# IMPROVING PREDICTION FROM STOCHASTIC SIMULATION VIA MODEL DISCREPANCY LEARNING

Henry Lam Xinyu Zhang

Department of Industrial Engineering and
Operations Research
Columbia University
500 W. 120th Street
New York, NY 10027, USA

Matthew Plumlee

Department of Industrial Engineering and Management Sciences Northwestern University 2145 Sheridan Road Evanston, IL 60208, USA

#### **ABSTRACT**

Stochastic simulation is an indispensable tool in operations and management applications. However, simulation models are only approximations to reality, and often bear discrepancies with the generating processes of real output data. We investigate a framework to statistically learn these discrepancies under the presence of data on past implemented system configurations, which allows us to improve prediction using simulation models. We focus on the case of general continuous output data that generalizes previous work. Our approach utilizes (a combination of) regression analysis and optimization formulations constrained on suitable summary statistics. We demonstrate our approach with a numerical example.

# 1 INTRODUCTION

Stochastic simulation is used ubiquitously in decision analytics in many operations and management applications. It describes system dynamics under alternate configurations that can be used for analytical tasks. While modelers often attempt to build models as realistic as possible, due to resource constraints and ignorance of unobserved system features, these models are at best approximations and in almost all practical cases bear discrepancies with reality. For example, in operational settings that are naturally cast as "queues", the serial structures of interarrival times (e.g., Livny et al. 1993), discretionary or strategic behaviors (e.g., balking, reneging, queue selection; Pazgal and Radas 2008, Veeraraghavan and Debo 2009), individualized and contextual behaviors of the customers and servers (e.g., Aksin et al. 2007), and structural server changes over time (e.g., Brown et al. 2005) are all difficult to capture. Inadequate reflection of these features can degrade the predictive power of the simulation model and potentially affect the reliability of decision-making.

In this paper, we investigate a general framework that builds on our previous work Plumlee and Lam (2016, 2017) to learn about the imperfectness of stochastic simulation models and subsequently use this information to improve prediction. To describe our investigation, we first introduce some notations and define the notion of model discrepancy. Imagine that we are interested in an output variable,  $Y \in \mathcal{Y}$ , from a stochastic system which has a true distribution  $\pi(\cdot)$ . The variable can represent waiting time, queue length, reward, other responses or any collections of them that are relevant to decision-making. A typical performance analysis consists of evaluating a quantity  $\rho(\pi)$ , which in this paper we focus on the expectation-type performance measure  $E_{\pi}[f(Y)]$  for some function  $f(\cdot): \mathcal{Y} \to \mathbb{R}$ .

When real-world observations from  $\pi$  are abundant, the performance analyses can be accurately approximated by replacing  $\pi$  with data, e.g., the empirical distribution. However, in decision analyses (e.g., optimality or feasibility tests), we are typically interested in system designs that are sparsely sampled or

even never adopted before. So, instead, model builders use "simulation models" that generate outputs with distribution, say,  $\tilde{\pi}(\cdot)$ , as an approximation to  $\pi(\cdot)$  that can be estimated much more easily via simulation replications. If  $\tilde{\pi}(\cdot) \neq \pi(\cdot)$ , then  $\rho(\tilde{\pi})$  may not equal  $\rho(\pi)$  and the outcomes of the analysis become erroneous. Our interest is to combine both the observations on  $\pi$  and the simulation model  $\tilde{\pi}$  to generate prediction for  $\rho(\pi)$  that is better than using either one alone.

Our approach hinges on estimating the difference between  $\tilde{\pi}(\cdot)$  and  $\pi(\cdot)$  and then integrating this bias estimate with the simulation model to predict  $\rho(\pi)$ . It requires the existence of some control variable or design point  $x \in \mathcal{X}$  that represents the set of system parameters, and that real system observations are available at some x. The value of a design point can include decision variables (e.g., number of servers, routing policy) and important components of the system (e.g., arrival rate of customers). The availability of observations from different design points provides an opportunity to learn and correct for the imperfectness of  $\tilde{\pi}$ , including at un- or under-observed points, by data-pooling.

We consider the class of mappings  $\mathscr{P}: \mathscr{X} \times \mathscr{Y} \to \mathbb{R}$ , such that if  $p \in \mathscr{P}$ , then  $p(x,\cdot)$  is a valid probability measure over  $\mathscr{Y}$ . We expand the definition of the true model  $\pi$  and the simulation model  $\tilde{\pi}$  to be in  $\mathscr{P}$ . In general, we have the pointwise relationship

$$\pi(x, y) = \tilde{\pi}(x, y) + \delta(x, y) \tag{1}$$

where we call  $\delta(\cdot,\cdot)$  the *model discrepancy*. When  $\tilde{\pi} \neq \pi$ , we have  $\delta \neq 0$  and model discrepancy is present. The relationship (1) can be viewed as a stochastic counterpart of a similar discrepancy notion in the literature of deterministic computer experiments (where no y in (1) is present). The latter is under the framework of model calibration that refers generally to the refinement or correction of a hypothetical simulation model through real-world data to enhance the model prediction capability. Starting with the seminal statistical article of Kennedy and O'Hagan (2001), calibration for deterministic simulation has been substantially investigated (e.g., Bayarri et al. 2012, Oakley and O'Hagan 2004, Higdon et al. 2004, Plumlee 2016). In the stochastic simulation literature, other than parameter calibration (implemented in some simulation software such as Anylogic), calibration on model structure is commonly conducted together with model validation, integrated in the model development process, through various techniques including running statistical tests (e.g., Balci and Sargent 1982, Kleijnen 2015, Sargent 2013, Kleijnen, Cheng, and Bettonvil 2001, Kleijnen and Sargent 2000). This paper detours from the common approach, and utilizes the idea of inferring model discrepancies in deterministic computer experiments to improve the prediction of stochastic simulation. Compared to the established practice, the avoidance of building increasingly sophisticated models could potentially save costs and time in refining the model, and prevent the unfortunate situation in case an ultimately satisfactory model is absent.

We will consider two methods in this paper. The first one estimates  $\rho(\delta(x,\cdot)) := \rho(\pi(x,\cdot)) - \rho(\tilde{\pi}(x,\cdot))$  directly by regression against the covariate x. The second method exploits further the fact that  $\pi$  and  $\tilde{\pi}$  are probability distributions, so that they satisfy standard conditions for valid probability measures that constrain the possible values of  $\rho(\pi)$ . This second method can be thought of as inferring the difference of distributions  $\delta(x,\cdot) = \pi(x,\cdot) - \tilde{\pi}(x,\cdot)$  and using it to calculate  $\rho(\pi)$ . However, since  $\delta(x,\cdot)$  can be a high- or even infinite-dimensional object (in the case of continuous output distributions) that is challenging to estimate, our second approach instead operates on a collection of summary statistics on  $\delta(x,\cdot)$ , and posits suitable optimization formulations to compute bounds that take into account the relation among these summary statistics induced by the probability distribution structure.

As mentioned before, this work builds on Plumlee and Lam (2016, 2017) and generalizes their investigation. Besides the difference regarding the representation of model discrepancy (additive versus multiplicative form), Plumlee and Lam (2016, 2017) focus on finite-valued outputs in a Bayesian framework. In this paper, we are primarily interested in outputs that are continuous, which pose substantial additional challenges since they are now infinite-dimensional objects. Our aforementioned second method novelly uses summary statistics to facilitate tractable statistical estimation, but with infinite-dimensional optimization programs posited over the space of probability distributions to recover bounds incurred by the continuous

distributions. With well-chosen functions to define the summary statistics, these optimization programs can be reformulated by duality as tractable finite-dimensional programs. These developments provide an implementable mechanism to take into account distributional information that is infinite-dimensional in nature.

In the rest of this paper, we first describe the two proposed methods in Sections 2 and 3. Then we provide an extensive numerical study and compare our performance with both simulation-only and data-only predictions in Section 4. We conclude this paper in Section 5.

#### 2 LEARNING MODEL DISCREPANCY VIA REGRESSION

Consider a collection of design point values, say  $\{1, \dots, s\}$ . For each value x = i, suppose there are i.i.d. real-world output observations  $Y_{ij} \sim \pi(i,\cdot)$ ,  $j=1,\ldots,n_i$ . The sample size  $n_i$  can be possibly zero for some i (so that no  $Y_{ij}$  exists).

