

All The BPs

CMPUT 328

Nilanjan Ray

OK, let's apply chain rule to a computational graph where all variables and parameters are scalars

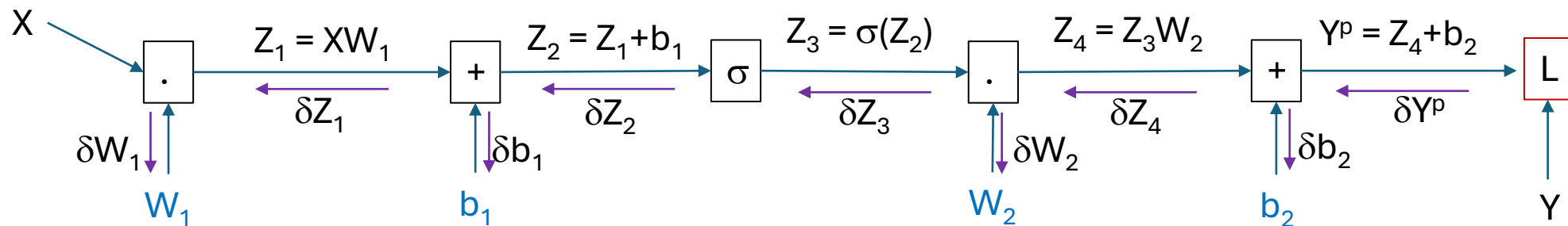
So, our scalar neural net is:

$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2$$

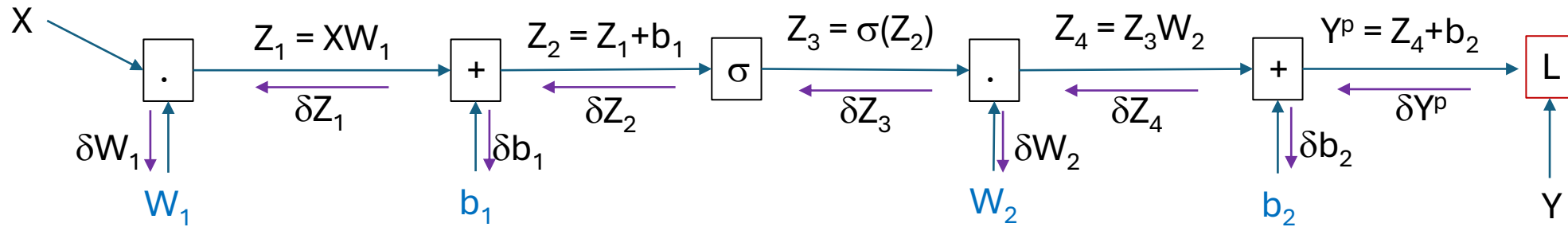
with a square (aka Euclidean) loss function: $L(Y^p, Y) = \frac{1}{2}(Y^p - Y)^2$

As usual, X is the input, Y^p is the output, and W_1, b_1, W_2, b_2 are parameters of the neural net, σ is a non-linear function.

The computational graph for this scalar neural net is (also showing loss gradient symbols):



Chain rule for a scalar neural net...



$$1 \quad \delta Y^p \equiv \frac{\partial L}{\partial Y^p} = \frac{\partial}{\partial Y^p} \left[\frac{1}{2} (Y^p - Y)^2 \right] = y^p - y$$

$$4 \quad \delta Z_2 \equiv \frac{\partial L}{\partial Z_2} = \frac{\partial Z_3}{\partial Z_2} \frac{\partial L}{\partial Z_3} = \sigma'(Z_2) \delta Z_3$$

$$2 \quad \delta Z_4 \equiv \frac{\partial L}{\partial Z_4} = \frac{\partial Y^p}{\partial Z_4} \frac{\partial L}{\partial Y^p} = \delta Y^p$$

Because $Z_3 = \sigma(Z_2)$, $\frac{\partial Z_3}{\partial Z_2} = \sigma'(Z_2)$

Because $Y^p = Z_4 + b_2$, $\frac{\partial Y^p}{\partial Z_4} = 1$

$$5 \quad \delta Z_1 \equiv \frac{\partial L}{\partial Z_1} = \frac{\partial Z_2}{\partial Z_1} \frac{\partial L}{\partial Z_2} = \delta Z_2$$

