

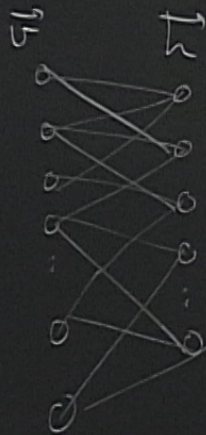
Title: PSI 2017/2018 - Machine Learning for Many Body Physics - Lecture 11

Date: Apr 20, 2018 09:00 AM

URL: <http://pirsa.org/18040063>

Abstract:

RBM



$$E_{\lambda} = - \sum_i b_i v_i - \sum_j c_j h_j - \sum_{ij} W_{ij} v_i h_j$$

$$Z_{\lambda} = \sum_{\mathbf{v}, \mathbf{h}} e^{-E_{\lambda}}$$

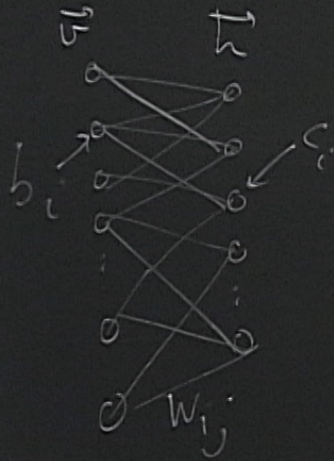
$$\lambda = (W, \vec{b}, \vec{c}) = \theta$$

Training & Sampling a RBM

First: Training

Assume you have a data set: $S = \{ \vec{x}_1, \vec{x}_2, \dots, \vec{x}_L \}$

RBM:



$$E_{\lambda} = - \sum_{i=1}^m b_i v_i - \sum_{j=1}^n c_j h_j - \sum_{ij} w_{ij} v_i h_j$$

$$Z_{\lambda} = \sum_{\mathbf{v}, \mathbf{h}} e^{-E_{\lambda}}$$

$$\lambda = (W, \vec{b}, \vec{c}) = \theta$$

Training &

First: Train

Assume y

Training & Sampling a RBM

First: Training

Assume you have a data set: $S = \{ \vec{x}_1, \vec{x}_2, \dots, \vec{x}_e \}$
which is drawn from the underlying "physical" distribution $q(\vec{x})$
'Training' is adjusting λ so that

$$P_{\lambda}(\vec{x})$$

Training & Sampling a RBM

$$S = \{\vec{v}_1, \vec{v}_2, \dots, \vec{v}_l\}$$

First: Training

Assume you have a data set:

$$= \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_l\}$$

which is drawn from the underlying "physical" distribution $q(\vec{x})$

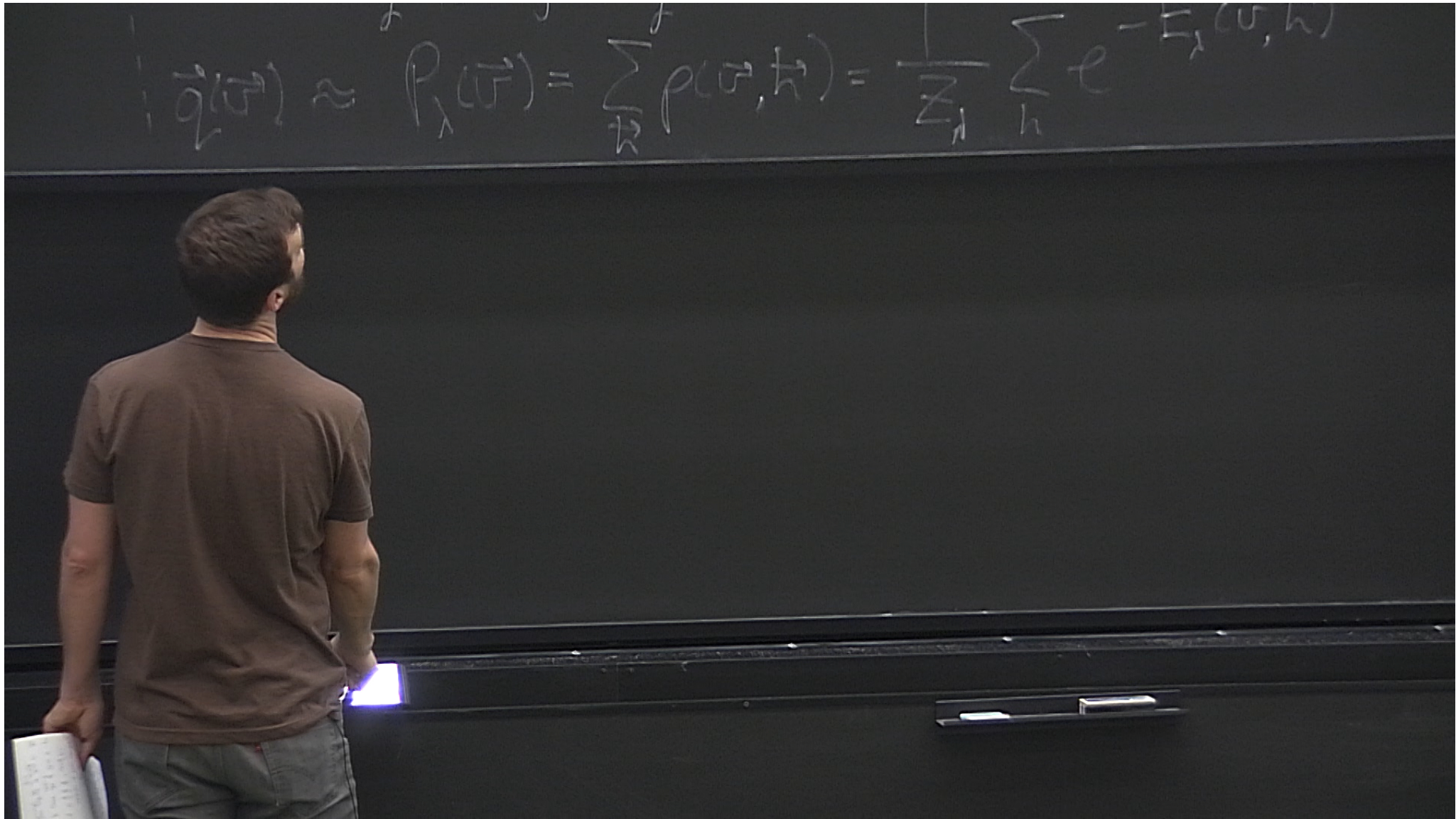
'Training' is adjusting λ so that

$$P_{\lambda}(\vec{v})$$

First: Training

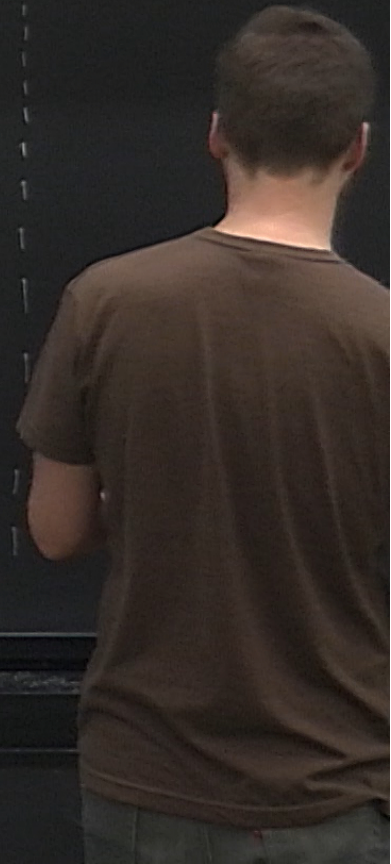
Assume you have a data set: $\mathcal{D} = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_l\}$
which is drawn from the underlying "physical" distribution $q(\vec{x})$
'learning' is adjusting λ so that

$$p_{\lambda}(\vec{w}) = \sum_{\vec{h}} p(\vec{w}, \vec{h}) = \frac{1}{Z_{\lambda}} \sum_{\vec{h}} e^{-E_{\lambda}(\vec{w}, \vec{h})}$$



Second: Sample / Reconstruct / generate

After training is complete, samples of \vec{v}
(and \vec{h}) can be drawn from the RBM.
this is done according to the joint $p(\vec{v}, \vec{h})$



$$\vec{q}(\vec{\omega}) \approx P_{\lambda}(\vec{\omega}) = \int_{\mathbb{R}} p(\vec{\omega}, h) = \frac{1}{Z_{\lambda}} \int_{\mathbb{R}} e^{-\dots}$$

What use are these generated samples?

