

On the Design of Channel Shortening Demodulators for Iterative Receivers in Linear Vector Channels

SHA HU¹, (Member, IEEE), AND FREDRIK RUSEK¹

Department of Electrical and Information Technology, Lund University, SE-22100 Lund, Sweden

Corresponding author: Sha Hu (sha.hu@eit.lth.se)

ABSTRACT We consider designing demodulators for linear vector channels that use reduced-size trellis descriptions for the received signal. We assume an iterative receiver and use interference cancellation (IC) with soft-information provided by an outer decoder to mitigate the signal part that is not covered by a reduced-size trellis description. In order to reach a trellis description, a linear filter is applied as a front end to compress the signal structure into a small trellis. This process requires three parameters to be designed: 1) the front-end filter; 2) the feedback filter through which the IC is done; and 3) a target response which specifies the trellis. Demodulators of this form have been studied before under the name *channel shortening* (CS), but the interplay between CS, IC, and the trellis-search processes has not been adequately addressed in the literature. In this paper, we analyze two types of CS demodulators that are based on the Forney and Ungerboeck detection models, respectively. The parameters are jointly optimized with a generalized mutual information (GMI) function. We also introduce a third type of CS demodulator that is, in general, suboptimal, which has closed-form solutions. Furthermore, signal-to-noise ratio asymptotic properties are analyzed, and we show that the third CS demodulator asymptotically converges to the optimal CS demodulator in the sense of GMI maximization.

INDEX TERMS Channel shortening (CS), intersymbol interference (ISI), multi-input multi-output (MIMO), front-end filter, feedback filter, target response, generalized mutual information (GMI), Forney model, Ungerboeck model, turbo equalization, demodulator, linear minimum mean square error (LMMSE), interference cancellation (IC), extrinsic information transfer (EXIT), block-error-rate (BLER), BCJR.

I. INTRODUCTION

Channel shortening (CS) demodulators have a long and rich history, see [3]–[16]. For intersymbol interference (ISI) channels, Forney, Jr., [17] showed that Viterbi Algorithm (VA) [19] implements maximum likelihood (ML) detection. However, the complexity of VA is exponential in the memory of the channel which prohibits its use in many cases of interest. As a remedy, Falconer and Magee [4] proposed in 1973 the concept of CS. The concept is to filter the received signal with a prefilter such that the effective channel has a much shorter duration than the original one, and then apply VA to the shortened channel.

Traditionally, CS demodulators have been optimized from various criteria such as minimum mean square error (MMSE) and signal-to-interference-plus-noise ratio (SINR) [5]–[13]. Venkataramani and Erden [14] attempted to minimize the error probability of an uncoded system, which leads to a new

notion of posterior equivalence between the target response and the filtered channel. But as [14] works with uncoded error probabilities, the analysis does not adequately address the case of coded systems and Shannon capacity properties. The first paper that works with capacity-related cost measures is [15] where the authors considered the achievable rate, in the form of generalized mutual information (GMI) [18], [20]–[23], that the transceiver system can achieve if a CS demodulator is adopted. However, [15] is limited to ISI channels only, and the design method in [15] of the CS demodulator is in fact not always possible to execute. The limitations of [15] were first dealt with in [18], which extended the CS concept to any linear vector channel such as multi-input multi-output (MIMO) and ISI channels, and resulted in a closed-form optimization procedure.

On the other hand, iterative receivers such as turbo equalization [24]–[28] followed as a natural extension to

turbo codes for designs of iterative detection and decoding receivers. When it comes to turbo equalization, common settings of the demodulator are [24] the maximum *a posteriori* (MAP) demodulator [43] and its suboptimal variants such as dimension-reduction and subspace based detections [29], [30], and linear MMSE (LMMSE) [31], [32]. The suboptimal demodulators replace the MAP with a linear equalizer or a decision feedback equalizer (DFE) to reduce the prohibitive complexity. One open problem in the area of turbo equalization is the development of other non-trellis-based detection methods that provide performance between MAP and LMMSE [24], [32].

Instead of fully removing the trellis-based detection, another approach is to reduce the memory-size of the original linear vector channel through an interference cancellation (IC) based prefiltering. To the best of our knowledge, there is limited literature [26], [33] on such a demodulator design that combines both IC based prefiltering and a memory-size shortened BCJR in an iterative receiver. A closely related concept is delayed-decision-feedback-sequence-estimation (DDFSE) [34], [35], which also reduces the number of states in the BCJR. However, in DDFSE the IC is done within a single iteration, and not between the iterations.

In this paper, we generalize the GMI-maximization based CS demodulator in [18] to cooperate with iterative receivers. With iterative receivers, it is reasonable to expect that better detection-performance can be attained by allowing the parameters of CS demodulator to change over each iteration. However, the CS demodulator in [18] does not take the prior information into account, rendering a static design in all iterations. We aim at constructing a CS demodulator that takes soft-information provided by the outer-decoder into account such that the parameters of CS demodulator are designed for a particular level of prior knowledge. This procedure includes an IC mechanism to deal with the signal part that are not handled by the trellis-search, i.e., the BCJR. Preliminary results for CS demodulators in iterative receivers are available in [2] and [36], but this paper non-trivially advances the state-of-the-art.

Although the trellis-search based detection is still utilized in CS demodulator, the memory-size ν of the linear vector channel has been reduced which results in significant complexity-reduction compared to MAP. Meanwhile, with different values of ν , CS demodulator provides trade-offs between the performance of LMMSE and MAP. As what will become clear later that CS demodulator is closely related to the concept of LMMSE with parallel interference cancellation (LMMSE-PIC) [37]–[39], which cooperates the soft-information into its filter-coefficients and IC process. By setting $\nu = 0$ in CS demodulator, it is identical to LMMSE-PIC whose trellis-search process is trivial as different layers are assumed to be independent after the front-end filtering. However, CS demodulator advances the LMMSE-PIC by including a trellis-search, where the parameters of the front-end filter, IC, and the trellis-search are jointly optimized. On the other hand, by setting ν to the original

memory-size of the linear vector channel, CS demodulator is identical to MAP. Therefore, CS demodulator provides a generalized framework that includes LMMSE-PIC and MAP as two extreme cases for designing an iterative receiver.

The rest of the paper is organized as follows: The linear vector channel model and the iterative receiver structure are introduced in Section II, while the general form of CS demodulator and the GMI are described in Section III. In Section IV we analyze three types of CS demodulators for MIMO channels. In Section V we deal with ISI channels as asymptotic versions of the results established in Section IV. The signal-to-noise ratio (SNR) asymptotic of the CS demodulators are discussed in Section VI. Numerical results are then shown in Section VII, and Section VIII summarizes the paper.

NOTATION

Throughout the paper, a capital bold letter such as \mathbf{A} represents a matrix, a lower case bold letter \mathbf{a} represents a vector, and a capital letter A represents a number. The expression $\mathbf{A} < \mathbf{0}$ means matrix \mathbf{A} is negative-definite, while $\mathbf{A} > \mathbf{0}$ means \mathbf{A} is positive-definite. We let \mathbf{I}_K represent a $K \times K$ identity matrix and the dimension is omitted when it can be understood from the context. The superscripts have the following meanings: $(\cdot)^*$ is complex conjugate, $(\cdot)^T$ is matrix transpose, $(\cdot)^\dagger$ denotes the conjugate transpose of a matrix, $(\cdot)^{-1}$ is matrix inverse. In addition, \propto means proportional to, $\mathbb{E}[\cdot]$ is the expectation operator, $\text{Tr}(\cdot)$ takes the trace of a matrix, $\mathcal{R}\{\cdot\}$ returns the real part of a variable, \otimes is the Kronecker multiplication operator, $\text{vec}(\mathbf{A})$ is a column vector containing the columns of matrix \mathbf{A} stacked on top of each other, and $[A, B]$ is the set of integers $\{k : A \leq k \leq B\}$.

Further, we say that a matrix \mathbf{A} is banded within diagonals $[-\nu_1, \nu_2]$ ($\nu_1, \nu_2 \geq 0$), if the (k, ℓ) th element $A(k, \ell)$ satisfies¹

$$A(k, \ell) = 0, \ell - k > \nu_1 \text{ or } k - \ell > \nu_2.$$

Moreover, we define two matrix notations $[\]_\nu$ and $[\]_{\setminus \nu}$ such that $\mathbf{A} = [\mathbf{A}]_\nu + [\mathbf{A}]_{\setminus \nu}$, and $[\mathbf{A}]_\nu$ is banded within diagonals $[-\nu, \nu]$ where $[\mathbf{A}]_{\setminus \nu}$ are constrained to zeros.

II. SYSTEM MODEL

We consider linear vector channels according to the received signal model

$$\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{n} \quad (1)$$

where \mathbf{y} is an $N \times 1$ vector of the received signal; \mathbf{x} is a $K \times 1$ vector comprising unit-energy coded symbols that belong to a constellation \mathcal{X} ; \mathbf{H} is an $N \times K$ matrix representing the communication channel which is perfectly known to the receiver; and \mathbf{n} is zero-mean complex-valued Gaussian noise vector with a covariance matrix $N_0\mathbf{I}$.

The model (1) may represent many different communication systems and we consider two typical cases, i.e., when

¹Note that ν_1 refers to the number of upper diagonals of \mathbf{A} that are non-zero. We have this convention in order to subsequently follow standard notation for Toeplitz matrices [40].

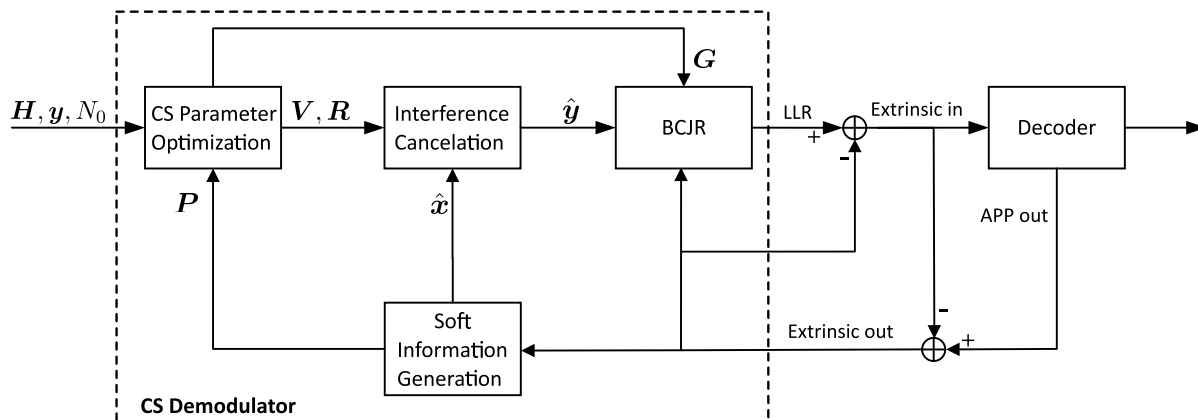


FIGURE 1. Iterative receiver structure with CS demodulator and outer-decoder. The target of the CS demodulator is to maximize the GMI through jointly optimizing the parameters V, R , and G , which are referred to as the front-end filter, IC matrix, and trellis-representation matrix, respectively.

H represents a multi-input multi-output (MIMO) or an ISI channel. In the MIMO case, the variables N and K are finite while they grow without bounds in the ISI case. For the MIMO case, a block-fading model is assumed to perform an analysis for a whole transmitted data-block.

In an iterative receiver, the feedbacks from the outer-decoder can be utilized to improve the performance. As the outer-decoder provides the demodulator with *a posteriori* probability (APP) and extrinsic information (in terms of bit log-likelihood ratio (LLR)) [41], [42], side-information is present about the symbols x and we represent this by the probability mass function $p_k(s) = P(x_k = s)$, ($0 \leq k \leq K - 1$). Note that the side-information does not consider the dependency among the symbols, but are symbol-wise marginal probabilities. This reflects the situation encountered in iterative receivers with perfect interleaving. In those cases, the prior probabilities provided from previous iterations are assumed independent, i.e., $P(x = s) = \prod p_k(s)$, and the demodulator can compute $\hat{x} = \mathbb{E}[x] = [\hat{x}_0, \hat{x}_1, \dots, \hat{x}_{K-1}]^T$ in a per-entry fashion as

$$\hat{x}_k = \sum_{s \in \mathcal{X}} s p_k(s),$$

where the expectations are computed with respect to the prior distribution $p_k(s)$.

With soft-information \hat{x} , we define a $K \times K$ diagonal matrix P reflecting the accuracy of the side-information as

$$P = \mathbb{E}_y[x \hat{x}^\dagger] = \mathbb{E}_y[\hat{x} \hat{x}^\dagger], \tag{2}$$

where x is the transmitted symbol for the received signal y , and the exception “ \mathbb{E}_y ” is taken by averaging $\hat{x} \hat{x}^\dagger$ through the whole transmitted data-block. Note that $0 \leq P \leq I$, and when there is no soft-information available we have $P = 0$, while with perfect feedback we get $P = I$.

For the ISI case, as there is only a single transmit-layer and K is the length of data-block, P can be simplified as

$$P = \alpha I, \tag{3}$$

where

$$\alpha = \frac{1}{K} \sum_{k=0}^{K-1} |\hat{x}_k|^2.$$

The task of the demodulator is to generate soft-information about x given the observable y and the side-information. The optimal demodulator is the MAP [43], [44] which evaluates the posterior probabilities $P(x_k = s|y)$. However, the number of leaves of the search-tree in MAP is in general $|\mathcal{X}|^K$ which is prohibitive for practical applications. With CS demodulator we force the signal model to be an lower-triangular matrix with only $\nu + 1$ ($0 \leq \nu < K - 1$) non-zero diagonals, by means of a linear filter,² where ν is referred to as the memory-size of the CS demodulator. Then, the BCJR [45] can be applied over a trellis with $|\mathcal{X}|^\nu$ states. Further, as there is side-information present about x , the parts of H that are outside the memory of BCJR can be partly eliminated by means of IC through the prior mean \hat{x} .

The structure of an iterative receiver utilizing a CS demodulator is depicted in Fig. 1. The extrinsic information from an outer-decoder is used to compute an estimate \hat{x} and a matrix P that reflects the feedback quality. Based on the updated P in each iteration, the optimal CS parameters are solved by maximizing the GMI. A prefiltering and IC process are then implemented on y with optimal V and R to obtain the signal \hat{y} , which is sent to a memory- ν BCJR module specified by an optimal G . Further, the extrinsic information iteratively exchanged between the BCJR and the outer-decoder is also used as *a priori* information for decoding the transmitted symbols. Note that if we set $\nu = K - 1$, the search-space of the CS demodulator is no longer a trellis but corresponds to the original tree and is equivalent to the MAP, while the LMMSE-PIC is a special case of the CS demodulation with $\nu = 0$.

²For finite length linear vector channels such as MIMO channel, “filtering” means matrix multiplication.

III. THE GENERAL FORM OF THE CS DEMODULATOR

We state two lemmas that are useful later, and Lemma 2 can be verified straightforwardly.

Lemma 1: Let \mathbf{A}_1 and \mathbf{A}_2 be two $K \times K$ matrices, where \mathbf{A}_1 is invertible and banded within diagonals $[-v, v]$. If $[\mathbf{A}_1^{-1}]_v = [\mathbf{A}_2]_v$, then

$$\text{Tr}(\mathbf{A}_1 \mathbf{A}_2) = \text{Tr}(\mathbf{I}).$$

Proof: Let $\mathbf{A}_3 = \mathbf{A}_2 - \mathbf{A}_1^{-1}$, then it holds that $[\mathbf{A}_3]_v = \mathbf{0}$ and $\mathbf{A}_3 = [\mathbf{A}_3]_{\setminus v}$. As $\mathbf{A}_1 = [\mathbf{A}_1]_v$, the elements along the main diagonal of $\mathbf{A}_1 \mathbf{A}_3$ are zeros. Hence, we have $\text{Tr}(\mathbf{A}_1 \mathbf{A}_2) = \text{Tr}(\mathbf{A}_1 (\mathbf{A}_1^{-1} + \mathbf{A}_3)) = \text{Tr}(\mathbf{I})$. ■

Lemma 2: Let \mathbf{A}_1 and \mathbf{A}_2 be two $K \times K$ matrices that are banded within diagonals $[-v_1, v_2]$ and $[-v_3, v_4]$, respectively. Then the product $\mathbf{A}_1 \mathbf{A}_2$ is banded within diagonals $[\max(-v_1 + v_3, 1 - K), \min(v_2 + v_4, K - 1)]$.

A. SYSTEM MODEL OF THE CS DEMODULATOR

The CS demodulators that we investigate operate on the basis of a mismatched³ function

$$\tilde{p}(\mathbf{y}|\mathbf{x}) = \exp(2\mathcal{R}\{\mathbf{x}^\dagger(\mathbf{V}\mathbf{y} - \mathbf{R}\hat{\mathbf{x}})\} - \mathbf{x}^\dagger \mathbf{G}\mathbf{x}) \quad (4)$$

for given feedbacks $\hat{\mathbf{x}}$, instead of the true conditional probability distribution function (pdf)

$$p(\mathbf{y}|\mathbf{x}) = \frac{1}{(\pi N_0)^N} \exp\left(-\frac{\|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2}{N_0}\right). \quad (5)$$

Note that $\tilde{p}(\mathbf{y}|\mathbf{x})$ may not be a valid pdf, but this is irrelevant for demodulation, see [46]. The matrices \mathbf{V} , \mathbf{R} , and \mathbf{G} are denoted as the front-end filter, IC matrix, and trellis-representation matrix as mentioned earlier, respectively. Without loss of generality, we have absorbed the noise power N_0 into them. Detection models (4) and (5) are equivalent if we set $\mathbf{V} = \mathbf{H}^\dagger/N_0$, $\mathbf{R} = \mathbf{0}$, and $\mathbf{G} = \mathbf{H}^\dagger \mathbf{H}/N_0$, in which case the CS demodulator represents the MAP.

The detection model (4) has its roots in Falconer and Magee’s paper [4] with an additional IC process, in which case it is described as

$$\tilde{T}(\mathbf{y}|\mathbf{x}) = \exp(-\|\mathbf{W}\mathbf{y} - \mathbf{T}\hat{\mathbf{x}} - \mathbf{F}\mathbf{x}\|^2) \quad (6)$$

By setting $\mathbf{T} = \mathbf{0}$, we obtain the same model in [4]. If identifying $\mathbf{V} = \mathbf{F}^\dagger \mathbf{W}$, $\mathbf{R} = \mathbf{F}^\dagger \mathbf{T}$, and $\mathbf{G} = \mathbf{F}^\dagger \mathbf{F}$, the model (6) is equivalent to (4) since

$$\begin{aligned} \tilde{T}(\mathbf{y}|\mathbf{x}) &\propto \exp(2\mathcal{R}\{\mathbf{x}^\dagger(\mathbf{F}^\dagger \mathbf{W}\mathbf{y} - \mathbf{F}^\dagger \mathbf{T}\hat{\mathbf{x}})\} - \mathbf{x}^\dagger \mathbf{F}^\dagger \mathbf{F}\mathbf{x}) \\ &= \exp(2\mathcal{R}\{\mathbf{x}^\dagger(\mathbf{V}\mathbf{y} - \mathbf{R}\hat{\mathbf{x}})\} - \mathbf{x}^\dagger \mathbf{G}\mathbf{x}). \end{aligned}$$

Detection model (6) is usually denoted as “Forney” model due to its Euclidean-distance form, while the general model (4) is called “Ungerboeck” model [47]–[49]. An advantage of the Ungerboeck model is that the parameter optimization through GMI-maximization is simpler [18] than

³By “mismatched” we mean that $\tilde{p}(\mathbf{y}|\mathbf{x})$ may not be a valid pdf and in general differs from the true pdf $p(\mathbf{y}|\mathbf{x})$ even with $\hat{\mathbf{x}} = \mathbf{0}$, but such a “mismatched” property is for the purpose of reducing the size of trellis-description in the BCJR.

the Forney model. However, as both models can be viewed as “natural” CS demodulators, we investigate both in CS demodulator design for iterative receivers in this paper.

In order to optimize $(\mathbf{V}, \mathbf{R}, \mathbf{G})$, we choose to work with the GMI which is an achievable rate for a receiver that operates on the basis of a mismatched version of the channel law. The GMI in nats per channel-use is defined as

$$I_{\text{GMI}} = -\mathbb{E}_{\mathbf{y}} [\log \tilde{p}(\mathbf{y})] + \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\log \tilde{p}(\mathbf{y}|\mathbf{x})] \quad (7)$$

where

$$\tilde{p}(\mathbf{y}) = \frac{1}{\pi^K} \int \tilde{p}(\mathbf{y}|\mathbf{x}) \exp(-\|\mathbf{x}\|^2) d\mathbf{x}$$

and the expectation is taken over the true pdfs $p(\mathbf{y})$ and $p(\mathbf{x}, \mathbf{y})$. Although finite constellations \mathcal{X} are used in practice, they are hard to analyze. In order to obtain a mathematically tractable problem, here we use a zero-mean complex-valued Gaussian constellation with an unit-variance for each entry in \mathbf{x} . With Gaussian inputs, the trellis discussed earlier has no proper meaning as the number of states is infinite. However, the Gaussian assumption is only made to design the CS parameters $(\mathbf{V}, \mathbf{R}, \mathbf{G})$ that can also be used for finite constellations.

We first state Theorem 1 that states the GMI for model (4).

