# Week 03 Lecture 05
# Artificial Neural Network

Wan Fang

Southern University of Science and Technology

# Linear Regression

- **w**
  - The **_weight_** determines the influence of each feature on our prediction, usually a vector form with $w_i$

- $b$
  - The **_bias_** says what value the predicted price should take when all features take 0

- Given a dataset, **our goal** is
  - To choose the weights **w** and bias $b$ such that on average, the predictions made based on our model best fit the true prices observed in the data.

$$\hat{y} = w_1 \cdot x_1 + \ldots + w_d \cdot x_d + b \longrightarrow \hat{y} = \mathbf{w}^T \mathbf{x} + b.$$

# Linear Regression

$$\hat{y} = w_1 \cdot x_1 + \ldots + w_d \cdot x_d + b \longrightarrow \hat{y} = \mathbf{w}^T \mathbf{x} + b$$

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w} + b$$

- Vectorization
  - All features into a vector **x** for a single data point
  - All weights into a vector **w**
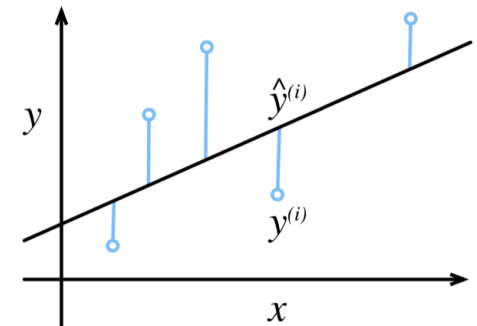  - Our entire dataset as the *design matrix* **X,** including one row for every example and one column for every feature

$$\mathbf{X} = \begin{bmatrix} x_1^{(1)} & \cdots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(i)} & \cdots & x_d^{(i)} \end{bmatrix} \quad \text{one row for every example}$$

one column
for every feature

# Loss Function

- To quantify the distance between the **predicted** and **real** value.
  - usually be a non-negative number where smaller values are better
  - perfect predictions incur a loss of 0

- The Sum of Squared Errors $l^{(i)}(\mathbf{w}, b) = \frac{1}{2}\left(\hat{y}^{(i)} - y^{(i)}\right)^2$
  - the empirical error is only
    a function of the model parameters

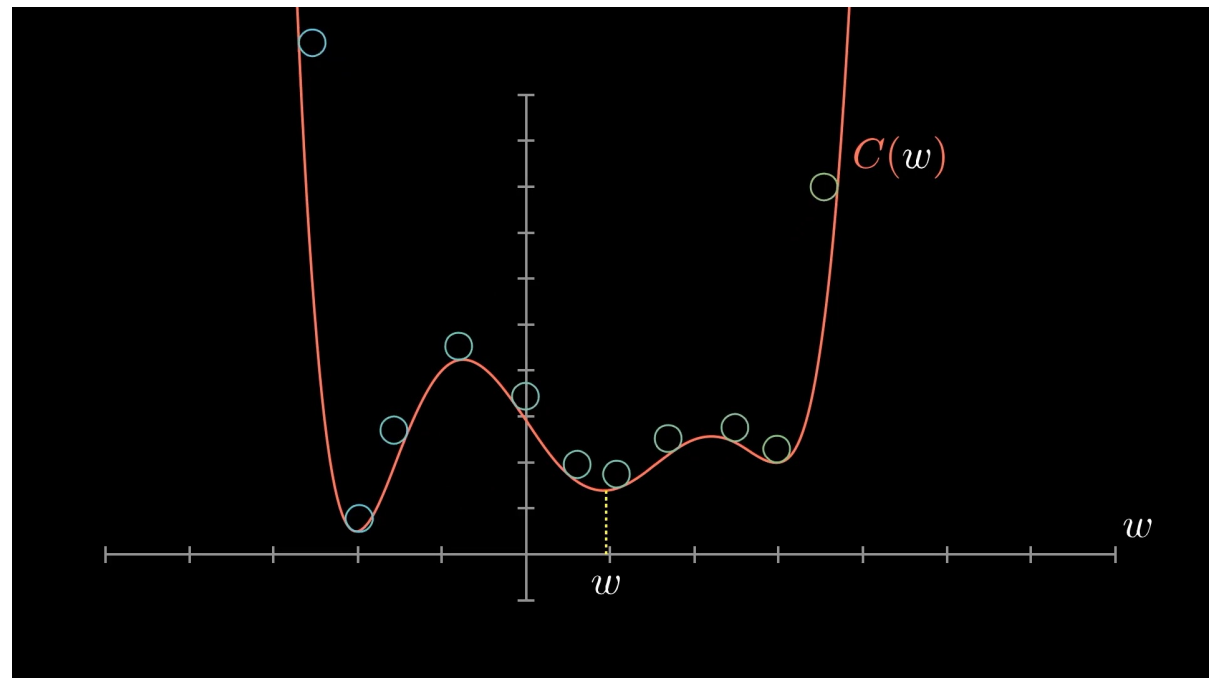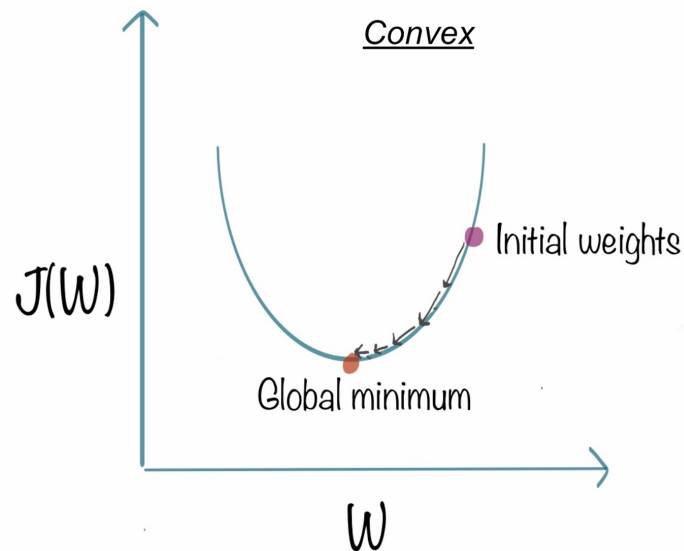- Loss Function as an averaged SSE

