

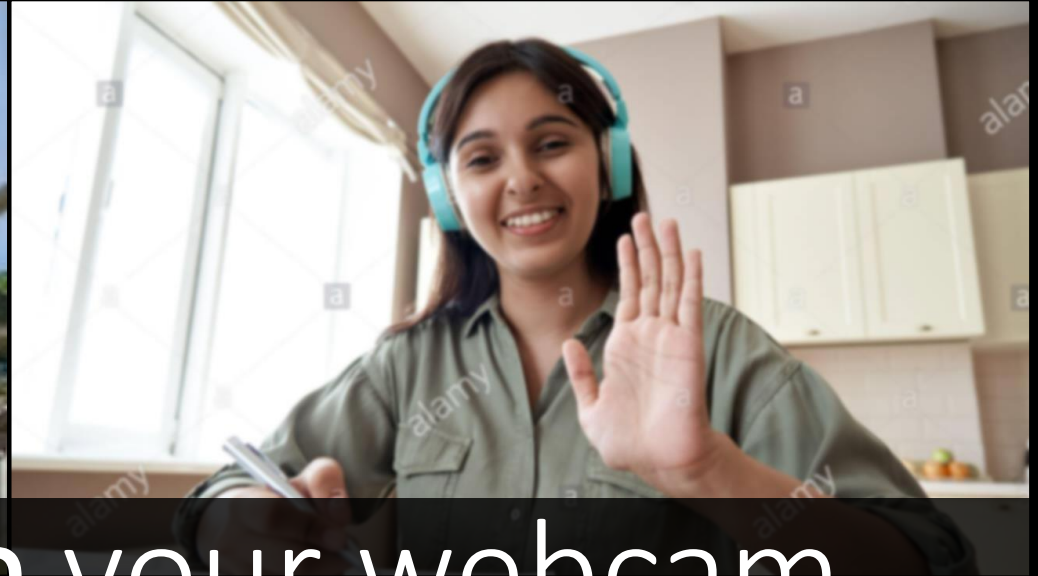
Ask Me Anything (AMA)

CS 3244
Machine Learning

9
B



NUS | Computing



Please turn on your webcam



Mystery Student

Week 09A: Lecture Outline

1. Math Notation Primer
2. Backprop Revisiting
3. AMA

Average difference metrics for test dataset

Mean Absolute Error (MAE)

$$MAE = \frac{1}{m} \sum_{j=1}^m |\hat{y}_j - y_j|$$

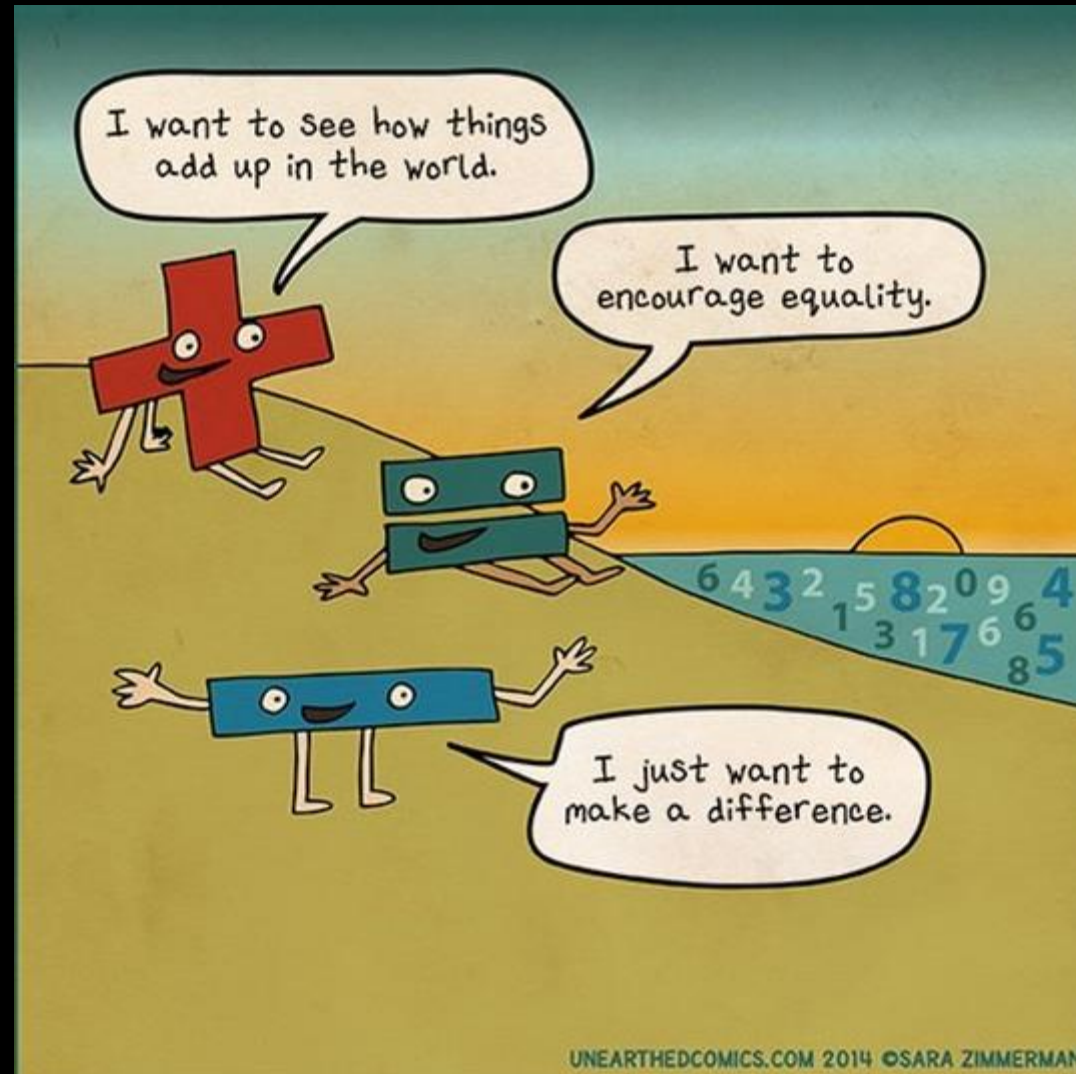
Mean Squared Error (MSE)

