

Time Series Prediction via Recurrent Neural Networks with the Information Bottleneck Principle

Duo Xu and Faramarz Fekri
{dxu301, faramarz.fekri}@ece.gatech.edu

Introduction

- **Background:**
 - Time series prediction is an important interdisciplinary topic in computer sciences, statistics, and econometrics.
 - Recently, deep learning methods have been applied successfully to the time series prediction.
- **Recurrent Neural Network:**
 - RNNs are particularly suitable for modeling time series as they operate on input information as well as the trace of previously acquired information.
 - The most successful models: LSTM and its variant GRU.
- **Probabilistic Modeling:**
 - We propose an stochastic RNN trained using the Recurrent Information Bottleneck (RIB) as objective function.
 - The method maximizes the mutual information between latent state and the target with low latent complexity.
 - Our model: the input is first encoded to latent state, then fed into the RNN cell to generate next hidden state. The decoder generates the distribution of the predicted value of time series with latent and hidden state together.
- **Two advantages:**
 - Stochastic latent state can augment RNN to better utilize recent observations.
 - RIB can model the temporal dependencies of high-dimensional time series with lower complexity.

Proposed Methods

- **Gated Recurrent Unit:** it has gating units that modulate the information flow inside the unit, but without having a separate memory cells:

$$\begin{aligned} r_t &= \sigma(\mathbf{W}_r \mathbf{x}_t + \mathbf{U}_r \mathbf{h}_{t-1} + \mathbf{b}_r) \\ \mathbf{u}_t &= \sigma(\mathbf{W}_u \mathbf{x}_t + \mathbf{U}_u \mathbf{h}_{t-1} + \mathbf{b}_u) \\ \mathbf{c}_t &= \sigma(\mathbf{W}_c \mathbf{x}_t + \mathbf{U}_c (r_t \odot \mathbf{h}_{t-1}) + \mathbf{b}_c) \\ \mathbf{h}_t &= \mathbf{u}_t \odot \mathbf{h}_{t-1} + (1 - \mathbf{u}_t) \odot \mathbf{c}_t \end{aligned}$$

with the input (\mathbf{x}_t), the hidden state of GRU cell (\mathbf{h}_t), the reset gate (r_t), the forgetting gate (\mathbf{u}_t), and candidate activation (\mathbf{c}_t).

- **Information Bottleneck:** It is information-theoretic view of deep neural networks. Assuming input random variable is X , latent state is Z and target is Y , the Markov chain $Y \leftrightarrow X \leftrightarrow Z$ holds.
- The aim is to get latent Z informative about the target Y , i.e. to maximize mutual information of Y and Z .
- In order to find a low-complexity representation of input, we can assume the mutual information between X and Z to be less than a pre-defined ϵ . Then we can formulate the optimization problem as

$$\max_{\theta} I(Z, Y; \theta) \quad \text{s.t.} \quad I(X, Z; \theta) < \epsilon$$

By introducing Lagrange multiplier $\beta \geq 0$, we have the objective:

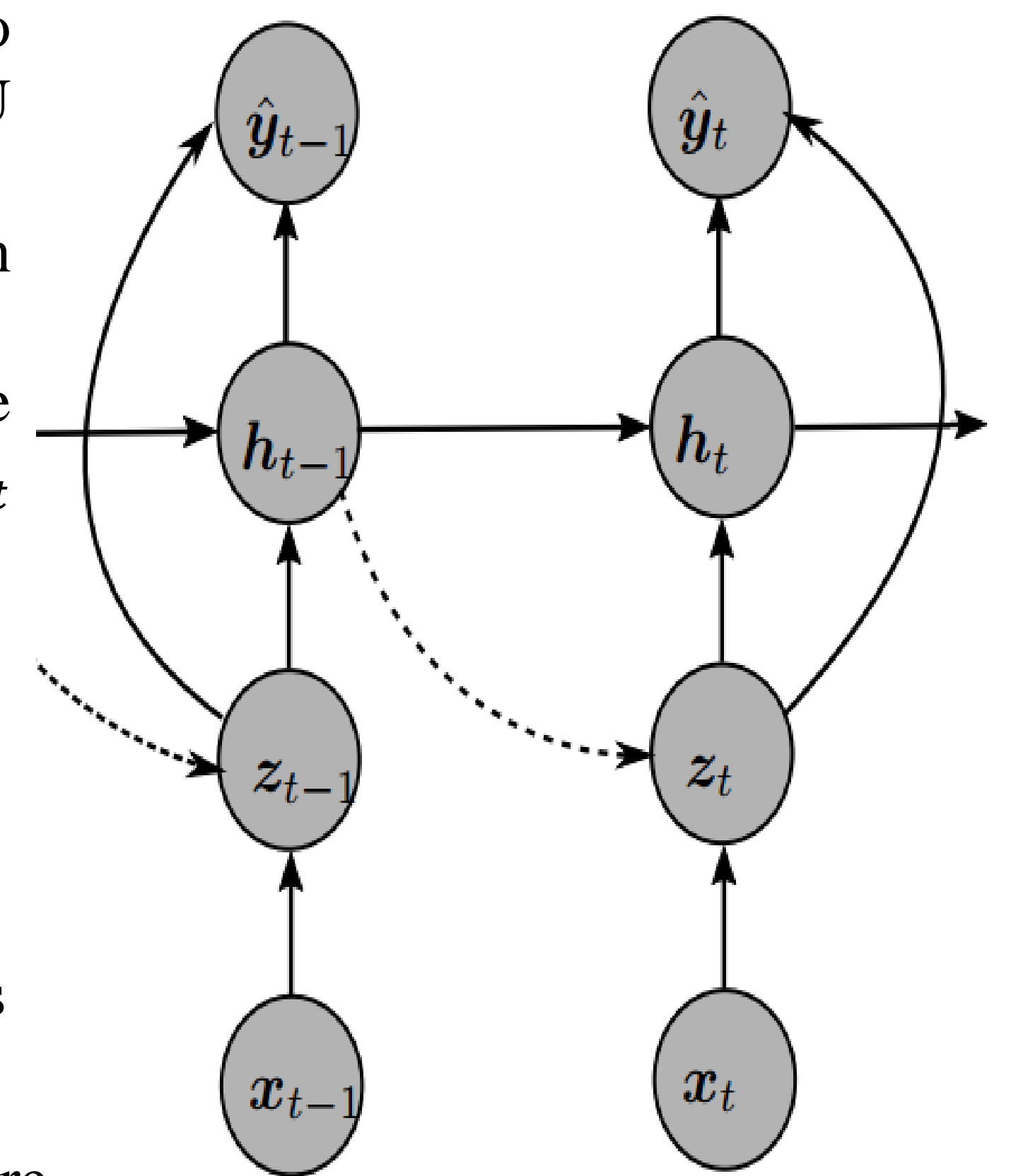
$$L_{IB}(\theta) = I(Z, Y; \theta) - \beta I(X, Z; \theta)$$

- **Recurrent Information Bottleneck:** Information bottleneck is extended to recurrent models. We use multi-variate Gaussian distribution as a variational approximation and model temporal relationships by hidden states of GRU. The objective function is,

$$\frac{1}{N} \sum_{n=1}^N \left[\sum_{t=1}^T \int p(z_{nt}|x_{nt}) \log q(y_{nt}|z_{nt}) dz_{nt} - \beta \int p(z_{nt}|x_{nt}) \log \frac{p(z_{nt}|x_{nt})}{\tilde{p}_t(z_{nt})} dz_{nt} \right]$$

Proposed Model

- In our model, the input is first encoded to latent state \mathbf{z}_t , which is then fed into GRU for state transition.
- \mathbf{z}_t is then decoded together with hidden state to generate output.
- The dashed lines in the diagram show the conditional dependency in both prior of \mathbf{z}_t and encoding of the input.
- The encoding follows distribution as $Z_t \sim \mathcal{N}(\mu_{\text{enc}}(\mathbf{x}_t, \mathbf{h}_{t-1}), \text{diag}(\sigma_{\text{enc}}^2(\mathbf{x}_t, \mathbf{h}_{t-1})))$
- The decoder follows distribution as $\hat{Y}_t \sim \mathcal{N}(\mu_{\text{dec}}(\mathbf{z}_t, \mathbf{h}_{t-1}), \text{diag}(\sigma_{\text{dec}}^2(\mathbf{z}_t, \mathbf{h}_{t-1})))$
- The prior of latent \mathbf{z}_t follows distribution as $\tilde{p}(\mathbf{z}_t|\mathbf{h}_{t-1}) = \mathcal{N}(\mu_{\text{prior}}(\mathbf{h}_{t-1}), \text{diag}(\sigma_{\text{prior}}^2(\mathbf{h}_{t-1})))$ where all mean and variance functions here are realized by multi-layer neural networks.
- In the GRU cell, the hidden state is translated following GRU operations.

