

# CNN for intelligent control

Wen Yu

# NN for image Classification



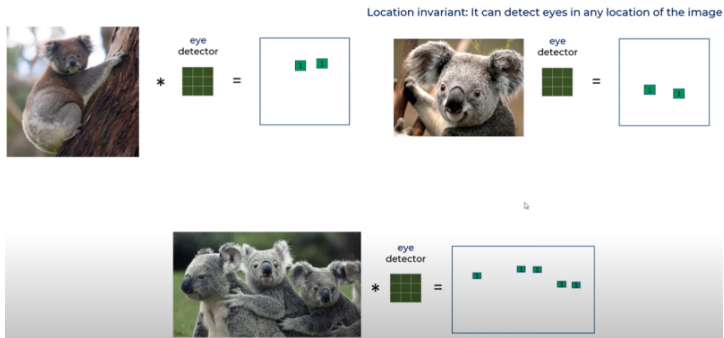
Image size =  $1920 \times 1080 \times 3$

First layer neurons =  $1920 \times 1080 \times 3 \sim 6$  million

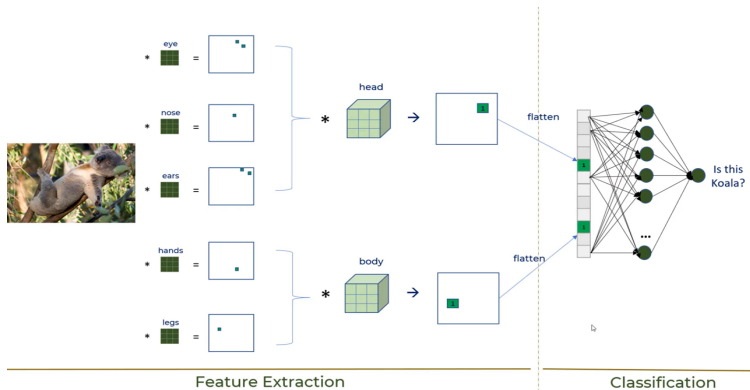
Hidden layer neurons = Let's say you keep it  $\sim 4$  million

Weights between input and hidden layer =  $6 \text{ mil} \times 4 \text{ mil}$   
=  $24$  million

