

# Lecture 5: Neural Networks

## Waitlist update

I was confused about the way waitlists work on Monday =(

We have set enrollment sizes of 35 / 85 for 498 / 598

Each day overrides will be sent automatically in waitlist order to fill up to capacity

If you don't enroll within a day of getting an override you will be dropped from the waitlist

# Assignment 1

Was due on Sunday

If you use all 3 late days then you can turn it in today with no penalty

If you enrolled late, your A1 will be due **one week from the time you enrolled**

## Assignment 2

Due Monday, September 30

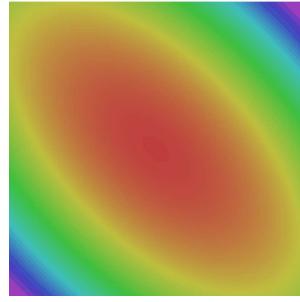
Much longer than A1 – Start early

Your submission **must** pass the validation script to be graded!

We will be lenient on A1 submissions, but starting with A2 we will not grade your assignment if it does not pass the validation script

Where we are:

1. Use **Linear Models** for image classification problems
  2. Use **Loss Functions** to express preferences over different choices of weights
  3. Use **Stochastic Gradient Descent** to minimize our loss functions and train the model
- $s = f(x; W) = Wx$
- 



```
v = 0  
for t in range(num_steps):  
    dw = compute_gradient(w)  
    v = rho * v + dw  
    w -= learning_rate * v
```

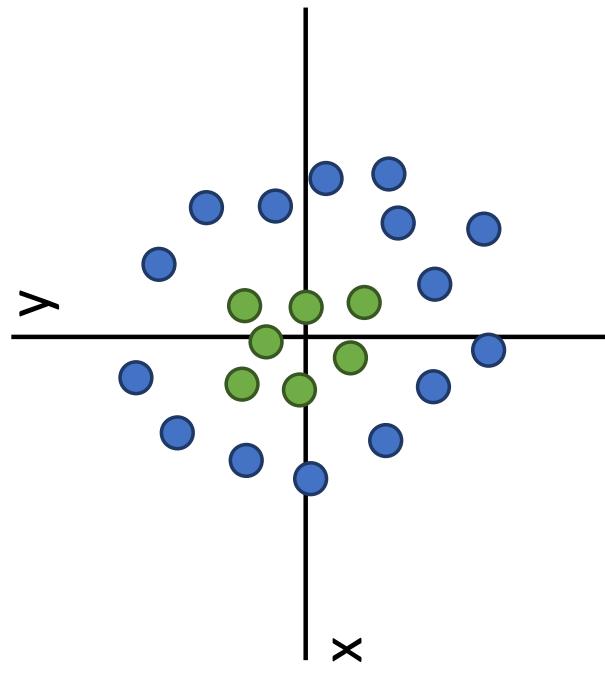
September 18, 2019

Lecture 5 - 5

Justin Johnson

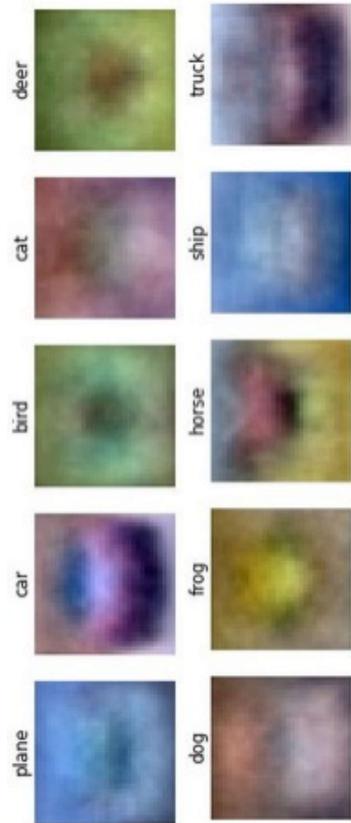
**Problem:** Linear Classifiers aren't that powerful

### Geometric Viewpoint

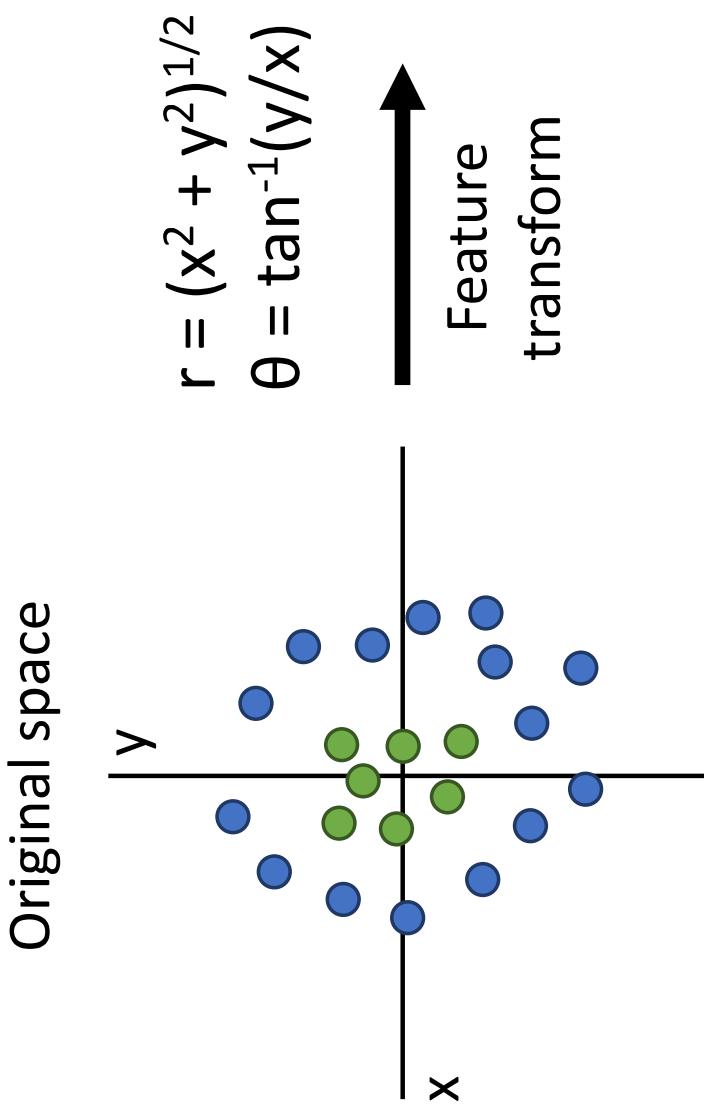


### Visual Viewpoint

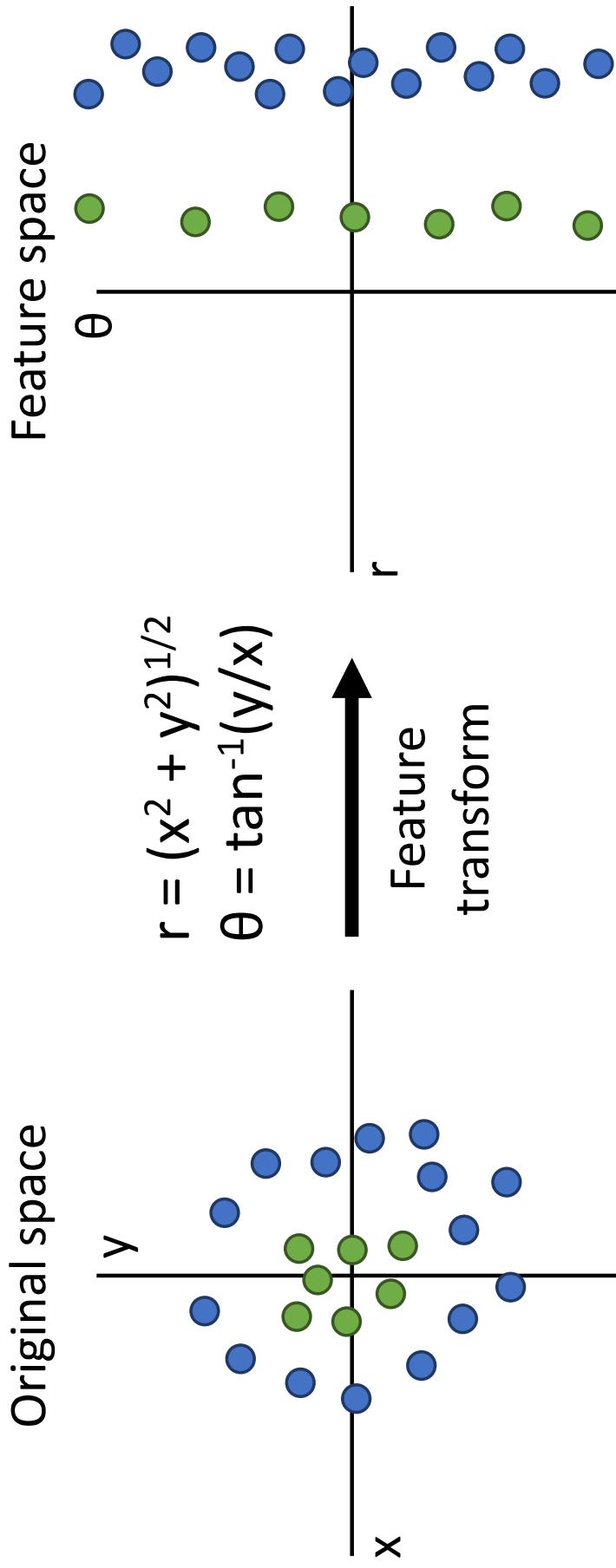
One template per class:  
Can't recognize different  
modes of a class



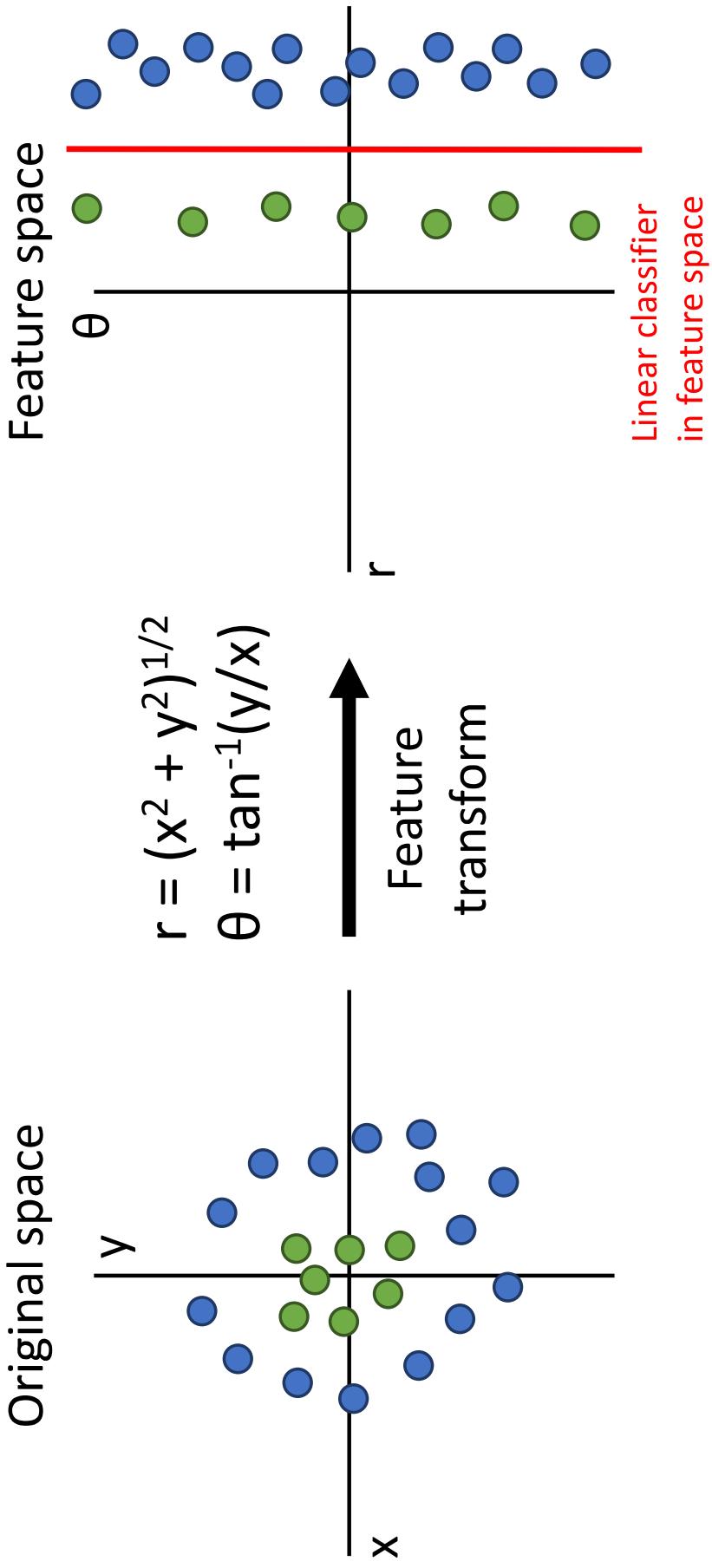
# One solution: Feature Transforms



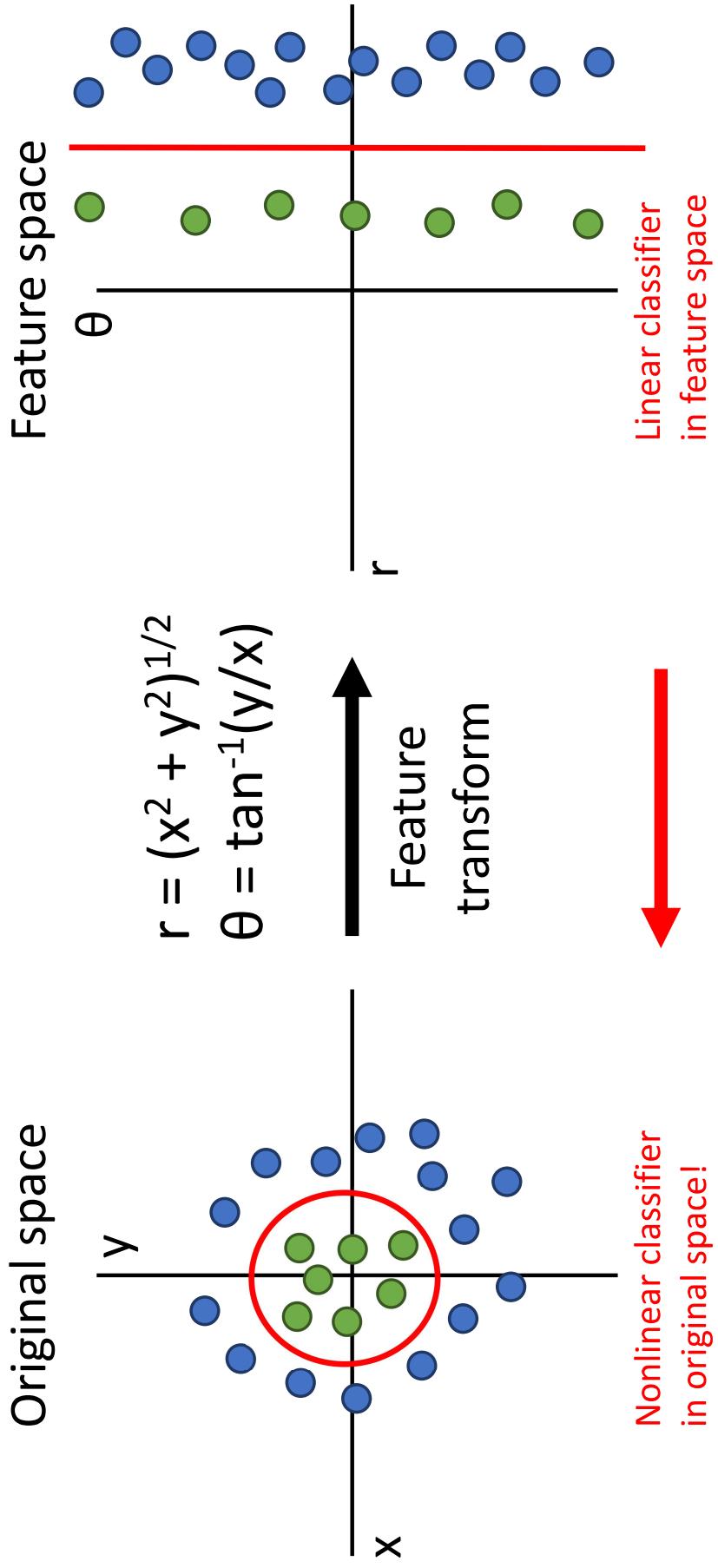
# One solution: Feature Transforms



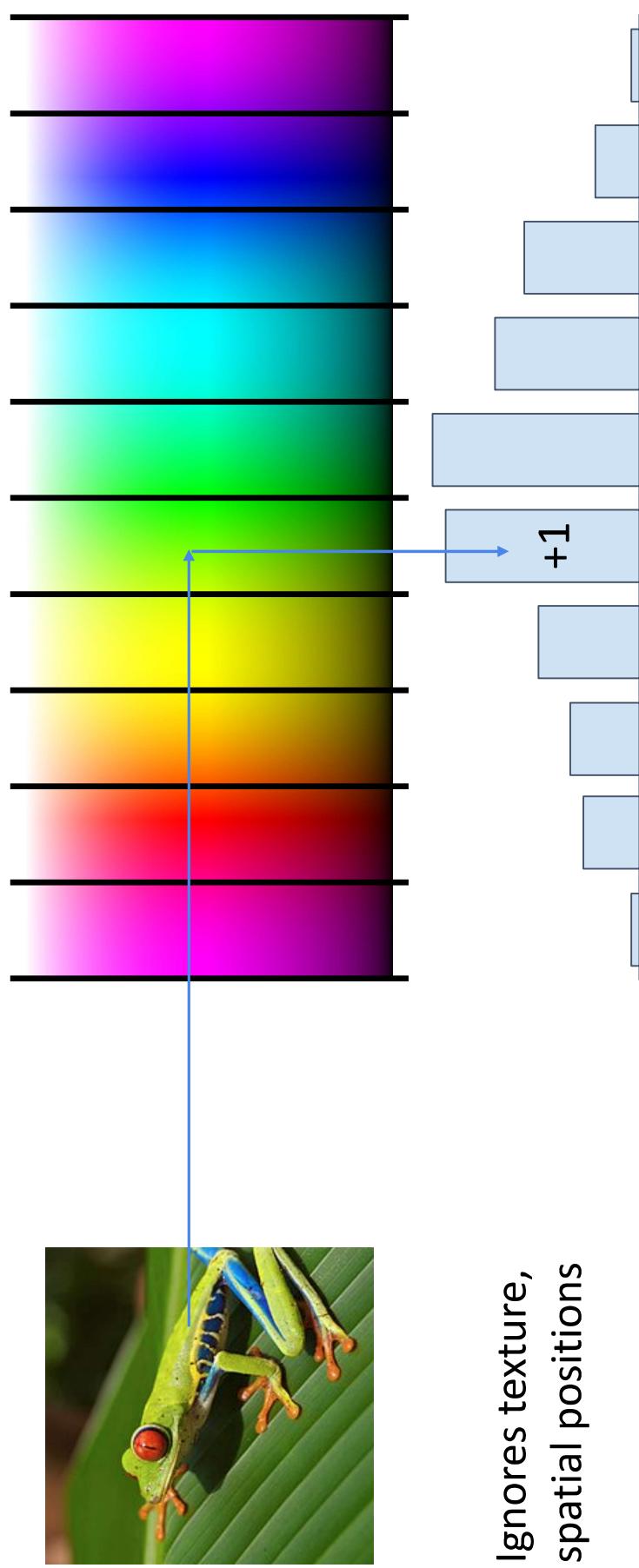
# One solution: Feature Transforms



# One solution: Feature Transforms



# Image Features: Color Histogram



Frog image is in the public domain

Justin Johnson

Lecture 5 - 11

September 18, 2019

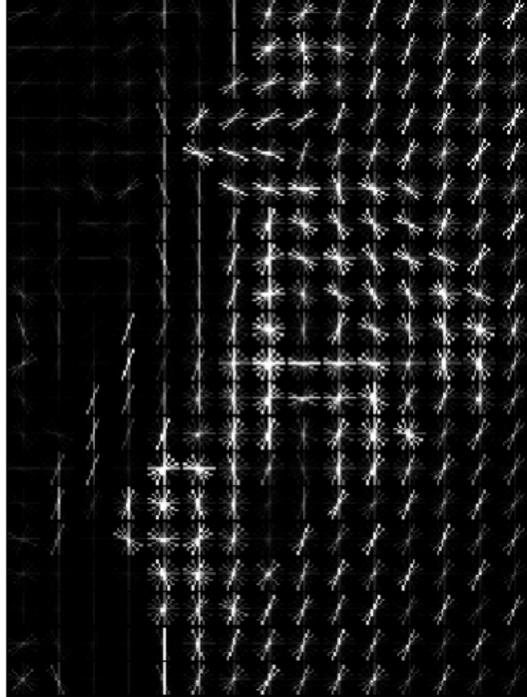
# Image Features: Histogram of Oriented Gradients (HoG)



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

Lowe, "Object recognition from local scale-invariant features", ICCV 1999  
Dalal and Triggs, "Histograms of oriented gradients for human detection," CVPR 2005

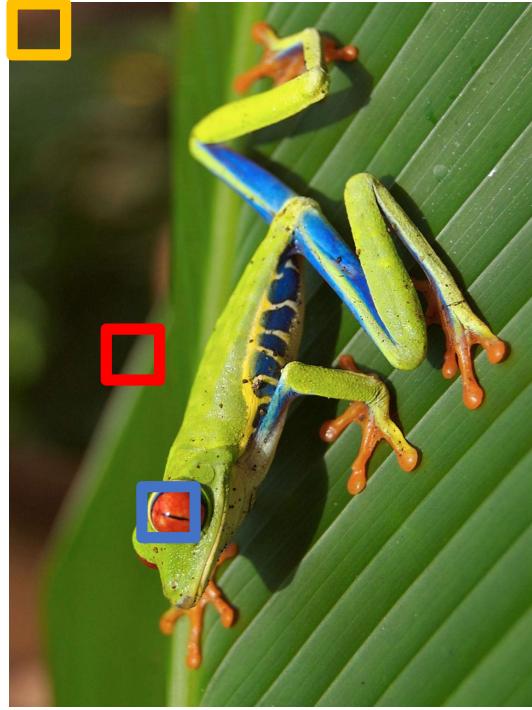
# Image Features: Histogram of Oriented Gradients (HoG)



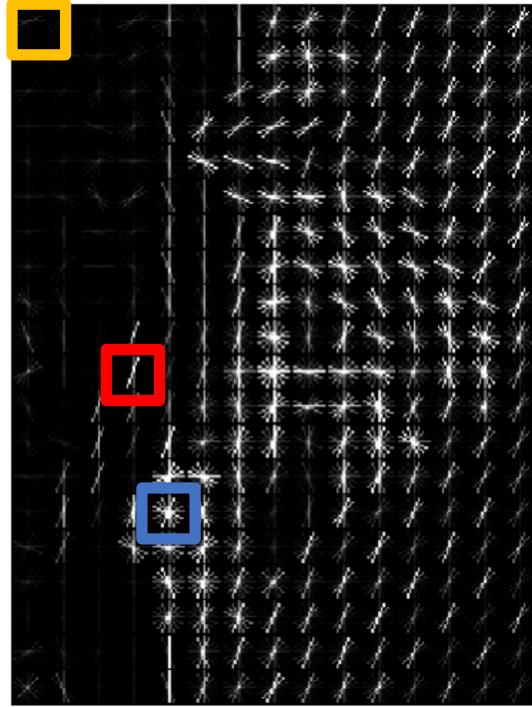
1. Compute edge direction / strength at each pixel
  2. Divide image into 8x8 regions
  3. Within each region compute a histogram of edge directions weighted by edge strength
- Example: 320x240 image gets divided into 40x30 bins; 8 directions per bin; feature vector has  $30 * 40 * 9 = 10,800$  numbers

Lowe, "Object recognition from local scale-invariant features", ICCV 1999  
Dalal and Triggs, "Histograms of oriented gradients for human detection," CVPR 2005

# Image Features: Histogram of Oriented Gradients (HoG)



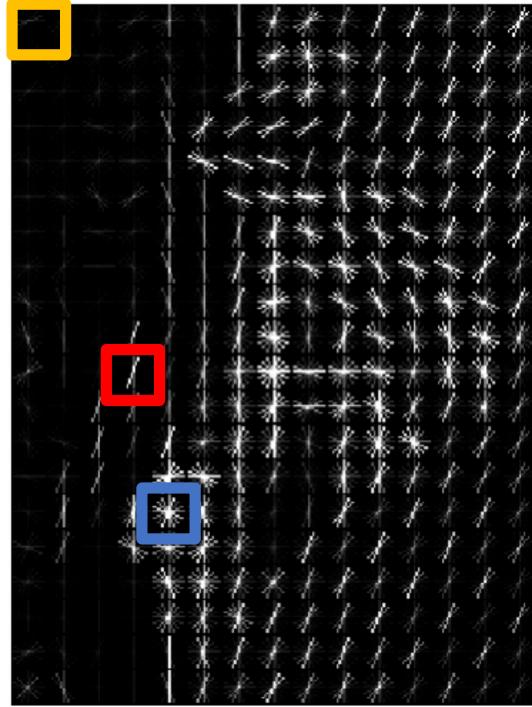
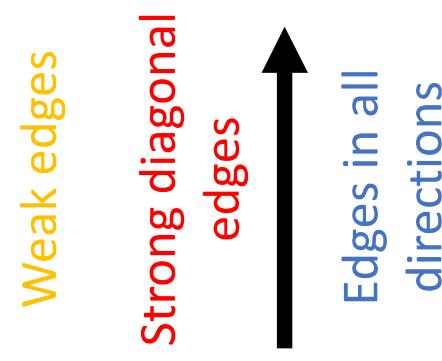
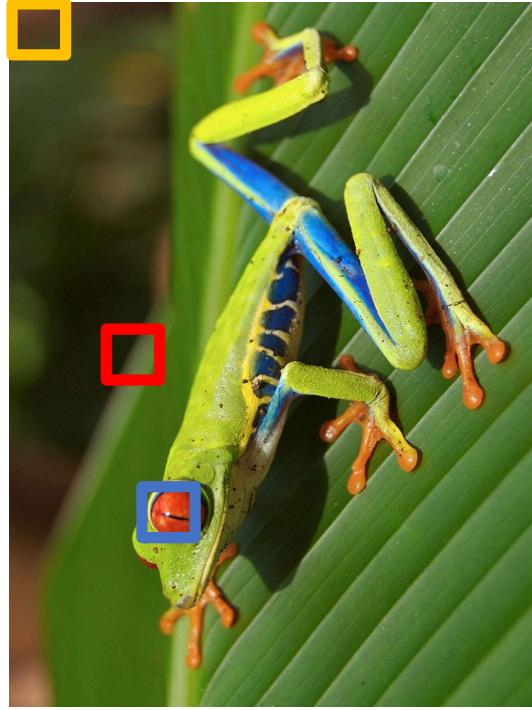
Weak edges  
Strong diagonal edges  
Edges in all directions



1. Compute edge direction / strength at each pixel
  2. Divide image into  $8 \times 8$  regions
  3. Within each region compute a histogram of edge directions weighted by edge strength
- Example: 320x240 image gets divided into  $40 \times 30$  bins; 8 directions per bin; feature vector has  $30 * 40 * 9 = 10,800$  numbers

Lowe, "Object recognition from local scale-invariant features", ICCV 1999  
Dalal and Triggs, "Histograms of oriented gradients for human detection," CVPR 2005

