# cs231n

Lecture 4,5 summary

# review



$$score = f(x; W)$$

$$L = \frac{1}{N}\sum L_i(f, y_i)$$

loss

$$L = \frac{1}{N}\sum L_i(f, y_i) + R(W)$$

R

$$R(W)$$

# review



J(W)

Initial weights

Global minimum
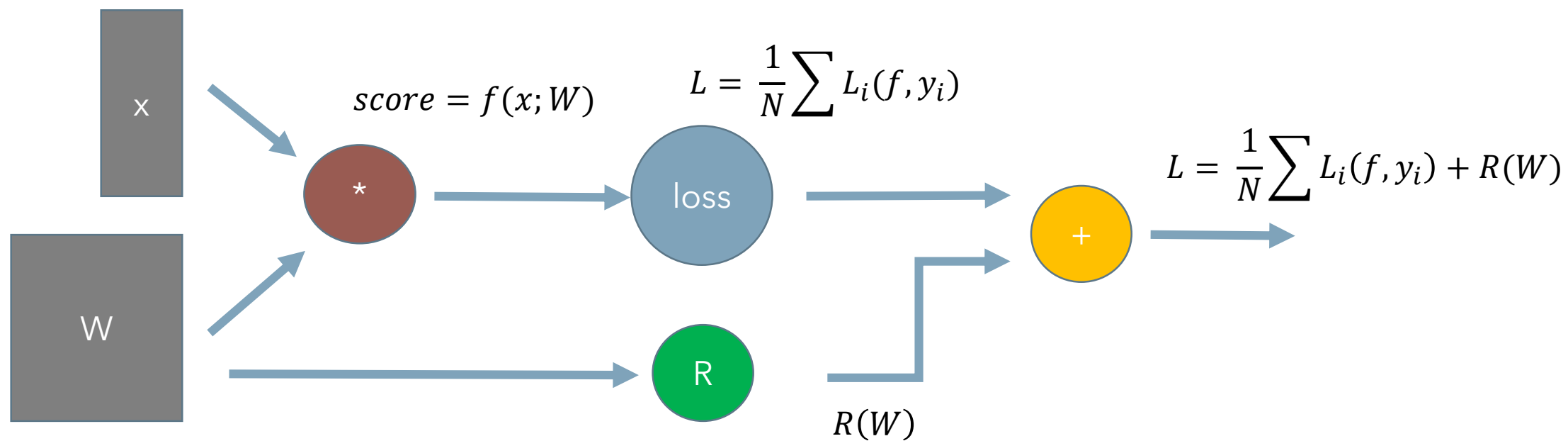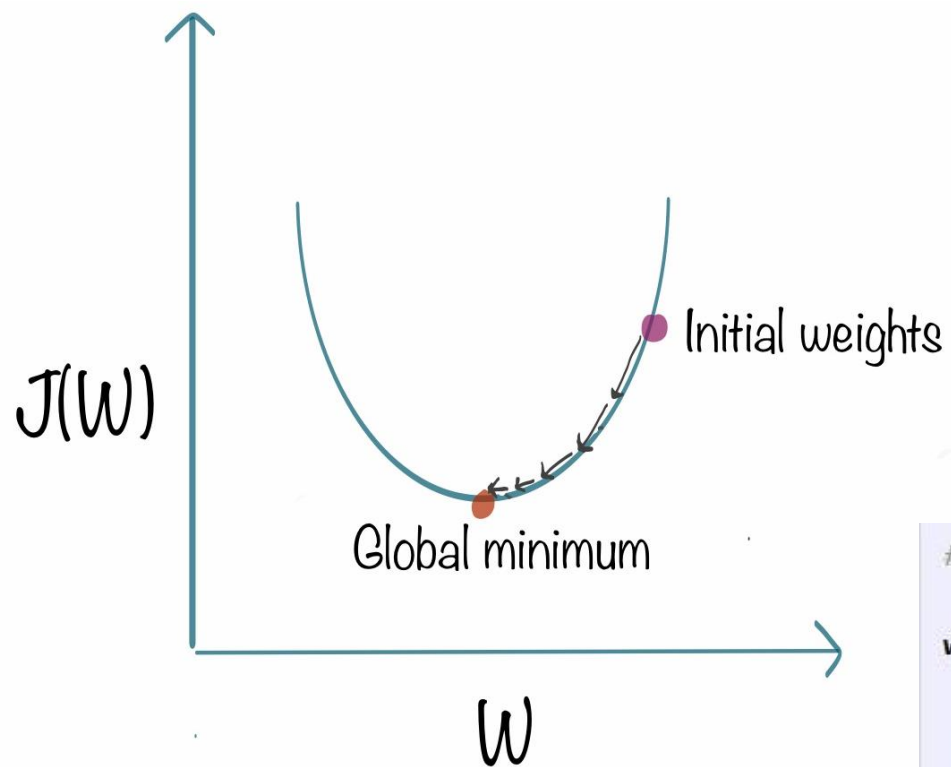
W

