

Rare-event analysis of modulated Ornstein-Uhlenbeck processes

H M Jansen, M Mandjes, Koen de Turck, Sabine Wittevrongel

► **To cite this version:**

H M Jansen, M Mandjes, Koen de Turck, Sabine Wittevrongel. Rare-event analysis of modulated Ornstein-Uhlenbeck processes. Performance Evaluation, Elsevier, 2017, 112, pp.1 - 14. 10.1016/j.peva.2017.02.002 . hal-01580377

HAL Id: hal-01580377

<https://hal.archives-ouvertes.fr/hal-01580377>

Submitted on 1 Sep 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

RARE-EVENT ANALYSIS OF MODULATED ORNSTEIN-UHLENBECK PROCESSES

H. M. JANSEN^{1,3}, M. MANDJES¹, K. DE TURCK², S. WITTEVRONGEL³

ABSTRACT. This paper studies Ornstein-Uhlenbeck (OU) processes in a random environment. The OU model has found widespread use in networking, as a Gaussian approximation of the user-level dynamics that allows explicit analysis; adding modulation to it allows incorporating phenomena in which the users' activity level is affected by exogenous factors. The focus lies on rare-event analysis: under a specific scaling of the parameters involved, we establish the large deviations asymptotics of the probability that the process reaches an extreme value. The decay rate of this probability is generally only implicitly available (as the solution to a variational problem), but specializing to the case of Markov modulation we succeed in devising efficient numerical procedures.

KEYWORDS. Large deviations asymptotics ◦ random environment ◦ Markov modulation ◦ Ornstein-Uhlenbeck process ◦ Hamilton-Jacobi-Bellman equations

¹ Korteweg-de Vries Institute for Mathematics, University of Amsterdam, Science Park 904, 1098 XH Amsterdam, the Netherlands.

² CentraleSupélec, Département de Télécommunications, Laboratory of Signals and Systems (L2S), UMR8506 Plateau de Moulon, 3 rue Joliot-Curie, 91192 Gif sur Yvette, France.

³ TELIN, Ghent University, Sint-Pietersnieuwstraat 41, B-9000 Ghent, Belgium.

{h.m.jansen|m.r.h.mandjes}@uva.nl, koen.deturck@centralesupelec.fr,
Sabine.Wittevrongel@ugent.be

1. INTRODUCTION

In the performance analysis of communication systems, various canonical elements can be distinguished. At a detailed time-scale, the amount of information stored at an individual network node is modeled as a (stable) queue. The simplest among these is the classical M/M/1 queue, but substantially more general queues are amenable for analysis. Under rather general circumstances — most notably, it is required that the variance of the amount of traffic entering in a time window of given length be finite — such systems can be approximated by Brownian motion reflected at 0; see e.g. [18, Chs. V, IX]. At a somewhat coarser time-scale, the number of clients simultaneously present is typically modeled by a mean-reverting process: users enter the system at a roughly constant rate, but the departure rate is roughly proportional to the number of users present. The simplest model in this class is the M/M/ ∞ queue, which can also be approximated by its Gaussian counterpart. As it turns out, this Gaussian process is the well-known Ornstein-Uhlenbeck (OU) process; see e.g. [16, Section 6.6] and [18, Section 10.4].

In this paper we concentrate on models of the latter kind: mean-reverting models, essentially corresponding to the user-level dynamics. The focus is on developing techniques

to analyze their rare-event behavior: we present results that facilitate the evaluation of the probability that at a given point in time the user-level congestion level is unusually high. Motivated by the reasoning above, we consider the situation that the underlying dynamics are of the OU type, but we add one distinguishing complication: we allow the system to react to a random environment. This means that the parameters of the OU process depend on the state of an external, independently evolving, process — this process is typically referred to as the *background process*. From a practical standpoint, incorporating such modulation offers substantial advantages: it enables the modeling of all sorts of phenomena in which the users' activity level is affected by exogenous factors.

To state our contributions, let us now introduce our model in somewhat greater detail. We consider the process $(M(t))_{t \geq 0}$, to be thought of as the aggregate user activity level. A 'normal' (non-modulated, that is) OU process is given, for parameters $\gamma, \sigma^2 > 0$ and $\alpha \in \mathbb{R}$, through the stochastic differential equation (SDE)

$$dM(t) = (\alpha - \gamma M(t)) dt + \sigma dB(t),$$

where $(B(t))_{t \geq 0}$ is a standard Brownian motion. This OU process, which can be used to approximate the number of customers in specific classes of infinite-server models, has been studied in detail. In particular it has been proven that $M(t)$ has a Normal distribution (with parameters that are explicit functions of α, γ, σ^2 , and t); the long term mean is given by α/γ , whereas the long term variance is $\sigma^2/(2\gamma)$. In addition, the distribution of the running maximum has been derived, albeit in terms of special functions [1, 7].

The *modulated* OU process, which is the object of interest of this paper, is represented by the SDE

$$dM(t) = (\alpha(J(t)) - \gamma(J(t))M(t)) dt + \sigma(J(t)) dB(t),$$

where $(J(t))_{t \geq 0}$ is the (independently evolving) background process. We allow a great level of generality with respect to the modulating process $J(\cdot)$; the assumptions imposed amount to requiring that certain functionals of the modulating process satisfy the large deviations principle. A concrete example to keep in mind, also featuring in our earlier work [13], is that of $J(\cdot)$ corresponding to a finite-state irreducible Markov process.

The main contributions of this paper are the following. In the first place we derive a large deviations principle with which we can assess the probability that $M(t)$ exceeds some (large) threshold; we do so in an asymptotic regime in which the parameters are scaled, comparably to how this was done in prior work [12]. However, to derive the large deviations principle, we take an approach that strongly differs from the approach in [12]. This different approach is inspired by the approach in [14] and leads to an expression of the large deviations rate function that is particularly amenable to numerical evaluation. The resulting expressions of the decay rate of the probability of interest are in terms of variational problems, with an insightful decomposition into (i) the impact of the background process and (ii) that of the driving Brownian motion (conditional on the background process).

Then we specifically focus on the numerical evaluation of these variational problems for the special case that the background process corresponds to an irreducible continuous-time finite-state Markov chain. We do so by setting up various numerical schemes which are provably equivalent (as they correspond to the same Hamilton-Jacobi-Bellman equations), but that differ in terms of numerical features. In particular, we managed to reduce the number of dimensions of the variational problem that needs to be solved from 2 to 1. We then consider the special case in which the background process is relatively slow, in which we find intuitively appealing results; notably, we prove that, along the most likely path, the background process jumps at most $2d - 2$ times, with d the number of states of the background process; this is in stark contrast with earlier findings for the infinite-server queue in [4], where the background process jumps at most $d - 1$ times.

In previous work on modulated mean-reverting processes there was a strong focus on the central limit regime [2, 5, 13]; under various conditions convergence to an ordinary (i.e., non-modulated) OU process has been established. Related results on large deviations for modulated infinite-server queues can be found in e.g. [4, 14, 17], whereas diffusion-type processes are considered in e.g. [12, 15].

This paper is organized as follows. Section 2 defines our model and proves the large deviations asymptotics of the probability of our interest, which are presented in terms of a so-called *large deviations principle* [8, Section I.2]. Section 3 then derives the rate function for the special case that the background process is an irreducible continuous-time finite-state Markov chain. Next, in Section 4, we focus on the numerical evaluation of the decay rate; in addition explicit properties are given for the case that the background process evolves slowly. The last section covers a number of illustrative examples.

2. LARGE DEVIATIONS PRINCIPLE

In this section, we prove the large deviations principle (LDP) for an appropriately scaled modulated OU process. Let us first define modulated OU processes and next recall the definition of the LDP.