# 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

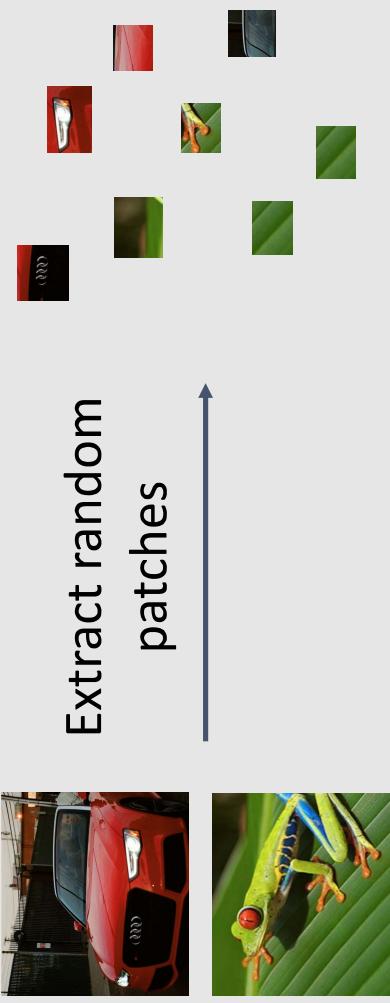
Captures texture and position, robust to small image changes

Example: 320x240 image gets divided into 40x30 bins; 8 directions per bin; feature vector has  $30 * 40 * 9 = 10,800$  numbers

Lowe, "Object recognition from local scale-invariant features", ICCV 1999  
Dalal and Triggs, "Histograms of oriented gradients for human detection," CVPR 2005

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

## Step 1: Build codebook



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

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

Car image is CC0 1.0 public domain

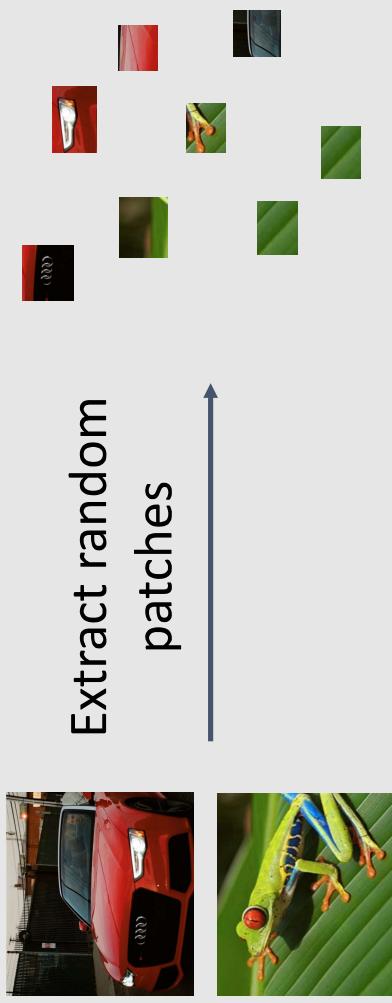
Justin Johnson

Lecture 5 - 16

September 18, 2019

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

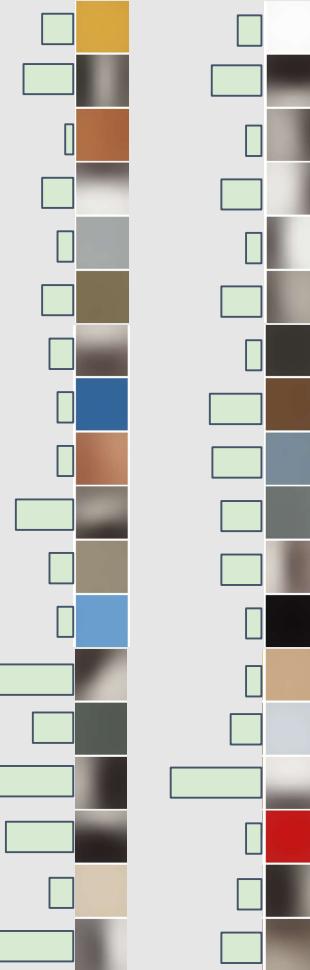
## Step 1: Build codebook



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

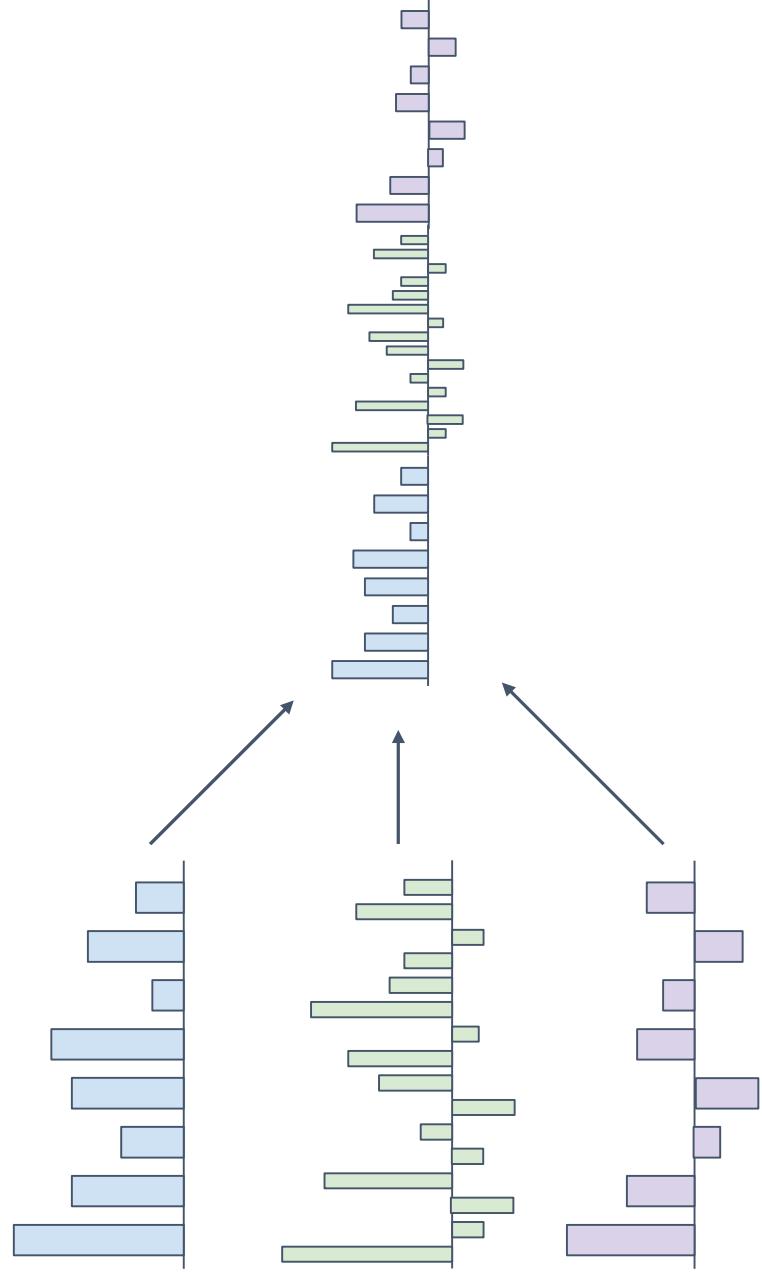
A horizontal arrow points from the text "Cluster patches to form ‘codebook’ of ‘visual words’" to the right.

## 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
- $\left. \begin{array}{l} \text{SIFT: 128-dim} \\ \text{color: 96-dim} \end{array} \right\}$  reduced to 64-dim with PCA

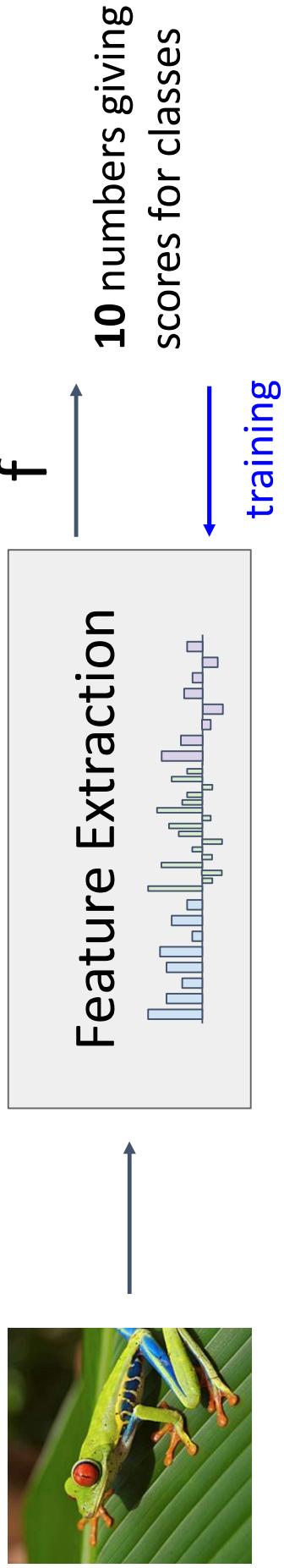
FV extraction and compression:

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

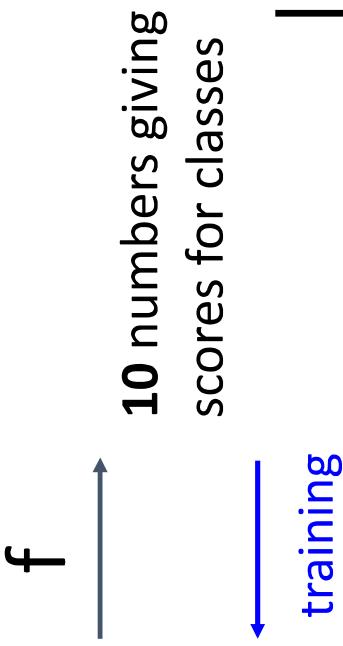
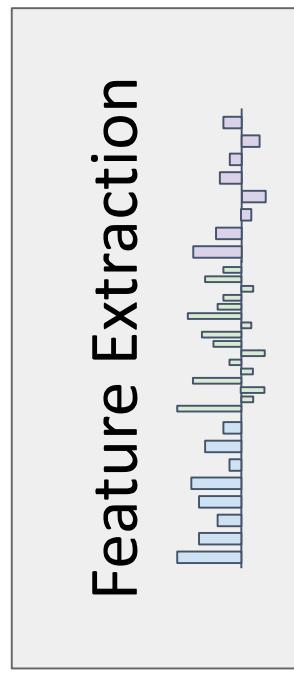
One-vs-all SVM learning with SGD

Late fusion of SIFT and color systems

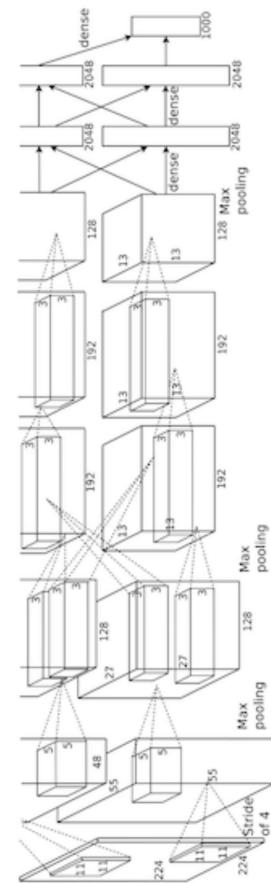
# Image Features



# Image Features vs Neural Networks



Krizhevsky, Sutskever, and Hinton, "Imagenet classification with deep convolutional neural networks", NIPS 2012.  
Figure copyright Krizhevsky, Sutskever, and Hinton, 2012.  
Reproduced with permission.



# 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:

$$\begin{aligned} f &= Wx \\ (\text{Now}) \quad &\begin{aligned} &\text{2-layer Neural Network} \qquad f = W_2 \max(0, W_1 x) \\ &\text{or 3-layer Neural Network} \qquad f = W_3 \max(0, W_2 \max(0, W_1 x)) \end{aligned} \end{aligned}$$

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

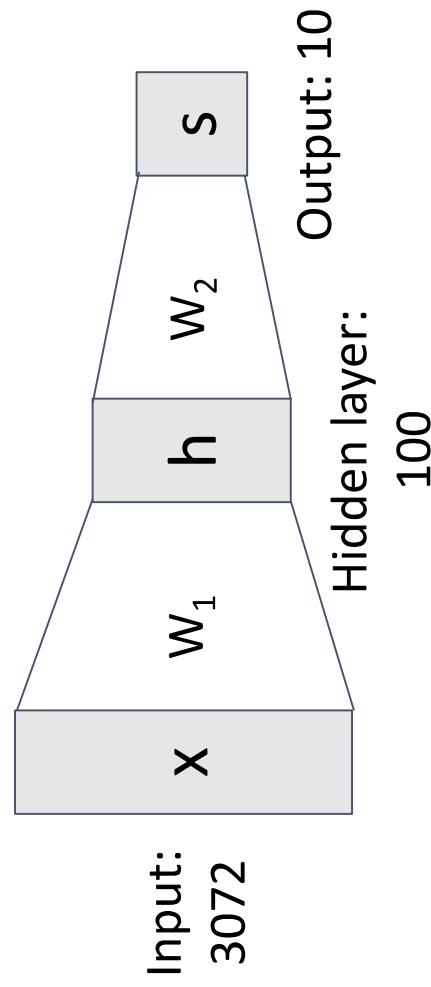
(In practice we will usually add a learnable bias at each layer as well)

# Neural Networks

(Before) Linear score function:

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

(Now) 2-layer Neural Network



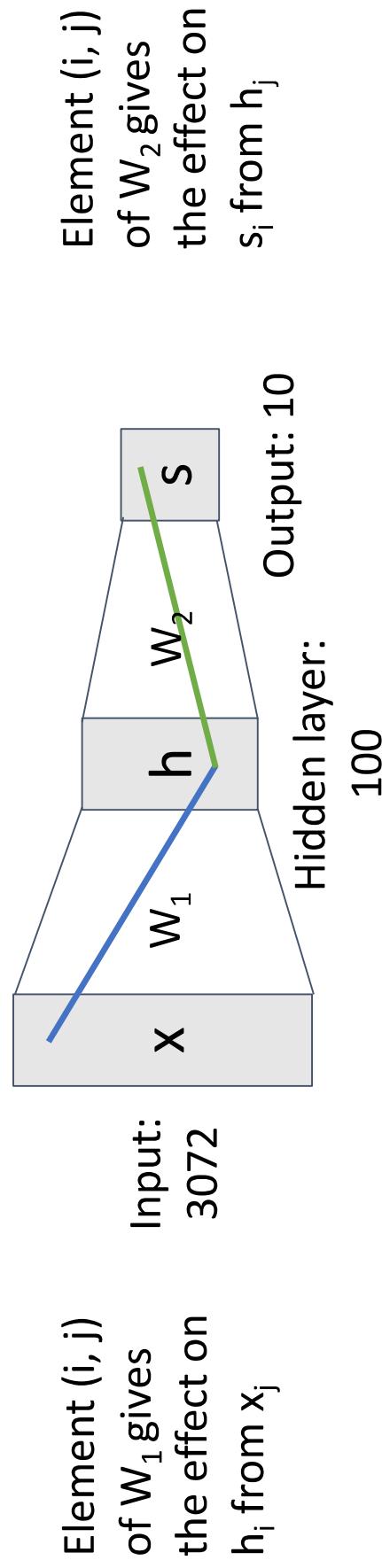
$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$

# Neural Networks

(Before) Linear score function:

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

(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

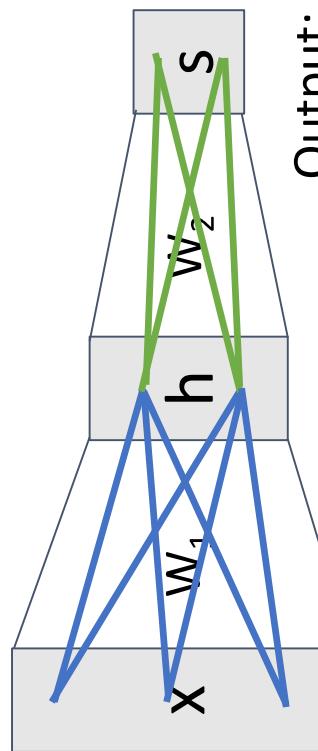
(Before) Linear score function:

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

(Now) 2-layer Neural Network

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

Input:  
3072



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

Output: 10

Hidden layer:  
100

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

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

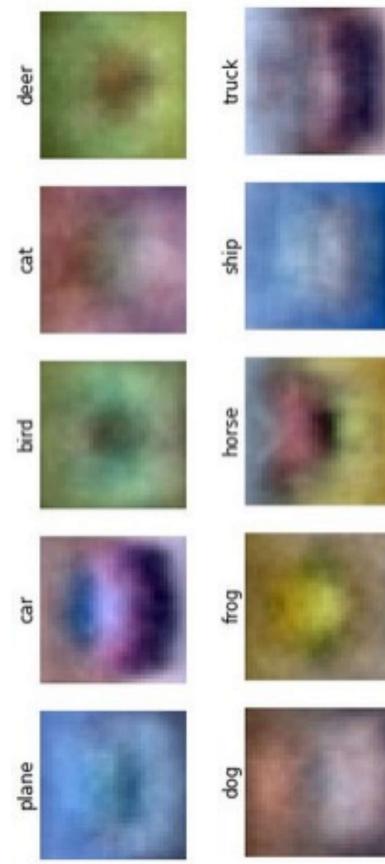
Fully-connected neural network

Also “Multi-Layer Perceptron” (MLP)

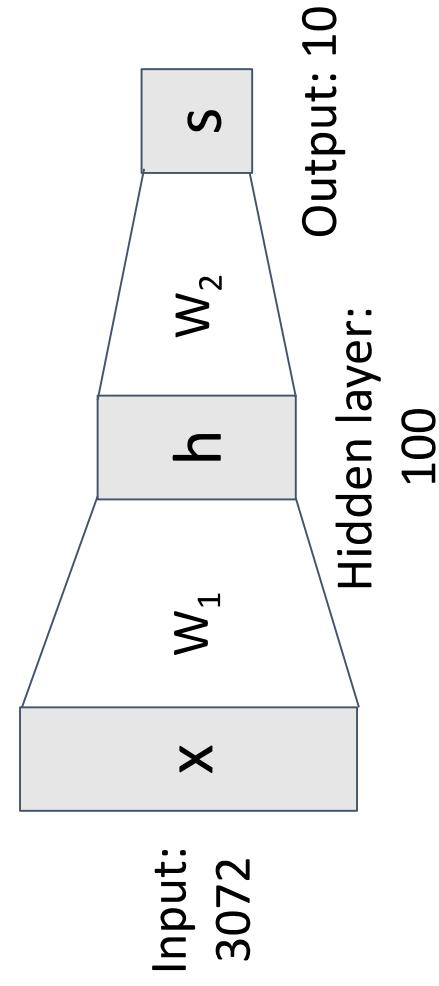
# Neural Networks

(Before) Linear score function:

Linear classifier: One template per class



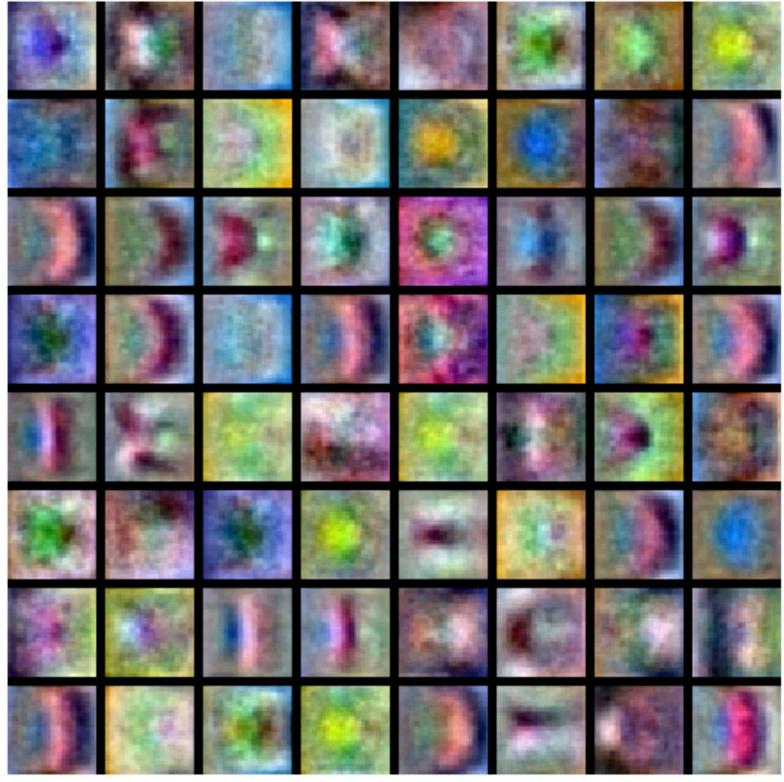
(Now) 2-layer Neural Network



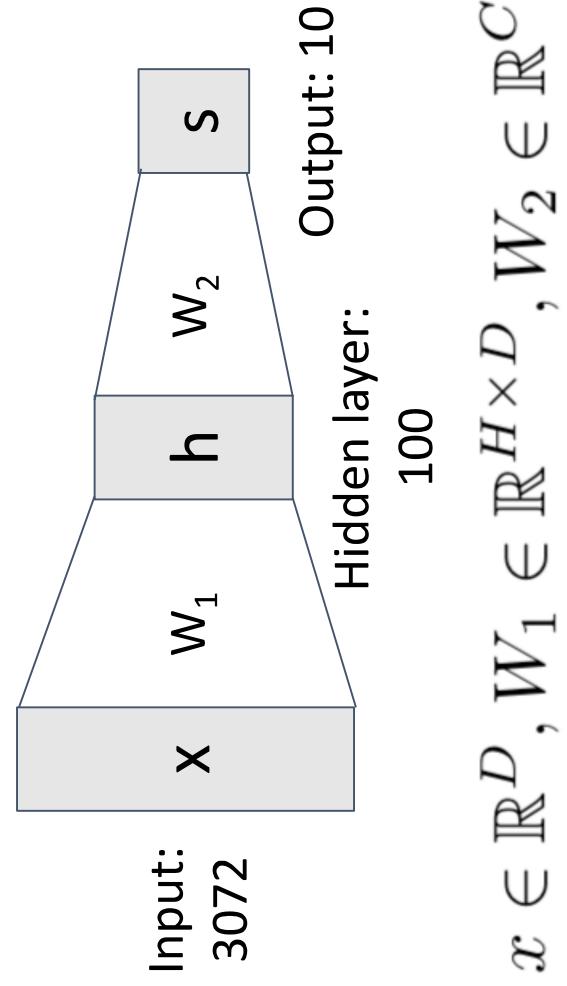
$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$

# Neural Networks

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



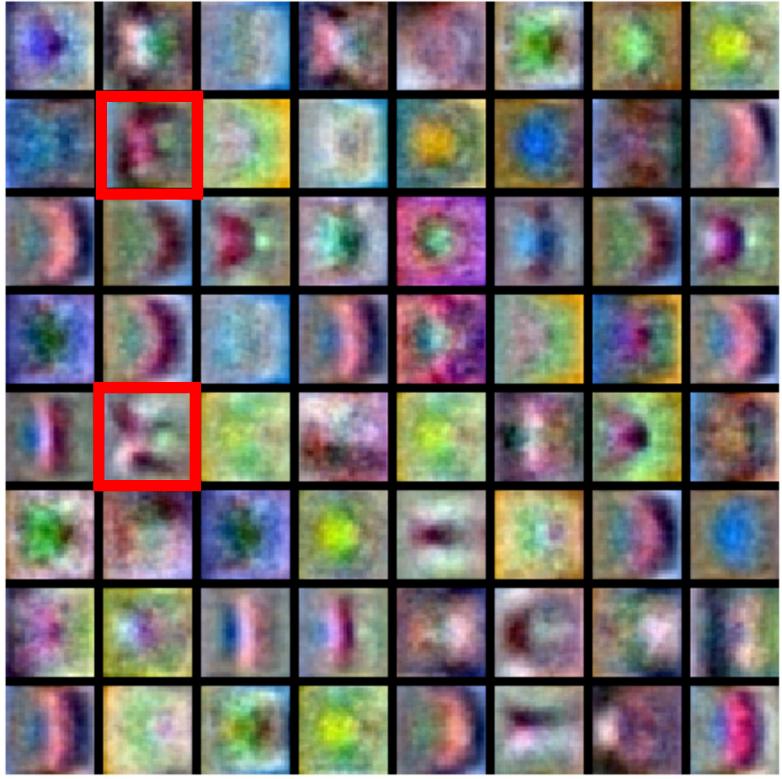
(Before) Linear score function:  
(Now) 2-layer Neural Network



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

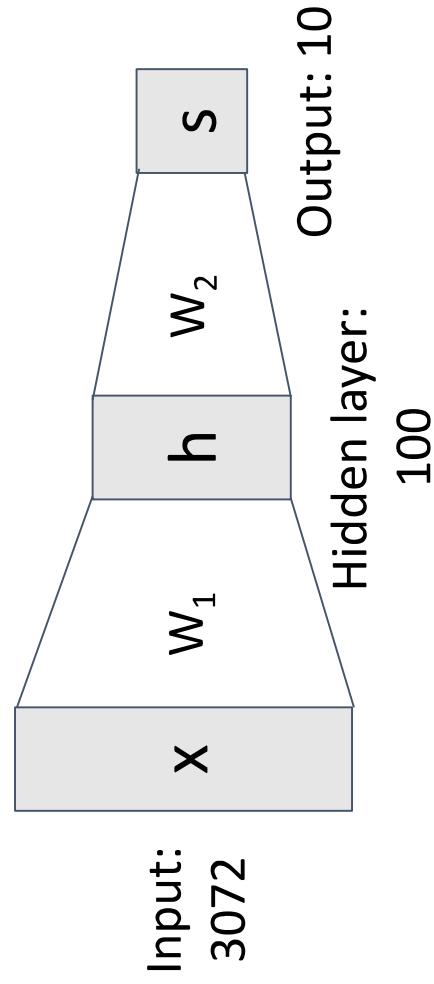
# Neural Networks

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



(Before) Linear score function:

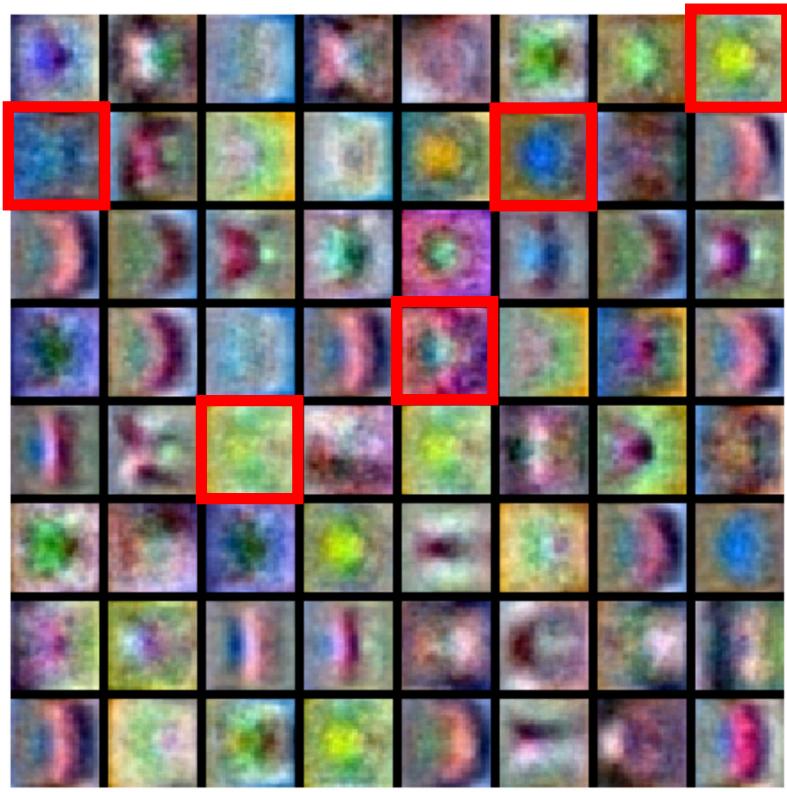
(Now) 2-layer Neural Network



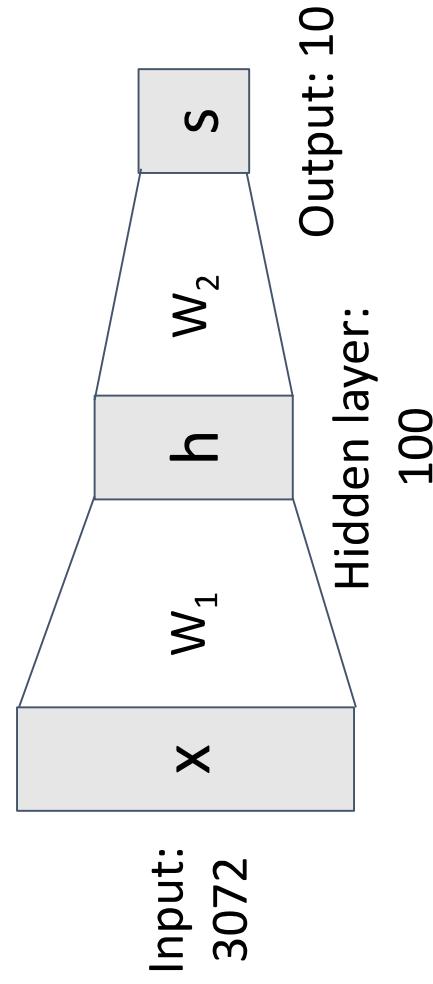
$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$

# Neural Networks

“Distributed representation”:  
Most templates not interpretable!