Estimators can be calculated, e.g. for observable Q
 $\langle Q \rangle$

Second: Sample / Reconstruct / generate

After training is complete, samples of \vec{v}
(and \vec{h}) can be drawn from the RBM.
this is done according to the joint $p(\vec{v}, \vec{h})$

$$q(\vec{v}) \approx$$

What use

Estimate

L

$$\vec{q}(\vec{v}) \approx P_{\lambda}(\vec{v}) = \sum_{\vec{k}} p(\vec{v}, \vec{k}) = \frac{1}{Z_{\lambda}} \sum_{\vec{k}} e^{-\beta E(\vec{v}, \vec{k})}$$

What use are these generated samples?

Estimators can be calculated, e.g. for observable Q

$$\langle Q \rangle_{P_{\lambda}(\vec{v}, \vec{k})} = \frac{1}{Z_{\lambda}} \sum_{\vec{v}} \sum_{\vec{k}} Q p(\vec{v}, \vec{k})$$

$$\vec{q}(\vec{\sigma}) \approx P_{\lambda}(\vec{\sigma}) = \sum_{\vec{h}} p(\vec{\sigma}, \vec{h}) = \frac{1}{Z_{\lambda}} \sum_{\vec{h}} e^{-\beta \mathcal{H}(\vec{\sigma}, \vec{h})}$$

What use are these generated samples?

Estimators can be calculated, e.g. for observable Q

$$\langle Q \rangle_{P_{\lambda}(\vec{\sigma}, \vec{h})} = \frac{1}{Z_{\lambda}} \sum_{\vec{\sigma}} \sum_{\vec{h}} Q p_{\lambda}(\vec{\sigma}, \vec{h})$$

e.g. imagine defining a magnetization, etc ($\sigma_i = 0, 1$, $h_i = 0, 1$)

$$\vec{q}(\vec{\sigma}) \approx P_{\lambda}(\vec{\sigma}) = \sum_{\vec{h}} p(\vec{\sigma}, \vec{h}) = \frac{1}{Z_{\lambda}} \sum_{\vec{h}} e^{-\beta \mathcal{H}(\vec{\sigma}, \vec{h})}$$

What use are these generated samples?

Estimators can be calculated, e.g. for observable Q

$$\langle Q \rangle_{P_{\lambda}(\vec{\sigma}, \vec{h})} = \frac{1}{Z_{\lambda}} \sum_{\vec{\sigma}} \sum_{\vec{h}} Q p(\vec{\sigma}, \vec{h})$$

e.g. imagine defining a magnetization, etc ($\sigma_i = 0, 1$, $h_i = 0, 1$)
OR restrict the definition of Q to act on visible:

$$\text{ie } Q = Q_U,$$

$$\langle Q_U \rangle_{P_1(\vec{U}, \vec{k})} = \frac{1}{Z_1} \sum_{\vec{U}} Q_U \underbrace{\sum_{\vec{k}} P_1(\vec{U}, \vec{k})}_{P_1(\vec{U})}$$

