

Deep Feedforward Networks

Lecture slides for Chapter 6 of *Deep Learning*

www.deeplearningbook.org

Ian Goodfellow

Last updated 2016-10-04

Adapted by m.n. for CMPS 392

Multilayer perceptrons (MLP)

- Approximate some function f^*
- A feedforward network defines a mapping
$$y = f(x; \theta)$$
- No feedback connections
- Functions are composed in a chain $f(x) = f^{(3)}(f^{(2)}(f^{(1)}(x)))$
 - $f^{(1)}$ is the first layer, $f^{(2)}$ is the second layer, and so on ...
- The training examples specify what the output layer must do
 - The other layers are called hidden layers

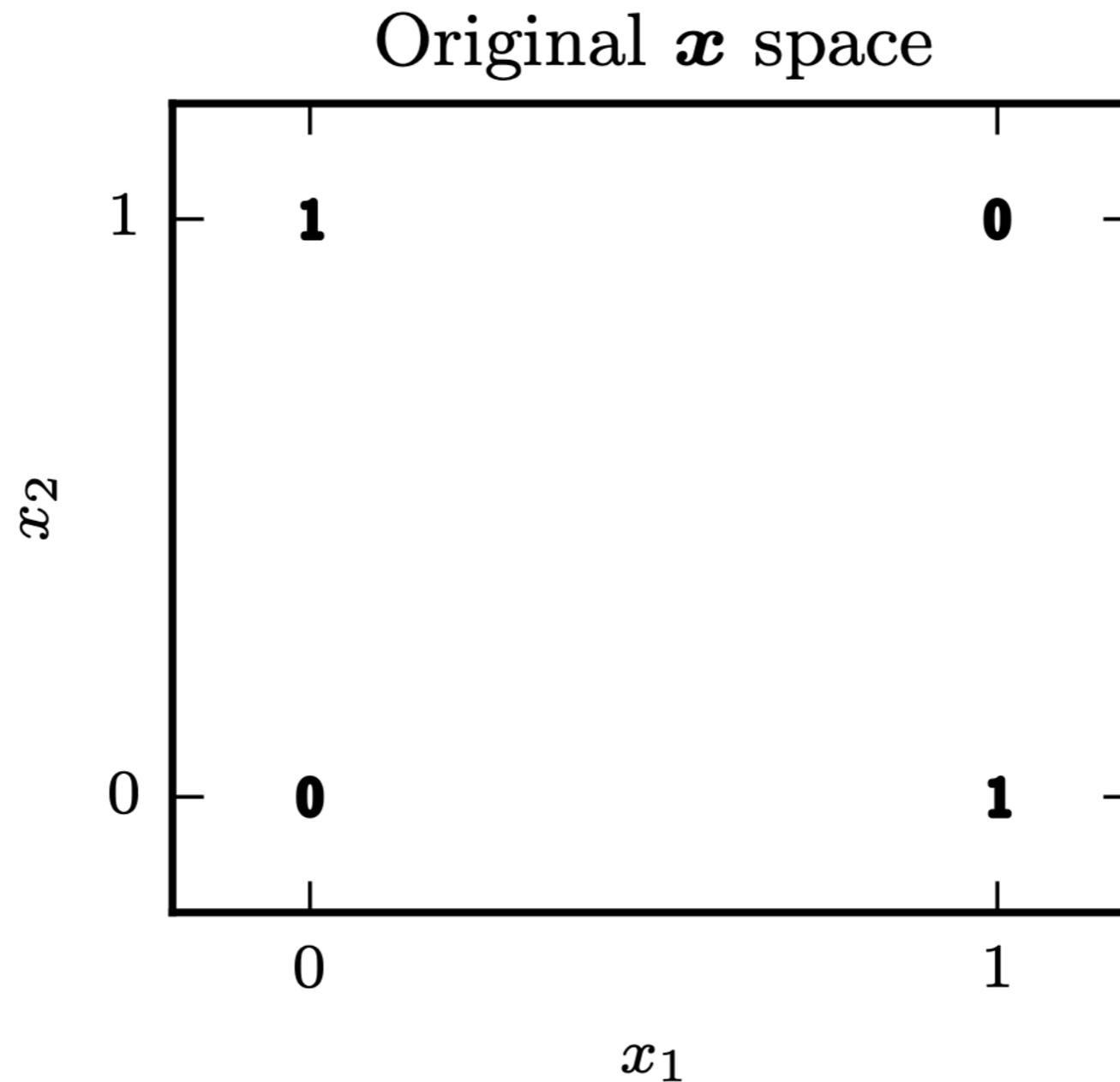
What's the idea?

- The strategy of deep learning is to learn a new representation of the data $\phi(\mathbf{x}, \boldsymbol{\theta})$. Think of this as the output of a hidden layer.
- Parameters w maps the output of hidden layers to the desired output: $\phi(\mathbf{x}, \boldsymbol{\theta})^T \mathbf{w}$
- We give up on the convexity of the training problem
- $\phi(\mathbf{x}, \boldsymbol{\theta})$ can be very generic, in the same time its design can be guided by human experts

Roadmap

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

XOR is not linearly separable



Linear model?

- Cost function:

$$J(\theta) = \frac{1}{4} \sum_{\mathbf{x} \in \mathbb{X}} (f^*(\mathbf{x}) - f(\mathbf{x}; \theta))^2$$

- Linear model

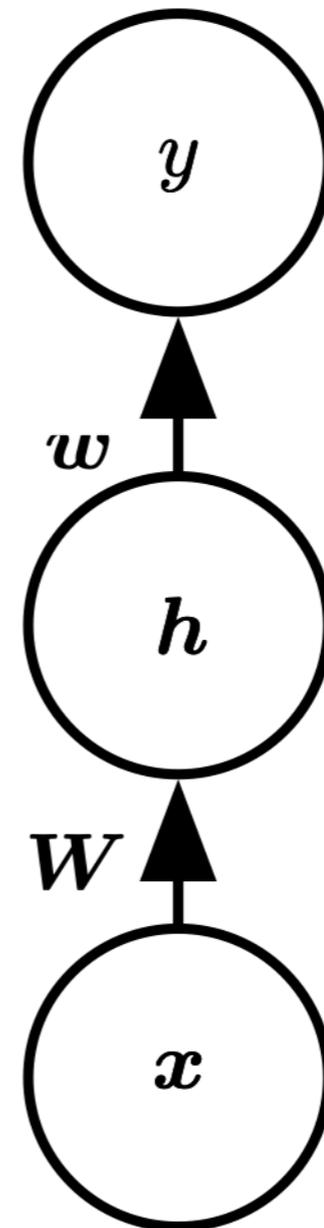
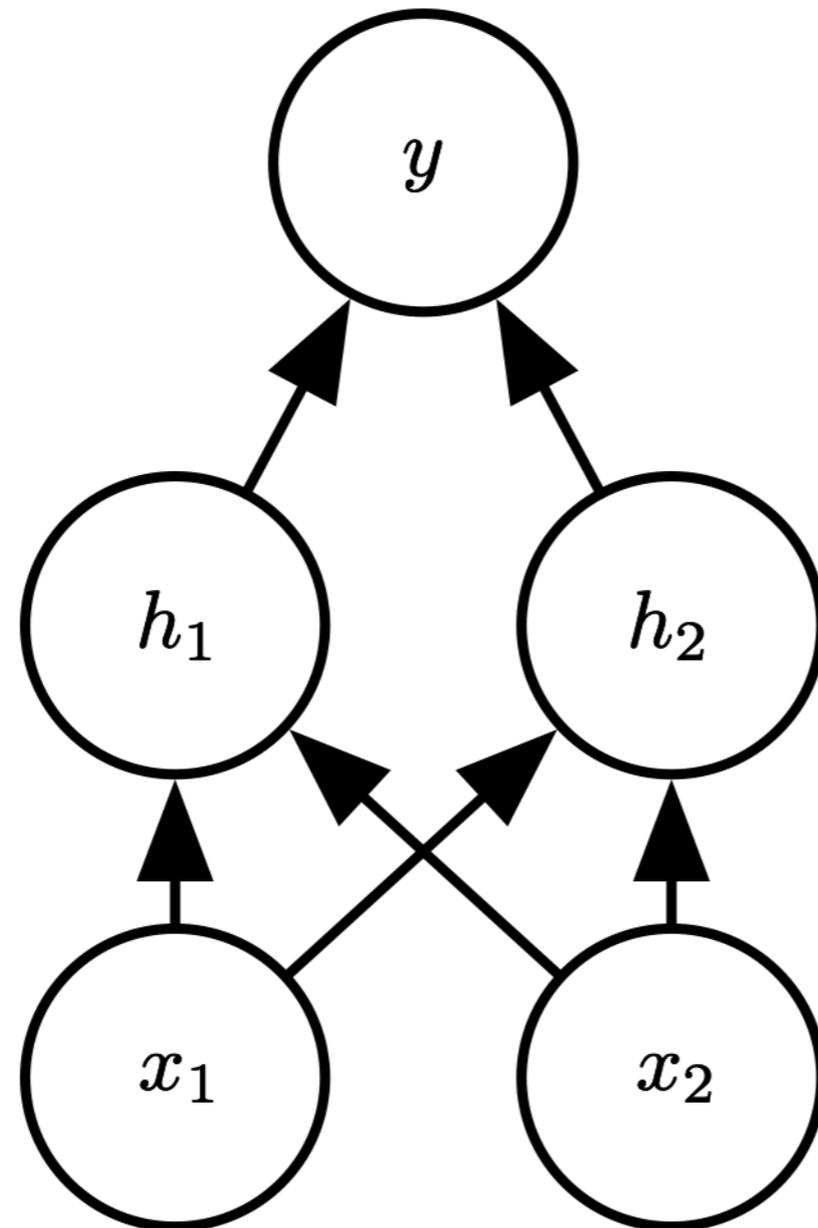
$$f(\mathbf{x}; \mathbf{w}, b) = \mathbf{x}^\top \mathbf{w} + b.$$

- Normal equations ?

