

COSI 165B

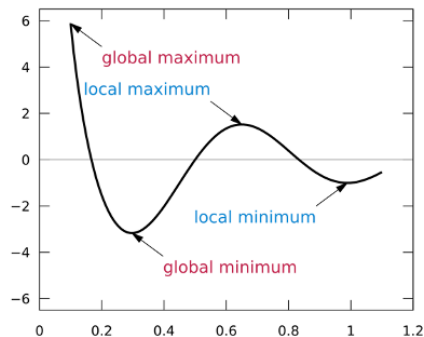
Deep Learning

Chuxu Zhang
Computer Science Department
Brandeis University

3/1/2021

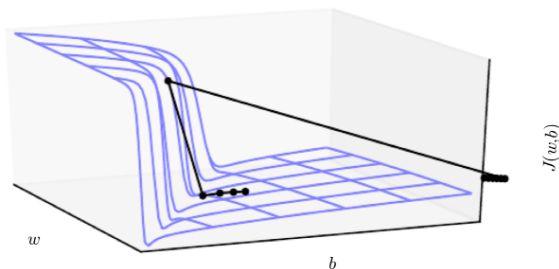
□ Optimization

Local Minima



non-convex function
(neural networks)

Exploding Gradients



cliff region (catapult parameters very far)
gradient clipping

Algorithm 8.1 Stochastic gradient descent (SGD) update at training iteration k

Require: Learning rate ϵ_k .

Require: Initial parameter θ

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Apply update: $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$

end while

Learning rate: it is necessary to gradually decrease the learning rate over time

$$\epsilon_k = (1 - \alpha)\epsilon_0 + \alpha\epsilon_{\tau} \quad \text{linear decay}$$

$$\alpha = \frac{k}{\tau}.$$

□ Optimization

Algorithm 8.2 Stochastic gradient descent (SGD) with momentum

Require: Learning rate ϵ , momentum parameter α .

Require: Initial parameter θ , initial velocity v .

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{x^{(1)}, \dots, x^{(m)}\}$ with corresponding targets $y^{(i)}$.

 Compute gradient estimate: $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$

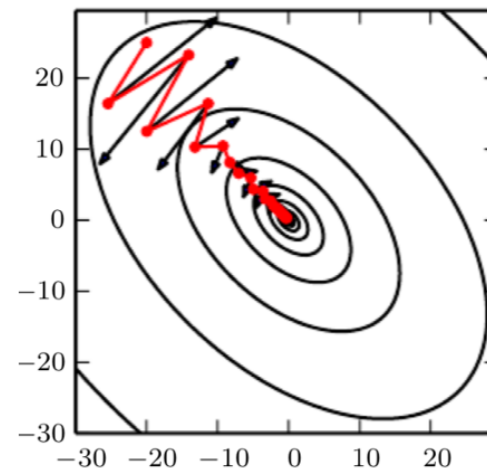
 Compute velocity update: $v \leftarrow \alpha v - \epsilon g$

 Apply update: $\theta \leftarrow \theta + v$

end while

$$v \leftarrow \alpha v - \epsilon \nabla_{\theta} \left(\frac{1}{m} \sum_{i=1}^m L(f(x^{(i)}; \theta), y^{(i)}) \right),$$

$$\theta \leftarrow \theta + v.$$



The velocity v accumulates the gradient elements

The larger α is, the more previous gradients affect the current direction

□ Optimization

Algorithm 8.4 The AdaGrad algorithm

Require: Global learning rate ϵ

Require: Initial parameter θ

Require: Small constant δ , perhaps 10^{-7} , for numerical stability

Initialize gradient accumulation variable $\mathbf{r} = \mathbf{0}$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Accumulate squared gradient: $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$

 Compute update: $\Delta \theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$. (Division and square root applied element-wise)

 Apply update: $\theta \leftarrow \theta + \Delta \theta$

end while

Algorithm 8.7 The Adam algorithm

Require: Step size ϵ (Suggested default: 0.001)

Require: Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1)$. (Suggested defaults: 0.9 and 0.999 respectively)

Require: Small constant δ used for numerical stabilization. (Suggested default: 10^{-8})

Require: Initial parameters θ

Initialize 1st and 2nd moment variables $\mathbf{s} = \mathbf{0}$, $\mathbf{r} = \mathbf{0}$

Initialize time step $t = 0$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient: $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

$t \leftarrow t + 1$

 Update biased first moment estimate: $\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$

 Update biased second moment estimate: $\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$

 Correct bias in first moment: $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$

 Correct bias in second moment: $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$

 Compute update: $\Delta \theta = -\epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \delta}}$ (operations applied element-wise)

 Apply update: $\theta \leftarrow \theta + \Delta \theta$

end while

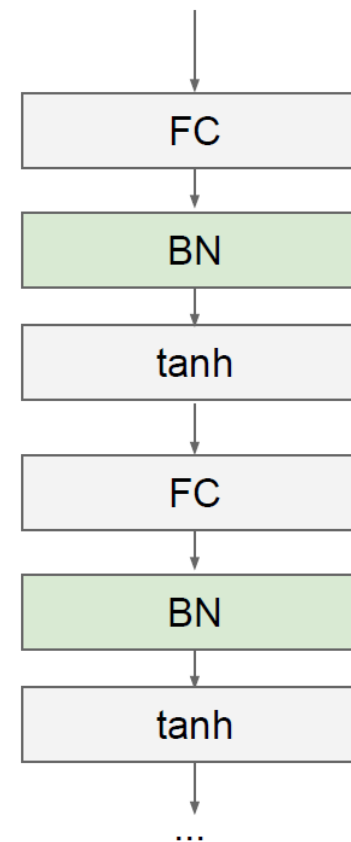
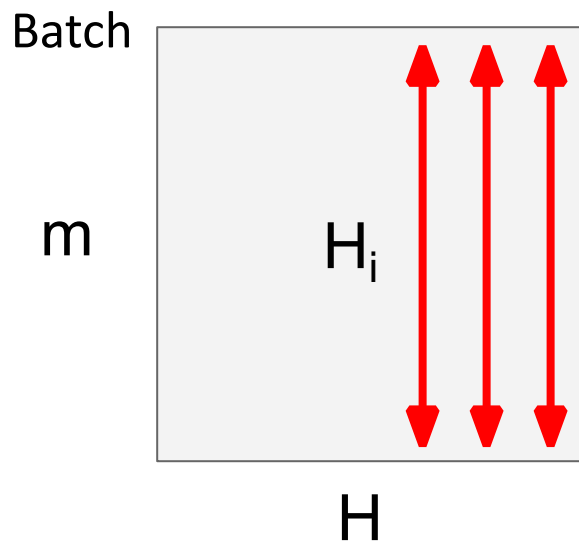
□ Optimization

