

ECE/SIOC 228 Machine Learning for Physical Applications

Lecture 5: Convolutional Neural Network

Yuanyuan Shi

Assistant Professor, ECE

University of California, San Diego

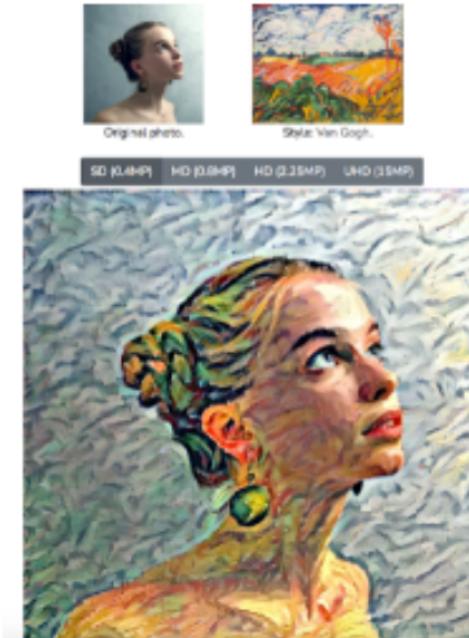
Computer Vision Problems

UC San Diego

Image classification



→Cat? (0/1)



SD (0.4MP) HD (0.8MP) HD (2.2MP) UHD (3.5MP)

<https://neuralstyle.art/>

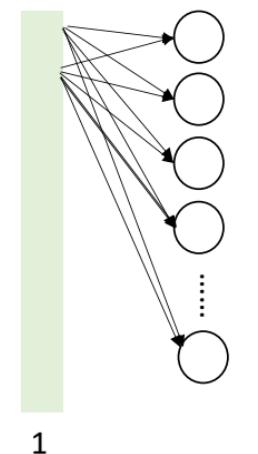
Using MLP?

UC San Diego

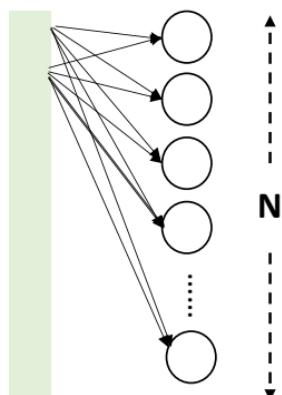


Reshape

$$100 \times 100 \times 3 = 30,000$$



1



1 Layer of MLP

| | N | # parameters |
|--|---------|---------------|
| | 1,000 | 30,000,000 |
| | 10,000 | 300,000,000 |
| | 100,000 | 3,000,000,000 |

of parameters become huge
even if we only have 1 layer !

Spatial Structure

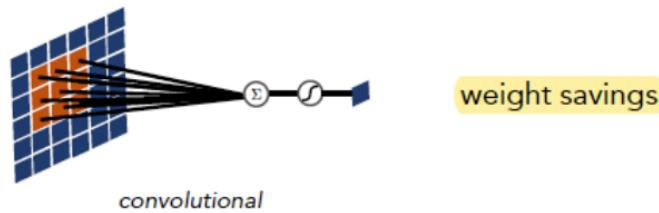
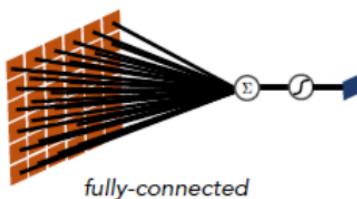
UC San Diego



- We can exploit the spatial structure of image data to reduce the amount of learning



- ✓ nearby areas tend to contain stronger patterns
- ✓ nearby pixels tend to share similar characteristics and are combined in certain ways

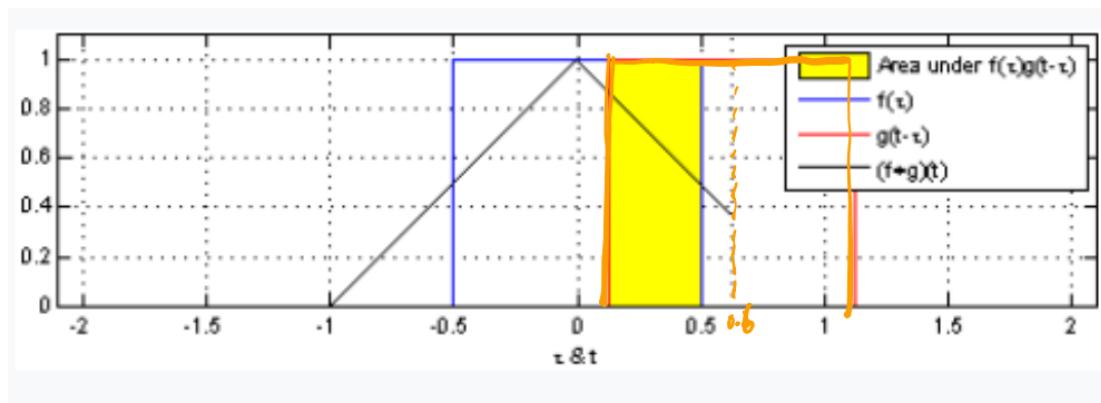


Convolution

UC San Diego

- **Convolution** is a mathematical operation on two functions (f and g) that produces a third function ($f * g$) that expresses how the shape of one is modified by the other.

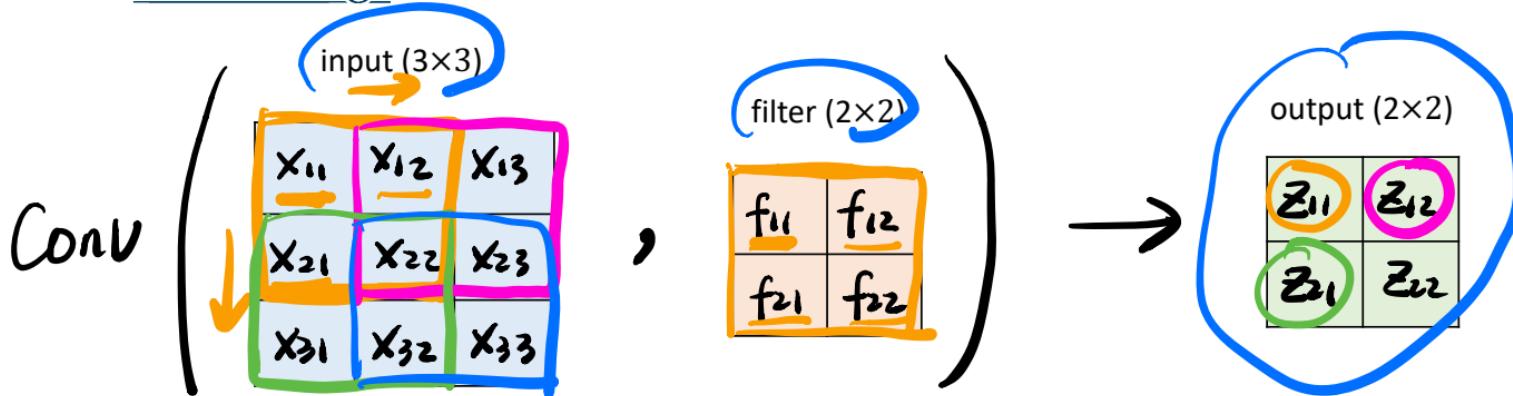
$$(f * g)(t) = \int_{-\infty}^{\infty} f(\tau)g(t - \tau)d\tau = (g * f)(t)$$



The amount of yellow is the area of $\int_{-\infty}^{\infty} f(\tau)g(t - \tau)d\tau$ for the current value of t .

Convolution in Neural Networks

UC San Diego



$$\underline{z_{11}} = f_{11}x_{11} + f_{12}x_{12} + f_{21}x_{21} + f_{22}x_{22}$$

$$\underline{z_{12}} = f_{11}x_{12} + f_{12}x_{13} + f_{21}x_{22} + f_{22}x_{23}$$

$$\underline{z_{21}} = f_{11}x_{21} + f_{12}x_{22} + f_{21}x_{31} + f_{22}x_{32}$$

$$\underline{z_{22}} = f_{11}x_{22} + f_{12}x_{23} + f_{21}x_{32} + f_{22}x_{33}$$

Convolution Example: Vertical Edge Detection

UC San Diego

Conv

Input (6×6)

| | | | | | |
|----|----------------------|----------------------|----------------------|---------------------|----------------------|
| 10 | 10 x ₁ | 10 x ₀ | 0 x ₋₁ | 0 | 0 |
| 10 | 10 x ₁ | 10 x ₀ | 0 x ₋₁ | 0 | 0 |
| 10 | 10 x ₁ | 10 x ₀ | 0 x ₋₁ | 0 | 0 |
| 10 | 10 | 10 | 0 x ₁ | 0 x ₀ | 0 x ₋₁ |
| 10 | 10 | 10 | 0 x ₁ | 0 x ₀ | 0 x ₋₁ |
| 10 | 10 | 10 | 0 x ₁ | 0 x ₀ | 0 x ₋₁ |

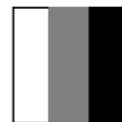
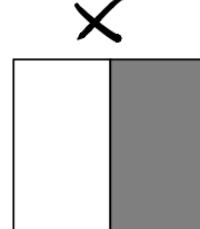
Filter (3×3)

| | | |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |
| 1 | 0 | -1 |

f

Output (4×4)

| | | | |
|---|----|----|---|
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |



=



Convolution Example: Horizontal Edge Detection

UC San Diego

Input (6×6)

| | | | | | |
|----|----|----|----|----|----|
| 10 | 10 | 10 | 10 | 10 | 10 |
| 10 | 10 | 10 | 10 | 10 | 10 |
| 10 | 10 | 10 | 10 | 10 | 10 |
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 |

Filter (3×3)

| | | |
|----|----|----|
| 1 | 1 | 1 |
| 0 | 0 | 0 |
| -1 | -1 | -1 |

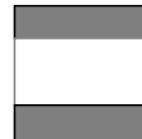
=

Output (4×4)

| | | | |
|----|----|----|----|
| 0 | 0 | 0 | 0 |
| 30 | 30 | 30 | 30 |
| 30 | 30 | 30 | 30 |
| 0 | 0 | 0 | 0 |

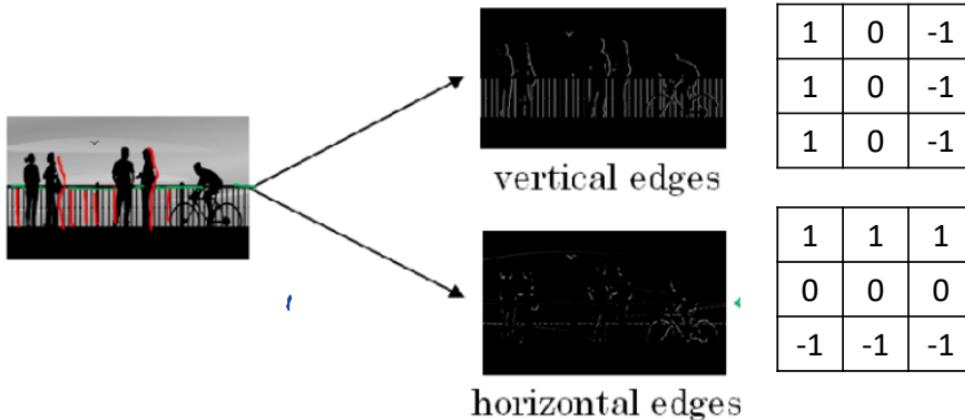


=



Learning to Detect Edges

UC San Diego

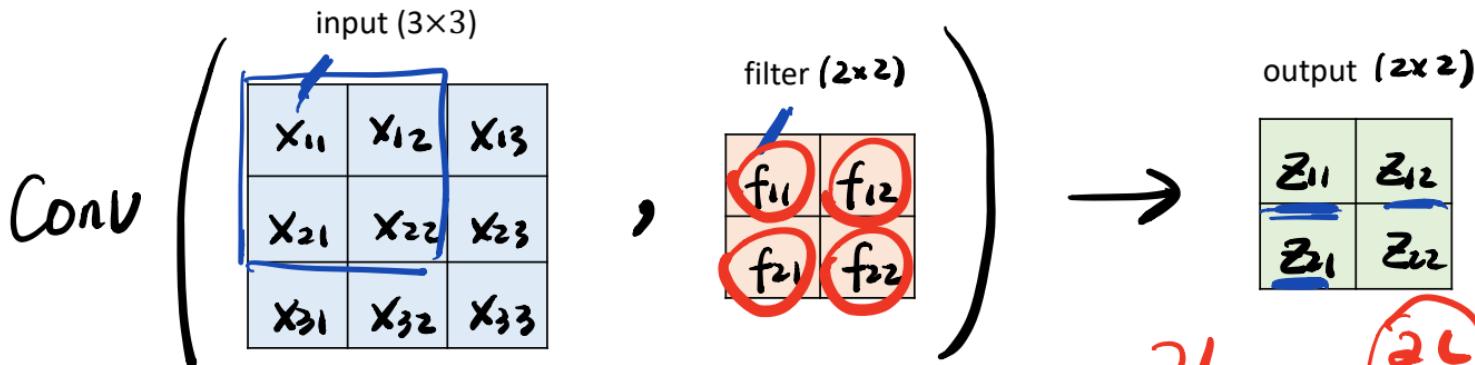


| | | |
|----------|----------|----------|
| f_{11} | f_{12} | f_{13} |
| f_{21} | f_{22} | f_{23} |
| f_{31} | f_{32} | f_{33} |

Filter (3×3)

Forward Computation of Convolutional Layer

UC San Diego



$$z_{11} = \underline{f_{11}x_{11}} + f_{12}x_{12} + f_{21}x_{21} + f_{22}x_{22}$$

$$z_{12} = \underline{f_{11}x_{12}} + f_{12}x_{13} + f_{21}x_{22} + f_{22}x_{23}$$

$$z_{21} = \underline{f_{11}x_{21}} + f_{12}x_{22} + f_{21}x_{31} + f_{22}x_{32}$$

$$z_{22} = \underline{f_{11}x_{22}} + f_{12}x_{23} + f_{21}x_{32} + f_{22}x_{33}$$

$$\frac{\partial L}{\partial f_{11}} = \cancel{\frac{\partial L}{\partial z_{11}}} \cdot \frac{\partial z_{11}}{\partial f_{11}}$$

$$\frac{\partial L}{\partial f_{12}} = \cancel{\frac{\partial L}{\partial z_{12}}} \cdot \frac{\partial z_{12}}{\partial f_{12}}$$

$$\frac{\partial L}{\partial f_{21}} = \cancel{\frac{\partial L}{\partial z_{21}}} \cdot \frac{\partial z_{21}}{\partial f_{21}}$$

$$\frac{\partial L}{\partial f_{22}} = \cancel{\frac{\partial L}{\partial z_{22}}} \cdot \frac{\partial z_{22}}{\partial f_{22}}$$

Backward Computation of Convolutional Layer

UC San Diego

- Now, Compute the derivatives of loss "L" with respect to filter "f", assume. $dZ = \frac{\partial L}{\partial Z}$ known

$$\frac{\partial L}{\partial f_{11}} = \frac{\partial L}{\partial Z_{11}} \cdot \frac{\partial Z_{11}}{\partial f_{11}} + \frac{\partial L}{\partial Z_{12}} \cdot \frac{\partial Z_{12}}{\partial f_{11}} + \frac{\partial L}{\partial Z_{21}} \cdot \frac{\partial Z_{21}}{\partial f_{11}} + \frac{\partial L}{\partial Z_{22}} \cdot \frac{\partial Z_{22}}{\partial f_{11}}$$

$$\frac{\partial L}{\partial f_{12}} = \frac{\partial L}{\partial Z_{11}} \cdot \frac{\partial Z_{11}}{\partial f_{12}} + \frac{\partial L}{\partial Z_{12}} \cdot \frac{\partial Z_{12}}{\partial f_{12}} + \frac{\partial L}{\partial Z_{21}} \cdot \frac{\partial Z_{21}}{\partial f_{12}} + \frac{\partial L}{\partial Z_{22}} \cdot \frac{\partial Z_{22}}{\partial f_{12}}$$

$$\frac{\partial L}{\partial f_{21}} = \frac{\partial L}{\partial Z_{11}} \cdot \frac{\partial Z_{11}}{\partial f_{21}} + \frac{\partial L}{\partial Z_{12}} \cdot \frac{\partial Z_{12}}{\partial f_{21}} + \frac{\partial L}{\partial Z_{21}} \cdot \frac{\partial Z_{21}}{\partial f_{21}} + \frac{\partial L}{\partial Z_{22}} \cdot \frac{\partial Z_{22}}{\partial f_{21}}$$

$$\frac{\partial L}{\partial f_{22}} = \frac{\partial L}{\partial Z_{11}} \cdot \frac{\partial Z_{11}}{\partial f_{22}} + \frac{\partial L}{\partial Z_{12}} \cdot \frac{\partial Z_{12}}{\partial f_{22}} + \frac{\partial L}{\partial Z_{21}} \cdot \frac{\partial Z_{21}}{\partial f_{22}} + \frac{\partial L}{\partial Z_{22}} \cdot \frac{\partial Z_{22}}{\partial f_{22}}$$

Backward Computation of Convolutional Layer

UC San Diego

- Which can be written as,

$$\frac{\partial L}{\partial f_{11}} = \frac{\partial L}{\partial z_{11}} \cdot X_{11} + \frac{\partial L}{\partial z_{12}} \cdot X_{12} + \frac{\partial L}{\partial z_{21}} \cdot X_{21} + \frac{\partial L}{\partial z_{22}} \cdot X_{22}$$

$$\frac{\partial L}{\partial f_{12}} = \frac{\partial L}{\partial z_{11}} \cdot X_{12} + \frac{\partial L}{\partial z_{12}} \cdot X_{13} + \frac{\partial L}{\partial z_{21}} \cdot X_{22} + \frac{\partial L}{\partial z_{22}} \cdot X_{23}$$