We assume there is enough computational resource, so that the simulation output distribution  $\tilde{\pi}$  can be exactly obtained.

Suppose we are interested in  $\rho(\pi(x_0,\cdot)) = E_{\pi(x_0,\cdot)}[f(Y)]$  for some  $x_0$ . Using data alone, we will output  $ar{Y}_{x_0\cdot}=(1/n_{x_0})\sum_{j=1}^{n_{x_0}}f(Y_{ij}).$  Clearly, if there is no data at  $x_0$ , this approach does not apply. On the other hand, the simulation model gives  $\rho(\tilde{\pi}(x_0,\cdot))=E_{\tilde{\pi}(x_0,\cdot)}[f(Y)]$  as the predicted value of  $\rho(\pi(x_0,\cdot))$ . To combine both  $\{Y_{ij}\}$  and  $\tilde{\pi}$  in our prediction, we consider the estimation of  $\rho(\delta(x_0,\cdot))$  by the

regression problem

$$f(Y_{ij}) - \rho(\tilde{\pi}(i,\cdot)) = h(i) + \varepsilon_{ij} \tag{2}$$

where  $h(\cdot)$  is an unknown function of the design point x, and  $\varepsilon_{ij}$ ,  $j=1,\ldots,n_i$  are mean-zero i.i.d. noises from the real data, for each i. The responses in the regression are  $f(Y_{ij}) - \rho(\tilde{\pi}(i,\cdot))$  whose mean values are the model discrepancies  $\rho(\delta(i,\cdot))$ .

Regression (2) can be conducted by a variety of tools; we use local linear regression with a standard R package in our experiments. This gives us a confidence interval (CI) for  $\rho(\delta(x_0,\cdot))$ , say [L,U]. Our prediction of  $\rho(\pi(x_0,\cdot))$  is then expressed as a CI  $[\rho(\tilde{\pi}(x_0,\cdot)) + L, \rho(\tilde{\pi}(x_0,\cdot)) + U]$ . Clearly, if [L,U] indeed covers the true  $\rho(\delta(x_0,\cdot))$  with probability  $1-\alpha$ , then  $[\rho(\tilde{\pi}(x_0,\cdot))+L,\rho(\tilde{\pi}(x_0,\cdot))+U]$  covers  $\rho(\pi(x_0,\cdot))$ with the same probability.

The above scheme can be modified with

$$f(Y_{ij}) = \beta \rho(\tilde{\pi}(i,\cdot)) + h(i) + \varepsilon_{ij}$$
(3)

The additional unknown scaling parameter  $\beta$  represents a multiplicative model discrepancy. The response in (3) now has mean  $\rho(\pi(i,\cdot))$ . The CI for  $\rho(\pi(x_0,\cdot))$  can be obtained directly from the regression (3).

### COMBINING REGRESSION OUTCOMES WITH OPTIMIZATION

We consider an enhancement of the method in Section 2 by introducing the auxiliary statistics  $\gamma_k(\pi(x,\cdot)) =$  $E_{\pi(x,\cdot)}[g_k(Y)]$  for  $k=1,\ldots,m$ . One of these statistics can be taken as  $\rho(\pi(x,\cdot))$ , e.g., let  $g_1=f$ . Then, instead of using (2) only, we impose the multivariate-output regression

$$g_k(Y_{ij}) - \gamma_k(\pi(i,\cdot)) = h_k(i) + \varepsilon_{ijk} \tag{4}$$

where  $h_k(\cdot)$  is the individual regression function for the k-th statistic. Similar to Section 2,  $\varepsilon_{ijk}$ ,  $j=1,\ldots,n_i$ are mean-zero i.i.d. noises, for each i, and  $g_k(Y_{ij}) - \gamma_k(\pi(i,\cdot))$  has mean equal to  $\gamma_k(\delta(i,\cdot))$ .

The idea is that estimates of  $\gamma_k(\delta(x_0,\cdot)), k=1,\ldots,m$  give more information than using  $\rho(\delta(x_0,\cdot))$ alone. For convenience, we denote  $\gamma = (\gamma_k)_{k=1,\dots,m}$  and  $g = (g_k)_{k=1,\dots,m}$ . To incorporate the new information, we first need a confidence region (CR) for  $\gamma(\delta(x_0,\cdot))$ . As in Section 2, this can be done by a variety of statistical tools, and here we use a basic scheme of running individual local regression at each k, making an additional assumption that  $\varepsilon_{ijk}$ ,  $k=1,\ldots,m$  are independent for each i,j. For each individual regression, we can obtain the point estimate of  $\gamma_k(\delta(x_0,\cdot))$ , say  $\hat{\gamma}_k$ , and its standard error, say  $s_k$ . Under standard assumption we have  $\hat{\gamma}_k \stackrel{approx.}{\sim} N(\gamma_k(\delta(x_0,\cdot)), \sigma_k^2)$ , independent among the k's, where  $\sigma_k^2$  is the sampling variance that can be estimated by  $s_k^2$ . This gives us a CR

$$\mathscr{U} = \left\{ (z_1, \dots, z_m) \in \mathbb{R}^m : \sum_{k=1}^m \frac{(z_k - \hat{\gamma}_k)^2}{s_k^2} \le \chi_{1-\alpha, m}^2 \right\}$$
 (5)

for  $\gamma(\delta(x_0,\cdot))$ , where  $\chi^2_{1-\alpha,m}$  is the  $(1-\alpha)$ -quantile of a  $\chi^2$ -distribution with degree of freedom m. We then impose the optimization problems

$$\max \ \rho(q) \ \text{ subject to } \ \gamma(q) - \gamma(\tilde{\pi}(x_0, \cdot)) \in \mathscr{U}$$
 (6)

and

min 
$$\rho(q)$$
 subject to  $\gamma(q) - \gamma(\tilde{\pi}(x_0, \cdot)) \in \mathcal{U}$  (7)

where the decision variable q is a probability distribution on the response space  $\mathscr{Y}$ . The constraints in (6) and (7) incorporate information on the possible values of  $\gamma(q)$  with probability  $1-\alpha$ . Given this information, the outcomes of (6) and (7) give the best upper and lower bounds on  $\rho(\pi(x_0,\cdot))$ . If the CR  $\mathscr{U}$  contains the true  $\gamma(\tilde{\pi}(x_0,\cdot))$  with confidence level  $1-\alpha$ , then (6) and (7) will form a CI for  $\rho(\pi(x_0,\cdot))$  with at least the same level. This idea originates from the literature of data-driven distributionally robust optimization (e.g., Delage and Ye 2010, Ben-Tal et al. 2013, Goh and Sim 2010, Wiesemann et al. 2014), where  $\mathscr{U}$  is often known as the uncertainty set or the ambiguity set. We also mention that the use of multiple simulation outputs has also been considered in the conventional model validation literature, in which Bonferroni correction is commonly suggested when conducting the relevant statistical tests (e.g., Sargent 2013, Kleijnen 2015).

Note that (6) and (7) are generalized moment problems, where the optimization contains an objective function and constraint functions that are all moments of a random variable. By using conic duality, (6) can be reformulated as the dual problem

$$\min_{\kappa, \nu \in \mathbb{R}, \lambda \in \mathbb{R}^{m}} \quad \kappa + (\gamma(\tilde{\pi}(x_{0}, \cdot)) + \hat{\gamma})^{\mathsf{T}} \lambda + \sqrt{\chi_{1-\alpha, m}^{2}} \nu$$
subject to
$$\kappa + g(y)^{\mathsf{T}} \lambda - f(y) \ge 0 \text{ for all } y \in \mathscr{Y}$$

$$\|s^{\mathsf{T}} \lambda\|_{2} \le \nu$$
(8)

where  $\hat{\gamma} = (\hat{\gamma}_k)_{k=1,\dots,m}$  and  $s = (s_k)_{k=1,\dots,m}$ , and similarly for (7) (by considering  $-\max\{-\rho(\pi)\}$ ). Strong duality holds if  $\gamma(q) - \gamma(\tilde{\pi}(x_0,\cdot))$  lies in the interior of  $\mathscr{U}$  and under appropriate topological assumptions on  $\mathscr{Y}$  (e.g., Shapiro 2001). Without the condition, weak duality still implies that (8) provides a conservative approximation for (6).

