

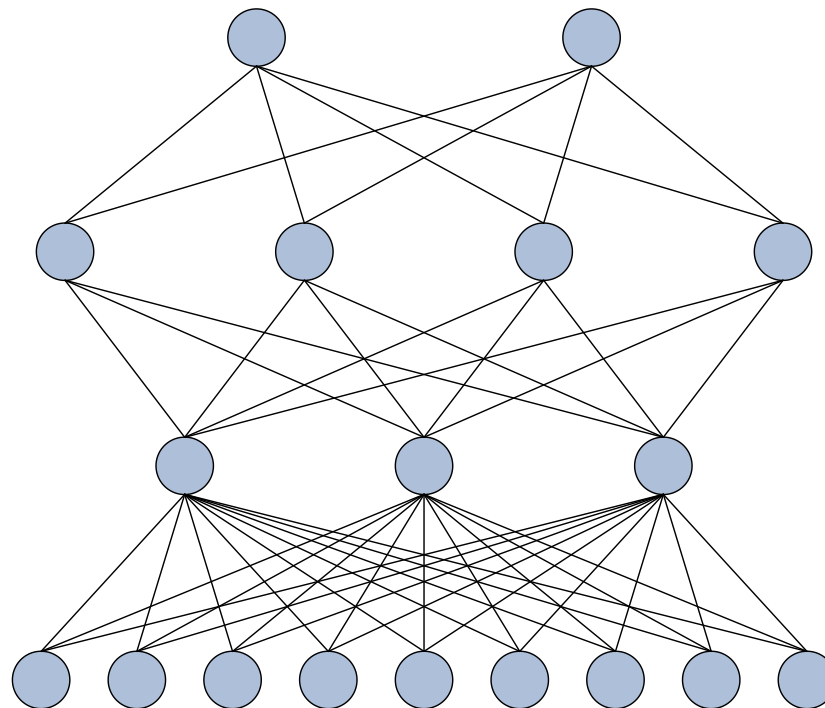
Lecture 4

Towards Deep Learning

(January 30, 2015)

Mu Zhu
University of Waterloo

Deep Network



Boltzmann Distribution

- probability distribution for a complex system

$$p(\mathbf{x}) = \frac{1}{Z} e^{f(\mathbf{x}; \boldsymbol{\theta})} \quad \text{with} \quad Z \equiv \sum_{\mathbf{x}} e^{f(\mathbf{x}; \boldsymbol{\theta})} \quad \left[\text{or} \quad \int e^{f(\mathbf{x}; \boldsymbol{\theta})} d\mathbf{x} \right]$$

- often,

$$f(\mathbf{x}; \boldsymbol{\theta}) = \frac{-u(\mathbf{x}; \boldsymbol{\vartheta})}{kT},$$

where

$u(\mathbf{x}; \boldsymbol{\vartheta})$ = energy function;

k = Boltzmann constant;

T = thermodynamic temperature

- e.g., lattice of particles, protein molecule

Boltzmann Distribution

$$\log[p(\mathbf{x}; \boldsymbol{\theta})] = f(\mathbf{x}; \boldsymbol{\theta}) - \log \left[\sum_{\mathbf{x}} e^{f(\mathbf{x}; \boldsymbol{\theta})} \right]$$

$$\begin{aligned} \frac{d}{d\boldsymbol{\theta}} \log[p(\mathbf{x}; \boldsymbol{\theta})] &= \frac{d}{d\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}) - \frac{1}{\sum_{\mathbf{x}} e^{f(\mathbf{x}; \boldsymbol{\theta})}} \left[\sum_{\mathbf{x}} e^{f(\mathbf{x}; \boldsymbol{\theta})} \cdot \frac{d}{d\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}) \right] \\ &= \frac{d}{d\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}) - \sum_{\mathbf{x}} \frac{e^{f(\mathbf{x}; \boldsymbol{\theta})}}{Z} \cdot \frac{d}{d\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}) \\ &= \frac{d}{d\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}) - \sum_{\mathbf{x}} p(\mathbf{x}) \cdot \frac{d}{d\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}) \\ &= \frac{d}{d\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}) - \mathbb{E} \left[\frac{d}{d\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}) \right] \end{aligned}$$

