ and $\exp\left(\int_0^T [\lambda_2(u) - \lambda_1(u)] du\right) \geq \frac{\omega_2}{\omega_1}$ if $N = 0$. Clearly, if the reverse inequality in (??) holds, then we classify X to ω_2 . The log transform of (??) gives the alternative convenient form of the rule ψ_T^* , i.e., $X \in \omega_1$ if

In the context of the classification problem the class 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 $\psi_T^* \lambda_2 \mathbf{X}_T \in d\omega_1 + i \log \left(\frac{\pi_2}{\pi_1} \right)$. 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

$$\prod_{i=1}^n \frac{\lambda_1(t_i)}{\lambda_2(t_i)} \exp \left(\int_0^T [\lambda_2(u) - \lambda_1(u)] du \right) \geq \frac{\pi_2}{\pi_1}. \quad (4)$$

$$\prod_{i=1}^n \log \left(\int_0^{t_i} \lambda_1(t) dt \right) dN(t) > \gamma. \quad (6)$$

assuming that $N \geq 1$ and $\exp\left(\int_0^T [\lambda_2(u) - \lambda_1(u)] du\right) \geq \frac{\log\left(\frac{N}{\lambda_1(t_2)}\right)}{\pi_1}$ 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

$$\lambda(t) \sum_{i=1}^N \log \begin{cases} \lambda_1(t_i) & \text{if } \mathbf{X} \in \omega_1 \\ \frac{\lambda_1(t_i)}{\lambda_2(t_i)} & \text{if } \mathbf{X} \in \omega_2 \end{cases}. \quad (7)$$

where our further considerations it is useful to represent the class intensity function fully written in terms of the stochastic intensity factor and shape function [?]. Thus, let $\lambda_1(t) = \tau_1 p_1(t)$, $\lambda_2(t) = \tau_2 p_2(t)$, where 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

$$\tau_i = \int_0^T \lambda_i(u) du, \quad \frac{\lambda_1(t)}{\lambda_2(t)} = \lambda_i(u)/\tau_i, \quad i = 1, 2. \quad (8)$$

Clearly $p_1(t)$, $p_2(t)$ are well-defined probability density functions on $[0, T]$. The representation in (??) allows us 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

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

$$\sum_{i=1}^N \log\left(\frac{\lambda_1(t_i)}{p_2(t_i)}\right) \geq \eta, \quad (9)$$

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 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 $\psi_T^*; \mathbf{X} \in \omega_1$. Clearly $p_1(t)$, $p_2(t)$ are well-defined probability density functions on $[0, T]$. The representation in (??) allows us 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 $\geq \log\left(\frac{\pi_2}{\pi_1}\right)$ as follows, $\mathbf{X} \in \omega_1$ if

provided that $N \geq 1$. In the case $N = 0$ this reads $\sum_{i=1}^N \log\left(\frac{p_1(t_i)}{p_2(t_i)}\right) \geq \log\left(\frac{\pi_2}{\pi_1}\right)$. The risk associated with the rule $\psi_T^*(\mathbf{x})$ in (??) (or (??)) is defined as $R_T^* = P(\psi_T^*(\mathbf{X}) \neq Y)$ and is referred as the Bayes risk. Here $Y \in \{\omega_1, \omega_2\}$ is the true class label of \mathbf{X} . For our future studies we express the Bayes risk in terms of the decision function $W_T(\mathbf{X})$, i.e., we write where $\eta = \tau_1 - \tau_2 + N \log\left(\frac{\pi_2}{\pi_1}\right) + \log\left(\frac{\pi_2}{\pi_1}\right)$. 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 $\lambda_2(t) = \lambda_2$, we have the homogeneous spike train data then the Bayes rule takes the following form $\psi_T^*; \mathbf{X} \in \omega_1$. 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 Sections ?? 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

provided that $N \geq 1$. In the case $N = 0$ this reads as $\lambda(t) dt \geq 1/\log\left(\frac{\pi_2}{\pi_1}\right)$. The risk associated with the rule $\psi_T^*(\mathbf{x})$ in (??) (or (??)) is defined as $R_T^* = P(\psi_T^*(\mathbf{X}) \neq Y)$ and is referred as the Bayes risk. Here $Y \in \{\omega_1, \omega_2\}$ is the true class label of \mathbf{X} . For our future studies we express the Bayes risk in terms of the decision function $W_T(\mathbf{X})$, i.e., The formula in (??) can be viewed as the local signal plus noise decomposition, where the noise process $dM(t)$ reveals the local martingale structure [?]. Appendix A gives the pertinent results concerning the martingale decomposition of the underlying spiking process.

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

The decomposition in (??) allows us to express $P(\text{Classification} | \mathbf{X} \in \omega_1)$ (or its version in (??)) in the convenient stochastic integral form. In fact, by virtue of (??) and (??) we write the left-hand side of (??) as R_T^* as 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 Sections ?? 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

$$dN(t) = \lambda(t) dt + dM(t) \frac{\lambda_1(t)}{\lambda_2(t)}, \quad (13)$$

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 local signal plus noise decomposition, where the noise process $dM(t)$ 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.

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

III. THE BAYES RULE AND RISK: BOUNDS AND ASYMPTOTIC BEHAVIOR

A. The Bayes Decision Function $\int_0^T \log\left(\frac{\lambda_1(t)}{\lambda_2(t)}\right) dN(t) = \int_0^T \log\left(\frac{\lambda_1(t)}{\lambda_2(t)}\right) \lambda(t) dt$

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

$$\int_0^T \log\left(\frac{\lambda_1(t)}{\lambda_2(t)}\right) dM(t) \geq \eta_T, \quad (14)$$

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. $\int_0^T g(t) dM(t)$. (15)

Here $g(t) = \log\left(\frac{\lambda_1(t)}{\lambda_2(t)}\right) = \text{Hd}\left(\frac{\lambda_1(t)}{\lambda_2(t)}\right)$ Bayes Rule and Risk: Bounds and Asymptotic Behavior

A. The Bayes Decision Function

In this section we examine the optimal decision function derived in (??) in 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 (??), $\mathbf{X} \in \omega_{16}$

$$U_T(\mathbf{X}) \geq \alpha_T + \log\left(\frac{p_2(t)}{p_1(t)}\right)\left(\frac{\lambda_2(t)}{\pi_1}\right) dt, \quad (14)$$

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 present under stochastic integration. Hence, since $\mathbb{E}[dM(t)] = 0$ the process

$$\text{Here } g(t) = \log\left(\frac{\lambda_1(t)}{\lambda_2(t)}\right) = \log\left(\frac{p_1(t)}{p_2(t)}\right) + \log\left(\frac{\tau_2}{\tau_1}\right) \text{ and } \begin{cases} U_t(\mathbf{X}) = \int_0^t g(u)dM(u), 0 \leq t \leq T \\ \alpha_T = \tau_1 - \tau_2 + \log\left(\frac{\tau_2}{\tau_1}\right) \int_0^T \lambda(t)dt \end{cases}$$

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

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 $\mathbb{E}[dM(t)] = 0$ the process. 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

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

It is known that $K_T(p \parallel q) \geq 0$ and $K_T(p \parallel q) = 0$ if $p = q$. 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

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.

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

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

If $\mathbf{X} \in \omega_1$ then

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

$$K_T(p \parallel q) = \int_0^T \alpha_T \log\left(\frac{p(t)}{q(t)}\right) K_T(p_t \parallel p_2) dt. \quad (19)$$

(b) If $\mathbf{X} \in \omega_2$ then $K_T(p \parallel q) \geq 0$ and $K_T(p \parallel q) = 0$ 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

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

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

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 $[\frac{\tau_1 K_T(p_1 \parallel p_2)}{\tau_2}, \frac{(\tau_1 - \tau_2)^2}{\tau_1}]$ if $\mathbf{X} \in \omega_1$ and $\mathbf{X} \in \omega_2$, respectively. It is also worth noting that $(\tau_1 - \tau_2)^2 \leq \left\{ \int_0^T (\lambda_1(t) - \lambda_2(t)) dt \right\}^2$ 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 $\alpha_T = -\tau_1 K_T(p_1 \parallel p_2)$ for $\mathbf{X} \in \omega_1$ and ω_2 then $K_T(p_2 \parallel p_1)$ for $\mathbf{X} \in \omega_2$.

The next result concerns 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 we evaluate the basic statistical feature of this term by deriving its variance.

Lemma 2. Let us consider the stochastic part $U_T(\mathbf{X})$ of the Bayes rule in (??). Then,

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

(a) As if 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 $\text{Var}[U_T(\mathbf{X})]$ 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. It is also worth noting that $\int_0^T \left\{ \log \left(\frac{p_1(t)}{p_2(t)} \right) + \log \left(\frac{\lambda_1(t)}{\lambda_2(t)} \right) \right\} p_2(t) dt$ represents the square of the difference of the average number of events on $[0, T]$ coming from classes $p_1(t)$ and ω_2 . As a result, if $\tau_1 = \tau_2$ then $\alpha_T = -\tau_1 K_T(p_1 \parallel p_2)$ for $\mathbf{X} \in \omega_1$ and $\alpha_T = \tau_2 K_T(p_2 \parallel p_1)$ for $\mathbf{X} \in \omega_2$.

(b) The next result concerns 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 we evaluate the basic statistical feature of this term by deriving its variance.

