

Deep Feedforward Networks

Lecture slides for Chapter 6 of *Deep Learning*

www.deeplearningbook.org

Ian Goodfellow

Last updated 2016-10-04

Roadmap

- Example: Learning XOR
- Gradient-Based Learning
- Hidden Units
- Architecture Design
- Back-Propagation

XOR is not linearly separable

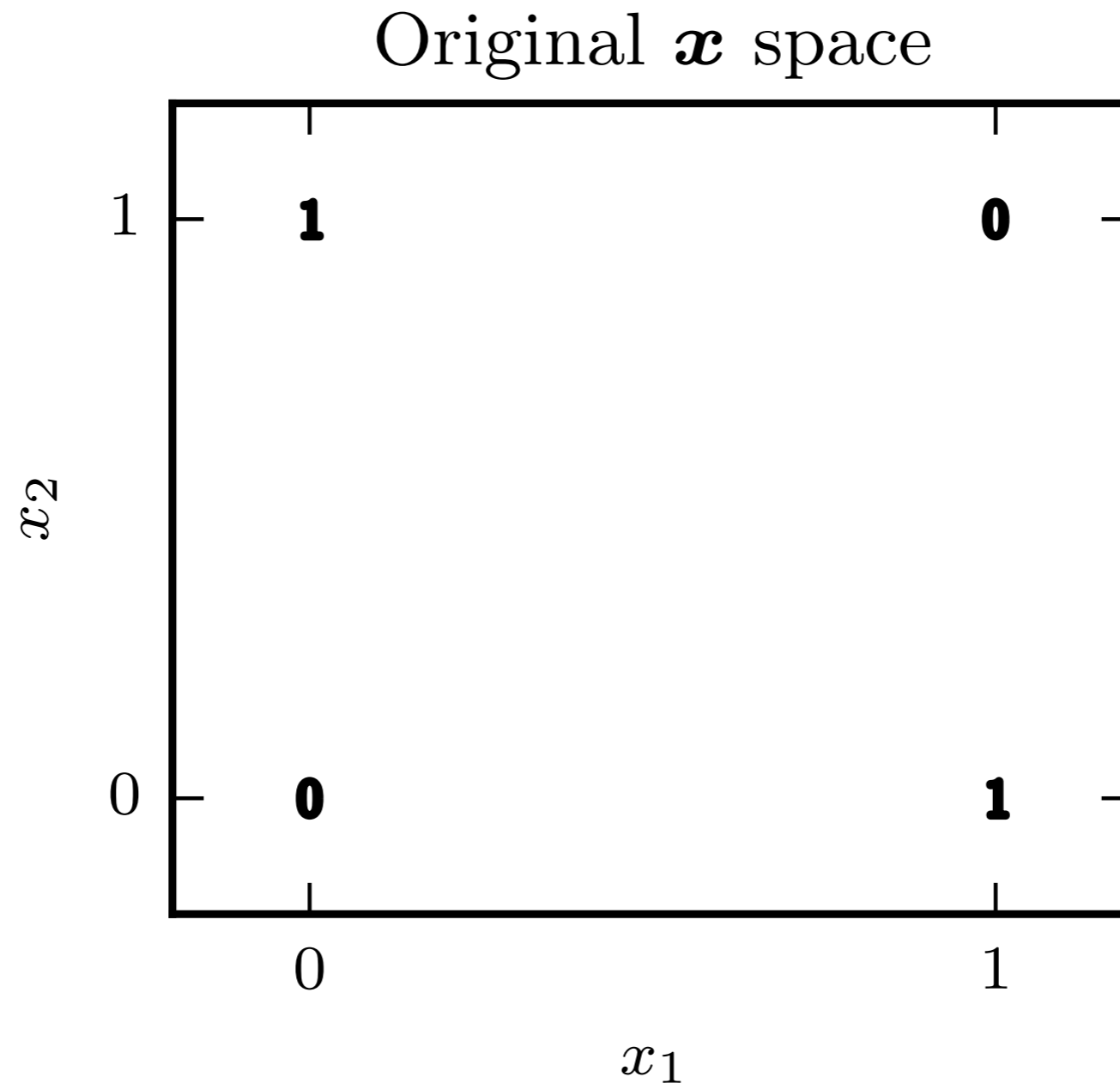


Figure 6.1, left

Rectified Linear Activation

