

The optimization behind deep neural networks

Marleen Balvert

Disclaimer

This is a tutorial,
not a research talk

Course material

- www.neuralnetworksanddeeplearning.com
- www.deeplearningbook.org

5 7 1 6 4

7 2 2 5 9

0 0 4 9 7

6 6 1 0 7

5 0 7 9 3

5 6 2 7 2

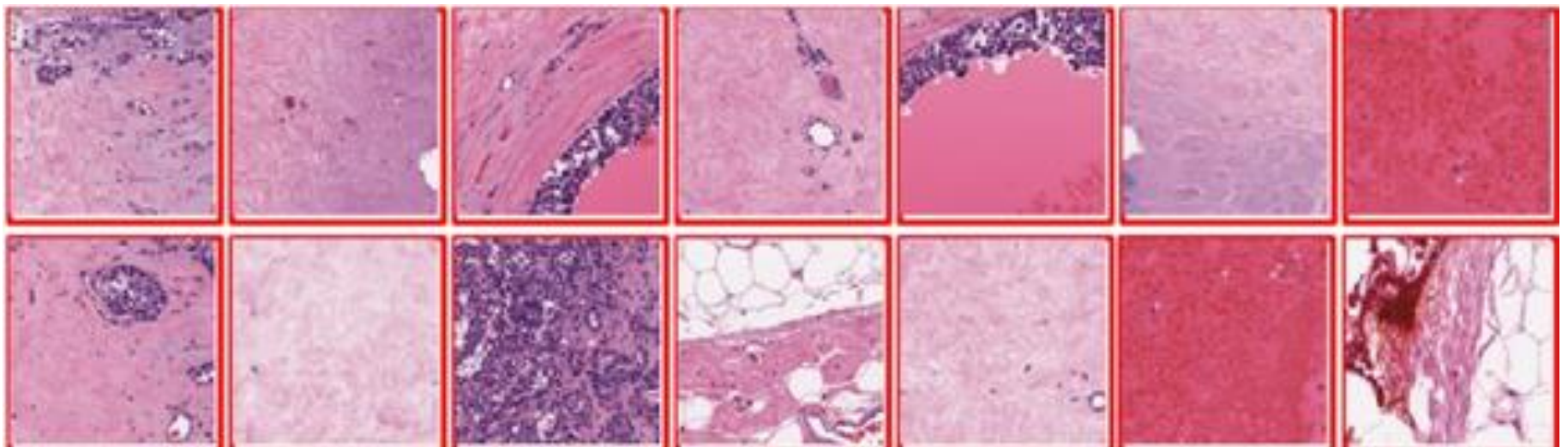
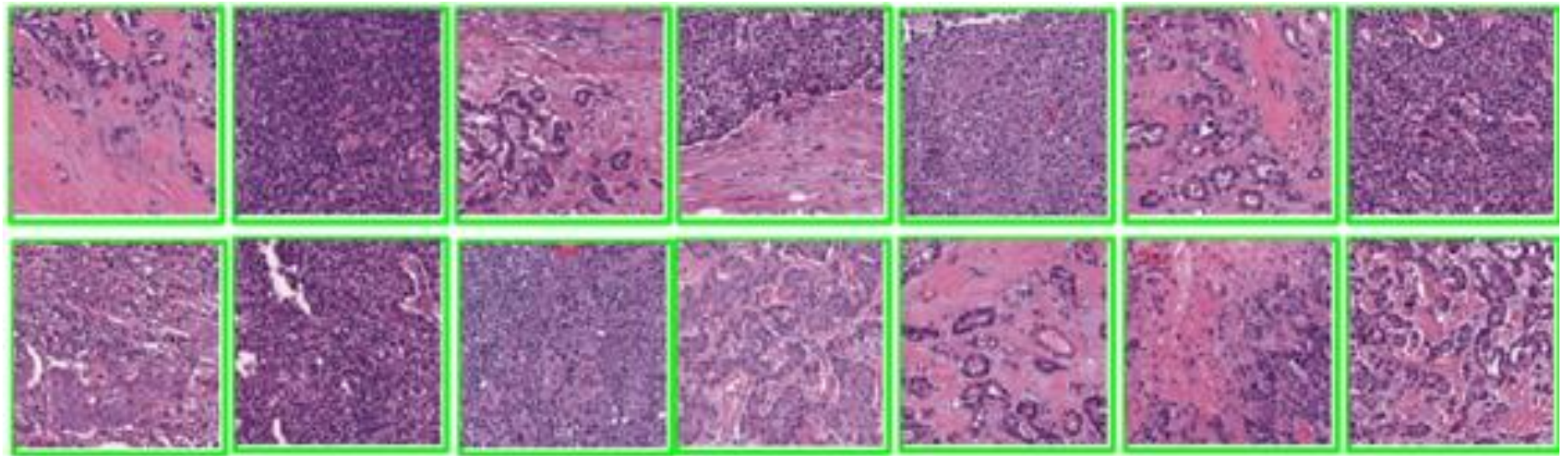
3 3 5 7 0

5 3 4 5 0

0 3 5 2 7

8 0 2 0 2







Deep neural network basics

The classification problem (recap)

Given data:

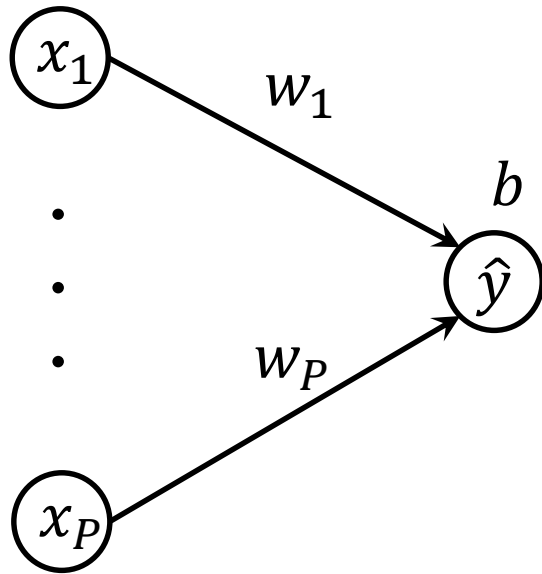
- Samples $1, \dots, N$
- Independent variables $x_1, \dots, x_N \in \mathbb{R}^P$
- Dependent variables $y_1, \dots, y_N \in \{0, \dots, K\}$
 - Often converted into K dummy variables

Goal: find the relation $y = f(x)$

Goal: make a model that determines the class of a new sample

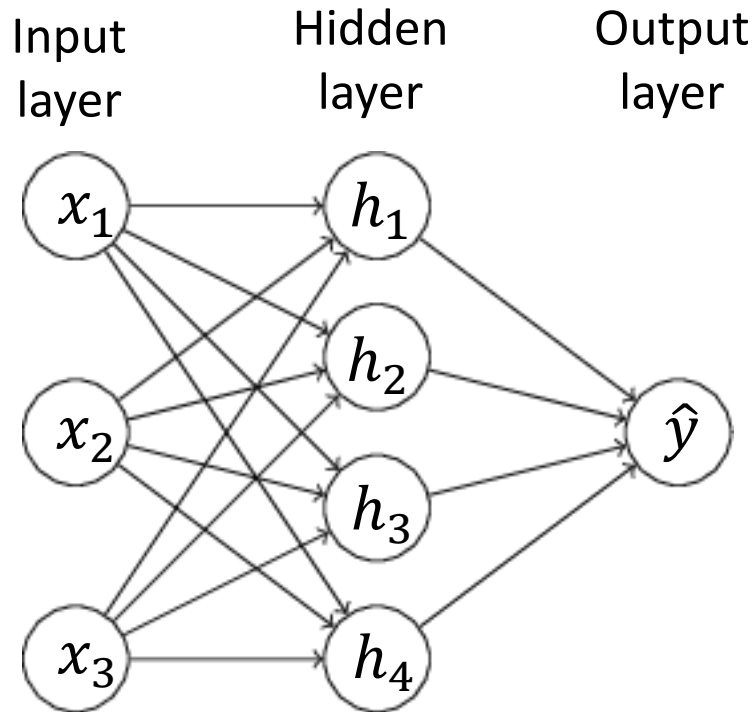
Basic element: neuron or perceptron

Notation: $x \in \mathbb{R}^P$ is feature vector of single sample



$$\hat{y} = \begin{cases} 1 & \text{if } w^T x + b \geq 0 \\ 0 & \text{if } w^T x + b < 0 \end{cases}$$

Neural network (NN)

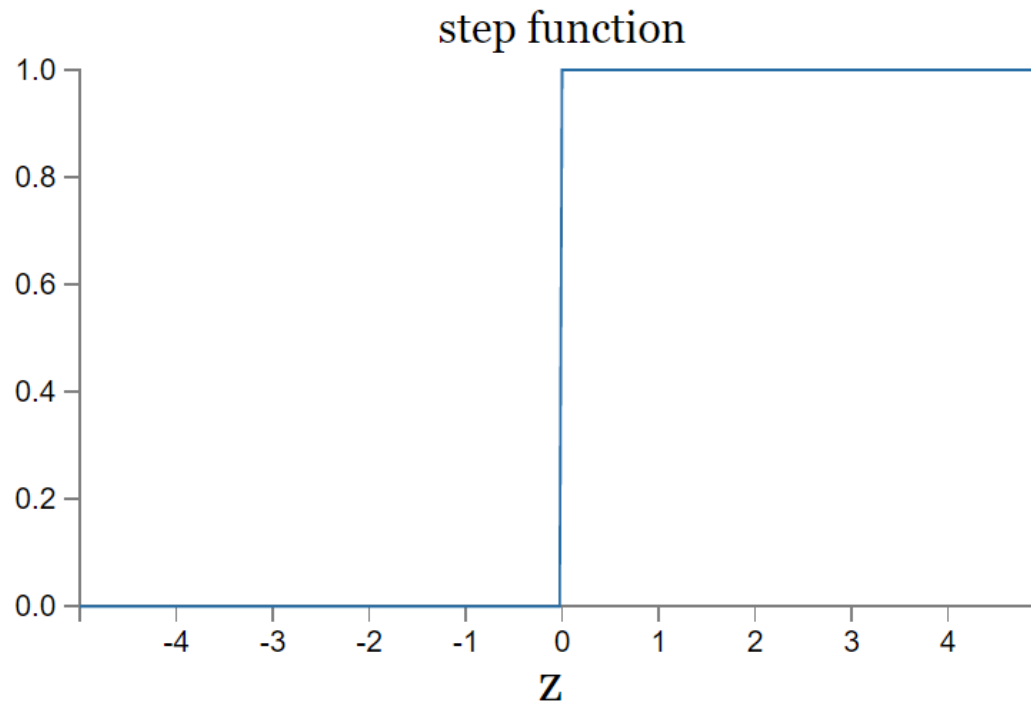


$$h_i = \begin{cases} 1 & \text{if } w^{iT} x + b_i \geq 0 \\ 0 & \text{if } w^{iT} x + b_i < 0 \end{cases}$$

$$\hat{y} = \begin{cases} 1 & \text{if } w^{5T} h + b_5 \geq 0 \\ 0 & \text{if } w^{5T} h + b_5 < 0 \end{cases}$$

