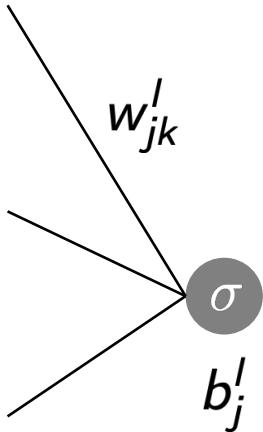


Computational MRI

Training neural networks with backpropagation

Neural networks



$$a_j^I = \sigma \left(\sum_k w_{jk}^I a_k^{I-1} + b_j^I \right)$$

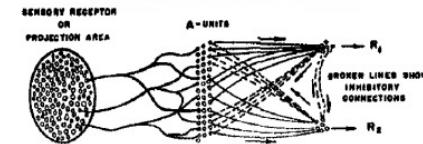


FIG. 2A. Schematic representation of connections in a simple perceptron.

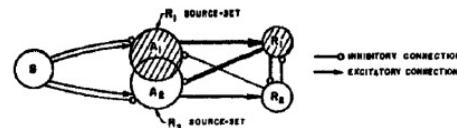
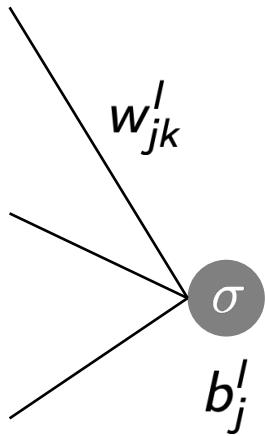


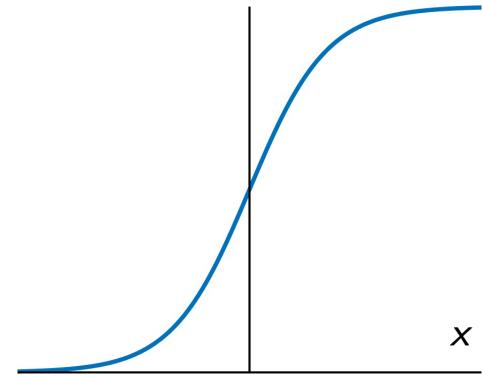
FIG. 2B. Venn diagram of the same perceptron (shading shows active sets for R_1 response).

Neural networks

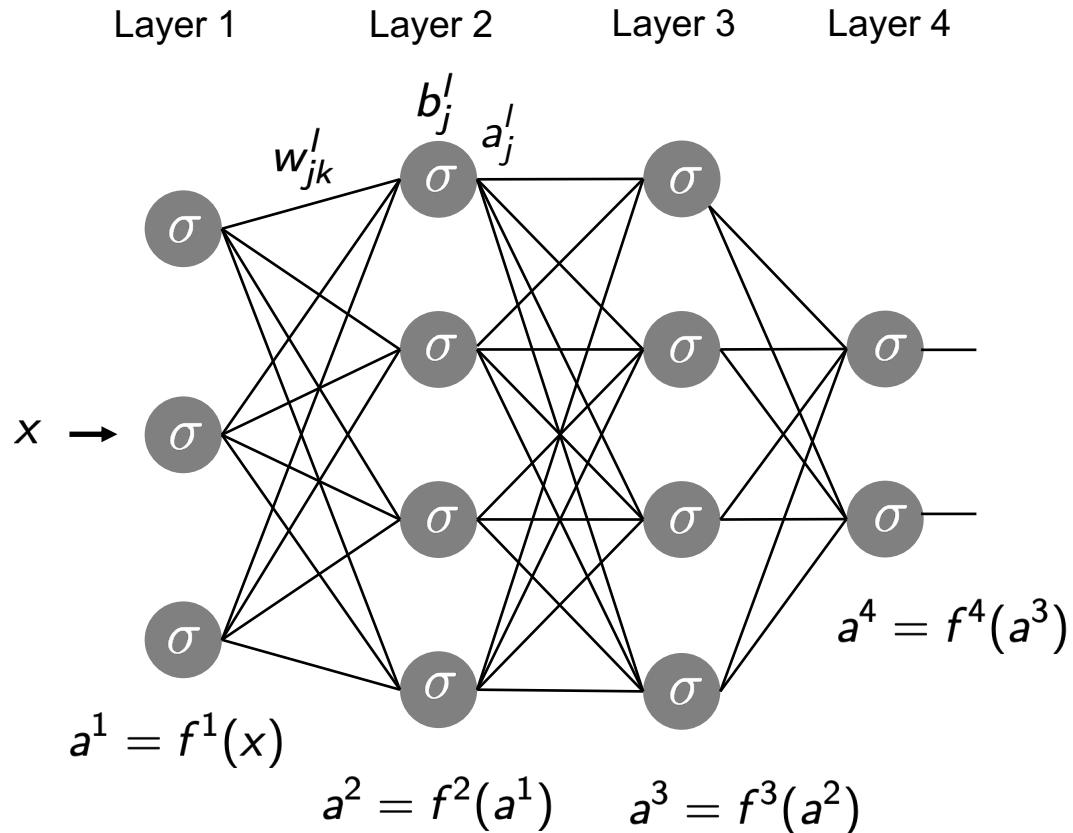


$$a_j^I = \sigma \left(\sum_k w_{jk}^I a_k^{I-1} + b_j^I \right)$$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



Neural networks



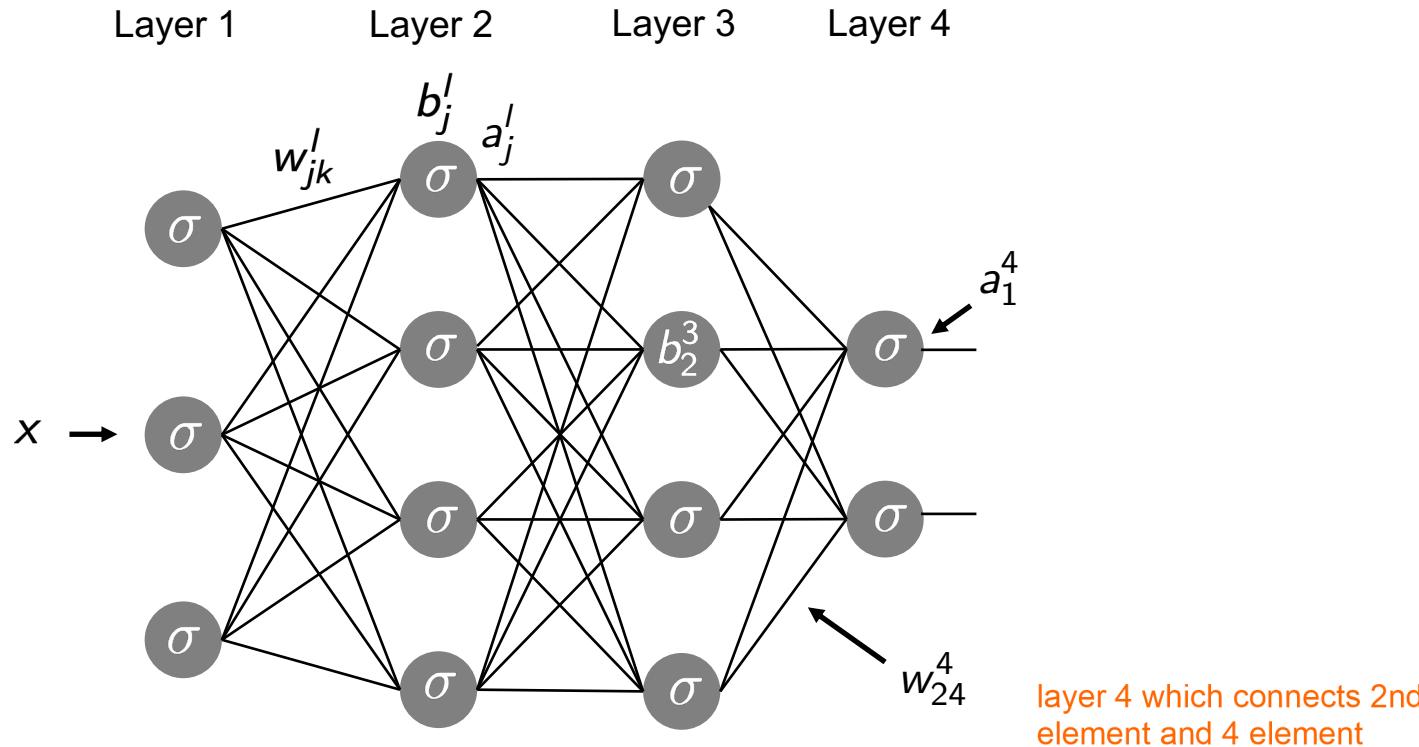
Large number of model parameters

High descriptive capacity

$$a^4 = F(x) = f^4(f^3(f^2(f^1(x))))$$

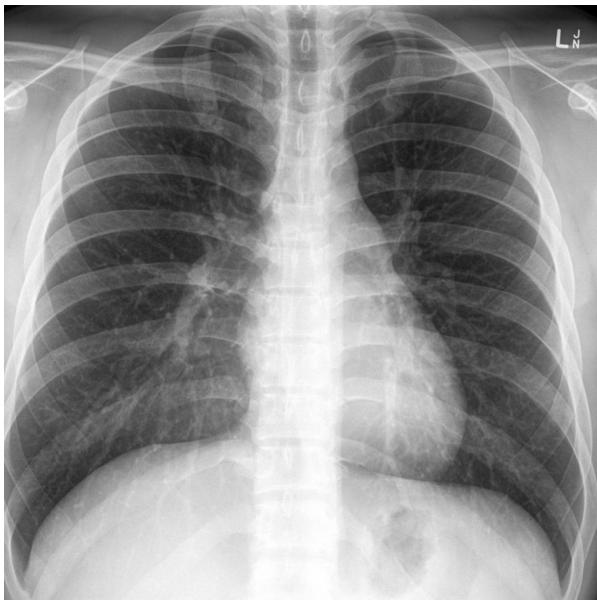
Cybenko: Math Control Sig Sys 1989

Neural networks: Notation



ER chest X-Ray diagnostic classification

Normal



COVID-19



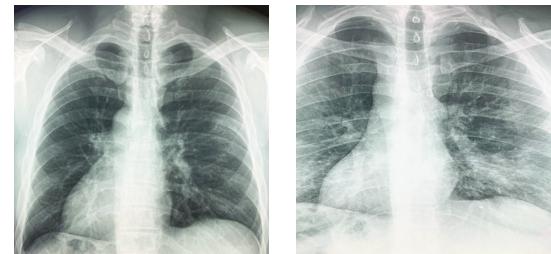
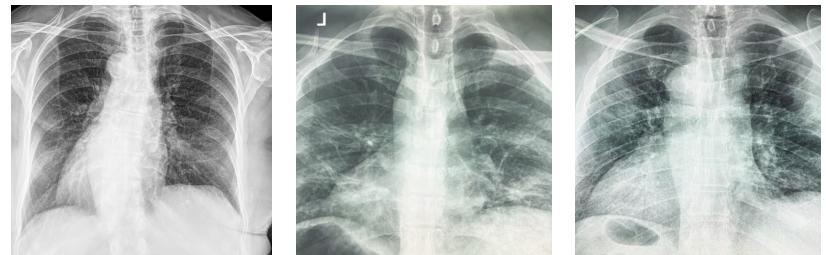
Chest X-Ray data set