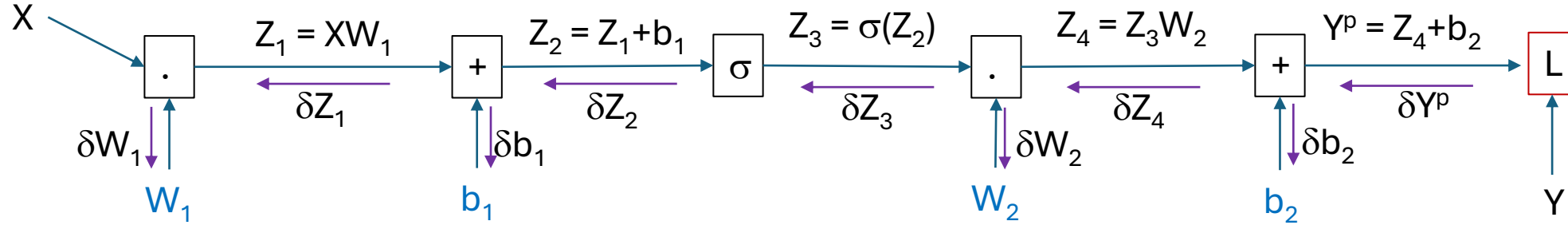
$$3 \quad \delta Z_3 \equiv \frac{\partial L}{\partial Z_3} = \frac{\partial Z_4}{\partial Z_3} \frac{\partial L}{\partial Z_4} = W_2 \delta Z_4$$

Because $Z_4 = Z_3 W_2$, $\frac{\partial Z_4}{\partial Z_3} = W_2$

Because $Z_2 = Z_1 + b_1$, $\frac{\partial Z_2}{\partial Z_1} = 1$

But we need loss derivatives with respect to parameters...

Loss derivatives w.r.t. parameters



$$1 \quad \delta W_1 \equiv \frac{\partial L}{\partial W_1} = \frac{\partial Z_1}{\partial W_1} \frac{\partial L}{\partial Z_1} = X \delta Z_1$$

$$\text{Because } Z_1 = XW_1, \frac{\partial Z_1}{\partial W_1} = X$$

$$2 \quad \delta b_1 \equiv \frac{\partial L}{\partial b_1} = \frac{\partial Z_2}{\partial b_1} \frac{\partial L}{\partial Z_2} = \delta Z_2$$

$$\text{Because } Z_2 = Z_1 + b_1, \frac{\partial Z_2}{\partial b_1} = 1$$

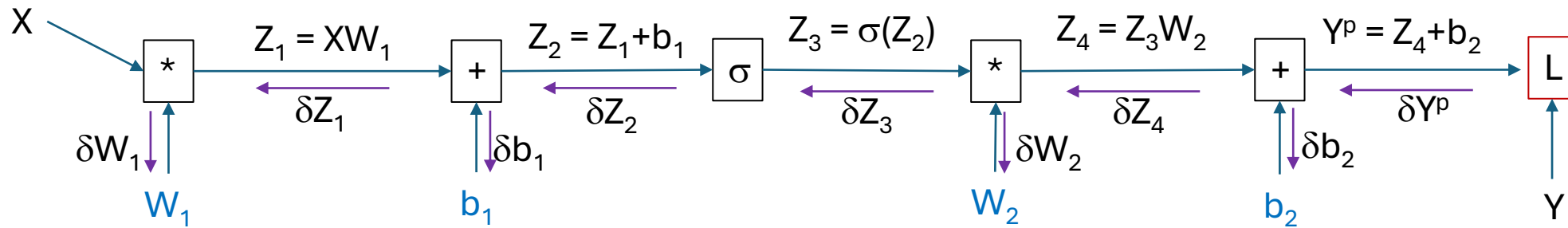
$$3 \quad \delta W_2 \equiv \frac{\partial L}{\partial W_2} = \frac{\partial Z_4}{\partial W_2} \frac{\partial L}{\partial Z_4} = Z_3 \delta Z_4$$

$$\text{Because } Z_4 = Z_3W_2, \frac{\partial Z_4}{\partial W_2} = Z_3$$

$$4 \quad \delta b_2 \equiv \frac{\partial L}{\partial b_2} = \frac{\partial Y^p}{\partial b_2} \frac{\partial L}{\partial Y^p} = \delta Y^p$$

$$\text{Because } Y^p = Z_4 + b_2, \frac{\partial Y^p}{\partial b_2} = 1$$

Chain rule for a (general) neural net



$$1 \quad \delta Y^p \equiv \frac{\partial L}{\partial Y^p} = \frac{\partial}{\partial Y^p} \left[\frac{1}{2} \|Y^p - Y\|^2 \right] = y^p - y$$

$$4 \quad \delta Z_2 \equiv \frac{\partial L}{\partial Z_2} = \frac{\partial Z_3}{\partial Z_2} \frac{\partial L}{\partial Z_3} = \sigma'(Z_2) \cdot \delta Z_3$$

$$2 \quad \delta Z_4 \equiv \frac{\partial L}{\partial Z_4} = \frac{\partial Y^p}{\partial Z_4} \frac{\partial L}{\partial Y^p} = \delta Y^p$$

