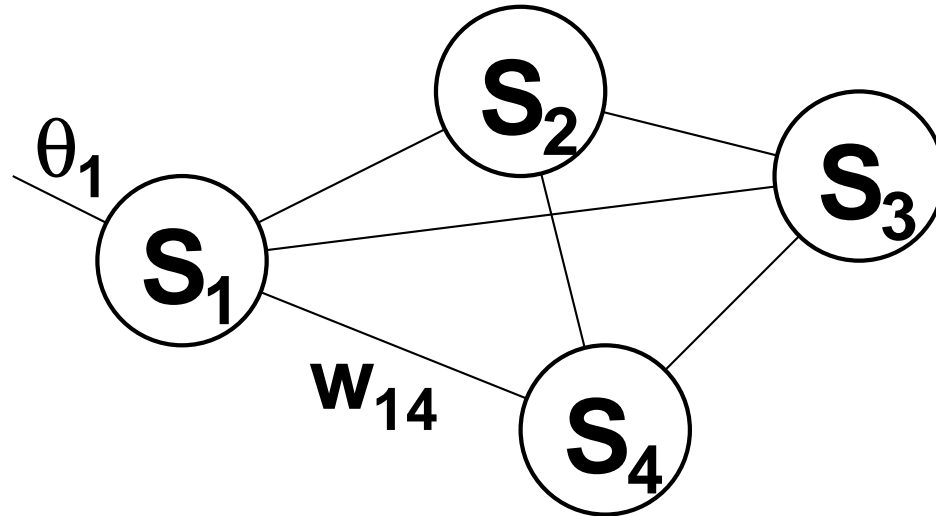


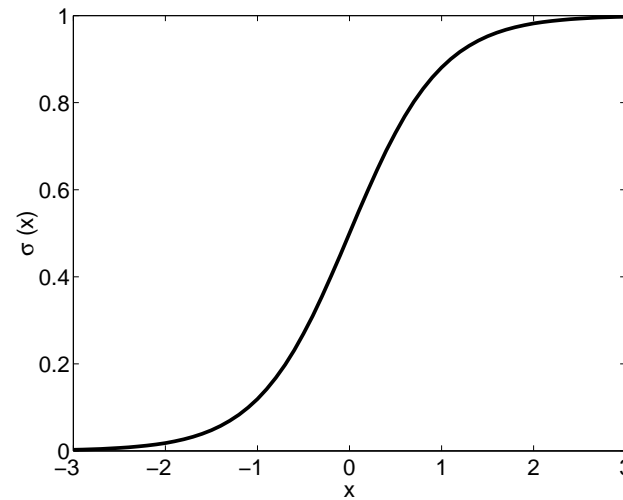
# Variational Methods and Bounds

- Introduction to Boltzmann machines
- Simple approximations
- Introduction to variational methods
- The search for bounds
- Summary

# The Boltzmann Machine



$$p(s_i = 1) = \sigma \left( \theta_i + \sum_j w_{ij} s_j \right)$$



## The Boltzmann Machine

The *energy* of the Boltzmann machine for a certain *state* is

$$-E(\vec{s}) = \sum_i \theta_i s_i + \frac{1}{2} \sum_{ij} w_{ij} s_i s_j$$

The *probability* to find the BM in state  $\vec{s}$ :