- $\mathbf{W} = [0,0]$, $b=0.5$

Network Diagrams

$$f(x) = w^\top W^\top x.$$



$$f^{(2)}(h) = h^\top w$$

$$f^{(1)}(x) = W^\top x$$

Figure 6.2

Network Diagrams

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

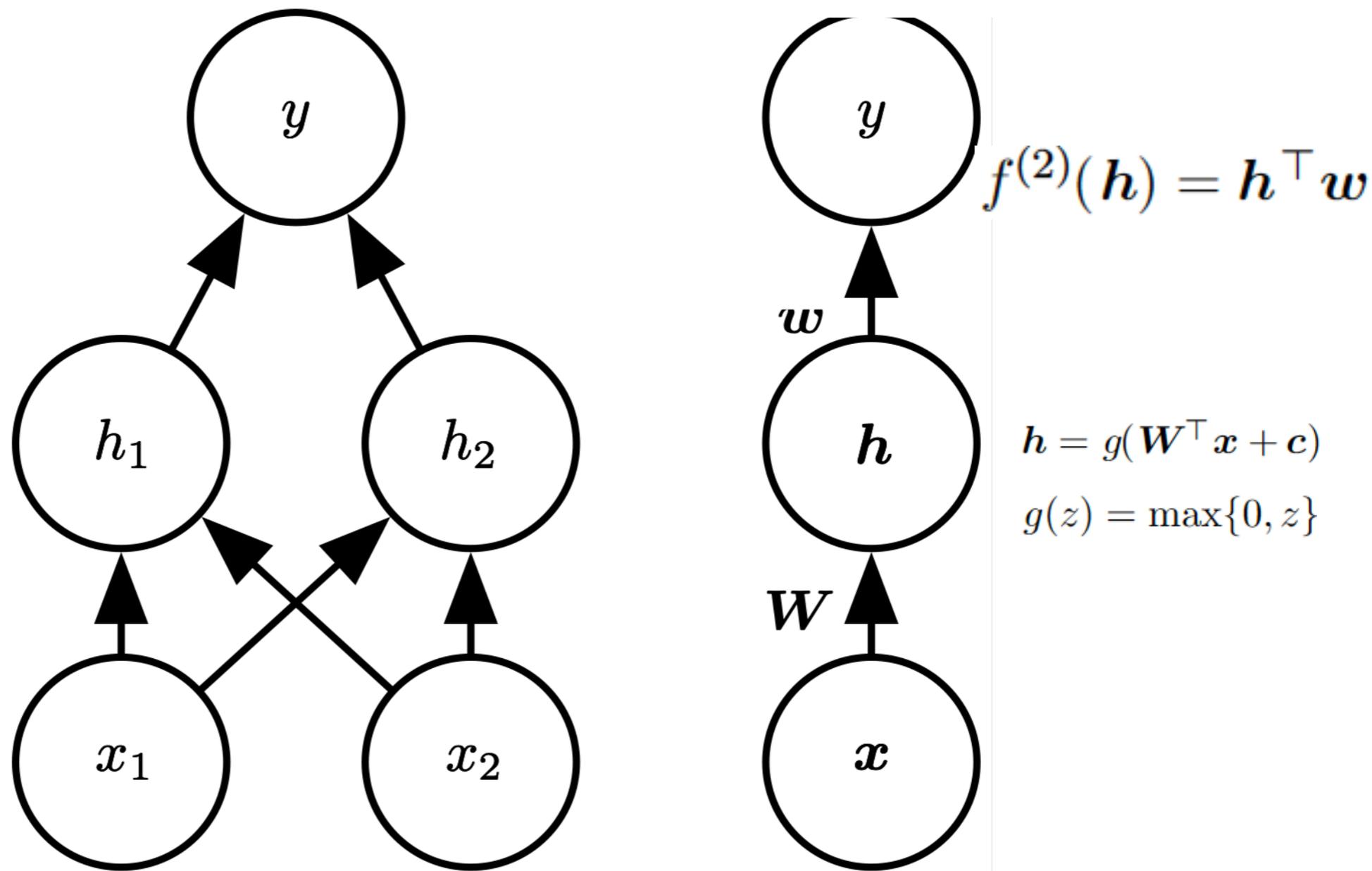


Figure 6.2

Rectified Linear Activation

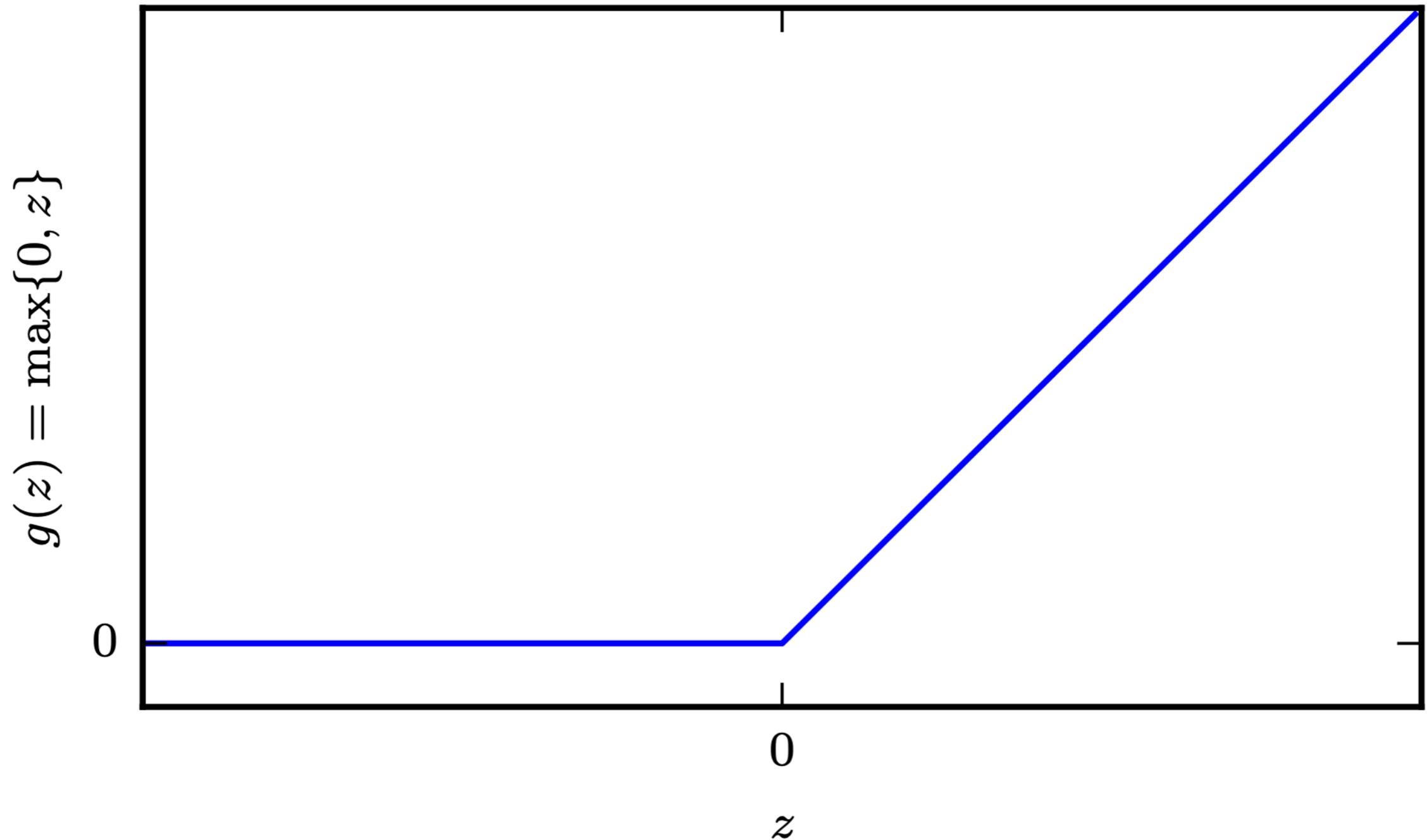


Figure 6.3

Solving XOR

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

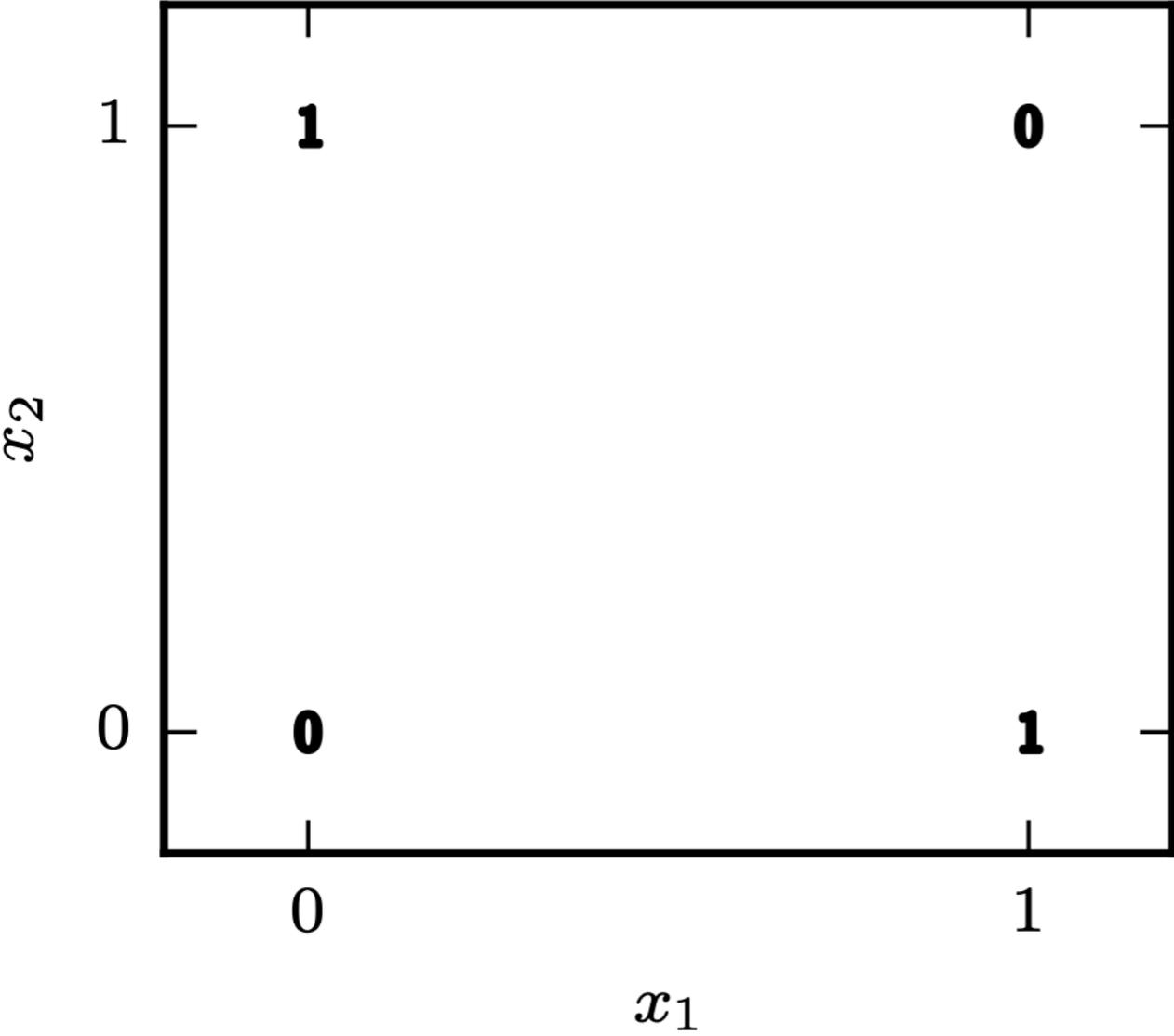
$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad (6.4)$$

