

Lecture 23

[Convolutional neural networks] (CNN)

NN for images

Continuous convolution:

$$[f \times g](\tilde{z}) = \int d\tilde{u} f(\tilde{u}) g(\tilde{z} - \tilde{u})$$

Discrete convolution:

if $f(\tilde{u})$ is defined at $\{-L, \dots, 0, \dots, L\}$:

$f(-L) = w_{-L}, \dots, f(L) = w_L$ filter/kernel

if $g(\tilde{u})$ is defined at $\{-N, \dots, 0, \dots, N\}$:

$g(-N) = x_{-N}, \dots, g(N) = x_N$ input

Then $[w \times x](i) = w_{-L} x_{i+L} + \dots + w_0 x_i + \dots + w_L x_{i-L}$

(1) $i+L$ & $i-L$ must be within the $[-N, N]$ range

Similarly, one can define a cross-correlation:

$$[w \otimes x](i) = w_{-L} x_{i-L} + \dots + w_{-1} x_{i-1} + \dots + w_L x_{i+L}$$

(2)

(1) & (2) are the same if $w_{-k} = w_k$.

One can restrict Eq.(2) to non-negative indices:

$$[w \otimes x](i) = \sum_{u=0}^{L-1} w_u x_{i+u} \quad \text{filter of length } L$$

Ex.

0 1 2 3 4 5 6	\times	1 2	=	2 5 8 11 14 17
input		kernel		output

Similarly, in 2D

$$[w \otimes x](i, j) = \sum_{u=0}^{H-1} \sum_{v=0}^{W-1} w_{u,v} x_{i+u, j+v}$$

2D $H \times W$ filter

Ex.

0	13	2
3	4	5
6	7	8

 \times

0	1
2	3

 $=$

89	25
37	43

input kernel
 (feature) (feature map)

In matrix notation,

$$\begin{pmatrix} w_1 & w_2 \\ w_3 & w_4 \end{pmatrix} \times \begin{pmatrix} x_1 & x_2 & x_3 \\ x_4 & x_5 & x_6 \\ x_7 & x_8 & x_9 \end{pmatrix} = \begin{pmatrix} w_1x_1 + w_2x_2 + w_3x_4 + w_4x_5 \dots \\ w_1x_4 + w_2x_5 + w_3x_7 + w_4x_8 \dots \end{pmatrix} \quad (1)$$

$$\textcircled{=} \begin{pmatrix} w_1 & w_2 & 0 & w_3 & w_4 & 0 & 0 & 0 & 0 \\ 0 & w_1 & w_2 & 0 & w_3 & w_4 & 0 & 0 & 0 \\ 0 & 0 & 0 & w_1 & w_2 & 0 & w_3 & w_4 & 0 \\ 0 & 0 & 0 & 0 & w_1 & w_2 & 0 & w_3 & w_4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_9 \end{pmatrix}$$

↑
 sparse weight matrices,
 shared weights

4x9 ↓
 flattened

In general, $f_h \times f_w$ filter applied to an image of size $x_h \times x_w$ produces an output of size $(x_h - f_h + 1) \times (x_w - f_w + 1)$: use zero padding of the image (borders of zeros) so that the output retains the $x_h \times x_w$ size.

With symmetric padding of size p_h, p_w
we have:

$$x_h \times x_w \xrightarrow{\text{padding}} (x_h + 2p_h) \times (x_w + 2p_w) \xrightarrow{\text{convolution}} (x_h + 2p_h - f_h + 1) \times (x_w + 2p_w - f_w + 1) \rightarrow$$

$$\rightarrow x_h \times x_w.$$

=====

$$\begin{cases} 2p_h = f_h - 1 \\ 2p_w = f_w - 1 \end{cases}$$

Strided convolution: apply filters
with step sizes > 1 or, equivalently,
~~steps every~~ keeps only $\frac{1}{(S_h S_w)}$
convolutions

$\uparrow \uparrow$
Step sizes

Output size:

$$\left\lfloor \frac{x_h + 2p_h - f_h + S_h}{S_h} \right\rfloor \times \left\lfloor \frac{x_w + 2p_w - f_w + S_w}{S_w} \right\rfloor$$

smaller than input

Multiple channels:

For ex., RGB channels in color images

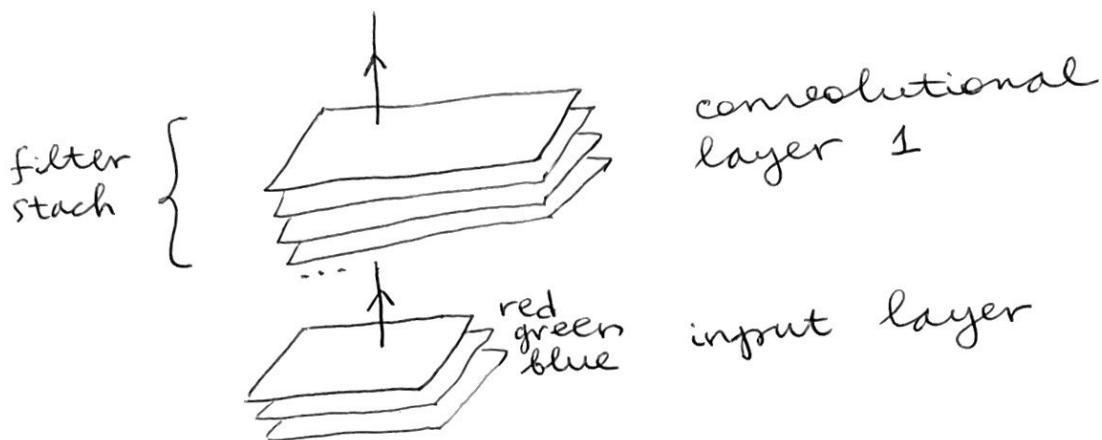
Then $z_{i,j} = b + \sum_{c=0}^{C-1} \sum_{u=0}^{H-1} \sum_{v=0}^{W-1} w_{u,v,c} x_{s_i+u, s_j+v, c}$

(*) bias $S = S_h = S_w \leftarrow \text{stride}$

Eq. (*) provides a single filter; with multiple filters, we have:

$$z_{i,j,d} = b_d + \sum_c \sum_u \sum_v w_{u,v,c,d} x_{s_i+u, s_j+v, c}$$

(**) filter index: $d = 0, \dots, D-1$



Pointwise convolution:

no "averaging" across locations

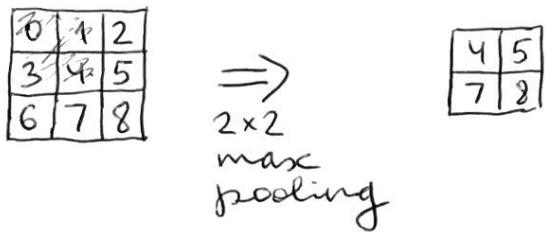
$$z_{i,j,d} = b_d + \sum_{c=0}^{C-1} x_{i,j,c} w_{0,0,c,d}$$

Same spatial dimensions: $H \times W$,
but the number of channels
goes from C to D .

Pooling layers: convolutions preserve
spatial information (at reduced resolution)

Idea: "pool" incoming values by taking
a max (max pooling) or an average
(ave pooling).

Ex.



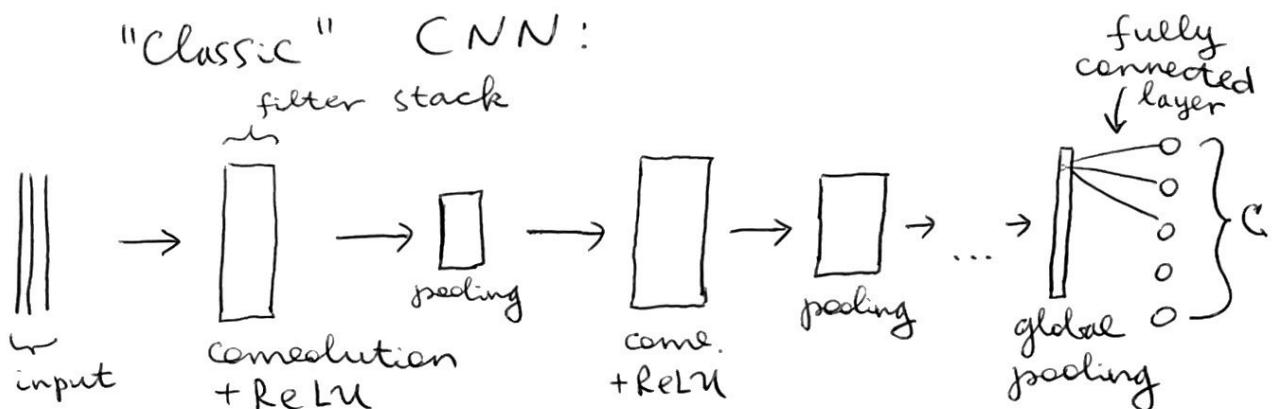
[Pooling is usually done ~~on~~ separately]
on each feature channel.

at the end of the CNN, one can do
global average pooling: (filter size = $H \times W$)

$$H \times W \times D \rightarrow 1 \times 1 \times D = D$$

↑
channels

This D -dim vector can be fed into a softmax output layer with C nodes.



→ Deep CNNs may suffer from vanishing or exploding grads.

Idea: add normalization layers to the CNN architecture.