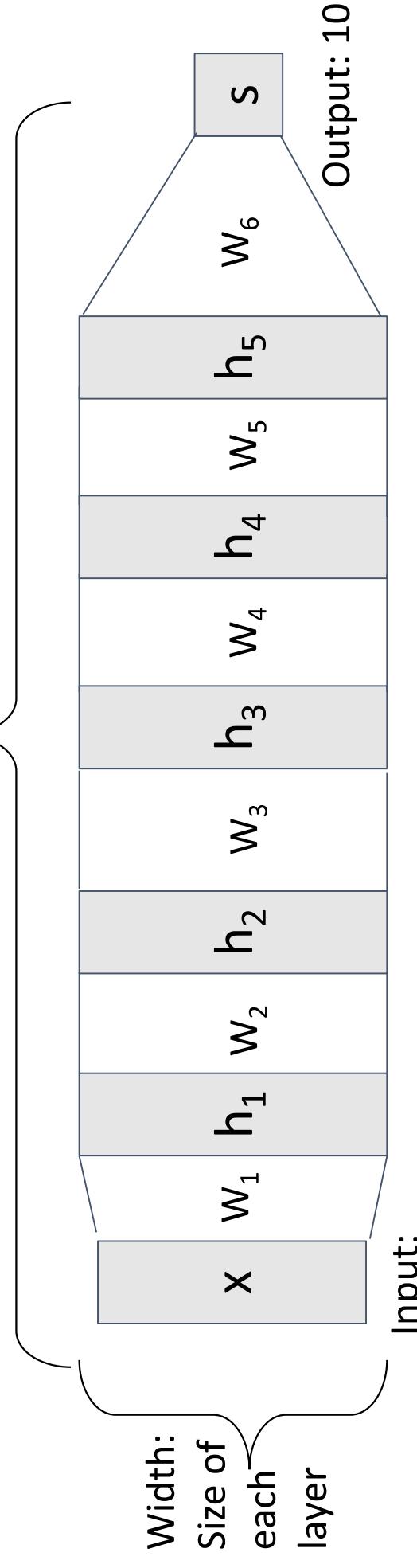
(Before) Linear score function:  
(Now) 2-layer Neural Network



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

# Deep Neural Networks

Depth = number of layers



Input:  
3072

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

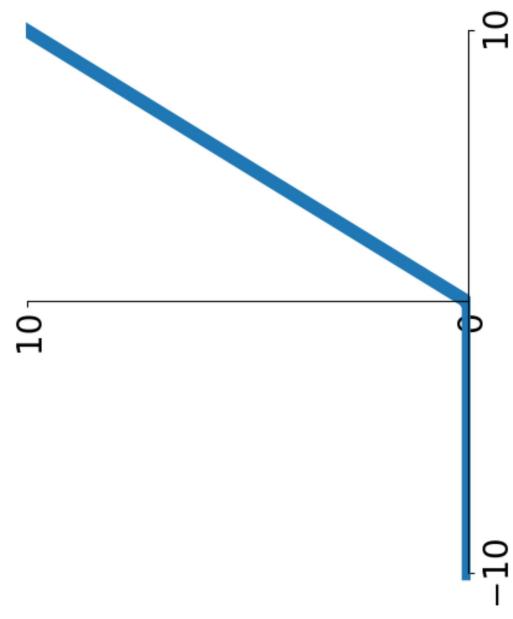
## Activation Functions

2-layer Neural Network

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

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

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



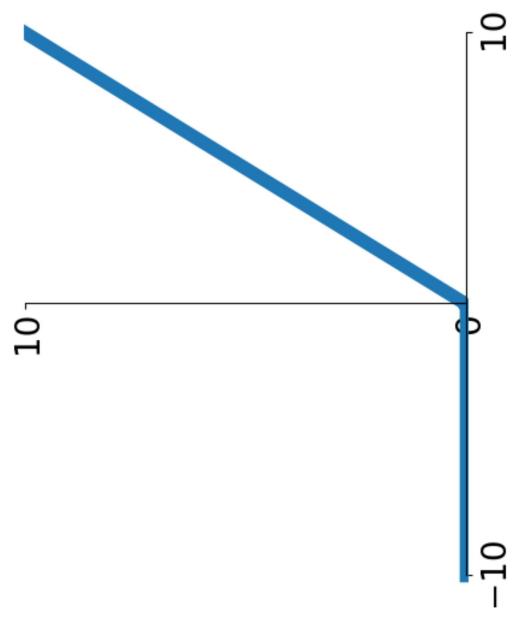
## Activation Functions

2-layer Neural Network

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

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

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



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

$$s = W_2 W_1 x$$

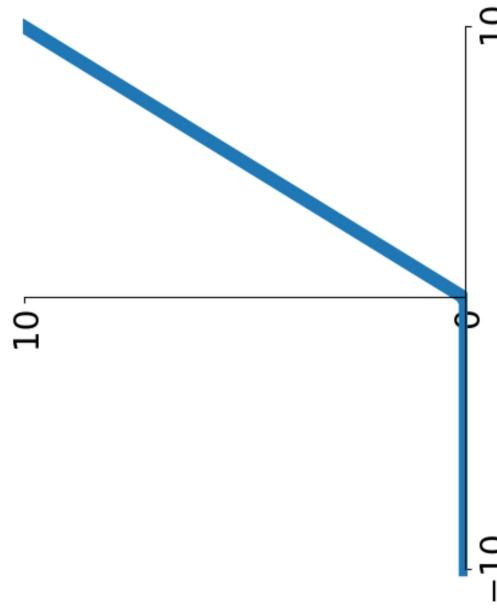
## Activation Functions

2-layer Neural Network

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

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

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



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

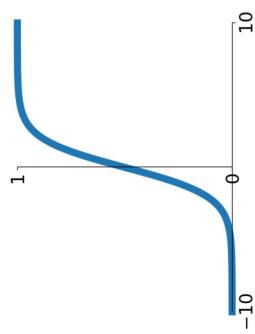
$$\begin{aligned} s &= W_2 W_1 x \\ W_3 &= W_2 W_1 \in \mathbb{R}^{C \times H} \quad s = W_3 x \end{aligned}$$

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

# Activation Functions

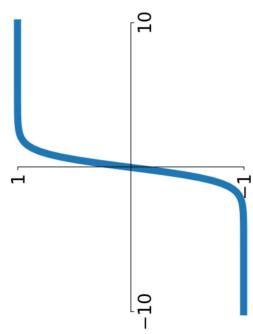
**Sigmoid**

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



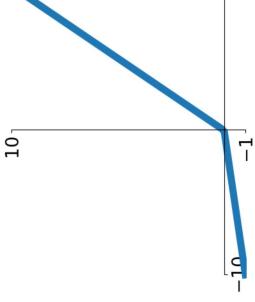
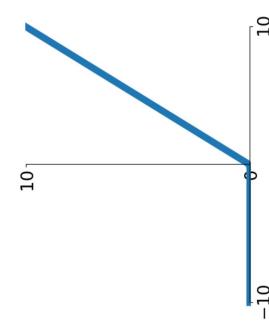
**tanh**

$$\tanh(x)$$



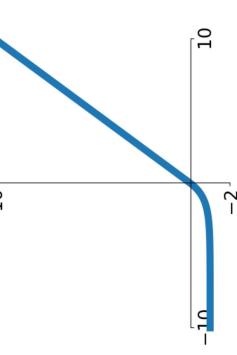
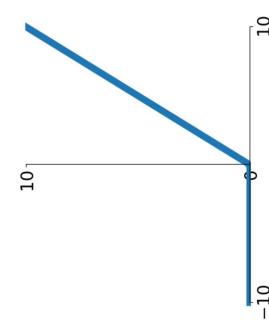
**ReLU**

$$\max(0, x)$$



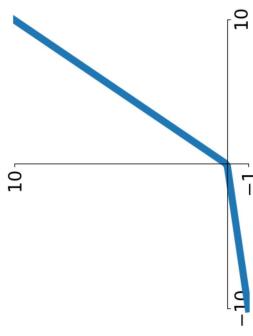
**Maxout**

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

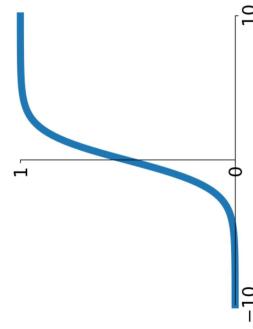


# Activation Functions

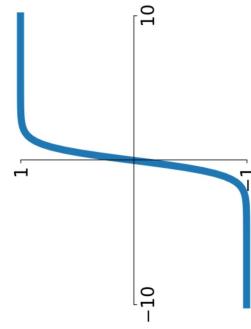
ReLU is a good default choice  
for most problems



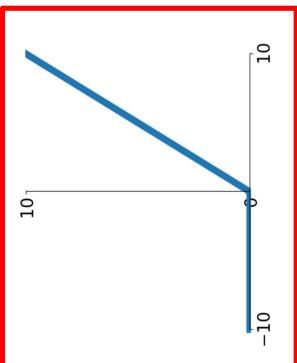
**Leaky ReLU**  
 $\max(0.1x, x)$



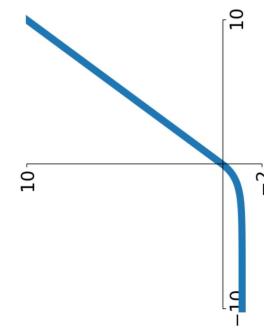
**Sigmoid**  
 $\sigma(x) = \frac{1}{1+e^{-x}}$



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



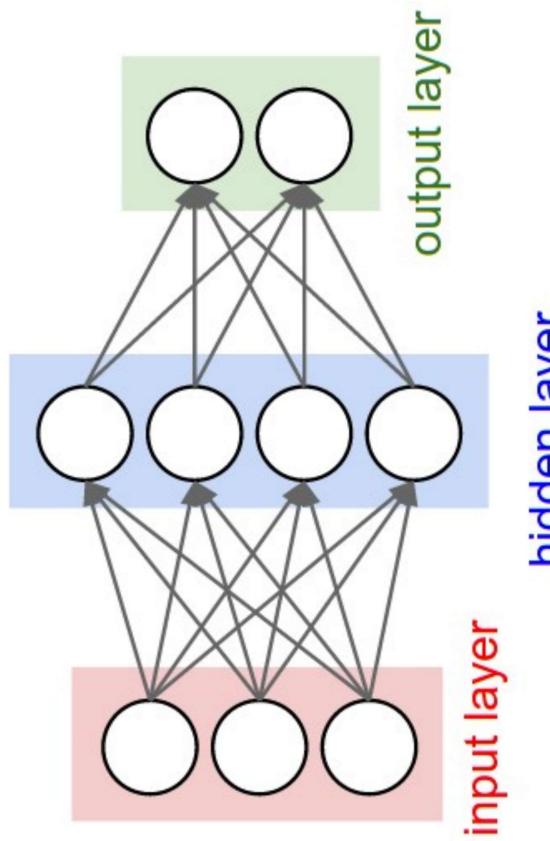
**ReLU**  
 $\max(0, x)$



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

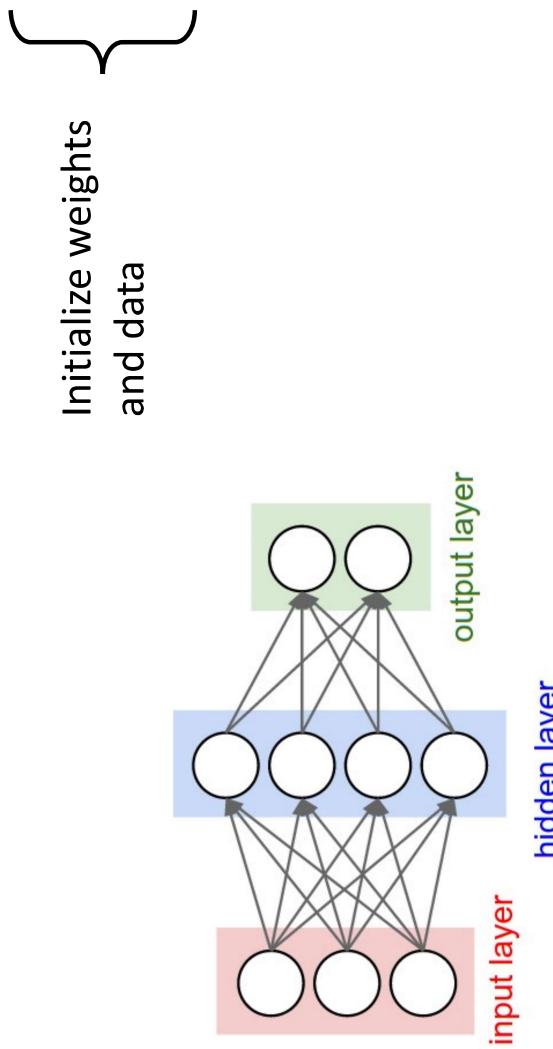
# Neural Net in <20 lines!

```
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
```



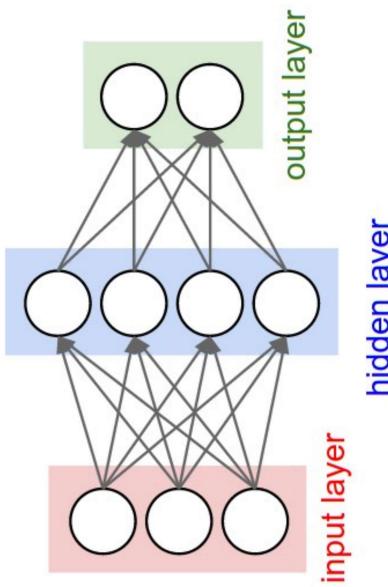
# Neural Net in <20 lines!

```
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
```



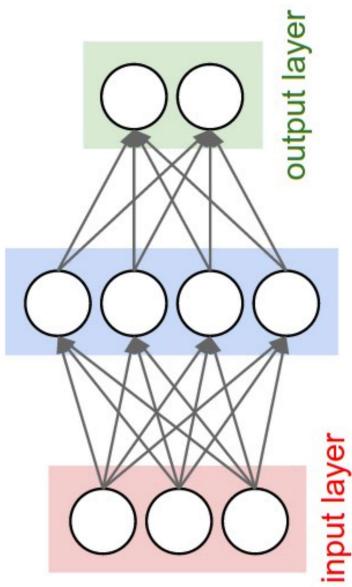
# Neural Net in <20 lines!

```
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
```



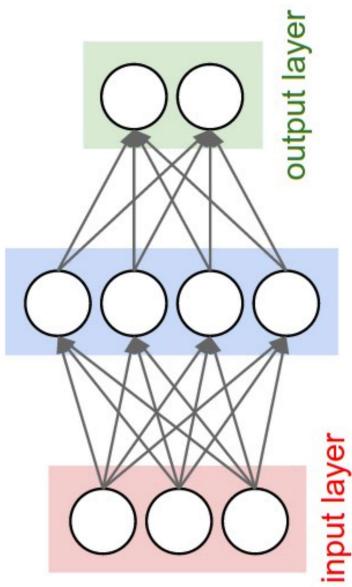
# Neural Net in <20 lines!

```
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    dh = dy_pred.dot(w2.T)
13    dw1 = x.T.dot(dh * h * (1 - h))
14    dw2 = h.T.dot(dy_pred)
15    w1 -= 1e-4 * dw1
16    w2 -= 1e-4 * dw2
```



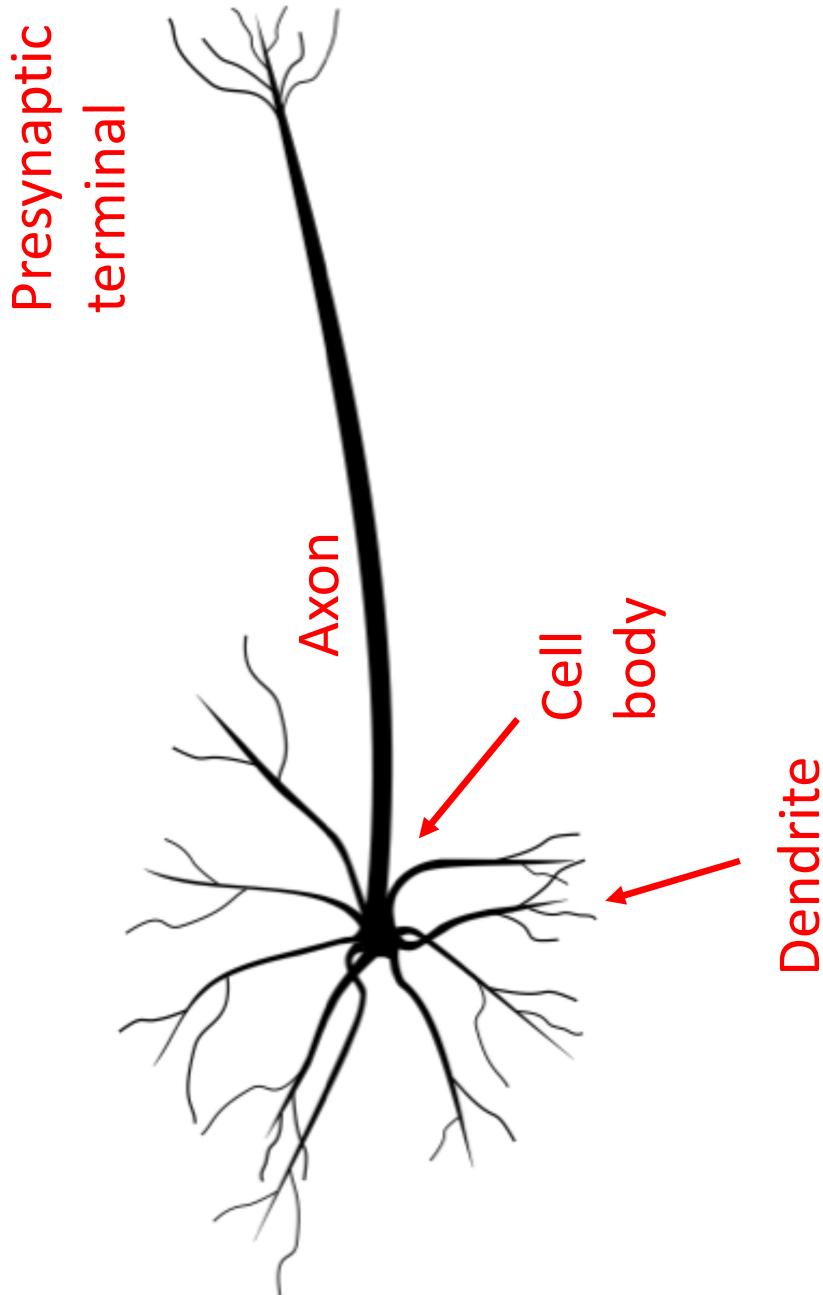
# Neural Net in <20 lines!

```
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    dh = dy_pred.dot(w2.T)
13    dw2 = h.T.dot(dy_pred)
14    dw1 = x.T.dot(dh * h * (1 - h))
15    w1 -= 1e-4 * dw1
16    w2 -= 1e-4 * dw2
```