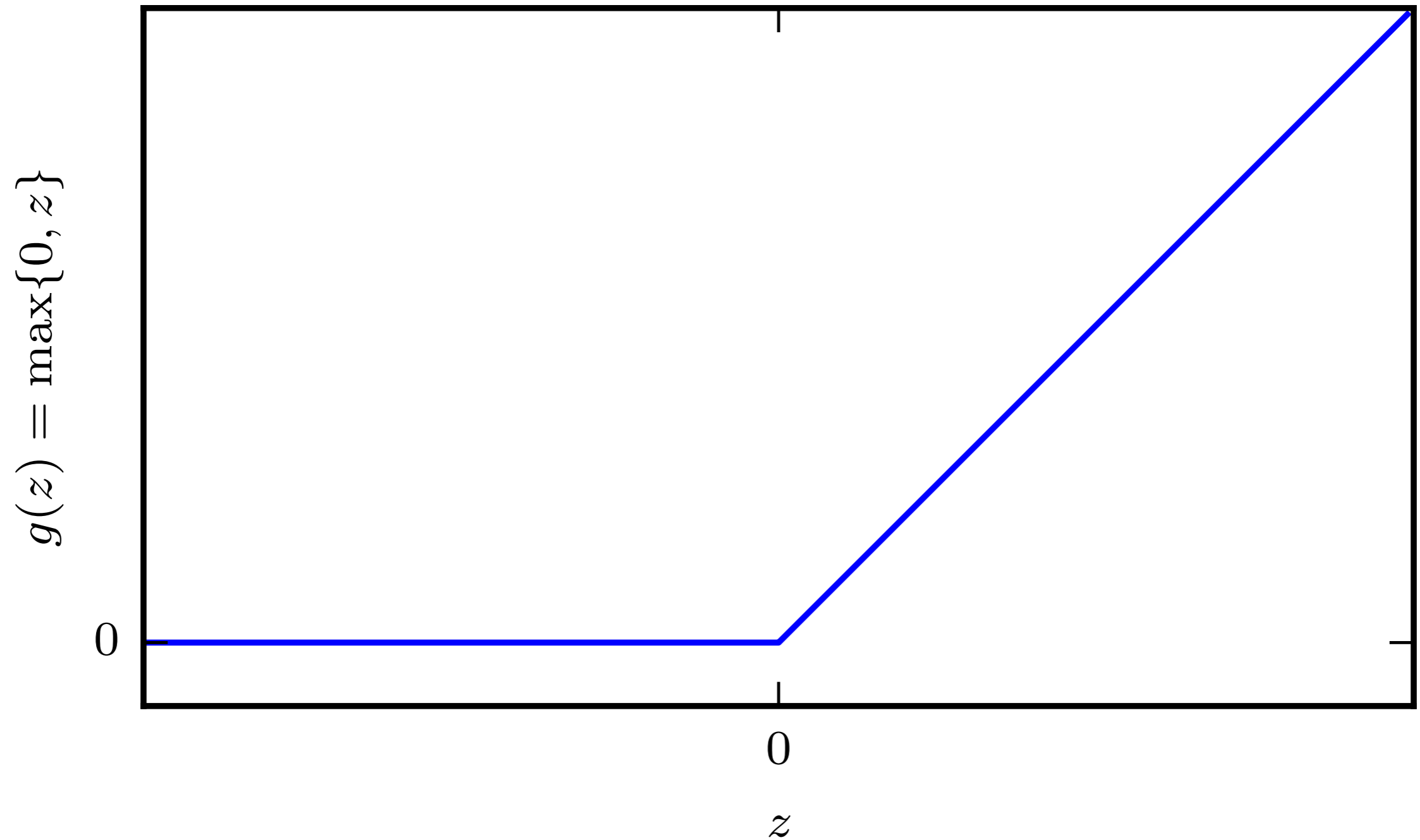


Figure 6.3

Network Diagrams

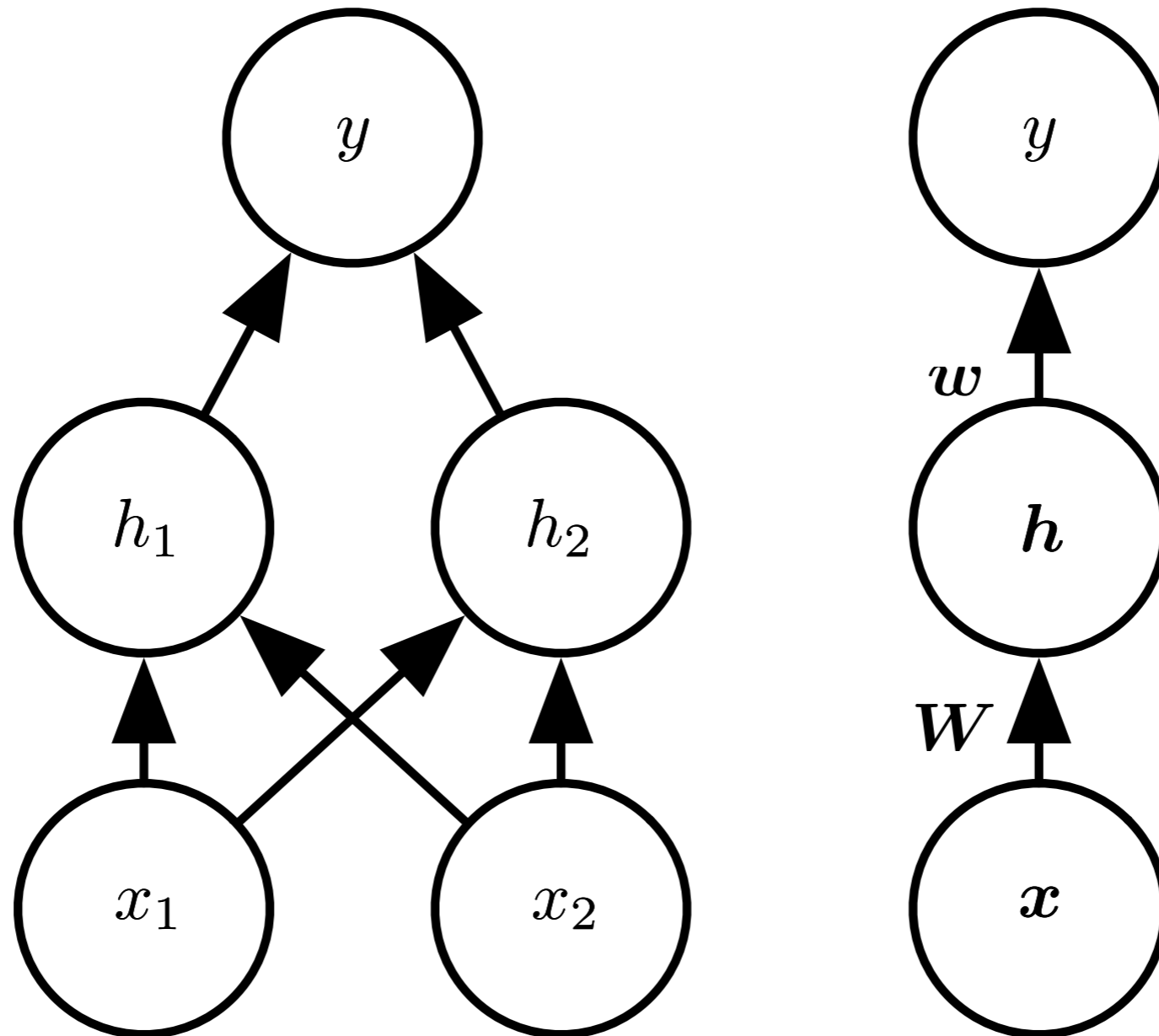


Figure 6.2

Solving XOR

$$f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^\top \max\{0, \mathbf{W}^\top \mathbf{x} + \mathbf{c}\} + b. \quad (6.3)$$

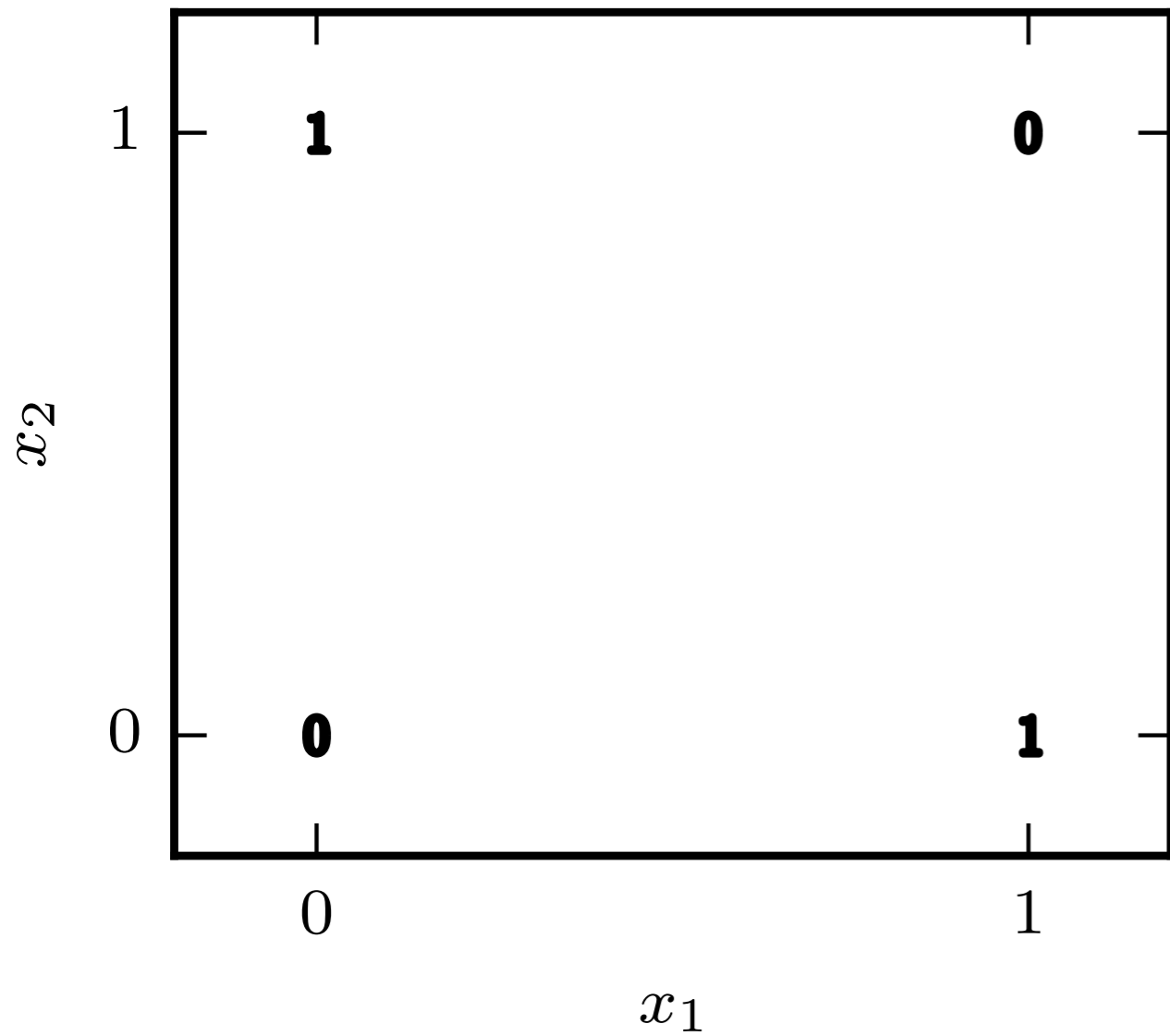
$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad (6.4)$$

