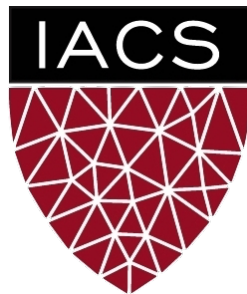


# Lecture 12-13: Basic Neural Nets

## Deep Feedforward Networks

CS 109B, STAT 121B, AC 209B, CSE 109B

Mark Glickman and Pavlos Protopapas



# Beyond Linear Models

- Linear models
  - Can be fit efficiently (via convex optimization)
  - Limited model capacity
- Alternative:

$$f(x) = w^T \phi(x)$$

where  $\phi$  is a *non-linear transform*

# Traditional ML

- Manually engineer  $\phi$ 
  - Domain specific, enormous human effort
- Generic transform
  - Maps to a higher-dimensional space
  - Kernel methods: e.g. RBF kernels
  - Over fitting: does not generalize well to test set
  - Cannot encode enough prior information

# Deep Learning

- Directly learn  $\phi$

$$f(x; \theta) = w^T \phi(x; \theta)$$

where  $\theta$  are parameters of the transform

- $\phi$  defines hidden layers
- Non-convex optimization
- Can encode prior beliefs, generalizes well

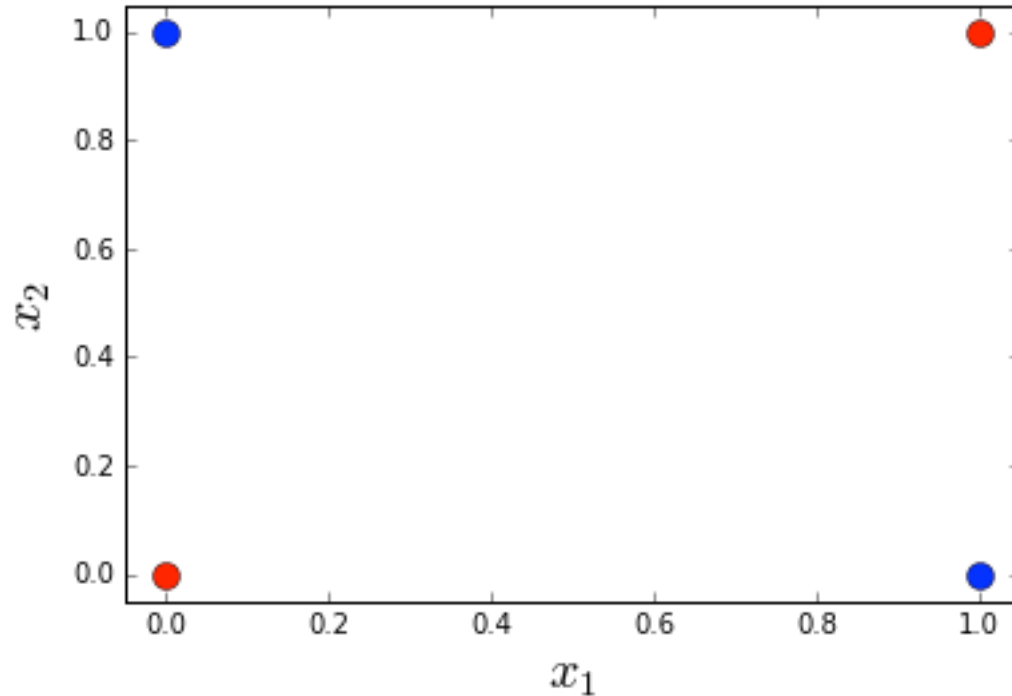
# SVM vs Neural Networks

- Hand-written digit recognition: MNIST data



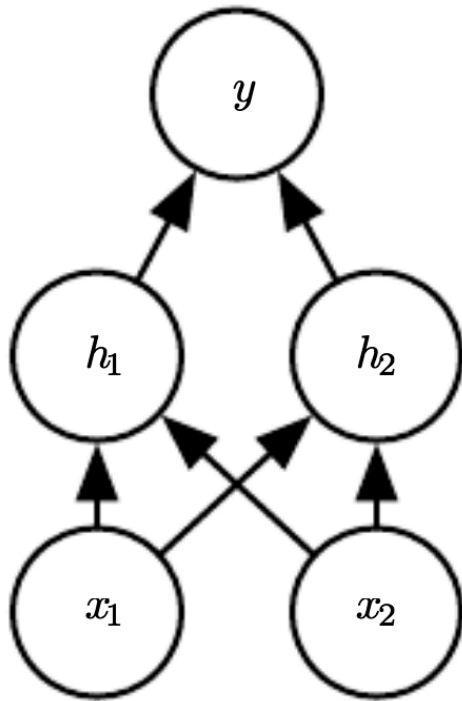
See illustration in notebook

# Example: Learning XOR



- Optimal linear model (sq. loss)
  - Predicts 0.5 on all points

# Example: Learning XOR



$$h_1 = \sigma(w_1^T x + c_1)$$

$$h_2 = \sigma(w_2^T x + c_2)$$

$$y = \sigma(w^T h + b)$$

where,

$$\sigma(z) = \max\{0, z\}$$

See illustration in notebook

# Design Choices

- Cost function
- Output units
- Hidden units
- Architecture
- Optimizer



# Cost Function

- Cross-entropy between training data and model distribution (i.e. **negative log-likelihood**)

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

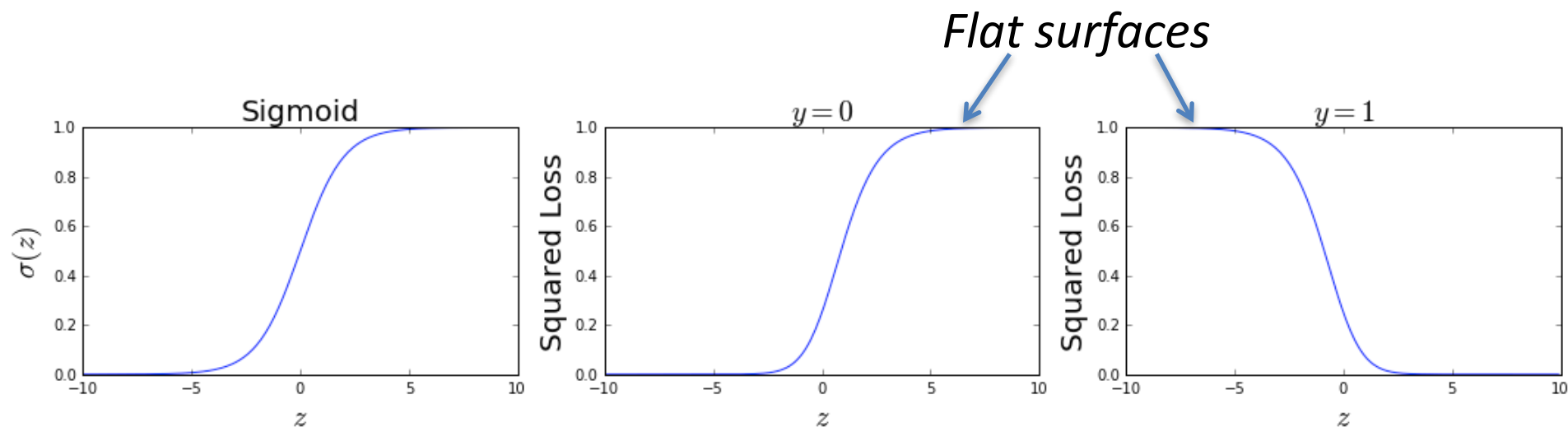
- Do not need to design separate cost functions
- Gradient of cost function must be large enough

# Cost Function

- Example: sigmoid output + squared loss

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

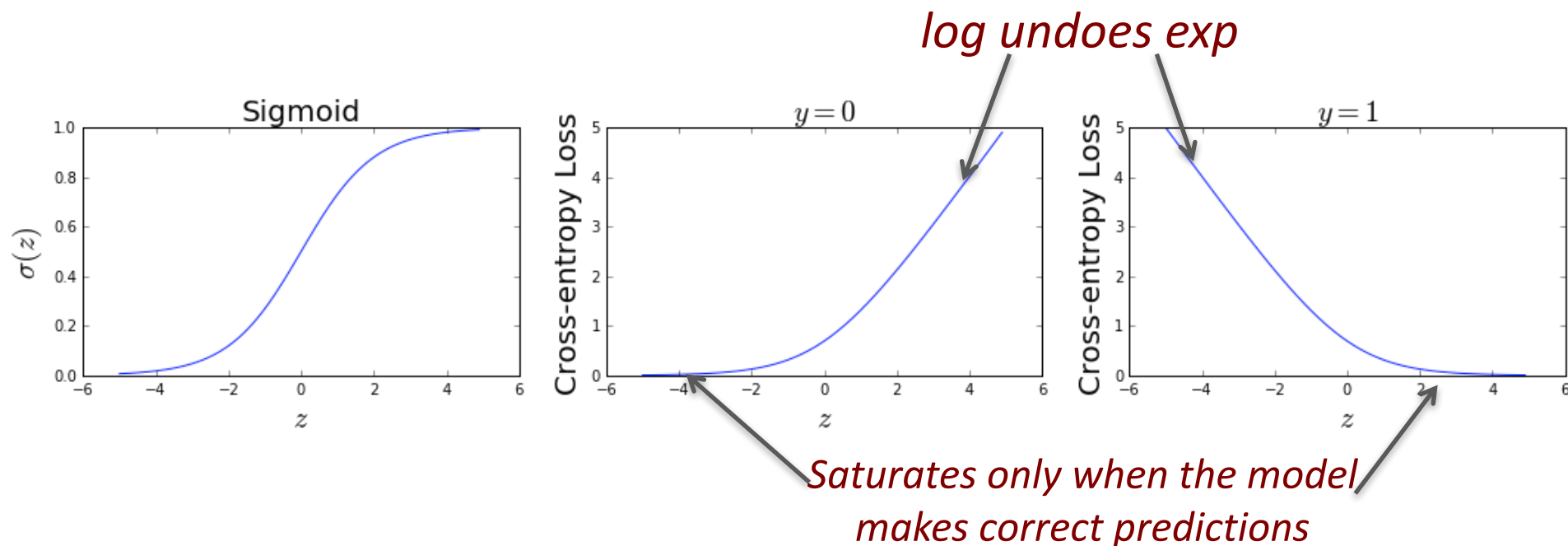
$$L_{sq}(y, z) = (y - \sigma(z))^2$$



# Cost Function

- Example: sigmoid output + cross-entropy loss

$$L_{ce}(y, z) = -(y \log(z) + (1 - y) \log(1 - z))$$



# Design Choices

- Cost function
- Output units
- Hidden units
- Architecture
- Optimizer

# Output Units

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

# Softmax Output

- Discrete / Multinoulli output distribution
- For output scores  $z_1, \dots, z_n$