Normal



$$y = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

COVID-19

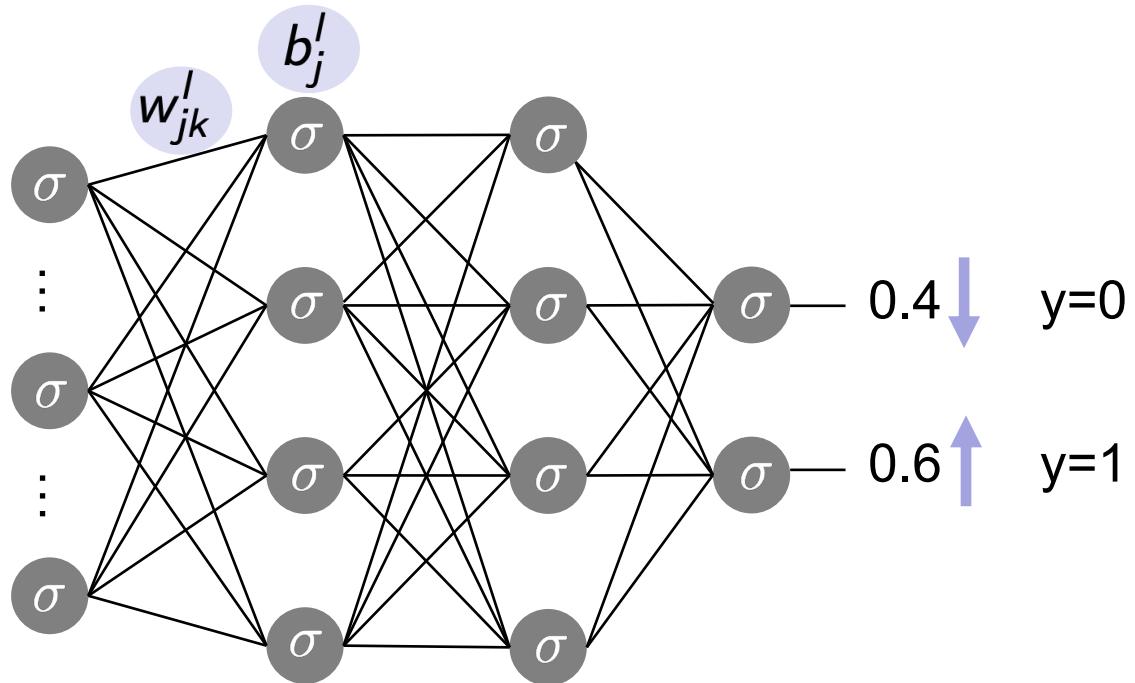


$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\{(x_1, y_1), \dots (x_N, y_N)\}$$

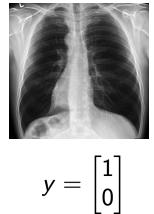
Neural network training

COVID-19



Change weights and biases to bring output closer to target

Neural network training: Cost function



$$y = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$



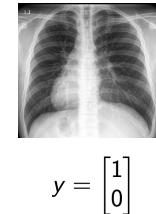
$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$



$$y = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$



$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$



$$y = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$



$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$



$$y = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$



$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$



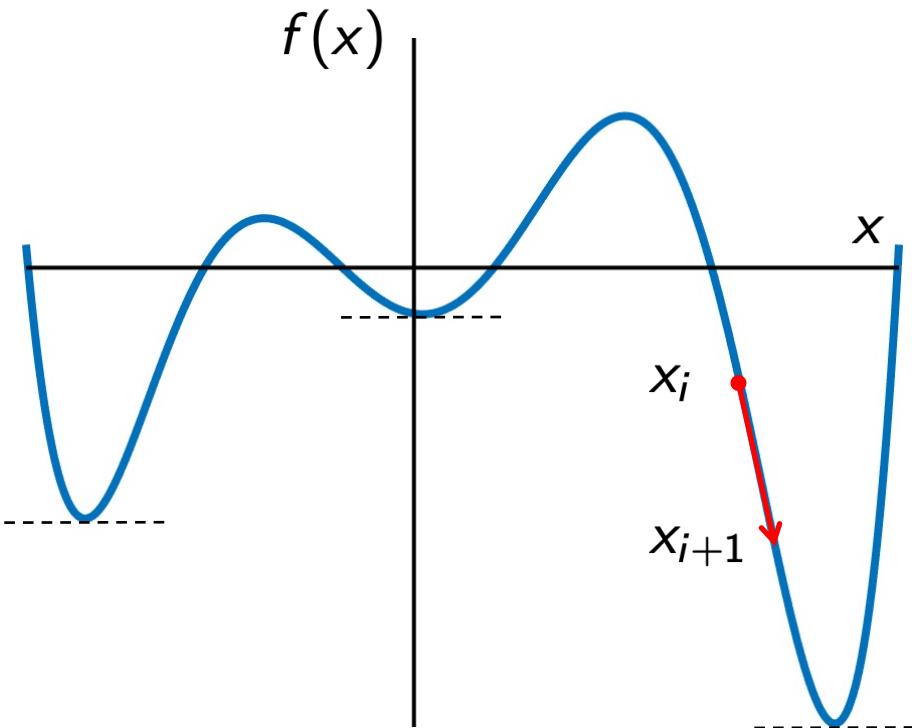
$$y = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$



$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

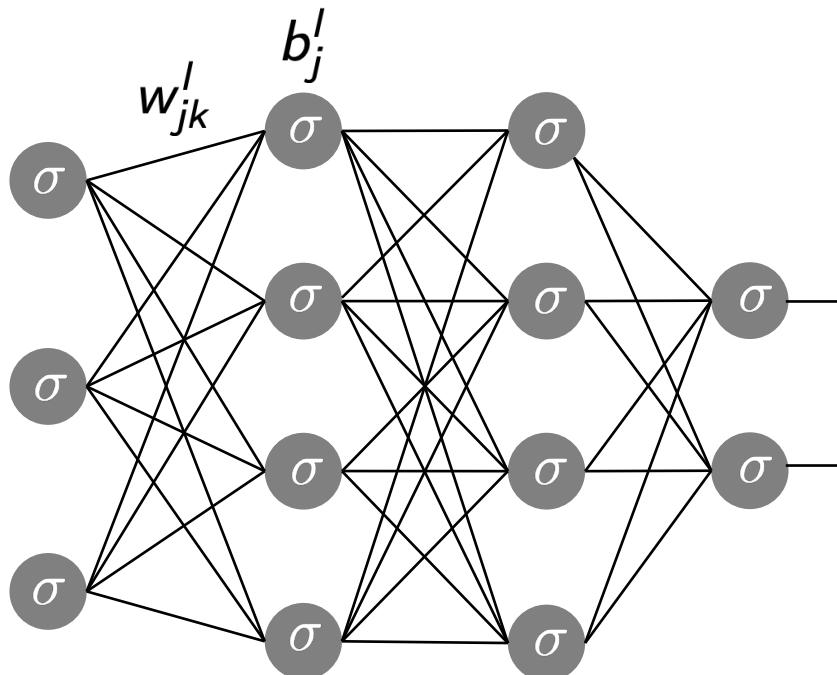
$$C(w, b) = \frac{1}{2N} \sum_{x=1}^N ||y_x - a_x||_2^2$$

Find minimum of function: Gradient descent



$$x_{i+1} = x_i - \alpha \frac{\partial f(x_i)}{\partial x_i}$$

Neural network training: Gradient descent



$$\min_{w,b} C(w, b)$$



$$\tilde{w}_{jk}^I = w_{jk}^I - \alpha \frac{\partial C}{\partial w_{jk}^I}$$

$$\tilde{b}_j^I = b_j^I - \alpha \frac{\partial C}{\partial b_j^I}$$

→ Chain rule

Backpropagation

**Learning representations
by back-propagating errors**

**David E. Rumelhart*, Geoffrey E. Hinton†
& Ronald J. Williams***

* Institute for Cognitive Science, C-015, University of California,
San Diego, La Jolla, California 92093, USA

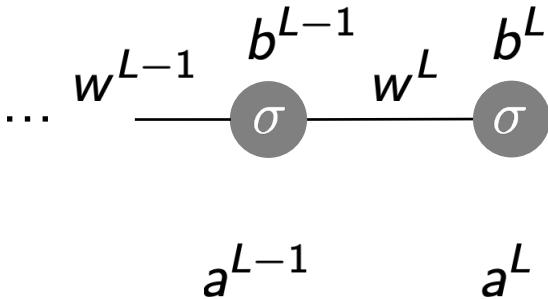
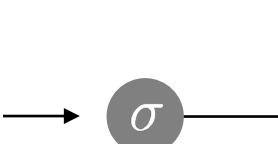
† Department of Computer Science, Carnegie-Mellon University,
Pittsburgh, Philadelphia 15213, USA

We describe a new learning procedure, back-propagation, for networks of neurone-like units. The procedure repeatedly adjusts the weights of the connections in the network so as to minimize a measure of the difference between the actual output vector of the net and the desired output vector. As a result of the weight adjustments, internal ‘hidden’ units which are not part of the input or output come to represent important features of the task domain, and the regularities in the task are captured by the interactions of these units. The ability to create useful new features distinguishes back-propagation from earlier, simpler methods such as the perceptron-convergence procedure¹.

$$\frac{\partial C}{\partial w_{jk}^l} \quad \frac{\partial C}{\partial b_j^l}$$

→ Efficient recursion to calculate these partial derivatives

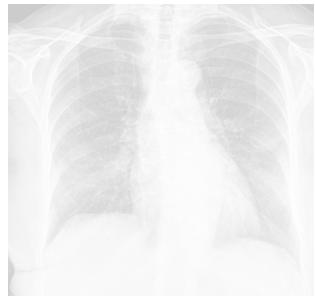
Backpropagation: Forward pass



$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$C_0 = \frac{1}{2}(y - a^L)^2$$

Backpropagation: Backward pass



$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$



$$a^1$$

...