Lemma 2. Let us consider the stochastic part $U_T(\mathbf{X})$ of the Bayes rule in (??). Then,

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

$$\mathbf{V}_T(p \parallel q) = \tau_1 \int_0^T \left\{ \log \left(\frac{p_1(t)}{p_2(t)} \right) + \log \left(\frac{\lambda_1(t)}{\lambda_2(t)} \right) \right\}^2 p_2(t) dt. \quad (21)$$

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

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

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

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

$$\mathbf{V}_T(p \parallel q) = \tau_1 \left\{ \mathbf{V}_T(p_1 \parallel p_2) \left(\frac{p(t)}{q(t)} \right) p(t) dt + 2 \log \left(\frac{\tau_1}{\tau_2} \right) \mathbf{K}_T(p_1 \parallel p_2) + \log^2 \left(\frac{\tau_1}{\tau_2} \right) \right\}. \quad (22)$$

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

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 of

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

$$T U_T(\mathbf{X}) = \frac{1}{T} \int_0^T g(t) dM(t) \quad (26)$$

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

$$A2 : \frac{1}{T} \int_0^T \lambda_i(u) du \rightarrow d, \quad i = 1, 2 \text{ as } T \rightarrow \infty. \quad (27)$$

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

The meaning of this condition is that the average number of events from each class increases linearly with T . It is worth noting that interesting questions 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 the limit behavior of $\text{Var}[U_T(\mathbf{X})]$ as $T \rightarrow \infty$. It is clear that the limit behavior 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 $\mathbf{X} \in \omega_1$

$$\frac{1}{T} U_T(\mathbf{X}) = \frac{1}{T} \int_0^T g(t) dM(t) \quad (26)$$

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

$$A2 : \frac{1}{T} \int_0^T \lambda_i(u) du \rightarrow d, \quad i = 1, 2 \text{ as } T \rightarrow \infty. \quad (27)$$

The meaning of this condition is that the average number of events from 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 $\text{Var}[U_T(\mathbf{X})]$ as $T \rightarrow \infty$. It is clear that such

The right-hand side of this inequality is equal to $\left(\int_0^T \log \left(\frac{\lambda_1(t)}{\lambda_2(t)} \right) p_1(t) dt \right)^2$ and the upper bound is not smaller than $\text{Var}[U_T^2(\mathbf{X})]$. In fact, recalling (22) and (23) we have that if $\mathbf{X} \in \omega_1$

$$\text{Var}[U_T^2(\mathbf{X})] \leq \text{Var} \int_0^T \left[C \log \left(\frac{\lambda_1(t)}{\lambda_2(t)} \right) \log \left(\frac{p_1(t)}{\delta} \right) \right]^2 dt. \quad (29)$$

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

$$d \log^2 \left(\frac{\delta}{C} \right) \leq \tau_1 \log^2 \left(\frac{C}{\delta} \right). \quad (28)$$

On the other hand by (22) and Lemma 3 we get

$$\text{Var}[U_T(\mathbf{X})] \leq \tau_2 \log^2 \left(\frac{C}{\delta} \right). \quad (30)$$

The bounds in (22), (23) and the $\text{Var}[U_T(\mathbf{X})]$ in (28) lead to the following limit behavior of $\text{Var}[U_T(\mathbf{X})]$.

Lemma 4. Let the assumptions **A1**, **A2** hold. Then for $\mathbf{X} \in \omega_1$ or $\mathbf{X} \in \omega_2$ we have

$$d \log^2 \left(\frac{\delta}{C} \right) \leq \lim_{T \rightarrow \infty} \text{Var} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right] + 2 \log \left(\frac{\tau_1}{\tau_2} \right) K_T(p_1 \parallel p_2) + \log^2 \left(\frac{\tau_1}{\tau_2} \right).$$

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 (23) this is not smaller than $\tau_1 \log^2 \left(\frac{\delta}{C} \right)$. Hence, if $\mathbf{X} \in \omega_1$ this gives the following bounds

The question whether the inferior and superior limits in (22) are equal remains open. It should be noted that if (22) 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 (22) 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 (22) 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

Lemma 4. Let the assumptions **A1**, **A2** hold. Then for $\mathbf{X} \in \omega_1$ or $\mathbf{X} \in \omega_2$ we have

$$d \log^2 \left(\frac{\delta}{C} \right) \leq \lim_{T \rightarrow \infty} \text{Var} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right], \quad (32)$$

whereas if $\mathbf{X} \in \omega_2$ then

$$\lim_{T \rightarrow \infty} \text{Var} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right] = d \log^2(\mu) \log^2 \left(\frac{C}{\delta} \right). \quad (33)$$

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 $\mu \rightarrow \infty$. It should be noted that if (22) 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 (22) let us consider the following simple example defined in (22).

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 the conditions of Lemma ?? (a)-(d). This implies that the condition in (22) or 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

$$\frac{1}{T} U_T(\mathbf{X}) = \frac{1}{T} \int_0^T \log \left(\frac{\lambda_1(t)}{\lambda_2(t)} \right) dM(t) \rightarrow 0 \quad (P) \quad (34)$$

$$\lim_{T \rightarrow \infty} \text{Var} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right] = d \log^2(\mu), \quad (32)$$

whereas if $\mathbf{X} \in \omega_2$ then

$$\text{The part of } \mathbf{X} \in \omega_2 \text{ then is a direct application of Lemma ?? and the Chebyshev inequality. In fact, let us consider the case } \mathbf{X} \in \omega_1. \text{ Then, for any } \epsilon > 0 \text{ we have } \lim_{T \rightarrow \infty} \text{Var} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right] = d \mu \log^2(\mu). \quad (33)$$

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 $\mu \rightarrow \infty$. The right-hand side of (22) is equal to $\text{Var} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right] / T \epsilon^2$, where due to (22) the limit superior of $\text{Var} \left[\frac{1}{\sqrt{T}} U_T(\mathbf{X}) \right]$ is bounded by a finite constant. This confirms the claim of Theorem ??.

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 (22). 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.

$$\frac{1}{T} U_T(\mathbf{X}) = \frac{1}{T} \int_0^T \log \left(\frac{\lambda_1(t)}{\lambda_2(t)} \right) dM(t) \rightarrow 0 \quad (P) \quad (34)$$

Lemma 5. Let $N(t)$ be the counting process allowing the decomposition in (22). 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

- (a) $|g(t)| \leq u_T$ for all $t \in [0, T]$. Then for \mathbf{X} coming either from class ω_1 or class ω_2 we have
- (b) $\int_0^T g(t) \lambda_1(t) dt \leq c_T$. Then for any $\epsilon > 0$ we have

$$\mathbb{P} \left(\frac{1}{T} |U_T(\mathbf{X})| \geq \epsilon \right) \leq \frac{\text{Var}[U_T(\mathbf{X})]}{T^2 \epsilon^2}. \quad (35)$$

The right-hand side of (??) is finite and bounded by a finite constant $\text{Var}\left[\frac{1}{\sqrt{T}}U_T(\mathbf{X})\right]$. Then, for each $\epsilon > 0$ we have due to (??) the limit superior of $\text{Var}\left[\frac{1}{\sqrt{T}}U_T(\mathbf{X})\right]$ is bounded by a finite constant. This confirms the claim of Theorem ??.

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 or 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 ?? can be directly applied for the evaluation of the stochastic integral in (??). In fact, with $g(t) = \log\left(\frac{\lambda_1(t)}{\lambda_2(t)}\right)$ and by the assumption **A1**, we have that $g(t)$ is a process (??). Hence, the decomposition in Lemma ?? is met with U_T met with $\int_0^T g(t)dM(t)$ (??) for the stochastic integrand of the probability (??) function in Appendix A with respect to the conditional measure $M(t)$. Let us now suppose that

- $|g(t)| \leq u_T$ for all $t \in [0, T]$.
 - $\int_0^T g^2(t)\lambda(t)dt \leq v_T$,
- $$\text{Var}[U_T(\mathbf{X})] = T\text{Var}\left[\frac{1}{\sqrt{T}}U_T(\mathbf{X})\right] = T\theta_T,$$

where u_T and v_T are some finite constants. Then, for each $\epsilon > 0$ we have

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 $T > 0$.

Lemma 6. Suppose that the assumption **A1** holds for any finite T . If \mathbf{X} coming either from class ω_1 or class ω_2 and every $\epsilon > 0$ we have

Lemma ?? can be directly applied for the evaluation of the stochastic integral in (??). In fact, with $g(t) = \log\left(\frac{\lambda_1(t)}{\lambda_2(t)}\right)$ and by the assumption **A1**, we have that $|g(t)| \leq \log\left(\frac{C}{\delta}\right)$. Hence, the condition (a) in Lemma ?? is met with $u_T = \log\left(\frac{C}{\delta}\right)$ for all $T > 0$. By virtue of the property (??) in Appendix A the integral in the condition (b) of Lemma ?? reads as

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 it by finding a sequence of times T_n , such that $T_n \rightarrow \infty$ and then employ the standard Borel-Cantelli lemma. We refer to [?] for details for such a discretization strategy. The bound is valid for any finite $T > 0$.

Theorem 2. Let the assumptions **A1** and **A2** hold. Then for \mathbf{X} coming either from class ω_1 or class ω_2 and every $\epsilon > 0$ we have

$$\frac{1}{T}\Pr\left(\frac{1}{T}|U_T(\mathbf{X})| \geq \epsilon\right) \leq 2\exp\left[-\frac{T\theta_T - \epsilon^2}{2\theta_T + \epsilon}\right] \quad (\text{a.s.}),$$

as $T \rightarrow \infty$, where $u = \log\left(\frac{C}{\delta}\right)$ and the factor θ_T is defined in (??).

B. The Bayes Risk 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 a discretization strategy.

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

$$\frac{1}{T}U_T(\mathbf{X}) = \frac{1}{T}\Pr\left(U_T(\mathbf{X}) \geq \alpha_T \left|\frac{\lambda_1(\tau_1)}{\lambda_2(\tau_1)}\right| \log\left(\frac{\pi_2}{\pi_1}\right) \mid \mathbf{X} \in \omega_2\right),$$

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

$$\alpha_T = \tau_1 - \tau_2 + \tau_2 \log\left(\frac{\tau_2}{\tau_1}\right) + \tau_2 \mathbf{K}_T(p_2 \parallel p_1). \quad (41)$$

B. The Bayes Risk

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 (??). Owing to (??) it suffices to consider the probability of misclassification $\Pr(W_T(\mathbf{X}) \geq \eta_T \mid \mathbf{X} \in \omega_2)$. The Hence, we have the following convergence result that also gives the upper bound for the Bayes risk analysis of the probability $\Pr(W_T(\mathbf{X}) < \eta_T \mid \mathbf{X} \in \omega_1)$ is analogous. By virtue of (??) we can write

Theorem 3. Let the assumptions **A1** and **A2** hold. Then we have

$$= \Pr\left(U_T(\mathbf{X}) \geq \alpha_T + \log\left(\frac{\pi_2}{\pi_1}\right) \mid \mathbf{X} \in \omega_2\right), \quad (40)$$

Furthermore,

$$\text{where } U_T(\mathbf{X}) \text{ is defined in (??) and } \alpha_T \text{ (under the fact that } \mathbf{X} \in \omega_2\text{)} \leq (\pi_1 a_T + \pi_2 b_T) \frac{1}{T}, \quad (42)$$

$$\text{for some finite constants } a_T, b_T. \quad \alpha_T = \tau_1 - \tau_2 + \tau_2 \log\left(\frac{\tau_2}{\tau_1}\right) + \tau_2 \mathbf{K}_T(p_2 \parallel p_1). \quad (41)$$

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 (??).