Let $B(\cdot)$ be a standard Brownian motion and let the background process $J(\cdot)$ be an independent càdlàg stochastic process taking values in \mathbb{R}^d . Additionally, let $\alpha: \mathbb{R}^d \rightarrow \mathbb{R}$, $\gamma: \mathbb{R}^d \rightarrow \mathbb{R}$, and $\sigma: \mathbb{R}^d \rightarrow [0, \infty)$ be continuous functions. Then we define the modulated OU process via the SDE

$$dM(t) = (\alpha(J(t)) - \gamma(J(t))M(t)) dt + \sigma(J(t)) dB(t).$$

Following the arguments of [13], it turns out that a solution to the above SDE exists and is unique. Moreover, the arguments in [13] also show that $M(t)$ has a ‘mixed Normal distribution’, i.e., it has a Normal distribution with *random* mean

$$m(t) \equiv m_{J(\cdot)}(t) := M(0) \exp\left(-\int_0^t \gamma(J(r)) dr\right) + \int_0^t \alpha(J(s)) \exp\left(-\int_s^t \gamma(J(r)) dr\right) ds$$

and *random* variance

$$v(t) \equiv v_{J(\cdot)}(t) := \int_0^t \sigma^2(J(s)) \exp\left(-2\int_s^t \gamma(J(r)) dr\right) ds.$$

Informally, the ‘mixed Normal’ property entails that, *conditional on the path of the background process* between 0 and t , the random variable $M(t)$ has a Normal distribution.

We scale $\alpha \mapsto n\alpha$ and $\sigma^2 \mapsto n\sigma^2$. To stress the dependence on the scaling parameter n , we call the resulting process $M_n(t)$; the corresponding random parameters are denoted by $m_n(t)$ and $v_n(t)$. We also allow for a scaling of the process $J(t)$: that is, $J(t) \mapsto J_n(\cdot)$, where conditions on the process $J_n(\cdot)$ will be specified below, but broadly speaking the scaling should be such that an appropriately scaled and centered version of the vector (m_n, v_n) admits an LDP.

Let us briefly recall the definition of an LDP. Given some Hausdorff topological space \mathcal{X} and a sequence of random variables $\{X_n\}_{n \in \mathbb{N}}$ taking values in \mathcal{X} , we say that X_n satisfies the LDP with speed $n^{1-2\beta}$ and rate function I if

$$\limsup_{n \rightarrow \infty} n^{2\beta-1} \log \mathbb{P}(X_n \in F) \leq - \inf_{a \in F} I(a)$$

for all closed sets $F \subset \mathcal{X}$ and

$$\liminf_{n \rightarrow \infty} n^{2\beta-1} \log \mathbb{P}(X_n \in G) \geq - \inf_{a \in G} I(a)$$

for all open sets $G \subset \mathcal{X}$. Here, I is a nonnegative lower semicontinuous function and $\beta \in [0, 1/2)$. If $\beta = 0$, then this definition coincides with the classical definition of an LDP. If $\beta \in (0, 1/2)$, then this definition gives the definition of an MDP (moderate deviations principle). If we do not specify the speed of an LDP, we mean the classical LDP. For more background on LDPs and MDPs, see e.g. [8] and [10].

The following theorem, which is the main result of this section, states that under the proviso that a scaled and centered version of $m_n(t)$ jointly with a scaled version of $v_n(t)$ satisfies the LDP, then so does an appropriately centered and scaled version of $M_n(t)$. This LDP is in terms of two functions \mathcal{L} and ψ , where $\mathcal{L}: \mathbb{R} \times [0, \infty) \rightarrow [0, \infty]$ is given by

$$\mathcal{L}(x, y) := \begin{cases} 0 & \text{if } x = 0, \\ \frac{1}{2} \frac{x^2}{y} & \text{if } x \neq 0, y > 0, \\ \infty & \text{if } x \neq 0, y = 0, \end{cases}$$

and $\psi: \mathbb{R} \times [0, \infty) \rightarrow [0, \infty]$ is some rate function. Given ψ , we define the set $\Psi := \{(m, v) : \psi(m, v) < \infty\}$.

Theorem 2.1. *Let $\rho(t)$ be a random variable and let $\beta \in [0, \frac{1}{2})$. Suppose that the random vector*

$$(1) \quad \left(n^\beta \left(\frac{m_n(t)}{n} - \rho(t) \right), \frac{v_n(t)}{n} \right)$$

satisfies the LDP with speed $n^{1-2\beta}$ in $\mathbb{R} \times [0, \infty)$ with rate function ψ . Then

$$n^\beta \left(\frac{M_n(t)}{n} - \rho(t) \right)$$

satisfies the LDP with speed $n^{1-2\beta}$, with the corresponding rate function I given by

$$(2) \quad I(x) := \inf_{(m, v) \in \Psi} [\mathcal{L}(x - m, v) + \psi(m, v)].$$

Proof. The proof consists of a lower bound and an upper bound.

We start with the lower bound. Let $G \subset \mathbb{R}$ be open and let $x \in G$. Take any $\delta > 0$ such that $B(x, \delta) \subset G$ (where $B(x, \delta)$ is an open ball of radius δ and center x) and fix $(m, v) \in \Psi$. Let $\epsilon = \delta/2$ and denote

$$\mathcal{A}_n \equiv \mathcal{A}_n(\epsilon) := \left\{ n^\beta \left(\frac{m_n(t)}{n} - \rho(t) \right) \in B(m, \epsilon), \frac{v_n(t)}{n} \in B(v, \epsilon) \right\}.$$

Because $\psi(m, v) < \infty$, there exists $N_\epsilon \in \mathbb{N}$ such that $\mathbb{P}(\mathcal{A}_n) > 0$ for all $n \geq N_\epsilon$. Now write

$$\begin{aligned} \mathbb{P} \left(n^\beta \left(\frac{M_n(t)}{n} - \rho(t) \right) \in G \right) &\geq \mathbb{P} \left(n^\beta \left(\frac{M_n(t)}{n} - \rho(t) \right) \in B(x, \delta) \right) \\ &\geq \mathbb{P} \left(\left\{ n^\beta \left(\frac{M_n(t)}{n} - \rho(t) \right) \in B(x, \delta) \right\} \cap \mathcal{A}_n \right) \\ (3) \quad &= \mathbb{P} \left(n^\beta \left(\frac{M_n(t)}{n} - \rho(t) \right) \in B(x, \delta) \mid \mathcal{A}_n \right) \mathbb{P}(\mathcal{A}_n) \end{aligned}$$

and observe that, as we took $\epsilon = \delta/2$,

$$\begin{aligned} &\mathbb{P} \left(n^\beta \left(\frac{M_n(t)}{n} - \rho(t) \right) \in B(x, \delta) \mid \mathcal{A}_n \right) \\ &= \mathbb{P} \left(n^{\beta-1} (M_n(t) - m_n(t)) + n^\beta \left(\frac{m_n(t)}{n} - \rho(t) \right) \in B(x, \delta) \mid \mathcal{A}_n \right) \\ &\geq \mathbb{P} \left(n^{\beta-1} (M_n(t) - m_n(t)) \in B(x - m, \delta/2) \mid \mathcal{A}_n \right). \end{aligned}$$

Recall that the random variable $n^{\beta-1} (M_n(t) - m_n(t))$ has a Normal distribution with mean 0 and (random) variance $n^{2\beta-1} v_n(t)/n$. As a consequence, it is easy to see that the following lower bound applies:

$$\begin{aligned} &\liminf_{n \rightarrow \infty} n^{2\beta-1} \log \mathbb{P} \left(n^{\beta-1} (M_n(t) - m_n(t)) \in B(x - m, \delta/2) \mid \mathcal{A}_n \right) \\ (4) \quad &\geq \liminf_{n \rightarrow \infty} \inf_{\sigma^2 \in B_+(v, \epsilon)} n^{2\beta-1} \log \mathbb{P} \left(\mathcal{N} \left(0, n^{2\beta-1} \sigma^2 \right) \in B(x - m, \delta/2) \right). \end{aligned}$$

Here, $B_+(v, \epsilon) := B(v, \epsilon) \cap [0, \infty)$ and $\mathcal{N}(a, b^2)$ represents a Normally distributed random variable with mean a and variance b^2 . Observe that $B_+(v, \epsilon)$ is nonempty, since $\mathbb{P}(\mathcal{A}_n) > 0$. We are left with evaluating the behavior of the lower bound (4).