Because $Z_3 = \sigma(Z_2)$, $\frac{\partial Z_3}{\partial Z_2} = \sigma'(Z_2)$

Because $Y^p = Z_4 + b_2$, $\frac{\partial Y^p}{\partial Z_4} = 1$

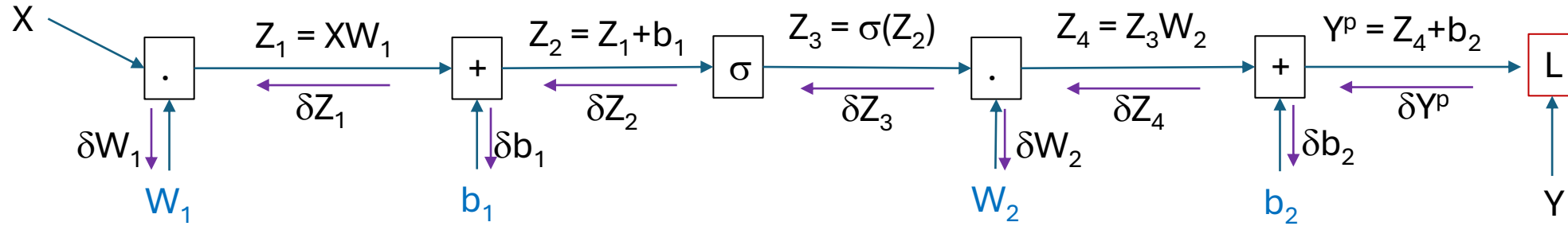
$$5 \quad \delta Z_1 \equiv \frac{\partial L}{\partial Z_1} = \frac{\partial Z_2}{\partial Z_1} \frac{\partial L}{\partial Z_2} = \delta Z_2$$

$$3 \quad \delta Z_3 \equiv \frac{\partial L}{\partial Z_3} = \frac{\partial Z_4}{\partial Z_3} \frac{\partial L}{\partial Z_4} = \delta Z_4 W_2^T$$

Because $Z_4 = Z_3 W_2$, $\frac{\partial Z_4}{\partial Z_3} = W_2^T$

Because $Z_2 = Z_1 + b_1$, $\frac{\partial Z_2}{\partial Z_1} = 1$

Loss derivatives w.r.t. **matrix or vector** parameters



$$1 \quad \delta W_1 \equiv \frac{\partial L}{\partial W_1} = \frac{\partial Z_1}{\partial W_1} \frac{\partial L}{\partial Z_1} = \mathbf{X}^T \delta \mathbf{Z}_1$$

$$\text{Because } Z_1 = XW_1, \frac{\partial Z_1}{\partial W_1} = \mathbf{X}^T$$

$$3 \quad \delta W_2 \equiv \frac{\partial L}{\partial W_2} = \frac{\partial Z_4}{\partial W_2} \frac{\partial L}{\partial Z_4} = \mathbf{Z}_3^T \delta \mathbf{Z}_4$$

$$\text{Because } Z_4 = Z_3W_2, \frac{\partial Z_4}{\partial W_2} = \mathbf{Z}_3^T$$

$$2 \quad \delta b_1 \equiv \frac{\partial L}{\partial b_1} = \frac{\partial Z_2}{\partial b_1} \frac{\partial L}{\partial Z_2} = \sum_k (\delta \mathbf{Z}_2)_{k,:}$$

$$\text{Because } Z_2 = Z_1 + b_1, \frac{\partial Z_2}{\partial b_1} = [\mathbf{1}, \dots, \mathbf{1}]$$

$$4 \quad \delta b_2 \equiv \frac{\partial L}{\partial b_2} = \frac{\partial Y^p}{\partial b_2} \frac{\partial L}{\partial Y^p} = \sum_k (\delta \mathbf{Y}^p)_{k,:}$$

$$\text{Because } Y^p = Z_4 + b_2, \frac{\partial Y^p}{\partial b_2} = [\mathbf{1}, \dots, \mathbf{1}]$$

