

DenseNet学习笔记

前言

之前的一些研究表明了在输入层和输出层之间添加一些跳接可以让网络架构更深，且训练更有效率。例如 [ResNet](#) [1]，解决了深层网络梯度消失的问题，而 [GoogleNet](#) [2] 则是让网络加宽。借鉴这两种思想，让网络中各层之间的信息传递，将**所有的层连接起来**，这就是 [DenseNet](#) [3] 的基本思想。

在传统的卷积神经网络中，第 L 层就有 L 个连接，每一层和其他的层相互连接，所以总共的跳接就有 $\frac{L(L+1)}{2}$ ，如 [Figure 1](#) 所示。对于每一层来说，所有此前的网络层的特征图作为输入，而其自身的特征图作为之后所有层的输入。[DenseNet](#) 有以下几个优点：

- 减轻了梯度消失问题 (vanishing-gradient)
- 加强了特征传播 (feature propagation)
- 更有效地利用特征
- 大大减少了参数数量

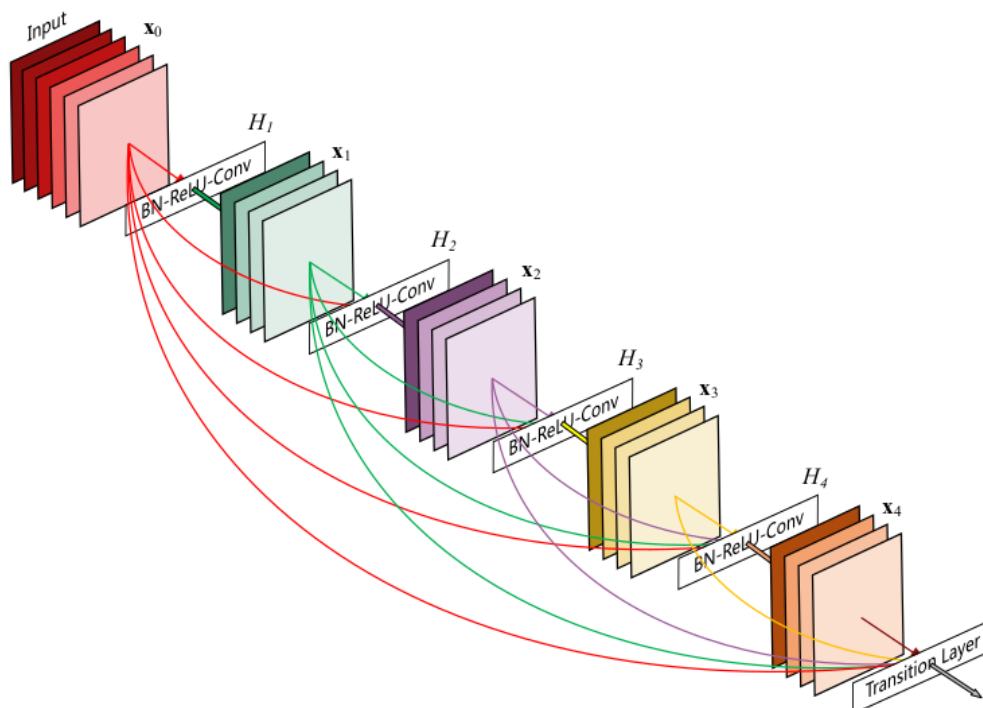


Figure 1: A 5-layer dense block with a growth rate of $k = 4$. Each layer takes all preceding feature-maps as input.