$$MSE = \frac{1}{m} \sum_{j=1}^m (\hat{y}_j - y_j)^2$$

Root Mean Squared Error (RMSE)

$$RMSE = \sqrt{\frac{1}{m} \sum_{j=1}^m (\hat{y}_j - y_j)^2}$$

MSE and RMSE penalize larger differences **more** than MAE



Why Math in Machine Learning?

Idea \rightarrow Math \rightarrow Coding

Math enables **deterministic** (algorithmic) execution (implementation)

Math helps to **check idea** is correct (logically consistent)

Math helps to **check coding** is correct (bug free)

Notation

n = Number of features in \mathbf{x}
 m = Number of instances in dataset

- **Scalar**: not bolded, lower case

x

- **Vector**: bolded, lower case

$$\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

- **Matrix**: bolded, upper case

$$\mathbf{X} = \begin{pmatrix} x_{11} & \cdots & x_{m1} \\ \vdots & \ddots & \vdots \\ x_{1n} & \cdots & x_{mn} \end{pmatrix}$$

Functions with Vectors and Matrices

- Scalar-by-scalar:
 - $y(x) = wx$ for scaling input
- Scalar-by-vector:
 - $y(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} = \mathbf{w}^\top \mathbf{x} = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = w_1x_1 + w_1x_2 + w_2x_1 + w_2x_2$ for weighted sum
- Vector-by-vector:
 - $\mathbf{y}(\mathbf{x}) = w\mathbf{x} = w \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} wx_1 \\ wx_2 \end{pmatrix}$ for scaled outputs (same weight)

Functions with Vectors and Matrices

- Vector-by-matrix:

- $\mathbf{y}(\mathbf{X}) = \mathbf{W} \circ \mathbf{X} = \begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix} \begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{pmatrix} = \begin{pmatrix} w_{11}x_{11} & w_{12}x_{12} \\ w_{21}x_{21} & w_{22}x_{22} \end{pmatrix}$

- Using **Hadamard** product \circ for element-wise multiplication

- For backpropagation (see later), convenient multiplication of corresponding elements between matrices

- Matrix-by-matrix:

- $\mathbf{Y}(\mathbf{X}) = \mathbf{W} * \mathbf{X} = \begin{pmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{pmatrix} \begin{pmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{pmatrix}$
$$= \begin{pmatrix} w_{11}x_{11} + w_{12}x_{12} + w_{21}x_{21} + w_{22}x_{22} & w_{11}x_{12} + w_{12}x_{13} + w_{21}x_{22} + w_{22}x_{23} \\ w_{11}x_{21} + w_{12}x_{22} + w_{21}x_{31} + w_{22}x_{32} & w_{11}x_{22} + w_{12}x_{23} + w_{21}x_{32} + w_{22}x_{33} \end{pmatrix}$$

- Using **Convolution** operator $*$ for element-wise multiplication then sum [W08b]

- For computer vision filters (kernels)

Weighted Sum

Summation Series = Scalar

$$\sum_{r=0}^n w_r x_r$$
$$w_1 x_1 + \cdots + w_r x_r + \cdots + w_n x_n$$

Vector Dot Product = Scalar

$$\mathbf{w} \cdot \mathbf{x} = \begin{pmatrix} w_1 \\ \vdots \\ w_r \\ \vdots \\ w_n \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ \vdots \\ x_r \\ \vdots \\ x_n \end{pmatrix}$$

Transposed Vector Multiplication = Scalar

$$\mathbf{w}^T \mathbf{x} = (w_1 \quad \cdots \quad w_r \quad \cdots \quad w_n) \begin{pmatrix} x_1 \\ \vdots \\ x_r \\ \vdots \\ x_n \end{pmatrix}$$

Transposed Matrix Multiplication = Vector

$$\mathbf{W}^T \mathbf{x} = \begin{pmatrix} w_{11} & \cdots & w_{1r} & \cdots & w_{1n} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ w_{r1} & \cdots & w_{rr} & \cdots & w_{rn} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ w_{n1} & \cdots & w_{nr} & \cdots & w_{nn} \end{pmatrix}^T \begin{pmatrix} x_1 \\ \vdots \\ x_r \\ \vdots \\ x_n \end{pmatrix}$$