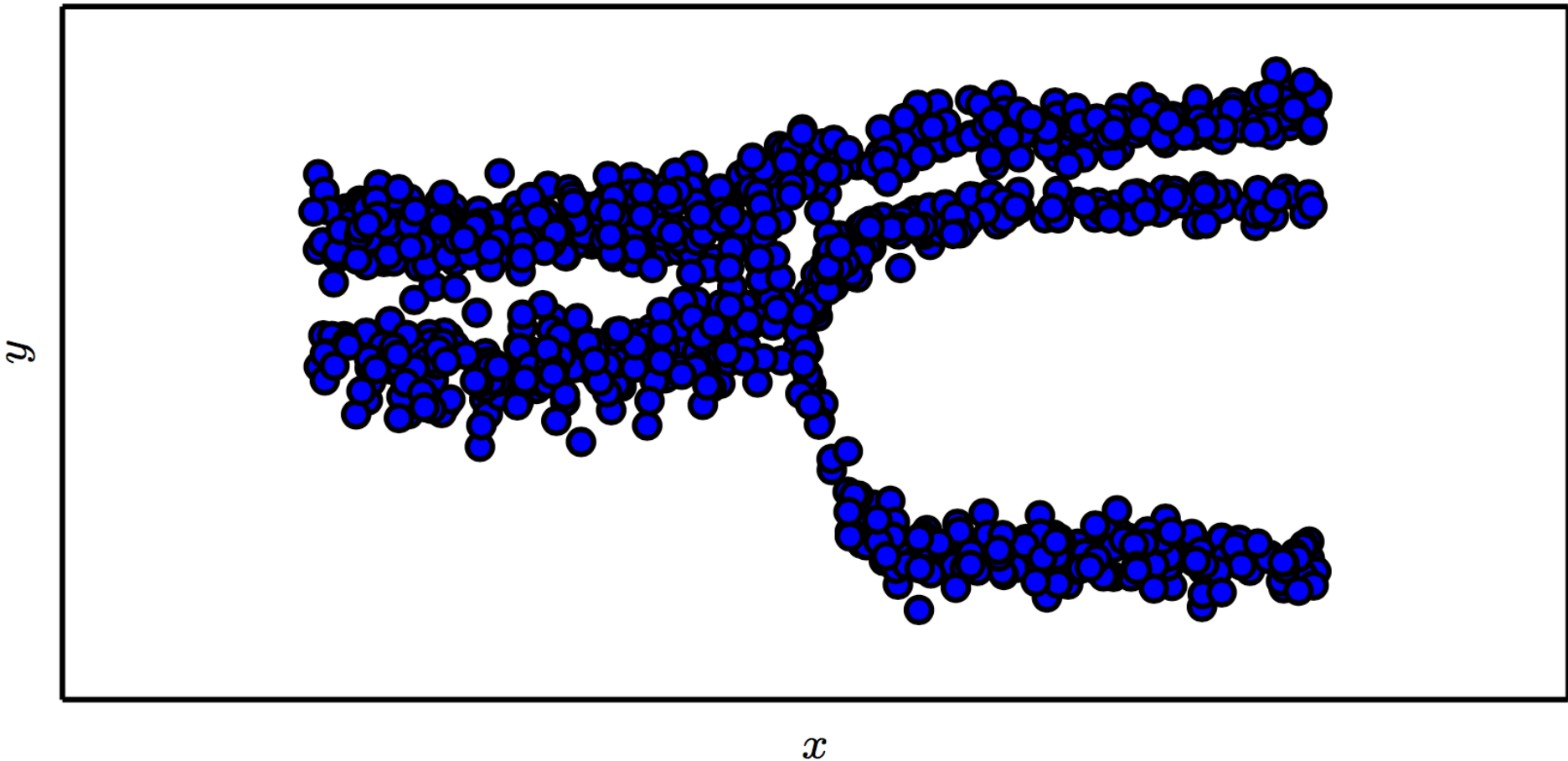
$$\text{softmax}(z)_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$

- Log-likelihood undoes exp

$$\begin{aligned}\log \text{softmax}(z)_i &= z_i - \log \sum_j \exp(z_j) \\ &\approx z_i - \max_j z_j\end{aligned}$$

(Score to target label – Maximum score)

# Mixture Density Output



(Goodfellow 2017)

# Design Choices

- Cost function
- Output units
- Hidden units
- Architecture
- Optimizer



# Hidden Units

$$\mathbf{h} = g(\mathbf{W}^T x + \mathbf{b})$$

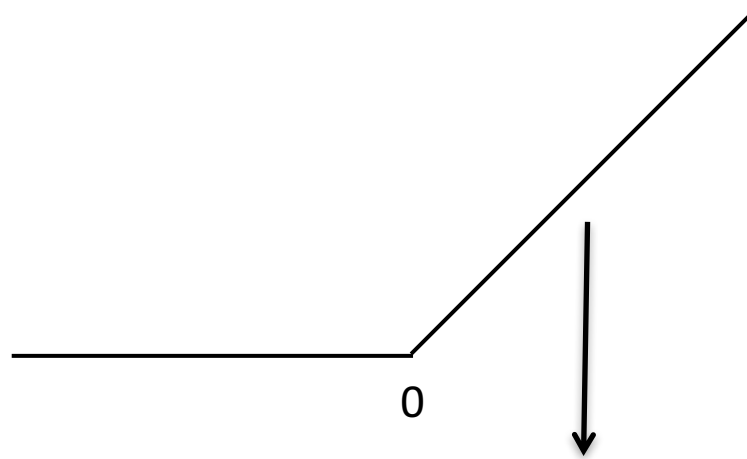
with activation function  $g$

- Ensure **gradients remain large** through hidden unit
- Preferred: **piece-wise linear activation**
- **Avoid sigmoid**/tanh activation
  - Do not provide useful gradient info when they saturate

# ReLU

- Rectified Linear Units

$$g(z) = \max\{0, z\}$$



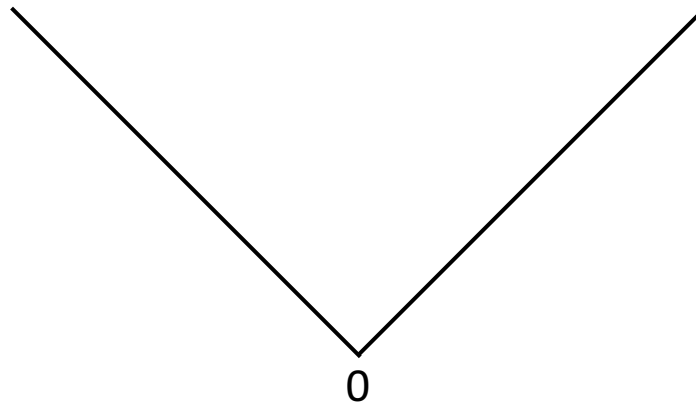
- Gradient is 1 whenever unit is active
  - More useful for learning compared to sigmoid
  - No useful gradient information when  $z < 0$

# Generalized ReLU

- Generalization: For  $\alpha_i > 0$ ,

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

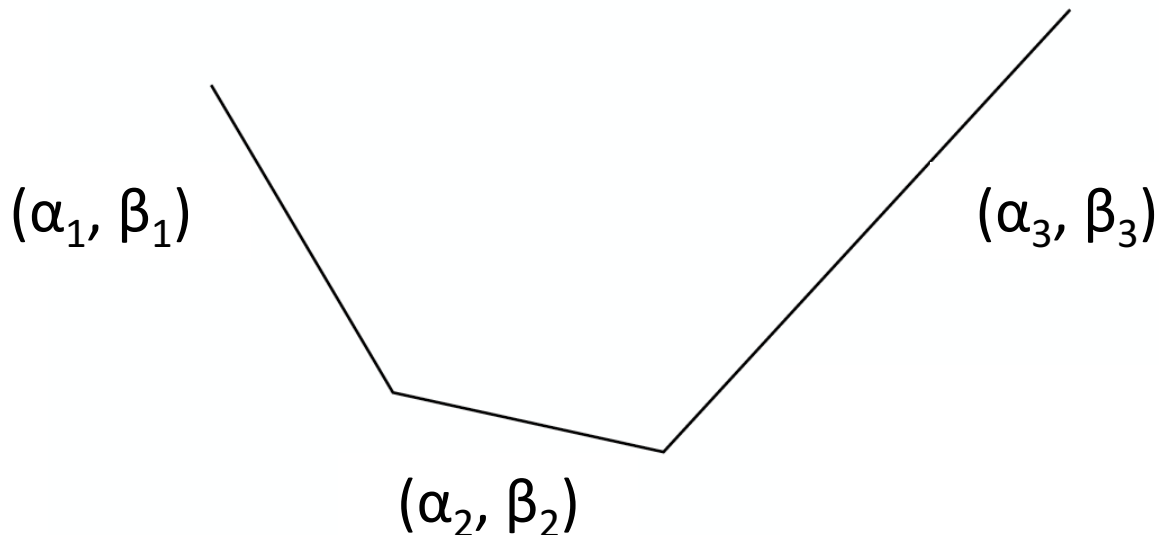
- E.g. Absolute value ReLU:  $\alpha_i = -1 \Rightarrow g(z) = |z|$



# Maxout

- Directly learn the activation function
  - Max of  $k$  linear functions

$$g(z) = \max_{i \in \{1, \dots, k\}} \alpha_i z_i + \beta_i$$

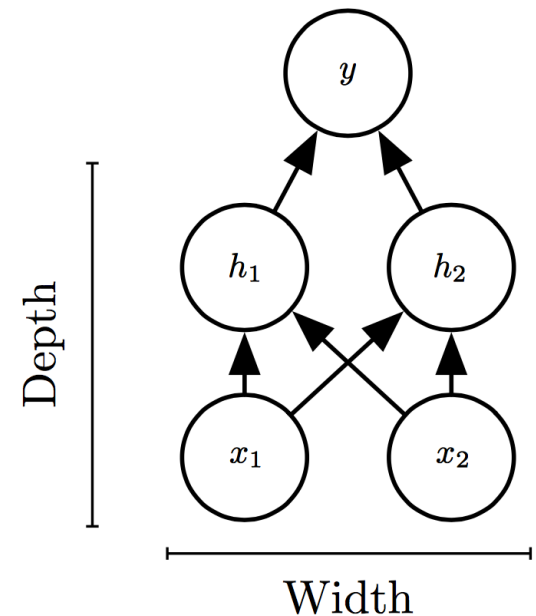


# Design Choices

- Cost function
- Output units
- Hidden units
- Architecture
- Optimizer

# Universal Approximation Theorem

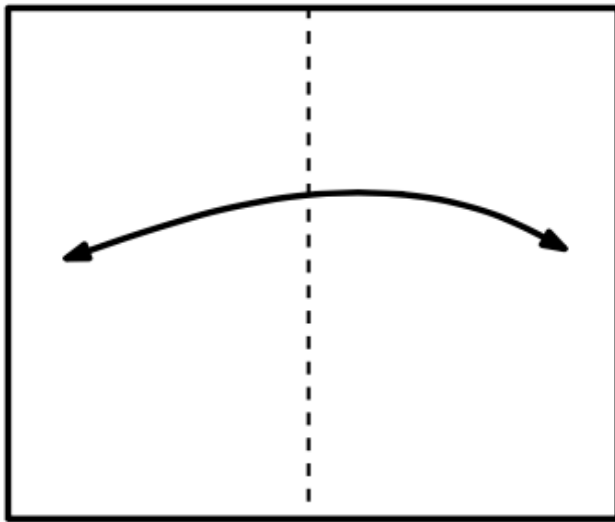
- One hidden layer is enough to *represent* an approximation of any function to an arbitrary degree of accuracy
- So why deeper?
  - Shallow net may need (exponentially) more width
  - Shallow net may overfit more



