

# Minimum-Information LQG Control Supplementary Material

Roy Fox<sup>†</sup> and Naftali Tishby<sup>†</sup>

## APPENDIX I PERFECTLY MATCHED CHANNEL

In this appendix we construct a channel that is perfectly matched to the sequential source code derived in Theorem 1, in Part I of this paper [1, Section III-B]. Recall that in a perfectly matched source-channel pair the optimal source coding and the optimal channel coding can be implemented jointly for single letters, without requiring longer blocks. This allows us to use them in a perception-action cycle, where we cannot accumulate a block of inputs before emitting an output.

The main results of [2], applied to our setting, can be summarized as follows. We wish to find a memoryless channel into which we can input an encoding  $w_t = g(\hat{x}_{y_t})$ , such that  $\hat{x}_{u_t} = h(\hat{w}_t)$  can be decoded from the channel output  $\hat{w}_t$ . Suppose that we are concerned with the power needed to transmit  $w_t$  and thus the input cost is  $w_t^\top w_t$ . Then the source  $\hat{x}_{y_t}$  and the channel  $w_t \rightarrow \hat{w}_t$  are perfectly matched if there exist an encoder and a decoder such that

- 1) The Kullback-Leibler divergence  $\mathbb{D}[f(\hat{w}_t|w_t)||f(\hat{w}_t)]$  between the conditional and marginal densities of  $\hat{w}_t$ , as a function of  $w_t$ , equals  $c_1 w_t^\top w_t + c_2$ , for some constants  $c_1 \geq 0$  and  $c_2$ ; and
- 2)  $f(\hat{x}_{u_t}|\hat{x}_{y_t})$  satisfies the conditions in Theorem 1.

To meet these conditions, we can choose the channel, the encoder and the decoder to have

$$\begin{aligned} w_t &= D^{1/2} V^\top \Sigma_{\hat{x}_y}^{1/2} \hat{x}_{y_t} \\ \hat{w}_t &= w_t + v_t; & v_t &\sim \mathcal{N}(0, I - D) \\ \hat{x}_{u_t} &= \Sigma_{\hat{x}_y}^{1/2} V D^{1/2} \hat{w}_t, \end{aligned}$$

with  $D$  and  $V$  as in Theorem 1. Then

$$\begin{aligned} \Sigma_w &= D \\ \Sigma_{\hat{w}} &= I \\ \Sigma_{\hat{x}_u} &= \Sigma_{\hat{x}_y}^{1/2} V D V^\top \Sigma_{\hat{x}_y}^{1/2} = \Sigma_{\hat{x}_u; \hat{x}_y}, \end{aligned}$$

and it can be verified that

$$\mathbb{D}[f(\hat{w}_t|w_t)||f(\hat{w}_t)] = \frac{1}{2} w_t^\top \Sigma_{\hat{w}}^{-1} w_t + \text{const},$$

as required.

The capacity of the additive Gaussian noise channel with noise covariance  $I - D$ , under the appropriate expected

<sup>†</sup>School of Computer Science and Engineering, The Hebrew University, {royf, tishby}@cs.huji.ac.il

\*This work was supported by the DARPA MSEE Program, the Gatsby Charitable Foundation, the Israel Science Foundation and the Intel ICRI-CI Institute

power constraint, is indeed achieved by a Gaussian input with covariance  $D$  and is equal to the information rate in Theorem 1. As shown in [2], this means that constraining the expected power  $\Sigma_w$  is equivalent to constraining the information rate  $\mathbb{I}[\hat{x}_{y_t}; \hat{x}_{u_t}]$ .

Note, however, that the matched channel noise covariance depends on the constraint, through the solution in Theorem 1. Moreover, this result is not applicable when the best channel available to the designer of the controller is not the matched channel above, in which case both the channel and the sequential source coding generally need to be adapted.

## APPENDIX II PROOF OF LEMMA 1 OF PART I

In this appendix we restate and prove Lemma 1 of Part I [1, Section IV-A].

*Lemma 1:* Let  $x$  and  $\hat{x}$  be 0-mean jointly Gaussian random variables. The following properties are equivalent:

- 1) There exists a random variable  $u$ , jointly Gaussian with  $x$ , such that  $\hat{x}(u) = \arg \min_{\hat{x}} \mathbb{E}[\|\hat{x} - x\|^2 | u] = \mathbb{E}[x | u]$ .
- 2)  $\Sigma_{\hat{x}; x} = \Sigma_{\hat{x}}$ .
- 3)  $\Sigma_{x|\hat{x}} = \Sigma_x - \Sigma_{\hat{x}}$ , where  $\Sigma_{x|\hat{x}}$  is the conditional covariance matrix of  $x$  given  $\hat{x}$ , implying  $\Sigma_x \succeq \Sigma_{\hat{x}}$ .
- 4)  $\hat{x} = \mathbb{E}[x | \hat{x}]$ .

