# Neural Networks

# +

# Backpropagation

Matt Gormley
Lecture 12
Feb. 24, 2023

# Reminders
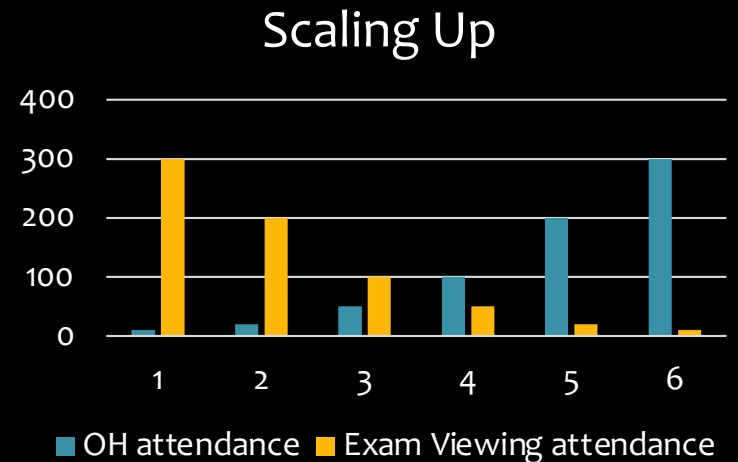
- **Post-Exam Followup:**
  - **Exam Viewing**
  - **Exit Poll: Exam 1**
  - **Grade Summary 1**
- **Homework 4: Logistic Regression**
  - **Out: Fri, Feb 17**
  - **Due: Sun, Feb 26 at 11:59pm**
- **Homework 5: Neural Networks**
  - **Out: Sun, Feb 26**
  - **Due: Fri, Mar 17 at 11:59pm**

Scaling Up



- OH attendance
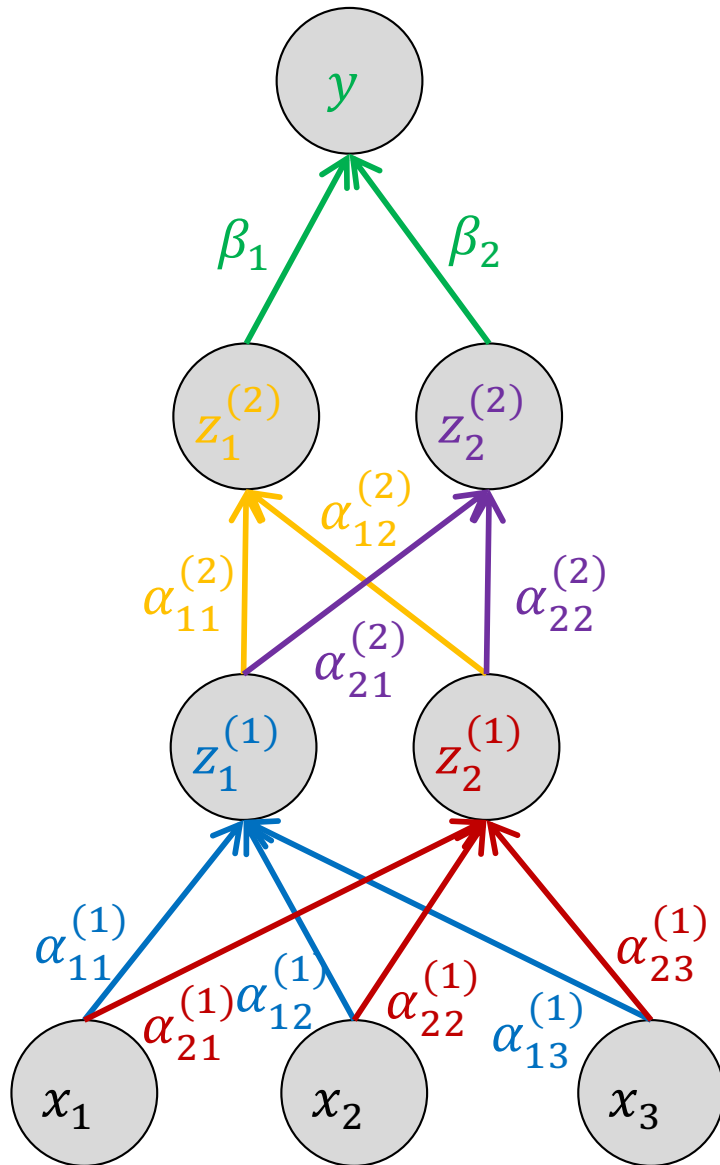- Exam Viewing attendance

# ARCHITECTURES

# Neural Network Architectures

Even for a basic Neural Network, there are many design decisions to make:

1. # of hidden layers (depth)

2. # of units per hidden layer (width)

3. Type of activation function (nonlinearity)

4. Form of objective function

5. How to initialize the parameters

# Neural Network



*Example: Neural Network with 2 Hidden Layers and 2 Hidden Units*

$$z_1^{(1)} = \sigma(\alpha_{11}^{(1)}x_1 + \alpha_{12}^{(1)}x_2 + \alpha_{13}^{(1)}x_3 + \alpha_{10}^{(1)})$$

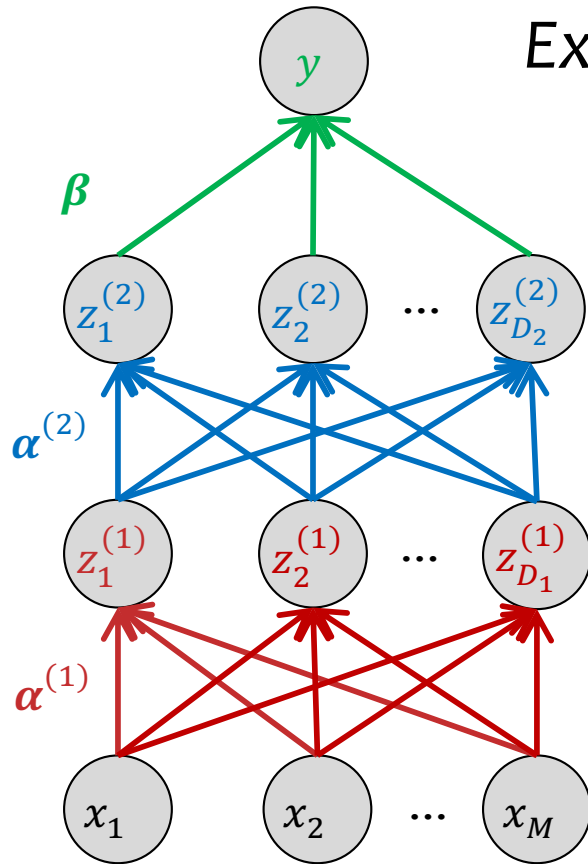$$z_2^{(1)} = \sigma(\alpha_{21}^{(1)}x_1 + \alpha_{22}^{(1)}x_2 + \alpha_{23}^{(1)}x_3 + \alpha_{20}^{(1)})$$

$$z_1^{(2)} = \sigma(\alpha_{11}^{(2)}z_1^{(1)} + \alpha_{12}^{(2)}z_2^{(1)} + \alpha_{10}^{(2)})$$

$$z_2^{(2)} = \sigma(\alpha_{21}^{(2)}z_1^{(1)} + \alpha_{22}^{(2)}z_2^{(1)} + \alpha_{20}^{(2)})$$

$$y = \sigma(\beta_1 z_1^{(2)} + \beta_2 z_2^{(2)} + \beta_0)$$

# Neural Network (Matrix Form)

*Example: Arbitrary Feed-forward Neural Network*



$$\boldsymbol{\beta} \in \mathbb{R}^{D_2}$$

$$\beta_0 \in \mathbb{R}$$

$$\boldsymbol{\alpha}^{(2)} \in \mathbb{R}^{M \times D_2}$$

$$\boldsymbol{b}^{(2)} \in \mathbb{R}^{D_2}$$

$$\boldsymbol{\alpha}^{(1)} \in \mathbb{R}^{M \times D_1}$$

$$\boldsymbol{b}^{(1)} \in \mathbb{R}^{D_1}$$

$$y = \sigma((\boldsymbol{\beta})^T \boldsymbol{z}^{(2)} + \beta_0)$$

$$\boldsymbol{z}^{(2)} = \sigma((\boldsymbol{\alpha}^{(2)})^T \boldsymbol{z}^{(1)} + \boldsymbol{b}^{(2)})$$

$$\boldsymbol{z}^{(1)} = \sigma((\boldsymbol{\alpha}^{(1)})^T \boldsymbol{x} + \boldsymbol{b}^{(1)})$$

# Neural Network (Vector Form)

*Neural Network with 1 Hidden Layers and 2 Hidden Units (Matrix Form)*

Output

$y$

$$y = \sigma(\boldsymbol{\beta}^T \mathbf{z})$$

Weights

$\beta_1$   $\beta_2$

Hidden Layer

$z_1$   $z_2$

$$z_2 = \sigma(\boldsymbol{\alpha}_{2,.}^T \mathbf{x})$$

$$z_1 = \sigma(\boldsymbol{\alpha}_{1,.}^T \mathbf{x})$$

Weights

$\alpha_{11}$   $\alpha_{21}$   $\alpha_{12}$   $\alpha_{22}$   $\alpha_{13}$   $\alpha_{23}$

Input

$x_1$   $x_2$   $x_3$

12

# ACTIVATION FUNCTIONS

# Activation Functions

Neural Network with sigmoid activation functions



Output

Hidden Layer

Input

(F) **Loss**
$$J = \frac{1}{2}(y - y^*)^2$$

(E) **Output (sigmoid)**
$$y = \frac{1}{1+\exp(-b)}$$

(D) **Output (linear)**
$$b = \sum_{j=0}^{D} \beta_j z_j$$

(C) **Hidden (sigmoid)**
$$z_j = \frac{1}{1+\exp(-a_j)}, \ \ \forall j$$

(B) **Hidden (linear)**
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \ \ \forall j$$

(A) **Input**
Given $x_i, \ \forall i$

14

# Activation Functions

Neural Network with arbitrary nonlinear activation functions



Output

Hidden Layer

Input

(F) **Loss**
$$J = \tfrac{1}{2}(y - y^*)^2$$

(E) **Output (nonlinear)**
$$y = \sigma(b)$$

(D) **Output (linear)**
$$b = \sum_{j=0}^{D} \beta_j z_j$$

(C) **Hidden (nonlinear)**
$$z_j = \sigma(a_j), \ \forall j$$

(B) **Hidden (linear)**
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \ \forall j$$

(A) **Input**
Given $x_i, \ \forall i$

# Activation Functions

So far, we've assumed that the activation function (nonlinearity) is always the sigmoid function...

**Sigmoid (aka. logistic) function**



$\sigma(x) = \frac{1}{1 + \exp(-x)}$

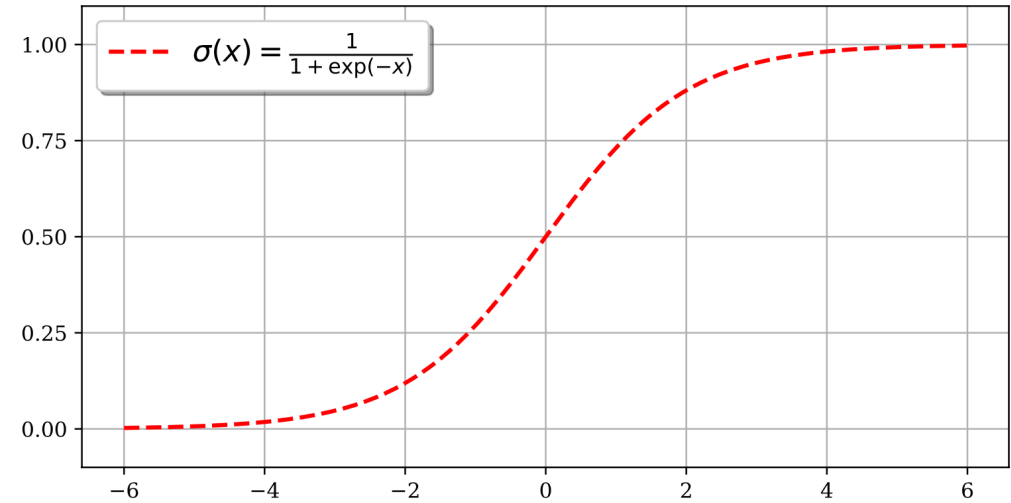... but the sigmoid is not widely used in modern neural networks
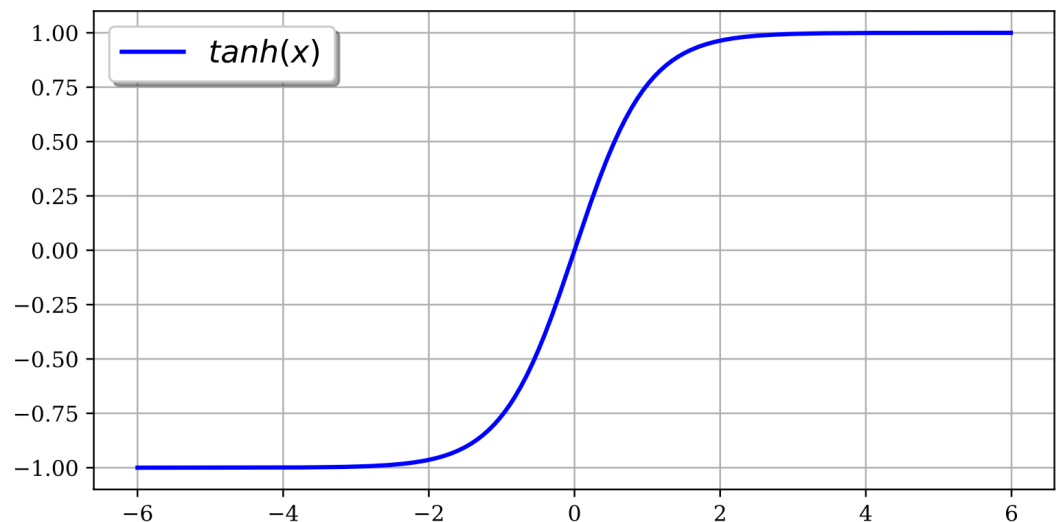
**Hyperbolic tangent function**



