

Learning Recurrent Models with Temporally Local Rules

Azwar Abdulsalam (abdulsal@purdue.edu)
Joseph G. Makin (jgmakin@purdue.edu)



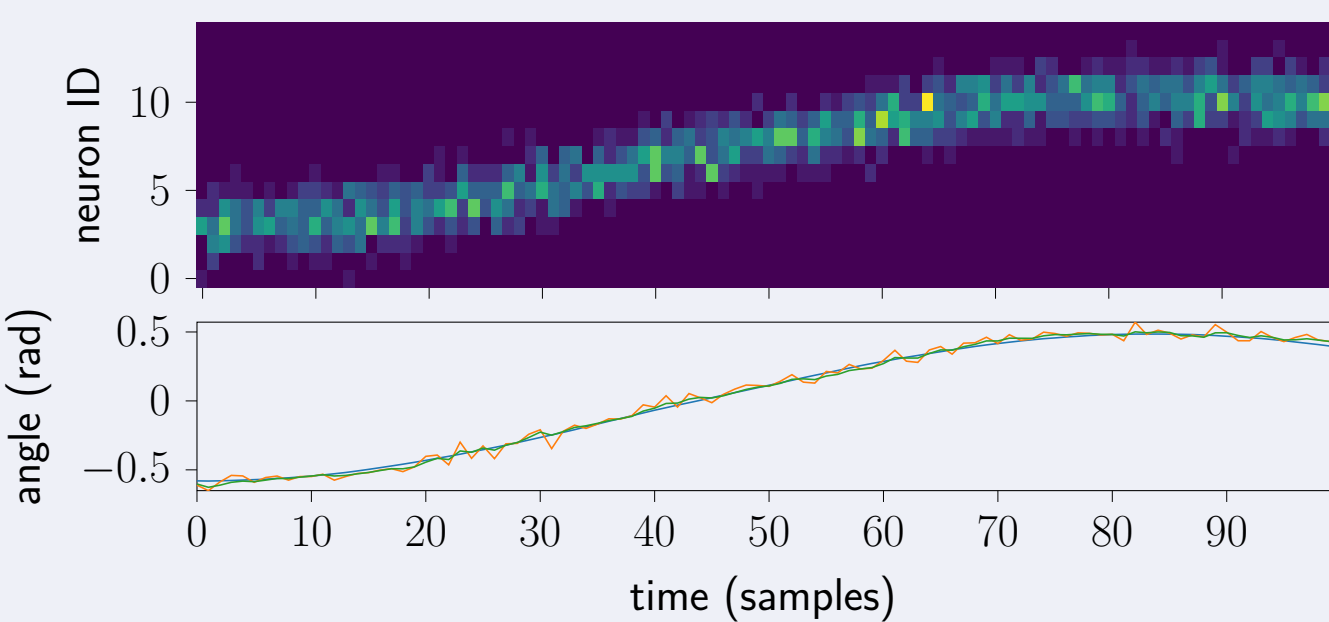
Elmore Family School of Electrical and
Computer Engineering, Purdue University

Abstract

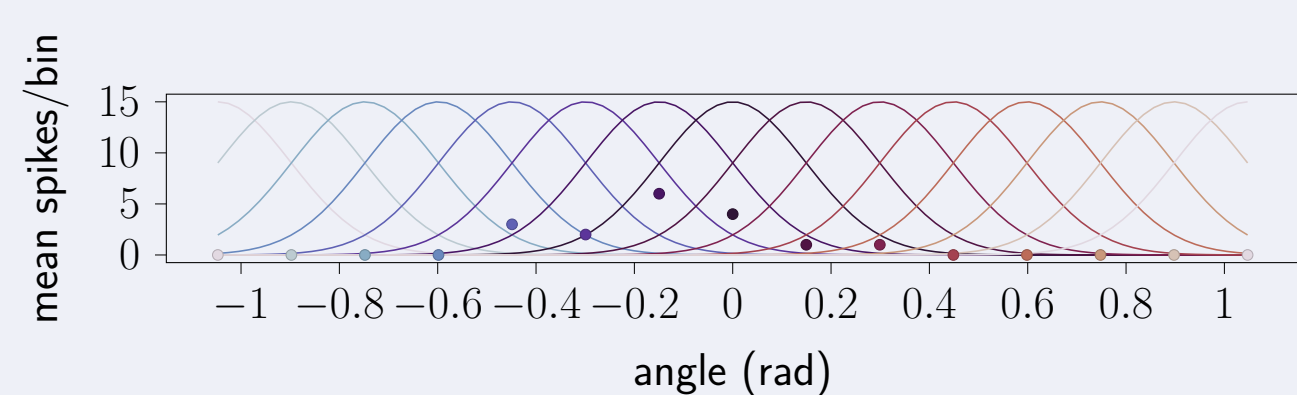
- Fitting **generative models** to sequential data typically requires **backprop-through-time**
- BPTT is biologically implausible and computationally expensive
- We investigate an alternative: **require the generative model to learn the joint distribution over current and previous states, rather than merely the forward-transition probabilities** [2, 3]
- Two architectures: **rEFH, rVAE**
- On toy datasets the procedure has the same effect as including BPTT

