Distributionally Robust Receive Beamforming

Shixiong Wang, Wei Dai, and Geoffrey Ye Li, Fellow, IEEE

Abstract—This article investigates signal estimation in wireless transmission (i.e., receive beamforming) from the perspective of statistical machine learning, where the transmit signals may be from an integrated sensing and communication system; that is, 1) signals may be not only discrete constellation points but also arbitrary complex values; 2) signals may be spatially correlated. Particular attention is paid to handling various uncertainties such as the uncertainty of the transmit signal covariance, the uncertainty of the channel matrix, the uncertainty of the channel noise covariance, the existence of channel impulse noises, and the limited sample size of pilots. To proceed, a distributionally robust machine learning framework that is insensitive to the above uncertainties is proposed, which reveals that channel estimation is not a necessary operation. For optimal linear estimation, the proposed framework includes several existing beamformers as special cases such as diagonal loading and eigenvalue thresholding. For optimal nonlinear estimation, estimators are limited in reproducing kernel Hilbert spaces and neural network function spaces, and corresponding uncertainty-aware solutions (e.g., kernelized diagonal loading) are derived. In addition, we prove that the ridge and kernel ridge regression methods in machine learning are distributionally robust against diagonal perturbation in feature covariance.

Index Terms—Wireless Transmission, Smart Antenna, Machine Learning, Robust Estimation, Robust Beamforming, Distributional Uncertainty, Channel Uncertainty, Limited Pilot.

I. INTRODUCTION

N wireless transmission, detection and estimation of transmitted signals is of high importance, and beamforming at array receivers serves as a key signal-processing technique to suppress interference and environmental noises. The earliest beamforming solutions rely on the use of phase shifters (e.g., phased arrays) to steer and shape wave lobes, while advanced beamforming methods allow the employment of digital signal processing units, which introduce additional structural freedom (e.g., fully digital, hybrid, nonlinear, wideband) in beamformer design and significant performance improvement in signal recovery [1]–[3].

In traditional communication systems, transmitted signals are discrete points from constellations. Therefore, signal recovery, commonly referred to as *signal detection*, can be cast into a classification problem from the perspective of statistical machine learning, and the number of candidate classes is determined by the number of points in the employed constellation. Research in this stream includes, e.g., [4]–[9] as well as references therein, and the performance measure

S. Wang, W. Dai, and G. Li are with the Department of Electrical and Electronic Engineering, Imperial College London, London SW7 2AZ, United Kingdom (E-mail: s.wang@u.nus.edu; wei.dai1@imperial.ac.uk; geoffrey.li@imperial.ac.uk).

This work is supported by the UK Department for Science, Innovation and Technology under the Future Open Networks Research Challenge project TUDOR (Towards Ubiquitous 3D Open Resilient Network). The views expressed are those of the authors and do not necessarily represent the project.

for signal detection is usually the misclassification rate (i.e., symbol error rate); representative algorithms encompass the maximum likelihood detector, the sphere decoding, etc. In another research stream, the signal recovery performance is evaluated using mean-squared errors (cf., signal-to-interferenceplus-noise ratio), and the resultant signal recovery problem is commonly known as signal estimation, which can be considered as a regression problem from the perspective of statistical machine learning. By comparing the estimated symbols with the constellation points afterward, the detection of discrete symbols can be realized. For this case, till now, typical beamforming solutions include zero-forcing receivers, Wiener receivers (i.e., linear minimum mean-squared error receivers), Capon receivers (i.e., minimum variance distortionless response receivers), and nonlinear receivers such as neuralnetwork receivers [10]-[12]. On the basis of these canonical approaches, variants such as robust beamformers working against the limited size of pilot samples and the uncertainty in steering vectors [13]–[18] have also been intensively reported; among these robust solutions, the diagonal loading method [19], [14, Eq. (11)] and the eigenvalue thresholding method [20], [14, Eq. (12)], etc., are popular due to their excellent balance in practical performance and technical simplicity.

Different from traditional paradigms, in emerging communication systems, e.g., integrated sensing and communication (ISAC) systems, transmitted signals may be arbitrary complex-valued symbols and spatially correlated [21]–[23]. As a result, mean-squared error is a preferred performance measure to investigate the *receive beamforming and estimation* problem of wireless signals, which is, therefore, the focus of this article.

Although a large body of problems have been attacked in the area, the following signal-processing problems of beamforming and estimation in wireless transmission remain unsolved.

- 1) What is the relation between the signal-model-based approaches (e.g., Wiener and Capon receivers) and the data-driven approaches (e.g., deep-learning receivers)? In other words, how can we build a mathematically unified modeling framework to interpret all the existing digital receive beamformers?
- 2) In addition to the limited pilot size and the uncertainty in steering vectors, there exist other uncertainties in the signal model: the uncertainty in the transmit signal covariance, the uncertainty of the communication channel matrix, the uncertainty of the channel noise covariance, and the presence of channel impulse noises (i.e., outliers). Therefore, how can we handle all these types of uncertainties in a unified solution framework?
- 3) Existing literature mainly studied the robustness theory of linear beamformers against limited pilot size and the uncertainty in steering vectors [13]–[18]. However, how can we develop the theory of robust nonlinear beamformers

against all the aforementioned uncertainties?

To this end, this article designs a unified modeling and solution framework for receive beamforming of wireless signals, in consideration of the scarcity of the pilot data and the different uncertainties in the signal model.

A. Contributions

The contributions of this article can be summarized from the aspects of machine learning theory and wireless transmission theory.

In terms of machine learning theory, we give a justification of the popular ridge regression and kernel ridge regression (i.e., quadratic loss function plus squared-F-norm regularization) from the perspective of distributional robustness against diagonal perturbation in feature covariance, which enriches the theory of trustworthy machine learning; see Theorems 2 and 3, as well as Corollaries 3 and 5.

In terms of wireless transmission theory, the contributions are outlined below.

- 1) We build a fundamentally theoretical framework for receive beamforming from the perspective of statistical machine learning. In addition to the linear estimation methods, nonlinear approaches (i.e., nonlinear beamforming) are also discussed in reproducing kernel Hilbert spaces and neural network function spaces. In particular, we reveal that channel estimation is not a necessary operation in receive beamforming. For details, see Subsection III-A.
- 2) The presented framework is particularly developed from the perspective of distributional robustness which can therefore combat the limited size of pilot data and several types of uncertainties in the wireless signal model such as the uncertainty in the transmit power matrix, the uncertainty in the communication channel matrix, the existence of channel impulse noises (i.e., outliers), the uncertainty in the covariance matrix of channel noises, etc. For details, see Subsection III-B, and the technical developments in Sections IV and V.
- 3) Existing methods such as diagonal loading and eigenvalue thresholding are proven to be distributionally robust against the limited pilot size and all the aforementioned uncertainties in the wireless signal model. Extensions of diagonal loading and eigenvalue thresholding are proposed as well. Moreover, the kernelized diagonal loading and the kernelized eigenvalue thresholding methods are put forward for nonlinear estimation cases. For details, see Corollary 1, Examples 4 and 5, and Subsections IV-B.
- 4) The distributionally robust receive beamforming and signal estimation problems across multiple frames, where channel conditions may change, are also investigated. For details, see Subsections IV-C and V-A2.

B. Notations

The N-dimensional real (coordinate) space and complex (coordinate) space are denoted as \mathbb{R}^N and \mathbb{C}^N , respectively. Lowercase symbols (e.g., \boldsymbol{x}) denote vectors (column by default) and uppercase ones (e.g., \boldsymbol{X}) denote matrices. We use

the Roman font for random quantities (e.g., \mathbf{x}, \mathbf{X}) and the italic font for deterministic quantities (e.g., \mathbf{x}, \mathbf{X}). Let $\operatorname{Re} \mathbf{X}$ be the real part of a complex quantity \mathbf{X} (a vector or matrix) and $\operatorname{Im} \mathbf{X}$ be the imaginary part of \mathbf{X} . For a vector $\mathbf{x} \in \mathbb{C}^N$, let

$$\underline{\boldsymbol{x}} \coloneqq \left[\begin{array}{c} \operatorname{Re} \boldsymbol{x} \\ \operatorname{Im} \boldsymbol{x} \end{array} \right] \in \mathbb{R}^{2N}$$

be the real-space representation of \boldsymbol{x} ; for a matrix $\boldsymbol{H} \in \mathbb{C}^{N \times M}$, let

$$\underline{\boldsymbol{H}} \coloneqq \left[\begin{array}{c} \operatorname{Re} \boldsymbol{H} \\ \operatorname{Im} \boldsymbol{H} \end{array} \right], \qquad \underline{\underline{\boldsymbol{H}}} \coloneqq \left[\begin{array}{cc} \operatorname{Re} \boldsymbol{H} & -\operatorname{Im} \boldsymbol{H} \\ \operatorname{Im} \boldsymbol{H} & \operatorname{Re} \boldsymbol{H} \end{array} \right]$$

be the real-space representations of \boldsymbol{H} where $\underline{\boldsymbol{H}} \in \mathbb{R}^{2N \times M}$ and $\underline{\boldsymbol{H}} \in \mathbb{R}^{2N \times 2M}$. The running index set induced by an integer N is defined as $[N] \coloneqq \{1,2,\ldots,N\}$. To concatenate matrices and vectors, MATLAB notations are used: i.e., $[\boldsymbol{A}, \boldsymbol{B}]$ for row stacking and $[\boldsymbol{A}; \boldsymbol{B}]$ for column stacking. We let $\Gamma_M \coloneqq [\boldsymbol{I}_M, \boldsymbol{J}_M] \in \mathbb{C}^{M \times 2M}$ where \boldsymbol{I}_M denotes the M-dimensional identity matrix, $\boldsymbol{J}_M \coloneqq j \cdot \boldsymbol{I}_M$, and j denotes the imaginary unit. Let $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ denote a real Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$. We use $\mathcal{C}\mathcal{N}(\boldsymbol{s}, \boldsymbol{P}, \boldsymbol{C})$ to denote a complex Gaussian distribution with mean \boldsymbol{s} , covariance \boldsymbol{P} , and pseudo-covariance \boldsymbol{C} ; if \boldsymbol{C} is not specified, we imply $\boldsymbol{C} = \boldsymbol{0}$.

II. PRELIMINARIES

We review two popular structured representation methods of nonlinear functions $\phi : \mathbb{R}^N \to \mathbb{R}^M$. More details can be seen in Appendix A of the online supplementary materials.

A. Reproducing Kernel Hilbert Spaces

A reproducing kernel Hilbert space (RKHS) \mathcal{H} induced by the kernel function $\ker: \mathbb{R}^N \times \mathbb{R}^N \to \mathbb{R}$ and a collection of points $\{x_1, x_2, \dots, x_L\} \subset \mathbb{R}^N$ is a set of functions from \mathbb{R}^N to \mathbb{R} ; L may be infinite. Every function $\phi: \mathbb{R}^N \to \mathbb{R}$ in the functional space \mathcal{H} can be represented by a linear combination [24, p. 539; Chap. 14]

$$\phi(\boldsymbol{x}) = \sum_{i=1}^{L} \omega_i \cdot \ker(\boldsymbol{x}, \boldsymbol{x}_i), \ \forall \boldsymbol{x} \in \mathbb{R}^N$$
 (1)

where $\{\omega_i\}_{i\in[L]}$ are the combination weights; $\omega_i\in\mathbb{R}$ for every $i\in[L]$. The matrix form of (1) for M-multiple functions are

$$\phi(oldsymbol{x})\coloneqq \left[egin{array}{c} \phi_1(oldsymbol{x}) \ \phi_2(oldsymbol{x}) \ dots \ \phi_M(oldsymbol{x}) \end{array}
ight] = oldsymbol{W}\cdotoldsymbol{arphi}(oldsymbol{x})\coloneqq \left[egin{array}{c} oldsymbol{\omega}_1 \ oldsymbol{\omega}_2 \ dots \ oldsymbol{\omega}_M \end{array}
ight] \cdotoldsymbol{arphi}(oldsymbol{x}), \ \ (2)$$

where $\omega_1, \omega_2, \dots, \omega_M \in \mathbb{R}^L$ are weight row-vectors for functions $\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots, \phi_M(\mathbf{x})$, respectively, and

$$oldsymbol{W}\coloneqq\left[egin{array}{c} oldsymbol{\omega}_1\ oldsymbol{\omega}_2\ dots\ oldsymbol{\omega}_M \end{array}
ight]\in\mathbb{R}^{M imes L}, \quad oldsymbol{arphi}(oldsymbol{x})\coloneqq\left[egin{array}{c} \ker(oldsymbol{x},oldsymbol{x}_1)\ \ker(oldsymbol{x},oldsymbol{x}_2)\ dots\ \ker(oldsymbol{x},oldsymbol{x}_L) \end{array}
ight]. \quad ext{(3)}$$

Since a kernel function is pre-designed (i.e., fixed) for an RKHS \mathcal{H} , (2) suggests a W-linear representation of x-nonlinear functions $\phi(x)$ in \mathcal{H}^M . Note that there exists a

one-to-one correspondence between ϕ and W: for every $\phi : \mathbb{R}^N \to \mathbb{R}^M$, there exists a $W \in \mathbb{R}^{M \times L}$, and vice versa.

B. Neural Networks

Neural networks (NN) are another powerful tool to represent (i.e., approximate) nonlinear functions. A neural network function space (NNFS) K characterizes (or parameterizes) a set of multi-input multi-output functions. Typical choices are multilayer feed-forward neural networks, recurrent neural networks, etc. For beamforming and estimation of wireless signals, the multi-layer feed-forward neural networks are standard [10]-[12]. Suppose that we have R-1 hidden layers (so in total R+1 layers including one input layer and one output layer) and each layer r = 0, 1, ..., R contains T_r neurons. To represent a function $\phi: \mathbb{R}^N \to \mathbb{R}^M$, for the input layer r=0and output layer r = R, we have $T_0 = N$ and $T_R = M$, respectively. Let the output of the r^{th} layer be $y_r \in \mathbb{R}^{T_r}$. For every layer r, we have $y_r = \sigma_r(W_r^{\circ} \cdot y_{r-1} + b_r)$ where $\boldsymbol{W}_r^{\circ} \in \mathbb{R}^{T_r \times T_{r-1}}$ is the weight matrix, $\boldsymbol{b}_r \in \mathbb{R}^{T_r}$ is the bias vector, and the multi-output function σ_r is the activation function which is entry-wise identical. Hence, every function $\phi: \mathbb{R}^N \to \mathbb{R}^M$ in a NNFS can be recursively expressed as [25, Chap. 5], [26]

$$\phi(\boldsymbol{x}) = \sigma_R(\boldsymbol{W}_R \cdot [\boldsymbol{y}_{R-1}(\boldsymbol{x}); 1])
\boldsymbol{y}_r(\boldsymbol{x}) = \sigma_r(\boldsymbol{W}_r \cdot [\boldsymbol{y}_{r-1}(\boldsymbol{x}); 1]), \quad r \in [R-1]$$

$$\boldsymbol{y}_0(\boldsymbol{x}) = \boldsymbol{x}.$$
(4)

where $W_r := [W_r^{\circ}, b_r]$ for $r \in [R]$. Note that the activation functions can vary from one layer to another.

III. PROBLEM FORMULATION

Consider a narrow-band wireless signal model

$$\mathbf{x} = \mathbf{H}\mathbf{s} + \mathbf{v} \tag{5}$$

where $\mathbf{x} \in \mathbb{C}^N$ is the received signal, $\mathbf{s} \in \mathbb{C}^M$ is the transmitted signal, $\mathbf{H} \in \mathbb{C}^{N \times M}$ is the channel matrix, and $\mathbf{v} \in \mathbb{C}^N$ is the zero-mean channel noise. The precoding operation (if exists) is integrated in \mathbf{H} . The transmitted symbols \mathbf{s} have zero means, which may be not only discrete values from constellations such as quadrature amplitude modulation but also arbitrary symbols such as integrated sensing and communication (ISAC) signals. We consider L pilots $\mathbf{S} \coloneqq (\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_L)$ in each frame, and the corresponding received symbols are $\mathbf{X} \coloneqq (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_L)$ under the noise $(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_L)$. We suppose that $\mathbf{R}_s \coloneqq \mathbb{E}\mathbf{s}^H$ and $\mathbf{R}_v \coloneqq \mathbb{E}\mathbf{v}\mathbf{v}^H$ may not be identity or diagonal matrices: i.e., the components of \mathbf{s} can be correlated (e.g., in ISAC), so can be these of \mathbf{v} . Consider the real-space representation of the signal model (5) by stacking the real and imaginary components:

$$\mathbf{x} = \mathbf{H} \cdot \mathbf{s} + \mathbf{v},\tag{6}$$

where $\underline{\mathbf{x}} \in \mathbb{R}^{2N}$, $\underline{\underline{H}} \in \mathbb{R}^{2N \times 2M}$, $\underline{\mathbf{s}} \in \mathbb{R}^{2M}$, and $\underline{\mathbf{v}} \in \mathbb{R}^{2N}$. The expressions of $\underline{R}_{\underline{x}} \coloneqq \mathbb{E}\underline{\mathbf{x}}\underline{\mathbf{x}}^\mathsf{T}$, $\underline{R}_{\underline{s}} \coloneqq \mathbb{E}\underline{\mathbf{s}}\underline{\mathbf{s}}^\mathsf{T}$, $\underline{R}_{\underline{x}\underline{s}} \coloneqq \mathbb{E}\underline{\mathbf{x}}\underline{\mathbf{s}}^\mathsf{T}$, and $\underline{R}_{\underline{v}} \coloneqq \mathbb{E}\underline{\mathbf{v}}\underline{\mathbf{v}}^\mathsf{T}$ can be readily obtained; see Appendix B of the online supplementary materials. Signal estimation in real spaces can be technically simpler than that in complex spaces.

