

Deep Generative Models for Sequence Learning

Zhe Gan

Advisor: Dr. Lawrence Carin

Ph.D. Qualifying Exam

December 2nd, 2015

Outline

- 1 Literature Review
 - State Space Models
 - Recurrent Neural Networks
 - Temporal Restricted Boltzmann Machines
- 2 Deep Temporal Sigmoid Belief Networks
 - Model Formulation
 - Scalable Learning and Inference
 - Experiments
- 3 Conclusion and Future work

Table of Contents

- 1 Literature Review
 - State Space Models
 - Recurrent Neural Networks
 - Temporal Restricted Boltzmann Machines
- 2 Deep Temporal Sigmoid Belief Networks
 - Model Formulation
 - Scalable Learning and Inference
 - Experiments
- 3 Conclusion and Future work

State Space Models

- We are interested in developing probabilistic models for sequential data.
- Hidden Markov Models (**HMMs**) [12]: a mixture model, multinomial latent variables
- Linear Dynamical Systems (**LDS**) [8]: continuous latent variables

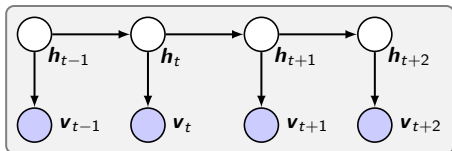


Figure: Graphical model for HMM and LDS.

Table of Contents

- 1 Literature Review
 - State Space Models
 - Recurrent Neural Networks
 - Temporal Restricted Boltzmann Machines
- 2 Deep Temporal Sigmoid Belief Networks
 - Model Formulation
 - Scalable Learning and Inference
 - Experiments
- 3 Conclusion and Future work

Recurrent Neural Networks

- Recurrent Neural Networks (RNNs) takes a sequence as input.
- Recursively processing each input \mathbf{v}_t while maintaining its internal hidden state \mathbf{h}_t .

$$\mathbf{h}_t = f(\mathbf{v}_{t-1}, \mathbf{h}_{t-1}) \quad (1)$$

where f is a deterministic non-linear transition function.

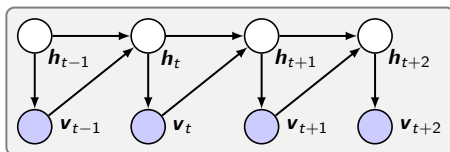


Figure: Illustration for RNNs.

Recurrent Neural Networks

- RNN defines the likelihood as

$$p(\mathbf{x}_1, \dots, \mathbf{x}_T) = \prod_{t=1}^T p(\mathbf{x}_t | \mathbf{x}_{<t}), \quad p(\mathbf{x}_t | \mathbf{x}_{<t}) = g(\mathbf{h}_t) \quad (2)$$

- How to define f ?

- ① Logistic function

$$\mathbf{h}_t = \sigma(\mathbf{W}\mathbf{x}_{t-1} + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{b}) \quad (3)$$

- ② Long Short-Term Memory (LSTM) [7]

- ③ Gated Recurrent Units (GRU) [3]

- All the randomness inside the model lies in the usage of the conditional probability $p(\mathbf{x}_t | \mathbf{x}_{<t})$.

Table of Contents

- 1 Literature Review
 - State Space Models
 - Recurrent Neural Networks
 - Temporal Restricted Boltzmann Machines

- 2 Deep Temporal Sigmoid Belief Networks
 - Model Formulation
 - Scalable Learning and Inference
 - Experiments

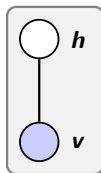
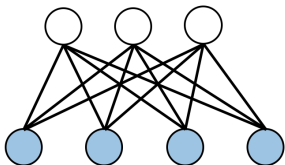
- 3 Conclusion and Future work

Restricted Boltzmann Machines

RBM [6] is an *undirected* graphical model that represents a joint distribution over binary visible units $\mathbf{v} \in \{0, 1\}^M$ and binary hidden units $\mathbf{h} \in \{0, 1\}^J$ as

$$P(\mathbf{v}, \mathbf{h}) = \exp\{-E(\mathbf{v}, \mathbf{h})\} / Z \quad (4)$$

$$E(\mathbf{v}, \mathbf{h}) = -(\mathbf{h}^\top \mathbf{W} \mathbf{v} + \mathbf{c}^\top \mathbf{v} + \mathbf{b}^\top \mathbf{h}) \quad (5)$$



Restricted Boltzmann Machines

- The posterior for \mathbf{v} and \mathbf{h} are factorized.

$$P(h_j = 1 | \mathbf{v}) = \sigma \left(\sum_m w_{jm} v_m + b_j \right) \quad (6)$$

$$P(v_i = 1 | \mathbf{h}) = \sigma \left(\sum_m w_{im} h_m + c_m \right) \quad (7)$$

- **Inference:** block Gibbs sampling.
- **Learning:** Contrastive Divergence (CD).
- **Application:**
 - 1 Deep Belief Networks (DBNs) and Deep Boltzmann Machines (DBMs)
 - 2 Learn the sequential dependencies in time-series data [9, 13, 14, 15, 16].

