

Minimum Entropy Rate Simplification of Stochastic Processes: Supplemental Material

Gustav Eje Henter, *Member, IEEE*, and W. Bastiaan Kleijn, *Fellow, IEEE*



THIS document contains supplemental material for the IEEE Transactions on Pattern Analysis and Machine Intelligence article “Minimum Entropy Rate Simplification of Stochastic Processes.” The supplement is divided into three appendices: the first on MERS for Gaussian processes, and the remaining two on, respectively, the theory and the experimental results of MERS for Markov chains.

A SUPPLEMENTAL DERIVATIONS FOR GAUSSIAN MERS

This section presents detailed derivations of the three solutions to MERS for Gaussian processes presented in Sections 4.1 through 4.3 of the main article, along with arguments for the translation identities in the article Section 4.4.

A.1 Gaussian MERS Solution

We first consider the purely nondeterministic case; the result is easily extended to arbitrary stationary and ergodic Gaussian processes using the Wold decomposition.

Let X and \tilde{X} be two discrete-time stationary and ergodic purely nondeterministic univariate Gaussian processes, with spectral power density functions $R_X(e^{i\omega})$ and $R_{\tilde{X}}(e^{i\omega})$ respectively. These are by necessity positive and defined on $(-\pi, \pi]$. The process \tilde{X} is considered given, and we are trying to find a corresponding X which minimizes the differential entropy rate $h_\infty(X)$ under a differential KL-divergence rate constraint $d_\infty(X \parallel \tilde{X}) \leq d$ and a constraint $R_X(e^{i\omega}) \geq r_{\min}$ on the minimum power at any frequency.

Following [1], the differential entropy rate of X can be calculated as

$$h_\infty(X) = \frac{1}{4\pi} \int_{-\pi}^{\pi} \log(4\pi^2 e^2 R_X(e^{i\omega})) d\omega. \quad (1)$$

Furthermore, the relative differential entropy rate is

$$d_\infty(X \parallel \tilde{X}) = \frac{1}{4\pi} \int_{-\pi}^{\pi} \left(\frac{R_X(e^{i\omega})}{R_{\tilde{X}}(e^{i\omega})} - 1 - \log \frac{R_X(e^{i\omega})}{R_{\tilde{X}}(e^{i\omega})} \right) d\omega. \quad (2)$$

To minimize the differential entropy rate under our divergence constraint we shall use variational calculus in a computation similar to that of Burg [2], but with an added constraint on the minimum power. Fix a Lagrange multiplier $\lambda > 0$ for the divergence constraint and a set (function) of Lagrange multipliers $(4\pi)^{-1} \mu(x) \geq 0$ for the power constraint, and seek the stationary points of the functional

$$\eta_\lambda(R_X(\cdot)) = \frac{1}{4\pi} \int_{-\pi}^{\pi} \log(4\pi^2 e^2 R_X(e^{i\omega})) d\omega + \lambda \frac{1}{4\pi} \int_{-\pi}^{\pi} \left(\frac{R_X(e^{i\omega})}{R_{\tilde{X}}(e^{i\omega})} - 1 - \log \frac{R_X(e^{i\omega})}{R_{\tilde{X}}(e^{i\omega})} \right) d\omega + \int_{-\pi}^{\pi} \frac{1}{4\pi} \mu(\omega) (r_{\min} - R_X(e^{i\omega})) d\omega \quad (3)$$

$$= \frac{1}{4\pi} \int_{-\pi}^{\pi} \left(\log(4\pi^2 e^2 R_X(e^{i\omega})) + \lambda \frac{R_X(e^{i\omega})}{R_{\tilde{X}}(e^{i\omega})} - \lambda - \lambda \log \frac{R_X(e^{i\omega})}{R_{\tilde{X}}(e^{i\omega})} + \mu(\omega) (r_{\min} - R_X(e^{i\omega})) \right) d\omega. \quad (4)$$

-
- G. E. Henter is with the Centre for Speech Technology Research at the University of Edinburgh, United Kingdom. A major portion of this research took place while he was with the Communication Theory laboratory, School of Electrical Engineering at KTH Royal Institute of Technology, Stockholm, Sweden. E-mail: gustav.henter@ee.kth.se.
 - W. B. Kleijn is with the Communications and Signal Processing Group at Victoria University of Wellington, New Zealand, and the Multimedia Computing Group at TU Delft, The Netherlands. E-mail: bastiaan.kleijn@ecs.vuw.ac.nz.
 - This research was supported by the LISTA (Listening Talker) project. The project LISTA acknowledges the financial support of the Future and Emerging Technologies (FET) programme within the Seventh Framework Programme for Research of the European Commission, under FET-Open grant number: 256230.

