



GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN

Deep Learning

Lecture 4: Multi-layer perceptron

Alexander Ecker
Institut für Informatik, Uni Göttingen



<https://alexanderecker.wordpress.com>

– Credit: some of the slides based on Ian Goodfellow's slides –

Agenda for this week

Motivation for deep learning

Perceptron

Multi-layer perceptron

Stochastic gradient descent

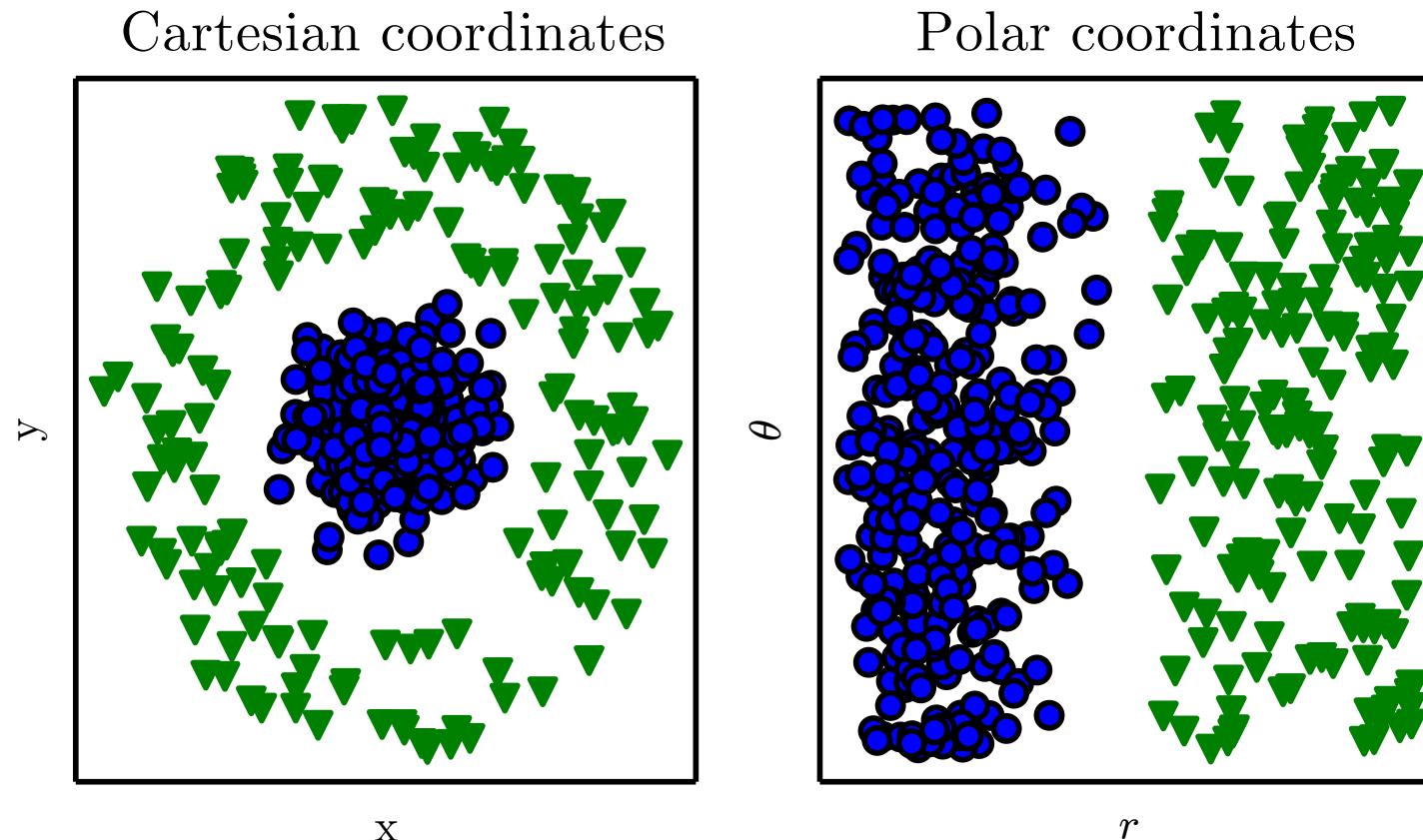
Regularization of deep nets

Architecture choices

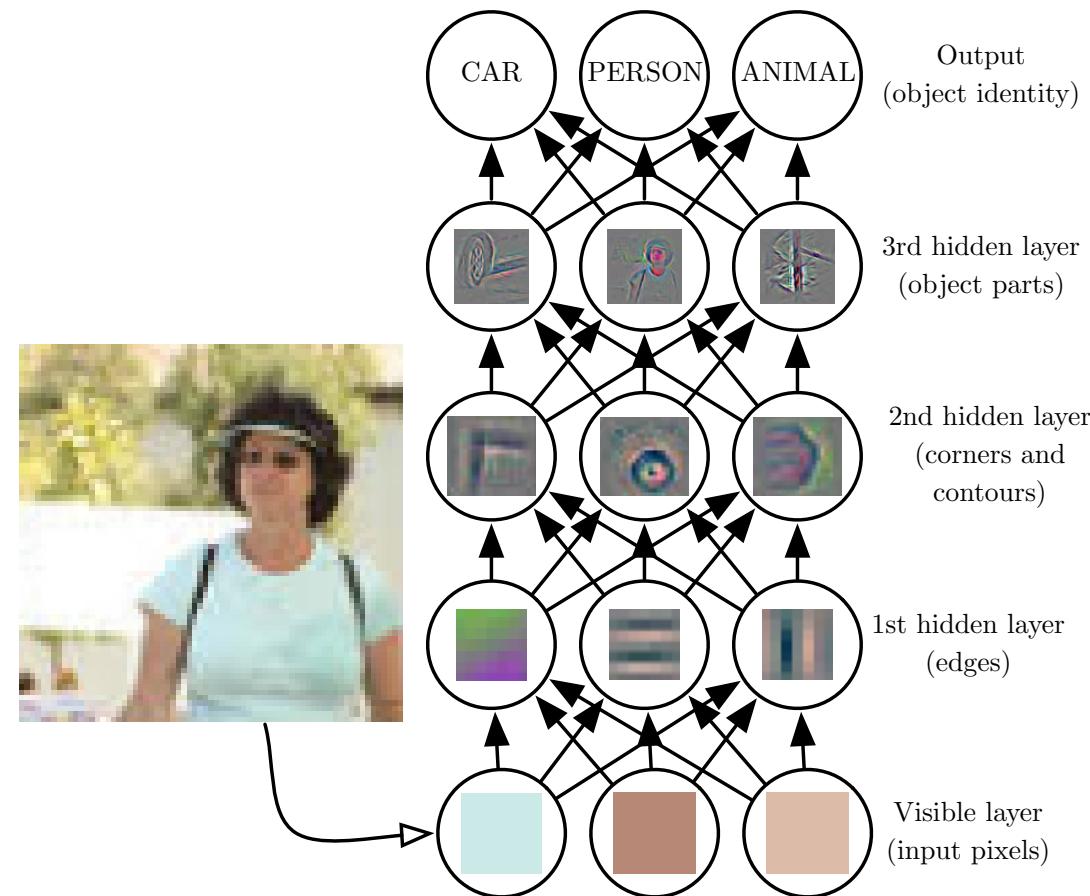
Deep learning frameworks

Why deep learning?

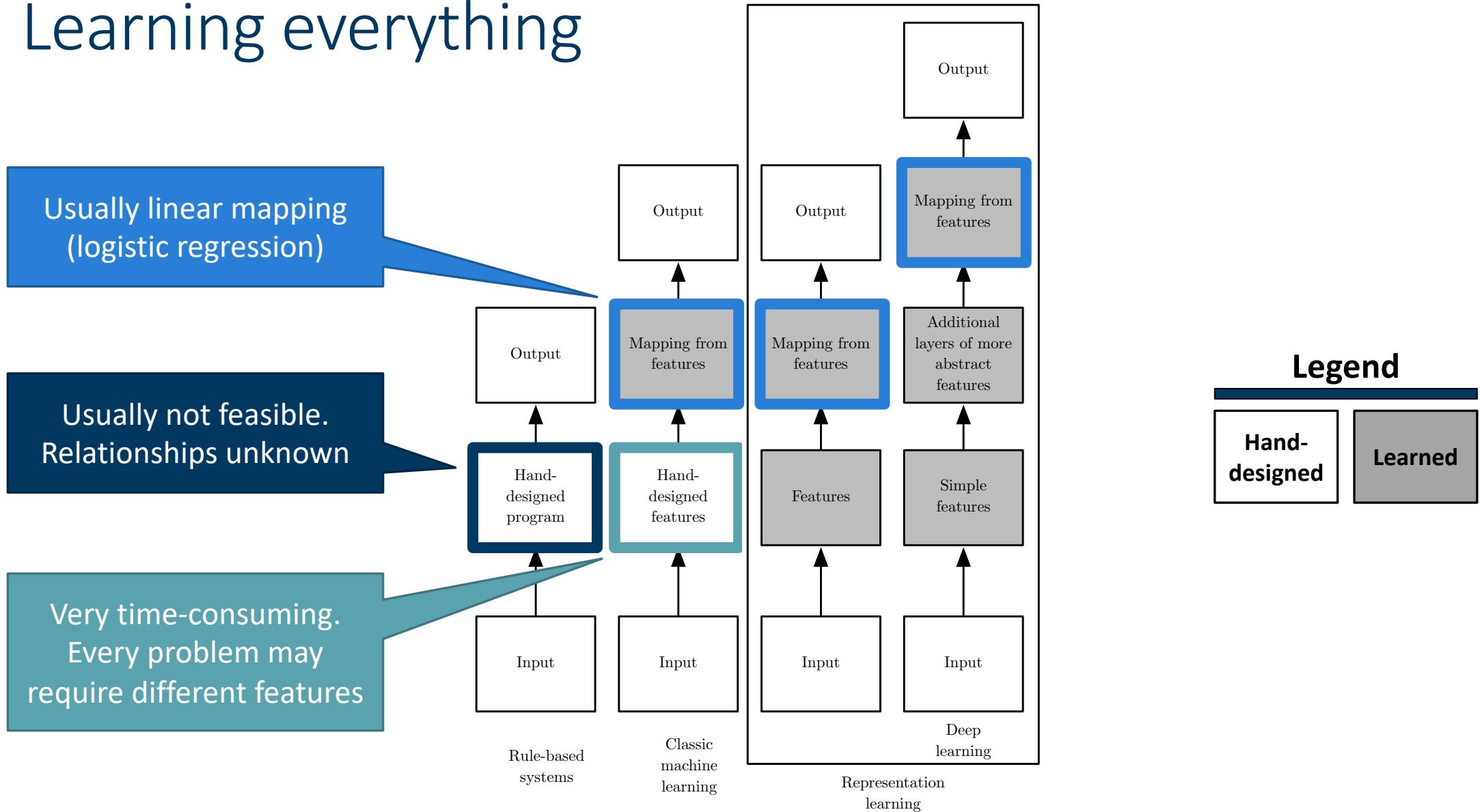
Representations matter



Depth: repeated composition

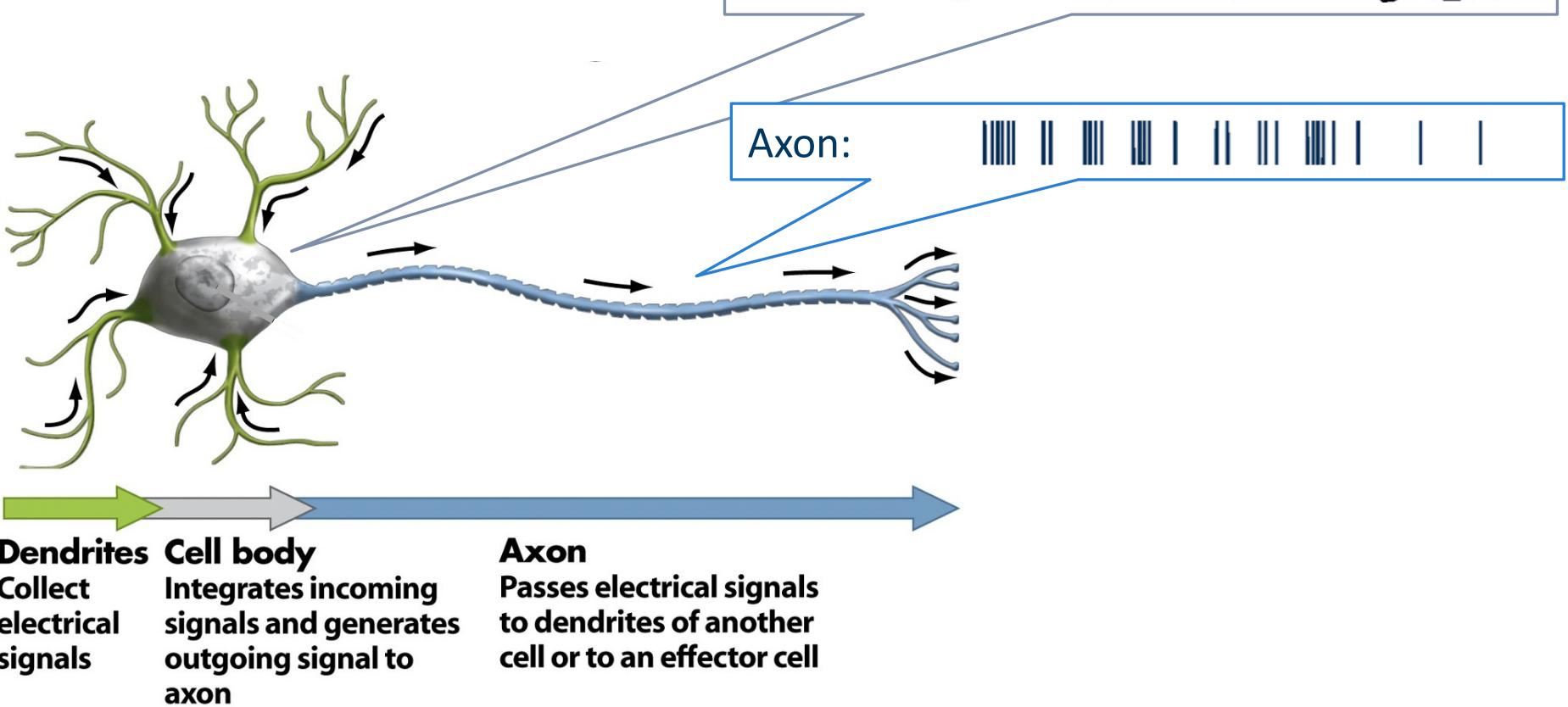


Learning everything

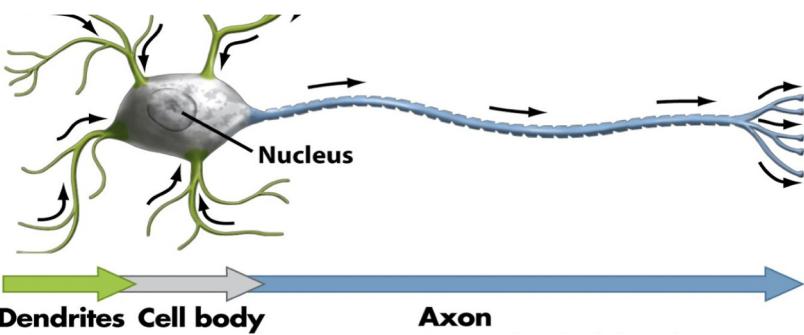
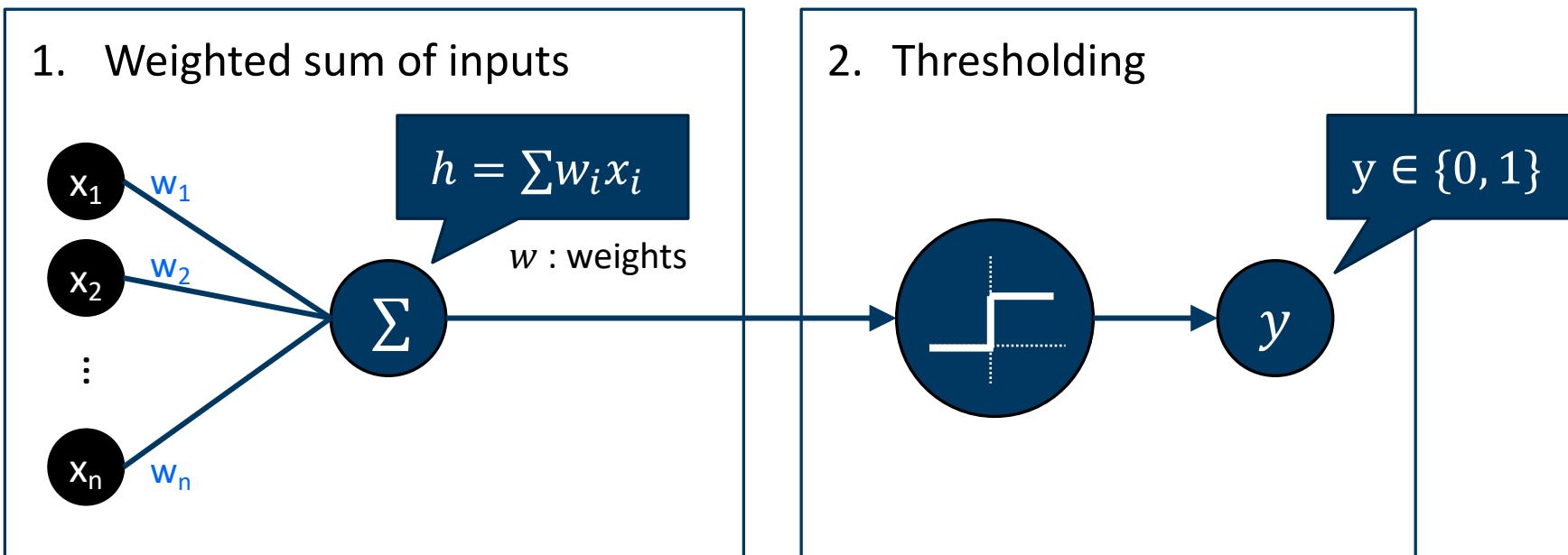