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

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

Solving XOR

Original x space



Learned h space

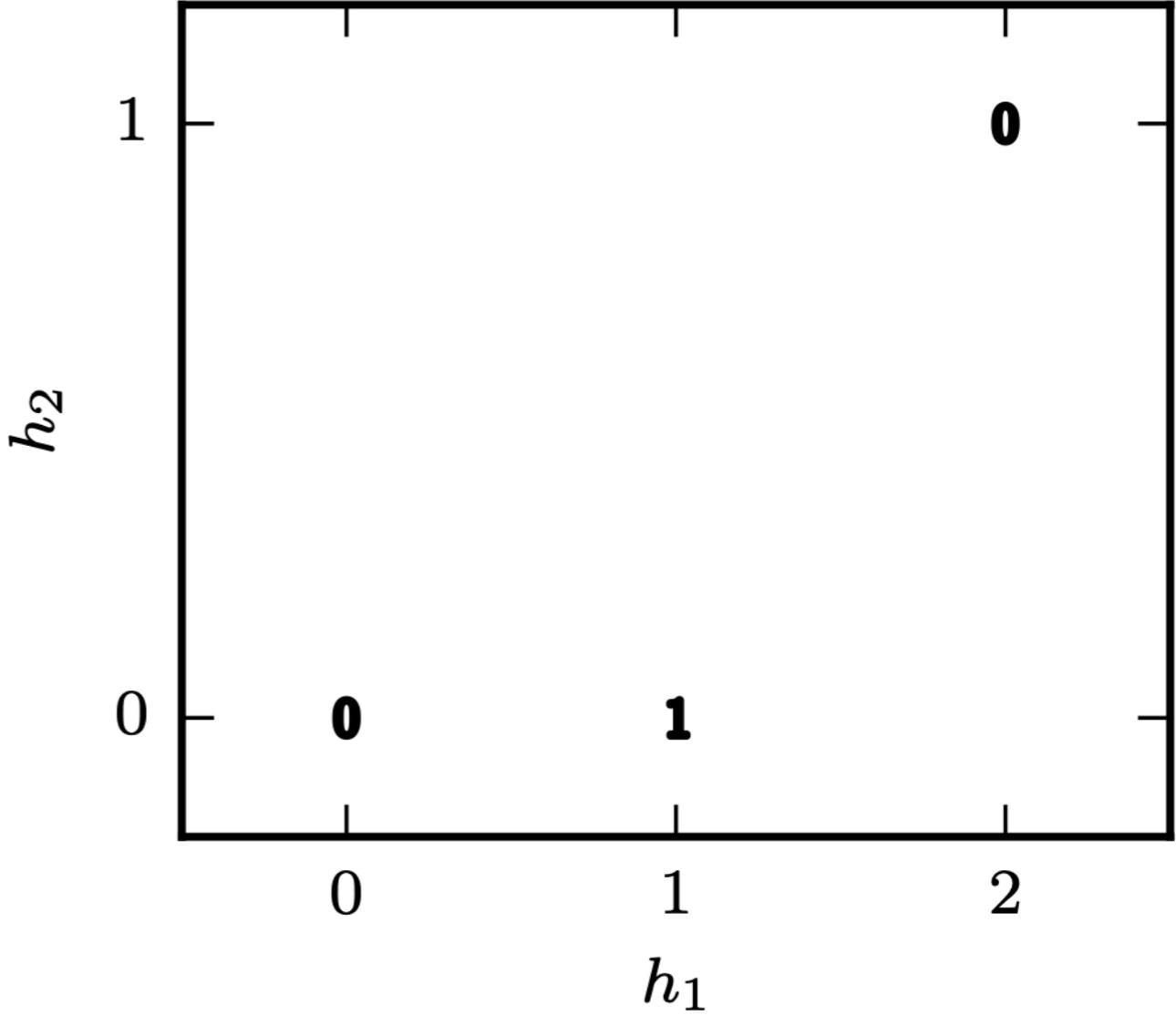


Figure 6.1

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
 - Rather than linear equations solvers used in *linear regression*
 - Or convex optimization algorithms used in *SVM and logistic regression*

Conditional Distributions and Cross-Entropy

$$J(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} L(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m L(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, \boldsymbol{\theta})$$

$$L(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) = -\log p(\mathbf{y} \mid \mathbf{x}; \boldsymbol{\theta}).$$

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

Stochastic gradient descent

$$J(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, y \sim \hat{p}_{\text{data}}} L(\mathbf{x}, y, \boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m L(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta})$$

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} L(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta}).$$

Computing the gradient is $O(m)$

Sample a mini-batch

$$\mathbb{B} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m')}\}$$

$$\mathbf{g} = \frac{1}{m'} \nabla_{\boldsymbol{\theta}} \sum_{i=1}^{m'} L(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta})$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \epsilon \mathbf{g}$$

Stochastic gradient descent

- No convergence guarantees
- Sensitive to the values of the initial parameters
- recipe
 - Initialize all weights to small random numbers
 - Biases must be initialized to zero or small positive values

Linear output units

- Suppose the network provides a set of hidden features $\mathbf{h} = f(\mathbf{x}, \boldsymbol{\theta})$
- Predict a vector of continuous variables \mathbf{y}
- Linear output unit: $\hat{\mathbf{y}} = \mathbf{W}^T \mathbf{h} + \mathbf{b}$
- It corresponds to produce the mean of a conditional Gaussian distribution:
 - $p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}; \hat{\mathbf{y}}; \mathbf{I})$
 - Minimizing cross entropy is equivalent to minimizing the mean squared error
- Linear units do not saturate

Sigmoid output units

- Predict a binary variable y (binary classification problem)
- $\hat{y} = P(y = 1 | \mathbf{x})$ based on: $\mathbf{h} = f(\mathbf{x}, \boldsymbol{\theta})$
- Bad idea: $P(y = 1 | \mathbf{x}) = \max \{ 0, \min \{ 1, \mathbf{w}^\top \mathbf{h} + b \} \}$
- Sigmoid output unit $\hat{y} = \sigma(\mathbf{w}^\top \mathbf{h} + b)$
 - $z = \mathbf{w}^\top \mathbf{h} + b$ is called logit
 - $\hat{y} = \sigma(z)$
- The log in the cost function undoes the exp of the sigmoid
 - No gradient saturation

Saturation

- $J(\theta) = -\log \sigma(z)$ if $y = 1$
- $J(\theta) = -\log(1 - \sigma(z)) = -\log(\sigma(-z))$ if $y = 0$
- One equation: $J(\theta) = -\log \sigma((2y - 1)z)$
 $= \zeta((1 - 2y)z)$
- Saturates only when $(1 - 2y)z$ is very negative
 - When we have converged to the solution
- In the limit of extremely incorrect z , the softplus function does not shrink the gradient at all

Softmax output units

- A generalization of the sigmoid to represent a distribution over n values
 - multi-label classification
 - In fact it is **soft arg-max**, winner-take-all
 - one of the outputs is nearly one and the others are nearly 0.

- $\hat{y}_i = P(y = i|x) = \text{softmax}(\mathbf{z})_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$

- $\mathbf{z} = \mathbf{W}^T \mathbf{h} + \mathbf{b}$ (we can think that \mathbf{z} are unnormalized log probabilities given by the network)

this term cannot saturate, learning can proceed

- Cost function: $-\log(\text{softmax}(\mathbf{z})_i) = -z_i + \log \sum_j \exp(z_j)$

- We can impose a requirement that one element of \mathbf{z} be fixed.

- But it is simpler to implement the overparametrized version.

Saturation

- If z_i is extremely positive \Rightarrow saturates to 1
- If z_i is extremely negative \Rightarrow saturates to 0

$$\text{softmax}(\mathbf{z}) = \text{softmax}(\mathbf{z} - \max_i z_i).$$

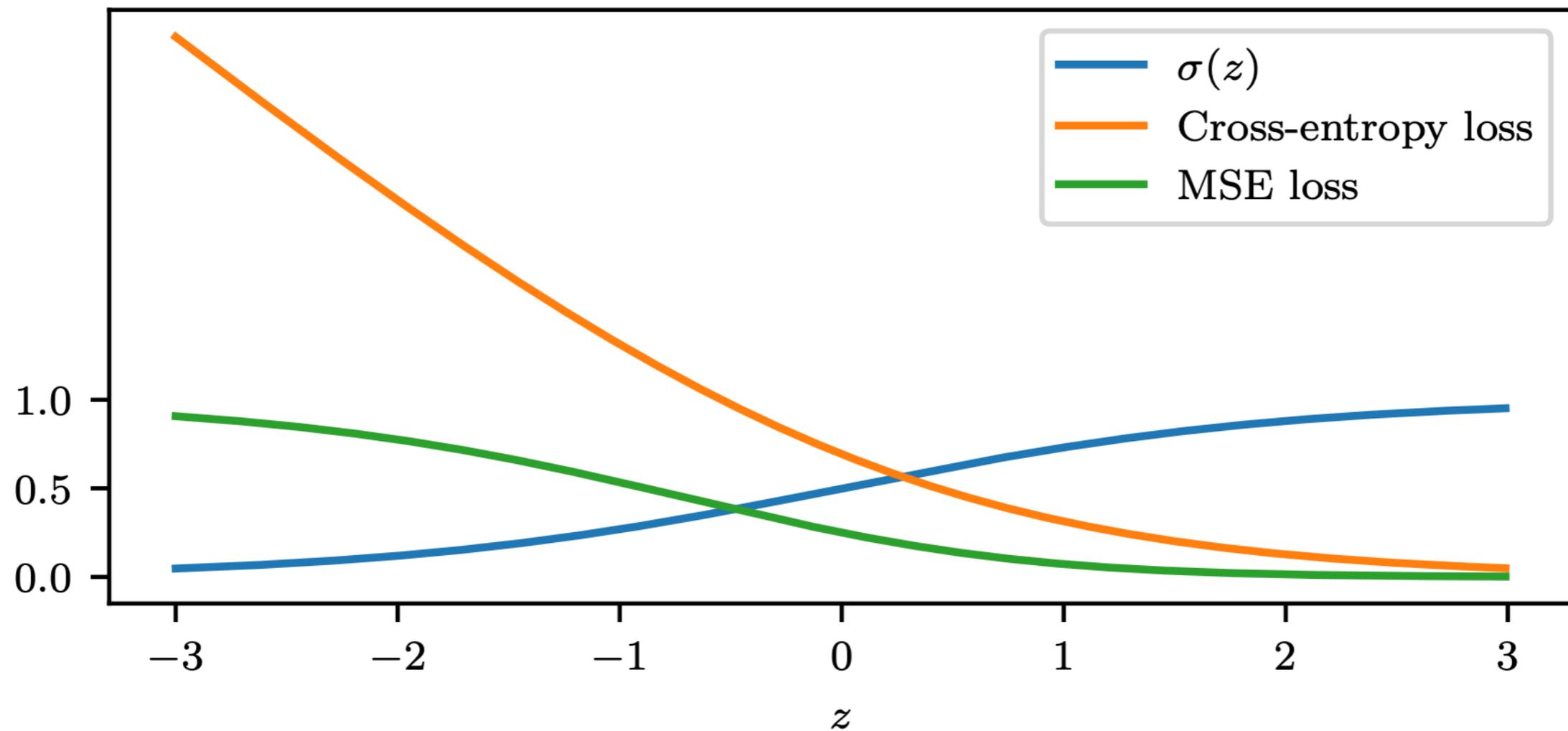
- An output $\text{softmax}(\mathbf{z})_i$ saturates to 1 when the corresponding input is maximal ($z_i = \max_i z_i$) and z_i is much greater than all the other inputs.
 - The output $\text{softmax}(\mathbf{z})_i$ can also saturate to 0 when z_i is not maximal and the maximum is much greater.
- MSE loss function will perform poorly because of these gradient vanishing problems