Note that (8) still has an infinite number of constraints, but for specific choices of f and g one can reduce (8) to finite-dimensional optimization problems. We focus here on the setting that  $\mathscr{Y} = [0, \mathrm{UB}] \subset \mathbb{R}$  where UB is the largest possible value of Y, and f(y) and  $g_k(y)$  are in the form  $y^r$  where r is a rational number. This allows us to reduce (8) to finite-dimensional semidefinite programs (SDP) by generalizing the technique in Bertsimas and Popescu (2005), which consider moment equalities only and does not contain the second order cone constraint in the dual formulation (8). For example, in the case that f(y) = y and  $g_k(y) = y^{k/\tilde{m}}$  for  $k = 1, 2, \dots, m$ , (8) becomes

$$\min_{\kappa, \nu \in \mathbb{R}, \lambda \in \mathbb{R}^{m}} \quad \kappa + (\gamma(\tilde{\pi}(x_{0}, \cdot)) + \hat{\gamma})^{\mathsf{T}}\lambda + \sqrt{\chi_{1-\alpha,m}^{2}}\nu$$
subject to
$$\kappa + \sum_{i=1}^{m} \lambda_{i} y^{i/\tilde{m}} - y \ge 0 \text{ for all } y \in [0, \mathsf{UB}]$$

$$\|s^{\mathsf{T}}\lambda\|_{2} \le \nu$$

$$(9)$$

By a change of variable  $y^{1/\tilde{m}} \to y$ , we can transform the algebraic constraints into positive semidefinite constraints, and get

$$\begin{aligned} \min_{\kappa,\nu\in\mathbb{R},\lambda\in\mathbb{R}^{m},U\in\mathbb{R}^{(m+1)\times(m+1)}} & \kappa + (\gamma(\tilde{\pi}(x_{0},\cdot)) + \hat{\gamma})^{\mathsf{T}}\lambda + \sqrt{\chi_{1-\alpha,m}^{2}}\nu \\ \text{subject to} & U \succeq 0 \\ & \sum\limits_{i,j:i+j=2l+1} u_{ij} = 0, \quad l=1,\cdots,m \\ & \sum\limits_{i,j:i+j=2l+2} u_{ij} = \sum\limits_{r=0}^{l} y_{r} {\tilde{m}-r \choose l-r} \mathsf{U}\mathsf{B}^{r/\tilde{m}}, \quad l=0,\cdots,m \\ & \|s^{\mathsf{T}}\lambda\|_{2} \leq \nu \end{aligned} \tag{10}$$

where  $y_0 = \kappa, y_i = \lambda_i \mathbf{1}_{i \neq \tilde{m}} + (\lambda_i - 1) \mathbf{1}_{i = \tilde{m}}, [U]_{ij} = u_{ij}$ . Similarly, for (7), the objective function can be written as  $-\max\{-\rho(\pi)\}$ . The dual formulation of  $\max\{-\rho(\pi)\}$  remains as (10) except that  $y_{\tilde{m}}$  is replaced by  $\lambda_{\tilde{m}} + 1$ .

We choose the power in  $g_k(\cdot)$  as  $k/\tilde{m}$  instead of integers so as to avoid a blow-up in the magnitude of the moment as k increases. Note that there are other plausible choices of statistics, e.g., quantile-type statistics in the form  $g_k(y) = I(y \le b_k)$  for some  $b_k$ , or a combination of these and the power moments, which can also be similarly converted into SDP.

As in Section 2, an alternative is to consider the regression

$$g_k(Y_{ij}) = \beta_k \gamma_k(\tilde{\pi}(i,\cdot)) + h_k(i) + \varepsilon_{ijk}, \tag{11}$$

with responses  $g_k(Y_{ij})$  and additional scaling parameters  $\beta_k$ . We can form a CR  $\mathscr{U}$  for  $\gamma(\pi(x_0,\cdot))$  in this case as (5), where  $\hat{\gamma}_k$  is now the point estimate of  $\gamma_k(\pi(x_0,\cdot))$  under the regression and  $s_k$  is the corresponding standard error. Then we use the optimization problems

max 
$$\rho(q)$$
 subject to  $\gamma(q) \in \mathcal{U}$ 

and

min 
$$\rho(q)$$
 subject to  $\gamma(q) \in \mathcal{U}$ 

The maximization can be dualized to (8) and in the special case considered above to (10) (without the  $\gamma(\tilde{\pi}(x_0,\cdot))$  term). Similar treatments for the minimization as discussed above also hold.

We note that the accuracy of our proposed approach, here and in Section 2, depends on the complexity of the simulation model and the reality, in that it determines how challenging it is to find an acceptable regression model for capturing the model discrepancy. In the next section, we will test our approach with some numerical studies on a simple queueing system.

# 4 NUMERICAL STUDY

We use a queueing example to test the approaches we propose in Sections 2 and 3. We first consider a simulation model and the "real" system generated as follows: The simulation model is a M/M/x queue with arrival rate  $\lambda = 5$  and service rate  $\mu = 0.05$ . The real system is a M/M/x queue with a random arrival rate as an absolute value of a normally distributed random variable with mean 5 and standard deviation  $5 \cdot 0.1$ , and a random service rate as an absolute value of a normally distributed random variable with mean 0.045 and standard deviation  $0.045 \cdot 0.1$  (Think of these as, e.g., some random daily characteristics). In addition, each customer has a probability 0.2 to abandon the queue if the waiting time is larger than an exponential random variable with mean  $\frac{450}{x}$  where x is the number of servers. The output quantity of interest Y is the average waiting time among the first 50 customers.

The design point x corresponds to the number of servers, which ranges on  $x = 10, 11, \dots, 25$ . For this example, we can easily run plenty of simulation, say  $10^5$  replications, to get  $\tilde{\pi}$  that can be considered as

exactly obtained. For each  $x = 10, \dots, 24$ , we generate 10 observations on the real system. We deliberately collect no observations for x = 25 as a test point.

We use f(y) = y and  $g_k(y) = y^{k/\tilde{m}}$ , so that we can use the SDP (10) when adopting the optimization-enhanced approaches. We use a confidence level  $1 - \alpha = 95\%$  in the experiment.

We consider five methods:

• Data-only method (Method D) only utilizes real data  $Y_{ij}$ . For i = 1, ..., 24, the CI is

$$\left[\bar{Y}_i - z_{1-\alpha/2} \frac{\hat{\sigma}_i}{\sqrt{10}}, \bar{Y}_i + z_{1-\alpha/2} \frac{\hat{\sigma}_i}{\sqrt{10}}\right],$$

where  $\bar{Y}_i = \frac{\sum\limits_{j=1}^{10} Y_{ij}}{10}$ ,  $\hat{\sigma}_i = \sqrt{\frac{\sum\limits_{i=1}^{10} (Y_{ij} - \bar{Y}_i)^2}{10 - 1}}$ , and  $z_{1-\alpha/2}$  is the  $(1 - \alpha/2)^{th}$  quantile of a standard normal distribution. For x = 25, there is no data so Method D does not apply.

- Regression-only method 1 (Method R1) is based on the regression problem (2) in Section 2.
- Optimization-enhanced method 1 (Method O1) is based on optimizing over the CR obtained from (4) in Section 3.
- Regression-only method 2 (Method R2) is based on (3) in Section 2.
- Optimization-enhanced method 2 (Method O2) is based on optimizing over the CR obtained from (11) in Section 3.
- Simulation-only method (Method S) uses the (accurate) point estimate of our simulation model.

Each method except Method S outputs a lower bound and a upper bound value for each experiment, thereby forming a  $(1-\alpha)$ -level confidence interval. For Methods R1, O1, R2, and O2, we use the R function "npreg" in the "np" package to make local linear regressions.

