# 5. Neural Networks
## GEV6135 Deep Learning for Visual Recognition and Applications

**Kibok Lee**

Assistant Professor of

Applied Statistics / Statistics and Data Science

Sep 29, 2022

# Assignment 3

- Due **Monday 10/10, 11:59pm KST**

- Training linear classifiers (Lec 3) with
  - SVM/Softmax loss (Lec 3)
  - SGD (Lec 4)

- If you feel difficult, consider to take **option 2**.

- Please read the instruction carefully!
  - Do **not write or modify any code outside** of the designated blocks.
  - Do **not add or delete cells** from the notebook.
  - Do **not import** additional libraries.
    - + Do not use torch.nn unless instructed.
  - **Run all cells**, and do **not clear out the outputs**, before submitting.
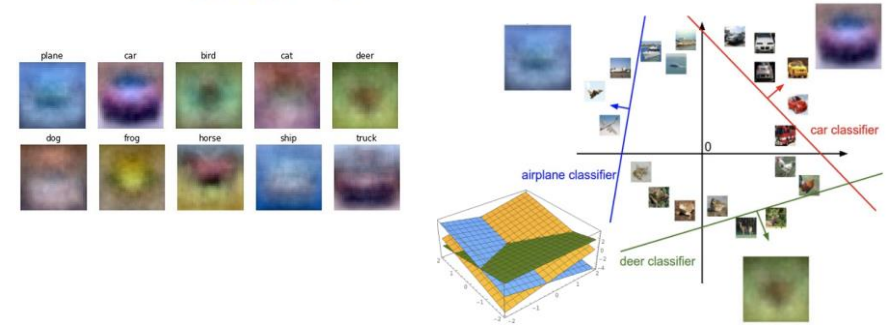  - Do **not zip by yourself**, run the provided code.

# Assignment 4

- Will be released around **Tuesday 10/4**

- Expected due **Monday 10/17**

- Training two-layer neural networks (Lec 5) with
  - Softmax loss (Lec 3)
  - SGD (Lec 4)

- A1 grading by this weekend?

# Where we are:

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

1. Use **Linear Models** for image classification problems

2. Use **Loss Functions** to express preferences over different choices of weights

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

3. Use **Regularization** to prevent overfitting to training data

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

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

4. Use **Stochastic Gradient Descent** to minimize our loss functions and train the model

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

# Problem: Linear Classifiers aren't that powerful

## Geometric Viewpoint



## Visual Viewpoint

One template per class:
Can't recognize different
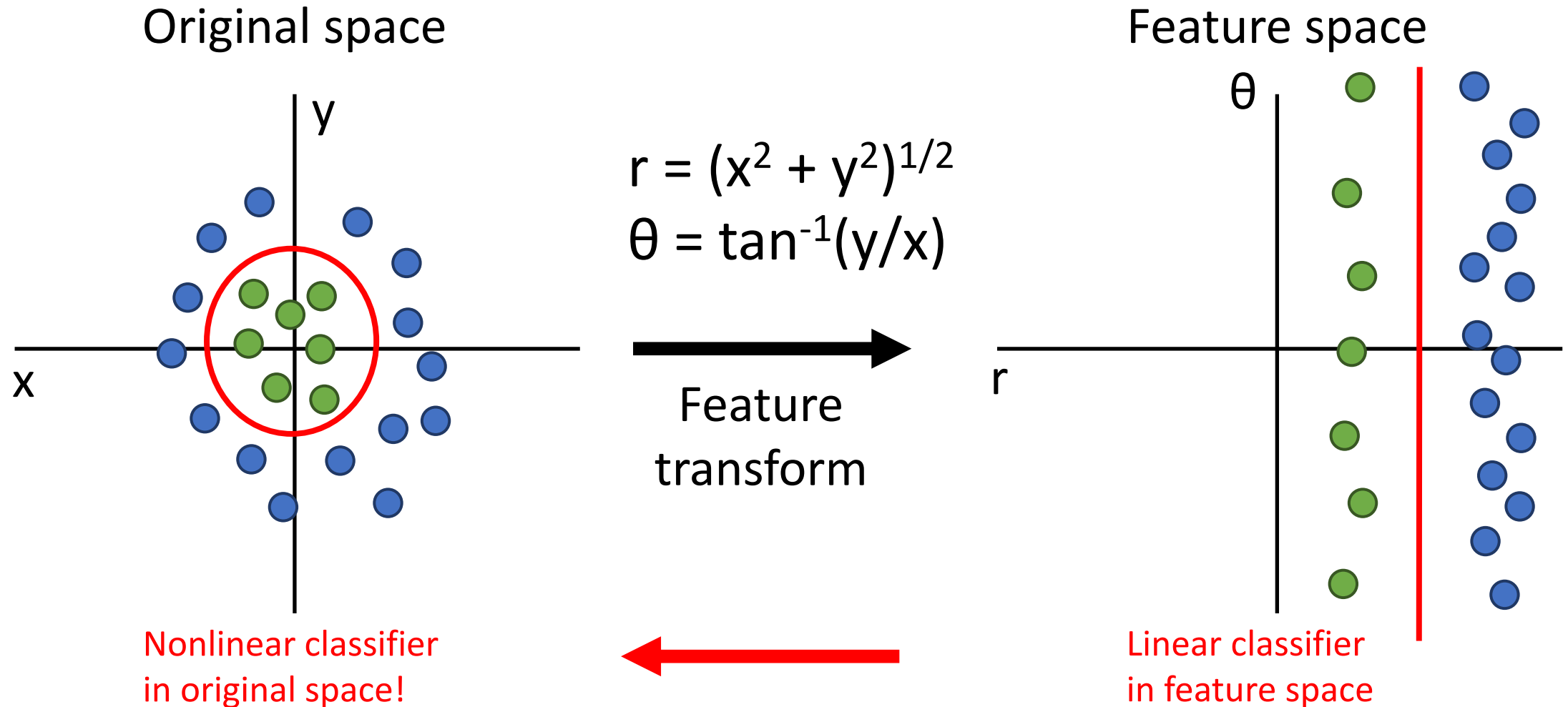modes of a class
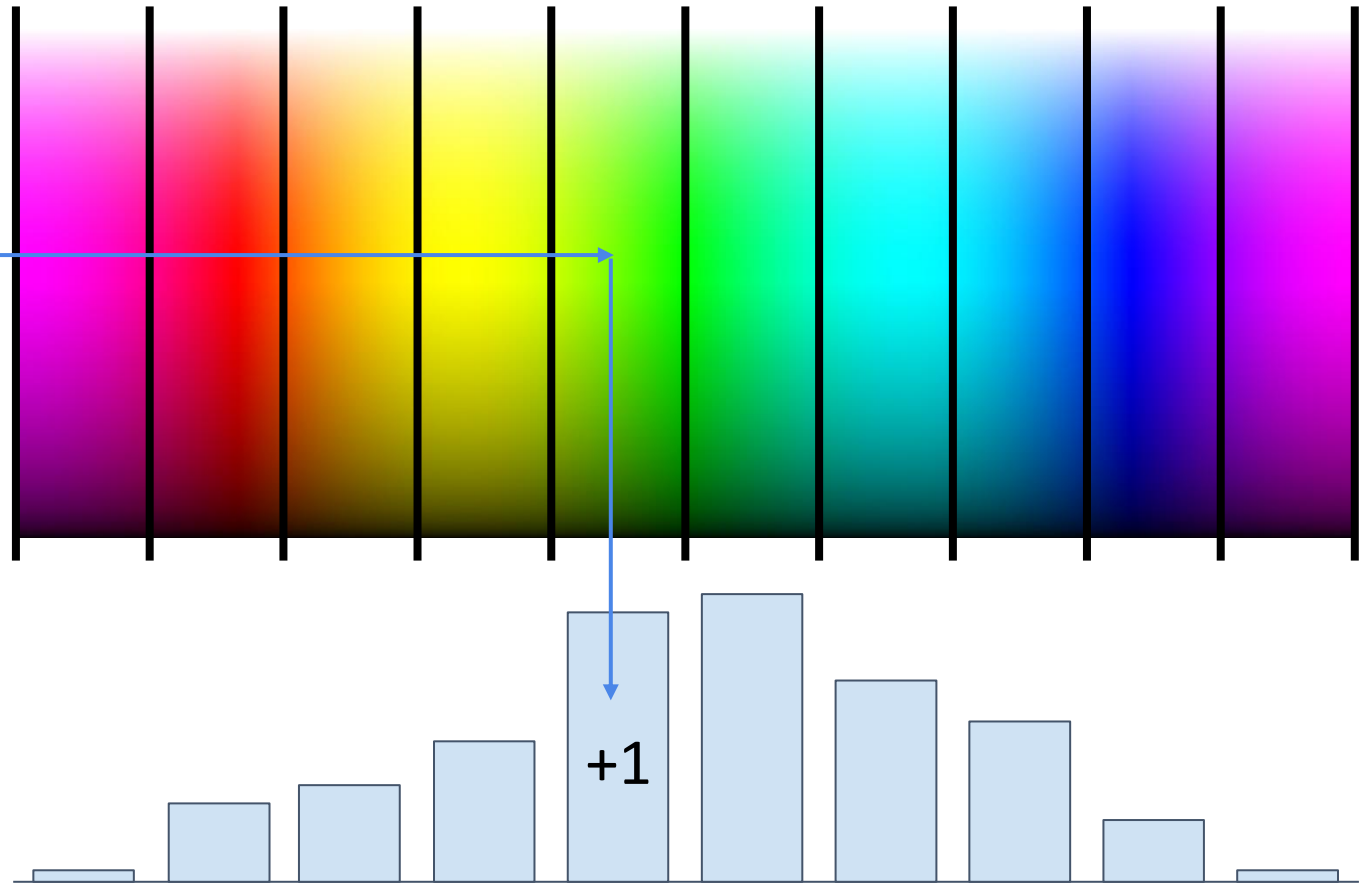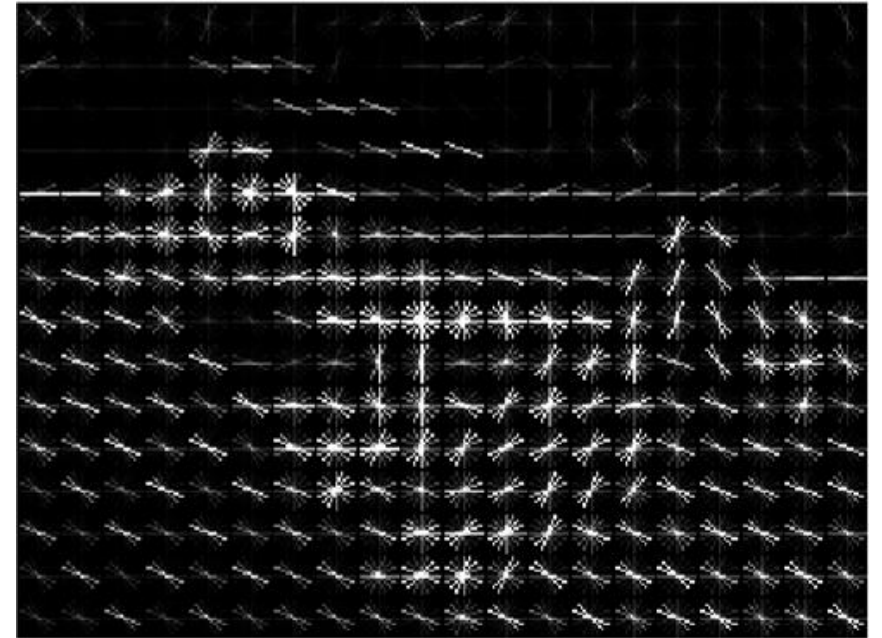
# One solution: Feature Transforms

