

Detailed Balance

When transitioning $\mu \rightarrow \nu$ to ensure we're sampling from the correct $P(\vec{x})$

$$P(\mu)T(\mu \rightarrow \nu) = P(\nu)T(\nu \rightarrow \mu)$$

$$\text{or } \frac{T(\mu \rightarrow \nu)}{T(\nu \rightarrow \mu)} = \frac{P(\nu)}{P(\mu)}$$

e.g.) for Boltzmann dist:

$$\frac{T(\mu \rightarrow \nu)}{T(\nu \rightarrow \mu)} = \frac{\sum e^{-\beta E_\nu}}{\sum e^{-\beta E_\mu}} = e^{-\beta(E_\nu - E_\mu)}$$

Difficult part: designing the update itself

Design trick: split the transition probability

$$T(\mu \rightarrow \nu) = g(\mu \rightarrow \nu) A(\mu \rightarrow \nu)$$

selection probability \uparrow acceptance ratio

e.g.) Ising model single-spin flip

propose an update $x_i \rightarrow -x_i$ (if $\neq 1$)

then $q(\mu \rightarrow \nu) = \frac{1}{N}$ for any μ, ν

$$\frac{T(\mu \rightarrow \nu)}{T(\nu \rightarrow \mu)} = \frac{\sum_{\nu} \frac{1}{N} A(\mu \rightarrow \nu)}{\sum_{\mu} \frac{1}{N} A(\nu \rightarrow \mu)} = \frac{P(\nu)}{P(\mu)}$$

ie.
$$\frac{A(\mu \rightarrow \nu)}{A(\nu \rightarrow \mu)} = e^{-\beta(E_\nu - E_\mu)}$$

the "Metropolis" algorithm chooses

$$A(\nu \rightarrow \mu) = 1 \quad \text{if } E_\mu < E_\nu$$

then

$$A(\mu \rightarrow \nu) = e^{-\beta(E_\nu - E_\mu)}$$

when the energy is raised.

Generalized version

$$A(\mu \rightarrow \nu) = \min \left\{ 1, \frac{P(\nu)}{P(\mu)} \cdot \frac{q(\nu \rightarrow \mu)}{q(\mu \rightarrow \nu)} \right\}$$

let's derive another algorithm used in RBMs called "Gibbs" sampling.

Some definitions:

$P(A|B)$ "prob of A given B"

$$P(A) = \sum_B P(A, B) \quad \text{"sum rule"}$$

marginal prob.

$P(A, B)$ "joint distribution"

$$P(A, B) = P(B|A) P(A)$$

"product rule"

let's use these to suggest an update
that replaces x_i ($\vec{x}^\mu = (\underbrace{x_1}_B, \underbrace{x_2}_A, \dots, \underbrace{x_i}_B, \dots, x_N)$)
with a value drawn
from $P(x_i | \vec{x}_{-i}^\mu)$

$\nwarrow x_i$ is omitted

ie.

$$\begin{aligned} q(\mu \rightarrow \nu) &= P(x_i^\nu | \vec{x}_{-i}^\mu) \\ &= P(x_i^\nu | \vec{x}_{-i}^\nu) \end{aligned}$$

let's use

$$\begin{aligned} P(\nu) &= P(\vec{x}^\nu) \\ &= P(x_i^\nu | \vec{x}_{-i}^\nu) P(\vec{x}_{-i}^\nu) \end{aligned}$$

$$\begin{aligned}
\frac{A(\mu \rightarrow \nu)}{A(\nu \rightarrow \mu)} &= \frac{P(\nu)}{P(\mu)} \frac{g(\mu \rightarrow \nu)}{g(\nu \rightarrow \mu)} \\
&= \frac{P(\nu)}{P(\mu)} \frac{P(x_i^\mu | \vec{x}_{-i}^\mu)}{P(x_i^\nu | \vec{x}_{-i}^\nu)} \\
&= \frac{P(x_i^\nu | \vec{x}_{-i}^\nu) P(\vec{x}_{-i}^\nu) P(x_i^\mu | \vec{x}_{-i}^\mu)}{P(x_i^\mu | \vec{x}_{-i}^\mu) P(\vec{x}_{-i}^\mu) P(x_i^\nu | \vec{x}_{-i}^\nu)} \\
&= 1 \quad \text{since } \vec{x}_{-i}^\mu = \vec{x}_{-i}^\nu
\end{aligned}$$

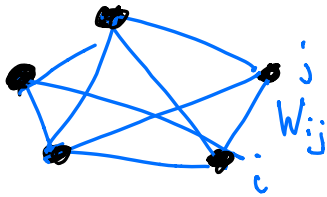
Therefore in Gibbs sampling the acceptance ratio is always 1.

Recall: our goal is to take
 $\mathcal{D} = \{\vec{x}\}$ and to find (an
approximate) $p(\vec{x}) \simeq P(\vec{x})$

This indicates a parametrization
 \rightarrow goal tune λ using \mathcal{D} .

The Hopfield Network (1982)

Prob. graphical model / Hopfield neural network is defined with



N nodes
for n variables

$$E(\vec{x}) = - \sum_{ij} W_{ij} x_i x_j - \sum_i b_i x_i$$

using

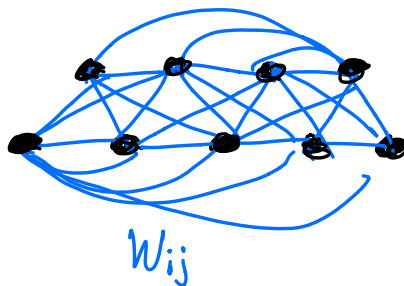
$$P_{\lambda}(\vec{x}) = \frac{1}{Z} e^{-E(\vec{x})}$$

"learn" $\lambda = (W, b)$ so that $p_{\lambda}(\vec{x}) \approx P(\vec{x})$

Boltzmann Machine (Ackley, Hinton, Sejnowski 1985)

Similar to Hopfield with an additional latent space

or
"hidden units"



\vec{h} hidden
 \vec{v} visible

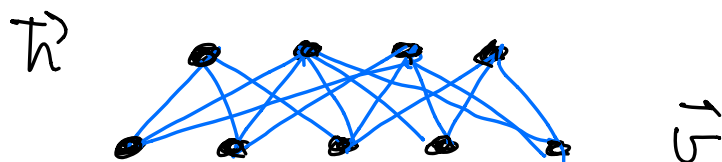
Here we have a graphical prob. dist.

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

"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)



no intra-layer couplings

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

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

and

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