

Lecture 5: Image Classification with CNNs

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 1

April 18, 2023

Image Classification: A core task in Computer Vision



(assume given a set of labels)
{dog, cat, truck, plane, ...}

→
cat
dog
bird
deer
truck

Pixel space



$$f(x) = Wx$$



Class scores

Image features



$$f(x) = Wx$$

Feature Representation

Class scores

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 16

April 18, 2023

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 17

April 18, 2023

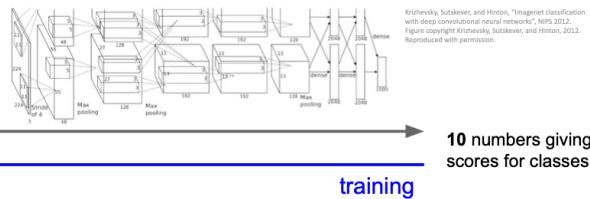
Image features vs. ConvNets



f

10 numbers giving scores for classes

training



10 numbers giving scores for classes

training

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 22

April 18, 2023

Last Time: Neural Networks

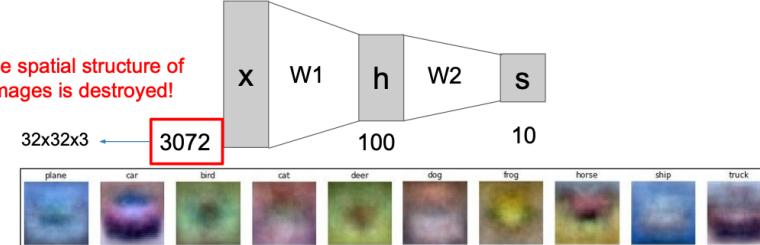
Linear score function:

$$f = Wx$$

2-layer Neural Network

$$f = W_2 \max(0, W_1 x)$$

The spatial structure of images is destroyed!



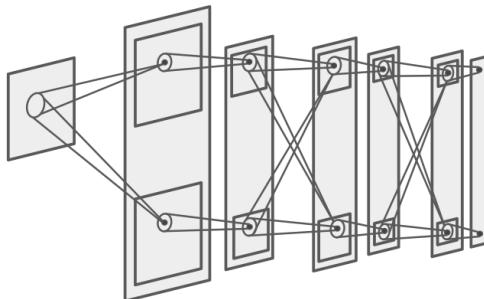
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 23

April 18, 2023

A bit of history:

Neocognitron
[Fukushima 1980]



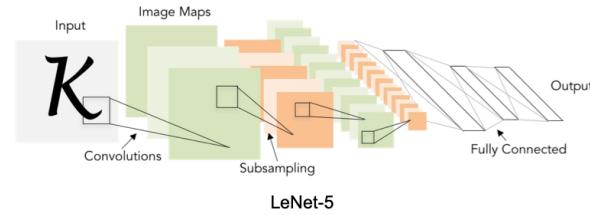
"sandwich" architecture (SCSCSC...)
simple cells: modifiable parameters
complex cells: perform pooling

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 33

April 18, 2023

A bit of history:
Gradient-based learning applied to document recognition
[LeCun, Bottou, Bengio, Haffner 1998]



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 34

April 18, 2023

A bit of history:

ImageNet Classification with Deep Convolutional Neural Networks
[Krizhevsky, Sutskever, Hinton, 2012]

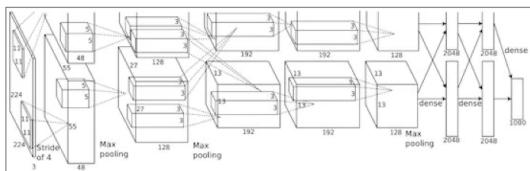


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

AlexNet

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 35

April 18, 2023

Fast-forward to today: ConvNets are everywhere

Classification



Figures copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Retrieval



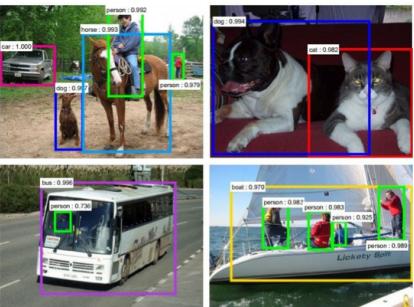
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 36

April 18, 2023

Fast-forward to today: ConvNets are everywhere

Detection



Figures copyright Shaqing Ren, Kaiming He, Ross Girshick, Jian Sun, 2015. Reproduced with permission.
[Faster R-CNN: Ren, He, Girshick, Sun 2015]

Segmentation



Figures copyright Clement Farabet, 2012.
Reproduced with permission.

[Farabet et al., 2012]

No errors



A white teddy bear sitting in the grass

Minor errors



A man in a baseball uniform throwing a ball

Somewhat related



A woman is holding a cat in her hand



A man riding a wave on top of a surfboard



A cat sitting on a suitcase on the floor



A woman standing on a beach holding a surfboard

Image Captioning

[Vinyals et al., 2015]
[Karpathy and Fei-Fei, 2015]

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 37

April 18, 2023

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 43

April 18, 2023

All images are CC0 Public domain:
<https://pixabay.com/en/teddy-bear-white-cute-toy/> (1643010)
<https://pixabay.com/en/baseball-player-sports-baseball/> (1667712)
<https://pixabay.com/en/cat-suitcase-animal-feline/> (1667967)
<https://pixabay.com/en/woman-meditation-lake-meditation/> (99008)
<https://pixabay.com/en/baseball-player-shortstop-field/> (1045263)

Captions generated by Justin Johnson using **NeuralTalk2**

Convolutional Neural Networks

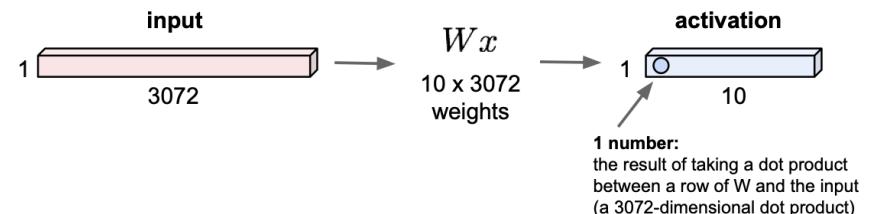
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 45

April 18, 2023

Fully Connected Layer

32x32x3 image \rightarrow stretch to 3072 x 1



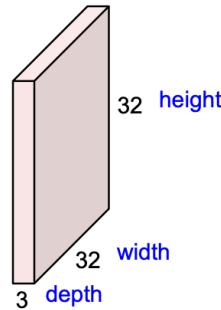
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 47

April 18, 2023

Convolution Layer

32x32x3 image \rightarrow preserve spatial structure



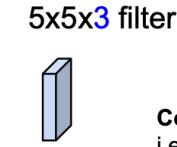
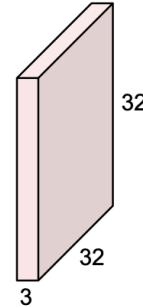
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 48