Remark 1. Hence we have that the assumptions **A1** and **A2** the Bayes risk tends under the assumption **T** to zero. The proof of Theorem ?? directs a set of useful tools of the framework of large deviation theory established in Theorem ??, see (??).

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

Theorem 3. Let the assumptions **A1** and **A2** hold. Then, we have $c_1 = \frac{1}{d} \left(\frac{\log(C/\delta)}{\log(\delta/C)} \right)^2$. (43)

Hence, for large T one can write $R_T^* \prec c_1 \frac{1}{T}$. $R_T^* \rightarrow 0$ as $T \rightarrow \infty$.

Furthermore, if the result of Lemma ?? we can substantially improve the bound in (??). Hence, we have the following result.

Theorem 4. Let the assumptions **A1** and **A2** hold. Then, we have $R_T^* \leq (\pi_1 a_T + \pi_2 b_T) \frac{1}{T}$, (42)

for some finite constants a_T, b_T . $R_T^* \leq \pi_1 \exp[-A_T T] + \pi_2 \exp[-B_T T]$, (44)

for the proof of Theorem ?? B_T 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 $U_T(X)$ in (??).

The proof of Theorem ?? is deferred to Appendix B, where also the explicit expressions for A_T and B_T are presented.

Remark 1. Hence under the assumptions **A1** and **A2** the Bayes risk tends to zero with the rate $1/T$. 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

$$c_2 = d \frac{1}{3} \left(\frac{\log(A/\delta)}{\log(\delta/C)} \right)^2, \quad (45)$$

Hence, for large T one can write $R_T^* \prec c_2 \frac{1}{T}$; 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 .

Example 4. Consider the assumptions **A1** and **A2** hold. Then, we have ?? Then, using the results in (??) and (??) and some algebra we can show the following counterpart of the result of Theorem ??

$$R_T^* \leq \pi_1 \exp[-A_T T] + \pi_2 \exp[-B_T T], \quad (44)$$

$$R_T^* \prec \pi_1 \exp[-c_1(\mu) d T] + \pi_2 \exp[-c_2(\mu) d T]. \quad (46)$$

for some finite constants A_T, B_T . The asymptotic constants $c_1(\mu), c_2(\mu)$ can be written in the explicit form and they obey the following properties

The proof of Theorem ?? is deferred to Appendix B, where also the explicit expressions for A_T and B_T are presented.

$$\lim_{\mu \rightarrow \infty} c_1(\mu) = \lim_{\mu \rightarrow \infty} c_2(\mu) = 0$$

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

$$c_2 = d \frac{1}{3} \left(\frac{\log(\delta/C)}{\log(\delta/C)} \right)^2 = \infty. \quad (45)$$

The former limit corresponds to the indistinguishable case, i.e. $\log(C/\delta) = \lambda_2(t)$. On the other hand, the latter limit exhibits that where δ/C then $R_T^* \rightarrow 0$. Again 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 .

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

$$R_T^* \leq \sqrt{\pi_1 \pi_2} \exp(-\beta(T)),$$

Example 2. Consider the classification problem discussed in Example ?? Then, using the results in (??) and (??) and some algebra we can show the following counterpart of the result of Theorem ?? is the classical Bhattacharya bound [?] 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 $c(\mu) d T$, where

The asymptotic constant $c(\mu)$ can be written in the explicit form appearing at the bottom of the following properties

Remark 4. 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 growth of $\lambda_1(t), \lambda_2(t)$ one could use the multiplicative model due to Aalen [?], i.e. we consider $\lim_{\mu \rightarrow \infty} c_1(\mu) = \lim_{\mu \rightarrow \infty} c_2(\mu) = \infty$.

The former limit corresponds to the indistinguishable case, i.e., $\lambda_i(t) = \lambda_1(t) = \lambda_2(t)$. On the other hand, the latter limit exhibits that if $\mu \rightarrow \infty$ then $R_T^* \rightarrow 0$. Again the assumption **A1** is not needed here.

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

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

Example 3. Let us model the class intensities in the following form, 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 factor. This is the classical Bhattacharya bound [?] 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 $c(\mu) d T$, where $c(\mu) = (\mu + 1)/2 - \sqrt{\mu}$. Interestingly $c(\mu) \geq c_1(\mu), c_2(\mu)$, where $c_1(\mu), c_2(\mu)$ appear in our bound in (??). (48)

$$+ 0.5 \cos \left(\frac{\pi}{3\sqrt{2}} t + \frac{\pi}{4} + \phi \right)$$

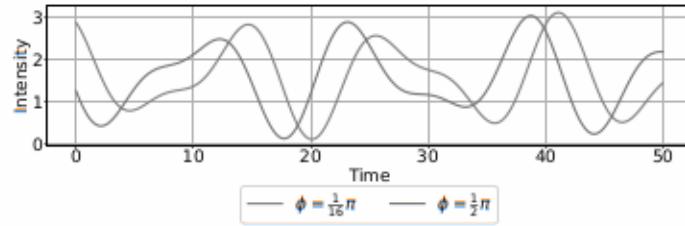


Figure 11. Intensity functions $\lambda_1(t) = (t\lambda_1^*(t) \frac{\pi}{16})^\phi$ and $\lambda_2(t) = (t\lambda_2^*(t) \frac{\pi}{2})^\phi$.

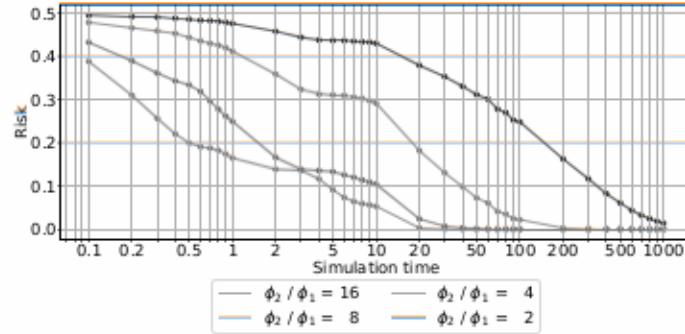


Figure 22. The Bayes risk R_T versus risk T for a two-class classification problem.

Remark 3. The order of convergence of the Bayes risk R_T depicted zero is determined by the condition in **A2**. This is due to the fact that the Bayes risk tends to zero as T gets larger. If **A2** does not hold intensities $\lambda_i(t)$, $i = 1, 2$ exponentially increase and the convergence of R_T to zero is impossible. This is in stark contrast to the convergence of intensities $\lambda_i(t)$ in one hand, whereas the fast rate of convergence is considered 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.

where $\gamma_i(t)$, $i = 1, 2$ are fixed functions and d is a parameter that is allowed to grow. It is an interesting alternative to derive the results obtained in this paper on nonparametric multiplicative class intensity model in (??).

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 is the empirical counterpart of the optimal Bayes rule in (??) or equivalently in (??). We have already pointed out that the single-sample based intensity function estimate cannot be consistent unless there is a certain mechanism that makes the intensity function increase, e.g., the multiplicative model in (??). In this paper we consider the intensity model based on the increasing number of replicates of the class spiking processes. Hence, contrary to the results of Section ?? the observation interval $[0, T]$ is kept constant.

Various choices of ϕ define $\lambda_i(t)$, $i = 1, 2$. Figure ?? depicts $\lambda_1(t) = \lambda_1^*(t; \frac{\pi}{16})^\phi$ and $\lambda_2(t) = \lambda_2^*(t; \frac{\pi}{2})^\phi$.

Figure ?? illustrates the fact that the Bayes risk tends to zero as T gets larger. The model of class intensities defined labeled spiking processes (\mathbf{X}_j, Y_j) . Here \mathbf{X}_j is the variable-length vector, i.e., $\mathbf{X}_j = [t^{(j)}_1, t^{(j)}_2, \dots, t^{(j)}_{N^{(j)}}]^T$ and $Y_j \in \{\omega_1, \omega_2\}$, in (??) is parametrized by ϕ , i.e., we set $\lambda_1(t) = \lambda_1^*(t; \phi_1)$ and $\lambda_2(t) = \lambda_2^*(t; \phi_2)$. 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.