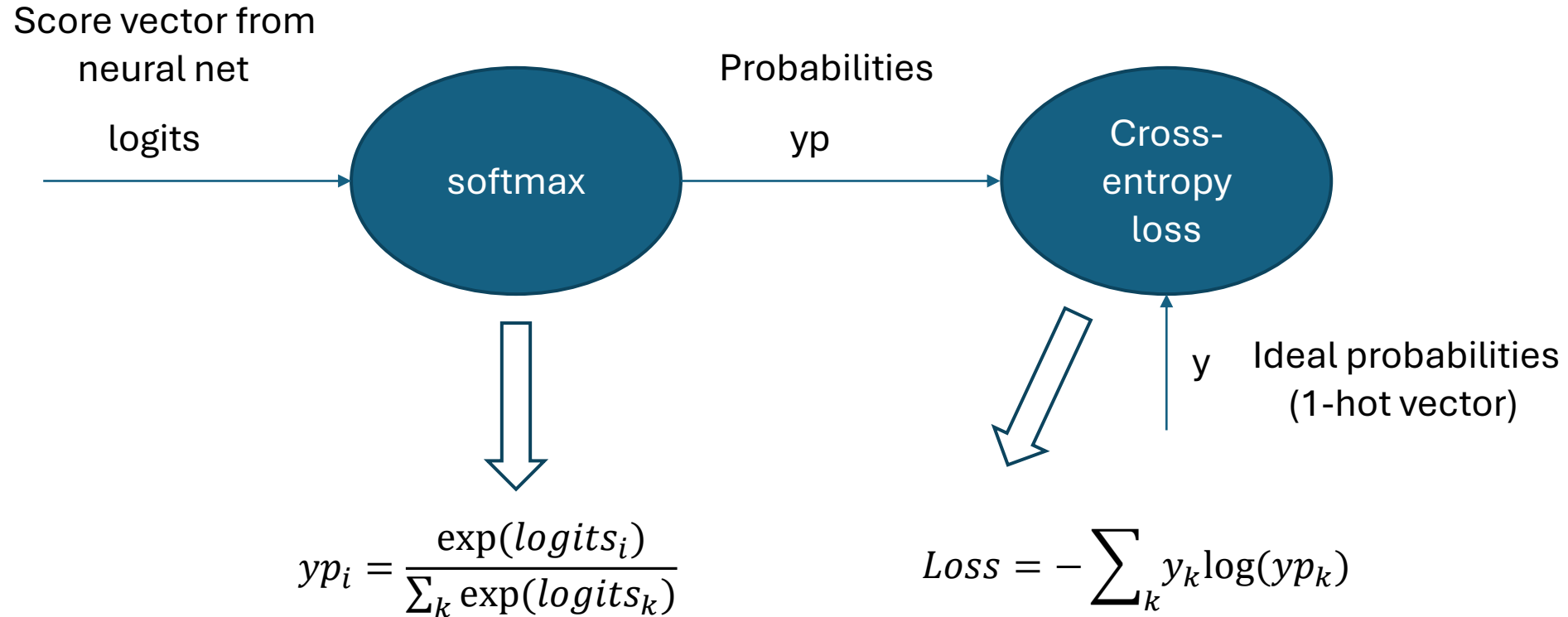
Backprop derivation for loss function

Euclidean loss function: $Loss(Y^p, Y) = \frac{1}{2} \|Y^p - Y\|^2 = \frac{1}{2} \sum_i (Y_i^p - Y_i)^2$

i^{th} component of δY^p vector: $\delta Y_i^p = \frac{\partial}{\partial Y_i^p} Loss(Y^p, Y) = \frac{\partial}{\partial Y_i^p} \frac{1}{2} \sum_k (Y_k^p - Y_k)^2 = Y_i^p - Y_i$

Using vector notation: $\delta Y^p = Y^p - Y$

Softmax and cross-entropy loss



To backpropagate error, we need to compute: $\delta(\text{logits})_i \equiv \frac{\partial(\text{Loss})}{\partial(\text{logits})_i}$

Softmax and cross-entropy loss: backprop

Score vector from
neural net

logits

softmax

Probabilities

yp

Cross-
entropy
loss

y Ideal probabilities
(1-hot vector)

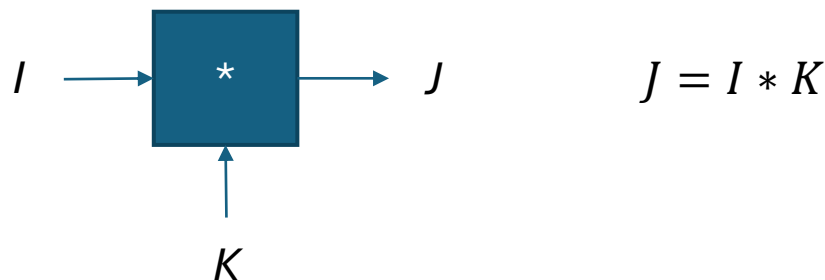
$$Loss = - \sum_k y_k \log(yp_k) \longrightarrow \frac{\partial(Loss)}{\partial(yp)_k} = -\frac{y_k}{yp_k}$$

$$yp_i = \frac{\exp(logits_i)}{\sum_k \exp(logits_k)} \longrightarrow \frac{\partial(yp)_k}{\partial(logits)_i} = \begin{cases} yp_i(1 - yp_i), & \text{if } i = k, \\ -yp_i yp_k, & \text{otherwise.} \end{cases}$$

Using the above two results in the chain rule, $\delta(logits)_i \equiv \frac{\partial(Loss)}{\partial(logits)_i} = \sum_k \frac{\partial(yp)_k}{\partial(logits)_i} \frac{\partial(Loss)}{\partial(yp)_k} = yp_i - y_i$

What if, instead of cross-entropy, we used L2 loss along with softmax?

Backpropagation (BP) for a conv layer



Input to convolution layer: I , a H -by- W matrix

Parameter of the layer: K , a h -by- w matrix

Output of the layer: J , a $(H-h+1)$ -by- $(W-w+1)$ matrix

Assume: $H \geq h$ and $W \geq w$

Given the gradient of loss function δJ with respect to J , BP tries to find answers to the following:

(1) What is the gradient of the loss function with respect to K ? Denote this gradient by δK .

(2) What is the gradient of the loss function with respect to I ? Denote this gradient by δI .

Why do we need δK ? Because, we want to adjust the parameter K by gradient descent: $K = K - (\text{learning rate})\delta K$

Why do we need δI ? Because, we want to apply BP to the layer that precedes this conv layer.