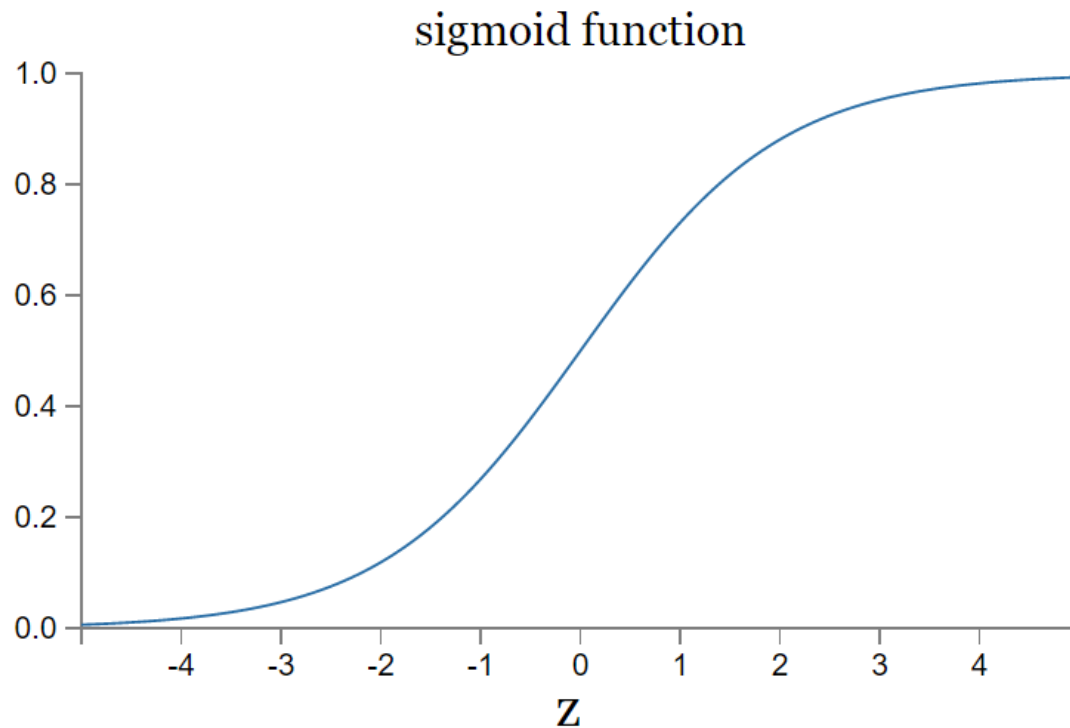
The activation function: step

$$\hat{y} = \begin{cases} 1 & \text{if } w^T x + b \geq 0 \\ 0 & \text{if } w^T x + b < 0 \end{cases}$$



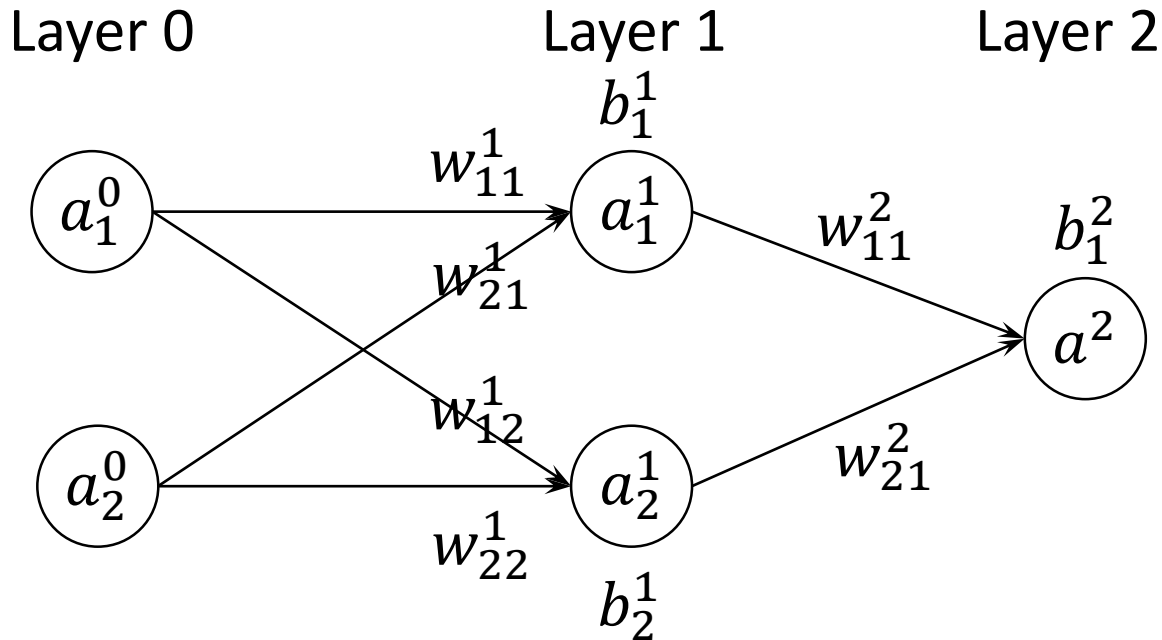
The activation function: sigmoid

$$\hat{y} = \frac{1}{1 + e^{-w^T x - b}} = \sigma(w^T x + b)$$



Optimizing a (deep) neural network

Notation



Where:

- w_{jk}^l weight from node j in layer $l - 1$ to node k in layer l
- b_j^l bias in node j of layer l
- $a_j^l = \sigma(w^l a^{l-1} + b^l), l = 2, 3$
- $a_j^0 = x_j$

Optimizing the weights and biases

- Cost function, e.g.

$$C(y, \hat{y}) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

Optimizing the weights and biases

- Cost function, e.g.

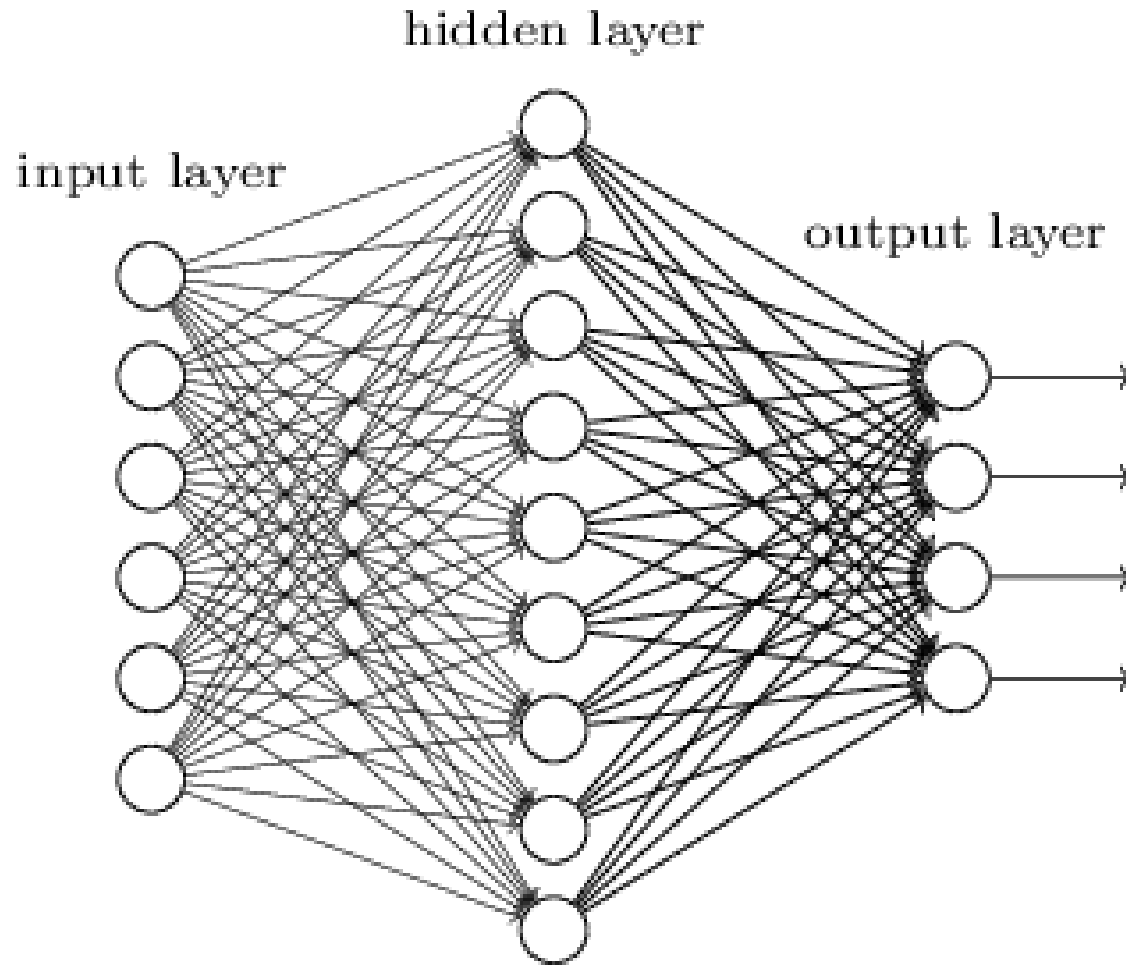
$$C(y, X, w, b) = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i(w, b, x_i))^2$$

- Minimize the cost function using gradient descent:

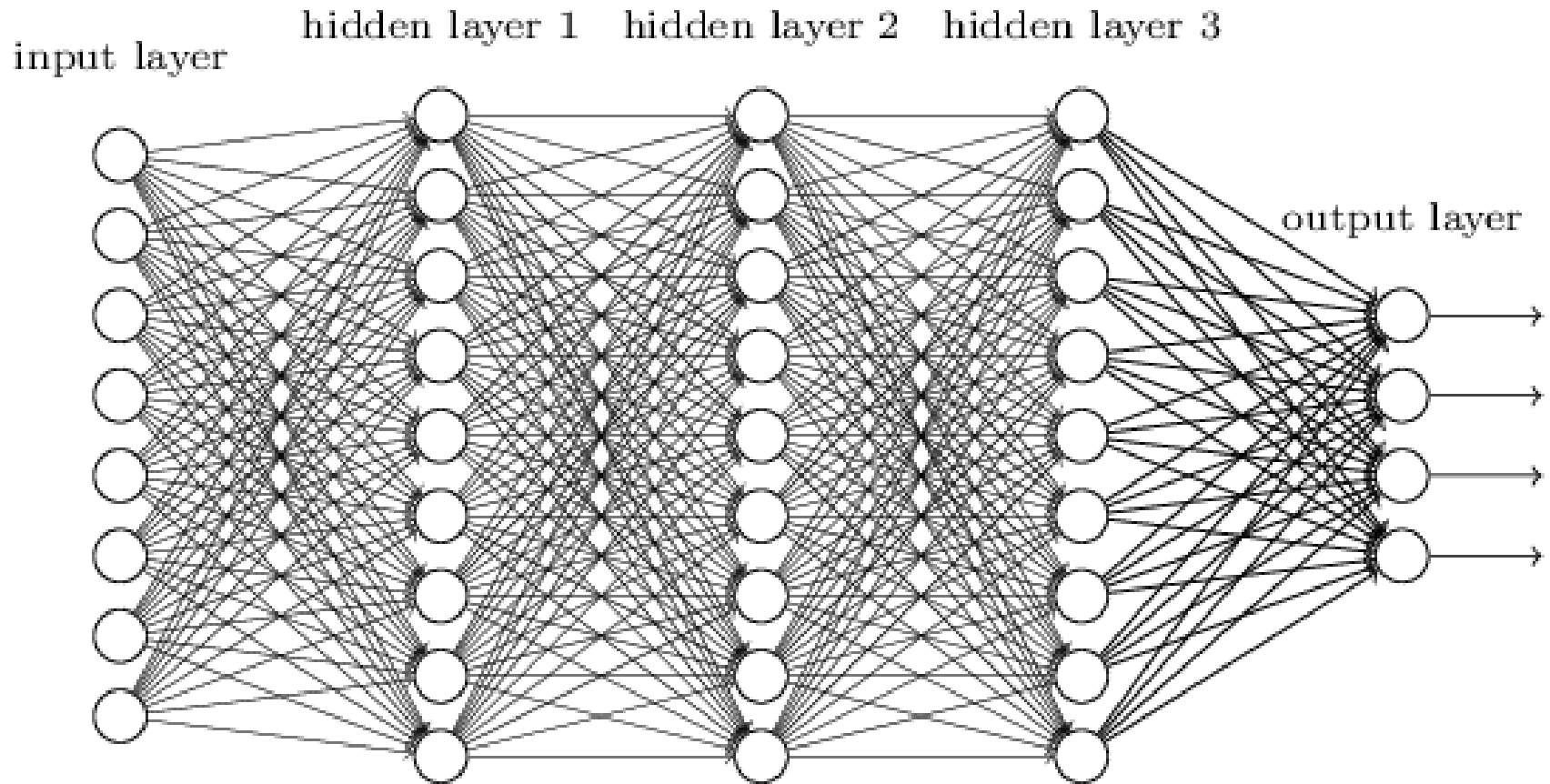
$$w_{jk}^l \leftarrow w_{jk}^l - \eta \left[\frac{\partial C}{\partial w_{jk}^l} \right]$$

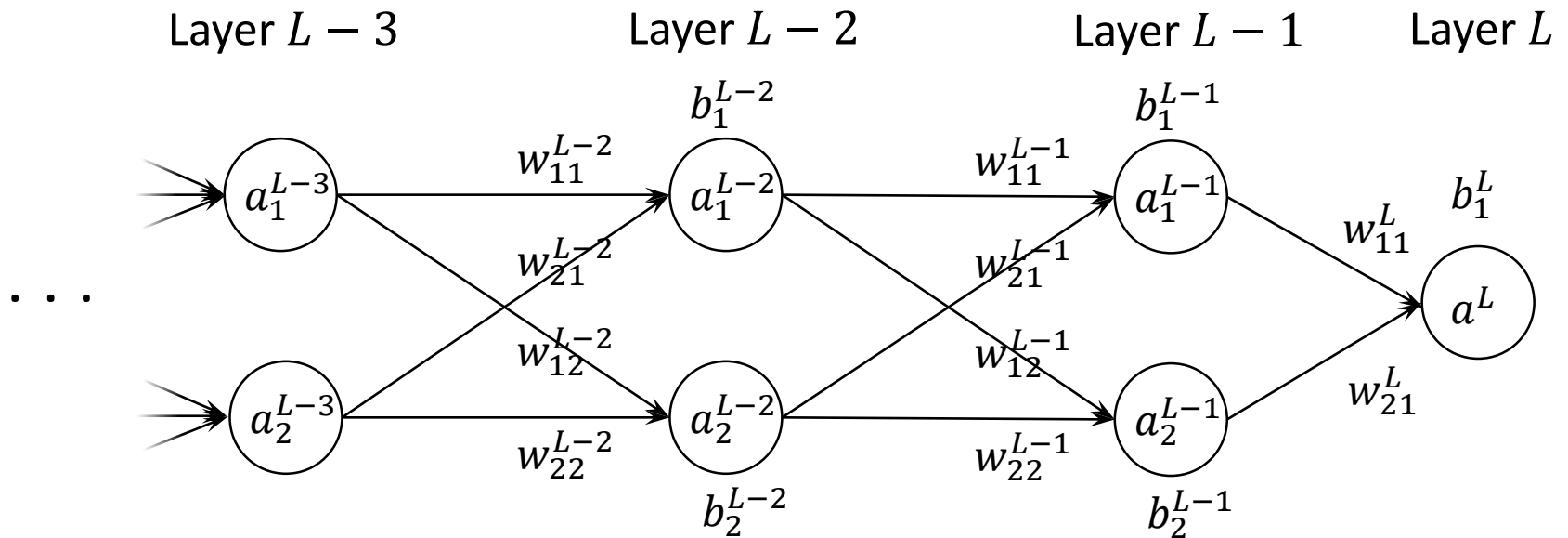
$$b_j^l \leftarrow b_j^l - \eta \left[\frac{\partial C}{\partial b_j^l} \right]$$

Neural networks (NNs)



Deep neural networks (DNNs)





Define:

- $a^l = \sigma(w^l a^{l-1} + b^l)$ **activation** of layer l
- $z^l = w^l a^{l-1} + b^l$ **weighted input** to layer l
- $\delta_j^l = \frac{\partial C}{\partial z_j^l}$ **error** in node j of layer l

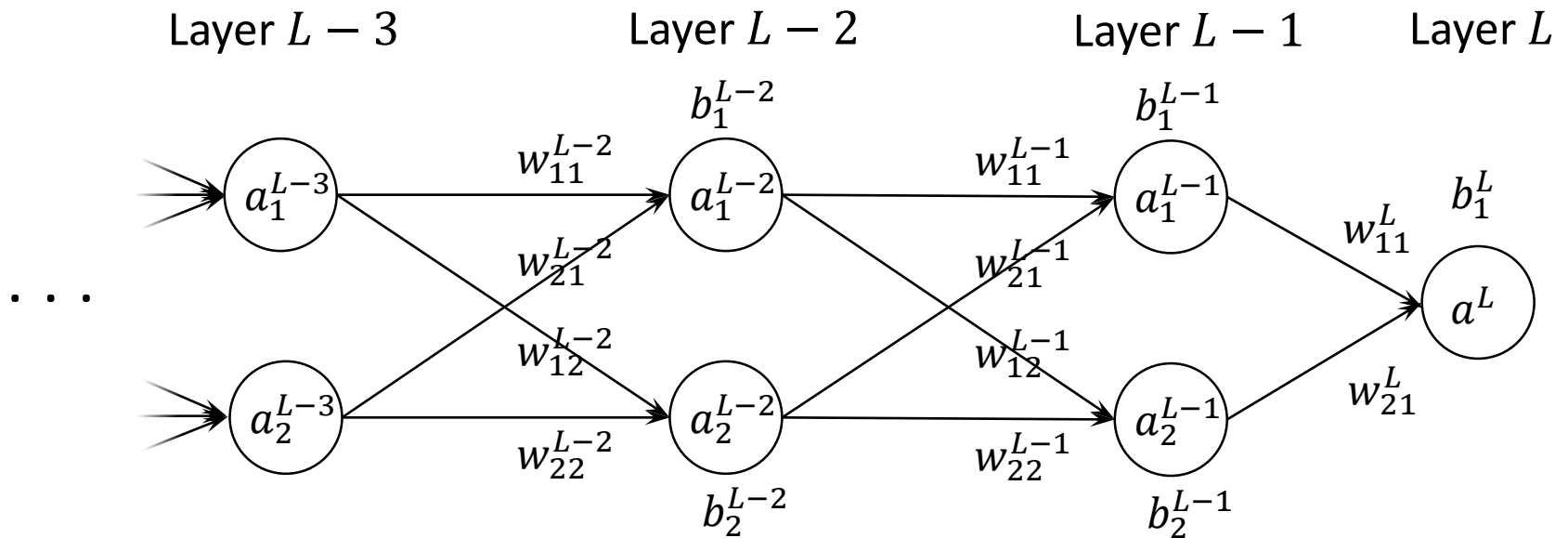
We can now compute:

- $\delta^L = \nabla_{a^L} C \cdot \sigma'(z^L)$
- $\delta^l = \left((w^{l+1})^T \delta^{l+1} \right) \odot \sigma'(z^l)$



$$\frac{\partial C}{\partial b_j^l} = \frac{\partial C}{\partial z_j^l} \cdot \frac{\partial z_j^l}{\partial b_j^l}$$

$$\frac{\partial C}{\partial w_{jk}^l} = \frac{\partial C}{\partial z_j^l} \cdot \frac{\partial z_j^l}{\partial w_{jk}^l}$$



Backpropagation

Set the input $a^0 = x_j$. Then:

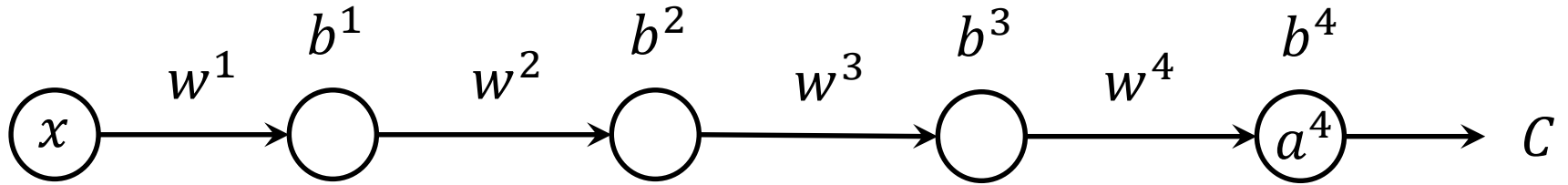
1. Feedforward: compute $z^l = w^l a^{l-1} + b^l$, $a^l = \sigma(z^l)$
2. Output error: compute $\delta^L = \nabla_a C \odot \sigma'(z^L)$
3. Backpropagate the error:

$$\delta^l = \left((w^{l+1})^T \delta^{l+1} \right) \odot \sigma'(z^l)$$

4. Compute the gradients: $\frac{\partial C}{\partial b_j^l} = \delta_j^l$ and $\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$

