

Deep Learning Crash Course



www.deeplearningcrashcourse.org

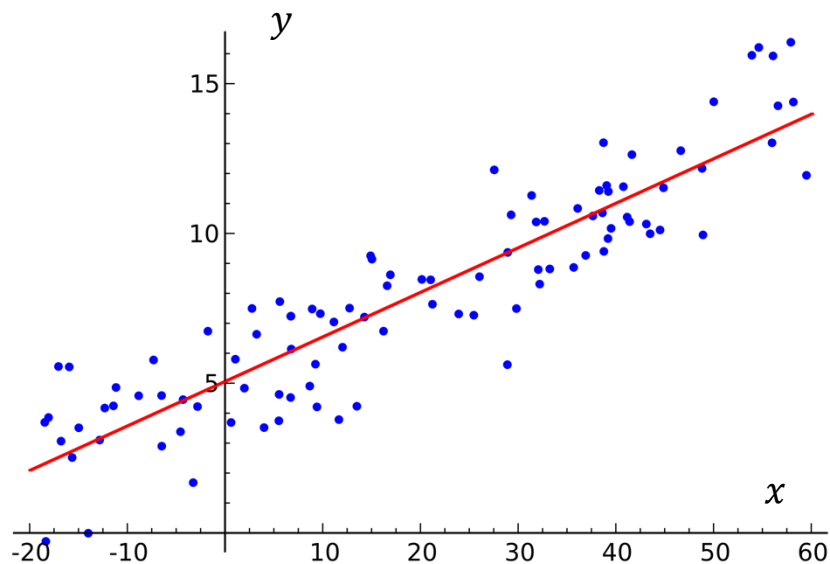
Hui Xue

Fall 2021

Lecture 1

- Deep learning as a data driven approach
- Binary and multi-class classification
- Multi-layer perceptron (MLP)

Linear regression



Model:

$$y = wx + b$$

Loss:

$$\ell(w, b) = \sum_{i=0}^{N-1} [y^{(i)} - (wx^{(i)} + b)]^2$$

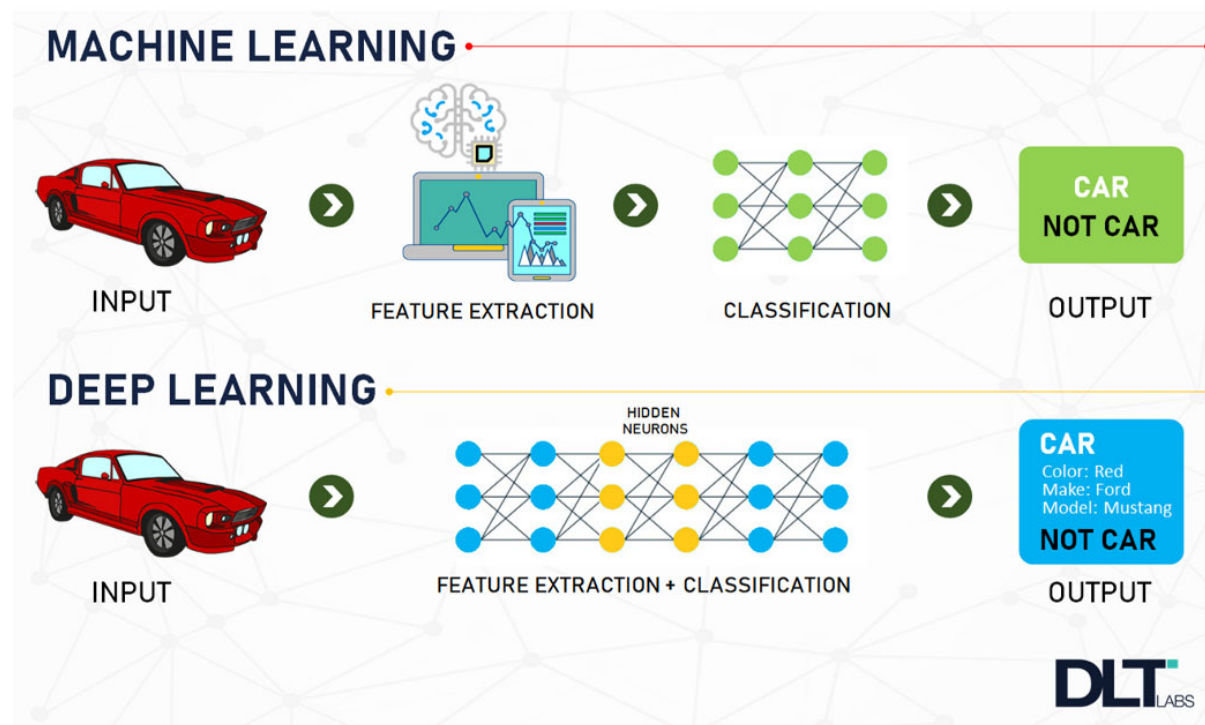
Optimization:

$$\min_{w, b} \ell(w, b)$$

Minimize the empirical risk,
computed on the measured data
sample $(x^{(i)}, y^{(i)})$

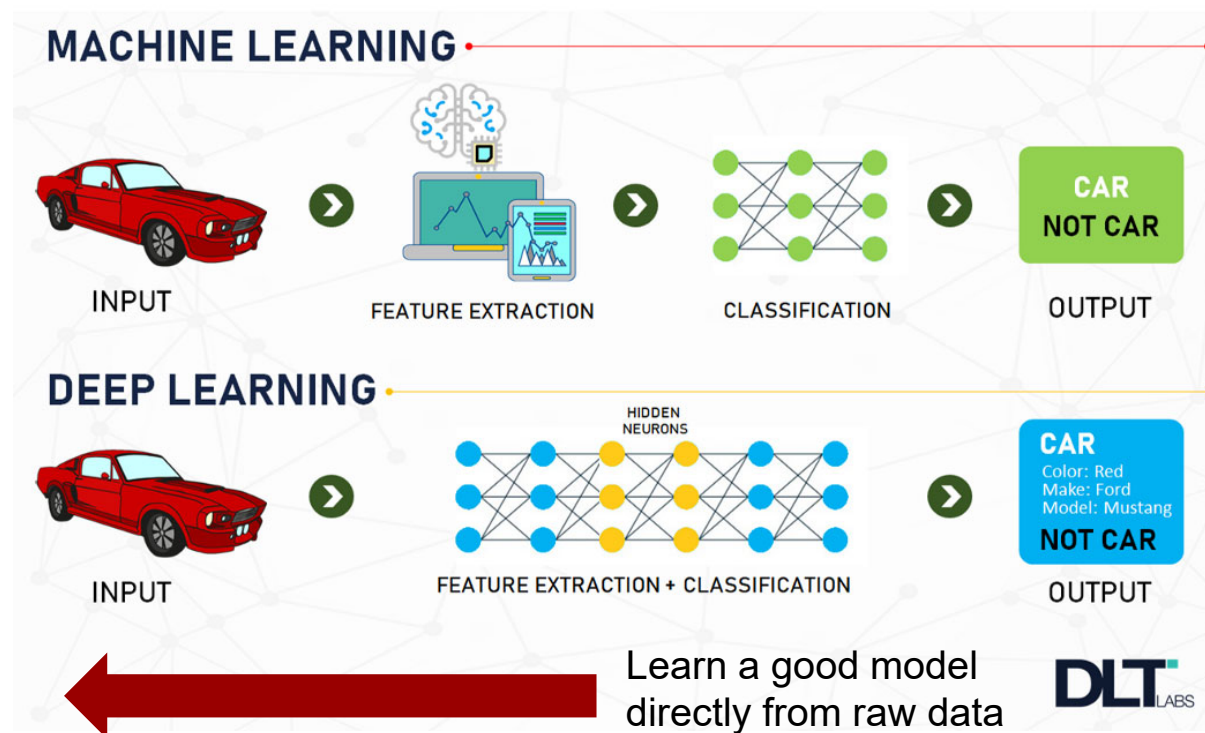
https://en.wikipedia.org/wiki/Linear_regression#/media/File:Linear_regression.svg

Deep Learning model is data-driven

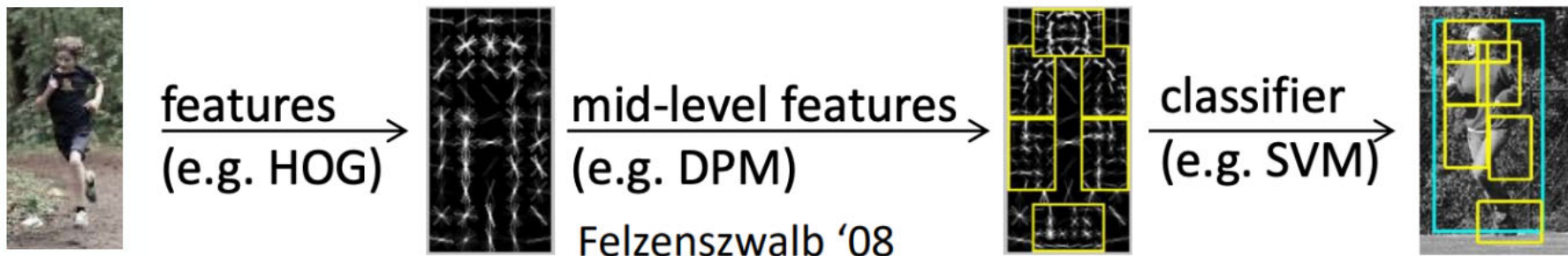


<https://dltlabs.medium.com/understanding-machine-learning-deep-learning-f5aa95264d61>

Deep Learning model is data-driven



Feature engineering



histogram of oriented gradients

- Human design and build feature extractor – very hard to scale, e.g. to 2000 object classes
- Model works on the feature vectors, instead of original data
- Limited by the human insights and amount of data

Object Detection with Discriminatively Trained Part Based Models. <https://cs.brown.edu/people/pfelzens/papers/lsvm-pami.pdf>

Data driven detection



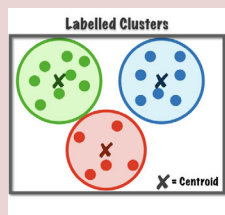
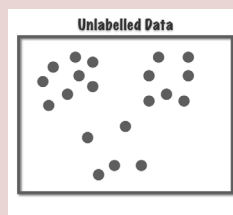
<https://youtu.be/VOC3huqHrss>

- Start from image to object detection (class and bounding box)
- End-to-end training, without human crafted features
- So powerful, can be done in real-time

Deep Learning Landscapes

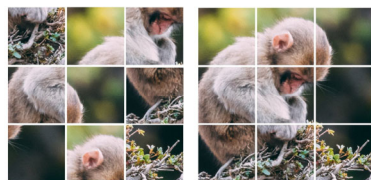
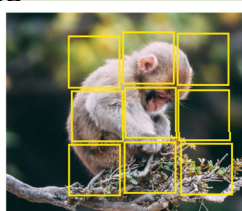
Unsupervised Learning

- Unlabeled X
- Learn representation in data



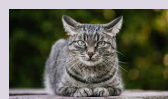
Self-supervised Learning

- Unlabeled X
- Generate self-supervisory signal
- Learn effective representation of data to help downstream apps



Supervised Learning

- Labelled X and Y
- Learn model: $X \rightarrow Y$
- Make prediction

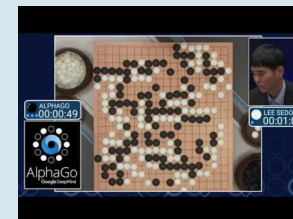
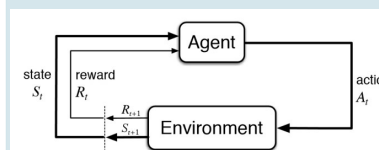


Is a cat?