Gradient

- **Derivative:** d

- $\frac{dy}{dx}$ is the derivative of y relative to x

- **Partial derivative:** ∂

- $\frac{\partial y}{\partial x_1}$ is the derivative of y relative to x_1
- But y also depends on other variables (e.g., x_2 so, we can also calculate $\frac{\partial y}{\partial x_2}$)

- **Gradient:** ∇

- To calculate the derivative relative to *all* x_1 and x_2 together
- $\nabla y(\mathbf{x})$ is the gradient of y relative to all variables $\mathbf{x} = (x_1, \dots, x_n)^\top$

Vector in denominator means
Derivative for each variable is
put in separate, corresponding
variable dimension

$$\nabla y(\mathbf{x}) = \frac{\partial y}{\partial \mathbf{x}} = \begin{pmatrix} \partial y / \partial x_1 \\ \vdots \\ \partial y / \partial x_n \end{pmatrix} = \left(\frac{\partial y}{\partial x_1} \quad \dots \quad \frac{\partial y}{\partial x_n} \right)^\top$$

- Assumes Cartesian coordinates (linear, orthogonal)

Matrix Calculus – not in exam

n = Number of features in \mathbf{x}
 m = Number of instances in dataset
 N = Number of y prediction tasks

Scalar-by-Vector (= 1D Vector)

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y}{\partial x_1} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{pmatrix}$$

Scalar-by-Matrix (= 2D Matrix)

$$\frac{\partial y}{\partial \mathbf{X}} = \begin{pmatrix} \frac{\partial y}{\partial x_{11}} & \dots & \frac{\partial y}{\partial x_{1m}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial x_{n1}} & \dots & \frac{\partial y}{\partial x_{nm}} \end{pmatrix}$$

Vector-by-Vector (= 2D Matrix)

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_N}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_n} & \dots & \frac{\partial y_N}{\partial x_n} \end{pmatrix}$$

Vector-by-Matrix (= 3D Matrix)

$$\frac{\partial \mathbf{y}}{\partial \mathbf{X}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_{11}} & \dots & \frac{\partial y_1}{\partial x_{1m}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_{n1}} & \dots & \frac{\partial y_1}{\partial x_{nm}} \end{pmatrix} \dots \begin{pmatrix} \frac{\partial y_N}{\partial x_{11}} & \dots & \frac{\partial y_N}{\partial x_{1m}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_N}{\partial x_{n1}} & \dots & \frac{\partial y_N}{\partial x_{nm}} \end{pmatrix}$$