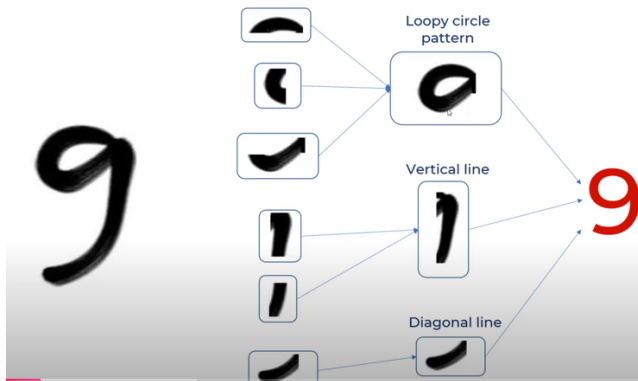
# Feature extraction



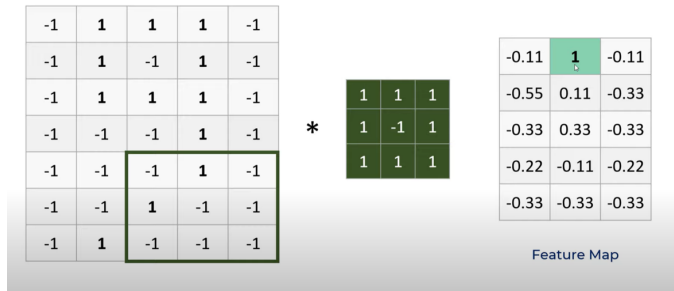
# Classification



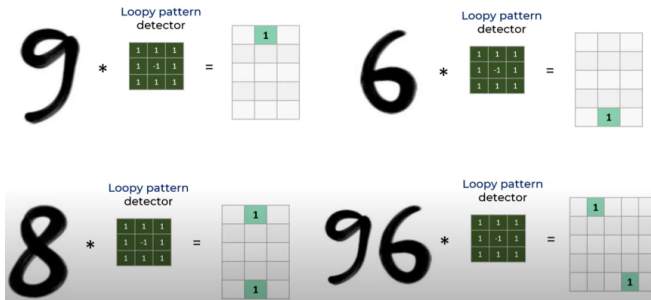
# Features



# Convolution operation

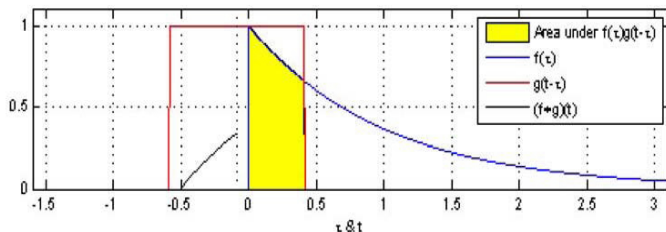


# Feature extraction



# Convolution in discrete-time

$$f * g = \sum_{k=-\infty}^{\infty} f(k) g(n-k)$$





# Convolution operation

- Convolution produces the area amount of two overlap functions
- multiplication of a original function and its translation.

$$h(t) = f * g = \int_{-\infty}^{\infty} f(\tau) g(t - \tau) d\tau$$

- The output of a linear time-invariant system (LTI) is the convolution of the input signal and the system (the impulse response) gives
- In probability theory, the probability distribution of the **sum** of two independent random variables is the convolution of their individual distributions.

The distribution of the sum  $Z = X + Y$  of two independent discrete variables is

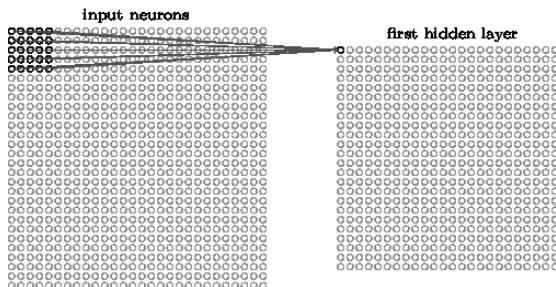
$$P(Z = z) = \sum P(X = k) P(y = z - k)$$
$$h(z) = \int_{-\infty}^{\infty} f(\tau) g(z - \tau) d\tau$$

where  $f(t)$  and  $g(t)$  are density functions

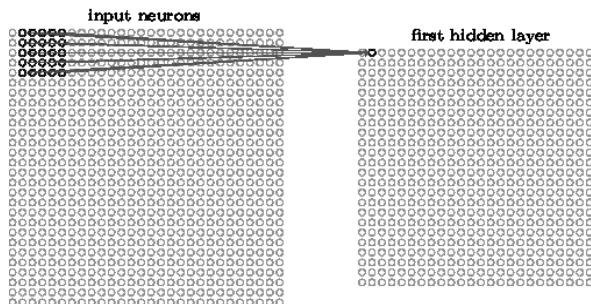
# Convolution: Local receptive fields

Every input pixel are not connected to every hidden neuron. Only localized regions of the input image is sent to hidden layers.

A image with  $28 \times 28$  pixels, only  $5 \times 5$  region is activated, this region is called the "local receptive field" for the hidden neuron.



# Convolution



Hidden layer has  $28 - 5 + 1 = 24$ ,  $24 \times 24$  pixels, when step length is 1.

We can also move the local receptive field 2 pixels (step is 2)

- The  $5 \times 5$  region is called: "filter", or neuron, or kernel
- the region which the filter passed called "receptive field"
- The values of the elements of the filter called "weights"
- $24 \times 24$  pixels are called "activation map", or feature map

