# **Deep Learning**

Kayhan Batmanghelich

**Slides: Matt Gormley** 

#### Motivation

# Why is everyone talking about Deep Learning?

- Because a lot of money is invested in it...
  - DeepMind: Acquired by Google for \$400
     million
  - DNNResearch: Three person startup (including Geoff Hinton) acquired by Google for unknown price tag
  - Ersatz, MetaMind, Nervana, Skylab:
     Deep Learning startups commanding millions
     of VC dollars
- Because it made the **front page** of the New York Times







#### Motivation

1980s

,1990s

2006

2016

# Why is everyone talking about Deep Learning?

#### C 1960s Deep learning:

- Has won numerous pattern recognition competitions
- Does so with minimal feature engineering

This wasn't always the case!

Since 1980s: Form of models hasn't changed much, but lots of new tricks...

- More hidden units
- Better (online) optimization
- New nonlinear functions (ReLUs)
- Faster computers (CPUs and GPUs)

#### Background

## A Recipe for Machine Learning

- 1. Given training data: $\{oldsymbol{x}_i,oldsymbol{y}_i\}_{i=1}^N$
- 2. Choose each of these:
  - Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

Loss function

 $\ell(\hat{oldsymbol{y}},oldsymbol{y}_i)\in\mathbb{R}$ 



**Examples:** Linear regression, Logistic regression, Neural Network

**Examples**: Mean-squared error, Cross Entropy

#### Background

## A Recipe for Machine Learning

- 1. Given training data: $\{oldsymbol{x}_i,oldsymbol{y}_i\}_{i=1}^N$
- 2. Choose each of these:
  - Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

Loss function

 $\ell(\hat{oldsymbol{y}},oldsymbol{y}_i)\in\mathbb{R}$ 

3. Define goal:  $\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^N \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$ 

4. Train with SGD:(take small stepsopposite the gradient)

 $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$ 



### A Recipe for Gradients

1. Given training dat $\{oldsymbol{x}_i,oldsymbol{y}_i\}_{i=1}^N$ 

#### 2. Choose each of t

Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

Loss function

 $\ell(\hat{oldsymbol{y}},oldsymbol{y}_i)\in\mathbb{R}$ 

**Backpropagation** can compute this gradient!

And it's a **special case of a more general algorithm** called reversemode automatic differentiation that can compute the gradient of any differentiable function efficiently!

> opposite the gradient)  $\boldsymbol{\theta}^{(t)} = -\eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$

## A Recipe for

## Goals for Today's Lecture

- 1. Explore a new class of decision functions (Deep Nets)
  - 2. Consider variants of this recipe for training

#### 2. choose each or these:

Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

Loss function

 $\ell(\hat{oldsymbol{y}},oldsymbol{y}_i)\in\mathbb{R}$ 

Train with SGD:
Ike small steps
opposite the gradient)

 $oldsymbol{ heta}^{(t+1)} = oldsymbol{ heta}^{(t)} - \eta_t 
abla \ell(f_{oldsymbol{ heta}}(oldsymbol{x}_i),oldsymbol{y}_i)$ 

# Outline

- Motivation
- Deep Neural Networks (DNNs)
  - Background: Decision functions
  - Background: Neural Networks
  - Three ideas for training a DNN
  - Experiments: MNIST digit classification

#### Deep Belief Networks (DBNs)

- Sigmoid Belief Network
- Contrastive Divergence learning
- Restricted Boltzman Machines (RBMs)
- RBMs as infinitely deep Sigmoid Belief Nets
- Learning DBNs
- Deep Boltzman Machines (DBMs)
  - Boltzman Machines
  - Learning Boltzman Machines
  - Learning DBMs

# Outline

Motivation

#### Deep Neural Networks (DNNs)

- Background: Decision functions
- Background: Neural Networks
- Three ideas for training a DNN
- Experiments: MNIST digit classification
- Deep Belief Networks (DBNs)
  - Sigmoid Belief Network
  - Contrastive Divergence learning
  - Restricted Boltzman Machines (RBMs)
  - RBMs as infinitely deep Sigmoid Belief Nets
  - Learning DBNs
- Deep Boltzman Machines (DBMs)
  - Boltzman Machines
  - Learning Boltzman Machines
  - Learning DBMs













## **Multi-Class Output**



## **Deeper Networks**

This lecture:



## **Deeper Networks**

#### This lecture:



## **Deeper Networks**



## Different Levels of Abstraction

- We don't know the "right" levels of abstraction
- So let the model figure it out!





