
Machine Learning 10-601

Tom M. Mitchell
Machine Learning Department
Carnegie Mellon University

September 27, 2017

This section:

- Backpropagation
- Representation learning
- Convolutional nets

Reading:

- Goodfellow: Convolutional nets

Training Deep Nets

1. Choose loss function $J(\theta)$ to optimize

- sum of squared errors for y continuous: $\sum (y - h(x; \theta))^2$
- maximize conditional likelihood: $\sum \log P(y|x; \theta)$
- MAP estimate: $\sum \log P(y|x; \theta) P(\theta)$
- ~~0/1 loss. Sum of classification errors: $\sum \delta(y = h(x; \theta))$~~
- ...

2. Design network architecture

- Network of layers (ReLU's, sigmoid, convolutions, ...)
- Widths of layers
- Fully or partly interconnected
- ...

3. Training algorithm

- Derive gradient components $\frac{\partial J(\theta)}{\partial \theta_i}$
- Choose gradient descent method, including stopping condition
- Experiment with alternative architectures

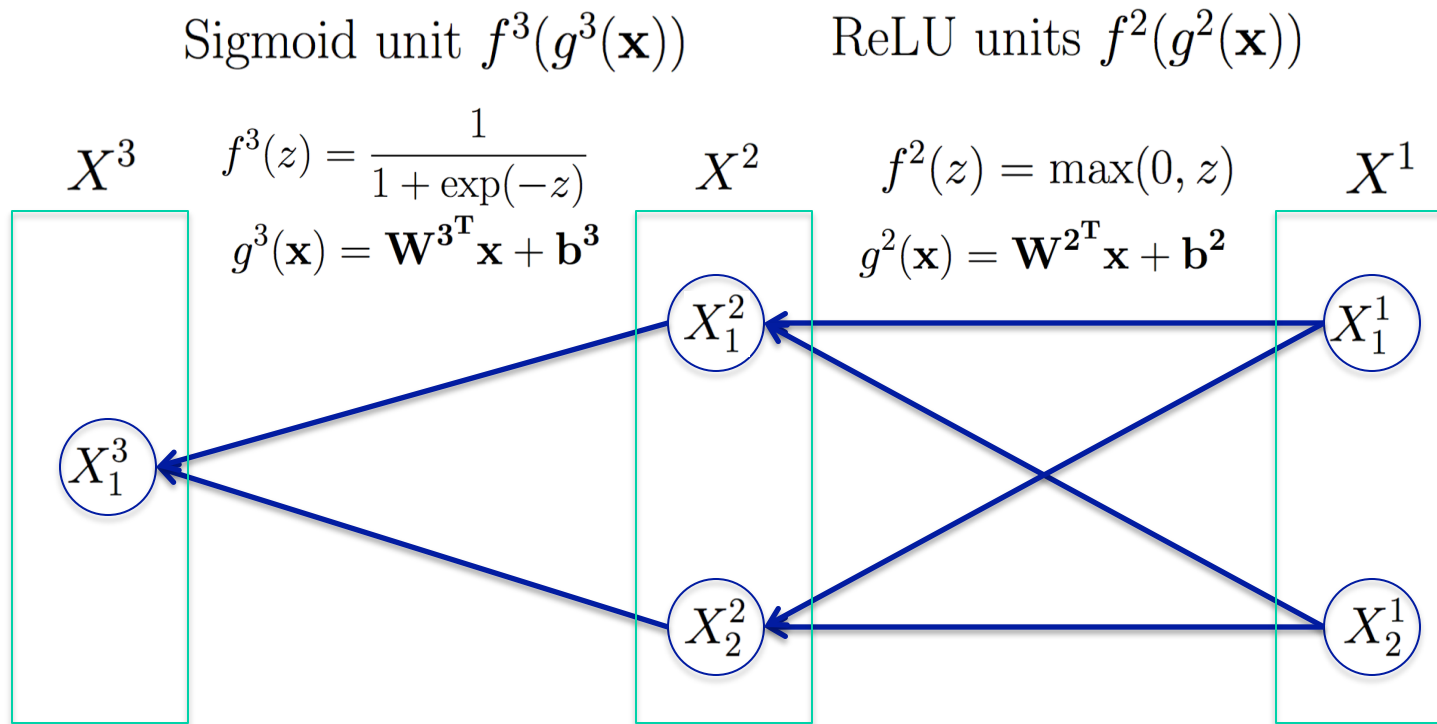
Example: Learn probabilistic XOR

- Given boolean Y , X_1 , X_2 learn $P(Y|X_1, X_2)$, where

$$P(Y = 0|X_1 = X_2) = 0.9$$

$$P(Y = 1|X_1 \neq X_2) = 0.9$$

- Can we learn this with logistic regression?



Loss function $J(\theta)$ to be minimized: negative log likelihood

$$J(\theta) = \sum_{\langle \mathbf{x}, y \rangle \in D} -\log P(Y = y | X = \mathbf{x})$$

where $X_1^3 = P(Y = 1 | X = \mathbf{X}^1)$, $\theta = \{\mathbf{W}^3, \mathbf{b}^3, \mathbf{W}^2, \mathbf{b}^2\}$

Example: Learn probabilistic XOR

- Given boolean Y , X_1 , X_2 learn $P(Y|X_1, X_2)$, where

$$P(Y = 0|X_1 = X_2) = 0.9$$

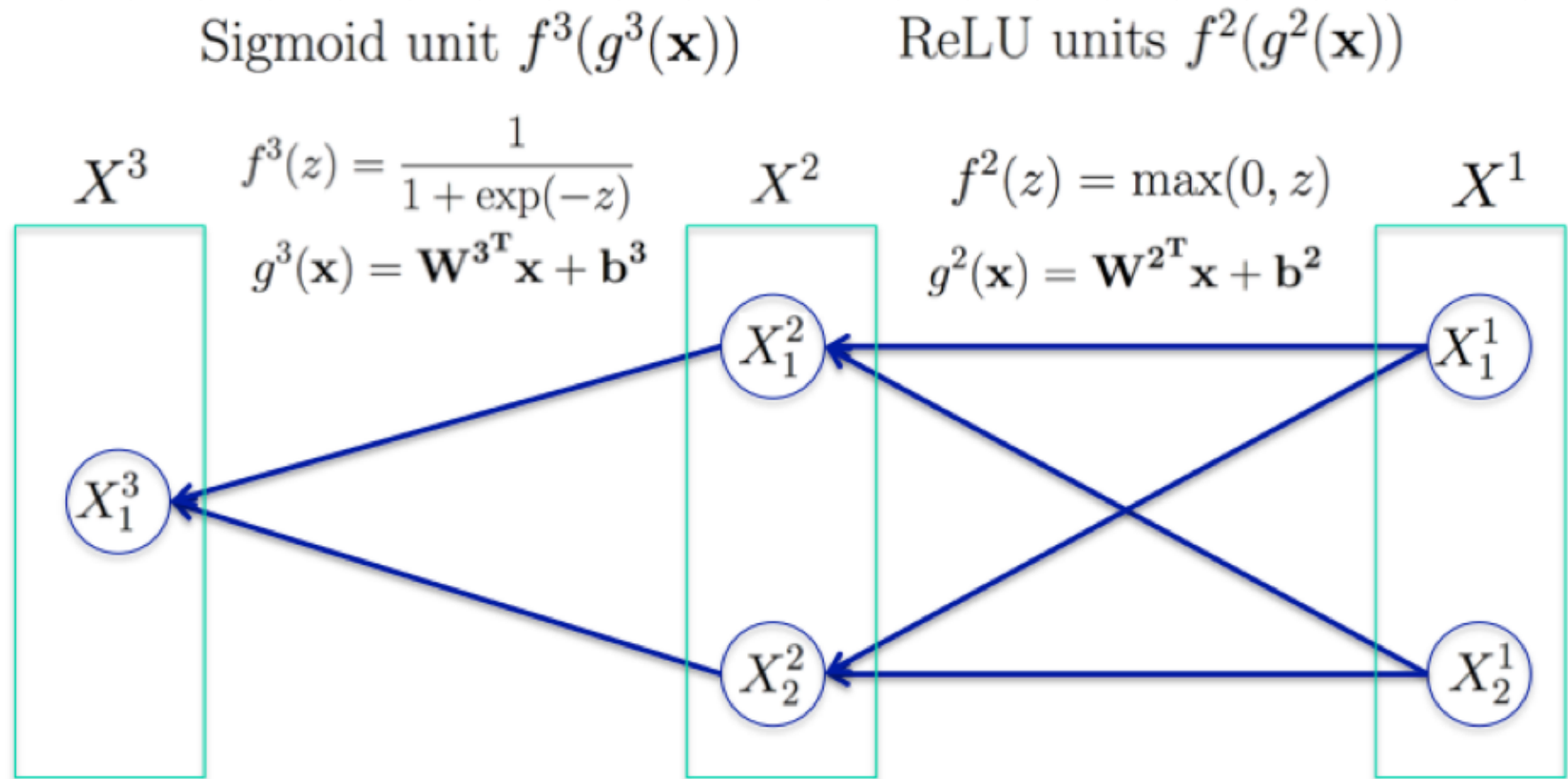
$$P(Y = 1|X_1 \neq X_2) = 0.9$$

- Choose max conditional likelihood,
equally, minimize negative log conditional likelihood

$$\begin{aligned} J(\theta) &= -\sum_k \log P(Y = y_k|X = \mathbf{x}_k) \\ &= -\sum_k y_k \log P(Y = 1|X = x_k) + (1 - y_k)(1 - P(Y = 1|X = x_k)) \end{aligned}$$