Batch Normalization

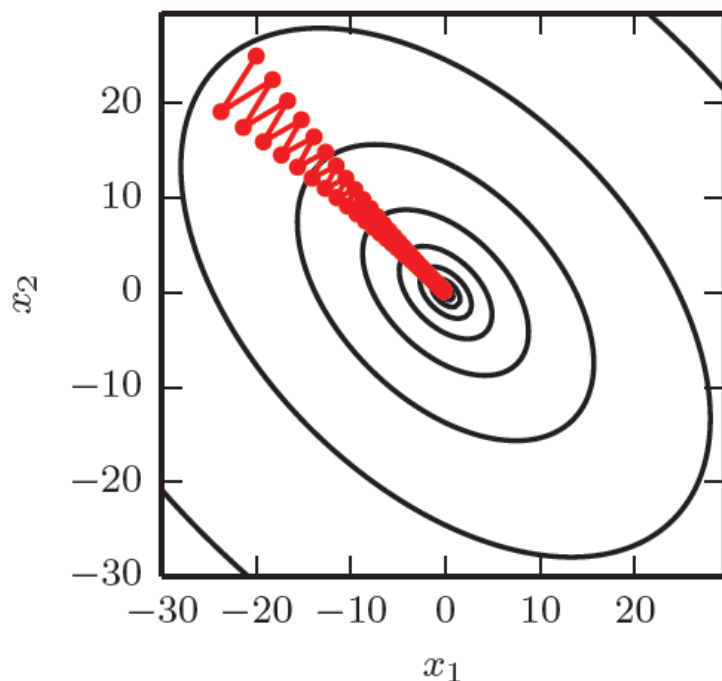
$$H' = \frac{H - \mu}{\sigma}$$

$$\mu = \frac{1}{m} \sum_i H_{i,:}$$

$$\sigma = \sqrt{\delta + \frac{1}{m} \sum_i (H - \mu)_i^2}$$



□ Optimization



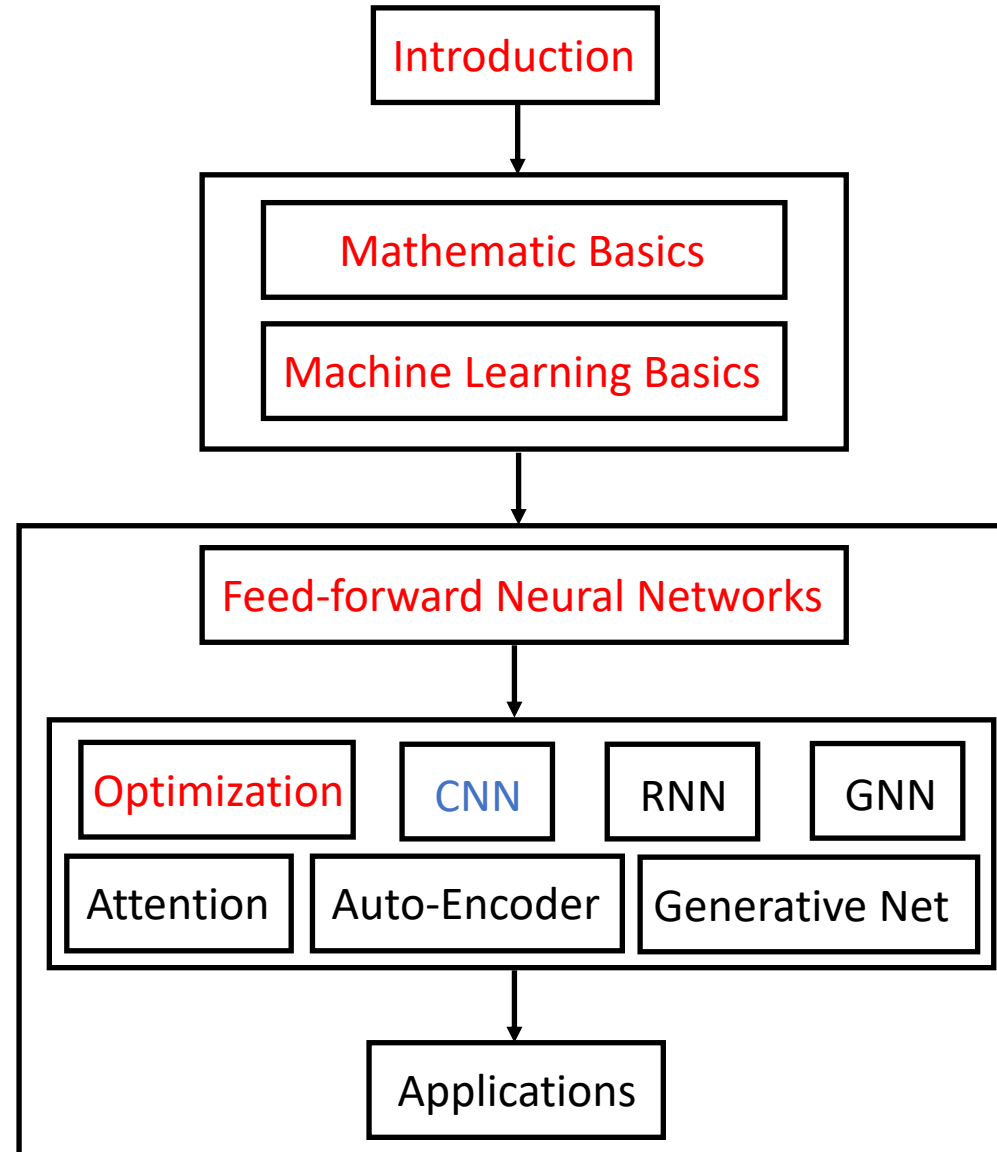
Ill-condition (Hessian matrix): the ratio of the largest singular value to the smallest singular value is large.

The loss function have contours that are very long ellipsoids.

The gradient directions away from the center are nearly orthogonal to the useful direction of descent.