Boltzmann Distribution

- given

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_i \stackrel{iid}{\sim} p(\mathbf{x}; \boldsymbol{\theta})$$

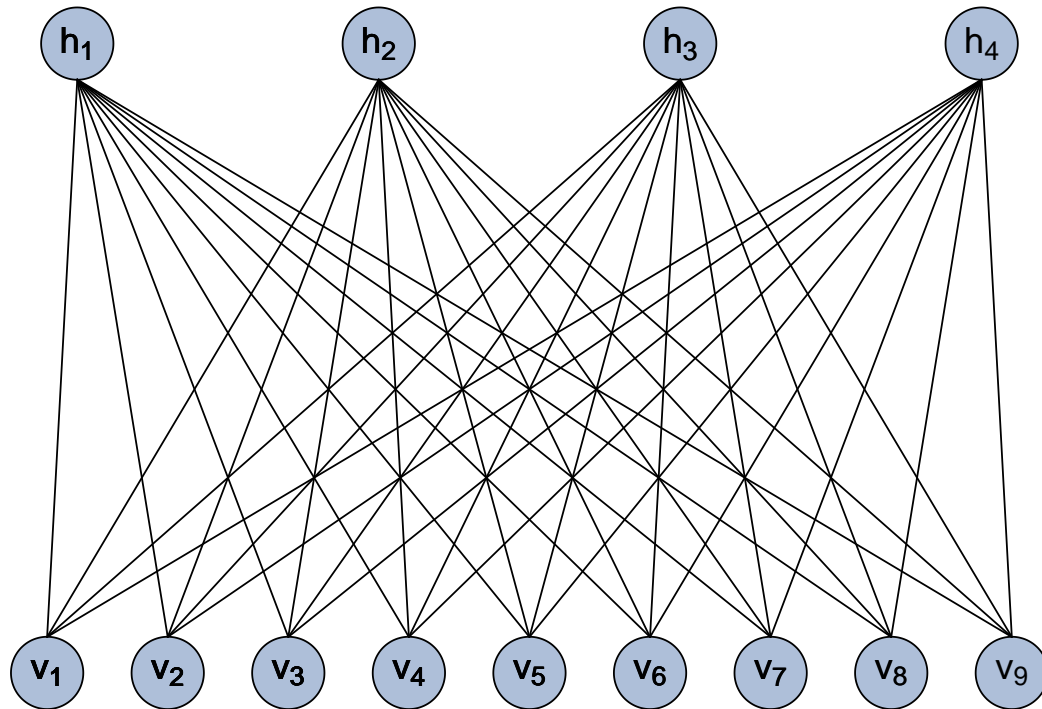
- log-likelihood is

$$\ell(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \log[p(\mathbf{x}_i; \boldsymbol{\theta})]$$

- its first derivative is

$$\begin{aligned} \frac{d}{d\boldsymbol{\theta}} \ell(\boldsymbol{\theta}) &= \frac{1}{n} \sum_{i=1}^n \left\{ \frac{d}{d\boldsymbol{\theta}} f(\mathbf{x}_i; \boldsymbol{\theta}) - \mathbb{E} \left[\frac{d}{d\boldsymbol{\theta}} f(\mathbf{x}_i; \boldsymbol{\theta}) \right] \right\} \\ &\equiv \widehat{\mathbb{E}} \left[\frac{d}{d\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}) \right] - \mathbb{E} \left[\frac{d}{d\boldsymbol{\theta}} f(\mathbf{x}; \boldsymbol{\theta}) \right] \end{aligned}$$

Restricted Boltzmann Machine



bottom nodes $\mathbf{v} = (v_1, v_2, \dots)^T$; top nodes $\mathbf{h} = (h_1, h_2, \dots)^T$