Output Types

Output Type	Output Distribution	Output Layer	Cost Function
Binary	Bernoulli	Sigmoid	Binary cross-entropy
Discrete	Multinoulli	Softmax	Discrete cross-entropy
Continuous		Linear	(MSE)

Don't mix and match

Sigmoid output with target of 1



Other output units

$$\mathcal{N}(x; \mu, \beta^{-1}) = \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{1}{2}\beta(x - \mu)^2\right). \quad (3.22)$$

- *Design a neural network to estimate the variance of a gaussian distribution representing $P(y|x)$*
- Loss function:
 - $-\log \mathcal{N}(x; \mu; \beta^{-1}) \propto -\frac{1}{2}\log \beta + \frac{1}{2}\beta(x - \mu)^2$
- Check notebook

Roadmap

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

Hidden units

- Most hidden units can be described as:
 - accepting a vector of inputs \mathbf{x} ,
 - computing an affine transformation $\mathbf{z} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$,
 - And applying an element-wise nonlinear function $g(\mathbf{z})$.
- Use ReLUs, 90% of the time
- For some research projects, get creative
- Many hidden units perform comparably to ReLUs. New hidden units that perform comparably are rarely interesting.

ReLU is not differentiable at 0!

- The functions used in the context of neural networks usually have defined left derivatives and defined right derivatives.
 - In the case of $g(z) = \max\{0, z\}$,
 - the left derivative at $z = 0$ is 0
 - and the right derivative is 1.
 - Software implementations of neural network training usually return one of the one-sided derivatives
 - rather than reporting that the derivative is undefined or raising an error.

When ReLu is active

- When the input is positive, we say that the rectifier is active:
 - the gradient is 1
 - The second derivative is 0, no second-order effects.
- When initializing the parameters of $W^T x + b$,
 - set all elements of b to a small, positive value, such as 0.1
 - This makes it very likely that the rectified linear units will be initially active for most inputs in the training set and allow the derivatives to pass through.

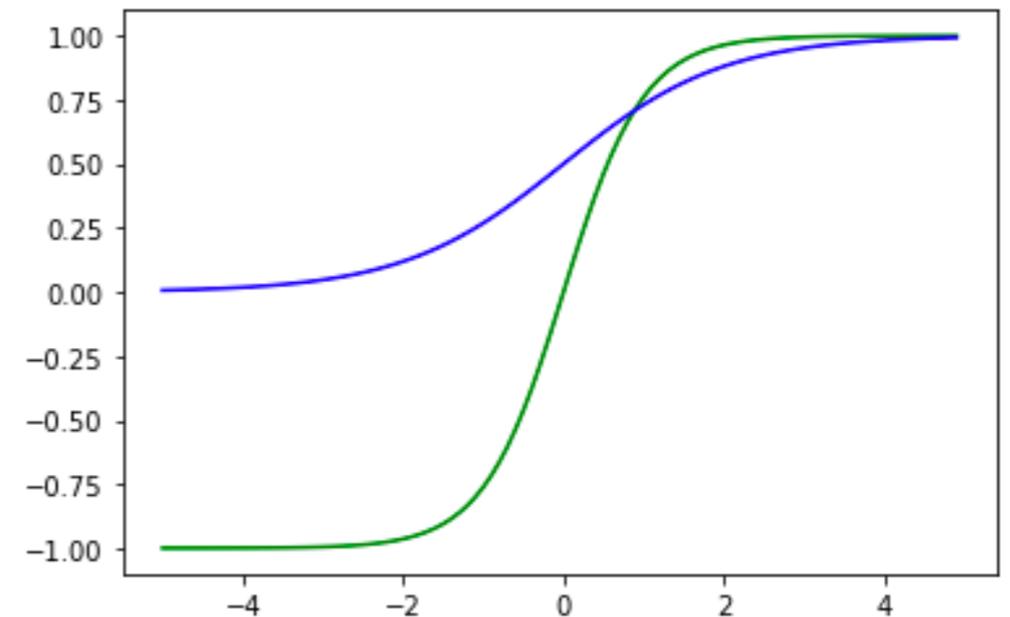
ReLU generalizations

$$h_i = g(\mathbf{z}, \boldsymbol{\alpha})_i = \max(0, z_i) + \alpha_i \min(0, z_i)$$

- Absolute value rectification: $\alpha_i = -1$
- Leaky ReLU: $\alpha_i = 0.01$
- Parametric ReLU: α_i is learnable
- Rectified linear units and its generalizations are based on the principle that models are easier to optimize if their behavior is closer to linear.

Other hidden units

- Sigmoid
 - $g(z) = \sigma(z)$
 - They can saturate
 - which makes learning difficult
- Hyperbolic tangent
 - $g(z) = \tanh(z)$.
 - Related to sigmoid: $\tanh(z) = 2\sigma(2z) - 1$.
 - It typically performs better
 - since it resembles the identity function
 - So long as the activations of the network can be kept small.



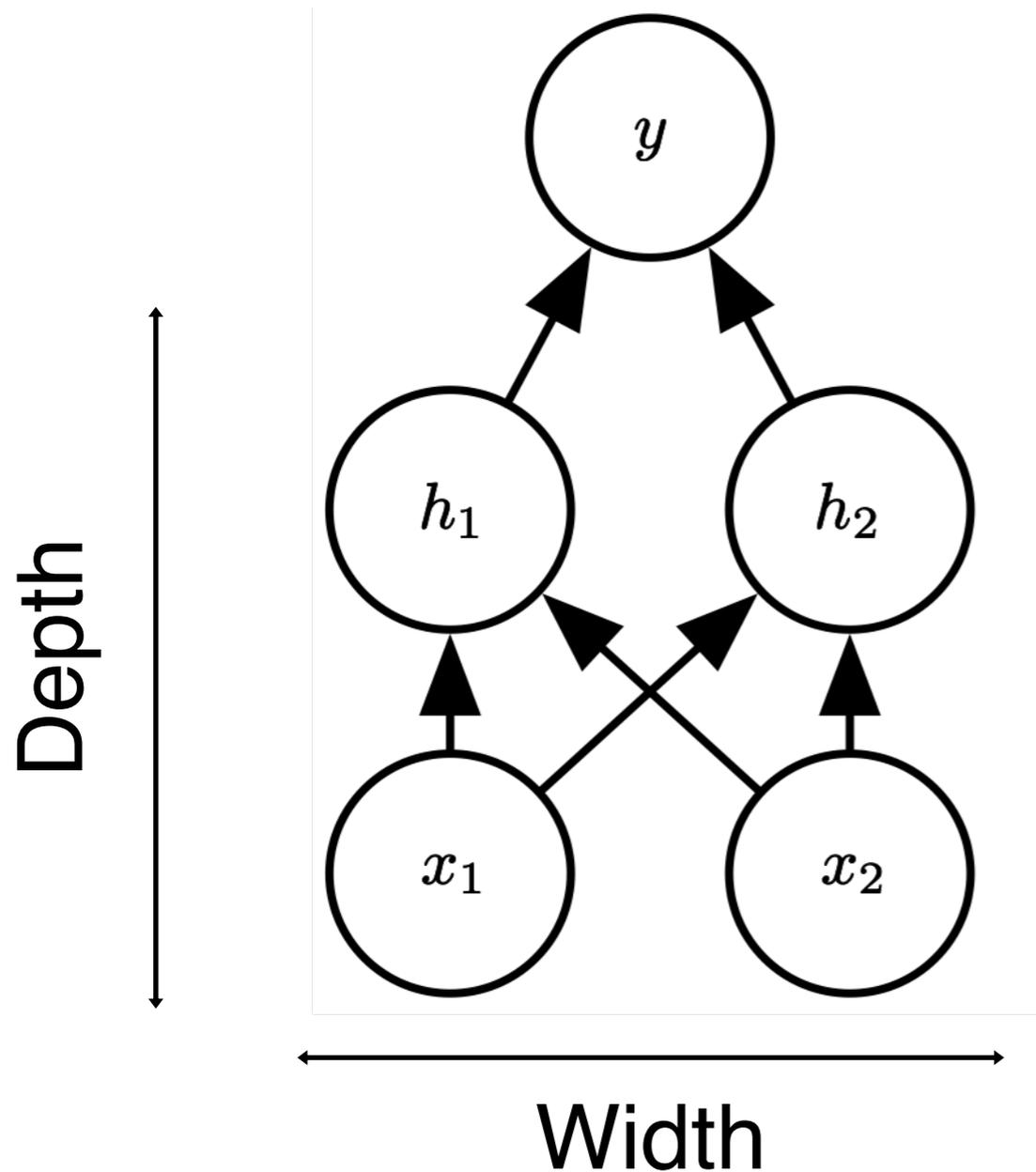
Other hidden units

- The authors tested a feedforward network using $\mathbf{h} = \cos(\mathbf{W}^T \mathbf{x} + \mathbf{b})$ on the MNIST dataset and obtained an error rate of less than 1%
- Linear hidden units offer an effective way of reducing the number of parameters in a network
 - Assume \mathbf{W}^T is $p \times n$
 - $\mathbf{W}^T = \mathbf{V}^T \mathbf{U}^T$ where \mathbf{V}^T is $p \times q$ and \mathbf{U}^T is $q \times n$
 - In total we have $q \times (n+p)$ trainable params
 - Much less than $\ll p \times n$ for small q

Roadmap

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

Architecture Basics



- $h^{(1)} = g^{(1)}(W^{(1)T}x + b^{(1)})$
- $h^{(2)} = g^{(2)}(W^{(2)T}h^{(1)} + b^{(2)})$
- ...