Original space

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

Feature space

Feature transform →

← 

Nonlinear classifier
in original space!

Linear classifier
in feature space

# Image Features: Color Histogram



Ignores texture, spatial positions

+1

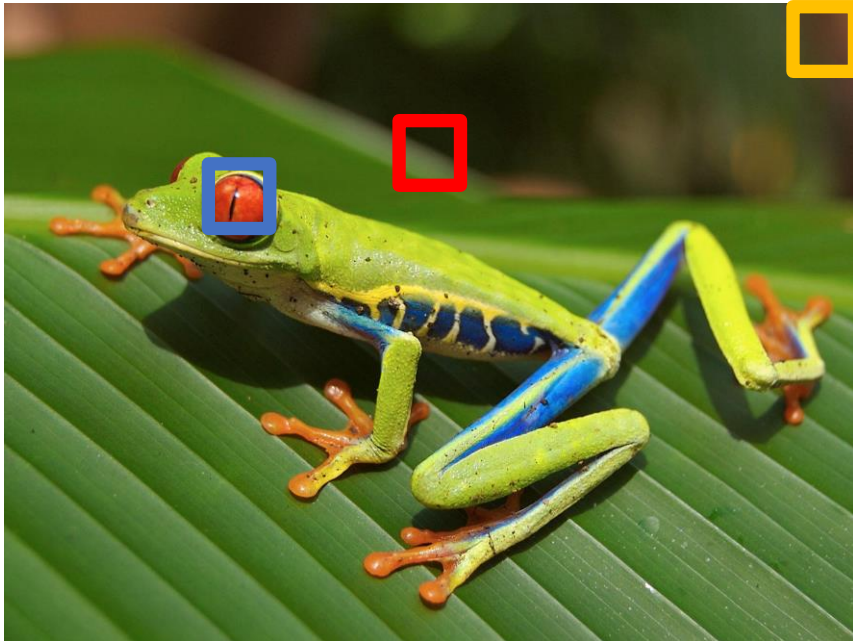# Image Features: Histogram of Oriented Gradients (HoG)



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

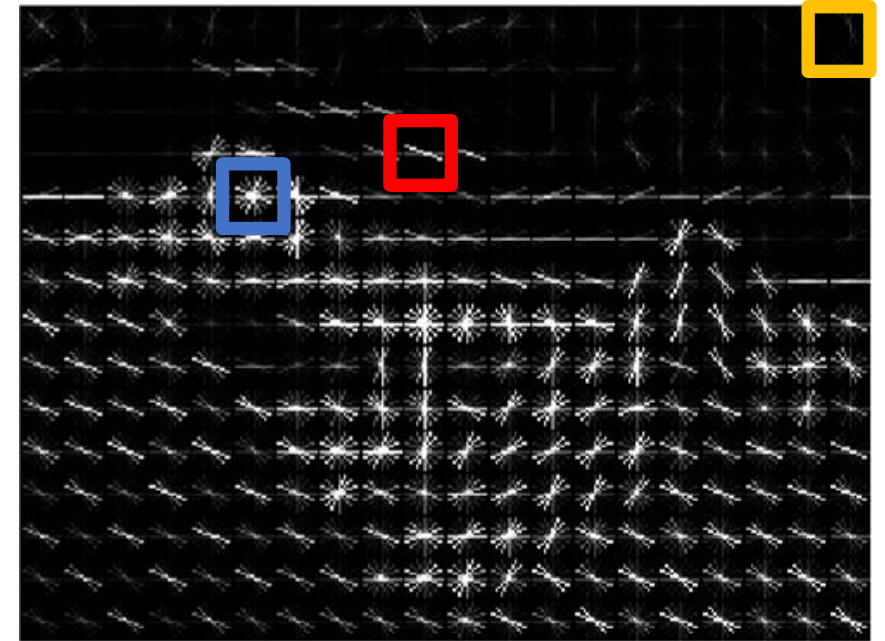# Image Features: Histogram of Oriented Gradients (HoG)



Weak edges

Strong diagonal edges

Edges in all directions

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; 8 directions per bin; feature vector has 30*40*9 = 10,800 numbers
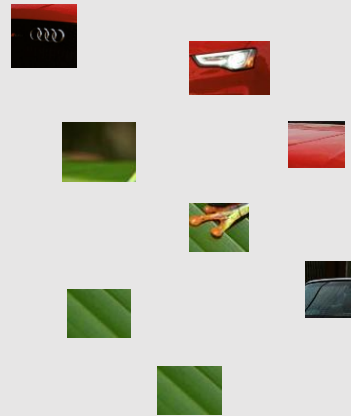
Lowe, "Object recognition from local scale-invariant features", ICCV 1999
Dalal and Triggs, "Histograms of oriented gradients for human detection," CVPR 2005
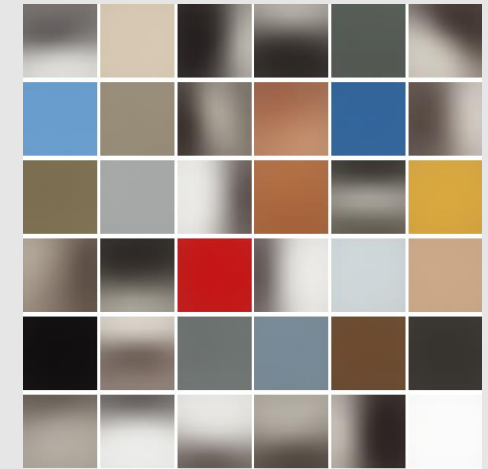
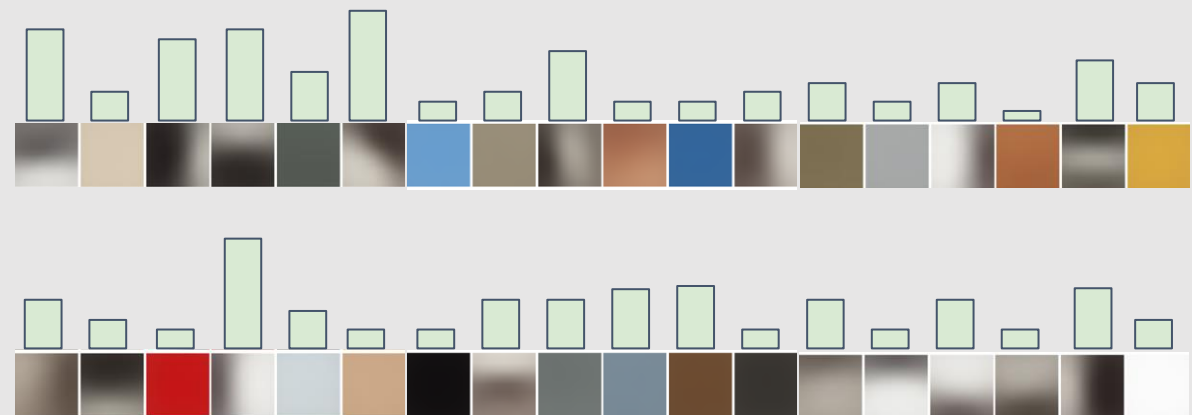# Image Features: Bag of Words (Data-Driven!)