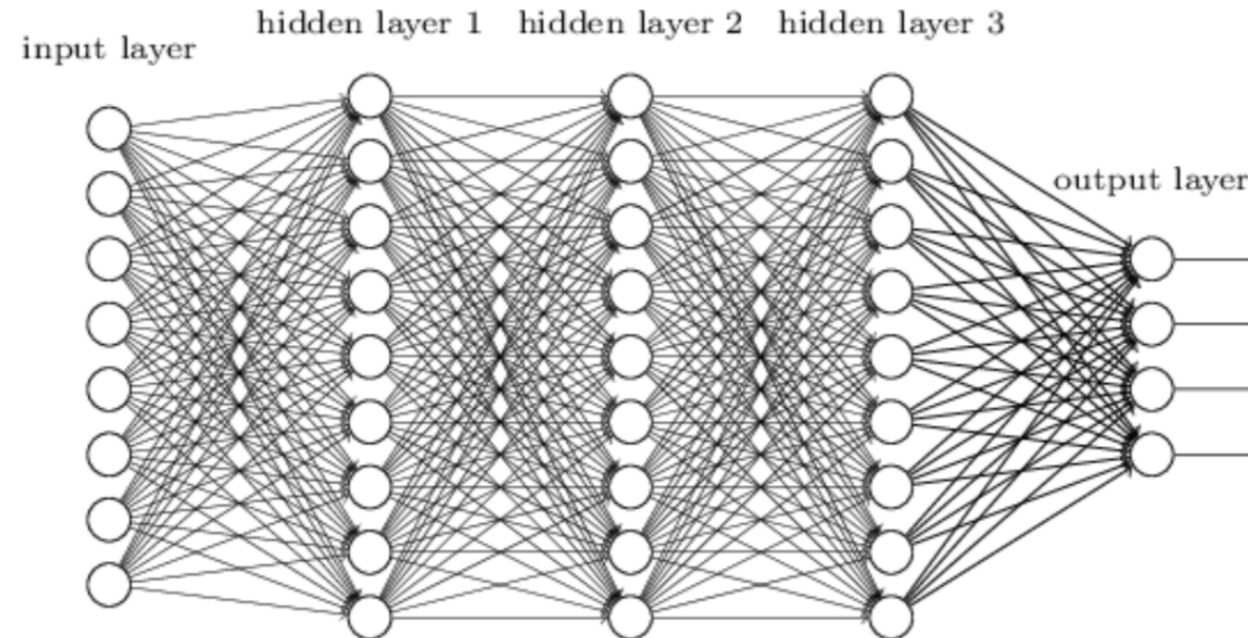
Zig-zag path: slow learning.

Structure of This Course



❑ Feed-forward neural network for image classification

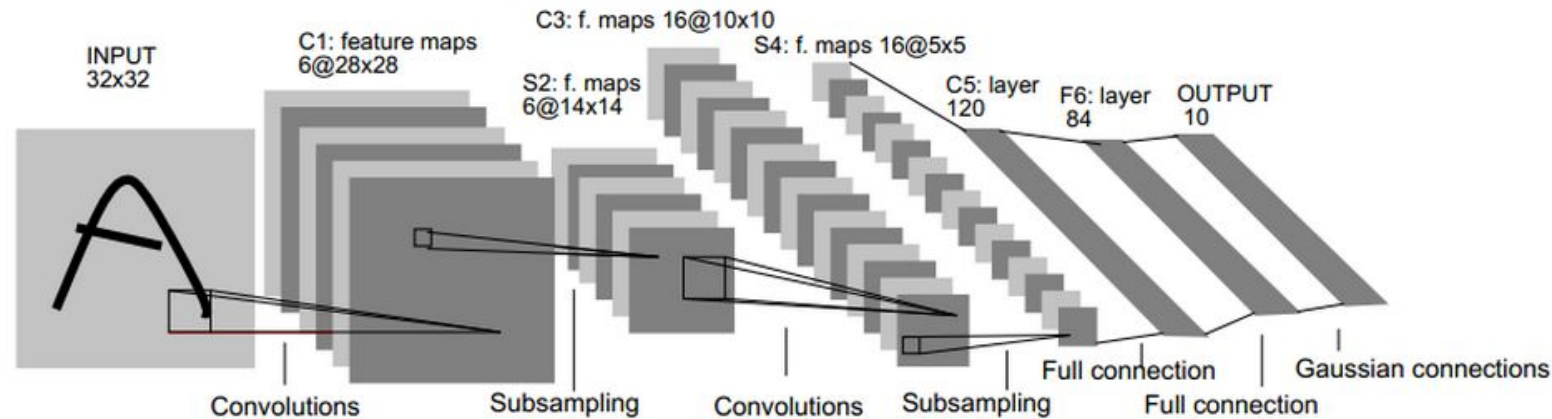
Input: vector of
image pixel values



No spatial information!
Many parameters!

How can we use **spatial structure** in the input to inform the architecture of the network?

□ What is CNN

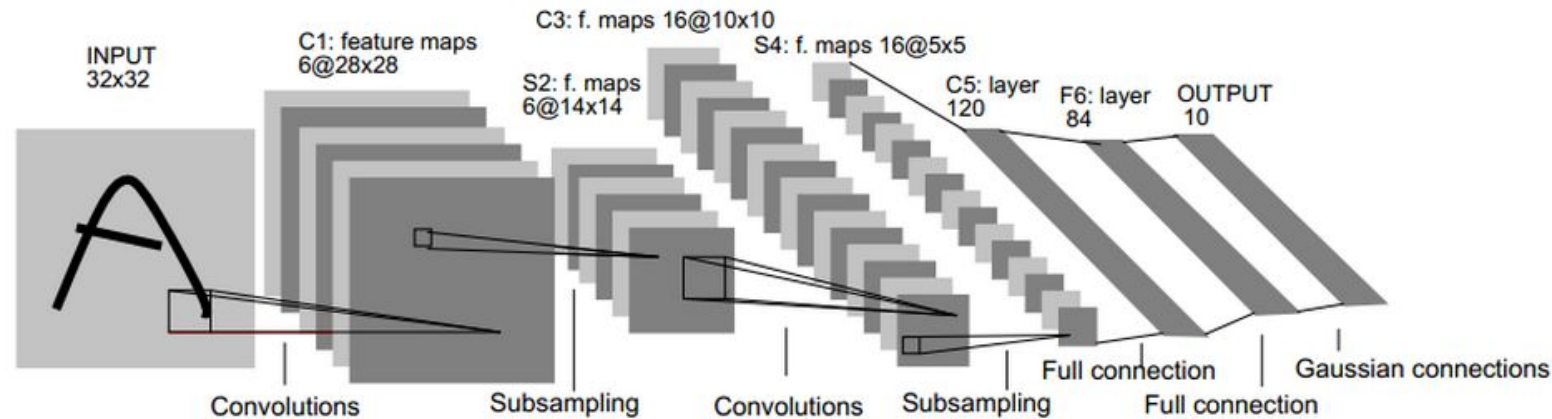


CNN is a specialized kind of neural network for processing data that has grid-like topology, such as image data which can be thought of as a 2D grid of pixels.