Theorem 1: The GMI for the detection model (4) equals

$$\begin{aligned} I_{\text{GMI}}(\mathbf{V}, \mathbf{R}, \mathbf{G}) &= \log(\det(\mathbf{I} + \mathbf{G})) - \text{Tr}(\mathbf{G}) + 2\mathcal{R}\{\text{Tr}(\mathbf{V}\mathbf{H} - \mathbf{R}\mathbf{P})\} \\ &\quad - \text{Tr}((\mathbf{I} + \mathbf{G})^{-1}(\mathbf{V}(N_0\mathbf{I} + \mathbf{H}\mathbf{H}^\dagger)\mathbf{V}^\dagger \\ &\quad - 2\mathcal{R}\{\mathbf{V}\mathbf{H}\mathbf{P}\mathbf{R}^\dagger\} + \mathbf{R}\mathbf{P}\mathbf{R}^\dagger)). \end{aligned} \quad (8)$$

The proof of Theorem 1 is given in Appendix A. Here we use the fact $(\mathbf{I} + \mathbf{G}) \succ \mathbf{0}$ shown in [18], otherwise the GMI is not well-defined. With any parameters $(\mathbf{V}, \mathbf{R}, \mathbf{G})$, the GMI can be calculated via (8), however, they may not be optimal in the sense of GMI-maximization.

We next illustrate Theorem 1 with two examples.

Example 1: Extended Zero-Forcing filter (EZF). We extend the zero-Forcing filter [50] to only partly invert the channel so that a trellis-search is necessary after the EZF. In view of the CS demodulator, we select the parameters in (4) as:

$$\mathbf{V} = (\mathbf{I} + \mathbf{G})(\mathbf{H}^\dagger \mathbf{H})^{-1} \mathbf{H}^\dagger, \quad \text{and } \mathbf{R} = \mathbf{0},$$

and then optimize (8) over \mathbf{G} . To satisfy the constraint of having a trellis with $|\mathcal{X}|^v$ states, it must hold $\mathbf{G} = [\mathbf{G}]_v$. The optimal \mathbf{G} , in the sense of maximizing (8), is shown (later in Theorem 2) to satisfy

$$[(\mathbf{I} + \mathbf{G})^{-1}]_v = N_0[(\mathbf{H}^\dagger \mathbf{H})^{-1}]_v.$$

Utilizing Lemma 1, the GMI in (8) for the optimal \mathbf{G} equals

$$\begin{aligned} I_{\text{GMI}} &= \log(\det(\mathbf{I} + \mathbf{G})) + \text{Tr}(\mathbf{I} - N_0(\mathbf{H}^\dagger \mathbf{H})^{-1}(\mathbf{I} + \mathbf{G})) \\ &= \log(\det(\mathbf{I} + \mathbf{G})). \end{aligned}$$

Example 2: Truncated Matched filter (TMF). As previously mentioned, the MAP demodulator (5) can be written in the form (4) by setting $\mathbf{V} = \mathbf{H}^\dagger/N_0$, $\mathbf{R} = \mathbf{0}$, and $\mathbf{G} = \mathbf{H}^\dagger \mathbf{H}/N_0$.

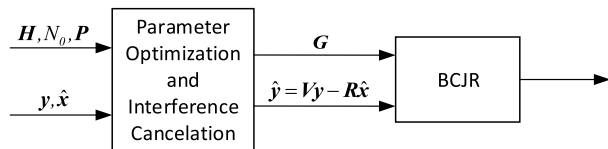


FIGURE 2. CS demodulator that maximizes the GMI with the tuple (x, \hat{y}) .

The front-end in this case is a matched filter [51] and the BCJR is implemented over the Ungerboeck model. To reach a trellis with $|\mathcal{X}|^\nu$ states, we truncate \mathbf{G} to its center $2\nu + 1$ diagonals, i.e., we use the following parameters in (4):

$$\mathbf{V} = \mathbf{H}^\dagger / N_0, \quad \mathbf{R} = \mathbf{0}, \quad \text{and} \quad \mathbf{G} = [\mathbf{H}^\dagger \mathbf{H} / N_0]_\nu.$$

With these settings in TMF, the GMI in (8) equals

$$I_{\text{GMI}} = \log(\det(\mathbf{I} + [\mathbf{H}^\dagger \mathbf{H} / N_0]_\nu)) - \text{Tr}(\mathbf{H}^\dagger \mathbf{H} (N_0 \mathbf{I} + [\mathbf{H}^\dagger \mathbf{H}]_\nu)^{-1} [\mathbf{H}^\dagger \mathbf{H}]_{\setminus \nu}).$$

B. CONSTRAINTS ON THE PARAMETER \mathbf{R} FOR CS DEMODULATORS

As illustrated in Fig. 2, our approach to design a CS demodulator consists of two steps:

- Construction of a signal $\hat{y} = \mathbf{V}\mathbf{y} - \mathbf{R}\hat{\mathbf{x}}$ based on the received signal \mathbf{y} and a prior mean $\hat{\mathbf{x}}$;
- BCJR demodulation of \hat{y} operating on a reduced number of states $|\mathcal{X}|^\nu$.

This procedure is analogous to LMMSE-PIC which first subtracts the interference followed by applying a Wiener filter, and then concludes by the BCJR that operates on a diagonal \mathbf{G} . The statistical behavior of $(\mathbf{x}, \hat{\mathbf{y}})$ can be superior to that of the original (\mathbf{x}, \mathbf{y}) , as the former tuple corresponds to a statistically different channel. As shown in Example 3 below, the GMI obtained with tuple $(\mathbf{x}, \hat{\mathbf{y}})$ with a perfect feedback $\hat{\mathbf{x}}$ can be infinitely large, which exceeds the channel capacity with the original tuple (\mathbf{x}, \mathbf{y}) . Therefore, the computed value of GMI may have little relevance for the performance of the transceiver system. In order for GMI to have bearing on performance, it is critical to put constraints on IC matrix \mathbf{R} .

Example 3: Let the system model be

$$\mathbf{y} = \mathbf{x} + \mathbf{n}.$$

and assume a perfect feedback, i.e., $\hat{\mathbf{x}} = \mathbf{x}$. The demodulator parameters are taken as $\mathbf{V} = \mathbf{0}$, $\mathbf{R} = -(1+\beta)\mathbf{I}$, and $\mathbf{G} = \beta\mathbf{I}$, where β is an arbitrary value. Then, we have

$$\hat{\mathbf{y}} = \mathbf{V}\mathbf{y} - \mathbf{R}\hat{\mathbf{x}} = (1+\beta)\mathbf{x},$$

and the GMI in (8) for the tuple $(\mathbf{x}, \hat{\mathbf{y}})$ equals

$$I_{\text{GMI}}(\mathbf{V}, \mathbf{R}, \mathbf{G}) = K(1 + \log(1 + \beta)).$$

In order to maximize the GMI, the demodulator will choose $\beta \rightarrow \infty$ to make I_{GMI} infinite. This is because, except for using the feedback information for IC, the demodulator uses the prior mean $\hat{\mathbf{x}}$ as a signal energy via \mathbf{R} . A demodulator equipped with these parameters will have significant error

propagation and does not have much operational meaning for an iterative receiver. Thus, we conclude that unless constraints are put on \mathbf{R} , the GMI value is not relevant. Three typical shapes of \mathbf{R} are specified in Fig. 3. All three have in common that rather than adding signal energy, the rationale of \mathbf{R} should be to remove interference. Therefore, at the very minimum the diagonal elements of \mathbf{R} should be constrained to zeros, so that the demodulation of each symbol in \mathbf{x} does not rely on its own prior mean $\hat{\mathbf{x}}$. Such a constraint is perfectly aligned with the operations of LMMSE-PIC, where \hat{x}_ℓ is not used for demodulation of x_ℓ . Further, the rationale of the constraints we impose on \mathbf{R} is to follow the principle of extrinsic information: The BCJR module should not rely on the prior information \hat{x}_ℓ when demodulating x_ℓ (this requires more than just the diagonal of \mathbf{R} to be zero).

We point out the fact that the GMI can exceed the channel capacity is a consequence of our choice not to include the side-information as a prior distribution on \mathbf{x} when evaluating the GMI. If we did, then the GMI is decaying with increasing quality of the side-information (due to the mutual information $I(\mathbf{x}, \mathbf{y}|\hat{\mathbf{x}})$ goes to 0 as $\hat{\mathbf{x}}$ becomes perfect). At last, we acknowledge the fact that a permutation of the columns of \mathbf{H} can boost the performance of the CS demodulator whenever $0 < \nu < K - 1$ for MIMO channels. However, minimum-phase conversions of ISI channels are not beneficial as we will solve for the optimal front-end filter.

IV. PARAMETER OPTIMIZATION FOR THE MIMO CASE

In this section, we elaborate the parameter optimization for MIMO channels. We introduce three different methods, namely, Method I, II, and III. We start with the classical Forney model (6) based demodulator, i.e., Method I, and then extend the model to the Ungerboeck model (4), i.e., Method II. As both Method I and Method II need gradient-based approaches for optimizations over target response, by carefully examining the properties of the CS demodulator with Ungerboeck model, we propose a suboptimal Method III which has an explicit construction and all CS parameters are in closed-forms.

A. METHOD I

In Method I, the CS demodulator is based on detection model (6) and the following structures of the CS parameters $(\mathbf{W}, \mathbf{T}, \mathbf{F})$ are imposed:

- The $K \times N$ matrix \mathbf{W} has no constraints.
- The $K \times K$ matrix \mathbf{F} is lower-triangular where only the main diagonal and the first ν (the memory-size of \mathbf{F}) lower diagonals are non-zero, i.e., \mathbf{F} is banded within diagonals $[0, \nu]$ ($0 \leq \nu < K - 1$). Moreover, the main diagonal of \mathbf{F} is constrained to have positive values.
- The $K \times K$ matrix \mathbf{T} is constrained to be zero wherever \mathbf{F} can take non-zero values.

The constraint on \mathbf{F} is to shorten the memory of the trellis in BCJR, while the purpose of the constraint on \mathbf{T} is to cancel the signal part that \mathbf{F} cannot handle. From Theorem 1 and by

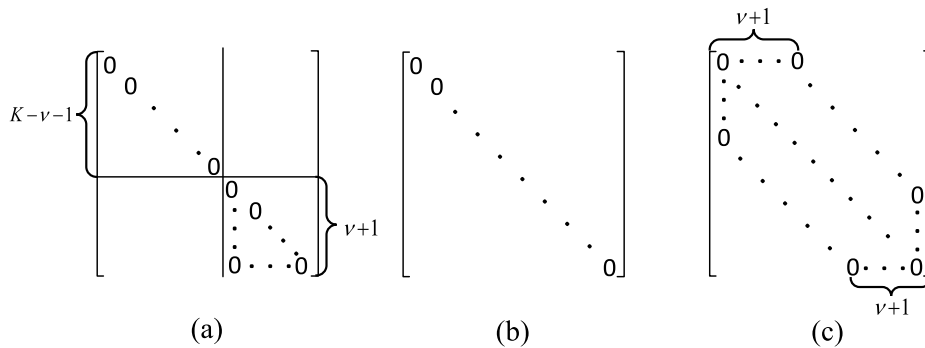


FIGURE 3. Three different types of shape of matrix R , where ν is the memory-size of F or G , i.e., the memory-size of the BCJR.

identifying $V = F^\dagger W$, $R = F^\dagger T$, and $G = F^\dagger F$, the GMI in (8) of Method I equals

$$I_{\text{GMI}}(\mathbf{W}, \mathbf{T}, \mathbf{F}) = \log(\det(\mathbf{I} + \mathbf{F}^\dagger \mathbf{F})) - \text{Tr}(\mathbf{F}^\dagger \mathbf{F}) + 2\mathcal{R}\{\text{Tr}(\mathbf{F}^\dagger(\mathbf{W}\mathbf{H} - \mathbf{T}\mathbf{P}))\} - \text{Tr}((\mathbf{I} + \mathbf{F}^\dagger \mathbf{F})^{-1} \mathbf{L}_1) \quad (9)$$

where

$$\mathbf{L}_1 = \mathbf{F}^\dagger \mathbf{W}(\mathbf{N}_0 \mathbf{I} + \mathbf{H}\mathbf{H}^\dagger) \mathbf{W}^\dagger \mathbf{F} - 2\mathcal{R}\{\mathbf{F}^\dagger \mathbf{W}\mathbf{H}\mathbf{P}\mathbf{T}^\dagger \mathbf{F}\} + \mathbf{F}^\dagger \mathbf{T}\mathbf{P}\mathbf{T}^\dagger \mathbf{F}.$$

With the aforementioned constraints on F and T , the matrix $R = F^\dagger T$ has a form of shape (a) in Fig. 3. That is, all diagonal elements are zero as well as the lower-triangular part of the $(\nu + 1) \times (\nu + 1)$ small matrix at the right-bottom corner.

In order to optimize (9) over (W, T, F) , we first introduce an $S \times K^2$ indication matrix Ω only consisting of ones and zeros,⁴ having a single 1 in each row, and S equals the number of elements in T that are allowed to be non-zero. Let $\mathbb{I}(\text{vec}(\mathbf{T}))$ be a vector that contains the positions where the vector $\text{vec}(\mathbf{T})$ is allowed to be non-zero. Then, the value of the k th entry in $\mathbb{I}(\text{vec}(\mathbf{T}))$ gives the column where row k of Ω is 1. That is, the $S \times 1$ vector $\Omega \text{vec}(\mathbf{T})$ stacks the columns of T on top of each other but with all elements that are constrained to zero removed.

With such a definition of Ω and defining two $K \times K$ matrices

$$\mathbf{M} = \mathbf{H}^\dagger (\mathbf{N}_0 \mathbf{I} + \mathbf{H}\mathbf{H}^\dagger)^{-1} \mathbf{H} - \mathbf{I}, \quad (10)$$

$$\tilde{\mathbf{M}} = \mathbf{P}(\mathbf{I} + \mathbf{M})\mathbf{P} - \mathbf{P}, \quad (11)$$

the GMI for the optimal W and T is given in Proposition 1 with the proof in Appendix B.

Proposition 1: Defining an $S \times K^2$ matrix

$$\mathbf{D} = \Omega((\mathbf{P}\mathbf{M}^*) \otimes \mathbf{I}_K),$$

the optimal W maximizing the GMI in (9) is

$$\mathbf{W}_{\text{opt}} = \mathbf{F}^{-\dagger} (\mathbf{I} + \mathbf{F}^\dagger \mathbf{F} + \mathbf{F}^\dagger \mathbf{T}\mathbf{P}) \mathbf{H}^\dagger (\mathbf{N}_0 \mathbf{I} + \mathbf{H}\mathbf{H}^\dagger)^{-1}, \quad (12)$$

⁴For instance, assuming $\mathbf{T} = \begin{bmatrix} 0 & 1 \\ 2 & 0 \end{bmatrix}$, then the indication matrix $\Omega = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$, and the vector $\Omega \text{vec}(\mathbf{T}) = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$.

and when $P \neq 0$ the optimal T is given by

$$\text{vec}(\mathbf{T}_{\text{opt}}) = -\Omega^T (\Omega (\tilde{\mathbf{M}}^* \otimes (\mathbf{F}(\mathbf{I} + \mathbf{F}^\dagger \mathbf{F})^{-1} \mathbf{F}^\dagger)) \Omega^T)^{-1} \text{vec}(\mathbf{F}). \quad (13)$$

With the optimal W and T , the GMI reads,

$$I_{\text{GMI}}(\mathbf{W}_{\text{opt}}, \mathbf{T}_{\text{opt}}, \mathbf{F}) = \begin{cases} I_1(\mathbf{F}), & \mathbf{P} = \mathbf{0}, \\ I_1(\mathbf{F}) + \delta_1(\mathbf{F}), & \mathbf{P} \neq \mathbf{0}, \end{cases} \quad (14)$$

where the functions $I_1(\mathbf{F})$ and $\delta_1(\mathbf{F})$ are defined as

$$I_1(\mathbf{F}) = K + \log(\det(\mathbf{I} + \mathbf{F}^\dagger \mathbf{F})) + \text{Tr}(\mathbf{M}(\mathbf{I} + \mathbf{F}^\dagger \mathbf{F})), \quad (15)$$

$$\delta_1(\mathbf{F}) = -\text{vec}(\mathbf{F})^\dagger \mathbf{D}^\dagger (\Omega (\tilde{\mathbf{M}}^* \otimes (\mathbf{F}(\mathbf{I} + \mathbf{F}^\dagger \mathbf{F})^{-1} \mathbf{F}^\dagger)) \Omega^T)^{-1} \mathbf{D} \text{vec}(\mathbf{F}). \quad (16)$$

Remark 1: From the definitions in (10) and (11), \mathbf{M} is the negative of MSE matrix and consequently, it holds that $\tilde{\mathbf{M}} \leq 0$. Hence, $\delta_1(\mathbf{F}) \geq 0$ represents the GMI increments from the soft-feedback.

Before maximization the GMI, we state Theorem 2 that deals with a general maximization problem.

Theorem 2: Define a scalar function I with respect to a $K \times K$ matrix \mathbf{G} as

$$I(\mathbf{G}) = K + \log(\det(\mathbf{I} + \mathbf{G})) + \text{Tr}(\mathbf{M}(\mathbf{I} + \mathbf{G})) \quad (17)$$

where \mathbf{G} satisfies $\mathbf{G} = [\mathbf{G}]_\nu$. Then, the optimal \mathbf{G} maximizing I is the unique solution that satisfies

$$[(\mathbf{I} + \mathbf{G}_{\text{opt}})^{-1}]_\nu = -[\mathbf{M}]_\nu. \quad (18)$$

With \mathbf{G}_{opt} , the maximal I equals

$$I(\mathbf{G}_{\text{opt}}) = \log(\det(\mathbf{I} + \mathbf{G}_{\text{opt}})). \quad (19)$$

Proof: Taking the first order differential of I with respect to \mathbf{G} and noticing that \mathbf{G} is banded within diagonals $[-\nu, \nu]$, yields (18) after some manipulations. The existence and uniqueness of such an optimal solution for (18) is proved in [52, Th. 2] and also illustrated in [18, Proposition 2]. By Lemma 1 and (18), it holds that $\text{Tr}([\mathbf{I} + \mathbf{G}_{\text{opt}}]^{-1} \mathbf{M}) = -K$, and then (19) follows. ■

Optimizing over F in (14) when $P \neq 0$ is difficult and cannot be carried out in closed-form. In Appendix C we show

by an example that (14) is in general non-concave. Therefore, a gradient based numerical optimization procedure is utilized to search for the optimal \mathbf{F} . In the i th iteration, we construct

$$\mathbf{F}^{(i)} = \mathbf{F}^{(i-1)} + \nabla_{\mathbf{F}^*} I_{\text{GMI}}(\mathbf{W}_{\text{opt}}, \mathbf{T}_{\text{opt}}, \mathbf{F}^{(i-1)})$$

where $\nabla_{\mathbf{F}^*} I_{\text{GMI}}(\mathbf{W}_{\text{opt}}, \mathbf{T}_{\text{opt}}, \mathbf{F})$ is the conjugate gradient of the GMI with respect to (the non-zero part of) \mathbf{F} , which is derived in Appendix D.

If replacing $\mathbf{F}^\dagger \mathbf{F}$ by \mathbf{G} , (15) has the same form as (17) when $\mathbf{P} = \mathbf{0}$, and \mathbf{G}_{opt} is in closed-form shown in Theorem 2. If $\mathbf{G}_{\text{opt}} \succ \mathbf{0}$, then the optimal \mathbf{F} is equal to its Cholesky decomposition. Whenever it is not, a gradient based numerical optimization is utilized to optimize (15), and \mathbf{G}_{opt} from Theorem 2 is used to initialize the starting point of \mathbf{F} , which has been observed to be reliable.

Next we establish a connection between the front-end filter \mathbf{W} and IC matrix \mathbf{T} in Method I.

Proposition 2: For $\mathbf{P} \neq \mathbf{0}$, and with the optimal \mathbf{W} and \mathbf{T} , the matrix $\mathbf{F}^\dagger(\mathbf{W}_{\text{opt}}\mathbf{H} - \mathbf{T}_{\text{opt}})$ is banded within diagonals $[-\nu, K-1]$.

Proof: Noting that $\mathbf{\Omega}^T \mathbf{\Omega} \text{vec}(\mathbf{T}_{\text{opt}}) = \text{vec}(\mathbf{T}_{\text{opt}})$ and using $\mathbf{\Omega} \mathbf{\Omega}^T = \mathbf{I}$, from (13) and (76) it holds that

$$\begin{aligned} & \mathbf{\Omega}(\tilde{\mathbf{M}}^* \otimes (\mathbf{F}(\mathbf{I} + \mathbf{F}^\dagger \mathbf{F})^{-1} \mathbf{F}^\dagger)) \mathbf{\Omega}^T \mathbf{\Omega} \text{vec}(\mathbf{T}_{\text{opt}}) \\ &= \mathbf{\Omega} \text{vec}(\mathbf{F}(\mathbf{I} + \mathbf{F}^\dagger \mathbf{F})^{-1} \mathbf{F}^\dagger \mathbf{T}_{\text{opt}} \tilde{\mathbf{M}}) \\ & \quad - \mathbf{\Omega} \text{vec}(\mathbf{FMP}), \end{aligned} \quad (20)$$

which shows that the elements of the matrix

$$\Delta = \mathbf{F}(\mathbf{I} + \mathbf{F}^\dagger \mathbf{F})^{-1} \mathbf{F}^\dagger \mathbf{T}_{\text{opt}} \tilde{\mathbf{M}} + \mathbf{FMP}$$

are zeros wherever \mathbf{T} can be non-zero. Hence Δ is banded within diagonals $[0, \nu]$. On the other hand, with the optimal \mathbf{W} in (12), and \mathbf{M} , $\tilde{\mathbf{M}}$ defined in (10) and (11), respectively, we have

$$\mathbf{F}^\dagger(\mathbf{W}_{\text{opt}}\mathbf{H} - \mathbf{T}_{\text{opt}}) - (\mathbf{I} + \mathbf{F}^\dagger \mathbf{F}) = (\mathbf{I} + \mathbf{F}^\dagger \mathbf{F}) \mathbf{F}^{-1} \Delta \mathbf{P}^{-1}. \quad (21)$$

By noting that \mathbf{F}^{-1} is lower-triangular, $\mathbf{I} + \mathbf{F}^\dagger \mathbf{F}$ is banded within diagonals $[-\nu, \nu]$ and \mathbf{P} is diagonal. Utilizing Lemma 2, the r.h.s of (21) is banded within diagonals $[-\nu, K-1]$. Therefore, $\mathbf{F}^\dagger(\mathbf{W}_{\text{opt}}\mathbf{H} - \mathbf{T}_{\text{opt}})$ is banded within diagonals $[-\nu, K-1]$. ■

Proposition 2 reveals an interesting and somewhat surprising fact that although the BCJR only has a memory-size ν , the interference outside the memory-size shall not be perfectly canceled with the optimal CS demodulator in Method I. As will be shown later, such a property also holds for the other two designs of CS demodulator, namely, Method II and III.