$$p(\vec{s}) = \frac{1}{Z} \exp(-E(\vec{s}))$$

The normalizing constant is the *partition function*

$$Z = \sum_{\text{all } \vec{s}} \exp(-E(\vec{s}))$$

## Derived Quantities

$$Z = \sum_{\text{all } \vec{s}} \exp(-E(\vec{s}))$$

*Means and correlations*

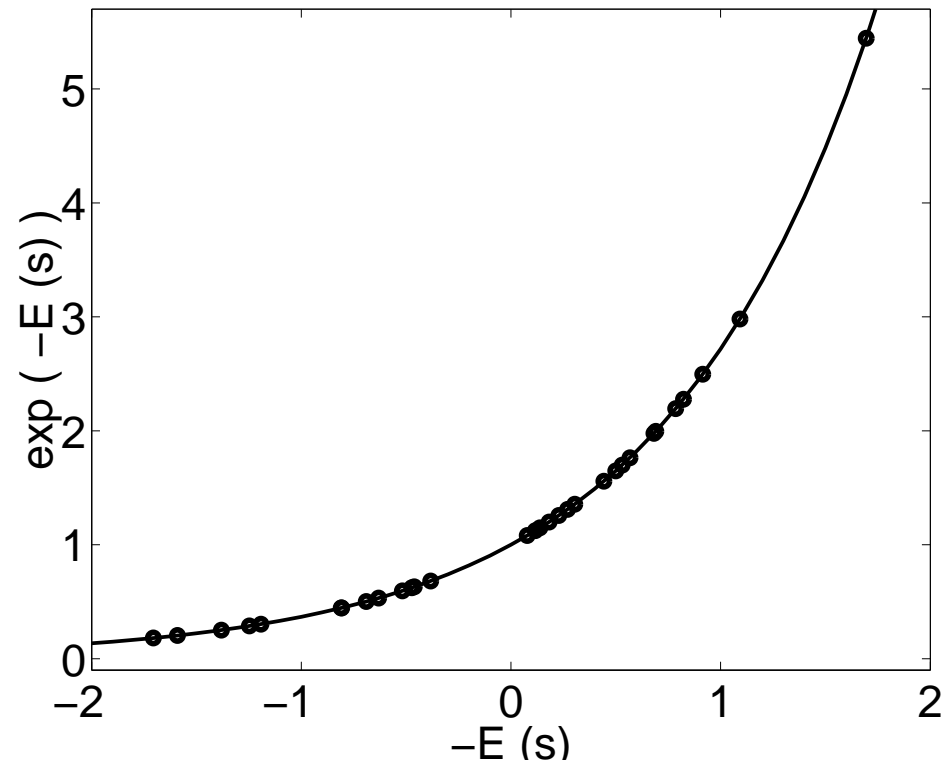
$$\begin{aligned} \langle s_i \rangle &= \frac{\partial}{\partial \theta_i} \log Z = \frac{1}{Z} \frac{\partial Z}{\partial \theta_i} = \frac{1}{Z} \sum_{\text{all } \vec{s}} \frac{\partial}{\partial \theta_i} \exp(-E(\vec{s})) \\ &= \frac{1}{Z} \sum_{\text{all } \vec{s}} \exp(-E(\vec{s})) s_i = \sum_{\text{all } \vec{s}} p(\vec{s}) s_i \end{aligned}$$

$$\langle s_i s_j \rangle = \frac{\partial}{\partial w_{ij}} \log Z$$

# Computing the Partition Function

$$Z = \sum_{\text{all } \vec{s}} \exp(-E(\vec{s}))$$

| State     | exp(Energy) |
|-----------|-------------|
| — — — — — | 3.979020    |
| — — — — + | 1.586260    |
| — — — + — | 3.480777    |
| — — — + + | 1.994746    |
| — — + — — | 0.566925    |
| — — + — + | 0.505540    |
| — — + + — | 1.466552    |
| — — + + + | 1.879924    |
| — + — — — | 1.604714    |
| — + — — + | 0.833811    |
| — + — + — | 2.247136    |
| — + — + + | 1.678465    |
| — + + — — | 0.795865    |
| — + + — + | 0.924999    |
| — + + + — | 3.295664    |
| — + + + + | 5.506271    |
| + — — — — | 0.642141    |
| + — — — + | 0.335461    |
| + — — + — | 0.606811    |
| + — — + + | 0.455699    |
| + — + — — | 0.157148    |
| + — + — + | 0.188888    |
| + — + + — | 0.439139    |
| + — + + + | 0.737664    |