The Euler-Lagrange equation is obtained by setting the derivative of the integrand with respect to $R_X(e^{i\omega})$ to zero:

$$\frac{\partial}{\partial R_X(e^{i\omega})} \left(\log(4\pi^2 e^2 R_X(e^{i\omega})) + \lambda \frac{R_X(e^{i\omega})}{R_{\tilde{X}}(e^{i\omega})} - \lambda - \lambda \log \frac{R_X(e^{i\omega})}{R_{\tilde{X}}(e^{i\omega})} + \mu(\omega)(r_{\min} - R_X(e^{i\omega})) \right) = 0 \quad (5)$$

$$\frac{\partial}{\partial R_X(e^{i\omega})} \left((1 - \lambda) \log R_X(e^{i\omega}) + \lambda \frac{R_X(e^{i\omega})}{R_{\tilde{X}}(e^{i\omega})} + \lambda \log R_{\tilde{X}}(e^{i\omega}) - \mu(\omega) R_X(e^{i\omega}) \right) = 0 \quad (6)$$

$$\frac{1 - \lambda}{R_X(e^{i\omega})} + \frac{\lambda}{R_{\tilde{X}}(e^{i\omega})} - \mu(\omega) = 0. \quad (7)$$

We can now distinguish two cases, the first being when the minimum-power constraint is not active for a particular ω . In this case complementary slackness requires $\mu(\omega) = 0$, and one obtains

$$\frac{1 - \lambda}{R_X(e^{i\omega})} + \frac{\lambda}{R_{\tilde{X}}(e^{i\omega})} = 0 \quad (8)$$

$$R_X(e^{i\omega}) = \frac{\lambda - 1}{\lambda} R_{\tilde{X}}(e^{i\omega}), \quad (9)$$

which is valid as long as

$$\frac{\lambda - 1}{\lambda} R_{\tilde{X}}(e^{i\omega}) \geq r_{\min}. \quad (10)$$

We see that the MERS simplification of X in this case is a simple scaling (contraction) of the spectral density $R_{\tilde{X}}(e^{i\omega})$ by a factor $\frac{\lambda-1}{\lambda} = \alpha^{-1}$. Obviously, only values $\lambda > 1$ are valid in practice; $\lambda \in (0, 1)$ gives nonsensical, negative spectral powers. In the corresponding interval $\alpha \in (1, \infty)$ we recover all MERS solutions between the original \tilde{X} and fully predictable X , for the case of processes \mathcal{X} with no lower bound on the spectral power ($r_{\min} = 0$).

If $r_{\min} > 0$, there is a risk that the right-hand side of Equation (9) falls below r_{\min} , in which case the minimum-power constraint becomes active, so that $R_X(e^{i\omega}) = r_{\min}$. For such ω , we obtain (see Equation (7))

$$\frac{1 - \lambda}{r_{\min}} + \frac{\lambda}{R_{\tilde{X}}(e^{i\omega})} - \mu(\omega) = 0 \quad (11)$$

$$\mu(\omega) = \frac{\lambda}{R_{\tilde{X}}(e^{i\omega})} + \frac{1 - \lambda}{r_{\min}} \quad (12)$$

$$\mu(\omega) = \frac{\frac{\alpha}{\alpha-1}}{R_{\tilde{X}}(e^{i\omega})} + \frac{1 - \frac{\alpha}{\alpha-1}}{r_{\min}} \quad (13)$$

$$\mu(\omega) = \frac{1}{\alpha - 1} \left((\alpha^{-1} R_{\tilde{X}}(e^{i\omega}))^{-1} - (r_{\min})^{-1} \right), \quad (14)$$

which is strictly positive as long as

$$\alpha^{-1} R_{\tilde{X}}(e^{i\omega}) < r_{\min}. \quad (15)$$

Putting the two cases together, we obtain the solution

$$R_X(e^{i\omega}) = \max(r_{\min}, \alpha^{-1} R_{\tilde{X}}(e^{i\omega})) \quad (16)$$

$$\mu(\omega) = \frac{1}{\alpha - 1} \max\left(0, (\alpha^{-1} R_{\tilde{X}}(e^{i\omega}))^{-1} - (r_{\min})^{-1}\right) \quad (17)$$

for $\alpha \in (1, \infty)$. It is straightforward to verify that this satisfies the requisite optimality criteria. The minimum differential entropy rate achievable is limited by the white noise process X_{\min} defined by

$$R_{X_{\min}}(e^{i\omega}) \equiv r_{\min}, \quad (18)$$

which has

$$h_{\infty}(X_{\min}) = \frac{1}{2} \log(4\pi^2 e^2 r_{\min}). \quad (19)$$

If, on the other hand, $r_{\min} \geq \min_{\omega} R_{\tilde{X}}(e^{i\omega})$, we have $\tilde{X} \notin \mathcal{X}$ (\tilde{X} is not in the feasible set), and the minimum achievable distortion is nonzero.