A. Optimal Estimation

1) Optimal Nonlinear Estimation (Nonlinear Beamforming): To recover \mathbf{s} using \mathbf{x} , we consider an estimator $\hat{\mathbf{s}} := \phi(\mathbf{x})$ at the receiver where $\phi: \mathbb{C}^N \to \mathbb{C}^M$ is a Borelmeasurable function; note that $\phi(\mathbf{x})$ may be nonlinear in general. The signal estimation problem at the receiver can be written as a statistical machine-learning problem under the joint data distribution $\mathbb{P}_{\mathbf{x},\mathbf{s}}$ of (\mathbf{x},\mathbf{s}) , that is,

$$\min_{\boldsymbol{\phi} \in \mathcal{B}_{\mathbb{C}^N \to \mathbb{C}^M}} \operatorname{Tr} \mathbb{E}_{\mathbf{x}, \mathbf{s}} [\boldsymbol{\phi}(\mathbf{x}) - \mathbf{s}] [\boldsymbol{\phi}(\mathbf{x}) - \mathbf{s}]^{\mathsf{H}}, \tag{7}$$

where $\mathcal{B}_{\mathbb{C}^N \to \mathbb{C}^M}$ contains all Borel-measurable estimators from \mathbb{C}^N to \mathbb{C}^M . In what follows, we omit the notational dependence on \mathbb{C}^N and \mathbb{C}^M , and use \mathcal{B} as a shorthand. The optimal estimator, in the sense of minimum mean-squared error, is known as the conditional mean of s given x, i.e.,

$$\hat{\mathbf{s}} = \boldsymbol{\phi}(\mathbf{x}) = \mathbb{E}(\mathbf{s}|\mathbf{x}). \tag{8}$$

Therefore, we have $s = \hat{s} + e$, that is,

$$\mathbf{s} = \phi(\mathbf{x}) + \mathbf{e} \tag{9}$$

where \mathbf{e} denotes the signal's estimation error; (9) implies that, in practice where $\mathbb{P}_{\mathbf{x},\mathbf{s}}$ is unknown, the datadriven estimation problem of \mathbf{s} given \mathbf{x} can be seen as a nonlinear regression problem of the collected data pairs $\{(x_1,s_1),(x_2,s_2),\ldots,(x_L,s_L)\}$. Usually, it is technically complicated to find the optimal $\phi(\cdot)$ from the whole space \mathcal{B} of Borel-measurable functions. Therefore, in practice, we may find the optimal approximation of $\phi(\cdot)$ in an RKHS \mathcal{H} or a NNFS \mathcal{K} ; note that \mathcal{H} and \mathcal{K} are two subspaces of \mathcal{B} . However, both \mathcal{H} and \mathcal{K} are sufficiently rich because they can be dense in the space of all continuous bounded functions.

2) Optimal Linear Estimation (Beamforming): When x and x are jointly Gaussian (e.g., when x and y are jointly Gaussian), the optimal estimator y is linear in x:

$$\hat{\mathbf{s}} = \mathbf{W}\mathbf{x},\tag{10}$$

where $\boldsymbol{W} \in \mathbb{C}^{M \times N}$ is called a smart antenna or a receive beamformer. In this linear case, (7) reduces to the usual Wiener-Hopf beamforming problem

$$\min_{\mathbf{W}} \operatorname{Tr} \mathbb{E}_{\mathbf{x}, \mathbf{s}} [\mathbf{W} \mathbf{x} - \mathbf{s}] [\mathbf{W} \mathbf{x} - \mathbf{s}]^{\mathsf{H}}, \tag{11}$$

that is,

$$\min_{\boldsymbol{W}} \operatorname{Tr} \left[\boldsymbol{W} \boldsymbol{R}_{x} \boldsymbol{W}^{\mathsf{H}} - \boldsymbol{W} \boldsymbol{R}_{xs} - \boldsymbol{R}_{xs}^{\mathsf{H}} \boldsymbol{W}^{\mathsf{H}} + \boldsymbol{R}_{s} \right], \quad (12)$$

where $\mathbf{R}_x := \mathbb{E}\mathbf{x}\mathbf{x}^{\mathsf{H}} \in \mathbb{C}^{N \times N}$ and $\mathbf{R}_{xs} := \mathbb{E}\mathbf{x}\mathbf{s}^{\mathsf{H}} \in \mathbb{C}^{N \times M}$. Since $\mathbf{R}_x = \mathbf{H}\mathbf{R}_s\mathbf{H}^{\mathsf{H}} + \mathbf{R}_v$ and $\mathbf{R}_{xs} = \mathbf{H}\mathbf{R}_s + \mathbb{E}\mathbf{v}\mathbf{s}^{\mathsf{H}} = \mathbf{H}\mathbf{R}_s$, the solution of (12), or (11), is

$$W_{\text{Wiener}}^{\star} = R_{xs}^{\mathsf{H}} R_x^{-1} = R_s H^{\mathsf{H}} [H R_s H^{\mathsf{H}} + R_v]^{-1},$$
 (13)

which is known as the Wiener beamformer. With an additional constraint $WH = I_M$ (i.e., distortionless response), (12) gives the Capon beamformer. Both the Wiener beamformer and the Capon beamformer maximize the output signal-to-interference-plus-noise ratio (SINR); hence, both are optimal in the sense of maximum output SINR. No matter whether

 $\mathbb{P}_{\mathbf{x},\mathbf{s}}$ is Gaussian or not, (11) or (12) gives the *optimal linear beamformer* in the sense of maximizing the output SINR among all linear beamformers.

- 3) Role of Channel Estimation: Eqs. (7) and (11) imply that channel estimation is not a necessary step in receive beamforming. The only necessary element, from the perspective of statistical machine learning, is the joint distribution $\mathbb{P}_{\mathbf{x},\mathbf{s}}$ of the received signal \mathbf{x} and the transmitted signal \mathbf{s} . Therefore, the following two points can be highlighted.
- 1) If the joint distribution $\mathbb{P}_{\mathbf{x},\mathbf{s}}$ is non-Gaussian, we just need to learn the mapping ϕ using (7).
- 2) If the joint distribution $\mathbb{P}_{\mathbf{x},\mathbf{s}}$ is (or assumed to be) Gaussian, we just need to learn the covariance matrices \mathbf{R}_{xs} and \mathbf{R}_{x} ; cf. (13). If, further, the channel matrix \mathbf{H} is known, \mathbf{R}_{xs} and \mathbf{R}_{x} can be expressed using \mathbf{H} .

B. Distributional Uncertainty and Distributional Robustness

For ease of conceptual illustration, we start with the following stationary-channel assumption in this subsection: The channel statistics remain unchanged within the communication frame so that the joint distribution $\mathbb{P}_{\mathbf{x},\mathbf{s}}$ is fixed over time. That is, pilot data $\{(x_1,s_1),(x_2,s_2),\ldots,(x_L,s_L)\}$ and nonpilot communication data are drawn from the same unknown distribution $\mathbb{P}_{\mathbf{x},\mathbf{s}}$. For the general case where the channel is not statistically stationary within a frame, see Appendix C of the online supplementary materials; the statistical non-stationarity of $\mathbb{P}_{\mathbf{x},\mathbf{s}}$ may be due to the time-selectivity of the transmit power matrix \mathbf{R}_s , of the channel matrix \mathbf{H} , and/or of the channel noise covariance \mathbf{R}_v .

1) Issue of Distributional Uncertainty: In practice, the true joint distribution $\mathbb{P}_{\mathbf{x},\mathbf{s}}$ is unknown but can be estimated by the pilot data. Hence, the estimation of wireless signals is a data-driven statistical inference (i.e., statistical machine learning) problem. We let

$$\hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}} := \frac{1}{L} \sum_{i=1}^{L} \delta_{(\boldsymbol{x}_i, \boldsymbol{s}_i)}$$
 (14)

denote the empirical distribution supported on the L collected data $\{(\boldsymbol{x}_i, \boldsymbol{s}_i)\}_{i \in [L]}$, where $\delta_{(\boldsymbol{x}_i, \boldsymbol{s}_i)}$ denotes the Dirac distribution (i.e., point-mass distribution) centered on $(\boldsymbol{x}_i, \boldsymbol{s}_i)$; note that $\hat{\mathbb{P}}_{\mathbf{x}, \mathbf{s}}$ is a discrete distribution. If we use the estimated joint distribution $\hat{\mathbb{P}}_{\mathbf{x}, \mathbf{s}}$ as a surrogate of the true joint distribution $\mathbb{P}_{\mathbf{x}, \mathbf{s}}$, (7) becomes the conventional empirical risk minimization (ERM) $\min_{\boldsymbol{\phi} \in \mathcal{B}} \operatorname{Tr} \mathbb{E}_{(\mathbf{x}, \mathbf{s}) \sim \hat{\mathbb{P}}_{\mathbf{x}, \mathbf{s}}}[\boldsymbol{\phi}(\mathbf{x}) - \mathbf{s}][\boldsymbol{\phi}(\mathbf{x}) - \mathbf{s}]^{\mathsf{H}}$, i.e.,

$$\min_{\boldsymbol{\phi} \in \mathcal{B}} \operatorname{Tr} \frac{1}{L} \sum_{i=1}^{L} [\boldsymbol{\phi}(\boldsymbol{x}_i) - \boldsymbol{s}_i] [\boldsymbol{\phi}(\boldsymbol{x}_i) - \boldsymbol{s}_i]^{\mathsf{H}}. \tag{15}$$

Likewise, (12) become the conventional beamforming problem

$$\min_{\boldsymbol{W}} \operatorname{Tr} \left[\boldsymbol{W} \hat{\boldsymbol{R}}_{x} \boldsymbol{W}^{\mathsf{H}} - \boldsymbol{W} \hat{\boldsymbol{R}}_{xs} - \hat{\boldsymbol{R}}_{xs}^{\mathsf{H}} \boldsymbol{W}^{\mathsf{H}} + \hat{\boldsymbol{R}}_{s} \right], \quad (16)$$

where \hat{R}_x , \hat{R}_{xs} , and \hat{R}_s are the training-sample-estimated (i.e., nominal) values of R_x , R_{xs} , and R_s , respectively.

There exists the distributional difference between the sample-defined nominal distribution $\hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}$ and true datagenerating distribution $\mathbb{P}_{\mathbf{x},\mathbf{s}}$ due to the limited size of the

training data set (i.e., limited pilot length) and the time-selectivity of $\mathbb{P}_{\mathbf{x},\mathbf{s}}$. From the perspective of applied statistics and machine learning, the distributional difference between $\hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}$ and $\mathbb{P}_{\mathbf{x},\mathbf{s}}$ (i.e., the distributional uncertainty of $\hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}$ compared to $\mathbb{P}_{\mathbf{x},\mathbf{s}}$) may cause significant performance degradation of (15) compared to (7), so is the performance deterioration of (16) compared to (12). For extensive reading on this point, see Appendix C of the online supplementary materials. Therefore, to reduce the adverse effect introduced by the distributional uncertainty in $\hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}$, a new surrogate of (7) rather than the sample-averaged approximation in (15) is expected.

2) Distibutionally Robust Estimation: To combat the distributional uncertainty in $\hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}$, we consider the distributionally robust counterpart of (7)

$$\min_{\boldsymbol{\phi} \in \mathcal{B}} \max_{\mathbb{P}_{\mathbf{x},\mathbf{s}} \in \mathcal{U}_{\mathbf{x},\mathbf{s}}} \operatorname{Tr} \mathbb{E}_{\mathbf{x},\mathbf{s}} [\boldsymbol{\phi}(\mathbf{x}) - \mathbf{s}] [\boldsymbol{\phi}(\mathbf{x}) - \mathbf{s}]^{\mathsf{H}},$$
 (17)

where $\mathcal{U}_{\mathbf{x},\mathbf{s}}$, called a distributional uncertainty set, contains a collection of distributions that are close to the nominal distribution (i.e., the sample-estimated distribution) $\hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}$;

$$\mathcal{U}_{\mathbf{x},\mathbf{s}} := \{ \mathbb{P}_{\mathbf{x},\mathbf{s}} | d(\mathbb{P}_{\mathbf{x},\mathbf{s}}, \hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}) \le \epsilon \}, \tag{18}$$

where $d(\cdot,\cdot)$ denotes a similarity measure (e.g., metric or divergence) between two distributions and $\epsilon \geq 0$ an uncertainty quantification level. Since $\hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}$ is discrete and $\mathbb{P}_{\mathbf{x},\mathbf{s}}$ is not, the Wasserstein distance [27, Def. 2] and the maximum mean discrepancy (MMD) distance [28, Def. 2.1] are the typical choices of $d(\cdot,\cdot)$ to construct $\mathcal{U}_{\mathbf{x},\mathbf{s}}$. When $\hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}$ and $\mathbb{P}_{\mathbf{x},\mathbf{s}}$ are parametric distributions (e.g., Gaussian, exponential family), divergences such as the Kullback—Leibler (KL) divergence, or more general ϕ -divergence, are also applicable to particularize $d(\cdot,\cdot)$ because parameters can be estimated using samples. When $\epsilon=0$, (17) reduces to (15).

If $\mathcal{U}_{\mathbf{x},\mathbf{s}}$ contains only Gaussian distributions, (17) is particularized to

$$\min_{\mathbf{W}} \max_{\mathbf{R}} \quad \text{Tr} \left[\mathbf{W} \mathbf{R}_{x} \mathbf{W}^{\mathsf{H}} - \mathbf{W} \mathbf{R}_{xs} - \mathbf{R}_{xs}^{\mathsf{H}} \mathbf{W}^{\mathsf{H}} + \mathbf{R}_{s} \right]
\text{s.t.} \quad d_{0}(\mathbf{R}, \ \hat{\mathbf{R}}) \leq \epsilon_{0},
\mathbf{R} \succeq \mathbf{0},$$
(19)

where

$$\boldsymbol{R} := \begin{bmatrix} \boldsymbol{R}_{x} & \boldsymbol{R}_{xs} \\ \boldsymbol{R}_{xs}^{\mathsf{H}} & \boldsymbol{R}_{s} \end{bmatrix}, \quad \hat{\boldsymbol{R}} := \begin{bmatrix} \hat{\boldsymbol{R}}_{x} & \hat{\boldsymbol{R}}_{xs} \\ \hat{\boldsymbol{R}}_{xs}^{\mathsf{H}} & \hat{\boldsymbol{R}}_{s} \end{bmatrix}, \quad (20)$$

because every zero-mean complex Gaussian distribution is uniquely characterized by its covariance and pseudo-covariance, but in receive beamforming, we do not consider pseudo-covariances; cf. (13); d_0 denotes the matrix similarity measures (e.g., matrix distances); $\epsilon_0 \geq 0$ is the uncertainty quantification parameter. When $\epsilon_0 = 0$, (19) reduces to (16).

IV. DISTRIBUTIONALLY ROBUST LINEAR ESTIMATION

Due to several practical benefits of linear estimation, for example, the simplicity of hardware structures, the clarity of physical meaning (i.e., constructive and destructive interference through beamforming), and the easiness of computations, investigating distributionally robust linear estimation problems is important. This section particularly studies Problem (19).

A. General Framework and Concrete Examples

The following lemma solves Problem (19).

Lemma 1: Suppose that the constraint $d_0(\mathbf{R}, \hat{\mathbf{R}}) \leq \epsilon_0$ is compact convex and \mathbf{R}_x is invertible. Let \mathbf{R}^* solve the problem below:

$$\begin{aligned} \max_{\boldsymbol{R}} & & \operatorname{Tr}\left[-\boldsymbol{R}_{xs}^{\mathsf{H}}\boldsymbol{R}_{x}^{-1}\boldsymbol{R}_{xs}+\boldsymbol{R}_{s}\right] \\ & \text{s.t.} & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & \\ & & & \\ & & \\ & & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & & \\ & \\ & & \\ & & \\ & \\ & & \\ & \\ & & \\$$

Construct W^{\star} using R^{\star} as follows:

$$W^* \coloneqq R_{xs}^{\mathsf{H}} R_x^{\mathsf{L}-1}. \tag{22}$$

Then (W^*, R^*) is a solution to Problem (19). On the other hand, if (W', R') solves Problem (19), then R' is a solution to (21) and (W', R') satisfies (22).

Proof: See Appendix D of the online supplementary materials. \Box

Let

$$f_1(\mathbf{R}) := \text{Tr}\left[-\mathbf{R}_{xs}^{\mathsf{H}}\mathbf{R}_x^{-1}\mathbf{R}_{xs} + \mathbf{R}_s\right]$$
 (23)

denote the objective function of (21). When R_s and R_{xs} are fixed, we define

$$f_2(\mathbf{R}_x) := \operatorname{Tr} \left[-\mathbf{R}_{xs}^{\mathsf{H}} \mathbf{R}_x^{-1} \mathbf{R}_{xs} + \mathbf{R}_s \right]. \tag{24}$$

The theorem below studies the properties of f_1 and f_2 .

Theorem 1: Consider the definition of \mathbf{R} in (20). The functions f_1 defined in (23) and f_2 defined in (24) are monotonically increasing in \mathbf{R} and \mathbf{R}_x , respectively. To be specific, if $\mathbf{R}_1 \succeq \mathbf{R}_2 \succeq \mathbf{0}$, $\mathbf{R}_{1,x} \succ \mathbf{0}$, and $\mathbf{R}_{2,x} \succ \mathbf{0}$, we have $f_1(\mathbf{R}_1) \geq f_1(\mathbf{R}_2)$. In addition, if $\mathbf{R}_{1,x} \succeq \mathbf{R}_{2,x} \succ \mathbf{0}$, we have $f_2(\mathbf{R}_{1,x}) \geq f_2(\mathbf{R}_{2,x})$.

Proof: See Appendix E of the online supplementary materials. \Box

To concretely solve (21), we need to particularize d_0 . This article investigates the following uncertainty sets.

Definition 1 (Additive Moment Uncertainty Set): The additive moment uncertainty set of R is constructed as

$$\hat{R} - \epsilon_0 E \prec R \prec \hat{R} + \epsilon_0 E, \tag{25}$$

for some $E \succeq \mathbf{0}$ and $\epsilon_0 \geq 0$ such that $\hat{\mathbf{R}} - \epsilon_0 \mathbf{E} \succeq \mathbf{0}$.

Definition 1 is motivated by the fact that the difference $R-\hat{R}$ is bounded by some threshold matrix E and error quantification level ϵ_0 : specifically, $-\epsilon_0 E \leq R - \hat{R} \leq \epsilon_0 E$. In practice, we can consider the threshold as an identity matrix because, for every non-identity $E \geq 0$, we have $E \leq \lambda_1 I_{N+M}$ where λ_1 is the largest eigenvalue of E.

Definition 2 (Diagonal-Loading Uncertainty Set): The diagonal-loading uncertainty set of ${\bf R}$ is constructed as

$$\hat{R} - \epsilon_0 I_{N+M} \leq R \leq \hat{R} + \epsilon_0 I_{N+M}, \tag{26}$$

for some $\epsilon_0 \geq 0$ such that $\hat{\boldsymbol{R}} - \epsilon_0 \boldsymbol{I}_{N+M} \succeq \boldsymbol{0}$.