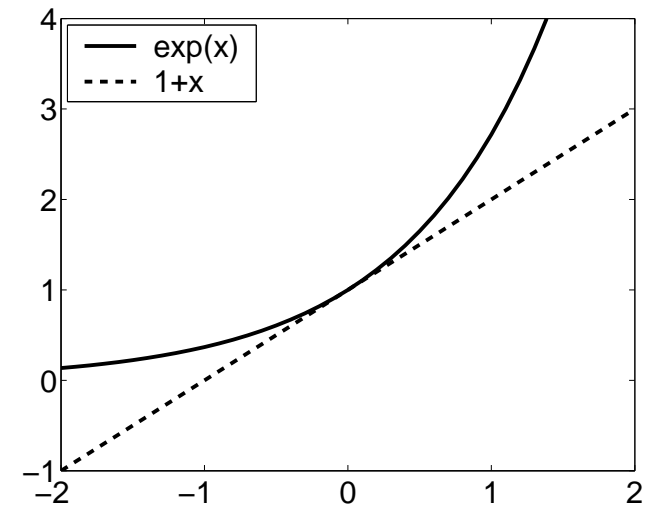
## A Simple Approximation

Approximate  $\exp(x)$  by  $1+x$

Thus

$$\begin{aligned} Z &= \sum_{\text{all } \vec{s}} \exp(-E(\vec{s})) \approx \sum_{\text{all } \vec{s}} (1 - E(\vec{s})) \\ &= 2^N - \sum_{\text{all } \vec{s}} E(\vec{s}) = 2^N \end{aligned}$$

which is a quite poor approximation.



## Variational Approaches

For all  $\mu$  we know

$$\exp(x) \geq e^\mu + e^\mu (x - \mu) = e^\mu (1 + x - \mu)$$

Thus

$$\begin{aligned} \forall_\mu Z &= \sum_{\text{all } \vec{s}} \exp(-E(\vec{s})) \geq \sum_{\text{all } \vec{s}} e^\mu (1 - E(\vec{s}) - \mu) \\ &= 2^N e^\mu (1 - \mu) = B(\mu) \end{aligned}$$

The best approximation is the maximum of  $B(\mu)$ .

$$\frac{\partial}{\partial \mu} B(\mu) = -\mu e^\mu = 0 \Rightarrow \mu = 0$$

and again we find

$$Z \geq 2^N$$

## Variational Approaches

$$\exp(x) \geq e^\mu (1 + x - \mu)$$

Thus

$$\forall_{\mu(\vec{s})} Z = \sum_{\text{all } \vec{s}} \exp(-E(\vec{s})) \geq \sum_{\text{all } \vec{s}} e^{\mu(\vec{s})} (1 - E(\vec{s}) - \mu(\vec{s}))$$

and we choose

$$\mu(\vec{s}) = \mu + \sum_i h_i s_i$$

which can be optimised with respect to  $\mu$  and  $h_i$ .



## Other Bounds

The Kullback-Leibler divergence is a bound:

$$K(q, p) = \sum_{\text{all } \vec{s}} q(\vec{s}) \log \frac{q(\vec{s})}{p(\vec{s})} \geq 0$$

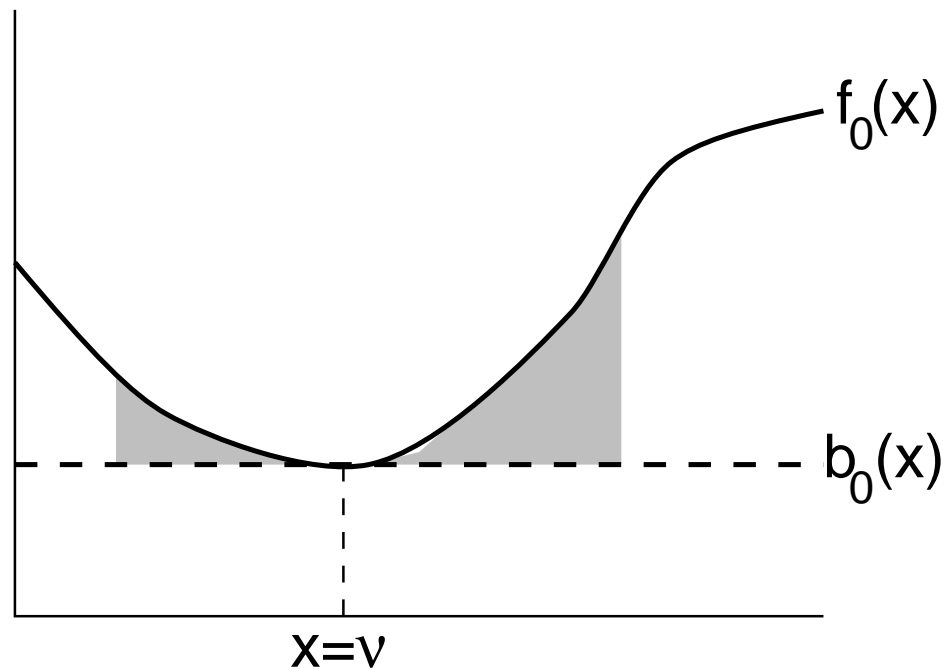
But we have also

$$\log x \leq \frac{x}{\mu} - 1 + \log \mu$$

and

$$\tanh x \leq \frac{1}{2\sqrt{2}} (x - \mu)^2 + (1 - \tanh^2 \mu) (x - \mu) + \tanh \mu$$