1

Reinforcement Learning

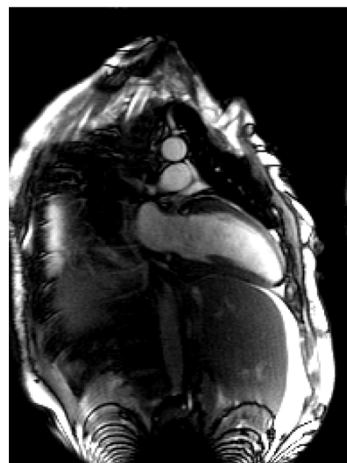
- Learn actions to reach the goal given an environment



Heart finding problem

A running example : heart finding from MRI images

$$X = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_N \end{bmatrix} =$$



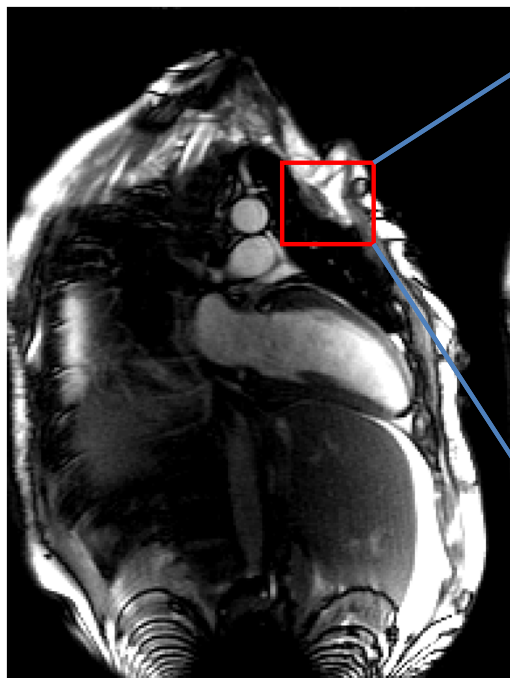
$Y = 1$, if imaging the heart



$Y = 0$

Computer sees the number, not content

What computer sees in this image:



967	965	967	975	984	993	1004	1015	1027	1041
1003	985	975	988	995	1001	1011	1021	1031	1042
838	944	1025	993	1001	1017	1017	1027	1037	1043
685	891	1053	994	1004	1026	1019	1029	1037	1042
737	912	1050	998	1006	1024	1019	1027	1033	1040
776	951	1074	1002	1006	1025	1018	1025	1030	1042
785	882	994	1010	1018	1024	1024	1027	1033	1055
825	789	886	1052	1043	1021	1040	1037	1046	1083
870	829	928	1092	1069	1041	1062	1059	1077	1131
937	932	956	1030	1099	1111	1099	1109	1143	1214
1068	1063	1014	1032	1179	1225	1193	1214	1262	1351
1240	1235	1195	1213	1330	1364	1336	1360	1419	1524
1390	1385	1350	1373	1481	1496	1462	1494	1567	1690
1528	1519	1488	1489	1546	1584	1593	1627	1713	1860
1708	1693	1681	1655	1625	1692	1769	1799	1892	2060
1921	1905	1892	1864	1830	1888	1964	2003	2107	2281
2106	2090	2072	2044	2017	2071	2145	2199	2317	2495
2256	2240	2219	2194	2177	2227	2303	2378	2510	2688
2404	2390	2355	2347	2385	2416	2459	2558	2702	2873
2555	2541	2501	2499	2561	2600	2642	2749	2896	3062
2724	2707	2672	2665	2700	2746	2825	2954	3102	3257
2895	2876	2846	2837	2850	2887	2993	3156	3303	3442
3044	3024	2995	2992	3014	3063	3176	3337	3480	3611
3202	3181	3153	3159	3196	3258	3372	3525	3656	3773
3358	3337	3313	3325	3367	3434	3545	3687	3815	3926
3491	3471	3452	3467	3514	3593	3689	3804	3944	4062
3635	3615	3598	3617	3670	3753	3843	3945	4082	4195
3771	3752	3739	3764	3821	3904	3996	4097	4220	4320
3881	3866	3860	3893	3955	4039	4129	4226	4343	4436
3996	3987	3991	4030	4095	4179	4266	4359	4471	4556
4121	4118	4132	4176	4241	4322	4409	4498	4591	4668
4248	4252	4273	4319	4384	4462	4548	4631	4705	4775
4376	4386	4414	4461	4525	4600	4682	4760	4829	4893
4499	4514	4548	4597	4660	4734	4811	4886	4952	5009

Model as a mapping function

First idea : design a mapping function to map image pixel values to the probability

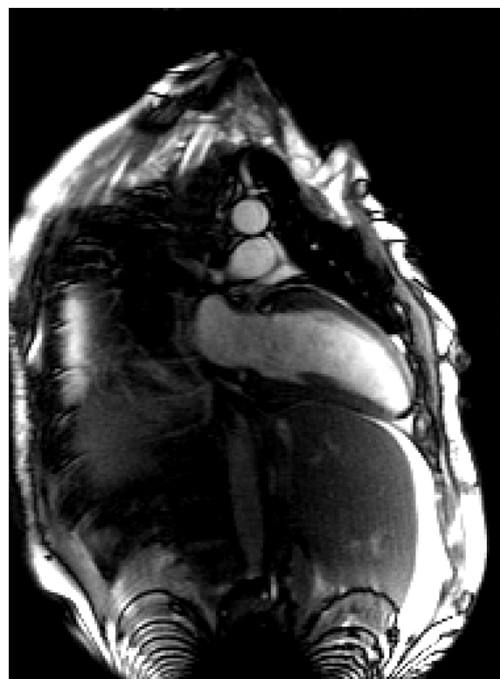
$$y = f(X)$$

y : a single floating value

lower case, regular X : $N \times 1$ vector

$N = H \times W$, number of pixels

upper case, bold



A scalar for the probability where heart was imaged

Linear mapping, weights and bias

Start with a linear mapping

$$z = \mathbf{W}X + b$$

The diagram illustrates the components of the linear mapping equation $z = \mathbf{W}X + b$. Three blue arrows point from descriptive labels below to the terms in the equation: one from 'z : 1x1 scalar' to 'z', one from '**W** : 1xN weights matrix' to '**W**', and one from 'b : 1x1 bias' to 'b'.

z : 1x1 scalar

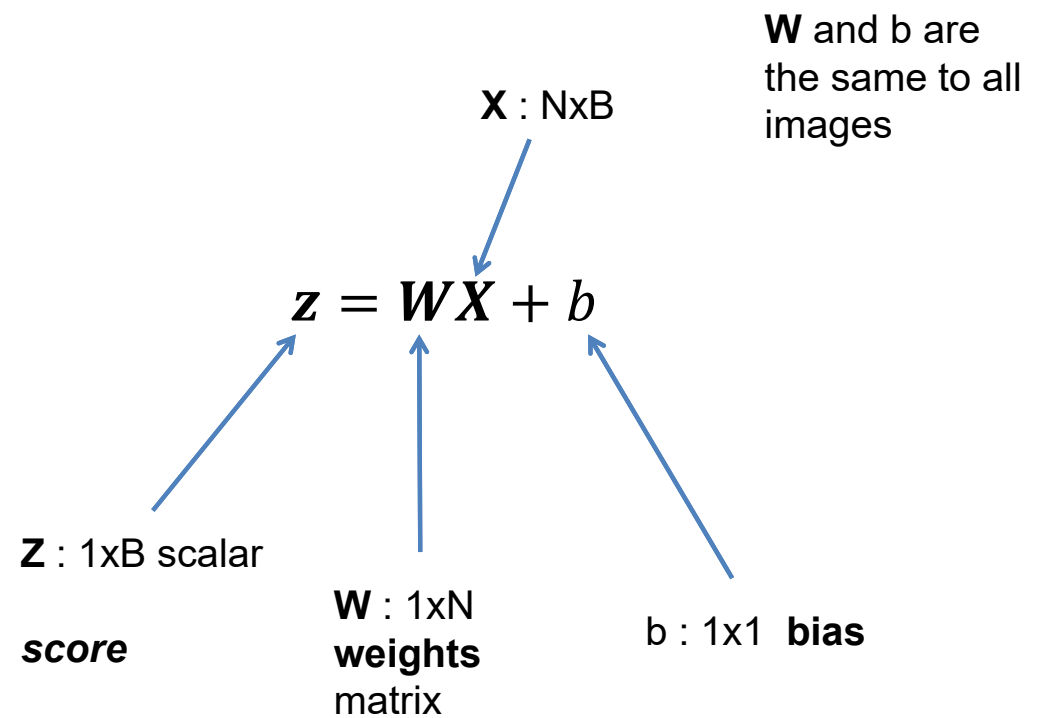
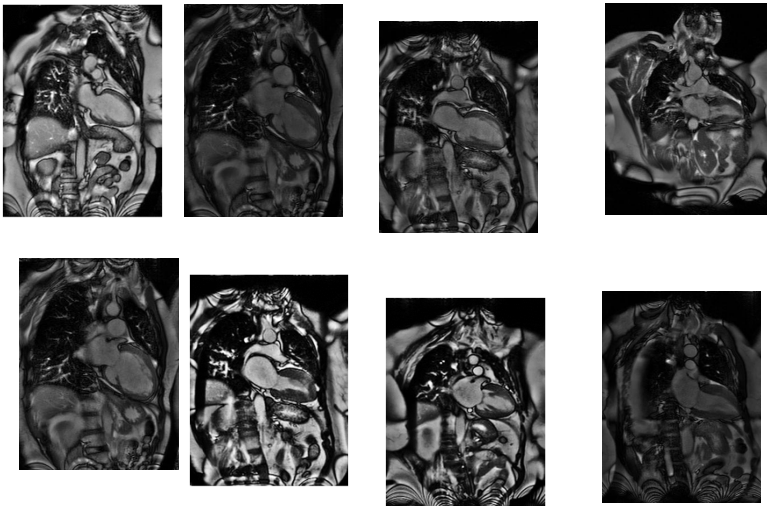
W : 1xN
weights
matrix

b : 1x1 **bias**

Process a *Batch*

We can process multiple images in one pass

B=8, a **batch**





Broadcasting

$$\mathbf{Z} = \mathbf{WX} + b$$

$W : 1 \times N$

$X : N \times B$

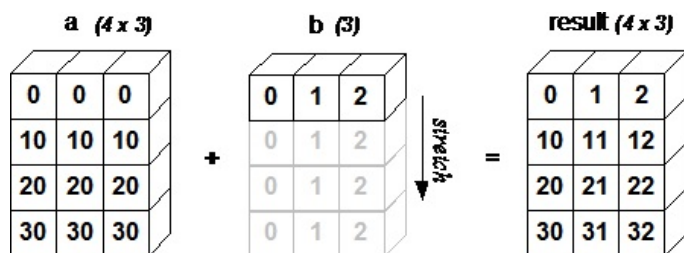
1×1 scalar

So Z is $1 \times B$

Add same b to every element of Z

Broadcasting : repeat the smaller array to make its shape compatible to larger array

- starts with the trailing (i.e. rightmost) dimensions and works its way left
- **Repeat any dimension of 1**