Temporal Restricted Boltzmann Machines

- TRBM [13] is defined as a sequence of RBMs

$$p(\mathbf{V}, \mathbf{H}) = p_0(\mathbf{v}_1, \mathbf{h}_1) \prod_{t=2}^T p(\mathbf{v}_t, \mathbf{h}_t | \mathbf{h}_{t-1}) \quad (8)$$

- Each $p(\mathbf{v}_t, \mathbf{h}_t | \mathbf{h}_{t-1})$ is defined as a conditional RBM.

$$p(\mathbf{v}_t, \mathbf{h}_t | \mathbf{h}_{t-1}) \propto \exp(\mathbf{v}_t^\top \mathbf{W} \mathbf{h}_t + \mathbf{v}_t^\top \mathbf{c} + \mathbf{h}_t^\top (\mathbf{b} + \mathbf{U}^\top \mathbf{h}_{t-1}))$$

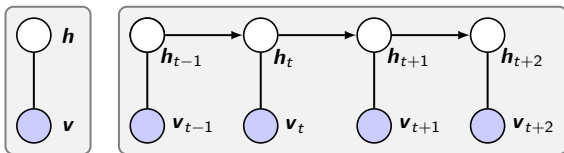


Figure: Graphical model for TRBM.

Model Reviews

| PROPERTIES | HMM | LDS | RNN | TRBM |
|--------------------------|-----|-----|-----|------|
| DIRECTED | ✓ | ✓ | ✓ | ✗ |
| LATENT | ✓ | ✓ | ✗ | ✓ |
| NONLINEAR | ✗ | ✗ | ✓ | ✓ |
| DISTRIBUTED ¹ | ✗ | — | — | ✓ |

Can we find a model that has all the properties listed above? A deep directed latent variable model.

¹Each hidden state in HMM is a one-hot vector. To represent 2^N distinct states, an HMM requires a length- 2^N one-hot vector, while for TRBM, only a length- N vector is needed.

Table of Contents

- 1 Literature Review
 - State Space Models
 - Recurrent Neural Networks
 - Temporal Restricted Boltzmann Machines
- 2 Deep Temporal Sigmoid Belief Networks
 - Model Formulation
 - Scalable Learning and Inference
 - Experiments
- 3 Conclusion and Future work

Overview

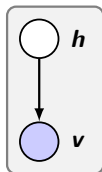
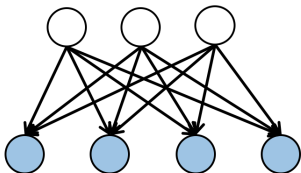
- **Problem of interest:** Developing **deep directed latent variable** models for sequential data.
- **Main idea:**
 - ① Constructing a hierarchy of Temporal Sigmoid Belief Networks (**TSBNs**).
 - ② TSBN is defined as a sequential stack of Sigmoid Belief Networks (**SBNs**).
- **Challenge:** Designing **scalable learning and inference** algorithms.
- **Solution:**
 - ① Stochastic Variational Inference (**SVI**).
 - ② Design a recognition model for fast inference.

Sigmoid Belief Networks

An SBN [11] models a visible vector $\mathbf{v} \in \{0, 1\}^M$, in terms of hidden variables $\mathbf{h} \in \{0, 1\}^J$ and weights $\mathbf{W} \in \mathbb{R}^{M \times J}$ with

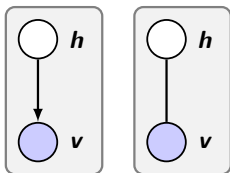
$$p(v_m = 1 | \mathbf{h}) = \sigma(\mathbf{w}_m^\top \mathbf{h} + c_m), \quad p(h_j = 1) = \sigma(b_j) \quad (9)$$

where $\sigma(x) \triangleq 1/(1 + e^{-x})$.