$$\mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \quad (6.5)$$

$$\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, \quad (6.6)$$

Solving XOR

Original \mathbf{x} space



Learned \mathbf{h} space

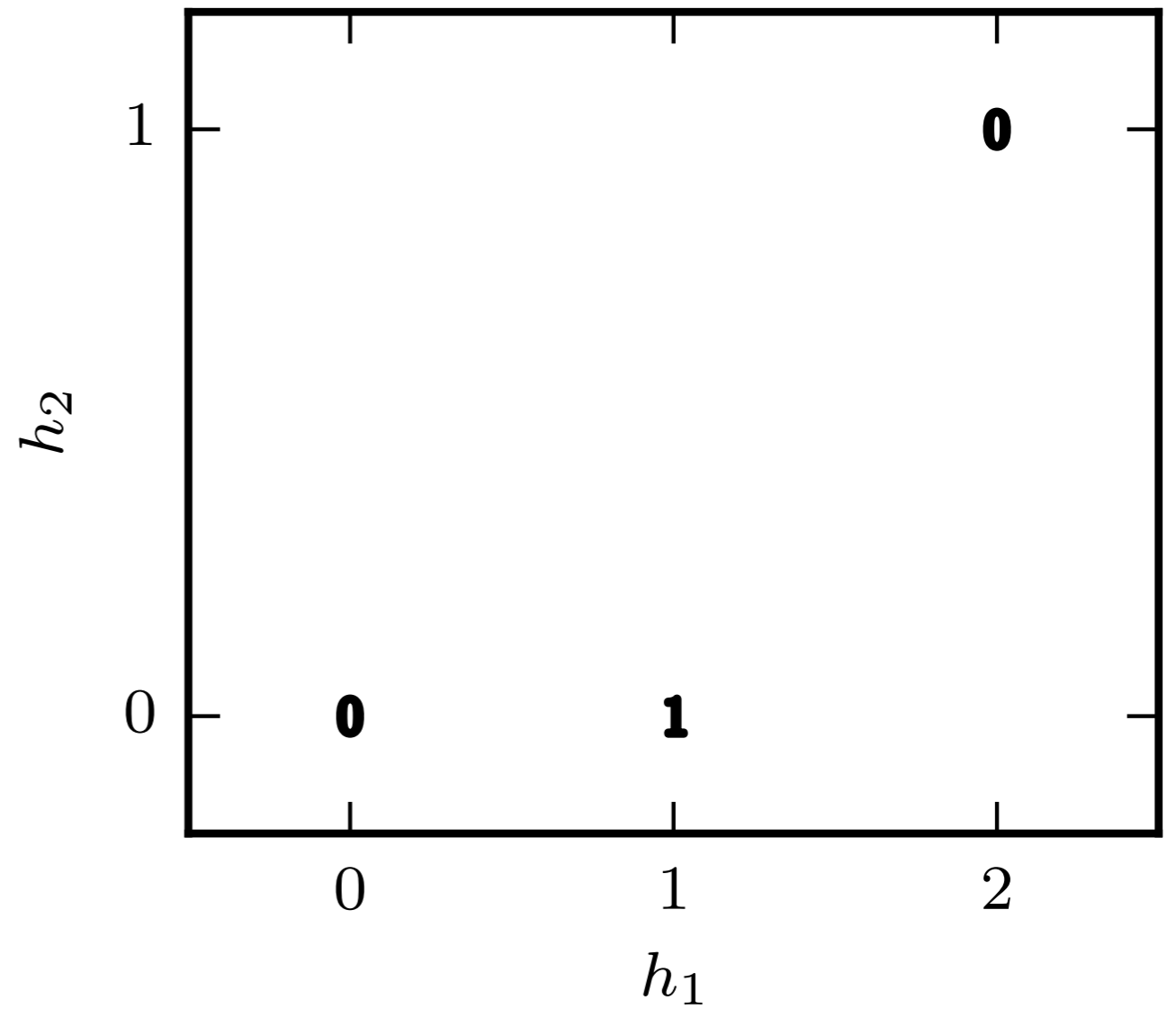


Figure 6.1

Roadmap

- Example: Learning XOR
- Gradient-Based Learning
- Hidden Units
- Architecture Design
- Back-Propagation

Gradient-Based Learning

- Specify
 - Model
 - Cost
- Design model and cost so cost is smooth
- Minimize cost using gradient descent or related techniques

Conditional Distributions and Cross-Entropy

$$J(\boldsymbol{\theta}) = -\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{y} \mid \mathbf{x}). \quad (6.12)$$

Output Types

Output Type	Output Distribution	Output Layer	Cost Function
Binary	Bernoulli	Sigmoid	Binary cross-entropy
Discrete	Multinoulli	Softmax	Discrete cross-entropy
Continuous	Gaussian	Linear	Gaussian cross-entropy (MSE)
Continuous	Mixture of Gaussian	Mixture Density	Cross-entropy
Continuous	Arbitrary	See part III: GAN, VAE, FVBN	Various