B. METHOD II

Method II origins from Ungerboeck's 1974 paper [47]. Different from Method I, an Ungerboeck detection model (4) instead of the Forney model (6) is applied. The Ungerboeck model has been extensively discussed in [48], [49], and [53]. The system model (4) has the following constraints:

- The $K \times N$ matrix \mathbf{V} has no constraints.

- The $K \times K$ matrix \mathbf{G} is Hermitian and satisfies $\mathbf{G} = [\mathbf{G}]_\nu$, and $\mathbf{I} + \mathbf{G} \succ \mathbf{0}$, where ν is the memory-size of \mathbf{G} .
- The $K \times K$ matrix \mathbf{R} can have specified shapes, but at a minimum the main diagonal is constrained to zeros.

Instead of $(\mathbf{W}, \mathbf{T}, \mathbf{F})$, in Method II we optimize $(\mathbf{V}, \mathbf{R}, \mathbf{G})$ for (8). The same definition of an indication matrix $\mathbf{\Omega}$ is used as in Method I, but now $\mathbf{\Omega}$ corresponds to \mathbf{R} instead of \mathbf{T} . We continue to let S denote the number of elements that are allowed to be non-zero in \mathbf{R} . That is, the $S \times 1$ vector $\mathbf{\Omega} \text{vec}(\mathbf{R})$ stacks the columns of \mathbf{R} on top of each other and with all elements that are constrained to zeros removed.

In Method II, we have Proposition 3 showing the GMI with optimal \mathbf{V} and \mathbf{R} , with the proof given in Appendix E.

Proposition 3: Define an $S \times 1$ vector $\mathbf{d} = \mathbf{\Omega} \text{vec}(\mathbf{MP})$, the optimal \mathbf{V} for the GMI in (8) is

$$\mathbf{V}_{\text{opt}} = (\mathbf{I} + \mathbf{G} + \mathbf{RP}) \mathbf{H}^\dagger (\mathbf{H} \mathbf{H}^\dagger + N_0 \mathbf{I})^{-1}, \quad (22)$$

and when $\mathbf{P} \neq \mathbf{0}$ the optimal \mathbf{R} is given by

$$\text{vec}(\mathbf{R}_{\text{opt}}) = -\mathbf{\Omega}^T (\mathbf{\Omega}(\tilde{\mathbf{M}}^* \otimes (\mathbf{I} + \mathbf{G})^{-1}) \mathbf{\Omega}^T)^{-1} \mathbf{d}. \quad (23)$$

With the optimal \mathbf{V} and \mathbf{R} , the GMI in (8) equals

$$I_{\text{GMI}}(\mathbf{V}_{\text{opt}}, \mathbf{R}_{\text{opt}}, \mathbf{G}) = \begin{cases} I_2(\mathbf{G}), & \mathbf{P} = \mathbf{0}, \\ I_2(\mathbf{G}) + \delta_2(\mathbf{G}), & \mathbf{P} \neq \mathbf{0}, \end{cases} \quad (24)$$

where the functions $I_2(\mathbf{G})$ and $\delta_2(\mathbf{G})$ are defined as

$$I_2(\mathbf{G}) = K + \log(\det(\mathbf{I} + \mathbf{G})) + \text{Tr}(\mathbf{M}(\mathbf{I} + \mathbf{G})), \quad (25)$$

$$\delta_2(\mathbf{G}) = -\mathbf{d}^\dagger (\mathbf{\Omega}(\tilde{\mathbf{M}}^* \otimes (\mathbf{I} + \mathbf{G})^{-1}) \mathbf{\Omega}^T)^{-1} \mathbf{d}. \quad (26)$$

Similar to $\delta_1(\mathbf{F})$ in Method I, $\delta_2(\mathbf{G}) \geq 0$ represents the GMI increments from the feedback. Further, when $\mathbf{P} \neq \mathbf{0}$, the optimization over \mathbf{G} in (24) also uses a gradient based numerical optimization, and the gradient of $I_{\text{GMI}}(\mathbf{V}_{\text{opt}}, \mathbf{R}_{\text{opt}}, \mathbf{G})$ with respect to (the non-zero part of) \mathbf{G} is provided in Appendix F. The closed-form \mathbf{G} from Theorem 2 can be used as the starting point. However, different from Method I, the optimization procedure in this case is concave whose proof is in Appendix G.

Although the optimal \mathbf{R} is solved in closed-form in (23), we shall specify the constraint (reflected by $\mathbf{\Omega}$) on it. We consider two types of \mathbf{R} in Method II. Firstly, as we are interested in the comparison between Method I and II, we consider the shape (a) in Fig. 3, which has the same shape as for $\mathbf{R} = \mathbf{F}^\dagger \mathbf{T}$ in Method I. Secondly, we consider a band-shaped \mathbf{R} with memory-size ν_R , where shape (b) and (c) in Fig. 3 are typical cases with $\nu_R = 0$ and $\nu_R = \nu$, respectively. With shape (b), we only limit the diagonal elements of \mathbf{R} to be zero and intend to eliminate the interference as much as possible. With shape (c), we limit \mathbf{R} to have the opposite form of \mathbf{G} , that is, the elements of \mathbf{R} are constrained to be zero wherever \mathbf{G} is non-zero. The intention is to only cancel the interference that the BCJR represented by \mathbf{G} cannot handle. Shape (c) is based on the same idea as Method I, but now operates on the Ungerboeck model.

The connection between the optimal front-end filter \mathbf{V} and IC matrix \mathbf{R} in Method II is established in Proposition 4.

Proposition 4: For $\mathbf{P} \neq \mathbf{0}$ and the optimal \mathbf{V} and \mathbf{R} ,

$$[\mathbf{V}_{\text{opt}}\mathbf{H}]_{\setminus(\nu+\nu_R)} = [\mathbf{R}_{\text{opt}}]_{\setminus(\nu+\nu_R)}. \quad (27)$$

That is, the elements of $\mathbf{V}_{\text{opt}}\mathbf{H}$ and \mathbf{R}_{opt} are equal outside the center $2(\nu+\nu_R)+1$ diagonals for any \mathbf{G} that is banded within diagonals $[-\nu, \nu]$, where $\nu_R = 0$ for \mathbf{R} with both shapes (a) and (b), and $\nu_R = \nu$ for \mathbf{R} with shape (c), respectively.

Proof: Following the similar steps in the proof of Proposition 2, (23) can be rewritten as

$$\mathbf{\Omega}\text{vec}((\mathbf{I} + \mathbf{G})^{-1}\mathbf{R}_{\text{opt}}\tilde{\mathbf{M}}) = -\mathbf{\Omega}\text{vec}(\mathbf{M}\mathbf{P}). \quad (28)$$

It shows that the elements of

$$\Delta = (\mathbf{I} + \mathbf{G})^{-1}\mathbf{R}_{\text{opt}}\tilde{\mathbf{M}}\mathbf{P}^{-1} + \mathbf{M}$$

are zeros wherever \mathbf{R} can be non-zero. On the other hand, with the optimal \mathbf{V} in (22) we have

$$\mathbf{V}_{\text{opt}}\mathbf{H} - \mathbf{R}_{\text{opt}} - (\mathbf{I} + \mathbf{G}) = (\mathbf{I} + \mathbf{G})\Delta. \quad (29)$$

As $\mathbf{I} + \mathbf{G}$ is banded within diagonals $[-\nu, \nu]$, using Lemma 2 and with the three shapes⁵ of \mathbf{R} in Fig. 3, it can be shown that the r.h.s of (29) is banded within diagonals $[-(\nu + \nu_R), \nu + \nu_R]$, where $\nu_R = 0$ for the shape (a) and (b), and $\nu_R = \nu$ for the shape (c). Therefore, $\mathbf{V}_{\text{opt}}\mathbf{H} - \mathbf{R}_{\text{opt}}$ on the l.h.s of (29) is banded within diagonals $[-(\nu + \nu_R), \nu + \nu_R]$. ■

The same as Proposition 2 for Method I, Proposition 4 shows that the signal part not considered in \mathbf{G} (the BCJR) shall not be perfectly canceled inside the center $2(\nu + \nu_R)+1$ diagonals, instead of the center $2\nu+1$ diagonals where \mathbf{G} is constrained to be non-zero. With LMMSE-PIC, we have $\nu = \nu_R = 0$ and Proposition 4 is natural and frequently used. However, when $\nu_R > 0$, a more general property is now revealed that $\mathbf{V}_{\text{opt}}\mathbf{H}$ and \mathbf{R} are only equal outside the center $2(\nu + \nu_R)+1$ diagonals.

C. METHOD III

So far we have discussed two types of CS demodulators based on Forney and Ungerboeck detection models. One disadvantage is that both methods need numerical optimizations to obtain the optimal target responses. Next we construct a third CS demodulator that has closed-form solutions for all CS parameters, whose GMI is slightly suboptimal.

Method III relies on the same operations as Method II with $\mathbf{P} = \mathbf{0}$. By inserting \mathbf{V}_{opt} in (22) back into (4) and setting $\mathbf{R} = \mathbf{0}$, the CS demodulator actually operates on the mismatched function

$$\begin{aligned} \tilde{p}(\mathbf{y}|\mathbf{x}) &= \exp(2\mathcal{R}\{\mathbf{x}^\dagger\mathbf{V}_{\text{opt}}\mathbf{y}\} - \mathbf{x}^\dagger\mathbf{G}\mathbf{x}) \\ &= \exp(2\mathcal{R}\{\mathbf{x}^\dagger(\mathbf{I} + \mathbf{G})\tilde{\mathbf{x}}\} - \mathbf{x}^\dagger\mathbf{G}\mathbf{x}) \end{aligned} \quad (30)$$

where

$$\tilde{\mathbf{x}} = \mathbf{H}^\dagger(\mathbf{H}\mathbf{H}^\dagger + N_0\mathbf{I})^{-1}\mathbf{y}$$

is the LMMSE estimate of \mathbf{x} . As can be seen from (30), the BCJR is based on $\tilde{\mathbf{x}}$. With soft-feedback we can therefore

⁵The case that \mathbf{R} with shape (a) is slightly different, but it can be verified straightforwardly.

replace $\tilde{\mathbf{x}}$ by LMMSE-PIC estimates $\tilde{\mathbf{x}}$. That is, instead of (30) we now operate on

$$\tilde{p}(\mathbf{y}|\mathbf{x}, \tilde{\mathbf{x}}) = \exp(2\mathcal{R}\{\mathbf{x}^\dagger(\mathbf{I} + \mathbf{G})\tilde{\mathbf{x}}\} - \mathbf{x}^\dagger\mathbf{G}\mathbf{x}) \quad (31)$$

where \mathbf{G} has the same banded-shape as the first two methods, but optimized according to $\tilde{\mathbf{x}}$.

The estimate $\tilde{\mathbf{x}}$ is constructed as follows. As we prefer to handle the interference through the trellis-search process, the IC should not be present within the memory-size ν . In other words, the signal vector after the IC that is used to form the k th symbol of $\tilde{\mathbf{x}}$ is denoted as

$$\tilde{\mathbf{y}}_k = \mathbf{y} - \sum_{n \in \mathcal{A}_k} \mathbf{h}_n \hat{x}_n \quad (32)$$

where

$$\mathcal{A}_k = \{0 \leq n \leq K-1 : n \notin [\max(0, k-\nu), \min(k+\nu, K-1)]\}.$$

Denoting p_n as the n th diagonal element of \mathbf{P} , the Wiener-filter coefficients [54] for the k th symbol are calculated as

$$\hat{\mathbf{w}}_k = \mathbf{h}_k^\dagger(\mathbf{H}^\dagger\mathbf{C}_k\mathbf{H} + N_0\mathbf{I})^{-1} \quad (33)$$

where \mathbf{C}_k is a diagonal whose n th diagonal element is

$$C_k(n) = \begin{cases} 1 - p_n, & k \in \mathcal{A}_k, \\ 1, & \text{otherwise.} \end{cases} \quad (34)$$

The estimate $\tilde{\mathbf{x}}$ is then obtained through

$$\begin{aligned} \tilde{\mathbf{x}} &= [\hat{\mathbf{w}}_1\tilde{\mathbf{y}}_1, \hat{\mathbf{w}}_2\tilde{\mathbf{y}}_2, \dots, \hat{\mathbf{w}}_K\tilde{\mathbf{y}}_K]^T \\ &= \hat{\mathbf{W}}\mathbf{y} - \hat{\mathbf{C}}\hat{\mathbf{x}} \end{aligned} \quad (35)$$

where the coefficient matrix $\hat{\mathbf{W}}$ and IC matrix $\hat{\mathbf{C}}$ defined as

$$\hat{\mathbf{W}} = [\hat{\mathbf{w}}_1^T, \hat{\mathbf{w}}_2^T, \dots, \hat{\mathbf{w}}_K^T]^T, \quad (36)$$

$$\hat{\mathbf{C}} = [\hat{\mathbf{W}}\mathbf{H}]_{\setminus\nu}. \quad (37)$$

Inserting $\tilde{\mathbf{x}}$ in (35) back into (31), the detection model we operate on is

$$\tilde{p}(\mathbf{y}|\mathbf{x}) = \exp(2\mathcal{R}\{\mathbf{x}^\dagger((\mathbf{I} + \mathbf{G})\hat{\mathbf{W}}\mathbf{y} - (\mathbf{I} + \mathbf{G})\hat{\mathbf{C}}\hat{\mathbf{x}})\} - \mathbf{x}^\dagger\mathbf{G}\mathbf{x}). \quad (38)$$

Note that (38) is a also special case of (4), by identifying

$$\begin{aligned} \mathbf{V} &= (\mathbf{I} + \mathbf{G})\tilde{\mathbf{W}}, \\ \mathbf{R} &= (\mathbf{I} + \mathbf{G})\tilde{\mathbf{C}}. \end{aligned}$$

The GMI in (8) in this case reads, after some manipulations,

$$I_{\text{GMI}}(\mathbf{G}) = K + \log(\det(\mathbf{I} + \mathbf{G})) + \text{Tr}(\hat{\mathbf{M}}(\mathbf{I} + \mathbf{G})), \quad (39)$$

with $\hat{\mathbf{M}}$ (the updated \mathbf{M} in Method II) defined as⁶

$$\begin{aligned} \hat{\mathbf{M}} &= \hat{\mathbf{W}}\mathbf{H}\mathbf{P}\hat{\mathbf{C}}^\dagger + \hat{\mathbf{W}}\mathbf{H} - \mathbf{P}\hat{\mathbf{C}}^\dagger \\ &\quad (\hat{\mathbf{W}}\mathbf{H}\mathbf{P}\hat{\mathbf{C}}^\dagger + \hat{\mathbf{W}}\mathbf{H} - \mathbf{P}\hat{\mathbf{C}}^\dagger)^\dagger \\ &\quad - \hat{\mathbf{W}}(\mathbf{H}\mathbf{H}^\dagger + N_0\mathbf{I})\hat{\mathbf{W}}^\dagger - \hat{\mathbf{C}}\mathbf{P}\hat{\mathbf{C}}^\dagger - \mathbf{I}, \end{aligned} \quad (40)$$

⁶It can also be shown that $\hat{\mathbf{M}}$ is the negative of the updated MSE matrix, i.e., $\hat{\mathbf{M}} = -\mathbb{E}[(\mathbf{x} - \tilde{\mathbf{x}})(\mathbf{x} - \tilde{\mathbf{x}})^\dagger] = -\mathbb{E}[(\mathbf{x} - \hat{\mathbf{W}}\mathbf{y} + \hat{\mathbf{C}}\hat{\mathbf{x}})(\mathbf{x} - \hat{\mathbf{W}}\mathbf{y} + \hat{\mathbf{C}}\hat{\mathbf{x}})^\dagger]$.

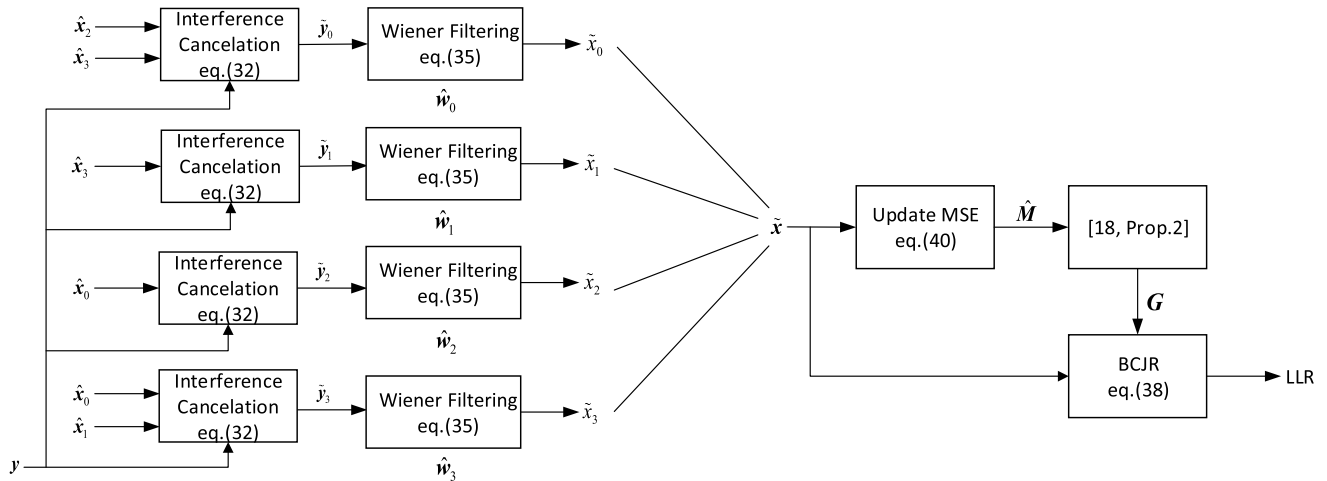


FIGURE 4. An graphical overview of Method III with $K = 4$ and $\nu = 1$.

The optimal \mathbf{G} for (39) is then obtained from Theorem 2, and the optimal GMI is

$$I_{\text{GMI}}(\mathbf{G}_{\text{opt}}) = \log(\det(\mathbf{I} + \mathbf{G}_{\text{opt}})).$$

An graphical overview of Method III for $K = 4$ and $\nu = 1$ is illustrated in Fig. 4. For any \mathbf{G} with memory-size ν , the IC matrix $(\mathbf{I} + \mathbf{G})\tilde{\mathbf{C}}$ is zero along the main diagonal, which guarantees that the extrinsic information will not be used for current symbols in the IC process. In GMI sense, Method III will not outperform Method II with \mathbf{R} in shape (b), but it can outperform the GMI of Method II with \mathbf{R} in shape (c), as it can be verified that \mathbf{R} in shape (c) has zeros at positions where $(\mathbf{I} + \mathbf{G})\hat{\mathbf{C}}$ are also zeros.

Remark 2: Since $\hat{\mathbf{W}}\mathbf{H} - \hat{\mathbf{C}} = [\hat{\mathbf{W}}\mathbf{H}]_{\nu}$ and by Lemma 2, $(\mathbf{I} + \mathbf{G})(\hat{\mathbf{W}}\mathbf{H} - \hat{\mathbf{C}})$ is banded within diagonals $[-2\nu, 2\nu]$. Therefore, it holds that $[(\mathbf{I} + \mathbf{G})\hat{\mathbf{W}}\mathbf{H}]_{\setminus 2\nu} = [(\mathbf{I} + \mathbf{G})\hat{\mathbf{C}}]_{\setminus 2\nu}$, and Proposition 4 is also true for Method III where $\nu_{\text{R}} = \nu$.

V. PARAMETER OPTIMIZATION FOR THE ISI CASE

In this section, we extend the CS demodulators to the ISI case with the block length K being sufficiently large and $\mathbf{P} = \alpha\mathbf{I}$. The formulas for the GMI in (8), (9), and (39) can be directly applied to (1), but as the rate I_{GMI} depends on K , we consider an asymptotic rate

$$\bar{I} = \lim_{K \rightarrow \infty} \frac{1}{K} I_{\text{GMI}}.$$

Ideally, in the ISI case the front-end matrix \mathbf{V} and IC matrix \mathbf{R} corresponding to linear filtering operations with infinite taps, but in practice they have finite tap-lengths. Hence, we analyze the properties of \mathbf{V} , \mathbf{R} (and also \mathbf{W} , \mathbf{T}) with finite numbers of taps and approximate them by band-shaped Toeplitz matrices. Further, the trellis-representation matrix \mathbf{G} (and \mathbf{F}), and channel matrix \mathbf{H} are also band-shaped Toeplitz matrices. Therefore, in the ISI case all matrices we consider are assumed to be band-shaped Toeplitz matrices,

whose dimensions are sufficiently large such that the asymptotic properties can be analyzed. In [55] a complete theoretic machinery for ISI channels is derived and a result is that as $K \rightarrow \infty$ the linear convolution in (1) can be replaced by a circular convolution.

