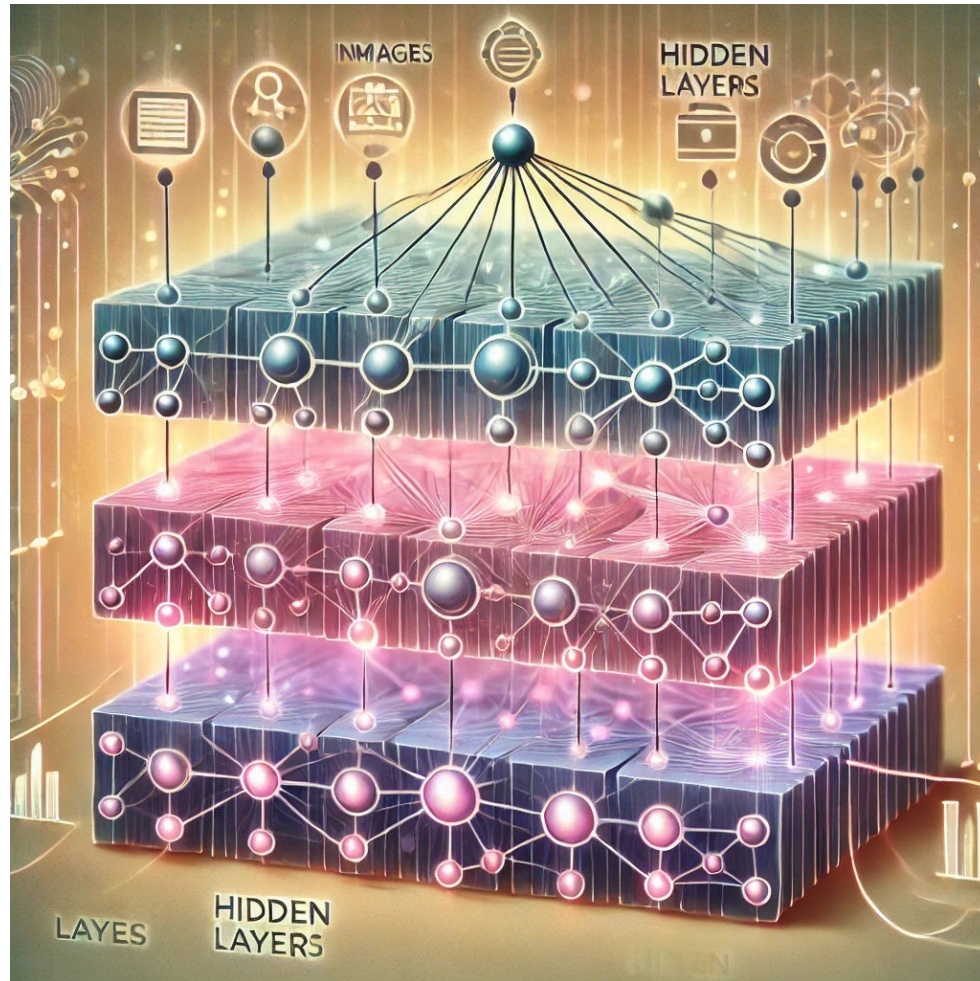


# Neural Networks

---



AI604 Deep Learning for Computer Vision

Prof. Hyunjung Shim

Slide credit: Justin Johnson, Fei-Fei Li, Ehsan Adeli

# Recap: Classification/Loss/Optimizer

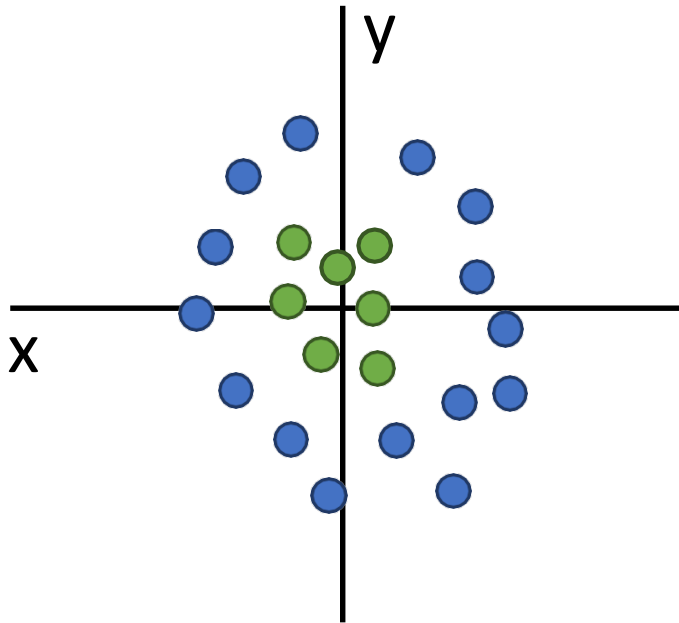
---

1. Use **Linear Models** for image classification problems (three perspectives)
2. Use **Loss Functions** (SVM, Softmax) to express preferences over different choices of weights
3. Use **Stochastic Gradient Descent** to minimize our loss functions and train the model

# Problem: Linear Classifiers aren't that powerful

---

## Geometric Viewpoint



## Visual Viewpoint

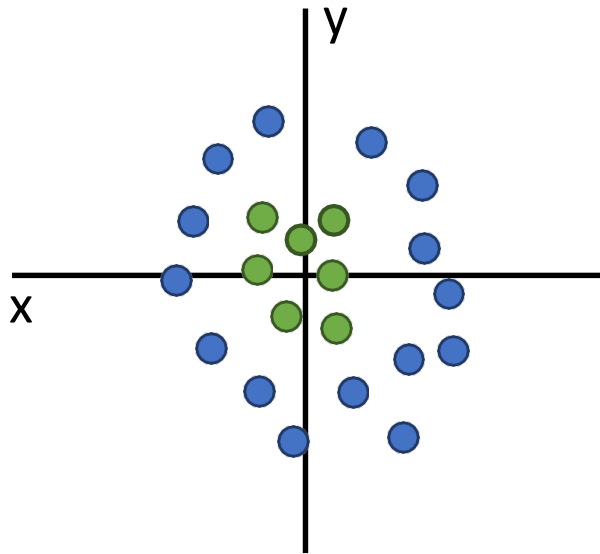
One template per class: Can't recognize different modes of a class



# One solution: Feature Transforms

---

Original space



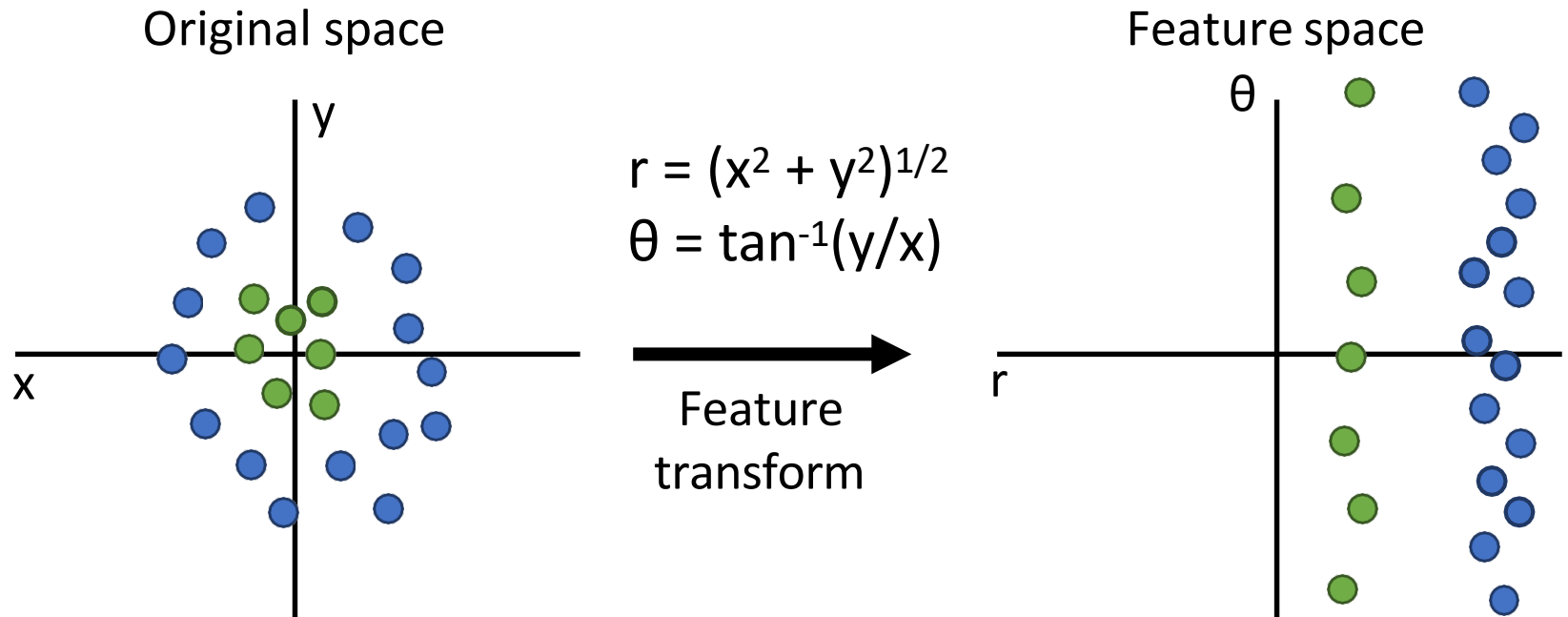
$$r = (x^2 + y^2)^{1/2}$$
$$\theta = \tan^{-1}(y/x)$$



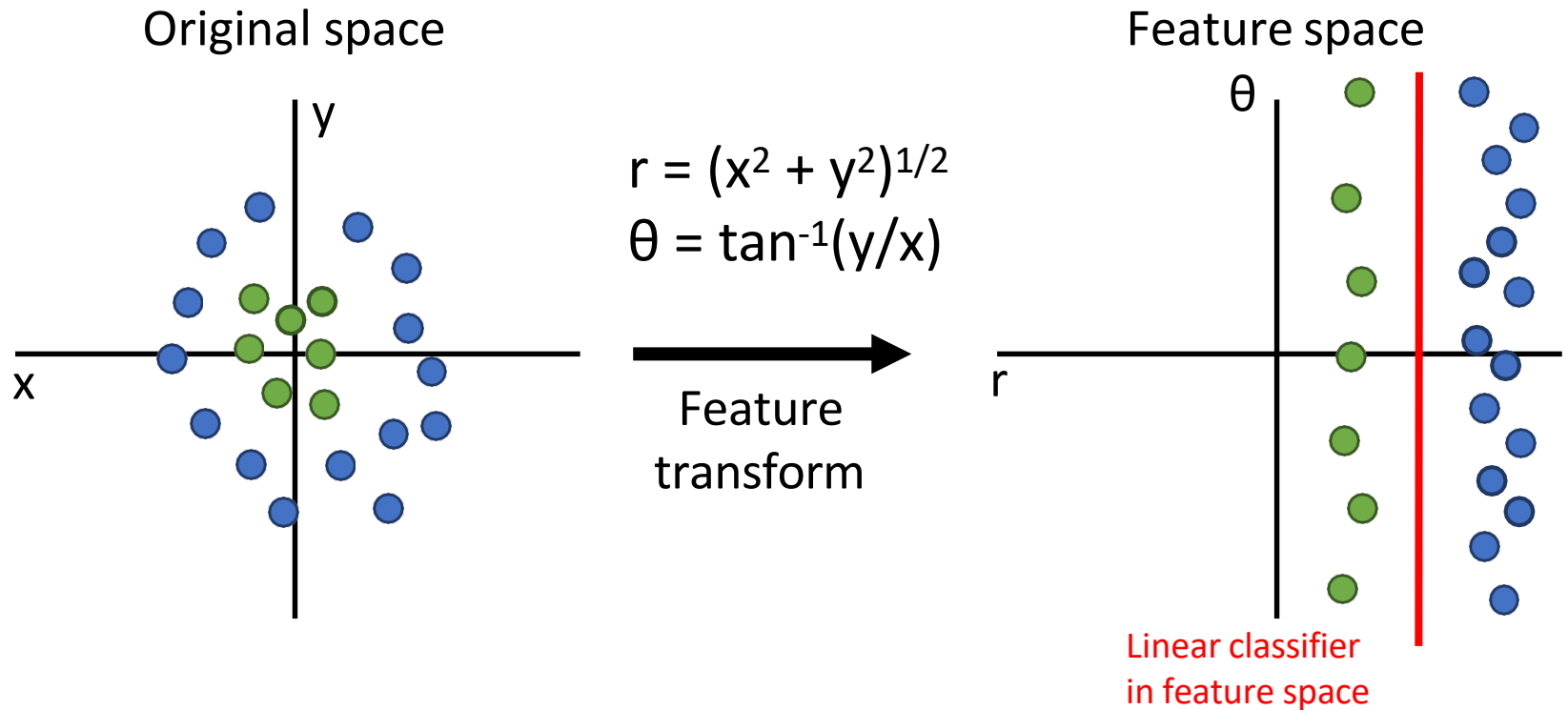
Feature  
transform

# One solution: Feature Transforms

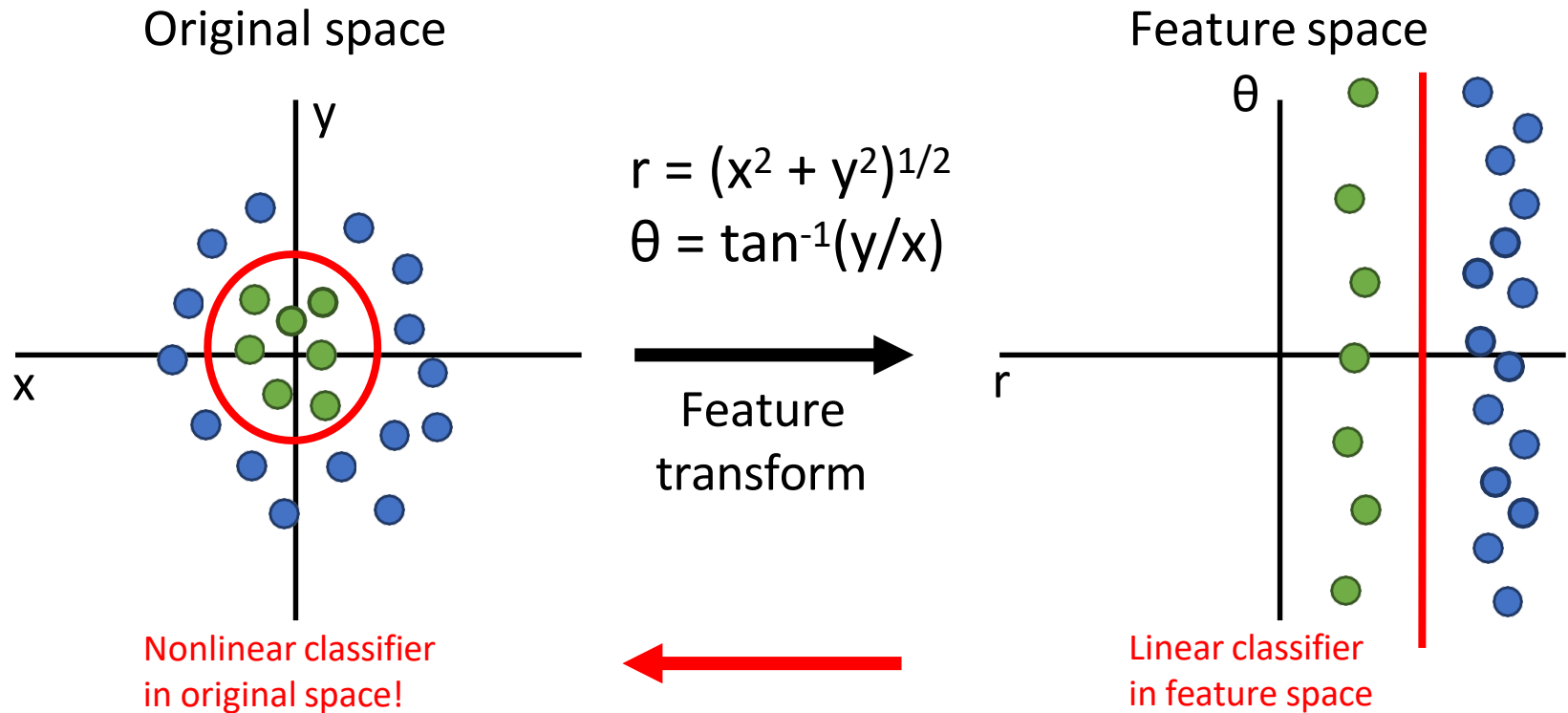
---



# One solution: Feature Transforms



# One solution: Feature Transforms

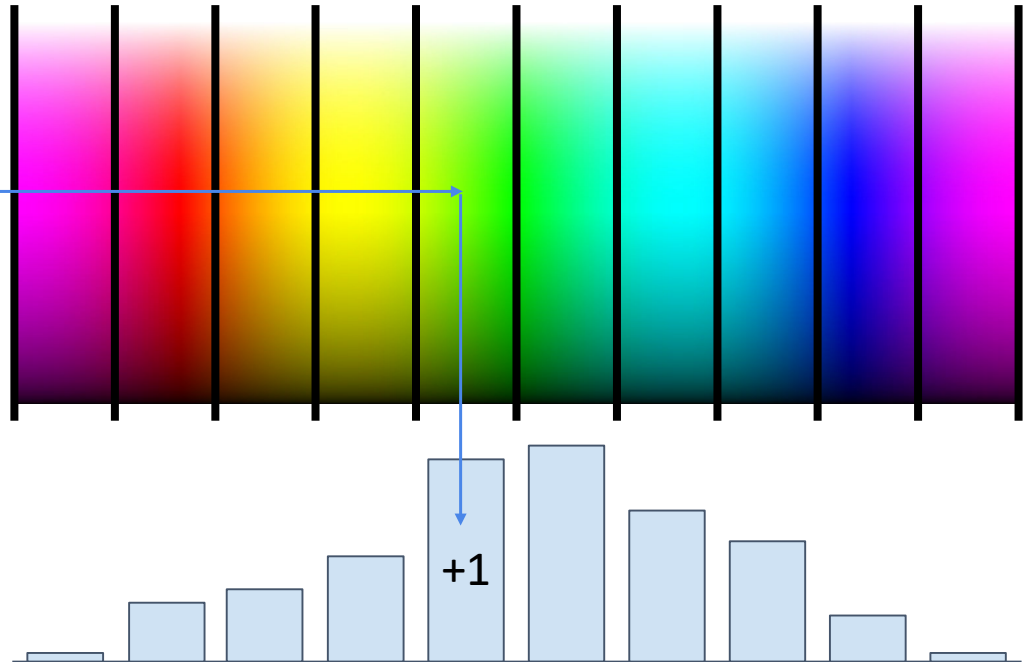


# Image Features: Color Histogram

---



Ignores texture,  
spatial positions



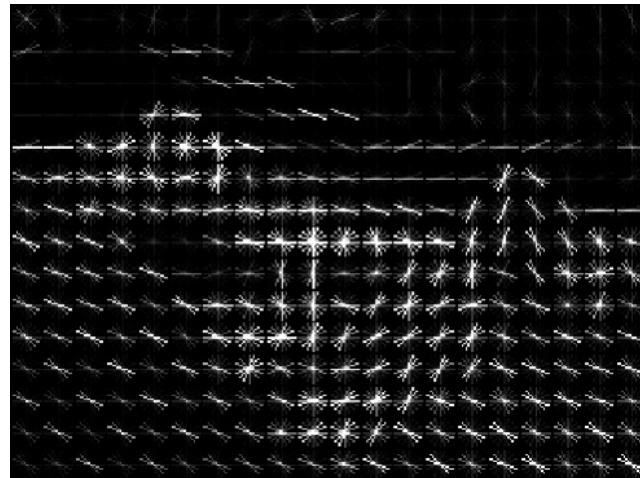


# Image Features: Histogram of Oriented Gradients (HoG)



1. Compute edge direction / strength at each pixel
2. Divide image into 8x8 regions
3. Within each region compute a histogram of edge directions weighted by edge strength

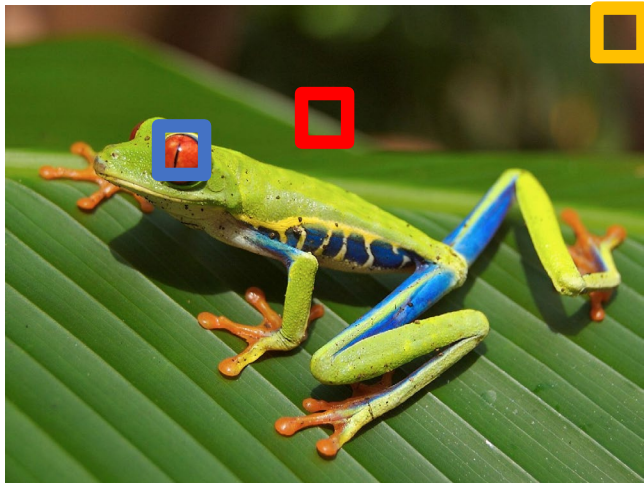
# Image Features: Histogram of Oriented Gradients (HoG)



1. Compute edge direction / strength at each pixel
2. Divide image into 8x8 regions
3. Within each region compute a histogram of edge directions weighted by edge strength

Example: 320x240 image gets divided into 40x30 bins; 9 directions per bin; feature vector has  $30 \times 40 \times 9 = 10,800$  numbers

# Image Features: Histogram of Oriented Gradients (HoG)

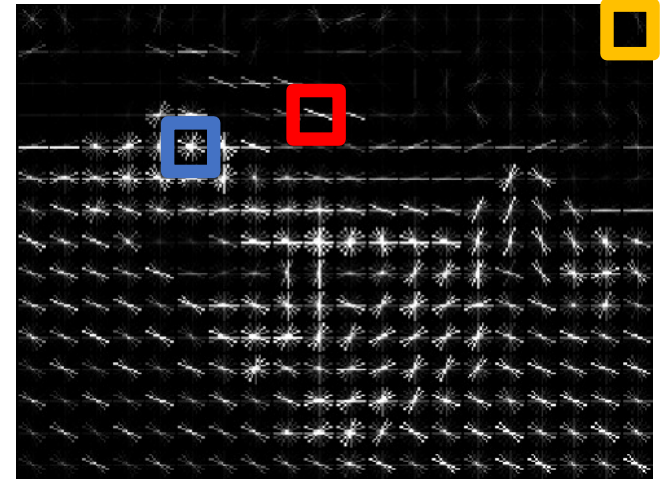


Weak edges

Strong diagonal  
edges



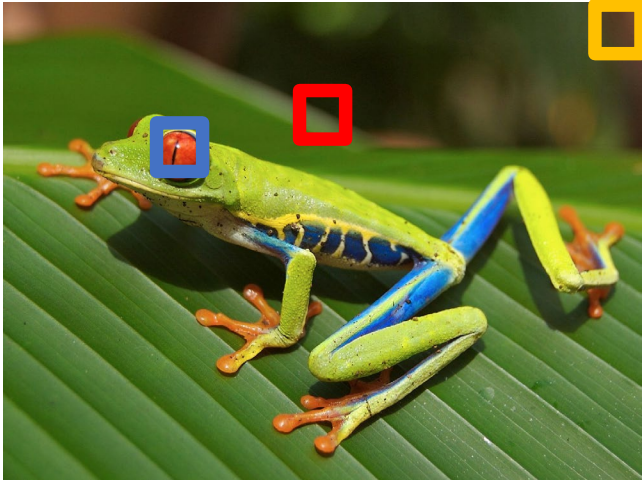
Edges in all  
directions



1. Compute edge direction / strength at each pixel
2. Divide image into 8x8 regions
3. Within each region compute a histogram of edge directions weighted by edge strength