3rd layer "Objects"





Pixels

#### Example from Honglak Lee (NIPS 2010)

## Different Levels of Abstraction

### Face Recognition:

- Deep Network
   can build up
   increasingly
   higher levels of
   abstraction
- Lines, parts, regions



3rd layer "Objects"

2nd layer "Object parts"



Pixels

Example from Honglak Lee (NIPS 2010)

## Different Levels of Abstraction





2nd layer "Object parts"

1st layer "Edges"

Pixels

#### Example from Honglak Lee (NIPS 2010)

## A Recipe for

## Goals for Today's Lecture

- Explore a new class of decision functions (Deep Neural Networks)
  - 2. Consider variants of this recipe for training

#### 2. Choose each of these:

Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

Loss function



4. Train with SGD:(take small steps opposite the gradient)

 $oldsymbol{ heta}^{(t+1)} = oldsymbol{ heta}^{(t)} - \eta_t 
abla \ell(f_{oldsymbol{ heta}}(oldsymbol{x}_i),oldsymbol{y}_i)$ 

## Idea #1: No pre-training

- Idea #1: (Just like a shallow network)
  - Compute the supervised gradient by backpropagation.
  - Take small steps in the direction of the gradient (SGD)

# Backpropagation

Backpropagation

is just repeated application of the **chain rule** from Calculus 101.

y = g(u) and u = h(x).

**Chain Rule:**  $dy_i \ du_j$  $dy_i$  $\forall i, k$  $\overline{du_j} \, \overline{dx_k}$  $dx_k$ 





# Backpropagation

Case 2: Forward Neural Network  $J = y^* \log q$  $q = \frac{1}{1+q}$  $b = \sum_{i=1}^{D} distance descent for the second sec$ 

$$J = y^* \log q + (1 - y^*) \log(1 - q)$$
$$q = \frac{1}{1 + \exp(-b)}$$
$$b = \sum_{j=0}^{D} \beta_j z_j$$
$$z_j = \frac{1}{1 + \exp(-a_j)}$$
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i$$

Backward  

$$\frac{dJ}{dq} = \frac{y^*}{q} + \frac{(1-y^*)}{q-1}$$

$$\frac{dJ}{db} = \frac{dJ}{dy}\frac{dy}{db}, \frac{dy}{db} = \frac{\exp(b)}{(\exp(b)+1)^2}$$

$$\frac{dJ}{d\beta_j} = \frac{dJ}{db}\frac{db}{d\beta_j}, \frac{db}{d\beta_j} = z_j$$

$$\frac{dJ}{dz_j} = \frac{dJ}{db}\frac{db}{dz_j}, \frac{db}{dz_j} = \beta_j$$

$$\frac{dJ}{da_j} = \frac{dJ}{dz_j}\frac{dz_j}{da_j}, \frac{dz_j}{da_j} = \frac{\exp(a_j)}{(\exp(a_j)+1)^2}$$

$$\frac{dJ}{d\alpha_{ji}} = \frac{dJ}{da_j}\frac{da_j}{d\alpha_{ji}}, \frac{da_j}{d\alpha_{ji}} = x_i$$

$$\frac{dJ}{dx_i} = \frac{dJ}{da_j}\frac{da_j}{dx_i}, \frac{da_j}{dx_i} = \sum_{j=0}^D \alpha_{ji}$$
30

## Idea #1: No pre-training

- Idea #1: (Just like a shallow network)
  - Compute the supervised gradient by backpropagation.
  - Take small steps in the direction of the gradient (SGD)

## **Comparison on MNIST**

- Results from Bengio et al. (2006) on MNIST digit classification task
- Percent error (lower is better)



## **Comparison on MNIST**

- Results from Bengio et al. (2006) on MNIST digit classification task
- Percent error (lower is better)



## Idea #1: No pre-training

- Idea #1: (Just like a shallow network)
  - Compute the supervised gradient by backpropagation.
  - Take small steps in the direction of the gradient (SGD)
- What goes wrong?
  - A. Gets stuck in local optima
    - Nonconvex objective
    - Usually start at a random (bad) point in parameter space
  - B. Gradient is progressively getting more dilute
    - "Vanishing gradients"

## Problem A: Nonconvexity

- Where does the nonconvexity come from?
- Even a simple quadratic z = xy objective is nonconvex:



## Problem A: Nonconvexity

- Where does the nonconvexity come from?
- Even a simple quadratic z = xy objective is nonconvex:



## Problem A: Nonconvexity

Stochastic Gradient Descent...