$$L(\mathbf{w}, b) = \frac{1}{n}\sum_{i=1}^{n} l^{(i)}(\mathbf{w}, b) = \frac{1}{n}\sum_{i=1}^{n} \frac{1}{2}\left(\mathbf{w}^\top \mathbf{x}^{(i)} + b - y^{(i)}\right)^2$$

$$\mathbf{w}^*, b^* = \underset{\mathbf{w}, b}{\operatorname{argmin}}\ L(\mathbf{w}, b)$$

# Gradient Descent

- **Iteratively reducing** the error by updating the parameters in the direction that incrementally lowers the loss function
  - On *convex* loss surfaces, it will eventually converge to a global minimum
  - For *nonconvex* surfaces, it will at least lead towards a (hopefully good) local minimum.



- The key technique for optimizing *nearly any* deep learning model

# Linear Classification

- Hypothesis
  - Acceptance depending on Test and Grade

- Data
  - $\left(x^{(i)}, y^{(i)}\right)$

- Input
  - $x_1^{(i)}$ as test scores and $x_2^{(i)}$ as test scores

- Output
  - $\hat{y}^{(i)}$ as a threshold decision of Accept or Reject

- Model
  - A linear boundary line to separate the data
    - $w_1 x_1 + w_2 x_2 + b = 0$
  - A threshold to activate a decision against the line
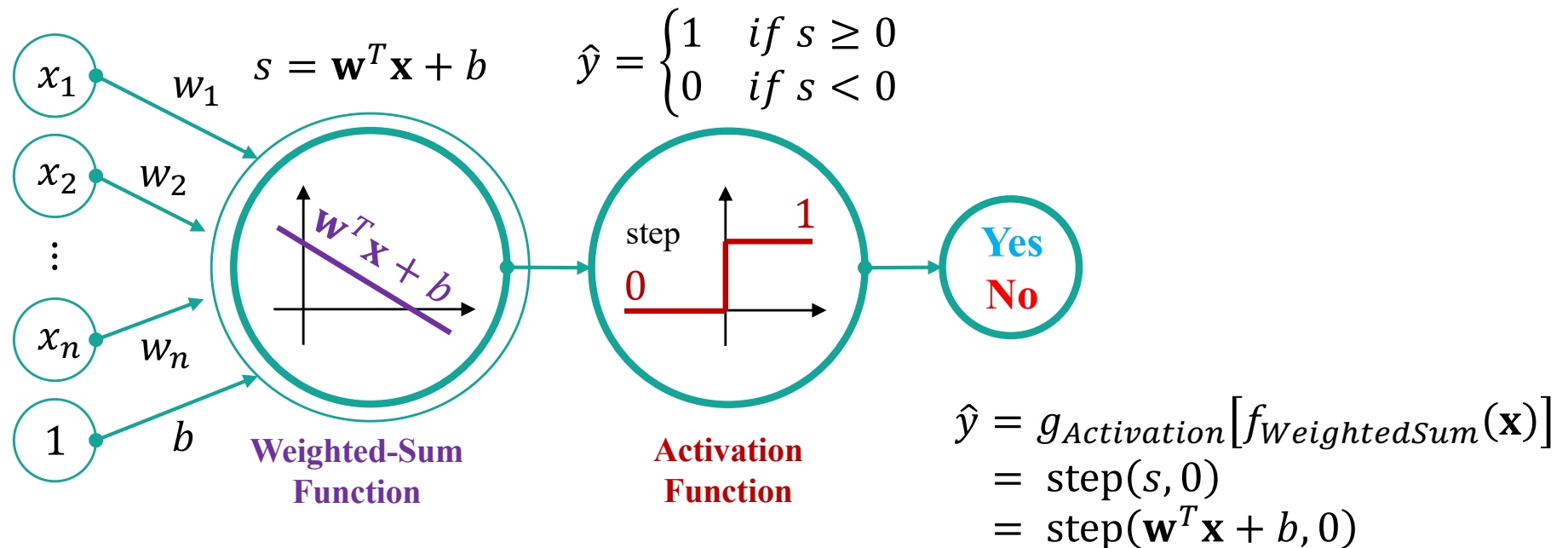    - $> 0$: Accept; $< 0$: Reject
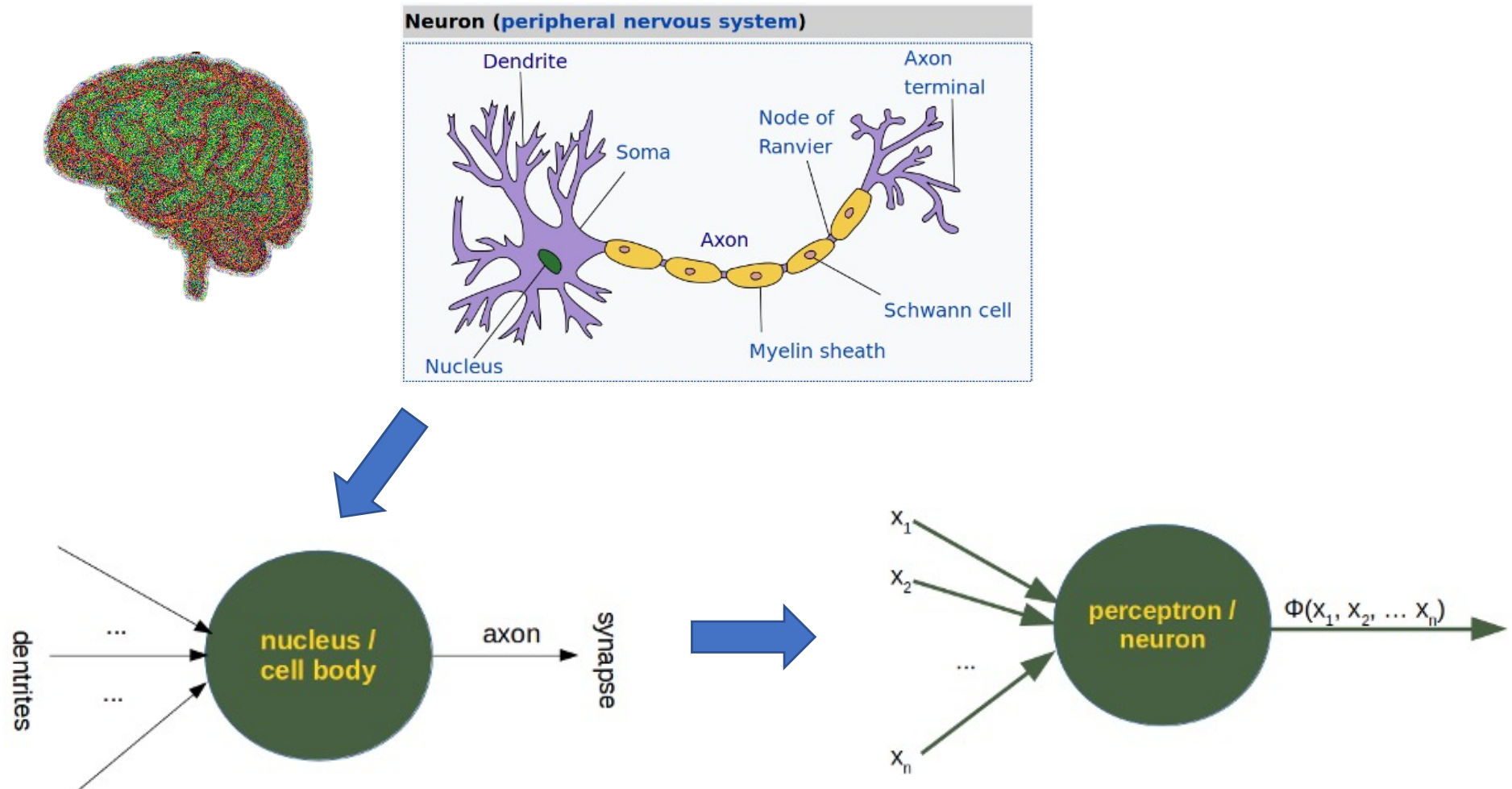
*An example of acceptance at a University*