$tanh(x)$

# Activation Functions

- ## sigmoid, σ(x)
  - output in range (0,1)
  - good for probabilistic outputs
- ## hyperbolic tangent, tanh(x)
  - similar shape to sigmoid, but output in range (-1,+1)
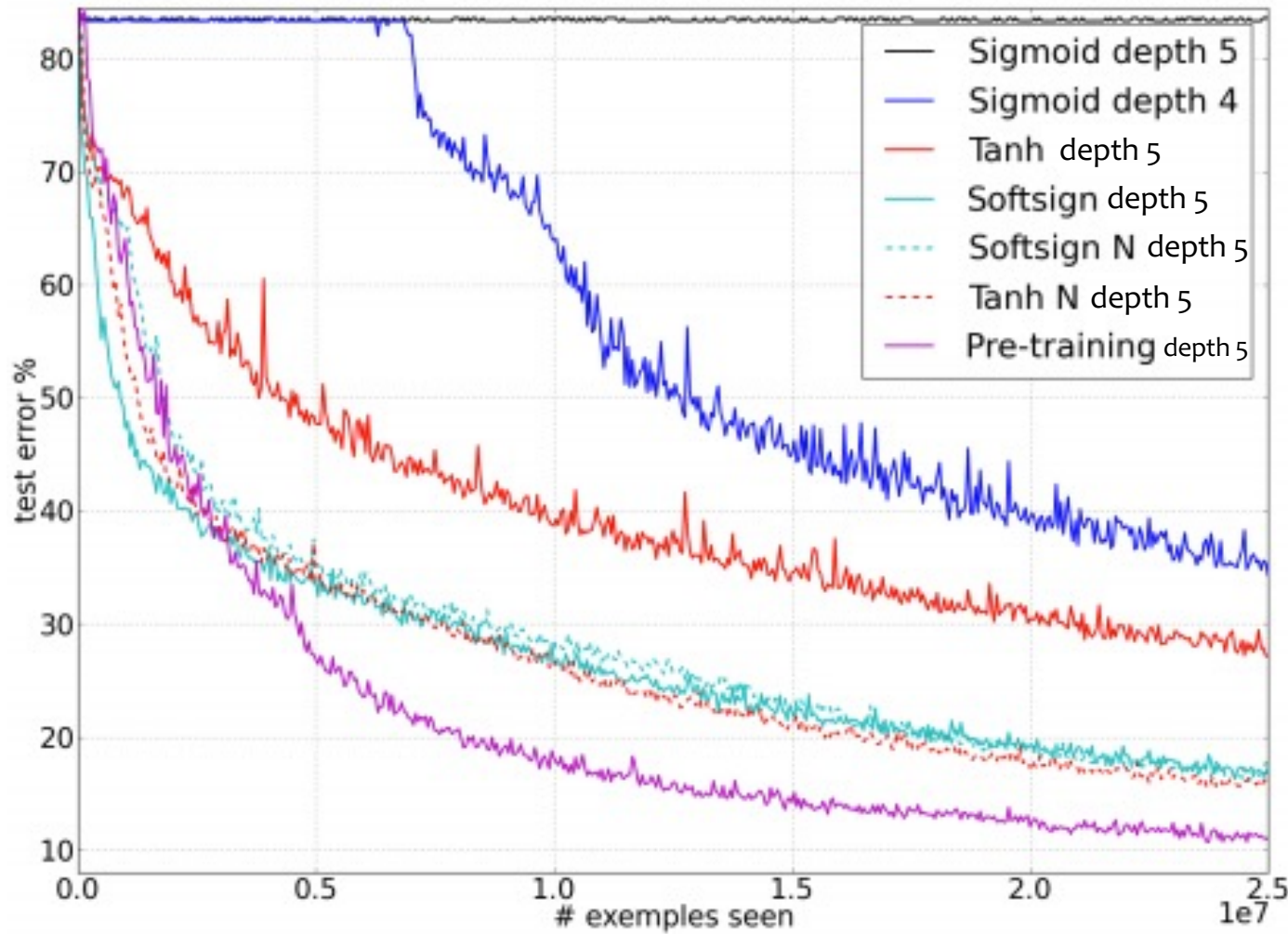
**Sigmoid (aka. logistic) function**



$\sigma(x) = \frac{1}{1 + \exp(-x)}$

**Hyperbolic tangent function**



$tanh(x)$

# Understanding the difficulty of training deep feedforward neural networks

AI Stats 2010
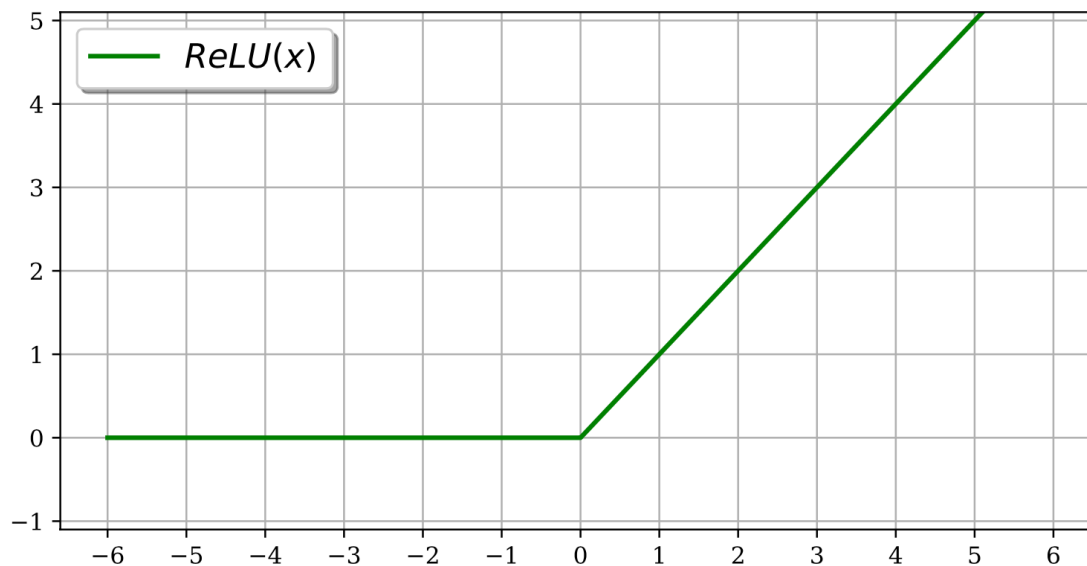


Figure from Glorot & Bentio (2010)

# Activation Functions

- Rectified Linear Unit (ReLU)
    - avoids the vanishing gradient problem
    - derivative is fast to compute

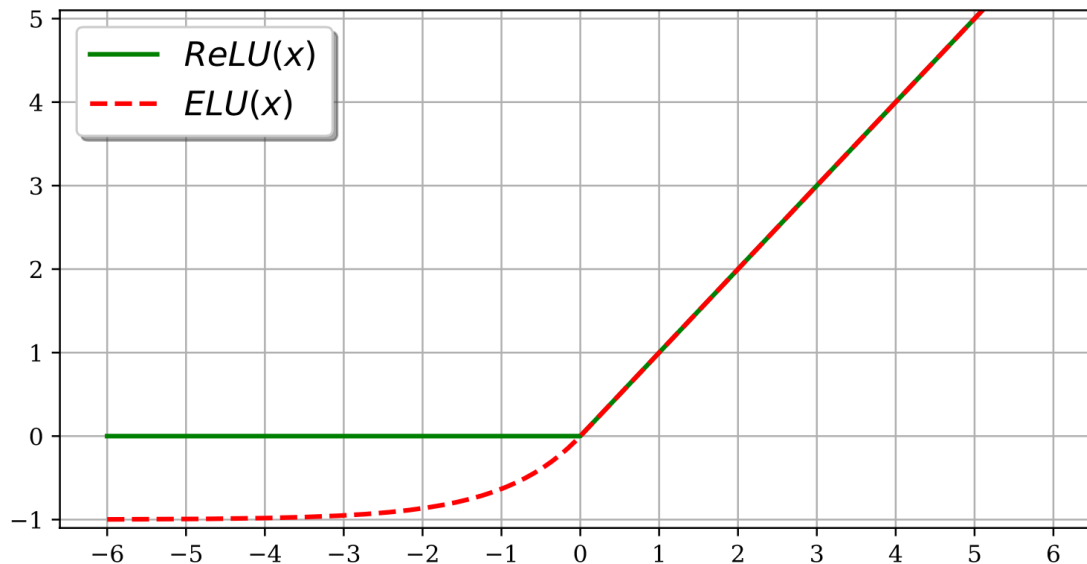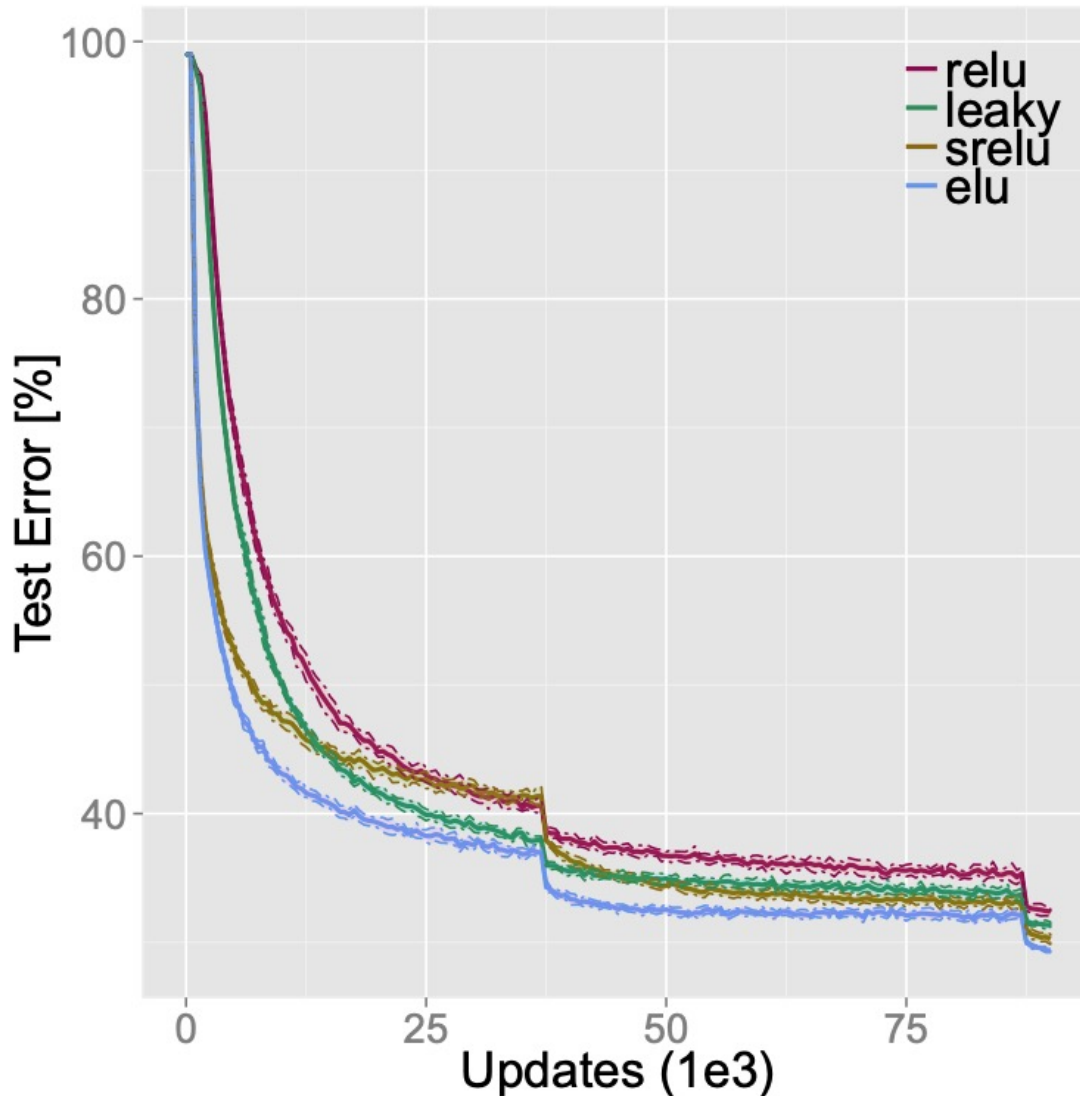$$\text{ReLU}(x) = max(0, x)$$

