# CS 4364/6364

# **Machine Learning**

Fall Semester 9/28/2023
Lecture 11
Back-Propagation 1

John Sipple
jsipple@gwu.edu

Overview and fundamentals

Info Theory, Calculus, Chain Rule, Optimization

Additive Models, Trees and Related Methods, Random Forests, SVM
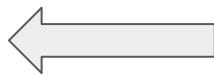
Regularization Optimization

Anomaly Detection Convolutional Nets

Attention, Transformer, Autoencoders, GANs

Thanksgiving

Reinforcement Learning

Final Exam

1  2  3  4  5  6  7  8  9  10  11  12  13  14  15  16  17

Linear Algebra, PCA, Probability & Info Theory

Feedforward Networks and Backprop

Similarity and Clustering

Explainable AI LLMs

AI & Society Project Presentations

Linear Regression, Linear Classification

Midterm + Break

Computer Vision Sequence Models

# General Architecture of Feedforward Nets



Output prediction
$\hat{y}$

weights & bias

Output layer $n$

Hidden layers $2,...,n-1$

Forward Prop

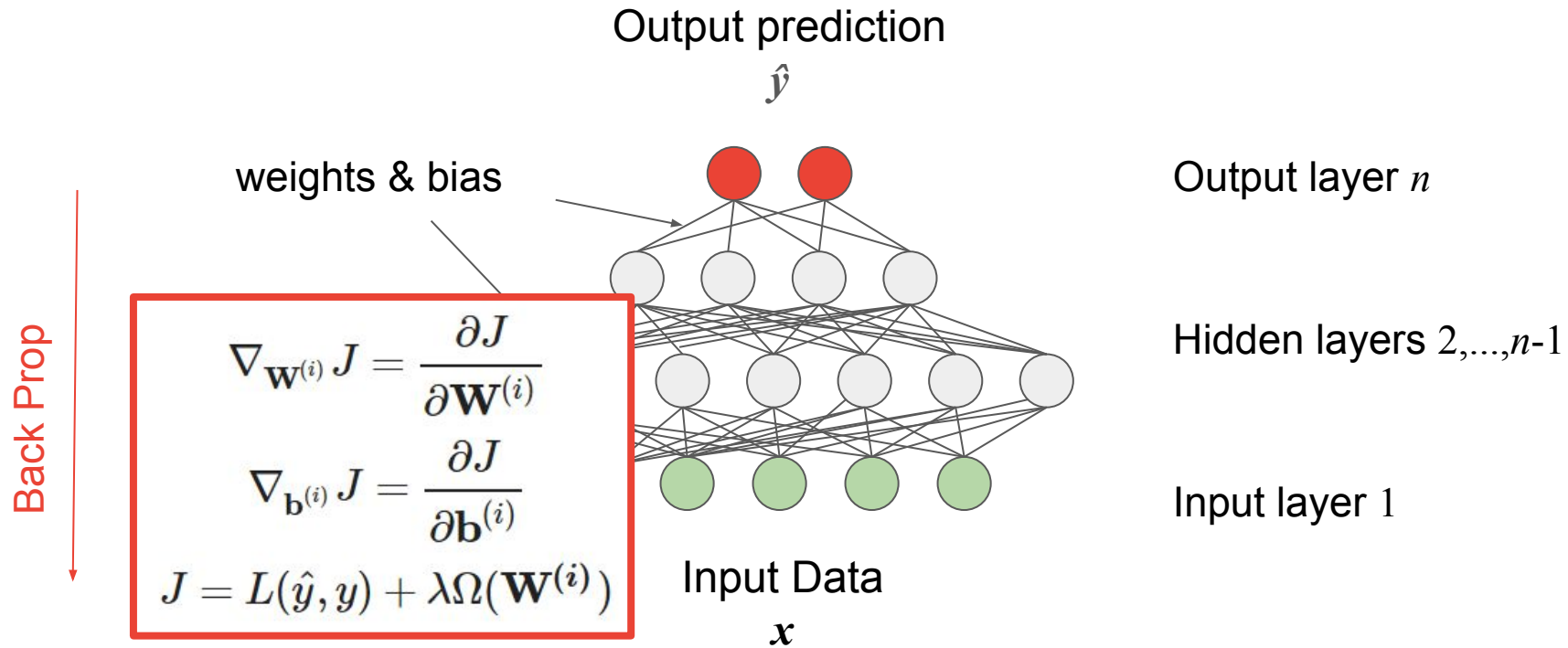Input layer $1$

Input Data
$x$

# Essential Parts of Training a Neural Net

- Loss function

- Optimization

- Gradient Computation  ⬅ **Back-propagation**

# Outputs of Back-Propagation

Output prediction
$\hat{y}$

weights & bias

Output layer $n$

$$\nabla_{\mathbf{W}^{(i)}} J = \frac{\partial J}{\partial \mathbf{W}^{(i)}}$$

Hidden layers $2,...,n\text{-}1$

$$\nabla_{\mathbf{b}^{(i)}} J = \frac{\partial J}{\partial \mathbf{b}^{(i)}}$$

Input layer $1$