**Step 1: Build codebook**

Extract random patches

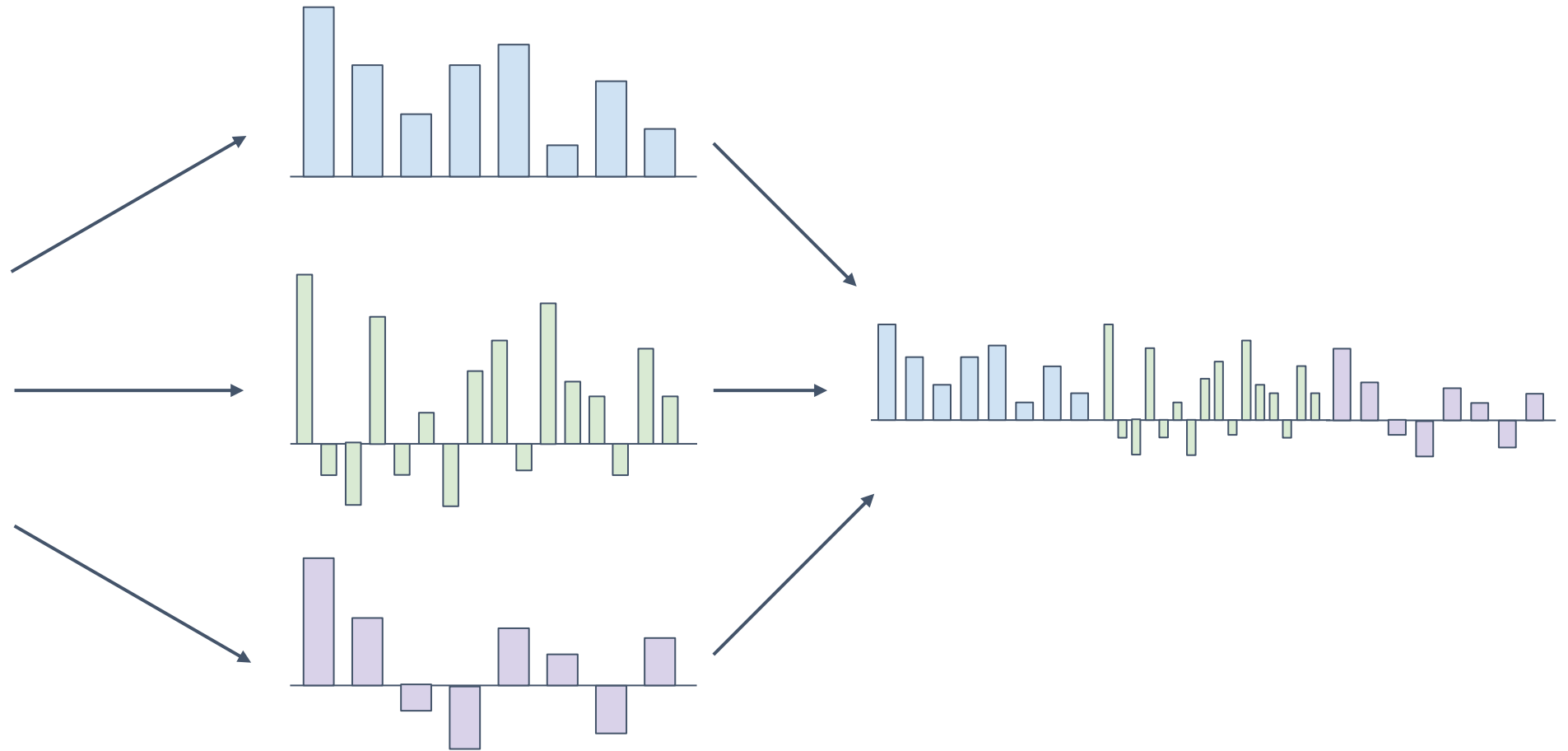Cluster patches to form "codebook" of "visual words"

**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 ≈ 10k patches per image
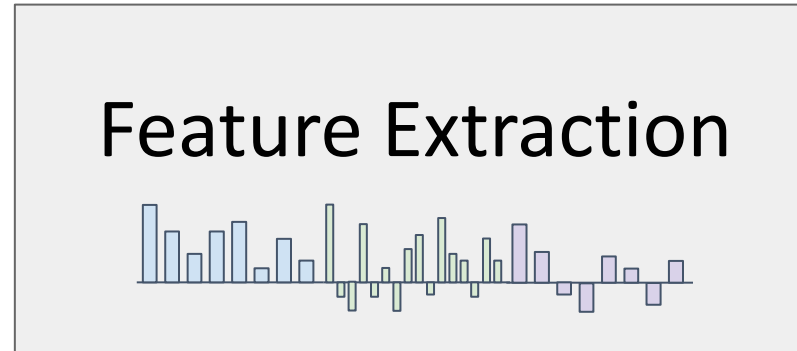- SIFT: 128-dim ⎫
- color: 96-dim ⎭ reduced to 64-dim with PCA

FV extraction and compression:
- N=1,024 Gaussians, R=4 regions ⇨ 520K dim x 2
- compression: G=8, b=1 bit per dimension

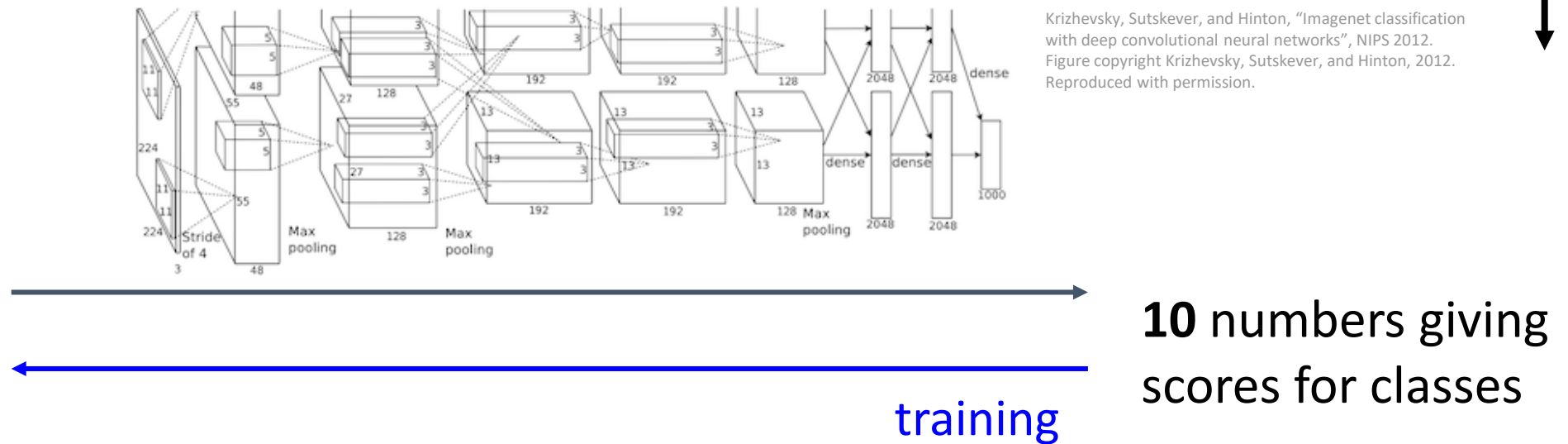One-vs-all SVM learning with SGD

Late fusion of SIFT and color systems

F. Perronnin, J. Sánchez, "Compressed Fisher vectors for LSVRC", PASCAL VOC / ImageNet workshop, ICCV, 2011.

# Image Features

# Image Features vs Neural Networks



Feature Extraction

**f**

**10** numbers giving scores for classes

training

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

**10** numbers giving scores for classes

training

# Neural Networks

Linear classifiers



This image is CC0 1.0 public domain

# Neural Networks

**Input**: $x \in \mathbb{R}^D$     **Output**: $s(x) \in \mathbb{R}^C$     **Activation function**: $f$

**Before**: Linear Classifier:    $s(x) = Wx + b$
Learnable parameters: $W \in \mathbb{R}^{C \times D}, b \in \mathbb{R}^C$

**Now:** Two-Layer Neural Network: $s(x) = W_2 \, f(W_1 x + b_1) + b_2$
Learnable parameters: $W_1 \in \mathbb{R}^{H \times D}, b_1 \in \mathbb{R}^H, W_2 \in \mathbb{R}^{C \times H}, b_2 \in \mathbb{R}^C$

# Neural Networks

**Input**: $x \in \mathbb{R}^D$    **Output**: $s(x) \in \mathbb{R}^C$    **Activation function**: $f$

**Before**: Linear Classifier:    $s(x) = Wx + b$
Learnable parameters:  $W \in \mathbb{R}^{C \times D}, b \in \mathbb{R}^C$

Feature Extraction
Linear Classifier

**Now:** Two-Layer Neural Network:  $s(x) = W_2 \, f(W_1 x + b_1) + b_2$
Learnable parameters: $W_1 \in \mathbb{R}^{H \times D}, b_1 \in \mathbb{R}^H, W_2 \in \mathbb{R}^{C \times H}, b_2 \in \mathbb{R}^C$

# Neural Networks

**Input**: $x \in \mathbb{R}^D$     **Output**: $s(x) \in \mathbb{R}^C$     **Activation function**: $f$

**Before**: Linear Classifier:   $s(x) = Wx + b$
Learnable parameters:  $W \in \mathbb{R}^{C \times D}, b \in \mathbb{R}^C$