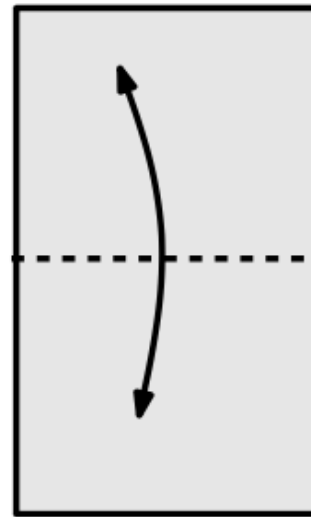
# Exponential Gain with Depth

- Each hidden layer **folds** the space of activations of the previous layer. E.g. abs activation  $g(z) = |z|$

Montúfar (2014)



1. Fold along the vertical axis



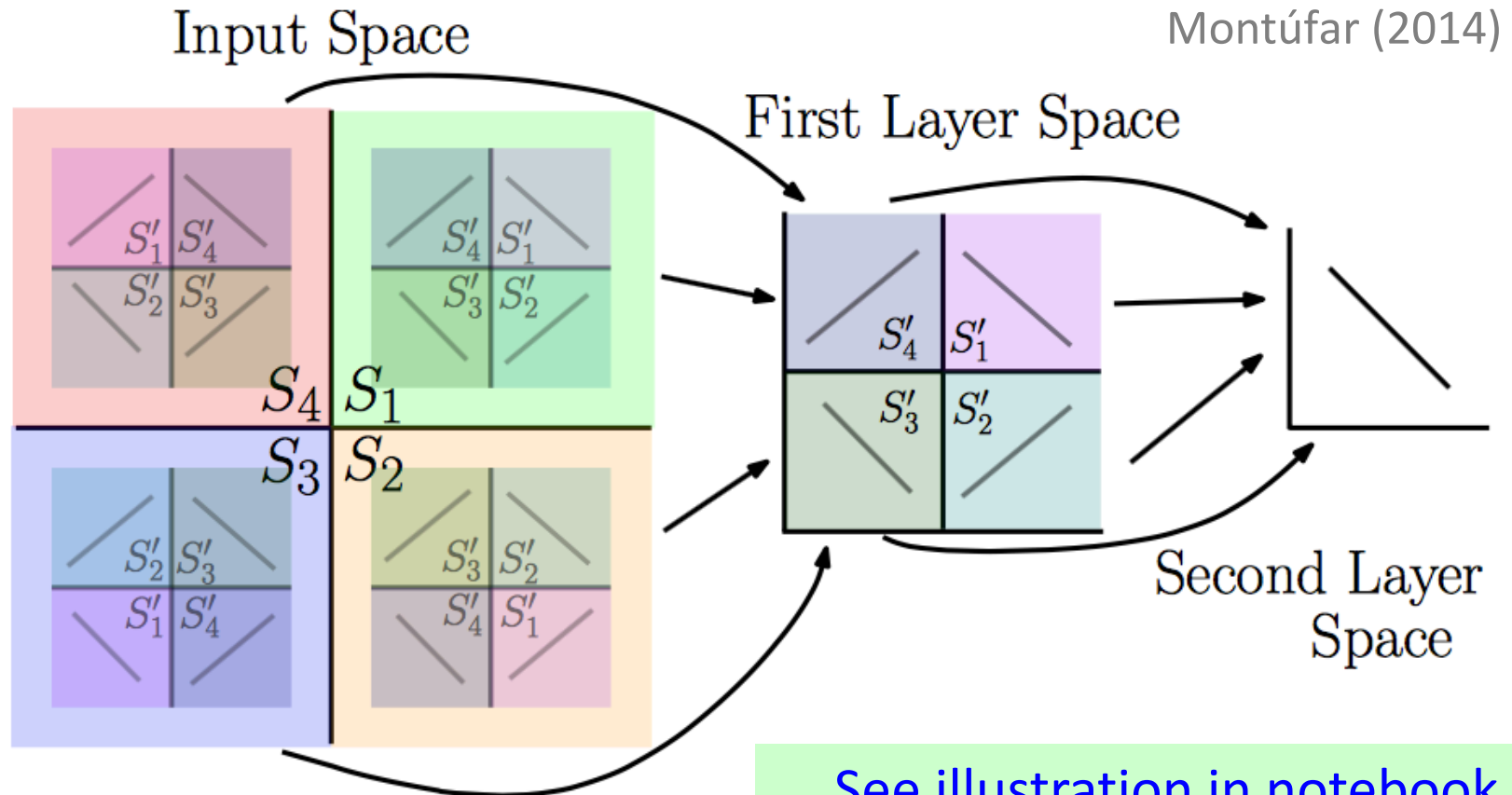
2. Fold along the horizontal axis



3.