# 4.1 Comparisons with Data-only and Simulation-only Methods

We first investigate the coverage probabilities of the prediction bounds generated by the methods. Table 1 shows the estimated coverage probabilities from repetition of 100 experiments with  $g_k(y) = y^{k/10}, k = 1, \dots, 10$ . For each entry in the table, we use  $a \pm b$  to denote the CI of the estimated coverage probability from the 100 experiments, i.e., a is the average number of success in covering the true value and b is the half-width of the associated CI for the coverage probability.

For Method D, the coverage probability ranges from 0.816 to 0.999. It is roughly around the theoretical coverage probability 95% but there are mild fluctuations due to the small number of repetitions. That being said, both regression-only and optimization-enhanced methods obtain more fluctuated values as the number of servers varies. For example, the coverage probability ranges from 0.372 to 0.999 for Method R1 and ranges from 0.493 to 1 for Method O1. Therefore, the data-only method shows a more reliable performance overall. Even though in some cases (e.g., when x = 20) the correction methods (i.e., Methods R1, O1, R2, O2) all obtain reasonable coverages, the data-only method has uniformly better performances and is the best in other cases (e.g., x = 10). However, if one considers the case x = 25 (i.e., the extreme case where there is no data), then Method D does not apply, but the correction methods still give reasonable coverage (though the performances vary). That the correction methods tend to have poorer coverage probabilities could be due to the crudeness of the regression models, whose assumptions and post-hoc analyses have not been carefully conducted in this experiment. In other words, the regression models used in this example could deviate from the true behavior of  $\rho(\delta(x,\cdot))$  or  $\rho(\pi(x,\cdot))$ .

# Lam, Plumlee, and Zhang

Table 1: Coverage probabilities for different methods.

| Number  | of | Method D           | Method R1         | Method O1        | Method R2          | Method O2        |
|---------|----|--------------------|-------------------|------------------|--------------------|------------------|
| servers |    |                    |                   |                  |                    |                  |
| 10      |    | $0.88 \pm 0.064$   | $0.47 \pm 0.098$  | $0.59\pm0.097$   | $0.57\pm0.098$     | $0.66\pm0.093$   |
| 11      |    | $0.91 \pm 0.056$   | $0.59 \pm 0.097$  | $0.63 \pm 0.095$ | $0.67 \pm 0.093$   | $0.77 \pm 0.083$ |
| 12      |    | $0.96 \pm 0.039$   | $0.67 \pm 0.093$  | $0.78 \pm 0.082$ | $0.81 \pm 0.077$   | $0.88 \pm 0.064$ |
| 13      |    | $0.88 {\pm} 0.064$ | $0.77 \pm 0.083$  | $0.83 \pm 0.074$ | $0.90 \pm 0.059$   | $0.93 \pm 0.050$ |
| 14      |    | $0.89 \pm 0.062$   | $0.76 \pm\ 0.084$ | $0.79 \pm 0.080$ | $0.90 \pm 0.059$   | $0.90 \pm 0.059$ |
| 15      |    | $0.91 \pm 0.056$   | $0.79 \pm 0.080$  | $0.78 \pm 0.082$ | $0.98 {\pm} 0.028$ | $0.96\pm0.039$   |
| 16      |    | $0.93 \pm 0.050$   | $0.85 \pm\ 0.070$ | $0.91 \pm 0.056$ | $0.94 \pm 0.047$   | $0.96 \pm 0.039$ |
| 17      |    | $0.88 {\pm} 0.064$ | $0.91 \pm\ 0.056$ | $0.95 \pm 0.043$ | $0.98{\pm}0.028$   | $0.97 \pm 0.034$ |
| 18      |    | $0.93 \pm 0.050$   | $0.93 \pm 0.050$  | $0.91 \pm 0.056$ | $0.95 \pm 0.043$   | $0.93 \pm 0.050$ |
| 19      |    | $0.92 \pm 0.053$   | $0.95 \pm\ 0.043$ | $0.97 \pm 0.034$ | $0.90 \pm 0.059$   | $0.94\pm0.047$   |
| 20      |    | $0.92 \pm 0.053$   | $0.96 \pm\ 0.039$ | $0.96 \pm 0.039$ | $0.99 \pm 0.020$   | $0.97 \pm 0.034$ |
| 21      |    | $0.95 \pm 0.043$   | $0.96 \pm\ 0.039$ | $0.97 \pm 0.034$ | $0.95 \pm 0.043$   | $0.94 \pm 0.047$ |
| 22      |    | $0.95 \pm 0.043$   | $0.90 \pm\ 0.059$ | $0.89 \pm 0.062$ | $0.99 \pm 0.020$   | $0.95 \pm 0.043$ |
| 23      |    | $0.91 \pm 0.056$   | $0.89 \pm\ 0.062$ | $0.88 \pm 0.064$ | $0.99 \pm 0.020$   | $0.95 \pm 0.043$ |
| 24      |    | $0.95 \pm 0.043$   | $0.89 \pm\ 0.062$ | $0.92 \pm 0.053$ | $0.87 \pm 0.066$   | $0.87 \pm 0.066$ |
| 25      |    | -                  | $0.81 \pm\ 0.077$ | $0.70\pm0.090$   | $0.87 {\pm} 0.066$ | $0.73 \pm 0.087$ |

Besides coverage probability, we are interested in gaining some understanding on how conservative the methods are. This can be measured by the difference between the lower and upper confidence bounds for each method. We call these differences the prediction gaps. Table 2 shows the statistics of this gap for each of the methods, among the 100 experiments with  $g_k(y) = y^{k/10}, k = 1, \dots, 10$ . For each entry in the table except the ones in the last column, we use  $a \pm b$  to denote the CI of the prediction gap from the 100 experiments (i.e., a is the average gap, and b is the CI half-width). The last column shows the absolute difference between the point estimate of the simulation model and the truth, which can be viewed as the magnitude of the model error made in using the simulation.

Table 2: Prediction gaps for different methods.

| Number  | of | Method D          | Method R1         | Method O1         | Method R2         | Method O2         | Method S |
|---------|----|-------------------|-------------------|-------------------|-------------------|-------------------|----------|
| servers |    |                   |                   |                   |                   |                   |          |
| 10      |    | $7.488 \pm 0.355$ | $1.820 \pm 0.263$ | $2.035 \pm 0.242$ | $2.660 \pm 0.205$ | $2.943 \pm 0.171$ | 0.413    |
| 11      |    | $6.005 \pm 0.346$ | $1.561 \pm 0.190$ | $1.670 \pm 0.155$ | $2.631 \pm 0.176$ | $2.884 \pm 0.161$ | 0.348    |
| 12      |    | $5.527 \pm 0.275$ | $1.325 \pm 0.150$ | $1.486 \pm 0.127$ | $2.322 \pm 0.125$ | $2.489 \pm 0.094$ | 0.236    |
| 13      |    | $4.751 \pm 0.249$ | $1.315 \pm 0.140$ | $1.467 \pm 0.133$ | $2.127 \pm 0.085$ | $2.360 \pm 0.077$ | 0.191    |
| 14      |    | $4.490 \pm 0.220$ | $1.334 \pm 0.134$ | $1.371 \pm 0.103$ | $1.979 \pm 0.040$ | $2.188 \pm 0.050$ | 0.137    |
| 15      |    | $3.878 \pm 0.208$ | $1.172 \pm 0.087$ | $1.247 \pm 0.083$ | $1.890 \pm 0.031$ | $2.035 \pm 0.041$ | 0.041    |
| 16      |    | $3.582 \pm 0.177$ | $1.143 \pm 0.087$ | $1.184 \pm 0.086$ | $1.797 \pm 0.030$ | $1.933 \pm 0.030$ | 0.027    |
| 17      |    | $2.987 \pm 0.144$ | $1.038 \pm 0.053$ | $1.060 \pm 0.049$ | $1.713 \pm 0.037$ | $1.808 \pm 0.031$ | 0.018    |
| 18      |    | $2.791 \pm 0.138$ | $1.000 \pm 0.057$ | $0.981 \pm 0.050$ | $1.627 \pm 0.051$ | $1.679 \pm 0.038$ | 0.059    |
| 19      |    | $2.533 \pm 0.125$ | $0.981 \pm 0.046$ | $0.973 \pm 0.047$ | $1.527 \pm 0.065$ | $1.557 \pm 0.046$ | 0.093    |
| 20      |    | $2.274 \pm 0.108$ | $0.958 \pm 0.043$ | $0.945 \pm 0.045$ | $1.579 \pm 0.067$ | $1.470 \pm 0.047$ | 0.101    |
| 21      |    | $2.112 \pm 0.098$ | $0.906 \pm 0.026$ | $0.886 \pm 0.029$ | $1.489 \pm 0.075$ | $1.383 \pm 0.054$ | 0.140    |
| 22      |    | $1.831 \pm 0.086$ | $0.893 \pm 0.025$ | $0.849 \pm 0.028$ | $1.452 \pm 0.081$ | $1.306 \pm 0.045$ | 0.136    |
| 23      |    | $1.673 \pm 0.088$ | $0.885 \pm 0.024$ | $0.832 \pm 0.030$ | $1.522 \pm 0.082$ | $1.272 \pm 0.048$ | 0.143    |
| 24      |    | $1.501 \pm 0.077$ | $0.898 \pm 0.019$ | $0.806 \pm 0.026$ | $1.462 \pm 0.084$ | $1.165 \pm 0.050$ | 0.164    |
| 25      |    | _                 | $0.999 \pm 0.084$ | $0.871 \pm 0.078$ | $1.514 \pm 0.080$ | $1.157 \pm 0.067$ | 0.162    |

