

# Data-Driven Prediction with Stochastic Data: Confidence Regions and Minimum Mean-Squared Error Estimates

Mingzhou Yin, Andrea Iannelli, and Roy S. Smith

**Abstract**—Recently, direct data-driven prediction has found important applications for controlling unknown systems, particularly in predictive control. Such an approach provides exact prediction using behavioral system theory when noise-free data are available. For stochastic data, although approximate predictors exist based on different statistical criteria, they fail to provide statistical guarantees of prediction accuracy. In this paper, confidence regions are provided for these stochastic predictors based on the prediction error distribution. Leveraging this, an optimal predictor which achieves minimum mean-squared prediction error is also proposed to enhance prediction accuracy. These results depend on some true model parameters, but they can also be replaced with an approximate data-driven formulation in practice. Numerical results show that the derived confidence region is valid and smaller prediction errors are observed for the proposed minimum mean-squared error estimate, even with the approximate data-driven formulation.

## I. INTRODUCTION

In dynamical system analysis, one of the fundamental problems is to predict system responses from given inputs and initial conditions. Conventionally, this is done by simulating a model of the system, derived from first principles and/or experimental data. However, increasing complexity of systems poses challenges to the modeling process. Direct approaches have therefore been widely pursued to obtain reliable predictions of system responses without an explicit model [1]. In what follows, the term ‘data-driven’ refers to such direct approaches.

A seminal result, known as the Willems’ fundamental lemma [2], shows that data-driven prediction can be conducted by linearly combining historical trajectory data with persistently exciting inputs for linear systems. A more general version of the lemma was recently given in [3]. This result enables model-based control design techniques to be adopted with direct data-driven formulations. This framework is especially suitable for predictive control, where multiple data-driven algorithms have been developed, including subspace predictive control [4], data-enabled predictive control [5], and behavioral input-output parametrization [6]. Successful applications have been described in [7], [8].

Recently, the extension of the fundamental lemma to stochastic data from a system identification point of view

has been drawing increasing interest [9]. Such work includes model predictive control based on the prediction error method [10], maximum likelihood signal matrix model [11], [12], and a Wasserstein distance minimization approach [7].

With stochastic data, both the historical trajectories and the prediction conditions are uncertain, which makes it difficult to obtain statistical guarantees of the predictors. This limits the application of data-driven predictors to control design, particularly when robustness requirements and safety constraints exist. As a result, to the best of our knowledge, existing work on robust data-driven control [13], [14], [15] is restricted to bounded noise models with often loose prediction error bounds.

In this paper, a statistical framework on the accuracy of the predicted response is established under the assumption of Gaussian output noise. With this framework, confidence regions are available for a general form of stochastic data-driven predictors. The confidence region depends on the extended observability matrix of the system, but it can also be approximated through a data-driven formulation of model properties without direct knowledge of model parameters. The validity of the derived confidence regions is verified by numerical examples.

In addition, this statistical framework allows computation of the mean-squared error (MSE) of the predictor. In this way, a novel stochastic data-driven predictor is designed to be optimal for prediction accuracy in terms of minimizing the MSE. This optimal algorithm can be obtained in practice with a data-driven model characterization. It is shown numerically that the proposed minimum MSE predictor obtains smaller prediction errors than existing stochastic predictors.

*Notation.* A Gaussian distribution with mean  $\mu$  and covariance  $\Sigma$  is indicated by  $\mathcal{N}(\mu, \Sigma)$ . The expectation and the covariance of a random vector  $x$  are denoted by  $\mathbb{E}(x)$  and  $\text{cov}(x)$  respectively. For a vector  $x$  and a positive definite matrix  $Q$ , the weighted Euclidean norm  $(x^T Q x)^{\frac{1}{2}}$  is denoted by  $\|x\|_Q$ . For a matrix  $X$ , the vectorization operator  $\text{vec}(X)$  stacks its columns in a single vector;  $X^\dagger$  indicates the Moore-Penrose pseudoinverse. For a sequence of matrices  $X_1, \dots, X_n$ , we denote  $[X_1^T \dots X_n^T]^T$  by  $\text{col}(X_1, \dots, X_n)$ .

## II. THE DATA-DRIVEN PREDICTION PROBLEM

### A. Problem Statement

Consider a discrete-time linear time-invariant (LTI) system with output noise, given by

$$\begin{cases} x_{t+1} &= Ax_t + Bu_t, \\ y_t &= Cx_t + Du_t + w_t, \end{cases} \quad (1)$$

This work was supported by the Swiss National Science Foundation under Grant 200021\_178890.

The authors are with the Automatic Control Laboratory, Swiss Federal Institute of Technology (ETH Zurich), Physikstrasse 3, 8092 Zurich, Switzerland, {myin, iannelli, rsmith}@control.ee.ethz.ch.

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

where  $x_t \in \mathbb{R}^{n_x}$ ,  $u_t \in \mathbb{R}^{n_u}$ ,  $y_t \in \mathbb{R}^{n_y}$ ,  $w_t \in \mathbb{R}^{n_y}$  are the states, inputs, outputs, and output noise respectively. In this paper, we assume that the system is observable with observability index (lag)  $l$ .

In data-driven prediction, the model parameters  $A, B, C, D$  are unknown, but  $M$  length- $L$  input-output trajectories

$$z_i^d = \text{col}\left(u_{t_i}^d, \dots, u_{t_i+L-1}^d, y_{t_i}^d, \dots, y_{t_i+L-1}^d\right) \in \mathbb{R}^{L(n_u+n_y)}, \quad (2)$$

where  $i = 0, \dots, M-1$ , have been collected. The matrix that concatenates these trajectories

$$Z = \begin{bmatrix} z_0^d & \dots & z_{M-1}^d \end{bmatrix} \in \mathbb{R}^{L(n_u+n_y) \times M} \quad (3)$$

is termed the signal matrix [11]. Depending on the construction, we can choose either  $t_{i+1} = t_i + 1$  for a mosaic Hankel signal matrix, or  $t_{i+1} = t_i + L$  for a Page signal matrix [16]. The trajectories can also come from independent experiments [17].