A Linear Boundary Line of $2x_1 + x_2 - 18 = 0$

as a decision criteria from regression to classification

# Perceptron with an Activation Function

- An Artificial Neuron with two nodes
  - **Weighted-sum node**
    - Calculate a linear equation $s(x)$ with inputs on the weights plus bias
  - **Activation node**
    - Apply the step function to get the predicted result $\hat{y}(s)$

$$s = \mathbf{w}^T\mathbf{x} + b \qquad \hat{y} = \begin{cases} 1 & if\ s \geq 0 \\ 0 & if\ s < 0 \end{cases}$$



Weighted-Sum Function

Activation Function

$$\hat{y} = g_{Activation}\left[f_{WeightedSum}(\mathbf{x})\right]$$
$$= \text{step}(s, 0)$$
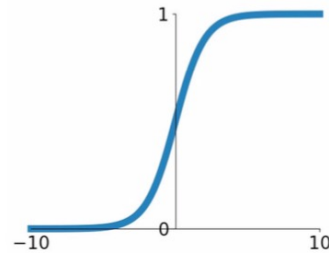$$= \text{step}(\mathbf{w}^T\mathbf{x} + b, 0)$$

# A Perceptron as an Artificial Neuron

# Activation Function
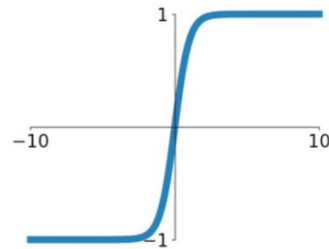
**Sigmoid**

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

**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**

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

# Limitation of Single Perceptron

*What if the boundary line is non-linear?*



- Unable to classify nonlinear scenarios

# Multi-Layer Perceptrons



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"

# Exercise I

- Build your first neural networks on the website
  - https://playground.tensorflow.org/


- Play with different data types, features, network structures. Can Neural networks separate nonlinear features?

- How does nonlinearity come?

- How important is input features

- How important is number of neurons and layers of neurons?

# Exercise II



Hello World: handwritten digits classification - MNIST

MNIST = Mixed National Institute of Standards and Technology - Download the dataset at http://yann.lecun.com/exdb/mnist/

# A Toy Example of Training a Neural Network



training digits and their labels

validation digits and their labels

**Tensors** as matrices storage of data
- [28, 28, 1]: A 28x28 pixel grayscale image (Gray)
- [128, 28, 28, 3]: A batch of 128 color images (RGB)

Loss

Accuracy

```
468/468 [==============================] - 2s 5ms/step - loss: 0.4290 - acc: 0.8862 - val_loss: 0.4003 - val_acc: 0.8963
Epoch 9/10
378/468 [======================>......] - ETA: 0s - loss: 0.4230 - acc: 0.8866
```

predictions from local fonts (bad predictions in red)

168 sample validation digits out of 10000 with bad predictions in red and sorted first