# Activation Functions

- Rectified Linear Unit (ReLU)
  - avoids the vanishing gradient problem
  - derivative is fast to compute

$$\textbf{ReLU}(x) = max(0, x)$$

- Exponential Linear Unit (ELU)
  - same as ReLU on positive inputs
  - unlike ReLU, allows negative outputs and smoothly transitions for x < 0

$$\textbf{ELU}(x) = \begin{cases} x, & \text{if } x > 0 \\ \alpha(\exp(x) - 1), & \text{if } x \leq 0 \end{cases}$$

# Activation Functions
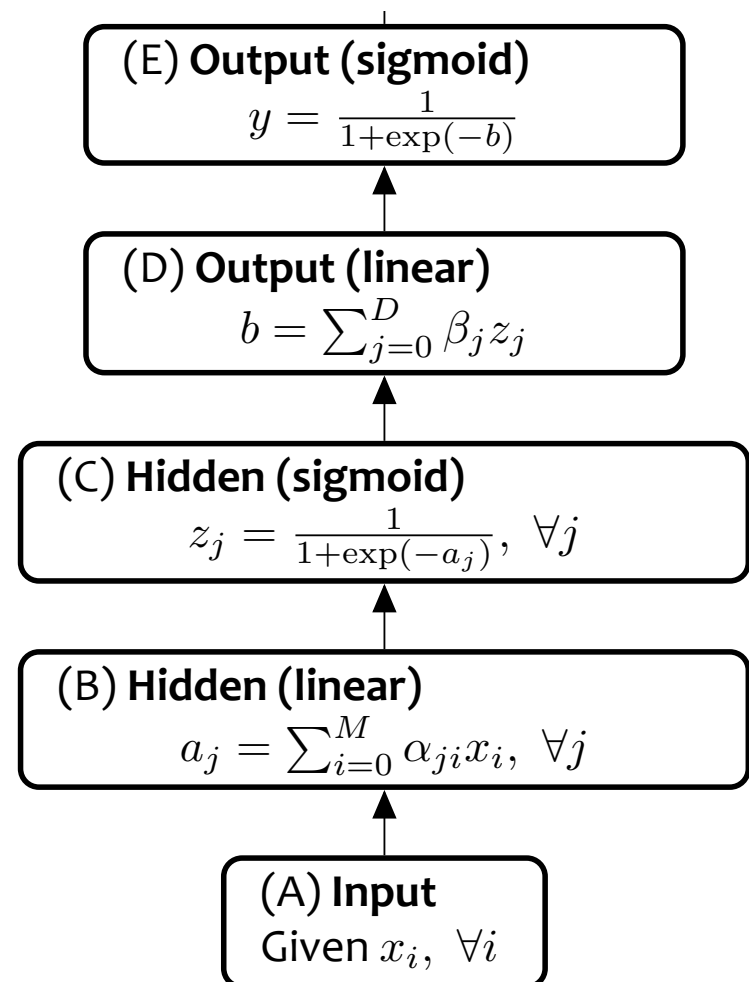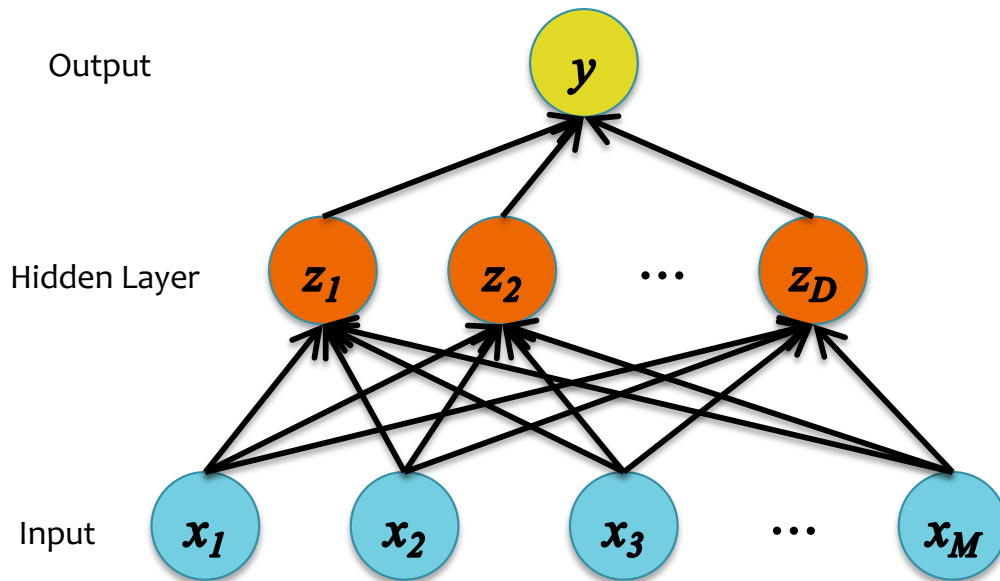
**Image Classification Benchmark (CIFAR-10)**



1. Training loss converges fastest with ELU

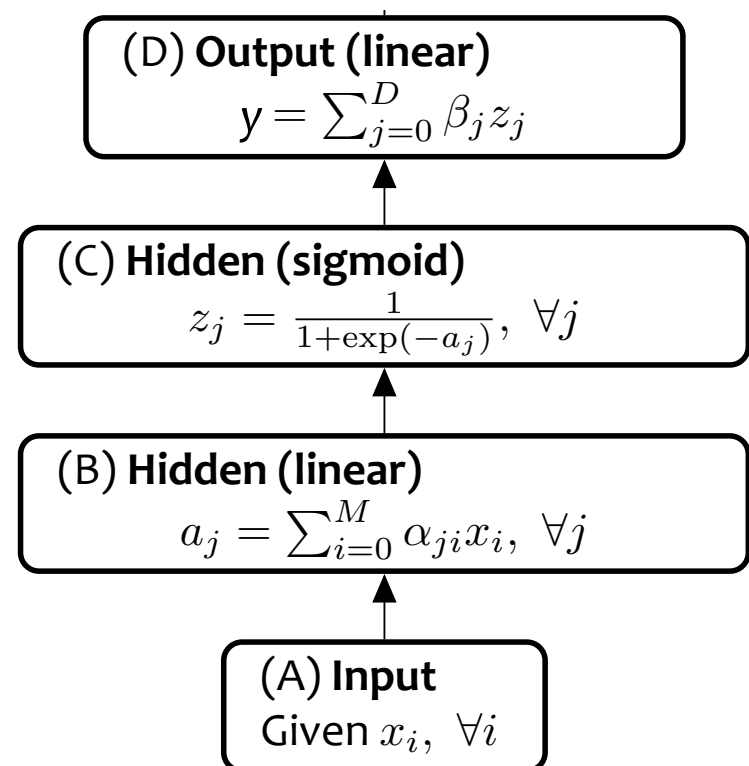2. ELU(x) yields lower test error than ReLU(x) on CIFAR-10

Figure from Clevert et al. (2016)

# LOSS FUNCTIONS & OUTPUT LAYERS

# Neural Network for Classification



**Output**

**Hidden Layer**

**Input**

(E) **Output (sigmoid)**
$$y = \frac{1}{1+\exp(-b)}$$

(D) **Output (linear)**
$$b = \sum_{j=0}^{D} \beta_j z_j$$

(C) **Hidden (sigmoid)**
$$z_j = \frac{1}{1+\exp(-a_j)}, \ \forall j$$

(B) **Hidden (linear)**
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \ \forall j$$

(A) **Input**
Given $x_i, \ \forall i$

23

# Neural Network for Regression



Output

Hidden Layer

Input

(D) **Output (linear)**
$$y = \sum_{j=0}^{D} \beta_j z_j$$

(C) **Hidden (sigmoid)**
$$z_j = \frac{1}{1+\exp(-a_j)}, \ \ \forall j$$

(B) **Hidden (linear)**
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \ \ \forall j$$

(A) **Input**
Given $x_i, \ \forall i$

# Objective Functions for NNs

1. Quadratic Loss:
   - the same objective as Linear Regression
   - i.e. mean squared error

$$J = \ell_Q(y, y^{(i)}) = \frac{1}{2}(y - y^{(i)})^2$$
$$\frac{dJ}{dy} = y - y^{(i)}$$

2. Binary Cross-Entropy:
   - the same objective as Binary Logistic Regression
   - i.e. negative log likelihood
   - This requires our output y to be a probability in [0,1]

$$J = \ell_{CE}(y, y^{(i)}) = -(y^{(i)} \log(y) + (1 - y^{(i)}) \log(1 - y))$$
$$\frac{dJ}{dy} = -\left(y^{(i)} \frac{1}{y} + (1 - y^{(i)}) \frac{1}{y - 1}\right)$$

