

# ① Energy models

①

- ① • definition
- ICA case, generative vs. energy model

- ② Learning
  - exact (hybrid sampling)
  - CD
  - score matching

## ③ DBN

-RBM,

- complementary prior  $\leftrightarrow$  infinite BN w. tied weights

$\leftrightarrow$  RBM

- greedy learning
- wake-sleep phase

④ PoE, FoE

⑤ Further reading

④ Learning in energy models

$p^0(x)$  real distribution of the observed data

$p^{\infty}(x; \theta)$  model distribution =  $\frac{1}{Z(\theta)} \exp(-E(x; \theta))$

we want to minimize

$$KL(p^0 \parallel p_{\theta}^{\infty}) = \int p^0(x) \log \frac{p^0(x)}{p^{\infty}(x)} dx$$

$$= cte + \cancel{\log Z(\theta)} + \langle E(x; \theta) \rangle_{p^0}$$

$$\frac{\partial KL(p^0 \parallel p_{\theta}^{\infty})}{\partial \theta} = \cancel{\frac{\partial \log Z(\theta)}{\partial \theta}} + \left\langle \frac{\partial E(x; \theta)}{\partial \theta} \right\rangle_{p^0}$$

$$\frac{\partial \log Z(\theta)}{\partial \theta} = \frac{1}{Z(\theta)} \frac{\partial}{\partial \theta} \int \exp(-E(x; \theta)) dx$$

$$= \int \frac{1}{Z(\theta)} \exp(-E(x; \theta)) \left( -\frac{\partial E(x; \theta)}{\partial \theta} \right) dx$$

$$= - \left\langle \frac{\partial E(x; \theta)}{\partial \theta} \right\rangle_{p_{\theta}^{\infty}}$$

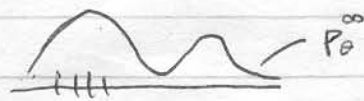
$$\Rightarrow \frac{\partial KL}{\partial \theta} = \left\langle \frac{\partial E(x; \theta)}{\partial \theta} \right\rangle_{p^0} - \left\langle \frac{\partial E(x; \theta)}{\partial \theta} \right\rangle_{p_{\theta}^{\infty}}$$

estimate using  
MCMC (Hamiltonian)

- + consistent
- + easily adaptable to more complex models
- it takes long to reach equilibrium dist.
- variance of MCMC estimator is usually high

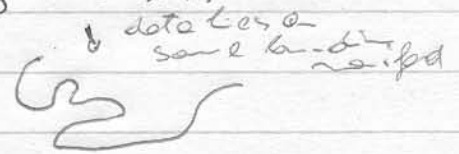
# 5) Catastrophic Divergence

Two ideas: 1) start the MC at the data distr.  $p^0$  rather than at some vague distr.



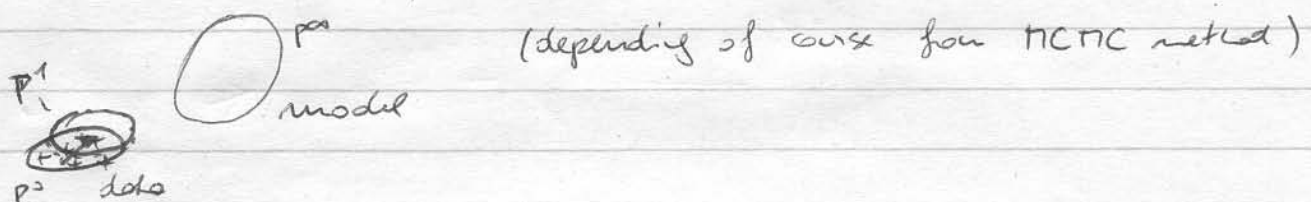
there is no pressure to remove empty modes

~~constant~~  
 => it's hard to estimate  
 the value of this mode (normalizing constant)