Restricted Boltzmann Machine

- bottom nodes $\mathbf{v} = (v_1, v_2, \dots)^\top$
- top nodes $\mathbf{h} = (h_1, h_2, \dots)^\top$
- Boltzmann distribution

$$p(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}) = \frac{1}{Z} e^{f(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta})}$$

with

$$\begin{aligned} f(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}) &= \mathbf{h}^\top \mathbf{W} \mathbf{v} + \boldsymbol{\alpha}^\top \mathbf{h} + \boldsymbol{\beta}^\top \mathbf{v} \\ &= \sum_t h_t \mathbf{w}_t^\top \mathbf{v} + \sum_t \alpha_t h_t + \boldsymbol{\beta}^\top \mathbf{v} \end{aligned}$$

i.e., $\boldsymbol{\theta} = \{\mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\beta}\}$

Restricted Boltzmann Machine

- if just one binary top node $h \in \{0, 1\}$
- get

$$p(\mathbf{v}, h = 1) = \frac{1}{Z} e^{\mathbf{w}^T \mathbf{v} + \alpha + \beta^T \mathbf{v}}, \quad p(\mathbf{v}, h = 0) = \frac{1}{Z} e^{\beta^T \mathbf{v}}$$

- so

$$\frac{p(\mathbf{v}, h = 1)}{p(\mathbf{v}, h = 0)} = \frac{\mathbb{P}(h = 1 | \mathbf{v}) f(\mathbf{v})}{\mathbb{P}(h = 0 | \mathbf{v}) f(\mathbf{v})}$$
$$\Rightarrow \log \frac{\mathbb{P}(h = 1 | \mathbf{v})}{\mathbb{P}(h = 0 | \mathbf{v})} = \mathbf{w}^T \mathbf{v} + \alpha$$

- hence, model for $h | \mathbf{v}$ is usual **logistic regression**

Restricted Boltzmann Machine

- if more than one binary top nodes $h_1, h_2, \dots \in \{0, 1\}$
- get

$$f(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}) = h_t \mathbf{w}_t^T \mathbf{v} + \alpha_t h_t + \sum_{s \neq t} h_s \mathbf{w}_s^T \mathbf{v} + \sum_{s \neq t} \alpha_s h_s + \boldsymbol{\beta}^T \mathbf{v}$$

- so

$$\log \frac{\mathbb{P}(h_t = 1 | \mathbf{v}, \mathbf{h}_{-t})}{\mathbb{P}(h_t = 0 | \mathbf{v}, \mathbf{h}_{-t})} = \mathbf{w}_t^T \mathbf{v} + \alpha_t$$

- in general, model for $h_t | \mathbf{v}, \mathbf{h}_{-t}$ is usual **logistic regression**
- notice **conditional independence** between h_t and \mathbf{h}_{-t} given \mathbf{v}

Fitting RBM

- given $(\mathbf{v}_i, \mathbf{h}_i)$, $i = 1, 2, \dots, n$, MLE by gradient ascent,

$$\mathbf{W}_{new} = \mathbf{W}_{old} + \varepsilon \left[\frac{d\ell}{d\mathbf{W}} \right]_{\boldsymbol{\theta}_{old}},$$

and likewise for $\boldsymbol{\alpha}$, $\boldsymbol{\beta}$

- gradients:

$$f(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}) = \mathbf{h}^T \mathbf{W} \mathbf{v} + \boldsymbol{\alpha}^T \mathbf{h} + \boldsymbol{\beta}^T \mathbf{v}$$

↓

$$\frac{d\ell}{d\mathbf{W}} = \hat{\mathbf{E}}[\mathbf{h}\mathbf{v}^T] - \mathbf{E}[\mathbf{h}\mathbf{v}^T]$$

$$\frac{d\ell}{d\boldsymbol{\alpha}} = \hat{\mathbf{E}}[\mathbf{h}] - \mathbf{E}[\mathbf{h}], \quad \frac{d\ell}{d\boldsymbol{\beta}} = \hat{\mathbf{E}}[\mathbf{v}] - \mathbf{E}[\mathbf{v}]$$