**Now:** Two-Layer Neural Network:  $s(x) = W_2 f(W_1 x + b_1) + b_2$
Learnable parameters: $W_1 \in \mathbb{R}^{H \times D}, b_1 \in \mathbb{R}^H, W_2 \in \mathbb{R}^{C \times H}, b_2 \in \mathbb{R}^C$

Or Three-Layer Neural Network:
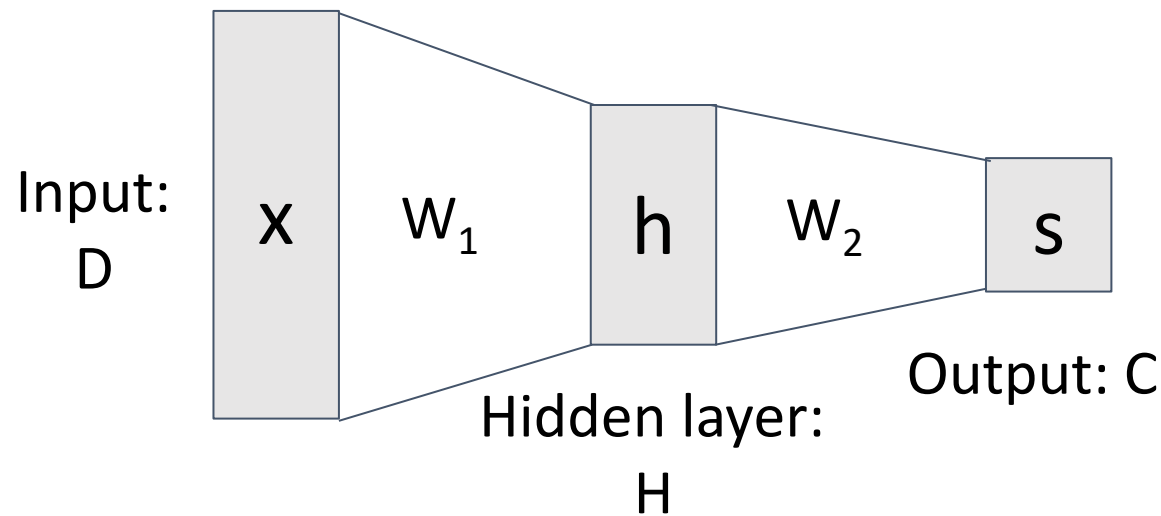$s(x) = W_3 f(W_2 f(W_1 x + b_1) + b_2) + b_3$

# Neural Networks

**Before**: Linear classifier $\qquad\qquad s(x) = Wx + b$

**Now**: 2-layer Neural Network $\quad s(x) = W_2\, f(W_1 x + b_1) + b_2$

Input:
D

x $\quad$ W$_1$ $\quad$ h $\quad$ W$_2$ $\quad$ s

Hidden layer:
H

Output: C

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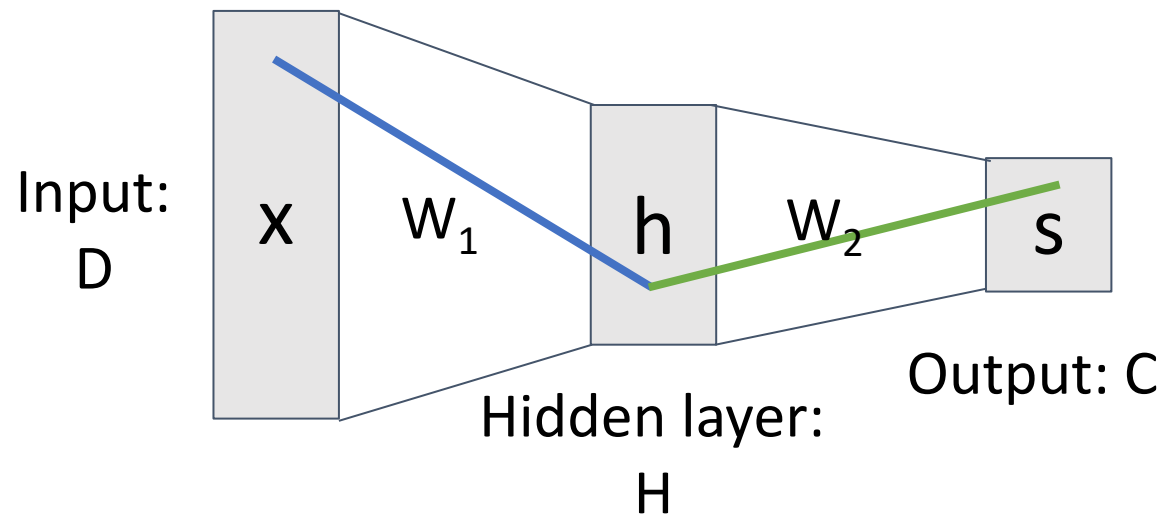
# Neural Networks

**Before**: Linear classifier $\qquad s(x) = Wx + b$

**Now**: 2-layer Neural Network $\quad s(x) = W_2\, f(W_1 x + b_1) + b_2$

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

Input: D

x    $W_1$    h    $W_2$    s

Hidden layer: H

Output: C

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
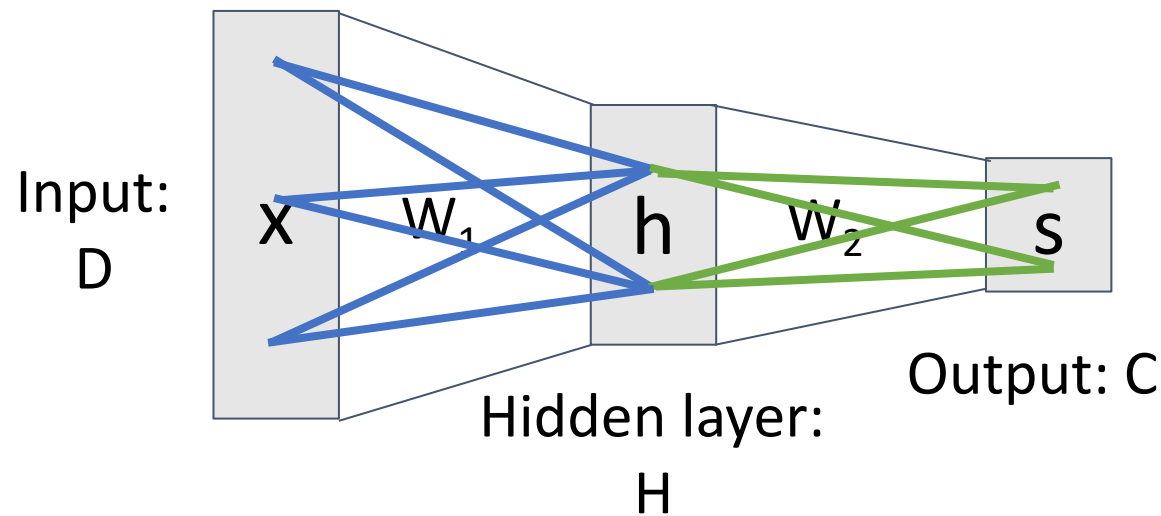
**Before**: Linear classifier $\qquad s(x) = Wx + b$

**Now**: 2-layer Neural Network $\quad s(x) = W_2\, f(W_1 x + b_1) + b_2$

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

All elements of x affect all elements of h

Input: D

x $\quad$ $W_1$ $\quad$ h $\quad$ $W_2$ $\quad$ s

Hidden layer: H

Output: C

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
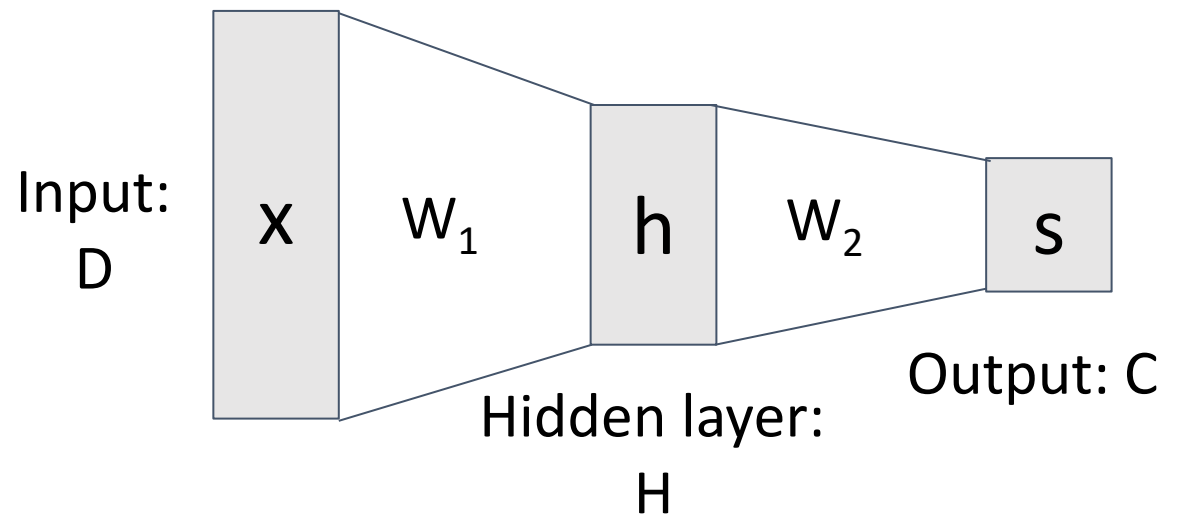Also "Multi-Layer Perceptron" (MLP)

