

Prediction-based Similarity Identification for Autoregressive Processes

Hanwei Wu, Qiwen Wang and Markus Flierl

KTH Royal Institute of Technology, Stockholm

School of Electrical Engineering and Computer Science



Introduction

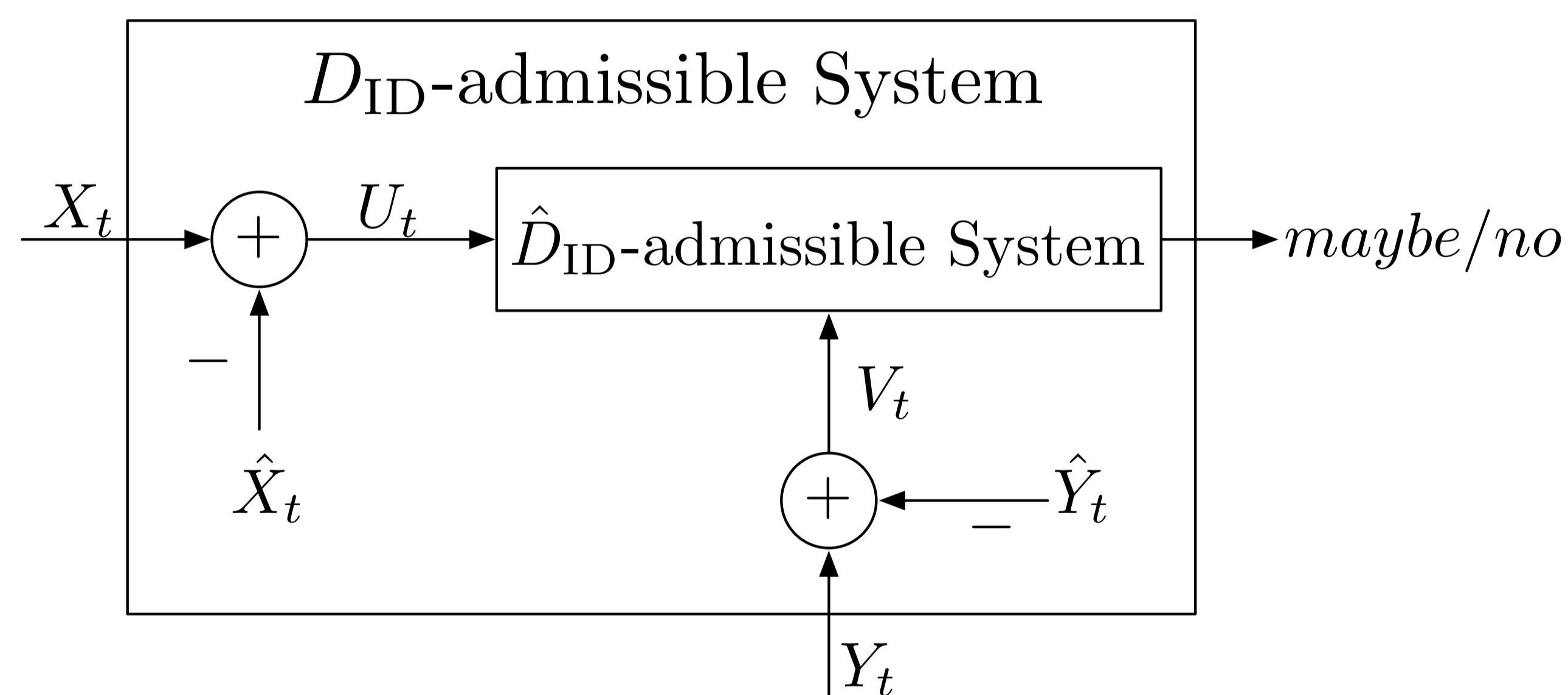
Problem

- Compression of autoregressive processes for similarity identification [1].

Goals

- Use prediction-based model to process autoregressive signals for similarity identification.
- Derive the identification rate of the autoregressive processes using prediction-based model.

Prediction-based Approach



- Database and the query are zero-mean autoregressive processes:

$$X_t = U_t + \mathbf{a}_m^T \mathbf{X}_{t-1}^{(m)} \text{ and } Y_t = V_t + \mathbf{a}_m^T \mathbf{Y}_{t-1}^{(m)},$$

- Optimal predictor for autoregressive processes:

$$\hat{X}_t = \mathbf{a}_{N_p}^T \mathbf{X}_{t-1}^{(N_p)}, \text{ where } \mathbf{a}_{N_p} = (a_1, \dots, a_m, 0, \dots, 0)^T.$$

- Similarity identification is conducted in the embedded \hat{D}_{ID} -admissible system.