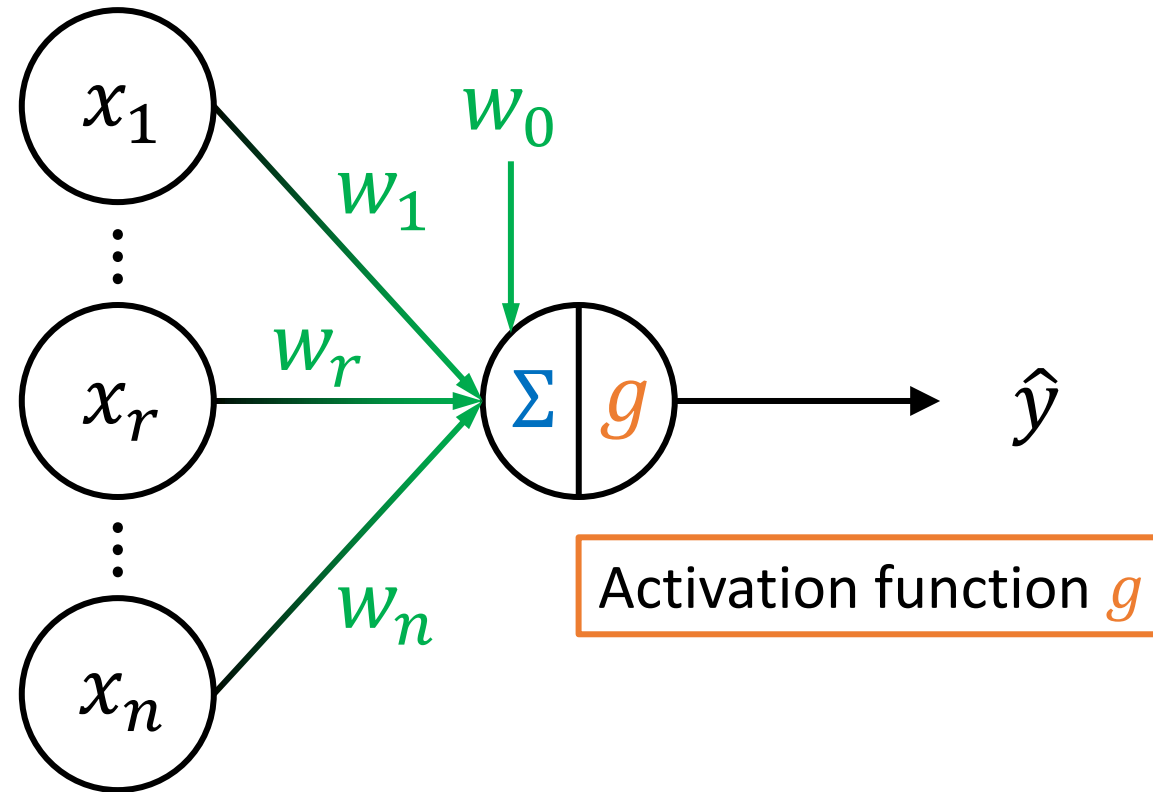
Along 3rd dimension

This math informs what matrix **shapes** you need to implement

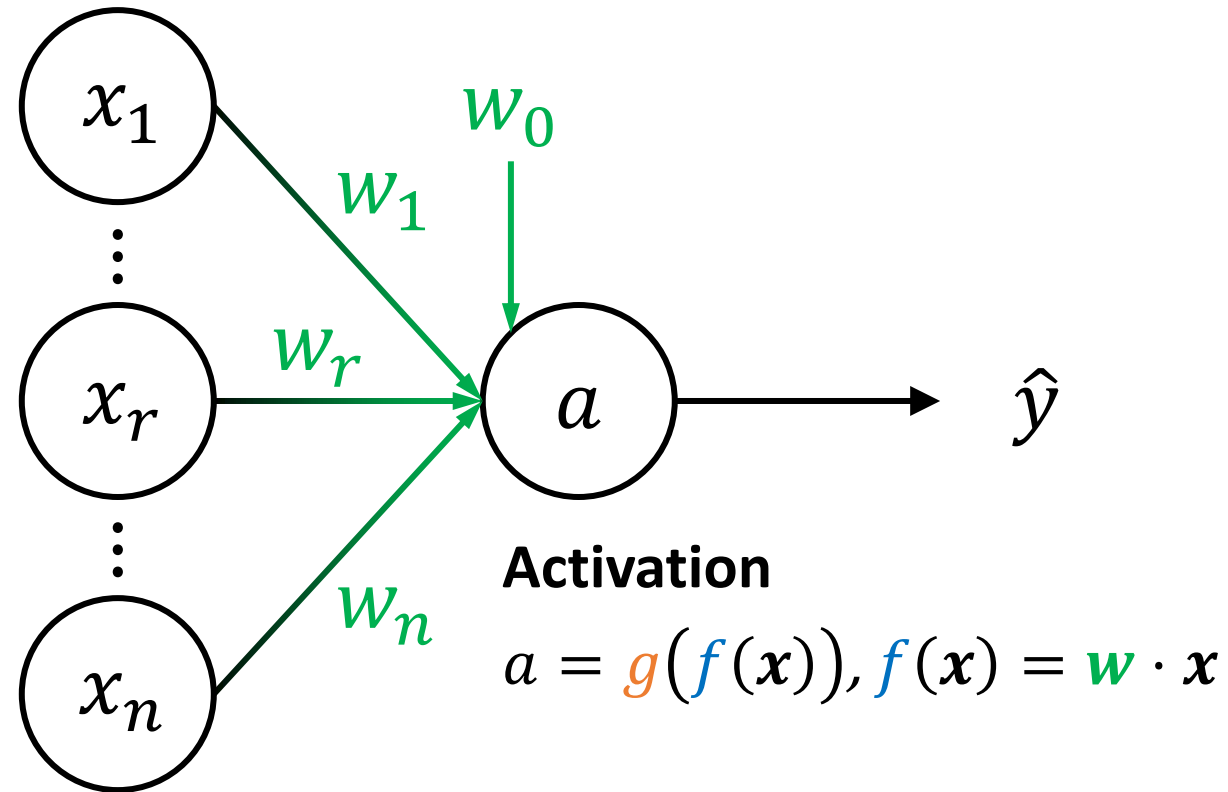


Backpropagation “backprop”

