# Convolutional neural network

2024年9月12日 22:23
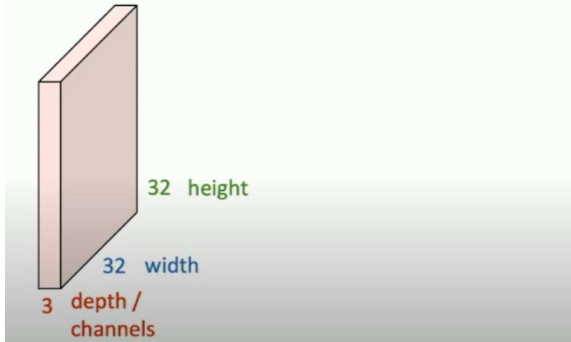
$32 \times 32 \times 3 \rightarrow 3072 \times 1$   fully connected
$\underbrace{}_{1024}$

3 Components: Conv layers, Pooling, Normalization

1 Convolution Layers

1)

1° ③x32x32 image: preserve spatial structure

32  height

32  width

3  depth / channels

③x5x5 filter

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

3x5x5 filter

slide over all spatial loc →

activation map

28

28

1

1×28×28

**1 number:**
the result of taking a dot product between the filter and a small 3x5x5 chunk of the image
(i.e. 3*5*5 = 75-dimensional dot product + bias)

$w^T x + b$

32

32

3

3×32×32

2° 可以有更多 filters, 提取不同特征 → act- maps, stack them together.

## Convolution Layer

3x32x32 image

Consider 6 filters, each 3x5x5
6-dim bias vec

Convolution Layer

32

32

3

6x3x5x5
(W) filters

6 activation maps, each 1x28x28

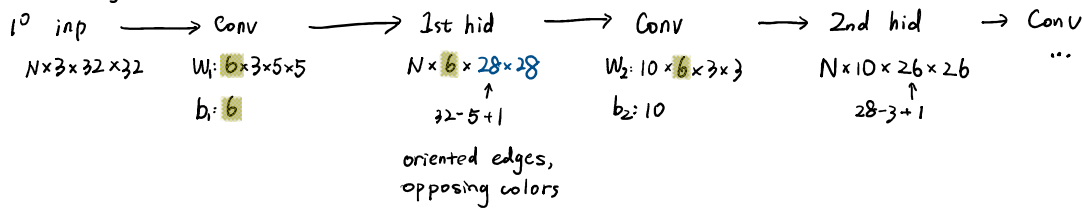Stack activations to get a 6x28x28 output image!

3° Batch of img  $2 \times 3 \times 32 \times 32 \xrightarrow{6 \times 3 \times 5 \times 5 \text{ filters}} 2 \times 6 \times 28 \times 28$  Batch of outp

$N \times C_{in} \times H \times W$   $C_{out} \times C_{in} \times K_w \times K_h$   $N \times C_{out} \times H' \times W'$

分区 CV 的第 1 页

$3°$ Batch of img $2\times3\times32\times32 \xrightarrow{6\times3\times5\times5 \text{ filters}} 2\times6\times28\times28$ Batch of outp

$\quad\quad N\times C_{in}\times H\times W \quad\quad C_{out}\times C_{in}\times K_w\times K_h \quad N\times C_{out}\times H'\times W'$

$\quad\quad\quad\quad\quad\uparrow$ inp channels $\quad\quad\uparrow$ mayb dif

## 2) Stacking Convolutions

$1°$ inp $\longrightarrow$ Conv $\longrightarrow$ 1st hid $\longrightarrow$ Conv $\longrightarrow$ 2nd hid $\longrightarrow$ Conv

$N\times3\times32\times32 \quad W_1: 6\times3\times5\times5 \quad N\times6\times28\times28 \quad W_2: 10\times6\times3\times3 \quad N\times10\times26\times26 \quad \cdots$

$\quad\quad\quad\quad\quad b_1: 6 \quad\quad\quad\quad \underset{32-5+1}{\uparrow} \quad\quad\quad b_2: 10 \quad\quad\quad \underset{28-3+1}{\uparrow}$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad$ oriented edges,
$\quad\quad\quad\quad\quad\quad\quad\quad\quad$ opposing colors

Stacking 2 conv $\rightarrow$ another conv $\quad y = W_2 W_1 x \quad$ linear classifier

## $2°$ Size

① In: W $\quad\quad\quad 7\times7$

$\quad$ Fil: K $\quad\quad\quad 3\times3$

$\quad$ Out W$-$K$+1 \quad 5\times5$

$\quad$ Feature maps shrink, lim #layers, sol.

