# *13. Matrix algebra and Backpropagation*

## M.A.Z. Chowdhury and M.A. Oehlschlaeger

Department of Mechanical, Aerospace and Nuclear Engineering
Rensselaer Polytechnic Institute
Troy, New York

*chowdm@rpi.edu*
*oehlsm@rpi.edu*

### MANE 4962 and 6962

## *Regular announcement*

☞ Quiz 4 Today.

☞ Quiz 4 is based on notes from Lecture 10.

☞ HW 4 is due March 2.

☞ Initial project proposal due March 2.

# Outline

- ☞ HW4, Linear Algebra, and backpropagation
- ☞ Further discussion about the initial and revised project proposal submission
- ☞ We solved regression problems using a fully connected neural network, $y \in \mathbb{R}$.
  outputs one number, we used mse.
- ☞ We solved a binary classification problem using a fully connected neural network
  $y \in [0, 1]$.
  outputs one number, we used binary crossentropy.
- ☞ Solve a multi-class classification problem using a fully connected neural networks
  $y \in \mathbb{R}^{K}$, K is the number of classes.
  outputs K numbers, we used categorical crossentropy.
- ☞ For project, We can solve regression problems using a fully connected neural network,
  $y \in \mathbb{R}^{u}$.
  outputs $u$ numbers. Use regression cost function, i.e., mse-like.

Representing Vectors and Matrices

Vector with $n$ elements, $x \in \mathbb{R}^n$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

Matrix with $m$ rows & $n$ columns, such that elements of it are real numbers, $A \in \mathbb{R}^{m \times n}$

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix} = \begin{bmatrix} | & | & & | \\ a^1 & a^2 & \cdots & a^n \\ | & | & & | \end{bmatrix} = \begin{bmatrix} - & a_1^T & - \\ - & a_2^T & - \\ & \vdots & \\ - & a_m^T & - \end{bmatrix}$$

# Identity & Diagonal Matrices

Identity matrix, $I \in \mathbb{R}^{n \times n}$, is a square matrix and zeros everywhere else.

$$I_{ij} = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}$$

Special Property For all $A \in \mathbb{R}^{m \times n}$

$$AI = A = IA$$

For a diagonal matrix, $D = \text{diag}(d_1, d_2, \ldots, d_n)$

$$D_{ij} = \begin{cases} d_i, & i = j \\ 0, & i \neq j \end{cases}$$

So, $I = \text{diag}(1, 1, \ldots, 1)$

Vector - Vector product

$$① \quad x^T y \in \mathbb{R} = \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \sum x_i y_i$$

$$② \quad x y^T \in \mathbb{R}^{m \times n} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \begin{bmatrix} y_1 & y_2 & \cdots & y_n \end{bmatrix} = \begin{bmatrix} x_1 y_1 & x_1 y_2 & \cdots & x_1 y_n \\ x_2 y_1 & x_2 y_2 & \cdots & x_2 y_n \\ \vdots & \vdots & & \vdots \\ x_m y_1 & x_m y_2 & \cdots & x_m y_n \end{bmatrix}$$

① is called vector-vector inner product,
it is also known as dot product.

② is called vector-vector dot product.

Matrix - Vector Product

$$\text{①} \quad y = Ax = \begin{bmatrix} - & a_1^T & - \\ - & a_2^T & - \\ & \vdots & \\ - & a_m^T & - \end{bmatrix} x = \begin{bmatrix} a_1^T x \\ a_2^T x \\ \vdots \\ a_m^T x \end{bmatrix}$$

in ① we multiplied by rows

$$y = Ax = \begin{bmatrix} | & | & & | \\ a^1 & a^2 & \cdots & a^n \\ | & | & & | \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} | \\ a^1 \\ | \end{bmatrix} x_1 + \begin{bmatrix} | \\ a^2 \\ | \end{bmatrix} x_2 \cdots + \begin{bmatrix} | \\ a^n \\ | \end{bmatrix} x_n$$

in ② we multiplied by columns

y is a linear combination of columns of A

Matrix-Vector product (contd.)