## Stochastic Gradient Descent (SGD)

$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(x_i, y_i, W) + \lambda R(W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^{N} \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

```
# Vanilla Minibatch Gradient Descent

while True:
    data_batch = sample_training_data(data, 256) # sample 256 examples
    weights_grad = evaluate_gradient(loss_fun, data_batch, weights)
    weights += - step_size * weights_grad # perform parameter update
```

# backpropagation

- How to update?
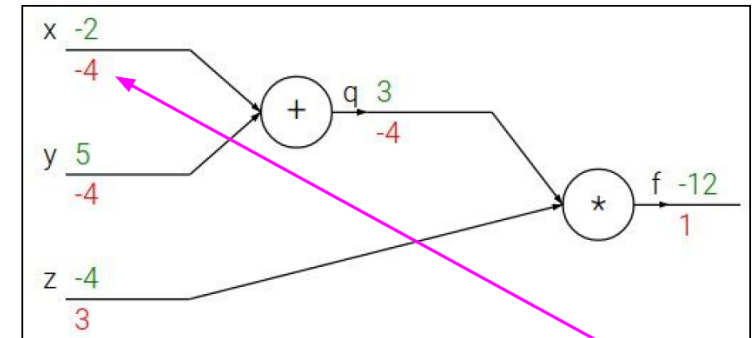
- A:using chain rule

Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$
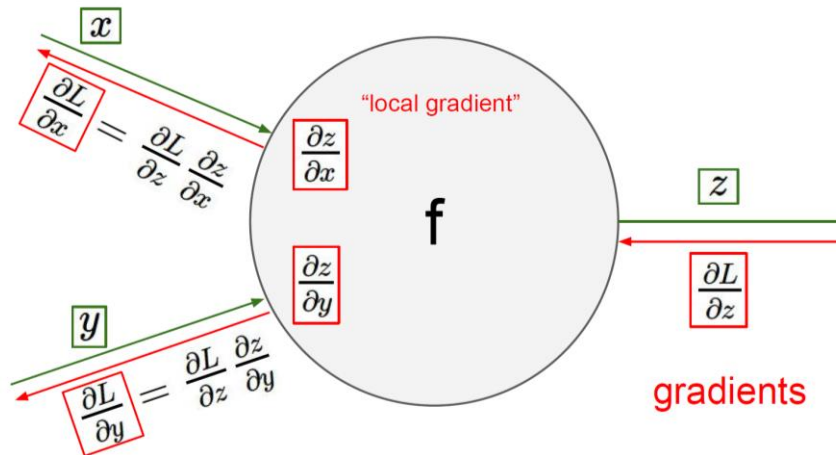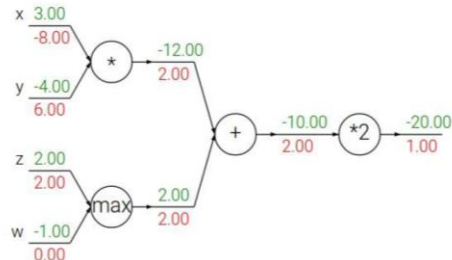


Chain rule:

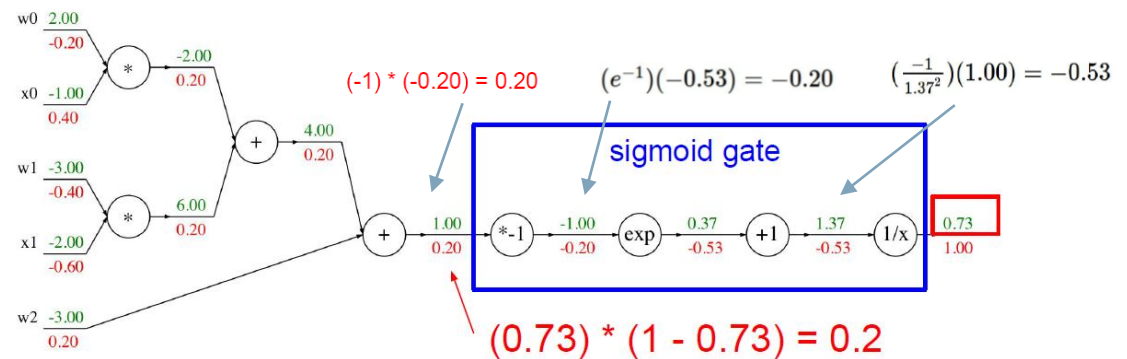$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

$$\frac{\partial f}{\partial x}$$

# backpropagation

# Neural Network

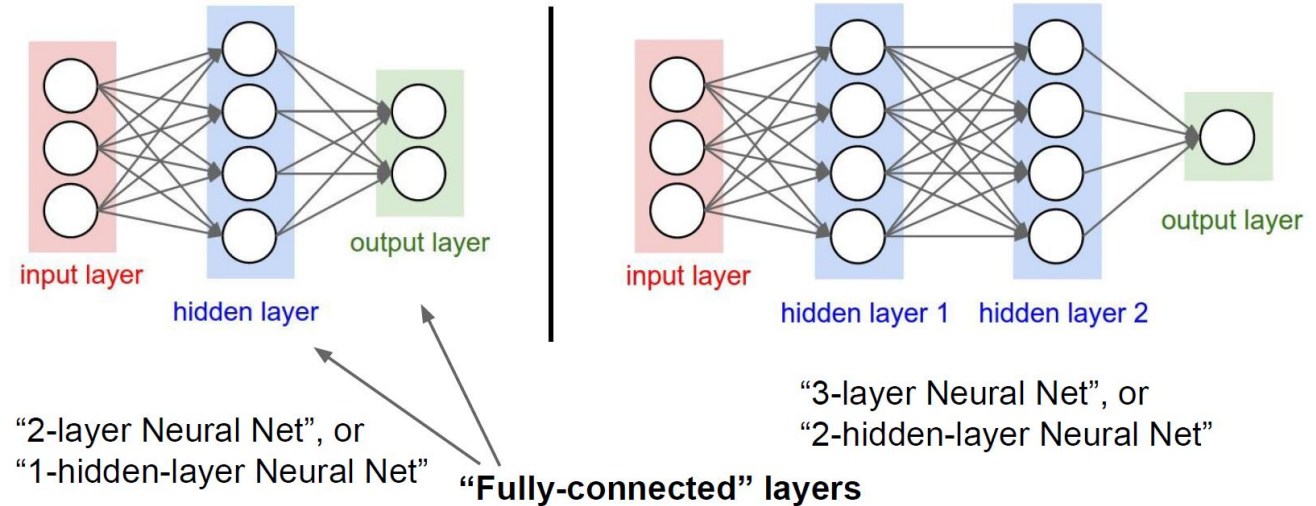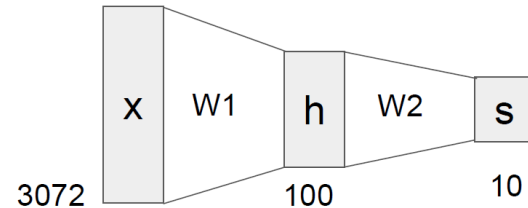

Neural networks: without the brain stuff

(**Before**) Linear score function: $f = Wx$

(**Now**) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$



3072    x   W1   h   W2   s    100    10



input layer

hidden layer

output layer

"2-layer Neural Net", or
"1-hidden-layer Neural Net"