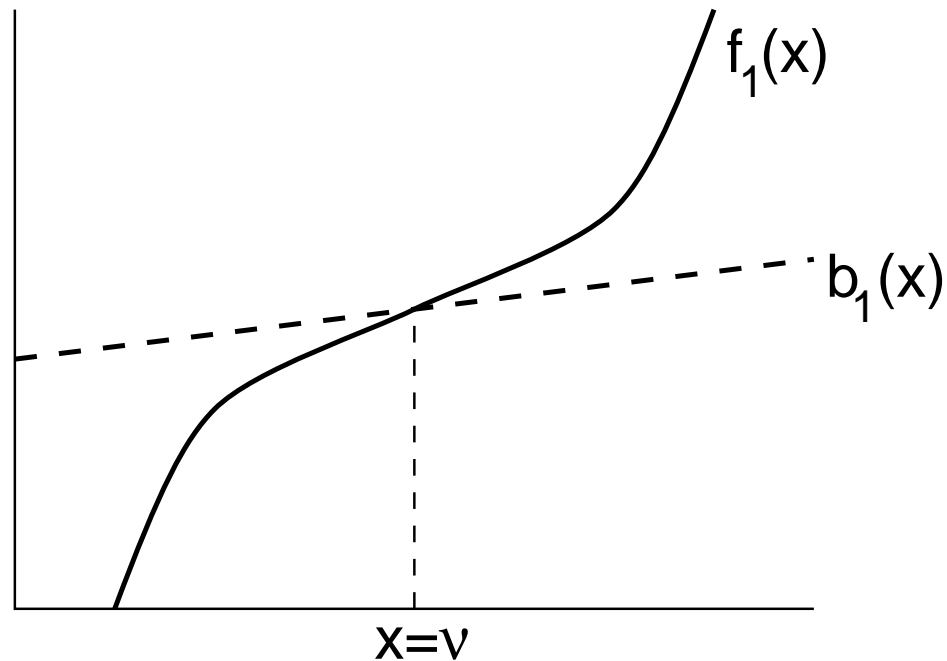
## Improving Bounds



$$f_0(x) \geq b_0(x)$$

$$f_1(x) = \int_{x'=\nu}^{x'=x} f_0(x') \, dx' \begin{cases} \geq & (x \geq \nu) \\ \leq & (x \leq \nu) \end{cases} \int_{x'=\nu}^{x'=x} b_0(x') \, dx' = b_1(x)$$

## Improving Bounds



$$f_1(x) \leq b_1(x)$$

$$f_2(x) = \int_{x'=\nu}^{x'=x} f_1(x') \, dx' \geq \int_{x'=\nu}^{x'=x} b_1(x') \, dx' = b_2(x)$$

## Improving Bounds

Given  $f_0(x) \geq b_0(x)$  we can derive

- $f_1 = \int f_0$  and  $b_1 = \int b_0$  with  $f_1(\nu) = b_1(\nu)$  for some  $\nu$
- $f_2 = \int f_1$  and  $b_2 = \int b_1$  with  $f_2(\nu) = b_2(\nu)$  for that  $\nu$