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 nets may need (exponentially) more width
 - Shallow nets may overfit more

Bad idea: One hidden layer

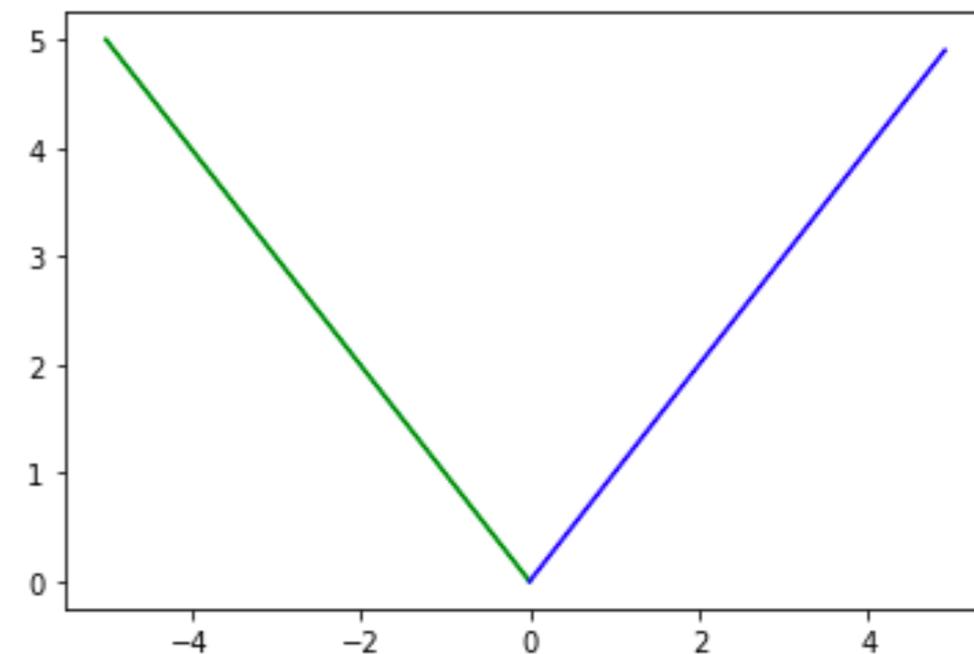
- Consider binary functions over $v \in \{0,1\}^n$
 - The image of any member of the possible 2^n inputs can be 0 or 1
 - In total, we have 2^{2^n} possible functions
 - Selecting one such function requires $O(2^n)$ degrees of freedom
- Working with a single layer is sufficient to represent ‘any’ function
 - Yet the width of the layer can be exponential, and prone to overfitting

Montufar et. al. (2014)

- A deep rectifier network can represent functions with a number of regions that is exponential in the depth of the network.

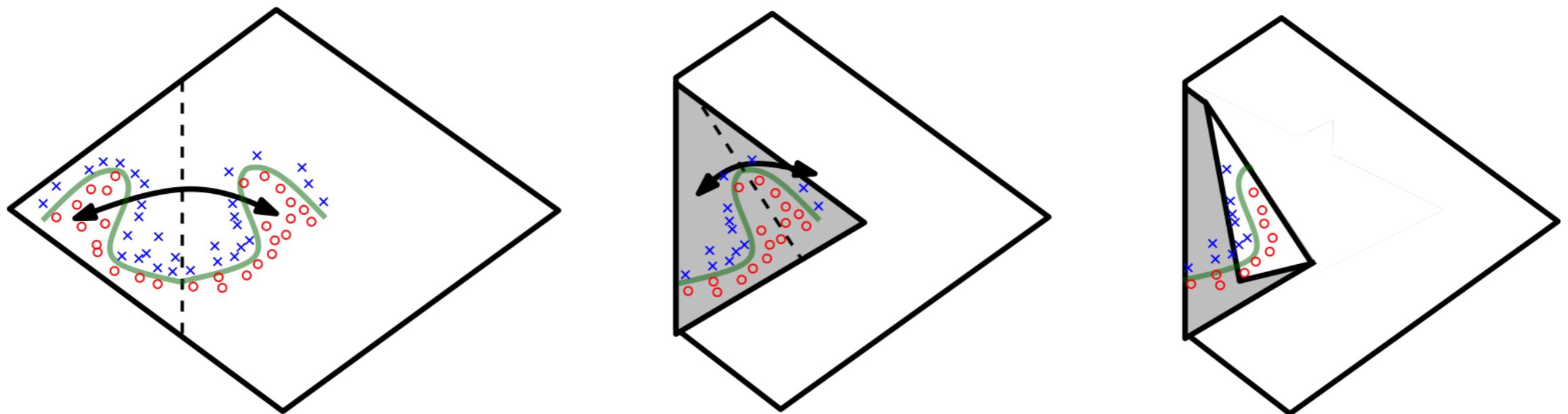
Example: Absolute value rectification

- Has the same output for every pair of mirror points
- The mirror axis of symmetry is given by $Wh + b$
- A function computed on top of that unit is simpler
- If we fold again, we get an even simpler function, and so on ..



Exponential Representation

Advantage of Depth

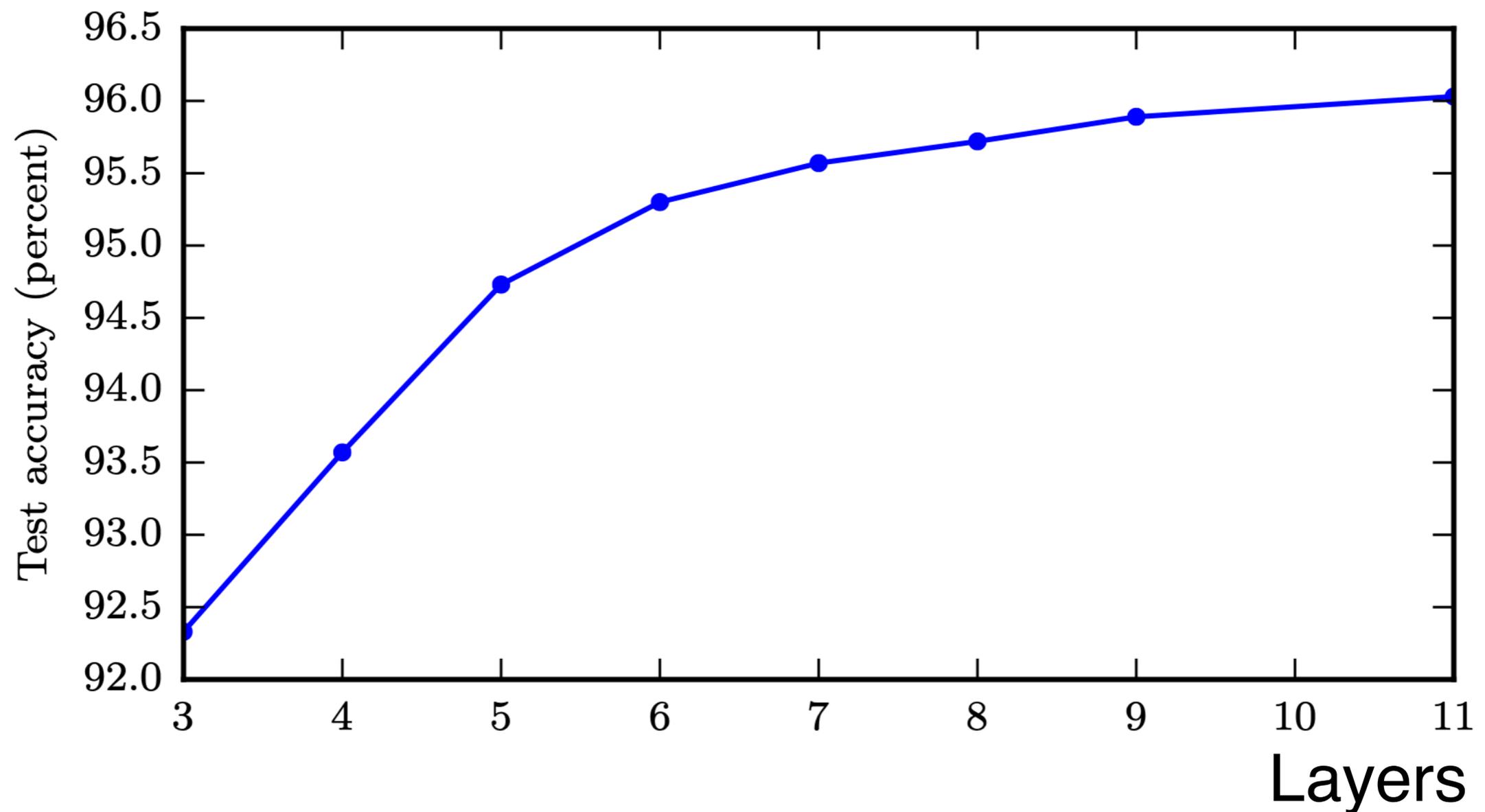


We obtain an exponentially large number of piecewise linear regions which can capture all kinds of repeating patterns.

Why deep?

- Choosing a deep model encodes a very general belief (or prior) that the function we want to learn should involve composition of several simpler functions
- Alternately, we can interpret the use of a deep architecture as expressing a belief that the function we want to learn is a computer program consisting of multiple steps, where each step makes use of the previous step's output
- Empirically, greater depth does seem to result in better generalization for a wide variety of tasks