**60,000 training images**
**10,000 testing images**

# A Single-layer Network of Image Classification

Straightened into an array of features ($x$)

*784 pixels*

28x28 pixels

Input Image

softmax

0   1   2   9

Through Softmax operation to predict 10 labels ($y$) with a normalized distrution

*weighted sum of all pixels + bias*

softmax

$$softmax(L_n) = \frac{e^{L_n}}{||e^L||}$$

*neuron outputs*

Vectorized dataset into a **batch** (of 100 images) for efficient operation

*X: 100 images, one per line, flattened*

*10 columns*

$$
\begin{array}{cccccc}
W_{0,0} & W_{0,1} & W_{0,2} & W_{0,3} & \cdots & W_{0,9} \\
W_{1,0} & W_{1,1} & W_{1,2} & W_{1,3} & \cdots & W_{1,9} \\
W_{2,0} & W_{2,1} & W_{2,2} & W_{2,3} & \cdots & W_{2,9} \\
W_{3,0} & W_{3,1} & W_{3,2} & W_{3,3} & \cdots & W_{3,9} \\
W_{4,0} & W_{4,1} & W_{4,2} & W_{4,3} & \cdots & W_{4,9} \\
W_{5,0} & W_{5,1} & W_{5,2} & W_{5,3} & \cdots & W_{5,9} \\
W_{6,0} & W_{6,1} & W_{6,2} & W_{6,3} & \cdots & W_{6,9} \\
W_{7,0} & W_{7,1} & W_{7,2} & W_{7,3} & \cdots & W_{7,9} \\
W_{8,0} & W_{8,1} & W_{8,2} & W_{8,3} & \cdots & W_{8,9} \\
& & \cdots & & & \\
W_{783,0} & W_{783,1} & W_{783,2} & \cdots & W_{783,9}
\end{array}
$$

*784 lines*

*broadcast*

$$L = X.W + b$$

8 4 6 6 4 3

*784 pixels*

$$
\begin{array}{cccccc}
L_{0,0} & L_{0,1} & L_{0,2} & L_{0,3} & \cdots & L_{0,9} \\
L_{1,0} & L_{1,1} & L_{1,2} & L_{1,3} & \cdots & L_{1,9} \\
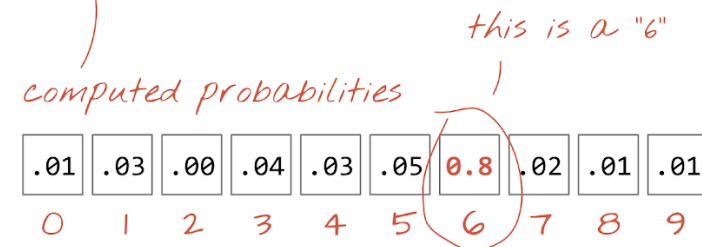L_{2,0} & L_{2,1} & L_{2,2} & L_{2,3} & \cdots & L_{2,9} \\
L_{3,0} & L_{3,1} & L_{3,2} & L_{3,3} & \cdots & L_{3,9} \\
L_{4,0} & L_{4,1} & L_{4,2} & L_{4,3} & \cdots & L_{4,9} \\
& & \cdots & & & \\
L_{99,0} & L_{99,1} & L_{99,2} & & \cdots & L_{99,9}
\end{array}
$$

$+$    $b_0$  $b_1$  $b_2$  $b_3$  $\cdots$  $b_9$

$+$ *Same 10 biases on all lines*

# Softmax on a Batch of Images

Predictions
Y[100, 10]

Images
X[100, 784]

Weights
W[784,10]

Biases
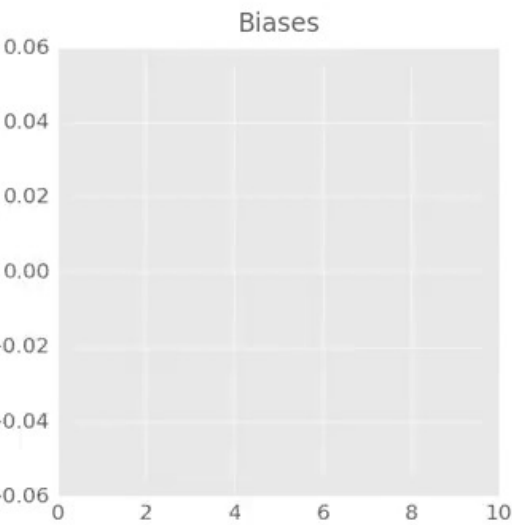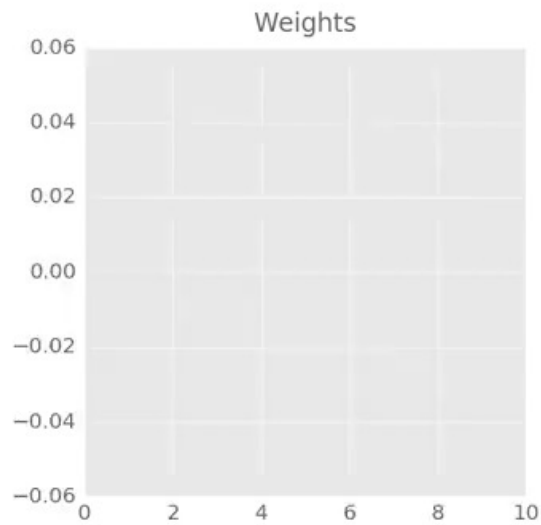b[10]

$$Y = softmax(X.W + b)$$

applied line
by line

matrix multiply

broadcast
on all lines

tensor shapes in [ ]

tensor shapes:   X[100, 784]   W[748,10]   b[10]