Fitting RBM

- given $\{(\mathbf{v}_i, \mathbf{h}_i)\}_{i=1}^n$, $\hat{\mathbb{E}}[\mathbf{h}\mathbf{v}^T]$, $\hat{\mathbb{E}}[\mathbf{h}]$, $\hat{\mathbb{E}}[\mathbf{v}]$ are easy to compute
 - just take empirical averages [definition of $\hat{\mathbb{E}}(\cdot)$]
- but $\mathbb{E}[\mathbf{h}\mathbf{v}^T]$, $\mathbb{E}[\mathbf{h}]$, $\mathbb{E}[\mathbf{v}]$ are hard to compute
 - estimate by drawing an **MCMC** sample from $p(\mathbf{v}, \mathbf{h})$
 - in particular, do **Gibbs Sampling** for just a few iterations

Gibbs Sampling Repeat

- given \mathbf{v} , sample \mathbf{h} from the conditional distribution of $\mathbf{h}|\mathbf{v}$;
- given \mathbf{h} , sample \mathbf{v} from the conditional distribution of $\mathbf{v}|\mathbf{h}$;

until convergence (“burn-in”).

Gibbs Sampling

- recall: model for $h_t | \mathbf{v}, \mathbf{h}_{-t}$ is usual logistic regression, so

$$h_t | \mathbf{v}, \mathbf{h}_{-t} \sim \text{Bernoulli} [\sigma (\alpha_t + \mathbf{w}_t^T \cdot \mathbf{v})]$$

- if v_1, v_2, \dots also **binary**, then symmetry gives

$$v_b | \mathbf{h}, \mathbf{v}_{-b} \sim \text{Bernoulli} [\sigma (\beta_b + \mathbf{h}^T \mathbf{w}_{\cdot b})]$$

- Gibbs sampling amounts to back-and-forth coin flips
- note: $\sigma(\cdot)$ above denotes the sigmoid function

Towards Deep Learning

- stack many RBMs on top of each other
- top nodes for layer ℓ becomes bottom nodes for layer $\ell + 1$, i.e.,

$$\mathbf{v}^{(\ell+1)} = \mathbf{h}^{(\ell)}$$

- fit the whole thing layer by layer
- note intermediate layers are latent (**hidden**, not **visible**)
- so for things like $\hat{\mathbb{E}}(\mathbf{h})$, can use $\hat{\mathbb{E}}(\tilde{\mathbf{h}})$ where $\tilde{\mathbf{h}}_i \equiv \mathbb{E}(\mathbf{h}|\mathbf{v}_i)$



kind of an EM procedure

Towards Deep Learning

- repeat

- with current parameters $\boldsymbol{\theta} = \{\boldsymbol{\alpha}, \boldsymbol{\beta}, \mathbf{W}\}$

- (a) estimate $\tilde{\mathbf{h}}_i \equiv \mathbb{E}(\mathbf{h}|\mathbf{v}_i) \Rightarrow \hat{\mathbb{E}}(\cdot)$

- (b) draw $\mathbf{h}_i \sim p(\mathbf{h}|\mathbf{v}_i)$ and $\mathbf{v}_i \sim p(\mathbf{v}|\mathbf{h}_i)$ [or repeat] $\Rightarrow \mathbb{E}(\cdot)$

- move along the gradient (a crude estimate of it)

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \varepsilon \left[\hat{\mathbb{E}}(\cdot) - \mathbb{E}(\cdot) \right]$$