Such  $\hat{x}$  is called a minimum mean square error (MMSE) estimator (of  $u$ ) for  $x$ .

*Proof:* (1  $\implies$  2) Assume without loss of generality that  $u$  has mean 0. Then

$$\hat{x} = \Sigma_{x;u} \Sigma_u^\dagger u,$$

implying

$$\Sigma_{\hat{x}; x} = \Sigma_{x;u} \Sigma_u^\dagger \Sigma_{u;x} = \Sigma_{\hat{x}}.$$

(2  $\implies$  3)

$$\Sigma_{x|\hat{x}} = \Sigma_x - \Sigma_{x;\hat{x}} \Sigma_{\hat{x}}^\dagger \Sigma_{\hat{x};x} = \Sigma_x - \Sigma_{\hat{x}}.$$

(3  $\implies$  4) Since  $x$  and  $\hat{x}$  are 0-mean and jointly Gaussian, we can write for some  $T$

$$x = T \hat{x} + \xi; \quad \xi \sim \mathcal{N}(0, \Sigma_{x|\hat{x}}),$$

implying

$$\Sigma_x = T \Sigma_{\hat{x}} T^\top + \Sigma_x - \Sigma_{\hat{x}},$$

thus without loss of generality  $T = I$ .

(4  $\implies$  1) Taking  $u = \hat{x}$ , we have

$$\begin{aligned} & \arg \min_{\hat{x}'} \mathbb{E}[\|\hat{x}' - x\|^2 | u] \\ &= \arg \min_{\hat{x}'} (\hat{x}'^\top \hat{x}' - 2\hat{x}'^\top \mathbb{E}[x|u]) + \mathbb{E}[x^\top x | u], \end{aligned}$$

which is optimized by  $\hat{x}' = \mathbb{E}[x|u]$ .  $\square$

### APPENDIX III

#### PROOF OF LEMMA 2 OF PART I

In this appendix we restate and prove Lemma 2 of Part I [1, Section IV-A].

*Lemma 2:* The bounded memoryless LTI controller optimization problem (Problem 1) is solved by a control law of the form

$$\hat{x}_{y_t} = Ky_t \quad (5a)$$

$$\hat{x}_{u_t} = W\hat{x}_{y_t} + \omega_t; \quad \omega_t \sim \mathcal{N}(0, \Sigma_\omega) \quad (5b)$$

$$u_t = L\hat{x}_{u_t}, \quad (5c)$$

where  $W \in \mathbb{R}^{n \times n}$ ,  $\Sigma_\omega \in \mathbb{R}^{n \times n}$ ,  $L \in \mathbb{R}^{\ell \times n}$ ,  $\omega_t$  is independent of  $y_t$ ,  $\hat{x}_{u_t}$  is a MMSE estimator for  $\hat{x}_{y_t}$  and

$$\mathbb{I}[y_t; u_t] = \mathbb{I}[\hat{x}_{y_t}; \hat{x}_{u_t}]. \quad (6)$$

*Proof:* Consider a LTI controller  $\pi$  of the form

$$u_t = Hy_t + \eta_t; \quad \eta_t \sim \mathcal{N}(0, \Sigma_\eta), \quad (III.1)$$

satisfying the Markov network

$$\begin{array}{ccc} x_t & - & y_t & - & u_t \\ & & | & & | \\ & & \hat{x}_{y_t} & & \hat{x}_{u_t}. \end{array}$$

We now construct a controller  $\pi'$  with control law  $u'_t$  based on the estimator  $\hat{x}'_{u_t}$  by defining the Markov chain

$$x_t - y_t - \hat{x}_{y_t} - u''_t - \hat{x}'_{u_t} - u'_t$$

such that each consecutive pair of variables has the same joint distribution as their unprimed namesakes. Since  $\hat{x}_{y_t}$  is a sufficient statistic of  $y_t$  for  $x_t$ , we have the Markov chain  $x_t - \hat{x}_{y_t} - y_t - u_t$ , implying that  $u''_t$  has the same joint distribution with  $x_t$  as  $u_t$  does. Likewise,  $\hat{x}'_{u_t}$  has the same joint distribution with  $x_t$  as  $\hat{x}_{u_t}$  does. Since  $\hat{x}_{u_t}$  is a sufficient statistic of  $u_t$  for  $x_t$ , we have that  $u'_t$  also has the same joint distribution with  $x_t$  as  $u_t$  does.