Convolutional networks have been tremendously successful in practical applications.

The name “convolutional neural network” indicates that the network employs a mathematical operation called convolution. Convolution is a specialized kind of linear operation.

□ What is CNN



Convolutional Layer: convolution operation

Pooling Layer: max/mean pooling

Fully Connected Layer: feed-forward neural network

□ Convolution

Example: tracking the location of a spaceship with a laser sensor

$$s(t) = \int x(a)w(t-a)da$$

↑
weighting
function

↓
single
measurement

$$\text{feature map} \leftarrow s(t) = (x * w)(t)$$

↑
kernel

↓
input

provide measurements at every
instant in time is not realistic

discrete convolution

$$s(t) = (x * w)(t) = \sum_{a=-\infty}^{\infty} x(a)w(t-a)$$

multidimensional array input/parameters
(tensors)

□ Convolution

1D Convolution

$$s(t) = (x * w)(t) = \sum_{a=-\infty}^{\infty} x(a)w(t-a)$$

2D Convolution

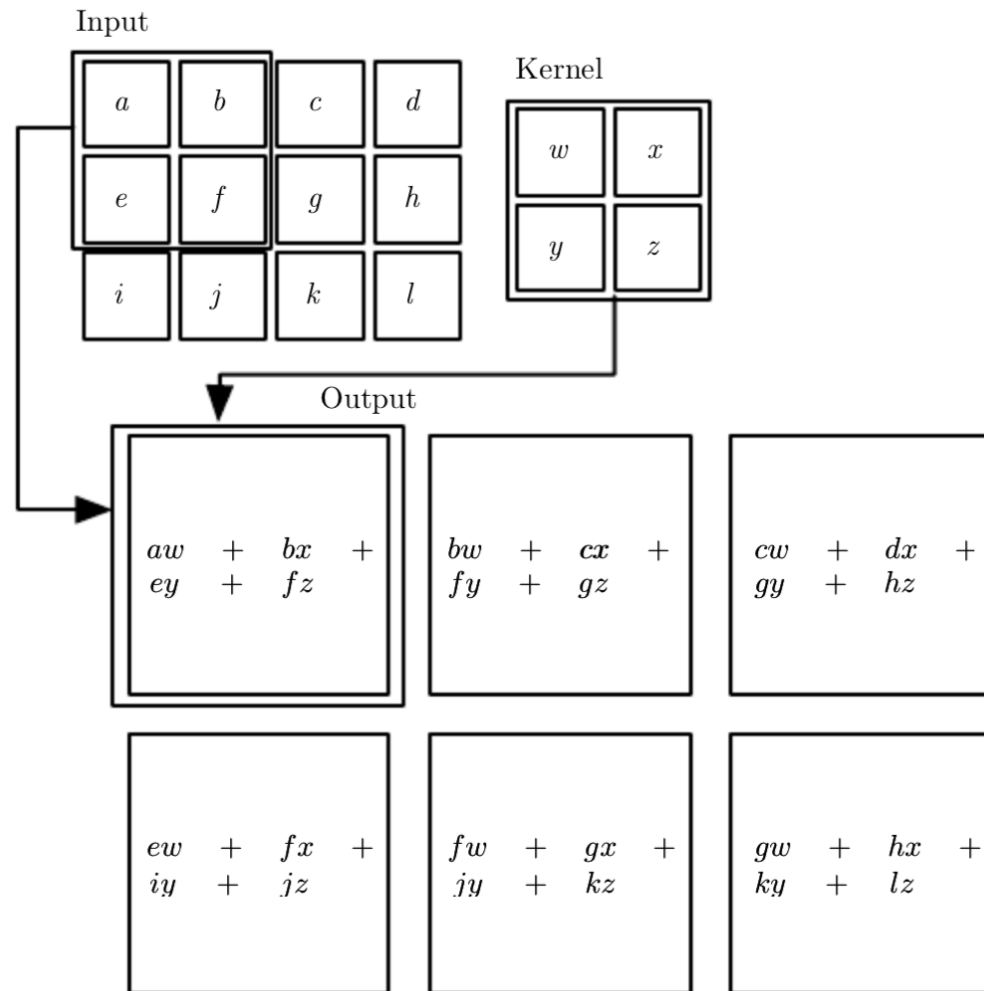
$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(m, n)K(i-m, j-n)$$