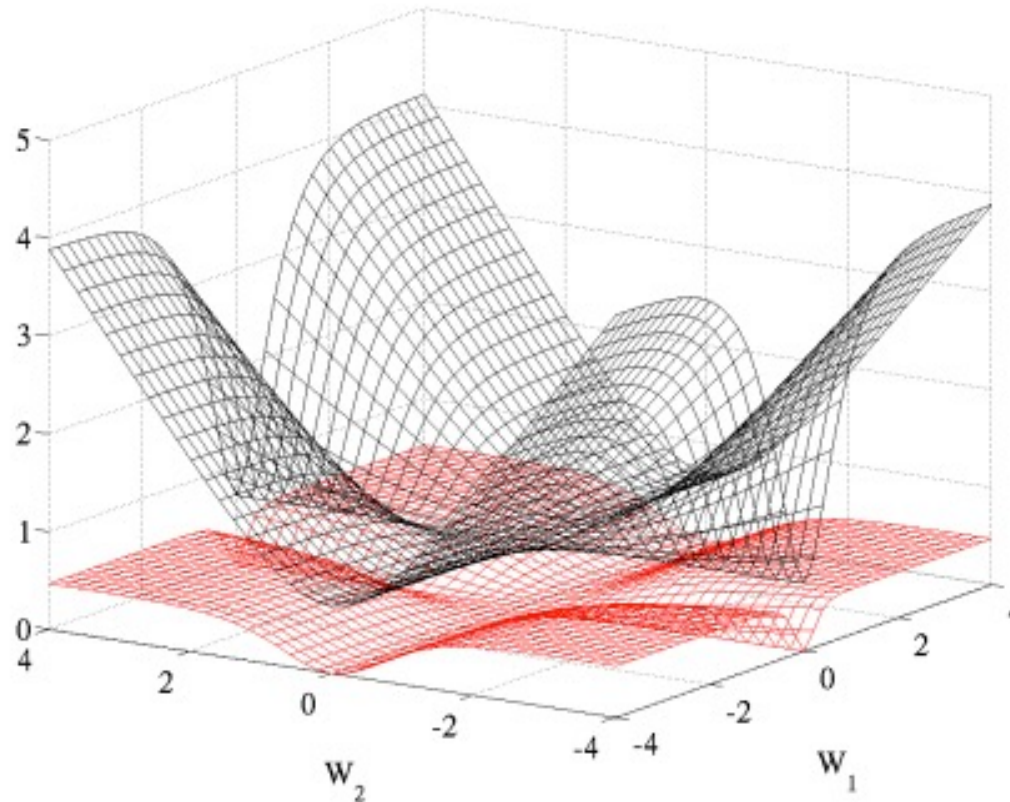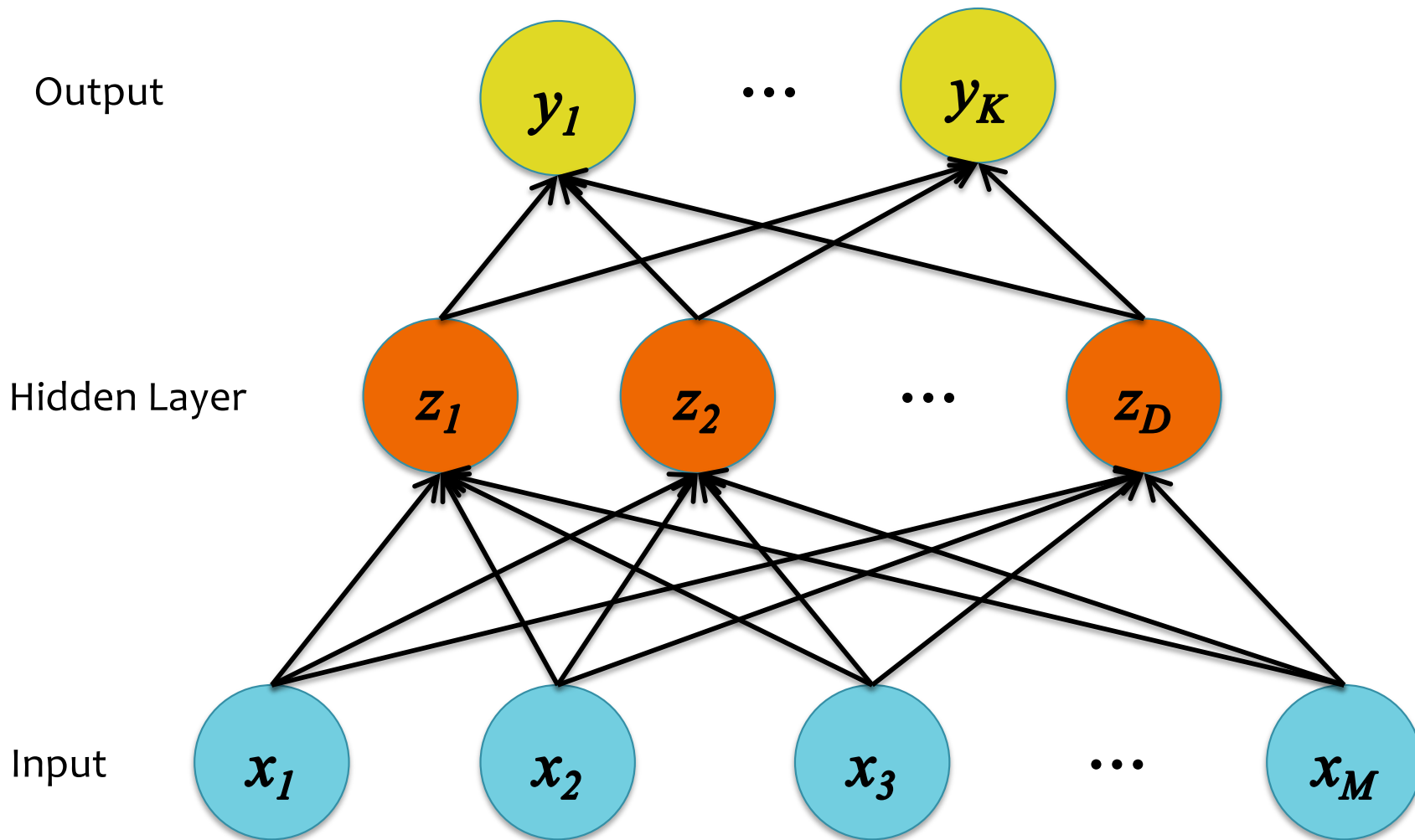# Objective Functions for NNs

**Cross-entropy vs. Quadratic loss**



Figure 5: Cross entropy (black, surface on top) and quadratic (red, bottom surface) cost as a function of two weights (one at each layer) of a network with two layers, $W_1$ respectively on the first layer and $W_2$ on the second, output layer.
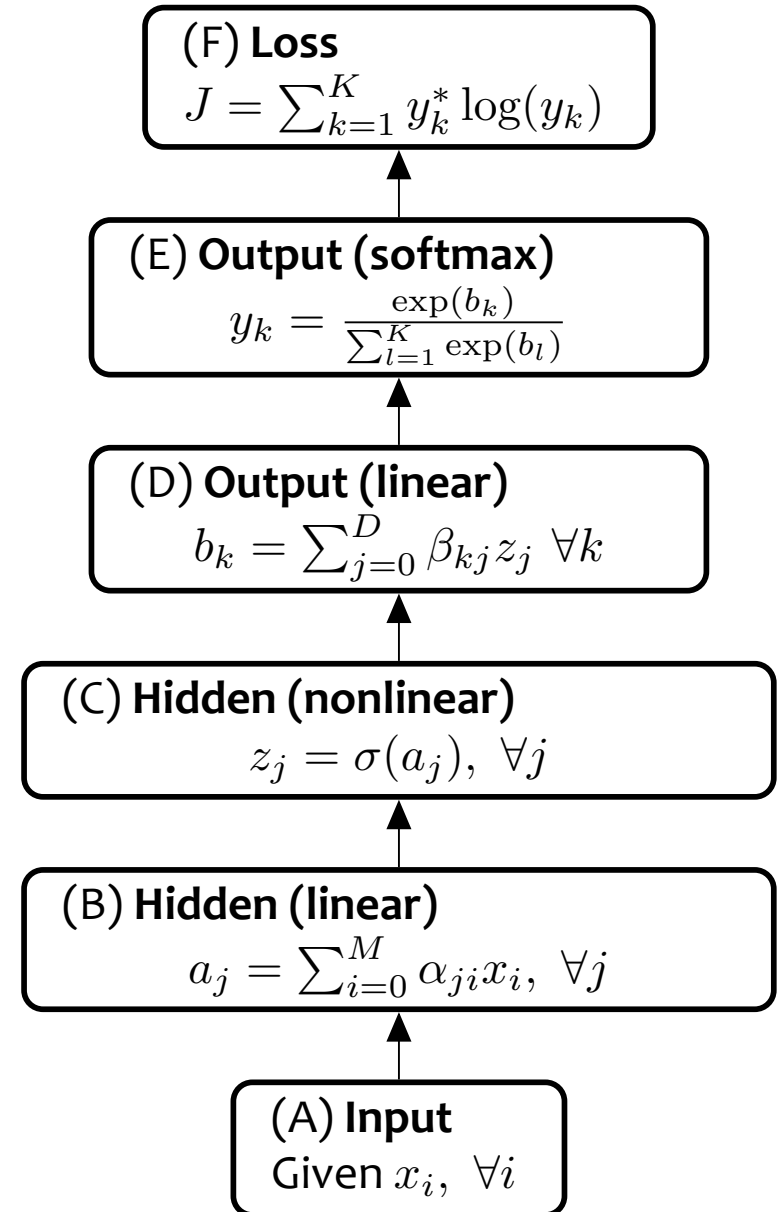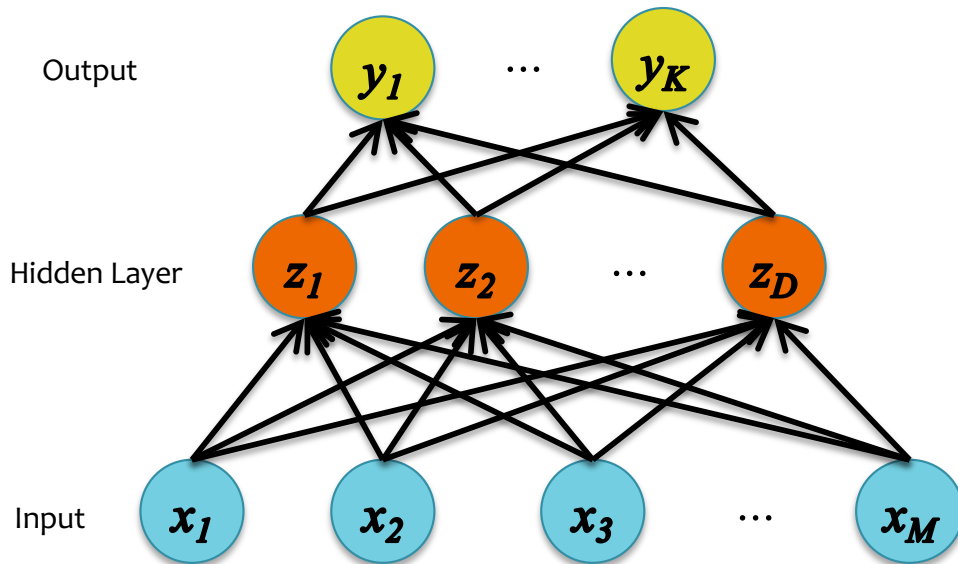
# Multiclass Output

Output

Hidden Layer

Input

# Multiclass Output

Softmax:

$$y_k = \frac{\exp(b_k)}{\sum_{l=1}^{K} \exp(b_l)}$$



(F) **Loss**
$$J = \sum_{k=1}^{K} y_k^* \log(y_k)$$

(E) **Output (softmax)**
$$y_k = \frac{\exp(b_k)}{\sum_{l=1}^{K} \exp(b_l)}$$

(D) **Output (linear)**
$$b_k = \sum_{j=0}^{D} \beta_{kj} z_j \ \forall k$$

(C) **Hidden (nonlinear)**
$$z_j = \sigma(a_j), \ \forall j$$

(B) **Hidden (linear)**
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \ \forall j$$

(A) **Input**
Given $x_i, \ \forall i$

28

# Objective Functions for NNs

3. Cross-Entropy for Multiclass Outputs:
   - i.e. negative log likelihood for multiclass outputs
   - Suppose output is a random variable Y that takes one of K values
   - Let $\mathbf{y}^{(i)}$ represent our true label as a one-hot vector:

$$\mathbf{y}^{(i)} = \begin{array}{|c|c|c|c|c|c|c|c|} \hline 0 & 0 & 0 & 1 & 0 & 0 & \ldots & 0 \\ \hline \end{array}$$

$$\quad\ \ 1 \quad\ 2 \quad\ 3 \quad\ 4 \quad\ 5 \quad\ 6 \quad \ldots \quad K$$

   - Assume our model outputs a length K vector of probabilities:

$$\mathbf{y} = \text{softmax}(f_{\text{scores}}(\mathbf{x}, \boldsymbol{\theta}))$$

   - Then we can write the log-likelihood of a single training example $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$ as:
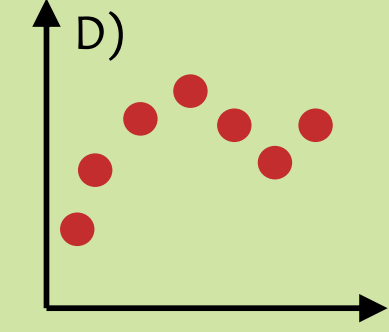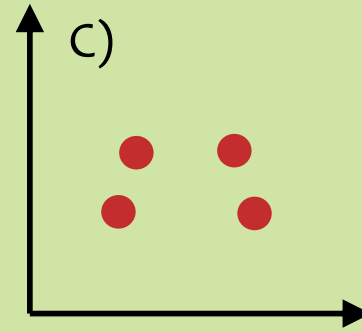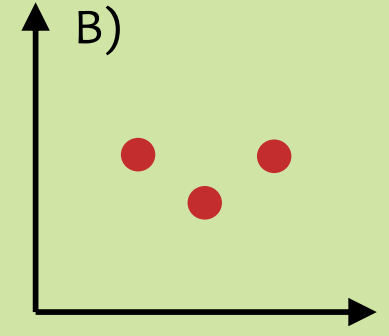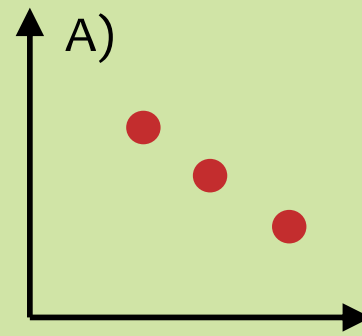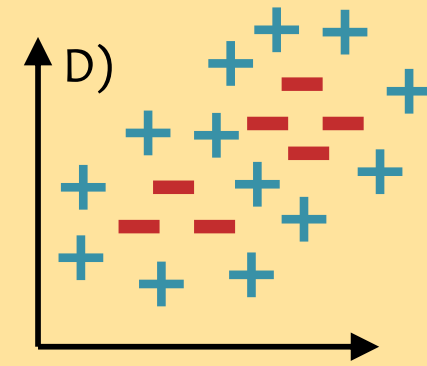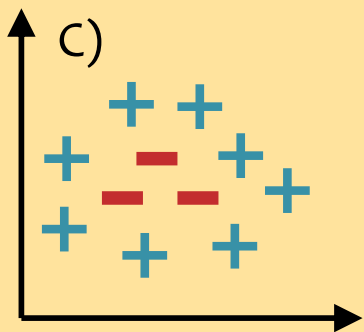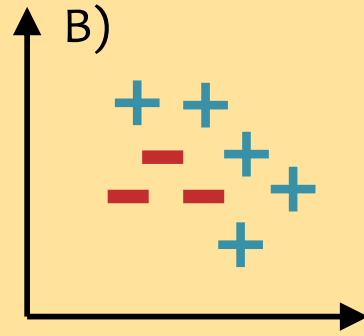
$$J = \ell_{CE}(\mathbf{y}, \mathbf{y}^{(i)}) = -\sum_{k=1}^{K} y_k^{(i)} \log(y_k)$$

