

Likelihood ratio Haar variance stabilization and normalization for Poisson and other non-Gaussian noise removal

Piotr Fryzlewicz*

June 16, 2017

Abstract

We propose a methodology for denoising, variance-stabilizing and normalizing signals whose varying mean and variance are linked via a single parameter, such as Poisson or scaled chi-squared. Our key observation is that the signed and square-rooted generalized log-likelihood ratio test for the equality of the local means is approximately distributed as standard normal under the null. We use these test statistics within the Haar wavelet transform at each scale and location, referring to them as the *likelihood ratio Haar (LRH) coefficients* of the data. In the denoising algorithm, the LRH coefficients are used as thresholding decision statistics, which enables the use of thresholds suitable for i.i.d. Gaussian noise. In the variance-stabilizing and normalizing algorithm, the LRH coefficients replace the standard Haar coefficients in the Haar basis expansion. We prove the consistency of our LRH smoother for Poisson counts with a near-parametric rate, and various numerical experiments demonstrate the good practical performance of our methodology.

Key words: variance-stabilizing transform, Haar-Fisz, Anscombe transform, log transform, Box-Cox transform, Gaussianization.

*Department of Statistics, London School of Economics, Houghton Street, London WC2A 2AE, UK. Email: p.fryzlewicz@lse.ac.uk. Work supported by the Engineering and Physical Sciences Research Council grant no. EP/L014246/1.

1 Introduction

The popularity of wavelets and their potential for useful applications in data science did not escape the attention of Peter Hall¹, who wrote, amongst others, on threshold choice in wavelet curve estimation (Hall and Patil, 1996a,b), wavelet methods for functions with many discontinuities (Hall et al., 1996), wavelets for regression with irregular design (Hall and Turlach, 1997) and block-thresholded wavelet estimators (Hall et al., 1999). I learned of Peter through wavelets by reading some of his papers on the topic during my doctoral study. I remember my surprise at discovering that both my then PhD supervisor, Guy Nason, and someone else I knew, Prakash Patil, had co-authored papers with Peter Hall. When I shared my surprise with Guy, he responded by saying that he did not know many people who were *not* Peter’s co-authors! Even though I can unfortunately count myself in this “minority” category, I have learned and am still learning a great deal from Peter, especially by appreciating the careful and elegant way in which he used mathematics to support his arguments.

Traditional wavelet transformations are orthonormal transformations of the input data into coefficients that carry information about the local behaviour of the data at a range of dyadic scales and locations. They tend to offer sparse representation of the input data, with a small number of wavelet coefficients often being able to encode much of the energy of the input signal, and are computable and invertible in linear time via recursive pyramid algorithms (Mallat, 1989; Daubechies, 1992). Reviews of the use of wavelets in statistics can be found, for example, in Vidakovic (1999) and Nason (2008). One canonical task facilitated by wavelets is the removal of noise from signals, which usually proceeds by taking a wavelet transform of the data, thresholding away the (typically many) wavelet coefficients that are small in magnitude, preserving those few that are large in magnitude, and taking the inverse wavelet transform. Since the seminal paper by Donoho and Johnstone (1994) in which the general idea was first proposed, several other methods for wavelet smoothing of one-dimensional signals have appeared, but the vast majority make the i.i.d. Gaussian noise assumption. By contrast, the focus

¹This article is to appear in a special issue of *Statistica Sinica* in memory of Prof. Peter Hall.

of this article is the treatment of signals in which the variance of the noise is a function of its mean; this includes Poisson- or scaled-chi-squared-distributed signals. (Throughout the paper, we refer to a distribution as a ‘scaled chi-squared’, or simply ‘chi-squared’, if it takes the form $\sigma^2 m^{-1} \chi_m^2$.)

The simplest example of a wavelet transform, and the focus of this article, is the Haar transform, which can be described as a sequence of symmetric scaled differences of consecutive local means of the data, computed at dyadic scales and locations and naturally forming a binary tree consisting of ‘parents’ and ‘children’. Its local difference mechanism means that it offers sparse representations for (approximately) piecewise-constant signals. Our starting point is the observation that testing whether or not each Haar coefficient of a signal exceeds a certain threshold (in the denoising task described above) can be interpreted as the *likelihood ratio test* for the equality of the corresponding local means of the signal in the i.i.d. Gaussian noise model. In this paper, we take this observation further and propose similar multiscale likelihood ratio tests for other distributions, most notably those in which the variance is a function of the mean, such as Poisson or scaled chi-squared. The proposed multiscale likelihood ratio tests will reduce to the traditional thresholding of Haar wavelet coefficients for Gaussian data, but will take entirely different and new forms for other distributions. This will lead to a new, unified class of algorithms useful for problems such as e.g. Poisson intensity estimation, Poisson image denoising, spectral density estimation in time series, or time-varying volatility estimation in finance. (Extension of our methodology to images is as straightforward as the extension of the standard one-dimensional Haar wavelet transform to two dimensions.)

The new multiscale likelihood ratio tests will naturally induce a new construction, *likelihood ratio (Haar) wavelets*, which have the benefit of producing (equivalents of) Haar wavelet coefficients that are asymptotically standard normal under the null hypothesis of the corresponding local means being equal, even for inhomogeneous non-Gaussian signals. This will (a) make it much easier to choose a single threshold parameter in smoothing these kinds of data and (b) serve as a basis for new normalizing transformations for these kinds of data, which bring their distribution close to Gaussianity. This

article demonstrates both these phenomena. The device that enables these results is the Wilks’ theorem, according to which the signed square-rooted likelihood ratio statistic will often be approximately distributed as standard normal, a fact that, we believe, has not been explored in a variance-stabilization context before.

Wavelet-based Poisson noise removal, with or without the use of a variance-stabilizing and/or normalizing transform, has a long history. For a Poisson variable X , the Anscombe (1948) transform $2(X + 3/8)^{1/2}$ brings its distribution to approximate normality with variance one. Donoho (1993) proposes to pre-process Poisson data via the Anscombe transform, and then use wavelet-based smoothing techniques suitable for i.i.d. Gaussian noise. This and a number of other wavelet-based techniques for denoising Poisson-contaminated signals are reviewed and compared in Besbeas et al. (2004). These include the translation-invariant multiscale Bayesian techniques by Kolaczyk (1999b) and Timmermann and Nowak (1997, 1999), shown to outperform earlier techniques in Kolaczyk (1997, 1999a) and Nowak and Baraniuk (1999). Willett and Nowak (2003) propose the use of “platelets” in Poisson image denoising. The Haar-Fisz methodology Fryzlewicz and Nason (2004), drawing inspiration from earlier work by Fisz (1955) outside the wavelet context, proceeds by decomposing the Poisson data via the standard Haar transform, then variance-stabilizing the Haar coefficients by dividing them by the MLE of their own standard deviation, and then using thresholds suitable for i.i.d. Gaussian noise with variance one. Closely related ideas appear in Luisier et al. (2010) and Reynaud-Bouret and Rivoirard (2010). Jansen (2006) extends the Haar-Fisz idea to other wavelets. As an alternative to Anscombe’s transform, which is known not to work well for low Poisson intensities, Zhang et al. (2008) introduce a more involved square-root-type variance-stabilizing transform for (filtered) Poisson data. Hirakawa and Wolfe (2012) propose Bayesian Haar-based shrinkage for Poisson signals based on the exact distribution of the difference of two Poisson variates (the Skellam distribution).

In multiplicative set-ups, such as signals distributed as $X_k = \sigma_k^2 m^{-1} \chi_m^2$, the logarithmic transform stabilizes the variance exactly, but does not bring the distribution of the transformed X_k close to

normality, especially not for small values of m such as 1 or 2. In the context of spectral density estimation in time series, in which the signal is approximately exponentially distributed, wavelet shrinkage for the logged (and hence variance-stabilized) periodogram is studied, amongst others, in Moulin (1994), Gao (1997), Pensky et al. (2007) and Freyermuth et al. (2010). An alternative route, via pre-estimation of the variance of the wavelet coefficients (rather than via variance stabilization) is taken in Neumann (1996). Haar-Fisz or wavelet-Fisz estimation for the periodogram or other (approximate) chi-squared models is developed in Fryzlewicz et al. (2006), Fryzlewicz and Nason (2006) and Fryzlewicz et al. (2008). In more general settings, wavelet estimation for exponential families with quadratic or cubic variance functions is considered in Antoniadis and Sapatinas (2001), Antoniadis et al. (2001) and Brown et al. (2010). The Haar-Fisz or wavelet-Fisz transformations for unknown distributions are studied in Fryzlewicz (2008), Fryzlewicz et al. (2007), Motakis et al. (2006) and Nason (2014). Variance-stabilizing transformations are reviewed in the (unpublished) manuscript by Foi (2009).

Our approach departs from the existing literature in that our variance-stabilization and normalization device does not involve either the pre-estimation of the variance (as, effectively, in the Haar-Fisz transform) or the application of a Box-Cox-type transform (as in the Anscombe variance stabilization for Poisson data or the logarithmic transform in multiplicative models). By contrast, we use the entire likelihood for the purpose of variance-stabilization and normalization. As a result, the thresholding decision in our proposed smoothing methodology is not based on the usual wavelet detail coefficients, but on the newly-proposed likelihood ratio Haar coefficients. For completeness, we mention that Kocaczyk and Nowak (2004) construct multiscale decompositions of the Poisson likelihood, which leads them to consider binomial likelihood ratio tests for the purpose of thresholding; however, this is done in a context that does not use the signed and square-rooted generalized log-likelihood ratio tests or utilize their variance-stabilizing or normalizing properties.

The paper is organized as follows. Section 2 motivates and introduces the concept of likelihood ratio Haar coefficients and outlines our general methodology for smoothing and variance stabiliza-

tion/normalization. Section 3 describes our method in two special cases, those of the Poisson and the scaled chi-squared distribution. Section 4 formulates and discusses a consistency result for the Poisson smoother. Section 5 provides a numerical study illustrating the practical performance of our smoothing and variance stabilization/normalization algorithms. The online supplement contains the proof of our theoretical result and further technical details.

