

Efficiency of the Girsanov Transformation Approach for Parametric Sensitivity Analysis of Stochastic Chemical Kinetics*

Ting Wang[†] and Muruhan Rathinam[‡]

Abstract. The most common Monte Carlo methods for sensitivity analysis of stochastic reaction networks are the finite difference (FD), Girsanov transformation (GT), and regularized pathwise derivative (RPD) methods. It has been numerically observed in the literature that the biased FD and RPD methods tend to have lower variance than the unbiased GT method and that centering the GT method (CGT) reduces its variance. We provide a theoretical justification for these observations in terms of system size asymptotic analysis under what is known as the classical scaling. Our analysis applies to GT, CGT, and FD and shows that the standard deviations of their estimators when normalized by the actual sensitivity scale as $\mathcal{O}(N^{1/2})$, $\mathcal{O}(1)$, and $\mathcal{O}(N^{-1/2})$, respectively, as system size $N \rightarrow \infty$. In the case of the FD methods, the $N \rightarrow \infty$ asymptotics are obtained keeping the finite difference perturbation h fixed. Our numerical examples verify that our order estimates are sharp and that the variance of the RPD method scales similarly to the FD methods. We combine our large N asymptotics with previously known small h asymptotics to obtain the best choice of h in terms of N and estimate the number N_s of simulations required to achieve a prescribed relative \mathcal{L}_2 error δ . This shows that N_s depends on δ and N as $\delta^{-2-\frac{\gamma_2}{\gamma_1}} N^{-1}$, δ^{-2} , and $N\delta^{-2}$ for FD, CGT, and GT, respectively. Here $\gamma_1 > 0$, $\gamma_2 > 0$ depend on the type of FD method used.

Key words. stochastic chemical kinetics, Girsanov transformation, asymptotic analysis, parametric sensitivity, finite difference, variance analysis

AMS subject classifications. Primary, 60H35, 65C99; Secondary, 92C42, 92C45

DOI. 10.1137/140998111

1. Introduction. Estimation of parametric sensitivities of dynamical systems is an essential part of the modeling and parameter estimation process. For instance, the problem of finding the set of parameters that best fit some observed data can be formulated as an optimization problem over the parameter space where the partial derivatives of the objective function depend on the parametric sensitivities defined as partial derivatives of some system output with respect to the parameters.

In deterministic dynamical systems governed by ordinary differential equations (ODEs), the sensitivities defined by the partial derivatives $\partial f(X(t))/\partial c_k$ of some function f of the state with respect to the parameters are essentially computed by numerical integration of an auxiliary system of evolution equations obtained by linearization of the original ODEs. In contrast, for stochastic dynamical systems, several vastly different approaches exist. We note

*Received by the editors December 3, 2014; accepted for publication (in revised form) August 25, 2016; published electronically November 1, 2016.

<http://www.siam.org/journals/juq/4/99811.html>

[†]Department of Mathematical Sciences, University of Delaware, Newark, DE 19716 (tingw@udel.edu).

[‡]Department of Mathematics and Statistics, University of Maryland, Baltimore County, Baltimore, MD 21250 (muruhan@umbc.edu).

that we shall treat the parameters c_k as deterministic and not as random quantities, while the dynamic behavior of the systems we consider is stochastic.

Our primary focus will be stochastically modeled chemical reaction systems. While the stochastic chemical kinetic model under the well stirred assumption [10] has been around for decades, it wasn't until the late nineties that the importance of stochastic chemical models in some applications was realized [4, 19]. Especially, intracellular chemical reactions systems often contain certain molecular species in small copy numbers, and as such, the deterministic model based on ODEs or partial differential equations (PDEs) for the concentrations of various molecular species is not appropriate. A more appropriate model, under the well stirred assumption, consists of a continuous time Markov process $X(t)$ with the nonnegative integer lattice \mathbb{Z}_+^n as state space.

While we focus on stochastic chemical kinetics which we describe in the next subsection, we note that analogous models appear in other fields such as epidemiology and predator-prey systems.

1.1. Stochastic chemical kinetics. As a simple example, let us consider the chemical reaction system



consisting of three species S_1, S_2 , and S_3 undergoing two reaction channels. The state space is the set \mathbb{Z}_+^3 of nonnegative three-dimensional integer vectors, where the state $x = (x_1, x_2, x_3)$ describes the copy numbers x_1 of S_1 , x_2 of S_2 , and x_3 of S_3 . When the first reaction channel fires, the state changes by $\nu_1 = (-1, -1, 1)^T$, and when the second reaction channel fires it changes by $\nu_2 = (1, 1, -1)^T$. The quantities ν_j are known as *stoichiometric vectors* and for chemical reaction systems the ν_j are parameters and state independent. The “probabilistic rate” at which these two reactions occur is given by the *intensity functions* $a_1(x, c)$ and $a_2(x, c)$ (where c is a vector of parameters). The precise meaning of the intensity functions is as follows. If $X(t) = (X_1(t), X_2(t), X_3(t))$ is the stochastic process of species counts, then given $X(t) = x$, the probability of at least one firing of the j th reaction channel during interval $(t, t + h]$ is $a_j(x, c)h + o(h)$ as $h \rightarrow 0^+$.

Stochastic mass action form. Under the well stirred model of Gillespie [10], the intensity functions take the *stochastic mass action* form: $a_1(x, c) = c_1 x_1 x_2$ and $a_2(x, c) = c_2 x_3$. The rationale for this specific form is based on the following considerations. The probability that a given pair of one S_1 molecule and one S_2 molecule comes together and react during time interval $(t, t + h]$ is given by $c_1 h + o(h)$, where c_1 is a constant. Given that there are $x_1 x_2$ different ways to choose the pair, we obtain the probability of $c_1 x_1 x_2 h + o(h)$ for any pair of S_1 and S_2 to react. Likewise, the probability that a given S_3 molecule gives rise to an S_1 and an S_2 via the second reaction during $(t, t + h]$ is given by $c_2 h + o(h)$, where c_2 is a constant. Given that there are x_3 different S_3 molecules, we obtain the probability of $c_2 x_3 h + o(h)$ for any of the S_3 to react.

General chemical system. More generally, a chemical reaction system consists of m reaction channels and n chemical species $\{S_1, \dots, S_n\}$. The n -dimensional state vector $X(t)$ characterizes the state of the system where each entry $X_i(t)$ represents the number of molecules of the species S_i at time t . The firing of a reaction channel $j \in \{1, \dots, m\}$ at time t causes

the state to be incremented by the stoichiometric vector ν_j . We assume that X is *càdlàg*, i.e., paths of X are right continuous with left-hand limits and hence, if reaction channel j fires at time t , then $X(t) = X(t-) + \nu_j$. For $j = 1, \dots, m$ we denote by $R_j(t)$ the number of firings of the j th reaction channel during $(0, t]$. Thus $X(t) = X(0) + \nu R(t)$ for $t \geq 0$, where ν is the *stoichiometric matrix* whose j th column is ν_j and $R(t) = (R_1(t), \dots, R_m(t))^T$. We note that $R(0) = 0$ and $R_j(t) - R_j(t-)$ is either 0 or 1. The process X is assumed to be Markovian, and associated with each reaction channel is an *intensity function* (also known as *propensity function* in the chemical kinetics literature) $a_j(x, c), j = 1, \dots, m$, which is such that, given $X(t) = x$, the probability of one or more firings of reaction channel j during $(t, t + h]$ is $a_j(x, c)h + o(h)$ as $h \rightarrow 0^+$. Here, c are parameters. Following the terminology of [8], we note that R_j are counting processes which admit the \mathcal{F}_t -predictable intensity process $a_j(X(t-), c)$, where \mathcal{F}_t is the filtration generated by X and R .

Random time change representation. Naturally, the probability laws of the stochastic processes X and R depend on the parameters c . For the purpose of analysis, it proves convenient to find a way to represent the processes X and R corresponding to different c values on the same sample space $(\Omega, \mathcal{F}, \mathbb{P})$. To this end, we use the *random time change representation* [9] to express X via the stochastic equation

$$(2) \quad X(t, c) = x_0 + \sum_{j=1}^m Y_j \left(\int_0^t a_j(X(s, c), c) ds \right) \nu_j,$$

where Y_j are independent unit rate Poisson processes. It follows that

$$(3) \quad R_j(t, c) = Y_j \left(\int_0^t a_j(X(s, c), c) ds \right), \quad j = 1, \dots, m,$$

where x_0 is the initial state assumed to be deterministic. We note that in this representation, we have a family of stochastic processes $X(t, c)$ and $R(t, c)$ on the same sample space $(\Omega, \mathcal{F}, \mathbb{P})$ where each element $\omega \in \Omega$ may be identified with a specific trajectory of $Y(t) = (Y_1(t), \dots, Y_m(t))$, the underlying unit rate independent Poissons. We note that the $Y_j(t)$ do not depend on the parameters. See [24] for a detailed explanation of how to compute $X(t, c)$ once a sample path of $Y(t)$ is generated.

1.2. Parametric sensitivity estimation. We consider parametric sensitivities of the stochastic process $X(t, c)$ with respect to an output function $f : \mathbb{Z}_+^n \rightarrow \mathbb{R}$, defined by the partial derivatives

$$\frac{\partial}{\partial c_k} \mathbb{E}(f(X(t, c))),$$

where c_k are scalar parameters, f is some suitable scalar function of the state space, \mathbb{E} is the expectation, and $t > 0$ is some fixed final time. For simplicity we shall focus on one scalar parameter c . When the number of species n is large (in several applications it is of the order of 10–100), due to the curse of dimensionality, Monte Carlo approaches are the most viable for both simulation of the process X as well as estimation of sensitivities. Monte Carlo simulation of exact sample paths of the process X is feasible and is provided by the well-known SSA or Gillespie algorithm [10]. In this context several different Monte Carlo approaches exist for the numerical computation of the parametric sensitivities as well.

We shall use $\mathcal{S}(t, c)$ to denote the exact sensitivity

$$(4) \quad \mathcal{S}(t, c) = \frac{\partial}{\partial c} \mathbb{E}(f(X(t, c))).$$

As we will see later in this section, all the Monte Carlo methods for computing the sensitivity involve the estimation of the expected value $\mathbb{E}(S(t, c))$ of some process $S(t, c)$ at time $t > 0$, via independent and identically distributed sample estimation, where $S(t, c)$ can be computed easily from the knowledge of system parameters, the function f , and the sample path of X on the time interval $[0, t]$. In other words, one generates N_s independent copies $X^{(i)}(t, c)$ of $X(t, c)$ for $i = 1, \dots, N_s$ and then computes the corresponding copies $S^{(i)}(t, c)$ of $S(t, c)$. Then the sensitivity is estimated by

$$\bar{S}(t, c) = \frac{1}{N_s} \sum_{i=1}^{N_s} S^{(i)}(t, c).$$

Since $\mathbb{E}(\bar{S}(t, c)) = \mathbb{E}(S(t, c))$ and $\text{Var}(\bar{S}(t, c)) = \text{Var}(S(t, c))/N_s$, the accuracy of this estimate depends on the error (known as bias) $\mathbb{E}(S(t, c) - \mathcal{S}(t, c))$, the variance $\text{Var}(S(t, c))$ of the underlying estimator $S(t, c)$, and the sample size N_s .

One way to quantify the error in estimation is via the mean square error:

$$(5) \quad \mathbb{E}(|\bar{S}(t, c) - \mathcal{S}(t, c)|^2) = \frac{\text{Var}(S(t, c))}{N_s} + (\mathbb{E}(S(t, c) - \mathcal{S}(t, c)))^2.$$

If $\text{Var}(S(t, c))$ is large, then one requires a greater number N_s of simulations, resulting in loss of efficiency. On the other hand if a biased estimator is used, increasing the number of simulations N_s does not help. It is often useful to consider the *relative error* (RE) defined by

$$(6) \quad \text{RE} = \sqrt{\mathbb{E}(|\bar{S}(t, c) - \mathcal{S}(t, c)|^2)} / |\mathcal{S}(t, c)|,$$

provided the true sensitivity $\mathcal{S}(t, c)$ is nonzero.

Throughout this paper, we shall refer to $S(t, c)$ as the *underlying estimator* or simply the *estimator* and $\bar{S}(t, c)$ as the *ultimate estimator*. As the properties of the latter depend directly on that of the former and N_s , the analysis of the variance of the underlying estimator $S(t, c)$ shall be our focus. We define the *relative standard deviation* (RSD) and the *relative bias* (RB) of the underlying estimator $S(t, c)$ by

$$(7) \quad \text{RSD} = \sqrt{\text{Var}(S(t, c))} / |S(t, c)|$$

and

$$(8) \quad \text{RB} = \mathbb{E}(S(t, c) - \mathcal{S}(t, c)) / |\mathcal{S}(t, c)|$$

when $\mathcal{S}(t, c) \neq 0$. We note that the RE is given by

$$(9) \quad \text{RE} = \sqrt{\frac{\text{RSD}^2}{N_s} + \text{RB}^2}.$$

Now, we turn our attention to the description of some common Monte Carlo sensitivity estimators. As a general reference on this topic we suggest [5, 12]. The Monte Carlo methods for sensitivity can broadly be categorized into finite difference (FD) methods [1, 5, 24], pathwise derivative (PD) methods [5, 26], and likelihood ratio or Girsanov transformation (GT) methods [5, 20].

