# RECURRENT NEURAL NETWORKS

{Grady Kurpasi}

{SSIE 616}

{Prof H. Lewis}

{https://github.com/GradyKurpasi/RNN-Tutorial}

## $x, y$

| x₁ | x₂ | x₃ | y₁ | y₂ |
|--------|--------|--------|--------|--------|
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |
| 0.1000 | -5.0000 | 3.0000 | 0.1221 | 0.0964 |
| 6.0000 | -5.5420 | 4.8970 | 0.1061 | 0.0702 |
| 4.0000 | 8.0000 | 9.0000 | 0.0996 | 0.0641 |
| 12.0000 | -2.0000 | 0.0063 | 0.1110 | 0.0732 |
| 6.0000 | -5.5000 | 4.8970 | 0.1060 | 0.0701 |

$x, y$ is titled as $x, y$ with columns $x_1, x_2, x_3, y_1, y_2$.

## $w_1[i, k]$

| i | w[i,1] | w[i,2] |
|---|--------|--------|
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

## $w_2[k, j]$

| k | w[k,1] | w[k,2] |
|---|--------|--------|
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

$a_2[j]$
$j = 1..n$

$a_1[k]$
$k = 1..l$

$a_0[i]$
$i = 1..m$

$n = 2$
\# output

$l = 2$
\# hidden

$m = 2$
\# input

$a_2[1]$  $a_2[2]$

$a_1[0]$  $a_1[1]$  $a_1[2]$

$a_0[0]$  $a_0[1]$  $a_0[2]$  $a_0[3]$

$x, y$

| x₁ | x₂ | x₃ | y₁ | y₂ |
|---|---|---|---|---|
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |
| 0.1000 | -5.0000 | 3.0000 | 0.1221 | 0.0964 |
| 6.0000 | -5.5420 | 4.8970 | 0.1061 | 0.0702 |
| 4.0000 | 8.0000 | 9.0000 | 0.0996 | 0.0641 |
| 12.0000 | -2.0000 | 0.0063 | 0.1110 | 0.0732 |
| 6.0000 | -5.5000 | 4.8970 | 0.1060 | 0.0701 |

$w_1[i,k]$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

$w_2[k,j]$

| k | w[k,1] | w[k,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

$a_2[j]$
j = 1..n

$a_1[k]$
k = 1..l

$a_0[i]$
i = 1..m

$subscript_2$

$subscript_1$

$subscript_0$

## $x, y$

| $x_1$ | $x_2$ | $x_3$ | $y_1$ | $y_2$ |
|---|---|---|---|---|
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |
| 0.1000 | -5.0000 | 3.0000 | 0.1221 | 0.0964 |
| 6.0000 | -5.5420 | 4.8970 | 0.1061 | 0.0702 |
| 4.0000 | 8.0000 | 9.0000 | 0.0996 | 0.0641 |
| 12.0000 | -2.0000 | 0.0063 | 0.1110 | 0.0732 |
| 6.0000 | -5.5000 | 4.8970 | 0.1060 | 0.0701 |

## $w_1[i, k]$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

## $w_2[k, j]$

| k | w[k,1] | w[k,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

$x, y$

| x₁ | x₂ | x₃ | y₁ | y₂ |
|---|---|---|---|---|
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |

$w_1[i, k]$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | **0.5** | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

$w_2[k, j]$

| k | w[k,1] | w[k,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

$a_2[j]$

$a_1[k]$

$a_0[i]$

$$s_1[k] = \sum_{i=0}^{m} w_1[i,k] * a_0 i \quad \text{for k} = 1..l$$

$$s_1[1] = w_1[0,1] * a_0[0] = .05 * 1$$

$$x, y$$

| x₁ | x₂ | x₃ | y₁ | y₂ |
|---|---|---|---|---|
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |

$$w_1[i,k]$$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | **0.1** | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

$$w_2[k,j]$$

| k | w[k,1] | w[k,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

$$s_1[k] = \sum_{i=0}^{3} w_1[i,k] * a_0 i \quad \text{for k} = 1, 2$$

$$s_1[1] = \quad w_1[0,1] * a_0[0] = \quad .05 * 1$$

$$w_1[1,1] * a_0[1] = \quad .1 * 1$$

## $x, y$

| x₁ | x₂ | x₃ | y₁ | y₂ |
|---|---|---|---|---|
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |

## $w_1[i, k]$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | **0.5** | 0.5 |
| 1 | **0.1** | 0.2 |
| 2 | **0.3** | 0.4 |
| 3 | **0.5** | 0.6 |

## $w_2[k, j]$

| k | w[k,1] | w[k,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |



$$s_1[k] = \sum_{i=0}^{3} w_1[i,k] * a_0 i \quad \text{for k} = 1, 2$$

$$s_1[1] = \quad w_1[0,1] * a_0[0] = \quad .05 \ * \ 1 \ +$$

$$w_1[1,1] * a_0[1] = \quad .1 \ * \ 1 \ +$$

$$w_1[2,1] * a_0[2] = \quad .3 \ * \ 4 \ +$$

$$w_1[3,1] * a_0[3] = \quad .5 \ * \ 5$$

$$= \quad 4.3$$

| $x, y$ | | | | |
|---|---|---|---|---|
| **x₁** | **x₂** | **x₃** | **y₁** | **y₂** |
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |

| $w_1[i, k]$ | | |
|---|---|---|
| **i** | **w[i,1]** | **w[i,2]** |
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

| $w_2[k, j]$ | | |
|---|---|---|
| **k** | **w[k,1]** | **w[k,2]** |
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

$a_2[j]$

$a_1[k]$

$a_0[i]$



$$a_1[k] = \frac{1}{1 + e^{-s_1[k]}} \quad \text{for k} = 1, 2$$

$$a_1[1] = \frac{1}{1 + e^{-4.3}} = .9866$$

**x, y**

| x₁ | x₂ | x₃ | y₁ | y₂ |
|---|---|---|---|---|
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |

$w_1[i,k]$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

$w_2[k,j]$

| k | w[k,1] | w[k,2] |
|---|---|---|
| 0 | **0.5** | 0.5 |
| 1 | **0.7** | 0.8 |
| 2 | **0.9** | 0.1 |

$$s_2[j] = \sum_{k=0}^{l} w_2[k,j] * a_1 i \quad \text{for j} = 1..\,\text{m}$$

$$a_2[j] = \frac{1}{1 + e^{-s_2[j]}} \quad \text{for j} = 1, 2$$
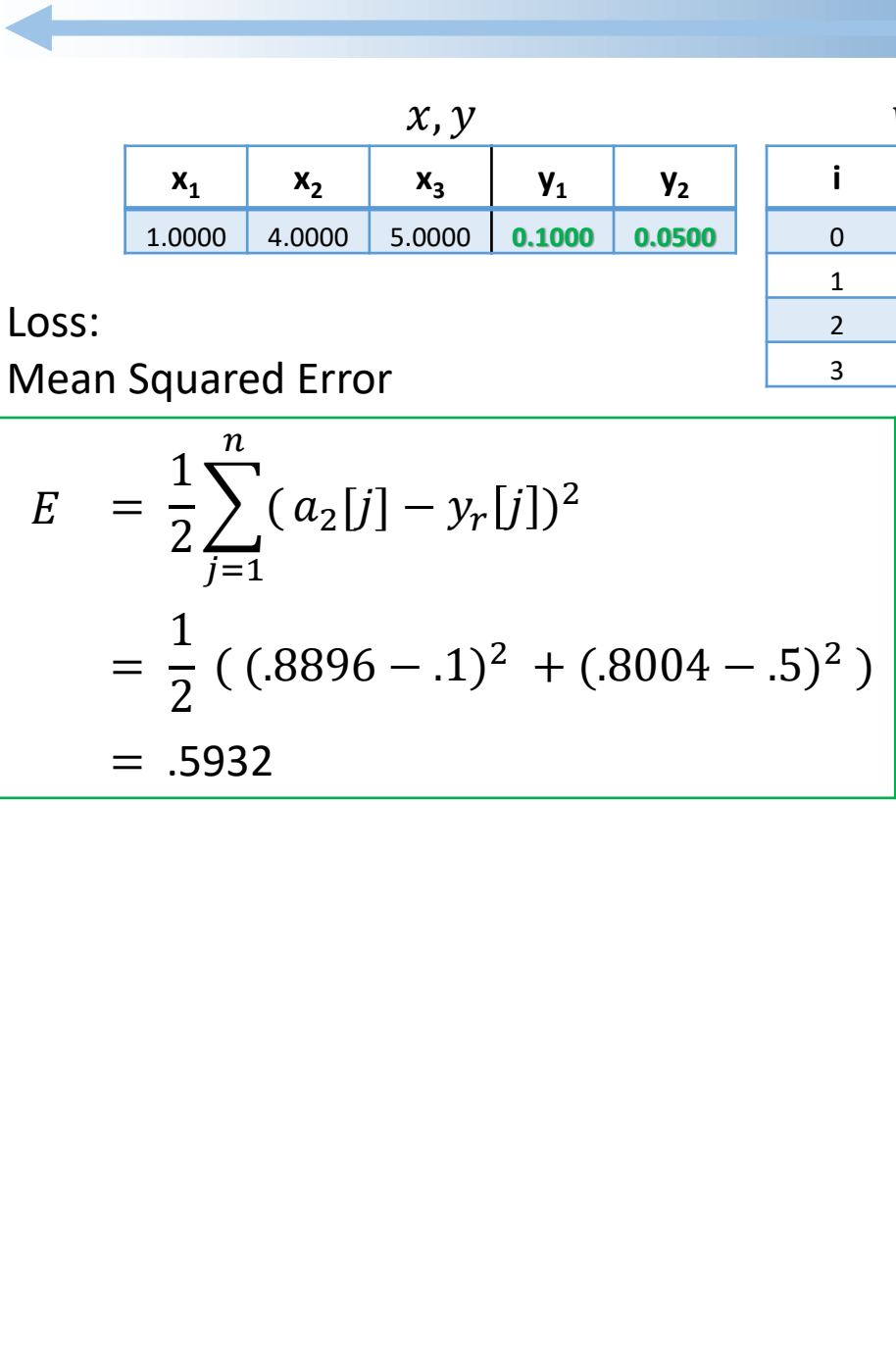
$s_2[1] =$