Due to the concentration property of the sample-covariance \hat{R} to the true covariance R when the true distribution $\mathbb{P}_{\mathbf{x},\mathbf{s}}$ is fixed within a frame, finite values of ϵ_0 exist for every sample size L; NB: $\epsilon_0 \to 0$ as $L \to \infty$. However, given L, the smallest ϵ_0 cannot be practically calculated because it depends on the true but unknown $\mathbb{P}_{\mathbf{x},\mathbf{s}}$. If E is block-diagonal, the generalized diagonal-loading uncertainty set can be motivated.

Definition 3 (Generalized Diagonal-Loading Uncertainty Set): The generalized diagonal-loading uncertainty set of R is constructed as

$$\begin{bmatrix}
\hat{R}_{x} & \hat{R}_{xs} \\
\hat{R}_{xs}^{\mathsf{H}} & \hat{R}_{s}
\end{bmatrix} - \epsilon_{0} \begin{bmatrix} F & 0 \\ 0 & G \end{bmatrix}$$

$$\preceq \begin{bmatrix} R_{x} & R_{xs} \\ R_{xs}^{\mathsf{H}} & R_{s} \end{bmatrix}$$

$$\preceq \begin{bmatrix} \hat{R}_{x} & \hat{R}_{xs} \\ R_{xs}^{\mathsf{H}} & \hat{R}_{s} \end{bmatrix} + \epsilon_{0} \begin{bmatrix} F & 0 \\ 0 & G \end{bmatrix},$$
(27)

for some $F, G \succeq 0$ and $\epsilon_0 \geq 0$ such that the matrix in the first line is positive semi-definite.

Definitions 1, 2, and 3 are introduced for the first time in this article. Another type of moment-based uncertainty set is popular in the literature, which we refer to as the multiplicative moment uncertainty set for differentiation.

Definition 4 (Multiplicative Moment Uncertainty Set [29]): The multiplicative moment uncertainty set of R is given as

$$\theta_1 \hat{\mathbf{R}} \preceq \mathbf{R} \preceq \theta_2 \hat{\mathbf{R}},\tag{28}$$

for some $\theta_2 \ge 1 \ge \theta_1 \ge 0$.

The following corollary shows the distributionally robust linear beamformers associated with the various uncertainty sets in Definitions 1, 2, 3, and 4.

Corollary 1 (of Theorem 1): Consider the moment-based uncertainty sets in Definitions 1, 2, 3, and 4. The distributionally robust linear beamforming (21) is analytically solved by the corresponding upper bounds of \mathbf{R} . To be specific,

C1) Under Definition 1, the additive-moment distributionally robust (DR-AM) beamformer is

$$\begin{aligned} \boldsymbol{W}_{\mathrm{DR-AM}}^{\star} &= (\hat{\boldsymbol{R}}_{xs} + \epsilon_0 \boldsymbol{E}_{xs})^{\mathsf{H}} (\hat{\boldsymbol{R}}_x + \epsilon_0 \boldsymbol{E}_x)^{-1} \\ &= (\hat{\boldsymbol{H}} \hat{\boldsymbol{R}}_s + \epsilon_0 \boldsymbol{E}_{xs})^{\mathsf{H}} \cdot \\ &= [\hat{\boldsymbol{H}} \hat{\boldsymbol{R}}_s \hat{\boldsymbol{H}}^{\mathsf{H}} + \hat{\boldsymbol{R}}_v + \epsilon_0 \boldsymbol{E}_x]^{-1}. \end{aligned}$$

C2) Under Definition 2, the diagnal-loading distributionally robust (DR-DL) beamformer is

$$\begin{aligned} \boldsymbol{W}_{\text{DR-DL}}^{\star} &= \hat{\boldsymbol{R}}_{xs}^{\mathsf{H}} [\hat{\boldsymbol{R}}_{x} + \epsilon_{0} \boldsymbol{I}_{N}]^{-1} \\ &= \hat{\boldsymbol{R}}_{s} \hat{\boldsymbol{H}}^{\mathsf{H}} [\hat{\boldsymbol{H}} \hat{\boldsymbol{R}}_{s} \hat{\boldsymbol{H}}^{\mathsf{H}} + \hat{\boldsymbol{R}}_{v} + \epsilon_{0} \boldsymbol{I}_{N}]^{-1}, \end{aligned} \tag{30}$$

which is also known as the loaded sample matrix inversion method [19], [14, Eq. (11)] and widely-used in the practice of wireless communications.

C3) Under Definition 3, the generalized diagonal-loading distributionally robust beamformer (DR-GDL) is

$$\begin{aligned} \boldsymbol{W}_{\text{DR-GDL}}^{\star} &= \hat{\boldsymbol{R}}_{xs}^{\mathsf{H}} [\hat{\boldsymbol{R}}_{x} + \epsilon_{0} \boldsymbol{F}]^{-1} \\ &= \hat{\boldsymbol{R}}_{s} \hat{\boldsymbol{H}}^{\mathsf{H}} [\hat{\boldsymbol{H}} \hat{\boldsymbol{R}}_{s} \hat{\boldsymbol{H}}^{\mathsf{H}} + \hat{\boldsymbol{R}}_{v} + \epsilon_{0} \boldsymbol{F}]^{-1}. \end{aligned} \tag{31}$$

C4) Under Definition 4, the multiplicative-moment (MM) distributionally robust beamformer is identical to the Wiener beamformer (13) at nominal values:

The corresponding estimation errors are simple to obtain.
Corollary 1 implies that, in the sense of the same induced robust beamformers, the diagonal-loading uncertainty set (26)

and the generalized diagonal-loading uncertainty set (27) are technically equivalent to the following trimmed versions.

Definition 5 (Trimmed Diagonal-Loading Uncertainty Sets): By setting G := 0 in (27), in terms of R_x , (27) reduces to the trimmed generalized diagonal-loading uncertainty set:

$$\hat{\mathbf{R}}_x - \epsilon_0 \mathbf{F} \le \mathbf{R}_x \le \hat{\mathbf{R}}_x + \epsilon_0 \mathbf{F}. \tag{33}$$

The trimmed diagonal-loading uncertainty set

$$\hat{R}_x - \epsilon_0 I_N \leq R_x \leq \hat{R}_x + \epsilon_0 I_N, \tag{34}$$

is obtained by letting $F := I_N$.

The robust beamformers corresponding to the trimmed uncertainty sets (33) and (34) remain the same as defined in (31) and (30), respectively; cf. Theorem 1.

As we can see from Corollary 1, the primary benefit of using the moment-based uncertainty sets is the computational simplicity due to the availability of closed-form solutions. If the uncertainty sets are constructed using the Wasserstein distance $\sqrt{\text{Tr}[\mathbf{R}+\hat{\mathbf{R}}-2(\hat{\mathbf{R}}^{1/2}\mathbf{R}\hat{\mathbf{R}}^{1/2})^{1/2}]} \leq \epsilon_0$ or the KL divergence $\frac{1}{2}[\text{Tr}[\hat{\boldsymbol{R}}^{-1}\boldsymbol{R} - \boldsymbol{I}_{N+M}] - \ln\det(\hat{\boldsymbol{R}}^{-1}\boldsymbol{R})] \leq \epsilon_0$ between $\mathcal{CN}(0,R)$ and $\mathcal{CN}(0,R)$, the induced distributionally robust linear beamforming problems have no closed-form solutions, and therefore, are computationally prohibitive in practice. In addition, Corollary 1 suggests that the distributionally robust beamformer under the multiplicative moment uncertainty set (28) is the same as the nominal beamformer $\hat{R}_{rs}^{\mathsf{H}}\hat{R}_{r}^{-1}$, which essentially do not introduce robustness in wireless signal estimation; this is another motivation why we construct new moment-based uncertainty sets in Definitions 1, 2, and 3. However, we can modify the multiplicative moment uncertainty set in Definition 4 to achieve robustness.

Definition 6 (Modified Multiplicative Moment Uncertainty Set): The modified multiplicative moment uncertainty set of \boldsymbol{R} is given as

$$\begin{bmatrix} \theta_{1}\hat{\mathbf{R}}_{x} & \hat{\mathbf{R}}_{xs} \\ \hat{\mathbf{R}}_{xs}^{\mathsf{H}} & \theta_{1}\hat{\mathbf{R}}_{s} \end{bmatrix} \preceq \begin{bmatrix} \mathbf{R}_{x} & \mathbf{R}_{xs} \\ \mathbf{R}_{xs}^{\mathsf{H}} & \mathbf{R}_{s} \end{bmatrix} \preceq \begin{bmatrix} \theta_{2}\hat{\mathbf{R}}_{x} & \hat{\mathbf{R}}_{xs} \\ \hat{\mathbf{R}}_{xs}^{\mathsf{H}} & \theta_{2}\hat{\mathbf{R}}_{s} \end{bmatrix}$$
(35)

for some $\theta_2 \ge 1 \ge \theta_1 \ge 0$ such that the left-most matrix is positive semi-definite. \Box

The robust beamformer under the modified multiplicative moment uncertainty set (35) is $W_{\text{DR-MMM}}^{\star} = \hat{R}_{xs}^{\text{H}} \cdot [\theta_2 \hat{R}_x]^{-1}$. In terms of the uncertainties of R_s and R_v , Problem (21)

In terms of the uncertainties of R_s and R_v , Problem (21) can be explicitly written as

$$\max_{\boldsymbol{R}_{s},\boldsymbol{R}_{v}} \operatorname{Tr}\left[\boldsymbol{R}_{s} - \boldsymbol{R}_{s}\boldsymbol{H}^{\mathsf{H}}(\boldsymbol{H}\boldsymbol{R}_{s}\boldsymbol{H}^{\mathsf{H}} + \boldsymbol{R}_{v})^{-1}\boldsymbol{H}\boldsymbol{R}_{s}\right]$$
s.t.
$$d_{1}(\boldsymbol{R}_{s},\hat{\boldsymbol{R}}_{s}) \leq \epsilon_{1},$$

$$d_{2}(\boldsymbol{R}_{v},\hat{\boldsymbol{R}}_{v}) \leq \epsilon_{2},$$

$$\boldsymbol{R}_{s} \succeq \boldsymbol{0}, \ \boldsymbol{R}_{v} \succeq \boldsymbol{0},$$
(36)

for some similarity measures d_1 and d_2 and nonnegative scalars ϵ_1 and ϵ_2 . For every given $(\mathbf{R}_s, \mathbf{R}_v)$, the associated beamformer is given in (13). When the uncertainty in the channel matrix must be investigated, we can consider

$$\max_{\mathbf{H}} \operatorname{Tr} \left[\mathbf{R}_{s} - \mathbf{R}_{s} \mathbf{H}^{\mathsf{H}} (\mathbf{H} \mathbf{R}_{s} \mathbf{H}^{\mathsf{H}} + \mathbf{R}_{v})^{-1} \mathbf{H} \mathbf{R}_{s} \right]$$
s.t.
$$d_{3}(\mathbf{H}, \hat{\mathbf{H}}) \leq \epsilon_{3},$$
(37)

which is not a semi-definite program. In addition, the gradient of the objective function with respect to \boldsymbol{H} is complicated to obtain. Hence, practically, we should avoid directly attacking Problem (37); this can be done by directly considering the uncertainties of \boldsymbol{R}_x and \boldsymbol{R}_{xs} (i.e., \boldsymbol{R}) because the uncertainties of \boldsymbol{R}_s , \boldsymbol{R}_v , and \boldsymbol{H} can be reflected in the uncertainties of \boldsymbol{R}_x and \boldsymbol{R}_{xs} ; cf. $\boldsymbol{R}_x = \boldsymbol{H}\boldsymbol{R}_s\boldsymbol{H}^H + \boldsymbol{R}_v$ and $\boldsymbol{R}_{xs} = \boldsymbol{H}\boldsymbol{R}_s$.

In addition to Corollary 1, below we provide other concrete examples to further showcase the usefulness and applications of the distributionally robust beamforming formulations (21) and (36), where the trimmed uncertainty sets are employed.

Example 1 (Distributionally Robust Capon Beamforming): We consider a distributionally robust Capon beamforming problem under the trimmed uncertainty set (34):

$$\begin{aligned} \min \max_{\pmb{W}} \max_{\pmb{R}_x} & \operatorname{Tr} \left[\pmb{W} \pmb{R}_x \pmb{W}^\mathsf{H} - 2 \pmb{R}_s + \pmb{R}_s \right] \\ \text{s.t.} & \pmb{W} \pmb{H} = \pmb{I}_M, \\ & \hat{\pmb{R}}_x - \epsilon_0 \pmb{I}_N \preceq \pmb{R}_x \preceq \hat{\pmb{R}}_x + \epsilon_0 \pmb{I}_N, \\ & \pmb{R}_x \succ \pmb{0}, \end{aligned}$$

which is equivalent, in the sense of the same minimizer, to

$$\begin{aligned} \min \max_{\pmb{W}} & \operatorname{Tr} \left[\pmb{W} \pmb{R}_x \pmb{W}^\mathsf{H} \right] \\ \text{s.t.} & \pmb{W} \pmb{H} = \pmb{I}_M, \\ & \hat{\pmb{R}}_x - \epsilon_0 \pmb{I}_N \preceq \pmb{R}_x \preceq \hat{\pmb{R}}_x + \epsilon_0 \pmb{I}_N, \\ & \pmb{R}_x \succ \pmb{0}. \end{aligned}$$

According to Theorem 1, the above display is equivalent to

$$egin{aligned} \min_{m{W}} & \operatorname{Tr}\left[m{W}\hat{m{R}}_xm{W}^\mathsf{H}
ight] + \epsilon_0 \cdot \operatorname{Tr}\left[m{W}m{W}^\mathsf{H}
ight] \ & ext{s.t.} & m{W}m{H} = m{I}_M. \end{aligned}$$

The above formulation is the squared-F-norm-regularized Capon beamformer [14, Eq. (10)] whose solution is

$$W_{\text{DR-Capon}}^{\star} = [\boldsymbol{H}^{\mathsf{H}} (\hat{\boldsymbol{R}}_x + \epsilon_0 \boldsymbol{I}_N)^{-1} \boldsymbol{H}]^{-1} \cdot \boldsymbol{H}^{\mathsf{H}} (\hat{\boldsymbol{R}}_x + \epsilon_0 \boldsymbol{I}_N)^{-1},$$
(38)

which is the diagonal-loading Capon beamformer.

Example 2 (Eigenvalue Thresholding): Suppose that R_x admits the eigenvalues of diag $\{\lambda_1, \lambda_2, \dots, \lambda_N\}$ in descending order and the eigenvectors in Q (columns). Let $0 \le \mu \le 1$ be a shrinking coefficient. If we assume $R_x \le \hat{R}_{x,\text{thr}}$ where

we have the distributionally robust beamformer

$$\boldsymbol{W}_{\text{DR-ET}}^{\star} = \boldsymbol{R}_{xs}^{\mathsf{H}} \hat{\boldsymbol{R}}_{x,\text{thr}}^{-1},\tag{40}$$

which is known as the eigenvalue thresholding method [20], [14, Eq. (12)].

Example 3 (Distributionally Robust Beamforming for Uncertain \mathbf{R}_s and \mathbf{R}_v): Consider Problem (36). Since the objective of (36) is increasing in both \mathbf{R}_s and \mathbf{R}_v , if

$$\hat{R}_s - \epsilon_1 I_M \prec R_s \prec \hat{R}_s + \epsilon_1 I_M$$

¹This claim can be routinely proven in analogy to Theorem 1 and a real-space case in [30, Thm. 1].

we have a distributionally robust beamformer

$$W_{\text{DR}}^{\star} = (\hat{\boldsymbol{R}}_s + \epsilon_1 \boldsymbol{I}_M) \boldsymbol{H}^{\mathsf{H}} [\boldsymbol{H} (\hat{\boldsymbol{R}}_s + \epsilon_1 \boldsymbol{I}_M) \boldsymbol{H}^{\mathsf{H}} + \boldsymbol{R}_v]^{-1}$$

$$= (\hat{\boldsymbol{R}}_s + \epsilon_1 \boldsymbol{I}_M) \boldsymbol{H}^{\mathsf{H}} [\boldsymbol{H} \hat{\boldsymbol{R}}_s \boldsymbol{H}^{\mathsf{H}} + \boldsymbol{R}_v + \epsilon_1 \boldsymbol{H} \boldsymbol{H}^{\mathsf{H}}]^{-1};$$
(41)

if instead

$$\hat{\mathbf{R}}_s - \epsilon_1 \mathbf{H}^{\mathsf{H}} (\mathbf{H} \mathbf{H}^{\mathsf{H}})^{-2} \mathbf{H} \leq \mathbf{R}_s \leq \hat{\mathbf{R}}_s + \epsilon_1 \mathbf{H}^{\mathsf{H}} (\mathbf{H} \mathbf{H}^{\mathsf{H}})^{-2} \mathbf{H},$$
(42)

we have

$$W_{\mathrm{DR}}^{\star} = [\hat{R}_{s}H^{\mathrm{H}} + \epsilon_{1}H^{\mathrm{H}}(HH^{\mathrm{H}})^{-1}] \cdot [H\hat{R}_{s}H^{\mathrm{H}} + R_{v} + \epsilon_{1}I_{M}]^{-1}, \tag{43}$$

which is a modified diagonal-loading beamformer. On the other hand, if

$$\hat{\mathbf{R}}_v - \epsilon_2 \mathbf{I}_N \prec \mathbf{R}_v \prec \hat{\mathbf{R}}_v + \epsilon_2 \mathbf{I}_N$$

we have

$$\boldsymbol{W}_{\mathrm{DR}}^{\star} = \boldsymbol{R}_{s} \boldsymbol{H}^{\mathsf{H}} [\boldsymbol{H} \boldsymbol{R}_{s} \boldsymbol{H}^{\mathsf{H}} + \hat{\boldsymbol{R}}_{v} + \epsilon_{2} \boldsymbol{I}_{N}]^{-1}, \quad (44)$$

which is also a diagonal-loading beamformer.

Motivated by Corollary 1 and Examples $1\sim3$, as well as the trimmed uncertainty sets in Definition 5, we have the following important theorem, which justifies the popular ridge regression in machine learning.

