

EXTENSION OF WALD'S FIRST LEMMA TO MARKOV PROCESSES

GEORGE V. MOUSTAKIDES,* *University of Patras*

Abstract

Let $\xi_0, \xi_1, \xi_2, \dots$ be a homogeneous Markov process and let S_n denote the partial sum $S_n = \theta(\xi_1) + \dots + \theta(\xi_n)$, where $\theta(\xi)$ is a scalar nonlinearity. If N is a stopping time with $\mathbb{E}N < \infty$ and the Markov process $\{\xi_n\}_{n=0}^{\infty}$ satisfies certain ergodicity properties, we then show that $\mathbb{E}S_N = [\lim_{n \rightarrow \infty} \mathbb{E}\theta(\xi_n)]\mathbb{E}N + \mathbb{E}\omega(\xi_0) - \mathbb{E}\omega(\xi_N)$. The function $\omega(\xi)$ is a well defined scalar nonlinearity directly related to $\theta(\xi)$ through a Poisson integral equation, with the characteristic that $\omega(\xi)$ becomes zero in the i.i.d. case. Consequently our result constitutes an extension to Wald's first lemma for the case of Markov processes. We also show that, when $\mathbb{E}N \rightarrow \infty$, the correction term is negligible as compared to $\mathbb{E}N$ in the sense that $\mathbb{E}\omega(\xi_0) - \mathbb{E}\omega(\xi_N) = o(\mathbb{E}N)$.

Keywords: Wald's lemma; stopping times

AMS 1991 Subject Classification: Primary 60G40

1. Introduction

Wald's lemmas constitute a very powerful tool in sequential analysis for evaluating the performance of sequential schemes. In particular they can be used in schemes where the test statistic is a random sum of i.i.d. random variables. This includes both popular tests of sequential analysis, namely the sequential probability ratio test (SPRT) used for the sequential hypotheses problem, and the CUSUM test used for the disruption problem. Both problems are known to have a large number of applications in the areas of signal, image and speech processing, communications, systems monitoring, economics and so on. Since, in such applications, the data encountered are rarely i.i.d., the extension of Wald's lemmas to more general data types is of primary interest. In this work we focus on one of Wald's lemmas and more specifically on the one known as *Wald's first lemma*.

Let X_1, X_2, \dots be a sequence of random variables with partial sum $S_n = X_1 + \dots + X_n$. Let N denote any stopping time with respect to the filtration generated by $\{X_n\}$. Wald's first lemma [18], extended by Blackwell [1], states that if $\{X_n\}$ is an i.i.d. sequence with $\mathbb{E}|X_1| < \infty$ and the stopping time N has $\mathbb{E}N < \infty$, then

$$\mathbb{E}S_N = \mathbb{E}X_1 \mathbb{E}N. \quad (1)$$

Generalizations of this lemma consider primarily the case $\mathbb{E}X_1 = 0$ and conditions on the stopping time N , and the sequence $\{X_n\}$ that ensures $\mathbb{E}S_N = 0$. Burkholder and Gandy [2] show that if $\mathbb{E}|X_1|^\alpha < \infty$ and $\mathbb{E}N^{1/\alpha} < \infty$ then $\mathbb{E}S_N = 0$ for the case $\alpha = 2$, for which a generalization to any $1 < \alpha \leq 2$ was made in [3]. Chow *et al.* [4] extended this result further, to a class of denormalized U -statistics.

Received 9 June 1997; revision received 4 February 1998.

* Postal address: Department of Computer Engineering and Informatics, University of Patras, 26 500 Patras, Greece.
Email address: moustaki@cti.gr.

Recent publications aim to obtain conditions that are necessary and sufficient for the validity of Wald's lemma. Gundy [7] considered the case $\mathbb{E}X_1^2 = 1$ and introduced necessary and sufficient conditions on N that guarantee $\lim_{n \rightarrow \infty} \mathbb{E}[|S_{N \wedge n} - S_N|] = 0$ and $\mathbb{E}S_N = 0$. Klass [11] showed that if $\{\tilde{X}_n\}$ is an independent copy of $\{X_n\}$ and independent of N such that $\mathbb{E} \max_{n \leq N} |\tilde{S}_n| < \infty$ then $\mathbb{E}S_N = 0$. La Pena [5] introduced a tail probability version of Wald's lemma for expectations of randomly stopped sums and showed that the work of Gundy and Klass was related. Roters [14] showed that whenever $\mathbb{E}S_N$ exists then (1) is valid provided that we do not simultaneously have $\mathbb{E}X_1 = 0$ and $\mathbb{E}N = \infty$, and that finiteness of $\mathbb{E}S_N$ is in fact a necessary condition for the validity of the lemma [15].

Generalizations of Wald's lemma in a different direction, namely the dependent data case, are also reported in the literature. Franken and Lisek [6] showed that if $\{X_n\}$ is a stationary sequence with distribution P then there exists a stationary sequence $\{[Y_n, X_n]\}$, where $Y_n \in \{0, 1\}$ with distribution Q such that $\mathbb{E}_Q[S_N | Y_0 = 1] = \mathbb{E}_Q[N | Y_0 = 1]\mathbb{E}_P X_1$.

Sadowsky [16] considered the case where the random variable $X_n = \theta(\xi_n)$, with $\theta(\xi)$ a scalar nonlinearity and $\{\xi_n\}$ a homogeneous Markov process. Using a generalization of Wald's other identity, involving products of random variables, he showed that

$$\mathbb{E}S_N = \mu'(0)\mathbb{E}N - \mathbb{E}r'(\xi_N, 0) + \mathbb{E}r'(\xi_0, 0), \tag{2}$$

where $\mu'(s)$ and $r'(\xi, s)$ are the derivatives of $\mu(s)$ and $r(\xi, s)$ with respect to s and $\mu(s)$ and $r(\xi, s)$ is a solution to the eigenvalue problem

$$e^{\mu(s)}r(\xi, s) = \mathbb{E}[e^{s\theta(\xi_1)}r(\xi_1, s) | \xi_0 = \xi], \tag{3}$$

where $e^{\mu(s)}$ is the maximum (in magnitude) eigenvalue. It must be noted that both quantities $\mu'(0)$ and $r'(\xi, 0)$ entering (2) can be computed only in cases where the eigenvalue problem in (3) can be explicitly solved and the maximum eigenvalue identified. It turns out that this is not easy most of the time, even for simple Markov processes (as for example finite state chains). The difficulty is mainly due to the parametrization of (3) with respect to the parameter s .