Challenge: Unstable gradients

A simple example

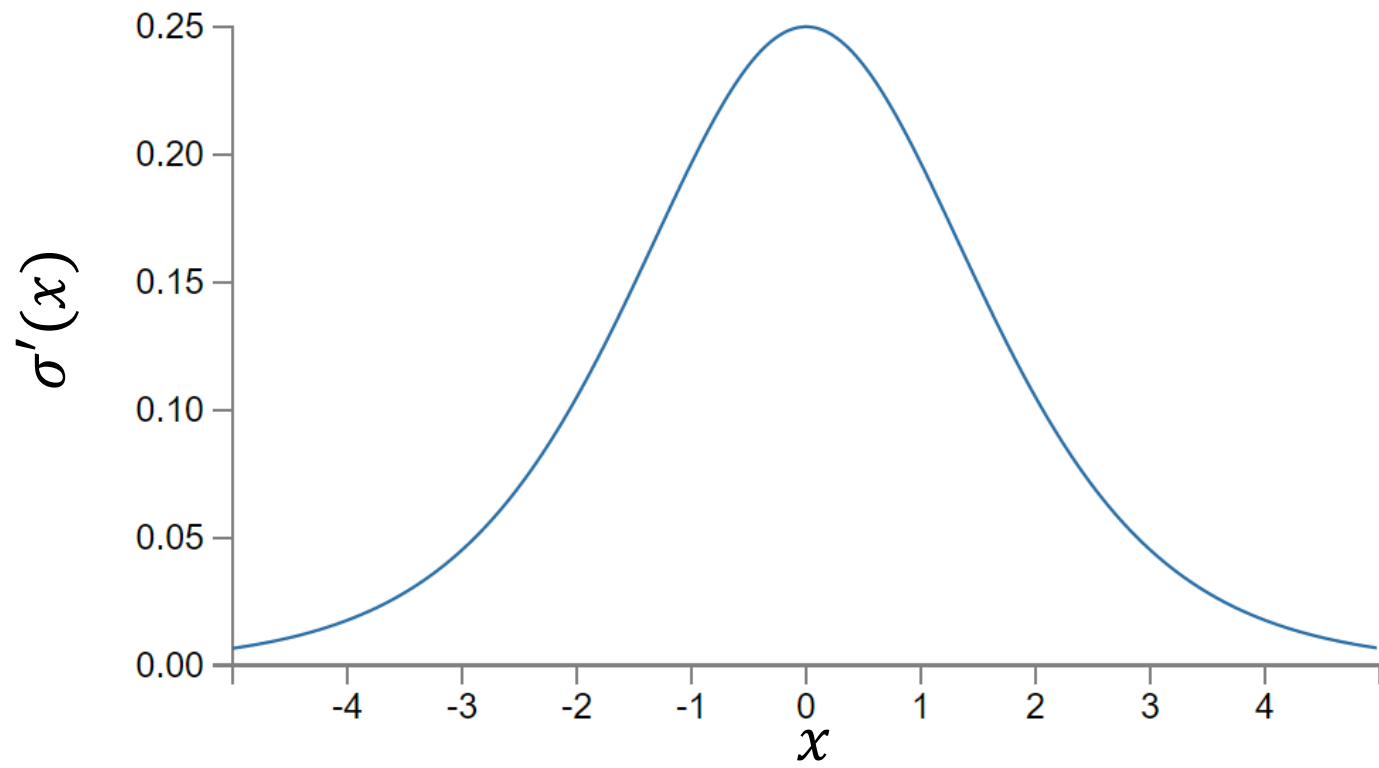


$$\frac{\partial C}{\partial b^4} = \sigma'(z^4) \nabla_{a^4} C$$

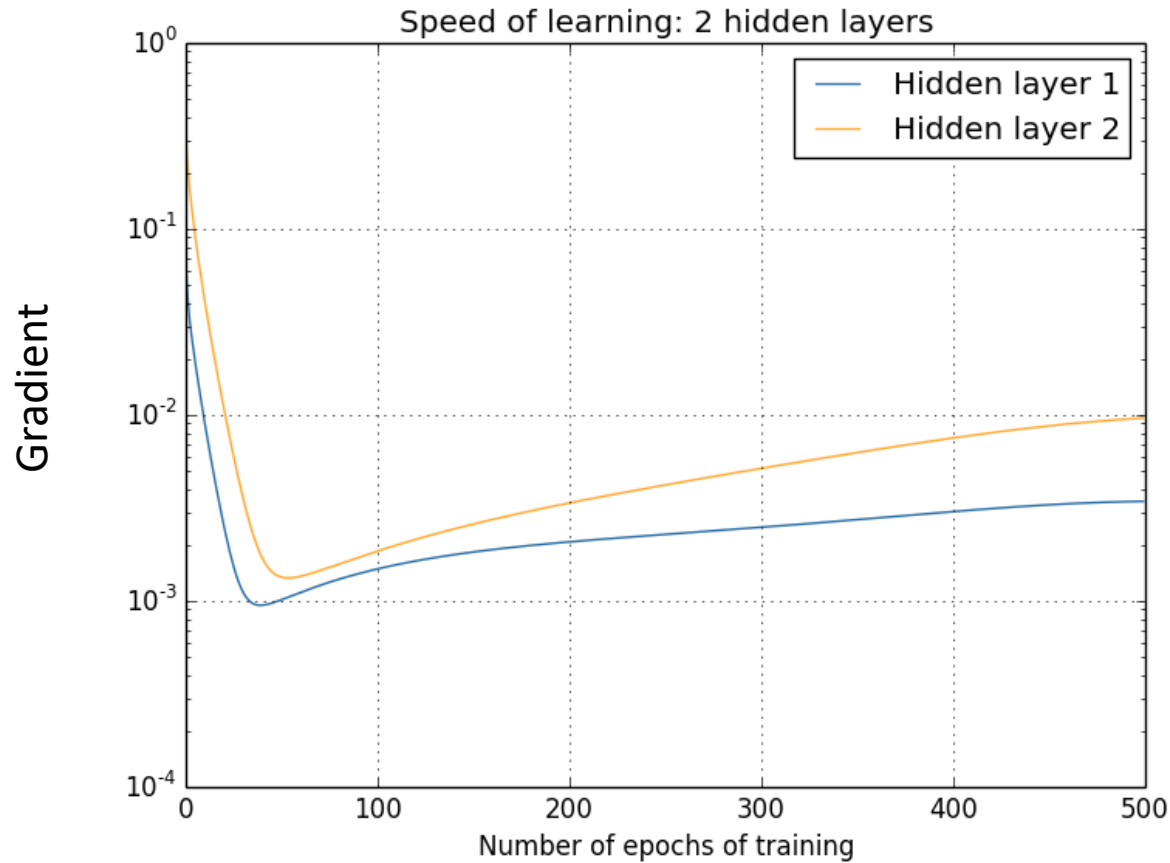
$$\frac{\partial C}{\partial b^3} = \sigma'(z^3) \cdot w^4 \cdot \sigma'(z^4) \nabla_{a^4} C$$

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

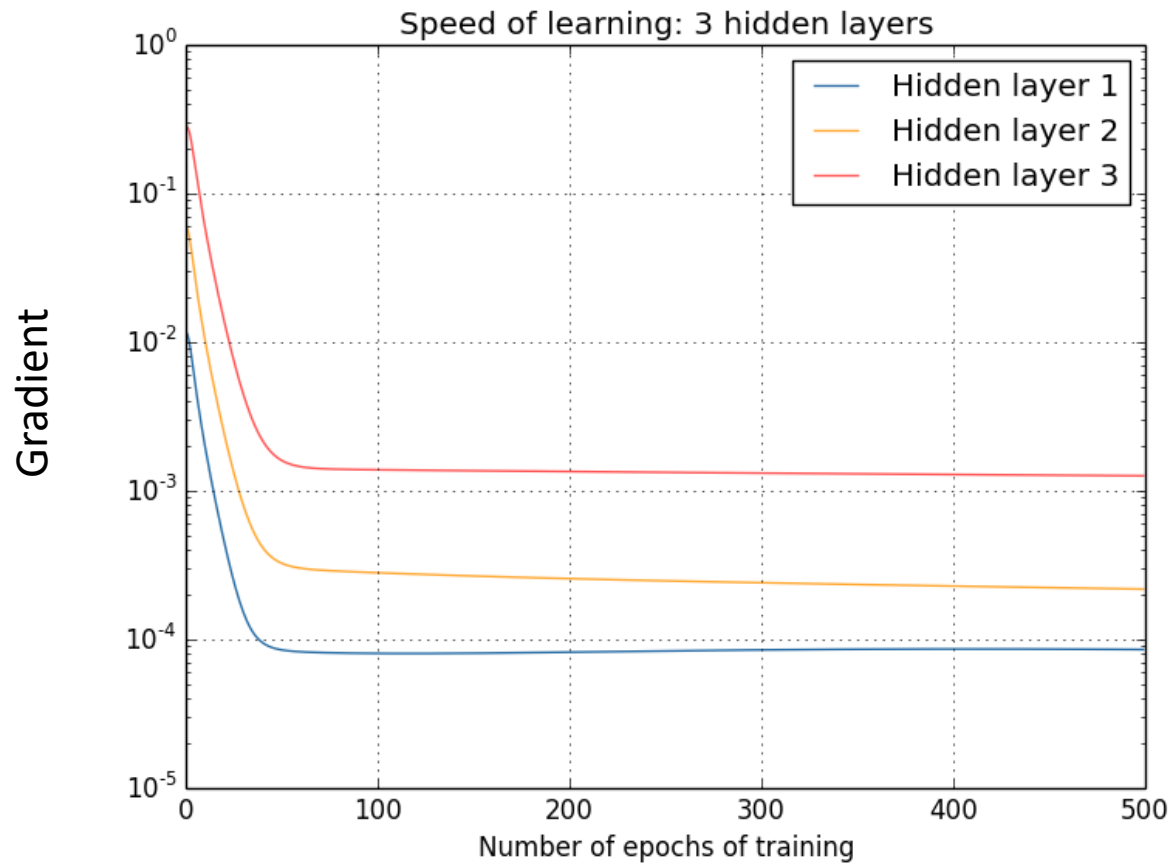
$$\frac{\partial C}{\partial b^1} = \sigma'(z^1) \cdot w^2 \cdot \sigma'(z^2) \cdot w^3 \cdot \sigma'(z^3) \cdot w^4 \cdot \sigma'(z^4) \nabla_{a^4} C$$