The Perceptron

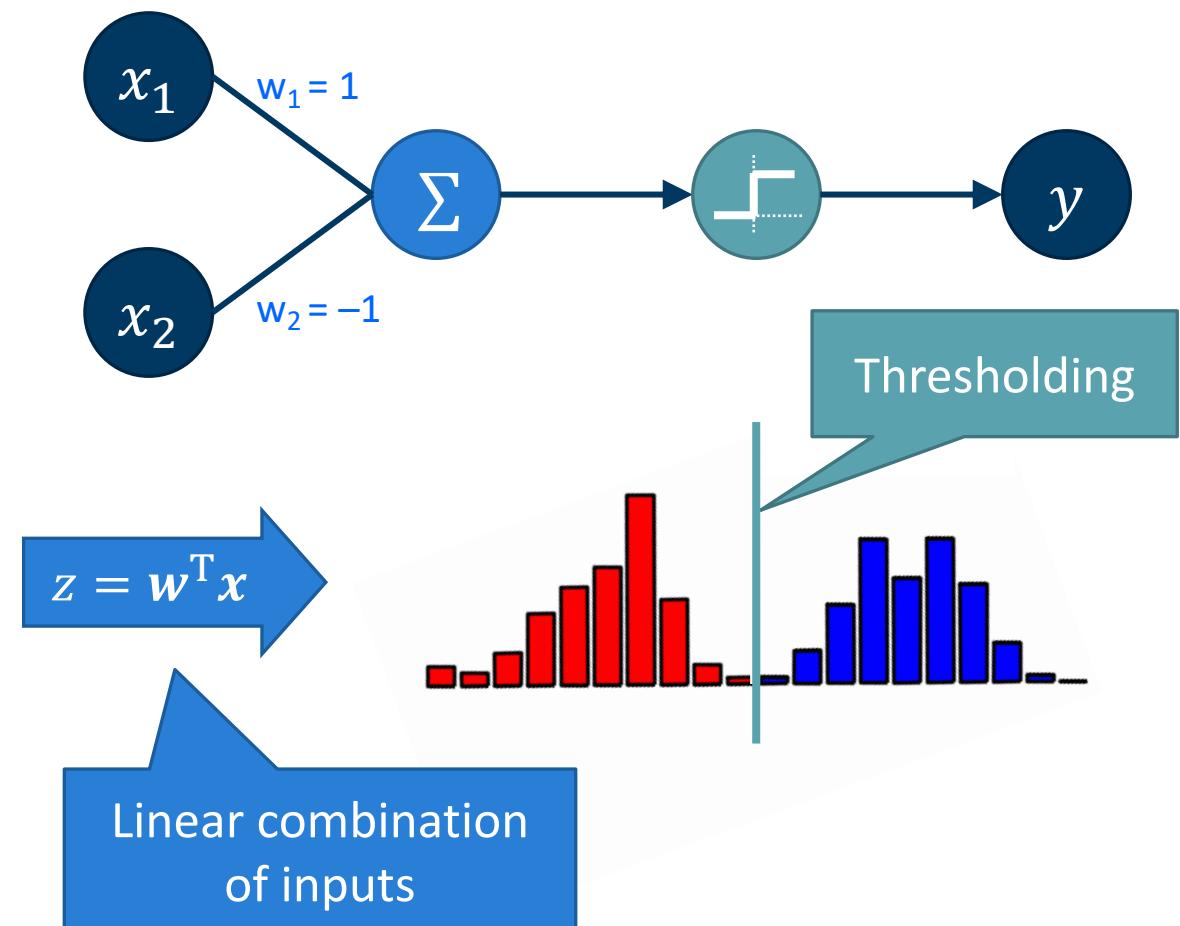
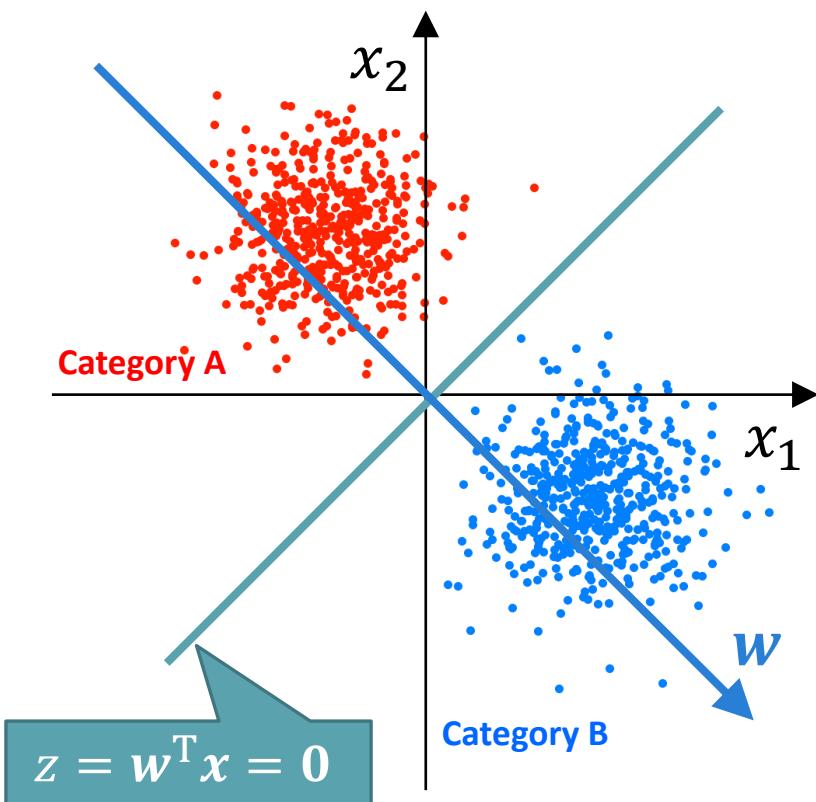
Neurons in the brain



The Perceptron (Rosenblatt 1958)



Perceptron solves linearly separable problems



Perceptron learning algorithm

Algorithm

$$\mathbf{w} \leftarrow \mathbf{0}$$

Iterate over training examples $(\mathbf{x}^{(i)}, y^{(i)})$ until convergence:

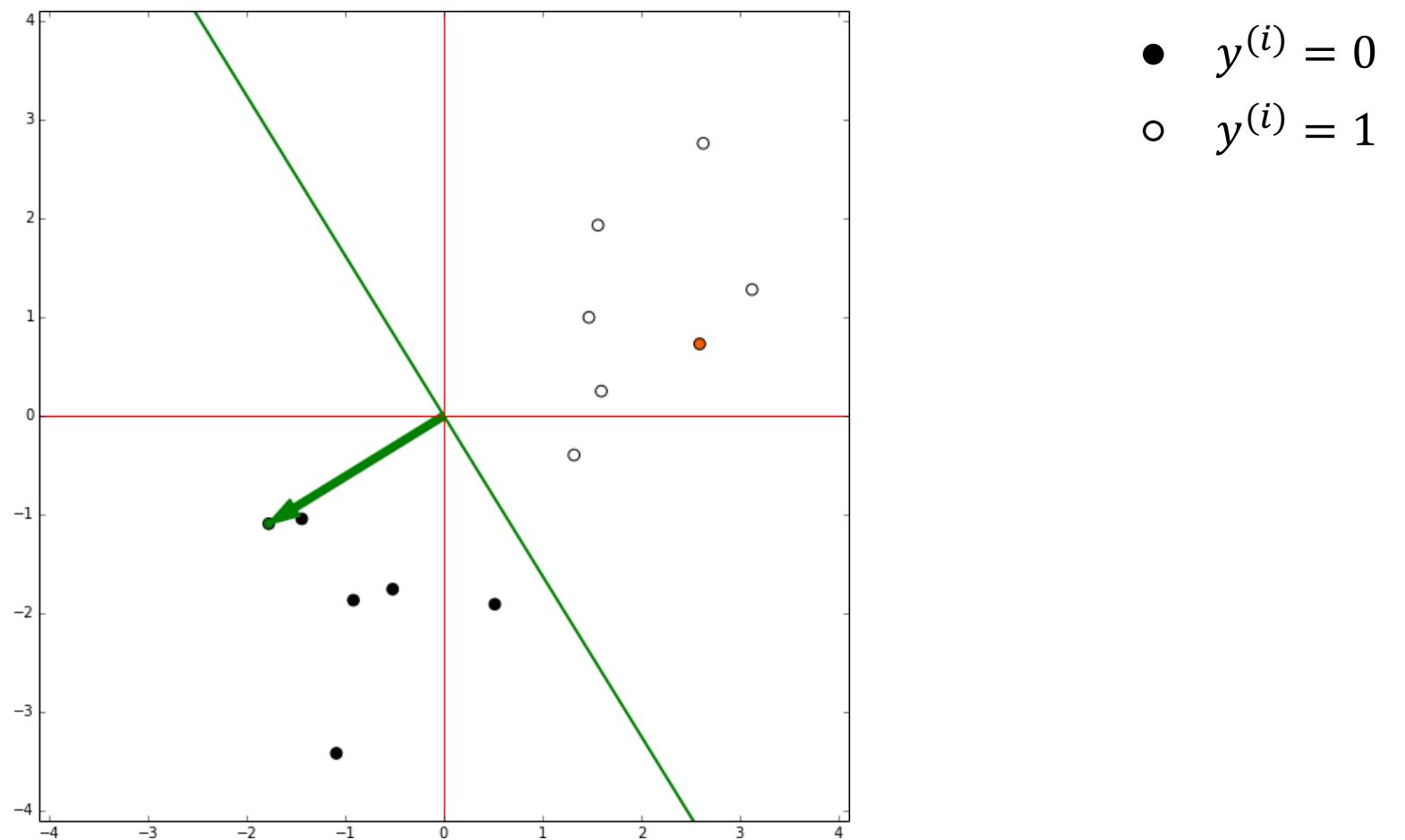
$$\hat{y}^{(i)} \leftarrow \mathbf{w}^T \mathbf{x}^{(i)}$$

$$e \leftarrow y^{(i)} - \hat{y}^{(i)}$$

$$\mathbf{w} \leftarrow \mathbf{w} + e \cdot \mathbf{x}$$

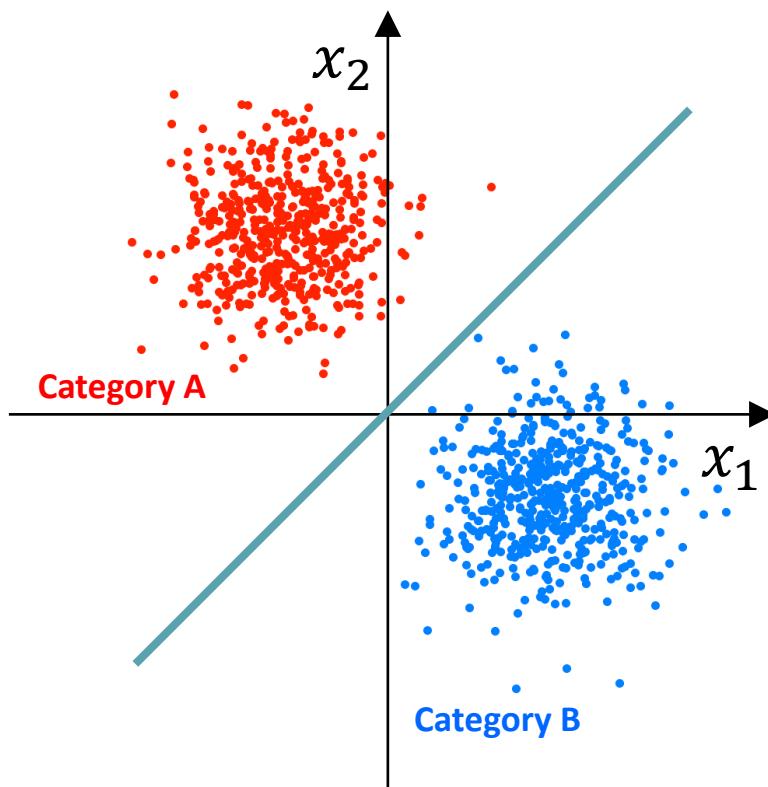
If all training examples can be classified correctly (problem is linearly separable), this algorithm provably converges to a correct solution in a finite number of steps.