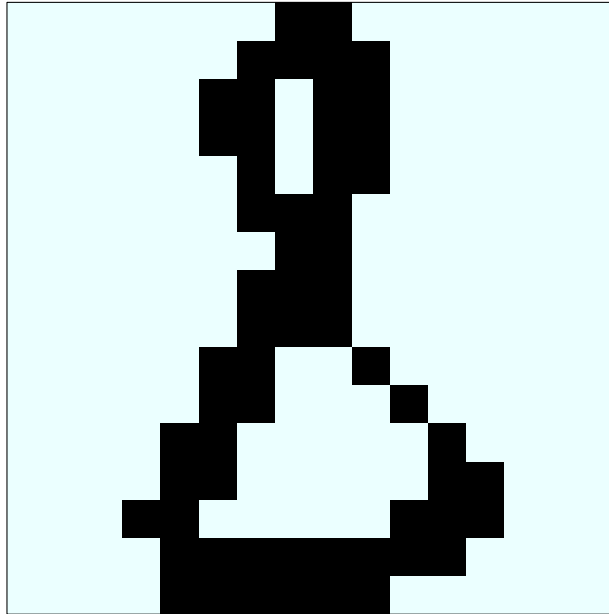
until some stopping criterion

- proceed to next layer

On v_1, v_2, \dots Being Binary

- not as restrictive as it appears
- can already handle **image** data and **text** data
 - to some extent [examples next two slides]
- can generalize to other inputs
 - make some changes to the model $f(\mathbf{v}, \mathbf{h})$
 - $p(\mathbf{v}|\mathbf{h})$ no longer logistic model [obviously]
 - ideally, want $p(\mathbf{v}|\mathbf{h})$ “nice” to sample from

Example: Image Data



each v_1, v_2, \dots a binary pixel

Example: Text Data

	Word 1	Word 2	Word 3
"hate"			
"I"	1		
"love"		1	
"math"			
"you"			1

each v_1, v_2, \dots an indicator for a particular word

The Gaussian-Bernoulli RBM

- $v_1, v_2, \dots \in \{0, 1\}$:

$$\begin{aligned} f(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}) &= \mathbf{h}^T \mathbf{W} \mathbf{v} + \boldsymbol{\alpha}^T \mathbf{h} + \boldsymbol{\beta}^T \mathbf{v} \\ &= \sum_b [\mathbf{h}^T \mathbf{w}_{\cdot b}] v_b + \boldsymbol{\alpha}^T \mathbf{h} + \sum_b \beta_b v_b \end{aligned}$$

- $v_1, v_2, \dots \in \mathbb{R}$:

$$f(\mathbf{v}, \mathbf{h}; \boldsymbol{\theta}) = \sum_b [\mathbf{h}^T \mathbf{w}_{\cdot b}] \left[\frac{v_b}{\tau_b} \right] + \boldsymbol{\alpha}^T \mathbf{h} - \sum_b \frac{(v_b - \beta_b)^2}{2\tau_b^2}$$

Exercise Let $\tilde{\mathbf{v}} = (v_1/\tau_1, v_2/\tau_2, \dots)^T$. Show that

- (a) $p(h_t | \mathbf{v}, \mathbf{h}_{-t}) \sim \text{Bernoulli}[\sigma(\alpha_t + \mathbf{w}_{t \cdot}^T \tilde{\mathbf{v}})]$; (easy)
- (b) $p(v_b | \mathbf{h}, \mathbf{v}_{-b}) \sim \text{N}(\beta_b + \tau_b(\mathbf{h}^T \mathbf{w}_{\cdot b}), \tau_b^2)$. (slightly harder)

Iterative Optimization

- consider an iterative rule, $x_{t+1} = m(x_t)$, for minimizing $f(x)$

- Newton's method:

$$m(x_t) = x_t - \frac{f'(x_t)}{f''(x_t)}$$

- gradient descent:

$$m(x_t) = x_t - f'(x_t)$$

- multivariate case ... same idea, with **gradient** and **Hessian**, i.e.,

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \varepsilon \mathbf{H}_t^{-1} \mathbf{g}_t, \quad \mathbf{x}_{t+1} = \mathbf{x}_t - \varepsilon \mathbf{g}_t;$$

will explain extra ε later

Iterative Optimization

- suppose iteration converges to a **local minimum** x^*
- must have $m(x^*) = x^*$ [i.e., a **fixed point**]
- then, in neighborhood near x^* ,

$$\begin{aligned} e_{t+1} &\equiv x_{t+1} - x^* = m(x_t) - m(x^*) \\ &= \left[m(x^*) + m'(x^*)(x_t - x^*) + \frac{m''(\xi)}{2}(x_t - x^*)^2 \right] - m(x^*) \\ &= m'(x^*)e_t + \frac{m''(\xi)}{2}e_t^2 \end{aligned}$$