The problem is then to predict output trajectory  $\mathbf{y} = \text{col}(y_0, \dots, y_{L-1})$  from any given input trajectory  $\mathbf{u} = \text{col}(u_0, \dots, u_{L-1})$  using only the collected historical trajectories. To obtain a unique output trajectory, the initial condition is also fixed by measuring the immediate past input-output trajectory  $\mathbf{u}_{\text{ini}} = \text{col}(u_{-L_0}, \dots, u_{-1})$  and  $\mathbf{y}_{\text{ini}} = \text{col}(y_{-L_0}, \dots, y_{-1})$ , where  $L_0 = L - L' \geq l$ . In other words, the data-driven prediction problem aims to find an input-output mapping in the following form:

$$\mathbf{y} = \mathcal{F}_Z(\mathbf{u}; \mathbf{u}_{\text{ini}}, \mathbf{y}_{\text{ini}}). \quad (4)$$

### B. Noise-Free Data-Driven Prediction

In the noise-free case, the following lemma provides a condition for the existence of an exact mapping.

*Lemma 1:* If  $w_t = \mathbf{0}$ , the exact mapping in the form of (4) exists if  $\text{rank}(Z) = n_u L + n_x$ .

*Proof:* According to Corollary 19 in [3], if  $\text{rank}(Z) = n_u L + n_x$ , for all  $(\mathbf{u}_{\text{ini}}, \mathbf{u}, \mathbf{y}_{\text{ini}}, \mathbf{y})$ , there exists  $g \in \mathbb{R}^M$ , such that  $\text{col}(\mathbf{u}_{\text{ini}}, \mathbf{u}, \mathbf{y}_{\text{ini}}, \mathbf{y}) = Zg$ . Note that the observability index  $l$  satisfies  $n_y l \geq n_x$ . The dimension of  $\text{col}(\mathbf{u}_{\text{ini}}, \mathbf{u}, \mathbf{y}_{\text{ini}})$  then satisfies  $n_u L + n_y L_0 \geq \text{rank}(Z)$ , so  $\mathbf{y}$  can be uniquely determined by  $(\mathbf{u}_{\text{ini}}, \mathbf{u}, \mathbf{y}_{\text{ini}})$ . ■

Define a partition of  $Z$  as

$$Z = \text{col}(U_p, U_f, Y_p, Y_f), \quad (5)$$

where  $U_p \in \mathbb{R}^{n_u L_0 \times M}$ ,  $U_f \in \mathbb{R}^{n_u L' \times M}$ ,  $Y_p \in \mathbb{R}^{n_y L_0 \times M}$ ,  $Y_f \in \mathbb{R}^{n_y L' \times M}$ . Following the proof of Lemma 1, the mapping can be obtained by first solving the linear system

$$\text{col}(U_p, U_f, Y_p) g = \text{col}(\mathbf{u}_{\text{ini}}, \mathbf{u}, \mathbf{y}_{\text{ini}}), \quad (6)$$

and then applying  $\mathbf{y} = Y_f g$ . Although any solution to (6) is applicable (Proposition 1 in [18]), the pseudo-inverse solution is the most commonly used. So the solution to the noise-free data-driven prediction problem is,

$$\mathcal{F}_Z(\cdot) = Y_f g_{\text{pinv}}, \quad g_{\text{pinv}} = \begin{bmatrix} U_p \\ U_f \\ Y_p \end{bmatrix}^\dagger \begin{bmatrix} \mathbf{u}_{\text{ini}} \\ \mathbf{u} \\ \mathbf{y}_{\text{ini}} \end{bmatrix}. \quad (7)$$

### C. Data-Driven Prediction with Stochastic Data

When the output noise  $w_t$  is no longer zero but a realization of a stochastic process, Lemma 1 no longer holds and the mapping (1) can only be estimated approximately. The output noise leads to uncertainties in both the output signal matrix  $\text{col}(Y_p, Y_f)$  and the output initial condition  $\mathbf{y}_{\text{ini}}$ . In this paper, the distribution of  $w_t$  is assumed to be zero-mean Gaussian. Then, the distributions of  $\mathbf{y}_{\text{ini}}$  and  $\text{col}(Y_p, Y_f)$  are also Gaussian. In what follows, the distributions are denoted by

$$\begin{aligned} \mathbf{y}_{\text{ini}} &\sim \mathcal{N}(\mathbf{y}_{\text{ini}}^0, \Sigma_{\mathbf{y}_{\text{ini}}}), \\ \text{vec}\left(\begin{bmatrix} Y_p \\ Y_f \end{bmatrix}\right) &\sim \mathcal{N}\left(\text{vec}\left(\begin{bmatrix} Y_p^0 \\ Y_f^0 \end{bmatrix}\right), \Sigma_Y\right), \end{aligned} \quad (8)$$

where  $\mathbf{y}_{\text{ini}}^0$ ,  $Y_p^0$ , and  $Y_f^0$  are noise-free versions of  $\mathbf{y}_{\text{ini}}$ ,  $Y_p$ , and  $Y_f$  respectively, and  $\mathbf{y}_{\text{ini}}$  is uncorrelated with  $\text{col}(Y_p, Y_f)$ .

Under this assumption, for a given  $g$ , the distribution of

$$\begin{bmatrix} Y_p \\ Y_f \end{bmatrix} g = (g^\top \otimes \mathbb{I}_{n_y L}) \text{vec}\left(\begin{bmatrix} Y_p \\ Y_f \end{bmatrix}\right) \quad (9)$$

is thus

$$\begin{bmatrix} Y_p \\ Y_f \end{bmatrix} g \mid g \sim \mathcal{N}\left(\begin{bmatrix} Y_p^0 \\ Y_f^0 \end{bmatrix} g, \underbrace{\begin{bmatrix} \Sigma_p & \Sigma_{pf} \\ \Sigma_{pf}^\top & \Sigma_f \end{bmatrix}}_{\Sigma_g}\right), \quad (10)$$

where  $\Sigma_g = (g^\top \otimes \mathbb{I}_{n_y L}) \Sigma_Y (g \otimes \mathbb{I}_{n_y L})$ .