Theorem 2 (Ridge Regression and Tikhonov Regularization): Consider a linear regression problem

$$s = Wx + e$$

and the distributionally robust linear estimator of W

$$\min_{\boldsymbol{W} \in \mathbb{C}^{M \times N}} \max_{\mathbb{P}_{\mathbf{x}, \mathbf{s}} \in \mathcal{U}_{\mathbf{x}, \mathbf{s}}} \operatorname{Tr} \mathbb{E}_{\mathbf{x}, \mathbf{s}} [\boldsymbol{W} \mathbf{x} - \mathbf{s}] [\boldsymbol{W} \mathbf{x} - \mathbf{s}]^{\mathsf{H}},$$

which can be particularized to (19). Supposing that the secondorder moment of x is uncertain and quantified as

$$\hat{R}_x - \epsilon_0 I_N \preceq R_x \preceq \hat{R}_x + \epsilon_0 I_N$$

then the distributionally robust linear estimator of \boldsymbol{W} becomes a ridge regression (i.e., squared-F-norm regularized) method

$$\min_{m{W}} \operatorname{Tr} \left[m{W} \hat{m{R}}_x m{W}^\mathsf{H} - m{W} \hat{m{R}}_{xs} - \hat{m{R}}_{xs}^\mathsf{H} m{W}^\mathsf{H} + \hat{m{R}}_s
ight] + \epsilon_0 \operatorname{Tr} \left[m{W} m{W}^\mathsf{H}
ight].$$

The regularization term becomes $\operatorname{Tr}[WFW^{H}]$, which is known as the Tikhonov regularizer, if

$$\hat{R}_x - \epsilon_0 F \prec R_x \prec \hat{R}_x + \epsilon_0 F$$

for some $F \succeq 0$.

Proof: This is due to Lemma 1 and Theorem 1. Just note that $\operatorname{Tr}\left[\boldsymbol{W}(\hat{\boldsymbol{R}}_x+\epsilon_0\boldsymbol{F})\boldsymbol{W}^{\mathsf{H}}-\boldsymbol{W}\hat{\boldsymbol{R}}_{xs}-\hat{\boldsymbol{R}}_{xs}^{\mathsf{H}}\boldsymbol{W}^{\mathsf{H}}+\hat{\boldsymbol{R}}_s\right]=\operatorname{Tr}\left[\boldsymbol{W}\hat{\boldsymbol{R}}_x\boldsymbol{W}^{\mathsf{H}}-\boldsymbol{W}\hat{\boldsymbol{R}}_{xs}-\hat{\boldsymbol{R}}_{xs}^{\mathsf{H}}\boldsymbol{W}^{\mathsf{H}}+\hat{\boldsymbol{R}}_s\right]+\epsilon_0\operatorname{Tr}\left[\boldsymbol{W}\boldsymbol{F}\boldsymbol{W}^{\mathsf{H}}\right].$ This completes the proof.

Note that in Theorem 2, the second-order moment of s is not considered because it does not influence the optimal solution of W: i.e., the optimal solution of W does not depend on the value of R_s . Theorem 2 gives a new theoretical interpretation of the popular ridge regression in machine learning from the perspective of distributional robustness against second-moment uncertainties of the feature vector \mathbf{x} ; another interpretation of ridge regression from the perspective of distributional

robustness under martingale constraints is identified in [31, Ex. 3.3]. When the uncertainty is quantified by the Wasserstein distance, a similar result can be seen in [32, Prop. 3], [33, Prop. 2], which however is not a ridge regression formulation because in [32, Prop. 3] and [33, Prop. 2], the loss function is square-rooted and the norm regularizer is not squared; cf. also [27, Rem. 18 and 19]. The corollary below justifies the rationale of any norm-regularized method.

Corollary 2: The following squared-norm-regularized beamforming formulation can combat the distributional uncertainty:

$$\min_{\boldsymbol{W}} \operatorname{Tr} \left[\boldsymbol{W} \hat{\boldsymbol{R}}_{x} \boldsymbol{W}^{\mathsf{H}} - \boldsymbol{W} \hat{\boldsymbol{R}}_{xs} - \hat{\boldsymbol{R}}_{xs}^{\mathsf{H}} \boldsymbol{W}^{\mathsf{H}} + \hat{\boldsymbol{R}}_{s} \right] + \lambda \| \boldsymbol{W} \|^{2},$$
(45)

where $\|\cdot\|$ denotes any matrix norm. This is because all norms on $\mathbb{C}^{M\times N}$ are equivalent; hence, there exists some $\lambda\geq 0$ such that $\lambda \|\boldsymbol{W}\|^2 \geq \epsilon_0 \|\boldsymbol{W}\|_F^2 = \epsilon_0 \operatorname{Tr}\left[\boldsymbol{W}\boldsymbol{W}^{\mathsf{H}}\right]$. As a result, (45) can upper bound the ridge cost in Theorem 2.

Motivated by Theorem 2, the following corollary is immediate, which gives another interpretation of ridge regression and Tikhonov regularization from the perspective of data augmentation through data perturbation (cf. noise injection in image [34] and speech [35] processing).

Corollary 3 (Data Augmentation for Linear Regression): Consider a linear regression problem on (\mathbf{x}, \mathbf{s}) with data perturbation vectors $(\boldsymbol{\Delta}_x, \boldsymbol{\Delta}_s)$

$$(\mathbf{s} + \boldsymbol{\Delta}_s) = \boldsymbol{W}(\mathbf{x} + \boldsymbol{\Delta}_x) + \mathbf{e},$$

and the distributionally robust linear estimator of W

$$\begin{split} & \min_{\boldsymbol{W} \in \mathbb{C}^{M \times N}} \max_{\mathbb{P}_{\boldsymbol{\Delta}_{x}, \boldsymbol{\Delta}_{s}} \in \mathcal{U}_{\boldsymbol{\Delta}_{x}, \boldsymbol{\Delta}_{s}}} \operatorname{Tr} \mathbb{E}_{(\mathbf{x}, \mathbf{s}) \sim \hat{\mathbb{P}}_{\mathbf{x}, \mathbf{s}}} \mathbb{E}_{\boldsymbol{\Delta}_{x}, \boldsymbol{\Delta}_{s}} \Big\{ \\ & [\boldsymbol{W}(\mathbf{x} + \boldsymbol{\Delta}_{x}) - (\mathbf{s} + \boldsymbol{\Delta}_{s})] [\boldsymbol{W}(\mathbf{x} + \boldsymbol{\Delta}_{x}) - (\mathbf{s} + \boldsymbol{\Delta}_{s})]^{\mathsf{H}} \Big\}. \end{split}$$

Suppose that Δ_x is uncorrelated with \mathbf{x} , with \mathbf{s} , and with Δ_s ; in addition, Δ_s is uncorrelated with \mathbf{x} . If the second-order moment of Δ_x is upper bounded as $\mathbb{E}\Delta_x\Delta_x^{\mathsf{H}} \leq \epsilon_0 \mathbf{I}_N$, then the distributionally robust linear estimator of \mathbf{W} becomes a ridge regression (i.e., squared-F-norm regularized) method

$$\min_{\mathbf{W}} \operatorname{Tr} \left[\mathbf{W} \hat{\mathbf{R}}_{x} \mathbf{W}^{\mathsf{H}} - \mathbf{W} \hat{\mathbf{R}}_{xs} - \hat{\mathbf{R}}_{xs}^{\mathsf{H}} \mathbf{W}^{\mathsf{H}} + \hat{\mathbf{R}}_{s} \right] + \epsilon_{0} \operatorname{Tr} \left[\mathbf{W} \mathbf{W}^{\mathsf{H}} \right].$$

The regularization term becomes $\operatorname{Tr}\left[\boldsymbol{W}\boldsymbol{F}\boldsymbol{W}^{\mathsf{H}}\right]$, which is known as the Tikhonov regularizer, if $\mathbb{E}\boldsymbol{\Delta}_{x}\boldsymbol{\Delta}_{x}^{\mathsf{H}} \leq \epsilon_{0}\boldsymbol{F}$, for some $\boldsymbol{F} \succeq \boldsymbol{0}$.

The second-order moment of Δ_s is not considered in Corollary 3 as it does not influence the optimal value of W.

B. Complex Uncertainty Sets

Below we remark on more general construction methods for the uncertainty set of R using the Wasserstein distance and the F-norm, beyond the moment-based methods in Definitions $1\sim6$. However, note that such complicated approaches are computationally prohibitive in practice when N or M is large.

1) Wasserstein Distributionally Robust Beamforming: We start with the Wasserstein distance:

$$\max_{\boldsymbol{R}} \quad \text{Tr} \left[-\boldsymbol{R}_{xs}^{\mathsf{H}} \boldsymbol{R}_{x}^{-1} \boldsymbol{R}_{xs} + \boldsymbol{R}_{s} \right]$$
s.t.
$$\quad \text{Tr} \left[\boldsymbol{R} + \hat{\boldsymbol{R}} - 2(\hat{\boldsymbol{R}}^{1/2} \boldsymbol{R} \hat{\boldsymbol{R}}^{1/2})^{1/2} \right] \leq \epsilon_{0}^{2} \qquad (46)$$

$$\boldsymbol{R} \succeq \boldsymbol{0}, \ \boldsymbol{R}_{x} \succ \boldsymbol{0}.$$

The first constraint in the above display is a particularization of the Wasserstein distance between $\mathcal{CN}(\mathbf{0}, \mathbf{R})$ and $\mathcal{CN}(\mathbf{0}, \hat{\mathbf{R}})$. Problem (46) is equivalent, in the sense of the same maximizer(s), to

$$\max_{\mathbf{R}} \operatorname{Tr} \mathbf{R}$$
s.t.
$$\operatorname{Tr} \left[\mathbf{R} + \hat{\mathbf{R}} - 2(\hat{\mathbf{R}}^{1/2} \mathbf{R} \hat{\mathbf{R}}^{1/2})^{1/2} \right] \leq \epsilon_0^2 \qquad (47)$$

$$\mathbf{R} \succeq \mathbf{0}, \ \mathbf{R}_x \succ \mathbf{0},$$

since the objective of (46) is increasing in \mathbf{R} ; cf. Theorem 1. Problem (47) is a nonlinear positive semi-definite program (P-SDP). However, we can give it a linear reformulation.

Proposition 1: Problem (47) can be equivalently reformulated into a linear P-SDP

$$\max_{\boldsymbol{R},\boldsymbol{U}} \operatorname{Tr} \boldsymbol{R}$$
s.t.
$$\operatorname{Tr} \left[\boldsymbol{R} + \hat{\boldsymbol{R}} - 2\boldsymbol{U} \right] \leq \epsilon_0^2$$

$$\left[\begin{array}{ccc} \hat{\boldsymbol{R}}^{1/2} \boldsymbol{R} \hat{\boldsymbol{R}}^{1/2} & \boldsymbol{U} \\ \boldsymbol{U} & \boldsymbol{I}_{N+M} \end{array} \right] \succeq \mathbf{0}$$

$$\boldsymbol{R} \succeq \mathbf{0}, \ \boldsymbol{R}_x \succ \mathbf{0}, \ \boldsymbol{U} \succeq \mathbf{0}.$$
(48)

Proof: This is by applying the Schur complement. \Box Complex-valued linear P-SDP can be solved using, e.g., the YALMIP solver.²

Suppose that \mathbf{R}^{\star} solves (48). The corresponding Wasserstein distributionally robust beamformer is given as

$$\boldsymbol{W}_{\text{DR-Wasserstein}}^{\star} = \boldsymbol{R}_{xs}^{\star \mathsf{H}} \boldsymbol{R}_{x}^{\star - 1}. \tag{49}$$

Next, we separately investigate the uncertainties in \hat{R}_s and \hat{R}_v . From (36), we have

$$\max_{\boldsymbol{R}_{s},\boldsymbol{R}_{v}} \operatorname{Tr} \left[\boldsymbol{R}_{s} - \boldsymbol{R}_{s} \boldsymbol{H}^{\mathsf{H}} (\boldsymbol{H} \boldsymbol{R}_{s} \boldsymbol{H}^{\mathsf{H}} + \boldsymbol{R}_{v})^{-1} \boldsymbol{H} \boldsymbol{R}_{s} \right]$$
s.t.
$$\operatorname{Tr} \left[\boldsymbol{R}_{s} + \hat{\boldsymbol{R}}_{s} - 2(\hat{\boldsymbol{R}}_{s}^{1/2} \boldsymbol{R}_{s} \hat{\boldsymbol{R}}_{s}^{1/2})^{1/2} \right] \leq \epsilon_{1}^{2}$$

$$\operatorname{Tr} \left[\boldsymbol{R}_{v} + \hat{\boldsymbol{R}}_{v} - 2(\hat{\boldsymbol{R}}_{v}^{1/2} \boldsymbol{R}_{v} \hat{\boldsymbol{R}}_{v}^{1/2})^{1/2} \right] \leq \epsilon_{2}^{2}$$

$$\boldsymbol{R}_{s} \succeq \boldsymbol{0}, \ \boldsymbol{R}_{v} \succeq \boldsymbol{0},$$

$$(50)$$

where we ignore the uncertainty of H for technical tractability. Problem (50) can be transformed into a linear P-SDP using a similar technique as in Proposition 1. One can just introduce an inequality $U \succeq R_s H^{\mathsf{H}} (HR_s H^{\mathsf{H}} + R)^{-1} HR_s$ and the objective function will become $\operatorname{Tr}[R_s - U]$.

Suppose that $(\mathbf{R}_s^{\star}, \mathbf{R}_v^{\star})$ solves (50). The corresponding Wasserstein distributionally robust beamformer is given as

$$W_{\text{DR-Wasserstein-Individual}}^{\star} = R_s^{\star} H^{\mathsf{H}} [H R_s^{\star} H^{\mathsf{H}} + R_v^{\star}]^{-1}. \quad (51)$$

2) F-Norm Distributionally Robust Beamforming: Under the F-norm, we just need to replace the Wasserstein distance. To be specific, for example, (47) becomes

$$\max_{\mathbf{R}} \operatorname{Tr} \mathbf{R}$$
s.t.
$$\operatorname{Tr} \left[(\mathbf{R} - \hat{\mathbf{R}})^{\mathsf{H}} (\mathbf{R} - \hat{\mathbf{R}}) \right] \leq \epsilon_0^2$$

$$\mathbf{R} \succeq \mathbf{0}, \ \mathbf{R}_x \succ \mathbf{0}.$$
(52)

The linear reformulation of the above display is given in the proposition below.

Proposition 2: The nonlinear P-SDP (52) can be equivalently reformulated into a linear P-SDP

$$\max_{\boldsymbol{R},\boldsymbol{U}} \quad \operatorname{Tr} \boldsymbol{R} \\
\text{s.t.} \quad \operatorname{Tr} [\boldsymbol{U}] \leq \epsilon_0^2, \\
\left[\begin{array}{cc} \boldsymbol{U} & (\boldsymbol{R} - \hat{\boldsymbol{R}})^{\mathsf{H}} \\ (\boldsymbol{R} - \hat{\boldsymbol{R}}) & \boldsymbol{I}_{N+M} \end{array} \right] \succeq \boldsymbol{0}, \\
\boldsymbol{R} \succeq \boldsymbol{0}, \ \boldsymbol{R}_x \succ \boldsymbol{0}, \ \boldsymbol{U} \succeq \boldsymbol{0}. \\$$
(53)

Proof: This is by applying the Schur complement.

C. Multi-Frame Case: Dynamic Channel Evolution

Each frame contains a pilot block used for beamformer design. Although the channel state information (CSI) may change from one frame to another, the CSI between the two consecutive frames is highly correlated. This correlation can benefit beamformer design across multiple frames. Suppose that $\{(s_1,x_1),(s_2,x_2),\ldots,(s_L,x_L)\}$ is the training data in the current frame and $\{(s'_1,x'_1),(s'_2,x'_2),\ldots,(s'_L,x'_L)\}$ is the history data in the immediately preceding frame. In such a case, the distributional difference between $\hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}$ and $\hat{\mathbb{P}}_{\mathbf{x}',\mathbf{s}'}$ is upper bounded, that is, $d(\hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}, \hat{\mathbb{P}}_{\mathbf{x}',\mathbf{s}'}) \leq \epsilon'$ for some proper distance d and a real number $\epsilon' \geq 0$ where $\hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}} \coloneqq \frac{1}{L} \sum_{i=1}^L \delta_{(x_i,s_i)}$ and $\hat{\mathbb{P}}_{\mathbf{x}',\mathbf{s}'} \coloneqq \frac{1}{L} \sum_{i=1}^L \delta_{(x_i',s_i')}$. Since a beamformer $\mathbf{W} = \mathcal{F}(\mathbb{P}_{\mathbf{x},\mathbf{s}})$ is a continuous function.

Since a beamformer $W = \mathcal{F}(\mathbb{P}_{\mathbf{x},\mathbf{s}})$ is a continuous functional $\mathcal{F}(\cdot)$ of data distribution $\mathbb{P}_{\mathbf{x},\mathbf{s}}$, cf. (11), we have $\|W - W'\|_F = \|\mathcal{F}(\hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}) - \mathcal{F}(\hat{\mathbb{P}}_{\mathbf{x}',\mathbf{s}'})\|_F \leq C \cdot d(\hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}},\,\hat{\mathbb{P}}_{\mathbf{x}',\mathbf{s}'}) \leq \epsilon$ for some positive constant $C \geq 0$ and upper bound $\epsilon \geq 0$ where W' is the beamformer associated with $\hat{\mathbb{P}}_{\mathbf{x}',\mathbf{s}'}$ in the previous frame. Therefore, the beamforming problem (12) becomes a constrained problem

$$egin{aligned} \min_{m{W}} & \operatorname{Tr}\left[m{W}m{R}_xm{W}^\mathsf{H} - m{W}m{R}_{xs} - m{R}_{xs}^\mathsf{H}m{W}^\mathsf{H} + m{R}_s
ight] \ & ext{s.t.} & \operatorname{Tr}[m{W} - m{W}'][m{W} - m{W}']^\mathsf{H} < \epsilon^2. \end{aligned}$$

By the Lagrange duality theory, it is equivalent to

$$\min_{\boldsymbol{W}} \operatorname{Tr} \left[\boldsymbol{W} \boldsymbol{R}_{x} \boldsymbol{W}^{\mathsf{H}} - \boldsymbol{W} \boldsymbol{R}_{xs} - \boldsymbol{R}_{xs}^{\mathsf{H}} \boldsymbol{W}^{\mathsf{H}} + \boldsymbol{R}_{s} \right] + \\
\lambda \cdot \operatorname{Tr} \left[\boldsymbol{W} - \boldsymbol{W}' \right] \left[\boldsymbol{W} - \boldsymbol{W}' \right]^{\mathsf{H}} \\
= \min_{\boldsymbol{W}} \operatorname{Tr} \left[\boldsymbol{W} (\boldsymbol{R}_{x} + \lambda \boldsymbol{I}_{N}) \boldsymbol{W}^{\mathsf{H}} - \boldsymbol{W} (\boldsymbol{R}_{xs} + \lambda \boldsymbol{W}'^{\mathsf{H}}) - \\
(\boldsymbol{R}_{xs} + \lambda \boldsymbol{W}'^{\mathsf{H}})^{\mathsf{H}} \boldsymbol{W}^{\mathsf{H}} + (\boldsymbol{R}_{s} + \lambda \boldsymbol{W}' \boldsymbol{W}'^{\mathsf{H}}) \right], \tag{54}$$

for some $\lambda \geq 0$. As a result, we have the Wiener beamformer for the multi-frame case, where we can treat \mathbf{W}' as a **prior knowledge** of \mathbf{W} .