we can multiply on the left by a row vector

$$y^T = x^T A = x^T \begin{bmatrix} 1 & 1 & & 1 \\ a^1 & a^2 & \cdots & a^n \\ 1 & 1 & & 1 \end{bmatrix} = \begin{bmatrix} x^T a^1 & x^T a^2 & \cdots & x^T a^n \end{bmatrix}$$

or

$$y^T = \begin{bmatrix} x_1 & x_2 & \cdots & x_m \end{bmatrix} \begin{bmatrix} - & a_1^T & - \\ - & a_2^T & - \\ & \vdots & \\ - & a_m & - \end{bmatrix} = x_1 \begin{bmatrix} - a_1^T - \end{bmatrix} + x_2 \begin{bmatrix} - a_2^T - \end{bmatrix} + \ldots + x_m \begin{bmatrix} - a_m^T - \end{bmatrix}$$

$y^T$ is a linear combination of

rows of $A$

Matrix - matrix multiplication

1. Using the concept of vector-vector inner (dot) product

$$C = AB = \begin{bmatrix} - a_1^T - \\ - a_2^T - \\ \vdots \\ - a_m^T - \end{bmatrix} = \begin{bmatrix} | & | & & | \\ b^1 & b^2 & \cdots & b^n \\ | & | & & | \end{bmatrix}$$

$$= \begin{bmatrix} a_1^T b^1 & a_1^T b^2 & \cdots & a_1^T b^n \\ a_2^T b^1 & a_2^T b^2 & \cdots & a_2^T b^n \\ \vdots & \vdots & \ddots & \vdots \\ a_m^T b^1 & a_m^T b^2 & \cdots & a_m^T b^n \end{bmatrix}$$

Matrix - matrix multiplication (contd.)

2. Using the concept of outer product

$$C = AB = \begin{bmatrix} | & | & & | \\ a^1 & a^2 & \cdots & a^n \\ | & | & & | \end{bmatrix} \begin{bmatrix} - & b_1^T & - \\ - & b_2^T & - \\ & \vdots & \\ - & b_n^T & - \end{bmatrix} = \sum_{i=1}^{n} a^i b_i^T$$

Matrix - matrix multiplication (contd.)

3. As a set of matrix - vector products

$$C = AB = A \begin{bmatrix} | & | & & | \\ b^1 & b^2 & \cdots & b^n \\ | & | & & | \end{bmatrix} = \begin{bmatrix} | & | & & | \\ Ab^1 & Ab^2 & \cdots & Ab^n \\ | & | & & | \end{bmatrix}$$

The $i$-th column of $C$ is given by the matrix-vector product with the vector on the right, $c_i = Ab_i$.

Matrix-matrix multiplication (contd.)

4. As a set of vector-matrix products

$$C = AB = \begin{bmatrix} -a_1^T- \\ -a_2^T- \\ \vdots \\ -a_m^T- \end{bmatrix} B = \begin{bmatrix} -a_1^T B- \\ -a_2^T B- \\ \vdots \\ -a_m^T B- \end{bmatrix}$$

Properties of matrix-matrix multiplication

① Associative : $(AB)C = A(BC)$

② Distributive : $A(B+C) = AB + AC$

③ Not necessarily commutative

For example, if $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{n \times a}$

and $m \neq n \neq a$, then $BA$ does not exist.

Hadamard or Elementwise Product

→ Matrix Product, $C = AB$

$$\Rightarrow C_{ij} = \sum_{k=1}^{n} A_{ik} B_{kj}$$

→ Hadamard Product, $C = A \odot B$

↳ $A, B, C$ are of the same size

↳ Multiply elements in $A$ and $B$ at same position. $(A \odot B)_{ij} = A_{ij} B_{ij}$

↳ Example

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \odot \begin{bmatrix} 1 & 5 & 6 \\ 1 & 7 & 8 \end{bmatrix} = \begin{bmatrix} 1 & 10 & 18 \\ 4 & 35 & 48 \end{bmatrix}$$

Transpose

The transpose operation flips the rows and columns. Given, $A \in \mathbb{R}^{m \times n}$, its transpose $A^T \in \mathbb{R}^{n \times m}$ whose entries are given by,

$$(A^T)_{ij} = A_{ji}$$

Properties

$$(A^T)^T = A$$
$$(AB)^T = B^T A^T$$
$$(A+B)^T = A^T + B^T$$

if $A = A^T$, then A is symmetric.