Example: 320x240 image gets divided into 40x30 bins; 9 directions per bin; feature vector has  $30 \times 40 \times 9 = 10,800$  numbers

# Image Features: Histogram of Oriented Gradients (HoG)



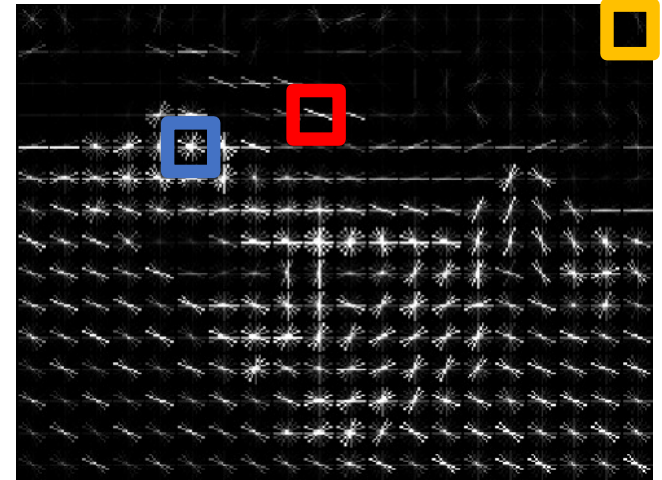
Weak edges

Strong diagonal edges



Edges in all directions

Captures texture and position, robust to small image changes



1. Compute edge direction / strength at each pixel
2. Divide image into 8x8 regions
3. Within each region compute a histogram of edge directions weighted by edge strength

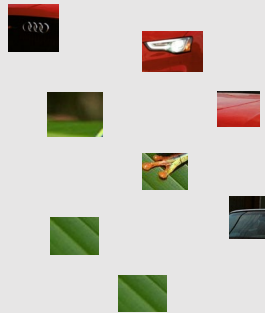
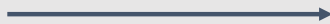
Example: 320x240 image gets divided into 40x30 bins; 9 directions per bin; feature vector has  $30 \times 40 \times 9 = 10,800$  numbers

# Image Features: Bag of Words (Data-Driven!)

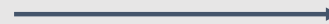
## Step 1: Build codebook



Extract random  
patches



Cluster patches to  
form “codebook”  
of “visual words”

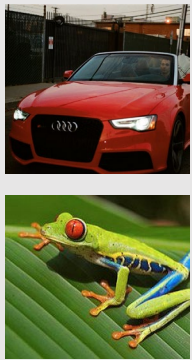


Fei-Fei and Perona, “A bayesian hierarchical model for learning natural scene categories”, CVPR 2005

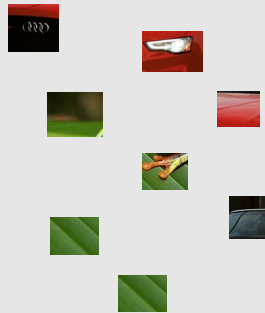


# Image Features: Bag of Words (Data-Driven!)

## Step 1: Build codebook



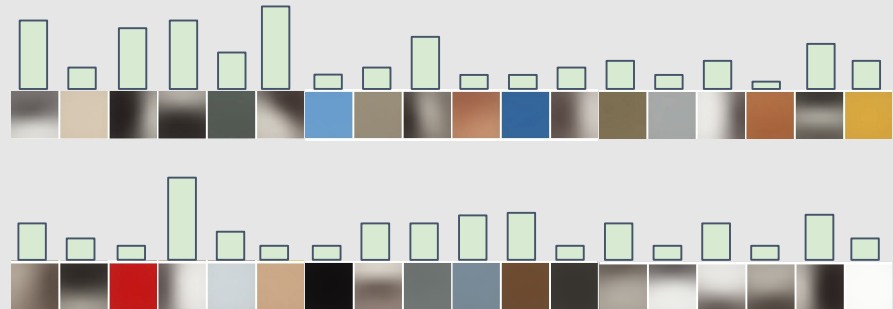
Extract random  
patches



Cluster patches to  
form “codebook”  
of “visual words”



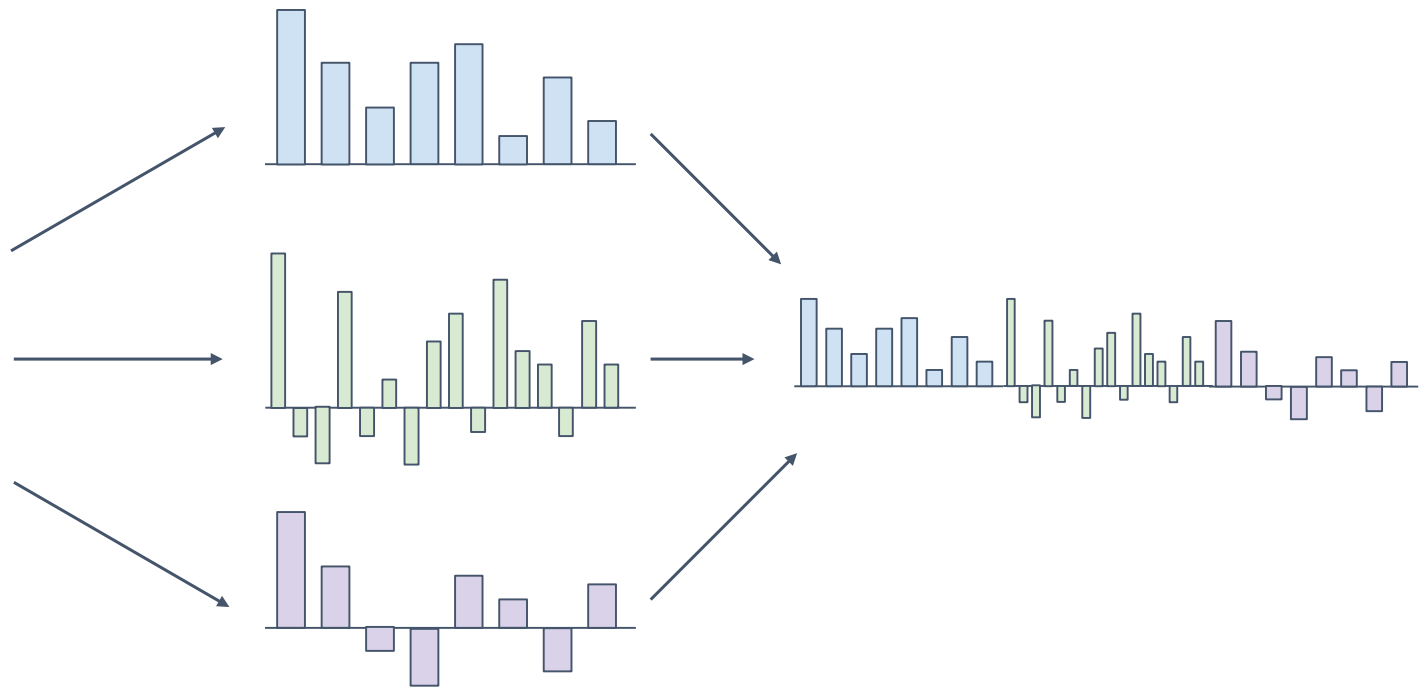
## Step 2: Encode images



Fei-Fei and Perona, “A bayesian hierarchical model for learning natural scene categories”, CVPR 2005

# Image Features

---



# Example: Winner of 2011 ImageNet challenge

---

Low-level feature extraction  $\approx$  10k patches per image

- SIFT: 128-dim
  - color: 96-dim
- } reduced to 64-dim with PCA

FV extraction and compression:

- $N=1,024$  Gaussians,  $R=4$  regions  $\Rightarrow$  520K dim x 2
- compression:  $G=8$ ,  $b=1$  bit per dimension

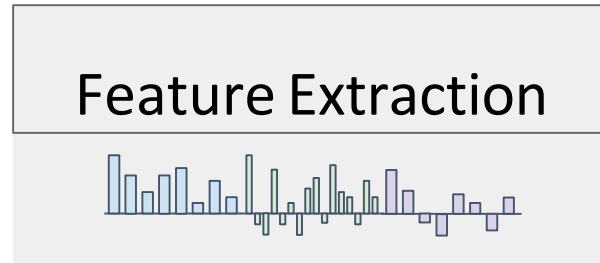
One-vs-all SVM learning with SGD

Late fusion of SIFT and color systems



# Image Features

---



$f$

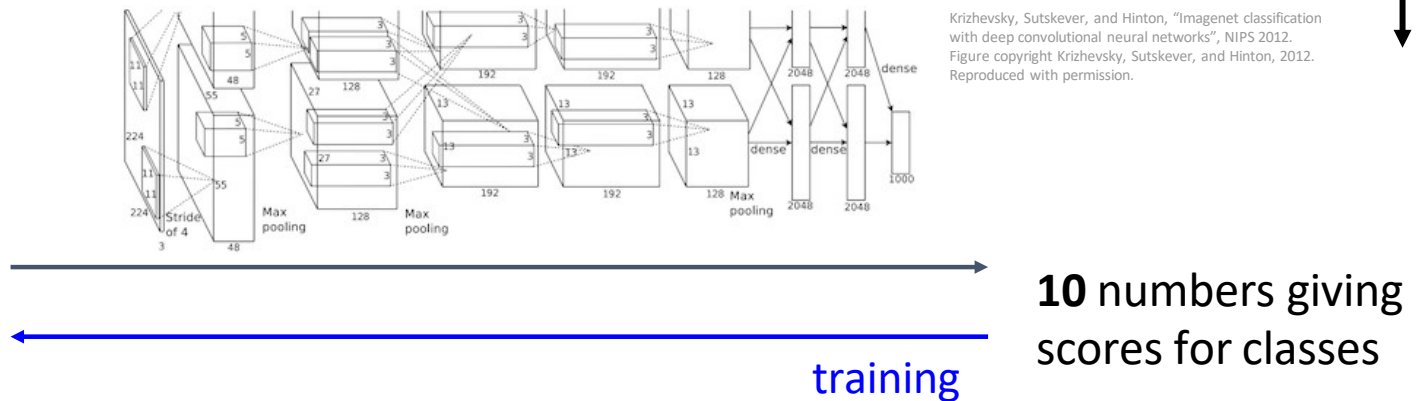
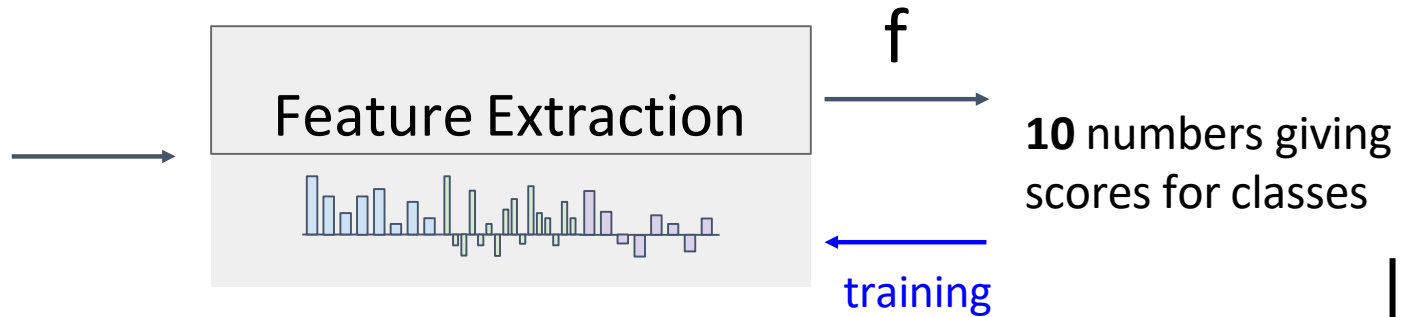


**10** numbers giving  
scores for classes



training

# Image Features



# Neural Networks

---

(**Before**) Linear score function:

$$f = Wx$$

$$x \in \mathbb{R}^D, W \in \mathbb{R}^{C \times D}$$

# Neural Networks

---

(**Before**) Linear score function:

$$f = Wx$$

(**Now**) 2-layer Neural Network

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

$$W_2 \in \mathbb{R}^{C \times H} \quad W_1 \in \mathbb{R}^{H \times D} \quad x \in \mathbb{R}^D$$

(In practice we will usually add a learnable bias at each layer as well)

# Neural Networks

(**Before**) Linear score function:

$$f = Wx$$

(**Now**) 2-layer Neural Network  
or 3-layer Neural Network

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

$$f = W_3 \max(0, W_2 \max(0, W_1 x))$$

$$W_3 \in \mathbb{R}^{C \times H_2} \quad W_2 \in \mathbb{R}^{H_2 \times H_1} \quad W_1 \in \mathbb{R}^{H_1 \times D} \quad x \in \mathbb{R}^D$$

(In practice we will usually add a learnable bias at each layer as well)

# Neural Networks

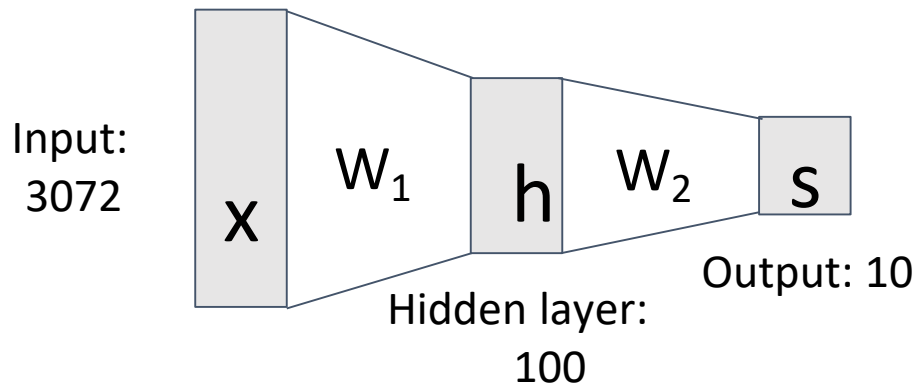
---

(**Before**) Linear score function:

$$f = Wx$$

(**Now**) 2-layer Neural Network

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



$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$

# Neural Networks

---

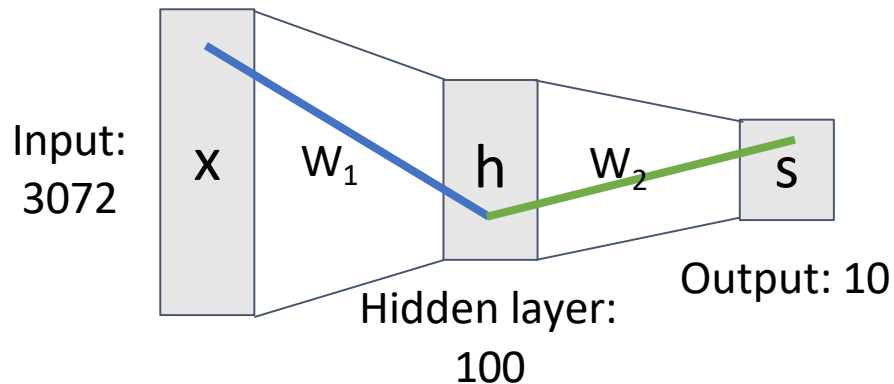
(**Before**) Linear score function:

$$f = Wx$$

(**Now**) 2-layer Neural Network

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

Element (i, j)  
of  $W_1$  gives  
the effect on  
 $h_i$  from  $x_j$



Element (i, j)  
of  $W_2$  gives  
the effect on  
 $s_i$  from  $h_j$

$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$

# Neural Networks

---

**(Before)** Linear score function:

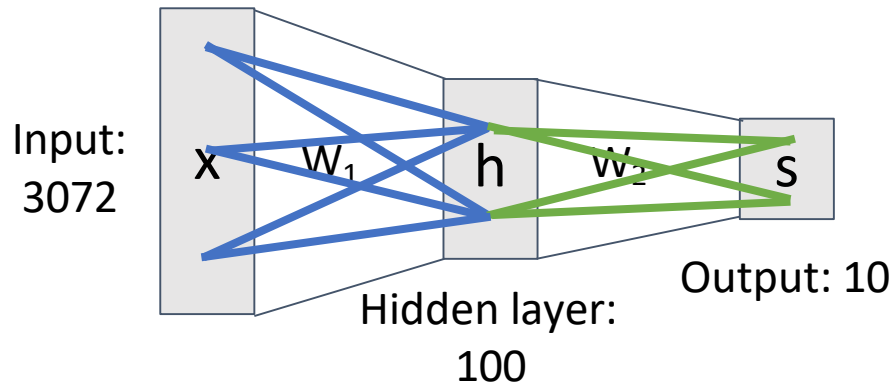
$$f = Wx$$

**(Now)** 2-layer Neural Network

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

Element  $(i, j)$  of  $W_1$   
gives the effect on  
 $h_i$  from  $x_j$

All elements  
of  $x$  affect all  
elements of  $h$



Element  $(i, j)$  of  $W_2$   
gives the effect on  
 $s_i$  from  $h_j$

All elements  
of  $h$  affect all  
elements of  $s$

Fully-connected neural network  
Also “Multi-Layer Perceptron” (MLP)



# Neural Networks

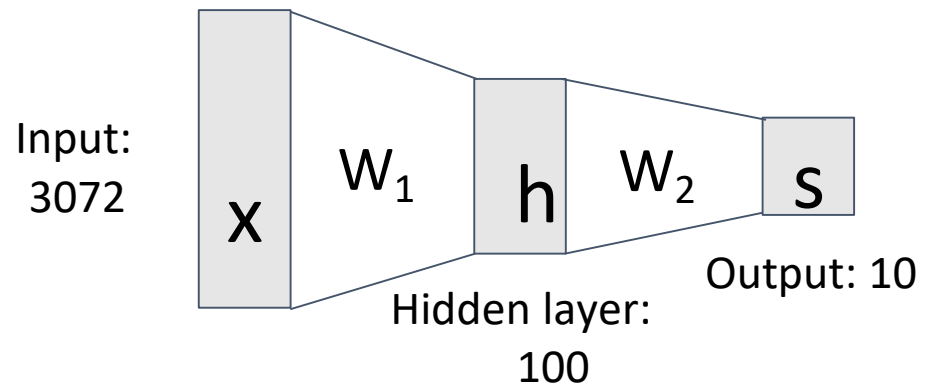
---

Linear classifier: One template per class



**(Before)** Linear score function:

**(Now)** 2-layer Neural Network



$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$

# Neural Networks

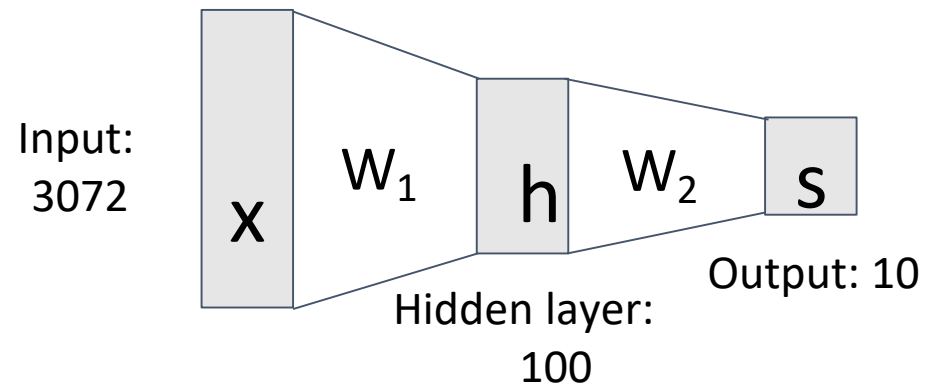
---

Neural net: first layer is bank of templates;  
Second layer recombines templates



**(Before)** Linear score function:

**(Now)** 2-layer Neural Network



$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$

# Neural Networks

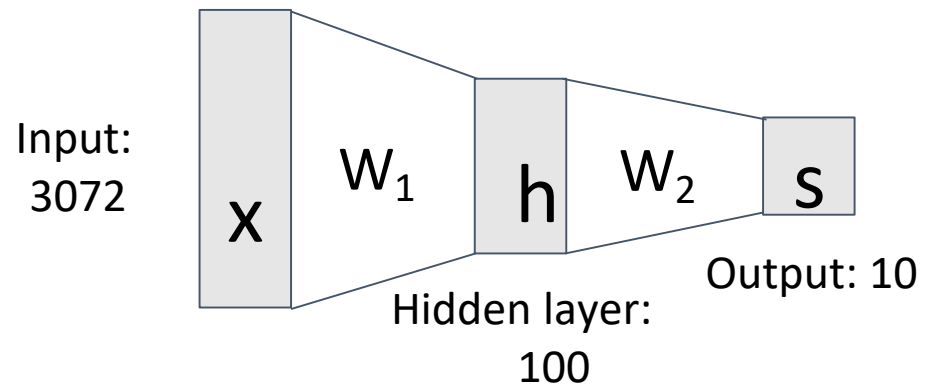
---

Can use different templates to cover multiple modes of a class!



**(Before)** Linear score function:

**(Now)** 2-layer Neural Network



$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$

# Neural Networks

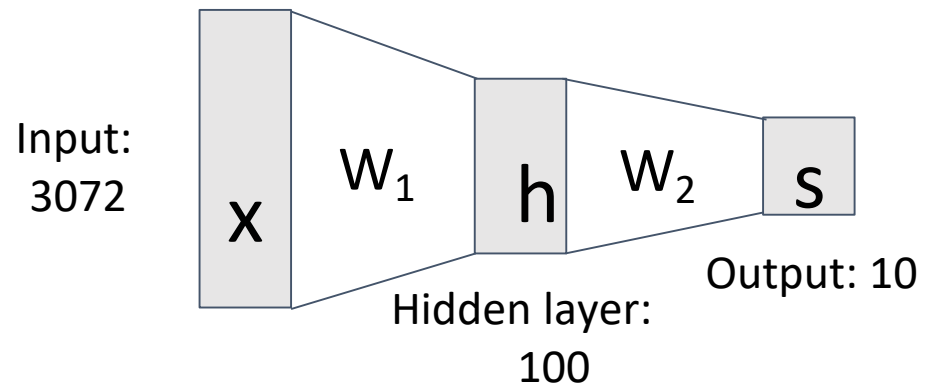
---

“Distributed representation”:  
Most templates not interpretable!