DenseNet 架构

假设 X_0 是输入卷积网络的单张图片，网络包括 L 层，每一层都实现了非线性变换 $H_l(\cdot)$ ，其中 l 表示的是第 l 层。 $H_l(\cdot)$ 是包含了批量归一化 (Batch Normalization, BN)、ReLU、池化和卷积的组合操作，将 l^{th} 层的输出命名为 X_l 。

ResNets

传统的卷积前馈网络将 l^{th} 的输出作为 $(l + 1)^{th}$ 层的输入，得到这个转换公式： $X_l = H_l(X_{l-1})$ 。而 **ResNet** 通过标识函数 (identity function) 添加了一个绕过非线性变换 $H_l(\cdot)$ 的跳接

$$X_l = H_l(x_{l-1}) + x_{l-1} \quad (1)$$

ResNet 的一个优点是梯度可以直接通过标识函数 (identity function) 从后面的层流向前面的层。但是，标识函数 (identity function) 和 H_l 层的输出通过求和进行组合，这可能会阻碍网络中信息的流动。

Dense 连接

为了进一步地层与层之间的信息流，**DenseNet** 提出了一个不同的连接模型：对于每一层，都添加一个跳接到其他所有之后的层。**Figure 1**表示了 **DenseNet** 连接的方式。因此， l^{th} 层网络接受了所有之前层的特征图 X_0, \dots, X_{l-1} 作为输入：

$$X_l = H_l([X_0, X_1, \dots, X_{l-1}]) \quad (2)$$

其中 $[X_0, X_1, \dots, X_{l-1}]$ 表示的是 $0, \dots, l - 1$ 层得到的特征图拼接的结果。

Composite function

$H_l(\cdot)$ 表示的是三个连续的操作：

- batch normalization (BN)
- rectified linear unit (ReLU)
- 3 x 3 Conv

池化层

当特征图尺寸变化时，**式2**中的拼接操作不可行。但是，卷积网络一个重要的部分就是降采样层，用于改变特征图的尺寸。为了在 **DenseNet** 架构中实现降采样，将网络分为多个紧密连接的 **dense blocks**，如**Figure 2**所示。

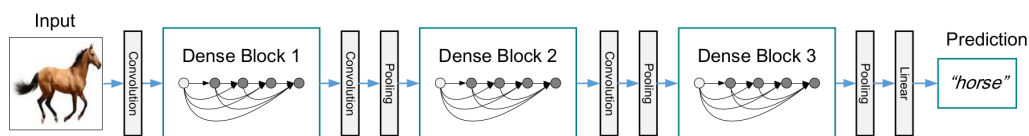


Figure 2: A deep DenseNet with three dense blocks. The layers between two adjacent blocks are referred to as transition layers and change feature-map sizes via convolution and pooling.

将 **dense block** 之间的层叫做过渡层，在这里做卷积和池化操作。过渡层包含批量归一层和 1×1 卷积层，紧跟一个 2×2 平均池化层

Growth rate

如果每个函数 H_l 产生 k 个特征图，之后的 l^{th} 层有 $k_0 + k \times (l - 1)$ 个输入特征图，其中 k_0 表示输入层的通道数。**DenseNet** 和现有的网络架构最重要的区别是 **DenseNet** 层数很窄，仅有 $k = 12$ 。将 k 定义为网络的增长率。

Bottleneck layers

尽管每一层都只产生 k 个输出特征图，仍然有许多输入。**ResNet** 中在 3×3 卷积前使用 1×1 卷积作为 **bottleneck** 层减少输入特征图的数量，可以提高计算效率。使用了 **Bottleneck** 的网络命名为 **DenseNet-B**。

Compression

为了进一步使模型更加紧凑，在过渡层减少特征图的数量。如果 **dense block** 包括 m 个特征图，让之后的过渡层产生 $[\theta_m]$ 输出特征图，其中 $0 < \theta \leq 1$ 表示压缩因子。如果 $\theta = 1$ ，表示特征图数量经过过渡层保持不变。在试验中设置 $\theta = 0.5$ 。将使用了 **bottleneck** 和过渡层设置 $\theta < 1$ 的网络命名为 **DenseNet-BC**