$w_2[0,2] * a_1[0] = \quad .5 * 1$

$w_2[1,2] * a_1[1] = \quad .7 * .9866$

$w_2[2,2] * a_1[2] = \quad .9 * .9950$

$= \quad 2.0862$

$$a_2[1] = \frac{1}{1 + e^{-2.0862}} = .8896 = \hat{y}$$

Desired -  y1: .1    y2: .05

Calculated - ŷ1: .8896    ŷ2: .8004

$a_2[j]$

$a_1[k]$

$a_0[i]$



## $x, y$

| $x_1$ | $x_2$ | $x_3$ | $y_1$ | $y_2$ |
|--------|--------|--------|--------|--------|
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |

## $w_1[i, k]$

| i | w[i,1] | w[i,2] |
|---|--------|--------|
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

## $w_2[k, j]$

| k | w[k,1] | w[k,2] |
|---|--------|--------|
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

Loss:

Mean Squared Error

$$E = \frac{1}{2}\sum_{j=1}^{n} (a_2[j] - y_r[j])^2$$

$$= \frac{1}{2}\left( (.8896 - .1)^2 + (.8004 - .5)^2 \right)$$

$$= .5932$$

$E = .5932$

y1
.1

y2
.05

Desired -  .1          .05

Calculated - .8896      .8004

$\hat{y}1$          $\hat{y}2$

$a_2[j]$

$a_2[1]$          $a_2[2]$

$w_2[0,1]$

$a_1[k]$

1

$a_1[0]$     $a_1[1]$     $a_1[2]$

$a_0[i]$

1      1      4      5

$a_0[0]$     $a_0[1]$     $a_0[2]$     $a_0[3]$

| | $x, y$ | | | |
|---|---|---|---|---|
| **x₁** | **x₂** | **x₃** | **y₁** | **y₂** |
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |

| | $w_1[i,k]$ | |
|---|---|---|
| **i** | **w[i,1]** | **w[i,2]** |
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

| | $w_2[k,j]$ | |
|---|---|---|
| **k** | **w[k,1]** | **w[k,2]** |
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

Partial Derivative of Error with respect to Weight$_2$[0,1]

Functional Dependencies

$$E = \frac{1}{2}(\hat{y} - y1)^2 \longrightarrow \text{E}(\hat{y})$$

$$\hat{y} = \sigma(s_2) \longrightarrow \text{f}(s_2)$$

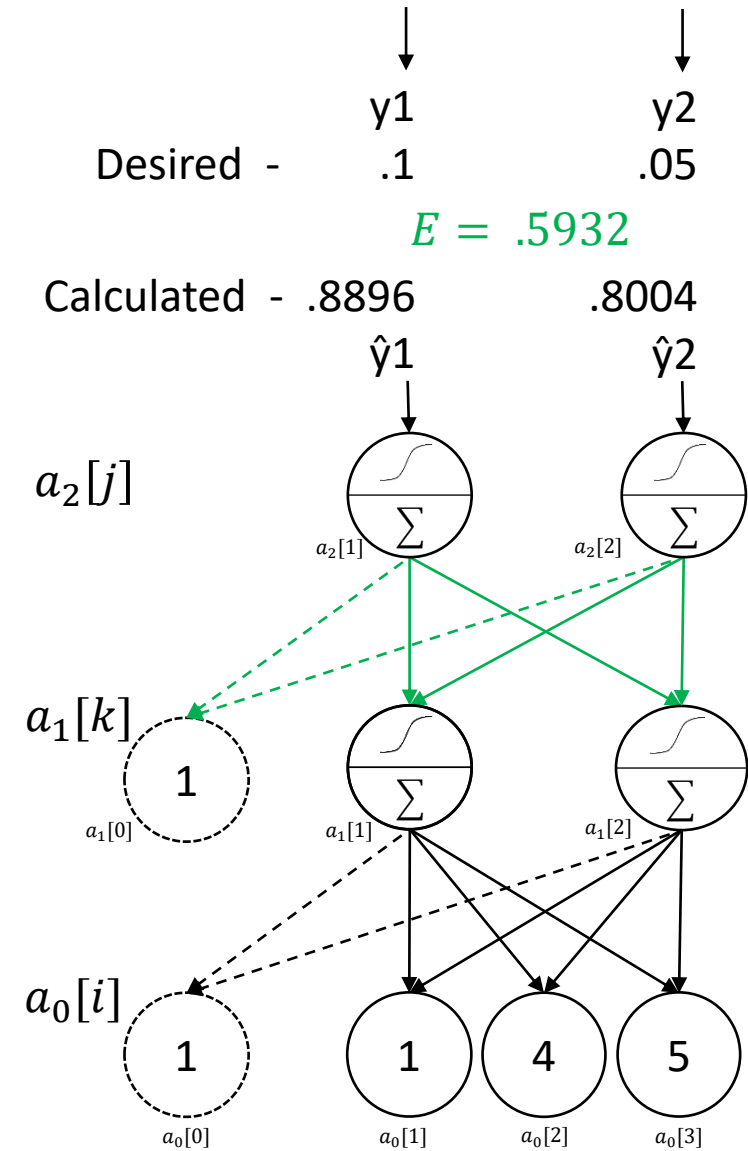$$s_2 = \sum a_1[k] * w_2[0,1] + .6906 + .8955 \longrightarrow \text{g}(w_2)$$

Chain Rule
$$E(f(g(w_2)))$$

y1      y2

Desired -    .1      .05

$E = .5932$

Calculated - .8896      .8004

$\hat{y}1$      $\hat{y}2$

$a_2[j]$

$a_1[k]$

$a_0[i]$



| $x, y$ | | | | |
|---|---|---|---|---|
| **x₁** | **x₂** | **x₃** | **Y₁** | **Y₂** |
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |

| $w_1[i, k]$ | | |
|---|---|---|
| **i** | **w[i,1]** | **w[i,2]** |
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

| $w_2[k, j]$ | | |
|---|---|---|
| **k** | **w[k,1]** | **w[k,2]** |
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

$$\frac{\partial E}{\partial w_2[k,j]} = \frac{\partial E}{\partial a_2[j]} * \frac{\partial a_2[j]}{\partial s_2[j]} * \frac{\partial s_2[j]}{\partial w_2[k,j]}$$

$$= \{\hat{y}_j - y_k\} * \{a_2[j](1 - a_2[j])\} * \{a_1[k]\}$$

y1
.1

Desired -  .1

y2
.05

Calculated - .8896

.8004

$\hat{y}1$

$\hat{y}2$

$a_2[j]$

$a_2[1]$

$a_2[2]$

$w_2[0,1]$

$a_1[k]$

1

$a_1[0]$

$a_1[1]$

$a_1[2]$

$a_0[i]$

1

$a_0[0]$

1

$a_0[1]$

4

$a_0[2]$

5

$a_0[3]$

## $x, y$

| x₁ | x₂ | x₃ | y₁ | y₂ |
|---|---|---|---|---|
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |

## $w_1[i, k]$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

## $w_2[k, j]$

| k | w[k,1] | w[k,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

$$\frac{\partial E}{\partial w_2[0,1]} = \frac{\partial E}{\partial a_2[1]} * \frac{\partial a_2[1]}{\partial s_2[1]} * \frac{\partial s_2[1]}{\partial w_2[0,1]}$$

$$= \{\hat{y}_j - y_k\} * \{a_2[1](1 - a_2[1])\} * \{a_1[0]\}$$

$$= .8896 - .1$$

# Backpropagation

Desired - y1 .1  y2 .05

Calculated - .8896   .8004

ŷ1   ŷ2

$a_2[j]$

$a_2[1]$  $a_2[2]$

$w_2[0,1]$

$a_1[k]$

1

$a_1[0]$   $a_1[1]$   $a_1[2]$

$a_0[i]$

1   1   4   5

$a_0[0]$   $a_0[1]$   $a_0[2]$   $a_0[3]$

## $x, y$

| x₁ | x₂ | x₃ | y₁ | y₂ |
|---|---|---|---|---|
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |

## $w_1[i, k]$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

## $w_2[k, j]$

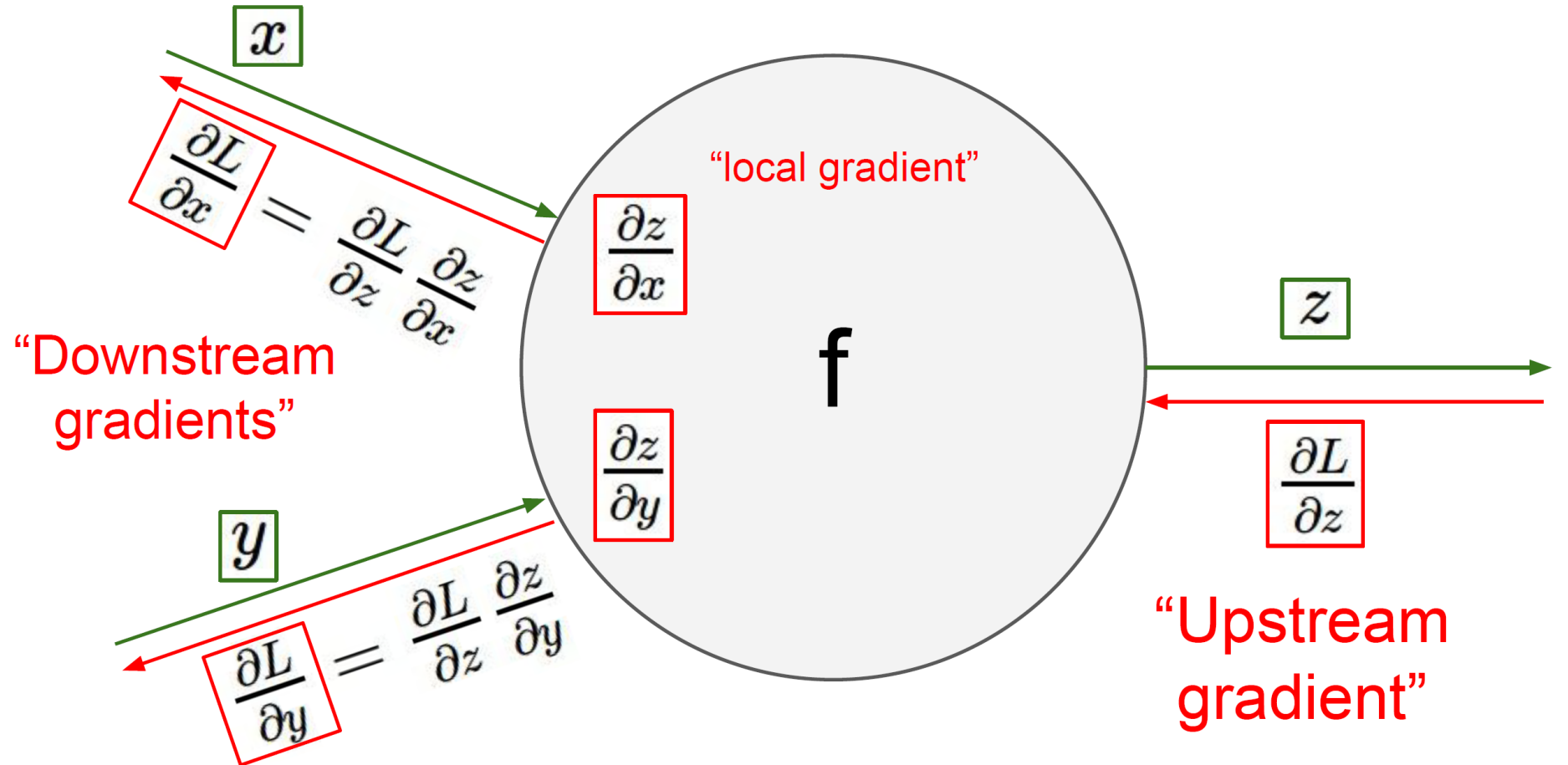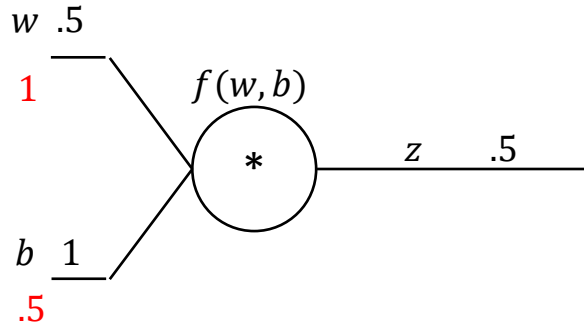| k | w[k,1] | w[k,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

$$\frac{\partial E}{\partial w_2[0,1]} = \frac{\partial E}{\partial a_2[1]} * \frac{\partial a_2[1]}{\partial s_2[1]} * \frac{\partial s_2[1]}{\partial w_2[0,1]}$$

$$= \{\hat{y}_j - y_k\} * \{a_2[1](1 - a_2[1])\} * \{a_1[0]\}$$

$$= .8896 - .1 * .8896(1 - .8896)$$

Desired - y1 .1    y2 .05

Calculated - .8896    .8004

$a_2[j]$

$\hat{y}1$    $\hat{y}2$

$a_2[1]$    $a_2[2]$

$w_2[0,1]$

$a_1[k]$

1

$a_1[0]$    $a_1[1]$    $a_1[2]$

$a_0[i]$

1    1    4    5

$a_0[0]$    $a_0[1]$    $a_0[2]$    $a_0[3]$

### $x, y$

| $x_1$ | $x_2$ | $x_3$ | $y_1$ | $y_2$ |
|---|---|---|---|---|
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |

### $w_1[i, k]$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

### $w_2[k, j]$

| k | w[k,1] | w[k,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

$$\frac{\partial E}{\partial w_2[0,1]} = \frac{\partial E}{\partial a_2[1]} * \frac{\partial a_2[1]}{\partial s_2[1]} * \frac{\partial s_2[1]}{\partial w_2[0,1]}$$