A (4d array): $8 \times 1 \times 6 \times 1$

B (3d array): $7 \times 1 \times 5$

Result (4d array): $8 \times 7 \times 6 \times 5$

<https://numpy.org/doc/stable/user/basics.broadcasting.html>

Enforce probability

$$z = WX + b$$

*z is not a value [0 1], so it
is not a probability*

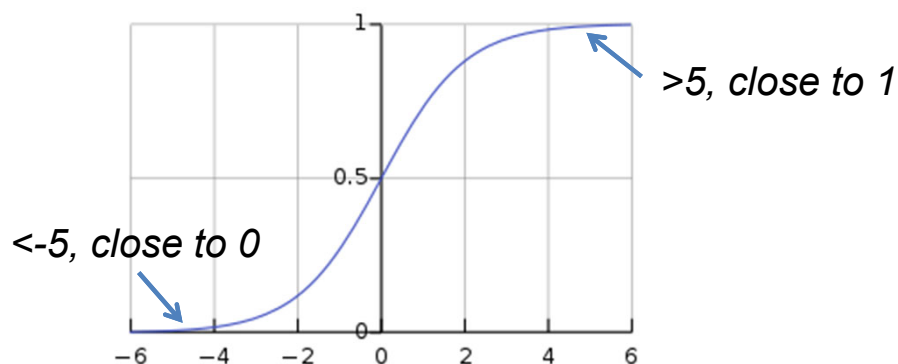


$$a = f(z) = f(WX + b)$$

f is applied element-wise

*A choice of f is the sigmoid
function*

Sigmoid function $\sigma(x) = \frac{1}{1 + \exp(-x)}$

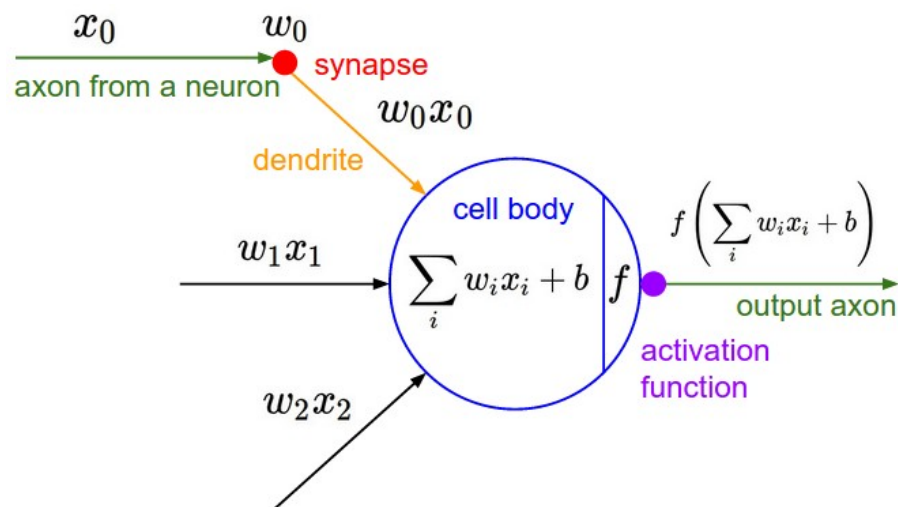


$$\frac{d\sigma(x)}{dx} = \sigma(x)[1 - \sigma(x)]$$

Active range -5 to 5

Accept inputs from $-\infty$ to $+\infty$

This is a “neuron” in neural network



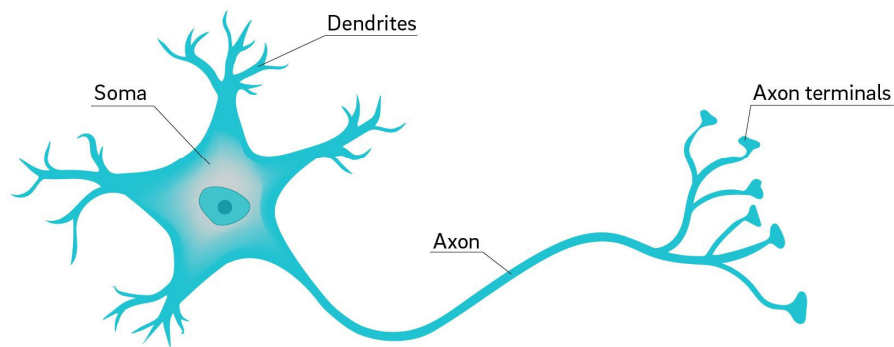
Neuron = linear mapping + nonlinear activation

- Historically, inspired by biological neuron
- Building block of DL
- **Logistic regression** if sigmoid nonlinear activation is used

Picture is from <https://cs231n.github.io/neural-networks-1/>

This concept was invented in 1958

Neuron



The New Yorker, December 6, 1958 P.44

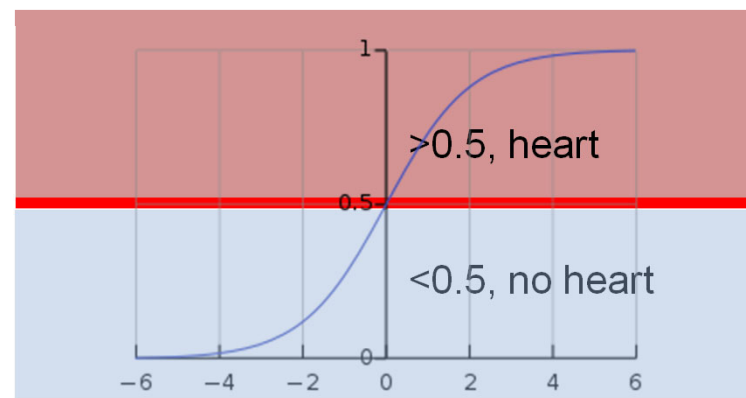
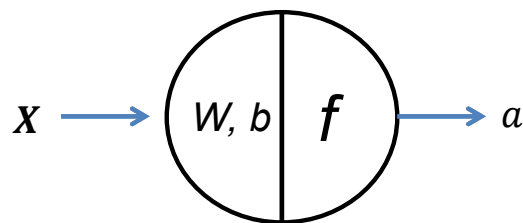
Talk story about the perceptron, a new electronic brain which hasn't been built, but which has been successfully simulated on the I.B.M. 704. Talk with Dr. Frank Rosenblatt, of the Cornell Aeronautical Laboratory, who is one of the two men who developed the prodigy; the other man is Dr. Marshall C. Yovits, of the Office of Naval Research, in Washington. Dr. Rosenblatt defined the perceptron as the first non-biological object which will achieve an organization of its external environment in a meaningful way. It interacts with its environment, forming concepts that have not been made ready for it by a human agent. If a triangle is held up, the perceptron's eye picks up the image & conveys it along a random succession of lines to the response units, where the image is registered. It can tell the difference betw. a cat and a dog,

<https://www.newyorker.com/magazine/1958/12/06/rival-2>

<https://medium.datadriveninvestor.com/introduction-to-neural-networks-a0fe9ec0a947>

Binary classification

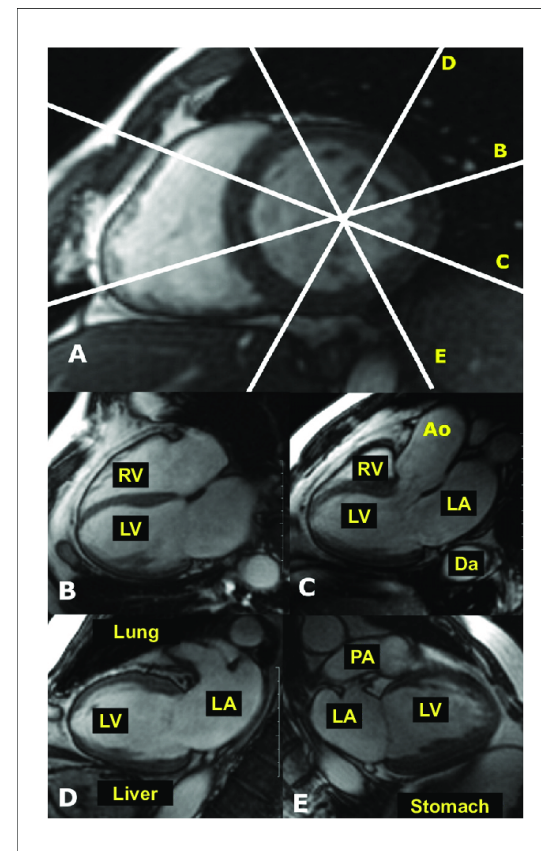
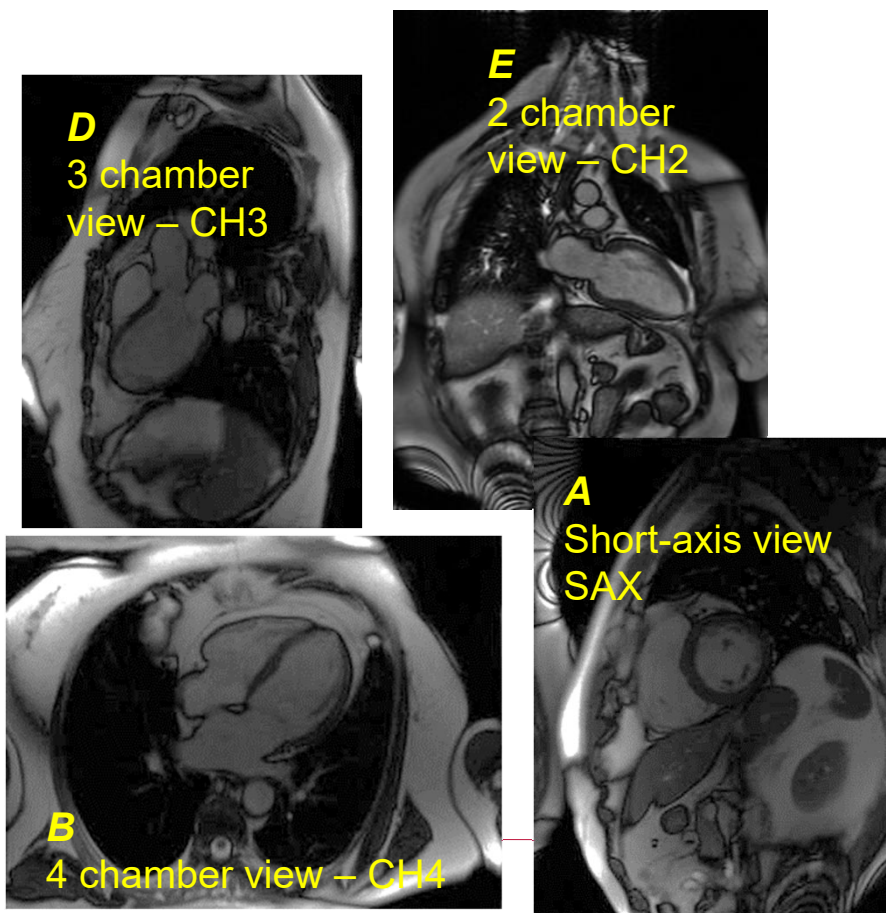
$$a = f(\mathbf{z}) = f(\mathbf{W}\mathbf{X} + b)$$



- Only two output status “Yes” or “No”
- One scalar is needed
- Binary decision

Multi-class classification

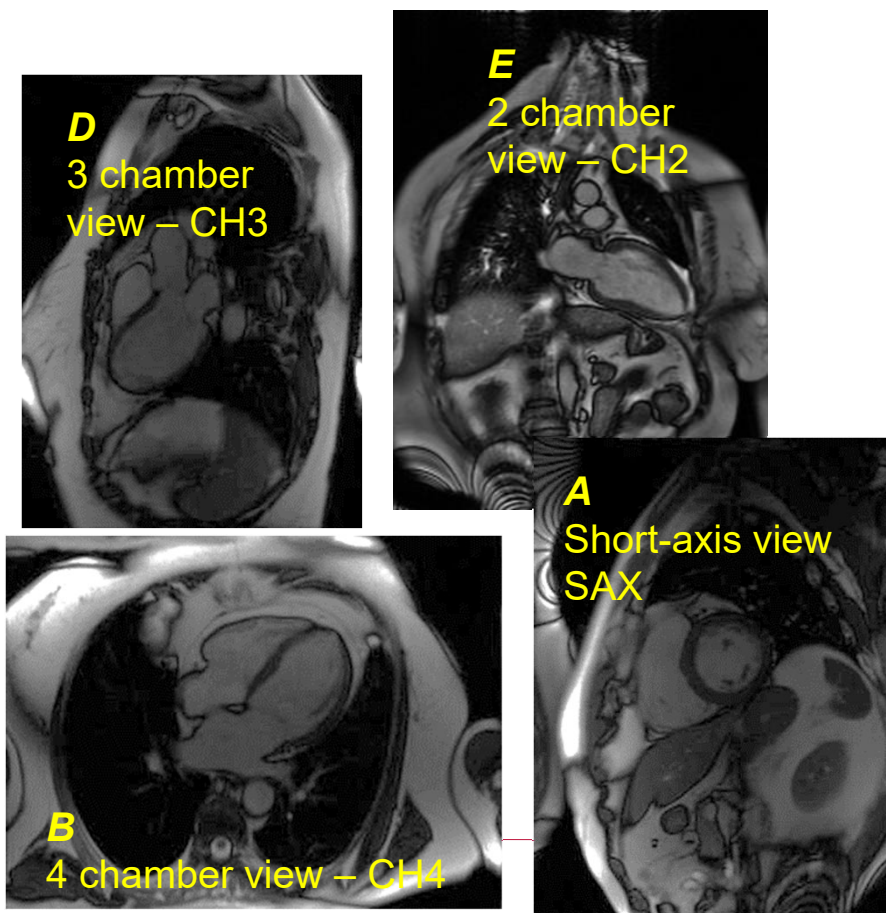
Cardiac view detection



https://www.scielo.br/scielo.php?script=sci_arttext&pid=S0066-782X2010001600014&lng=pt&nrm=iso&tng=pt