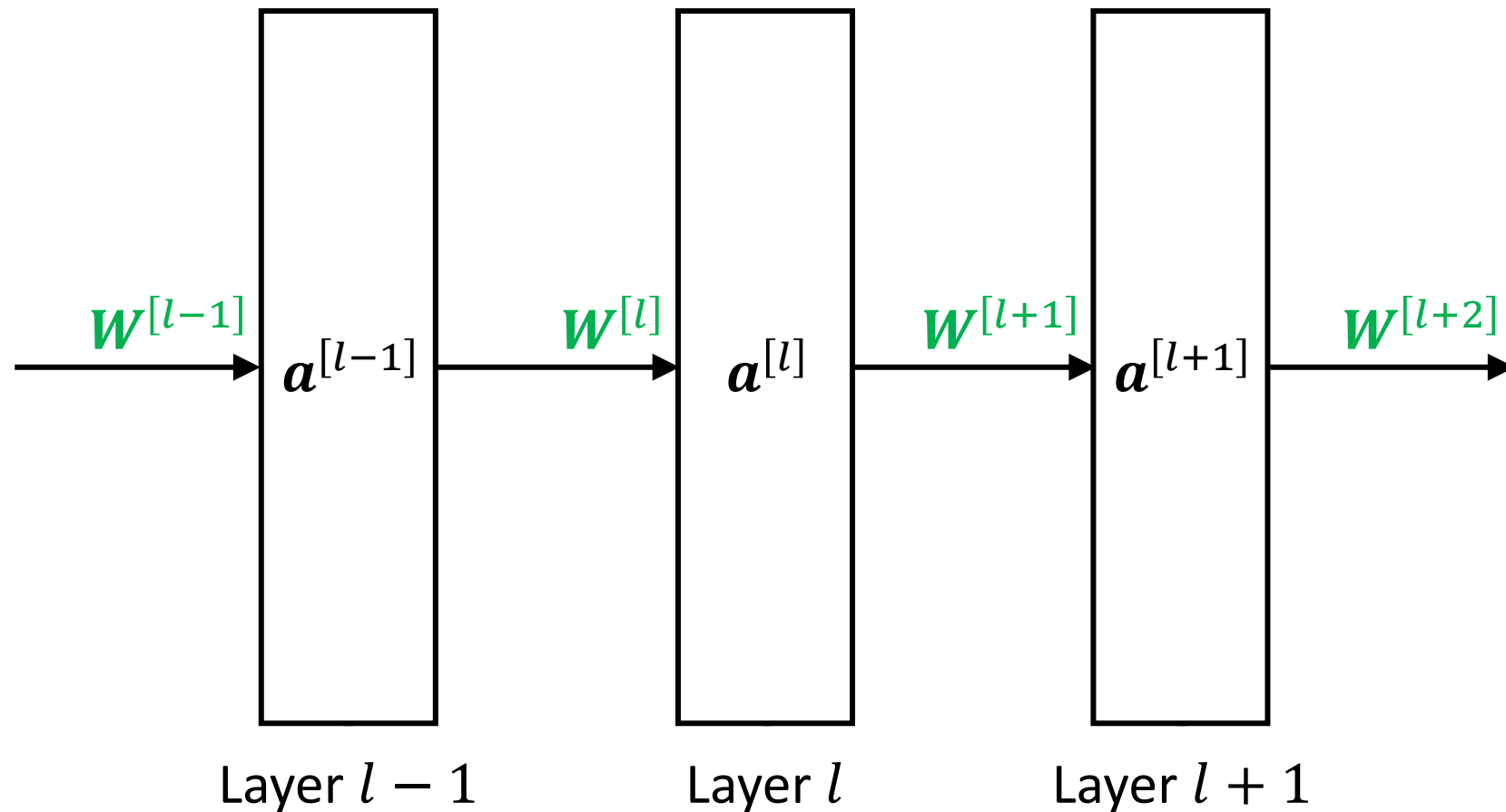
Single-Layer Perceptron



Single-Layer Perceptron



Neural Network



Layer Activation

$$a = g(f(x)), f(x) = w \cdot x$$

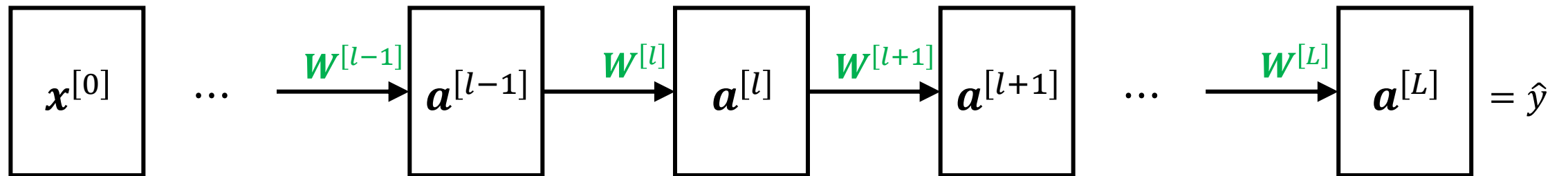
Single-Layer
Perceptron

$$\underset{\substack{\text{Layer } l \\ \text{Activations}}}{\mathbf{a}^{[l]}} = \underset{\substack{\text{Layer } l \\ \text{Activation} \\ \text{Function}}}{g^{[l]}} \left(\left(\underset{\substack{\text{Layer } l \\ \text{Weights}}}{\mathbf{W}^{[l]}} \right)^T \underset{\substack{\text{Layer } l - 1 \\ \text{Activations}}}{\mathbf{a}^{[l-1]}} \right)$$

Layer l in
Neural Network

Forward Propagation

$$g^{[1]} \left((W^{[1]})^\top x^{[0]} \right) = a^{[1]} \quad g^{[l]} \left((W^{[l]})^\top a^{[l-1]} \right) = a^{[l]} \quad g^{[L]} \left((W^{[L]})^\top a^{[L-1]} \right) = a^{[L]}$$



$$\hat{y}(x) = g^{[L]} \left((W^{[L]})^\top g^{[L-1]} \left(\dots \left(g^{[l]} \left((W^{[l]})^\top g^{[l-1]} \left(\dots \left(g^{[1]} \left((W^{[1]})^\top x^{[0]} \right) \right) \right) \right) \right) \right) \right)$$