Thus the controller  $\pi'$  induces the same stochastic process  $\{x_t, u'_t\}$  and the same external cost. Note that  $u'_t$  may not have the same joint distribution with  $y_t$  as  $u_t$  does and due to the data-processing inequality [3]

$$\begin{aligned} \mathbb{I}[y_t; u_t] &\geq \mathbb{I}[\hat{x}_{y_t}; u_t] = \mathbb{I}[\hat{x}_{y_t}; u''_t] \\ &\geq \mathbb{I}[\hat{x}_{y_t}; \hat{x}'_{u_t}] \geq \mathbb{I}[y_t; u'_t]. \end{aligned}$$

Therefore  $\pi'$  performs at least as well as  $\pi$  and equally well when  $\pi$  is optimal, proving (6).

$\hat{x}'_{u_t}$  is a MMSE estimator for  $\hat{x}_{y_t}$  since

$$\begin{aligned} \mathbb{E}[\hat{x}_{y_t} | \hat{x}'_{u_t}] &= \mathbb{E}[\mathbb{E}[x_t | y_t] | \hat{x}'_{u_t}] \\ &= \mathbb{E}[x_t | \hat{x}'_{u_t}] = \hat{x}'_{u_t}, \end{aligned}$$

where the second equality follows from  $x_t - y_t - \hat{x}'_{u_t}$ .

Finally, it may not be clear from the above analysis that  $u'_t$  is optimally deterministic in  $\hat{x}'_{u_t}$ . If  $u_t$  has covariance  $\Sigma_\nu$  given  $\hat{x}'_{u_t}$ , the Lagrangian of the optimization problem (9) in Part I) depends on  $\Sigma_\nu$  only through the terms

$$\frac{1}{2}(\text{tr}(R\Sigma_\nu) + \text{tr}(SB\Sigma_\nu B^\top)).$$

Since  $R + B^\top SB \succeq 0$  is positive semidefinite, we can take  $\Sigma_\nu = 0$  without loss of performance, recovering the structure (5). Intuitively, the argument is that any noise added to  $u'_t$ , beyond  $\hat{x}'_{u_t}$ , is not helpful in compressing  $x_t$  and can only increase the external cost without saving any communication cost.

In the other direction, let  $u_t$  satisfy the form of Lemma 2. We can rewrite  $u_t$  in the form (III.1), with

$$H = LWK$$

$$\Sigma_\eta = L\Sigma_\omega L^\top. \quad \square$$

### APPENDIX IV

#### PROOF OF THEOREM 1 OF PART I

In this appendix we restate and prove Theorem 1 of Part I [1, Section IV-A], which relies on the following Lagrangian developed there.

$$\begin{aligned} \mathcal{F}_{\Sigma_x, \Sigma_{\hat{x}_u}, L, S; \beta} &= \frac{1}{2}(\beta^{-1}(\log |\Sigma_{\hat{x}_y}|_\dagger - \log |\Sigma_{\hat{x}_y | \hat{x}_u}|_\dagger) \quad (9) \\ &+ \text{tr}(Q\Sigma_x) + \text{tr}(RL\Sigma_{\hat{x}_u}L^\top) \\ &+ \text{tr}(S((A + BL)\Sigma_{\hat{x}_u}(A + BL)^\top \\ &+ A\Sigma_{x|\hat{x}_u}A^\top + \Sigma_\xi - \Sigma_x))). \end{aligned}$$