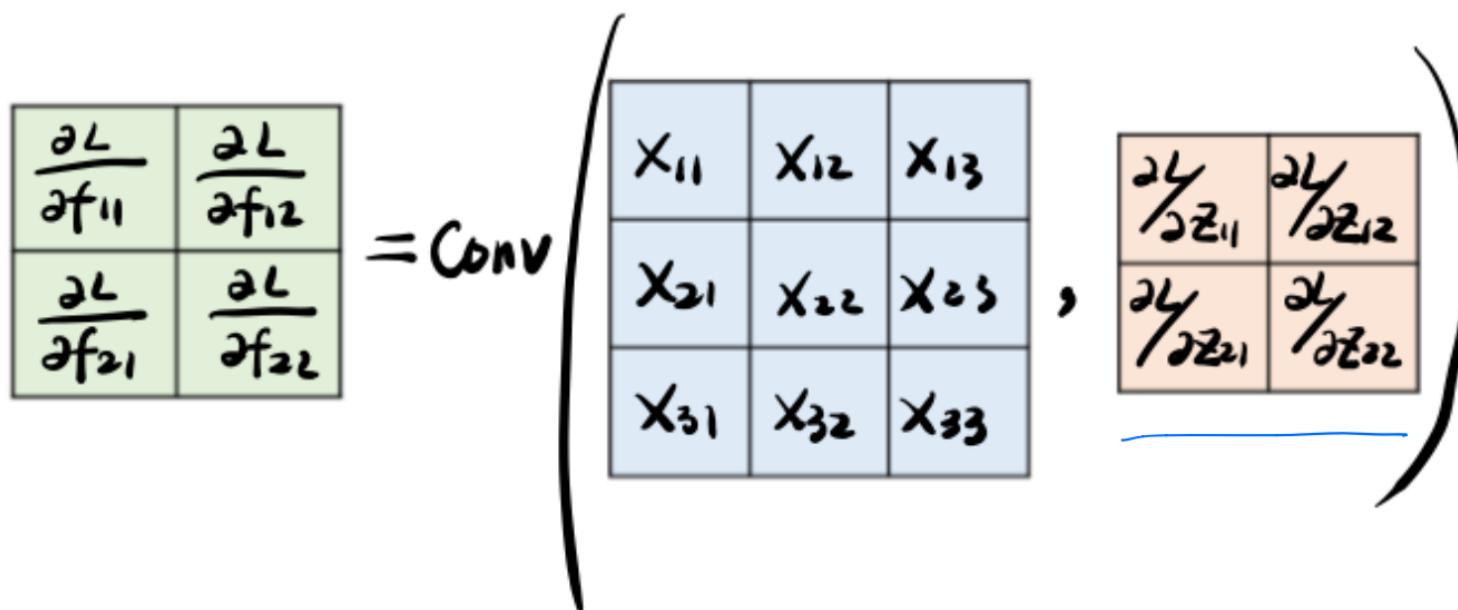
$$\frac{\partial L}{\partial f_{21}} = \frac{\partial L}{\partial z_{11}} \cdot X_{21} + \frac{\partial L}{\partial z_{12}} \cdot X_{22} + \frac{\partial L}{\partial z_{21}} \cdot X_{31} + \frac{\partial L}{\partial z_{22}} \cdot X_{32}$$

$$\frac{\partial L}{\partial f_{22}} = \frac{\partial L}{\partial z_{11}} \cdot X_{22} + \frac{\partial L}{\partial z_{12}} \cdot X_{23} + \frac{\partial L}{\partial z_{21}} \cdot X_{32} + \frac{\partial L}{\partial z_{22}} \cdot X_{33}$$

Backward Computation of Convolutional Layer

UC San Diego

- The derivative can also be written in the form of convolution!



Backward Computation of Convolutional Layer

UC San Diego

Similarly, we can find the derivatives of the input matrix X w.r.t. error L

$$\frac{\partial L}{\partial X_{11}} = \boxed{\frac{\partial L}{\partial Z_{11}}} \cdot \boxed{\frac{\partial Z_{11}}{\partial X_{11}}} + \frac{\partial L}{\partial Z_{12}} \cdot \frac{\partial Z_{12}}{\partial X_{11}} + \frac{\partial L}{\partial Z_{21}} \cdot \frac{\partial Z_{21}}{\partial X_{11}} + \frac{\partial L}{\partial Z_{22}} \cdot \frac{\partial Z_{22}}{\partial X_{11}}$$
$$= \frac{\partial L}{\partial Z_{11}} \cdot f_{11} + \frac{\partial L}{\partial Z_{12}} \cdot 0 + \frac{\partial L}{\partial Z_{21}} \cdot 0 + \frac{\partial L}{\partial Z_{22}} \cdot 0$$

$$\frac{\partial L}{\partial X_{12}} = \frac{\partial L}{\partial Z_{11}} \cdot f_{12} + \frac{\partial L}{\partial Z_{21}} f_{11}$$

Backward Computation of Convolutional Layer

UC San Diego

$$\frac{\partial L}{\partial x_{13}} = \frac{\partial L}{\partial z_{12}} \cdot f_{12}$$

$$\frac{\partial L}{\partial x_{21}} = \frac{\partial L}{\partial z_{11}} \cdot f_{21} + \frac{\partial L}{\partial z_{21}} \cdot f_{11}$$

$$\frac{\partial L}{\partial x_{22}} = \frac{\partial L}{\partial z_{11}} \cdot f_{22} + \frac{\partial L}{\partial z_{12}} \cdot f_{21} + \frac{\partial L}{\partial z_{21}} \cdot f_{12} + \frac{\partial L}{\partial z_{22}} \cdot f_{11}$$

$$\frac{\partial L}{\partial x_{23}} = \frac{\partial L}{\partial z_{12}} \cdot f_{22} + \frac{\partial L}{\partial z_{22}} \cdot f_{11}$$

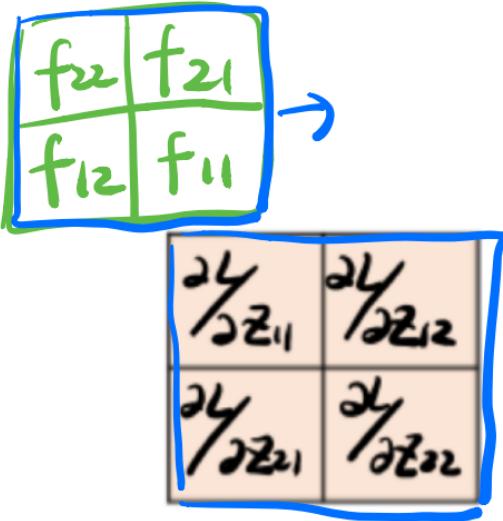
$$\frac{\partial L}{\partial x_{31}} = \frac{\partial L}{\partial z_{21}} \cdot f_{21}$$

$$\frac{\partial L}{\partial x_{32}} = \frac{\partial L}{\partial z_{21}} \cdot f_{22} + \frac{\partial L}{\partial z_{22}} \cdot f_{21}$$

$$\frac{\partial L}{\partial x_{33}} = \frac{\partial L}{\partial z_{22}} \cdot f_{22}$$

Backward Computation of Convolutional Layer

UC San Diego



(full
convolution)

$$\frac{\partial L}{\partial x_{11}} = f_{11} \cdot \frac{\partial L}{\partial z_{11}}$$

$$\frac{\partial L}{\partial x_{12}} = f_{12} \cdot \frac{\partial L}{\partial z_{11}} + f_{11} \cdot \frac{\partial L}{\partial z_{12}}$$

$$\frac{\partial L}{\partial x_{13}} = f_{12} \cdot \frac{\partial L}{\partial z_{12}}$$

$$\frac{\partial L}{\partial x_{21}} = f_{21} \cdot \frac{\partial L}{\partial z_{11}} + f_{11} \cdot \frac{\partial L}{\partial z_{21}}$$

$$\frac{\partial L}{\partial x_{22}} = f_{22} \cdot \frac{\partial L}{\partial z_{11}} + f_{21} \cdot \frac{\partial L}{\partial z_{21}} + f_{12} \cdot \frac{\partial L}{\partial z_{12}} + f_{11} \cdot \frac{\partial L}{\partial z_{22}}$$

$$\frac{\partial L}{\partial x_{23}} = f_{22} \cdot \frac{\partial L}{\partial z_{12}} + f_{12} \cdot f_{12} \cdot \frac{\partial L}{\partial z_{22}}$$

$$\frac{\partial L}{\partial x_{31}} = f_{21} \cdot \frac{\partial L}{\partial z_{21}}$$

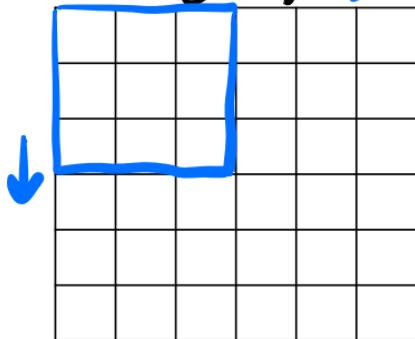
$$\frac{\partial L}{\partial x_{32}} = f_{22} \cdot \frac{\partial L}{\partial z_{21}} + f_{21} \cdot \frac{\partial L}{\partial z_{22}}$$

$$\frac{\partial L}{\partial x_{33}} = f_{22} \cdot \frac{\partial L}{\partial z_{22}}$$

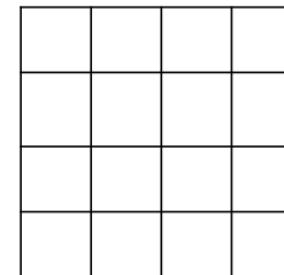
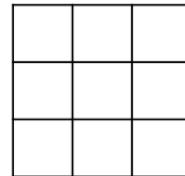
Padding