Better Generalization with Greater Depth



Task: transcribe multi-digit numbers from photographs of addresses

Large, Shallow Models Overfit More

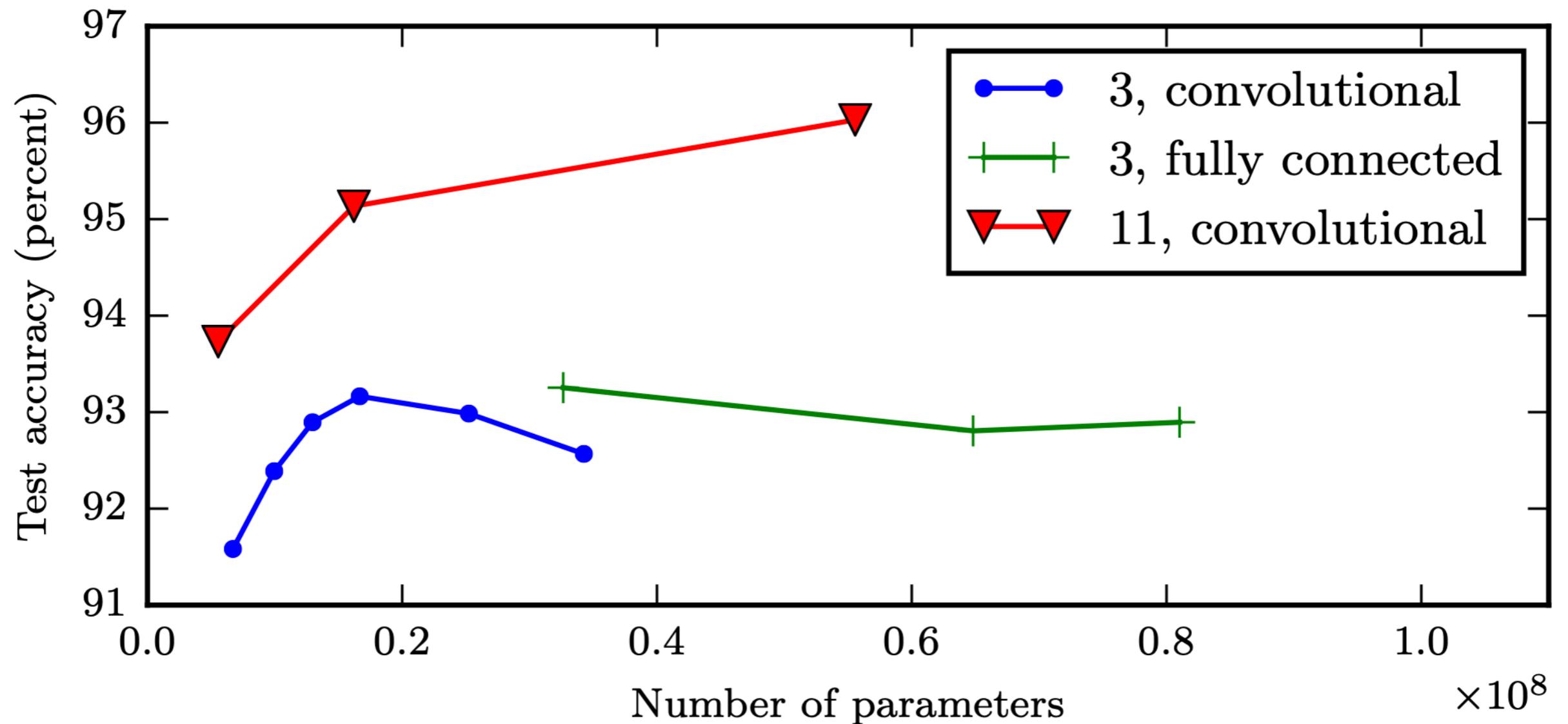


Figure 6.7

Other architectural considerations

- Other architectures: CNN, RNN
- Skip connections: going from layer i to layer $i + 2$ or higher.
 - make it easier for the gradient to flow from output layers to layers nearer the input.

Design some neural networks for fun:

<https://playground.tensorflow.org/>

Roadmap

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

Compute the gradient

- Forward propagation: $x \rightarrow \hat{y} \rightarrow J(\theta)$
- Backprop: $\frac{\partial J(\theta)}{\partial \theta} = ?$
 - Computing derivatives by propagating information backward through the network
 - Simple and inexpensive
 - Based on the chain rule:

$$\theta = \mathbf{W}, \mathbf{c}, \mathbf{w}, b$$

$$\nabla_{\theta} J = ?$$

- $\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{c}$

- $\mathbf{h} = g(\mathbf{z})$

- $\hat{y} = \mathbf{w}^T \mathbf{h} + b$

- $J(\hat{y}) = (y - \hat{y})^2$

$$\frac{\partial J}{\partial c_1} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h_1} \frac{\partial h_1}{\partial z_1} \frac{\partial z_1}{\partial c_1}$$

$$\frac{\partial J}{\partial W_{11}} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h_1} \frac{\partial h_1}{\partial z_1} \frac{\partial z_1}{\partial W_{11}}$$

Chain Rule

- More generally

- $\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^m, z \in \mathbb{R}$

- $z = f(\mathbf{y}) = f(g(\mathbf{x}))$

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

- Vector notation:

- $\nabla_{\mathbf{x}} z = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T \nabla_{\mathbf{y}} z$

- Or: $\nabla_{\mathbf{x}} z = \sum_j \frac{\partial z}{\partial y_j} \nabla_{\mathbf{x}} y_j$

$n \times m$
Jacobian
matrix

- Backprop consists of performing Jacobian-gradient product for each operation in the graph.

Back-Propagation

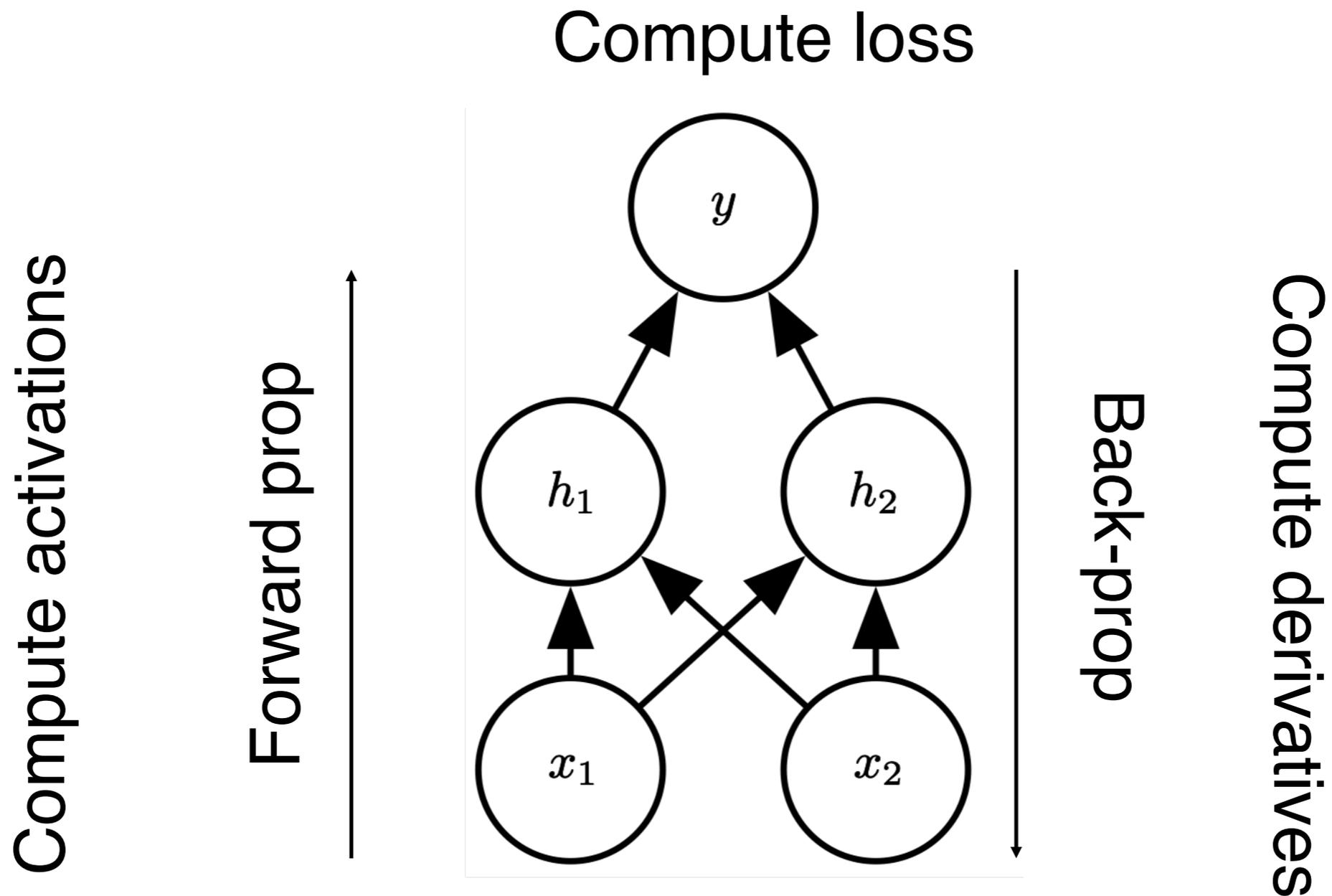
- Back-propagation is “**just the chain rule**” of calculus

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

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

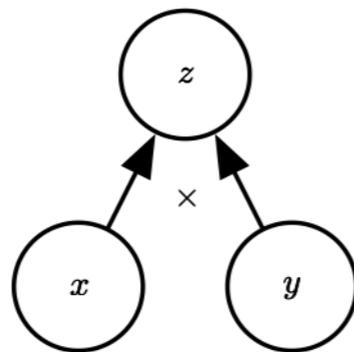
- But it’s a particular implementation of the chain rule
 - Uses dynamic programming (table filling)
 - Avoids recomputing repeated subexpressions
 - Speed vs memory tradeoff

Simple Back-Prop Example

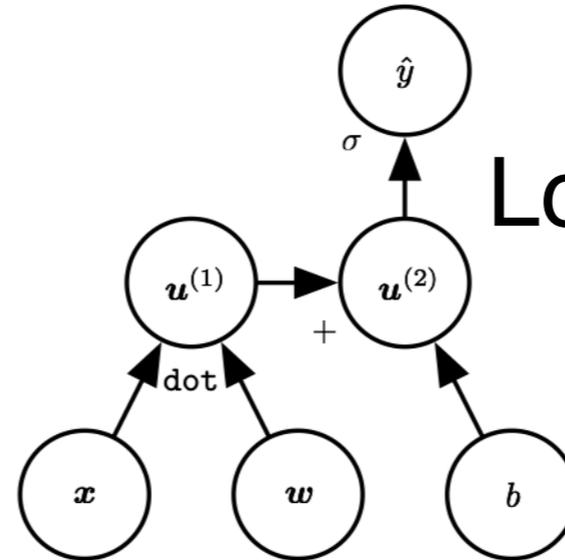


General implementation Computation Graphs

Multiplication



(a)



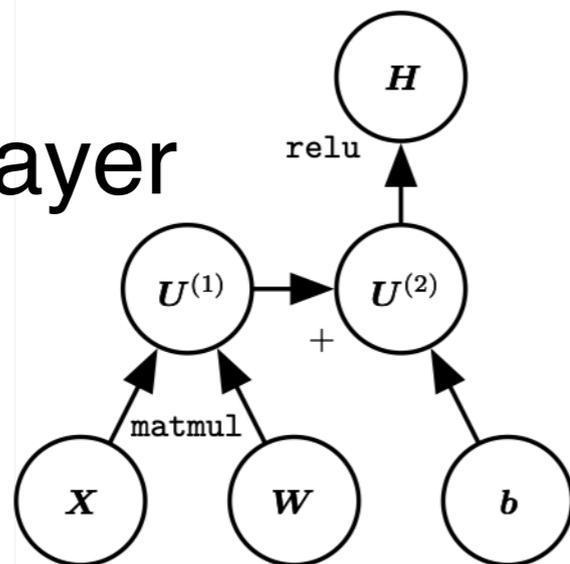
(b)