Table 2 shows that the prediction gaps using correction methods are significantly shorter than that using data-only method for every x. For example, at x = 18, Method D has a prediction gap about 2.791, but Method R1 has 1 and Method O1 has 0.981. Moreover, at some particular x, both high coverage probability and small prediction gap are obtained simultaneously for correction methods. For instance, the coverage probabilities at x = 18 for both data-only and correction methods are roughly around 95% and the probabilities only differ by  $5\% \sim 10\%$  for some other choices of x such as 15, 16, 23, 24. This shows that, while the coverage probabilities are more reliable when using only data, their predictions are generally

much more conservative. The correction methods use information on the simulation model to reduce the conservativeness.

Considering the last column in Table 2, we see that the simulation model errors are in overall smaller than the prediction gaps in this example. However, simulation method in this example is, in a sense, always erroneous as the Monte Carlo errors of an inaccurate model are washed away with abundant simulation runs, leaving point estimates that are different from the truth. We note that different simulation model will give different model error, and there is no direct guarantee whether the model error as we define will be smaller or larger than the prediction gaps from the correction methods.

# 4.2 Comparisons Among Regression-only and Optimization-enhanced Methods

Next we repeat the experiment, using a smaller number, but higher orders, of moments  $g_k(y) = y^{k/2}$ ,  $k = 1, \dots, 3$ . Tables 3 and 4 give the coverage probabilities and prediction gaps respectively under this new scheme. The data are all newly generated in this experiment, and hence the results for Method D differ slightly from the previous experiment. We also realize that the R function "npreg" for local linear regression seems to give slightly different prediction value each time even with identical input data, possibly due to numerical issues. For regression-only method R1, these fluctuations are negligible and consequently the difference of results with the previous experiment mainly arises from the random data generation. For method R2, the changes range from 0.01 to 0.05 in both the coverage probability and the prediction gap, which arises from a combination of random data generation and the R function "npreg".

Table 3 shows that the coverage probabilities for regression-only methods and optimization-enhanced methods are comparable to each other for every x. For example, for x=16, the coverage probabilities differ by only roughly 0.02. Table 4 shows that the prediction gaps obtained from optimization-enhanced methods are generally shorter than that from regression-only methods. For instance, at x=14, Method R1 has a prediction gap about 1.223 and Method O1 has 1.142, whereas at say x=19 Method R2 has 1.670 and Method O2 has 1.478. Besides the shorter gaps, optimization-enhanced methods give prediction gaps that have smaller variability. For example, the CI length of the prediction gap is shorter for Method O2 than for Method R2 for almost every x.

Table 3: Coverage probabilities for different methods under a different set of moments.

| Number of servers | of | Method D        | Method R1       | Method O1       | Method R2       | Method O2       |
|-------------------|----|-----------------|-----------------|-----------------|-----------------|-----------------|
| 10                |    | $0.89 \pm 0.06$ | $0.39 \pm 0.10$ | $0.41 \pm 0.10$ | $0.54 \pm 0.10$ | $0.47 \pm 0.10$ |
| 11                |    | $0.93 \pm 0.05$ | $0.59 \pm 0.10$ | $0.59 \pm 0.10$ | $0.76 \pm 0.08$ | $0.71 \pm 0.09$ |
| 12                |    | $0.87 \pm 0.07$ | $0.70 \pm 0.09$ | $0.66 \pm 0.09$ | $0.87 \pm 0.07$ | $0.86 \pm 0.07$ |
| 13                |    | $0.91 \pm 0.06$ | $0.79 \pm 0.08$ | $0.75 \pm 0.09$ | $0.92 \pm 0.05$ | $0.90 \pm 0.06$ |
| 14                |    | $0.97 \pm 0.03$ | $0.82\pm0.08$   | $0.79 \pm 0.08$ | $0.94 \pm 0.05$ | $0.86 \pm 0.07$ |
| 15                |    | $0.90 \pm 0.06$ | $0.82\pm0.08$   | $0.79\pm0.08$   | $0.95\pm0.04$   | $0.93 \pm 0.05$ |
| 16                |    | $0.92 \pm 0.05$ | $0.88\pm0.06$   | $0.86 \pm 0.07$ | $0.97 \pm 0.03$ | $0.95 \pm 0.04$ |
| 17                |    | $0.92\pm0.05$   | $0.89\pm0.06$   | $0.85\pm0.07$   | $0.95\pm0.04$   | $0.94 \pm 0.05$ |
| 18                |    | $0.94 \pm 0.05$ | $0.99\pm0.02$   | $0.99\pm0.02$   | $0.95\pm0.04$   | $0.92\pm0.05$   |
| 19                |    | $0.90 \pm 0.06$ | $0.98 \pm 0.03$ | $0.98 \pm 0.03$ | $0.99 \pm 0.02$ | $0.97 \pm 0.03$ |
| 20                |    | $0.84 \pm 0.07$ | $0.93\pm0.05$   | $0.94 \pm 0.05$ | $0.99\pm0.02$   | $0.98 \pm 0.03$ |
| 21                |    | $0.93 \pm 0.05$ | $0.97 \pm 0.03$ | $0.97 \pm 0.03$ | $0.95\pm0.04$   | $0.94 \pm 0.05$ |
| 22                |    | $0.89\pm0.06$   | $0.93\pm0.05$   | $0.93 \pm 0.05$ | $0.98\pm0.03$   | $0.97 \pm 0.03$ |
| 23                |    | $0.87 \pm 0.07$ | $0.93 \pm 0.05$ | $0.93 \pm 0.05$ | $0.98 \pm 0.03$ | $0.97 \pm 0.03$ |
| 24                |    | $0.94 \pm 0.05$ | $0.83\pm0.07$   | $0.87\pm0.07$   | $0.95\pm0.04$   | $0.95 \pm 0.04$ |
| 25                |    | _               | $0.74\pm0.09$   | $0.68\pm0.01$   | $0.77\pm0.08$   | $0.78\pm0.01$   |

| TC 1 1 1 TO 1' 1'    | C 1'CC 4             | 41 1 1           | 1'CC 4 4           | • .      |
|----------------------|----------------------|------------------|--------------------|----------|
| Table /I: Prediction | gane for different   | methods under a  | different cet of   | momente  |
| Table 4: Prediction  | gains for difficient | incurvas unaci a | i uniforchi sei oi | moments. |
|                      |                      |                  |                    |          |