Y = tf.nn.softmax(tf.matmul(X, W) + b)

matrix multiply

broadcast
on all lines

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | **1** | 0 | 0 | 0 |

actual probabilities, "one-hot" encoded

$$-\sum Y_i'.log(Y_i)$$ Cross entropy

this is a "6"

computed probabilities

| .01 | .03 | .00 | .04 | .03 | .05 | **0.8** | .02 | .01 | .01 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

# Training Process

# Adding Layers



$$\frac{1}{1 + e^{-x}}$$

sigmoid

**Getting flat**
- The gradient can become very small and training get slower and slower.

Simply adding more layers with sigmoid activations does not give us the expected results …

# Special Care for Deep Networks

784

200 relu

100 relu

60 relu

30 relu

10 softmax

0 1 2 ... 9

$$Y = \textbf{relu}(\sum_i \textbf{W}_i X_i + \textbf{b})$$

inputs
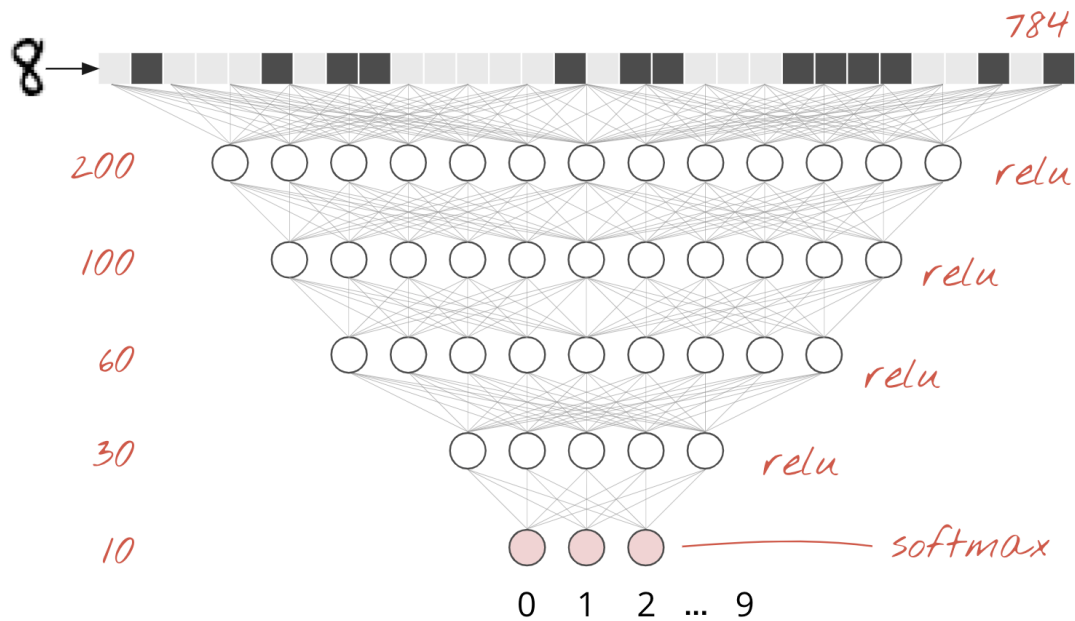
activation

weights

bias

relu

cross-entropy

**HANDS ON:**

Replace the `'sgd'` optimizer with a better one, for example `'adam'` and train again.

97%

**HANDS ON:**

Replace all `activation='sigmoid'` with `activation='relu'` in your layers and train again.

https://github.com/GoogleCloudPlatform/tensorflow-without-a-phd/blob/master/tensorflow-mnist-tutorial/keras_02_mnist_dense.ipynb

- DS323: AI in Design

# Convolutional Networks

# Convolutional Networks Applications

# A Design Challenge with Increasing Dimensions



$512 \times 512 \times 3 = 765,432$

Gray-scaled     MNIST            CIFAR-10     Color (RGB)

$28 \times 28 \times 1 = 784$      **Input** *Features*      $32 \times 32 \times 3 = 3{,}072$

Weighted-Sum
$(10 \times 784 \leftarrow W \rightarrow 10 \times 3072)$

Activation
(?)

**Output** *Label*

0    1    2      9         airplane   cat    frog      ship

Regular Neural Nets don't scale well to full images

# Convolutional Operation

**Input**

| | | | |
|---|---|---|---|
| a | b | c | d |
| e | f | g | h |
| i | j | k | l |

Kernel

| | |
|---|---|
| w | x |
| y | z |

**Output**

| aw + bx + ey + fz | bw + **cx** + fy + gz | cw + dx + gy + hz |
|---|---|---|
| ew + fx + iy + jz | fw + gx + jy + kz | gw + hx + ky + lz |

An Image
with 3
channels
of RGB

Input Volume (+pad 1) (7x7x3)

x[:,:,0]

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

x[:,:,1]

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

x[:,:,2]

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

Filter W0 (3x3x3)

w0[:,:,0]

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

-2

w0[:,:,1]

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

1

w0[:,:,2]

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

5

Bias b0 (1x1x1)

b0[:,:,0]

| 1 |
|---|

1

Filter W1 (3x3x3)

w1[:,:,0]

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

w1[:,:,1]

