- VGG (Visual Geometry Group)
  Very Deep Convolutional Networks For Large-Scale Image Recognition

- 亮点: ① 通过堆叠多个 3×3 的卷积核来替代大尺度卷积核 (减少所需参数)
  ② 论文中提到, 可以通过堆叠 两个 3×3 的卷积核替代 5×5 的卷积核,
     堆叠 三个 3×3 的卷积核替代 7×7 的卷积核
  ⇒ 拥有相同的感受野

Feature map: $F=1$
Conv3×3 (3): $F=(1-1)×1+3=3$
conv 3×3 (2): $F=(3-1)×1+3=5$
Conv 3×3 (1): $F=(5-1)×1+3=7$
#所需参数    <Input channel = C>
$7×7×C×C = 49C^2$
$3×3×C×C + 3×3×C×C + 3×3×C×C = 27C^2$

一般使用D

receptive field

决定某一层输出结果中 一个元素
所对应的输入层的区域大小.
<输出 feature map 上的一个单元
对应输入层上区域的大小>

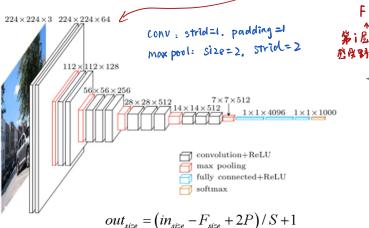| ConvNet Configuration | | | | | |
|---|---|---|---|---|---|
| A | A-LRN | B | C | D | E |
| 11 weight layers | 11 weight layers | 13 weight layers | 16 weight layers | 16 weight layers | 19 weight layers |
| input (224 × 224 RGB image) | | | | | |
| conv3-64 | conv3-64 LRN | conv3-64 conv3-64 | conv3-64 conv3-64 | conv3-64 conv3-64 | conv3-64 conv3-64 |
| maxpool | | | | | |
| conv3-128 | conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 |
| maxpool | | | | | |
| conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 conv1-256 | conv3-256 conv3-256 conv3-256 | conv3-256 conv3-256 conv3-256 conv3-256 |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 conv1-512 | conv3-512 conv3-512 conv3-512 | conv3-512 conv3-512 conv3-512 conv3-512 |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 conv1-512 | conv3-512 conv3-512 conv3-512 | conv3-512 conv3-512 conv3-512 conv3-512 |
| maxpool | | | | | |
| FC-4096 | | | | | |
| FC-4096 | | | | | |
| FC-1000 | | | | | |
| soft-max | | | | | |



输出
2×2×1    Max Pool
         Size: 2×2
         Stride: 2
4×4×1
         Conv1
         Size: 3×3
输入      Stride = 2
9×9×1
$out_{size} = (in_{size} - F_{size} + 2P)/S + 1$

$F(i) = (F(i+1) -1) × Stride + Ksize$
第i层          kernal size
感受野

eg. Feature map: $F=1$
    Pool1:  $F = (1-1) × 2 + 2 = 2$
    Conv1:  $F = (2-1) × 2 + 3 = 5$

224×224×3   224×224×64
conv: strid=1, padding=1
max pool: size=2, strid=2

112×112×128
56×56×256
28×28×512
14×14×512
7×7×512
1×1×4096   1×1×1000

convolution+ReLU
max pooling
fully connected+ReLU
softmax

$out_{size} = (in_{size} - F_{size} + 2P)/S + 1$