The FD methods involve approximation of the partial derivative by the simple finite difference $\mathbb{E}[f(X(t, c+h)) - f(X(t, c))]/h$ or some higher order finite difference. In the case of the simple finite difference above,

$$(10) \quad S_{\text{FD}}(t, c) = h^{-1}[f(X(t, c+h)) - f(X(t, c))].$$

Thus, $\mathbb{E}(S_{\text{FD}}(t, c)) \neq \frac{\partial}{\partial c} \mathbb{E}(f(X(t, c)))$ in general, and the bias is decreased by decreasing h . On the other hand,

$$\text{Var}(S_{\text{FD}}(t, c)) = h^{-2} \{ \text{Var}(f(X(t, c+h))) + \text{Var}(f(X(t, c))) - 2\text{Cov}(f(X(t, c+h)), f(X(t, c))) \}.$$

In general the numerator does not vanish as fast as h^2 when $h \rightarrow 0$, showing that small h leads to large variance. When $f(X(t, c+h))$ and $f(X(t, c))$ are strongly positively correlated, one may expect the variance to be small. If the processes $X(t, c)$ and $X(t, c+h)$ are taken to be independent, which is accomplished by the use of two independent streams of random numbers in the simulation, the resulting FD method is known as the *independent random number* (IRN) method. If the processes $X(t, c)$ and $X(t, c+h)$ are strongly coupled, which is accomplished by the use of a common random number stream, the resulting approach is known as *common random number* (CRN) method. In general, the CRN FD methods have much lower variance than the IRN FD methods. Moreover, different approaches to couple the processes $X(t, c+h)$ and $X(t, c)$ lead to different covariances and hence different variances for the FD estimators. See [1, 24] for some approaches.

In the PD method one takes

$$S_{\text{PD}}(t, c) = \frac{\partial}{\partial c} f(X(t, c)),$$

and the method is applicable provided the derivative exists, analytical computation of the derivative is possible, and the commutation

$$(11) \quad \mathbb{E} \left(\frac{\partial}{\partial c} f(X(t, c)) \right) = \frac{\partial}{\partial c} \mathbb{E}(f(X(t, c)))$$

holds. In the context of stochastic chemical kinetics, direct application of the PD method is not valid as the commutation in (11) does not hold. To see this, note that $f(X(t, c, \omega))$ is piecewise constant in c for fixed t and ω and hence the derivative is 0, while the sensitivity $\partial \mathbb{E}(f(X(t, c)))/\partial c$ is in general nonzero, showing that the commutation in (11) is not valid (see [26] for details). It is possible to regularize the problem by replacing $\partial f(X(t, c))/\partial c$ with

$$(12) \quad S_{\text{RPD}}(t, c) = \frac{\partial}{\partial c} \left(\frac{1}{2w} \int_{t-w}^{t+w} f(X(s, c)) ds \right)$$

to obtain the *regularized pathwise derivative* (RPD) estimator for which the commutation of derivative with expectation holds for a restricted class of examples [26]. This, however,

results in a bias which increases with large w . Also see [12] for similar work in the context of computing the sensitivity of path integrals.

The GT approach may be motivated in different ways. For the purpose of our analysis based on the random time change representation, it is natural to start with the family of processes $X(t, c)$ parametrized by c that are all defined on $(\Omega, \mathcal{F}, \mathbb{P})$ as mentioned before. Suppose the sensitivity is required at a specific parameter value $c = c_0$. Under certain regularity conditions, a family of new probability measures $P(c)$ may be constructed on the same sample space (Ω, \mathcal{F}) for a range of c values in a neighborhood of c_0 so that $P(c_0) = \mathbb{P}$, i.e., coincides with the original probability measure (see [8], for instance). Moreover, the probability measures $P(c)$ are absolutely continuous with respect to $P(c_0)$ and the $P(c)$ -law of the process $X(t, c_0)$ is the same as the $P(c_0) (= \mathbb{P})$ -law of the process $X(t, c)$. In other words, for all suitable functions f ,

$$\int_{\Omega} f(X(t, c)) dP(c_0) = \int_{\Omega} f(X(t, c_0)) dP(c).$$

We observe that the left-hand side is $\mathbb{E}(f(X(t, c)))$. If we denote by $L(t, c, c_0)$ the Radon–Nikodym derivative $dP(c)/dP(c_0)$, then we have

$$\begin{aligned} (13) \quad \frac{\partial}{\partial c} \Big|_{c=c_0} \mathbb{E}(f(X(t, c))) &= \frac{\partial}{\partial c} \Big|_{c=c_0} \int_{\Omega} f(X(t, c_0)) L(t, c, c_0) dP(c_0) \\ &= \int_{\Omega} f(X(t, c_0)) \frac{\partial}{\partial c} \Big|_{c=c_0} L(t, c, c_0) dP(c_0) \end{aligned}$$

provided the differentiation inside the integral is valid. It turns out that

$$(14) \quad Z(t, c_0) = \frac{\partial}{\partial c} \Big|_{c=c_0} L(t, c, c_0)$$

is analytically tractable and the required sensitivity is given by

$$\frac{\partial}{\partial c} \Big|_{c=c_0} \mathbb{E}(f(X(t, c))) = \mathbb{E}[f(X(t, c_0)) Z(t, c_0)],$$

thus the sensitivity estimator $S(t, c_0) = f(X(t, c_0)) Z(t, c_0)$.

In the context of stochastic chemical kinetics, the weight process Z defined by (14) is given by [20, 26]

$$(15) \quad Z(t, c) = \sum_{j=1}^m \int_0^t \frac{\partial a_j}{\partial c}(X(s-, c), c) dR_j(s, c) - \sum_{j=1}^m \int_0^t \frac{\partial a_j}{\partial c}(X(s, c), c) ds.$$

We have dropped c_0 in favor of c for notational ease; however, it must be noted that all computations are carried out at the specific parameter value c at which the sensitivity is required.

We also investigate a modified GT method inspired by the work in [28], which we call the *centered Girsanov transformation* (CGT) method, in which we replace the estimator

$f(X(t, c))Z(t, c)$ with $(f(X(t, c)) - \mathbb{E}(f(X(t, c))))Z(t, c)$. Since $Z(t, c)$ has zero mean this new estimator has the same mean as the original one and hence is also unbiased. In practice $\mathbb{E}(f(X(t, c)))$ is not known and needs to be estimated as well. One approach would be to generate N_s independent copies $X^i(t, c)$ of $X(t, c)$, then use

$$\overline{f(X(t, c))} = \frac{1}{N_s} \sum_{i=1}^{N_s} f(X^i(t, c))$$

to estimate $\mathbb{E}(f(X(t, c)))$, and then use

$$\bar{S}_{\text{CGT}} = \frac{1}{N_s} \sum_{i=1}^{N_s} \left(f(X^i(t, c)) - \overline{f(X(t, c))} \right) Z^{(i)}(t, c)$$

as the ultimate estimator. In this case $\mathbb{E}(\bar{S}_{\text{CGT}}) \neq \mathbb{E}(f(X(t, c))Z(t, c))$ and the estimator is biased. However, when N_s is large the bias is $\mathcal{O}(1/N_s)$. Also

$$\text{Var}(\bar{S}_{\text{CGT}}) = \text{Var}(S_{\text{CGT}})/N_s + \mathcal{O}(1/N_s^2),$$

where $S_{\text{CGT}} = (f(X(t, c)) - \mathbb{E}(f(X(t, c))))Z(t, c)$ is the underlying CGT estimator. So it is adequate to study the variance of $(f(X(t, c)) - \mathbb{E}(f(X(t, c))))Z(t, c)$. In the formula used in [28] for the ultimate estimator, $Z^{(i)}$ above were replaced by $Z^{(i)} - \bar{Z}$, where \bar{Z} was the sample mean of $Z^{(i)}$. When the sample size N_s is large, both ultimate estimators are similar. For the purpose of analysis, we shall focus on the underlying CGT estimator

$$(16) \quad S_{\text{CGT}} = f(X(t, c))Z(t, c) - \mathbb{E}(f(X(t, c)))Z(t, c).$$

We note that the variances of the GT and CGT estimators are given by the following formulae:

$$(17) \quad \begin{aligned} \text{Var}(S_{\text{GT}}) &= \mathbb{E}((f(X(t, c)))^2 Z^2(t, c)) - \mathbb{E}^2(f(X(t, c))Z(t, c)), \\ \text{Var}(S_{\text{CGT}}) &= \text{Var}(S_{\text{GT}}) - 2\mathbb{E}(f(X(t, c))Z^2(t, c)) + \mathbb{E}^2(f(X(t, c)))\mathbb{E}(Z^2(t, c)). \end{aligned}$$

It must be noted that it is not always the case that $\text{Var}(S_{\text{GT}})$ is greater than or equal to $\text{Var}(S_{\text{CGT}})$. Thus, one cannot conclude that CGT is always superior to GT. However, it was observed in [28] as well as in our simulations that CGT tends to have lower variance than GT in most examples.

Recently introduced methods, the *auxiliary path algorithm* [14] and the *Poisson path algorithm* [15], do not strictly belong to the three categories mentioned above. While they are closely related to the FD and PD methods, they provide unbiased estimators similar to the GT. We do not investigate these methods in this paper.

It has been observed that the PD method, when applicable, yields an estimator with lower variance than the GT estimator, which is applicable in most situations [5, 26]. In the context of stochastic chemical kinetics, the RPD method is applicable only to a limited class of examples and results in a biased estimator [26]. The FD methods also result in biased estimators. Both the FD and RPD methods also involve the use of method parameters, h or w , and the smaller

these are the less the bias of these methods. However, decreasing h or w results in an increase in the variance of the FD or RPD estimators, respectively. The GT estimator, on the other hand, is unbiased and does not involve method parameters to be determined. However, it has been observed that in many situations, the GT estimator has much larger variance compared to the FD and RPD estimators [5, 20, 24, 26]. To our knowledge, no theoretical explanation has been presented for the large variance of the GT method observed in many applications. In this paper, we provide a theoretical explanation for the large variance.

Remark 1. If a coefficient $c_j = 0$ in the stochastic mass action form of intensity functions, then reaction channel j is absent. However, one may want to compute the sensitivity at $c_j = 0$ to see the effect of “turning off” a reaction channel. In this case the GT or CGT method does not work; in fact the weight process Z is undefined. However, the FD methods work. Given the dependence of Z on c_j , one also expects the variance of Z to approach infinity as $c_j \rightarrow 0$. This was numerically examined in [14].

1.3. System size dependence in stochastic mass action. In stochastic chemical kinetics as well as other population models, there is a “system size parameter” N and in the $N \rightarrow \infty$ these systems behave deterministically (see Chapter 11 of [9], for instance). Our analysis shows that the variance of the GT method grows much faster in N than the variances of the FD methods.

We describe the general stochastic mass action form of intensities that commonly arise in stochastic chemical kinetics [10] and describe how system size N enters into the model. If we divide the stoichiometric vector ν_j into two parts, such that $\nu_j = \nu'_j - \nu''_j$, where

ν'_j = the vector number of molecules of each species that are created in the j th reaction,
 ν''_j = the vector number of molecules of each species that are consumed in the j th reaction,
 then the intensity of the j th reaction is

$$(18) \quad a_j^N(x, c) = \frac{c_j}{N^{|\nu''_j|-1}} \prod_i^n \binom{x_i}{\nu''_{ij}},$$

where $|\nu''_j| = \sum_{i=1}^n \nu''_{ij}$, N is the volume of the system times Avogadro’s number, and c_j is a constant specifying the rate of the reaction. We note that the term $\binom{x_i}{\nu''_{ij}}$ represents the number of ways to choose ν''_{ij} molecules from x_i molecules of the i th species. The term $1/N^{|\nu''_j|-1}$ also plays a critical role. To understand this, let us return to the example in (1). Let us relabel the parameters as c'_1 and c'_2 . As $c'_1 h + o(h)$ is the probability that a given pair of S_1 and S_2 interact during $(t, t+h]$, one expects c'_1 to depend on the system volume or equivalently on system size N in inverse proportion: $c'_1 = c_1/N$. Here, the newly defined c_1 is independent of system size N . On the other hand, for the monomolecular reaction, the probability $c'_2 h + o(h)$ of a given S_3 molecule reacting during $(t, t+h]$ is independent of system size N . In general, when $|\nu''_j|$ number of molecules come together to react, the term c'_j will depend on system size N as

$$(19) \quad c'_j = c_j / N^{|\nu''_j|-1}.$$

See [10] for more details. It must be noted that it is often useful to model “pure production” reactions, represented by an abstract chemical equation as $\emptyset \rightarrow S$, and the stochastic chemical

models in literature often utilize such reactions. In this case, the stochastic mass action form of intensity function is a constant c' and it is natural to take its dependence on N to be proportional: $c' = cN$, still satisfying the formula $c'_j = c_j/N^{|\nu''_j|-1}$.

Thus the intensity functions a_j^N depend on N and x in a specific manner, referred to as *density dependence* (see Chapter 11 of [9]). This density dependence leads to a deterministic limiting behavior in the large system size ($N \rightarrow \infty$) when the initial conditions are also scaled by N so that the initial species counts per volume (concentration) is held constant. The relevant theorem from [9] will be restated in the next section.

The parameters c'_j and c_j . We note that the parameters c'_j (which depend on N) are sometimes referred to as the *stochastic parameters*, while c_j are referred to as the *deterministic parameters*. In practice, one works with c'_j , and hence the sensitivities with respect to c'_j will be relevant. The sensitivities with respect to c_j are related to those with respect to c'_j via

$$(20) \quad S_j(t, c) = \frac{\partial}{\partial c_j} \mathbb{E}(f(X(t))) = \frac{\partial}{\partial c'_j} \mathbb{E}(f(X(t))) N^{1-|\nu''_j|} = S'_j(t, c) N^{1-|\nu''_j|}.$$

Moreover, if S is a sensitivity estimator for the sensitivity with respect to the deterministic parameter c_j , then $S' = S N^{|\nu''_j|-1}$ is a sensitivity estimator for the sensitivity with respect to the stochastic parameter c'_j . While the variances and biases of the stochastic and deterministic sensitivity estimators scale differently with system size N , the relative quantities RE, RSD, and RB will scale the same way. Therefore, without loss of generality, in the rest of the paper, we shall only concern ourselves with sensitivities with respect to the deterministic parameters c_j .