# Trace

Trace of a square matrix $A \in \mathbb{R}^{n \times n}$ is denoted by $\text{tr } A$ or $\text{tr}(A)$ is the sum of diagonal elements of $A$.

$$\text{tr } A = \sum_{i=1}^{n} A_{ii}$$

Properties : Consider $A \in \mathbb{R}^{n \times n}$ & $B \in \mathbb{R}^{n \times n}$

① $\text{tr } A = \text{tr } A^T$

② $\text{tr}(A+B) = \text{tr } A + \text{tr } B$

③ $\text{tr}(\alpha A) = \alpha \text{ tr } A, \alpha \in \mathbb{R}$

④ $\text{tr}(AB) = \text{tr}(BA)$, for $A, B$ such that $AB$ is square

⑤ $\text{tr}(ABC) = \text{tr}(BCA) = \text{tr}(CAB)$, such that $ABC$ is square. Can be expanded more

Rank of a matrix

Column rank of $A \in \mathbb{R}^{m \times n}$ is the largest number of columns of A that constitute a linearly independent set.

Row rank of $A \in \mathbb{R}^{m \times n}$ is the largest number of rows of A that constitute a linearly independent set.

For $A \in \mathbb{R}^{m \times n}$, column rank and row rank are equal, so generally rank is represented by $\text{rank}(A)$.

Properties of the rank

- For $A \in \mathbb{R}^{m \times n}$, $\text{rank}(A) \leq \min(m, n)$

- If $\text{rank}(A) = \min(m, n)$, then $A$ is <u>full rank</u>.

- For $A \in \mathbb{R}^{m \times n}$, $\text{rank}(A) = \text{rank}(A^T)$

- For $A \in \mathbb{R}^{m \times p}$, $B \in \mathbb{R}^{p \times n}$, $\text{rank}(AB) \leq \min(\text{rank}(A), \text{rank}(B))$

- For $A, B \in \mathbb{R}^{m \times n}$, $\text{rank}(A+B) \leq \text{rank}(A) + \text{rank}(B)$

## Inverse of A (A is a square matrix)

→ inverse of a square matrix $A \in \mathbb{R}^{n \times n}$ is denoted by $A^{-1}$, and is the unique matrix such that $A^{-1}A = I = AA^{-1}$

→ A is <u>invertible</u> or <u>non-singular</u> if $A^{-1}$ exists and <u>non-invertible</u> or <u>singular</u> otherwise.

→ For square matrix A to have an inverse $A^{-1}$, A must be <u>full rank</u>.

Properties of inverse of A (A is a square matrix)

$$\left(A^{-1}\right)^{-1} = A$$

$$(AB)^{-1} = B^{-1} A^{-1}$$

$$\left(A^{-1}\right)^{T} = \left(A^{T}\right)^{-1} \quad , \text{ sometimes denoted } A^{-T}$$

# Orthogonal Matrices

→ Two vectors are orthogonal, if $x^T y = 0$

→ A vector is normalized if $||x||_2 = 1$

→ A square matrix, $A \in R^{n \times n}$ is orthogonal if all its columns are orthogonal to each other and normalized.

→ Such columns are called <u>orthonormal</u>

# Properties of orthogonal matrices

→ The inverse of an orthogonal matrix is its transpose.

$$A^{-1} = A^T \quad \text{(A is orthogonal)}$$

$$A^T A = I = A A^T$$

→ Operating on a vector with an orthogonal matrix does not change its Euclidean norm,

$$\|U x\|_2 = \|x\|_2 \quad, \quad \begin{array}{l} U \text{ is orthogonal} \\ U \in \mathbb{R}^{n \times n} \\ x \in \mathbb{R}^n \end{array}$$

# Backpropagation (simplified)

Consider this network, with bias set to zero.



$a$ = activation

$\sigma$ = activation function

# *Backpropagation (simplified)*



$$a^{(1)} = x$$

$$z^{(2)} = \omega^{(1)} a^{(1)}$$

$$a^{(2)} = \sigma(z^{(2)})$$

$$z^{(3)} = \omega^{(2)} a^{(2)}$$

$$a^{(3)} = \sigma(z^{(3)})$$

$$z^{(4)} = \omega^{(3)} a^{(3)}$$

$$a^{(4)} = \sigma(z^{(4)})$$

$$\hat{y} = a^{(4)}$$

Forward
step

Consider a single data point

$$J = -\{y \ln \hat{y} + (1-y) \ln(1-\hat{y})\}$$

then,

$$\frac{\partial J}{\partial \omega^{(3)}} = \frac{\partial J}{\partial a^{(4)}} \cdot \frac{\partial a^{(4)}}{\partial z^{(4)}} \cdot \frac{\partial z^{(4)}}{\partial \omega^{(3)}} \quad \sigma = \text{sigmoid function}$$