$$= \{\hat{y}_j - y_k\} * \{a_2[1](1 - a_2[1])\} * \{a_1[0]\}$$
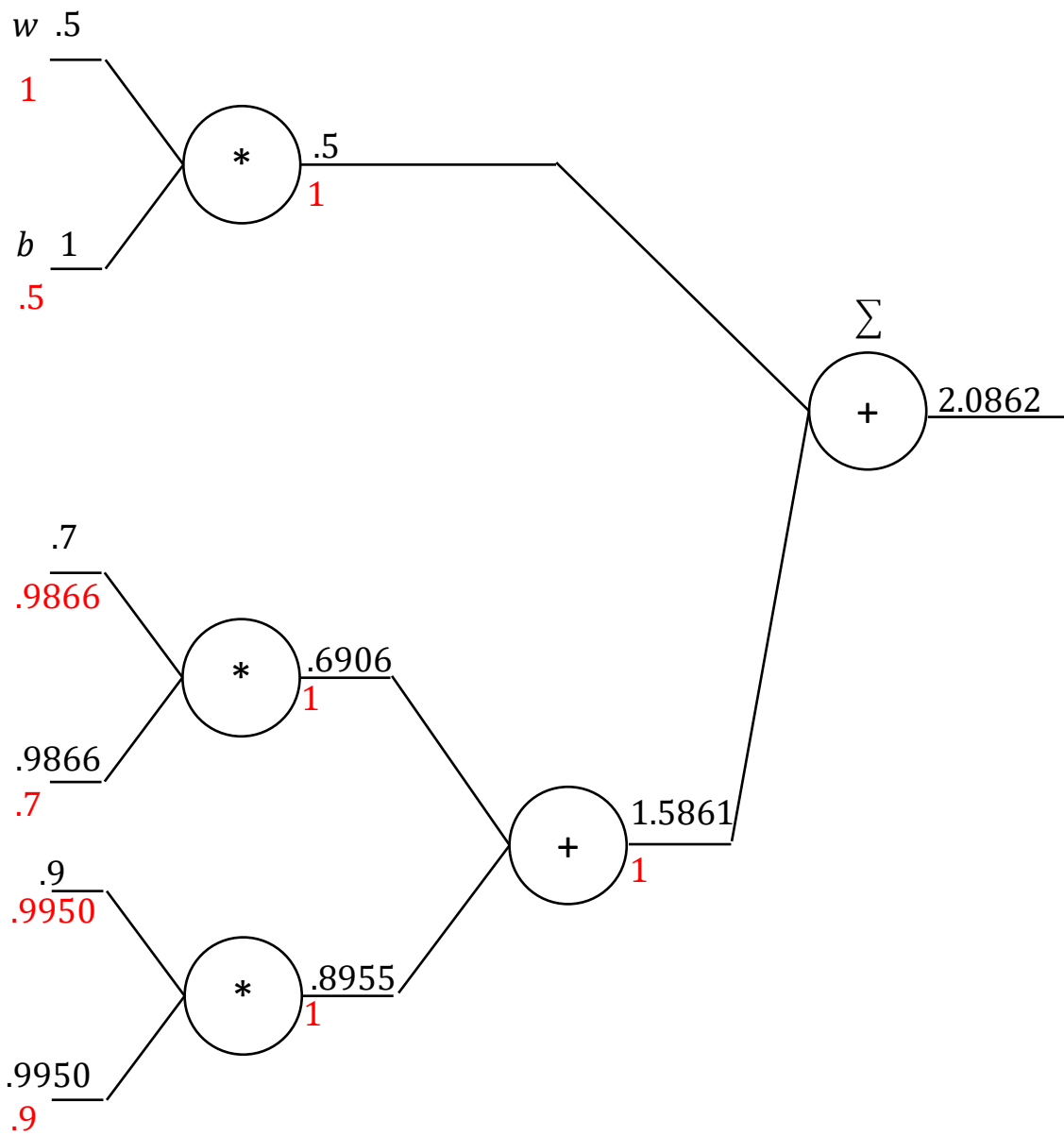
$$= .8896 - .1 * .8896(1 - .8896) * 1$$

$$= .0775$$

### $c_2[k, j]$

| k | c[k,1] | c[k,2] |
|---|---|---|
| 0 | .0775 | 0 |
| 1 | 0 | 0 |
| 2 | 0 | 0 |

$w$ .5

1

$f(w, b)$

*

$z$ .5

$b$ 1

.5

$w$ .5
1
$f(w, b)$
*
$z$ .5
$b$ 1
.5



$x$

$$\boxed{\frac{\partial L}{\partial x}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x}$$

"local gradient"

$$\frac{\partial z}{\partial x}$$

**f**

$z$

"Downstream gradients"

$$\frac{\partial z}{\partial y}$$

$y$

$$\boxed{\frac{\partial L}{\partial y}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial y}$$

$$\frac{\partial L}{\partial z}$$

"Upstream gradient"

cs231n_2019_lecture04.pdf (stanford.edu)

$w$ .5

1

$b$ 1

.5

$*$ .5

1

$\Sigma$

$+$ 2.0862

.7

.9866

.9866

.7

$*$ .6906

1

.9

.9950

.9950

.9

$*$ .8955

1

$+$ 1.5861

1

Sigmoid Function

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

Sigmoid Function

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

$w$  .5
1

*  .5
1

$b$  1
.5

$\Sigma$

+  2.0862  *  −2.0862  exp  .1242  +  1.1242  1/x  .8896
                                        −.7913 $x$ 1        1.0000

−1                                1

.7
.9866

*  .6906
1

.9866
.7

+  1.5861
1

.9
.9950

*  .8955
1

.9950
.9

$$\frac{1}{1.1242^2} = -.7913$$

$$f(x) = ax, \qquad \frac{df}{dx} = a \qquad\qquad f(x) = c + x, \qquad \frac{df}{dx} = 1$$

$$f(x) = e^x, \qquad \frac{df}{dx} = e^x \qquad\qquad f(x) = \frac{1}{x}, \qquad \frac{df}{dx} = -\frac{1}{x^2}$$

Sigmoid Function

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

$w$ .5

1

.5

1

$b$ 1

.5

$\Sigma$

+   2.0862   *   $-2.0862$   exp   .1242   +   1.1242   1/x   .8896

$-1\,x - .0982$
$=.0982$

$.1241\,x - .7913$
$= -.0982$

$1\,x - .7913$

$-.7913\,x\ 1$

1.0000

$-1$

1

$1\,x - .7913$

$e^{-2.0862} = .1241$

.7

.9866

*

.6906

1

.9866

.7

+

1.5861

1

.9

.9950

*

.8955

1

.9950

.9

$$f(x) = ax, \qquad \frac{df}{dx} = a \qquad\qquad f(x) = c + x, \qquad \frac{df}{dx} = 1$$

$$f(x) = e^{x}, \qquad \frac{df}{dx} = e^{x} \qquad\qquad f(x) = \frac{1}{x}, \qquad \frac{df}{dx} = -\frac{1}{x^{2}}$$

$w$ .5
1

.5
1

$b$ 1
.5

**Sigmoid Function**

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

$\sum$

+    2.0862    *    $-2.0862$    exp    .1242    +    1.1242    1/x    .8896

$-1\,x - .0982$
$=.0982$

$-1$

.1241 x − .7913
= -.0982

1 x − .7913

−.7913 x 1

1.0000

1

1 x − .7913

.7
.9866

.9866
.7

.9
.9950

.6906
1

*

1.5861
1

+

$(1 - .8896).8896 = .0982$

.9950
.9

.8955
1

*

**Sigmoid Derivative**
$$= \big(1 - \sigma(x)\big) * \sigma(x)$$

$w$ .5

$1 * .0775$
$= .0775$

$*$ .5
$1*.0775$
$=.0775$

$b$ 1
.5

.7
.9866

$*$ .6906
1

.9866
.7

.9
.9950

$*$ .8955
1

.9950
.9

$+$ 1.5861
1

$\Sigma$

$+$ 2.0862
$.0982x.7896$
$=.0775$

$-1$

$\int$ $\hat{y}$ .8896
$.8896 - .1$
$= .7986 \, x \, 1$
$=.7986$

$y$ .1

E .3117
1

Desired -

| | y1 | y2 |
|---|---|---|
| | .1 | .05 |

Calculated - .8896 .8004

$\hat{y}1$ $\hat{y}2$

$a_2[j]$

$w_2[0,1]$

$a_1[k]$

$a_0[i]$

### $x, y$

| $x_1$ | $x_2$ | $x_3$ | $y_1$ | $y_2$ |
|---|---|---|---|---|
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |

### $w_1[i,k]$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

### $w_2[k,j]$

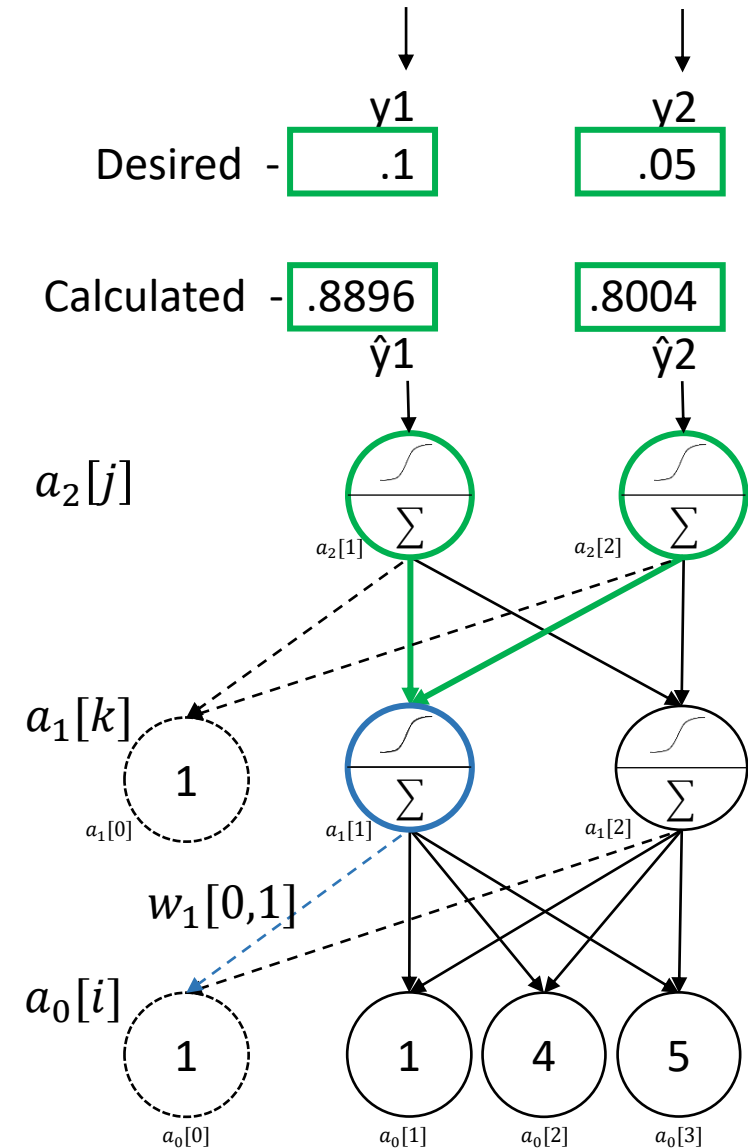| k | w[k,1] | w[k,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