**"Fully-connected" layers**

input layer

hidden layer 1   hidden layer 2

output layer

"3-layer Neural Net", or
"2-hidden-layer Neural Net"
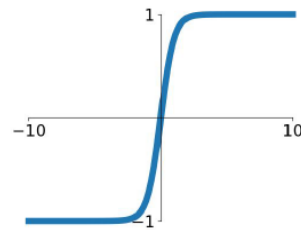
# Activation function
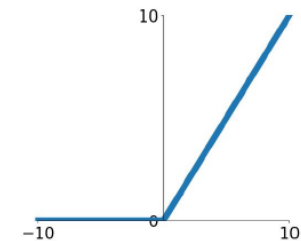
**Sigmoid**

$$\sigma(x) = \frac{1}{1+e^{-x}}$$
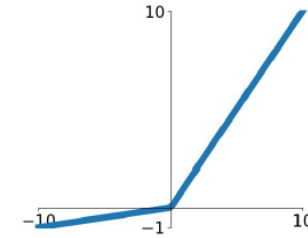
**tanh**

$$\tanh(x)$$

**ReLU**

$$\max(0, x)$$

**Leaky ReLU**

$$\max(0.1x, x)$$

**Maxout**

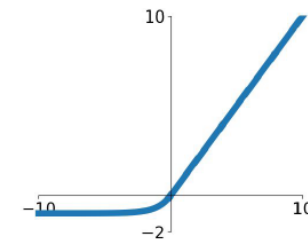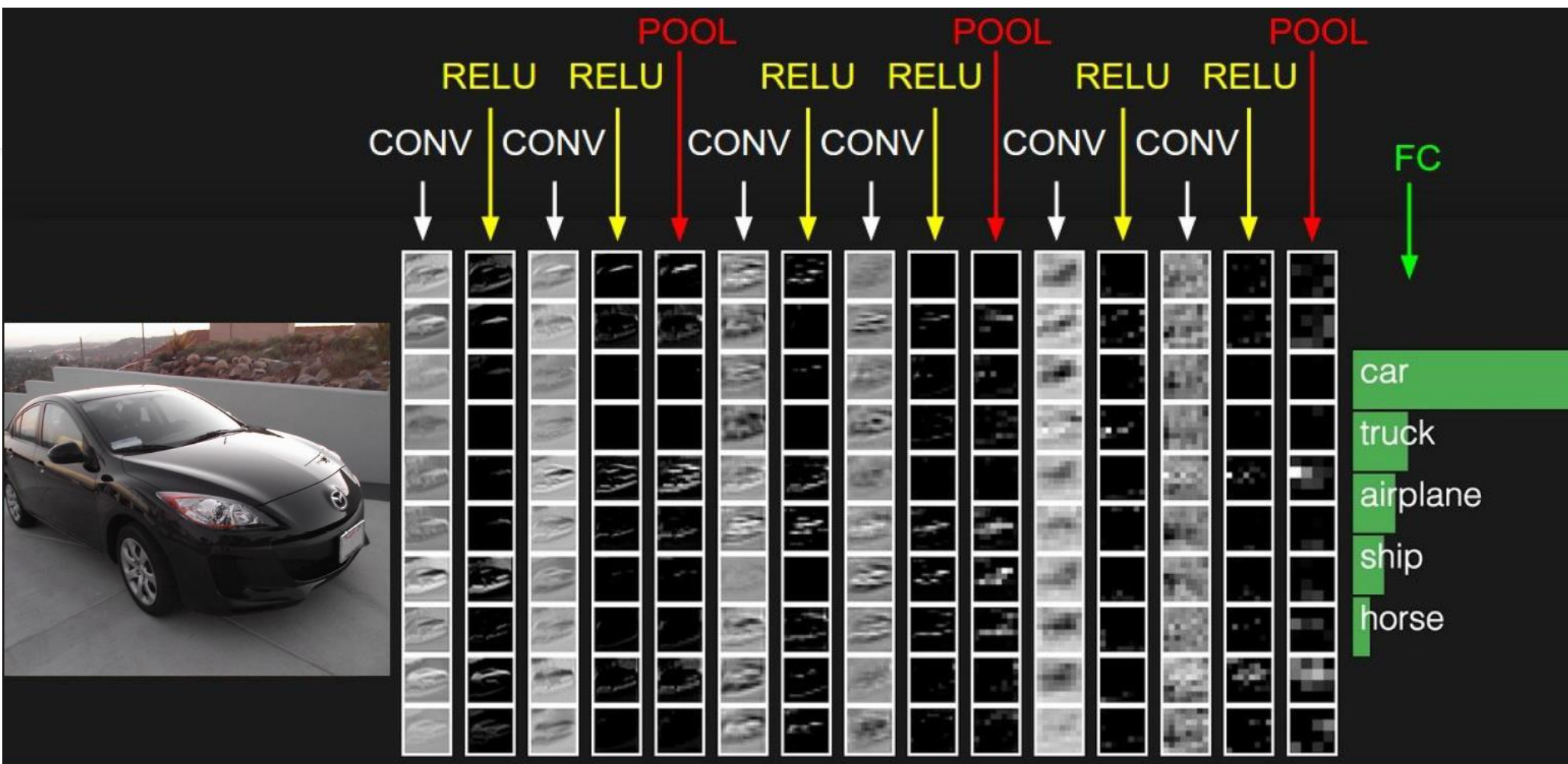$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

