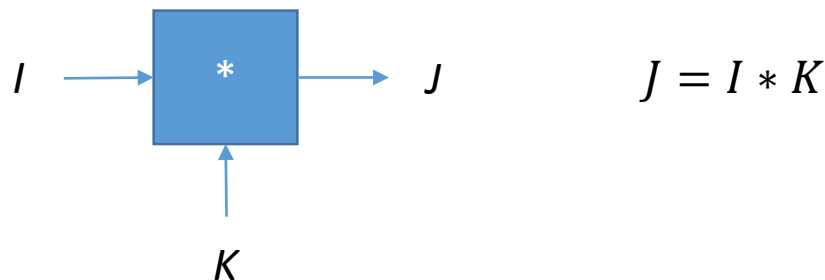


Training Convnet with Backprop

Nilanjan Ray

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 \mathcal{J} 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, $\delta K = I * \delta J$

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

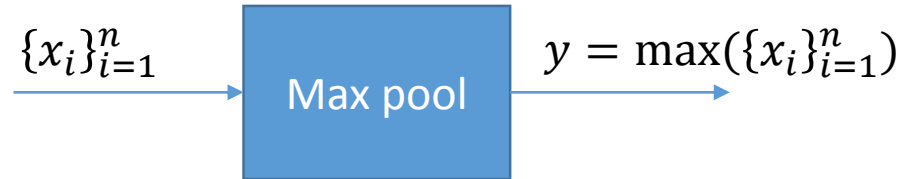
“pad” function adds $(h-1)$ 0 rows at the top and bottom and also adds $(w-1)$ 0 columns at the left and at the right of a matrix.

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

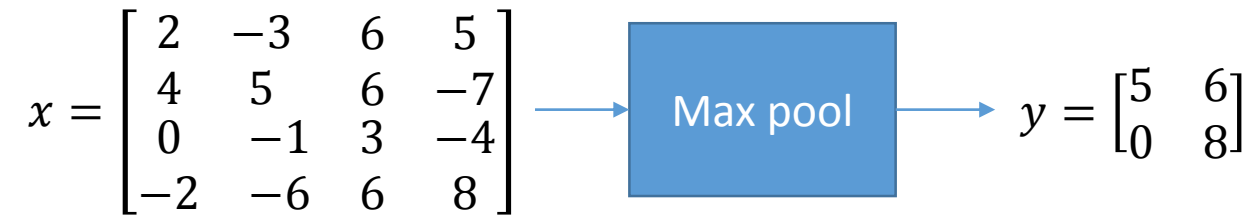


Note that in case of a tie, only a single index is chosen for the following operation:

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

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

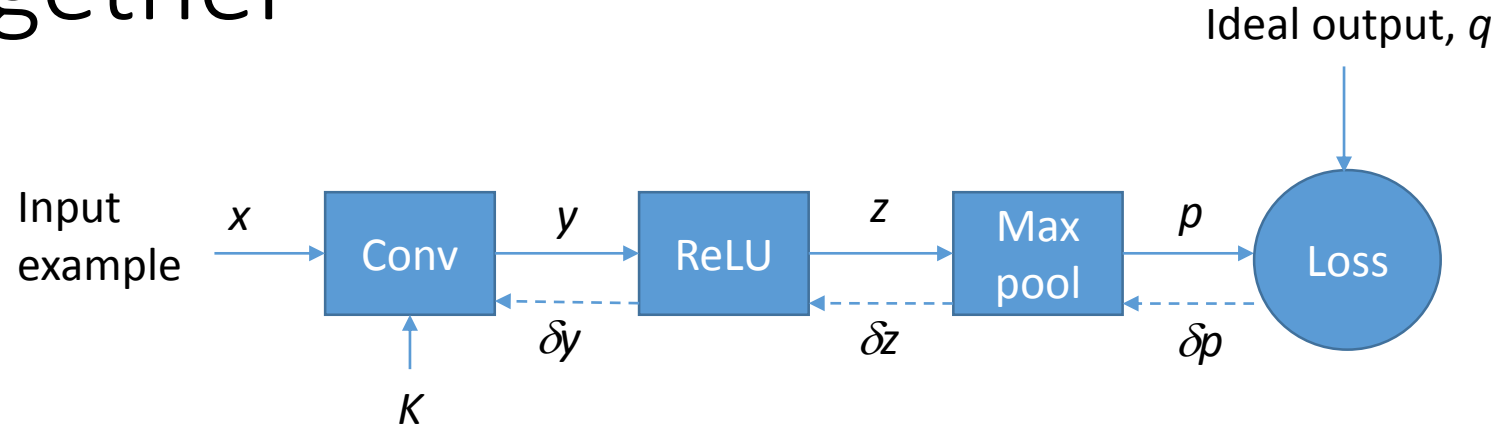


By chain rule:
$$\delta x = \frac{\partial y}{\partial x} \delta y = \begin{cases} \delta y, & \text{if } x > 0, \\ 0, & \text{otherwise.} \end{cases}$$

Example: $x = \begin{bmatrix} -5 & 4 \\ 0 & 2 \end{bmatrix}$  $y = \begin{bmatrix} 0 & 4 \\ 0 & 2 \end{bmatrix}$

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:

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 preceding layer before conv in the example above.