# Neural Networks

Linear classifier: One template per class

| | | | | |
|---|---|---|---|---|
| plane | car | bird | cat | deer |
| dog | frog | horse | ship | truck |

(**Before**) Linear score function:

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

Input:
D

x   $W_1$   h   $W_2$   s

Hidden layer:
H

Output: C

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

# Neural Networks

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



(**Before**) Linear score function:

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

Input:
D

x

$W_1$

h

$W_2$

s

Hidden layer:
H

Output: C

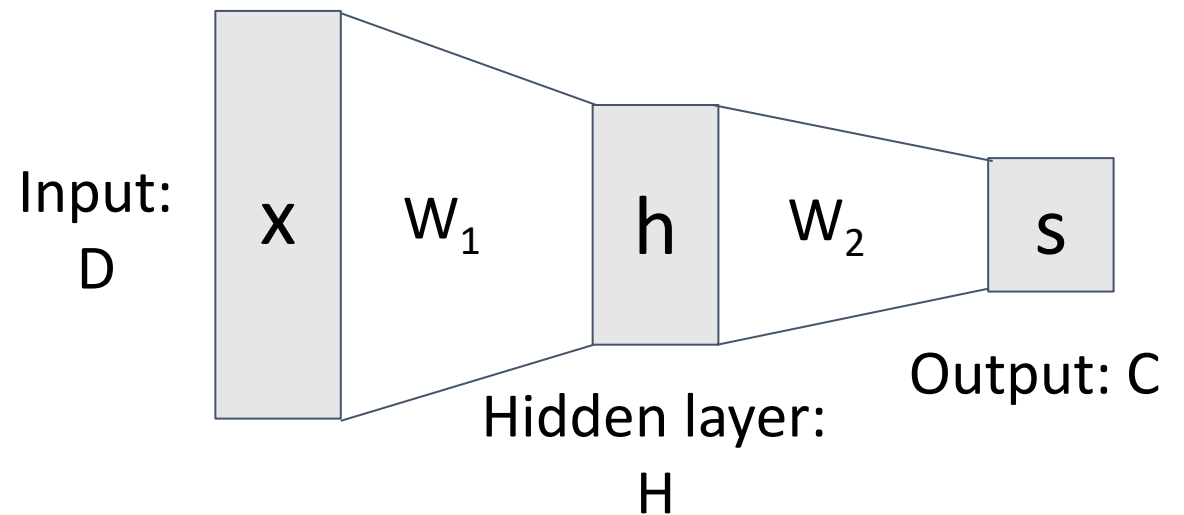$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$

# Neural Networks

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



(**Before**) Linear score function:

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



Input: D

x    $W_1$    h    $W_2$    s

Hidden layer: H

Output: C

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

# Neural Networks

"Distributed representation":
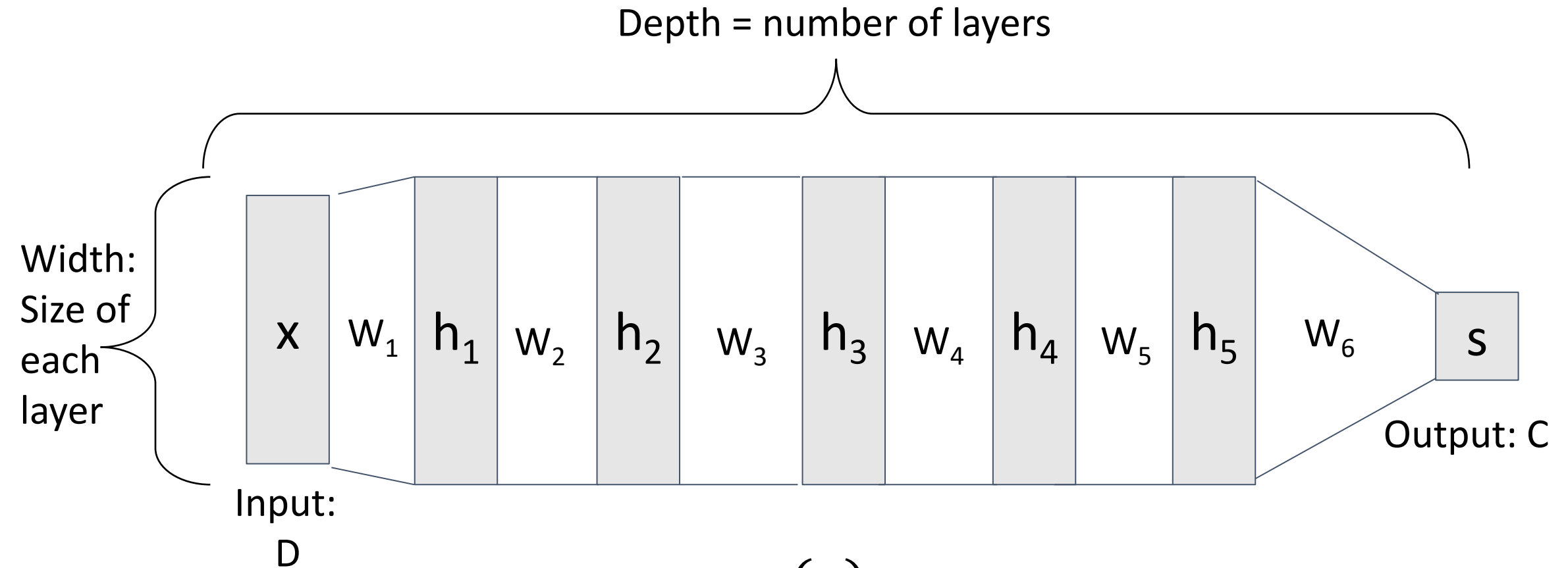Most templates not interpretable!



(**Before**) Linear score function:

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



Input: D

x   $W_1$   h   $W_2$   s

Hidden layer: H

Output: C

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

# Deep Neural Networks



Depth = number of layers

Width: Size of each layer

Input: D

Output: C

$$s(x) =$$
$$W_6\,f(W_5\,f(W_4\,f(W_3\,f(W_2\,f(W_1 x + b_1) + b_2) + b_3) + b_4) + b_5) + b_6$$

# Activation Functions

## 2-layer Neural Network

$$s(x) = W_2 \boxed{f} (W_1 x + b_1) + b_2$$

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?

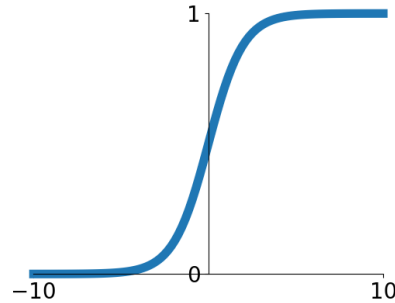$$s(x) = W_2 (W_1 x + b_1) + b_2$$

# Activation Functions

## 2-layer Neural Network

$$s(x) = W_2 \boxed{f}(W_1 x + b_1) + b_2$$

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?

$$s(x) = W_2(W_1 x + b_1) + b_2$$
$$= (W_1 W_2)x + (W_2 b_1 + b_2)$$

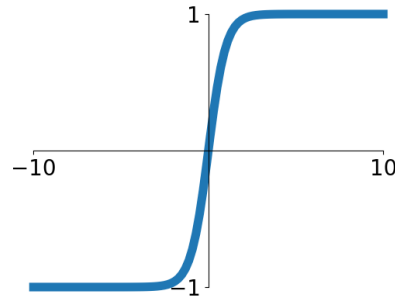**A**: We end up with a linear classifier!

# Activation Functions
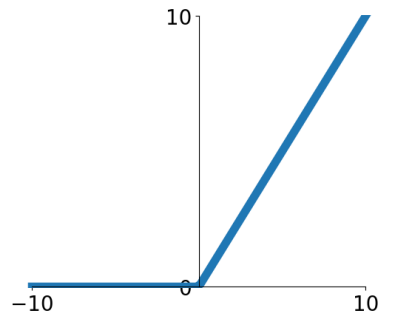
**Sigmoid**

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

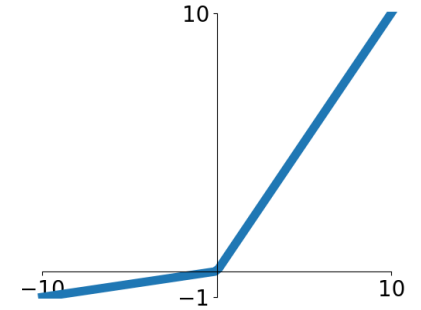**tanh**

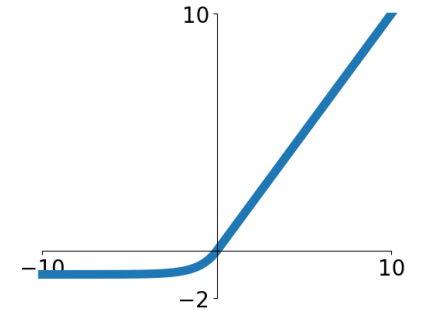$$\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1}$$