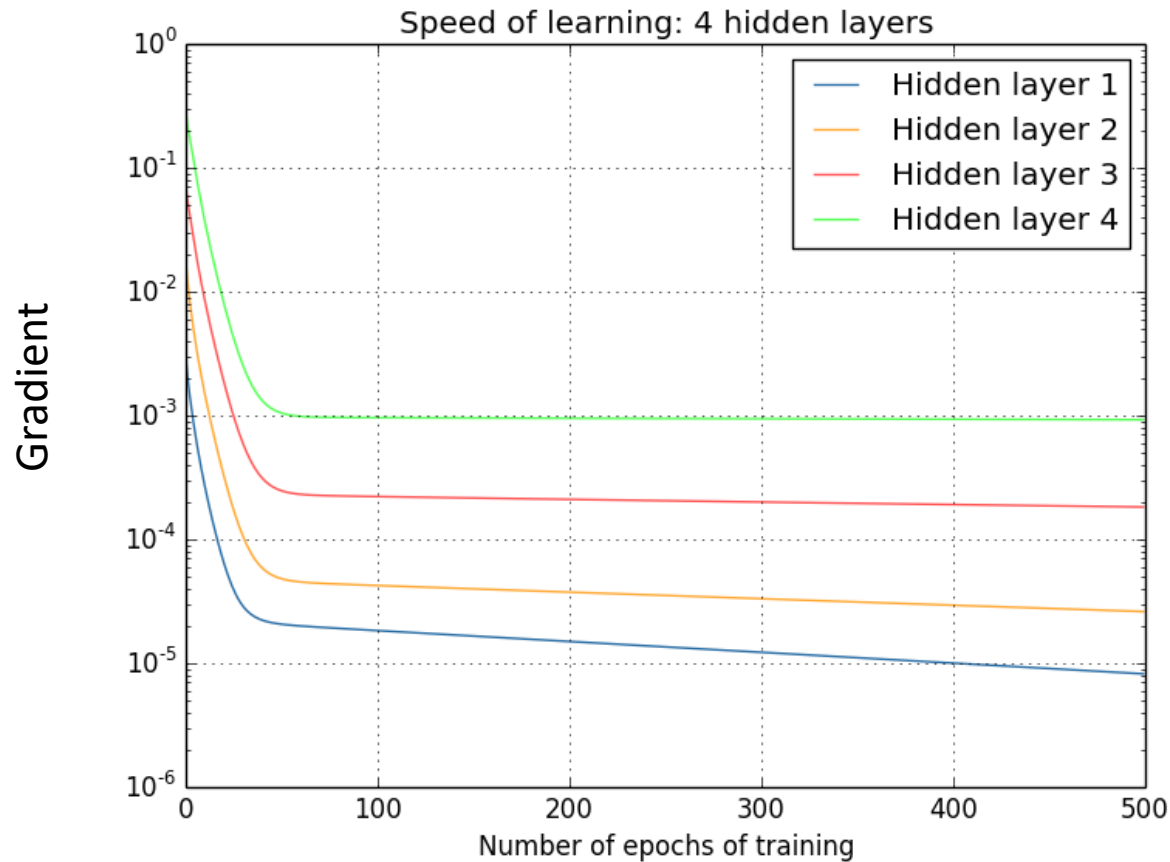
Vanishing gradients



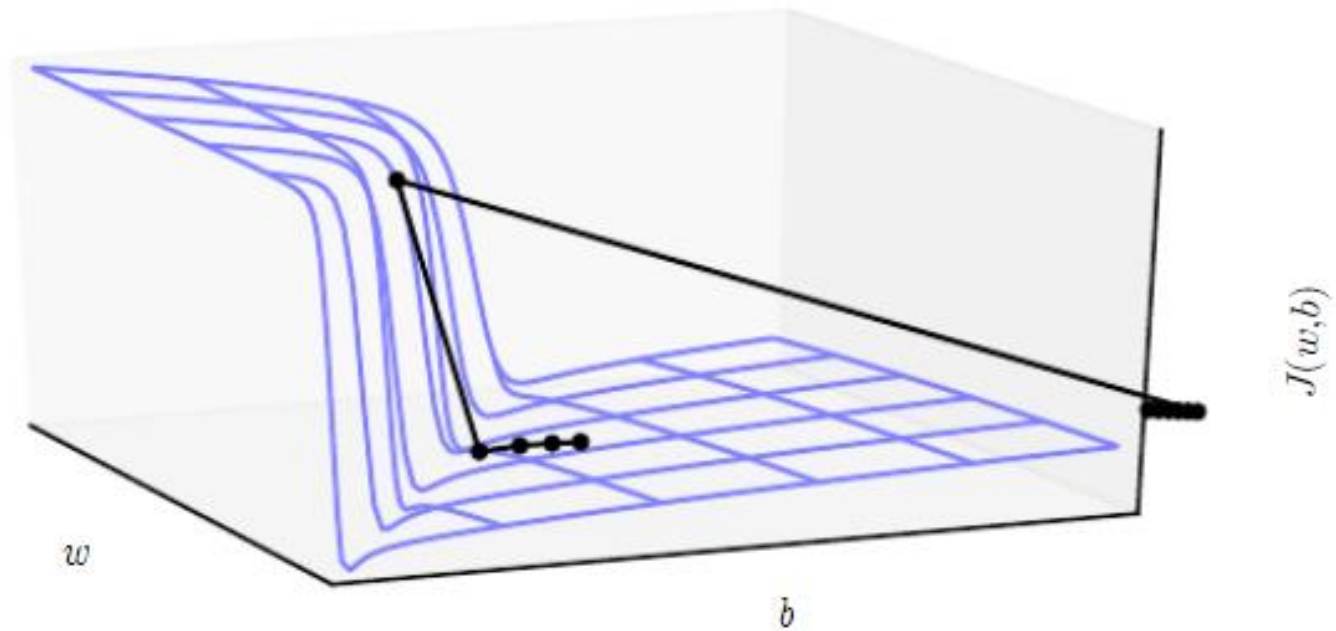
Vanishing gradients



Vanishing gradients



Exploding gradients



Solutions

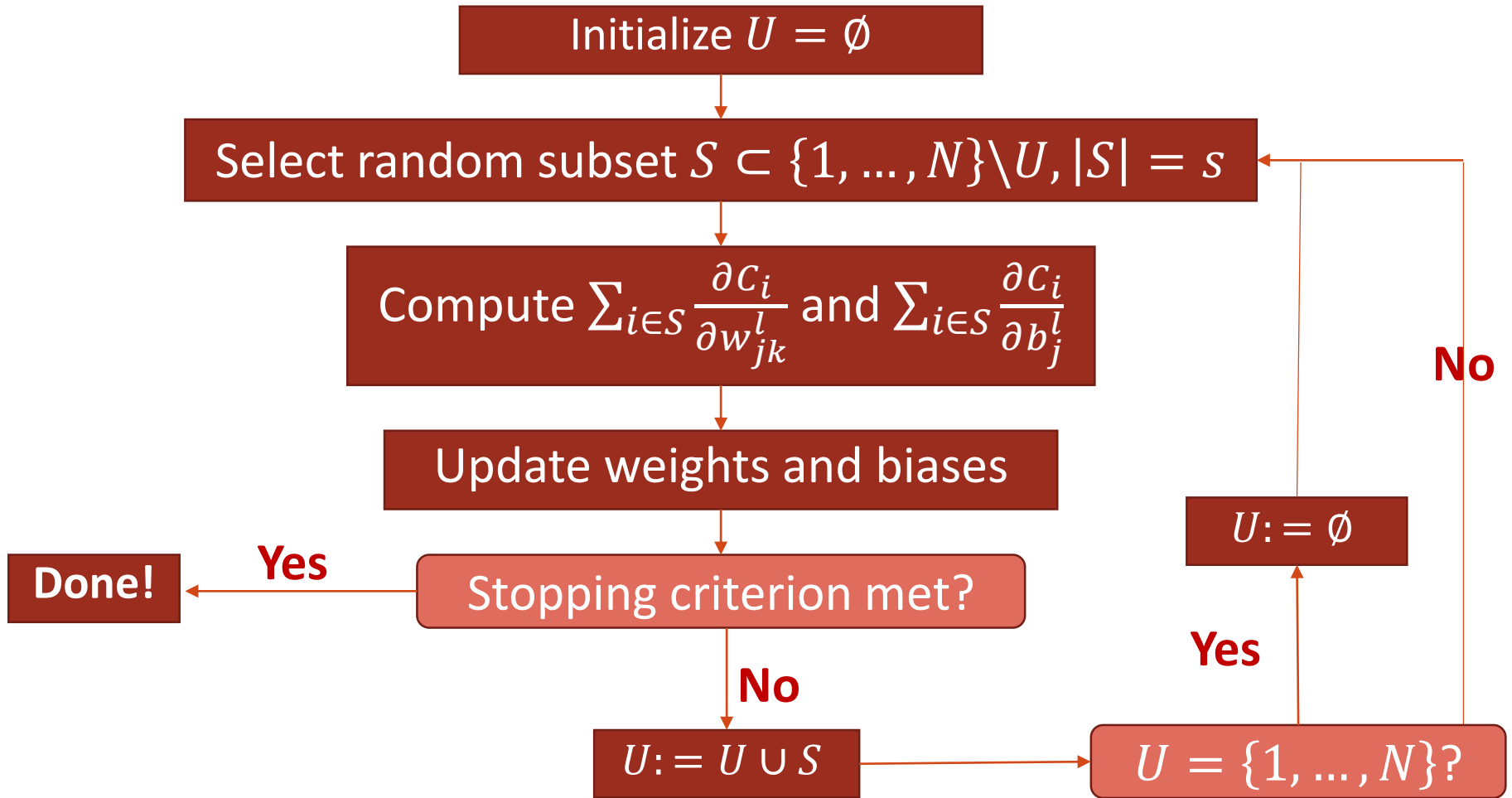
- Other activation functions
- Batch normalization
- Gradient clipping (for exploding gradients)
- Other?

Challenge: large number of samples

$$C(\cdot) = \sum_{i=1}^N C_i(\cdot)$$

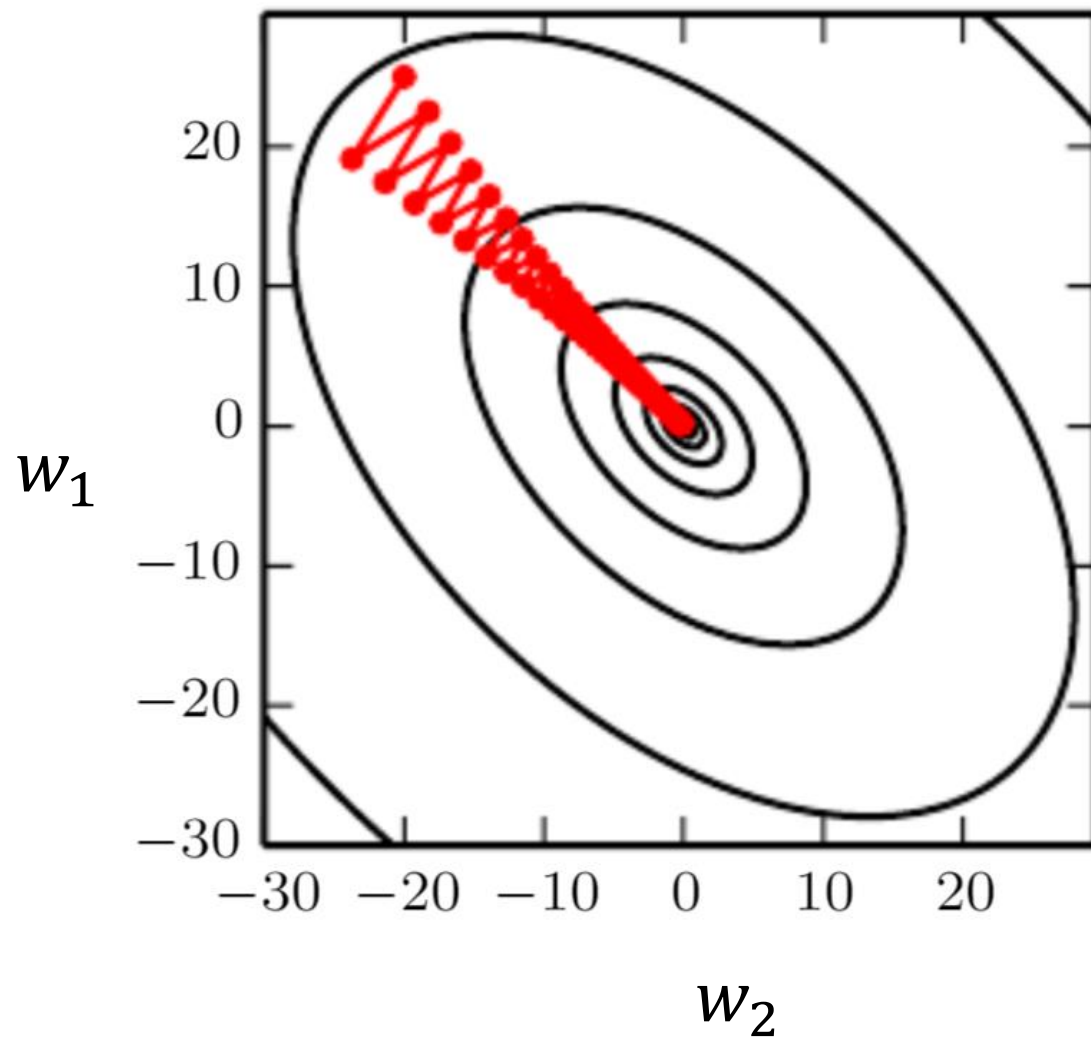
$$\frac{\partial C}{\partial w_{jk}^l} = \sum_{i=1}^N \frac{\partial C_i}{\partial w_{jk}^l}$$