2 General methodology

Let X_1, \dots, X_n be a sequence of independent univariate random variables such that $X_k \sim F(\theta_k)$, where $F(\theta)$ is a family of distributions parameterized by a scalar parameter θ such that $\mathbb{E}(X_k) = \theta_k$. Our two running examples are: $X_k \sim \text{Pois}(\lambda_k)$, and $X_k \sim \sigma_k^2 m^{-1} \chi_m^2$ (throughout the paper, we refer to the latter example as ‘scaled chi-squared’ or simply ‘chi-squared’). Extensions to higher-dimensional parameters are possible, but certain aspects of the asymptotic normality are then lost.

We recall the traditional Haar transform and the fundamentals of signal smoothing via (Haar) wavelet thresholding. In the following, we assume that $n = 2^J$, where J is an integer. Extensions to non-dyadic n are possible, see e.g. Wickerhauser (1994). Given the input data $\mathbf{X} = (X_1, \dots, X_n)$, we define $\mathbf{s}_0 = (s_{0,1}, \dots, s_{0,n}) = \mathbf{X}$. The Haar transform recursively performs the following steps

$$s_{j,k} = 2^{-1/2}(s_{j-1,2k-1} + s_{j-1,2k}), \quad d_{j,k} = 2^{-1/2}(s_{j-1,2k-1} - s_{j-1,2k}), \quad (1)$$

for $j = 1, \dots, J$ and $k = 1, \dots, 2^{J-j}$. The indices j and k are thought of as “scale” and “location” parameters, respectively, and the coefficients $s_{j,k}$ and $d_{j,k}$ as the “smooth” and “detail” coefficients (respectively) at scale j , location k . It is easy to express $s_{j,k}$ and $d_{j,k}$ as explicit functions of \mathbf{X} :

$$s_{j,k} = 2^{-j/2} \sum_{i=(k-1)2^j+1}^{k2^j} X_i, \quad d_{j,k} = 2^{-j/2} \left(\sum_{i=(k-1)2^j+1}^{(k-1)2^j+2^j-1} X_i - \sum_{i=(k-1)2^j+2^j-1+1}^{k2^j} X_i \right).$$

Defining $\mathbf{d}_j = (d_{j,k})_{k=1}^{2^{J-j}}$, the Haar transform H of \mathbf{X} is $H(\mathbf{X}) = (\mathbf{d}_1, \dots, \mathbf{d}_J, s_{J,1})$. The ‘‘pyramid’’ algorithm in formulae (1) enables the computation of $H(\mathbf{X})$ in $O(n)$ operations. $H(\mathbf{X})$ is an orthonormal transform of \mathbf{X} and can be inverted by undoing steps (1). If the mean signal $\Theta = (\theta_1, \dots, \theta_n)$ is piecewise-constant, then those coefficients $d_{j,k}$ that correspond to the locally constant segments of Θ will be zero-centered. This justifies the following procedure for estimating the mean vector Θ : take the Haar transform of \mathbf{X} , retain those coefficients $d_{j,k}$ for which $|d_{j,k}| > t$ for a certain threshold t and set the others to zero, then take the inverse Haar transform of the thus-‘‘hard’’-thresholded vector. In the i.i.d. Gaussian noise model, in which $X_k = \theta_k + \varepsilon_k$, where $\varepsilon \sim N(0, \sigma^2)$ with σ^2 assumed known, the operation $|d_{j,k}| > t$ is the likelihood ratio test for the local constancy of Θ in the following sense.

1. Assume $(\theta_u)_{u=(k-1)2^j+1}^{(k-1)2^j+2^{j-1}} = \theta^{(1)}$ for all u , and $(\theta_v)_{v=(k-1)2^j+2^{j-1}+1}^{k2^j} = \theta^{(2)}$ for all v . The indices u (respectively v) are the same as those corresponding to the X_u 's (X_v 's) with positive (negative) weights in $d_{j,k}$.
2. Test $H_0 : \theta^{(1)} = \theta^{(2)}$ against $H_1 : \theta^{(1)} \neq \theta^{(2)}$; the Gaussian likelihood ratio test reduces to $|d_{j,k}| > t$, where t is naturally related to the desired significance level. H_0 can alternatively be phrased as $E(d_{j,k}) = 0$, and H_1 – as $E(d_{j,k}) \neq 0$.

Because under each H_0 , the variable $d_{j,k}$ is distributed as $N(0, \sigma^2)$ due to the orthonormality of the Haar transform, the same t can meaningfully be used across different scales and locations (j, k) .

In models other than Gaussian, the operation $|d_{j,k}| > t$ can typically no longer be interpreted as the likelihood ratio test for the equality of $\theta^{(1)}$ and $\theta^{(2)}$. Moreover, the distribution of $d_{j,k}$ will not generally be the same under each H_0 but will, in many models, vary with the local (unknown) parameters $(\theta_i)_{i=(k-1)2^j+1}^{k2^j}$, which makes the selection of t in the operation $|d_{j,k}| > t$ challenging. This is, for example, the case in our running examples, $X_k \sim \text{Pois}(\lambda_k)$ and $X_k \sim \sigma_k^2 m^{-1} \chi_m^2$, both of which are such that $\text{Var}(X_k)$ is a non-trivial function of $E(X_k)$, which translates into the dependence of $d_{j,k}$ on the local means vector $(\theta_i)_{i=(k-1)2^j+1}^{k2^j}$, even under the null hypothesis $E(d_{j,k}) = 0$.

In the (non-Gaussian) model under consideration, our proposal is to remedy this by replacing the operation $|d_{j,k}| > t$ with a likelihood ratio test for $H_0 : \theta^{(1)} = \theta^{(2)}$ against $H_1 : \theta^{(1)} \neq \theta^{(2)}$ suitable for the distribution at hand. More specifically, denoting by $L(\theta | X_{k_1}, \dots, X_{k_2})$ the likelihood of the constant parameter θ given the data X_{k_1}, \dots, X_{k_2} and by $\hat{\theta}^{(1)}, \hat{\theta}^{(2)}$ the MLEs of $\theta^{(1)}, \theta^{(2)}$, respectively, we design a new Haar-like transform, in which we replace the “test statistic” $d_{j,k}$ by

$$g_{j,k} = \text{sign}(\hat{\theta}^{(1)} - \hat{\theta}^{(2)}) \left\{ 2 \log \frac{\sup_{\theta^{(1)}} L(\theta^{(1)} | X_{(k-1)2^j+1}, \dots, X_{(k-1)2^j+2^j-1}) \sup_{\theta^{(2)}} L(\theta^{(2)} | X_{(k-1)2^j+2^j-1+1}, \dots, X_{k2^j})}{\sup_{\theta} L(\theta | X_{(k-1)2^j+1}, \dots, X_{k2^j})} \right\}^{1/2}, \quad (2)$$

the signed and square-rooted generalized log-likelihood ratio statistic for testing H_0 against H_1 . The rationale is that by Wilks’ theorem, under H_0 , this quantity will asymptotically be distributed as $N(0, 1)$ for a class of models that includes, amongst others, our two running examples. We refer to $g_{j,k}$ as the *likelihood ratio Haar coefficient* of \mathbf{X} at scale j and location k . By performing this replacement, we tailor-make a new Haar transform suitable for the distribution of the input vector.

2.1 General methodology for smoothing

We now outline the general methodology for signal smoothing (denoising) involving likelihood ratio Haar wavelets. The problem is to estimate Θ from \mathbf{X} . Let \mathbb{I} be the indicator function. The basic smoothing algorithm proceeds as follows.

1. With \mathbf{X} on input, compute the coefficients $s_{j,k}$, $d_{j,k}$ and $g_{j,k}$ as defined by (1) and (2).
2. Estimate each $\mu_{j,k} := E(d_{j,k})$ by

$$\hat{\mu}_{j,k} = \begin{cases} 0 & j = 1, \dots, J_0, \\ d_{j,k} \mathbb{I}(|g_{j,k}| > t) & j = J_0 + 1, \dots, J. \end{cases} \quad (3)$$

3. Defining $\hat{\mu}_j = (\hat{\mu}_{j,k})_{k=1}^{2^{J-j}}$, compute the inverse Haar transform of the vector $(\hat{\mu}_1, \dots, \hat{\mu}_J, s_{J,1})$ and use it as the estimate $\hat{\Theta}$ of Θ .

We set $\hat{\mu}_{j,k} = 0$ at the finest scales because of a certain strong asymptotic normality argument; see the proof of Theorem 4.1. This theorem also specifies the permitted magnitude of J_0 . The operation in the second line of (3) is referred to as hard thresholding; soft thresholding, in which the surviving coefficients are shrunk towards zero, is also possible. The threshold t is a tuning parameter of the procedure and we discuss its selection later. The above algorithm differs from the standard smoothing using Haar wavelets in that we use $g_{j,k}$, rather than $d_{j,k}$, as the thresholding statistic.

2.2 General methodology for variance stabilization and normalization

Due to the fact that $g_{j,k}$ will typically be distributed as close to $N(0,1)$ under each H_0 (that is, for the majority of scales j and locations k), replacing the coefficients $d_{j,k}$ with $g_{j,k}$ can be viewed as “normalizing” or “Gaussianizing” the input signal in the Haar wavelet domain. The standard inverse Haar transform will then yield a normalized version of the input signal. We outline the basic algorithm below.

1. With \mathbf{X} on input, compute the coefficients $s_{j,k}$ and $g_{j,k}$ as defined by (1) and (2).
2. Defining $\mathbf{g}_j = (g_{j,k})_{k=1}^{2^{J-j}}$, compute the inverse Haar transform of the vector $(\mathbf{g}_1, \dots, \mathbf{g}_J, s_{J,1})$ and denote the resulting vector by $G(\mathbf{X})$.

Throughout the paper, we will be referring to $G(\mathbf{X})$ as the *likelihood ratio Haar transform* of X . In the online supplement, we show that the likelihood Haar transform is invertible, at least in the Poisson and chi-squared cases. An invertible variance-stabilization transformation such as $G(\mathbf{X})$ is useful as it enables the smoothing of X in a modular way: (i) apply $G(X)$, (ii) use any smoother suitable for i.i.d. standard normal noise, (iii) take the inverse of $G(X)$ to obtain a smoothed version of X .