Finally, we note that in the stochastic mass action form of intensity functions, there is precisely one (deterministic) parameter c_j for each intensity function a_j and the parameters enter multiplicatively. Hence $\frac{\partial a_j}{\partial c_j} / a_j = 1/c_j$. This leads to the simple form for the weight process $Z(t, c)$ for the sensitivity with respect to c_j

$$(21) \quad Z(t, c) = \frac{1}{c_j} \left(R_j(t, c) - \int_0^t a_j(X(s, c)) ds \right).$$

1.4. An illustrative example. To investigate the estimator variance for the GT, CGT, and FD methods, we consider the analytically tractable birth death model from population dynamics, which also appears in gene regulatory networks where mRNA is produced at a constant probabilistic rate and decays at a rate proportional to the number of mRNA. The model is described by



The intensity functions are $a_1^N(x, c) = Nc_1$ and $a_2^N(x, c) = c_2x$. We consider the output function $f(x) = x$. Denoting by X^N the system size dependence of the process, it can be shown that

$$(23) \quad \mathbb{E}(X^N(t, c)) = Nx_0 e^{-c_2 t} + \frac{Nc_1}{c_2} (1 - e^{-c_2 t}),$$

where we have chosen a deterministic initial condition $X^N(0) = Nx_0$. The sensitivities with respect to c_1 and c_2 are given by

$$\begin{aligned}\frac{\partial}{\partial c_1} \mathbb{E}(X^N(t, c)) &= \frac{N}{c_2} (1 - e^{-c_2 t}), \\ \frac{\partial}{\partial c_2} \mathbb{E}(X^N(t, c)) &= -Nx_0 t e^{-c_2 t} - \frac{Nc_1}{c_2^2} (1 - e^{-c_2 t}) + \frac{Nc_1}{c_2} t e^{-c_2 t}.\end{aligned}$$

We observe that the both sensitivities are $\mathcal{O}(N)$ as $N \rightarrow \infty$. Also, in terms of t both sensitivities are $\mathcal{O}(1)$ as $t \rightarrow \infty$.

To study the variance of the GT and CGT estimators, first we consider the sensitivity $\frac{\partial}{\partial c_1} \mathbb{E}(X^N(t, c))$. The population process $X^N(t, c)$ and the weight process $Z^N(t, c)$ in this case can be written as

$$\begin{aligned}(24) \quad X^N(t, c) &= Nx_0 - \int_{(0, t]} dR_1^N(s, c) + \int_{(0, t]} dR_2^N(s, c), \\ Z^N(t, c) &= \int_{(0, t]} \frac{1}{c_1} dR_1^N(s, c) - N \int_0^t ds,\end{aligned}$$

where R_j^N and Z^N show dependence on N . One can use the Ito formula for processes driven by finite variation processes (see [25]) to write down the stochastic equations for $(X^N)^\alpha(t, c)(Z^N)^\beta(t, c)$, for the integer powers $0 \leq \alpha, \beta \leq 2$, and then take expectations to obtain a coupled system of linear ODEs for $\mathbb{E}((X^N)^\alpha(t, c)(Z^N)^\beta(t, c))$. Then the variance of GT and CGT estimators can be computed by the relations (17) with $f(x) = x$.

After lengthy calculations with the aid of Maple symbolic software one can show that

$$\begin{aligned}(25) \quad \text{Var}(S_{\text{GT}}) &= \frac{Ne^{-2c_2 t}}{c_1 c_2^2} (e^{2c_2 t} N^2 c_1^2 t + Nc_1 t c_2 e^{2c_2 t} + 2e^{c_2 t} N^2 c_1 c_2 t x_0 + e^{c_2 t} c_2^2 t N x_0 \\ &\quad + c_2^2 t N^2 x_0^2 - 2e^{c_2 t} N^2 c_1^2 t - e^{c_2 t} Nc_1 c_2 t - 2N^2 c_1 t c_2 x_0 \\ &\quad - Nx_0 t c_2^2 + 3Nc_1 e^{2c_2 t} + e^{2c_2 t} c_2 + 2Nx_0 e^{c_2 t} c_2 \\ &\quad + N^2 c_1^2 t - 6e^{c_2 t} Nc_1 - e^{c_2 t} c_2 - 2Nx_0 c_2 + 3Nc_1)\end{aligned}$$

and

$$\begin{aligned}(26) \quad \text{Var}(S_{\text{CGT}}) &= \frac{Ne^{-2c_2 t}}{c_1 c_2^2} (Nc_1 t c_2 e^{2c_2 t} + e^{c_2 t} c_2^2 t N x_0 - e^{c_2 t} Nc_1 c_2 t - Nx_0 t c_2^2 + Nc_1 e^{2c_2 t} \\ &\quad + e^{2c_2 t} c_2 - 2e^{c_2 t} Nc_1 - e^{c_2 t} c_2 + Nc_1).\end{aligned}$$

We observe that the variance of the GT estimator is $\mathcal{O}(N^3)$ while that of the CGT estimator is $\mathcal{O}(N^2)$, as $N \rightarrow \infty$. On the other hand, both estimators have $\mathcal{O}(t)$ variance as $t \rightarrow \infty$. Hence, in the $N \rightarrow \infty$ limit, the RSD of the GT estimator is $\mathcal{O}(N^{1/2})$ and the RSD of the CGT estimator is $\mathcal{O}(1)$. We can also conclude that in the $t \rightarrow \infty$ limit, the RSD is $\mathcal{O}(\sqrt{t})$ for both methods.

Second, we consider the sensitivity $\frac{\partial}{\partial c_2} \mathbb{E}(X^N(t, c))$. The weight process $Z^N(t, c)$ in this case can be written as

$$(27) \quad Z^N(t, c) = \int_{(0,t]} \frac{1}{c_2} dR_2^N(s, c) - \int_0^t X^N(s, c) ds,$$

and the analysis, while possible, is more complicated. For simplicity, we choose $c_1 = 0$, so the process now corresponds to a pure death process. In this case, the variances of GT and CGT estimators can be shown to be

$$(28) \quad \begin{aligned} \text{Var}(S_{\text{GT}}) = & \frac{1}{c_2^2} (e^{-2c_2t} N^3 x_0^3 - 4e^{-2c_2t} N^2 x_0^2 + 3e^{-2c_2t} N x_0 + 3e^{-2c_2t} N^2 x_0^2 t^2 c_2^2 \\ & - 2e^{-3c_2t} N x_0 + 3e^{-3c_2t} N^2 x_0^2 + e^{-c_2t} N^2 x_0^2 - e^{-c_2t} N x_0 \\ & + e^{-c_2t} N x_0 t^2 c_2^2 - 4e^{-2c_2t} t^2 c_2^2 N x_0 - e^{-3c_2t} N^2 x_0^3) \end{aligned}$$

and

$$(29) \quad \begin{aligned} \text{Var}(S_{\text{CGT}}) = & \frac{1}{c_2^2} (-2e^{-2c_2t} N^2 x_0^2 + 3e^{-2c_2t} N x_0 + e^{-2c_2t} N^2 x_0^2 t^2 c_2^2 \\ & - 2e^{-3c_2t} N x_0 + e^{-3c_2t} N^2 x_0^2 + e^{-c_2t} N^2 x_0^2 \\ & - e^{-c_2t} N x_0 + e^{-c_2t} N x_0 t^2 c_2^2 - 4e^{-2c_2t} N x_0 t^2 c_2^2). \end{aligned}$$

When dependence on system size N is concerned, the variance of GT estimator is $\mathcal{O}(N^3)$ while that of the CGT estimator is only $\mathcal{O}(N^2)$. As in the case of the parameter c_1 , we again obtain that the RSD of the GT method is $\mathcal{O}(N^{1/2})$ while that of CGT is $\mathcal{O}(1)$, as $N \rightarrow \infty$. Finally, we note that large t behavior is uninteresting as the system enters the absorbing state 0 eventually.

Now we consider any FD estimator, and we can bound its variance as

$$(30) \quad \begin{aligned} \text{Var}(S_{\text{FD}}) &= h^{-2} \text{Var}(X^N(t, c+h) - X^N(t, c)) \\ &\leq 2h^{-2} \{ \text{Var}(X^N(t, c+h)) + \text{Var}(X^N(t, c)) \}. \end{aligned}$$

We also note that [23]

$$(31) \quad \text{Var}(X^N(t, c)) = N x_0 (1 - e^{-c_2t}) e^{-c_2t} + \frac{N c_1}{c_2} (1 - e^{-c_2t}).$$

In our analysis we shall treat the finite difference perturbation h of the parameter as independent of system size N so that we consider the variance and bias of the FD estimator as a function of the two variables h and N . From the above equation, we see that for any fixed h , the variance of an FD estimator is $\mathcal{O}(N)$ and hence the RSD of the FD estimator is $\mathcal{O}(N^{-1/2})$ as $N \rightarrow \infty$. Finally, we note that for fixed N , as $t \rightarrow \infty$, the variance of the FD estimator is $\mathcal{O}(1)$.

We note here that the above upper bound for $\text{Var}(S_{\text{FD}})$ is exactly twice the variance of the IRN FD method. If a CRN FD method is used, the variance is in general much smaller. Nevertheless, our numerical results show that the asymptotic order in N is sharp even for CRN.

From the expression for $\mathbb{E}(X^N(t, c))$ in (23) it can easily be shown that the RB defined by (8) of any FD method is $\mathcal{O}(1)$ as $N \rightarrow \infty$ (with h fixed) when sensitivity of $\mathbb{E}(X^N(t, c))$ with respect to c_1 or c_2 is considered.

To summarize, we note that when computing the sensitivity of $\mathbb{E}(X^N(t, c))$ in this example, with respect to either of the parameters c_1 or c_2 , we observe that the RSDs of the GT, CGT, and FD estimators scale with system size N as $\mathcal{O}(N^{1/2})$, $\mathcal{O}(1)$, and $\mathcal{O}(N^{-1/2})$, respectively. If N is modestly large (say, 10–100), a significant amount of reduction in the RSD can be expected using CGT over GT. On the other hand FD methods will have even lower variance when compared to both GT and CGT as system size increases. However, the FD methods are biased, and for fixed h the RB remains $\mathcal{O}(1)$ as $N \rightarrow \infty$.

1.5. Contributions of this paper. Our analysis will show that the observations made about the RSD and the RB of the GT, CGT, and FD estimators in the context of the particular example of the previous subsection generalize to a large class of stochastic reaction networks. These general results are provided in section 4. While our analysis does not apply to the RPD method, our numerical simulations show that RPD has system size dependence similar to the FD methods. While our RSD analysis in the cases of CGT and CRN FD estimators is not proven to be sharp, the numerical simulations show that the estimates in terms of large system size N are sharp.

Our analysis thus provides theoretical evidence that centering (to obtain the CGT method) significantly improves the efficiency of the GT methods. Since the FD methods are biased while the GT and CGT methods are not, an efficiency comparison must be based on variance and bias. In the case of the FD estimators which depend on system size N as well as the perturbation parameter h , our analysis in section 4 treats h and N as independent variables and provides the large N behavior for fixed h . The small h behavior of the FD methods (for fixed N) is well known [5]. In section 6, we combine our large N results with the existing small h results for the FD methods in order to decide the optimal choice of h as a function of N , and we provide an estimate of efficiency (as measured by the number N_s of trajectories needed to achieve a given value δ for the RE) of the GT, CGT, and FD methods.

2. General setup and running assumptions. As mentioned in the previous section, the system size shall be the key to our analytical explanation for the larger variance of the GT estimator. In this section we set the stage for the system size analysis and state some assumptions that shall be carried throughout the rest of the paper. We shall use the notation $|x|$ for the norm of a vector (any norm in \mathbb{R}^n would do) and $\|\nu\|$ for the corresponding induced norm of a matrix.

Remark 2. Our analysis will focus on processes X , R , and Z corresponding to different system sizes N ; however, the deterministic parameter value c is fixed at a specific value at which the sensitivity is sought. For notational ease and readability, we shall not show the dependence of these processes and intensity functions on c and display c only when it explicitly appears outside these.

We will study the family of processes X^N indexed by $N \geq 1$ corresponding to the family of intensity functions a_j^N that are represented on the same sample space via the stochastic equation

$$(32) \quad X^N(t) = Nx_0 + \sum_{j=1}^m Y_j \left(\int_0^t a_j^N(X^N(s)) ds \right) \nu_j, \quad N \geq 1,$$

where Y_j are independent unit rate Poisson processes and we have taken $X^N(0) = Nx_0$, where $x_0 \in \mathbb{R}_+^n$ is fixed (deterministic). We also define the corresponding family of vector reaction count processes $R^N(t)$ whose j th component $R_j^N(t)$ counts the number of reaction events of type j that occurred during $(0, t]$. Thus

$$R_j^N(t) = Y_j \left(\int_0^t a_j^N(X^N(s)) ds \right), \quad N \geq 1, \quad j = 1, \dots, m.$$

We also define the centered processes $M^N(t) = (M_1^N(t), \dots, M_m^N(t))$ by

$$M_j^N(t) = R_j^N(t) - \int_0^t a_j^N(X^N(s)) ds, \quad N \geq 1, \quad j = 1, \dots, m.$$

We shall state five running assumptions under which the rest of the analysis in this paper is carried out. We note that Assumptions 1–3 are assumptions on the intensity functions and their dependence on parameters and system size. These assumptions are satisfied by the stochastic mass action form of intensity functions and are intended to generalize certain key properties of the stochastic mass action form of intensity functions. Not all stochastic models of intensity functions in the literature follow the stochastic mass action form. In such cases, our analysis will still apply provided these assumptions are met.

Assumption 1. We assume the following form of parameter dependence on the intensity function. For each $j = 1, \dots, m$ and $N \geq 1$,

$$(33) \quad a_j^N(x, c) = c_j b_j^N(x),$$

where $b_j^N : \mathbb{R}^n \rightarrow \mathbb{R}$ are such that b_j^N restricted to \mathbb{Z}_+^n are nonnegative. This also implies that there are precisely m parameters, one for each reaction j .

For the analysis in this paper we need not assume the stochastic mass action form but merely the density dependence which is stated by our Assumption 2.