*Theorem 1:* Given  $\beta$ , the Lagrangian (9) is minimized by a controller satisfying the forward equations

$$\Sigma_x = (A + BL)\Sigma_{\hat{x}_u}(A + BL)^\top \quad (10a)$$

$$+ A\Sigma_{x|\hat{x}_u}A^\top + \Sigma_\xi$$

$$\Sigma_y = C\Sigma_x C^\top + \Sigma_\epsilon \quad (10b)$$

$$K = \Sigma_x C^\top \Sigma_y^\dagger \quad (10c)$$

$$\Sigma_{\hat{x}_y} = K\Sigma_y K^\top, \quad (10d)$$

the backward equations

$$M = \beta^{-1}C^\top K^\top (\Sigma_{\hat{x}_y | \hat{x}_u}^\dagger - \Sigma_{\hat{x}_y}^\dagger) K C \quad (10e)$$

$$S = Q + A^\top S A - M, \quad (10f)$$

$$L = -(R + B^\top S B)^\dagger B^\top S A \quad (10g)$$

$$N = L^\top (R + B^\top S B) L \quad (10h)$$

and the control-based estimator covariance

$$\Sigma_{\hat{x}_u} = \Sigma_{\hat{x}_y}^{1/2} V D V^\top \Sigma_{\hat{x}_y}^{1/2}, \quad (10i)$$

the latter determined by the eigenvalue decomposition (EVD)

$$V \Lambda V^\top = \Sigma_{\hat{x}_y}^{1/2} N \Sigma_{\hat{x}_y}^{1/2} \quad (10j)$$

having  $V$  orthogonal with  $n - \text{rank}(\Sigma_{\hat{x}_y})$  columns spanning the kernel of  $\Sigma_{\hat{x}_y}$  and  $\Lambda = \text{diag}\{\lambda_i\}$  and by the active mode coefficient matrix

$$D = \text{diag} \left\{ \begin{array}{ll} 1 - \beta^{-1} \lambda_i^{-1} & \lambda_i > \beta^{-1} \\ 0 & \lambda_i \leq \beta^{-1} \end{array} \right\}. \quad (10k)$$

*Proof:* The minimum of the Lagrangian (9) must satisfy the first-order optimality conditions, i.e. that the gradient with respect to each parameter is 0 at the optimum. We start by differentiating  $\mathcal{F}$  by the feedback gain  $L$

$$\partial_L \mathcal{F}_{\Sigma_x, \Sigma_{\hat{x}_u}, L, S; \beta} = RL \Sigma_{\hat{x}_u} + B^\top S(A + BL) \Sigma_{\hat{x}_u} = 0,$$

which we rewrite as

$$(R + B^\top SB)L \Sigma_{\hat{x}_u} = -B^\top SA \Sigma_{\hat{x}_u}.$$

As this equation shows,  $L$  is underdetermined in the kernel of  $\Sigma_{\hat{x}_u}$ , since these modes are always 0 in  $\hat{x}_{u_t}$  and have no effect on  $u_t$ .  $L$  is also underdetermined in the kernel of  $R + B^\top SB$ , since these modes have no cost (immediate or future) and can be controlled in any way without affecting the solution's performance. Thus without loss of performance we can take

$$L = -(R + B^\top SB)^\dagger B^\top SA.$$

We substitute this solution back into the Lagrangian, to get

$$\mathcal{F}_{\Sigma_x, \Sigma_{\hat{x}_u}, S; \beta} = \frac{1}{2}(\beta^{-1}(\log |\Sigma_{\hat{x}_y}|_\dagger - \log |\Sigma_{\hat{x}_y|\hat{x}_u}|_\dagger) + \text{tr}(M \Sigma_x) - \text{tr}(N \Sigma_{\hat{x}_u}) + \text{tr}(S \Sigma_\xi)), \quad (\text{IV.1})$$

with

$$\begin{aligned} M &= Q + A^\top SA - S \\ N &= L^\top (R + B^\top SB)L \\ &= A^\top SB(R + B^\top SB)^\dagger B^\top SA. \end{aligned}$$

The problem of optimizing over  $\Sigma_{\hat{x}_u}$  given the other parameters can now be written, up to constants, as the semidefinite program (SDP)

$$\begin{aligned} \max_{\Sigma_{\hat{x}_u}} \quad & \log |\Sigma_{\hat{x}_y} - \Sigma_{\hat{x}_u}|_\dagger + \beta \text{tr}(N \Sigma_{\hat{x}_u}) \\ \text{s.t.} \quad & 0 \preceq \Sigma_{\hat{x}_u} \preceq \Sigma_{\hat{x}_y}. \end{aligned}$$

By Lemma V.1 in Appendix V, the optimum is achieved when  $\Sigma_{\hat{x}_u}$  satisfies (10i)–(10k).