UC San Diego

(6x6) →

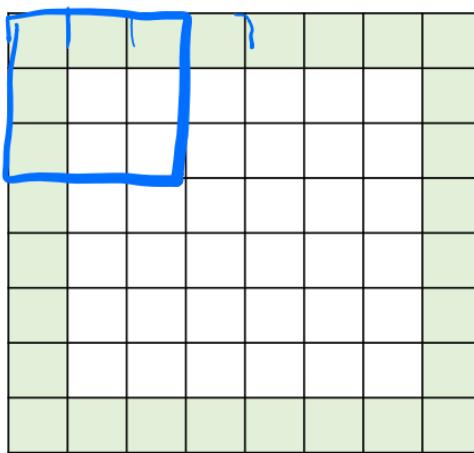


3x3

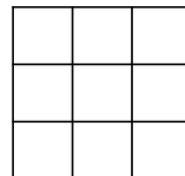


- dimension shrinks

(4x4)

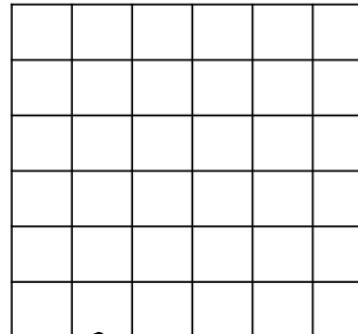


(8x8)



3x3

- padding zero



(8x8)

Same padding:
dimension stays the same after
Convolution

Strided Convolution

UC San Diego

- Only compute output at some integer interval to downsample image

The diagram illustrates a strided convolution operation. An input matrix of size 7x7 is shown with values ranging from 0 to 8. A 3x3 filter is applied with a stride of 2. The filter's receptive field is highlighted with colored boxes: green for the top-left element (17), orange for the middle element (29), and blue for the bottom-right element (9). The input values are labeled with subscripts indicating their position relative to the filter's receptive field. For example, the value 17 has subscripts 0, 0, 0, while the value 29 has subscripts 1, 0, 2. The bottom row of the input matrix is also labeled with its stride value, 2.

| | | | | | | |
|-----------------|----------------|------------------|----------------|-----------------|----------------|----------------|
| 3 ₀ | 0 ₀ | 1 ₀ | 2 ₀ | 2 ₀ | 7 ₀ | 4 ₀ |
| 1 ₁ | 5 ₀ | 8 ₁ | 9 ₀ | 9 ₁₂ | 3 ₀ | 1 ₂ |
| 2 ₋₁ | 7 ₀ | 2 ₋₁₁ | 5 ₀ | 5 ₋₁ | 1 ₀ | 3 ₁ |
| 0 | 1 | 3 | 1 | 1 | 7 | 8 |
| 4 | 2 | 1 | 6 | 6 | 2 | 8 |
| 2 | 4 | 5 | 2 | 2 | 3 | 9 |
| 1 | 3 | 4 | 5 | 7 | 8 | 3 |

Input (7x7)

Filter (3x3)

| | | |
|----|---|---|
| 0 | 0 | 0 |
| 1 | 0 | 2 |
| -1 | 0 | 1 |

Stride 2

=

Output (3x3)

| | | |
|----|----|---|
| 17 | 29 | 9 |
| | | |
| | | |

Summary of Convolutions

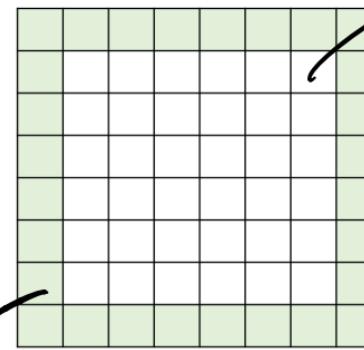
UC San Diego

$n \times n$ image

padding p

$f \times f$ filter

stride s



6x6 image

| | | |
|----|---|---|
| 0 | 0 | 0 |
| 1 | 0 | 2 |
| -1 | 0 | 1 |

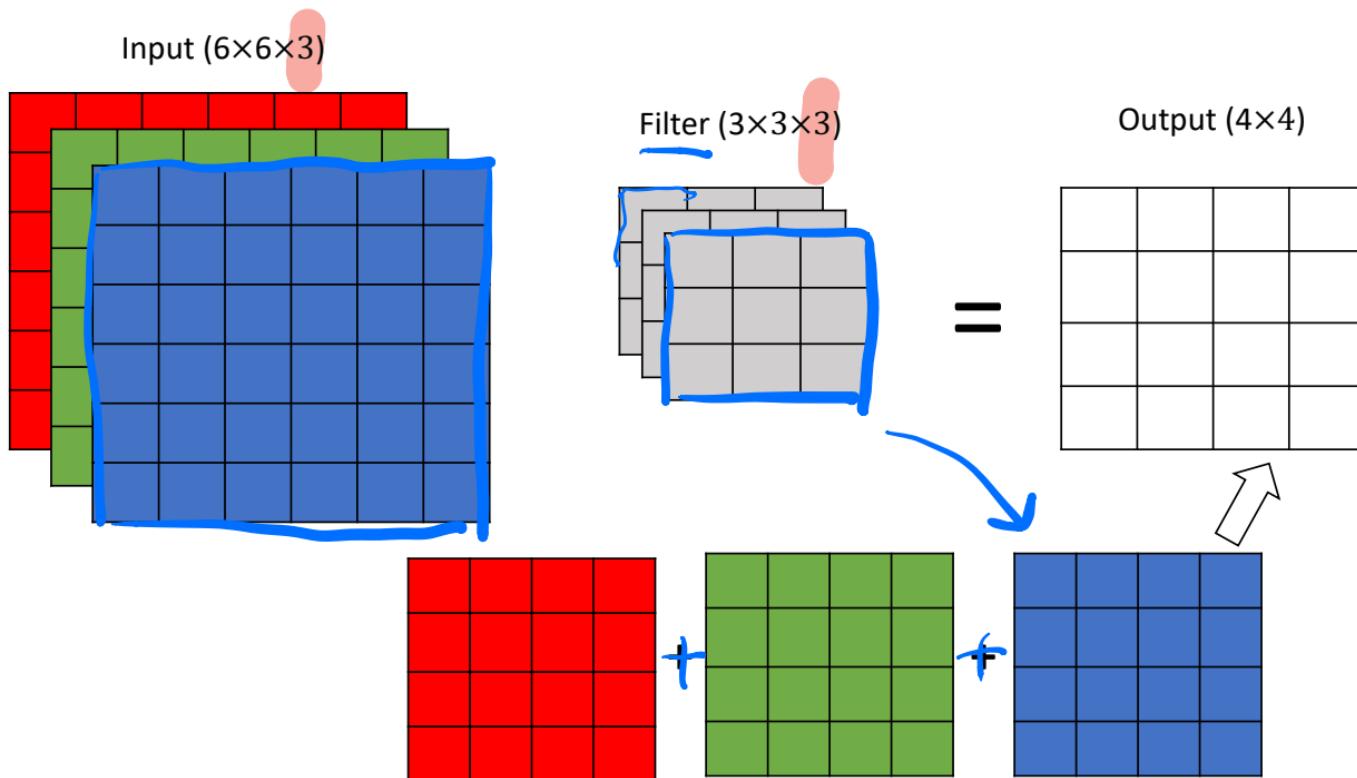
Filter (3x3)

$$\Rightarrow \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor$$

floor function

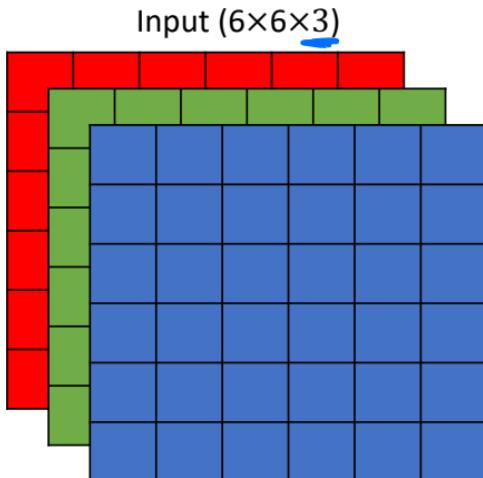
Convolutions on RGB images

UC San Diego

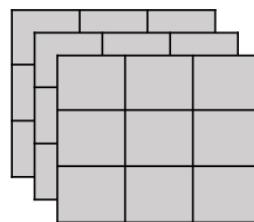


More Filters

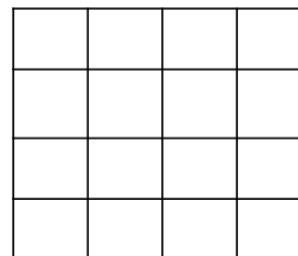
UC San Diego



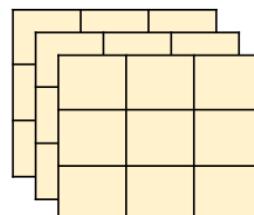
f_xf_y
Filter1 ($3 \times 3 \times 3$)



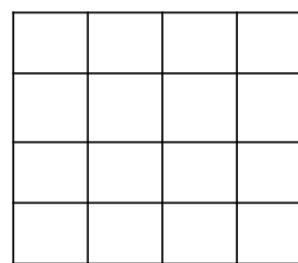
Output1 (4×4)



=



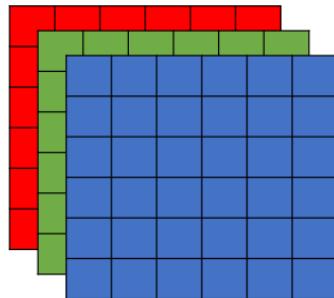
Filter2 ($3 \times 3 \times 3$)



