

Lecture 12: Convolutional Neural Networks

Introduction to Machine Learning [25737]

Sajjad Amini

Sharif University of Technology

- 1 Approach Definition
- 2 Common Layers
- 3 All Together

Except explicitly cited, the reference for the material in slides is:

- Murphy, K. P. (2022). *Probabilistic machine learning: an introduction*. MIT press.

Section 1

Approach Definition

MLPs

Assume transformation $\mathbf{z} = \varphi(\mathbf{W}\mathbf{x})$. Then j -th element in \mathbf{z} can be represented as:

$$z_j = \varphi(\mathbf{w}_j^T \mathbf{x})$$

We can interpret this equation as computing the similarity between \mathbf{x} and \mathbf{w}_j . When working with 2D images, this structure can lead to severe problems.

Approach Definition

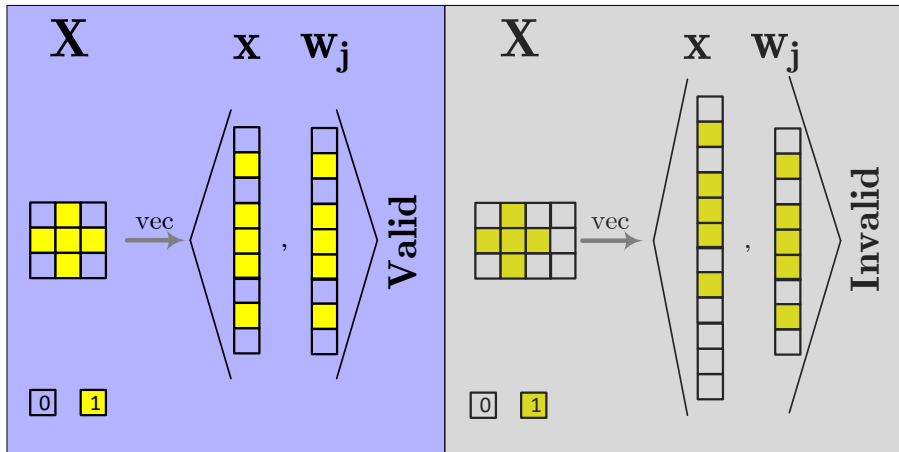


Figure: Not applicable when image size is changed

Approach Definition

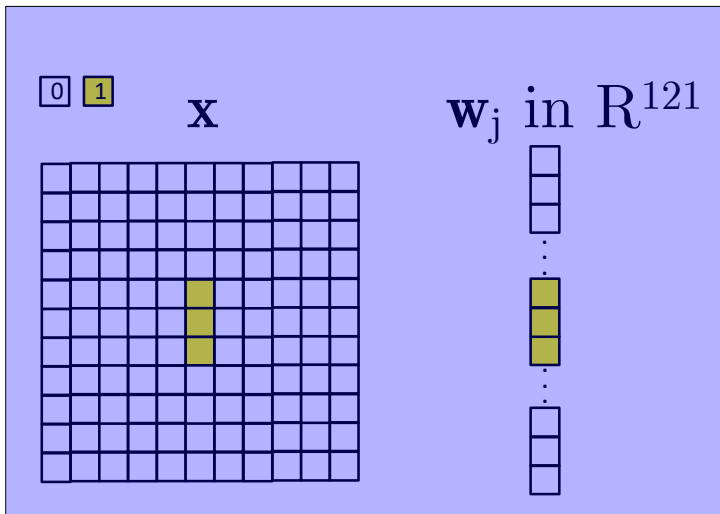


Figure: Highly redundant

Approach Definition

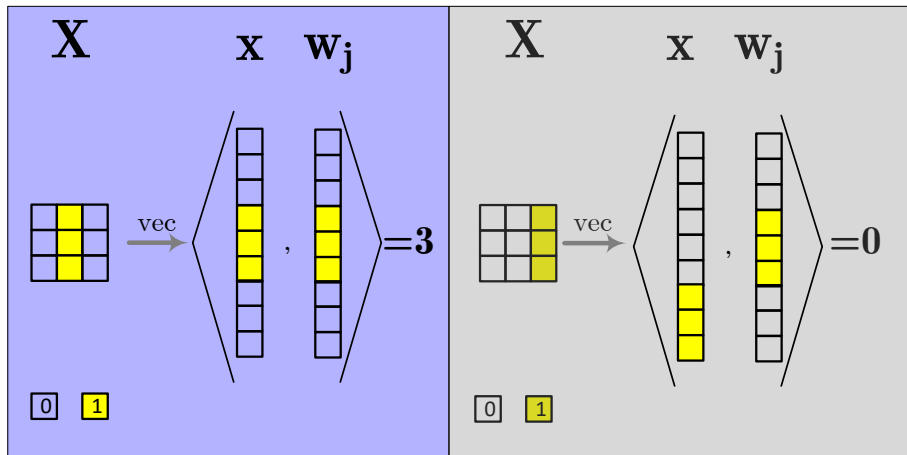


Figure: Not exhibiting translation invariance

Convolutional Neural Networks (CNN)