A special case of the noise model is when the noise is i.i.d. with  $w_t \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbb{I}_{n_y})$ , and the signal matrix  $Z$  is constructed as a Page matrix or from independent trajectories. In this case, we have  $\Sigma_{\mathbf{y}_{\text{ini}}} = \sigma^2 \mathbb{I}_{n_y L_0}$ ,  $\Sigma_Y = \sigma^2 \mathbb{I}_{n_y L M}$ , and thus  $\Sigma_g = \sigma^2 \|g\|_2^2 \mathbb{I}_{n_y L}$ .

Different algorithms have been developed under this noise model, most of which share the following form:

$$\mathcal{F}_Z(\cdot) = Y_f g, \quad (11a)$$

$$\begin{bmatrix} U_p \\ U_f \\ Y_p \end{bmatrix} g = \begin{bmatrix} \mathbf{u}_{\text{ini}} \\ \mathbf{u} \\ \mathbf{y}_{\text{ini}} + \delta \end{bmatrix}. \quad (11b)$$

The slack variable  $\delta$  is introduced to compensate for the error in both  $Y_p$  and  $\mathbf{y}_{\text{ini}}$ . The algorithms then propose different strategies for balancing the magnitude of  $g$  and the slack variable  $\delta$ . The algorithms are summarized as follows.

**Subspace predictor** [10], [19]: the solution of the algorithm is exactly the same as that for the noise-free case (7). However, the interpretation here is different. It corresponds to the least-squares estimate of a linear mapping:

$$\mathcal{F}_Z(\cdot) = F_Z \text{col}(\mathbf{u}_{\text{ini}}, \mathbf{u}, \mathbf{y}_{\text{ini}}), \quad (12)$$

where

$$F_Z = \underset{F}{\text{argmin}} \left\| Y_f - F \text{col}(U_p, U_f, Y_p) \right\|_F^2. \quad (13)$$

This coincides with finding the vector  $g$  that minimizes  $\|g\|_2^2$  subject to (11b) and  $\delta = \mathbf{0}$ .

**Signal matrix model** [11], [12]: this algorithm uses maximum likelihood estimation to find the vector  $g$  that

maximizes the conditional probability of  $\text{col}(\delta, Y_f g)$  given  $g$ :

$$\min_{g, \delta} \log \det \left( \Sigma_g + \begin{bmatrix} \Sigma_{y_{\text{ini}}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \right) + \delta^\top (\Sigma_p + \Sigma_{y_{\text{ini}}})^{-1} \delta, \quad (14)$$

subject to (11b). When  $\Sigma_{y_{\text{ini}}} = \sigma^2 \mathbb{I}_{n_y L_0}$  and  $\Sigma_g = \sigma^2 \|g\|_2^2 \mathbb{I}_{n_y L}$ , an approximate quadratic program of (14) has been derived as

$$\begin{aligned} \min_{g, \delta} \quad & \|\delta\|_2^2 + n_y \left( L' \sigma^2 / \|g_{\text{pinv}}\|_2^2 + L \sigma^2 \right) \|g\|_2^2 \\ \text{s.t.} \quad & \end{aligned} \quad (15)$$

**Wasserstein distance minimization** [7]: this algorithm finds the vector  $g$  that minimizes the Wasserstein distance between the stochastic distribution of  $\mathbf{y}_{\text{ini}}$  and that of  $Y_p g$ :

$$\min_{g, \delta} \|\delta\|_2^2 + \text{tr} \left( \Sigma_{y_{\text{ini}}} + \Sigma_p - 2 (\Sigma_{y_{\text{ini}}} \Sigma_p)^{1/2} \right), \quad (16)$$

subject to (11b). When  $\Sigma_{y_{\text{ini}}} = \sigma^2 \mathbb{I}_{n_y L_0}$  and  $\Sigma_g = \sigma^2 \|g\|_2^2 \mathbb{I}_{n_y L}$ , an approximate quadratic program of (16) has been derived as

$$\begin{aligned} \min_{g, \delta} \quad & \|\delta\|_2^2 + n_y L_0 \sigma^2 \|g\|_2^2 \\ \text{s.t.} \quad & \end{aligned} \quad (17)$$

It is noted that the algorithms (12), (15), and (17) can be expressed in the following unified form:

$$\begin{aligned} \mathcal{F}_Z(\cdot) = Y_f \underset{g}{\text{argmin}} \quad & \|\delta\|_2^2 + \lambda \|g\|_2^2 \\ \text{s.t.} \quad & \end{aligned} \quad (18)$$

where  $\lambda \rightarrow 0$  for (12),  $\lambda = n_y \left( L' \sigma^2 / \|g_{\text{pinv}}\|_2^2 + L \sigma^2 \right)$  for (15), and  $\lambda = n_y L_0 \sigma^2$  for (17). With an abuse of notation,  $\underset{g}{\text{argmin}}_g$  denotes the optimal solution of  $g$  for the program depending on both  $g$  and  $\delta$ . The optimization problem in (18) is a strongly convex quadratic program with only equality constraints. It admits a closed-form solution that is linear with respect to  $\text{col}(\mathbf{u}_{\text{ini}}, \mathbf{u}, \mathbf{y}_{\text{ini}})$ .

### III. CONFIDENCE REGION ANALYSIS

In this section, confidence regions are established for the stochastic data-driven prediction algorithms discussed in Section II-C. The result first exploits information from the underlying state-space model. Then, a data-driven approximation of the model information is proposed.

#### A. Derivation of the Confidence Region

For any stochastic data-driven predictor in the form of (11), the output estimate (4) differs from the true output  $\mathbf{y}_0$  due to the following two sources of error: 1) the output part of the signal matrix  $Y_f$  is noisy, 2) the predictor estimates a trajectory whose output initial condition is  $Y_p^0 g$ , which differs from the trajectory to be predicted whose output initial condition is  $\mathbf{y}_{\text{ini}}^0$ . By characterizing the distributions of these two sources of error for a particular estimate of  $g$  and  $\delta$ , we obtain the following confidence region for stochastic data-driven prediction.