$$w^{L-1}$$

$$\sigma$$

$$a^{L-1}$$

$$\sigma$$

$$w^L$$

$$a^L$$

$$a^L = \sigma \underbrace{\left(w^L a^{L-1} + b^L \right)}_{z^L}$$

$$C_0 = \frac{1}{2}(y - a^L)^2$$

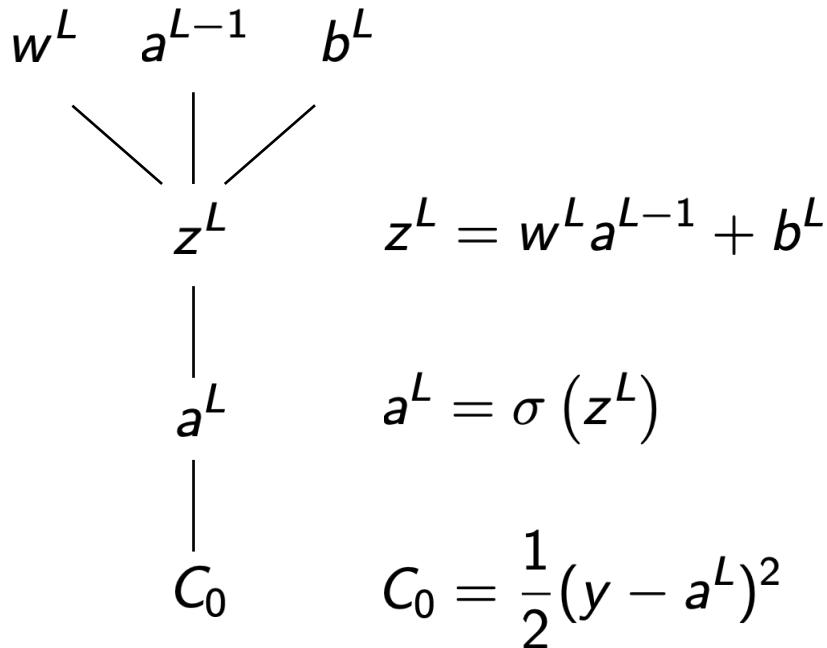
Backpropagation: Backward pass

$$\begin{array}{ccc} w^L & a^{L-1} & b^L \\ \backslash & | & / \\ z^L & & z^L = w^L a^{L-1} + b^L \\ | & & \\ a^L & & a^L = \sigma(z^L) \\ | & & \\ C_0 & & C_0 = \frac{1}{2}(y - a^L)^2 \end{array}$$

Partial derivatives

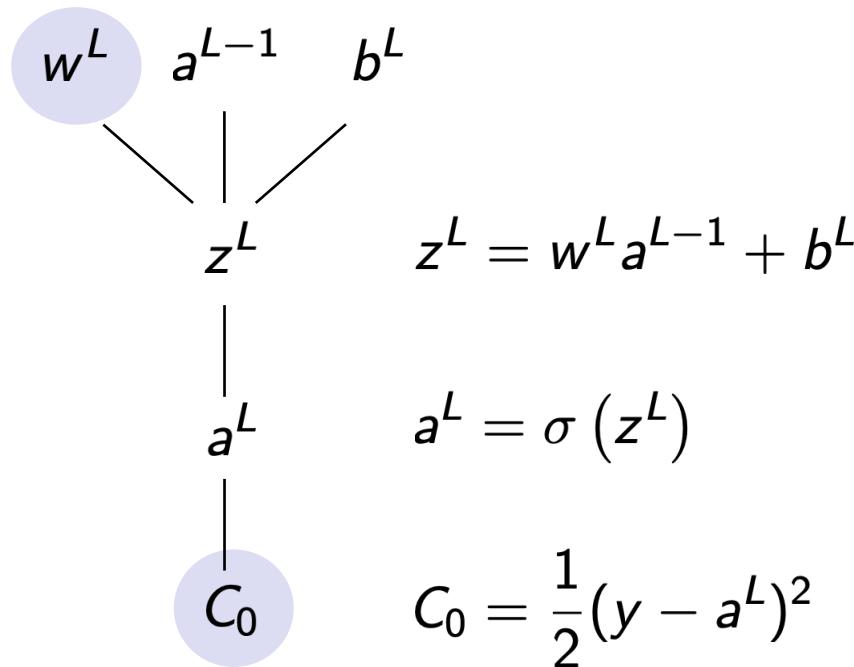
$$\frac{\partial C_0}{\partial w^L} = ?$$

$$\frac{\partial C_0}{\partial b^L} = ?$$



Backpropagation: Weights

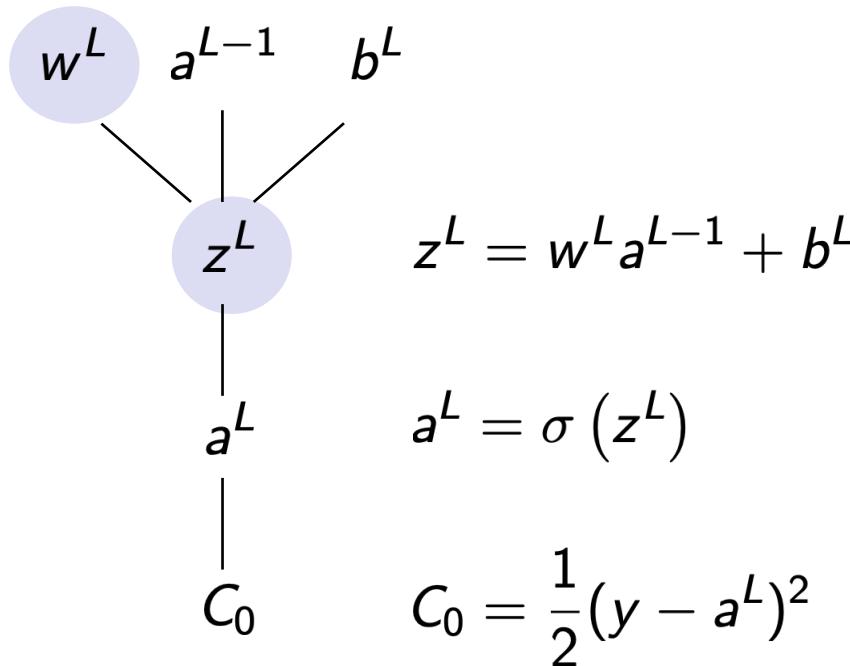
$$\frac{\partial C_0}{\partial w^L}$$



Backpropagation: Weights

$$\frac{\partial C_0}{\partial w^L} = \frac{\partial z^L}{\partial w^L}$$

$$\frac{\partial z^L}{\partial w^L} = a^{L-1}$$

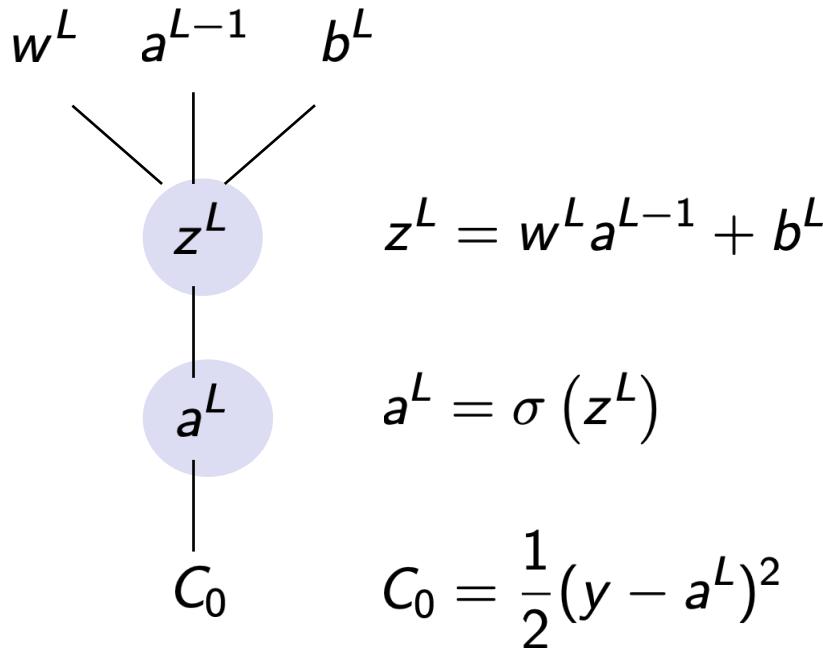


Backpropagation: Weights

$$\frac{\partial C_0}{\partial w^L} = \frac{\partial z^L}{\partial w^L} \frac{\partial a^L}{\partial z^L}$$

$$\frac{\partial z^L}{\partial w^L} = a^{L-1}$$

$$\frac{\partial a^L}{\partial z^L} = \sigma'(z^L)$$



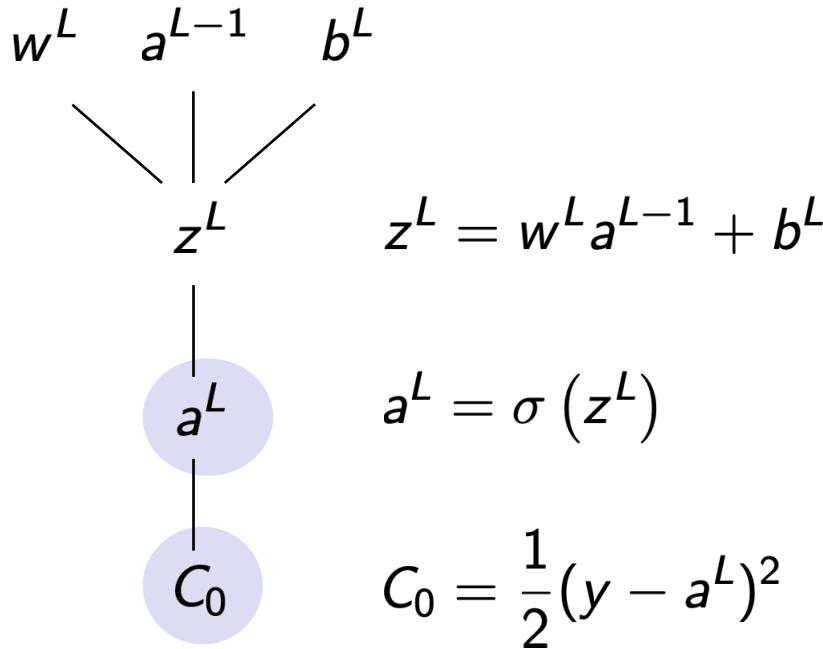
Backpropagation: Weights

$$\frac{\partial C_0}{\partial w^L} = \frac{\partial z^L}{\partial w^L} \frac{\partial a^L}{\partial z^L} \frac{\partial C_0}{\partial a^L}$$