Perceptron learning animated

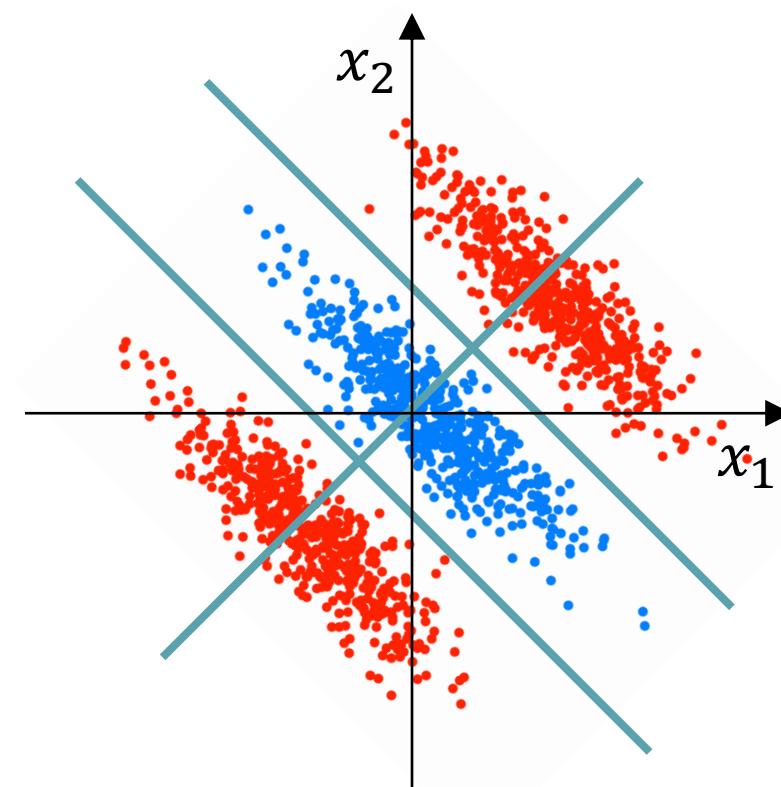


Linear (in)separability

Linearly separable

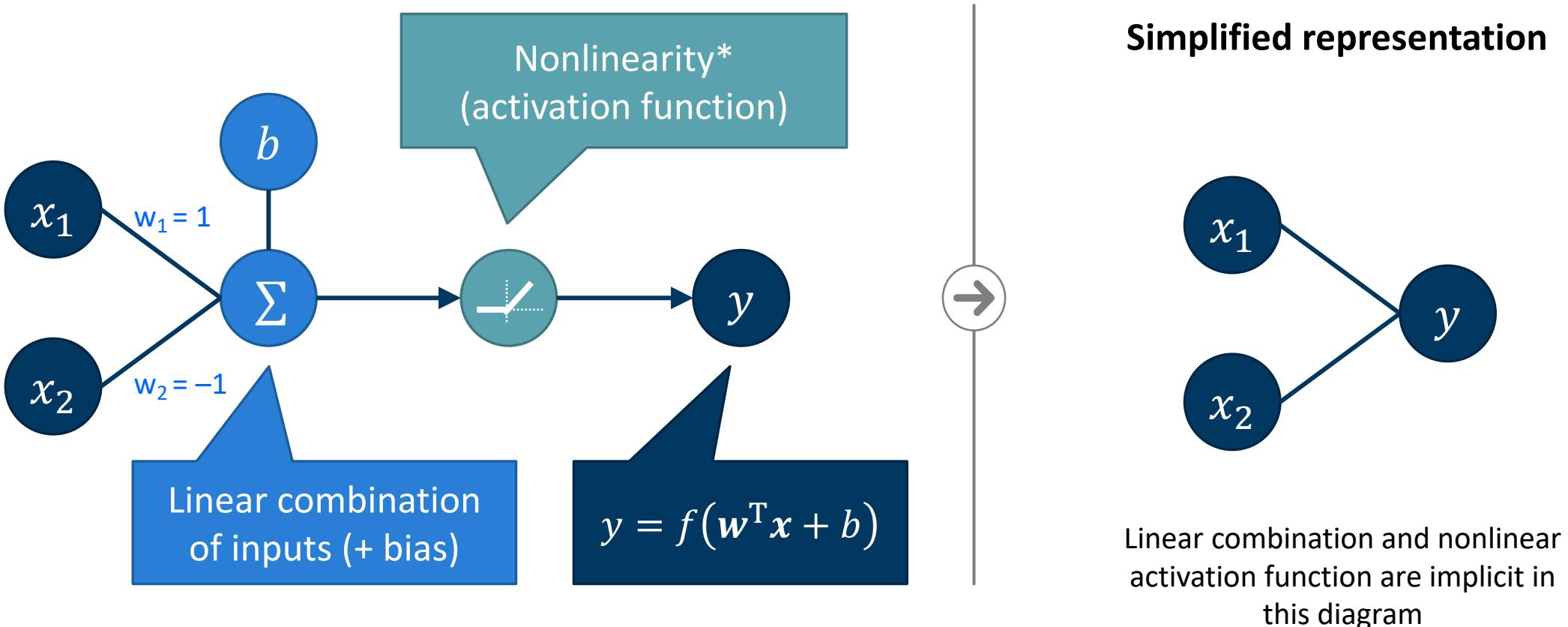


Not linearly separable



Multi-layer perceptron

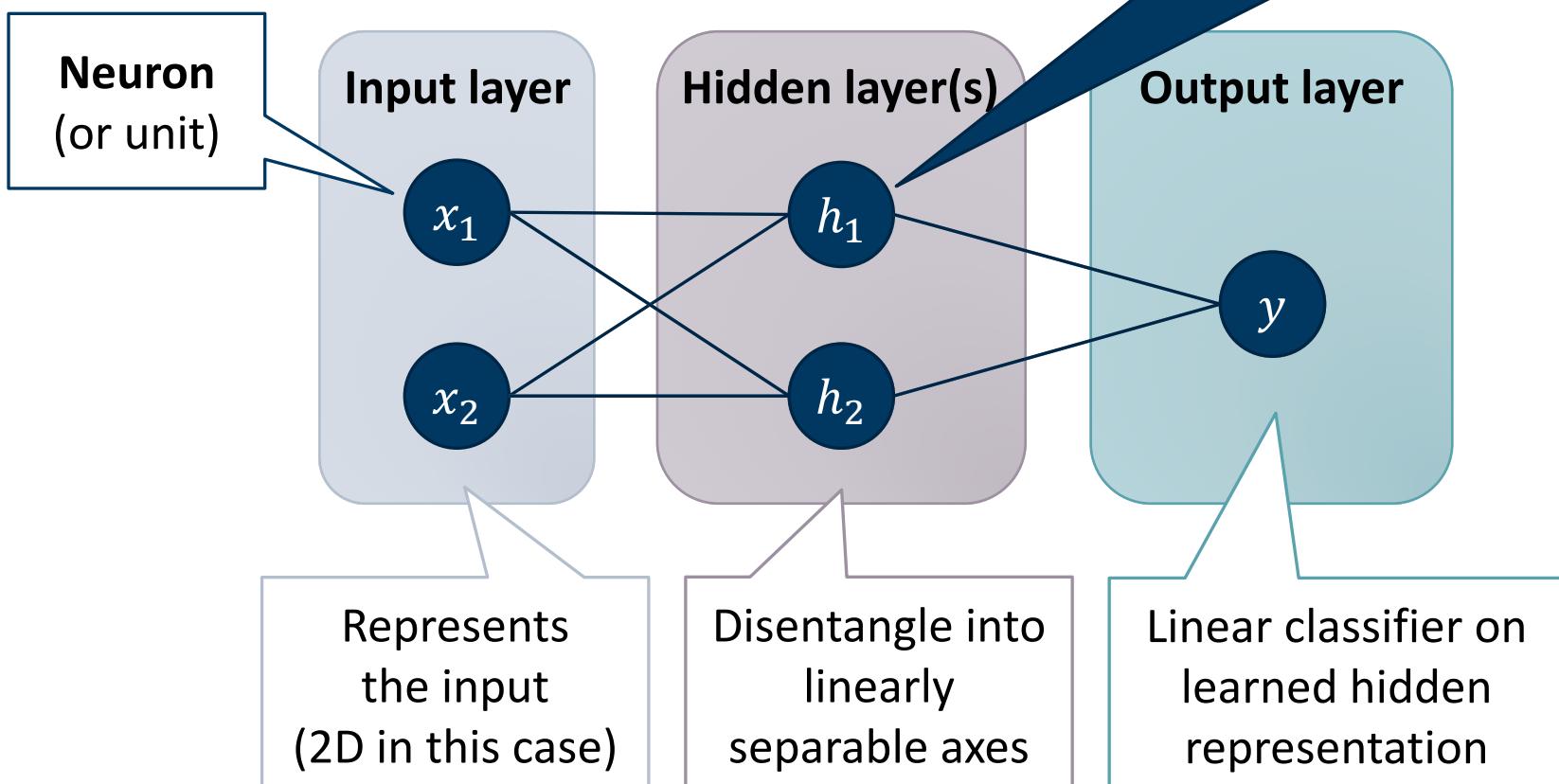
Perceptron: more generally



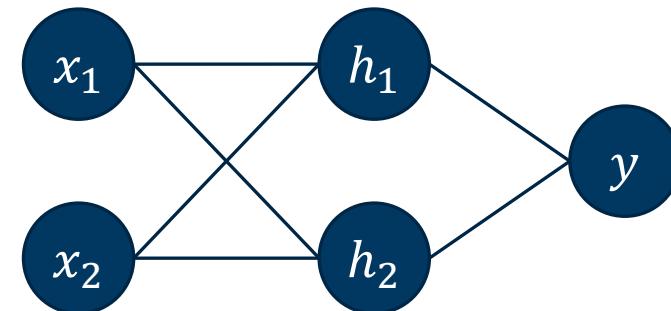
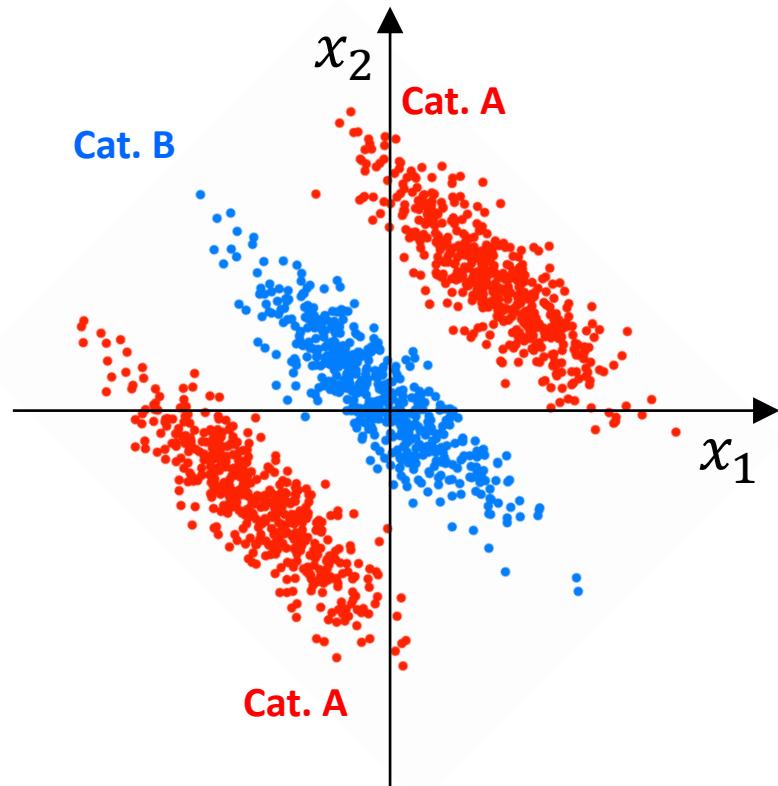
* The original perceptron and its learning rule used the thresholding operation as activation function. The term is now being used more broadly (see e.g. multi-layer perceptron in a few slides).

Multi-layer perceptron

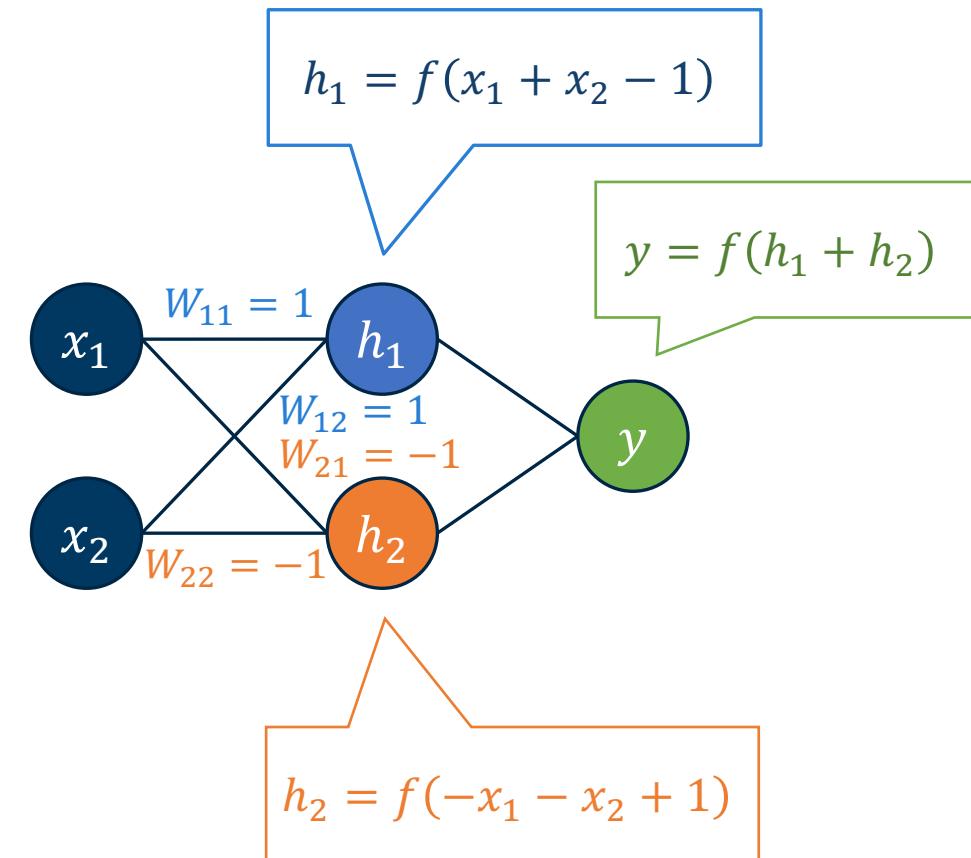
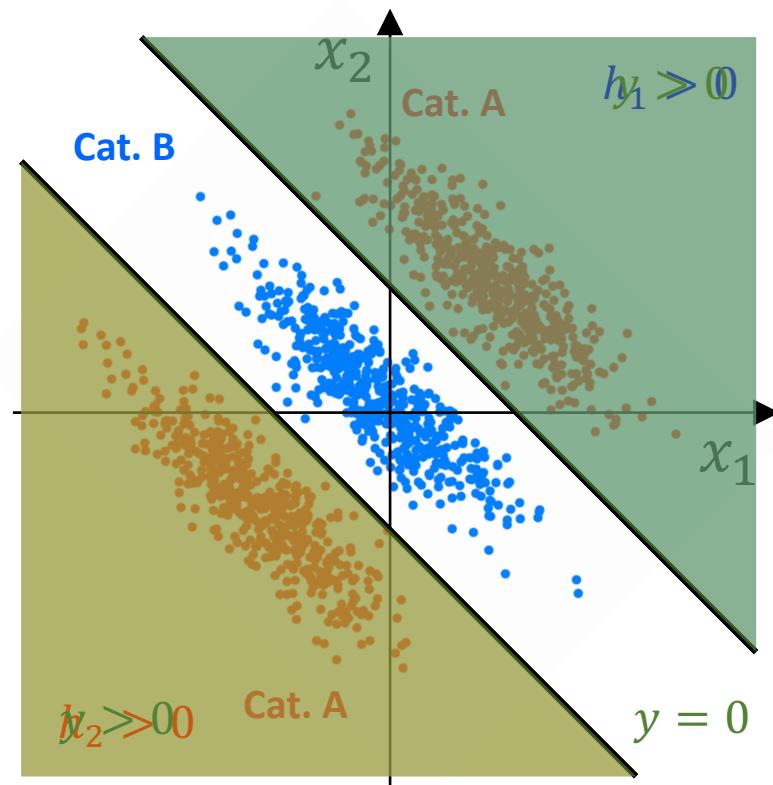
Each unit is a Perceptron:
$$h_1 = f(W_{1:}x + b_1)$$



Solving the toy problem with a two-layer perceptron

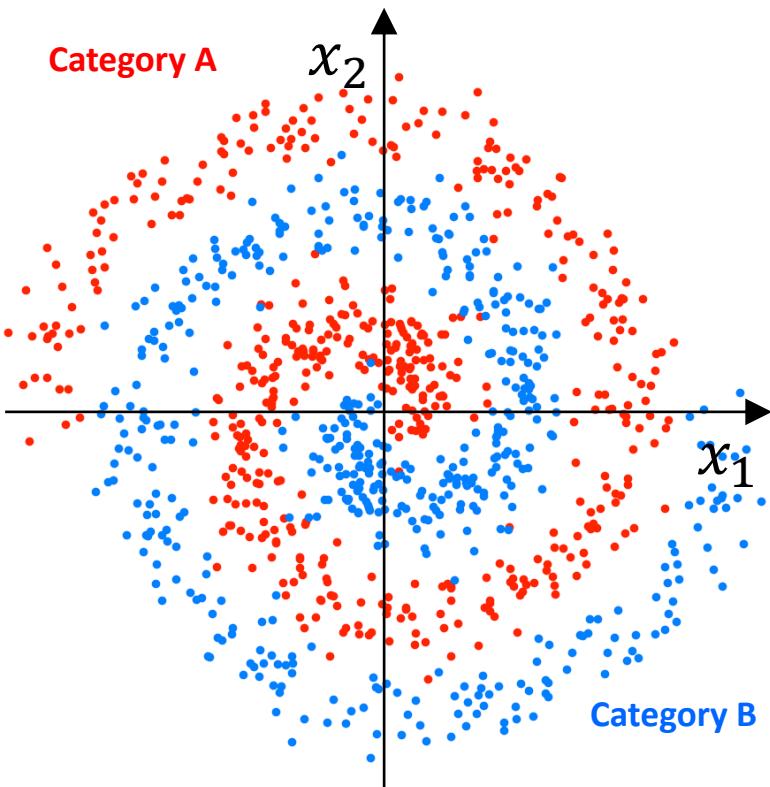


Solving the toy problem with a two-layer perceptron

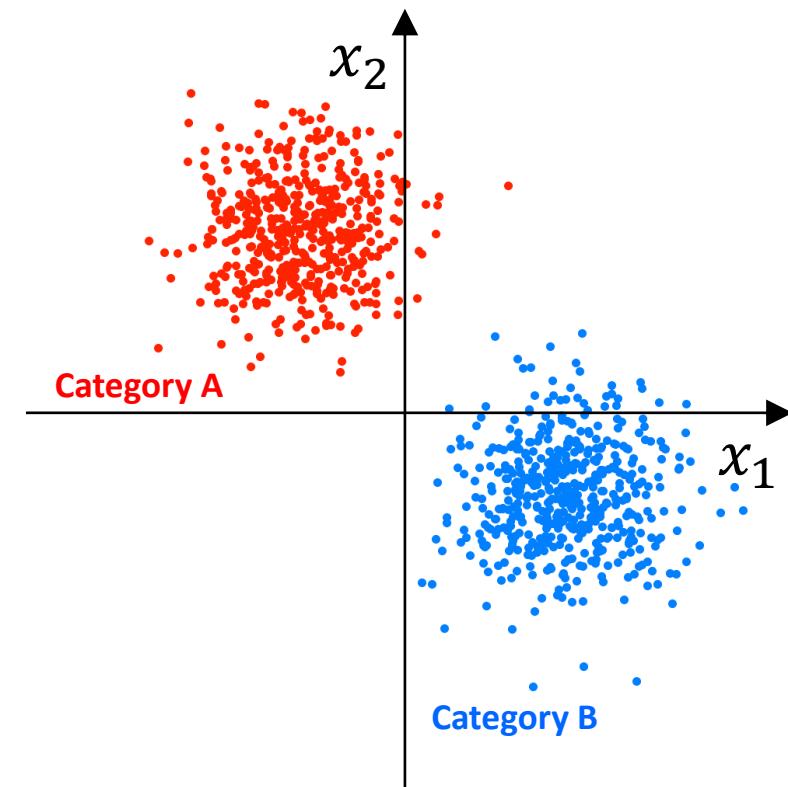


Untangling the manifold: More complex case

Entangled \rightarrow Not linearly separable



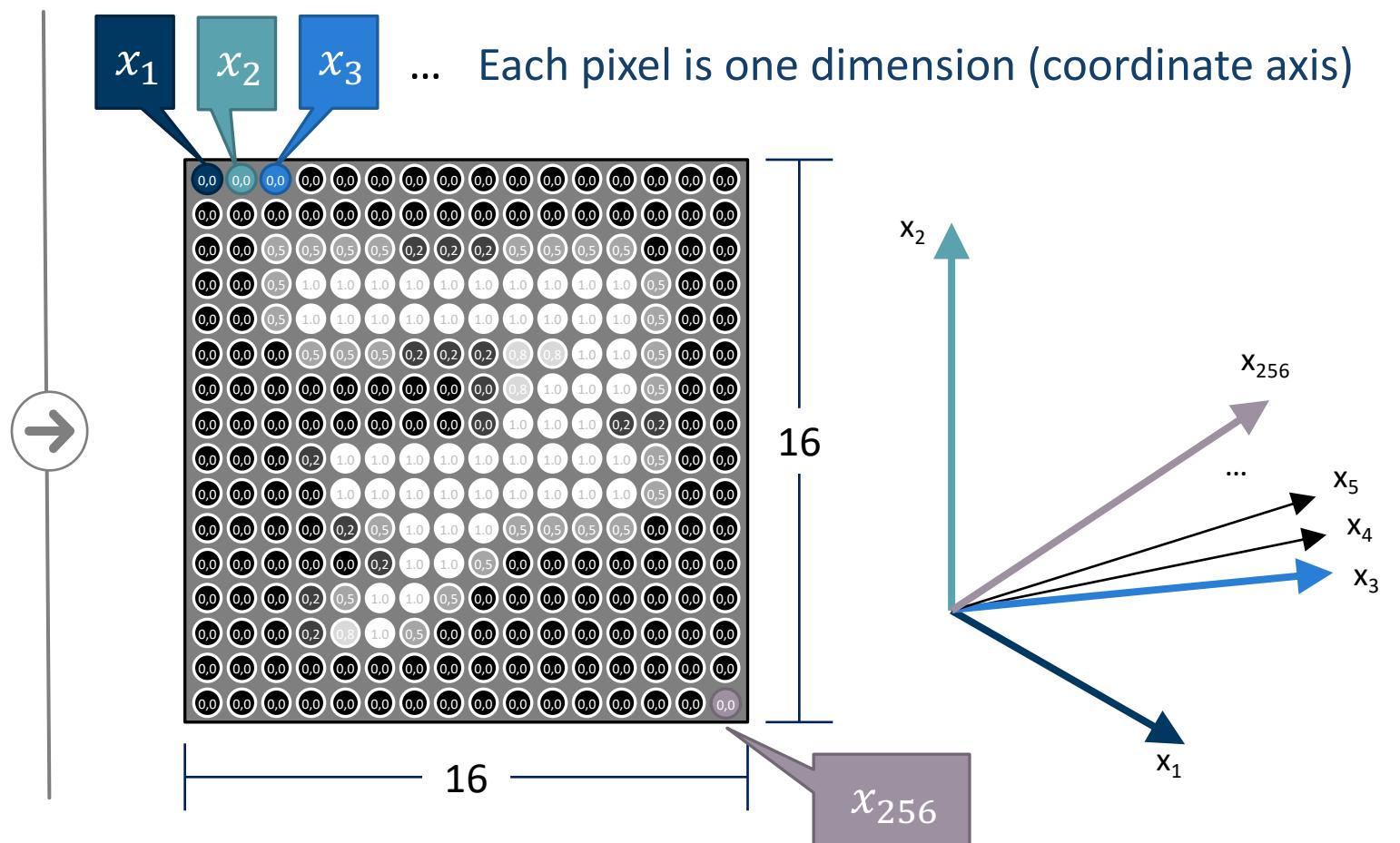
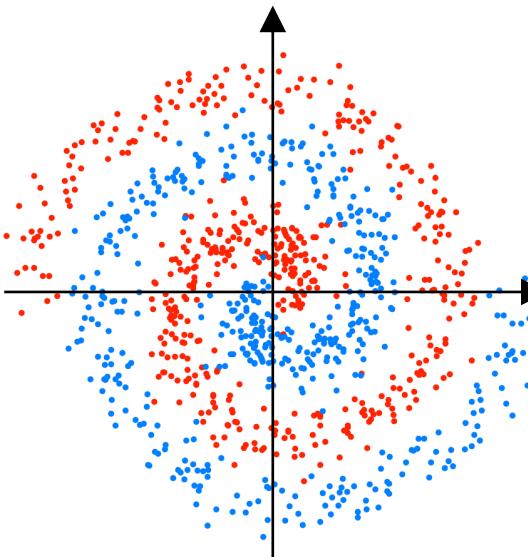
Disentangled \rightarrow Linearly separable



