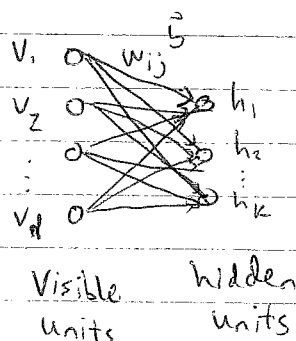


## Deep Belief Networks

Recall ~~a~~ a standard layer of an MLP:

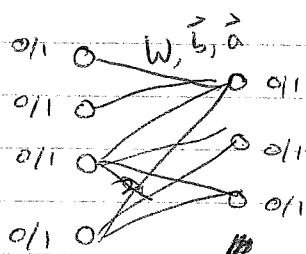


$d \times k$  weight matrix  $W$

Sigmoid activation

$$h_i = \sigma(W_i^T V + b_i)$$

## Restricted Boltzmann Machines



Standard RBM: assume  
visible and hidden units are  
binary.

Note  $d+k$  units  $\Rightarrow 2^{d+k}$  different settings of the variables.

Let  $\Theta = \{W, \vec{b}, \vec{a}\}$

Define energy function : 
$$E(v, h) = - \sum_{i=1}^d a_i v_i - \sum_{j=1}^k b_j h_j - \sum_{i=1}^d \sum_{j=1}^k v_i w_{ij} h_j$$

$$= -\vec{a}^T \vec{v} - \vec{b}^T \vec{h} - \vec{v}^T W \vec{h}$$

Then write

$$p(v, h; \Theta) = \frac{1}{Z} \cdot e^{-E(v, h)} \propto e^{-E(v, h)}$$

$$Z = \sum_v \sum_h e^{-E(v, h)} \leftarrow \text{sum has } 2^{d+k} \text{ terms in it.}$$

Once you have joint prob, can mess around with it:

$$p(\mathbf{v}; \theta) = \frac{1}{Z} \sum_h e^{-E(\mathbf{v}, \mathbf{h})}$$

$$p(\mathbf{v}|\mathbf{h}; \theta) \quad p(\mathbf{h}|\mathbf{v}; \theta)$$

Not entirely trivial fact  $\rightarrow$

$$p(\vec{v}|\vec{h}; \theta) = \prod_{i=1}^d p(v_i|h_i; \theta)$$

For example...

$$p\left(\begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \middle| \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}\right) = p(v_1=1 | \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}) \cdot p(v_2=1 | \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}) \cdot p(v_3=0 | \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix})$$

$$\text{Similarly - } p(\vec{h}|\vec{v}; \theta) = \prod_{j=1}^k p(h_j|\vec{v}; \theta) \cdot p(v_4=1 | \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix})$$

And, further...  $p(h_j=1|\vec{v}) = \sigma(w_j^T \mathbf{v} + b_j)$

$$p(v_i=1|\vec{h}) = \sigma(\mathbf{h}^T \mathbf{w}_i + a_i)$$

$\rightarrow$  compare to a feed-forward network

- ~~output~~ output of sigmoid treated as an actual value, **not** as prob. value in FFM
- In the probabilistic case, we can "go backwards" from hidden to visible.



Can sample data once we know  $\theta$  and  $\mathbf{h}$ !

See derivation page at end, or look in textbook for

How to learn parameters  $\theta$ ?

Maximum Likelihood!

For a single visible unit  $\vec{v}$

$$p(\vec{v}) = \frac{1}{Z} \sum_h e^{-E(\vec{v}, \vec{h})}$$

For a sequence of visible units  $\vec{v}^{(1)}, \dots, \vec{v}^{(n)}$

$$p(V) = \prod_{p=1}^n \frac{1}{Z} \sum_h e^{-E(\vec{v}^{(p)}, \vec{h})}$$

Treat this as a function of  $\theta$  and maximize w.r.t.  $\theta$ .