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

w1[:,:,2]

| -1 | 1 | -1 |
|---|---|---|
| -1 | -1 | 1 |
| -1 | -1 | 1 |

Bias b1 (1x1x1)

b1[:,:,0]

| 0 |
|---|

Output Volume (3x3x2)

o[:,:,0]

| 5 | 5 | 0 |
|---|---|---|
| 0 | -5 | -4 |
| -4 | -5 | -2 |

-2+1+5+1

o[:,:,1]

| -6 | -8 | -3 |
|---|---|---|
| -6 | -6 | -1 |
| -5 | -3 | -1 |

toggle movement

**Motivations**
- Sparse interactions
- Parameter sharing
- Equivariant representations
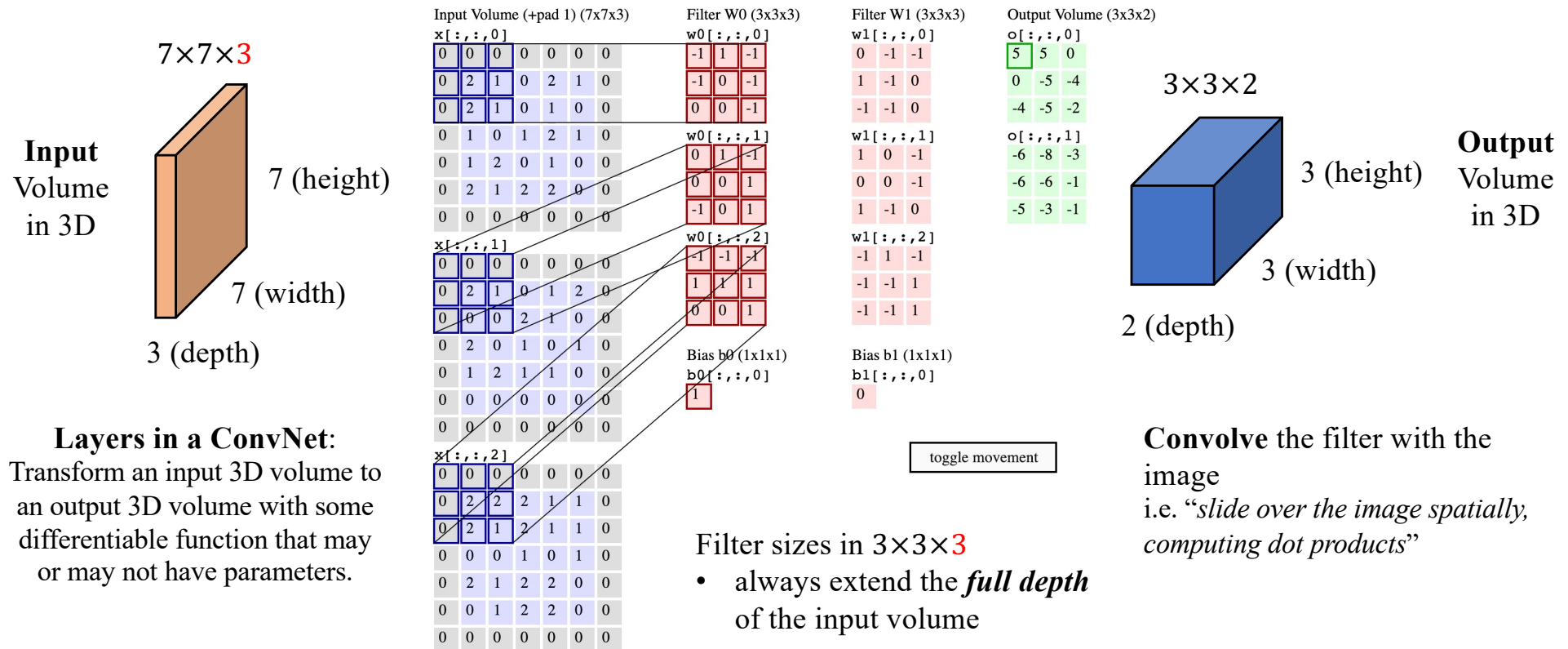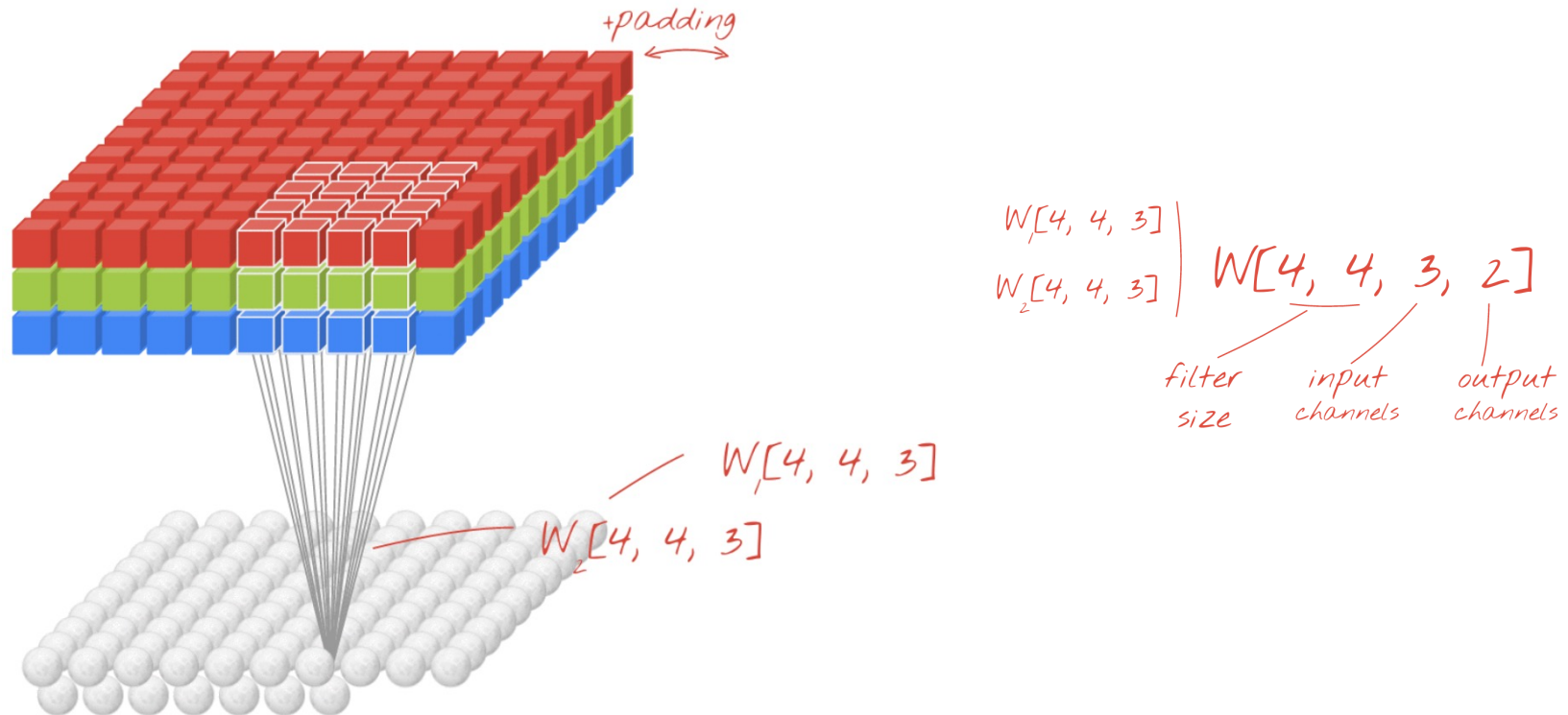
# Convolution in 3D Volumes

