

Boltzmann Machines

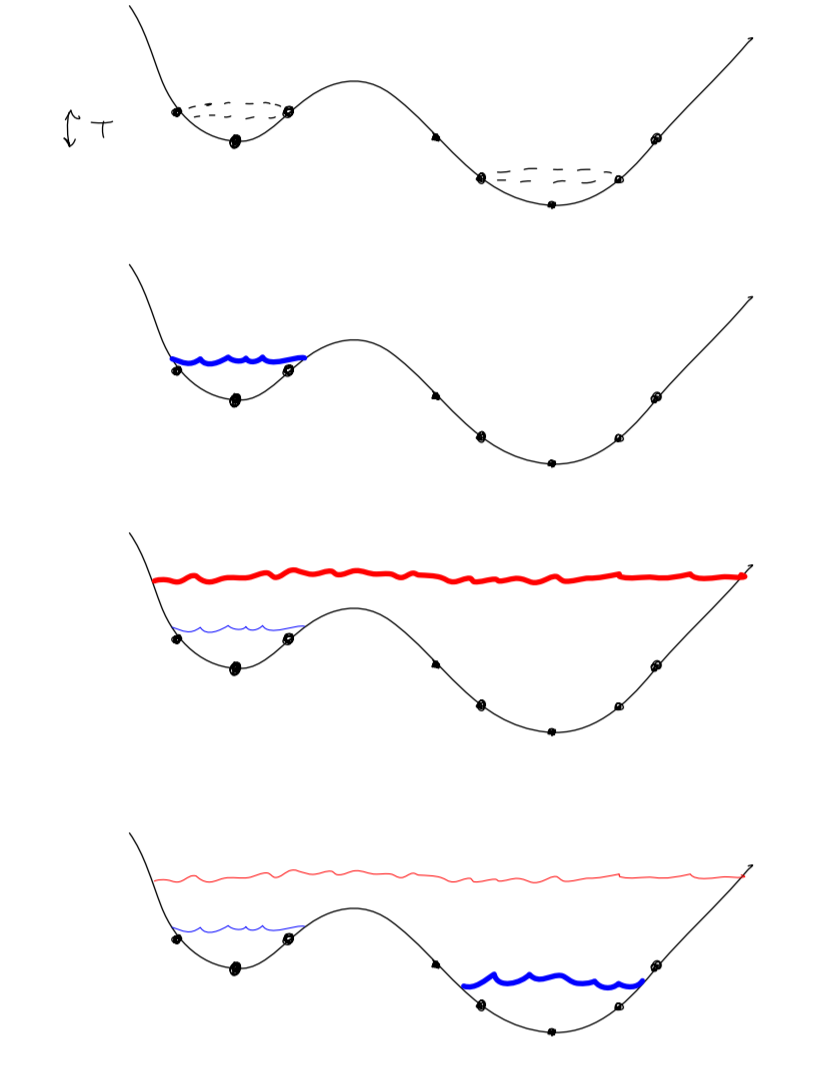
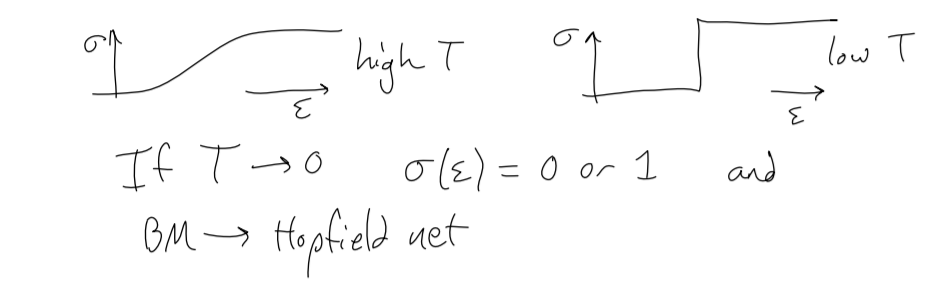
Stochastic generalization of Hopfield net.

