

Lecture 11

(Artificial) Neural Networks

Song Chaoyang

Assistant Professor

Department of Mechanical and Energy Engineering

songcy@sustech.edu.cn

Tinker With a Neural Network

TensorBoard: Visualizing Learning

Epoch: 000,000 | Learning rate: 0.03 | Activation: Tanh | Regularization: None | Regularization rate: 0 | Problem type: Classification

DATA
Which dataset do you want to use?
Ratio of training to test data: 50%
Noise: 0
Batch size: 10
REGENERATE

FEATURES
Which properties do you want to feed in?
 X_1
 X_2
 X_1^2
 X_2^2
 $X_1 X_2$
 $\sin(X_1)$
 $\sin(X_2)$

2 HIDDEN LAYERS
4 neurons | 2 neurons

This is the output from one neuron. Hover to see it larger.

The outputs are mixed with varying weights, shown by the thickness of the lines.

OUTPUT
Test loss 0.515
Training loss 0.507

Colors shows data, neuron and weight values.

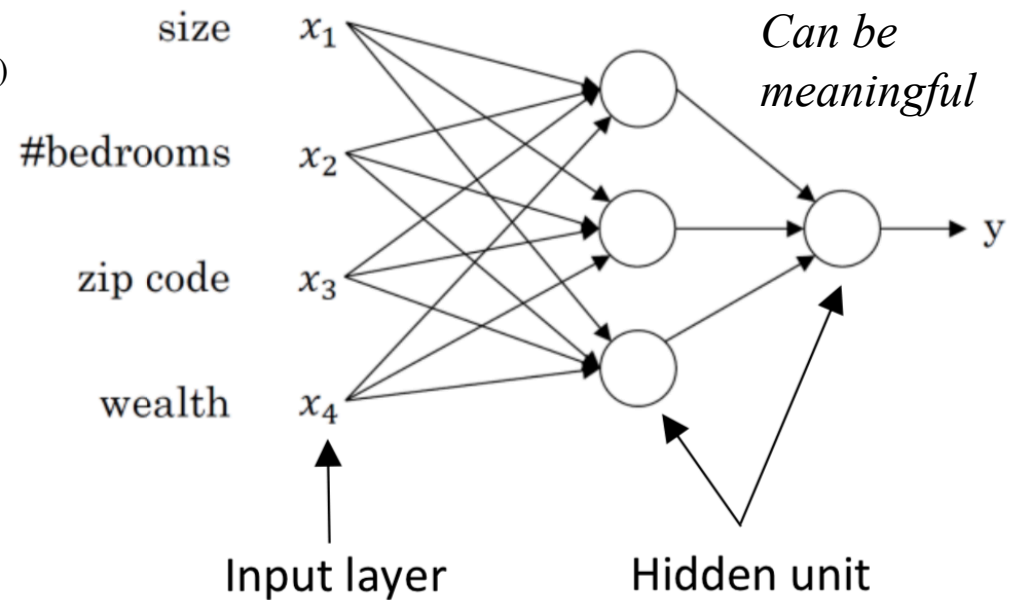
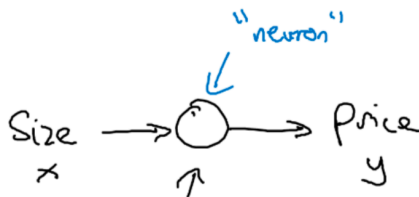
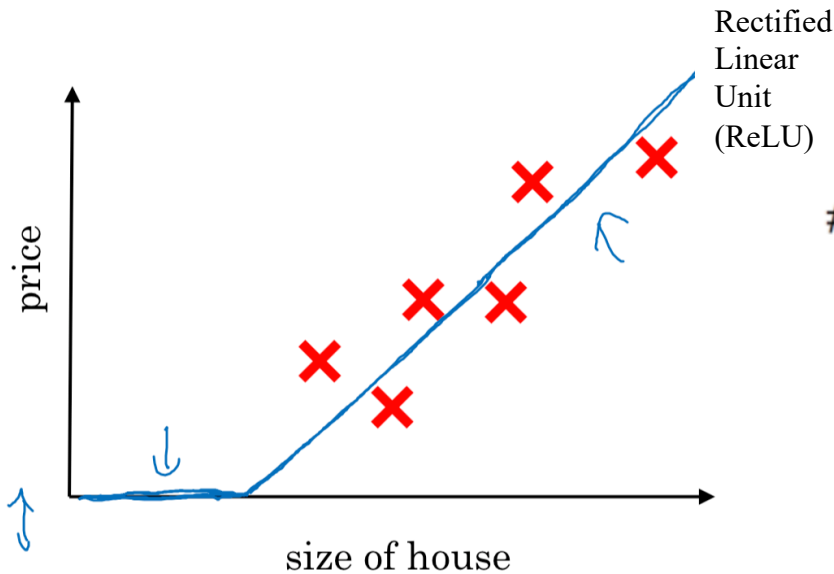
Show test data Discretize output

https://www.tensorflow.org/guide/summaries_and_tensorboard

What is a Neural Network?

It is a powerful learning algorithm inspired by how the brain works.

Housing Price Prediction



Supervised Learning with NNs

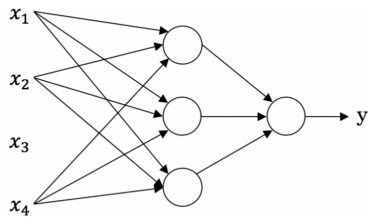
Examples, NNs & Data Types

Input(x)	Output (y)	Application
Home features	Price	Real Estate } <i>Standard NN</i>
Ad, user info	Click on ad? (0/1)	Online Advertising }
Image	Object (1,...,1000)	Photo tagging } <i>CNN</i>
<u>Audio</u>	Text transcript	Speech recognition } <i>RNN</i>
<u>English</u>	Chinese	Machine translation }
<u>Image, Radar info</u>	Position of other cars	Autonomous driving } <i>Custom/ Hybrid</i>