# Neural Network Errors

**Question X:** For which of the datasets below does there exist a one-hidden layer neural network that achieves zero *classification* error? **Select all that apply.**

**Question Y:** For which of the datasets below does there exist a one-hidden layer neural network for *regression* that achieves *nearly* zero MSE? **Select all that apply.**

# Neural Networks Objectives

*You should be able to…*

- Explain the biological motivations for a neural network
- Combine simpler models (e.g. linear regression, binary logistic regression, multinomial logistic regression) as components to build up feed-forward neural network architectures
- Explain the reasons why a neural network can model nonlinear decision boundaries for classification
- Compare and contrast feature engineering with learning features
- Identify (some of) the options available when designing the architecture of a neural network
- Implement a feed-forward neural network

Computing Gradients

# APPROACHES TO DIFFERENTIATION

**Background**

# A Recipe for Machine Learning

**1. Given training data:**

$$\{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^{N}$$

**2. Choose each of these:**

– Decision function

$$\hat{\boldsymbol{y}} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

– Loss function

$$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}_i) \in \mathbb{R}$$

**3. Define goal:**

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \sum_{i=1}^{N} \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

**4. Train with SGD:**

(take small steps opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

## Gradients

**Backpropagation** can compute this gradient!

And it's a **special case of a more general algorithm** called reverse-mode automatic differentiation that can compute the gradient of any differentiable function efficiently!

1. Given training dat

$$\{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^N$$

$$\boldsymbol{y}_i)$$

2. Choose each of th

– Decision functio

$$\hat{y} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

opposite the gradient)

– Loss function

$$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}_i) \in \mathbb{R}$$

$$\boldsymbol{\theta}^{(t} \qquad {}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

# Approaches to Differentiation

- **Question 1:**
  When can we compute the gradients for an arbitrary neural network?

- **Question 2:**
  When can we make the gradient computation efficient?

# Approaches to Differentiation

1. Finite Difference Method
   - Pro: Great for testing implementations of backpropagation
   - Con: Slow for high dimensional inputs / outputs
   - Required: Ability to call the function f(**x**) on any input **x**

2. Symbolic Differentiation
   - Note: The method you learned in high-school
   - Note: Used by Mathematica / Wolfram Alpha / Maple
   - Pro: Yields easily interpretable derivatives
   - Con: Leads to exponential computation time if not carefully implemented
   - Required: Mathematical expression that defines f(**x**)

Given $f : \mathbb{R}^A \to \mathbb{R}^B, f(\mathbf{x})$

Compute $\dfrac{\partial f(\mathbf{x})_i}{\partial x_j} \forall i, j$

# Approaches to Differentiation

3. Automatic Differentiation - Reverse Mode
   – Note: Called *Backpropagation* when applied to Neural Nets
   – Pro: Computes partial derivatives of one output $f(\mathbf{x})_i$ with respect to all inputs $x_j$ in time proportional to computation of $f(\mathbf{x})$
   – Con: Slow for high dimensional outputs (e.g. vector-valued functions)
   – Required: Algorithm for computing $f(\mathbf{x})$
4. Automatic Differentiation - Forward Mode
   – Note: Easy to implement. Uses dual numbers.
   – Pro: Computes partial derivatives of all outputs $f(\mathbf{x})_i$ with respect to one input $x_j$ in time proportional to computation of $f(\mathbf{x})$
   – Con: Slow for high dimensional inputs (e.g. vector-valued $\mathbf{x}$)
   – Required: Algorithm for computing $f(\mathbf{x})$

Given $f : \mathbb{R}^A \to \mathbb{R}^B, f(\mathbf{x})$

Compute $\dfrac{\partial f(\mathbf{x})_i}{\partial x_j} \forall i, j$

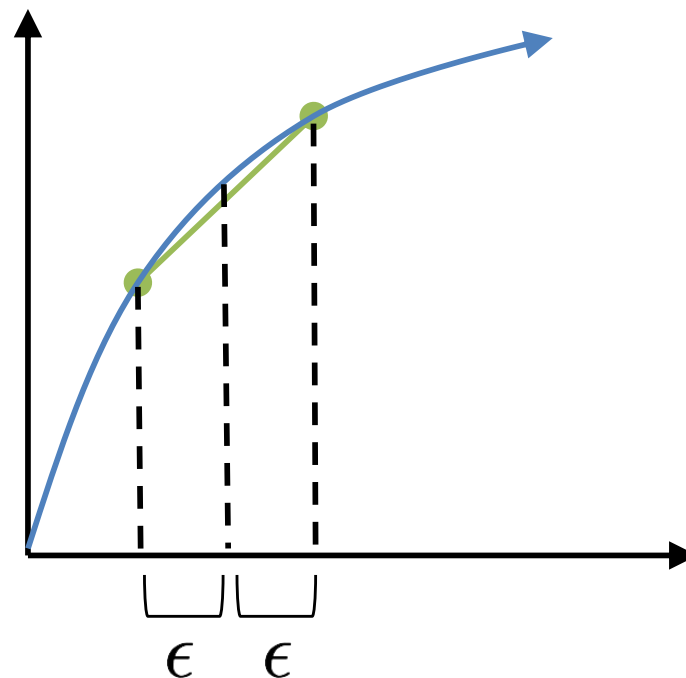# THE FINITE DIFFERENCE METHOD

# Finite Difference Method

The *centered* finite difference approximation is:

$$\frac{\partial}{\partial \theta_i} J(\boldsymbol{\theta}) \approx \frac{(J(\boldsymbol{\theta} + \epsilon \cdot \boldsymbol{d}_i) - J(\boldsymbol{\theta} - \epsilon \cdot \boldsymbol{d}_i))}{2\epsilon} \qquad (1)$$

where $\boldsymbol{d}_i$ is a 1-hot vector consisting of all zeros except for the $i$th entry of $\boldsymbol{d}_i$, which has value 1.

**Notes:**

- Suffers from issues of floating point precision, in practice
- Typically only appropriate to use on small examples with an appropriately chosen epsilon

# Differentiation Quiz

**Speed Quiz: 2 minute time limit.**

**Differentiation Quiz #1:**

Suppose x = 2 and z = 3, what are dy/dx and dy/dz for the function below? **Round your answer to the nearest integer**.

$$y = \exp(xz) + \frac{xz}{\log(x)} + \frac{\sin(\log(x))}{xz}$$

**Answer:** *Answers below are in the form [dy/dx, dy/dz]*

A.  [42, -72]

B.  [72, -42]

C.  [100, 127]

D.  [127, 100]

E.  [1208, 810]

F.  [810, 1208]

G.  [1505, 94]

H.  [94, 1505]

# Differentiation Quiz

## Differentiation Quiz #2:

A neural network with 2 hidden layers can be written as:

$$y = \sigma(\boldsymbol{\beta}^T \sigma((\boldsymbol{\alpha}^{(2)})^T \sigma((\boldsymbol{\alpha}^{(1)})^T \mathbf{x})))$$

where $y \in \mathbb{R}$, $\mathbf{x} \in \mathbb{R}^{D^{(0)}}$, $\boldsymbol{\beta} \in \mathbb{R}^{D^{(2)}}$ and $\boldsymbol{\alpha}^{(i)}$ is a $D^{(i)} \times D^{(i-1)}$ matrix. Nonlinear functions are applied elementwise:

$$\sigma(\mathbf{a}) = [\sigma(a_1), \ldots, \sigma(a_K)]^T$$

Let $\sigma$ be sigmoid: $\sigma(a) = \frac{1}{1 + exp - a}$

What is $\frac{\partial y}{\partial \beta_j}$ and $\frac{\partial y}{\partial \alpha_j^{(i)}}$ for all $i$, $j$.
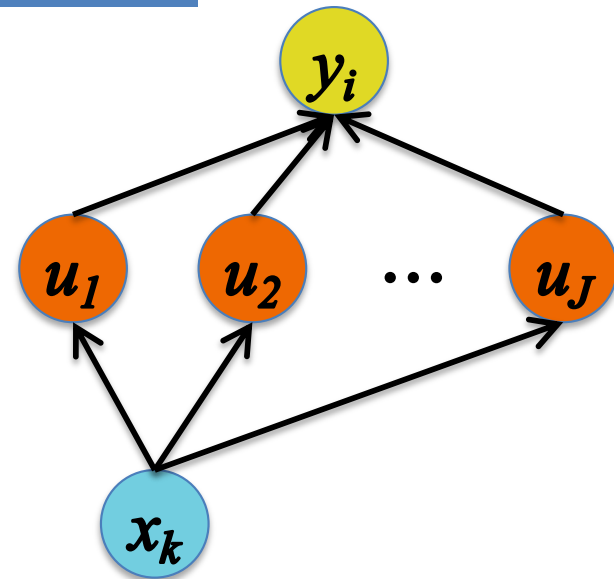
# THE CHAIN RULE OF CALCULUS

# Chain Rule

*Whiteboard*

– Chain Rule of Calculus

# Chain Rule

**Given:** $\boldsymbol{y} = g(\boldsymbol{u})$ and $\boldsymbol{u} = h(\boldsymbol{x})$.

**Chain Rule:**

$$\frac{dy_i}{dx_k} = \sum_{j=1}^{J} \frac{dy_i}{du_j} \frac{du_j}{dx_k}, \quad \forall i, k$$
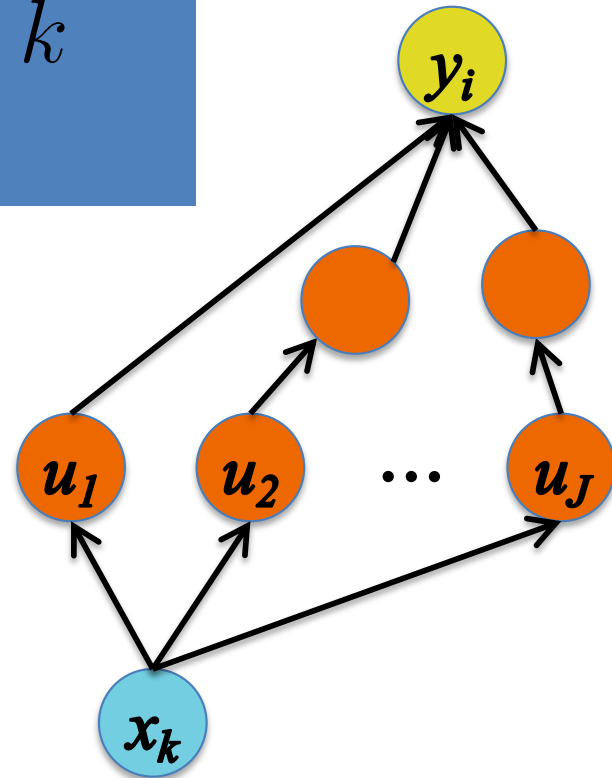
# Chain Rule

**Given:** $y = g(u)$ and $u = h(x)$.

**Chain Rule:**

$$\frac{dy_i}{dx_k} = \sum_{j=1}^{J} \frac{dy_i}{du_j}\frac{du_j}{dx_k}, \quad \forall i, k$$