Experiments

LTI system + PPC [1]

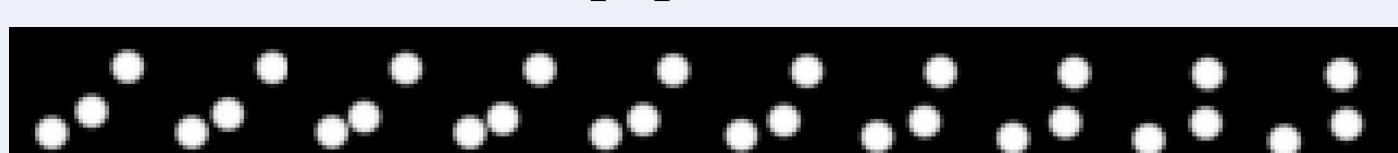


$$p(z_{t+1}|z_t) = \mathcal{N}(\mathbf{A}z_t, \Sigma_{\hat{x}}) \quad p(\mathbf{y}_t|z_t) = \prod_{i=1}^{15} \text{Pois}(y_t^i | g_t f_i(z_t^1))$$



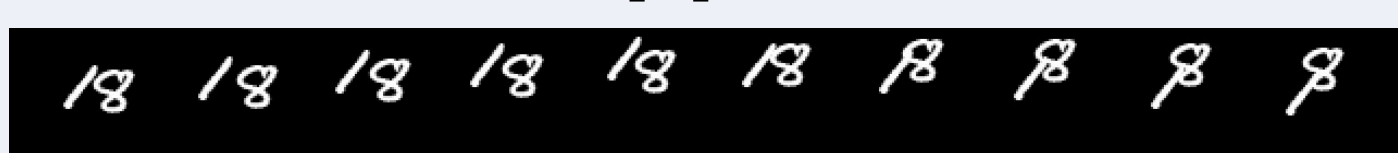
- A crude biophysical model of neurons reporting a stimulus.
- Second-order (oscillatory!) dynamics, but only position is "reported" by the neurons.
- Observation model is *nonlinear*, but a closed-form solution is still available (KF). This allows us to determine what order model was learned.
- What algorithms can learn **2nd-order** dynamics?

Bouncing balls [4]



- Video sequences of 3 balls bouncing off each other and walls with complete energy conservation
- constant velocities (2nd-order, linear) + collisions (nonlinear)

Moving MNIST [4]



- Video sequences of moving MNIST digits bouncing off walls and passing through one another
- constant velocities (2nd-order, linear) + overlap (nonlinear) + occlusions (nonlinear)

Models

$$\hat{p}(\hat{\mathbf{u}}_{t-1}|\hat{\mathbf{x}}_t; \theta) = \text{Bern}(\sigma(\mathbf{W}_u^T \hat{\mathbf{x}}_t + \mathbf{b}_u))$$

$$\hat{p}(\hat{\mathbf{x}}_t|\hat{\mathbf{u}}_{t-1}, \hat{\mathbf{y}}_t; \theta) = \text{Bern}(\sigma(\mathbf{W}_x \hat{\mathbf{u}}_{t-1} + \mathbf{W}_y \hat{\mathbf{y}}_t + \mathbf{b}_x))$$