In this paper we intend to consider a similar data case as in [16], consequently we will obtain the same generalized form of Wald's lemma as in (2). However, because of the different approach we are going to follow, we will identify in a more precise way the two parameters entering in (2). Specifically for stopping times with $\mathbb{E}N < \infty$ and under suitable conditions on the Markov process $\{\xi_n\}$ and the function $\theta(\xi)$, we are going to show that

$$\mathbb{E}S_N = [\lim_{n \rightarrow \infty} \mathbb{E}\theta(\xi_n)]\mathbb{E}N + \mathbb{E}\omega(\xi_0) - \mathbb{E}\omega(\xi_N),$$

where the scalar function $\omega(\xi)$ is related to $\theta(\xi)$ through a Poisson integral equation involving the transition probability of the process. Regarding this last equation we must note that it is significantly simpler to solve when compared to (3) (actually we will present a solution in the form of a series) because it is not parametrized with respect to any parameter. Furthermore for several interesting Markov processes and nonlinearities $\theta(\xi)$ we will be able to present closed form solutions for the function $\omega(\xi)$ (one of the examples being finite state chains). One last point we must make is that when we consider Markov chains on countable state spaces and bounded $\theta(\xi)$ functions, then we can recover the proposed extension of Wald's first lemma from Dynkin's identity as we will see in Section 2.1.

Before going into any detail, we need to call upon certain definitions and results regarding Markov processes. Thus let $\{\xi_n\}$ be a discrete time Markov process that evolves on some state

space X equipped with a σ -field $\mathcal{B}(X)$. Let $\mathbb{P}^n(\xi, A)$ denote the n th step transition probability of $\{\xi_n\}$, that is

$$\mathbb{P}^n(\xi, A) = \text{Prob}(\xi_n \in A \mid \xi_0 = \xi), \quad n \in \mathbb{Z}_+, \xi \in X, A \in \mathcal{B}(X),$$

where $\mathbb{Z}_+ = \{0, 1, 2, \dots\}$. Let also $\pi(A)$ denote the invariant probability measure of the process.

We will be mainly interested in processes where $\mathbb{P}^n(\xi, A)$ converges to $\pi(A)$ at an exponential rate. The theory of φ -irreducible processes described in [12] or [13] will provide us with criteria that can guarantee the proper form of convergence required by our analysis, and also with various background results that will facilitate our proofs. Let us first proceed with the necessary definitions. A φ -irreducible process is a process for which there exists a non-trivial measure φ such that

$$\varphi(A) > 0 \Rightarrow \sum_n \mathbb{P}^n(\xi, A) > 0, \quad \xi \in X.$$

It turns out (see [12]) that if an invariant probability measure π exists then it is the perfect candidate as a φ measure, in the sense that the process is π -irreducible and the measure π is a maximal irreducibility measure. Of course this property is seldom used in practice since one of the main advantages of φ -irreducibility is, along with some additional conditions, to guarantee existence of the π measure. However, it will turn out to be very helpful in our examples.

We also need to define the concepts of *small sets* and *aperiodicity*. A set $C \in \mathcal{B}(X)$ is *small* if for $n > 0$ there exists a non-trivial measure ν_n on $\mathcal{B}(X)$ such that

$$\mathbb{P}^n(\xi, A) \geq \nu_n(A), \quad \xi \in C, A \in \mathcal{B}(X). \quad (4)$$

It is known (see [12, Chapter 5]), that for φ -irreducible processes any set $A \in \mathcal{B}(X)$ with $\varphi(A) > 0$ contains a small set. A process is called *aperiodic* if for a small set C with $\varphi(C) > 0$ the g.c.d. of the n satisfying (4) is 1. (If this is true for a single small set it is also true for all small sets.) We must note that in many interesting cases it is relatively easy to verify these notions. Several examples are presented in [12].

Let $\mathbb{E}_\xi\{\cdot\}$ denote conditional expectation given that state $\xi_0 = \xi$, furthermore if $g(\xi)$ is a scalar function, let $\mathbb{P}^n g, \pi g$ denote $\mathbb{P}^n g = \int g(\tau) \mathbb{P}^n(\xi, d\tau)$ and $\pi g = \int g(\tau) \pi(d\tau)$. We then have the following theorem from Chapters 15 and 16 of [12] that is concerned with the exponential convergence of \mathbb{P}^n toward π .

Theorem 1.1. *Let $\{\xi_n\}$ be φ -irreducible and aperiodic, then the following two conditions are equivalent.*

(i) *There exists a function $V(\xi) \geq 1$, finite at least for one ξ , and a small set C such that for some constants $0 \leq \lambda < 1$ and $b < \infty$, the following condition (known as drift condition) is satisfied*

$$\mathbb{P} V \leq \lambda V + b \mathbb{1}_C. \quad (5)$$

(ii) *The process is geometrically ergodic in the sense that there exists a probability measure π and a function $V(\xi) \geq 1$, finite π -a.s. ($V(\xi)$ can be the function defined in (i)), and constants $0 \leq \rho < 1$ and $R < \infty$, such that for all $n \in \mathbb{Z}_+$ and for all ξ for which $V(\xi)$ is finite we have*

$$\sup_{|g| \leq V} |\mathbb{P}^n g - \pi g| \leq \rho^n R V(\xi). \quad (6)$$

Although for finite state chains (and processes with compact state space X) the exponential convergence toward the invariant probability measure is uniform, this is not the case for general processes where there is dependence on the initial state. It is this dependence that Theorem 1.1 reveals and controls efficiently for φ -irreducible aperiodic processes.