To solve the challenges mentioned above, CNNs are introduced where matrix multiplication is replaced with convolution operator.

- We can compute the convolution of different size images with the same filter.
- The size of convolution filter is smaller than the size of input features.
- Convolution is a template matching operator and can present translation invariance.

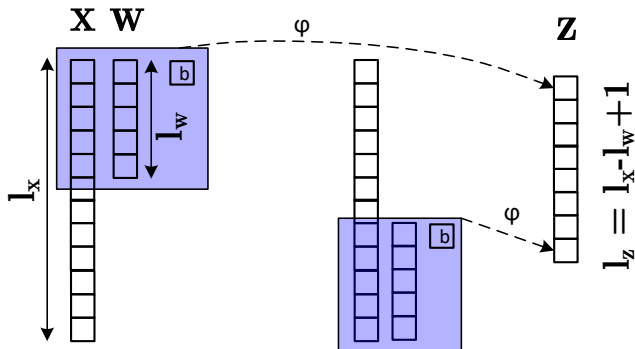
Section 2

Common Layers

1D Convolution Layer

1D convolution (valid)

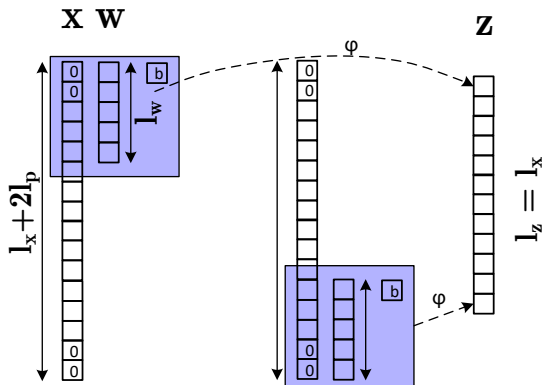
$$\begin{cases} \mathbf{x} \in \mathbb{R}^{l_x} \\ \mathbf{w} \in \mathbb{R}^{l_w} \end{cases} \Rightarrow z_p = \varphi \left(b + \sum_{j=0}^{K-1} w_j x_{p+j} \right), \quad 0 \leq p \leq l_x - l_w$$



1D Convolution Layer

1D convolution (same)

- If we pad each side of input feature vector $\mathbf{x} \in \mathbb{R}^{l_x}$ with p elements and $\mathbf{w} \in \mathbb{R}^{l_w}$, then the output size will be $l_x + 2p - l_w + 1$.
- If we select $p = \frac{l_w - 1}{2}$, then input and output sizes are the same.



1D Convolution Layer

Convolution to Matrix Multiplication

Assume we have the following convolution operator:

$$z = \varphi(\mathbf{x} \circledast \mathbf{w} + b), \quad \begin{cases} \mathbf{x} \in \mathbb{R}^{l_x} \\ \mathbf{w} \in \mathbb{R}^3 \\ \mathbf{z} \in \mathbb{R}^{l_x-2} \end{cases}$$

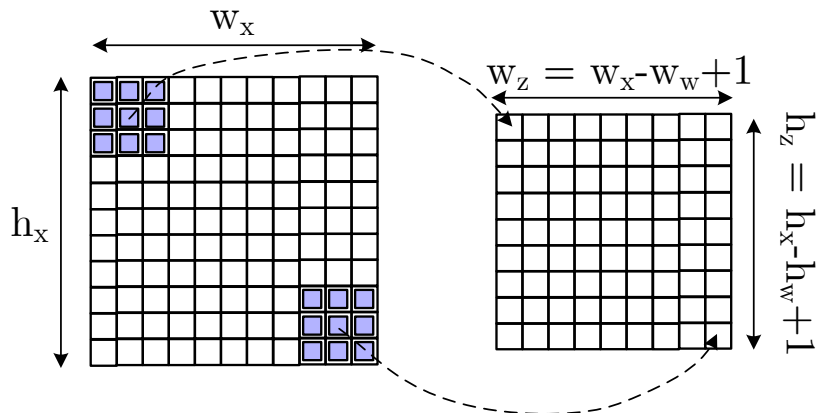
then we can write above mapping in matrix multiplication as:

$$\begin{bmatrix} z_0 \\ z_1 \\ \vdots \\ z_{l_x-4} \\ z_{l_x-3} \end{bmatrix} = \varphi \left(\begin{bmatrix} w_0 & w_1 & w_2 & 0 & \dots & 0 & 0 & 0 & 0 \\ 0 & w_0 & w_1 & w_2 & \dots & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \dots & w_0 & w_1 & w_2 & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 & w_0 & w_1 & w_2 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_{l_x-2} \\ x_{l_x-1} \end{bmatrix} + \begin{bmatrix} b \\ b \\ \vdots \\ b \\ b \end{bmatrix} \right)$$

2D Convolution Layer

2D convolution (valid)

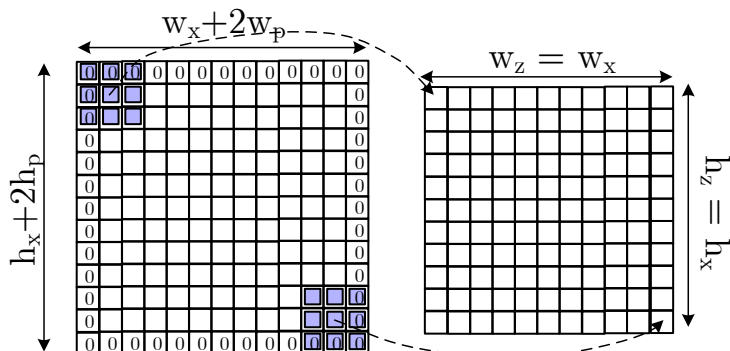
$$\begin{cases} \mathbf{X} \in \mathbb{R}^{h_x \times w_x} \\ \mathbf{W} \in \mathbb{R}^{h_w \times w_w} \end{cases} \Rightarrow z_{pq} = \varphi \left(b + \sum_{i=0}^{h_w-1} \sum_{j=0}^{w_w-1} w_{ij} x_{(p+i)(q+j)} \right), \quad \begin{cases} 0 \leq p \leq h_x - h_w \\ 0 \leq q \leq w_x - w_w \end{cases}$$



2D Convolution Layer

2D convolution (same)

- If we pad each side of input feature matrix $\mathbf{X} \in \mathbb{R}^{h_x \times w_x}$ with h_p and w_p elements and $\mathbf{w} \in \mathbb{R}^{h_w \times w_w}$, then the output size will be $(h_x + 2h_p - h_w + 1) \times (w_x + 2w_p - w_w + 1)$.
- If we select $h_p = \frac{h_w - 1}{2}$ and $w_p = \frac{w_w - 1}{2}$, then input and output sizes are the same.



2D Convolution Layer

Convolution to Matrix Multiplication

Assume we have the following convolution operator:

$$\mathbf{Z} = \varphi(\mathbf{X} \circledast \mathbf{W} + b), \quad \begin{cases} \mathbf{X} \in \mathbb{R}^{3 \times 3} \\ \mathbf{W} \in \mathbb{R}^{2 \times 2} \\ \mathbf{Z} \in \mathbb{R}^{2 \times 2} \end{cases}$$

Assume column-wise vectorizing. Then we can write above mapping in matrix multiplication as:

$$\begin{bmatrix} z_{00} \\ z_{10} \\ z_{01} \\ z_{11} \end{bmatrix} = \varphi \left(\begin{bmatrix} w_{00} & w_{10} & 0 & w_{01} & w_{11} & 0 & 0 & 0 & 0 \\ 0 & w_{00} & w_{10} & 0 & w_{01} & w_{11} & 0 & 0 & 0 \\ 0 & 0 & 0 & w_{00} & w_{10} & 0 & w_{01} & w_{11} & 0 \\ 0 & 0 & 0 & 0 & w_{00} & w_{10} & 0 & w_{01} & w_{11} \end{bmatrix} \begin{bmatrix} x_{00} \\ x_{10} \\ x_{20} \\ x_{01} \\ x_{11} \\ x_{21} \\ x_{02} \\ x_{12} \\ x_{22} \end{bmatrix} + \begin{bmatrix} b \\ b \\ b \\ b \end{bmatrix} \right)$$

Receptive Field

For each element in the output features, its receptive field are the elements in the input features that form the it. In the above example, $\{x_{00}, x_{10}, x_{01}, x_{11}\}$ are receptive field of z_{00} .