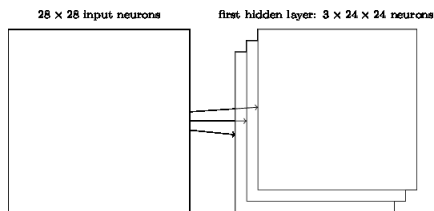
# Convolution

The values of hidden neurons are

$$y = b + \sum_{i=1}^5 \sum_{j=1}^5 w_{i,j} x_{i,j}$$

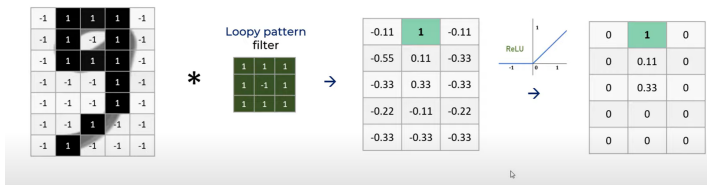
where  $b$  is the shared value for the bias,  $w_{i,j}$  is a  $5 \times 5$  array of shared weights,  $x_{i,j}$  denotes the input activation at position  $i, j$

The convolution is regarded as local features extraction. If we need more than one feature map, we need several different filters to convolute the image.



Training parameter:

# Activation function: ReLU



# Activation function: ReLU

$$y = \phi \left( b + \sum_{i=1}^5 \sum_{j=1}^5 w_{i,j} x_{i,j} \right)$$

where  $\phi$  is the neural activation function - perhaps the sigmoid function or tanh.

Rectified linear unit (ReLU)

$$y = \max \left[ 0, \left( b + \sum_{i=1}^5 \sum_{j=1}^5 w_{i,j} x_{i,j} \right) \right]$$

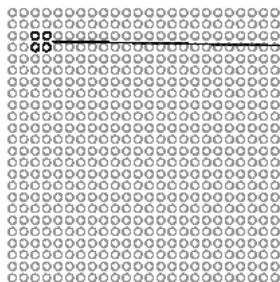
This activation function enables better training of deeper networks, the most popular activation function for deep neural networks.

Convolutional networks are well adapted to the translation invariance of images:

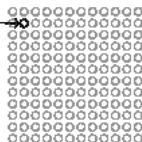
- All the neurons in the first hidden layer detect exactly the same feature
- The weights and bias correspond a particular local receptive field,

Simplify the information from the convolutional layer. It is a condensed feature map

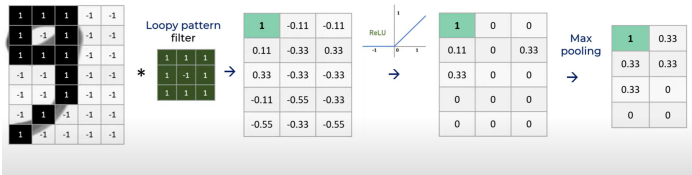
hidden neurons (output from feature map)



max-pooling units



Shifted 9 at  
different position



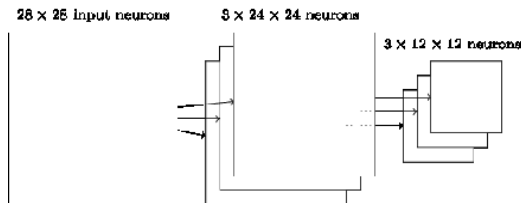


Each unit in the pooling layer summarizes a region of  $2 \times 2$  in the convolutional layer. Many fewer pooled features, and so this helps reduce the number of parameters needed in later layers.  $24 \times 24 \rightarrow 12 \times 12$

Operation of the "summarization"

- max-pooling: the maximum of the  $2 \times 2$  input region: where a given feature is found, then throws away the exact positional information. Once a feature has been found, its exact location isn't as important
- $L_2$  pooling: square root of the sum of the  $2 \times 2$  input region: is a way of condensing information from the convolutional layer.
- average:  $\frac{x_j}{\max(x_j)}$

# Fully-connected layer

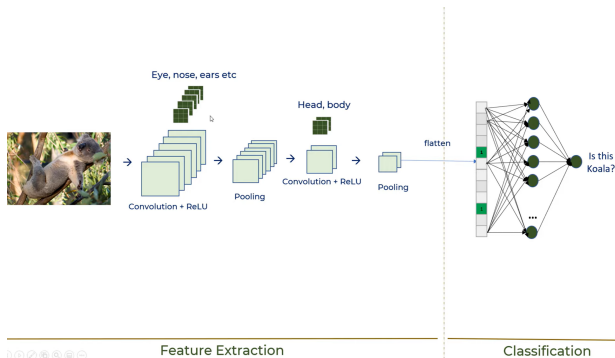


If we want to classify 10 digits, they are 10 output neurons.  
This layer connects every neuron from pooling layer to every one of the 10 output neurons.