Stochastic gradient descent



Challenge: valleys

Or: the Hessian is ill-conditioned



Solution: momentum

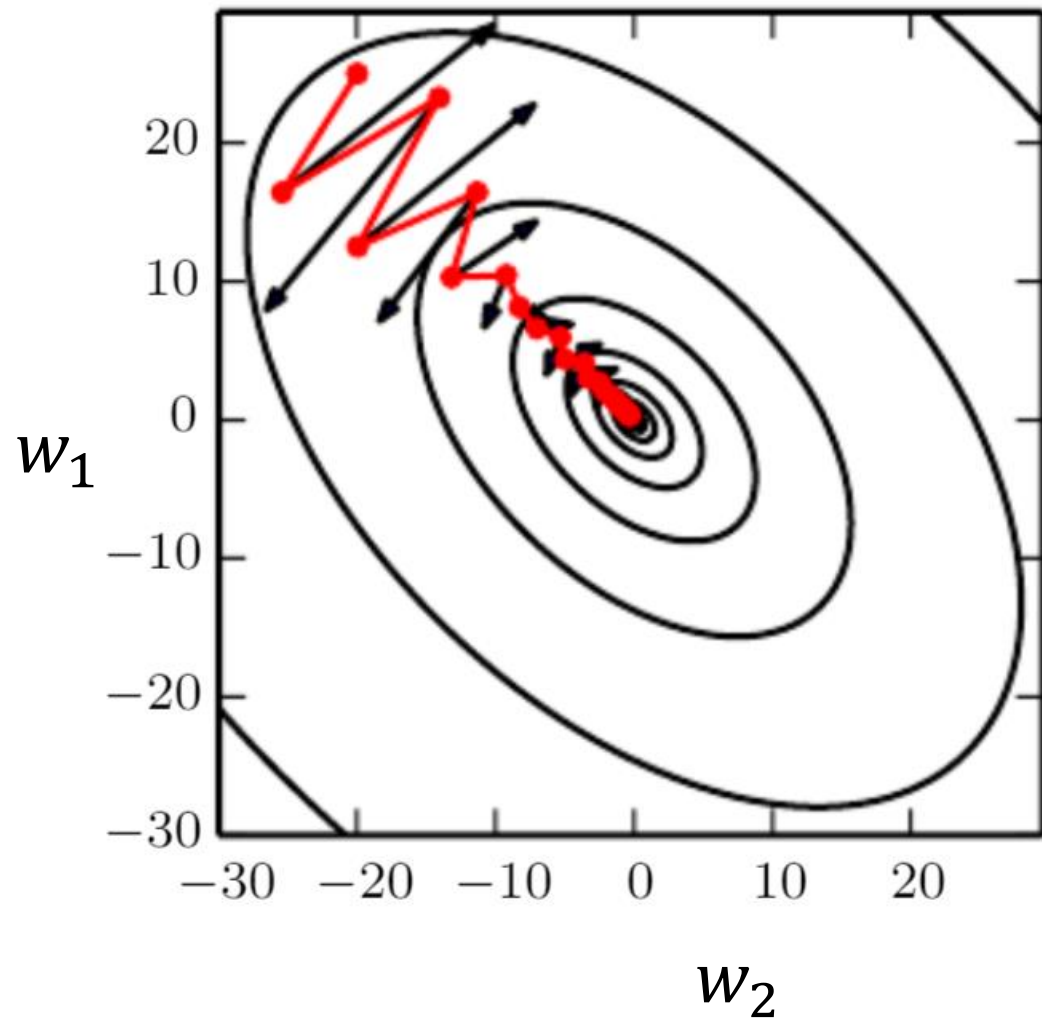
- Gradient descent:

$$w \leftarrow w - \eta \nabla_w C$$

- Momentum:

$$v \leftarrow \alpha \cdot v - \eta \nabla_w C$$

$$w \leftarrow w + v$$



Solution: AdaGrad

- Gradient descent:

$$w \leftarrow w - \eta \nabla_w C$$

- AdaGrad:

$$r \leftarrow r + \nabla_w C^T \nabla_w C$$

$$w \leftarrow w - \frac{\eta}{\sqrt{\delta + r}} \nabla_w C$$

Solution: RMSProp

- Gradient descent:

$$w \leftarrow w - \eta \nabla_w C$$

- RMSProp:

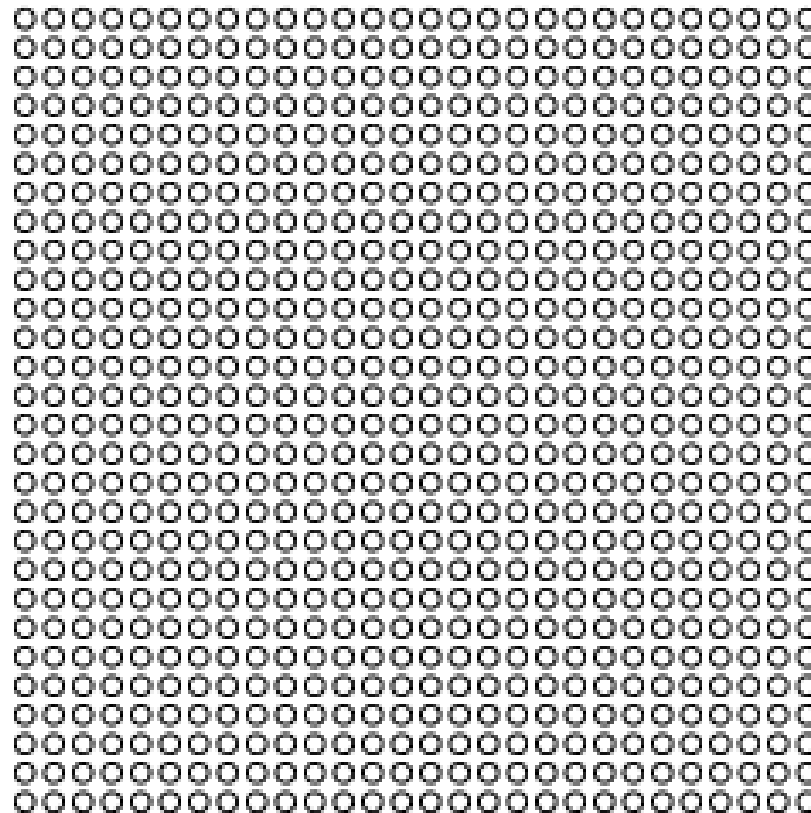
$$r \leftarrow \alpha r + (1 - \alpha) \nabla_w C^T \nabla_w C$$