$$\frac{\partial z^L}{\partial w^L} = a^{L-1}$$

$$\frac{\partial a^L}{\partial z^L} = \sigma'(z^L)$$

$$\frac{\partial C_0}{\partial a^L} = a^L - y$$



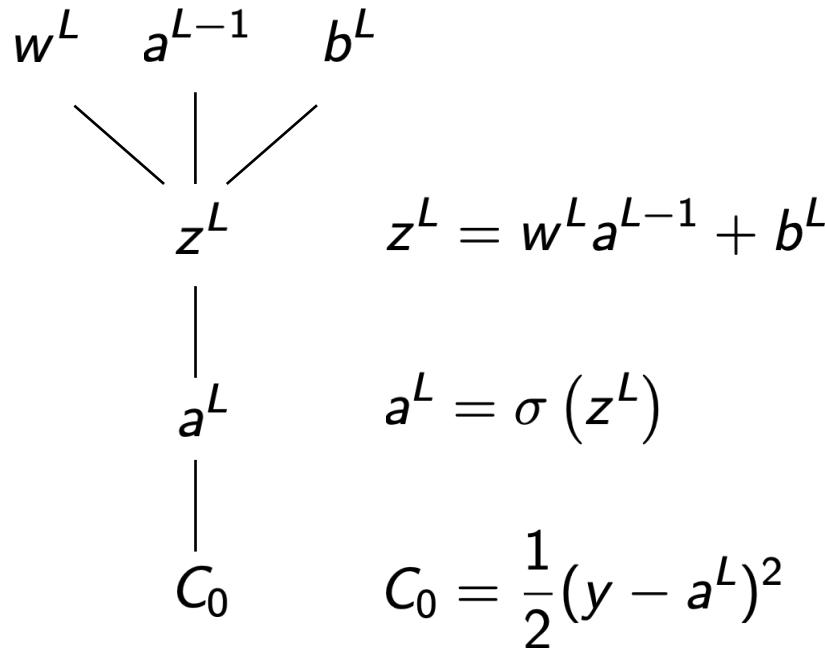
Backpropagation: Weights

$$\frac{\partial C_0}{\partial w^L} = \frac{\partial z^L}{\partial w^L} \frac{\partial a^L}{\partial z^L} \frac{\partial C_0}{\partial a^L} = a^{L-1} \sigma'(z^L) (a^L - y)$$

$$\frac{\partial z^L}{\partial w^L} = a^{L-1}$$

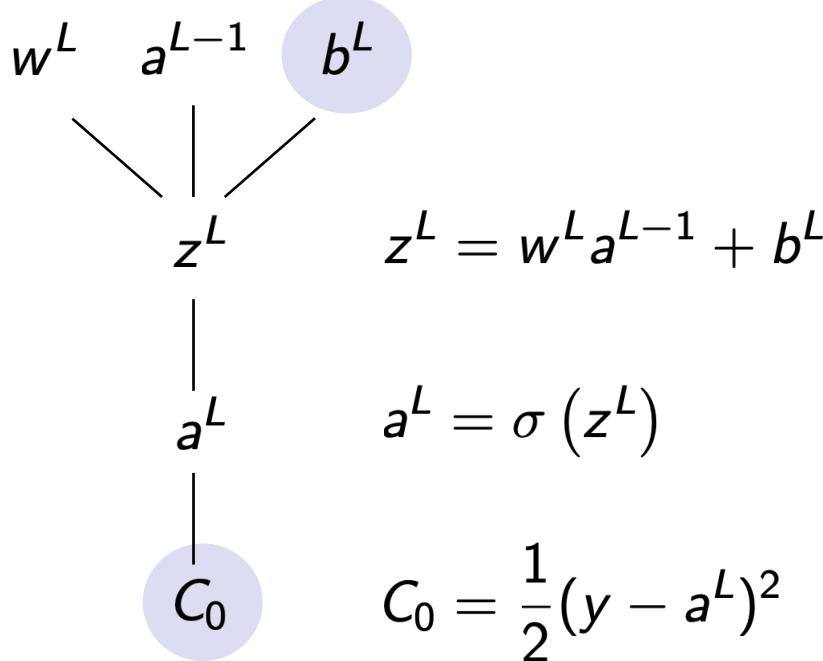
$$\frac{\partial a^L}{\partial z^L} = \sigma'(z^L)$$

$$\frac{\partial C_0}{\partial a^L} = a^L - y$$



Backpropagation: Bias

$$\frac{\partial C_0}{\partial b^L} = \frac{\partial a^L}{\partial z^L} \frac{\partial C_0}{\partial a^L} = \sigma'(z^L)(a^L - y)$$

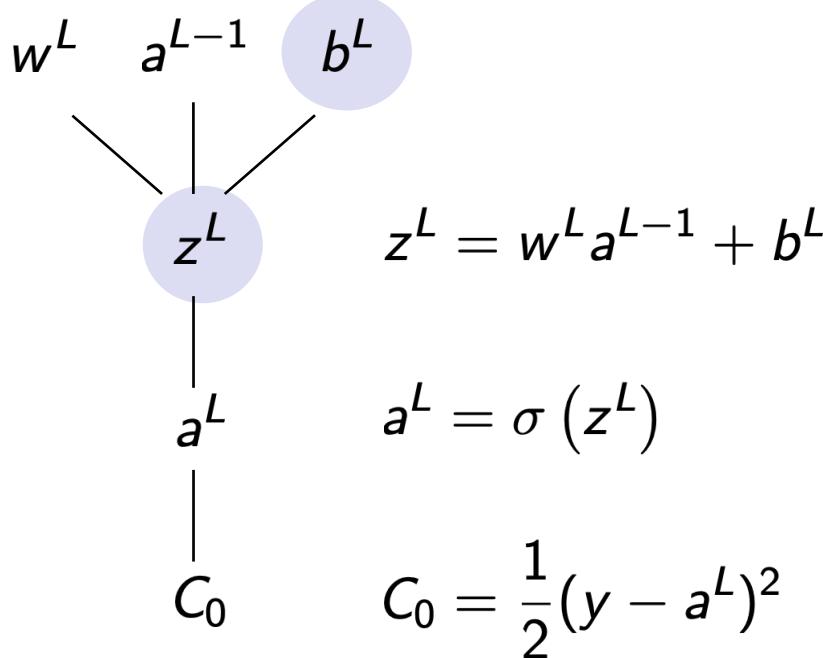


$$\frac{\partial a^L}{\partial z^L} = \sigma'(z^L)$$

$$\frac{\partial C_0}{\partial a^L} = a^L - y$$

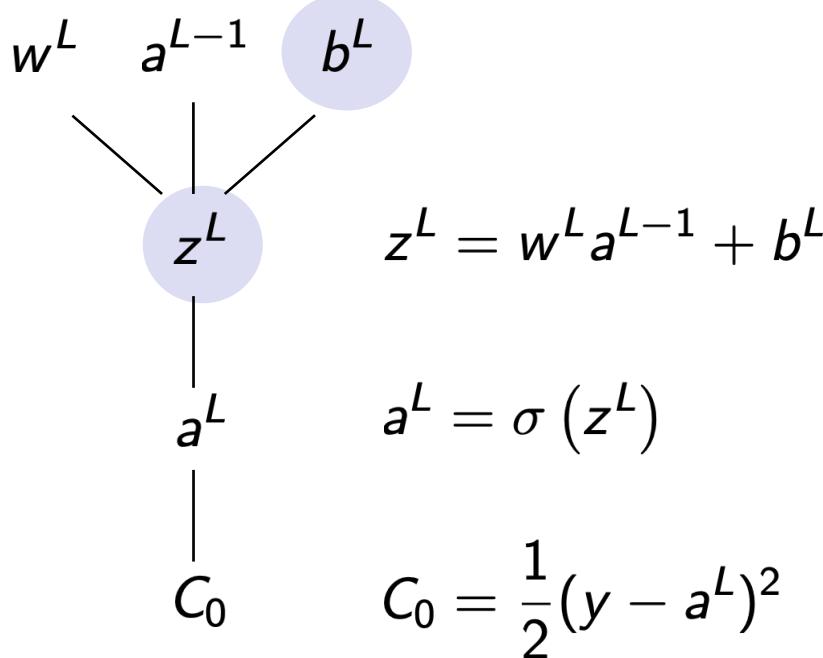
Backpropagation: Bias

$$\frac{\partial C_0}{\partial b^L} = \frac{\partial z^L}{\partial b^L} \frac{\partial a^L}{\partial z^L} \frac{\partial C_0}{\partial a^L} = \sigma'(z^L)(a^L - y)$$



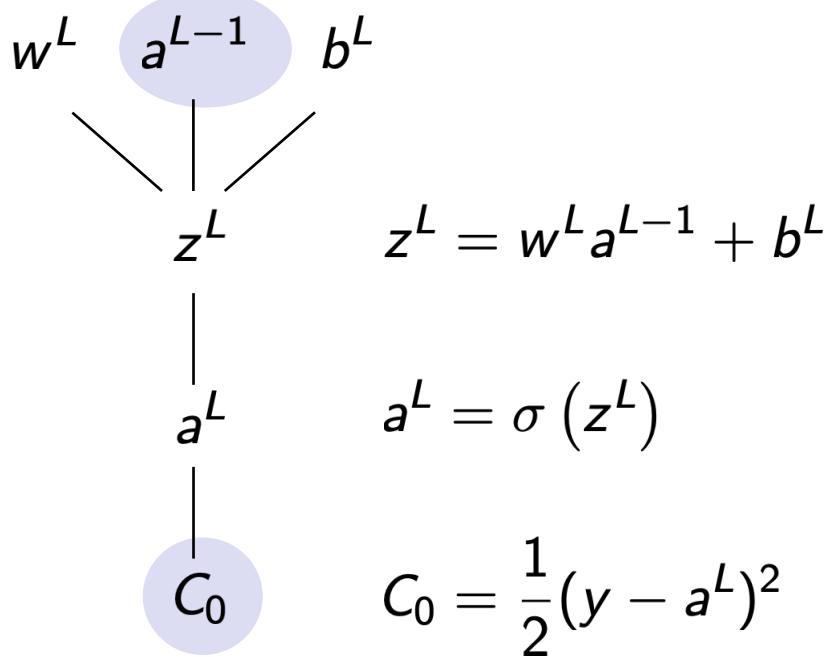
Backpropagation: Bias

$$\frac{\partial C_0}{\partial b^L} = \frac{\partial z^L}{\partial b^L} \frac{\partial a^L}{\partial z^L} \frac{\partial C_0}{\partial a^L} = 1\sigma'(z^L)(a^L - y)$$