April 18, 2023

Convolution Layer

32x32x3 image



Filters always extend the full depth of the input volume

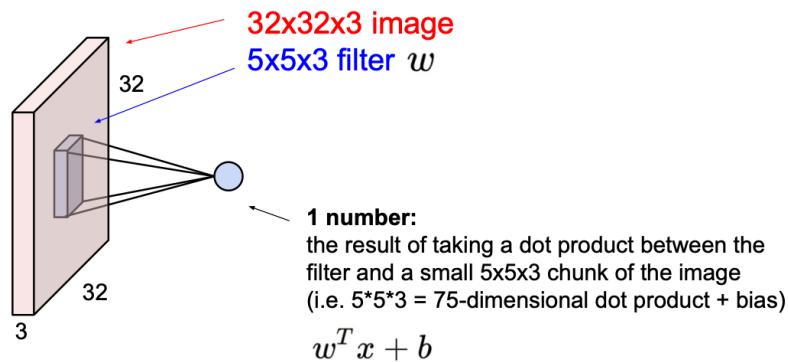
Convolve the filter with the image
i.e. "slide over the image spatially,
computing dot products"

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 50

April 18, 2023

Convolution Layer

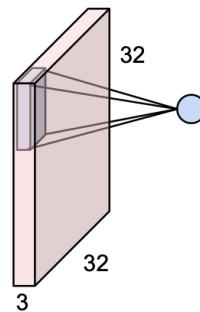


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 51

April 18, 2023

Convolution Layer

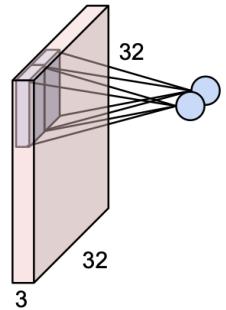


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 52

April 18, 2023

Convolution Layer

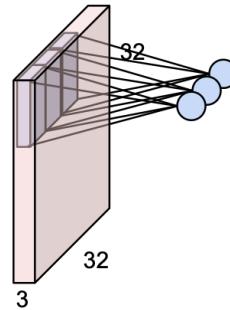


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 53

April 18, 2023

Convolution Layer

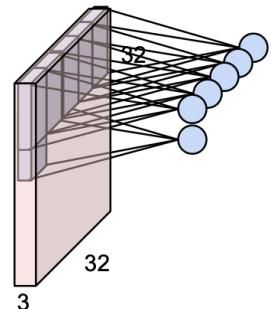


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 54

April 18, 2023

Convolution Layer

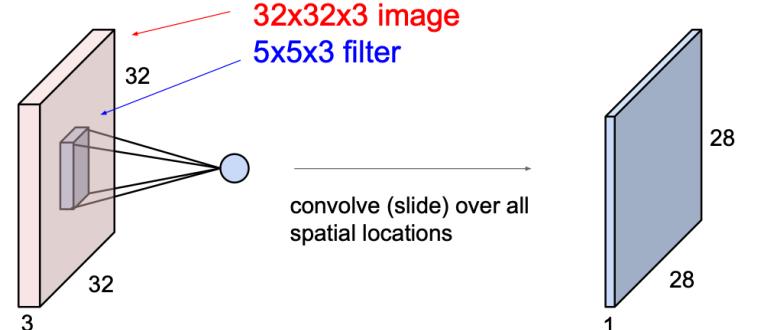


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 55

April 18, 2023

Convolution Layer



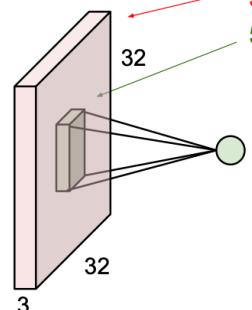
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 56

April 18, 2023

Convolution Layer

consider a second, green filter



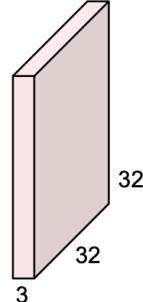
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 57

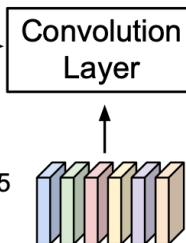
April 18, 2023

Convolution Layer

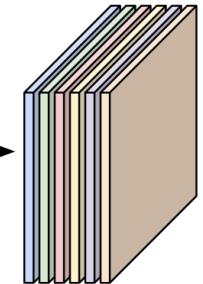
3x32x32 image



Consider 6 filters,
each 3x5x5



6 activation maps,
each 1x28x28



Stack activations to get a
6x28x28 output image!

Slide inspiration: Justin Johnson

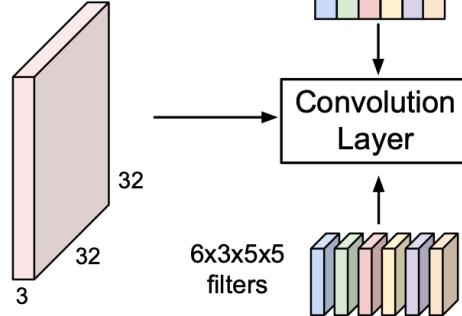
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 58

April 18, 2023

Convolution Layer

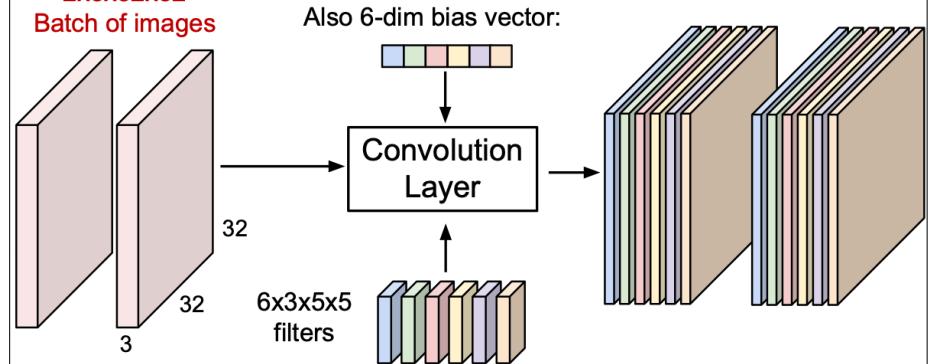
3x32x32 image



6 activation maps,
each 1x28x28

Convolution Layer

2x3x32x32
Batch of images



Slide inspiration: Justin Johnson

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 59

April 18, 2023

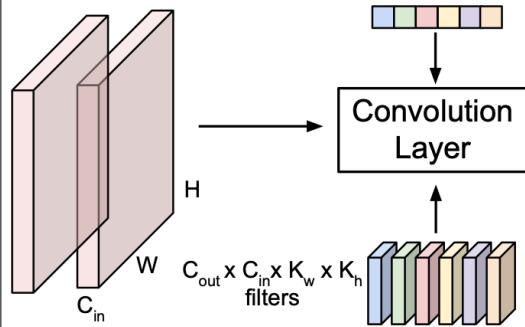
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 61