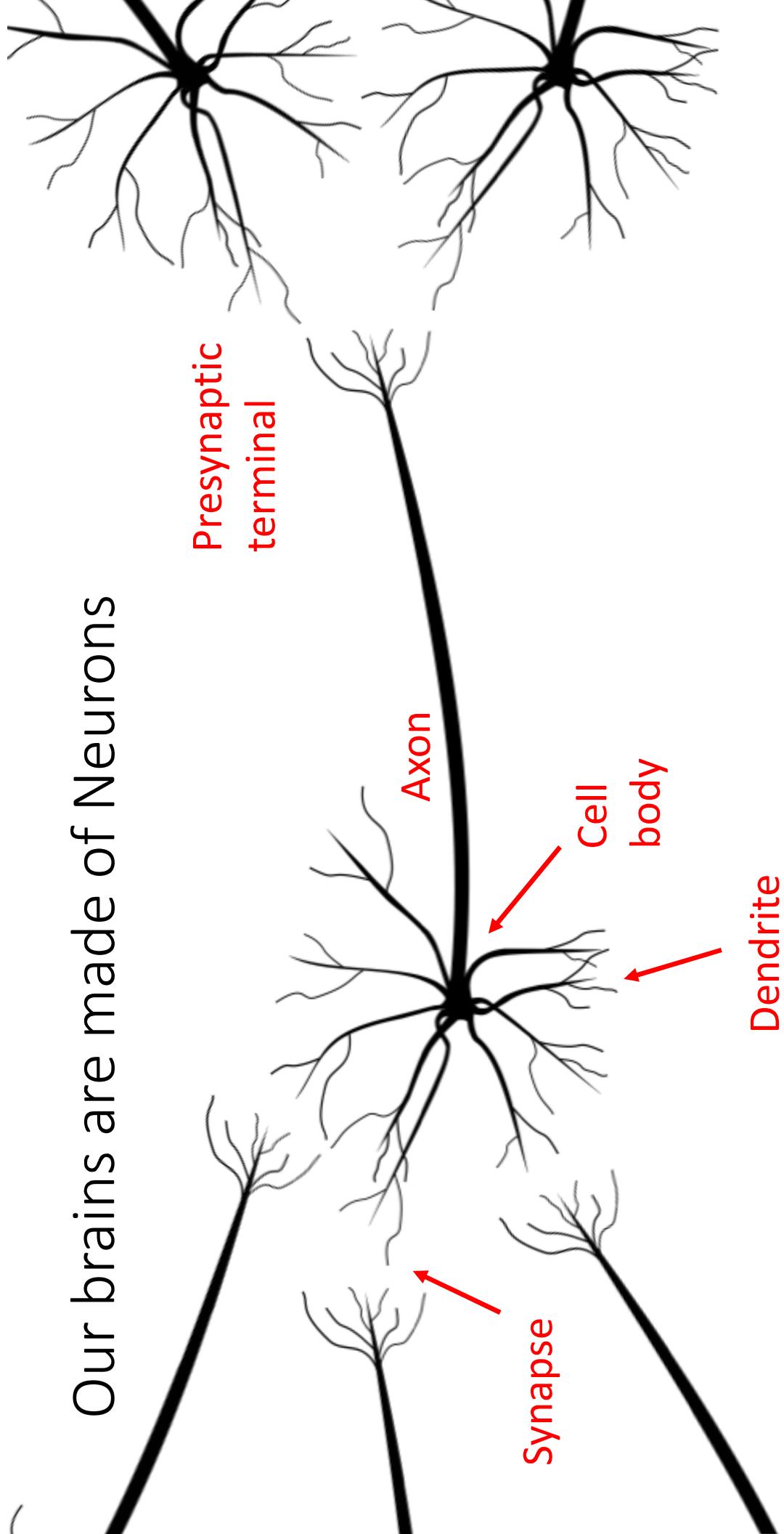




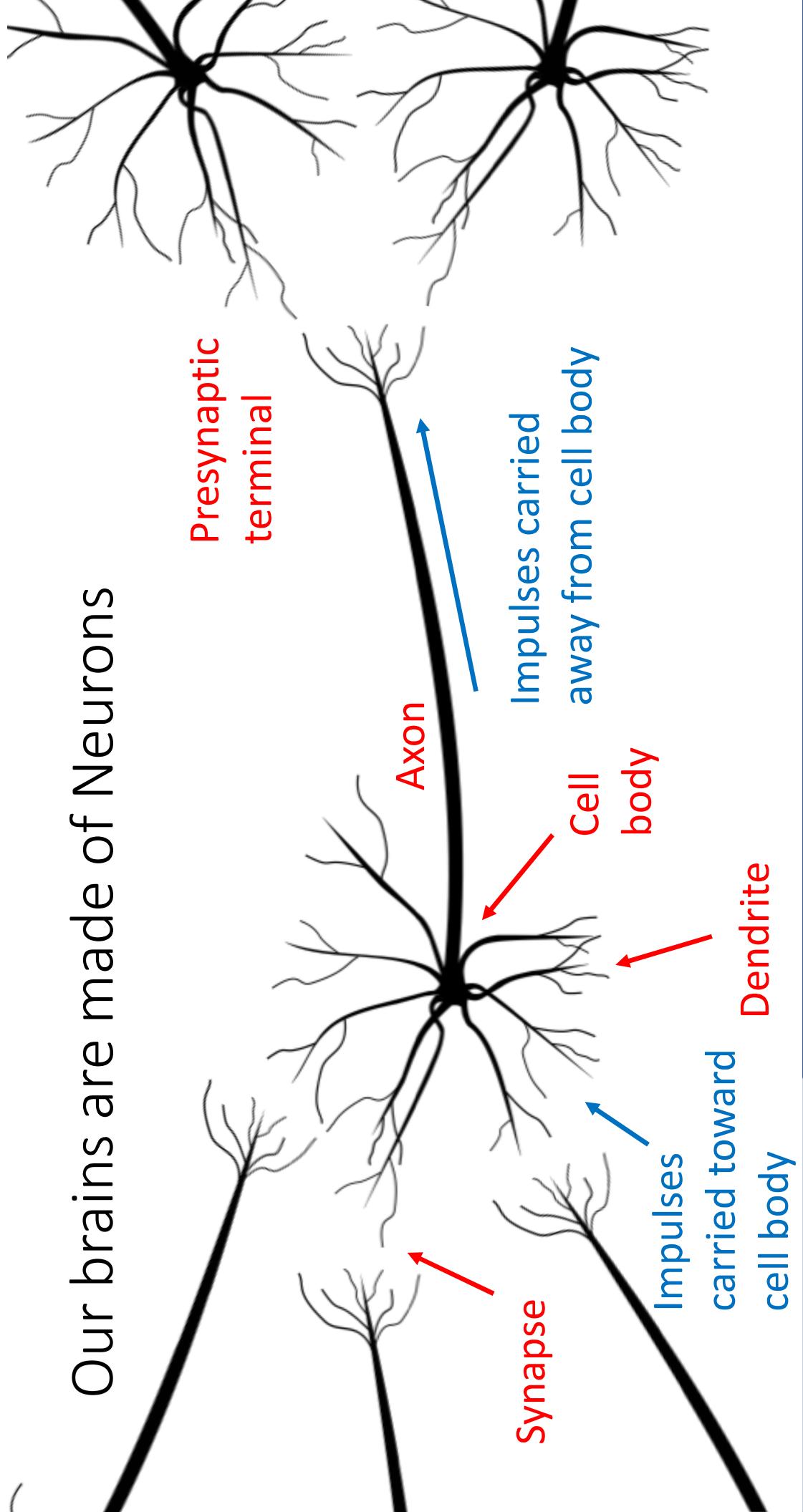
Our brains are made of Neurons



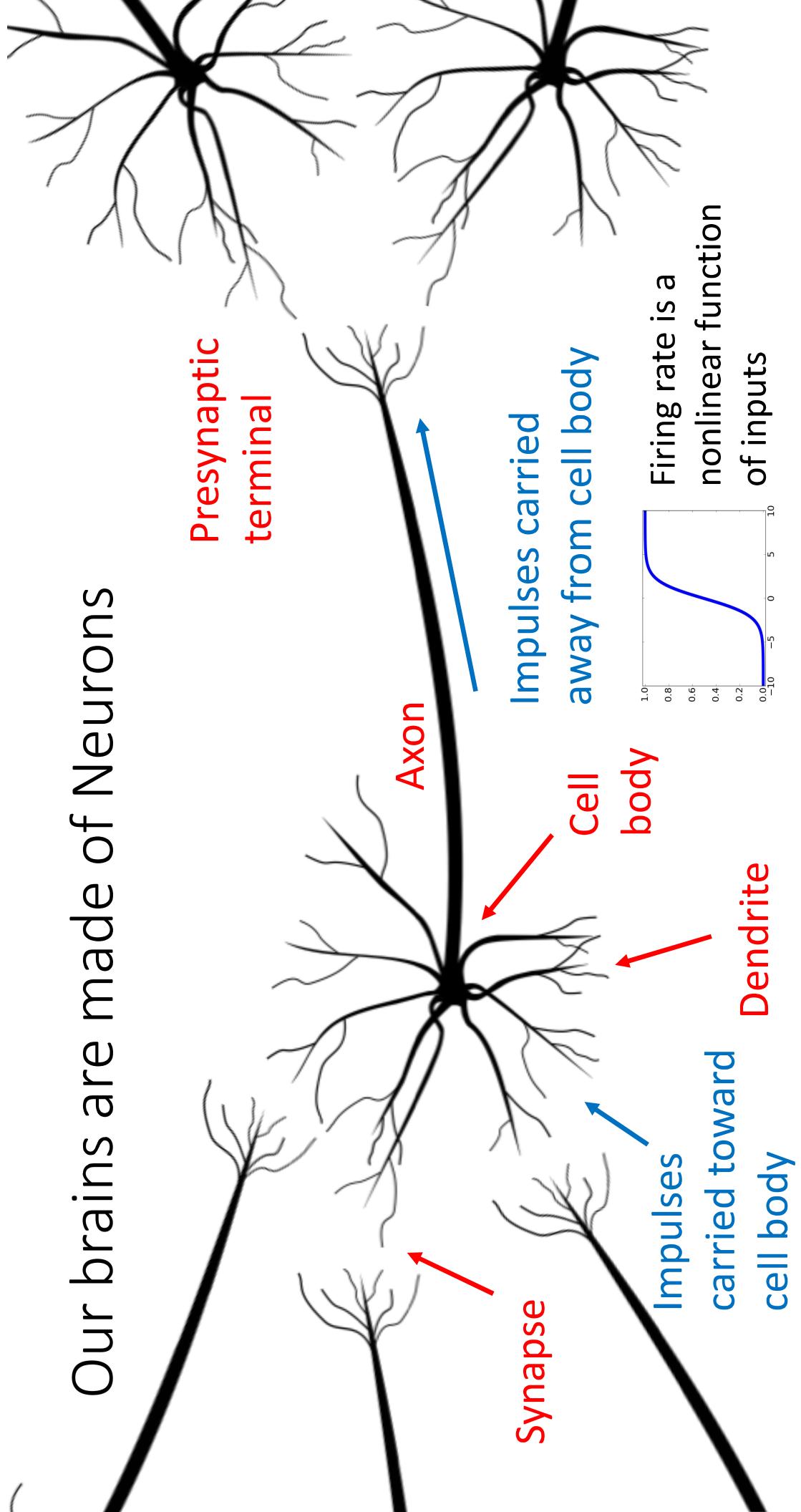
Our brains are made of Neurons



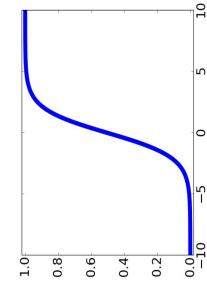
Our brains are made of Neurons

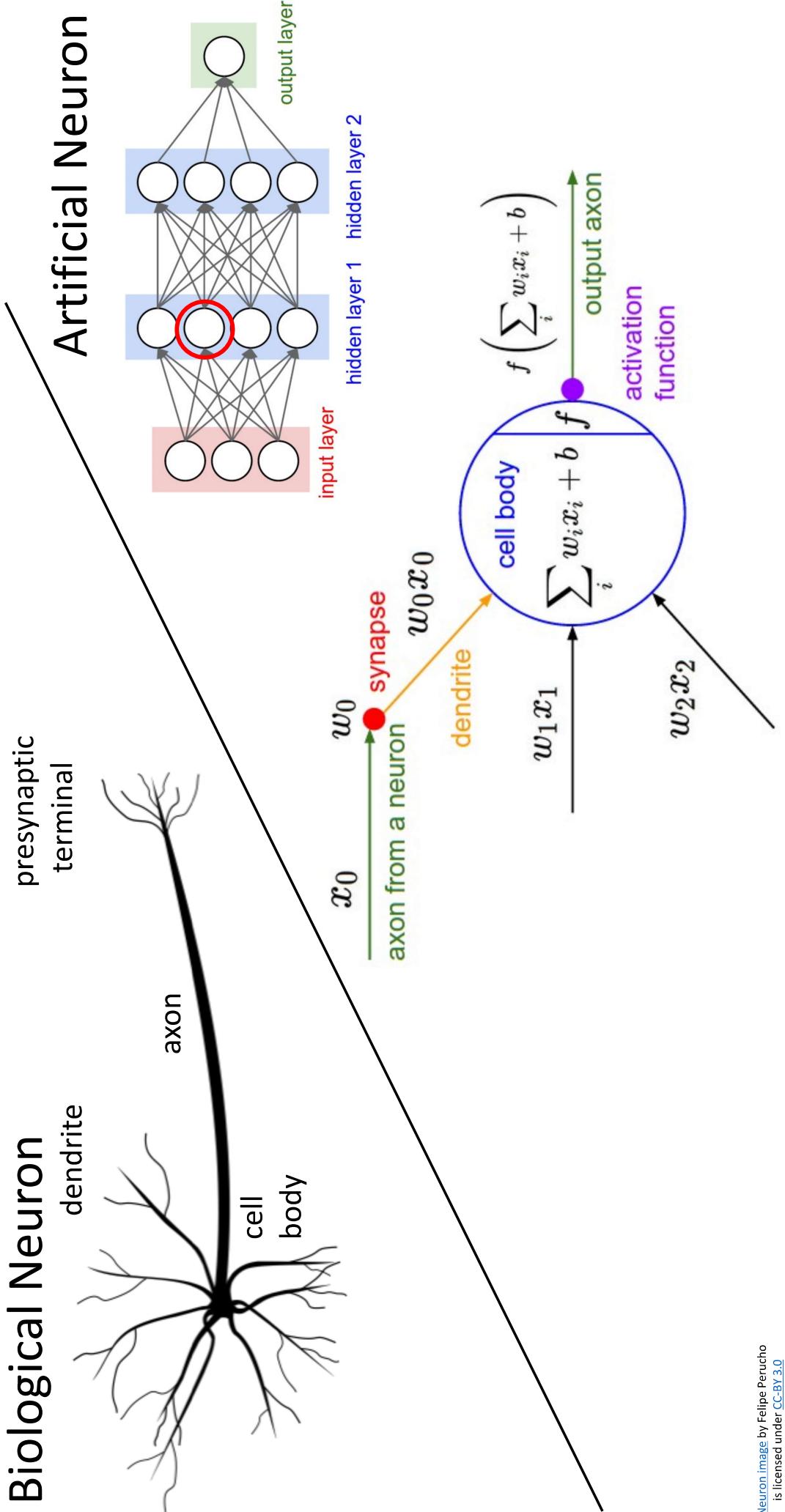


Our brains are made of Neurons

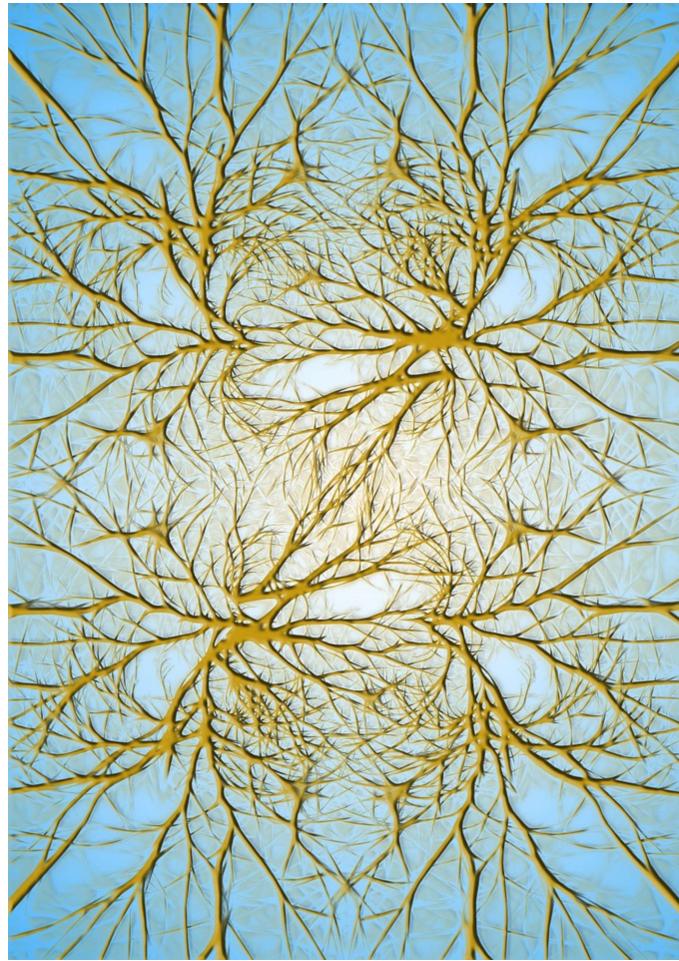


Firing rate is a  
nonlinear function  
of inputs



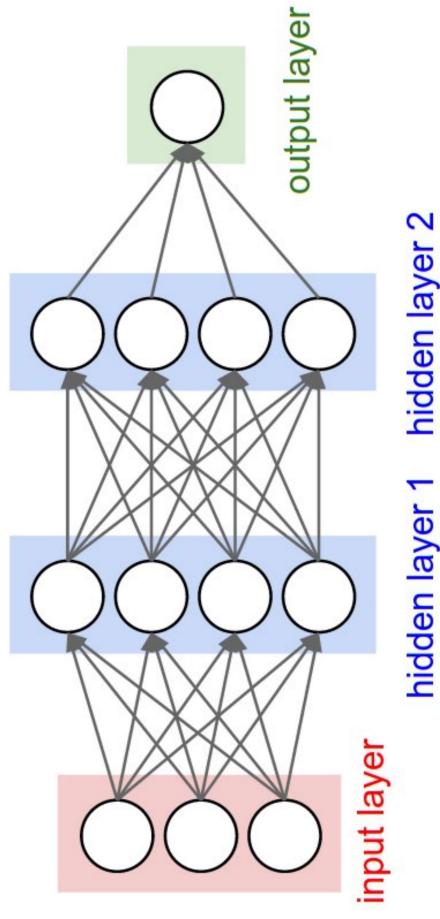


## Biological Neurons: Complex connectivity patterns

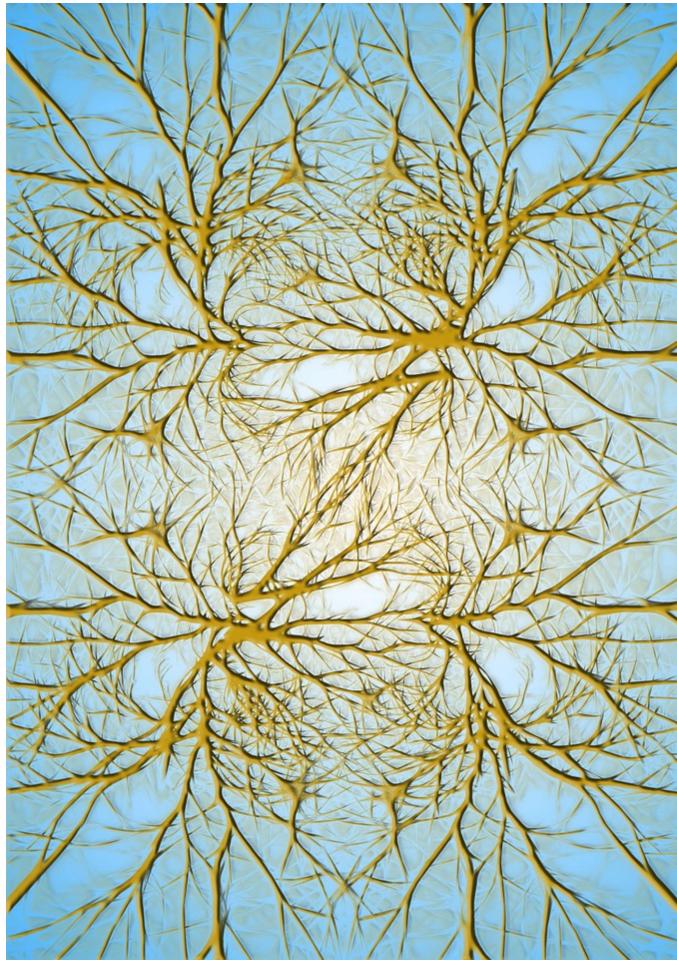


This image is CC0 Public Domain

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

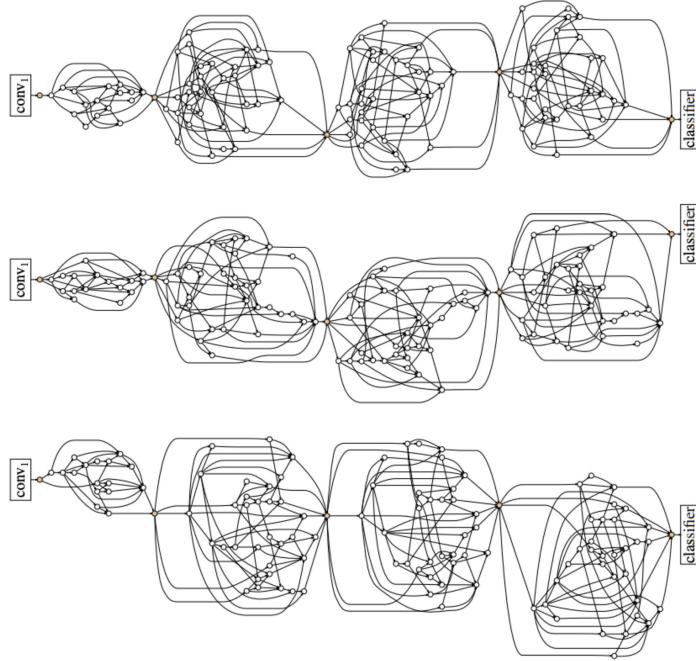


Biological Neurons:  
Complex connectivity patterns



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

[Dendritic Computation. London and Häusser]

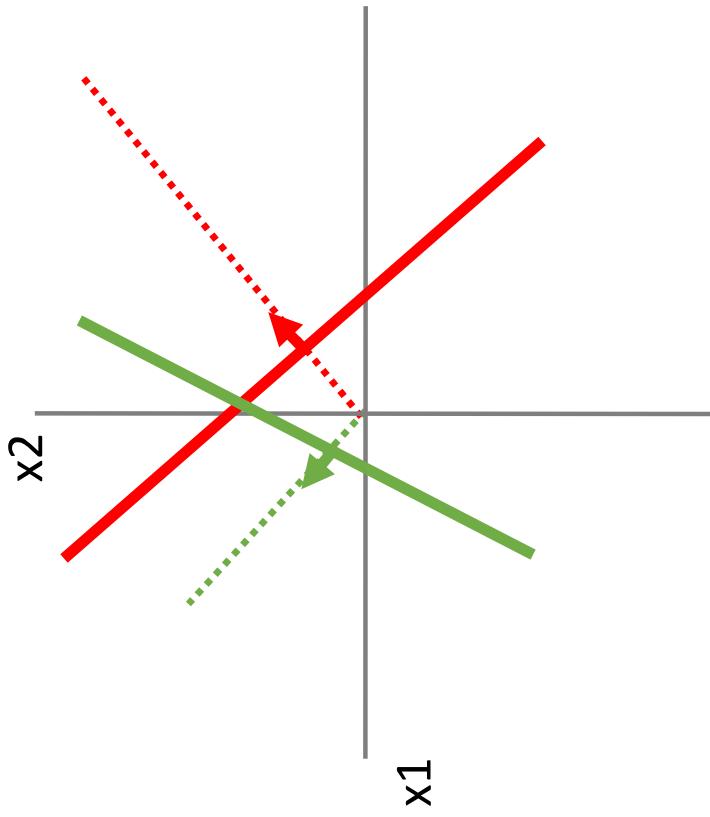
Justin Johnson

Lecture 5 - 51

September 18, 2019

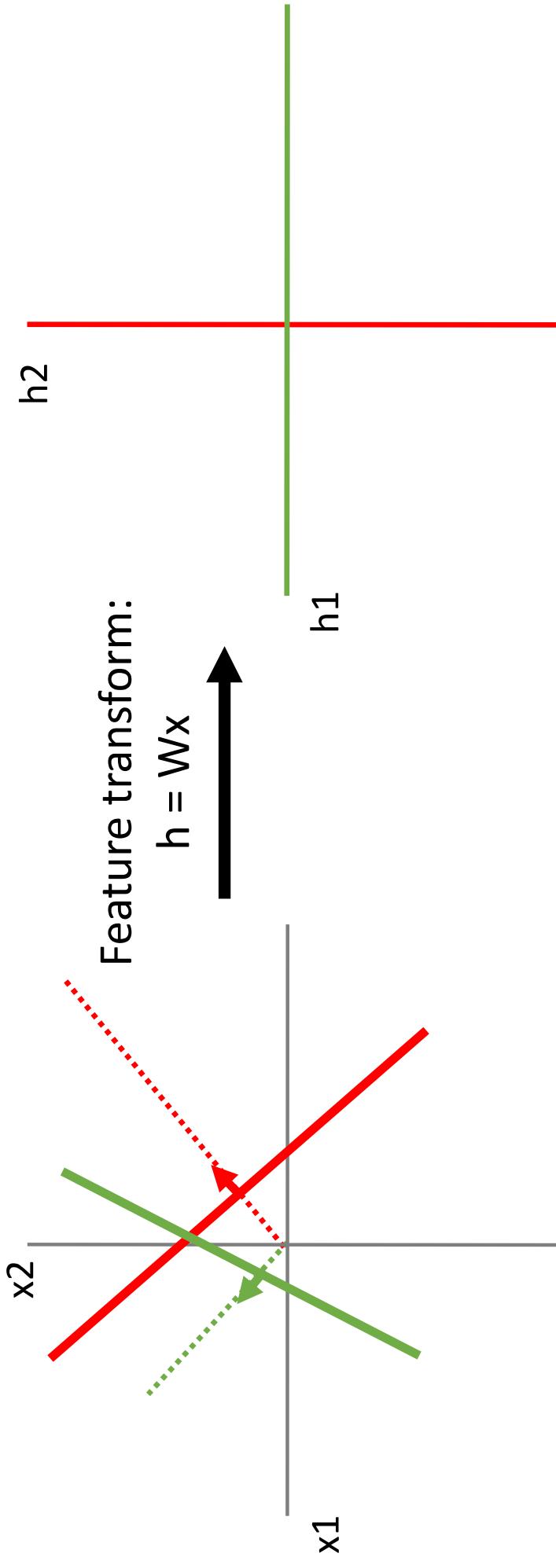
# Space Warping

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



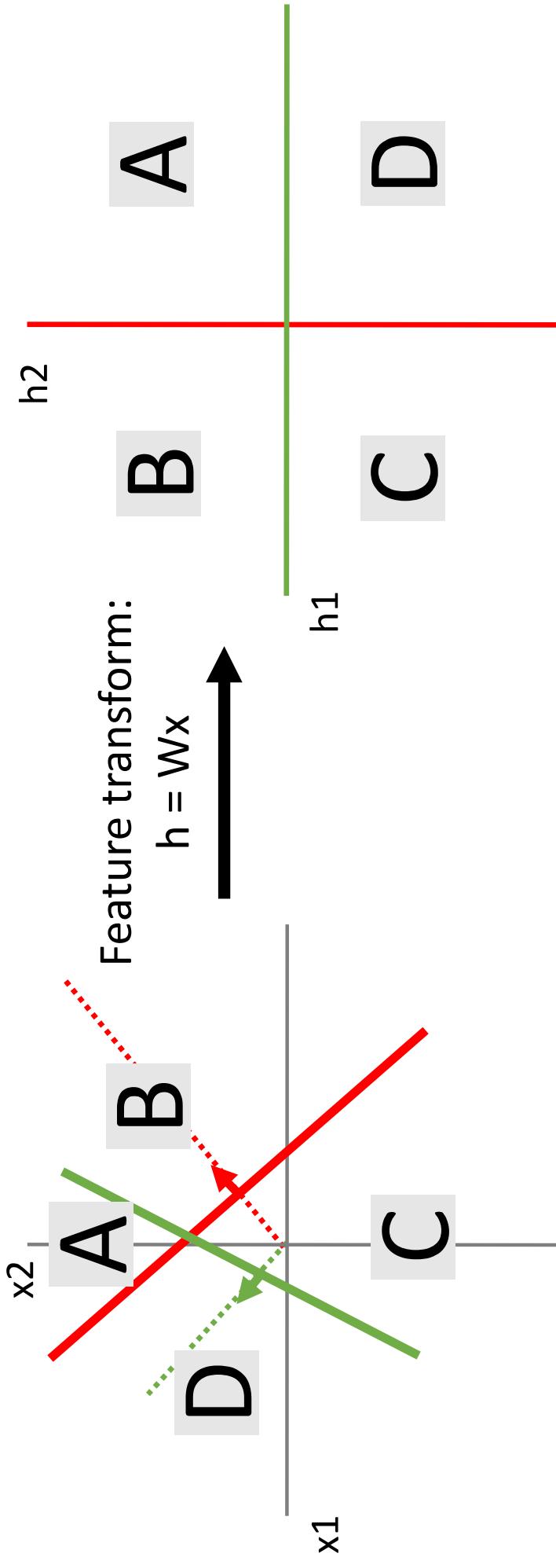
# Space Warping

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



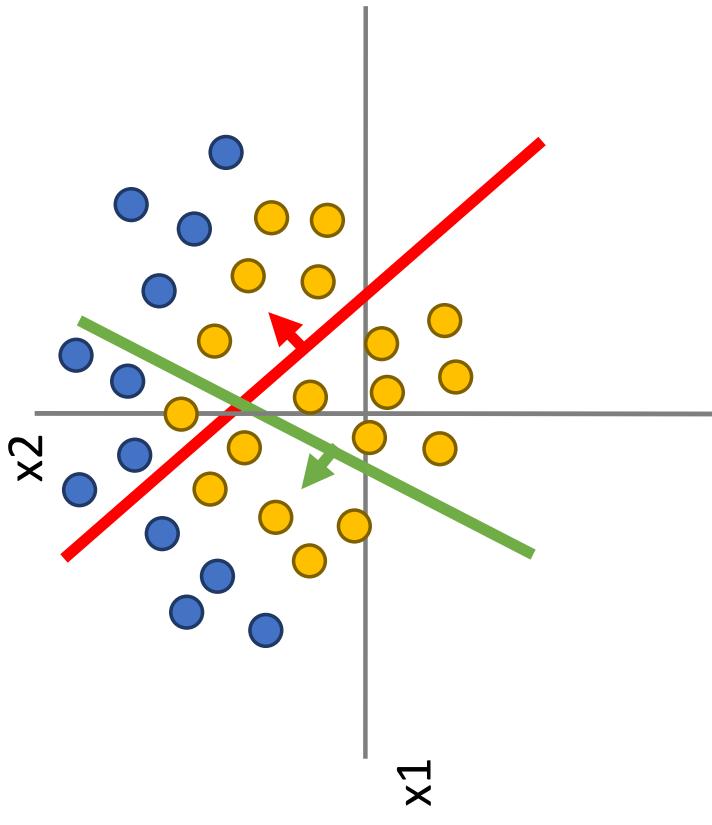
# Space Warping

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



# Space Warping

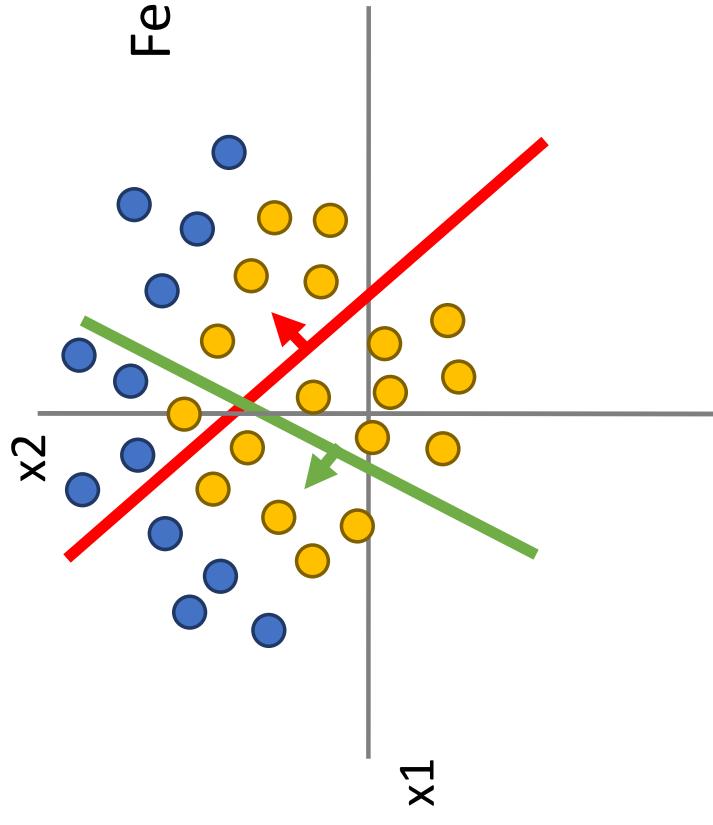
Points not linearly  
separable in original space



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

# Space Warping

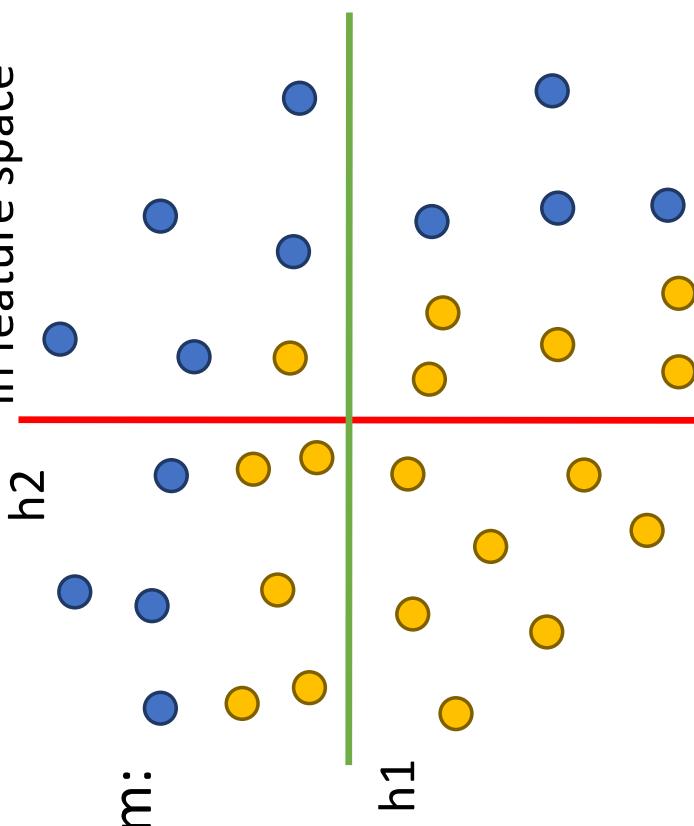
Points not linearly  
separable in original space