Redundancy in the Output Features

As the receptive field for neighboring features in the output of convolution layer are highly overlapped, there exist redundancy in the output features. Strided convolution is designed to eliminate this redundancy.

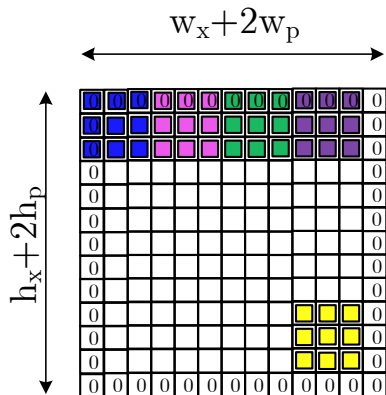
Strided Convolution

Strided Convolution is ordinary convolution while we skip every s_h and s_w in vertical and horizontal shift.

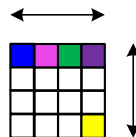
Output Dimensions

$$\left\{ \begin{array}{l} \mathbf{X} \in \mathbb{R}^{h_x \times w_x} \\ \mathbf{W} \in \mathbb{R}^{h_w \times w_w} \\ \text{Height Padding : } h_p \\ \text{Width Padding : } w_p \\ \text{Height Stride : } h_s \\ \text{Width Stride : } w_s \end{array} \right. \Rightarrow \dim(\mathbf{Z}) : \left\lfloor \frac{h_x + 2h_p - h_w + h_s}{h_s} \right\rfloor \times \left\lfloor \frac{w_x + 2w_p - w_w}{w_s} \right\rfloor$$

Strided Convolution



$$W_z = \text{floor}((w_x + 2w_p - w_w + w_s) / w_s)$$



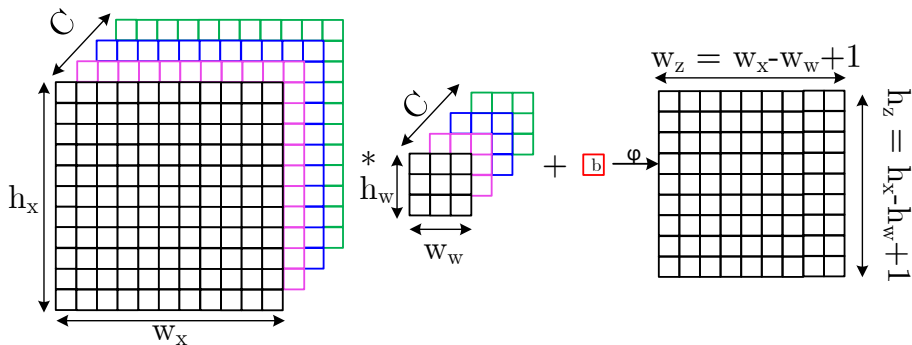
$$h_z = \text{floor}((h_x + 2h_p - h_w + h_s) / h_s)$$

Figure: 2D convolution (same)

Extension 1: Multiple Input Channel

Convolution with multi-channel input

$$\begin{aligned} \begin{cases} \mathbf{X} \in \mathbb{R}^{h_x \times w_x \times C} \\ \mathbf{W} \in \mathbb{R}^{h_w \times w_w \times C} \end{cases} &\Rightarrow z_{pq} = \varphi \left(b + \sum_{k=0}^{C-1} \sum_{i=0}^{h_w-1} \sum_{j=0}^{w_w-1} w_{ijk} x_{(p+i)(q+j)k} \right), \quad \begin{cases} 0 \leq p \leq h_x - h_w \\ 0 \leq q \leq w_x - w_w \end{cases} \\ &\Rightarrow \mathbf{Z} \in \mathbb{R}^{(h_x - h_w + 1) \times (w_x - w_w + 1)} \end{aligned}$$