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

Input Data
$x$

Back Prop

# The Gradient

General gradient of $f$ with respect to $\boldsymbol{x}$, but ignoring $\boldsymbol{y}$

$$\nabla_{\boldsymbol{x}} f(\boldsymbol{x}, \boldsymbol{y})$$

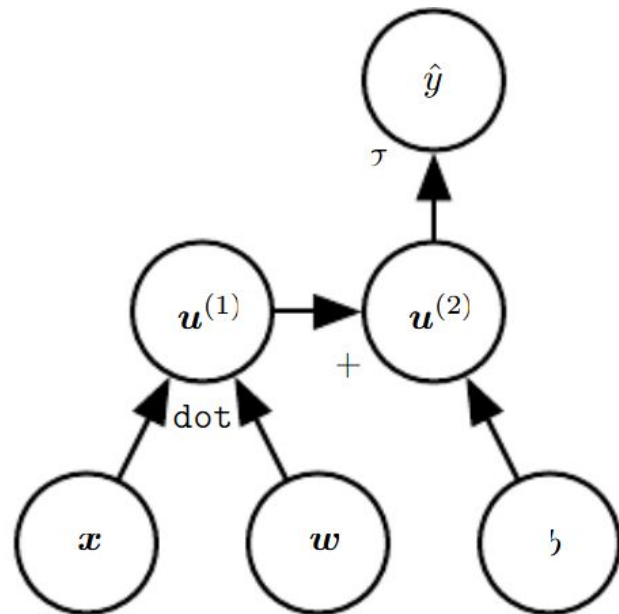Gradient of the loss $J$ with respect to parameters $\boldsymbol{\theta}$:

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$
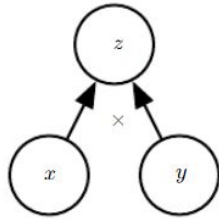
# Computational Graphs

# Computational Graph

- Efficient method of decomposing forward/backward propagations
- Graph of nodes and edges, *G(V,E)*
- Variables are nodes
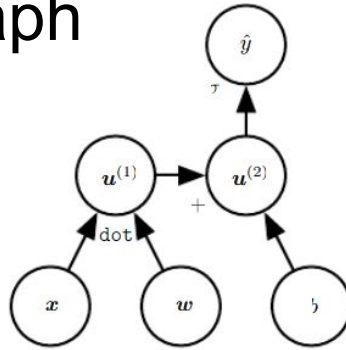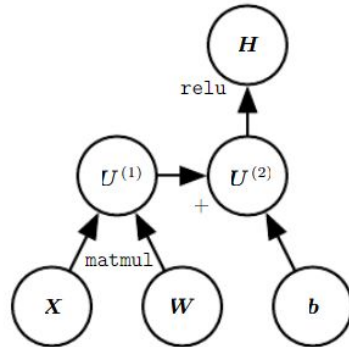- Operations are edges

$$\hat{y} = \sigma(\boldsymbol{w}^\mathsf{T}\boldsymbol{x} + b)$$

# Computational Graph

(a) $z = xy$

(b) $\hat{y} = \sigma(\boldsymbol{w}^\mathsf{T}\boldsymbol{x} + b)$

(c) $\boldsymbol{H} = \max\{0, \boldsymbol{X}\boldsymbol{W} + \boldsymbol{b}\}$

(d) Two outputs:

- $\hat{y} = \boldsymbol{x}^\mathsf{T}\boldsymbol{w}$
- $\lambda\sum_i w_i^2$

# Chain Rule, Jacobians, and Gradients

# Chain Rule of Calculus

$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$$

Knowing the instantaneous rate of change of $z$ relative to $y$ and that of $y$ relative to $x$ allows one to calculate the instantaneous rate of change of $z$ relative to $x$ as the product of the two rates of change.