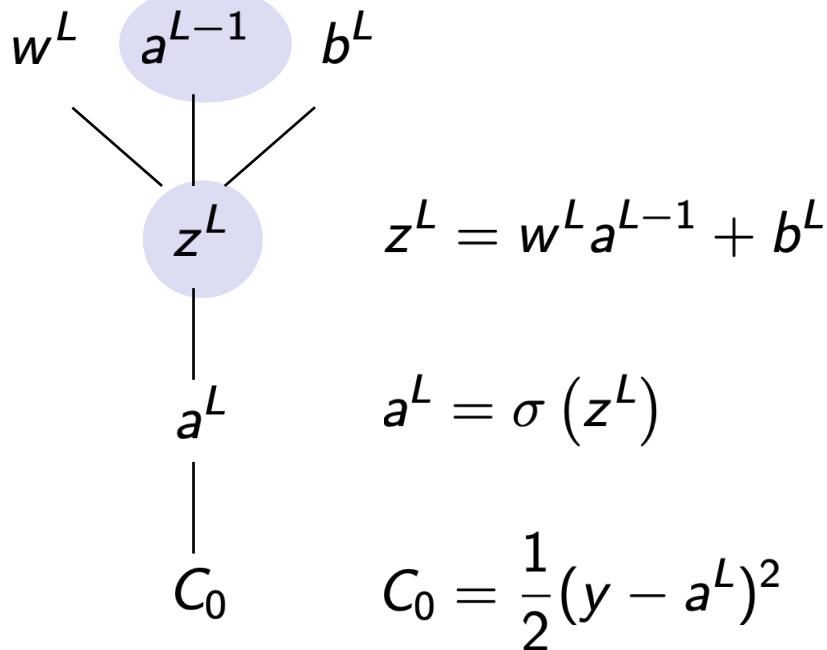
Backpropagation: Activation of previous layer

$$\frac{\partial C_0}{\partial a^{L-1}} = \frac{\partial a^L}{\partial z^L} \frac{\partial C_0}{\partial a^L} = \sigma'(z^L)(a^L - y)$$



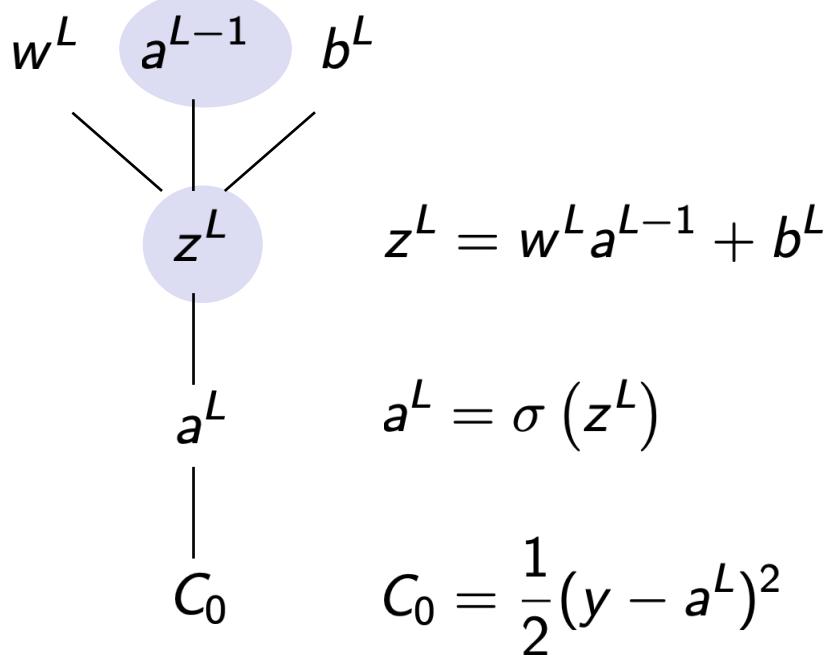
Backpropagation: Activation of previous layer

$$\frac{\partial C_0}{\partial a^{L-1}} = \frac{\partial z^L}{\partial a^{L-1}} \frac{\partial a^L}{\partial z^L} \frac{\partial C_0}{\partial a^L} = \sigma'(z^L)(a^L - y)$$



Backpropagation: Activation of previous layer

$$\frac{\partial C_0}{\partial a^{L-1}} = \frac{\partial z^L}{\partial a^{L-1}} \frac{\partial a^L}{\partial z^L} \frac{\partial C_0}{\partial a^L} = w^L \sigma'(z^L) (a^L - y)$$



Backpropagation

$$\frac{\partial C_0}{\partial w^L} = a^{L-1} \sigma'(z^L) \frac{\partial C_0}{\partial a^L}$$

$$\frac{\partial C_0}{\partial b^L} = \sigma'(z^L) \frac{\partial C_0}{\partial a^L}$$

$$\frac{\partial C_0}{\partial a^L} = a^L - y$$



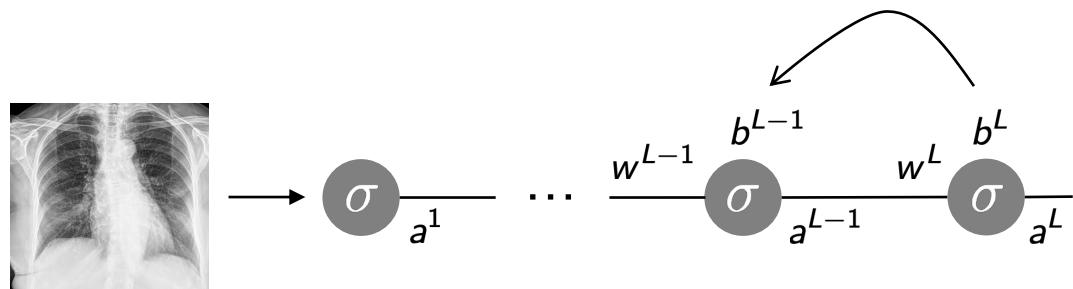
$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

Backpropagation

$$\frac{\partial C_0}{\partial w^{L-1}} = a^{L-2} \sigma'(z^{L-1}) \frac{\partial C_0}{\partial a^{L-1}}$$

$$\frac{\partial C_0}{\partial b^{L-1}} = \sigma'(z^{L-1}) \frac{\partial C_0}{\partial a^{L-1}}$$

$$\frac{\partial C_0}{\partial a^L} = w^L \sigma'(z^L) \frac{\partial C_0}{\partial a^L}$$



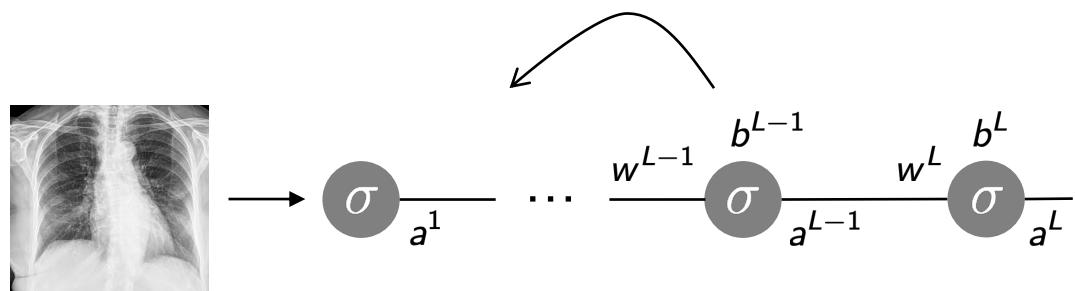
$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

Backpropagation

$$\frac{\partial C_0}{\partial w^l} = a^{l-1} \sigma'(z^l) \frac{\partial C_0}{\partial a^l}$$

$$\frac{\partial C_0}{\partial b^l} = \sigma'(z^l) \frac{\partial C_0}{\partial a^l}$$

$$\frac{\partial C_0}{\partial a^l} = w^{l+1} \sigma'(z^{l+1}) \frac{\partial C_0}{\partial a^{l+1}}$$



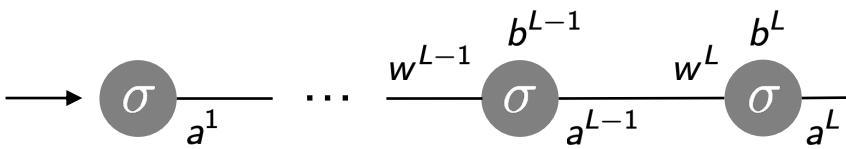
$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

Plug partial derivatives into gradient descent

$$\frac{\partial C_0}{\partial w^l} = a^{l-1} \sigma'(z^l) \frac{\partial C_0}{\partial a^l}$$

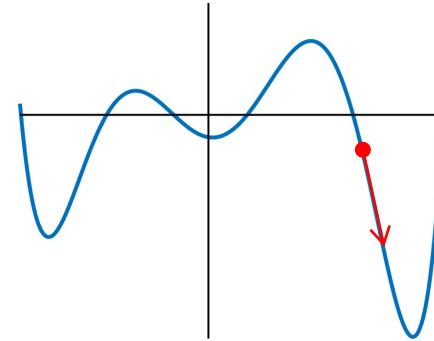
$$\frac{\partial C_0}{\partial b^l} = \sigma'(z^l) \frac{\partial C_0}{\partial a^l}$$

$$\frac{\partial C_0}{\partial a^l} = w^{l+1} \sigma'(z^{l+1}) \frac{\partial C_0}{\partial a^{l+1}}$$



$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\min_{w,b} C(w, b)$$



Gradient descent

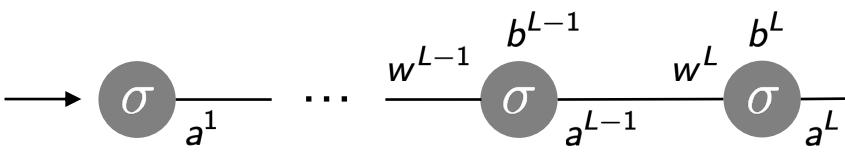
Loop over training examples

$$\frac{\partial C_0}{\partial w^l} = a^{l-1} \sigma'(z^l) \frac{\partial C_0}{\partial a^l}$$

$$\frac{\partial C_0}{\partial b^l} = \sigma'(z^l) \frac{\partial C_0}{\partial a^l}$$

$$\frac{\partial C_0}{\partial a^l} = w^{l+1} \sigma'(z^{l+1}) \frac{\partial C_0}{\partial a^{l+1}}$$

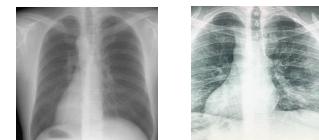
$$\min_{w,b} C(w, b)$$



$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$



$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$



$$y = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

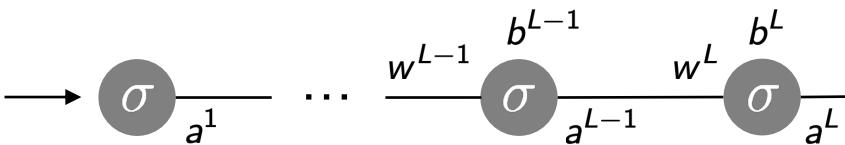
...