Multi-class classification

Cardiac view detection



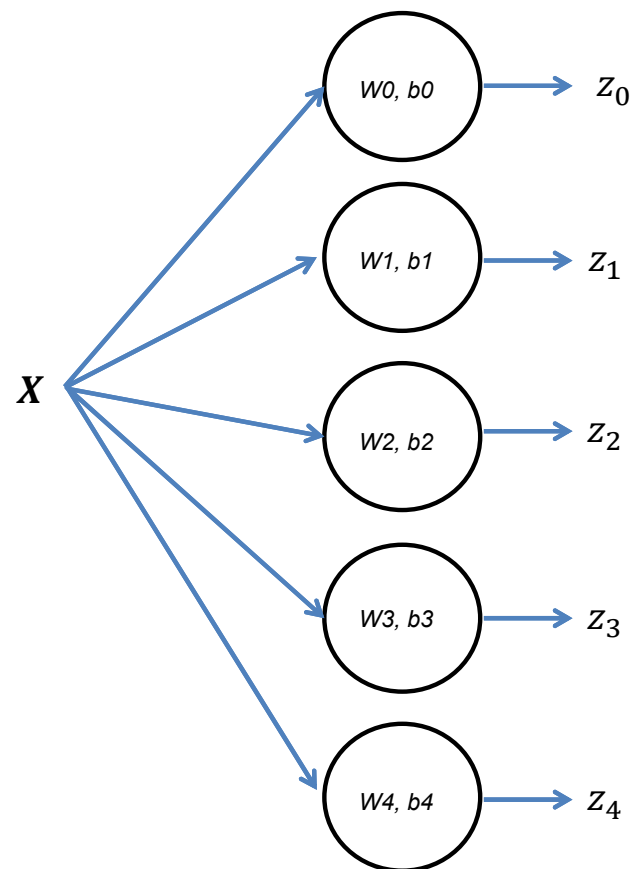
Given an image, decide which imaging view it is

$$\longrightarrow Y = \begin{cases} CH3, & 1 & 0 & 0 & 0 & 0 \\ CH2, & 0 & 1 & 0 & 0 & 0 \\ CH4, & 0 & 0 & 1 & 0 & 0 \\ SAX, & 0 & 0 & 0 & 1 & 0 \\ Other, & 0 & 0 & 0 & 0 & 1 \end{cases}$$

Now we need 5 numbers for K=5 classes

This type of encoding y is “one-hot encoding”

Multi-class classification



Generate 5 numbers as the score for each class

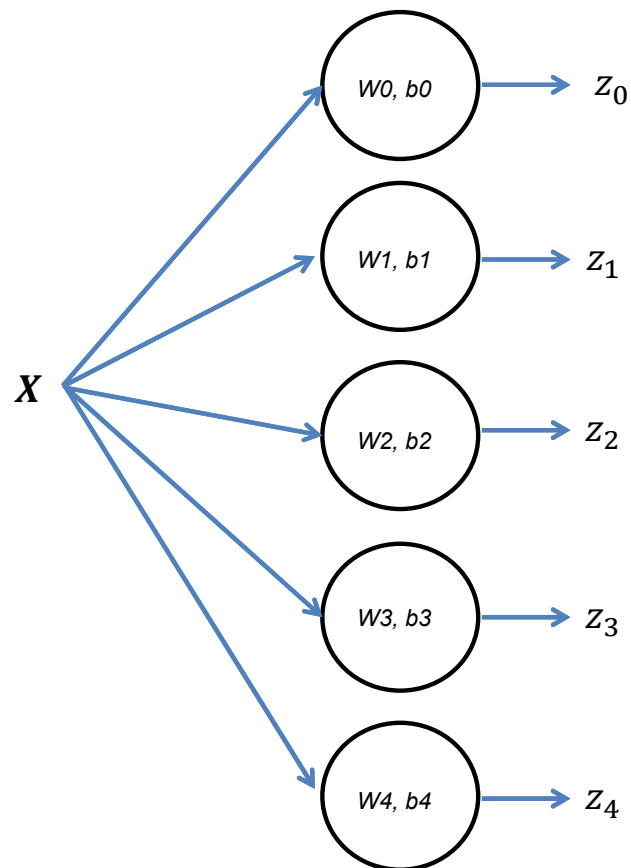
$$\mathbf{Z} = [z_0, z_1, z_2, z_3, z_4]$$

X : $N \times B$, images

W_0, W_1, \dots, W_4 : $1 \times N$ weights

b_0, b_1, \dots, b_4 : 1×1 bias

Multi-class classification



X : $N \times B$, images

W_0, W_1, \dots, W_4 : $1 \times N$ weights

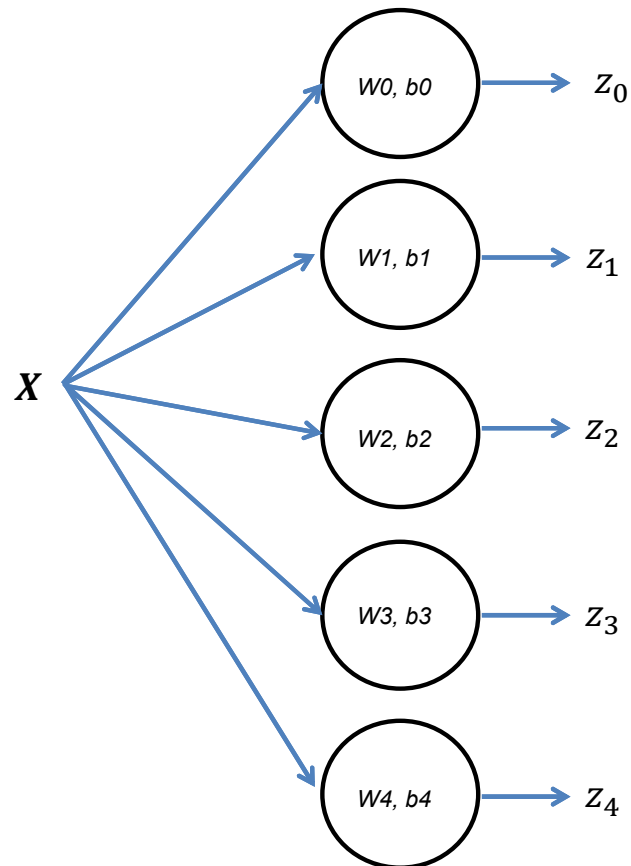
b_0, b_1, \dots, b_4 : 1×1 bias

We can put all weights in one matrix and bias

$$W = \begin{bmatrix} - & W_0 & - \\ - & W_1 & - \\ - & W_2 & - \\ - & W_3 & - \\ - & W_4 & - \end{bmatrix} \quad K \times N \text{ matrix, every row for each class}$$

$$b = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix} \quad K \times 1 \text{ vector}$$

Multi-class classification



Matrix representation:

$$z = WX + b$$

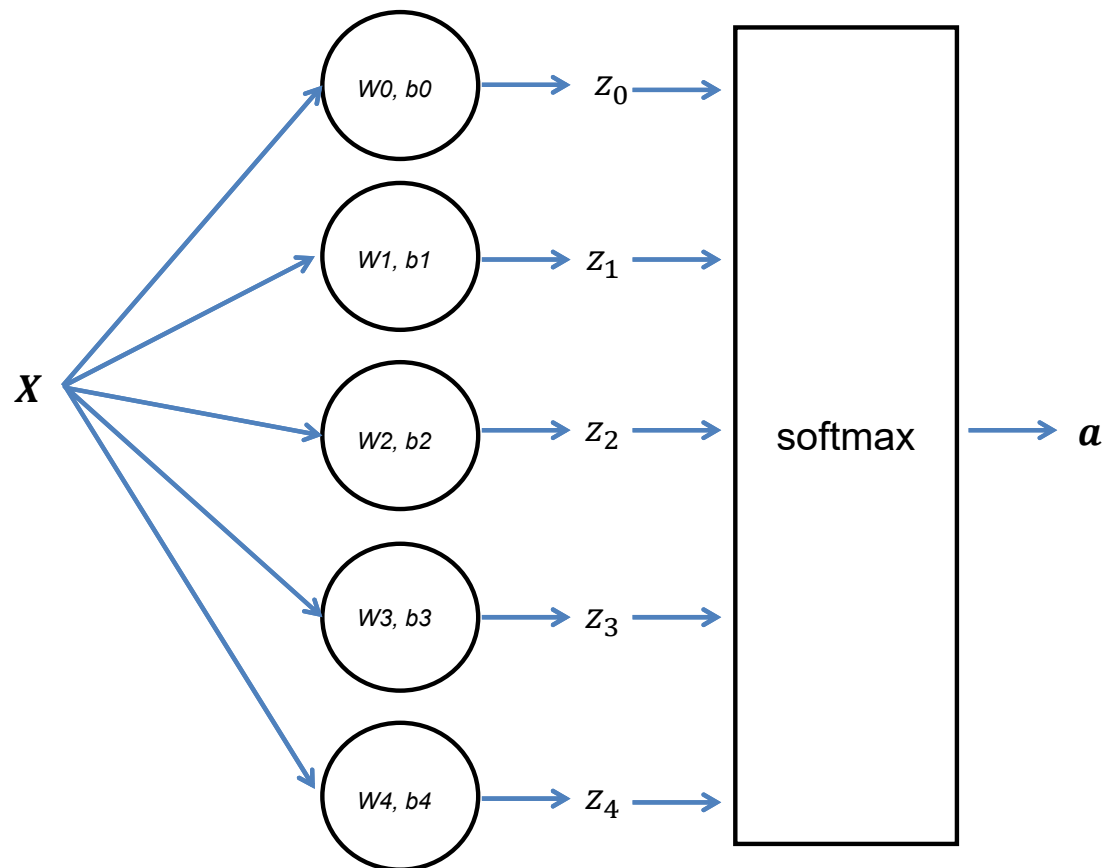
X : $N \times B$, input images

W : $K \times N$, weights

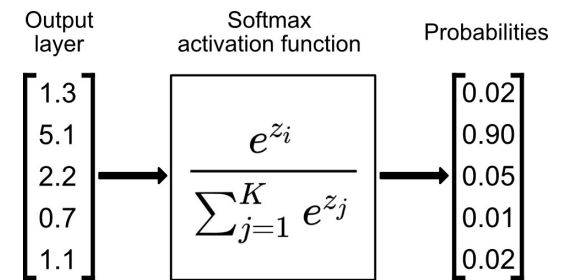
b : $K \times 1$, bias

Z : $K \times B$, score

Softmax



$$\text{softmax}(\mathbf{Z}) = \frac{\exp(z_i)}{\sum_{j=0}^{K-1} \exp(z_j)}$$



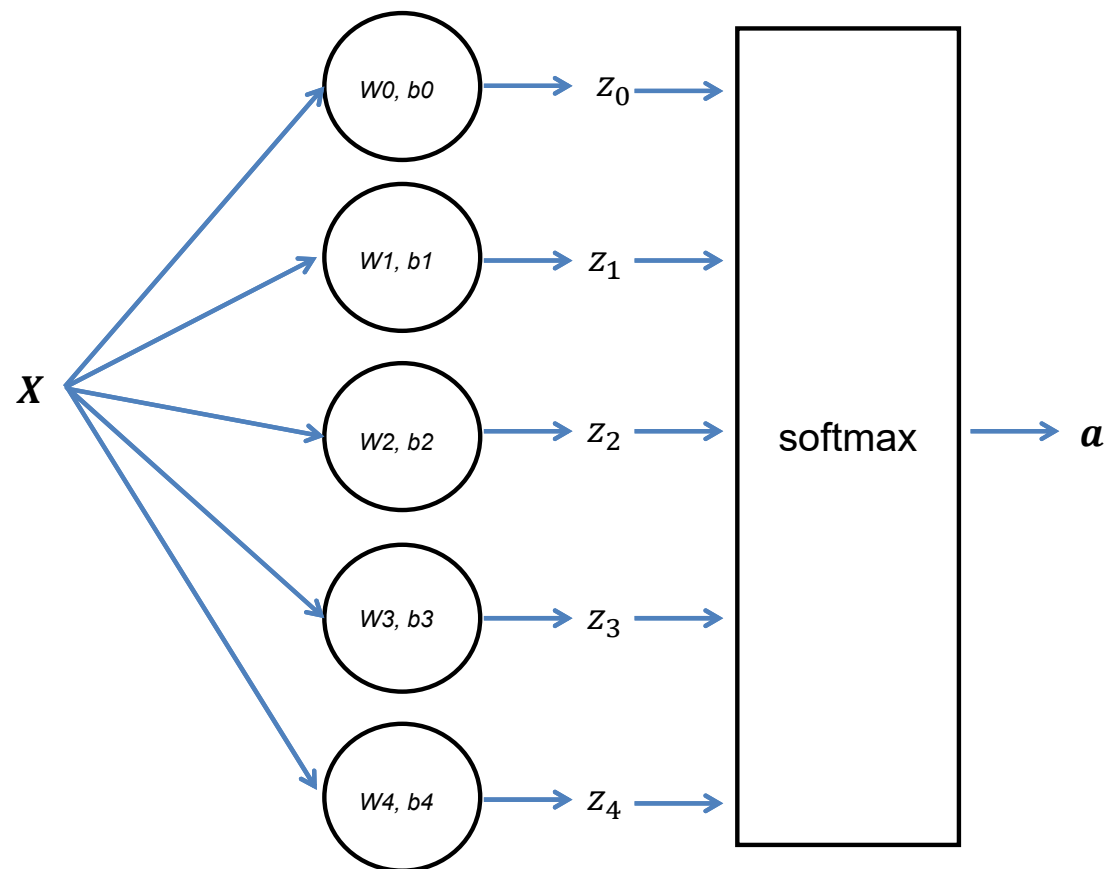
Softmax trick:

$$\text{softmax}(\mathbf{Z} - c) = \text{softmax}(\mathbf{Z})$$

compute $\text{softmax}(\mathbf{Z})$ **as** $\text{softmax}(\mathbf{Z} - \max(\mathbf{Z}))$

Better numeric stability

Logits

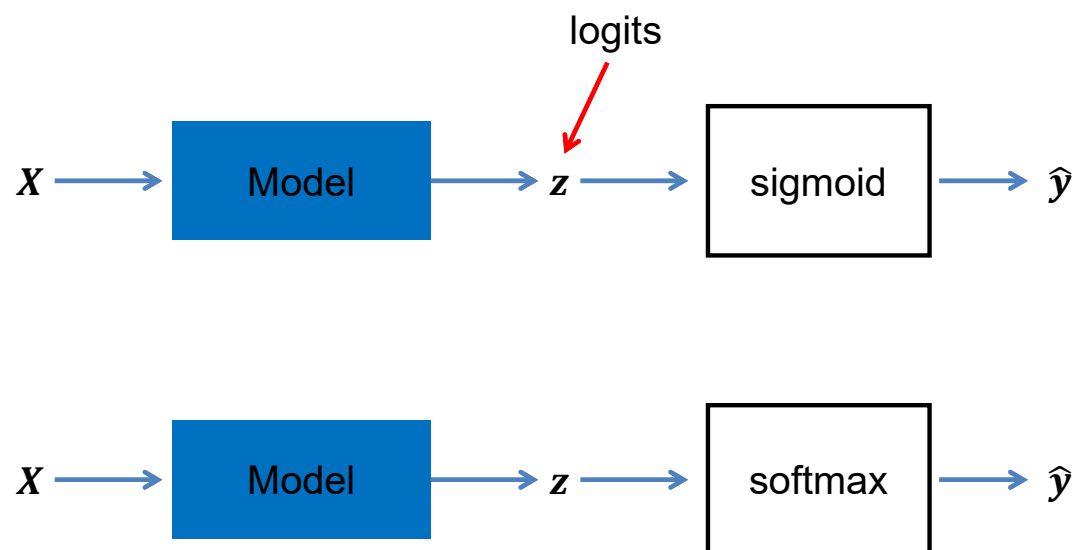


When score \mathbf{Z} is used as inputs to softmax, it is also called “logits”

In the machine learning, logits can be defined as a vector of raw predictions that are inputted into normalization function to generate probabilities.

$$\text{logits}(\mathbf{a}) = \log\left(\frac{a}{1-a}\right)$$
$$a \in [0, 1], \text{logits} \in (-\infty, +\infty)$$

Logits



- The term “logits” is often used in DL frameworks
- Same concept for binary classification
- When models are illustrated, often the sigmoid/softmax are included in the model plotting
- Loss is computed on $L(\mathbf{y}, \hat{\mathbf{y}})$; every sample is (\mathbf{x}, \mathbf{y}) ; model outputs $\hat{\mathbf{y}}$

Binary vs. multi-class classification

Binary

$$\mathbf{a} = \sigma(\mathbf{WX} + \mathbf{b})$$

$\mathbf{X} : \mathbf{N} \times \mathbf{B}$
 $\mathbf{W} : \mathbf{1} \times \mathbf{N}$
 $\mathbf{b} : \mathbf{1} \times \mathbf{1}$
 $\mathbf{a} : \mathbf{1} \times \mathbf{B}$

Apply sigmoid element-wise

Multi-class

$$\mathbf{a} = \textit{softmax}(\mathbf{WX} + \mathbf{b})$$

$\mathbf{X} : \mathbf{N} \times \mathbf{B}$
 $\mathbf{W} : \mathbf{K} \times \mathbf{N}$
 $\mathbf{b} : \mathbf{K} \times \mathbf{1}$
 $\mathbf{a} : \mathbf{K} \times \mathbf{B}$

Apply softmax along
class K

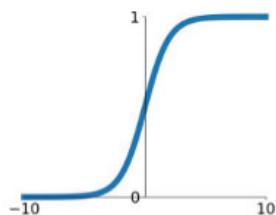
We used one layer of linear mapping and one nonlinear activation function



Activation functions

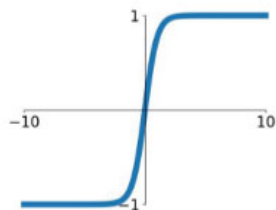
Sigmoid

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



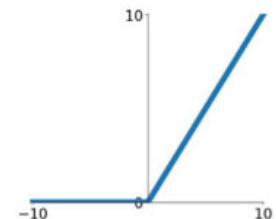
tanh

$$\tanh(x)$$



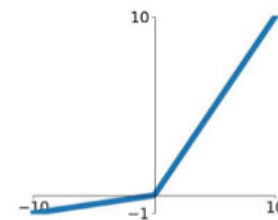
ReLU

$$\max(0, x)$$



Leaky ReLU

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

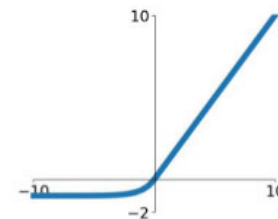


Maxout

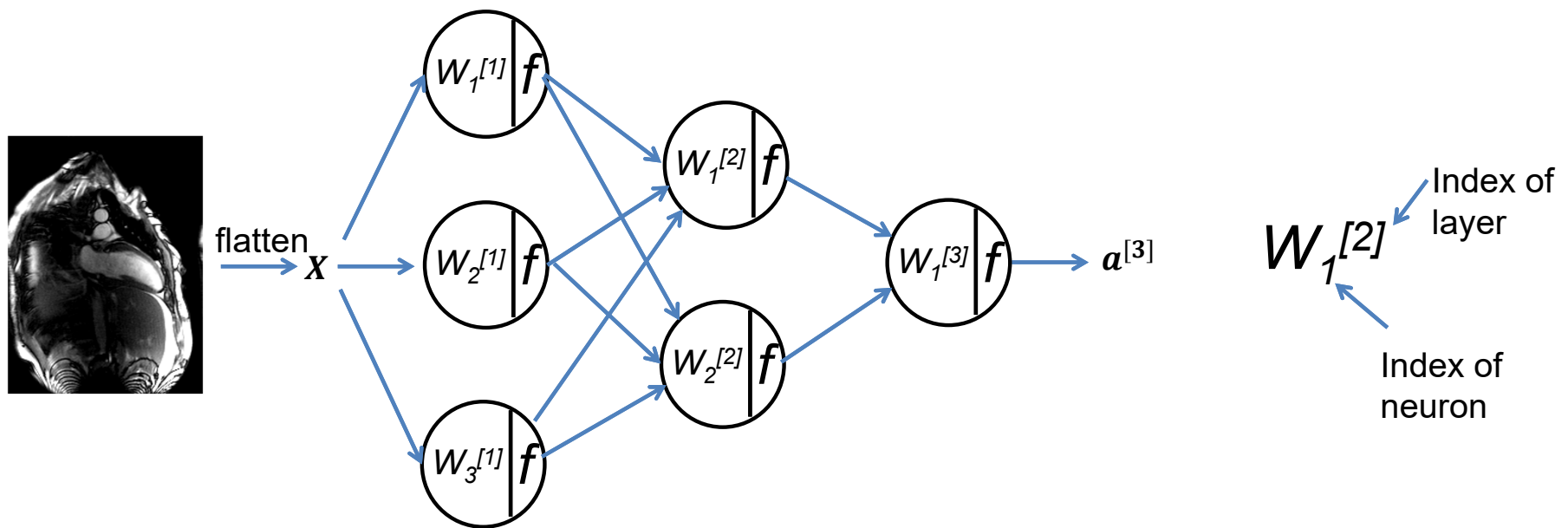
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