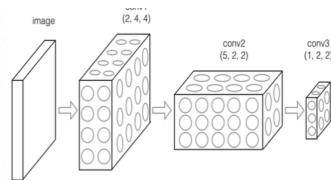
Structured Data

Size	#bedrooms	...	Price (1000\$s)
2104	3		400
1600	3		330
2400	3		369
⋮	⋮		⋮
3000	4		540

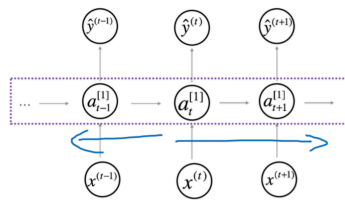
User Age	Ad Id	...	Click
41	93242		1
80	93287		0
18	87312		1
⋮	⋮		⋮
27	71244		1



Standard NN

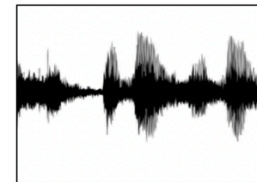


Convolutional NN

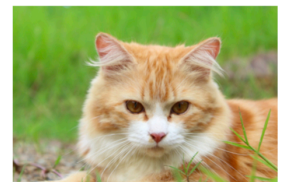


Recurrent NN

Unstructured Data



Audio



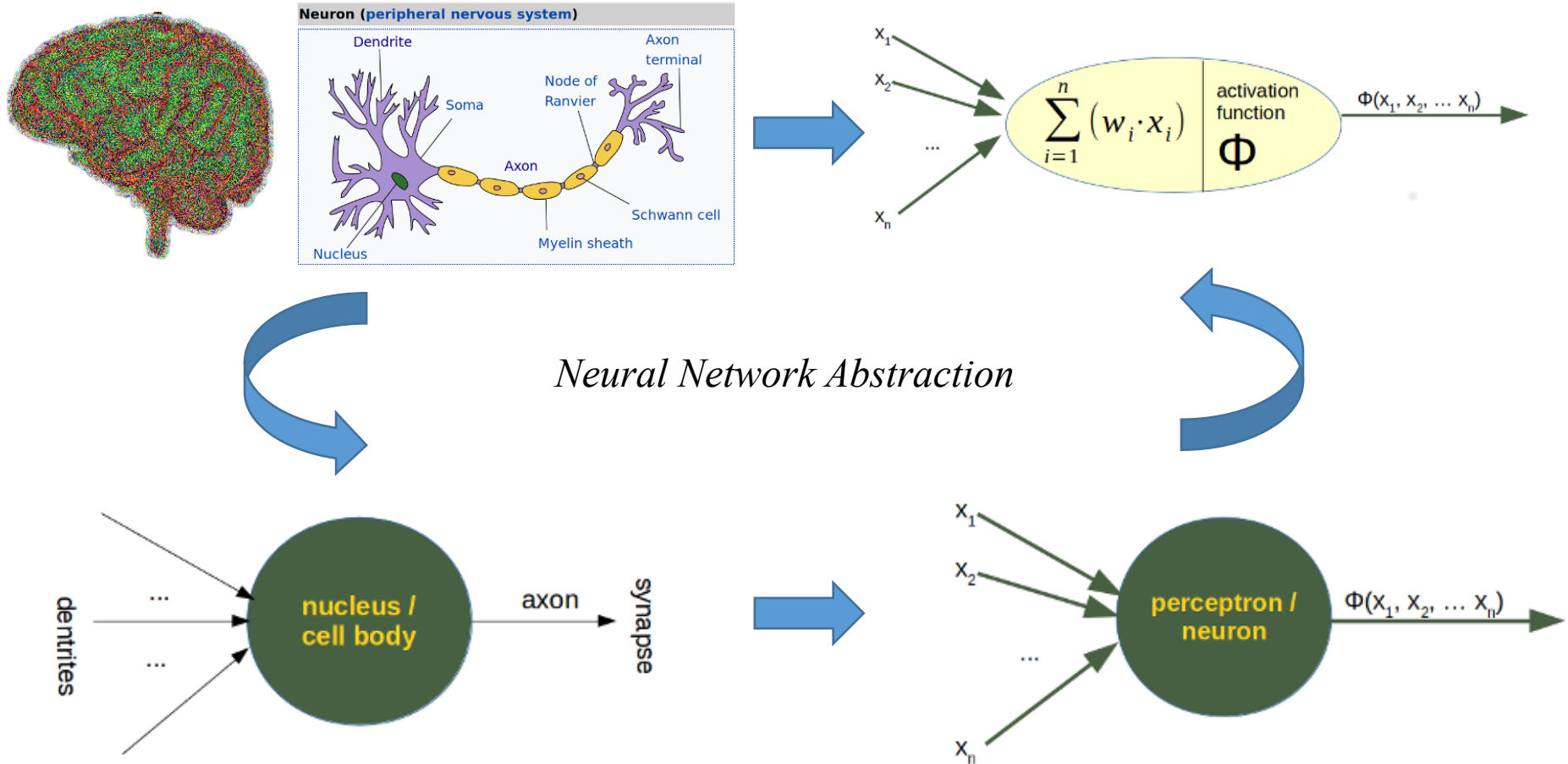
Image

Four scores and seven years ago...

Text

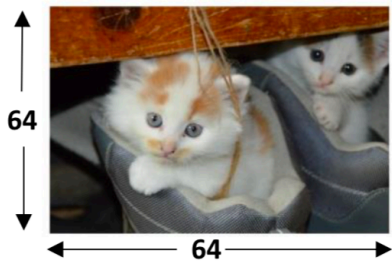
What is a Neural Network?

From biological inspiration to mathematical modeling



Example: Binary Classification

0 or 1 (*discrete value output*), it is a question ...