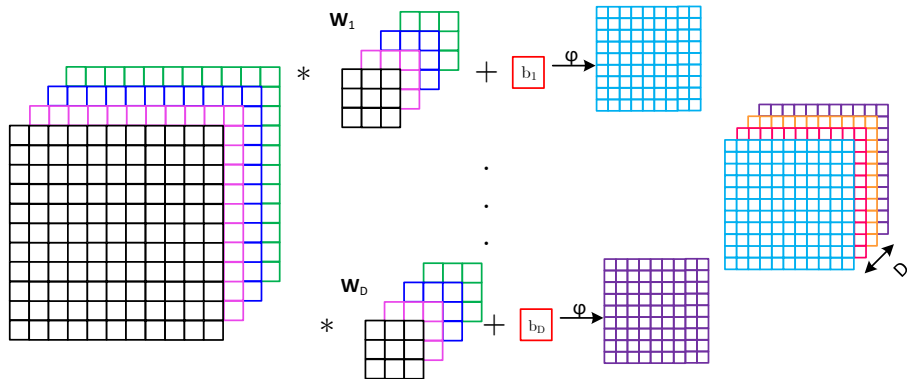
Extension 2: Multiple Output Channel

Convolution with multi-channel output

$$\begin{cases} \mathbf{X} \in \mathbb{R}^{h_x \times w_x \times C} \\ \{(\mathbf{W}_d, b_d)\}_{d=0}^{D-1} \\ \mathbf{W}_d \in \mathbb{R}^{h_w \times w_w \times C} \\ b_d \in \mathbb{R} \end{cases} \Rightarrow \begin{cases} z_{pq0} = \varphi \left(b_0 + \sum_{k=0}^{C-1} \sum_{i=0}^{h_w-1} \sum_{j=0}^{w_w-1} w_{ijk0} x_{(p+i)(q+j)k} \right) \\ \vdots \\ z_{pq(D-1)} = \varphi \left(b_{D-1} + \sum_{k=0}^{C-1} \sum_{i=0}^{h_w-1} \sum_{j=0}^{w_w-1} w_{ijk(D-1)} x_{(p+i)(q+j)k} \right) \\ 0 \leq p \leq h_x - h_w \\ 0 \leq q \leq w_x - w_w \end{cases}$$

We can concatenate matrices $\{\mathbf{Z}_0, \dots, \mathbf{Z}_{D-1}\}$ which results in $\mathbf{Z} \in \mathbb{R}^{(h_x - h_w + 1) \times (w_x - w_w + 1) \times D}$

Convolution with multi-channel output



Convolutional Layer

All Together

Assume:

$$\mathbf{X} \in \mathbb{R}^{h_x \times w_x \times C}$$

Input feature tensor

$$\mathbf{W} \in \mathbb{R}^{h_w \times w_w \times C \times D}$$

Weight tensor

$$\mathbf{b} \in \mathbb{R}^D$$

Bias vector

$$h_p, w_p$$

height and width of padding, respectively

$$h_s, w_s$$

height and width of stride, respectively

Then the output $\mathbf{Z} = \mathbf{W} \circledast \mathbf{X} + \mathbf{b}$ is of the following dimensions:

$$\underbrace{\left\lfloor \frac{h_x + 2h_p - h_w + h_s}{h_s} \right\rfloor}_{h_z} \times \underbrace{\left\lfloor \frac{w_x + 2w_p - w_w + w_s}{w_s} \right\rfloor}_{w_z} \times D$$

and

$$z_{pqd} = \varphi \left(b_c + \sum_{c=0}^{C-1} \sum_{i=0}^{h_w-1} \sum_{j=0}^{w_w-1} w_{ijkc} \hat{x}_{(h_s \times p+i)(w_s \times q+j)c} \right), \begin{cases} 0 \leq p \leq h_z - 1 \\ 0 \leq q \leq w_z - 1 \\ 0 \leq d \leq D - 1 \end{cases}$$

Pooling Layers

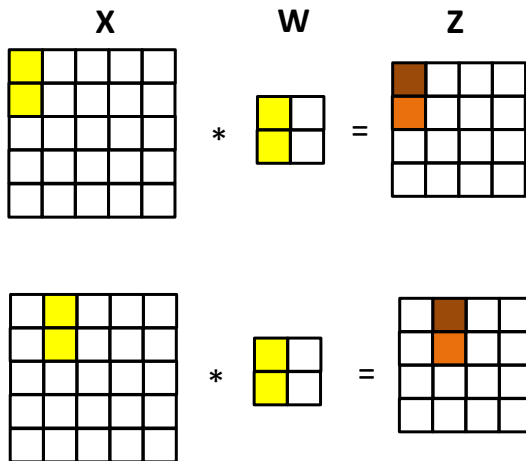


Figure: Output tensor carry information about the location

Pooling Layers

Assume:

$$\mathbf{X} \in \mathbb{R}^{h_x \times w_x \times C}$$

$$h_f, w_f$$

$$h_s, w_s$$

$$p(\cdot)$$

Input feature tensor

height and width of pooling, respectively

height and width of stride, respectively

Pooling operator

Then the output of pooling layer is:

$$z_{pqc} = p(x_{(h_s \times p:h_s \times p+h_f)(w_s \times q:w_s \times q+w_f)(c)})$$