Mixture Density Outputs

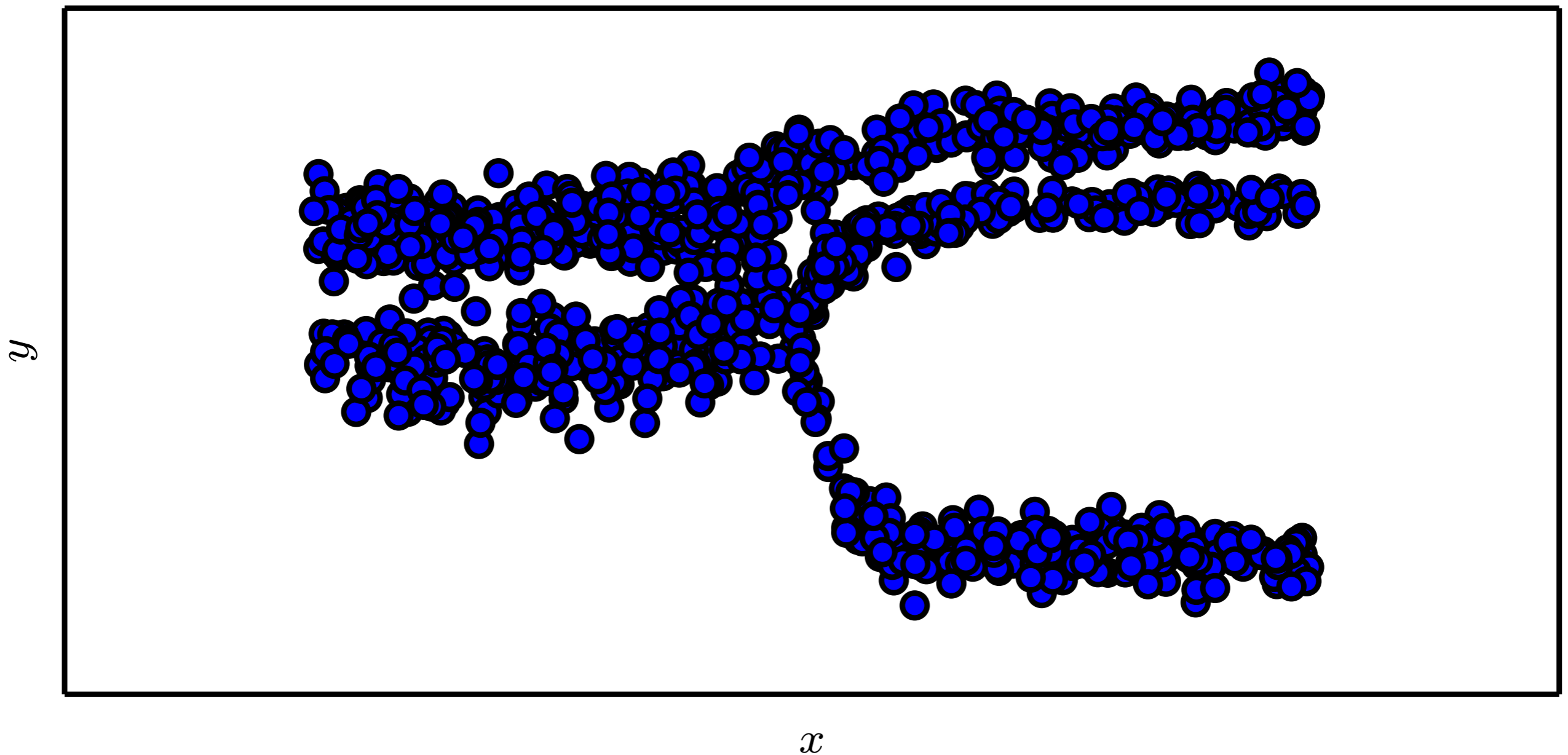
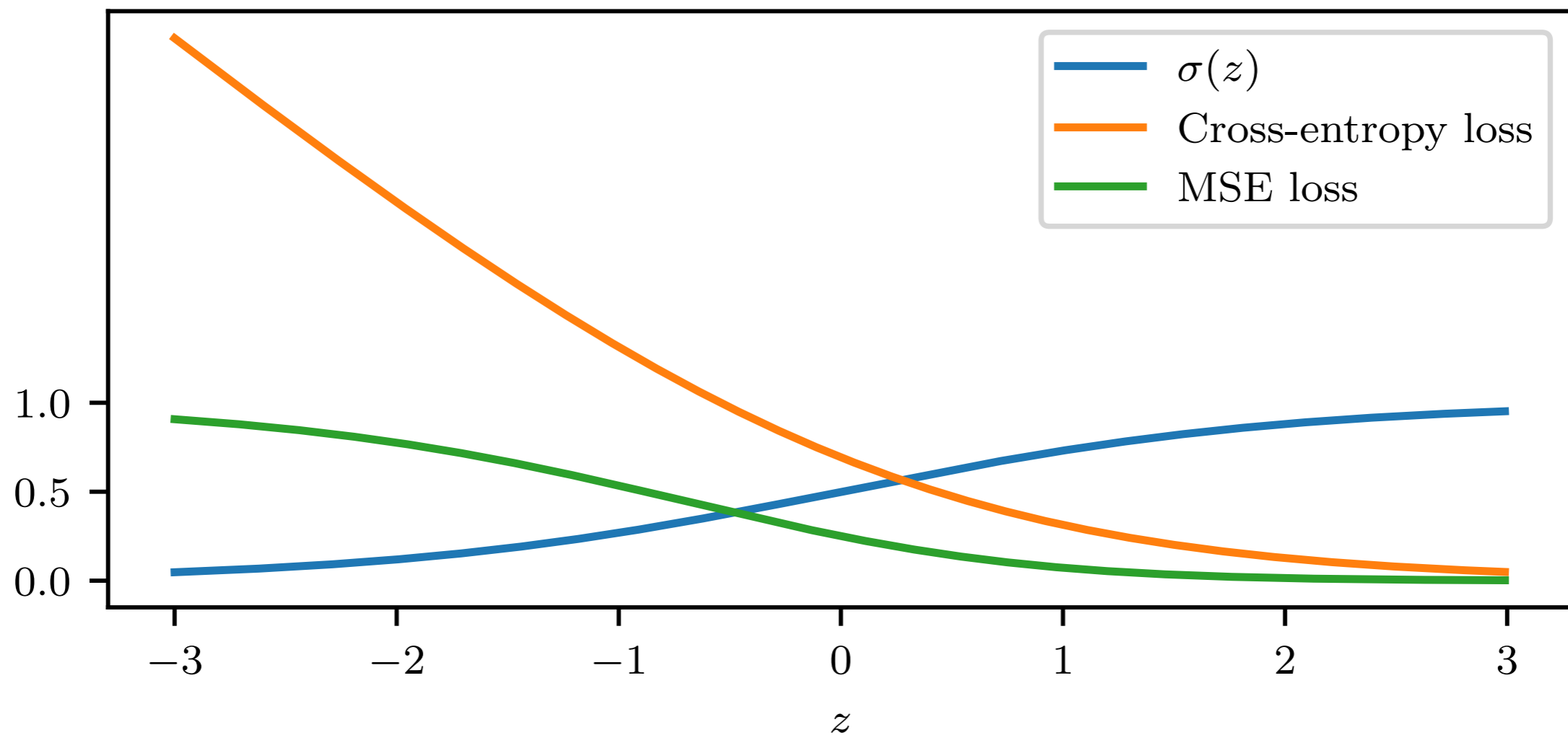


Figure 6.4

Don't mix and match

Sigmoid output with target of 1



Roadmap

- Example: Learning XOR
- Gradient-Based Learning
- Hidden Units
- Architecture Design
- Back-Propagation