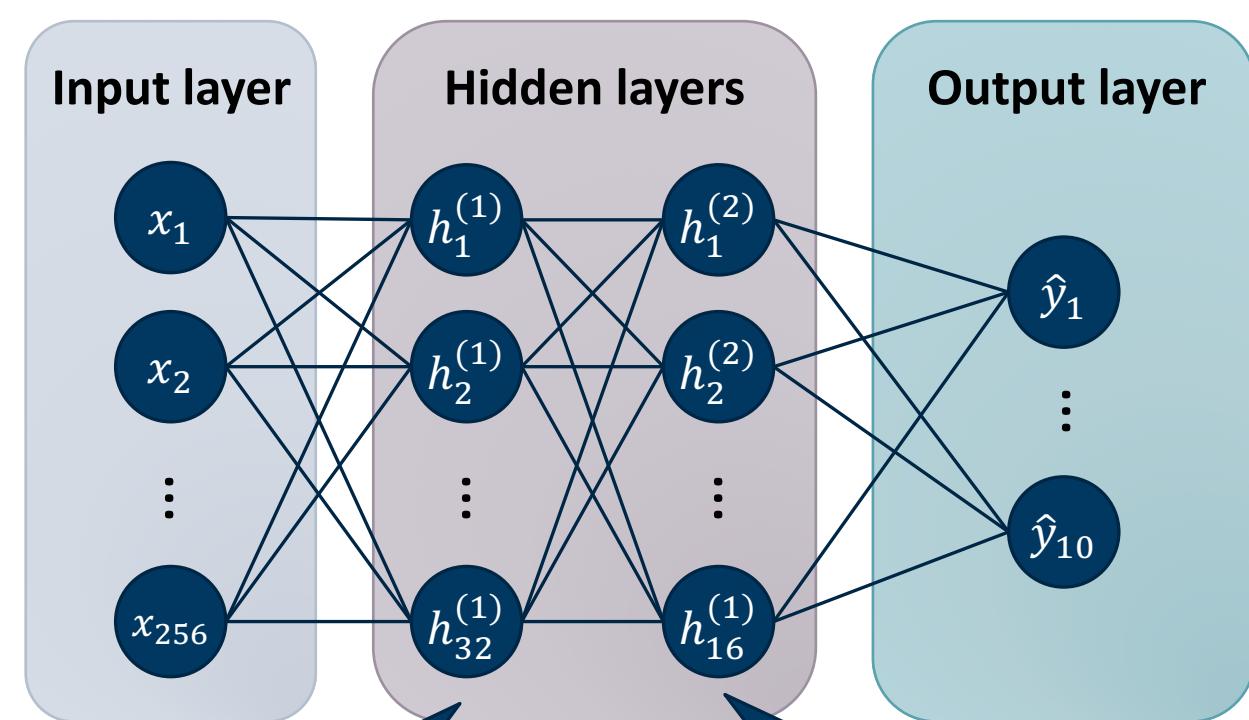
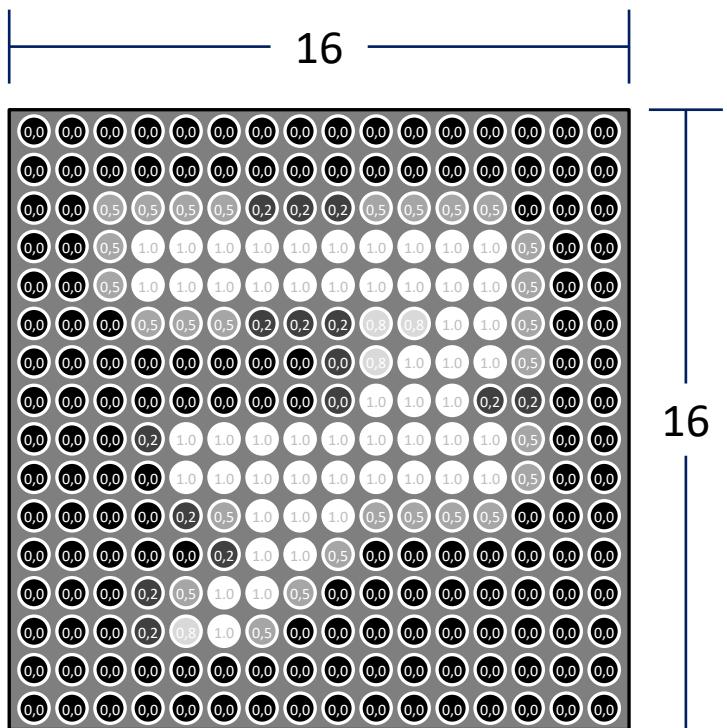
Disentangle
through multiple
hidden layers

2D to many-D: same principle, harder to visualize

2D toy example so far



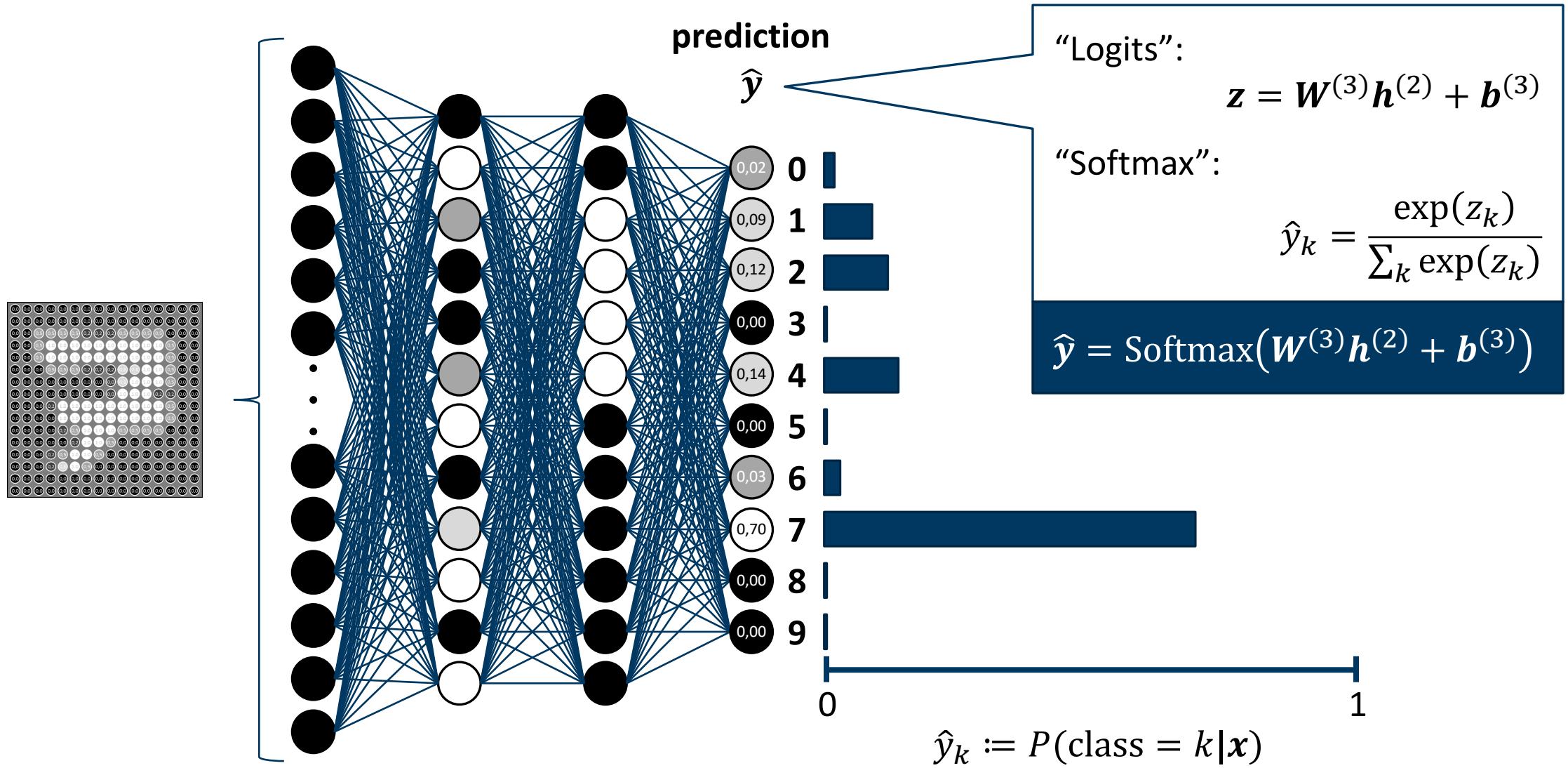
Multi-layer perceptron (MLP) for classifying handwritten digits



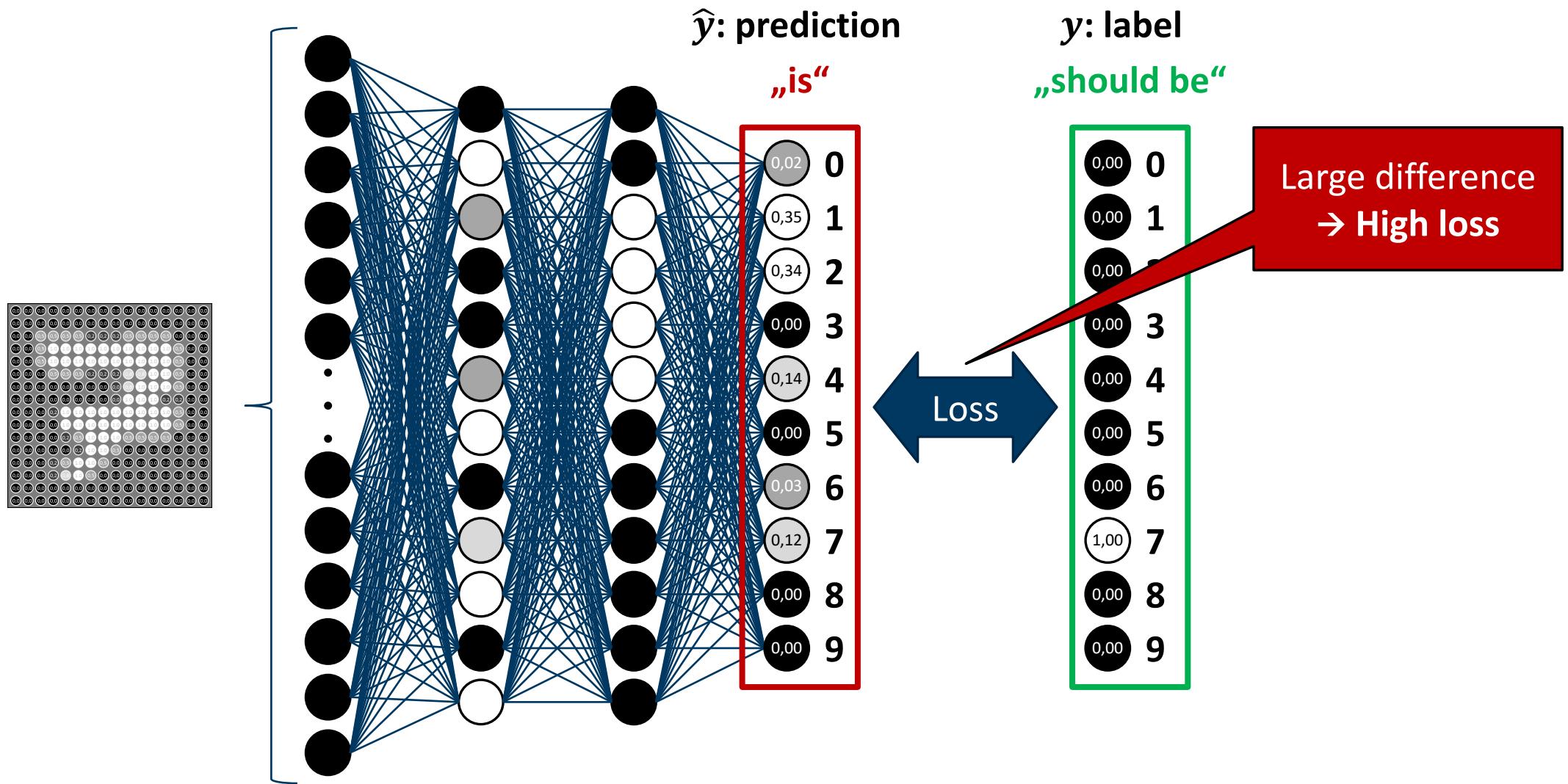
$$\mathbf{h}^{(1)} = f(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$$

$$\mathbf{h}^{(2)} = f(\mathbf{W}^{(2)}\mathbf{h}^{(1)} + \mathbf{b}^{(2)})$$

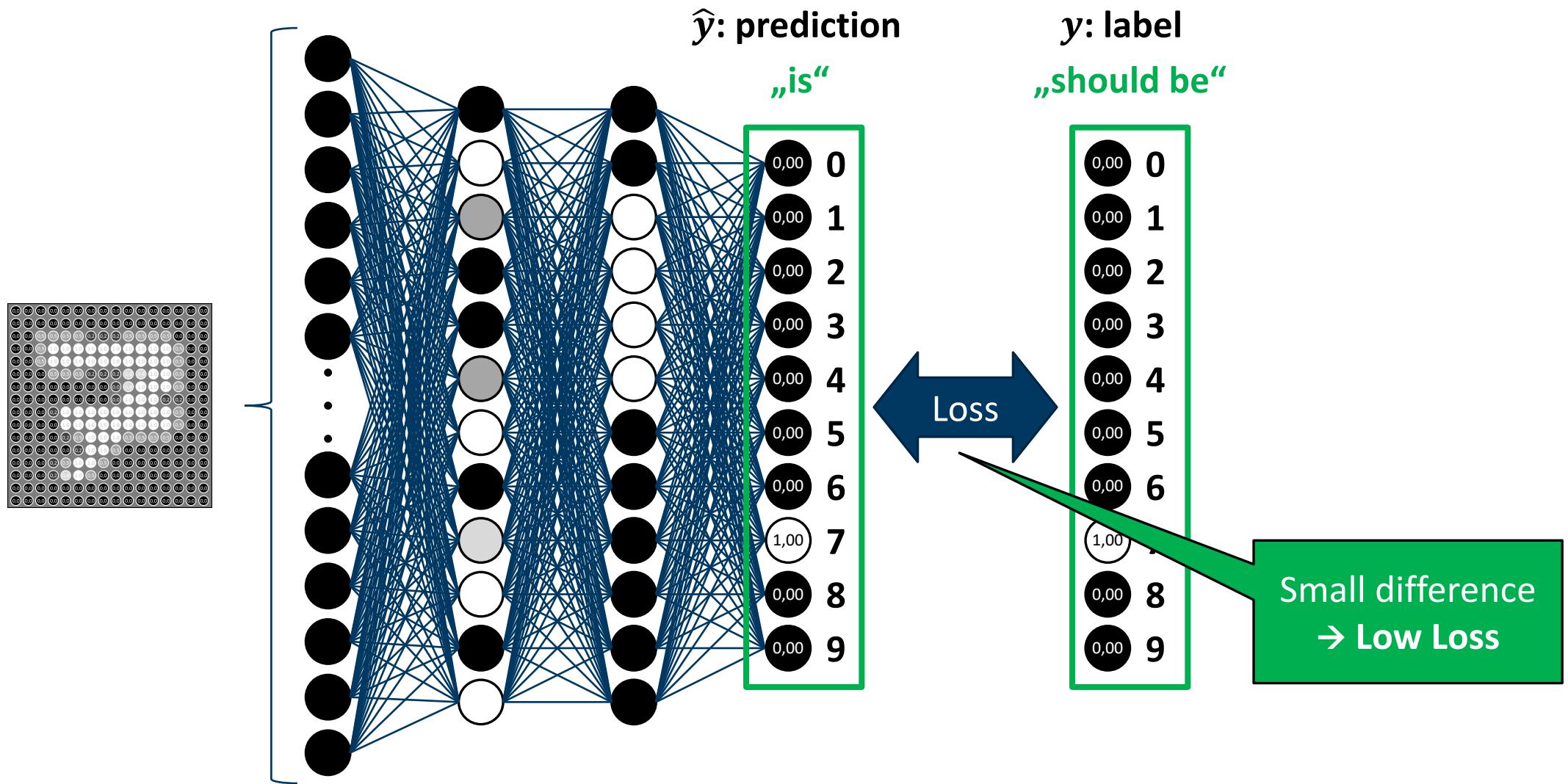
Output layer: probability distribution over classes