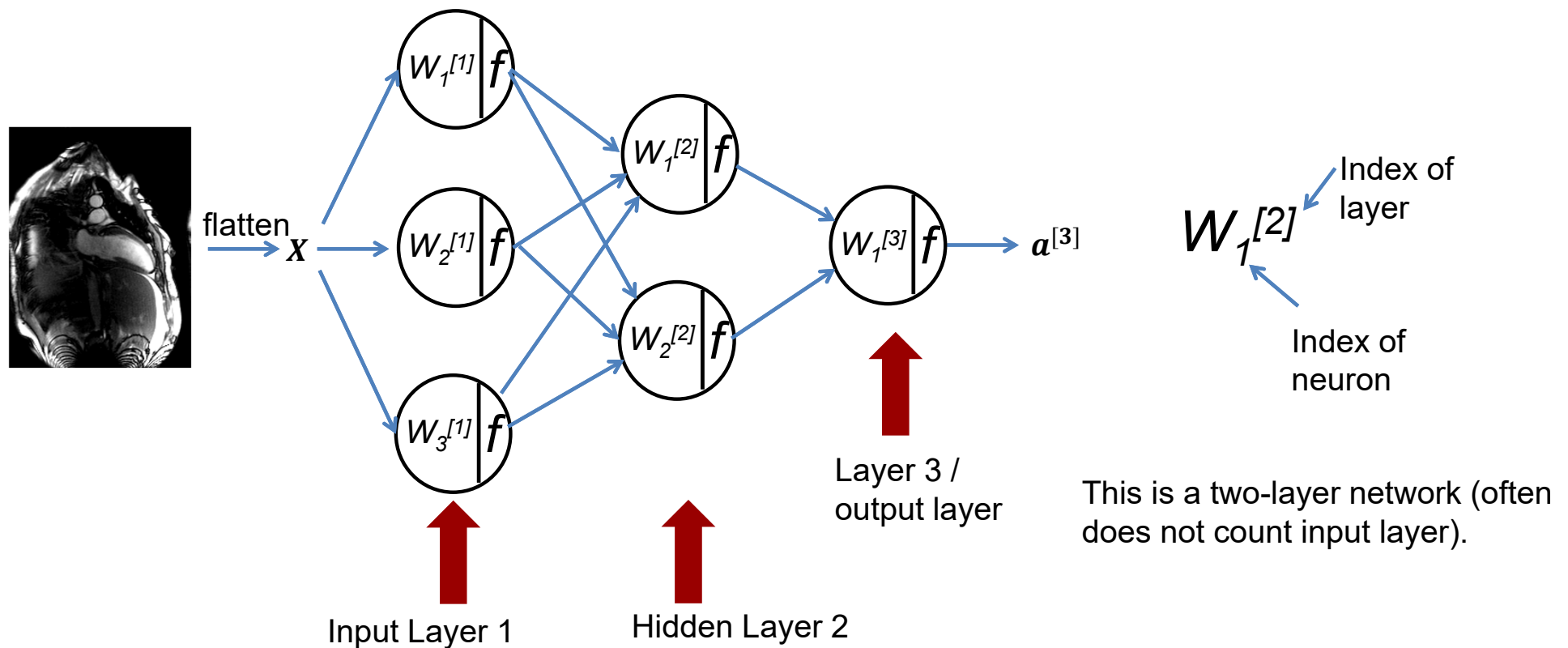
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



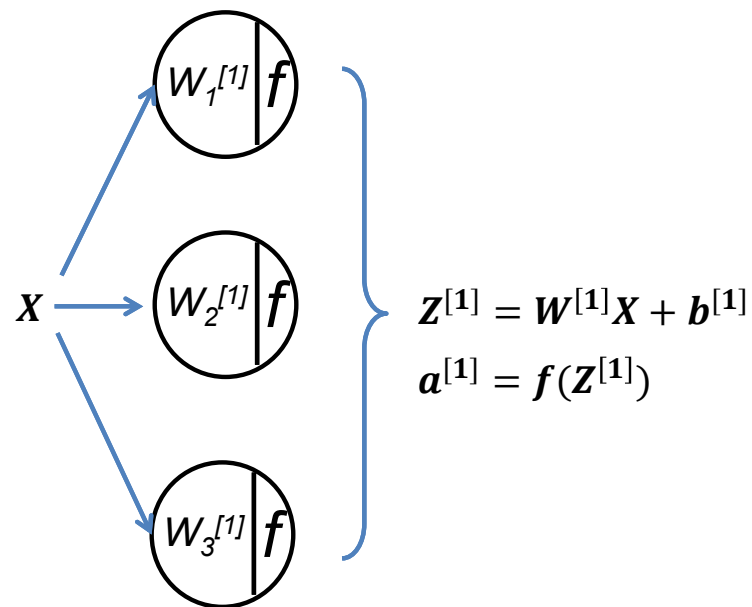
Multi-layer Perceptron (MLP)



Multi-layer Perceptron (MLP)



Multi-layer Perceptron (MLP)



$X : N \times 1$

$W^{[1]} : 3 \times N$

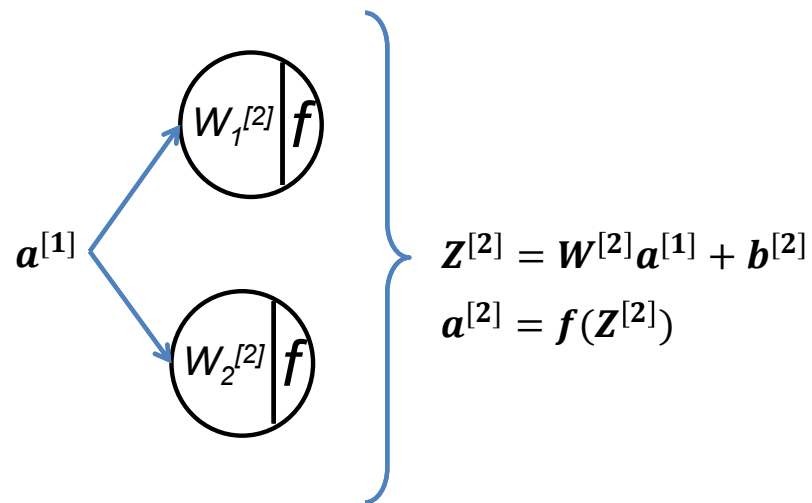
$b^{[1]} : 3 \times 1$

$Z^{[1]} : 3 \times 1$, score

$a^{[1]} : 3 \times 1$, activation

f : applied element-wise

Multi-layer Perceptron (MLP)



$$a^{[1]} : 3 \times 1$$

$$W^{[2]} : 2 \times 3$$

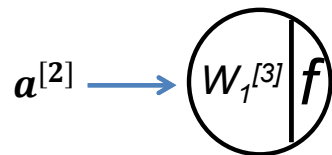
$$b^{[2]} : 2 \times 1$$

$$Z^{[2]} : 2 \times 1$$

$$a^{[2]} : 2 \times 1$$

f : applied element-wise

Multi-layer Perceptron (MLP)



$$\mathbf{Z}^{[3]} = \mathbf{W}^{[3]} \mathbf{a}^{[2]} + \mathbf{b}^{[3]}$$

$$\mathbf{a}^{[3]} = f(\mathbf{Z}^{[3]})$$

$$\mathbf{a}^{[2]} : 2 \times 1$$

$$\mathbf{W}^{[3]} : 1 \times 2$$

$$\mathbf{b}^{[3]} : 1 \times 1$$

$$\mathbf{Z}^{[3]} : 1 \times 1$$

$$\mathbf{a}^{[3]} : 1 \times 1$$

f : applied element-wise

Forward Pass

$$\mathbf{z}^{[1]} = \mathbf{W}^{[1]}\mathbf{X} + \mathbf{b}^{[1]}$$

$$\mathbf{a}^{[1]} = f(\mathbf{z}^{[1]})$$

$$\mathbf{z}^{[2]} = \mathbf{W}^{[2]}\mathbf{a}^{[1]} + \mathbf{b}^{[2]}$$

$$\mathbf{a}^{[2]} = f(\mathbf{z}^{[2]})$$

$$\mathbf{z}^{[3]} = \mathbf{W}^{[3]}\mathbf{a}^{[2]} + \mathbf{b}^{[3]}$$

$$\hat{\mathbf{y}} = \mathbf{a}^{[3]} = f(\mathbf{z}^{[3]})$$

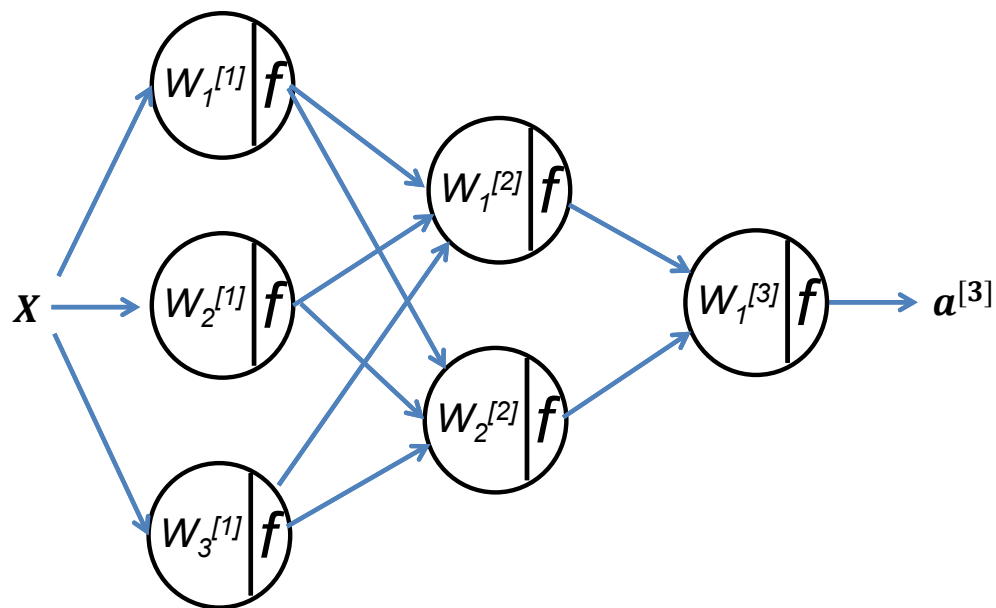
Forward pass: from input \mathbf{X} to output $\hat{\mathbf{y}}$

For every layer,

\mathbf{W} .shape is [number of output, number of input]

\mathbf{b} .shape is [number of output, 1]

Number of parameters?



$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

$$a^{[1]} = f(Z^{[1]})$$

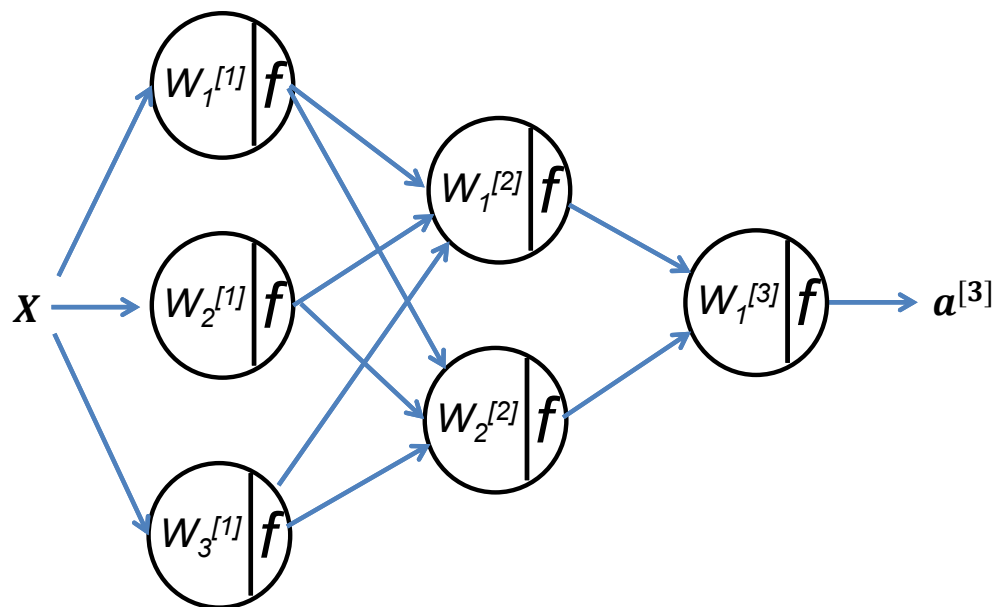
$$Z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = f(Z^{[2]})$$

$$Z^{[3]} = W^{[3]}a^{[2]} + b^{[3]}$$

$$\hat{y} = a^{[3]} = f(Z^{[3]})$$

Number of parameters?