Claim 1 (Multi-Frame Beamforming): The Wiener beamformer for the multi-frame case is given by

$$W_{\text{Wiener-MF}}^{\star} = [\mathbf{R}_{xs} + \lambda \mathbf{W}^{\prime \mathsf{H}}][\mathbf{R}_{x} + \lambda \mathbf{I}_{N}]^{-1}$$
$$= [\mathbf{R}_{s}\mathbf{H}^{\mathsf{H}} + \lambda \mathbf{W}^{\prime \mathsf{H}}][\mathbf{H}\mathbf{R}_{s}\mathbf{H}^{\mathsf{H}} + \mathbf{R}_{v} + \lambda \mathbf{I}_{N}]^{-1},$$
(55)

where $\lambda \geq 0$ is a tuning parameter to control the similarity between W and W'. Specifically, if λ is large, W must be close to W'; if λ is small, W can be far away from W'. \square With the result in Claim 1, (21) becomes

$$\max_{\mathbf{R}} \operatorname{Tr} \left[-(\mathbf{R}_{xs} + \lambda \mathbf{W}^{\prime \mathsf{H}})^{\mathsf{H}} \cdot (\mathbf{R}_{x} + \lambda \mathbf{I}_{N})^{-1} \cdot (\mathbf{R}_{xs} + \lambda \mathbf{W}^{\prime \mathsf{H}}) + (\mathbf{R}_{s} + \lambda \mathbf{W}^{\prime \mathsf{W}})^{\mathsf{H}} \right]$$
s.t.
$$d_{0}(\mathbf{R}, \hat{\mathbf{R}}) \leq \epsilon_{0},$$

$$\mathbf{R} \succeq \mathbf{0},$$
(56)

²See https://yalmip.github.io/inside/complexproblems/.

whose objective function is monotonically increasing in R.

The remaining distributional robustness modeling and analyses against the uncertainties in R are technically straightforward, and therefore, we omit them here. Upon using the diagonal-loading method on R, a distributionally robust beamformer for the multi-frame case is

$$\boldsymbol{W}_{\mathrm{DR-Wiener-MF}}^{\star} = [\hat{\boldsymbol{R}}_{xs} + \lambda \boldsymbol{W}^{\prime \mathsf{H}}] \cdot [\hat{\boldsymbol{R}}_{x} + \lambda \boldsymbol{I}_{N} + \epsilon_{0} \boldsymbol{I}_{N}]^{-1},$$

where ϵ_0 is an uncertainty quantification parameter for R.

V. DISTRIBUTIONALLY ROBUST NONLINEAR ESTIMATION

For the convenience of the technical treatment, we study the estimation problem in real spaces. Nonlinear estimators are to be limited in reproducing kernel Hilbert spaces and feedforward multi-layer neural network function spaces.

A. Reproducing Kernel Hilbert Spaces

1) General Framework and Concrete Examples: As a standard treatment in machine learning, we use the partial pilot data $\{\underline{x}_1,\underline{x}_2,\ldots,\underline{x}_L\}$ to construct the reproducing kernel Hilbert spaces, and use the whole pilot data $\{(\underline{x}_1,\underline{s}_1),(\underline{x}_2,\underline{s}_2),\ldots,(\underline{x}_L,\underline{s}_L)\}$ to train the optimal estimator in an RKHS.

With the **W**-linear representation of $\phi(\cdot)$ in (2), i.e., $\phi(\cdot)$ = $\boldsymbol{W}\boldsymbol{\varphi}(\cdot)$, the distributionally robust estimation problem (17)

$$\min_{\boldsymbol{W} \in \mathbb{R}^{2M \times L}} \max_{\mathbb{P}_{\underline{\mathbf{x}},\underline{\mathbf{s}}} \in \mathcal{U}_{\underline{\mathbf{x}},\underline{\mathbf{s}}}} \operatorname{Tr} \mathbb{E}_{\underline{\mathbf{x}},\underline{\underline{\mathbf{s}}}} [\boldsymbol{W} \cdot \boldsymbol{\varphi}(\underline{\mathbf{x}}) - \underline{\mathbf{s}}] [\boldsymbol{W} \cdot \boldsymbol{\varphi}(\underline{\mathbf{x}}) - \underline{\mathbf{s}}]^{\mathsf{T}}.$$
(57)

The proposition below reformulates and solves (57).

Proposition 3: Let K denote the kernel matrix associated with the kernel function $\ker(\cdot,\cdot)$ whose (i,j)-entry is defined as

$$\boldsymbol{K}_{i,j} := \ker(\underline{\boldsymbol{x}}_i, \underline{\boldsymbol{x}}_j), \quad \forall i, j \in [L].$$

Let $\underline{\mathbf{z}} := \varphi(\underline{\mathbf{x}})$. Then, the distributionally robust $\underline{\mathbf{x}}$ -nonlinear estimation problem (57) can be rewritten as a distributionally robust z-linear beamforming problem as

where $\hat{R}_{\underline{z}} = \frac{1}{L}K^2$, $\hat{R}_{\underline{z}\underline{s}} = \frac{1}{L}K\underline{S}^{\mathsf{T}}$, $\hat{R}_{\underline{s}} = \frac{1}{L}\underline{S}\underline{S}^{\mathsf{T}}$, and $\underline{S} := [\operatorname{Re} S; \operatorname{Im} S] = [\underline{s}_1, \underline{s}_2, \dots, \underline{s}_L]$. In addition, the strong minmax property holds for (58): i.e., the order of min and max can be exchanged provided that the first constraint is compact convex. As a result, given every pair of $(\mathbf{R}_z, \mathbf{R}_{zs}, \mathbf{R}_s)$, the optimal Wiener beamformer is³

$$W_{\text{RKHS}}^{\star} = \hat{R}_{\underline{z}\underline{s}}^{\mathsf{T}} \cdot \hat{R}_{\underline{z}}^{-1} = \underline{S}KK^{-2} = \underline{S}K^{-1},$$
 (59)

which transforms (58) to

$$\max_{\mathbf{R}_{\underline{z}}, \mathbf{R}_{\underline{z}\underline{s}}, \mathbf{R}_{\underline{s}}} \quad \operatorname{Tr} \left[-\mathbf{R}_{\underline{z}\underline{s}}^{\mathsf{T}} \mathbf{R}_{\underline{z}\underline{s}}^{-1} \mathbf{R}_{\underline{z}\underline{s}} + \mathbf{R}_{\underline{s}} \right]
s.t. \quad d_{0} \left(\begin{bmatrix} \mathbf{R}_{\underline{z}} & \mathbf{R}_{\underline{z}\underline{s}} \\ \mathbf{R}_{\underline{z}\underline{s}}^{\mathsf{T}} & \mathbf{R}_{\underline{s}} \end{bmatrix}, \begin{bmatrix} \hat{\mathbf{R}}_{\underline{z}} & \hat{\mathbf{R}}_{\underline{z}\underline{s}} \\ \hat{\mathbf{R}}_{\underline{z}\underline{s}}^{\mathsf{T}} & \hat{\mathbf{R}}_{\underline{s}} \end{bmatrix} \right) \leq \epsilon_{0},
\begin{bmatrix} \mathbf{R}_{\underline{z}} & \mathbf{R}_{\underline{z}\underline{s}} \\ \mathbf{R}_{\underline{z}\underline{s}}^{\mathsf{T}} & \mathbf{R}_{\underline{s}} \end{bmatrix} \succeq \mathbf{0}, \quad \mathbf{R}_{\underline{z}} \succ \mathbf{0}.$$
(60)

Proof: Treating [z; s] as, or approximating [z; s] using, a joint Gaussian random vector due to the linear estimation relation $\hat{\mathbf{g}} = \mathbf{W}\mathbf{z}$ in RKHS [cf. (57)], then the results in Lemma 1 apply. For details, see Appendix F of the online supplementary materials.

In (58), d_0 defines a matrix similarity measure to quantify the uncertainty of the covariance matrix of $[\underline{\mathbf{z}};\underline{\mathbf{s}}]$, and $\epsilon_0 \geq 0$ quantifies the uncertainty level. Proposition 3 reveals the benefit of the kernel trick (2), that is, the possibility to represent a nonlinear estimation problem as a linear one.

The claim below summarizes the solution of (17) in the RKHS induced by the kernel function $ker(\cdot, \cdot)$.

Claim 2: Suppose that $(R_z^{\star}, R_{zs}^{\star}, R_s^{\star})$ solves (60). Then the optimal estimator of (17) in the \overline{RKHS} induced by $\ker(\cdot,\cdot)$ is given by

$$\phi^{\star}(\mathbf{x}) = \mathbf{\Gamma}_{M} \cdot \mathbf{R}_{\underline{z}\underline{s}}^{\star \mathsf{T}} \cdot \mathbf{R}_{\underline{z}}^{\star - 1} \cdot \boldsymbol{\varphi}(\underline{\mathbf{x}}), \tag{61}$$

where $\underline{\mathbf{x}} = [\text{Re } \mathbf{x}; \text{ Im } \mathbf{x}]$ is the real-space representation of \mathbf{x} , $\Gamma_M \coloneqq [I_M, J_M]$ is defined in Subsection I-B, and

$$oldsymbol{arphi}(\underline{\mathbf{x}}) \coloneqq \left[egin{array}{c} \ker(\underline{\mathbf{x}},\underline{oldsymbol{x}}_1) \ \ker(\underline{\mathbf{x}},\underline{oldsymbol{x}}_2) \ dots \ \ker(\underline{\mathbf{x}},\underline{oldsymbol{x}}_L) \end{array}
ight].$$

In addition, the corresponding worst-case estimation error covariance is

$$\Gamma_{M} \cdot \left[-R_{\underline{z}\underline{s}}^{\star \mathsf{T}} R_{\underline{z}}^{\star -1} R_{\underline{z}\underline{s}}^{\star} + R_{\underline{s}}^{\star} \right] \cdot \Gamma_{M}^{\mathsf{H}}, \tag{62}$$

which upper bounds the true estimation error covariance. Concrete examples of Claim 2 are given as follows.

Example 4 (Kernelized Diagonal Loading): By using the trimmed diagonal-loading uncertainty set for R_z , i.e.,

$$\hat{\boldsymbol{R}}_{\underline{z}} - \epsilon_0 \boldsymbol{I}_L \preceq \boldsymbol{R}_{\underline{z}} \preceq \hat{\boldsymbol{R}}_{\underline{z}} + \epsilon_0 \boldsymbol{I}_L$$

$$\phi^{\star}(\mathbf{x}) = \mathbf{\Gamma}_{M} \cdot \frac{1}{L} \underline{\mathbf{S}} \mathbf{K} \cdot \left(\frac{1}{L} \mathbf{K}^{2} + \epsilon_{0} \mathbf{I}_{L}\right)^{-1} \cdot \boldsymbol{\varphi}(\underline{\mathbf{x}}), \quad (63)$$

which is obtained at the upper bound of R_z . Furthermore, in this case, the distributionally robust formulation (57) is equivalent to a squared-F-norm-regularized formulation

$$\min_{\boldsymbol{W}} \operatorname{Tr} \mathbb{E}_{(\underline{\mathbf{x}},\underline{\mathbf{s}}) \sim \hat{\mathbb{P}}_{\underline{\mathbf{x}},\underline{\mathbf{s}}}} [\boldsymbol{W} \cdot \boldsymbol{\varphi}(\underline{\mathbf{x}}) - \underline{\mathbf{s}}] [\boldsymbol{W} \cdot \boldsymbol{\varphi}(\underline{\mathbf{x}}) - \underline{\mathbf{s}}]^{\mathsf{T}} + \epsilon_0 \cdot \operatorname{Tr} [\boldsymbol{W} \boldsymbol{W}^{\mathsf{T}}],$$
(64)

which can be proven by replacing R_z in (58) with its upper bound.

Example 5 (Kernelized Eigenvalue Thresholding): The kernelized eigenvalue thresholding method can be designed in

³For numerical stability in the inversion of K, L should not be too large.

analogy to Example 2. The two key steps are to obtain the eigenvalue decomposition of $\hat{R}_{\underline{z}} = K^2/L$ and then lift the eigenvalues; cf. (39).

In addition, Example 4 motivates the following important theorem for statistical machine learning.

Theorem 3 (Kernel Ridge Regression and Kernel Tikhonov Regularization): Consider the nonlinear regression problem

$$s = \phi(x) + e$$

and the distributionally robust estimator of $\phi(\underline{\mathbf{x}}) = W \cdot \varphi(\underline{\mathbf{x}})$ in the RKHS induced by the kernel function $\ker(\cdot, \cdot)$, i.e.,

$$\min_{\boldsymbol{W} \in \mathbb{R}^{2M \times L}} \max_{\mathbb{P}_{\underline{\mathbf{x}},\underline{\mathbf{s}}} \in \mathcal{U}_{\underline{\mathbf{x}},\underline{\mathbf{s}}}} \operatorname{Tr} \mathbb{E}_{\underline{\mathbf{x}},\underline{\underline{\mathbf{s}}}} [\boldsymbol{W} \cdot \boldsymbol{\varphi}(\underline{\mathbf{x}}) - \underline{\mathbf{s}}] [\boldsymbol{W} \cdot \boldsymbol{\varphi}(\underline{\mathbf{x}}) - \underline{\mathbf{s}}]^{\mathsf{T}}.$$

Supposing that only the second-order moment of $\underline{\mathbf{z}} := \varphi(\underline{\mathbf{x}})$ is uncertain and quantified as

$$\hat{\boldsymbol{R}}_{\underline{z}} - \epsilon_0 \boldsymbol{I}_L \preceq \boldsymbol{R}_{\underline{z}} \preceq \hat{\boldsymbol{R}}_{\underline{z}} + \epsilon_0 \boldsymbol{I}_L,$$

then the distributionally robust estimator of W becomes a kernel ridge regression method (64). The regularization term in (64) becomes the Tikhonov regularizer $\text{Tr}[WFW^{\mathsf{T}}]$ if

$$\hat{\boldsymbol{R}}_z - \epsilon_0 \boldsymbol{F} \preceq \boldsymbol{R}_z \preceq \hat{\boldsymbol{R}}_z + \epsilon_0 \boldsymbol{F}$$

for some $F \succ 0$.

Proof: See Example 4; cf. Theorem 2.
$$\Box$$

Theorem 3 gives the kernel ridge regression an interpretation of distributional robustness. The usual choice of F in Theorem 3 is the L-divided kernel matrix K/L; see, e.g., [36, Eq. (4)], [24, Eqs. (15.110) and (15.113)]. As a result, from (63), we have

$$\phi^{\star}(\mathbf{x}) = \mathbf{\Gamma}_{M} \cdot \underline{\mathbf{S}} \cdot (\mathbf{K} + \epsilon_{0} \mathbf{I}_{L})^{-1} \cdot \varphi(\underline{\mathbf{x}}), \tag{65}$$

which is another type of kernel ridge regression (i.e., a new kernelized diagonal-loading method).

In analogy to Corollary 2, the following corollary motivated from (64) is immediate.

Corollary 4: The following squared-norm-regularized method in RKHSs can combat the distributional uncertainty:

$$\min_{\boldsymbol{W}} \operatorname{Tr} \mathbb{E}_{(\underline{\mathbf{x}},\underline{\mathbf{s}}) \sim \hat{\mathbb{P}}_{\underline{\mathbf{x}},\underline{\mathbf{s}}}} [\boldsymbol{W} \cdot \boldsymbol{\varphi}(\underline{\mathbf{x}}) - \underline{\mathbf{s}}] [\boldsymbol{W} \cdot \boldsymbol{\varphi}(\underline{\mathbf{x}}) - \underline{\mathbf{s}}]^{\mathsf{T}} + \\
\lambda \cdot ||\boldsymbol{W}||^{2}, \tag{66}$$

for any matrix norm $\|\cdot\|$; cf. Corollary 2.

Moreover, in analogy to Corollary 3, the following corollary is immediate.

Corollary 5 (Data Augmentation for Kernel Regression): Consider the nonlinear regression problem in Theorem 3. Its data-perturbed counterpart can be constructed by taking into account the data perturbation vectors $(\Delta_{\underline{s}}, \Delta_{\underline{z}})$. Suppose that $\Delta_{\underline{s}}$ is uncorrelated with \underline{z} , with \underline{s} , and with $\Delta_{\underline{s}}$; in addition, $\Delta_{\underline{s}}$ is uncorrelated with \underline{z} . If the second-order moment of $\Delta_{\underline{s}}$ is upper bounded by $\epsilon_0 I_L$, then the distributionally robust estimator of W becomes a kernel ridge regression (i.e., squared-F-norm regularized) method (64). The regularization term becomes $\mathrm{Tr}\left[WFW^{\mathrm{H}}\right]$, which is known as the Tikhonov regularizer, if the second-order moment of $\Delta_{\underline{s}}$ is upper bounded by $\epsilon_0 F$ for some $F \succeq 0$.

General uncertainty sets using the Wasserstein distance or the F-norm, beyond the diagonal ϵ_0 -perturbation (cf. Example

4), can be straightforwardly employed and the distributional robustness modeling and analyses remain routine; cf. Subsection IV-B. Hence, we omit them here. However, such complicated approaches are computationally prohibitive in practice when L or M is large.