2. Generalization of Wald's first lemma

Let us now assume that the process $\{\xi_n\}$ is φ -irreducible and aperiodic satisfying either of the two conditions of Theorem 1.1 for some given function $V(\xi) \geq 1$. Furthermore let us assume that any stopping time we use is adapted to the filtration $\{\mathcal{F}_n\}$ generated by the process $\{\xi_n\}$.

To this end consider the set \mathcal{L}_V^∞ of all functions $g(\xi)$ that satisfy $\sup_\xi |g(\xi)|/V(\xi) < \infty$. The space \mathcal{L}_V^∞ equipped with the norm $\|g\|_V = \sup_\xi |g(\xi)|/V(\xi)$ is a Banach space. In a sense \mathcal{L}_V^∞ contains all the functions for which we can *a priori* guarantee existence of their expectation under the invariant measure and also geometric ergodicity in the form of Theorem 1.1. With the help of Theorem 1.1, we can now show the following lemma.

Lemma 2.1. *Let $g(\xi) \in \mathcal{L}_V^\infty$ then*

$$\|\mathbb{P}^n g - \pi g\|_V \leq \rho^n \|g\|_V R,$$

where ρ and R are the constants defined in Theorem 1.1.

Proof. The statement of this lemma is equivalent to (6), since we can observe that for any function $g(\xi) \in \mathcal{L}_V^\infty$ we have from the definition of the norm that $|g(\xi)|/\|g\|_V \leq V(\xi)$.

Let us now present an important property for functions in \mathcal{L}_V^∞ , involving randomly stopped sums, that will be used in the following.

Lemma 2.2. *Let $g(\xi) \in \mathcal{L}_V^\infty$ then for any stopping time N we have that*

$$\mathbb{E} \left[\sum_{n=0}^{N-1} |g(\xi_n)| \right] \leq \frac{\|g\|_V}{1-\lambda} [\mathbb{E} V(\xi_0) + b\mathbb{E} N],$$

where λ and b are the constants defined in the drift condition (5).

Proof. The proof is a special case of Proposition 11.3.2 of [12] with $Z_k = \|g\|_V V(\xi_k)/(1-\lambda)$, $f_k(\xi) = |g(\xi)|$ and $s_k(\xi) = \|g\|_V b/(1-\lambda)$.

We call upon our last lemma before the presentation of our main theorem. This lemma introduces the function $\omega(\xi)$ that enters in the correction term of the generalized form of Wald's lemma.

Lemma 2.3. *Let $\theta(\xi) \in \mathcal{L}_V^\infty$, and consider the following Poisson integral equation with respect to the unknown function $\omega(\xi)$, i.e.*

$$\mathbb{P} \omega = \omega - (\mathbb{P} \theta - \pi \theta), \tag{7}$$

with the constraint

$$\pi \omega = 0, \tag{8}$$

then the constrained Poisson equation has a unique solution in \mathcal{L}_V^∞ given by the series

$$\omega = \sum_{n=1}^{\infty} \mathbb{P}^n \theta - \pi \theta, \quad (9)$$

which is convergent in norm in \mathcal{L}_V^∞ .

Proof. Notice that if $\omega(\xi)$ is a solution to (7) so is $c + \omega(\xi)$ for any constant c . Setting the constraint $\pi \omega = 0$ results in a unique solution as we will soon see. That the series defined in (9) converges is an immediate consequence of Lemma 2.1. That $\omega(\xi)$ is a solution to the Poisson equation (7) satisfying the constraint in (8) can be easily verified using the Markov property of the process. Finally, to show that $\omega(\xi)$ is unique assume that there exist two functions $\omega_1(\xi)$ and $\omega_2(\xi)$, both in \mathcal{L}_V^∞ , satisfying (7) and (8). Define $\delta(\xi) = \omega_1(\xi) - \omega_2(\xi)$, then $\delta(\xi) \in \mathcal{L}_V^\infty$, $\pi \delta = 0$ and $\mathbb{P} \delta = \delta$. Applying the last relation repeatedly, we conclude that for any $n \in \mathbb{Z}_+$ we have $\mathbb{P}^n \delta = \delta$. This last equality combined with Lemma 2.1 yields $\|\delta\|_V \leq \rho^n \|\delta\|_V R$. Since $0 \leq \rho < 1$ we have $\|\delta\|_V = 0$ and this concludes the proof.

We are now in a position to prove our main result.

Theorem 2.1. *Let $\mathbb{E} V(\xi_0) < \infty$ and $\theta(\xi) \in \mathcal{L}_V^\infty$. Define $S_n = \theta(\xi_1) + \dots + \theta(\xi_n)$, then for any stopping time N with $\mathbb{E} N < \infty$ we have*

$$\mathbb{E} S_N = [\lim_{n \rightarrow \infty} \mathbb{E} \theta(\xi_n)] \mathbb{E} N + \mathbb{E} \omega(\xi_0) - \mathbb{E} \omega(\xi_N), \quad (10)$$

with the function $\omega(\xi)$ defined in Lemma 2.3.

Proof. As in the proof for the i.i.d. case (see [17]) we need to introduce a proper martingale. Thus let us consider

$$U_n = S_n - (\pi \theta)n + \omega(\xi_n);$$

using (7) it is straightforward to show that U_n is indeed a martingale. To complete our proof we also need to show that

$$\lim_{n \rightarrow \infty} \mathbb{E}[|U_n| \mathbb{1}_{\{N > n\}}] = 0. \quad (11)$$

We have

$$|U_n| \leq \sum_{k=1}^n |\theta(\xi_k)| + |\pi \theta|n + |\omega(\xi_n)|,$$

consequently on the set $\{N > n\}$ we can write

$$|U_n| \leq \sum_{k=0}^{N-1} |\theta(\xi_k)| + |\pi \theta|N + \sum_{k=0}^{N-1} |\omega(\xi_k)|. \quad (12)$$

Since both functions $\theta(\xi)$ and $\omega(\xi)$ belong to \mathcal{L}_V^∞ , with the help of Lemma 2.2 and the fact that we assumed that $\mathbb{E} V(\xi_0) < \infty$, we conclude that all three terms on the right-hand side of (12) have finite expectation. Thus (11) is satisfied.