2) run the MC for only a few iterations rather than until equilibrium

idea: if the model is correct  $p_\theta^* = p^0$ , then MCMC does not change the distr. interpretation: remove any consistent tendency to move away from the data distribution



$$\Delta \theta \propto - \left\langle \frac{\partial E(x)}{\partial \theta} \right\rangle_{p^0} + \left\langle \frac{\partial E(x)}{\partial \theta} \right\rangle_{p^n}$$

$n$  typically very small (1)

## CO learning

1. compute  $\frac{\partial E(x)}{\partial \theta}$ , average over data set  $\underline{x}_k$
2. run MCMC sampler for  $m$  steps starting at  $\underline{x}_k \rightarrow \underline{z}_k = T^m \underline{x}_k$
3. compute  $\frac{\partial E(x)}{\partial \theta}$ , average over  $\underline{z}_k$
4. update parameters using

$$\Delta w_{ij} = - \frac{\eta}{N} \left( \sum_{\text{data}} \frac{\partial E(x_k)}{\partial \theta} - \sum_{\text{samples}} \frac{\partial E(z_k)}{\partial \theta} \right)$$

If  $p^0$  is in the space of  $p_\theta^*$ , and the MC mixes good enough  $\Rightarrow$  fixed point of PL solution

In general, it has been shown that there is a bias, but in practice it is small.

6

The update rule correspond to an approximate gradient descent step on obj. function

$$CD = KL(p^0 || p^0) - KL(p^n || p^0)$$

$$\frac{\partial CD}{\partial \theta} = \left\langle \frac{\partial E(x)}{\partial \theta} \right\rangle_{p^0} - \left\langle \frac{\partial E(x)}{\partial \theta} \right\rangle_{p^n} - \frac{\partial KL(p^n || p^0)}{\partial p^n} \frac{\partial p^n}{\partial \theta}$$

effect of change of parameters on the KL

→ empirically found to be small

one can do the same also if there are latent vars!

Example: CD learning of RBM parameters  
give solution for RBM

Score Matching

An easy way to learn a subset of energy models without performing MCMC

$$p(x; \theta) = \frac{1}{Z(\theta)} q(x; \theta)$$

$$\psi(x; \theta) := \nabla_x \log p(x; \theta) = \begin{pmatrix} \frac{\partial \log p(x; \theta)}{\partial x_1} \\ \vdots \\ \frac{\partial \log p(x; \theta)}{\partial x_n} \end{pmatrix}$$

"score function" in paper term

$$= \nabla_x \log q(x; \theta) \rightarrow \text{does not depend on } Z(\theta)$$

$$\psi_0(x) := \nabla_x \log p^0(x)$$

Score matching:

$$\text{minimize } J(\theta) := \frac{1}{2} \int dx p^0(x) \|\psi(x; \theta) - \psi_0(x)\|^2$$

$$\hat{\theta} = \arg \min_{\theta} J(\theta)$$

this term looks difficult to compute

\* makes score function of model similar to that of data → not clear what does that mean if  $\psi_0$  is not in  $\psi(\cdot, \theta)$  space

⑦ Theorem 1

If  $\psi(x; \theta)$  differentiable + other reg. conditions  
 $p^0$  diff,  $E_x(\|\psi(x; \theta)\|^2)$  finite  
 and  $E_x(\|\psi_0(x)\|^2)$  finite  
 $p^0(x) \psi(x; \theta) \xrightarrow{\|\theta\| \rightarrow \infty} 0$   $\forall \theta$

$$(2) \Rightarrow J(\theta) = \int p^0(x) \sum_{i=1}^m \left( \underbrace{\partial_i \psi_i(x; \theta)}_{\frac{\partial^2 \log q(x; \theta)}{\partial x_i^2}} + \frac{1}{2} \psi_i(x; \theta)^2 \right) dx + c\theta$$

$J(\theta)$  depends only on simple expectations of certain functions of the non-normalized pdf

