

Course ID 80250993 Dates: 9/16-12/23/2021 @ I-205
Tencent/VooV Meeting 521 4678 6257 (passwd 1205)



Chapter 16

Unsupervised Learning Neural Networks (part 2)

Xuegong Zhang
Dec. 2, 2021



16.2

Hopfield Network



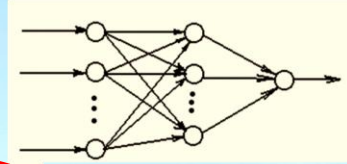
Xuegong Zhang

2

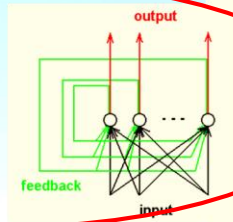


Three Major Types of ANN (in 1980s)

- Feedforward NN
 - Multi-Layer Perceptron

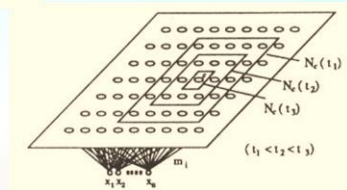


- Feedback NN
 - Hopfield NN



Energy-based neural network

- Competitive Learning NN
 - Self-organizing map



Xuegang Zhang

3

Xuegang Zhang

4

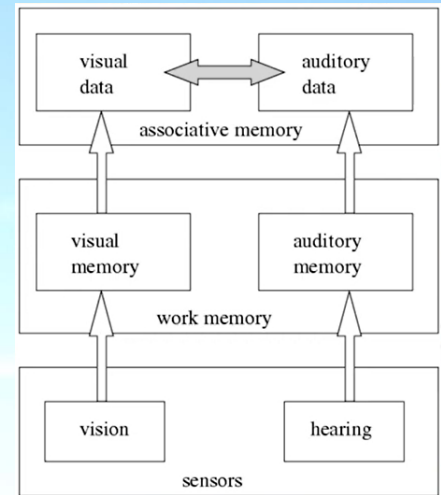
Associate Memory



- Associate the part with the whole



- Associate one type of memory with another type



Xuegong Zhang

5

Hopfield Network



- First by Little in 1974 and popularized in 1982 by John Hopfield
 - W.A. Little, The existence of persistent states in the brain, *Mathematical Biosciences* 19: 101-120, 1974
 - John J. Hopfield, Neural networks and physical systems with emergent collective computational abilities, *PNAS* 79: 2554-2558, 1982
- For “associative memory” or “content-addressable memory (CAM)”

MATHEMATICAL BIOSCIENCES 19, 101-120 (1974)

The Existence of Persistent States in the Brain

W. A. LITTLE
Department of Physics, Stanford University, Stanford, California

Communicated by S. M. Ulam

ABSTRACT

We show that given certain plausible assumptions the existence of persistent states in a neural network can occur only if a certain transfer matrix has degenerate maximum eigenvalues. The existence of such states of persistent order is directly analogous to the existence of long range order in an Ising spin system; while the transition to the state of persistent order is analogous to the transition to the ordered phase of the spin system. It is shown that the persistent state is also characterized by correlations between neurons throughout the brain. It is suggested that these persistent states are associated with short term memory while the eigenvectors of the transfer matrix are a representation of long term memory. A numerical example is given that illustrates certain of these features.

Proc. Natl. Acad. Sci. USA
Vol. 79, pp. 2554-2558, April 1982
Biophysics

Neural networks and physical systems with emergent collective computational abilities

(associative memory/parallel processing/categorization/content-addressable memory/full soft devices)

J. J. HOPFIELD

ABSTRACT Computational properties of use to biological organisms or to the construction of computers can emerge as collective properties of systems having a large number of simple equivalent components (or neurons). The physical meaning of content-addressable memory is described by an appropriate phase space flow of the state of a system. A model of such a system is given, based on aspects of neurobiology but readily adapted to integrated circuits. The collective properties of this model produce a content-addressable memory which correctly yields an entire memory from any subpart of sufficient size. The algorithm for the time evolution of the state of the system is based on asynchronous parallel processing. Additional emergent collective properties include some capacity for generalization, familiarity recognition, categorization, error correction, and time sequence retention. The collective properties are only weakly sensitive to details of the modeling or the failure of individual devices.



John J. Hopfield
(July 15, 1933 -)

Xuegong Zhang

6

Biology is a Dynamic System



- John J. Hopfield:

“... If you look at the feedforward artificial neural nets, they are done as: first of all I spend some time learning, not performing, then turn off learning and I perform.

That's NOT biology.

... As I look more deeply at neurobiology, even as associate memory, I got to face the fact that the dynamics of synapse changes is going on all the time, ”



Xuegong Zhang

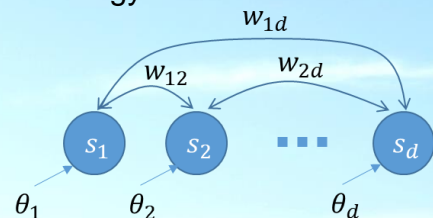
7

Hopfield Net as a Dynamic System



- An array of **binary threshold units** (McCulloch-Pitts neurons) fully connected with **symmetric weights**
- Each binary “configuration” of the network has an energy
- Energy (Lyapunov function)

$$E = -\frac{1}{2} \sum_{i,j} w_{ij} s_i s_j + \sum_i \theta_i s_i$$



Xuegong Zhang

8

Hopfield Net as a Dynamic System

- Min Energy

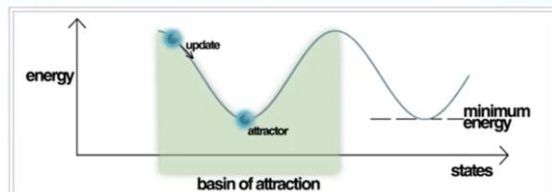
$$E = -\frac{1}{2} \sum_{i,j} w_{ij} s_i s_j + \sum_i \theta_i s_i$$

- For each node

$$\Delta E_i = E(s_i = 0) - E(s_i = 1) = \sum_j w_{ij} s_j - \theta_i \quad \text{--- the energy gap}$$

If $\Delta E_i \geq 0$, $s_i = 1$ lowers the energy; otherwise $s_i = 0$ lowers the energy.

- The “binary threshold decision rule” makes the network to settle to a minimum of energy function



Energy Landscape of a Hopfield Network, highlighting the current state of the network (up the hill), an attractor state to which it will eventually converge, a minimum energy level and a basin of attraction shaded in green. Note how the update of the Hopfield Network is always going down in Energy.