Energy: $E = -\sum_{ij} J_{ij} x_i x_j - \sum_i h_i x_i$

Boltzmann distribution: $Q(x) = \frac{e^{-\frac{1}{T}E(x)}}{Z}$
 $Z = \sum_x e^{-\frac{1}{T}E(x)}$ $\beta = \frac{1}{T}$
 "partition function" inverse temperature

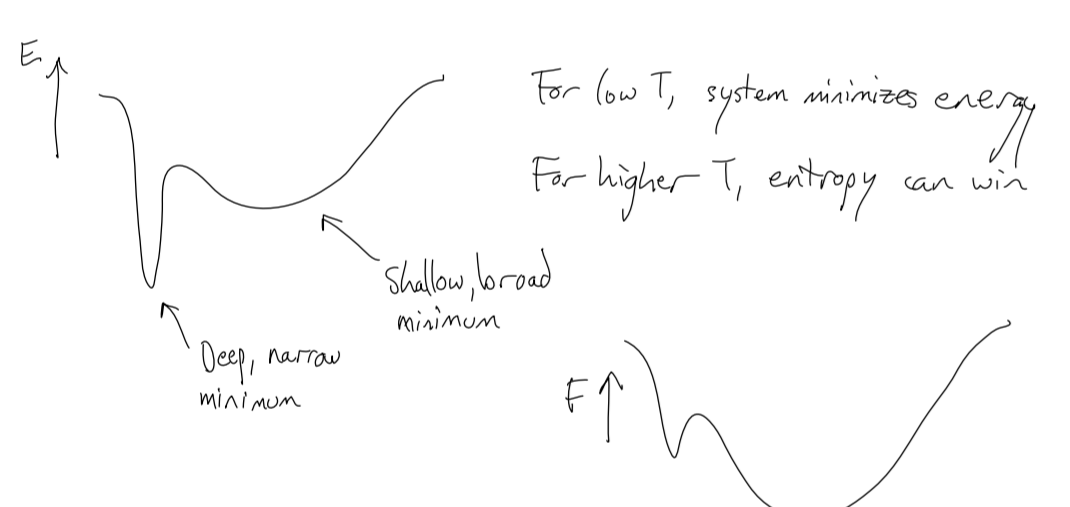
Hopfield update: $x_i^{t+1} = \text{argmin}_i E(x_i, x_i^t)$

Boltzmann update: $Q(x_i^{t+1}) = \sigma(\frac{1}{T}E(x_i^t, x_i^t))$
 "glazier dynamics" or "Gibbs sampling" $\sigma(\epsilon) = \frac{1}{1+e^{-\epsilon}}$



With $T > 0$, Energy is no longer minimized. But Free Energy is!

$F = U - TS$
 Hidden Free Energy Total energy $\langle E \rangle$ Entropy Temperature



Interpretation: a system has a higher probability of exploring many higher energy states than staying in a few lower energy states.
 $Q(x) = e^{-\beta E(x)}$ so if there are more than $e^{\beta \Delta E}$ higher E states, the system will usually live there.