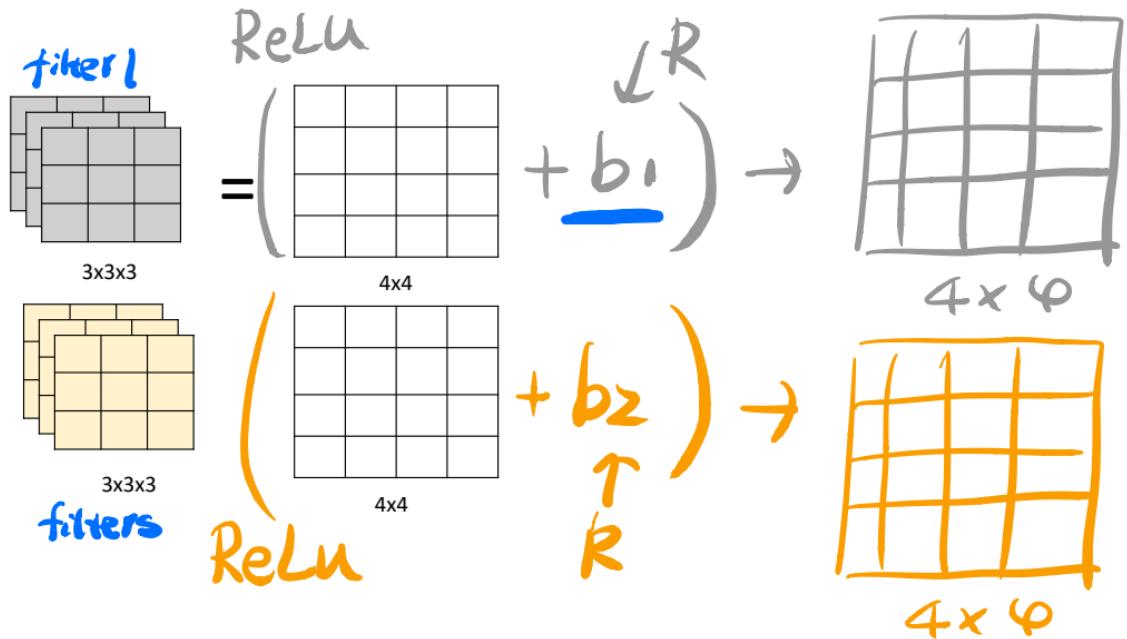
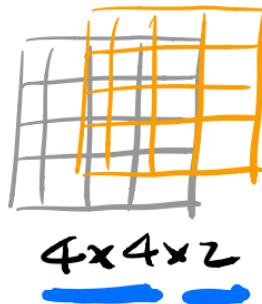
Output2 (4×4)

One layer of Convolutional Neural Nets

UC San Diego



↓
1 layer
of CNN



use as input for
next convolutional layer

Summary of Notation

UC San Diego

For 1 Convolutional layer:

$f^{[l]}$ = filter size

$p^{[l]}$ = padding

$s^{[l]}$ = stride

$n_c^{[l]}$ = number of conv filters

Input: $\underline{n_H^{[l-1]}} \times \underline{n_W^{[l-1]}} \times \underline{n_C^{[l-1]}}$

Output: $n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

Parameters in each filter

$$\underline{f^{[l]} \times f^{[l]} \times n_c^{[l-1]}}$$

Total # of parameters in 1conv layer

$$f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times \underline{n_c^{[l]}}$$

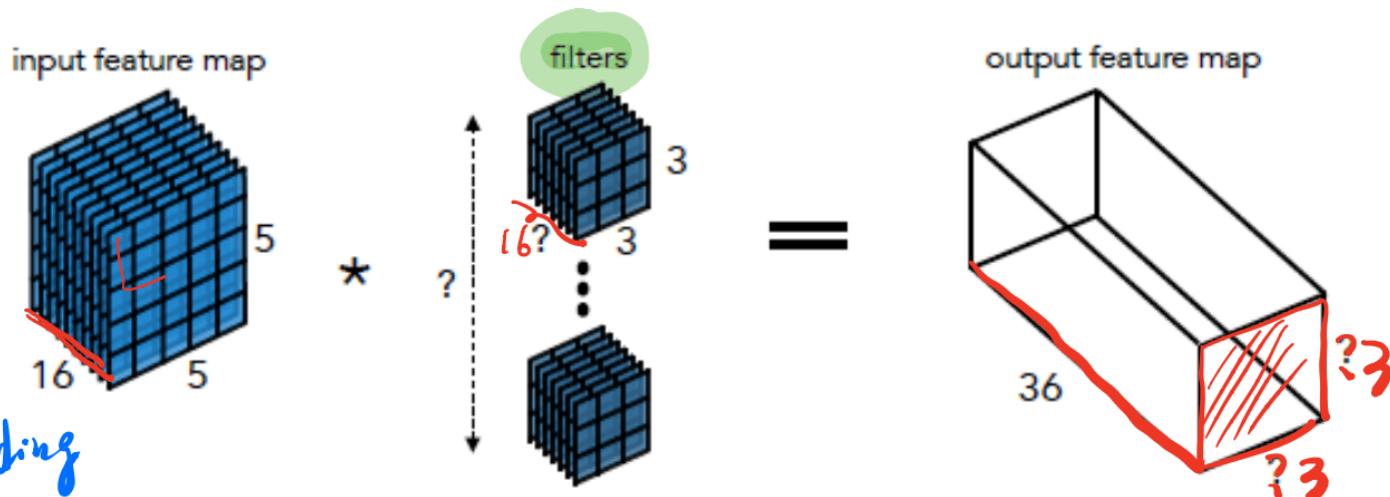
number of filters

$$\text{~~~} n_H^{[l]} = \left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$$

$$n_W^{[l]} = \left\lfloor \frac{n_W^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$$

Pop Quiz

UC San Diego



*no padding
stride = 1*

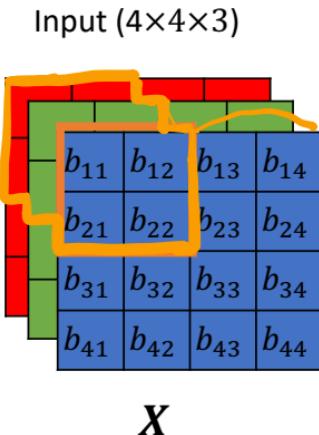
- What is the size of each filter?
- How many filters do we have?

3x3x16

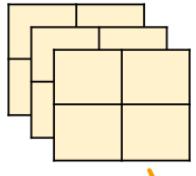
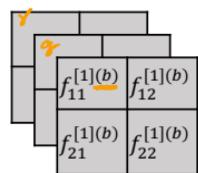
36

Implementation

UC San Diego



Filter1 ($2 \times 2 \times 3$)



Filter2 ($2 \times 2 \times 3$)

stride = 2

- Take each $[2 \times 2 \times 3]$ blocks of pixels in the input and stretch each block into a column vector of size 12
- Iterating this process with $s = 2$ gives $(4-2)/2+1 = 2$ locations along both width and height, leading to a input matrix with size $[12 \times 4]$
- The weights of the CONV layer are similarly stretched out into rows. Since we have 2 filters of size $[2 \times 2 \times 3]$, this gives us a matrix W of size $[7 \times 12]$
- We can express convolution as WX
- Output of this operation is $[2 \times 4]$
- We can reshape the output to $[2 \times 2 \times 2]$

$f_{11}^{[1](b)}, f_{12}^{[1](b)}, f_{21}^{[1](b)}, f_{22}^{[1](b)}, f_{11}^{[1](g)}, f_{12}^{[1](g)}, f_{21}^{[1](g)}, f_{22}^{[1](g)}, f_{11}^{[1](r)}, f_{12}^{[1](r)}, f_{21}^{[1](r)}, f_{22}^{[1](r)}$
 $f_{11}^{[2](b)}, f_{12}^{[2](b)}, f_{21}^{[2](b)}, f_{22}^{[2](b)}, f_{11}^{[2](g)}, f_{12}^{[2](g)}, f_{21}^{[2](g)}, f_{22}^{[2](g)}, f_{11}^{[2](r)}, f_{12}^{[2](r)}, f_{21}^{[2](r)}, f_{22}^{[2](r)}$

W

X_{col}

$b_{11}, b_{13}, b_{31}, b_{33}$
 $b_{12}, b_{14}, b_{32}, b_{34}$
 $b_{21}, b_{23}, b_{41}, b_{43}$
 $b_{22}, b_{24}, b_{42}, b_{44}$
 $g_{11}, g_{13}, g_{31}, g_{33}$
 $g_{12}, g_{14}, g_{32}, g_{34}$
 $g_{21}, g_{23}, g_{41}, g_{43}$
 $g_{22}, g_{24}, g_{42}, g_{44}$
 $r_{11}, r_{13}, r_{31}, r_{33}$
 $r_{12}, r_{14}, r_{32}, r_{34}$
 $r_{21}, r_{23}, r_{41}, r_{43}$
 $r_{22}, r_{24}, r_{42}, r_{44}$

Convolution v.s. Fully Connected Layer

UC San Diego

Input (consider 1 channel)

| | | | |
|----------|----------|----------|----------|
| b_{11} | b_{12} | b_{13} | b_{14} |
| b_{21} | b_{22} | b_{23} | b_{24} |
| b_{31} | b_{32} | b_{33} | b_{34} |
| b_{41} | b_{42} | b_{43} | b_{44} |