April 18, 2023

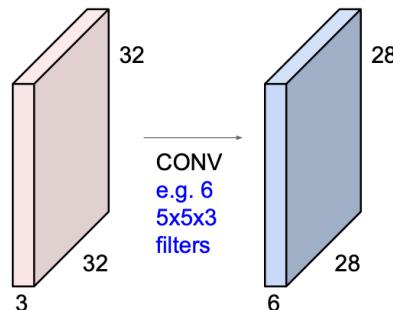
Convolution Layer

$N \times C_{in} \times H \times W$
Batch of images



$N \times C_{out} \times H' \times W'$
Batch of outputs

Preview: ConvNet is a sequence of Convolution Layers



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 62

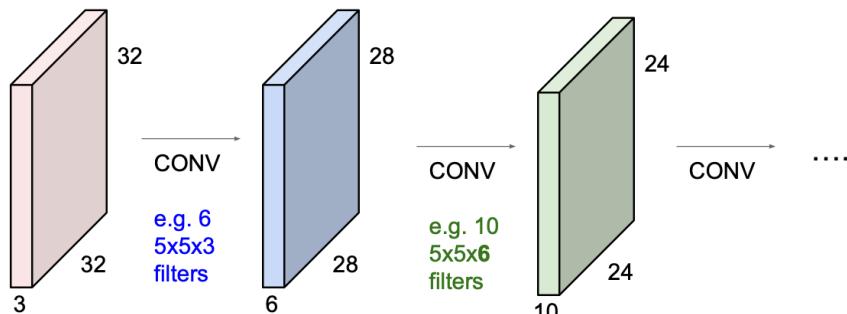
April 18, 2023

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 63

April 18, 2023

Preview: ConvNet is a sequence of Convolution Layers

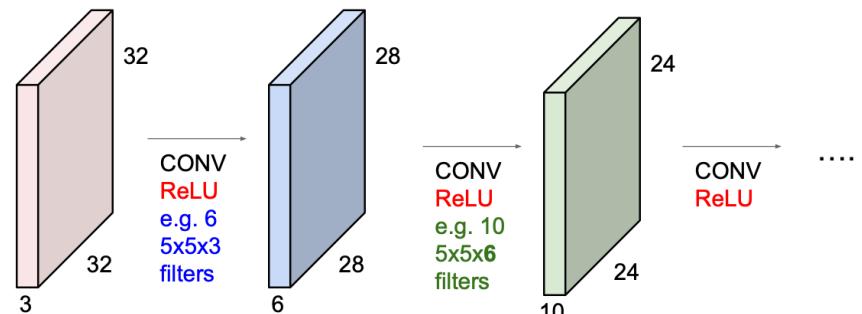


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 64

April 18, 2023

Preview: ConvNet is a sequence of Convolution Layers, interspersed with activation functions

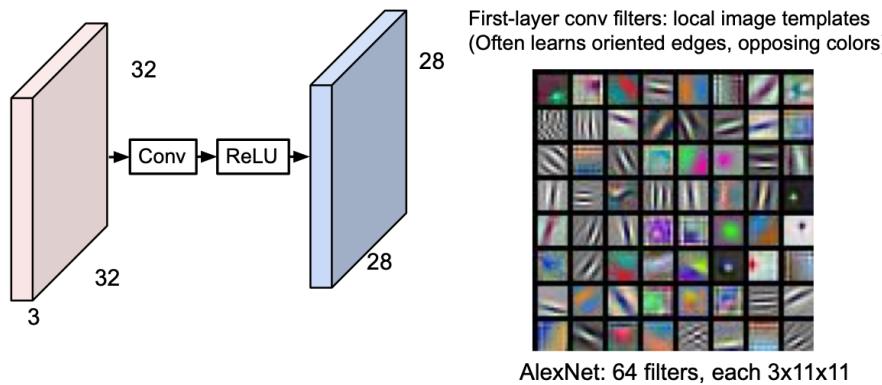


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 65

April 18, 2023

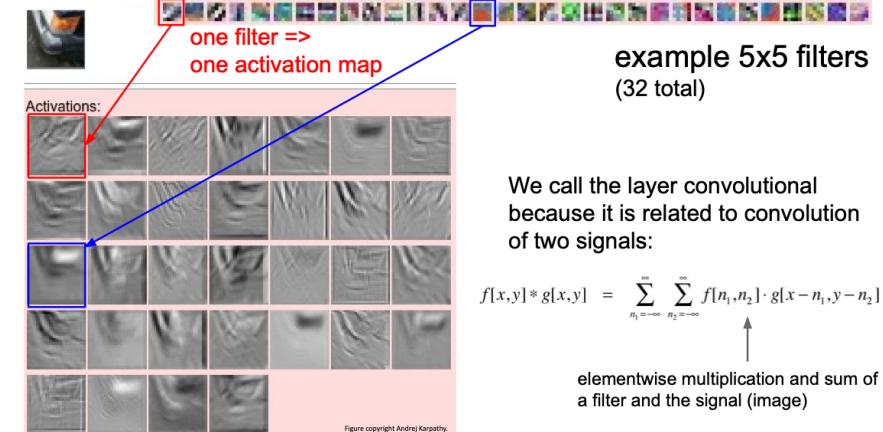
Preview: What do convolutional filters learn?



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 68

April 18, 2023

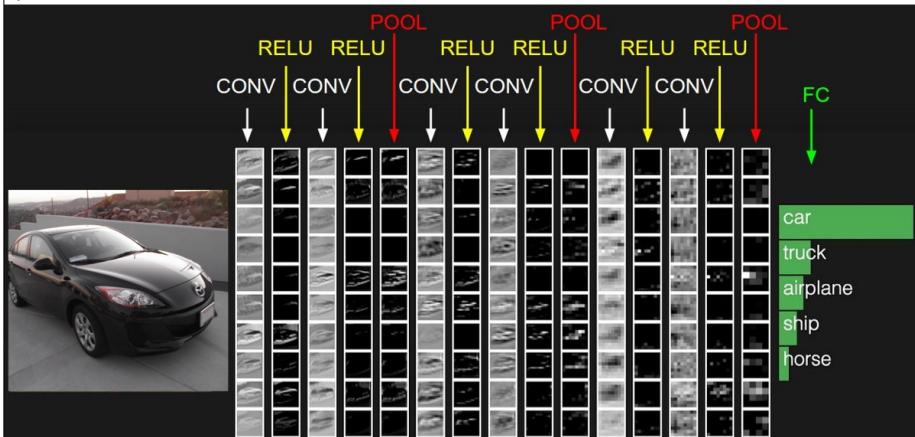


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 69

April 18, 2023

preview:

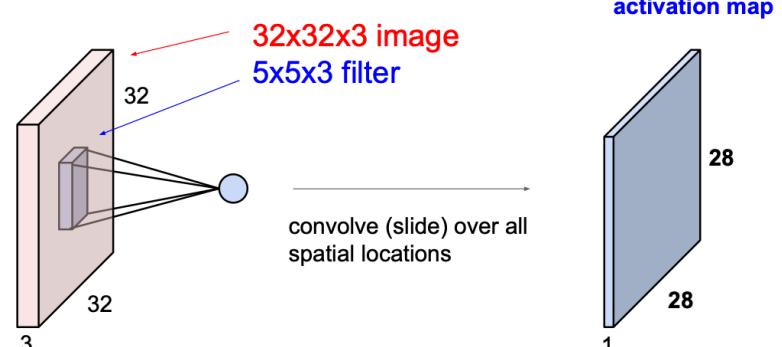


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 70

April 18, 2023

A closer look at spatial dimensions:

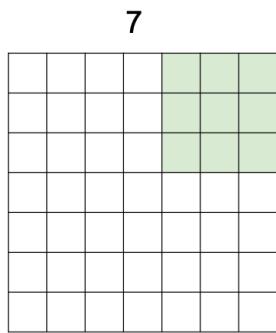


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 71

April 18, 2023

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

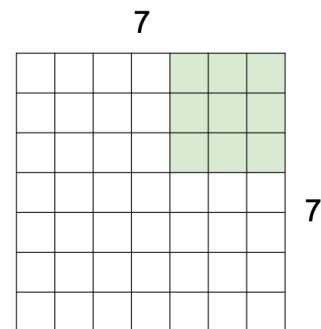
=> 5x5 output

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 76

April 18, 2023

A closer look at spatial dimensions:



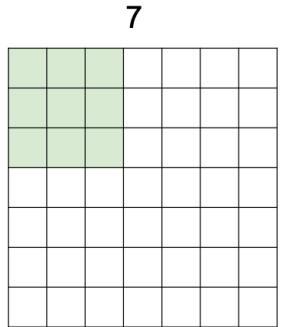
7x7 input (spatially)
assume 3x3 filter
applied with stride 2
=> 3x3 output!

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 79

April 18, 2023

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

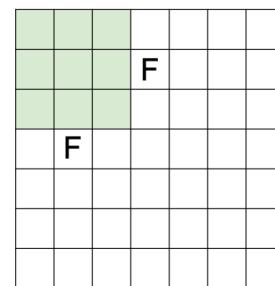
doesn't fit!
cannot apply 3x3 filter on
7x7 input with stride 3.

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 81

April 18, 2023

N



Output size:
(N - F) / stride + 1

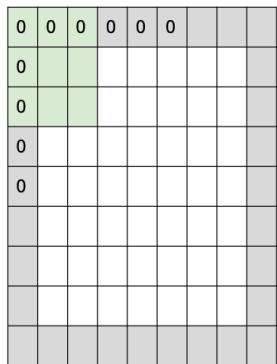
e.g. N = 7, F = 3:
stride 1 => $(7 - 3)/1 + 1 = 5$
stride 2 => $(7 - 3)/2 + 1 = 3$
stride 3 => $(7 - 3)/3 + 1 = 2.33 \vdots$

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 82

April 18, 2023

In practice: Common to zero pad the border



e.g. input 7x7
3x3 filter, applied with **stride 1**
pad with 1 pixel border => what is the output?

7x7 output!

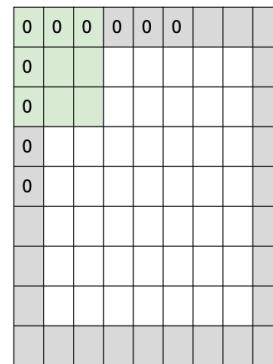
(recall):
 $(N + 2P - F) / \text{stride} + 1$

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 84

April 18, 2023

In practice: Common to zero pad the border



e.g. input 7x7
3x3 filter, applied with **stride 1**
pad with 1 pixel border => what is the output?

7x7 output!

in general, common to see CONV layers with
stride 1, filters of size FxF, and zero-padding with
(F-1)/2. (will preserve size spatially)

e.g. F = 3 => zero pad with 1
F = 5 => zero pad with 2
F = 7 => zero pad with 3

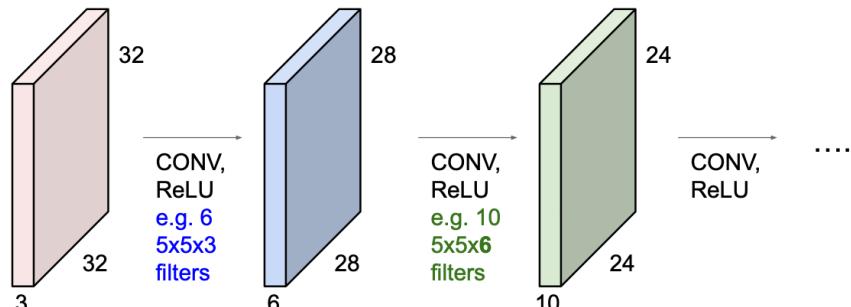
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 85

April 18, 2023

Remember back to...

E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially!
(32 -> 28 -> 24 ...). Shrinking too fast is not good, doesn't work well.



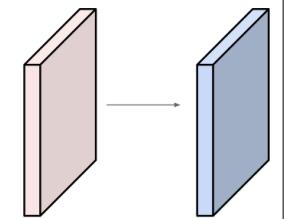
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 86

April 18, 2023

Examples time:

Input volume: $32 \times 32 \times 3$
10 5×5 filters with stride 1, pad 2



Output volume size: ?

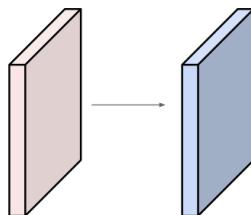
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 87

April 18, 2023

Examples time:

Input volume: $32 \times 32 \times 3$
10 5×5 filters with stride 1, pad 2



Number of parameters in this layer?

Receptive Fields

For convolution with kernel size K, each element in the output depends on a $K \times K$ **receptive field** in the input



Input

Output

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 89

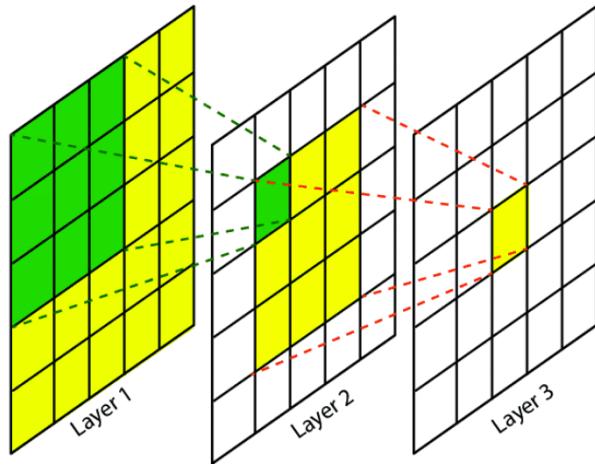
April 18, 2023

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 91

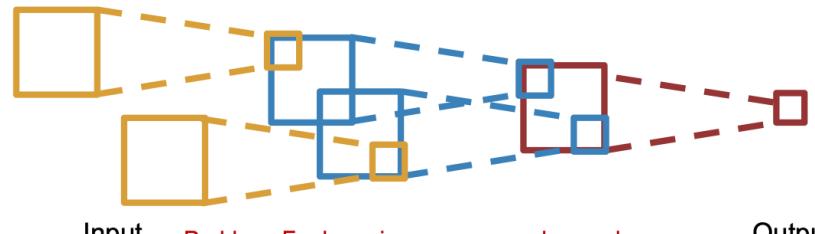
April 18, 2023

Slide inspiration: Justin Johnson



Receptive Fields

Each successive convolution adds $K - 1$ to the receptive field size
With L layers the receptive field size is $1 + L * (K - 1)$



Problem: For large images we need many layers for each output to "see" the whole image image

Solution: Downsample inside the network

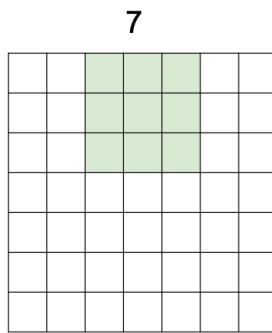
Slide inspiration: Justin Johnson

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 94

April 18, 2023

Solution: Strided Convolution



7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**
=> 3x3 output!