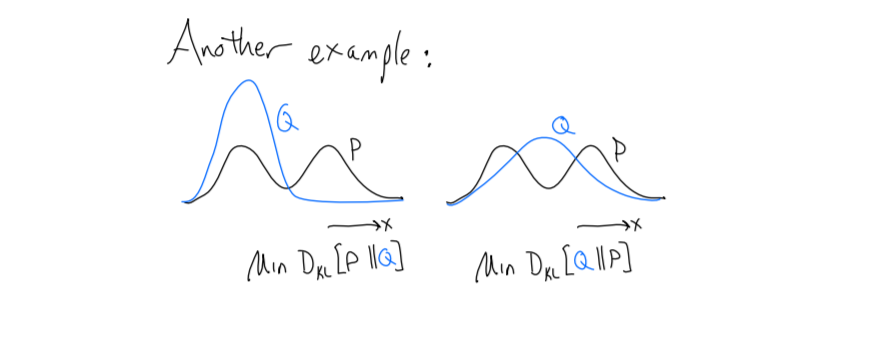
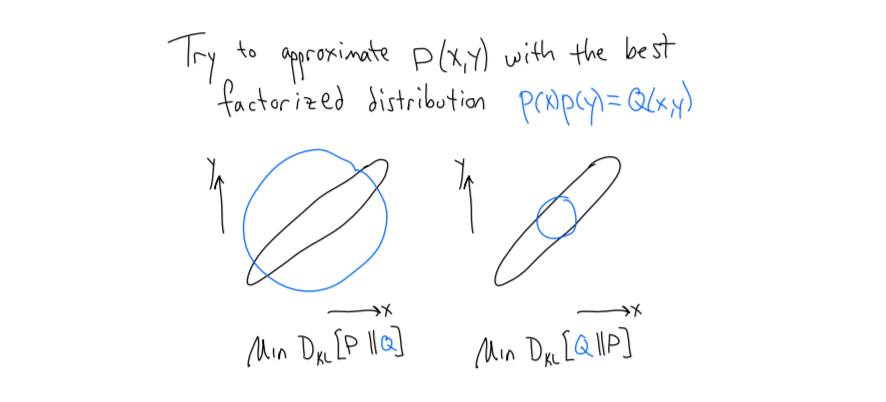
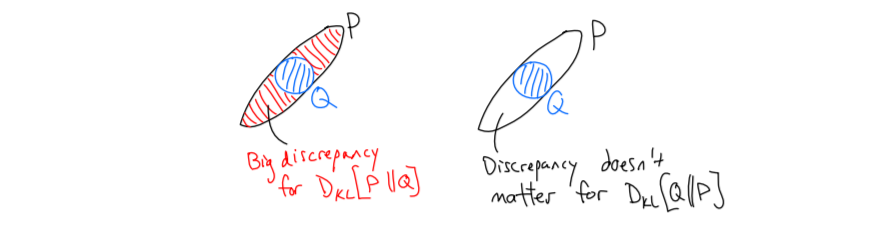
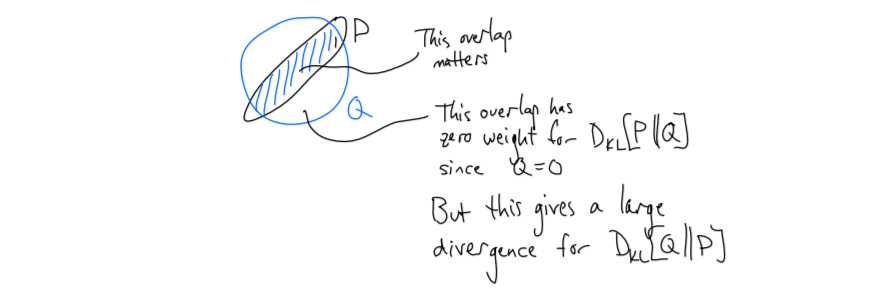
Learning: Try to make $Q(x)$ match a target $P(x)$.

Define a cost function: Kullback-Leibler divergence $D_{KL}[P||Q] = \sum_x P(x) \log \frac{P(x)}{Q(x)}$

$D_{KL} \geq 0$ with $=$ iff $P=Q$.