... climbs to the top of the nearest hill...



## Problem A: Nonconvexity

Stochastic Gradient Descent...

... climbs to the top of the nearest hill...


### Problem A: Nonconvexity

Stochastic Gradient Descent...

... climbs to the top of the nearest hill...



### Problem A: Nonconvexity

Stochastic Gradient Descent...

... climbs to the top of the nearest hill...



### Problem A: Nonconvexity

Stochastic Gradient Descent...

... climbs to the top of the nearest hill...

... which might not lead to the top of the mountain



# Problem B: Vanishing Gradients

The gradient for an edge dat the base of the network depends on the gradients of many edges above it

The chain rule multiplies many of these partial derivatives together



# Problem B: Vanishing Gradients

The gradient for an edge dat the base of the network depends on the find gradients of many edges above it

The chain rule multiplies many of these partial derivatives together



# Problem B: Vanishing Gradients

The gradient for an edge at the base of the network depends on the gradients of many edges above it

The chain rule multiplies many of these partial derivatives together



# Idea #2: Supervised Pre-training

- Idea #2: (Two Steps)
  - Train each level of the model in a greedy way
  - Then use our original idea
- 1. Supervised Pre-training
  - Use labeled data
  - Work bottom-up
    - Train hidden layer 1. Then fix its parameters.
    - Train hidden layer 2. Then fix its parameters.
    - ...
    - Train hidden layer n. Then fix its parameters.
- 2. Supervised Fine-tuning
  - Use labeled data to train following "Idea #1"
  - Refine the features by backpropagation so that they become tuned to the end-task

# Idea #2: Supervised Pre-training

- Idea #2: (Two Steps)
  - Train each level of the model in a greedy way
  - Then use our **original idea**



# Idea #2: Supervised Pre-training

#### Idea #2: (Two Steps)

- Train each level of the model in a greedy way
- Then use our original idea



# Idea #2: Supervised Pre-training

#### Idea #2: (Two Steps) У Hidden Layer 3 **C**<sub>1</sub> **C**<sub>2</sub> CF b₁ b<sub>E</sub> Hidden Layer 2 b<sub>2</sub> Hidden Layer 1 a<sub>1</sub> **a**<sub>2</sub> ••• a<sub>D</sub> ••• **X**2 X<sub>1</sub> X<sub>3</sub> X<sub>M</sub> Input

Training

# Idea #2: Supervised Pre-training





# **Comparison on MNIST**

- Results from Bengio et al. (2006) on MNIST digit classification task
- Percent error (lower is better)



# **Comparison on MNIST**

- Results from Bengio et al. (2006) on MNIST digit classification task
- Percent error (lower is better)



# Idea #3: Unsupervised Pre-training

- Idea #3: (Two Steps)
  - Use our original idea, but **pick a better starting point**
  - Train each level of the model in a greedy way
- 1. Unsupervised Pre-training
  - Use unlabeled data
  - Work bottom-up
    - Train hidden layer 1. Then fix its parameters.
    - Train hidden layer 2. Then fix its parameters.
    - .
    - Train hidden layer n. Then fix its parameters.
- 2. Supervised Fine-tuning
  - Use labeled data to train following "Idea #1"
  - Refine the features by backpropagation so that they become tuned to the end-task

### Unsupervised pretraining of the first layer:

- What should it predict?
- What else do we observe?
- The input!



### Unsupervised pretraining of the first layer:

- What should it predict?
- What else do we observe?
- The input!

# This topology defines an Auto-encoder.



# Auto-Encoders

Key idea: Encourage z to give small reconstruction error:

- x' is the reconstruction of x
- Loss =  $|| x DECODER(ENCODER(x)) ||^2$
- Train with the same backpropagation algorithm for 2-layer Neural Networks with  $x_{\rm m}$  as both input and output.

```
DECODER: x' = h(W'z)
```

ENCODER: z = h(Wx)



Slide adapted from Raman Arora

### Unsupervised pretraining

- Work bottom-up
  - Train hidden layer 1.
     Then fix its parameters.
  - Train hidden layer 2.
     Then fix its parameters.
  - ...
  - Train hidden layer n.
     Then fix its parameters.