Example

Feed Forward



\mathbf{x}^3	$g^3(\mathbf{x}^2)$	\mathbf{W}^{3T}	\mathbf{b}^3
0.53	0.12	0.10 -0.09	0.10

\mathbf{x}^2	$g^2(\mathbf{x}^1)$	\mathbf{W}^{2T}	\mathbf{b}^2
0.20 0.00 1	0.20 -0.15	0.10 -0.10 -0.20 0.10	0.10 0.05

\mathbf{x}^1
1 0 1

Derive the gradient we need

$$\begin{aligned} J(\theta) &= -\sum_k \log P(Y = y_k | X = \mathbf{x}_k) \\ &= -\sum_k y_k \log P(Y = 1 | X = x_k) + (1 - y_k)(1 - P(Y = 1 | X = x_k)) \end{aligned}$$

simplify notation by considering just one training example

$$\begin{aligned} \frac{\partial J(\theta)}{\partial X_1^3} &= \frac{\partial (-Y \log P(Y = 1 | X) - (1 - Y)(1 - P(Y = 1 | X)))}{\partial X_1^3} \\ &= \frac{\partial (-Y \log X_1^3 - (1 - Y) \log(1 - X_1^3))}{\partial X_1^3} \\ &= \frac{-Y}{X_1^3} - (1 - Y) \frac{1}{1 - X_1^3} (-1) \\ &= \frac{-Y}{X_1^3} + \frac{(1 - Y)}{1 - X_1^3} \end{aligned}$$

recall $\frac{\partial \ln z}{\partial z} = \frac{1}{z}$

Derive the gradient we need

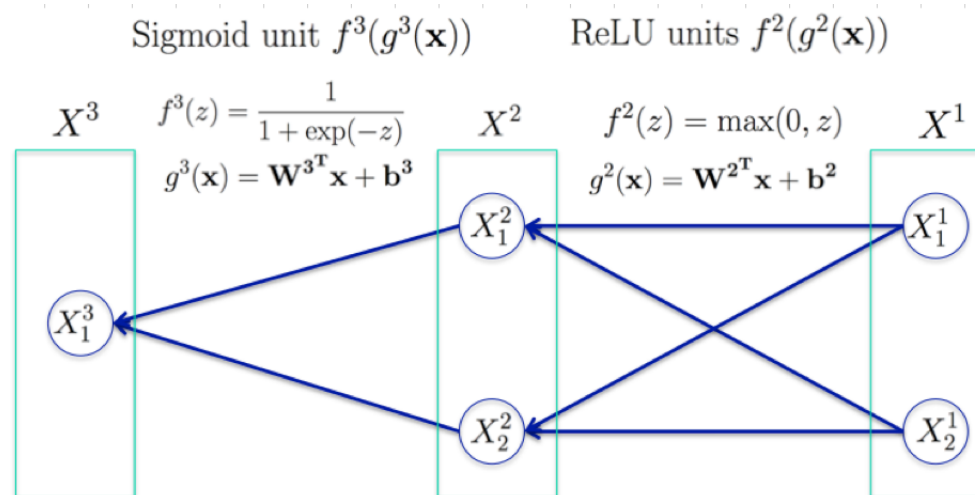
$$\begin{aligned} J(\theta) &= -\sum_k \log P(Y = y_k | X = \mathbf{x}_k) \\ &= -\sum_k y_k \log P(Y = 1 | X = x_k) + (1 - y_k) \log P(Y = 0 | X = x_k) \end{aligned}$$

simplify notation by considering just one training example:

$$\begin{aligned} \frac{\partial J(\theta)}{\partial X_1^3} &= \frac{\partial (-Y \log P(Y = 1|X) - (1 - Y)(1 - P(Y = 1|X)))}{\partial X_1^3} \\ &= \frac{\partial (-Y \log X_1^3 - (1 - Y) \log(1 - X_1^3))}{\partial X_1^3} \\ &= \frac{-Y}{X_1^3} - (1 - Y) \frac{1}{1 - X_1^3} (-1) \\ &= \frac{-Y}{X_1^3} + \frac{(1 - Y)}{1 - X_1^3} \end{aligned}$$

recall $\frac{\partial \ln z}{\partial z} = \frac{1}{z}$

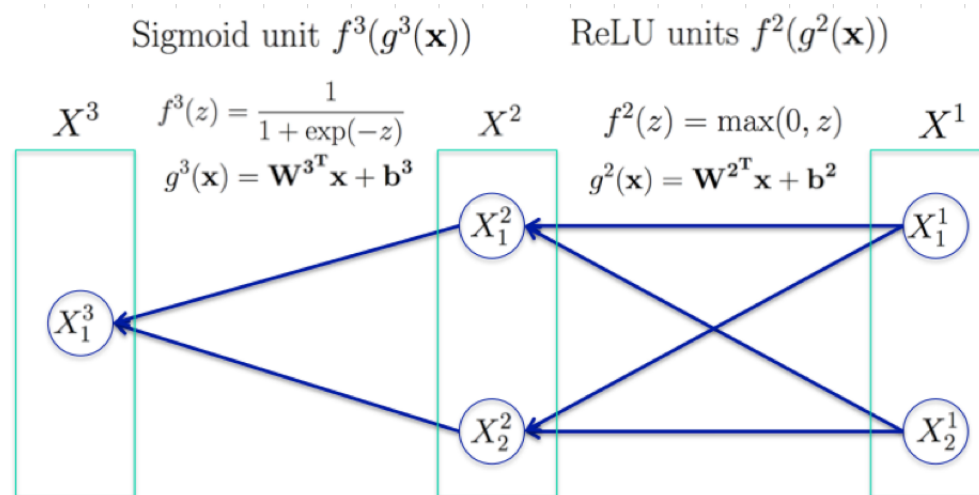
	x^3	$g^3(x^2)$	w^{3T}	b^3		x^2	$g^2(x^1)$	w^{2T}	b^2		x^1
$Y_{\text{true}}=1$	0.53	0.12	0.10 -0.09 0.10			0.20 0.00	0.20 -0.15	0.10 -0.10 0.10 -0.20 0.10 0.05			1 0
						1					1



$$\frac{\partial J(\theta)}{\partial g_1^3} = \frac{\partial J(\theta)}{\partial X_1^3} \frac{\partial X_1^3}{\partial g^3} =$$

recall $\frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1 - \sigma(z))$
 where $\sigma(z) = \frac{1}{1 + \exp(-z)}$

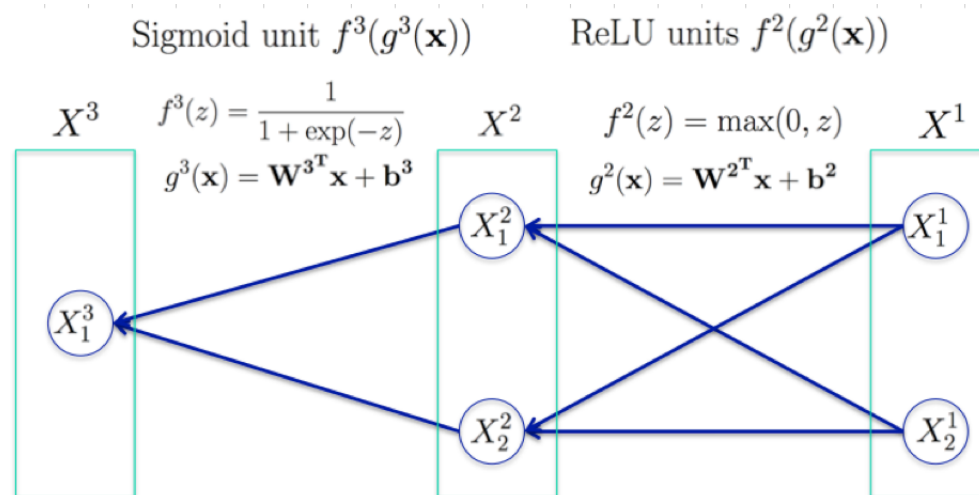
	\mathbf{x}^3	$\mathbf{g}^3(\mathbf{x}^2)$	\mathbf{W}^{3T}	\mathbf{b}^3		\mathbf{x}^2	$\mathbf{g}^2(\mathbf{x}^1)$	\mathbf{W}^{2T}	\mathbf{b}^2		\mathbf{x}^1
$Y_{\text{true}}=1$	0.53	0.12	0.10 -0.09 0.10		0.20 0.00 1	0.20 -0.15	0.10 -0.10 0.10 -0.20 0.10 0.05		1 0 1		



$$\frac{\partial J(\theta)}{\partial g_1^3} = \frac{\partial J(\theta)}{\partial X_1^3} \frac{\partial X_1^3}{\partial g^3} = \frac{\partial J(\theta)}{\partial X_1^3} X_1^3 (1 - X_1^3) =$$

recall $\frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1 - \sigma(z))$
 where $\sigma(z) = \frac{1}{1 + \exp(-z)}$

