

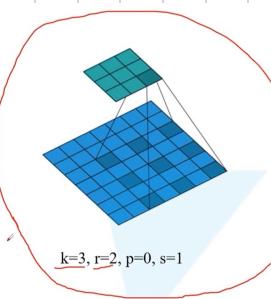
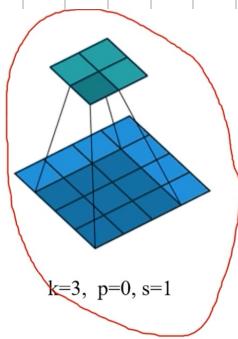
Dilated Convolution

240906

(1) 也叫法 Atrous convolution

(2) function { (1) 增大感受野

} (2) 保持原输入特征图 W, H



✓ 什么是 dilated convolution?

○ How dilated Convolution?
(怎么使用)

(3) Why Dilated convolution? 为什么使用膨胀卷积?

① 在语义分割中，如FCN 通常使用分类网络作为网络的 backbone
使用 backbone 对图片进行一系列的下采样，再使用一系列的上采样
还原原图的大小

② 通常使用的分类网络将高度 & 宽度下采样 32倍，后面再使用上采样
还原原图尺寸，而下采样倍率太大 我们再还原为原因有很大影响
如 VGG 由 maxpooling Layer 池化、下采样 首先降低图片高 & 宽度
其次丢失细节信息 以及比较小的目标 这些丢失的细节及目标无法
通过上采样还原 这会导致语义分割效果不理想

③ 那去掉 maxpooling 呢？ 确实不会降低高 & 宽度了。但这样的话
特征图对应的感受野却变小了

④ So dilated convolution

(4) 那么如此，那全用 dilated convolution?

Ans = No

《Understanding convolution for Semantic Segmentation》

gridding effect

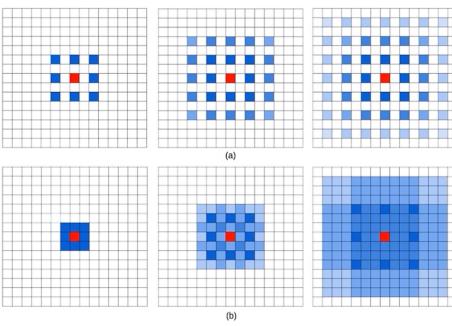
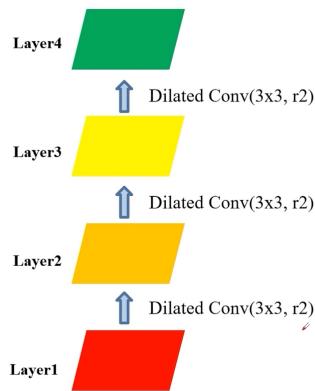


Figure 2. Illustration of the gridding problem. Left to right: the pixels (marked in blue) contributes to the calculation of the center pixel (marked in red) through three convolution layers with kernel size 3×3 . (a) all convolutional layers have a dilation rate $r = 2$. (b) subsequent convolutional layers have dilation rates of $r = 1, 2, 3$, respectively.

Experiment 1

(1) 什么是 gridding effect?

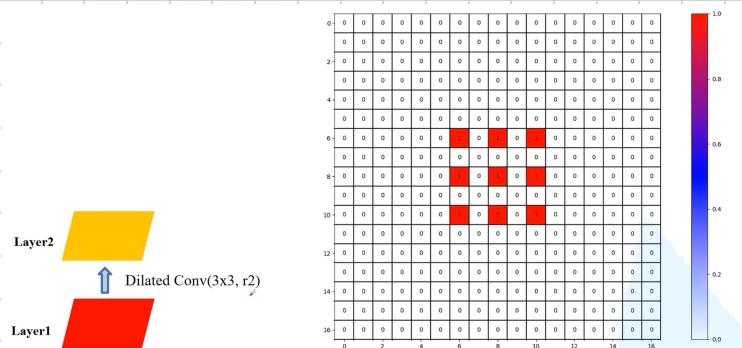


3个 3×3 的dilated convolution

$r=2$

$r=1$ ordinary convolution

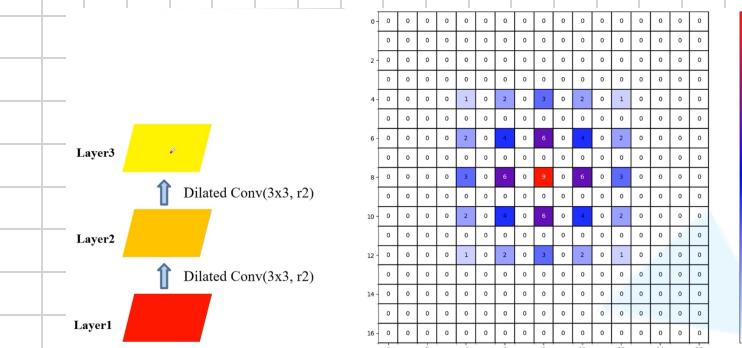
(1) 首先看 Layer1 \rightarrow Layer2 用到了 Layer1 的哪些元素