Take logs -

$$\log p(V) = \sum_{p=1}^n \left[ \log \left( \sum_h e^{-E(\vec{v}^{(p)}, \vec{h})} \right) - \log \left( \sum_v \sum_h e^{-E(v, h)} \right) \right]$$

$Z$  is not a "constant" w.r.t.  $\theta$ ! Problematic!

~~1st term~~

$$\frac{\partial}{\partial w_{ij}} E(v, h) = ? \quad -v_i h_j$$

1st Term of Log Likelihood:

$$\frac{\partial}{\partial w_{ij}} \sum_{i=1}^n \log \left( \sum_h e^{-E(\vec{v}^{(i)}, \vec{h})} \right)$$

$$= \sum_{p=1}^n \frac{\partial}{\partial w_{ij}} \log \left( \sum_h e^{-E(\vec{v}^{(p)}, \vec{h})} \right)$$

$$= \sum_{p=1}^n \left[ \frac{e^{-E(\vec{v}^{(p)}, \vec{h})}}{\sum_h e^{-E(\vec{v}^{(p)}, \vec{h})}} \cdot \sum_h \left[ \frac{\partial}{\partial w_{ij}} E(\vec{v}^{(p)}, \vec{h}) \right] \right]$$

(ordered thing is exactly  $-p(\vec{h}|\vec{v}^{(p)})$ :

$$p(\vec{h}|\vec{v}^{(p)}) = \frac{p(\vec{v}^{(p)}, \vec{h})}{p(\vec{v})} = \frac{\frac{1}{Z} \cdot e^{-E(\vec{v}^{(p)}, \vec{h})}}{\frac{1}{Z} \sum_{\vec{h}} e^{-E(\vec{v}^{(p)}, \vec{h})}}$$

So... 1st term is exactly

$$\sum_{p=1}^n \sum_{\vec{h}} p(\vec{h}|\vec{v}^{(p)}) \cdot v_i h_j \quad \leftarrow \text{This is an expectation of } v_i h_j \text{ wrt. } p(\vec{h}|\vec{v}^{(p)})$$

Similarly, 2nd term is

$$- \sum_{p=1}^n \sum_{\vec{v}, \vec{h}} p(\vec{v}, \vec{h}) \cdot v_i h_j \quad \leftarrow \text{also an expectation but wrt. joint probability}$$

Scale by  $n$ .

Hinton writes this as

$$E_{\text{data}}(v_i h_j) - E_{\text{model}}(v_i h_j)$$

GD. would operate as

$$W_{ij} \leftarrow W_{ij} - \epsilon [E_{\text{data}}(v_i h_j) - E_{\text{model}}(v_i h_j)]$$

Question is how to compute these quantities. Shorter answer - you cannot.

We can get an unbiased sample of this value. How?

Given  $\vec{v}^{(p)}$ , sampling  $p(\vec{h}|\vec{v}^{(p)})$  and then multiplying by  $v_i h_j$  gives you an unbiased sample.

Sample  $h_j$  via  $p(h_j=1|\vec{v}^{(p)})$  ~~and~~ and multiply by  $v_i$ . Average over all  $\vec{v}^{(p)}$ .

Now need sample of  $E_{\text{model}}(v|h_j)$ . This is a little trickier.  
Need to sample from  $p(v, h)$ .

Gibbs Sampling: To sample from a distribution  $p(v, h)$ , repeat the following:

Start with some value of ~~new~~  $v$  (in our case,  $\vec{v}^{(p)}$  is known)

Sample <sup>new  $h$</sup>  from  $p(h|v)$

~~Sample~~ Sample new  $v$  from  $p(v|h)$

repeat for a "long time"

Eventually  $v$  and  $h$  given by this algorithm will be a sample from  $p(v, h)$

[Need to establish 2 things — i. that once  $(v, h)$  samples are from the true  $p(v, h)$ , ~~the~~ <sup>any subsequent</sup> samples are also from  $p(v, h)$ ; ii. irrespective of the initial choice of  $v$ , the distribution of  $(v, h)$  approaches the desired distribution:  $p(v, h)$ .]