2.1. Ergodic chains on countable spaces and Dynkin's identity

It is worth noting that the expression $\mathbb{E}\omega(\xi_0) - \mathbb{E}\omega(\xi_N)$ also appears in another popular identity that has been in the literature for several years—we refer to Dynkin's well-known identity. As we will see, it is possible to recover our proposed extension through this identity as well. We must point out, however, that, although not explicitly stated, Dynkin's identity was indirectly used in our previous proof since it is part of the proof of [12, Proposition 11.3.2], which we employed in showing Lemma 2.2.

Although the extension of Wald's first lemma was introduced for Markov processes $\{\xi_n\}$ evolving on general state spaces, the processes in question, due to Theorem 1.1, were characterized by geometric ergodicity. This of course suggests that our analysis excludes Markov processes that are ergodic but not geometrically ergodic. As our referee pointed out, for the special case of countable state spaces and bounded $\theta(\xi)$ functions it is possible to show the extension of Wald's lemma, for general ergodic Markov chains, through Dynkin's identity. Let us present this fact in the following theorem

Theorem 2.2. *Let the Markov chain $\{\xi_n\}$ be aperiodic, ergodic and evolving on a countable state space X . For any bounded function $\theta(\xi)$ define $\omega(\xi)$ to satisfy (9), then the corresponding series is absolutely convergent. If $\omega(\xi)$ is bounded then relation (10) is also valid.*

Proof. It is possible to modify the proof of Theorem 2.1 by using alternative drift conditions guaranteeing more general forms of ergodicity (see [12]). In fact we could even go without the requirement of boundedness for the $\omega(\xi)$ function. We prefer, however, to proceed along a different line in order to emphasize the association of the proposed extension with Dynkin's identity.

If $\omega(\xi)$ satisfies (9) then by the standard coupling proof of the ergodic theorem for Markov chains [8] we have that the series converges absolutely, moreover we can easily show that (7) is also true.

We must note that although $\theta(\xi)$ is (by assumption) bounded this is not necessarily the case for $\omega(\xi)$, which can be unbounded even for countable state spaces. If $\omega(\xi)$ turns out to be bounded (as in the case of finite state spaces) then we can apply Dynkin's identity and we obtain

$$\mathbb{E}_\xi \omega(\xi_N) - \omega(\xi) = \mathbb{E}_\xi \left[\sum_{n=1}^N \mathbb{P} \omega(\xi_{n-1}) - \omega(\xi_{n-1}) \right]. \tag{13}$$

Now using Doob's optional sampling theorem we can write

$$\mathbb{E}_\xi \left[\sum_{n=1}^N \theta(\xi_n) - \mathbb{P} \theta(\xi_{n-1}) \right] = 0. \tag{14}$$

Substituting (7) and (14) into (13) one recovers the proposed extension.

3. Study of the correction term

The generalized form of Wald's lemma reduces to the usual form (1) when $\omega(\xi) = 0$. From Lemma 2.3 we can see that this can happen if and only if $\mathbb{P} \theta = \pi \theta$ or equivalently

$$\mathbb{P} [\theta - \pi \theta] = 0. \tag{15}$$

Clearly this condition is always true when $\theta(\xi)$ is a constant but also for any function $\theta(\xi)$ with finite expectation if $\{\xi_n\}$ is i.i.d. Equation (15) may also hold, except in these two cases, for special combinations of Markov processes and nonlinearities. Notice that \mathbb{P} can be regarded as a linear operator from \mathcal{L}_V^∞ to \mathcal{L}_V^∞ . Condition (15) is equivalent to saying that this operator has a non-trivial null space and that $\theta - \pi\theta$ belongs to this null space. Therefore, if $\theta(\xi) = c + g(\xi)$ with $g(\xi)$ in the null space of \mathbb{P} and c is a constant, (15) is valid.

Let us now examine the correction term $\mathbb{E}\omega(\xi_0) - \mathbb{E}\omega(\xi_N)$ and specifically its rate of growth as $\mathbb{E}N$ tends to infinity for the case where $\omega(\xi) \neq 0$. If $\omega(\xi)$ is not π -a.s. bounded then this term can grow without bound as $\mathbb{E}N \rightarrow \infty$ (consider for instance N as the time of first entry of the process in the set $\{\omega(\xi) > M\}$ or the set $\{\omega(\xi) < -M\}$ for arbitrarily large positive M). Even for cases where the correction term can become arbitrarily large, it turns out that its size relative to $\mathbb{E}N$ becomes negligible as $\mathbb{E}N \rightarrow \infty$. This suggests an asymptotic validity of Wald's original lemma in a sense described by the following theorem.

Theorem 3.1. *If $V(\xi) \geq 1$ is a drift function satisfying either of the two conditions of Theorem 1.1 and such that $\mathbb{E}V(\xi_0) < \infty$, then for any function $\theta(\xi) \in \mathcal{L}_V^\infty$ and stopping times with $\mathbb{E}N < \infty$, we have*

$$\lim_{\mathbb{E}N \rightarrow \infty} \frac{\mathbb{E}S_N}{\mathbb{E}N} = \pi\theta. \quad (16)$$

Proof. In Lemma 2.1 we have seen that for $\theta(\xi) \in \mathcal{L}_V^\infty$ we also have $\omega(\xi) \in \mathcal{L}_V^\infty$. This means that $|\mathbb{E}\omega(\xi_0) - \mathbb{E}\omega(\xi_N)| \leq \|\omega\|_V[\mathbb{E}V(\xi_0) + \mathbb{E}V(\xi_N)]$. Since $\|\omega\|_V$ is independent of the stopping time N , in order to prove (16), it suffices to show

$$\lim_{\mathbb{E}N \rightarrow \infty} \frac{\mathbb{E}V(\xi_N)}{\mathbb{E}N} = 0. \quad (17)$$