### Unsupervised pretraining

- Work bottom-up
  - Train hidden layer 1.
     Then fix its parameters.
  - Train hidden layer 2.
     Then fix its parameters.
  - Train hidden layer n.
     Then fix its parameters.



### Unsupervised pretraining

- Work bottom-up
  - Train hidden layer 1.
     Then fix its parameters.
  - Train hidden layer 2.
     Then fix its parameters.
  - ...
     Train hidden layer n.
     Then fix its parameters.



### Unsupervised pretraining

- Work bottom-up
  - Train hidden layer 1.
     Then fix its parameters.
  - Train hidden layer 2. Hidden La Then fix its parameters.
    - ...
  - Train hidden layer n. Hidden Then fix its parameters.

#### Supervised fine-tuning Backprop and update all

parameters



# Deep Network Training

#### • Idea #1:

1. Supervised fine-tuning only

#### • Idea #2:

- 1. Supervised layer-wise pre-training
- 2. Supervised fine-tuning

#### • Idea #3:

- 1. Unsupervised layer-wise pre-training
- 2. Supervised fine-tuning

# **Comparison on MNIST**

- Results from Bengio et al. (2006) on MNIST digit classification task
- Percent error (lower is better)



# **Comparison on MNIST**

- Results from Bengio et al. (2006) on MNIST digit classification task
- Percent error (lower is better)



Is layer-wise pre-training always necessary?

**In 2010,** a record on a hand-writing recognition task was set by standard supervised backpropagation (our Idea #1).

**HOW?** A very fast implementation on GPUs.

See Ciresen et al. (2010)

# Deep Learning

- Goal: learn features at different levels of abstraction
- Training can be tricky due to...
  - Nonconvexity
  - Vanishing gradients
- Unsupervised layer-wise pre-training can help with both!

# Outline

- Motivation
- Deep Neural Networks (DNNs)
  - Background: Decision functions
  - Background: Neural Networks
  - Three ideas for training a DNN
  - Experiments: MNIST digit classification

#### Deep Belief Networks (DBNs)

- Sigmoid Belief Network
- Contrastive Divergence learning
- Restricted Boltzman Machines (RBMs)
- RBMs as infinitely deep Sigmoid Belief Nets
- Learning DBNs
- Deep Boltzman Machines (DBMs)
  - Boltzman Machines
  - Learning Boltzman Machines
  - Learning DBMs

Question:

# How does this relate to Graphical Models?

The first "Deep Learning" papers in 2006 were innovations in training a particular flavor of Belief Network.

Those models happen to also be neural nets.

### DBNs

# MNIST Digit Generation

- This section: Suppose you want to build a generative model capable of explaining handwritten digits
- Goal:
  - To have a model p(x)
     from which we can
     sample digits that look
     realistic
  - Learn **unsupervised** hidden representation of an image

0	0	0	0	0	0	0	0	0	9
1	1	1	1	1	1	1	1	1	l
2	2	2	2	ī	2	2	2	2	2
3	3	3	3	٦	3	3	3	ъ	3
ч	4	4	4	4	Ч	ч	4	4	1
5	5	5	5	5	5	5	5	5	3
6	Ь	6	6	G	6	6	6	6	6
1	7	7	7	7	7	7	7	7	7
8	8	в	8	8	8	8	8	ર્ટ	е
4	9	9	5	9	9	9	9	9	9

### DBNs

# Sigmoid Belief Networks

- Directed graphical model of binary variables in fully connected layers
- Only bottom layer is observed
- Specific parameterization of the conditional probabilities:

$$p(x_i | \text{parents}(x_i)) = \frac{1}{1 + \exp(-\sum_j w_{ij} x_j)}$$
Unknown Params



# A bit of (relevant) digression: Contrastive Divergence



# Contrastive Divergence Training

**Contrastive Divergence** is a general tool for learning a generative distribution, where the derivative of the log partition function is intractable to compute.

Max likelihood principle to train the model:



### DBNs

Fill in

# Contrastive Divergence Training



### DBNs

# Contrastive Divergence Training

**Contrastive Divergence** is a general tool for learning a generative distribution, where the derivative of the log partition function is intractable to compute.