Assumption 2. We suppose that for each $j = 1, \dots, m$ and each $x \in \mathbb{R}_+^n$, the limit $\lim_{N \rightarrow \infty} a_j^N(Nx)/N = a_j(x)$ exists and, moreover, for each compact $K \subset \mathbb{R}_+^n$, the collection of functions $a_j^N(Nx) - Na_j(x)$ is uniformly bounded for $x \in K$ and $N \geq 1$. We note that this implies that for each compact set $K \subset \mathbb{R}_+^n$ there exists a constant $B_K > 0$ such that

$$(34) \quad \left| \frac{a_j^N(Nx)}{N} - a_j(x) \right| \leq \frac{B_K}{N}, \quad x \in K, \quad j = 1, \dots, m, \quad N \geq 1.$$

Defining $X_N(t) = N^{-1}X^N(t)$, we note that X_N can be interpreted as the *concentration* of molecules at time t for system size N . We note that X_N are coupled via the following stochastic equations:

$$(35) \quad X_N(t) = x_0 + \sum_{j=1}^m N^{-1} Y_j \left(\int_0^t a_j^N(NX_N(s)) ds \right) \nu_j.$$

We state the following theorem regarding the limiting behavior of X_N (see [9] for details). The deterministic limit X of X_N is also referred to as the *fluid limit*.

Theorem 1 (Theorem 2.1 of Chapter 11 in [9]). *Suppose that Assumption 2 holds. Moreover, assume that for each compact $K \subset \mathbb{R}^n$,*

$$\sum_{j=1}^m |\nu_j| \sup_{x \in K} a_j(x) < \infty$$

and that $F(x) = \sum_{j=1}^m \nu_j a_j(x)$ is Lipschitz on K , that is, for each $x, y \in K$, there exists some constant M_K such that

$$|F(x) - F(y)| \leq M_K |x - y|.$$

Suppose $t > 0$ is in the forward maximal interval of existence of solution X for the ODE initial value problem

$$X(t) = x_0 + \int_0^t F(X(s)) ds.$$

Then

$$\limsup_N \sup_{s \leq t} |X_N(s) - X(s)| = 0 \quad \text{a.s.,}$$

where the deterministic limit X satisfies the ODE above.

Remark 3. We note that with fixed initial condition $X_N(0) = x_0$ we want $X^N(0) = Nx_0$ to belong to \mathbb{Z}_+^n , which may not hold for all $N \geq 1$ but we assume that it holds for a sequence of N values tending to ∞ . For instance, if x_0 is rational this is true. This is adequate for our purposes.

In order to satisfy the conditions stated in Theorem 1 we shall assume the following.

Assumption 3. For each $j = 1, \dots, m$, the functions $a_j(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ are continuously differentiable. This automatically implies the Lipschitz condition in Theorem 1.

The following assumption is used to facilitate the analysis in this paper. Several, but not all, examples in applications satisfy this assumption.

Assumption 4. We assume that the sequence of concentration processes X_N is uniformly bounded, that is, there exists a constant Γ such that for all $t \geq 0$,

$$(36) \quad |X_N(t)| \leq \Gamma \quad \text{a.s.}$$

for all $N \geq 1$.

We note that if there exists a strictly positive vector $\gamma \in \mathbb{R}_+^m$ so that $\gamma^T \nu_j \leq 0$ for each j , then this assumption is satisfied. We note that a form of converse of this statement is also true [22].

Now we turn our attention to the sensitivity. Given $f : \mathbb{R}^n \rightarrow \mathbb{R}$, we are interested in computing the sensitivity

$$\frac{\partial}{\partial c} \mathbb{E}(f(X^N(t))),$$

where $c \in (0, \infty)$ is a parameter. In view of Assumption 1, without loss of generality, we shall take $c = c_1$. Then we note that the GT sensitivity estimator is $f(X^N(t))Z^N(t)$ and the CGT estimator is $[f(X^N(t)) - \mathbb{E}(f(X^N(t)))]Z^N(t)$, where we note that $Z^N(t) = M_1^N(t)/c_1$ in this case.

As we are concerned with families of processes indexed by N , it makes sense to consider a corresponding family of functions $f^N : \mathbb{R}^n \rightarrow \mathbb{R}$ instead of one function f and make reasonable assumptions on f^N and f .

To motivate the assumption we make on f^N and f , we note that we shall be concerned with $f^N(X^N(t)) = f^N(NX_N(t))$ which we wish to compare with $f(X(t))$. When $f^N(x) = x_i$, one of the components of x , we have

$$f^N(NX_N(t))/N = X_{N_i}(t) \rightarrow X_i(t) = f(X(t))$$

with $f(x) = x_i$. Alternatively, if $f^N(x) = x_i^\alpha$ for some $\alpha > 0$ we have

$$f^N(NX_N(t))/N^\alpha = (X_{N_i}(t))^\alpha \rightarrow (X_i(t))^\alpha = f(X(t))$$

with $f(x) = x_i^\alpha$. If, however, $f^N(x) = x_i^2 + x_i$, then we have

$$f^N(NX_N(t))/N^2 = (X_{N_i}(t))^2 + X_{N_i}(t)/N \rightarrow (X_i(t))^2 = f(X(t)),$$

where $f(x) = x_i^2$. In this case we note that $f^N(Nx)/N^2 - f(x) = x_i/N$, which tends to 0 as $1/N$, uniformly for x in a compact set. Motivated by this, we impose the following assumption.

Assumption 5. We assume that there exist a function f and a constant $\alpha > 0$ such that for each compact set $K \subset \mathbb{R}_+^n$,

$$(37) \quad |f^N(Nx)/N^\alpha - f(x)| \leq \frac{L_K}{\sqrt{N}}, \quad x \in K, \quad N \geq 1,$$

for some constant $L_K > 0$.

We remark that the $\mathcal{O}(1/\sqrt{N})$ behavior is adequate for our proofs.

We note that the running assumptions 1–5 will be assumed throughout the rest of the paper.

3. Large N behavior. In this section we derive results concerning the $N \rightarrow \infty$ limit for the various relevant processes. Throughout the rest of the paper $X(t)$ will denote the solution of the equation

$$(38) \quad X(t) = x_0 + \sum_{j=1}^m \nu_j \int_0^t a_j(X(s)) ds,$$

where $x_0 \in \mathbb{R}_+^n$ is fixed.

Lemma 2. For each $j = 1, \dots, m$, there exists $A_j > 0$ such that, for all $t > 0$,

$$\frac{a_j^N(NX_N(t))}{N} \leq A_j \quad a.s.$$

for all $N \geq 1$.

Proof. By Assumption 4, the processes X_N are contained in a compact set of \mathbb{R}^n , say, K ; therefore, for each j we have the estimation

$$\sup_{t \geq 0} \frac{a_j^N(NX_N(t))}{N} \leq \sup_{x \in K} \frac{a_j^N(Nx)}{N}.$$

Since $N^{-1}a_j^N(Nx)$ converges uniformly to $a_j(x)$ for $x \in K$ by (34) in Assumption 2, it is apparent that $\sup_{x \in K} N^{-1}a_j^N(Nx)$ is bounded by continuity of a_j . Hence $\sup_{t \geq 0} N^{-1}a_j^N(NX_N(t))$ is bounded by a constant A_j . ■

Lemma 3. For each $j = 1, \dots, m$, and $t > 0$, we have

$$\sup_{s \leq t} \left| \frac{a_j^N(NX_N(s))}{N} - a_j(X(s)) \right| \rightarrow 0 \quad a.s.$$

as $N \rightarrow \infty$.

Proof. We may write

$$\begin{aligned} & \left| \frac{a_j^N(NX_N(s))}{N} - a_j(X(s)) \right| \\ & \leq \left| \frac{a_j^N(NX_N(s))}{N} - a_j(X_N(s)) \right| + |a_j(X_N(s)) - a_j(X(s))|. \end{aligned}$$

The first part on the right-hand side converges to zero uniformly for s in $[0, t]$ because of Assumptions 2 and 4. To see that the second part on the right-hand side converges uniformly to 0 on $[0, t]$, note that by Assumptions 3 and 4, a_j is Lipschitz continuous on the compact set K (which contains X_N and X); hence the result follows by Theorem 1. ■

We define a family of scaled reaction count processes $R_N(t)$ by $R_N(t) = R^N(t)/N$.

Lemma 4. For each $j = 1, 2, \dots, m$ and $t > 0$,

$$\sup_{s \leq t} \left| R_{Nj}(s) - \int_0^s a_j(X(u)) du \right| \rightarrow 0 \quad a.s.$$

as $N \rightarrow \infty$.

Proof. Recall that $R_j^N(t) = Y_j(\int_0^t a_j^N(NX_N(s)) ds)$. For each $j = 1, \dots, m$,

$$\begin{aligned} & \sup_{s \leq t} \left| \frac{1}{N} Y_j \left(\int_0^s a_j^N(NX_N(u)) du \right) - \int_0^s a_j(X(u)) du \right| \\ & \leq \sup_{s \leq t} \left| \frac{1}{N} Y_j \left(\int_0^s a_j^N(NX_N(u)) du \right) - \frac{1}{N} \int_0^s a_j^N(NX_N(u)) du \right| \\ & \quad + \int_0^t \left| \frac{1}{N} a_j^N(NX_N(u)) - a_j(X(u)) \right| du. \end{aligned}$$

The second term on the right-hand side converges to zero by Lemma 3. Setting $\tilde{Y}(t) = Y(t) - t$, the first term on the right can be written and then bounded as

$$\sup_{s \leq t} \left| \frac{1}{N} \tilde{Y}_j \left(\int_0^s a_j^N(NX_N(u)) du \right) \right| \leq \sup_{s \leq t} \left| \frac{1}{N} \tilde{Y}_j(NA_j s) \right| \quad \text{a.s.},$$

where the last term converges to zero by the law of large numbers for Poisson processes (see Theorem 1.2 in [3]). ■

Lemma 5. For a given $t > 0$, suppose that f is continuous at $X(t)$. Then

$$(39) \quad \lim_{N \rightarrow \infty} |f^N(NX_N(t))/N^\alpha - f(X(t))| = 0 \quad \text{a.s.}$$

Proof. Write

$$\begin{aligned} |f^N(NX_N(t))/N^\alpha - f(X(t))| &\leq |f^N(NX_N(t))/N^\alpha - f(X_N(t))| \\ &\quad + |f(X_N(t)) - f(X(t))|. \end{aligned}$$

The first term converges to zero almost surely by Assumption 4 and (37) in Assumption 5. The second term converges to zero by the continuity assumption on f since $X_N(t)$ converges to $X(t)$ almost surely. ■

Recall the definition of M^N ,

$$M^N(t) = R^N(t) - \int_0^t a^N(NX_N(s)) ds.$$

Note that in general, $M^N(t)$ is an m -dimensional local martingale (see [21, 16] for a definition) for each N , but by Lemma 2 it follows that $\mathbb{E}[R_j^N(t)] \leq NA_j t$ for all $t > 0$, which makes $M^N(t)$ a martingale. We define the scaled processes $M_N = N^{-1}M^N$ and $Z_N = N^{-1}Z^N$. We note that $Z^N(t) = M_1^N(t)/c_1$ and $Z_N(t) = M_{N1}(t)/c_1$.

Let us denote by $D^m[0, \infty)$ the space of càdlàg functions mapping from $[0, \infty)$ into \mathbb{R}^m , endowed with the Skorohod topology (see [7] for definitions). We provide a lemma on the weak convergence of M_N .

Lemma 6. Let $C(t) = (c_{ij}(t))$ be the $m \times m$ matrix-valued function, where

$$(40) \quad c_{ij}(t) = \begin{cases} \int_0^t a_j(X(s)) ds, & i = j, \\ 0, & i \neq j. \end{cases}$$

Then $\sqrt{N}M_N \Rightarrow \bar{M}$ on $D^m[0, \infty)$, where $\bar{M}(t)$ is an m -dimensional Gaussian process with independent increments, having mean vector and covariance matrix

$$(41) \quad \mathbb{E}[\bar{M}(t)] = (0, \dots, 0), \quad \mathbb{E}[\bar{M}(t)\bar{M}(t)^T] = C(t).$$

In particular, the scaled Girsanov sensitivity (or weight) process $\sqrt{N}Z_N \Rightarrow U$ on $D[0, \infty)$, where

$$(42) \quad U(t) = \frac{1}{c_1} \bar{M}_1(t).$$

Also since U has continuous sample paths, for each $t > 0$, we have

$$\sqrt{N}Z_N(t) \Rightarrow U(t).$$

Proof. The proof relies on the martingale functional central limit theorem (FCLT) proved in [29]. Note that each jump of $\sqrt{N}M_N$ has size $1/\sqrt{N}$; therefore,

$$\lim_{N \rightarrow \infty} \mathbb{E} \left[\sup_{s \leq t} \left| \sqrt{N}M_N(s) - \sqrt{N}M_N(s-) \right| \right] = 0.$$

Also, for each pair (i, j) with $i, j = 1, \dots, m$, and each $t > 0$, since the jump size for M_{Nj} is always N^{-1} and there are no simultaneous jumps, we have the following quadratic covariation:

$$(43) \quad \left[\sqrt{N}M_{Ni}, \sqrt{N}M_{Nj} \right] (t) = \begin{cases} R_{Nj}, & i = j, \\ 0, & i \neq j. \end{cases}$$

By Lemma 4, $R_{Nj}(t)$ converges almost surely to $c_{jj}(t) = \int_0^t a_j(X(s))ds$. Then, for each pair (i, j) ,

$$\left[\sqrt{N}M_{Ni}, \sqrt{N}M_{Nj} \right] (t) \rightarrow c_{ij}(t)$$

almost surely and hence in probability. Thus, the weak convergence of M_N follows from the martingale FCLT. ■

Lemma 7. For each $p \geq 1$, there exists a constant $\beta(p)$ such that for all $t > 0$

$$(44) \quad \limsup_N \mathbb{E} \left(\sup_{s \leq t} \left| \sqrt{N}M_N(s) \right| \right)^p \leq \beta(p)t^{p/2}.$$

Proof. Observe that the quadratic variation (see [21] for a definition) of $\sqrt{N}M_N$ is

$$\left[\sqrt{N}M_N, \sqrt{N}M_N \right] (t) = N^{-1} \sum_{j=1}^m Y_j \left(\int_0^t a_j^N(NX_N(s))ds \right).$$

By the Burkholder–Davis–Gundy inequality (see [21]), there exists a constant $C(p)$ (depends on p) such that

$$\begin{aligned} \mathbb{E} \left(\sup_{s \leq t} \left| \sqrt{N}M_N(s) \right| \right)^p &\leq C(p) \mathbb{E} \left(\frac{1}{N} \sum_{j=1}^m Y_j \left(\int_0^t a_j^N(NX_N(s))ds \right) \right)^{p/2} \\ &\leq C(p) \mathbb{E} \left(\frac{1}{N} \sum_{j=1}^m Y_j (NA_j t) \right)^{p/2} \\ &\leq C(p) N^{-p/2} \left(\mathbb{E} \left(\sum_{j=1}^m Y_j (NA_j t) \right)^p \right)^{1/2}, \end{aligned}$$

where we have used Lemma 2.