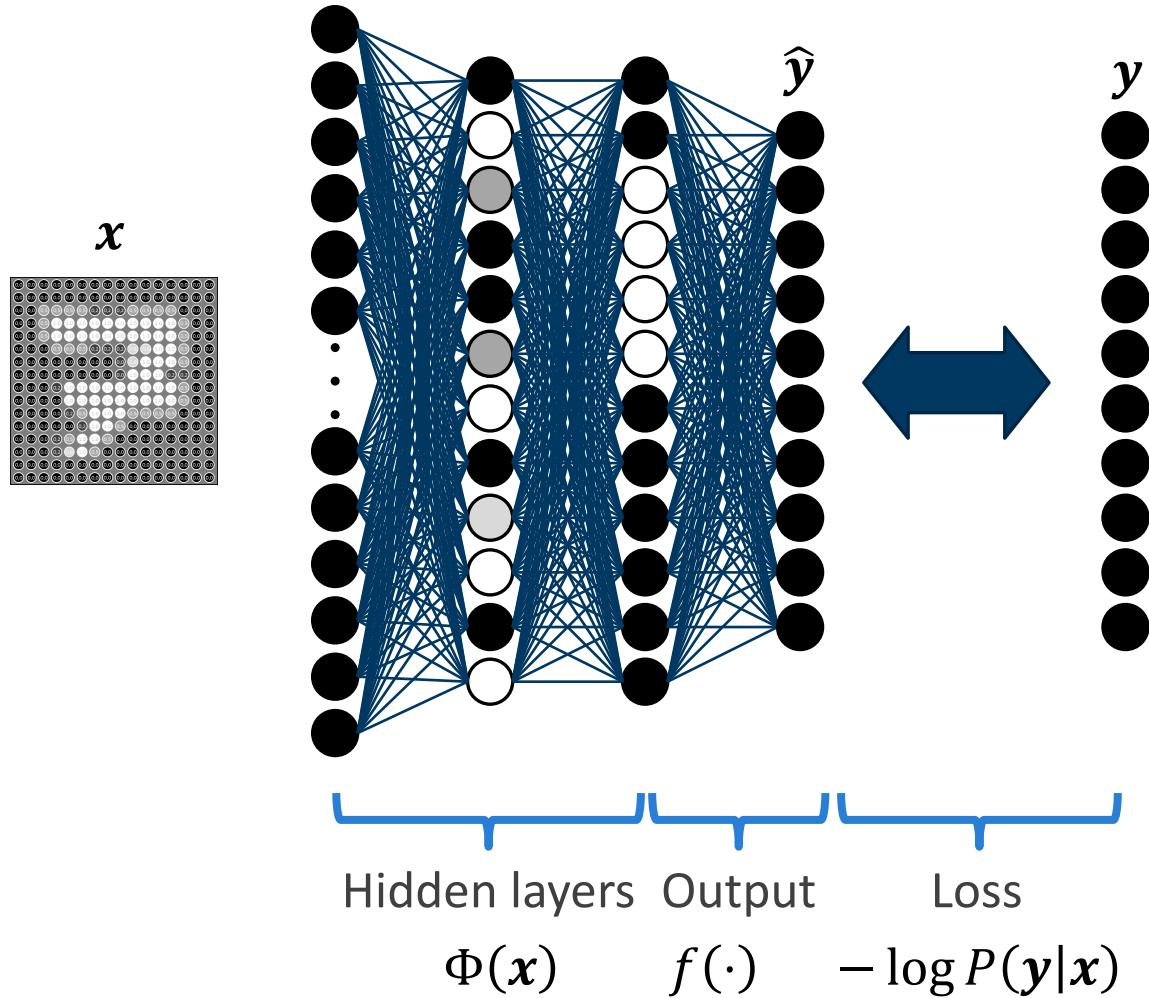
Loss function measures how well the net performs



Loss function measures how well the net performs



More generally: output layers and loss functions



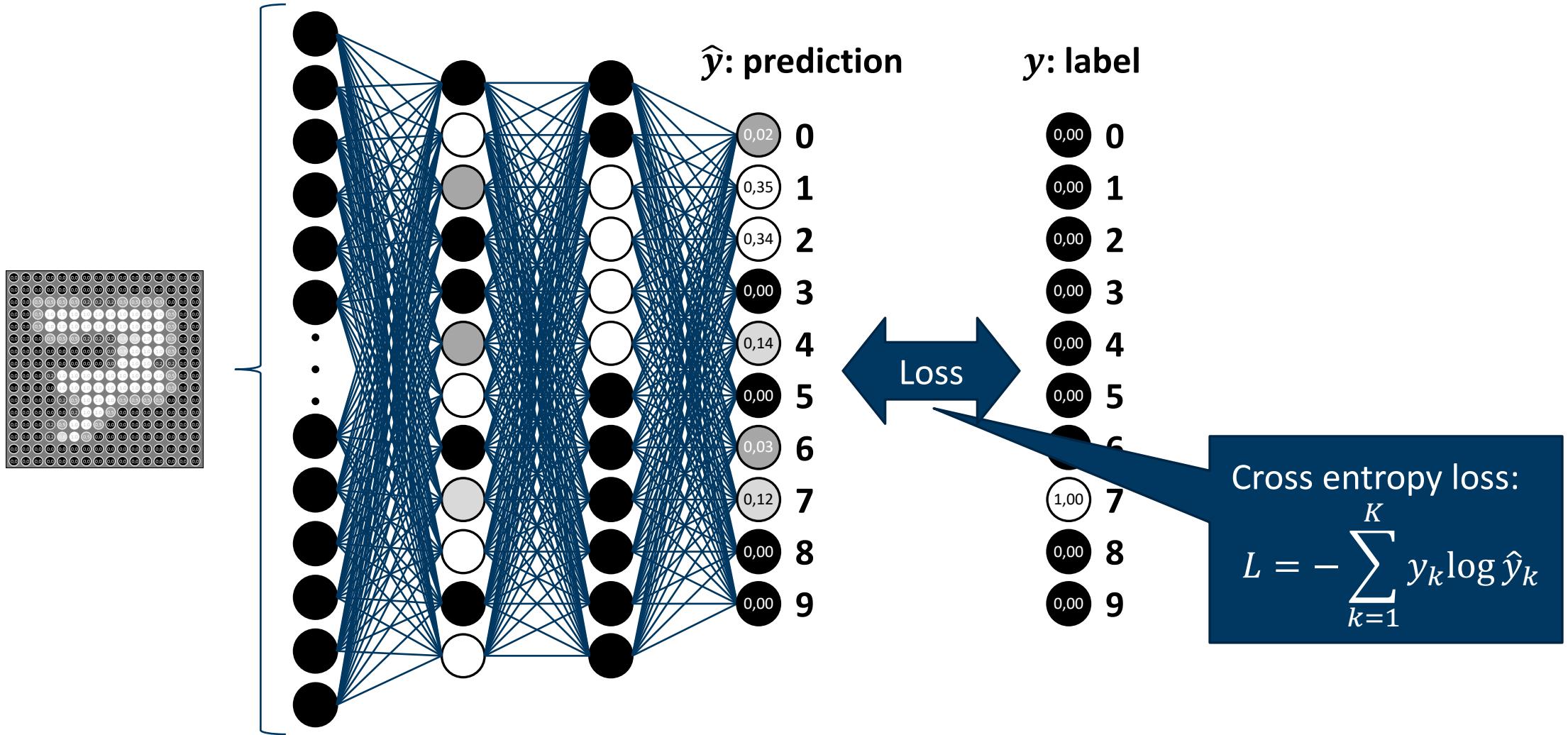
Model

$$P(y|x) = f(\Phi(x))$$

Loss

$$L = -\log P(y|x)$$

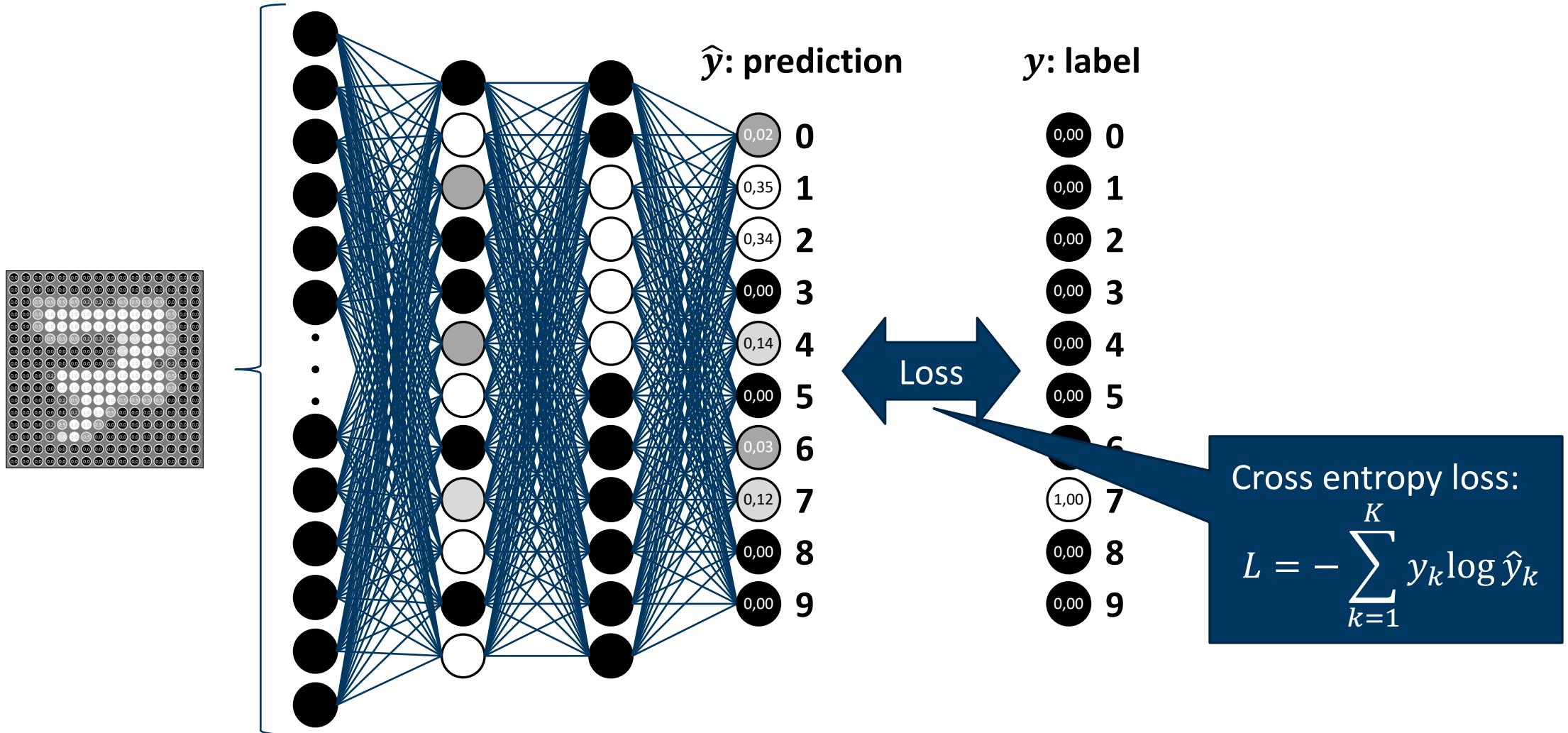
Loss function for classification: cross entropy



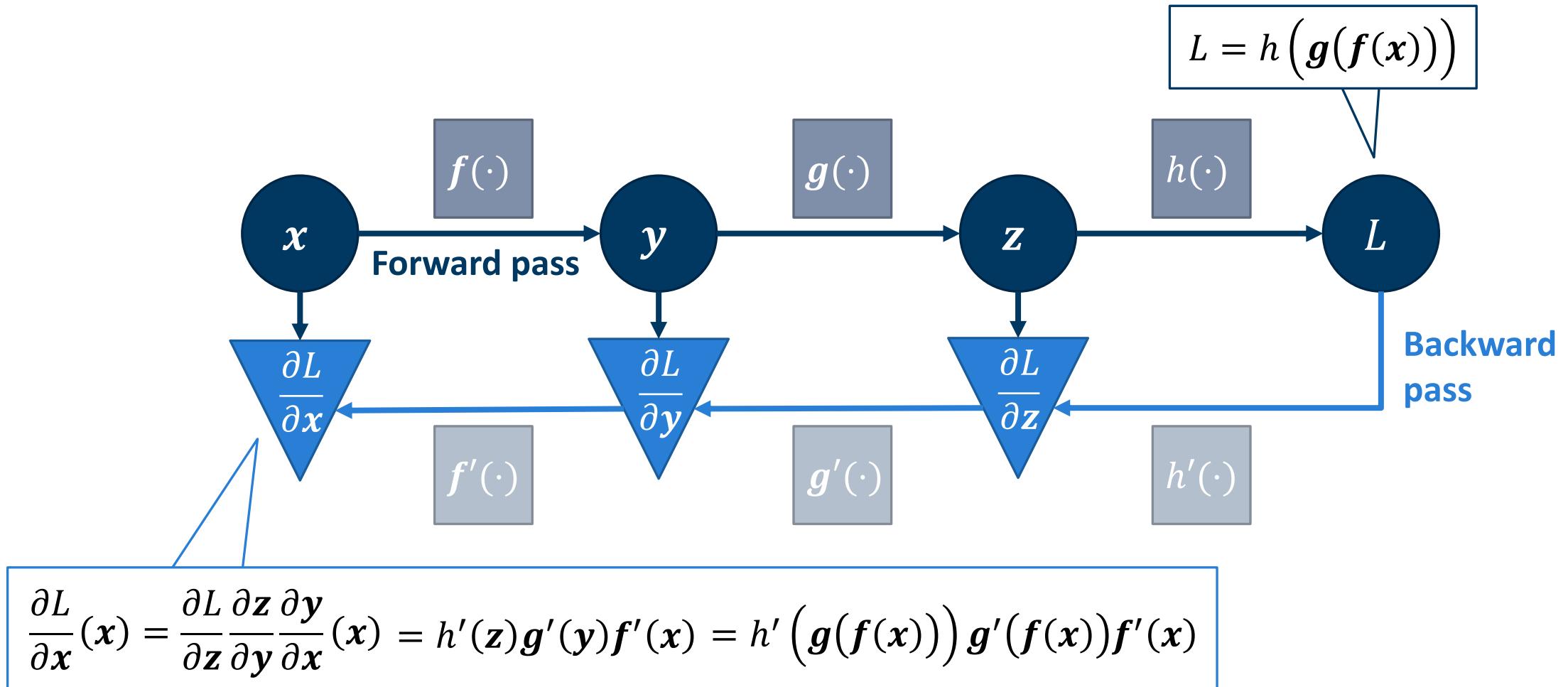
Types of output layer and loss functions

Type of variable	Distribution	Output unit	Activation function	Loss function
y	$P(y x)$		$f(\cdot)$	$-\log P(y x)$
Continuous	Gaussian	Linear	$f(z) = z$	$(\hat{y} - y)^2$
Binary	Bernoulli	Sigmoid	$f(z) = \frac{1}{1 + \exp(-z)}$	$y \log \hat{y}$
Categorical	Multinomial	Softmax	$f(z_i) = \frac{\exp(z_i)}{\sum_k \exp(z_k)}$	$\sum_k y_k \log \hat{y}_k$

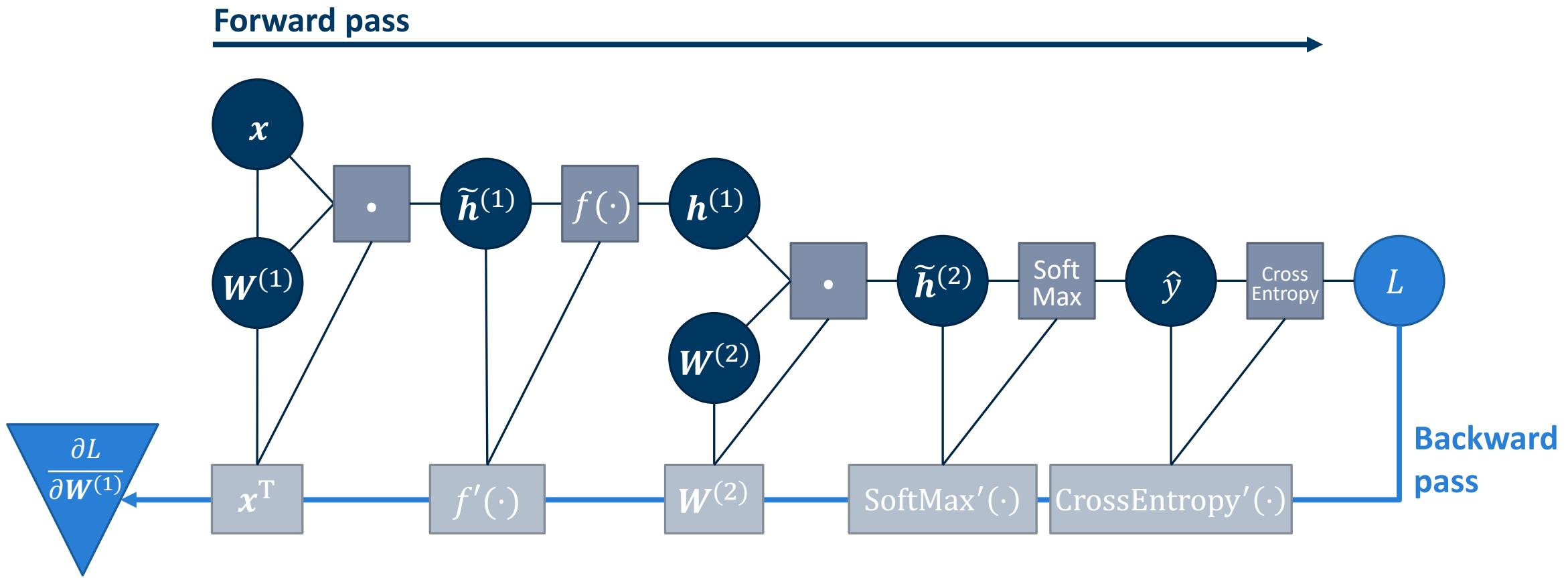
How to optimize the loss function?



