

Learning Long Sequences with TRBM

James Bergstra

Aug 7-11 2007

I'm researching recurrent networks as probability models of sequential data.

Data

- ultimately : audio, video, language, music, robotic control
- currently: toy bit sequences, symbolic melodies

Models

- recurrent nets with regularization to improve BPTT
- temporally hierarchical nets
- partly linear nets

How to detect/model long-term and multi-scale patterns???

Temporal Restricted Boltzmann Machines

- 1 What they are
- 2 How they are trained
- 3 One way things can go wrong when learning:
 - long sequences
 - many non-overlapping sequences
 - sequences with long-term statistical dependencies
- 4 Ideas to get around this problem (some results)

The ones with temporal connections between hidden units.

Ask for picture!

3 Phase Learning Algo

- 1 Optimize W as RBM of non-temporal $P(v)$
 - 1 Project sequence into observation space.
 - 2 Learn (RBM) density model of projected points.
- 2 Optimize U for given $z^{1..T}$
 - Choosing $z^{(t)} = E_W[h|v^{(t)}]$
- 3 Continue to optimize W, U to model $v^{(1..t)}$
 - W by Contrastive Divergence
 - U by backprop of bias gradient

How United States of America Goes Wrong

$v =$ "united states of america united states of america ..."

united \rightarrow 00 *america* \rightarrow 01
of \rightarrow 10 *states* \rightarrow 11

- Phase 1 Suppose the RBM learns the mapping above
- Phase 2 First component of z must perform XOR with 1-layer net ... **not possible!**
- Phase 3 W, U are nowhere near a solution, better to restart with joint optimization.

Semantics of z_t in TRBM

$$\begin{aligned} & \log P(v^{(1)}, v^{(2)}, \dots, v^{(t)}, \dots) \\ &= \sum_t \log P(v^{(t)} | z^{(t-1)}) \quad \textit{hidden markov assumption} \\ &\propto \sum_t \log \sum_h e^{-(v^{(t)}W + z^{(t-1)}U)h} \quad \textit{RBM} \\ &\propto \sum_t e^{-\textit{freeEnergy}(v^{(t)} | z^{(t-1)})} \end{aligned}$$

z_t must be predictive of **FUTURE** $v^{(t+1)}$, $z^{(t+1)}$
(within constraints imposed by functional form of z_t)

Why is Phase 2 hard?

Phase 2 begins with $z^{(t)} = \mathbb{E}_W[h|v^{(t)}]$, and **tries** to solve for U :

$$\text{sigm}(Uz^{(1)} + b) = z^{(2)}$$

$$\text{sigm}(Uz^{(2)} + b) = z^{(3)}$$

...

$$\text{sigm}(Uz^{(T-1)} + b) = z^{(T)}$$

Shallow: When T is greater than the number of dimensions of $z^{(t)}$, then no solution generally exists (linear separability).

Deep: *Many* (most?) trajectories of length $> K$ through $\{0, 1\}^K$ are not possible with an iterated system of form $z^{(t)} = \text{sgn}(Uz_{(t-1)} + b)$, though long trajectories exist.

What to do?

Avoid the problem:

- Use sufficiently large Z vectors, jump straight to phase 3
- Add fresh units to the system for phase 3

[Try to] solve the problem:

- Decouple $z^{(t)}$ from $P(h|v^{(t)})$, so that $z^{(t)} = f(v^{(t)}, z^{(t-1)})$
 - Differentiable f enables optimization by BPTT (really!)
 - I have some ideas (Yoshua too!) to improve BPTT
 - some preliminary results... (ongoing work)

Questions? Comments?