If $v - \epsilon > 0$, then (4) equals

$$\begin{aligned} &\liminf_{n \rightarrow \infty} \inf_{\sigma^2 \in B_+(v, \epsilon)} n^{2\beta-1} \log \int_{x-m-\delta/2}^{x-m+\delta/2} \frac{1}{\sqrt{2\pi n^{2\beta-1} \sigma^2}} e^{-\frac{1}{2} \frac{z^2}{n^{2\beta-1} \sigma^2}} dz \\ &= \liminf_{n \rightarrow \infty} \inf_{\sigma^2 \in B_+(v, \epsilon)} n^{2\beta-1} \log \int_{x-m-\delta/2}^{x-m+\delta/2} e^{-\frac{1}{2} \frac{z^2}{n^{2\beta-1} \sigma^2}} dz \\ &= \liminf_{n \rightarrow \infty} n^{2\beta-1} \log \int_{x-m-\delta/2}^{x-m+\delta/2} e^{n^{1-2\beta} \left(-\frac{1}{2} \frac{z^2}{v-\epsilon} \right)} dz \\ &= - \inf_{a \in B(x-m, \delta/2)} \frac{1}{2} \frac{a^2}{v-\epsilon} = - \inf_{a \in B(x, \delta/2)} \mathcal{L}(a - m, \max\{0, v - \epsilon\}). \end{aligned}$$

The penultimate equality follows from Varadhan's Lemma [8, Thm. 4.3.1].

The case that $v - \epsilon \leq 0$ should be dealt with separately. If in addition to $v - \epsilon \leq 0$ also $0 \notin B(x - m, \delta/2)$, then (4) equals

$$\begin{aligned} \liminf_{n \rightarrow \infty} n^{2\beta-1} \log \mathbb{P}(\mathcal{N}(0, 0) \in B(x - m, \delta/2)) &= -\infty \\ &= - \inf_{a \in B(x, \delta/2)} \mathcal{L}(a - m, \max\{0, v - \epsilon\}). \end{aligned}$$

We finally consider the case that $v - \epsilon \leq 0$ and $0 \in B(x - m, \delta/2)$. It is readily seen that in this case there exists $\delta' > 0$ such that $B(0, \delta') \subset B(x - m, \delta/2)$. As a consequence, (4) is bounded below by

$$\begin{aligned} &\liminf_{n \rightarrow \infty} \inf_{\sigma^2 \in B_+(v, \epsilon)} n^{2\beta-1} \log \mathbb{P}(\mathcal{N}(0, n^{2\beta-1}\sigma^2) \in B(0, \delta')) \\ &= \liminf_{n \rightarrow \infty} \inf_{\sigma^2 \in (0, v+\epsilon)} n^{2\beta-1} \log \mathbb{P}(\mathcal{N}(0, n^{2\beta-1}\sigma^2) \in B(0, \delta')) \\ &= \liminf_{n \rightarrow \infty} \inf_{\sigma^2 \in (0, v+\epsilon)} n^{2\beta-1} \log \mathbb{P}(\mathcal{N}(0, 1) \in B(0, \sqrt{n^{1-2\beta}} \delta' / \sqrt{\sigma^2})) \\ &= \liminf_{n \rightarrow \infty} n^{2\beta-1} \log \mathbb{P}(\mathcal{N}(0, 1) \in B(0, \sqrt{n^{1-2\beta}} \delta' / \sqrt{v + \epsilon})) \\ &= 0 = - \inf_{a \in B(x, \delta/2)} \mathcal{L}(a - m, \max\{0, v - \epsilon\}). \end{aligned}$$

Now recalling the definition of \mathcal{A}_n , lower bound (3), and using that $\epsilon = \delta/2$, we thus obtain for any open G ,

$$\begin{aligned} &\liminf_{n \rightarrow \infty} n^{2\beta-1} \log \mathbb{P}\left(n^\beta \left(\frac{M_n(t)}{n} - \rho(t)\right) \in G\right) \\ &\geq \liminf_{n \rightarrow \infty} n^{2\beta-1} \log \mathbb{P}\left(n^\beta \left(\frac{M_n(t)}{n} - \rho(t)\right) \in B(x, \delta) \mid \mathcal{A}_n\right) + \liminf_{n \rightarrow \infty} n^{2\beta-1} \log \mathbb{P}(\mathcal{A}_n) \\ &\geq - \inf_{a \in B(x, \delta/2)} \mathcal{L}(a - m, \max\{0, v - \delta/2\}) - \inf_{(\tilde{m}, \tilde{v}) \in B(m, \delta/2) \times B(v, \delta/2)} \psi(\tilde{m}, \tilde{v}). \end{aligned}$$

The next step is to take $\delta \downarrow 0$ and to use the lower semicontinuity of both \mathcal{L} and ψ , so as to obtain that

$$\liminf_{n \rightarrow \infty} n^{2\beta-1} \log \mathbb{P}\left(n^\beta \left(\frac{M_n(t)}{n} - \rho(t)\right) \in G\right) \geq -\mathcal{L}(x - m, v) - \psi(m, v).$$

Since this holds for any $x \in G$ and $(m, v) \in \Psi$, it follows that

$$\liminf_{n \rightarrow \infty} n^{2\beta-1} \log \mathbb{P}\left(n^\beta \left(\frac{M_n(t)}{n} - \rho(t)\right) \in G\right) \geq - \inf_{x \in G} \inf_{(m, v) \in \Psi} [\mathcal{L}(x - m, v) + \psi(m, v)],$$

as required.

We now turn to the upper bound. To prove this large deviations upper bound, take any closed set $F \subset \mathbb{R}$. Let μ_n denote the probability measure on $\mathbb{R} \times [0, \infty)$ induced by the

random vector (1). Then

$$\begin{aligned}
& \mathbb{P}\left(n^\beta\left(\frac{M_n(t)}{n} - \rho(t)\right) \in F\right) \\
&= \mathbb{P}\left(n^{\beta-1}(M_n(t) - m_n(t)) + n^\beta\left(\frac{m_n(t)}{n} - \rho(t)\right) \in F\right) \\
&= \int_{\mathbb{R} \times [0, \infty)} \mathbb{P}\left(\mathcal{N}\left(0, n^{2\beta-1}v\right) + m \in F\right) d\mu_n(m, v) \\
&\leq \int_{\mathbb{R} \times [0, \infty)} 2 \exp\left(-\inf_{x \in F} \mathcal{L}\left(x - m, n^{2\beta-1}v\right)\right) d\mu_n(m, v) \\
&= \int_{\mathbb{R} \times [0, \infty)} 2 \exp\left(n^{1-2\beta}\left[-\inf_{x \in F} \mathcal{L}(x - m, v)\right]\right) d\mu_n(m, v),
\end{aligned}$$

where the inequality follows from [9, Lemma 4.1]. Consequently, we have, as an immediate application of Varadhan's Lemma [8, Thm. 4.3.1],

$$\begin{aligned}
& \limsup_{n \rightarrow \infty} n^{2\beta-1} \log \mathbb{P}\left(n^\beta\left(\frac{M_n(t)}{n} - \rho(t)\right) \in F\right) \\
&\leq \sup_{(m, v) \in \mathbb{R} \times [0, \infty)} \left[-\inf_{x \in F} \mathcal{L}(x - m, v) - \psi(m, v)\right] \\
&= -\inf_{x \in F} \inf_{(m, v) \in \Psi} [\mathcal{L}(x - m, v) + \psi(m, v)].
\end{aligned}$$

This completes the proof. \square

Remark 2.2. The decomposition featuring in (2) lends itself to an appealing interpretation. The term $\psi(m, v)$ can be interpreted as the 'cost' of the background behaving such that $M_n(t)$ has mean m and variance v , whereas $\mathcal{L}(x - m, v)$ is the 'cost' of a Gaussian random variable with mean m and variance v attaining the value x . Similar decomposition results can be found in e.g. [12, 15, 17].

3. THE MARKOV-MODULATED ORNSTEIN-UHLENBECK PROCESS

In this section we derive the LDP for the Markov-modulated OU (MMOU) process $M(\cdot)$, i.e., the OU process that is modulated by an irreducible Markov chain $J(\cdot)$ with state space $\{1, \dots, d\}$ and generator matrix Q . For ease of exposition, we assume that $M(0) = 0$, but we note that this assumption can be easily lifted. The scaling $J(\cdot) \mapsto J_n(\cdot)$ we consider first corresponds to a linear scaling of the generator matrix, which means that $J_n(\cdot)$ is the Markov-chain generated by nQ ; the section is concluded with a few reflections on the case that Q is scaled sublinearly in n .

We use the following strategy to establish the LDP. First, we show that the empirical measure induced by the Markov chain satisfies the LDP in an appropriate space. Then, we observe that (1) may be viewed as a continuous map defined on this space. Third, we invoke the Contraction Principle [8, Thm. 4.2.1] to derive the LDP for (1) and thus obtain from Thm. 2.1 the LDP for the scaled MMOU process together with the corresponding rate function.