SBN vs. RBM

- Graphical Model



- Energy function

$$-E_{SBN}(\mathbf{v}, \mathbf{h}) = \mathbf{v}^\top \mathbf{c} + \mathbf{v}^\top \mathbf{W} \mathbf{h} + \mathbf{h}^\top \mathbf{b} - \sum_m \log(1 + e^{\mathbf{w}_m^\top \mathbf{h} + c_m})$$

$$-E_{RBM}(\mathbf{v}, \mathbf{h}) = \mathbf{v}^\top \mathbf{c} + \mathbf{v}^\top \mathbf{W} \mathbf{h} + \mathbf{h}^\top \mathbf{b}$$

- How to generate data
 - SBN**: ancestral sampling
 - RBM**: iterative Gibbs sampling

SBN vs. RBM

- Inference methods
 - ① **SBN**: Gibbs sampling, mean-field VB, Recognition model
 - ② **RBM**: Gibbs sampling
- Learning methods
 - ① **SBN**: SGD, Gibbs sampling, mean-field VB, MCEM (Polya-Gamma data augmentation)
 - ② **RBM**: Contrastive Divergence
- SBN has been shown potential to build deep models [4, 5].
 - ① Binary image modeling
 - ② Topic modeling

Temporal Sigmoid Belief Networks

- TSBN is defined as a sequence of SBNs.

$$p_{\theta}(\mathbf{V}, \mathbf{H}) = p(\mathbf{h}_1)p(\mathbf{v}_1|\mathbf{h}_1) \cdot \prod_{t=2}^T p(\mathbf{h}_t|\mathbf{h}_{t-1}, \mathbf{v}_{t-1}) \cdot p(\mathbf{v}_t|\mathbf{h}_t, \mathbf{v}_{t-1})$$

- For $t = 1, \dots, T$, each conditional distribution is expressed as

$$p(h_{jt} = 1 | \mathbf{h}_{t-1}, \mathbf{v}_{t-1}) = \sigma(\mathbf{w}_{1j}^{\top} \mathbf{h}_{t-1} + \mathbf{w}_{3j}^{\top} \mathbf{v}_{t-1} + b_j) \quad (10)$$

$$p(v_{mt} = 1 | \mathbf{h}_t, \mathbf{v}_{t-1}) = \sigma(\mathbf{w}_{2m}^{\top} \mathbf{h}_t + \mathbf{w}_{4m}^{\top} \mathbf{v}_{t-1} + c_m) \quad (11)$$

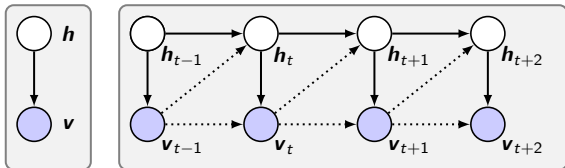



Figure: Graphical model for TSBN.

TSBN Properties

- Our TSBN model can be viewed as :
 - ① an HMM with distributed hidden state representations²;
 - ② a LDS with non-linear dynamics;
 - ③ a probabilistic construction of RNN;
 - ④ a directed-graphical-model counterpart of TRBM.

²Each hidden state in the HMM is represented as a *one-hot* length- J vector, while in the TSBN, the hidden states can be any length- J binary vector. 

TSBN Variants: Modeling Different Data Types

- **Real-valued data:** $p(\mathbf{v}_t | \mathbf{h}_t, \mathbf{v}_{t-1}) = \mathcal{N}(\boldsymbol{\mu}_t, \text{diag}(\boldsymbol{\sigma}_t^2))$, where

$$\mu_{mt} = \mathbf{w}_{2m}^\top \mathbf{h}_t + \mathbf{w}_{4m}^\top \mathbf{v}_{t-1} + c_m, \quad (12)$$

$$\log \sigma_{mt}^2 = (\mathbf{w}'_{2m})^\top \mathbf{h}_t + (\mathbf{w}'_{4m})^\top \mathbf{v}_{t-1} + c'_m \quad (13)$$

- **Count data:** $p(\mathbf{v}_t | \mathbf{h}_t, \mathbf{v}_{t-1}) = \prod_{m=1}^M y_{mt}^{v_{mt}}$, where

$$y_{mt} = \frac{\exp(\mathbf{w}_{2m}^\top \mathbf{h}_t + \mathbf{w}_{4m}^\top \mathbf{v}_{t-1} + c_m)}{\sum_{m'=1}^M \exp(\mathbf{w}_{2m'}^\top \mathbf{h}_t + \mathbf{w}_{4m'}^\top \mathbf{v}_{t-1} + c_{m'})}. \quad (14)$$

TSBN Variants: Boosting Performance

- **High Order:** $\mathbf{h}_t, \mathbf{v}_t$ depends on $\mathbf{h}_{t-n:t-1}, \mathbf{v}_{t-n:t-1}$.
- **Going Deep:**
 - 1 adding stochastic hidden layers;
 - 2 adding deterministic hidden layers.
- Conditional independent “becomes” conditional dependent.