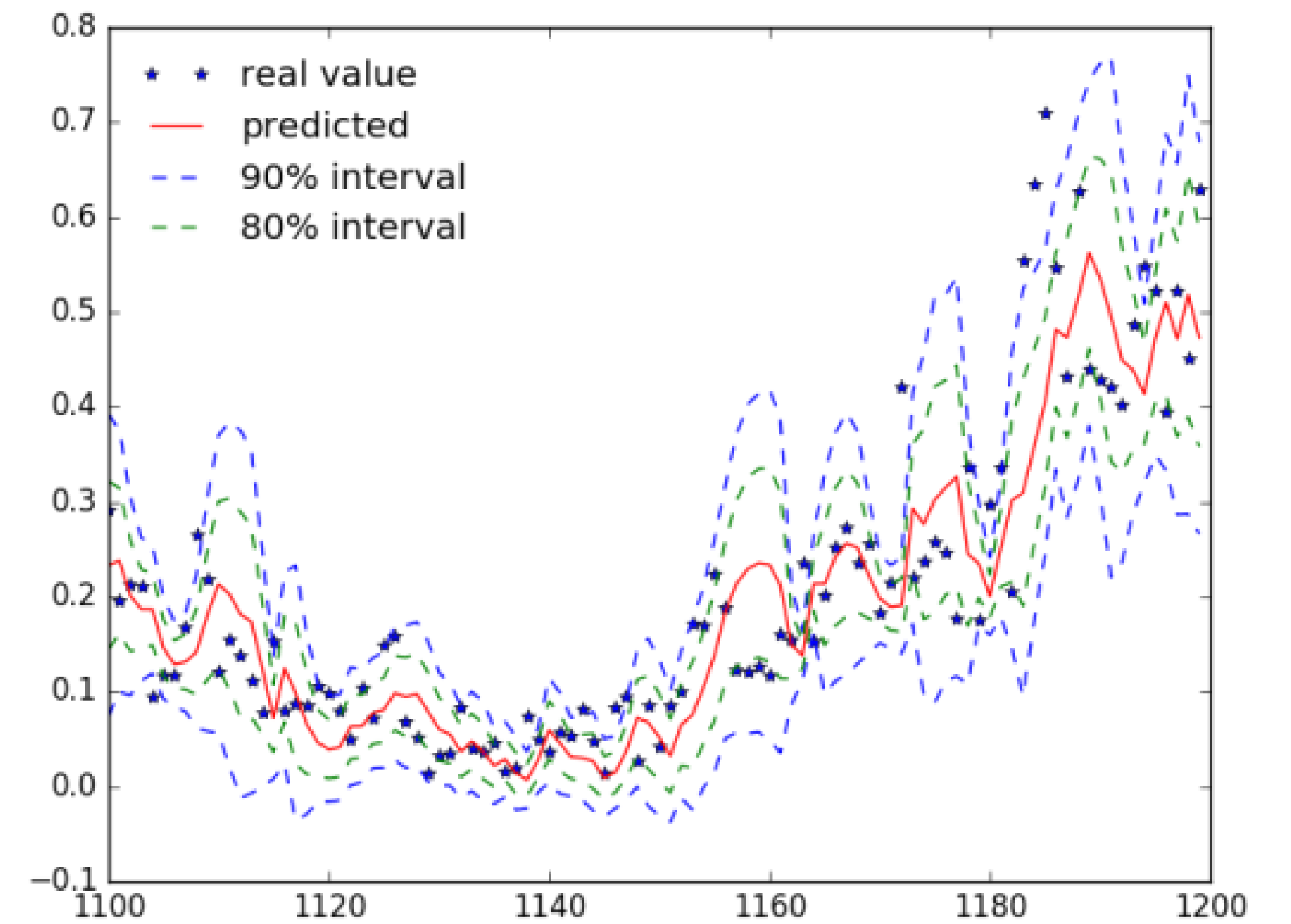


Experiments

- The first experiment uses sunspot dataset. It contains sunspot number collected in Zurich from Jan. 1749 to Dec. 1983. We use first 1000 data points for training and the rest for testing. The $1 - \alpha$ level confidence interval of prediction is

$$[\mu_{\text{dec}} - z_{1-\alpha/2} \sqrt{\sigma_{\text{dec}}^2},$$

$$\mu_{\text{dec}} + z_{1-\alpha/2} \sqrt{\sigma_{\text{dec}}^2}]$$



	Proposed	RBM	SAE	LSTM
MAE	0.0453	0.0911	0.116	0.0524
RMSE	0.0612	0.132	0.150	0.0739

- The second experiment is on traffic data, a collection of 15 months of daily data from California's Department of Transportation, describing occupancy rate of 963 freeways in bay area.

	Proposed	MatFact	VRNN	SAE
NMAE	0.1127	0.1935	0.2103	0.2234
NMRSE	0.3608	0.4263	0.4312	0.4566

Conclusion

- We propose to use information bottleneck to model time series probabilistically, and obtain better prediction performance than previous methods.

Acknowledgement

- This work is partially supported by the Center for Energy and Geo Processing (CeGP) at Georgia Tech and by King Fahd University of Petroleum and Minerals (KFUPM),