

Computer Vision

Montag, 20. Juli 2020 09:50

Convolution - filter

-- elementwise multiplication

to create filter in: Python -- conv_forward; tensorflow -- tf.nn.conv2d; Keras -- conv2D

Edge detection:

vertical

$$\begin{matrix} 6 \times 6 \\ \text{grid} \end{matrix} * \begin{matrix} 3 \times 3 \\ \begin{bmatrix} 1 & 0 & -1 \\ -1 & 0 & 1 \end{bmatrix} \end{matrix} = \begin{matrix} 4 \times 4 \\ \text{output grid} \end{matrix} \rightarrow \text{vertical edges}$$

horizontal

$$\begin{matrix} 6 \times 6 \\ \text{grid} \end{matrix} * \begin{matrix} 3 \times 3 \\ \begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix} \end{matrix} = \begin{matrix} 4 \times 4 \\ \text{output grid} \end{matrix} \rightarrow \text{horizontal edges}$$

Padding

$$\begin{matrix} n \times n \\ \text{grid} \end{matrix} \xrightarrow{p} \begin{matrix} (n+2p) \times (n+2p) \\ \text{padded grid} \end{matrix} * \begin{matrix} 3 \times 3 \\ \text{filter} \end{matrix} = \begin{matrix} (n+2p) \times (n+2p) \\ \text{output grid} \end{matrix}$$

the same dimension

valid convolution -> no padding

same convolution -> $p = (f-1)/2$ -> output size == input size

Stride

$$\begin{matrix} 5 \times 5 \\ \text{grid} \end{matrix} * \begin{matrix} 3 \times 3 \\ \text{filter} \end{matrix} = \begin{matrix} 2 \times 2 \\ \text{output grid} \end{matrix}$$

stride = 2

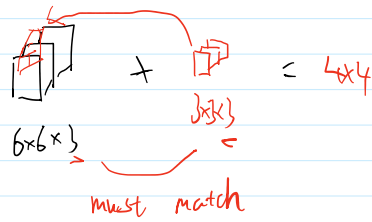
$$\left(\frac{n+2p-f}{s} + 1 \right)$$

if not integer:
 $\text{floor}((n+2p-f)/s) + 1$

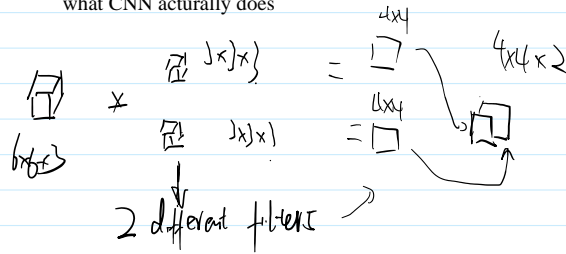
Convolution on RGB image -- convolution on each channel

what CNN actually does

what



what CNN actually does



$$\underbrace{n \times n \times n_c}_{\# \text{ channels}} * \underbrace{f \times f \times n_c}_{\# \text{ filters}} = (n - f + 1) \times (n - f + 1) \times n_c$$

for convolutional layer l :

$f^{[l]}$ - filter size

$p^{[l]}$ - padding

$s^{[l]}$ - stride

$n_c^{[l]}$ - # filters

$$Z_{\text{input}}^{[l-1]}: n_H^{[l-1]} \times n_W^{[l-1]} \times n_c^{[l-1]}$$

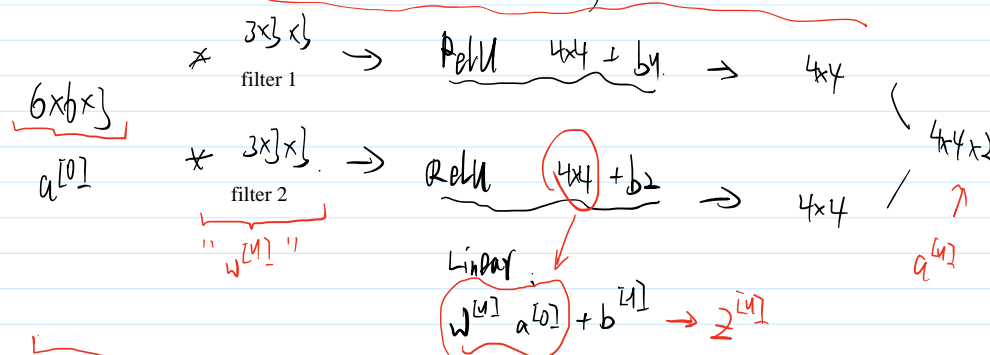
$$Z_{\text{output}}^{[l]}: n_H^{[l]} \times n_W^{[l]} \times n_c^{[l]}$$

$$n_H^{[l]} = \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1$$

$$a^{[l]} \text{ activation } A^{[l]} = m * a^{[l]}$$

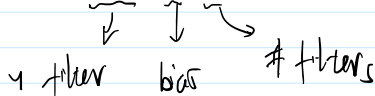
$$\text{weight } \underbrace{f^{[l]} \times f^{[l]} \times n_c^{[l-1]}}_{\text{dimension of one filter}} \times n_c^{[l]}$$