Derivation of δK

$$J(i, j) = \sum_{l=1}^h \sum_{m=1}^w I(i + l - 1, j + m - 1) K(l, m) \quad \longrightarrow \quad \frac{\partial J(i, j)}{\partial K(p, q)} = I(i + p - 1, j + q - 1)$$

Using chain rule of derivative:

$$\delta K(p, q) = \sum_{i=1}^{H-h+1} \sum_{j=1}^{W-w+1} \frac{\partial J(i, j)}{\partial K(p, q)} \delta J(i, j) = \sum_{i=1}^{H-h+1} \sum_{j=1}^{W-w+1} I(i + p - 1, j + q - 1) \delta J(i, j)$$

Thus, $\boxed{\delta K = I * \delta J}$

Derivation of δI

$$J(i, j) = \sum_{l=1}^h \sum_{m=1}^w I(i + l - 1, j + m - 1) K(l, m) \quad \longrightarrow$$

$$\frac{\partial J(i, j)}{\partial I(p, q)} = \begin{cases} K(p - i + 1, q - j + 1), & \text{if } 0 \leq p - i \leq h - 1 \text{ and } 0 \leq q - j \leq w - 1, \\ 0, & \text{otherwise.} \end{cases}$$

Using chain rule of derivative:

$$\begin{aligned} \delta I(p, q) &= \sum_{i=1}^{H-h+1} \sum_{j=1}^{W-w+1} \frac{\partial J(i, j)}{\partial I(p, q)} \delta J(i, j) = \sum_{i=\max(1, p-h+1)}^{\min(p, H-h+1)} \sum_{j=\max(1, q-w+1)}^{\min(q, W-w+1)} K(p - i + 1, q - j + 1) \delta J(i, j) \\ &= \sum_{l=\max(1, p+h-H)}^{\min(p, h)} \sum_{m=\max(1, q+w-W)}^{\min(q, w)} K(l, m) \delta J(p - l + 1, q - m + 1) \end{aligned}$$

Thus, $\delta I = \text{pad}(\delta J) * \text{flip}(K)$

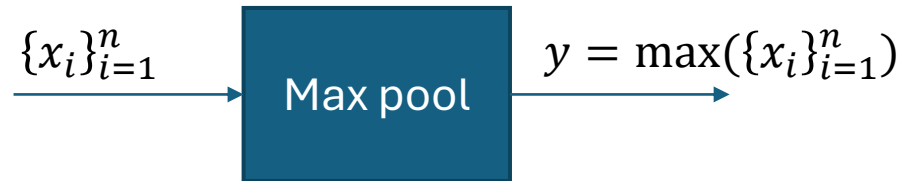
“pad” function adds $(h-1)$ 0 rows at the top and bottom and also adds $(w-1)$ 0 columns at the left and at the right of a matrix.

Size of $\text{pad}(\delta J)$ is $(H+h-1)$ -by- $(W+w-1)$.

$\text{flip}(K)$ is best understood by an example:

$$K = \begin{bmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \end{bmatrix} \quad \text{flip}(K) = \begin{bmatrix} 6 & 4 & 2 \\ 5 & 3 & 1 \end{bmatrix}$$

BP for a max pooling layer

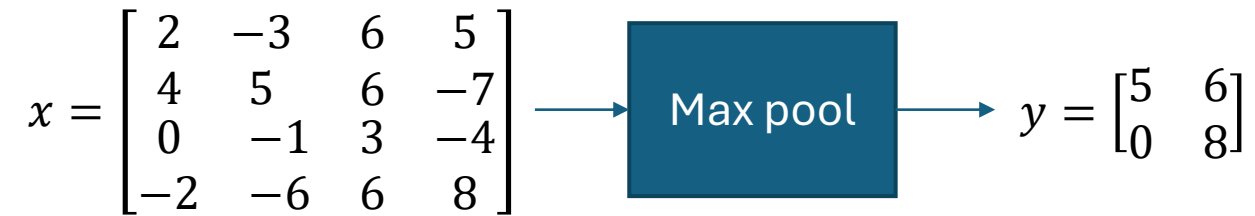


Note that in case of a tie, only a single index is chosen for the following operation:

$$i = \operatorname{argmax}_k \{x_k\}_{k=1}^n$$

By chain rule: $\delta x_i = \frac{\partial y}{\partial x_i} \delta y = \begin{cases} \delta y, & \text{if } i = \operatorname{argmax}_k \{x_k\}_{k=1}^n, \\ 0, & \text{otherwise.} \end{cases}$

Example of a 2-by-2, stride 2 max pooling:



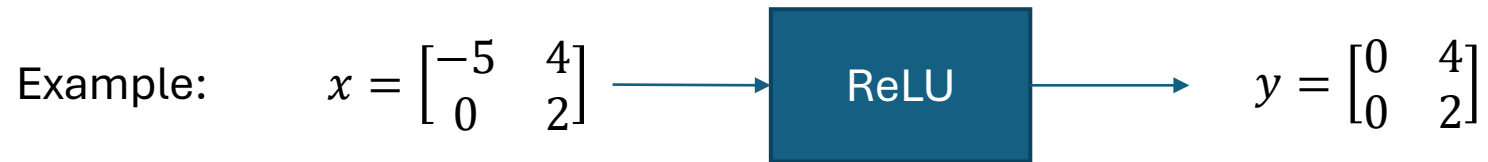
Suppose, $\delta y = \begin{bmatrix} \delta y_1 & \delta y_3 \\ \delta y_2 & \delta y_4 \end{bmatrix},$

then, $\delta x = \begin{bmatrix} 0 & 0 & \delta y_3 & 0 \\ 0 & \delta y_1 & 0 & 0 \\ \delta y_2 & 0 & 0 & 0 \\ 0 & 0 & 0 & \delta y_4 \end{bmatrix}.$

BP for a ReLU layer

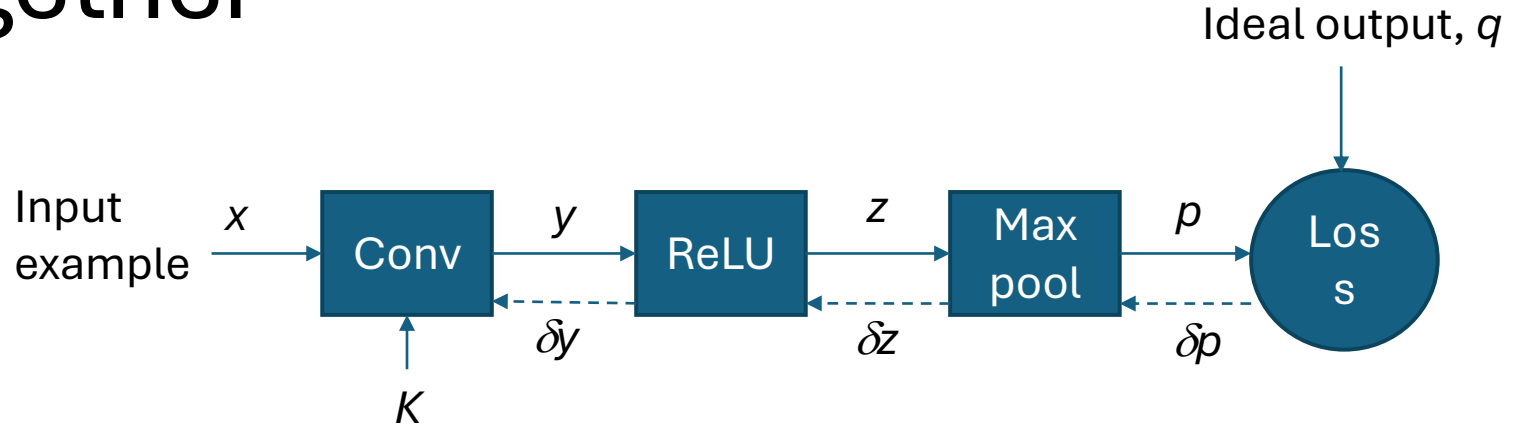


By chain rule:
$$\delta x = \frac{\partial y}{\partial x} \delta y = \begin{cases} \delta y, & \text{if } x > 0, \\ 0, & \text{otherwise.} \end{cases}$$



Suppose, $\delta y = \begin{bmatrix} \delta y_1 & \delta y_3 \\ \delta y_2 & \delta y_4 \end{bmatrix}$, then $\delta x = \begin{bmatrix} 0 & \delta y_3 \\ 0 & \delta y_4 \end{bmatrix}$.

Putting it all together



Convnet training algorithm:

An example network with 3 layers and a loss function

Initialize parameter K .

Iterate:

Step 1: (Forward pass)

Step 1a: Randomly choose a training example x and its corresponding ideal output q .

Step 1b: Pass x through “Conv” to get y ; pass y through ReLU to get z ; pass z through Max pool to get p .

Step 2: Compute “Loss” function for diagnostic purposes. /* Loss function measures deviation of p from q . */

Step 3: (Backward pass aka backpropagation)

Step 3a: Compute gradient of Loss function with respect to p . Denote this gradient by δp .