This is still slow, unfortunately. Hinton introduced a further simplification (contrastive divergence) that only takes a single step. So, the full alg. looks like this:

Average over a minibatch:

- pick a visible point  $v^{(p)}$
- Sample  $h$  via  $p(h|v^{(p)})$
- Sample  $\tilde{v}$  via  $p(\tilde{v}|h)$
- $\Delta w_{ij} = v_i^{(p)} \cdot h_j - \tilde{v}_i \cdot h_j$

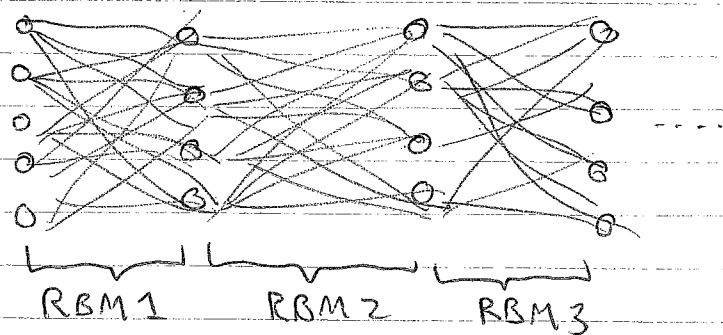
[Sometimes, instead of  $p(\tilde{v})$  use the actual probability  $p(\tilde{v}|h)$  when computing  $\Delta w_{ij}$ ]

Note: this procedure is known to not converge to a <sup>local</sup> max of the likelihood.

## Next Step: Deep Belief Networks

Idea: Stack several layers of RBMs on top of one another, and train layer-by-layer.

Picture:



Idea is this ... train weights of RBM 1 using visible/training data.

- Freeze weights of ~~RB~~ RBM 2.
  - Treat the hidden units of RBM 1 as the visible units of RBM 2.
  - Pass each training pt through RBM 2, sampling the visible units for RBM 2.
  - Train RBM 2 using these sampled data.
  - Freeze weights of RBM 2.
  - etc...
- "Pretraining"

Next, use the weights learned by this procedure as an initialization to a feedforward sigmoidal network.

"Fine-tuning": Run SGD with supervised data

## Extensions to basic RBMs

What if weights are non-binary?

Extension 1: binary  $\rightarrow$  1-of-k

For binary, we have something like

$$\sigma(x) = \frac{e^x}{1+e^x} = \frac{1}{1+e^{-x}}$$

Generalize to 1-of-k via softmax units

$$p_j = \frac{e^{x_j}}{\sum_{i=1}^k e^{x_i}}$$

Don't really have to change ~~anything~~ much at all.  
Energy function same. Only changes how samples are taken and probabilities computed.

$$p(h_j=1 | \vec{v}; \theta) = \sigma(\vec{v}^T \vec{w}_j + b_j)$$

$$\Rightarrow p(h_j=c | \vec{v}; \theta) = \frac{e^{\vec{v}^T \vec{w}_j + b_j}}{\sum_{j=1}^k e^{\vec{v}^T \vec{w}_j + b_j}}$$

Real visible units / Binary hidden units

$$E(v, h) = \sum_{i \in \text{vis}} \frac{(v_i - a_i)^2}{2\sigma_i^2} - \sum_{j \in \text{hid}} b_j h_j - \sum_{ij} \frac{v_i}{\sigma_i} \cdot h_j \cdot w_{ij}$$

Usual = assume data is whitened first, so  $\sigma_i = 1$

Then: we have

$$p(h_j = 1 | \vec{v}; \theta) = \text{same as before} \\ = \sigma(\vec{v}^T \mathbf{w}_j + b_j)$$

$$p(v_i = 1 | \vec{h}; \theta) \sim ?$$

Gaussian distribution - Exercise:

What is the precise Gaussian?  
(Need to compute for HW 5.)