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 $\hat{\pi}_i = L_i/L$ and $\hat{\pi}_i = I_i/I$. In order to estimate $\{(\tau_i, p_i(t)), i = 1, 2\}$ 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 one can form the unbiased estimate

A. Plug-in Classifiers

of τ_i as $\hat{\tau}_i^{(j)} = N^{(j)}$. However, $\text{Var}[N^{(j)}|Y_j = \omega_i] = \tau_i$ and this is an inconsistent estimate of τ_i . The latter fact results from In practice one does not know the true class intensities functions and must rely on some training data in order to form the local Poisson behavior of the spiking process, see Appendix A. Nevertheless, the aggregation of $\{\hat{\tau}_i^{(j)}\}$ leads to consistent a data-driven classification rule. In this paper we apply the plug-in strategy to design a classifier, i.e. the classifier that estimate of τ_i for the increased size of the training set. Hence, let

is the empirical counterpart of the optimal Bayes rule in (??) or equivalently in (??). We have already pointed out that the single-sample based intensity function estimate cannot be consistent unless there is a certain mechanism that makes the intensity function increase, e.g., the multiplicative model in (??). In this paper we consider the intensity model based on the increasing number of replicates of the class spiking processes. Hence, contrary to the results of Section ?? the observation interval $[0, T]$ is kept constant.

be to estimate $\mathbf{D}_L = \{(\mathbf{X}_i, Y_i)\}_{i=1}^L$ in the analogous way the learning with the probing of estimating L independent observations of the nonparametric estimate of $p_i(t)$ based on the single sample \mathbf{X}_j from the class ω_i . Then, the aggregated estimate of $p_i(t)$ takes the following form

$$\widehat{p}_i(t) = \frac{1}{L} \sum_{j=1}^L p_i^{[j]}(t) \mathbf{1}(Y_j = \omega_i), \quad i = 1, 2. \quad (50)$$

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 $\widehat{\pi}_L = L/N$ and $\widehat{\pi}_2 = L_2/L$. In order to estimate $\{(\tau_i, p_i(t)), i = 1, 2\}$ 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 one can form the unbiased estimate of τ_i as $\widehat{\tau}_i^{[j]} = N^{[j]}/L$. However, $\text{Var}[N^{[j]}|Y_j = \omega_i] = \tau_i$ and this is an inconsistent estimate of τ_i . The latter fact results from the local Poisson behavior of the InSituIn process, see Appendix A. Nevertheless, the aggregation of this leads to consistent estimate of τ_i for the increased size of the training set. Hence 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 first result reveals that this probability tends to one if the training set tends to infinity with the problem of estimating $p_i(t)$. Let $\widehat{p}_i^{[j]}(t)$ 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 of $p_i(t)$ takes the following form

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

Then,

$$\widehat{p}_i(t) = \frac{1}{L_i} \sum_{j=1}^{L_i} \widehat{p}_i^{[j]}(t) \mathbf{1}(Y_j = \omega_i), \quad i = 1, 2. \quad (50)$$

$$\mathbf{P}(\widehat{\psi}_{L,T}(\mathbf{x}) = \psi_T^*(\mathbf{x})) \rightarrow 1$$

Plugging (??) and (??) into (??) gives us the following empirical classification rule $\widehat{\psi}_{L,T}$: classify $\mathbf{X} = [t_1, \dots, t_N; N] \in \omega_1$ as $L \rightarrow \infty$. The proof of Theorem ?? is given in Appendix C. 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 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 first result reveals that this probability tends to one if the size of the training set tends to infinity.

The condition in (??) of Theorem ?? assures that the decision 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.

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

Then,

$$\mathbf{P}(\widehat{W}_{L,T}(\mathbf{x}) \rightarrow W_T(\mathbf{x})) \rightarrow 1 \quad (54)$$

as $L \rightarrow \infty$.

as $L \rightarrow \infty$. The proof of Theorem ?? is given in Appendix C. 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

$$\rho(\widehat{\psi}_{L,T}(\mathbf{x}), \psi_T^*(\mathbf{x})) \rightarrow 0 \quad (P) \quad (55)$$

as $L \rightarrow \infty$, where

The local consistency of $\widehat{\psi}_{L,T}$ leads to the global convergence characterized by the conditional risk. Hence, let $\mathbf{P}(\widehat{\psi}_{L,T}(\mathbf{x}) \neq Y | \mathbf{D}_L) = \mathbb{E}[\mathbf{1}(\widehat{\psi}_{L,T}(\mathbf{x}) \neq Y) | \mathbf{D}_L]$ be the conditional risk associated with the rule $\widehat{\psi}_{L,T}$. Since the condition in (??) of Theorem ?? assures that the decision 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.

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, the above is bounded by

$$\mathbb{E}[\widehat{W}_{L,T}(\mathbf{x}), W_T(\mathbf{x}) | \mathbf{D}_L]. \quad (54)$$

Owing to (??) and Lebesgue's dominated convergence theorem we obtain the main result of this section.

Theorem 6. Consider the class of plug-in classifiers defined in Theorem ?? Suppose that (B) conditions of the respective method. Then, we have the following Bayes risk consistency characterization by the threshold value $\hat{\eta}_{L,T}$. 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 as $L \rightarrow \infty$.

$$\hat{\eta}_{L,T} \rightarrow \eta_T \quad (P) \quad (55)$$

B. Kernel Classifiers

The local consistency of $\hat{\psi}_{L,T}$ leads to the global convergence characterized by the conditional risk. Hence, let it be 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 $R_T^* = E[1(\psi_T^*(X) \neq Y)]$, then one can write

$$0 \leq R_{L,T} - \hat{R}_T^*(t) = \sum_{l=1}^{N^{[j]}} K_h(t - t_l^{[j]}) , \quad (57)$$

$$= E[1(\hat{\psi}_{L,T}(X) \neq Y) - 1(\psi_T^*(X) \neq Y) | D_L] .$$

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

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

Owing to (??) and Lebesgue's dominated convergence theorem we obtain the main result of this section.

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

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

$$R_{L,T} \rightarrow R_T^*(P) \quad (56)$$

as $L \rightarrow \infty$.

$$\hat{p}_i^{[j]}(t) = \frac{1}{N^{[j]}} \sum_{l=1}^{N^{[j]}} K_h(t - t_l^{[j]}) .$$

B. Kernel Classifiers

As we have already pointed in Section ?? the estimates $\hat{\lambda}_i^{[j]}(t)$, $\hat{p}_i^{[j]}(t)$ 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 $\{\hat{\lambda}_i^{[j]}\}$, $\{\hat{p}_i^{[j]}\}$. This leads to the following aggregated kernel estimate of $p_i(t)$

$$\hat{\lambda}_i^{[j]}(t) = \sum_{l=1}^{N^{[j]}} K_h(t - t_l^{[j]}) , \quad (57)$$

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

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

Moreover, the aggregated estimate $\hat{\tau}_i$ of τ_i is defined in (??). Plugging $\hat{p}_i(t)$ and $\hat{\tau}_i$, $i = 1, 2$ into (??) we obtain the kernel classification rule. The aggregated kernel estimate $\hat{\lambda}_i(t)$ 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 $\hat{p}_i(t)$ and $\hat{\lambda}_i(t)$ is the same and therefore we can verify the requirement in (??) for the kernel intensity estimate. Hence, with some abuse of the notation let $\{X_1, X_2, \dots, X_L\}$ be the data set from the fixed

These crucial tuning parameters h are called (the characterized by the controls the intensity function of this) Thus, one observes the kernel $K_h(t)$, $\{N^{[j]}(t)\}$ of the counting process $N(t)$, where $N^{[j]}(t)$ is represented by the feature vector $X_j = [t_1^{[j]}, \dots, t_{N^{[j]}}^{[j]}; N^{[j]}]$

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

$$dN^{[j]}(t) = \lambda(t) dt + dM^{[j]}(t), \quad j = 1, \dots, L .$$

$$\hat{p}_i^{[j]}(t) = \frac{1}{N^{[j]}} \sum_{l=1}^{N^{[j]}} K_h(t - t_l^{[j]}) .$$

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

As we have already pointed in Section ?? the estimates $\hat{\lambda}_i^{[j]}(t)$, $\hat{p}_i^{[j]}(t)$ cannot be consistent by merely increasing T (??) overcome this problem one can utilize the observed multiple training vectors and aggregate the single-sample estimates $\{\hat{\lambda}_i^{[j]}\}$, $\{\hat{p}_i^{[j]}\}$. This leads to the following aggregated kernel estimate of $p_i(t)$

$$d\bar{N}_L(t) = \frac{1}{L} \sum_{l=1}^L dN^{[j]}(t) ,$$

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

$$d\bar{M}_L(t) = \frac{1}{L} \sum_{l=1}^L dM^{[j]}(t) .$$

Moreover, the aggregated estimate $\hat{\tau}_i$ of τ_i is defined in (??). Plugging $\hat{p}_i(t)$ and $\hat{\tau}_i$, $i = 1, 2$ into (??) we obtain the kernel classification rule. The aggregated kernel estimate $\hat{\lambda}_i(t)$ 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 $\hat{p}_i(t)$ and $\hat{\lambda}_i(t)$ is the same and therefore we can verify the

require important (??) for that the aggregated residual process $\bar{N}_L(t)$ meets all the properties listed in Appendix A. Hence, the data set from the properties in (??) and 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 $\text{Var}[d\bar{M}_N(t)] = \lambda(t)dt$. The local martingale decomposition in (??) for $N^{[j]}(t)$ reads

$$\text{Var}\left[\int_0^T g(u)\lambda(u)du + dM^{[j]}(t)\right] = jg^2(u)\lambda(u)du. \quad (60)$$

This gives the analogous decomposition for the aggregated counting process, i.e., we have
The single-sample kernel estimate of $\lambda(t)$ is as in (??), whereas its aggregated version takes the form

$$d\bar{N}_L(t) = \lambda(t)dt + d\bar{M}_L(t), \quad (59)$$

where

$$\hat{\lambda}(t) = \frac{1}{L} \sum_{j=1}^L \hat{\lambda}^{[j]}(t). \quad (61)$$

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

$$\hat{\lambda}(t) = \int_0^T K_1(t-s) d\bar{N}_L(s), \quad d\bar{M}_L(t) = \frac{1}{L} \sum_{j=1}^L dM^{[j]}(t). \quad (62)$$

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 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

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

$$\begin{aligned} \text{Var}[d\bar{M}_L(t)] &= \frac{1}{L} \lambda(t)dt, \\ \text{Var}\left[\int_0^T g(u)d\bar{M}_L(u)\right] &= \frac{1}{L} \int_0^T \frac{1}{h} K^2\left(\frac{t-s}{h}\right) \lambda(s)ds. \end{aligned} \quad (64)$$

These formulas and the standard analysis developed in the context of kernel estimates [?], [?] reveal that if
The single-sample kernel estimate of $\lambda(t)$ is as in (??) and whereas its aggregated version takes the form

then

$$\hat{\lambda}(t) \xrightarrow{L \rightarrow \infty} \lambda(t) \text{ as } L \rightarrow \infty. \quad (65)$$

for $t \in (0, T)$ 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 $K_h(t) = \int_0^T K(t-s) d\bar{N}_L(s)$. 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 along with the properties of the estimate $\hat{\lambda}(t)$ in (??) and the properties of $d\bar{M}_L(t)$ (see (??)) yield the following identities for the bias and the variance of $\hat{\lambda}(t)$ in [0, ∞). Let the kernel function $K(t)$ be Lipschitz continuous on $[-1, 1]$. 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

$$h(E)[\hat{\lambda}(t)] \text{ and } \int_0^T \frac{1}{h} K\left(\frac{t-s}{h}\right) \lambda(s)ds \rightarrow \infty. \quad (66)$$