Batch normalization (BN):

ensures the distribution of activations has mean = 0 & var = 1 when averaged over a minibatch.

$$\bar{\mu}_B = \frac{1}{|B|} \sum_{\vec{z} \in B} \vec{z}, \quad \sigma_B^2 = \frac{1}{|B|} \sum_{\vec{z} \in B} (\vec{z} - \bar{\mu}_B)^2$$

vector of activations in a layer

B is the minibatch of size $|B|$ containing example n
(datapoint)

$$\vec{z}_n' = \frac{\vec{z}_n - \bar{\mu}_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad \begin{matrix} \text{standardized} \\ \epsilon > 0 = \text{const} \end{matrix}$$

$$\vec{z}_n'' = \gamma \vec{z}_n' + \beta \quad \begin{matrix} \text{shifted \& scaled} \\ \text{element-wise} \quad \left[\begin{matrix} \text{output of the BN} \\ \text{layer} \end{matrix} \right] \end{matrix}$$

γ & β are fitted prms.

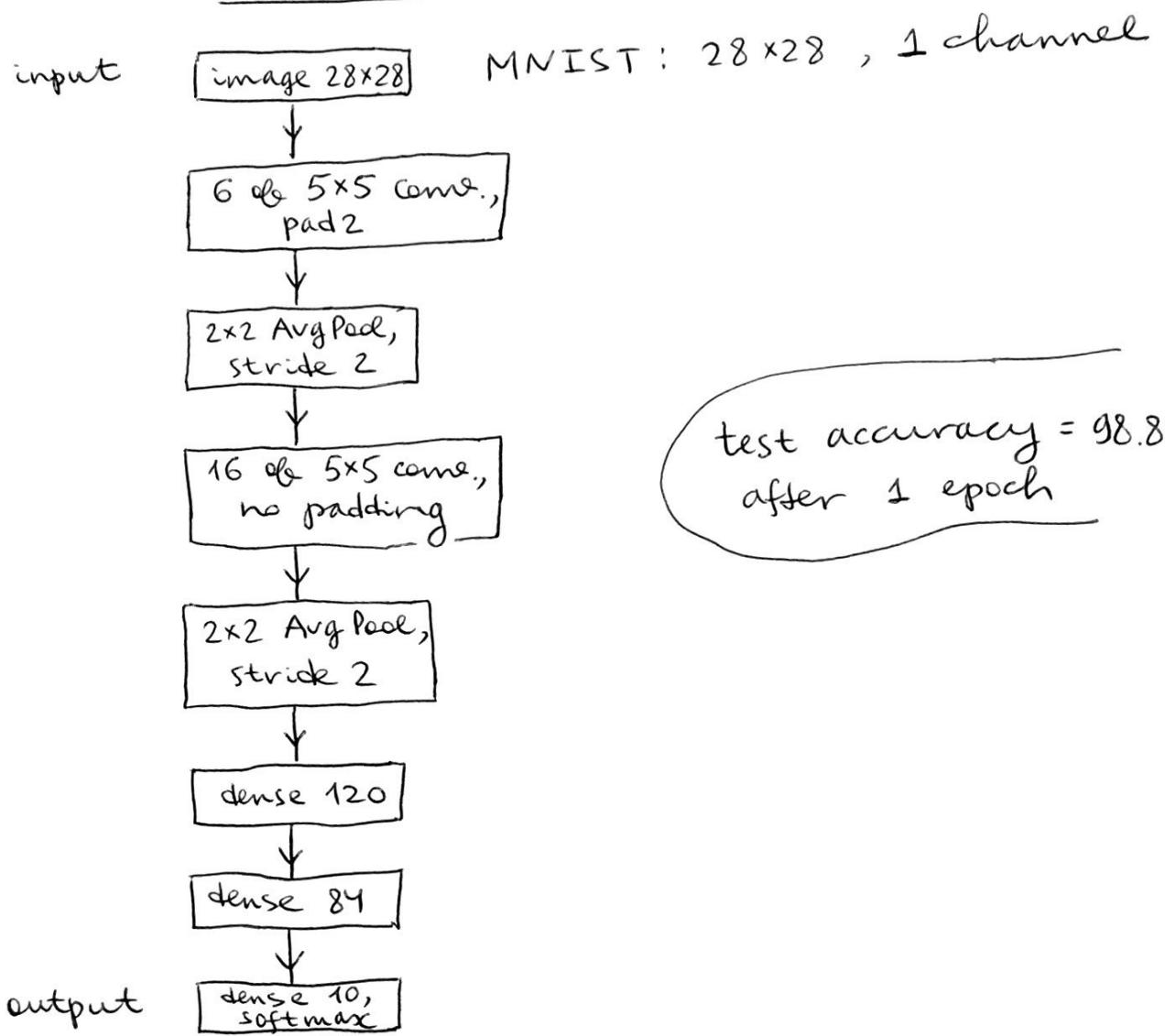
In the context of CNN, $\bar{\mu}_B$ is averaged across both spatial locations and examples $\Rightarrow \dim\{\bar{\mu}_B\} = \# \text{ channels}$.