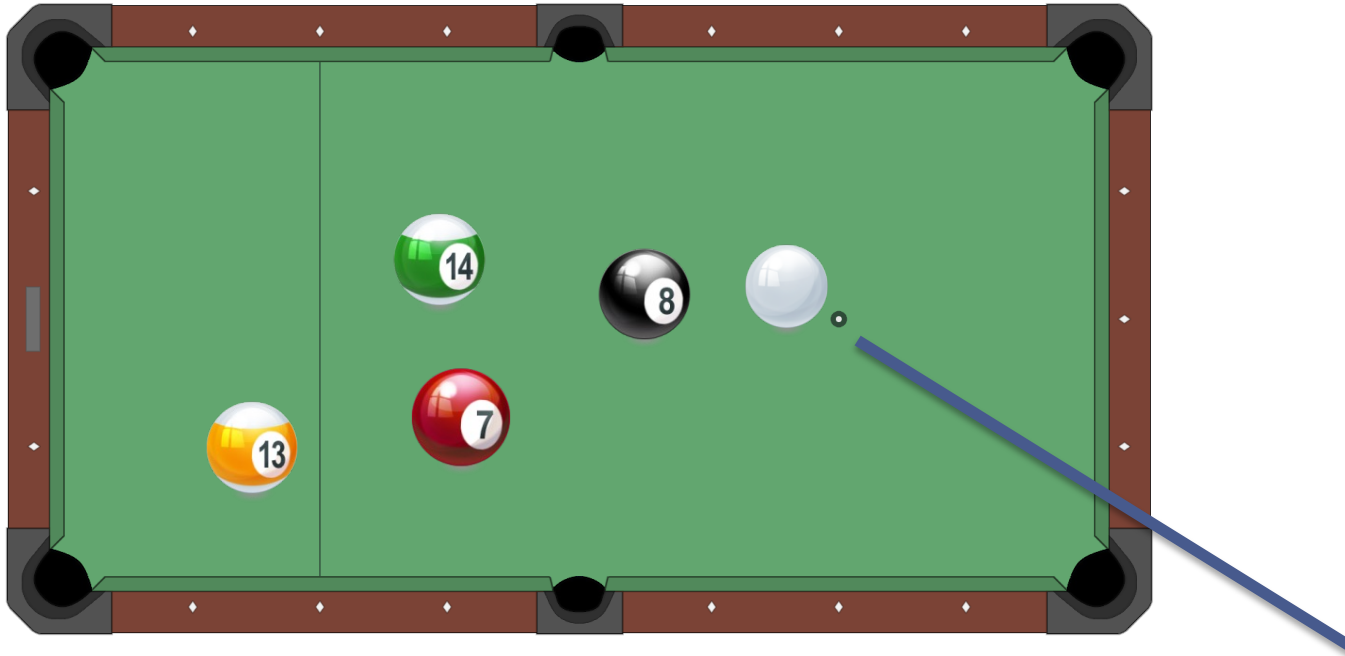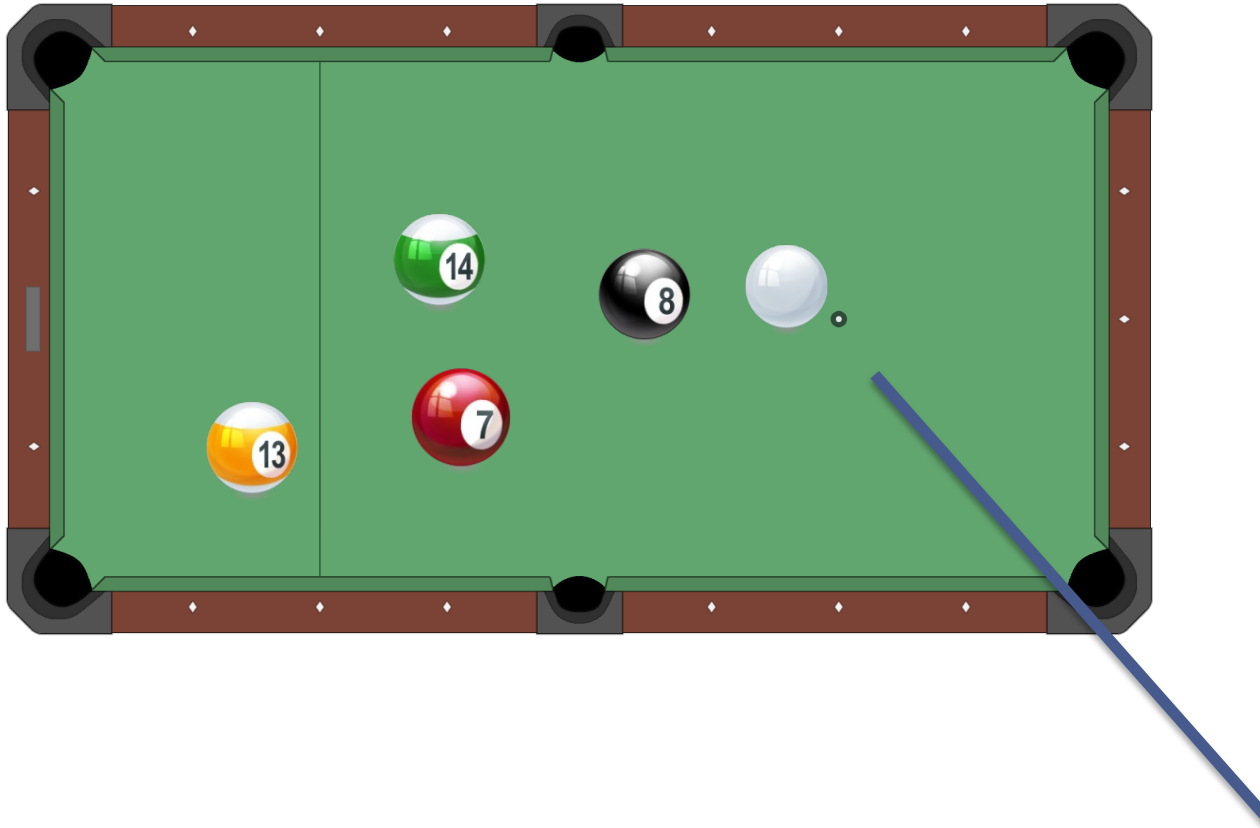**Backpropagation** is just repeated application of the **chain rule** from Calculus 101.
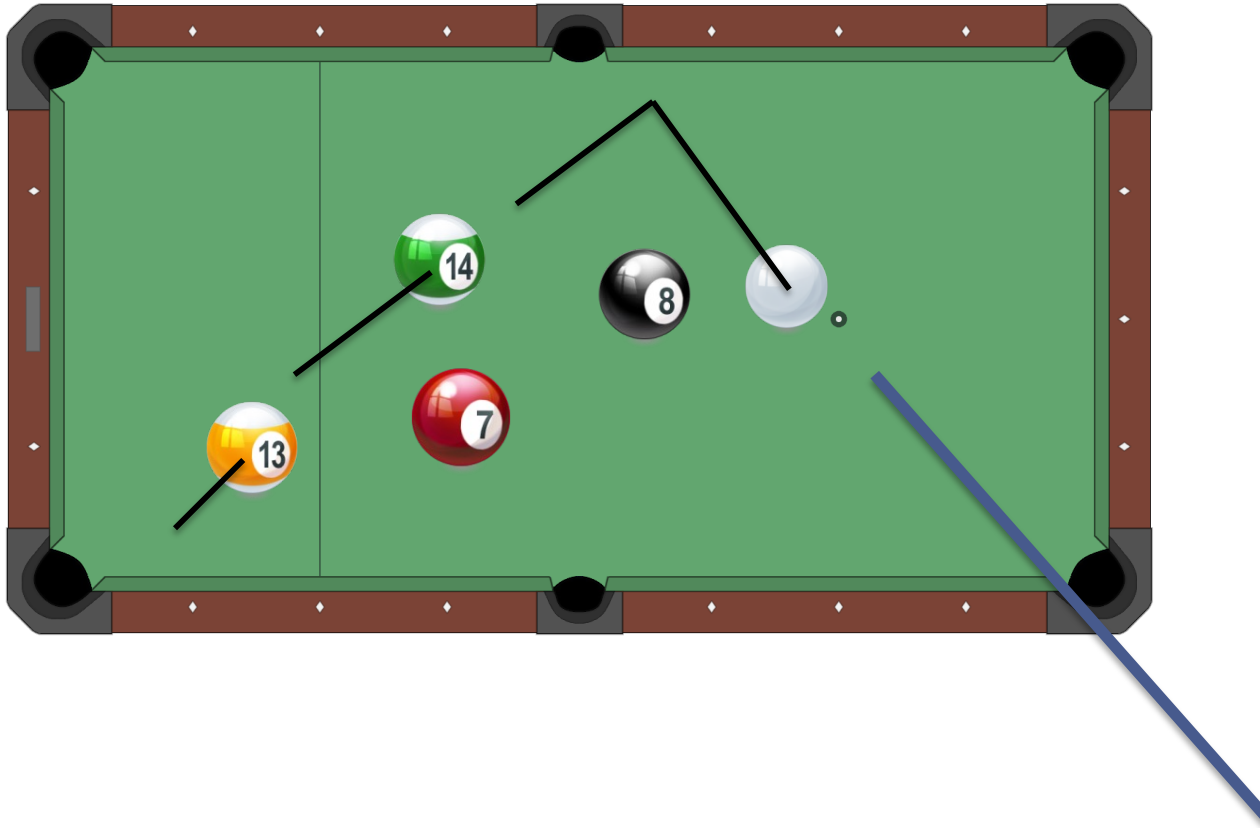
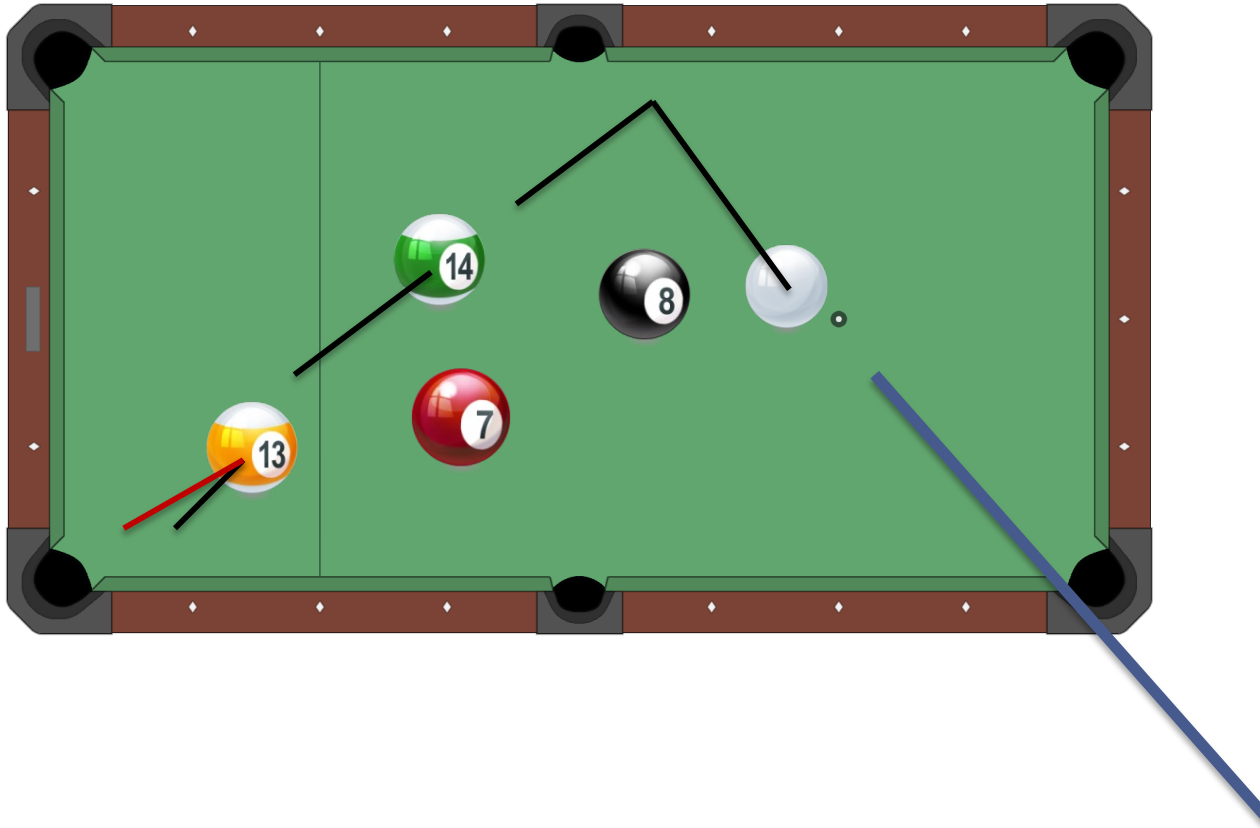Intuitions