There are  $12 \times 12 \times 3 \times 10 = 4,320$  weights

# CNN-structure

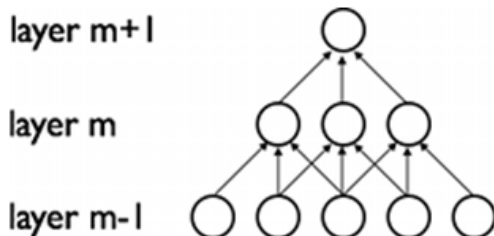
[Convolution ReLU → Pooling (Sub-sampling) → Convolution ReLU → Pooling (Sub-sampling) → Full connection]



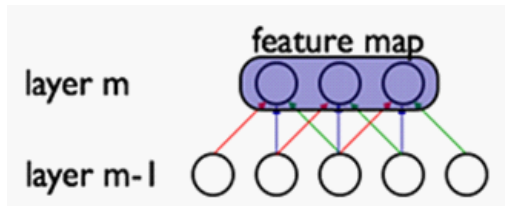
# Properties of CNN

- 1 local receptive fields
- 2 shared weights
- 3 pooling

# Properties of CNN: sparse Connectivity - local receptive fields



# Properties of CNN: Shared weights



# Why does the convolutional neural network have higher accuracy

Machine learning is largely about two things

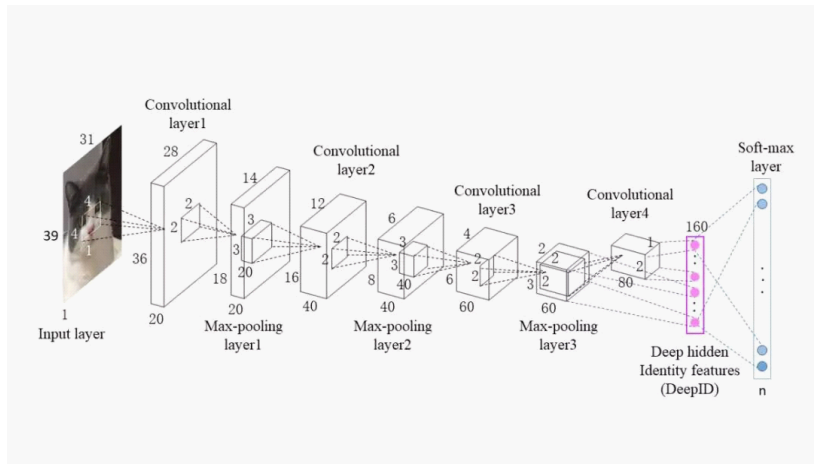
- Optimization: How you fit your model to the data
- Trade-offs: Precision/recall;  
Bias/variance; Powerful/interpretable; Complexity/training speed/data needed

# Advantage of CNN

- CNNs utilize the spatial information, weight sharing to reduce number of parameters
- Each convolutional kernel scans the entire input feature map to produce an output map: a feature is the same no matter where in the receptive field that feature is located. This means they learn more meaningful features
- Pooling builds up tolerance to severe distortions, such as lens distortions, camera view angle. This pooling takes place at multiple levels, the low-level features are allowed to shift. This makes the CNNs very robust.
- Context: Spatial pooling also makes the receptive fields larger and hence context gets to be increasingly considered after each pooling operation. Hence at high-level layers, the receptive fields will include a huge area of the input stimuli such as an image. It is a well known fact that context improves performance in many recognition tasks.
- Hierarchical feature learning.