Xuegong Zhang

9

Hopfield Net as a Dynamic System

- The “binary threshold decision rule”**

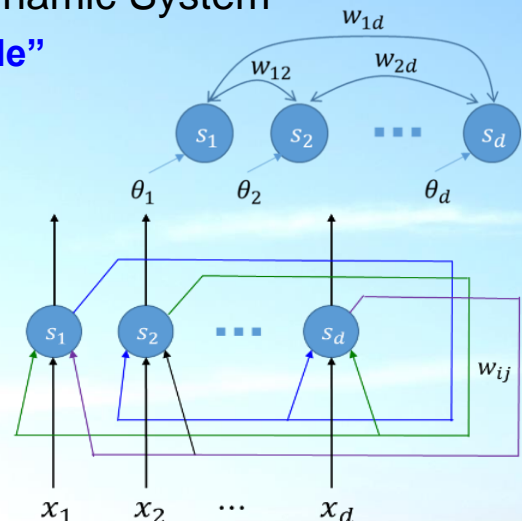
- At time 0, $s = [s_1, \dots, s_d]^T = x$,

$$x = [x_1, \dots, x_d]^T \in \{1, -1\}^d$$

- $s_i(t) = f_w(s_i(t-1))$, $s_i(0) = x_i$

- At time t , updating:

$$s_i = \begin{cases} +1 & \text{if } \sum_j w_{ij} s_j \geq \theta_i \\ -1 & \text{else} \end{cases}$$



Xuegong Zhang

10

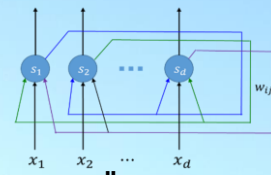
Hopfield Network

- Min Energy (Lyapunov function)

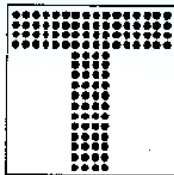
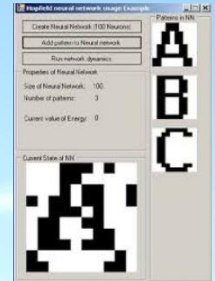
$$E = -\frac{1}{2} \sum_{i,j} w_{ij} s_i s_j + \sum_i \theta_i s_i$$

- The dynamic procedure converges to a “memory”.

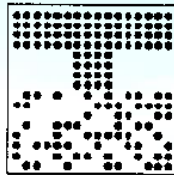
- Memory \leftrightarrow energy minima of neural net



Associative memory



Original 'T'



half of image
corrupted by
noise

“It’s like reconstructing a dinosaur from a few bones.”
— Hinton

--- A good tutorial: <http://web.cs.ucla.edu/~rosen/161/notes/hopfield.html>

Xuegong Zhang

11

Hopfield Network

- How to store the memory in?

- Learning Algorithm:

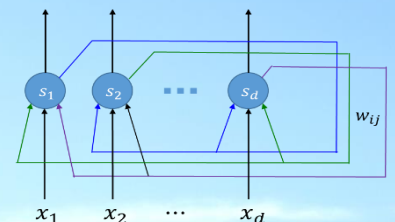
- Hebb rule: “Neurons that fire together wire together”

if $x_i, x_j \in \{-1, 1\}$,

$$w_{ij} = \frac{1}{n} \sum_{k=1}^n x_i^k x_j^k$$

if $x_i, x_j \in \{0, 1\}$,

$$w_{ij} = \frac{1}{n} \sum_{k=1}^n (2x_i^k - 1)(2x_j^k - 1)$$



Hebb's 1949

“When an axon of cell A . . . excites cell B and repeatedly or persistently takes part in firing it, some growth process or metabolic change takes place in one or both cells so that A's efficiency as one of the cells firing B is increased.”

The Organization of Behavior: A Neuropsychological Theory



--- A good tutorial: <http://web.cs.ucla.edu/~rosen/161/notes/hopfield.html>

Xuegong Zhang

12



The storage capacity of a Hopfield Net

- Using Hopfield's storage rule the capacity of a totally connected net with N units is only about $0.15N$ memories.
 - At N bits per memory this is only $0.15 N^2$ bits.
 - This does not make efficient use of the bits required to store the weights.
- The net has N^2 weights and biases.
- After storing M memories, each connection weight has an integer value in the range $[-M, M]$.
- So the number of bits required to store the weights and biases is: $N^2 \log(2M + 1)$

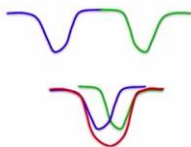


Avoiding spurious minima by unlearning

- Hopfield, Feinstein and Palmer suggested the following strategy:
 - Let the net settle from a random initial state and then do **unlearning**.
 - This will get rid of deep, spurious minima and increase memory capacity.
- They showed that this worked.
 - But they had no analysis.
- Crick and Mitchison proposed unlearning as a model of what dreams are for.
 - That's why you don't remember them (unless you wake up during the dream)
- But how much unlearning should we do?
 - Can we derive unlearning as the right way to minimize some cost function?

Spurious minima limit capacity

- Each time we memorize a configuration, we hope to create a new energy minimum.
 - But what if two nearby minima merge to create a minimum at an intermediate location?
 - This limits the capacity of a Hopfield net.



The state space is the corners of a hypercube. Showing it as a 1-D continuous space is a misrepresentation.

— Slide from Hinton's lecture

Xuegong Zhang

13



10-second break



Xuegong Zhang

14



16.3

Restrictive Boltzmann Machines

Xuegong Zhang

15

Restricted Boltzmann Machines (RBM)



COGNITIVE SCIENCE 9, 147-169 (1985)

A Learning Algorithm for Boltzmann Machines*

DAVID H. ACKLEY
GEOFFREY E. HINTON
*Computer Science Department
Carnegie-Mellon University*
TERRENCE J. SEJNOWSKI
*Biophysics Department
The Johns Hopkins University*

The computational power of massively parallel networks of simple processing elements resides in the communication bandwidth provided by the hardware connections between elements. These connections can allow a significant



Reducing the Dimensionality of Data with Neural Networks

G. E. Hinton* and R. R. Salakhutdinov

High-dimensional data can be converted to low-dimensional codes by training a multilayer neural network with a small central layer to reconstruct high-dimensional input vectors. Gradient descent can be used for fine-tuning the weights in such "autoencoder" networks, but this works well only if the initial weights are close to a good solution. We describe an effective way of initializing the weights that allows deep autoencoder networks to learn low-dimensional codes that work much better than principal components analysis as a tool to reduce the dimensionality of data.

red SHG
x in Fig.
we find
he noise
l closely
it power
mission
gle with
ms from
Rs (see
detuning
velength
he SHG
For ex-
vertical
a small
xcitation
incident
gnificant
gain po-

Dimensionality reduction facilitates the classification, visualization, communication, and storage of high-dimensional data. A simple and widely used method is principal components analysis (PCA), which

finds the directions of greatest variance in the data set and represents each data point by its coordinates along each of these directions. We describe a nonlinear generalization of PCA that uses an adaptive, multilayer "encoder" network

28 JULY 2006 VOL 313 SCIENCE www.sciencemag.org

Xuegong Zhang

16

What are "Patterns"?

