

JACOBIAN (OPTIONAL MATERIAL)

JACOBIAN IS THE DERIVATIVE OF MULTI
INPUT - MULTI OUTPUT FUNCTIONS

$$\bar{y} = f(\bar{x})$$

DEFINITION (FOR VECTOR IN - VECTOR OUT)

$$\frac{\partial f(\bar{x})}{\partial \bar{x}} = \begin{bmatrix} \frac{\partial f_1(x)}{\partial x_1}, \frac{\partial f_1(x)}{\partial x_2}, \dots, \frac{\partial f_1(x)}{\partial x_n} \\ \vdots \\ \frac{\partial f_m(x)}{\partial x_1}, \frac{\partial f_m(x)}{\partial x_2}, \dots, \frac{\partial f_m(x)}{\partial x_n} \end{bmatrix}$$

NO OF
OUTPUTS

NO. OF INPUTS

SOMETIMES THE JACOBIAN IS DEFINED
AS A TRANSPOSE OF THIS. IT DOESN'T MATTER
UNTIL IT IS CONSISTENT.

FOR MATRIX-VECTOR MULTIPLICATION

$$\bar{y} = \underline{A} \bar{x} \rightarrow \bar{y} = f(\bar{x})$$

NO OF INPUTS; NO. OF ELEMENTS OF \bar{x}
NO OF OUTPUTS; NO OF ELEMENTS OF \bar{y}

$$\frac{\partial \underline{A} \bar{x}}{\partial \bar{x}} = \underbrace{\left[\quad \right]}_{\text{NO OF } x} \left\{ \text{NO OF } y \right\}$$

$$\underline{A} \bar{x} = \underbrace{\begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \vdots & & & \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}}_{\text{NO OF } y} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}}_{\text{NO OF } x} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

$$y_i = a_{i1} x_1 + a_{i2} x_2 + \dots + a_{in} x_n$$

$$\frac{\partial y_i}{\partial x_j} = a_{ij}$$



$$\frac{\partial \underline{A} \bar{x}}{\partial \bar{x}} = \underline{A}$$

JACOBIAN WITH RESPECT TO THE MATRIX IS A 3D TENSOR. FOR A MATRIX INPUT - MATRIX OUTPUT FUNCTIONS, IT IS A 4D TENSOR.

$$\bar{y}_{[m]} = \bar{A}_{[m \times n]} \bar{x}_{[n]}$$

$$\frac{\partial \bar{y}}{\partial \bar{A}} = \text{TENSOR OF } \underbrace{m \times m}_{\text{NO OF OUTPUTS}} \times \underbrace{n}_{\text{NO OF INPUTS}}$$

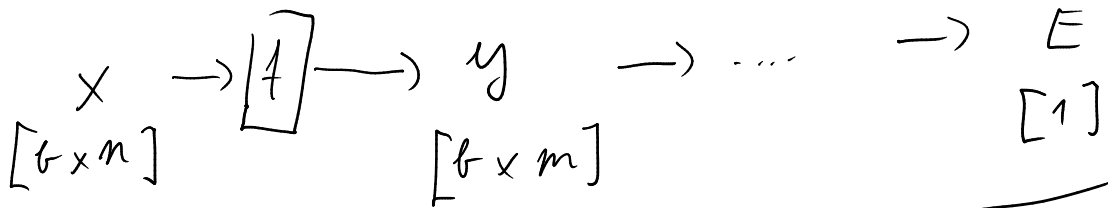
$$\frac{\partial \bar{y}}{\partial \bar{A}}_{ijk} = \frac{\partial y_i}{\partial a_{jk}}$$

BUT FORTUNATELY WE ALWAYS HAVE A SCALAR OUTPUT FUNCTION, BECAUSE THE LOSS IS ALWAYS A SCALAR.

NOTE THAT WE ARE ALWAYS CALCULATING THE GRADIENT WITH RESPECT TO THE LOSS, WHICH IS SCALAR.

FOR BATCHED TRAINING, DATA IS A $[b \times m]$ MATRIX. EACH TRANSFORMATION CAN CHANGE THE NUMBER OF CHANNELS, BUT NOT THE NUMBER OF BATCHES, THUS THE OUTPUT IS $[b \times m]$.