Identification rate R_{ID}^P

- Vector representation of the autoregressive process:

$$\mathbf{x} = \mathbf{M}_t \mathbf{u} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_t \end{bmatrix} = \begin{bmatrix} 1 & & & \\ m_1 & 1 & & \\ & m_1 & 1 & \\ \vdots & & \ddots & \\ m_{t-1} & \dots & m_1 & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_t \end{bmatrix}.$$

- Relation of the similarity measures : data space v.s. residual space

$$d(\mathbf{x}, \mathbf{y}) = \frac{1}{t} \sum_{i=1}^t \lambda_{t,i} d(\tilde{u}_i, \tilde{v}_i),$$

where $\tilde{u} = \mathbf{Q}\mathbf{u}$, $\tilde{v} = \mathbf{Q}\mathbf{v}$, \mathbf{Q} and $\lambda_{t,i}$ are the eigenmatrix and eigenvalues of $\mathbf{P}_t = \mathbf{M}_t^T \mathbf{M}_t$.

- Assume for each time step an ideal identification system for Gaussian data, i.e., $d(\tilde{u}_i, \tilde{v}_i) \leq \tilde{D}_{ID}^{(i)}$ with identification rate $\tilde{R}_{ID}^{(i)}$. The identification rate based on the prediction-model is obtained by the following constrained optimization:

$$\begin{aligned} \max_{\tilde{D}_{ID}^{(1)}, \dots, \tilde{D}_{ID}^{(t)}} D_{ID} &= \frac{1}{t} \sum_{i=1}^t \lambda_{t,i} \tilde{D}_{ID}^{(i)} \\ \text{s.t. } \frac{1}{t} \sum_{i=1}^t \tilde{R}_{ID}^{(i)} &\leq R_{ID}, \\ \text{s.t. } \tilde{R}_{ID}^{(i)} &\geq 0. \end{aligned}$$

- Identification rate of the prediction-based model for autoregressive Gaussian processes is given by

$$\begin{aligned} R_{ID}^P &= \frac{1}{t} \sum_{i=1}^t \max \left\{ \log_2 \left(\frac{2 \ln(2) \lambda_{t,i}}{v} \right), 0 \right\} \\ D_{ID}^P &= \frac{1}{t} \sum_{i=1}^t \lambda_{t,i} 2 \left(1 - 2^{-R_{ID}^{(i)}} \right). \end{aligned}$$

Special Case

R_{ID}^{PS} for Autoregressive Processes

- Only the smallest eigenvalue of $\mathbf{P}_t = \mathbf{M}_t^T \mathbf{M}_t$ is known.
- Relation of the similarity measures : data space v.s. residual space

$$d(\mathbf{x}, \mathbf{y}) \geq \lambda_{\min} d(\mathbf{u}, \mathbf{v})$$

- Similarity threshold for embedded similarity identification system

$$d(\mathbf{u}, \mathbf{v}) \leq \frac{d(\mathbf{x}, \mathbf{y})}{\lambda_{\min}} \leq \frac{D_{ID}}{\lambda_{\min}} := \hat{D}_{ID}.$$

- Identification rate of Gaussian autoregressive processes for the special case

$$R_{ID}^{PS} = \log_2 \left(\frac{2 \lambda_{\min}}{2 \lambda_{\min} - D_{ID}} \right).$$

Asymptotic Upper Bound of R_{ID}^{PS}

- \mathbf{P}_t is asymptotically equivalent to a Toeplitz matrix $T_t(g)$, where $g(\omega)$ is

$$g(\omega) = \sum_{k=-\infty}^{\infty} e^{-jk\omega} \left\{ \sum_{i=0}^{\infty} m_i m_{i+k} \right\} = \left| \sum_{k=0}^{\infty} m_k e^{-jk\omega} \right|^2.$$