Proof

$$J(\theta) = \frac{1}{2} \int p^0(x) \left( \|\psi(x; \theta)\|^2 + \underbrace{\|\psi_0(x)\|^2}_{\text{does not depend on } \theta} - 2 \psi(x; \theta)^T \psi_0(x) \right) dx$$

$$= \int p^0(x) \left( \underbrace{\frac{1}{2} \sum_i \psi_i(x; \theta)^2}_{\text{second term above}} - \underbrace{\psi(x; \theta)^T \psi_0(x)}_* \right) dx$$

$$* = \sum_i \int p^0(x) \psi_i(x; \theta) \psi_i^0(x) dx$$

$$= \sum_i \int p^0(x) \frac{\partial \log p^0(x)}{\partial x_i} \psi_i(x; \theta) dx$$

$$= \sum_i \int \frac{p^0(x)}{p^0(x)} \frac{\partial p^0(x)}{\partial x_i} \psi_i(x; \theta) dx \quad \int u'v = uv - \int uv'$$

$$= - \int p^0(x) \partial_i \psi_i(x; \theta) dx \quad \square$$

partial integration, one index of the line

For energy models  $\psi_i(x; \theta) = \frac{\partial}{\partial x_i} \log \exp(-E(x))$

$$\Rightarrow J(\theta) = \left\langle \sum_i \frac{\partial^2 E(x)}{\partial x_i^2} + \frac{1}{2} \left( \frac{\partial E(x)}{\partial x_i} \right)^2 \right\rangle_{p^0}$$

⑧

Theorem 2 Assume 1)  $p^\circ(\cdot) = p(\cdot; \theta^*)$  for some  $\theta^*$   
(almost everywhere)

2) the model is non-degenerate, i.e.

$$\nexists \theta^{**} : p(\cdot; \theta^{**}) = p(\cdot; \theta^*)$$

3)  $\forall x, \theta : q(x; \theta) > 0$

$$\Rightarrow J(\theta) = 0 \Leftrightarrow \theta = \theta^*$$

Proof:  $J(\theta) = 0 \Rightarrow \psi^\circ(\cdot) = \psi(\cdot; \theta)$   
 $\stackrel{q^\circ(x) > 0}{\Rightarrow p^\circ(x) > 0}$

$$\Rightarrow \log p^\circ(\cdot) = \log p(\cdot; \theta) + \underbrace{c}_{=0 \text{ (pdfs)}} \quad \square$$

$\Rightarrow$  (local) existence

- In which sense  $p_\theta$  becomes more similar to  $p^\circ$  if  $p^\circ$  does not lie in the  $p_\theta$ 's space?

Eq. 2  $\Rightarrow$  To minimize  $J$ , the first term should be negative  $\Rightarrow$  maximum of  $\log p(x; \theta)$   
 { second term  $\Rightarrow$  as steep a maximum as possible

- CD is more general since it is applicable to latent variable models

extensions: binary vars ("ratio matching")  
 non negative data

## Example

① Multivariate Gaussian density

$$p(x; \Pi, \mu) = \frac{1}{Z(\Pi, \mu)} \exp -\frac{1}{2} (x-\mu)^T \Pi (x-\mu)$$

$$\psi(x; \Pi, \mu) = \nabla_x -\frac{1}{2} (x-\mu)^T \Pi (x-\mu)$$

$$= -\Pi (x-\mu)$$

$$\partial_i \psi(x) = -m_{ii}$$

$$\Rightarrow \tilde{J}(\Pi, \mu) = \frac{1}{T} \sum_{t=1}^T \left( \sum_i -m_{ii} + \frac{1}{2} (x_t - \mu)^T \Pi^T \Pi (x_t - \mu) \right)$$

$$\nabla_{\mu} \tilde{J} = \frac{1}{T} \sum_{t=1}^T \Pi (x_t - \mu)$$

$$\Pi \mu - \Pi \frac{1}{T} \sum_{t=1}^T x_t$$

$$= 0 \Leftrightarrow \mu = \frac{1}{T} \sum_{t=1}^T x_t$$

$\Pi$  pos def

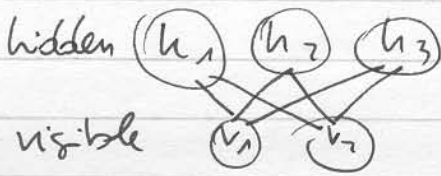
$$\nabla_{\Pi} \tilde{J} = -I + \frac{1}{2T} \sum_{t=1}^T (x_t - \mu)(x_t - \mu)^T$$
$$+ \frac{1}{2T} \left( \sum_{t=1}^T (x_t - \mu)(x_t - \mu)^T \right) \Pi$$

$$= 0 \Leftrightarrow \Pi^{-1} = \frac{1}{T} \sum_{t=1}^T (x_t - \mu)(x_t - \mu)^T$$

estimators are the the same as ML  
for any sample, not just asymptotically

② ICA (model, results)