Backpropagation: Efficiency and insights

$$\frac{\partial C_0}{\partial w^l} = a^{l-1} \sigma'(z^l) \frac{\partial C_0}{\partial a^l}$$

$$\frac{\partial C_0}{\partial b^l} = \sigma'(z^l) \frac{\partial C_0}{\partial a^l}$$

$$\frac{\partial C_0}{\partial a^l} = w^{l+1} \sigma'(z^{l+1}) \frac{\partial C_0}{\partial a^{l+1}}$$



$$y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

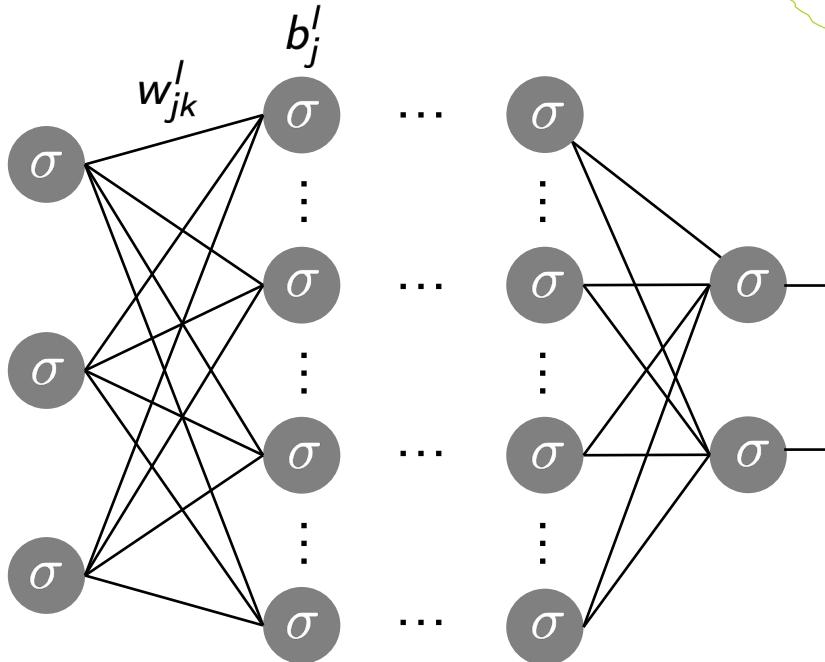
- Each computation involves just two layers
- Avoid recomputing identical expressions in the chain rule
- Gradients provide insight into what determines speed of learning

Backpropagation: General formulation

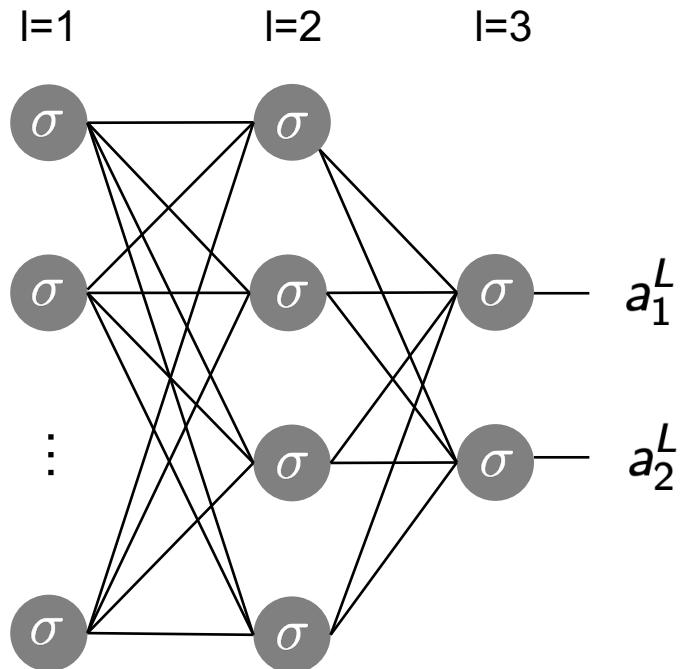
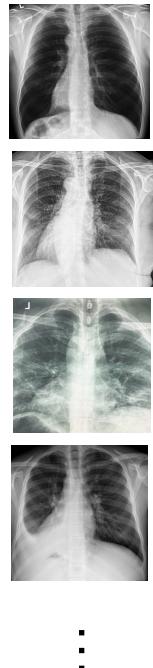
$$\frac{\partial C_0}{\partial w_{jk}^I} = a_k^{I-1} \sigma'(z_j^I) \frac{\partial C_0}{\partial a_j^I}$$

$$\frac{\partial C_0}{\partial b_j^I} = \sigma'(z_j^I) \frac{\partial C_0}{\partial a_j^I}$$

$$\frac{\partial C_0}{\partial a_j^I} = \sum_j w_{jk}^{I+1} \sigma'(z_j^{I+1}) \frac{\partial C_0}{\partial a_j^{I+1}}$$



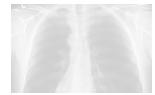
Back to our chest X-ray problem



$$\min_{w,b} C(w, b)$$

Accuracy

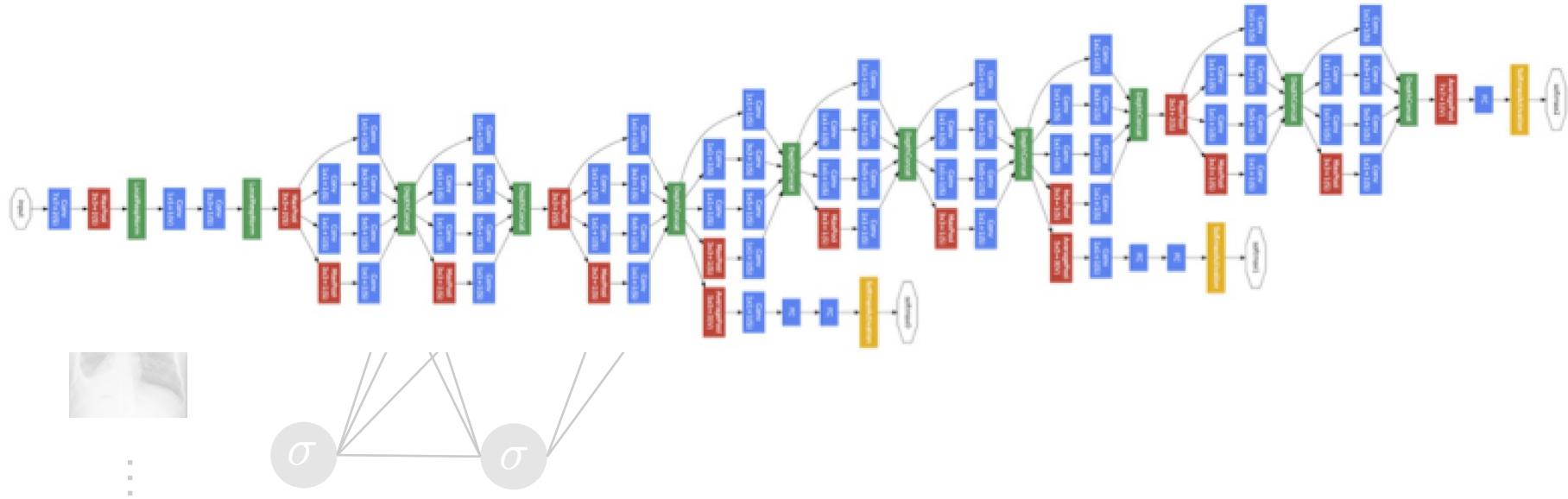
Back to our chest X-ray problem



|=1

|=2

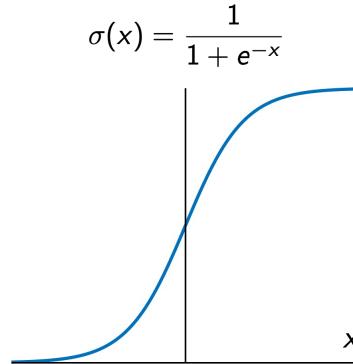
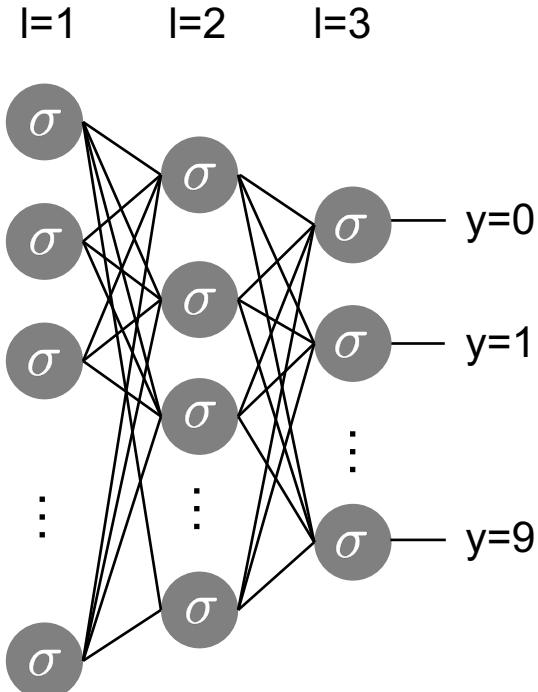
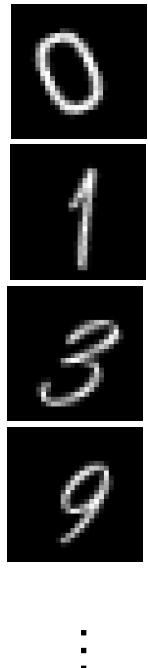
|=3



Handwritten character recognition: MNIST

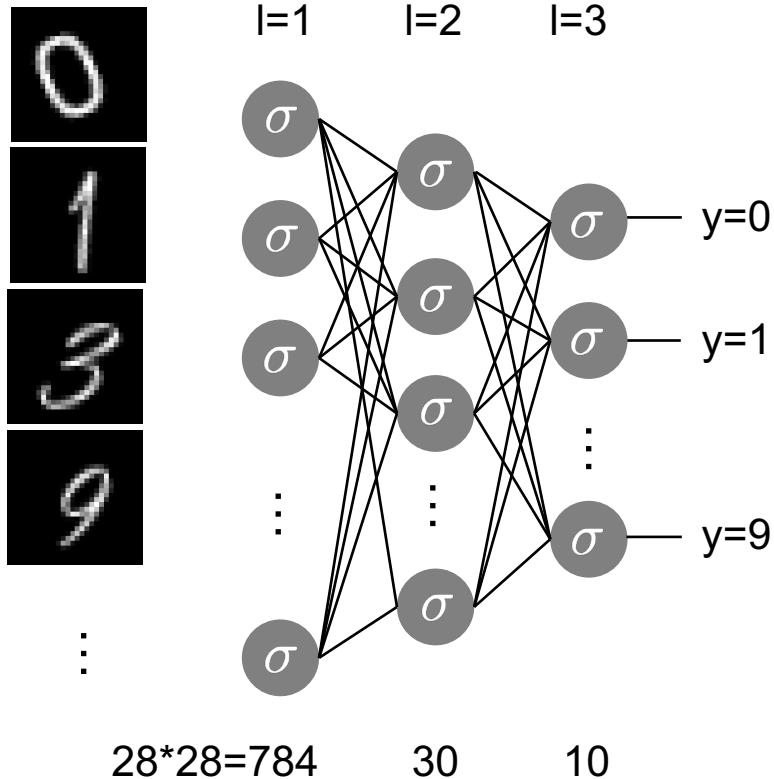
2 1 9 0 6 9 8 0 / 1
0 9 9 0 4 9 8 4 6 0
7 8 8 3 1 1 2 4 7 4
1 3 8 6 3 5 1 5 9 8
0 5 3 1 4 1 3 0 0 2
1 1 3 9 2 3 8 6 5 7

Adding more layers to our network?

