# Lec 7. COMPUTATION GRAPH     why we have Back Prop

Forward Pass → Forward Propagation

↙

Then Backward Propagation.

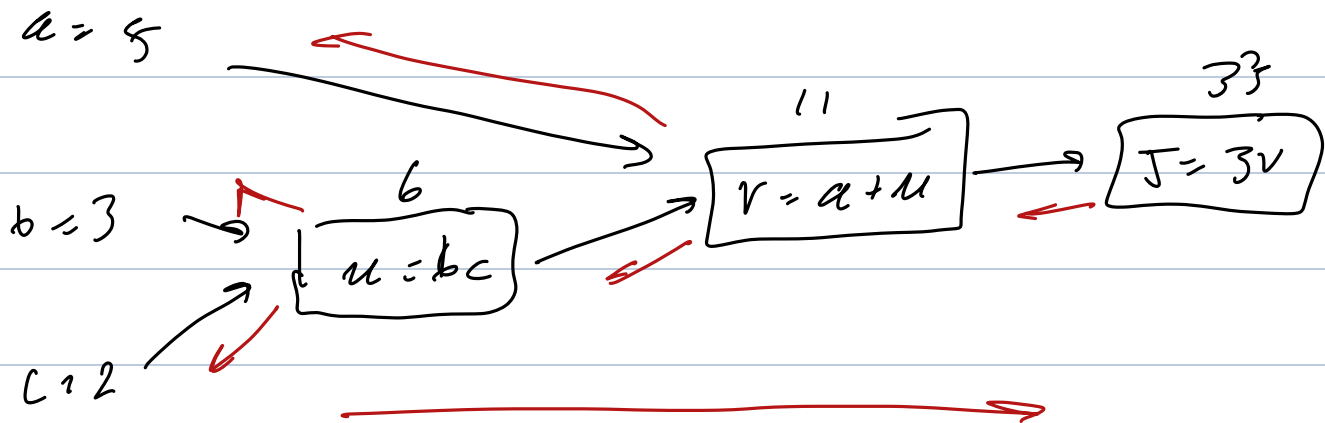## Computation Graph    example

$$J(a,b,c) = 3\underbrace{(a + \underbrace{bc}_{u})}_{V}$$

$$\underbrace{\phantom{3(a+bc)}}_{J}$$

$$u = bc$$

$$V = a + u$$

$$J = 3V$$

$a = 5$

$b = 3$

$c = 2$

$u = bc$   6

$V = a + u$   11

$J = 3v$   33

Left to Right, can use to compute J Pass.

Compute Derivate do a Right to Left Pass

$a = 5$

$b = 3$

$c = 2$



$6$

$u = bc$

$11$

$v = a + u$

$33$

$J = 3v$

$$\frac{dJ}{dv} = ? \qquad J = 3v \qquad \frac{dJ}{dv} = 3$$

$$\frac{dJ}{da} = \frac{\partial v}{\partial a} = 3(a + u)$$

Chain Rule.    $a \to v \to J.$

$$\frac{dJ}{da} = \frac{\partial J}{\partial v} \frac{\partial v}{\partial a}.$$

$$\frac{dJ}{dc} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial u} , \frac{\partial u}{\partial c}$$

$$= (3)(1)(b)$$

$$= b$$

For Code $\to$ Final Var we care about $\to$. Last node in computation graph is $J$.

$$\frac{d \text{ Final output var}}{d \text{ var}} = $$

$$\frac{dJ}{du} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial u} = (3)(1) = 3.$$

$$\frac{dJ}{db} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial u} \cdot \frac{\partial u}{\partial b}$$

$$= (3)(1)(c)$$

$$= 3c \quad \text{if } c = 2 \Rightarrow 6.$$

$$\frac{dJ}{dc} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial u} \cdot \frac{\partial u}{\partial c} = \begin{array}{l} (3)(1)(b) \\ = (3)(1)(3) = 9. \end{array}$$

<span style="color:red">So backwards prop allows for the chain rule to compute the derivatives</span>