Preserved spatial structure between the input and output volumes in width, height, number of channels

$7 \times 7 \times 3$

**Input** Volume in 3D

7 (height)

7 (width)

3 (depth)

**Layers in a ConvNet**:
Transform an input 3D volume to an output 3D volume with some differentiable function that may or may not have parameters.



Input Volume (+pad 1) (7x7x3)
x[:,:,0]

x[:,:,1]

x[:,:,2]

Filter W0 (3x3x3)
w0[:,:,0]

w0[:,:,1]

w0[:,:,2]

Bias b0 (1x1x1)
b0[:,:,0]
1

Filter W1 (3x3x3)
w1[:,:,0]

w1[:,:,1]

w1[:,:,2]

Bias b1 (1x1x1)
b1[:,:,0]
0

Output Volume (3x3x2)
o[:,:,0]
| 5 | 5 | 0 |
| 0 | -5 | -4 |
| -4 | -5 | -2 |

o[:,:,1]
| -6 | -8 | -3 |
| -6 | -6 | -1 |
| -5 | -3 | -1 |

toggle movement

Filter sizes in $3 \times 3 \times 3$
• always extend the ***full depth*** of the input volume

$3 \times 3 \times 2$

**Output** Volume in 3D

3 (height)

3 (width)

2 (depth)

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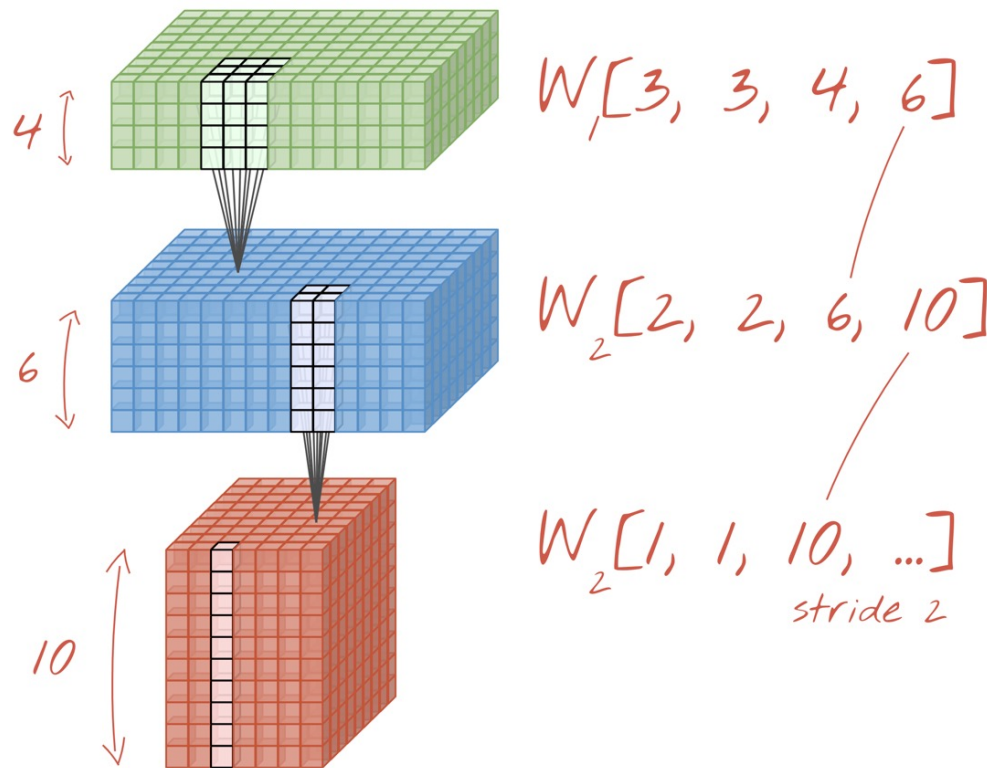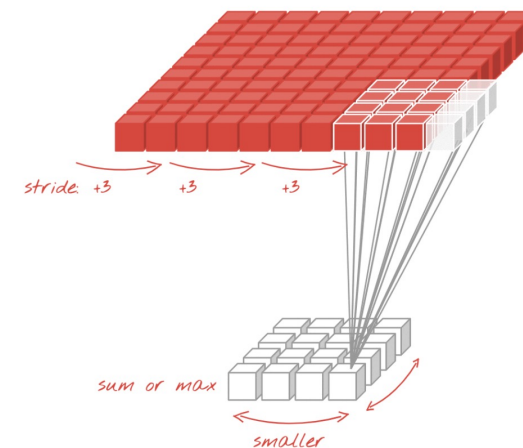i.e. "*slide over the image spatially, computing dot products*"