In the following, we denote the Fourier series introduced by a band-shaped Toeplitz matrix \mathbf{E} with infinitely large dimensions as $E(\omega)$, where \mathbf{E} is constrained to zero outside the middle $2N_{\text{E}} + 1$ diagonals, and N_{E} is referred to as the tap-length of $E(\omega)$. The Fourier series $E(\omega)$ is defined as

$$E(\omega) = \sum_{k=-N_{\text{E}}}^{N_{\text{E}}} e_k \exp(jk\omega)$$

and specified by a vector

$$\mathbf{e} = [e_{-N_{\text{E}}}, \dots, e_{-1}, e_0, e_1, \dots, e_{N_{\text{E}}}]$$

where e_0 is the element on the main diagonal and e_k is the element on k th lower ($k > 0$) or upper ($k < 0$) diagonal. As all quantities are evaluated as K grows large, $E(\omega)$ approaches the eigenvalue distribution of \mathbf{E} (see [40], [56] for a precise statement of this result).

We first state Theorem 3, which is an asymptotic version of Theorem 2 for ISI channels.

Theorem 3: Assume two band-shaped Toeplitz matrices \mathbf{G} and \mathbf{M} with infinitely large dimensions satisfying the conditions: $[\mathbf{G}]_{\setminus \nu} = \mathbf{0}$, $\mathbf{I} + \mathbf{G} > \mathbf{0}$, and $\mathbf{M} < \mathbf{0}$, and define a scalar function

$$\bar{I} = 1 + \frac{1}{2\pi} \int_{-\pi}^{\pi} (\log(1 + G(\omega)) + M(\omega)(1 + G(\omega))) d\omega. \quad (41)$$

Then, the optimal $G(\omega)$ that maximizes \bar{I} is

$$G_{\text{opt}}(\omega) = |u_0 + \hat{\mathbf{u}}\varphi(\omega)|^2 - 1,$$

where the $1 \times \nu$ vector

$$\varphi(\omega) = [\exp(j\omega), \exp(j2\omega), \dots, \exp(j\nu\omega)]^T,$$

and

$$u_0 = \frac{1}{\sqrt{\boldsymbol{\tau}_1^\dagger \boldsymbol{\tau}_2^{-1} \boldsymbol{\tau}_1 - \tau_0}},$$

$$\hat{\mathbf{u}} = -u_0 \boldsymbol{\tau}_1^\dagger \boldsymbol{\tau}_2^{-1}. \quad (42)$$

The real scalar τ_0 , $\nu \times 1$ vector $\boldsymbol{\tau}_1$, and $\nu \times \nu$ matrix $\boldsymbol{\tau}_2$ are defined as

$$\tau_0 = \frac{1}{2\pi} \int_{-\pi}^{\pi} M(\omega) d\omega,$$

$$\boldsymbol{\tau}_1 = \frac{1}{2\pi} \int_{-\pi}^{\pi} M(\omega) \boldsymbol{\varphi}(\omega) d\omega,$$

$$\boldsymbol{\tau}_2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} M(\omega) \boldsymbol{\varphi}(\omega) \boldsymbol{\varphi}(\omega)^\dagger d\omega,$$

respectively. Further, with $G_{\text{opt}}(\omega)$ the optimal \bar{I} is

$$\bar{I} = 2 \log(u_0). \quad (43)$$

Proof: As $\mathbf{I} + \mathbf{G} > \mathbf{0}$, we have $1 + G(\omega) = |U(\omega)|^2$ where $U(\omega) = u_0 + \hat{\mathbf{u}}\boldsymbol{\varphi}(\omega)$ and $\hat{\mathbf{u}} = [u_1, u_2, \dots, u_\nu]$. Then, \bar{I} in (41) can be written as

$$\bar{I} = 1 + 2 \log(u_0) + \frac{1}{2\pi} \int_{-\pi}^{\pi} M(\omega) (u_0^2 + 2\mathcal{R}\{u_0 \hat{\mathbf{u}}\boldsymbol{\varphi}(\omega)\} + \hat{\mathbf{u}}\boldsymbol{\varphi}(\omega)\boldsymbol{\varphi}^\dagger(\omega)\hat{\mathbf{u}}^\dagger) d\omega. \quad (44)$$

Taking the first order derivatives with respect to both u_0 and $\hat{\mathbf{u}}$ and optimizing them directly yields the optimal solutions in (42). Inserting (42) back into (44) and after some manipulations, the optimal asymptotic rate is in (43). ■

A. METHOD I

The structures of $(\mathbf{W}, \mathbf{T}, \mathbf{F})$ are the same as in Sec. IV-A, except that now the matrices have infinite dimensions. Applying Szegő's eigenvalue distribution theorem [56] to (9), the asymptotic rate equals

$$\bar{I}(W(\omega), T(\omega), F(\omega))$$

$$= \lim_{K \rightarrow \infty} \frac{1}{K} J_{\text{GMI}}(\mathbf{W}, \mathbf{T}, \mathbf{F})$$

$$= \frac{1}{2\pi} \int_{-\pi}^{\pi} \left(\log(1 + |F(\omega)|^2) - |F(\omega)|^2 - \frac{L_1(\omega)}{1 + |F(\omega)|^2} \right) d\omega$$

$$+ \frac{1}{\pi} \int_{-\pi}^{\pi} \mathcal{R}\{F^*(\omega)(W(\omega)H(\omega) - \alpha T(\omega))\} d\omega \quad (45)$$

where

$$L_1(\omega) = |F(\omega)W(\omega)|^2 (N_0 + |H(\omega)|^2) + \alpha |F(\omega)T(\omega)|^2 - 2\alpha |F(\omega)|^2 \mathcal{R}\{H(\omega)W(\omega)T^*(\omega)\}.$$

Note that, the Fourier series associated to \mathbf{M} and $\tilde{\mathbf{M}}$ in (10) and (11) are

$$M(\omega) = \frac{|H(\omega)|^2}{N_0 + |H(\omega)|^2} - 1, \quad (46)$$

$$\tilde{M}(\omega) = \alpha^2 (M(\omega) + 1) - \alpha. \quad (47)$$

Further, define a $(2N_T - \nu) \times 1$ vector

$$\boldsymbol{\phi}(\omega) = [\exp(-jN_T\omega), \dots, \exp(-j(\nu+1)\omega), \exp(j(\nu+1)\omega), \dots, \exp(jN_T\omega)]^T, \quad (48)$$

a $(2N_T - \nu) \times 1$ vector $\boldsymbol{\epsilon}_1$, and a $(2N_T - \nu) \times (2N_T - \nu)$ Hermitian matrix $\boldsymbol{\epsilon}_2$ as

$$\boldsymbol{\epsilon}_1 = \frac{\alpha}{2\pi} \int_{-\pi}^{\pi} M(\omega) F^*(\omega) \boldsymbol{\phi}(\omega) d\omega,$$

$$\boldsymbol{\epsilon}_2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{\tilde{M}(\omega) |F(\omega)|^2 \boldsymbol{\phi}(\omega) \boldsymbol{\phi}(\omega)^\dagger}{1 + |F(\omega)|^2} d\omega, \quad (49)$$

respectively, where N_T is the tap-length of $T(\omega)$, and $\nu + 1$ is the band-size that matrix \mathbf{T} is constrained to zero. Then, we have Proposition 5 with its proof⁷ given in Appendix H.

Proposition 5: The optimal $W(\omega)$ for the asymptotic rate in (45) is

$$W_{\text{opt}}(\omega) = \frac{H^*(\omega)}{F^*(\omega)(N_0 + |H(\omega)|^2)} (1 + |F(\omega)|^2 + \alpha F^*(\omega) T_{\text{opt}}(\omega)), \quad (50)$$

and when $0 < \alpha \leq 1$, the optimal $T(\omega)$ reads

$$T_{\text{opt}}(\omega) = -\boldsymbol{\epsilon}_1^\dagger \boldsymbol{\epsilon}_2^{-1} \boldsymbol{\phi}(\omega). \quad (51)$$

With the optimal $W(\omega)$ and $T(\omega)$, the asymptotic rate equals

$$\bar{I}(W_{\text{opt}}(\omega), T_{\text{opt}}(\omega), F(\omega))$$

$$= \begin{cases} \bar{I}_1(F(\omega)), & \alpha = 0, \\ \bar{I}_1(F(\omega)) + \bar{\delta}_1(F(\omega)), & 0 < \alpha \leq 1. \end{cases} \quad (52)$$

The functions $\bar{I}_1(F(\omega))$ and $\bar{\delta}_1(F(\omega))$ are defined as,⁸

$$\bar{I}_1(F(\omega)) = 1 + \frac{1}{2\pi} \int_{-\pi}^{\pi} \left(\log(1 + |F(\omega)|^2) + M(\omega)(1 + |F(\omega)|^2) \right) d\omega, \quad (53)$$

$$\bar{\delta}_1(F(\omega)) = -\boldsymbol{\epsilon}_1^\dagger \boldsymbol{\epsilon}_2^{-1} \boldsymbol{\epsilon}_1. \quad (54)$$

In ISI case, Method I is not concave and an example is also provided in Appendix C. Hence, a gradient based optimization is used to optimize $F(\omega)$ with the optimal $G_{\text{opt}}(\omega)$ from Theorem 3 utilized for initialization. The connection between the optimal front-end filter $W(\omega)$ and the IC filter $T(\omega)$ in Proposition 2 also holds for ISI channels. An asymptotic version of Proposition 2 is stated in Proposition 6.

Proposition 6: When $0 < \alpha \leq 1$, $a_k = b_k$ holds for $k < -(\nu + 1)$, where

$$a_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} F^*(\omega) W_{\text{opt}}(\omega) H(\omega) \exp(-jk\omega) d\omega,$$

$$b_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} F^*(\omega) T_{\text{opt}}(\omega) \exp(-jk\omega) d\omega.$$

⁷Note that, Proposition 5 is the same as [33, Th. 1] that has been derived for hard feedback symbols.

⁸Similar to MIMO channels, $\bar{\delta}_1(F(\omega))$ in (54) is only defined for $\alpha \neq 0$ which represents the rate increment with soft-information. The same holds for Method II.

Proof: As shown in Appendix H, the optimal $\tilde{\mathbf{t}}$ in (90) satisfies

$$\tilde{\mathbf{t}}_{\text{opt}} \mathbf{e}_2 = -\mathbf{e}_1^\dagger.$$

With the definitions of $\mathbf{e}_1, \mathbf{e}_2$ in (49), this is equivalent to

$$\begin{aligned} & \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{\tilde{M}(\omega) |F(\omega)|^2 T_{\text{opt}}(\omega) \phi(\omega)^\dagger}{1 + |F(\omega)|^2} d\omega \\ &= -\frac{\alpha}{2\pi} \int_{-\pi}^{\pi} F(\omega) M(\omega) \phi(\omega)^\dagger d\omega. \end{aligned} \quad (55)$$

On the other hand, with W_{opt} in (50) and $M(\omega), \tilde{M}(\omega)$ defined in (46) and (47), we have

$$\begin{aligned} & \frac{1}{2\pi} \int_{-\pi}^{\pi} (F^*(\omega) W_{\text{opt}}(\omega) H(\omega) - F^*(\omega) T_{\text{opt}} \\ & \quad - (1 + |F(\omega)|^2)) \exp(-jk\omega) d\omega \\ &= \frac{1}{2\pi} \int_{-\pi}^{\pi} \left(\frac{\tilde{M}(\omega) F^*(\omega) T_{\text{opt}}(\omega)}{\alpha} \right. \\ & \quad \left. + (1 + |F(\omega)|^2) M(\omega) \right) \exp(-jk\omega) d\omega. \end{aligned} \quad (56)$$

Transforming (55) and (56) back into matrix forms, we have that (20) and (21) hold. Following the same arguments as in the proof of Proposition 2, $\mathbf{F}^\dagger(\mathbf{W}_{\text{opt}} \mathbf{H} - \mathbf{R}_{\text{opt}})$ is banded within diagonals $[-\nu, K-1]$. Therefore, we have

$$\begin{aligned} & \frac{1}{2\pi} \int_{-\pi}^{\pi} (F^*(\omega) W_{\text{opt}}(\omega) H(\omega) \\ & \quad - F^*(\omega) T_{\text{opt}}(\omega)) \exp(-jk\omega) d\omega = 0 \end{aligned}$$

whenever $k < -(\nu+1)$, which proves Proposition 6. ■

B. METHOD II

The matrices $(\mathbf{V}, \mathbf{R}, \mathbf{G})$ have the same constraints as in Sec. IV-B while the dimensions of these matrices are infinitely large. However, as the shape (a) of \mathbf{R} in Fig. 3 is not meaningful when $N, K \rightarrow \infty$, it is not considered for ISI case. Applying Szegő's eigenvalue distribution theorem to (8), the asymptotic rate of Method II reads

$$\begin{aligned} & \bar{I}(\mathbf{V}(\omega), \mathbf{R}(\omega), \mathbf{G}(\omega)) \\ &= \lim_{K \rightarrow \infty} \frac{1}{K} I_{\text{GMI}}(\mathbf{V}, \mathbf{R}, \mathbf{G}) \\ &= \frac{1}{2\pi} \int_{-\pi}^{\pi} \left(\log(1 + G(\omega)) - G(\omega) - \frac{L_2(\omega)}{1 + G(\omega)} \right) d\omega \\ & \quad + \frac{1}{\pi} \int_{-\pi}^{\pi} \mathcal{R}\{(V(\omega)H(\omega) - \alpha R(\omega))\} d\omega \end{aligned} \quad (57)$$

where

$$\begin{aligned} L_2(\omega) &= |V(\omega)|^2 (N_0 + |H(\omega)|^2) \\ & \quad + \alpha |R(\omega)|^2 - 2\alpha \mathcal{R}\{H(\omega)V(\omega)R^*(\omega)\}. \end{aligned} \quad (58)$$

Define a $2(N_R - \nu_R) \times 1$ vector

$$\psi(\omega) = [\exp(-jN_R\omega), \dots, \exp(-j(\nu_R+1)\omega), \exp(j(\nu_R+1)\omega), \dots, \exp(jN_R\omega)]^T, \quad (59)$$

a $2(N_R - \nu_R) \times 1$ vector ξ_1 , and a $2(N_R - \nu_R) \times 2(N_R - \nu_R)$ Hermitian matrix ξ_2 as

$$\begin{aligned} \xi_1 &= \frac{\alpha}{2\pi} \int_{-\pi}^{\pi} M(\omega) \psi(\omega) d\omega, \\ \xi_2 &= \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{\tilde{M}(\omega) \psi(\omega) \psi(\omega)^\dagger}{1 + G(\omega)} d\omega, \end{aligned} \quad (60)$$

where N_R denotes the tap-length of $R_{\text{opt}}(\omega)$, and $2\nu_R+1$ is the band-size that \mathbf{R} is constrained to zero, and $M(\omega)$ and $\tilde{M}(\omega)$ are in (46) and (47), respectively. Then, we have Proposition 7 with the proof in Appendix I, where we also show that $R(\omega)$ is real and correspondingly, matrix \mathbf{R} has Hermitian symmetry.

Proposition 7: The optimal $V(\omega)$ for (57) is

$$V_{\text{opt}}(\omega) = \frac{H^*(\omega)}{N_0 + |H(\omega)|^2} (1 + G(\omega) + \alpha R_{\text{opt}}(\omega)), \quad (61)$$

and when $0 < \alpha \leq 1$ the optimal $R(\omega)$ reads

$$R_{\text{opt}}(\omega) = -\xi_1^\dagger \xi_2^{-1} \psi(\omega). \quad (62)$$

With the optimal $V(\omega)$ and $R(\omega)$, the asymptotic rate equals

$$\begin{aligned} & \bar{I}(V_{\text{opt}}(\omega), R_{\text{opt}}(\omega), G(\omega)) \\ &= \begin{cases} \bar{I}_2(G(\omega)), & \alpha = 0, \\ \bar{I}_2(G(\omega)) + \bar{\delta}_2(G(\omega)), & 0 < \alpha \leq 1. \end{cases} \end{aligned} \quad (63)$$

The functions $\bar{I}_1(G(\omega))$ and $\bar{\delta}_2(G(\omega))$ are defined as

$$\begin{aligned} \bar{I}_2(G(\omega)) &= 1 + \frac{1}{2\pi} \int_{-\pi}^{\pi} \left(\log(1 + G(\omega)) \right. \\ & \quad \left. + M(\omega)(1 + G(\omega)) \right) d\omega, \end{aligned} \quad (64)$$

$$\bar{\delta}_2(G(\omega)) = -\xi_1^\dagger \xi_2^{-1} \xi_1. \quad (65)$$

When $0 < \alpha \leq 1$, it still needs a gradient based optimization to find the optimal $G(\omega)$ for (64), and the closed-form solution in Theorem 3 is utilized as the starting point. Similar to the MIMO case, the asymptotic rate $\bar{I}(V_{\text{opt}}(\omega), R_{\text{opt}}(\omega), G(\omega))$ is concave with respect to $G(\omega)$, which is shown in Appendix J.

Proposition 8: When $0 < \alpha \leq 1$ it holds that $a_k = b_k$ for $|k| > \nu + \nu_R$, where

$$a_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} V_{\text{opt}}(\omega) H(\omega) \exp(-jk\omega) d\omega, \quad (66)$$

$$b_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} R_{\text{opt}}(\omega) \exp(-jk\omega) d\omega. \quad (67)$$

The connection of the optimal $V(\omega)$ and $R(\omega)$ stated in Proposition 8 is an asymptotic version of Proposition 4, and the proof follows the similar approach as Proposition 6 which is omitted. We show an example in Fig. 5 to illustrate Proposition 8 with Method II and $\nu = \nu_R = 1$. The Proakis-C [57] channel is tested at an SNR of 10 dB and α equals 0.1, 0.4 and 0.8, respectively. As $\nu_R = 1$, b_k as defined in (67) is constrained to zero for $0 \leq k \leq 1$. As can be seen, a_k as defined in (66) equals b_k only for $|k| > 2$, and when $|k| = 2$, a_k and b_k are not identical. This shows that with the optimal $V(\omega)$ and $R(\omega)$, the signal part along the second upper and lower diagonals that is not considered in $G(\omega)$ shall not be

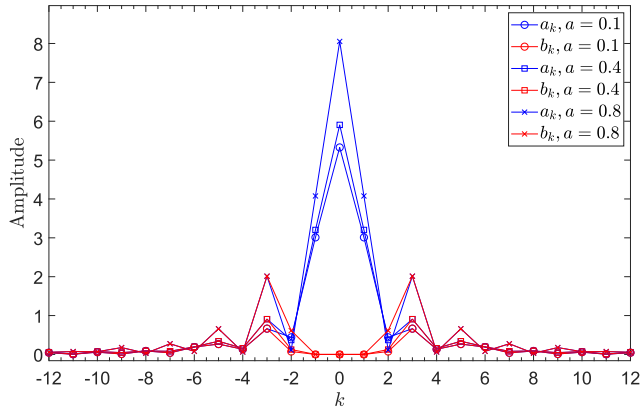


FIGURE 5. Comparison between a_k and b_k for Method II under Proakis-C channel $h = [0.227 \ 0.46 \ 0.688 \ 0.46 \ 0.227]$.

perfectly canceled out. This behavior cannot be seen in [32] which treats LMMSE-PIC for ISI channels due to $\nu = \nu_R = 0$.

C. METHOD III

In Method III, from (39) the asymptotic rate reads

$$\begin{aligned} \bar{I}(G(\omega)) &= \lim_{K \rightarrow \infty} \frac{1}{K} I_{\text{GMI}}(\mathbf{G}) \\ &= 1 + \frac{1}{2\pi} \int_{-\pi}^{\pi} (\log(1 + G(\omega)) + \hat{M}(\omega)(1 + G(\omega))) d\omega \end{aligned} \tag{68}$$

where according to (40) it holds that

$$\begin{aligned} \hat{M}(\omega) &= 2\mathcal{R}\{\alpha \hat{W}(\omega)H(\omega)\hat{C}^*(\omega) + \hat{W}(\omega)H(\omega) - \alpha \hat{C}^*(\omega)\} \\ &\quad - \frac{|\hat{W}(\omega)|^2}{N_0 + |H(\omega)|^2} - \alpha |\hat{C}(\omega)|^2 - 1. \end{aligned}$$

Replacing $M(\omega)$ by $\hat{M}(\omega)$, the optimal $G(\omega)$ and asymptotic rate \bar{I} follow from Theorem 3.

Remark 3: Proposition 8 also holds for Method III with $\nu_R = \nu$, due to the fact that $[(\mathbf{I} + \mathbf{G})(\hat{\mathbf{W}}\mathbf{H} - \hat{\mathbf{C}})]_{2\nu} = \mathbf{0}$.

VI. SNR ASYMPTOTICS

In this section, we analyze asymptotic properties of the elaborated CS demodulators, and show that as N_0 goes to 0 and ∞ , Method II and III are asymptotically equivalent. As Method I is inferior to Method II in GMI sense, we limit the analysis to Method II and III, and start with the analysis for the MIMO case.