Then for arbitrarily small $\epsilon > 0$

$$\text{Var}\left[\sup_{t \in [\epsilon, T-\epsilon]} \left| \frac{1}{L} \int_0^T \frac{1}{h} K^2\left(\frac{t-s}{h}\right) \lambda(s)ds \right| \right] \rightarrow 0 \text{ as } L \rightarrow \infty. \quad (67)$$

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

It is worth noting that the uniform convergence holds under the condition $Lh^3(L) \rightarrow \infty$. This is the stronger restriction than the one required for the pointwise convergence, where one needs that $Lh(L)Lh'(L) \rightarrow \infty$. We conjecture that (??) can be replaced by the weaker condition $Lh(L)/\log(L) \rightarrow \infty$. 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 have This for the pointwise convergence that holds at interior points of $[0, T]$ is known [?], [?] that the convergence fails at the boundary points near $t = 0$, $t = T$. This enforces us to confine the Theorem 7. Let the class intensities $\lambda_1(t), \lambda_2(t)$ satisfy the conditions of Lemma ?? If (??) holds, then the kernel classification rule is Bayes risk consistent, i.e., we have

R_{L,T} → R_T^{*}(P) that are able to restore the convergence property at the boundary points. The following lemma gives the sufficient conditions for the uniform convergence property of the estimate $\hat{\lambda}(t)$ in (??).

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

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.

Then for arbitrarily small $\epsilon > 0$

$$\sup_{t \in [\epsilon, T-\epsilon]} |\hat{\lambda}(t) - \lambda(t)| \rightarrow 0 \text{ (P) as } L \rightarrow \infty. \quad (68)$$

The selection of the bandwidth h of the nonparametric classifier under training on the $Lh(L)$ sample affords the stronger classification rule than standard $\text{rule}(P)$ if it is compared with the expression in (22) and (23) so that $Lh(L)/L \rightarrow 0$. We conjecture that (22) holds for $b \in (0, L)$ if only the weaker condition $Lh(L)/\log(L) \rightarrow \infty$ 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.

Theorem 7. Let the class intensities $\lambda_1(t), \lambda_2(t)$ satisfy the conditions of Lemma ???. If (??) holds, then the kernel classification rule is Bayes risk consistent, i.e., we have

$$\mathbb{E} [\hat{\lambda}(t) - \lambda(t)]^2 \leq R_{L,T} \left(\frac{1}{Lh} \right)^2 + O(h^4), \quad t \in (0, T).$$

as $L \rightarrow \infty$. The minimum of the error yields the asymptotically optimal choice of the bandwidth, i.e., $h^* = cL^{-1/5}$ for some positive constant c . Theoretically, the asymptotically optimal choice of the bandwidth for the kernel classifier may be quite different from the practical application rule. Then specifying the bandwidth according to the resampling technique (22) shows that if (22) holds our experiments lead us to choose a 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 and

$$\text{CV}(h) = \prod_{l=1}^p \prod_{r=1}^{N^{[l]}} \hat{p}_i(t_r^{[l]}; h), \quad (68)$$

This leads to the following asymptotical formula for the mean squared error

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 $L_i = O(q/Lh)$ size training set. Then, the bandwidth is selected as the one that maximizes $\text{CV}(h)$ in (??). This is equivalent to the following choice

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 as 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 [?], [?].

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 $p_i(t; h)$ be the kernel estimate in (??) specified by the bandwidth h . Then, the likelihood function of $p_i(t; h)$ specified by test data is given by $R_{L,T}$ with respect to the training set size L and the observation window size T .

In all experiments the kernel classifier is given by (??) with the estimated $\hat{\tau}_i, \hat{p}_i(t; h)$ 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} with the test data. Finally we denote $\mathbb{E}[\text{CV}(h)]$ as the empirically evaluated risk averaged over ten simulation runs with a testing set size of 10^4 . Then, the bandwidth is selected as we shall focus on the size of $\text{CV}(h)$ results obtained for the intensity function specified by (??). Unless noted otherwise, we refer to the intensity function pair parametrized by $\phi_1 = \pi/16$ and $\phi_2 = \pi/4$.

Figure ?? depicts the average risk versus T for the size of training data ranging from $L = 10$ to $L = 200$. The Bayes risk R_T^* is also plotted for comparison. The convergence of $\mathbb{E}[\text{CV}(h)]$ to zero analogous as it was observed for the Bayes risk (see Figure ??) is seen. Also the small value of the difference $\mathbb{E}[R_{L,T}] - R_T^*$ for all T should be noted. We also observe 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 in subsequent analysis, i.e., with the value of T fixed.

In order to assess the proposed methodology, we conduct a simulated data 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 observation window size T .

Next, we analyze the value of the optimal bandwidth selected according to the log-likelihood method versus T . For brevity, in Figure ??a we show only the results for h_1 , noting that the curves obtained for h_2 are analogous. We observe an increase in behavior of $R_{L,T}$ with T , which aligns with the notion that as the observation window increases, the distribution of events in time becomes sparser, yielding the larger bandwidth. On the other hand, the obtained results also show that $h(L) \rightarrow 0$ as 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 the intensity function pair parameters ϕ_1, ϕ_2 versus L and fixed $T = 10$. Clearly, as the difficulty of the problem increases, i.e., data. Finally, we denote $\mathbb{E}[R_{L,T}]$ as an empirically evaluated risk averaged over ten simulation runs with a testing set size of 10^4 . Bayes risk is higher for more difficult classification problems.

We shall focus on the simulation results obtained for the intensity function specified by (??). Unless noted otherwise, we refer to the intensity function pair parametrized by $\phi_1 = \pi/16$ and $\phi_2 = \pi/4$.

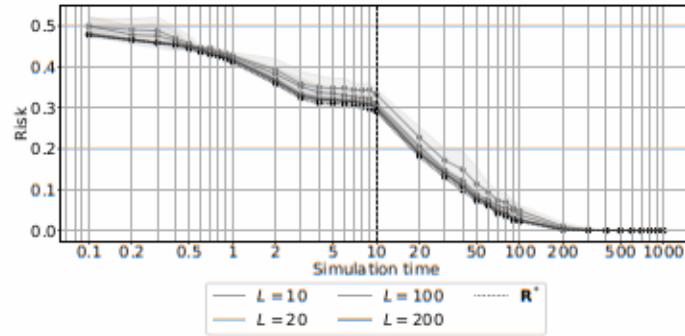


Figure 33. The average risks $\mathbb{E}[\mathbf{R}_L | \mathbf{R}_T^*]$ versus T for different values of L . The vertical dashed line denotes the simulation space slice in Figure ??-b-??.

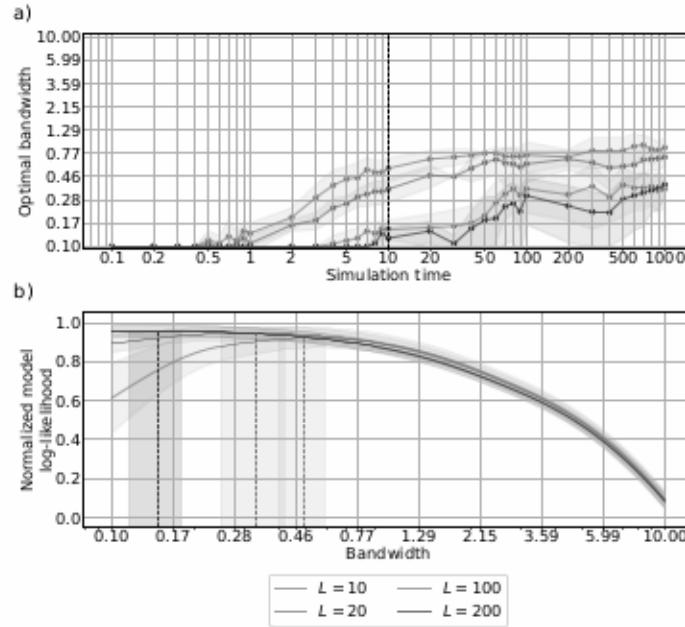


Figure 44. a) The average optimal bandwidth versus T for different values of L . The vertical dashed line at $T = 10$ denotes the simulation space slice in Figure ??-a-??b. b) The average normalized model log-likelihood for different values of L . The vertical dashed line at $h = 0.77$ denotes the maximum of the normalized log-likelihood function for $L = 10$. Note that the curves for $L = 100$ and $L = 200$ overlap.

Figure ?? depicts the average risk versus T for the size of training data ranging from $L = 10$ to $L = 200$. The Bayes risk \mathbf{R}_T^* is also plotted for comparison. The convergence of $\mathbb{E}[\mathbf{R}_{L,T}]$ to zero analogous as it was observed for the Bayes risk (see Figure ??) is seen. Also the small value of the difference $\mathbb{E}[\mathbf{R}_{L,T}] - \mathbf{R}_T^*$ for all T should be noted. We also observe 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 in subsequent analysis, i.e., with the value of T fixed.

Next, we analyze the value of the optimal bandwidth selected according to the log-likelihood method versus T . For brevity, in Figure ??-a we show only the results for \hat{h}_1 , noting that the curves obtained for \hat{h}_2 are analogous. We observe an increase in \hat{h}_1 with T , which aligns with the notion that as the observation window increases, the distribution of events in time becomes sparser, yielding the larger bandwidth. On the other hand, the obtained results also show that $h(L) \rightarrow 0$ as 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 to the Bayes risk for different values of the intensity function pair parameters ϕ_1, ϕ_2 versus L and fixed $T = 10$. Clearly, as the difficulty of the problem increases, i.e., when the two intensity functions become more similar to one another, the rate of convergence decreases. Also note that the Bayes risk is higher for more difficult classification problems.