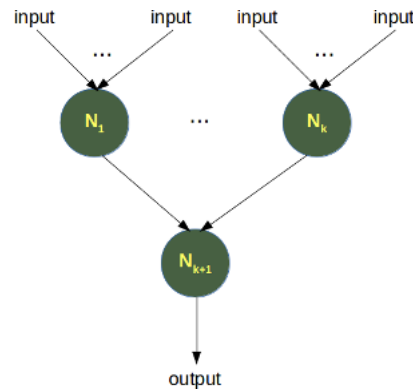
digitization

		Blue			
Green		255	134	93	22
Red		255	134	202	22
	Red	255	231	42	22
		123	94	83	2
		34	44	187	92
		34	76	232	124
		67	83	194	202

vectorization

$$x = \begin{bmatrix} 255 \\ 231 \\ 42 \\ \vdots \\ 255 \\ 134 \\ 202 \\ \vdots \\ 255 \\ 134 \\ 93 \\ \vdots \end{bmatrix}$$

} red
} green
} blue



A simple NN implementing AND

Input1	Input2	Output
0	0	0
0	1	0
1	0	0
1	1	1

```

0 0 0
1 [0 0] 0
2 [0 1] 0
3 [1 0] 0
4 [1 1] 1
    
```

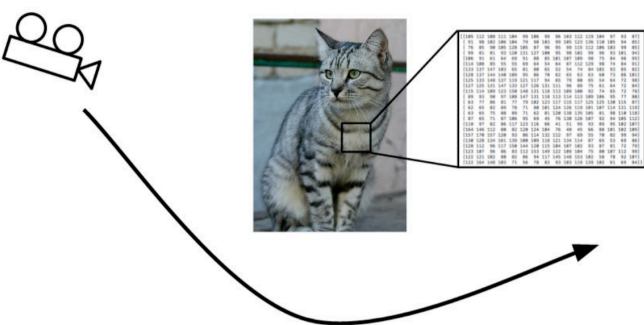
```

1 #!/usr/bin/python3
2 # adapted from https://www.python-course.eu/neural_networks.php
3 # use for ME336 CoBot Learning @ SUSTech by Dr. Song Chaoyang
4
5 import numpy as np
6 class Perceptron: # define the perceptron
7
8     def __init__(self, input_length, weights = None):
9         if weights == None:
10             self.weights = np.ones(input_length) * 0.5
11         else:
12             self.weights = weights
13
14     @staticmethod
15     def unit_step_function(x):
16         if x > 0.5:
17             return 1
18         return 0
19
20     def __call__(self, in_data):
21         weighted_input = self.weights * in_data
22         weighted_sum = weighted_input.sum()
23         return Perceptron.unit_step_function(weighted_sum)
24
25 p = Perceptron(2, np.array([0.5, 0.5])) # config the perceptron
26
27 for x in [np.array([0, 0]), np.array([0, 1]), # input
28           np.array([1, 0]), np.array([1, 1])]:
29     y = p(np.array(x)) # output
30     print(x, y)
    
```

Challenges of Recognition

Image Classification: A core task in computer vision

Viewpoint



All pixels change when the camera moves!

Illumination



This image is [CC0 1.0](#) public domain

Deformation



This image by [Umberto Salvagnin](#) is licensed under [CC-BY 2.0](#)

Occlusion

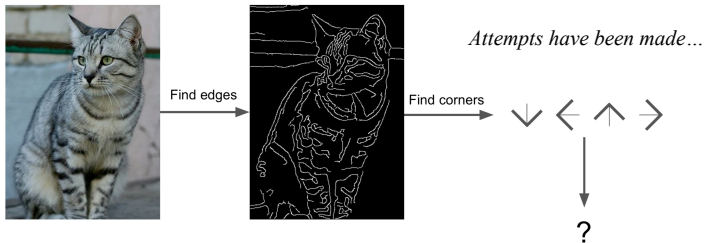


This image by [jonsson](#) is licensed under [CC-BY 2.0](#)

An image classifier

```
def classify_image(image):  
    # Some magic here?  
    return class_label
```

no obvious way to hard-code the algorithm for recognizing a cat, or other classes.



Clutter



This image is [CC0 1.0](#) public domain

Intraclass Variation



This image is [CC0 1.0](#) public domain

Machine Learning: Data-driven approach

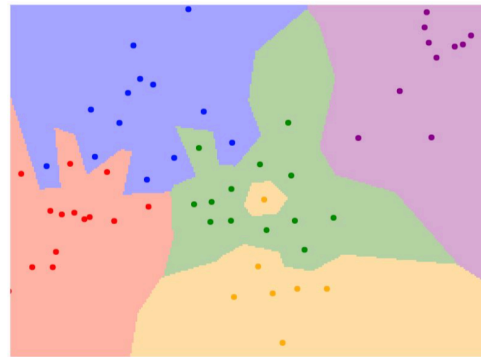
First classifier: Nearest Neighbor

airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

Example Dataset: CIFAR10
 10 classes
 50,000 training images
 10,000 testing images

1. Collect a dataset of images and labels
2. Use Machine Learning to train a classifier
3. Evaluate the classifier on new images

1-NN classifier



```
def train(images, labels):
    # Machine learning!
    return model
```

Memorize all data and labels

```
def predict(model, test_images):
    # Use model to predict labels
    return test_labels
```

Predict the label of the most similar training image

Distance Metric to compare images

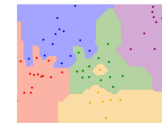
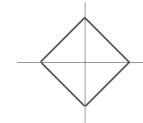
L1 distance: $d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$

test image				training image				pixel-wise absolute value differences			
56	32	10	18	10	20	24	17	46	12	14	1
90	23	128	133	8	10	89	100	82	13	39	33
24	26	178	200	12	16	178	170	12	10	0	30
2	0	255	220	4	32	233	112	2	32	22	108