每2行 & 2列之间都间隔了1行 / 1列 0

用到了 layer1 中 9 个 pixel 信息

(2) 当我们连续使用2个dilated convolution layer3 使用了 layer1 哪些元素?



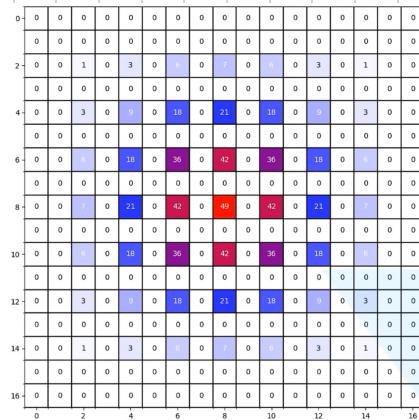
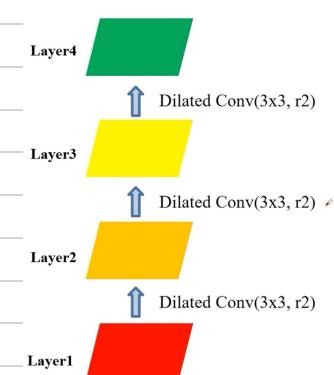
用到了 layer1 中 25 个 pixel 信息

颜色越红 次数越多

colorbar

(3) When we stack 3 dilated convolutions

layer4 uses pixels from layer1



$$RF = 13 \times 13$$

$$RF = 13 \times 13$$

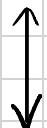
But many pixels are not used

Compared to Experiment 1

Highly regional usage of low-level receptive fields

(1) First, layer4 uses pixels from layer1, but they are not continuous.

Each non-zero element has a certain interval



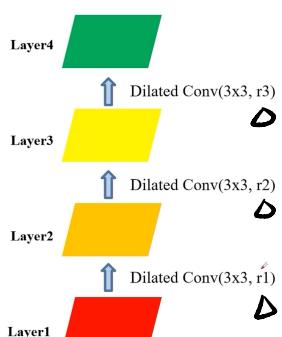
gridding effect

i.e., layer4 did not use all pixels in layer1's range. It only used part of them, leading to loss of information.

So we should avoid such situations.

Now, let's look at another experiment.

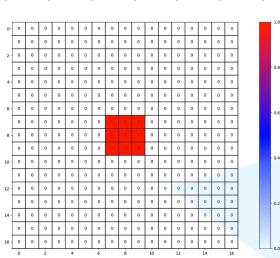
Experiment 2



● Observe dilation rates $r=1, 2, 3$

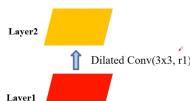
$r=1$ ordinary convolution

(1) Layer2 used which pixels from layer1?

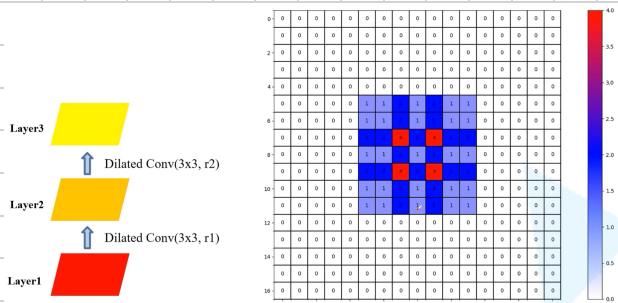


(1) Used 9 pixels

(2) Continuous



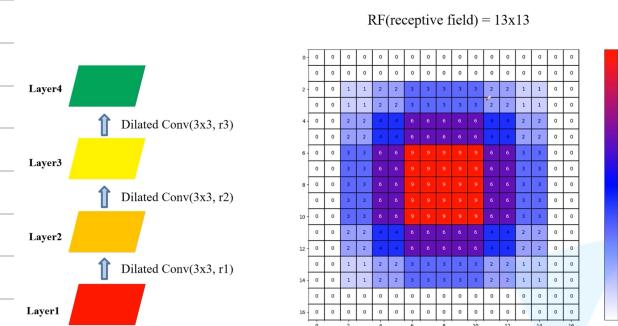
(2) Layer3 上的 pixel 对应 Layer1 上的哪些 pixel?