② Padding, $+0s$ around the input

$\quad$ In: W

$\quad$ Fil: K

$\quad$ Pad: P

$\quad$ Out: W$-$K$+1+2$P

(common: P $= (K-1)/2$ to make size in $=$ out)

## $3°$ Receptive Fields



| Layer | 3 | 2 | 1 | 0 |
|---|---|---|---|---|
| Field size | $1+3\times2=7$ | $1+2\times2=5$ | $1+1\times2=3$ | 1 |

$\quad$ L Layer receptive field size $= 1 + L\times(K-1)$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \underset{\text{fil size}}{\uparrow}$

e.g. input $1000\times1000$, $K=3$, $L=?$

$\quad\quad 1000 = 1 + L\times(3-1)$

$\quad\quad\quad L = 999\div2 = 499.5$

Large imgs need many layers for each outps to "see" the whole img, sol:
$\quad\quad\quad\quad\quad\quad$ (gloabal context)

## $4°$ Downsample inside the network

$\quad$ Controlling the stride: indirectly downsampling $\because$ fewer data pt are being processed.

$\quad$ ① Strided conv:

$\quad\quad$ In W $\quad\quad\quad\quad\quad\quad\quad 3\times32\times32$

$\quad\quad$ Fil K $\quad\quad\quad\quad\quad\quad\quad 10\times5\times5$

In  W                                 $3 \times 32 \times 32$

Fil  K                                $10 \times 5 \times 5$

Pad  P                                        2
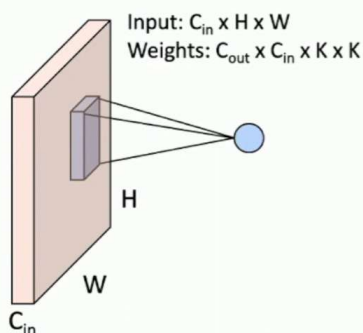
Str  S                                        1

Out  $(W - K + 2P)/S + 1$      $(32-5+2 \times 2)/1 + 1 = 32$    $10 \times 32 \times 32$

\# learnable para $\begin{cases} \text{① para/fl} & 3 \times 5 \times 5 + 1 = 76 \\ & \qquad\qquad\qquad \uparrow \\ & \qquad\qquad\quad bias \\ \text{② } \times \text{\#fl} & 76 \times 10 = 760 \end{cases}$

\# multiply-add oper $\begin{cases} \text{\# outs} & 10 \times 32 \times 32 = 10240 \\ \text{inner product/out} & 3 \times 5 \times 5 = 75 \\ total & 75 \times 10240 = 768 \text{ K} \end{cases}$



$32 \times 64 \times 1 \times 1$

1x1 CONV
with 32 filters
$\longrightarrow$

(each filter has size 1x1x64,
and performs a 64-
dimensional dot product)

Stacking 1x1 conv layers
gives MLP operating on
each input position

Lin et al, "Network in Network", ICLR 2014

5° Other conv



So far: 2D Convolution
Input: $C_{in}$ x H x W
Weights: $C_{out}$ x $C_{in}$ x K x K

1D Convolution
Input: $C_{in}$ x W
Weights: $C_{out}$ x $C_{in}$ x K

Audio...

3D Convolution
Input: $C_{in}$ x H x W x D
Weights: $C_{out}$ x $C_{in}$ x K x K x K

3D data

$C_{in}$-dim vector
at each point
in the volume

2 Pooling Layer: downsample 缩-压缩



224x224x64

pool
$\longrightarrow$

112x112x64

Hyperparameters:
Kernel Size
Stride
Pooling function

**Hyperparameters:**
Kernel Size
Stride
Pooling function

1° Max Pooling



Max pooling with 2x2 kernel size and stride 2

Introduces **invariance** to small spatial shifts
No learnable parameters!
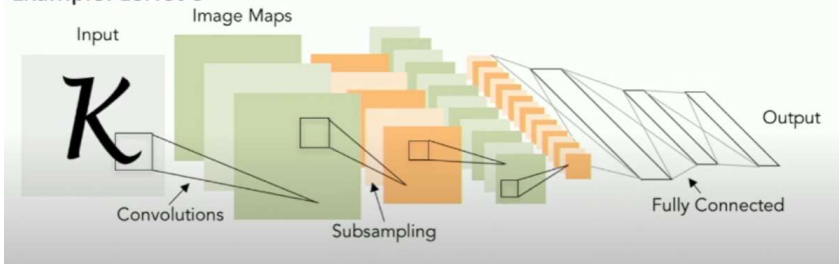
# 3 Convolutional Networks