# 9e Restricted Boltzmann Machine



bipartite: no hid-hid  
or vis-vis connections  
⇒ conditional independence

$$E(\underline{v}, \underline{h}) = - \sum_{ij} v_i h_j w_{ij} - \sum_i b_i^v v_i - \sum_j b_j^h h_j$$

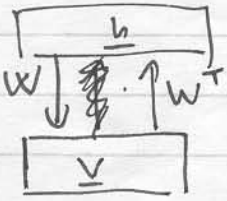
activation rules:

$$p(h_j = 1) = \sigma(b_j^h + \sum_i v_i w_{ij}), \quad \sigma(x) = \frac{1}{1 + e^{-x}} \text{ sigmoid}$$

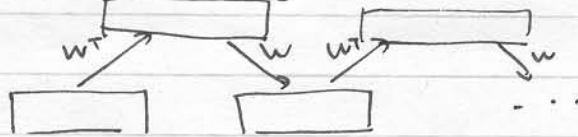
$$p(v_i = 1) = \sigma(h_i^v + \sum_j h_j w_{ij})$$

implement Gibbs sampling on  $\frac{1}{Z} \exp(-E(x))$

Inference: given  $\underline{v}$ , one can generate  $h_i$  independently because of conditional indep  
same  $w, b$



sample from equilibrium distribution  $p_w(\underline{v}, \underline{h})$   
by alternating Gibbs



Learning → arg max  $\log \prod_n p(\underline{v}^{(n)} | W, b) = \arg \max \log \prod_n \sum_{\underline{h}} p(\underline{v}^{(n)}, \underline{h} | W, b)$

$$\text{ML learning: } \frac{\partial}{\partial w_{ij}} \sum_n \ln p(\underline{v}^{(n)} | W, b) = \langle v_i h_j \rangle_{p^0} - \langle v_i h_j \rangle_{p^{\infty}}$$

$$\frac{\partial}{\partial b_i^v} \dots = \langle v_i \rangle_{p^0} - \langle v_i \rangle_{p^{\infty}}$$

$$\frac{\partial}{\partial b_j^h} \dots = \langle h_j \rangle_{p^0} - \langle h_j \rangle_{p^{\infty}}$$

can be efficiently learning using CD



$$E(\underline{v}, \underline{h}) = - \sum_i b_i^v v_i - \sum_j b_j^h h_j - \sum_{ij} v_i h_j w_{ij}$$

$$p(\underline{v}, \underline{h}) = \frac{1}{Z} \exp(-E(\underline{v}, \underline{h}))$$

$$p(\underline{v}) = \frac{1}{Z} \sum_{\underline{h}} \exp(-E(\underline{v}, \underline{h}))$$

ML learning

$$\max_{\underline{b}, \underline{w}} \log \prod_n p(\underline{v}^{(n)} | \underline{w}, \underline{b})$$

$$= \sum_n \ln \sum_{\underline{h}} p(\underline{v}^{(n)}, \underline{h} | \underline{w}, \underline{b})$$

$$= \sum_n \left( \ln \sum_{\underline{h}} \exp(-E(\underline{v}, \underline{h})) - \ln Z(\underline{w}) \right)$$

$$\frac{\partial}{\partial w_{ij}} \ln Z(\underline{w}) = \frac{1}{Z(\underline{w})} \sum_{\underline{v}, \underline{h}} \exp(-E(\underline{v}, \underline{h})) \cdot \left( - \frac{\partial}{\partial w_{ij}} E(\underline{v}, \underline{h}) \right)$$

$$= P(\underline{v}, \underline{h} | \underline{w}, \underline{b})$$

$$v_i h_j$$

$$= \sum_{\underline{v}, \underline{h}} P(\underline{v}, \underline{h} | \underline{w}, \underline{b}) \cdot v_i h_j$$

$$= \langle v_i h_j \rangle_{\text{model}}$$

$$\frac{\partial}{\partial w_{ij}} \ln \sum_{\underline{h}} \exp(-E(\underline{v}^{(n)}, \underline{h}))$$