# Training



$$y = \phi(Wx), \quad \phi(z) = \frac{e^{z_i}}{\sum e^{z_i}}$$

$y = \phi(z)$  is probability of  $z_i$ . The loss function of soft-max (Cross Entropy) is

$$J = \frac{1}{n} \sum^n [y_i \log p_i + (1 - y_i) \log (1 - p_i)]$$

where  $p_i$  is the probability of  $y_i$

# Training process

- 1 Full connection layer
- 2 Convolution layer
- 3 Polling layer

If the pattern  $x(k)$  belongs to the class  $i$ , for simple case the class is  $y_i(k)$ . The  $i$ -th element of the output  $y_i(k)$ , the training number is  $k = 1, \dots, N$

$$L = \sum_{i=1}^{10} \sum_{k=1}^N [y_i(k) - \hat{y}_i(k)]^2$$

$y_i(k)$  depends on the choice of the output activation function (soft-max).  
For the individual errors on each pattern  $i$

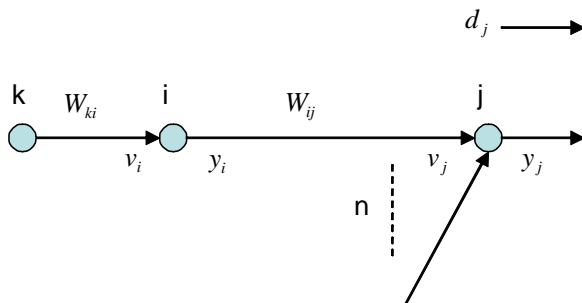
$$J_i = \sum_{k=1}^N [y_i(k) - \hat{y}_i(k)]^2$$

Gradient for "Full connection"

$$\hat{y}(k) = W\varphi(x(k))$$

$$W(k+1) = W(k) - \eta \frac{\partial J}{\partial W}$$

# Full connection layer training

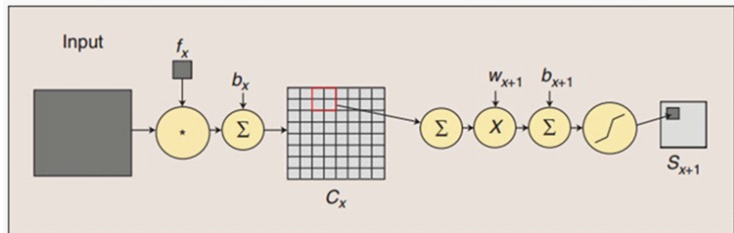


$$w_{ki}(k+1) = w_{ki}(k) - \eta e_j y_k$$

where

$$e(k) = y_j(k) - d_j(k), \quad e_i = e_j W_{ij} \phi_i'$$

# Convolution layer: feedforward

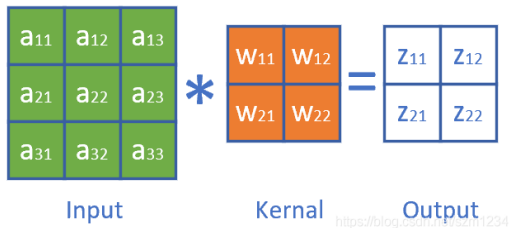


Feedforward calculation (feature map),

$$y = \phi \left( b + \sum_{i=1}^5 \sum_{j=1}^5 w_{i,j} x_{i,j} \right) = \phi \left( b + \sum_{i,j \in M} w_{i,j} x_{i,j} \right) = WX + b$$

where  $\phi$  is the active function (ReLU),  $w_{i,j}$   $b$  also called kernel and bias.

# Convolution layer: feedforward



$$z_{11} = w_{11}a_{11} + w_{12}a_{12} + w_{21}a_{21} + w_{22}a_{22} + b$$

$$z_{12} = w_{11}a_{12} + w_{12}a_{13} + w_{21}a_{22} + w_{22}a_{23} + b$$

some gradients are

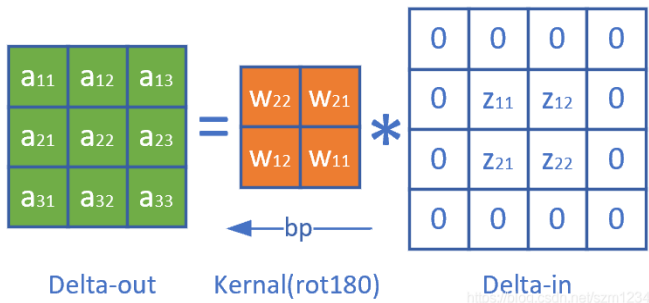
$$\frac{\partial J}{\partial a_{11}} = \frac{\partial J}{\partial z_{11}} \frac{\partial z_{11}}{\partial a_{11}} = e_{z_{11}} w_{11}$$

