

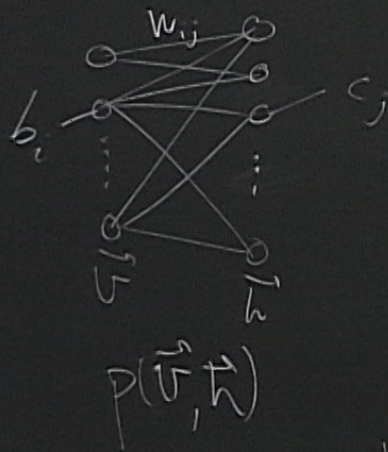
Title: PSI 2017/2018 - Machine Learning for Many Body Physics - Lecture 12

Date: Apr 23, 2018 09:00 AM

URL: <http://pirsa.org/18040064>

Abstract:

RBM:



$$E_{\lambda} = -\vec{b}^T \vec{v} - \vec{c}^T \vec{h} - \vec{v}^T \mathbf{W} \vec{h}$$
$$= - \sum_{i=1}^n b_i v_i - \sum_{j=1}^m c_j h_j - \sum_{ij} w_{ij} v_i h_j$$

$$\frac{\partial \log \mathcal{L}}{\partial \lambda} = - \sum_{\vec{h}} p(\vec{h} | \vec{v}) \frac{\partial E}{\partial \lambda} + \sum_{\vec{v}} p(\vec{v} | \vec{h}) \frac{\partial E}{\partial \lambda}$$

$$\lambda = (w, b, c)$$

$w_{ij} v_i h_j$

$$\sum_{\vec{v}, \vec{h}} p(\vec{v}, \vec{h}) \frac{\partial E}{\partial \lambda}$$

$$\begin{aligned} \frac{\partial \log \mathcal{L}}{\partial w_{ij}} &= \sum_{\vec{h}} p(\vec{h} | \vec{v}) h_i v_j - \sum_{\vec{v}} p(\vec{v}) \sum_{\vec{h}} p(\vec{h} | \vec{v}) h_i v_j \\ &= p(h_i = 1 | \vec{v}) v_j - \sum_{\vec{v}} p(\vec{v}) p(h_i = 1 | \vec{v}) v_j \end{aligned}$$

$$\lambda = (w, b, r)$$

In practice this is averaged over a mini-batch, size  $b \ll l$

$CD_k$  can be used

$$\frac{1}{b} \sum_{\vec{v} \in \text{mini batch}} \frac{\partial \mathcal{L}}{\partial w_{ij}} = \langle v_i h_j \rangle_{p(\vec{h}|\vec{v})} - \langle v_i h_j \rangle_{p(\vec{h}, \vec{v})}$$

$v_i$  is "clamped" to the data

↑  
expectation value for the stationary distribution of the model

This calculation has to be completed for all weights & biases

$$\nabla_{\lambda}(\log_2 p) = \begin{cases} \langle v h^T \rangle_{p(h|v)} - \langle v h^T \rangle_{p(v,h)} & \lambda = W \\ \langle v \rangle_{p(h|v)} - \langle v \rangle_{p(v,h)} & \lambda = b \\ \langle h \rangle_{p(h|v)} - \langle h \rangle_{p(v,h)} & \lambda = c \end{cases}$$

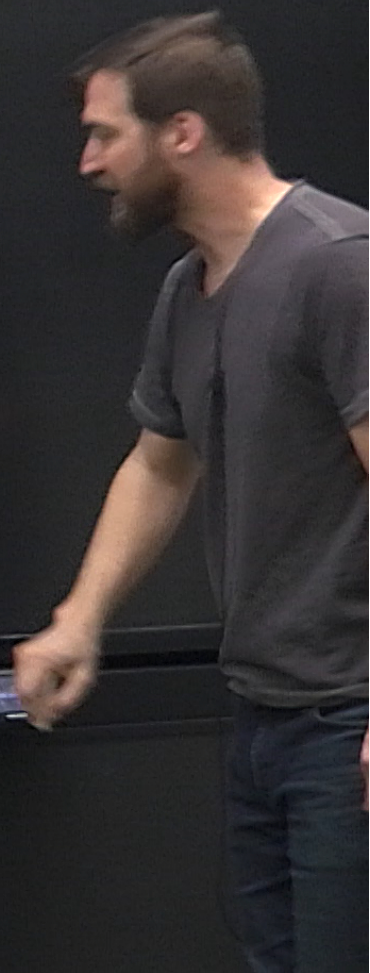
$$\text{and } \lambda' = \lambda - \eta \nabla_{\lambda} \text{KL}$$