# BACKPROPAGATION OF ERRORS
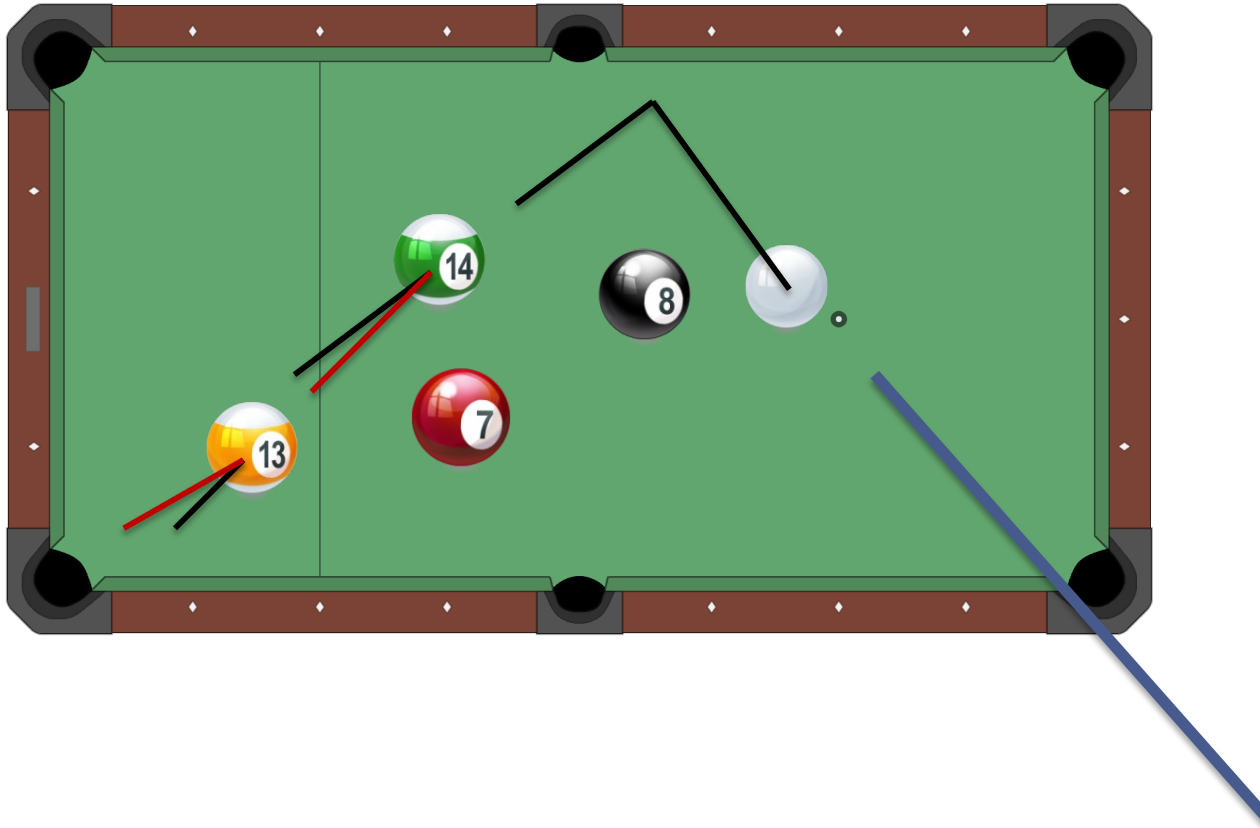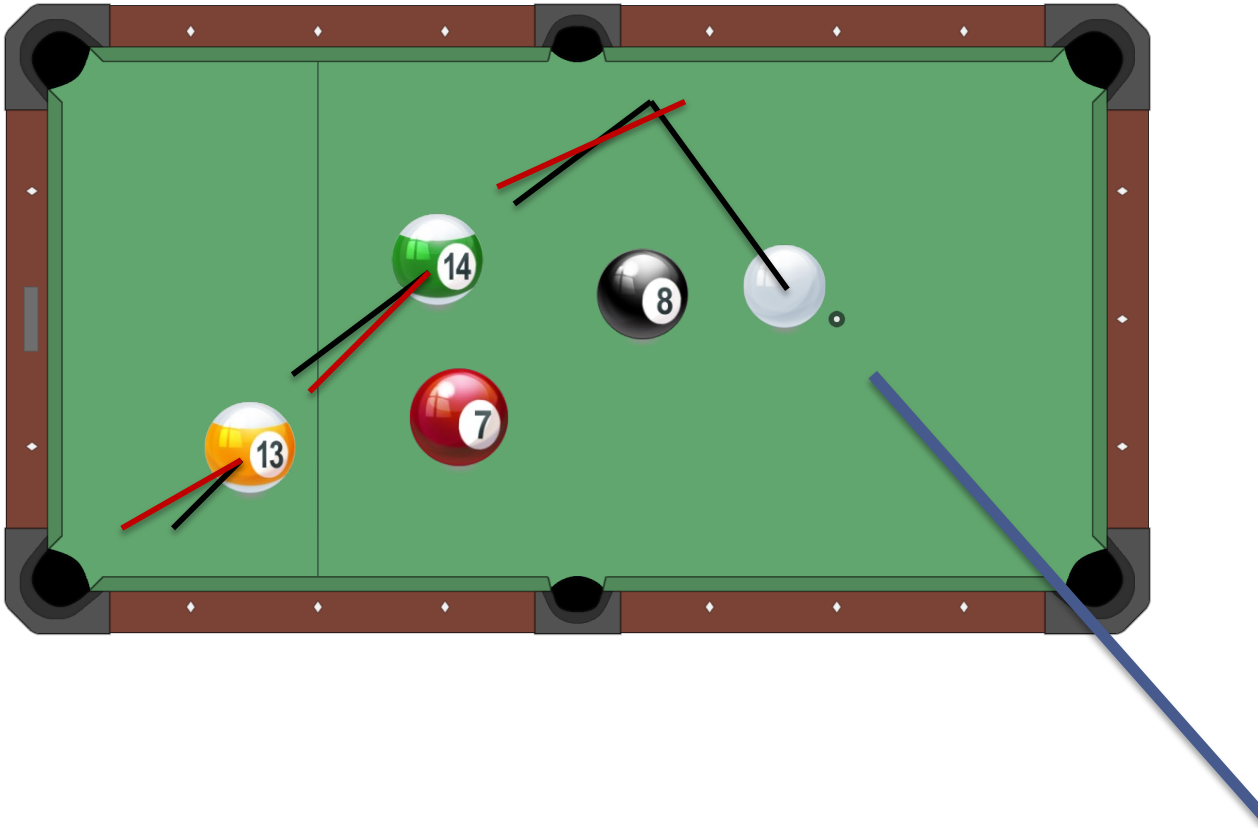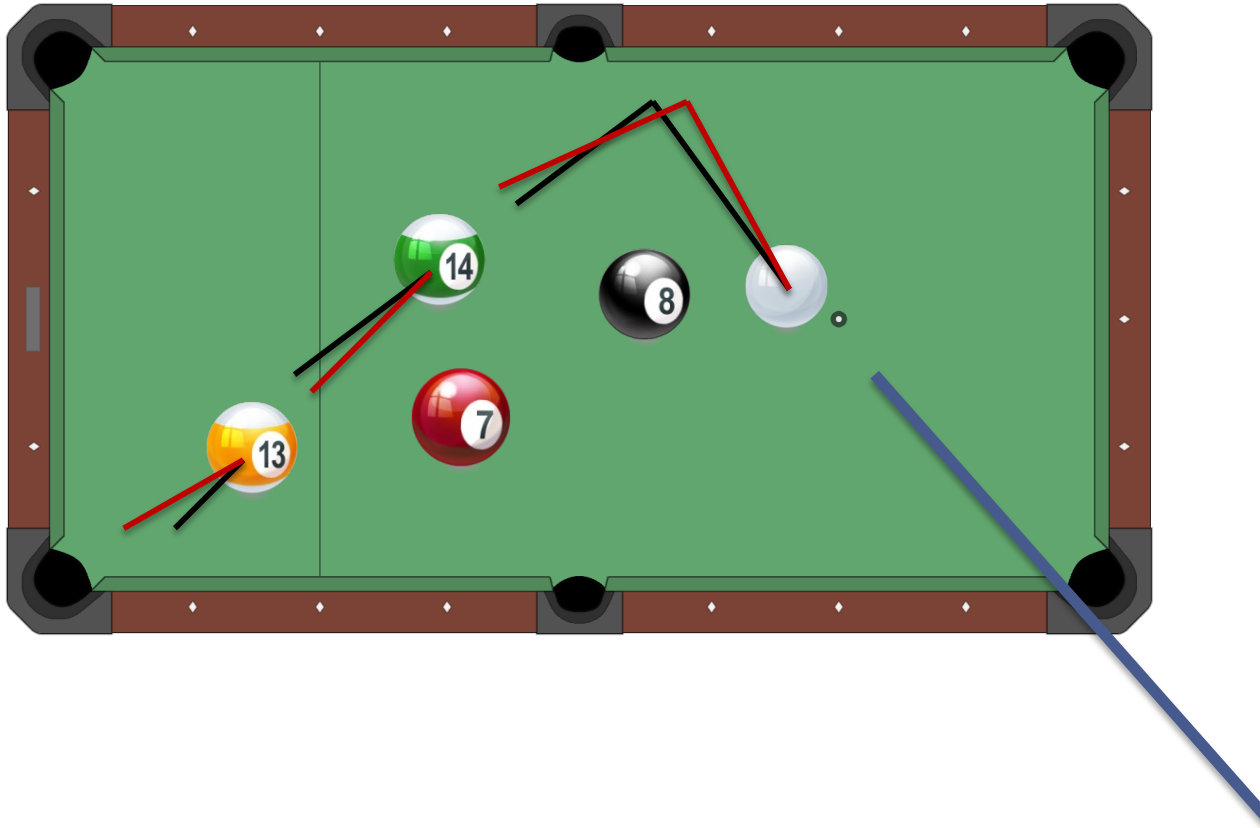
# Error Back-Propagation