Max likelihood principle to train the model:



- **A hurdle**: many MCMC cycles required to compute the second term.
- Hinton et al. assert that only a few MCMC cycles would be needed to calculate an approximate gradient.

# Contrastive Divergence Training



Slide from Marcus Frean, MLSS Tutorial 2010

**DBNs** 

# Contrastive Divergence Training

#### Another view:

**DBNs** 

$$\frac{\partial}{\partial w} \log Z = \frac{1}{Z} \frac{\partial}{\partial w} \sum_{\mathbf{v}} \sum_{\mathbf{h}} P^{\star}(\mathbf{v}, \mathbf{h})$$

$$= \frac{1}{Z} \sum_{\mathbf{v}} \sum_{\mathbf{h}} \frac{\partial}{\partial w} P^{\star}(\mathbf{v}, \mathbf{h})$$

$$= \frac{1}{Z} \sum_{\mathbf{v}} \sum_{\mathbf{h}} P^{\star}(\mathbf{v}, \mathbf{h}) \frac{\partial}{\partial w} \log P^{\star}(\mathbf{v}, \mathbf{h})$$

$$= \sum_{\mathbf{v}} \sum_{\mathbf{h}} P(\mathbf{v}, \mathbf{h}) \frac{\partial}{\partial w} \log P^{\star}(\mathbf{v}, \mathbf{h})$$
average over joint!

Slide from Marcus Frean, MLSS Tutorial 2010
## Contrastive Divergence Training

#### Another view:



Slide from Marcus Frean, MLSS Tutorial 2010

## Back to Sigmoid BN

## Contrastive Divergence Training

 $(x_i - p_i)x_j$ 

For a belief net the joint is automatically normalised: Z is a constant 1

- 2nd term is zero!
- for the weight  $w_{ij}$  from j into i, the gradient  $\frac{\partial \log L}{\partial w_{ij}}$
- stochastic gradient ascent:

$$\Delta w_{ij} \propto \underbrace{(x_i - p_i)x_j}_{\text{the "delta rule"}}$$

#### Xi

Pi

So this is a stochastic version of the EM algorithm, that you may have heard of. We iterate the following two steps:

#### E step: get samples from the posterior

M step: apply the learning rule that makes them more likely

## Sigmoid Belief Networks

- In practice, applying CD to a Deep Sigmoid Belief Nets fails
- Sampling from the posterior of many (deep) hidden layers doesn't approach the equilibrium distribution quickly enough
- Take home summary: Sigmoid BN are easy to sample from as a generative model, but hard to learn

# Note: this is a GM diagram not a NN!



### How about undirected models?

## **Boltzman Machines**

- Undirected graphical model of binary variables with pairwise potentials
- Parameterization of the potentials:

 $\psi_{ij}(x_i, x_j) =$ 



(In English: higher value of parameter  $W_{ij}$  leads to higher correlation between  $X_i$  and  $X_j$  on value 1)

 $\exp(x_i W_{ij})$ 

### Restricted Boltzman Machines

- Assume visible units are one layer, and hidden units are another.
- Throw out all the connections within each layer.



- $h_j \perp h_k \mid \mathbf{v}$
- the posterior  $P(\mathbf{h} \mid \mathbf{v})$  factors *c.f.* in a belief net, the *prior*  $P(\mathbf{h})$  factors

### DBNs Restricted Boltzman Machines Alternating Gibbs sampling Since none of the units within a layer are interconnected, we can do Gibbs sampling by updating the whole layer at a time.



### Restricted Boltzman Machines

#### learning in an RBM



then alternate between updating all the hidden units in parallel and updating all the visible units in parallel

$$\Delta w_{ij} = \eta \left[ \langle v_i h_j \rangle^0 - \langle v_i h_j \rangle^\infty \right]$$

#### restricted connectivity is trick #1:

it saves waiting for equilibrium in the clamped phase.

#### Slide from Marcus Frean, MLSS Tutorial 2010

### Restricted Boltzman Machines

#### trick # 2: curtail the Markov chain during learning



Repeat for all data:

- start with a training vector on the visible units
- update all the hidden units in parallel
- update all the visible units in parallel to get a "reconstruction"
- update the hidden units again

$$\Delta w_{ij} = \eta \left[ \langle v_i h_j \rangle^0 - \langle v_i h_j \rangle^1 \right]$$

This is not following the correct gradient, but works well in practice. Geoff Hinton calls it learning by "contrastive divergence".