The following limits can be verified straightforwardly:

$$\begin{aligned} \lim_{N_0 \rightarrow 0} \mathbf{M}/N_0 &= -(\mathbf{H}^\dagger \mathbf{H})^{-1}, \\ \lim_{N_0 \rightarrow \infty} N_0(\mathbf{I} + \mathbf{M}) &= \mathbf{H}^\dagger \mathbf{H}. \end{aligned} \tag{69}$$

Further, it also holds that

$$\begin{aligned} \lim_{N_0 \rightarrow 0} \tilde{\mathbf{M}} &= \mathbf{P}^2 - \mathbf{P}, \\ \lim_{N_0 \rightarrow \infty} \tilde{\mathbf{M}} &= -\mathbf{P}. \end{aligned} \tag{70}$$

As $\tilde{\mathbf{M}}$ should be invertible from the definition of $\delta_2(\mathbf{G})$ in (26), we restrict that $\mathbf{P} \prec \mathbf{I}$.

Lemma 3: When $N_0 \rightarrow 0$ and ∞ , the optimal \mathbf{G} for (24) in Method II satisfies (18), and the following limits hold,

$$\lim_{N_0 \rightarrow 0} [(N_0(\mathbf{I} + \mathbf{G}_{\text{opt}}))^{-1}]_{\nu} = [(\mathbf{H}^\dagger \mathbf{H})^{-1}]_{\nu}, \tag{71}$$

$$\lim_{N_0 \rightarrow \infty} [N_0 \mathbf{G}_{\text{opt}}]_{\nu} = [\mathbf{H}^\dagger \mathbf{H}]_{\nu}. \tag{72}$$

Proof: When $\mathbf{P} = \mathbf{0}$, from Theorem 2 the optimal \mathbf{G} for (24) satisfies (18). From (69), when $N_0 \rightarrow 0$ it holds that $\mathbf{M} \rightarrow \mathbf{0}$, and when $N_0 \rightarrow \infty$ we have $\mathbf{M} \rightarrow -\mathbf{I}$. Therefore, by the definition of Ω it holds that

$$\lim_{N_0 \rightarrow 0, \infty} \mathbf{d} = \Omega \text{vec}(\mathbf{M}\mathbf{P}) = \mathbf{0}.$$

This implies that the gradient $d_G(\delta_2)$ (showing in (85) in Appendix F) converges to zero. Hence the differentials of $I_{\text{GMI}}(\mathbf{V}_{\text{opt}}, \mathbf{R}_{\text{opt}}, \mathbf{G})$ in (24) with $\mathbf{P} \neq \mathbf{0}$ converges to the those with $\mathbf{P} = \mathbf{0}$. From (18) and (69), the limit (71) follows.

Next, since

$$\begin{aligned} \lim_{N_0 \rightarrow \infty} [N_0 (\mathbf{I} - (\mathbf{I} + \mathbf{G}_{\text{opt}})^{-1})]_{\nu} &= \lim_{N_0 \rightarrow \infty} [N_0(\mathbf{I} + \mathbf{M})]_{\nu} \\ &= [\mathbf{H}^\dagger \mathbf{H}]_{\nu}, \end{aligned} \tag{73}$$

which shows that $\mathbf{I} - (\mathbf{I} + \mathbf{G}_{\text{opt}})^{-1} \rightarrow \mathbf{0}^9$ as $N_0 \rightarrow \infty$. By the matrix inversion lemma [58], $\mathbf{I} - (\mathbf{I} + \mathbf{G}_{\text{opt}})^{-1} \rightarrow \mathbf{G}_{\text{opt}}$ as $N_0 \rightarrow \infty$, and combining it with (73) proves (72). ■

Lemma 4: In Method II, with the optimal \mathbf{G} , when $N_0 \rightarrow 0$ the GMI increment $\delta_2(\mathbf{G})$ in (26) converges to zero with speed $\mathcal{O}(1/N_0)$,¹⁰ and when $N_0 \rightarrow \infty$ it converges to zero with speed $\mathcal{O}(N_0^2)$, respectively.

Proof: As $N_0 \rightarrow 0$, from (69) we have

$$\begin{aligned} \lim_{N_0 \rightarrow 0} \mathbf{d}/N_0 &= \lim_{N_0 \rightarrow 0} \Omega \text{vec}(\mathbf{M}\mathbf{P}/N_0) \\ &= -\Omega \text{vec}((\mathbf{H}^\dagger \mathbf{H})^{-1} \mathbf{P}). \end{aligned}$$

Based on (70) and Lemma 3, the below equalities hold,

$$\begin{aligned} \delta_2(\mathbf{G}_{\text{opt}}) &= N_0 \frac{\mathbf{d}^\dagger}{N_0} \left(\Omega (\tilde{\mathbf{M}}^* \otimes \frac{(\mathbf{I} + \mathbf{G}_{\text{opt}})^{-1}}{N_0}) \Omega^T \right)^{-1} \frac{\mathbf{d}}{N_0} \\ &= \mathcal{O}(N_0). \end{aligned}$$

On the other hand, as $N_0 \rightarrow \infty$, and by the definition of Ω , from (69) we also have

$$\begin{aligned} \lim_{N_0 \rightarrow \infty} N_0 \mathbf{d} &= \lim_{N_0 \rightarrow \infty} \Omega \text{vec}(N_0 \mathbf{M}\mathbf{P}) \\ &= \lim_{N_0 \rightarrow \infty} \Omega \text{vec}(N_0(\mathbf{I} + \mathbf{M})\mathbf{P}) \\ &= \Omega \text{vec}(\mathbf{H}^\dagger \mathbf{H}\mathbf{P}). \end{aligned}$$

⁹A matrix $\mathbf{A} \rightarrow \mathbf{B}$ or a vector $\mathbf{a} \rightarrow \mathbf{b}$ means the non-zero elements of $\mathbf{A} - \mathbf{B}$ or $\mathbf{a} - \mathbf{b}$ converges to zero.

¹⁰Two scalars A and B as functions of a variable n converging to each other with speed $\mathcal{O}(n)$ means that there exists a constant C such that $\lim_{n \rightarrow \infty} n|A - B| < C$.

Again utilizing (70) and Lemma 3, the below equalities hold,

$$\begin{aligned}\delta_2(\mathbf{G}_{\text{opt}}) &= \frac{1}{N_0^2} (N_0 \mathbf{d}^\dagger) (\boldsymbol{\Omega} (\tilde{\mathbf{M}}^* \otimes (\mathbf{I} + \mathbf{G}_{\text{opt}})^{-1}) \boldsymbol{\Omega}^T)^{-1} (N_0 \mathbf{d}) \\ &= \mathcal{O}(1/N_0^2).\end{aligned}$$

Therefore, Lemma 4 holds. \blacksquare

Lemma 5: When both $N_0 \rightarrow 0$ and ∞ , the optimal GMI in Method III is independent of \mathbf{P} and converges to the optimal GMI with $\mathbf{P} = \mathbf{0}$. Further, (71) and (72) also hold for Method III.

The proof is given in Appendix K. Combining Lemmas 3-5, and using the fact that Method III and Method II are equivalent with $\mathbf{P} = \mathbf{0}$, we have the following Theorem 4.

Theorem 4: Assume that $\mathbf{P} \prec \mathbf{I}$, when $N_0 \rightarrow 0$ and ∞ , the optimal GMI in Method III converges to the optimal GMI in Method III with $\mathbf{P} = \mathbf{0}$. Further, the optimal GMI in Method II also converges to the same, with speed $\mathcal{O}(1/N_0)$ as SNR increase, and $\mathcal{O}(N_0^2)$ as SNR decreases, respectively. The optimal \mathbf{G} for both methods has the limits (71) and (72).

From Theorem 4 we know that, except for the case where one of the elements in the diagonal matrix \mathbf{P} is 1, the soft-feedback information becomes asymptotically insignificant for the design of the CS parameters. The reason is that, when $N_0 \rightarrow 0$, $\hat{\mathbf{x}}$ is overwhelmed by the noise, while when $N_0 \rightarrow \infty$, the optimal front-end filter will null out $\hat{\mathbf{x}}$ since the receiver can perfectly reconstruct the transmitted symbols without the side-information.

Remark 4: When $N_0 \rightarrow 0$ and ∞ , the optimal CS demodulator is the EZF demodulator defined in Example 1, and the TMF defined in Example 2, respectively.

With ISI channels, as the same constraint $\mathbf{P} = \alpha \mathbf{I} \prec \mathbf{I}$ shall hold, we make the restriction that $\alpha < 1$. The asymptotic properties for ISI channels are presented in Corollary 1, which is an asymptotic version of Theorem 4 when the channel matrix \mathbf{H} and CS parameters are band-shaped Toeplitz matrices with infinite dimensions. The detailed proof follows the same analysis as for the MIMO case and is omitted.

Corollary 1: Assume that $0 \leq \alpha < 1$, when $N_0 \rightarrow 0$ and ∞ , the optimal GMI in Method III converges to the optimal GMI in Method III with $\alpha = 0$. Further, the optimal GMI in Method II also converges to the same, with speed $\mathcal{O}(1/N_0)$ as SNR increase, and $\mathcal{O}(N_0^2)$ as SNR decreases, respectively. The optimal \mathbf{G} for both methods has the following asymptotic properties hold for $|k| \leq \nu$:

$$\begin{aligned}\lim_{N_0 \rightarrow 0} \int_{-\pi}^{\pi} \frac{1}{N_0(1 + G_{\text{opt}}(\omega))} \exp(-jk\omega) d\omega \\ &= \int_{-\pi}^{\pi} \frac{1}{|H(\omega)|^2} \exp(-jk\omega) d\omega, \\ \lim_{N_0 \rightarrow \infty} \int_{-\pi}^{\pi} N_0 G_{\text{opt}}(\omega) \exp(-jk\omega) d\omega \\ &= \int_{-\pi}^{\pi} |H(\omega)|^2 \exp(-jk\omega) d\omega.\end{aligned}$$

VII. NUMERICAL RESULTS

In this section, we provide numerical results to elaborate the behaviors of CS demodulators in an iterative detection and decoding receiver designs. In the considered MIMO channels, all channel elements are assumed to be independent and identically distributed (i.i.d.) complex-valued Gaussian variables with zero-means and unit-variance. Further, the transmit-power at each transmit-antenna is normalized to unity. For the ISI case, we use the typical Proakis-C channel as in Fig. 5.

A. GMI EVALUATION

We first evaluate the GMI under 5×5 MIMO channels with memory-size $\nu = 1$ for all CS demodulators. We simulate 10000 channel realizations for each SNR point. The GMIs are compared with those of the static CS demodulator [18], which is equivalent to Method II with $\mathbf{P} = \mathbf{0}$. The channel capacity is also presented for comparisons. The results of GMI are shown in Fig. 6. As the quality of soft-information improves beyond $\mathbf{P} = \mathbf{0}$, Method II with $\nu_R = 0$ performs the best among all CS demodulators, as it has the most degrees of freedom (DoF) in designing \mathbf{R} . Method II with $\nu_R = \nu$ is the worst among Method I and II, while Method I is slightly worse than Method II with \mathbf{R} of shape (a) in Fig. 3. This is because although the IC matrix \mathbf{R} is of shape (a) in both cases, \mathbf{R} in Method II is more general than in Method I, since the latter one is constrained to $\mathbf{R} = \mathbf{F}^\dagger \mathbf{T}$. Further, the GMI of Method III is inferior to Method II as expected.

The results show consistent GMI increments for all CS demodulators when the feedback quality improves. When \mathbf{P} increases from $\mathbf{P} = \mathbf{0}$ to the ideal case $\mathbf{P} = \mathbf{I}$, the channel capacity becomes inferior to the GMI as the pair $(\mathbf{x}, \hat{\mathbf{y}})$ is superior to (\mathbf{x}, \mathbf{y}) for information-transfer.

B. SNR ASYMPTOTIC OF THE GMI

Next, we evaluate the asymptotic properties described in Theorem 4 under 5×5 MIMO channels. As shown in Fig. 7, the GMIs of both Method II and III converge to that of Method III with $\mathbf{P} = \mathbf{0}$. In addition, the GMI converges to the EZF in Example 1 at high SNR, and the TMF in Example 2 at low SNR, respectively, which are aligned with the limits in Theorem 4.

C. EXIT CHARTS OF CS DEMODULATORS

In order to predict the dynamics of iterative receivers, we extrinsic information transfer (EXIT) charts [41], [59] to analyse iterative receivers. The demodulators and the decoder measure the output extrinsic information I_E based on a sequence of observations \mathbf{y} and *a priori* mutual information I_A .

In Fig. 8, we evaluate the EXIT charts for CS demodulators under 4×6 MIMO channels with $\nu = 2$ for \mathbf{F} and \mathbf{G} at an SNR of 10 dB. With Method II, we test different values of ν_R . As can be seen, when $\nu_R > \nu$, the demodulation-performance is inferior to $\nu_R \leq \nu$. This is because, the

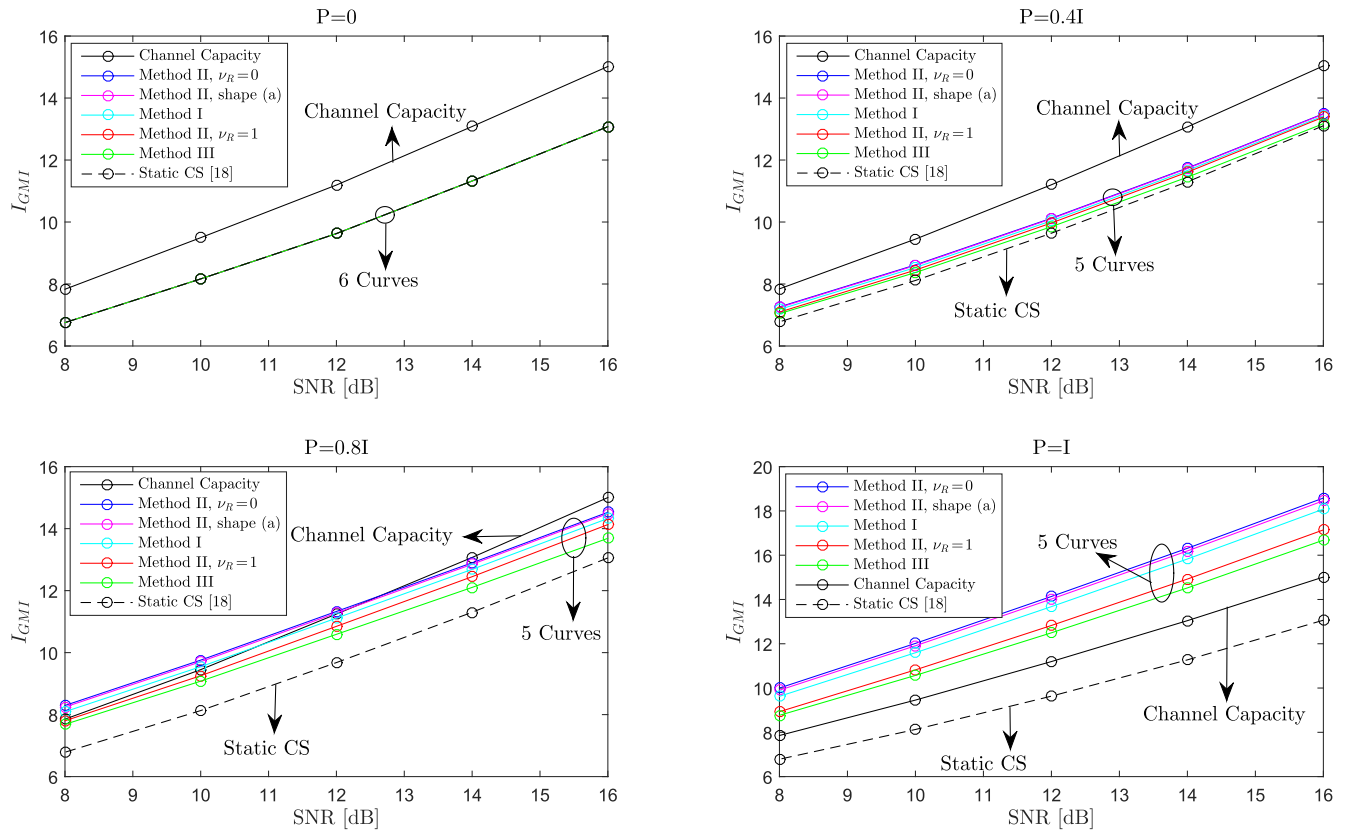


FIGURE 6. GMI of the CS demodulators under 5×5 MIMO i.i.d. complex-valued Gaussian channels with $\nu = 1$.

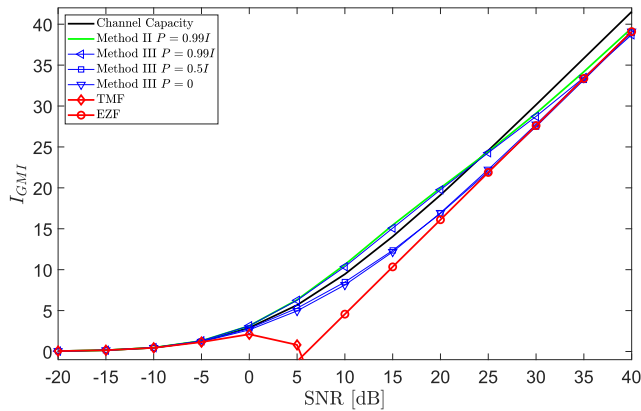


FIGURE 7. SNR asymptotic under 5×5 MIMO i.i.d. complex-valued Gaussian channels with $\nu = 1$.

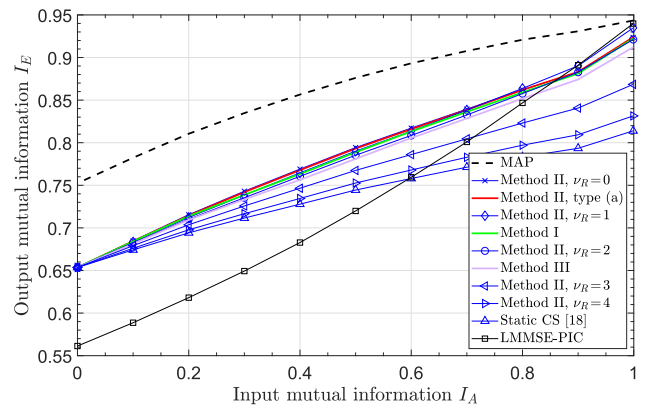


FIGURE 8. EXIT charts under 4×6 i.i.d. complex-valued Gaussian MIMO channels with $\nu = 2$ and different values of ν_R .

interference outside memory-size ν and inside memory-size ν_R is neither considered in the IC nor in the BCJR. Moreover, with $\nu_R \leq \nu$ the CS demodulators with Method II performs quite close to each other as well as Method I and III. For Method II with $\nu_R < \nu$, the interference inside memory-size ν and outside ν_R are considered both in the IC and BCJR processes. However, an interesting observation is that, with a larger input I_A , Method II with $\nu_R = 0$ is inferior to $\nu_R = 1$ and $\nu_R = 2$. Therefore, a conservative and robust approach

with Method II is to set $\nu_R = \nu$ such that the interference is either removed in IC or dealt with in BCJR, to get rid of potential error-propagation caused by redundant processing of the same signal part.

In Fig. 9, we show the iterative detection and decoding trajectories for CS demodulators under Proakis-C channel with $\nu = 2$ and at an SNR of 10 dB. We use a $(7, 5)$ -convolutional code [24] with a coded block-length $K = 2004$, and a random permutation is applied to the coded-bits. As can

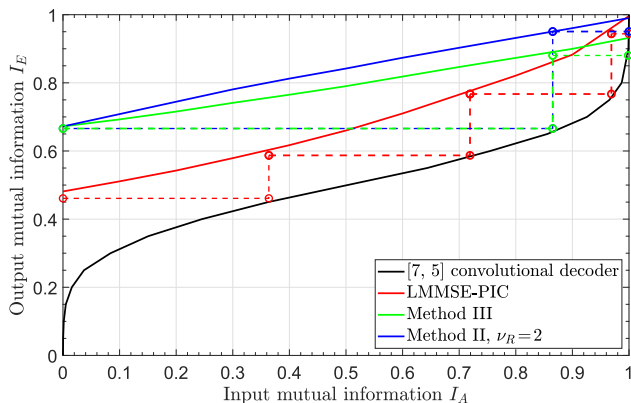


FIGURE 9. Iterative detection and decoding trajectories under Proakis-C channel at an SNR of 10 dB. The outer code is a (7, 5)-convolutional code with generator polynomials $g_0(D)=1+D^2$ and $g_1(D)=1+D+D^2$. The black curve represents the decoder and the dashed lines are the iterative detection and decoding trajectories for LMMSE-PIC, Method III and II, respectively.

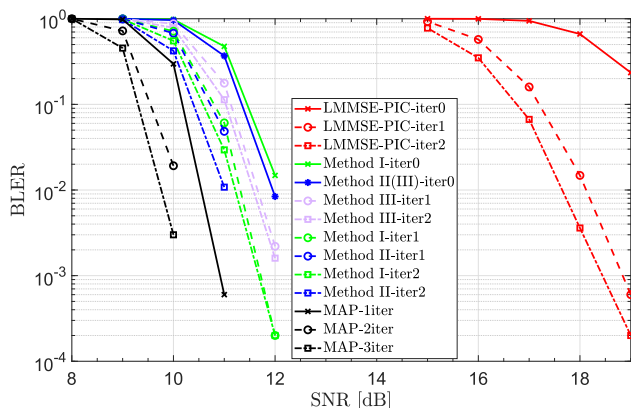


FIGURE 10. BLER performance of the LMMSE-PIC, Method I, II, and III, and MAP under Proakis-C channel with QPSK modulation.

be seen, the CS demodulators with Method II and III are superior to the LMMSE-PIC demodulator, and the iterative detection and decoding trajectories are well aligned with the EXIT charts.