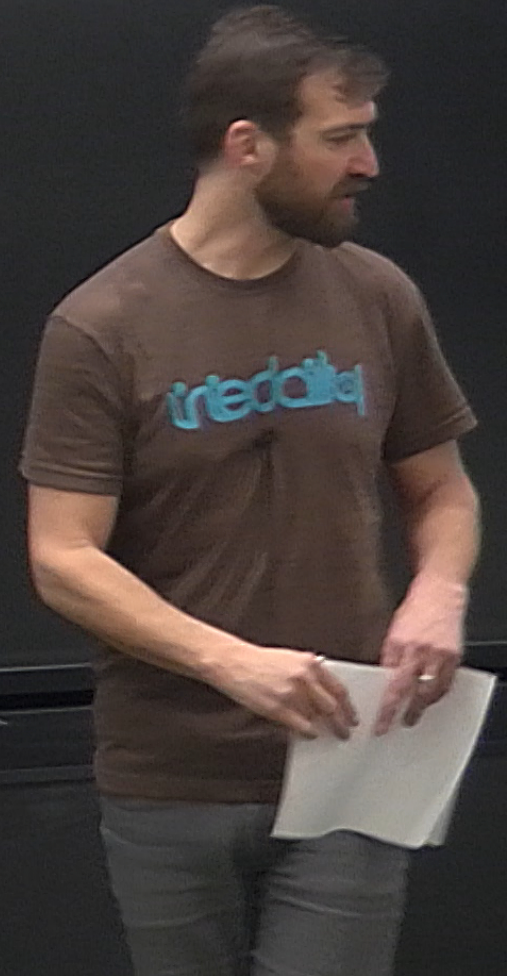
$$\text{ie } Q = Q_U,$$

$$\langle Q_U \rangle_{P_H(\vec{U}, \vec{h})} = \frac{1}{Z_H} \sum_{\vec{U}} Q_U \underbrace{\sum_{\vec{h}} P_H(\vec{U}, \vec{h})}_{P_H(\vec{U}) \approx q(\vec{U})}$$

ie $Q = Q_U$ \leftarrow Q a Hamiltonian over the variables
 a magnetization, c_u, χ , etc

$$\langle Q_U \rangle_{P_U(\vec{U}, \vec{h})} = \frac{1}{Z_U} \sum_{\vec{U}} Q_U \underbrace{\sum_{\vec{h}} P_U(\vec{U}, \vec{h})}_{P(\vec{U}) \approx q(\vec{U})}$$

since $Z_U = \sum_{\vec{U}} \sum_{\vec{h}} p(\vec{U}, \vec{h})$
 $= \sum_{\vec{U}} P_U(\vec{U})$



ie $Q = Q_U$ ← Q a Hamiltonian over the variables
 a magnetization, C_U , χ , etc

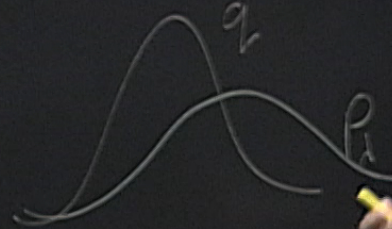
$$\langle Q_U \rangle_{P_U(\vec{U}, \beta)} = \frac{1}{Z_U} \sum_{\vec{U}} Q_U \underbrace{\sum_{\vec{h}} P_U(\vec{U}, \vec{h})}_{P_U(\vec{U}) \approx q(\vec{U})} \approx \langle Q_U \rangle_q$$

since $Z_U = \sum_{\vec{U}} \sum_{\vec{h}} p(\vec{U}, \vec{h})$
 $= \sum_{\vec{U}} P_U(\vec{U})$

Training a RBM

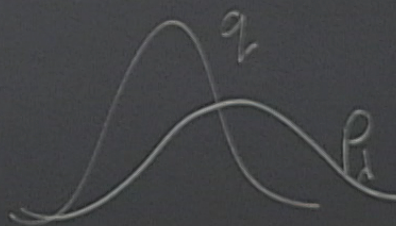
We want to optimize d to minimize the "distance" between

$p_1(\vec{v})$ and $q_2(\vec{v})$



Training a RBM

We want to optimize λ to minimize the "distance" between $p_{\lambda}(\vec{v})$ and $q(\vec{v})$



We do this by minimizing

the Kullback-Leibler (KL) divergence: $KL = \sum_{\vec{x}} q(\vec{x}) \log \frac{q(\vec{x})}{p(\vec{x})}$
for any q, p

$$KL = \sum_{\vec{x}} q(\vec{x}) \log q(\vec{x}) - \sum_{\vec{x}} q(\vec{x}) \log p(\vec{x})$$

for

\vec{x}) for us $\vec{x} = \vec{v}$ drawn from S
 q unknown (underlies \vec{v}), $p = p_{\lambda}(\vec{v})$



$$KL = \sum_{\vec{x}} q(\vec{x}) \log q(\vec{x}) - \sum_{\vec{x}} q(\vec{x}) \log p(\vec{x})$$

- non-symmetric measure of the "distance" between q and p
- always positive; zero if $p=q$
- first term is the entropy of the data set
- second term contains parameters θ

\vec{x}) for us $\vec{x} = \vec{v}$ drawn from S
 q unknown (underlies \vec{v}), $p = p_{\lambda}(\vec{v})$

and p Approximate $\langle \log p \rangle_q$
by the training samples $S = \{\vec{v}_1, \vec{v}_2, \dots, \vec{v}_n\}$

$$\sum_{\vec{v}} \log p(\vec{v}) = \log \mathcal{L} \quad \text{"log likelihood"}$$