*Theorem 1:* Consider a stochastic data-driven predictor  $\mathbf{y} = \mathcal{F}_Z(\mathbf{u}; \mathbf{u}_{\text{ini}}, \mathbf{y}_{\text{ini}}) = Y_f g$  satisfying (11). The true output  $\mathbf{y}_0$  is in the following ellipsoidal set w.p.  $p$ :

$$\mathcal{Y} = \left\{ \tilde{\mathbf{y}} \mid (\mathbf{y} - \tilde{\mathbf{y}} - \Gamma \delta)^\top \Sigma^{-1} (\mathbf{y} - \tilde{\mathbf{y}} - \Gamma \delta) \leq \mu_p \right\}, \quad (19)$$

where

$$\Gamma = \text{col}(CA^{L_0}, \dots, CA^{L-1}) \text{col}(C, \dots, CA^{L_0-1})^\dagger, \quad (20)$$

$$\Sigma = \begin{bmatrix} -\Gamma & \mathbb{I}_{n_y L'} \\ \mathbb{I}_{n_y L'} & \end{bmatrix} \Sigma_g \begin{bmatrix} -\Gamma^\top \\ \mathbb{I}_{n_y L'} \end{bmatrix} + \Gamma \Sigma_{y_{\text{ini}}} \Gamma^\top, \quad (21)$$

and  $\mu_p$  satisfies  $F_{\chi^2(L')}(\mu_p) = p$ , where  $F_{\chi^2(d)}(\cdot)$  is the cumulative distribution function of the  $\chi^2$ -distribution with  $d$  degrees of freedom.

*Proof:* Let the stochastic noise in  $Y_p$ ,  $Y_f$ , and  $\mathbf{y}_{\text{ini}}$  be  $E_p$ ,  $E_f$ , and  $\boldsymbol{\varepsilon}_{\text{ini}}$  respectively, i.e.,

$$E_p = Y_p - Y_p^0, \quad E_f = Y_f - Y_f^0, \quad \boldsymbol{\varepsilon}_{\text{ini}} = \mathbf{y}_{\text{ini}} - \mathbf{y}_{\text{ini}}^0. \quad (22)$$

The estimation error can be decomposed as follows, according to the two aforementioned sources of error

$$\mathbf{y} - \mathbf{y}_0 = E_f g + \mathbf{y}^-, \quad (23)$$

where  $\mathbf{y}^-$  is the error due to the discrepancy  $(Y_p^0 g - \mathbf{y}_{\text{ini}}^0)$  in the output initial condition. The initial condition error  $\mathbf{y}^-$  can be seen as the free response from initial condition  $\mathbf{u}_{\text{ini}}^- = \mathbf{0}$ ,  $\mathbf{y}_{\text{ini}}^- = Y_p^0 g - \mathbf{y}_{\text{ini}}^0$ . From (11b) and (22), we have

$$Y_p^0 g = \mathbf{y}_{\text{ini}} + \delta - E_p g, \quad \mathbf{y}_{\text{ini}}^0 = \mathbf{y}_{\text{ini}} - \boldsymbol{\varepsilon}_{\text{ini}}, \quad (24)$$

$$\mathbf{y}_{\text{ini}}^- = (\mathbf{y}_{\text{ini}} + \delta - E_p g) - (\mathbf{y}_{\text{ini}} - \boldsymbol{\varepsilon}_{\text{ini}}) = \delta + \boldsymbol{\varepsilon}_{\text{ini}} - E_p g. \quad (25)$$

Let the state of the trajectory at time  $-L_0$  be  $x^-$ . Then we have

$$\mathbf{y}_{\text{ini}}^- = \begin{bmatrix} C \\ \vdots \\ CA^{L_0-1} \end{bmatrix} x^-, \quad \mathbf{y}^- = \begin{bmatrix} CA^{L_0} \\ \vdots \\ CA^{L-1} \end{bmatrix} x^-. \quad (26)$$

Since  $L_0 \geq l$ ,  $\text{col}(C, \dots, CA^{L_0-1})$  has full column rank. Thus, we have  $x^- = \text{col}(C, \dots, CA^{L_0-1})^\dagger \mathbf{y}_{\text{ini}}^-$ . This directly leads to  $\mathbf{y}^- = \Gamma \mathbf{y}_{\text{ini}}^-$ . From (23)-(25), the estimation error is then

$$\mathbf{y} - \mathbf{y}_0 = E_f g + \Gamma (\delta + \boldsymbol{\varepsilon}_{\text{ini}} - E_p g). \quad (27)$$

Recall that  $\boldsymbol{\varepsilon}_{\text{ini}} \sim \mathcal{N}(\mathbf{0}, \Sigma_{y_{\text{ini}}})$ ,  $\text{col}(E_p, E_f) g \mid g \sim \mathcal{N}(\mathbf{0}, \Sigma_g)$ , and they are uncorrelated. The distribution of  $(\mathbf{y} - \mathbf{y}_0)$  given  $g$  and  $\delta$  is Gaussian with

$$\begin{aligned} \mathbb{E}(\mathbf{y} - \mathbf{y}_0) &= \Gamma \delta, \\ \text{cov}(\mathbf{y} - \mathbf{y}_0) &= \mathbb{E} \left( \begin{bmatrix} -\Gamma & \mathbb{I}_{n_y L'} \end{bmatrix} \begin{bmatrix} E_p \\ E_f \end{bmatrix} g + \Gamma \boldsymbol{\varepsilon}_{\text{ini}} \right) \\ &\quad \left( \begin{bmatrix} -\Gamma & \mathbb{I}_{n_y L'} \end{bmatrix} \begin{bmatrix} E_p \\ E_f \end{bmatrix} g + \Gamma \boldsymbol{\varepsilon}_{\text{ini}} \right)^\top \\ &= \begin{bmatrix} -\Gamma & \mathbb{I}_{n_y L'} \end{bmatrix} \Sigma_g \begin{bmatrix} -\Gamma^\top \\ \mathbb{I}_{n_y L'} \end{bmatrix} + \Gamma \Sigma_{y_{\text{ini}}} \Gamma^\top = \Sigma. \end{aligned} \quad (28)$$

Therefore,  $(\mathbf{y} - \mathbf{y}_0 - \Gamma \delta)^\top \Sigma^{-1} (\mathbf{y} - \mathbf{y}_0 - \Gamma \delta)$  is subject to the  $\chi^2$ -distribution with  $L'$  degrees of freedom. This directly leads to (19).  $\blacksquare$

*Remark 1:* Theorem 1 stills holds when the system is not observable by replacing  $A$ ,  $C$ , and  $l$  with those for the observable part of the system.