Following [12], we define \mathbb{M}_t as the space of d -dimensional functions $\nu: [0, t] \times \{1, \dots, d\} \rightarrow \mathbb{R}$ that may be represented as $\nu_i(s) = \int_0^s K_\nu(i; r) dr$, where $K_\nu(i; \cdot): [0, t] \rightarrow [0, 1]$ is Borel measurable and satisfies $\sum_{i=1}^d K_\nu(i; r) = 1$ for all $r \in [0, t]$. The function K_ν is referred to as the kernel of ν . We may interpret ν as a function measuring the amount of time a process has spent in each state, where $K_\nu(i; r)$ is the infinitesimal fraction of time the process spends in state i at time r . As \mathbb{M}_t is a subspace of the space of continuous functions, we equip \mathbb{M}_t with the supremum norm.

Let \mathbf{Z}_n be the empirical measure corresponding to the Markov chain $J_n(\cdot)$, i.e.,

$$Z_n(i; s) := \int_0^s 1_{\{J_n(r)=i\}} dr$$

for $i \in \{1, \dots, d\}$ and $s \in [0, t]$. Then \mathbf{Z}_n is a random element of \mathbb{M}_t and [12, Cor. 3.3] asserts that \mathbf{Z}_n satisfies the LDP in \mathbb{M}_t with the corresponding good rate function given by

$$(5) \quad \int_0^t \tilde{I}(K_\nu(s)) ds$$

for $\nu \in \mathbb{M}_t$, where \tilde{I} is the rate function defined as

$$(6) \quad \tilde{I}(K_\nu(s)) := \sup_{\mathbf{u} > \mathbf{0}} \left[- \sum_{i=1}^d \frac{(Q\mathbf{u})_i}{u_i} K_\nu(i; s) \right],$$

where $\mathbf{u} > \mathbf{0}$ is meant coordinatewise; for further background on the underlying LDP see e.g. [11, Thm. IV.14].

We know that $M_n(t)$ has a Normal distribution with random mean $m_n(t)$ and random variance $v_n(t)$. Now the crucial insight is that the joint behavior of $m_n(t)$ and $v_n(t)$ (as captured by (1)) is a continuous function of \mathbf{Z}_n . Indeed, defining

$$\zeta_m(\nu)(t) := \int_0^t \sum_{i=1}^d \alpha(i) K_\nu(i; s) \exp \left(- \int_s^t \sum_{j=1}^d \gamma(j) K_\nu(j; r) dr \right) ds$$

and

$$\zeta_v(\nu)(t) := \int_0^t \sum_{i=1}^d \sigma^2(i) K_\nu(i; s) \exp \left(-2 \int_s^t \sum_{j=1}^d \gamma(j) K_\nu(j; r) dr \right) ds,$$

it is readily verified that $\nu \mapsto (\zeta_m(\nu), \zeta_v(\nu))$ constitutes a continuous map from \mathbb{M}_t to $\mathbb{R} \times [0, \infty)$ and that (1) coincides with $(\zeta_m(\mathbf{Z}_n), \zeta_v(\mathbf{Z}_n))$ when $\beta = 0$ and $\rho(t) = 0$.

Since \mathbf{Z}_n satisfies an LDP in \mathbb{M}_t , it follows from the Contraction Principle that (1) satisfies the LDP in $\mathbb{R} \times [0, \infty)$ with corresponding rate function

$$\psi(m, v) = \inf_{\nu \in \mathbb{M}_t: (\zeta_m(\nu), \zeta_v(\nu)) = (m, v)} \int_0^t \tilde{I}(K_\nu(s)) ds.$$

Then Thm. 2.1 implies that $\frac{1}{n} M_n(t)$ satisfies the LDP with rate function

$$(7) \quad I(x) = \inf_{(m, v) \in \Psi} [\mathcal{L}(x - m, v) + \psi(m, v)].$$

We can also consider the case where the background chain is scaled at a slower speed, that is $Q \mapsto n^\alpha Q$, for $\alpha \in [0, 1)$. With slight abuse of notation, we denote by $J_{n^\alpha}(\cdot)$ a Markov chain with generator $n^\alpha Q$ and by \mathbf{Z}_{n^α} its empirical measure. For the moment, assume that $\alpha > 0$. Then we know from [12, Cor. 3.3] that \mathbf{Z}_{n^α} satisfies the LDP in \mathbb{M}_t with speed n^α and rate function given by Eqn. (5). Since this rate function is finite everywhere (cf. [11, Lem. IV.22]), it follows that

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}(\mathbf{Z}_{n^\alpha} \in G) = \liminf_{n \rightarrow \infty} \frac{1}{n^{1-\alpha}} \frac{1}{n^\alpha} \log \mathbb{P}(\mathbf{Z}_{n^\alpha} \in G) = 0$$

for every nonempty open set $G \subset \mathbb{M}_t$. Consequently, \mathbf{Z}_{n^α} satisfies the LDP in \mathbb{M}_t with speed n and rate function 0 if $\alpha > 0$. This also trivially holds for the case $\alpha = 0$. It follows that, under the scaling $Q \mapsto n^\alpha Q$ for some $\alpha \in [0, 1)$, the random vector (1) satisfies the LDP in $\mathbb{R} \times [0, \infty)$ with speed n ; the corresponding rate function $\psi(m, v)$ is 0 when there exists a ν such that $\zeta_m(\nu)(t) = m$ and $\zeta_v(\nu)(t) = v$, and ∞ elsewhere.

4. REFORMULATION OF THE VARIATIONAL PROBLEM

In this section, we reformulate the variational problem (7) in which the LDP corresponding to MMOU was presented, in a form that is easier to solve numerically. As it will turn out, the new formulation offers additional insight as well.

To this end, we transform the problem that we found in the previous section into the control problem (A) and derive its Hamilton-Jacobi-Bellman equation. By doing so, we reduce the dimension of the control problem from 2 dimensions to just 1. Next, we point out an equivalent (but at first appearance very different) control problem (B), followed by some observations on the numerical procedure. We conclude this section by demonstrating some properties of the model in the slow regime. Note that in this section, we will make use of various concepts in optimal control and calculus of variations. We refer to [3] for an excellent introduction to all of these topics.

4.1. The control problem and its HJB equation. We formulate a variational problem that yields the decay rate of the (rare) event that the Markov-modulated OU process reaches value a at time T , conditional on it starting off from level x at time t ; the scaling considered corresponds to $\beta = 0$ in the terminology of Section 2.

From (7), we have immediately have that

$$I(a) = \inf_{\nu \in \mathbb{M}_t} \left[\mathcal{L}(a - \zeta_m(\nu), \zeta_v(\nu)) + \int_t^T \tilde{I}(K_\nu(s)) ds \right].$$

For numerical convenience, instead of dealing with the vectors ν that record the occupation measures, we will work with their kernels $K_\nu(\cdot)$, which we denote as $\pi(\cdot)$ in this section. They consist of non-negative entries that add up to 1, and can be viewed as local ‘frequencies’ to be in a certain state averaged over a small interval of time. We can show that an optimization over π can be translated into one over ν and vice versa.

Then, by rewriting $\zeta_m(\nu)$ and $\zeta_v(\nu)$, we straightforwardly transform the LDP into variational problem (A): for $t < T$,

$$(A) \quad J^A(x, t) = \inf_{\pi(\cdot)} \left(\frac{(m(T) - a)^2}{2v(T)} + \int_t^T \tilde{I}(\pi(s)) ds \right)$$

where, for $s \in (t, T]$,

$$\begin{aligned} \dot{m}(s) &= \sum_{i=1}^d \pi_i(s) (\alpha_i - \gamma_i m(s)), & m(t) &= x, \\ \dot{v}(s) &= \sum_{i=1}^d \pi_i(s) (\sigma_i^2 - 2\gamma_i v(s)), & v(t) &= 0. \end{aligned}$$

Here $\tilde{I}(\cdot)$ is the large-deviations rate function of Eq. (6).

We now derive the HJB equation for this control problem (which is slightly non-standard due to the occurrence of final costs etc.). In this derivation, we tacitly assume that the value function J^A is differentiable with respect to t and x without proving it, as we feel that a fully rigorous treatment would be beyond the scope of the paper.