$$\frac{\partial E}{\partial w_2[0,1]} = \frac{\partial E}{\partial a_2[1]} * \frac{\partial a_2[1]}{\partial s_2[1]} * \frac{\partial s_2[1]}{\partial w_2[0,1]}$$

$$= \{\hat{y}_j - y_k\} * \{a_2[1](1 - a_2[1])\} * \{a_1[0]\}$$

$$= .8896 - .1 * .8896(1 - .8896) * 1$$

$$= .0775$$

### $c_2[k,j]$

| k | c[k,1] | c[k,2] |
|---|---|---|
| 0 | .0775 | 0 |
| 1 | 0 | 0 |
| 2 | 0 | 0 |

$w$ .5

$1 * .0775$
$= .0775$

$*$ .5
$1*.0775$
$=.0775$

$b$ 1
.5

$\Sigma$

$+$ 2.0862
$.0982 x .7896$
$=.0775$

$\int$ $\hat{y}$ .8896
$.8896 - .1$
$= .7986 \, x \, 1$
$=.7986$

E .3117
1

$-1$

$y$ .1

.7
.9866

$*$ .6906
1

.9866
.7

$+$ 1.5861
1

.9
.9950

$*$ .8955
1

.9950
.9

f 1

$1 + 2 = 3$

2

Desired - 
y1 .1
y2 .05

$E = .5932$

Calculated - .8896    .8004

$\hat{y}1$    $\hat{y}2$

$a_2[j]$

$a_{2}[1]$    $a_{2}[2]$

$a_1[k]$

1
$a_{1}[0]$

$a_{1}[1]$    $a_{1}[2]$

$a_0[i]$

1
$a_{0}[0]$

1
$a_{0}[1]$

4
$a_{0}[2]$

5
$a_{0}[3]$

### $x, y$

| $x_1$ | $x_2$ | $x_3$ | $y_1$ | $y_2$ |
|---|---|---|---|---|
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |

### $w_1[i, k]$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

### $w_2[k, j]$

| k | w[k,1] | w[k,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

$$\frac{\partial E}{\partial w_1[i,k]} = \frac{\partial E}{\partial a_1[k]} * \frac{\partial a_1[k]}{\partial s_1[k]} * \frac{\partial s_1[k]}{\partial w_1[i,k]}$$

$$= \frac{\partial E}{\partial a_1[k]} * \{a_1[j](1 - a_1[j])\} * \{a_0[i]\}$$

y1
Desired - .1

y2
.05

Calculated - .8896

.8004

$\hat{y}1$

$\hat{y}2$

$a_2[j]$

$a_2[1]$   $a_2[2]$

$a_1[k]$

1

$a_1[0]$   $a_1[1]$   $a_1[2]$

$w_1[0,1]$

$a_0[i]$

1     1     4     5

$a_0[0]$   $a_0[1]$   $a_0[2]$   $a_0[3]$

### $x, y$

| x₁ | x₂ | x₃ | y₁ | y₂ |
|---|---|---|---|---|
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |

### $w_1[i, k]$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

### $w_2[k, j]$

| k | w[k,1] | w[k,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

$$\frac{\partial E}{\partial w_1[i, k]} = \frac{\partial E}{\partial a_1[k]} * \frac{\partial a_1[k]}{\partial s_1[k]} * \frac{\partial s_1[k]}{\partial w_1[i, k]}$$

$$= \frac{\partial E}{\partial a_1[k]} * \{a_1[k](1 - a_1[k])\} * \{a_0[i]\}$$

$$\frac{\partial E}{\partial a_1[k]} = \sum_{j=1}^{n} \frac{\partial E}{\partial a_2[j]} * \frac{\partial a_2[j]}{\partial s_2[j]} * \frac{\partial s_2[j]}{\partial a_1[k, j]}$$

$$= \sum_{j=1}^{n} \{\hat{y}_j - y_k\} * \{a_2[j](1 - a_2[j])\} * \{w_2[k, j]\}$$

Desired - y1 `.1`   y2 `.05`

Calculated - `.8896` $\hat{y}1$   `.8004` $\hat{y}2$

$a_2[j]$

$a_1[k]$

$w_1[0,1]$

$a_0[i]$



| | $x, y$ | | | |
|---|---|---|---|---|
| $x_1$ | $x_2$ | $x_3$ | $y_1$ | $y_2$ |
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |

| | $w_1[i,k]$ | |
|---|---|---|
| i | w[i,1] | w[i,2] |
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

| | $w_2[k,j]$ | |
|---|---|---|
| k | w[k,1] | w[k,2] |
| 0 | 0.5 | 0.5 |
| 1 | **0.7** | **0.8** |
| 2 | 0.9 | 0.1 |

$$\frac{\partial E}{\partial w_1[0,1]} = \frac{\partial E}{\partial a_1[1]} * \frac{\partial a_1[1]}{\partial s_1[1]} * \frac{\partial s_1[1]}{\partial w_1[0,1]}$$

$$= \frac{\partial E}{\partial a_1[1]} * \{a_1[1](1-a_1[1])\} * \{a_0[0]\}$$

$$\frac{\partial E}{\partial a_1[k]} = \sum_{j=1}^{n} \frac{\partial E}{\partial a_2[j]} * \frac{\partial a_2[j]}{\partial s_2[j]} * \frac{\partial s_2[j]}{\partial a_1[k,j]}$$

$$= \sum_{j=1}^{n} \{\hat{y}_j - y_k\} * \{a_2[j](1-a_2[j])\} * \{w_2[k,j]\}$$

$$= .8896 - .1 \quad * \quad .8896(1 - .8896) \quad * \quad .7$$

$$+ .8004 - .05 \quad * \quad .8004(1 - .8004) \quad * \quad .8$$

$$= .1502$$

y1
.1

y2
.05

Desired -   .1        .05

Calculated -  .8896        .8004

$\hat{y}1$        $\hat{y}2$

$a_2[j]$



$a_1[k]$

$a_0[i]$

$w_1[0,1]$

### $x, y$

| $x_1$ | $x_2$ | $x_3$ | $y_1$ | $y_2$ |
|---|---|---|---|---|
| 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |

### $w_1[i,k]$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

### $w_2[k,j]$

| k | w[k,1] | w[k,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

$$\frac{\partial E}{\partial w_1[0,1]} = \frac{\partial E}{\partial a_1[1]} \quad * \quad \frac{\partial a_1[1]}{\partial s_1[1]} \quad * \quad \frac{\partial s_1[1]}{\partial w_1[0,1]}$$

$$= \frac{\partial E}{\partial a_1[1]} \quad * \quad \{a_1[1](1-a_1[1])\} \quad * \quad \{a_0[0]\}$$

$$= .1502 \quad * \quad .9866(1-.9866) \quad * \quad 1$$

$$= .0020$$

### $c_1[i,k]$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | .0020 | 0 |
| 1 | 0 | 0 |
| 2 | 0 | 0 |
| 3 | 0 | 0 |

Desired - y1    y2

Calculated - ŷ1    ŷ2

$a_2[j]$

$a_1[k]$

$a_0[i]$

### $x, y$

|  | $x_1$ | $x_2$ | $x_3$ | $y_1$ | $y_2$ |
|---|---|---|---|---|---|
| r =1 | 1.0000 | 4.0000 | 5.0000 | 0.1000 | 0.0500 |
| ⋮ | 0.1000 | -5.0000 | 3.0000 | 0.1221 | 0.0964 |
|  | 6.0000 | -5.5420 | 4.8970 | 0.1061 | 0.0702 |
| r =4 | 4.0000 | 8.0000 | 9.0000 | 0.0996 | 0.0641 |
|  | 12.0000 | -2.0000 | 0.0063 | 0.1110 | 0.0732 |
|  | 6.0000 | -5.5000 | 4.8970 | 0.1060 | 0.0701 |

### $w_1[i,k]$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.1 | 0.2 |
| 2 | 0.3 | 0.4 |
| 3 | 0.5 | 0.6 |

### $w_2[k,j]$

| k | w[k,1] | w[k,2] |
|---|---|---|
| 0 | 0.5 | 0.5 |
| 1 | 0.7 | 0.8 |
| 2 | 0.9 | 0.1 |