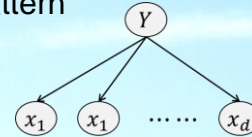


- Simple patterns

- Single feature: distribution \rightarrow pattern



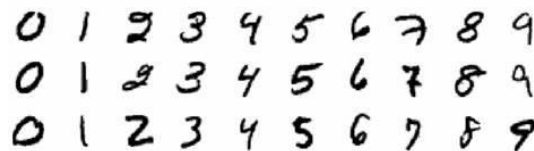
- Multiple independent features: distributions \rightarrow pattern



- Complex patterns

- Multiple related features:

relations among elements/features \rightarrow pattern



Xuegong Zhang

17

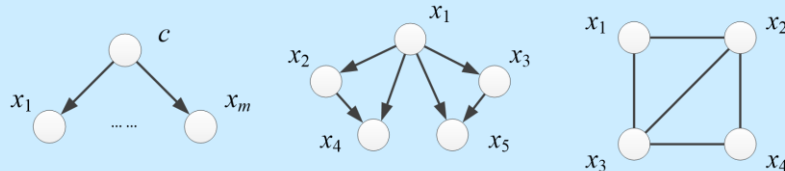
Question: How can we represent/store relations among elements?



Probabilistic Graphical Models



- Present the joint distribution of a set of random variables using a graph
 - Node: random variable
 - Edge: conditional dependence
 - Non-edge: conditional independence
- Applications:
 - computer vision, speech recognition
 - gene regulatory network inference
 - ...



18

What are “Patterns”?



• How can we capture relations among features?

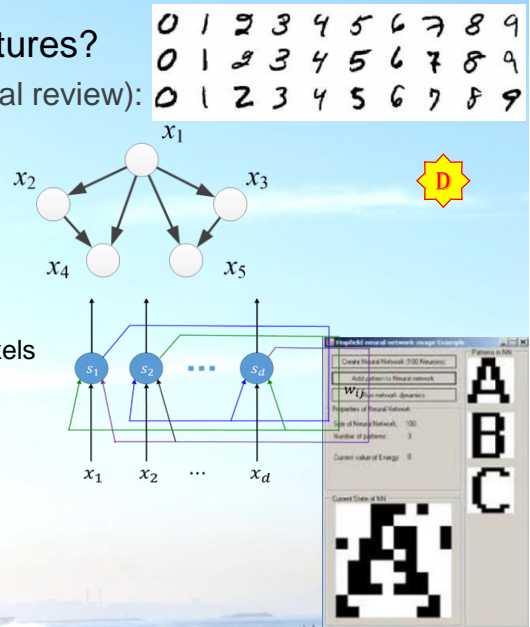
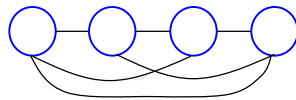
- Questions (to be thought about during the final review):
 - How does MLP capture relations among pixels?
 - How does SVM capture relations among pixels?
 - How does CNN capture relations among pixels?

• Knowledge → model → Bayesian net

• How does Hopfield Net capture relations?

- Hebb-rule training: finding associations among pixels

• Can we build a full joint probabilistic model?



Xuegong Zhang



16.3.1 Boltzmann Machines

Xuegong Zhang

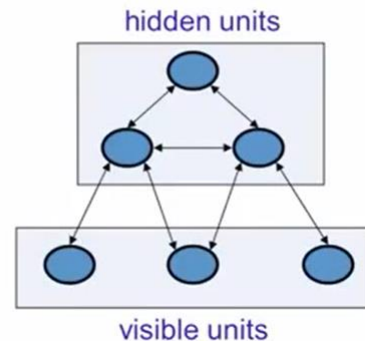
20

What are “Patterns”?



A different computational role for Hopfield nets

- Instead of using the net to store memories, use it to construct interpretations of sensory input.
 - The input is represented by the visible units.
 - The interpretation is represented by the states of the hidden units.
 - The badness of the interpretation is represented by the energy.



— Slide from Hinton's lecture

Xuegong Zhang

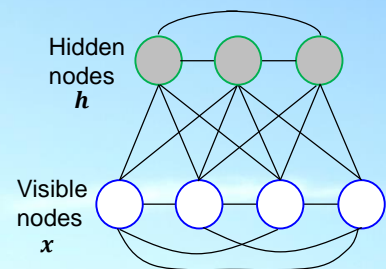
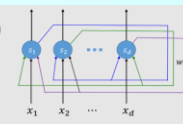
21

Boltzmann Machines (G.E. Hinton et al, 1983)



- Binary features (0 or 1 pixels) → visible nodes
 - Binary vector x
- Underlying relations → binary hidden nodes
 - Binary vector h
- State of the system: $\{x, h\}$
- Recall
 - Min Energy (Lyapunov function)

$$E = -\frac{1}{2} \sum_{i,j} w_{ij} s_i s_j + \sum_i \theta_i s_i$$



- What is a proper energy for Boltzmann Machines?

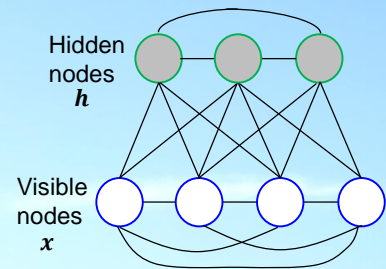
Xuegong Zhang

22

Boltzmann Machines (G.E. Hinton et al, 1983)



- A joint configuration of the network: (x, h)
- Two equivalent views:
 - Define the probability of a network configuration as $p(x, h) \propto e^{-E(x, h)}$
 - Define the probability as the probability of finding the network in that configuration after the stochastic binary units are updated many times



Xuegong Zhang

23

Boltzmann Machines (G.E. Hinton et al, 1983)

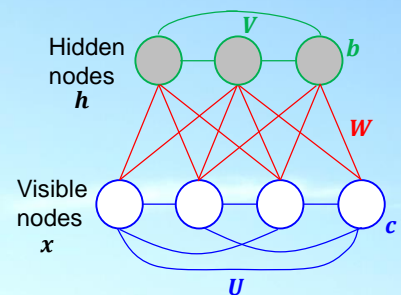


- Define the energy function as

$$E(x, h) = -\mathbf{c}^T \mathbf{x} - \mathbf{b}^T \mathbf{h} - \mathbf{h}^T \mathbf{W} \mathbf{x} - \mathbf{x}^T \mathbf{U} \mathbf{x} - \mathbf{h}^T \mathbf{V} \mathbf{h}$$

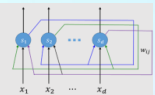
Offsets

Connection Weights



- Min Energy (Lyapunov function)

$$E = -\frac{1}{2} \sum_{i,j} w_{ij} s_i s_j + \sum_i \theta_i s_i$$



Xuegong Zhang

24

Boltzmann Machines (G.E. Hinton et al, 1983)



- The probability of a joint configuration

$$p(\mathbf{x}, \mathbf{h}) \propto e^{-E(\mathbf{x}, \mathbf{h})}$$

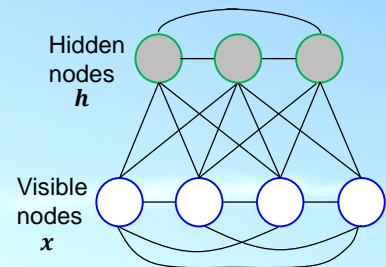
$$p(\mathbf{x}, \mathbf{h}) = \frac{e^{-E(\mathbf{x}, \mathbf{h})}}{\sum_{\mathbf{x}, \mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})}}$$

- The probability of a configuration of visible nodes \mathbf{x}

$$p(\mathbf{x}) = \frac{e^{-E(\mathbf{x})}}{Z} = \frac{\sum_{\mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})}}{\sum_{\mathbf{x}, \mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})}}$$