3 Specific examples: Poisson and chi-squared

For $X_i \sim \text{Pois}(\lambda)$, we have $P(X_i = k) = \exp(-\lambda) \frac{\lambda^k}{k!}$ for $k = 0, 1, \dots$, and if $X_s, \dots, X_e \sim \text{Pois}(\lambda)$, then the MLE $\hat{\lambda}$ of λ is $\bar{X}_s^e = \frac{1}{e-s+1} \sum_{i=s}^e X_i$. This, after straightforward algebra, leads to

$$g_{j,k} = \text{sign}(\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^j-1} - \bar{X}_{(k-1)2^j+2^j-1+1}^{k2^j}) 2^{j/2} \times \left\{ \log(\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^j-1}) \bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^j-1} + \log(\bar{X}_{(k-1)2^j+2^j-1+1}^{k2^j}) \bar{X}_{(k-1)2^j+2^j-1+1}^{k2^j} - 2 \log(\bar{X}_{(k-1)2^j+1}^{k2^j}) \bar{X}_{(k-1)2^j+1}^{k2^j} \right\}^{1/2}, \quad (4)$$

using the convention $0 \log 0 = 0$. For $X_i \sim \sigma_i^2 m^{-1} \chi_m^2 = \Gamma(m/2, m/(2\sigma_i^2))$, if $X_s, \dots, X_e \sim \Gamma(m/2, m/(2\sigma^2))$, then the MLE $\hat{\sigma}^2$ of σ^2 is $\bar{X}_s^e = \frac{1}{e-s+1} \sum_{i=s}^e X_i$. This gives

$$g_{j,k} = \text{sign}(\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^j-1} - \bar{X}_{(k-1)2^j+2^j-1+1}^{k2^j}) 2^{j/2} \times \left\{ m \left[\log(\bar{X}_{(k-1)2^j+1}^{k2^j}) - \frac{1}{2} \log(\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^j-1}) - \frac{1}{2} \log(\bar{X}_{(k-1)2^j+2^j-1+1}^{k2^j}) \right] \right\}^{1/2}. \quad (5)$$

Up to the multiplicative factor $m^{1/2}$, the form of the transform in (5) is the same for any m , which means that the chi-squared likelihood ratio Haar coefficients $g_{j,k}$ (computed with an arbitrary m) also achieve variance stabilization if m is unknown (but possibly to a constant different from one). In both the Poisson and the chi-squared cases, $g_{j,k}$ is a function of the local means $\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^j-1}$ and $\bar{X}_{(k-1)2^j+2^j-1+1}^{k2^j}$, which is unsurprising as these are sufficient statistics for the corresponding population means in both these distributions. These local means and, therefore, the coefficients $g_{j,k}$, can be computed in computational time $O(n)$ using the pyramid algorithm in formulae (1).

4 L_2 theory for the likelihood ratio Haar Poisson smoother

In this section, we provide a theoretical mean-square analysis of the performance of the signal smoothing algorithm involving likelihood ratio Haar wavelets, described in Section 2.1. Although we focus on the Poisson distribution, the statement of the result and the mechanics of the proof will be similar for certain other distributions, including scaled chi-squared. The following result holds.

Theorem 4.1 *Let $\Lambda = (\lambda_1, \dots, \lambda_n)$ be a positive piecewise-constant vector, i.e. let there exist up to N indices η_1, \dots, η_N for which $\lambda_{\eta_i} \neq \lambda_{\eta_{i-1}}$. Let $n = 2^J$, where J is a positive integer. Assume Λ is bounded from above and away from zero, and denote $\bar{\Lambda} := \max_i \lambda_i$, $\underline{\Lambda} := \min_i \lambda_i$, $\Lambda' = \bar{\Lambda} - \underline{\Lambda}$ and $\bar{\lambda}_s^e = \frac{1}{e-s+1} \sum_{i=s}^e \lambda_i$. Let $X_k \sim \text{Pois}(\lambda_k)$ for $k = 1, \dots, n$. Let $\hat{\Lambda}$ be obtained as in the algorithm of Section 2.1, using threshold t and with a fixed $\beta \in (0, 1)$ such that $J_0 = \lfloor \log_2 n^\beta \rfloor$. Then, with $d_{j,k}$ and $\mu_{j,k}$ defined in the algorithm of Section 2.1 and with $\bar{X}_s^e = \frac{1}{e-s+1} \sum_{i=s}^e X_i$, on set $\mathcal{A} \cap \mathcal{B}$, where*

$$\begin{aligned} \mathcal{A} &= \{\forall j = J_0 + 1, \dots, J, k = 1, \dots, 2^{J-j} \quad (\bar{\lambda}_{(k-1)2^{j+1}}^{k2^j})^{-1/2} |d_{j,k} - \mu_{j,k}| < t_1\}, \\ \mathcal{B} &= \{\forall j = J_0, \dots, J, k = 1, \dots, 2^{J-j} \quad 2^{j/2} (\bar{\lambda}_{(k-1)2^{j+1}}^{k2^j})^{-1/2} |\bar{X}_{(k-1)2^{j+1}}^{k2^j} - \bar{\lambda}_{(k-1)2^{j+1}}^{k2^j}| < t_2\}, \end{aligned}$$

whose probability approaches 1 as $n \rightarrow \infty$ if $t_1 = C_1 \log^{1/2} n$ and $t_2 = C_2 \log^{1/2} n$ with $C_1 > \{2(1-\beta)\}^{1/2}$ and $C_2 > \{2(1-\beta)\}^{1/2}$, if threshold t is such that

$$t \geq \frac{t_1}{(1 - t_2 2^{-\frac{J_0+1}{2}} \underline{\Lambda}^{-1/2})^{1/2}}, \quad (6)$$

we have

$$\begin{aligned} n^{-1} \|\hat{\Lambda} - \Lambda\|^2 &\leq \\ &\frac{1}{2} n^{-1} N (\Lambda')^2 (n^\beta - 1) + 2n^{-1} N \bar{\Lambda}^{1/2} \{(J - J_0)(t^2 + t_1^2) \bar{\Lambda}^{1/2} + t^2 t_2 (2 + 2^{1/2}) n^{-\beta/2}\} + n^{-1} t_1^2 \bar{\lambda}_1^n, \end{aligned}$$

where $\|\cdot\|$ is the l_2 -norm of its argument.

Bearing in mind the magnitudes of t_2 and J_0 , we can see that the term $t_2 2^{-\frac{J_0+1}{2}} \underline{\Lambda}^{-1/2}$ becomes arbitrarily close to zero for large n , and therefore, from formula (6), the threshold constant t can become arbitrarily close to t_1 . In particular, it is safe to set t to be the ‘‘universal’’ threshold suitable for iid $N(0, 1)$ noise (Donoho and Johnstone, 1994), that is $t = \{2 \log n\}^{1/2}$. It is in this sense that our

likelihood ratio Haar construction achieves variance stabilization and normalization: in order to denoise Poisson signals in which the variance of the noise depends on the local mean, we make it possible to use the universal Gaussian threshold, as if the noise were Gaussian with variance one. In classical Haar wavelet thresholding with $|d_{j,k}| > \tilde{t}$ as the thresholding decision, \tilde{t} would have to depend on the level of the Poisson intensity Λ over the support of $d_{j,k}$, which is unknown; our approach circumvents this.

If the number N of change-points does not increase with the sample size n , then the dominant term in the mean-square error is of order $O(n^{\beta-1})$. This suggests that β should be set to be “arbitrarily small”, in which case the MSE becomes arbitrarily close to the parametric rate $O(n^{-1})$.

5 Practical performance

In the online supplement, we demonstrate that the likelihood ratio Haar coefficients appear to offer better normalization and variance stabilization than the Fisz coefficients. In this section, we show that this translates into better MSE properties of the likelihood ratio Haar smoother than the analogous Haar-Fisz smoother, in both the Poisson and the exponential models, on the examples considered. For comprehensive comparison of the performance of the Haar-Fisz smoother to that of other techniques, see Fryzlewicz and Nason (2004), Besbeas et al. (2004) and Fryzlewicz (2008), among others. Our test signals are [1] Donoho and Johnstone’s (1994) `blocks` and [2] `bumps` functions, scaled to have (min, max) of [1] (0.681, 27.029) and [2] (1, 12.565), both of length $n = 2048$. We consider the following models: **(1a)**, **(2a)**: Poisson models, in which the signals [1], [2] (respectively) play the role of the Poisson intensity Λ , so that $X_k \sim \text{Pois}(\lambda_k)$; **(1b)**, **(2b)**: Exponential models, in which the signals [1], [2] (respectively) play the role of the exponential mean σ^2 , so that $X_k \sim \sigma_k^2 \text{Exp}(1) = \sigma_k^2 2^{-1} \chi_2^2$.

For all models, we compare the MSE performance of “like-for-like” likelihood ratio Haar and Haar-Fisz smoothers, both constructed as described in Section 2.1, except the Haar-Fisz smoother uses the corresponding coefficients $f_{j,k}$ in place of $g_{j,k}$. We use the non-decimated (translation invariant,

Method \ Model	(1a)	(1b)	(2a)	(2b)
Haar-Fisz	0.615	8.647	0.357	1.053
Likelihood ratio Haar	0.605	7.958	0.341	0.905

Table 1: MSE over 1000 simulations for the two methods and four models described in Section 5.

stationary, maximum overlap) Haar transform (Nason and Silverman, 1995) to achieve fast averaging over all possibly cyclic shifts of the input data. For better comparison of the effects of thresholding alone, we use $J_0 = 0$. We use the universal threshold $t = \{2 \log n\}^{1/2}$. Figures 1 and 2 show sample reconstructions for the likelihood ratio Haar method in the Poisson models (1a), (2a).

Table 1 shows that the likelihood ratio Haar smoother outperforms Haar-Fisz across all the models considered. For the Poisson models, the improvement is fairly modest (2% for `blocks`, 4% for `bumps`) but for the exponential models, it is more significant (8% for `blocks`, 14% for `bumps`). One important reason for this improved performance is that as demonstrated in the online supplement, the likelihood ratio Haar coefficients have a higher magnitude than the corresponding Fisz coefficients, and therefore more easily survive thresholding. This implies that the likelihood ratio Haar smoother lets through “more signal” compared to the Haar-Fisz smoother if both use the same threshold, however chosen. Another possible reason is that as shown in the online supplement, the likelihood ratio Haar coefficients are closer to variance-one normality than the Fisz coefficients and therefore the use of thresholds designed for standard normal noise may be more suitable for them.