(1) 1个 7×7 的 region

(2) 相连

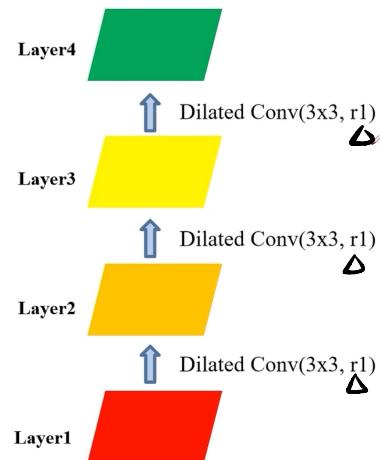
(3) Layer4 上的 pixel 对应 Layer1 上的哪些 pixel?



(1) RF = 13×13

(2) 连续 没有间隙

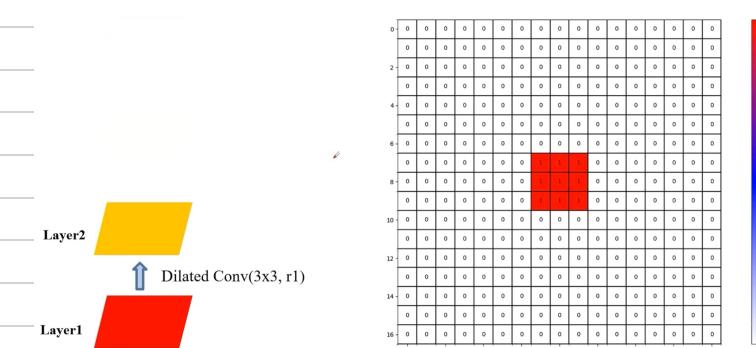
Experiment 3



当我们全部 $r=1$

连续使用3个 3×3 的 ordinary conv

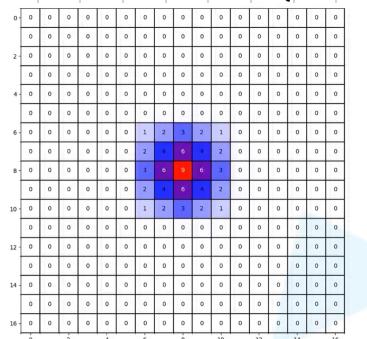
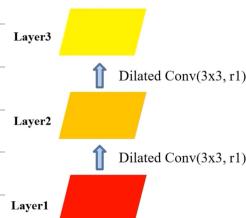
(1) Layer2 使用 Layer1 哪些 pixel , RF=?



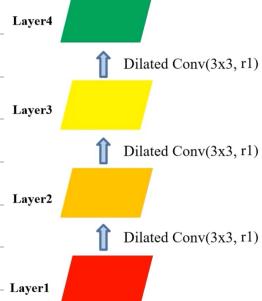
RF

= 3

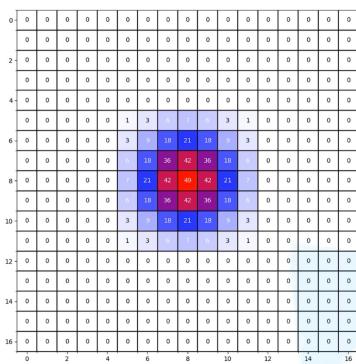
(2) Layer3 使用 Layer1



(3) Layer4 使用 Layer1



RF(receptive field) = 7x7



Layer4 RF = 7x7

首先 注意到 3组实验 全部是 3x3 的 conv \Rightarrow 参数数量是一样的

case1: $r=1$ 3个3x3 conv \rightarrow RF = 7x7

且连续 consecutive

case2: $r=1, 2, 3 \rightarrow RF = 13 \times 13$

summary: Hybrid Dilated Convolution

感受野可以变大且不损失信息

ref <<Understanding convolution for semantic segmentation>>

讨论如何设置 dilated rate

(6) About HDC & dilated rate (如何设计dilated rate)