$$\hat{p}(\hat{\mathbf{y}}_t|\hat{\mathbf{x}}_t; \theta) = \begin{cases} \text{Bern}(\sigma(\mathbf{W}_y^T \hat{\mathbf{x}}_t + \mathbf{b}_y)) \\ \text{Pois}(\exp(\mathbf{W}_y^T \hat{\mathbf{x}}_t + \mathbf{b}_y)) \end{cases}$$

$$\frac{d\mathcal{L}_{\text{rEFH}}}{d\theta} = \sum_{t=1}^T \frac{dH_{pp}[U_{t-1}, Y_t; \theta]}{d\theta} \approx \sum_{t=1}^T \mathbb{E}_{\hat{\mathbf{x}}_t, Y_t, U_{t-1}} \left[\frac{dE}{d\theta} \right] - \mathbb{E}_{\hat{\mathbf{x}}_t, \hat{Y}_t, \hat{U}_{t-1}} \left[\frac{dE}{d\theta} \right]$$

rEFH

$$\hat{p}(\hat{\mathbf{u}}_{t-1}|\hat{\mathbf{x}}_t; \theta) = \mathcal{N}(\mu_{\hat{x}}(\hat{\mathbf{x}}_t, \theta), \sigma_x^2 \mathbf{I})$$

$$\hat{p}(\hat{\mathbf{x}}_t; \theta) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\hat{p}(\hat{\mathbf{y}}_t|\hat{\mathbf{x}}_t; \theta) = \mathcal{N}(\mu_{\hat{y}}(\hat{\mathbf{x}}_t, \theta), \sigma_y^2 \mathbf{I})$$

$$\check{p}(\check{\mathbf{x}}_t|\mathbf{u}_{t-1}, \mathbf{y}_t; \phi) = \mathcal{N}(\nu_{\check{x}}(\mathbf{u}_{t-1}, \mathbf{y}_t), \Upsilon_{\check{x}}(\mathbf{u}_{t-1}, \mathbf{y}_t))$$

$$\frac{d\mathcal{L}_{\text{rVAE}}}{d\theta^T} = \frac{d}{d\theta^T} \sum_{t=1}^T \mathbb{E}_{\check{\mathbf{x}}_t, U_{t-1}, Y_t} [\log \check{p}(\check{\mathbf{X}}_t|U_{t-1}, Y_t; \phi) - \log \hat{p}(\check{\mathbf{X}}_t, U_{t-1}, Y_t; \theta)]$$

$$\approx \sum_{t=1}^T \mathbb{E}_{\check{\mathbf{x}}_t, U_{t-1}, Y_t} \left[\frac{d \log \check{p}(\check{\mathbf{X}}_t|U_{t-1}, Y_t; \phi)}{d\theta^T} - \frac{d \log \hat{p}(\check{\mathbf{X}}_t, U_{t-1}, Y_t; \theta)}{d\theta^T} \right]$$

rVAE

Letting the derivative pass through the expectations

- ignores the dependence of U_{t-1} on the parameters θ .
- amounts to discarding BPTT

