

Here we have a graphical prob. dist.

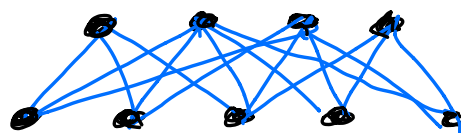
$P_{\lambda}(\vec{v}, \vec{h})$ ie. a joint distribution

"learning" is adjusting parameters λ so that

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

Restricted Boltzmann Machine (Hinton, Smolensky '86)

n of \vec{h}



\vec{v} m of these

no intra-layer couplings

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

$$v_i, h_j = 0, 1$$

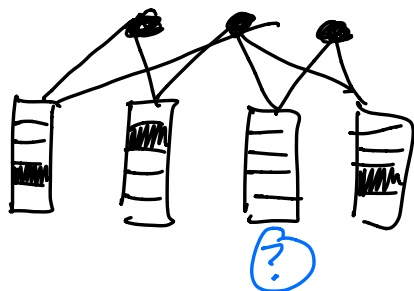
parameters are $\lambda = (W, \vec{b}, \vec{c})$

and

$$P_{\lambda}(\vec{v}, \vec{h}) = \frac{1}{Z_{\lambda}} e^{-E_{\lambda}(\vec{v}, \vec{h})}$$

Industry applications

- Netflix competition: collaborative filtering



- topic modelling: interpreting latent units as "semantic structure"
- Unsupervised "pre-training" for deep learning (pre-2012)

Now: Training an RBM: $\mathcal{D} = \{\vec{v}\}$
drawn from an unknown $P(\vec{v})$
adjusting λ so that

$$P_{\lambda}(\vec{v}) = \sum_k p(\vec{v}, k) \approx P(\vec{v})$$

Next: Sampling: after training, generate new \vec{v}, k drawn from $p(\vec{v}, k)$

and calculate $\langle \mathcal{O} \rangle_{p_\lambda(\vec{v}, \vec{h})} = \frac{1}{Z} \sum_{\vec{v}} \sum_{\vec{h}} \mathcal{O} p_\lambda(\vec{v}, \vec{h})$

note: if we define $\mathcal{O} = \mathcal{O}_v$ (e.g. E_v)

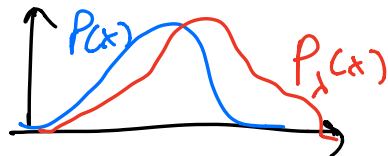
$$\langle \mathcal{O}_v \rangle_{p_\lambda(\vec{v}, \vec{h})} = \frac{1}{Z} \sum_v \mathcal{O}_v \underbrace{\sum_{\vec{h}} p_\lambda(\vec{v}, \vec{h})}_{p_\lambda(\vec{v}) \approx P(\vec{v})}$$

note

$$Z = \sum_v \sum_h p_\lambda(\vec{v}, \vec{h}) = \sum_{\vec{v}} p_\lambda(\vec{v})$$

TRAINING: Define a cost function

$$KL(P \parallel p_\lambda) = \sum_{\vec{x}} P(\vec{x}) \log \frac{P(\vec{x})}{p_\lambda(\vec{x})}$$



$$KL(P \parallel p_\lambda) = \underbrace{\sum_x P(\vec{x}) \log P(\vec{x})}_{\text{entropy of } P \equiv -H_P} - \underbrace{\sum_x P(\vec{x}) \log p_\lambda(\vec{x})}_{\text{depends on } \lambda}$$

recall $\mathcal{D} = \{\vec{x}\}$, $P_{\text{data}}(\vec{x}) = \frac{1}{|\mathcal{D}|} \sum_{\vec{x}_k \in \mathcal{D}} \delta_{\vec{x}, \vec{x}_k}$

$$KL(P \| P_\lambda) \simeq -H_D - \underbrace{\frac{1}{\|p\|} \sum_{\vec{x} \in \mathcal{B}} \log P_\lambda(\vec{x})}_{\text{use this for the cost fn } C_\lambda}$$