$$Z^{[1]} = W^{[1]}X + b^{[1]} \quad 3N+3$$

$$a^{[1]} = f(Z^{[1]})$$

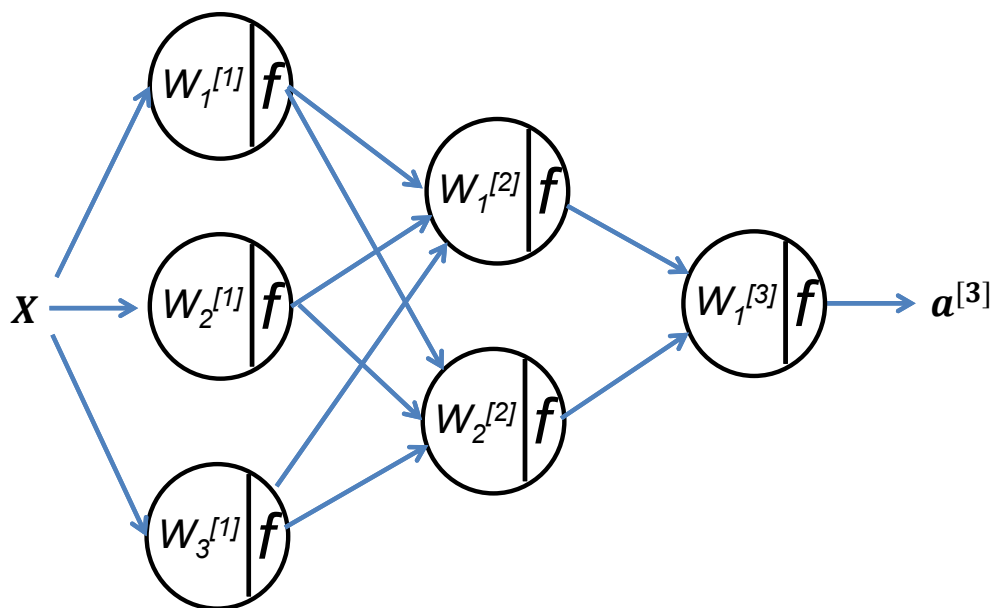
$$Z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = f(Z^{[2]})$$

$$Z^{[3]} = W^{[3]}a^{[2]} + b^{[3]}$$

$$\hat{y} = a^{[3]} = f(Z^{[3]})$$

Number of parameters?



$$Z^{[1]} = W^{[1]}X + b^{[1]} \quad 3N+3$$

$$a^{[1]} = f(Z^{[1]}) \quad 0$$

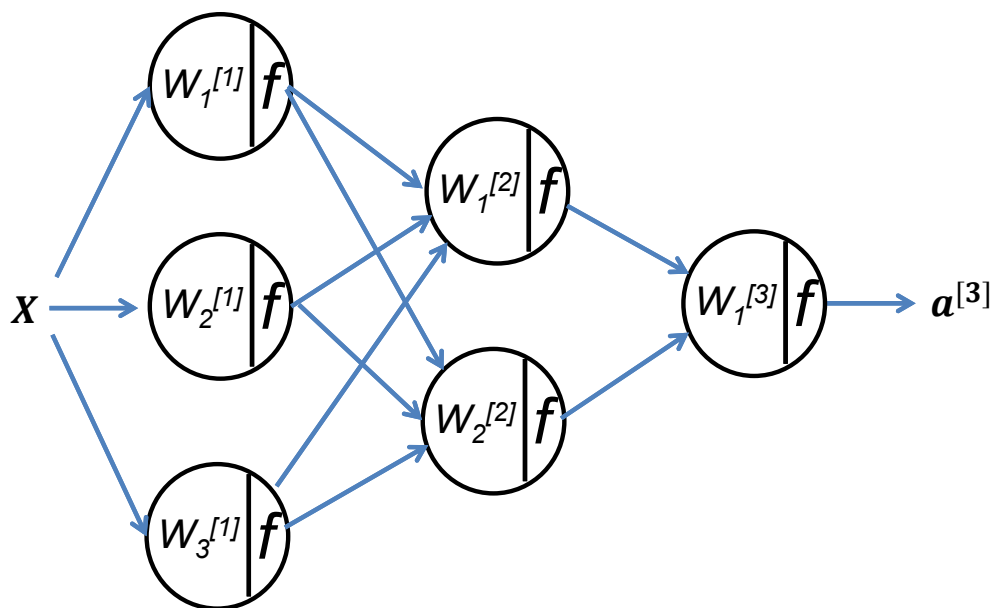
$$Z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = f(Z^{[2]})$$

$$Z^{[3]} = W^{[3]}a^{[2]} + b^{[3]}$$

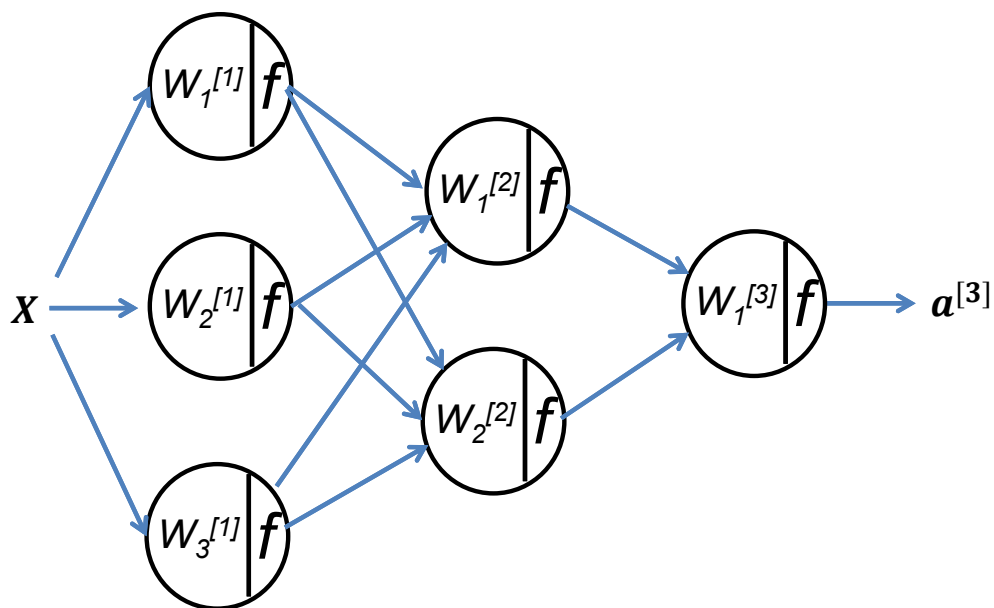
$$\hat{y} = a^{[3]} = f(Z^{[3]})$$

Number of parameters?



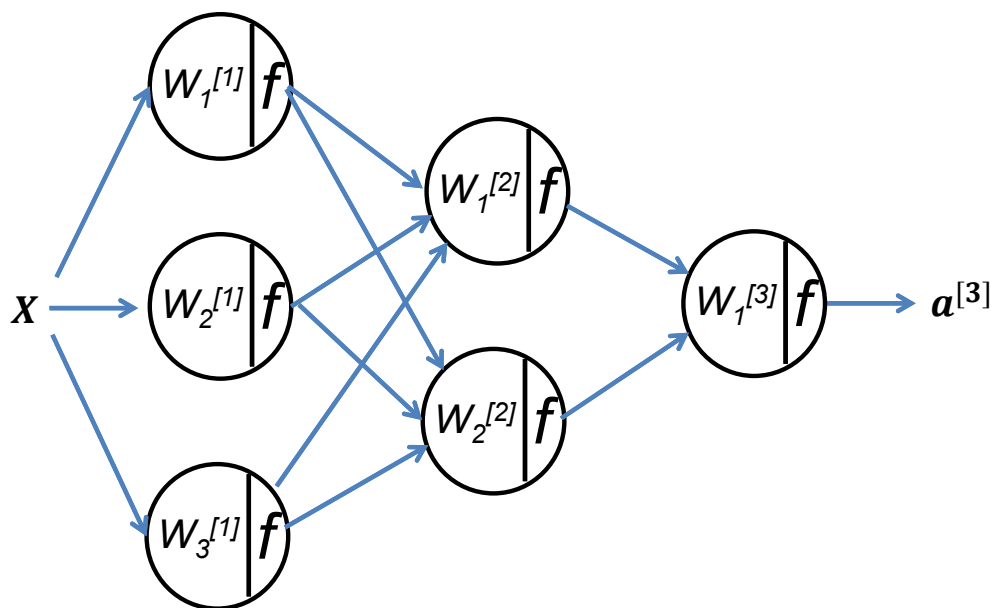
$$\begin{aligned} Z^{[1]} &= W^{[1]}X + b^{[1]} && 3N+3 \\ a^{[1]} &= f(Z^{[1]}) && 0 \\ Z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} && 2 \times 3 + 2 = 8 \\ a^{[2]} &= f(Z^{[2]}) \\ Z^{[3]} &= W^{[3]}a^{[2]} + b^{[3]} \\ \hat{y} &= a^{[3]} = f(Z^{[3]}) \end{aligned}$$

Number of parameters?



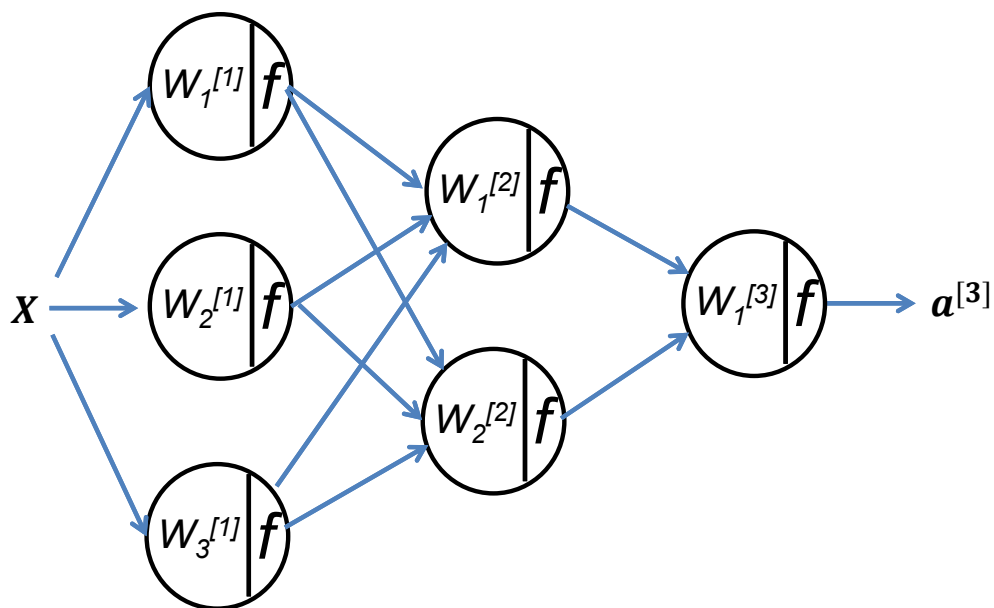
$$\begin{aligned} Z^{[1]} &= W^{[1]}X + b^{[1]} & 3N+3 \\ a^{[1]} &= f(Z^{[1]}) & 0 \\ Z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} & 2 \times 3 + 2 = 8 \\ a^{[2]} &= f(Z^{[2]}) & 0 \\ Z^{[3]} &= W^{[3]}a^{[2]} + b^{[3]} \\ \hat{y} &= a^{[3]} = f(Z^{[3]}) \end{aligned}$$

Number of parameters?