Understanding Convolution for Semantic Segmentation

Here we propose a simple solution- hybrid dilated convolution (HDC), to address this theoretical issue. Suppose we have N convolutional layers with kernel size $K \times K$ that have dilation rates of $[r_1, \dots, r_i, \dots, r_n]$, the goal of HDC is to let the final size of the RF of a series of convolutional operations fully covers a square region without any holes or missing edges. We define the "maximum distance between two nonzero values" as

$$M_i = \max[M_{i+1} - 2r_i, M_{i+1} - 2(M_{i+1} - r_i), r_i], \quad (2)$$

with $M_n = r_n$. The design goal is to let $M_2 \leq K$. For example, for kernel size $K = 3$, an $r = [1, 2, 5]$ pattern works as $M_2 = 2$; however, an $r = [1, 2, 9]$ pattern does not work as $M_2 = 5$. Practically, instead of using the same dilation rate for all layers after the downsampling occurs, we

Hybrid Dilated Convolution (HDC)

$$M_2 \ll K$$

HDC的目的

HDC Target

(第一个 $r = 1$, $r = 1$ 为普通卷积)

① 壓清什麼是 maximum distance

0	0	0	0	0	0	0	0	0	0	0	0
0	0	1	3	6	7	6	3	1	0	0	0
0	0	3	9	18	21	18	9	3	0	0	0
0	0	6	18	36	42	36	18	6	0	0	0
0	0	7	21	42	49	42	21	7	0	0	0
0	0	6	18	36	42	36	18	6	0	0	0
0	0	3	9	18	21	18	9	3	0	0	0
0	0	1	3	6	7	6	3	1	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0

maximum distance = 1

0	0	0	0	0	0	0	0	0	0	0	0
0	1	0	3	0	0	7	0	8	0	3	0
0	0	0	0	0	0	0	0	0	0	0	0
0	3	0	0	0	18	0	21	0	18	0	3
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	1	0	3	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0

maximum distance = 2

间隔1行/1列的零

② 看什么

相反数

$$M_i = \max [M_{i+1} - 2k_i, M_{i+1} - 2(M_{i+1} - k_i), k_i]$$

meaning:

M_i : 第 i 层 两个非 0 元素 (pixel) 之间最大距离 (maximum distance)

r_i : i th 的 dilated rate

M_n : 最后一层 maximum distance = r_n n th dilated distance

Target:

$$M_2 \leq k$$

k : kernel-size

区分 M_i & r_i

问题: 那 T_n 是 M_i or r_i ?

For example (举例说明)

e.g. kernel-size = 3 $r = [1, 2, 5]$ $M_2 = 2$ $M_n = M_B = 5$

$$M_2 = 2 \leq k = \text{kernel-size} = 3$$

Now discuss M_2 从哪里来的?

$$\text{从 } M_i = \max [M_{i+1} - 2k_i, M_{i+1} - 2(M_{i+1} - k_i), k_i]$$

$$M_2 = \max [M_3 - 2r_2, M_3 - 2(M_3 - r_2), r_2]$$

$$\left| \begin{array}{l} M_3 = r_3 = 5 \\ r_2 = 2 \end{array} \right. \quad (M_n = r_n) \Rightarrow M_3 - 2r_2 = 5 - 4 = 1$$

$$M_3 - 2(M_3 - r_2) = 5 - 2(5 - 2) = -1$$

$$r_2 = 2$$

$$\therefore M_2 = 2$$

me: r 是我们设计的 设计出来的 r (决定 M)

st. $M_2 \leq k$ (先有 r 再有 M)

eq2. $T_n = \{1, 2, 9\}$

out: $M_2 = ?$

deal:

$$M_1 = \max[M_1 + 2r_1, M_1 - 2(M_1 - r_1), r_1]$$

$$M_2 = \max[M_2 - 2r_2, M_2 - 2(M_2 - r_2), r_2]$$

$$\begin{cases} M_3 = r_3 = 9 \\ r_2 = 2 \end{cases} \quad M_3 - 2r_2 = 9 - 4 = 5$$

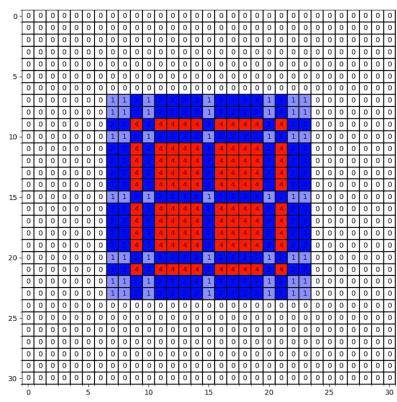
$$\begin{cases} M_3 = 9 \\ M_3 - r_2 = 9 - 2 = 7 \end{cases} \quad M_3 - 2(M_3 - r_2) = 9 - 14 = -5$$