### $c_1[i,k]$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | .0604 | .0332 |
| 1 | .1134 | .0427 |
| 2 | -.2929 | -.1658 |
| 3 | .2193 | .1129 |

### $c_2[k,j]$

| k | c[k,1] | c[k,2] |
|---|---|---|
| 0 | .3447 | .4816 |
| 1 | .2929 | .4176 |
| 2 | .2930 | .4191 |

### $w_1{}^*[i,k]$

| i | w[i,1] | w[i,2] |
|---|---|---|
| 0 | .4994 | .4997 |
| 1 | .0989 | .1996 |
| 2 | .3029 | .4017 |
| 3 | .4978 | .5989 |

### $w_2{}^*[k,j]$

| k | w[k,1] | w[k,2] |
|---|---|---|
| 0 | .4966 | .0039 |
| 1 | .4971 | .0038 |
| 2 | .0021 | .0004 |

Optimizer:

Stochastic Gradient Descent

Learning rate α = .01

$$w^* = w - \alpha * c$$

$w$ .5

$1 * .0775$
$= .0775$

$*$   .5
$1*.0775$
$=.0775$

$b$  1

.5

$\Sigma$

$+$

2.0862
$.0982x.7896$
$=.0775$

$\hat{y}$   .8896

$\int$

$.8896 - .1$
$= .7986 \ x \ 1$
$=.7986$

E   .3117

1

$-1$

$y$  .1

.7

.9866

$*$   .6906
1

.9866

.7

$+$   1.5861
1

.9

.9950

$*$   .8955
1

.9950

.9

## Rectified Linear Unit

| Sigmoid | Tanh | RELU |
|---|---|---|
| $g(z) = \dfrac{1}{1 + e^{-z}}$ | $g(z) = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$ | $g(z) = \max(0, z)$ |

CS 230 - Recurrent Neural Networks Cheatsheet (stanford.edu)

"hello"

Vocabulary

| "h" | [1, 0, 0, 0] |
| "e" | [0, 1, 0, 0] |
| "l" | [0, 0, 1, 0] |
| "o" | [0, 0, 0, 1] |

|  | y1 | y2 | y3 | y4 |
|---|---|---|---|---|
| "e"→ | 0 | 1 | 0 | 0 |

"hello"

Vocabulary

| "h" | [1, 0, 0, 0] |
|---|---|
| "e" | [0, 1, 0, 0] |
| "l" | [0, 0, 1, 0] |
| "o" | [0, 0, 0, 1] |

$a_2[j]$

$\hat{y}1 \quad \hat{y}2 \quad \hat{y}3 \quad \hat{y}4$

$a_2[0] \quad a_2[1] \quad a_2[2] \quad a_2[3]$

$a_1[k]$

$a_1[0] \quad a_1[1] \quad a_1[2]$

$a_0[i]$

x1  x2  x3  x4

$a_0[0] \quad a_0[1] \quad a_0[2] \quad a_0[3]$

"h"→   1   0   0   0

"h"→ "e"

"e"→ "l"

"l" → "l"
"l" → "o" **?**

Ambiguous Targets

"hel"→ "p"

"hell"→ "o"

"heliu" → "m"

Fixed Width Input/Output

"hello"

Vocabulary

| "h" | [1, 0, 0, 0] |
| "e" | [0, 1, 0, 0] |
| "l" | [0, 0, 1, 0] |
| "o" | [0, 0, 0, 1] |

"h" → "e"

"e" → "l"

"l" → "l" **?**
"l" → "o"

Ambiguous Targets

| | y1 | y2 | y3 | y4 |
|---|---|---|---|---|
| "e"→ | 0 | 1 | 0 | 0 |



$a_2[j]$

$a_1[k]$

$a_0[i]$

$\hat{y}1 \quad \hat{y}2 \quad \hat{y}3 \quad \hat{y}4$

| | x1 | x2 | x3 | x4 |
|---|---|---|---|---|
| "h"→ | 1 | 0 | 0 | 0 |

y

h

x

t

$\hat{y} = \text{softmax}(s_2)$
$s_2 = g(h, w_2)$

$h = \tanh(s_1)$
$s_1 = f(x, h_{t-1}, w_1)$

$$W_{yh} \quad\quad H$$
$$Y \begin{bmatrix} y_1h_1, y_2h_2 \ldots y_1h_k \\ y_2h_1, y_2h_2 \ldots y_2h_k \\ \\ y_nh_1, y_n\ddot{h}_2 \ldots y_nh_k \end{bmatrix}$$

$$W_{hh} \quad\quad H(t-1)$$
$$H \begin{bmatrix} h_1h_1, h_1h_2 \ldots h_1h_k \\ h_2h_1, h_2h_2 \ldots h_2h_k \\ \\ h_kh_1, h_k\ddot{h}_2 \ldots h_kh_k \end{bmatrix}$$

$$W_{xh} \quad\quad X$$
$$H \begin{bmatrix} h_1x_1, h_1x_2 \ldots h_1x_m \\ h_2x_1, h_2x_2 \ldots h_2x_m \\ \\ h_kx_1, h_k\ddot{x}_2 \ldots h_kx_m \end{bmatrix}$$

$$\hat{y} = \text{softmax}(s_2)$$
$$s_2 = g(h, w_2)$$

$$h = \tanh(s_1)$$
$$s_1 = f(x, h_{t-1}, w_1)$$

t=1

$$w_{yh} \qquad h \qquad s_2 \qquad \qquad \hat{y}$$

$$\begin{bmatrix} y_1h_1, y_2h_2 \dots y_1h_k \\ y_2h_1, y_2h_2 \dots y_2h_k \\ \\ y_nh_1, y_n\ddot{h}_2 \dots y_nh_k \end{bmatrix} * \begin{bmatrix} h_1 \\ h_2 \\ \dots \\ h_k \end{bmatrix} = \begin{bmatrix} s_1 \\ s_2 \\ \dots \\ s_k \end{bmatrix}, \text{softmax}(s_2) = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_k \end{bmatrix}$$

$$\text{softmax}(s_2) = \frac{e^{s_{2j}}}{\sum_{i=1}^{n} e^{s_{2i}}}$$

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{pmatrix}$$

y

$$w_{xh} \qquad\qquad x \qquad\qquad w_{hh} \qquad\qquad h_{t-1} \qquad s_1 \qquad\qquad h$$

$$\begin{bmatrix} h_1x_1, h_1x_2 \dots h_1x_m \\ h_2x_1, h_2x_2 \dots h_2x_m \\ \\ h_kx_1, h_k\ddot{x}_2 \dots h_kx_m \end{bmatrix} * \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_m \end{bmatrix} + \begin{bmatrix} h_1h_1, h_1h_2 \dots h_1h_k \\ h_2h_1, h_2h_2 \dots h_2h_k \\ \\ h_kh_1, h_k\ddot{h}_2 \dots h_kh_k \end{bmatrix} * \begin{bmatrix} h_1 \\ h_2 \\ \dots \\ h_k \end{bmatrix} = \begin{bmatrix} s_1 \\ s_2 \\ \dots \\ s_k \end{bmatrix}, \tanh(s_1) = \begin{bmatrix} h_1 \\ h_2 \\ \dots \\ h_k \end{bmatrix}$$

$$\begin{pmatrix} h_1 \\ h_2 \\ h_3 \end{pmatrix} \qquad \tanh(s_1) = \frac{e^{2x}-1}{e^{2x}+1}$$

h

x

$$w_{xh} \qquad\qquad x$$

$$\begin{bmatrix} h_1x_1, h_1x_2 \dots h_1x_m \\ h_2x_1, h_2x_2 \dots h_2x_m \\ \\ h_kx_1, h_k\ddot{x}_2 \dots h_kx_m \end{bmatrix} * \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_m \end{bmatrix}$$