$$= -\{y - a^{(4)}\} a^{(3)} \quad \text{[logistic regression] equivalent}$$

# *Backpropagation (simplified)*



$a^{(1)} = x$

$z^{(2)} = \omega^{(1)} a^{(1)}$

$a^{(2)} = \sigma(z^{(2)})$

$z^{(3)} = \omega^{(2)} a^{(2)}$

$a^{(3)} = \sigma(z^{(3)})$

$z^{(4)} = \omega^{(3)} a^{(3)}$
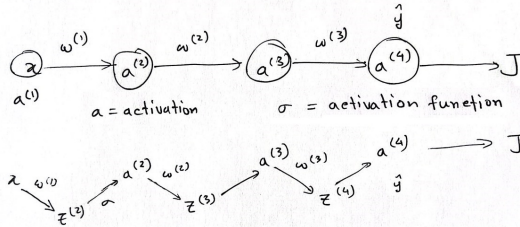
$a^{(4)} = \sigma(z^{(4)})$

$\hat{y} = a^{(4)}$

Forward
step

Let's define, $\dfrac{\partial J}{\partial z^{(4)}} = \delta^{(4)}$

generally, $\dfrac{\partial J}{\partial z^{(\ell)}} = \delta^{(\ell)}$

$$\therefore \quad \frac{\partial J}{\partial \omega^{(3)}} = -\{y - a^{(4)}\} a^{(3)}$$

# Backpropagation (simplified)



$a^{(1)} = x$

$z^{(2)} = \omega^{(1)} a^{(1)}$

$a^{(2)} = \sigma(z^{(2)})$

$z^{(3)} = \omega^{(2)} a^{(2)}$

$a^{(3)} = \sigma(z^{(3)})$

$z^{(4)} = \omega^{(3)} a^{(3)}$
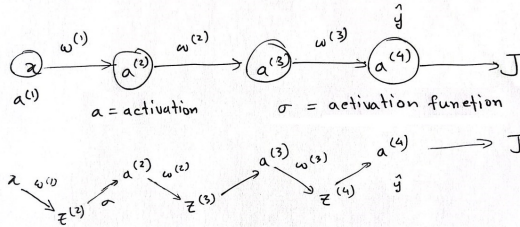
$a^{(4)} = \sigma(z^{(4)})$

$\hat{y} = a^{(4)}$

Forward
step

$$\frac{\partial J}{\partial \omega^{(2)}} = \frac{\partial J}{\partial a^{(4)}} \cdot \frac{\partial a^{(4)}}{\partial z^{(4)}} \cdot \frac{\partial z^{(4)}}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(3)}} \cdot \underbrace{\frac{\partial z^{(3)}}{\partial \omega^{(2)}}}$$

$$\frac{\partial J}{\partial z^{(3)}} = S^{(3)}$$

$$\therefore \frac{\partial J}{\partial \omega^{(2)}} = S^{(3)} a^{(2)}$$

Similarly, $\dfrac{\partial J}{\partial \omega^{(1)}} = S^{(2)} a^{(1)}$

Generally, $\boxed{\dfrac{\partial J}{\partial \omega^{(l)}} = S^{(l+1)} a^{(l)}}$

# Backpropagation (simplified)



$a^{(1)} = x$

$z^{(2)} = \omega^{(1)} a^{(1)}$

$a^{(2)} = \sigma(z^{(2)})$

$z^{(3)} = \omega^{(2)} a^{(2)}$

$a^{(3)} = \sigma(z^{(3)})$

$z^{(4)} = \omega^{(3)} a^{(3)}$

$a^{(4)} = \sigma(z^{(4)})$

$\hat{y} = a^{(4)}$

Forward step

Network with $L$ hidden layers, will have error,

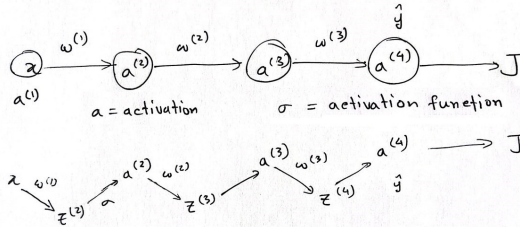$$\delta^{(L)} = -\{ y - a^{(L)} \} = \frac{\partial J}{\partial z^{(L)}}$$