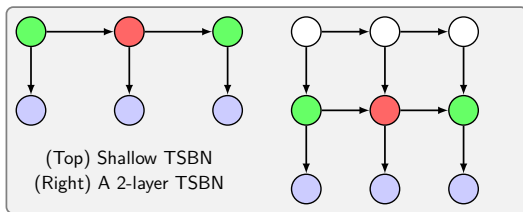


Table of Contents

- 1 Literature Review
 - State Space Models
 - Recurrent Neural Networks
 - Temporal Restricted Boltzmann Machines
- 2 Deep Temporal Sigmoid Belief Networks
 - Model Formulation
 - Scalable Learning and Inference
 - Experiments
- 3 Conclusion and Future work

Scalable Learning and Inference

- By using Polya-Gamma, Gibbs sampling or mean-field VB can be implemented [5, 11], but are inefficient.
- To allow for
 - ① tractable and scalable inference and parameter learning
 - ② without loss of the flexibility of the variational posterior

We apply recent advances in **stochastic variational Bayes** for **non-conjugate** inference [10].

Recognition model

- Variational Lower Bound

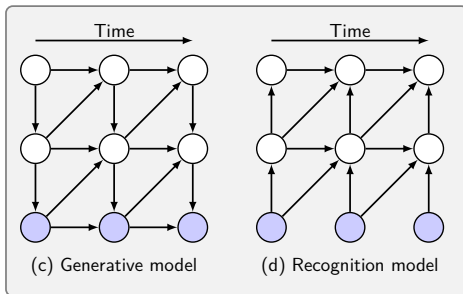
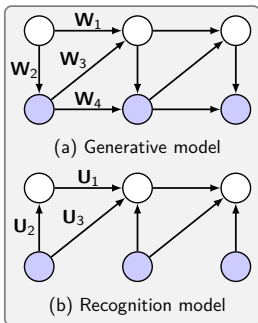
$$\mathcal{L}(\mathbf{V}, \theta, \phi) = \mathbb{E}_{q_{\phi}(\mathbf{H}|\mathbf{V})}[\log p_{\theta}(\mathbf{V}, \mathbf{H}) - \log q_{\phi}(\mathbf{H}|\mathbf{V})]. \quad (15)$$

- How to design $q_{\phi}(\mathbf{H}|\mathbf{V})$? No mean-field VB!
- **Recogniton model**: introduce a fixed-form distribution $q_{\phi}(\mathbf{H}|\mathbf{V})$ to approximate the true posterior $p(\mathbf{H}|\mathbf{V})$.
 - 1 Fast inference
 - 2 Utilization of a global set of parameters
 - 3 Potential better fit of the data

Recognition model

Recognition model is also defined as a TSBN.

$$q_{\phi}(\mathbf{H}|\mathbf{V}) = q(\mathbf{h}_1|\mathbf{v}_1) \cdot \prod_{t=2}^T q(\mathbf{h}_t|\mathbf{h}_{t-1}, \mathbf{v}_t, \mathbf{v}_{t-1}), \quad (16)$$



Parameter Learning

- Gradients:

$$\nabla_{\theta} \mathcal{L}(\mathbf{V}) = \mathbb{E}_{q_{\phi}} [\nabla_{\theta} \log p_{\theta}(\mathbf{V}, \mathbf{H})]$$

$$\nabla_{\phi} \mathcal{L}(\mathbf{V}) = \mathbb{E}_{q_{\phi}} [(\log p_{\theta}(\mathbf{V}, \mathbf{H}) - \log q_{\phi}(\mathbf{H}|\mathbf{V})) \times \nabla_{\phi} \log q_{\phi}(\mathbf{H}|\mathbf{V})]$$

- Variance Reduction:

- 1 centering the learning signal
- 2 variance normalization

Table of Contents

- 1 Literature Review
 - State Space Models
 - Recurrent Neural Networks
 - Temporal Restricted Boltzmann Machines
- 2 Deep Temporal Sigmoid Belief Networks
 - Model Formulation
 - Scalable Learning and Inference
 - Experiments
- 3 Conclusion and Future work

Experiments

- Datasets:

| | |
|--------------------|-------------|
| BOUNCING BALLS | BINARY |
| MOTION CAPTURE | REAL-VALUED |
| POLYPHONIC MUSIC | BINARY |
| STATE OF THE UNION | COUNT |

- Notations:

- ① *HMSBN*: TSBN model with $\mathbf{W}_3 = \mathbf{0}$ and $\mathbf{W}_4 = \mathbf{0}$;
- ② *DTSBN-S*: Deep TSBN with stochastic hidden layer;
- ③ *DTSBN-D*: Deep TSBN with deterministic hidden layer.