Let us briefly examine a counter-example when the proposed algorithm fails to converge. Consider the following Gaussian type intensity function

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

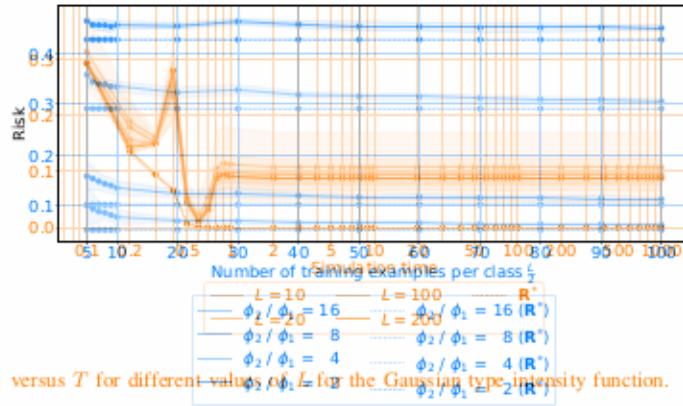
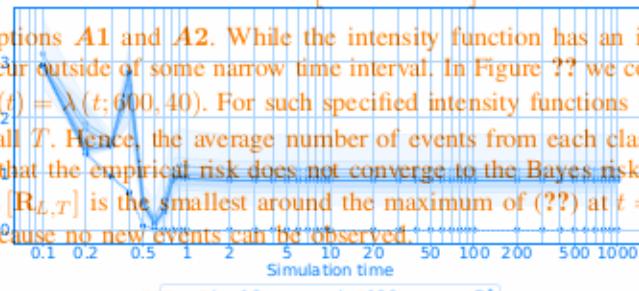


Figure 6. The average risk $E[R_{L,T}]$ versus T for different values of L for the Gaussian type-intensity function.

Figure 5. The average risk $E[R_{L,T}]$ 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 R_T^* for the associated intensity function pair.

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

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 $\lambda_1(t) = \lambda(t; 300, 20)$ and $\lambda_2(t) = \lambda(t; 600, 40)$. For such specified intensity functions we can evaluate that $\int_0^T \lambda_1(t) dt < 119$ and $\int_0^T \lambda_2(t) dt < 169$ 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 $E[R_{L,T}]$ 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.



VI. CONCLUDING REMARKS

In this paper we have developed the rigorous asymptotic analysis for the classification problem applied to spike trains data characterized by non-random intensity functions. The optimal Bayes rule was derived and its finite and asymptotic (with respect to the length of the observation interval) properties were established. This includes the exponential bound for the Bayes risk. Our asymptotic theory is relied on the martingale representation of counting processes. We then introduced a general class of plug-in empirical classification rules and formulated a generalization of the Bayes risk. This optimality property of some natural and interesting classifiers is derived from the plug-in rule with $\lambda_1(t) = \lambda(t; 300, 20)$ and $\lambda_2(t) = \lambda(t; 600, 40)$. For such specified intensity functions we can evaluate that $\int_0^T \lambda_1(t) dt < 119$ and $\int_0^T \lambda_2(t) dt < 169$ 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 $E[R_{L,T}]$ 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.

VI. Concluding Remarks

In this paper we have developed the rigorous asymptotic analysis for the classification problem applied to spike trains data characterized by non-random intensity functions. The optimal Bayes rule was derived and its finite and asymptotic (with respect to the length of the observation interval) properties were established. This includes the exponential bound for the Bayes risk. Our asymptotic theory is relied on the martingale representation of counting processes. We then introduced a general class of plug-in empirical classification rules and formulated a generalization of the Bayes risk. This optimality property is confirmed for 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.

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 The asymptotic theory of the classification problem discussed in this paper is based on the martingale method. This approach gives a generalization of the classification function within the assumptions provided the λ_k and their martingale representation, see (??) for the full classification theory. 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

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

where F_{ik} denotes the history of $N(t)$ in the interval $[0, t]$. Here we note that M is a general function due to the dependence probabilities. Utilizing prior to the gained decoupling position (see (??)) implies that the residual process allow us to generalize our asymptotic 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. (71)

satisfies the property

$$\mathbb{E}[\text{Appendix A 0.}] \quad (72)$$

This confirms the fact that the process $M(t) = N(t) - \int_0^t \lambda(u)du$ is a zero mean local martingale. 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

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

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

$$\text{Var}[dM(t)|\mathcal{F}_t] = \lambda(t)dt. \quad (71)$$

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

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

$$M_I(t) = N(t) - \int_0^t g(u)dM(u) \quad (72)$$

is a zero mean local martingale. 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 $I(t)$ is a zero mean martingale with respect to the history of the counting process $N(t)$. The variance of $I(t)$ is given by

$$\text{Var}[I(t)|\mathcal{F}_t] = \lambda(t)dt + dM(t). \quad (73)$$

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

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

$$\text{Var}[dM(t)|\mathcal{F}_t] = \lambda(t)dt. \quad (74)$$

Also $\text{Var}[dM(t)|\mathcal{F}_t] = \text{Var}[dN(t)|\mathcal{F}_t]$. $\text{Cov}\left[\int_0^t g_1(u)dM(u), \int_0^t g_2(u)dM(u), |\mathcal{F}_t\right]$

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

$$\text{where } g_1(t), g_2(t) \text{ are measurable functions. } I(t) = \int_0^t g(u)dM(u)$$

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 $I(t)$ is a zero mean martingale with respect to the history of the counting process $N(t)$. The variance of $I(t)$ is given by

$$\frac{x-1}{x+1} \leq \log(x) \leq \frac{2x-1}{2x}, \quad x > 0. \quad (75)$$

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

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

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

$$2\frac{x-1}{x+1} \leq \log(x) \leq \frac{x^2-1}{2x}, \quad x \geq 1. \quad (77)$$

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

$$\alpha_T = \tau_1 - \int_0^T g_1(u)g_2(u)\lambda(u)du - \tau_1 \int_0^T \log\left(\frac{p_1(t)}{p_2(t)}\right)p_1(t)dt. \quad (78)$$

Here $K_{g_1}(t), g_2(t)$ are measurable functions. $K_{g_1}(t)$ is the Kullback-Leibler divergence between the densities $p_1(t)$ and $p_2(t)$. Then, by virtue of (??) we have

$$\begin{aligned} \alpha_T &\leq \tau_1 - \tau_2 + \tau_1 \left\{ \frac{\tau_2}{\tau_1} - 1 \right\} - \tau_1 K_T(p_1 \| p_2) \\ &= -\tau_1 K_T(p_1 \| p_2) \end{aligned}$$

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,

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

$$\begin{aligned} \alpha_T &\geq \tau_1 - \tau_2 + \tau_1 \left\{ 1 - \frac{\tau_1}{\tau_2} \right\} - \tau_1 \mathbf{K}_T(p_1 \parallel p_2) \\ \frac{x - 1}{x + 1} &\leq \log(x) \leq \frac{x - 1}{x}, \quad x > 0. \end{aligned} \quad (78)$$

The tighter version of this inequality for $x \geq \frac{\tau_2}{\tau_1}$ reads as follows

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

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

$$\text{Then, the application of (??) gives } \alpha_T = \tau_1 - \tau_2 + \tau_1 \log \left(\frac{\tau_2}{\tau_1} \right) + \tau_1 \int_0^T \log \left(\frac{p_1(t)}{p_2(t)} \right) p_1(t) dt. \quad (80)$$

Proof of Lemma ??. The result of Lemma ?? is implied by the straightforward application of the identity in (??) to the stochastic integral $\int_0^T \log \left(\frac{p_1(t)}{p_2(t)} \right) p_1(t) dt$ is the Kullback-Leibler divergence between the densities $p_1(t)$ and $p_2(t)$. Then, by virtue of (??) we have

$$\int_0^T \log \left(\frac{\lambda_1(t)}{\lambda_2(t)} \right) dM(t).$$

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

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

$$\alpha_T \geq \tau_1 - \tau_2 + \tau_1 \left\{ 2 \exp \left[\frac{\tau_1}{\tau_2} \right] \frac{v_T}{u_T^2} J_1 \left(\frac{u_T}{v_T} p_1 \right) \parallel p_2 \right\} \quad (81)$$

where $J(x) = (1+x) \log(1+x) - x$, $x \geq \frac{\tau_1 - \tau_2}{\tau_2}$. Using the inequalities $\mathbf{K}_T(p_1 \parallel p_2) \geq 0$ we can easily obtain that

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

Proof of Lemma ??. The result of Lemma ?? is implied by the straightforward application of the identity in (??) to the stochastic integral

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

where $\kappa = \log(\pi_2/\pi_1)$. Since $\mathbf{X} \in \omega_2$ then

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

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

$$\mathbf{P}(|U_T| \geq \epsilon) \leq 2 \exp \left[- \frac{b_T}{u_T^2} J \left(\epsilon \frac{u_T}{v_T} \right) \right], \quad (84)$$

where

where $J(x) = (1+x) \log(1+x) - x$, $x \geq 0$. Using the Van Hove inequalities in (??) we can easily obtain that
 $b_T = \frac{1}{v_T} \sqrt{\frac{1}{T} \mathbf{H}_T(\mathbf{X})}$. \square

By the assumptions **A1** and **A2** we have

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

Then by the result of Lemma ?? and (??) $\mathbf{P} \left(\frac{1}{T} \mathbf{g}_T(\mathbf{X}) \geq \frac{1}{T} \alpha_T + \frac{1}{T} \kappa | \mathbf{X} \in \omega_2 \right)$, \square

where $\kappa = \log(\pi_2/\pi_1)$. Since $\mathbf{X} \in \omega_2$ then $b_T \leq \frac{d \log^2(\frac{C}{\delta})}{d^2 \log^2(\frac{\delta}{C})} = \frac{1}{d} \left(\frac{\log(C/\delta)}{\log(\delta/C)} \right)^2$.