add → 456

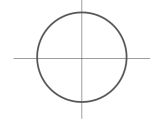
L1 (Manhattan) distance

$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$



L2 (Euclidean) distance

$$d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$



Hyperparameters

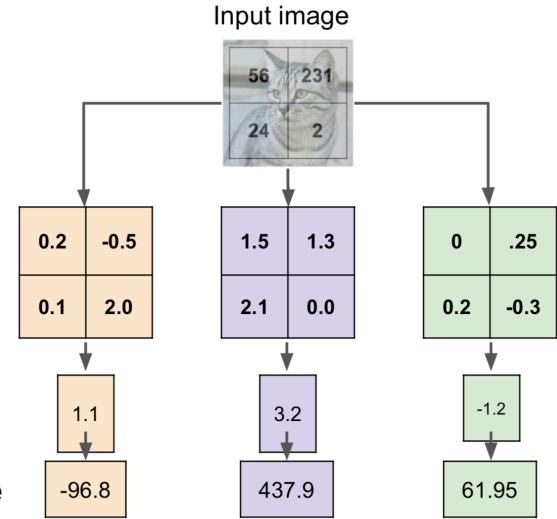
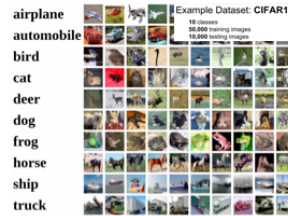
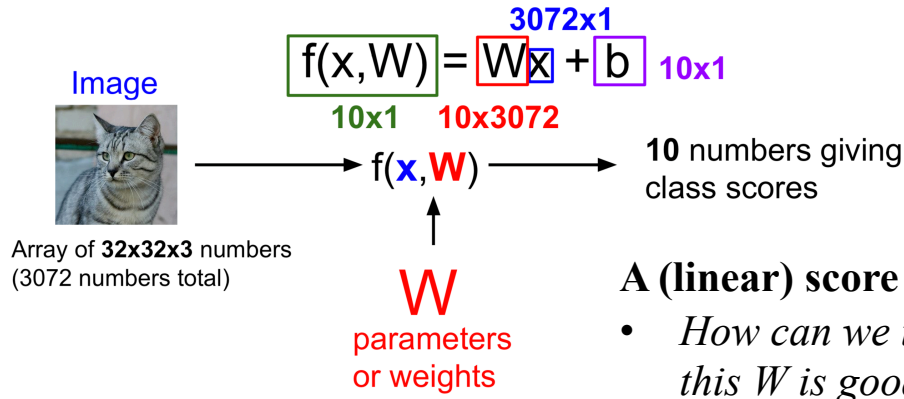
- What is the best value of k to use?
- What is the best distance to use?

Choices about the algorithm that we *set* rather than *learn*

- *Very problem-dependent.*
- *Must try them all out and see what works best.*



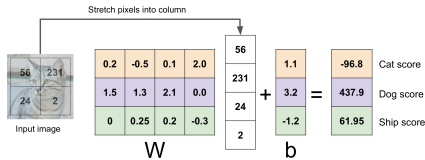
Parametric Approach: Linear Classifier



Viewpoint

Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

Algebraic



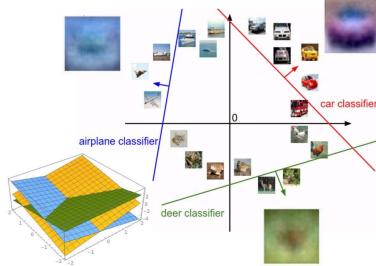
$$f(x, W) = Wx$$

Visual



One template per class

Geometric

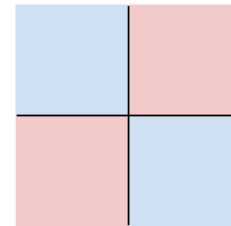


Hyperplanes cutting up space

Hard cases for a linear classifier

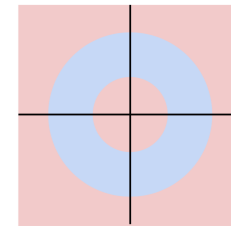
Class 1:
First and third quadrants

Class 2:
Second and fourth quadrants



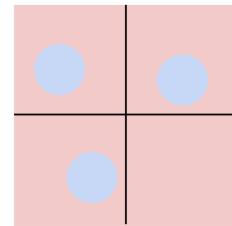
Class 1:
 $1 \leq L2 \text{ norm} \leq 2$

Class 2:
Everything else



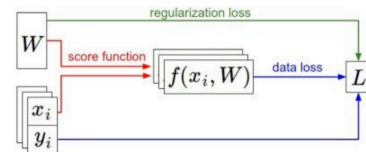
Class 1:
Three modes

Class 2:
Everything else



Loss Function

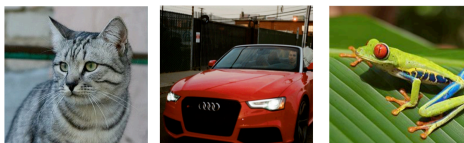
how good our current classifier is



Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:

Given a dataset of examples $\{(x_i, y_i)\}_{i=1}^N$

Where x_i is image and y_i is (integer) label



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1
Losses:	2.9	0	12.9

Loss over the dataset is a sum of loss over examples:

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

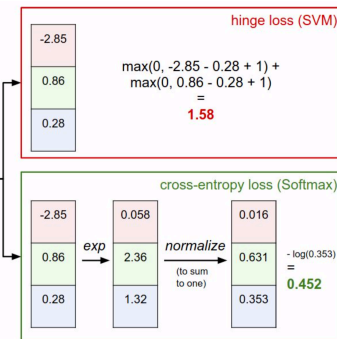
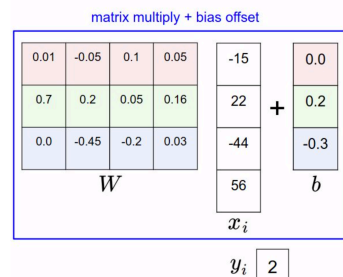
e.g. Multiclass SVM loss

$$L_i = \sum_{j \neq y_i} \begin{cases} 0 & \text{if } s_{y_i} \geq s_j + 1 \\ s_j - s_{y_i} + 1 & \text{otherwise} \end{cases}$$

$$= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

We only care about the relative distance (property of the loss), not the absolute difference (property of the data)

Softmax vs. SVM



Regularization

λ = regularization strength (hyperparameter)

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss}} + \underbrace{\lambda R(W)}_{\text{Regularization}}$$

Data loss: Model predictions should match training data

Regularization: Prevent the model from doing *too* well on training data

Simple examples

L2 regularization: $R(W) = \sum_k \sum_l W_{k,l}^2$

L1 regularization: $R(W) = \sum_k \sum_l |W_{k,l}|$

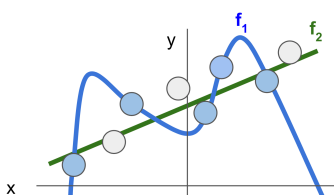
Elastic net (L1 + L2): $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