- Experimental Setup:

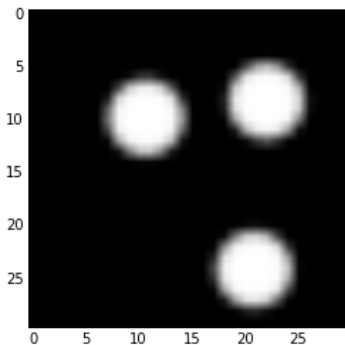
- Random initialization
- RMSprop optimization
- Monte Carlo sampling using only a single sample
- Weight decay regularization

Experiments

- **Prediction** of \mathbf{v}_t given $\mathbf{v}_{1:t-1}$
 - ① first obtain a sample from $q_\phi(\mathbf{h}_{1:t-1}|\mathbf{v}_{1:t-1})$;
 - ② calculate the conditional posterior $p_\theta(\mathbf{h}_t|\mathbf{h}_{1:t-1}, \mathbf{v}_{1:t-1})$ of the current hidden state ;
 - ③ make a prediction for \mathbf{v}_t using $p_\theta(\mathbf{v}_t|\mathbf{h}_{1:t}, \mathbf{v}_{1:t-1})$.
- **Generation**: ancestral sampling.

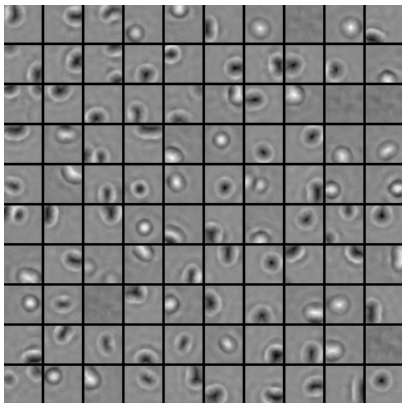
Bouncing balls dataset

Each video is of length 100 and of resolution 30×30 .



Bouncing balls dataset: Dictionaries

Our learned dictionaries are spatially localized.



Bouncing balls dataset: Prediction

- 1 Our approach achieves state-of-the-art.
- 2 A high-order TSBN reduces the prediction error significantly.
- 3 Using multiple layers improve performance.

| MODEL | DIM | ORDER | PRED. ERR. | NEG. LOG. LIKE. |
|--------------------|---------|-------|------------------------|-------------------------|
| DTsBN-S | 100-100 | 2 | 2.79 \pm 0.39 | 69.29 \pm 1.52 |
| DTsBN-D | 100-100 | 2 | 2.99 \pm 0.42 | 70.47 \pm 1.52 |
| TSBN | 100 | 4 | 3.07 \pm 0.40 | 70.41 \pm 1.55 |
| TSBN | 100 | 1 | 9.48 \pm 0.38 | 77.71 \pm 0.83 |
| RTRBM $^\diamond$ | 3750 | 1 | 3.88 \pm 0.33 | — |
| SRTRBM $^\diamond$ | 3750 | 1 | 3.31 \pm 0.33 | — |

Motion capture dataset: Prediction

We used 33 running and walking sequences of subject 35 in the CMU motion capture dataset.

TSBN-based models improves over the RBM-based models significantly.

| MODEL | WALKING | RUNNING |
|----------------------|------------------------|------------------------|
| DTSBN-s | 4.40 \pm 0.28 | 2.56 \pm 0.40 |
| DTSBN-D | 4.62 \pm 0.01 | 2.84 \pm 0.01 |
| TSBN | 5.12 \pm 0.50 | 4.85 \pm 1.26 |
| HMSBN | 10.77 \pm 1.15 | 7.39 \pm 0.47 |
| SS-SRTRBM \diamond | 8.13 \pm 0.06 | 5.88 \pm 0.05 |
| G-RTRBM \diamond | 14.41 \pm 0.38 | 10.91 \pm 0.27 |

Motion capture dataset: Generation

- 1 These generated data are readily produced from the model and demonstrate realistic behavior.
- 2 The smooth trajectories are walking movements, while the vibrating ones are running.