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 96

April 18, 2023

Convolution layer: summary

Common settings:

Let's assume input is $W_1 \times H_1 \times C$ $K = (\text{powers of 2, e.g. } 32, 64, 128, 512)$
Conv layer needs 4 hyperparameters:

- Number of filters K
- The filter size F
- The stride S
- The zero padding P

This will produce an output of $W_2 \times H_2 \times K$
where:

- $W_2 = (W_1 - F + 2P)/S + 1$
- $H_2 = (H_1 - F + 2P)/S + 1$

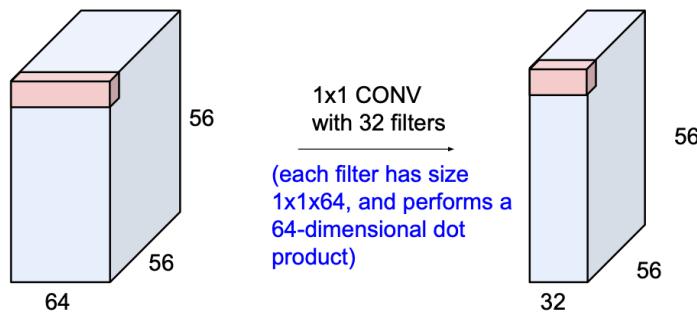
Number of parameters: F^2CK and K biases

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 98

April 18, 2023

(btw, 1x1 convolution layers make perfect sense)



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 100

April 18, 2023

Example: CONV layer in PyTorch

Conv2d

CLASS: torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True)

[SOURCE]

Applies a 2D convolution over an input signal composed of several input planes.

In the simplest case, the output value of the layer with input size (N, C_{in}, H, W) and output $(N, C_{out}, H_{out}, W_{out})$ can be precisely described as:

$$\text{out}(N_i, C_{out,j}) = \text{bias}(C_{out,j}) + \sum_{k=0}^{C_{in}-1} \text{weight}(C_{out,j}, k) * \text{input}(N_i, k)$$

where $*$ is the valid 2D cross-correlation operator, N is a batch size, C denotes a number of channels, H is a height of input planes in pixels, and W is width in pixels.

- `stride` controls the stride for the cross-correlation, a single number or a tuple.
- `padding` controls the amount of implicit zero-paddings on both sides for padding, number of points for each dimension.
- `dilation` controls the spacing between the kernel points; also known as the atrous algorithm, it is harder to describe, but this link is a nice visualization of what dilation does.
- `groups` controls the connections between inputs and outputs. `in_channels` and `out_channels` must both be divisible by `groups`. For example:
 - At `groups=1`, all inputs are convolved to all outputs.
 - At `groups=2`, the operation becomes equivalent to having two convolution layers side by side, each seeing half the input channels, and producing half the output channels, and both subsequently concatenated.
 - At `groups=in_channels`, each input channel is convolved with its own set of filters, of size $\left[\frac{C_{out}}{C_{in}}\right]$.

The parameters `kernel_size`, `stride`, `padding`, `dilation` can either be:

- a single int – in which case the same value is used for the height and width dimension
- a tuple of two ints – in which case, the first int is used for the height dimension, and the second int for the width dimension

PyTorch is licensed under BSD 3-clause.

Conv layer needs 4 hyperparameters:

- Number of filters K
- The filter size F
- The stride S
- The zero padding P

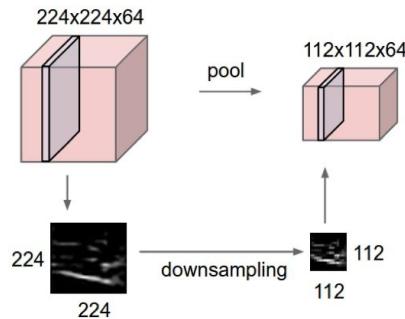
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 101

April 18, 2023

Pooling layer

- makes the representations smaller and more manageable
- operates over each activation map independently



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 107

April 18, 2023

MAX POOLING

Single depth slice

1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

max pool with 2x2 filters and stride 2

6	8
3	4

- No learnable parameters
- Introduces spatial invariance

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 5 - 109

April 18, 2023

Pooling layer: summary

Let's assume input is $W_1 \times H_1 \times C$

Conv layer needs 2 hyperparameters:

- The spatial extent **F**
- The stride **S**

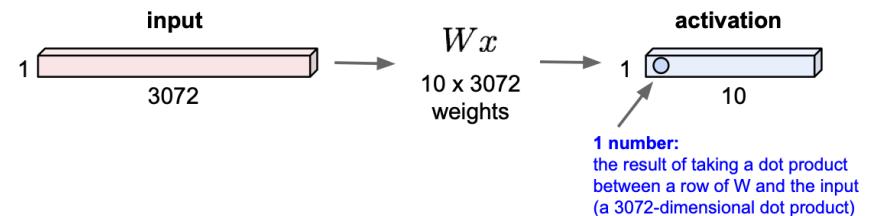
This will produce an output of $W_2 \times H_2 \times C$ where:

- $W_2 = (W_1 - F)/S + 1$
- $H_2 = (H_1 - F)/S + 1$

Number of parameters: 0

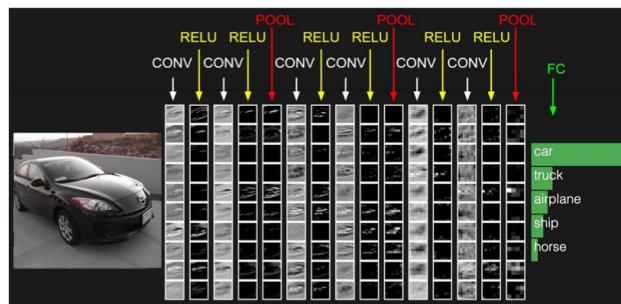
Reminder: Fully Connected Layer

32x32x3 image \rightarrow stretch to 3072 x 1

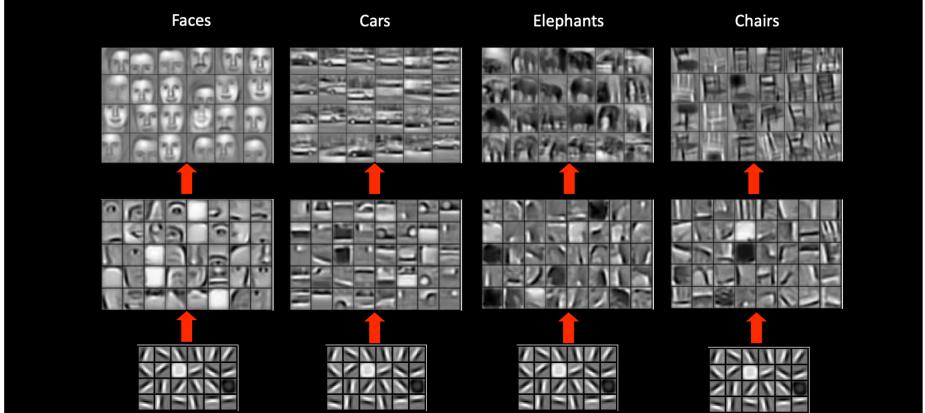


Fully Connected Layer (FC layer)