*Remark 2:* The derivation is inspired by the prediction error bound presented in Section IV.C of [14]. However, the results of [14] consider a bounded non-stochastic noise model and provide a deterministic but admittedly non-tight bound on  $\|\mathbf{y} - \mathbf{y}_0\|$ .

### B. Data-Driven Formulation of System Parameter $\Gamma$

The confidence region given in Theorem 1 is not available in practice since  $\Gamma$  is dependent on the unknown model parameters  $A$  and  $C$ . However, this system parameter matrix can be alternatively formulated by another data-driven prediction scheme offline.

As can be seen from the proof of Theorem 1, the matrix  $\Gamma$  can be considered as a linear data-driven predictor with  $\mathbf{u} = \mathbf{0}$  and  $\mathbf{u}_{\text{ini}} = \mathbf{0}$ . Supposing we have a noise-free signal matrix, the following lemma gives a data-driven version of Theorem 1 without knowledge of  $A$  and  $C$ .

*Lemma 2:* Let  $\bar{Z} = \text{col}(\bar{U}_p, \bar{U}_f, \bar{Y}_p, \bar{Y}_f)$  be a noise-free signal matrix with  $\text{rank}(\bar{Z}) = n_u L + n_x$ . If  $\Gamma$  is replaced by  $\Gamma_Z = \bar{Y}_f P$ , where  $P$  is the last  $n_y L_0$  columns of  $\text{col}(\bar{U}_p, \bar{U}_f, \bar{Y}_p)^\dagger$ , then Theorem 1 holds.

*Proof:* According to Lemma 1, for any output initial condition  $\mathbf{y}_{\text{ini}}$ ,  $\mathbf{y} = \Gamma_Z \mathbf{y}_{\text{ini}}$  is the unique free response with  $\mathbf{u}_{\text{ini}} = \mathbf{0}$ . So we have  $\mathbf{y}^- = \Gamma_Z \mathbf{y}_{\text{ini}}^-$ . The rest of the proof of Theorem 1 remains the same. ■

*Remark 3:* In general,  $\Gamma_Z \neq \Gamma$ . This is because when  $n_y L_0 > n_x$ , the valid  $\Gamma$  in the proof of Theorem 1 is not unique. The pseudo-inverse solution (20) gives only one possibility.

In practice, the noisy signal matrix  $Z$  can be used to find an approximation of the data-driven system parameter  $\Gamma_Z$ . Recall that the estimated mappings in the form of (18) admit linear solutions. So they can be employed to find an estimate of the linear mapping  $\Gamma_Z$  by setting  $\mathbf{u} = \mathbf{0}$ ,  $\mathbf{u}_{\text{ini}} = \mathbf{0}$ . The closed-form solution is given by

$$\hat{\Gamma}_Z = Y_f (F^{-1} - F^{-1} U^\top (U F^{-1} U^\top)^{-1} U F^{-1}) Y_p^\top, \quad (29)$$

where  $F = \lambda \mathbb{I}_M + Y_p^\top Y_p$  and  $U = \text{col}(U_p, U_f)$  as derived in [11]. The hyperparameter  $\lambda$  can be selected as approaching 0 (subspace predictor),  $n_y L \sigma^2$  (signal matrix model), or  $n_y L_0 \sigma^2$  (Wasserstein distance minimization), and the corresponding  $\hat{\Gamma}_Z$  estimates are denoted by  $\hat{\Gamma}_{\text{Sub}}$ ,  $\hat{\Gamma}_{\text{SMM}}$ , and  $\hat{\Gamma}_{\text{WD}}$  respectively. Note that in this estimation, the output initial condition  $\mathbf{y}_{\text{ini}}^-$  is known exactly without noise. This leads to a slight change in the hyperparameter of the signal matrix model solution. When  $\Gamma$  is replaced by  $\hat{\Gamma}_Z$ , Theorem 1 only holds approximately. The validity of the approximation will be investigated in Section V.

## IV. MINIMUM MEAN-SQUARED ERROR ALGORITHM

In the proof of Theorem 1, the distribution of the estimation error has been derived in order to quantify the confidence region for a given estimate of  $g$  and  $\delta$  with the algorithms discussed in Section II-C. In this section, this distribution

is used to propose a novel optimal predictor in the form of (11), which directly targets maximum prediction accuracy, instead of the statistical properties as in Section II-C. This algorithm finds  $g$  and  $\delta$  in the mapping by minimizing the expected estimation error subject to (28), which leads to the following proposition.

*Proposition 1:* The minimum MSE estimate of the mapping in the form of (11) is given by

$$\begin{aligned} \mathcal{F}_Z(\cdot) = Y_f \underset{g}{\text{argmin}} \delta^\top \Gamma^\top \Gamma \delta + \text{tr} \left( \begin{bmatrix} -\Gamma & \mathbb{I}_{n_y L'} \end{bmatrix} \Sigma_g \begin{bmatrix} -\Gamma^\top \\ \mathbb{I}_{n_y L'} \end{bmatrix} \right) \\ \text{s.t.} \quad (11\text{b}). \end{aligned} \quad (30)$$

*Proof:* From (28), we have

$$\begin{aligned} \text{MSE}(\mathbf{y} - \mathbf{y}_0) &= \mathbb{E}(\mathbf{y} - \mathbf{y}_0)^\top (\mathbf{y} - \mathbf{y}_0) \\ &= \text{tr} \left( \text{cov}(\mathbf{y} - \mathbf{y}_0) + \mathbb{E}(\mathbf{y} - \mathbf{y}_0) \mathbb{E}(\mathbf{y} - \mathbf{y}_0)^\top \right) \\ &= \text{tr}(\Sigma + \Gamma \delta \delta^\top \Gamma^\top) = \text{tr}(\Sigma) + \delta^\top \Gamma^\top \Gamma \delta. \end{aligned} \quad (31)$$

where the third equality comes from (28). From the definition of  $\Sigma$  in (21), it is observed that since  $\Gamma \Sigma_{\mathbf{y}_{\text{ini}}} \Gamma^\top$  does not depend on the optimization variables  $g$  and  $\delta$ , minimizing the MSE is equivalent to the optimization problem in (30). ■