FINALLY, WE HAVE THE REST OF THE NETWORK, THAT REDUCE ALL DATA TO A SINGLE SCALAR.



$$y = f(x)$$

$$E = g(y)$$

— THE REST OF THE NETWORK, INCLUDING LOSS.

$$\frac{\partial f(x)}{\partial x} : \begin{matrix} i & j \\ b \times m & \times & b \times m \\ \underbrace{\hspace{1cm}} & & \underbrace{\hspace{1cm}} \\ \text{OUTPUT} & & \text{INPUT} \end{matrix}$$

$$\frac{\partial E}{\partial y} : \begin{matrix} & i & j \\ 1 \times & b \times m \\ \underbrace{\hspace{1cm}} & & \underbrace{\hspace{1cm}} \\ \text{OUTPUT} & & \text{INPUT} \end{matrix} \sim b \times m$$

WE CAN IGNORE THE 1, IT IS ALWAYS 1 FOR SCALAR OUTPUT

$$\frac{\partial E}{\partial x} = \frac{\partial E}{\partial y} \otimes \underbrace{(2)}_{\text{GENERALIZED MATRIX MULTIPLICATION}} \frac{\partial y}{\partial x}$$

GENERALIZED MATRIX MULTIPLICATION

HOW MANY DIMENSIONS TO REDUCE (SUM OVER)

$$[i, j] \otimes_2 [i, j, k, l] \Rightarrow [k, l]$$

IN OUR CASE

$$\underbrace{X}_{[l \times n]} \underbrace{A}_{[n \times m]} = \underbrace{y}_{[l \times m]} \quad E = g(y)$$

WE ARE INTERESTED IN $\frac{\partial E}{\partial A}$ AND $\frac{\partial E}{\partial X}$.

$$\frac{\partial E}{\partial y} : \underbrace{1 \times l \times m}_{\text{OUT}} \underbrace{\sim l \times m}_{\text{IN}} \text{ MATRIX}$$

$$\frac{\partial y}{\partial A} : \underbrace{l \times m \times n \times m}_{\text{OUT}} \underbrace{\sim n \times m}_{\text{IN}} \text{ TENSOR}$$

$$\frac{\partial E}{\partial A} : \underbrace{1 \times n \times m}_{\text{OUT}} \underbrace{\sim n \times m}_{\text{IN}} \leftarrow \text{WE ARE INTERESTED IN THIS}$$

CHAIN RULE : $\frac{\partial E}{\partial A} = \frac{\partial E}{\partial y} \otimes_2 \frac{\partial y}{\partial A}$

$$l \begin{bmatrix} \dots \\ X \\ \dots \end{bmatrix}_n \begin{bmatrix} \dots \\ y \\ \dots \end{bmatrix}_m = \begin{bmatrix} \dots \\ A \\ \dots \end{bmatrix}_m$$

$$y_{ij} = \sum_k x_{ik} A_{kj}$$

DEF OF MATMUL

$$\frac{\partial y}{\partial A} \text{ ijk l} = \frac{\partial y_{ij}}{\partial A_{kl}} = \begin{cases} 0 & j \neq l \\ x_{ik} & j = l \end{cases}$$

$$\frac{\partial y}{\partial A}{}_{ij}{}_{kl} = \begin{cases} 0 & j \neq l \\ x_{il} & j = l \end{cases}$$

$$\frac{\partial E}{\partial A} = \frac{\partial E}{\partial y} \otimes \frac{\partial y}{\partial A}$$

$$[kl] [kl \quad ij]$$

$$\frac{\partial E}{\partial A}{}_{ij} = \sum_k \sum_l \underbrace{\frac{\partial E}{\partial y}{}_{kl} \frac{\partial y}{\partial A}{}_{kl \quad ij}}_{0 \text{ if } l \neq j \Rightarrow l = j}$$

$$\frac{\partial E}{\partial A}{}_{ij} = \sum_k \frac{\partial E}{\partial y}{}_{kj} \frac{\partial y}{\partial A}{}_{kj \quad ij}$$

$$\frac{\partial E}{\partial A}{}_{ij} = \sum_k \frac{\partial E}{\partial y}{}_{kj} x_{ki}$$

DEF OF MATMUL

$$\left(\frac{\partial E}{\partial y} \right)^T \times : [j, i]$$

MATRIX

BUT WE

WANT i AS

FIRST INDEX

TRANSPOSE!

$$\frac{\partial E}{\partial A} = \left(\left(\frac{\partial E}{\partial y} \right)^T \times \right)^T$$

$$\boxed{\frac{\partial E}{\partial A} = X^T \frac{\partial E}{\partial y}}$$

FOR MORE DETAILS, CHECK:
MML BOOK, PAGE 149 ([MML-1300K.GITHUB.IO](#)).