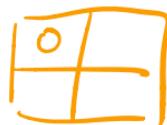
4×4

X

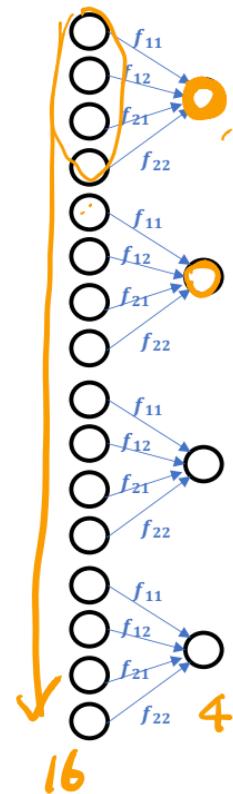
| | |
|----------|----------|
| f_{11} | f_{12} |
| f_{21} | f_{22} |

2×2

consider 1 filter



Convolution



local receptive field

In convolution layers, we enforce the following prior on learning:

- Connections are restricted to nearby regions
- Same filters (weights) are applied throughout the input
- We only have 4 rather than 64 parameters to learn!

weight sharing

$$16 \times 4 = 64$$

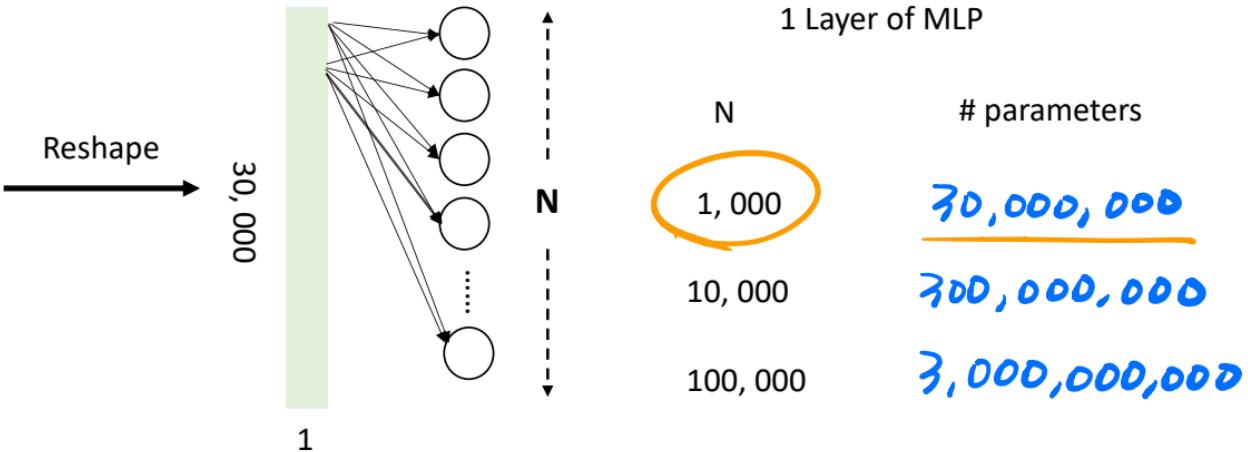
Convolution v.s. Fully Connected Layer

UC San Diego



Reshape

$$100 \times 100 \times 3 = 30,000$$



- If you have 100 filters that are $3 \times 3 \times 3$ in one layer of a neural network, how many parameters does that layer have?

$$(27+1) \times 100 = \underline{\underline{2800}}$$

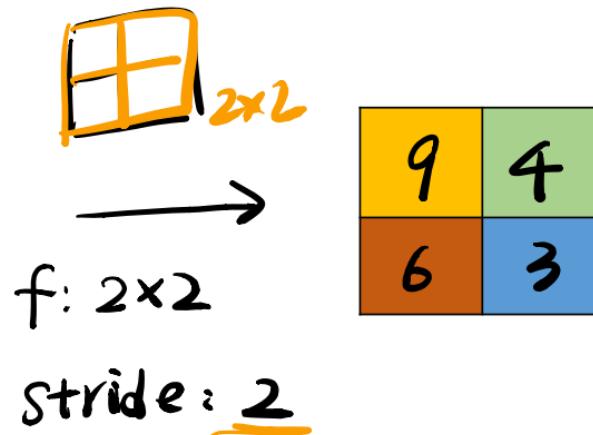
Pooling Layer- Max pooling

UC San Diego

Why pooling?

- Reduce feature dimension \Rightarrow Reduce # of learnable parameters
- Summarize features in a local region (Robust to variations)

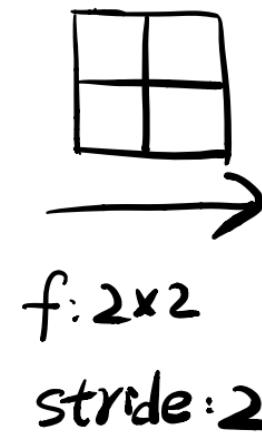
| | | | |
|---|---|---|---|
| 1 | 3 | 2 | 4 |
| 2 | 9 | 1 | 1 |
| 1 | 3 | 2 | 3 |
| 5 | 6 | 1 | 2 |



Pooling Layer- Average pooling

UC San Diego

| | | | |
|---|---|---|---|
| 1 | 3 | 2 | 4 |
| 2 | 9 | 1 | 1 |
| 1 | 3 | 2 | 3 |
| 5 | 6 | 1 | 2 |



| | |
|------|------|
| 3.75 | 1.25 |
| 4 | 2 |

Summary of Pooling

UC San Diego

★ No parameters to be learned

Hyperparameters:

f : filter size

✓ $f=2, S=2$ (shrink input size by 2)

s : stride

✓ $f=3, S=2$

Max or average pooling

$$\text{Input: } n_H \times n_W \times n_C \xrightarrow[\substack{\text{after} \\ \text{pooling}}]{\longrightarrow} \left\lfloor \frac{n_H - f}{S} + 1 \right\rfloor \times \left\lfloor \frac{n_W - f}{S} + 1 \right\rfloor \times n_C$$

Backward Computation of Pooling Layer

UC San Diego

| | | | |
|---|---|---|---|
| 1 | 3 | 2 | 4 |
| 2 | 9 | 1 | 1 |
| 1 | 3 | 2 | 3 |
| 5 | 6 | 1 | 2 |



max pooling

$$Z_{11} = \max(X_{11}, X_{12}, X_{21}, X_{22}) \\ = 9$$

$$\frac{\partial Z_{11}}{\partial X_{11}} = 0, \quad \frac{\partial Z_{11}}{\partial X_{12}} = 0$$

$$\frac{\partial Z_{11}}{\partial X_{21}} = 0, \quad \frac{\partial Z_{11}}{\partial X_{22}} = 1$$

Suppose $Z = \max(X_1, \dots, X_n)$

$$\frac{\partial Z}{\partial x} = \begin{cases} 1 & \text{if } \underline{x_i} = \max(X_1, \dots, X_n) \\ 0 & \text{otherwise} \end{cases}$$

Backward Computation of Pooling Layer

UC San Diego

| | | | |
|---|---|---|---|
| 1 | 3 | 2 | 4 |
| 2 | 9 | 1 | 1 |
| 1 | 3 | 2 | 3 |
| 5 | 6 | 1 | 2 |

average pooling



$$Z_{11} = \frac{1}{4} (x_{11} + x_{12} + x_{21} + x_{22})$$

$$\frac{\partial Z_{11}}{\partial x_{11}} = \frac{1}{4}$$

$$\frac{\partial Z_{11}}{\partial x_{12}} = \frac{1}{4}$$

$$\frac{\partial Z_{11}}{\partial x_{21}} = \frac{1}{4}$$

$$\frac{\partial Z_{11}}{\partial x_{22}} = \frac{1}{4}$$

Classic Networks

UC San Diego

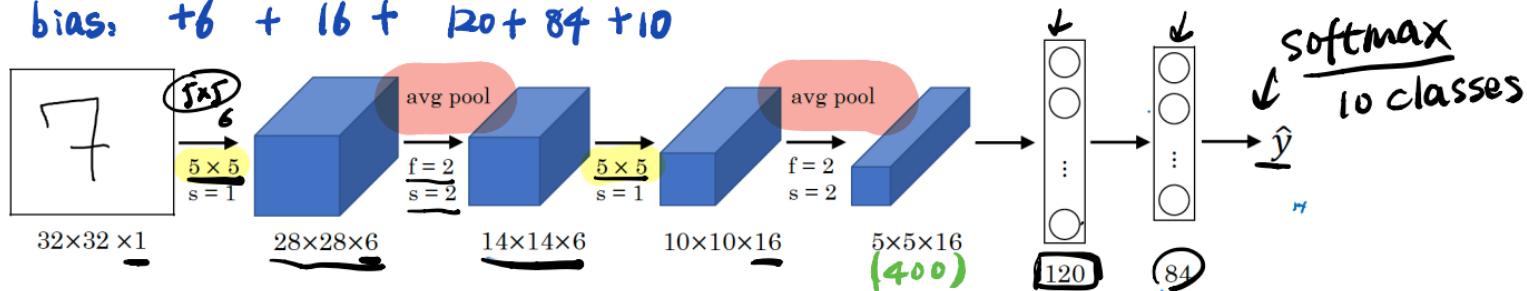
- Classic Networks
 - LeNet – 5
 - AlexNet
 - VGG -16
- ResNet (152 layer)