$$w \leftarrow w - \frac{\eta}{\sqrt{\delta + r}} \nabla_w C$$

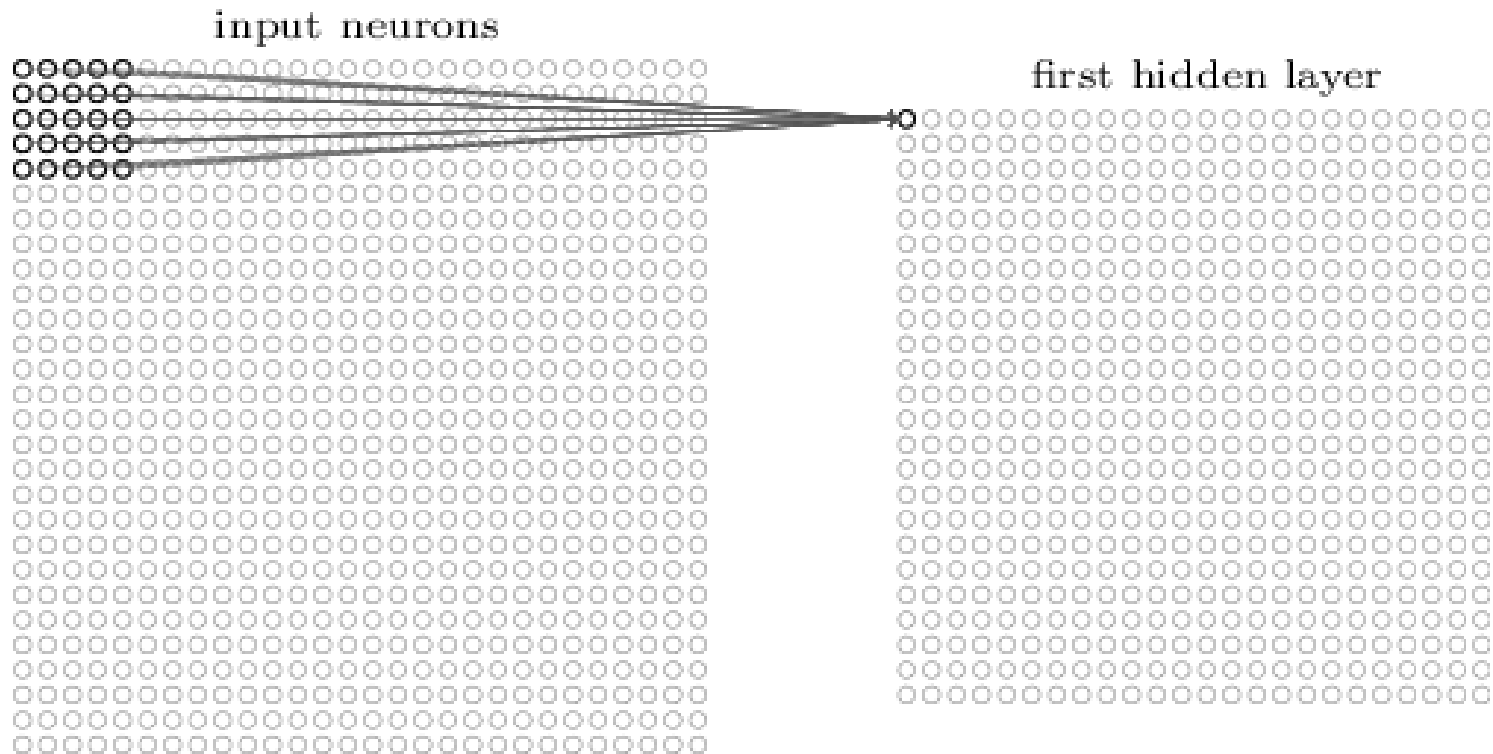
Challenge: architecture design

Convolution layers

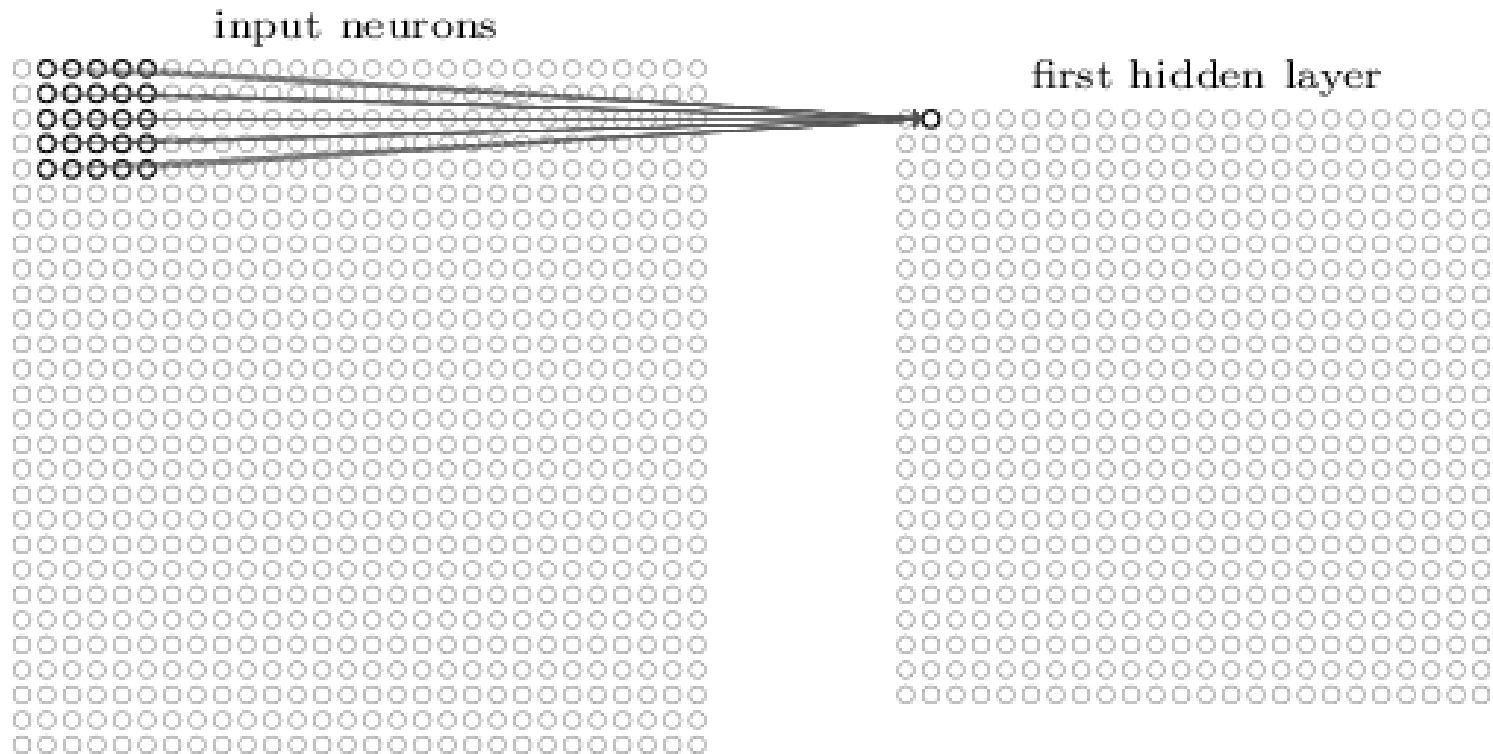
input neurons



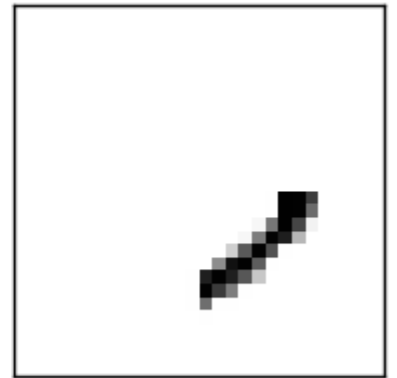
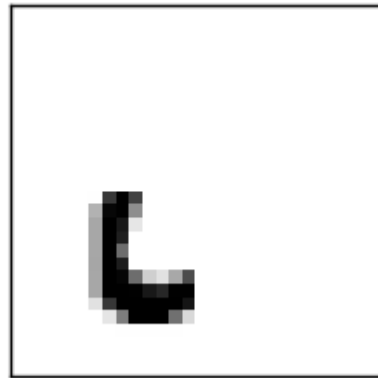
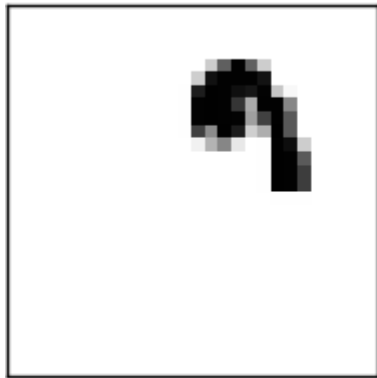
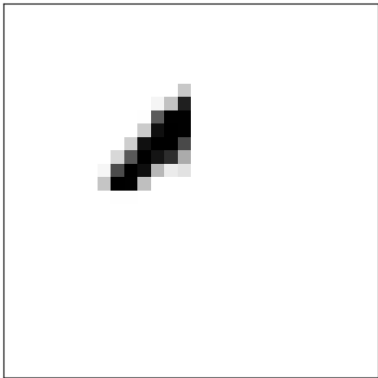
Convolution layers: Local receptive fields

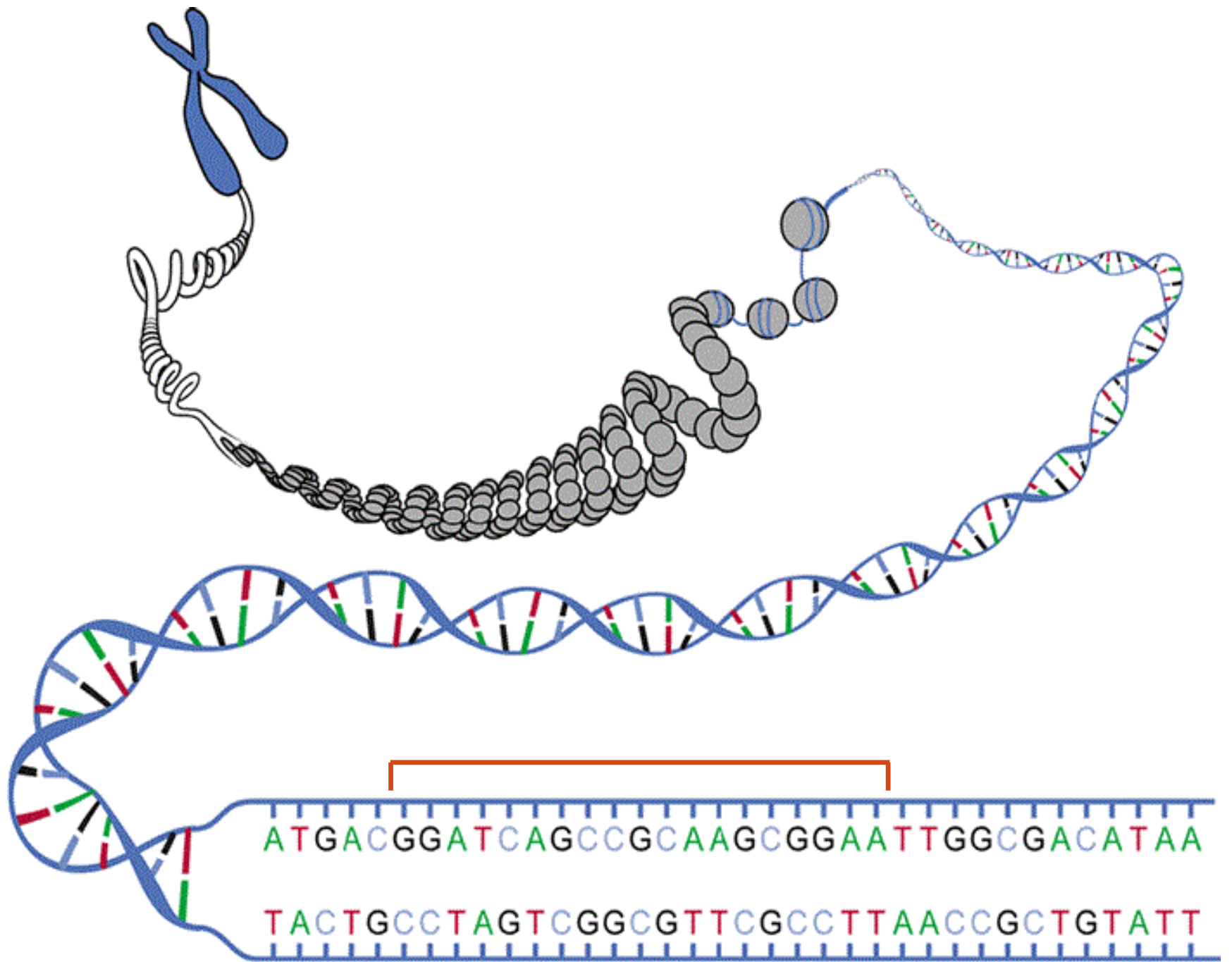


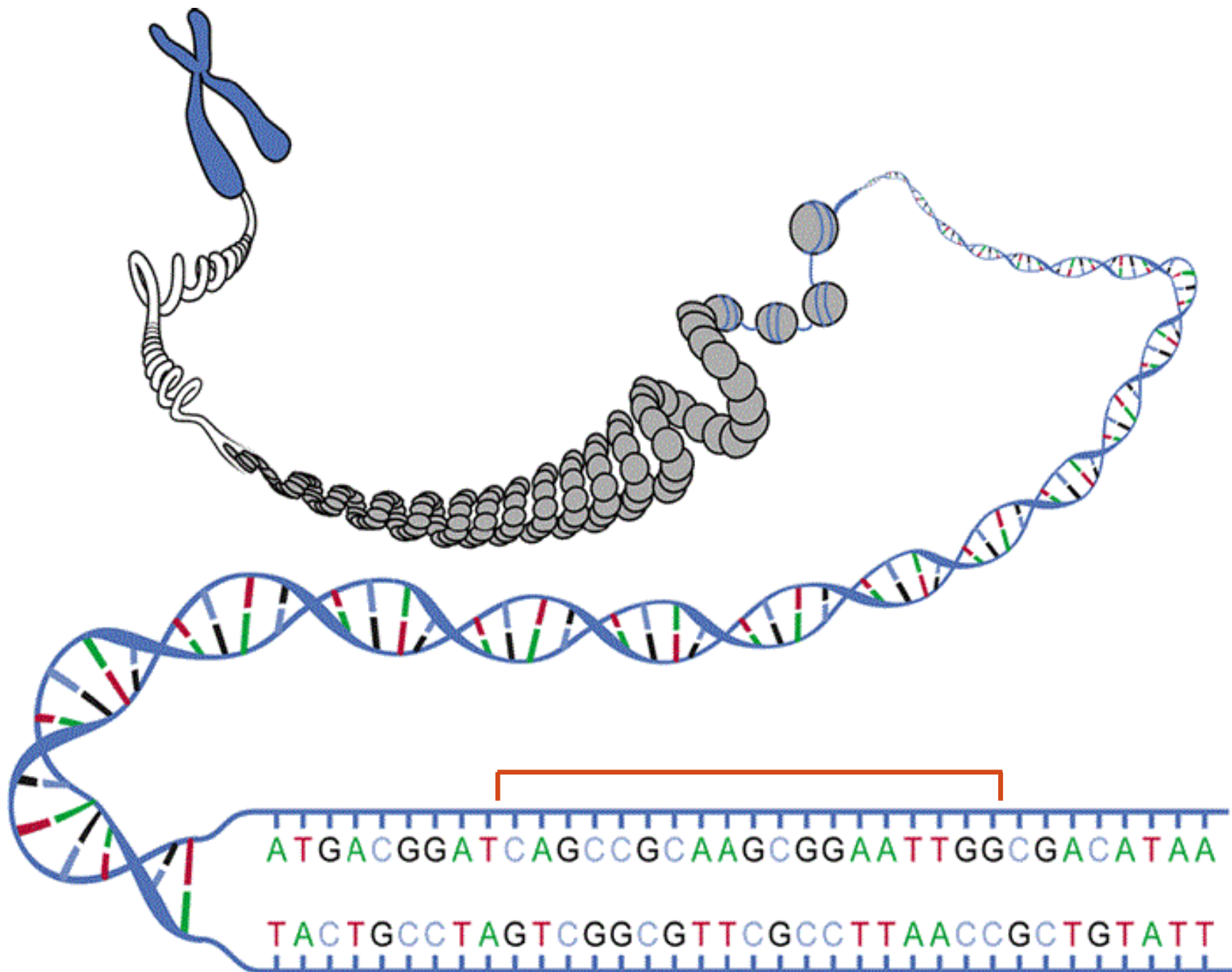
Convolution layers: Local receptive fields