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

Not linearly separable  
in feature space

Feature transform:  
 $h = Wx$

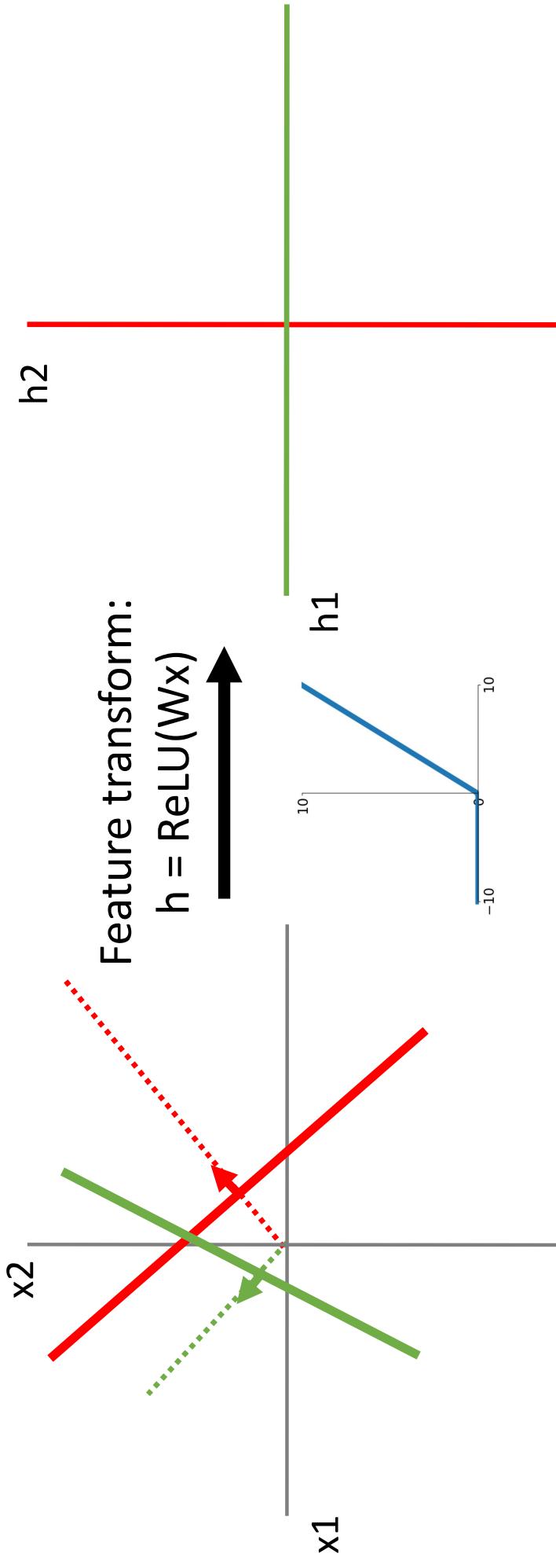


# Space Warping

Consider a neural net hidden layer:

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

Where  $x, h$  are both 2-dimensional

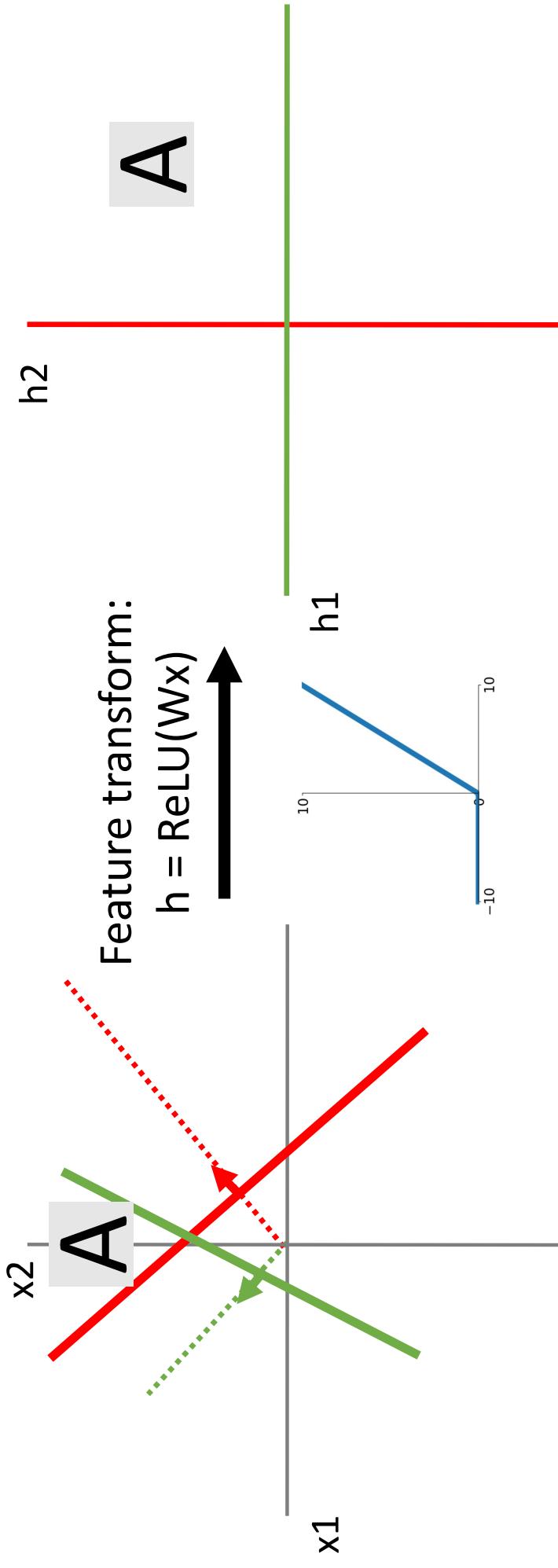


# Space Warping

Consider a neural net hidden layer:

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

Where  $x, h$  are both 2-dimensional

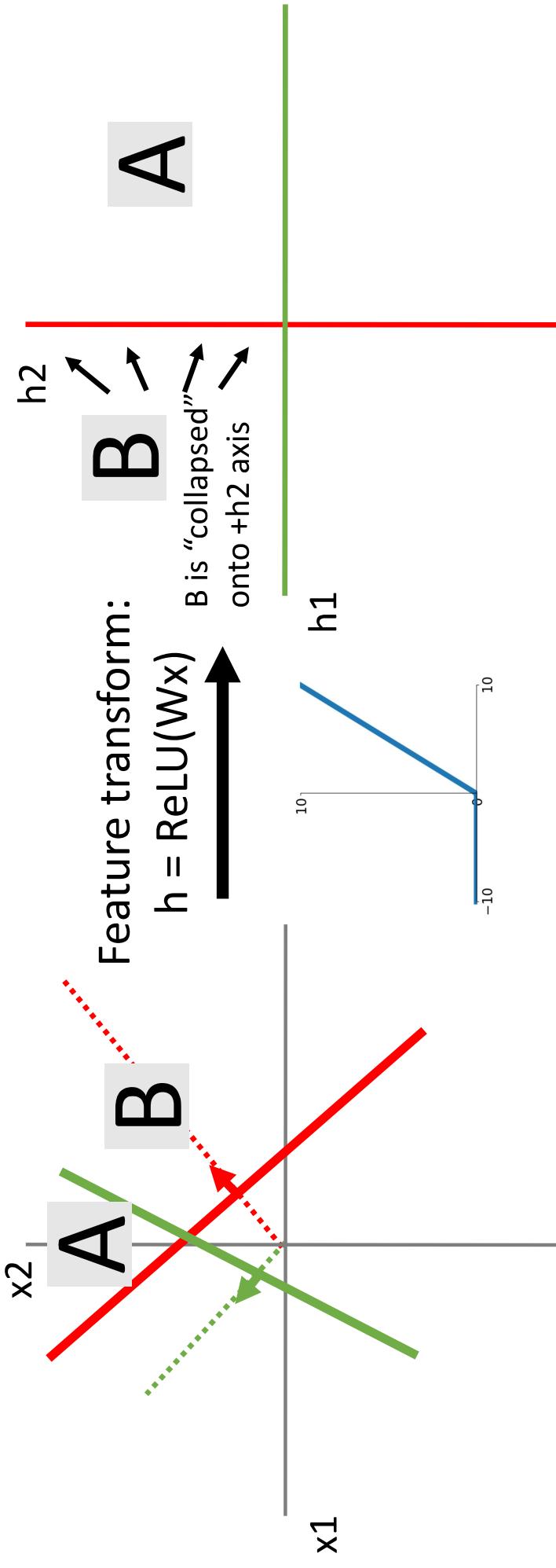


# Space Warping

Consider a neural net hidden layer:

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

Where  $x, h$  are both 2-dimensional

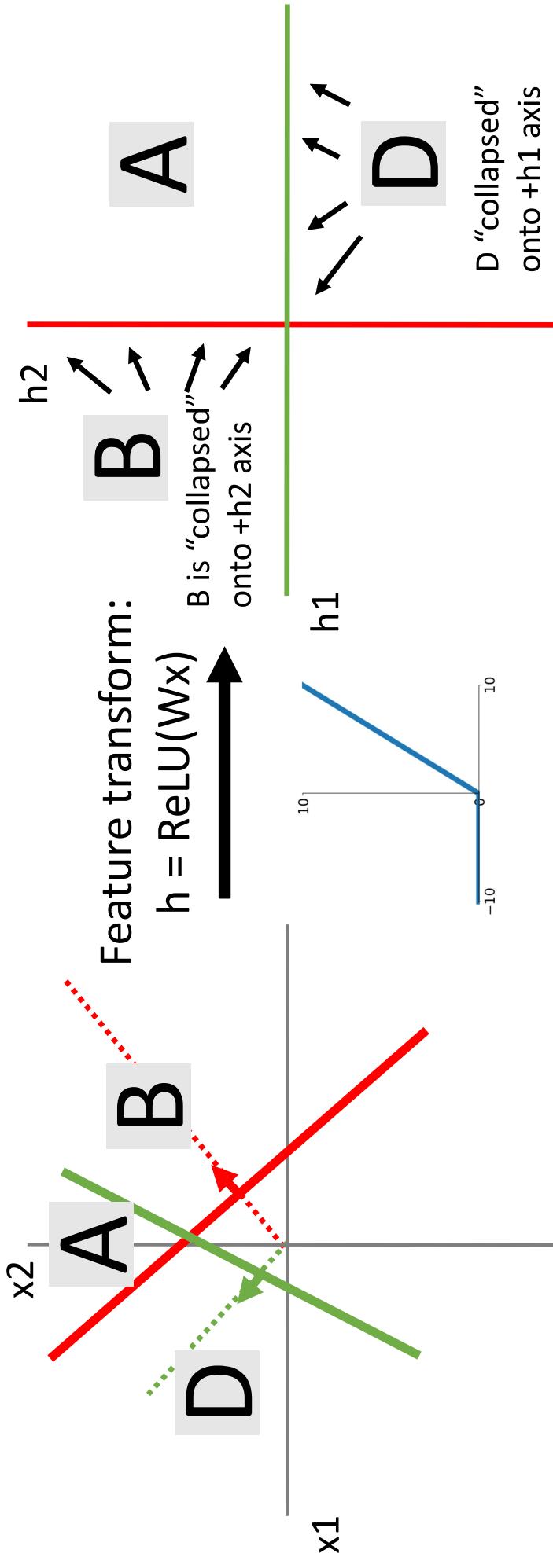


# Space Warping

Consider a neural net hidden layer:

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

Where  $x, h$  are both 2-dimensional

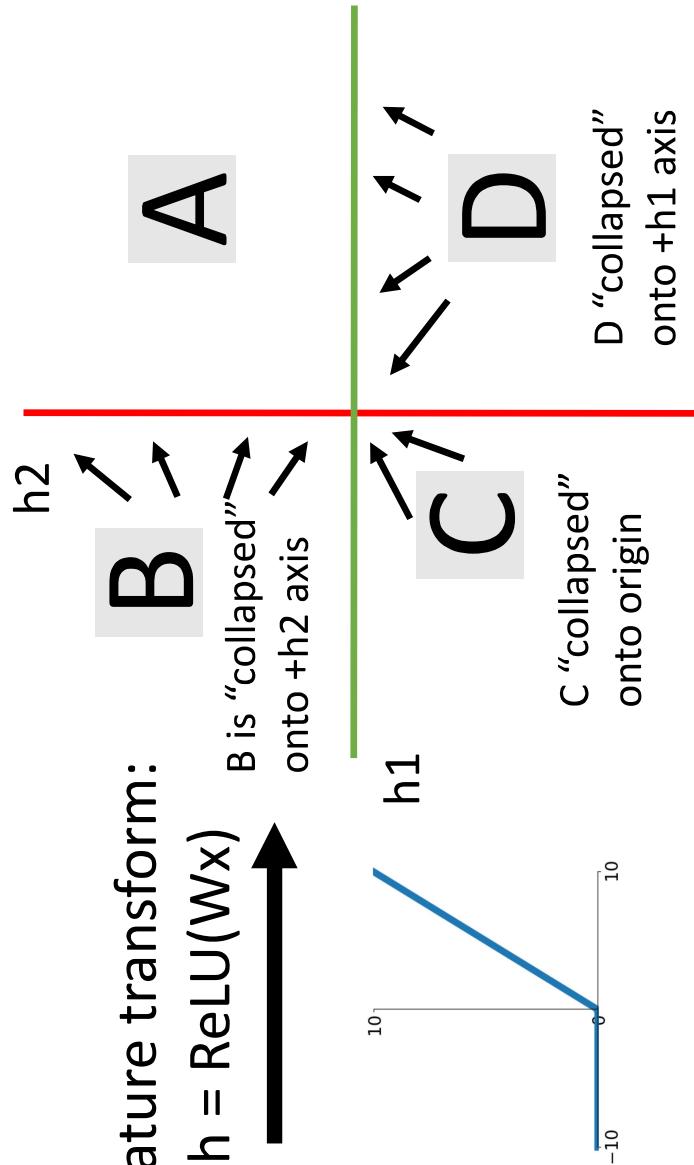
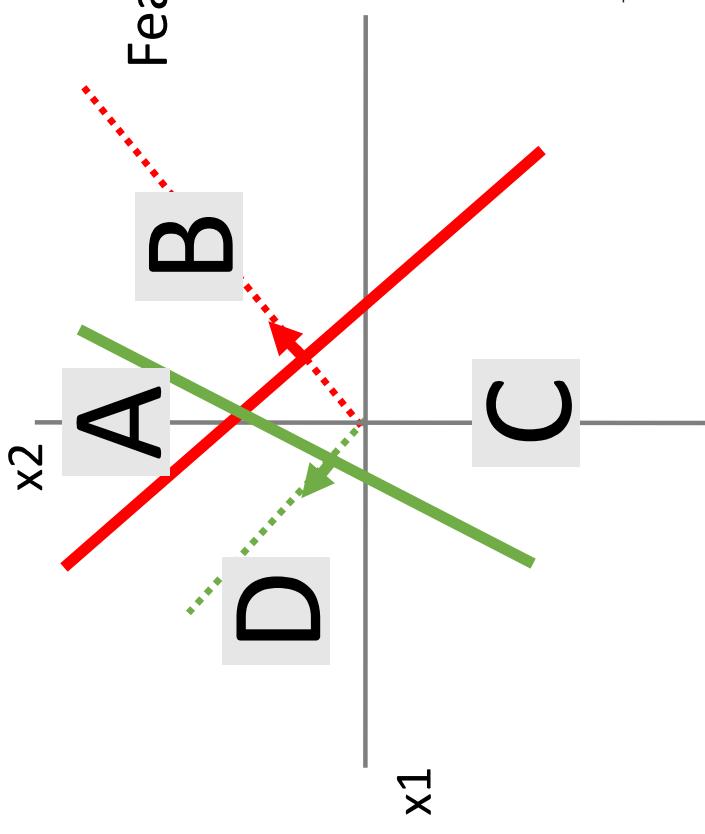


# Space Warping

Consider a neural net hidden layer:

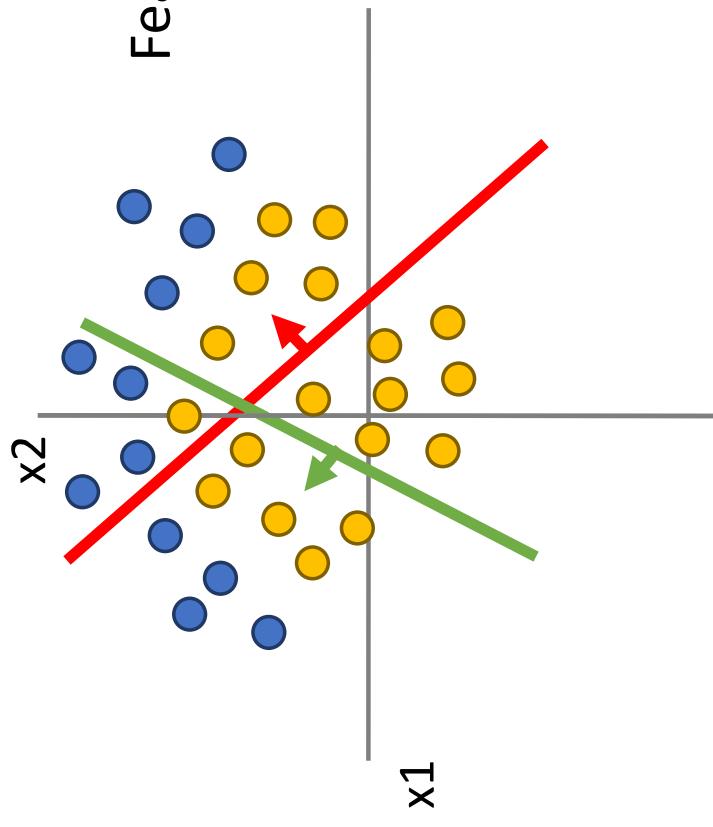
$$h = \text{ReLU}(Wx) = \max(0, Wx)$$

Where  $x, h$  are both 2-dimensional

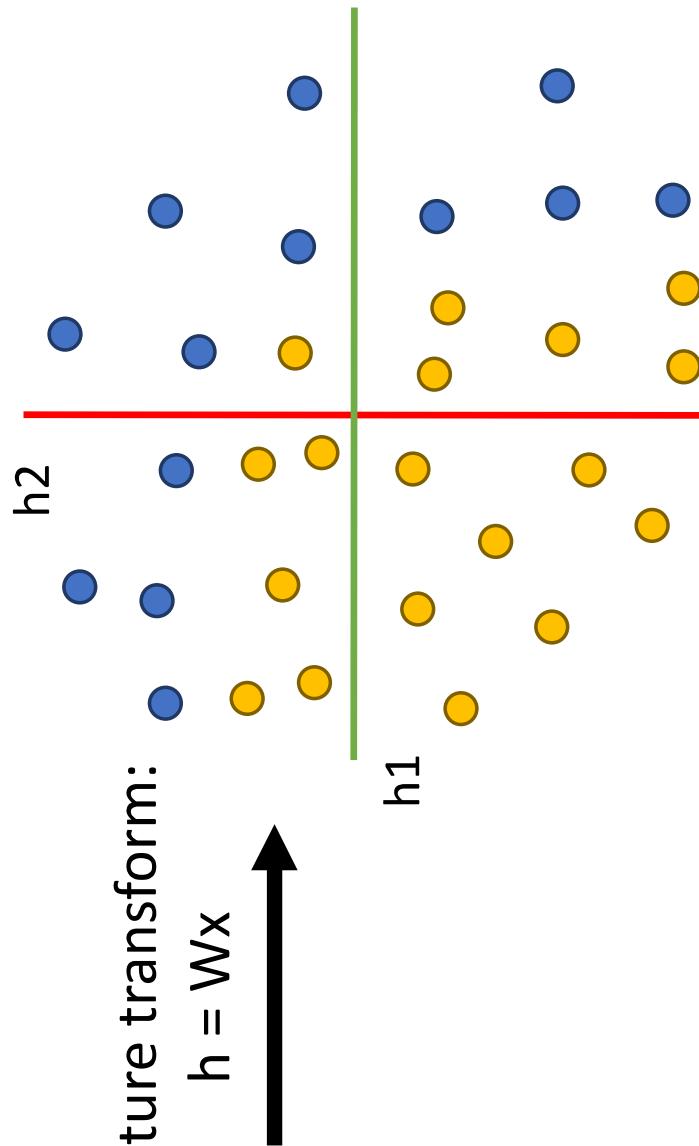


# Space Warping

Points not linearly  
separable in original space

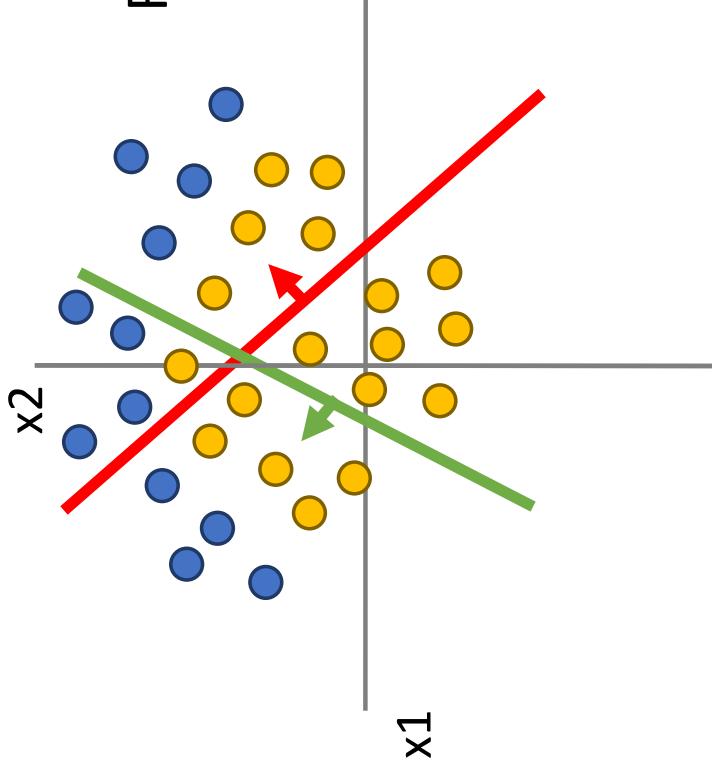


Consider a neural net hidden layer:  
 $h = \text{ReLU}(Wx) = \max(0, Wx)$   
Where  $x, h$  are both 2-dimensional

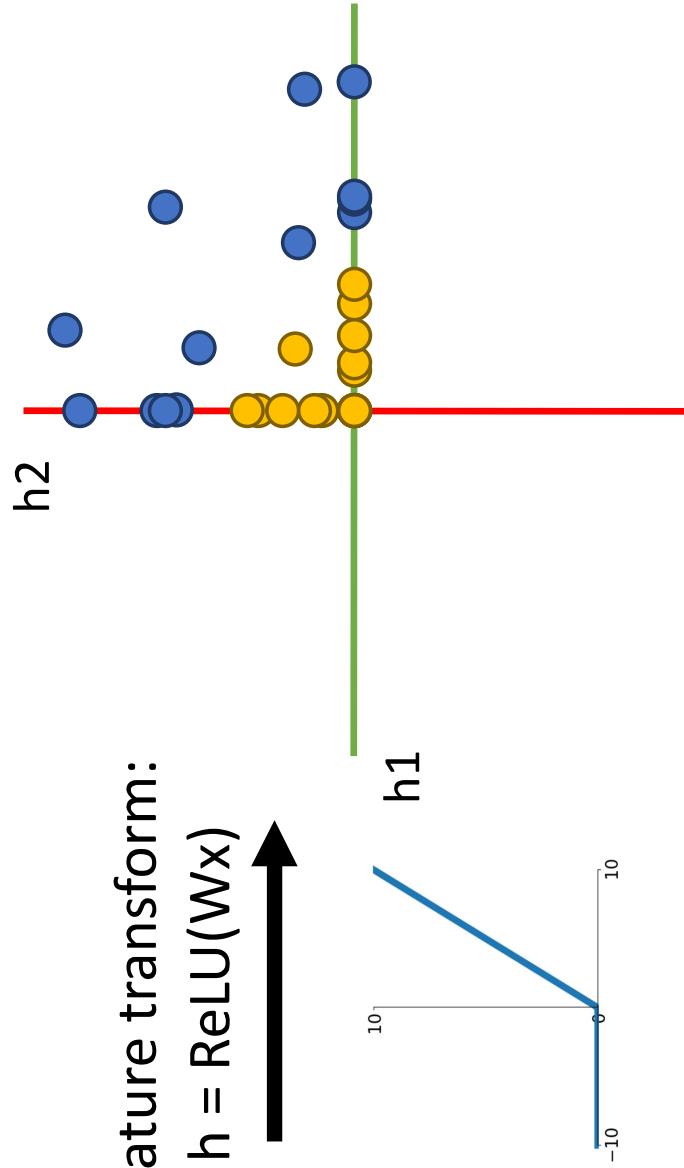


# Space Warping

Points not linearly  
separable in original space

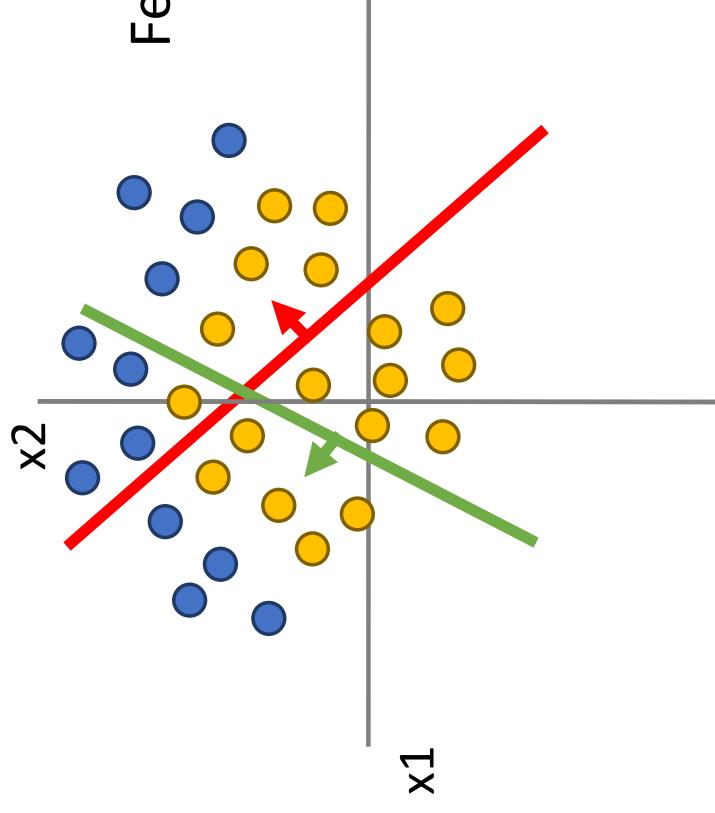


Consider a neural net hidden layer:  
 $h = \text{ReLU}(Wx) = \max(0, Wx)$   
Where  $x, h$  are both 2-dimensional