We now briefly illustrate the normalizing and variance-stabilizing properties of the likelihood ratio Haar transform $G(\cdot)$ described in Section 2.2, using data simulated from models (1a) and (1b). We use the non-decimated version of the Haar transform.

Figure 3 illustrates the results for the Poisson case. In both the Poisson and the exponential examples, the likelihood ratio Haar transform is a very good normalizer and variance-stabilizer: the transformed data minus the transformed signal shows good agreement with an i.i.d. normal sample; its sample variance equals 1.07 for the Poisson model and 1.14 for the exponential model. Particularly

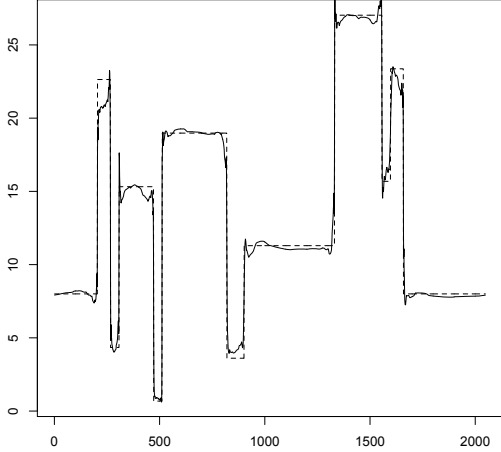


Figure 1: Sample likelihood ratio Haar reconstruction in model (1a), see Section 5 for details.

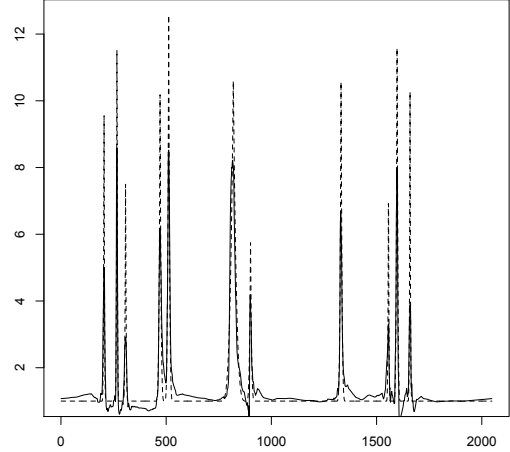


Figure 2: Sample likelihood ratio Haar reconstruction in model (2a), see Section 5 for details.

for the exponential model, the likelihood ratio Haar transform is a significantly better normalizer than the Haar-Fisz transform (not shown here).

5.1 California earthquake data

In this section, we revisit the Northern California earthquake dataset, analysed in Fryzlewicz and Nason (2004) and available from <http://quake.geo.berkeley.edu/ncedc/catalog-search.html>. We analyze the time series N_k , $k = 1, \dots, 1024$, where N_k is the number of earthquakes of magnitude 3.0 or more which occurred in the k th week, the first week under consideration starting April 22nd, 1981 and the final ending December 5th, 2000. We assume $N_k \sim \text{Pois}(\lambda_k)$ and estimate Λ using our likelihood ratio Haar smoother, used as described in Section 5.

The estimate and the data are shown in Figure 4. The appearance of the estimator reveals an interesting phenomenon, not necessarily easily seen in the noisy data: for many of the intensity spikes observed in this dataset, the intensity in the time period just before the spike appears to be much lower

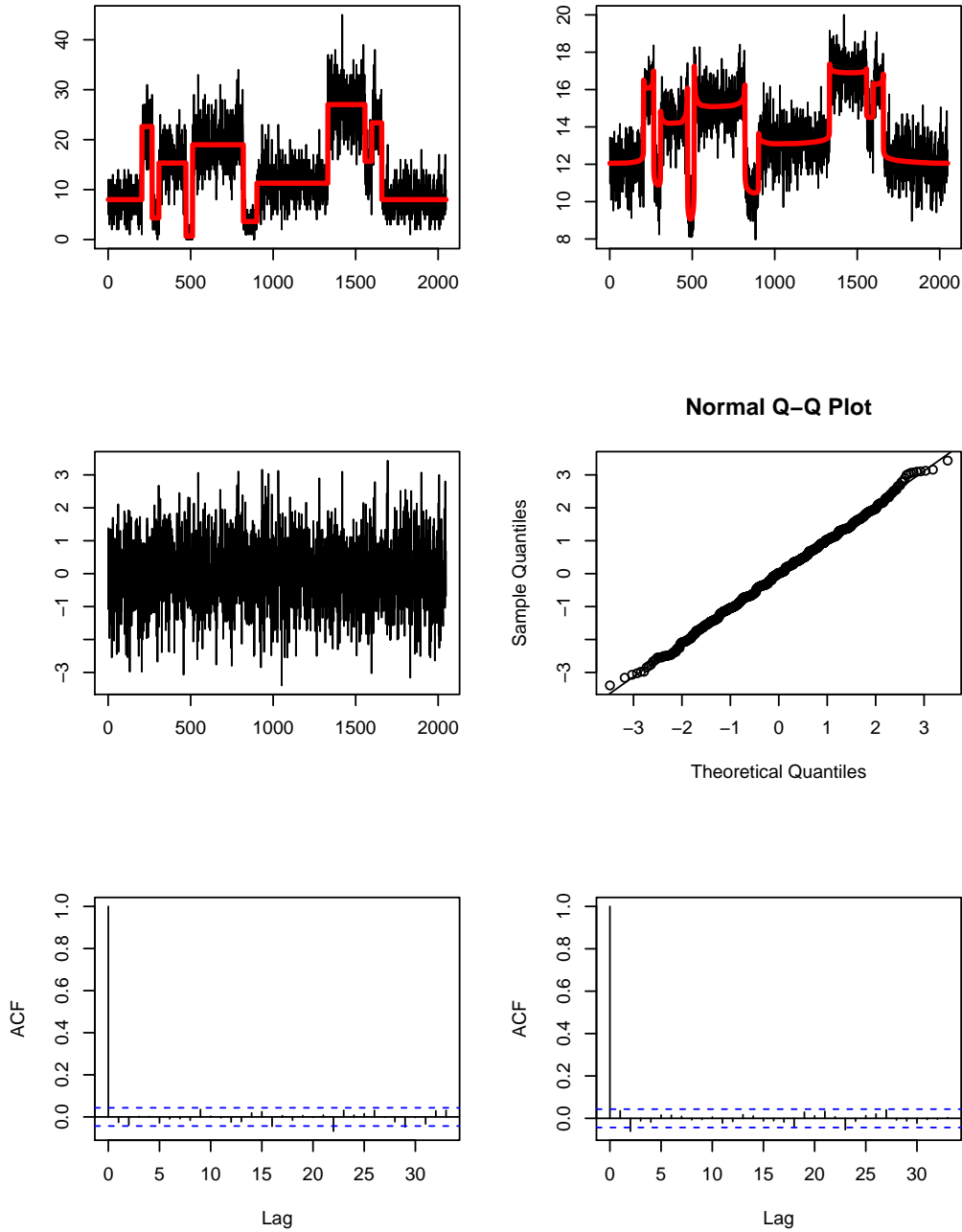


Figure 3: The Poisson model. Top left: Poisson intensity Λ (red) and simulated data \mathbf{X} (black). Top right: the likelihood ratio transform $G(\Lambda)$ (red) and $G(\mathbf{X})$ (black). Middle left: $G(\mathbf{X}) - G(\Lambda)$. Middle right: Q-Q plot of $G(\mathbf{X}) - G(\Lambda)$ against the normal quantiles. Bottom left: sample acf plot of $G(\mathbf{X}) - G(\Lambda)$. Bottom right: sample acf plot of $(G(\mathbf{X}) - G(\Lambda))^2$.

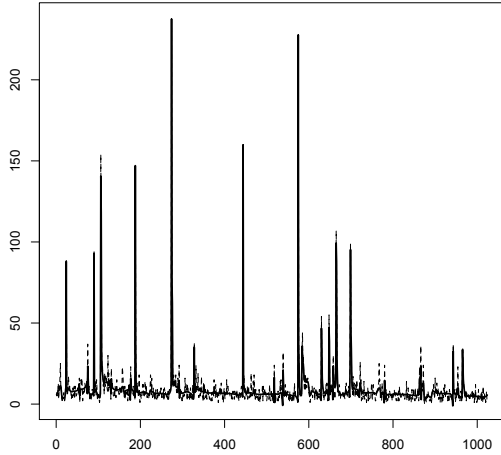


Figure 4: Northern California earthquake data: N_k (dashed) and the likelihood ratio Haar estimate (thick solid). See Section 5.1 for details.

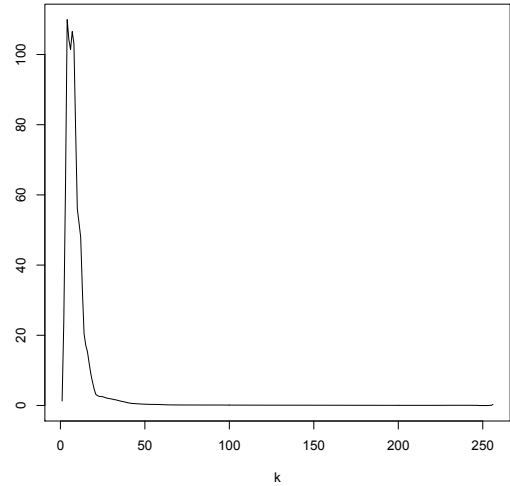


Figure 5: The likelihood ratio Haar smooth of M_k under the Poisson assumption.

than the intensity in the period following the spike, which may point to a degree of persistence in the seismic activity following the major spikes in activity observed in these data.

Further, we analyse the histogram of counts M_k , $k = 0, \dots, 255$, defined as the number of weeks in which k earthquakes of magnitude 3.0 or more which occurred. The raw data (not shown here) show an apparent bimodality with modes at 4 and 6. To verify whether this is a spurious or “real” effect, we smooth M_k using our likelihood ratio Haar smoother suitable for Poisson data (note that M_k , being a histogram, can approximately be modelled as Poisson-distributed). Figure 5 reveals that our fit preserves the bimodality, which gives support to the argument that this is a genuine, rather than spurious, effect. This finding points towards a mixture model with two components: one corresponding to “quieter” periods (i.e. those with a low intensity of earthquakes of magnitude 3.0 or more) and the other to periods with high earthquake intensity.