Training is done by stochastic gradient descent
(pick a mini-batch out of S)

The gradient of the (negative) log likelihood
is used:

$$\lambda' = \lambda - \eta \nabla KL$$

↑ learning rate

Consider a single training example \vec{v}

$$\log \mathcal{L} = \log p(\vec{v})$$

Consider a single training example \vec{v}

$$\log \mathcal{L} = \log p_{\lambda}(\vec{v}) = \log p(\vec{v})$$

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

$$= \log \sum_{\vec{h}} e^{-E(\vec{v}, \vec{h})} - \log \sum_{\vec{v}, \vec{h}} e^{-E(\vec{v}, \vec{h})}$$

What use are these generated samples?

Estimators can be calculated, e.g. for observable Q

$\lambda = \lambda - \eta \nabla \mathcal{L}$

↑ learning rate

the gradient

$$\begin{aligned} \frac{\partial \log \mathcal{L}}{\partial \lambda} &= \frac{\partial}{\partial \lambda} \left(\log \sum_h e^{-E} \right) - \frac{\partial}{\partial \lambda} \left(\log \sum_{v,h} e^{-E} \right) \\ &= \frac{-1}{\sum e^{-E}} \sum_h e^{-E} \frac{\partial E}{\partial \lambda} + \frac{1}{\sum_{v,h} e^{-E}} \sum_{v,h} e^{-E} \frac{\partial E}{\partial \lambda} \end{aligned}$$

note: $p(h|\nu) = \frac{p(\nu, h)}{p(\nu)}$

$$= \frac{\frac{1}{2} e^{-E}}{\frac{1}{2} \sum_n e^{-E}} = \frac{e^{-E}}{\sum_n e^{-E}}$$

note: $p(h|\nu) = \frac{p(\nu, h)}{p(\nu)}$

$$= \frac{\frac{1}{2} e^{-E}}{\frac{1}{2} \sum_n e^{-E}} = \frac{e^{-E}}{\sum_n e^{-E}}$$

$$\frac{d \log p}{d\lambda} = - \sum_n p(h|\nu) \frac{dE}{d\lambda} + \sum_{\nu, h} p(\nu, h) \frac{dE}{d\lambda} \quad \text{⊗}$$

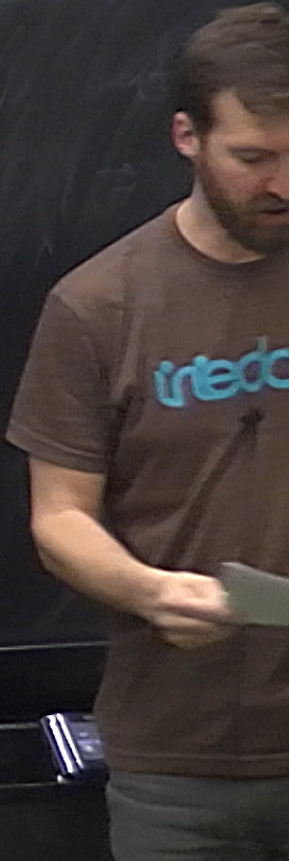
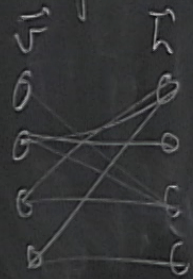
$$\frac{2 \log \mathcal{L}}{2d} = - \left\langle \frac{2E}{2d} \right\rangle$$

conditional
dist of h
given u

$$\frac{2 \log \mathcal{L}}{2d} = - \left\langle \frac{2E}{2d} \right\rangle_{\text{conditional dist of } h \text{ given } \mathcal{U}} + \left\langle \frac{2E}{2d} \right\rangle_{\text{full joint distribution}}$$

$$\frac{\partial \log \mathcal{L}}{\partial \lambda} = - \left\langle \frac{\partial E}{\partial \lambda} \right\rangle_{\text{conditional dist of } h \text{ given } \vec{v}} + \left\langle \frac{\partial E}{\partial \lambda} \right\rangle_{\text{full joint distribution}}$$