Computing gradients: chain rule → backpropagation

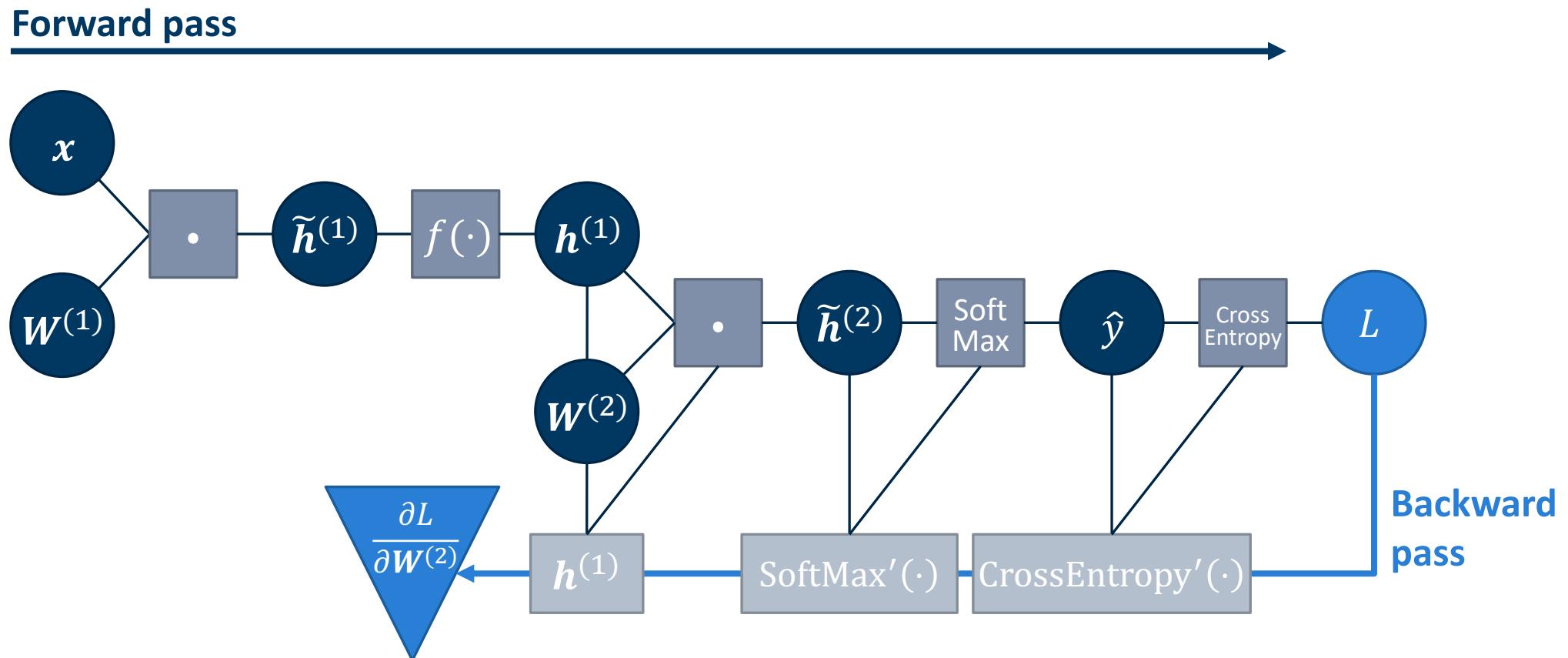


Computational graph



* For simplicity we're ignoring the bias terms in all layers

Computational graph



* For simplicity we're ignoring the bias terms in all layers

Automatic differentiation

Each operation „knows“ its gradient w.r.t. its inputs

You can compose arbitrary chains of operations (computational graph)

Gradients get computed automatically via chain rule

Modern deep learning frameworks
like *PyTorch*, *Tensorflow* etc. implement automatic differentiation

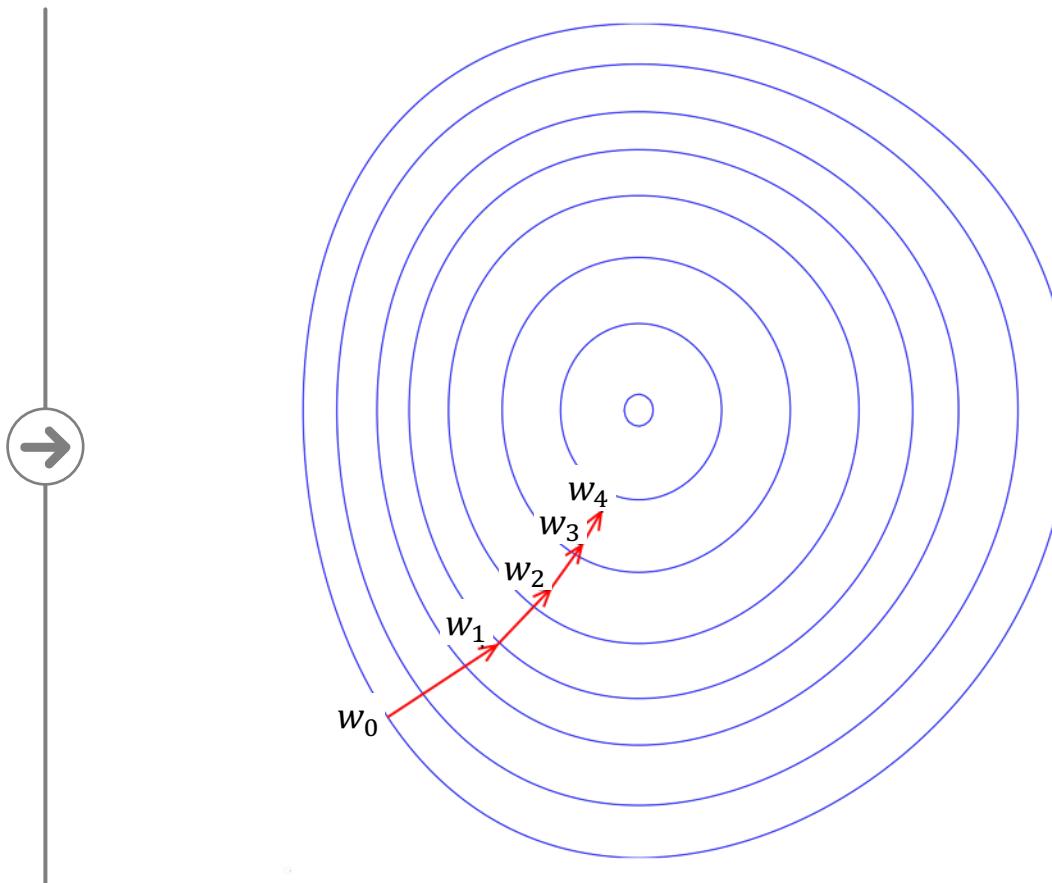
Stochastic gradient descent

Starting point: gradient descent to minimize the loss function

Start with a random guess for the parameters w

Iteratively update w by a small step in the direction of the gradient of the loss function

$$w \leftarrow w - \lambda \nabla_w L(w)$$



Stochastic gradient descent

Problem: Computing gradients for large datasets is expensive

BUT: Loss functions typically decompose into sums of losses per data point

$$L(\{x_n\}_{n=1}^N, w) = \sum_{n=1}^N L(x_n, w) \quad \text{where often} \quad L(x_n, w) = -\log p(y_n|x_n, w)$$

$$\nabla_w L(\{x_n\}_{n=1}^N, w) = \sum_{n=1}^N \nabla_w L(x_n, w)$$

↑
Entire dataset ↑
One data point

Alternative: minibatch gradient descent only uses a subset of all data points in each iteration

Idea: Use only a subset of data points in each iteration

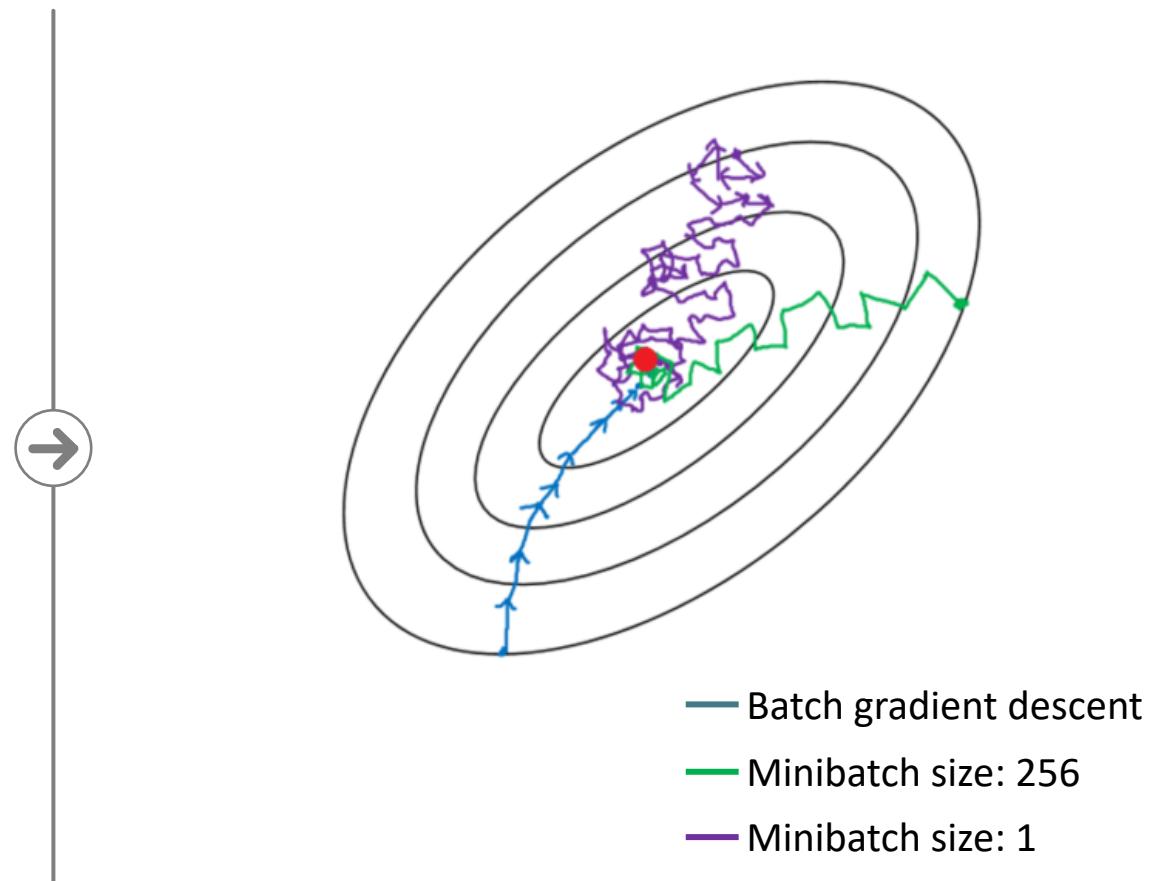
“Pure” stochastic gradient descent
(minibatch size 1)