Then we know

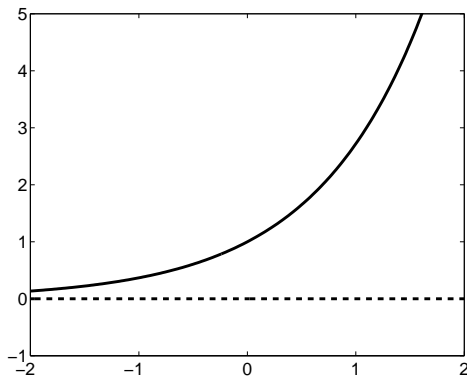
- $\forall_{x,\nu} f_2(x) \geq b_2(x)$

## Example: Exponential function

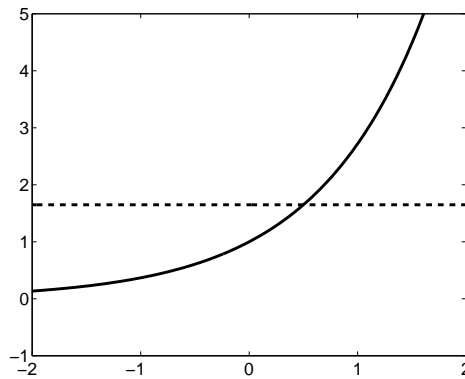
$$f_0(x) = e^x \geq 0 = b_0(x)$$

$$f_1(x) = e^x \leq e^\nu = b_1(x)$$

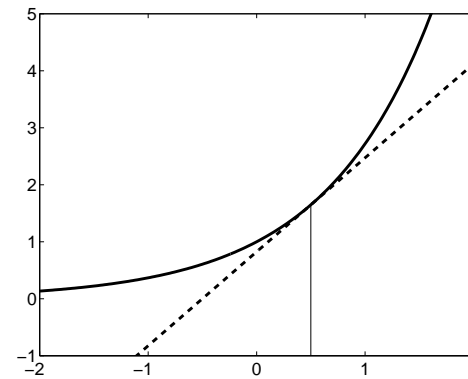
$$\forall x, \nu \quad f_2(x) = e^x \geq e^\nu (1 + x - \nu) = b_2(x)$$



$$f_0(x) \geq b_0(x)$$



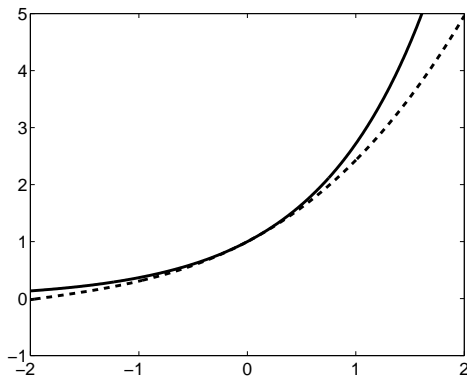
$$f_1(x) \leq b_1(x)$$



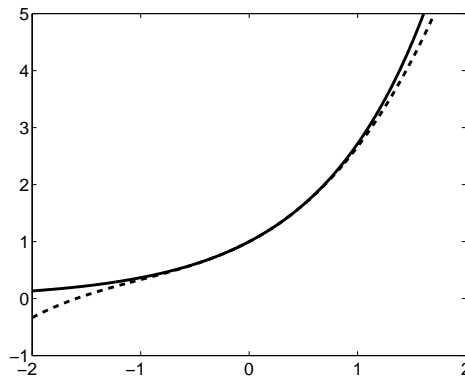
$$f_2(x) \geq b_2(x)$$

## Example: Exponential function

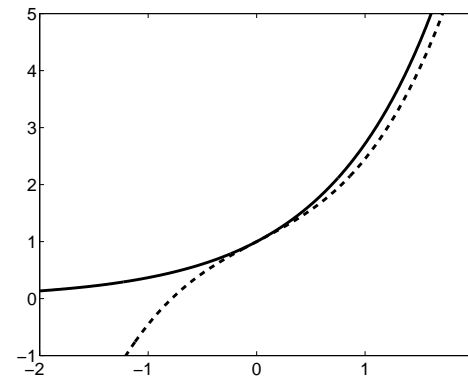
$$\forall x, \mu, \lambda \quad e^x \geq e^\mu \left\{ 1 + x - \mu + e^\lambda \left( \frac{1-\lambda}{2} (x-\mu)^2 + \frac{1}{6} (x-\mu)^3 \right) \right\}$$



$$\mu = 0$$
$$\lambda = -1$$



$$\mu = 0$$
$$\lambda = 0$$



$$\mu = 0$$
$$\lambda = +1$$

## Summary

- You have a function that is intractable (e.g.  $\sum_{\text{all } \vec{s}} \exp(-E(\vec{s}))$ )
- You can derive a (large) class of bounding functions ( $f(x) \geq b(x, \mu)$ )
- These functions are parametrized by  $\mu$ , the *variational parameters*
- The larger this class is, the better your estimate
- Optimize the class of bounding function with respect to  $\mu$  to find the tightest bound

## Take Home Message

1. I have seen a way to do approximate computations for a Boltzmann machine. But  $p(\text{I will use a BM}) = 0$ .
2. I have an intuition of what can be done with variational methods and the Boltzmann machine was just one of the applications.