**(Before)** Linear score function:

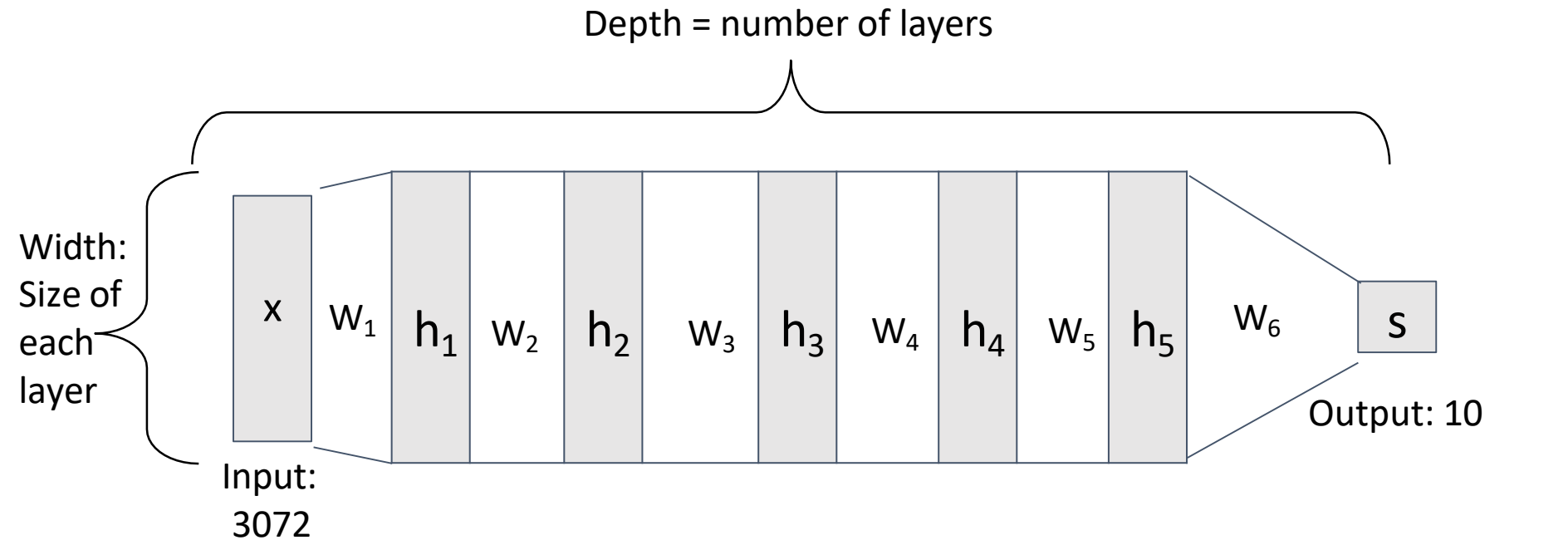
**(Now)** 2-layer Neural Network



$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$

# Deep Neural Networks

---



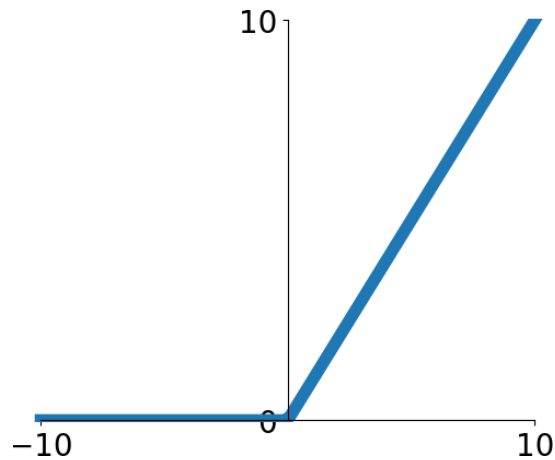
$$s = W_6 \max(0, W_6 \max(0, W_5 \max(0, W_4 \max(0, W_3 \max(0, W_2 \max(0, W_1 x))))))$$

# Activation Functions

---

## 2-layer Neural Network

The function  $ReLU(z) = \max(0, z)$  is called “Rectified Linear Unit”



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

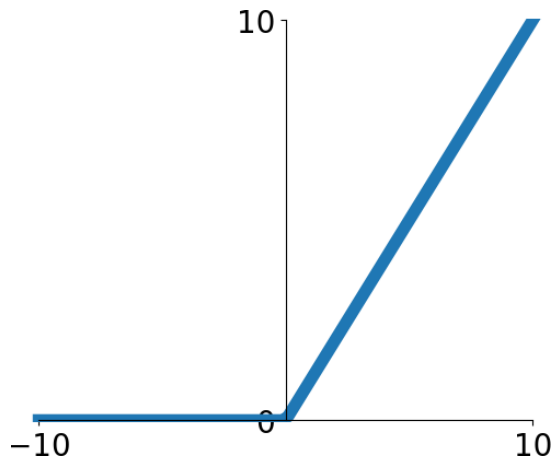
This is called the **activation function** of the neural network

# Activation Functions

---

## 2-layer Neural Network

The function  $ReLU(z) = \max(0, z)$  is called “Rectified Linear Unit”



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

This is called the **activation function** of the neural network

**Q:** What happens if we build a neural network with no activation function?

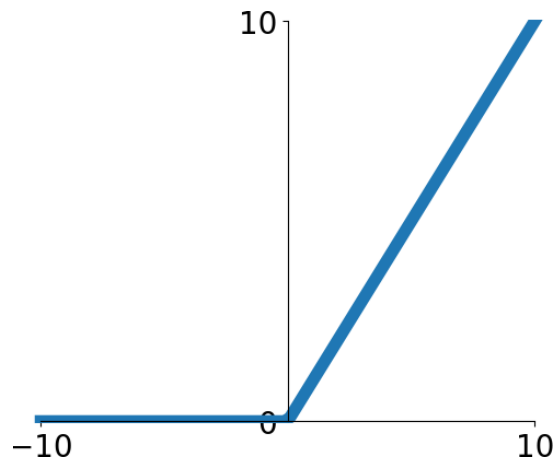
$$s = W_2 W_1 x$$

# Activation Functions

---

## 2-layer Neural Network

The function  $ReLU(z) = \max(0, z)$  is called “Rectified Linear Unit”



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

This is called the **activation function** of the neural network

**Q:** What happens if we build a neural network with no activation function?

$$s = W_2 W_1 x$$

$$W_3 = W_2 W_1 \in \mathbb{R}^{C \times H} \quad s = W_3 x$$

**A:** We end up with a linear classifier!

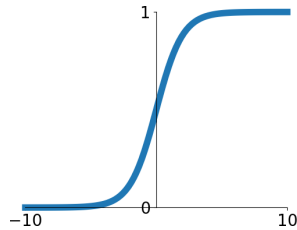


# Activation Functions

---

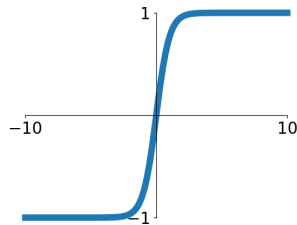
## Sigmoid

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



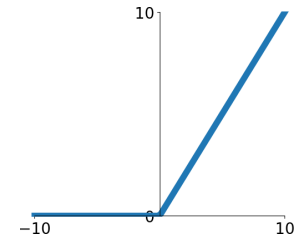
## tanh

$$\tanh(x)$$



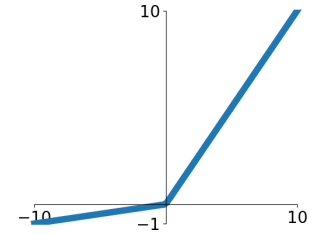
## ReLU

$$\max(0, x)$$



## Leaky ReLU

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

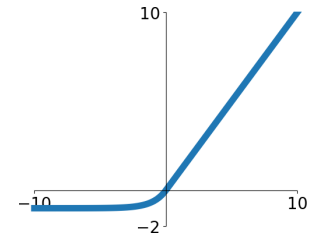


## Maxout

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

## ELU

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



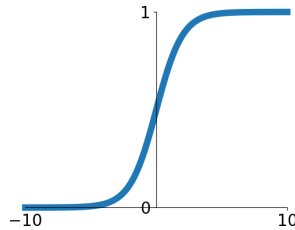
# Activation Functions

---

ReLU is a good default choice  
for most problems

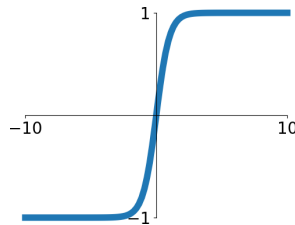
## Sigmoid

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



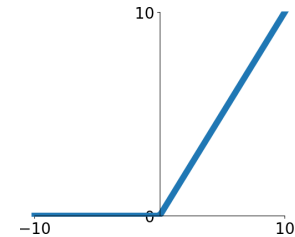
## tanh

$$\tanh(x)$$



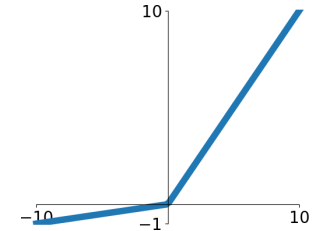
## ReLU

$$\max(0, x)$$



## Leaky ReLU

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

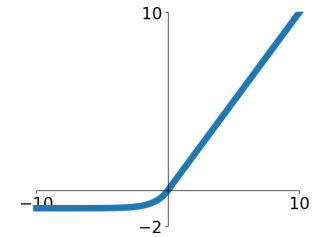


## Maxout

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

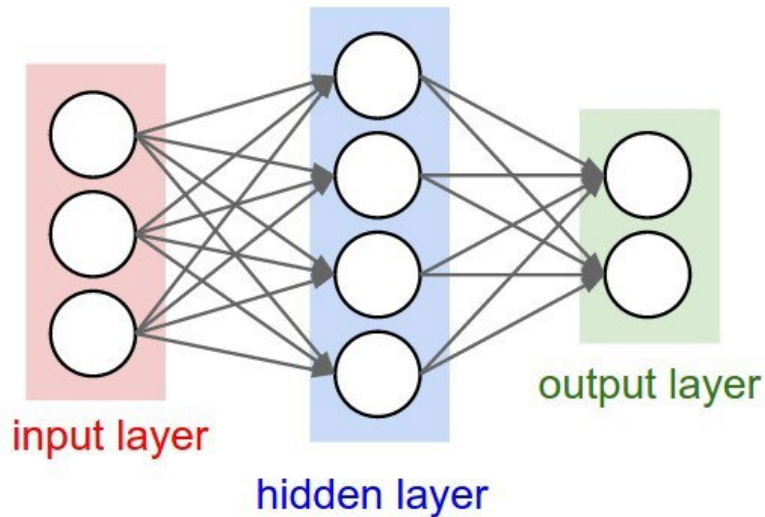
## ELU

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



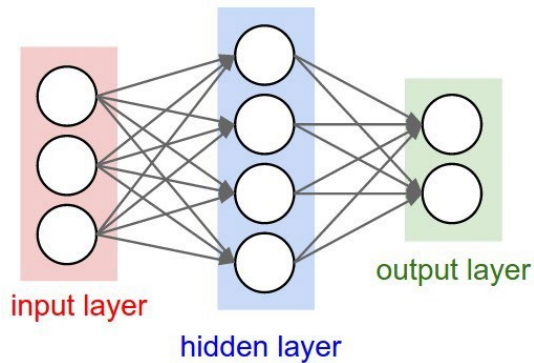
# Neural Net in <20 lines!

---



```
1  import numpy as np
2  from numpy.random import randn
3
4  N, Din, H, Dout = 64, 1000, 100, 10
5  x, y = randn(N, Din), randn(N, Dout)
6  w1, w2 = randn(Din, H), randn(H, Dout)
7  for t in range(10000):
8      h = 1.0 / (1.0 + np.exp(-x.dot(w1)))
9      y_pred = h.dot(w2)
10     loss = np.square(y_pred - y).sum()
11     dy_pred = 2.0 * (y_pred - y)
12     dw2 = h.T.dot(dy_pred)
13     dh = dy_pred.dot(w2.T)
14     dw1 = x.T.dot(dh * h * (1 - h))
15     w1 -= 1e-4 * dw1
16     w2 -= 1e-4 * dw2
```

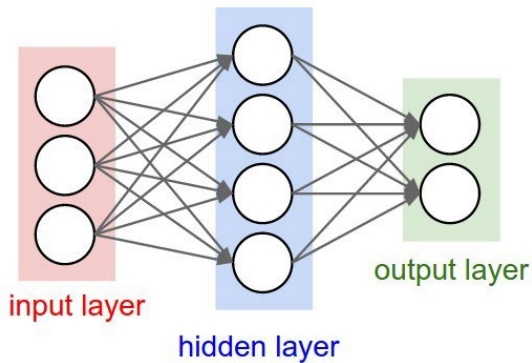
# Neural Net in <20 lines!



Initialize weights  
and data

```
1  import numpy as np
2  from numpy.random import randn
3
4  N, Din, H, Dout = 64, 1000, 100, 10
5  x, y = randn(N, Din), randn(N, Dout)
6  w1, w2 = randn(Din, H), randn(H, Dout)
7  for t in range(10000):
8      h = 1.0 / (1.0 + np.exp(-x.dot(w1)))
9      y_pred = h.dot(w2)
10     loss = np.square(y_pred - y).sum()
11     dy_pred = 2.0 * (y_pred - y)
12     dw2 = h.T.dot(dy_pred)
13     dh = dy_pred.dot(w2.T)
14     dw1 = x.T.dot(dh * h * (1 - h))
15     w1 -= 1e-4 * dw1
16     w2 -= 1e-4 * dw2
```

# Neural Net in <20 lines!

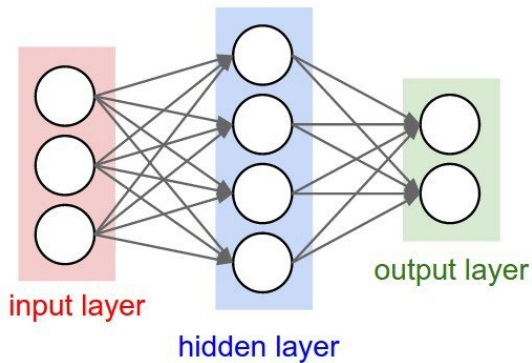


Initialize weights  
and data

Compute loss  
(sigmoid activation,  
L2 loss)

```
1  import numpy as np
2  from numpy.random import randn
3
4  N, Din, H, Dout = 64, 1000, 100, 10
5  x, y = randn(N, Din), randn(N, Dout)
6  w1, w2 = randn(Din, H), randn(H, Dout)
7  for t in range(10000):
8      h = 1.0 / (1.0 + np.exp(-x.dot(w1)))
9      y_pred = h.dot(w2)
10     loss = np.square(y_pred - y).sum()
11     dy_pred = 2.0 * (y_pred - y)
12     dw2 = h.T.dot(dy_pred)
13     dh = dy_pred.dot(w2.T)
14     dw1 = x.T.dot(dh * h * (1 - h))
15     w1 -= 1e-4 * dw1
16     w2 -= 1e-4 * dw2
```

# Neural Net in <20 lines!



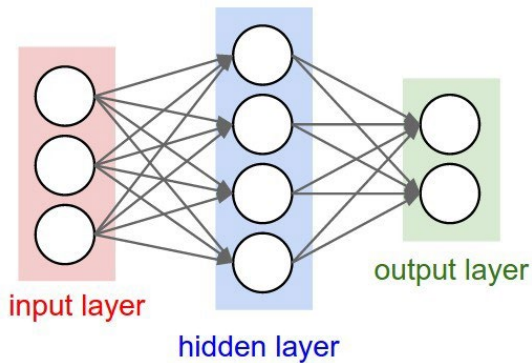
Initialize weights  
and data

Compute loss  
(sigmoid activation,  
L2 loss)

Compute  
gradients

```
1  import numpy as np
2  from numpy.random import randn
3
4  N, Din, H, Dout = 64, 1000, 100, 10
5  x, y = randn(N, Din), randn(N, Dout)
6  w1, w2 = randn(Din, H), randn(H, Dout)
7  for t in range(10000):
8      h = 1.0 / (1.0 + np.exp(-x.dot(w1)))
9      y_pred = h.dot(w2)
10     loss = np.square(y_pred - y).sum()
11     dy_pred = 2.0 * (y_pred - y)
12     dw2 = h.T.dot(dy_pred)
13     dh = dy_pred.dot(w2.T)
14     dw1 = x.T.dot(dh * h * (1 - h))
15     w1 -= 1e-4 * dw1
16     w2 -= 1e-4 * dw2
```

# Neural Net in <20 lines!



Initialize weights  
and data

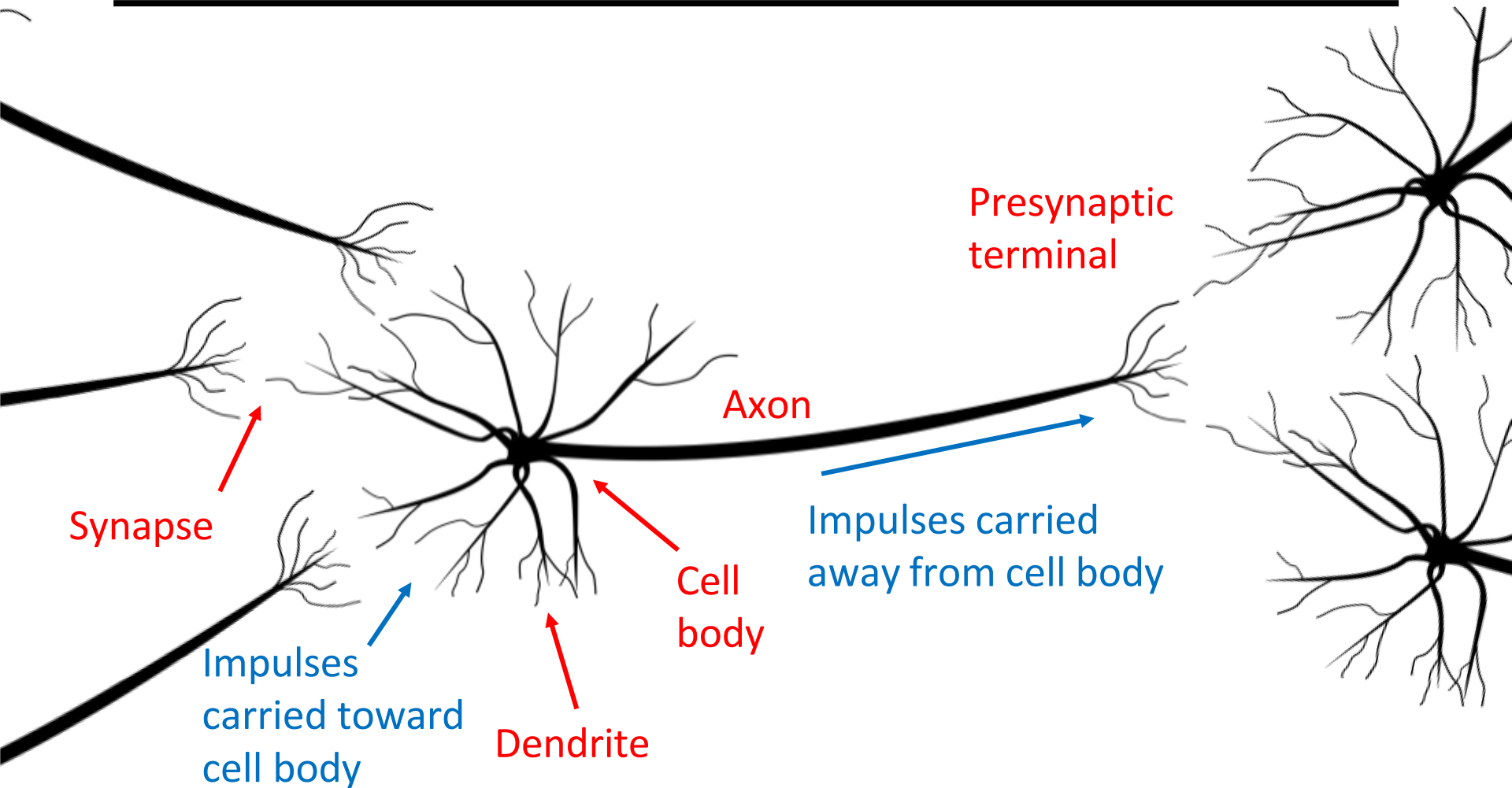
Compute loss  
(sigmoid activation,  
L2 loss)

Compute  
gradients

SGD  
step

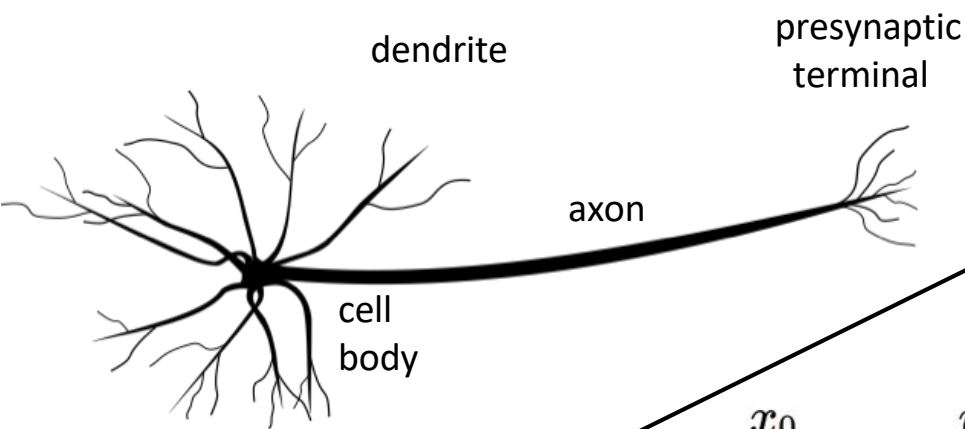
```
1  import numpy as np
2  from numpy.random import randn
3
4  N, Din, H, Dout = 64, 1000, 100, 10
5  x, y = randn(N, Din), randn(N, Dout)
6  w1, w2 = randn(Din, H), randn(H, Dout)
7  for t in range(10000):
8      h = 1.0 / (1.0 + np.exp(-x.dot(w1)))
9      y_pred = h.dot(w2)
10     loss = np.square(y_pred - y).sum()
11     dy_pred = 2.0 * (y_pred - y)
12     dw2 = h.T.dot(dy_pred)
13     dh = dy_pred.dot(w2.T)
14     dw1 = x.T.dot(dh * h * (1 - h))
15     w1 -= 1e-4 * dw1
16     w2 -= 1e-4 * dw2
```

# Our brains are made of Neurons

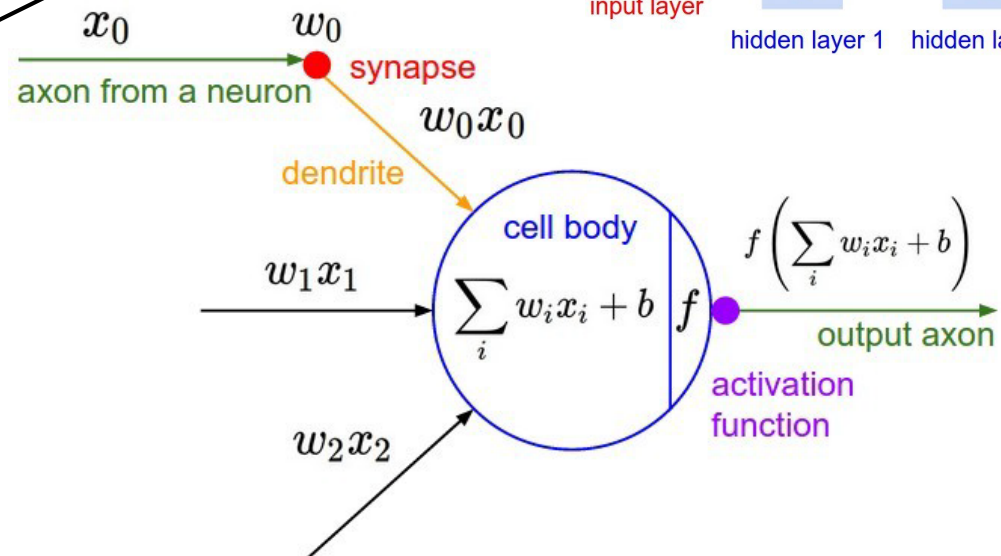
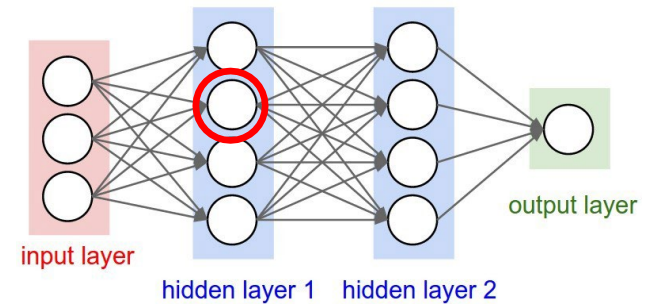




