Laurens van der Maaten



- Deep learning is what people called "neural networks" in the 80s and 90s
- If you think about a kernel SVM, it basically does the following:



- Deep learning is what people called "neural networks" in the 80s and 90s
- If you think about a kernel SVM, it basically does the following:



• This is a *shallow* model: there's only one latent feature representation

- Deep learning is what people called "neural networks" in the 80s and 90s
- If you think about a kernel SVM, it basically does the following:



- This is a *shallow* model: there's only one latent feature representation
- A deep model learns multiple, increasingly high-level feature representations:



· Central idea: build increasingly high-level representations of the data



· Central idea: build increasingly high-level representations of the data



• Simple feedforward networks have the following general structure:



- Simple feedforward networks have the following general structure:
- At every layer, there is some non-linear *activation function*





- Simple feedforward networks have the following general structure:
- At every layer, there is some non-linear *activation function*





\* Nair & Hinton, 2010

- Simple feedforward networks have the following general structure:
- Network is trained to minimize some loss between the output and the target:

$$C = \frac{1}{2}(y - \tilde{y})^2$$



• Backpropagation performs gradient descent w.r.t. the network weights:



$$\mathbf{h}_1 = f_1(\mathbf{W}_1^{\mathrm{T}}\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{W}_1^{\mathrm{T}}\mathbf{x})}$$

\* Rumelhart et al., 1986

• Backpropagation performs gradient descent w.r.t. the network weights:



$$\mathbf{h}_2 = f_2(\mathbf{W}_2^{\mathrm{T}}\mathbf{h}_1) = \frac{1}{1 + \exp(-\mathbf{W}_2^{\mathrm{T}}\mathbf{h}_1)}$$

$$\mathbf{h}_1 = f_1(\mathbf{W}_1^{\mathrm{T}}\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{W}_1^{\mathrm{T}}\mathbf{x})}$$

\* Rumelhart et al., 1986

• Backpropagation performs gradient descent w.r.t. the network weights:



$$\tilde{y} = f_3(\mathbf{W}_3^{\mathrm{T}}\mathbf{h}_2) = \frac{1}{1 + \exp(-\mathbf{W}_3^{\mathrm{T}}\mathbf{h}_2)}$$

$$\mathbf{h}_2 = f_2(\mathbf{W}_2^{\mathrm{T}}\mathbf{h}_1) = \frac{1}{1 + \exp(-\mathbf{W}_2^{\mathrm{T}}\mathbf{h}_1)}$$

$$\mathbf{h}_1 = f_1(\mathbf{W}_1^{\mathrm{T}}\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{W}_1^{\mathrm{T}}\mathbf{x})}$$

\* Rumelhart et al., 1986

• Often, we want the network to minimize the sum of squared errors:



$$C = \frac{1}{2}(y - \tilde{y})^2$$

- *Backpropagation* proceeds in two steps:
  - 1. Propagate the error signal over weights
  - 2.Compute the associated weight update

• Backpropagation performs *gradient descent* w.r.t. the network weights:



$$e_{3} = (y - \tilde{y})\frac{\partial f_{3}}{\partial \mathbf{W}_{3}} = (y - \tilde{y})\tilde{y}(1 - \tilde{y})$$
$$\frac{C}{\mathbf{V}_{3}} = e_{3}\mathbf{h}_{2}$$

• Backpropagation performs *gradient descent* w.r.t. the network weights:



$$e_{3} = (y - \tilde{y}) \frac{\partial f_{3}}{\partial \mathbf{W}_{3}} = (y - \tilde{y}) \tilde{y}(1 - \tilde{y})$$
$$\frac{\partial C}{\partial \mathbf{W}_{3}} = e_{3} \mathbf{h}_{2}$$
$$\mathbf{e}_{2} = \mathbf{h}_{2} \circ (1 - \mathbf{h}_{2}) \circ (\mathbf{W}_{3} e_{3})$$
$$\frac{\partial C}{\partial \mathbf{W}_{2}} = \mathbf{e}_{2}^{\mathrm{T}} \mathbf{h}_{1}$$

• Backpropagation performs *gradient descent* w.r.t. the network weights:



$$e_{3} = (y - \tilde{y}) \frac{\partial f_{3}}{\partial \mathbf{W}_{3}} = (y - \tilde{y}) \tilde{y}(1 - \tilde{y})$$
$$\frac{\partial C}{\partial \mathbf{W}_{3}} = e_{3} \mathbf{h}_{2}$$
$$\mathbf{e}_{2} = \mathbf{h}_{2} \circ (1 - \mathbf{h}_{2}) \circ (\mathbf{W}_{3} e_{3})$$
$$\frac{\partial C}{\partial \mathbf{W}_{2}} = \mathbf{e}_{2}^{\mathrm{T}} \mathbf{h}_{1}$$
$$\mathbf{e}_{1} = \mathbf{h}_{1} \circ (1 - \mathbf{h}_{1}) \circ (\mathbf{W}_{2} \mathbf{e}_{2})$$
$$\frac{\partial C}{\partial \mathbf{W}_{1}} = \mathbf{e}_{1}^{\mathrm{T}} \mathbf{x}$$

- Training a feedforward net on large images leads to huge numbers of weights
  - Do you expect to see some structure in those weights?