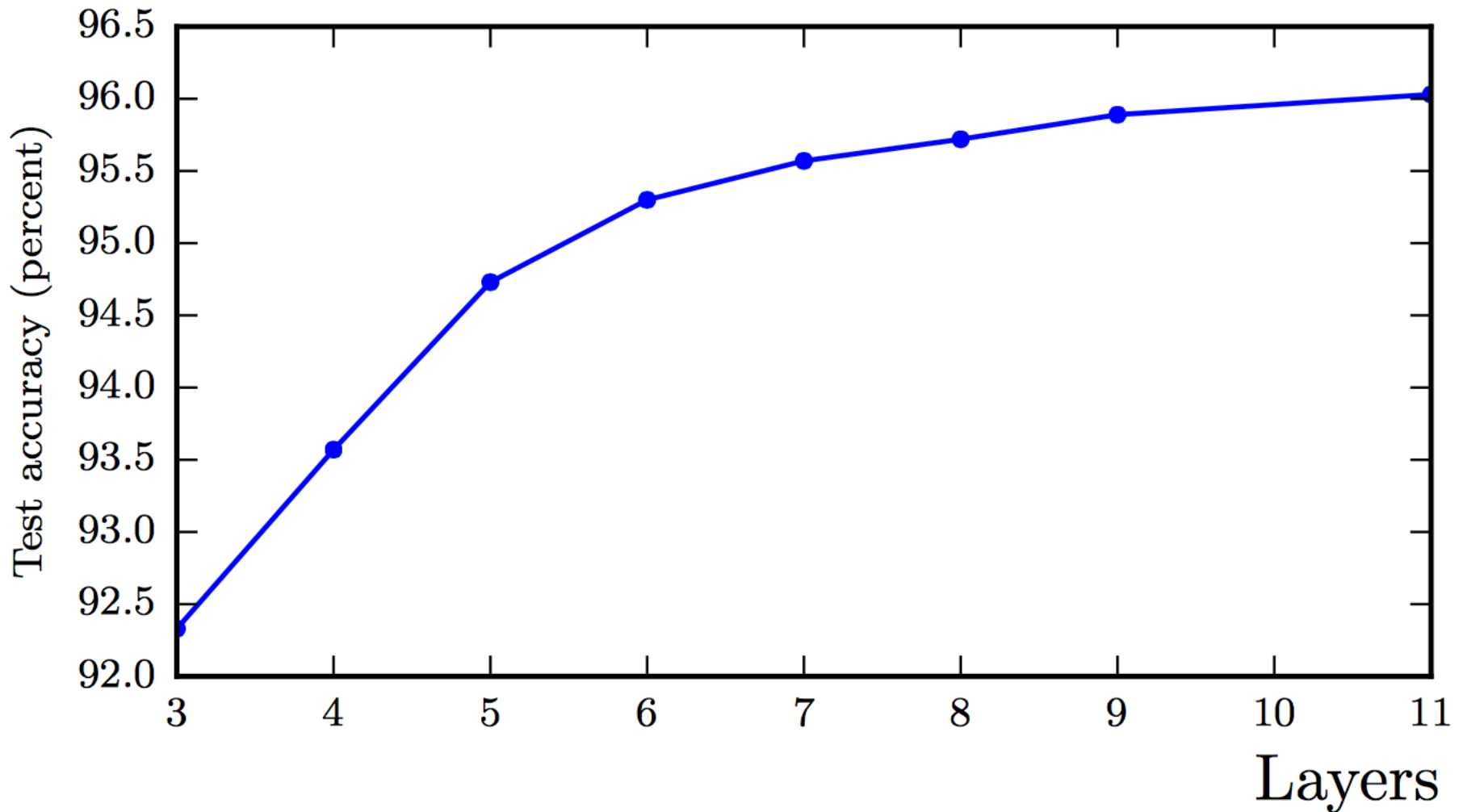
# Exponential Gain with Depth

- With  $N$  hidden layers, there are  $O(4^N)$  piecewise linear regions



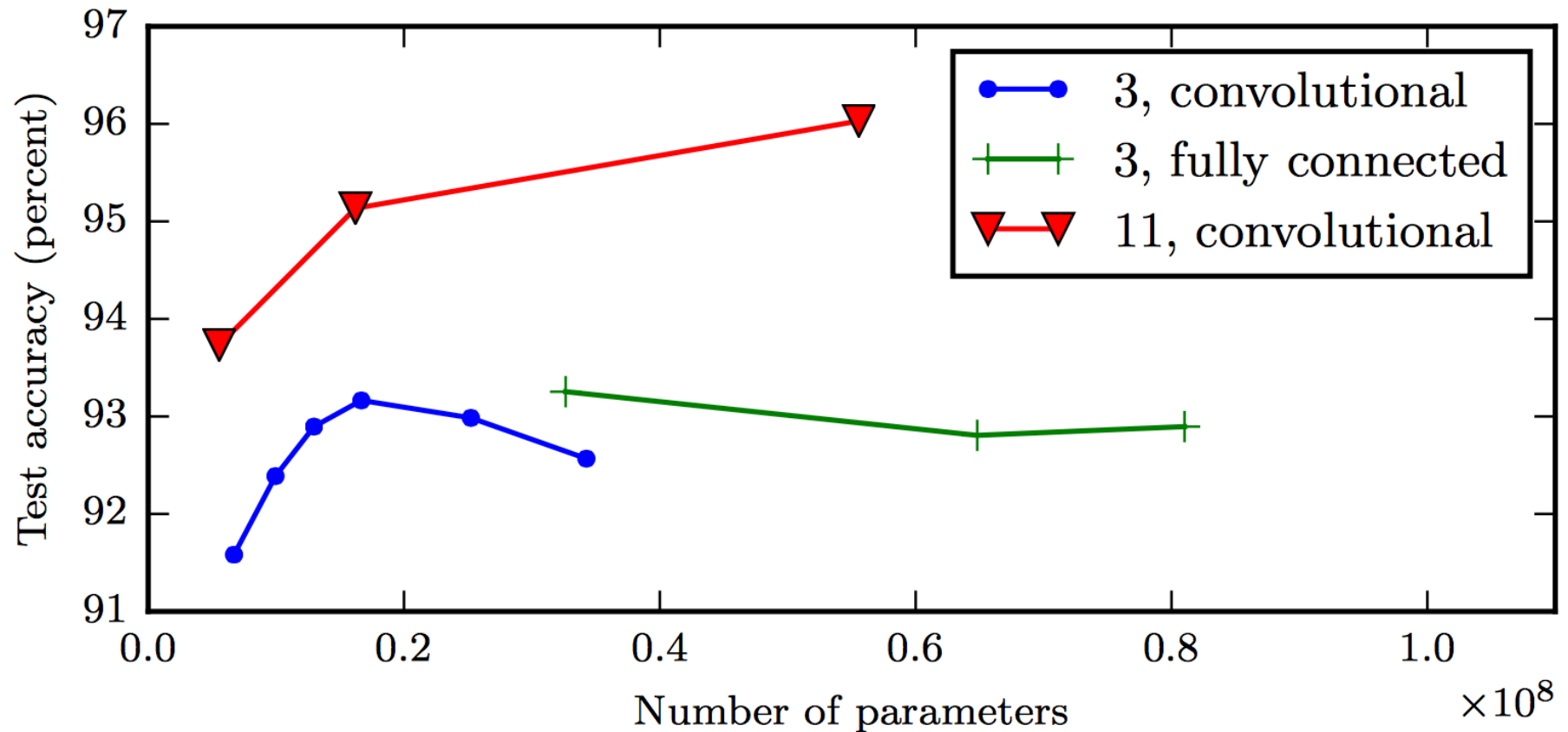


# Better Generalization with Depth



(Goodfellow 2017)

# Large, Shallow Nets Overfit More



(Goodfellow 2017)

# Design Choices

- Cost function
- Output units
- Hidden units
- Architecture
- Optimizer

# Gradient-based Optimizer

(e.g. stochastic gradient descent)

- “Chain rule” for computing gradients:

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

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

- For deeper networks

$$\frac{\partial z}{\partial x_i} = \sum_{j_1} \dots \sum_{j_m} \frac{\partial z}{\partial y_{j_1}} \dots \frac{\partial y_{j_m}}{\partial x_i}$$

Naïve computation  
takes exponential time

# Backpropagation

- Avoids repeated sub-expressions
- Uses dynamic programming (table filling)
- Trades-off memory for speed

# Backprop: Arithmetic

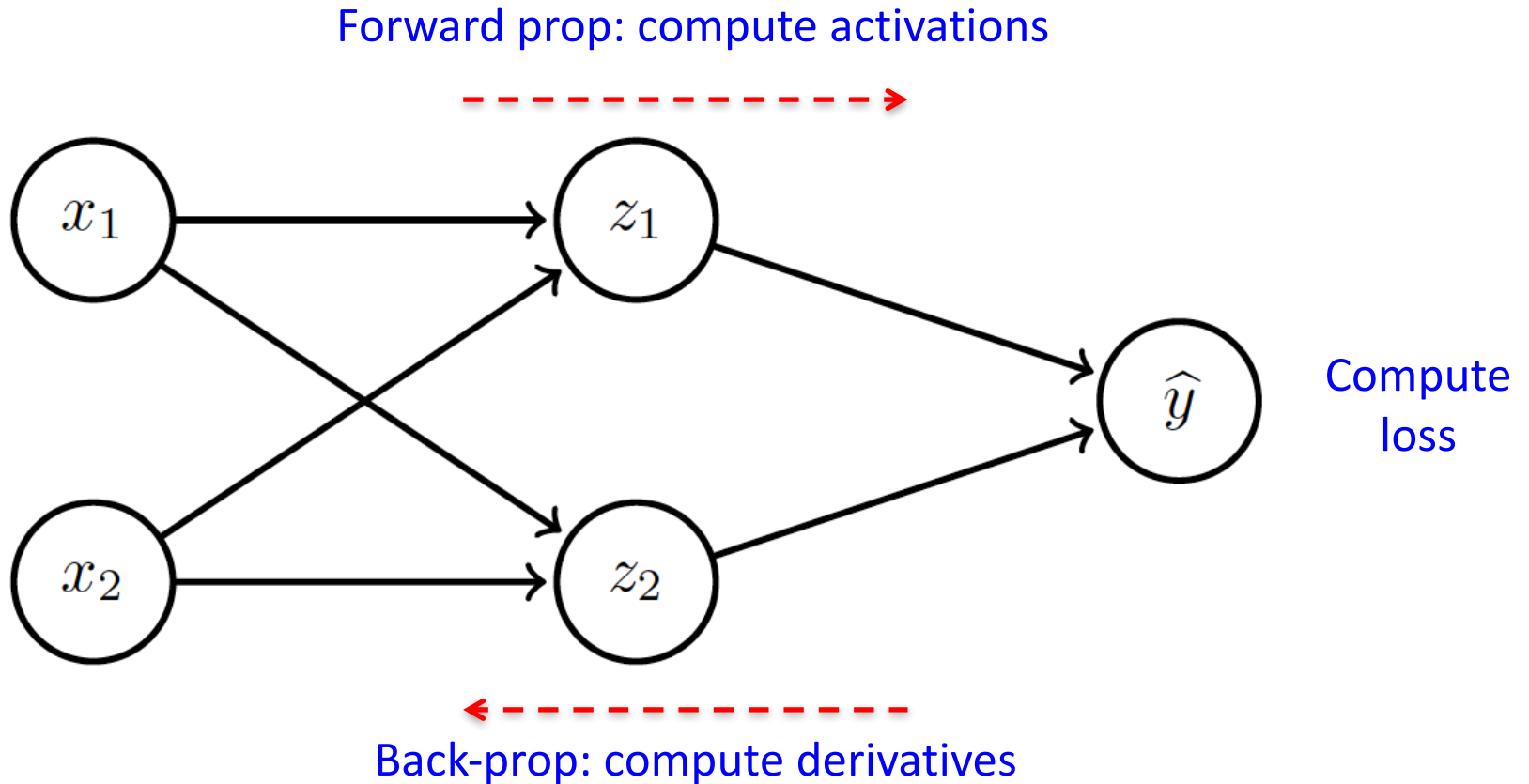
- Jacobian-gradient products

$$\begin{array}{l} \mathbf{z} = g(\mathbf{x}) \\ y = f(\mathbf{z}) \end{array} \quad \begin{array}{c} \left[ \begin{array}{c} \frac{\partial y}{\partial x_1} \\ \vdots \\ \frac{\partial y}{\partial x_m} \end{array} \right] \\ \text{grad w.r.t. } \mathbf{x} \end{array} = \begin{array}{c} \left[ \begin{array}{ccc} \frac{\partial z_1}{\partial x_1} & \dots & \frac{\partial z_m}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_1}{\partial x_n} & \dots & \frac{\partial z_m}{\partial x_n} \end{array} \right] \\ \text{Jacobian of 'g' } \end{array} \times \begin{array}{c} \left[ \begin{array}{c} \frac{\partial y}{\partial z_1} \\ \vdots \\ \frac{\partial y}{\partial z_m} \end{array} \right] \\ \text{grad w.r.t. } \mathbf{z} \end{array}$$

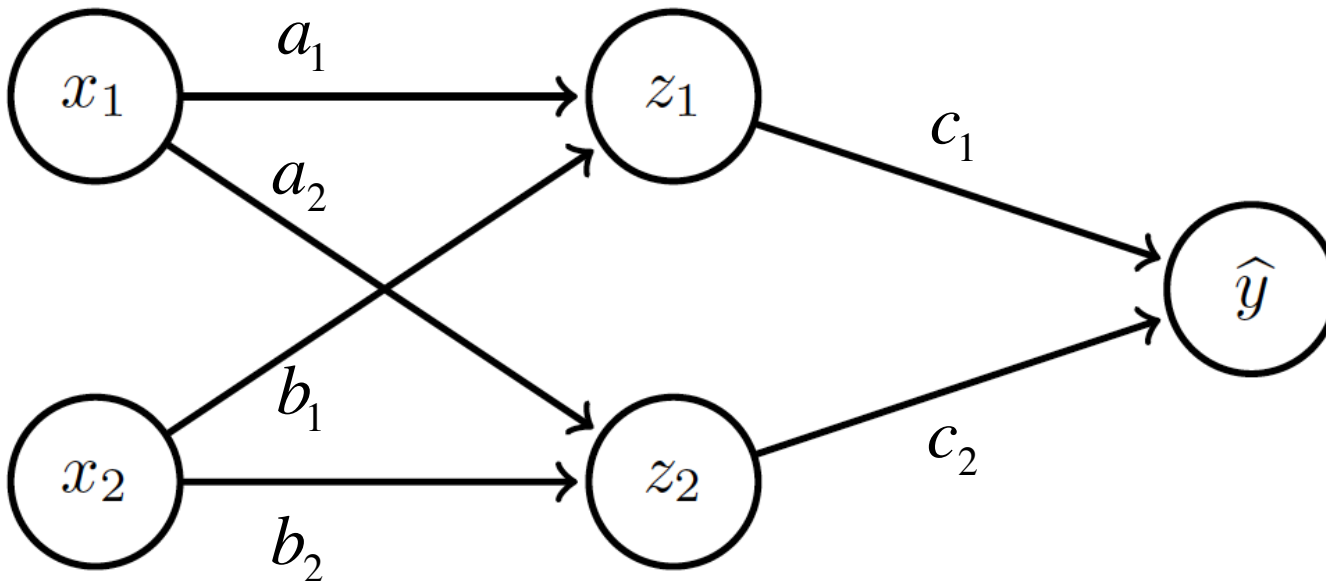
$$\nabla_{\mathbf{x}} y = \left( \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right)^T \nabla_{\mathbf{z}} y$$

Apply  
recursively!