# Biological Neuron

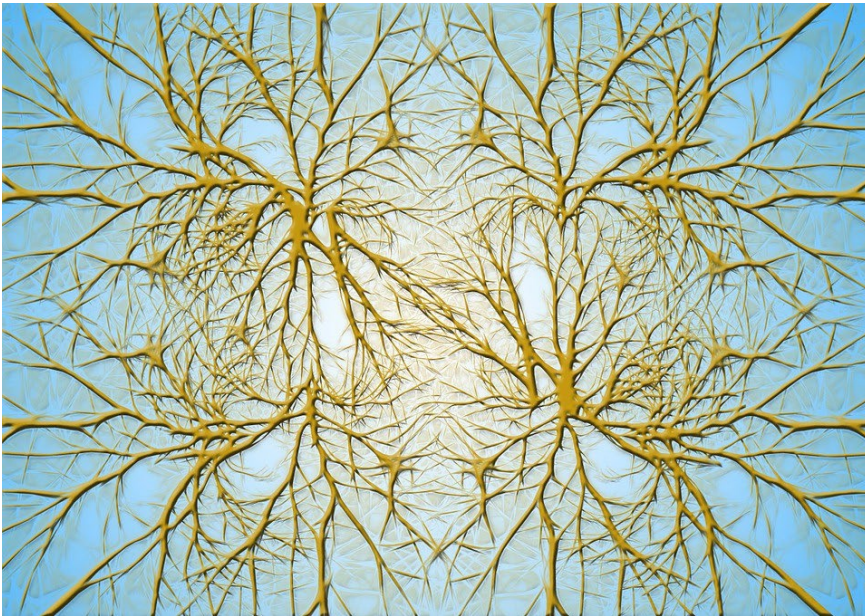


# Artificial Neuron



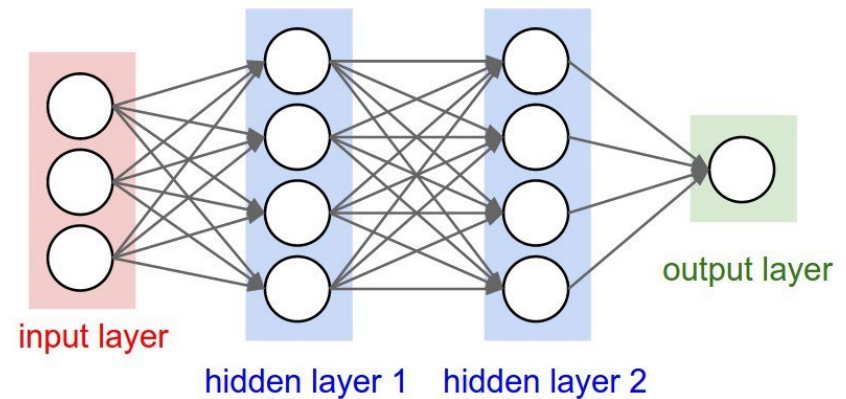
---

## Biological Neurons: Complex connectivity patterns



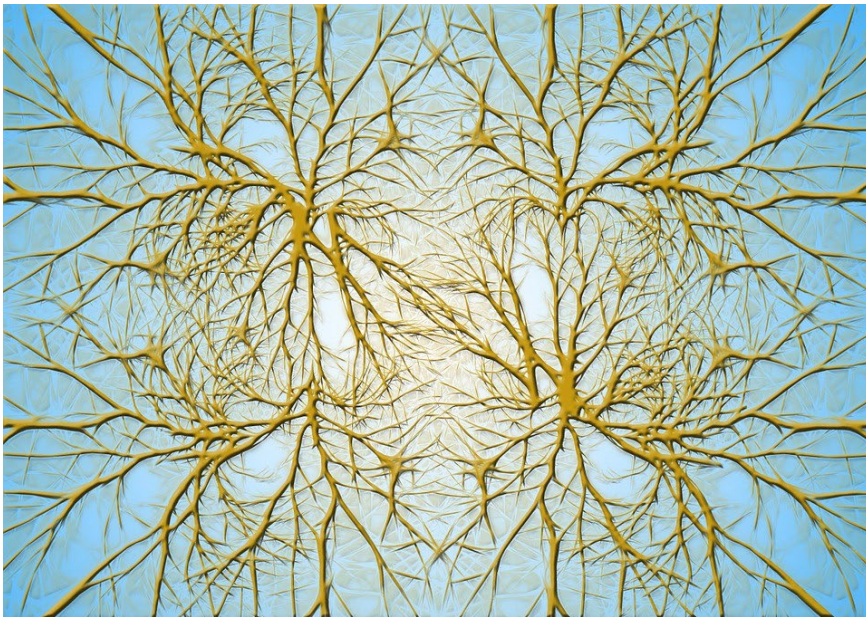
[This image is CC0 Public Domain](#)

## Neurons in a neural network: Organized into regular layers for computational efficiency



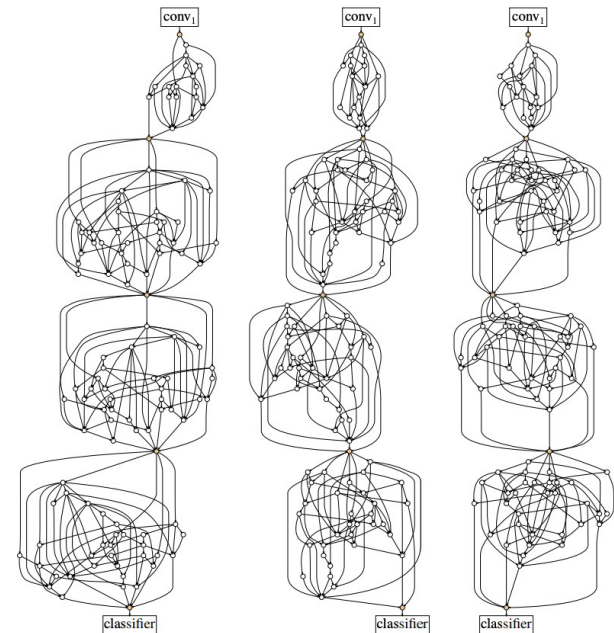
---

## Biological Neurons: Complex connectivity patterns



[This image is CC0 Public Domain](#)

But neural networks with random connections can work too!



Xie et al, "Exploring Randomly Wired Neural Networks for Image Recognition", ICCV 2019

# Be very careful with brain analogies!

---

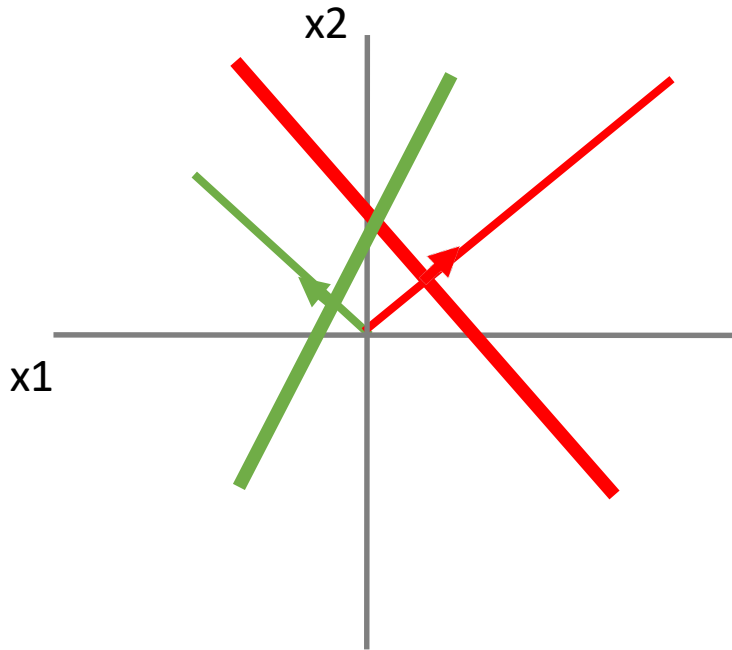
## **Biological Neurons:**

- Many different types
- Dendrites can perform complex non-linear computations
- Synapses are not a single weight but a complex non-linear dynamical system
- Rate code may not be adequate

# Space Warping

---

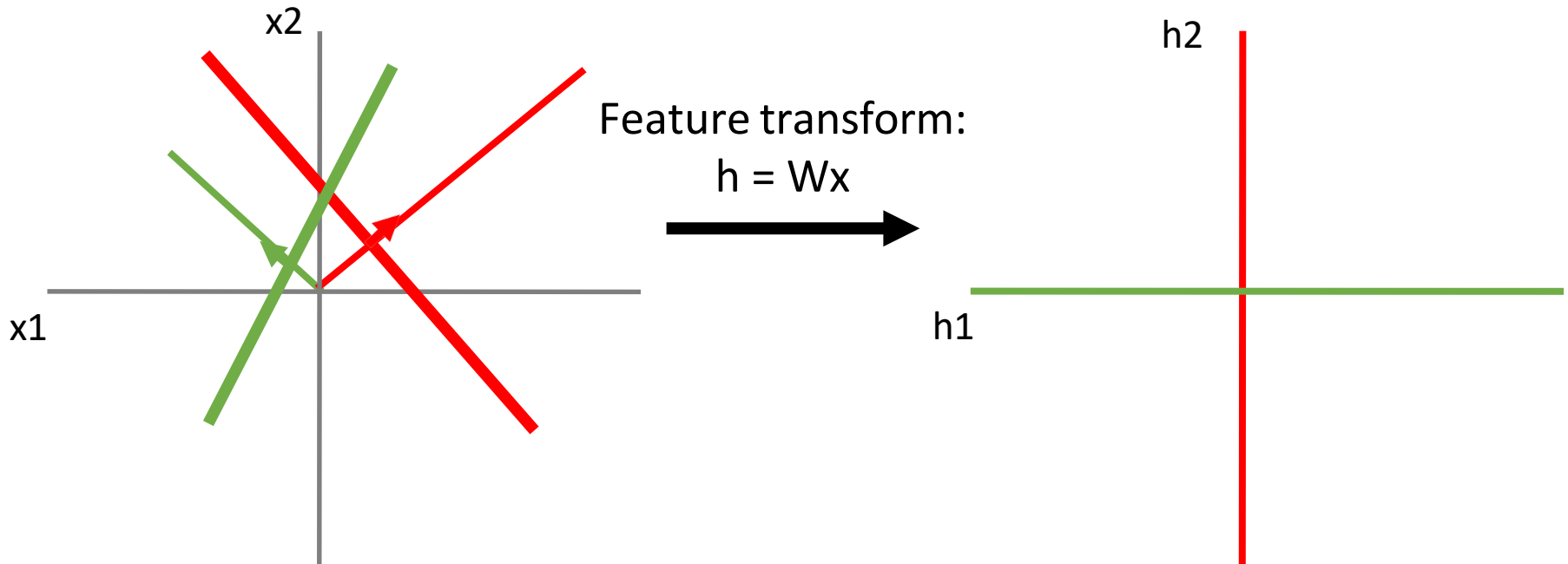
Consider a linear transform:  $h = Wx$   
Where  $x, h$  are both 2-dimensional



# Space Warping

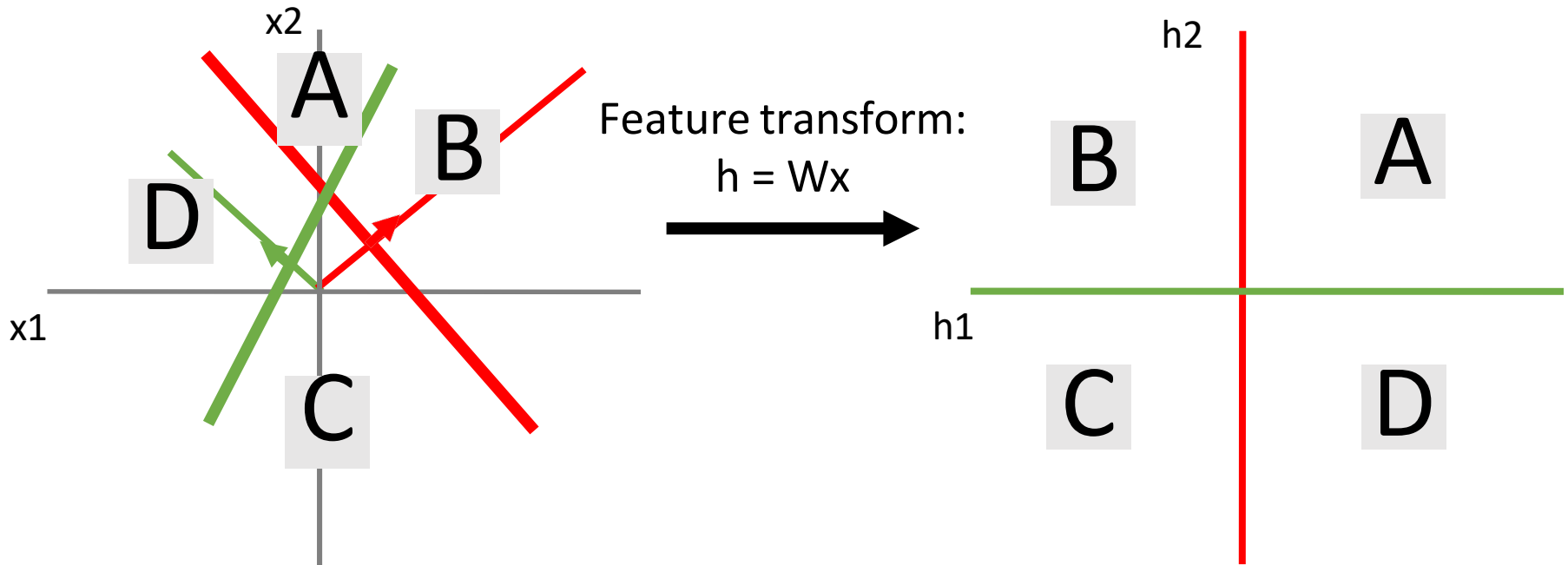
---

Consider a linear transform:  $h = Wx$   
Where  $x, h$  are both 2-dimensional



# Space Warping

Consider a linear transform:  $h = Wx$   
Where  $x, h$  are both 2-dimensional

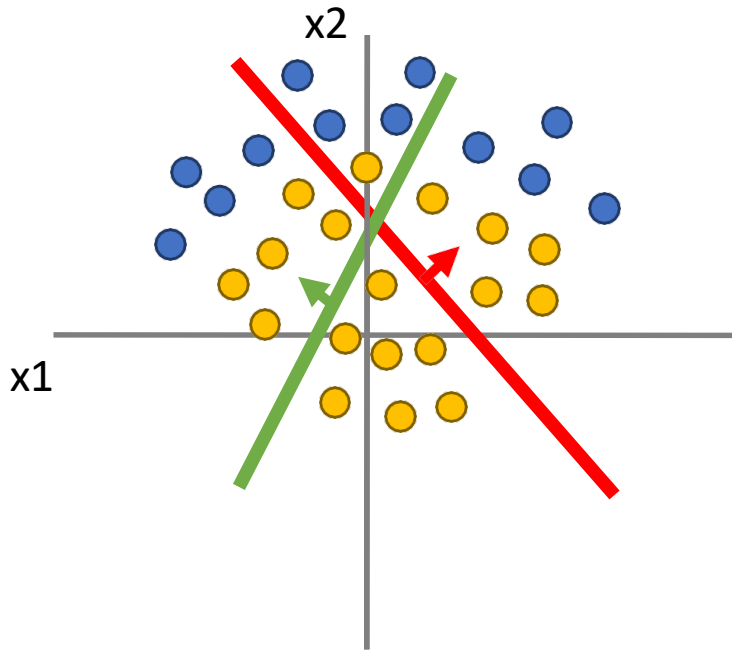




# Space Warping

---

Points not linearly separable in original space

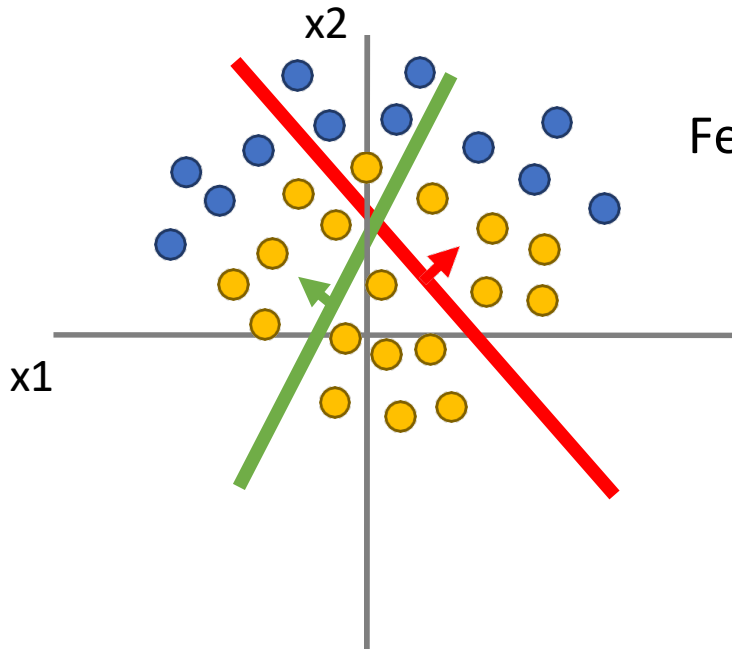


Consider a linear transform:  $h = Wx$   
Where  $x$ ,  $h$  are both 2-dimensional



# Space Warping

Points not linearly separable in original space

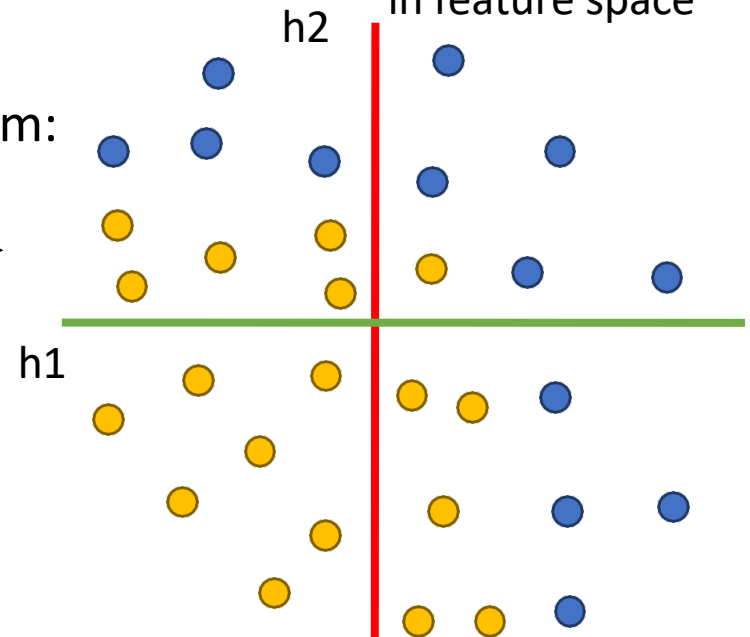


Feature transform:  
 $h = Wx$



Consider a linear transform:  $h = Wx$   
Where  $x, h$  are both 2-dimensional

Not linearly separable  
in feature space



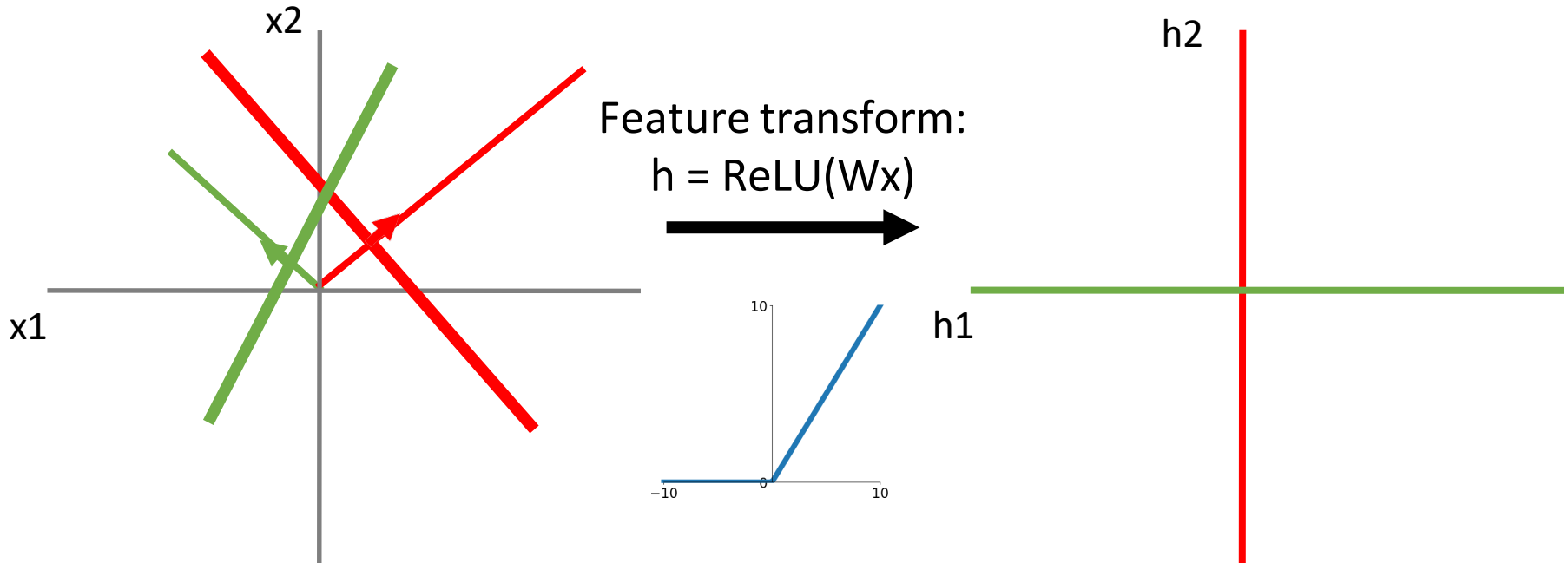
# Space Warping

---

Consider a neural net hidden layer:

$$h = \text{ReLU}(Wx) = \max(0, Wx)$$

Where  $x$ ,  $h$  are both 2-dimensional

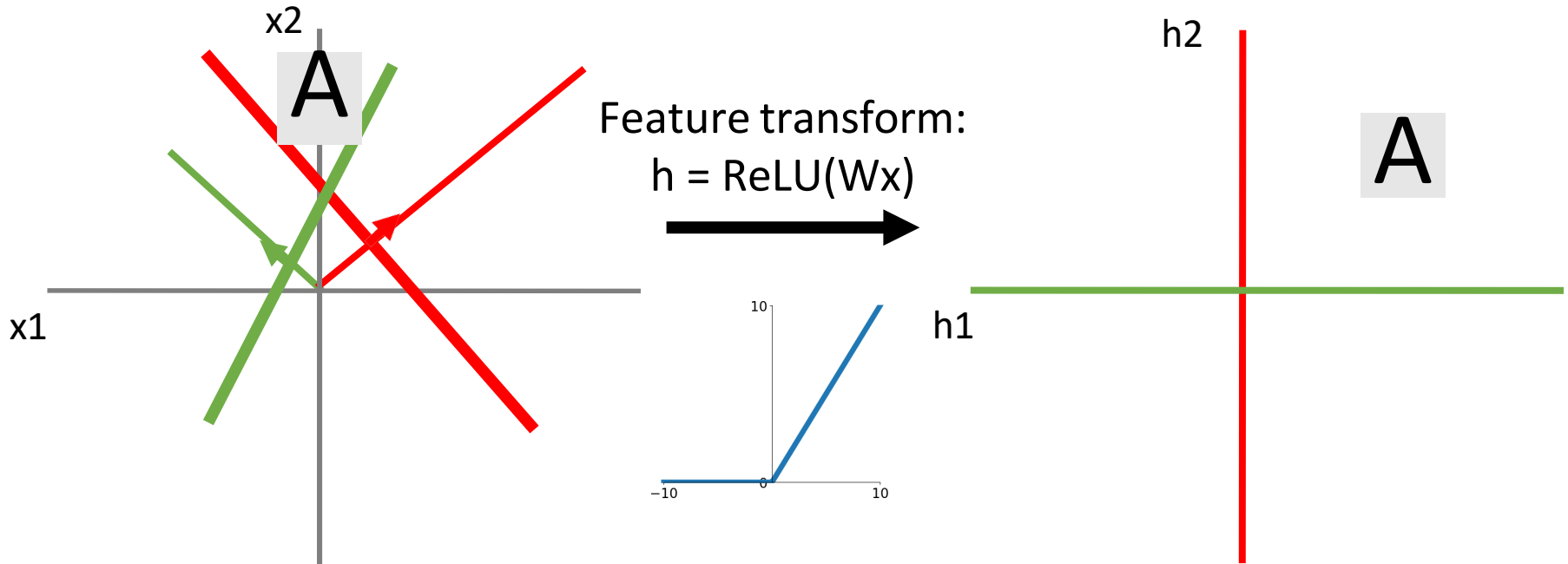


# Space Warping

Consider a neural net hidden layer:

$$h = \text{ReLU}(Wx) = \max(0, Wx)$$

Where  $x$ ,  $h$  are both 2-dimensional

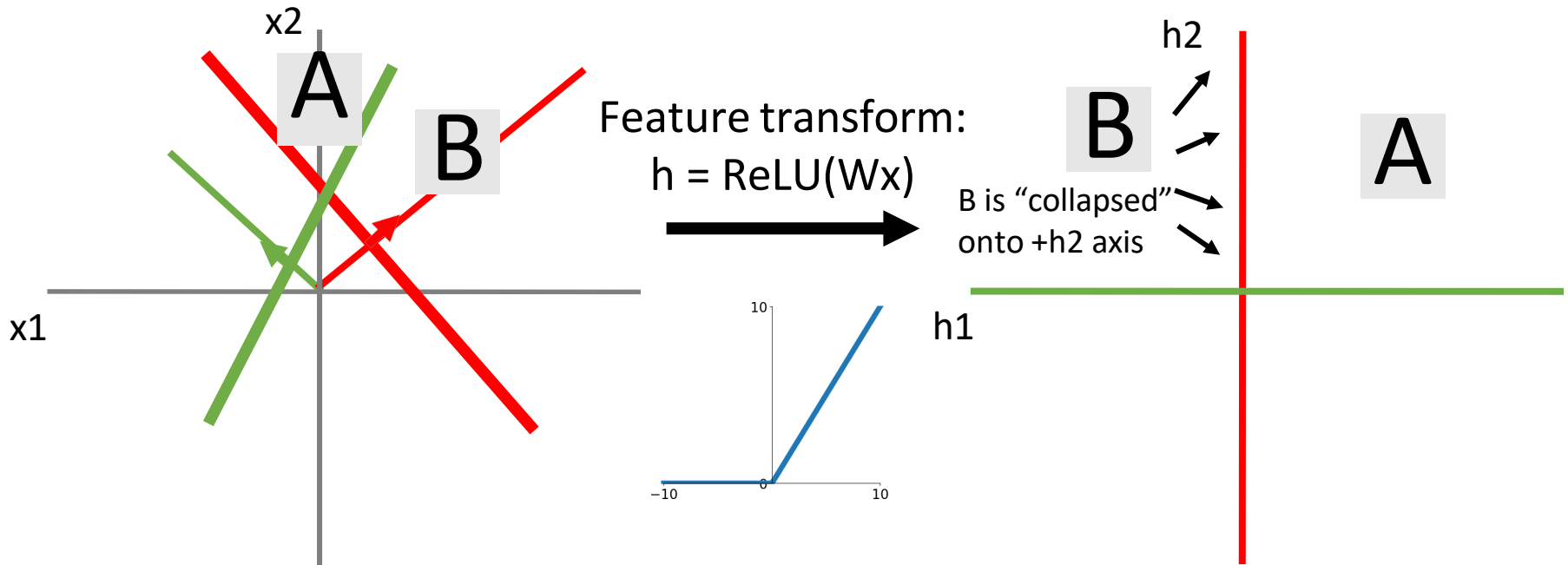


# Space Warping

Consider a neural net hidden layer:

$$h = \text{ReLU}(Wx) = \max(0, Wx)$$

Where  $x$ ,  $h$  are both 2-dimensional

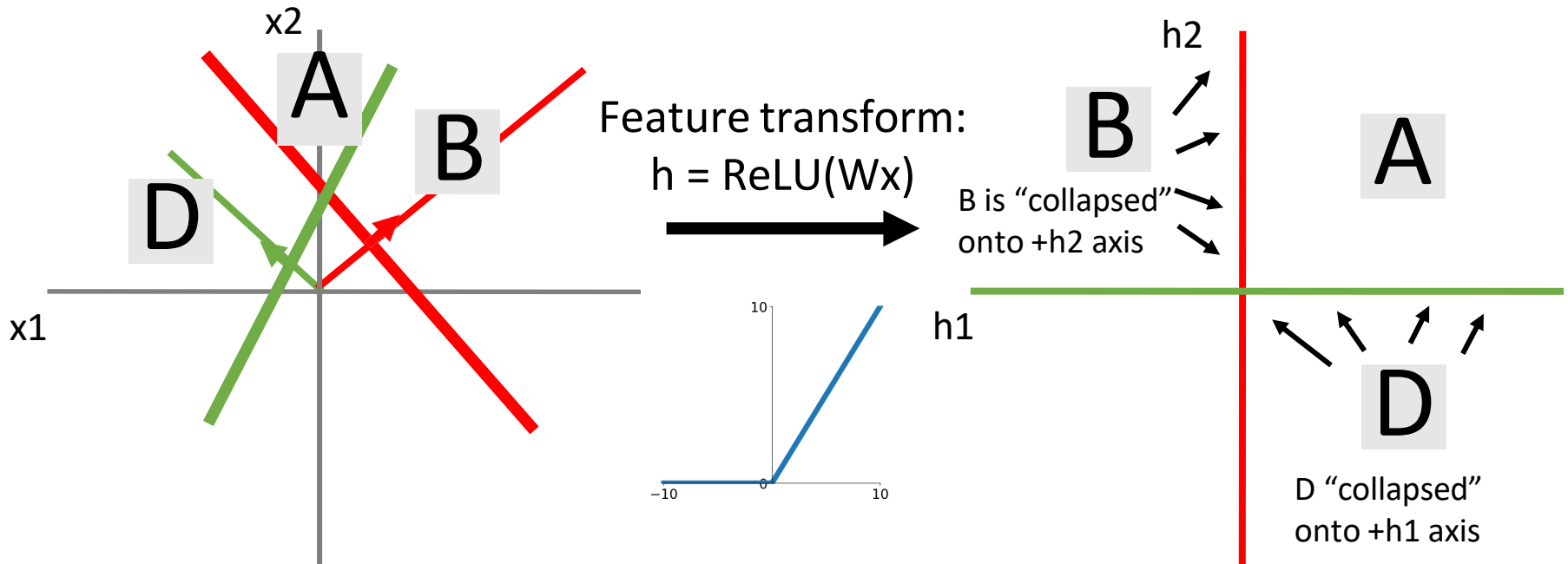


# Space Warping

Consider a neural net hidden layer:

$$h = \text{ReLU}(Wx) = \max(0, Wx)$$

Where  $x$ ,  $h$  are both 2-dimensional

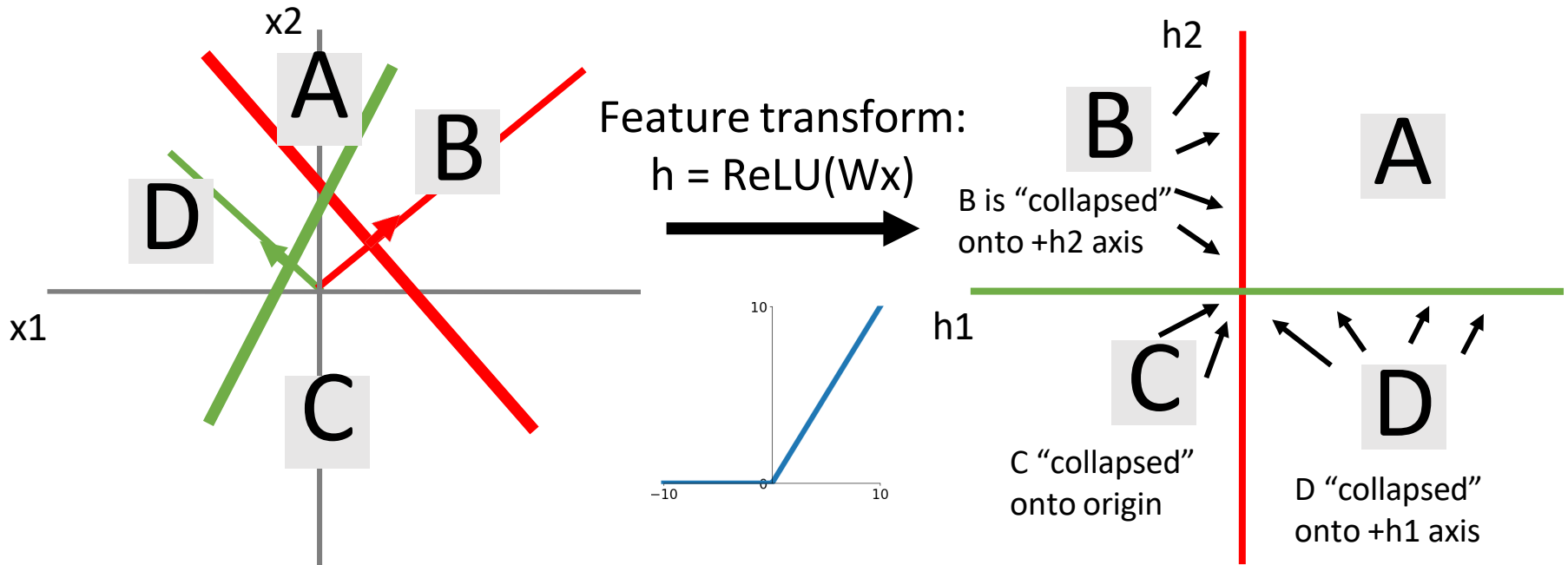


# Space Warping

Consider a neural net hidden layer:

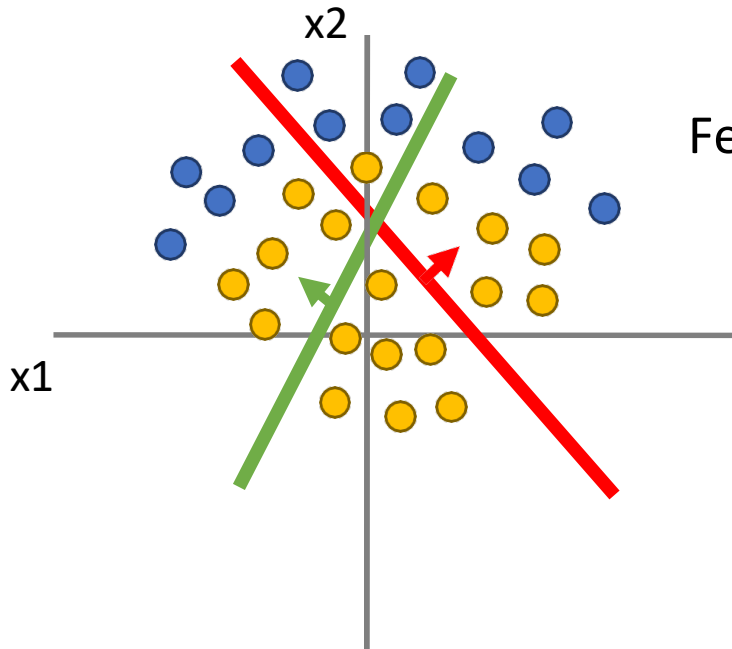
$$h = \text{ReLU}(Wx) = \max(0, Wx)$$

Where  $x$ ,  $h$  are both 2-dimensional



# Space Warping

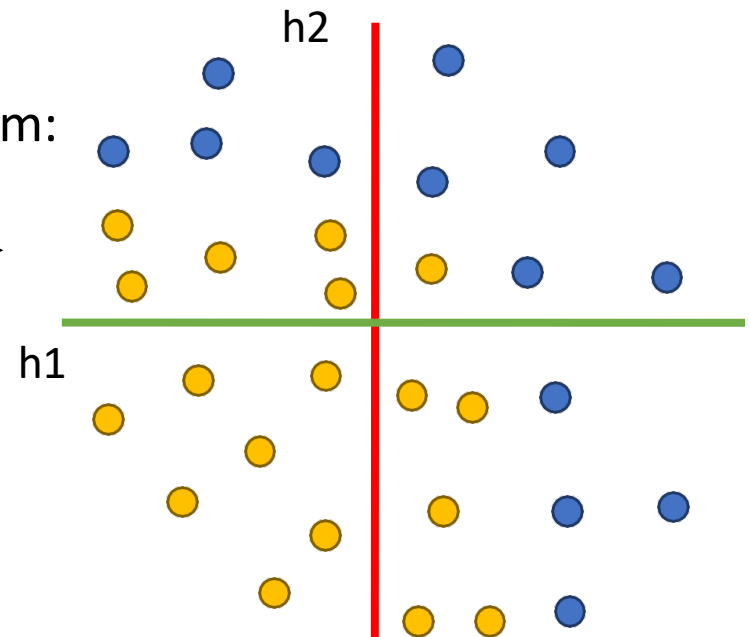
Points not linearly separable in original space



Feature transform:  
 $h = Wx$

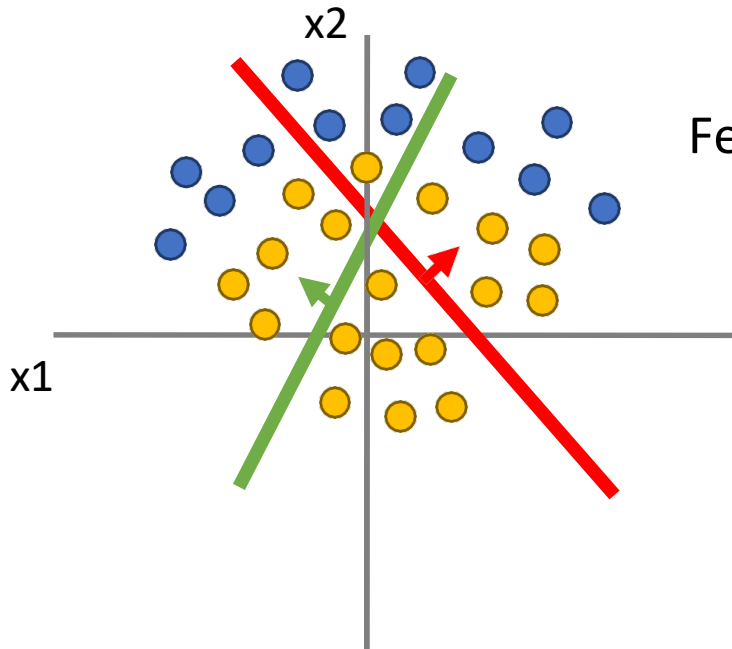


Consider a neural net hidden layer:  
 $h = \text{ReLU}(Wx) = \max(0, Wx)$   
Where  $x$ ,  $h$  are both 2-dimensional

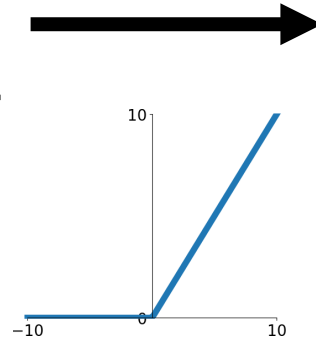


# Space Warping

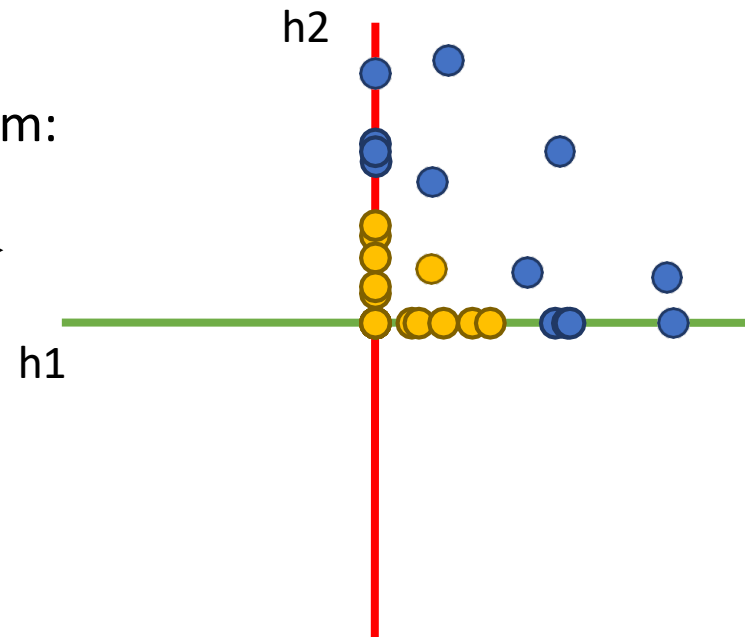
Points not linearly separable in original space



Feature transform:  
 $h = \text{ReLU}(Wx)$



Consider a neural net hidden layer:  
 $h = \text{ReLU}(Wx) = \max(0, Wx)$   
Where  $x$ ,  $h$  are both 2-dimensional



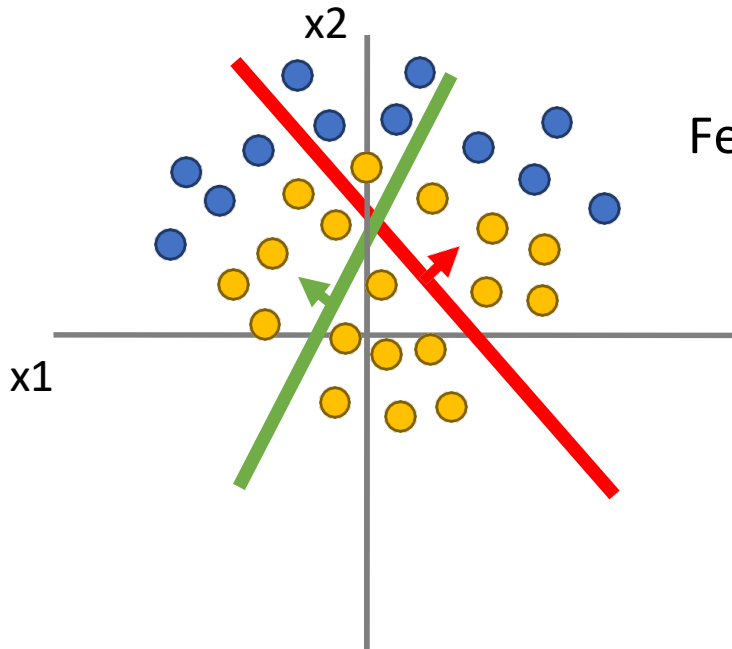


# Space Warping

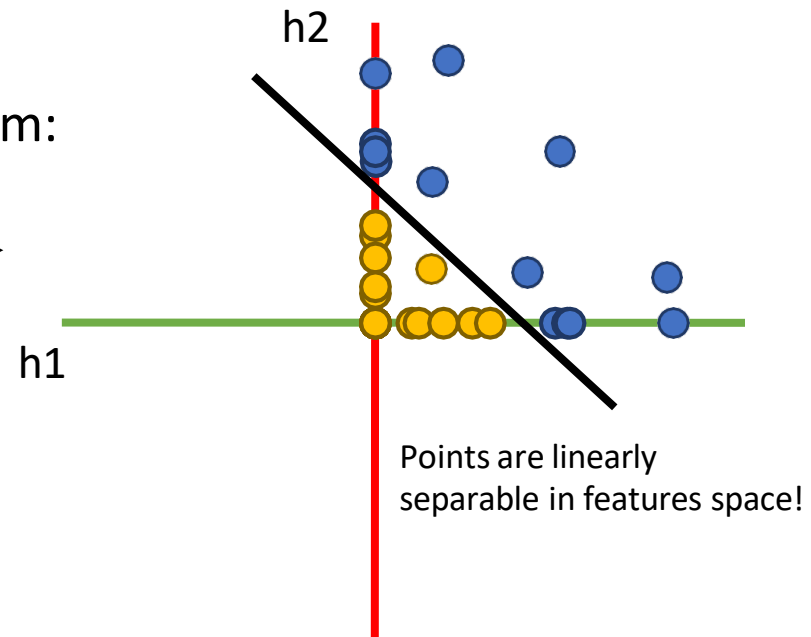
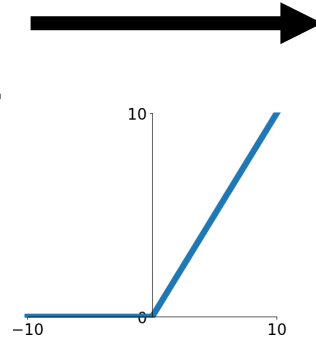
Consider a neural net hidden layer:  
 $h = \text{ReLU}(Wx) = \max(0, Wx)$