To calculate $p(\vec{h} | \vec{v})$ recall

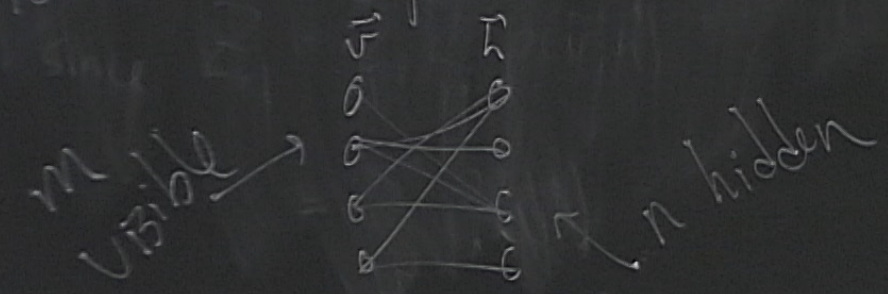


The hidden units are independent of one another
(similar for the visibles), and the conditional probs. factor

$$\text{ie. } p(\vec{h} | \vec{v}) = \prod_{i=1}^n p(h_i | \vec{v})$$

$$\frac{\partial \log \mathcal{L}}{\partial \mathbf{d}} = - \left\langle \frac{\partial E}{\partial \mathbf{d}} \right\rangle_{\text{conditional dist of } h \text{ given } \mathbf{v}} + \left\langle \frac{\partial E}{\partial \mathbf{d}} \right\rangle_{\text{full joint distribution}}$$

To calculate $p(\mathbf{h} | \mathbf{v})$ recall



The hidden unit
 (similar for
 ie. $p(\mathbf{h} | \mathbf{v})$
 $p(\mathbf{v} | \mathbf{h})$

The hidden units are independent of one another
(similar for the visibles), and the conditional probs. factor

$$\text{ie. } p(\vec{h} | \vec{v}) = \prod_{i=1}^n p(h_i | \vec{v})$$

$$p(\vec{v} | \vec{h}) = \prod_{i=1}^m p(v_i | \vec{h})$$

Next time we'll calculate

$$P(h_i=1 | \vec{v}) = \sigma \left(\sum_{j=1}^m W_{ij} v_j + c_i \right)$$

$$P(v_j=1 | \vec{h}) = \sigma \left(\sum_{i=1}^n W_{ij} h_i + b_j \right)$$

where

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

let's calculate the first term $\left\langle \frac{2E}{2h} \right\rangle_{p(h/v)}$

2d

2d conditional dist of h given v

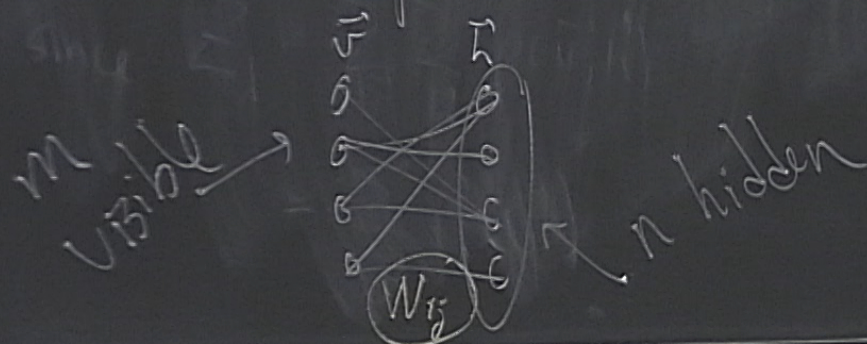
2d full joint distribution

(similar for

ie. $p(\vec{h} | \vec{v})$

$p(\vec{v} | \vec{h})$

To calculate $p(\vec{h} | \vec{v})$ recall



Next time we'll calculate

$$p(h_i=1 | \vec{v}) = \sigma \left(\sum_{j=1}^m W_{ij} v_j + c_i \right)$$

$$p(v_i=1 | \vec{h}) = \sigma \left(\sum_{j=1}^n W_{ij} h_j + b_i \right)$$

let's calculate the fir

let's do the example o

lets calculate the first term = $\left\langle \frac{2E}{2\hbar} \right\rangle_{p(\hbar|\vec{v})}$

lets do the example of $d = W_{ij}$

$$\sum_{\vec{h}} p(\vec{h}|\vec{v}) \frac{2E}{2W_{ij}} = \sum_{\vec{h}} p(\vec{h}|\vec{v}) h_i v_j = \sum_{\vec{h}} \prod_{k=1}^n p(h_k|\vec{v}) h_i v_j$$

lets calculate the first term $\left\langle \frac{\partial E}{\partial h} \right\rangle_{p(h|\vec{v})}$