References

- F. Anscombe. The transformation of Poisson, binomial and negative-binomial data. *Biometrika*, 35: 246–254, 1948.
- A. Antoniadis and T. Sapatinas. Wavelet shrinkage for natural exponential families with quadratic variance functions. *Biometrika*, 88:805–820, 2001.
- A. Antoniadis, P. Besbeas, and T. Sapatinas. Wavelet shrinkage for natural exponential families with cubic variance functions. *Sankhya Ser. A*, 63:309–327, 2001.
- P. Besbeas, I. De Feis, and T. Sapatinas. A comparative study of wavelet shrinkage estimators for Poisson counts. *Int. Statist. Review*, 72:209–237, 2004.
- L. Brown, T. Cai, and H. Zhou. Nonparametric regression in exponential families. *Annals of Statistics*, 38:2005–2046, 2010.
- I. Daubechies. *Ten Lectures on Wavelets*. SIAM, Philadelphia, Pa., 1992.
- D. Donoho. Nonlinear wavelet methods for recovery of signals, densities, and spectra from indirect and noisy data. In *Proc. Sympos. Appl. Math.*, pages 173–205. Amer. Math. Soc., 1993.
- D. Donoho and I. Johnstone. Ideal spatial adaptation by wavelet shrinkage. *Biometrika*, 81:425–455, 1994.
- M. Fisz. The limiting distribution of a function of two independent random variables and its statistical application. *Colloquium Mathematicum*, 3:138–146, 1955.
- A. Foi. Optimization of variance-stabilizing transformations. *Preprint*, 2009.
- J.-M. Freyermuth, H. Ombao, and R. von Sachs. Tree-structured wavelet estimation in a mixed effects model for spectra of replicated time series. *Journal of the American Statistical Association*, 105: 634–646, 2010.
- P. Fryzlewicz. Data-driven wavelet-Fisz methodology for nonparametric function estimation. *Electronic Journal of Statistics*, 2:863–896, 2008.

- P. Fryzlewicz and G. Nason. A Haar-Fisz algorithm for Poisson intensity estimation. *Journal of Computational and Graphical Statistics*, 13:621–638, 2004.
- P. Fryzlewicz and G. Nason. Haar-Fisz estimation of evolutionary wavelet spectra. *Journal of the Royal Statistical Society Series B*, 68:611–634, 2006.
- P. Fryzlewicz, T. Sapatinas, and S. Subba Rao. A Haar-Fisz technique for locally stationary volatility estimation. *Biometrika*, 93:687–704, 2006.
- P. Fryzlewicz, V. Delouille, and G. Nason. GOES-8 X-ray sensor variance stabilization using the multiscale data-driven Haar-Fisztransform. *J. Roy. Statist. Soc. C*, to appear, 2007.
- P. Fryzlewicz, G. Nason, and R. von Sachs. A wavelet-Fisz approach to spectrum estimation. *Journal of Time Series Analysis*, 29:868–880, 2008.
- H-Y. Gao. Choice of thresholds for wavelet shrinkage estimate of the spectrum. *J. Time Ser. Anal.*, 18:231–252, 1997.
- P. Hall and P. Patil. On the choice of smoothing parameter, threshold and truncation in nonparametric regression by non-linear wavelet methods. *Journal of the Royal Statistical Society Series B*, 58:361–377, 1996a.
- P. Hall and P. Patil. Effect of threshold rules on performance of wavelet-based curve estimators. *Statistica Sinica*, 6:331–345, 1996b.
- P. Hall and B. Turlach. Interpolation methods for nonlinear wavelet regression with irregularly spaced design. *Annals of Statistics*, 25:1912–1925, 1997.
- P. Hall, I. McKay, and B. Turlach. Performance of wavelet methods for functions with many discontinuities. *Annals of Statistics*, 24:2462–2476, 1996.
- P. Hall, G. Kerkycharian, and D. Picard. On the minimax optimality of block thresholded wavelet estimators. *Statistica Sinica*, 9:33–49, 1999.
- K. Hirakawa and P. Wolfe. Skellam shrinkage: Wavelet-based intensity estimation for inhomogeneous Poisson data. *IEEE Transactions on Information Theory*, 58:1080–1093, 2012.

- M. Jansen. Multiscale Poisson data smoothing. *J. R. Statist. Soc. B*, 68:27–48, 2006.
- E. Kolaczyk. Non-parametric estimation of Gamma-ray burst intensities using Haar wavelets. *The Astrophysical Journal*, 483:340–349, 1997.
- E. Kolaczyk. Wavelet shrinkage estimation of certain Poisson intensity signals using corrected thresholds. *Statistica Sinica*, 9:119–135, 1999a.
- E. Kolaczyk. Bayesian multiscale models for Poisson processes. *Journal of the American Statistical Association*, 94:920–933, 1999b.
- E. Kolaczyk and R. Nowak. Multiscale likelihood analysis and complexity penalized estimation. *Ann. Statist.*, 32:500–527, 2004.
- F. Luisier, C. Vonesch, T. Blu, and M. Unser. Fast interscale wavelet denoising of Poisson-corrupted images. *Signal Processing*, 90:415–427, 2010.
- S. Mallat. A theory for multiresolution signal decomposition: the wavelet representation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 11:674–693, 1989.
- E. Motakis, G. Nason, P. Fryzlewicz, and G. Rutter. Variance stabilization and normalization for one-color microarray data using a data-driven multiscale approach. *Bioinformatics*, 22:2547–2553, 2006.
- P. Moulin. Wavelet thresholding techniques for power spectrum estimation. *IEEE Trans. Sig. Proc.*, 42:3126–3136, 1994.
- G. Nason. *Wavelet Methods in Statistics with R*. Springer, New York, 2008.
- G. Nason. Multiscale variance stabilization via maximum likelihood. *Biometrika*, 101:499–504, 2014.
- G. Nason and B. Silverman. The stationary wavelet transform and some statistical applications. In A. Antoniadis and G. Oppenheim, editors, *Lecture Notes in Statistics*, vol. 103, pages 281–300. Springer-Verlag, 1995.
- M. Neumann. Spectral density estimation via nonlinear wavelet methods for stationary non-Gaussian time series. *Journal of Time Series Analysis*, 17:601–633, 1996.

- R. Nowak and R. Baraniuk. Wavelet domain filtering for photon imaging systems. *IEEE Transactions on Image Processing*, 8:666–678, 1999.
- M. Pensky, B. Vidakovic, and D. De Canditis. Bayesian decision theoretic scale-adaptive estimation of a log-spectral density. *Statistica Sinica*, 17:635–666, 2007.
- P. Reynaud-Bouret and V. Rivoirard. Near optimal thresholding estimation of a Poisson intensity on the real line. *Electronic Journal of Statistics*, 4:172–238, 2010.
- K. Timmermann and R. Nowak. Multiscale Bayesian estimation of Poisson intensities. In *Proceedings of the Asilomar Conference on Signals, Systems and Computers*, pages 95–90. IEEE Computer Press, Pacific Grove, CA, 1997.
- K. Timmermann and R. Nowak. Multiscale modeling and estimation of Poisson processes with application to photon-limited imaging. *IEEE Transactions on Information Theory*, 45:846–862, 1999.
- B. Vidakovic. *Statistical Modeling by Wavelets*. Wiley, New York, 1999.
- M. Wickerhauser. *Adapted Wavelet Analysis from Theory to Software*. AK Peters, 1994.
- R. Willett and R. Nowak. Platelets: a multiscale approach for recovering edges and surfaces in photon-limited medical imaging. *IEEE Transactions on Medical Imaging*, 22:332–350, 2003.
- B. Zhang, J. Fadili, and J.-L. Starck. Wavelets, ridgelets, and curvelets for Poisson noise removal. *IEEE Transactions on Image Processing*, 17:1093–1108, 2008.

Likelihood ratio Haar variance stabilization and normalization for Poisson and other non-Gaussian noise removal – online supplement

Piotr Fryzlewicz*

June 15, 2017

1 Links between likelihood ratio Haar wavelets and the Haar-Fisz methodology (with numerical examples)

This section compares the likelihood ratio Haar coefficients $g_{j,k}$, defined in the general, Poisson and chi-squared cases in formulae (2), (4) and (5) of the main paper, respectively, to the Fisz coefficients $f_{j,k}$ (Fryzlewicz and Nason, 2004), which the above work defines as the Haar coefficients $d_{j,k}$ divided by the maximum likelihood estimates of their own standard deviation under the null hypothesis $E(d_{j,k}) = 0$. We start with the Poisson case and note that by Fryzlewicz and Nason (2004), $f_{j,k}$ is then expressed as

$$f_{j,k} = 2^{j/2-1} \frac{\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^j-1} - \bar{X}_{(k-1)2^j+2^j-1+1}^{k2^j}}{\sqrt{\bar{X}_{(k-1)2^j+1}^{k2^j}}}.$$

We first note that $\text{sign}(g_{j,k}) = \text{sign}(f_{j,k})$ and that Lemma 3.2, used with $f(u) = u \log u$; $f(0) = 0$ in the notation of that lemma, reduces to $|g_{j,k}| \geq |f_{j,k}|$. Moreover, since the inequality in Lemma 3.2 arises as a simple application of Jensen's inequality to the convex function $f(\cdot)$, it is intuitively apparent that the less convexity in $f(\cdot)$, the closer $g_{j,k}$ will be to $f_{j,k}$. Noting that $f''(u) = u^{-1}$ and therefore the degree of convexity in $f(u)$ decreases as u increases, it can heuristically be observed that $g_{j,k}$ and $f_{j,k}$ should be closer to each other for larger values of $\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^j-1}$ and $\bar{X}_{(k-1)2^j+2^j-1+1}^{k2^j}$ (i.e. for high Poisson intensities), and further apart otherwise.

To illustrate this phenomenon and other interesting similarities and differences between the Fisz and the likelihood ratio Haar coefficients in the Poisson case, consider the following two numerical examples, in which we simulate 1000 realisations of $\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^j-1}$

*Department of Statistics, London School of Economics, Houghton Street, London WC2A 2AE, UK. Email: p.fryzlewicz@lse.ac.uk. Work supported by the Engineering and Physical Sciences Research Council grant no. EP/L014246/1.

and $\bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j}$ and compute the corresponding 1000 realisations of $\{g_{j,k}^{(i)}\}_{i=1}^{1000}$ and $\{f_{j,k}^{(i)}\}_{i=1}^{1000}$.