Where  $x$ ,  $h$  are both 2-dimensional

Points not linearly  
separable in original space



Feature transform:  
 $h = \text{ReLU}(Wx)$

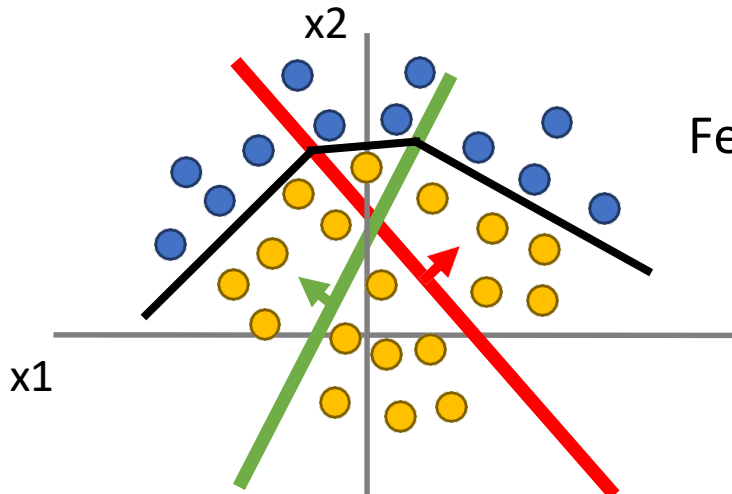


# Space Warping

Consider a neural net hidden layer:  
 $h = \text{ReLU}(Wx) = \max(0, Wx)$

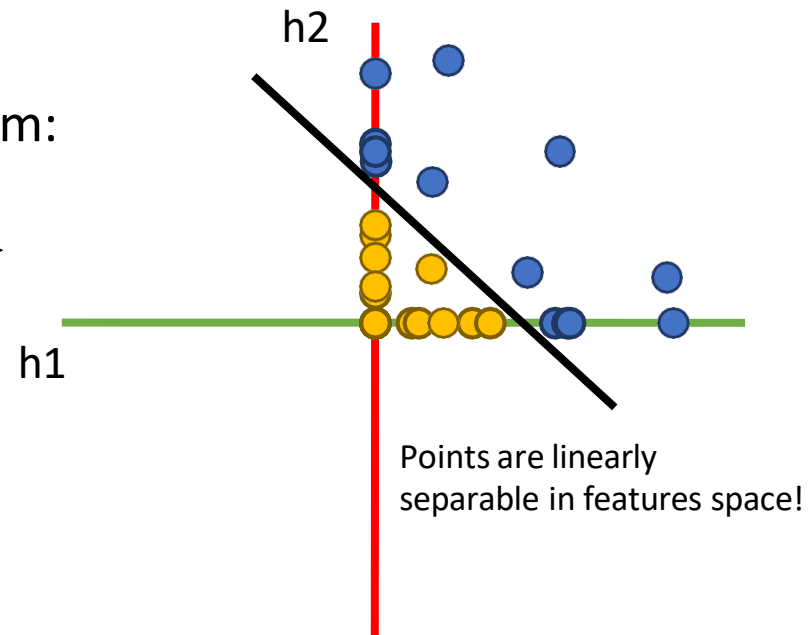
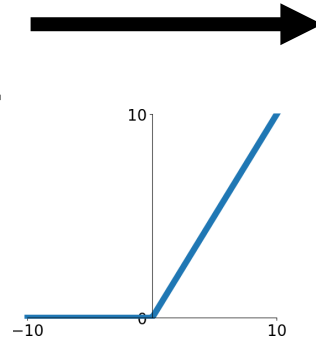
Where  $x$ ,  $h$  are both 2-dimensional

Points not linearly  
separable in original space



Linear classifier in feature  
space gives nonlinear  
classifier in original space

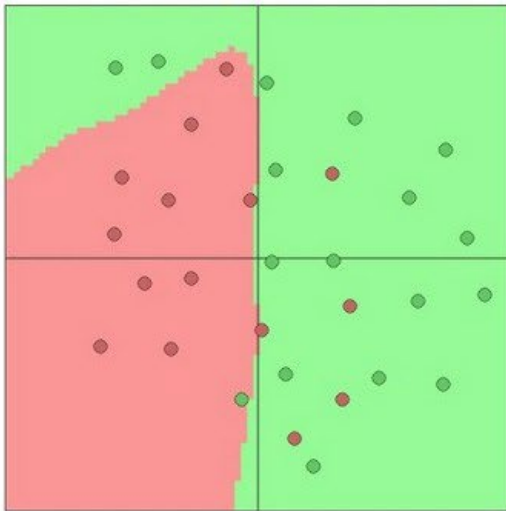
Feature transform:  
 $h = \text{ReLU}(Wx)$



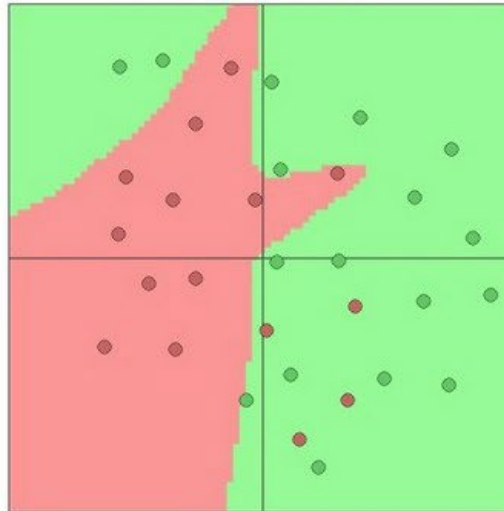
Points are linearly  
separable in features space!

# Setting the number of layers and their sizes

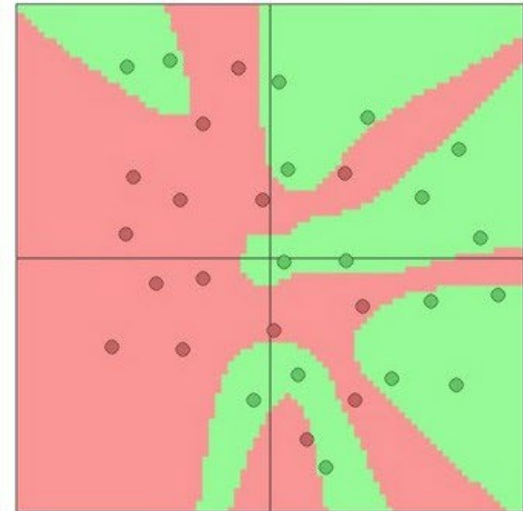
3 hidden units



6 hidden units



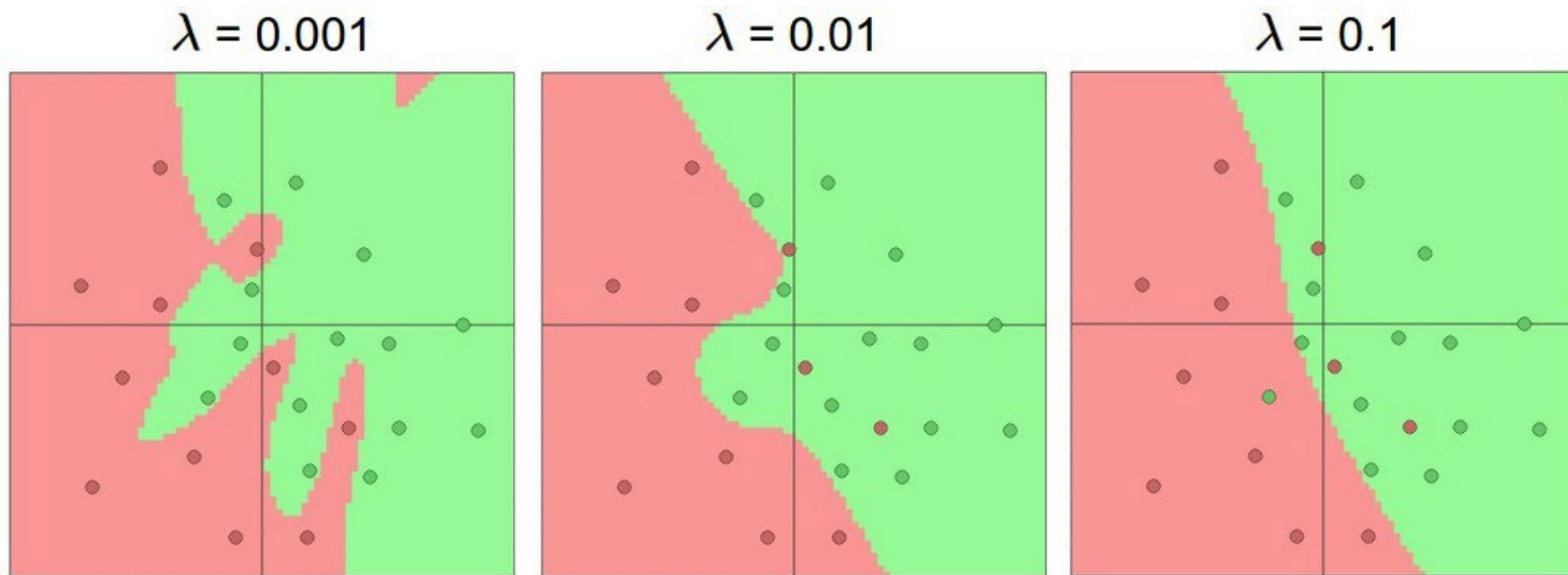
20 hidden units



More hidden units = more capacity

# Don't regularize with size; instead use stronger L2

---



(Web demo with ConvNetJS:

<http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>)

# Universal Approximation

---

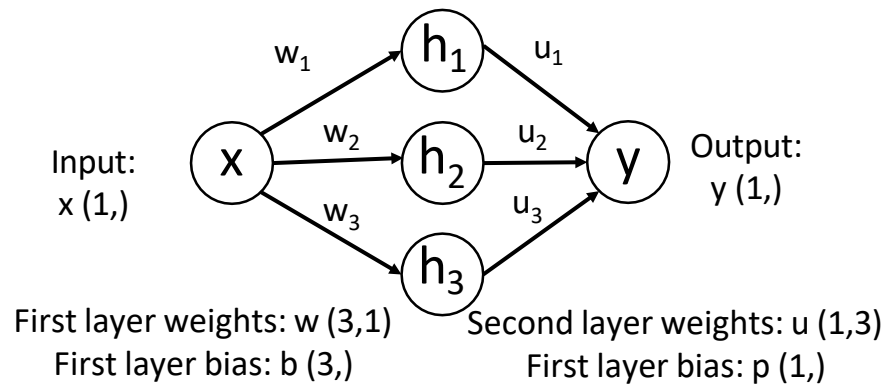
A neural network with one hidden layer can approximate any function  $f: \mathbb{R}^N \rightarrow \mathbb{R}^M$  with arbitrary precision\*

\*Many technical conditions: Only holds on compact subsets of  $\mathbb{R}^N$ ; function must be continuous; need to define “arbitrary precision”; etc

# Universal Approximation

---

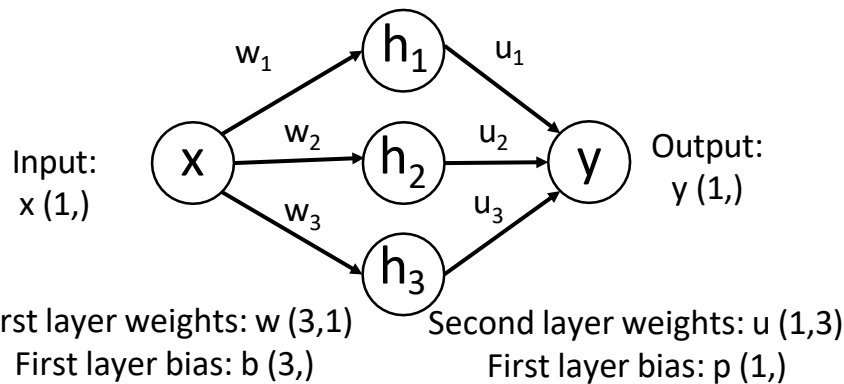
Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



# Universal Approximation

---

Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



$$h_1 = \max(0, w_1 * x + b_1)$$

$$h_2 = \max(0, w_2 * x + b_2)$$

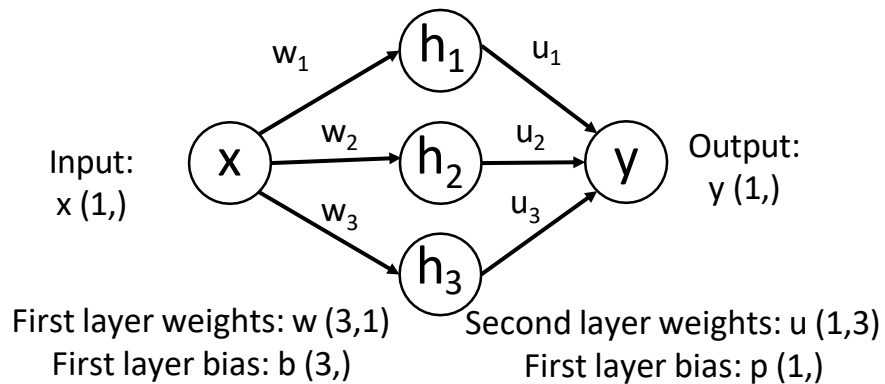
$$h_3 = \max(0, w_3 * x + b_3)$$

$$y = u_1 * h_1 + u_2 * h_2 + u_3 * h_3 + p$$

# Universal Approximation

---

Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



$$h_1 = \max(0, w_1 * x + b_1)$$

$$h_2 = \max(0, w_2 * x + b_2)$$

$$h_3 = \max(0, w_3 * x + b_3)$$

$$y = u_1 * h_1 + u_2 * h_2 + u_3 * h_3 + p$$

$$y = u_1 * \max(0, w_1 * x + b_1)$$

$$+ u_2 * \max(0, w_2 * x + b_2)$$

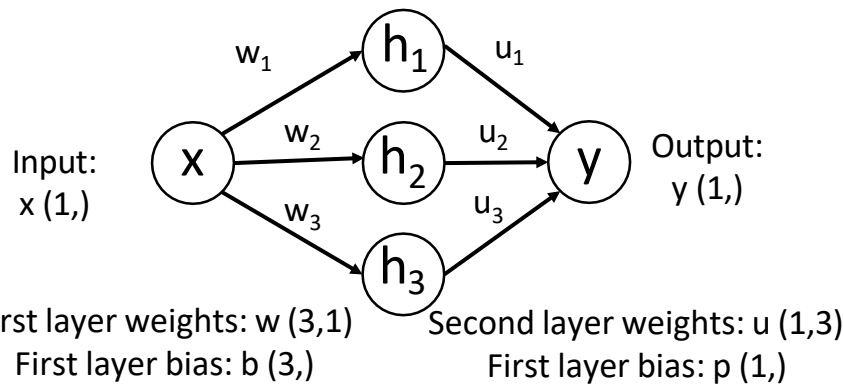
$$+ u_3 * \max(0, w_3 * x + b_3)$$

$$+ p$$



# Universal Approximation

Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



$$h_1 = \max(0, w_1 * x + b_1)$$

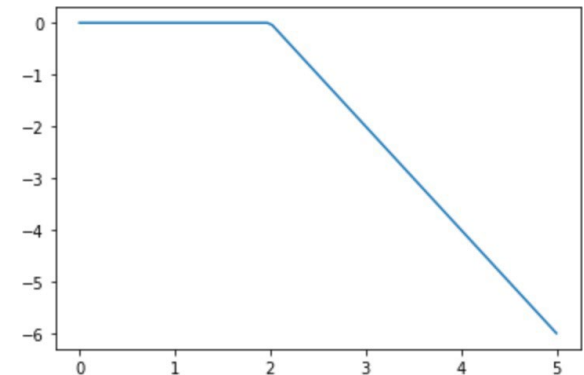
$$h_2 = \max(0, w_2 * x + b_2)$$

$$h_3 = \max(0, w_3 * x + b_3)$$

$$y = u_1 * h_1 + u_2 * h_2 + u_3 * h_3 + p$$

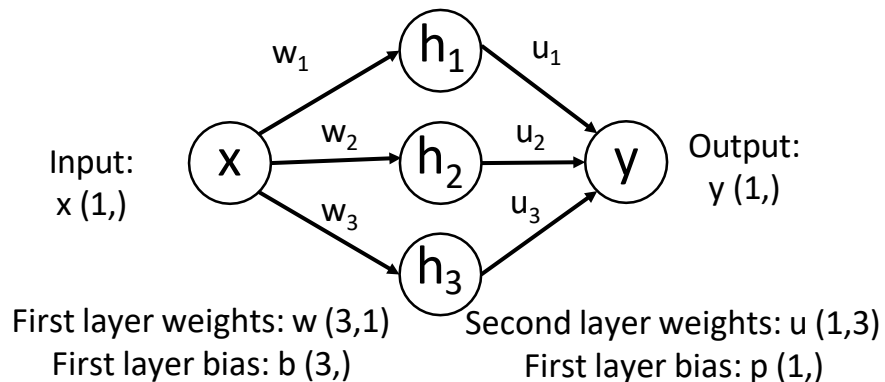
$$y = u_1 * \max(0, w_1 * x + b_1) + u_2 * \max(0, w_2 * x + b_2) + u_3 * \max(0, w_3 * x + b_3) + p$$

Output is a sum of shifted, scaled ReLUs:



# Universal Approximation

Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



$$h_1 = \max(0, w_1 * x + b_1)$$

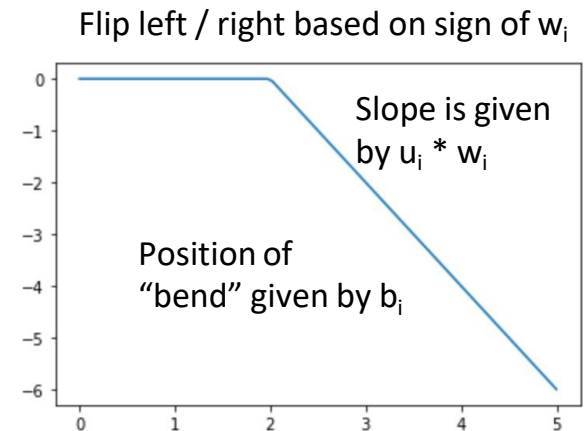
$$h_2 = \max(0, w_2 * x + b_2)$$

$$h_3 = \max(0, w_3 * x + b_3)$$

$$y = u_1 * h_1 + u_2 * h_2 + u_3 * h_3 + p$$

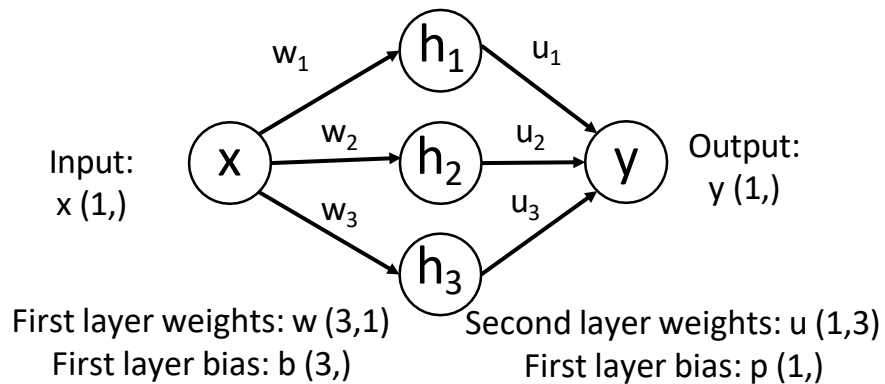
$$y = u_1 * \max(0, w_1 * x + b_1) + u_2 * \max(0, w_2 * x + b_2) + u_3 * \max(0, w_3 * x + b_3) + p$$

Output is a sum of shifted, scaled ReLUs:



# Universal Approximation

Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



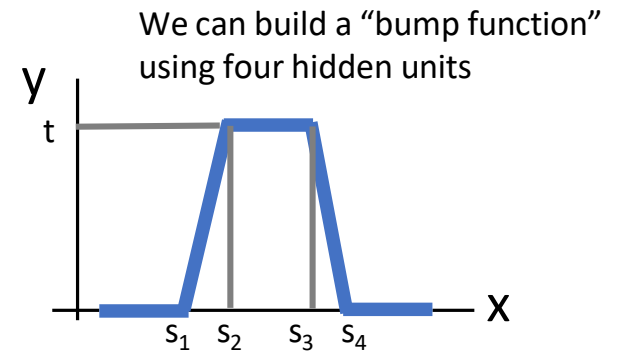
$$h_1 = \max(0, w_1 * x + b_1)$$

$$h_2 = \max(0, w_2 * x + b_2)$$

$$h_3 = \max(0, w_3 * x + b_3)$$

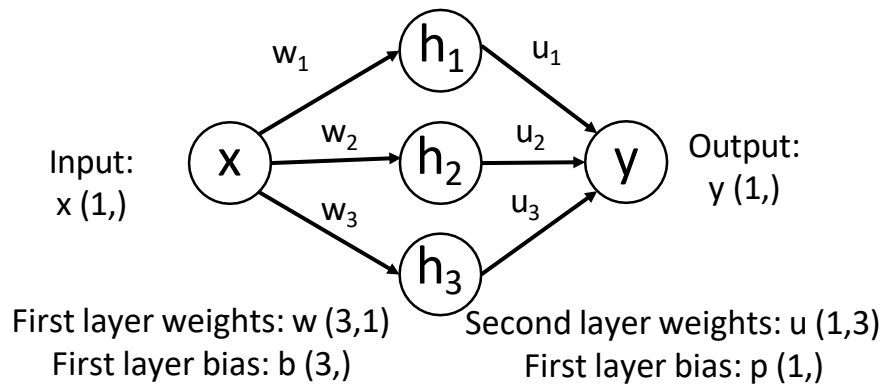
$$y = u_1 * h_1 + u_2 * h_2 + u_3 * h_3 + p$$

$$y = u_1 * \max(0, w_1 * x + b_1) + u_2 * \max(0, w_2 * x + b_2) + u_3 * \max(0, w_3 * x + b_3) + p$$