By invoking Bellman's principle on time $t + \Delta t$ and position $x + a'\Delta t$ (without loss of generality) we find

$$J^A(x, t) = \inf_{a', \pi(\cdot)} \left(\int_t^{t+\Delta t} \tilde{I}(\pi(s)) ds + \frac{(m(t+\Delta t) - x - a'\Delta t)^2}{2v(t+\Delta t)} + J^A(x + a'\Delta t, t + \Delta t) \right).$$

After 'Tayloring', dividing by Δt and ignoring $o(\Delta t)$ -terms, it turns out that

$$0 = \inf_{a', \pi} \left(\tilde{I}(\pi) + \frac{\sum_i \pi_i(t) (\alpha_i - \gamma_i a) - a'^2}{2 \sum_i \pi_i(t) \sigma_i^2} + a' \partial_x J^A + \partial_t J^A \right),$$

this is due to the fact that

$$m(t + \Delta t) - x - a'\Delta t = m(t) + \Delta t \cdot \sum_i \pi_i(t) (\alpha_i - \gamma_i x) - x - a'\Delta t,$$

while

$$v(t + \Delta t) = \Delta t \sum_i \pi_i(t) \sigma_i^2.$$

As there are no constraints on a' , we can find the minimum by differentiating and finding the zero of the resulting expression:

$$\partial_x J^A = \frac{\sum_i \pi_i (\alpha_i - \gamma_i x) - a'}{\sum_i \pi_i \sigma_i^2},$$

leading to

$$a' = \partial_x J^A \sum_i \pi_i \sigma_i^2 + \sum_i \pi_i (\alpha_i - \gamma_i x).$$

We conclude that we have found the relation

$$(8) \quad \partial_t J^A + \inf_{\pi} \left(\tilde{I}(\pi) + \partial_x J^A \sum_i \pi_i \left(\alpha_i - \gamma_i x - \partial_x J^A \frac{1}{2} \sigma_i^2 \right) \right) = 0.$$

In order to gain some intuition for this HJB equation, we define control problem (B) that has —as it will turn out— the same HJB equation:

$$(B) \quad J^B(x, t) = \inf_{\pi(\cdot), x(\cdot)} \left(\int_t^T \frac{n(s)^2}{2 \sum_i \pi_i(s) \sigma_i^2} ds + \int_t^T \tilde{I}(\pi(s)) ds \right),$$

where, for a given path $x(\cdot)$,

$$n(s) := \dot{x}(s) - \sum_i \pi_i(s) (\alpha_i - \gamma_i x(s)); \quad x(t) = x,$$

for $s \in (t, T]$. Evidently, the value of $x(T)$ should match the target level a . In approach (B) the decay rate is expressed as an integral over ‘local’ costs; it is a special case of the sample-path LDP that was derived in [12, Thm. 3.1].

Problem (B) can be considered as the standard HJB use case:

$$\partial_t J^B + \inf_{\pi, n} \left(\tilde{I}(\pi) + \partial_x J^B \left(\sum_i \pi_i (\alpha_i - \gamma_i a) + n \right) + \frac{n^2}{2 \sum_i \pi_i \sigma_i^2} \right) = 0.$$

The optimization over n can be done explicitly, and doing so, we eventually find the same equation for J^B as for J^A . We thus conclude equivalence of both HJB equations.

Observe that the variational problems (A) and (B) have a markedly different appearance: both involve an optimization over paths $\pi(\cdot)$ (recording the evolution of the state frequencies of the background process), but the other term in the minimization looks significantly different.

4.2. Solving the HJB equation with the Pontryagin minimum principle. We now point out how the HJB equation (8) can be solved numerically. There are several widely different solving strategies (see e.g. [3]). We mention the direct approach (that is, discretizing the partial differential equation and translating the problem into a nonlinear program), dynamical programming approaches, and approaches based on Pontryagin’s minimum principle. We have opted for the latter because as we will see, it is an excellent fit with a minimal numerical burden. Corresponding to the HJB equation (8), we identify the Hamiltonian ¹:

$$H(x, \theta, \pi) = \tilde{I}(\pi) + \theta \sum_i \pi_i \left(\alpha_i - \gamma_i x - \frac{1}{2} \sigma_i^2 \theta \right).$$

From this we readily derive the Pontryagin equations. Following the standard procedure, we find the *state equation*

$$\dot{x}(t) = \partial_\theta H = \sum_i \pi_i(t) (\alpha_i - \gamma_i x(t) - \sigma_i^2 \theta(t)),$$

and the *co-state equation*

$$\dot{\theta}(t) = -\partial_x H = \sum_i \pi_i(t) \gamma_i \theta(t).$$

¹The Hamiltonian in control theory was conceived of by Lev Pontryagin and was inspired by (but different from) the Hamiltonian in mechanics. It can be viewed as an extension of the idea of Langrange multipliers. See e.g. [3].

Due to the local optimality principle, we thus obtain

$$\pi(t) = \arg \min_{\pi} H(x(t), \theta(t)),$$

upon which the boundary conditions $x(0) = x$ and $x(T) = a$ are imposed. As we have two boundary conditions on $x(\cdot)$, but none on $\theta(\cdot)$, we have a so-called *two-point boundary problem*. Also for this subproblem, there are different approaches, but one of the simplest is the shooting method, which consists of repeatedly guessing and refining the value of $\theta(0)$ which leads the system to the desired target a . As $\theta(0)$ is a scalar in this case, and the reached final point is monotone in $\theta(0)$, we can solve it efficiently by what boils down to a secant method.

Once we have solved this system of two ordinary differential equations and one optimization problem, we find the decay rate by plugging the optimal paths $x(\cdot)$ and $\pi(\cdot)$ back into the formula (B).

Note that we can also solve the original control problem (A) directly by Pontryagin's method. However, as this constitutes a control problem with two state dimensions (i.e., m and v), it is numerically more cumbersome than the approach taken here (in particular it requires a two-dimensional shooting method, which is numerically more fragile).

4.3. The slow regime. In this subsection we consider the (limiting) regime where the background chain is *slow*. What we mean by this, is that $\alpha \mapsto n\alpha$ and $\sigma^2 \mapsto n\sigma^2$ as before, but the background process (and thus the entries of the Q matrix) remains unscaled, or is sublinearly scaled (e.g. by a factor n^α for some $\alpha \in [0, 1)$). Intuitively speaking, this effectively entails that the probability for the background chain to follow a deviant path dominates the probability of the OU process doing so, and hence the contribution corresponding to the integral over $I(\pi(t))$ in the control problems (A) and (B) vanishes. We refer to the remarks at the end of Section 3 that indicate how this intuition can be made rigorous.

The Hamiltonian of the control problem under consideration therefore has the following shape:

$$(9) \quad H(x, \theta, \pi) = \sum_{i=1}^d \pi_i \theta \left(\alpha_i - \gamma_i x - \frac{1}{2} \sigma_i^2 \theta \right).$$

We note that the Hamiltonian is stationary along the optimal path $(x^*(t), \theta^*(t), \pi^*(t))$. Indeed, when the final time t of the control problem is fixed and the Hamiltonian does not depend explicitly on time ($\partial_t H \equiv 0$), then

$$H(x^*(t), \theta^*(t), \pi^*(t)) \equiv \text{constant}.$$

Also, according to the Pontryagin principle, the optimal costate $\theta^*(t)$ satisfies

$$\dot{\theta}^*(t) = -\partial_x H = \sum_{i=1}^d \pi_i^*(t) \gamma_i \theta^*(t).$$

We see that $\theta^*(t)$ is a monotone increasing (decreasing, respectively) function under the proviso that $\theta^*(0) > 0$ ($\theta^*(0) < 0$, respectively). We concentrate for now on $\theta^*(0) > 0$, but

an analogous reasoning can be made for $\theta^*(0) < 0$; it is obvious that $\theta^*(0) = 0$ is a trivial case.

Without loss of generality, suppose that the Hamiltonian along the optimal path equals the number H . Then, the optimal path of the state $x^*(\theta)$ (considered as a function of the co-state, which we can do, as the co-state is monotonously increasing), is given by

$$(10) \quad x^*(\theta) = \frac{1}{\theta} \min_{i \in \{1, \dots, d\}} \left(-\frac{\sigma_i^2}{2\gamma_i} \theta^2 + \frac{\alpha_i}{\gamma_i} \theta - \frac{H}{\gamma_i} \right),$$

where we used the fact that $\gamma_i > 0$ for all $i \in \{1, \dots, d\}$. Notice that in this expression the minimum is taken over d concave parabolas (observe that all coefficients $\sigma_i^2/(2\gamma_i)$ are non-negative). In case some σ_i are zero, the corresponding parabolas degenerate into straight lines.