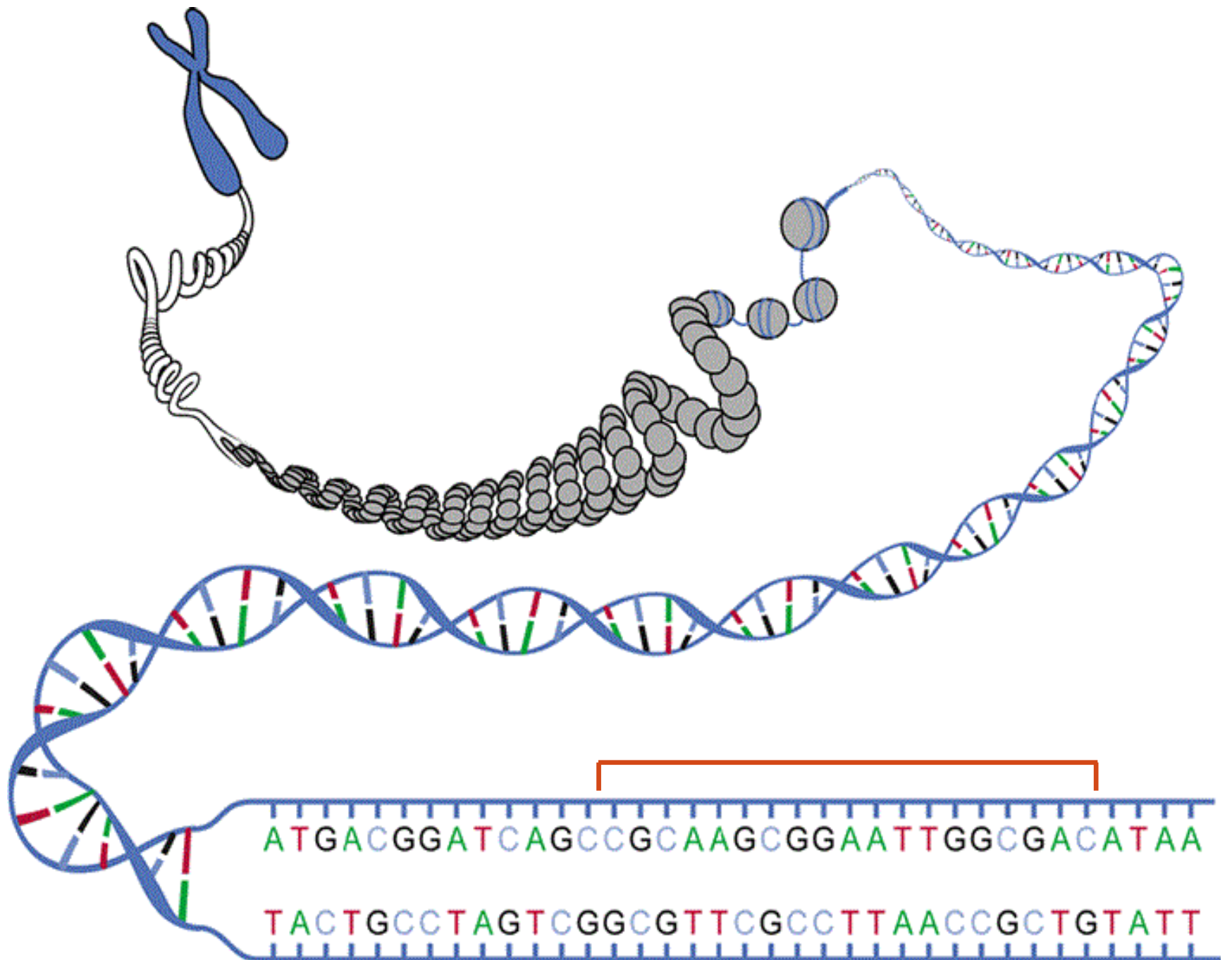


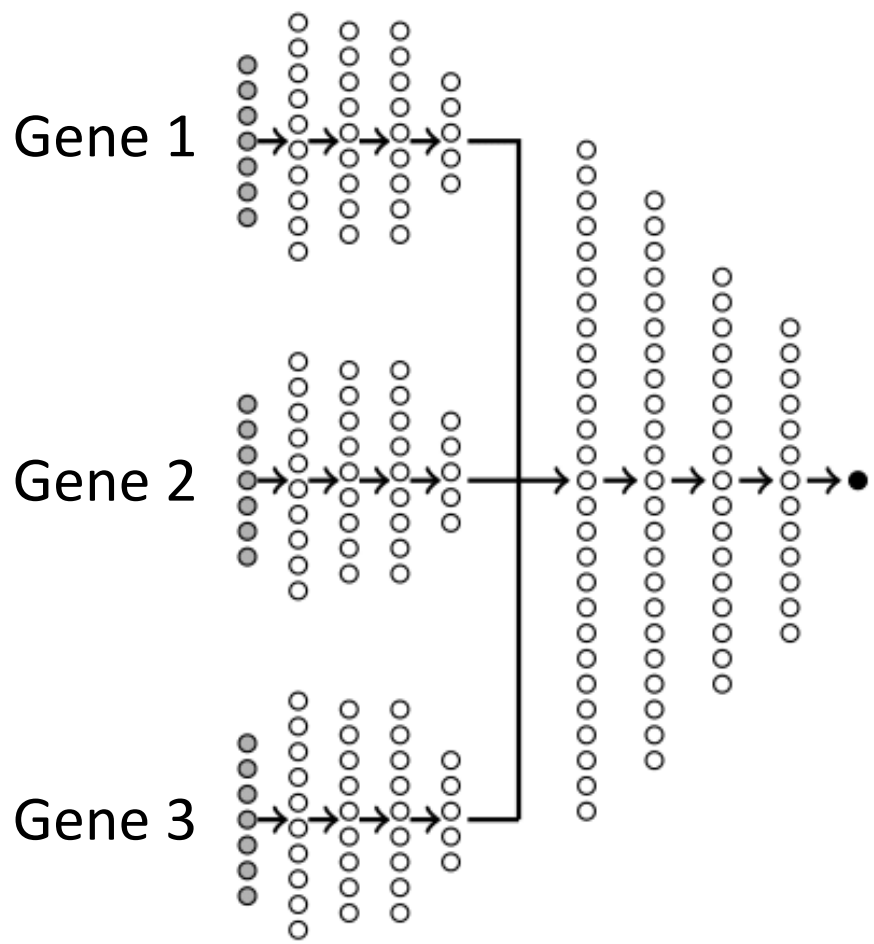
Learning features



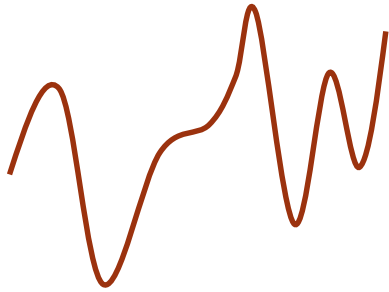








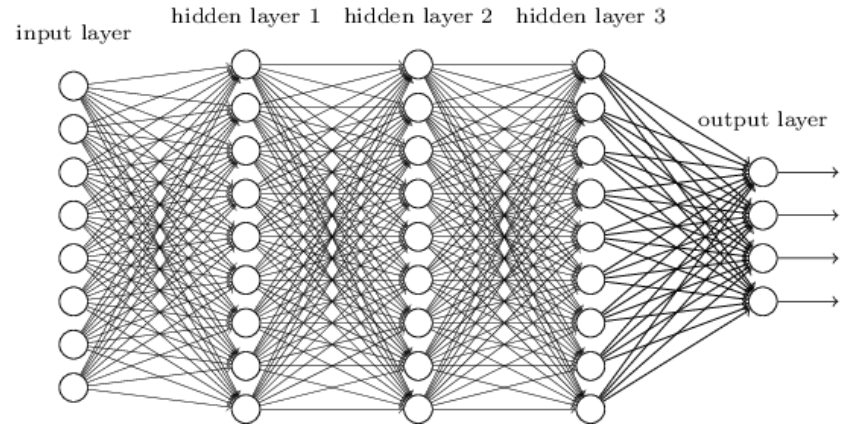
Some more challenges...



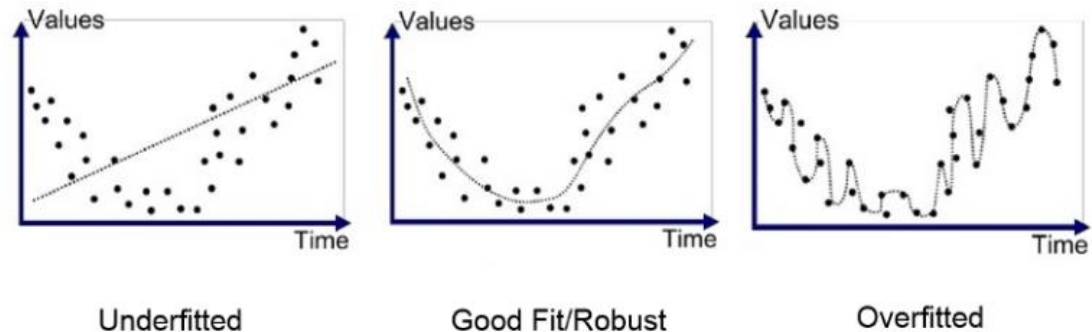
Local minima and saddle points

$$H_x = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_1 x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n x_1} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

Second order methods do not scale well



Network design is a manual process



Risk of overfitting

To conclude