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

Then, by the probability of misclassification $P(W_T(\mathbf{X}) < \eta_T | \mathbf{X} \in \omega_1)$ in (??) 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 $P(W_T(\mathbf{X}) < \eta_T | \mathbf{X} \in \omega_1)$ 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 $P(W_T(\mathbf{X}) \geq \eta_T | \mathbf{X} \in \omega_2)$ or equivalently the probability in (??). We wish to use the exponential inequality in (??) of Lemma ??.

$$\text{Then, the probability } P\left[\frac{1}{\sqrt{T}}U_T(\mathbf{X})\right] \text{ is bounded by} \quad (85)$$

$$b_T = \frac{\exp\left[\frac{T^{-1}\alpha_T + \frac{2}{T}\kappa T^{-1}\kappa}{2\theta_T + u\epsilon_T}\right]^2}{\exp\left[\frac{T^{-1}\alpha_T + \frac{2}{T}\kappa T^{-1}\kappa}{2\theta_T + u\epsilon_T}\right]}, \quad (89)$$

By the assumptions **A1** and **A2** we have

where $u = \log\left(\frac{C}{\delta}\right)$ characterizes the assumption **A1**, $\epsilon_T = \frac{1}{T}\alpha_T + \frac{1}{T}\kappa$, and $\theta_T = \text{Var}\left[\frac{1}{\sqrt{T}}U_T(\mathbf{X})\right]$. This defines the exponential factor

and also

$$B_T = \frac{\epsilon_T^2}{2\theta_T + u\epsilon_T}\left(\frac{\delta}{C}\right). \quad (87)$$

Owing to Lemma ??, (??) and (??), the limit inferior of B_T is not smaller than

$$\begin{aligned} \liminf_{T \rightarrow \infty} B_T &\geq \frac{d^2 \log^2(\delta/C)}{d \log^2\left(\frac{C}{\delta}(C/\delta)\right) + \frac{u d \log(C/\delta)}{2 \log(C/\delta)} \cdot \frac{2}{d} \log^2\left(\frac{C}{\delta}\right)} \\ \liminf_{T \rightarrow \infty} b_T &\leq \frac{d^2 \log^2\left(\frac{C}{\delta}\right) \delta / C}{d^2 \log^2\left(\frac{C}{\delta}\right) \delta / C} = \frac{2}{d} \log\left(\frac{C}{\delta}\right). \end{aligned} \quad (88)$$

Hence, the probability of misclassification $P(W_T(\mathbf{X}) \geq \eta_T | \mathbf{X} \in \omega_2)$ is bounded by b_T/T , where the superior limit of b_T is given in (??) with the required bound. Since the probability of misclassification $P(W_T(\mathbf{X}) > \eta_T | \mathbf{X} \in \omega_1)$ is bounded by $P(W_T(\mathbf{X}) < \eta_T | \mathbf{X} \in \omega_1)$ the superior limit of this also give Theorem ?? This completes the proof of Theorem ??.

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

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

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

$$\begin{aligned} \exp\left[-\frac{\epsilon_T^2}{2\theta_T + u\epsilon_T}\right] &\geq 1 - \delta/2, \\ \text{where } u = \log\left(\frac{C}{\delta}\right) \text{ characterizes the assumption } \mathbf{A1}, \epsilon_T = \frac{1}{T}\alpha_T + \frac{1}{T}\kappa, \text{ and } \theta_T = \text{Var}\left[\frac{1}{\sqrt{T}}U_T(\mathbf{X})\right]. \end{aligned} \quad (90)$$

Let $\psi_T^*(\mathbf{x}) = \omega_1$, i.e., we have $W_T(\mathbf{x}) > \eta_T$. Then, $B_T = \frac{\epsilon_T^2}{2\theta_T + u\epsilon_T}$.

Owing to Lemma ??, (??) and (??) $P(\widehat{W}_{L,T}(\mathbf{x}) - W_T(\mathbf{x}) < \epsilon) > 1 - \delta/2$.

$$\begin{aligned} \text{The right-hand side of this equality is not smaller than} \quad &\frac{d^2 \log^2(\delta/C)}{2d \log^2(C/\delta) + ud \log(C/\delta)} \\ P\left(\left|\widehat{W}_{L,T}(\mathbf{x}) - \widehat{\eta}_{L,T}\right| - \left|W_T(\mathbf{x}) - \eta_T\right| < 2\epsilon\right) &= \frac{2}{d} \log\left(\frac{C}{\delta}\right) \end{aligned} \quad (91)$$

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

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

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

Appendix C

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

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

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

Proof of Lemma ??. The proof will be based on the following version of Helly's theorem [?] for the Stieltjes integral. Let $\psi_T^*(\mathbf{x}) = \omega_1$, i.e., we have $W_T(\mathbf{x}) > \eta_T$. Then,

$$P\left(f(x) \rightarrow f(x) \text{ uniformly on } [a, b] \text{ as } L \rightarrow \infty\right) = P\left(\widehat{W}_{L,T}(\mathbf{x}) \rightarrow \widehat{\eta}_{L,T}\right).$$

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

$$P\left(\left|\int_a^b (\widehat{W}(x)dx - \widehat{\eta}_{L,T})\right| < 2\epsilon\right) \quad (93)$$

Consider the $\frac{1}{2}$ optimal decision. Moreover, $W_T(x)$ is bounded from below by its empirical counterpart $\widehat{W}_{L,T}(x)$ in (??). Then, we can write (see ??))

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

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

$$\mathbf{P} \left(\left| \widehat{W}_{L,T}(x) - W_T(x) \right| < \epsilon \right) + \mathbf{P} (|\widehat{\eta}_{L,T} - \eta_T| < \epsilon) - 1. \quad (94)$$

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

$$\left| \log \left(\frac{\widehat{p}_1(t)}{p_2(t)} \right) - \log \left(\frac{p_1(t)}{p_2(t)} \right) \right| \rightarrow 0 \text{ (P)}. \quad (95)$$

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

$$f_L(x) \rightarrow f(x), \text{ uniformly on } [a, b] \text{ as } L \rightarrow \infty.$$

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

$$= \left| \int_a^b \frac{(p_1(t) - p_1(t))p_2(t) + (p_2(t) - \widehat{p}_2(t))p_1(t)}{f_L(x)(p_2(t) - p_2(t))p_1(t)} dg(x) \right| \rightarrow 0 \text{ (P) as } L \rightarrow \infty. \quad (93)$$

This is not greater than. Consider the optimal decision function $W_T(x)$ in (??) and its empirical counterpart $\widehat{W}_{L,T}(x)$ in (??). Then, we can write (see ??))

$$p_1(t)p_2(t)$$

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 not exceed

$$= \left| \frac{C}{\delta} \int_0^T \left[\log \left(\frac{\widehat{p}_1(t)}{p_2(t)p_1(t)} \right) - \log \left(\frac{p_1(t)}{p_2(t)p_1(t)} \right) \right] dN(t) \right|. \quad (94)$$

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

$$\left| \log \left(\frac{\widehat{p}_1(t)}{p_2(t)} \right) \right| \rightarrow 0 \text{ (P)}. \quad (95)$$

Proof of Lemma ??. We wish to show that

$$\text{uniformly on } [0, T] \text{ as } L \rightarrow \infty$$

Observe that the left-hand side of (??) is equal to $\left| \sup_{t \in T_e} \left| \widehat{\lambda}(t) - \lambda(t) \right| \right| \log \left(\frac{p_1(t)p_2(t)}{p_2(t)p_1(t)} \right)$ as $L \rightarrow \infty$. Then, using (??) this is bounded by

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

$$\left| \widehat{\lambda}(t) - \lambda(t) \right| \leq \left| \widehat{\lambda}(t) - \mathbb{E}[\widehat{\lambda}(t)] \right| + \left| \mathbb{E}[\widehat{\lambda}(t)] - \lambda(t) \right|. \quad (97)$$

Owing to (??) the bias term is equal to

$$\left| \mathbb{E} \left[\widehat{\lambda}(t) \right] - \lambda(t) \right| \leq \left| \int_{t/h}^{(T-t)/h} K(s) \lambda(p_2(t)/p_1(t)) ds - \lambda(t) \right|$$

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. 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 not exceed

$$\left| \mathbb{E} \left[\widehat{\lambda}(t) \right] \right| = \left| \int_0^1 K(s) |\lambda(t+hs) - \lambda(t)| ds \right|$$

$$\leq M_\lambda h \int_0^1 K(s) |s| ds$$

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

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., $T_\epsilon = \bigcup_{j=1}^J U_j$, where the size of U_j is denoted as $\Delta(L)$. Clearly the number of intervals is of order

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

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

$$|\widehat{\lambda}(t) - \lambda(t)| \leq \left| \widehat{\lambda}(t) - \mathbb{E}[\widehat{\lambda}(t)] \right| + \left| \mathbb{E}[\widehat{\lambda}(t)] - \lambda(t) \right|. \quad (97)$$

Owing to (I) to (III) the bias term is indeed point of \mathbf{U}_j . Then, the uniform norm of the stochastic term in (??) can be bounded as follows

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

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

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

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 Consider first the term A_1 . By virtue of $q(L)$ we have

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 (??) 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

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

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 $\text{Var} [\bar{N}_L(T)] = \frac{1}{L} \int_0^T \lambda(t)^2 dt$. This proves that

$$\begin{aligned} &+ \max_{1 \leq j \leq q(L)} \sup_{t \in T_\epsilon \setminus \mathbf{U}_j} \left| \mathbb{E} [\widehat{\lambda}(t)] - \mathbb{E} [\widehat{\lambda}(u_j)] \right| \cdot \\ &+ \max_{1 \leq j \leq q(L)} \left| \widehat{\lambda}(u_j) - \mathbb{E} [\widehat{\lambda}(u_j)] \right| \end{aligned} \quad (98)$$

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

uniformly in $t \in T_\epsilon$. Concerning the term A_2 in (??) we can use (??). Then, we obtain

$$= A_1 + A_2 + A_3.$$

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

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

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

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

$$|\widehat{\lambda}(t) - \bar{\lambda}(u_j)| \leq M_K \frac{\Delta(L)}{h^2} \int_0^T d\bar{N}_L(s). \quad (99)$$

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

Finally, let us consider the term A_3 in (??). First we note that $\mathbb{E} [\bar{N}_L(T)] = \int_0^T \lambda(t) dt$ and $\text{Var} [\bar{N}_L(T)] = \frac{1}{L} \int_0^T \lambda(t)^2 dt$. This proves that

$$\mathbb{P} \left(\max_{t \in T_\epsilon} |\widehat{\lambda}(u_j) - \mathbb{E} [\widehat{\lambda}(u_j)]| \geq \delta \right) = \mathcal{O} \left(\frac{\Delta(L)}{h^2} \right) \quad (a.s.) \quad (100)$$

uniformly in $t \in T_\epsilon$. Concerning the term A_2 in (??) we can use (??). Then, we obtain

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

$$\begin{aligned} & \mathbb{E} [\widehat{\lambda}(t)] - \mathbb{E} [\widehat{\lambda}(u_j)] \\ &= \widehat{\lambda}(t) \int_0^T \mathbb{E} [\widehat{\lambda}(s)] ds = \int_0^T K_h(t-s) \int_0^s K_h(u_j-s) \lambda(s) ds ds. \end{aligned}$$

This, (??) and Chebyshev inequality yield

This gives

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

Note that the right-hand side of this inequality is of order $\mathcal{O}(1/L \Delta(L))$ uniformly in $t \in T_\epsilon$. This, (??) 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 $\delta > 0$

$$\mathbb{P} (A_3 \geq \delta) = \mathcal{O}(1/\Delta(L)Lh)$$

or equivalently $A_3 = \mathcal{O}_P(1/\sqrt{\Delta(L)Lh})$. Hence, balancing $A_3 = \mathcal{O}_P(1/\sqrt{\Delta(L)Lh})$ 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 $\mathbb{E} [\widehat{\lambda}(t)] \geq \delta$

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

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.

