# Recurrent Generative Stochastic Networks for Sequence Prediction

## Abstract

We present a new generative model for unsupervised learning of sequence representations, the recurrent generative stochastic network (RNN-GSN).

## 1. Introduction

Unsupervised sequence learning is an important problem in machine learning given that most information (ranging from speech to video to even consumer behavior) is often unlabeled and has a sequential structure. Most of these sequences consist of high-dimensional, complex objects such as words in text, images in video, or chords in music. Recently, recurrent neural networks (RNN) (Rumelhart et al., 1986) have become state-of-the-art for sequence representation because they have an internal memory that can learn long-term, temporal dependencies. Further advances with Hessian-free optimization and Long Short-term Memory for RNNs (Martens & Sutskever, 2011; Hochreiter & Schmidhuber, 1996) have enabled learning highly complex temporal dependencies.

While traditional RNNs can learn complex sequences through maintaining an internal memory, they don't deal well when modeling complex input distributions, where the conditional distribution at each time step is highly multimodal. In real-world data, we often care more about this conditional distribution rather than the expected value. In music, for example, the existence of a particular note can highly change the probabilities with which other notes occur at the same time in a chord. In consumer behavior, the existence of a particular action taken can highly change the probabilities of other actions made within the same time window. We would prefer to reason over these distributions at each time step to form a better representation of the sequences.

Deep RNN architectures have been proposed to alleviate the traditional architecture shortcomings with complex input distributions (Pascanu et al., 2013). One proposed framework involves using deep input-to-hidden and

hidden-to-output functions to reduce input complexity, making the RNN's task easier.

(Boulanger-Lewandowski et al., 2012)

In Sections 2, 3, and 4, we introduce the GSN, RNN, and RNN-GSN architectures. We then provide experimental validation of the RNN-GSN on sequences of MNIST images and midi datasets in Section 5.

## 2. Generative Stochastic Networks

Generative stochastic networks (GSN) are a generalization of the denoising auto-encoder and help solve the problem of mixing between many major modes of the input data distribution.

### 2.1. Denoising auto-encoder

Denoising auto-encoders use a Markov chain to learn a reconstruction distribution $P(X|\widetilde{X})$ given a corruption process $C(\widetilde{X}|X)$ for some data $X$. Denoising auto-encoders have been shown as generative models (Bengio et al., 2013b), where the Markov chain can be iteratively sampled from:

$$X_t \sim P_\Theta(X|\widetilde{X}_{t-1})$$
$$\widetilde{X}_t \sim C(\widetilde{X}|X_t)$$

As long as the learned distribution $P_{\Theta_n}(X|\widetilde{X})$ is a consistent estimator of the true conditional distribution $P(X|\widetilde{X})$ and the Markov chain is ergodic, then as $n \to \infty$, the asymptotic distribution $\pi_n(X)$ of the generated samples from the denoising auto-encoder converges to the data-generating distribution $P(X)$ (Bengio et al., 2013b)).

### 2.2. Easing restrictive conditions on the denoising auto-encoder

A few restrictive conditions are necessary to guarantee ergodicity of the Markov chain - requiring $C(\widetilde{X}|X) > 0$ everywhere that $P(X) > 0$. Particularly, a large region $V$ containing any possible $X$ is defined such that the probability of moving between any two points in a single jump $C(\widetilde{X}|X)$ must be greater than 0. This restriction requires that $P_{\Theta_n}(X|\widetilde{X})$ has the ability to model every mode of $P(X)$, which is a problem this model was meant to avoid.

To ease this restriction, Bengio et al. (Bengio et al., 2013a) prove that using a $C(\widetilde{X}|X)$ that only makes small jumps allows $P_\Theta(X|\widetilde{X})$ to model a small part of the space $V$ around each $\widetilde{X}$. This weaker condition means that modeling the reconstruction distribution $P(X|\widetilde{X})$ would be easier since it would probably have fewer modes.

However, the jump size $\sigma$ between points must still be large enough to guarantee that one can jump often enough between the major modes of $P(X)$ to overcome the regions of low probability: $\sigma$ must be larger than half the largest distance of low probability between two nearby modes, such that $V$ has at least a single connected component between modes. This presents a tradeoff between the difficulty of learning $P_\Theta(X|\widetilde{X})$ and the ease of mixing between modes separated by this low probability region.

### 2.3. Generalizing to GSN

While denoising auto-encoders can rely on $X_t$ alone through a deterministic procedure for the state of the Markov chain, GSNs introduce a latent variable $H_t$ that acts as an additional state variable in the Markov chain along with the visible $X_t$ (Bengio et al., 2013a):

$$H_{t+1} \sim P_{\Theta_1}(H|H_t, X_t)$$
$$X_{t+1} \sim P_{\Theta_2}(X|H_{t+1})$$

The latent state variable $H$ can be equivalently defined as $H_{t+1} = f_{\Theta_1}(X_t, Z_t, H_t)$, a learned function $f$ with an independent noise source $Z_t$ such that $X_t$ cannot be reconstructed exactly from $H_{t+1}$. If $X_t$ could be recovered from $H_{t+1}$, the reconstruction distribution would simply converge to the Dirac at $X$. Denoising auto-encoders are therefore a special case of GSNs, where $f$ is fixed instead of learned.