**ReLU**

$$\max(0, x)$$

**Leaky ReLU**

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

**ELU**

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

**GELU**

$$= 0.5x\left[1 + \mathrm{erf}(x/\sqrt{2})\right]$$
$$\approx x\sigma(1.702x)$$

# Neural Net in <20 lines!



input layer

hidden layer

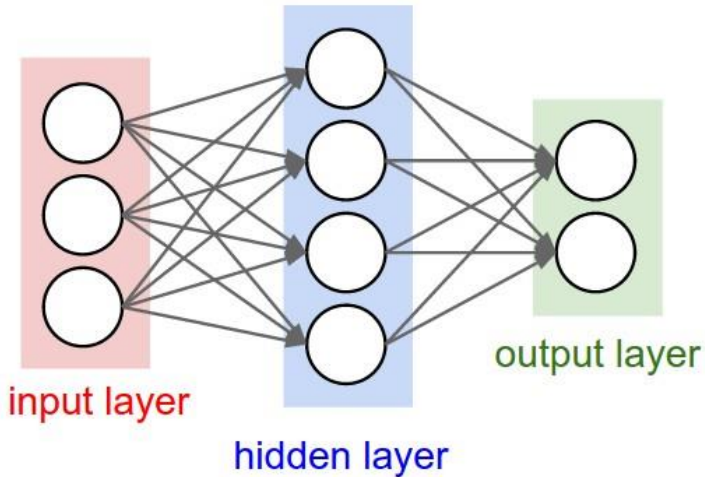output layer

```
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)

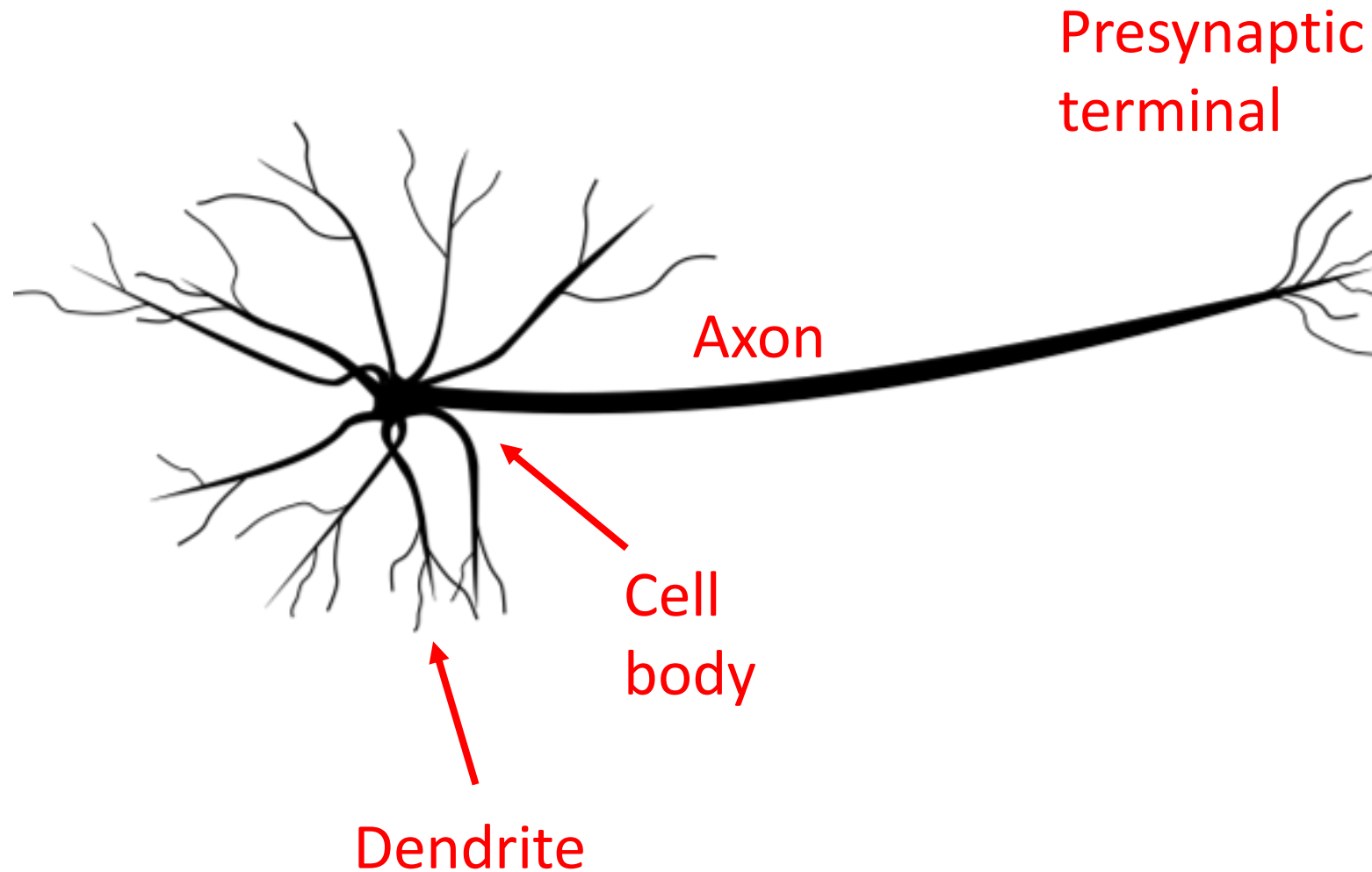Compute gradients

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
```
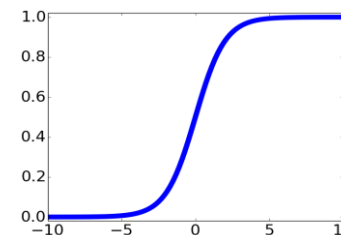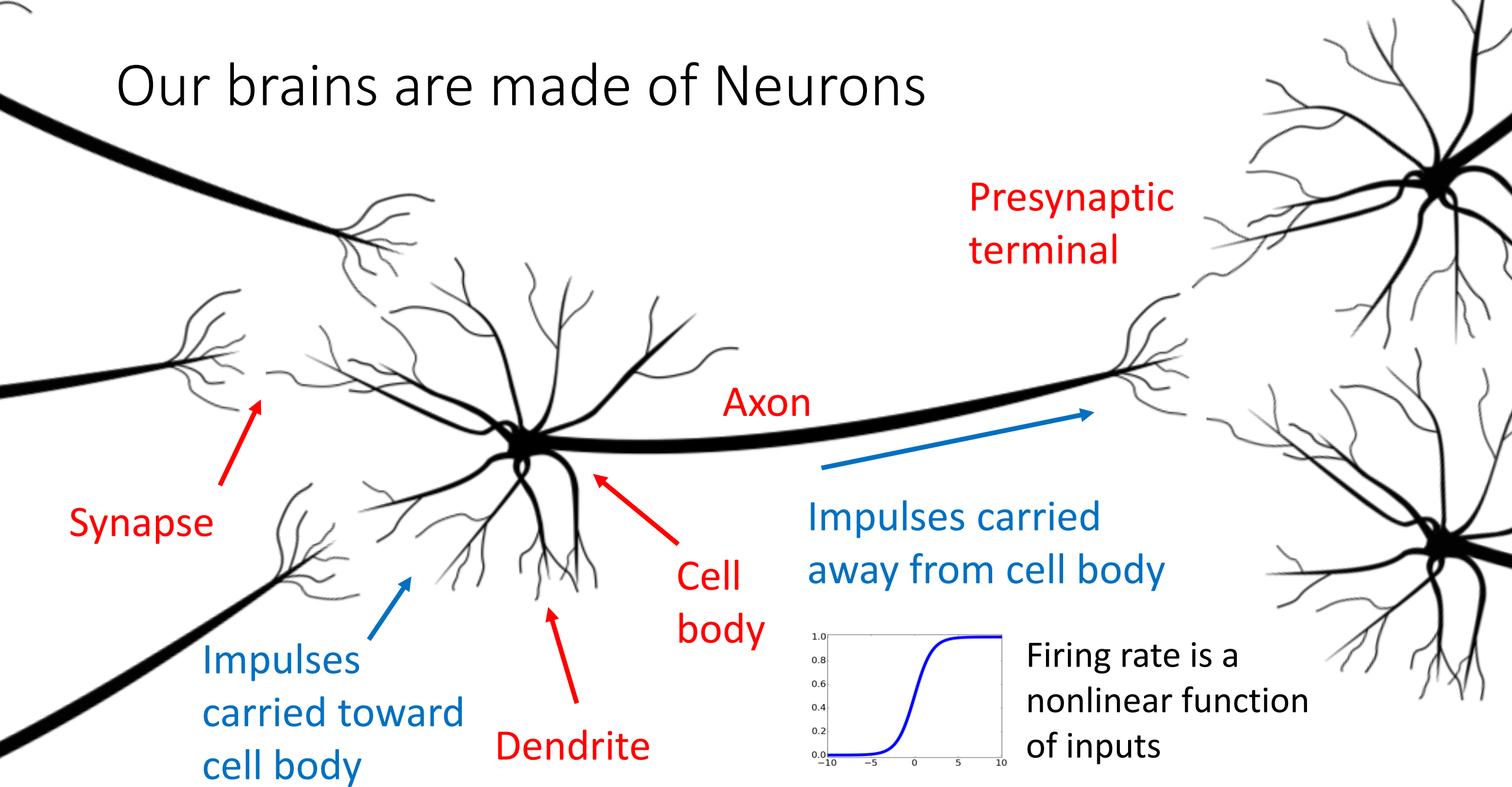
# Attendance Check