Pooling Layers

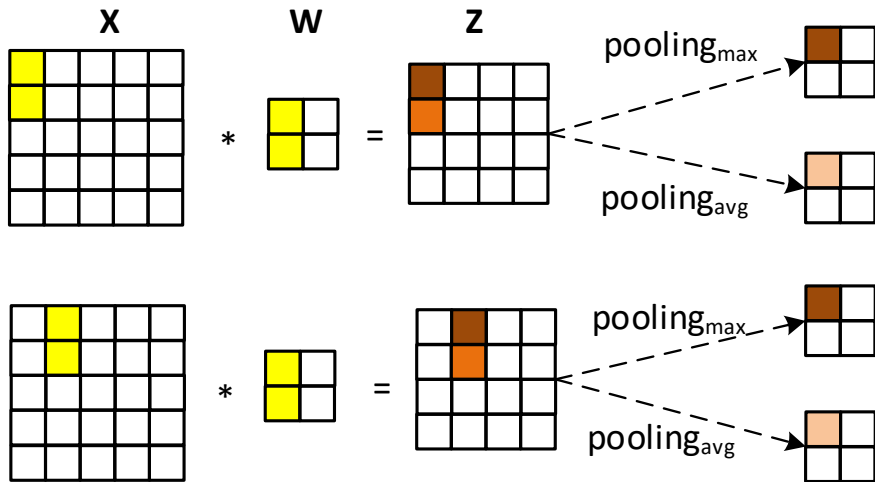


Figure: Pooling layer provide local invariance

Flattening Layer

Frame Title

Flattening layers are used to reshape the input feature tensor $\mathbf{X} \in \mathbb{R}^{h_x \times w_x \times C}$ into output feature vector $\mathbf{z} \in \mathbb{R}^{(h_x \times w_x \times C)}$ using vectorizing operator.

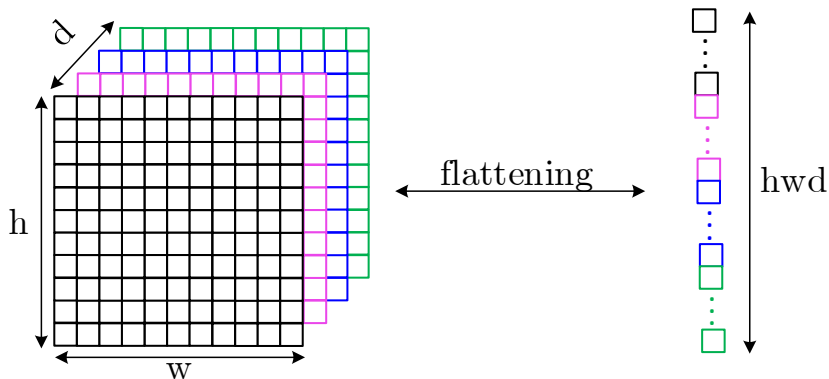


Figure: Flattening Layer

Section 3

All Together

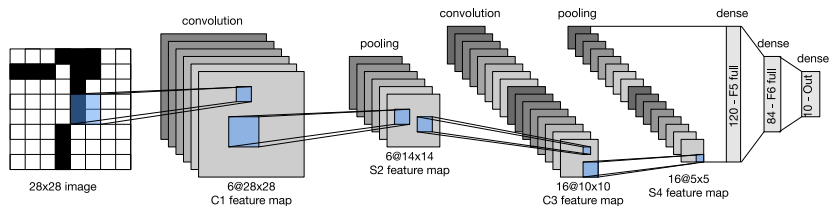


Figure: LeNet5 for MNIST classification (Test accuracy: 98.8% after 1 epoch)

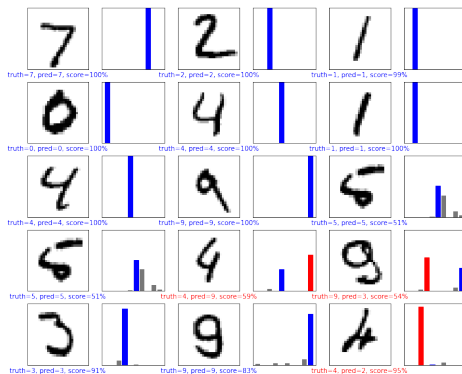


Figure: Result of LeNet5 for MNIST classification