

Q1. a) Denote filter size = $k \times k$, stride = S , current receptive field size = $r \times r$, distance (stride) between 2 adjacent features = j .

We have the receptive field size of the output feature map is equal to the area covered by k input features $(k-1) \times j_{in}$, plus the extra area covered by the receptive field of the input feature on the border.

Also, the stride of output receptive field is equal to the stride in the input map times the number of input features that got "jumped over" when applying the convolution.

Thus, we have

$$\begin{cases} j_{out} = j_{in} \times S \\ r_{out} = r_{in} + (k-1)j_{in} \end{cases}$$

For the input image (the actual input to the network), $r=1, j=1$

b) Pooling operation acts like a convolution given a different stride S' and filter size k' . The above expression still holds.

$$j'_{out} = j_{out} \times S', \quad r'_{out} = r_{out} + (k'-1)j_{out}$$

c) VGG16 architecture: Input $\rightarrow 2$ conv-64 \rightarrow maxpool $\rightarrow 2$ conv-128 \rightarrow maxpool $\rightarrow 3$ conv-256 \rightarrow maxpool $\rightarrow 3$ conv-512 \rightarrow maxpool $\rightarrow 3$ conv-512 \rightarrow maxpool \rightarrow FC-4096 \rightarrow FC-4096 \rightarrow FC-1000 \rightarrow softmax. Also, $k=3, S=1, k'=2, S'=2$

Applying the derived expression we have:

$j_0=1, r_0=1$ (input) $\rightarrow j_1=j_0 S=1, r_1=r_0+(k-1)j_0=1+2=3; j_2=j_1 S=1, r_2=r_1+(k-1)j_1=3+2=5$ (conv-64)
 $\rightarrow j_3=j_2 S'=2, r_3=r_2+(k'-1)j_2=5+1=6$ (maxpool) $\rightarrow j_4=j_3 S=2, r_4=r_3+(k-1)j_3=6+2=8; j_5=2, r_5=14$ (conv-128)
 $\rightarrow j_6=4, r_6=16$ (maxpool 2) $\rightarrow j_7=4, r_7=24; j_8=4, r_8=32; j_9=4, r_9=40$ (conv-256)
 $\rightarrow j_{10}=8, r_{10}=44$ (maxpool 3) $\rightarrow j_{11}=8, r_{11}=60; j_{12}=8, r_{12}=76$ (conv-512)
 $\rightarrow j_{13}=8, r_{13}=92$ (conv-512) $\rightarrow j_{14}=16, r_{14}=100$ (maxpool 4) $\rightarrow j_{15}=16, r_{15}=132$ (conv-512)
 $\rightarrow j_{16}=16, r_{16}=164; j_{17}=16, r_{17}=196$ (conv-512) $\rightarrow j_{18}=32, r_{18}=212$ (maxpool 5) \rightarrow FC.
 So the receptive field of VGG16 right before the first fully connected layer has a size of 212×212 , and stride 32.