2) Multi-Frame Case: Dynamic Channel Evolution: As in (54), the multi-frame formulation in RKHSs is

$$\min_{\boldsymbol{W} \in \mathbb{R}^{2M \times L}} \operatorname{Tr} \mathbb{E}_{\underline{\mathbf{x}},\underline{\mathbf{s}}} [\boldsymbol{W} \cdot \boldsymbol{\varphi}(\underline{\mathbf{x}}) - \underline{\mathbf{s}}] [\boldsymbol{W} \cdot \boldsymbol{\varphi}(\underline{\mathbf{x}}) - \underline{\mathbf{s}}]^{\mathsf{T}} + \\
\lambda \cdot \operatorname{Tr} [\boldsymbol{W} - \boldsymbol{W}'] [\boldsymbol{W} - \boldsymbol{W}']^{\mathsf{T}}, \tag{67}$$

where W' denotes the beamformer in the immediately preceding frame and serves as a *prior knowledge* of W.

Claim 3 (Multi-Frame Estimation in RHKS): The solution of (67) is given by (cf. (59))

$$W_{\text{RKHS-MF}}^{\star} = [\mathbf{R}_{\underline{z}\underline{s}} + \lambda \mathbf{W}^{\prime\mathsf{T}}][\mathbf{R}_{\underline{z}} + \lambda \mathbf{I}_{L}]^{-1}$$

$$= \left(\frac{1}{L}\underline{\mathbf{S}}\mathbf{K} + \lambda \mathbf{W}^{\prime\mathsf{T}}\right) \cdot \left(\frac{1}{L}\mathbf{K}^{2} + \lambda \mathbf{I}_{L}\right)^{-1},$$
(68)

where $\lambda \geq 0$ is a tuning parameter to control the similarity between W and W'; cf. Claim 1.

The remaining distributional robustness modeling and analyses on (67) against the uncertainties in $\hat{R}_{\underline{z}}$, $\hat{R}_{\underline{x}\underline{z}}$, and $\hat{R}_{\underline{s}}$ are technically straightforward; cf. Subsection IV-C. Therefore, we omit them here.

B. Neural Networks

With the W-parameterization $\phi_{W_{[R]}}(\underline{\mathbf{x}})$ of $\phi(\underline{\mathbf{x}})$ in feedforward multi-layer neural networks, i.e., (4), the distributionally robust estimation problem (17) becomes

$$\min_{\boldsymbol{W}_{[R]}} \max_{\mathbb{P}_{\underline{\mathbf{x}},\underline{\mathbf{s}}} \in \mathcal{U}_{\underline{\mathbf{x}},\underline{\mathbf{s}}}} \operatorname{Tr} \mathbb{E}_{\underline{\mathbf{x}},\underline{\mathbf{s}}} [\phi_{\boldsymbol{W}_{[R]}}(\underline{\mathbf{x}}) - \underline{\mathbf{s}}] [\phi_{\boldsymbol{W}_{[R]}}(\underline{\mathbf{x}}) - \underline{\mathbf{s}}]^{\mathsf{T}}, (69)$$

where $W_{[R]} := \{W_1, W_2, \dots, W_R\}$ and $\phi_{W_{[R]}}(\underline{\mathbf{x}})$ is defined in (4). Problem (69) is highly nonlinear in both argument $\underline{\mathbf{x}}$ and parameter $W_{[R]}$, which is different from the case in reproducing kernel Hilbert spaces where the $W_{[R]}$ -linearization features. Hence, problem (69) is too complicated to solve to global optimality. According to [27, Cor. 33], under several technical conditions (plus the boundedness of the feasible region of $W_{[R]}$), (69) is upper bounded by a spectral-norm-regularized empirical risk minimization problem

$$\min_{\boldsymbol{W}_{[R]}} \frac{1}{L} \sum_{i=1}^{L} \operatorname{Tr}[\boldsymbol{\phi}_{\boldsymbol{W}_{[R]}}(\underline{\boldsymbol{x}}_{i}) - \underline{\boldsymbol{s}}_{i}][\boldsymbol{\phi}_{\boldsymbol{W}_{[R]}}(\underline{\boldsymbol{x}}_{i}) - \underline{\boldsymbol{s}}_{i}]^{\mathsf{T}} + \\
\lambda' \cdot \sum_{r=1}^{R} \|\boldsymbol{W}_{r}\|_{2}, \tag{70}$$

for some regularization coefficient $\lambda' \geq 0$, where $\|\cdot\|_2$ denotes the spectral norm of a matrix (i.e., the induced 2-norm). Eq. (70) rigorously justifies the popular norm regularization method in training neural networks: By reducing the upper bound (70) of (69), the value of (69) can be diminished as well. The regularized ERM problem (70) is reminiscent of the ridge regression and the kernel ridge regression methods in Theorems 2 and 3 for distributional robustness in linear regression and RKHS linear regression, respectively. Supposing that $\boldsymbol{W}_{[R]}^{\star}$ is an (approximated, or sub-optimal) solution⁴ of

⁴Neural networks are hard to be globally optimized.

(70), then the distributionally robust optimal estimator of the transmitted signal s can be obtained as

$$\hat{\mathbf{s}} = \mathbf{\Gamma}_M \cdot \boldsymbol{\phi}_{\mathbf{W}_{[R]}^{\star}}(\underline{\mathbf{x}}).$$

Therefore, in training a neural network for wireless signal estimation, it is recommended to apply the norm regularization methods. Since norms on real spaces are equivalent, (70) can be further upper bounded by

$$\min_{\boldsymbol{W}_{[R]}} \frac{1}{L} \sum_{i=1}^{L} \operatorname{Tr}[\boldsymbol{\phi}_{\boldsymbol{W}_{[R]}}(\underline{\boldsymbol{x}}_{i}) - \underline{\boldsymbol{s}}_{i}][\boldsymbol{\phi}_{\boldsymbol{W}_{[R]}}(\underline{\boldsymbol{x}}_{i}) - \underline{\boldsymbol{s}}_{i}]^{\mathsf{T}} + \\
\lambda \cdot \sum_{r=1}^{R} \|\boldsymbol{W}_{r}\|, \tag{71}$$

for any matrix norm $\|\cdot\|$ and some $\lambda \geq 0$; λ depends on λ' and $\|\cdot\|$. As a result, to achieve distributional robustness in training a neural network, any-norm-regularized learning method in (71) can be considered.

VI. EXPERIMENTS

We consider a point-to-point multiple-input-multiple-output (MIMO) wireless communication problem where the transmitter is located at [0,0] and the receiver is at $[500\mathrm{m},450\mathrm{m}]$. We randomly sample 25 points according to the uniform distribution on the square of $[0,500\mathrm{m}] \times [0,500\mathrm{m}]$ to denote the scatters' positions; i.e., there exist 25 radio paths. All the source data and codes are available online at GitHub with thorough implementation comments: https://github.com/Spratm-Asleaf/Beamforming. In this section, we only present major experimental setups and results; readers can use the shared source codes to explore (or verify) minor ones.

The following eleven methods are implemented in the experiments: 1) Wiener: Wiener beamformer (13), upper expression; 2) Wiener-DL: Wiener beamformer with diagonal loading (30), upper expression; 3) Wiener-DR: Distributionally robust Wiener beamformer (49) and (53); 4) Wiener-CE: Channel-estimation-based Wiener beamformer (13), lower expression; 5) Wiener-CE-DL: Channel-estimation-based Wiener beamformer with diagonal loading (30), lower expression; 6) Wiener-CE-DR: Distributionally robust channelestimation-based Wiener beamformer (41) and (31); 7) Capon: Capon beamformer (38) for $\epsilon_0 = 0$; 8) **Capon-DL**: Capon beamformer with diagonal loading (38); 9) ZF: Zero-forcing beamformer where $W_{\rm ZF}\coloneqq (\hat{H}^{\sf H}\hat{H})^{-1}\hat{H}^{\sf H}$ and \hat{H} denotes the estimated channel matrix; 10) Kernel: Kernel receiver (61) with $\epsilon_0 = 0$ in (60); and 11) **Kernel-DL**: Kernel receiver with diagonal loading (65). Note that the diagonal-loadingbased methods are particular cases of distributionally robust beamformers; see, e.g., Corollary 1 and Example 4. The deeplearning-based (DL-based) methods in Subsection V-B are not implemented in this section because they have been deeply studied in our previous publications, e.g., [10], [12]; we only comment on the advantages and disadvantages of DL-based methods compared with the listed eleven methods in Section VII (Conclusions).

When covariance matrix R_s of transmitted signal s is unknown for the receiver (e.g., in ISAC systems, R_s needs to

vary from one frame to another for sensing), R_s is estimated by the sample covariance matrix $\hat{R}_s = SS^H/L$. The channel matrix H is estimated using the minimum mean-squared error method, i.e., $\hat{H} = XS^H(SS^H)^{-1}$. Covariance matrix R_v of channel noise v is estimated using the least-square method, i.e., $\hat{R}_v = (X - \hat{H}S)(X - \hat{H}S)^H/L$. The matrices \hat{R}_s , \hat{H} , and \hat{R}_v are therefore uncertain compared to their true (but unknown; possibly time-varying) values R_s , H, and R_v , respectively. The matrices \hat{R}_s , \hat{H} , and \hat{R}_v are used in beamformers such as the channel-estimation-based Wiener beamformer (30), the Capon beamformer, and the zero-forcing beamformer.

The beamformers are determined on the training data set (i.e., pilot data). The performance evaluation method of beamformers is mean-squared estimation error (MSE) on the test data set (i.e., non-pilot communication data): to be specific, $\|S_{\text{test}} - \hat{S}_{\text{test}}\|_F^2/(M \times L_{\text{test}})$ where $S_{\text{test}} \in \mathbb{C}^{M \times L_{\text{test}}}$ is the test data block, \hat{S}_{test} is its estimate, and L_{test} is the length of non-pilot test data units. As data-driven machine learning methods, all parameters (e.g., uncertainty quantification coefficients ϵ 's) of beamformers can be tuned using the popular cross-validation (e.g., one-shot cross-validation) method. The parameters can also be empirically tuned to save training times because cross-validation imposes a significant computational burden. This article mainly uses the empirical tuning method (i.e., trial-and-error) to tune each beamformer to achieve its best average performance. For each test case, the MSE performances are averaged on 250 Monte-Carlo episodes.

We consider an experimental scenario where impulse channel noises exist; i.e., the channel is non-Gaussian so linear beamformers are no longer sufficient. (Complementary experimental setups and results can be seen in Appendix G of the online supplementary materials.) The detailed setups are as follows. The transmitter has four antennas (i.e., M=4) with unit transmit power; without loss of generality, each antenna is assumed to emit continuous-valued complex Gaussian signals. The receiver has eight antennas (i.e., N=8). The SNR is -10dB, which is a challenging situation. The channel has impulse noises: i.e., in L received signals (i.e., $[\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_L]$) that are contaminated by usual complex Gaussian channel noises, 10% of them are also contaminated by uniform noises with the maximum amplitude of 1.5, which is a relatively large value compared to the amplitude of the usual Gaussian channel noises. We assume that a communication frame contains 500 non-pilot data units; i.e., $L_{\text{test}} = 500$. The experimental results are shown in Tables I~VI, from which the following main points can be outlined.

- 1) A larger number of pilot data benefits the estimation performances of wireless signals.
- The diagonal loading operation can significantly improve the estimation performances especially when the pilot data size is relatively small.
- 3) Since the signal model under impulse channel noises is no longer linear Gaussian, the optimal beamformer in the MSE sense must be nonlinear. Therefore, the Kernel and the Kernel-DL methods have the potential to outperform other linear beamformers, i.e., to *suppress outliers*. However, in practice, the non-robust Kernel method may

- undergo numerical instability in calculating the inverse of the kernel matrix K. Therefore, its actual MSEs are not necessarily smaller than those of linear beamformers. Nevertheless, the robust Kernel-DL method consistently outperforms all other beamformers.
- 4) Distributionally robust beamformers (including diagonal-loading ones) can combat the adverse effect introduced by the limited pilot size and several types of uncertainties in the signal model (e.g., outliers). To be specific, for example, all diagonal-loading beamformers can outperform their original non-diagonal-loading counterparts; cf. the Wiener and the Wiener-DL methods, the Wiener-CE and the Wiener-CE-DL methods, the Capon and the Capon-DL methods, and the Kernel and the Kernel-DL methods. In addition, the Wiener-DR beamformer (53) using the *F*-norm uncertainty set has the potential to outperform the Wiener-DL beamformer (30) that employs the simple uncertainty set (26).
- 5) Although the Wiener-DR beamformer has the potential to work better than the Wiener-DL beamformer, it has a significant computational burden, which may not be suitable for timely use in practice especially when the computing resources are limited. Hence, the Wiener-DL beamformer is practically promising because it can provide an excellent balance between the computational burden and the actual performance.

Remarks on Parameter Tuning: From experiments, we find that the uncertainty quantification coefficients ϵ 's (e.g., in diagonal loading) can be neither too large nor too small. When ϵ 's are too large, the beamformers become overly conservative, while when ϵ 's are too small, the beamformers cannot offer sufficient robustness against data scarcity and model uncertainties. In both cases of inappropriate ϵ 's, the performances of beamformers degrade significantly. Therefore, ϵ 's must be carefully tuned in practice, and a rigorous method to tune ϵ 's can be the cross-validation method on the training data set (i.e., the pilot data set). If practitioners just pursue satisfaction rather than optimality, the empirical tuning method is recommended to save training times.

TABLE I
EXPERIMENTAL RESULTS (PILOT SIZE = 10)

Beamformer	MSE	Time	Beamformer	MSE	Time
Wnr	3.30	1.49e-04	Wnr-DL	2.11	9.81e-06
Wnr-DR	1.97	3.16e+00	Wnr-CE	3.30	4.59e-05
Wnr-CE-DL	2.50	2.17e-05	Wnr-CE-DR	3.31	4.63e-05
Capon	5.44	4.42e-05	Capon-DL	4.52	2.50e-05
ZF	2.12	2.54e-05	Kernel	1.07	1.60e-04
Kernel-DL	0.80	5.59e-05			

Wnr: The abbreviation for Wiener.

Time: The training time averaged on 250 Monte-Carlo episodes.

VII. CONCLUSIONS

This article introduces a unified mathematical framework for receive beamforming of wireless signals from the perspective of data-driven machine learning, which reveals that channel estimation is not a necessary operation. To combat the limited

TABLE II
EXPERIMENTAL RESULTS (PILOT SIZE = 15)

Beamformer	MSE	Time	Beamformer	MSE	Time
Wnr	1.38	1.65e-04	Wnr-DL	1.23	1.10e-05
Wnr-DR	1.07	3.21e+00	Wnr-CE	1.38	4.44e-05
Wnr-CE-DL	1.30	2.12e-05	Wnr-CE-DR	1.39	4.28e-05
Capon	4.48	4.31e-05	Capon-DL	4.34	2.42e-05
ZF	2.97	2.44e-05	Kernel	1.12	1.94e-04
Kernel-DL	0.70	9.23e-05			

TABLE III
EXPERIMENTAL RESULTS (PILOT SIZE = 20)

Beamformer	MSE	Time	Beamformer	MSE	Time
Wnr	1.12	1.86e-04	Wnr-DL	1.05	1.87e-05
Wnr-DR	0.93	7.19e+00	Wnr-CE	1.12	5.78e-05
Wnr-CE-DL	1.08	3.14e-05	Wnr-CE-DR	1.13	6.01e-05
Capon	5.01	5.93e-05	Capon-DL	4.94	3.81e-05
ZF	3.82	3.56e-05	Kernel	1.20	4.48e-04
Kernel-DL	0.66	3.11e-04			

TABLE IV
EXPERIMENTAL RESULTS (PILOT SIZE = 25)

Beamformer	MSE	Time	Beamformer	MSE	Time
Wnr	0.92	1.41e-04	Wnr-DL	0.88	1.11e-05
Wnr-DR	0.80	4.22e+00	Wnr-CE	0.92	5.02e-05
Wnr-CE-DL	0.90	2.44e-05	Wnr-CE-DR	0.92	4.78e-05
Capon	4.94	4.93e-05	Capon-DL	4.89	2.85e-05
ZF	4.06	2.72e-05	Kernel	1.14	4.26e-04
Kernel-DL	0.60	2.95e-04			

TABLE V EXPERIMENTAL RESULTS (PILOT SIZE = 50)

Beamformer	MSE	Time	Beamformer	MSE	Time
Wnr	0.69	1.75e-04	Wnr-DL	0.68	1.85e-05
Wnr-DR	0.65	6.10e+00	Wnr-CE	0.69	5.81e-05
Wnr-CE-DL	0.68	3.03e-05	Wnr-CE-DR	0.70	5.90e-05
Capon	6.95	5.97e-05	Capon-DL	6.93	3.75e-05
ZF	6.36	3.38e-05	Kernel	0.92	1.81e-03
Kernel-DL	0.53	1.67e-03			

TABLE VI EXPERIMENTAL RESULTS (PILOT SIZE = 100)

Beamformer	MSE	Time	Beamformer	MSE	Time
Wnr	0.57	3.41e-04	Wnr-DL	0.57	3.64e-05
Wnr-DR	0.55	4.96e+00	Wnr-CE	0.57	6.35e-05
Wnr-CE-DL	0.57	2.93e-05	Wnr-CE-DR	0.58	6.07e-05
Capon	9.89	6.88e-05	Capon-DL	9.88	3.99e-05
ZF	9.45	3.27e-05	Kernel	0.72	5.93e-03
Kernel-DL	0.49	5.83e-03			

pilot size and several types of uncertainties in the signal model, the distributionally robust (DR) receive beamforming framework is then suggested. We prove that the diagonalloading (DL) methods are distributionally robust against the scarcity of pilot data and the uncertainties in the signal model. In addition, we generalize the diagonal-loading methods to achieve better estimation performance (e.g., the DR Wiener beamformer using F-norm for uncertainty quantification), at the cost of significantly higher computational burdens. Experiments suggest that nonlinear beamformers such as the Kernel and the Kernel-DL methods have the potential when the pilot size is small and/or the signal model is not linear Gaussian. Compared with the Kernel and the Kernel-DL beamformers, neural-network-based solutions [10], [12] have the stronger expressive capability of nonlinearities, which however are unscalable in the numbers of transmit and receive antennas, and significantly more time-consuming in training and more troublesome in tuning hyper-parameters (e.g., the number of layers and the number of neurons in each layer) than the studied eleven beamformers.