- Contains neurons that connect to the entire input volume, as in ordinary Neural Networks



Features learned from training on different object classes.



Alexnet

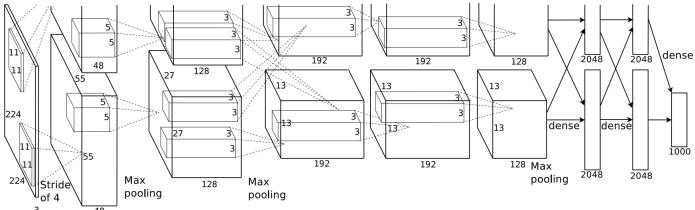


Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

"ImageNet Classification with Deep Convolutional Neural Networks", NIPS, 2012

Pytorch

```
class AlexNet(nn.Module):
    def __init__(self, num_classes: int = 1000, dropout: float = 0.5) -> None:
        super().__init__()
        self.features = nn.Sequential(
            nn.Conv2d(3, 64, kernel_size=11, stride=4, padding=2),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
            nn.Conv2d(64, 192, kernel_size=5, padding=2),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
            nn.Conv2d(192, 384, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.Conv2d(384, 256, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.Conv2d(256, 256, kernel_size=3, padding=1),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=3, stride=2),
        )
        self.avgpool = nn.AdaptiveAvgPool2d((6, 6))
        self.classifier = nn.Sequential(
            nn.Dropout(p=dropout),
            nn.Linear(256 * 6 * 6, 4096),
            nn.ReLU(inplace=True),
            nn.Dropout(p=dropout),
            nn.Linear(4096, 4096),
            nn.ReLU(inplace=True),
            nn.Linear(4096, num_classes),
        )

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        x = self.features(x)
        x = self.avgpool(x)
        x = torch.flatten(x, 1)
        x = self.classifier(x)
        return x
```

AlexNet was originally introduced in "ImageNet Classification with Deep Convolutional Neural Networks". This implementation is based instead on the "One Weird Trick for Parallelizing Convolutional Neural Networks" paper.

single-column model with 64, 192, 384, 384, 256 filters in the five convolutional layers, respectively

```
import torch
from PIL import Image
from torchvision import transforms

model = torch.hub.load('pytorch/vision:v0.10.0', 'alexnet', pretrained=True)
model.eval()

input_image = Image.open(filename)
preprocess = transforms.Compose([
    transforms.Resize(256),
    transforms.CenterCrop(224),
    transforms.ToTensor(),
    transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]),
])
input_tensor = preprocess(input_image)
# create a mini-batch as expected by the model
input_batch = input_tensor.unsqueeze(0)

if torch.cuda.is_available():
    input_batch = input_batch.to('cuda')
    model.to('cuda')

with torch.no_grad():
    output = model(input_batch)
    probabilities = torch.nn.functional.softmax(output[0], dim=0)

print(torch.topk(probabilities, 5))
# https://raw.githubusercontent.com/pytorch/hub/master/imagenet_classes.txt
```

Pytorch

Table 1: ConvNet configurations (shown in columns). The depth of the configurations increases from the left (A) to the right (E), as more layers are added (the added layers are shown in bold). The convolutional layer parameters are denoted as "conv/receptive field size)-(number of channels)". The ReLU activation function is not shown for brevity.

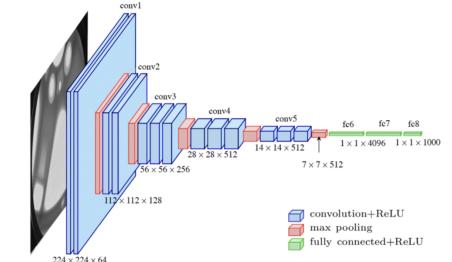
ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers		13 weight layers	15 weight layers	16 weight layers	19 weight layers
input (224x224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 conv3-64 maxpool	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
conv3-128	conv3-128	conv3-128 conv3-128 maxpool	conv3-128 conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
conv3-256	conv3-256	conv3-256 conv3-256 maxpool	conv3-256 conv3-256 conv1-256	conv3-256 conv3-256 conv3-256	conv3-256 conv3-256 conv3-256
conv3-512	conv3-512	conv3-512 conv3-512 conv1-512 maxpool	conv3-512 conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 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 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512
conv3-512	conv3-512	conv3-512 conv3-512 conv3-512 conv3-512 conv3-512 FC-4096 maxpool	conv3-512 conv3-512 conv3-512 conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512 conv3-512 fc6 fc7 fc8
conv3-512	conv3-512	conv3-512 conv3-512 conv3-512 conv3-512 conv3-512 conv3-512 FC-4096 FC-4096 FC-1000 soft-max	conv3-512 conv3-512 conv3-512 conv3-512 conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512 conv3-512 conv3-512 conv3-512	conv3-512 conv3-512 conv3-512 conv3-512 conv3-512 fc6 fc7 fc8

Table 2: Number of parameters (in millions).

Network	A,A-LRN	B	C	D	E
Number of parameters	133	133	134	138	144

"Very Deep Convolutional Networks for Large-Scale Image Recognition", ICLR, 2015

VGG 16



```

def VGG16():
    model = Sequential()

    model.add(Conv2D(filters=64, kernel_size=(3,3), padding='same', activation='relu', input_shape=(224,224,3)))
    model.add(Conv2D(filters=64, kernel_size=(3,3), padding='same', activation='relu'))
    model.add(MaxPooling2D(pool_size=(2,2), strides=2))

    filters_convs = [(128, 2), (256, 3), (512, 3), (512,3)]
    for n_filters, n_convs in filters_convs:
        for _ in np.arange(n_convs):
            model.add(Conv2D(filters=n_filters, kernel_size=(3,3), padding='same', activation='relu'))
    model.add(MaxPooling2D(pool_size=(2,2), strides=2))

    model.add(Flatten())
    model.add(Dense(4096, activation = 'relu'))
    model.add(Dropout(0.5))
    model.add(Dense(4096, activation = 'relu'))
    model.add(Dropout(0.5))
    model.add(Dense(1000,activation = 'softmax'))

    opt = SGD(lr = 0.01)
    model.compile(loss = categorical_crossentropy, optimizer = opt, metrics = ['accuracy'])

    return model

https://medium.com/mlearning-ai/an-overview-of-vgg16-and-nin-models-96e4bf398484