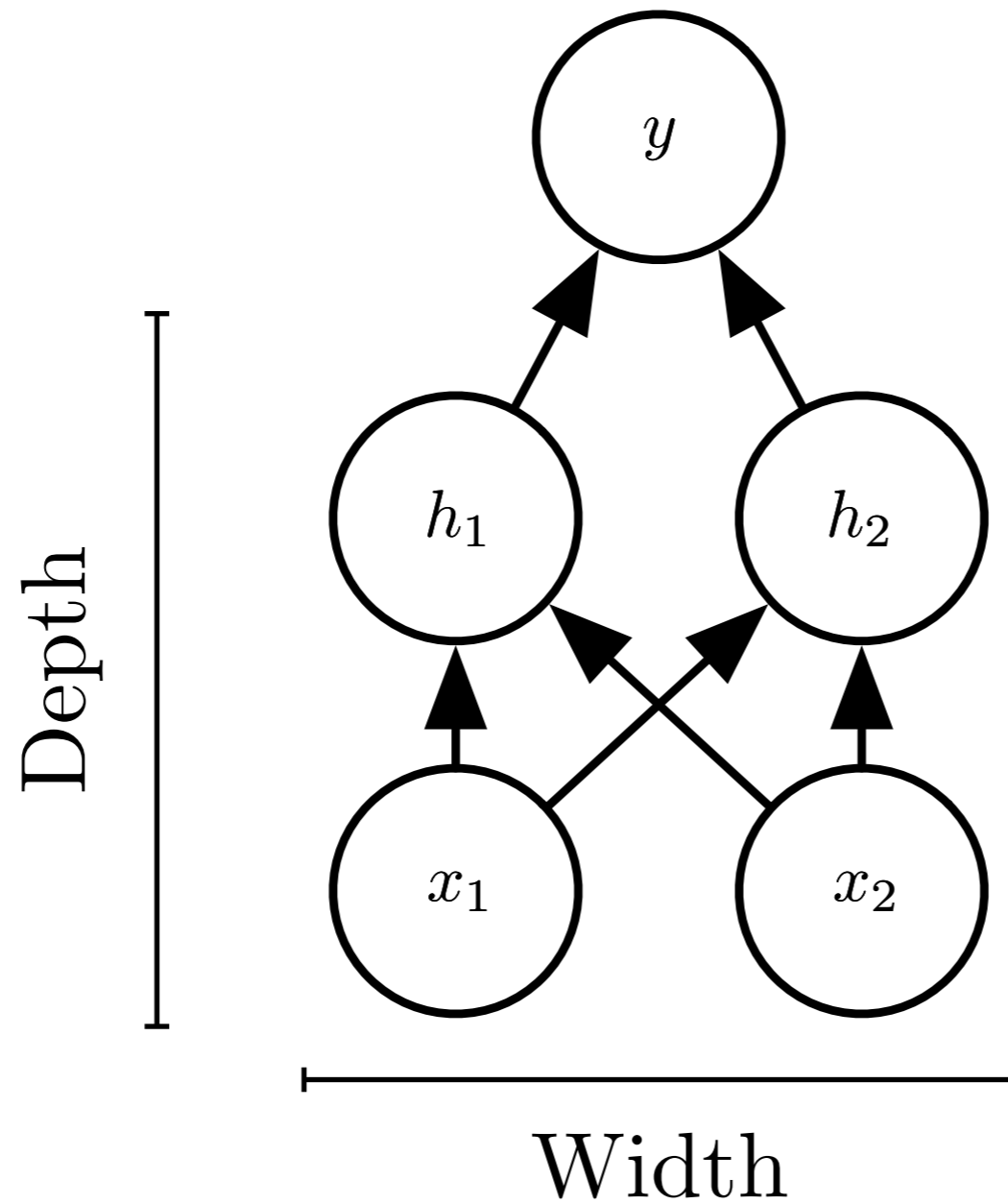
Hidden units

- Use ReLUs, 90% of the time
- For RNNs, see Chapter 10
- For some research projects, get creative
- Many hidden units perform comparably to ReLUs. New hidden units that perform comparably are rarely interesting.

Roadmap

- Example: Learning XOR
- Gradient-Based Learning
- Hidden Units
- Architecture Design
- Back-Propagation

Architecture Basics



Universal Approximator Theorem

- One hidden layer is enough to *represent* (not *learn*) an approximation of any function to an arbitrary degree of accuracy
- So why deeper?
 - Shallow net may need (exponentially) more width
 - Shallow net may overfit more