$$w \leftarrow w - \lambda \nabla_w L(x_i, w)$$

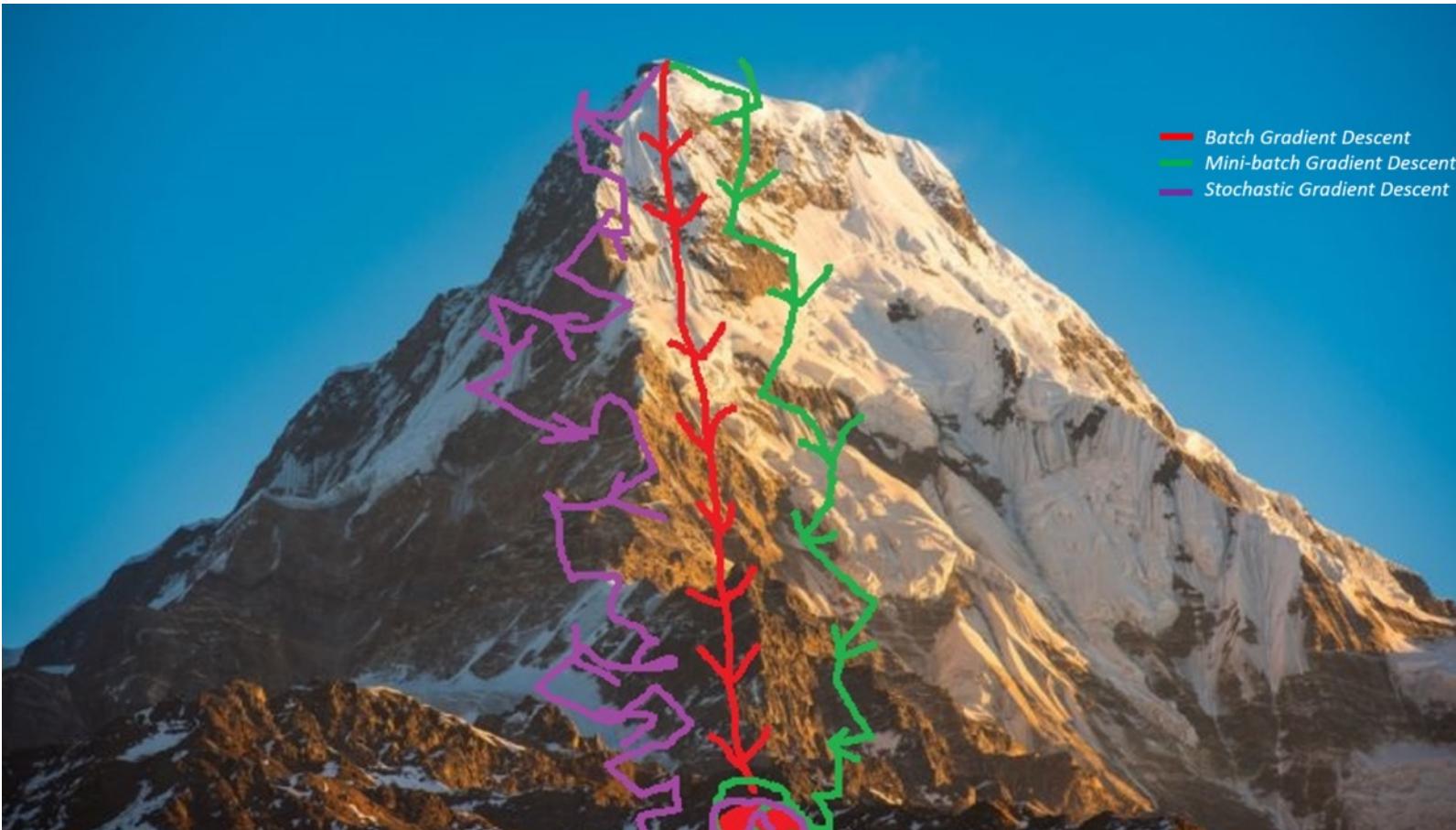
Minibatch gradient descent

$$w \leftarrow w - \lambda \nabla_w L(\{x_i\}_{i \in \mathcal{B}}, w)$$

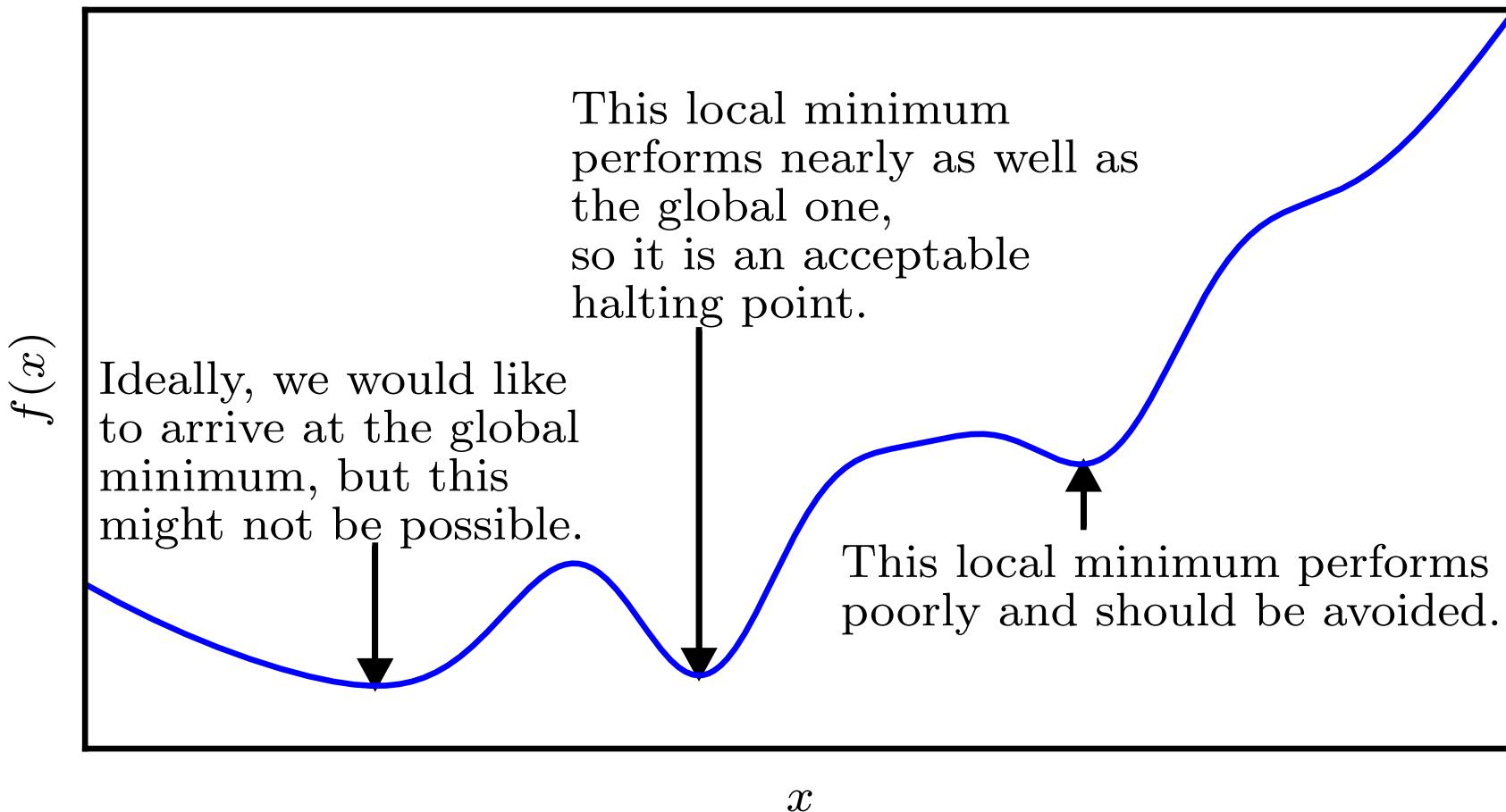
↑
Minibatch



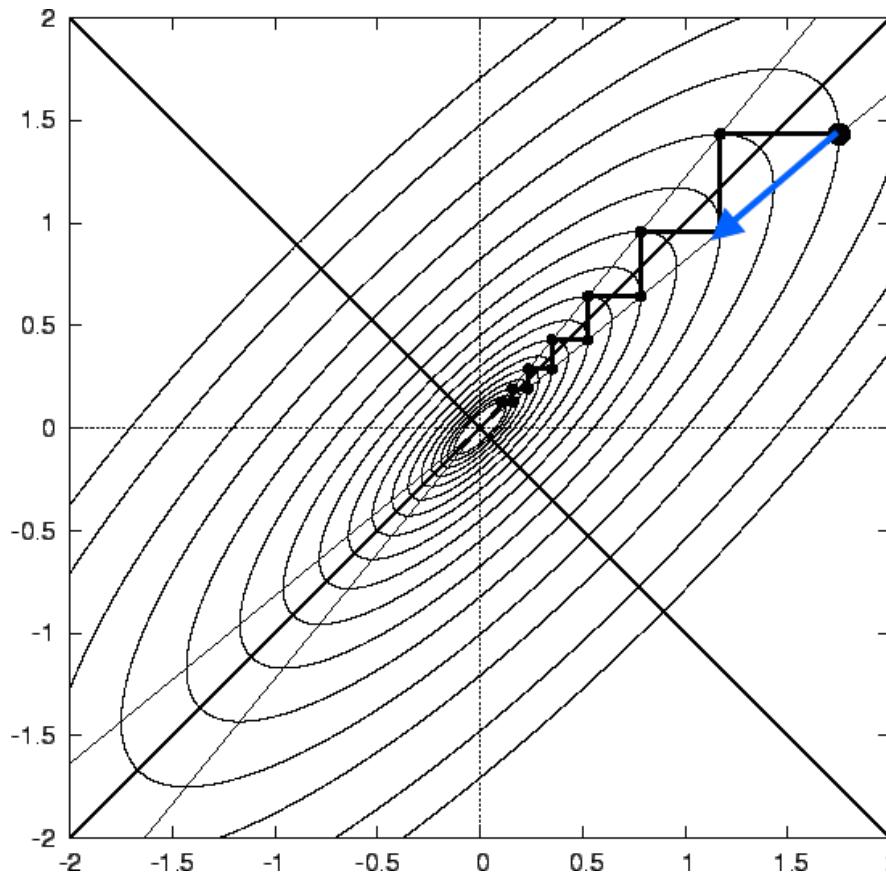
Different gradient descent strategies



Gradient-based optimization does neither always find the absolute minimum...



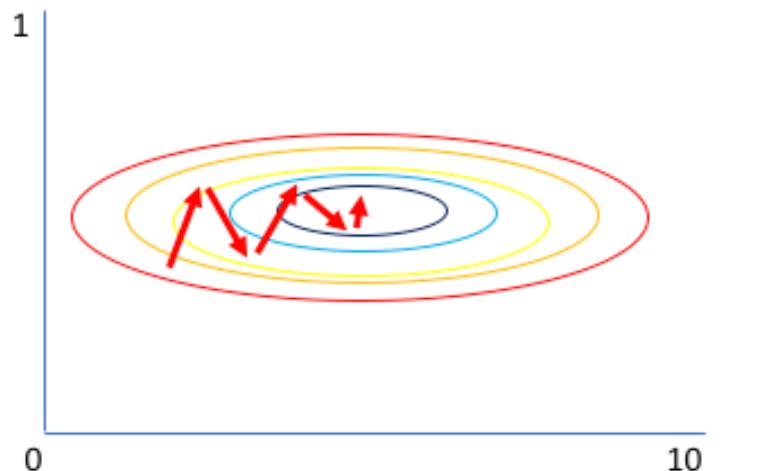
...nor does it find the optimal direction



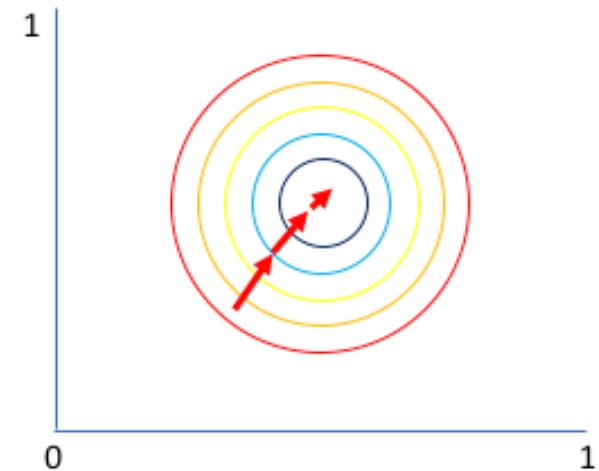
Steepest descent
better direction

Second order optimizers

Would find a better direction by taking into account the curvature (2nd derivative, Hessian)



Normalize by inverse Hessian



BUT: Are generally too expensive and don't work with minibatches

Stochastic gradient descent with momentum

Idea: keep a history of gradient steps to get smoother updates

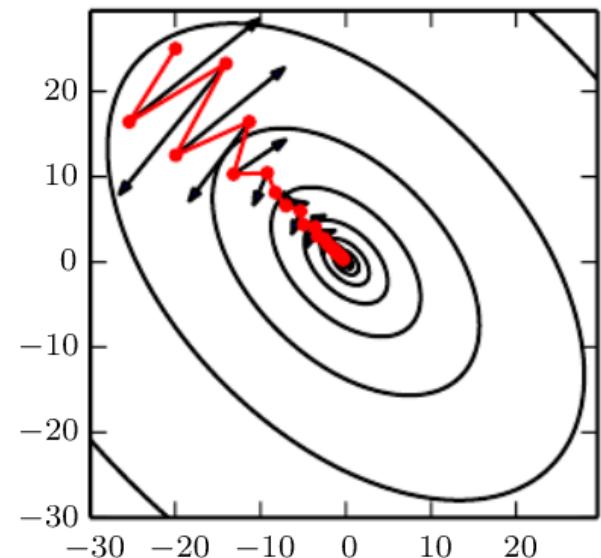
Update velocity v :

$$v \leftarrow \alpha v + \lambda \nabla_w E$$

$\alpha \in [0,1)$ controls how quickly the effect of past gradients decays

Update parameters:

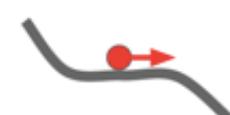
$$w \leftarrow w - v$$



Local Minima



Saddle points



AdaGrad

AdaGrad = “adaptive gradient algorithm”

Per-parameter learning rates according to:

$$w_i \leftarrow w_i - \frac{\eta}{\sqrt{G_i}} g_i^{(t)} \quad \text{where} \quad G_i = \sum_{\tau=1}^t (g_i^{(\tau)})^2$$

$\sqrt{G_i}$ is the L_2 norm of previous gradients

- Extreme parameter updates get damped
- Parameters with small updates get larger learning rates

Gradient of w_i
in iteration τ

RMSProp

RMSProp = “root mean square propagation”

Per-parameter learning rates according to:

$$w_i \leftarrow w_i - \frac{\eta}{\sqrt{v^{(t)}}} g_i^{(t)}$$

$$v^{(t)} = \gamma v^{(t-1)} + (1 - \gamma) (g_i^{(t)})^2$$

$v^{(t)}$ is an exponential moving average of the squared gradient

→ Similar effects to AdaGrad, but “forgets” earlier gradient magnitudes over time

Most popular optimizer to achieve results quickly and robustly

Adam

Adam = “Adaptive Moment Estimation”

Combination of Momentum and RMSProp optimizer

$$\begin{aligned} m_w^{(t)} &\leftarrow \beta_1 m_w^{(t-1)} + (1 - \beta_1) g^{(t)} \\ v_w^{(t)} &\leftarrow \beta_2 v_w^{(t-1)} + (1 - \beta_2) (g^{(t)})^2 \end{aligned}$$

Exponential moving average of mean and variance of gradient (biased estimate)

$$\hat{m}_w = \frac{m_w^{(t)}}{1 - \beta_1^t} \quad \hat{v}_w = \frac{v_w^{(t)}}{1 - \beta_2^t}$$

Bias correction for moments

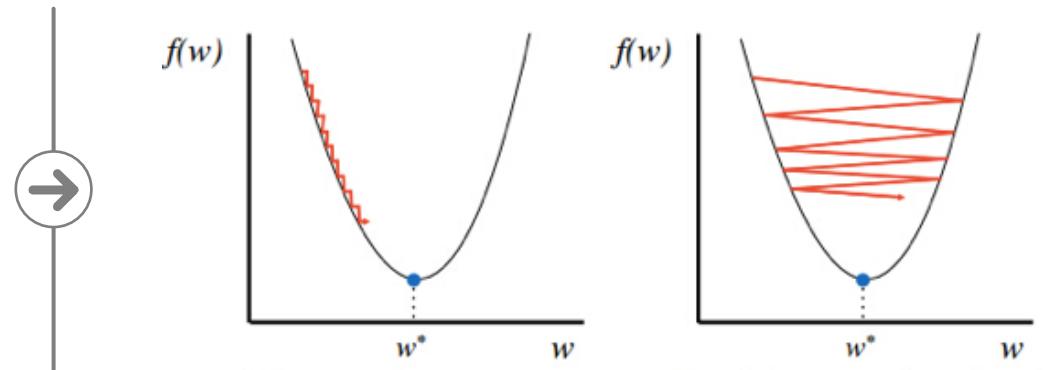
$$w^{(t)} \leftarrow w^{(t-1)} - \eta \frac{\hat{m}_w}{\sqrt{\hat{v}_w} + \epsilon}$$

Parameter update

Optimization results depend on hyperparameters

Choice of step size (learning rate) is important

- Too small: converge slowly or find local minimum
- Too large: oscillations or even divergence



Momentum affects convergence

- Too small: get stuck in local minima or saddle points
- Too large: overshoot optimum or spiral around it

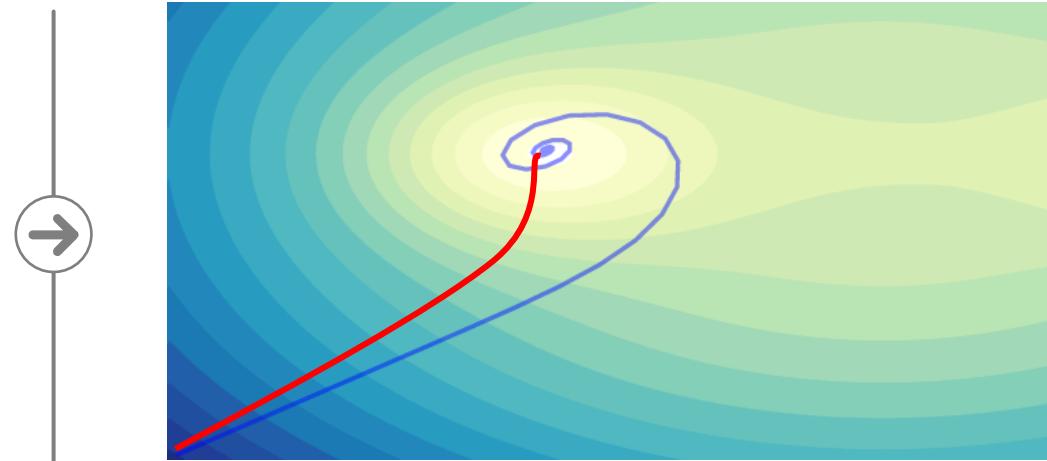
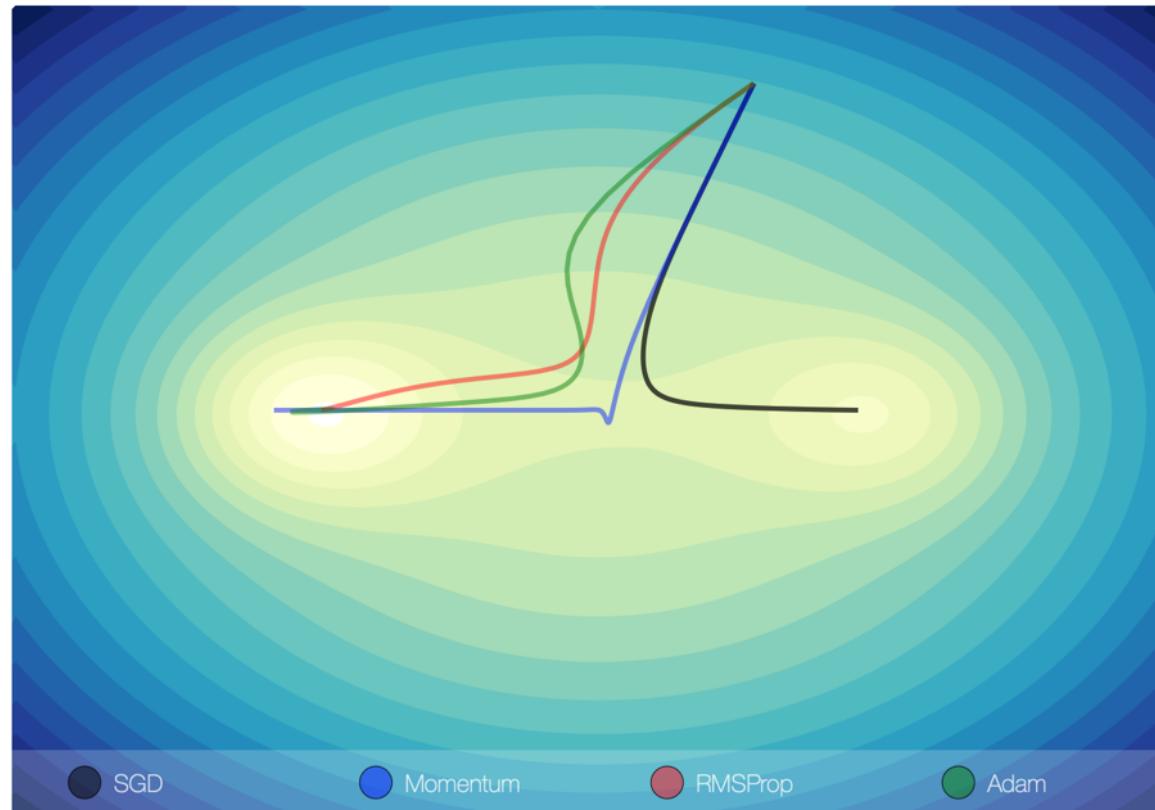
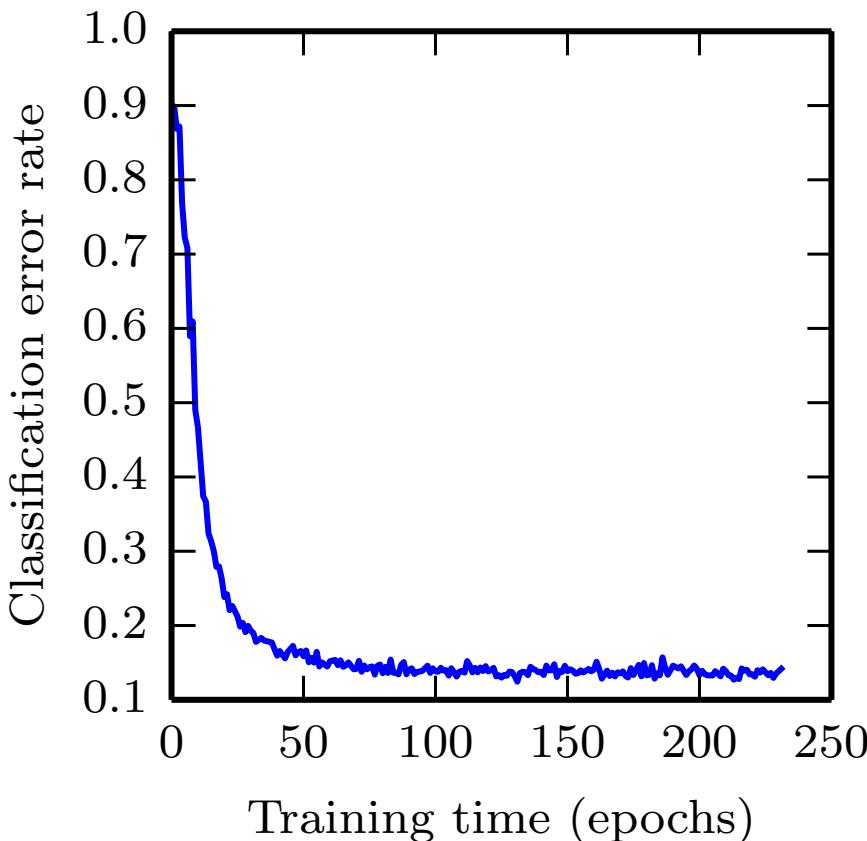
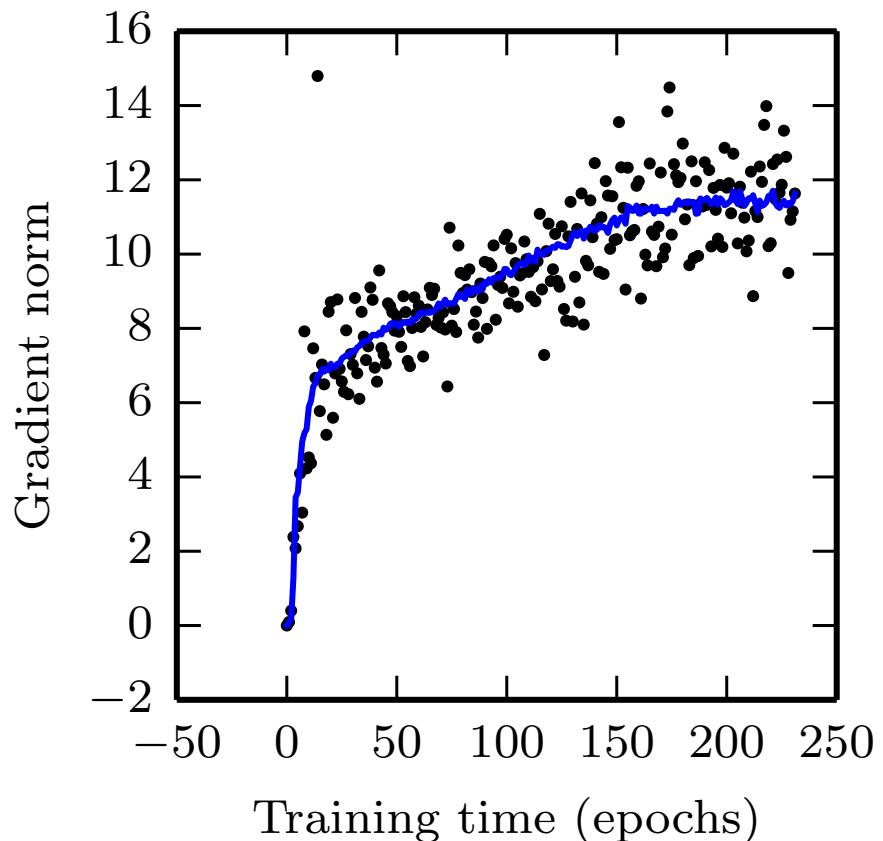