# Backprop: Overview



# Backprop: Example



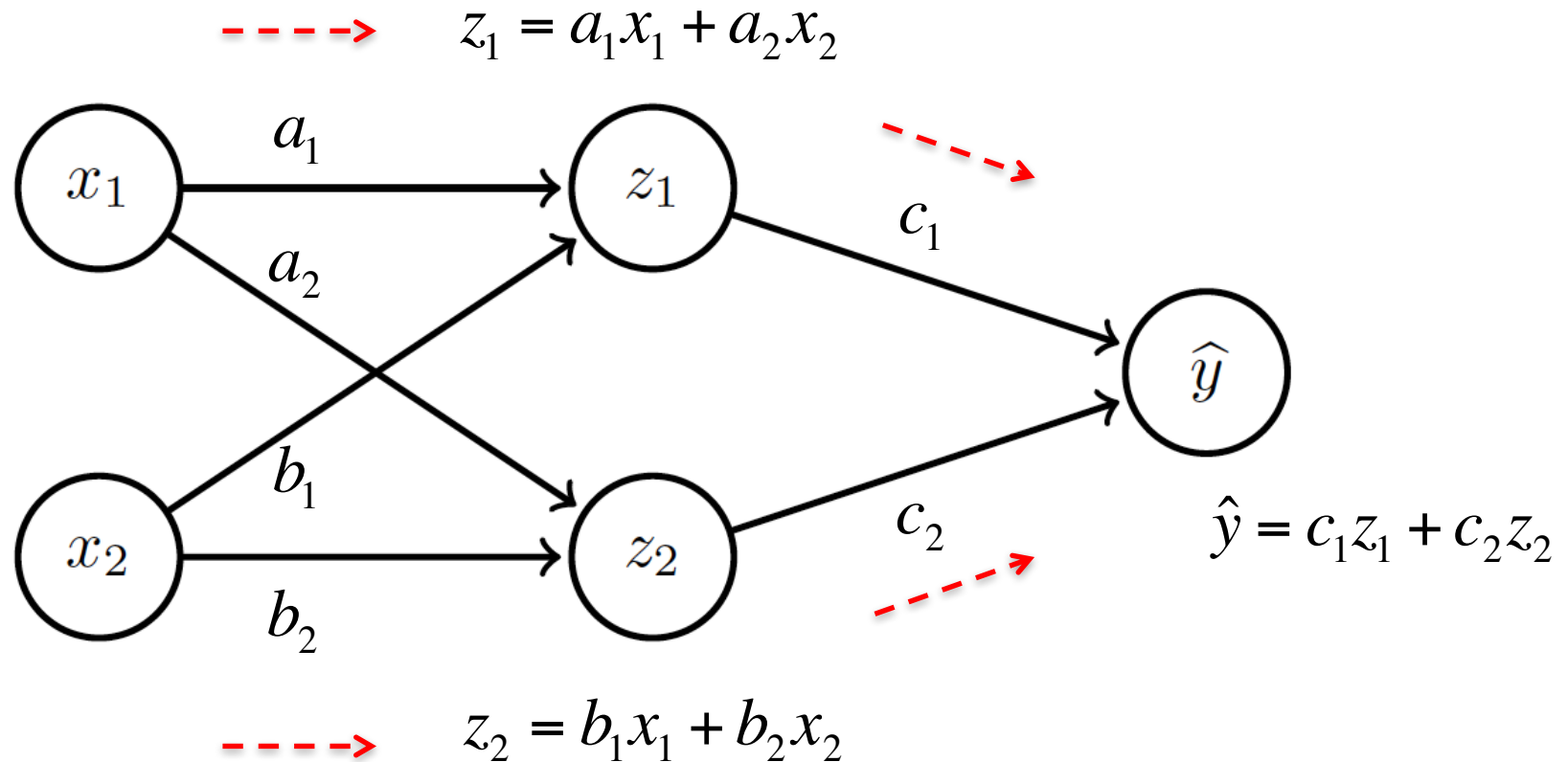
Linear activation functions

No bias

Squared loss

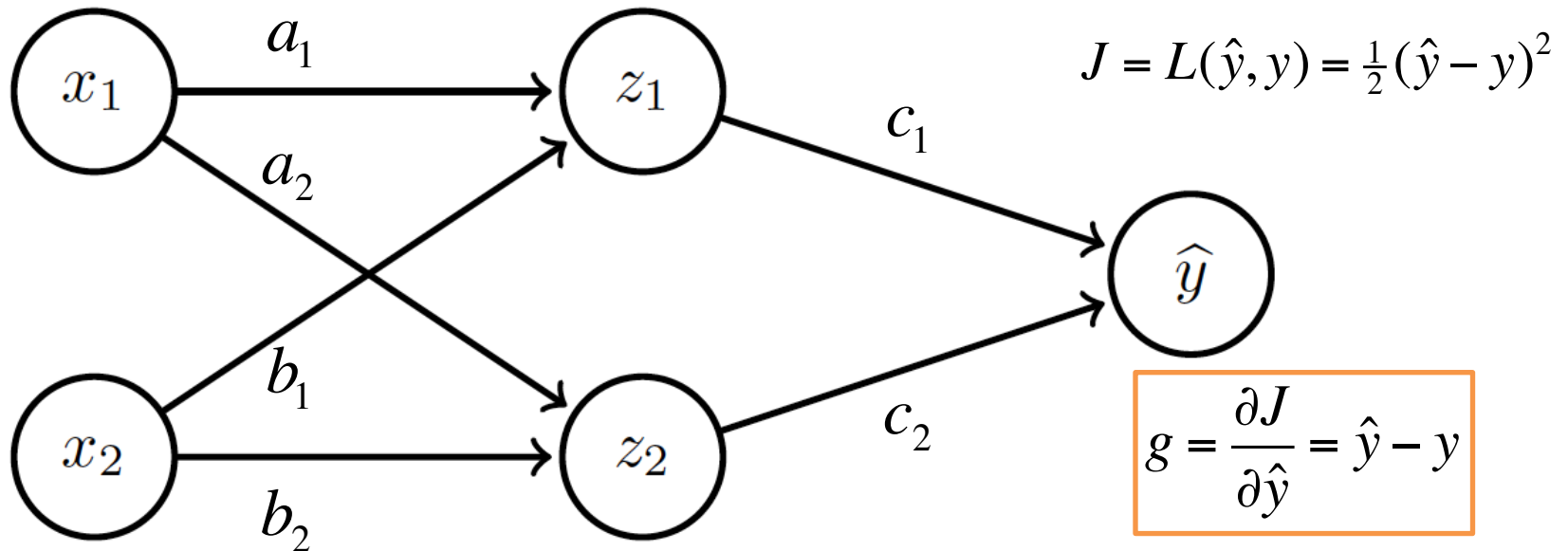


# Backprop: Example



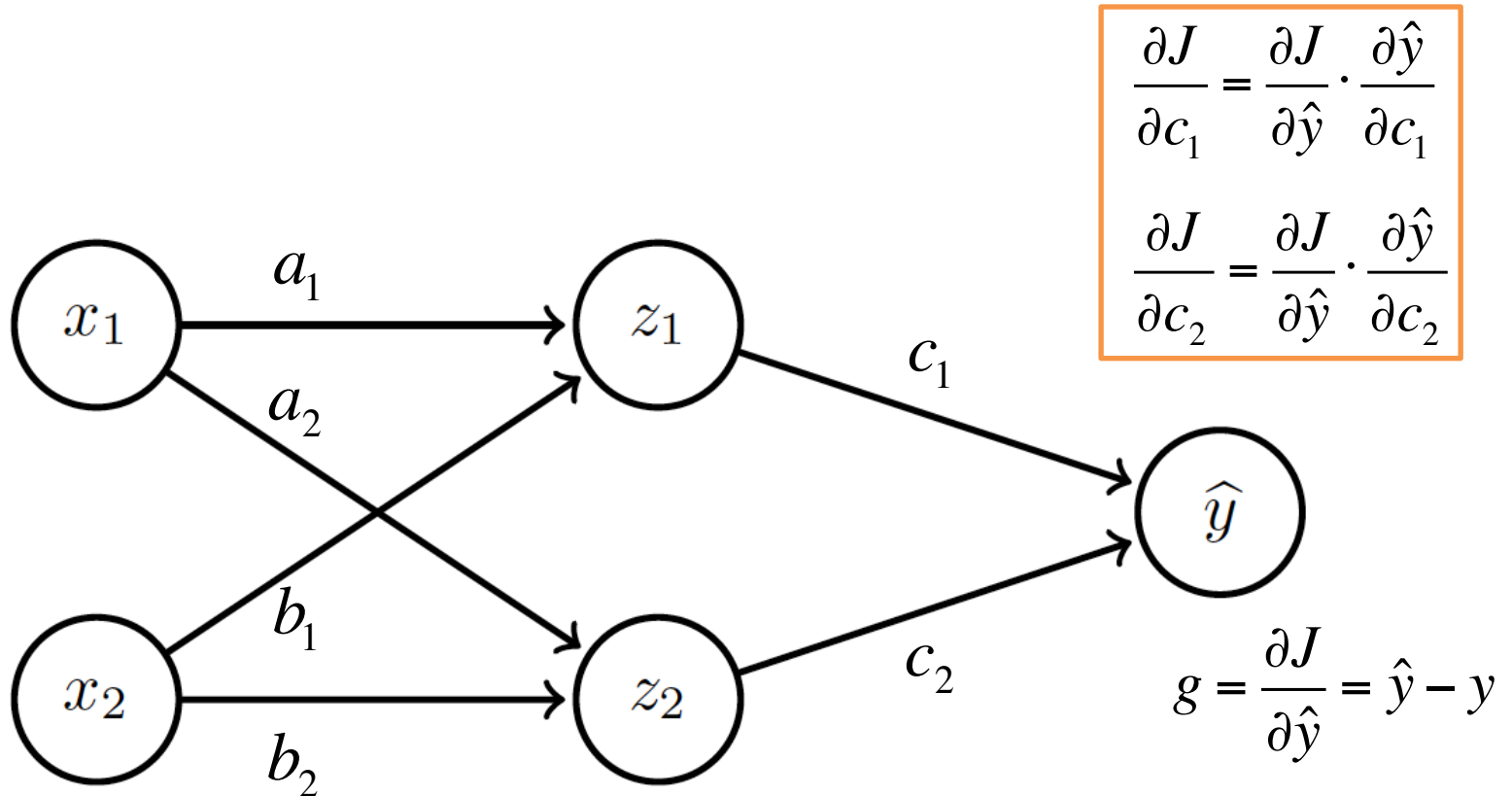
Forward prop: Propagate activations to output layer