George F. Simmons: "*if a car travels twice as fast as a bicycle and the bicycle is four times as fast as a walking man, then the car travels 2 × 4 = 8 times as fast as the man.*"

# The Jacobian (from before…)

If the function has $m$ inputs and $n$ outputs:

$$f : \mathbb{R}^m \to \mathbb{R}^n$$

The derivative $m$-by-$n$ matrix is called the Jacobian $J \in \mathbb{R}^{n \times m}$ of $f$:

$$\mathbf{J} = \begin{bmatrix} \frac{\partial}{\partial x_1} f(x)_1 & \frac{\partial}{\partial x_2} f(x)_1 & \cdots & \frac{\partial}{\partial x_m} f(x)_1 \\ \frac{\partial}{\partial x_1} f(x)_2 & \frac{\partial}{\partial x_2} f(x)_2 & \cdots & \frac{\partial}{\partial x_m} f(x)_2 \\ \cdots & & & \\ \frac{\partial}{\partial x_1} f(x)_n & \frac{\partial}{\partial x_2} f(x)_n & \cdots & \frac{\partial}{\partial x_m} f(x)_n \end{bmatrix}$$

# Jacobian with matrix input output

Let's take a 2 x 2 matrix $\mathbf{X}$:

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix}$$

On which an elementwise operation is perfomed in the forward pass: $a_{ij} = \sigma(x_{ij})$

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

# Jacobian Example

For the backward pass, we need the Jacobian $\frac{\partial A}{\partial x_{ij}}$:

$$\frac{\partial A}{\partial X} = \begin{bmatrix} \dfrac{\partial a_{11}}{\partial x_{11}} & \dfrac{\partial a_{12}}{\partial x_{11}} & \dfrac{\partial a_{21}}{\partial x_{11}} & \dfrac{\partial a_{22}}{\partial x_{11}} \\[2ex] \dfrac{\partial a_{11}}{\partial x_{12}} & \dfrac{\partial a_{12}}{\partial x_{12}} & \dfrac{\partial a_{21}}{\partial x_{12}} & \dfrac{\partial a_{22}}{\partial x_{12}} \\[2ex] \dfrac{\partial a_{11}}{\partial x_{21}} & \dfrac{\partial a_{12}}{\partial x_{21}} & \dfrac{\partial a_{21}}{\partial x_{21}} & \dfrac{\partial a_{22}}{\partial x_{21}} \\[2ex] \dfrac{\partial a_{11}}{\partial x_{22}} & \dfrac{\partial a_{12}}{\partial x_{22}} & \dfrac{\partial a_{21}}{\partial x_{22}} & \dfrac{\partial a_{22}}{\partial x_{22}} \end{bmatrix}$$

# Jacobian Example

For most operations, the non-diagonal terms reduce to 0.

$$\frac{\partial A}{\partial X} = \begin{bmatrix} \frac{\partial a_{11}}{\partial x_{11}} & 0 & 0 & 0 \\ 0 & \frac{\partial a_{12}}{\partial x_{12}} & 0 & 0 \\ 0 & 0 & \frac{\partial a_{21}}{\partial x_{21}} & 0 \\ 0 & 0 & 0 & \frac{\partial a_{22}}{\partial x_{22}} \end{bmatrix}$$