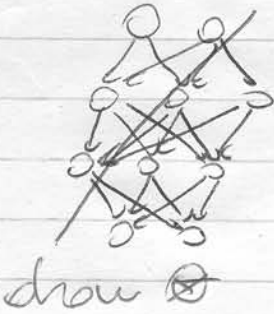
$$= \sum_{\underline{h}} \frac{1}{\sum_{\underline{h}} \exp(-E(\underline{v}^{(n)}, \underline{h}))} \exp(-E(\underline{v}^{(n)}, \underline{h})) \cdot v_i h_j$$

$$= \sum_{\underline{h}} P(\underline{h} | \underline{w}, \underline{b}, \underline{v}^{(n)})$$

$$\Rightarrow \begin{cases} \frac{\partial}{\partial w_{ij}} \sum_n \ln p(\underline{v}^{(n)} | \underline{w}, \underline{b}) = \langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{model}} \\ \frac{\partial}{\partial b_i^v} \text{---} = \langle v_i \rangle_{\text{data}} - \langle v_i \rangle_{\text{model}} \\ \frac{\partial}{\partial b_j^h} \text{---} = \langle h_j \rangle_{\text{data}} - \langle h_j \rangle_{\text{model}} \end{cases}$$

# (10) Deep Belief Networks

Energy models and CD can be used to learn efficiently directed belief networks with many layers



logistic belief network

binary units

generate:

$$p(A_i^{(e)} = 1 \mid \mathbf{h}_m^{(e+1)}) = \frac{1}{1 + \exp(-b_i - \sum_j h_j^{(e+1)} w_{ij}^e)}$$

Learning is very difficult:

1 - The probability of the observed data is a very complicated fct of the params

2 - inference is a nightmare because of explaining away  $\rightarrow$  the posterior dist. is typically intractable

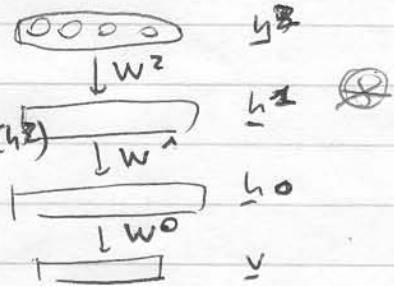
(MCMC, variational methods)

It would be nice to learn one layer at a time in a greedy way, but:

$$\textcircled{1} p(\mathbf{x}, \mathbf{h}^0) = p(\mathbf{x} \mid \mathbf{h}^0) p(\mathbf{h}^0)$$

$$\sum_{h_2, h_3} p(h^0 \mid h^1) p(h^1 \mid h^2) p(h^2)$$

need to integrate over all possible configurations of higher layers to get the prior  $\rightarrow$  p. 11

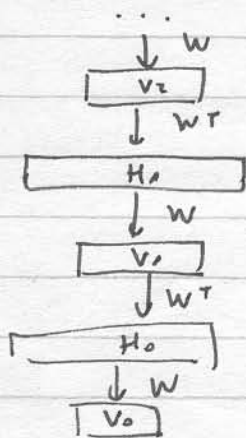


② the weights interact

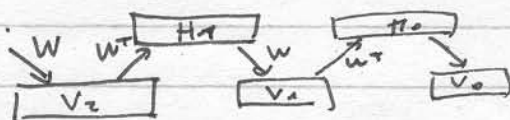
③ the posterior is going to be intractable, anyway

(12) Let's consider a simpler architecture where this problems are not an issue, and then relax that ~~constraints~~ back to the general case. Maybe that can give us a clever way to initialize ~~the~~ the model's parameters.