To construct an optimal path $x^*(t)$ satisfying $x^*(0) = x$ and $x^*(T) = a$, we first identify the θ_0 and θ_T satisfying $x^*(\theta_0) = x$ and $x^*(\theta_T) = a$, and then construct the accompanying $\theta^*(t)$. Note that by fixing H , we fix the time T at which the desired end point a is reached, so we have to vary H to find the path that reaches a at time T (which can be done by an elementary bisection procedure).

We proceed by establishing structural properties of the resulting optimal path. The following lemma identifies an upper bound on the number of jumps in the background process along the optimal path.

Lemma 4.1. *In the slow regime, the optimal path of the background process of a OU process modulated by a Markov chain with d states is at most $2d - 2$, and for every d there is a set of parameters for which this maximum number of jumps is reached.*

Proof. We order the parabolas according to the times they become the minimal parabola for the first time. Parabola 1 is necessarily the parabola with the smallest quadratic coefficient $\sigma_i^2/(2\gamma_i)$, parabola 2 is the one that becomes minimal after that, etc. Now, if at a certain point θ_1 the minimal parabola changes from $n \in \{1, \dots, d\}$ back to $m \in \{1, \dots, d\}$ (with $m < n$), then this means that parabolas $m+1, m+2, \dots, n$ cannot be minimal parabolas in (θ_1, ∞) , as they have intersected already twice with parabola m : once on their way to becoming the minimal parabola, and once giving back the top spot to parabola m (at θ_1 for parabola n , before that for the parabolas $m+1, \dots, n-1$). We say that the parabolas $m+1, m+2, \dots, n$ become *inactive* at θ_1 , or remain inactive, in case a previous jump instant made them already so.

A maximal number of jumps is thus achieved when at each jump instant either a fresh parabola becomes minimal, or causes only 1 parabola to become inactive. In such a scheme, each parabola except the first one has two jumps associated with it (when it becomes minimal for the first time and when it becomes inactive). Hence the total number of jumps is at most $2d - 2$. \square

Remark 4.2. We can contrast the result of Lemma 4.1 with the corresponding infinite-server case, where it was found that for d background states there are at most $d - 1$ jumps. This was shown by proving that along the optimal path all states are visited at most once;

see [6] for the variant of the Markov-modulated infinite-server queue in which the service times are sampled upon arrival, and [4] for the variant in which the hazard rate of leaving the system changes is determined by the *current* state of the background process.

Compared to Markov-modulated infinite server queue, in the MMOU case the set of possibilities is much richer. As two striking examples, we mention the ‘necklace’, which follows a path of the type

$$1 \circ 2 \circ 1 \circ 3 \circ 1 \circ 4 \circ 1 \circ 5 \circ 1 \circ \dots \circ 1 \circ d \circ 1$$

(with d visits to state 1, and one visit to the other $d - 1$ states, corresponding with $2d - 2$ jumps) and the ‘lotus’

$$1 \circ 2 \circ 3 \circ \dots \circ (d - 1) \circ d \circ (d - 1) \circ \dots \circ 1$$

(with every state except the d -th being visited twice, also corresponding with $2d - 2$ jumps).

5. NUMERICAL EXPERIMENTS

In this section, we illustrate the numerical procedures detailed in the previous section, as well as the observation of the maximum number of jumps. Specifically we consider Ornstein-Uhlenbeck processes modulated by a two-state background Markov process; to keep the number of parameters limited, we assume that the transition rate matrix is of the form

$$Q = \begin{pmatrix} -q & q \\ q & -q \end{pmatrix}$$

for some $q > 0$.

In the first set of experiments, we choose the parameters $\alpha_1 = -2$, $\alpha_2 = 2$, $\gamma_1 = 1$, $\gamma_2 = 4$, $\sigma_1^2 = 2$, $\sigma_2^2 = 2$, and $q = 0.01$. If an OU process with parameters $\alpha_1, \gamma_1, \sigma_1^2$ were active all the time, it would converge (as $t \rightarrow \infty$) to -2 , if an OU process with parameters $\alpha_2, \gamma_2, \sigma_2^2$ were active all the time, it would converge (as $t \rightarrow \infty$) to $\frac{1}{2}$. On average, the background process is (due to the symmetry of Q) half of the time in state 1, and half of the time in state 2. We fix the starting level at $a_0 = 3$.

In Fig. 1 we plot, as a function of the target level $a \equiv a_t$,

$$J(a) := J(a_t, t),$$

for $t = 1$. The graph is obtained by solving for each value of a the associated HJB equations. We observe that $J(\cdot)$ is convex (as expected), and has the value 0 around the mean position at time $t = 1$. The green dot corresponds with $a = 0.15$, the red dot with $a = -2.13$, and the blue dot with $a = 1.56$. For these three values of a , we show in Fig. 2 the optimal paths of the Markov modulated OU process (which is the function $a(\cdot)$, as was introduced in Section 4.2), and in Fig. 3 the optimal path of the state frequencies $\pi(\cdot)$ of the background process (which in terms of the ‘local fraction of time’ spent in state 1, i.e., $\pi_1(\cdot)$).

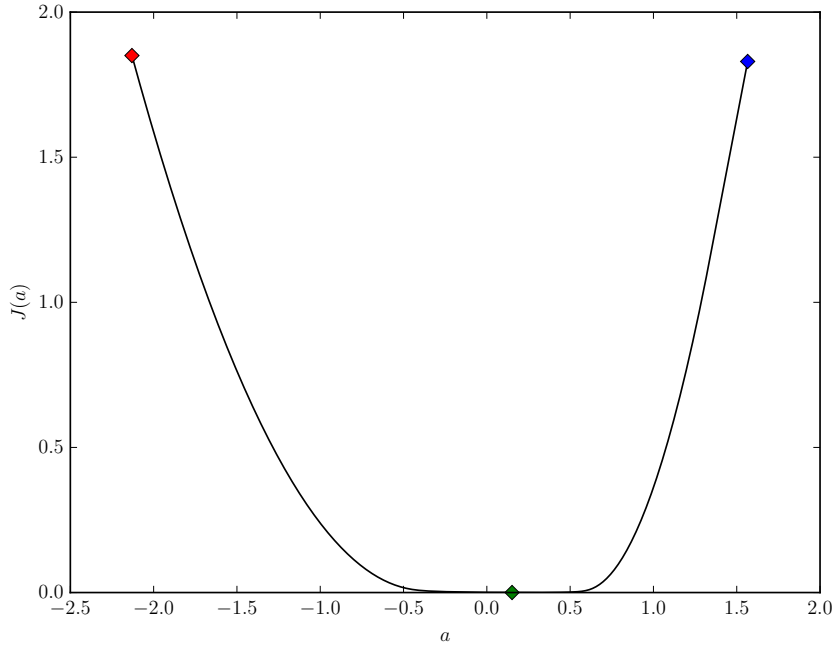


FIGURE 1. Decay rate versus target a reached at time $t = 1$. As explained in the text, the colored dots link this plot to the next two plots.

- As $a = 0.15$ corresponds to the process' expected value at time 1, hardly any effort is needed to reach this value, explaining that the state frequencies are consistently close to $\frac{1}{2}$.
- The target value $a = -2.13$ is substantially smaller than the mean. Fig. 3 indicates that the most likely way of reaching this target value is (roughly) by staying in state 2 until $t = 0.2$, and then staying in state 1 until $t = 1$; this makes sense, as state 2 corresponds with a long term of -2 , as we observed above. Because q is relatively small we are 'close to' the slow regime described in Section 4.3; this explains why the transition of the path in Fig. 3 from the value 0 to the value 1 is rather sharp. For the same parameter set and $q = 0$, we would find a background path with 1 transition.
- The target value $a = 1.56$ is larger than the mean. The graph shows that in this case it is apparently optimal to be in state 1 till roughly $t = 0.6$, then visit state 2 between 0.6 and 0.85 (taking advantage of the higher variance), and then to be again in state 1 until time 1 (taking advantage of the higher mean); again the transitions along this path are rather sharp, and correspond to a path with 2 jumps, which is the maximal amount for a 2-state Markov chain.