If we assume that  $\Sigma_Y = \sigma^2 \mathbb{I}_{n_y L M}$ , (30) becomes

$$\begin{aligned} \mathcal{F}_Z(\cdot) = Y_f \underset{g}{\text{argmin}} \|\delta\|_Q^2 + \lambda_{\text{MSE}} \|g\|_2^2 \\ \text{s.t.} \quad (11\text{b}), \end{aligned} \quad (32)$$

where  $Q = \Gamma^\top \Gamma$  and  $\lambda_{\text{MSE}} = \sigma^2 n_y L' + \sigma^2 \text{tr}(Q)$ . This optimization problem is very similar to the unified form (18) for existing algorithms, except that the Euclidean norm of  $\delta$  is now weighted by  $Q$ . The solution (32) is also linear with respect to  $\text{col}(\mathbf{u}_{\text{ini}}, \mathbf{u}, \mathbf{y}_{\text{ini}})$ .

The implications of Proposition 1 are twofold. On the one hand, it provides the optimal solution to the data-driven prediction problem with output noise in terms of minimizing the MSE. Although the optimal solution relies on the unknown extended observability matrix to formulate  $\Gamma$ , it can be used with a preliminary model or a model set via minimax approaches.

On the other hand, similar to establishing the confidence region, the parameter  $\Gamma$  used in the minimum MSE solution (30) can be replaced by the data-driven estimate  $\hat{\Gamma}_Z$  (29) derived from the same signal matrix for an approximate solution. This leads to the minimum-MSE data-driven predictor, denoted as Algorithm 1.

## V. NUMERICAL EXAMPLES

In this section, numerical tests are conducted to illustrate the validity of the derived confidence region and the effectiveness of the proposed minimum-MSE algorithm. In the examples, stochastic data with i.i.d. noise are collected from one single experiment and used in  $Z$  with a Page matrix construction. Unit Gaussian input sequences are used to generate the data.

**Algorithm 1** The minimum-MSE data-driven predictor with stochastic data

- 1: **Given:** signal matrix  $Z$ , noise model  $\Sigma_Y, \Sigma_{y_{ini}}$ , confidence level  $p$ .
- 2: **Input:**  $\mathbf{u}_{ini}, \mathbf{y}_{ini}, \mathbf{u}$ .
- 3: Calculate  $\hat{\Gamma}_Z$  by (29).
- 4: Find  $\mathbf{y} = \mathcal{F}_Z(\mathbf{u}; \mathbf{u}_{ini}, \mathbf{y}_{ini})$  by solving (30) with  $\Gamma = \hat{\Gamma}_Z$ .
- 5: Find  $p$ -confidence region  $\mathcal{Y}$  by (19) with  $\Gamma = \hat{\Gamma}_Z$ .
- 6: **Output:**  $\mathbf{y}, \mathcal{Y}$ .

First, we consider a simple two-dimensional example for illustration purposes. The prediction problem is to find the first two points ( $L' = 2$ ) in the step response of the following fourth-order system

$$G_1(z) = \frac{0.1059(0.1z^4 + z^3 + 0.5z^2)}{z^4 - 2.2z^3 + 2.42z^2 - 1.87z + 0.7225}. \quad (33)$$

The prediction conditions are  $\mathbf{u}_{ini} = \mathbf{0}$ ,  $\mathbf{y}_{ini} = \mathbf{0}$ , and  $\mathbf{u} = [1 \ 1]^T$ . The following parameters are used:  $L = 10$ ,  $L_0 = 8$ ,  $M = 80$ , and noise level  $\sigma^2 = 0.1$ . A confidence level of  $p = 0.90$  is used in the following figures.

Figure 1 compares the confidence regions obtained using model-based  $\Gamma$  (20) (*CR-MB*), data-driven  $\hat{\Gamma}_{Sub}$  (*CR-Sub*),  $\hat{\Gamma}_{SMM}$  (*CR-SMM*), and  $\hat{\Gamma}_{WD}$  (*CR-WD*). The confidence regions are tested on the minimum-MSE predictor with data-driven  $\hat{\Gamma}_{SMM}$  (*MSE-SMM*). 10 different realizations of the stochastic data are plotted. The results show that the data-driven formulations (*CR-Sub*, *CR-SMM*, and *CR-WD*) obtain similar confidence regions, but are different from the model-based formulation. This is because the data-driven formulations with  $\hat{\Gamma}_Z$  estimate the noise-free  $\Gamma_Z$  that is different from the model-based  $\Gamma$ . Nevertheless, all the confidence regions are valid for this problem, since the true trajectory lies in the regions with high probability.

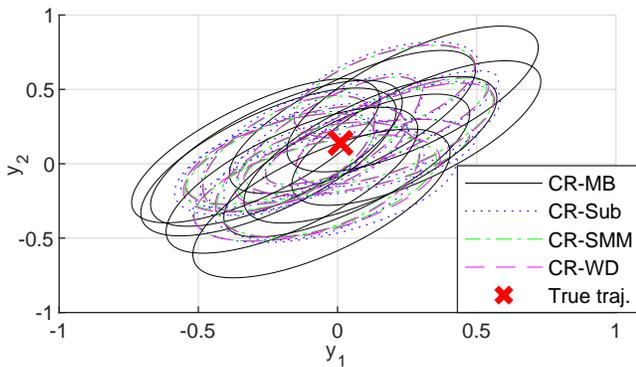


Fig. 1. Comparison of different confidence region formulations ( $p = 0.90$ ) tested on the *MSE-SMM* predictor with 10 different realizations of the stochastic data.