# Universal Approximation

Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



$$h_1 = \max(0, w_1 * x + b_1)$$

$$h_2 = \max(0, w_2 * x + b_2)$$

$$h_3 = \max(0, w_3 * x + b_3)$$

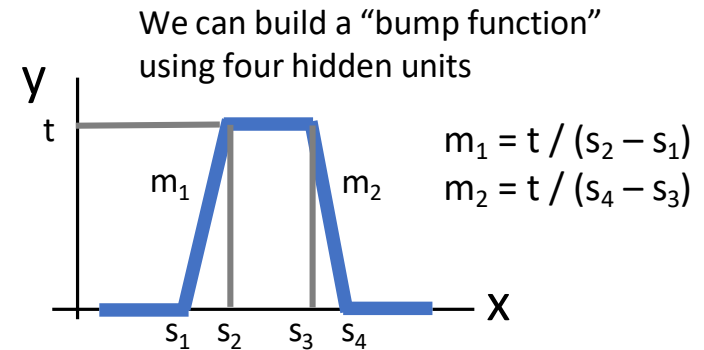
$$y = u_1 * h_1 + u_2 * h_2 + u_3 * h_3 + p$$

$$y = u_1 * \max(0, w_1 * x + b_1)$$

$$+ u_2 * \max(0, w_2 * x + b_2)$$

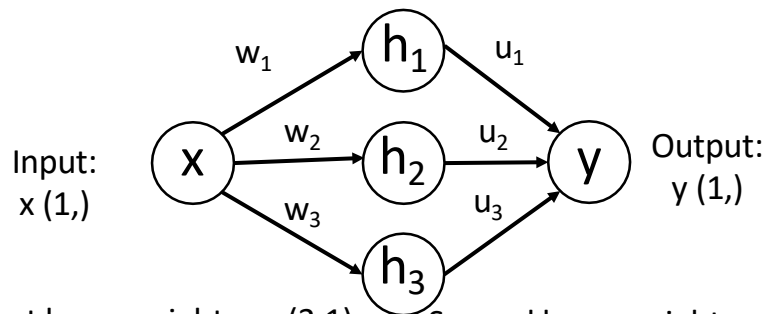
$$+ u_3 * \max(0, w_3 * x + b_3)$$

$$+ p$$



# Universal Approximation

Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



First layer weights:  $w (3,1)$

First layer bias:  $b (3,)$

Second layer weights:  $u (1,3)$

First layer bias:  $p (1,)$

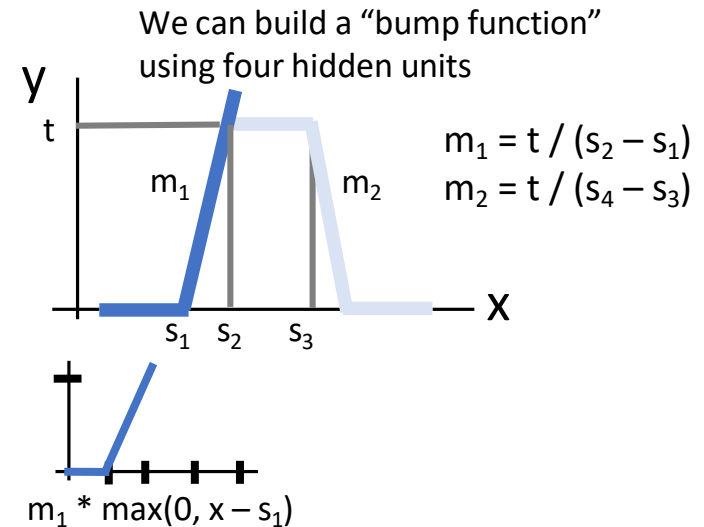
$$h_1 = \max(0, w_1 * x + b_1)$$

$$h_2 = \max(0, w_2 * x + b_2)$$

$$h_3 = \max(0, w_3 * x + b_3)$$

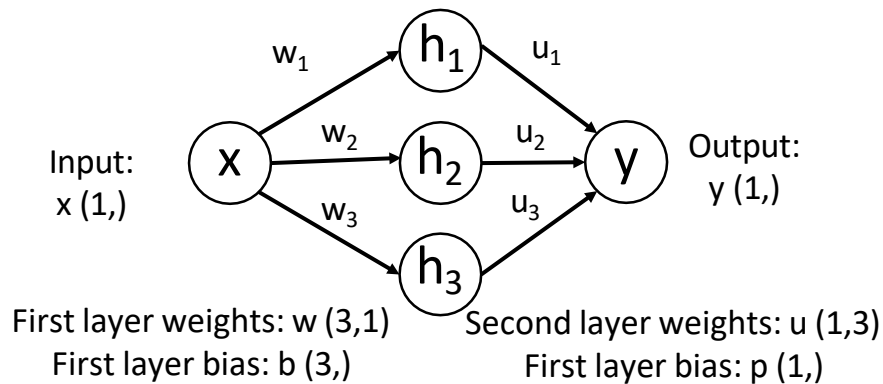
$$y = u_1 * h_1 + u_2 * h_2 + u_3 * h_3 + p$$

$$y = u_1 * \max(0, w_1 * x + b_1) + u_2 * \max(0, w_2 * x + b_2) + u_3 * \max(0, w_3 * x + b_3) + p$$



# Universal Approximation

Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



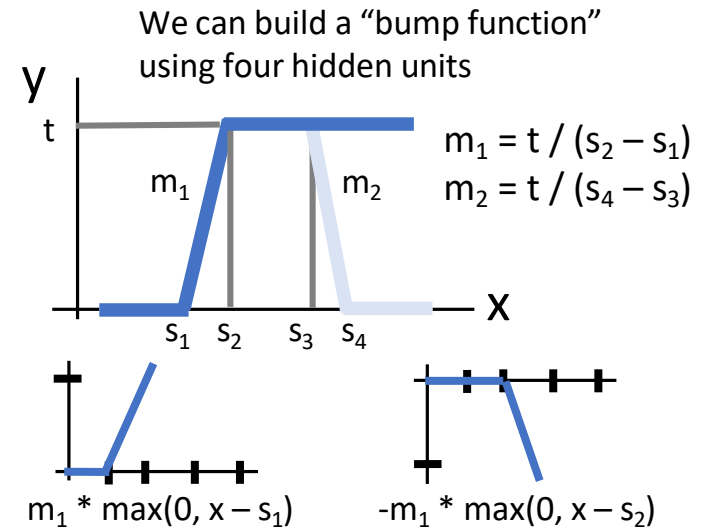
$$h_1 = \max(0, w_1 * x + b_1)$$

$$h_2 = \max(0, w_2 * x + b_2)$$

$$h_3 = \max(0, w_3 * x + b_3)$$

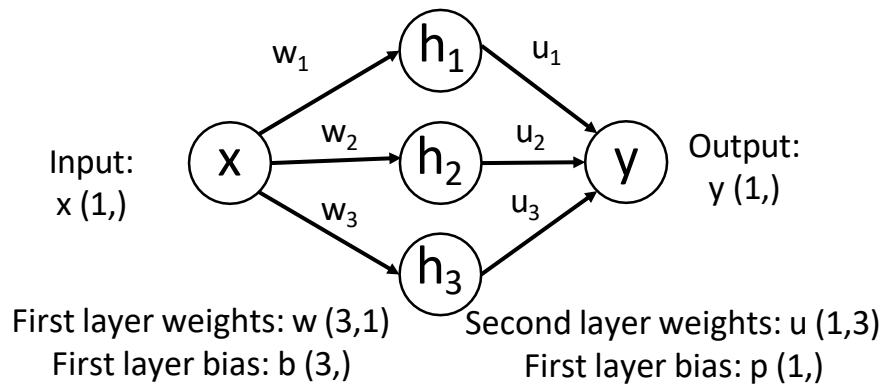
$$y = u_1 * h_1 + u_2 * h_2 + u_3 * h_3 + p$$

$$y = u_1 * \max(0, w_1 * x + b_1) + u_2 * \max(0, w_2 * x + b_2) + u_3 * \max(0, w_3 * x + b_3) + p$$



# Universal Approximation

Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



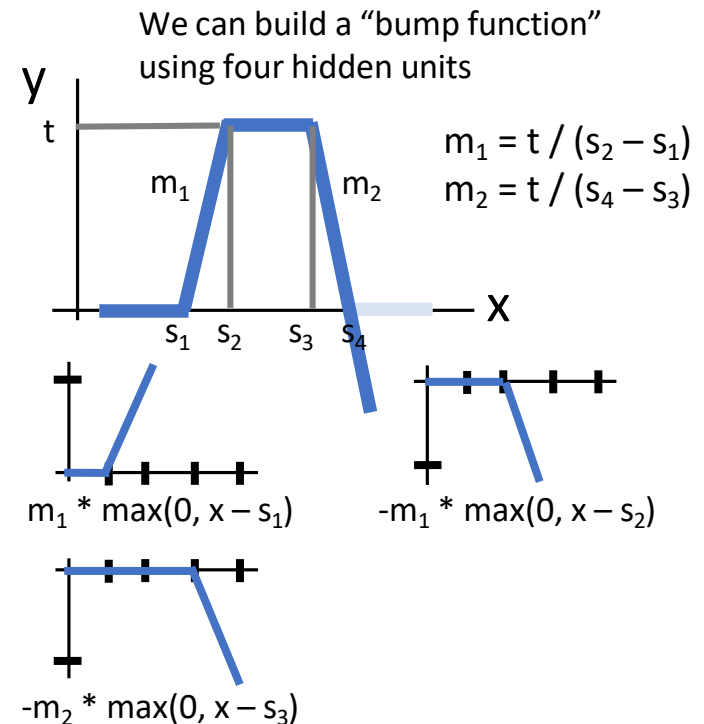
$$h_1 = \max(0, w_1 * x + b_1)$$

$$h_2 = \max(0, w_2 * x + b_2)$$

$$h_3 = \max(0, w_3 * x + b_3)$$

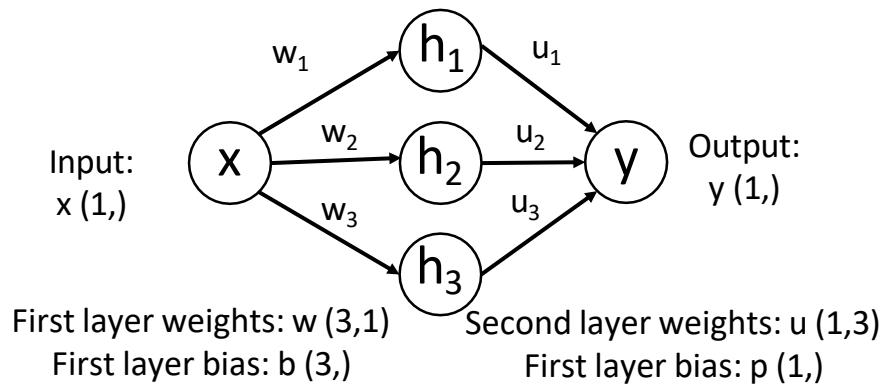
$$y = u_1 * h_1 + u_2 * h_2 + u_3 * h_3 + p$$

$$y = u_1 * \max(0, w_1 * x + b_1) + u_2 * \max(0, w_2 * x + b_2) + u_3 * \max(0, w_3 * x + b_3) + p$$



# Universal Approximation

Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



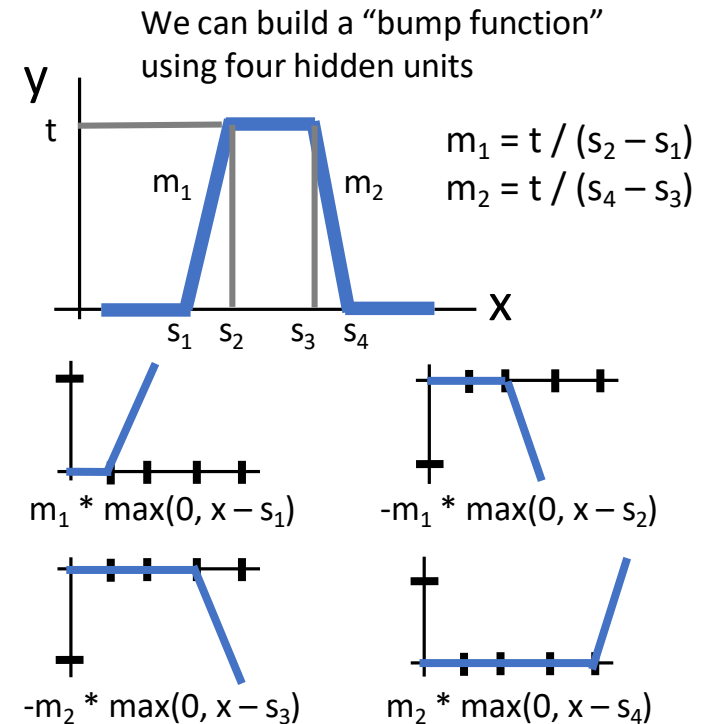
$$h_1 = \max(0, w_1 * x + b_1)$$

$$h_2 = \max(0, w_2 * x + b_2)$$

$$h_3 = \max(0, w_3 * x + b_3)$$

$$y = u_1 * h_1 + u_2 * h_2 + u_3 * h_3 + p$$

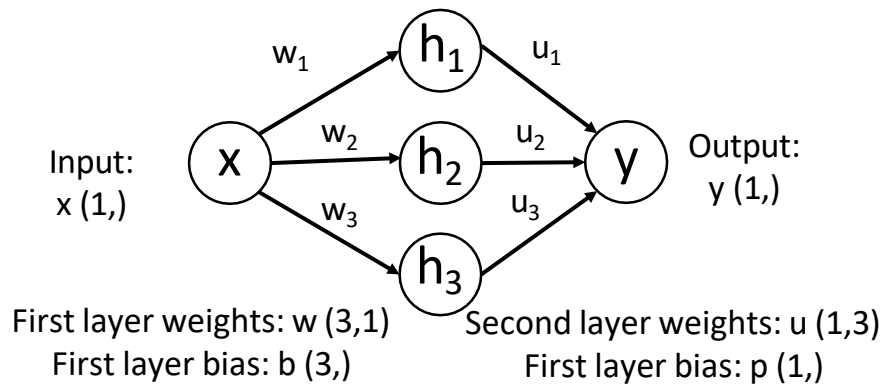
$$y = u_1 * \max(0, w_1 * x + b_1) + u_2 * \max(0, w_2 * x + b_2) + u_3 * \max(0, w_3 * x + b_3) + p$$





# Universal Approximation

Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



$$h_1 = \max(0, w_1 * x + b_1)$$

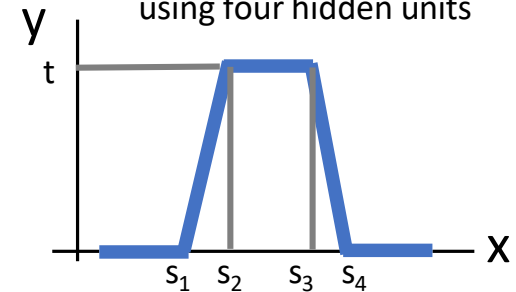
$$h_2 = \max(0, w_2 * x + b_2)$$

$$h_3 = \max(0, w_3 * x + b_3)$$

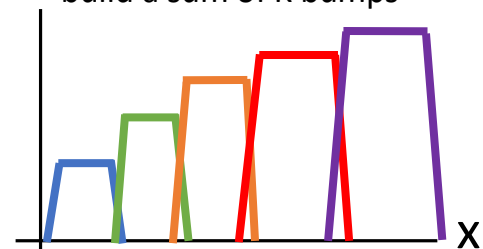
$$y = u_1 * h_1 + u_2 * h_2 + u_3 * h_3 + p$$

$$y = u_1 * \max(0, w_1 * x + b_1) + u_2 * \max(0, w_2 * x + b_2) + u_3 * \max(0, w_3 * x + b_3) + p$$

We can build a “bump function” using four hidden units

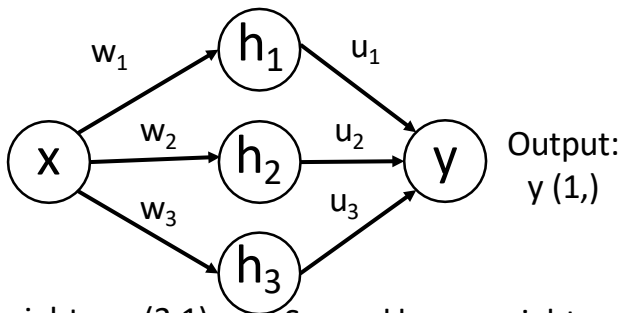


With  $4K$  hidden units we can build a sum of  $K$  bumps



# Universal Approximation

Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



First layer weights:  $w(3,1)$

First layer bias:  $b(3,)$

Second layer weights:  $u(1,3)$

First layer bias:  $p(1,)$

$$h_1 = \max(0, w_1 * x + b_1)$$

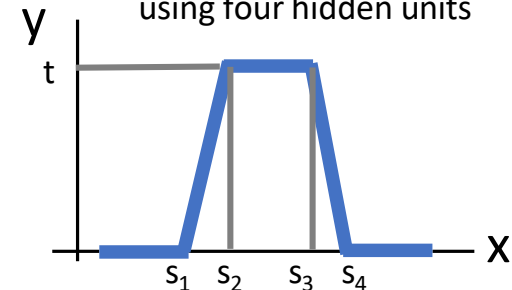
$$h_2 = \max(0, w_2 * x + b_2)$$

$$h_3 = \max(0, w_3 * x + b_3)$$

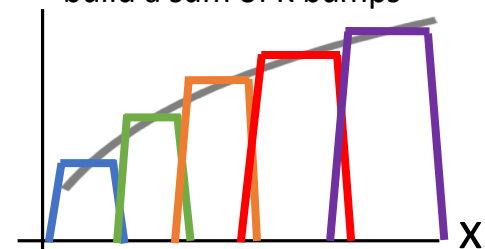
$$y = u_1 * h_1 + u_2 * h_2 + u_3 * h_3 + p$$

$$y = u_1 * \max(0, w_1 * x + b_1) + u_2 * \max(0, w_2 * x + b_2) + u_3 * \max(0, w_3 * x + b_3) + p$$

We can build a “bump function” using four hidden units



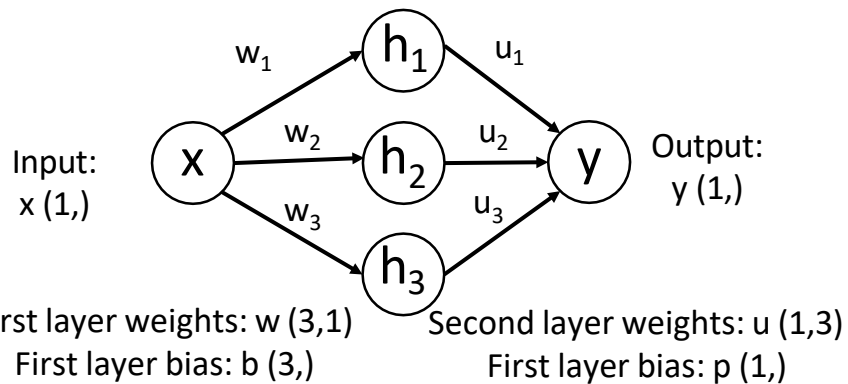
With 4K hidden units we can build a sum of K bumps



Approximate functions with bumps!

# Universal Approximation

Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



$$h_1 = \max(0, w_1 * x + b_1)$$

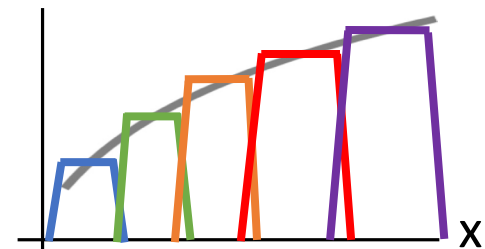
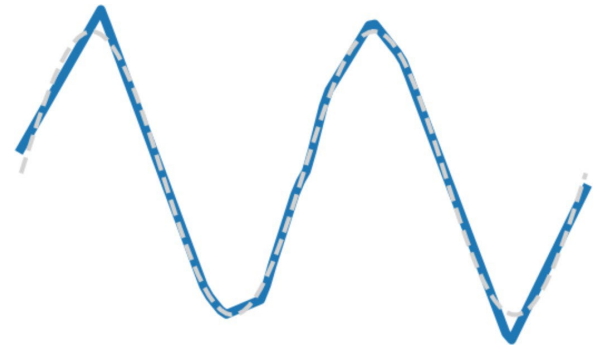
$$h_2 = \max(0, w_2 * x + b_2)$$

$$h_3 = \max(0, w_3 * x + b_3)$$

$$y = u_1 * h_1 + u_2 * h_2 + u_3 * h_3 + p$$

$$y = u_1 * \max(0, w_1 * x + b_1) + u_2 * \max(0, w_2 * x + b_2) + u_3 * \max(0, w_3 * x + b_3) + p$$

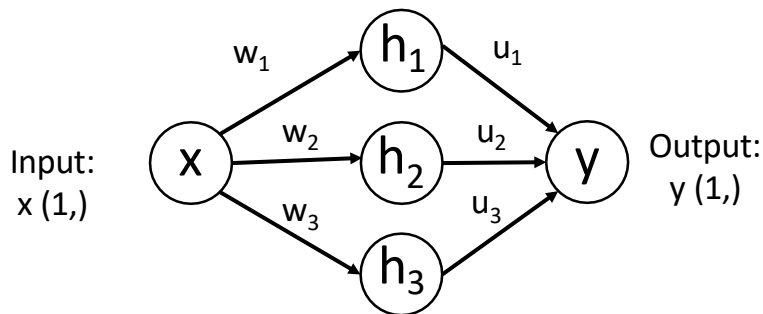
Reality check: Networks don't really learn bumps!



Approximate functions with bumps!

# Universal Approximation

Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



Universal approximation tells us:

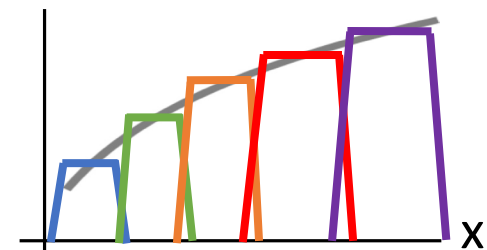
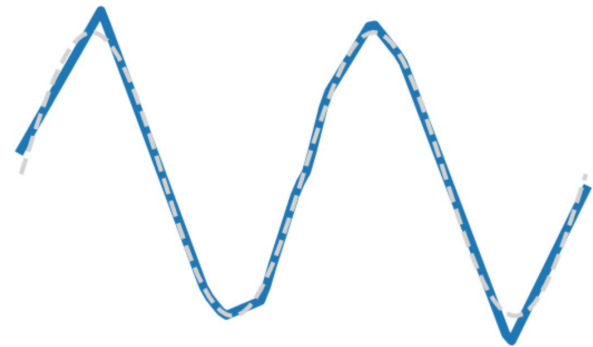
- Neural nets can represent any function

Universal approximation DOES NOT tell us:

- Whether we can actually learn any function with SGD
- How much data we need to learn a function

Remember: kNN is also a universal approximator!

Reality check: Networks don't really learn bumps!



Approximate functions with bumps!