- Normalization factor: the partition function

$$Z = \sum_{\mathbf{x}} e^{-E(\mathbf{x})} = \sum_{\mathbf{x}, \mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})}$$



— The Boltzmann distribution

Xuegong Zhang

25

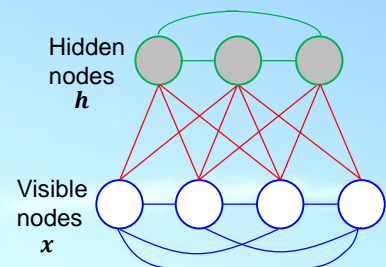
Boltzmann Machines (G.E. Hinton et al, 1983)



- Energy with \mathbf{x} : consider all possible hidden vectors \mathbf{h}

$$E(\mathbf{x}) = -\log \sum_{\mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})}$$

$$E(\mathbf{x}, \mathbf{h}) = -\mathbf{c}^T \mathbf{x} - \mathbf{b}^T \mathbf{h} - \mathbf{h}^T \mathbf{W} \mathbf{x} - \mathbf{x}^T \mathbf{U} \mathbf{x} - \mathbf{h}^T \mathbf{V} \mathbf{h}$$



→ A fully-connected Boltzmann machine

- Can represent complex probabilistic relations
- But parameter estimation is hard

Xuegong Zhang

26



Recall: what do we do for a model

Three Central Problems in HMM



- The Evaluation Problem
 - Given an HMM, complete with transition and emission probabilities, how to determine **the probability of a particular sequence of symbols** being generated by that model?
- The Decoding Problem
 - Given an HMM as well as a set of observations, how to determine the **most likely sequence of hidden states** that led to those observations?
- The Learning Problem
 - Given the coarse structure of the model but **not** the probabilities, and given a set of training observations of symbols, how to determine **the probabilities**?



Xuegong Zhang

27



16.3.2 Restrictive Boltzmann Machines (RBM)

Xuegong Zhang

28

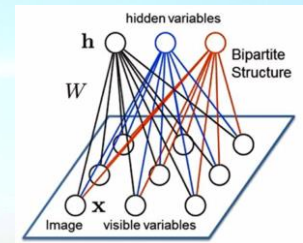
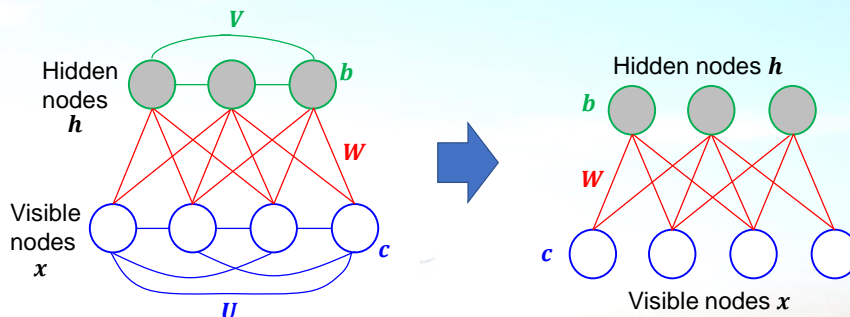
Restricted Boltzmann Machines (RBM)



Harmony Network or Harmoniums (Smolensky, 1986)

- “Restricted” connections:

- Restricted to connections between visible and hidden nodes
- Bipartite graph
- Visible nodes $x \in \{0,1\}^D$, Hidden nodes $h \in \{0,1\}^H$



Xuegong Zhang

29

Restricted Boltzmann Machines (RBM)



Harmony Network or Harmoniums (Smolensky, 1986)

- Energy function

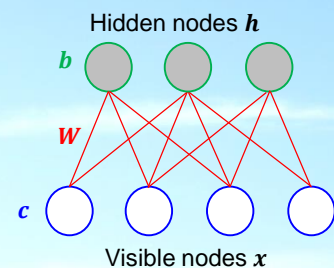
$$E(x, h) = -c^T x - b^T h - h^T W x$$

- Joint probability

$$p(x, h) = e^{-E(x, h)} / Z$$

- Normalization factor (partition function)

$$Z = \sum_{x, h} e^{-E(x, h)}$$



Xuegong Zhang

30

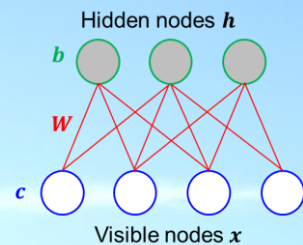
Restricted Boltzmann Machines (RBM)



- Energy function

$$E(\mathbf{x}, \mathbf{h}) = -\mathbf{h}^T \mathbf{W} \mathbf{x} - \mathbf{c}^T \mathbf{x} - \mathbf{b}^T \mathbf{h}$$

$$= -\sum_j \sum_k \mathbf{W}_{j,k} h_j x_k - \sum_k \mathbf{c}_k x_k - \sum_j \mathbf{b}_j h_j$$



- Joint probability

$$p(\mathbf{x}, \mathbf{h}) = e^{-E(\mathbf{x}, \mathbf{h})} / Z$$

$$= \frac{1}{Z} \exp(\sum_j \sum_k \mathbf{W}_{j,k} h_j x_k + \sum_k \mathbf{c}_k x_k + \sum_j \mathbf{b}_j h_j)$$

$$= \frac{1}{Z} \prod_j \prod_k \exp(\mathbf{W}_{j,k} h_j x_k) \prod_k \exp(\mathbf{c}_k x_k) \prod_j \exp(\mathbf{b}_j h_j)$$

Correction in (12-23) on p.305 of the textbook

—— Product of Experts

Xuegong Zhang

31

Restricted Boltzmann Machines (RBM)



$$E(\mathbf{x}, \mathbf{h}) = -\mathbf{h}^T \mathbf{W} \mathbf{x} - \mathbf{c}^T \mathbf{x} - \mathbf{b}^T \mathbf{h}$$

$$= -\sum_j \sum_k \mathbf{W}_{j,k} h_j x_k - \sum_k \mathbf{c}_k x_k - \sum_j \mathbf{b}_j h_j$$

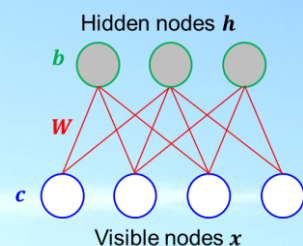
$$p(\mathbf{x}, \mathbf{h}) = e^{-E(\mathbf{x}, \mathbf{h})} / Z$$

$$= \frac{1}{Z} \prod_j \prod_k \exp(\mathbf{W}_{j,k} h_j x_k) \prod_k \exp(\mathbf{c}_k x_k) \prod_j \exp(\mathbf{b}_j h_j)$$

$$p(\mathbf{x}) = \frac{1}{Z} e^{-E(\mathbf{x})} = \frac{1}{Z} \sum_{\mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})}$$

- Conditional probability

$$p(\mathbf{h}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{h})}{p(\mathbf{x})} = \frac{1}{Z'} \prod_j \exp(\mathbf{b}_j h_j + h_j \mathbf{W}_{j.} \mathbf{x})$$