REFERENCES

- [1] T. Lo, H. Leung, and J. Litva, "Nonlinear beamforming," *Electronics Letters*, vol. 4, no. 27, pp. 350–352, 1991.
- [2] S. Yang and L. Hanzo, "Fifty years of MIMO detection: The road to large-scale MIMOs," *IEEE Commun. Surveys Tuts.*, vol. 17, no. 4, pp. 1941–1988, 2015.
- [3] A. M. Elbir, K. V. Mishra, S. A. Vorobyov, and R. W. Heath, "Twenty-five years of advances in beamforming: From convex and nonconvex optimization to learning techniques," *IEEE Signal Processing Mag.*, vol. 40, no. 4, pp. 118–131, 2023.
- [4] S. Chen, S. Tan, L. Xu, and L. Hanzo, "Adaptive minimum error-rate filtering design: A review," *Signal Processing*, vol. 88, no. 7, pp. 1671– 1697, 2008.
- [5] S. Chen, A. Wolfgang, C. J. Harris, and L. Hanzo, "Symmetric RBF classifier for nonlinear detection in multiple-antenna-aided systems," *IEEE Trans. Neural Networks*, vol. 19, no. 5, pp. 737–745, 2008.
- [6] A. Navia-Vazquez, M. Martinez-Ramon, L. E. Garcia-Munoz, and C. G. Christodoulou, "Approximate kernel orthogonalization for antenna array processing," *IEEE Trans. Antennas Propagat.*, vol. 58, no. 12, pp. 3942–3350, 2010.
- [7] M. Neinavaie, M. Derakhtian, and S. A. Vorobyov, "Lossless dimension reduction for integer least squares with application to sphere decoding," *IEEE Trans. Signal Processing*, vol. 68, pp. 6547–6561, 2020.
- [8] J. Liao, J. Zhao, F. Gao, and G. Y. Li, "Deep learning aided low complex breadth-first tree search for MIMO detection," *IEEE Trans. Wireless Commun.*, 2023.
- [9] D. A. Awan, R. L. Cavalcante, M. Yukawa, and S. Stanczak, "Robust online multiuser detection: A hybrid model-data driven approach," *IEEE Trans. Signal Processing*, 2023.
- [10] H. Ye, G. Y. Li, and B.-H. Juang, "Power of deep learning for channel estimation and signal detection in OFDM systems," *IEEE Wireless Commun. Lett.*, vol. 7, no. 1, pp. 114–117, 2017.
- [11] H. He, C.-K. Wen, S. Jin, and G. Y. Li, "Model-driven deep learning for MIMO detection," *IEEE Trans. Signal Processing*, vol. 68, pp. 1702– 1715, 2020.
- [12] N. Van Huynh and G. Y. Li, "Transfer learning for signal detection in wireless networks," *IEEE Wireless Commun. Lett.*, vol. 11, no. 11, pp. 2325–2329, 2022.
- [13] J. Li, P. Stoica, and Z. Wang, "On robust Capon beamforming and diagonal loading," *IEEE Trans. Signal Processing*, vol. 51, no. 7, pp. 1702–1715, 2003.
- [14] R. G. Lorenz and S. P. Boyd, "Robust minimum variance beamforming," IEEE Trans. Signal Processing, vol. 53, no. 5, pp. 1684–1696, 2005.
- [15] X. Zhang, Y. Li, N. Ge, and J. Lu, "Robust minimum variance beamforming under distributional uncertainty," in 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2015, pp. 2514–2518.

- [16] B. Li, Y. Rong, J. Sun, and K. L. Teo, "A distributionally robust minimum variance beamformer design," *IEEE Signal Processing Lett.*, vol. 25, no. 1, pp. 105–109, 2017.
- [17] Y. Huang, W. Yang, and S. A. Vorobyov, "Robust adaptive beamforming maximizing the worst-case SINR over distributional uncertainty sets for random inc matrix and signal steering vector," in ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2022, pp. 4918–4922.
- [18] Y. Huang, H. Fu, S. A. Vorobyov, and Z.-Q. Luo, "Robust adaptive beamforming via worst-case SINR maximization with nonconvex uncertainty sets," *IEEE Trans. Signal Processing*, vol. 71, pp. 218–232, 2023.
- [19] H. Cox, R. Zeskind, and M. Owen, "Robust adaptive beamforming," IEEE Trans. Acoust., Speech, Signal Processing, vol. 35, no. 10, pp. 1365–1376, 1987.
- [20] K. Harmanci, J. Tabrikian, and J. L. Krolik, "Relationships between adaptive minimum variance beamforming and optimal source localization," *IEEE Trans. Signal Processing*, vol. 48, no. 1, pp. 1–12, 2000.
- [21] F. Liu, L. Zhou, C. Masouros, A. Li, W. Luo, and A. Petropulu, "To-ward dual-functional radar-communication systems: Optimal waveform design," *IEEE Trans. Signal Processing*, vol. 66, no. 16, pp. 4264–4279, 2018.
- [22] J. A. Zhang, F. Liu, C. Masouros, R. W. Heath, Z. Feng, L. Zheng, and A. Petropulu, "An overview of signal processing techniques for joint communication and radar sensing," *IEEE J. Select. Topics Signal Processing*, vol. 15, no. 6, pp. 1295–1315, 2021.
- [23] Y. Xiong, F. Liu, Y. Cui, W. Yuan, T. X. Han, and G. Caire, "On the fundamental tradeoff of integrated sensing and communications under Gaussian channels," *IEEE Trans. Inform. Theory*, 2023.
- [24] K. P. Murphy, Machine Learning: A Probabilistic Perspective. MIT Press, 2012.
- [25] C. M. Bishop and N. M. Nasrabadi, Pattern Recognition and Machine Learning. Springer, 2006, vol. 4, no. 4.
- [26] G. Li and J. Ding, "Towards understanding variation-constrained deep neural networks," *IEEE Trans. Signal Processing*, vol. 71, pp. 631–640, 2023.
- [27] S. Shafieezadeh-Abadeh, D. Kuhn, and P. M. Esfahani, "Regularization via mass transportation," *Journal of Machine Learning Research*, vol. 20, no. 103, pp. 1–68, 2019.
- [28] M. Staib and S. Jegelka, "Distributionally robust optimization and generalization in kernel methods," Advances in Neural Information Processing Systems, vol. 32, 2019.
- [29] E. Delage and Y. Ye, "Distributionally robust optimization under moment uncertainty with application to data-driven problems," *Operations Research*, vol. 58, no. 3, pp. 595–612, 2010.
- [30] S. Wang, "Distributionally robust state estimation for jump linear systems," *IEEE Trans. Signal Processing*, 2023.
- [31] J. Li, S. Lin, J. Blanchet, and V. A. Nguyen, "Tikhonov regularization is optimal transport robust under martingale constraints," *Advances in Neural Information Processing Systems*, vol. 35, pp. 17677–17689, 2022.
- [32] D. Kuhn, P. M. Esfahani, V. A. Nguyen, and S. Shafieezadeh-Abadeh, "Wasserstein distributionally robust optimization: Theory and applications in machine learning," in *Operations Research & Management Science in the Age of Analytics*. Informs, 2019, pp. 130–166.
- [33] J. Blanchet, Y. Kang, and K. Murthy, "Robust Wasserstein profile inference and applications to machine learning," *Journal of Applied Probability*, vol. 56, no. 3, pp. 830–857, 2019.
- [34] C. Shorten and T. M. Khoshgoftaar, "A survey on image data augmentation for deep learning," *Journal of Big Data*, vol. 6, 2019.
- [35] G. Saon, Z. Tüske, K. Audhkhasi, and B. Kingsbury, "Sequence noise injected training for end-to-end speech recognition," in *ICASSP 2019-*2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2019, pp. 6261–6265.
- [36] K. Vu, J. C. Snyder, L. Li, M. Rupp, B. F. Chen, T. Khelif, K.-R. Müller, and K. Burke, "Understanding kernel ridge regression: Common behaviors from simple functions to density functionals," *International Journal of Quantum Chemistry*, vol. 115, no. 16, pp. 1115–1128, 2015.

Supplementary Materials

APPENDIX A STRUCTURED REPRESENTATION OF NONLINEAR FUNCTIONS

In Section II, we have reviewed two popular frameworks for representing (nonlinear) functions: reproducing kernel Hilbert spaces (RKHS) and neural network function spaces (NNFS). Typical kernel functions $\ker(\cdot,\cdot)$ to define RKHSs include Gaussian kernel, Matern kernel, Linear kernel, Laplacian kernel, and Polynomial kernel. Mathematical details of these kernel functions can be found in [24, Subsec. 14.2], [27, Ex. 1]. Typical activation functions $\sigma(\cdot)$ to define NNFSs include Hyperbolic tangent (i.e, tanh) function, Softmax function, Sigmoid function, Rectified linear unit (ReLU) function, and Exponential linear unit (ELU) function. Mathematical details of these activation functions can be found in [27, Ex. 2].

APPENDIX B

DETAILS ON REAL-SPACE SIGNAL REPRESENTATION

Let $C_x \coloneqq \mathbb{E}\mathbf{x}\mathbf{x}^\mathsf{T}$, $C_s \coloneqq \mathbb{E}\mathbf{s}\mathbf{s}^\mathsf{T}$, and $C_v \coloneqq \mathbb{E}\mathbf{v}\mathbf{v}^\mathsf{T} = \mathbf{0}$. We have

$$\begin{split} & \boldsymbol{R}_{\underline{x}} \coloneqq \mathbb{E}\underline{\mathbf{x}}^\mathsf{T} = \frac{1}{2} \left[\begin{array}{ccc} \operatorname{Re}(\boldsymbol{R}_x + \boldsymbol{C}_x) & \operatorname{Im}(-\boldsymbol{R}_x + \boldsymbol{C}_x) \\ \operatorname{Im}(\boldsymbol{R}_x + \boldsymbol{C}_x) & \operatorname{Re}(\boldsymbol{R}_x - \boldsymbol{C}_x) \end{array} \right]. \\ & \boldsymbol{R}_{\underline{s}} \coloneqq \mathbb{E}\underline{\mathbf{s}}^\mathsf{T} & = \frac{1}{2} \left[\begin{array}{ccc} \operatorname{Re}(\boldsymbol{R}_s + \boldsymbol{C}_s) & \operatorname{Im}(-\boldsymbol{R}_s + \boldsymbol{C}_s) \\ \operatorname{Im}(\boldsymbol{R}_s + \boldsymbol{C}_s) & \operatorname{Re}(\boldsymbol{R}_s - \boldsymbol{C}_s) \end{array} \right], \end{split}$$

and

$$\boldsymbol{R}_{\underline{v}} \coloneqq \mathbb{E} \underline{\mathbf{v}} \underline{\mathbf{v}}^{\mathsf{T}} = \frac{1}{2} \begin{bmatrix} \operatorname{Re} \boldsymbol{R}_{v} & \operatorname{Im} - \boldsymbol{R}_{v} \\ \operatorname{Im} \boldsymbol{R}_{v} & \operatorname{Re} \boldsymbol{R}_{v} \end{bmatrix}.$$

Note that the following identities hold: $R_x = HR_sH^H + R_v$, $C_x = HC_sH^T$, $R_{\underline{x}} = \underline{\underline{H}} \cdot R_{\underline{s}} \cdot \underline{\underline{H}}^T + R_v$, and $R_{\underline{xs}} = \underline{\underline{H}} \cdot R_{\underline{s}}$.

APPENDIX C

EXTENSIVE READING ON DISTRIBUTIONAL UNCERTAINTY

A. Generalization Error and Distributional Robustness

We use (7) and (15) as examples to illustrate the concepts. Supposing that ϕ^* solves the true problem (7) and ϕ_{ERM}^* solves the surrogate problem (15), we have

$$\begin{aligned} & \min_{\boldsymbol{\phi}} \operatorname{Tr} \mathbb{E}_{(\mathbf{x}, \mathbf{s}) \sim \mathbb{P}_{\mathbf{x}, \mathbf{s}}} [\boldsymbol{\phi}(\mathbf{x}) - \mathbf{s}] [\boldsymbol{\phi}(\mathbf{x}) - \mathbf{s}]^{\mathsf{H}} \\ & = \operatorname{Tr} \mathbb{E}_{(\mathbf{x}, \mathbf{s}) \sim \mathbb{P}_{\mathbf{x}, \mathbf{s}}} [\boldsymbol{\phi}^{\star}(\mathbf{x}) - \mathbf{s}] [\boldsymbol{\phi}^{\star}(\mathbf{x}) - \mathbf{s}]^{\mathsf{H}} \\ & \leq \operatorname{Tr} \mathbb{E}_{(\mathbf{x}, \mathbf{s}) \sim \mathbb{P}_{\mathbf{x}, \mathbf{s}}} [\boldsymbol{\phi}^{\star}_{\mathsf{ERM}}(\mathbf{x}) - \mathbf{s}] [\boldsymbol{\phi}^{\star}_{\mathsf{ERM}}(\mathbf{x}) - \mathbf{s}]^{\mathsf{H}}. \end{aligned}$$
(72)

To clarify further, the testing error in the last line (evaluated at the true distribution $\mathbb{P}_{\mathbf{x},\mathbf{s}}$) of the learned estimator $\phi_{\text{ERM}}^{\star}$ may be (much) larger than the optimal error in the first two lines, although $\phi_{\text{ERM}}^{\star}$ has the smallest training error (evaluated at the nominal distribution $\hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}$), i.e.,

$$\min_{\boldsymbol{\phi}} \operatorname{Tr} \mathbb{E}_{(\mathbf{x},\mathbf{s}) \sim \hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}} [\boldsymbol{\phi}(\mathbf{x}) - \mathbf{s}] [\boldsymbol{\phi}(\mathbf{x}) - \mathbf{s}]^{\mathsf{H}} \\
= \min_{\boldsymbol{\phi}} \operatorname{Tr} \frac{1}{L} \sum_{i=1}^{L} [\boldsymbol{\phi}(\boldsymbol{x}_{i}) - \boldsymbol{s}_{i}] [\boldsymbol{\phi}(\boldsymbol{x}_{i}) - \boldsymbol{s}_{i}]^{\mathsf{H}} \\
= \operatorname{Tr} \frac{1}{L} \sum_{i=1}^{L} [\boldsymbol{\phi}_{\mathsf{ERM}}^{\mathsf{ERM}}(\boldsymbol{x}_{i}) - \boldsymbol{s}_{i}] [\boldsymbol{\phi}_{\mathsf{ERM}}^{\mathsf{ERM}}(\boldsymbol{x}_{i}) - \boldsymbol{s}_{i}]^{\mathsf{H}} \\
\leq \operatorname{Tr} \frac{1}{L} \sum_{i=1}^{L} [\boldsymbol{\phi}^{\mathsf{*}}(\boldsymbol{x}_{i}) - \boldsymbol{s}_{i}] [\boldsymbol{\phi}^{\mathsf{*}}(\boldsymbol{x}_{i}) - \boldsymbol{s}_{i}]^{\mathsf{H}}.$$
(73)

In the terminologies of machine learning, the difference between the testing error and the training error, i.e.,

$$\begin{aligned} &\operatorname{Tr} \, \mathbb{E}_{(\mathbf{x},\mathbf{s}) \sim \mathbb{P}_{\mathbf{x},\mathbf{s}}} [\phi_{\text{ERM}}^{\star}(\mathbf{x}) - \mathbf{s}] [\phi_{\text{ERM}}^{\star}(\mathbf{x}) - \mathbf{s}]^{\mathsf{H}} - \\ &\operatorname{Tr} \, \mathbb{E}_{(\mathbf{x},\mathbf{s}) \sim \hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}} [\phi_{\text{ERM}}^{\star}(\mathbf{x}) - \mathbf{s}] [\phi_{\text{ERM}}^{\star}(\mathbf{x}) - \mathbf{s}]^{\mathsf{H}} \\ &= \operatorname{Tr} \, \mathbb{E}_{\mathbf{x},\mathbf{s}} [\phi_{\text{ERM}}^{\star}(\mathbf{x}) - \mathbf{s}] [\phi_{\text{ERM}}^{\star}(\mathbf{x}) - \mathbf{s}]^{\mathsf{H}} - \\ &\operatorname{Tr} \, \frac{1}{L} \sum_{i=1}^{L} [\phi_{\text{ERM}}^{\star}(x_i) - s_i] [\phi_{\text{ERM}}^{\star}(x_i) - s_i]^{\mathsf{H}} \end{aligned}$$

is called the *generalization error* of ϕ_{ERM}^{\star} ; the difference between the testing error and the optimal error, i.e.,

$$\begin{aligned} \operatorname{Tr} \mathbb{E}_{\mathbf{x},\mathbf{s}} [\phi_{\text{ERM}}^{\star}(\mathbf{x}) - \mathbf{s}] [\phi_{\text{ERM}}^{\star}(\mathbf{x}) - \mathbf{s}]^{\mathsf{H}} - \\ \operatorname{Tr} \mathbb{E}_{\mathbf{x},\mathbf{s}} [\phi^{\star}(\mathbf{x}) - \mathbf{s}] [\phi^{\star}(\mathbf{x}) - \mathbf{s}]^{\mathsf{H}} \end{aligned}$$

is called the *excess risk* of ϕ_{ERM}^{\star} . In machine learning practice, we want to reduce both the generalization error and the excess risk. Most attention in the literature has been particularly paid to reducing generalization errors. Specifically, an upper bound of the true cost $\operatorname{Tr} \mathbb{E}_{(\mathbf{x},\mathbf{s})\sim \mathbb{P}_{\mathbf{x},\mathbf{s}}}[\phi(\mathbf{x})-\mathbf{s}][\phi(\mathbf{x})-\mathbf{s}]^H$ is first found and then minimize the upper bound: by minimizing the upper bound, the true cost can also be reduced.

