#### **Deep Neural Networks** CS 550: Machine Learning

### **Deep Architectures**

 They are hard to train by backpropagation due to the vanishing gradient problem

$$\frac{\partial loss}{\partial W_{ij}} = \frac{\partial loss}{\partial net_j} \frac{\partial net_j}{\partial W_{ij}} = \delta_j x_i$$
$$\delta_j = \left[\sum_{(j+1)} \delta_{(j+1)} W_{j(j+1)}\right] \sigma'(net_j)$$

- However, when the initial weights are good enough, backpropagation works well
- Layerwise pretraining
  - Restricted Boltzmann machines
  - Autoencoders



#### Layerwise Pretraining

First, train one layer at a time, optimizing P(x)



#### Layerwise Pretraining

 Then, fine-tune weights, optimizing P(y|x) by backpropagation



*Fine-tune weights by backpropagation* 

#### Restricted Boltzmann Machines (RBMs)

RBM is a simple energy-based model

$$p(x,h) = \frac{1}{Z_{\theta}} \exp(-E_{\theta}(x,h))$$

$$E_{\theta}(x,h) = -x^{T} W h - b^{T} x - d^{T} h$$

$$Z_{\theta} = \sum_{(x,h)} \exp(-E_{\theta}(x,h))$$
normalizer

It only allows h-x interactions

- Stacked RBMs can be used to construct a deep belief net, which is a probabilistic generative model
- Stacked RBMs can also be used to initialize a deep neural network

#### Training RBMs to optimize P(x)

Maximize the log-likelihood of data

$$\partial_{W_{ij}} \log P_W(x = x^{(m)}) = \partial_{W_{ij}} \log \sum_h P_W(x = x^{(m)}, h)$$

$$= -\partial_{W_{ij}} \log Z_W + \partial_{W_{ij}} \log \sum_h \exp(-E_W(x^{(m)}, h))$$

$$= -E_{p(x,h)} [x_i h_j] + E_{p(h|x=x^{(m)})} [x_i^{(m)} h_j]$$
Negative phase

Negative phase *comes from the model comes from the data* 

*Positive phase* 

The negative phase term is expensive to calculate since it requires sampling (x, h) from the model. Contrastive divergence is a faster solution.

#### **Contrastive Divergence Algorithm**

# Initialize with the training sample and wait only a few sampling steps (usually 1 step)

1. Let  $x^{(m)}$  be a training sample and  $w_{ij}$ ,  $b_i$ , and  $d_j$  be the current weights

2. Sample 
$$\hat{h}_j \in \{0, 1\}$$
 from  $p\left(h_j \mid x = x^{(m)}\right) = \sigma\left(\sum_i w_{ij} x_i^{(m)} + d_j\right), \forall j$ 

3. Sample  $\tilde{x}_i \in \{0, 1\}$  from  $p\left(x_i \mid h = \hat{h}\right) = \sigma\left(\sum_j w_{ij} \hat{h}_j + b_i\right), \forall i$ 

4. Sample 
$$\tilde{h}_j \in \{0, 1\}$$
 from  $p(h_j \mid x = \tilde{x}) = \sigma(\sum_i w_{ij} \tilde{x}_i + d_j), \forall j$   
h's and x's are assumed  
to be binary variables

$$w_{ij} = w_{ij} + \gamma (x_i^{(m)} \hat{h}_j - \tilde{x}_i \tilde{h}_j)$$
  

$$b_i = b_i + \gamma (x_i^{(m)} - \tilde{x}_i)$$
  

$$d_j = d_j + \gamma (\hat{h}_j - \tilde{h}_j)$$

#### Autoencoders

- They learn to "compress" and "reconstruct" input data
- Learn the weights to minimize the reconstruction loss

Encoder: 
$$h = \sigma(W x + b)$$
  
Decoder:  $x' = \sigma(W' h + d)$ 

$$loss = \sum_{m} \left( x^{(m)} - x' \right)^2$$

- This is the same backpropagation for a network with one hidden layer, where x<sup>(m)</sup> is both input and output
- They can be stacked to form a deep neural network
  - Cheaper alternatives to RBMs
  - However, unlike RBMs, they are deterministic and cannot form a deep generative model

#### **Denoising Autoencoders**

Perturb input sample by adding noise to it

Encoder:  $h = \sigma(W \tilde{x} + b)$  $\tilde{x} = x + noise$ Decoder:  $x' = \sigma(W' h + d)$ 

 Learn the weights to minimize the reconstruction loss with respect to the original input sample

$$loss = \sum_{m} \left( x^{(m)} - x' \right)^2$$

#### Layerwise Pretraining

#### Is it always necessary?

- Answer in 2006: Yes!
- Answer in 2014: No!
  - If initialization is done well by design (e.g., sparse connections and convolutional nets), there may not a vanishing gradient problem
  - If the net is trained on an extremely large dataset, it may not overfit

### **Convolutional Neural Networks**

 A CNN consists of a number of convolutional and subsampling (pooling) layers optionally followed by fully connected layers



### **Convolutional Neural Networks**

 When the input data is an image, a fully connected layer will produce a huge number of weights (parameters) to be learned



Example: 200x200 image 25K hidden units → ~1B parameters

Slide credit: M. A. Ranzato

- However, spatial correlation is local and statistics is similar at different locations
- Thus, small kernels are defined and their parameters are shared by all pixels
- It is convolution with learned kernels