Example: $D_{KL}[P(x,y)||P(x)P(y)] = I(x,y)$
 $D_{KL}[S(x)||P(x)] = I(y,P)$

$D_{KL}[Q||P] \neq D_{KL}[P||Q] \Rightarrow$ not a distance
 $\sum_x P(x)P(x) - \sum_x P(x)Q(x)$



Minimize D_{KL} with respect to parameters of $E(x; \{J, h\})$

$0 = \frac{\partial D_{KL}[P||Q]}{\partial J_{ij}} = \frac{\partial}{\partial J_{ij}} \left[\sum_x P(x) \log \frac{P(x)}{Q(x)} \right]$

$= \frac{\partial}{\partial J_{ij}} \left[\sum_x P(x) (-\log P(x) + \log Q(x)) \right]$

$= \frac{\partial}{\partial J_{ij}} \left[\sum_x P(x) \left(\log \frac{1}{Z} + \sum_{ij} J_{ij} x_i x_j + \sum_i h_i x_i \right) \right]$

$= \frac{\partial}{\partial J_{ij}} \left[\sum_x P(x) \left(\log \frac{1}{Z} + \sum_{ij} J_{ij} x_i x_j + \sum_i h_i x_i \right) \right]$

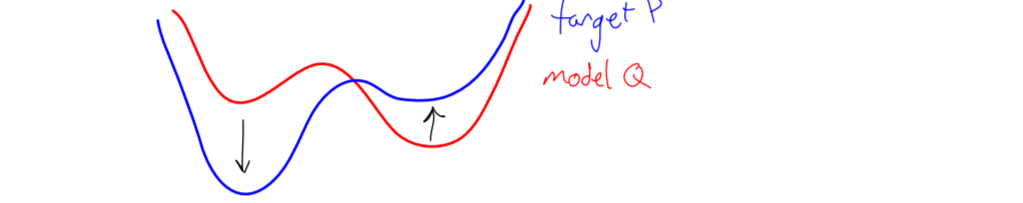
$\frac{\partial D_{KL}}{\partial J_{ij}} = \frac{1}{T} (\langle x_i x_j \rangle_Q - \langle x_i \rangle_Q \langle x_j \rangle_Q)$

$\frac{\partial D_{KL}}{\partial h_i} = \frac{1}{T} (\langle x_i \rangle_Q - \langle x_i \rangle_P)$

To minimize cost, we could thus set $\langle x_i x_j \rangle_Q = \langle x_i x_j \rangle_P$

It's hard to specify J_{ij} to produce a given $\langle x_i x_j \rangle_P$. Instead, we can descent the gradient of D_{KL} : $J_{ij}^{t+1} = -\alpha (\langle x_i x_j \rangle_Q - \langle x_i x_j \rangle_P)$

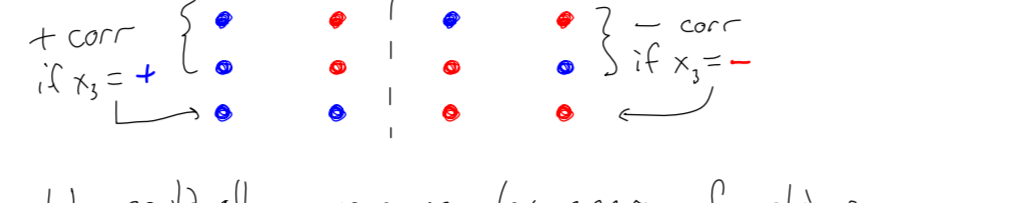
This lowers the energy of x where $Q < P$ and raises it where $Q > P$.



We still need to estimate $\langle x_i x_j \rangle_Q$, which is also hard... This is accomplished by sampling, or sometimes by approximation (e.g. mean field assuming neurons are independent)

Why is it hard to just calculate $\langle x_i x_j \rangle_Q$ from $Q(x)$?