LeNet - 5

UC San Diego

parameters: $150 + 2400 + 48000 + 10080 + 840$

bias: $+6 + 16 + 120 + 84 + 10$



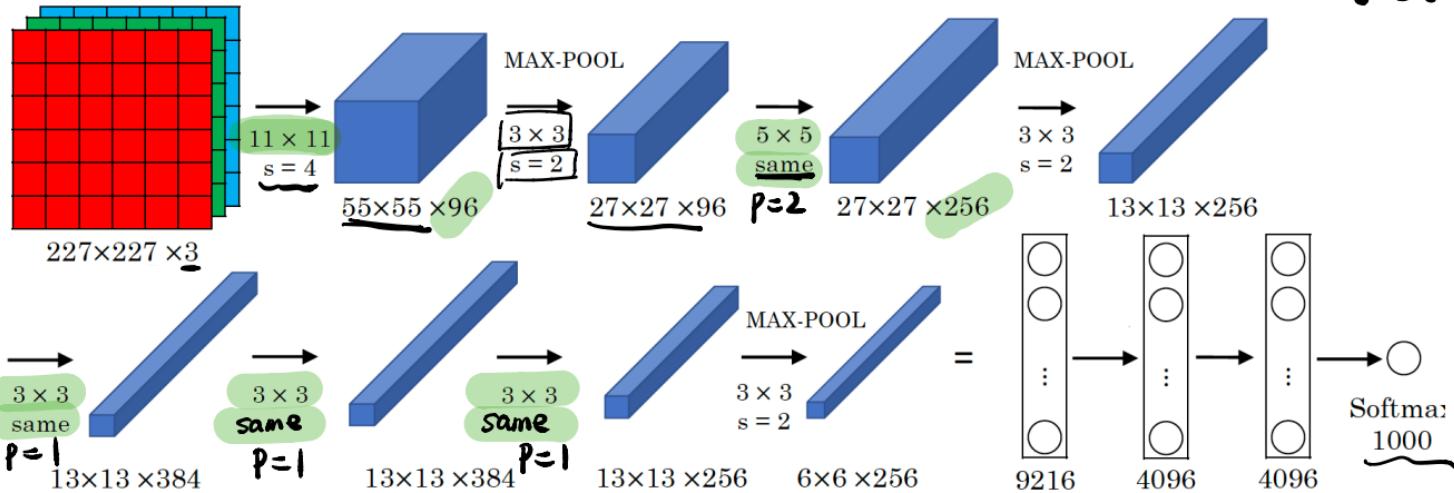
- $(n_H, n_W) \downarrow n_C \uparrow$, no padding (with padding)
- Conv pool conv pool fc fc output
- Activation function: tanh (V.S. ReLU)
- Average pooling (V.S. max pooling)

~60k parameters

AlexNet

UC San Diego

- 90 epoch for 6 days in 2 GTX 580 Gpus. • 2012 Winner: 17% top-5 error on test data



- bigger than LeNet-5
- ReLU, max pooling, padding, + Dropout

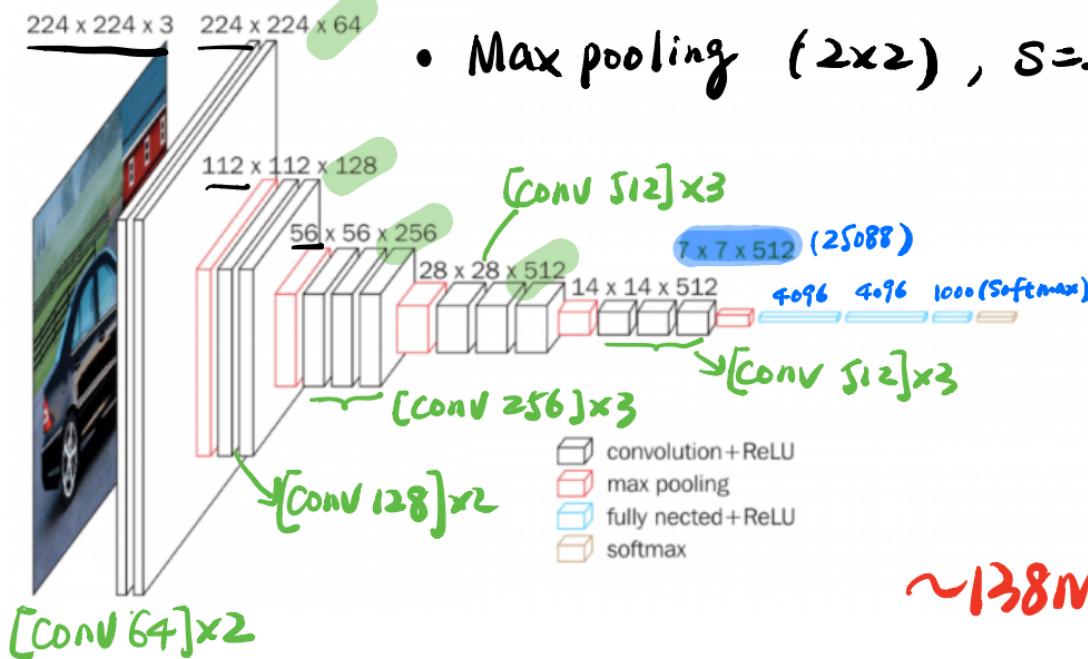
(~60m parameter)

15 conv + 3 FC layers

VGG -16

UC San Diego

1st place
for 2014
ImageNet
challenge
(7.32% top-5
error on
test data)



Residual Network (ResNet)

UC San Diego

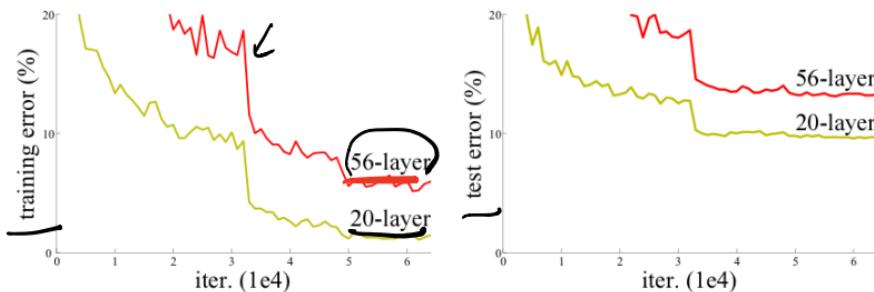


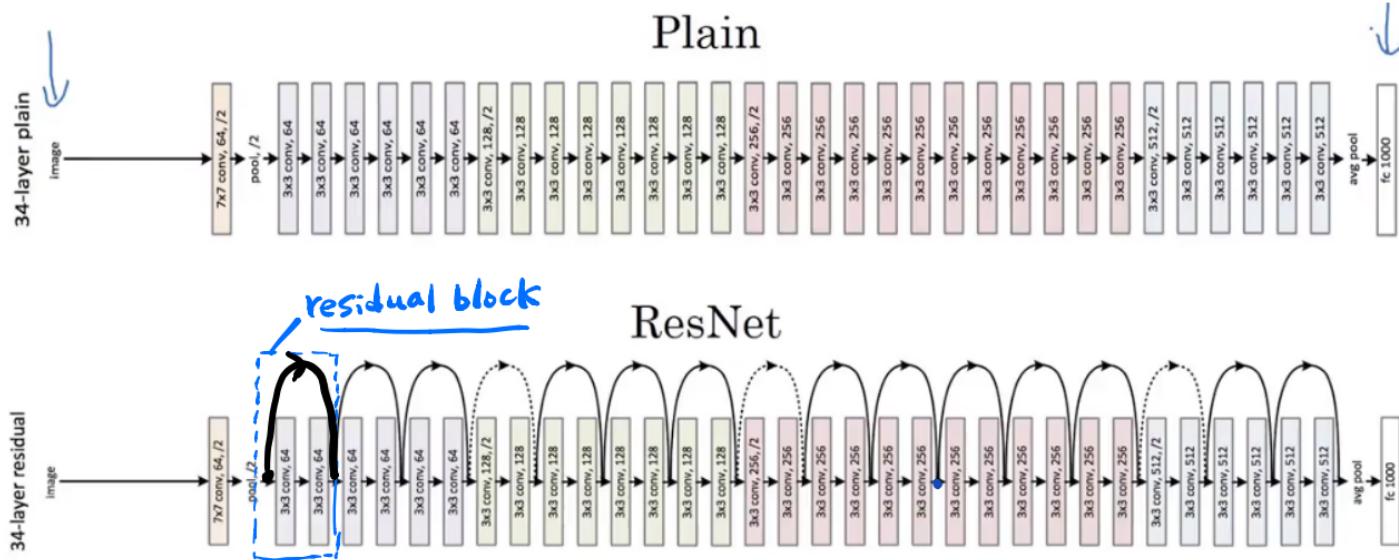
Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

- Before ResNet, deeper network does not guarantee better performance
- Gradient vanishing after many layers



Residual Network (ResNet)