A.2 Weighted Itakura-Saito Divergence MERS

By multiplying the Itakura-Saito criterion by the observed signal spectrum $R_{\tilde{X}}(e^{i\omega})$, we obtain a weighted dissimilarity measure

$$d_{\text{IS}}^q(X \parallel \tilde{X}) = \frac{1}{4\pi} \int_{-\pi}^{\pi} \left(R_X(e^{i\omega}) - R_{\tilde{X}}(e^{i\omega}) + R_{\tilde{X}}(e^{i\omega}) \log \left(\frac{R_{\tilde{X}}(e^{i\omega})}{R_X(e^{i\omega})} \right) \right) d\omega \quad (20)$$

which emphasizes correctness of spectral peaks more than spectral valleys. This new measure is nonnegative, since it can be written

$$d_{\text{IS}}^q(X \parallel \tilde{X}) = \frac{1}{4\pi} \int_{-\pi}^{\pi} R_{\tilde{X}}(e^{i\omega}) f \left(\frac{R_{\tilde{X}}(e^{i\omega})}{R_X(e^{i\omega})} \right) d\omega, \quad (21)$$

where $f(t) = t - 1 - \log t$ and $R_{\tilde{X}}(e^{i\omega})$ are both nonnegative. The measure evaluates to zero when $R_X(e^{i\omega}) = R_{\tilde{X}}(e^{i\omega})$, and is strictly positive whenever $R_X(e^{i\omega}) \neq R_{\tilde{X}}(e^{i\omega})$ on a set of nonzero measure, assuming $R_{\tilde{X}}(e^{i\omega}) > 0$ everywhere (which is required for ergodicity). Expression (20) is superficially similar to the generalized KL-divergence in [3], but with the signs of the $R_X(e^{i\omega})$ and $R_{\tilde{X}}(e^{i\omega})$ -terms reversed. It is not a Bregman divergence or an f -divergence.

To solve the MERS problem using the weighted Itakura-Saito divergence in (20), we introduce a Lagrange multiplier $\lambda > 0$ for the new dissimilarity constraint, along with $\mu(\omega) \geq 0$ for the lower bound $R_X(e^{i\omega}) \geq r_{\min}$ as before. This yields the functional

$$\begin{aligned} \eta_{\lambda}(R_X(\cdot)) = & \frac{1}{4\pi} \int_{-\pi}^{\pi} \log(4\pi^2 e^2 R_X(e^{i\omega})) d\omega + \int_{-\pi}^{\pi} \mu(\omega) (r_{\min} - R_X(e^{i\omega})) d\omega \\ & + \lambda \frac{1}{4\pi} \int_{-\pi}^{\pi} \left(R_X(e^{i\omega}) - R_{\tilde{X}}(e^{i\omega}) + R_{\tilde{X}}(e^{i\omega}) \log \left(\frac{R_{\tilde{X}}(e^{i\omega})}{R_X(e^{i\omega})} \right) \right) d\omega. \end{aligned} \quad (22)$$

As usual, the Euler-Lagrange equation is obtained by setting the derivative of the integrand with respect to $R_X(e^{i\omega})$ to zero. This produces

$$\frac{\partial}{\partial R_X(e^{i\omega})} \left(\log R_X(e^{i\omega}) + \lambda R_X(e^{i\omega}) - \lambda R_{\tilde{X}}(e^{i\omega}) + \lambda R_{\tilde{X}}(e^{i\omega}) \log \left(\frac{R_{\tilde{X}}(e^{i\omega})}{R_X(e^{i\omega})} \right) \right) = 0 \quad (23)$$

$$\frac{1}{R_X(e^{i\omega})} + (\lambda - \mu(\omega)) - \lambda \frac{R_{\tilde{X}}(e^{i\omega})}{R_X(e^{i\omega})} = 0 \quad (24)$$

$$\frac{1}{\lambda} + \left(1 - \frac{1}{\lambda} \mu(\omega) \right) R_X(e^{i\omega}) - R_{\tilde{X}}(e^{i\omega}) = 0, \quad (25)$$

since $R_X(e^{i\omega})$ and λ are both nonnegative.

We once again must consider two distinct cases. The first is that the lower-bound constraint is not active for a particular ω , $\mu(\omega) = 0$, and so

$$\frac{1}{\lambda} + R_X(e^{i\omega}) - R_{\tilde{X}}(e^{i\omega}) = 0 \quad (26)$$

$$R_X(e^{i\omega}) = R_{\tilde{X}}(e^{i\omega}) - \frac{1}{\lambda}, \quad (27)$$

which is valid whenever

$$R_{\tilde{X}}(e^{i\omega}) - \frac{1}{\lambda} \geq r_{\min}. \quad (28)$$

This is a simple subtraction from the power spectral density. For processes with no lower bound on the spectral power (i.e., $r_{\min} = 0$), we see that $R_X(e^{i\omega}) \rightarrow R_{\tilde{X}}(e^{i\omega})$ as $\lambda^{-1} \rightarrow 0$ (no simplification), while the entropy rate of X approaches zero as λ^{-1} increases towards $\min_{\omega} R_{\tilde{X}}(e^{i\omega})$ (complete predictability, under certain continuity assumptions on $R_{\tilde{X}}(e^{i\omega})$).