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix}$$

$$
\overset{w_{yh}}{\begin{bmatrix} y_1h_1, y_2h_2 \ldots y_1h_k \\ y_2h_1, y_2h_2 \ldots y_2h_k \\ \\ y_nh_1, y_n\ddot{h}_2 \ldots y_nh_k \end{bmatrix}} * \overset{h}{\begin{bmatrix} h_1 \\ h_2 \\ \ldots \\ h_k \end{bmatrix}} = \overset{s_2}{\begin{bmatrix} s_1 \\ s_2 \\ \ldots \\ s_k \end{bmatrix}}, \mathrm{softmax}(s_2) = \overset{\hat{y}}{\begin{bmatrix} y_1 \\ y_2 \\ \ldots \\ y_k \end{bmatrix}}
$$

$$
\overset{w_{xh}}{\begin{bmatrix} h_1x_1, h_1x_2 \ldots h_1x_m \\ h_2x_1, h_2x_2 \ldots h_2x_m \\ \\ h_kx_1, h_k\ddot{x}_2 \ldots h_kx_m \end{bmatrix}} * \overset{x}{\begin{bmatrix} x_1 \\ x_2 \\ \ldots \\ x_m \end{bmatrix}} + \overset{w_{hh}}{\begin{bmatrix} h_1h_1, h_1h_2 \ldots h_1h_k \\ h_2h_1, h_2h_2 \ldots h_2h_k \\ \\ h_kh_1, h_k\ddot{h}_2 \ldots h_kh_k \end{bmatrix}} * \overset{h_{t-1}}{\begin{bmatrix} h_1 \\ h_2 \\ \ldots \\ h_k \end{bmatrix}} = \overset{s_1}{\begin{bmatrix} s_1 \\ s_2 \\ \ldots \\ s_k \end{bmatrix}}, \tanh(s_1) = \overset{h}{\begin{bmatrix} h_1 \\ h_2 \\ \ldots \\ h_k \end{bmatrix}}
$$

$$
\overset{w_{xh}}{\begin{bmatrix} h_1x_1, h_1x_2 \ldots h_1x_m \\ h_2x_1, h_2x_2 \ldots h_2x_m \\ \\ h_kx_1, h_k\ddot{x}_2 \ldots h_kx_m \end{bmatrix}} * \overset{x}{\begin{bmatrix} x_1 \\ x_2 \\ \ldots \\ x_m \end{bmatrix}}
$$

$$L_{total} = -\sum_{t=1}^{n} y_t \log(\hat{y})$$

$$L_t = -y_t \log(\hat{y})$$

Multi-class Cross Entropy Loss

$$L_{total} = -\sum_{t=1}^{n} y_t \log(\hat{y})$$

$$L_t = -y_t \log(\hat{y})$$

"e"  "l"  "l"

L1  L2  **L3**

$\hat{y}$  $\hat{y}$  $\hat{y}$

$$\frac{\partial L3}{\partial W_{yh}} = \frac{\partial L3}{\partial \hat{Y}3} * \frac{\partial \hat{Y}3}{\partial W_{yh}}$$
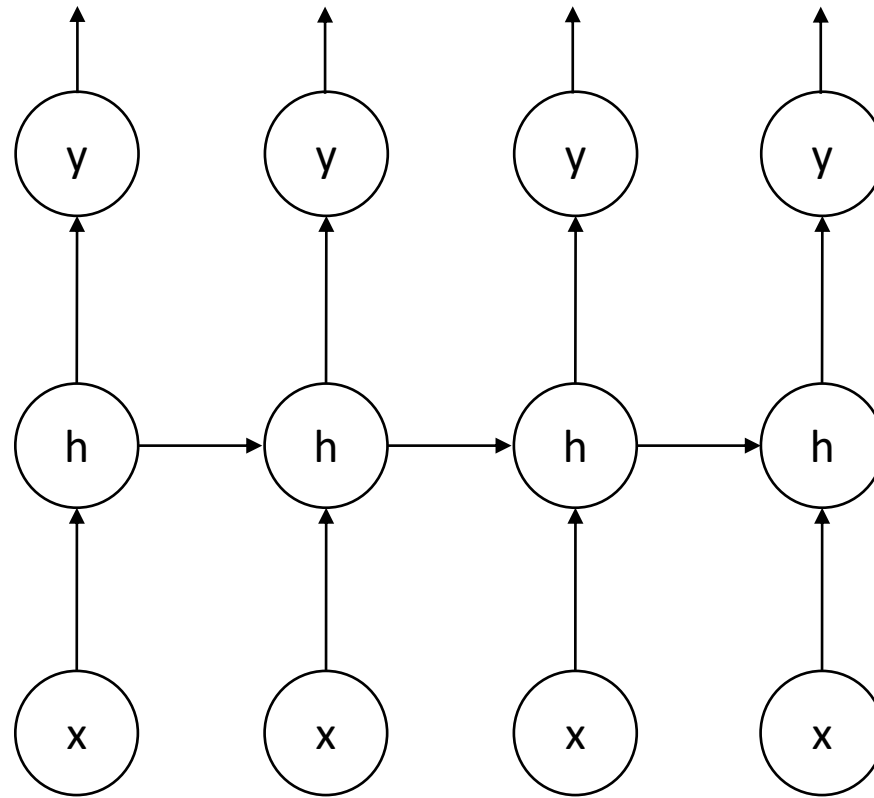
$$\frac{\partial LT}{\partial W_{yh}} = \frac{\partial LT}{\partial \hat{Y}T} * \frac{\partial \hat{Y}T}{\partial W_{yh}}$$

$\boldsymbol{W_{yh}}$

$h_0$  h  h  **h**

"h"  "e"  "l"

t=1  t=2  t=3

For derivatives of Softmax and Cross Entropy Loss
Part 2: Softmax Regression (saitcelebi.com)
How to compute the derivative of softmax and cross-entropy – Charlee Li

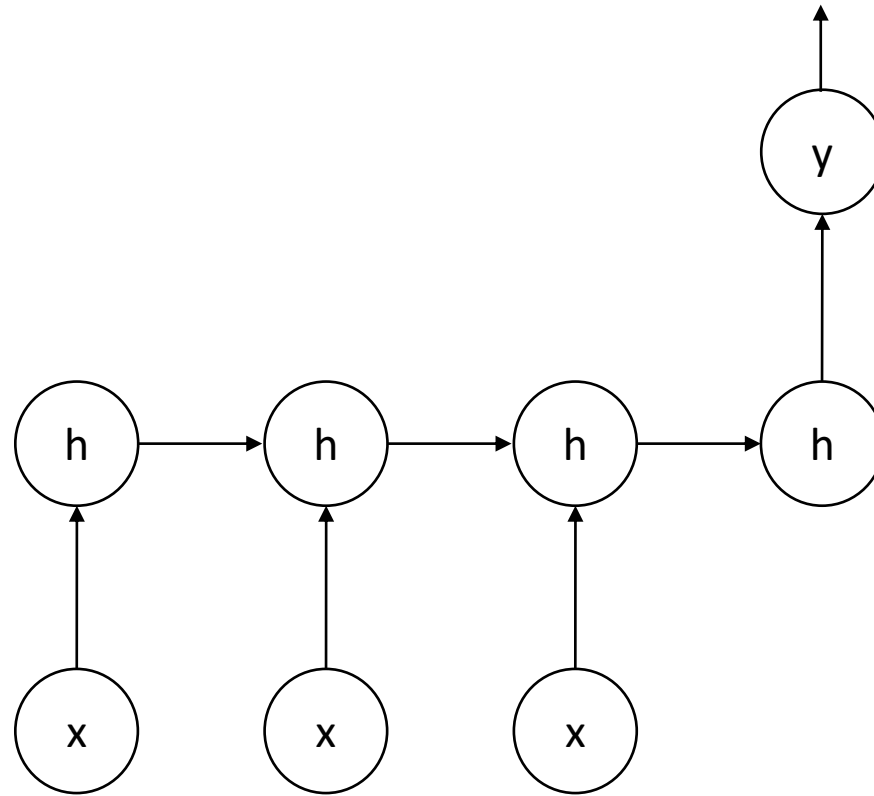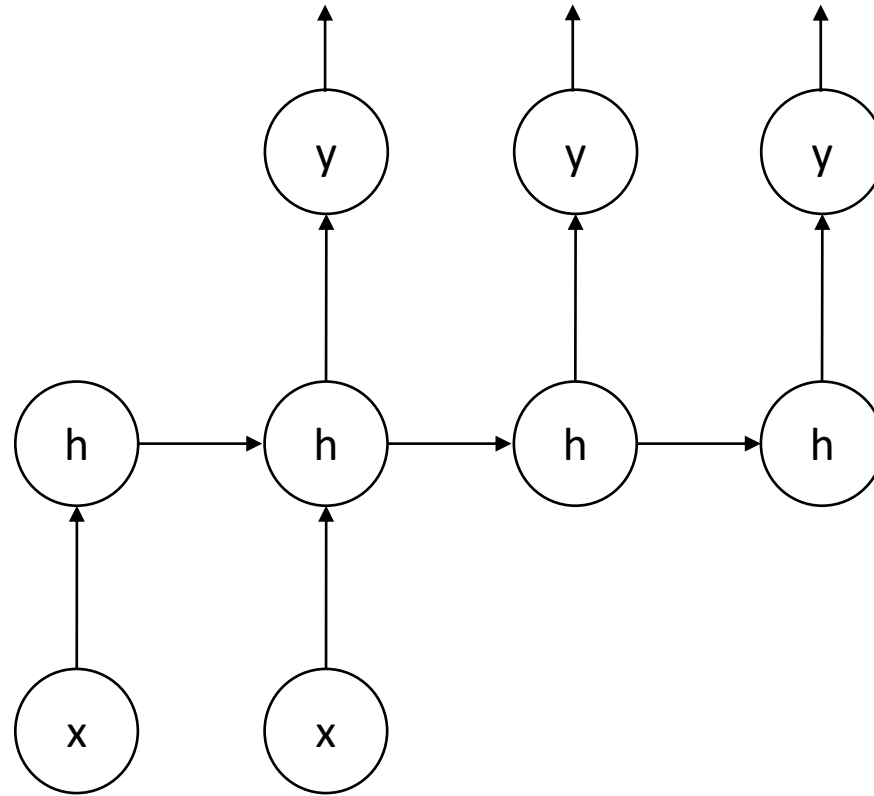$$L_{total} = -\sum_{t=1}^{n} y_t \log(\hat{y})$$

$$L_t = -y_t \log(\hat{y})$$

"e"   "l"   "l"