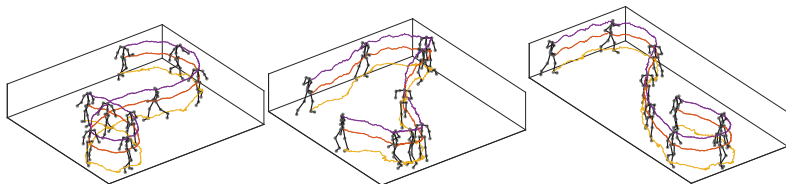


Figure: Generated motion trajectories. (Left) Walking. (Middle) Running-running-walking. (Right) Running-walking.

Polyphonic music dataset: Prediction

Four polyphonic music sequences of piano [2]: Piano-midi.de, Nottingham, MuseData and JSB chorales.

Table: Test log-likelihood for music datasets. (\diamond) taken from [2].

| MODEL | PIANO. | NOTT. | MUSE. | JSB. |
|---------------------|--------------|--------------|--------------|--------------|
| TSBN | -7.98 | -3.67 | -6.81 | -7.48 |
| RNN-NADE \diamond | -7.05 | -2.31 | -5.60 | -5.56 |
| RTRBM \diamond | -7.36 | -2.62 | -6.35 | -6.35 |
| RNN \diamond | -8.37 | -4.46 | -8.13 | -8.71 |

Polyphonic music dataset: Generation

We can generate different styles of music based on different training datasets.

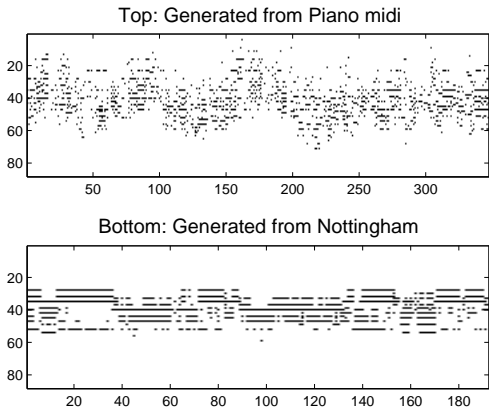


Figure: Generated samples..

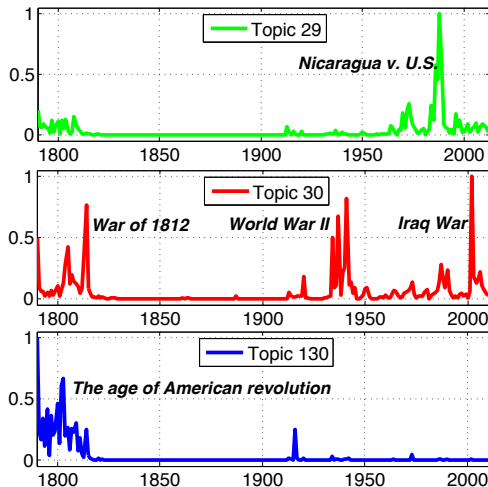
State of the Union dataset: Prediction

- Prediction is concerned with estimating the held-out words.
- MP: *mean precision* over all the years that appear.
- PP: *predictive precision* for the final year.

| MODEL | DIM | MP | PP |
|--------------------|-------|--------------------------|--------------------------|
| HMSBN | 25 | 0.327 \pm 0.002 | 0.353 \pm 0.070 |
| DHMSBN-s | 25-25 | 0.299 \pm 0.001 | 0.378 \pm 0.006 |
| GP-DPFA \diamond | 100 | 0.223 \pm 0.001 | 0.189 \pm 0.003 |
| DRFM \diamond | 25 | 0.217 \pm 0.003 | 0.177 \pm 0.010 |

State of the Union dataset: Dynamic Topic Modeling

The learned trajectory exhibits different temporal patterns across the topics.



State of the Union dataset: Dynamic Topic Modeling

Table: Top 10 most probable words associated with the STU topics.

| Topic #29 | Topic #30 | Topic #130 | Topic #64 | Topic #70 | Topic #74 |
|------------|-----------|------------|-------------|------------|-------------|
| family | officer | government | generations | Iraqi | Philippines |
| budget | civilized | country | generation | Qaida | islands |
| Nicaragua | warfare | public | recognize | Iraq | axis |
| free | enemy | law | brave | Iraqis | Nazis |
| future | whilst | present | crime | Al | Japanese |
| freedom | gained | citizens | race | Saddam | Germans |
| excellence | lake | united | balanced | ballistic | mines |
| drugs | safety | house | streets | terrorists | sailors |
| families | American | foreign | college | hussein | Nazi |
| God | militia | gentlemen | school | failures | hats |

Conclusion and Future work

- Conclusion
 - 1 Proposed Temporal Sigmoid Belief Networks
 - 2 Developed an efficient variational optimization algorithm
 - 3 Extensive applications to diverse data types
- Future work
 - 1 Controlled style transitioning
 - 2 Language modeling
 - 3 Multi-modality learning

References I



A. Acharya, J. Ghosh, and M. Zhou.

Nonparametric Bayesian factor analysis for dynamic count matrices.
In *AISTATS*, 2015.