In the alternative case, the condition in Equation (28) is violated, in which case the lower bound gives $R_X(e^{i\omega}) = r_{\min}$. Slotting this into Equation (25) and assuming $r_{\min} > 0$, we obtain

$$\frac{1}{\lambda} + \left(1 - \frac{1}{\lambda} \mu(\omega) \right) r_{\min} - R_{\tilde{X}}(e^{i\omega}) = 0 \quad (29)$$

$$-r_{\min} \frac{1}{\lambda} \mu(\omega) = R_{\tilde{X}}(e^{i\omega}) - \frac{1}{\lambda} - r_{\min} \quad (30)$$

$$\mu(\omega) = \frac{\lambda}{r_{\min}} \left(r_{\min} + \frac{1}{\lambda} - R_{\tilde{X}}(e^{i\omega}) \right), \quad (31)$$

which is strictly positive whenever

$$R_{\tilde{X}}(e^{i\omega}) - \frac{1}{\lambda} < r_{\min}. \quad (32)$$

Combining the two cases, we arrive at the solution

$$R_X(e^{i\omega}) = \max(r_{\min}, R_{\tilde{X}}(e^{i\omega}) - \lambda^{-1}) \quad (33)$$

$$\mu(\omega) = \frac{\lambda}{r_{\min}} \max(0, r_{\min} + \lambda^{-1} - R_{\tilde{X}}(e^{i\omega})) \quad (34)$$

for $r_{\min} > 0$, valid for all $\lambda \in (0, \infty)$. (For $r_{\min} = 0$ we have the previous solution derived in Equation (27).) Exactly like the general solution in Section A.1, there is a nonzero minimum rate $h_{\infty}(X_{\min})$, and the minimum distortion can be nonzero as well.

A.3 Variance-Constrained Gaussian MERS

Another interesting variation of MERS is obtained by constraining the standard relative entropy rate (regular Itakura-Saito divergence), but preventing the X -process from simply shrinking away by only considering processes that satisfy the additional variance constraint

$$\text{Var}(X_t) = \text{Var}(\tilde{X}_t). \quad (35)$$

For Gaussian processes the variance may be written in terms of the spectral density function as

$$\text{Var}(X_t) = \frac{1}{2\pi} \int_{-\pi}^{\pi} R_X(e^{i\omega}) d\omega, \quad (36)$$

so the constraint (35) may be expressed

$$\frac{1}{2\pi} \int_{-\pi}^{\pi} (R_X(e^{i\omega}) - R_{\tilde{X}}(e^{i\omega})) d\omega = 0. \quad (37)$$

Introducing an additional Lagrange multiplier μ for the variance leads to a modified functional

$$\begin{aligned} \eta_{\lambda, \mu}(R_X(\cdot)) = \frac{1}{4\pi} \int_{-\pi}^{\pi} \left(\log(4\pi^2 e^2 R_X(e^{i\omega})) + \lambda \frac{R_X(e^{i\omega})}{R_{\tilde{X}}(e^{i\omega})} - \lambda - \lambda \log \frac{R_X(e^{i\omega})}{R_{\tilde{X}}(e^{i\omega})} \right) d\omega \\ + 2\mu \frac{1}{4\pi} \int_{-\pi}^{\pi} (R_X(e^{i\omega}) - R_{\tilde{X}}(e^{i\omega})) d\omega. \end{aligned} \quad (38)$$

We again derive the Euler-Lagrange equation as

$$\frac{\partial}{\partial R_X(e^{i\omega})} \left((1 - \lambda) \log R_X(e^{i\omega}) + \lambda \frac{R_X(e^{i\omega})}{R_{\tilde{X}}(e^{i\omega})} + 2\mu R_X(e^{i\omega}) \right) = 0 \quad (39)$$

$$\frac{1 - \lambda}{R_X(e^{i\omega})} + \frac{\lambda}{R_{\tilde{X}}(e^{i\omega})} + 2\mu = 0 \quad (40)$$

$$\frac{\lambda - 1}{R_X(e^{i\omega})} = \frac{2\mu R_{\tilde{X}}(e^{i\omega}) + \lambda}{R_{\tilde{X}}(e^{i\omega})}, \quad (41)$$

which gives

$$R_X(e^{i\omega}) = \frac{(\lambda - 1) R_{\tilde{X}}(e^{i\omega})}{2\mu R_{\tilde{X}}(e^{i\omega}) + \lambda} \quad (42)$$

$$= \frac{1}{\nu} \frac{R_{\tilde{X}}(e^{i\omega})}{R_{\tilde{X}}(e^{i\omega}) + \frac{\alpha}{\nu}}. \quad (43)$$

The new Lagrange multipliers $\nu = \frac{2\mu}{\lambda - 1}$ and α here can take on any nonzero real values, though the result only makes sense if $\nu R_{\tilde{X}}(e^{i\omega}) + \alpha > 0$ everywhere.

By considering the behavior of the function $f(r) = \frac{1}{\nu} \frac{r}{r + c}$ for a subset $r \in R = [R_{\min}, R_{\max}] \subset (0, \infty)$ of the positive real line, we see that we that $\nu < 0, c < -R_{\max}$ corresponds to positive and increasing concave $f(r)$ with a derivative $f'(r) = \frac{1}{\nu} \frac{c}{(r + c)^2}$ greater than $-\frac{1}{\nu} > 0$. This implies that spectral peaks are emphasized, similar to reverse water-filling.