More complex:

Dropout

Batch normalization

Stochastic depth, fractional pooling, etc



Regularization pushes against fitting the data too well so we don't fit noise in the data

Softmax Classifier (Multinomial Logistic Regression)

Want to interpret raw classifier scores as **probabilities**



$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{Softmax Function}$$

Probabilities must be ≥ 0

Probabilities must sum to 1

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

cat
car
frog

cat	3.2
car	5.1
frog	-1.7

exp

cat	24.5
car	164.0
frog	0.18

normalize

cat	0.13
car	0.87
frog	0.00

cat	1.00
car	0.00
frog	0.00

Unnormalized log-probabilities / logits

unnormalized probabilities

probabilities

Correct

compare \leftarrow

Cross Entropy $H(P, Q) = H(p) + D_{KL}(P||Q)$

Optimization using Gradient Descent

Vanilla Gradient Descent

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

In 1-dimension, the derivative of a function:

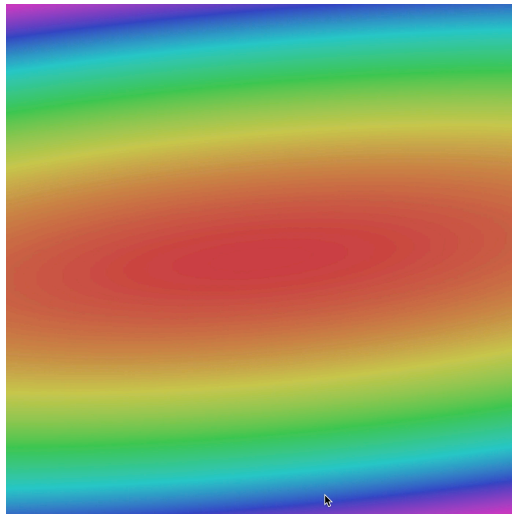
$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

Numerical gradient: slow (:, approximate :, easy to write :)
Analytic gradient: fast (:, exact :, error-prone :)

In practice: Derive analytic gradient, check your implementation with numerical gradient

In multiple dimensions, the **gradient** is the vector of (partial derivatives) along each dimension

The slope in any direction is the **dot product** of the direction with the gradient
 The direction of steepest descent is the **negative gradient**