# Convex Functions

---

A function  $f : X \subseteq \mathbb{R}^N \rightarrow \mathbb{R}$  is **convex** if for all  $x_1, x_2 \in X, t \in [0, 1]$ ,

$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2)$$

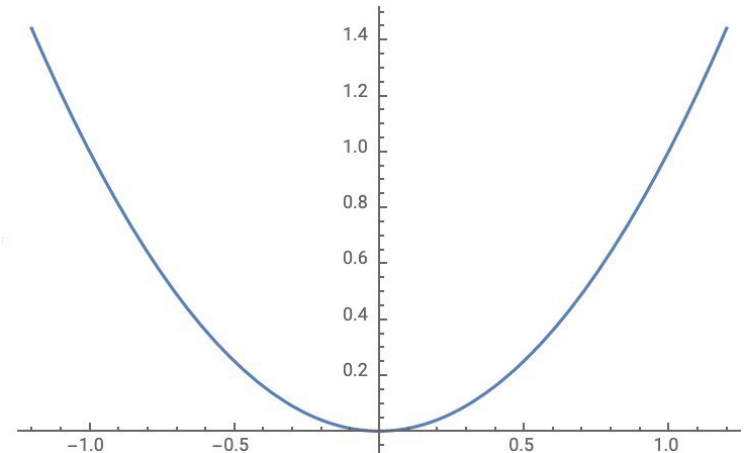
# Convex Functions

---

A function  $f : X \subseteq \mathbb{R}^N \rightarrow \mathbb{R}$  is **convex** if for all  $x_1, x_2 \in X, t \in [0, 1]$ ,

$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2)$$

Example:  $f(x) = x^2$  is convex:



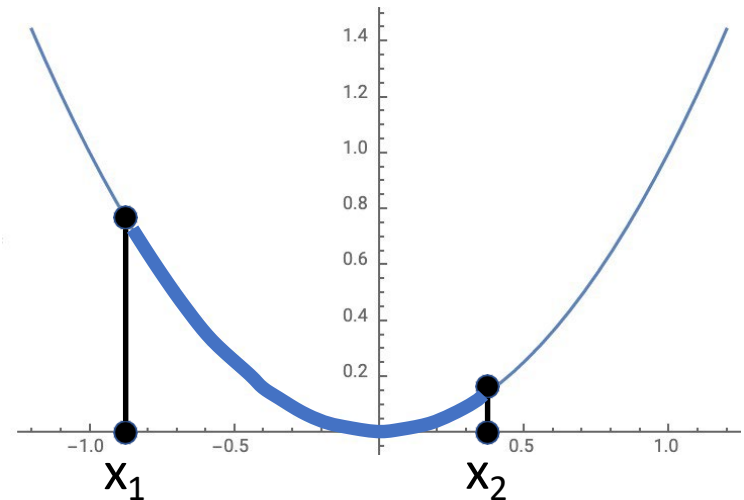
# Convex Functions

---

A function  $f : X \subseteq \mathbb{R}^N \rightarrow \mathbb{R}$  is **convex** if for all  $x_1, x_2 \in X, t \in [0, 1]$ ,

$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2)$$

Example:  $f(x) = x^2$  is convex:



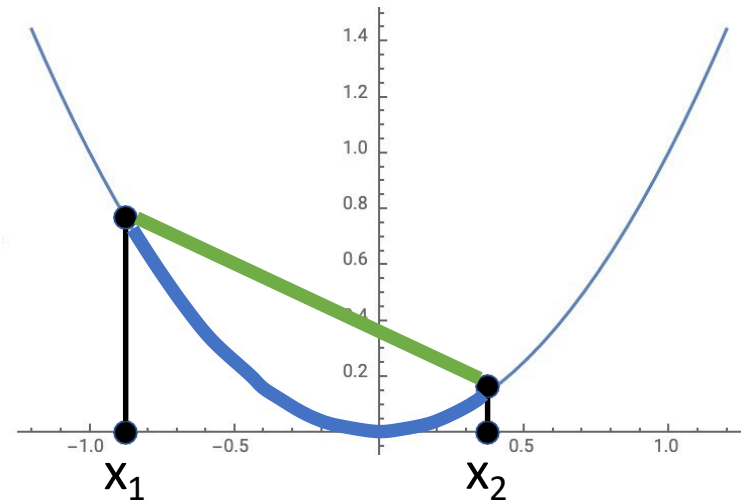
# Convex Functions

---

A function  $f : X \subseteq \mathbb{R}^N \rightarrow \mathbb{R}$  is **convex** if for all  $x_1, x_2 \in X, t \in [0, 1]$ ,

$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2)$$

Example:  $f(x) = x^2$  is convex:





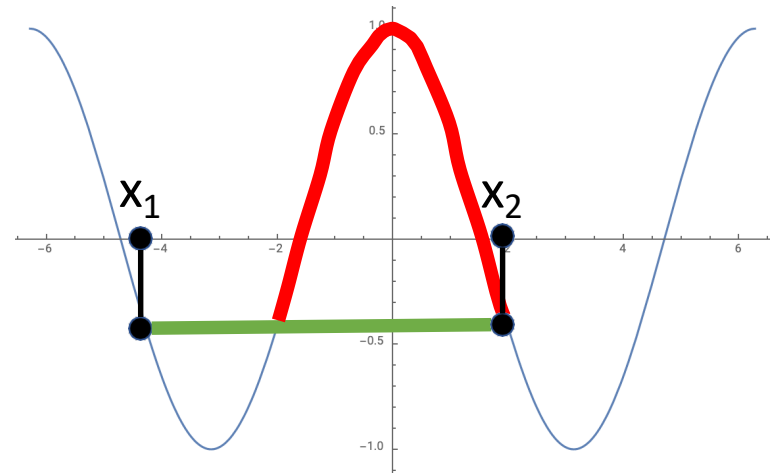
# Convex Functions

---

A function  $f : X \subseteq \mathbb{R}^N \rightarrow \mathbb{R}$  is **convex** if for all  $x_1, x_2 \in X, t \in [0, 1]$ ,

$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2)$$

Example:  $f(x) = \cos(x)$   
is not convex:



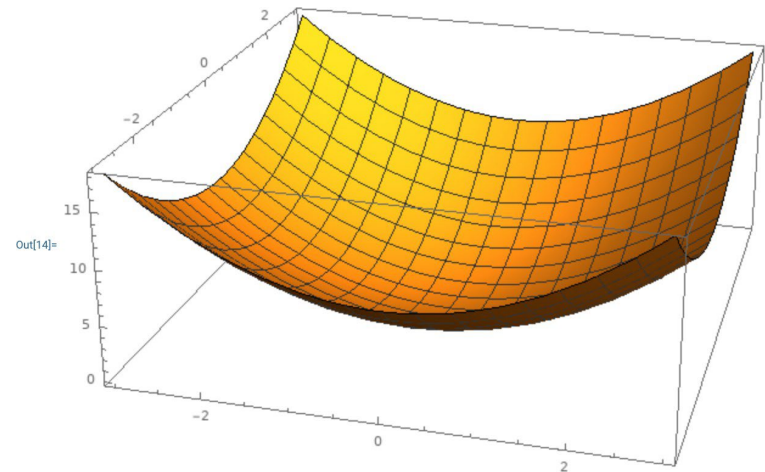
# Convex Functions

---

A function  $f : X \subseteq \mathbb{R}^N \rightarrow \mathbb{R}$  is **convex** if for all  $x_1, x_2 \in X, t \in [0, 1]$

$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2) \quad ,$$

**Intuition:** A convex function is  
a (multidimensional) bowl



\*Many technical details! See e.g. IOE 661 / MATH 663

# Convex Functions

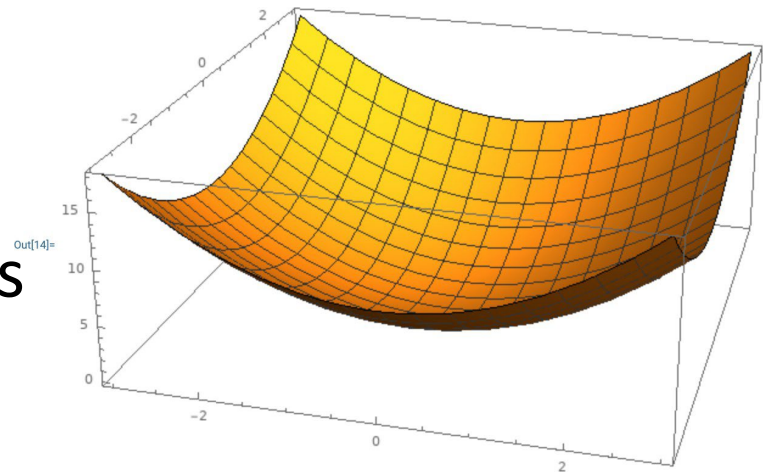
---

A function  $f : X \subseteq \mathbb{R}^N \rightarrow \mathbb{R}$  is **convex** if for all  $x_1, x_2 \in X, t \in [0, 1]$ ,

$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2)$$

**Intuition:** A convex function is a (multidimensional) bowl

Generally speaking, convex functions are **easy to optimize**: can derive theoretical guarantees about **converging to global minimum\***



# Convex Functions

---

A function  $f : X \subseteq \mathbb{R}^N \rightarrow \mathbb{R}$  is **convex** if for all  $x_1, x_2 \in X, t \in [0, 1]$ ,

$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2)$$

**Intuition:** A convex function is a (multidimensional) bowl

Generally speaking, convex functions are **easy to optimize**: can derive theoretical guarantees about **converging to global minimum\***

Linear classifiers optimize a **convex function!**

$$s = f(x; W) = Wx$$

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

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

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

$R(W)$  = L2 or L1 regularization

# Convex Functions

---

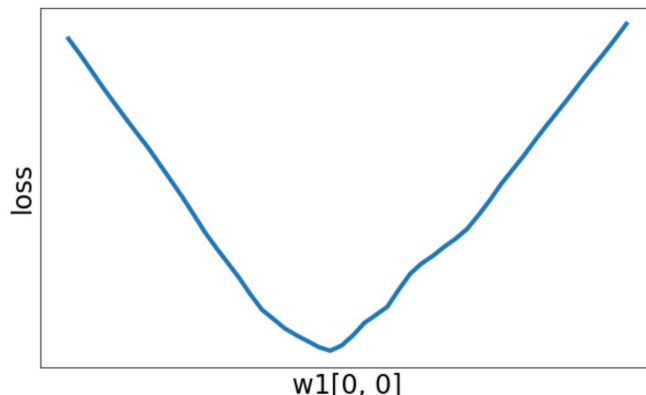
A function  $f : X \subseteq \mathbb{R}^N \rightarrow \mathbb{R}$  is **convex** if for all  $x_1, x_2 \in X, t \in [0, 1]$ ,

$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2)$$

**Intuition:** A convex function is a (multidimensional) bowl

Generally speaking, convex functions are **easy to optimize**: can derive theoretical guarantees about **converging to global minimum\***

Neural net losses sometimes look convex-ish:



1D slice of loss landscape for a 4-layer ReLU network with 10 input features, 32 units per hidden layer, 10 categories, with softmax loss

# Convex Functions

---

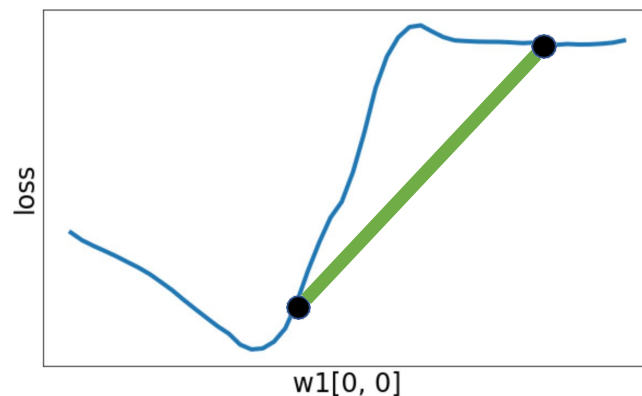
A function  $f : X \subseteq \mathbb{R}^N \rightarrow \mathbb{R}$  is **convex** if for all  $x_1, x_2 \in X, t \in [0, 1]$ ,

$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2)$$

**Intuition:** A convex function is a (multidimensional) bowl

Generally speaking, convex functions are **easy to optimize**: can derive theoretical guarantees about **converging to global minimum\***

But often clearly nonconvex:



1D slice of loss landscape for a 4-layer ReLU network with 10 input features, 32 units per hidden layer, 10 categories, with softmax loss

# Convex Functions

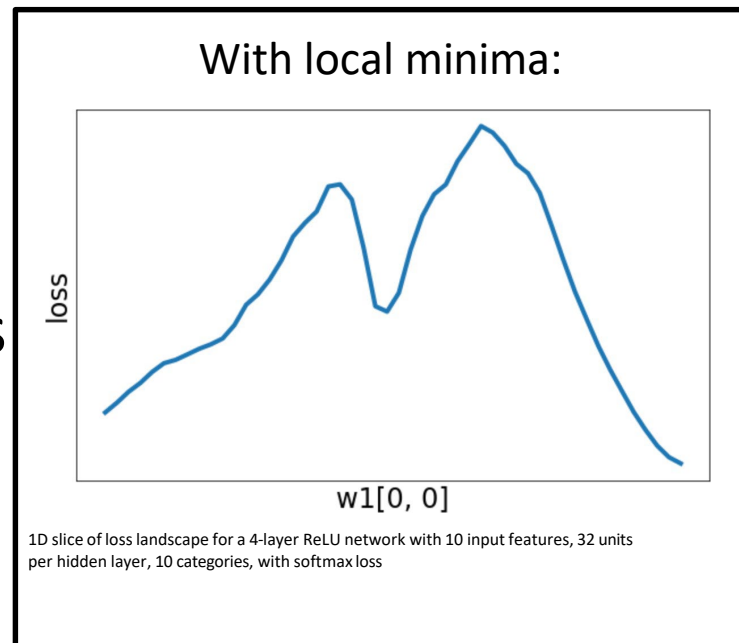
---

A function  $f : X \subseteq \mathbb{R}^N \rightarrow \mathbb{R}$  is **convex** if for all  $x_1, x_2 \in X, t \in [0, 1]$ ,

$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2)$$

**Intuition:** A convex function is a (multidimensional) bowl

Generally speaking, convex functions are **easy to optimize**: can derive theoretical guarantees about **converging to global minimum\***



# Convex Functions

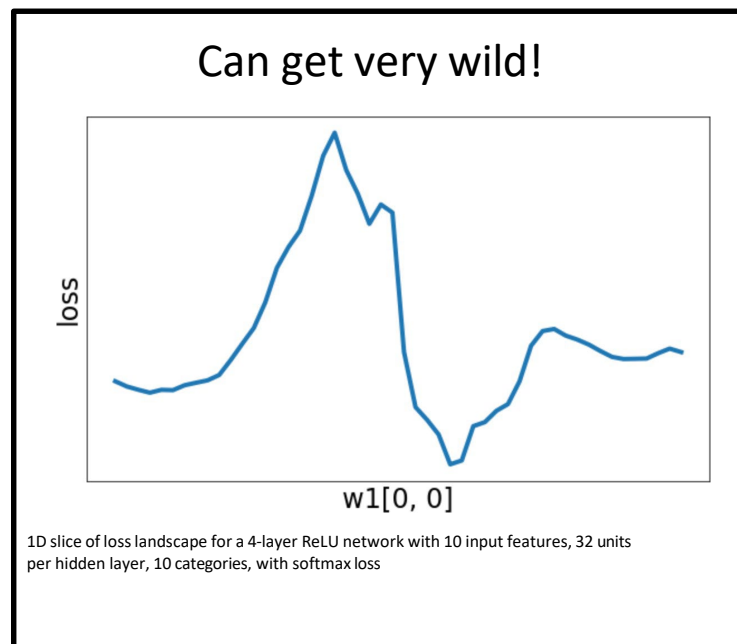
---

A function  $f : X \subseteq \mathbb{R}^N \rightarrow \mathbb{R}$  is **convex** if for all  $x_1, x_2 \in X, t \in [0, 1]$ ,

$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2)$$

**Intuition:** A convex function is a (multidimensional) bowl

Generally speaking, convex functions are **easy to optimize**: can derive theoretical guarantees about **converging to global minimum**\*





# Convex Functions

---

A function  $f : X \subseteq \mathbb{R}^N \rightarrow \mathbb{R}$  is **convex** if for all  $x_1, x_2 \in X, t \in [0, 1]$ ,

$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2)$$

**Intuition:** A convex function is a (multidimensional) bowl

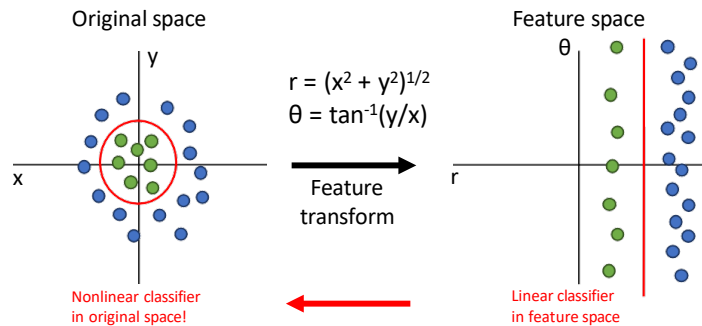
Generally speaking, convex functions are **easy to optimize**: can derive theoretical guarantees about **converging to global minimum\***

Most neural networks need **nonconvex optimization**

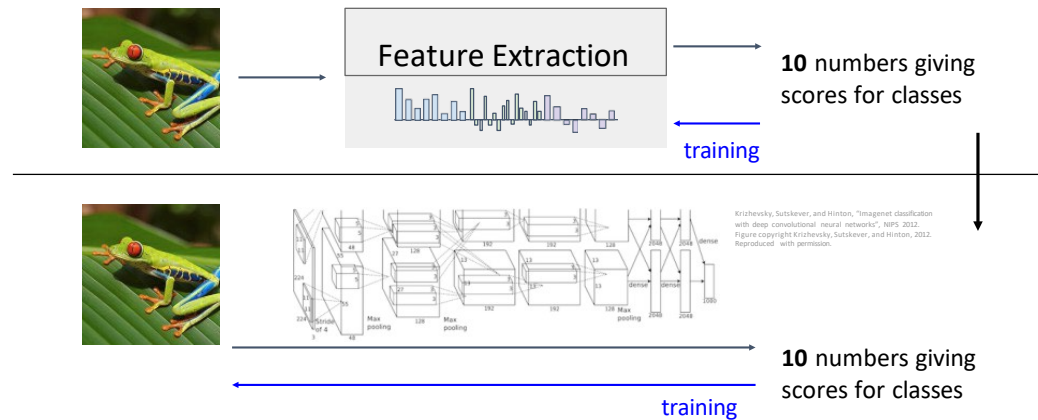
- Few or no guarantees about convergence
- Empirically it seems to work anyway
- Active area of research

# Summary

Feature transform + Linear classifier  
allows nonlinear decision boundaries



Neural Networks as learnable feature transforms

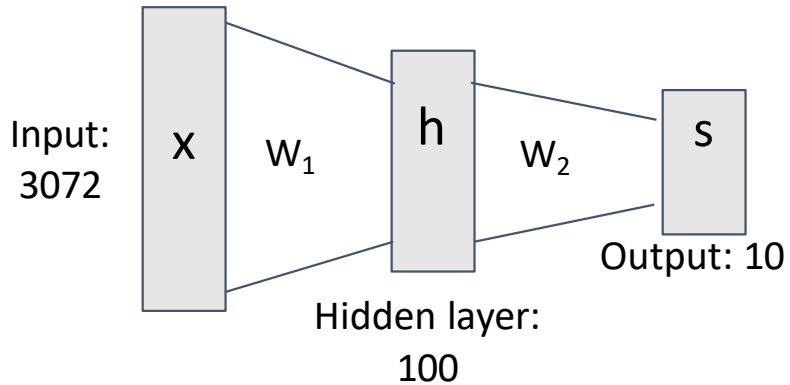


# Summary

---

From linear classifiers to  
fully-connected networks

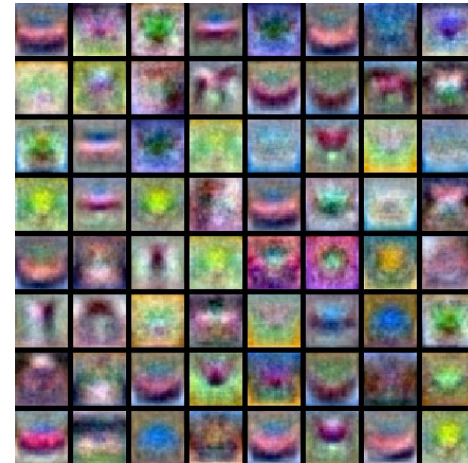
$$f = W_2 \max(0, W_1 x)$$



Linear classifier: One template per class



Neural networks: Many reusable templates

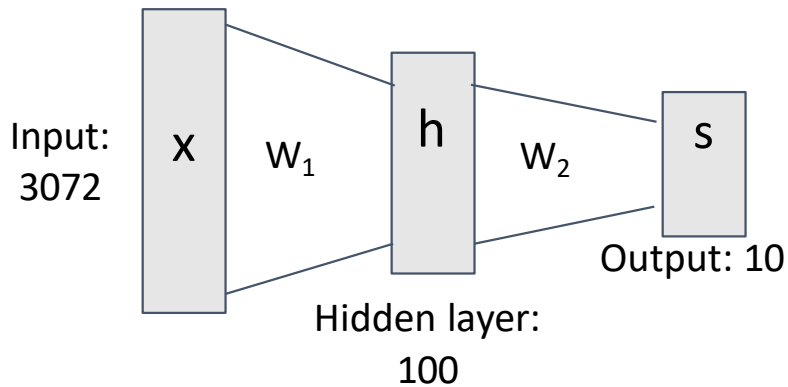


# Summary

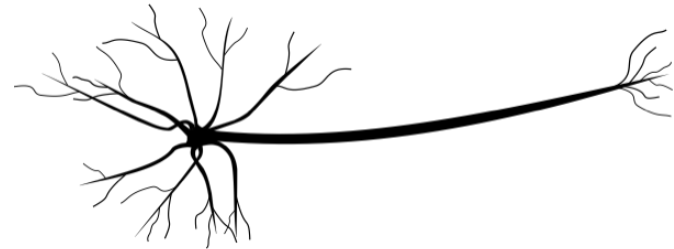
---

From linear classifiers to  
fully-connected networks

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



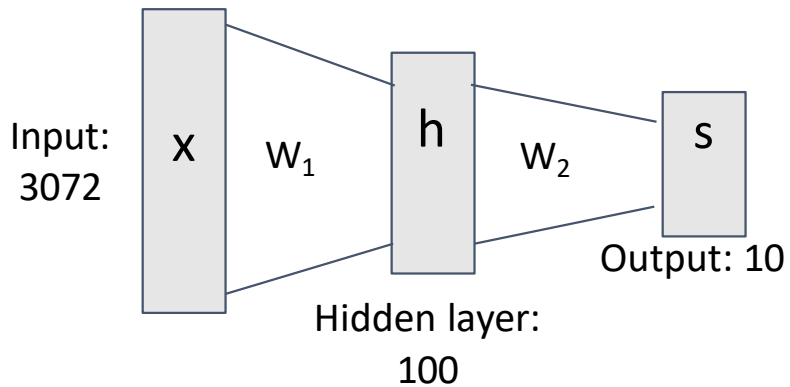
Neural networks loosely inspired by biological  
neurons but be careful with analogies



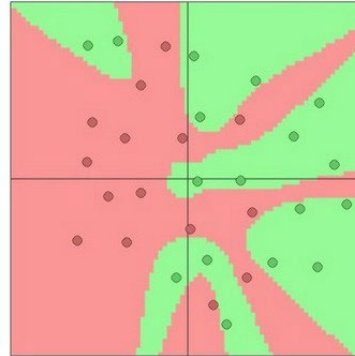
# Summary

From linear classifiers to  
fully-connected networks

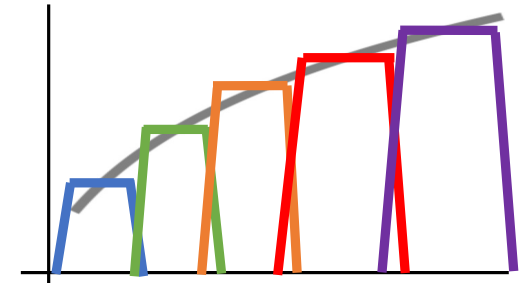
$$f = W_2 \max(0, W_1 x)$$



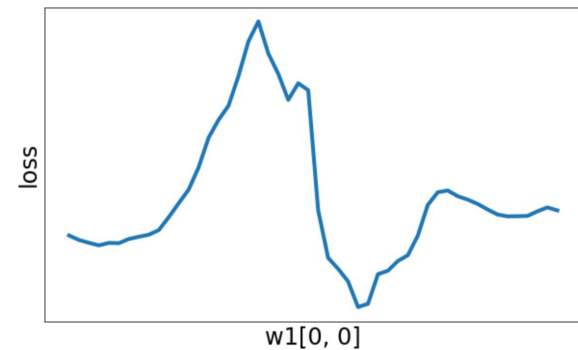
Space Warping



Universal Approximation



Nonconvex



# Problem: How to compute gradients?

---

$$s = f(x; W_1, W_2) = W_2 \max(0, W_1 x) \quad \text{Nonlinear score function}$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \quad \text{SVM Loss on predictions}$$

$$R(W) = \sum_k W_k^2 \quad \text{Regularization}$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i + \lambda R(W_1) + \lambda R(W_2) \quad \text{Total loss: data loss + regularization}$$

If we can compute  $\frac{\partial L}{\partial W_1}, \frac{\partial L}{\partial W_2}$  we can learn  $W_1$  and  $W_2$