The region $\nu, c > 0$, meanwhile, corresponds to a positive but concave f with slope between zero and $\frac{1}{\nu}$. This mapping smooths out spectral peaks and instead maximizes entropy under the constraints.

A.4 Translation Invariance

The fact that $h_\infty(\mu + X) = h_\infty(X)$ when μ is deterministic is theorem 2.4.1 in [1]. For

$$d_\infty(\mu + X \parallel \tilde{\mu} + \tilde{X}) \geq d_\infty(\tilde{\mu} + X \parallel \tilde{\mu} + \tilde{X}), \quad (44)$$

it suffices to prove that

$$d_\infty(\Delta\mu + X \parallel \tilde{X}) \geq d_\infty(X \parallel \tilde{X}), \quad (45)$$

where $\Delta\mu = \mu - \tilde{\mu}$. Since $\Delta\mu + X$ and X have the same spectral density functions (and covariance matrices), the general theorem 2.4.5 from [1] applies. This establishes

$$d_\infty(\Delta\mu + X \parallel \tilde{X}) = d_\infty(\Delta\mu + X \parallel X) + d_\infty(X \parallel \tilde{X}) \quad (46)$$

$$\geq d_\infty(X \parallel \tilde{X}), \quad (47)$$

which is the desired result. Equality occurs when $d_\infty(\Delta\mu + X \parallel X) = 0 \Leftrightarrow \Delta\mu + X = X \Leftrightarrow \mu = \tilde{\mu}$.

B SUPPLEMENTAL DERIVATIONS FOR MARKOV CHAIN MERS

This section presents an alternative representation of the MERS problem for Markov chains, and subsequently shows how the optimal solution presented in Section 5.1 of the main article can be derived from this representation.

B.1 Bigram Matrix MERS Formulation

The MERS problem for first-order Markov chains is to find a stationary, ergodic Markov chain X that minimizes the entropy rate $H_\infty(X)$ under a divergence constraint $D_\infty(X \parallel \tilde{X}) \leq D$. Here, \tilde{X} is another, given stationary and ergodic Markov chain over the same state space as X . We may represent X and \tilde{X} by their transition matrices \mathbf{A} and $\tilde{\mathbf{A}}$ with elements $a_{ij} = P(X_{t+1} = j \mid X_t = i)$, and similarly for $\tilde{\mathbf{A}}$. We also have the stationary distribution vector $\boldsymbol{\pi}$ with elements $\pi_i = P(X_t = i) > 0$, as described in the article.

With our notation, the entropy rate and relative entropy rate of a first order Markov chain X can be expressed as

$$H_\infty(X) = - \sum_{ij} \pi_i a_{ij} \log a_{ij} \quad (48)$$

$$D_\infty(X \parallel \tilde{X}) = \sum_{ij} \pi_i a_{ij} \log \frac{a_{ij}}{\tilde{a}_{ij}}. \quad (49)$$

However, these formulas are cumbersome in practice since they involve $\boldsymbol{\pi}$, which is a normalized version of the leading eigenvector of \mathbf{A}^\top and thus a complicated function of the matrix elements a_{ij} . A better option is to represent the Markov chain X by its *bigram probability matrix* \mathbf{B} , having elements $b_{ij} = P(X_t = i \wedge X_{t+1} = j)$. \mathbf{B} is related to the transition matrix through the formula $\mathbf{B} = \text{diag}(\boldsymbol{\pi}) \mathbf{A}$, a matrix version of the well-known relationship $P(X_t = i \wedge X_{t+1} = j) = P(X_t = i) P(X_{t+1} = j \mid X_t = i)$. Standard requirements for probability distributions give $\mathbf{B} \geq \mathbf{0}$ and $\mathbf{1}^\top \mathbf{B} \mathbf{1} = 1$. The fact that the process is stationary, so that $\pi_i = P(X_t = i) = P(X_{t+1} = i)$, also contributes a constraint $\mathbf{B} \mathbf{1} = \mathbf{B}^\top \mathbf{1}$ on the marginal distributions. It is easy to see that there is a one-to-one correspondence between \mathbf{A} and \mathbf{B} in our case.

Using the bigram probability matrix representation, the MERS problem for first-order Markov chains can be written as

$$\min_{\mathbf{B} \in \mathbb{R}^{k \times k}} - \sum_{ij} b_{ij} \log \frac{b_{ij}}{\sum_{j'} b_{ij'}} \quad (50)$$

subject to

$$\sum_{ij} b_{ij} \log \frac{b_{ij}}{\tilde{a}_{ij} \sum_{j'} b_{ij'}} \leq D \quad (51)$$

$$\mathbf{1}^\top \mathbf{B} \mathbf{1} = 1 \quad (52)$$

$$\mathbf{B} \mathbf{1} - \mathbf{B}^\top \mathbf{1} = \mathbf{0} \quad (53)$$

$$\mathbf{B} \geq \mathbf{0}. \quad (54)$$

This formulation is preferable for practical purposes since it does not include any implicit functions of the variables b_{ij} .