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

This, (??) and Chebyshev inequality yield

ACKNOWLEDGMENT

- [1] W. Gersiner, W. M. Kistler, R. Naud, 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 Note that the right hand side of this inequality is of order $O(1/Lh)$ uniformly in $t \in T_t$. This, (??) and the fact that $qL = O(1/\Delta(L))$ lead to the following uniform bound
- [3] O. Sohn, A. Türkmen, T. Januszkiewicz, and S. Gimpermann, "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. 1837–1849, 2008.

[6] N. Merhav, "Optimal correlators for detection and estimation in optical receivers," *IEEE Transactions on Information Theory*, vol. 67, pp. 5200–5210, or equivalently $A_3 = \mathcal{O}_P(1/\sqrt{\Delta(L)Lh})$. Hence, balancing $A_3 = \mathcal{O}_P(1/\sqrt{\Delta(L)Lh})$ versus $A_1, A_2 = \mathcal{O}(\Delta(L)/h^2)$

[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. Cholakidis, 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.

The proof of Lemma ?? has been completed.

[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, RTbust: Exploiting temporal patterns for botnet detection on twitter," in *Proceedings of the 10th ACM conference on web science*, 2019, pp. 183–192.

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

[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] O. Aalen, "Nonparametric inference for a family of counting processes," *The Annals of Statistics*, pp. 701–726, 1978.

[17] M. Pawlak, 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*, 2023, pp. 1–5.

[18] 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.

[19] D. Stoyan and H.-G. Müller, Cox process regression," *IEEE Transactions on Information Theory*, vol. 68, no. 2, pp. 1133–1156, 2022.

[20] L. Birgé and P. Massart, "Rates of convergence for minimum contrast estimators," *Probability Theory and Related Fields*, vol. 97, pp. 112–150, 1993.

[21] S. Van de Geer, "Exponential inequalities for martingales, with application to maximum likelihood estimation for coupling processes," *The Annals of Probability*, pp. 779–801, 1996.

[22] S. Ghosal and A. Van der Vaart, *Fundamentals of Nonparametric Bayesian Inference*. Cambridge University Press, 2017.

[23] S. Ghosal and A. Van der Vaart, "Fundamentals of Nonparametric Bayesian Inference," *Comptes Rendus. Mathématique*, vol. 359, no. 8, pp. 767–782, 2021.

[24] D. J. Daley and D. Vere-Jones, "On the estimation of frequency in point-process data," *Journal of Applied Probability*, vol. 19, pp. 383–394, 1982.

[25] M. P. Wand and M. C. Jones, *Kernel Smoothing*. CRC Press, 1994.

[26] P. K. Andersen, O. Borgan, R. D. Gill, and N. Keiding, *Statistical Models Based on Counting Processes*. Springer, 2012.

[27] W. Greblicki and M. Pawlak, *Nonparametric System Identification*. Cambridge University Press, 2008.

[28] M. Devroye, L. Györfi, and G. Lugosi, *A Probabilistic Theory of Pattern Recognition*. Springer, 1996.

[29] M. DeGroot, M. Jansson, and X. Ma, "Simple local polynomial density estimators," *Journal of the American Statistical Association*, vol. 115, no. 531, pp. 1449–1455, 2020.

[30] E. Giné and A. Guillou, "Rates of strong uniform consistency for multivariate kernel density estimators," *Annales de l'Institut Henri Poincaré (B) Probabilités et Statistiques*, vol. 38, no. 6, pp. 907–921, 2002.

[31] E. Giné and A. Guillou, "On the error rate for classifying point processes," in *60th IEEE Conference on Decision and Control (CDC)*. IEEE, 2021, pp. 120–125.

[32] E. Giné and A. Guillou, "Rates of strong uniform consistency for multivariate kernel density estimators," *Annales de l'Institut Henri Poincaré (B) Probabilités et Statistiques*, vol. 38, no. 6, pp. 907–921, 2002.

[33] E. Giné and A. Guillou, "Exponential inequalities for the supremum of some counting processes and their square martingales," *Comptes Rendus. Mathématique*, vol. 353, pp. 1–15, 2017.

[34] E. Giné and A. Guillou, "Rates of strong uniform consistency for multivariate kernel density estimators," *Annales de l'Institut Henri Poincaré (B) Probabilités et Statistiques*, vol. 38, no. 6, pp. 907–921, 2002.

[35] E. Giné and A. Guillou, "Exponential inequalities for the supremum of some counting processes and their square martingales," *Stochastic Processes and their Applications*, vol. 65, no. 1, pp. 81–101, 1996.

[36] M. Mazza, S. Cresci, M. Avvenuti, W. Quattrociocchi, and M. Tesconi, RTbust: Exploiting temporal patterns for botnet detection on twitter," in *Proceedings of the 10th ACM conference on web science*, 2019, pp. 183–192.

[37] J. Dedecker and A. B. Tsybakov, "Fast learning rates for plug-in classifiers," *The Annals of Statistics*, vol. 35, pp. 608–633, 2007.

[38] H. Fairhurst and J. Gama, "Event labeling combining ensemble detectors and background knowledge," *Progress in Artificial Intelligence*, vol. 1, no. 2, pp. 113–127, 2023.

[39] R. Lim, "Hawkes processes modeling, inference, and control: an overview," *SIAM Review*, vol. 65, pp. 331–374, 2023.

[40] W. Greblicki, *Mathematical Analysis*. Addison-Wesley, 1974.

[41] W. Greblicki, "Asymptotically optimal pattern recognition procedures with density estimates," *IEEE Transactions on Information Theory*, vol. 24, no. 2, pp. 250–251, 1978.

[42] 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.

[43] O. Aalen, "Nonparametric inference for a family of counting processes," *The Annals of Statistics*, pp. 701–726, 1978.

[44] M. Pawlak, 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*, 2023, pp. 1–5.

[45] Á. Gajardo and H.-G. Müller, "Cox point process regression," *IEEE Transactions on Information Theory*, vol. 68, no. 2, pp. 1133–1156, 2022.

[46] L. Birgé and P. Massart, "Convergence of minimum contrast estimators: Probability Theory and Related Fields", vol. 97, pp. 113–150, 1993. 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 at The Annals of Statistics, pp. 177–187, 1995.

[47] S. Ghosal and A. Van der Vaart, *Fundamentals of Nonparametric Bayesian Inference*. Cambridge University Press, 2017.

[48] R. Le Guével, "Exponential inequalities for the supremum of some counting processes and their square martingales," *Comptes Rendus. Mathématique*, vol. 349, no. 8, pp. 369–382, 2011. His interests include statistical signal processing, machine learning, and nonparametric modeling.

[49] D. Vere-Jones, "On the estimation of frequency in point-process data," *Journal of Applied Probability*, vol. 19, pp. 383–394, 1982.

[50] M. P. Wand and M. C. Jones, *Kernel Smoothing*. CRC Press, 1994.

[51] W. Greblicki and M. Pawlak, *Nonparametric System Identification*. Cambridge University Press, 2008.

[52] M. D. Cattaneo, M. Jansson, and X. Ma, "Simple local polynomial density estimators," *Journal of the American Statistical Association*, vol. 115, no. 531, pp. 1449–1455, 2020.

[53] E. Giné and A. Guillou, "Rates of strong uniform consistency for multivariate kernel density estimators," *Annales de l'Institut Henri Poincaré (B) Probabilités et Statistiques*, vol. 38, no. 6, pp. 907–921, 2002.

- [28] E. Masry, "Multi-Mateusz Pabian received his 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 2017 he was with the Wireless Sensor and Control Networks Group, AGH University of Science and Technology, where he was involved in the design of 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 Postdoctoral 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.
- [29] J.-Y. Audibert and A. Boulle, "Event-based learning," in *Advances in Statistical Signal Processing: Trends and Applications*, vol. 65, no. 1, pp. 1–60, 2016.
- [30] R. Lima, "Hawkes processes method for signal design and signal acquisition via noisy SISAs," *The Rendiconti del Circolo Matematico di Palermo*, vol. 63, pp. 311–374, 2023.
- [31] T. Apostol, *Mathematical Analysis: A Collection of Problems and Solutions*, Mir, Moscow, 1974.



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 Electrical and Computer Engineering from the AGH University of Science and Technology, Kraków, Poland, in 2017. He completed a scholarship from NKTAGH University of Science and Technology, Kraków, Poland, in 2016–2017. The program mainly focused on the expertise of signal processing and signal processing. AGH University of Science and Technology and its potential where from 2015 to 2022 he worked on the design of machine learning for medical signal processing and signal processing. In 2011, he joined the Event-Based Control and Signal Processing Group at AGH University of Science and Technology, where he currently works on signal processing and machine learning systems. He was a Visiting Student and Postdoctoral 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 Postdoctoral 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.