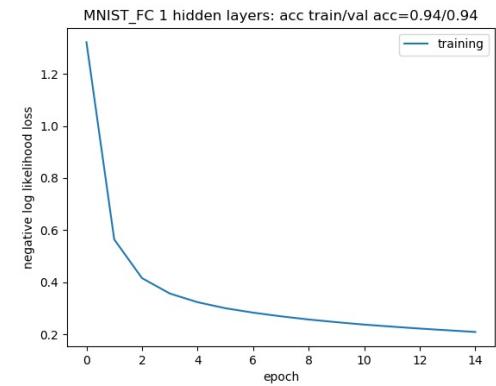


Sigmoid nonlinearity

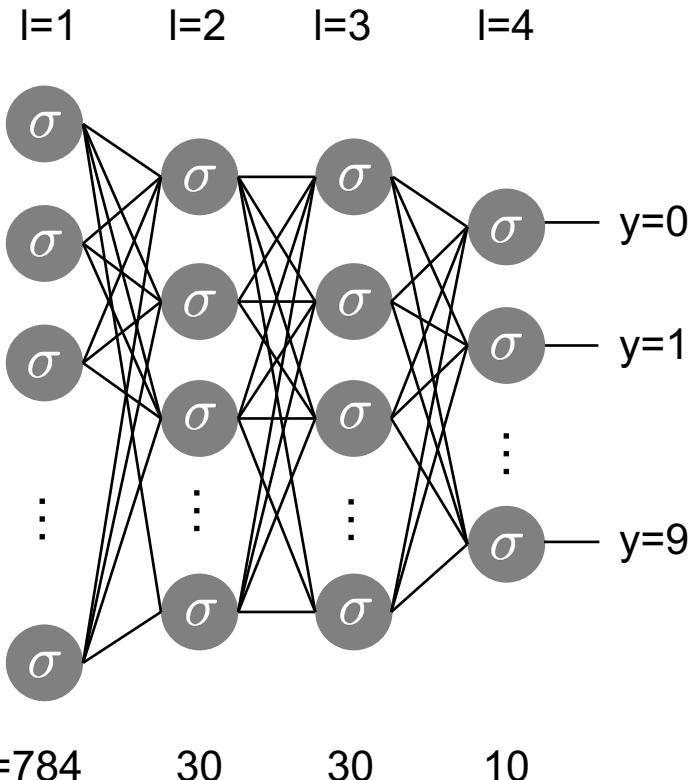
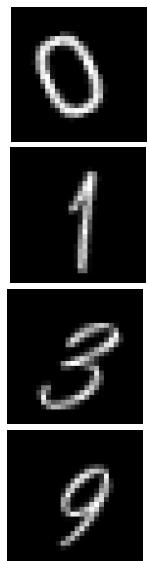
Adding more layers to our network?



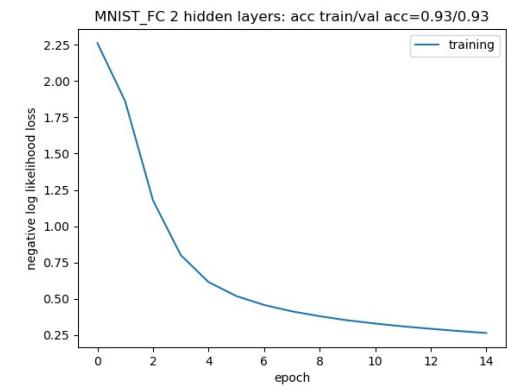
$$\min_{w,b} C(w, b)$$



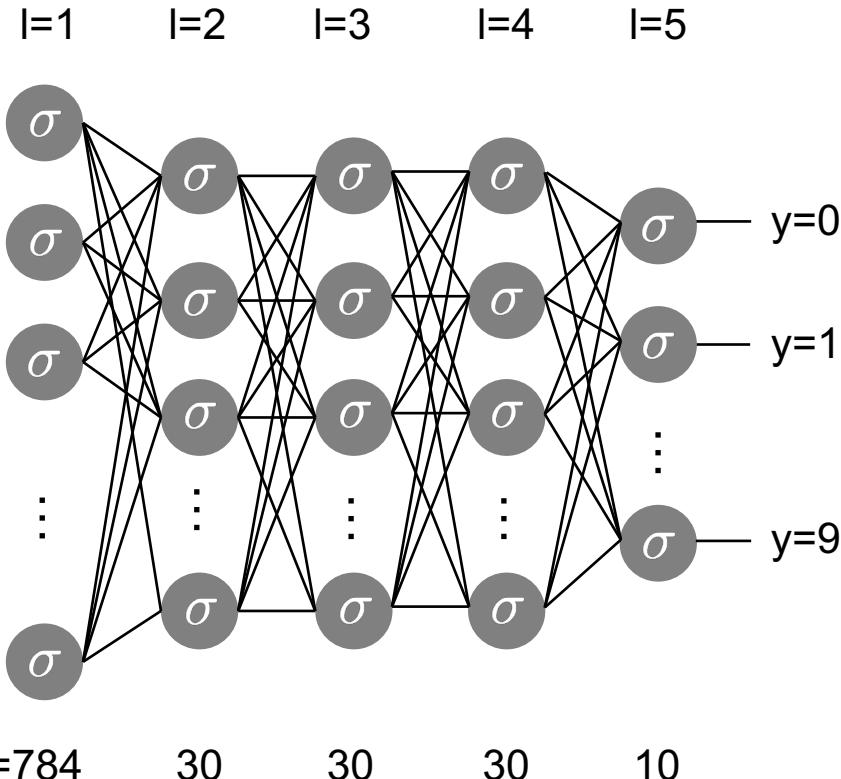
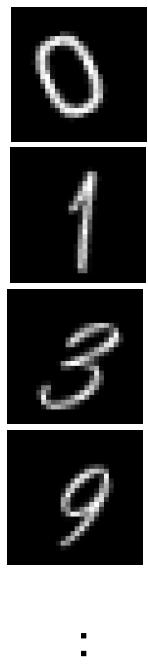
Adding more layers to our network?



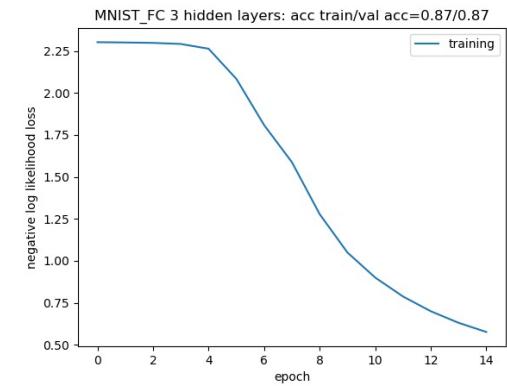
$$\min_{w,b} C(w, b)$$



Adding more layers to our network?

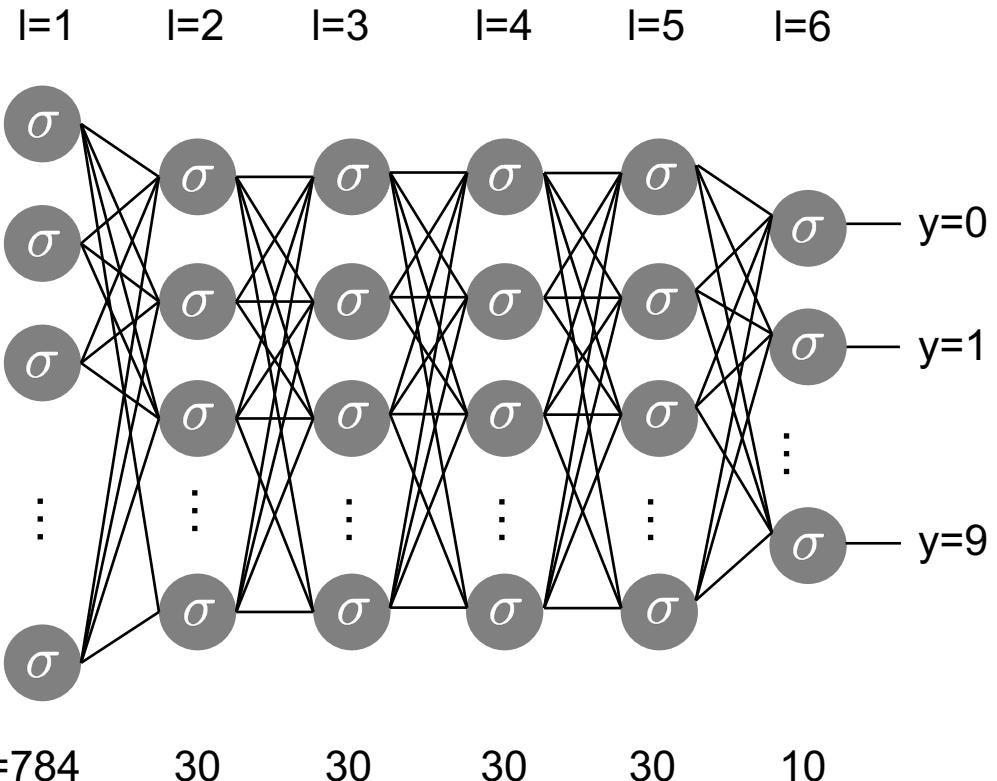
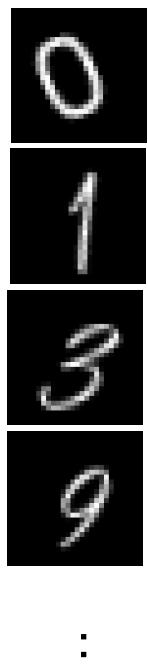


$$\min_{w,b} C(w, b)$$

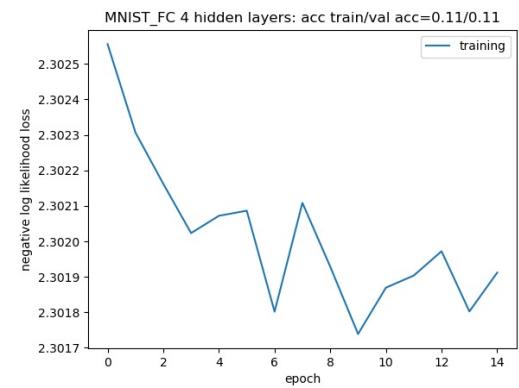


$$28 \times 28 = 784$$

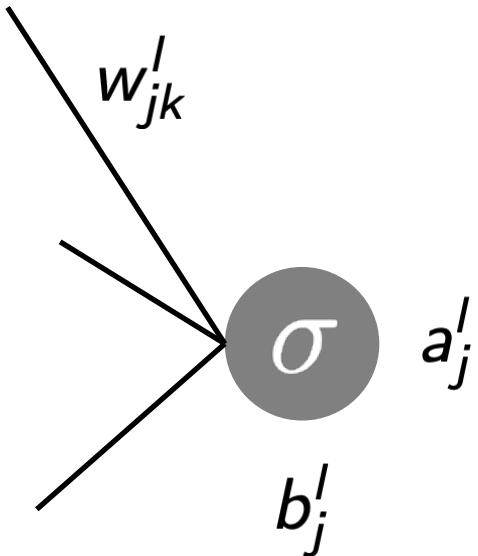
Adding more layers to our network?