Then, the sizes of the confidence regions are analyzed for different stochastic data-driven predictors. The following predictors are compared: 1) subspace predictor (12) (*Sub*), 2) signal matrix model (15) (*SMM*), 3) Wasserstein distance minimization (17) (*WD*), and 4) minimum-MSE predictor

(Algorithm 1) using model-based  $\Gamma$  (20) (*MSE-MB*), data-driven  $\hat{\Gamma}_{Sub}$  (*MSE-Sub*),  $\hat{\Gamma}_{SMM}$  (*MSE-SMM*), and  $\hat{\Gamma}_{WD}$  (*MSE-WD*). Figure 2 shows the confidence regions of these stochastic predictors with model-based  $\Gamma$  (*CR-MB*). As can be seen from the figure, the existing algorithms (*Sub*, *SMM*, and *WD*) have larger confidence regions compared to the minimum-MSE algorithms (*MSE-MB* and *MSE-SMM*). This illustrates the effectiveness of the proposed algorithm in improving prediction accuracy. In this example, the confidence regions of *MSE-Sub* and *MSE-WD* are very close to that of *MSE-SMM*, so they are omitted in Figure 2.

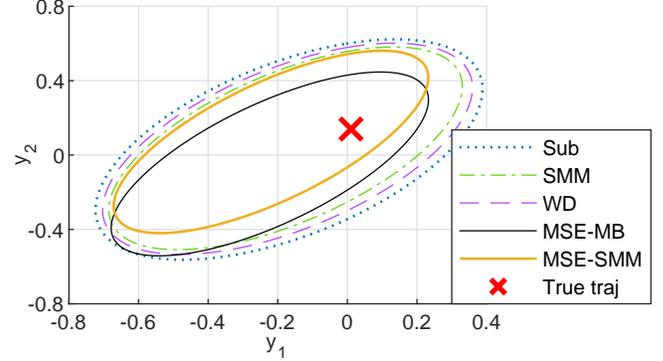


Fig. 2. Comparison of different stochastic data-driven predictors in terms of the confidence regions ( $p = 0.90$ ) with model-based  $\Gamma$  (*CR-MB*).

To quantitatively assess the derived confidence region and the minimum-MSE prediction algorithm, the following campaign of 1000 Monte Carlo simulations is set up. A bank of 1000 single-input, single-output systems are randomly generated by the `drss` command in MATLAB with random numbers of states between 3 and 8. These random systems are normalized to have an  $\mathcal{H}_2$ -gain of 1. The prediction problem uses the following parameters:  $L = 20$ ,  $L_0 = 8$ ,  $L' = 12$ , and  $M = 320$ . The input  $\mathbf{u}$  and the initial condition ( $\mathbf{u}_{ini}, \mathbf{y}_{ini}$ ) are selected randomly with a unit Gaussian distribution.

Table I compares the percentage of the simulations where the true response is in the confidence region, i.e.,  $\mathbf{y}_0^i \in \mathcal{Y}^i$  for the  $i$ -th simulation, for the model-based and different data-driven formulations. Two confidence levels  $p = 0.95$  and  $p = 0.99$  are selected. The noise level is selected as  $\sigma^2 = 0.1$ . The rows in Table I correspond to different predictors, whereas the columns correspond to different formulations of the confidence region. It can be seen from the table that the empirical confidence levels match the targeted  $p$ -value well with the model-based  $\Gamma$  (*CR-MB*) for all three predictors, where Theorem 1 is satisfied exactly. With the data-driven estimates  $\hat{\Gamma}_Z$ , the confidence regions become marginally more conservative as the empirical confidence levels are slightly larger in Table I. The results of the three data-driven estimates (*CR-Sub*, *CR-SMM*, *CR-WD*) are similar, which indicates that the confidence region is not very sensitive to the choice of  $\hat{\Gamma}_Z$  estimation method.

Table II compares the empirical MSE of the predictors in the Monte Carlo simulations to the MSE estimated by (31) with the approximate data-driven confidence regions. The

TABLE I  
EMPIRICAL CONFIDENCE LEVELS OF THE CONFIDENCE REGIONS

$p = 0.95$	CR-MB	CR-Sub	CR-SMM	CR-WD
Sub	97.1%	98.7%	98.4%	98.7%
SMM	96.8%	97.4%	97.3%	97.3%
MSE-SMM	95.2%	96.4%	96.2%	96.4%
$p = 0.99$	CR-MB	CR-Sub	CR-SMM	CR-WD
Sub	99.3%	100%	99.8%	99.9%
SMM	99.2%	99.7%	99.7%	99.7%
MSE-SMM	99.0%	99.3%	99.2%	99.3%

empirical MSE is computed as

$$\text{MSE}_{\text{emp}}(\mathbf{y} - \mathbf{y}_0) = \frac{1}{N_s} \sum_{i=1}^{N_s} \|\mathbf{y}^i - \mathbf{y}_0^i\|_2^2, \quad (34)$$

where  $\mathbf{y}^i$  and  $\mathbf{y}_0^i$  are the predicted and the true responses of the  $i$ -th simulation respectively, and  $N_s = 1000$ . Two different noise levels of  $\sigma^2 = 0.1$  and  $\sigma^2 = 1$  are considered. Similar to the observation from Table I, the estimated MSE is shown to be more conservative compared to the empirical ones for all three predictors. It is also observed that the region *CR-SMM* is the less conservative among those tested here. However, the estimated MSE can correctly predict the relative error magnitudes of different predictors. This illustrates that the estimated MSE can be a good indicator of prediction accuracy, which motivates its use as the objective function in Algorithm 1. Only three representative predictors are shown in Table I and Table II for clarity. The results of the other algorithms are similar.