**ELU**

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

# Convolutional NN

## Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1

**input**

1 _____ 3072

$Wx$

10 x 3072 weights

**activation**

1 [○] _____ 10

**1 number:**
the result of taking a dot product between a row of W and the input (a 3072-dimensional dot product)

## Convolution Layer

32x32x3 image

Filters always extend the full depth of the input volume

5x5x3 filter

32

32

3

**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"

WHY? Preserve spatial inform.(공간정보)

# Convolutional NN

32x32x3 image
5x5x3 filter

32

**activation map**

28

convolve (slide) over all
spatial locations

32

**1 number:**
the result of taking a dot product between the
filter and a small 5x5x3 chunk of the image
(i.e. 5*5*3 = 75-dimensional dot product + bias)

28

3

1

$$w^T x + b$$

# Stride

STRIDE = 1

STRIDE = 2

STRIDE = 3 는 맞지 않아 쓸 수 없다!

# Padding

## In practice: Common to zero pad the border

| 0 | 0 | 0 | 0 | 0 | 0 | | | |
|---|---|---|---|---|---|---|---|---|
| 0 | | | | | | | | |
| 0 | | | | | | | | |
| 0 | | | | | | | | |
| 0 | | | | | | | | |
| | | | | | | | | |
| | | | | | | | | |
| | | | | | | | | |
| | | | | | | | | |

e.g. input 7x7
**3x3** filter, applied with **stride 1**
**pad with 1 pixel** border => what is the output?

**7x7 output!**
in general, common to see CONV layers with
stride 1, filters of size FxF, and zero-padding with
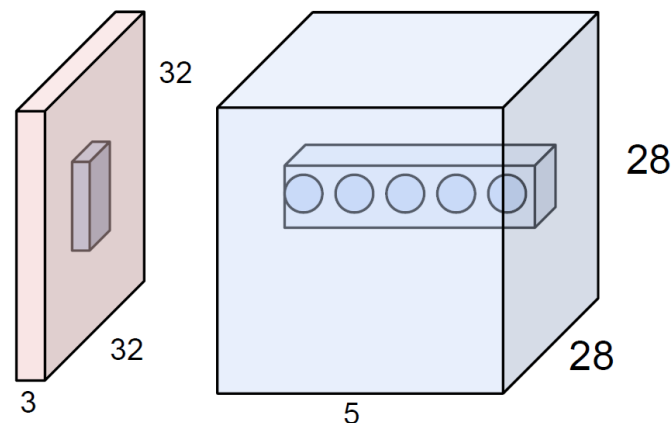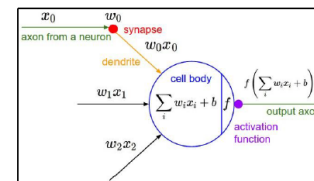(F-1)/2. (will preserve size spatially)
e.g. F = 3 => zero pad with 1
    F = 5 => zero pad with 2
    F = 7 => zero pad with 3