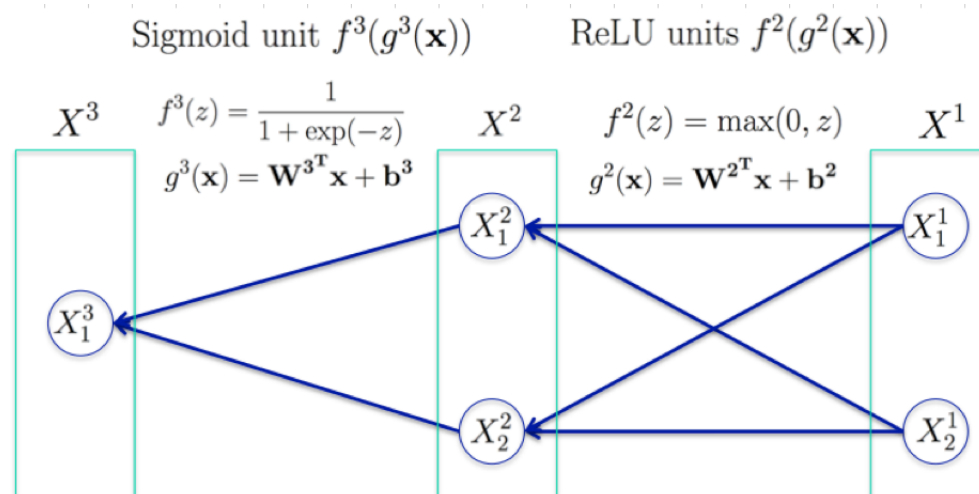
	X^3	$g^3(X^2)$	W^{3T}	b^3		X^2	$g^2(X^1)$	W^{2T}	b^2		X^1
$Y_{\text{true}}=1$	0.53	0.12	0.10 -0.09 0.10			0.20 0.00 1	0.20 -0.15	0.10 -0.10 0.10 -0.20 0.10 0.05			1 0 1



$$\frac{\partial J(\theta)}{\partial g_1^3} = \frac{\partial J(\theta)}{\partial X_1^3} \frac{\partial X_1^3}{\partial g^3} = \frac{\partial J(\theta)}{\partial X_1^3} X_1^3 (1 - X_1^3) = -0.47$$

$$\frac{\partial J(\theta)}{\partial w_i^3} = \frac{\partial J(\theta)}{\partial g^3} \frac{\partial g^3}{\partial w_i^3} = \frac{\partial J(\theta)}{\partial g^3} X_i^2$$

	\mathbf{X}^3	$\mathbf{g}^3(\mathbf{X}^2)$	\mathbf{W}^{3T}	\mathbf{b}^3		\mathbf{X}^2	$\mathbf{g}^2(\mathbf{X}^1)$	\mathbf{W}^{2T}	\mathbf{b}^2		\mathbf{X}^1
$Y_{\text{true}}=1$	0.53	0.12	0.10 -0.09 0.10		0.20 0.00 1	0.20 -0.15	0.10 -0.10 0.10 -0.20 0.10 0.05		1 0 1		



$$\frac{\partial J(\theta)}{\partial g_1^3} = \frac{\partial J(\theta)}{\partial X_1^3} \frac{\partial X_1^3}{\partial g^3} = \frac{\partial J(\theta)}{\partial X_1^3} X_1^3 (1 - X_1^3)$$

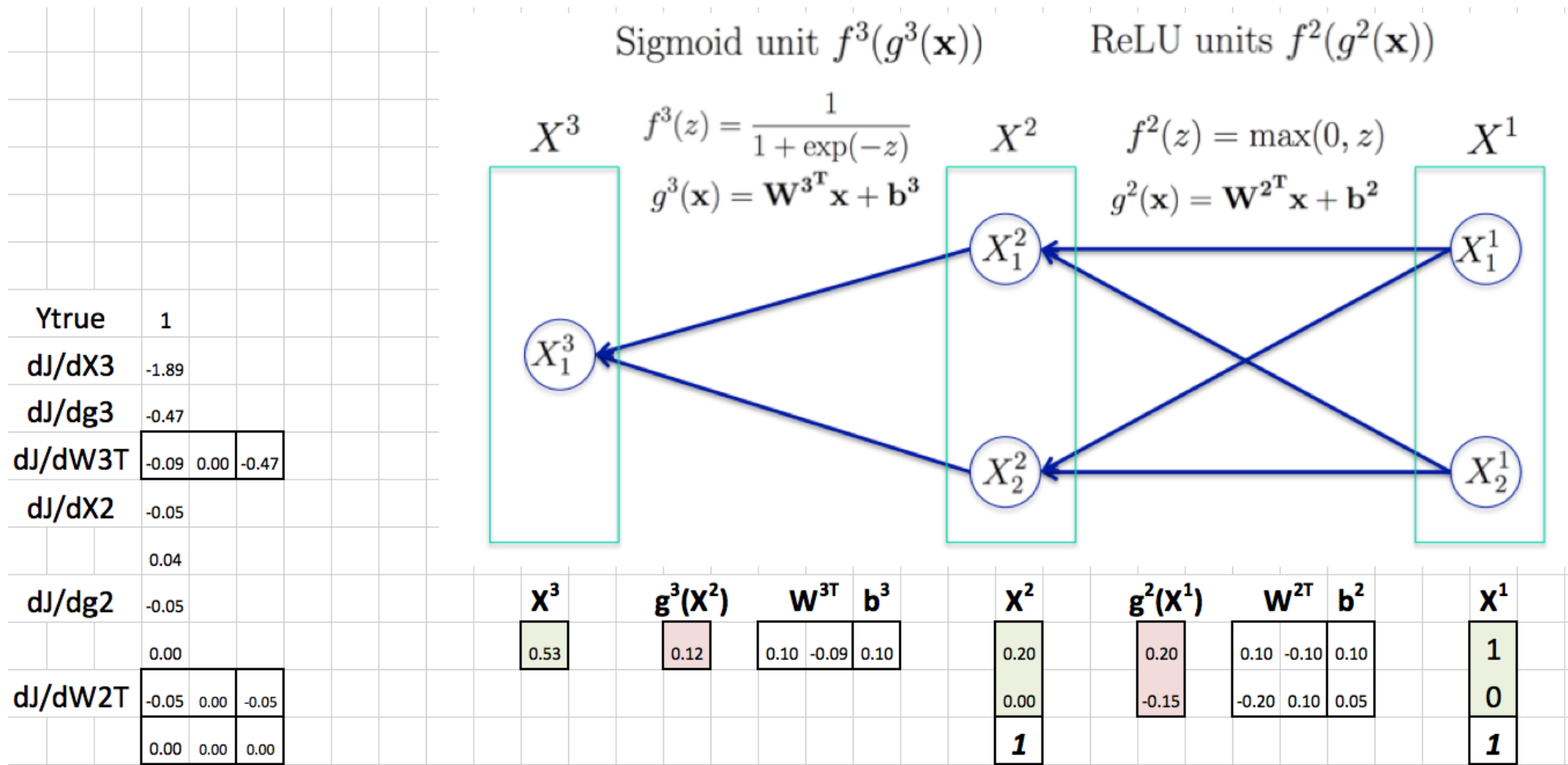
$$\frac{\partial J(\theta)}{\partial w_i^3} = \frac{\partial J(\theta)}{\partial g^3} \frac{\partial g^3}{\partial w_i^3} = \frac{\partial J(\theta)}{\partial g^3} X_i^2$$

$$\frac{\partial J(\theta)}{\partial X_i^2} = \frac{\partial J(\theta)}{\partial g^3} \frac{\partial g^3}{\partial X_i^2} = \frac{\partial J(\theta)}{\partial g^3} w_i^3$$

$$\frac{\partial J(\theta)}{\partial g_i^2} = \frac{\partial J(\theta)}{\partial X_i^2} \frac{\partial X_i^2}{\partial g_i^2} = \frac{\partial J(\theta)}{\partial X_i^2} \times [\text{if } g_i^2 > 0 \text{ then } 1 \text{ else } 0]$$

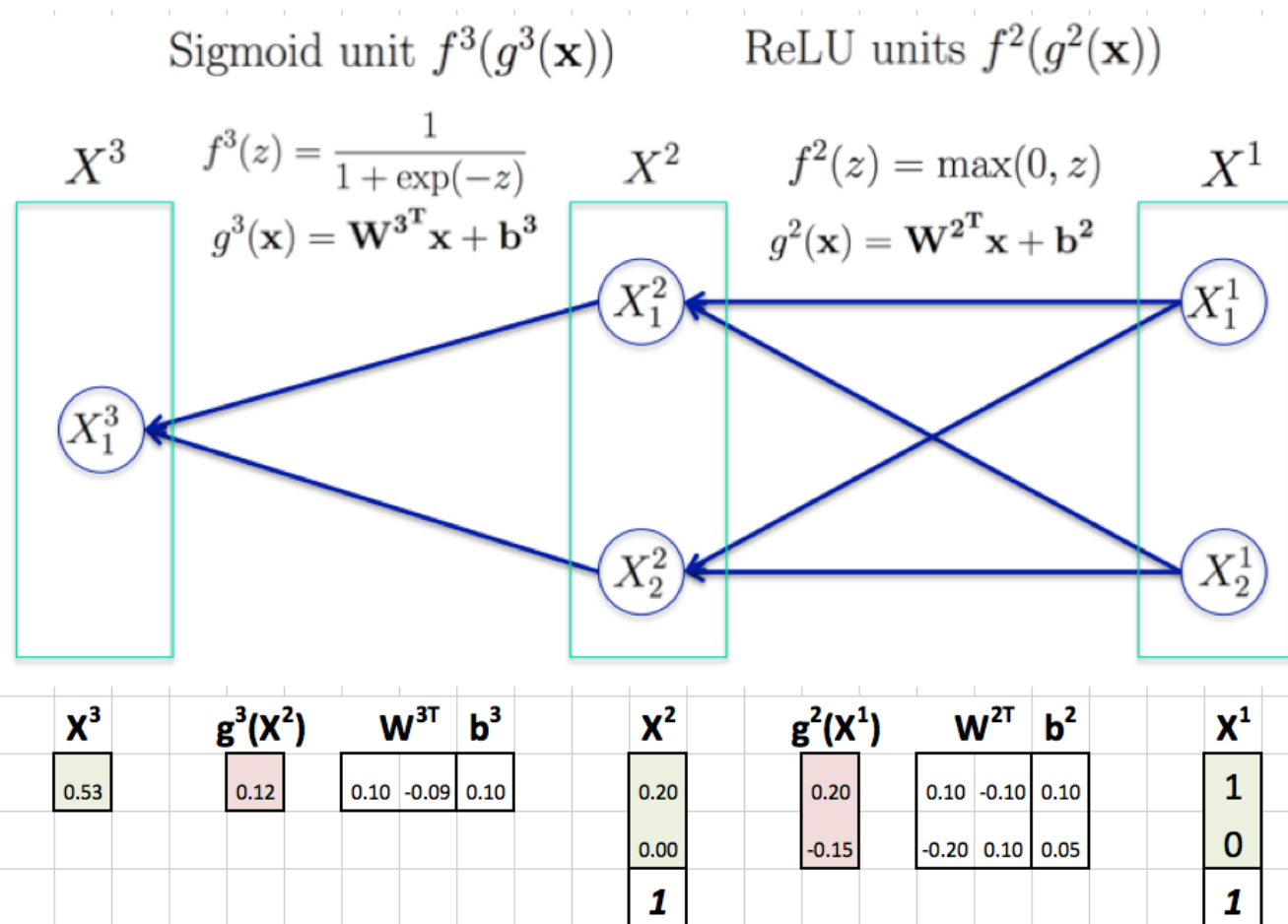
$$\frac{\partial J(\theta)}{\partial w_{ik}^2} = \frac{\partial J(\theta)}{\partial g_i^2} \frac{\partial g_i^2}{\partial X_k^1} = \frac{\partial J(\theta)}{\partial g^3} X_k^1$$