$$\frac{\partial J}{\partial a_{12}} = \frac{\partial J}{\partial z_{11}} \frac{\partial z_{11}}{\partial a_{12}} + \frac{\partial J}{\partial z_{12}} \frac{\partial z_{12}}{\partial a_{12}} = e_{z_{11}} w_{12} + e_{z_{12}} w_{11}$$

$$\frac{\partial J}{\partial a_{22}} = \frac{\partial J}{\partial z_{11}} \frac{\partial z_{11}}{\partial a_{22}} + \frac{\partial J}{\partial z_{12}} \frac{\partial z_{12}}{\partial a_{22}} + \frac{\partial J}{\partial z_{21}} \frac{\partial z_{21}}{\partial a_{22}} + \frac{\partial J}{\partial z_{22}} \frac{\partial z_{22}}{\partial a_{22}}$$



# Convolution layer training: BP



$$\begin{aligned}\frac{\partial J}{\partial a_{22}} &= e_{z_{11}} w_{22} + e_{z_{12}} w_{21} + e_{z_{21}} w_{12} + e_{z_{22}} w_{11} \\ &\rightarrow z_{11} w_{22} + z_{12} w_{21} + z_{21} w_{12} + z_{22} w_{11}\end{aligned}$$

The error is propagated as: the original kernel  $W$  is rotated by 180 degrees, and then do the convolution operation.

In the corner the padding is needed (fill zero in corners)

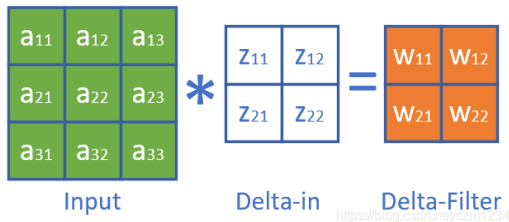
# Convolution layer training: BP

The error in the convolution layer is propagated as

$$e_{out} = e_{in} * W^{rot(180)}$$

$$W(k+1) = W(k) - \eta \frac{\partial J}{\partial W} = W(k) - \eta A e_{out}$$

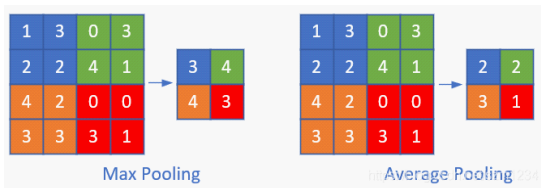
It can be regarded as



So the convolution layers updated as

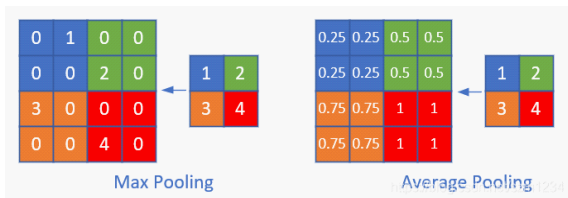
$$W(k+1) = W(k) - \eta A e_{in} * W^{rot(180)}$$

# Pooling layer training



Maximum pooling, the error will be propagated to the position of the original maximum value, and the other positions are all 0.

Average pooling, the error are propagated to all position with the averaged values.



# Pooling layer training

Max pooling: if

$$y = \max(a, b, c, d) = b$$

$$e_{out,a} = 0, \quad e_{out,b} = e_{in}, \quad e_{out,c} = 0, \quad e_{out,d} = 0$$

Average pooling: if

$$y = \frac{1}{4}(a + b + c + d)$$

$$e_{out,a} = \frac{1}{4}e_{in}, \quad e_{out,b} = \frac{1}{4}e_{in}, \quad e_{out,c} = \frac{1}{4}e_{in}, \quad e_{out,d} = \frac{1}{4}e_{in}$$

Up-sampling of average pooling is

$$e_{out} = up[e_{in}], \quad up(z) = z \otimes \mathbf{1}_{n \times m}$$

where  $\otimes$  is the Kronecker product.