D. LINK PERFORMANCE

We next turn to link-level simulations with turbo codes [60] where the outer-decoder uses 8 internal iterations. A single code-block over all transmit symbols is used. At each SNR point 20000 data-blocks are simulated and the block-error-ratio (BLER) is measured. In all simulations, maximal three global iterations are used between the demodulators. The tap-length of both the front-end filter and IC filter are set to $8L$, and $\nu_R = \nu$ for Method II.

In Fig. 10, we evaluate the BLER under Proakis-C channel with QPSK symbols and $\nu = 2$ for all CS demodulators. A (1064, 1600) turbo code is used. Note that, at the first iteration when there is no soft-information, Method II and III are overlapped with each other. With CS demodulators, the gap to the MAP demodulator is less than 1 dB,

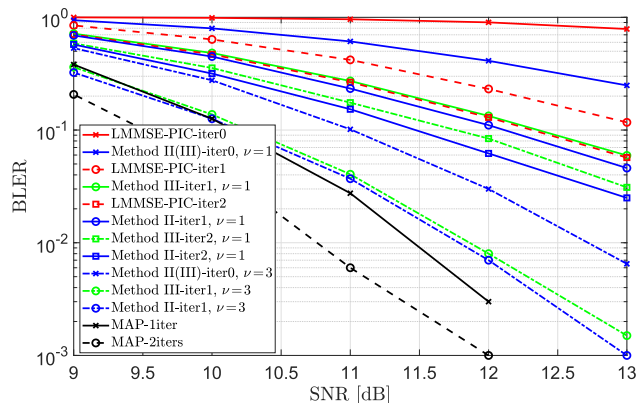


FIGURE 11. BLER performance of the LMMSE-PIC, Method II, and III, and MAP under 4×6 MIMO channels with QPSK modulation.

while the LMMSE-PIC has a gap to the MAP that is up to 10 dB. Moreover, Method II performs slightly better than Method I, and Method III is slightly inferior to both methods. However, Method III has the advantage of less computational complexity than the others since all parameters are in closed-forms.

In Fig. 11, we evaluate the BLER under 4×6 MIMO channels with QPSK symbols and $\nu = 3$ for all CS demodulators. A (1064,1800)-turbo code is used. As $N < K$, the LMMSE-PIC fails [61] at the first iteration due to the lack of sufficient receive diversity. However, the CS demodulators with $\nu = 3$ significantly improve the performance and with less than 1 dB gap to the MAP at 10% BLER. Further, the CS demodulators with $\nu = 1$ after three iterations is less than 2 dB away from the MAP. More interestingly, with less complexity, Method III performs close to Method II.

Finally we remark that for the sake of complexity savings, the CS parameters do not need to be updated through all iterations. Once the feedback information quality is good enough and the changes in \mathbf{P} or α are small, the CS parameters can be kept unchanged in successive iterations.

VIII. SUMMARY

In this paper we have considered the design of CS demodulators for linear vector channels that use a trellis-representation of the received signal, in combination with interference cancellation (IC) of the signal part that is not appropriately modeled by the trellis. In order to reach a trellis-representation, a linear filter is applied as front-end. It is an extension of the well studied CS demodulators to iterative receivers and a generalization of the LMMSE-PIC demodulator to cooperate with trellis-search in turbo equalization.

We have analyzed the properties of three different methods for designing optimal CS demodulators, as all of them may come across as natural CS demodulators. In the used framework, there are three parameters that need to be optimized: the front-end filter, the IC filter, and the target response. Based on maximizing the GMI, the first two are solved in closed-forms, while the third one needs numerical optimizations and simple

gradient based approaches have been shown to perform well. In particular, in Method III we have constructed all three CS parameters explicitly at a cost of small GMI losses.

Numerical results have been provided to illustrate the behaviors of the proposed CS demodulators. In general, Method II based on the Ungerboeck model is superior to Method I based on the Forney model. Method II has the advantage over Method I that the optimization procedure is concave. Further, the suboptimal Method III performs close to Method I and II. An interesting result is that the IC of the CS demodulators should not perfectly cancel the effective channel outside the memory-size ν , a property that cannot be seen in the LMMSE-PIC demodulator as $\nu = 0$. Moreover, we have also analyzed asymptotic properties of the CS demodulators and showed that Method III converges to Method II asymptotically when the noise power goes either zero or infinity.

**APPENDIX A
DERIVATION OF THE GMI**

By making the eigenvalue decomposition $\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\dagger = \mathbf{G}$ and letting $\mathbf{s} = \mathbf{Q}^\dagger\mathbf{x}$. As \mathbf{x} is assumed to be zero mean complex Gaussian random vector with covariance matrix \mathbf{I} , we can write $\tilde{p}(\mathbf{y}|\mathbf{x})$ in (4) as

$$\tilde{p}(\mathbf{y}|\mathbf{x}) = \exp(2\mathcal{R}\{\mathbf{s}^\dagger\mathbf{d}\} - \mathbf{s}^\dagger\mathbf{\Lambda}\mathbf{s}), \tag{74}$$

where $\mathbf{d} = \mathbf{Q}^\dagger(\mathbf{V}\mathbf{y} - \mathbf{R}\hat{\mathbf{x}})$. We can now evaluate

$$\begin{aligned} \tilde{p}(\mathbf{y}) &= \int \tilde{p}(\mathbf{y}|\mathbf{x})p(\mathbf{x})d\mathbf{x} \\ &= \frac{1}{\pi^K} \int \exp(2\mathcal{R}\{\mathbf{s}^\dagger\mathbf{d}\} - \mathbf{s}^\dagger\mathbf{\Lambda}\mathbf{s}) \exp(-\mathbf{s}^\dagger\mathbf{s})d\mathbf{s} \\ &= \prod_{k=1}^N \frac{1}{1+\lambda_k} \exp\left(\frac{|d_k|^2}{1+\lambda_k}\right), \end{aligned}$$

where λ_k is the k th diagonal element of $\mathbf{\Lambda}$ and d_k is the k th entry of \mathbf{d} . Taking the average over \mathbf{y} gives

$$-\mathbb{E}_{p(\mathbf{y})}[\log(\tilde{p}(\mathbf{y}))] = \log(\det(\mathbf{I} + \mathbf{G})) - \text{Tr}(\mathbf{L}(\mathbf{I} + \mathbf{G})^{-1})$$

where the matrix $\mathbf{L} = \mathbb{E}[\mathbf{Q}\mathbf{d}\mathbf{d}^\dagger\mathbf{Q}^\dagger]$ is given by

$$\mathbf{L} = \mathbf{V}(N_0\mathbf{I} + \mathbf{H}\mathbf{H}^\dagger)\mathbf{V}^\dagger - \mathbf{V}\mathbf{H}\mathbf{P}\mathbf{R}^\dagger - \mathbf{R}\mathbf{P}\mathbf{H}^\dagger\mathbf{V}^\dagger + \mathbf{R}\mathbf{P}\mathbf{R}^\dagger.$$

On the other hand, we have

$$-\mathbb{E}_{p(\mathbf{y},\mathbf{x})}[\log(\tilde{p}(\mathbf{y}|\mathbf{x}))] = \text{Tr}(\mathbf{G}) - 2\mathcal{R}\{\text{Tr}(\mathbf{V}\mathbf{H} - \mathbf{R}\mathbf{P})\}.$$

Combining the two expectations, the GMI reads

$$\begin{aligned} I_{\text{GMI}}(\mathbf{V}, \mathbf{R}, \mathbf{G}) &= \log(\det(\mathbf{I} + \mathbf{G})) - \text{Tr}(\mathbf{L}(\mathbf{I} + \mathbf{G})^{-1}) - \text{Tr}(\mathbf{G}) \\ &\quad + 2\mathcal{R}\{\text{Tr}(\mathbf{V}\mathbf{H} - \mathbf{R}\mathbf{P})\} \\ &= \log(\det(\mathbf{I} + \mathbf{G})) - \text{Tr}(\mathbf{G}) + 2\mathcal{R}\{\text{Tr}(\mathbf{V}\mathbf{H} - \mathbf{R}\mathbf{P})\} \\ &\quad - \text{Tr}((\mathbf{I} + \mathbf{G})^{-1}(\mathbf{V}[\mathbf{H}\mathbf{H}^\dagger + N_0\mathbf{I}]\mathbf{V}^\dagger \\ &\quad - 2\mathcal{R}\{\mathbf{V}\mathbf{H}\mathbf{P}\mathbf{R}^\dagger\} + \mathbf{R}\mathbf{P}\mathbf{R}^\dagger)). \end{aligned}$$

**APPENDIX B
THE PROOF OF PROPOSITION 1**

As the formula of GMI in (9) is quadratic in \mathbf{W} and no constraints apply to \mathbf{W} , taking the gradient of $I_{\text{GMI}}(\mathbf{W}, \mathbf{T}, \mathbf{F})$ with respect to \mathbf{W} and setting it to zero, the optimal \mathbf{W} is then in (12). Inserting \mathbf{W}_{opt} back into (9) yields, after some manipulations,

$$\begin{aligned} I_{\text{GMI}}(\mathbf{W}_{\text{opt}}, \mathbf{T}, \mathbf{F}) &= K + \log(\det(\mathbf{I} + \mathbf{F}^\dagger\mathbf{F})) + \text{Tr}(\mathbf{T}^\dagger\mathbf{F}(\mathbf{I} + \mathbf{F}^\dagger\mathbf{F})^{-1}\mathbf{F}^\dagger\mathbf{T}\tilde{\mathbf{M}}) \\ &\quad + \text{Tr}(\mathbf{M}(\mathbf{I} + \mathbf{F}^\dagger\mathbf{F})) + 2\mathcal{R}\{\text{Tr}(\mathbf{P}\mathbf{M}\mathbf{F}^\dagger\mathbf{T})\}. \end{aligned} \tag{75}$$

where \mathbf{M} and $\tilde{\mathbf{M}}$ are defined in (10) and (11), respectively. When $\mathbf{P} = \mathbf{0}$, (75) equals

$$I_1(\mathbf{F}) = K + \log(\det(\mathbf{I} + \mathbf{F}^\dagger\mathbf{F})) + \text{Tr}(\mathbf{M}(\mathbf{I} + \mathbf{F}^\dagger\mathbf{F})).$$

In this case, there is no soft-information available and the matrix \mathbf{T} is not included in the formula. When $\mathbf{P} \neq \mathbf{0}$, the terms of I_{GMI} in (75) related to \mathbf{T} are

$$f(\mathbf{T}) = \text{Tr}(\mathbf{T}^\dagger\mathbf{F}(\mathbf{I} + \mathbf{F}^\dagger\mathbf{F})^{-1}\mathbf{F}^\dagger\mathbf{T}\tilde{\mathbf{M}}) + 2\mathcal{R}\{\text{Tr}(\mathbf{P}\mathbf{M}\mathbf{F}^\dagger\mathbf{T})\}.$$

Let \mathbf{t}_k denote the k th column of \mathbf{T} , but with all elements in rows $[k, \min(k + \nu, K - 1)]$ removed, and define the column vector $\mathbf{t} = [\mathbf{t}_0^T, \mathbf{t}_1^T, \dots, \mathbf{t}_{K-1}^T]^T$. Then, by the definition of indication matrix $\mathbf{\Omega}$, we have

$$\mathbf{t} = \mathbf{\Omega}\text{vec}(\mathbf{T}).$$

Similarly, let \mathbf{z}_k denote the k th column of the matrix $\mathbf{F}\mathbf{M}\mathbf{P}$ with all elements in rows $[k, \min(k + \nu, K - 1)]$ removed, and define a row vector $\mathbf{z} = [\mathbf{z}_0^T, \mathbf{z}_1^T, \dots, \mathbf{z}_{K-1}^T]^T$. Then, we have

$$\mathbf{z} = \mathbf{\Omega}\text{vec}(\mathbf{F}\mathbf{M}\mathbf{P}) = \mathbf{\Omega}((\mathbf{P}\mathbf{M}^*) \otimes \mathbf{I}_K)\text{vec}(\mathbf{F}).$$

Finally, defining a Hermitian matrix

$$\hat{\mathbf{B}}_1 = \mathbf{\Omega}(\tilde{\mathbf{M}}^* \otimes (\mathbf{F}(\mathbf{I} + \mathbf{F}^\dagger\mathbf{F})^{-1}\mathbf{F}^\dagger))\mathbf{\Omega}^T,$$

and with that we can rewrite $f(\mathbf{T})$ as

$$f(\mathbf{T}) = \mathbf{t}^\dagger \hat{\mathbf{B}}_1 \mathbf{t} + 2\mathcal{R}\{\mathbf{z}^\dagger \mathbf{t}\}.$$

Taking the gradient of $f(\mathbf{T})$ with respect to \mathbf{t} and setting it to zero yields

$$\mathbf{t}_{\text{opt}} = -\hat{\mathbf{B}}_1^{-1} \mathbf{z}. \tag{76}$$

Transferring \mathbf{t}_{opt} back into \mathbf{T}_{opt} gives the optimal \mathbf{T} in (13), and inserting it back into $f(\mathbf{T})$ gives

$$f(\mathbf{T}_{\text{opt}}) = -\mathbf{z}^\dagger \hat{\mathbf{B}}_1^{-1} \mathbf{z}.$$

Thus, with the optimal \mathbf{W} and \mathbf{T} and when $\mathbf{P} \neq \mathbf{0}$ the GMI equals

$$\begin{aligned} I_{\text{GMI}}(\mathbf{W}_{\text{opt}}, \mathbf{T}_{\text{opt}}, \mathbf{F}) &= K + \log(\det(\mathbf{I} + \mathbf{F}^\dagger\mathbf{F})) + \text{Tr}(\mathbf{M}(\mathbf{I} + \mathbf{F}^\dagger\mathbf{F})) \\ &\quad - \text{vec}(\mathbf{F})^\dagger \mathbf{D}^\dagger (\mathbf{\Omega}(\tilde{\mathbf{M}}^* \otimes (\mathbf{F}(\mathbf{I} + \mathbf{F}^\dagger\mathbf{F})^{-1}\mathbf{F}^\dagger))\mathbf{\Omega}^T)^{-1} \mathbf{D}\text{vec}(\mathbf{F}). \end{aligned}$$

where

$$\mathbf{D} = \mathbf{\Omega}((\mathbf{P}\mathbf{M}^*) \otimes \mathbf{I}_K).$$

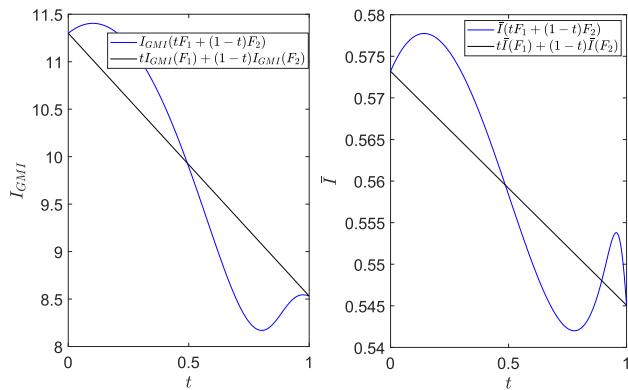


FIGURE 12. Non-concaveness of Method I under 5 × 5 MIMO channel (left figure) and Proakis-C ISI channel (right figure).

**APPENDIX C
NON-CONCAVITY EXAMPLES OF METHOD I**

We give two examples to demonstrate the non-concavity of Method I for both MIMO and ISI cases by assuming that $\mathbf{P} = \mathbf{I}$ and $\alpha = 1$, respectively. The memory-size is $\nu = 1$ and the noise power N_0 equals 1 in both cases. A 5 × 5 MIMO channel and the Proakis-C ISI channel are used.

Example 4: MIMO case:

$$\mathbf{H} = \begin{bmatrix} 2 & 0 & -3 & 5 & 4 \\ -5 & 2 & -1 & 0 & 2 \\ 2 & -4 & 3 & 3 & 3 \\ -1 & -5 & -4 & 1 & 2 \\ 0 & -2 & 0 & 5 & 5 \end{bmatrix},$$

$$\mathbf{F}_1 = \begin{bmatrix} 4.94 & 4.45 & 0 & 0 & 0 \\ 0 & 0.21 & 3.85 & 0 & 0 \\ 0 & 0 & 5.56 & 1.76 & 0 \\ 0 & 0 & 0 & 0.61 & 7.10 \\ 0 & 0 & 0 & 0 & 2.79 \end{bmatrix},$$

$$\mathbf{F}_2 = \begin{bmatrix} 2.03 & 6.17 & 0 & 0 & 0 \\ 0 & 5.22 & 3.56 & 0 & 0 \\ 0 & 0 & 7.43 & 0.73 & 0 \\ 0 & 0 & 0 & 4.98 & 4.32 \\ 0 & 0 & 0 & 0 & 10.11 \end{bmatrix}.$$

Example 5: ISI case:

$$\mathbf{h} = [0.227 \ 0.460 \ 0.688 \ 0.460 \ 0.227],$$

$$\mathbf{f}_1 = [0.1606 \ 0.9009], \mathbf{f}_2 = [0.2230 \ 0.2035].$$

The $I_{\text{GMI}}(\mathbf{W}_{\text{opt}}, \mathbf{T}_{\text{opt}}, \mathbf{F})$ in (14) as a function \mathbf{F} is plotted on the left of Fig. 12, while the $\bar{I}(\mathbf{W}_{\text{opt}}(\omega), \mathbf{T}_{\text{opt}}(\omega), \mathbf{F}(\omega))$ in (52) as a function of $\mathbf{F}(\omega)$ is plotted on the right. If $I_{\text{GMI}}(\mathbf{W}_{\text{opt}}, \mathbf{T}_{\text{opt}}, \mathbf{F})$ and $\bar{I}(\mathbf{W}_{\text{opt}}(\omega), \mathbf{T}_{\text{opt}}(\omega), \mathbf{F}(\omega))$ are concave or convex, the blue curves lie above or below the black curves, which clearly does not hold in the examples.

**APPENDIX D
THE GRADIENT IN METHOD I FOR THE MIMO CASE**

In this section we derive the first order differential of the GMI given in (14) with respect to \mathbf{F} . In order to utilize the differential with respect to a matrix, we use the α -differential

as defined in [62]. Assume a matrix $\mathbf{Y}_{N,K}$ with dimension $N \times K$ and a matrix $\mathbf{X}_{M,S}$ with dimension $M \times S$, define $d_{\mathbf{X}}\mathbf{Y}$ as the α -differential of \mathbf{Y} with respect to \mathbf{X} . Further, defining y_ℓ and x_ℓ as $[y_1, y_2, \dots, y_{NK}] = \text{vec}(\mathbf{Y})^T$ and $[x_1, x_2, \dots, x_{MS}] = \text{vec}(\mathbf{X})^T$, the α -differential $d_{\mathbf{X}}\mathbf{Y}$ is

$$d_{\mathbf{X}}\mathbf{Y} = \frac{\partial \text{vec}(\mathbf{Y})}{\partial \text{vec}(\mathbf{X})^T} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_{MS}} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_2}{\partial x_{MS}} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial y_{NK}}{\partial x_1} & \frac{\partial y_{NK}}{\partial x_2} & \dots & \frac{\partial y_{NK}}{\partial x_{MS}} \end{bmatrix}.$$

The reason for adopting the α -differential is because it keeps both the chain rule and the product rule. We introduce an $NK \times NK$ permutation matrix $\mathbf{Z}_{N,K}$, which satisfies the condition $\text{vec}(\mathbf{Y}^T) = \mathbf{Z}_{N,K}\text{vec}(\mathbf{Y})$. It is easy to verify that $\mathbf{Z}_{N,K}^{-1} = \mathbf{Z}_{K,N}$. In addition, when $N = 1$ or $K = 1$, \mathbf{Y} is a vector and it holds that $\text{vec}(\mathbf{Y}^T) = \text{vec}(\mathbf{Y})$, hence, we have $\mathbf{Z}_{N,1} = \mathbf{I}_N$ and $\mathbf{Z}_{1,K} = \mathbf{I}_K$. Furthermore, by definition it holds that $d_{\mathbf{F}}(\mathbf{F}) = d_{\mathbf{F}}(\text{vec}(\mathbf{F})) = \mathbf{I}$, and $d_{\mathbf{F}}(\mathbf{F}^\dagger) = d_{\mathbf{F}}(\text{vec}(\mathbf{F}^\dagger)) = \mathbf{0}$.