Step 3b: Compute δz given δp . /* Look at “BP for Max pooling.” */

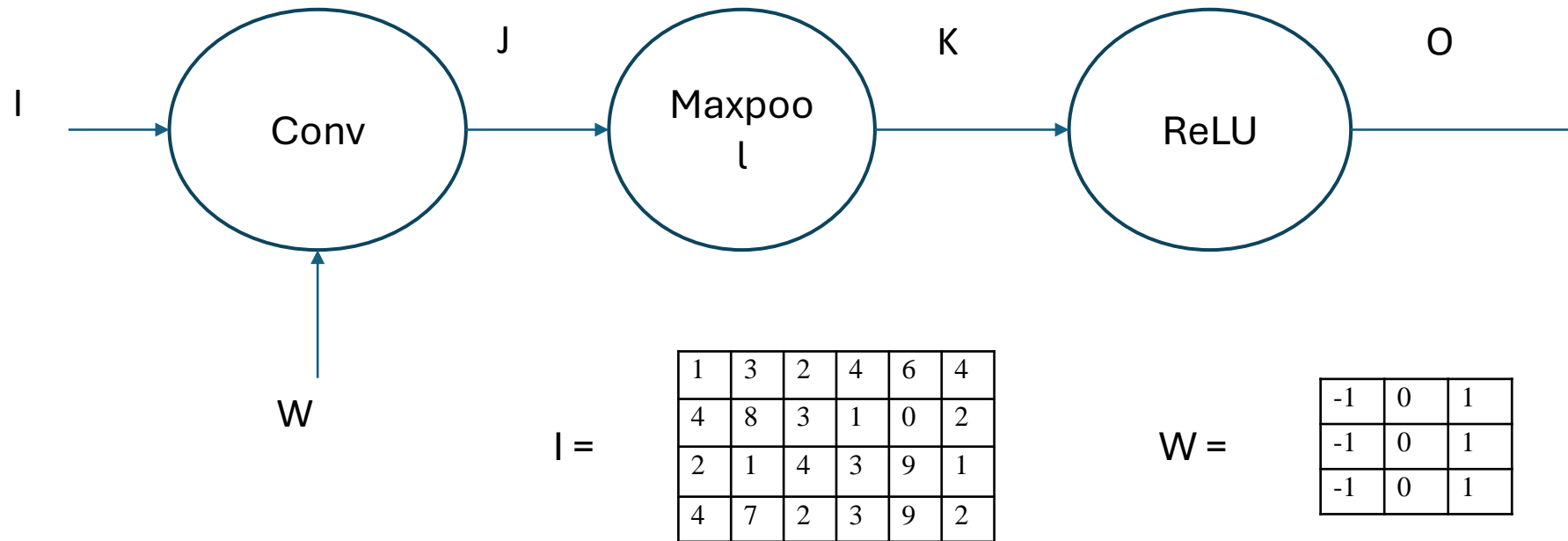
Step 3c: Compute δy given δz . /* Look at “BP for ReLU.” */

Step 4c: Compute δK given δy . /* Look at “BP for Conv.” */

Step 4: (Update parameter K by gradient descent) $K = K - (\text{learning rate}) \delta K$.

Note: We don’t have to compute δx , because there is no layer preceding conv layer in the example above.

Example w/PyTorch



Compute J , K and O . You are not doing any zero padding while doing the convolution. Assume stride size 1 for the convolution. Assume a 2-by-2 max pooling with stride 2. Write J , K and O below.

Now Assume that ideal output $IO = [7, 10]$. Also assume the loss as 0.5 times the square of Euclidean distance between IO and O . Compute back-propagation for O , K , J and W . Verify using pytorch autograd!

Example 2

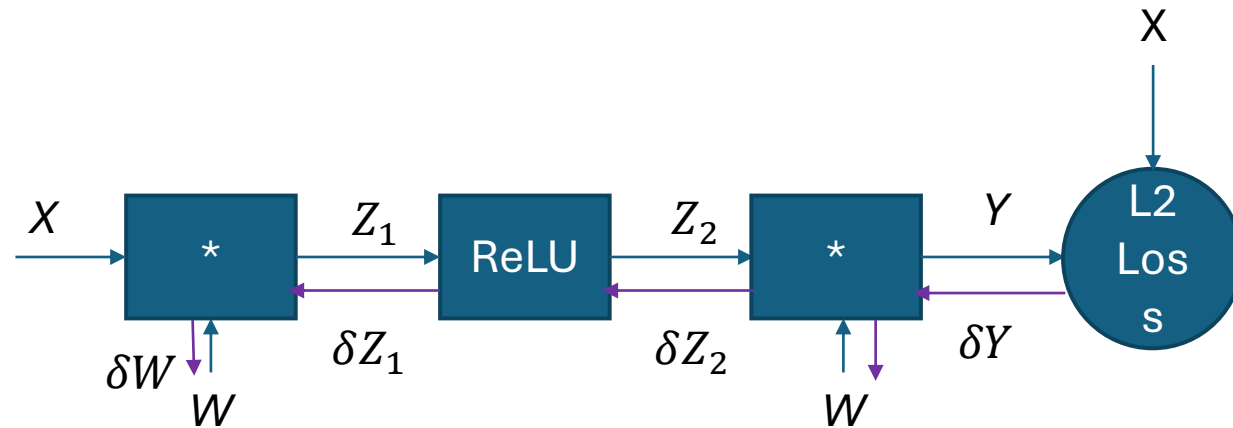
Consider a fully connected neural net:

$$Y = \text{ReLU}(X*W)*W$$

The loss is L2 between Y and X. Compute the gradient of the loss with respect to W. Here '*' refers to matrix multiplication.

Note that the same parameter matrix W appears in two computational nodes – called shared parameters/weights.