# The Design of a Convolutional Layer

A convolutional layer is defined by the filter (or kernel) size, the number of filters applied and the stride

+padding

$W_1[4, 4, 3]$
$W_2[4, 4, 3]$

$W[4, 4, 3, 2]$

filter size    input channels    output channels

$W_1[4, 4, 3]$
$W_2[4, 4, 3]$

# Output Volume Size



$W_1[3, 3, 4, 6]$

$W_2[2, 2, 6, 10]$

$W_2[1, 1, 10, ...]$

*stride 2*

- Depth (number of channels):
  - *adjusted by using more or fewer filters*

- Width & Height:
  - *adjusted by using a stride >1*
  - *(or with a max-pooling operation)*

stride: +3   +3   +3

sum or max

smaller

Defined by the filter (or kernel) size, the number of filters applied and the stride

# The Last Layer
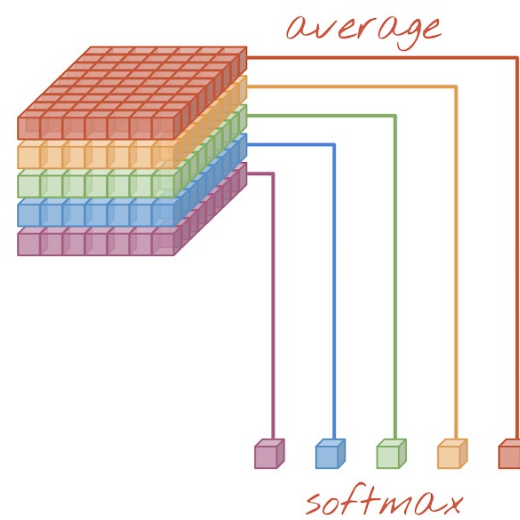
## From a Cubic Volume in 3D to predicted labels

Fully connected layer

Global average pooling

Similar like a normal neural network

Expensive in #weights

But preseves the location data (*x, y*)
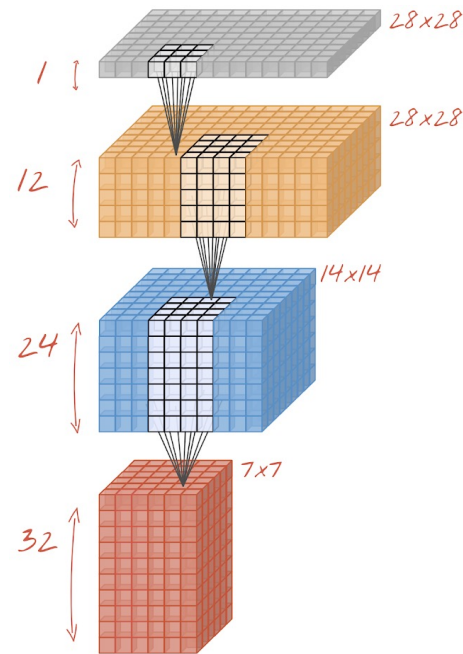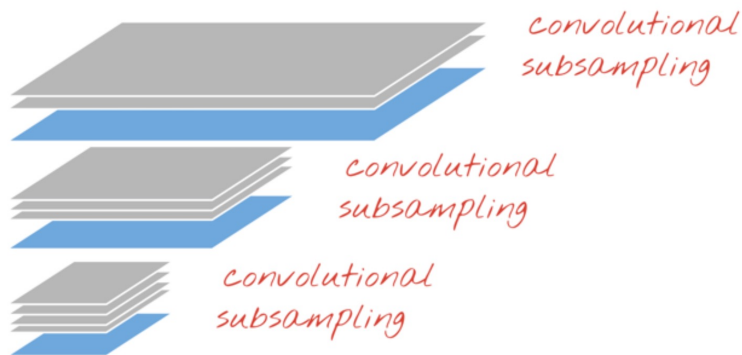
7

7

5

245

flatten

*W* [245, 5]

softmax

average

softmax

Much lighter in calculation

The average pooling explicitly discards all location data
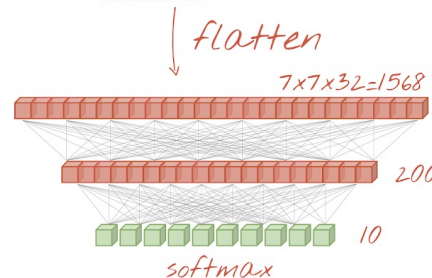
**1225** weights   *cheaper*   **0** weights

Stacking Up a ConvNet

*Layer-by-layer*



convolutional subsampling

convolutional subsampling

convolutional subsampling

28x28

1

Convolutional 3x3 filters=12
$W_1[3, 3, 1, 12]$

28x28

12

Convolutional 6x6 filters=24
$W_2[6, 6, 12, 24]$ stride 2

14x14

24

Convolutional 6x6 filters=32
$W_3[6, 6, 24, 32]$ stride 2

7x7

32

strided convolution

flatten

7x7x32=1568

Dense layer
$W_4[1568, 200]$

200

10

softmax

Softmax dense layer
$W_5[200, 10]$

# Exercise: Handwritten digits classification

- MNIST forward Neural Network:
  - https://ml4a.github.io/demos/forward_pass_mnist/
- MNIST confusion matrix
  - https://ml4a.github.io/demos/confusion_mnist/
- MNIST Convolutional Neural Network: filters
  - https://ml4a.github.io/demos/convolution_all/
- crfm.stanford.edu/2023/03/13/alpaca.html
- MNIST TensorFlow.js playground: https://cs.stanford.edu/people/karpathy/convnetjs/demo/mnist.html
- Explainable AI Demos: https://lrpserver.hhi.fraunhofer.de/image-classification

# Thank you~

Wan Fang

Southern University of Science and Technology