B.2 Markov Chain MERS Solution

General nonconvex optimization problems like the above can be demanding to solve. However, it is here possible to solve the problem analytically, using a technique derived from the Blahut-Arimoto algorithm [4], [5] in rate-distortion theory. The first step is to fix a Lagrange multiplier $\lambda \geq 0$ for the dissimilarity constraint (51), and minimize the combined objective function

$$f_\lambda(\mathbf{B}) = (\lambda - 1) \sum_{ij} b_{ij} \log \frac{b_{ij}}{\sum_{j'} b_{ij'}} - \lambda \sum_{ij} b_{ij} \log \tilde{a}_{ij} - \lambda D \quad (55)$$

under the remaining constraints. By considering λ fixed, we define a particular trade-off between simplicity and divergence, corresponding to a point on the rate-distortion curve. Like for the Gaussian case, we shall see that all relevant solutions correspond to $\lambda > 1$, where we can define the factor $\alpha = \frac{\lambda}{\lambda-1} \in (1, \infty)$ as before. In this region we may divide by $\lambda - 1$, and instead maximize

$$f_\alpha(\mathbf{B}) = \sum_{ij} b_{ij} \log b_{ij} - \sum_{ij} b_{ij} \log \sum_{j'} b_{ij'} - \alpha \sum_{ij} b_{ij} \log \tilde{a}_{ij} - \alpha D. \quad (56)$$

Next, we use the same trick as for the Blahut-Arimoto algorithm, transforming the problem to a minimization over two sets of variables, where each set is straightforward to optimize if the other is fixed. We do this by noting that

$$\max_{\mathbf{q} \geq \mathbf{0}} \sum_{ij} b_{ij} \log q_i \text{ subject to } \mathbf{1}^\top \mathbf{q} = 1 \quad (57)$$

is solved by $\mathbf{q}^*(\mathbf{B}) = \mathbf{B}\mathbf{1}$ for nonnegative \mathbf{B} . (Superscripted stars here signify optimal solutions, and are not related to the underlying process X^* defined within the main article text.) This lets us replace the difficult term $b_{ij} \log \sum_{j'} b_{ij'}$, containing a sum within a logarithm, and instead define an augmented objective function

$$f_\alpha(\mathbf{B}, \mathbf{q}) = \sum_{ij} b_{ij} \log b_{ij} - \sum_{ij} b_{ij} \log q_i - \alpha \sum_{ij} b_{ij} \log \tilde{a}_{ij}. \quad (58)$$

Maximizing $f_\alpha(\mathbf{B}, \mathbf{q})$ subject to constraints (52) through (54), as well as $\mathbf{q} \geq \mathbf{0}$ and $\mathbf{1}^\top \mathbf{q} = 1$, yields a \mathbf{B} that is also optimal for the original problem.

We already know how to optimize $f_\alpha(\mathbf{B}, \mathbf{q})$ for a given, constant \mathbf{B} . If we instead fix \mathbf{q} and introduce real Lagrange multipliers $(\log \nu - 1)$ (for the sum constraint (52)) and $\log \mu_i$ (for the stationarity constraints (53)), with ν and all μ_i positive, we can find the optimal $\mathbf{B}^*(\mathbf{q})$ from the stationary points of the function

$$f_{\alpha, \mu, \nu}(\mathbf{B}, \mathbf{q}) = \sum_{ij} b_{ij} \log b_{ij} - \sum_{ij} b_{ij} \log q_i - \alpha \sum_{ij} b_{ij} \log \tilde{a}_{ij} + (\log \nu - 1) \sum_{ij} b_{ij} - (\log \nu - 1) + \sum_i \log \mu_i \sum_j (b_{ij} - b_{ji}). \quad (59)$$

The stationary points satisfy $\nabla_{\mathbf{B}} f_{\alpha, \mu, \nu}(\mathbf{B}^*, \mathbf{q}) = \mathbf{0}$, giving

$$\frac{\partial}{\partial b_{ij}} f_{\alpha, \mu, \nu}(\mathbf{B}^*, \mathbf{q}) = 0 \quad (60)$$

$$\log \frac{b_{ij}^*}{q_i} + 1 - \alpha \log \tilde{a}_{ij} + \log \nu - 1 + \log \mu_i - \log \mu_j = 0 \quad (61)$$

$$\frac{b_{ij}^*}{q_i} (\tilde{a}_{ij})^{-\alpha} \nu \frac{\mu_i}{\mu_j} = 1 \quad (62)$$

$$b_{ij}^*(\mathbf{q}) = \frac{1}{\nu} \frac{\mu_j}{\mu_i} q_i (\tilde{a}_{ij})^\alpha. \quad (63)$$

This solution naturally satisfies nonnegativity, with ν being a simple normalization. The vector μ is more difficult to compute, and will be derived later. For now, we simply assume μ known.