We start by reviewing a few properties [62], [63] of α -differential below that will be used later, where both matrices \mathbf{X} and \mathbf{Y} are functions of \mathbf{F} and the dimensions are specified by subscripts associated to them:

$$d_{\mathbf{F}}(\mathbf{X}_{K,K}^{-1}) = -(\mathbf{X}_{K,K}^{-T} \otimes \mathbf{X}_{K,K}^{-1})d_{\mathbf{F}}\mathbf{X}_{K,K}$$

$$d_{\mathbf{F}}(\mathbf{Y}_{N,K}\mathbf{X}_{K,S}) = (\mathbf{X}_{K,S}^T \otimes \mathbf{I}_N)d_{\mathbf{F}}\mathbf{Y}_{N,K} + (\mathbf{I}_S \otimes \mathbf{Y}_{N,K})d_{\mathbf{F}}\mathbf{X}_{K,S}$$

$$d_{\mathbf{F}}(\log(\det(\mathbf{X}_{K,K}))) = \text{vec}(\mathbf{X}_{K,K}^{-T})^T d_{\mathbf{F}}\mathbf{X}_{K,K}$$

$$d_{\mathbf{F}}(\mathbf{Y}_{N,K} \otimes \mathbf{X}_{M,S}) = (\mathbf{I}_K \otimes \mathbf{Z}_{S,N} \otimes \mathbf{I}_M) \times (\mathbf{I}_{NK} \otimes \text{vec}(\mathbf{X}))d_{\mathbf{F}}\mathbf{Y}_{N,K} + (\mathbf{I}_K \otimes \mathbf{Z}_{S,N} \otimes \mathbf{I}_M) \times (\text{vec}(\mathbf{Y}) \otimes \mathbf{I}_{MS})d_{\mathbf{F}}\mathbf{X}_{M,S}.$$

The α -differential of $I_1(\mathbf{F})$ with respect to \mathbf{F} is

$$d_{\mathbf{F}}(I_1) = \text{vec}((\mathbf{I} + \mathbf{F}^\dagger\mathbf{F})^{-T})^T (\mathbf{I}_K \otimes \mathbf{F}^\dagger) + \text{vec}(\mathbf{F}^*\mathbf{M}^T)^T = \text{vec}(\mathbf{F}\mathbf{M} + \mathbf{F}(\mathbf{I} + \mathbf{F}^\dagger\mathbf{F})^{-1})^\dagger. \quad (77)$$

Defining a $K \times K$ matrix $\mathbf{B} = \mathbf{F}(\mathbf{I} + \mathbf{F}^\dagger\mathbf{F})^{-1}\mathbf{F}^\dagger$ and an $S \times S$ matrix $\mathbf{\Pi} = (\mathbf{\Omega}(\tilde{\mathbf{M}}^* \otimes \mathbf{B})\mathbf{\Omega}^T)^{-1}$, the α -differential of $\delta_1(\mathbf{F})$ with respect to \mathbf{F} is

$$d_{\mathbf{F}}(\delta_1) = -\text{vec}(\mathbf{F})^\dagger \mathbf{D}^\dagger \left((\text{vec}(\mathbf{F})^T \mathbf{D}^T) \otimes \mathbf{I}_S \right) d_{\mathbf{F}}(\mathbf{\Pi}) - \text{vec}(\mathbf{F})^\dagger \mathbf{D}^\dagger \mathbf{\Pi} \mathbf{D}, \quad (78)$$

where

$$d_{\mathbf{F}}(\mathbf{\Pi}) = -(\mathbf{\Pi}^T \otimes \mathbf{\Pi})d_{\mathbf{F}}(\mathbf{\Omega}(\tilde{\mathbf{M}}^* \otimes \mathbf{B})\mathbf{\Omega}^T) = -((\mathbf{\Pi}^T \mathbf{\Omega}) \otimes (\mathbf{\Pi} \mathbf{\Omega}))(\mathbf{I}_K \otimes \mathbf{Z}_{K,K} \otimes \mathbf{I}_K) \cdot (\text{vec}(\tilde{\mathbf{M}}^*) \otimes \mathbf{I}_{K^2})d_{\mathbf{F}}\mathbf{B} \quad (79)$$

and

$$\begin{aligned} d_F(\mathbf{B}) &= d_F(\mathbf{I} - (\mathbf{I} + \mathbf{F}\mathbf{F}^\dagger)^{-1}) \\ &= ((\mathbf{I} + \mathbf{F}\mathbf{F}^\dagger)^{-T}) \otimes ((\mathbf{I} + \mathbf{F}\mathbf{F}^\dagger)^{-1}) d_F(\mathbf{I} + \mathbf{F}\mathbf{F}^\dagger) \\ &= ((\mathbf{I} + \mathbf{F}\mathbf{F}^\dagger)^{-T}) \otimes ((\mathbf{I} + \mathbf{F}\mathbf{F}^\dagger)^{-1}) (\mathbf{F}^* \otimes \mathbf{I}_K) \\ &= (\mathbf{F}^* (\mathbf{I} + \mathbf{F}\mathbf{F}^\dagger)^{-T}) \otimes (\mathbf{I} - \mathbf{B}). \end{aligned} \quad (80)$$

Then, defining a $K \times K$ matrix

$$\tilde{\mathbf{F}} = (\mathbf{I} + \mathbf{F}^\dagger \mathbf{F})^{-1} \mathbf{F}^\dagger$$

and a $K^4 \times K^2$ matrix

$$\Psi = (\mathbf{I}_K \otimes \mathbf{Z}_{K,K} \otimes \mathbf{I}_K) (\text{vec}(\tilde{\mathbf{M}}^*) \otimes \mathbf{I}_{K^2}), \quad (81)$$

and by combing (77)-(81), we finally have when $\mathbf{P} \neq \mathbf{0}$ that

$$\begin{aligned} d_F(I_{\text{GMI}}(\mathbf{W}_{\text{opt}}, \mathbf{T}_{2\text{pt}}, \mathbf{F})) &= d_F(I_1) + d_F(\delta_1) \\ &= \text{vec}(\mathbf{F}\mathbf{M} + \tilde{\mathbf{F}}^\dagger)^\dagger - \text{vec}(\mathbf{F})^\dagger \mathbf{D}^\dagger \Pi \mathbf{D} \\ &\quad + \text{vec}(\mathbf{F})^\dagger \mathbf{D}^\dagger ((\Omega^T \Pi \mathbf{D} \text{vec}(\mathbf{F}))^T \otimes (\Pi \Omega)) \Psi (\tilde{\mathbf{F}}^T \otimes (\mathbf{I} - \mathbf{B})). \end{aligned}$$

APPENDIX E THE PROOF OF PROPOSITION 3

As the formula of GMI in (8) is quadratic in \mathbf{V} and no constraints apply to \mathbf{V} , taking the gradient of $I_{\text{GMI}}(\mathbf{V}, \mathbf{R}, \mathbf{G})$ with respect to \mathbf{V} and setting it to zero, yields the optimal \mathbf{V} in (22). Inserting \mathbf{V}_{opt} back into (8) gives, after some manipulations,

$$\begin{aligned} I_{\text{GMI}}(\mathbf{V}_{\text{opt}}, \mathbf{R}, \mathbf{G}) &= K + \log(\det(\mathbf{I} + \mathbf{G})) + 2\mathcal{R}\{\text{Tr}(\mathbf{P}\mathbf{M}\mathbf{R})\} \\ &\quad + \text{Tr}(\mathbf{M}(\mathbf{I} + \mathbf{G})) + \text{Tr}((\mathbf{I} + \mathbf{G})^{-1} \mathbf{R} \tilde{\mathbf{M}} \mathbf{R}^\dagger). \end{aligned} \quad (82)$$

When $\mathbf{P} = \mathbf{0}$, (82) equals

$$I_2(\mathbf{G}) = K + \log(\det(\mathbf{I} + \mathbf{G})) + \text{Tr}(\mathbf{M}(\mathbf{I} + \mathbf{G})).$$

Further, when $\mathbf{P} \neq \mathbf{0}$ the terms of $I_{\text{GMI}}(\mathbf{V}_{\text{opt}}, \mathbf{R}, \mathbf{G})$ in (82) related to \mathbf{R} are

$$g(\mathbf{R}) = 2\mathcal{R}\{\text{Tr}(\mathbf{P}\mathbf{M}\mathbf{R})\} + \text{Tr}((\mathbf{I} + \mathbf{G})^{-1} \mathbf{R} \tilde{\mathbf{M}} \mathbf{R}^\dagger).$$

Let \mathbf{r}_k denote the k th column of \mathbf{R} , with all elements in rows $[\max(0, k - \nu_R), \min(k + \nu_R, K - 1)]$ removed, and define the column vector $\mathbf{r} = [\mathbf{r}_0^T, \mathbf{r}_1^T, \dots, \mathbf{r}_{K-1}^T]^T$. Then, we have

$$\mathbf{r} = \Omega \text{vec}(\mathbf{R}).$$

Further, let \mathbf{d}_k denote the k th column of the matrix $\mathbf{M}\mathbf{P}$ with all elements in rows $[\max(0, k - \nu_R), \min(k + \nu_R, K - 1)]$ removed, and define the vector $\mathbf{d} = [\mathbf{d}_0^T, \mathbf{d}_1^T, \dots, \mathbf{d}_{K-1}^T]^T$. From the definition of \mathbf{d} , we have

$$\mathbf{d} = \Omega \text{vec}(\mathbf{M}\mathbf{P}).$$

Defining a Hermitian matrix $\hat{\mathbf{B}}_2$ as

$$\hat{\mathbf{B}}_2 = \Omega (\tilde{\mathbf{M}}^* \otimes (\mathbf{I} + \mathbf{G})^{-1}) \Omega^T,$$

we can write $f(\mathbf{R})$ as

$$g(\mathbf{R}) = \mathbf{r}^\dagger \hat{\mathbf{B}}_2 \mathbf{r} + 2\mathcal{R}\{\mathbf{d}^\dagger \mathbf{r}\}.$$

Therefore, the optimal \mathbf{r} is

$$\mathbf{r}_{\text{opt}} = -\hat{\mathbf{B}}_2^{-1} \mathbf{d}. \quad (83)$$

Transferring \mathbf{r}_{opt} back into \mathbf{R}_{opt} gives the optimal \mathbf{R} in (23), and inserting it back into $g(\mathbf{R})$ gives

$$g(\mathbf{R}_{\text{opt}}) = -\mathbf{d}^\dagger \hat{\mathbf{B}}_2^{-1} \mathbf{d}.$$

Thus, with the optimal \mathbf{V} and \mathbf{R} , when $\mathbf{P} \neq \mathbf{0}$ the GMI equals

$$\begin{aligned} I_{\text{GMI}}(\mathbf{V}_{\text{opt}}, \mathbf{R}_{\text{opt}}, \mathbf{G}) &= K + \log(\det(\mathbf{I} + \mathbf{G})) + \text{Tr}(\mathbf{M}(\mathbf{I} + \mathbf{G})) \\ &\quad - \mathbf{d}^\dagger (\Omega (\tilde{\mathbf{M}}^* \otimes (\mathbf{I} + \mathbf{G})^{-1}) \Omega^T)^{-1} \mathbf{d}. \end{aligned}$$

APPENDIX F THE GRADIENT IN METHOD II FOR THE MIMO CASE

Taking the α -differential of $I_2(\mathbf{G})$ with respect to \mathbf{G} yields

$$d_G(I_2) = \text{vec}((\mathbf{I} + \mathbf{G})^{-1} + \mathbf{M})^\dagger. \quad (84)$$

Defining an $S \times S$ Hermitian matrix

$$\Phi = (\Omega (\tilde{\mathbf{M}}^* \otimes (\mathbf{I} + \mathbf{G})^{-1}) \Omega^T)^{-1}$$

and taking the α -differential of $\delta_2(\mathbf{G})$ with respect to \mathbf{G} yields

$$\begin{aligned} d_G(\delta_2) &= -(\mathbf{d}^T \otimes \mathbf{d}^\dagger) d_G(\Phi) \\ &= (\mathbf{d}^T \otimes \mathbf{d}^\dagger) (\Phi^T \otimes \Phi) d_G(\Omega (\tilde{\mathbf{M}}^* \otimes (\mathbf{I} + \mathbf{G})^{-1}) \Omega^T) \\ &= ((\mathbf{d}^T \Phi^T) \otimes (\mathbf{d}^\dagger \Phi)) (\Omega \otimes \Omega) d_G(\tilde{\mathbf{M}}^* \otimes (\mathbf{I} + \mathbf{G})^{-1}) \\ &= ((\mathbf{d}^T \Phi^T \Omega) \otimes (\mathbf{d}^\dagger \Phi \Omega)) \Psi d_G((\mathbf{I} + \mathbf{G})^{-1}) \\ &= -((\mathbf{d}^T \Phi^T \Omega) \otimes (\mathbf{d}^\dagger \Phi \Omega)) \Psi ((\mathbf{I} + \mathbf{G})^{-T} \otimes (\mathbf{I} + \mathbf{G})^{-1}) \end{aligned} \quad (85)$$

where Ψ is defined in (81). Combining (84) and (85), we obtain

$$\begin{aligned} d_G(I_{\text{GMI}}(\mathbf{V}_{\text{opt}}, \mathbf{R}_{\text{opt}}, \mathbf{G})) &= d_G(I_2) + d_G(\delta_2) \\ &= \text{vec}((\mathbf{I} + \mathbf{G})^{-1} + \mathbf{M})^\dagger \\ &\quad - ((\mathbf{d}^T \Phi^T \Omega) \otimes (\mathbf{d}^\dagger \Phi \Omega)) \Psi ((\mathbf{I} + \mathbf{G})^{-T} \otimes (\mathbf{I} + \mathbf{G})^{-1}). \end{aligned}$$

APPENDIX G THE CONCAVITY PROOF OF METHOD II FOR THE MIMO CASE

When $\mathbf{P} = \mathbf{0}$, as $\log(\det(\mathbf{I} + \mathbf{G}))$ is concave [64] and $\text{Tr}(\mathbf{M}(\mathbf{I} + \mathbf{G}))$ is linear in \mathbf{G} , the function $I_2(\mathbf{G})$ in (25) is concave with respect to \mathbf{G} whenever $\mathbf{I} + \mathbf{G} \succ \mathbf{0}$.

The concavity when $\mathbf{P} \neq \mathbf{0}$ can be deduced from the composition theorem in [64, Ch. 3.6]. For a positive-definite matrix \mathbf{X} , $\mathbf{d}^\dagger \mathbf{X}^{-1} \mathbf{d}$ is convex and non-increasing (with respect to the generalized inequality for positive-definite Hermitian matrices, see [64], [65]) for any column vector \mathbf{d} . Furthermore, since $\mathbf{I} + \mathbf{G} \succ \mathbf{0}$, $(\mathbf{I} + \mathbf{G})^{-1}$ is convex. Since $\tilde{\mathbf{M}} \prec \mathbf{0}$ $\mathbf{X} = \Omega (\tilde{\mathbf{M}}^* \otimes (\mathbf{I} + \mathbf{G})^{-1}) \Omega^T$ is concave in \mathbf{G} . By the composition theorem, $\mathbf{d}^\dagger (\Omega (\tilde{\mathbf{M}}^* \otimes (\mathbf{I} + \mathbf{G})^{-1}) \Omega^T)^{-1} \mathbf{d}$ is convex, and $\delta_2(\mathbf{G})$ is then concave. Therefore, the function $I_{\text{GMI}}(\mathbf{V}_{\text{opt}}, \mathbf{R}_{\text{opt}}, \mathbf{G})$ in (24) is concave with respect to \mathbf{G} whenever $\mathbf{I} + \mathbf{G} \succ \mathbf{0}$.

APPENDIX H THE PROOF OF PROPOSITION 5

The Fourier series associated to the Toeplitz matrix \mathbf{W} is

$$W(\omega) = \sum_{k=-\infty}^{\infty} w_k \exp(jk\omega),$$

and the differential of $\bar{I}(W(\omega), T(\omega), F(\omega))$ in (45) with respect to w_k (where ω is fixed) is

$$\begin{aligned} \frac{\partial \bar{I}}{\partial w_k} = & -\frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{|F(\omega)|^2 (N_0 + |H(\omega)|^2) W^*(\omega)}{1 + |F(\omega)|^2} \exp(jk\omega) d\omega \\ & + \frac{1}{\pi} \int_{-\pi}^{\pi} \left(F^*(\omega) H(\omega) + \frac{\alpha |F(\omega)|^2 H(\omega) T^*(\omega)}{1 + |F(\omega)|^2} \right) \\ & \times \exp(jk\omega) d\omega. \end{aligned} \quad (86)$$

Since (86) should equal zero for all k , the optimal $W(\omega)$ is given in (50). Inserting $W_{\text{opt}}(\omega)$ back into (45) yields

$$\begin{aligned} \bar{I}(W_{\text{opt}}(\omega), T(\omega), F(\omega)) \\ = 1 + \frac{\alpha}{\pi} \int_{-\pi}^{\pi} \mathcal{R}\{F^*(\omega) T(\omega) M(\omega)\} d\omega \\ + \frac{1}{2\pi} \int_{-\pi}^{\pi} \left(\log(1 + |F(\omega)|^2) \right. \\ \left. + \frac{\tilde{M}(\omega) |T(\omega) F(\omega)|^2}{1 + |F(\omega)|^2} + M(\omega) (1 + |F(\omega)|^2) \right) d\omega. \end{aligned} \quad (87)$$

where $M(\omega)$ and $\tilde{M}(\omega)$ are defined in (46) and (47), respectively.

When $\alpha = 0$, the GMI in (87) equals (53) while when $0 < \alpha \leq 1$, the terms related to $T(\omega)$ in (87) are

$$\begin{aligned} f(T(\omega)) = & \frac{\alpha}{\pi} \int_{-\pi}^{\pi} \mathcal{R}\{F^*(\omega) T(\omega) M(\omega)\} d\omega \\ & + \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{\tilde{M}(\omega) |T(\omega) F(\omega)|^2}{1 + |F(\omega)|^2} d\omega. \end{aligned} \quad (88)$$

As the elements of the main diagonal and the first ν lower diagonals of matrix \mathbf{T} are constrained to zeros, we define the vector $\tilde{\mathbf{t}}$ that specifies the Toeplitz matrix \mathbf{T} as

$$\tilde{\mathbf{t}} = [t_{-N_T}, \dots, t_{-1}, t_{\nu+1}, \dots, t_{N_T}],$$

and with $\phi(\omega)$ defined in (48), the Fourier series $T(\omega)$ with a finite tap-length N_T is

$$T(\omega) = \sum_{-N_T \leq k \leq N_T, k \notin [0, \nu]} t_k \exp(jk\omega) = \tilde{\mathbf{t}} \phi(\omega). \quad (89)$$

Further, with $\mathbf{e}_1, \mathbf{e}_2$ defined in (49), (88) can be rewritten as

$$f(T(\omega)) = \tilde{\mathbf{t}} \mathbf{e}_2 \tilde{\mathbf{t}}^\dagger + 2\mathcal{R}\{\tilde{\mathbf{t}} \mathbf{e}_1\}.$$

Therefore, the optimal $\tilde{\mathbf{t}}$ is

$$\tilde{\mathbf{t}}_{\text{opt}} = -\mathbf{e}_1^\dagger \mathbf{e}_2^{-1}. \quad (90)$$

Inserting $\tilde{\mathbf{t}}_{\text{opt}}$ back into (87)-(89), the optimal $T(\omega)$ is then in (51), and $\bar{I}(W(\omega), T(\omega), F(\omega))$ with the optimal $W(\omega)$ and $T(\omega)$ is given in (52) after some manipulations.

APPENDIX I THE PROOF OF PROPOSITION 7

The Fourier series associated to the Toeplitz matrix \mathbf{V} is

$$V(\omega) = \sum_{k=-\infty}^{\infty} v_k \exp(jk\omega)$$

and the differential of $\bar{I}(V(\omega), R(\omega), G(\omega))$ in (57) with respect to v_k (where ω is fixed) is

$$\begin{aligned} \frac{\partial \bar{I}}{\partial v_k} = & -\frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{(N_0 + |H(\omega)|^2) V^*(\omega)}{1 + G(\omega)} \exp(jk\omega) d\omega \\ & + \frac{1}{\pi} \int_{-\pi}^{\pi} \left(H(\omega) + \frac{\alpha H(\omega) R^*(\omega)}{1 + G(\omega)} \right) \exp(jk\omega) d\omega. \end{aligned} \quad (91)$$

Since (91) equals zero for all k , the optimal $V(\omega)$ is given in (61). Putting $V_{\text{opt}}(\omega)$ in (61) back into (57) yields

$$\begin{aligned} \bar{I}(V_{\text{opt}}(\omega), R(\omega), G(\omega)) \\ = 1 + \frac{\alpha}{\pi} \int_{-\pi}^{\pi} \mathcal{R}\{M(\omega) R(\omega)\} d\omega + \frac{1}{2\pi} \int_{-\pi}^{\pi} \left(\log(1 + G(\omega)) \right. \\ \left. + \frac{\tilde{M}(\omega) |R(\omega)|^2}{1 + G(\omega)} + M(\omega) (1 + G(\omega)) \right) d\omega. \end{aligned} \quad (92)$$

When $\alpha = 0$, the GMI in (92) equals (64), and when $0 < \alpha \leq 1$, the terms of $\bar{I}(V_{\text{opt}}(\omega), R(\omega), G(\omega))$ related to $R(\omega)$ in (92) are

$$\begin{aligned} g(R(\omega)) = & \frac{\alpha}{\pi} \int_{-\pi}^{\pi} \mathcal{R}\{M(\omega) R(\omega)\} d\omega \\ & + \frac{1}{2\pi} \int_{-\pi}^{\pi} \frac{\tilde{M}(\omega) |R(\omega)|^2}{1 + G(\omega)} d\omega. \end{aligned} \quad (93)$$

Defining the vector $\tilde{\mathbf{r}}$ that specifies the Toeplitz matrix \mathbf{R} as

$$\tilde{\mathbf{r}} = [r_{-N_R}, \dots, r_{-\nu_R-1}, r_{\nu_R+1}, \dots, r_{N_R}],$$

and with $\psi(\omega)$ defined in (59), the Fourier series $R(\omega)$ with a finite tap-length N_R is

$$R(\omega) = \sum_{-N_R \leq k \leq N_R, k \notin [-\nu_R, \nu_R]} r_k \exp(jk\omega) = \tilde{\mathbf{r}} \psi(\omega) \quad (94)$$

where $2\nu_R + 1$ is the band-size that \mathbf{R} is constrained to zero. With $\boldsymbol{\zeta}_1$ and $\boldsymbol{\zeta}_2$ defined in (60), (93) can be written as

$$g(R(\omega)) = \tilde{\mathbf{r}} \boldsymbol{\zeta}_2 \tilde{\mathbf{r}}^\dagger + 2\mathcal{R}\{\tilde{\mathbf{r}} \boldsymbol{\zeta}_1\}.$$

Therefore, the optimal $\tilde{\mathbf{r}}$ is

$$\tilde{\mathbf{r}}_{\text{opt}} = -\boldsymbol{\zeta}_1^\dagger \boldsymbol{\zeta}_2^{-1}. \quad (95)$$

This shows that $\tilde{\mathbf{r}}_{\text{opt}}$ has Hermitian symmetry, since $G(\omega)$, $M(\omega)$, and $\tilde{M}(\omega)$ are all real-valued and $R_{\text{opt}}(\omega)$ is thusly real-valued. Putting $\tilde{\mathbf{r}}_{\text{opt}}$ back into (92)-(94), the optimal $R(\omega)$ is in (62) and $\bar{I}(V(\omega), R(\omega), G(\omega))$ with the optimal $V(\omega)$ and $R(\omega)$ is in (63), respectively.