- thus, a “basic requirement” is $|m'(x^*)| < 1$ (and nearby x^*)

Gradient Descent

$$m(x_t) = x_t - f'(x_t) \quad \Rightarrow \quad m'(x^*) = 1 - f''(x^*)$$

$$x^* \text{ local minimum} \quad \Rightarrow \quad f''(x^*) > 0$$

so, need $f''(x^*) < 1$ (and nearby x^*)

ensure by letting $m(x_t) = x_t - \varepsilon f'(x_t)$

Newton's Method

$$m(x_t) = x_t - \frac{f'(x_t)}{f''(x_t)}$$

$$m'(x^*) = 1 - \frac{f''(x^*)f''(x^*) - f'''(x^*)f'(x^*)}{[f''(x^*)]^2} = \frac{f'''(x^*)f'(x^*)}{[f''(x^*)]^2} = 0$$

get $e_{t+1} \sim O(e_t^2)$, much faster local convergence

but second derivative $\mathbf{H}_{d \times d}$ is not easy for large d

Quasi-Newton

- use local **curvature** information \Rightarrow fast local convergence
- want to “do some of this” without computing \mathbf{H}_t
- key idea: construct a sequence \mathbf{B}_t to “mimic” \mathbf{H}_t
- example [**symmetric rank-1 (SR1)**]: construct \mathbf{B}_{t+1} so that
 - (a) $\mathbf{g}_{t+1} = \mathbf{g}_t + \mathbf{B}_{t+1}(\mathbf{x}_{t+1} - \mathbf{x}_t)$ [**\mathbf{B}_{t+1} “like a Hessian”**]
 - (b) $\mathbf{B}_{t+1} = \mathbf{B}_t + \mathbf{u}\mathbf{v}^T$ [**update “just a little”**]
 - (c) $\mathbf{u} \propto \mathbf{v}$ [**ensures symmetry**]
- \exists many variations ...

Some Details of SR1

$$(\mathbf{B}_t + \mathbf{u}\mathbf{v}^\top) \underbrace{(\mathbf{x}_{t+1} - \mathbf{x}_t)}_{\Delta\mathbf{x}_t} = \underbrace{\mathbf{g}_{t+1} - \mathbf{g}_t}_{\Delta\mathbf{g}_t}$$

$$\underbrace{\mathbf{u} [\mathbf{v}^\top(\Delta\mathbf{x}_t)]}_{\text{scalar}} = (\Delta\mathbf{g}_t) - \mathbf{B}_t(\Delta\mathbf{x}_t) \quad \Rightarrow \quad \mathbf{u} = \frac{(\Delta\mathbf{g}_t) - \mathbf{B}_t(\Delta\mathbf{x}_t)}{\mathbf{v}^\top(\Delta\mathbf{x}_t)}$$

taking $\mathbf{u} = \mathbf{v} = \gamma [(\Delta\mathbf{g}_t) - \mathbf{B}_t(\Delta\mathbf{x}_t)]$ leads to

$$\mathbf{B}_{t+1} = \mathbf{B}_t + \frac{[(\Delta\mathbf{g}_t) - \mathbf{B}_t(\Delta\mathbf{x}_t)] [(\Delta\mathbf{g}_t) - \mathbf{B}_t(\Delta\mathbf{x}_t)]^\top}{[(\Delta\mathbf{g}_t) - \mathbf{B}_t(\Delta\mathbf{x}_t)]^\top (\Delta\mathbf{x}_t)}$$

Summary

- key ideas:
 - RBMs (building block for deep learning)
 - gradient descent; Gibbs sampling
 - local convergence behavior (gradient *vs* Newton)
- specific techniques:
 - quasi-Newton (SR1)
 - Gaussian-Bernoulli RBM

Next ...

- a short, 10-minute break
- lecture by Dr. R. Grosse on some current work about [RBMs](#)