Exponential Representation

Advantage of Depth

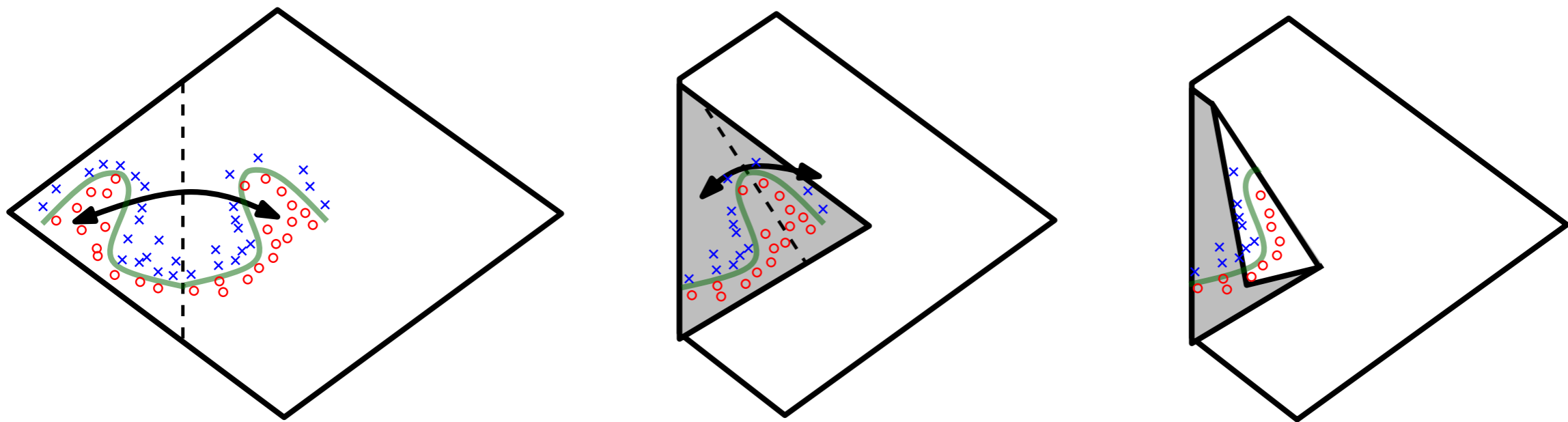


Figure 6.5

Better Generalization with Greater Depth

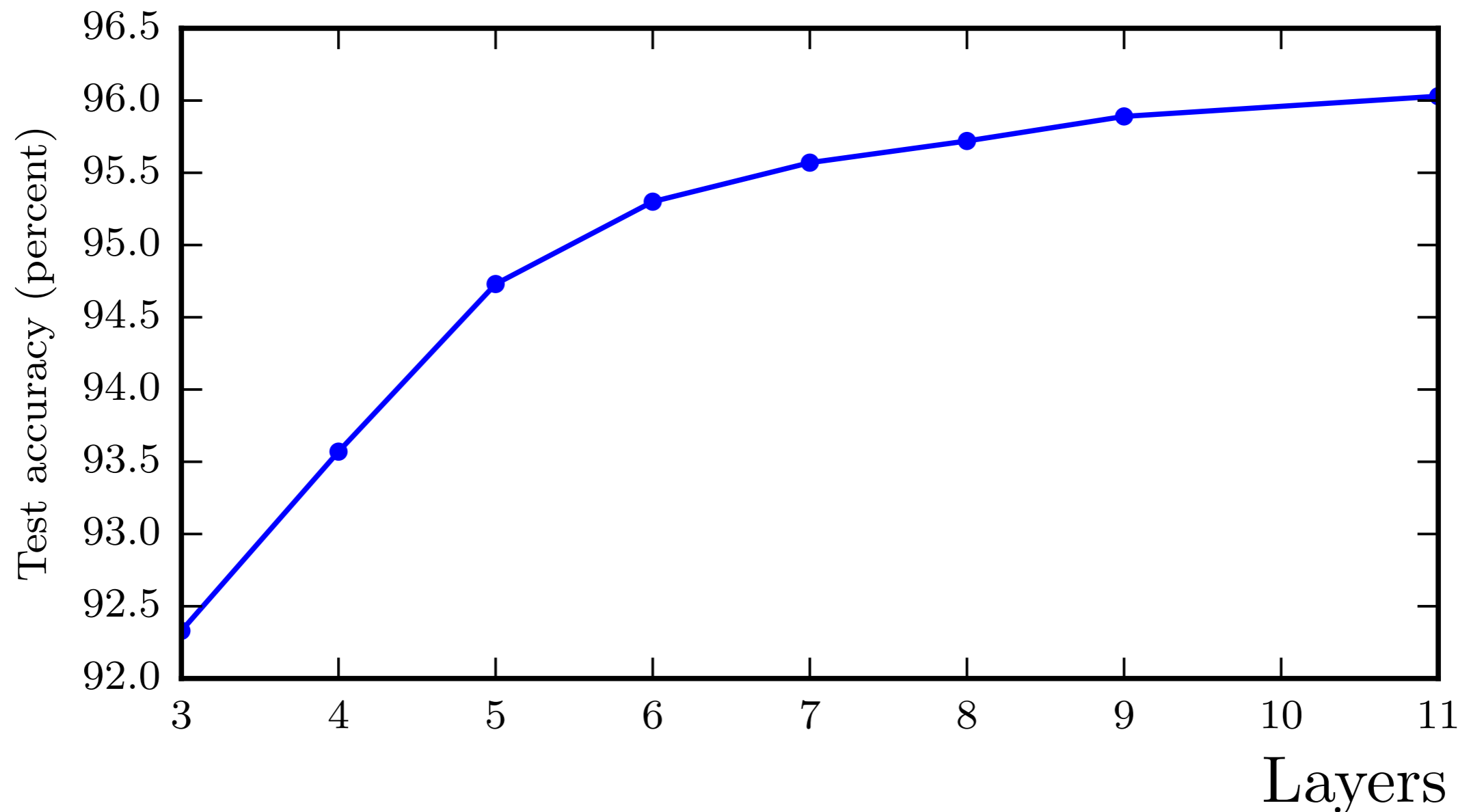


Figure 6.6

Large, Shallow Models Overfit More

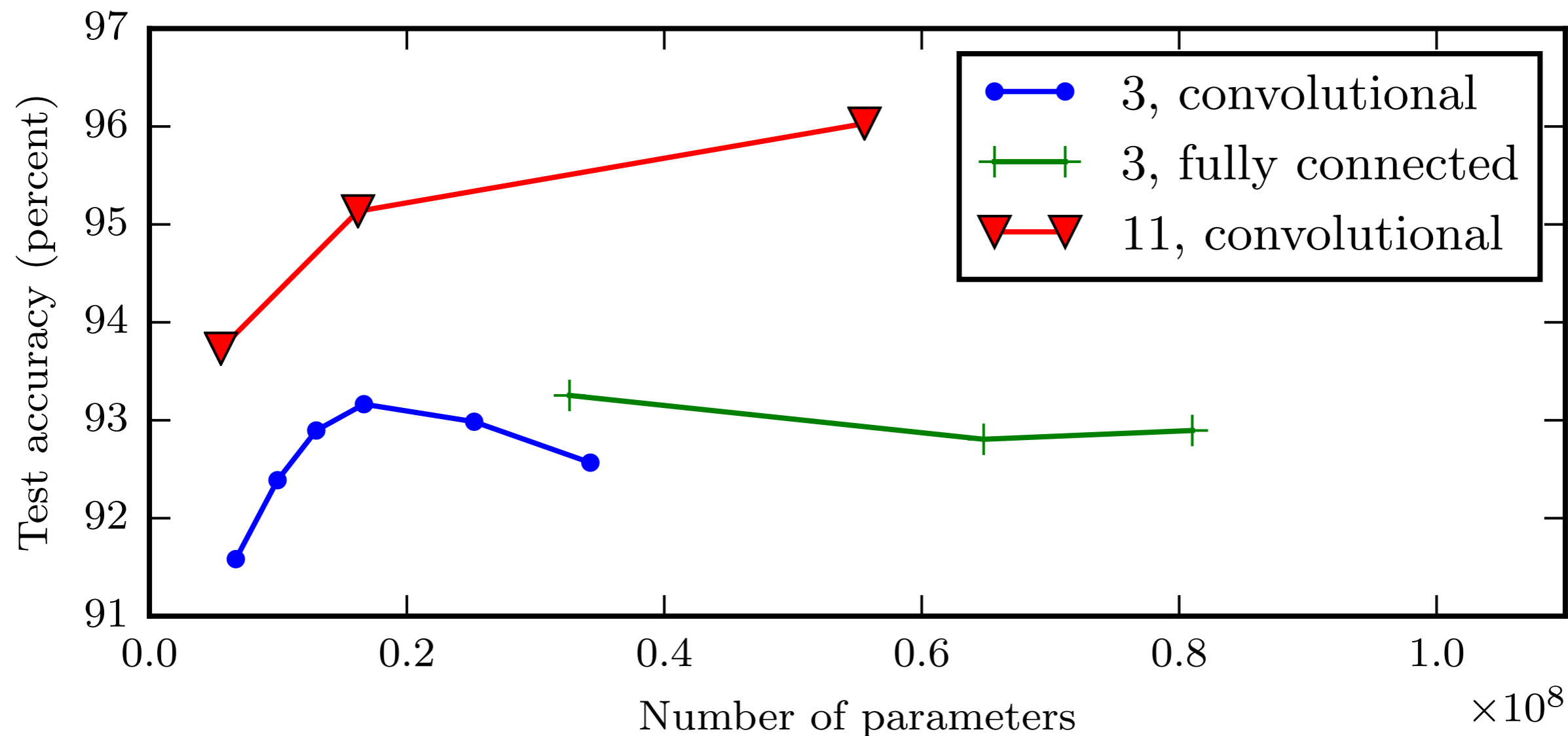


Figure 6.7

Roadmap

- Example: Learning XOR
- Gradient-Based Learning
- Hidden Units
- Architecture Design
- Back-Propagation

Back-Propagation

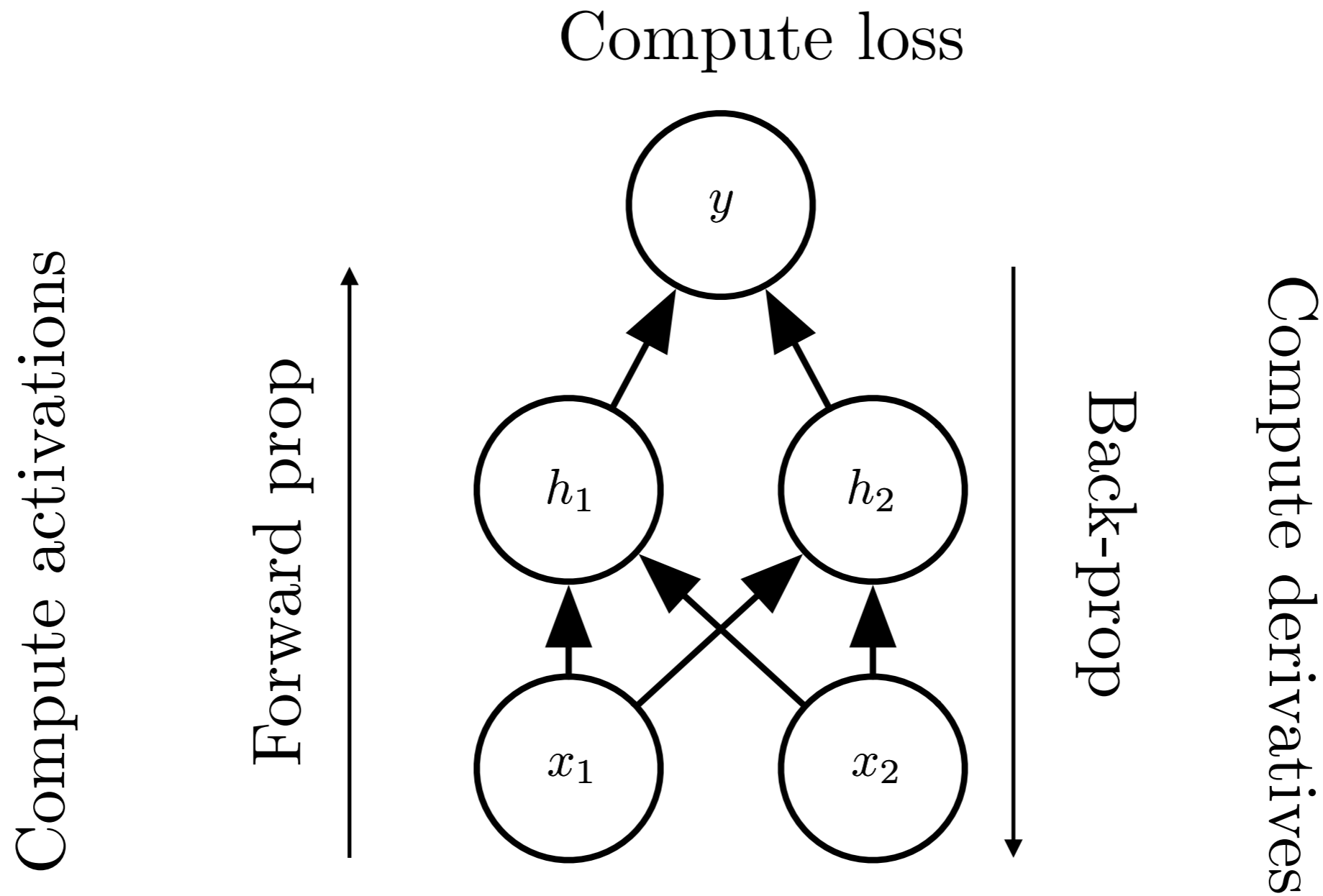
- Back-propagation is “just the chain rule” of calculus