Microsoft Visual Studio 디버그 콘솔

```
7
FILETER : 3x3   STRIDE : 1
NO PADDING
1 2 3 3 3 2 1
2 4 6 6 6 4 2
3 6 9 9 9 6 3
3 6 9 9 9 6 3
3 6 9 9 9 6 3
2 4 6 6 6 4 2
1 2 3 3 3 2 1
PADDING
4 6 6 6 6 6 4
6 9 9 9 9 9 6
6 9 9 9 9 9 6
6 9 9 9 9 9 6
6 9 9 9 9 9 6
6 9 9 9 9 9 6
6 9 9 9 9 9 6
4 6 6 6 6 6 4
```
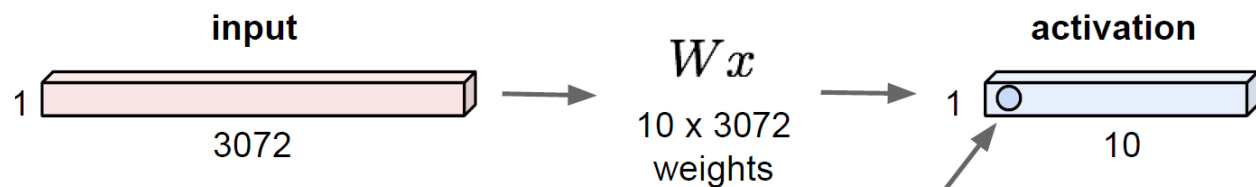
The brain/neuron view of CONV Layer



E.g. with 5 filters,
CONV layer consists of
neurons arranged in a 3D grid
(28x28x5)