N. Boulanger-Lewandowski, Y. Bengio, and P. Vincent.

Modeling temporal dependencies in high-dimensional sequences: Application to polyphonic music generation and transcription.
In *ICML*, 2012.



Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio.

Learning phrase representations using rnn encoder-decoder for statistical machine translation.
In *arXiv:1406.1078*, 2014.



Z. Gan, C. Chen, R. Henao, D. Carlson, and L. Carin.

Scalable deep poisson factor analysis for topic modeling.
In *ICML*, 2015.



Z. Gan, R. Henao, D. Carlson, and L. Carin.

Learning deep sigmoid belief networks with data augmentation.
In *AISTATS*, 2015.



G. Hinton.

Training products of experts by minimizing contrastive divergence.
In *Neural computation*, 2002.

References II



Sepp Hochreiter and Jürgen Schmidhuber.

Long short-term memory.
In *Neural computation*, 1997.



R. Kalman.

Mathematical description of linear dynamical systems.
In *J. the Society for Industrial & Applied Mathematics, Series A: Control*, 1963.



R. Mittelman, B. Kuipers, S. Savarese, and H. Lee.

Structured recurrent temporal restricted boltzmann machines.
In *ICML*, 2014.



A. Mnih and K. Gregor.

Neural variational inference and learning in belief networks.
In *ICML*, 2014.



R. Neal.

Connectionist learning of belief networks.
In *Artificial intelligence*, 1992.



L. Rabiner and B. Juang.

An introduction to hidden markov models.
In *ASSP Magazine, IEEE*, 1986.



I. Sutskever and G. Hinton.

Learning multilevel distributed representations for high-dimensional sequences.
In *AISTATS*, 2007.

References III



I. Sutskever, G. Hinton, and G. Taylor.

The recurrent temporal restricted boltzmann machine.
In *NIPS*, 2009.



G. Taylor and G. Hinton.

Factored conditional restricted boltzmann machines for modeling motion style.
In *ICML*, 2009.



G. Taylor, G. Hinton, and S. Roweis.

Modeling human motion using binary latent variables.
In *NIPS*, 2006.

Backup I

Contrastive Divergence:

$$\frac{\partial P(\mathbf{v})}{\partial \theta} = \mathbb{E}_{p(\mathbf{v}, \mathbf{h})} \left[\frac{\partial}{\partial \theta} E(\mathbf{v}, \mathbf{h}) \right] - \mathbb{E}_{p(\mathbf{h}|\mathbf{v})} \left[\frac{\partial}{\partial \theta} E(\mathbf{v}, \mathbf{h}) \right] \quad (17)$$

Backup II

The contribution of \mathbf{v} to the log-likelihood can be lower-bounded as follows

$$\log p_{\theta} = \log \sum_{\mathbf{h}} p_{\theta}(\mathbf{v}, \mathbf{h}) \quad (18)$$

$$\geq \sum_{\mathbf{h}} q_{\phi}(\mathbf{h}|\mathbf{v}) \log \frac{p_{\theta}(\mathbf{v}, \mathbf{h})}{q_{\phi}(\mathbf{h}|\mathbf{v})} \quad (19)$$

$$= \mathbb{E}_{q_{\phi}(\mathbf{h}|\mathbf{v})} [\log p_{\theta}(\mathbf{v}, \mathbf{h}) - \log q_{\phi}(\mathbf{h}|\mathbf{v})] \quad (20)$$

$$= \mathcal{L}(\mathbf{v}, \theta, \phi) \quad (21)$$

We can also rewrite the bound as

$$\mathcal{L}(\mathbf{v}, \theta, \phi) = \log p_{\theta}(\mathbf{v}) - KL(q_{\phi}(\mathbf{h}|\mathbf{v}), p_{\theta}(\mathbf{h}|\mathbf{v})) \quad (22)$$