Hence,

$$\limsup_N \mathbb{E} \left(\sup_{s \leq t} \left| \sqrt{N} M_N(s) \right| \right)^p \leq \limsup_N C(p) N^{-p/2} \left(\mathbb{E} \left(\sum_{j=1}^m Y_j(N A_j t) \right)^p \right)^{1/2}.$$

First we observe that for $j = 1, \dots, m$, the p th moment of the Poisson random variable $Y_j(N A_j t)$ is a polynomial of degree p in $N A_j t$. Also, noting that Y_j are independent, we obtain that the right-hand side is bounded by a term $\beta(p) t^{p/2}$, where $\beta(p)$ is a constant. ■

Since $Z^N(t) = c_1^{-1} M_1^N(t)$, we immediately have the following property regarding the process Z_N .

Lemma 8. *For each $p \geq 1$, there exists a constant $\gamma(p)$ such that for all $t > 0$,*

$$(45) \quad \limsup_N \mathbb{E} \left(\sup_{s \leq t} \sqrt{N} |Z_N(s)| \right)^p \leq \gamma(p) t^{p/2}.$$

Define the process $V_N(t) = \sqrt{N}(X_N(t) - X(t))$. Let us consider the moment of this process on a compact time interval.

Lemma 9. *For each $p \geq 1$, there exist constants $\bar{\beta}(p), K(p)$ such that for all $t > 0$*

$$\limsup_N \sup_{s \leq t} \mathbb{E} (|V_N(s)|^p) \leq \bar{\beta}(p) t^{p/2} e^{K(p)t^p}.$$

Proof. Recall that

$$X_N(s) = x_0 + \nu R_N(s)$$

and

$$X(s) = x_0 + \int_0^s \nu a(X(u)) du,$$

where ν is the n by m dimensional stoichiometric matrix. One can write V_N as

$$\begin{aligned} V_N(s) &= \sqrt{N} \nu R_N(s) - \sqrt{N} \int_0^s \nu a(X(u)) du \\ &= \sqrt{N} \nu \left(R_N(s) - \int_0^s \frac{a^N(N X_N(u))}{N} du \right) \\ &\quad + \sqrt{N} \nu \left(\int_0^s \frac{a^N(N X_N(u))}{N} - a(X(u)) du \right). \end{aligned}$$

Note that we denote $M_N(s) = R_N(s) - \int_0^s N^{-1} a^N(N X_N(u)) du$, and hence

$$|V_N(s)| \leq \|\nu\| \left| \sqrt{N} M_N(s) \right| + \|\nu\| \int_0^s \sqrt{N} \left| \frac{a^N(N X_N(u))}{N} - a(X(u)) \right| du.$$

To estimate the second term on the right-hand side of the last inequality, we note that

$$\begin{aligned} \sqrt{N} \left| \frac{a^N(N X_N(u))}{N} - a(X(u)) \right| &\leq \sqrt{N} \left| \frac{a^N(N X_N(u))}{N} - a(X_N(u)) \right| \\ &\quad + \sqrt{N} |a(X_N(u)) - a(X(u))|. \end{aligned}$$

Since X_N lies in a compact set K according to Assumption 4, we have for all $u > 0$,

$$\left| \frac{a^N(NX_N(u))}{N} - a(X_N(u)) \right| \leq \frac{\tilde{B}_K}{N},$$

where we have used Assumption 2 and \tilde{B}_K is related to B_K from (34).

On the other hand, for each $j = 1, \dots, m$, by Assumption 3, a_j is continuously differentiable and hence it is Lipschitz continuous on the compact set K . Hence, there exists a Lipschitz constant C_j such that for all $u > 0$,

$$|a_j(X_N(u)) - a_j(X(u))| \leq C_j |X_N(u) - X(u)|.$$

It follows that there exists a constant C such that

$$|a(X_N(u)) - a(X(u))| \leq C |X_N(u) - X(u)|,$$

where $\|\cdot\|$ can be norm on \mathbb{R}^m . Therefore,

$$\begin{aligned} |V_N(s)| &\leq \|\nu\| \left(\left| \sqrt{N} M_N(s) \right| + N^{-1/2} \tilde{B}_K s + C \int_0^s \sqrt{N} |X_N(u) - X(u)| du \right) \\ &= \|\nu\| \left(\left| \sqrt{N} M_N(s) \right| + N^{-1/2} \tilde{B}_K s + C \int_0^s |V_N(u)| du \right). \end{aligned}$$

In virtue of the inequality $(a + b + c)^p \leq 3^p(a^p + b^p + c^p)$ and the Holder's inequality, we obtain

$$|V_N(s)|^p \leq (3\|\nu\|)^p \left(\left| \sqrt{N} M_N(s) \right|^p + N^{-p/2} (\tilde{B}_K s)^p + C^p s^{p-1} \int_0^s |V_N(u)|^p du \right).$$

Taking expected value of both sides, for $s \in [0, t]$,

$$\begin{aligned} \mathbb{E}|V_N(s)|^p &\leq (3\|\nu\|)^p \left(\mathbb{E} \left| \sqrt{N} M_N(s) \right|^p + N^{-p/2} (\tilde{B}_K t)^p \right) \\ &\quad + (3\|\nu\|)^p C^p s^{p-1} \left(\int_0^s \mathbb{E}|V_N(u)|^p du \right). \end{aligned}$$

To estimate the first term of the right-hand side, recall that in the proof of Lemma 7,

$$\mathbb{E} \left(\sup_{s \leq t} \left| \sqrt{N} M_N(s) \right| \right)^p \leq C(p) N^{-p/2} \left(\mathbb{E} \left(\sum_{j=1}^m Y_j(N A_j t) \right)^p \right)^{1/2}.$$

For convenience, let us denote

$$\Phi_N(t) = C(p) N^{-p/2} \left(\mathbb{E} \left(\sum_{j=1}^m Y_j(N A_j t) \right)^p \right)^{1/2}.$$

Therefore,

$$\mathbb{E}|V_N(s)|^p \leq (3\|\nu\|)^p \left(\Phi_N(t) + N^{-p/2}(\tilde{B}_K t)^p + C^p s^{p-1} \left(\int_0^s \mathbb{E}|V_N(u)|^p du \right) \right).$$

We note that $\mathbb{E}|V_N(s)|^p$ is continuous in s , and applying the Gronwall inequality, we obtain that, for $s \leq t$,

$$\mathbb{E}|V_N(s)|^p \leq (3\|\nu\|)^p \left(\Phi_N(t) + N^{-p/2}(\tilde{B}_K t)^p \right) e^{(3\|\nu\|)^p C^p t^p}.$$

Taking supremum over $s \in [0, t]$ and then taking \limsup_N , the result follows from the same considerations as in the proof of Lemma 7. ■

4. Scaling of sensitivity, estimator bias, and estimator variance. In this section, we study the system size dependence of the sensitivity

$$s^N = \frac{\partial}{\partial c} \mathbb{E}(f^N(X^N(t)))$$

and the biases as well as the variances of the GT, CGT, and FD estimators. In the case of the FD estimators, the parameter perturbation h is fixed when $N \rightarrow \infty$. As mentioned earlier, the difference between the sensitivity with respect to the stochastic parameter and with respect to the deterministic parameter is merely a scaling factor $N^{|\nu_j''|-1}$ and hence the RSD, RB, and RE are unchanged regardless of whether one considers the sensitivity with respect to the stochastic parameter or the deterministic parameter. From an analytical point of view, it is convenient to study the sensitivity with respect to the deterministic parameter.

Recall that the sensitivity estimator of the GT method is

$$f^N(X^N(t, c))Z^N(t, c),$$

where $f^N : \mathbb{R}^n \rightarrow \mathbb{R}$. We remind the reader that f^N satisfies Assumption 5, that is, there exist a function f and a constant α such that

$$\left| \frac{f^N(Nx)}{N^\alpha} - f(x) \right| \leq \frac{L_K}{\sqrt{N}}.$$

Theorem 10. *In addition to our running assumptions, we assume that f in (37) is continuously differentiable. Then for each $t \geq 0$*

$$\sup_{s \leq t} \mathbb{E}(f^N(X^N(s))Z^N(s)) = \mathcal{O}(N^\alpha).$$

That is, the true sensitivity is asymptotically $\mathcal{O}(N^\alpha)$ uniformly on $[0, t]$.

Proof. It is sufficient to show that $\sup_{s \leq t} \mathbb{E}(f^N(X^N(s))Z^N(s))/N^\alpha$ is bounded in N . Instead of working with $\mathbb{E}(f^N(X^N(s))Z^N(s))/N^\alpha$, we use

$$\mathbb{E} \left(\frac{f^N(X^N(s))}{N^\alpha} Z^N(s) - f(X(s))Z^N(s) \right)$$

because they are equal but the latter is easier to work with.

Note that f is continuously differentiable hence Lipschitz on the compact set K corresponding to Assumption 4. Denote by C_K the Lipschitz constant for f . Using the assumptions on f^N and f and writing X^N in terms of V_N as

$$X^N(s) = NX(s) + \sqrt{N}V_N(s),$$

which leads to

$$\begin{aligned} & \left| \frac{f^N(NX(s) + \sqrt{N}V_N(s))}{N^\alpha} - f(X(s)) \right| |Z^N(s)| \\ & \leq \left| \frac{f^N(NX(s) + \sqrt{N}V_N(s))}{N^\alpha} - f\left(X(s) + \frac{V_N(s)}{\sqrt{N}}\right) \right| |Z^N(s)| \\ & \quad + \left| f\left(X(s) + \frac{V_N(s)}{\sqrt{N}}\right) - f(X(s)) \right| |Z^N(s)| \\ & \leq \frac{L_K}{\sqrt{N}} |Z^N(s)| + C_K |V_N(s)| \frac{|Z^N(s)|}{\sqrt{N}} \\ & \leq L_K \sqrt{N} |Z_N(s)| + \frac{1}{2} C_K (|V_N(s)|^2 + N |Z_N(s)|^2), \end{aligned}$$

where L_K is as defined in Assumption 5. Taking expectation on both sides, the result follows from Lemmas 8 and 9. \blacksquare

Remark 4. While the proof above does not show that the order $\mathcal{O}(N^\alpha)$ is sharp, it can be shown to be sharp if, under the N^α scaling, the sensitivity of the stochastic process is shown to limit to the sensitivity of the deterministic limit $f(X(t))$ as $N \rightarrow \infty$. In fact, under additional assumptions, this limit can be shown [13]. Our numerical results in section 5 also show $\mathcal{O}(N^\alpha)$ behavior.

Recall that the FD estimator is defined in (10) as

$$S_{\text{FD}}^N(t, c) = h^{-1} [f^N(X^N(t, c+h)) - f^N(X^N(t, c))].$$

Based on the last theorem, with a little more effort we conclude the following corollary regarding the bias of FD estimator.

Corollary 11. *In addition to the running assumptions, if we assume that f is continuously differentiable, then for each $t > 0$, we have*

$$\mathbb{E}(S_{\text{FD}}^N(t) - \mathcal{S}^N(t)) = \mathcal{O}(N^\alpha),$$

where $\mathcal{S}^N(t)$ represents the true sensitivity at t . That is, the bias of FD estimator is asymptotically $\mathcal{O}(N^\alpha)$.

Proof. Since we have shown that the true sensitivity scales like $\mathcal{O}(N^\alpha)$, it suffices to show that $\mathbb{E}(f^N(X^N(t, c)))$ is asymptotically of order $\mathcal{O}(N^\alpha)$ for any c . In fact, by Lemma 5, $f^N(X^N(t))/N^\alpha$ converges almost surely to $f(X(t))$. To apply the dominate convergence theorem, note that Assumption 5 implies

$$\frac{|f^N(X^N(t))|}{N^\alpha} \leq |f(X(t))| + \frac{L_K}{\sqrt{N}}.$$

By virtue of Assumption 4, the right-hand side of the above equality is bounded in N and hence it is integrable. Finally, the dominate convergence theorem gives the result. ■

Next, we investigate the variance of the GT estimator in terms of the system size N . The following lemma concerning the weak convergence of joint distribution is crucial for the proof of Theorem 13.

Lemma 12. *Let X_n and Y_n be \mathbb{R}^m valued and \mathbb{R}^k valued sequences of random variables, respectively. Suppose X_n converges to X in probability (where X is deterministic) and $Y_n \Rightarrow Y$. Then $(X_n, Y_n) \Rightarrow (X, Y)$ in \mathbb{R}^{m+k} .*

Proof. Let $x \in \mathbb{R}^m$ be such that $X = x$ almost surely. First we show that $(X, Y_n) \Rightarrow (X, Y)$. If $f : \mathbb{R}^{m+k} \rightarrow \mathbb{R}$ is bounded and continuous, then so is $g : \mathbb{R}^k \rightarrow \mathbb{R}$ defined by $g(y) = f(x, y)$. Since $Y_n \Rightarrow Y$ we have that

$$\mathbb{E}(f(X, Y_n)) = \mathbb{E}(g(Y_n)) \rightarrow \mathbb{E}(g(Y)) = \mathbb{E}(f(X, Y)).$$

Now $\|(X_n, Y_n) - (X, Y_n)\| = \|X_n - X\|$ and since $X_n \rightarrow X$ in probability, $\|X_n - X\| \rightarrow 0$ in probability (implies convergence in distribution). Thus by Theorem 3.1 in [7] we have that $(X_n, Y_n) \Rightarrow (X, Y)$. ■

Theorem 13. *In addition to our running assumptions, we assume that f in (37) is bounded on every compact set and for a given $t > 0$, f is continuous at $X(t)$. Then we have*

$$(46) \quad N^{-1-2\alpha} \mathbb{E} \{ (f^N(X^N(t)))^2 (Z^N(t))^2 \} \rightarrow (f(X(t)))^2 \frac{1}{c_1} \int_0^t a_1(X(s)) ds$$

as $N \rightarrow \infty$. Furthermore, for each $t > 0$,

$$\sup_{s \leq t} \mathbb{E} ((f^N(X^N(s))) Z^N(s))^2 = \mathcal{O}(N^{2\alpha+1}).$$

Proof. Lemma 8 implies the uniform integrability of $N^{-1}(Z^N(t))^2$. By Assumption 4 and (37) we have that $(f^N(X^N(t)))^2/N^{2\alpha}$ is a uniformly bounded sequence. Thus $N^{-1-2\alpha}(f^N(X^N(t)))^2(Z^N(t))^2$ is uniformly integrable.