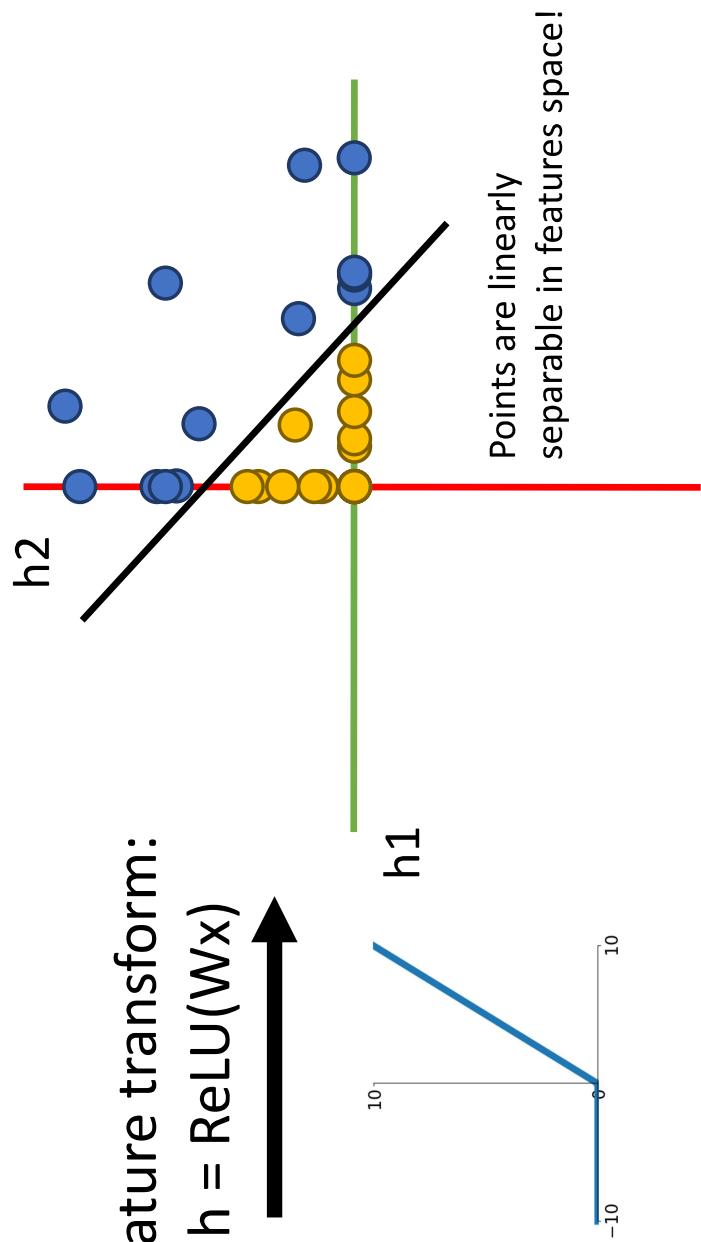


# Space Warping

Points not linearly  
separable in original space

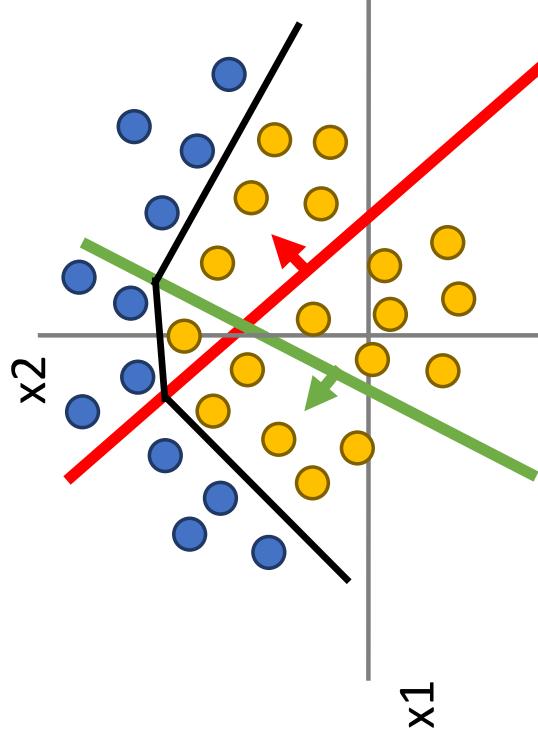


Consider a neural net hidden layer:  
 $h = \text{ReLU}(Wx) = \max(0, Wx)$   
Where  $x, h$  are both 2-dimensional



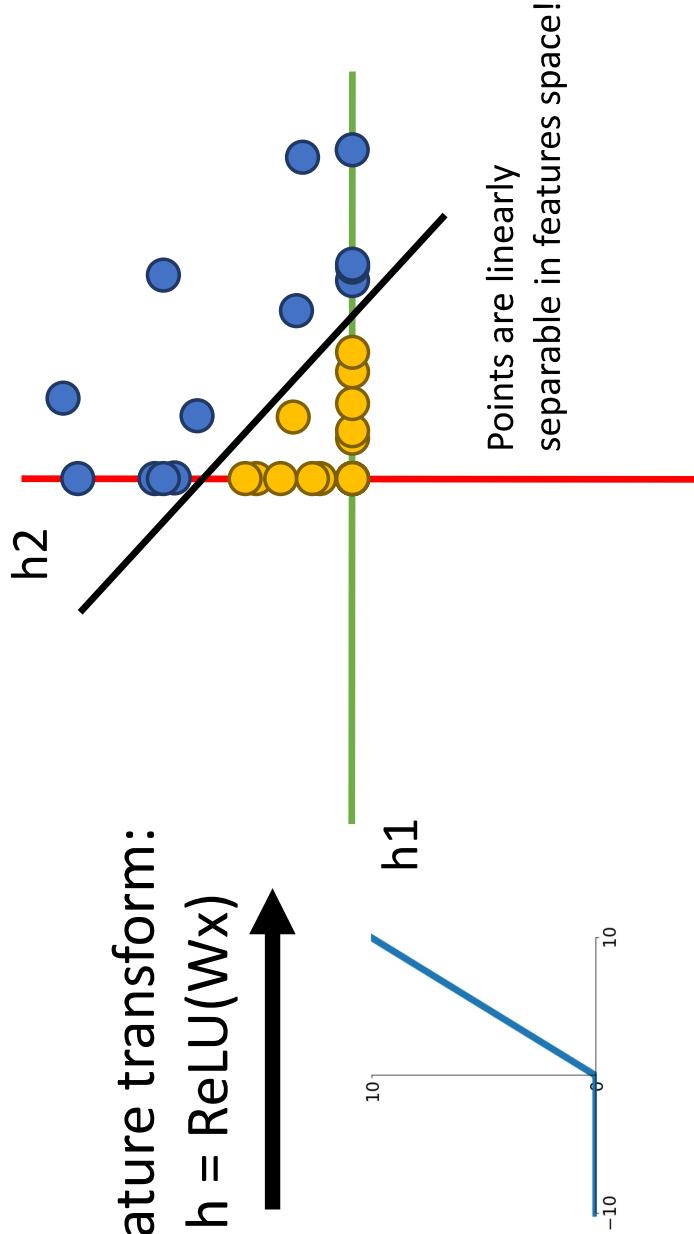
# Space Warping

Points not linearly  
separable in original space



Linear classifier in feature  
space gives nonlinear  
classifier in original space

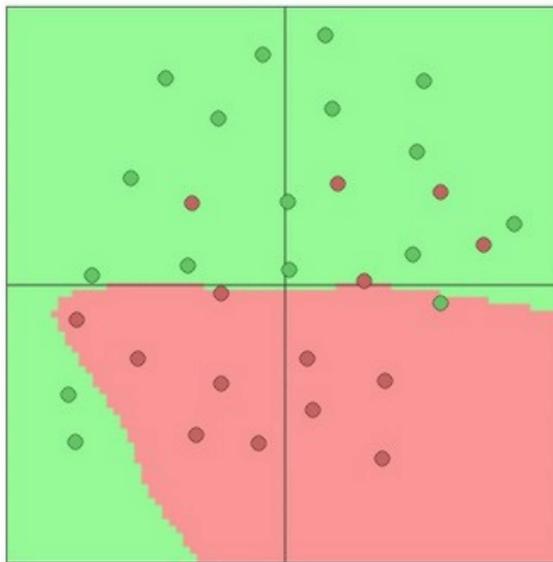
Consider a neural net hidden layer:  
 $h = \text{ReLU}(Wx) = \max(0, Wx)$   
Where  $x, h$  are both 2-dimensional



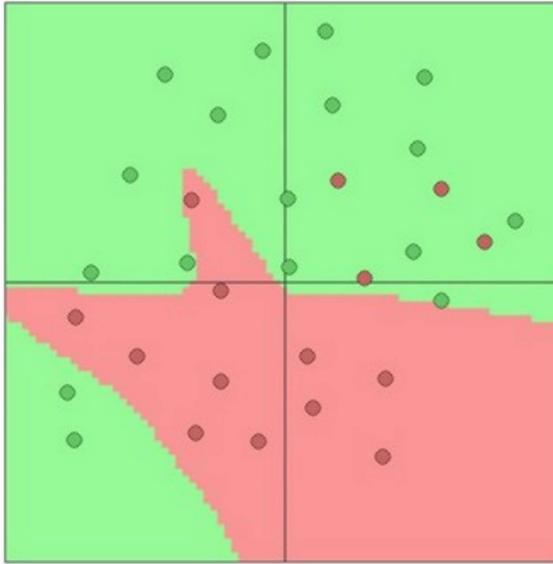
Points are linearly  
separable in features space!

# Setting the number of layers and their sizes

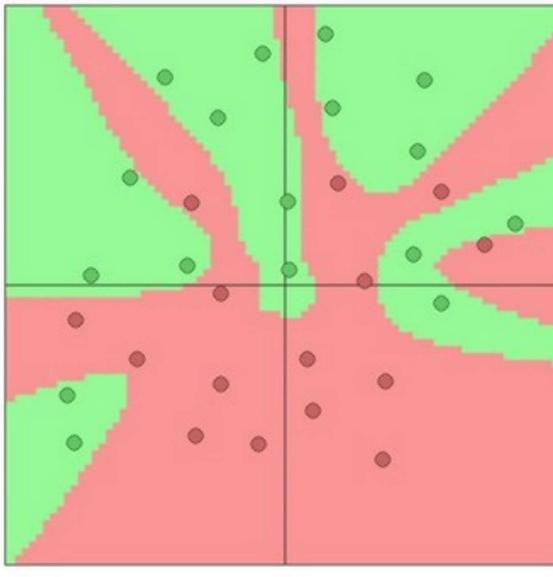
3 hidden units



6 hidden units



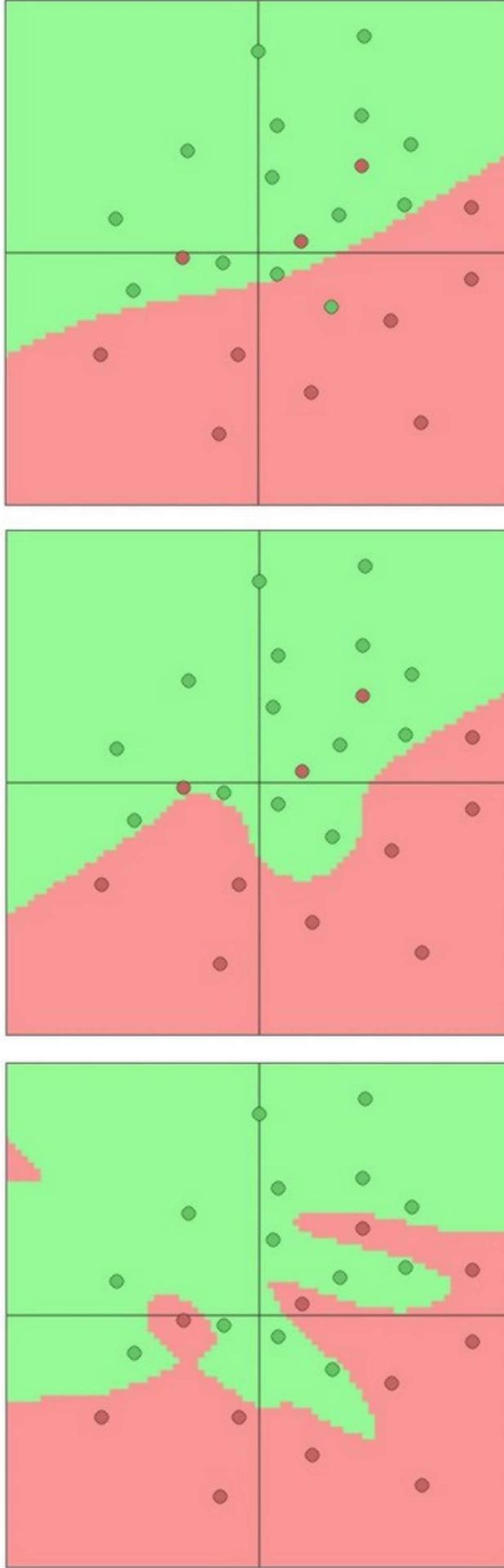
20 hidden units



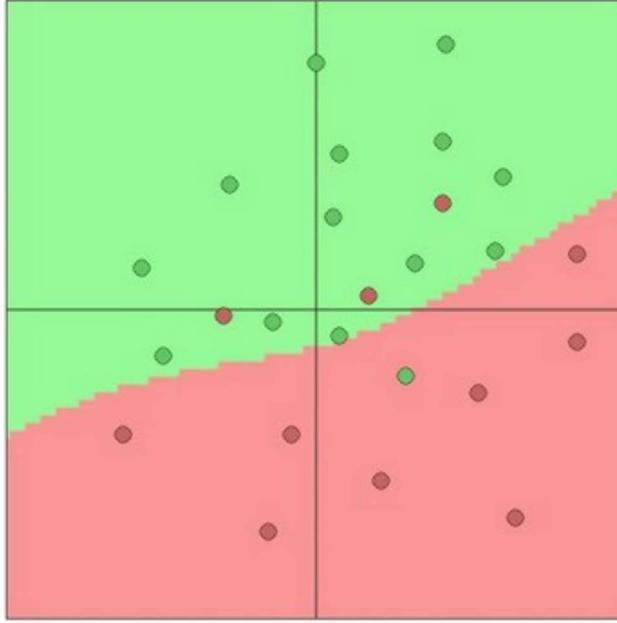
More hidden units = more capacity

Don't regularize with size; instead use stronger L2

$$\lambda = 0.001$$


$$\lambda = 0.01$$

$$\lambda = 0.1$$



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

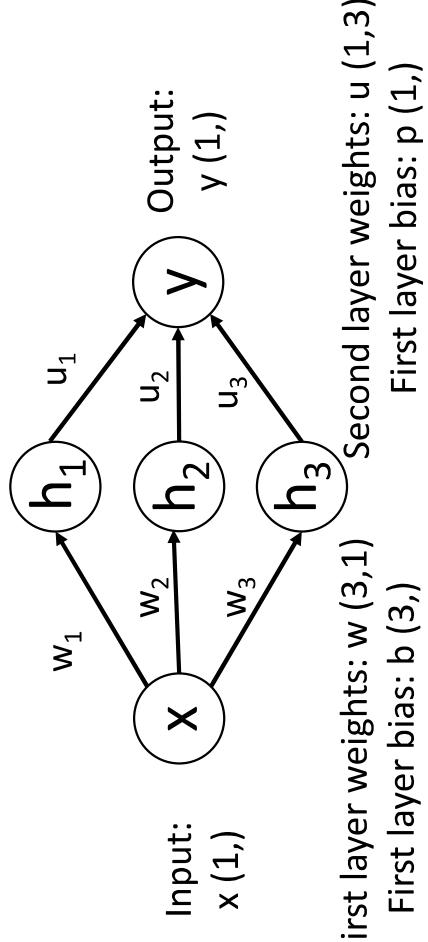
# Universal Approximation

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

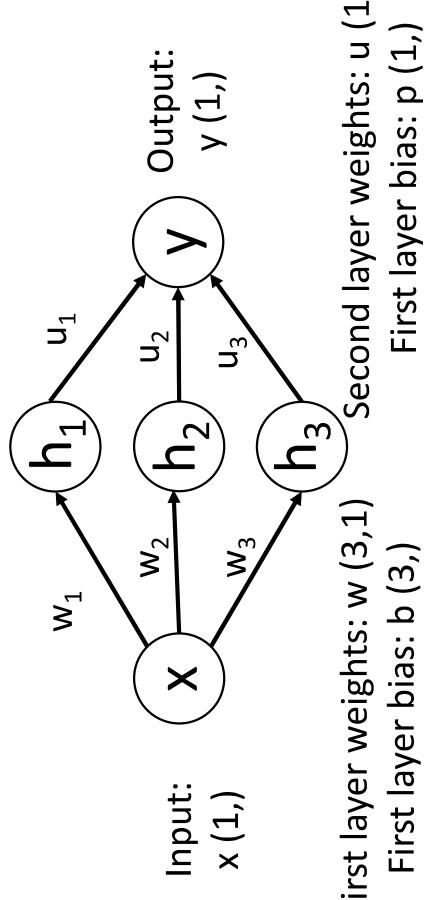
# Universal Approximation

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



# Universal Approximation

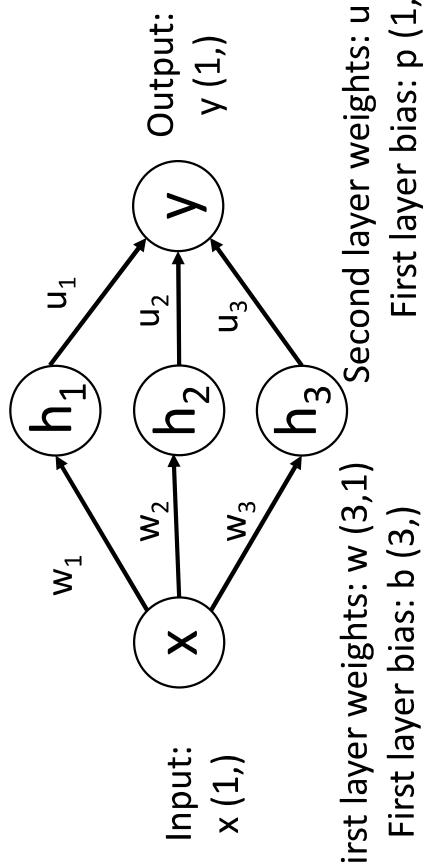
Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



$$\begin{aligned}h1 &= \max(0, w1 * x + b1) \\h2 &= \max(0, w2 * x + b2) \\h3 &= \max(0, w3 * x + b3) \\y &= u1 * h1 + u2 * h2 + u3 * h3 + p\end{aligned}$$

# Universal Approximation

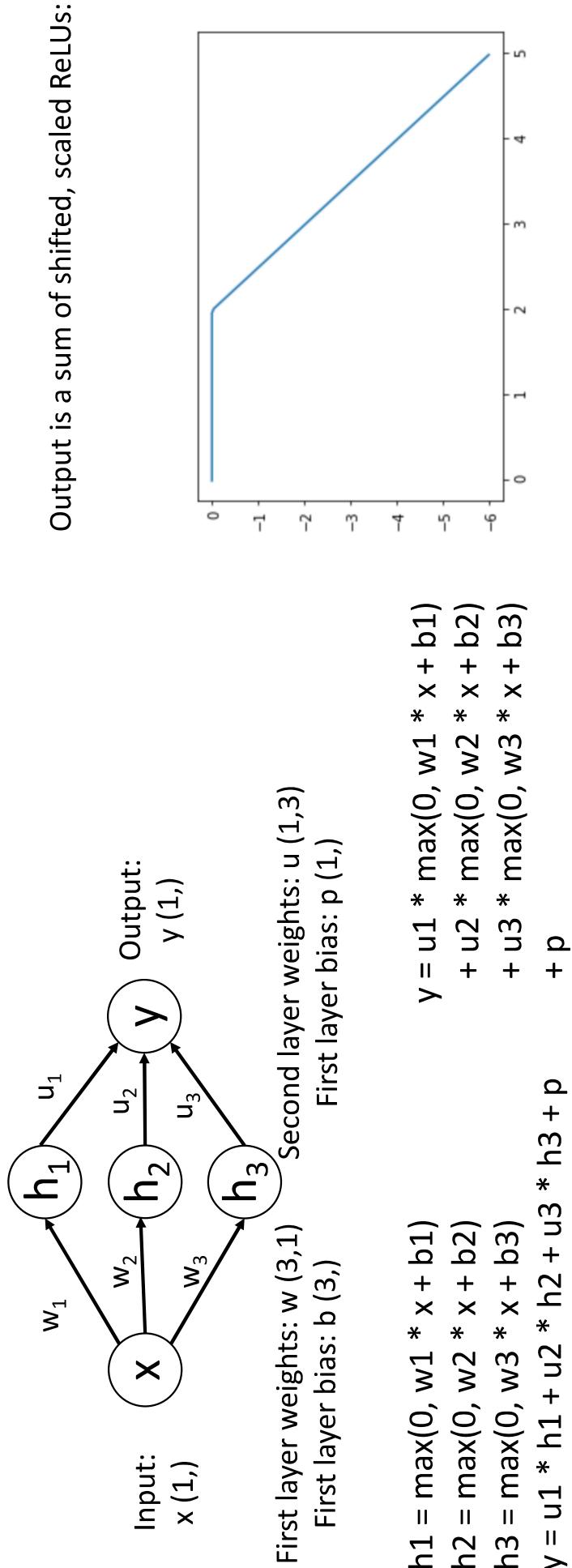
Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



$$\begin{aligned} h1 &= \max(0, w1 * x + b1) \\ h2 &= \max(0, w2 * x + b2) \\ h3 &= \max(0, w3 * x + b3) \\ y &= u1 * h1 + u2 * h2 + u3 * h3 + p \end{aligned}$$
$$\begin{aligned} y &= u1 * \max(0, w1 * x + b1) \\ &\quad + u2 * \max(0, w2 * x + b2) \\ &\quad + u3 * \max(0, w3 * x + b3) \\ &\quad + p \end{aligned}$$

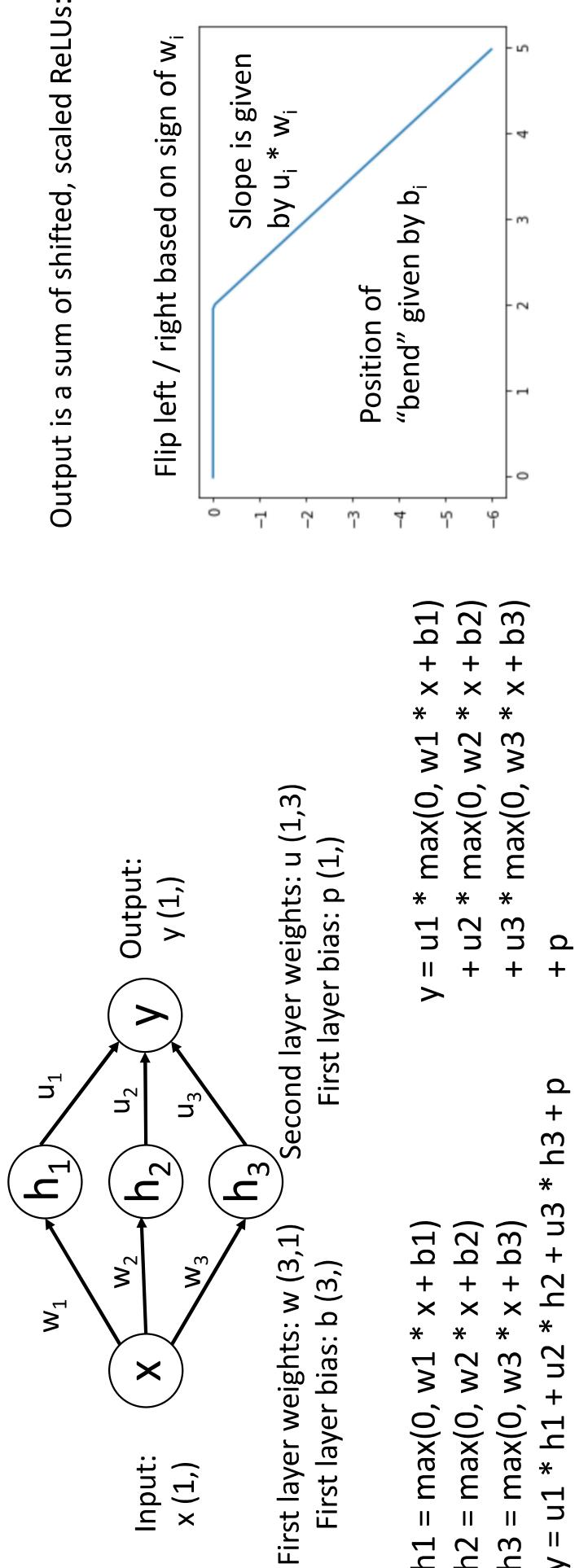
# Universal Approximation

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



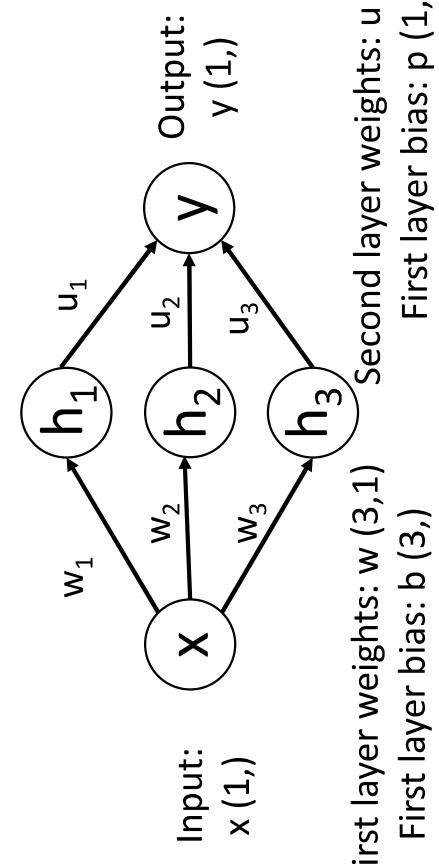
# Universal Approximation

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

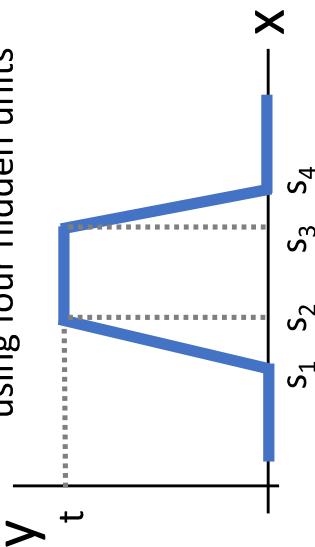


# Universal Approximation

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



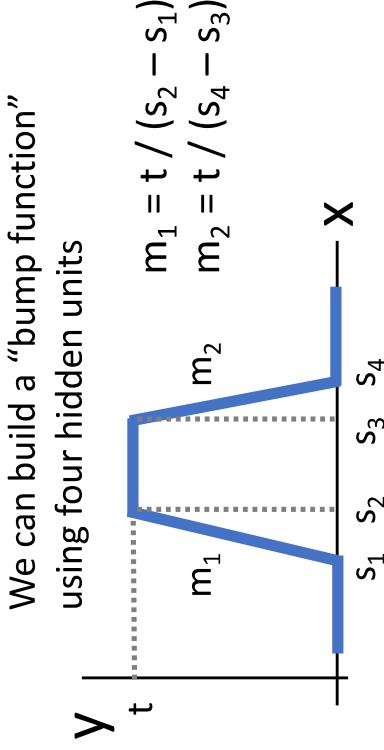
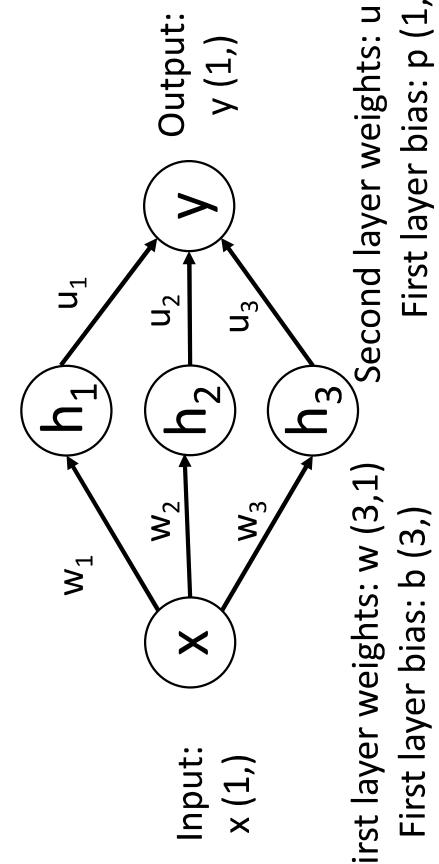
We can build a “bump function”  
using four hidden units



$$\begin{aligned} h1 &= \max(0, w1 * x + b1) \\ h2 &= \max(0, w2 * x + b2) \\ h3 &= \max(0, w3 * x + b3) \\ y &= u1 * h1 + u2 * h2 + u3 * h3 + p \end{aligned}$$
$$y = u1 * \max(0, w1 * x + b1) + u2 * \max(0, w2 * x + b2) + u3 * \max(0, w3 * x + b3) + p$$