- The minimum eigenvalue of a Toeplitz matrix converges to the essential infimum N_l of $g(\omega)$

$$\lim_{t \rightarrow \infty} \min_i \tau_{t,i} = N_l.$$

Simulation Results

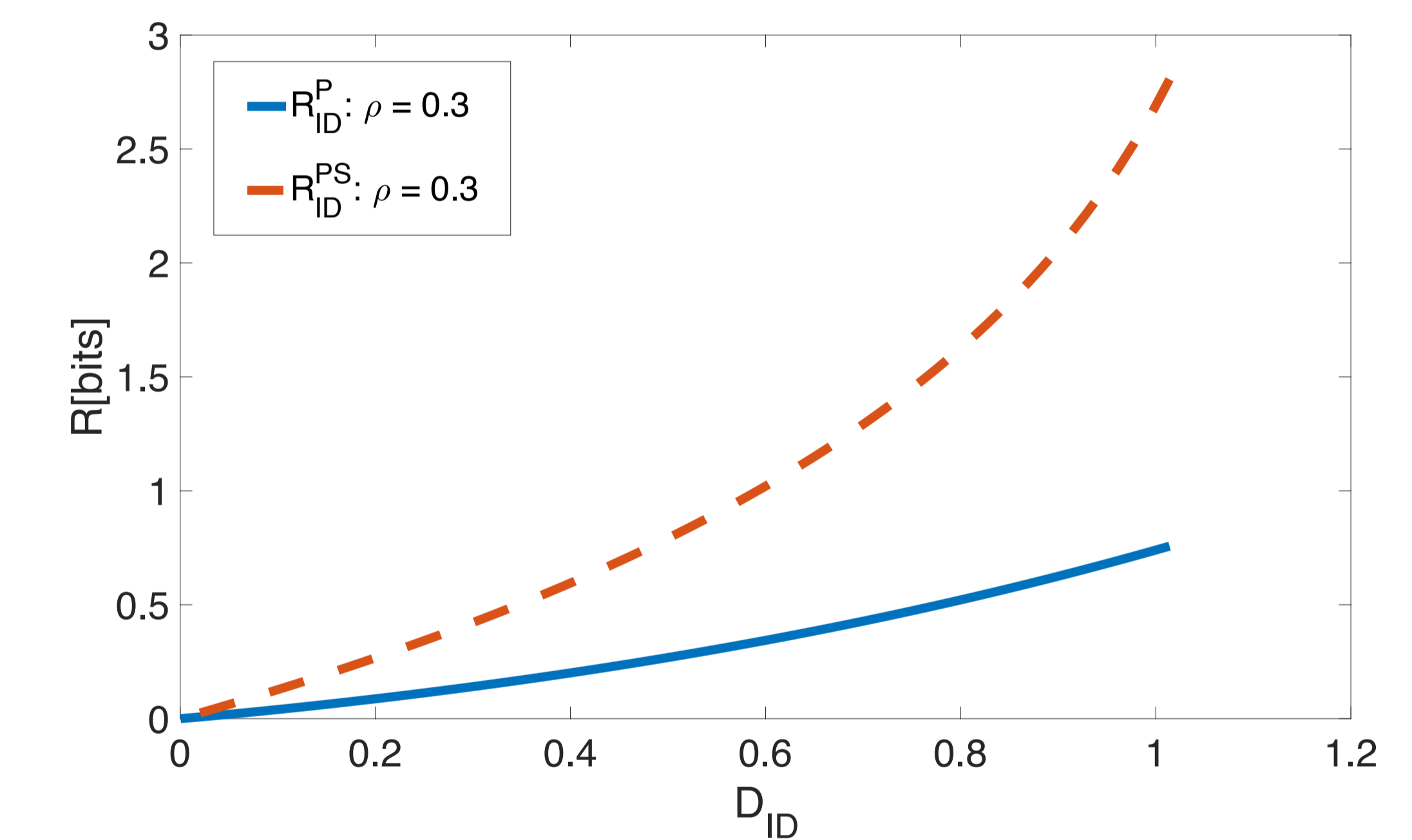


Figure 1: R_{ID}^P and R_{ID}^{PS} for AR(1) sequences with $\rho = 0.3$, and variance $\sigma_X^2 = \frac{1}{1-\rho^2}$

Conclusions

- Propose a prediction-based model for computing the identification rate of Gaussian autoregressive processes.
- The identification rate depends on a sequence of eigenvalues that we derive from our prediction model.
- The identification rate for the special case depends only on minimum eigenvalue of the Toeplitz matrix.

References

- [1] A. Ingber, T. Courtade, and T. Weissman, "Compression for quadratic similarity queries", *IEEE Trans. on Information Theory*, vol. 61, no. 5, pp. 2729-2747, May 2015.