Then G.D. works as follows:

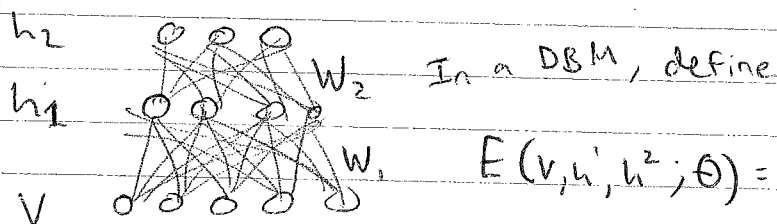
Average over minibatch:

- Given  $v^{(p)}$  sample  $h$  via  $p(h | v^{(p)})$  (sigmoid)
- Given  $h$  sample  $\tilde{v}$  via  $p(\tilde{v} | h)$  (Gaussian)
- $\Delta w_{ij} = v_i^{(p)} h_j - \tilde{v}_i h_j$

## ~~Convolutional Deep Boltzmann Networks (Deep DBNs)~~ Deep Boltzmann Machines

- When designing a joint probability distribution over a deep network, stacks of RBMs are not the cleanest solution.
- For instance, they don't have a simple joint distribution.

Alternative - Deep Boltzmann Machine  
Suppose there are 2 hidden layers



$$E(v, h_1, h_2; \theta) = -v^T w_1 h_1 - h_1^T w_2 h_2$$

$$\text{i.e. } p(v, \theta) = \frac{1}{Z} \sum_{h_1, h_2} \exp(-E(v, h_1, h_2; \theta))$$



$$p(h_j^1 = 1 | v, h^2) = \sigma\left(\sum_i W_{ij}^1 v_i + \sum_m W_{jm}^2 h_m^2\right)$$

$$p(h_m^2 = 1 | h^1) = \sigma\left(\sum_j W_{jm}^2 h_j^1\right)$$

$$p(v_i = 1 | h^1) = \sigma\left(\sum_j W_{ij}^1 h_j^1\right)$$

Can use sample-based methods for parameter learning, but  
 can also use Variational inference (won't discuss)

$$p(h|v; \theta) = \frac{p(v, h; \theta)}{p(v; \theta)}$$

Some Derivations

$$= \frac{1}{p(v; \theta)} \cdot \frac{1}{Z} \cdot e^{a^T v + b^T h + v^T W h}$$

$$\propto e^{b^T h + v^T W h}$$

$$= \exp \left( \sum_{j=1}^K b_j h_j + \sum_{j=1}^K v^T W_{:,j} h_j \right) \quad W_{:,j} = W_j$$

$$= \prod_{j=1}^K \exp(b_j h_j + v^T W_{:,j} h_j)$$

Follows that  $p(h|v; \theta)$  will actually factorize over the hidden units. In particular,

$$p(h_j = 1 | \vec{v}; \theta) = \frac{\tilde{p}(h_j = 1 | \vec{v}; \theta)}{\tilde{p}(h_j = 0 | \vec{v}; \theta) + \tilde{p}(h_j = 1 | \vec{v}; \theta)}$$

$$= \frac{\exp(b_j + v^T W_{:,j})}{1 + \exp(b_j + \vec{v}^T W_{:,j})}$$

$$= \sigma(w_j^T \vec{v} + b_j)$$

Now pull out  $e^{a^T v}$  term on num/denom

Another way:

$$\frac{p(v, h; \theta)}{p(v; \theta)} = \frac{\frac{1}{Z} \cdot e^{-E(v, h)}}{\frac{1}{Z} \cdot \sum_h e^{-E(v, h)}} = \frac{e^{-E(v, h)}}{\sum_h e^{-E(v, h)}} \quad \uparrow$$