# Universal Approximation

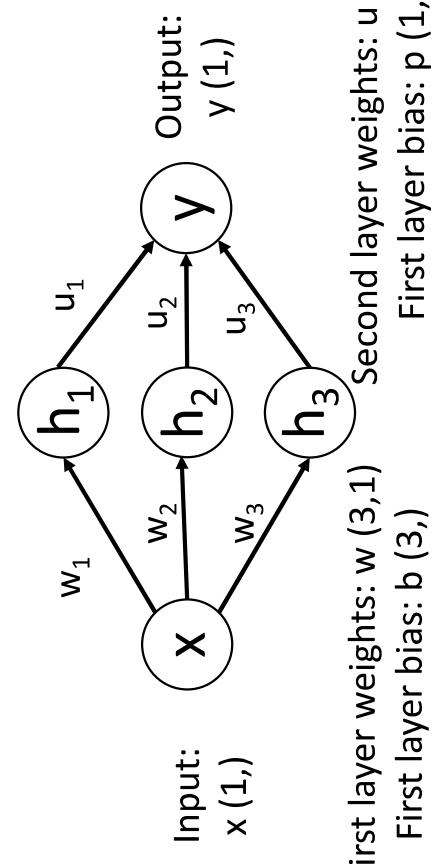
Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



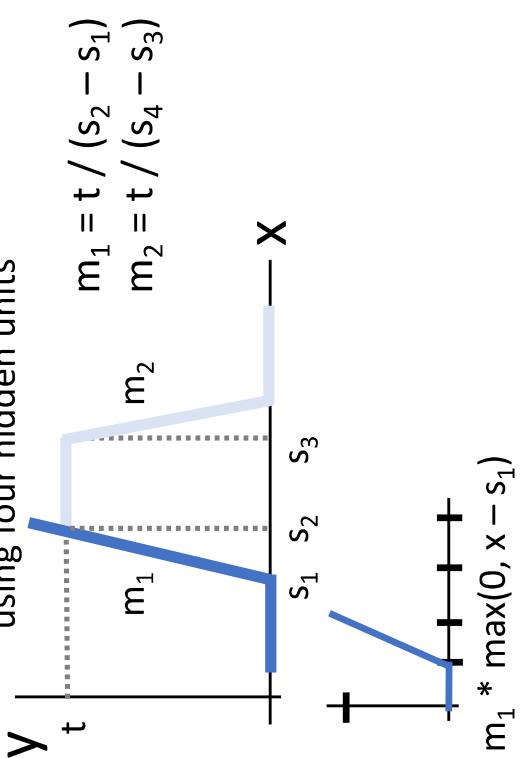
```
h1 = max(0, w1 * x + b1)
h2 = max(0, w2 * x + b2)
h3 = max(0, w3 * x + b3)
y = u1 * h1 + u2 * h2 + u3 * h3 + p
y = u1 * max(0, w1 * x + b1)
      + u2 * max(0, w2 * x + b2)
      + u3 * max(0, w3 * x + b3)
      + p
```

# Universal Approximation

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



We can build a "bump function" using four hidden units



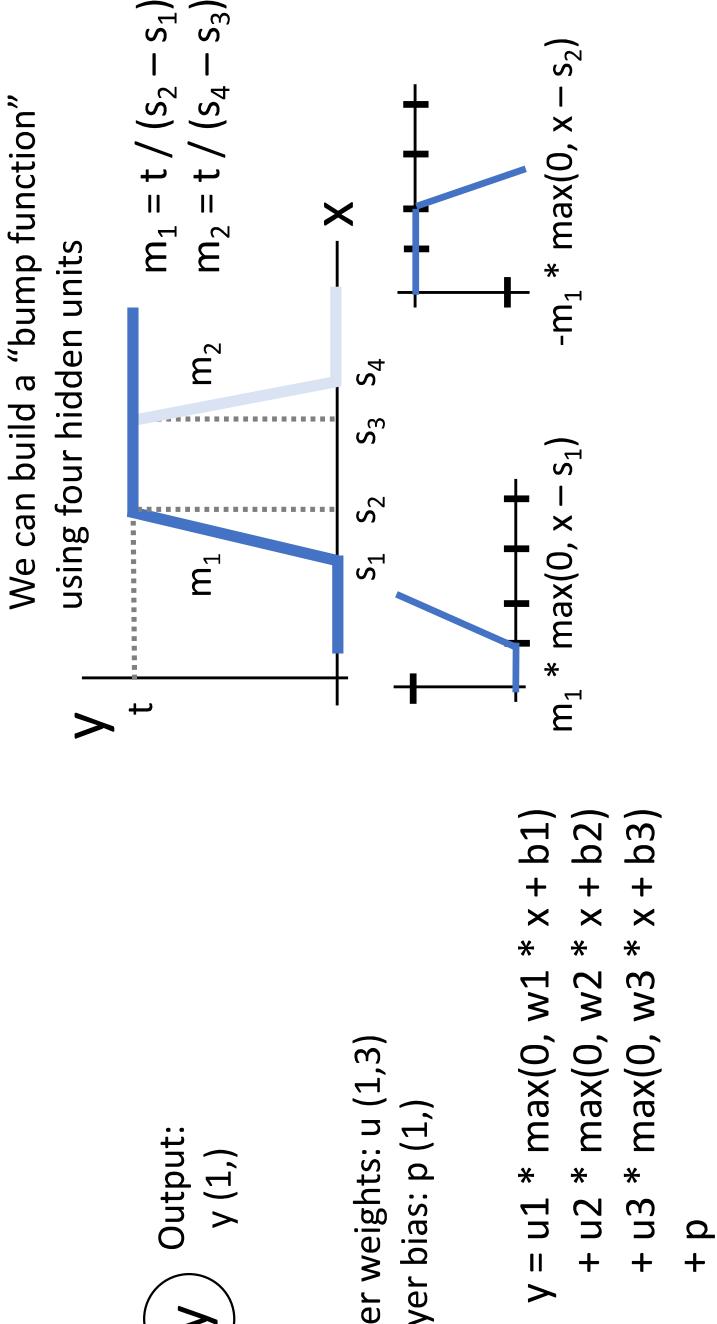
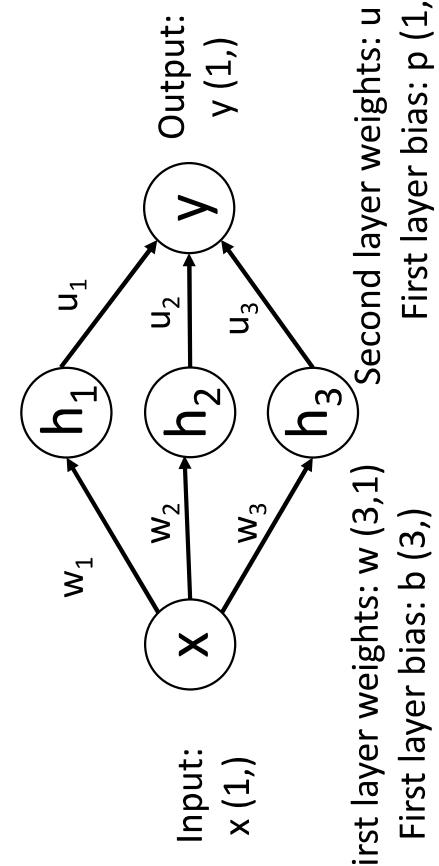
$$\begin{aligned}
 h1 &= \max(0, w1 * x + b1) \\
 h2 &= \max(0, w2 * x + b2) \\
 h3 &= \max(0, w3 * x + b3) \\
 y &= u1 * h1 + u2 * h2 + u3 * h3 + p
 \end{aligned}$$

$$\begin{aligned}
 y &= u1 * \max(0, w1 * x + b1) \\
 &\quad + u2 * \max(0, w2 * x + b2) \\
 &\quad + u3 * \max(0, w3 * x + b3) \\
 &\quad + p
 \end{aligned}$$

$$\begin{aligned}
 m_1 * \max(0, x - s_1) \\
 m_1 * \max(0, x - s_1)
 \end{aligned}$$

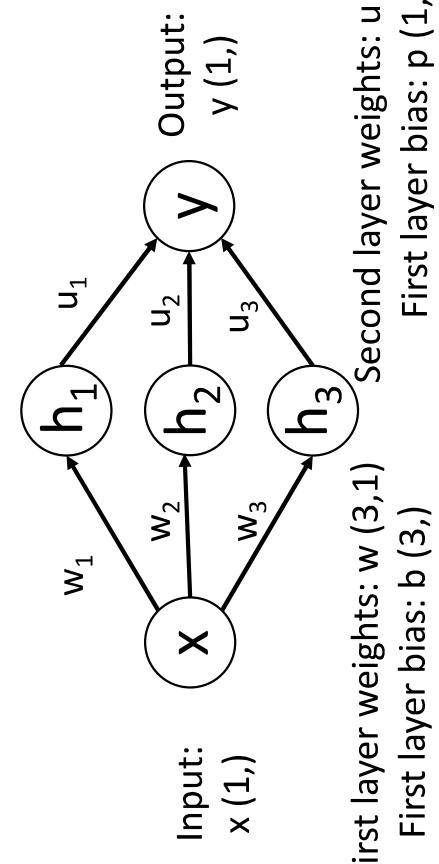
# Universal Approximation

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

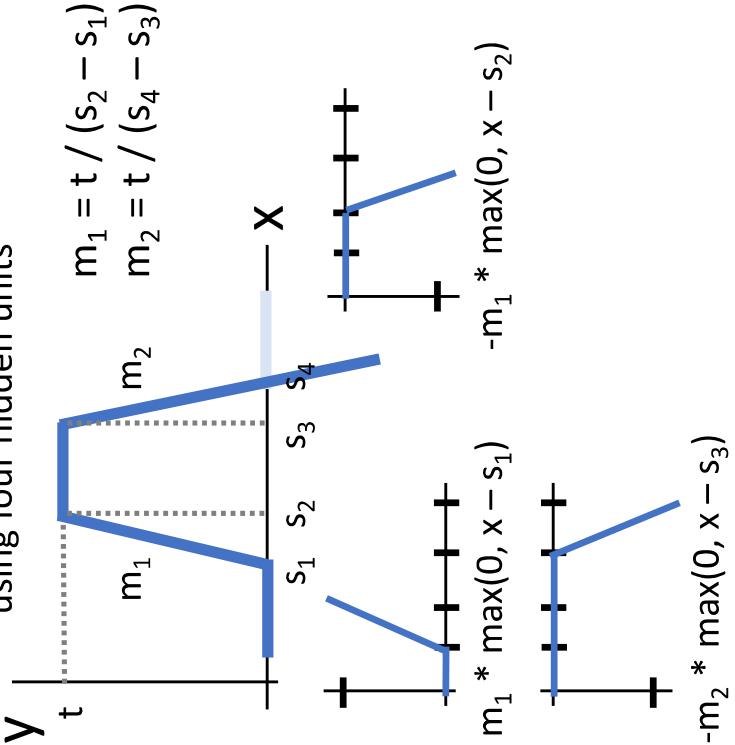


# Universal Approximation

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



We can build a "bump function"  
using four hidden units

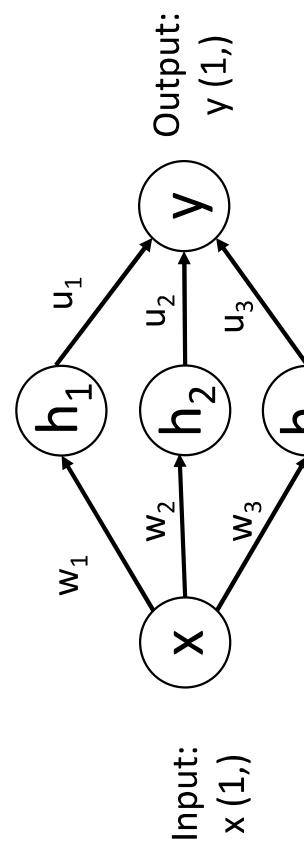


$$\begin{aligned}
 h1 &= \max(0, w1 * x + b1) \\
 h2 &= \max(0, w2 * x + b2) \\
 h3 &= \max(0, w3 * x + b3) \\
 h4 &= \max(0, w4 * x + b4) \\
 y &= h1 + h2 + h3 + h4
 \end{aligned}$$

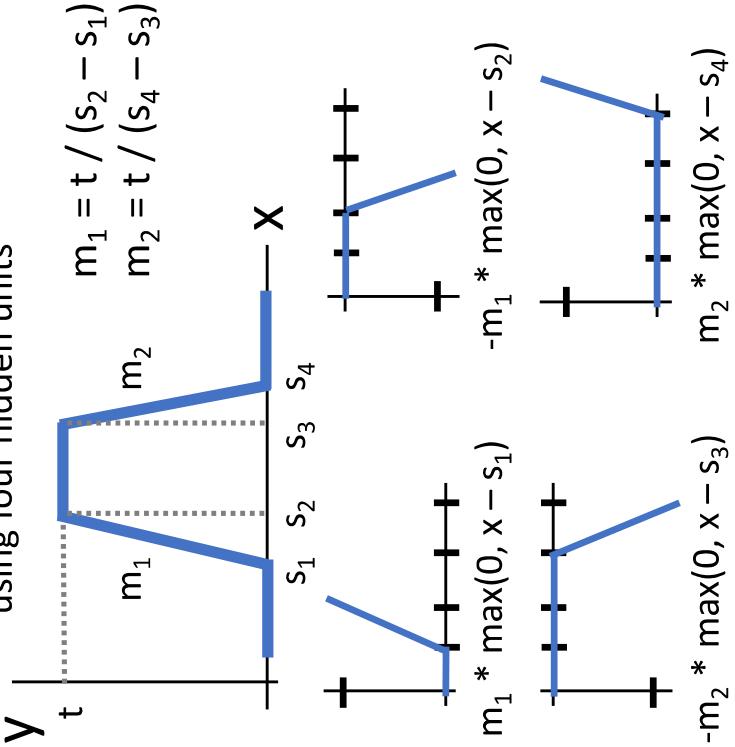
$$\begin{aligned}
 y &= u1 * \max(0, w1 * x + b1) \\
 &\quad + u2 * \max(0, w2 * x + b2) \\
 &\quad + u3 * \max(0, w3 * x + b3) \\
 &\quad + u4 * \max(0, x - s_1)
 \end{aligned}$$

# Universal Approximation

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



We can build a "bump function" using four hidden units



First layer weights:  $w (3,1)$   
First layer bias:  $b (3,)$   
Second layer weights:  $u (1,3)$   
Second layer bias:  $p (1,)$

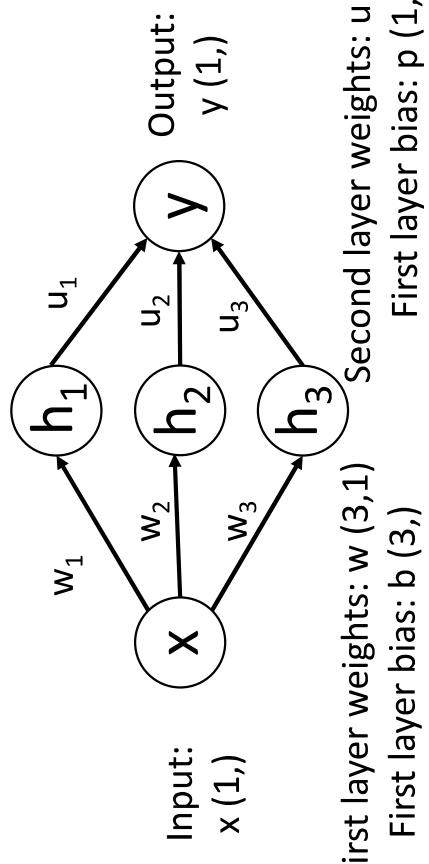
$$\begin{aligned} h1 &= \max(0, w1 * x + b1) \\ h2 &= \max(0, w2 * x + b2) \\ h3 &= \max(0, w3 * x + b3) \\ y &= u1 * h1 + u2 * h2 + u3 * h3 + p \end{aligned}$$

$$\begin{aligned} y &= u1 * \max(0, w1 * x + b1) \\ &\quad + u2 * \max(0, w2 * x + b2) \\ &\quad + u3 * \max(0, w3 * x + b3) \\ &\quad + p \end{aligned}$$

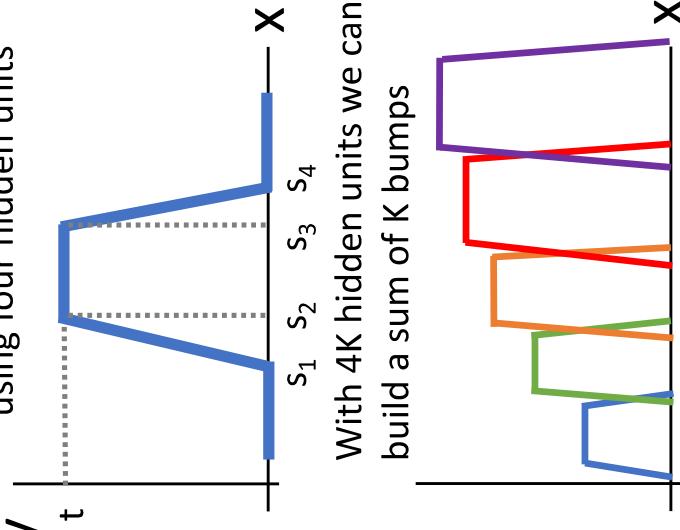
$$\begin{aligned} h1 &= \max(0, w1 * x + b1) \\ h2 &= \max(0, w2 * x + b2) \\ h3 &= \max(0, w3 * x + b3) \\ y &= u1 * \max(0, x - s_1) \\ &\quad + u2 * \max(0, x - s_2) \\ &\quad + u3 * \max(0, x - s_3) \end{aligned}$$

# Universal Approximation

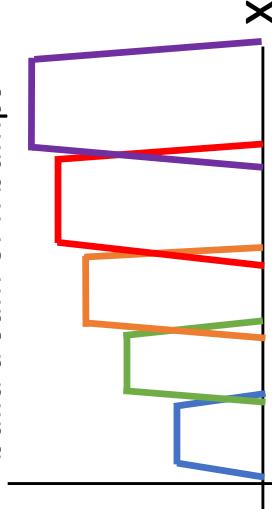
Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



We can build a "bump function" using four hidden units



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

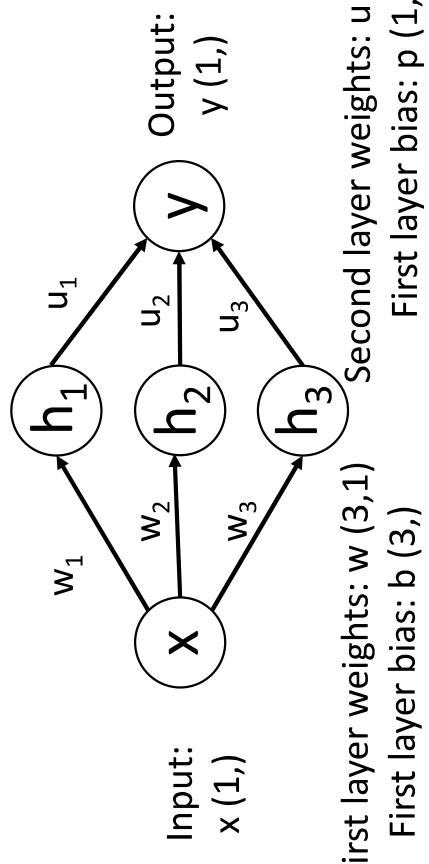


$$\begin{aligned}y &= u1 * \max(0, w1 * x + b1) \\&\quad + u2 * \max(0, w2 * x + b2) \\&\quad + u3 * \max(0, w3 * x + b3) \\&\quad + p\end{aligned}$$

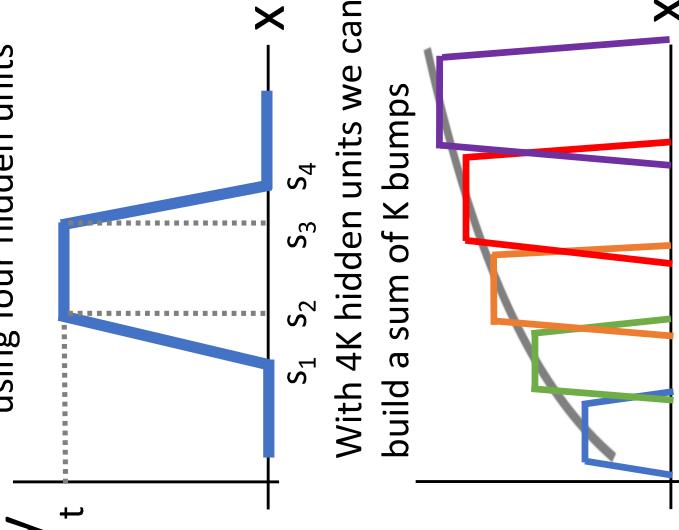
$$\begin{aligned}h1 &= \max(0, w1 * x + b1) \\h2 &= \max(0, w2 * x + b2) \\h3 &= \max(0, w3 * x + b3) \\y &= u1 * h1 + u2 * h2 + u3 * h3 + p\end{aligned}$$

# Universal Approximation

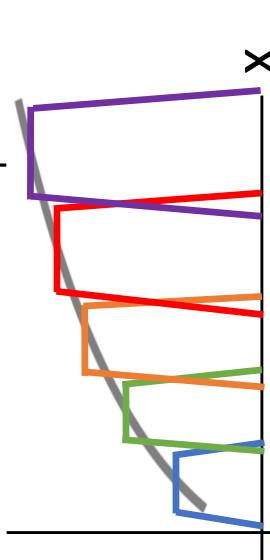
Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



We can build a "bump function"  
using four hidden units



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



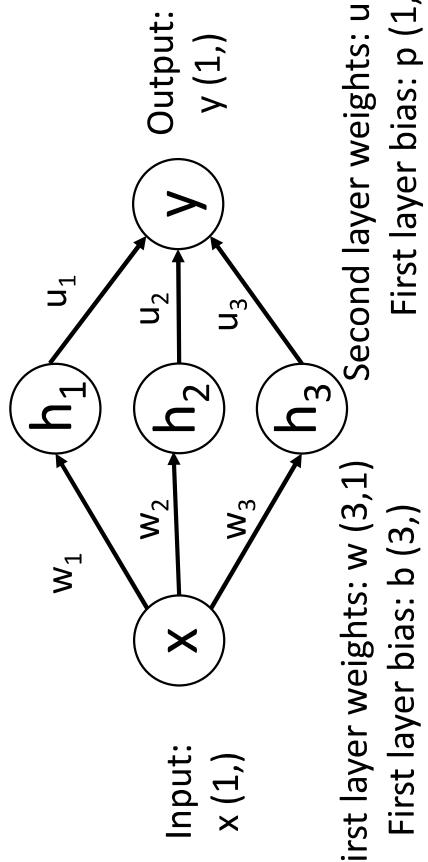
Approximate functions with bumps!

$$\begin{aligned} h1 &= \max(0, w1 * x + b1) \\ h2 &= \max(0, w2 * x + b2) \\ h3 &= \max(0, w3 * x + b3) \\ y &= u1 * h1 + u2 * h2 + u3 * h3 + p \end{aligned}$$

$$y = u1 * \max(0, w1 * x + b1) + u2 * \max(0, w2 * x + b2) + u3 * \max(0, w3 * x + b3) + p$$

# Universal Approximation

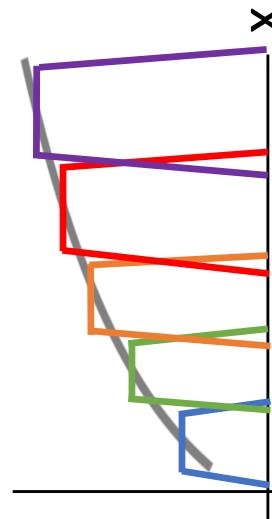
Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



What about...

- Gaps between bumps?
- Other nonlinearities?
- Higher-dimensional functions?

See [Nielsen, Chapter 4](#)



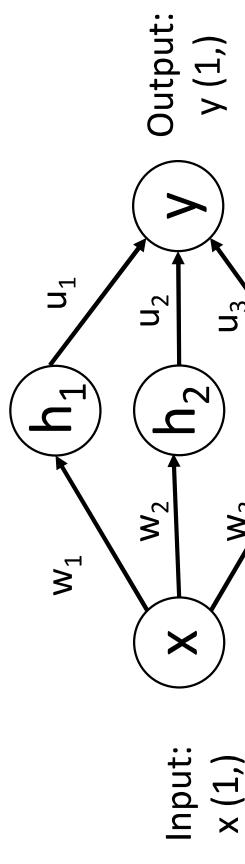
Approximate functions with bumps!

$$\begin{aligned} h1 &= \max(0, w1 * x + b1) \\ h2 &= \max(0, w2 * x + b2) \\ h3 &= \max(0, w3 * x + b3) \\ y &= u1 * h1 + u2 * h2 + u3 * h3 + p \end{aligned}$$

$$\begin{aligned} y &= u1 * \max(0, w1 * x + b1) \\ &\quad + u2 * \max(0, w2 * x + b2) \\ &\quad + u3 * \max(0, w3 * x + b3) \\ &\quad + p \end{aligned}$$

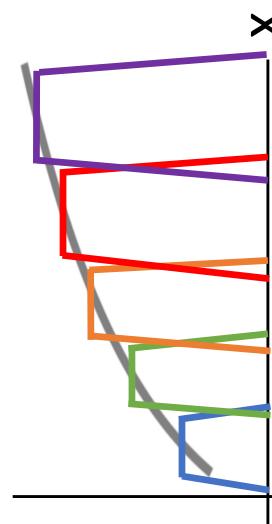
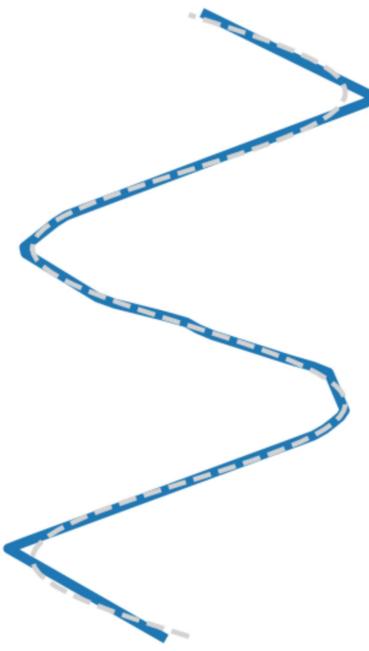
# Universal Approximation

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:  $p(1,)$

Reality check: Networks don't really learn bumps!