GSNs also use the notion of walkbacks to aid training. Walkbacks are the process of generating samples by iteratively sampling from $P_{\Theta_1}(H|H_t, X_t)$ and $P_{\Theta_2}(X|H_{t+1})$ for a given input. By using walkbacks, the model is more likely to seek out spurious modes in the data distribution and correct for them (Bengio et al., 2013b). The resulting Markov chain of a GSN is inspired by Gibbs sampling, but with stochastic units at each layer that can be backpropagated (Rezende et al., 2014).

Experimental results with GSNs show that their latent states mix well between the major modes of the data - mixing faster at higher layers in the model (Bengio et al., 2013a). Using this property, we tested a simple temporal GSN model to predict sequences of inputs. The temporal GSN uses a linear transformation $H \rightarrow H$ to predict $P(H_t|H_{t-1}, ..., H_{t-n})$ with an input history of size $n$. Using this predicted $H_t$, an expected input $x_t$ was sampled from the model's learned reconstuction distribution
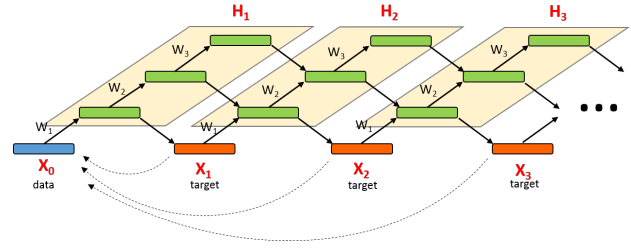


*Figure 1.* Unrolled GSN Markov chain.

$P_{\Theta_2}(X|H)$.

Qualitatively, this model appears to predict temporal dependencies well within its history window $n$ (Figure 3). This result provided motivation to use the latent state $H$ as the emitted parameter in a better-suited temporal model, such as an RNN.
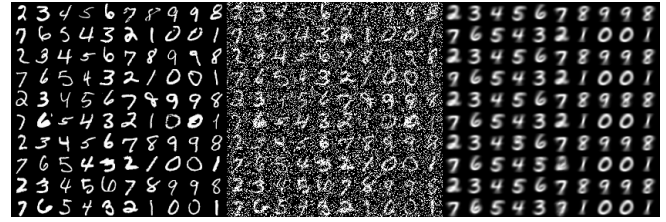


*Figure 2.* Temporal GSN with history $n = 2$ trained on an arbitrary MNIST sequence. Original input sequence is on the left, a noisy version of the input fed into the model is in the middle, and the predicted input based on the history is on the right.

## 3. Recurrent Neural Networks

This is the RNN section.

## 4. The RNN-GSN

The Recurrent GSN (RNN-GSN)

## 5. Experiments

This is the experiments section.

### 5.1. Sequences of MNIST digits

arbitrary sequences of images.

- Sequence1 is a simple linear sequence of digits 0-9 repeating.

- Sequence2 introduces one bit of parity by alternating sequences 0-9 and 9-0 repeating, where the next value depends on whether the sequence is ascending or descending.

- Sequence3 creates a longer, more complex sequence by introducing a second bit of parity. It is formed by:

### 5.2. Sequences of polyphonic music

midi stuff.

- Piano-midi.de .......

- Nottingham .....

- MuseData ....

- JSB chorales .....

## 6. Conclusion

This is the conclusion.

## References

Bengio, Yoshua, Thibodeau-Laufer, Eric, and Yosinski, Jason. Deep generative stochastic networks trainable by backprop. *CoRR*, abs/1306.1091, 2013a.

Bengio, Yoshua, Yao, Li, Alain, Guillaume, and Vincent, Pascal. Generalized denoising auto-encoders as generative models. *CoRR*, abs/1305.6663, 2013b.

Boulanger-Lewandowski, Nicolas, Bengio, Yoshua, and Vincent, Pascal. Modeling temporal dependencies in high-dimensional sequences: Application to polyphonic music generation and transcription. In *Proceedings of the 29th International Conference on Machine Learning*, 2012.

Hochreiter, Sepp and Schmidhuber, Jrgen. Bridging long time lags by weight guessing and "long short term memory". In *Spatiotemporal models in biological and artificial systems*, pp. 65–72. IOS Press, 1996.

Martens, James and Sutskever, Ilya. Learning recurrent neural networks with hessian-free optimization. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pp. 1033–1040, 2011.

Pascanu, Razvan, Gülçehre, Çaglar, Cho, Kyunghyun, and Bengio, Yoshua. How to construct deep recurrent neural networks. *CoRR*, abs/1312.6026, 2013. URL http://arxiv.org/abs/1312.6026.

Rezende, Danilo J., Mohamed, Shakir, and Wierstra, Daan. Stochastic backpropagation and approximate inference in deep generateive models. In *Proceedings of the 31st International Conference on Machine Learning*, 2014.

Rumelhart, David E., Hinton, Geoffrey E., and Williams, Ronald J. Learning representations by back-propagating errors. In *Parallel Dist. Proc.*, pp. 318–362. MIT Press, 1986.