#### Slide from Marcus Frean, MLSS Tutorial 2010

## Deep Belief Networks (DBNs)

RBMs are equivalent to infinitely deep belief networks



sampling from this is the same as sampling from the network on the right.



## Deep Belief Networks (DBNs)

#### RBMs are equivalent to infinitely deep belief networks



- So when we train an RBM, we're really training an  $\infty^{ly}$  deep sigmoid belief net!
- It's just that the weights of all layers are tied.

Slide from Marcus Frean, MLSS Tutorial 2010

### Let's apply it on MNIST

### Unsupervised Learning of DBNs

Setting A: DBN Autoencoder

- I. Pre-train a stack of RBMs in greedy layerwise fashion
- II. Unroll the RBMs to create an autoencoder (i.e. bottom-up and top-down weights are untied)
- III. Fine-tune the parameters using backpropagation

### Unsupervised Learning of DBNs

### Setting A: DBN Autoencoder

DBNs

- I. Pre-train a stack of RBMs in greedy layerwise fashion
- II. Unroll the RBMs to create an autoencoder (i.e. bottom-up and top-down weights are untied)
- III. Fine-tune the parameters using backpropagation



## Unsupervised Learning of DBNs

Setting A: DBN Autoencoder

DBNs

- I. Pre-train a stack of RBMs in greedy layerwise fashion
- II. Unroll the RBMs to create an autoencoder (i.e. bottom-up and top-down weights are untied)
- III. Fine-tune the parameters using backpropagation





### Unsupervised Learning of DBNs

Setting A: DBN Autoencoder

- Pre-train a stack of RBMs in greedy layerwise fashion
- II. Unroll the RBMs to create an autoencoder (i.e. bottom-up and top-down weights are untied)
- III. Fine-tune the parameters using backpropagation



### Supervised Learning of DBNs

Setting B: DBN classifier

- Pre-train a stack of RBMs in greedy layerwise fashion (unsupervised)
- II. Fine-tune the parameters using backpropagation by minimizing classification error on the training data

## DBNs MNIST Digit Generation



- Comparison of deep autoencoder, logistic PCA, and PCA
- Each method projects the real data down to a vector of 30 real numbers
- Then reconstructs the data from the low-dimensional projection

## MNIST Digit Recognition

Experimental evaluation of DBN with greedy layerwise pretraining and fine-tuning via the wakesleep algorithm Examples of correctly recognized handwritten digits that the neural network had never seen before

## MNIST Digit Recognition

Experimental evaluation of DBN with greedy layerwise pretraining and fine-tuning via the wakesleep algorithm How well does it discriminate on MNIST test set with no extra information about geometric distortions?

- Generative model based on RBM's 1.25%
  Support Vector Machine (Decoste et. al.) 1.4%
  Backprop with 1000 hiddens (Platt) ~1.6%
  Backprop with 500 -->300 hiddens ~1.6%
  K-Nearest Neighbor ~3.3%
- See Le Cun et. al. 1998 for more results
- Its better than backprop and much more neurally plausible because the neurons only need to send one kind of signal, and the teacher can be another sensory input.

## Outline

- Motivation
- Deep Neural Networks (DNNs)
  - Background: Decision functions
  - Background: Neural Networks
  - Three ideas for training a DNN
  - Experiments: MNIST digit classification
- Deep Belief Networks (DBNs)
  - Sigmoid Belief Network
  - Contrastive Divergence learning
  - Restricted Boltzman Machines (RBMs)
  - RBMs as infinitely deep Sigmoid Belief Nets
  - Learning DBNs
- Deep Boltzman Machines (DBMs)
  - Boltzman Machines
  - Learning Boltzman Machines
  - Learning DBMs

### Deep Boltzman Machines

- DBNs are a hybrid directed/undi rected graphical model
- DBMs are a purely undirected graphical model



### Deep Boltzman Machines

# Can we use the same techniques to train a DBM?



## Learning Standard Boltzman Machines

- Undirected graphical model of binary variables with pairwise potentials
- Parameterization of the potentials:

$$\psi_{ij}(x_i, x_j) =$$

$$\exp(x_i W_{ij} x_j)$$

(In English: higher value of parameter W<sub>ij</sub> leads to higher correlation between X<sub>i</sub> and X<sub>j</sub> on value 1)