```

Keras

```

# load vgg model
from keras.applications.vgg16 import VGG16

# load the model
model = VGG16()

# prints summary
model.summary()
# Total params: 138,357,544
# Trainable params: 138,357,544

```

```

tf.keras.applications.VGG16(
    include_top=True,
    weights="imagenet",
    input_tensor=None,
    input_shape=None,
    pooling=None,
    classes=1000,
    classifier_activation="softmax",
)

```

[ConvNetJS demo: training on CIFAR-10]

ConvNetJS CIFAR-10 demo

Description

This demo trains a Convolutional Neural Network on the [CIFAR-10 dataset](#) in your browser, with nothing but Javascript. The state of the art on this dataset is about 90% accuracy and human performance is at about 94% (not perfect as the dataset can be a bit ambiguous). I used [this python script](#) to prepare the [original files](#) (python version) into batches of images that can be easily loaded into page DOM with img tags.

This dataset is more difficult and it takes longer to train a network. Data augmentation includes random flipping and random image shift by up to 2px horizontally or vertically.

By default, in this demo we're using Adadelta which is one of per-parameter adaptive step size methods, so we don't have to worry about changing learning rates or momentum over time. However, I still included the text fields for changing these if you'd like to play around with SGD+Momentum trainer.

Report questions/bugs/suggestions to [@karpathy](#).



<http://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

Summary

- ConvNets stack CONV,POOL,FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Historically architectures looked like
 $[(CONV-RELU)*N-POOL?] * M - (FC-RELU)*K, SOFTMAX$
where N is usually up to ~5, M is large, $0 \leq K \leq 2$.
- But recent advances such as ResNet/GoogLeNet have challenged this paradigm

Transfer learning

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 7 -

129

April 25, 2023

You need a lot of data if you want to train/use CNNs?

Fei-Fei Li, Yunzhu Li, Ruohan Gao

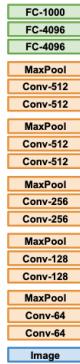
Lecture 7 -

130

April 25, 2023

Transfer Learning with CNNs

1. Train on Imagenet



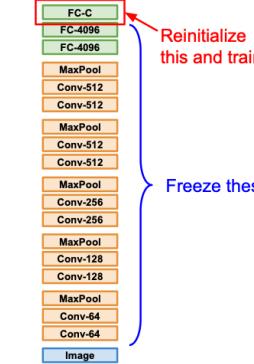
Donahue et al. "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al. "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

Transfer Learning with CNNs

1. Train on Imagenet



2. Small Dataset (C classes)



Donahue et al. "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al. "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 7 -

131

April 25, 2023

Fei-Fei Li, Yunzhu Li, Ruohan Gao

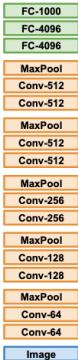
Lecture 7 -

132

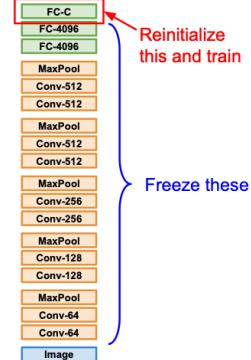
April 25, 2023

Transfer Learning with CNNs

1. Train on Imagenet

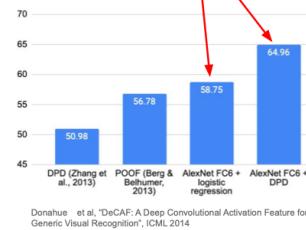


2. Small Dataset (C classes)



Donahue et al., "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al., "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

Finetuned from AlexNet



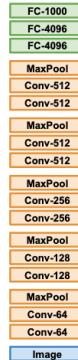
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 7 -

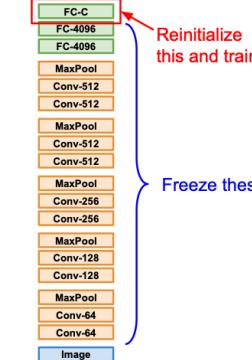
133 April 25, 2023

Transfer Learning with CNNs

1. Train on Imagenet

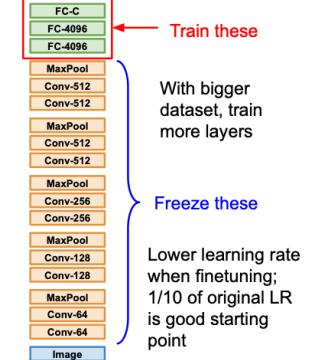


2. Small Dataset (C classes)



Donahue et al., "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014
Razavian et al., "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

3. Bigger dataset



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 7 -

134 April 25, 2023

More specific

More generic

	very similar dataset	very different dataset
very little data	Use Linear Classifier on top layer	You're in trouble... Try linear classifier from different stages
quite a lot of data	Finetune a few layers	Finetune a larger number of layers

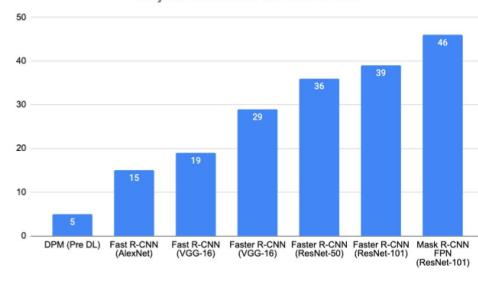
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 6 -

143 April 20, 2023

Transfer learning with CNNs - Architecture matters

Object detection on MSCOCO

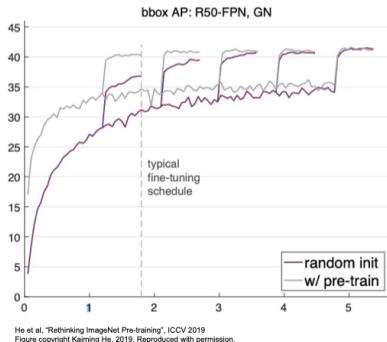


Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 6 - 147

April 20, 2023

Transfer learning with CNNs is pervasive... But recent results show it might not always be necessary!



Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 6 - 148 April 20, 2023

Batch Normalization

Consider a single layer $y = Wx$

The following could lead to tough optimization:

- Inputs x are not *centered around zero* (need large bias)
- Inputs x have different scaling per-element (entries in W will need to vary a lot)

Idea: force inputs to be “nicely scaled” at each layer!

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 6 - 4

April 20, 2023

Batch Normalization

[Ioffe and Szegedy, 2015]

“you want zero-mean unit-variance activations? just make them so.”

consider a batch of activations at some layer. To make each dimension zero-mean unit-variance, apply:

$$\hat{x}^{(k)} = \frac{x^{(k)} - \text{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

this is a vanilla
differentiable function...