No parameter gradient, only error propagation

# Application of CNN: nonlinear system identification

Consider the following unknown discrete-time nonlinear system

$$y(k) = \Gamma [x(k)]$$

where

$$x(k) = [y(k-1), \dots, y(k-n_y), u(k), \dots, u(k-n_u)]^T \quad (1)$$

Simulation model (parallel)

$$\hat{y}(k) = F[u(k), \dots, u(k-n)] \quad (2)$$

also

$$\hat{y}(k) = N[\hat{y}(k-1), \dots, \hat{y}(k-m), u(k), \dots, u(k-n)] \quad (3)$$

Prediction model (series-parallel)

$$\hat{y}(k) = N[y(k-1), \dots, y(k-m), u(k), \dots, u(k-n)] \quad (4)$$

# CNN structure

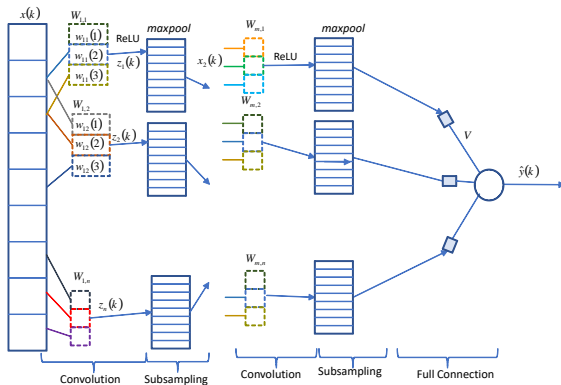


Figure: CNN model for system identification

$$\hat{y}(k) = V\Phi[x(k)] \quad (5)$$

the convoluted output is

$$z = x * W_{ij} \quad z_{zq} = \sum_{k=1}^p x[p(q-1) + k] w_{ij}[k], \quad z = [z_{z1} \cdots z_{zN-p}] \quad (6)$$

we use the rectified linear unit (ReLU)

$$Z(k) = \max[0, z(k)], \quad z = [z_1 \cdots z_n] \quad (7)$$

For system identification,

$$x_2(k) = \max_p [Z(k), s] \quad (8)$$

where  $s$  is the shrink parameter

- In the pre-training, the weights of the convolutional filters  $W_{ij}$  are randomly assigned first, then we use a small dataset and the Backpropagation of CNN to update them.
- In the fine tuning, the weights of the convolutional filters are fixed as the values of the pre-training, we only update the weight  $V$  in the last layer.



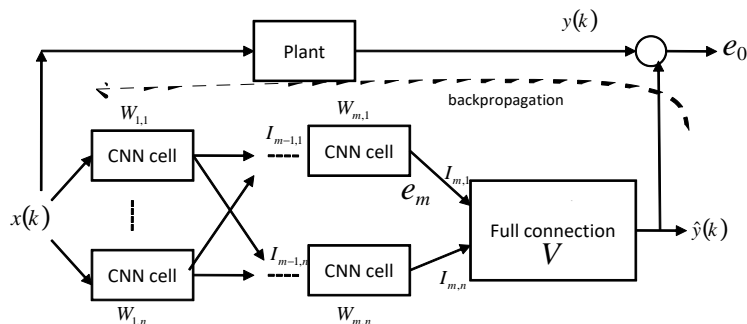


Figure: The scheme of the system identification with CNN

$$J = \frac{1}{2} e_0^2 \quad e_o = \hat{y} - y \quad (9)$$

$$V(k+1) = V(k) - \eta I_m(k) e_o(k) \quad (10)$$

The error is back-propagated to the layer  $m$  as

$$e_m(k) = e_o(k) \frac{\partial \sigma}{\partial t} V = e_o(k) V(k) \quad (11)$$

where  $\sigma$  is the active function of the full connection layer. Here we use linear function,  $\frac{\partial \sigma}{\partial t} = I$ .

For each CNN cell, we need to update the filter weight  $W_{ij}$  in (6)

$$W_{ij}^{(k+1)} = W_{ij}^{(k)} - \eta \frac{\partial J}{\partial W_{ij}}$$

where  $J$  is defined in (9).