| Number  | of | Method D          | Method R1         | Method O1         | Method R2         | Method O2         | Method S |
|---------|----|-------------------|-------------------|-------------------|-------------------|-------------------|----------|
| servers |    |                   |                   |                   |                   |                   |          |
| 10      |    | $7.165 \pm 0.390$ | $1.753 \pm 0.238$ | $1.581 \pm 0.190$ | $2.583 \pm 0.176$ | $2.377 \pm 0.149$ | 0.413    |
| 11      |    | $6.156 \pm 0.293$ | $1.498 \pm 0.187$ | $1.345 \pm 0.132$ | $2.314 \pm 0.150$ | $2.079 \pm 0.124$ | 0.348    |
| 12      |    | $5.544 \pm 0.285$ | $1.358 \pm 0.157$ | $1.288 \pm 0.129$ | $2.304 \pm 0.115$ | $2.074 \pm 0.097$ | 0.236    |
| 13      |    | $4.850 \pm 0.260$ | $1.295 \pm 0.122$ | $1.202 \pm 0.098$ | $2.146 \pm 0.084$ | $1.950 \pm 0.068$ | 0.191    |
| 14      |    | $4.468 \pm 0.208$ | $1.223 \pm 0.117$ | $1.142 \pm 0.095$ | $2.048 \pm 0.056$ | $1.812 \pm 0.041$ | 0.137    |
| 15      |    | $3.729 \pm 0.162$ | $1.229 \pm 0.097$ | $1.150 \pm 0.077$ | $1.885 \pm 0.026$ | $1.665 \pm 0.026$ | 0.041    |
| 16      |    | $3.637 \pm 0.176$ | $1.104 \pm 0.082$ | $1.041 \pm 0.064$ | $1.803 \pm 0.030$ | $1.583 \pm 0.034$ | 0.027    |
| 17      |    | $3.145 \pm 0.184$ | $1.082 \pm 0.079$ | $1.018 \pm 0.060$ | $1.754 \pm 0.031$ | $1.585 \pm 0.039$ | 0.018    |
| 18      |    | $2.920 \pm 0.138$ | $1.048 \pm 0.066$ | $0.999 \pm 0.050$ | $1.659 \pm 0.047$ | $1.492 \pm 0.051$ | 0.059    |
| 19      |    | $2.563 \pm 0.126$ | $0.977 \pm 0.044$ | $0.941 \pm 0.035$ | $1.670 \pm 0.061$ | $1.478 \pm 0.061$ | 0.093    |
| 20      |    | $2.244 \pm 0.120$ | $0.962 \pm 0.039$ | $0.922 \pm 0.029$ | $1.523 \pm 0.065$ | $1.362 \pm 0.064$ | 0.101    |
| 21      |    | $2.053 \pm 0.102$ | $0.917 \pm 0.031$ | $0.904 \pm 0.026$ | $1.557 \pm 0.076$ | $1.383 \pm 0.071$ | 0.140    |
| 22      |    | $1.879 \pm 0.089$ | $0.889 \pm 0.025$ | $0.875 \pm 0.024$ | $1.526 \pm 0.083$ | $1.373 \pm 0.075$ | 0.136    |
| 23      |    | $1.690 \pm 0.077$ | $0.867 \pm 0.019$ | $0.873 \pm 0.025$ | $1.523\pm0.082$   | $1.376 \pm 0.073$ | 0.143    |
| 24      |    | $1.565\pm0.072$   | $0.909 \pm 0.026$ | $0.897 \pm 0.028$ | $1.532\pm0.078$   | $1.378 \pm 0.071$ | 0.164    |
| 25      |    | _                 | $1.017 \pm 0.066$ | $0.966 \pm 0.005$ | $1.519 \pm 0.080$ | $1.345 \pm 0.008$ | 0.162    |

In comparing different regression models (2) and (3), we see that correction methods based on (3) generally give higher coverage probabilities compared with those based on (2). For example, at x=12, Method R1 obtains coverage probability  $0.67\pm0.093$  and Method R2 obtains  $0.81\pm0.077$  in Table 1, and  $0.70\pm0.09$  and  $0.87\pm0.07$  respectively in Table 3. Likewise, for optimization-enhanced methods, Method O1 obtains coverage probability  $0.78\pm0.082$  and Method O2 obtains  $0.88\pm0.064$  in Table 1, and  $0.66\pm0.09$  and  $0.86\pm0.07$  respectively in Table 3. However, higher coverage probabilities seem to come with the price of wider prediction gaps. For instance, Method R1 gives prediction gaps  $1.325\pm0.015$  and Method R2 gives  $2.322\pm0.0125$  in Table 2, and  $1.358\pm0.157$  and  $2.304\pm0.115$  respectively in Table 4. Likewise, for optimization-enhanced methods, Method O1 gives prediction gaps  $1.486\pm0.127$  and Method O2 gives  $2.489\pm0.094$  in Table 2, and  $1.288\pm0.129$  and  $2.074\pm0.097$  respectively in Table 4. From these observations, there does not seem to be a concrete conclusion which regression model is better. One who prefers a higher coverage probability at the cost of a conservative prediction gap may plausibly choose the second regression model, and vice versa.

Next we analyze the sensitivity to the choice of moments in the constraints. We choose four sets of parameters as follows:  $(m, \tilde{m}) = (3,3), (6,6), (10,10), (3,2)$ . We classify two categories in the choices of moment constraints with given functions  $f(y) = y, g_k(y) = y^{k/\tilde{m}}, k = 1, \cdots, m$ . The first category is  $m \le \tilde{m}$  and the second category is  $m > \tilde{m}$ . Hence,  $(m, \tilde{m}) = (3,3), (6,6), (10,10)$  are under the first category while  $(m, \tilde{m}) = (3,2)$  belongs to the second category. In the first category, the order of moments in the constraints does not exceed that of the objective function. On the other hand, the second category allows higher-order moments in the constraints. In this experiment, we focus on Methods R1 and O1.

Tables 5 and 6 give the coverage probabilities and prediction gaps respectively under the above four sets of parameters. Table 5 shows that optimization-enhanced methods under the first category generally give higher coverage probabilities than regression-only methods, while vice versa under the second category. For example, at x = 16, Method R1 outputs coverage probability 0.85. Under the first category,  $(m, \tilde{m}) = (3,3), (6,6), (10,10)$  give coverage probabilities 0.88,0.88,0.87 respectively. Under the second category, namely  $(m, \tilde{m}) = (3,2)$ , the coverage probability is 0.83. Similar to an observation made before, lower coverage probability generally comes with a shorter prediction gap. Table 6 shows that optimization-enhanced methods under the first category generally give wider prediction gaps than regression-only methods, and vice versa under the second category. For instance, at x = 16, Method R1 outputs a prediction gap 1.214. Under the first category,  $(m, \tilde{m}) = (3,3), (6,6), (10,10)$  output prediction gaps 1.327, 1.306, 1.273 respectively. Under the second category, namely  $(m, \tilde{m}) = (3,2)$ , the prediction gap is 1.122. We attribute such phenomenon to a stronger impact from moments that are higher order, i.e., in the second category, on constraining the possible values of the performance measure, which shrinks

the prediction gap. Under the first category, the effect of moments is weaker, and the combined additional statistical errors in estimating these moments lead to wider prediction gaps than regression-only methods.

Table 5: Coverage probabilities for different methods and different settings.