$$\min_{w,b} C(w, b)$$

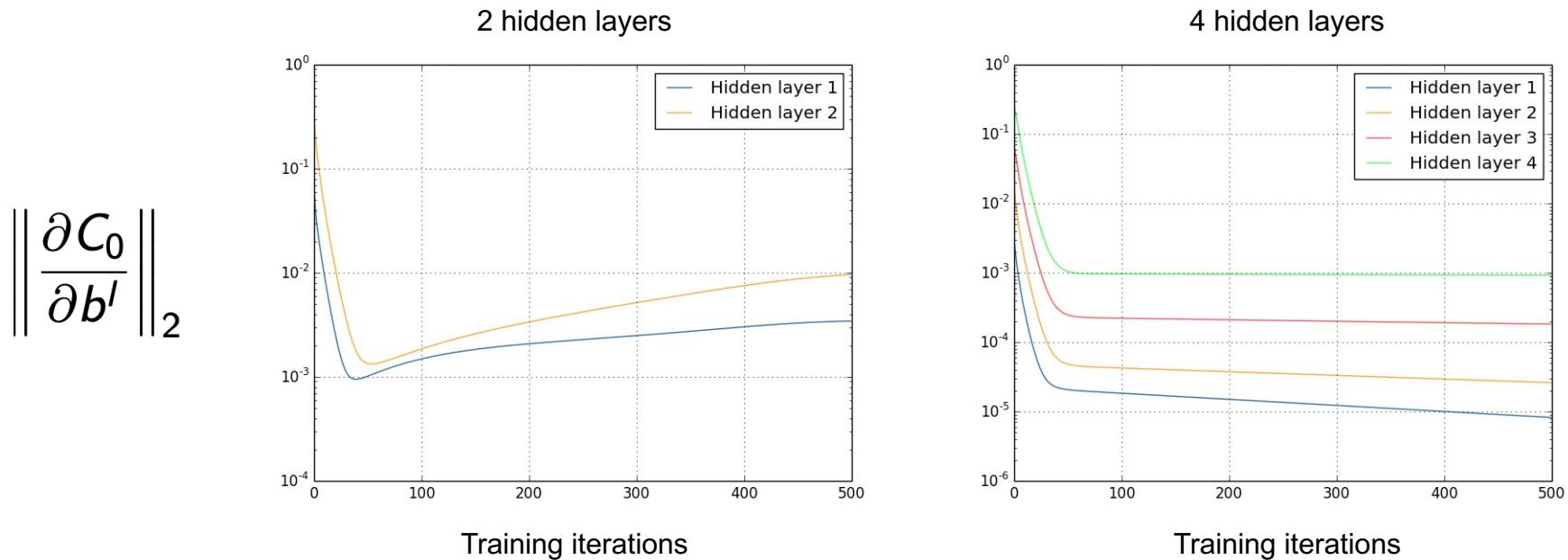


Insights from backpropagation gradients

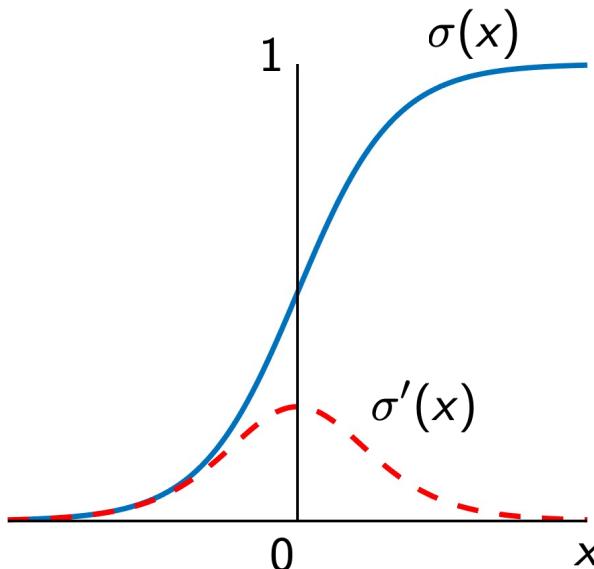


$$\left\| \frac{\partial C_0}{\partial b'} \right\|_2 = \left\| \sigma'(z') \frac{\partial C_0}{\partial a'} \right\|_2$$

Strength of gradients determines how fast neurons learn

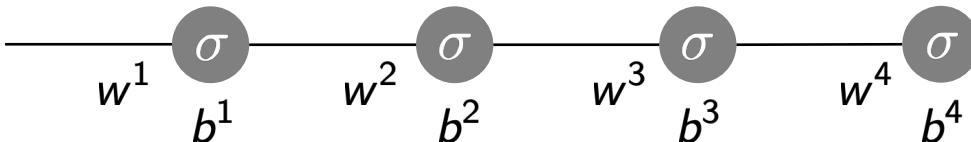


Vanishing gradient problem



$$\frac{\partial C_0}{\partial b_j^l} = \sigma'(z_j^l) \frac{\partial C_0}{\partial a_j^l}$$

$$\frac{\partial C_0}{\partial b^1} = \sigma'(z^1) \cdot w^2 \cdot \sigma'(z^2) \cdot w^3 \cdot \sigma'(z^3) \cdot w^4 \cdot \sigma'(z^3) \cdot \frac{\partial C_0}{\partial a^4}$$

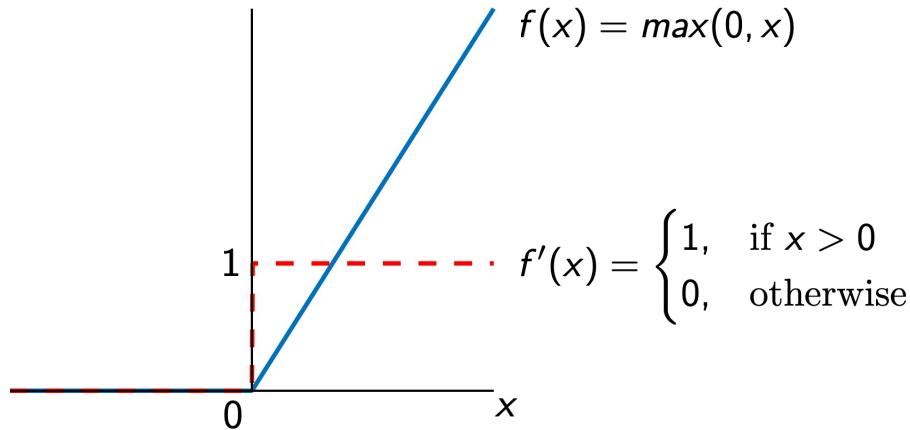


max(sigmoid gradient)=1/4, 太小了,
层数越多, 累乘起来越接近0, 导致梯度消失

Neurons get adjusted at very different rates → Gradient instability

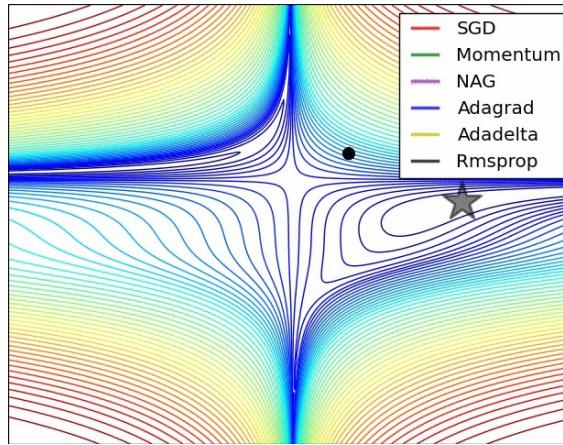
Ongoing research

- Nonlinearities: ReLU and variants,...



Ongoing research

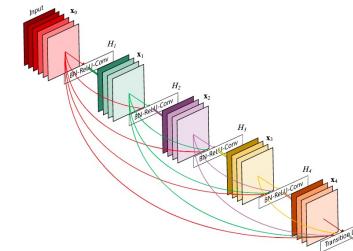
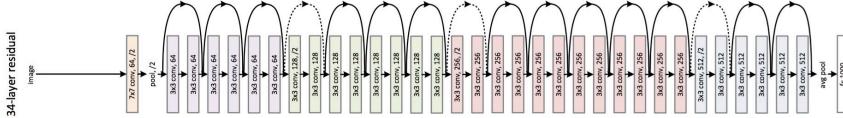
- Nonlinearities
- Optimization algorithms (ADAM, Adagrad, RMSProp,...)



Kingma 2014
Duchi 2011
Hinton 2013
Ruder 2016

Ongoing research

- Nonlinearities
- Optimization algorithms (ADAM, Adagrad, RMSProp,...)
- Network architectures



Summary

- Short recap of neural networks
- Training neural networks with gradient descent
- Backpropagation: Efficient implementation of chain rule
- Deep neural networks: Gradient instability

Exercise

Computational MR imaging

Laboratory 10: Training neural networks with backpropagation

Report is due on Wednesday the week after the lab session at 23:59. Send your report by email to Bruno Riemenschneider (bruno.riemenschneider@fau.de) and Florian Knoll (florian.knoll@fau.de).

Learning objectives

- Train a fully connected network for MNIST classification in Pytorch
- Change the number of hidden layers from 1 to 5
- Observe the
- Examine effects of dropping data on neural network performance
- Learn about effects of over and underfitting

1. Training a neural network for MNIST classification

Open the script `MNIST_classification_pytorch_1p5_cpu.py`. This script serves as a template that handles file I/O and plotting.

- i. Familiarize yourself with the script, the data, and the labels. Plot an example of selected training images.
- ii. Define a fully connected neural network architecture with one hidden layer with 30 neurons and train it. Use the sigmoid activation function for input and the hidden layers, and log softmax for the output layer. Recommended training settings are:
 - i. Batch size: 64
 - ii. Loss function: NLLLoss
 - iii. Optimizer: SGD
 - iv. Training epochs: 15
 - v. Learning rate: 0.003
 - vi. Momentum: 0.9
- iii. Plot the training loss over the epochs

2. Investigate the impact of the number of hidden layers and the activation function

- i. Increase the number of hidden layers to [1,5].
- ii. Change the activation function from sigmoid to ReLU for the input and hidden layers.
- iii. Plot the training loss over the epochs for all experiments and discuss your findings.