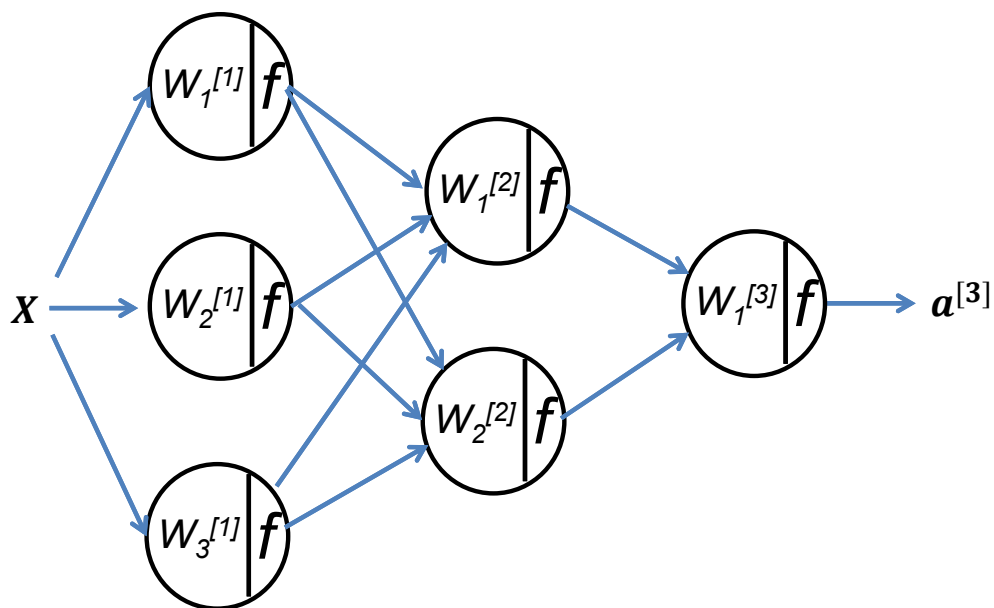
$$\begin{aligned} Z^{[1]} &= W^{[1]}X + b^{[1]} && 3N+3 \\ a^{[1]} &= f(Z^{[1]}) && 0 \\ Z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} && 2 \times 3 + 2 = 8 \\ a^{[2]} &= f(Z^{[2]}) && 0 \\ Z^{[3]} &= W^{[3]}a^{[2]} + b^{[3]} && 1 \times 2 + 1 = 3 \\ \hat{y} &= a^{[3]} = f(Z^{[3]}) \end{aligned}$$

Number of parameters?



$$\begin{aligned} Z^{[1]} &= W^{[1]}X + b^{[1]} & 3N+3 \\ a^{[1]} &= f(Z^{[1]}) & 0 \\ Z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} & 2 \times 3 + 2 = 8 \\ a^{[2]} &= f(Z^{[2]}) & 0 \\ Z^{[3]} &= W^{[3]}a^{[2]} + b^{[3]} & 1 \times 2 + 1 = 3 \\ \hat{y} = a^{[3]} &= f(Z^{[3]}) & 0 \end{aligned}$$

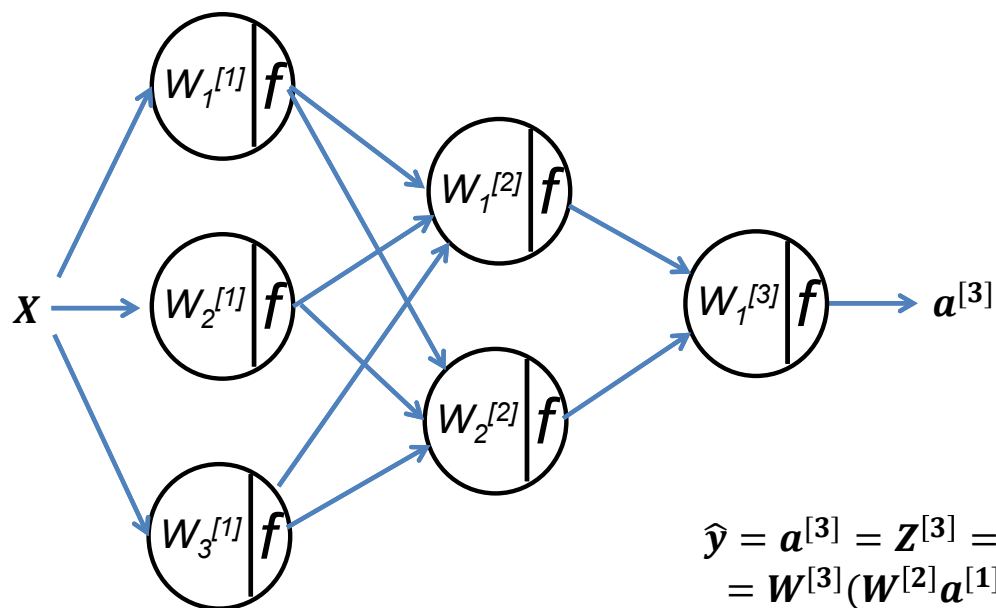
Number of parameters?



$$\begin{aligned}
 Z^{[1]} &= W^{[1]}X + b^{[1]} & 3N+3 \\
 a^{[1]} &= f(Z^{[1]}) & 0 \\
 Z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} & 2 \times 3 + 2 = 8 \\
 a^{[2]} &= f(Z^{[2]}) & 0 \\
 Z^{[3]} &= W^{[3]}a^{[2]} + b^{[3]} & 1 \times 2 + 1 = 3 \\
 \hat{y} = a^{[3]} &= f(Z^{[3]}) & 0
 \end{aligned}$$

$$3N+14$$

Nonlinear activation is necessary



if $f(x) = x$

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

$$a^{[1]} = Z^{[1]}$$

$$Z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = Z^{[2]}$$

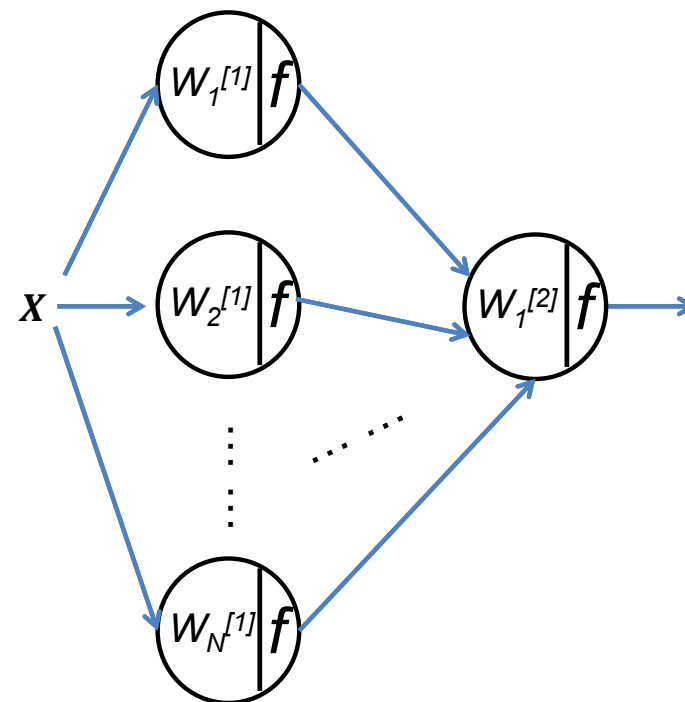
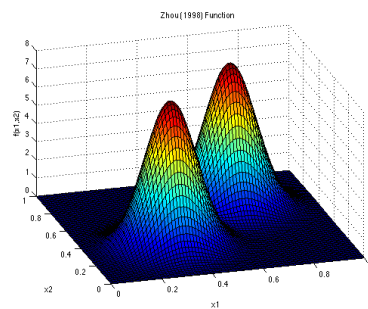
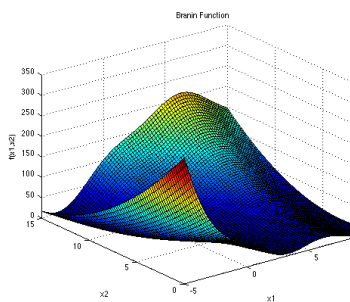
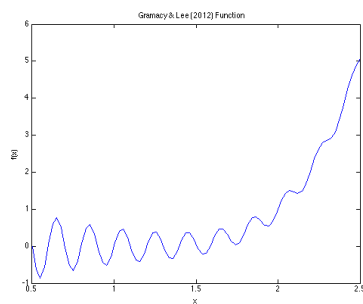
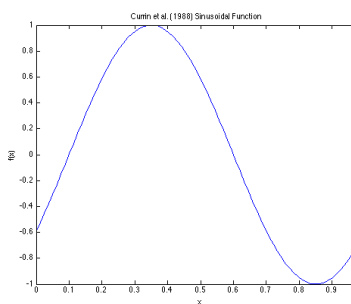
$$Z^{[3]} = W^{[3]}a^{[2]} + b^{[3]}$$

$$\begin{aligned}\hat{y} &= a^{[3]} = Z^{[3]} = W^{[3]}a^{[2]} + b^{[3]} \\ &= W^{[3]}(W^{[2]}a^{[1]} + b^{[2]}) + b^{[3]} \\ &= W^{[3]}(W^{[2]}(W^{[1]}X + b^{[1]}) + b^{[2]}) + b^{[3]} \\ &= W^{[3]}W^{[2]}W^{[1]}X + W^{[3]}W^{[2]}b^{[1]} + W^{[3]}b^{[2]} + b^{[3]}\end{aligned}$$

A linear model ...



Universal Approximation

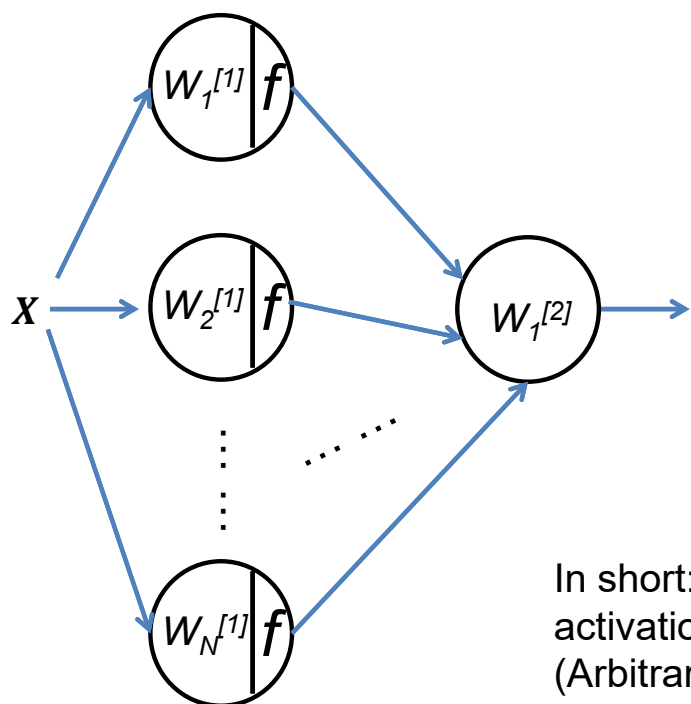


Can a MLP approximate any functions?



Universal Approximation

Can a MLP approximate any functions?



Approximation Capabilities of Multilayer Feedforward Networks

KURT HORNIK

Technische Universität Wien, Vienna, Austria

(Received 30 January 1990; revised and accepted 25 October 1990)

Abstract—We show that standard multilayer feedforward networks with as few as a single hidden layer and arbitrary bounded and nonconstant activation function are universal approximators with respect to $L^p(\mu)$ performance criteria, for arbitrary finite input environment measures μ , provided only that sufficiently many hidden units are available. If the activation function is continuous, bounded and nonconstant, then continuous mappings can be learned uniformly over compact input sets. We also give very general conditions ensuring that networks with sufficiently smooth activation functions are capable of arbitrarily accurate approximation to a function and its derivatives.

[https://doi.org/10.1016/0893-6080\(91\)90009-T](https://doi.org/10.1016/0893-6080(91)90009-T)

In short: A MLP with at least one hidden layer and nonlinear activation function can approximate any continuous function (Arbitrary Width).

Similar conclusions exist for arbitrary depth, CNN and other conditions

https://en.wikipedia.org/wiki/Universal_approximation_theorem