# Backprop: Example



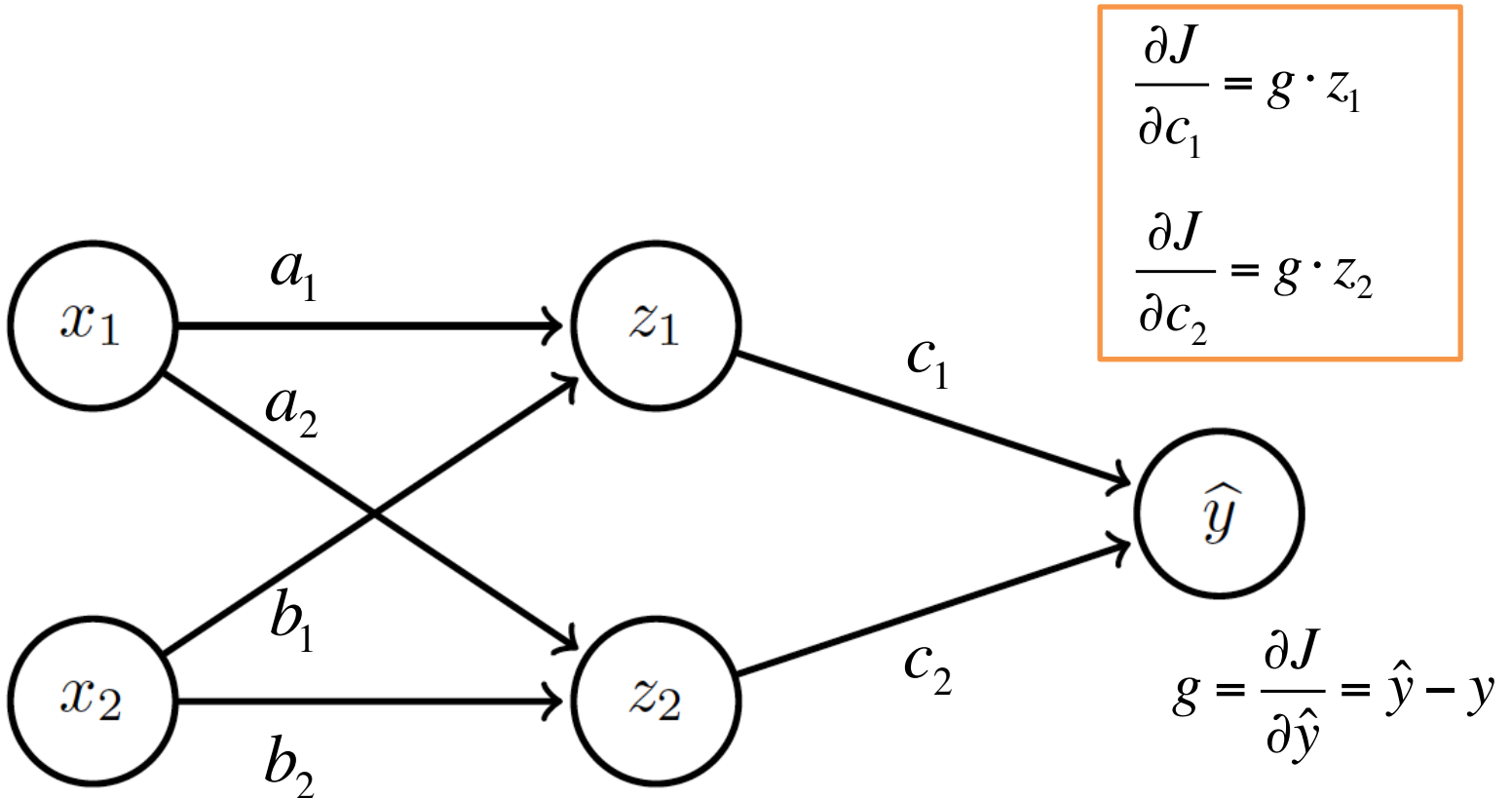
Backward prop: Compute loss and its derivative

# Backprop: Example



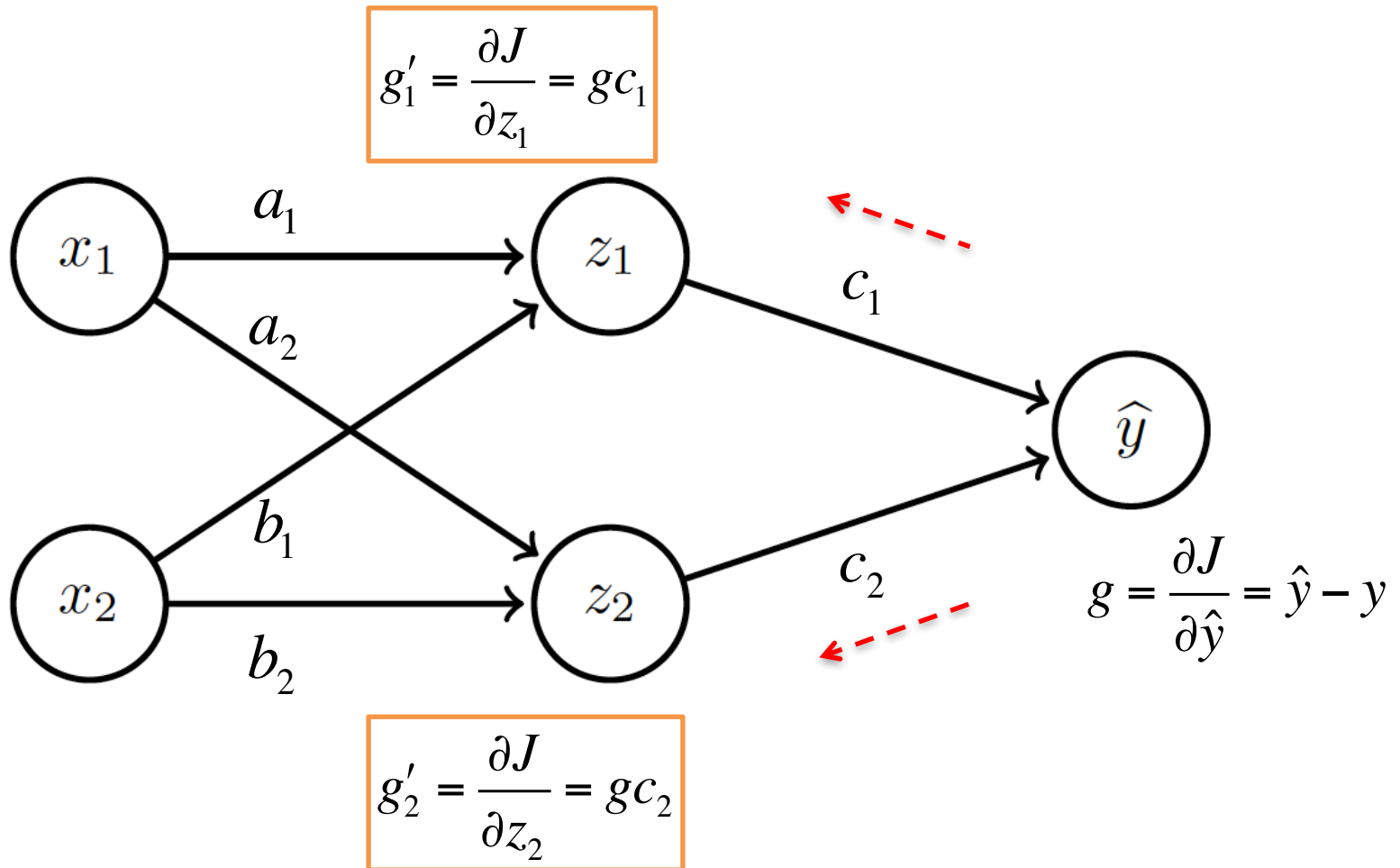
Backward prop: Compute derivatives w.r.t. weights  $c_1$  and  $c_2$

# Backprop: Example



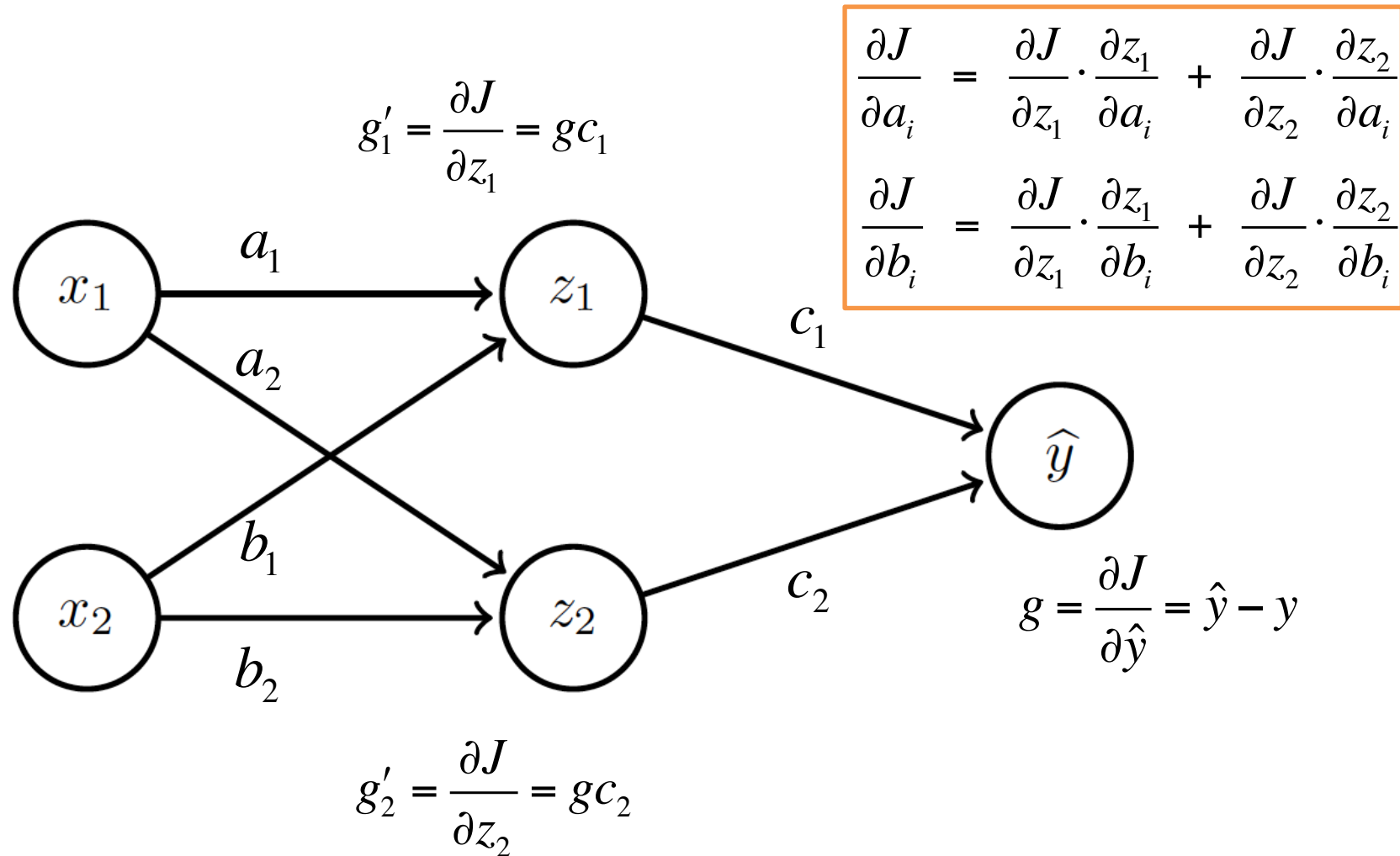
Backward prop: Compute derivatives w.r.t. weights  $c_1$  and  $c_2$