| Number  | of | Method D        | Method R1       | Method O1           | Method O1           | Method O1            | Method O1         |
|---------|----|-----------------|-----------------|---------------------|---------------------|----------------------|-------------------|
| servers |    |                 |                 | $m = \tilde{m} = 3$ | $m = \tilde{m} = 6$ | $m = \tilde{m} = 10$ | $m=3 \tilde{m}=2$ |
| 10      |    | $0.91 \pm 0.06$ | $0.45 \pm 0.10$ | $0.56 \pm 0.10$     | $0.52 \pm 0.10$     | $0.50 \pm 0.10$      | $0.40 \pm 0.10$   |
| 11      |    | $0.91 \pm 0.05$ | $0.59\pm0.10$   | $0.67 \pm 0.09$     | $0.67 \pm 0.09$     | $0.66 \pm 0.09$      | $0.58\pm0.10$     |
| 12      |    | $0.87 \pm 0.07$ | $0.7 \pm 0.09$  | $0.71 \pm 0.09$     | $0.70 \pm 0.09$     | $0.69 \pm 0.09$      | $0.65 \pm 0.09$   |
| 13      |    | $0.92\pm0.06$   | $0.72\pm0.09$   | $0.73 \pm 0.09$     | $0.74 \pm 0.09$     | $0.72\pm0.09$        | $0.66\pm0.09$     |
| 14      |    | $0.93 \pm 0.03$ | $0.76\pm0.08$   | $0.80\pm0.08$       | $0.79 \pm 0.08$     | $0.78\pm0.08$        | $0.72\pm0.09$     |
| 15      |    | $0.92\pm0.06$   | $0.81\pm0.08$   | $0.84 \pm 0.07$     | $0.83\pm0.07$       | $0.81\pm0.08$        | $0.79\pm0.08$     |
| 16      |    | $0.91 \pm 0.05$ | $0.85\pm0.07$   | $0.88 \pm 0.06$     | $0.88\pm0.06$       | $0.87 \pm 0.07$      | $0.83 \pm 0.07$   |
| 17      |    | $0.94\pm0.05$   | $0.89\pm0.06$   | $0.92 \pm 0.05$     | $0.91 \pm 0.06$     | $0.90\pm0.06$        | $0.87\pm0.07$     |
| 18      |    | $0.94\pm0.05$   | $0.95\pm0.04$   | $0.97 \pm 0.03$     | $0.96\pm0.04$       | $0.96\pm0.04$        | $0.92\pm0.05$     |
| 19      |    | $0.9\pm0.06$    | $0.94 \pm 0.05$ | $0.98 \pm 0.03$     | $0.97 \pm 0.03$     | $0.97 \pm 0.03$      | $0.94 \pm 0.05$   |
| 20      |    | $0.91 \pm 0.07$ | $0.97 \pm 0.03$ | $0.97 \pm 0.03$     | $0.95\pm0.04$       | $0.95\pm0.04$        | $0.92\pm0.05$     |
| 21      |    | $0.94\pm0.05$   | $0.97\pm0.03$   | $0.95\pm0.04$       | $0.94\pm0.05$       | $0.94\pm0.05$        | $0.92\pm0.05$     |
| 22      |    | $0.91\pm0.06$   | $0.93 \pm 0.05$ | $0.96 \pm 0.04$     | $0.96\pm0.04$       | $0.96\pm0.04$        | $0.94\pm0.05$     |
| 23      |    | $0.95\pm0.07$   | $0.92\pm0.05$   | $0.93 \pm 0.05$     | $0.93 \pm 0.05$     | $0.93 \pm 0.05$      | $0.93 \pm 0.05$   |
| 24      |    | $0.93\pm0.05$   | $0.89\pm0.06$   | $0.92\pm0.05$       | $0.88\pm0.06$       | $0.86\pm0.07$        | $0.91\pm0.06$     |
| 25      |    | _               | $0.73 \pm 0.09$ | $0.84 \pm 0.01$     | $0.81 \pm 0.01$     | $0.76 \pm 0.01$      | $0.73 \pm 0.01$   |

Table 6: Prediction gaps for different methods and different settings.

| Number<br>servers | of | Method D          | Method R1                           | Method O1 $m = \tilde{m} = 3$          | Method O1 $m = \tilde{m} = 6$ | Method O1 $m = \tilde{m} = 10$      | Method O1 $m = 3 \tilde{m} = 2$     | Method S |
|-------------------|----|-------------------|-------------------------------------|--|-------------------------------|-------------------------------------|-------------------------------------|----------|
| 10                |    | $7.375 \pm 0.396$ | $1.993 \pm 0.252$                   | $\frac{m-m-3}{2.207+0.247}$            | $\frac{m-m-0}{2.102+0.226}$   | $\frac{m-m-10}{2.049 \pm 0.218}$    | $\frac{m-3m-2}{1.714+0.188}$        | 0.413    |
| 11                |    | $6.154 \pm 0.337$ | $1.747 \pm 0.196$                   | $2.207 \pm 0.247$<br>$2.010 \pm 0.205$ | $1.980 \pm 0.189$             | $1.943 \pm 0.179$                   | $1.714 \pm 0.166$ $1.588 \pm 0.152$ | 0.348    |
| 12                |    | $5.470 \pm 0.267$ | $1.747 \pm 0.190$ $1.578 \pm 0.156$ | $1.797 \pm 0.159$                      | $1.781 \pm 0.150$             | $1.749 \pm 0.179$                   | $1.368 \pm 0.132$ $1.442 \pm 0.119$ | 0.236    |
| 13                |    | $4.815 \pm 0.237$ | $1.462 \pm 0.130$                   | $1.635 \pm 0.127$                      | $1.624 \pm 0.121$             | $1.749 \pm 0.144$ $1.589 \pm 0.117$ | $1.332 \pm 0.096$                   | 0.191    |
| 14                |    | $4.560 \pm 0.234$ | $1.370 \pm 0.115$                   | $1.522 \pm 0.114$                      | $1.506 \pm 0.109$             | $1.470 \pm 0.105$                   | $1.250 \pm 0.084$                   | 0.137    |
| 15                |    | $3.961 \pm 0.190$ | $1.282 \pm 0.095$                   | $1.411 \pm 0.094$                      | $1.389 \pm 0.090$             | $1.355 \pm 0.087$                   | $1.177 \pm 0.070$                   | 0.041    |
| 16                |    | $3.413 \pm 0.160$ | $1.214 \pm 0.084$                   | $1.327 \pm 0.085$                      | $1.306 \pm 0.083$             | $1.273 \pm 0.080$                   | $1.122 \pm 0.060$                   | 0.027    |
| 17                |    | $3.307 \pm 0.181$ | $1.143 \pm 0.075$                   | $1.245 \pm 0.082$                      | $1.213 \pm 0.075$             | $1.177 \pm 0.072$                   | $1.073 \pm 0.058$                   | 0.018    |
| 18                |    | $2.748 \pm 0.133$ | $1.071 \pm 0.059$                   | $1.160 \pm 0.067$                      | $1.125 \pm 0.063$             | $1.092 \pm 0.061$                   | $1.011 \pm 0.044$                   | 0.059    |
| 19                |    | $2.578 \pm 0.132$ | $1.011 \pm 0.048$                   | $1.088 \pm 0.054$                      | $1.056 \pm 0.051$             | $1.026 \pm 0.049$                   | $0.965 \pm 0.036$                   | 0.093    |
| 20                |    | $2.384 \pm 0.132$ | $0.956 \pm 0.034$                   | $1.017 \pm 0.037$                      | $0.980 \pm 0.035$             | $0.951 \pm 0.035$                   | $0.915 \pm 0.024$                   | 0.101    |
| 21                |    | $2.149 \pm 0.119$ | $0.920 \pm 0.027$                   | $0.981 \pm 0.033$                      | $0.948 \pm 0.032$             | $0.921 \pm 0.032$                   | $0.889 \pm 0.025$                   | 0.14     |
| 22                |    | $1.904 \pm 0.102$ | $0.900 \pm 0.026$                   | $0.950 \pm 0.030$                      | $0.912 \pm 0.031$             | $0.884 \pm 0.031$                   | $0.871 \pm 0.027$                   | 0.136    |
| 23                |    | $1.716 \pm 0.077$ | $0.885 \pm 0.020$                   | $0.926 \pm 0.024$                      | $0.886 \pm 0.024$             | $0.858 \pm 0.025$                   | $0.862\pm0.024$                     | 0.143    |
| 24                |    | $1.558 \pm 0.079$ | $0.907 \pm 0.020$                   | $0.921\pm0.030$                        | $0.870 \pm 0.031$             | $0.837 \pm 0.032$                   | $0.870 \pm 0.027$                   | 0.164    |
| 25                |    | _                 | $1.021\pm0.056$                     | $0.998 \pm 0.006$                      | $0.928 \pm 0.006$             | $0.894 \pm 0.006$                   | $0.946 \pm 0.005$                   | 0.162    |

Under the first category, the coverage probability and prediction gap is generally smaller as the number of moments increases for every x. For example, at x=14, as  $(m,\tilde{m})$  ranges from (3,3) to (10,10), the coverage probability decreases from  $0.80\pm0.08$  to  $0.78\pm0.08$ , and the prediction gap shortens from  $1.522\pm0.114$  to  $1.470\pm0.105$  respectively. However, this phenomenon may not always occur. In general, as the number of moments increases, we gain more related information about the true distribution, thereby shrinking the feasible region. For example, in setting  $(m,\tilde{m})=(6,6)$  we capture information not only about  $EY^{1/3}$ ,  $EY^{2/3}$ ,  $EY^{3/3}$  but also  $EY^{1/6}$ ,  $EY^{3/6}$ ,  $EY^{5/6}$ , the latter not utilized in the case  $(m,\tilde{m})=(3,3)$ . On the other hand, as we add more constraints, the inequalities on the moments may become looser due to more estimation requirements. More moments means a smaller feasible region while the additional estimation on the ranges of these moments enlarges the feasible region, so it remains unclear which factor dominates.