Back propagation



update each parameter according to $\theta_i \leftarrow \theta_i - \eta \frac{\partial J(\theta)}{\partial \theta_i}$

Ytrue	1
dJ/dX3	-1.89
dJ/dg3	-0.47
dJ/dW3T	-0.09 0.00 -0.47
dJ/dX2	-0.05
	0.04
dJ/dg2	-0.05
	0.00
dJ/dW2T	-0.05 0.00 -0.05
	0.00 0.00 0.00

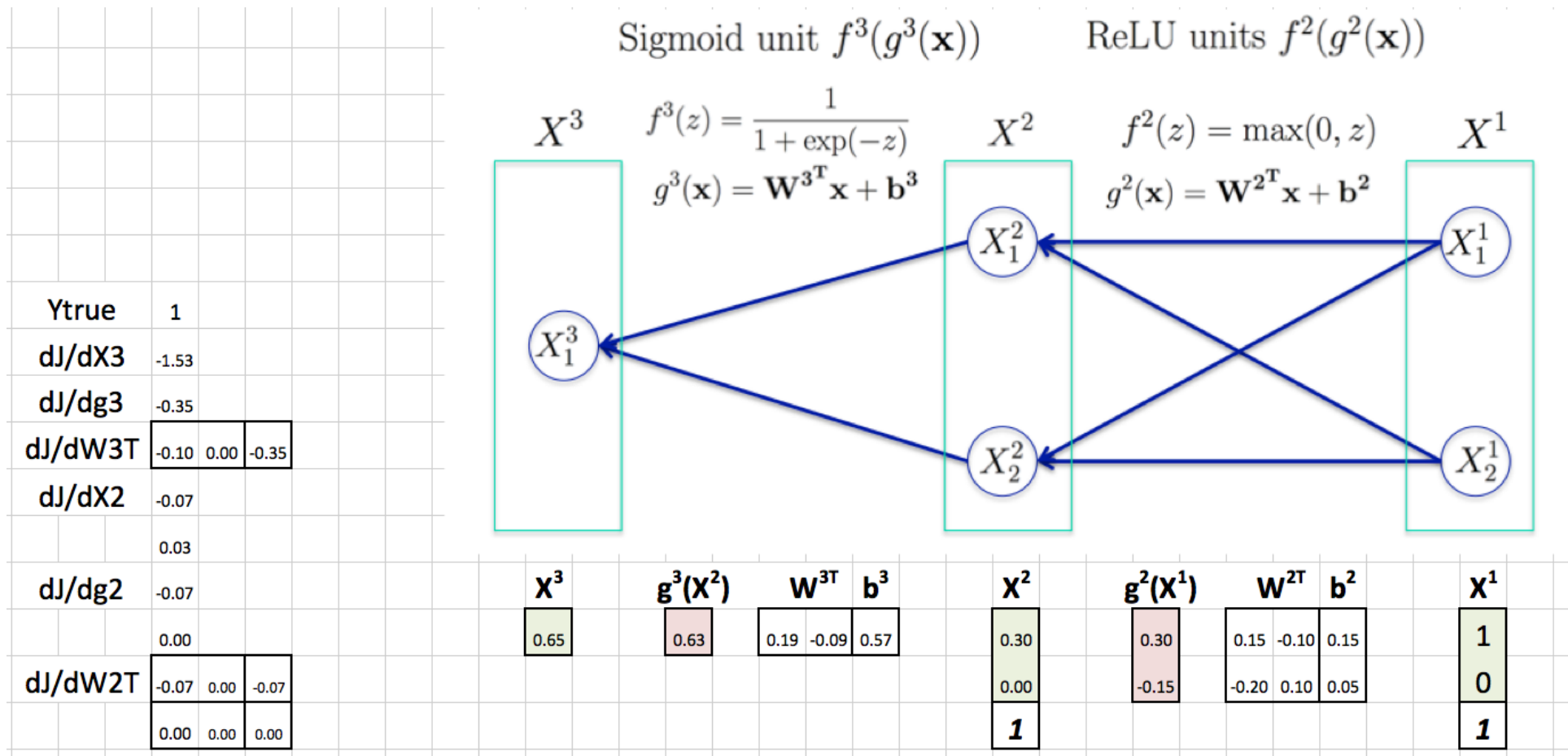


update each parameter according to $\theta_i \leftarrow \theta_i - \eta \frac{\partial J(\theta)}{\partial \theta_i}$

using $\eta = 1$:

\mathbf{x}^3	$\mathbf{g}^3(\mathbf{X}^2)$	\mathbf{W}^{3T}	\mathbf{b}^3	\mathbf{x}^2	$\mathbf{g}^2(\mathbf{X}^1)$	\mathbf{W}^{2T}	\mathbf{b}^2	\mathbf{x}^1
0.65	0.63	0.19 -0.09 0.57		0.30	0.30	0.15 -0.10 0.15		1
				0.00	-0.15	-0.20 0.10 0.05		0
				1				1

Next training example.. next stochastic gradient step...



Backpropagation Algorithm

for sigmoid netwk, minimizing $\sum_d (t_d - o_d)^2$

Initialize all weights to small random numbers.
Until satisfied, Do

- For each training example, Do
 1. Input the training example to the network and compute the network outputs
 2. For each output unit k

$$\delta_k \leftarrow o_k(1 - o_k)(t_k - o_k)$$

3. For each hidden unit h

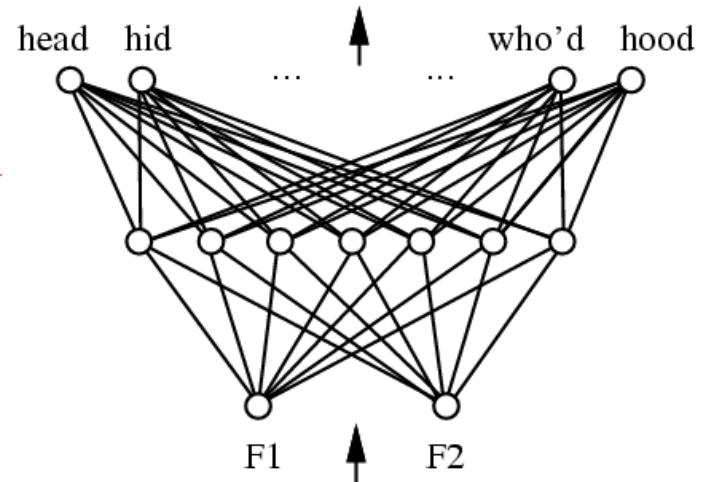
$$\delta_h \leftarrow o_h(1 - o_h) \sum_{k \in \text{outputs}} w_{h,k} \delta_k$$

4. Update each network weight $w_{i,j}$

$$w_{i,j} \leftarrow w_{i,j} + \Delta w_{i,j}$$

where

$$\Delta w_{i,j} = \eta \delta_j x_i$$



x_d = input
 o_i = observed i^{th} unit output
 t_i = target output
 w_{ij} = wt from i to j
 δ_i = error term backpropagated to unit k :
 $dJ(\theta) / dg^k(X^{k-1})$