50

# Error Back-Propagation
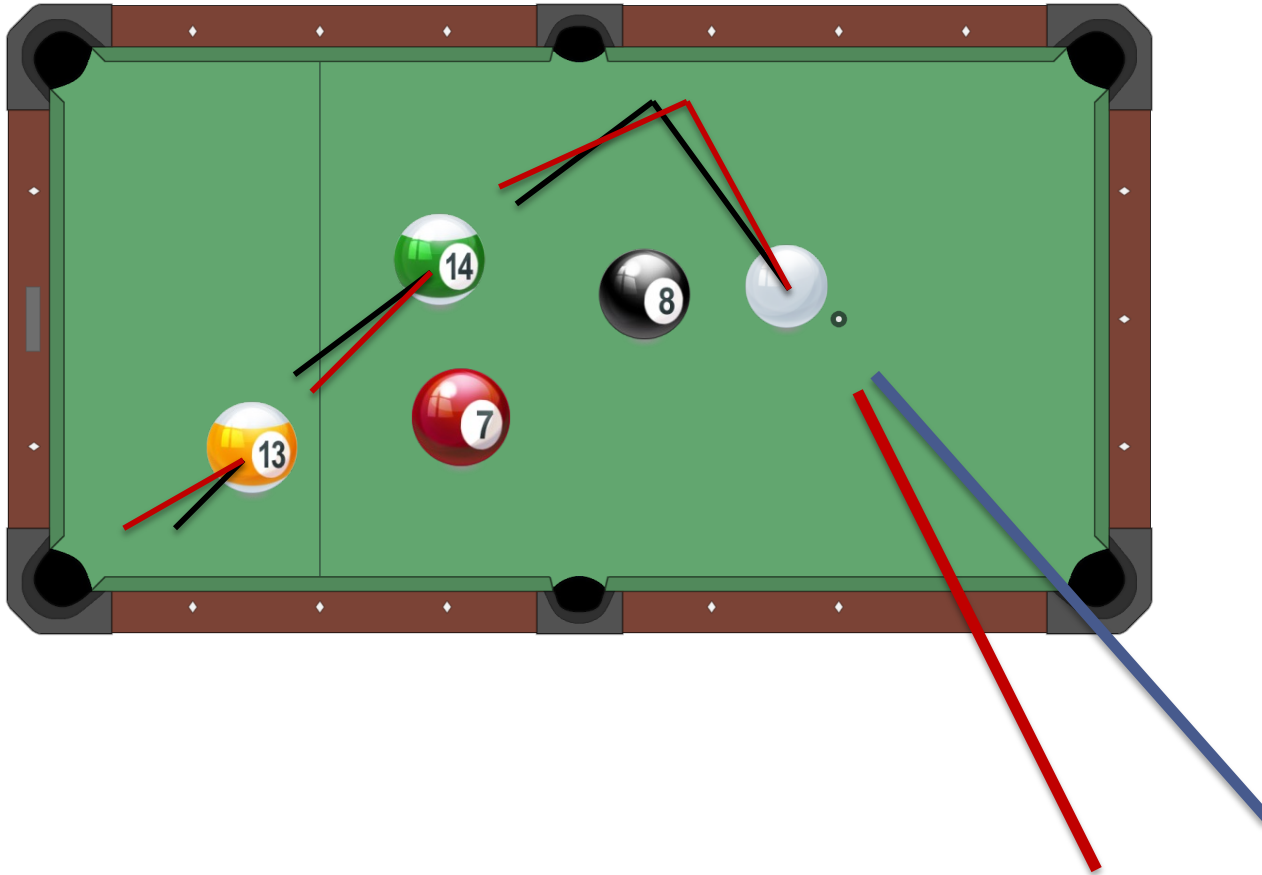
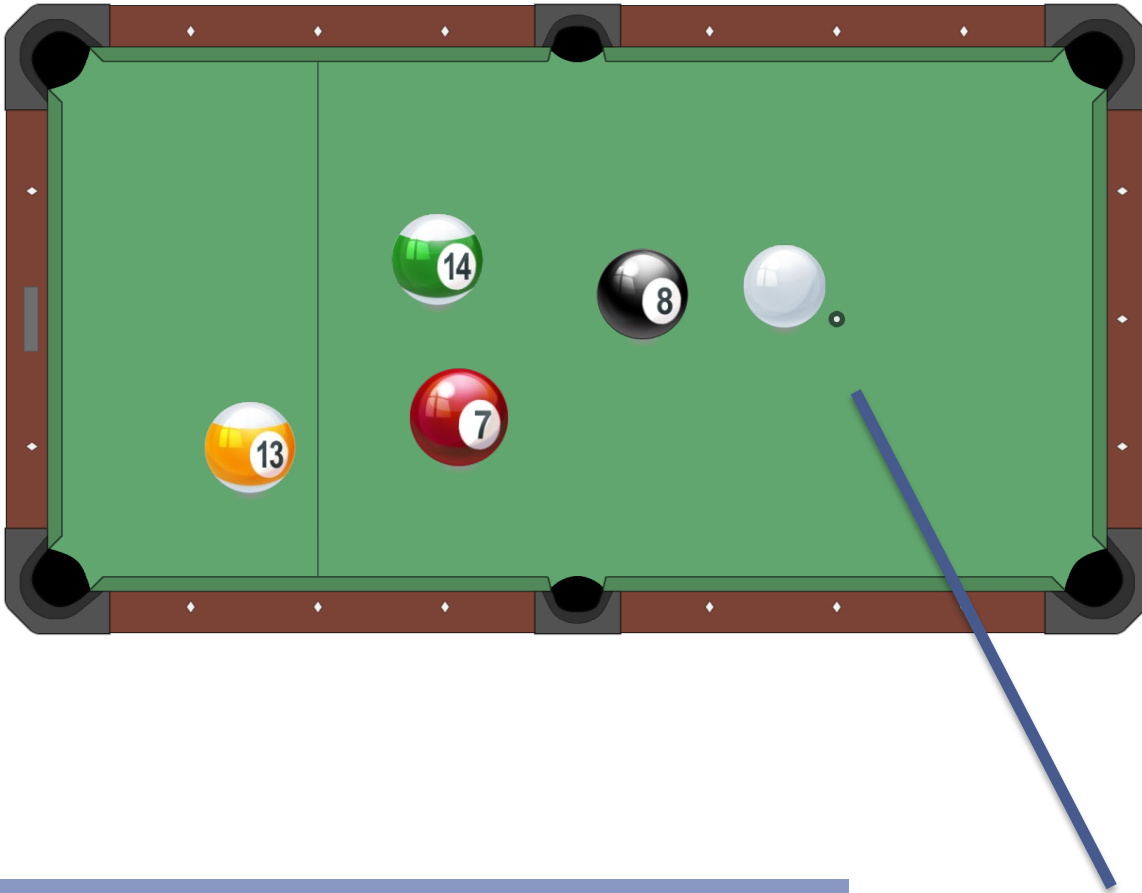# Error Back-Propagation

52

# Error Back-Propagation

# Error Back-Propagation

54

# Error Back-Propagation

# Error Back-Propagation

Slide from (Stoyanov & Eisner, 2012)

# Error Back-Propagation

# Error Back-Propagation

58

# Error Back-Propagation



$p(y|\mathbf{x}^{(i)})$

$\mathbf{z}$

$\theta$

$y^{(i)}$

59

# FORWARD COMPUTATION FOR A COMPUTATION GRAPH

# Backpropagation

*Whiteboard*

– From equation to forward computation

– Representing a simple function as a computation graph

**Differentiation Quiz #1:**

Suppose x = 2 and z = 3, what are dy/dx and dy/dz for the function below? **Round your answer to the nearest integer.**

$$y = \exp(xz) + \frac{xz}{\log(x)} + \frac{\sin(\log(x))}{xz}$$

Algorithm

# BACKPROPAGATION FOR A COMPUTATION GRAPH

# Backpropagation

*Whiteboard*

– Backprogation on a simple computation graph

**Differentiation Quiz #1:**

Suppose x = 2 and z = 3, what are dy/dx and dy/dz for the function below? **Round your answer to the nearest integer.**

$$y = \exp(xz) + \frac{xz}{\log(x)} + \frac{\sin(\log(x))}{xz}$$

# Backpropagation

**Simple Example:** The goal is to compute $J = \cos(\sin(x^2) + 3x^2)$ on the forward pass and the derivative $\frac{dJ}{dx}$ on the backward pass.

Forward

$J = cos(u)$

$u = u_1 + u_2$

$u_1 = sin(t)$

$u_2 = 3t$

$t = x^2$

# Backpropagation

**Simple Example:** The goal is to compute $J = \cos(\sin(x^2) + 3x^2)$ on the forward pass and the derivative $\frac{dJ}{dx}$ on the backward pass.

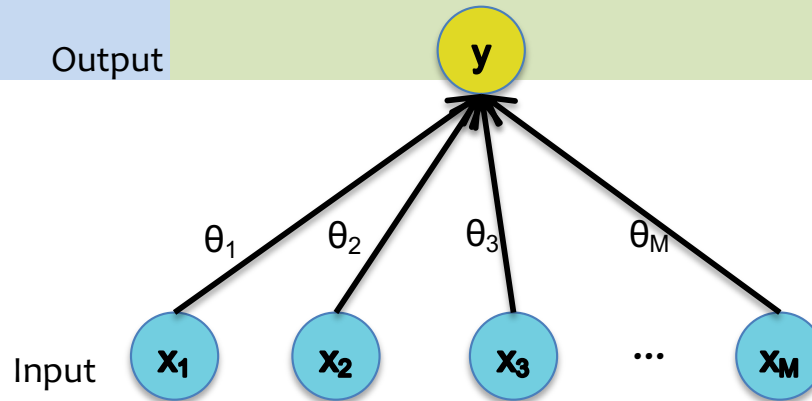| Forward | Backward |
|---|---|
| $J = cos(u)$ | $\dfrac{dJ}{du} \mathrel{+}= -sin(u)$ |
| $u = u_1 + u_2$ | $\dfrac{dJ}{du_1} \mathrel{+}= \dfrac{dJ}{du}\dfrac{du}{du_1},\quad \dfrac{du}{du_1} = 1 \qquad\qquad \dfrac{dJ}{du_2} \mathrel{+}= \dfrac{dJ}{du}\dfrac{du}{du_2},\quad \dfrac{du}{du_2} = 1$ |
| $u_1 = sin(t)$ | $\dfrac{dJ}{dt} \mathrel{+}= \dfrac{dJ}{du_1}\dfrac{du_1}{dt},\quad \dfrac{du_1}{dt} = \cos(t)$ |
| $u_2 = 3t$ | $\dfrac{dJ}{dt} \mathrel{+}= \dfrac{dJ}{du_2}\dfrac{du_2}{dt},\quad \dfrac{du_2}{dt} = 3$ |
| $t = x^2$ | $\dfrac{dJ}{dx} \mathrel{+}= \dfrac{dJ}{dt}\dfrac{dt}{dx},\quad \dfrac{dt}{dx} = 2x$ |

# Backpropagation

**Case 1:**
**Logistic**
**Regression**

Output

y

$\theta_1$ $\theta_2$ $\theta_3$ $\theta_M$

Input x₁ x₂ x₃ ... xₘ

Forward

$$J = y^* \log y + (1 - y^*) \log(1 - y)$$

$$y = \frac{1}{1 + \exp(-a)}$$

$$a = \sum_{j=0}^{D} \theta_j x_j$$

68

# Backpropagation

Output

**y**

**Case 1:**
**Logistic**
**Regression**

$\theta_1$  $\theta_2$  $\theta_3$  $\theta_M$

Input  $x_1$  $x_2$  $x_3$  ...  $x_M$

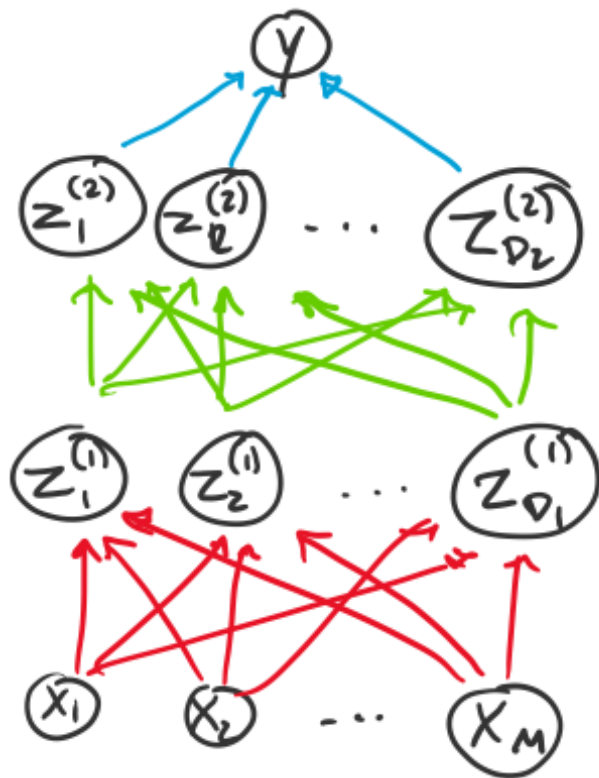| Forward | Backward |
|---|---|
| $J = y^* \log y + (1 - y^*) \log(1 - y)$ | $\dfrac{dJ}{dy} = \dfrac{y^*}{y} + \dfrac{(1 - y^*)}{y - 1}$ |
| $y = \dfrac{1}{1 + \exp(-a)}$ | $\dfrac{dJ}{da} = \dfrac{dJ}{dy}\dfrac{dy}{da}, \ \dfrac{dy}{da} = \dfrac{\exp(-a)}{(\exp(-a) + 1)^2}$ |
| $a = \displaystyle\sum_{j=0}^{D} \theta_j x_j$ | $\dfrac{dJ}{d\theta_j} = \dfrac{dJ}{da}\dfrac{da}{d\theta_j}, \ \dfrac{da}{d\theta_j} = x_j$ $\dfrac{dJ}{dx_j} = \dfrac{dJ}{da}\dfrac{da}{dx_j}, \ \dfrac{da}{dx_j} = \theta_j$ |

A 2-Hidden Layer Neural Network

# TRAINING / FORWARD COMPUTATION / BACKWARD COMPUTATION

# Backpropagation

**Recall:** Our 2-Hidden Layer Neural Network
**Question:** How do we train this model?



$$\beta \in \mathbb{R}^{D_2}$$
$$\beta_0 \in \mathbb{R}$$

$$y = \sigma\left(\vec{\beta}^T \vec{z}^{(2)} + \beta_0\right)$$

$$\alpha^{(2)} \in \mathbb{R}^{D_1 \times D_2}$$
$$\vec{b}^{(2)} \in \mathbb{R}^{D_2}$$

$$\vec{z}^{(2)} = \sigma\left((\alpha^{(2)})^T \vec{z}^{(1)} + \vec{b}^{(2)}\right)$$

$$\alpha^{(1)} \in \mathbb{R}^{M \times D_1}$$
$$\vec{b}^{(1)} \in \mathbb{R}^{D_1}$$

$$\vec{z}^{(1)} = \sigma\left((\alpha^{(1)})^T \vec{x} + \vec{b}^{(1)}\right)$$

# Backpropagation

*Whiteboard*

– Example: Backpropagation for Neural Network with 2-Hidden Layers

- SGD Training
- Forward Computation
- Computation Graph
- Backward Computation