By Lemma 5 we have that $N^{-2\alpha}(f^N(X^N(t)))^2$ converges to $(f(X(t)))^2$ almost surely. We also have that $N^{-1/2}Z^N(t)$ converges weakly to $U(t)$. Thus by Lemma 12 and the continuous mapping theorem we have that

$$N^{-1-2\alpha}(f^N(X^N(t)))^2(Z^N(t))^2 \Rightarrow (f(X(t)))^2 U^2(t).$$

By Theorem 3.5 from [7], we note that if a uniformly integrable sequence converges weakly, then it converges in the mean, and hence the result (46) follows.

Also, recall that $(f^N(X^N(t)))^2/N^{2\alpha}$ is uniformly bounded, hence

$$N^{-2\alpha-1} \sup_{s \leq t} \mathbb{E} ((f^N(X^N(s))) Z^N(s))^2 \leq \tilde{C} \mathbb{E} (\sup_{s \leq t} \sqrt{N} |Z_N(s)|)^2.$$

Taking \limsup_N and applying Lemma 8 yields the second result. ■

Note that the above theorem does not assume f is continuously differentiable. However, to state the result regarding the estimator variance for the GT method, we still need to assume continuous differentiability on f so that we can use Theorem 10.

Corollary 14. *In addition to our running assumptions, we assume that f in (37) is continuously differentiable. Then for given $t > 0$, the estimator variance of the GT method is asymptotically $\mathcal{O}(N^{2\alpha+1})$ uniformly on $[0, t]$.*

Next, we will explore the variance of the centered GT approach.

Theorem 15. *In addition to our running assumptions, we assume that f in (37) is continuously differentiable. Then for each $t > 0$,*

$$\sup_{s \leq t} \mathbb{E} \left[(f^N(X^N(s)) - \mathbb{E}[f^N(X^N(s))]) Z^N(s) \right]^2 = \mathcal{O}(N^{2\alpha}).$$

Proof. Write

$$\begin{aligned} & \mathbb{E} \left(\left| \frac{f^N(X^N(s))}{N^\alpha} - \mathbb{E} \left(\frac{f^N(X^N(s))}{N^\alpha} \right) \right|^2 (Z^N(s))^2 \right) \\ & \leq 2\mathbb{E} \left(\left| \frac{f^N(X^N(s))}{N^\alpha} - f(X(s)) \right|^2 (Z^N(s))^2 \right) \\ & \quad + 2\mathbb{E} \left(\left| f(X(s)) - \mathbb{E} \left(\frac{f^N(X^N(s))}{N^\alpha} \right) \right|^2 (Z^N(s))^2 \right) \\ & \leq 2\mathbb{E} \left(\left| \frac{f^N(X^N(s))}{N^\alpha} - f(X(s)) \right|^2 (Z^N(s))^2 \right) \\ & \quad + 2\mathbb{E} \left(\left| \frac{f^N(X^N(s))}{N^\alpha} - f(X(s)) \right|^2 \right) \mathbb{E}(Z^N(s))^2, \end{aligned}$$

where the last inequality is true due to the fact that $f(X(s))$ is deterministic. Using a similar argument as in the proof of Theorem 10, the first term on the right-hand side can be bounded by

$$4L_K^2 \mathbb{E} \left(|\sqrt{N} Z_N(s)| \right)^2 + 4C_K^2 \mathbb{E} \left(|V_N(s)| \sqrt{N} |Z_N(s)| \right)^2.$$

Similarly, the second term on the right-hand side can be bounded by

$$4L_K^2 \mathbb{E} \left(\sqrt{N} |Z_N(s)| \right)^2 + 4C_K^2 \mathbb{E} |V_N(s)|^2 \mathbb{E} \left(\sqrt{N} |Z_N(s)| \right)^2.$$

Both of the above terms are bounded in N uniformly on $[0, t]$ by Lemmas 8 and 9. ■

Combining this result with Theorem 10, the following corollary is immediate.

Corollary 16. *For any given $t > 0$, the estimator variance of the CGT method is asymptotically $\mathcal{O}(N^{2\alpha})$ uniformly on $[0, t]$.*

Theorem 17. *In addition to our running assumptions, we assume that f in (37) is continuously differentiable. Then for each $t > 0$ and $h > 0$,*

$$\sup_{s \leq t} \text{Var}(f^N(X^N(s, c+h)) - f^N(X^N(s, c))) = \mathcal{O}(N^{2\alpha-1}).$$

That is, the estimator variance of the FD method is asymptotically $\mathcal{O}(N^{2\alpha-1})$.

Proof. Note that

$$\begin{aligned} & \text{Var}(f^N(X^N(s, c+h)) - f^N(X^N(s, c))) \\ & \leq 2\text{Var}(f^N(X^N(s, c+h))) + 2\text{Var}(f^N(X^N(s, c))). \end{aligned}$$

Hence it is sufficient to show that $\text{Var}(f^N(X^N(t, c))) = \mathcal{O}(N^{2\alpha-1})$ for any c . We write

$$\frac{1}{N^{2\alpha-1}} \text{Var}(f^N(X^N(s, c))) = N \mathbb{E} \left(\left| \frac{f^N(X^N(s, c))}{N^\alpha} - \mathbb{E} \left(\frac{f^N(X^N(s, c))}{N^\alpha} \right) \right|^2 \right).$$

One can estimate the right-hand side by using the same argument as in Theorem 15 to obtain an upper bound $8L_K^2 + 8C_K^2 \mathbb{E}(|V_N(s)|)^2$, which is bounded in N uniformly on $[0, t]$ by Lemma 9. ■

Remark 5. Based on Theorem 10, Corollary 14, Corollary 16, and Theorem 17, we may expect the RSDs of the GT, CGT, and FD methods to scale as $\mathcal{O}(N^{1/2})$, $\mathcal{O}(1)$ and $\mathcal{O}(N^{-1/2})$, respectively. Since in Theorem 10 we do not have an exact limit for the sensitivity itself, this conclusion is not rigorously proven. As mentioned in Remark 4, under additional assumptions, this conclusion will be true. Our numerical results in the next section also support this statement. Moreover, we note that the $\mathcal{O}(N^{2\alpha+1})$ estimates in Theorem 13 and Corollary 14 are sharp.

5. Numerical examples. We illustrate the dependence of RSD of various sensitivity estimators (with respect to the deterministic parameter) on the system size N via numerical examples. When comparing the GT or CGT method with the FD or RPD method, we must bear in mind that while GT and CGT do not have method parameters, the FD method has a perturbation parameter h and the RPD method has a window size parameter w , making the comparison not straightforward. Moreover, the FD and RPD methods are biased. A proper practical comparison involves choosing parameters h and w to obtain an acceptable bias. We do not pursue such a detailed comparison here as we are focused solely on the dependence on system size N . In the case of FD or RPD methods, we fix h or w , respectively, and vary N . We also use the CRN FD method instead of the IRN FD, as that is the more commonly used approach. Moreover, since our variance estimates for FD methods were derived based on an upper bound which is twice that of the IRN FD method, it is important to compare the performance of CRN FD to see if the order estimate $\mathcal{O}(N^{-1/2})$ for the RSD is sharp.

We note that in the very large system size limit, the stochastic system behaves nearly deterministically and hence none of these stochastic sensitivity methods are needed; traditional ODE sensitivity methods would do. However, when the system size N is modestly large, say, $N = 100$, the system may not be approximated by the ODE and our asymptotic analysis may be relevant in this regime. Our numerical results below show this.

5.1. Numerical example 1. The reversible isomerization model consists of two species S_1 and S_2 and involves the following two reactions:



In the model with system size N , the intensity functions for processes R_1^N and R_2^N are

$$a_1^N(X^N(t), c) = c_1 X_1^N(t),$$

$$a_2^N(X^N(t), c) = c_2 X_2^N(t),$$

respectively. The stoichiometric vectors are $\nu_1 = [-1, 1]^T$ and $\nu_2 = [1, -1]^T$.

In this example, the expectation of the population of species at a fixed time t can be computed analytically:

$$(48) \quad E[X_1^N(t)] = X_1^N(0) + \frac{1 - e^{-(c_1+c_2)t}}{c_1 + c_2} (c_2 X_2^N(0) - c_1 X_1^N(0)),$$

$$(49) \quad E[X_2^N(t)] = X_2^N(0) - \frac{1 - e^{-(c_1+c_2)t}}{c_1 + c_2} (c_2 X_2^N(0) - c_1 X_1^N(0)),$$

where $X_1^N(0)$ and $X_2^N(0)$ are assumed to be deterministic. One can compute the exact sensitivities by differentiating (48) and (49) with respect to parameters. In the numerical tests considered here, we choose parameters $c_1 = 0.3$ and $c_2 = 0.2$ and the initial population $X_1^N(0) = N$ and $X_2^N(0) = N$, where N is the system size parameter. We set the terminal time $T = 10$ and compute the sensitivity for $N = 1, 2, 5, 10, 50, 100, 200, 300, 400, 500, 600, 700, 800, 900$, and 1000. We use four different methods here, namely, GT, CGT, CRN FD, and RPD. We note that by CRN FD, we mean the common random number and one-sided finite difference method in conjunction with Gillespie's SSA [24]. The perturbation parameter for the CRN FD method is $h = 0.01$ for parameter c_1 and the window size parameter $w = 1.0$ for the RPD method for terminal time $T = 10$. The number of trajectories for simulation is $N_s = 10^6$ for each system size N . We consider sensitivities with respect to c_1 of the expected values of four different output functions.

The first output function we consider here is $f^N(x) = x_1$ for all N , that is, we compute the sensitivity of $\mathbb{E}(X_1^N(T))$ with respect to parameter c_1 . Obviously, conditions in Assumption 5 are satisfied with $\alpha = 1$ and $f(x) = x_1$. We examine the growth of sensitivity of $\mathbb{E}(X_1^N(T))$ with respect to c_1 in terms of N using 10^6 independent trajectories. The computed sensitivity and the error in the sensitivity estimate are shown in Figure 1(a), and Figure 1(b) shows the loglog plot of RSD of all four methods.

The second output function we use for testing is $f^N(x) = x_1^2$ for all N . By (37), $f(x) = x_1^2$ and $\alpha = 2$ in Assumption 5. Similar to the case of output function $f^N(x) = x_1$, the exact sensitivity in this case can be calculated and hence we show the error in the sensitivity estimate as an inset plot. See Figure 2 for sensitivity and RSD. The third output function we consider is $f^N(x) = \sin(x_1/N)$ and so $f(x) = \sin x_1$. It can be seen that for this case, $\alpha = 0$ in Assumption 5. The plot for the numerical result is shown in Figure 3.

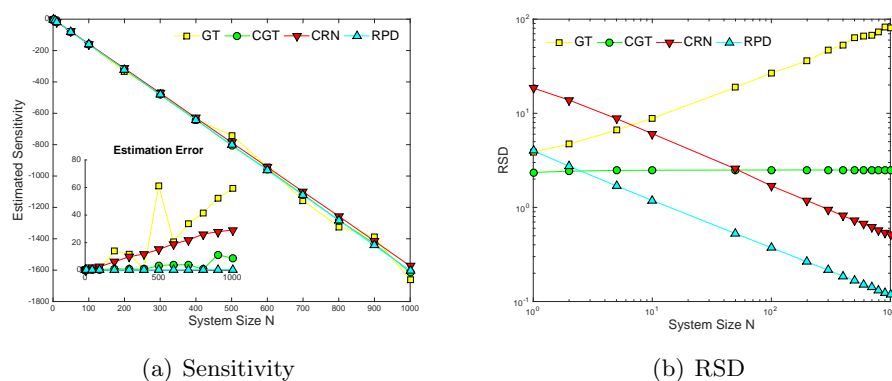


Figure 1. Estimated sensitivity (left) and error in the sensitivity estimate (inset) of $\mathbb{E}(X_1^N(T))$ with respect to c_1 and RSD (right) at terminal time $T = 10$ for the reversible isomerization model.

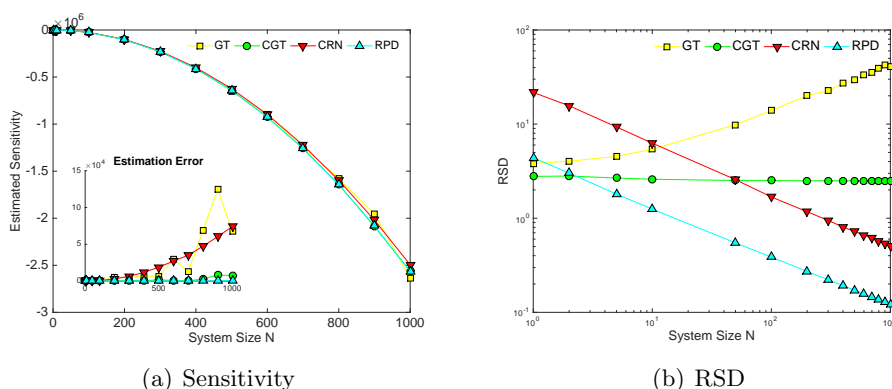


Figure 2. Estimated sensitivity (left) and error in the sensitivity estimate (inset) of $\mathbb{E}(X_1^N(T))^2$ with respect to c_1 and RSD (right) at terminal time $T = 10$ for the reversible isomerization model.