Many modifications to gradient descent

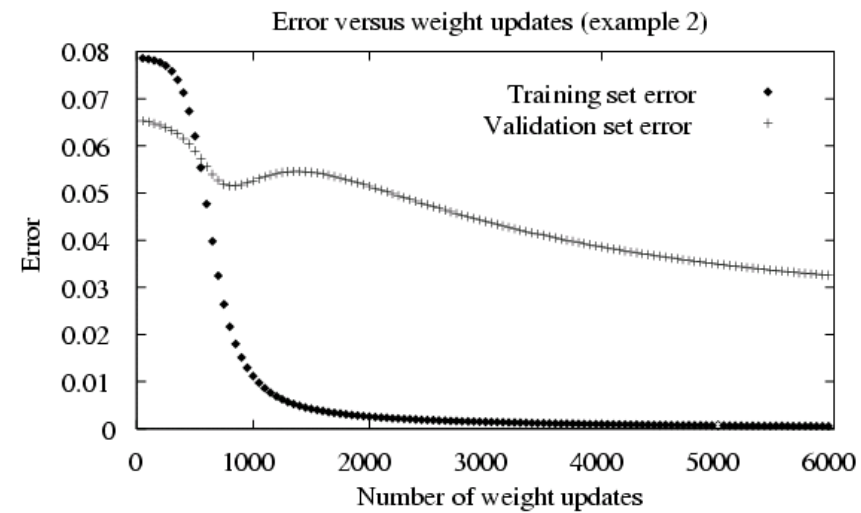
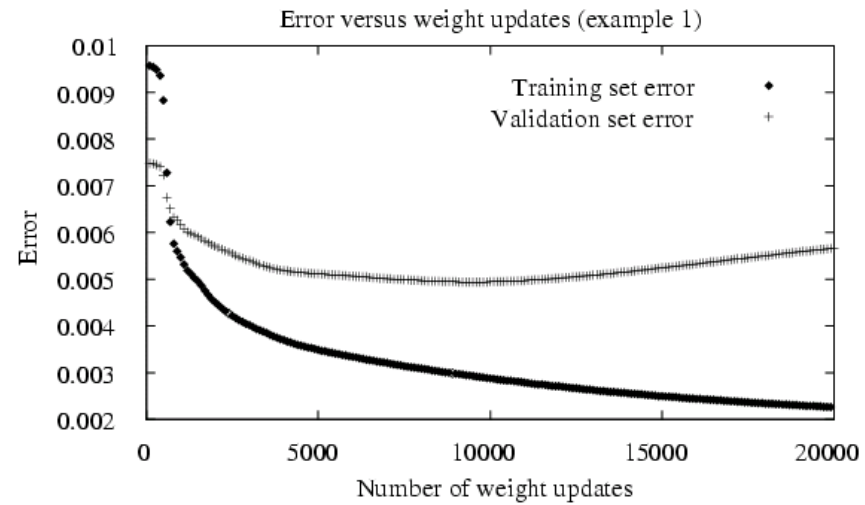
- Stochastic vs. Batch gradient descent (and mini-batches)
- Momentum
- Weight decay (MAP estimate with zero-mean prior)
- Gradient clipping
- Batch normalization
- Dropout
- Adagrad
- Adam
- no end in sight

See ML Department course on Optimization Methods

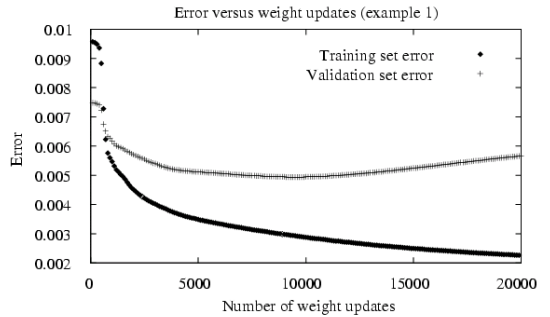
Gradient Descent and Backpropagation

- Updates every network parameter simultaneously, each iteration
- Easily generalized to arbitrary directed acyclic graphs
- Finds local minimum in $J(\theta)$, not necessarily global min
- Minimizes $J(\theta)$ over training examples, not necessarily future...
- Training can require hours, days, weeks, GPUs
- Applying network after training is relatively very fast

Overfitting in ANNs



Dealing with Overfitting



Our learning algorithm involves a parameter
 n = number of gradient descent iterations

How do we choose n to optimize future error or loss?

- Separate available data into training and validation set
- Use training to perform gradient descent
- $n \leftarrow$ number of iterations that optimizes validation set error

→ *gives unbiased estimate of optimal n*
(but still an optimistically biased estimate of true error)

Expressive Capabilities of ANNs

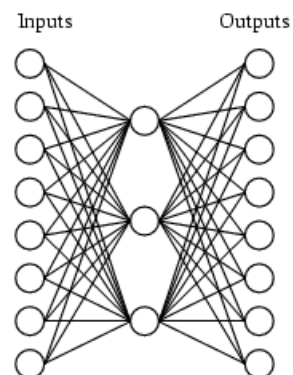
Boolean functions:

- Every boolean function can be represented by network with single hidden layer
- but might require exponential (in number of inputs) hidden units

Continuous functions:

- Every bounded continuous function can be approximated with arbitrarily small error, by network with one hidden layer [Cybenko 1989; Hornik et al. 1989]
- Any function can be approximated to arbitrary accuracy by a network with two hidden layers [Cybenko 1988].

Learning Hidden Layer Representations



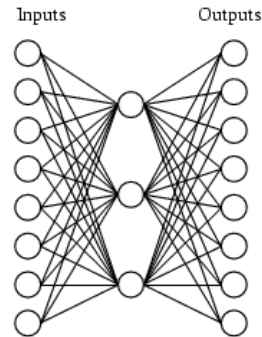
A target function:

Input	Output
10000000 →	10000000
01000000 →	01000000
00100000 →	00100000
00010000 →	00010000
00001000 →	00001000
00000100 →	00000100
00000010 →	00000010
00000001 →	00000001

Can this be learned??

Learning Hidden Layer Representations

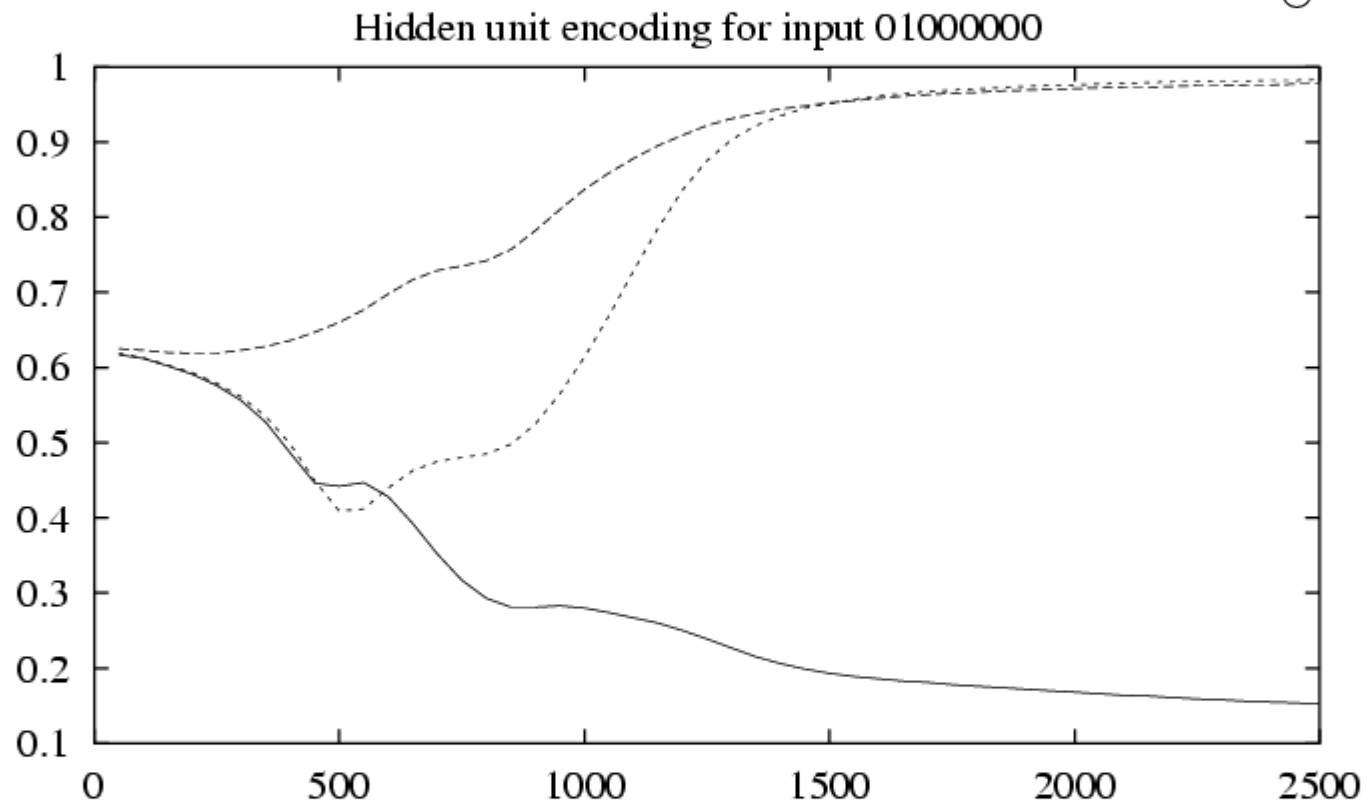
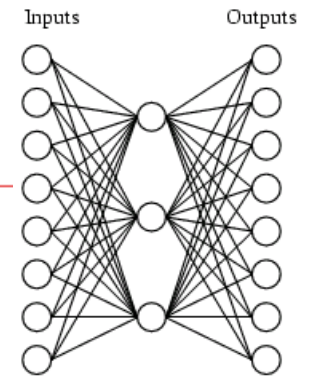
A network:



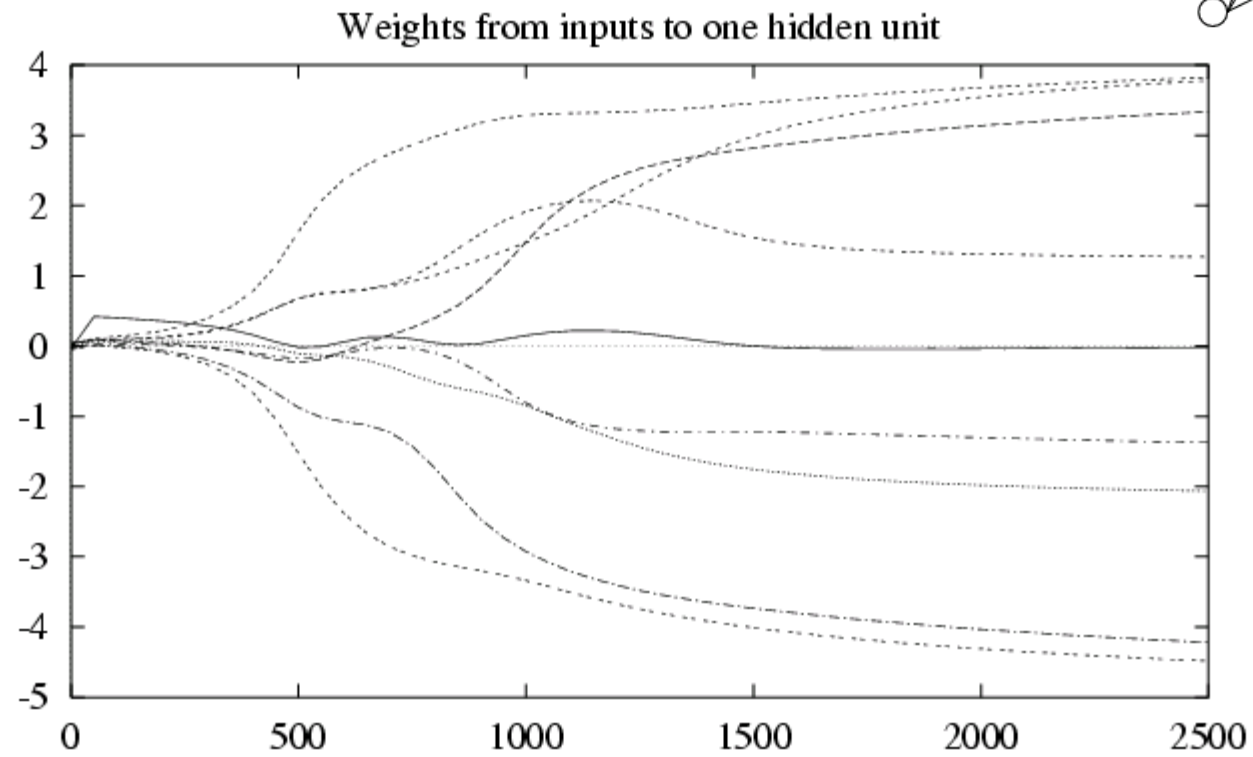
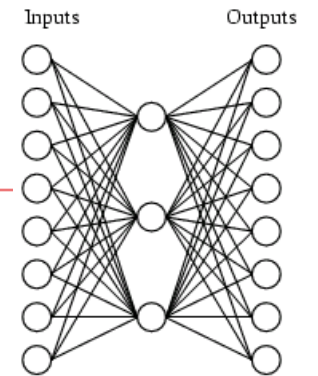
Learned hidden layer representation:

Input		Hidden Values		Output
10000000	→	.89 .04 .08	→	10000000
01000000	→	.01 .11 .88	→	01000000
00100000	→	.01 .97 .27	→	00100000
00010000	→	.99 .97 .71	→	00010000
00001000	→	.03 .05 .02	→	00001000
00000100	→	.22 .99 .99	→	00000100
00000010	→	.80 .01 .98	→	00000010
00000001	→	.60 .94 .01	→	00000001