Finally, we compare the two categories by keeping the number of moments the same, i.e., compare  $(m, \tilde{m}) = (3,3)$  with (3,2). The second category seems to offer smaller coverage probability and prediction gap for every x. For instance, at x = 12, the second category has a coverage probability that is roughly 0.06 smaller than the first, and the prediction gap is roughly 0.355 shorter.

#### 5 CONCLUSION

We investigate an approach based on a combination of regression and optimization to learn and correct for the model discrepancy of a simulation model with real-world observations. Our approach targets at continuous output variables that are not directly handleable by previous work due to their infinite-dimensional nature. One version of our approach regresses the objective function of interest on the observed model discrepancies, or the observations themselves, against the design points. Another version further regresses the moments of the model discrepancies or observations and uses a moment optimization to compute bounds that account for the distributional nature of the outputs. We present the regression details and how to reformulate and subsequently solve the optimization in the second version via duality and semidefinite programming. Our experimental results show that the data-only approach could give better coverage probability in the availability of data, but could be more conservative in terms of wider confidence interval length than our correction methods. Moreover, the data-only approach is only applicable for historically observed design points but trivially breaks down in the new design points. Our results also show that our correction methods improve pure simulation models that are misspecified, in that the simulation-only estimates are systematically different from the truth while our correction methods generate interval estimates that cover or are closer to it.

#### **ACKNOWLEDGMENTS**

We gratefully acknowledge support from the National Science Foundation under grants CMMI-1542020, CMMI-1523453 and CAREER CMMI-1653339.

#### REFERENCES

- Aksin, Z., M. Armony, and V. Mehrotra. 2007. "The Modern Call Center: A Multi-Disciplinary Perspective on Operations Management Research". *Production and Operations Management* 16 (6): 665–688.
- Balci, O., and R. G. Sargent. 1982. "Some Examples of Simulation Model Validation Using Hypothesis Testing". In *Proceedings of the 1982 Winter Simulation Conference-Volume 2*, edited by H. J. Highland, O. Madrigal, and Y. W. Chao, 621–629. Piscataway, New Jersey: Institute of Electrical and Electronics Engineers, Inc.
- Bayarri, M. J., J. O. Berger, R. Paulo, J. Sacks, J. A. Cafeo, J. Cavendish, C.-H. Lin, and J. Tu. 2012. "A Framework for Validation of Computer Models". *Technometrics*.
- Ben-Tal, A., D. Den Hertog, A. De Waegenaere, B. Melenberg, and G. Rennen. 2013. "Robust Solutions of Optimization Problems Affected by Uncertain Probabilities". *Management Science* 59 (2): 341–357.
- Bertsimas, D., and I. Popescu. 2005. "Optimal Inequalities in Probability Theory: A Convex Optimization Approach". *SIAM Journal on Optimization* 15 (3): 780–804.
- Brown, L., N. Gans, A. Mandelbaum, A. Sakov, H. Shen, S. Zeltyn, and L. Zhao. 2005. "Statistical Analysis of a Telephone Call Center: A Queueing-Science Perspective". *Journal of the American statistical association* 100 (469): 36–50.
- Delage, E., and Y. Ye. 2010. "Distributionally Robust Optimization under Moment Uncertainty with Application to Data-Driven Problems". *Operations research* 58 (3): 595–612.
- Goh, J., and M. Sim. 2010. "Distributionally Robust Optimization and its Tractable Approximations". *Operations research* 58 (4-part-1): 902–917.
- Higdon, D., M. Kennedy, J. C. Cavendish, J. A. Cafeo, and R. D. Ryne. 2004. "Combining Field Data and Computer Simulations for Calibration and Prediction". *SIAM Journal on Scientific Computing* 26 (2): 448–466.
- Kennedy, M. C., and A. O'Hagan. 2001. "Bayesian Calibration of Computer Models". *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 63 (3): 425–464.
- Kleijnen, J. P. 2015. Design and Analysis of Simulation Experiments. Second ed, Volume 20. Springer.

- Kleijnen, J. P., and R. G. Sargent. 2000. "A Methodology for Fitting and Validating Metamodels in Simulation". *European Journal of Operational Research* 120 (1): 14–29.
- Kleijnen, J. P. C., R. C. H. Cheng, and B. Bettonvil. 2001. "Validation of Trace-Driven Simulation Models: Bootsrap Tests". *Management Science* 47 (11): 1533–1538.
- Livny, M., B. Melamed, and A. K. Tsiolis. 1993. "The Impact of Autocorrelation on Queuing Systems". *Management Science* 39 (3): 322–339.
- Oakley, J. E., and A. O'Hagan. 2004. "Probabilistic Sensitivity Analysis of Complex Models: A Bayesian Approach". *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 66 (3): 751–769.
- Pazgal, A. I., and S. Radas. 2008. "Comparison of Customer Balking and Reneging Behavior to Queueing Theory Predictions: An Experimental Study". *Computers & Operations Research* 35 (8): 2537–2548.
- Plumlee, M. 2016. "Bayesian Calibration of Inexact Computer Models". *Journal of the American Statistical Association* (just-accepted).
- Plumlee, M., and H. Lam. 2016. "Learning Stochastic Model Discrepancy". In *Proceedings of the 2016 Winter Simulation Conference*, edited by P. Frazier, T. Roeder, R. Szechtman, and E. Zhou, 413–424. Piscataway, New Jersey: Institute of Electrical and Electronics Engineers, Inc.
- Plumlee, M., and H. Lam. 2017. "An Uncertainty Quantification Method for Inexact Simulation Models". *arXiv preprint arXiv:1707.06544*.
- Sargent, G. R. 2013. "Verification and Validation of Simulation Models". *Journal of Simulation* 7 (1): 12–24.
- Shapiro, A. 2001. "On Duality Theory of Conic Linear Problems". In *Semi-infinite programming*, 135–165. Springer.
- Veeraraghavan, S., and L. Debo. 2009. "Joining Longer Queues: Information Externalities in Queue Choice". *Manufacturing & Service Operations Management* 11 (4): 543–562.
- Wiesemann, W., D. Kuhn, and M. Sim. 2014. "Distributionally Robust Convex Optimization". *Operations Research* 62 (6): 1358–1376.

#### **AUTHOR BIOGRAPHIES**

**HENRY LAM** is an Associate Professor in the Department of Industrial Engineering and Operations Research at Columbia University. He received his Ph.D. degree in statistics at Harvard University, and was a faculty member in the Department of Mathematics and Statistics at Boston University and the Department of Industrial and Operations Engineering at the University of Michigan before joining Columbia. His research focuses on stochastic simulation, risk analysis, and simulation optimization. His email address is khl2114@columbia.edu.

MATTHEW PLUMLEE is an Assistant Professor at the Department of Industrial Engineering and Management Sciences at Northwestern University. He received his Ph.D. degree from the School of Industrial and Systems Engineering at Georgia Tech, and was an Assistant Professor at the University of Michigan before joining Northwestern. His interests lie in the interface of data and modeling, specifically in methods for experimentation and uncertainty quantification for complex systems. His email address is matthew.plumlee@northwestern.edu.

XINYU ZHANG is a Ph.D. student in the Department of Industrial Engineering and Operations Research at Columbia University. He obtained his bachelor's degree in the University of Michigan, majoring in physics and applied mathematics, and was a Ph.D. student in the Department of Industrial and Operations Engineering there before joining Columbia. His research interests include stochastic optimization and extreme event analysis. His email address is zhang.xinyu@columbia.edu.