DBMs	Learning Standard Boltzman Machines
Visible units: $\mathbf{v} \in \{0\}$ Hidden units: $\mathbf{h} \in \{0\}$	$[0,1]^D$ $P(\mathbf{v},\mathbf{h}) = $
Likelihood: $E(\mathbf{v}, \mathbf{h}; \theta) = -\frac{1}{2} \mathbf{v}^{\top} \mathbf{L} \mathbf{v}$ $\theta_{p}(\mathbf{v}; \theta) = \frac{p^{*}(\mathbf{x}; \theta)}{p^{*}(\mathbf{x}; \theta)} \int 1$	$-\frac{1}{2}\mathbf{h}^{T}\mathbf{J}\mathbf{h} - \mathbf{v}^{T}\mathbf{W}\mathbf{h}$ $-\sum \exp\left(-E(\mathbf{v}, \mathbf{h}; \theta)\right),$
$Z(\theta) = \sum_{i=1}^{N} Z(\theta)$	$\sum_{\mathbf{v}} \sum_{\mathbf{h}} \exp\left(-E(\mathbf{v}, \mathbf{h}; \theta)\right),$
2	$\mathcal{P}(x_i, y_j) \ge 2x^2$ $\log \phi(.) = \mathcal{K}_i \mathcal{W}_i \mathcal{J}_j$ 10

### Learning Standard Boltzman Machines

(Old) idea from Hinton & Sejnowski (1983): For each iteration of optimization, run a separate MCMC chain for each of the data and model expectations to approximate the parameter updates.

Delta updates to each of model parameters:

$$\Delta \mathbf{W} = \alpha \left( \mathbf{E}_{P_{\text{data}}} [\mathbf{v} \mathbf{h}^{\top}] - \mathbf{E}_{P_{\text{model}}} [\mathbf{v} \mathbf{h}^{\top}] \right),$$
  

$$\Delta \mathbf{L} = \alpha \left( \mathbf{E}_{P_{\text{data}}} [\mathbf{v} \mathbf{v}^{\top}] - \mathbf{E}_{P_{\text{model}}} [\mathbf{v} \mathbf{v}^{\top}] \right),$$
  

$$\Delta \mathbf{J} = \alpha \left( \mathbf{E}_{P_{\text{data}}} [\mathbf{h} \mathbf{h}^{\top}] - \mathbf{E}_{P_{\text{model}}} [\mathbf{h} \mathbf{h}^{\top}] \right),$$

Full conditionals for Gibbs sampler:

DBMs

$$p(h_j = 1 | \mathbf{v}, \mathbf{h}_{-j}) = \sigma \Big( \sum_{i=1}^{D} W_{ij} v_i + \sum_{m=1 \setminus j}^{P} J_{jm} h_j \Big),$$
$$p(v_i = 1 | \mathbf{h}, \mathbf{v}_{-i}) = \sigma \Big( \sum_{j=1}^{P} W_{ij} h_j + \sum_{k=1 \setminus i}^{D} L_{ik} v_j \Big),$$

111

### Learning Standard Boltzman Machines



Full conditionals for Gibbs sampler:

DBMs

$$p(h_j = 1 | \mathbf{v}, \mathbf{h}_{-j}) = \sigma \Big( \sum_{i=1}^{D} W_{ij} v_i + \sum_{m=1 \setminus j}^{P} J_{jm} h_j \Big),$$
$$p(v_i = 1 | \mathbf{h}, \mathbf{v}_{-i}) = \sigma \Big( \sum_{j=1}^{P} W_{ij} h_j + \sum_{k=1 \setminus i}^{D} L_{ik} v_j \Big),$$

112

## Learning Standard Boltzman Machines

(New) idea from Salakhutinov & Hinton (2009):

- Step 1) Approximate the data distribution by variational inference.
- Step 2) Approximate the model distribution with a "persistent" Markov chain (from iteration to iteration)