APPENDIX J

THE CONCAVITY PROOF OF METHOD II WITH ISI CHANNELS

In order to prove that $\bar{I}(V_{\text{opt}}(\omega), R_{\text{opt}}(\omega), G(\omega))$ in (63) is concave with respect to $G(\omega)$, it is sufficient to prove that $\xi_1^\dagger \xi_2^{-1} \xi_1$ is convex with respect to $G(\omega)$. For a positive-definite matrix ξ_2 , $\xi_1^\dagger \xi_2^{-1} \xi_1$ is convex and non-increasing (with respect to a generalized inequality for positive-definite Hermitian matrices) in $G(\omega)$ for any vector ξ_1 and with arbitrary finite tap-length N_R . As matrix $\tilde{\mathbf{M}} \prec \mathbf{0}$, ξ_2 in (60) is concave with respect to $G(\omega)$ under the constraint that $\mathbf{I} + \mathbf{G} \succ \mathbf{0}$ is positive-definite. Hence, $\xi_1^\dagger \xi_2^{-1} \xi_1$ is convex in $G(\omega)$ by the composition theorem [64].

APPENDIX K

THE PROOF OF LEMMA 5

From Theorem 2, the optimal \mathbf{G} in Method III satisfies

$$[(\mathbf{I} + \mathbf{G}_{\text{opt}})^{-1}]_v = -[\hat{\mathbf{M}}]_v.$$

Note that, when $\mathbf{P} = \mathbf{0}$ Method II and III are equivalent as $\hat{\mathbf{M}} = \mathbf{M}$. Hence, in order to prove Lemma 4, it is sufficient to show that $[\hat{\mathbf{M}}]_v$ converges to $[\mathbf{M}]_v$ as $N_0 \rightarrow 0$ and ∞ in Method III. When $\mathbf{P} \prec \mathbf{I}$, $\mathbf{C}_k \succ \mathbf{0}$ in (34), and when $N_0 \rightarrow 0$,

$$\begin{aligned} & \mathbf{H}^\dagger (\mathbf{H} \mathbf{C}_k \mathbf{H}^\dagger + N_0 \mathbf{I})^{-1} \mathbf{H} \\ &= \mathbf{C}_k^{-1} (\mathbf{H}^\dagger \mathbf{H} + N_0 \mathbf{C}_k^{-1})^{-1} \mathbf{H}^\dagger \mathbf{H} \\ &= \mathbf{C}_k^{-1} (\mathbf{I} - N_0 \mathbf{C}_k^{-1} (\mathbf{H}^\dagger \mathbf{H})^{-1}) + \mathcal{O}(N_0^2). \end{aligned} \quad (96)$$

Therefore, with $\hat{\mathbf{W}}$ and $\hat{\mathbf{C}}$ defined through (33)-(37) and using (96), it holds that

$$\begin{aligned} \hat{\mathbf{W}} \mathbf{H} &= \mathbf{I} - N_0 (\mathbf{H}^\dagger \mathbf{H})^{-1} + \mathcal{O}(N_0^2), \\ \hat{\mathbf{C}} &= [\hat{\mathbf{W}} \mathbf{H}]_{\setminus v} = -N_0 [(\mathbf{H}^\dagger \mathbf{H})^{-1}]_{\setminus v} + \mathcal{O}(N_0^2). \end{aligned} \quad (97)$$

With (97) and $\hat{\mathbf{M}}$ in (40), it can be verified that

$$\lim_{N_0 \rightarrow 0} [\hat{\mathbf{M}}/N_0]_v = -[(\mathbf{H}^\dagger \mathbf{H})^{-1}]_v.$$

On the other hand, when $N_0 \rightarrow \infty$, from (33)-(40) we have

$$\begin{aligned} N_0 \hat{\mathbf{W}} &= \mathbf{H}^\dagger (\mathbf{H} \mathbf{C}_k \mathbf{H}^\dagger / N_0 + \mathbf{I})^{-1} = \mathbf{H}^\dagger + \mathcal{O}(1/N_0), \\ N_0 \hat{\mathbf{C}} &= [\hat{\mathbf{W}} \mathbf{H}]_{\setminus v} = [\mathbf{H}^\dagger \mathbf{H}]_{\setminus v} + \mathcal{O}(1/N_0). \end{aligned} \quad (98)$$

With (98), it can be verified that

$$\lim_{N_0 \rightarrow \infty} [N_0(\mathbf{I} + \hat{\mathbf{M}})]_v = [\mathbf{H}^\dagger \mathbf{H}]_v.$$

Therefore, from (69) it shows that $[\hat{\mathbf{M}}]_v$ converges to $[\mathbf{M}]_v$ when $N_0 \rightarrow 0$ and ∞ , which completes the proof.

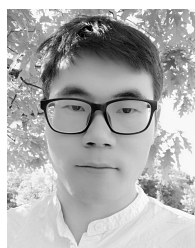
ACKNOWLEDGMENT

This work has been included my Ph.D. dissertation [1]. This paper was presented in part at the IEEE 26th Annual International Symposium on Personal, Indoor and Mobile Radio Communications, Hong Kong, September 2015 [2].

REFERENCES

- [1] S. Hu, "Channel shortening in wireless communication," Ph.D. dissertation, Dept. Elect. Inf. Technol., Lund Univ., Lund, Sweden, Dec. 2017.
- [2] S. Hu and F. Rusek, "On the design of reduced state demodulators with interference cancellation for iterative receivers," in *Proc. IEEE 26th Annu. Int. Symp. Pers., Indoor, Mobile Radio Commun. (PIMRC)*, Sep. 2015, pp. 981–985.
- [3] N. Al-Dhahir, "FIR channel-shortening equalizers for MIMO ISI channels," *IEEE Trans. Commun.*, vol. 49, no. 2, pp. 213–218, Feb. 2001.
- [4] D. D. Falconer and F. R. Magee, Jr., "Adaptive channel memory truncation for maximum likelihood sequence estimation," *Bell Syst. Tech. J.*, vol. 51, no. 9, pp. 1541–1562, Nov. 1973.
- [5] S. A. Fredricsson, "Joint optimization of transmitter and receiver filters in digital PAM systems with a Viterbi detector," *IEEE Trans. Inf. Theory*, vol. IT-22, no. 2, pp. 200–210, Mar. 1976.
- [6] C. T. Beare, "The choice of the desired impulse response in combined linear-Viterbi algorithm equalizers," *IEEE Trans. Commun.*, vol. COM-26, pp. 1301–1307, Aug. 1978.
- [7] N. Sundström, O. Edfors, P. Ödling, H. Eriksson, T. Koski, and P. O. Börjesson, "Combined linear-Viterbi equalizers—A comparative study and a minimax design," in *Proc. IEEE Veh. Technol. Conf. (VTC)*, Stockholm, Sweden, vol. 2, Jun. 1994, pp. 1263–1267.
- [8] N. Al-Dhahir and J. M. Cioffi, "Efficiently computed reduced-parameter input-aided MMSE equalizers for ML detection: A unified approach," *IEEE Trans. Inf. Theory*, vol. 42, no. 3, pp. 903–915, May 1996.
- [9] M. Lagunas, A. I. Perez-Neia, and J. Vidal, "Joint beamforming and Viterbi equalizer in wireless communications," in *Proc. Conf. Rec. 31st Asilomar Conf. Signals, Syst. Comput.*, vol. 1, Nov. 1997, pp. 915–919.
- [10] S. A. Aldosari, S. A. Alshebeili, and A. M. Al-Sanie, "A new MSE approach for combined linear-Viterbi equalizers," in *Proc. IEEE Veh. Technol. Conf. (VTC)*, Tokyo, Japan, vol. 3, May 2000, pp. 1707–1711.
- [11] R. Venkataramani and S. Sankaranarayanan, "Optimal channel shortening equalization for MIMO ISI channels," in *Proc. IEEE Global Telecomm. (GLOBECOM)*, New Orleans, LO, USA, Dec. 2008, pp. 1–5.
- [12] A. Shaheem, "Iterative detection for wireless communications," Ph.D. dissertation, School Elect., Electron. Comput. Eng., Univ. Western Australia, Perth, WA, USA, 2008.
- [13] U. L. Dang, W. H. Gerstacker, and D. T. M. Slock, "Maximum SINR prefiltering for reduced-state trellis-based equalization," in *Proc. IEEE Int. Conf. Commun. (ICC)*, Kyoto, Japan, Jun. 2011, pp. 1–6.
- [14] R. Venkataramani and M. F. Erden. (2007). "A posteriori equivalence: A new perspective for design of optimal channel shortening equalizers." [Online]. Available: <https://arxiv.org/abs/0710.3802>
- [15] I. Abov-Paycal and A. Lapidoth, "On the capacity of reduced-complexity receivers for intersymbol interference channels," in *Proc. Conf. Inf. Sci. Syst. (CISS)*, Princeton, NJ, USA: Princeton Univ., Mar. 2000, pp. 263–266.
- [16] D. Darsena and F. Verde, "Minimum-mean-output-energy blind adaptive channel shortening for multicarrier SIMO transceivers," *IEEE Trans. Signal Process.*, vol. 55, no. 12, pp. 5755–5771, Dec. 2007.
- [17] G. D. Forney, Jr., "Maximum-likelihood sequence estimation of digital sequences in the presence of intersymbol interference," *IEEE Trans. Inf. Theory*, vol. IT-18, no. 3, pp. 363–378, May 1972.
- [18] F. Rusek and A. Prlja, "Optimal channel shortening for MIMO and ISI channels," *IEEE Trans. Wireless Commun.*, vol. 11, no. 2, pp. 810–818, Feb. 2012.
- [19] A. Viterbi, "Error bounds for convolutional codes and an asymptotically optimum decoding algorithm," *IEEE Trans. Inf. Theory*, vol. IT-13, no. 2, pp. 260–269, Apr. 1967.
- [20] N. Merhav, G. Kaplan, A. Lapidoth, and S. Shamai (Shitz), "On information rates for mismatched decoders," *IEEE Trans. Inf. Theory*, vol. 40, no. 6, pp. 1953–1967, Nov. 1994.
- [21] A. Ganti, A. Lapidoth, and İ. E. Telatar, "Mismatched decoding revisited: General alphabets, channels with memory, and the wide-band limit," *IEEE Trans. Inf. Theory*, vol. 46, no. 7, pp. 2315–2328, Nov. 2000.
- [22] M. R. McKay, I. B. Collings, and A. M. Tulino, "Achievable sum rate of MIMO MMSE receivers: A general analytic framework," *IEEE Trans. Inf. Theory*, vol. 56, no. 1, pp. 396–410, Jan. 2010.
- [23] H. Weingarten, Y. Steinberg, and S. Shamai (Shitz), "Gaussian codes and weighted nearest neighbor decoding in fading multiple-antenna channels," *IEEE Trans. Inf. Theory*, vol. 50, no. 8, pp. 1665–1686, Aug. 2004.
- [24] M. Tüchler and A. C. Singer, "Turbo equalization: An overview," *IEEE Trans. Inf. Theory*, vol. 57, no. 2, pp. 920–952, Feb. 2011.

- [25] S.-J. Lee, A. C. Singer, and N. R. Shanbhag, "Linear turbo equalization analysis via BER transfer and EXIT charts," *IEEE Trans. Signal Process.*, vol. 53, no. 8, pp. 2883–2897, Aug. 2005.
- [26] A. Shaheem, H.-J. Zepernick, and M. Caldera, "Enhanced channel shortened turbo equalization," in *Proc. Int. Conf. Adv. Technol. Commun. (ATC)*, Oct. 2008, pp. 8–11.
- [27] A. Glavieux, C. Laot, and J. Labat, "Turbo equalization over a frequency selective channel," in *Proc. Int. Symp. Turbo Codes Rel. Topics*, Brest, France, Sep. 1997, pp. 96–102.
- [28] R. R. Lopes and J. R. Barry, "The soft-feedback equalizer for turbo equalization of highly dispersive channels," *IEEE Trans. Commun.*, vol. 54, no. 5, pp. 783–788, May 2006.
- [29] J. W. Choi, B. Shim, A. C. Singer, and N. I. Cho, "Low-complexity decoding via reduced dimension maximum-likelihood search," *IEEE Trans. Signal Process.*, vol. 58, no. 3, pp. 1780–1793, Mar. 2010.
- [30] J. W. Choi, B. Lee, and B. Shim, "Iterative group detection and decoding for large MIMO systems," *J. Commun. Netw.*, vol. 17, no. 6, pp. 609–621, Dec. 2015.
- [31] J. W. Choi, A. C. Singer, J. W. Lee, and N. I. Cho, "Improved linear soft-input soft-output detection via soft feedback successive interference cancellation," *IEEE Trans. Commun.*, vol. 58, no. 3, pp. 986–996, Mar. 2010.
- [32] M. Tüchler, A. C. Singer, and R. Kötter, "Minimum mean squared error equalization using a priori information," *IEEE Trans. Signal Process.*, vol. 50, no. 3, pp. 673–683, Mar. 2000.
- [33] S. Hu, H. Kröll, Q. Huang, and F. Rusek, "Optimal channel shortener design for reduced-state soft-output Viterbi equalizer in single-carrier systems," *IEEE Trans. Commun.*, vol. 65, no. 6, pp. 2568–2582, Jun. 2017.
- [34] A. Berthet, R. Visoz, and P. Tortelier, "Sub-optimal turbo-detection for coded 8-PSK signals over ISI channels with application to EDGE advanced mobile system," in *Proc. IEEE Veh. Technol. Conf. (VTC)*, Sep. 2000, pp. 151–157.
- [35] A. Duel-Hallen and C. Heegard, "Delayed decision-feedback sequence estimation," *IEEE Trans. Commun.*, vol. 37, no. 5, pp. 428–436, May 1989.
- [36] F. Rusek, N. Al-Dhahir, and A. Goma, "A rate-maximizing channel-shortening detector with soft feedback side information," in *Proc. IEEE Global Telecom. (GLOBECOM)*, Anaheim, CA, USA, Dec. 2012, pp. 1–6.
- [37] C. Studer, S. Fateh, and D. Seethaler, "ASIC implementation of soft-input soft-output MIMO detection using MMSE parallel interference cancellation," *IEEE J. Solid-State Circuits*, vol. 46, no. 7, pp. 1754–1765, Jul. 2011.
- [38] M. Witzke, S. Bärö, F. Schreckenbach, and J. Hagenauer, "Iterative detection of MIMO signals with linear detectors," in *Proc. Asilomar Conf. Signals, Syst. Comput. (ACSSC)*, Monterey, CA, USA, Nov. 2002, pp. 289–293.
- [39] J. Zhang, H. Nguyen, and G. Mandyam, "LMMSE-based iterative and turbo equalization methods for CDMA downlink channels," in *Proc. IEEE 6th Workshop Signal Process. Adv. Wireless Commun.*, Jun. 2005, pp. 231–235.
- [40] R. M. Gray, "Toeplitz and circulant matrices: A review," *Found. Trends Commun. Inf. Theory*, vol. 2, no. 3, pp. 155–239, 2006.
- [41] S. ten Brink, "Convergence behavior of iteratively decoded parallel concatenated codes," *IEEE Trans. Commun.*, vol. 49, no. 10, pp. 1727–1737, Oct. 2001.
- [42] J. Hagenauer, "Source-controlled channel decoding," *IEEE Trans. Commun.*, vol. 43, no. 9, pp. 2449–2457, Sep. 1995.
- [43] S. M. Kay, *Fundamentals of Statistical Signal Processing: Estimation Theory*, vol. 1. Upper Saddle River, NJ, USA: Prentice-Hall, 1993.
- [44] B. M. Hochwald and S. ten Brink, "Achieving near-capacity on a multiple-antenna channel," *IEEE Trans. Commun.*, vol. 51, no. 3, pp. 389–399, Mar. 2003.
- [45] L. Bahl, J. Cocke, F. Jelinek, and J. Raviv, "Optimal decoding of linear codes for minimizing symbol error rate," *IEEE Trans. Inf. Theory*, vol. IT-20, no. 2, pp. 284–287, Mar. 1974.
- [46] F. Rusek and D. Fertonani, "Bounds on the information rate of intersymbol interference channels based on mismatched receivers," *IEEE Trans. Inf. Theory*, vol. 58, no. 3, pp. 1470–1482, Mar. 2012.
- [47] G. Ungerboeck, "Adaptive maximum-likelihood receiver for carrier-modulated data-transmission systems," *IEEE Trans. Commun.*, vol. COM-22, no. 5, pp. 624–636, May 1974.
- [48] F. Rusek, G. Colavolpe, C. E. W. Sundberg, "40 years with the Ungerboeck model: A look at its potentialities [Lecture Notes]," *Signal Process. Mag.*, vol. 32, no. 3, pp. 156–161, May 2015.
- [49] F. Rusek, M. Loncar, and A. Prlija, "A comparison of ungerboeck and Forney models for reduced-complexity ISI equalization," in *Proc. IEEE Global Telecom. (GLOBECOM)*, Washington, DC, USA, Dec. 2007, pp. 1431–1436.
- [50] G. Caire and S. Shamai (Shitz), "On the achievable throughput of a multiantenna Gaussian broadcast channel," *IEEE Trans. Inf. Theory*, vol. 49, no. 7, pp. 1691–1706, Jul. 2003.
- [51] G. L. Turin, "An introduction to digital matched filters," *Proc. IEEE*, vol. 64, no. 7, pp. 1092–1112, Jul. 1976.
- [52] A. Kavčić and J. M. F. Moura, "Matrices with banded inverses: inversion algorithms and factorization of Gauss–Markov processes," *IEEE Trans. Inf. Theory*, vol. 46, no. 4, pp. 1495–1509, Jul. 2000.
- [53] G. Colavolpe and A. Barbieri, "On MAP symbol detection for ISI channels using the Ungerboeck observation model," *IEEE Commun. Lett.*, vol. 9, no. 8, pp. 720–722, Aug. 2005.
- [54] O. Edfors, M. Sandell, J. J. van de Beek, S. K. Wilson, and P. O. Börjesson, "OFDM channel estimation by singular value decomposition," *IEEE Trans. Commun.*, vol. 46, no. 7, pp. 931–939, Jul. 1998.
- [55] W. Hirt, "Capacity and information rates of discrete-time channels with memory," Ph.D. dissertation, Inst. Signal Inf. Process., Swiss Federal Inst. Technol., Zürich, Switzerland, 1988.
- [56] U. Grenander and G. Szegő, *Toeplitz Forms and Their Applications*. Berkeley, CA, USA: Univ. California Press, 1958.
- [57] J. G. Proakis and M. Salehi, *Digital Communications*, 5th ed. New York, NY, USA: McGraw-Hill, 2008.
- [58] G. H. Golub and C. F. Van Loan, *Matrix Computations*, 3rd ed. Baltimore, MD, USA: The Johns Hopkins Univ. Press, 1996.
- [59] S. ten Brink, "Convergence of iterative decoding," *Electron. Lett.*, vol. 35, no. 10, pp. 806–808, May 1999.
- [60] *Evolved Universal Terrestrial Radio Access (E-UTRA); Multiplexing and Channel Coding, Release 12*, document TS 36.212, 3GPP, Mar. 2015.
- [61] F. Rusek and O. Edfors, "An information theoretic characterization of channel shortening receivers," in *Proc. 47th Asilomar Conf. Signals, Syst. Comput.*, Pacific Grove, CA, USA, Nov. 2013, pp. 2108–2112.
- [62] J. R. Magnus, "On the concept of matrix derivative," *J. Multivariate Anal.*, vol. 101, no. 9, pp. 2200–2206, Oct. 2001.
- [63] P. L. Fackler, *Notes on Matrix Calculus*. Raleigh, NC, USA: North Carolina State Univ., Sep. 2005.
- [64] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, 2004.
- [65] C. Davis, "Notions generalizing convexity for functions defined on spaces of matrices," in *Proc. Symp. Pure Math.*, vol. 7. Providence, RI, USA: AMS, 1963, pp. 187–201.



SHA HU (S'15–M'18) was born in Hubei, China, in 1985. He received the Ph.D. degree in electrical engineering from Lund University, Lund, Sweden, in 2018, and the M.S. and B.S. degrees in pure mathematics from Wuhan University, China, in 2008 and 2006, respectively. Since 2008, he has been with Huawei Technologies, where he was involved in baseband algorithm research. He has co-authored 18 patents in these fields. Since 2015, he joined the Department of Electrical and Information Technology, Lund University. His research interests include applied information theory, signal processing, multi-input multi-output detection, channel shortening, and precoder design.



FREDRIK RUSEK was born in Lund, Sweden, in 1978. He received the M.S. and Ph.D. degrees in electrical engineering from Lund University, Lund, in 2003 and 2007, respectively. Since 2012, he has been an Associate Professor with the Department of Electrical and Information Technology, Lund Institute of Technology. His research interests include modulation theory, equalization, wireless communications, and applied information theory.