# Lecture 12
# Network Tuning

Song Chaoyang

Assistant Professor

Department of Mechanical and Energy Engineering

songcy@sustech.edu.cn

AncoraSIR.com

SUSTech
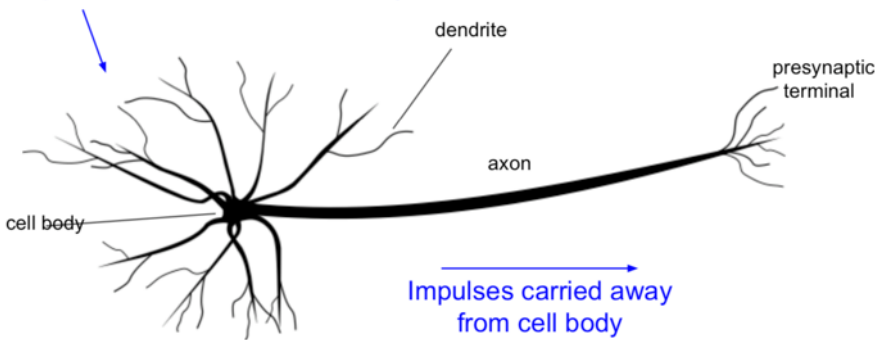Southern University
of Science and Technology

# Neural Network

## *Biological Inspiration*

Impulses carried toward cell body

dendrite

presynaptic terminal

axon

cell body

Impulses carried away from cell body

This image by Felipe Perucho is licensed under CC-BY 3.0

Linear score function:  $f = Wx$

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

x  W1  h  W2  s

3072  100  10

plane  car  bird  cat  deer  dog  frog  horse  ship  truck

## Neural networks: Architectures

input layer  hidden layer  output layer

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

input layer  hidden layer 1  hidden layer 2  output layer

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

"Fully-connected" layers

```
# forward-pass of a 3-layer neural network:
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function (use sigmoid)
x = np.random.randn(3, 1) # random input vector of three numbers (3x1)
h1 = f(np.dot(W1, x) + b1) # calculate first hidden layer activations (4x1)
h2 = f(np.dot(W2, h1) + b2) # calculate second hidden layer activations (4x1)
out = np.dot(W3, h2) + b3 # output neuron (1x1)
```

## Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1

input  $Wx$  activation

1  3072  10 x 3072 weights  1  10

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

$x_0$  $w_0$  synapse

axon from a neuron

$w_0 x_0$

dendrite

$w_1 x_1$

cell body

$\sum_i w_i x_i + b$   $f$

$f\left(\sum_i w_i x_i + b\right)$

output axon

activation function

$w_2 x_2$

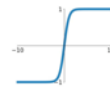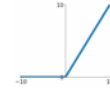## Activation functions

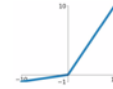**Sigmoid**  $\sigma(x) = \frac{1}{1+e^{-x}}$

**tanh**  $\tanh(x)$

**ReLU**  $\max(0, x)$
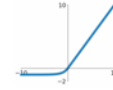
**Leaky ReLU**  $\max(0.1x, x)$

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

**ELU**  $\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$
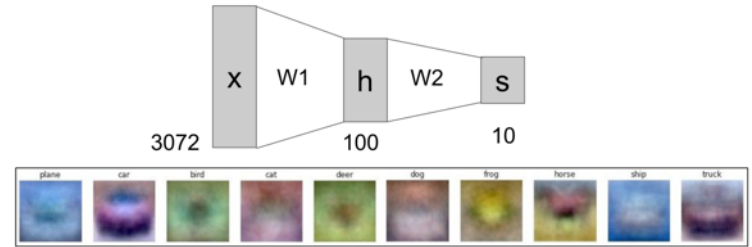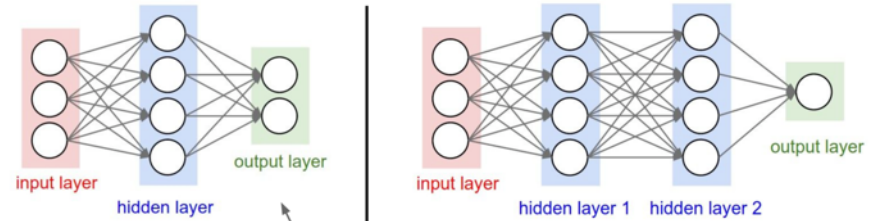
- Use ReLU. Be careful with your learning rates
- Try out Leaky ReLU / Maxout / ELU
- Try out tanh but don't expect much
- Don't use sigmoid

AncoraSIR.com

# Convolutional Neural Networks

## Convolve (slide) over all spatial locations



Filters always extend the full depth of the input volume

32x32x3 image -> preserve spatial structure

32x32x3 image
5x5x3 filter

activation maps

convolve (slide) over all spatial locations

**Preview:** ConvNet is a sequence of Convolution Layers, interspersed with activation functions

CONV, ReLU e.g. 6 5x5x3 filters

CONV, ReLU e.g. 10 5x5x6 filters

CONV, ReLU ....