$$\begin{cases} r_2 = 2 \end{cases}$$

$\therefore M_2 = 5 > k = 3$ | 不满足设计要求
这组参数不合适

eg1 & eg2 visualization

observe 高层 conv 用到底层 conv 哪些 pixel

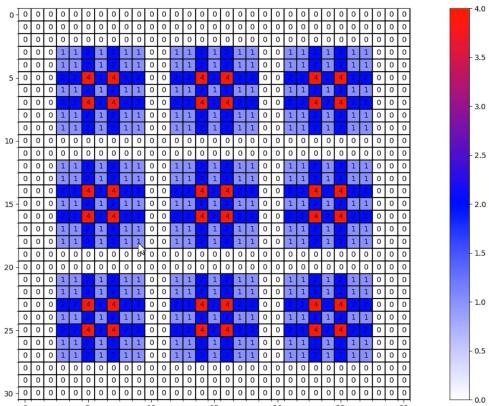


$$r = [1, 2, 5]$$

整个区域所有 pixel

consecutive

没有 gridding effect 问题



$$r = [1, 2, 9]$$

从图中可以看到 maximum distance = 3

(即 row / column 之间有 2 行 12 列 0)

由 eg2 $M_2 = 5$

求 M_1

$$M_1 = \max[M_2 - 2r_1, -(M_2 - 2r_1), r_1]$$

$$M_2 - 2r_1 = 5 - 2 \times 1 = 3, -3, 1$$

$\therefore M_1 = 3 = \text{maximum distance}$

(7) discuss [1, 2, 5] [1, 2, 9] 为什么都从1开始?

$r=1 \Leftrightarrow$ ordinary convolution

(8) principle₁ : $M \leq k$

principle₂ dilation rate 锯齿结构 重叠

Understanding Convolution for Semantic Segmentation

use a different dilation rate for each layer. In our network, the assignment of dilation rate follows a sawtooth wave-like heuristic: a number of layers are grouped together to form the “rising edge” of the wave that has an increasing dilation rate, and the next group repeats the same pattern. For example, for all layers that have dilation rate $r = 2$, we form 3 succeeding layers as a group, and change their dilation rates to be 1, 2, and 3, respectively. By doing this, the top layer can access information from a broader range of pixels, in the same region as the original configuration (Figure 2(b)). This process is repeated through all layers, thus making the receptive field unchanged at the top layer.

Hybrid Dilated Convolution (HDC)

锯齿形状

将 dilation rates 设置成锯齿结构, 例如:
[1, 2, 3, 1, 2, 3]



{ 锯齿结构
公约数≤1

(9) principle₃ 公约数≤1

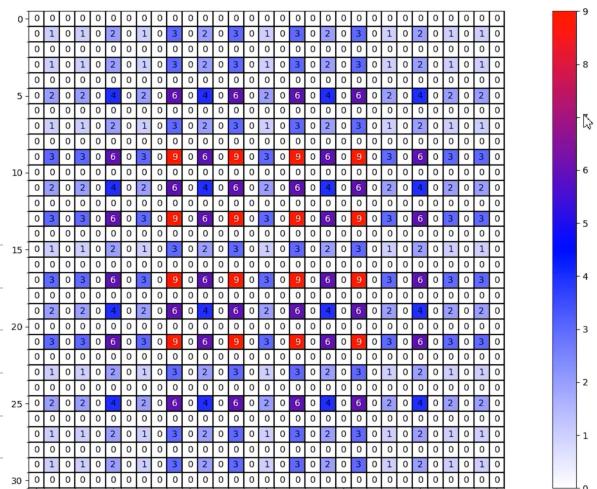
Understanding Convolution for Semantic Segmentation

Another benefit of HDC is that it can use arbitrary dilation rates through the process, thus naturally enlarging the receptive fields of the network without adding extra modules [29], which is important for recognizing objects that are relatively big. One important thing to note, however, is that the dilation rate within a group should not have a common factor relationship (like 2, 4, 8, etc.), otherwise the gridding problem will still hold for the top layer. This is a key difference between our HDC approach and the atrous spatial pyramid pooling (ASPP) module in [3], or the context aggregation module in [29], where dilation factors that have common factor relationships are used. In addition, HDC is naturally integrated with the original layers of the network, without any need to add extra modules as in [29, 3].

Hybrid Dilated Convolution (HDC)

公约数不能大于1

otherwise gridding problem



VISUALIZATION: $r = [2, 4, 8]$

Summary HDC 3个 principle

(1) $M_2 \leq K$
(2) 锯齿状
(3) 公约数 ≤ 1