Delta updates to each of model parameters:

$$\Delta \mathbf{W} = \alpha \left( \left\langle \mathbf{v} \mathbf{h}^T \right\rangle_{\mathbf{v} \in \mathcal{D}, \mathbf{h} \sim p(\mathbf{h} | \mathbf{v})} - \left\langle \mathbf{v} \mathbf{h}^T \right\rangle_{\mathbf{v}, \mathbf{h} \sim p(\mathbf{h}, \mathbf{v})} \right)$$
$$\Delta \mathbf{L} = \alpha \left( \left\langle \mathbf{v} \mathbf{v}^T \right\rangle_{\mathbf{v} \in \mathcal{D}, \mathbf{h} \sim p(\mathbf{h} | \mathbf{v})} - \left\langle \mathbf{v} \mathbf{v}^T \right\rangle_{\mathbf{v}, \mathbf{h} \sim p(\mathbf{h}, \mathbf{v})} \right)$$
$$\Delta \mathbf{J} = \alpha \left( \left\langle \mathbf{h} \mathbf{h}^T \right\rangle_{\mathbf{v} \in \mathcal{D}, \mathbf{h} \sim p(\mathbf{h} | \mathbf{v})} - \left\langle \mathbf{h} \mathbf{h}^T \right\rangle_{\mathbf{v}, \mathbf{h} \sim p(\mathbf{h}, \mathbf{v})} \right)$$



## Learning Standard Boltzman Machines

(New) idea from Salakhutinov & Hinton (2009):

- Step 1) Approximate the data distribution by variational inference.
- Step 2) Approximate the model distribution with a "persistent" Markov chain (from iteration to iteration)

Delta updates to each of model parameters:

$$\Delta \mathbf{W} = \alpha \left( \left\langle \mathbf{v} \mathbf{h}^T \right\rangle_{\mathbf{v} \in \mathcal{D}, \mathbf{h} \sim p(\mathbf{h} | \mathbf{v})} - \left\langle \mathbf{v} \mathbf{h}^T \right\rangle_{\mathbf{v}, \mathbf{h} \sim p(\mathbf{h}, \mathbf{v})} \right.$$

Step 1) Approximate the data distribution...

Mean-field approximation:

$$q(\mathbf{h};\mu) = \prod_{j=1}^{P} q(h_i)$$

 $q(h_i = 1) = \mu_i$ 

Variational lower-bound of log-likelihood:

 $\ln p(\mathbf{v}; \theta) \geq \sum_{\mathbf{h}} q(\mathbf{h} | \mathbf{v}; \mu) \ln p(\mathbf{v}, \mathbf{h}; \theta) + \mathcal{H}(q)$ 

Fixed-point equations for variational params:

$$\mu_j \leftarrow \sigma \big(\sum_i W_{ij} v_i + \sum_{m \setminus j} J_{mj} \mu_m\big)$$

### Learning Standard Boltzman Machines

(New) idea from Salakhutinov & Hinton (2009):

- Step 1) Approximate the data distribution by variational inference.
- Step 2) Approximate the model distribution with a "persistent" Markov chain (from iteration to iteration)

Delta updates to each of model parameters:

$$\Delta \mathbf{W} = \alpha \left( \left\langle \mathbf{v} \mathbf{h}^T \right\rangle_{\mathbf{v} \in \mathcal{D}, \mathbf{h} \sim p(\mathbf{h} | \mathbf{v})} - \left\langle \mathbf{v} \mathbf{h}^T \right\rangle_{\mathbf{v}, \mathbf{h} \sim p(\mathbf{h}, \mathbf{v})} \right.$$

Step 2) Approximate the model distribution...

Why not use variational inference for the model expectation as well?

Difference of the two mean-field approximated expectations above would cause learning algorithm to **maximize** divergence between true and mean-field distributions.

Persistent CD adds correlations between successive iterations, but not an issue.

### Deep Boltzman Machines

- DBNs are a hybrid directed/undi rected graphical model
- DBMs are a purely undirected graphical model



## Learning Deep Boltzman Machines

Can we use the same techniques to train a DBM?

- Pre-train a stack of RBMs in greedy layerwise fashion (requires some caution to avoid double counting)
- II. Use those parameters to initialize two step meanfield approach to learning full Boltzman machine (i.e. the full DBM)



### Document Clustering and Retrieval

#### **Clustering Results**

- Goal: cluster related documents
- Figures show projection to 2 dimensions
- Color shows true categories



Figure from (Salakhutdinov and Hinton, 2009)

#### DBN



## Deep Learning

Lots to explore:

- Other nonlinear functions
  - Rectified Linear Units (ReLUs)
- Popular (classic) architectures:
  - Convolutional Neural Networks (CNN)
  - Long-term Short-term Memory (LSTM)
- Modern architectures
  - Stacked SVMs with random projections
  - Sum-product Networks