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}. \quad (6.44)$$

$$\nabla_{\mathbf{x}} z = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^\top \nabla_{\mathbf{y}} z, \quad (6.46)$$

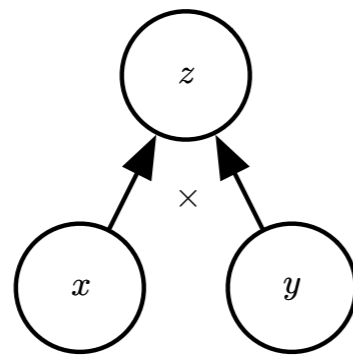
- But it’s a particular implementation of the chain rule
 - Uses dynamic programming (table filling)
 - Avoids recomputing repeated subexpressions
 - Speed vs memory tradeoff

Simple Back-Prop Example



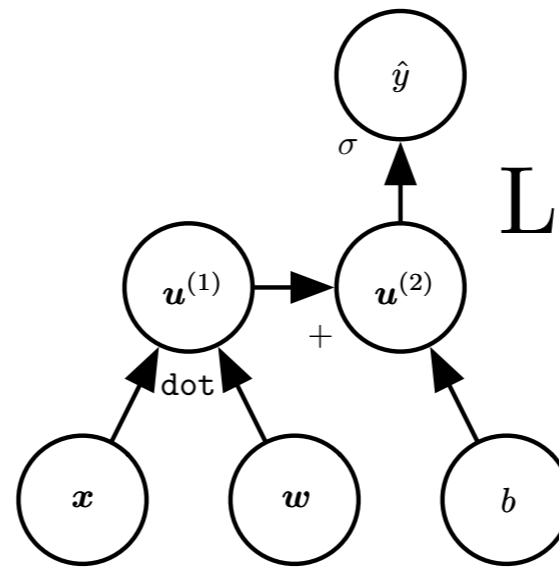
Computation Graphs

Multiplication



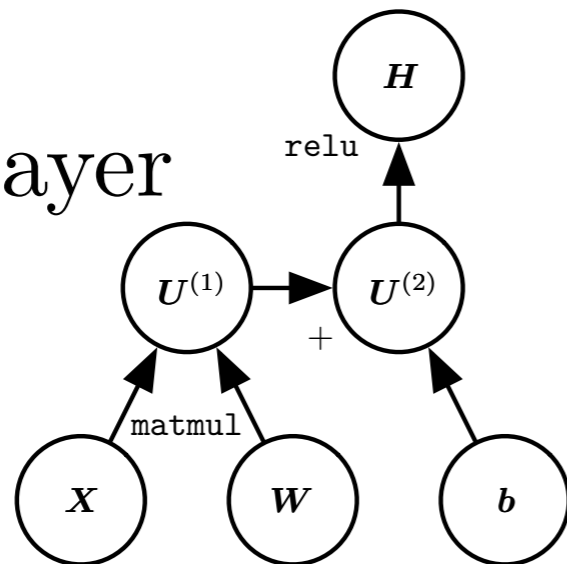
(a)

Logistic regression



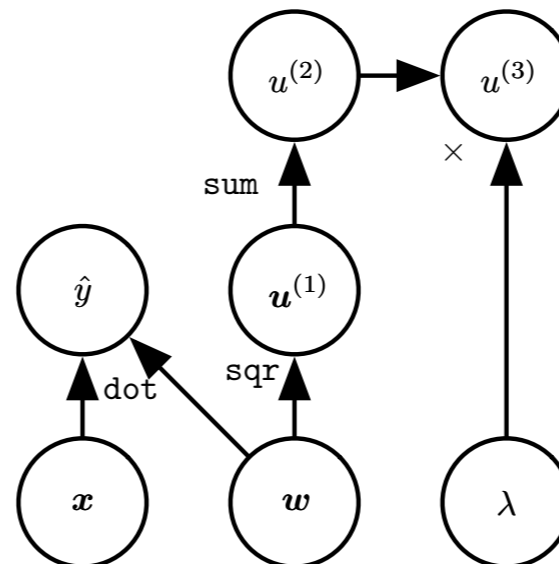
(b)

ReLU layer



(c)

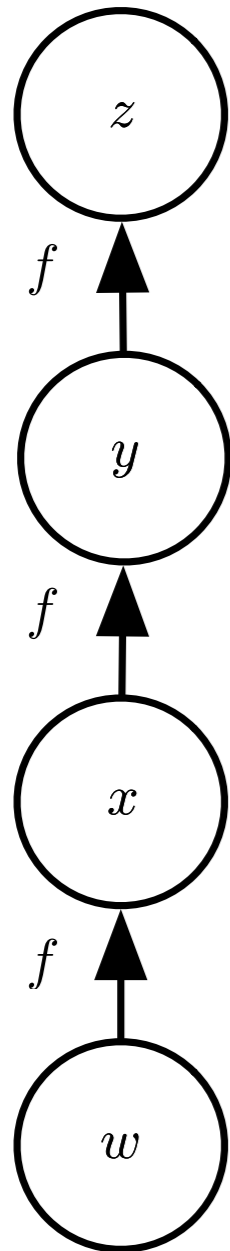
Linear regression
and weight decay



(d)