To show (17), let m be an integer greater than zero and let M be a positive constant. Let C_M denote the set $C_M = \{V(\xi) > M\}$ then, using monotone convergence, we can write

$$\begin{aligned} \mathbb{E}V(\xi_N) &= \mathbb{E}[V(\xi_N) \mathbb{1}_{C_M^c}(\xi_N)] + \mathbb{E}[V(\xi_N) \mathbb{1}_{C_M}(\xi_N)] \\ &\leq M + \sum_{n=0}^{\infty} \mathbb{E}[V(\xi_n) \mathbb{1}_{C_M}(\xi_n) \mathbb{1}(N = n)] \\ &\leq M + \sum_{n=0}^{m-1} \mathbb{E}[V(\xi_n) \mathbb{1}_{C_M}(\xi_n)] + \sum_{n=m}^{\infty} \mathbb{E}[V(\xi_n) \mathbb{1}_{C_M}(\xi_n) \mathbb{1}(N > n - m)] \\ &= M + \sum_{n=0}^{m-1} \mathbb{E}[V(\xi_n) \mathbb{1}_{C_M}(\xi_n)] + \sum_{n=m}^{\infty} \mathbb{E}[\mathbb{E}[V(\xi_n) \mathbb{1}_{C_M}(\xi_n) \mid \xi_{n-m}] \mathbb{1}(N > n - m)], \end{aligned} \quad (18)$$

where the last equality is true because the event $\{N > n - m\}$ is \mathcal{F}_{n-m} measurable. Notice that $V(\xi) \mathbb{1}_{C_M}(\xi) \in \mathcal{L}_V^\infty$ with $\|V \mathbb{1}_{C_M}\|_V \leq 1$. Applying Lemma 2.1 for $g(\xi) = V(\xi) \mathbb{1}_{C_M}(\xi)$ and using the Markov property of the process, we have for $n > k$ that

$$\mathbb{E}[V(\xi_n) \mathbb{1}_{C_M}(\xi_n) \mid \xi_k] \leq \pi[V \mathbb{1}_{C_M}] + \rho^{n-k} R V(\xi_k).$$

This inequality, if used in (18), yields

$$\mathbb{E}[V(\xi_N)] \leq M + m(\pi[V \mathbb{1}_{C_M}] + R' \mathbb{E}V(\xi_0)) + \pi[V \mathbb{1}_{C_M}] \mathbb{E}N + \rho^m R \mathbb{E} \left[\sum_{n=0}^{N-1} V(\xi_n) \right],$$

where $R' = \max\{R, 1\}$. Using Lemma 2.2 to bound the last sum, then dividing by $\mathbb{E}N$ and taking the limit as $\mathbb{E}N \rightarrow \infty$, we conclude that

$$\limsup_{\mathbb{E}N \rightarrow \infty} \frac{\mathbb{E}V(\xi_N)}{\mathbb{E}N} \leq \pi[V\mathbb{1}_{C_M}] + \rho^m \frac{Rb}{1-\lambda}.$$

Since $\pi V < \infty$ and $0 \leq \rho < 1$, the right-hand side of the previous inequality can become arbitrarily small by tending M and m to infinity. This concludes the proof.

4. Examples

In this section we are going to present several examples and identify explicitly all the terms entering in the generalized form of Wald's lemma. For every example, in order to apply our results, we also need to show the existence of a function $V(\xi)$ that satisfies the drift condition (5). At this point we must stress that the drift condition does not have a unique solution, and the different solutions are not necessarily equivalent, in the sense that they produce the same space of functions \mathcal{L}_V^∞ . Consequently if we are interested in applying our results to some given function $\theta(\xi)$, it is clear that it suffices to find one solution to the drift condition such that $\theta(\xi) \in \mathcal{L}_V^\infty$.

4.1. Finite state Markov chains

Let the chain $\{\xi_n\}$ have K states and let P be the corresponding transition matrix with the i th row of P denoting the transition probabilities of state i . We know that any transition matrix P has a unit eigenvalue whereas all other eigenvalues have magnitude that cannot exceed unity. Let us assume that the unit eigenvalue is simple and that all other eigenvalues have magnitude strictly less than unity. Then the invariant probability vector π exists and it is the left eigenvector to the unit eigenvalue, $J = [1 \cdot \cdot \cdot 1]^t$ being the corresponding right eigenvector. We can now verify (by induction) that $P^n - J\pi^t = (P - J\pi^t)^n$. The matrix $P - J\pi^t$ has the same eigenvalues as P except the unit eigenvalue that has become zero. This means that all eigenvalues of $P - J\pi^t$ have a magnitude that is strictly smaller than unity. Since for any square matrix A we have $\lim_{n \rightarrow \infty} \sqrt[n]{\|A^n\|} = \max_i |a_i|$, where a_i represents the eigenvalues of A and $\|\cdot\|$ represents any matrix norm (see for example [10, pp. 36–38]), we conclude that there exist R and $0 \leq \rho < 1$ such that $\|P^n - J\pi^t\| \leq \rho^n R$. This in turn suggests [12, Theorem 16.0.2] that the chain is aperiodic and the drift condition is satisfied by an everywhere bounded function $V(\xi)$. Finally, as we stated in Section 1, the chain is also π -irreducible.

For this example, any function $\theta(\xi)$ can in fact be regarded as a vector θ of length K and the space \mathcal{L}_V^∞ coincides with the space \mathbb{R}^K . To find the two quantities entering in the generalized form of the lemma, notice first that the expectation of θ under the invariant measure is simply $\pi^t \theta$. To find the vector ω we have the following linear system of equations with respect to ω (which is the analogue of the Poisson equation (7) and the constraint (8)),

$$\begin{aligned} (P - I)\omega &= -(P - J\pi^t)\theta, \\ \pi^t \omega &= 0. \end{aligned}$$

It should also be noted that a sufficient condition for the assumption regarding the eigenvalues of the matrix P to hold is that there exists $m \in \mathbb{Z}_+$ such that all the entries of the matrix P^m are positive (see [9, Chapter IV]).