Now: use gradient descent

$$\lambda' = \lambda - \eta \nabla_\lambda C_\lambda$$

Let's consider a single training example \vec{v}

$$\begin{aligned} \log p_\lambda(\vec{v}) &= \log \frac{1}{Z} \sum_h e^{-E(v, h)} \\ &= \log \sum_h e^{-E} - \log \sum_{v, h} e^{-E} \end{aligned}$$

The gradient:

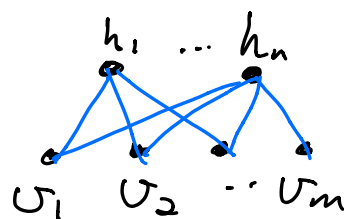
$$\begin{aligned} \frac{\partial \log p_\lambda}{\partial \lambda} &= \frac{\partial}{\partial \lambda} \left(\log \sum_h e^{-E} \right) - \frac{\partial}{\partial \lambda} \left(\log \sum_{v, h} e^{-E} \right) \\ &= -\frac{1}{\sum_h e^{-E}} \sum_h e^{-E} \cdot \frac{\partial E}{\partial \lambda} + \frac{1}{\sum_{v,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

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

$$\frac{2 \log p_{\lambda}}{2\lambda} = - \sum_h p(h|\vec{v}) \frac{2E}{2\lambda} + \sum_{v|h} p(v|h) \frac{2E}{2\lambda}$$

and recall $\frac{2E}{2W_{ij}} = -v_i h_j$ etc

Recall the RBM graph



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

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

This gives "Block" Gibbs sampling.

Simple to calculate an analytical expression

$$p(h_j=1|\vec{v}) = \sigma\left(\sum_{i=1}^m W_{ij} v_i + c_j\right)$$

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

↑ sigmoid function

Back to $\textcircled{*}$, let's calculate first term for $\lambda = W_{ij}$

$$\begin{aligned}
 \sum_{\vec{h}} p(\vec{h} | \vec{v}) \frac{2E}{2W_{ij}} &= \sum_{\vec{h}} \underbrace{p(\vec{h} | \vec{v})}_{p(h_i | \vec{v}) p(\vec{h}_{-i} | \vec{v})} h_i v_j \\
 &= \sum_{h_i} p(h_i | \vec{v}) h_i v_j \sum_{\vec{h}_{-i}} p(\vec{h}_{-i} | \vec{v}) \\
 &\quad \xrightarrow{\text{trace}} \sum_{\{h_i = 0, 1\}} \\
 &\quad \text{ie } (p(h_i = 0 | \vec{v}) + p(h_i = 1 | \vec{v})) \cdot () \\
 &= \sum_{h_i} p(h_i | \vec{v}) h_i v_j
 \end{aligned}$$

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

$$= \sigma \left(\sum_{k=1}^m W_{ik} v_k + c_k \right) v_j$$

Next, look at second term in $\textcircled{*}$

$$\sum_{\vec{v}, \vec{h}} p(\vec{v}, \vec{h}) \frac{2E}{2\lambda} = \sum_{\vec{v}} p(\vec{v}) \sum_{\vec{h}} p(\vec{h} | \vec{v}) \frac{2E}{2\lambda}$$

Inner sum is tractable, but $\sum_{\vec{v}} \dots$ is a trace over 2^m states, and is not.

How is this sum handled? With something called "contrastive divergence"

just fudge the estimator with a very short Markov Chain:

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

get a rough estimate with a small ($k=1 \dots 5$) number of steps in the Markov chain.

For completeness, remember we want to use the entire training set $\{\vec{v}\} = \mathcal{D}$

$$C_\lambda = -\frac{1}{\|\mathcal{D}\|} \sum_{\vec{v} \in \mathcal{D}} \log p_\lambda(\vec{v})$$

$$\nabla_\lambda C_\lambda = -\frac{1}{\|\mathcal{D}\|} \sum_{\vec{v} \in \mathcal{D}} \frac{2}{2\lambda} \log p_\lambda(\vec{v})$$

$$\frac{2 \log p_\lambda(\vec{v})}{2\lambda} = -\sum_{\vec{h}} p(\vec{h}|\vec{v}) \frac{2E(\vec{v}|\vec{h})}{2\lambda} + \sum_{\vec{h}} p(\vec{h}|\vec{v}^*) \frac{2E(\vec{v}^*|\vec{h})}{2\lambda}$$