Logistic regression

$$\hat{y} = x^T w$$

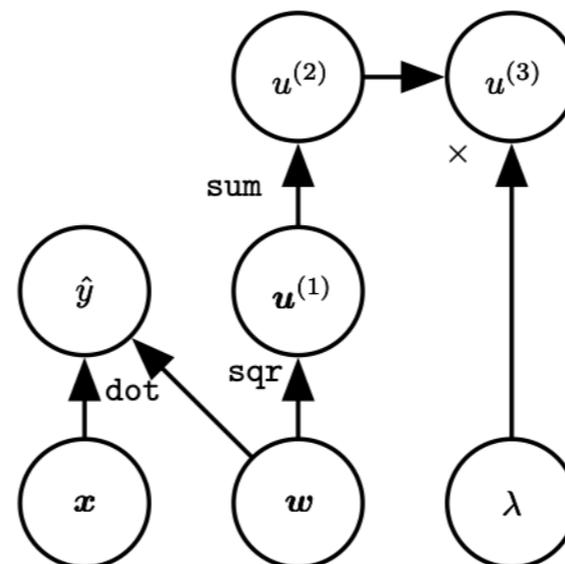
$$u^3 = \lambda \sum w_i^2$$

ReLU layer



(c)

Mini-batch

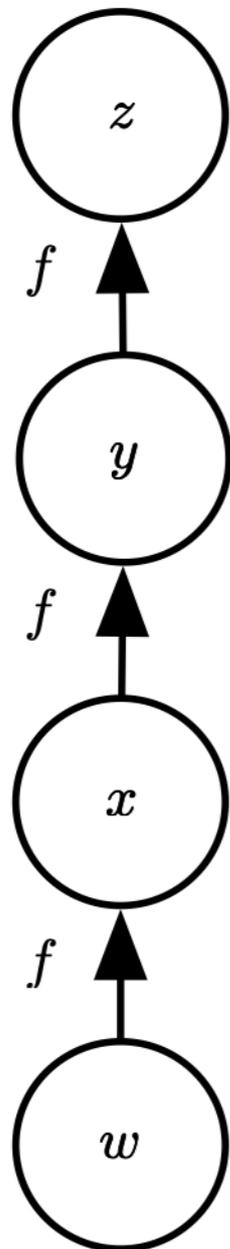


(d)

Linear regression
and weight decay

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

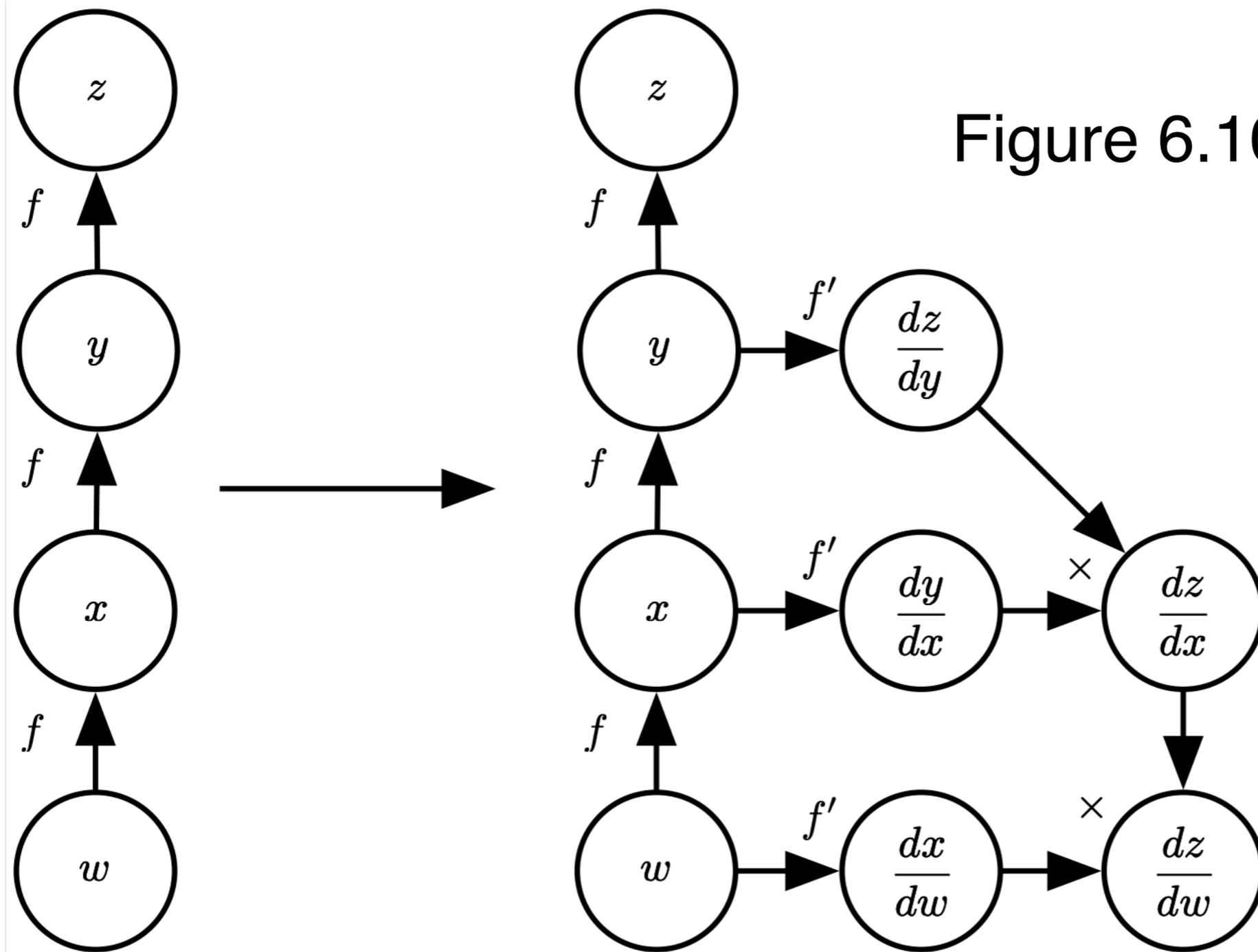


Figure 6.10

Implementing the chain rule

- Given a computation graph, we wish to compute

$$\frac{\partial u^{(n)}}{\partial u^{(i)}}, \text{ for } i = 1, \dots, n_i$$

- $u^{(n)}$ can be the cost function

- $u^{(i)}$ are the trainable parameters of the network

- $u^{(i)}$ is associated with an operation f

- $u^{(i)} = f(\mathbb{A}_i)$ where \mathbb{A}_i is the set of parents of $u^{(i)}$

$$\frac{\partial u^{(n)}}{\partial u^{(i)}} = \sum_{j:i \in Pa(u^{(j)})} \frac{\partial u^{(n)}}{\partial u^{(j)}} \frac{\partial u^{(j)}}{\partial u^{(i)}}$$

Forward propagation

Algorithm 6.1 A procedure that performs the computations mapping n_i inputs $u^{(1)}$ to $u^{(n_i)}$ to an output $u^{(n)}$. This defines a computational graph where each node computes numerical value $u^{(i)}$ by applying a function $f^{(i)}$ to the set of arguments $\mathbb{A}^{(i)}$ that comprises the values of previous nodes $u^{(j)}$, $j < i$, with $j \in Pa(u^{(i)})$. The input to the computational graph is the vector \mathbf{x} , and is set into the first n_i nodes $u^{(1)}$ to $u^{(n_i)}$. The output of the computational graph is read off the last (output) node $u^{(n)}$.

```
for  $i = 1, \dots, n_i$  do
     $u^{(i)} \leftarrow x_i$ 
end for
for  $i = n_i + 1, \dots, n$  do
     $\mathbb{A}^{(i)} \leftarrow \{u^{(j)} \mid j \in Pa(u^{(i)})\}$ 
     $u^{(i)} \leftarrow f^{(i)}(\mathbb{A}^{(i)})$ 
end for
return  $u^{(n)}$ 
```

Backward propagation

Algorithm 6.2 Simplified version of the back-propagation algorithm for computing the derivatives of $u^{(n)}$ with respect to the variables in the graph. This example is intended to further understanding by showing a simplified case where all variables are scalars, and we wish to compute the derivatives with respect to $u^{(1)}, \dots, u^{(n_i)}$. This simplified version computes the derivatives of all nodes in the graph. The computational cost of this algorithm is proportional to the number of edges in the graph, assuming that the partial derivative associated with each edge requires a constant time. This is of the same order as the number of computations for the forward propagation. Each $\frac{\partial u^{(i)}}{\partial u^{(j)}}$ is a function of the parents $u^{(j)}$ of $u^{(i)}$, thus linking the nodes of the forward graph to those added for the back-propagation graph.

Run forward propagation (algorithm 6.1 for this example) to obtain the activations of the network

Initialize `grad_table`, a data structure that will store the derivatives that have been computed. The entry `grad_table[u(i)]` will store the computed value of $\frac{\partial u^{(n)}}{\partial u^{(i)}}$.

`grad_table[u(n)] ← 1`

for $j = n - 1$ down to 1 do

The next line computes $\frac{\partial u^{(n)}}{\partial u^{(j)}} = \sum_{i:j \in Pa(u^{(i)})} \frac{\partial u^{(n)}}{\partial u^{(i)}} \frac{\partial u^{(i)}}{\partial u^{(j)}}$ using stored values:

`grad_table[u(j)] ← $\sum_{i:j \in Pa(u^{(i)})} \text{grad_table}[u^{(i)}] \frac{\partial u^{(i)}}{\partial u^{(j)}}$`

end for

return {`grad_table[u(i)]` | $i = 1, \dots, n_i$ }

$$\frac{\delta u^n}{\delta u^n} = 1$$

Efficiency

- The amount of computation required for performing the back-propagation scales linearly with the number of edges in G ,
 - the computation for each edge corresponds to computing a partial derivative (of one node with respect to one of its parents)
 - as well as performing one multiplication and one addition.
- Naïve approach:

$$\frac{\partial u^{(n)}}{\partial u^{(j)}} = \sum_{\substack{\text{path } (u^{(\pi_1)}, u^{(\pi_2)}, \dots, u^{(\pi_t)}), \\ \text{from } \pi_1=j \text{ to } \pi_t=n}} \prod_{k=2}^t \frac{\partial u^{(\pi_k)}}{\partial u^{(\pi_{k-1})}}.$$

Two implementation approaches

- Take a computational graph and a set of numerical values for the inputs to the graph, then return a set of numerical values describing the gradient at those input values. We call this approach “symbol-to-number” differentiation. This is the approach used by **PyTorch**.
- Another approach is to take a computational graph and add additional nodes to the graph that provide a symbolic description of the desired derivatives. This is the approach taken **TensorFlow**.
 - Because the derivatives are just another computational graph, it is possible to run back-propagation again, derivatives.