Convolution is commutative, meaning we can equivalently write

$$S(i, j) = (K * I)(i, j) = \sum_m \sum_n I(i-m, j-n)K(m, n)$$

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(i+m, j+n)K(m, n)$$

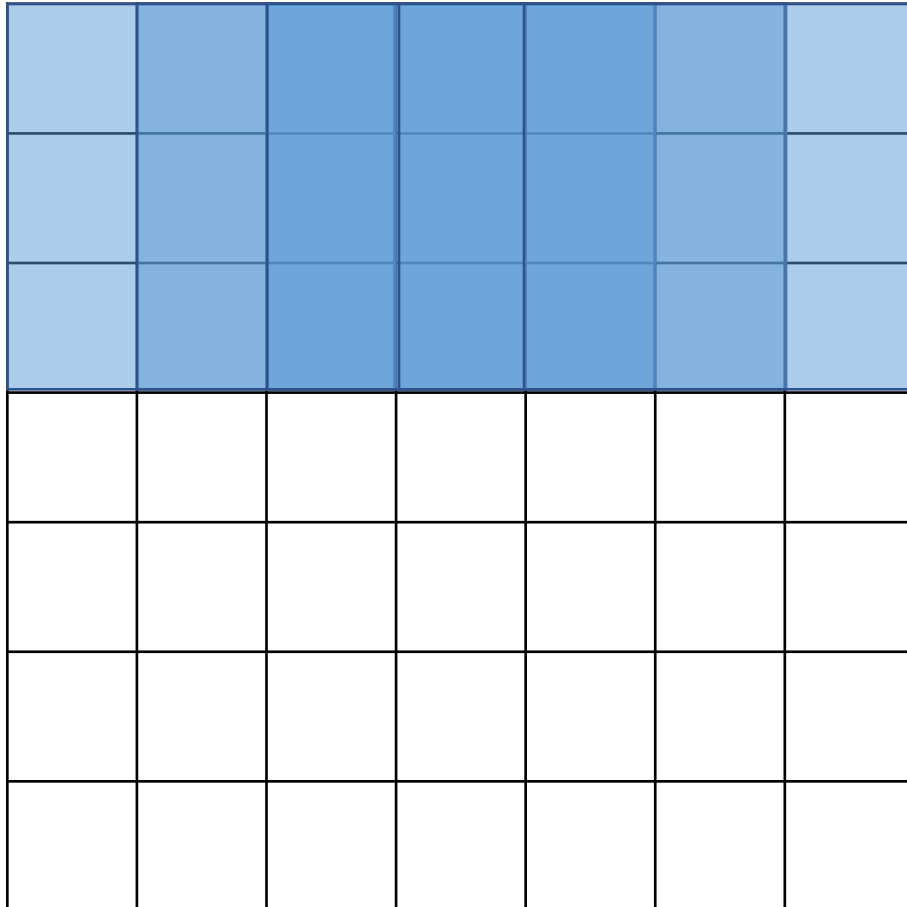
□ Convolution: example



$$S(i, j) = (K * I)(i, j) = \sum_m \sum_n I(i - m, j - n) K(m, n)$$

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(i + m, j + n) K(m, n)$$

□ Convolution: more examples



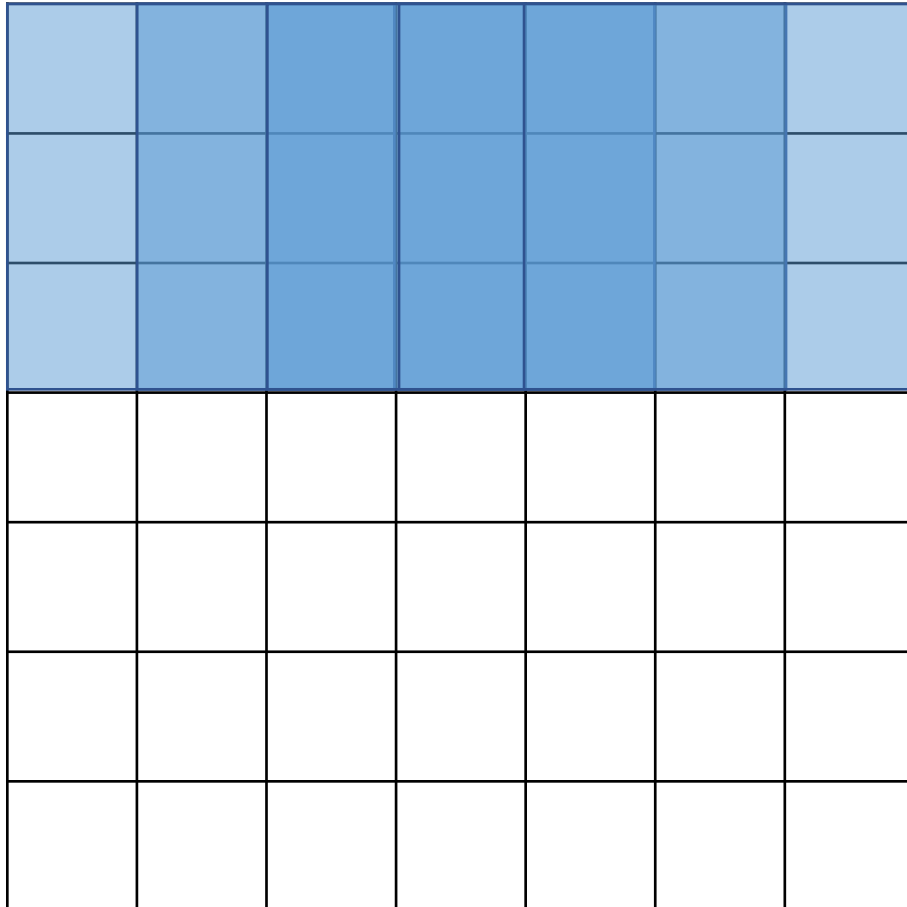
Input: 7x7

Kernel: 3x3

Stride: 1x1

Output: 5x5

□ Convolution: more examples



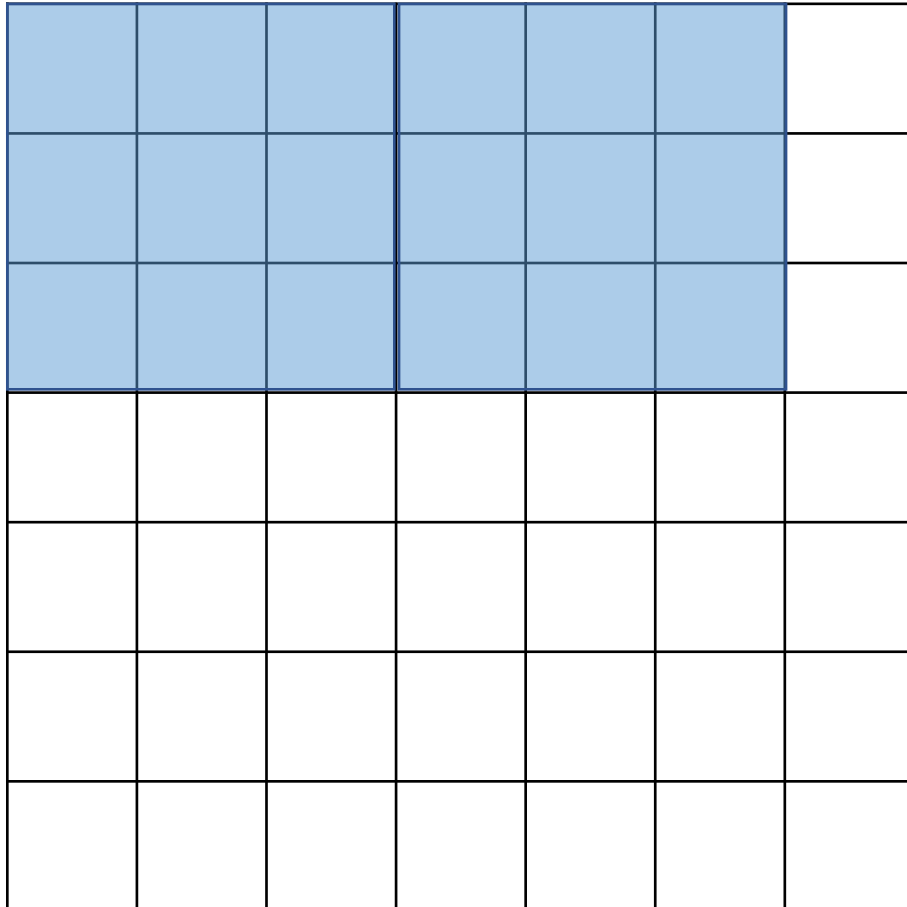
Input: 7x7

Kernel: 3x3

Stride: 2x2

Output: 3x3

□ Convolution: more examples



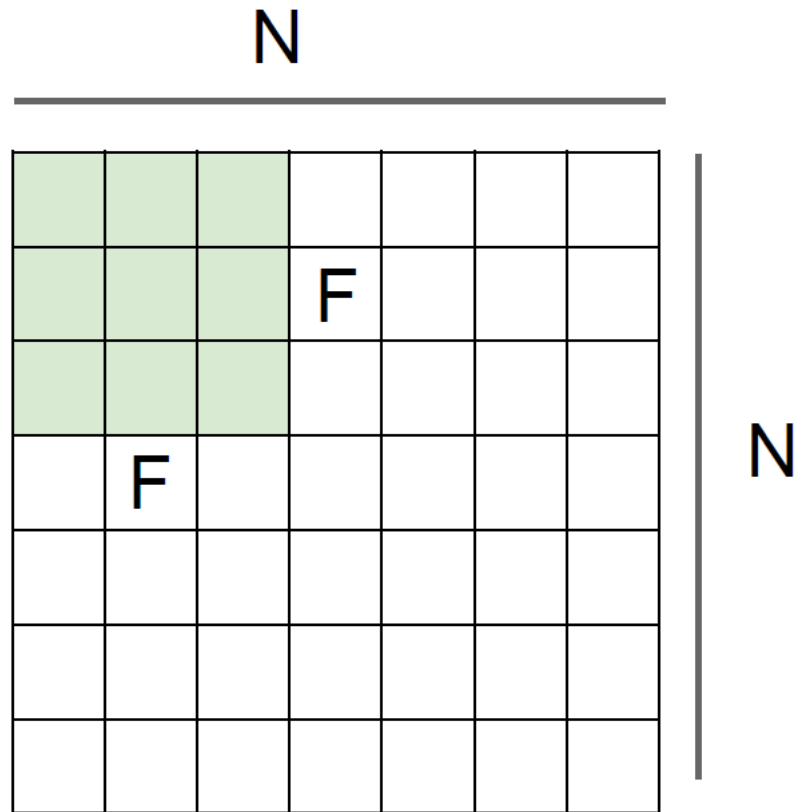
Input: 7x7

Kernel: 3x3

Stride: 3x3

Output: invalid dimension

□ Convolution: output size



$$\text{Output Size} = (N - F) / S + 1$$

$$N = 7, F = 3, S = 1, \text{Output} = 5$$

$$N = 7, F = 3, S = 2, \text{Output} = 3$$

$$N = 7, F = 3, S = 3, \text{Output} = 2.3$$

How to fix?

□ Convolution: zero padding