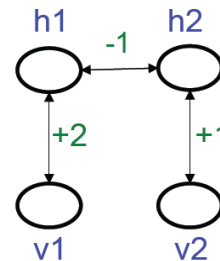
Xuegong Zhang

32

v	h	-E	e^{-E}	$p(\mathbf{v}, \mathbf{h})$	$p(\mathbf{v})$
1 1	1 1	2	7.39	.186	0.466
1 1	1 0	2	7.39	.186	
1 1	0 1	1	2.72	.069	
1 1	0 0	0	1	.025	
1 0	1 1	1	2.72	.069	0.305
1 0	1 0	2	7.39	.186	
1 0	0 1	0	1	.025	
1 0	0 0	0	1	.025	
0 1	1 1	0	1	.025	0.144
0 1	1 0	0	1	.025	
0 1	0 1	1	2.72	.069	
0 1	0 0	0	1	.025	
0 0	1 1	-1	0.37	.009	0.084
0 0	1 0	0	1	.025	
0 0	0 1	0	1	.025	
0 0	0 0	0	1	.025	

39.70

An example of how weights define a distribution



— Slide from Hinton's lecture

Xuegong Zhang

33

Restricted Boltzmann Machines (RBM)

- Conditional probability

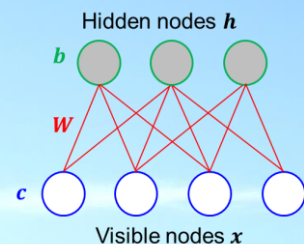
$$p(\mathbf{h}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{h})}{p(\mathbf{x})} = \frac{1}{Z'} \prod_j \exp(b_j h_j + h_j \mathbf{W}_j \cdot \mathbf{x})$$

$$P(h_j = 1|\mathbf{x}) = \frac{\exp(b_j + \mathbf{W}_j \cdot \mathbf{x})}{1 + \exp(b_j + \mathbf{W}_j \cdot \mathbf{x})} = \text{sigm}(b_j + \mathbf{W}_j \cdot \mathbf{x})$$

$$P(h_j = 0|\mathbf{x}) = 1 - \text{sigm}(b_j + \mathbf{W}_j \cdot \mathbf{x}) = \text{sigm}(-b_j - \mathbf{W}_j \cdot \mathbf{x})$$

$$p(\mathbf{h}|\mathbf{x}) = \prod_j p(h_j|\mathbf{x}) = \prod_j \text{sigm}((2h_j - 1)(b_j + \mathbf{W}_j \cdot \mathbf{x}))$$

$$p(\mathbf{x}|\mathbf{h}) = \prod_k p(x_k|\mathbf{h}) = \prod_k \text{sigm}((2x_k - 1)(c_k + \mathbf{h} \cdot \mathbf{W}_{\cdot k}))$$



sigm(): sigmoid function

Xuegong Zhang

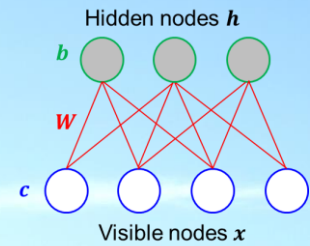
34

Restricted Boltzmann Machines (RBM)



$$p(\mathbf{h}|\mathbf{x}) = \prod_j \text{sigm}((2h_j - 1)(b_j + \mathbf{W}_j \cdot \mathbf{x}))$$

$$p(\mathbf{x}|\mathbf{h}) = \prod_k \text{sigm}((2x_k - 1)(c_k + \mathbf{h}\mathbf{W}_{\cdot k}))$$



• A generative model

- Knowing the model and data, we can sample the hidden states
- Giving the model and hidden states, we can sample the data

Xuegong Zhang

35

单选题 1分

设置

4-min break

- ☒ A Yes, I'm back to my computer.
- ☐ B Sorry, not yet.



Xuegong Zhang

提交

36

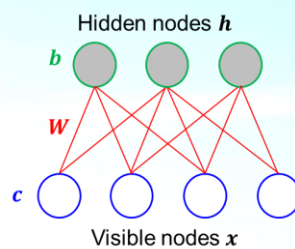
16.3.3 Learning of RBM

Xuegong Zhang

37

The Learning Problem

- Knowing the model and data, we can sample the hidden states.
- Giving the model and hidden states, we can sample the data.
- **Giving the data, how can we learn the model (weights, offsets)?**
 - Q: What is the basic idea for parametric learning?



Xuegong Zhang

38



Learning the parameters of RBM

$$p(\mathbf{x}, \mathbf{h}) = e^{-E(\mathbf{x}, \mathbf{h})} / Z = \frac{1}{Z} \prod_j \prod_k \exp(W_{j,k} h_j x_k) \prod_k \exp(c_k x_k) \prod_j \exp(b_j h_j)$$

Likelihood: $P(\mathbf{x}; \boldsymbol{\theta}) = \sum_{\mathbf{h}} p(\mathbf{x}, \mathbf{h}) = \sum_{\mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})} / Z$

The general idea of EM algorithms

- A statistical model for X with unobserved latent data or missing values Z and a vector of unknown parameters $\boldsymbol{\theta}$.

$$p(X, Z; \boldsymbol{\theta})$$

- Likelihood of the observed data

$$L(\boldsymbol{\theta}; X) = p(X|\boldsymbol{\theta}) = \int p(X, Z|\boldsymbol{\theta}) dZ$$

- However, we don't know Z and enumerating all possibilities is often infeasible.
- The idea is to make the estimation in two iterative steps:
 - Get the expected value of the log likelihood of $\boldsymbol{\theta}$ based on some estimate of Z given X with the current estimation $\boldsymbol{\theta}^{(t)}$ Expectation (E-step)
 - Find the next estimation $\boldsymbol{\theta}^{(t+1)}$ that maximizes this expected log likelihood Maximization (M-step)



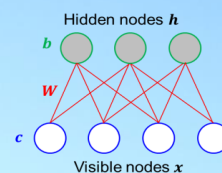
39

Learning Procedures for RBM

- MLE (Maximum Likelihood Estimation):

$$\max_{\boldsymbol{\theta}} P(\mathbf{x}; \boldsymbol{\theta})$$

$$P(\mathbf{x}; \boldsymbol{\theta}) = \sum_{\mathbf{h}} p(\mathbf{x}, \mathbf{h}) = \sum_{\mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})} / Z, \quad Z = \sum_{\mathbf{x}, \mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})}$$



- Minimum Loss: $\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}}(-\log P(\mathbf{x}))$

$$\frac{\partial(-\log P(\mathbf{x}))}{\partial \boldsymbol{\theta}} = \frac{\partial(-\log \sum_{\mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})} + \log \sum_{\mathbf{x}, \mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})})}{\partial \boldsymbol{\theta}}$$