Example 2 *solution*



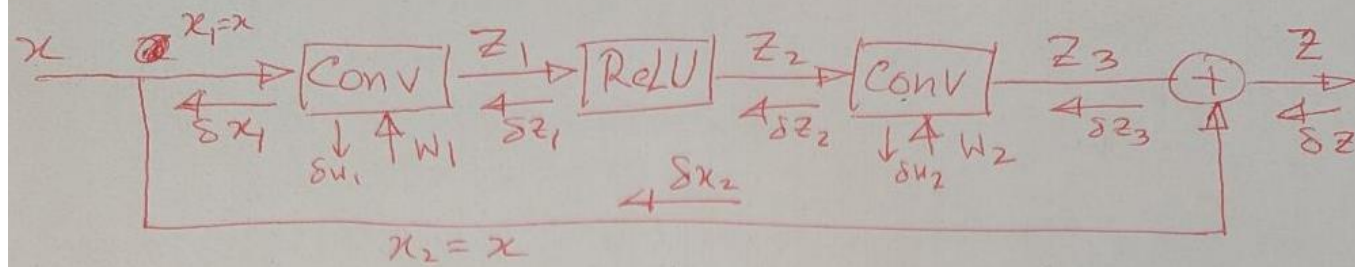
$$\delta Y = Y - X$$

$$\delta Z_2 = \delta Y * W^T$$

$$\delta Z_1 = ReLU'(Z_1) * \delta Z_2$$

$$\delta W = Z_2^T * \delta Y + X^T * \delta Z_1$$

BP in residual connection



computational graph for $z = \text{Conv}(\text{ReLU}(\text{Conv}(x, w_1)), w_2) + x$

$$z = \text{Conv}(\text{ReLU}(\text{Conv}(x, w_1)), w_2) + x$$

$$\delta x_2 = \delta z$$

$$\delta z_3 = \delta z$$

$$\delta z_2 = \text{Conv}(\text{Pad}(\delta z_3), \text{flip}(w_2))$$

$$\delta z_1 = \text{ReLU}'(z_1) \delta z_2$$

$$\delta x_1 = \text{Conv}(\text{Pad}(\delta z_1), \text{flip}(w_1))$$

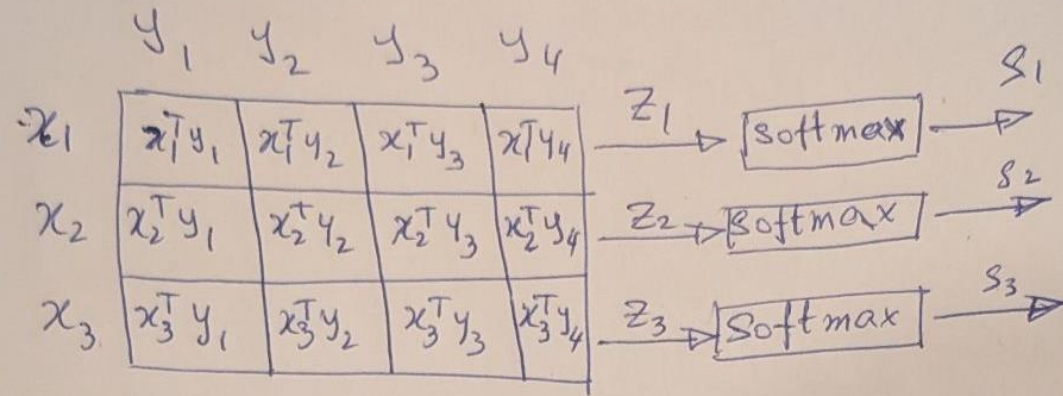
$$\delta x = \frac{\partial(\text{loss})}{\partial x} = \frac{\partial(\text{loss})}{\partial x_1} \frac{\partial x_1}{\partial x} + \frac{\partial(\text{loss})}{\partial x_2} \frac{\partial x_2}{\partial x}$$

$$= \delta x_1 + \delta x_2 \quad (\text{Because } x_1 = x \text{ and } x_2 = x)$$

$$\delta w_1 = \text{Conv}(x_1, \delta z_1)$$

$$\delta w_2 = \text{Conv}(z_2, \delta z_3)$$

BP in transformer



Given $\delta s_i, i=1,2,3$, compute $\delta x_i, i=1,2,3$ and $\delta y_j, j=1,2,3,4$.

Solution $\boxed{\delta z_i = J^i \delta s_i}, i=1,2,3.$

where J^i is the 4×4 Jacobian matrix for z_i :

$$(J^i)_{k,l} = \begin{cases} (s_i)_k (1 - (s_i)_k) & \text{if } k=l \\ -(s_i)_k (s_i)_l & \text{if } k \neq l. \end{cases}$$

$$\boxed{\delta x_i = \sum_{j=1}^4 y_j (\delta z_i)_j}, i=1,2,3$$

$$\boxed{\delta y_j = \sum_{i=1}^3 x_i (\delta z_i)_j}, j=1,2,3,4$$