At test time, first compute $\bar{\mu}_e$ & σ_e^2
on the entire training set layer index
and then 'freeze' these prms for testing.

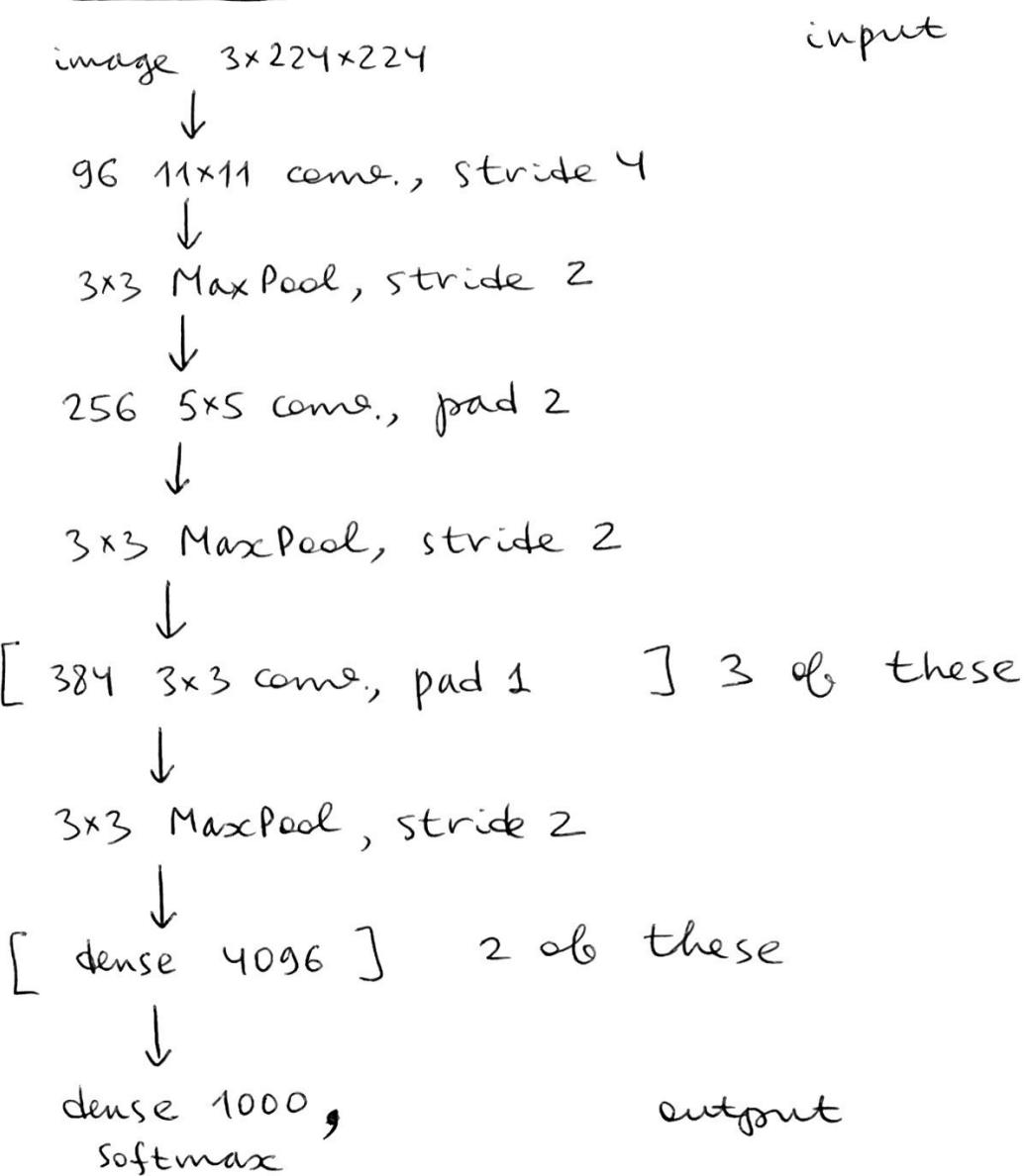
adaptive grad clipping: instead of BN layers, use grad clipping with $C = \text{fitting pm}$.

CNN architectures

LeNet5



AlexNet



~ 60 M prms (mostly due to dense layers)

ImageNet : ~ 1 M labeled examples

Top 5 error rate was reduced
from $\sim 26\%$ (SOTA) to $\sim 15\%$.