- $j = 2$, $E(\bar{X}_{(k-1)2^j+2^{j-1}+1}^{(k-1)2^j+2^{j-1}}) = 10$, $E(\bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j}) = 10.5$. As is apparent from Figure 1, the values of $g_{j,k}^{(i)} - f_{j,k}^{(i)}$ are close to zero. Figure 2 provides further evidence that the empirical distributions of $f_{j,k}^{(i)}$ and $g_{j,k}^{(i)}$ are difficult to distinguish by the naked eye. Q-q plots (not shown) exhibit good agreement for both $g_{j,k}^{(i)}$ and $f_{j,k}^{(i)}$ with the normal distribution, and we have $\widehat{\text{Var}}(g_{j,k}^{(i)}) = 1.06$ and $\widehat{\text{Var}}(f_{j,k}^{(i)}) = 1.05$, which provides evidence that both the likelihood ratio Haar coefficients and the Fisz coefficients achieve good variance stabilization in this high-intensity case.
- $j = 2$, $E(\bar{X}_{(k-1)2^j+2^{j-1}+1}^{(k-1)2^j+2^{j-1}}) = 0.2$, $E(\bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j}) = 0.7$. Figures 3 and 4 demonstrate that in this low-intensity case, the distributions of $f_{j,k}^{(i)}$ and $g_{j,k}^{(i)}$ are now further apart. The Fisz coefficients and the likelihood ratio Haar coefficients seem to be similarly close to the normal distribution, with the empirical skewness and kurtosis for $f_{j,k}^{(i)}$ being 0.39 and 2.52 (respectively) and those for $g_{j,k}^{(i)}$ being 0.35 and 2.53 (respectively). However, the likelihood ratio Haar coefficients achieve far better variance stabilization in this low-intensity example: we have $\widehat{\text{Var}}(g_{j,k}^{(i)}) = 0.92$ versus $\widehat{\text{Var}}(f_{j,k}^{(i)}) = 0.68$.

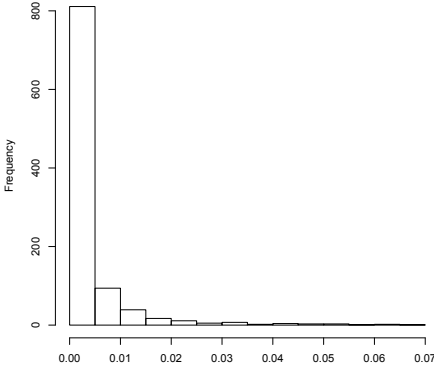


Figure 1: The Poisson case. Histogram of the empirical distribution of $\{|g_{j,k}^{(i)}| - |f_{j,k}^{(i)}|\}_{i=1}^{1000}$ with $j = 2$, $E(\bar{X}_{(k-1)2^j+2^{j-1}+1}^{(k-1)2^j+2^{j-1}}) = 10$, $E(\bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j}) = 10.5$.

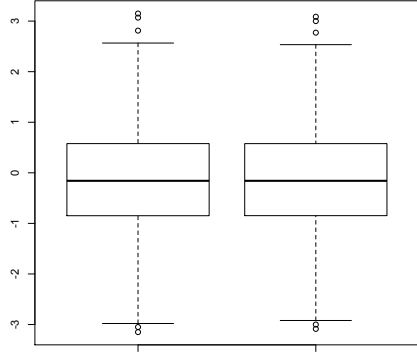


Figure 2: The Poisson case. Boxplots of the empirical distributions of $\{g_{j,k}^{(i)}\}_{i=1}^{1000}$ (left) and $\{f_{j,k}^{(i)}\}_{i=1}^{1000}$ (right) with $j = 2$, $E(\bar{X}_{(k-1)2^j+2^{j-1}+1}^{(k-1)2^j+2^{j-1}}) = 10$, $E(\bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j}) = 10.5$.

We now turn to the chi-squared distribution. The Fisz coefficients for the $\sigma^2\chi_1^2$ distribution are derived in Fryzlewicz et al. (2006), those for the exponential distribution ($\sigma^22^{-1}\chi_2^2$) appear in Fryzlewicz et al. (2008) and the general case $\sigma^2m^{-1}\chi_m^2$ is covered in Fryzlewicz (2008). In the general case of the $\sigma^2m^{-1}\chi_m^2$ distribution, using the notation from the current

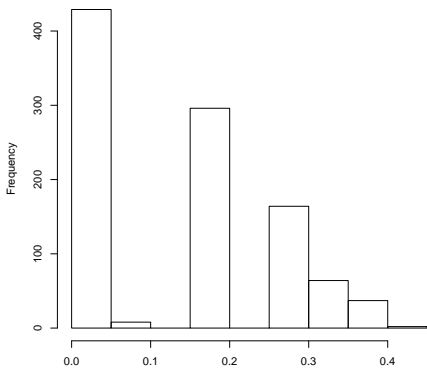


Figure 3: The Poisson case. Histogram of the empirical distribution of $\{|g_{j,k}^{(i)}| - |f_{j,k}^{(i)}|\}_{i=1}^{1000}$ with $j = 2$, $E(\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}}) = 0.2$, $E(\bar{X}_{(k-1)2^j+1}^{k2^j}) = 0.7$.

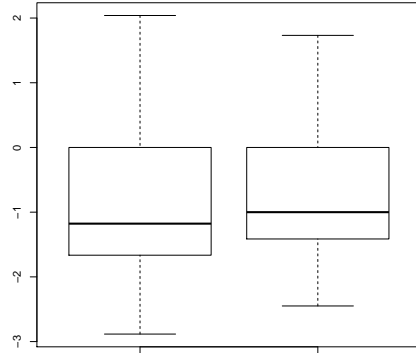


Figure 4: The Poisson case. Boxplots of the empirical distributions of $\{|g_{j,k}^{(i)}|\}_{i=1}^{1000}$ (left) and $\{f_{j,k}^{(i)}\}_{i=1}^{1000}$ (right) with $j = 2$, $E(\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}}) = 0.2$, $E(\bar{X}_{(k-1)2^j+1}^{k2^j}) = 0.7$.

paper, the Fisz coefficients $f_{j,k}$ are expressed as

$$f_{j,k} = 2^{\frac{j-3}{2}} m^{1/2} \frac{\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}} - \bar{X}_{(k-1)2^j+1}^{k2^j}}{\bar{X}_{(k-1)2^j+1}^{k2^j}}. \quad (1)$$

As in the Poisson case, we obviously have $\text{sign}(g_{j,k}) = \text{sign}(f_{j,k})$. Lemma 3.2, used with $f(u) = -\log u$ in the notation of that lemma, reduces to $|g_{j,k}| \geq |f_{j,k}|$. Moreover, by the same convexity argument as in the Poisson case, $g_{j,k}$ and $f_{j,k}$ will be closer to each other for larger values of $\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}}$ and $\bar{X}_{(k-1)2^j+1}^{k2^j}$.

A major difference between the Poisson and the chi-square cases is that in the chi-square case, $f_{j,k}$ is a compactly supported random variable (see formula (1)), whereas $g_{j,k}$ is not. This difference does not apply in the Poisson case, in which neither $f_{j,k}$ nor $g_{j,k}$ are compactly supported. This has implications for how quickly $f_{j,k}$ and $g_{j,k}$ approach the normal distribution (with increasing j or m) in the chi-square case, and we illustrate this numerically below.

As before, we simulate 1000 realisations of $\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}}$ and $\bar{X}_{(k-1)2^j+1}^{k2^j}$ and compute the corresponding 1000 realisations of $\{g_{j,k}^{(i)}\}_{i=1}^{1000}$ and $\{f_{j,k}^{(i)}\}_{i=1}^{1000}$. We consider the following four cases.

- $m = 1$, $j = 2$, $E(\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}}) = 10$, $E(\bar{X}_{(k-1)2^j+1}^{k2^j}) = 10.5$. In this case, the likelihood ratio Haar coefficients provide far better variance stabilization and normalization than the Fisz coefficients. For $f_{j,k}^{(i)}$, we have the following empirical values: variance 0.67, skewness 0.03, kurtosis 1.81. For $g_{j,k}^{(i)}$, we have variance 1.29, skewness 0.03, kurtosis 3.06. Figure 5 confirms the superiority of the likelihood ratio

Haar coefficients over the Fisz coefficients as regards their closeness to the normal distribution.

- $m = 1, j = 2, E(\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^j-1}) = 0.2, E(\bar{X}_{(k-1)2^j+2^j-1+1}^{k2^j}) = 0.7$. This low-sigma case differs from the previous one mainly in that both the likelihood ratio Haar coefficients and the Fisz coefficients are skewed to the right, although the Fisz coefficients (much) more so. For $f_{j,k}^{(i)}$, we have the following empirical values: variance 0.59, skewness 0.89, kurtosis 2.70. For $g_{j,k}^{(i)}$, we have variance 1.23, skewness 0.46, kurtosis 3.1. Figure 6 provides further visual evidence of the higher degree of symmetry in the likelihood ratio Haar coefficients and its closeness to the normal distribution.
- $m = 2, j = 2, E(\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^j-1}) = 10, E(\bar{X}_{(k-1)2^j+2^j-1+1}^{k2^j}) = 10.5$. As m increases, both the likelihood ratio Haar coefficients and the Fisz coefficients move closer towards variance-one normality, although again the likelihood ratio Haar coefficients beat Fisz. For $f_{j,k}^{(i)}$, we have the following empirical values: variance 0.81, skewness 0.05, kurtosis 2.19. For $g_{j,k}^{(i)}$, we have variance 1.16, skewness 0.03, kurtosis 2.97. Figure 7 shows both empirical distributions.
- $m = 2, j = 2, E(\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^j-1}) = 0.2, E(\bar{X}_{(k-1)2^j+2^j-1+1}^{k2^j}) = 0.7$. In this low-sigma case also, the likelihood ratio Haar coefficients appear to be far closer to variance-one normality than the Fisz coefficients. For $f_{j,k}^{(i)}$, we have the following empirical values: variance 0.57, skewness 1.15, kurtosis 4.08. For $g_{j,k}^{(i)}$, we have variance 1.04, skewness 0.45, kurtosis 3.64. Figure 8 shows both empirical distributions.