# Our brains are made of Neurons

Presynaptic terminal

Axon

Cell body

Dendrite

# Our brains are made of Neurons

Presynaptic terminal

Axon

Synapse

Impulses carried away from cell body

Cell body

Impulses carried toward cell body

Dendrite

Firing rate is a nonlinear function of inputs

# Biological Neuron

dendrite

presynaptic terminal

axon

cell body

# Artificial Neuron



input layer

hidden layer 1    hidden layer 2

output layer

$x_0$    $w_0$

axon from a neuron

synapse

$w_0 x_0$

dendrite

$w_1 x_1$

cell body

$\sum_i w_i x_i + b$    $f$

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

output axon

activation function

$w_2 x_2$

Biological Neurons:
Complex connectivity patterns

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



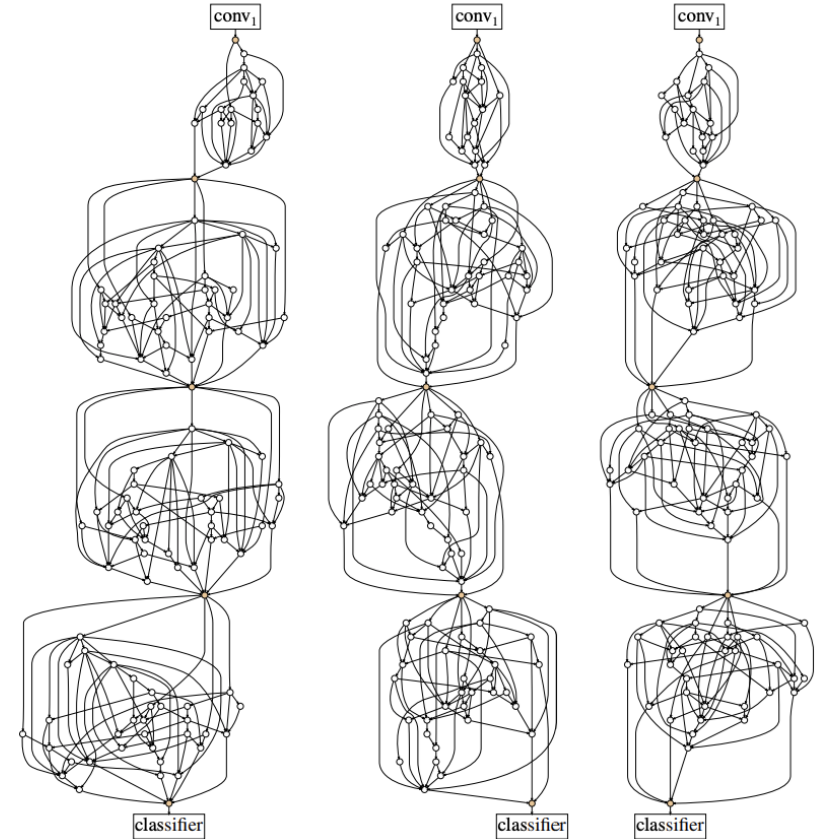This image is CC0 Public Domain

# 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
- Abstracting a neuron by "firing rate" isn't enough; temporal sequences of activations matter too (spiking neural networks)

[Dendritic Computation. London and Hausser]

# Space Warping

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



Feature transform:
h = Wx

# Space Warping

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

Points not linearly
separable in original space



Feature transform:
h = Wx

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



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

B is "collapsed" onto +h2 axis

C "collapsed" onto origin

D "collapsed" onto +h1 axis

# 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:
$$h = Wx$$

# Space Warping

Points not linearly
separable in original space

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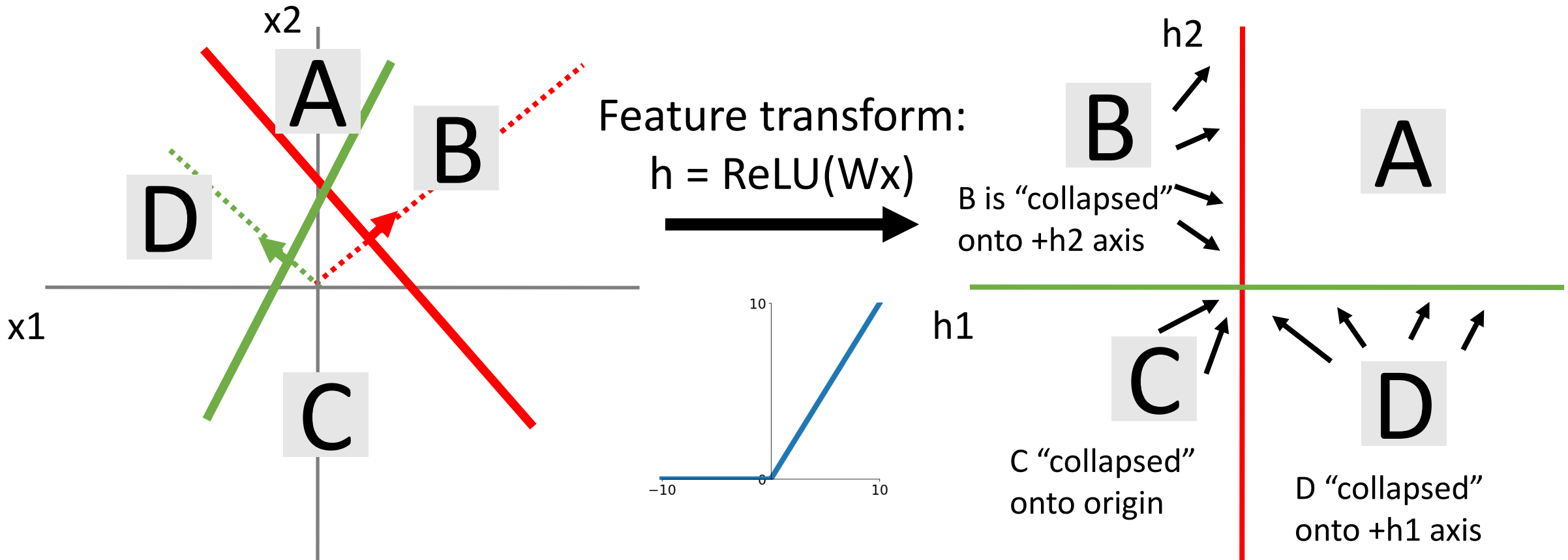h = ReLU(Wx) = max(0, Wx)
Where x, h are both 2-dimensional

Feature transform:
h = ReLU(Wx)

Linear classifier in feature
space gives nonlinear
classifier in original space

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 $f: R^N \rightarrow R^M$ with arbitrary precision*

*Many technical conditions: Only holds on compact subsets of $R^N$; function must be continuous; need to define "arbitrary precision"; etc

# Universal Approximation

Example: Approximating a function f: R -> R with a two-layer ReLU network



Input:
x (1,)

Output:
y (1,)

First layer weights: w (3,1)
First layer bias: b (3,)

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

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

Output is a sum of shifted, scaled ReLUs:

Flip left / right based on sign of $w_i$

Slope is given by $u_i * w_i$

Position of "bend" given by $b_i$

# Universal Approximation

Example: Approximating a function f: R -> R with a two-layer ReLU network



Input:
x (1,)

First layer weights: w (3,1)
First layer bias: b (3,)

Output:
y (1,)

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

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

$m_1 = t / (s_2 - s_1)$
$m_2 = t / (s_4 - s_3)$

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: R -> R with a two-layer ReLU network



Input:
x (1,)

$w_1$  $h_1$  $u_1$

$w_2$  $h_2$  $u_2$

$w_3$  $h_3$  $u_3$

y  Output:
y (1,)

First layer weights: w (3,1)
First layer bias: b (3,)

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

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

$m_1 = t / (s_2 - s_1)$
$m_2 = t / (s_4 - s_3)$

$m_1 * max(0, x - s_1)$

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: R -> R with a two-layer ReLU network

Input:
x (1,)

$w_1$ → $h_1$ → $u_1$

$w_2$ → $h_2$ → $u_2$

$w_3$ → $h_3$ → $u_3$

x → y

Output:
y (1,)

First layer weights: w (3,1)
First layer bias: b (3,)

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

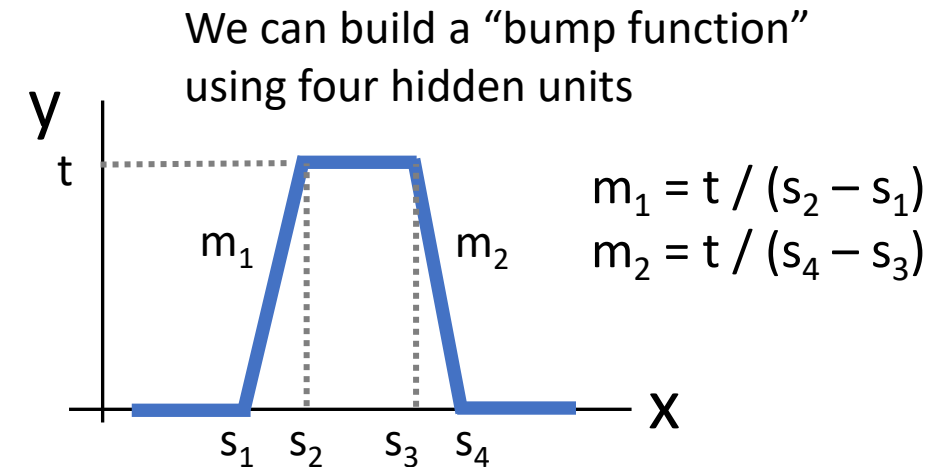h1 = max(0, w1 * x + b1)
h2 = max(0, w2 * x + b2)
h3 = max(0, w3 * x + b3)
y = u1 * h1 + u2 * h2 + u3 * h3 + p