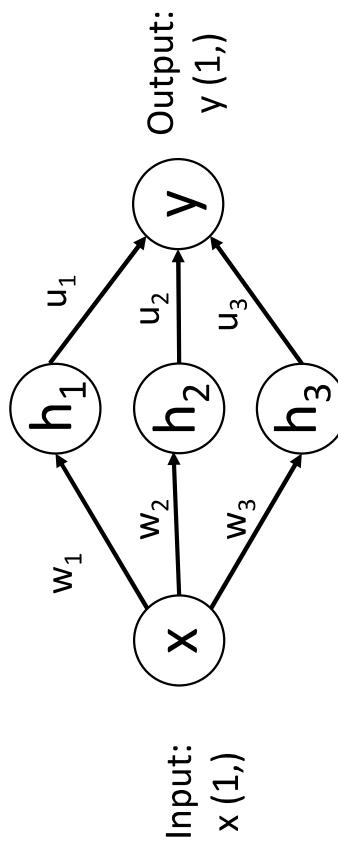
Approximate functions with bumps!

```
h1 = max(0, w1 * x + b1)
h2 = max(0, w2 * x + b2)
h3 = max(0, w3 * x + b3)
y = u1 * h1 + u2 * h2 + u3 * h3 + p
```

$$y = u1 * \max(0, w1 * x + b1) + u2 * \max(0, w2 * x + b2) + u3 * \max(0, w3 * x + b3) + p$$

# Universal Approximation

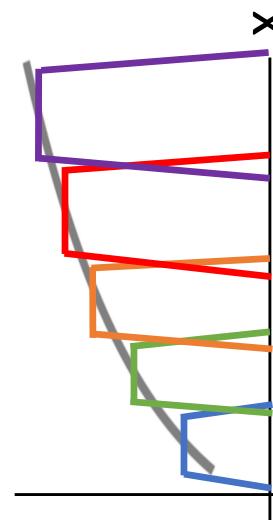
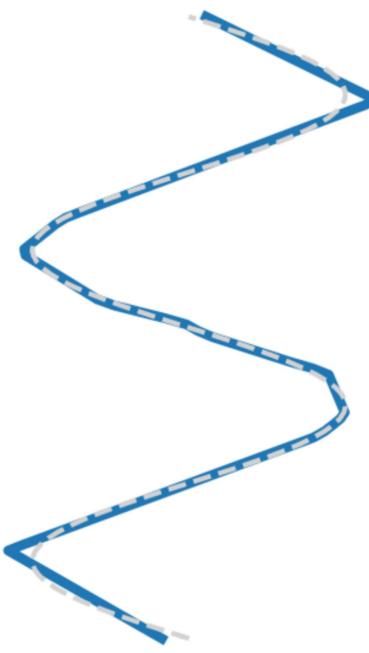
Example: Approximating a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  with a two-layer ReLU network



Universal approximation tells us:

- Neural nets can represent any function

Reality check: Networks don't really learn bumps!



Remember: kNN is also a universal approximator!

Approximate functions with bumps!

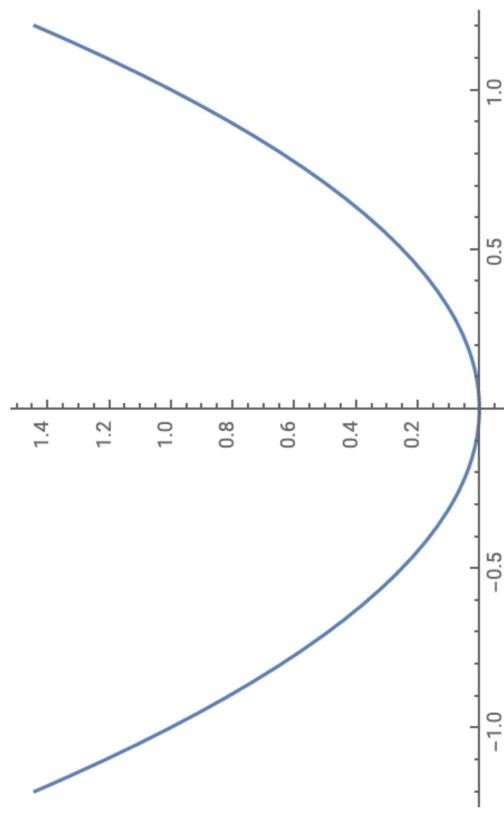
# Convex Functions

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  $x_1, x_2 \in X, t \in [0, 1]$ ,

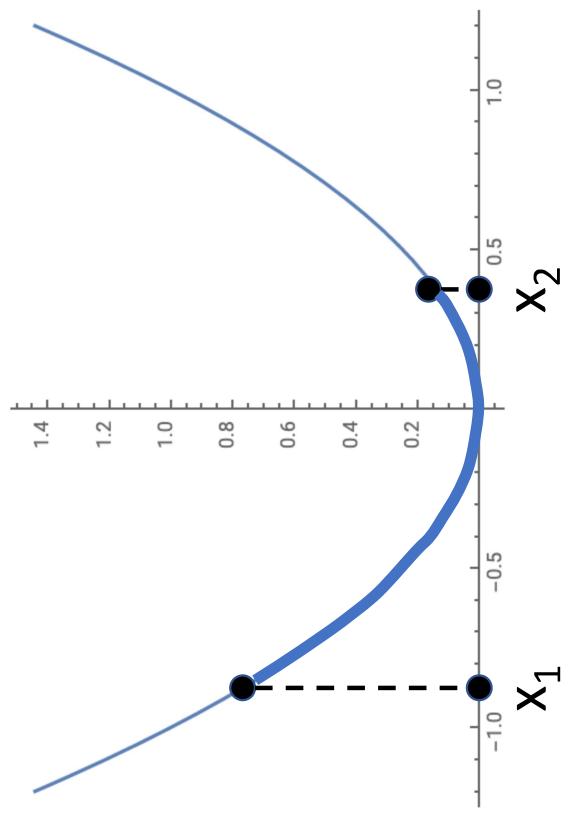
$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2)$$


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  $x_1, x_2 \in X, t \in [0, 1]$ ,

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

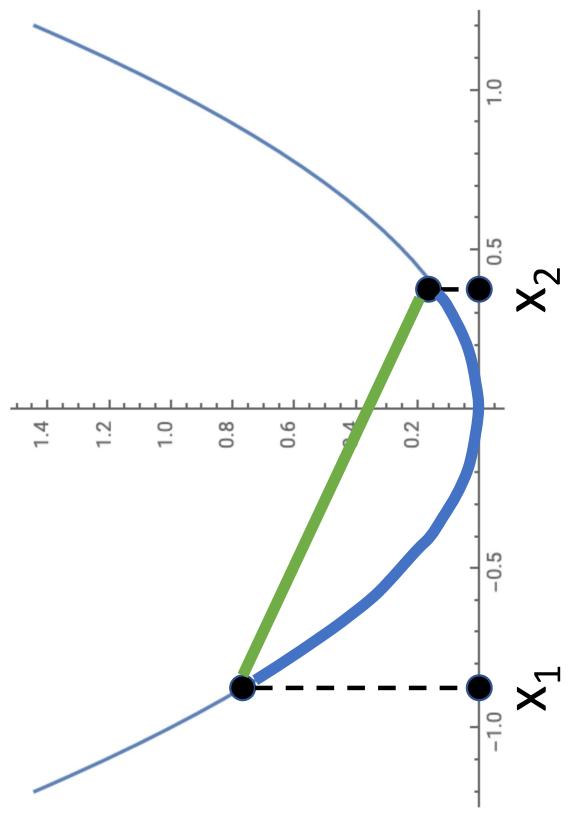


Example:  $f(x) = x^2$  is convex:

# Convex Functions

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

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

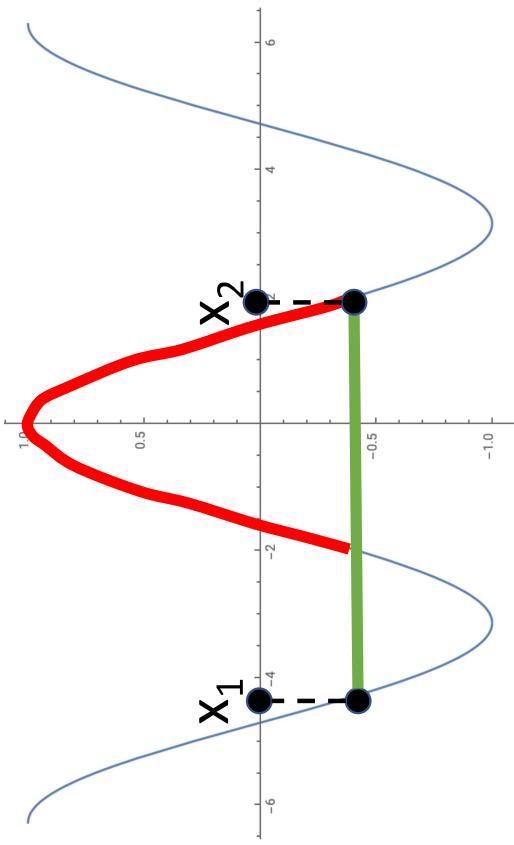


Example:  $f(x) = x^2$  is convex:

# Convex Functions

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

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



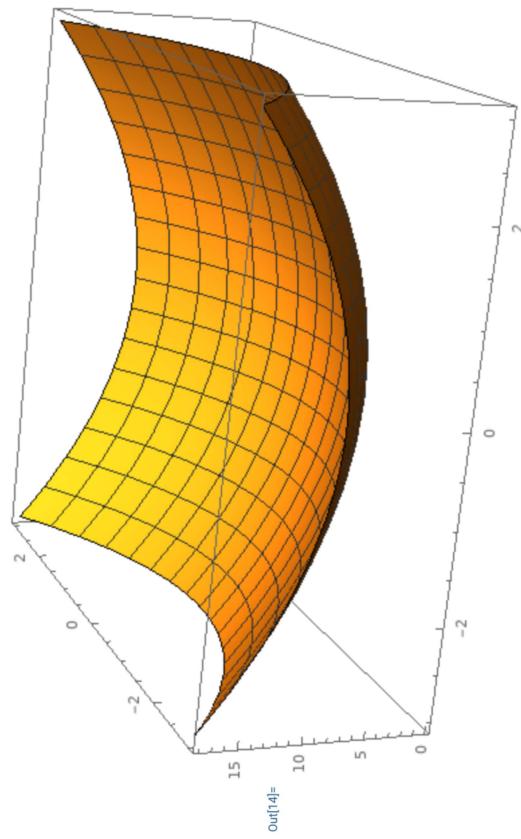
Example:  $f(x) = \cos(x)$   
is not convex:

# Convex Functions

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

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

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



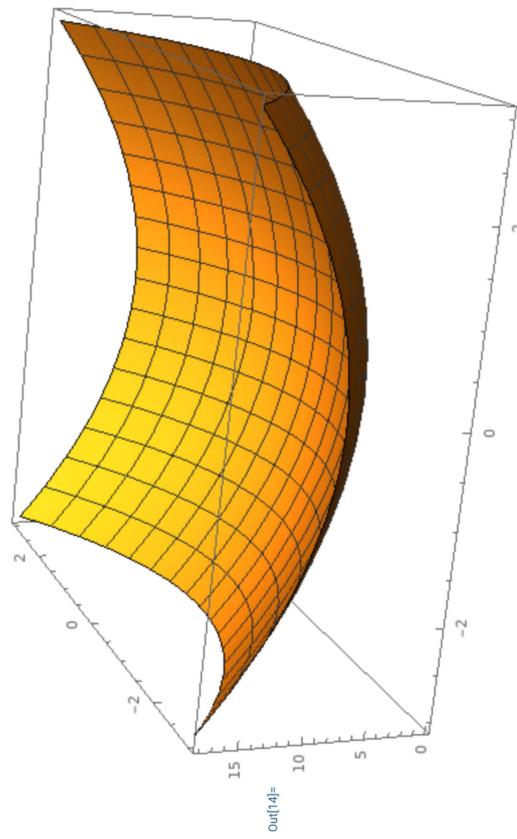
\*Many technical details! See e.g. IOE 661 / MATH 663

# 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**\*

\*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

Linear classifiers optimize a **convex function**!

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

$$\begin{aligned}s &= f(x; W) = Wx \\ L_i &= -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right) \text{ Softmax} \\ L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \text{ SVM} \\ L &= \frac{1}{N} \sum_{i=1}^N L_i + R(W) \\ R(W) &= \text{L2 or L1 regularization}\end{aligned}$$

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

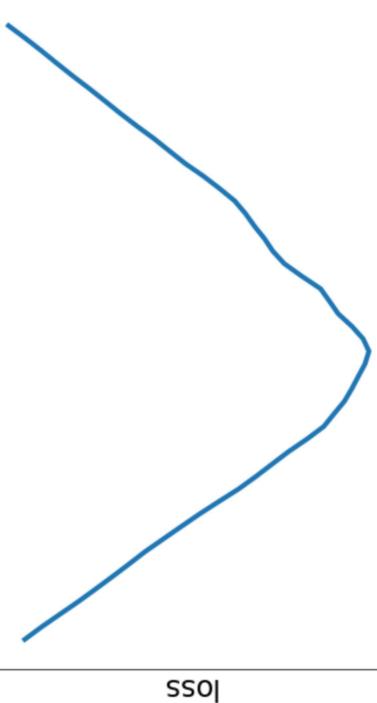
# Convex Functions

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\***



W1[0, 0]  
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

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

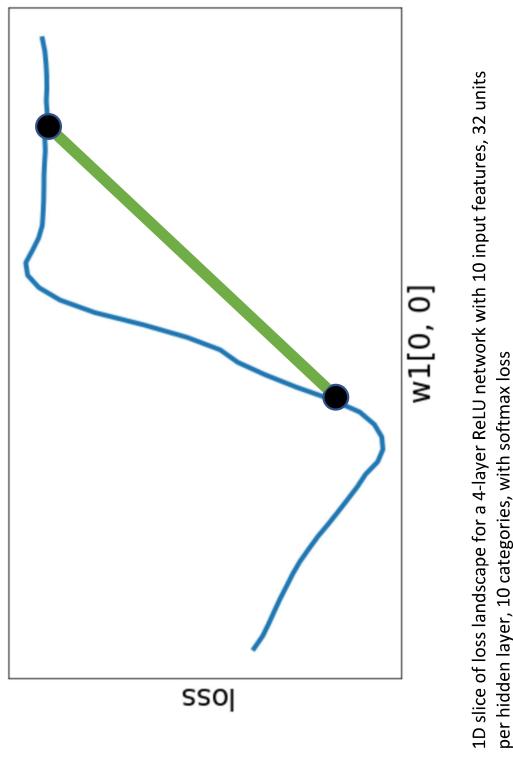
# Convex Functions

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

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

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

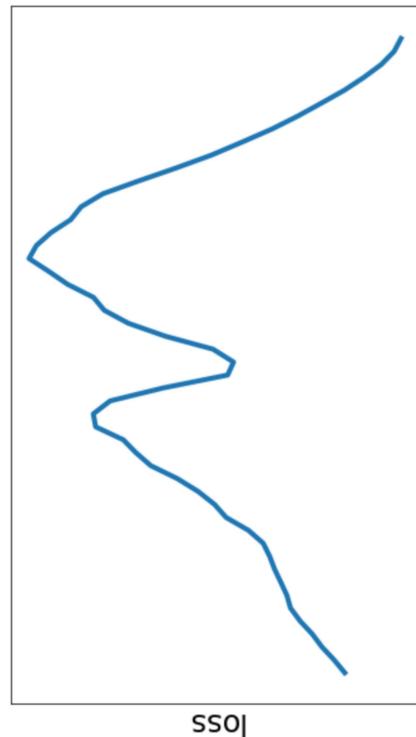
# Convex Functions

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

With local minima:



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

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

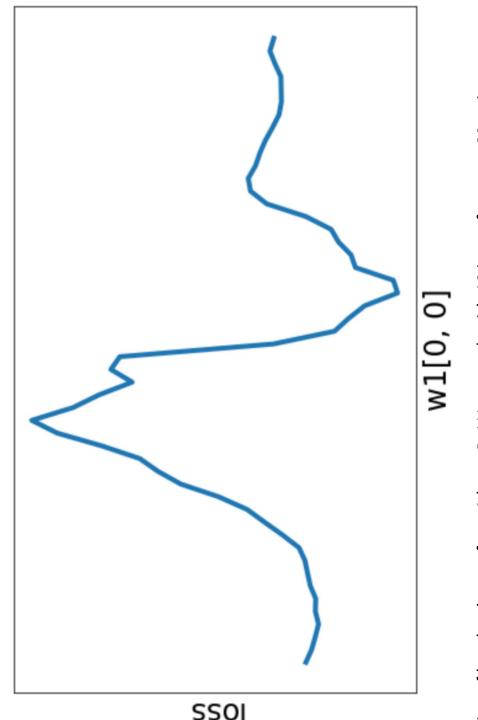
# Convex Functions

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

Can get very wild!



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

\*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**\*

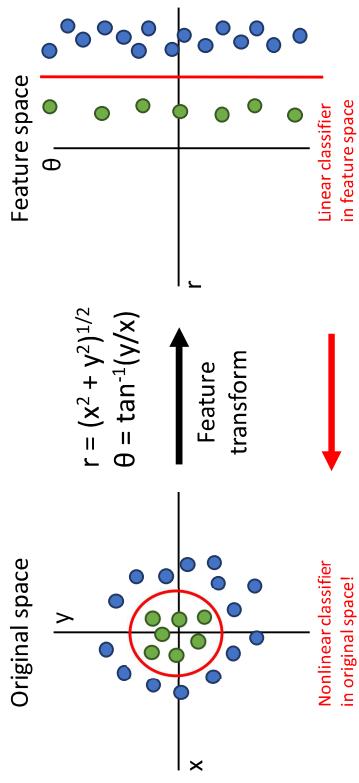
Most neural networks need **nonconvex optimization**

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

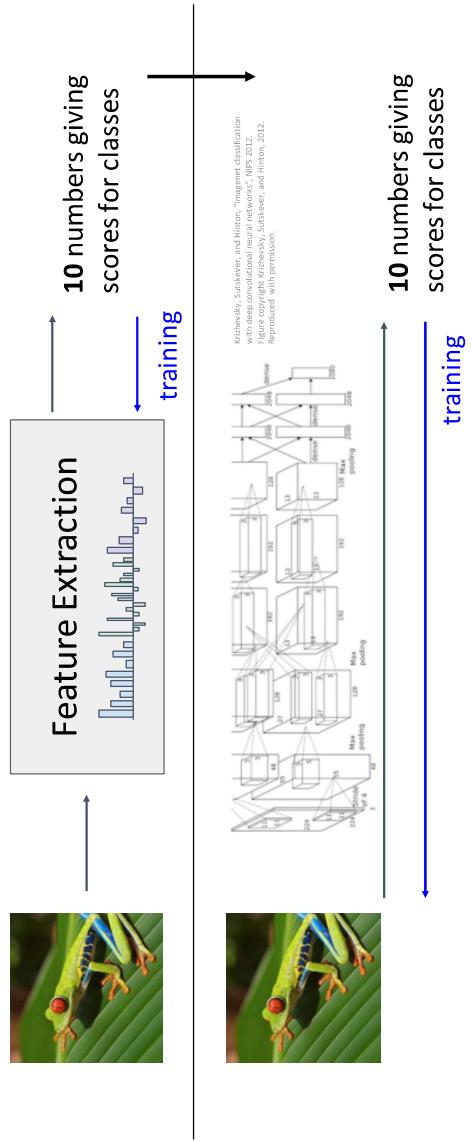
\*Many technical details! See e.g. IOE 661 / MATH 663

# Summary

Feature transform + Linear classifier  
allows nonlinear decision boundaries



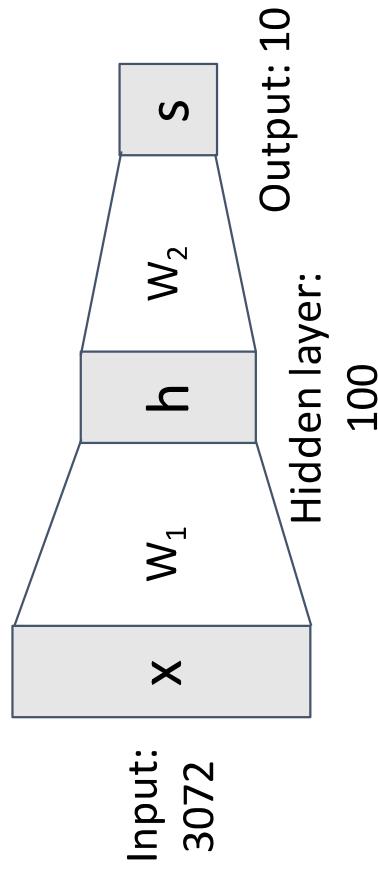
Neural Networks as learnable feature transforms



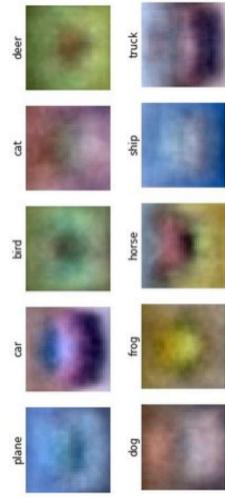
# Summary

From linear classifiers to  
fully-connected networks

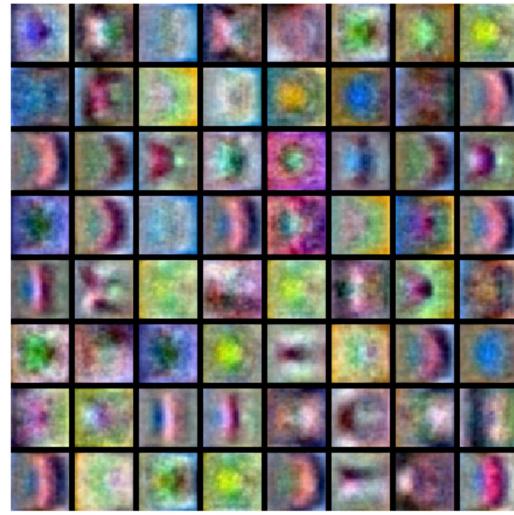
$$f = W_2 \max(0, W_1 x)$$



Linear classifier: One template per class



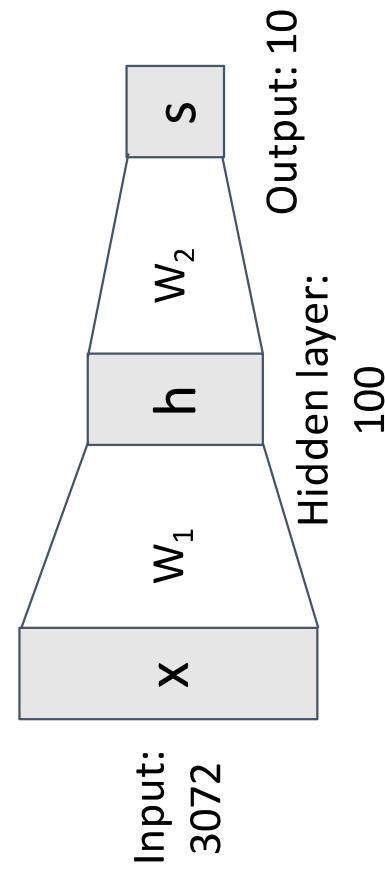
Neural networks: Many reusable templates



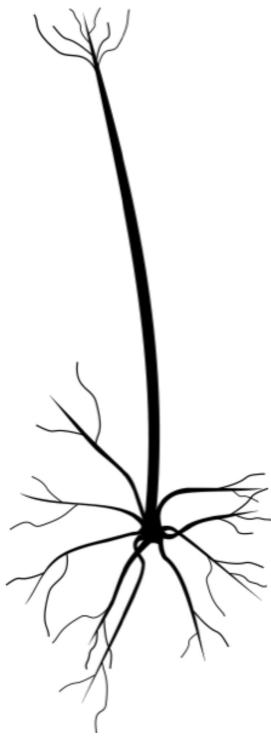
# Summary

From linear classifiers to  
fully-connected networks

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



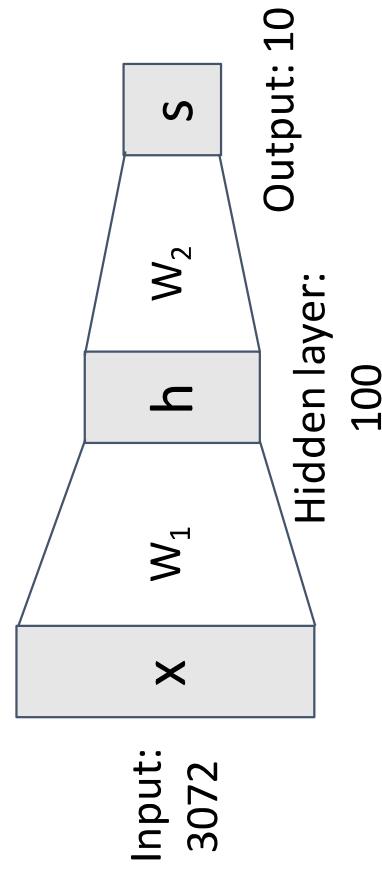
Neural networks loosely inspired by biological neurons but be careful with analogies



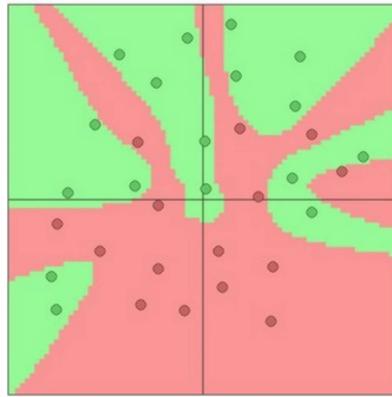
# Summary

From linear classifiers to  
fully-connected networks

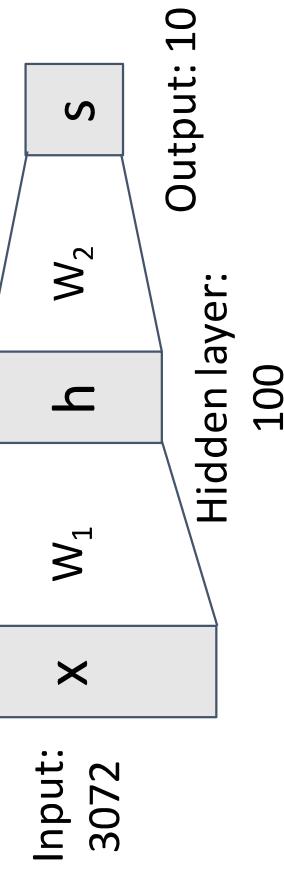
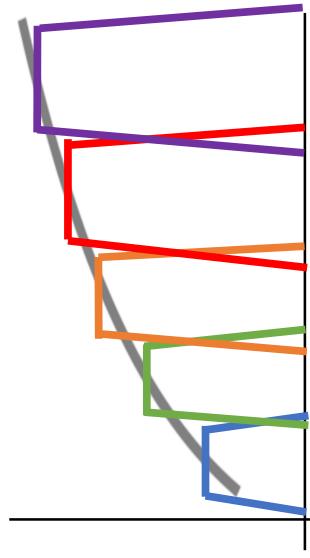
$$f = W_2 \max(0, W_1 x)$$



## Space Warping



## Universal Approximation



## Nonconvex

## Problem: How to compute gradients?

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

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

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

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

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

Next time:  
Backpropagation