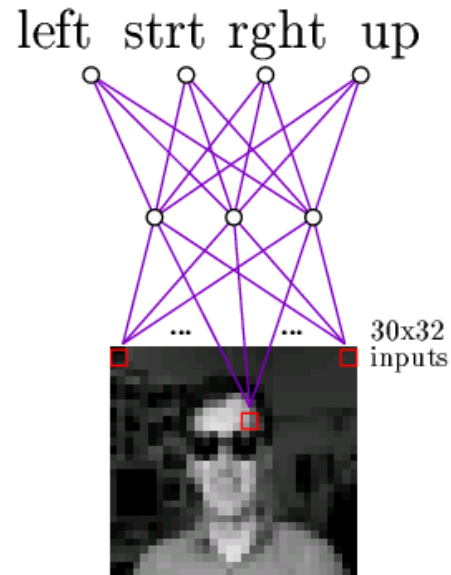
Training



Training



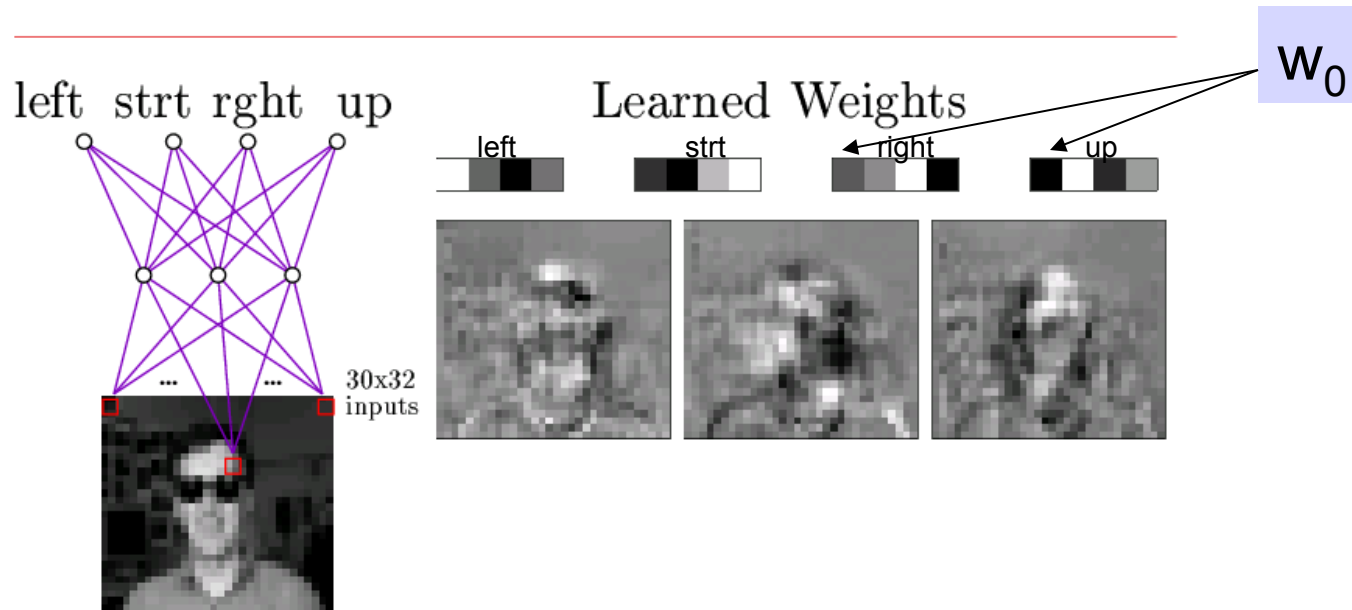
Neural Nets for Face Recognition



Typical input images

90% accurate learning head pose, and recognizing 1-of-20 faces

Learned Hidden Unit Weights

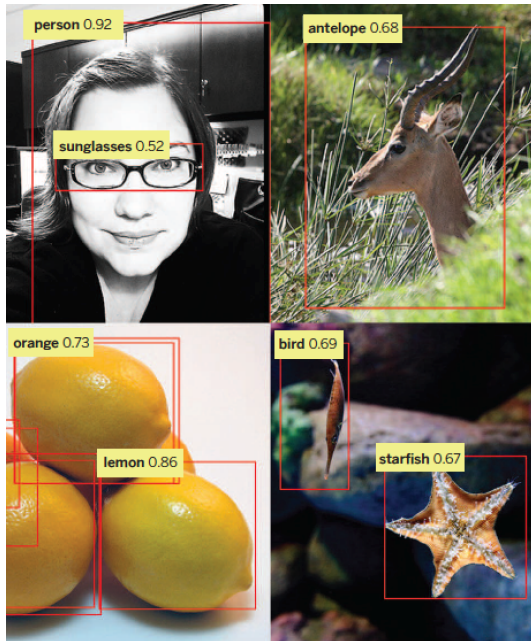


Typical input images

<http://www.cs.cmu.edu/~tom/faces.html>

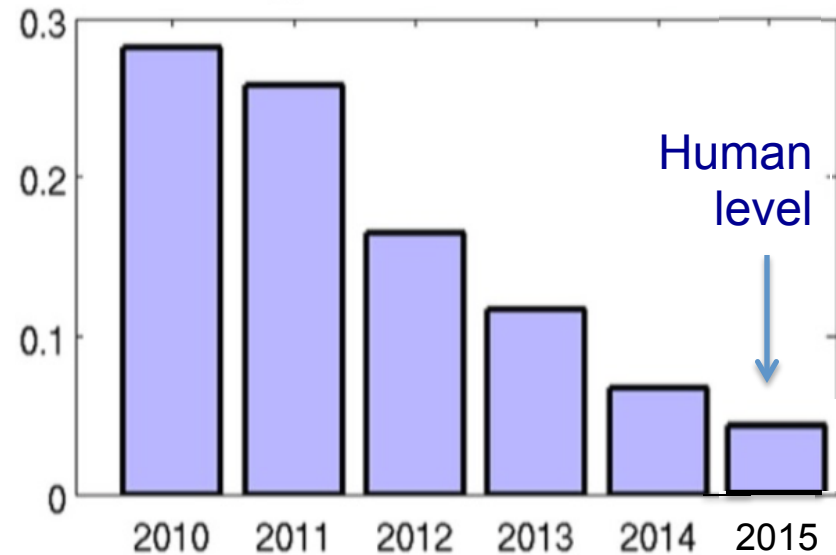
Convolutional Neural Nets

Computer Vision



Imagenet Visual Recognition

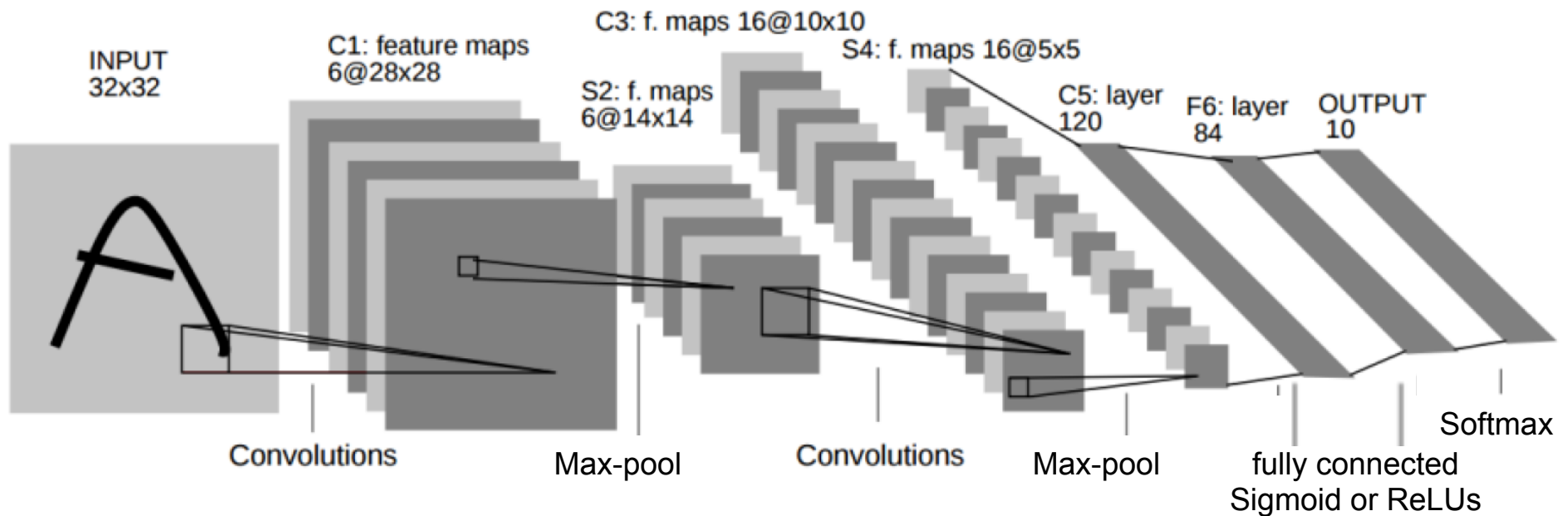
Error rate



Challenges:

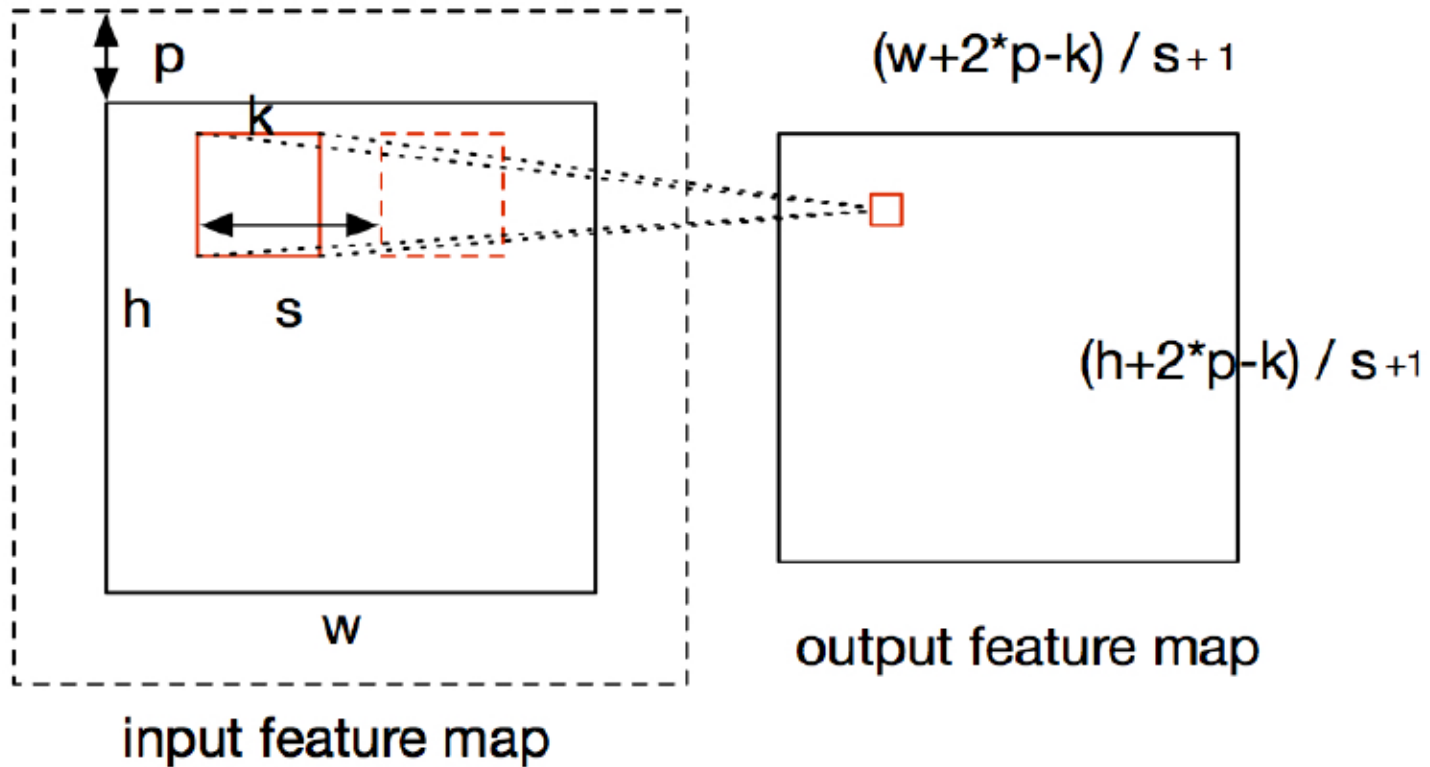
- invariance to translation, lighting, scaling, ...
- manual engineering of low-level image features
- ...