<u>Lec 9</u>  <u>Logistic Regression Gradient Descent</u>
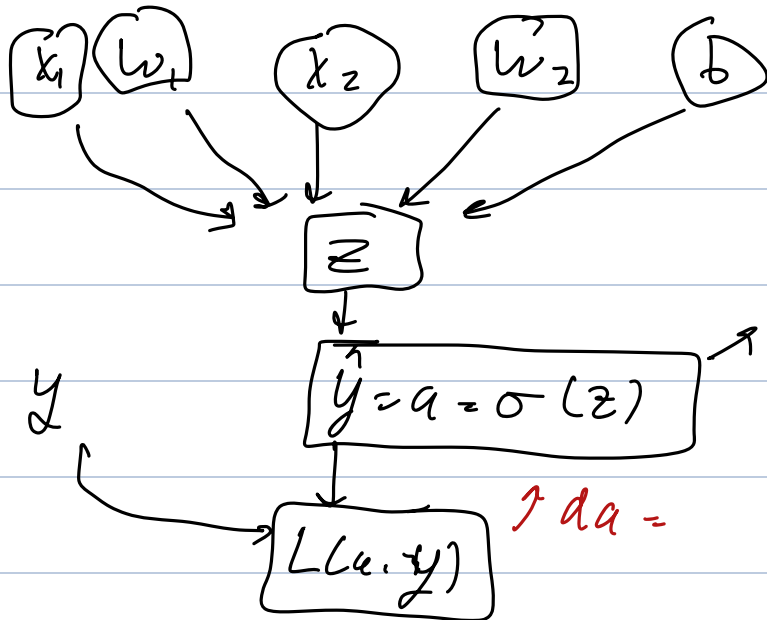
$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z)$$

$$L(a,y) = - (y \log(a) + (1-y) \log(1-a))$$

2 featured $x_1$ and $x_2$.

$x_1$  $x_2$

$$\Rightarrow z = w_1 x_1 + w_2 x_2 + b$$



$$\frac{da}{dz} = \frac{1}{1+e^{-z}} \Rightarrow a(1-a)$$

$\uparrow da =$

Modify $w$ and $b$ to minimize $L(a,y)$.

$da$.

$$\frac{\partial L(a,y)}{\partial w_1} = \frac{\partial (L(a,y))}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial w_1}$$

the weight input depend on + activation of last neuron

$$= -\left( \frac{y}{a} + \frac{(1-y)}{1-a}(-1) \right) \left( \sigma'(z) \right) (x_1)$$

$$= \left( \frac{-y}{a} + \frac{1-y}{1-a} \right) \left( \sigma'(z) \right) (x_1)$$

$$= (a-y)(x_1)$$

$$\frac{dL}{\partial w_1} = "dw_1" = x_1 \cdot dz$$

$$\frac{dL}{dw_2} = u \, dw_2^i = x_2 \cdot dz$$

$$db = dz.$$

Then $\rightarrow$

$$w_1 := w_1 - \alpha \, dw_1$$

One step of
gradient dscnt
w.r.t one
example

$$w_2 := w_2 - \alpha \, dw_2$$

$$b := b - \alpha \, db$$

---

## Lec9. Gradient Descent on $M$ Examples.

$\leftarrow$ Prediction $\hat{y}^{(i)}$

$$J(w,b) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(a^{(i)}, y^{(i)})$$

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

$$dw^{(i)}, \, dw_2^{(i)}, \, db^{(i)} \quad \underline{\text{from before}} \text{ just on } (x^{(i)}, y^{(i)})$$

$$\therefore \frac{\partial}{\partial w_1} J(w,b) = \frac{1}{m} \sum_{i=1}^{m} \underbrace{\frac{\partial}{\partial w_1} \mathcal{L}(a^{(i)}, y^{(i)})}_{\partial w_1^{(i)}}$$

$\hookrightarrow$ So this
is just the avg of these steps.

$\hookrightarrow$ take into acnt all of them.

Log Reg on M examples

$J = 0, \quad dw_1 = 0, \quad dw_2 = 0, \quad db = 0.$

For $i=1$ to $m$. → $O(m)$ (samples)

$$z^{(i)} = w^T x^{(i)} + b$$

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

$$J += (-[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log(1 - a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 \;+= \; x_1^{(i)} dz^{(i)}$$

$$dw_2 \;+= \; x_2 dz^{(i)}$$

$$db \;+= \; dz^{(i)}$$

$n=2 \uparrow$ 2 features

↳ for loop $O(n)$ (features)

$J /= m; \quad$ avg.

$dw_1 /= m$

$dw_2 /= m$

$db /= m$

$$\therefore \Rightarrow dw_1 = \frac{\partial J}{\partial w_1}$$

to sum over the entire set

1 step of gradient dcent.

$$w_1 := w_1 - \alpha \, dw_1$$

$$w_2 := w_2 - \alpha \, dw_2$$

$$b := b - \alpha \, db$$

One step of gradient step, would have to repeat all of this for <u>multiple steps</u>, to lower the <u>cost</u>

VECTORIZATION TECHNIQUES TO GET RID OF
EXPLICIT FOR LOOPS TO SPEED UP, ELSE
WILL TAKE <u>TOO LONG</u>

§§ Check note on <u>the derivative.</u>