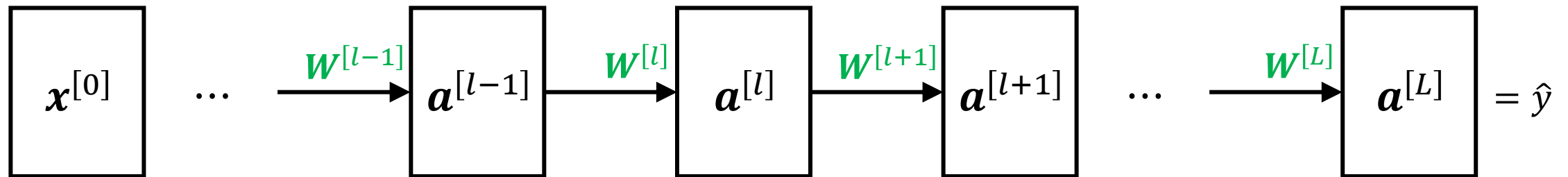
Forward Propagation

$$\mathbf{a}^{[l]} \equiv g^{[l]}(f^{[l]}), \quad f^{[l]} \equiv (W^{[l]})^\top \mathbf{a}^{[l]}$$

$$g^{[1]}(f^{[1]}(\mathbf{x}^{[0]})) = \mathbf{a}^{[1]}$$

$$g^{[l]}(f^{[l]}(\mathbf{a}^{[l-1]})) = \mathbf{a}^{[l]}$$

$$g^{[L]}(f^{[L]}(\mathbf{a}^{[L-1]})) = \mathbf{a}^{[L]}$$



$$\hat{y}(\mathbf{x}) = g^{[L]}(f^{[L]}(g^{[L-1]}(\dots(g^{[l]}(f^{[l]}(g^{[l-1]}(\dots(g^{[1]}(f^{[1]}(\mathbf{x}^{[0]}))))))))))$$

Even more Chain Rule:

Gradient of Neural Network

$$\hat{y}(\mathbf{x}) = g^{[L]}(f^{[L]}(g^{[L-1]}(\dots(g^{[l]}(f^{[l]}(g^{[l-1]}(\dots(g^{[1]}(f^{[1]}(\mathbf{x}^{[0]}))))))))))$$

Gradient relative to \mathbf{W}

$$\hat{y}'(\mathbf{W}^{[L-1]}) = \frac{\partial g^{[L]}}{\partial \mathbf{W}^{[L-1]}} = \frac{\partial f^{[L]}}{\partial \mathbf{W}^{[L-1]}} \boxed{\frac{\partial g^{[L]}}{\partial f^{[L]}}} - \delta^{[L-1]}$$

Reference

$$\mathbf{a}^{[l]} = g^{[l]}(f^{[l]})$$

$$\hat{y}'(\mathbf{W}^{[l+1]}) = \frac{\partial g^{[L]}}{\partial \mathbf{W}^{[l+1]}} = \frac{\partial f^{[l+1]}}{\partial \mathbf{W}^{[l+1]}} \boxed{\frac{\partial g^{[l+1]}}{\partial f^{[l+1]}} \dots \frac{\partial f^{[L]}}{\partial g^{[L-1]}} \frac{\partial g^{[L]}}{\partial f^{[L]}}} - \delta^{[l+1]}$$

$$f^{[l]} = (\mathbf{W}^{[l]})^T \mathbf{a}^{[l-1]}$$

$$\hat{y}'(\mathbf{W}^{[l]}) = \frac{\partial g^{[L]}}{\partial \mathbf{W}^{[l]}} = \frac{\partial f^{[l]}}{\partial \mathbf{W}^{[l]}} \frac{\partial g^{[l]}}{\partial f^{[l]}} \underbrace{\frac{\partial f^{[l+1]}}{\partial g^{[l]}} \frac{\partial g^{[l+1]}}{\partial f^{[l+1]}} \dots \frac{\partial f^{[L]}}{\partial g^{[L-1]}} \frac{\partial g^{[L]}}{\partial f^{[L]}}}_{\delta^{[l+1]}}$$

Recursive

$$\hat{y}'(\mathbf{W}^{[l]}) = \frac{\partial g^{[L]}}{\partial \mathbf{W}^{[l]}} = \frac{\partial f^{[l]}}{\partial \mathbf{W}^{[l]}} \frac{\partial g^{[l]}}{\partial f^{[l]}} \frac{\partial f^{[l+1]}}{\partial g^{[l]}} \delta^{[l+1]}$$

$$\hat{y}'(\mathbf{W}^{[1]}) = \frac{\partial g^{[L]}}{\partial \mathbf{W}^{[1]}} = \frac{\partial f^{[1]}}{\partial \mathbf{W}^{[1]}} \frac{\partial g^{[1]}}{\partial f^{[1]}} \dots \frac{\partial g^{[l]}}{\partial f^{[l]}} \frac{\partial f^{[l+1]}}{\partial g^{[l]}} \frac{\partial g^{[l+1]}}{\partial f^{[l+1]}} \dots \frac{\partial f^{[L]}}{\partial g^{[L-1]}} \frac{\partial g^{[L]}}{\partial f^{[L]}}$$

Even more Chain Rule:

Gradient of Neural Network

$$\hat{y}(\mathbf{x}) = g^{[L]}(f^{[L]}(g^{[L-1]}(\dots(g^{[l]}(f^{[l]}(g^{[l-1]}(\dots(g^{[1]}(f^{[1]}(\mathbf{x}^{[0]}))))))))))$$