Let's derive  $\sigma$  by considering  $p(\vec{h}|\vec{v})$   
 $p(\vec{h}|\vec{v}) = \frac{p(\vec{v}, \vec{h})}{p(\vec{v})}$  (product rule of Bayes)

Start with

$$p(\vec{v}) = \sum_{\vec{h}} p(\vec{v}, \vec{h}) = \frac{1}{Z} \sum_{\vec{h}} e^{(b^T \vec{v} + c^T \vec{h} + \vec{v}^T W \vec{h})}$$

$$= \frac{1}{Z} e^{b^T \vec{v}} \sum_{h_1=\{0,1\}} \sum_{h_2=\{0,1\}} \prod_{i=1}^n e^{(c^T + \vec{v}^T W) h_i}$$



$$= \frac{1}{Z} e^{b^T v} \sum_{h_1 \in \{0,1\}} e^{h_1 (c_1 + \sum_{j=1}^m W_{1j} v_j)} \sum_{h_2 \in \{0,1\}} e^{h_2 (\dots)} \dots \sum_{h_n \in \{0,1\}} \dots$$

$$= \frac{1}{Z} \prod_{j=1}^m e^{b_j v_j} \prod_{i=1}^n (1 + e^{(c_i + \sum_{j=1}^m W_{ij} v_j)})$$

$$= \frac{1}{Z} \exp \left[ b^T v + \sum_i \log (1 + e^{(c^T + v^T W)_i}) \right]$$

RBM is fitting  $F(\vec{v})$  to the "physical"  
Hamiltonian of the data set  $S$ .

Next:  $p(h|v) = \frac{p(v, h)}{p(v)} = \frac{1}{Z} \frac{1}{Z} e^{b^T v + c^T h + v^T W h}$



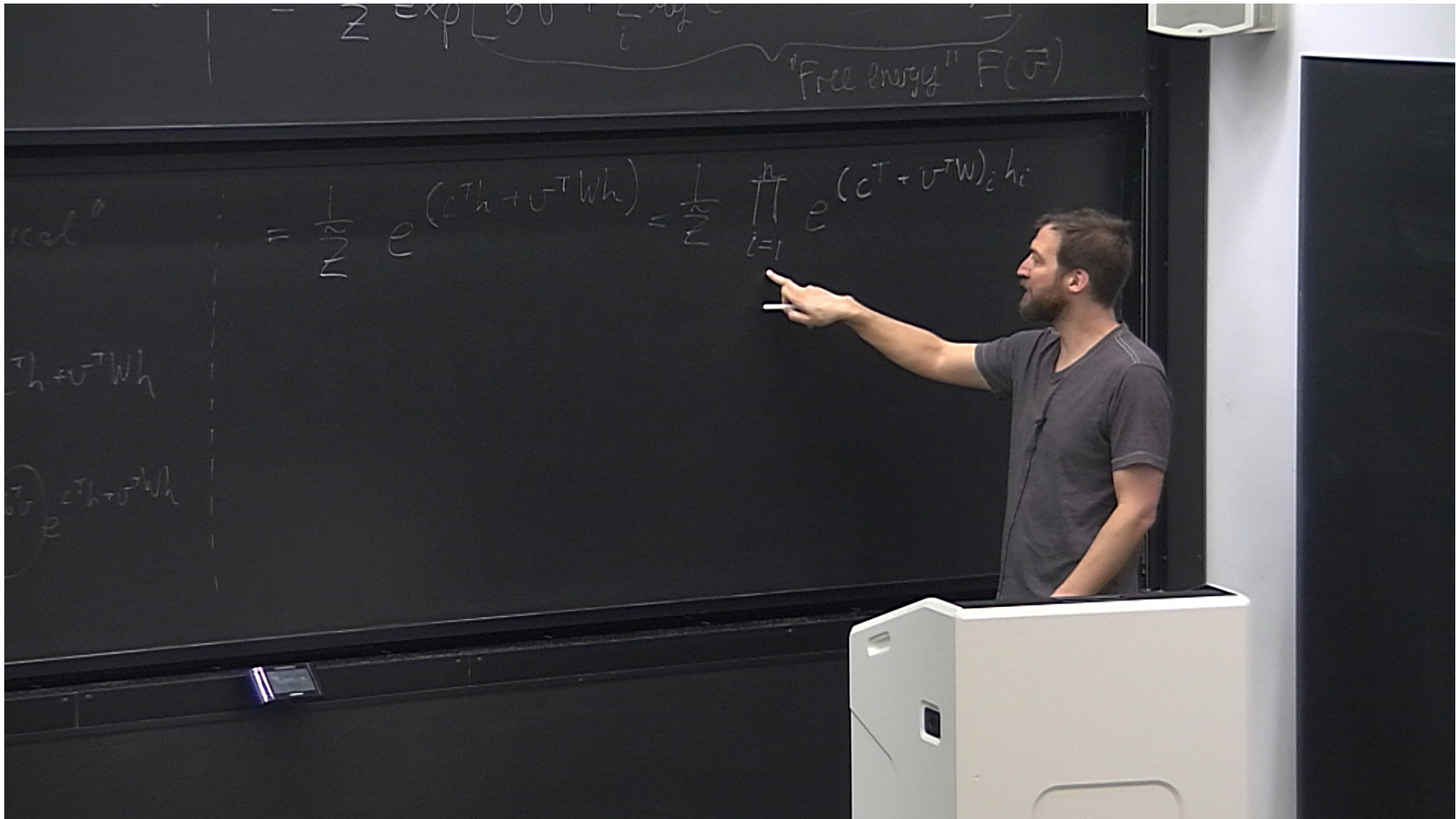
RBM is fitting  $F(\vec{v})$  to the "physical"  
Hamiltonian of the data set  $S$ .