Illustration of trajectories for different optimizers



<https://emiliendupont.github.io/2018/01/24/optimization-visualization/>

We usually don't even reach a local minimum



Deep Learning optimization way of life

PURE MATH WAY OF LIFE

Find literally the smallest value of $f(x)$

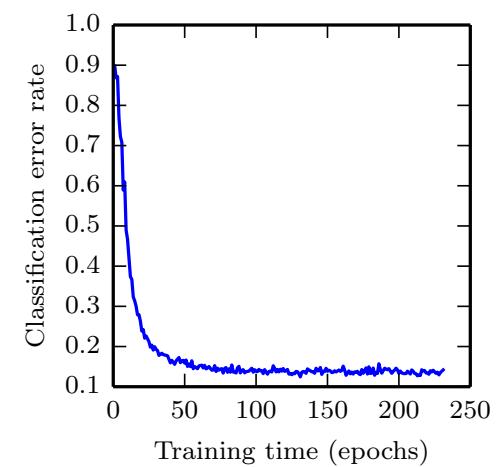
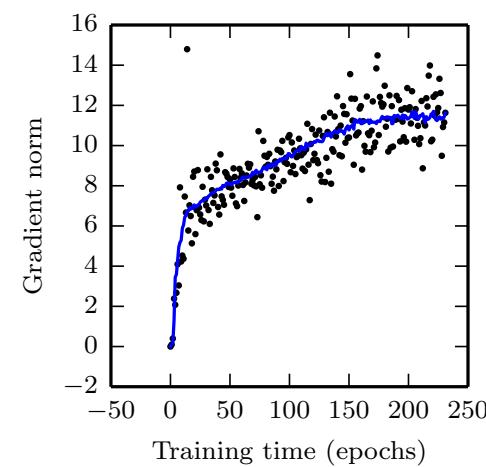
Alternatively: find some critical point of $f(x)$ where the value is locally smallest

(also true for convex methods like logistic regression/SVM)

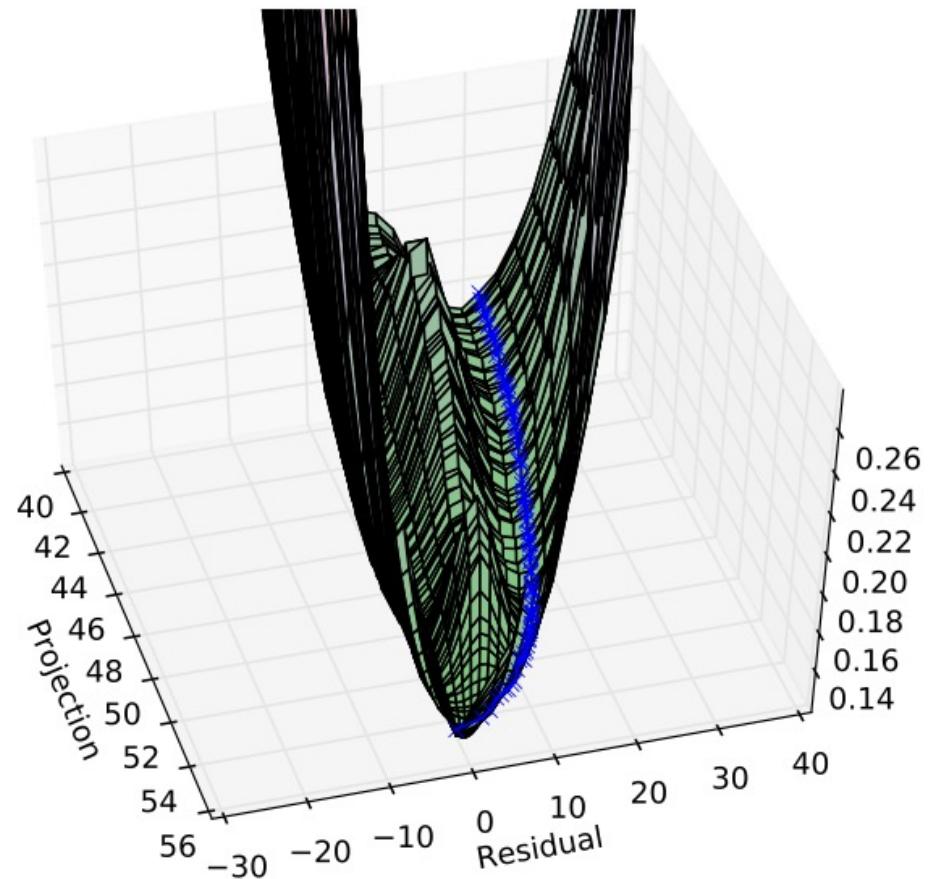


DEEP LEARNING WAY OF LIFE

Decrease the value of $f(x)$ a lot.



Neural network visualization



At end of learning:
- gradient is still large
- curvature is huge

(From “Qualitatively Characterizing Neural Network Optimization Problems”)

Regularization of MLPs

L_2 regularization (weight decay)

Assume Gaussian prior on weights θ

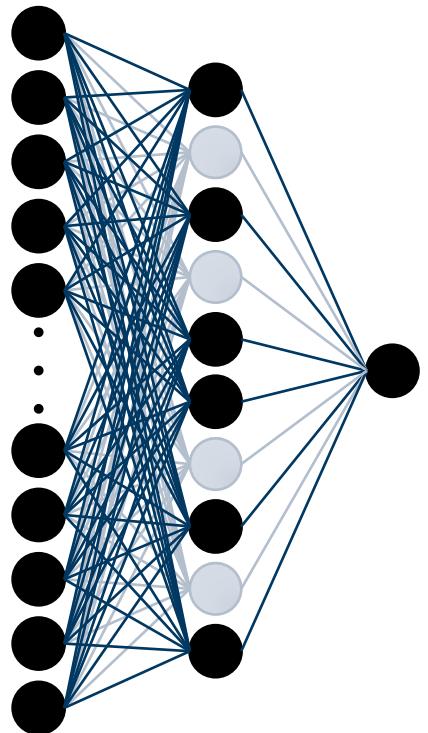
$$L = -\log P(\hat{y}|x, \theta) - \frac{\lambda}{2} \|\theta\|_2^2$$

$$\nabla_{\theta_i} = -\lambda \theta_i$$

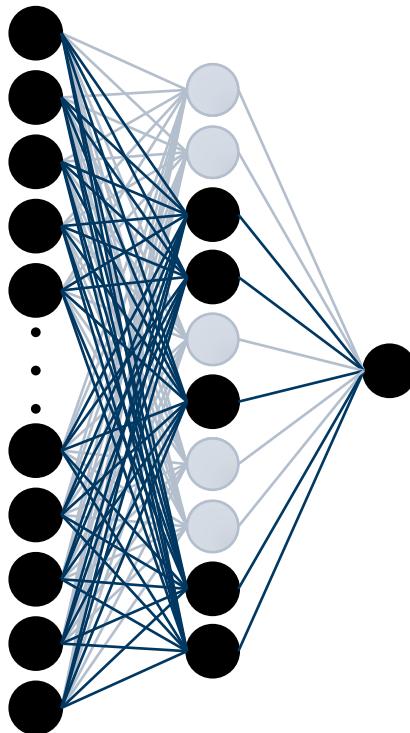
A.k.a **weight decay** because the weights decay towards zero in every iteration.

Dropout:

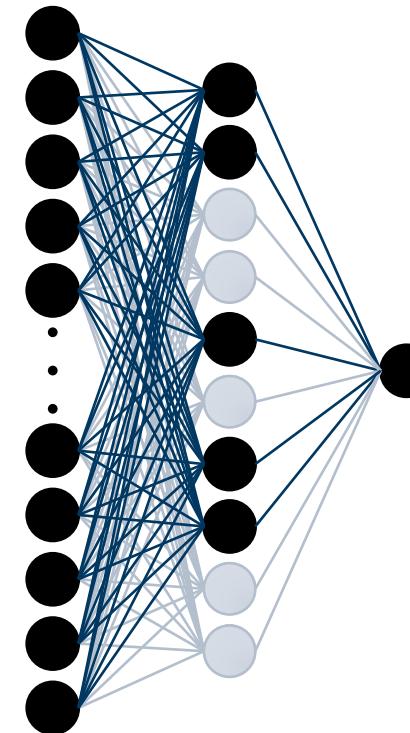
turn off hidden units randomly during training



Iteration 1



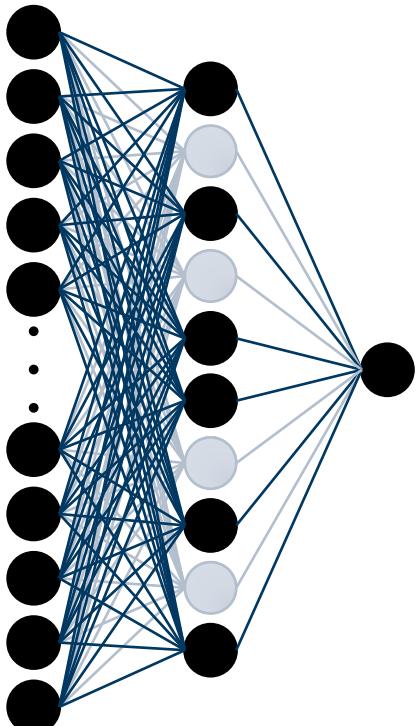
Iteration 2



Iteration 3

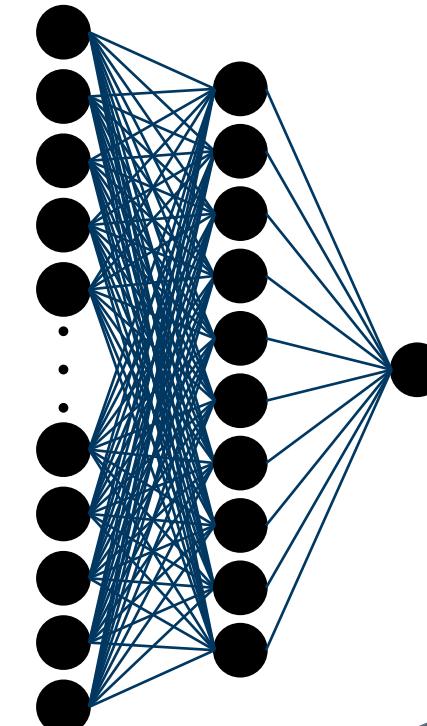
...

Dropout: training vs. inference



Training:
Train an ensemble
of models

$$h_i^{(n)} = \begin{cases} f(\mathbf{W}_{i:}^{(n)} \mathbf{h}^{(n-1)}) & \text{with prob. } p \\ 0 & \text{otherwise} \end{cases}$$

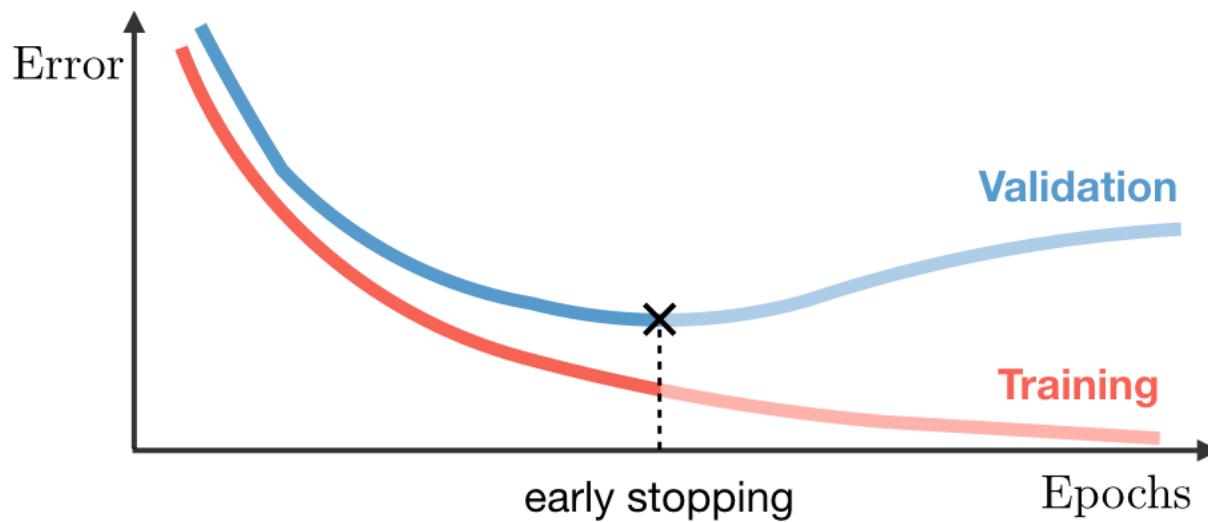


Inference:
Perform model
averaging

Adjust gain to keep
average activation
level constant

$$h_i^{(n)} = p \cdot f(\mathbf{W}_{i:}^{(n)} \mathbf{h}^{(n-1)})$$

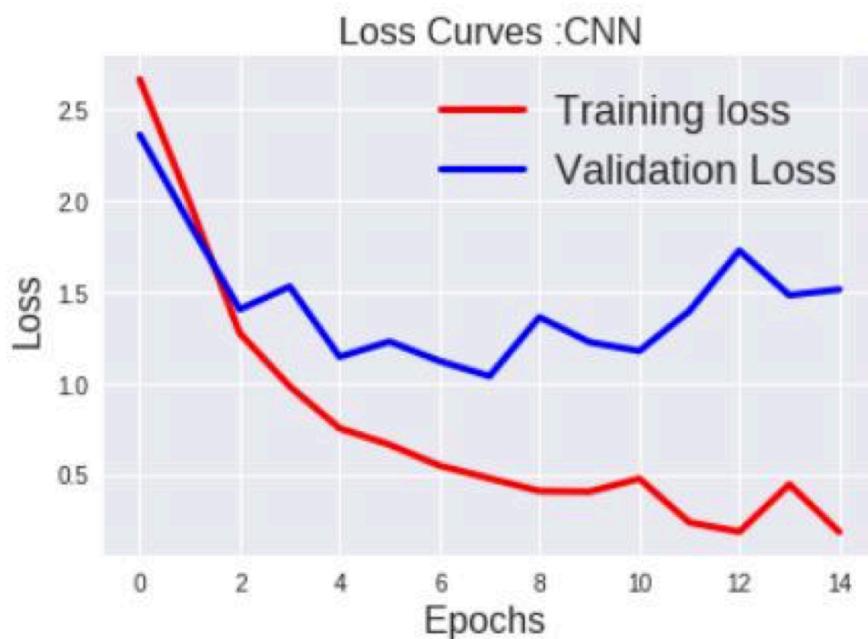
Early stopping



Idea: Stop training when generalization error increases

Early stopping: in practice

PROBLEM: LOSS FLUCTUATES



POSSIBLE SOLUTIONS

Smooth validation loss by computing running average (e.g. exponential moving average)



Introduce “patience” parameter τ : stop when loss has not decreased for τ steps/epochs in a row

Questions!