Gradient relative to $\mathbf{W}^{[l]}$

$$\hat{y}'(\mathbf{W}^{[l]}) = \frac{\partial g^{[L]}}{\partial \mathbf{W}^{[l]}} = \frac{\partial f^{[l]}}{\partial \mathbf{W}^{[l]}} \frac{\partial g^{[l]}}{\partial f^{[l]}} \frac{\partial f^{[l+1]}}{\partial g^{[l]}} \delta^{[l+1]}$$

Reference

$$\mathbf{a}^{[l]} = g^{[l]}(f^{[l]})$$

$$f^{[l]} = (\mathbf{W}^{[l]})^\top \mathbf{a}^{[l-1]}$$

$$\frac{\partial f^{[l]}}{\partial \mathbf{W}^{[l]}} = \mathbf{a}^{[l-1]} \quad \frac{\partial g^{[l]}}{\partial f^{[l]}} = g'^{[l]}(f^{[l]}) \quad \frac{\partial f^{[l+1]}}{\partial g^{[l]}} = \frac{\partial f^{[l+1]}}{\partial \mathbf{a}^{[l]}} = \mathbf{W}^{[l+1]}$$

$$\hat{y}'(\mathbf{W}^{[l]}) = \mathbf{a}^{[l-1]} \left[g'^{[l]}(f^{[l]}) \mathbf{W}^{[l+1]} \delta^{[l+1]} \right] = \mathbf{a}^{[l-1]} \delta^{[l]}$$

$$\hat{y}'(\mathbf{W}^{[l]}) = \mathbf{a}^{[l-1]} \delta^{[l]}$$

$$\delta^{[l]} = g'^{[l]}(f^{[l]}) \mathbf{W}^{[l+1]} \delta^{[l+1]}$$

Recursive

Matrix multiplication to match shape (not in exam)