AncoraSIR.com

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
```

Stochastic Gradient Descent (SGD)

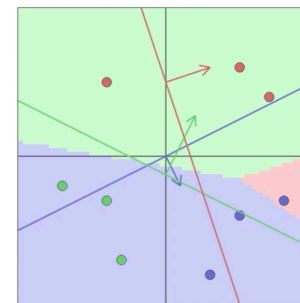
$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$$

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

Full sum expensive when N is large!

Approximate sum using a **minibatch** of examples
 32 / 64 / 128 common

Data points are shown as circles colored by their class (red/gree/blue). The background regions are colored by whichever class is most likely at any point according to the current weights. Each classifier is visualized by a line that indicates its zero score level set. For example, the blue classifier computes scores as $W_{0,0}x_0 + W_{0,1}x_1 + b_0$ and the blue line shows the set of points (x_0, x_1) that give score of zero. The blue arrow draws the vector $(W_{0,0}, W_{0,1})$, which shows the direction of score increase and its length is proportional to how steep the increase is. Note: you can drag the datapoints.



Parameters W, b are shown below. The value is in bold and its gradient (computed with backprop) is in *red, italic* below. Click the triangles to control the parameters.

$w[0,0]$	$w[0,1]$	$b[0]$
1.00 <i>-0.38</i>	2.00 <i>0.07</i>	0.00 <i>0.11</i>
$w[1,0]$	$w[1,1]$	$b[1]$
2.00 <i>0.51</i>	-4.00 <i>-0.58</i>	0.50 <i>-0.11</i>
$w[2,0]$	$w[2,1]$	$b[2]$
3.00 <i>0.17</i>	-1.00 <i>0.36</i>	-0.50 <i>0.00</i>

Step size: 0.10000
 Single parameter update
 Start repeated update
 Stop repeated update
 Randomize parameters

Visualization of the data loss computation. Each row is loss due to one datapoint. The first three columns are the 2D data x_i and the label y_i . The next three columns are the three class scores from each classifier $f(x_i; W, b) = Wx_i + b$ (E.g. $s[0] = x[0] * W[0,0] + x[1] * W[0,1] + b[0]$). The last column is the data loss for a single example, L_i .

$x[0]$	$x[1]$	y	$s[0]$	$s[1]$	$s[2]$	L	
0.50	0.40	0	1.30	-0.10	0.60	0.30	
0.80	0.30	0	1.40	0.90	1.60	1.70	
0.30	0.80	0	1.90	-2.10	-0.40	0.00	
-0.40	0.30	1	0.20	-1.50	-2.00	3.20	
-0.30	0.70	1	1.10	-2.90	-2.10	6.80	
-0.70	0.20	1	-0.30	-1.70	-2.80	2.40	
0.70	-0.40	2	-0.10	3.50	2.00	2.50	
0.50	-0.60	2	-0.70	3.90	1.60	3.30	
-0.40	-0.50	2	-1.40	1.70	-1.20	4.70	
						mean:	2.77

Total data loss: 2.77
 Regularization loss: 3.50
 Total loss: 6.27

L2 Regularization strength: 0.10000

Multiclass SVM loss formulation:
 Weston Watkins 1999
 One vs. All
 Structured SVM
 Softmax

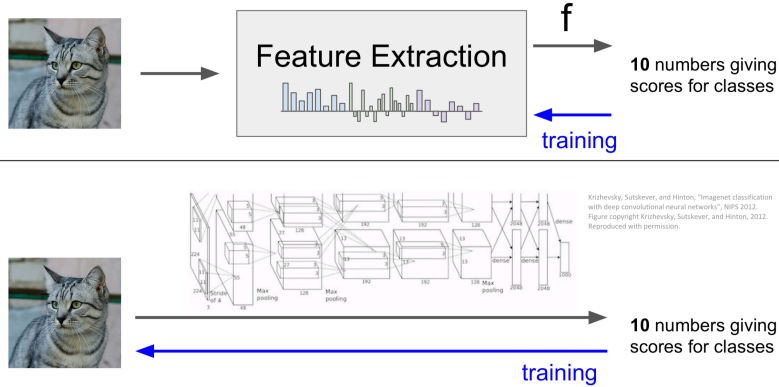
<http://vision.stanford.edu/teaching/cs231n-demos/linear-classify/>



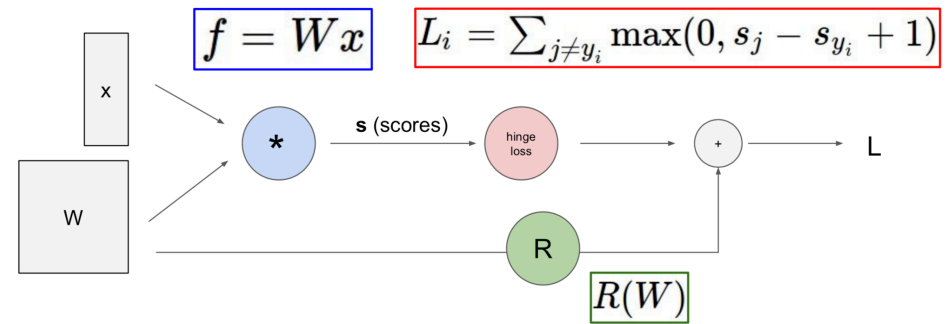
SUSTech
 Southern University
 of Science and Technology

Convolutional Neural Network

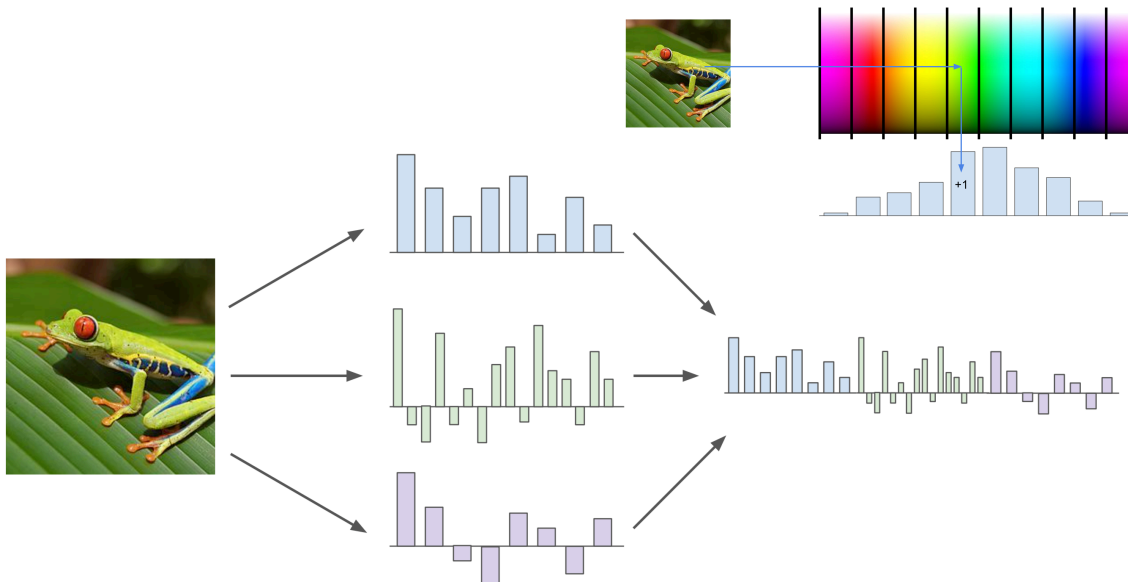
Image features vs ConvNets



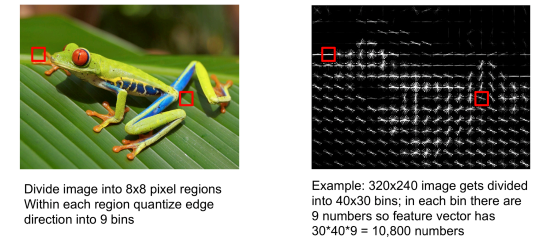
Computational Graph



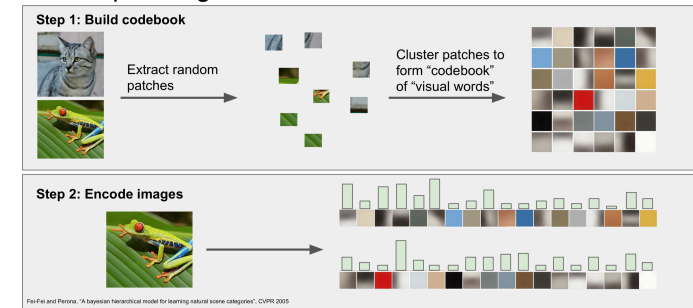
Example: Color Histogram



Example: Histogram of Oriented Gradients (HoG)



Example: Bag of Words



Backpropagation

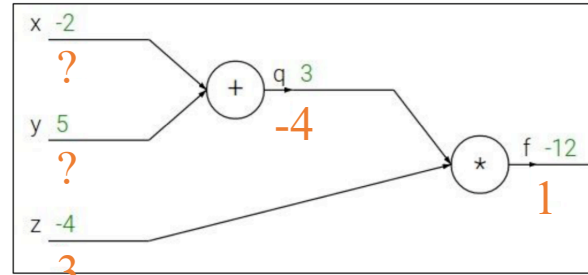
backward propagation of errors

Calculates the gradient of a loss function with respect to all the weights in the network.

Backpropagation: a simple example

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

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



Chain rule:

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

Upstream gradient

Local gradient

Chain rule:

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

Upstream gradient

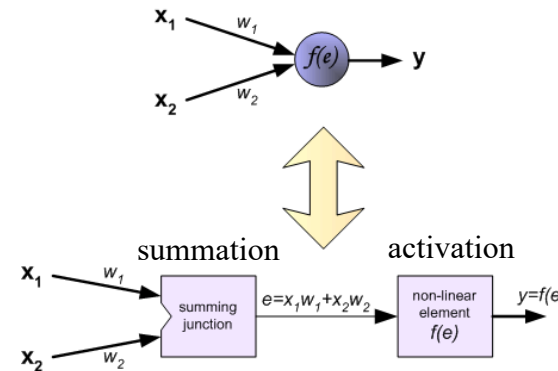
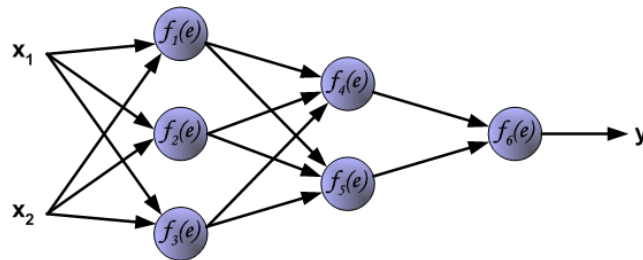
Local gradient

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

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

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

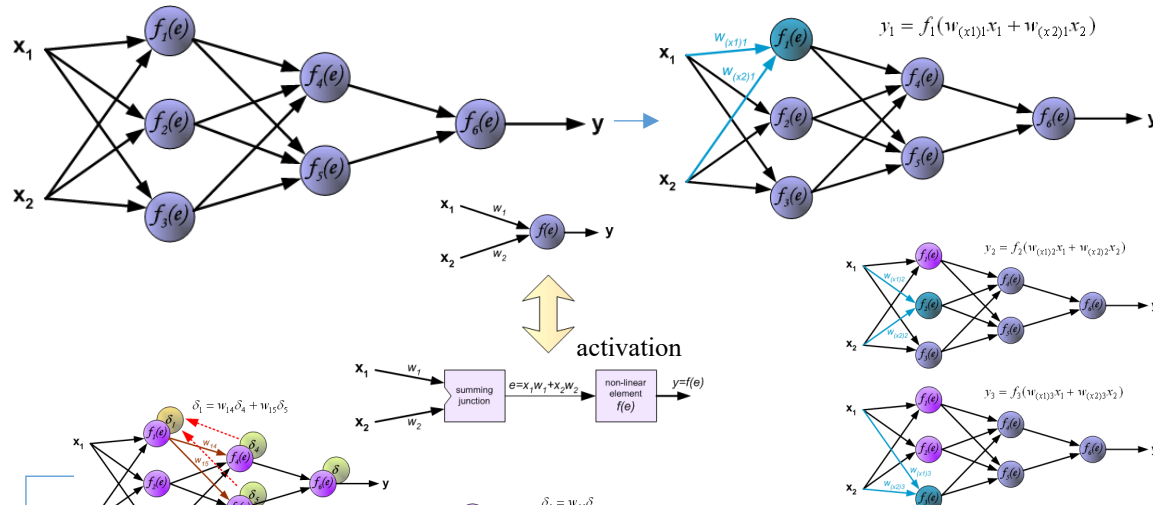
- neural nets will be very large: impractical to write down gradient formula by hand for all parameters
- backpropagation** = recursive application of the chain rule along a computational graph to compute the gradients of all inputs/parameters/intermediates
- implementations maintain a graph structure, where the nodes implement the **forward()** / **backward()** API
- forward**: compute result of an operation and save any intermediates needed for gradient computation in memory
- backward**: apply the chain rule to compute the gradient of the loss function with respect to the inputs



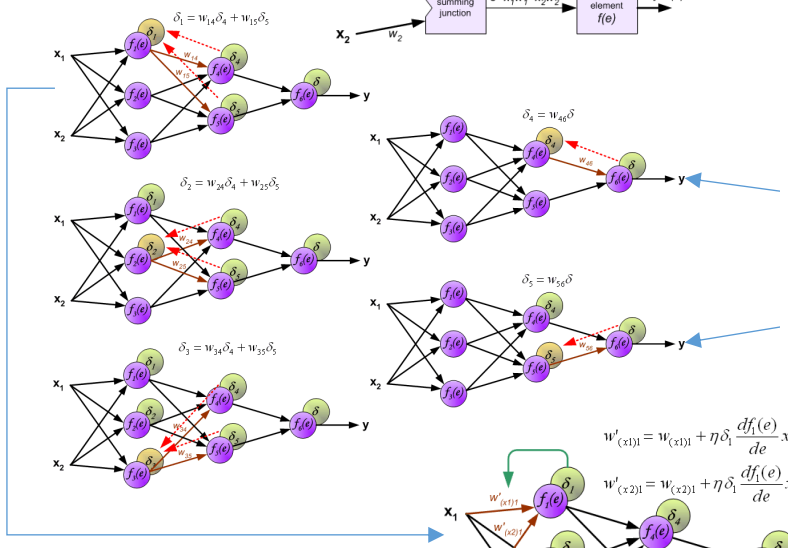
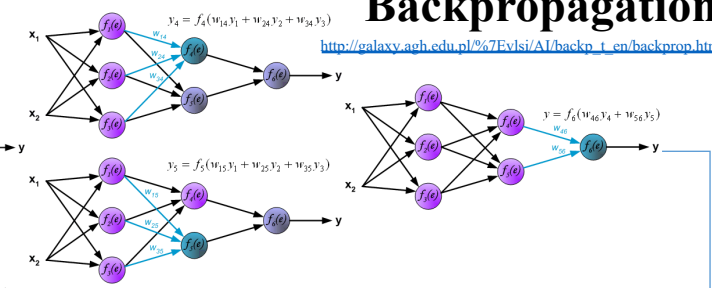
<http://galaxy.agh.edu.pl/%7Evlsi/AI/backprop.html>

Principles of Training Multi-layer Neural Network Using Backpropagation

http://galaxy.agh.edu.pl/%7Evlsl/Al/backp_t_en/backprop.html



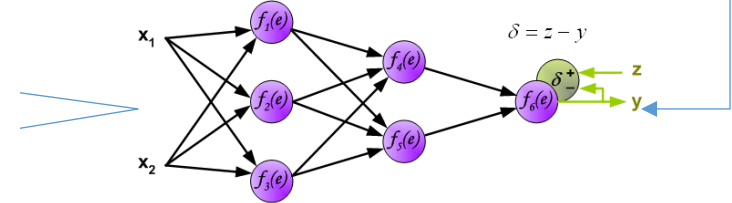
The idea is to propagate error signal d (computed in single teaching step) back to all neurons, which output signals were input for discussed neuron.