In the second set of experiments, we take the same parameters as before (i.e., $\alpha_1 = -2$, $\alpha_2 = 2$, $\gamma_1 = 1$, $\gamma_2 = 4$, $\sigma_1^2 = 2$, $\sigma_2^2 = 2$), but we fix the target value (at $a = 1.56$;

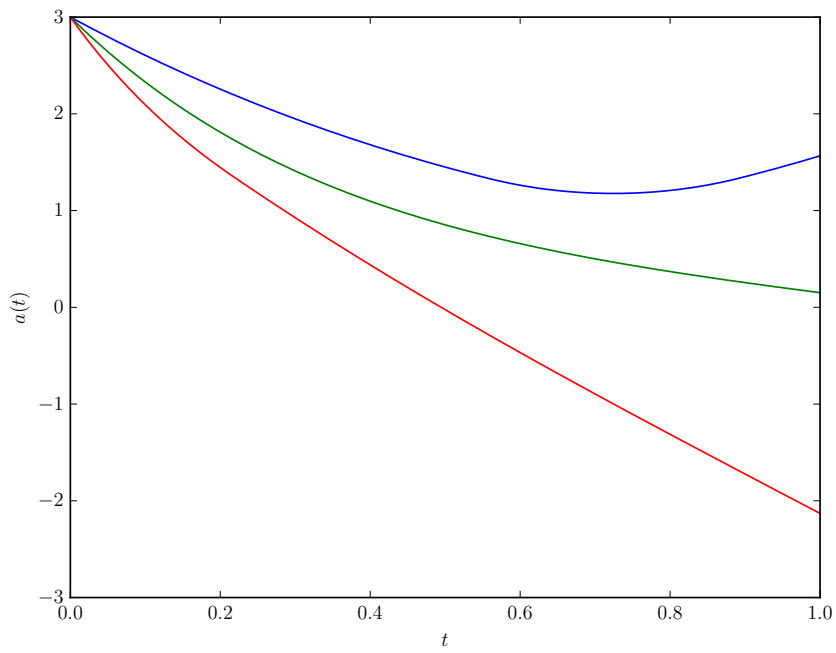


FIGURE 2. The optimal paths corresponding to the three scenarios in Fig. 1.

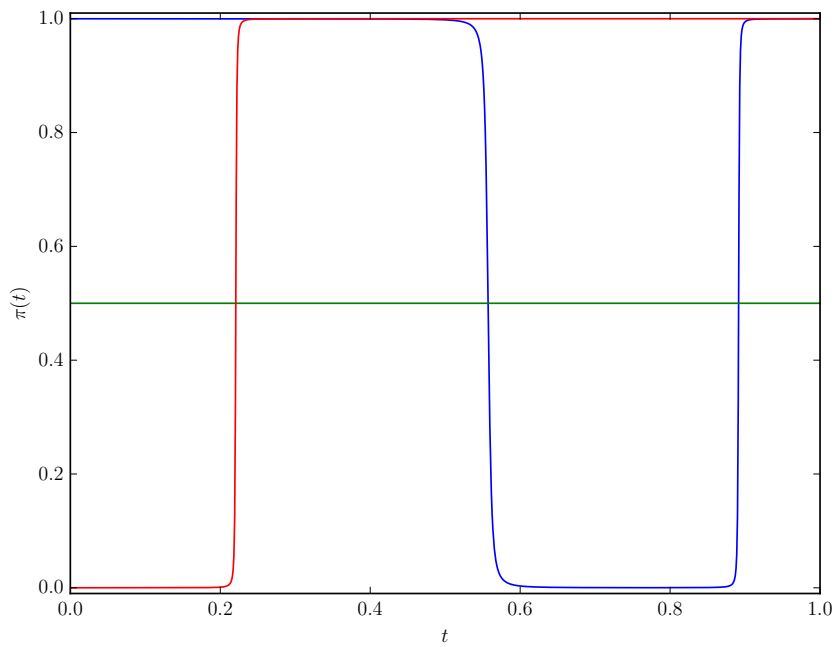


FIGURE 3. The background paths corresponding to the three scenarios in Fig. 1.

corresponding to the blue dot in the first set of experiments) and vary the value of q . We observe that the decay rate takes on the form of an ‘S’. This can be explained as follows.

- When q is small, then the cost of varying the background path is small as well, and the decay rate is dominated by the ‘OU’ term (which explains why we have an almost flat section in the curve), and we find the familiar two jump background path (see blue curve).
- As q and thus the cost of altering the background path increases, the optimal background path gets more rounded and follows the sharp transitions we found earlier only approximatively (green curve).
- For large q (red curve), the cost of altering the background path dominates, so the background path keeps close to the stationary value. As this path carries no ‘cost’, the decay rate is again almost flat for large q .
- The fact that the decay rate is increasing when q is increasing can be intuited as follows. As we already observed in Remark 2.2, the rate function can be thought of as the cost of deviating from the mean behavior. The rate function featuring in the LDP for a modulated OU process consists of two cost terms: the cost of a normal distribution deviating from its mean and the cost of the background process deviating from its mean to produce parameters m and v . When the background process J is the Markov chain described above, then a larger q implies that J converges faster to its mean (equilibrium) behavior. Consequently, a larger q makes it more difficult to deviate from the mean behavior and thus leads to a larger cost. This explains why the decay rate is increasing as q is increasing.

The intuitive reasoning above is, in fact, already reflected in the form of the rate function \tilde{I} in Eqn. (6). Indeed, this rate function is a linear function of the generator matrix of J , so increasing q will increase the decay rate.

To conclude this section, let us consider the slow regime in detail (for the same set of parameters and targetting $a = 1.56$ at time $t = 1$ as before). From our previous numerical calculations, we determine that this target corresponds with $\theta(0) = -0.4$, while $\theta(1) = -20$ and from this we have from (9) that $H = 1.84$. In the final figure, we plot the two parabolas of Eqn. (10) and see that indeed two bumps are encountered.

REFERENCES

- [1] L. ALILI, P. PATIE, and J. PEDERSEN (2005). Representations of the first hitting time density of an Ornstein-Uhlenbeck process. *Stochastic Models*, Vol. 21, pp. 967-980.
- [2] D. ANDERSON, J. BLOM, M. MANDJES, H. THORSODDOTTIR, and K. DE TURCK (2016). A functional central limit theorem for a Markov-modulated infinite-server queue. *Methodology and Computing in Applied Probability*, Vol. 18, pp. 153-168.
- [3] D. BERTSEKAS (1995). *Dynamic programming and optimal control*. Athena Scientific.
- [4] J. BLOM, K. DE TURCK, O. KELLA, and M. MANDJES (2014). Tail asymptotics of a Markov-modulated infinite-server queue. *Queueing Systems*, Vol. 78, pp. 337-357.
- [5] J. BLOM, K. DE TURCK, and M. MANDJES (2015). Analysis of Markov-modulated infinite-server queues in the central-limit regime. *Probability in the Engineering and Informational Sciences*, Vol. 29, pp. 433-459.
- [6] J. BLOM and M. MANDJES (2013). A large-deviations analysis of Markov-modulated infinite-server queues. *Operations Research Letters*, Vol. 41, pp. 220-225.

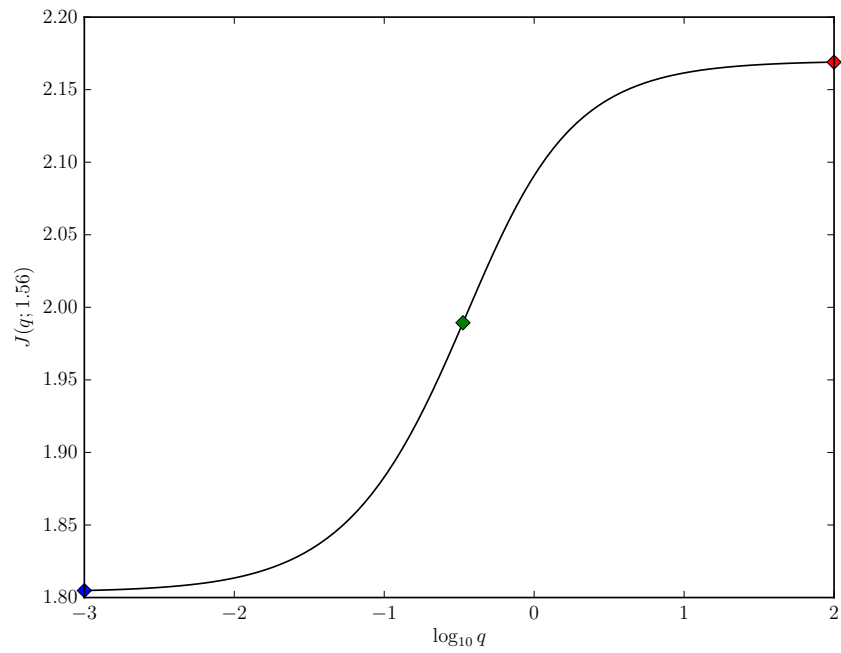


FIGURE 4. Decay rate versus the speed q of the Markov chain. As explained in the text, the colored dots link this plot to the next two plots.