For this network, $\delta^{(4)} = \frac{\partial J}{\partial z^{(4)}}$ is known

Now, $\delta^{(4)} = \frac{\partial J}{\partial a^{(4)}} \cdot \frac{\partial a^{(4)}}{\partial z^{(4)}}$

$$\delta^{(3)} = \frac{\partial J}{\partial z^{(3)}}$$

$$= \underbrace{\frac{\partial J}{\partial a^{(4)}} \cdot \frac{\partial a^{(4)}}{\partial z^{(4)}}}_{\delta^{(4)}} \cdot \underbrace{\frac{\partial z^{(4)}}{\partial a^{(3)}}}_{\omega^{(3)}} \cdot \underbrace{\frac{\partial a^{(3)}}{\partial z^{(3)}}}_{\sigma'(z^{(3)})}$$

# Backpropagation (simplified)



$$a^{(1)} = x$$
$$z^{(2)} = \omega^{(1)} a^{(1)}$$
$$a^{(2)} = \sigma(z^{(2)})$$
$$z^{(3)} = \omega^{(2)} a^{(2)}$$
$$a^{(3)} = \sigma(z^{(3)})$$
$$z^{(4)} = \omega^{(3)} a^{(3)}$$
$$a^{(4)} = \sigma(z^{(4)})$$
$$\hat{y} = a^{(4)}$$

Forward
step

$$\therefore \delta^{(3)} = \delta^{(4)} \omega^{(3)} \sigma'\left(z^{(3)}\right)$$

$$\boxed{S^{(l)} = \delta^{(l+1)} \omega^{(l)} \sigma'\left(z^{(l)}\right)}$$

Combine with $\quad \dfrac{\partial J}{\partial \omega^{(l)}} = \delta^{(l+1)} a^{(l)}$
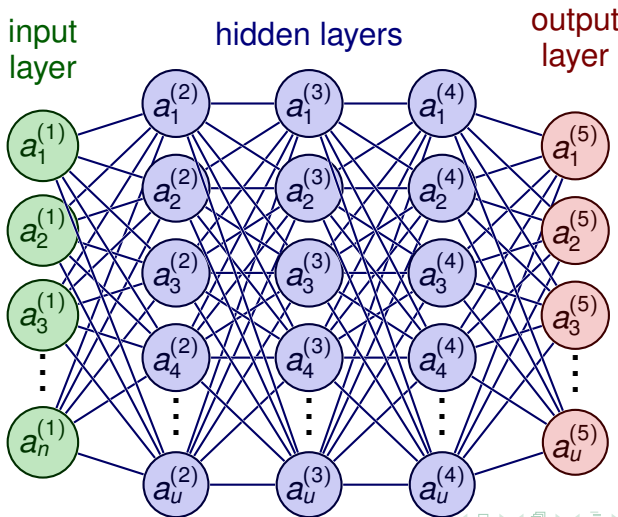
# A neural network

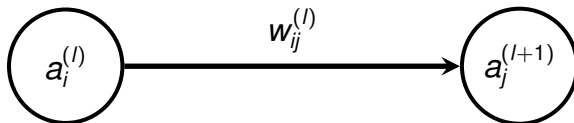$u$ represents number of units or neurons per layer

# A neural network: Only activations

*u* represents number of units or neurons per layer

# Forward Step: Calculate activations



$$a_j^{(l+1)} = \sigma(\sum_{i=0} w_{ij}^{(l)} a_i^{(l)})$$

Example : For the first neuron in the first hidden layer or second layer of the network, is

$$a_1^{(2)} = \sigma(\sum_{i=0} w_{i1}^{(1)} a_i^{(1)})$$

When i=0, w01 is taken into account. Similarly with every value of i you will progressively add the contribution of every neuron, to calculate the activation of the neuron of interest.

*Backpropagation: Update model parameters to reduce cost*

$$\delta^{(l)} = \delta^{(l+1)} w^{(l)} \sigma^{'}(z^{(l)})$$

$$\frac{\partial J}{\partial w^{(l)}} = \delta^{(l+1)} a^{(l)}$$

$\delta^{(L)}$ is known for the L-th layer.

*Backpropagation: Update model parameters to reduce cost*

$$\frac{\partial J}{\partial w_{ij}^{(l)}} = \delta_j^{(l+1)} a_i^{(l)}$$

$$\delta^{\vec{(l)}} = W^{(l)} \delta^{\vec{(l)}} \odot \sigma'(z^{\vec{(l)}})$$

$\odot$ is elementwise product or Hadamard product