To examine whether it is possible to have validity of Wald's lemma in its original form, notice that relation (15) takes the form $P[\theta - (\pi^t \theta)J] = 0$. This can be true if the matrix P has a zero eigenvalue.

4.2. Finite dependence processes

As a second example consider the sequence $\{\zeta_n\}_{n=-m+1}^{\infty}$ of i.i.d. random variables with probability measure μ . We are interested in nonlinear transformations of ξ_n where ξ_n is the m -tuple $\xi_n = (\zeta_n, \dots, \zeta_{n-m+1})$. Let us illustrate this class by limiting ourselves to the special case $m = 2$ (an extension to the general m -dependent case is straightforward). Since $\xi_n = (\zeta_n, \zeta_{n-1})$, the limiting measure π exists and $\pi = \mu \times \mu$. Furthermore for any integer $n \geq 2$ we have that $\mathbb{P}^n = \pi$. As stated in Section 1, if the invariant probability measure π exists then the process is π -irreducible. To test for aperiodicity, notice that if we select $\nu_n = \pi$ then (4) is valid with equality for any $n \geq 2$. Consequently the process is aperiodic.

Let us now identify a function $V(\xi)$. For this example it seems more appropriate to use condition (ii) of Theorem 1.1. Notice that we need to satisfy (6) only for $n = 0, 1$, since for $n \geq 2$ it is trivially satisfied by any function $V(\xi) \geq 1$. If $G(\xi) \geq 0$ is a function such that $\pi G < \infty$ let $V = 1 + G + \mathbb{P}G$. We can then show that (6) is satisfied for $n = 0, 1$, with $\rho = 0.5$, $R = 2(1 + 3\pi G)$, and observing that for $|g(\xi)| \leq V(\xi)$ we have $|\mathbb{P}^n g - \pi g| \leq \mathbb{P}^n V + \pi V$. Let us fix such a function $V(\xi)$ and consider $\theta(\xi) \in \mathcal{L}_V^{\infty}$. To compute the corresponding function $\omega(\xi)$ we observe that for the series (9) only the first term is non-zero, thus $\omega(\zeta_0, \zeta_{-1}) = \mathbb{E}[\theta(\zeta_1, \zeta_0) \mid \zeta_0, \zeta_{-1}] - \pi\theta$. From this last relation we also note that $\omega(\zeta_0, \zeta_{-1})$ depends only on ζ_0 and not on ζ_{-1} . Summarizing our result,

$$\mathbb{E} \left[\sum_{n=1}^N \theta(\zeta_n, \zeta_{n-1}) \right] = \mathbb{E}\theta(\zeta_1, \zeta_0)\mathbb{E}N + \mathbb{E}\omega(\zeta_0) - \mathbb{E}\omega(\zeta_N),$$

where $\omega(\zeta) = \mathbb{E}[\theta(\zeta_1, \zeta_0) \mid \zeta_0 = \zeta] - \mathbb{E}\theta(\zeta_1, \zeta_0)$.

Here, it is easy to find functions for which the correction term is zero and which thus satisfy Wald's lemma in its original form. Let $g(\zeta_1, \zeta_0) \in \mathcal{L}_V^{\infty}$ then define $\theta(\zeta_1, \zeta_0) = g(\zeta_1, \zeta_0) - \mathbb{E}[g(\zeta_1, \zeta_0) \mid \zeta_0] + c$, with c any constant. If the function $\theta(\zeta_1, \zeta_0)$ has this form then it satisfies relation (15), and consequently its corresponding function $\omega(\zeta)$ is equal to zero.

4.3. AR processes

We limit this presentation to the scalar case and consider processes of the form

$$\xi_n = a\xi_{n-1} + w_n,$$

where $\{w_n\}$ is an i.i.d. sequence and $|a| < 1$. Here, although it is not necessary, we assume for simplicity that the sequence $\{w_n\}$ has an everywhere positive density. Under this last assumption the process becomes φ -irreducible and aperiodic with φ being the Lebesgue measure and with compact sets being small sets (see [12, Chapter 16]).

We can now show that $V(\xi) = 1 + |\xi|^p$ satisfies the drift condition provided that $\mathbb{E}|w_1|^p < \infty$ for some $p \geq 1$. Indeed notice that by using Hölder's inequality we have the following for $\epsilon > 0$ and $p^{-1} + q^{-1} = 1$:

$$|\xi_1|^p = |a\xi + w_1|^p \leq \left(|a||\xi| + \epsilon \frac{|w_1|}{\epsilon} \right)^p \leq (|a|^q + \epsilon^q)^{p/q} \left(|\xi|^p + \frac{|w_1|^p}{\epsilon^p} \right).$$

Taking conditional expectation we then conclude that for small enough ϵ we can find $0 < \lambda' < 1$ and L' such that $\mathbb{P}V \leq \lambda'V + L'$. Selecting λ such that $\lambda' < \lambda < 1$, $L = L'/(\lambda - \lambda')$ and $C = \{V(\xi) \leq L\}$, the drift condition in (5) is then satisfied.

With function $V(\xi)$ in the form we have just introduced, we can apply the results of the previous sections to any function $\theta(\xi)$ that can be bounded by a polynomial of degree p .

However, for cases where all moments of w_1 exist, it is tempting to examine whether a function $V(\xi)$ can be constructed such that \mathcal{L}_V^∞ contains all polynomials. Such a function can indeed be proposed if for some $p \geq 1$ we can find a constant $c > 0$ such that $\mathbb{E}e^{c|w_1|^p} < \infty$ (this is true for example for any $1 \leq p \leq 2$ when w_1 is Gaussian). If such a condition holds then there exists a constant $\delta > 0$ such that $V(\xi) = e^{\delta|\xi|^p}$ satisfies the drift condition for properly selected λ , L and C . The proof is very similar to the polynomial $V(\xi)$ case presented above, thus we avoid giving any further details.