$$\begin{aligned} \text{Next: } p(h|v) &= \frac{p(v, h)}{p(v)} = \frac{1}{p(v)} \frac{1}{Z} e^{b^T v + c^T h + v^T W h} \\ &= \left( \frac{1}{Z} e^{b^T v} \sum_h e^{c^T h + v^T W h} \right)^{-1} \frac{1}{Z} e^{b^T v} e^{c^T h + v^T W h} \end{aligned}$$

$h_1 = (20, 1)$   $h_2 = (20, 1)$

RBM is fitting  $F(\vec{v})$  to the "physical" Hamiltonian of the data set  $S$ .

Next: 
$$p(h|v) = \frac{p(v, h)}{p(v)} = \frac{1}{p(v)} \frac{1}{Z} e^{b^T v + c^T h + v^T W h}$$
$$= \left( \frac{1}{Z} e^{b^T v} \sum_h e^{c^T h + v^T W h} \right)^{-1} \frac{1}{Z} e^{b^T v} e^{c^T h + v^T W h}$$



"Free energy"  $F(\vec{v})$

$$= \frac{1}{Z} e^{(c^T h + v^T W h)} = \frac{1}{Z} \prod_{i=1}^n e^{(c^T + v^T W)_i h_i}$$

this is consistent with  $p(h_i | \vec{v}) = \prod_{i=1}^n p(h_i | \vec{v})$

note:  $p(h_i=1 | \vec{v}) + p(h_i=0 | \vec{v}) = 1$  gives normalization

$$A(e^{(c^T + v^T W)_i} + 1) = 1, \quad A = \frac{1}{1 + e^{(c^T + v^T W)_i}}$$

$$\text{so } p(h_i = 1 | \vec{v}) = \frac{1}{1 + e^{-(C^T + U^T W)_i}}$$
$$= \sigma\left(C_i + \sum_{j=1}^J W_{ij} V_j\right)$$

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

And similarly for  $p(v_i = 1 | \vec{h})$

"Block Gibbs" sampling one "block" ( $\vec{v}$  or  $\vec{h}$ ) conditioned on the other.

$$\vec{v}_0 \rightarrow \vec{h}_0 \rightarrow \vec{v}_1 \rightarrow \vec{h}_1 \rightarrow \vec{v}_2 \rightarrow \vec{h}_2 \dots$$

TRAINING : Contrastive divergence - used to approximate  $\langle \cdot \rangle_{p(\vec{v}, \vec{h})}$  in the stochastic gradient descent.

"Block Gibbs" sampling one "block" ( $\vec{v}$  or  $\vec{h}$ ) conditioned on the other.

$$\vec{v}_0 \rightarrow \vec{h}_0 \rightarrow \vec{v}_1 \rightarrow \vec{h}_1 \rightarrow \vec{v}_2 \rightarrow \vec{h}_2 \dots \vec{v}_k \rightarrow \vec{h}_k$$

TRAINING : Contrastive divergence - used to approximate  $\langle \cdot \rangle_{p(\vec{v}, \vec{h})}$  in the stochastic gradient descent.  
 Can calculate this by equilibrating the RBM and sampling  $\vec{v}, \vec{h}$

CD<sub>k</sub> can be used : This calculation has to be completed for all weights & biases

$$\langle v h^T \rangle_{p(\vec{v}, \vec{h})} - \langle v h^T \rangle_{p(v, h)} \quad \lambda = W$$

Much simpler: run this Gibbs chain only  
for  $k \sim \mathcal{O}(1)$  steps ( $k=1$  sometimes works)