40



Learning Procedures for RBM

- Minimum Loss: $\theta^* = \underset{\theta}{\operatorname{argmin}}(-\log P(x))$
 $P(x) = \sum_{\mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})} / \sum_{\mathbf{x}, \mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})}$
 - Gradient descent: Updating weights along the **negative** direction of the gradient $\frac{\partial(-\log P(x))}{\partial \theta}$
- $$\frac{\partial(-\log P(x))}{\partial \theta} = \frac{\partial(-\log \sum_{\mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})} + \log \sum_{\mathbf{x}, \mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})})}{\partial \theta}$$
- $$= \frac{1}{\sum_{\mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})}} \sum_{\mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})} \frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta} - \frac{1}{\sum_{\mathbf{x}, \mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})}} \sum_{\mathbf{x}, \mathbf{h}} e^{-E(\mathbf{x}, \mathbf{h})} \frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta} = \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{x}) \frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta} - \sum_{\mathbf{x}, \mathbf{h}} p(\mathbf{x}, \mathbf{h}) \frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta}$$
- i.e. $\frac{\partial(-\log P(x))}{\partial \theta} = \mathbb{E}_{P(\mathbf{h}|\mathbf{x})} \frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta} - \mathbb{E}_{P(\mathbf{x}, \mathbf{h})} \frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta}$
- To calculate the gradient, we need to sample from $P(\mathbf{h}|\mathbf{x})$ and $P(\mathbf{x}, \mathbf{h})$.

Intractable

41

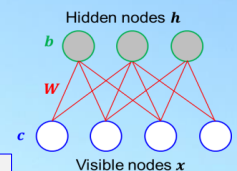
The Gradients



$$\frac{\partial(-\log P(x))}{\partial \theta} = \mathbb{E}_{P(\mathbf{h}|\mathbf{x})} \frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta} - \mathbb{E}_{P(\mathbf{x}, \mathbf{h})} \frac{\partial E(\mathbf{x}, \mathbf{h})}{\partial \theta}$$

Expectation over observed data
Negative of the direction:
energy ↓, likelihood ↑

Expectation over all values (random data)
Negative of the negative direction:
energy ↑, likelihood ↓



$$E(\mathbf{x}, \mathbf{h}) = -\mathbf{h}^T \mathbf{W} \mathbf{x} - \mathbf{c}^T \mathbf{x} - \mathbf{b}^T \mathbf{h} = -\sum_j \sum_k \mathbf{W}_{j,k} h_j x_k - \sum_k \mathbf{c}_k x_k - \sum_j \mathbf{b}_j h_j$$

$$\frac{\partial(-\log P(x))}{\partial \mathbf{W}} = -\mathbb{E}_{\mathbf{h}}(\mathbf{h}(\mathbf{x}(t))^T \mathbf{x}(t)) + \mathbb{E}_{\mathbf{x}, \mathbf{h}} \mathbf{h}^T \mathbf{x}$$

$$\frac{\partial(-\log P(x))}{\partial \mathbf{c}} = -\mathbf{x}(t) + \mathbb{E}_{\mathbf{x}, \mathbf{h}} \mathbf{x}$$

$$\frac{\partial(-\log P(x))}{\partial \mathbf{b}} = -\mathbb{E}_{\mathbf{h}} \mathbf{h}(\mathbf{x}(t)) + \mathbb{E}_{\mathbf{x}, \mathbf{h}} \mathbf{h}$$

$\mathbb{E}_{P(\mathbf{x}, \mathbf{h})}$

Intractable partition function

$\mathbb{E}_{P(\mathbf{h}|\mathbf{x})}$

$$p(\mathbf{h}|\mathbf{x}) = \prod_j \operatorname{sigm}((2h_j - 1)(b_j + \mathbf{W}_{j \cdot} \mathbf{x}))$$

Xuegong Zhang

42

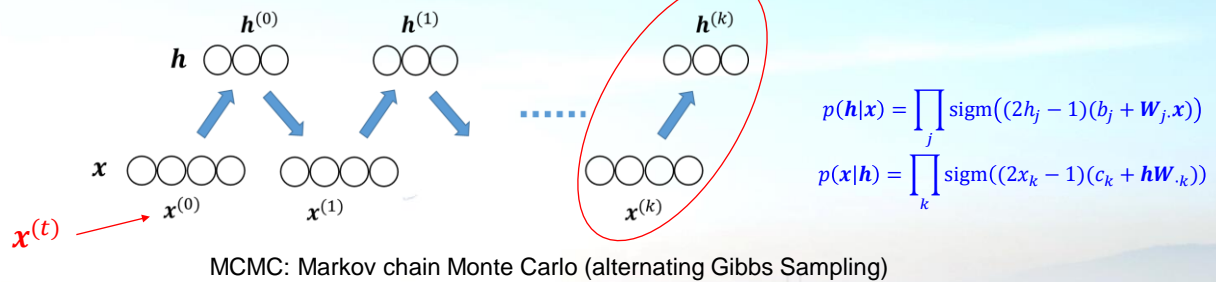


The Contrastive Divergence (CD) Algorithm

$$\frac{\partial(-\log P(x))}{\partial \theta} = \mathbb{E}_{P(h|x)} \frac{\partial E(x, h)}{\partial \theta} - \underbrace{\mathbb{E}_{P(x, h)}}_{\text{Estimation of the expectation } \mathbb{E}_{P(x, h)}()}$$

Idea for solution

- Replace the expectation by a **point estimate** at \tilde{x}
- Obtain the point \tilde{x} by Gibbs sampling
- Start sampling chain at a sample $x^{(t)}$



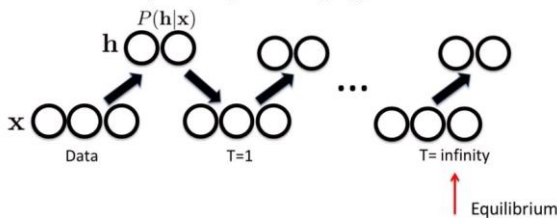
Xuegong Zhang

43

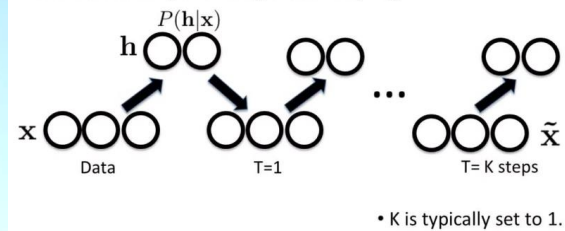


The Contrastive Divergence (CD) Algorithm

Run Markov chain (alternating Gibbs Sampling):