Rationale

Model requirements

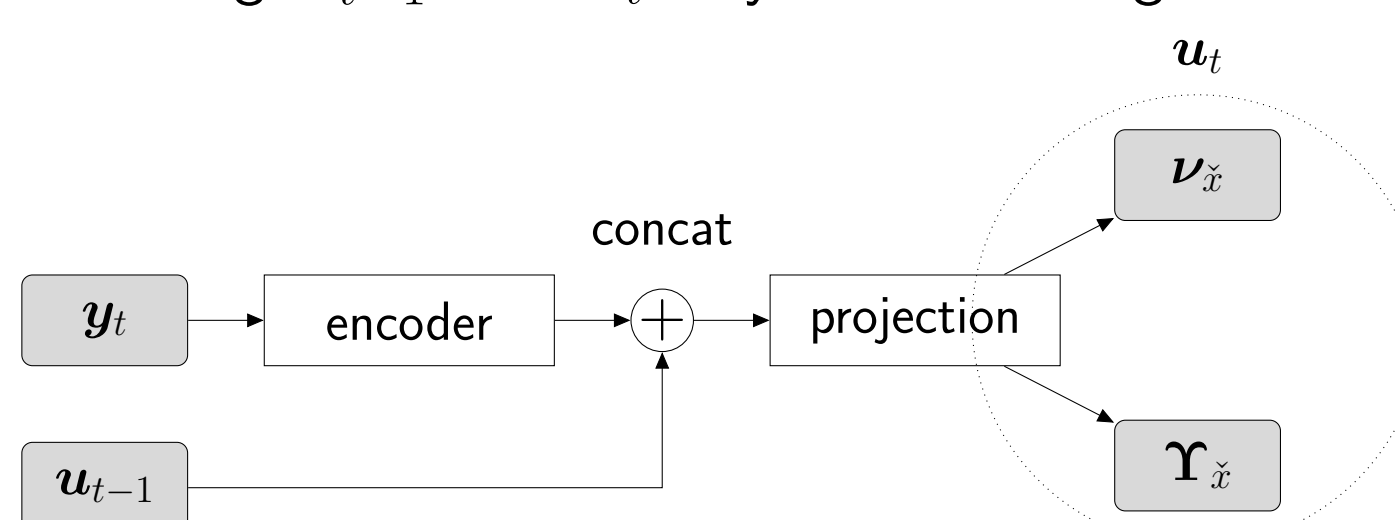
- 1) Generative model with latent variables, $\hat{\mathbf{X}}$
- 2) Latent variables inferable: $\hat{p}(\hat{\mathbf{x}}_t|\mathbf{u}_{t-1}, \mathbf{y}_t; \theta)$ or $\check{p}(\check{\mathbf{x}}_t|\mathbf{u}_{t-1}, \mathbf{y}_t; \phi)$
- 3) compression: $\dim(\hat{\mathbf{X}}) < \dim(\mathbf{Y})$

Candidate generative models

- ✗ GAN (cannot infer latent variables)
- ✗ Energy-Based Models (no latent variables)
- ✗ Diffusion ($\dim(\hat{\mathbf{X}}) = \dim(\mathbf{Y})$)
- ✗ Flow ($\dim(\hat{\mathbf{X}}) = \dim(\mathbf{Y})$)
- ✓ Variational Auto-Encoder
- ✓ Exponential-Family Harmonium

Architecture/Training

- The rVAE generative and recognition models make different independence statements. We try to minimize the discrepancy by concatenating U_{t-1} with Y_t only after encoding:



- All models trained with stochastic gradient descent and AdaM optimization. The learning rate was configured to 1e-4, with β_1 and β_2 values set to 0.9 and 0.999, respectively.

Results

Quantitative Results

MODEL	MSE
ORDER 0	12×10^{-4}
TVAE	9.5×10^{-4}
TRBM*	6.0×10^{-4}
KF-1	5.8×10^{-4}
rVAE	5.3×10^{-4}
rEFH	3.3×10^{-4}
RTRBM*	3.1×10^{-4}
KF-2	2.2×10^{-4}

Mean squared errors (MSE) for recovery of position information on the PPC experiment.

MODEL	MSE
ORDER 0	0.0120
TRBM	0.0124
rEFH	0.0067
RTRBM	0.0059

Mean squared errors (MSE) for bouncing-ball one-step predictions (w/clamped Gibbs sampling)

Qualitative Results



Bouncing-ball frame sequence generated by rVAE.



MovingMNIST frame sequence generated by rVAE.

Conclusions

- Learning joint over current and previous states seems to obviate BPTT
- rVAE learns 2nd-order dynamics from position observations, but worse than rEFH
- The procedure, though intuitive, requires a mathematical basis and scaling to handle more challenging datasets.

References

- [1] W. J. Ma, J. M. Beck, P. E. Latham, and A. Pouget. Bayesian Inference with Probabilistic Population Codes. *Nature Neuroscience*, 9(11):1423–1438, 2006.
- [2] J. G. Makin, B. K. Dichter, and P. N. Sabes. Learning to Estimate Dynamical State with Probabilistic Population Codes. *PLoS Computational Biology*, 11(11):1–28, 2015.
- [3] J. G. Makin, B. K. Dichter, and P. N. Sabes. Recurrent Exponential-Family Harmoniums without Backprop-Through-Time. <https://arxiv.org/abs/1605.05799>, 2016.
- [4] I. Sutskever. *Training Recurrent Neural Networks*. PhD thesis, University of Toronto, 2013.