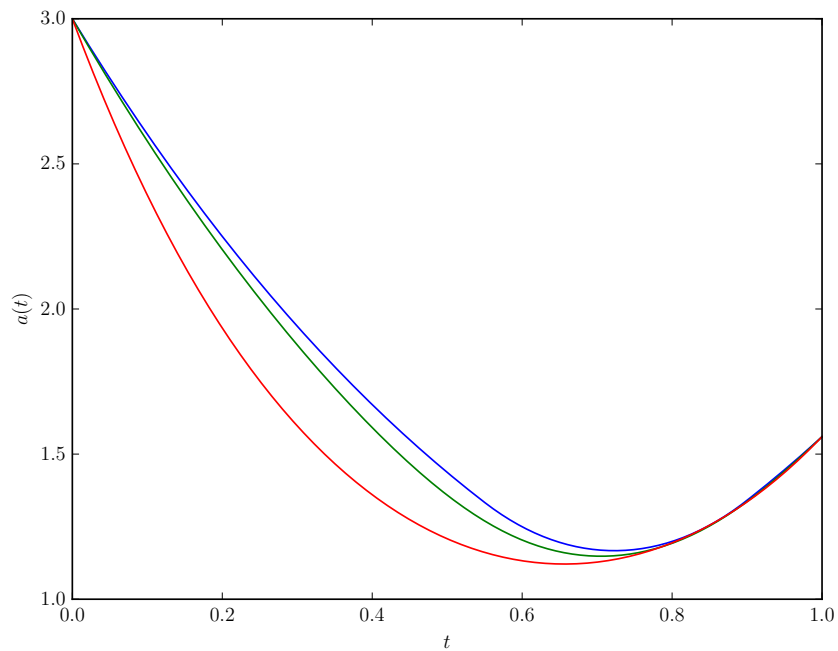


FIGURE 5. The optimal paths corresponding to the three scenarios in Fig. 4.

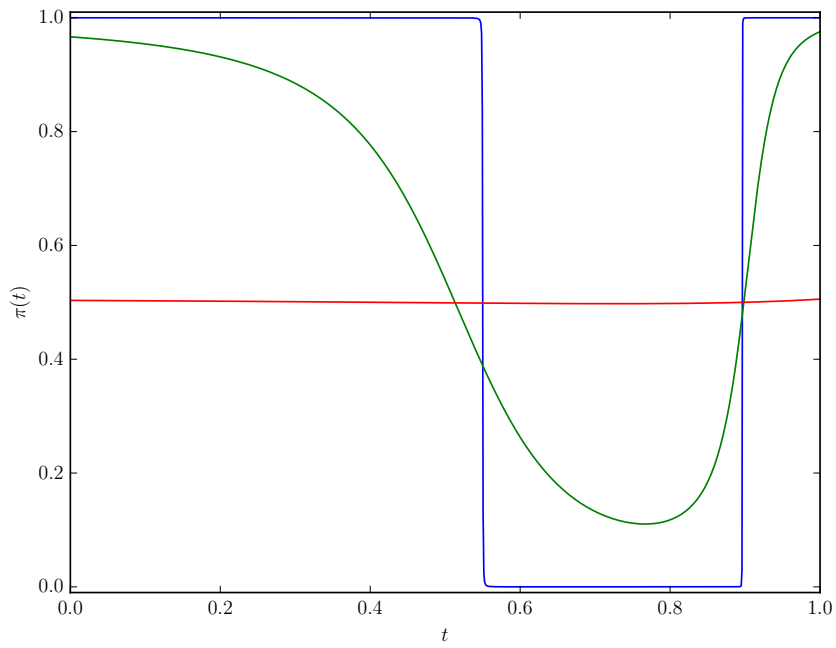


FIGURE 6. The background paths corresponding to the three scenarios in Fig. 4.

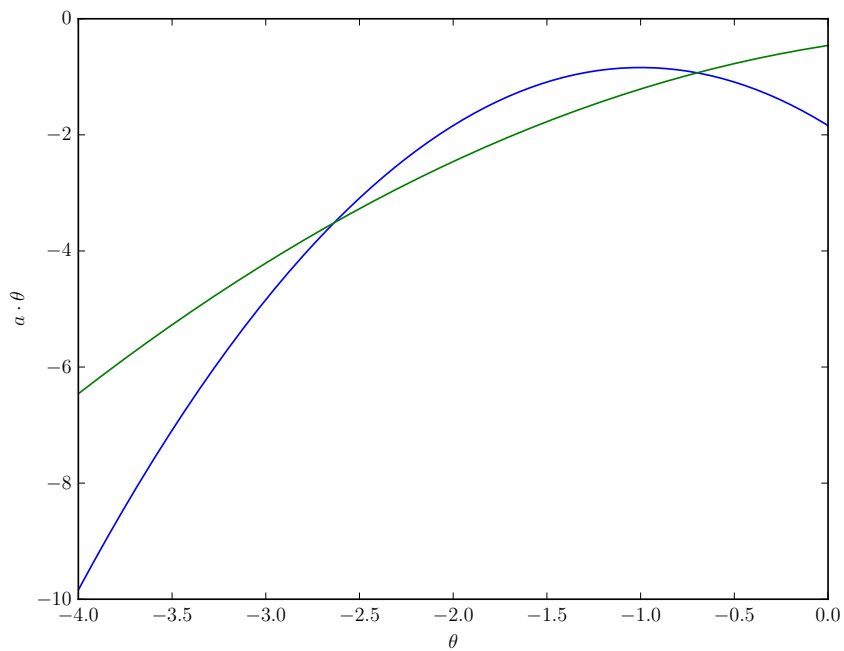


FIGURE 7. A plot of the optimal path in the θ - $x\theta$ plane.

- [7] L. BO, Y. WANG, and X. YANG (2011). First passage times of (reflected) Ornstein-Uhlenbeck processes over random jump boundaries. *Journal of Applied Probability*, Vol. 48, pp. 723-732.
- [8] A. DEMBO and O. ZEITOUNI (1998). *Large Deviations Techniques and Applications*, Springer, New York.
- [9] I. DINWOODIE and S. ZABELL (1992). Large deviations for exchangeable random vectors. *Annals of Probability*, Vol. 20, pp. 1147-1166.
- [10] A. GANESH, N. O'CONNELL, and D. WISCHIK (2004). *Big Queues*, Springer, New York.
- [11] F. DEN HOLLANDER (2000). *Large Deviations*. Fields Institute Monographs 14. AMS, Providence.
- [12] G. HUANG, M. MANDJES, and P. SPREIJ (2016). Large deviations for Markov-modulated diffusion processes with rapid switching. *Stochastic Processes and their Applications*, Vol. 126, pp. 1785-1818.
- [13] G. HUANG, H.M. JANSEN, M. MANDJES, P. SPREIJ, and K. DE TURCK (2016). Markov-modulated Ornstein-Uhlenbeck processes. *Advances in Applied Probability*, Vol. 48, pp 235-254.
- [14] H. M. JANSEN, M. MANDJES, K. DE TURCK, and S. WITTEVRONGEL (2016). A large deviations principle for infinite-server queues in a random environment. *Queueing Systems*, Vol. 82, pp 199-235.
- [15] R. LIPTSER (1996). Large deviations for two scaled diffusions. *Probability Theory and Related Fields*, Vol. 106, pp. 71-104.
- [16] P. ROBERT (2003). *Stochastic Networks and Queues*. Springer, New York.
- [17] K. DE TURCK and M. MANDJES (2014). Large deviations of an infinite-server system with a linearly scaled background process. *Performance Evaluation*, Vol. 75-76, pp. 36-49.
- [18] W. WHITT (2002). *Stochastic-process Limits*. Springer, New York.