L1    L2    **L3**

$$\frac{\partial L3}{\partial W_{hh}} = \frac{\partial L3}{\partial \hat{Y}3} * \frac{\partial \hat{Y}3}{\partial H3} * \frac{\partial H3}{\partial W_{hh}} +$$

$$\frac{\partial L3}{\partial \hat{Y}} * \frac{\partial \hat{Y}3}{\partial H3} * \frac{\partial H3}{\partial H2} * \frac{\partial H2}{\partial W_{hh}} +$$

$$\frac{\partial L3}{\partial \hat{Y}} * \frac{\partial \hat{Y}3}{\partial H3} * \frac{\partial H3}{\partial H2} * \frac{\partial H2}{\partial H1} * \frac{\partial H1}{\partial W_{hh}}$$

$\hat{y}$   $\hat{y}$   $\hat{y}$

$h_0$   $w_{hh}$   h   $w_{hh}$   h   $w_{hh}$   h

$$\frac{\partial LT}{\partial W_{hh}} = \sum_{t=1}^{T} \frac{\partial LT}{\partial \hat{Y}T} * \frac{\partial \hat{Y}T}{\partial H_t} * \frac{\partial H_t}{\partial W_{hh}}$$

"h"   "e"   "l"

t=1   t=2   t=3

For derivatives of Softmax and Cross Entropy Loss
Part 2: Softmax Regression (saitcelebi.com)
How to compute the derivative of softmax and cross-entropy – Charlee Li

$$L_{total} = -\sum_{t=1}^{n} y_t \log(\hat{y})$$

$$L_t = -y_t \log(\hat{y})$$

"e"      "l"      "l"

L1      L2      **L3**

$$\frac{\partial L3}{\partial W_{hh}} = \frac{\partial L3}{\partial \hat{Y}3} * \frac{\partial \hat{Y}3}{\partial H3} * \frac{\partial H3}{\partial W_{hx}} +$$

$$\frac{\partial L3}{\partial \hat{Y}} * \frac{\partial \hat{Y}3}{\partial H3} * \frac{\partial H3}{\partial H2} * \frac{\partial H2}{\partial W_{hx}} +$$

$$\frac{\partial L3}{\partial \hat{Y}} * \frac{\partial \hat{Y}3}{\partial H3} * \frac{\partial H3}{\partial H2} * \frac{\partial H2}{\partial H1} * \frac{\partial H1}{\partial W_{hx}}$$

ŷ      ŷ      ŷ

$h_0$      h      h      h

$w_{xh}$      $w_{xh}$      $w_{xh}$

$$\frac{\partial LT}{\partial W_{xh}} = \sum_{t=1}^{T} \frac{\partial LT}{\partial \hat{Y}T} * \frac{\partial \hat{Y}T}{\partial H_t} * \frac{\partial H_t}{\partial W_{xh}}$$

"h"      "e"      "l"

t=1      t=2      t=3

For derivatives of Softmax and Cross Entropy Loss
Part 2: Softmax Regression (saitcelebi.com)
How to compute the derivative of softmax and cross-entropy – Charlee Li

## One-to-One

**One-to-Many**

## Many-to-One

## Many-to-Many

Generation

# Vanishing Gradients



The      **bird**      that      I      really      really      really      liked      was      **flying**

0 = forget
1 = retain

= Signal Strength

$c_{t-1}$

X

+

c

X

+

X

X

$h_{t-1}$

+

W

h

x

Illustrated Guide to LSTM's and GRU's: A step by step explanation | by Michael Phi | Towards Data Science
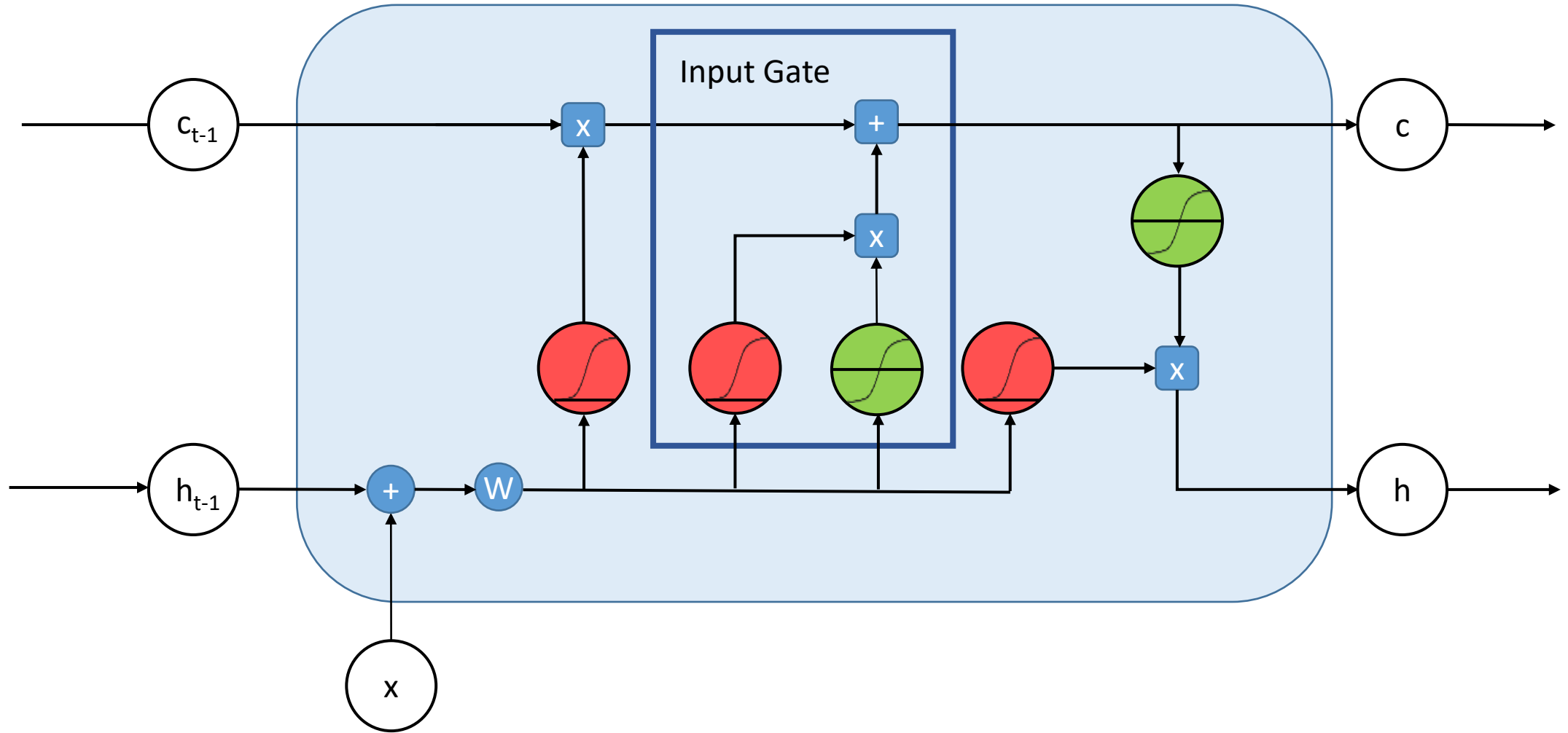
Long Short Term Memory

0 = forget
1 = retain

= Signal Strength

Vector Concatenation

Weight Multiplication

$c_{t-1}$

$c$

$h_{t-1}$

$h$

Long Short Term Memory

0 = forget
1 = retain

= Signal Strength

Previous Cell State

$c_{t-1}$

$h_{t-1}$

x

c

h

Long Short Term Memory

quote detection cell



line length tracking cell

CS 230 - Recurrent Neural Networks Cheatsheet (stanford.edu)