UC San Diego

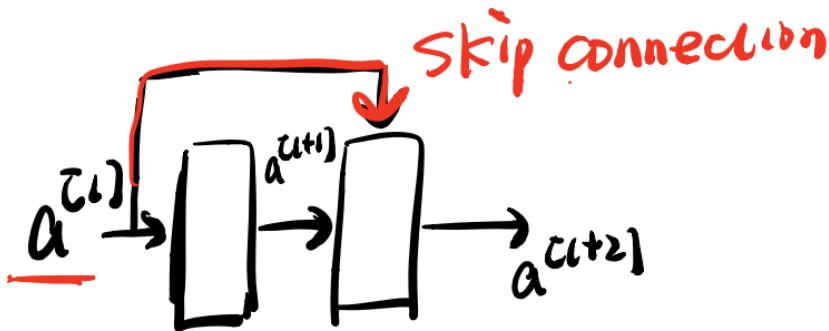


- * 1st place in 2015 ImageNet (3.57% test error (ensemble of 5 models))
in ImageNet test set using a 152-layer ResNet)
- * Object detection = 28% improvement over VGG-16

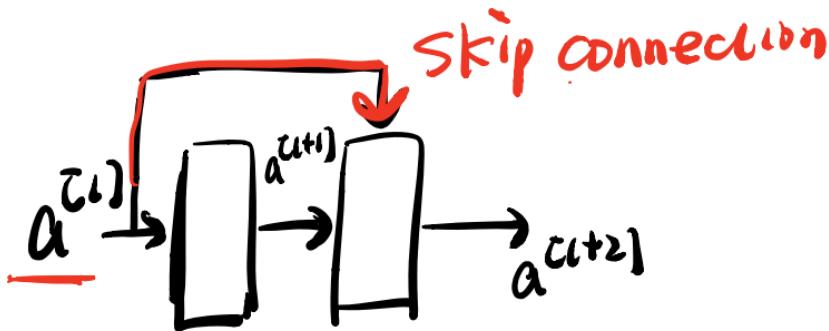
He et al., 2015. Deep Residual Learning for Image Recognition, <https://arxiv.org/abs/1512.03385>

Residual Network (ResNet)

UC San Diego



directly back-propagate
gradients to early layers!



$$\begin{aligned} z^{[l+1]} &= W^{[l+1]} a^{[l]} + b^{[l+1]} \\ a^{[l+1]} &= \sigma(z^{[l+1]}) \end{aligned}$$

$$\begin{aligned} z^{[l+2]} &= W^{[l+2]} a^{[l+1]} + b^{[l+2]} \\ a^{[l+2]} &= \sigma(z^{[l+2]}) \end{aligned}$$

↓
Residual block

$$\begin{aligned} z^{[l+2]} &= W^{[l+2]} a^{[l+1]} + b^{[l+2]} \\ a^{[l+2]} &= \sigma(z^{[l+2]} + a^{[l]}) \end{aligned}$$

Residual Network (ResNet)

UC San Diego

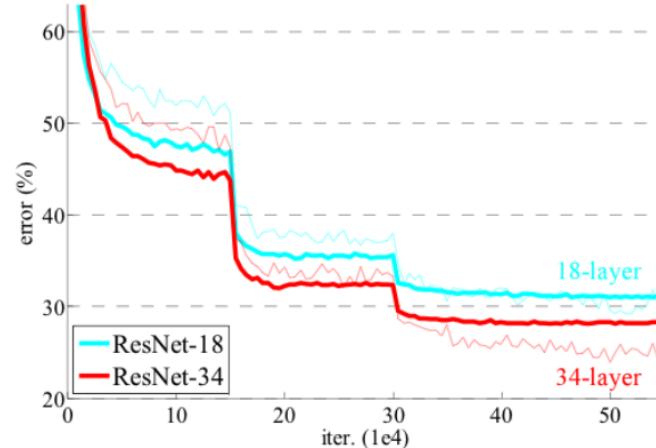
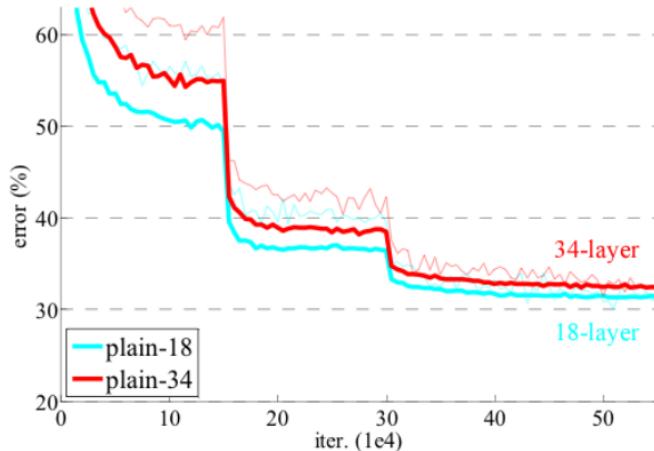


Figure 4. Training on **ImageNet**. Thin curves denote training error, and bold curves denote validation error of the center crops. Left: plain networks of 18 and 34 layers. Right: ResNets of 18 and 34 layers. In this plot, the residual networks have no extra parameter compared to their plain counterparts.

* ResNet enables training very deep neural networks.

Practical Tips

UC San Diego

- Use the classic architectures and architecture that works best on ImageNet
- Often times, you can download a pretrained model and finetune/do transfer learning with your own data (e.g., medical imaging or some applications with limited amount of data).

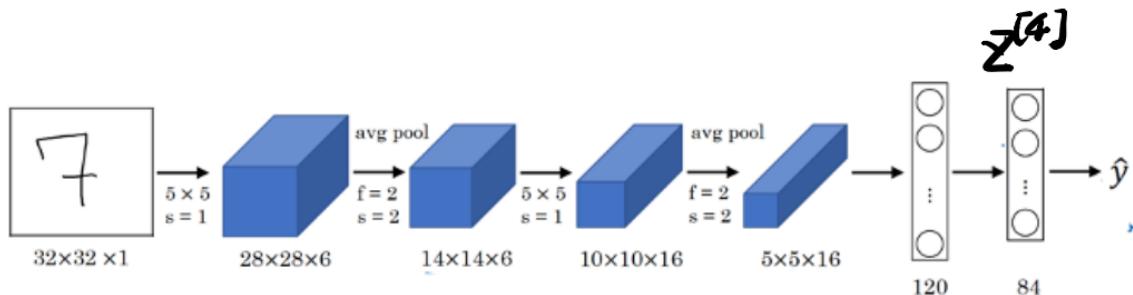
We provide pre-trained models, using the PyTorch `torch.utils.model_zoo`. These can be constructed by passing `pretrained=True`:

```
import torchvision.models as models
resnet18 = models.resnet18(pretrained=True)
alexnet = models.alexnet(pretrained=True)
squeezenet = models.squeezezenet1_0(pretrained=True)
vgg16 = models.vgg16(pretrained=True)
densenet = models.densenet161(pretrained=True)
inception = models.inception_v3(pretrained=True)
googlenet = models.googlenet(pretrained=True)
shufflenet = models.shufflenet_v2_x1_0(pretrained=True)
mobilenet_v2 = models.mobilenet_v2(pretrained=True)
```

Multi-class Classification

UC San Diego

3 6 8
6 7 5
2 1 7
4 8 1
7 6 1
7 5 9
2 2 2
0 2 3
0 1 4
7 1 2



$$\mathbf{z}^{[4]} \in \mathbb{R}^{84} \quad \mathbf{w}^{[5]} \in \mathbb{R}^{10 \times 84}$$

$$\mathbf{z}^{[5]} = \mathbf{w}^{[5]} \mathbf{z}^{[4]} + \mathbf{b}^{[5]}, \quad \mathbf{z}^{[5]} \in \mathbb{R}^{10}$$

$$\hat{y} = \mathbf{a}^{[5]} = \text{Softmax}(\mathbf{z}^{[5]})$$

$$\text{Softmax: } \delta(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

- \mathbf{z} : input vector to softmax
- K: number of elements (class) in \mathbf{z}

Multi-class Classification

UC San Diego

(Continue on softmax function)

$$\mathbf{z} = \begin{bmatrix} 9 \\ 5 \\ 0 \end{bmatrix} \quad e^{z_1} = e^9 \approx 8103$$
$$e^{z_2} = e^5 \approx 148$$

$$e^{z_3} = e^0 = 1$$

$$\sum_{j=1}^K e^{z_j} = e^{z_1} + e^{z_2} + e^{z_3} = 8252$$

$$\text{Softmax}, \sigma(\mathbf{z})_1 = \frac{8103}{8252} \approx 0.982, \sigma(\mathbf{z})_2 = \frac{148}{8252} \approx 0.0179$$

$$\sigma(\mathbf{z})_3 = \frac{1}{8252} \approx 0.0001$$

Multi-class Classification

UC San Diego

Binary classification

Sigmoid activation

$$\hat{y} = a^{[L]} = \text{sigmoid}(z^{[L]}) \\ = \frac{1}{1 + e^{-z^{[L]}}}, z^{[L]} \in \mathbb{R}$$

LOSS:

$$L = -[y \log \hat{y} + (1-y) \log(1-\hat{y})]$$

Multi-class classification

Softmax activation

$$\hat{y} = a^{[L]} = \text{softmax}(z^{[L]}) \\ = \frac{e^{z_i^{[L]}}}{\sum_{j=1}^K e^{z_j^{[L]}}}, z^{[L]} \in \mathbb{R}^K$$

LOSS

$$L = -\sum_{j=1}^K y_j \log \hat{y}_j$$