Finding closed form expressions for the quantities entering into the generalized form of Wald's lemma is not as straightforward as it was in the previous examples, although an interesting class of nonlinearities for which this is possible is the class of polynomials. Notice that if we have the moments $\mathbb{E}w_1^j$, $j = 0, 1, \dots, p$ available, then we can define polynomials $s_j(\xi)$, $j = 0, 1, \dots, p$, with $s_j(\xi)$ of degree j and the coefficient of the highest power equal to unity, satisfying

$$\mathbb{P} s_j = a^j s_j. \tag{19}$$

It must be noted that the computation of the coefficients of each polynomial involves the solution of a linear system of equations (when w_1 is zero mean Gaussian this process leads to the generation of a normalized version of the Hermite polynomials). Any polynomial $\theta(\xi)$ of degree $k \leq p$ can then be written as a linear combination of the polynomials $s_j(\xi)$, $j = 0, \dots, p$, namely $\theta(\xi) = \sum_{j=0}^k \theta_j s_j(\xi)$. Due to (19) we have $\mathbb{P}^n \theta = \sum_{j=0}^k a^{nj} \theta_j s_j(\xi)$ and thus $\pi \theta = \lim_{n \rightarrow \infty} \mathbb{P}^n \theta = \theta_0$. Applying these results in the series (9) yields the following form for $\omega(\xi)$,

$$\omega(\xi) = \sum_{j=1}^k \theta_j \frac{a^j}{1 - a^j} s_j(\xi), \tag{20}$$

and thus we have identified both quantities entering the generalized form of Wald's lemma. To see whether a $\theta(\xi)$ exists such that Wald's lemma is valid in its original form, notice that if $\theta(\xi)$ is a polynomial then the correction term cannot be zero (when $a \neq 0$) except in the trivial case of a constant $\theta(\xi)$. This is because $\omega(\xi)$ is zero if and only if $\theta_j = 0$, $j = 1, \dots, p$, which is only true when $\theta(\xi)$ is a constant.

The previous results can be extended, without much difficulty, to ARMA(p, q) processes and, more generally, to processes of the form $\xi_n = A\xi_{n-1} + Bw_n$, where ξ_n is a vector of length k , $\{w_n\}$ is an i.i.d. vector sequence with w_n of length r and with an everywhere positive density, A is a $k \times k$ matrix with all its eigenvalues inside the unit circle and B is a $k \times r$ matrix such that the pair (A, B) is controllable, i.e. the matrix $[B \ AB \ \dots \ A^{k-1}B]$ is of full rank (for more details see [12, Chapter 16]).

5. Application

Perhaps the most important application of Wald's identity can be considered to be the computation of the expectation of a randomly stopped log-likelihood ratio function. Specifically, for a Markov process $\{\xi_n\}_{n=0}^\infty$ that can be characterized by two different transition measures, $\mathbb{P}_1, \mathbb{P}_2$, we are interested in computing the expectation $\mathbb{E}[\sum_{n=1}^N \log(l(\xi_n, \xi_{n-1}))]$, where

$$l(\xi_1, \xi_0) = \frac{d\mathbb{P}_2(\xi_1, \xi_0)}{d\mathbb{P}_1(\xi_1, \xi_0)}$$

denotes the Radon–Nikodym derivative. The desired expectation is taken with respect to a probability measure induced by a transition measure \mathbb{P}_0 that is not necessarily equal to either of the two measures $\mathbb{P}_1, \mathbb{P}_2$.

It is clear that the required expectation can be evaluated using the theory we developed in previous sections. Indeed, by defining an new Markov process $\{\zeta_n\}$ with $\zeta_n = (\xi_n, \xi_{n-1})$ we can apply our results to $\{\zeta_n\}$. However, we wish to propose a slightly different approach that will reveal the special form of the quantities entering in the generalized identity in a more natural way. Let us consider the expectation $\mathbb{E}[\sum_{n=1}^N \theta(\xi_n, \xi_{n-1})]$ where $\theta(\xi_1, \xi_0)$ is some nonlinear function depending on two consecutive states of the process $\{\xi_n\}$.

Notice that since the stopping time N is $\{\mathcal{F}_n\}$ adapted we can write

$$\mathbb{E}\left[\sum_{n=1}^N \theta(\xi_n, \xi_{n-1})\right] = \mathbb{E}\left[\sum_{n=1}^N \mathbb{E}[\theta(\xi_n, \xi_{n-1}) \mid \xi_{n-1}]\right] = \mathbb{E}\left[\sum_{n=1}^N \tilde{\theta}(\xi_n)\right] + \mathbb{E}\tilde{\theta}(\xi_0) - \mathbb{E}\tilde{\theta}(\xi_N),$$

where we write $\tilde{\theta}(\xi_0) = \mathbb{E}[\theta(\xi_1, \xi_0) \mid \xi_0]$. We note that the last sum is in the form considered in the previous sections and therefore can be treated accordingly. On the other hand the two remaining terms are consistent with the form of the correction term of the generalized identity introduced in Theorem 2.1. Combining the different parts, and using the fact that $\lim_{n \rightarrow \infty} \mathbb{E}\tilde{\theta}(\xi_n) = \lim_{n \rightarrow \infty} \mathbb{E}\theta(\xi_n, \xi_{n-1})$, we conclude that

$$\mathbb{E}\left[\sum_{n=1}^N \theta(\xi_n, \xi_{n-1})\right] = \left[\lim_{n \rightarrow \infty} \mathbb{E}\theta(\xi_n, \xi_{n-1})\right] \mathbb{E}N + \mathbb{E}\omega(\xi_0) - \mathbb{E}\omega(\xi_N),$$

where $\omega(\xi) = \tilde{\omega}(\xi) + \tilde{\theta}(\xi)$ and $\tilde{\omega}(\xi)$ satisfies a Poisson integral equation as in Lemma 2.3 (see (7)), with $\theta(\xi)$ replaced by $\tilde{\theta}(\xi)$.