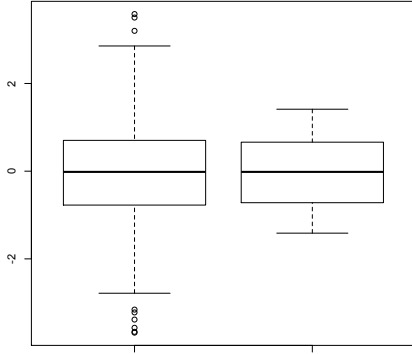


Figure 5: The chi-squared case. Boxplots of the empirical distributions of $\{g_{j,k}^{(i)}\}_{i=1}^{1000}$ (left) and $\{f_{j,k}^{(i)}\}_{i=1}^{1000}$ (right) with $m = 1, j = 2, E(\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^j-1}) = 10, E(\bar{X}_{(k-1)2^j+2^j-1+1}^{k2^j}) = 10.5$.

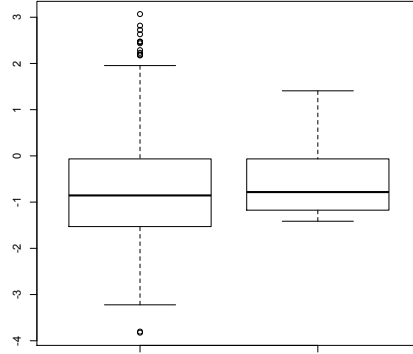


Figure 6: The chi-squared case. Boxplots of the empirical distributions of $\{g_{j,k}^{(i)}\}_{i=1}^{1000}$ (left) and $\{f_{j,k}^{(i)}\}_{i=1}^{1000}$ (right) with $m = 1, j = 2, E(\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^j-1}) = 0.2, E(\bar{X}_{(k-1)2^j+2^j-1+1}^{k2^j}) = 0.7$.

Overall, our empirical observations from the above (and other unreported) numerical exercises are as follows. For fine scales (i.e. those for which j is small) and for low degrees

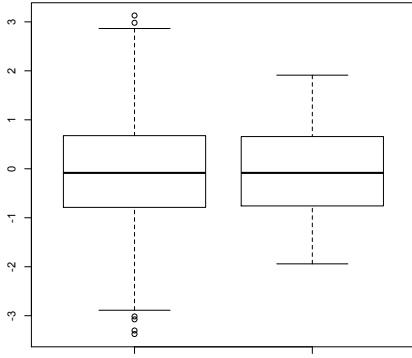


Figure 7: The chi-squared case. Boxplots of the empirical distributions of $\{g_{j,k}^{(i)}\}_{i=1}^{1000}$ (left) and $\{f_{j,k}^{(i)}\}_{i=1}^{1000}$ (right) with $m = 2$, $j = 2$, $E(\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}}) = 10$, $E(\bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j}) = 10.5$.

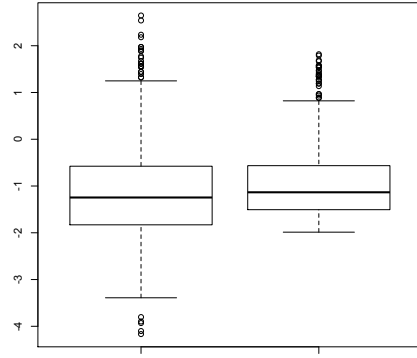


Figure 8: The chi-squared case. Boxplots of the empirical distributions of $\{g_{j,k}^{(i)}\}_{i=1}^{1000}$ (left) and $\{f_{j,k}^{(i)}\}_{i=1}^{1000}$ (right) with $m = 2$, $j = 2$, $E(\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}}) = 0.2$, $E(\bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j}) = 0.7$.

of freedom m , the likelihood ratio Haar coefficients are much closer to a normal variable with variance one than the corresponding Fisz coefficients. From the properties of the chi-squared distribution, the effect of increasing j while keeping m constant is similar to the effect of increasing m while keeping j constant. As m or j increases, the likelihood ratio Haar coefficients appear to move closer to the normal distribution with variance one. However, for the same to happen with Fisz coefficients, the two means, $E(\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}})$ and $E(\bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j})$, need to be relatively close to each other. The latter phenomenon can also be observed in the Poisson case for increasing j . This is not unexpected as the results from Fisz (1955) suggest that the asymptotic normality with variance one arises when the two means approach each other asymptotically; no results are provided in Fisz (1955) on the case in which the two means diverge.

We end this section with an interesting interpretation of Lemmas 3.2 and 3.4 in the case of the Poisson distribution. Note that together, they imply

$$2^{j/2-1} \frac{|\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}} - \bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j}|}{\sqrt{\frac{2}{\frac{1}{\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}}} + \frac{1}{\bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j}}}}} \geq |g_{j,k}| \geq 2^{j/2-1} \frac{|\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}} - \bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j}|}{\sqrt{\frac{1}{2} \left(\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}} + \bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j} \right)}}$$

on in other words, the magnitude of the likelihood ratio Haar coefficient is bounded from below by the magnitude of the corresponding Fisz coefficient and from above by the magnitude of a ‘‘Fisz-like’’ coefficient in which the arithmetic mean in the denominator has been replaced by the harmonic mean.

2 Invertibility of the likelihood Haar transform

Inverting the standard Haar transform proceeds by transforming each pair of coefficients $(s_{j,k}, d_{j,k})$ into $(s_{j-1,2k-1}, s_{j-1,2k})$, hierarchically for $j = J, \dots, 1$ (note that $s_{0,k} = X_k$). Similarly, to demonstrate that the likelihood Haar transform is invertible, we need to show that it is possible to transform $(s_{j,k}, g_{j,k})$ into $(s_{j-1,2k-1}, s_{j-1,2k})$.

We first show the invertibility of the Poisson likelihood ratio Haar transform. Denoting for brevity $u = \bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}}$, $v = \bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j}$ and ignoring some multiplicative constants and the square-root operation in $g_{j,k}$, which are irrelevant for invertibility, this amounts to showing that (u, v) can be uniquely determined from $(u+v)/2$ and $\text{sign}(u-v)\{u \log u + v \log v - (u+v) \log((u+v)/2)\}$. The term $\text{sign}(u-v)$ determines whether $u \leq v$ or vice versa, so assume that $u \leq v$ w.l.o.g. Denoting by a the known value of $u+v$, observe that the function $(a-v) \log(a-v) + v \log v$ is strictly increasing for $v \in [a/2, a]$, which means that v can be determined uniquely and therefore that the Poisson likelihood ratio Haar transform is invertible.

We now show the invertibility of the chi-squared likelihood ratio Haar transform. We denote $u = \bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}}$, $v = \bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j}$ and ignore some multiplicative constants and the square-root operation in $g_{j,k}$ which are irrelevant for invertibility. Assume that $u \leq v$ w.l.o.g. Denoting by a the known value of $u+v$, observe that the function $-\log(a-v) - \log v$ is strictly increasing for $v \in [a/2, a)$, which means that v can be determined uniquely and therefore that the chi-squared likelihood ratio Haar transform is invertible.

3 Technical results including proof of Theorem 4.1 from the main paper

Lemma 3.1 *Let function $f : [u, v] \rightarrow \mathbb{R}$ be such that f' is continuous on $[u, v]$ and f'' is continuous on (u, v) . There exists a point $\xi \in (u, v)$ such that*

$$f(u) - 2f\left(\frac{u+v}{2}\right) + f(v) = \frac{(u-v)^2}{4} f''(\xi).$$

Proof. Let $z = (u+v)/2$ and $\delta = (v-u)/2$, then

$$f(u) - 2f\left(\frac{u+v}{2}\right) + f(v) = f(z-\delta) - 2f(z) + f(z+\delta).$$

Defining $g(x) = f(z-x) - 2f(z) + f(z+x)$, Taylor's theorem yields

$$g(\delta) = g(0) + \delta g'(0) + \frac{\delta^2}{2} g''(\xi') = \frac{\delta^2}{2} g''(\xi') = \frac{\delta^2}{2} \{f''(z+\xi') + f''(z-\xi')\}, \quad (2)$$

where $\xi' \in (0, \delta)$. By the intermediate value theorem, there exists a $\xi \in (z-\xi', z+\xi') \subset [u, v]$ such that $\{f''(z+\xi') + f''(z-\xi')\}/2 = f''(\xi)$, which by (2) completes the result.

Lemma 3.2 *Let function $f : [u, v] \rightarrow \mathbb{R}$ be such that f' is continuous on $[u, v]$ and f'' is convex on (u, v) . Then*

$$f(u) - 2f\left(\frac{u+v}{2}\right) + f(v) \geq \frac{(u-v)^2}{4} f''\left(\frac{u+v}{2}\right).$$

Proof. Straightforward from the convexity of f'' and (2).

Lemma 3.3 *Let function $f : [u, v] \rightarrow \mathbb{R}$ be such that f' is continuous on $[u, v]$ and f'' is nonincreasing on $[u, v]$. Then*

$$f(u) - 2f\left(\frac{u+v}{2}\right) + f(v) \leq \frac{(u-v)^2}{8} \left\{ f''\left(\frac{u+v}{2}\right) + f''(u) \right\}.$$

Proof. Straightforward from (2) and the fact that f'' is nonincreasing on $[u, v]$.

Lemma 3.4 *Let function $f : [u, v] \rightarrow \mathbb{R}$ be such that f' is continuous on $[u, v]$ and f'' is convex on $[u, v]$. Then*

$$f(u) - 2f\left(\frac{u+v}{2}\right) + f(v) \leq \frac{(u-v)^2}{8} \{f''(v) + f''(u)\}.$$

Proof. Straightforward from the convexity of f'' and (2).