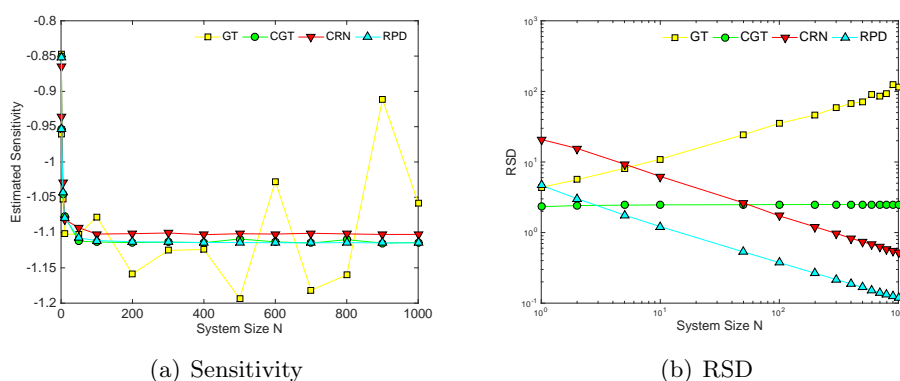


Figure 3. Estimated sensitivity of $\mathbb{E}(\sin(X_1^N(T)/N))$ with respect to c_1 (left) and RSD (right) at terminal time $T = 10$ for the reversible isomerization model.

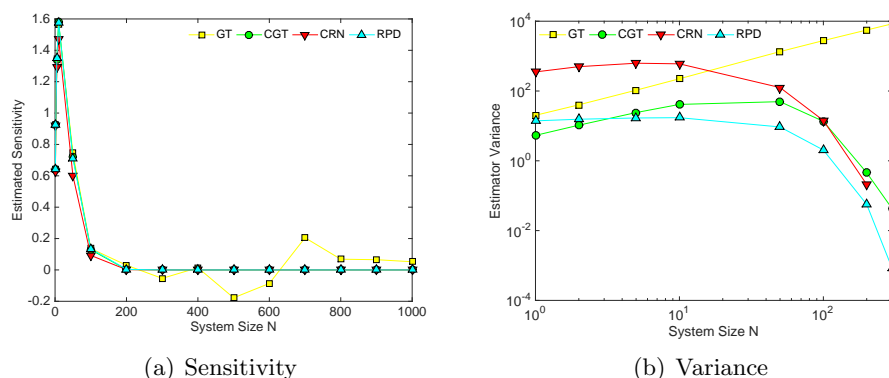


Figure 4. Estimated sensitivity of $\mathbb{P}(X_1^N(T) \leq X_2^N(T))$ with respect to c_1 (left) and variance (right) at terminal time $T = 10$ for the reversible isomerization model.

Table 1

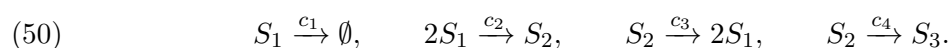
Observed slopes (via regression for large N) for the loglog plots of RSD for reversible isomerization model, that is, R_1 , R_2 , and R_3 are the observed asymptotic order of the estimator RSD (as a power of N) for $\mathbb{E}(X_1^N(T))$, $\mathbb{E}(X_1^N(T))^2$, and $\mathbb{E}(\sin(X_1^N(T)/N))$, respectively.

	R_1	R_2	R_3
GT	0.4992	0.4895	0.5724
CGT	-0.0004	-0.0008	0.0009
CRN FD	-0.5156	-0.5160	-0.5162
RPD	-0.5005	-0.5000	-0.5000

The last output function we consider here is the indicator function $f^N(x) = 1_{\{x_1 \leq x_2\}}(x)$, which does not satisfy the conditions in our theorems since $f = 1_{\{x_1 \leq x_2\}}$ is not continuously differentiable. However, numerical tests still show similar behavior as indicated by our theorems. Note that the sensitivity approaches zero as N increases to ∞ and hence RSD is not well defined for large N . Instead, we plot the estimator variance against N in Figure 4(b).

Additionally, Table 1 summarizes the rate of growth (as a power of N) of the numerically estimated RSD for the different estimators considered above. The results are in agreement with the theory.

5.2. Numerical example 2. As a second numerical example, let us consider the decaying-dimerizing model [11]



The stoichiometric vectors are $\nu_1 = [-1, 0, 0]^T$, $\nu_2 = [-2, 1, 0]^T$, $\nu_3 = [2, -1, 0]^T$, and $\nu_4 = [0, -1, 1]^T$. We set the initial population to be $X_1^N(0) = 10N$, $X_2^N(0) = 0$, $X_3^N(0) = 0$. Using the stochastic mass action form (18), the intensity for processes R_1^N , R_2^N , R_3^N , and R_4^N is

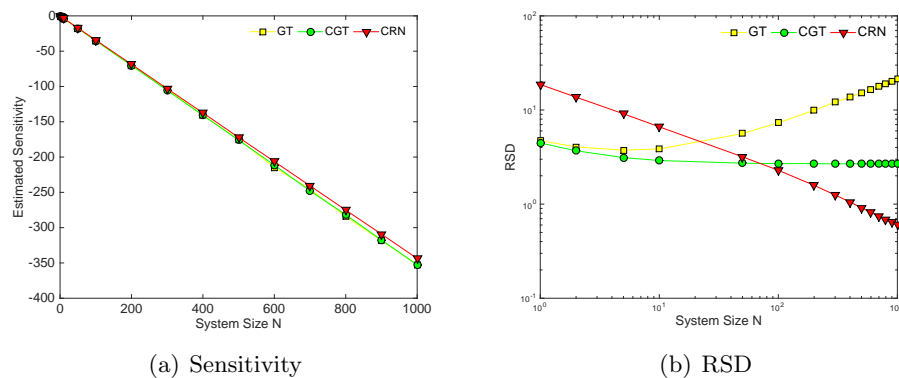


Figure 5. Estimated sensitivity of $\mathbb{E}[X_1^N(T)]$ with respect to c_1 and RSD at terminal time $T = 5$ for the decaying-dimerizing model.

$$\begin{aligned} a_1^N(X^N(t), c) &= c_1 X_1^N(t), \\ a_2^N(X^N(t), c) &= \frac{c_2}{2N} X_1^N(t)(X_1^N(t) - 1), \\ a_3^N(X^N(t), c) &= c_3 X_2^N(t), \\ a_4^N(X^N(t), c) &= c_4 X_2^N(t). \end{aligned}$$

We set the parameters as follows: $c_1 = 1.0$, $c_2 = 0.002$, $c_3 = 0.5$, and $c_4 = 0.04$. Note that the intensity for the second reaction is not linear, hence an analytical formula for the sensitivity is not attainable. We test the sensitivity and RSD for $\mathbb{E}[f^N(X_1^N)]$ with respect to c_1 . For the CRN FD method, we use the one-sided FD scheme and perturb the parameter c_1 by $h = 0.01$. Note that RPD is not applicable for this example since the firing of the first reaction will prevent the second reaction from happening when the population of S_1 is 1 (see [26]). Therefore, we only examine the RSDs of GT, CGT, and CRN FD here. For each system size N , the number of trajectories we use for simulation is $N_s = 10^6$. Plots of the sensitivity and RSD are shown in Figures 5, 6, and 7 for $\mathbb{E}(X_1^N(T))$, $\mathbb{E}(X_1^N(T))^2$, and $\mathbb{E}(\sin(X_1^N(T)/N))$, respectively. The rate of growth (as a power of N) of the numerically estimated RSD is summarized in Table 2.

5.3. Numerical example 3. In this numerical example, we revisit the reversible isomerization network to illustrate the asymptotic behavior of various estimators in terms of the terminal time T . Note that in this example, the deterministic parameters c_j and the stochastic parameters c'_j are the same. For ease of notation, we suppress N because we fix $N = 10$ and only let T change in this simulation. The initial population is $X_1(0) = 10$ and $X_2(0) = 10$. Parameters are taken to be $c_1 = 0.3$ and $c_2 = 0.2$ as before. In this case, the exact sensitivity can be obtained by taking derivative with respect to c_1 for (48). Figure 8(a) shows the sensitivities estimated by GT, CGT, and CRN FD against the true sensitivity as a function of T . Figure 8(b) shows the estimator variances as a function of T . It can be seen that all three estimators show a variance that grows linearly in T for the range of values of T considered here.

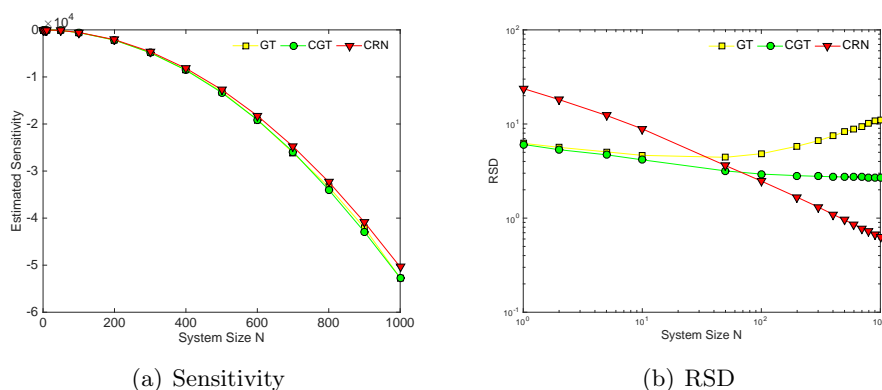


Figure 6. Estimated sensitivity of $\mathbb{E}(X_1^N(T))^2$ with respect to c_1 and RSD at terminal time $T = 5$ for the decaying-dimerizing model.

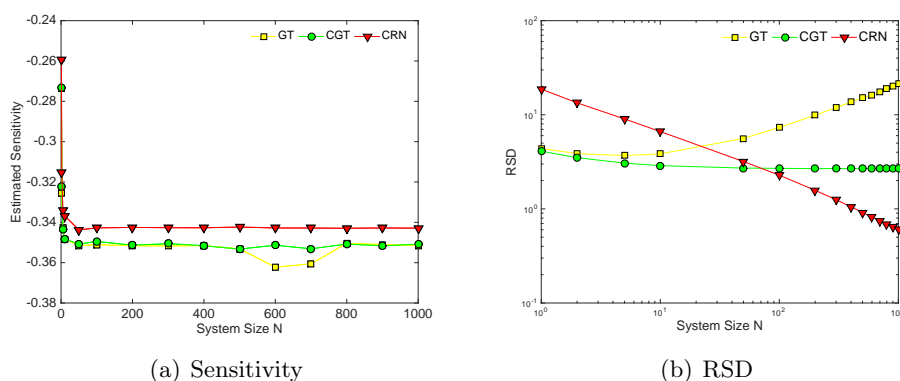


Figure 7. Estimated sensitivity of $\mathbb{E}(\sin(X_1^N(T)/N))$ with respect to c_1 and RSD at terminal time $T = 5$ for the decaying-dimerizing model.

Table 2

Observed slopes (via regression) for the loglog plots for RSD for the decaying-dimerizing model, that is, R_1 , R_2 , and R_3 are the observed asymptotic order of the estimator RSD (as a power of N) for $\mathbb{E}(X_1^N(T))$, $\mathbb{E}(X_1^N(T))^2$, and $\mathbb{E}(\sin(X_1^N(T)/N))$, respectively.

	R_1	R_2	R_3
GT	0.4689	0.4100	0.4737
CGT	-0.0040	-0.0257	-0.0008
CRN FD	-0.6022	-0.6068	-0.6009

In fact, this observation can be justified for the GT and CGT methods as follows. Recall the definition of the centered processes $M_j(t) = R_j(t) - \int_0^t a_j(X(s))ds$, $j = 1, \dots, m$. Since $X_j(t)$ are bounded in this network, one can show that

$$\mathbb{E}M_j^2(t) = \mathbb{E}([M_j, M_j](t)) = \mathbb{E}R_j(t) = c_j \int_0^t \mathbb{E}X_j(s)ds = \mathcal{O}(t),$$

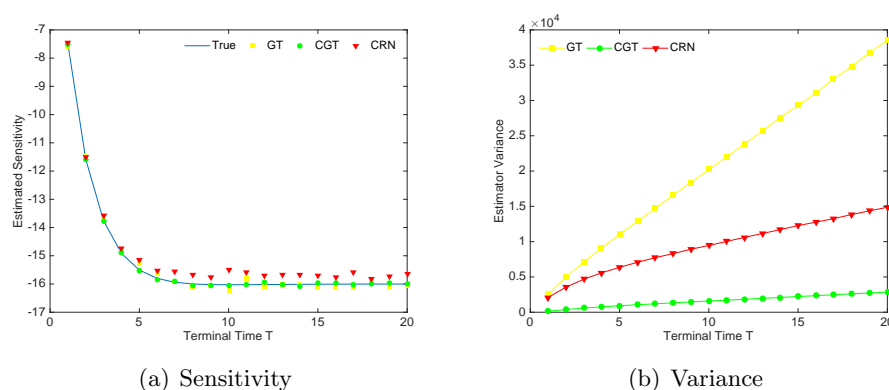


Figure 8. Estimated and true sensitivities (left) of $\mathbb{E}X_1(t)$ with respect to c_1 and the estimator variances (right) for the reversible isomerization model. The terminal time T (x -axis) ranges from 1 to 20.

where the first equality holds since $M_j(t), j = 1, 2$, is a L^2 -bounded martingale (see [21]). Therefore, we conclude that $\mathbb{E}Z^2(t) = \mathcal{O}(t)$ because in this case $Z(t) = c_1^{-1}M_1(t)$ and hence the variances of both GT and CGT are of $\mathcal{O}(t)$.

As for the variance of the FD estimator, the observed growth is approximately linear in t in the range of 10 to 20. However, from the upper bound used in the proof of Theorem 17, it is easy to see that the estimator variance remains bounded as $t \rightarrow \infty$.