For  $m$ -th CNN cell

$$W_{mj}^{(k+1)} = W_{mj}^{(k)} - \eta I_{m-1,j}(k) e_m(k)$$

In the max-pooling,

$$e_{mp}(k) = up[e_m(k)]$$

When the error  $e_{mp}$  goes through the convolution operation

$$e_{mc} = e_{mp} * W_{ij}$$

we use discrete-time de-convolution

$$e_{mc} = e_{mp} * [rot180(W_{ij})]$$

$$e_{m-1} = \text{up}[e_m(k)] * [\text{rot180}(W_{ij})] \quad (12)$$

For any CNN cell, the training law of the filters is

$$\begin{aligned} W_{ij}^{(k+1)} &= W_{ij}^{(k)} - \eta l_{i-1,j}(k) e_i(k) \\ e_i(k) &= \text{up}[e_{i+1}(k)] * [\text{rot180}(W_{ij})] \end{aligned} \quad (13)$$

The object to train the weight  $V$  is

$$J = \min_V \sum_k \|y(k) - \hat{y}(k)\|^2 \quad (14)$$

where

$$\hat{y}(k) = V\Phi[x(k)]$$

The training data are  $y(k)$  and  $\Phi[k]$ .

Consider all data  $k = 1 \dots N$ ,

$$\begin{aligned} \hat{Y} &= [ V\Phi(1) \quad \dots \quad V\Phi(N) ] = V\Psi \\ Y &= [ V\Phi(1) + e_o(1) \quad \dots \quad V\Phi(N) + e_o(N) ] \end{aligned} \quad (15)$$

where  $e_o(k)$  is the modeling error which defined. (15) in matrix form is

$$Y = V\Psi + E \quad (16)$$

where  $E = [e_o(1), \dots, e_o(N)]$ . (16) is a linear-in-parameter system.

When

$$\hat{V} = Y\Psi^T (\Psi\Psi^T)^{-1} = Y\Psi^+ \quad (17)$$

where  $\Psi^+$  is the pseudoinverse of  $\Psi$

$$J = \min_V \sum_k \|y(k) - \hat{y}(k)\|^2 \quad (18)$$

arrives the minimum value.

In each convolutional layer, there are  $n$  filters, the convoluted output is

$$z = x * W_{ij}, \quad z = [zz_1 \cdots zz_{N-p}], \quad zz_q = \sum_{k=1}^p x [p(q-1) + k] w_{ij} [k] \quad (19)$$

where  $p$  is the filter length,  $N$  is the length of  $x(k)$ ,  
 $z(k) = [zz_1 \cdots zz_{N-p}]$  is the feature map generated in this layer.

Linear unit (ReLU) as

$$Z(k) = \max [0, z(k)], \quad z = [z_1 \cdots z_n] \quad (20)$$

We use the max-pooling method in the sub-sampling layer

$$x_2(k) = \max_p [Z(k), s] \quad (21)$$

where  $s$  is the shrink parameter, which depends on the layer,  $x_2(k)$  is a new input for the second convolutional layer



If we use discrete Fourier transform (DFT), the convolution operation in (6) can be transformed into multiplication,

$$\mathcal{F}(x * W) = \mathcal{F}(x) \times \mathcal{F}(W) = \mathcal{F}(x) \mathcal{F}(W) \quad (22)$$

where  $F$  represents the DFT for the corresponding vector,  $*$  represents the convolution operation,  $F(x)$  and  $F(W)$  are vectors in the frequency domain,  $\times$  represents element-wise product.

## Theorem

The input data  $x(k)$ ,  $k = 1 \cdots N$ , pass through the filter  $W_{ij} = [w_{ij}(1), \cdots, w_{ij}(p)]$ ,  $p < N$ . There is the optimal input, which maximizes the CNN cell with the ReLU and the max-pooling as

$$x^*(k) = \frac{\sqrt{2}}{N} \cos\left(\frac{2\pi f_0}{N}k + \phi\right) \quad (23)$$

where  $f_0$  is the optimal frequency,  $\phi$  is the phase. .

- This theorem shows that for any inputs (random or determined), it corresponds to the optimal filter. The frequency of the optimum input is the frequency of maximum magnitude in the filter.
- Any random filter containing some frequencies of moderate magnitude can generate the maximum magnitude of the best input frequency.