Lemma 3.5 *The Poisson distribution satisfies Cramer's conditions.*

Proof. The Poisson distribution is log-concave, and Schudy and Sviridenko (2011), Lemma 7.4, show that all log-concave random variables Z are central moment bounded with real parameter $L > 0$, that is, satisfy for any integer $i \geq 1$,

$$E|Z - E(Z)|^i \leq i L E|Z - E(Z)|^{i-1}.$$

Moreover, again by Schudy and Sviridenko (2011), Lemma 7.5, we have

$$L = 1 + \max(E(|Z - E(Z)| \mid Z \geq E(Z)), E(|Z - E(Z)| \mid Z < E(Z))),$$

which for the $\text{Pois}(\lambda)$ distribution gives $L = O(\lambda^{1/2})$. But

$$\begin{aligned} E|Z - E(Z)|^i &\leq i L E|Z - E(Z)|^{i-1} \\ &\leq i! L^{i-2} E(Z - E(Z))^2, \end{aligned}$$

which completes the proof of the lemma.

Proof of Theorem 4.1 from the main paper.

We first show that $P(\mathcal{A} \cap \mathcal{B}) \rightarrow 1$. We have

$$P(\mathcal{A}^c) \leq \sum_{j=J_0+1}^J \sum_{k=1}^{2^{J-j}} P((\bar{\lambda}_{(k-1)2^j+1}^{k2^j})^{-1/2} |d_{j,k} - \mu_{j,k}| \geq t_1). \quad (3)$$

Since by Lemma 3.5, the Poisson distribution satisfies Cramer's conditions, Λ is bounded from above and away from zero, and $2^{J_0} = O(n^\beta)$ for $\beta \in (0, 1)$, the strong asymptotic normality from the Corollary underneath the proof of Theorem 1 in Rudzki et al. (1978) can be used, which in our context implies that if $t_1 = O(\log^{1/2} n)$, then

$$P((\bar{\lambda}_{(k-1)2^j+1}^{k2^j})^{-1/2} |d_{j,k} - \mu_{j,k}| \geq t_1) \leq C\Phi(t_1), \quad (4)$$

where $\Phi(\cdot)$ is the cdf of the standard normal distribution and C is a universal constant. Using (4), Mills' ratio inequality and the fact that $t_1 = C_1 \log^{1/2} n$, we bound (3) from above by $\tilde{C} \log^{-1/2} n n^{1-\beta-C_1^2/2}$, where \tilde{C} is a constant, which proves that $P(\mathcal{A}) \rightarrow 1$. The proof that $P(\mathcal{B}) \rightarrow 1$ is identical.

We now turn to the estimator. Due to the orthonormality of the Haar transform, we have

$$n^{-1} \|\hat{\Lambda} - \Lambda\|^2 = n^{-1} \sum_{j=1}^J \sum_{k=1}^{2^{J-j}} (\hat{\mu}_{j,k} - \mu_{j,k})^2 + n^{-1} (s_{J,1} - \tilde{\lambda})^2, \quad (5)$$

where $\tilde{\lambda} = n^{-1/2} \sum_{k=1}^n \lambda_k$.

We first consider scales $j = 1, \dots, J_0$, for which $\hat{\mu}_{j,k} = 0$. At each scale j , there are at most N indices k for which $\mu_{j,k} \neq 0$. From the definition of $d_{j,k}$, for those $\mu_{j,k}$, we have $\mu_{j,k} \leq 2^{j/2-1} \Lambda'$, which gives

$$\sum_{j=1}^{J_0} \sum_{k=1}^{2^{J-j}} (\hat{\mu}_{j,k} - \mu_{j,k})^2 \leq N(\Lambda')^2 \sum_{j=1}^{J_0} 2^{j-2} = N(\Lambda')^2 (2^{J_0-1} - \frac{1}{2}). \quad (6)$$

We now consider the remaining scales $j = J_0 + 1, \dots, J$ and first take an arbitrary index (j, k) for which λ_i is not constant for $i = (k-1)2^j + 1, \dots, k2^j$. For such a (j, k) , we have (using Lemma 3.2 in the second inequality)

$$\begin{aligned} (\hat{\mu}_{j,k} - \mu_{j,k})^2 &= (d_{j,k} \mathbb{I}(|g_{j,k}| > t) - \mu_{j,k})^2 \\ &\leq 2d_{j,k}^2 \mathbb{I}(|g_{j,k}| \leq t) + 2(d_{j,k} - \mu_{j,k})^2 \\ &\leq 2d_{j,k}^2 \mathbb{I}(|d_{j,k}| \leq t(\bar{X}_{(k-1)2^j+1}^{k2^j})^{1/2}) + 2(d_{j,k} - \mu_{j,k})^2 \\ &\leq 2t^2 \bar{X}_{(k-1)2^j+1}^{k2^j} + 2(d_{j,k} - \mu_{j,k})^2 \\ &\leq 2t^2 (\bar{\lambda}_{(k-1)2^j+1}^{k2^j} + t_2 2^{-j/2} (\bar{\lambda}_{(k-1)2^j+1}^{k2^j})^{1/2}) + 2t_1^2 \bar{\lambda}_{(k-1)2^j+1}^{k2^j}. \end{aligned}$$

Summing the bound over the at most N indices k within each scale for which λ_i is not constant for $i = (k-1)2^j + 1, \dots, k2^j$, as well as over scales $j = J_0 + 1, \dots, J$, and noting that $\bar{\lambda}_{(k-1)2^j+1}^{k2^j} \leq \bar{\Lambda}$, gives the upper bound of

$$2N\bar{\Lambda}^{1/2} \left\{ (J - J_0)(t^2 + t_1^2) \bar{\Lambda}^{1/2} + t^2 t_2 (1 + 2^{-1/2}) 2^{-\frac{J_0+1}{2}} \right\}. \quad (7)$$

We finally consider again the scales $j = J_0 + 1, \dots, J$ and those indices (j, k) for which λ_i is constant for $i = (k-1)2^j + 1, \dots, k2^j$, which implies $\mu_{j,k} = 0$. For each such (j, k) , we have

$$(\hat{\mu}_{j,k})^2 = d_{j,k}^2 \mathbb{I}(|g_{j,k}| > t).$$

Consider the following sequence of inequalities, with the first one being implied by Lemma

3.4, and the second using the fact that $\bar{\lambda}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}} = \bar{\lambda}_{(k-1)2^j+2^{j-1}+1}^{k2^j} = \bar{\lambda}_{(k-1)2^j+1}^{k2^j}$.

$$\begin{aligned}
|g_{j,k}| > t &\Rightarrow |d_{j,k}|2^{-1/2} \left| \frac{1}{\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}}} + \frac{1}{\bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j}} \right|^{1/2} > t \\
&\Rightarrow \frac{|d_{j,k}|}{(\bar{\lambda}_{(k-1)2^j+1}^{k2^j} - \delta)^{1/2}} > t \quad \vee \quad |\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}} - \bar{\lambda}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}}| \geq \delta \\
&\quad \vee \quad |\bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j} - \bar{\lambda}_{(k-1)2^j+2^{j-1}+1}^{k2^j}| \geq \delta \\
&\Leftrightarrow \frac{|d_{j,k}|}{(\bar{\lambda}_{(k-1)2^j+1}^{k2^j})^{1/2}} > t \left(1 - \frac{\delta}{\bar{\lambda}_{(k-1)2^j+1}^{k2^j}} \right)^{1/2} \\
&\quad \vee \quad 2^{j/2} (\bar{\lambda}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}})^{-1/2} |\bar{X}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}} - \bar{\lambda}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}}| \geq \delta 2^{j/2} (\bar{\lambda}_{(k-1)2^j+1}^{(k-1)2^j+2^{j-1}})^{-1/2} \\
&\quad \vee \quad 2^{j/2} (\bar{\lambda}_{(k-1)2^j+2^{j-1}+1}^{k2^j})^{-1/2} |\bar{X}_{(k-1)2^j+2^{j-1}+1}^{k2^j} - \bar{\lambda}_{(k-1)2^j+2^{j-1}+1}^{k2^j}| \geq \\
&\quad \delta 2^{j/2} (\bar{\lambda}_{(k-1)2^j+2^{j-1}+1}^{k2^j})^{-1/2}. \tag{8}
\end{aligned}$$

Let us set $\delta = t_2 2^{-j/2} (\bar{\lambda}_{(k-1)2^j+1}^{k2^j})^{1/2}$, then if

$$t_1 \leq t(1 - t_2 2^{-j/2} (\bar{\lambda}_{(k-1)2^j+1}^{k2^j})^{-1/2})^{1/2}, \tag{9}$$

then the right-hand side of the implication (8) is negated on $\mathcal{A} \cap \mathcal{B}$, which implies that so is the left-hand side, and therefore $\hat{\mu}_{j,k} = 0$. Note (9) is satisfied if (6) from the main paper holds.

Putting together (6) and (7) and noting that $n^{-1}(s_{J,1} - \tilde{\lambda})^2 \leq n^{-1}t_1^2 \bar{\lambda}_1^n$ on \mathcal{A} , we bound (5) by

$$\frac{1}{2} n^{-1} N(\Lambda')^2 (n^\beta - 1) + 2n^{-1} N \bar{\Lambda}^{1/2} \left\{ (J - J_0)(t^2 + t_1^2) \bar{\Lambda}^{1/2} + t^2 t_2 (2 + 2^{1/2}) n^{-\beta/2} \right\} + n^{-1} t_1^2 \bar{\lambda}_1^n$$

on condition that (6) from the main paper holds, which completes the proof.

References

- M. Fisz. The limiting distribution of a function of two independent random variables and its statistical application. *Colloquium Mathematicum*, 3:138–146, 1955.
- P. Fryzlewicz. Data-driven wavelet-Fisz methodology for nonparametric function estimation. *Electronic Journal of Statistics*, 2:863–896, 2008.
- P. Fryzlewicz and G. Nason. A Haar-Fisz algorithm for Poisson intensity estimation. *Journal of Computational and Graphical Statistics*, 13:621–638, 2004.
- P. Fryzlewicz, T. Sapatinas, and S. Subba Rao. A Haar-Fisz technique for locally stationary volatility estimation. *Biometrika*, 93:687–704, 2006.
- P. Fryzlewicz, G. Nason, and R. von Sachs. A wavelet-Fisz approach to spectrum estimation. *Journal of Time Series Analysis*, 29:868–880, 2008.

- R. Rudzkis, L. Saulis, and V. Statulevicius. A general lemma on probabilities of large deviations. *Lithuanian Math. J.*, 18:226–238, 1978.
- W. Schudy and M. Sviridenko. Bernstein-like concentration and moment inequalities for polynomials of independent random variables: Multilinear case. *Preprint*, 2011.