Finally, with  $P = \Sigma_{\hat{x}_y} \Sigma_{\hat{x}_y}^\dagger$  the projection onto the support of  $\hat{x}_{y_t}$  and since the range of  $\Sigma_{\hat{x}_u}$  is contained in that subspace, we have

$$\begin{aligned} \partial_{(\Sigma_x)_{i,j}} (\log |\Sigma_{\hat{x}_y}|_\dagger - \log |\Sigma_{\hat{x}_y|\hat{x}_u}|_\dagger) &= -\partial_{(\Sigma_x)_{i,j}} \log |P - \Sigma_{\hat{x}_u} \Sigma_{\hat{x}_y}^\dagger|_\dagger \\ &= -\partial_{(\Sigma_x)_{i,j}} \log |I - \Sigma_{\hat{x}_u} (P \Sigma_{\hat{x}_y} P)^\dagger| \\ &= \text{tr}((I - \Sigma_{\hat{x}_u} \Sigma_{\hat{x}_y}^\dagger)^{-1} \Sigma_{\hat{x}_u} \partial_{(\Sigma_x)_{i,j}} (P \Sigma_{\hat{x}_y} P)^\dagger). \end{aligned}$$

The purpose of introducing  $P$  is to notice that even if the range of  $\Sigma_{\hat{x}_y}$  is increased, this has no effect on the Lagrangian, because these modes are orthogonal to the range of  $\Sigma_{\hat{x}_u}$ . This allows us to treat  $P$  as constant, so that the range of  $P \Sigma_{\hat{x}_y} P$  is constant in a neighborhood of the solution, and the derivative of the pseudoinverse is simplified in this case to

$$\begin{aligned} \partial_{(\Sigma_x)_{i,j}} (P \Sigma_{\hat{x}_y} P)^\dagger &= -\Sigma_{\hat{x}_y}^\dagger (\partial_{(\Sigma_x)_{i,j}} \Sigma_{\hat{x}_y}) \Sigma_{\hat{x}_y}^\dagger \\ &= -\Sigma_{\hat{x}_y}^\dagger K C J_{i,j} C^\top K^\top \Sigma_{\hat{x}_y}^\dagger, \end{aligned}$$

with  $J_{i,j}$  the matrix with 1 in position  $(i, j)$  and 0 elsewhere. This yields

$$\begin{aligned} \partial_{\Sigma_x} \mathcal{F}_{\Sigma_x, \Sigma_{\hat{x}_u}, S; \beta} &= \frac{1}{2}(M - \beta^{-1} C^\top K^\top \Sigma_{\hat{x}_y}^\dagger (I - \Sigma_{\hat{x}_u} \Sigma_{\hat{x}_y}^\dagger)^{-1} \Sigma_{\hat{x}_u} \Sigma_{\hat{x}_y}^\dagger K C) \\ &= \frac{1}{2}(M - \beta^{-1} C^\top K^\top \Sigma_{\hat{x}_y}^\dagger ((I - \Sigma_{\hat{x}_u} \Sigma_{\hat{x}_y}^\dagger)^{-1} - I) K C) \\ &= \frac{1}{2}(M - \beta^{-1} C^\top K^\top (\Sigma_{\hat{x}_y|\hat{x}_u}^\dagger - \Sigma_{\hat{x}_y}^\dagger) K C) = 0, \end{aligned}$$

implying (10e).  $\square$

## APPENDIX V

### SEMIDEFINITE PROGRAM SOLUTION

In this appendix we state and prove the following solution to our SDP problem.

*Lemma V.1:* The semidefinite program

$$\begin{aligned} \max_{X \in \mathbb{S}_+^n} \quad & \log |M_1 - X|_\dagger + \text{tr}(M_2 X) \\ \text{s.t.} \quad & X \preceq M_1, \end{aligned}$$

with  $M_1, M_2 \succeq 0$ , is optimized by

$$X = M_1^{1/2} V D V^\top M_1^{1/2},$$

with the eigenvalue decomposition (EVD)

$$V \Lambda V^\top = M_1^{1/2} M_2 M_1^{1/2},$$

such that  $V$  is orthogonal with  $n - \text{rank}(M_1)$  columns spanning the kernel of  $M_1$  and  $\Lambda = \text{diag}\{\lambda_i\}$  and with

$$D = \text{diag} \left\{ \begin{array}{ll} 1 - \lambda_i^{-1} & \lambda_i > 1 \\ 0 & \lambda_i \leq 1 \end{array} \right\}.$$