Classic architecture: [Conv, ReLU, Pool] x N, flatten, [FC, ReLU] x N, FC

Example: LeNet-5



| Layer | Output Size | Weight Size |
|---|---|---|
| Input | 1 x 28 x 28 | |
| Conv ($C_{out}$=20, K=5, P=2, S=1) | 20 x 28 x 28 | 20 x 1 x 5 x 5 |
| ReLU | 20 x 28 x 28 | |
| MaxPool(K=2, S=2) | 20 x 14 x 14 | |
| Conv ($C_{out}$=50, K=5, P=2, S=1) | 50 x 14 x 14 | 50 x 20 x 5 x 5 |
| ReLU | 50 x 14 x 14 | |
| MaxPool(K=2, S=2) | 50 x 7 x 7 | |
| Flatten | 2450 | |
| Linear (2450 -> 500) | 500 | 2450 x 500 |
| ReLU | 500 | |

spatial size↓ (pl & strided)

# channels ↑ (total volume is preserved)

ReLU 不一定需要.

Deep NN: hard to train (converge), sol:

# 4 Normalization

1) Batch Nor in fully-connected

1° Idea: "Normalize" the outputs of a layer so they have zero mean and unit variance

Why? Helps reduce "internal covariate shift", improves optimization

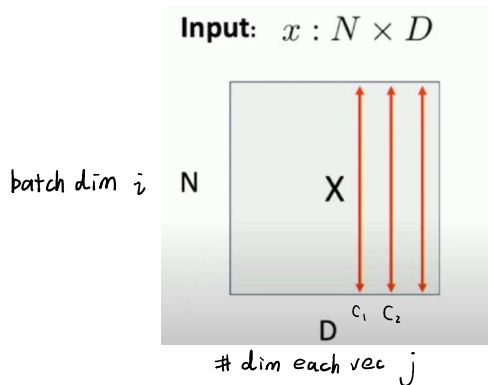We can normalize a batch of activations like this:

Why? Helps reduce "internal covariate shift", improves optimization

We can normalize a batch of activations like this:

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

This is a **differentiable function**, so we can use it as an operator in our networks and backprop through it!

convert inp → more standardized dist

**Input:** $x : N \times D$

batch dim $i$  N

X

$C_1$  $C_2$

D

# dim each vec $j$

per channel mean, shape $D$

$$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j} \qquad ①$$

~ std, shape $D$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^{N} (x_{i,j} - \mu_j)^2 \qquad ②$$

Normalized X,  $N \times D$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}} \qquad ③$$

<span style="color:red">What if $\mu = 0$, unit vec: too hard of a constrained, sol:</span>

$2°$ + Learnable scale & shift para: $\gamma, \beta : D$

Learning $\gamma = \sigma$, $\beta = \mu$, recover identity func.   ① ② ③

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j \qquad ④ \qquad N \times D$$

<span style="color:red">① ~③ Estimated depend on minibatch; ✗ do this at test-time! sol:</span>

$3°$ Test - Time use $\mu_j$, $\sigma_j^2$ got from training

## Batch Normalization: Test-Time

**Input:** $x : N \times D$

**Learnable scale and shift parameters:**

$\gamma, \beta : D$

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

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

$\sigma_j^2 =$ (Running) average of values seen during training — Per-channel std, shape is D

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

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

During testing batchnorm becomes a linear operator! Can be <u>fused</u> with the previous fully-connected or conv layer   v. 融合.

2) Batch N~ in Conv

## Batch Normalization for **fully-connected networks**

$$x: \quad N \times D$$

Normalize $\downarrow$

$$\mu, \sigma: \quad 1 \times D$$
$$\gamma, \beta: \quad 1 \times D$$
$$y = \gamma(x-\mu)/\sigma + \beta$$

*only on batch*

## Batch Normalization for **convolutional networks** (Spatial Batchnorm, BatchNorm2D)

$$x: \quad N \times C \times H \times W$$

Normalize $\downarrow \quad \downarrow \quad \downarrow$

$$\mu, \sigma: \quad 1 \times C \times 1 \times 1$$
$$\gamma, \beta: \quad 1 \times C \times 1 \times 1$$
$$y = \gamma(x-\mu)/\sigma + \beta$$

*on batch, spatial dims*

3)

| FC |
| BN |
| tanh |
| FC |
| BN |
| tanh |

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

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

⊕
- Makes deep networks **much** easier to train!
- 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!

ImageNet accuracy

- - - Inception
- - - BN-Baseline
...... BN-x5
—— BN-x30
- - BN-x5-Sigmoid
◆ Steps to match Inception

Training iterations

⊖
- Not well-understood theoretically (yet)
- Behaves differently during training and testing: this is a very common source of bugs!

4) Others Norm



Batch Norm     Layer Norm     Instance Norm     **Group Norm**