Let's do the example of $d = W_{ij}$

$$\sum_{\vec{h}} p(\vec{h}|\vec{v}) \frac{\partial E}{\partial W_{ij}} = \sum_{\vec{h}} p(\vec{h}|\vec{v}) h_i v_j = \sum_{\vec{h}} \prod_{k=1}^n p(h_k|\vec{v}) (h_i v_j)$$

introduce \vec{h}_{-i} is all the elements of \vec{h} minus h_i removed

$$= \sum_{h_i} \sum_{\vec{h}_{-i}} p(h_i | \vec{\sigma}) p(\vec{h}_{-i} | \vec{\sigma}) h_i \sigma_j$$

$$= \sum_{h_i} \sum_{t_{-i}} p(h_i | \vec{\sigma}) p(t_{-i} | \vec{\sigma}) h_i \sigma_j = \sum_{h_i} p(h_i | \vec{\sigma}) h_i \sigma_j$$

TRACE: $\sum_{h_i=0,1} \sum_{t_{-i}=\{0,1\}}$

$$= \sum_{h_i} p(h_i | \vec{v}) h_i v_i \underbrace{\sum_{h_i} p(h_i | \vec{v})}_{[\rho(h_i=0 | \vec{v}) + \rho(h_i=1 | \vec{v})] \dots}$$

$$= \sum_{h_i} p(h_i | \vec{v}) h_i v_i \underbrace{\sum_{h_i} p(h_i | \vec{v})}_{\left[p(h_i=0 | \vec{v}) + p(h_i=1 | \vec{v}) \right] \dots = 1}$$

$$\sum_{h_i} p(h_i | \vec{v}) h_i v_i \quad \underbrace{\sum_{h_i} p(h_i | \vec{v})}_{[p(h_i=0 | \vec{v}) + p(h_i=1 | \vec{v})] \dots = 1}$$

$$\sum_{h_i} p(h_i | \vec{v}) h_i v_j = (p(h_i=1 | \vec{v}) \cdot 1 + p(h_i=0 | \vec{v}) \cdot 0) v_j = p(h_i=1 | \vec{v}) v_j$$

$$\text{So } \sum_{\vec{h}} p(\vec{h}|\vec{v}) \frac{2E}{2W_{ij}} = \sigma \left(\sum_{j=1}^M W_{ij} v_j + c_i \right) v_j$$

The other term: $\left\langle \frac{2E}{2\lambda} \right\rangle_{p(\vec{v}, \vec{h})} = \sum_{\vec{v}, \vec{h}} p(\vec{v}, \vec{h}) \frac{2E}{2\lambda} = \sum_{\vec{v}, \vec{h}}$

$$\begin{aligned}
 & \sum_{\vec{h}} p(\vec{h}) \sum_{\vec{v}} p(\vec{v} | \vec{h}) \frac{2\pi}{2\lambda} = \sum_{\vec{h}} p(\vec{h}) \sum_{\vec{v}} p(\vec{v} | \vec{h}) \frac{2\pi}{2\lambda} \\
 & \sum_{\vec{h}} p(\vec{h}) \sum_{\vec{v}} p(\vec{v} | \vec{h}) \frac{2\pi}{2\lambda} = \sum_{\vec{h}} p(\vec{h}) \sum_{\vec{v}} p(\vec{v} | \vec{h}) \frac{2\pi}{2\lambda}
 \end{aligned}$$

$= (P(\vec{h}_1 | \vec{h}_0) + \dots + P(\vec{h}_n | \vec{h}_0))$

the inner sum can be factorized

$\vec{h}_i | \vec{h}_j$

$$\sum_{\vec{h}} \frac{2E}{2^n} = \sum_{\vec{h}} p(\vec{h}) \sum_{\vec{h}'} p(\vec{h}' | \vec{h}) \frac{2E}{2^n} = \sum_{\vec{h}} p(\vec{h}) \sum_{\vec{h}'} p(\vec{h}' | \vec{h}) \frac{2E}{2^n}$$

But you are left with $\sum_{\vec{h}'} p(\vec{h}' | \vec{h})$ or $\sum_{\vec{h}'} 1$ (ie 2^m or 2^n)