*Proof:* Let the EVD of  $M_1$  be

$$U \Psi U^\top = M_1,$$

with  $U$  orthogonal and  $\Psi$  diagonal, having

$$\Psi = \begin{bmatrix} \Psi_+ & 0 \\ 0 & 0_{(n-m) \times (n-m)} \end{bmatrix},$$

with  $m = \text{rank}(M_1)$ . Let

$$\Psi^\ddagger = \Psi^\dagger + I - \Psi^\dagger \Psi = \begin{bmatrix} \Psi_+^{-1} & 0 \\ 0 & I \end{bmatrix}.$$

By changing the variable to

$$Y = \Psi^{\ddagger/2} U^\top X U \Psi^{\ddagger/2},$$

the constraint of the SDP becomes

$$Y \preceq I_{m,n} = \begin{bmatrix} I_{m \times m} & 0 \\ 0 & 0_{(n-m) \times (n-m)} \end{bmatrix}.$$

$Y$  must therefore be 0 outside the upper-left  $m \times m$  block, and the SDP is equivalent, up to constants, to

$$\begin{aligned} \max_{Y \in \mathbb{S}_+^m} \quad & \log |I_{m,n} - Y|_\dagger + \text{tr}(\Psi^{1/2} U^\top M_2 U \Psi^{1/2} Y) \\ \text{s.t.} \quad & Y \preceq I_{m,n}. \end{aligned}$$

Let the EVD of the linear coefficient be

$$\bar{V} \bar{\Lambda} \bar{V}^\top = \Psi^{1/2} U^\top M_2 U \Psi^{1/2},$$

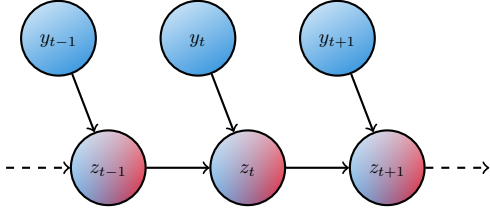


Fig. VI.1. Bayesian network of online inference from a sequence of independent observations

with

$$\bar{V} = \begin{bmatrix} \bar{V}_+ & 0 \\ 0 & I_{(n-m) \times (n-m)} \end{bmatrix}$$

orthogonal and preserving the kernel of  $\Psi$  and  $\Lambda = \text{diag}\{\lambda_i\}$ . We can again change the variable to

$$D = \bar{V}^\top Y \bar{V},$$

to get

$$\begin{aligned} \max_{D \in \mathbb{S}_+^n} \quad & \log |I_{m,n} - D|_{\dagger} + \text{tr}(\Lambda D) \\ \text{s.t.} \quad & D \preceq I_{m,n}, \end{aligned}$$

which can easily be solved using Hadamard's inequality [3], to find

$$D = \text{diag} \left\{ \begin{array}{ll} 1 - \lambda_i^{-1} & \lambda_i > 1 \\ 0 & \lambda_i \leq 1 \end{array} \right\}.$$

Finally, the lemma follows by unmaking the variable changes and taking

$$V = U \bar{V}. \quad \square$$

#### APPENDIX VI

##### PROPERTIES OF THE RETENTIVE DIRECTED INFORMATION

In this appendix we show how the retentive directed information (Definition 6 of Part II [4, Section III-A]) relates to the multi-information of Bayesian networks [5].

Consider the Bayesian network in Figure VI.1, which describes the process of online inference from a sequence of independent observations. The multi-information of this network, for horizon  $T$ , is equal to the retentive directed information

$$\begin{aligned} \mathbb{I}[y^T, z^T] &= \mathbb{E} \left[ \log \frac{f(y^T, z^T)}{\prod_{t=1}^T f(y_t) f(z_t)} \right] \\ &= \sum_{t=1}^T \mathbb{E} \left[ \log \frac{f(z_t | z^{t-1}, y^t)}{f(z_t)} \right] = \mathbb{I}[y^T \rightarrow z^T]. \end{aligned}$$

An important property of the directed information is that the mutual information between two sequences can be decomposed into the sum of directed information in both directions [6]

$$\mathbb{I}[x^T; z^T] = \mathbb{I}[x^T \rightarrow z^T] + \mathbb{I}[z^T \rightarrow x^T].$$

Interestingly, retentive directed information extends this property to the retentive control process (Figure 1 in Part II). This process can be thought of as consisting of four phases:

observation, inference, control and state transition. Its multi-information can accordingly be decomposed [7] into the sum

$$\begin{aligned} \mathbb{I}[x^T, y^T, z^T, u^T] &= \mathbb{I}[x^T \rightarrow y^T] + \mathbb{I}[y^T \rightarrow z^T] \\ &\quad + \mathbb{I}[z^T \rightarrow u^T] + \mathbb{I}[u^T \rightarrow x^T]. \end{aligned}$$

#### APPENDIX VII

##### STRUCTURE OF THE OPTIMAL RETENTIVE CONTROLLER

In this appendix we derive the structure of the optimal retentive controller summarized in Part II [4, Section III-C].