Hence, the Jacobian can be written as:

$$\frac{\partial A}{\partial X} = \text{diag}(f'(X))$$

where:

$$A = f(X)$$

# Jacobian

For the nonlinear activation function:

$$f(\mathbf{X}) = \tanh(\mathbf{X})$$

And its derivative:

$$f'(\mathbf{X}) = 1 - \tanh^2(\mathbf{X})$$

$$\frac{\partial \mathbf{A}}{\partial \mathbf{X}} = \begin{bmatrix} 1 - \tanh^2(x_{11}) & 0 & 0 & 0 \\ 0 & 1 - \tanh^2(x_{12}) & 0 & 0 \\ 0 & 0 & 1 - \tanh^2(x_{21}) & 0 \\ 0 & 0 & 0 & 1 - \tanh^2(x_{22}) \end{bmatrix}$$

# Generalizing the chain rule to vectors

Suppose

- $x \in \mathbb{R}^m$ and $y \in \mathbb{R}^n$ and
- $y = g(x) : \mathbb{R}^m \to \mathbb{R}^n$
- $z = f(y) : \mathbb{R}^n \to \mathbb{R}$ ← Loss function

Then

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

In vector notation:

$$\nabla_x z = \left( \frac{\partial y}{\partial x} \right)^{\mathsf{T}} \nabla_y z$$

Gradient of loss
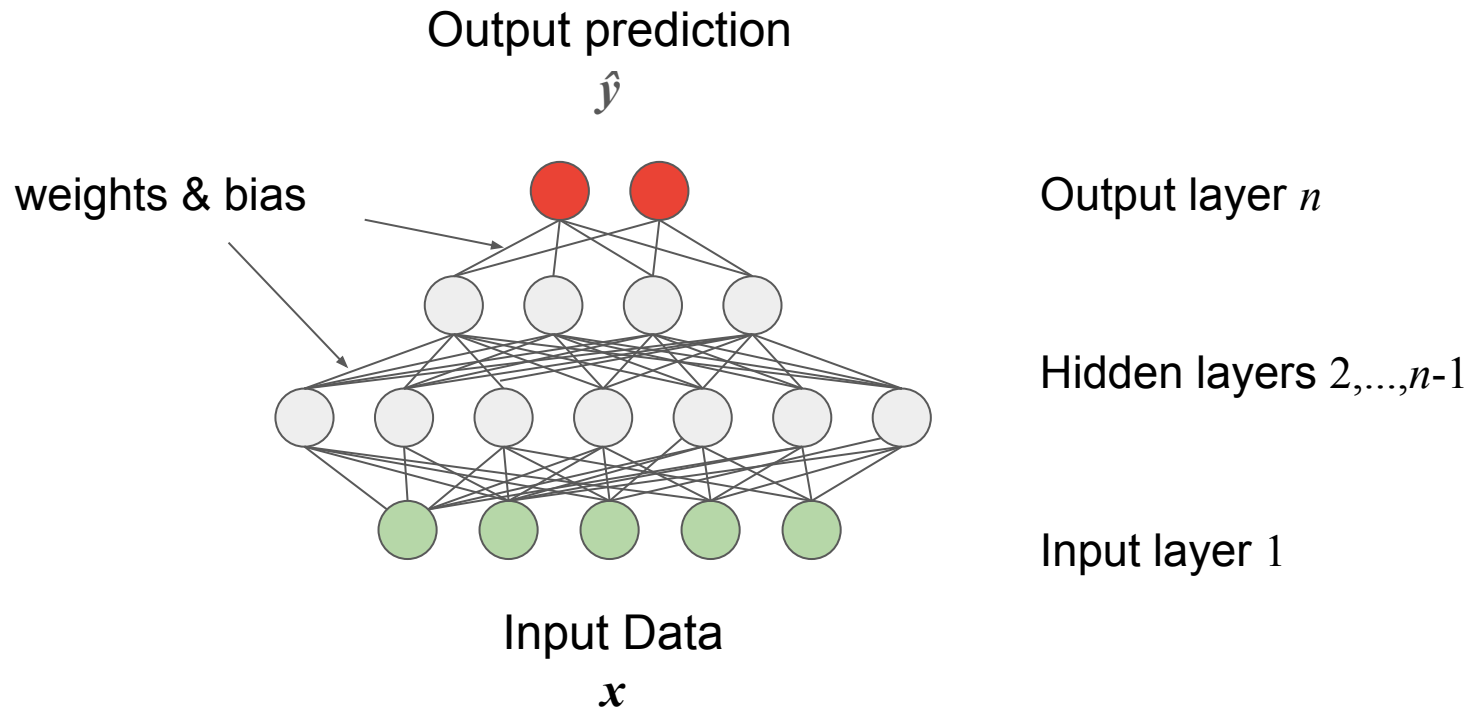
$n \times m$ Jacobian of $g$

# Generalizing the chain rule to tensors

- Conceptually the same idea as vectors, but arbitrary dimensionality: **multiply the Jacobian by the Gradient**.