Infinite directed model with tied weights



generating from this model is equivalent to let an RBM with connections  $W$  reach its equilibrium distribution:



$\Rightarrow$  the class of distribution that you can model with both is the same

Learning on infinite model w. tied weight  $\equiv$  learn a RBM

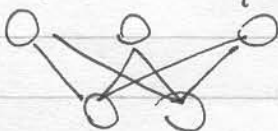
~~then~~  $\Rightarrow$  prior on  $h_0$  is complementary, inference is easy (use  $W^T$ )

alternative derivation: let's say we want to have a complementary prior on  $h_0$ .

$$\Rightarrow p(h_0 | v_0) = \prod_i p(h_i | v_0)$$

$$p(v_0 | h_0) = \prod_i p(v_i | h_0) \quad \text{by construction}$$

as an undirected graphical model these conditional indep. correspond to an RBM



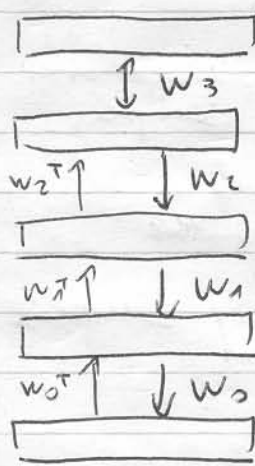
start from that, unroll it  $\rightarrow$  infinite model

(13)

we've seen: infinite  $\leftrightarrow$  RBM

## Greedy learning for DBN

consider hybrid model



} equivalent to assuming  $\infty$ -many layers with tied weights

1. Learn  $W_0$  assuming all weight matrices are tied (RBM)
2. Freeze  $W_0$  and commit ourselves to using  $W_0^T$  to infer ~~the state~~ the state in 1st hidden layer even if subsequent changes in higher-level weights mean that this inference method is no longer correct
3. Keeping all higher-weight matrices tied to each other, learn a model for higher-level data that is represented in  $h_0$

## "learning by rerepresenting"

Back-fitting with up-down algorithm

greedy learning is efficient but not optimal.

$\Rightarrow$  untie the recognition weights from the generative weights, retain the restriction that the posterior must be approx. by a factorial dist.

Refine the ~~the~~ weights using a CD variant of the wake-sleep algorithm:

- ① up-pass: use recognition weights to pick states for hidden variables; adjust generative weights:

$$\Delta W_{ij}^{n+1} \propto \left\langle h_j^{n+1} (h_j^n - \hat{h}_j^n) \right\rangle$$

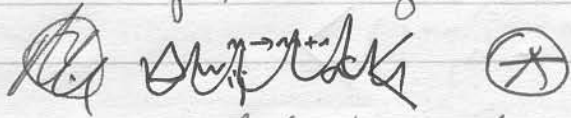
collected in up-pass

average over sample states  
reconstructed by sampling  $\rightarrow$  probability (\*)

The weight at the top are learned as before using RBM

110

② down-pass: start with state of the top-level associative memory, use generative weights to get states for hidden and visible layers; adjust recognition weights

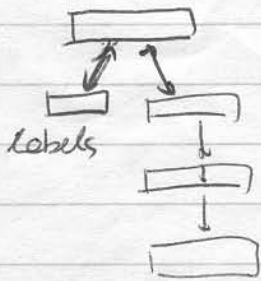


equivalent to wake-sleep if state taken from equilibrium dist. of top RBM. Here: init top RBM w. up-pass, run only a few iterations of gibbs sampling, then start down step

+ faster

+ eliminate node averaging

### Classification using DBN



softmax labels: exactly one unit is 1  

$$P_i = \frac{\exp(x_i)}{\sum_j \exp(x_j)} \rightarrow \text{probability of picking } i$$

The learning rules are unaffected by the competition

→ website demo!

performance on MNIST database: 1,25%

SVM 1,4% now even FDA 1,4%

best: 0,5% (I think)