For the structured feedback gain  $L$  we find using the Schur complement that

$$\begin{aligned} (R + B^\top S B)^\dagger &= \begin{bmatrix} R_u + B_{x;u}^\top S_x B_{x;u} & B_{x;u}^\top S_{x;m} \\ S_{m;x} B_{x;u} & S_m \end{bmatrix}^\dagger \\ &= \begin{bmatrix} S_{u|m}^\dagger & -S_{u|m}^\dagger B_{x;u}^\top S_{x;m} S_m^\dagger \\ -S_m^\dagger S_{m;x} B_{x;u} S_{u|m}^\dagger & S_m^\dagger \end{bmatrix}, \end{aligned}$$

with

$$S_{m|u}^\dagger = S_m^\dagger + S_m^\dagger S_{m;x} B_{x;u} S_{u|m}^\dagger B_{x;u}^\top S_{x;m} S_m^\dagger,$$

and so

$$\begin{aligned} L &= -(R + B^\top S B)^\dagger B^\top S A \\ &= -(R + B^\top S B)^\dagger \begin{bmatrix} B_{x;u}^\top S_x A_x & 0 \\ S_{m;x} A_x & 0 \end{bmatrix} \\ &= - \begin{bmatrix} S_{u|m}^\dagger B_{x;u}^\top S_{x|m} A_x & 0 \\ S_m^\dagger S_{m;x} (I - B_{x;u} S_{u|m}^\dagger B_{x;u}^\top S_{x|m}) A_x & 0 \end{bmatrix} \\ &= \begin{bmatrix} L_{u;x|m} & 0 \\ -S_m^\dagger S_{m;x} (A_x + B_{x;u} L_{u;x|m}) & 0 \end{bmatrix}, \end{aligned}$$

with

$$L_{u;x|m} = -S_{u|m}^\dagger B_{x;u}^\top S_{x|m} A_x.$$

We also have

$$\begin{aligned} N &= L^\top (R + B^\top S B) L \\ &= A^\top S B (R + B^\top S B)^\dagger B^\top S A = \begin{bmatrix} N_{x|m} & 0 \\ 0 & 0 \end{bmatrix} \\ N_{x|m} &= \begin{bmatrix} B_{x;u}^\top S_x A_x \\ S_{m;x} A_x \end{bmatrix}^\top L \begin{bmatrix} I \\ 0 \end{bmatrix} \\ &= A_x^\top (S_x - S_{x|m} + S_{x|m} B_{x;u} S_{u|m}^\dagger B_{x;u}^\top S_{x|m}) A_x. \end{aligned}$$

Dually, for the structured Kalman gain  $K$  we find that

$$\begin{aligned} \Sigma_y^\dagger &= \begin{bmatrix} \Sigma_y & \Sigma_{y;m} \\ \Sigma_{m;y} & \Sigma_m \end{bmatrix}^\dagger \\ &= \begin{bmatrix} \Sigma_{y|m}^\dagger & -\Sigma_{y|m}^\dagger \Sigma_{y;m} \Sigma_m^\dagger \\ -\Sigma_m^\dagger \Sigma_{m;y} \Sigma_{y|m}^\dagger & \Sigma_m^\dagger + \Sigma_m^\dagger \Sigma_{m;y} \Sigma_{y|m}^\dagger \Sigma_{y;m} \Sigma_m^\dagger \end{bmatrix}, \end{aligned}$$

and so

$$\begin{aligned} K &= \Sigma_x C^\top \Sigma_y^\dagger \\ &= [\Sigma_x C_{y;x}^\top \quad \Sigma_{x;m}] \begin{bmatrix} \Sigma_y & \Sigma_{y;m} \\ \Sigma_{m;y} & \Sigma_m \end{bmatrix}^\dagger \\ &= [K_{x;y|m} \quad (I - K_{x;y|m} C_{y;x}) \Sigma_{x;m} \Sigma_m^\dagger], \end{aligned}$$