# Backprop: Example



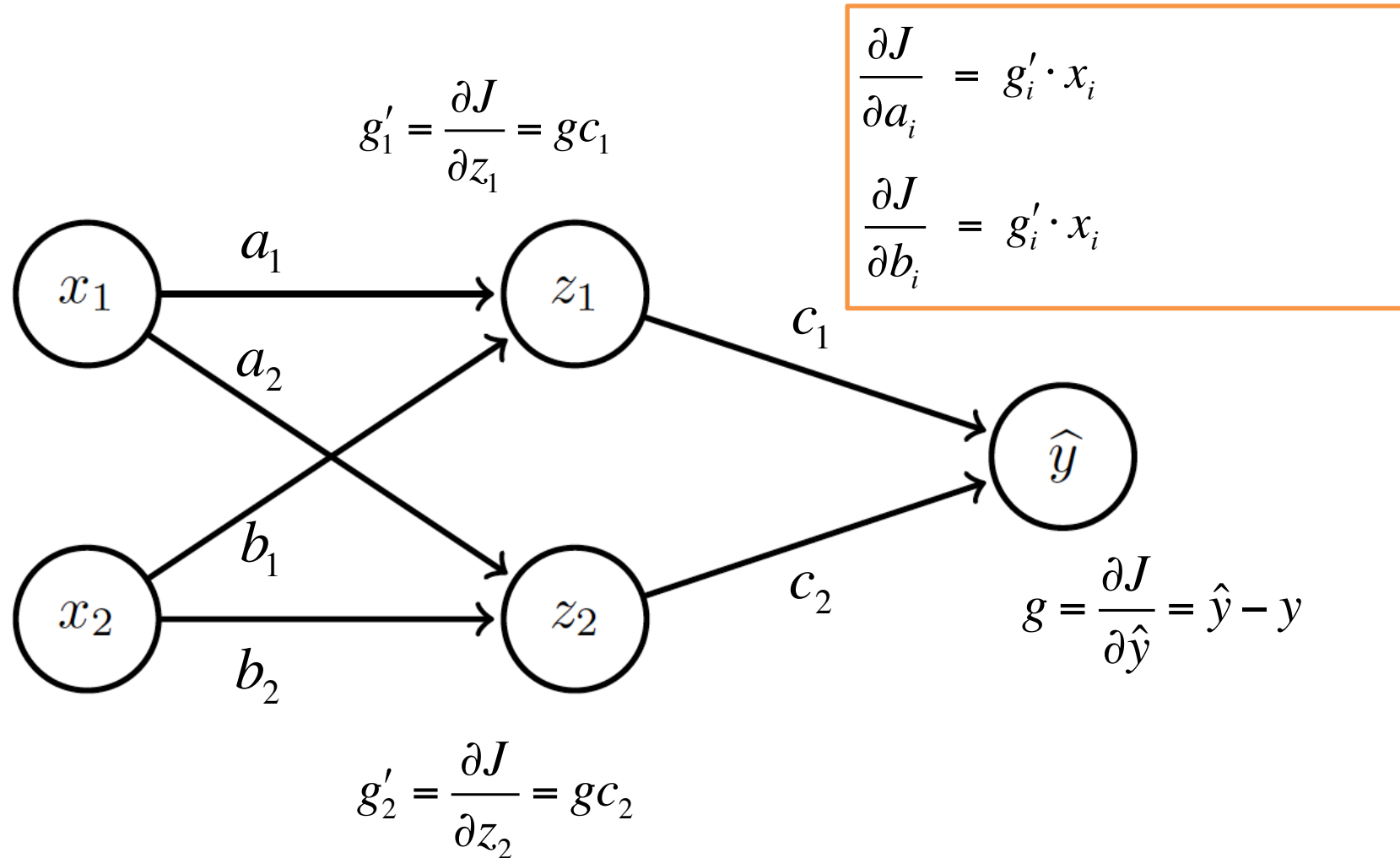
Backward prop: Propagate derivative to hidden layer

# Backprop: Example



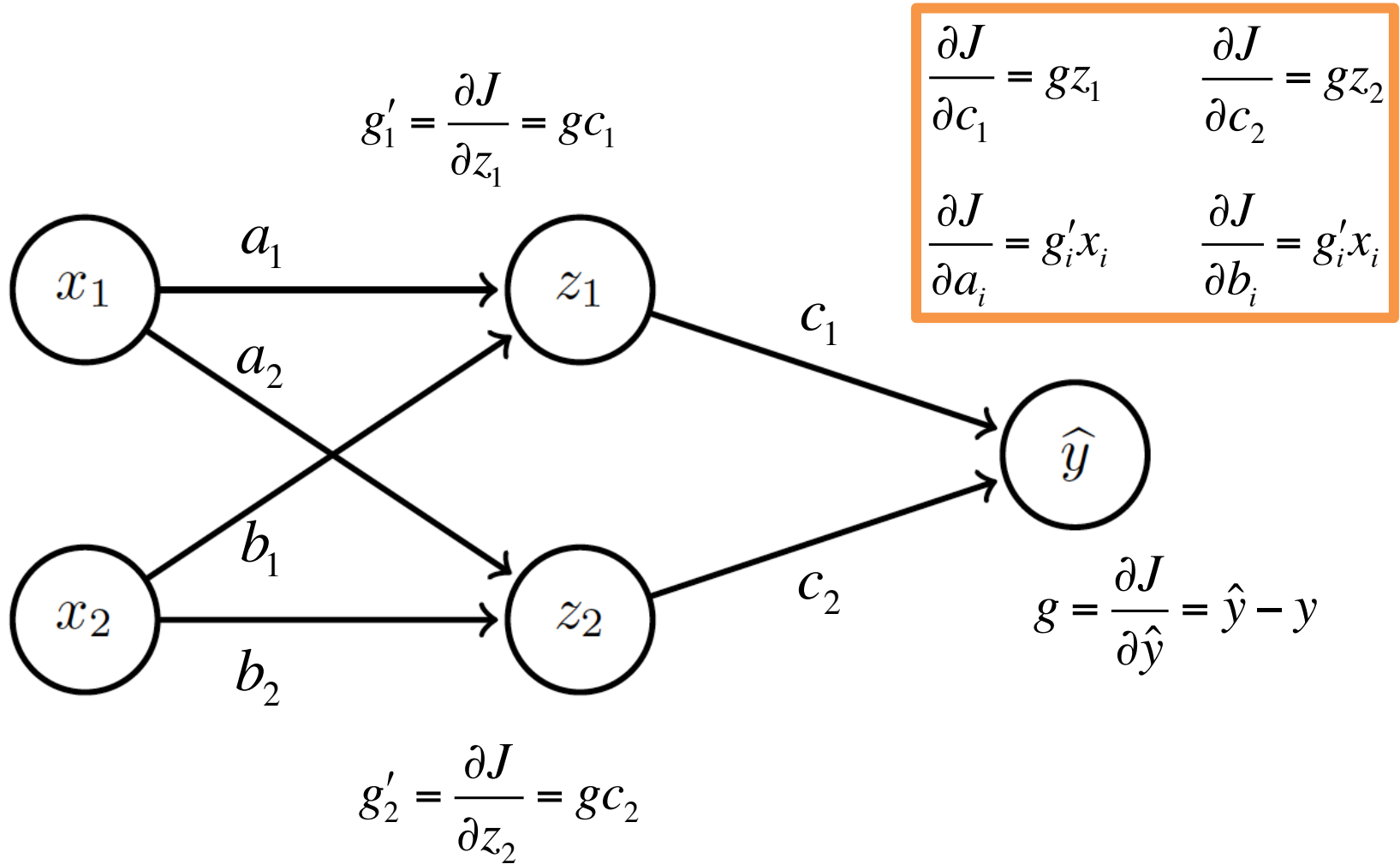
Backward prop: Compute derivatives w.r.t. weights  $a_1$ ,  $a_2$ ,  $b_1$  and  $b_2$

# Backprop: Example



Backward prop: Compute derivatives w.r.t. weights  $a_1, a_2, b_1$  and  $b_2$

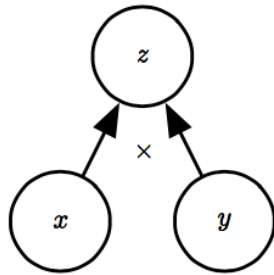
# Backprop: Example





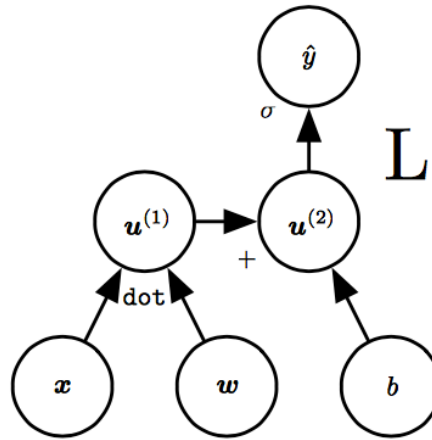
# Computation Graphs

Multiplication



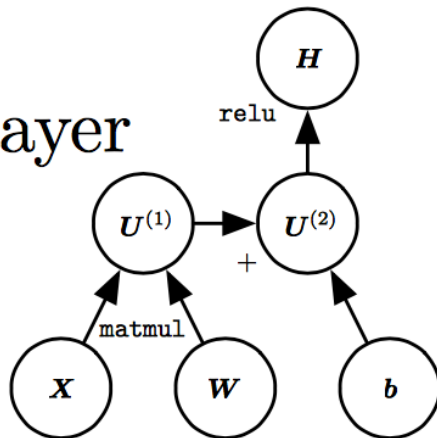
(a)