There will be 5 different
neurons all looking at the same
region in the input volume
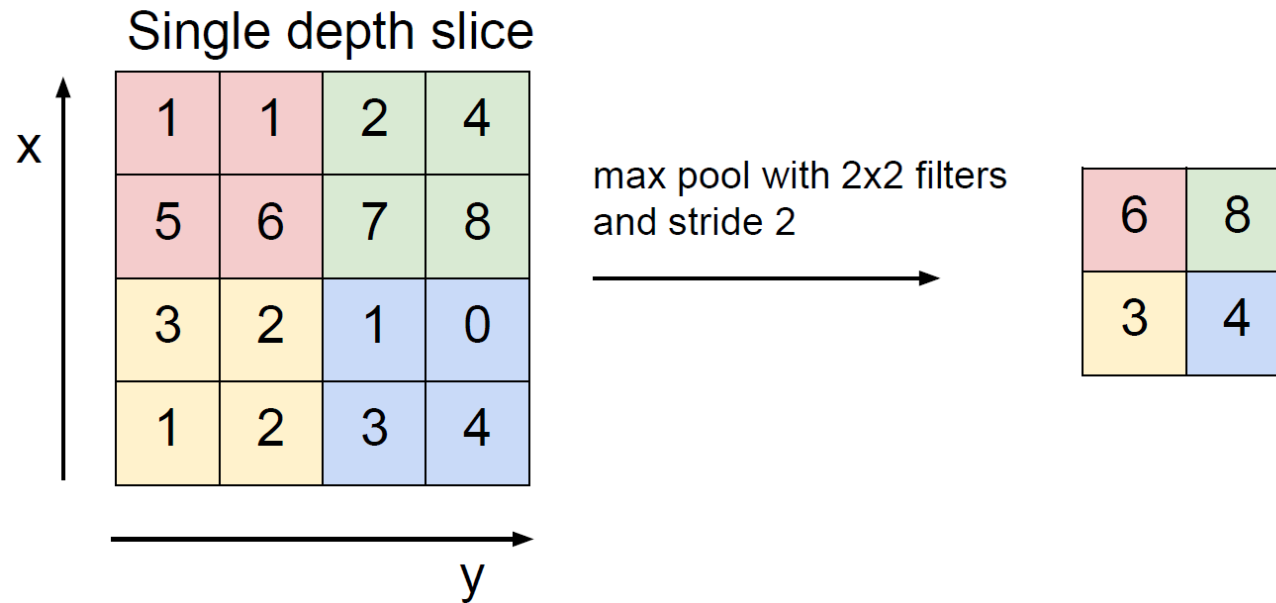
32x32x3 image -> stretch to 3072 x 1
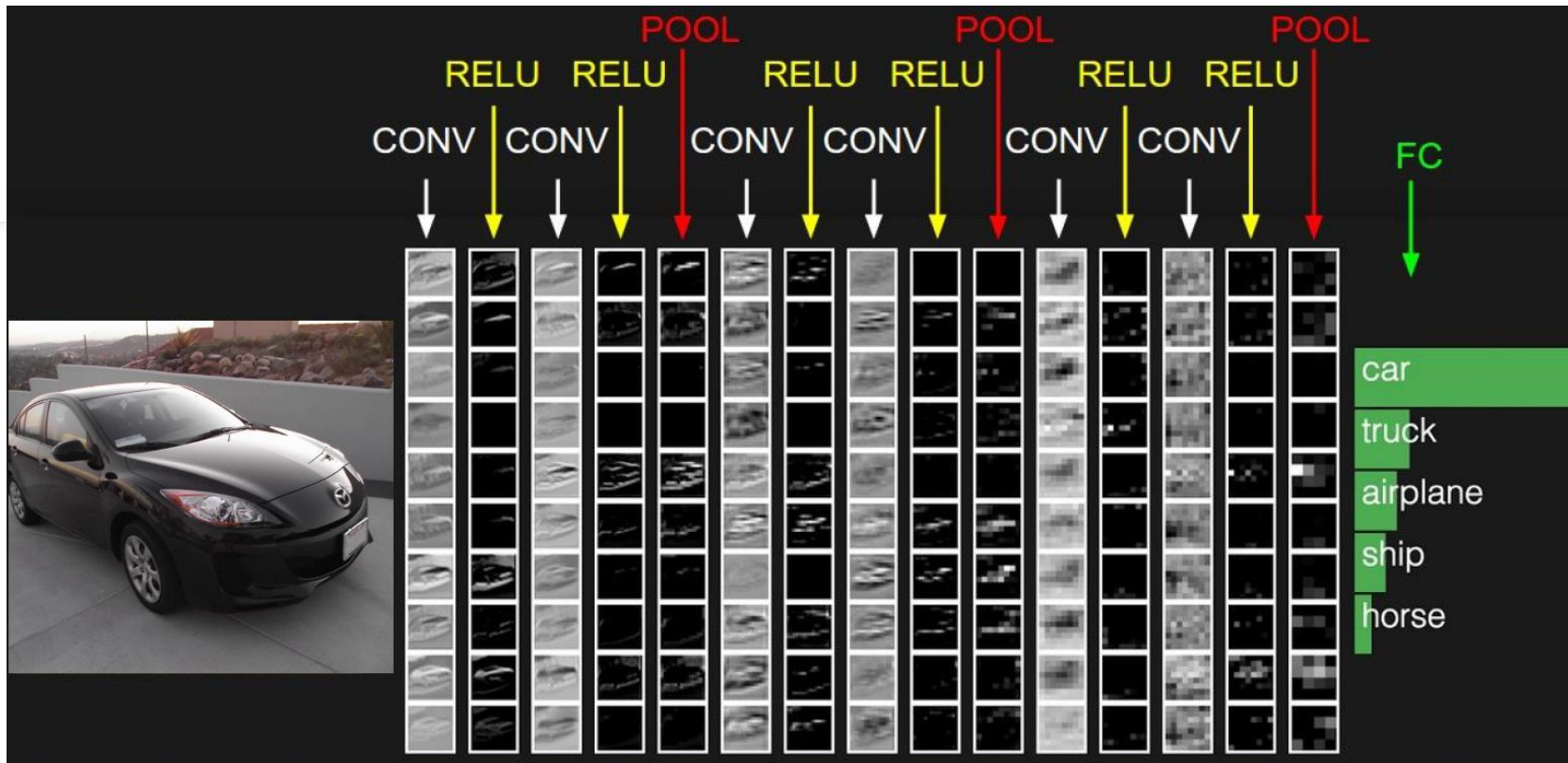
Each neuron
looks at the full
input volume

**input**

$$Wx$$

10 x 3072
weights

**activation**

1 number:
the result of taking a dot product
between a row of W and the input
(a 3072-dimensional dot product)

# Pooling ( downsampling)

## MAX POOLING

Single depth slice

마지막은 Fully connected layer를 통한 score를 계산

# Thanks