Backup III

Differentiating the variational lower bound w.r.t. to the recognition model parameters gives

$$\nabla_{\phi} \mathcal{L}(\mathbf{v}) = \nabla_{\phi} \mathbb{E}_q[\log p_{\theta}(\mathbf{v}, \mathbf{h}) - \log q_{\phi}(\mathbf{h}|\mathbf{v})] \quad (23)$$

$$\begin{aligned} &= \nabla_{\phi} \sum_{\mathbf{h}} q_{\phi}(\mathbf{h}|\mathbf{v}) \log p_{\theta}(\mathbf{v}, \mathbf{h}) - \nabla_{\phi} \sum_{\mathbf{h}} q_{\phi}(\mathbf{h}|\mathbf{v}) \log q_{\phi}(\mathbf{h}|\mathbf{v}) \\ &= \sum_{\mathbf{h}} \log p_{\theta}(\mathbf{v}, \mathbf{h}) \nabla_{\phi} q_{\phi}(\mathbf{h}|\mathbf{v}) - \sum_{\mathbf{h}} (\log q_{\phi}(\mathbf{h}|\mathbf{v}) + 1) \nabla_{\phi} q_{\phi}(\mathbf{h}|\mathbf{v}) \\ &= \sum_{\mathbf{h}} (\log p_{\theta}(\mathbf{v}, \mathbf{h}) - \log q_{\phi}(\mathbf{h}|\mathbf{v})) \nabla_{\phi} q_{\phi}(\mathbf{h}|\mathbf{v}) \end{aligned} \quad (24)$$

where we have used the fact that

$\sum_{\mathbf{h}} \nabla_{\phi} q_{\phi}(\mathbf{h}|\mathbf{v}) = \nabla_{\phi} \sum_{\mathbf{h}} q_{\phi}(\mathbf{h}|\mathbf{v}) = \nabla_{\phi} 1 = 0$. Using the identity $\nabla_{\phi} q_{\phi}(\mathbf{h}|\mathbf{v}) = q_{\phi}(\mathbf{h}|\mathbf{v}) \nabla_{\phi} \log q_{\phi}(\mathbf{h}|\mathbf{v})$, then gives

$$\nabla_{\phi} \mathcal{L}(\mathbf{v}) = \mathbb{E}_q[(\log p_{\theta}(\mathbf{v}, \mathbf{h}) - \log q_{\phi}(\mathbf{h}|\mathbf{v})) \times \nabla_{\phi} \log q_{\phi}(\mathbf{h}|\mathbf{v})] \quad (25)$$

Backup IV

$$\mathbb{E}_q[\nabla_{\phi} \log q_{\phi}(\mathbf{h}|\mathbf{v})] = \mathbb{E}_q \left[\frac{\nabla_{\phi} q_{\phi}(\mathbf{h}|\mathbf{v})}{q_{\phi}(\mathbf{h}|\mathbf{v})} \right] \quad (26)$$

$$= \sum_{\mathbf{h}} \nabla_{\phi} q_{\phi}(\mathbf{h}|\mathbf{v}) \quad (27)$$

$$= \nabla_{\phi} \sum_{\mathbf{h}} q_{\phi}(\mathbf{h}|\mathbf{v}) \quad (28)$$

$$= \nabla_{\phi} 1 \quad (29)$$

$$= 0 \quad (30)$$

Backup V

Algorithm 1 Compute gradient estimates.

$$\Delta\theta \leftarrow 0, \Delta\phi \leftarrow 0, \Delta\lambda \leftarrow 0, \mathcal{L} \leftarrow 0$$
for $t \leftarrow 1$ **to** T **do**
$$\mathbf{h}_t \sim q_\phi(\mathbf{h}_t|\mathbf{v}_t)$$
$$l_t \leftarrow \log p_\theta(\mathbf{v}_t, \mathbf{h}_t) - \log q_\phi(\mathbf{h}_t|\mathbf{v}_t)$$
$$\mathcal{L} \leftarrow \mathcal{L} + l_t$$
$$l_t \leftarrow l_t - C_\lambda(\mathbf{v}_t)$$
end for
$$c_b \leftarrow \text{mean}(l_1, \dots, l_T)$$
$$v_b \leftarrow \text{variance}(l_1, \dots, l_T)$$
$$c \leftarrow \alpha c + (1 - \alpha)c_b$$
$$v \leftarrow \alpha v + (1 - \alpha)v_b$$
for $t \leftarrow 1$ **to** T **do**
$$l_t \leftarrow \frac{l_t - c}{\max(1, \sqrt{v})}$$
$$\Delta\theta \leftarrow \Delta\theta + \nabla_\theta \log p_\theta(\mathbf{v}_t, \mathbf{h}_t)$$
$$\Delta\phi \leftarrow \Delta\phi + l_t \nabla_\phi \log q_\phi(\mathbf{h}_t|\mathbf{v}_t)$$
$$\Delta\lambda \leftarrow \Delta\lambda + l_t \nabla_\lambda C_\lambda(\mathbf{v}_t)$$
end for