0	0	0	0	0	0			
0								
0								
0								
0								

Input Size = $(N-F)/S + 1$

$N = N + 2$ (zero padding at boundary)

$N = 9, F = 3, S = 3, \text{Output} = 3$

□ Convolution: zero padding

0	0	0	0	0	0			
0								
0								
0								
0								

How many zero at boundary

$$Z = [(M-1)S + F - N]/2$$

When $S = 1$, $F-S/2$

$$N = 7, F = 3, S = 3, M = 3, Z = 1$$

□ Motivation/Strengths

Sparse Interactions

MLP: fully connected, every output unit interacts with every input unit. $O(m \times n)$

Convolution: making the kernel smaller than the input, sparse connection $O(k \times n)$

Parameter Sharing

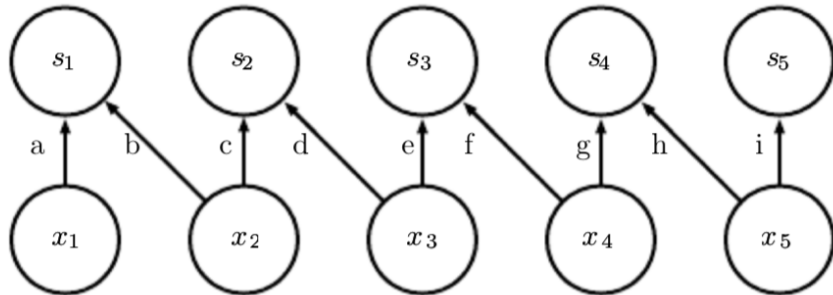
MLP: each element of the weight matrix is used exactly once when computing the output of a layer

Convolution: each member of the kernel is used at every position of the input, learn only one set of parameters for every position

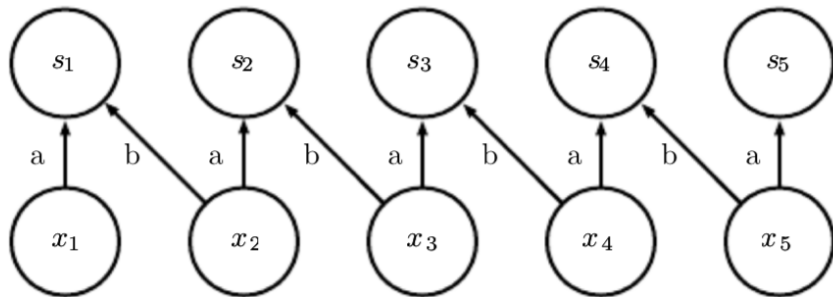
Equivariant Representations

If the input changes, the output changes in the same way: shift-convolution, convolution-shift