Run Markov chain (alternating Gibbs Sampling):



$$\frac{\partial(-\log P(x))}{\partial \mathbf{W}} = -\mathbb{E}_h(\mathbf{h}(x(t))^T \mathbf{x}(t)) + \underbrace{\mathbb{E}_{x, h} \mathbf{h}^T \mathbf{x}}_{\text{fake sample \& hidden state: weaken their connections}}$$

$$\frac{\partial(-\log P(x))}{\partial \mathbf{c}} = -\mathbf{x}(t) + \underbrace{\mathbb{E}_{x, h} \mathbf{x}}_{\text{true sample \& most probable hidden state: strengthen their connections}}$$

$$\frac{\partial(-\log P(x))}{\partial \mathbf{b}} = -\mathbb{E}_h \mathbf{h}(x(t)) + \underbrace{\mathbb{E}_{x, h} \mathbf{h}}_{\text{true sample \& most probable hidden state: strengthen their connections}}$$

$$\Delta \mathbf{W} = -\frac{\partial(-\log P(x))}{\partial \mathbf{W}} = \mathbf{h}(x(t))^T \mathbf{x}(t) - \tilde{\mathbf{h}}^T \tilde{\mathbf{x}}$$

$$\Delta \mathbf{c} = \mathbf{x}(t) - \tilde{\mathbf{x}}$$

$$\Delta \mathbf{b} = \mathbf{h}(x(t)) - \tilde{\mathbf{h}}$$

fake sample & hidden state: **weaken their connections**

true sample & most probable hidden state: **strengthen their connections**

Xuegong Zhang

44



CD-k Algorithm

- For each training example $\mathbf{x}^{(t)}$
 - Generate a **negative sample** $\tilde{\mathbf{x}}$ using k steps of Gibbs sampling, starting at the data point $\mathbf{x}^{(t)}$
 - Update model parameters:

$$\mathbf{W} \leftarrow \mathbf{W} + \alpha \left(\mathbf{h}(\mathbf{x}^{(t)}) \mathbf{x}^{(t)\top} - \mathbf{h}(\tilde{\mathbf{x}}) \tilde{\mathbf{x}}^\top \right)$$

$$\mathbf{b} \leftarrow \mathbf{b} + \alpha \left(\mathbf{h}(\mathbf{x}^{(t)}) - \mathbf{h}(\tilde{\mathbf{x}}) \right)$$

$$\mathbf{c} \leftarrow \mathbf{c} + \alpha \left(\mathbf{x}^{(t)} - \tilde{\mathbf{x}} \right)$$
 - Go back to 1 until **stopping criteria**

Xuegong Zhang

45

CD-1 Algorithm



Set learning rate η

Parameters to be learned: weights \mathbf{W} , offset vectors \mathbf{b} and \mathbf{c}

(1) Initialization

(2) Input a training example \mathbf{x}_0

Sampling:

- For all hidden nodes, calculate $P(h_j = 1 | \mathbf{x}_0) = \text{sigm}(b_j + \mathbf{W}_{j \cdot} \mathbf{x}_0)$, sample \mathbf{h}_0
- For all visible nodes, calculate $P(x_k = 1 | \mathbf{h}_0) = \text{sigm}(c_k + \mathbf{h}_0 \mathbf{W}_{\cdot k})$, sample \mathbf{x}_1
- For all hidden nodes, calculate $P(h_j = 1 | \mathbf{x}_1) = \text{sigm}(b_j + \mathbf{W}_{j \cdot} \mathbf{x}_1)$, sample \mathbf{h}_1

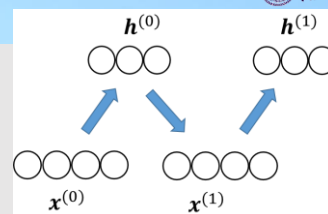
Updating:

$$\mathbf{W} \leftarrow \mathbf{W} + \eta (\mathbf{h}_0^\top \mathbf{x}_0 - \mathbf{h}_1^\top \mathbf{x}_1)$$

$$\mathbf{c} \leftarrow \mathbf{c} + \eta (\mathbf{x}_0 - \mathbf{x}_1)$$

$$\mathbf{b} \leftarrow \mathbf{b} + \eta (\mathbf{h}_0 - \mathbf{h}_1)$$

(3) Check for convergence: If converged or reached pre-set training rounds, stop; Otherwise go to (2).



Xuegong Zhang

46



10-second break



Xuegong Zhang

47



16.3.4 Using RBM for Recommendation

Xuegong Zhang

48



Hidden Nodes

Visible Nodes

Movie 1 Movie 2 Movie 3 Movie 4 Movie 5 Movie 6

Genre A Genre B Actor X Award Y Director Z

Movie 1 Movie 2 Movie 3 Movie 4 Movie 5 Movie 6

	Movie 1	Movie 2	Movie 3	Movie 4	Movie 5	Movie 6
User 1	1	0		1	1	1
User 2	0	1	0	0	1	0
User 3		1	1	0	0	
User 4	1	0	1	1	0	1
User 5	0		1	1		1
User 6	0	0	0	0	1	
User 7	1	0	1	1	0	1
User 8	0	1	1		0	1
User 9		0	1	1	1	1
User 10	1		0	0		0
User 11	0	1	1	1	0	1

Training

Genre A Genre B Actor X Award

Movie 1 Movie 2 Movie 3 Movie 4 Movie 5 Movie 6

	Movie 1	Movie 2	Movie 3	Movie 4	Movie 5	Movie 6
User 1	1	0		1	1	1
User 2	0	1	0	0	1	0
User 3		1	1	0	0	
User 4	1	0	1	1	0	1
User 5	0		1	1		1
User 6	0	0	0	0	1	
User 7	1	0	1	1	0	1
User 8	0	1	1		0	1
User 9		0	1	1	1	1
User 10	1		0	0		0
User 11	0	1	1	1	0	1

Xuegong Zhang

Slides from: Deep Learning A-Z_Hands-On Artificial Neural Networks Restricted Boltzmann Machine by SuperDataScience

49



0 1 0 1

The Matrix Fight Club Forrest Gump Pulp Fiction Titanic The Departed

Prediction

0 0 1 0 1 1

The Matrix Fight Club Forrest Gump Pulp Fiction Titanic The Departed

Xuegong Zhang

Slides from: Deep Learning A-Z_Hands-On Artificial Neural Networks Restricted Boltzmann Machine by SuperDataScience

50



16.3.5 RBM with Different Data Modalities

Xuegong Zhang

51

Gaussian Bernoulli RBMs

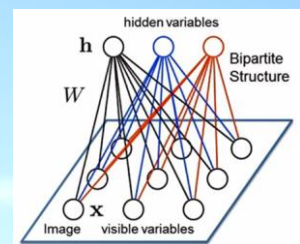


- Real-valued visible variables x

$$E(x, h) = -h^T W x - c^T x - b^T h + \frac{1}{2} x^T x$$

- $p(x|h)$ becomes a Gaussian distribution:

$$p(x|h) \sim N(c + W^T h, I)$$

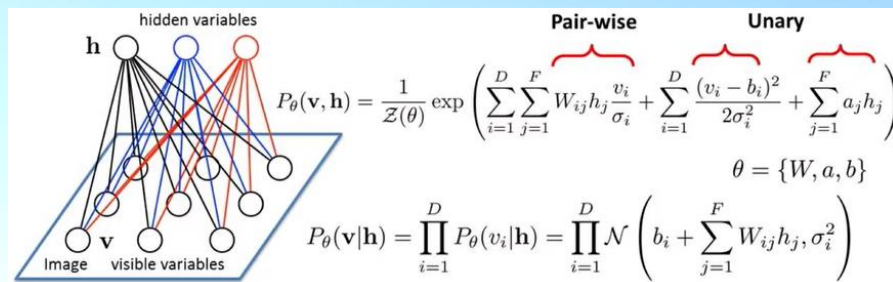


Xuegong Zhang

Slides from: Russ Salakhutdinov, Course 10707 at CMU ML Department
<http://www.cs.cmu.edu/~rsalakhu/10707>