Figure 6.8

Repeated Subexpressions



$$\frac{\partial z}{\partial w} \tag{6.50}$$

$$= \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w} \tag{6.51}$$

$$= f'(y) f'(x) f'(w) \tag{6.52}$$

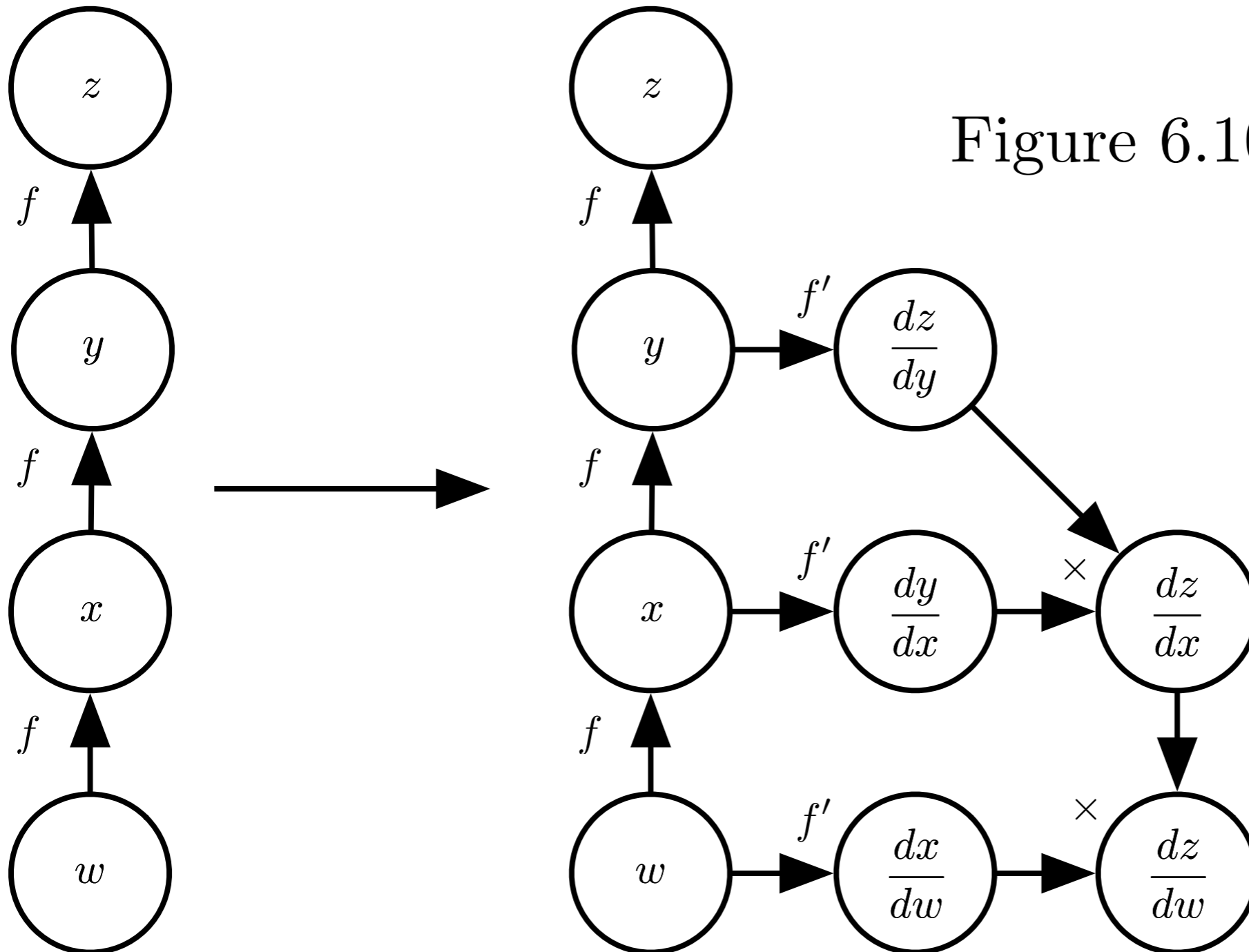
$$= f'(f(f(w))) f'(f(w)) f'(w) \tag{6.53}$$

Back-prop avoids computing this twice

Figure 6.9

Symbol-to-Symbol Differentiation

Figure 6.10



Neural Network Loss Function

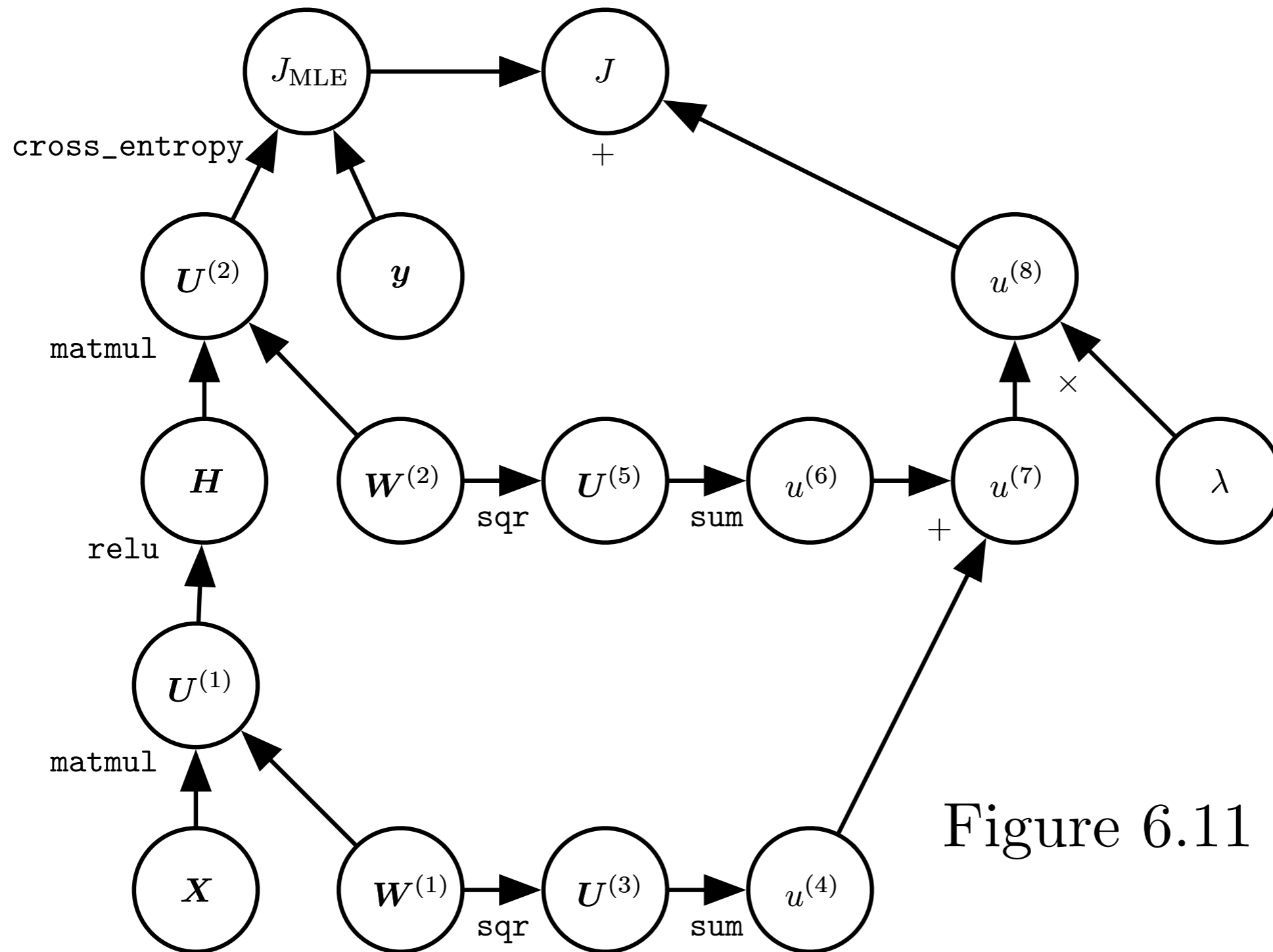


Figure 6.11

Hessian-vector Products

$$\mathbf{H}\mathbf{v} = \nabla_{\mathbf{x}} \left[(\nabla_{\mathbf{x}} f(\mathbf{x}))^\top \mathbf{v} \right]. \quad (6.59)$$

Questions