The formulas for $\mathbf{q}^*(\mathbf{B})$ and $\mathbf{B}^*(\mathbf{q})$ can, given an initial guess $\mathbf{B}^{(0)}$, be used for alternating updates of \mathbf{q} and \mathbf{B} guaranteed to decrease $f_{\alpha, \mu, \nu}$ unless a fixed point or a cycle is reached. At a fixed point of these iterations we have $q_i^* = \sum_j b_{ij}^*$. Slotting this into (63) yields

$$a_{ij}^* = \frac{b_{ij}^*}{\sum_{j'} b_{ij'}^*} = \frac{1}{\nu} \frac{\mu_j}{\mu_i} (\tilde{a}_{ij})^\alpha. \quad (64)$$

This can be written on matrix form as

$$\mathbf{A}^* = \frac{1}{\nu} (\text{diag } \mu)^{-1} \tilde{\mathbf{A}}^{(\alpha)} (\text{diag } \mu) \quad (65)$$

where $\tilde{\mathbf{A}}^{(\alpha)}$ denotes *Hadamard power* α of $\tilde{\mathbf{A}}$, the matrix resulting from taking all entries of $\tilde{\mathbf{A}}$ to the power α , i.e., $(\tilde{\mathbf{A}}^{(\alpha)})_{ij} = (\tilde{a}_{ij})^\alpha$.

To solve for $\boldsymbol{\mu}$ and ν we use the fact that the optimal transition matrix \mathbf{A} must satisfy the standard relation $\mathbf{A}\mathbf{1} = \mathbf{1}$ in order to properly define a Markov chain X . Hence

$$\frac{1}{\nu} (\text{diag } \boldsymbol{\mu})^{-1} \tilde{\mathbf{A}}^{(\alpha)} (\text{diag } \boldsymbol{\mu}) \mathbf{1} = \mathbf{1} \quad (66)$$

$$\tilde{\mathbf{A}}^{(\alpha)} \boldsymbol{\mu} = \nu \boldsymbol{\mu}, \quad (67)$$

so $\boldsymbol{\mu}$ is a right eigenvector of $\tilde{\mathbf{A}}^{(\alpha)}$ associated with the eigenvalue ν . In addition, we have the requirement that ν and μ_i are all positive, from before; this guarantees that all required inverses exist.

Since the original Markov chain \tilde{X} is ergodic, $\tilde{\mathbf{A}}$ and hence $\tilde{\mathbf{A}}^{(\alpha)}$ are irreducible. As the matrix elements are nonnegative, aperiodicity of $\tilde{\mathbf{A}}$ also establishes that $\tilde{\mathbf{A}}^{(\alpha)}$ is aperiodic. The Perron-Frobenius theorem then ensures that there is a unique eigenvector of $\tilde{\mathbf{A}}^{(\alpha)}$ with only real, positive elements, and that it is associated with the greatest eigenvalue, which is also real. This establishes that the unique optimal solution to the original MERS problem is the Markov chain X with transition matrix

$$\mathbf{A} = \frac{1}{\nu} (\text{diag } \boldsymbol{\mu})^{-1} \tilde{\mathbf{A}}^{(\alpha)} (\text{diag } \boldsymbol{\mu}) \quad (68)$$

where ν is the leading eigenvalue of the Hadamard power $\tilde{\mathbf{A}}^{(\alpha)}$ and $\boldsymbol{\mu}$ is the corresponding right eigenvector with only positive entries. We see that exponent values $\alpha \in (1, \infty)$ recover all solutions between the original \tilde{X} and complete predictability.

As an aside, Equation (68) shows that $\tilde{\mathbf{A}}^{(\alpha)}$ and \mathbf{A} have the same eigenvalues and eigenvectors, except for a scaling by ν and $\text{diag } \boldsymbol{\mu}$.

C SUPPLEMENTAL MATERIAL FOR MARKOV CHAIN EXPERIMENTS

This section presents, in Table 1 on page 8, excerpts from Swedish text data, along with sampled output from word-level Markov models based on this material. Errors, shown in red with wavy underline, are seen to decrease in frequency in the output of models simplified with MERS. Table 2 in Section 6.4 of the main article was derived by word-for-word translation of these Swedish text examples. Please consult the main article for more information.