Fei-Fei Li, Yunzhu Li, Ruohan Gao

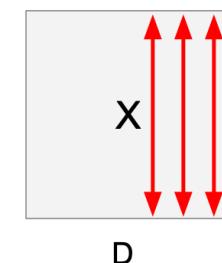
Lecture 6 - 5

April 20, 2023

Batch Normalization

[Ioffe and Szegedy, 2015]

Input: $x : N \times D$



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j} \quad \text{Per-channel mean, shape is D}$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2 \quad \text{Per-channel var, shape is D}$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}} \quad \text{Normalized x, Shape is } N \times D$$

Problem: What if zero-mean, unit variance is too hard of a constraint?

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 6 - 7

April 20, 2023

Batch Normalization

[Ioffe and Szegedy, 2015]

Input: $x : N \times D$

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j} \quad \text{Per-channel mean, shape is } D$$

Learnable scale and shift parameters:

$\gamma, \beta : D$

Learning $\gamma = \sigma$,
 $\beta = \mu$ will recover the identity function!

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2 \quad \text{Per-channel var, shape is } D$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}} \quad \text{Normalized x, Shape is } N \times D$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j \quad \text{Output, Shape is } N \times D$$

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 6 - 8

April 20, 2023

Estimates depend on minibatch; can't do this at test-time!

Batch Normalization: Test-Time

Input: $x : N \times D$

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j} \quad \text{Per-channel mean, shape is } D$$

Learnable scale and shift parameters:

$\gamma, \beta : D$

Learning $\gamma = \sigma$,
 $\beta = \mu$ will recover the identity function!

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2 \quad \text{Per-channel var, shape is } D$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}} \quad \text{Normalized x, Shape is } N \times D$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j \quad \text{Output, Shape is } N \times D$$

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 6 - 9

April 20, 2023

Batch Normalization: Test-Time

Input: $x : N \times D$

$$\mu_j = \text{(Running) average of values seen during training} \quad \text{Per-channel mean, shape is } D$$

Learnable scale and shift parameters:

$\gamma, \beta : D$

During testing batchnorm becomes a linear operator!
 Can be fused with the previous fully-connected or conv layer

$$\sigma_j^2 = \text{(Running) average of values seen during training} \quad \text{Per-channel var, shape is } D$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}} \quad \text{Normalized x, Shape is } N \times D$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j \quad \text{Output, Shape is } N \times D$$

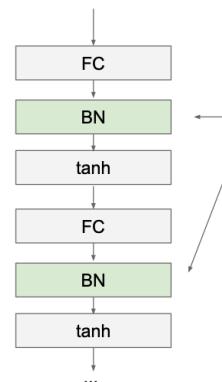
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 6 - 10

April 20, 2023

Batch Normalization

[Ioffe and Szegedy, 2015]



Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.

$$\hat{x}^{(k)} = \frac{x^{(k)} - E[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

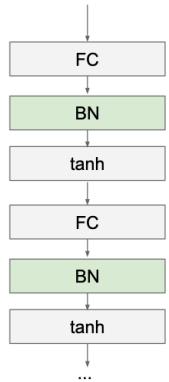
Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 6 - 11

April 20, 2023

Batch Normalization

[Ioffe and Szegedy, 2015]



- Makes deep networks **much** easier to train!
- Improves gradient flow
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Zero overhead at test-time: can be fused with conv!
- Behaves differently during training and testing: this is a very common source of bugs!

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 6 - 12

April 20, 2023

Batch Normalization for ConvNets

Batch Normalization for
fully-connected networks

$$\mathbf{x}: N \times D$$

Normalize

$$\mu, \sigma: 1 \times D$$

$$\gamma, \beta: 1 \times D$$

$$y = \gamma(x - \mu) / \sigma + \beta$$

Batch Normalization for
convolutional networks
(Spatial Batchnorm, BatchNorm2D)

$$\mathbf{x}: N \times C \times H \times W$$

Normalize

$$\mu, \sigma: 1 \times C \times 1 \times 1$$

$$\gamma, \beta: 1 \times C \times 1 \times 1$$

$$y = \gamma(x - \mu) / \sigma + \beta$$

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 6 - 13

April 20, 2023

Layer Normalization

Batch Normalization for
fully-connected networks

$$\mathbf{x}: N \times D$$

Normalize

$$\mu, \sigma: 1 \times D$$

$$\gamma, \beta: 1 \times D$$

$$y = \gamma(x - \mu) / \sigma + \beta$$

Layer Normalization for
fully-connected networks

Same behavior at train and test!
Can be used in recurrent networks

$$\mathbf{x}: N \times D$$

Normalize

$$\mu, \sigma: N \times 1$$

$$\gamma, \beta: 1 \times D$$

$$y = \gamma(x - \mu) / \sigma + \beta$$

Ba, Kiros, and Hinton, "Layer Normalization", arXiv 2016

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 6 - 14

April 20, 2023

Instance Normalization

Batch Normalization for
convolutional networks

$$\mathbf{x}: N \times C \times H \times W$$

Normalize

$$\mu, \sigma: 1 \times C \times 1 \times 1$$

$$\gamma, \beta: 1 \times C \times 1 \times 1$$

$$y = \gamma(x - \mu) / \sigma + \beta$$

Instance Normalization for
convolutional networks
Same behavior at train / test!

$$\mathbf{x}: N \times C \times H \times W$$

Normalize

$$\mu, \sigma: N \times C \times 1 \times 1$$

$$\gamma, \beta: 1 \times C \times 1 \times 1$$

$$y = \gamma(x - \mu) / \sigma + \beta$$

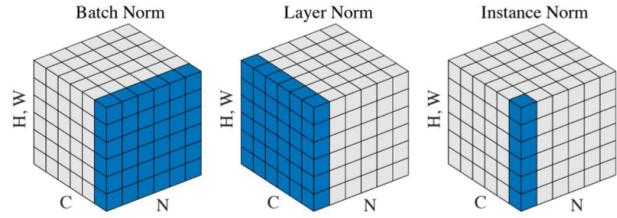
Ulyanov et al, Improved Texture Networks: Maximizing Quality and Diversity in Feed-forward Stylization and Texture Synthesis, CVPR 2017

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 6 - 15

April 20, 2023

Comparison of Normalization Layers



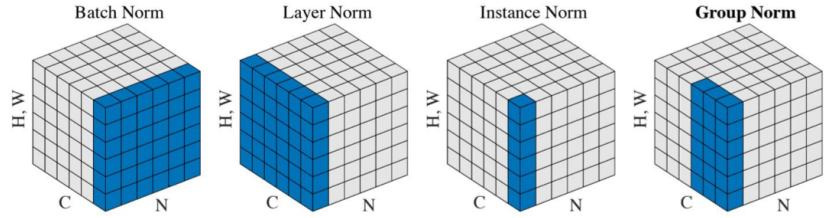
Wu and He, "Group Normalization", ECCV 2018

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 6 - 16

April 20, 2023

Group Normalization



Wu and He, "Group Normalization", ECCV 2018

Fei-Fei Li, Yunzhu Li, Ruohan Gao

Lecture 6 - 17

April 20, 2023