Logistic regression



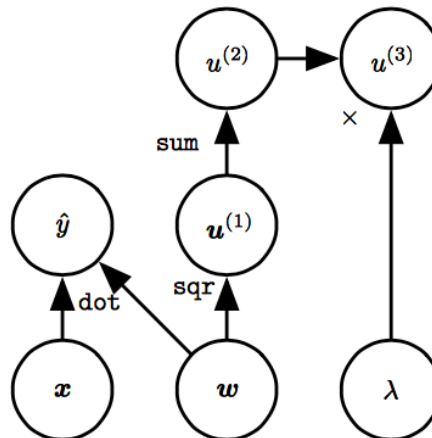
(b)

ReLU layer



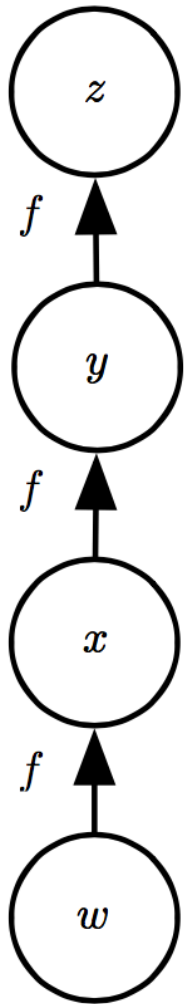
(c)

Linear regression  
and weight decay



(d)

# Repeated Sub-expressions



$$\begin{aligned} & \frac{\partial z}{\partial w} \\ &= \frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \frac{\partial x}{\partial w} \\ &= f'(y) f'(x) f'(w) \\ &= f'(f(f(w))) f'(f(w)) f'(w) \end{aligned}$$

Back-prop avoids computing this twice

# Backprop on Computation Graph

Maintain grad table

1: Initialize  $\mathbf{g} \in R^n$  where  $g_i$  denotes  $\frac{\partial u^n}{\partial u^i}$

2: for  $j = n - 1$  to 1 do:

3: 
$$g_j = \sum_{i: j \in Pa(u^i)} g_i \frac{\partial u^i}{\partial u^j}$$

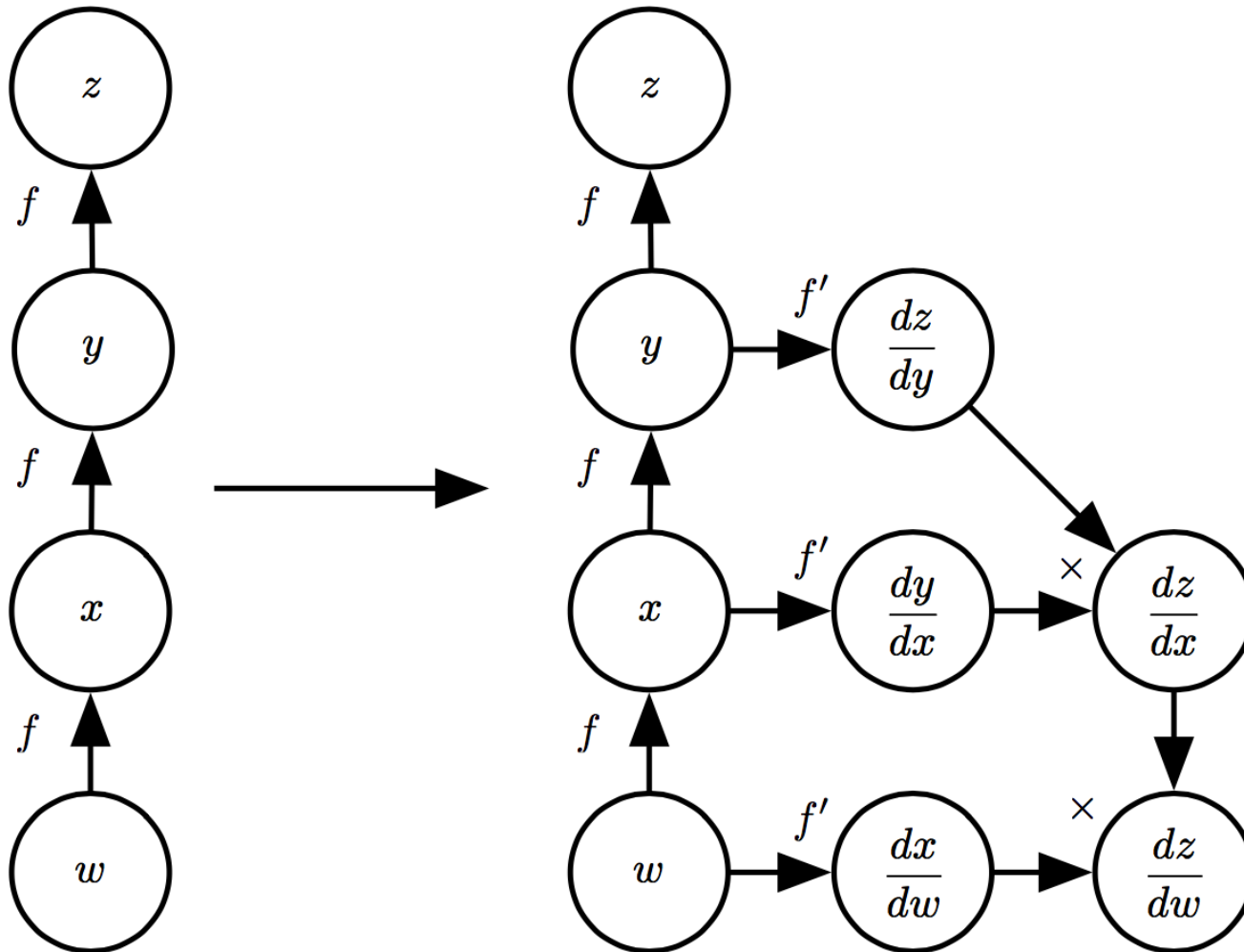
4: return  $\mathbf{g}$

Parents of  $u^i$

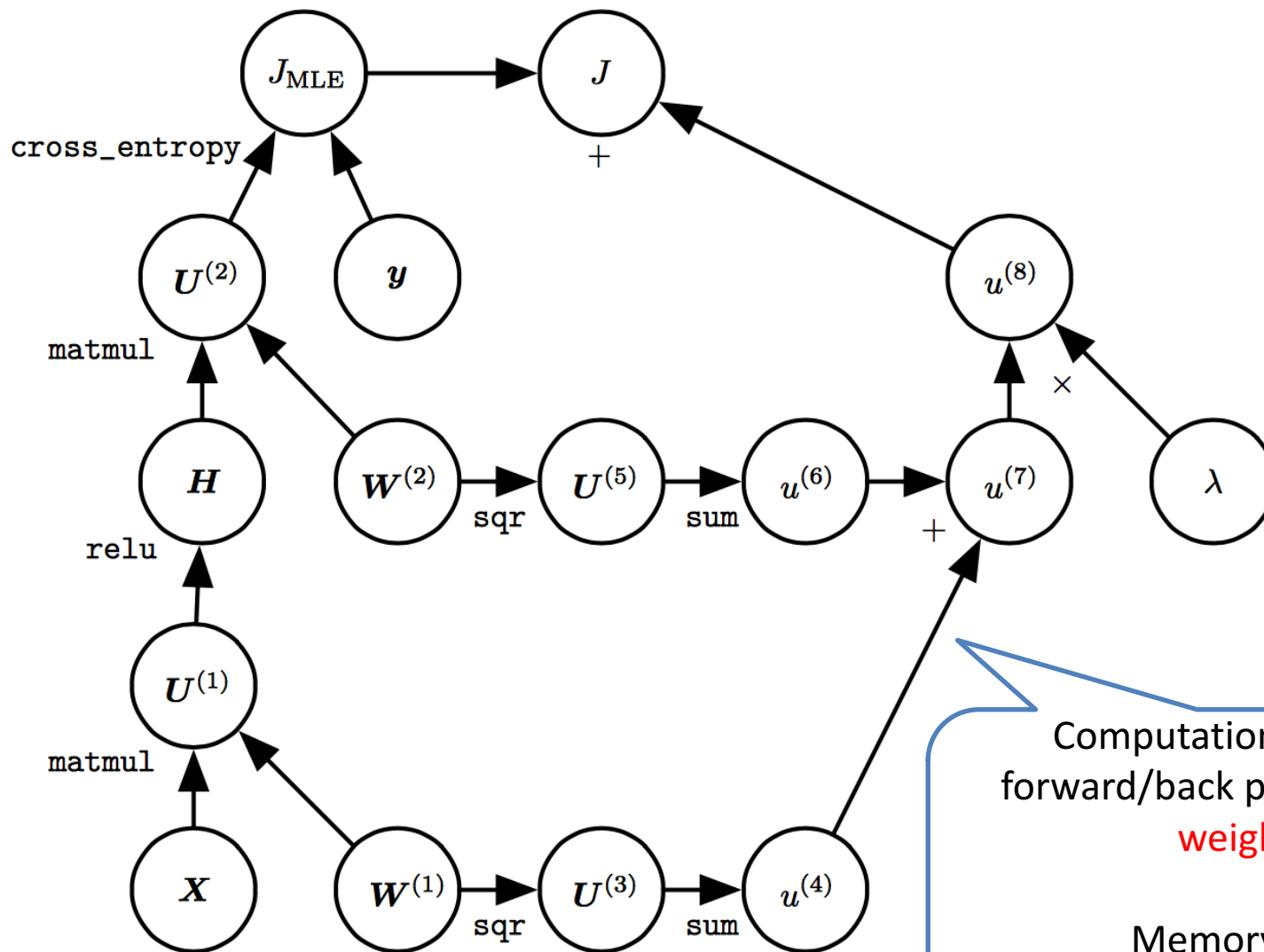
# Symbol-to-symbol Differentiation

- Derivatives as computation graphs
  - Same language for both forward and back-propagation
- During execution, replace symbolic inputs with numeric value
- Used by [Theano](#) and [TensorFlow](#)
- Symbol-to-number differentiation: e.g. Torch and Caffe

# Symbol-to-symbol Differentiation



# Training Feed-forward Nets



(Goodfellow et al. 2017)

Computational cost for  
forward/back prop:  $O(\text{\#num-weights})$

Memory cost:  
 $O(\text{\#num-layers} \times \text{minibatch-size})$