A Convolutional Neural Net for Handwritten Digit recognition: LeNet



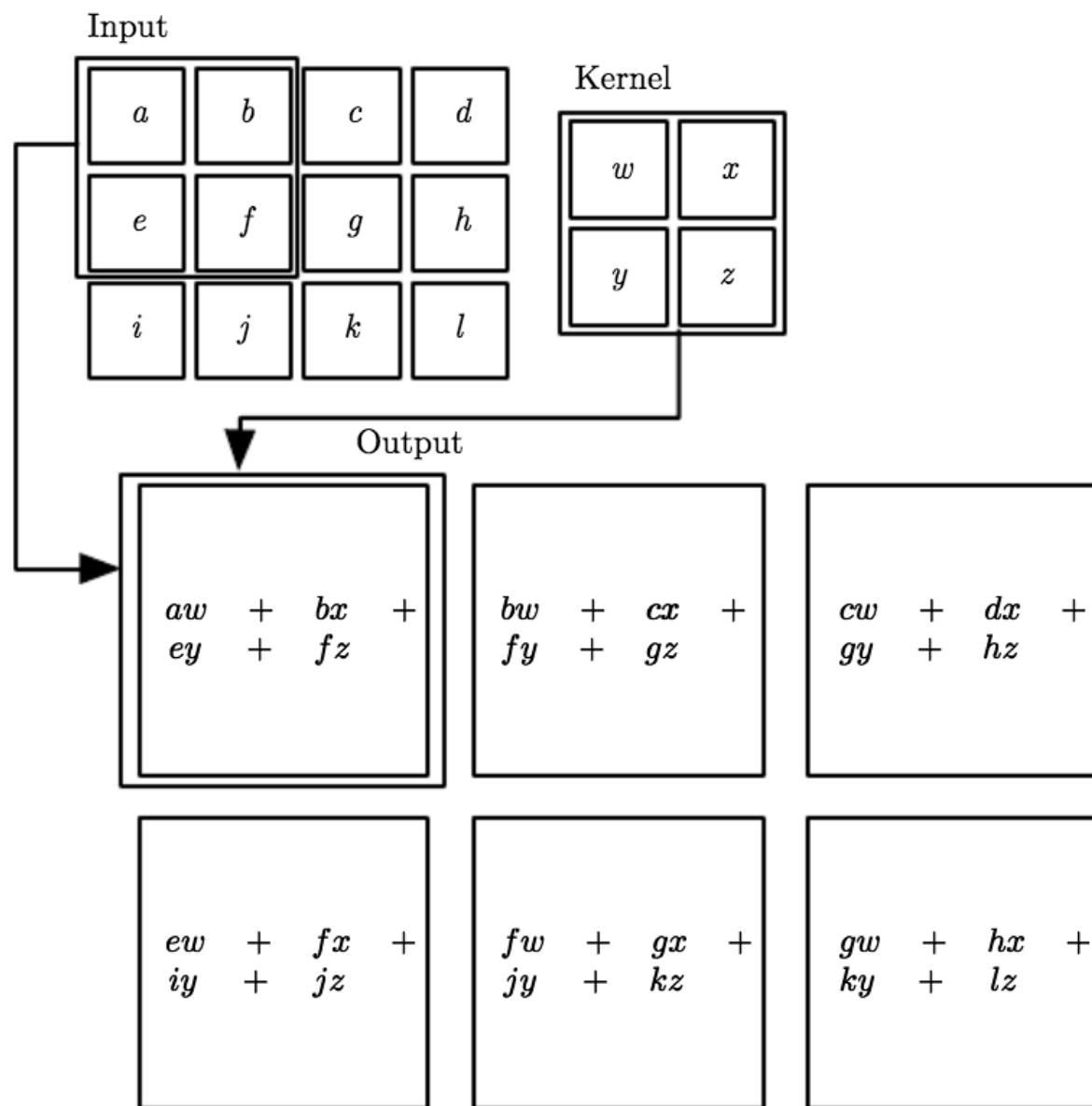
Convolution layer

p = padding
s = stride



$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(i + m, j + n) K(m, n)$$

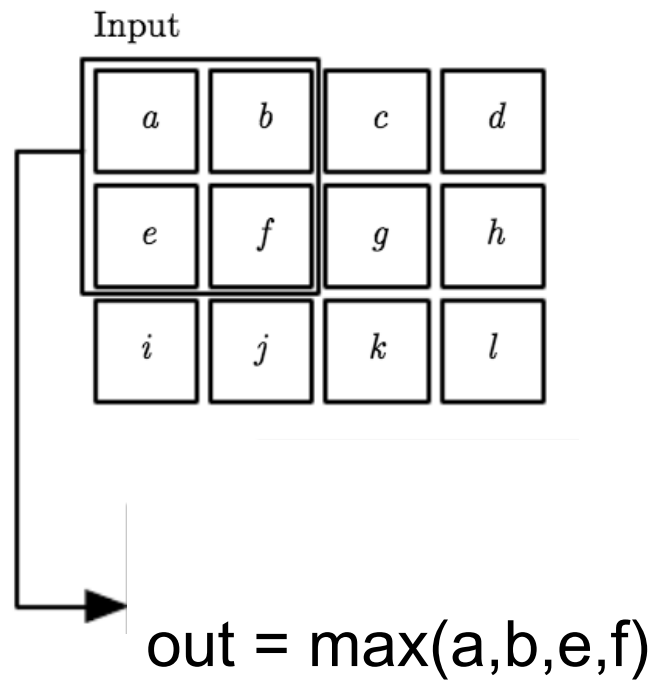
Convolution



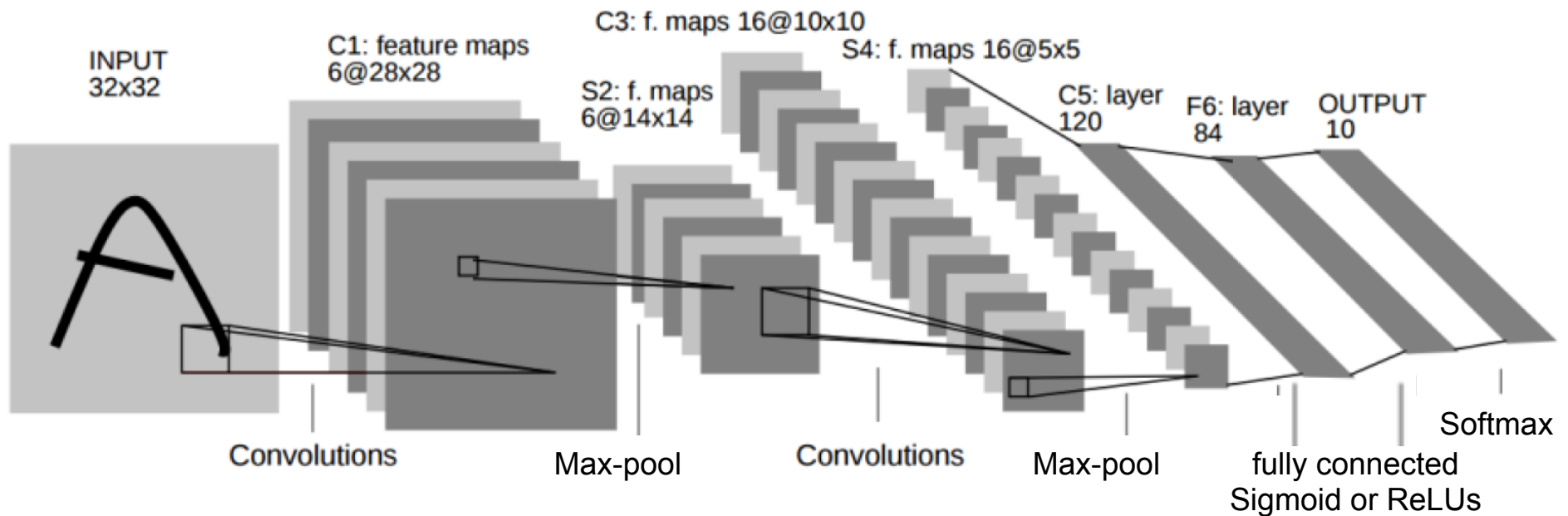
$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(i + m, j + n) K(m, n)$$

[from Goodfellow et al.]

Maxpool



A Convolutional Neural Net for Handwritten Digit recognition: LeNet



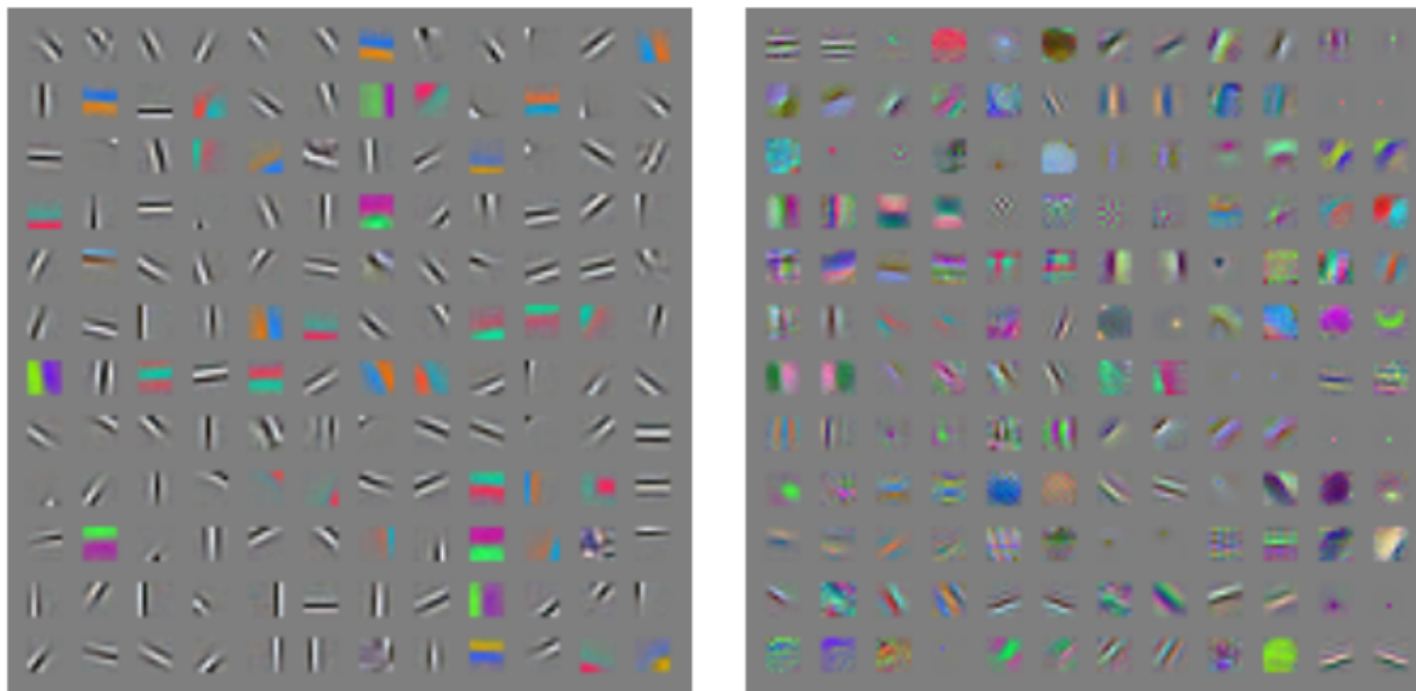


Figure 9.19: Many machine learning algorithms learn features that detect edges or specific colors of edges when applied to natural images. These feature detectors are reminiscent of the Gabor functions known to be present in primary visual cortex. (*Left*)Weights learned by an unsupervised learning algorithm (spike and slab sparse coding) applied to small image patches. (*Right*)Convolution kernels learned by the first layer of a fully supervised convolutional maxout network. Neighboring pairs of filters drive the same maxout unit.

[from Goodfellow et al.]

What you should know:

- Backpropagation algorithm
 - training network with gradient descent
 - how to derive gradient for simple cost functions, units
 - relationship to logistic regression
- Neural nets learn internal representations
- Convolutional neural networks
 - convolution operation
 - see new Homework

Many types of parameterized units

- Sigmoid activation units
- tanh activation units $\tanh(x) = \frac{(e^{2x} - 1)}{(e^{2x} + 1)}$
- ReLU
- Leaky ReLU (fixed non-zero slope for input < 0)
- Parametric ReLU (trainable slope for input < 0)
- Max Pool
- Inner Product
- GRU's
- LSTM's
- Matrix multiply
- no end in sight

Any unit $h(X; W)$ that is differentiable w.r.t. X and W