- Replace scalar indices with tuples e.g.,
$$i = \{(0,0,0),(0,0,1),\ldots\}$$

- For all possible index tuples, $i$, $(\nabla_X z)_i$ gives $\frac{\partial z}{\partial X_i}$, is equivalent to integer index $i$, $(\nabla_x z)_i$ gives $\frac{\partial z}{\partial x_i}$

- Given $Y = g(X)$, and $z = f(Y)$, then

$$\nabla_X z = \sum_j (\nabla_X Y_j) \frac{\partial z}{\partial Y_j}$$

Jacobian        Gradient

# Backprop: Recursive Chain Rule

# General Architecture of Feedforward Nets



Output prediction
$\hat{y}$

weights & bias

Output layer $n$

Hidden layers $2,...,n\text{-}1$

Input layer $1$

Input Data
$x$

# Recursive Chain Rule

**Forward prop**: Apply the function $f$ progressively from the input $w$ to the output $z$

Output prediction

$z$

$f$

$y$

$f$

$x$

$f$

$w$

Input value

$z = f(y)$

$y = f(x)$

$x = f(w)$

$z = f(f(f(w)))$
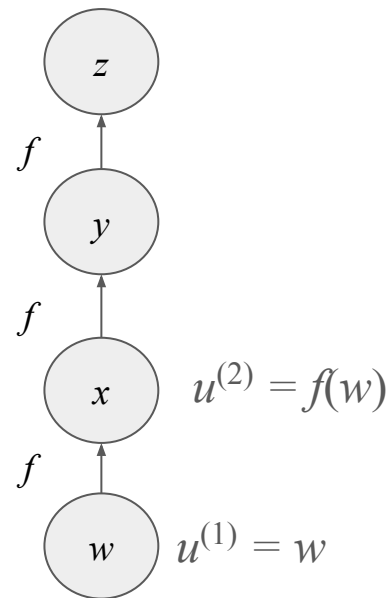
# Forward Propagation Algorithm

**Algorithm 6.2** Basic Forward-propagation

---

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

---

$Pa(x) = \{w\}$

$u^{(2)} = f(w)$

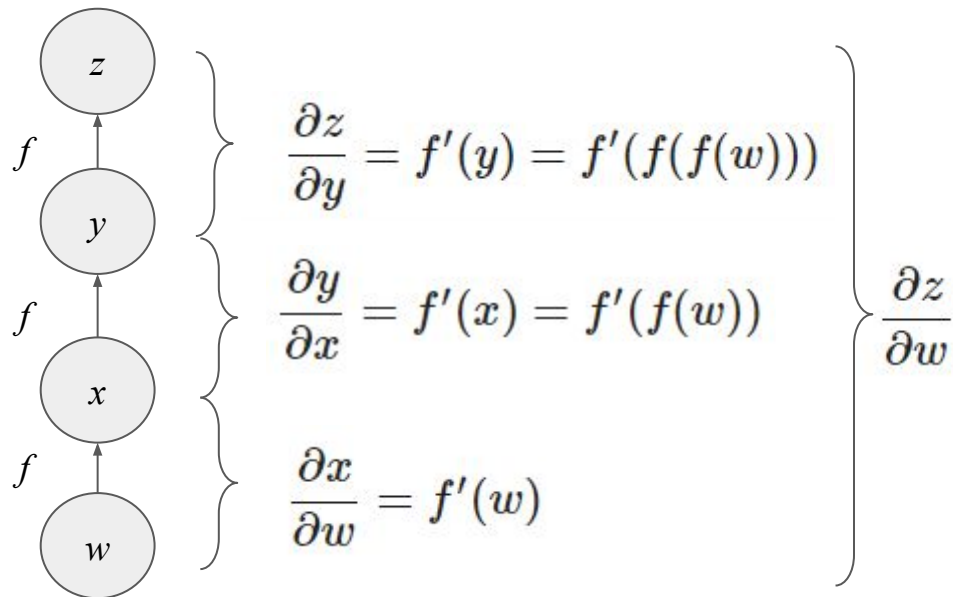$u^{(1)} = w$

# Recursive Chain Rule

**Backprop**: Obtain the gradient of the weights and biases by recursively applying chain rule from the Loss function down to the inputs