Let us apply the above result to Gaussian autoregressive processes. Assume that the Markov process satisfies $\xi_n = \alpha_i \xi_{n-1} + w_n$ with $\{w_n\}$ i.i.d. Gaussian $\mathcal{N}(0, \sigma_i^2)$, and $i = 0, 1, 2$, corresponding to the three different transition measures $\mathbb{P}_0, \mathbb{P}_1, \mathbb{P}_2$. We also assume that $|\alpha_i| < 1$, $i = 0, 1, 2$, for the corresponding processes to be stable. The log of the Radon–Nikodym derivative has the form $\theta(\xi_1, \xi_0) = c_1 \xi_1^2 + c_2 \xi_1 \xi_0 + c_3 \xi_0^2 + c_4$ with $c_1 = 0.5(\sigma_1^{-2} - \sigma_2^{-2})$, $c_2 = \alpha_2 \sigma_2^{-2} - \alpha_1 \sigma_1^{-2}$, $c_3 = 0.5(\alpha_1^2 \sigma_1^{-2} - \alpha_2^2 \sigma_2^{-2})$ and $c_4 = \log(\sigma_1/\sigma_2)$. Consequently we have that

$$\lim_{n \rightarrow \infty} \mathbb{E}\theta(\xi_n, \xi_{n-1}) = \frac{\sigma_0^2}{1 - \alpha_0^2} (c_1 + \alpha_0 c_2 + c_3) + c_4,$$

and $\tilde{\theta}(\xi)$ has the form

$$\tilde{\theta}(\xi) = (\alpha_0^2 c_1 + \alpha_0 c_2 + c_3) \xi^2 + (c_1 \sigma_0^2 + c_4).$$

To solve the Poisson integral equation for $\tilde{\omega}(\xi)$ we are going to use the method presented in Section 4.3. The required polynomials are $s_0(\xi) = 1$, $s_1(\xi) = \xi$ and $s_2(\xi) = \xi^2 - \sigma_0^2/(1 - \alpha_0^2)$. We can write $\tilde{\theta}(\xi) = \tilde{\theta}_2 s_2(\xi) + \tilde{\theta}_0 s_0(\xi)$, where $\tilde{\theta}_0 = (c_1 + \alpha_0 c_2 + c_3) \sigma_0^2 / (1 - \alpha_0^2) + c_4$ and $\tilde{\theta}_2 = \alpha_0^2 c_1 + \alpha_0 c_2 + c_3$ ($\tilde{\theta}_1 = 0$). To find $\tilde{\omega}(\xi)$ we use (20) and obtain

$$\tilde{\omega}(\xi) = \frac{\alpha_0^2}{1 - \alpha_0^2} (\alpha_0^2 c_1 + \alpha_0 c_2 + c_3) \left(\xi^2 - \frac{\sigma_0^2}{1 - \alpha_0^2} \right).$$

The final correction term that enters in the generalized identity is, as we stated previously, $\omega(\xi) = \tilde{\omega}(\xi) + \tilde{\theta}(\xi)$.

Acknowledgement

We would like to thank an anonymous reviewer for pointing out the relation of the proposed extension with Dynkin's identity.

References

- [1] BLACKWELL, D. (1946). On an equation by Wald. *Ann. Math. Statist.* **17**, 84–87.
- [2] BURKHOLDER, D. L. AND GUNDY, R. F. (1970). Extrapolation and interpolation of quasilinear operators on martingales. *Acta Math.* **124**, 249–304.
- [3] CHOW, Y. S., ROBBINS, H. AND SIEGMUND, D. (1971). *Great Expectations. The Theory of Optimal Stopping*. Houghton Mifflin, New York.
- [4] CHOW, Y. S., DE LE PENA, V. H. AND TEICHER, H. (1993). Wald's equation for a class of denormalized U -statistics. *Ann. Prob.* **21**, 1151–1158.
- [5] DE LA PENA, V. H. (1993). Inequalities for tails of adapted processes with an application to Wald's lemma. *J. Theoret. Prob.* **6**, 285–302.
- [6] FRANKEN, P. AND LISEK, B. (1982). On Wald's identity for dependent variables. *Z. Wahrscheinlichkeitsthe.* **60**, 143–150.
- [7] GUNDY, R. F. (1981). On a theorem of F. and M. Riesz and an equation of A. Wald. *Indiana Univ. Math. J.* **30**, 589–605.
- [8] HOEL, P. G., PORT, S. C. AND STONE, C. J. (1987). *Introduction to Stochastic Processes*. Waveland Press, Prospect Heights, IL.
- [9] ISAACSON, D. L. AND MADSEN, R. W. (1976). *Markov Chains Theory and Applications*. John Wiley, New York.
- [10] KATO, T. (1966). *Perturbation Theory for Linear Operators*. Springer, New York.
- [11] KLASS, M. J. (1988). A best possible improvement of Wald's equation. *Ann. Prob.* **16**, 840–853.
- [12] MEYN, S. P. AND TWEEDIE, R. L. (1993). *Markov Chains and Stochastic Stability*. Springer, New York.
- [13] NUMMELIN, E. (1984). *General Irreducible Markov Chains and Non-Negative Operators*. CUP, Cambridge, UK.
- [14] ROTERS, M. (1994). On the validity of Wald's equation. *J. Appl. Prob.* **31**, 949–957.
- [15] ROTERS, M. (1995). An optimal stopping problem for random walks with non-zero drift. *J. Appl. Prob.* **32**, 956–959.
- [16] SADOWSKY, J. S. (1989). A dependent data extension of Wald's identity and its application to sequential test performance computation. *IEEE Trans. Inform. Theory* **35**, 834–842.
- [17] SHIRYAYEV, A. N. (1978). *Optimal Stopping Rules*. Springer, New York.
- [18] WALD, A. (1945). Sequential tests of statistical hypotheses. *Ann. Math. Statist.* **16**, 117–186.