TABLE II  
COMPARISON OF THE ESTIMATED AND THE EMPIRICAL MSE

$\sigma^2 = 0.1$	Empirical	CR-Sub	CR-SMM	CR-WD
Sub	0.115	0.153	0.149	0.152
SMM	0.099	0.142	0.137	0.140
MSE-SMM	0.096	0.136	0.131	0.134
$\sigma^2 = 1$	Empirical	CR-Sub	CR-SMM	CR-WD
Sub	1.106	1.529	1.485	1.511
SMM	0.915	1.391	1.344	1.372
MSE-SMM	0.897	1.335	1.286	1.317

Finally, we compare the prediction accuracy of the predictors by the empirical MSE, under three different noise levels  $\sigma^2 = 0.1$ ,  $\sigma^2 = 0.5$ , and  $\sigma^2 = 1$ . The results are shown in Table III. For all three noise levels, the minimum-MSE predictor with model-based  $\Gamma$  (*MSE-MB*) achieves the minimum empirical MSE. This is expected as *MSE-MB* exactly optimizes for this objective as demonstrated in Proposition 1. However, the model-based  $\Gamma$  is not available in practice. Among the other practical algorithms, Algorithm 1 with  $\hat{\Gamma}_Z$  based on the signal matrix model (*MSE-SMM*) has the smallest empirical MSE, with slightly better performance than the direct signal matrix model approach (*SMM*). This result shows numerically that, with approximate data-driven formulations of  $\hat{\Gamma}_Z$ , the proposed minimum-MSE predictor still obtains a more accurate prediction than the existing algorithms.

TABLE III  
COMPARISON OF THE EMPIRICAL MSE FOR DIFFERENT PREDICTORS

	$\sigma^2 = 0.1$	$\sigma^2 = 0.5$	$\sigma^2 = 1$
Sub	0.115	0.558	1.106
SMM	0.099	0.476	0.915
WD	0.113	0.548	1.091
MSE-MB	0.094	0.435	0.833
MSE-Sub	0.097	0.464	0.908
MSE-SMM	0.096	0.460	0.897
MSE-WD	0.097	0.462	0.902

## VI. CONCLUSIONS

In this paper, the prediction error of data-driven predictors with stochastic data is characterized statistically. The framework provides ellipsoidal confidence regions for various predictors. It also offers a novel optimal predictor that minimizes the mean-squared prediction error directly. In practice, both the confidence region and the minimum-MSE predictor can be implemented with data-driven approximations that show good accuracy numerically.

Both the derived confidence region and the minimum-MSE predictor can contribute to more reliable and effective applications of stochastic data-driven predictors to predictive control design with robustness guarantees on the satisfaction of safety-critical constraints.

## REFERENCES

- [1] I. Markovsky and F. Dörfler, "Behavioral systems theory in data-driven analysis, signal processing, and control," *preprint*, 2021, Available: <http://homepages.vub.ac.be/~imarkovs/publications/overview-ddctr.pdf>.
- [2] J. C. Willems, P. Rapisarda, I. Markovsky, and B. L. M. De Moor, "A note on persistency of excitation," *Systems & Control Letters*, vol. 54, no. 4, pp. 325–329, 2005.
- [3] I. Markovsky and F. Dörfler, "Identifiability in the behavioral setting," *preprint*, 2020, Available: <http://homepages.vub.ac.be/~imarkovs/publications/identifiability.pdf>.
- [4] S. Sedghizadeh and S. Beheshti, "Data-driven subspace predictive control: Stability and horizon tuning," *Journal of the Franklin Institute*, vol. 355, no. 15, pp. 7509–7547, 2018.
- [5] J. Coulson, J. Lygeros, and F. Dörfler, "Data-enabled predictive control: In the shallows of the DeePC," in *2019 18th European Control Conference (ECC)*, 2019, pp. 307–312.
- [6] L. Frieri, B. Guo, A. Martin, and G. Ferrari-Trecate, "A behavioral input-output parametrization of control policies with suboptimality guarantees," *arXiv preprint arXiv:2102.13338*, 2021.
- [7] Y. Lian, J. Shi, M. P. Koch, and C. N. Jones, "Adaptive robust data-driven building control via bi-level reformulation: an experimental result," *arXiv preprint arXiv:2106.05740*, 2021.
- [8] L. Huang, J. Coulson, J. Lygeros, and F. Dörfler, "Decentralized data-enabled predictive control for power system oscillation damping," *IEEE Transactions on Control Systems Technology*, 2021.
- [9] F. Dörfler, J. Coulson, and I. Markovsky, "Bridging direct & indirect data-driven control formulations via regularizations and relaxations," *arXiv preprint arXiv:2101.01273*, 2021.
- [10] L. Huang, J. Coulson, J. Lygeros, and F. Dörfler, "Data-enabled predictive control for grid-connected power converters," in *IEEE Conference on Decision and Control (CDC)*, 2019, pp. 8130–8135.
- [11] M. Yin, A. Iannelli, and R. S. Smith, "Maximum likelihood estimation in data-driven modeling and control," *arXiv preprint arXiv:2011.00925*, 2020.
- [12] —, "Maximum likelihood signal matrix model for data-driven predictive control," in *Proceedings of the 3rd Conference on Learning for Dynamics and Control*, ser. Proceedings of Machine Learning Research, vol. 144, 2021, pp. 1004–1014.

- [13] A. Alanwar, Y. Stürz, and K. H. Johansson, “Robust data-driven predictive control using reachability analysis,” *arXiv preprint arXiv:2103.14110*, 2021.
- [14] J. Berberich, J. Kohler, M. A. Muller, and F. Allgower, “Data-driven model predictive control with stability and robustness guarantees,” *IEEE Transactions on Automatic Control*, vol. 66, no. 4, pp. 1702–1717, 2021.
- [15] —, “Robust constraint satisfaction in data-driven MPC,” in *2020 59th IEEE Conference on Decision and Control (CDC)*, 2020.
- [16] A. Damen, P. Van den Hof, and A. Hajdasinski, “Approximate realization based upon an alternative to the Hankel matrix: the Page matrix,” *Systems & Control Letters*, vol. 2, no. 4, pp. 202–208, 1982.
- [17] H. J. van Waarde, C. De Persis, M. K. Camlibel, and P. Tesi, “Willems’ fundamental lemma for state-space systems and its extension to multiple datasets,” *IEEE Control Systems Letters*, vol. 4, no. 3, pp. 602–607, 2020.
- [18] I. Markovsky and P. Rapisarda, “Data-driven simulation and control,” *International Journal of Control*, vol. 81, no. 12, pp. 1946–1959, 2008.
- [19] F. Fiedler and S. Lucia, “On the relationship between data-enabled predictive control and subspace predictive control,” *arXiv preprint arXiv:2011.13868*, 2021.