52

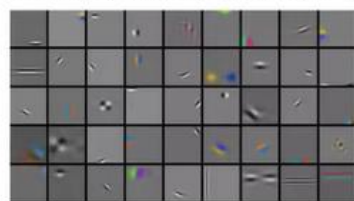
Gaussian Bernoulli RBMs



4 million unlabelled images



Learned features (out of 10,000)



$$\text{New Image} = p(h_7 = 1|v) = 0.9 * \text{feature}_1 + 0.8 * \text{feature}_2 + 0.6 * \text{feature}_3 + \dots$$

Xuegong Zhang

Department

53

Gaussian Bernoulli RBMs

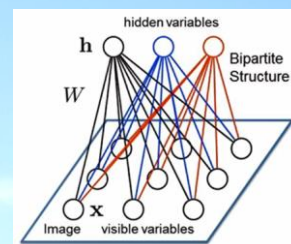


- Interpretation:
Mixture of exponential number of Gaussians

$$p(x) = \sum_h p(x|h)P(h)$$

where $P(h) = \int_x p(x, h) dx$

$$\text{and } p(x_i|h) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(x_i - b_i - \sigma_i \sum_j W_{ij} h_j)^2}{2\sigma_i^2}\right)$$



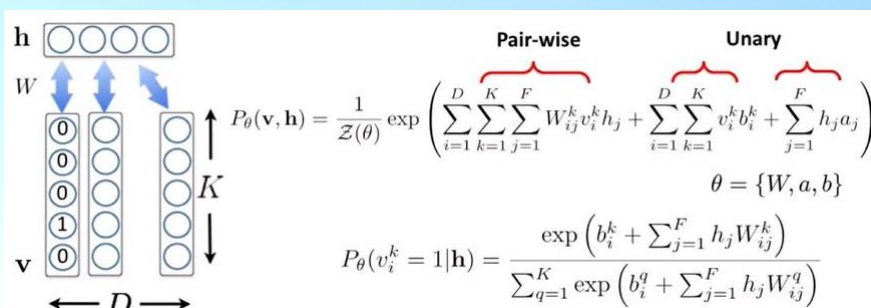
Xuegong Zhang

Slides from: Russ Salakhutdinov, Course 10707 at CMU ML Department
<http://www.cs.cmu.edu/~rsalakhu/10707>

54



RBM for Word Counts



- Replicated Softmax Model: undirected topic model:
 - Stochastic 1-of-K visible variables
 - Stochastic binary hidden variables

Salakhutdinov & Hinton, NIPS 2010

Xuegong Zhang

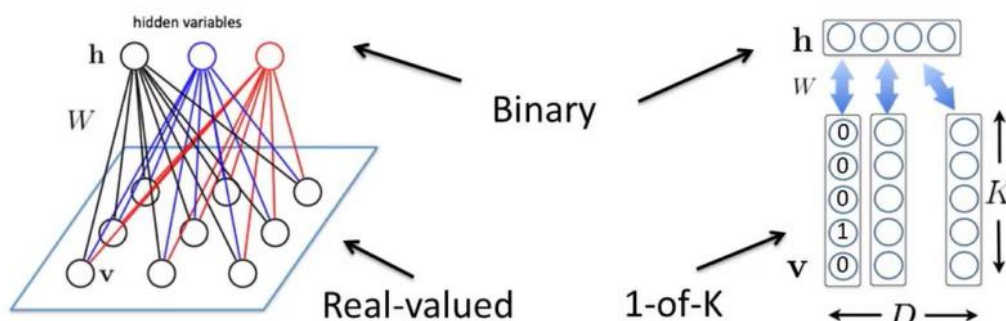
Slides from: Russ Salakhutdinov, Course 10707 at CMU ML Department
<http://www.cs.cmu.edu/~rsalakhu/10707>

55

Different Data Modalities



- Binary/Gaussian/Softmax RBMs: All have binary hidden variables but use them to model different kinds of data.



Xuegong Zhang

Slides from: Russ Salakhutdinov, Course 10707 at CMU ML Department
<http://www.cs.cmu.edu/~rsalakhu/10707>

56

Building a Probabilistic Model

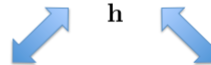
- Learn a joint density model:

$$P(\mathbf{h}, \mathbf{v}_{\text{image}}, \mathbf{v}_{\text{text}}).$$

- \mathbf{h} : “fused” representation for classification, retrieval.

$$P(\mathbf{h}|\mathbf{v}_{\text{image}}, \mathbf{v}_{\text{text}})$$

“Concept”



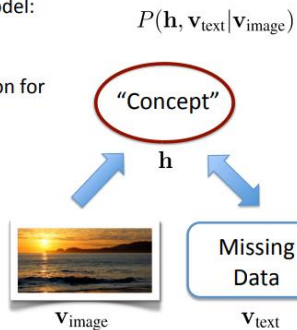
- Learn a joint density model:

$$P(\mathbf{h}, \mathbf{v}_{\text{image}}, \mathbf{v}_{\text{text}}).$$

- \mathbf{h} : “fused” representation for classification, retrieval.

- Generate data from conditional distributions for

- Image Annotation

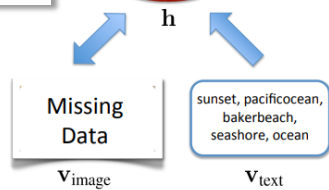


- Generate data from conditional distributions for

- Image Annotation
- Image Retrieval

$$P(\mathbf{h}, \mathbf{v}_{\text{image}}|\mathbf{v}_{\text{text}})$$

“Concept”



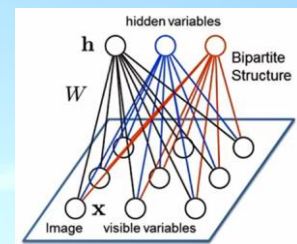
Xuegong Zhang

Slides from: Russ Salakhutdinov, Course 10707 at CMU ML Department
<http://www.cs.cmu.edu/~rsalakhu/10707>

57

Multiple Application Domains

- Natural Images
- Texts/Documents
- Recommendation/Collaborative filtering
- Videos
- Motion Capture
- Speech Perception
- ...



Xuegong Zhang

Slides from: Russ Salakhutdinov, Course 10707 at CMU ML Department
<http://www.cs.cmu.edu/~rsalakhu/10707>

58

Homework



- Problem sets and Computer exercises (PrEx8): choose at least one task among the following 4 optional tasks:
 - (optional 1) Write an essay on Hopfield Network with full mathematical developments.
 - (optional 2) Write an essay on RBM with real-valued and/or count variables with key mathematical developments.
 - (optional 3) Find a package of Hopfield Net and a demo task with real data, and do the computer exercise on the task.
 - (optional 4) Find a package of RBM and a demo task with real data, and do the computer exercise on the task.
- Deadline:
 - Dec. 15 (Wednesday), 23:00



Xuegong Zhang

59

单选题 1分

⚙ 设置

4-min break

- ☒ A Yes, I'm back to my computer.
- ☐ B Sorry, not yet.



Xuegong Zhang

提交

60