- Training a feedforward net on large images leads to huge numbers of weights
  - Do you expect to see some structure in those weights? Yes!

- Training a feedforward net on large images leads to huge numbers of weights
  - Do you expect to see some structure in those weights? Yes!
- Convolutional networks extract the same features from different image parts
- This restricts the number of weights
  we need to learn

The red connections all have the same weight.



\* LeCun et al., 1989

• Local *pooling* of features increases invariance to small image variations

- Local *pooling* of features increases invariance to small image variations
  - Such local pooling also happens in SIFT, HOG, BoW, etc.

- Local *pooling* of features increases invariance to small image variations
  - Such local pooling also happens in SIFT, HOG, BoW, etc.
- Convolutional networks generally also include such *pooling layers*:



\* LeCun et al., 1989

• LeNet5 was one of the first convolutional nets (for reading cheques):



• Erroneously recognized digits by LeNet5:



• Convolutional networks indeed appear to learn increasingly complex features:



# The ILSVRC-2012 Competition

• Substantial improvement of the state-of-the-art in object recognition:

#### ImageNet Challenge 2012

	Submission	Method	Error rate	
	SuperVision	Convolutional net	0.16422	<b>9</b> .8%
	ISI		0.26172	
	XRCE/INRIA	Other stuff (SVMs)	0.27058	ې 1.1%-
	OXFORD_VGG		0.27302	3 3S
X				

### **Recognition examples**



# The ILSVRC-2012 Competition

- The winning model was a single, large convolutional network:
  - Seven hidden layers not counting the max-pooling layers
  - First five layers were convolutional; the last two fully connected
  - Rectified linear units were used as activation function everywhere
  - Model used a *competitive normalization* trick: suppress hidden activities when nearby units have stronger activities to deal with intensity variations
  - Total number of 80 million parameters

# The ILSVRC-2012 Competition





• Convolutional networks have tons of parameters: prone to overfitting

- Convolutional networks have tons of parameters: prone to overfitting
- Simple solution is to use *weight decay*: penalize L2-norm of network weights

- Convolutional networks have tons of parameters: prone to overfitting
- Simple solution is to use *weight decay*: penalize L2-norm of network weights
- Alternative, better way to regularize such networks is via *dropout*:
  - Randomly switch off some hidden units during forward pass

- Convolutional networks have tons of parameters: prone to overfitting
- Simple solution is to use *weight decay*: penalize L2-norm of network weights
- Alternative, better way to regularize such networks is via *dropout*:
  - Randomly switch off some hidden units during forward pass
- Moreover, train model on slightly warped positive examples:



#### Some more examples



# Retrieval using convolutional net



# Implementing a convolutional net

- Convolutional networks perform many convolutions and matrix multiplications
  - These are very fast on GPUs (up to 100 times faster)
- GPU implementations (and Big Data) are essential in the recent successes of deep learning and convolutional nets

# Implementing a convolutional net

- Convolutional networks perform many convolutions and matrix multiplications
  - These are very fast on GPUs (up to 100 times faster)
- GPU implementations (and Big Data) are essential in the recent successes of deep learning and convolutional nets

• Libraries for building deep (convolutional) nets: Torch7, eblearn, Theano, Caffe

## Demonstration

#### • http://horatio.cs.nyu.edu



#### **Demo Notes**

- If your images have objects that are not in the 1,000 categories of ImageNet, the model will not know about them.
- Other objects can be added from all 20,000+ ImageNet categories (it may be slow to load the autocomplete results...just wait a little).
- The maximum file size for uploads in this demo is 10 MB.
- Only image files (JPEG, JPG, GIF, PNG) are allowed in this demo .
- You can drag & drop files from your desktop on this webpage with Google Chrome, Mozilla Firefox and Apple Safari.
- · Some mobile browsers are known to work, others will not. Try updating your browser or contact us with the problem.
- All images for your current IP and browsing session are shown above and not shown to others.
- This demo is powered by research out of New York University. Click here to find out more
- If you encounter problems, please contact zeiler@cs.nyu.edu

#### Demo created by: Matthew Zeiler

# Questions?