- $\vec{U}_0$  initialized from a training example
- $\vec{h}_0$  obtained from  $p(\vec{h} | \vec{U}_0)$
- repeat for  $k$  steps



CD<sub>k</sub> expression for  $\nabla_{KL} B$

$$\frac{\partial \log p}{\partial \lambda} \approx - \sum_k p(k | \vec{v}_0) \frac{\partial E_0}{\partial \lambda} + \sum_k p(k | \vec{v}_k) \frac{\partial E_k}{\partial \lambda}$$

giving for example

$$W' = W + \frac{\eta}{b} \sum_{\vec{v}_k \text{ mini batch}} (\vec{v}_0 \vec{h}_0^T - \vec{v}_k \vec{h}_k^T)$$

## The trained RBM & sampling

Once trained we have a "model"

$P_{\lambda}(\vec{v}, \vec{h})$  such that

$$P_{\lambda}(\vec{v}) = \sum_{\vec{h} \in \mathcal{H}} P_{\lambda}(\vec{v}, \vec{h}) \approx q(\vec{v})$$

underlies data in  $S$

We sample through Block Gibbs sampling  $\vec{v}$  and  $\vec{h}$   
 (we are interested in  $\vec{v}$ ,  $\langle O_v \rangle$ )

Recall:  $T(\mu \rightarrow \nu) = g(\mu \rightarrow \nu) A_{\mu \rightarrow \nu}$

"transition prob."      "selection prob."      "acceptance ratio"

detailed balance:  $\frac{T(\mu \rightarrow \nu)}{T(\nu \rightarrow \mu)} = \frac{P_\nu}{P_\mu}$  ← stationary state of the RBM

$$L = \frac{n}{T} (c^T + v^T W) c h$$

$$P_{\mu}(\vec{v}) = \sum_{\vec{h}} P_{\mu}(\vec{v}, \vec{h}) \approx q(\vec{v})$$

↑  
underlies data in S

detailed balance:  $\frac{T(\mu \rightarrow \nu)}{T(\nu \rightarrow \mu)} = \frac{P_{\nu}}{P_{\mu}}$  ← stationary state of the RBM

Recall Metropolis:  $A(\mu \rightarrow \nu) = \min \left\{ 1, \frac{P_{\nu}}{P_{\mu}} \times \frac{q(\nu \rightarrow \mu)}{q(\mu \rightarrow \nu)} \right\}$

consider RBM when the visibles are updated:

$$A((\vec{v}, \vec{h}) \rightarrow (\vec{v}', \vec{h})) = \min \left\{ 1, \frac{p(\vec{v}', \vec{h})}{p(\vec{v}, \vec{h})} \times \frac{p(\vec{v}, \vec{h})}{p(\vec{v}', \vec{h})} \right\}$$

since  $\frac{p(\vec{v}, \vec{h})}{p(\vec{v}, \vec{h})} = \frac{\frac{p(\vec{v}, \vec{h})}{p(\vec{v})}}{\frac{p(\vec{v}', \vec{h})}{p(\vec{v}', \vec{h})}} = \frac{p(\vec{v}, \vec{h})}{p(\vec{v}', \vec{h})}$

So  $A((\vec{v}, \vec{h}) \rightarrow (\vec{v}', \vec{h})) = 1$

$$P_{\lambda}(\vec{v}) = \sum_{\vec{h}} P_{\lambda}(\vec{v}, \vec{h}) \approx q(\vec{v})$$

↑  
involves data in S

detailed balance:

$$\frac{T(u \rightarrow v)}{T(v \rightarrow u)} = \frac{P_v}{P_u}$$

stationary state  
of the RBM

Recall Metropolis:

$$A(u \rightarrow v) = \min \left\{ 1, \frac{P_v}{P_u} \times \frac{q(v \rightarrow u)}{q(u \rightarrow v)} \right\}$$

consider RBM when the visibles are updated:

$$A((\vec{v}, \vec{h}) \rightarrow (\vec{v}', \vec{h})) = \min \left\{ 1, \frac{p(v'; h)}{p(v, h)} \times \frac{p(v|h)}{p(v'|h)} \right\}$$

since  $\frac{p(v|h)}{p(v'|h)} = \frac{\frac{p(v, h)}{p(h)}}{\frac{p(v', h)}{p(h)}} = \frac{p(v, h)}{p(v', h)}$

so  $A((v, h) \rightarrow (v', h)) = 1$

Rest is just like regular MCMC

- warm up (equilibration)
- collect statistics

$$\langle Q_i \rangle \approx \frac{1}{M} \sum_{i=1}^M Q_i$$