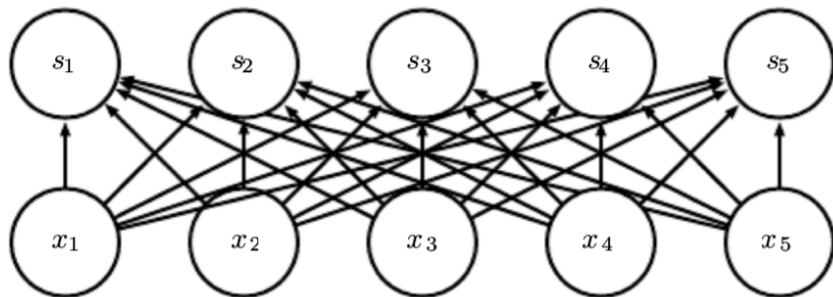
Comparison



Local Connection: patch size = 2, different parameters



Convolution: patch size = 2, same parameters



Fully Connection: path size = 5, different parameters

□ Pooling

Use a pooling function to modify the output of the layer.

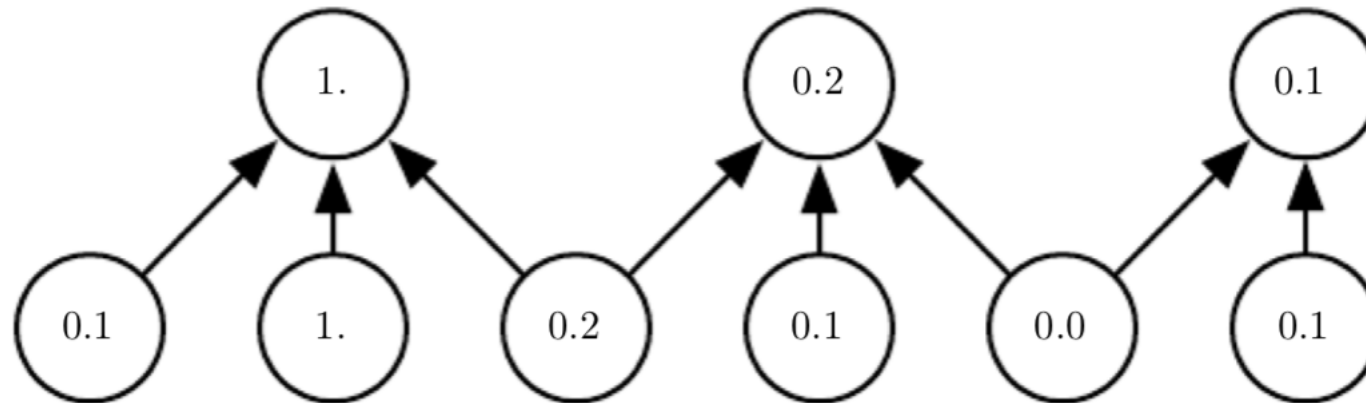
A pooling function replaces the output of the net at a certain location with a summary statistic of the nearby outputs.

e.g., max pooling: reports the maximum output within a rectangular neighborhood

In all cases, pooling helps to make the representation become approximately invariant to small translations of the input. Invariance to translation means that if we translate the input by a small amount, the values of most of the pooled outputs do not change.

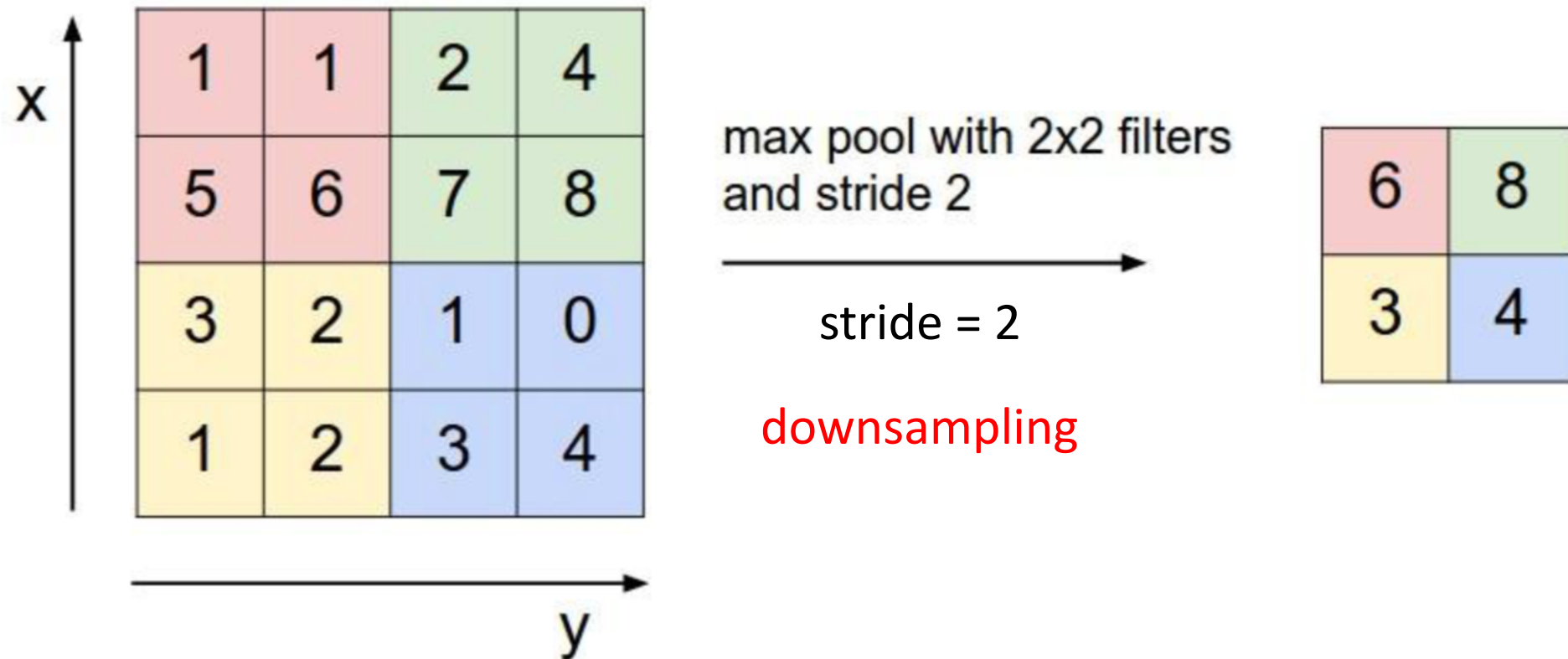
□ Pooling: example

Report summary statistics for pooling regions spaced k pixels apart rather than 1 pixel apart. This improves the computational efficiency of the network because the next layer has roughly k times fewer inputs to process.