Example: 200x200 image 25K hidden units 10x10 kernels → ~2.5M parameters (instead of 1B parameters)



Slide credit: M. A. Ranzato



Slide credit: M. A. Ranzato



Slide credit: M. A. Ranzato

Learn multiple filters





Rectified linear unit (ReLU) provides nonlinearity: u = max(0, x)

- fast to compute
- reduces the likelihood of the gradient to vanish
- better sparsity

Learn multiple filters





Choosing the architecture (the number of feature maps, size of kernels, and number of convolutional layers) is task dependent

# **Pooling Layer**



- By pooling filter responses at different locations
  - We gain robustness to the exact location of features
  - Receptive field becomes
     larger for the next layer (the next layer will look at larger spatial regions)





Slide credit: M. A. Ranzato

#### **Pooling Layer**



Slide credit: M. A. Ranzato

### Local Contrast Normalization

Equalizes feature maps/responses

$$h^{n+1}(x, y) = \frac{h^n(x, y) - m^n(N(x, y))}{\max(\varepsilon, \sigma^n(N(x, y)))}$$



Slide credits: R. Fergus and M. A. Ranzato

# **CNNs: Typical Architecture**





All layers are differentiable so that standard backpropagation can be used

#### After one stage

- Number of feature maps is usually increased (conv. layer)
- Spatial resolution is usually decreased (pooling layer and stride in conv. layer)
- Receptive field gets larger

#### After several stages

- Spatial resolution is greatly reduced and number of feature maps is large so convolution would not make any sense
- Next layer(s) will consist of fully connected layers (with or without hidden layers)

#### **CNN Examples**



ImageNet by Krizhevsky et al., 2012 LeNet-5 by LeCun et al., 1998

# Fully Convolutional Networks (FCNs)

- An FCN is designed for <u>semantic image segmentation</u> which predicts a label for each pixel of an image
  - Image (input) and its segmentation map (output) have the same dimensions
- As opposed to a CNN designed for <u>image classification</u> which predicts a class label for the entire image
  - Input has MxN dimensions and output is a single class label



# **CNNs for Image Classification**

- A CNN compresses an image into a set of feature maps to capture semantic/contextual information from the image
- This compression corresponds to downsampling the image using <u>convolution and pooling</u> layers
- Then it puts fully connected layers on the top of the feature maps to predict a class for the entire image



# FCNs for Image Segmentation

- An FCN recovers a larger-size segmentation map from the compressed image by <u>upsampling via deconvolution</u>
  - Downsampling path captures semantic/contextual information
  - Upsampling path recovers spatial information
  - No fully connected layer is used on the top
  - Skip connections (concatenations) from downsampling to upsampling layers are often used to recover the fine-grained spatial information lost in the downsampling path



#### **Recurrent Neural Networks**

- Feedforward neural networks assume that all inputs/outputs are independent
  - However, it is not true for sequential data (speech recognition, translation, etc)
- Recurrent neural networks do not have this assumption
  - They perform the same task for every element of a sequence, with the output being dependent on previous computations
  - They might be considered to have a "memory" that captures information about what has been calculated so far



#### **Recurrent Neural Networks**



- All steps of an RNN share the same weights (U, V, W)
  - This reflects the fact that each step performs the same task just with a different input
  - Unlike a deep neural network (which uses different weights at each different layer)

#### **Recurrent Neural Networks**



- Training an RNN also uses backpropagation (called backpropagation through time – BPTT)
  - In theory, RNNs can learn arbitrarily long sequences
  - But, in practice, they have difficulties in learning long-term dependencies

#### **Training Recurrent Neural Networks**



#### Long Short Term Memory Networks

- LSTM network is a type of RNN, which is explicitly designed to avoid long-term dependency problem
- It introduces an additional cell state c<sub>t</sub> that controls the flow of information over time



Standard RNN

LSTM network

# Long Short Term Memory Networks

#### It contains four layers

- **1.** Forget gate layer  $f_t$  to control how much to remember from the previous time steps
- **2.** Input gate layer  $i_t$  to control how much to use from the current time step
- 3. Output gate layer  $\tilde{O}_t$
- 4. Tanh layer  $\tilde{c}_t$

$$\begin{aligned} f_t &= \sigma \left( U_f x_t + W_f s_{t-1} \right) \\ i_t &= \sigma \left( U_i x_t + W_i s_{t-1} \right) \\ \tilde{o}_t &= \sigma \left( U_o x_t + W_o s_{t-1} \right) \\ \tilde{c}_t &= \tanh \left( U_c x_t + W_c s_{t-1} \right) \end{aligned} \right\}$$
 sigmoid in between  $C_t$  in between  $C_t$  is a sigmoid in  $C_t$  is a sigmoid in  $C_t$  in  $C_t$  is a sigmoid in  $C_t$  is a sigmoid in  $C_t$  in  $C_t$  in  $C_t$  is a sigmoid in  $C_t$  in  $C_t$  is a sigmoid in  $C_t$  in  $C_t$  is a sigmoid in  $C_t$  in  $C_t$ 

sigmoid gives values
in between 0 and 1
0: completely forget
1: completely keep
tanh gives values in
between -1 and 1

 $c_{t} = f_{t} \circ c_{t-1} + i_{t} \circ \tilde{c}_{t}$  $s_{t} = \tilde{o}_{t} \circ \tanh(c_{t})$ 

o is an entry-wise product