6. Discussion and concluding remarks. Our primary goal in this paper was to provide an analytical explanation of the phenomenon of larger estimator variance of the GT method compared to the FD (as well as RPD in the context of chemical kinetics) methods reported frequently in the literature [5, 20, 24, 26]. This was accomplished by our analysis in terms of system size N . The system size N was taken to be proportional to system volume in the context of stochastic chemical kinetics. Our analysis showed that the RSD (see (7) for a definition) of the GT, CGT, and FD sensitivity estimators is $\mathcal{O}(N^{1/2})$, $\mathcal{O}(1)$, and $\mathcal{O}(N^{-1/2})$, respectively, as $N \rightarrow \infty$. The numerical examples provided also illustrate this point. Additionally, our numerical examples suggested that the RSD of the RPD method also scales as $\mathcal{O}(N^{-1/2})$. We also showed that the RB (see (8) for a definition) of any FD method was asymptotically $\mathcal{O}(1)$ as $N \rightarrow \infty$. We note that in our analysis of the FD methods, we kept h fixed and considered $N \rightarrow \infty$ limit. Now we discuss, at least in theory, how h may be chosen in terms of system size N to obtain the best performance for the FD methods.

Number of simulations required to achieve a given RE. Since the FD methods are biased while the GT and CGT methods are not, we shall use the RE to compare the efficiencies of the GT, CGT, and FD estimators. More precisely, we shall estimate the number of trajectory simulations N_s required to achieve a given tolerance δ for the RE in the mean square sense, which includes RB and RSD (see (6) for the exact definition).

Our analysis for the FD methods was carried out so that large N behavior for fixed h was obtained. We may combine our large N analysis with small h behavior of the FD methods already studied in the literature [5]. In general, the bias of the one-sided FD estimator is $\mathcal{O}(h)$ as $h \rightarrow 0$, so we may expect the RB of an FD estimator to be given by $RB \approx C_2 h$ for small

h and large N , where C_2 does not depend on N or h . If higher order FD is used, then one expects $\text{RB} \approx C_2 h^{\gamma_1}$, where $\gamma_1 \geq 1$ in general. For instance, for the two-sided FD estimator we have that $\gamma_1 = 2$.

Moreover, when using the IRN FD method, the variance is $\mathcal{O}(1/h^2)$ as $h \rightarrow 0$, which is similar to behavior of the upper bound used in our proof of Theorem 17. However, when using CRN FD methods, one may typically expect $\mathcal{O}(1/h)$ dependence [5, 24]. This is because $\text{Cov}(f(X(t, c + h)), f(X(t, c)))$ is typically $\mathcal{O}(h)$ as $h \rightarrow 0$. Hence we may write $\text{RSD}^2 \approx C_1/(Nh^{\gamma_2})$ for small h and large N , where typically $\gamma_2 = 1$ or 2 depending on whether CRN or IRN is used, and C_1 is independent of h and N . Combining the bias and the variance, and using (9), we expect that for an FD method

$$(51) \quad \text{RE}^2 = \frac{\text{RSD}^2}{N_s} + \text{RB}^2 \approx \frac{C_1}{N_s N h^{\gamma_2}} + C_2^2 h^{2\gamma_1}.$$

At this point, we must remark that in order for the above approximation to hold rigorously, one must establish the joint limit as $(N, h) \rightarrow (\infty, 0)$. We believe that this could be done under additional regularity assumptions, but we shall not pursue this in this paper.

Extending the idea in [5] to include system size N , we look for the optimal choice of h (the one that minimizes RE), for a given system size N and number of simulations N_s . With some effort, one can see that the optimal h is given by

$$h \propto N^{\frac{-1}{2\gamma_1 + \gamma_2}} N_s^{\frac{-1}{2\gamma_1 + \gamma_2}},$$

and hence the minimal square RE for an FD method has the proportionality

$$(52) \quad \text{RE}^2 \propto N^{\frac{-2\gamma_1}{2\gamma_1 + \gamma_2}} N_s^{\frac{-2\gamma_1}{2\gamma_1 + \gamma_2}}.$$

On the other hand, for the CGT method, $\text{RE}^2 = \text{RSD}^2/N_s = C_3/N_s$ for large N , where C_3 is independent of N and N_s . Likewise, for the GT method, $\text{RE}^2 = \text{RSD}^2/N_s = C_4 N/N_s$ for large N , where C_4 is independent of N and N_s . Hence, for a specified value of δ for RE and a given system size N , the numbers of simulations required for the different methods are given by

$$(53) \quad N_s^{\text{FD}} \propto \delta^{-2 - \frac{\gamma_2}{\gamma_1}} N^{-1}, \quad N_s^{\text{CGT}} \propto \delta^{-2}, \quad N_s^{\text{GT}} \propto N \delta^{-2}.$$

We note that, as observed in [5], the optimal dependence of N_s on δ is δ^{-2} , which is achieved for an unbiased method. The biased FD methods have suboptimal dependence on δ , unless $\gamma_2 = 0$, which is typically not the case in the context of discrete state systems, as $\gamma_2 = 0$ implies the validity of the (unregularized) PD method [5]. However, when N is much larger than $\delta^{-\gamma_2/\gamma_1}$, we expect the FD method to be more efficient than the CGT or GT. For instance, for $\delta = 0.01$, if $N \gg 10$, say, $N = 50$, for instance, we may expect the two-sided CRN FD ($\gamma_1 = 2, \gamma_2 = 1$) to be more efficient than CGT, which will be more efficient than GT. If one-sided CRN FD is used ($\gamma_1 = \gamma_2 = 1$) or two-sided IRN FD is used ($\gamma_1 = \gamma_2 = 2$), we expect FD to be more efficient when $N \gg 100$, say, $N = 500$. If one-sided IRN FD is used ($\gamma_1 = 1, \gamma_2 = 2$) we expect FD to be more efficient than CGT only for $N \gg 10^4$.

Since the constants of proportionality that appear in the above discussion are not known in practice and are typically harder to estimate than the sensitivity itself, one may not expect to choose h in a straightforward manner based on the above discussion. Nevertheless, the above discussion provides some idea of the optimal efficiency that could be expected.

We also note that the comparison of an unbiased estimator with a biased one is more nuanced and qualitative. This is because, while one can estimate the variance of an estimator from the simulation, its bias cannot be estimated reliably unless one knows the exact quantity to be estimated! As a consequence, an unbiased estimator is preferable to a biased one, unless the unbiased estimator has exceedingly larger variance compared to the biased one. In this context, we mention that multilevel Monte Carlo approaches (see [2], for instance) may be used to combine a biased low variance estimator with an unbiased high variance estimator to obtain an efficient and unbiased estimator.

Factors other than system size that affect the RSD. We note that factors other than system size also affect the RSD of an estimator. One factor to study will be the dependence on t as $t \rightarrow \infty$. Our numerical simulations showed linear growth in t behavior for GT, CGT, and even FD methods for a practical range of t values (up to a few multiples of the time to stationarity). However, from a simple upper bound for the variance of the FD methods, we expect this growth to reach a finite maximum for systems that are ergodic. The $\mathcal{O}(t)$ behavior (as $t \rightarrow \infty$) for the variance of the GT and CGT methods can be justified theoretically, as explained in section 5.3. Thus, dependence on time does not explain the greater variance of GT compared to CGT.

Extension of the variance analysis. Our analysis made special use of the deterministic limit in the large system size under what is known as the classical scaling which was used by Ethier and Kurtz [9]. In other words, after suitable scaling, $f^N(X^N(t))$ converges to the deterministic limit $f(X(t))$ almost surely. However, the scaled weight processes $Z^N(t)/\sqrt{N}$ converge weakly to a Gaussian process $U(t)$. Our analysis combined the two limits to obtain the desired results. Our results were proven under Assumptions 1–5 stated in section 2. The first assumption assumes that the parameters enter multiplicatively: $a_j(x, c) = c_j b_j(x)$. This is satisfied by the stochastic mass action form of intensities. In some literature on chemical kinetics, there are some other forms of intensity functions that are used. Relaxing Assumption 1 to a general form will make the weight process Z^N more complicated, and it will be given by a stochastic integral where both the integrand and the integrator are stochastic processes indexed by N . To obtain convergence of $N^{-1-2\alpha} \mathbb{E}[(f^N(X^N(t)))^2 (Z^N(t))^2]$ one may need the result from [18] which analyzes the limit of a sequence of stochastic integrals. We speculate that Assumption 4 may be relaxed using stopping time arguments and sufficient integrability assumptions on the process.

In many practical systems some species are present in small numbers while others are present in large numbers, and some reaction parameters are much larger than the others, making the system “stiff.” The classical scaling studied here does not capture this. The more general scaling proposed in [6, 17] (again by Kurtz and collaborators) involves introducing a parameter N which appears with different exponents in both the stochastic parameters c'_j and the scaling of species and time itself. These analyses often provide stochastic limits to the scaled processes X^N . One could extend our current analysis along these lines to explore more subtle dependencies of the estimator variances. A related earlier work which scales all

“species” by the same factor ϵ , and scales time differently $\epsilon^{-\alpha}$, in the context of processes driven by Levy measure can be found in [27].

Acknowledgment. We would like to thank the anonymous referees for the comments that helped improve the manuscript.

REFERENCES

- [1] D. F. ANDERSON, *An efficient finite difference method for parameter sensitivities of continuous time Markov chains*, SIAM J. Numer. Anal., 50 (2012), pp. 2237–2258.
- [2] D. F. ANDERSON, D. J. HIGHAM, AND Y. SUN, *Complexity of multilevel monte carlo tau-leaping*, SIAM J. Numer. Anal., 52 (2014), pp. 3106–3127.
- [3] D. F. ANDERSON AND T. G. KURTZ, *Continuous time Markov chain models for chemical reaction networks*, in Design and Analysis of Biomolecular Circuits, Springer, New York, 2011, pp. 3–42.
- [4] A. P. ARKIN, J. ROSS, AND H. H. MCADAMS, *Stochastic kinetic analysis of developmental pathway bifurcation in phage λ -infected escherichia coli cells*, Genetics, 149 (1998), pp. 1633–1648.
- [5] S. ASMUSSEN AND P. W. GLYNN, *Stochastic Simulation: Algorithms and Analysis*, Stoch. Model. Appl. Probab. 57, Springer, New York, 2007.
- [6] K. BALL, T. G. KURTZ, L. POPOVIC, AND G. REMPALA, *Asymptotic analysis of multiscale approximations to reaction networks*, Ann. Appl. Probab., 16 (2006), pp. 1925–1961.
- [7] P. BILLINGSLEY, *Convergence of Probability Measures*, 2nd ed., Wiley Ser. Probab. Stat. Wiley, New York, 1999.
- [8] P. BRÉMAUD, *Point Processes and Queues: Martingale Dynamics*, Springer, New York, 1981.
- [9] S. N. ETHIER AND T. G. KURTZ, *Markov Processes: Characterization and Convergence*, 2nd ed., Wiley, New York, 2005.
- [10] D. T. GILLESPIE, *Exact stochastic simulation of coupled chemical reactions*, J. Phys. Chem., 81 (1977), pp. 2340–2361.
- [11] D. T. GILLESPIE, *Approximate accelerated stochastic simulation of chemically reacting systems*, J. Chem. Phys., 115 (2001), pp. 1716–1733.
- [12] P. GLASSERMAN, *Gradient Estimation via Perturbation Analysis*, Springer, New York, 1990.
- [13] A. GUPTA, *personal communication*, 2016.
- [14] A. GUPTA AND M. KHAMMASH, *Unbiased estimation of parameter sensitivities for stochastic chemical reaction networks*, SIAM J. Sci. Comput., 35 (2013), pp. A2598–A2620.
- [15] A. GUPTA AND M. KHAMMASH, *An efficient and unbiased method for sensitivity analysis of stochastic reaction networks*, J. R. Soc. Interface (2014), 20140979.
- [16] J. JACOD AND A. N. SHIRYAEV, *Limit Theorems for Stochastic Processes*, 2nd ed., Probab. Theory Stoch. Process. 288, Springer, New York, 2003.
- [17] H. KANG AND T. G. KURTZ, *Separation of time-scales and model reduction for stochastic reaction networks*, Ann. Appl. Probab., 23 (2013), pp. 529–583.
- [18] T. G. KURTZ AND P. PROTTER, *Weak limit theorems for stochastic integrals and stochastic differential equations*, Ann. Probab., 19 (1991), pp. 1035–1070.
- [19] H. H. MCADAMS AND A. P. ARKIN, *It’s a noisy business! Genetic regulation at the nanomolar scale*, Trends Genetics, 15 (1999), pp. 65–69.
- [20] S. PLYASUNOV AND A. P. ARKIN, *Efficient stochastic sensitivity analysis of discrete event systems*, J. Comput. Phys., 221 (2007), pp. 724–738.
- [21] P. PROTTER, *Stochastic Integration and Differential Equations*, 2nd ed., Springer, New York, 2005.
- [22] M. RATHINAM, *Moment growth bounds on continuous time Markov processes on non-negative integer lattices*, Quart. Appl. Math., 73 (2015), pp. 347–364.
- [23] M. RATHINAM AND H. EL-SAMAD, *Reversible-equivalent-monomolecular tau: A leaping method for “small number and stiff” stochastic chemical systems*, J. Comput. Phys., 224 (2007), pp. 897–923.
- [24] M. RATHINAM, P. W. SHEPPARD, AND M. KHAMMASH, *Efficient computation of parameter sensitivities of discrete stochastic chemical reaction networks*, J. Chem. Phys., 132 (2010), 034103.

- [25] L. C. G. ROGERS AND D. WILLIAMS, *Diffusions, Markov processes, and Martingales*. Vol. 2, Cambridge Math. Lib., Cambridge University Press, Cambridge, UK, 2000.
- [26] P. W. SHEPPARD, M. RATHINAM, AND M. KHAMMASH, *A pathwise derivative approach to the computation of parameter sensitivities in discrete stochastic chemical systems*, J. Chem. Phys., 136 (2012), 034115.
- [27] M. TOMISAKI, *Homogenization of càdlàg processes*, J. Math. Soc. Japan, 44 (1992), pp. 281–305.
- [28] P. B. WARREN AND R. J. ALLEN, *Steady-state parameter sensitivity in stochastic modeling via trajectory reweighting*, J. Chem. Phys., 136 (2012), 104106.
- [29] W. WHITT, *Proofs of the martingale FCLT*, Probab. Surv., 4 (2007), pp. 268–302.