Forward propagation for a neural network

Algorithm 6.3 Forward propagation through a typical deep neural network and the computation of the cost function. The loss $L(\hat{\mathbf{y}}, \mathbf{y})$ depends on the output $\hat{\mathbf{y}}$ and on the target \mathbf{y} (see section 6.2.1.1 for examples of loss functions). To obtain the total cost J , the loss may be added to a regularizer $\Omega(\theta)$, where θ contains all the parameters (weights and biases). Algorithm 6.4 shows how to compute gradients of J with respect to parameters \mathbf{W} and \mathbf{b} . For simplicity, this demonstration uses only a single input example \mathbf{x} . Practical applications should use a minibatch. See section 6.5.7 for a more realistic demonstration.

Require: Network depth, l

Require: $\mathbf{W}^{(i)}, i \in \{1, \dots, l\}$, the weight matrices of the model

Require: $\mathbf{b}^{(i)}, i \in \{1, \dots, l\}$, the bias parameters of the model

Require: \mathbf{x} , the input to process

Require: \mathbf{y} , the target output

$$\mathbf{h}^{(0)} = \mathbf{x}$$

for $k = 1, \dots, l$ **do**

$$\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}$$

$$\mathbf{h}^{(k)} = f(\mathbf{a}^{(k)})$$

end for

$$\hat{\mathbf{y}} = \mathbf{h}^{(l)}$$

$$J = L(\hat{\mathbf{y}}, \mathbf{y}) + \lambda \Omega(\theta)$$

Gradients of biases and weights

- We have $\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}$

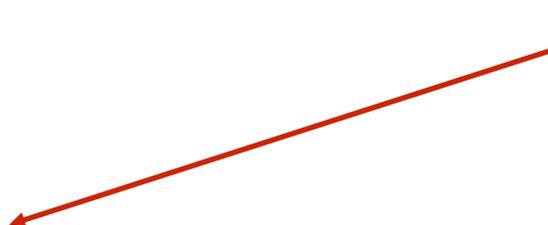
- $\frac{\partial J}{\partial b_i^{(k)}} = \frac{\partial L}{\partial a_i^{(k)}} \cdot 1 + \lambda \frac{\partial \Omega}{\partial b_i^{(k)}}$

- Vectorized: $\nabla_{\mathbf{b}^{(k)}} J = \nabla_{\mathbf{a}^{(k)}} L + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega$

- $\frac{\partial J}{\partial W_{i,j}^{(k)}} = \frac{\partial L}{\partial a_i^{(k)}} \cdot h_j^{(k-1)} + \lambda \frac{\partial \Omega}{\partial W_{i,j}^{(k)}}$

- Vectorized: $\nabla_{\mathbf{W}^{(k)}} J = (\nabla_{\mathbf{a}^{(k)}} L) \mathbf{h}^{(k-1)T} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega$

Outer Product



Gradient of representations

- We have $\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}$
 - $\nabla_{\mathbf{h}^{(k-1)}} L = \left(\frac{\partial \mathbf{a}^{(k)}}{\partial \mathbf{h}^{(k-1)}} \right)^T \nabla_{\mathbf{a}^{(k)}} L = \mathbf{W}^{(k)T} \nabla_{\mathbf{a}^{(k)}} L$
- $\mathbf{h}^{(k)} = f(\mathbf{a}^{(k)})$
 - $\nabla_{\mathbf{a}^{(k)}} L = \nabla_{\mathbf{h}^{(k)}} L \odot f'(\mathbf{a}^{(k)})$

Backward propagation for a neural network

Algorithm 6.4 Backward computation for the deep neural network of algorithm 6.3, which uses in addition to the input \mathbf{x} a target \mathbf{y} . This computation yields the gradients on the activations $\mathbf{a}^{(k)}$ for each layer k , starting from the output layer and going backwards to the first hidden layer. From these gradients, which can be interpreted as an indication of how each layer's output should change to reduce error, one can obtain the gradient on the parameters of each layer. The gradients on weights and biases can be immediately used as part of a stochastic gradient update (performing the update right after the gradients have been computed) or used with other gradient-based optimization methods.

After the forward computation, compute the gradient on the output layer:

$$\mathbf{g} \leftarrow \nabla_{\hat{\mathbf{y}}} J = \nabla_{\hat{\mathbf{y}}} L(\hat{\mathbf{y}}, \mathbf{y})$$

for $k = l, l - 1, \dots, 1$ do

Convert the gradient on the layer's output into a gradient into the pre-nonlinearity activation (element-wise multiplication if f is element-wise):

$$\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)})$$

Compute gradients on weights and biases (including the regularization term, where needed):

$$\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\theta)$$

$$\nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)\top} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\theta)$$

Propagate the gradients w.r.t. the next lower-level hidden layer's activations:

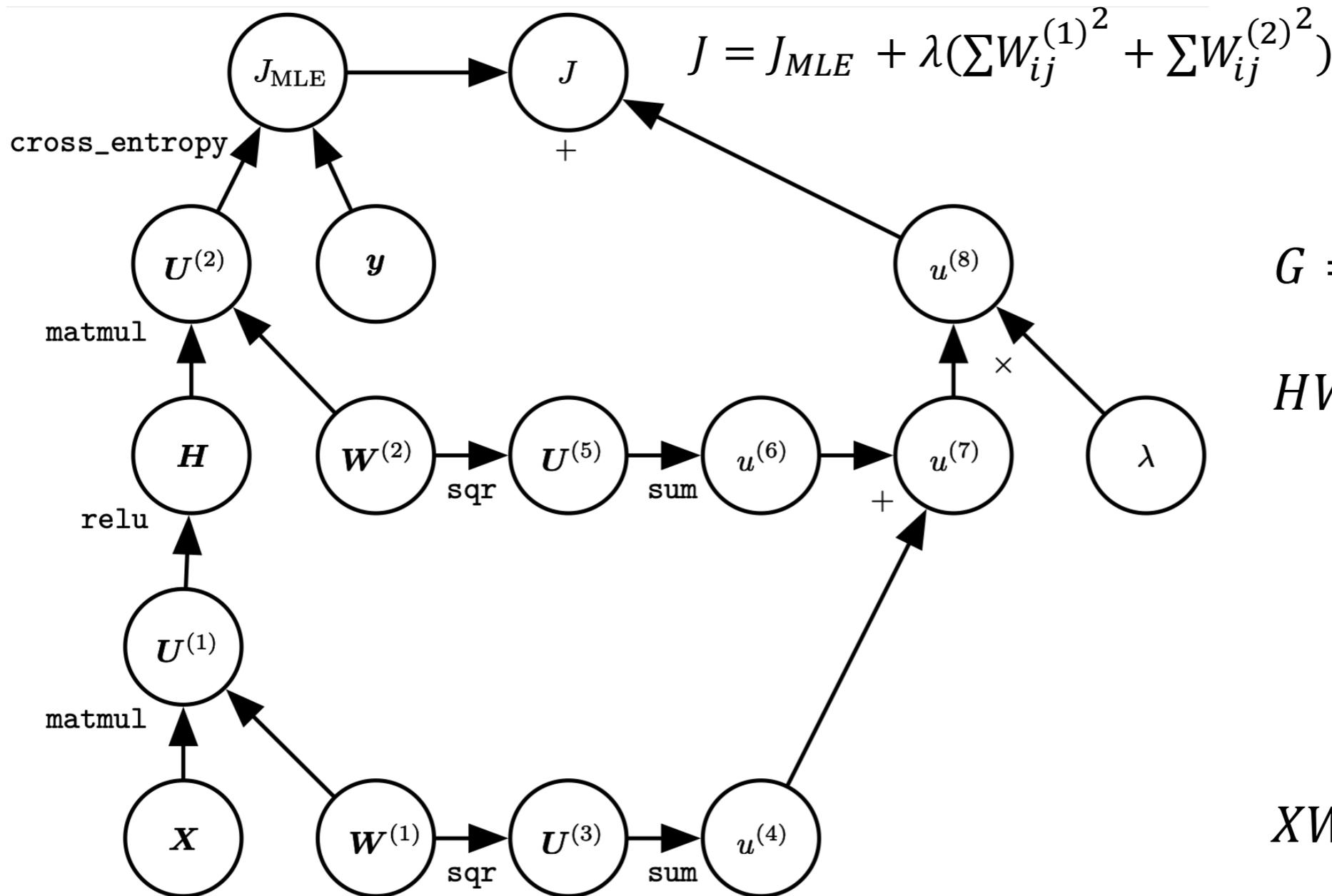
$$\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)\top} \mathbf{g}$$

end for

Gradients for matrix multiplication

- $HW = U$
 - $\nabla_H f = \nabla_U f W^T$
 - $\nabla_W f = H^T \nabla_U f$
- Useful to backprop over a batch of samples
- A layer receives G :
 - GW^T is the update for W
 - $H^T G$ is back propagated

Backprop for a batch of samples



$$G = \nabla_{U^{(2)}} J_{MLE}$$

$$HW^{(2)} = U^{(2)}$$

$$\nabla_{W^{(2)}} J = H^T G + 2\lambda W^{(2)}$$

$$\nabla_H J_{MLE} = GW^T$$

$$\nabla_{U^{(1)}} J_{MLE} = G'$$

$$XW^{(1)} = U^{(1)}$$

$$\nabla_{W^{(1)}} J = X^T G' + 2\lambda W^{(1)}$$

Cost forward propagation: $O(w)$ multiply-adds
 Cost backward propagation: $O(w)$ multiply-adds
 Memory-cost: $O(mn_h)$

Backprop in deep learning vs. the general field of automatic differentiation

- Reverse mode accumulation
- Forward mode accumulation
- Simplifying algebraic expressions (e.g. Theano)
 - $q_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$
 - $J = -\sum p_j \log(q_j)$
 - $\frac{\partial J}{\partial z_i} = q_i - p_i$
- Specialized libraries vs. automatic generation
 - Data structures with *bprop()* method
- Backprop is not the only way or the optimal way of computing the gradient, but it is practical

High-order derivatives

- Properties of the Hessian may guide the learning
- $J(\theta)$ is a function $\mathbb{R}^n \rightarrow \mathbb{R}$
 - The hessian of $J(\theta)$ is in $\mathbb{R}^{n \times n}$
- In typical deep learning, n can be in the billions
 - Computing the Hessian is not practical
- **Krylov methods** are a set of iterative techniques for performing various operations like:
 - approximately inverting a matrix
 - or finding eigenvectors or eigenvalues,
 - without using any operation other than **matrix-vector** products.

Hessian-vector Products

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

- Compute $\mathbf{H}\mathbf{e}^{(i)}$ for all $i = 1, \dots, n$
- Where $\mathbf{e}^{(i)}$ is the one hot vector:
 - $e_i^{(i)} = 1$
 - $e_j^{(i)} = 0$

Questions

- Watch Ian's lecture:
 - <https://drive.google.com/file/d/0B64011x02sIkRExCY0FDVXFcoHM/view>