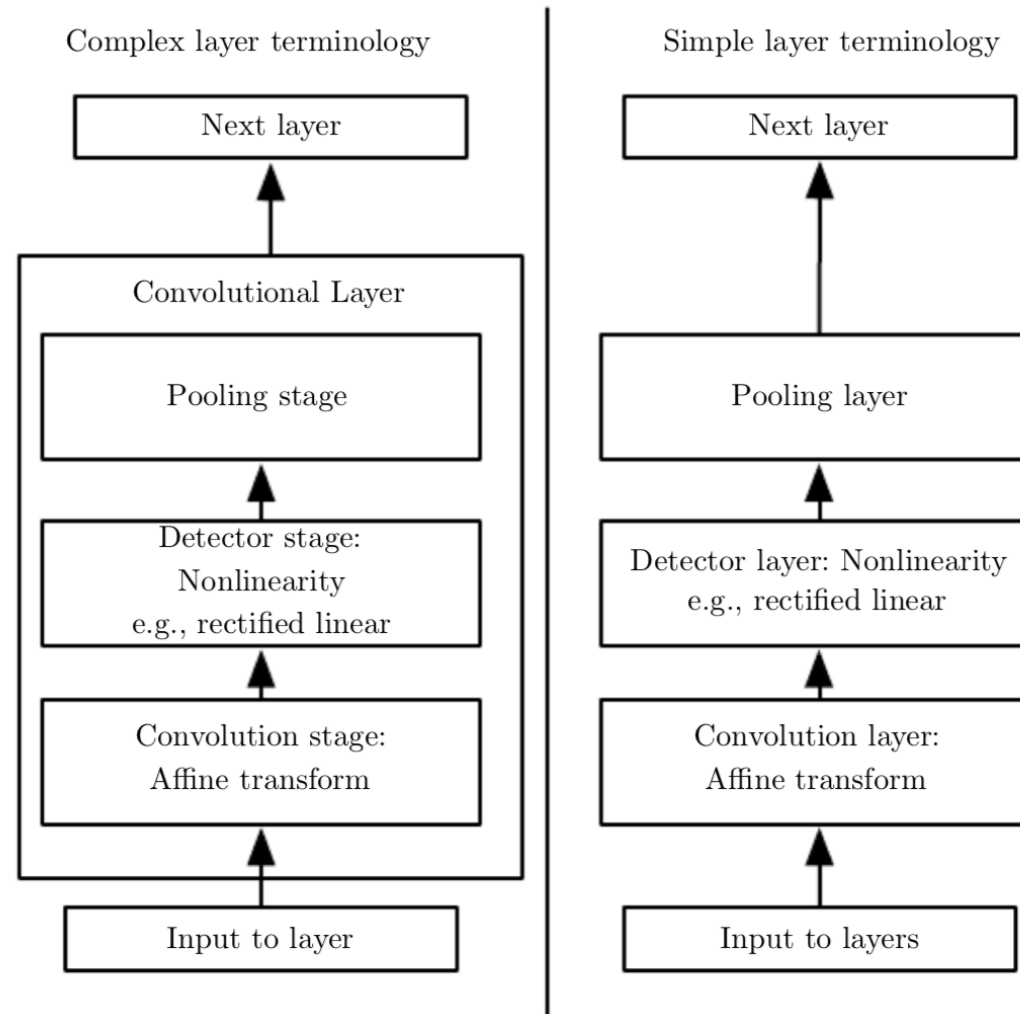


Pooling: max-pooling with a pool width of three and a stride between pools of two.

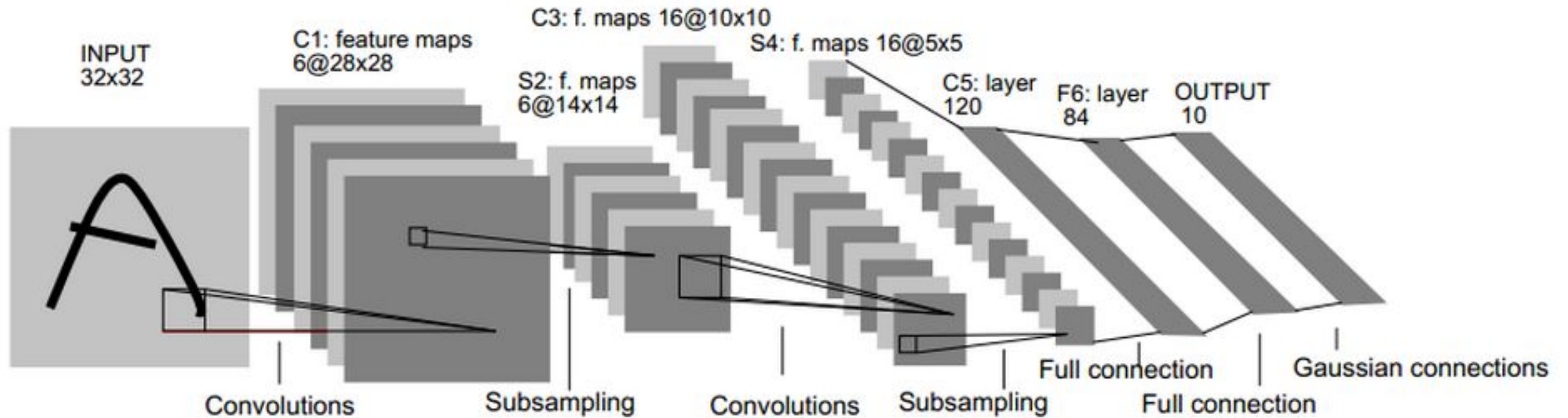
□ Pooling: 2D example



□ The components of a convolutional layer



□ A Full Model



Q & A