Common settings:

K = (powers of 2, e.g. 32, 64, 128, 512)
- F = 3, S = 1, P = 1
- F = 5, S = 1, P = 2
- F = 5, S = 2, P = ? (whatever fits)
- F = 1, S = 1, P = 0

**Summary**. To summarize, the Conv Layer:
- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
  - Number of filters $K$,
  - their spatial extent $F$,
  - the stride $S$,
  - the amount of zero padding $P$.
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
  - $W_2 = (W_1 - F + 2P)/S + 1$
  - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. width and height are computed equally by symmetry)
  - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and $K$ biases.
- In the output volume, the $d$-th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the $d$-th filter over the input volume with a stride of $S$, and then offset by $d$-th bias.

## Pooling layer
- makes the representations smaller and more manageable
- operates over each activation map independently:

224x224x64

pool

112x112x64

224

downsampling

112

**MAX POOLING**

Single depth slice

max pool with 2x2 filters and stride 2

- ConvNets stack CONV,POOL,FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Typical architectures look like
  **[(CONV-RELU)*N-POOL?]*M-(FC-RELU)*K,SOFTMAX**
  where N is usually up to ~5, M is large, 0 <= K <= 2.
  - but recent advances such as ResNet/GoogLeNet challenge this paradigm

# [ConvNetJS demo: training on CIFAR-10]

https://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html

# Construct of a Convolutional Layer

# python mnist_3.0_convolutional.py => 98.91%

## Convolutional neural network

28x28x1

convolutional layer, 4 channels
W1[5, 5, 1, 4] stride 1

28x28x4

convolutional layer, 8 channels
W2[4, 4, 4, 8] stride 2

14x14x8

convolutional layer, 12 channels
W3[4, 4, 8, 12] stride 2

7x7x12

200
10

fully connected layer    W4[7x7x12, 200]
softmax readout layer    W5[200, 10]

+ biases on all layers

## Tensorflow - initialisation

filter size    input channels    output channels

K=4
L=8
M=12

```
W1 = tf.Variable(tf.truncated_normal([5, 5, 1, K] ,stddev=0.1))
B1 = tf.Variable(tf.ones([K])/10)
W2 = tf.Variable(tf.truncated_normal([5, 5, K, L] ,stddev=0.1))
B2 = tf.Variable(tf.ones([L])/10)
W3 = tf.Variable(tf.truncated_normal([4, 4, L, M] ,stddev=0.1))
B3 = tf.Variable(tf.ones([M])/10)

N=200

W4 = tf.Variable(tf.truncated_normal([7*7*M, N] ,stddev=0.1))
B4 = tf.Variable(tf.ones([N])/10)
W5 = tf.Variable(tf.truncated_normal([N, 10] ,stddev=0.1))
B5 = tf.Variable(tf.zeros([10])/10)
```

weights initialised with random values

## Tensorflow - the model

input image batch    weights    stride    biases
X[100, 28, 28, 1]

```
Y1 = tf.nn.relu(tf.nn.conv2d(X, W1, strides=[1, 1, 1, 1], padding='SAME') + B1)
Y2 = tf.nn.relu(tf.nn.conv2d(Y1, W2, strides=[1, 2, 2, 1], padding='SAME') + B2)
Y3 = tf.nn.relu(tf.nn.conv2d(Y2, W3, strides=[1, 2, 2, 1], padding='SAME') + B3)

YY = tf.reshape(Y3, shape=[-1, 7 * 7 * M])

Y4 = tf.nn.relu(tf.matmul(YY, W4) + B4)
Y  = tf.nn.softmax(tf.matmul(Y4, W5) + B5)
```

flatten all values for fully connected layer

Y3 [100, 7, 7, 12]

YY [100, 7x7x12]



98.91%

Southern University of Science and Technology

# python mnist_3.1_convolutional_bigger_dropout.py => 99.3%

# Data Preprocessing

# Batch Normalization Theory

# Batch Normalization Done Right

# Convolutional Batch Normalization



AncoraSIR.com

# Example: Airplane or not?

# Transfer Learning

You need a lot of data if you want to train/use CNNs, really?



Transfer learning with CNNs is pervasive…
(it's the norm, not an exception)



Next class:
- Transfer learning for arcade claw machine robot

AncoraSIR.com

# Thank you!

Prof. Song Chaoyang

- Dr. Wan Fang (sophie.fwan@hotmail.com)

AncoraSIR.com