Fact 1: Suppose that the true distribution $\mathbb{P}_{0,\mathbf{x},\mathbf{s}}$ of (\mathbf{x},\mathbf{s}) is included in $\mathcal{U}_{\mathbf{x},\mathbf{s}}$; for notational clarity, we hereafter distinguish $\mathbb{P}_{0,\mathbf{x},\mathbf{s}}$ from $\mathbb{P}_{\mathbf{x},\mathbf{s}}$. The true objective function evaluated at $\mathbb{P}_{0,\mathbf{x},\mathbf{s}}$, i.e.,

$$\operatorname{Tr} \mathbb{E}_{(\mathbf{x},\mathbf{s}) \sim \mathbb{P}_{0,\mathbf{x},\mathbf{s}}} [\phi(\mathbf{x}) - \mathbf{s}] [\phi(\mathbf{x}) - \mathbf{s}]^{\mathsf{H}}, \quad \forall \phi \in \mathcal{B},$$
 (74)

is upper bounded by the worst-case objective function of (17), i.e.,

$$\max_{\mathbb{P}_{\mathbf{x},\mathbf{s}} \in \mathcal{U}_{\mathbf{x},\mathbf{s}}} \operatorname{Tr} \mathbb{E}_{(\mathbf{x},\mathbf{s}) \sim \mathbb{P}_{\mathbf{x},\mathbf{s}}} [\phi(\mathbf{x}) - \mathbf{s}] [\phi(\mathbf{x}) - \mathbf{s}]^{\mathsf{H}}, \quad \forall \phi \in \mathcal{B}. \tag{75}$$

Therefore, by diminishing the upper bound in (75), the true estimation error evaluated at $\mathbb{P}_{0,\mathbf{x},\mathbf{s}}$ can also be reduced. However, the conventional empirical estimation error evaluated at $\hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}$ cannot upper bound the true estimation error (74). This performance guarantee is the benefit of considering the distributionally robust method (17). Due to the weak convergence property of the empirical distribution to the true data-generating distribution, that is, $d(\mathbb{P}_{0,\mathbf{x},\mathbf{s}}, \hat{\mathbb{P}}_{\mathbf{x},\mathbf{s}}) \to 0$ as the sample size $L \to \infty$, there exists ϵ in (18) for every L, such that $\mathbb{P}_{0,\mathbf{x},\mathbf{s}}$ is included in $\mathcal{U}_{\mathbf{x},\mathbf{s}}$ in $\mathbb{P}^L_{0,\mathbf{x},\mathbf{s}}$ -probability (L-fold product measure of $\mathbb{P}_{0,\mathbf{x},\mathbf{s}}$).

B. Non-Stationary Channel Statistics

In the main body of the paper (see also Fact 1), we assume that the true data-generating distribution $\mathbb{P}_{0,\mathbf{x},\mathbf{s}}$ is time-invariant within a frame. In real-world operations, however, this assumption might be untenable.

As shown in Fig. 1, the frame contains eight data units; we suppose that the first four units are pilot symbols and the rest four units are communication-data symbols.

Let $\mathbb{P}_{0,\mathbf{x},\mathbf{s},i}$ denote the true data-generating distribution at time point t_i where $i=0,1,2,\ldots,8$. Specifically, we have $(\mathbf{x}_i,\mathbf{s}_i) \sim \mathbb{P}_{0,\mathbf{x},\mathbf{s},i}$ for every i. Therefore, the pilot data set (i.e.,

1

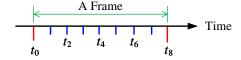


Fig. 1. True data-generating distributions might be time-varying in a frame.

the training data set) $\{(x_1,s_1),(x_2,s_2),(x_3,s_3),(x_4,s_4)\}$ can be seen as realizations of the mean distribution $\mathbb{P}_{\text{train},0,\mathbf{x},\mathbf{s}}$ of underlying true training-data distributions where $\mathbb{P}_{\text{train},0,\mathbf{x},\mathbf{s}} = \sum_{i=1}^4 h_i \mathbb{P}_{0,\mathbf{x},\mathbf{s},i}$, which is a mixture distribution with mixing weights $0 \leq h_1,h_2,h_3,h_4 \leq 1; \sum_{i=1}^4 h_i = 1$. Similarly, the communication data set (i.e., the testing data set) $\{(x_5,s_5),(x_6,s_6),(x_7,s_7),(x_8,s_8)\}$ can be seen as realizations of the mean $\mathbb{P}_{\text{test},0,\mathbf{x},\mathbf{s}}$ of the underlying true testing-data distributions where $\mathbb{P}_{\text{test},0,\mathbf{x},\mathbf{s}} = \sum_{i=5}^8 h_i \mathbb{P}_{0,\mathbf{x},\mathbf{s},i}$, with mixing weights $0 \leq h_5,h_6,h_7,h_8 \leq 1; \sum_{i=5}^8 h_i = 1$.

Suppose that

$$d(\hat{\mathbb{P}}_{\text{train},\mathbf{x},\mathbf{s}}, \mathbb{P}_{\text{train},0,\mathbf{x},\mathbf{s}}) \leq \epsilon_1,$$

where $\hat{\mathbb{P}}_{\text{train},\mathbf{x},\mathbf{s}} \coloneqq \frac{1}{4} \sum_{i=1}^{4} \delta_{(\boldsymbol{x}_i,\boldsymbol{s}_i)}$ is the data-driven estimate of $\mathbb{P}_{\text{train},0,\mathbf{x},\mathbf{s}}$ and

$$d(\mathbb{P}_{\text{train},0,\mathbf{x},\mathbf{s}}, \mathbb{P}_{\text{test},0,\mathbf{x},\mathbf{s}}) \leq \epsilon_2,$$

for some $\epsilon_1, \epsilon_2 \geq 0$. We have the uncertainty quantification

$$d(\mathbb{P}_{\text{test},0,\mathbf{x},\mathbf{s}}, \ \hat{\mathbb{P}}_{\text{train},\mathbf{x},\mathbf{s}}) \leq \epsilon \coloneqq \epsilon_1 + \epsilon_2.$$

Therefore, the distributionally robust modeling and solution framework is still valid to hedge against the distributional uncertainty in the nominal distribution $\hat{\mathbb{P}}_{\text{train},\mathbf{x},\mathbf{s}}$ compared to the underlying true distribution $\mathbb{P}_{\text{test},0,\mathbf{x},\mathbf{s}}$. When $\mathbb{P}_{\text{train},0,\mathbf{x},\mathbf{s}} = \mathbb{P}_{\text{test},0,\mathbf{x},\mathbf{s}}$, as assumed in the main body of the paper, we have $\epsilon_1 \to 0$ and $\epsilon \to \epsilon_2 = 0$ as the pilot size tends to infinity; however, when $\mathbb{P}_{\text{train},0,\mathbf{x},\mathbf{s}} \neq \mathbb{P}_{\text{test},0,\mathbf{x},\mathbf{s}}$, the radius $\epsilon \to \epsilon_2 \neq 0$ although $\epsilon_1 \to 0$.

Another justification for the DRO method is as follows. Suppose that there exists $\epsilon \geq 0$ such that

$$d(\mathbb{P}_{0,\mathbf{x},\mathbf{s},i}, \ \hat{\mathbb{P}}_{\text{train},\mathbf{x},\mathbf{s}}) \le \epsilon, \quad \forall i \in \{0,1,2,\dots,8\}.$$

It means that, at every snapshot in the frame, the true datagenerating distribution is included in the uncertainty set. Hence, the DRO cost can still upper bound the true cost even though the true distribution is time-varying; cf. Fact 1.

APPENDIX D PROOF OF LEMMA 1

Proof: The objective function of Problem (19) equals to

$$\left\langle \begin{bmatrix} \mathbf{W}^{\mathsf{H}}\mathbf{W} & -\mathbf{W}^{\mathsf{H}} \\ -\mathbf{W} & \mathbf{I}_{M} \end{bmatrix}, \begin{bmatrix} \mathbf{R}_{x} & \mathbf{R}_{xs} \\ \mathbf{R}_{xs}^{\mathsf{H}} & \mathbf{R}_{s} \end{bmatrix} \right\rangle, \quad (76)$$

where $\langle A, B \rangle := \operatorname{Tr} A^H B$ for two matrices A and B. Therefore, the objective function of (19) is convex in W and linear (thus concave) in the matrix variable R. Note also that complex Euclidean spaces are linear topological

spaces. Hence, due to Sion's minimax theorem, Problem (19) is equivalent to

$$egin{aligned} \max \min_{m{R}} & \operatorname{Tr}\left[m{W}m{R}_xm{W}^\mathsf{H} - m{W}m{R}_{xs} - m{R}_{xs}^\mathsf{H}m{W}^\mathsf{H} + m{R}_s
ight] \ & ext{s.t.} & d_0(m{R},\ \hat{m{R}}) \leq \epsilon_0, \ & m{R} \succeq \mathbf{0}. \end{aligned}$$

For every given R, the inner minimization sub-problem of (77) is solved by the Wiener beamformer $W_{\text{Wiener}}^{\star} = R_{xs}^{\text{H}} R_x^{-1}$, which transforms (77) to (21). This completes the proof. \square

APPENDIX E PROOF OF THEOREM 1

Proof: Consider the following optimization problem

$$\max_{\boldsymbol{R}} \quad \text{Tr} \left[-\boldsymbol{R}_{xs}^{\mathsf{H}} \boldsymbol{R}_{x}^{-1} \boldsymbol{R}_{xs} + \boldsymbol{R}_{s} \right]$$
s.t. $\boldsymbol{R} \succeq \boldsymbol{R}_{2},$ (78)
$$\boldsymbol{R}_{x} \succ \boldsymbol{0},$$

which, due to Lemma 1, is equivalent [in the sense of the same optimal objective value and maximizer(s) R^*] to

$$\begin{array}{ccc} \min \max \limits_{\substack{\boldsymbol{W} & \boldsymbol{R} \\ \boldsymbol{W} & \boldsymbol{R}}} & \left\langle \left[\begin{array}{cc} \boldsymbol{W}^{\mathsf{H}} \boldsymbol{W} & -\boldsymbol{W}^{\mathsf{H}} \\ -\boldsymbol{W} & \boldsymbol{I}_{M} \end{array} \right], \; \left[\begin{array}{ccc} \boldsymbol{R}_{x} & \boldsymbol{R}_{xs} \\ \boldsymbol{R}_{xs}^{\mathsf{H}} & \boldsymbol{R}_{s} \end{array} \right] \right\rangle \\ \text{s.t.} & \boldsymbol{R} \succeq \boldsymbol{R}_{2}, \\ \boldsymbol{R}_{x} \succ \boldsymbol{0}. \end{array}$$

Note that $egin{bmatrix} m{W}^{\mathsf{H}}m{W} & -m{W}^{\mathsf{H}} \\ -m{W} & m{I}_M \end{bmatrix} \succeq \mathbf{0}$, because for all $m{x} \in \mathbb{C}^N$ and $m{y} \in \mathbb{C}^M$, we have

$$[oldsymbol{x}^{\mathsf{H}}, \ oldsymbol{y}^{\mathsf{H}}] \left[egin{array}{cc} oldsymbol{W}^{\mathsf{H}} oldsymbol{W} & -oldsymbol{W}^{\mathsf{H}} \ -oldsymbol{W} & oldsymbol{I}_M \end{array}
ight] \left[egin{array}{cc} oldsymbol{x} \ oldsymbol{y} \end{array}
ight] = \|oldsymbol{W} oldsymbol{x} - oldsymbol{y}\|_F^2 \geq 0.$$

Therefore, for every given W, the objective function of (79) is increasing in R. As a result, the objective value of (78) is lower-bounded at R_2 : To be specific, $\forall R \succeq R_2$, we have

$$\text{Tr}\left[-\boldsymbol{R}_{xs}^{\mathsf{H}}\boldsymbol{R}_{x}^{-1}\boldsymbol{R}_{xs} + \boldsymbol{R}_{s}\right] \geq \text{Tr}\left[-\boldsymbol{R}_{2,xs}^{\mathsf{H}}\boldsymbol{R}_{2,x}^{-1}\boldsymbol{R}_{2,xs} + \boldsymbol{R}_{2,s}\right],$$
 i.e., $f_{1}(\boldsymbol{R}) \geq f_{1}(\boldsymbol{R}_{2})$, which proves the first part.

On the other hand, if $\mathbf{R}_{1,x} \succeq \mathbf{R}_{2,x} \succ \mathbf{0}$, we have $\mathbf{R}_{2,x}^{-1} \succeq \mathbf{R}_{1,x}^{-1}$. As a result, $f_2(\mathbf{R}_{1,x}) - f_2(\mathbf{R}_{2,x}) = \operatorname{Tr}\left[\mathbf{R}_{xs}^{\mathsf{H}}(\mathbf{R}_{2,x}^{-1} - \mathbf{R}_{1,x}^{-1})\mathbf{R}_{xs}\right] \geq 0$, completing the proof. \square

APPENDIX F PROOF OF PROPOSITION 3

Proof: Letting $\mathbf{z} := \varphi(\mathbf{x})$, (57) can be rewritten as

$$\min_{\boldsymbol{W} \in \mathbb{R}^{2M \times 2N}} \max_{\mathbb{P}_{\underline{\mathbf{z}},\underline{\mathbf{s}}} \in \mathcal{U}_{\underline{\mathbf{z}},\underline{\mathbf{s}}}} \operatorname{Tr} \mathbb{E}_{\underline{\mathbf{z}},\underline{\mathbf{s}}} [\boldsymbol{W}\underline{\mathbf{z}} - \underline{\mathbf{s}}] [\boldsymbol{W}\underline{\mathbf{z}} - \underline{\mathbf{s}}]^{\mathsf{T}}.$$
 (80)

Tantamount to the distributionally robust beamforming problem (19), Problem (80) reduces to (58) where

$$\hat{m{R}}_{m{\underline{z}}} \coloneqq rac{1}{L} \sum_{i=1}^L m{\underline{z}}_i m{\underline{z}}_i^\mathsf{T} = rac{1}{L} \sum_{i=1}^L m{arphi}(m{\underline{x}}_i) m{arphi}^\mathsf{T}(m{\underline{x}}_i) = rac{1}{L} m{K}^2,$$

$$\hat{m{R}}_{m{z}m{s}} \coloneqq rac{1}{L} \sum_{i=1}^{L} m{z}_i m{\underline{s}}_i^{\mathsf{T}} = rac{1}{L} \sum_{i=1}^{L} m{arphi}(m{x}_i) \cdot m{\underline{s}}_i^{\mathsf{T}} = rac{1}{L} m{K} m{\underline{S}}^{\mathsf{T}},$$

$$\hat{m{R}}_{\underline{s}}\coloneqq rac{1}{L}\sum_{i=1}^{L}m{\underline{s}}_{i}m{\underline{s}}_{i}^{\mathsf{T}} = rac{1}{L}\sum_{i=1}^{L}m{\underline{s}}_{i}\cdotm{\underline{s}}_{i}^{\mathsf{T}} = rac{1}{L}m{\underline{S}}m{\underline{S}}^{\mathsf{T}},$$

and

$$K \coloneqq [\boldsymbol{\varphi}(\underline{\boldsymbol{x}}_1), \boldsymbol{\varphi}(\underline{\boldsymbol{x}}_2), \dots, \boldsymbol{\varphi}(\underline{\boldsymbol{x}}_L)] \in \mathbb{R}^{L \times L}.$$

The rest claims are due to Lemma 1; NB: K is invertible. \square

APPENDIX G ADDITIONAL EXPERIMENTAL RESULTS

Complementary to experimental setups in Section VI, we consider pure complex Gaussian channel noises. First, we suppose that the transmit antennas emit continuous-valued complex signals; without loss of generality, Gaussian signals are used in experiments. The performance evaluation measure is therefore the mean-squared error (MSE). The experimental results are shown in Fig. 2.

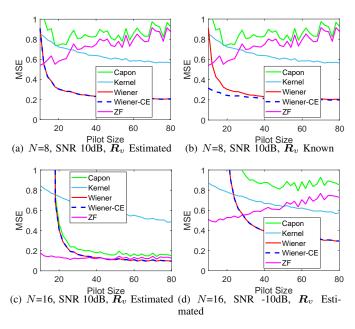


Fig. 2. Testing MSE against training pilot sizes under different numbers of receive antennas; only non-robust beamformers including non-diagonal-loading ones are considered. The true value of \mathbf{R}_v can be unknown and estimated using pilot data. The signal-to-noise ratio (SNR) is 10dB or -10dB.

From Fig. 2, the following main points can be outlined.

- For a fixed number M of transmit antennas, the larger the number N of receive antennas, the smaller the MSE; cf. Figs. 2(a) and 2(c). This fact is well-established and is due to the benefit of antenna diversity. In addition, for fixed N and M, the higher the SNR, the smaller the MSE; cf. Figs. 2(c) and 2(d); this is also well believed.
- 2) As the pilot size increases, the Wiener beamformer tends to have the best performance because the Wiener beamformer is optimal for the linear Gaussian signal model. When \mathbf{R}_v is accurately known, the Wiener-CE beamformer outperforms the general Wiener beamformer (cf. Fig. 2(b)) because the former also exploits the information of the linear signal model in addition to the pilot data, while the latter only utilizes the pilot data. However, when \mathbf{R}_v is estimated using the pilot data, the performances of the general Wiener beamformer and the Wiener-CE beamformer have no significant difference; cf. Figs. 2(a) and 2(c). Therefore, Fig. 2 validates our claim

- that channel estimation is not a necessary operation in receive beamforming and estimation of wireless signals; recall Subsection III-A3.
- 3) The ZF beamformer tends to be more efficient as N increases; cf. Figs. 2(a) and 2(c). However, the ZF beamformer becomes less satisfactory when the SNR decreases; cf. Figs. 2(c) and 2(d). The Capon beamformer is also unsatisfactory when N is small or the SNR is low.
- 4) The kernel beamformer, as a nonlinear method, cannot outperform linear beamformers because, for a linear Gaussian signal model, the optimal beamformer is linear. From the perspective of machine learning, nonlinear methods tend to overfit the limited training samples.

Second, we suppose that the transmit antennas emit discrete-valued symbols from a constellation that is modulated using quadrature phase-shift keying (QPSK). The performance evaluation measure is therefore the symbol error rate (SER). The experimental results are shown in Fig. 3. We find that all the conclusive main points from Fig. 2 can be obtained from Fig. 3 as well: this validates that *minimizing MSE reduces SER*. In addition, Figs. 3(c) and 3(d) reveal that the Wiener beamformer even slightly works better than the Wiener-CE beamformer when the pilot size is smaller than 15 because the uncertainty in the estimated \hat{R}_v , on the contrary, misleads the latter. Nevertheless, as the pilot size increases, the Wiener-CE beamformer tends to overlap the Wiener beamformer quickly.

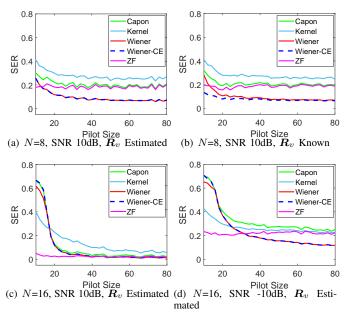


Fig. 3. Testing SER against training pilot sizes under different numbers of receive antennas; only non-robust beamformers including non-diagonal-loading ones are considered. The true value of \mathbf{R}_v can be unknown and estimated using pilot data. The signal-to-noise ratio (SNR) is 10dB or -10dB.