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

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

$$= f'(y)f'(x)f'(w)$$

$$= f'(f(f(w)))f'(f(w))f'(w)$$



$z$

$f$

$y$

$\dfrac{\partial z}{\partial y} = f'(y) = f'(f(f(w)))$

$f$

$x$

$\dfrac{\partial y}{\partial x} = f'(x) = f'(f(w))$

$\dfrac{\partial z}{\partial w}$

$f$

$w$

$\dfrac{\partial x}{\partial w} = f'(w)$

# Simplified Backprop Algo

**Algorithm 6.2** Simplified Back-propagation

---

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`$[u^{(n)}] \leftarrow 1$

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

**end for**

**return** $\{$`grad_table`$[u^{(i)}] \mid i = 1, \ldots, n_i\}$

---

# Improved FP/BP on a fully connected network

Apply the following enhancements to the previous algorithms:

- Add in the Loss function $L(\hat{y}, y)$
- Output prediction $\hat{y}$
- Include regularization $\lambda\Omega(\theta)$
- Compute weights and biases at each layer $W^{(i)}, b^{(i)}$

Regularized Loss function

$J = L(\hat{y}, y) - \lambda\Omega(\theta)$
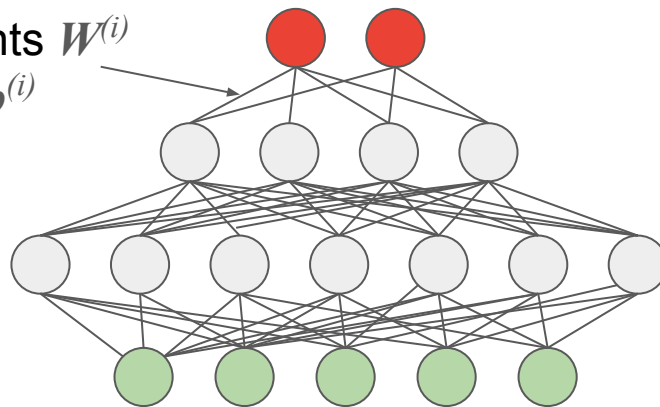
# General Architecture of Feedforward Nets

True label
$y$

Output prediction
$\hat{y}$

Weights $W^{(i)}$
Bias $b^{(i)}$

Output layer $n$

Hidden layers $2,...,n$-1

Regularized Loss function

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

Input layer $1$

Input Data
$x$

# Improved forward propagation

**Algorithm 6.3** Improved Forward-propagation

---

**Require:** Network depth, $l$

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

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

**Require:** $\boldsymbol{x}$, the input to process

**Require:** $\boldsymbol{y}$, the target output

$\quad \boldsymbol{h}^{(0)} = \boldsymbol{x}$

$\quad$ **for** $k = 1, \ldots, l$ **do**

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

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

$\quad$ **end for**

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

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

---

# Improved Back-propagation

**Algorithm 6.4** Improved Back-propagation

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

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

**for** $k = l, l-1, \ldots, 1$ **do**

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

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

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

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

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

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

**end for**

---

Linear Activation:

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

Layer output with non-linear activation:

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

# Structuring Backprop with a Computational Graph

- To reduce the runtime complexity from exponential to linear time, expand the computational graph with additional nodes for back propagation
- Note the use of the chain rule
- Graph is populated with values as soon as the parent nodes are available



$$\frac{dz}{dw} = \text{Jacobian} \times \text{Gradient} = \frac{dz}{dx} \times \frac{dx}{dw}$$

# Readings

- Goodfellow - Chapter 7
- Goodfellow - Chapter 8