When the error signal for each neuron is computed, the weights coefficients of each neuron input node may be modified.

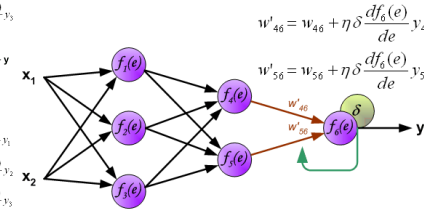
In formulas on the right, $df(e)/de$ represents derivative of neuron activation function (which weights are modified).

AncoraSIR.com



Coefficient η affects network teaching speed, to select this parameter:

- The first method is to start teaching process with large value of the parameter. While weights coefficients are being established the parameter is being decreased gradually.



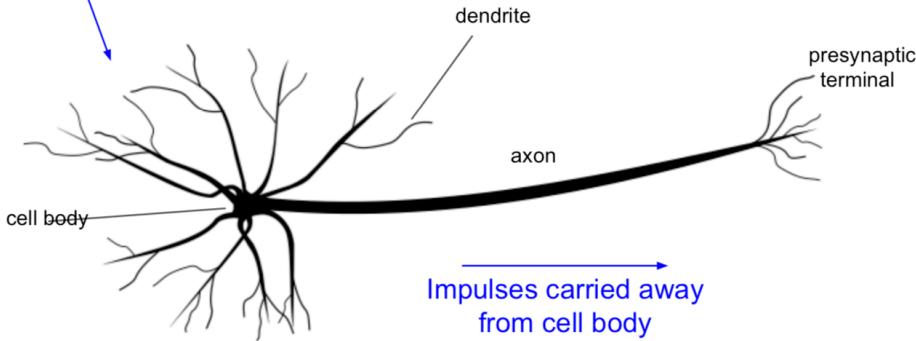
Finish one iteration of computation

- The second, more complicated, method starts teaching with small parameter value. During the teaching process the parameter is being increased when the teaching is advanced and then decreased again in the final stage. Starting teaching process with low parameter value enables to determine weights coefficients signs.

Neural Network

Biological Inspiration

Impulses carried toward cell body

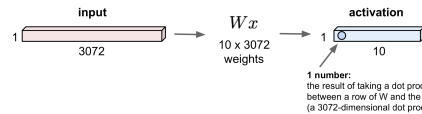
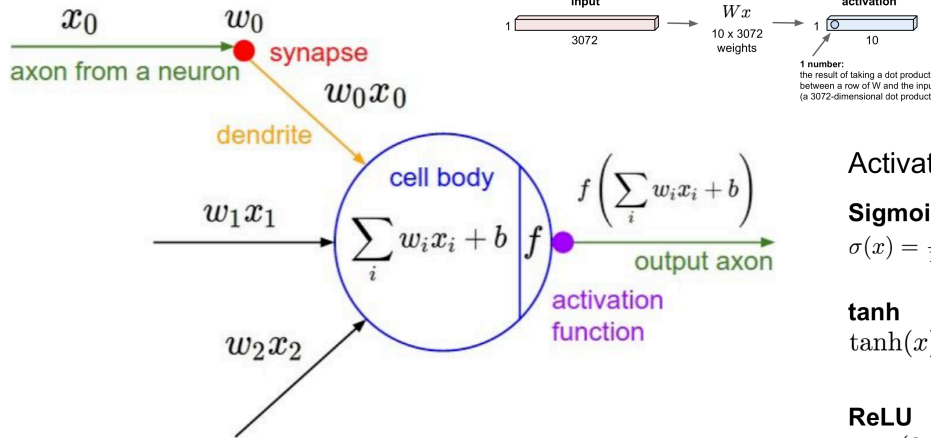


This image by Felipe Peruchio is licensed under CC-BY 3.0

Impulses carried away from cell body

Fully Connected Layer

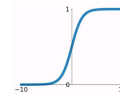
32x32x3 image -> stretch to 3072 x 1



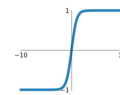
1 number: the result of taking a dot product between a row of W and the input (a 3072-dimensional dot product)

Activation functions

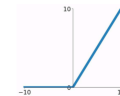
Sigmoid
 $\sigma(x) = \frac{1}{1+e^{-x}}$



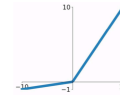
tanh
 $\tanh(x)$



ReLU
 $\max(0, x)$

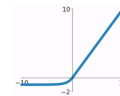


Leaky ReLU
 $\max(0.1x, x)$



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

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

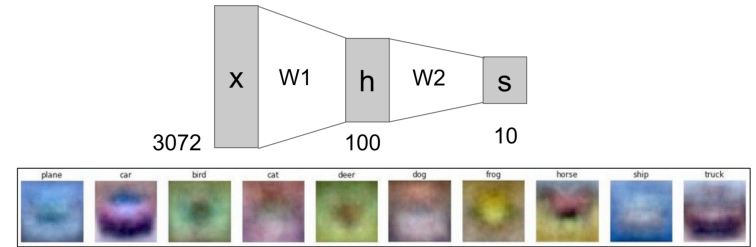


Linear score function:

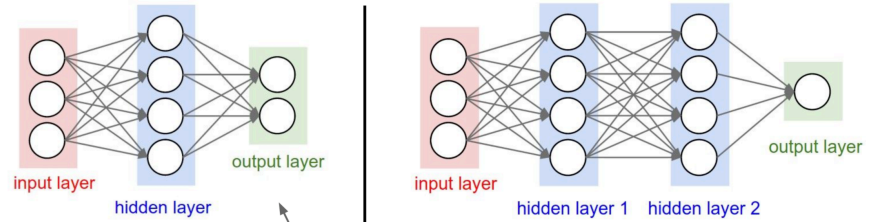
$$f = Wx$$

2-layer Neural Network

$$f = W_2 \max(0, W_1 x)$$



Neural networks: Architectures



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

"Fully-connected" layers

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

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)
```

Thank you!

Prof. Song Chaoyang

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