REFERENCES

- [1] S. Ihara, *Information Theory for Continuous Systems*. World Scientific Publishing Company, 1993.
- [2] J. P. Burg, *A New Analysis Technique for Time Series Data*, pp. 42–48. IEEE Press, 1978.
- [3] D. D. Lee and H. S. Seung, “Learning the parts of objects by non-negative matrix factorization,” *Nature*, vol. 401, pp. 788–791, Oct. 1999.
- [4] R. E. Blahut, “Computation of channel capacity and rate distortion function,” *IEEE Trans. Inf. Theory*, vol. IT-18, pp. 460–473, 1972.
- [5] S. Arimoto, “An algorithm for calculating the capacity of an arbitrary discrete memoryless channel,” *IEEE Trans. Inf. Theory*, vol. IT-18, pp. 14–20, 1972.

ett bad [paus] titta mamman [paus] ta en sko [paus]
 det är pappan [paus] titta olov [paus] ta ett bad
 [paus] var är olov nu [paus] ta ett bad [paus] var
 är olov nu [paus] det är en bil [paus] ta en sko
 [paus] ta en bil [paus] det är en bil [paus] titta
 bilen [paus] titta skon [paus] det är en bok [paus]
 ta ett bad [paus] var är badet nu [paus] titta bilen
 [paus] titta olov [paus] var är bilen nu [paus]
 titta skon [paus] ta en bil [paus] hej olov [paus]
 hej olov [paus] var är olov

(a) Sample from error-free corpus.

en bok [paus] hallå olov [paus] var är olov [paus]
 var är bilen nu [paus] hallå olov [paus] det är en
 bok [paus] det är olov [paus] titta pappan [paus]
 titta badet [paus] det är en bok [paus] var är skon
 nu [paus] ta en bil [paus] hej olov [paus] ta en sko
 [paus] det är en bok [paus] var är badet nu [paus]
 titta bilen [paus] titta bilen [paus] det är olov
 [paus] det är en bok [paus] ta en bil [paus] titta
 mamman [paus] hej olov [paus] ta en bil [paus] titta
 mamman [paus] titta skon [paus] titta

(c) Sample from X^* fit to error-free data.

hej olov [paus] hallå olov [paus] ööh ta en bil
 [paus] det är olov nu [paus] ta en bok [paus] det
 är pappan nu [paus] ta en en bok [paus] ta en bil
 [paus] det är en sko [paus] var är ööh en bil [paus]
 hej olov [paus] det är en bil [paus] ta ett en sko
 [paus] det är en bil [paus] ta en bil [paus] ta ett
 bad [paus] det är en bok [paus] ööh hej olov [paus]
 var är olov nu [paus] det är ett bad [paus] ööh ta
 en bok var är skon nu [paus] det

(e) Sample from MERS $X(R)$ at optimum denoising.

var är badet nu [paus] var är skon nu [paus] var är boken nu [paus] det är en bil [paus] det är en bil
 [paus] var är badet nu [paus] var är skon nu [paus] det är en sko [paus] det är en sko [paus] det är en
 bok [paus] var är bilen nu [paus] var är badet nu [paus] det är en bok [paus] det är ett bad [paus] var är
 boken nu [paus] var är badet nu [paus] det är ett bad [paus] det är ett bad [paus] det är ett bad [paus]
 det är en bok [paus]

(g) Sample from MERS $X(R)$ at low rate ($\alpha = 46$, $R \approx 0.41$).

Table 1

Excerpts from the clean and disturbed data, along with random samples from the fitted models X^* and \tilde{X} , the optimally denoised models from MERS and thresholding, and from low entropy-rate MERS (100 symbols each).

pappan nu hej [paus] olov [paus] [paus] hallå olov
 [paus] [paus] det är pappan [paus] titta ööh skon
 [paus] [paus] det är ööh olov [paus] [paus] hallå
 olov [paus] titta ööh boken [paus] [paus] var är
 olov nu ööh ta en bok [paus] ööh ööh hallå olov
 [paus] det är olov det är ööh en bil [paus] hallå
 olov [paus] ööh var är olov nu [paus] ööh var är
 olov nu [paus] hallå olov [paus] [paus] ta ett bad
 [paus] titta boken boken hej olov [paus] ööh titta
 mamman [paus] det är en bil [paus] ööh ta en bil
 hallå

(b) Sentences disturbed by random speech errors.

badet [paus] hej ööh olov [paus] ta en bok [paus]
 ta en sko en sko titta boken en bil [paus] det är
 pappan nu [paus] [paus] var är skon nu [paus] det är
 olov [paus] det är en sko [paus] hej olov [paus] ööh
 ta en sko [paus] var är boken nu [paus] det är en
 bil [paus] ööh ta ett bad ett bad [paus] titta bilen
 [paus] ööh titta bilen [paus] olov [paus] titta skon
 [paus] titta pappan var är ööh ett en bil [paus]
 titta olov det är ett bad [paus] var är är en bil
 [paus] [paus]

(d) Sample from \tilde{X} fit to corrupted data.

är ett bad [paus] titta mamman [paus] det är en bil
 [paus] det är [paus] pappan nu [paus] det är olov
nu [paus] hej olov [paus] hallå olov [paus] var är
 badet nu [paus] det är mamman nu [paus] ta en sko
 [paus] [paus] ta en bok [paus] hej olov [paus] det
 är mamman [paus] ööh det är ett bad [paus] ta ett
 bad [paus] hej olov [paus] ööh ta en bil [paus]
 hallå olov [paus] titta pappan [paus] ta en bok
 [paus] titta boken [paus] titta badet [paus] det är
 mamman nu [paus] [paus] hej olov [paus] titta boken

(f) Sample from thresholded $X'(R)$ at optimum denoising.