$$\text{bias } n_c^{[l]} \rightarrow (y_1, y_2, \dots, y_{n_c^{[l]}})$$

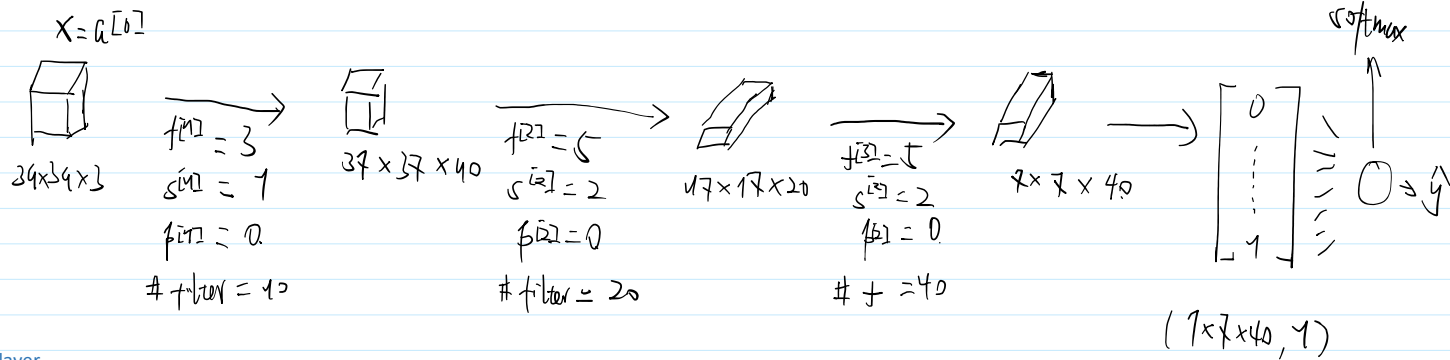


one layer, from $a^{[0]} \rightarrow a^{[1]}$ with 2 filters

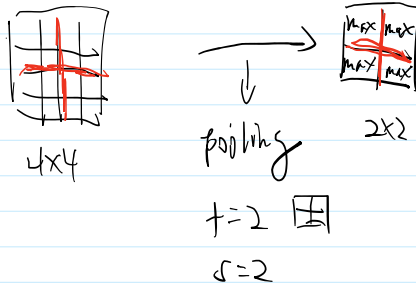
In this layer there are $(3*3*3 + 1) * 2 = 56$ parameters to learning.



Example:



Pooling layer



Intuition: reduce the dimension by keeping only the strongest features

Average pooling layer: instead of taking max, taking average.

Parameters to be learned:

Why using Conv net instead of conventional deep learning network? \rightarrow input size (features) is huge with respect to image/computer vision problem. To get same output size for each layer (how many features can be computed), conv net need to learn less parameters. Because, as to image convolution, actually a lot of pixels (information) are used more than once, so using conv net a lot of parameters(learned) are shared. So it is faster.

$$f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]} + \underbrace{n_c^{[l]}}_{\text{bias}}$$

Take away from Programming Assignment:

Build from scratch with numpy:

- Input shape = (m, n_H_prev, n_W_prev, n_C_prev)
- Output Z shape = (m, n_H, n_W, n_C)
- $n_H = \text{int}((n_H_prev + 2 \times \text{pad} - f) / \text{stride} + 1)$, same applied to n_W
- Weights shape (W) = (f, f, n_C_prev, n_C)
- biases (b) shape = (1, 1, 1, n_C)
- $A[i, h, w, c] = \text{activation}(Z[i, h, w, c])$
- Pooling layer: reduces the height and width of the input. It helps reduce computation, as well as helps make feature detectors more invariant to its position in the input.

Build with TensorFlow:

- `tf.nn.conv2d(X, W, strides = [1, s, s, 1], padding = 'SAME')`: given an input `XX` and a group of filters `WW`, this function convolves `WW`'s filters on `X`. The third parameter `([1, s, s, 1])` represents the strides for each dimension of the input (m, n_H_prev, n_W_prev, n_C_prev).
- `tf.nn.max_pool(A, ksize = [1, f, f, 1], strides = [1, s, s, 1], padding = 'SAME')`: given an input `A`, this function uses a window of size (f, f) and strides of size (s, s) to carry out max pooling over each window.
- `tf.nn.relu(Z)`: computes the elementwise ReLU of `Z`
- `tf.contrib.layers.flatten(P)`: given a tensor "P", this function takes each training (or test) example in the batch and flattens it into a 1D vector. (batch_size, k), where $k = h \times w \times c = h \times w \times c$.
- `tf.contrib.layers.fully_connected(F, num_outputs)`: given the flattened input `F`, it returns the output computed using a fully connected layer, automatically initializing the weights on this layer.
- `forward_propagation` function: CONV2D \rightarrow RELU \rightarrow MAXPOOL \rightarrow CONV2D \rightarrow RELU \rightarrow MAXPOOL \rightarrow \rightarrow FLATTEN \rightarrow FULLYCONNECTED.
- `tf.nn.softmax_cross_entropy_with_logits(logits = Z, labels = Y)`: computes the softmax entropy loss. This function both computes the softmax activation function as well as the resulting loss. Note that, the input of softmax is not `A` is `Z`.
- `tf.reduce_mean`: computes the mean of elements across dimensions of a tensor. Use this to calculate the sum of the losses over all the examples to get the overall cost.