实现细节

在所有除了 **ImageNet** 的数据集中，实验使用的 **DenseNet** 有三个 **dense block**，每个块的层数相等。在第一个 **dense block** 之前，对输入图像进行一个带有16（或者是 **DenseNet-BC** 增长率两倍）个输出通道的卷积操作。对于卷积核大小为 3×3 的卷积层，输入的每一侧都用一个像素进行零填充以修正特征图尺寸。在两个连续的 **dense block** 之间使用一个 1×1 的卷积接着一个 2×2 的池化层组成的过渡层。在最后一个 **dense block**，使用一个全局平均池化层和一个 **softmax** 函数。在这三个 **dense block** 中的特征图分别为 32×32 、 16×16 和 8×8 。

基本的 **DenseNet** 架构使用了以下的参数配置：

- $L = 40, k=12$
- $L = 100, k=12$
- $L = 100, k=24$

对于 DenseNet-BC，使用了以下的参数：

- $L = 100, k=12$
- $L = 250, k=24$
- $L = 190, k=40$

在 ImageNet 数据集的实验中，使用了 DenseNet-BC 结构，输入图像尺寸为 224×224 ，dense block 有4个。初始的卷积层包含 $2k$ 个步长为2的 7×7 卷积；其他层的特征图数量遵循设置 k 。ImageNet 配置如Table 1 所示

Layers	Output Size	DenseNet-121	DenseNet-169	DenseNet-201	DenseNet-264
Convolution	112×112	7×7 conv, stride 2			
Pooling	56×56	3×3 max pool, stride 2			
Dense Block (1)	56×56	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 6$
Transition Layer (1)	56×56	1×1 conv			
	28×28	2×2 average pool, stride 2			
Dense Block (2)	28×28	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 12$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 12$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 12$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 12$
Transition Layer (2)	28×28	1×1 conv			
	14×14	2×2 average pool, stride 2			
Dense Block (3)	14×14	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 24$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 32$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 48$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 64$
Transition Layer (3)	14×14	1×1 conv			
	7×7	2×2 average pool, stride 2			
Dense Block (4)	7×7	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 16$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 32$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 32$	$\begin{bmatrix} 1 \times 1 \text{ conv} \\ 3 \times 3 \text{ conv} \end{bmatrix} \times 48$
Classification Layer	1×1	7×7 global average pool			
		1000D fully-connected, softmax			

Table 1: DenseNet architectures for ImageNet. The growth rate for all the networks is $k = 32$. Note that each “conv” layer shown in the table corresponds the sequence BN-ReLU-Conv.

实验

数据集

CIFAR

训练集-50,000张图片，测试集10,000张图片，从训练集中选 5,000 张图片作为验证集。

- 使用了标准的数据增强，镜像，平移等
- 预处理使用了标准化

SVHN

训练集 73,257张图片，测试集26,032图片，还有531,131张图片作为额外的训练，从训练集中挑选6,000张图片作为验证集

- 没有使用任何数据增强

ImageNet

训练集使用了1.2m张图片，50,000张图片作为验证

- 使用了标准的数据增强
- 在测试的使用应用了 **single-crop** 和 **10-crop**

训练

- 使用的SGD方法训练
- **CIFAR**
 - batch size 64
 - epoch 300
- **SVHN**
 - batch size 64
 - epoch 40
- 初始学习率设置为0.1，在50%和75%训练进度除以10
- **ImageNet**
 - epoch 90
 - batch size 256
 - lr 0.1, 在30和60 epoch除以10