with

$$K_{x;y|m} = \Sigma_{x|m} C_{y;x}^\top \Sigma_{y|m}^\dagger.$$

Now constraining the controller to be MMSE, we have the structure

$$\Sigma_{\hat{x}} = \begin{bmatrix} \Sigma_{x|m} + \Sigma_m & \Sigma_m \\ \Sigma_m & \Sigma_m \end{bmatrix}$$

$$K = \begin{bmatrix} K_{x;y|m} & I - K_{x;y|m} C_{y;x} \end{bmatrix},$$

which we employ in differentiating  $\mathcal{F}$  (IV.1), to get

$$\begin{aligned} \partial_{\Sigma_{x|m}} \mathcal{F}_{\Sigma_{x|m}, \Sigma_m, \Sigma_{\hat{x}_u}, S; \beta} &= \begin{bmatrix} I \\ 0 \end{bmatrix}^\top \partial_{\Sigma_{\hat{x}}} \mathcal{F}_{\Sigma_{\hat{x}}, \Sigma_{\hat{x}_u}, S; \beta} \begin{bmatrix} I \\ 0 \end{bmatrix} \\ &= \frac{1}{2} \begin{bmatrix} I \\ 0 \end{bmatrix}^\top (M - \beta^{-1} C^\top K^\top Z K C) \begin{bmatrix} I \\ 0 \end{bmatrix} \\ &= \frac{1}{2} \left( \begin{bmatrix} I \\ 0 \end{bmatrix}^\top M \begin{bmatrix} I \\ 0 \end{bmatrix} - \beta^{-1} C_{y;x}^\top K_{x;y|m}^\top Z K_{x;y|m} C_{y;x} \right) = 0 \\ \partial_{\Sigma_m} \mathcal{F}_{\Sigma_{x|m}, \Sigma_m, \Sigma_{\hat{x}_u}, S; \beta} &= \begin{bmatrix} I \\ I \end{bmatrix}^\top \partial_{\Sigma_{\hat{x}}} \mathcal{F}_{\Sigma_{\hat{x}}, \Sigma_{\hat{x}_u}, S; \beta} \begin{bmatrix} I \\ I \end{bmatrix} \\ &= \frac{1}{2} \begin{bmatrix} I \\ I \end{bmatrix}^\top (M - \beta^{-1} C^\top K^\top Z K C) \begin{bmatrix} I \\ I \end{bmatrix} \\ &= \frac{1}{2} \left( \begin{bmatrix} I \\ I \end{bmatrix}^\top M \begin{bmatrix} I \\ I \end{bmatrix} - \beta^{-1} Z \right) = 0, \end{aligned}$$

with

$$Z = \Sigma_{\hat{x}_y|\hat{x}_u}^\dagger - \Sigma_{\hat{x}_y}^\dagger.$$

This leaves  $M$  overparameterized and we can choose to give it the structure

$$M = \begin{bmatrix} M_{x|m} + M_m & -M_m \\ -M_m & M_m \end{bmatrix}$$

with

$$M_{x|m} = \beta^{-1} Z$$

$$M_m = \beta^{-1} (C_{y;x}^\top K_{x;y|m}^\top Z K_{x;y|m} C_{y;x} - Z).$$

## REFERENCES

- [1] R. Fox and N. Tishby, "Minimum-information LQG control — Part I: Memoryless controllers," in *Proceedings of the 55th IEEE Conference on Decision and Control (CDC)*, 2016. [Online]. Available: <https://arxiv.org/abs/1606.01946>
- [2] M. Gastpar, B. Rimoldi, and M. Vetterli, "To code, or not to code: Lossy source-channel communication revisited," *IEEE Transactions on Information Theory*, vol. 49, no. 5, pp. 1147–1158, 2003.
- [3] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. Wiley, 2012.
- [4] R. Fox and N. Tishby, "Minimum-information LQG control — Part II: Retentive controllers," in *Proceedings of the 55th IEEE Conference on Decision and Control (CDC)*, 2016. [Online]. Available: <https://arxiv.org/abs/1606.01947>
- [5] N. Friedman, O. Mosenzon, N. Slonim, and N. Tishby, "Multivariate information bottleneck," in *Proceedings of the 17th Conference in Uncertainty in Artificial Intelligence*, 2001, pp. 152–161.
- [6] J. L. Massey, "Causality, feedback and directed information," in *Proceedings of the International Symposium on Information Theory and Its Applications*, 1990, pp. 303–305.
- [7] N. Tishby and D. Polani, "Information theory of decisions and actions," in *Perception-Action Cycle*, ser. Cognitive and Neural Systems. Springer, 2011, pp. 601–636.