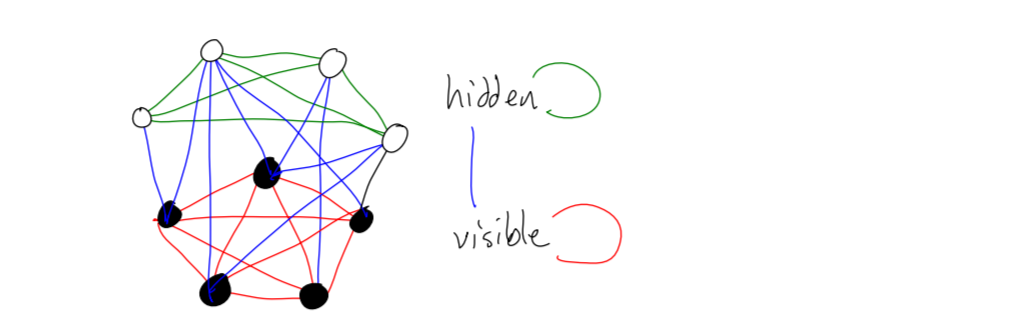
The second-order BM is limited in its expressiveness, and can't capture higher-order interactions:



We could allow more complex energy functions, e.g. $E(x) = -\sum_{ijk} K_{ijk} x_i x_j x_k$

to capture such ^{3rd} effects, but direct triple interactions are not generally neurally possible.

Instead, we can introduce HIDDEN NODES to capture more complex dependencies.



This architecture can accommodate any distribution over the 2^V possible states of the V visible units

How can we learn this?

Same rule as before, except that the hidden units are not specified by the target distribution:

$\Delta J_{ij} = -\alpha [\langle x_i x_j \rangle_{\text{clamped}} - \langle x_i x_j \rangle_{\text{free}}]$

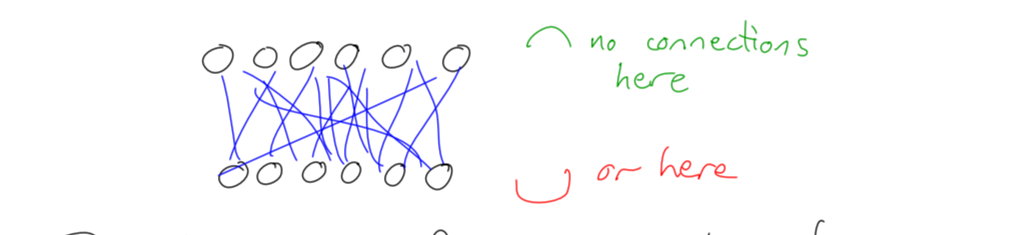
clamped: visible units are set to $x \sim P(x)$ and for each sample, hidden units are given by $Q(x, \text{visible} | \text{hidden})$.

free: all $x \sim Q(x)$

This is UNSUPERVISED learning: no teacher specifies the correct answers for the hidden units. The network finds hidden unit connections (latent variables) that best explain the distribution of the input data.

The architecture, however, must be specified, as fully general BMs (fully connected) are intractably slow.

\Rightarrow Restricted Boltzmann Machine



This makes sampling faster, at a loss of generality. Still need MANY trials to train a deep BM.

One interesting trick is to stop before equilibrium.

- Start at a point from the data, x^t , and only sample for a short time.
- Contrastive divergence
- Minimum probability flow

Another perspective: What is the distribution $Q(x)$ that matches some specified moments, e.g. $\langle x_i x_j \rangle = A_{ij}$ but is otherwise minimally uncertain?

Find maximum entropy subject to constraints. $H = -\sum Q \log Q$

$C = H - \sum_{ij} \beta_{ij} (\sum_x x_i x_j - A_{ij}) - \sum_i \beta_i (\sum_x x_i - 1)$

$\frac{\partial C}{\partial \beta_{ij}} = 0 = -(\log Q) - \sum_{ij} \beta_{ij} x_i x_j - \sum_i \beta_i x_i - 1$

$Q = e^{-(H + \sum_{ij} \beta_{ij} x_i x_j + \sum_i \beta_i x_i)}$

$= \frac{e^{-\sum_{ij} \beta_{ij} x_i x_j - \sum_i \beta_i x_i}}{Z} \geq e^{-H}$

Again the Boltzmann machine! (if $x \in \pm 1$; if $x \in \mathbb{R}$, then this is gaussian)

Example: <http://www.cs.toronto.edu/~hinton/digits.html>