结果

CIFAR和**SVHN**主要的结果如table 2所示

Method	Depth	Params	C10	C10+	C100	C100+	SVHN
Network in Network [22]	-	-	10.41	8.81	35.68	-	2.35
All-CNN [32]	-	-	9.08	7.25	-	33.71	-
Deeply Supervised Net [20]	-	-	9.69	7.97	-	34.57	1.92
Highway Network [34]	-	-	-	7.72	-	32.39	-
FractalNet [17]	21	38.6M	10.18	5.22	35.34	23.30	2.01
with Dropout/Drop-path	21	38.6M	7.33	4.60	28.20	23.73	1.87
ResNet [11]	110	1.7M	-	6.61	-	-	-
ResNet (reported by [13])	110	1.7M	13.63	6.41	44.74	27.22	2.01
ResNet with Stochastic Depth [13]	110	1.7M	11.66	5.23	37.80	24.58	1.75
	1202	10.2M	-	4.91	-	-	-
Wide ResNet [42]	16	11.0M	-	4.81	-	22.07	-
	28	36.5M	-	4.17	-	20.50	-
with Dropout	16	2.7M	-	-	-	-	1.64
ResNet (pre-activation) [12]	164	1.7M	11.26*	5.46	35.58*	24.33	-
	1001	10.2M	10.56*	4.62	33.47*	22.71	-
DenseNet ($k = 12$)	40	1.0M	7.00	5.24	27.55	24.42	1.79
DenseNet ($k = 12$)	100	7.0M	5.77	4.10	23.79	20.20	1.67
DenseNet ($k = 24$)	100	27.2M	5.83	3.74	23.42	19.25	1.59
DenseNet-BC ($k = 12$)	100	0.8M	5.92	4.51	24.15	22.27	1.76
DenseNet-BC ($k = 24$)	250	15.3M	5.19	3.62	19.64	17.60	1.74
DenseNet-BC ($k = 40$)	190	25.6M	-	3.46	-	17.18	-

Table 2: Error rates (%) on CIFAR and SVHN datasets. k denotes network's growth rate. Results that surpass all competing methods are **bold** and the overall best results are **blue**. "+" indicates standard data augmentation (translation and/or mirroring). * indicates results run by ourselves. All the results of DenseNets without data augmentation (C10, C100, SVHN) are obtained using Dropout. DenseNets achieve lower error rates while using fewer parameters than ResNet. Without data augmentation, DenseNet performs better by a large margin.

在**ImageNet**分类的结果和 **ResNet** 的对比如table 3和Figure 4所示。

Model	top-1	top-5
DenseNet-121	25.02 / 23.61	7.71 / 6.66
DenseNet-169	23.80 / 22.08	6.85 / 5.92
DenseNet-201	22.58 / 21.46	6.34 / 5.54
DenseNet-264	22.15 / 20.80	6.12 / 5.29

Table 3: The top-1 and top-5 error rates on the ImageNet validation set, with single-crop / 10-crop testing.

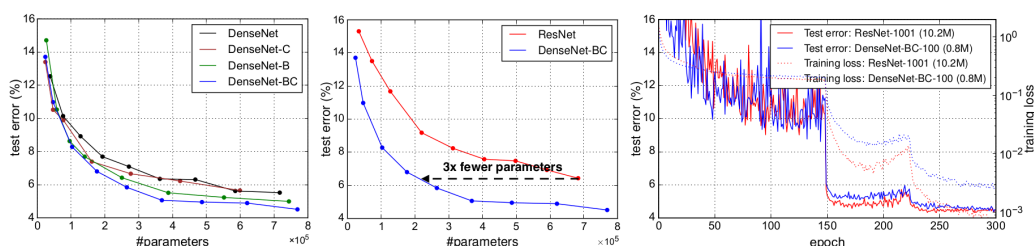


Figure 4: *Left:* Comparison of the parameter efficiency on C10+ between DenseNet variations. *Middle:* Comparison of the parameter efficiency between DenseNet-BC and (pre-activation) ResNets. DenseNet-BC requires about 1/3 of the parameters as ResNet to achieve comparable accuracy. *Right:* Training and testing curves of the 1001-layer pre-activation ResNet [12] with more than 10M parameters and a 100-layer DenseNet with only 0.8M parameters.

参考

1. He K , Zhang X , Ren S , et al. Deep Residual Learning for Image Recognition[J]. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016.
2. Szegedy C , Liu W , Jia Y , et al. Going Deeper with Convolutions[J]. IEEE Computer Society, 2014.
3. Huang G , Liu Z , Laurens V , et al. Densely Connected Convolutional Networks[J]. IEEE Computer Society, 2016.