nRows × nCols

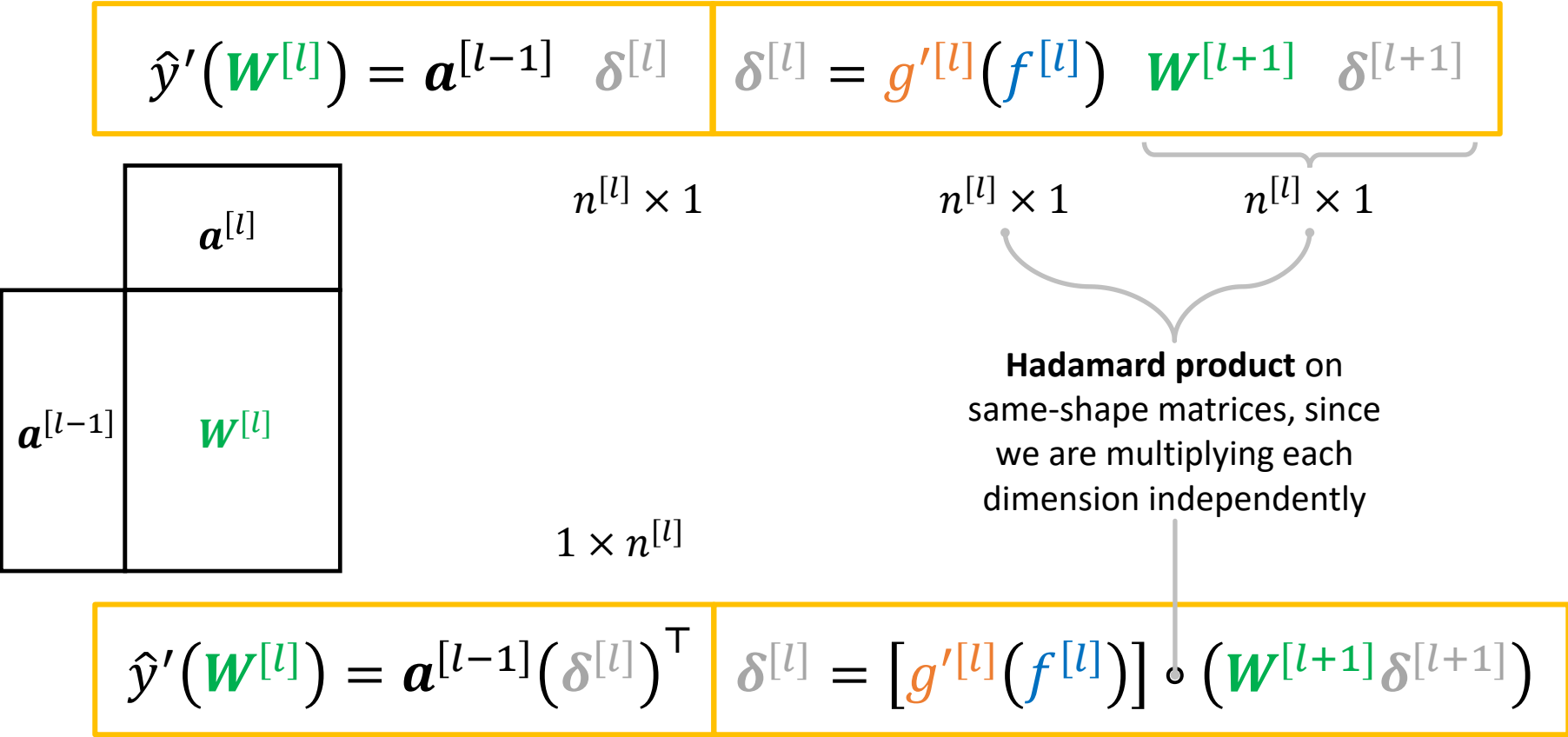
$n^{[l-1]} \times n^{[l]}$

$n^{[l-1]} \times 1$

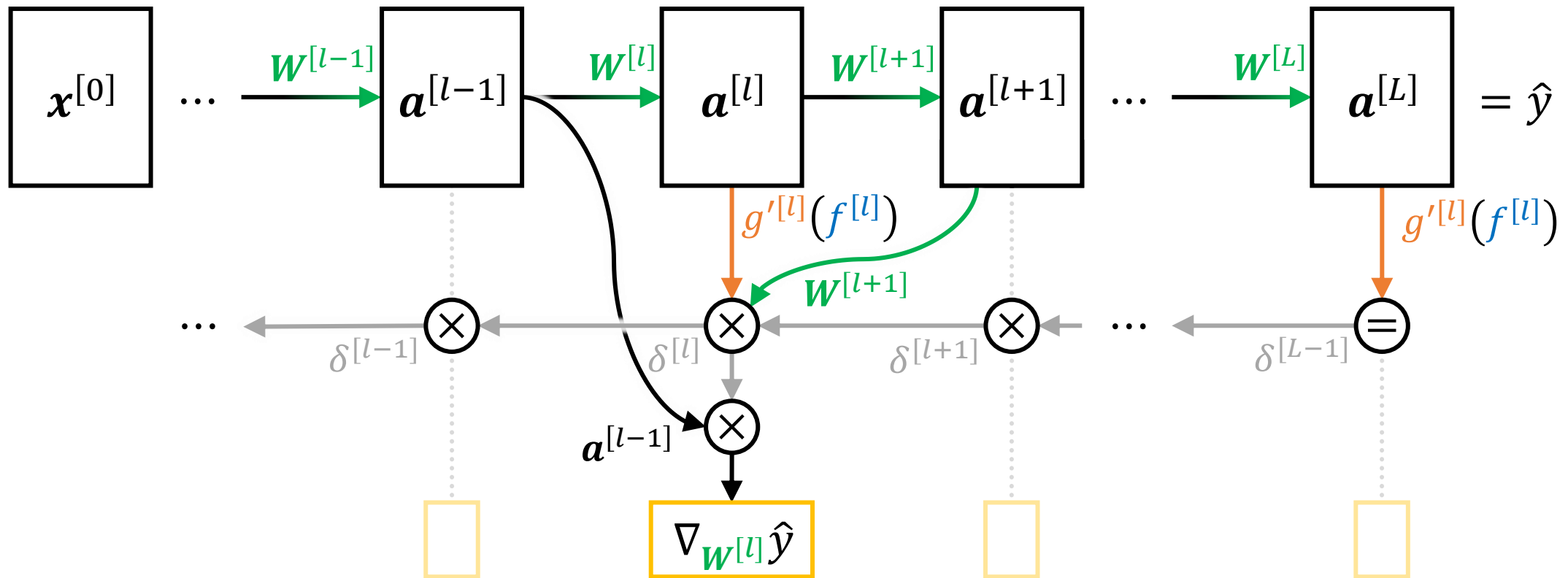
$n^{[l]} \times 1$

$n^{[l]} \times n^{[l+1]}$

$n^{[l+1]} \times 1$



Backward Propagation



$$\hat{y}'(W^{[l]}) = a^{[l-1]}(\delta^{[l]})^\top \quad \delta^{[l]} = [g'^{[l]}(f^{[l]})] \circ (W^{[l+1]} \delta^{[l+1]})$$

Backpropagation

Backpropagation **efficiently** computes the **gradient** by

- Avoiding **duplicate** calculations
- Not computing **unnecessary intermediate values**,
- Computing the **gradient** of ***each* layer**

Specifically, the gradient of the weighted input of each layer is calculated from back $[l + 1]$ to front $[l]$):

$$\hat{y}'(\mathbf{W}^{[l]}) = \mathbf{a}^{[l-1]}(\boldsymbol{\delta}^{[l]})^\top \quad \boldsymbol{\delta}^{[l]} = [g'^{[l]}(f^{[l]})] \circ (\mathbf{W}^{[l+1]} \boldsymbol{\delta}^{[l+1]})$$

Adapted from: <https://en.wikipedia.org/wiki/Backpropagation>

Auto Differentiation for Backprop

- Even with backprop, implementing the gradients is tedious
- Deep learning APIs have automated differentiation.
 - Tensor Flow [autodiff](#)
 - PyTorch [autograd](#)
 - Implement derivatives of many common functions
 - You just need to implement your layers and neurons; API will handle gradients
- **Caution**
 - If you want to implement **custom functions/layers** (not simple weighted sum)
 - They need to be **differentiable** to be able to calculate their **gradients**
 - Otherwise, backprop **cannot update** weights accurately



Questions!