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

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

y

t

$m_1$

$m_2$

$s_1$  $s_2$  $s_3$  $s_4$

x

$m_1 = t / (s_2 - s_1)$
$m_2 = t / (s_4 - s_3)$

$m_1 * max(0, x - s_1)$

$-m_1 * max(0, x - s_2)$

# Universal Approximation

Example: Approximating a function f: R -> R with a two-layer ReLU network



Input:
x (1,)

Output:
y (1,)

First layer weights: w (3,1)
First layer bias: b (3,)

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

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

$m_1 = t / (s_2 - s_1)$
$m_2 = t / (s_4 - s_3)$
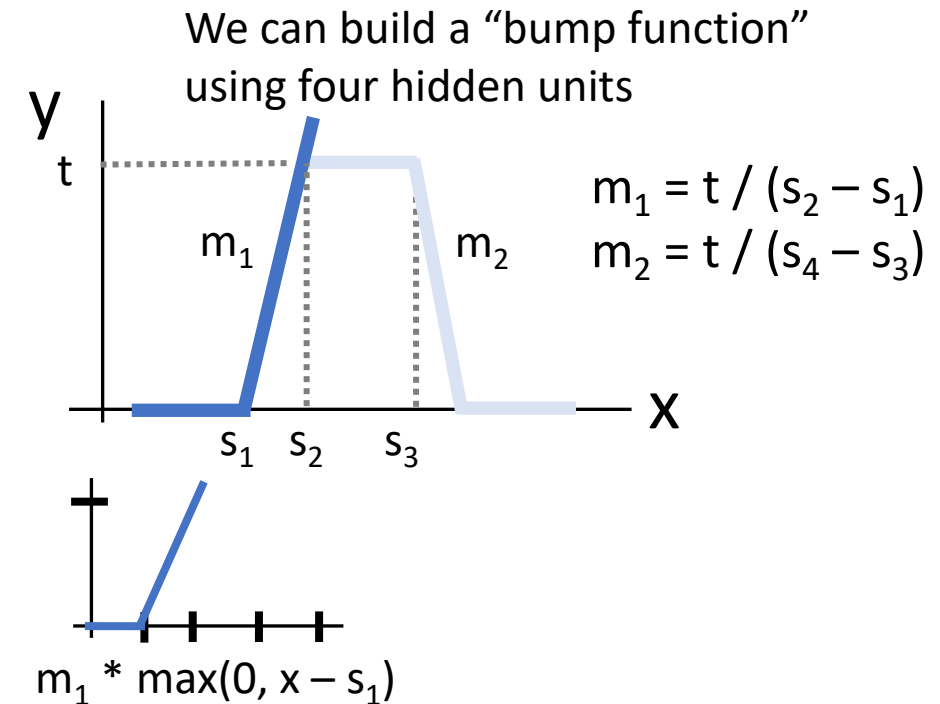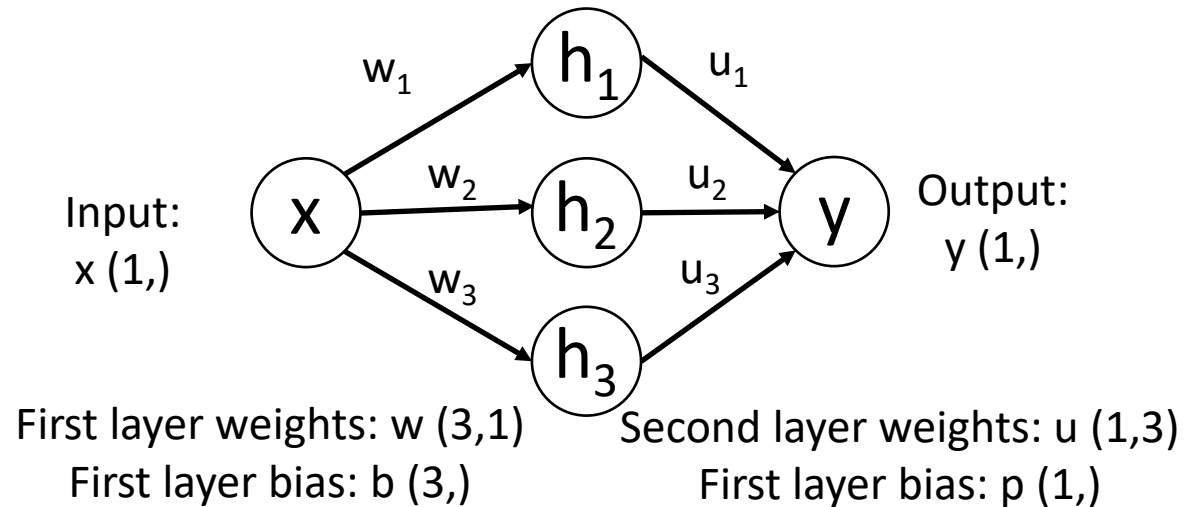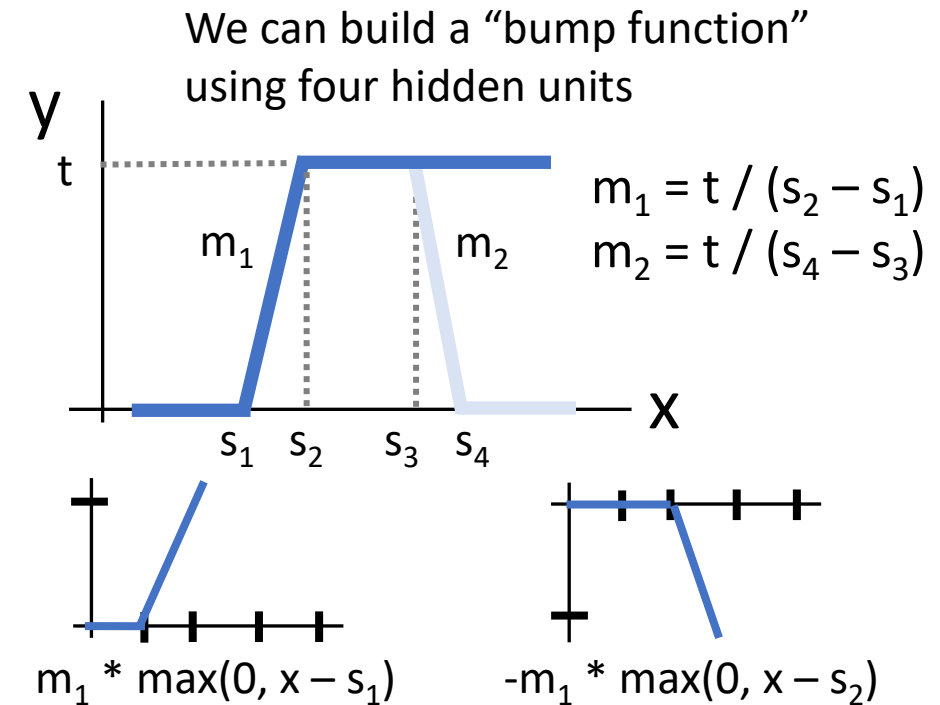
h1 = max(0, w1 * x + b1)
h2 = max(0, w2 * x + b2)
h3 = max(0, w3 * x + b3)
y = u1 * h1 + u2 * h2 + u3 * h3 + p

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

$m_1$ * max(0, x – $s_1$)

-$m_1$ * max(0, x – $s_2$)

-$m_2$ * max(0, x – $s_3$)

# Universal Approximation

Example: Approximating a function f: R -> R with a two-layer ReLU network



Input:
x (1,)

Output:
y (1,)

First layer weights: w (3,1)
First layer bias: b (3,)

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

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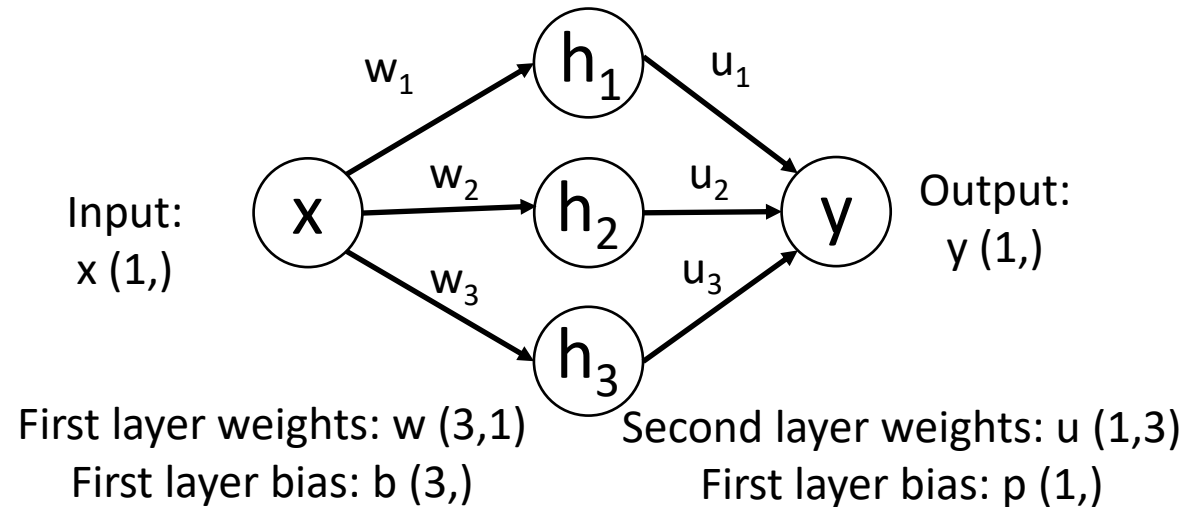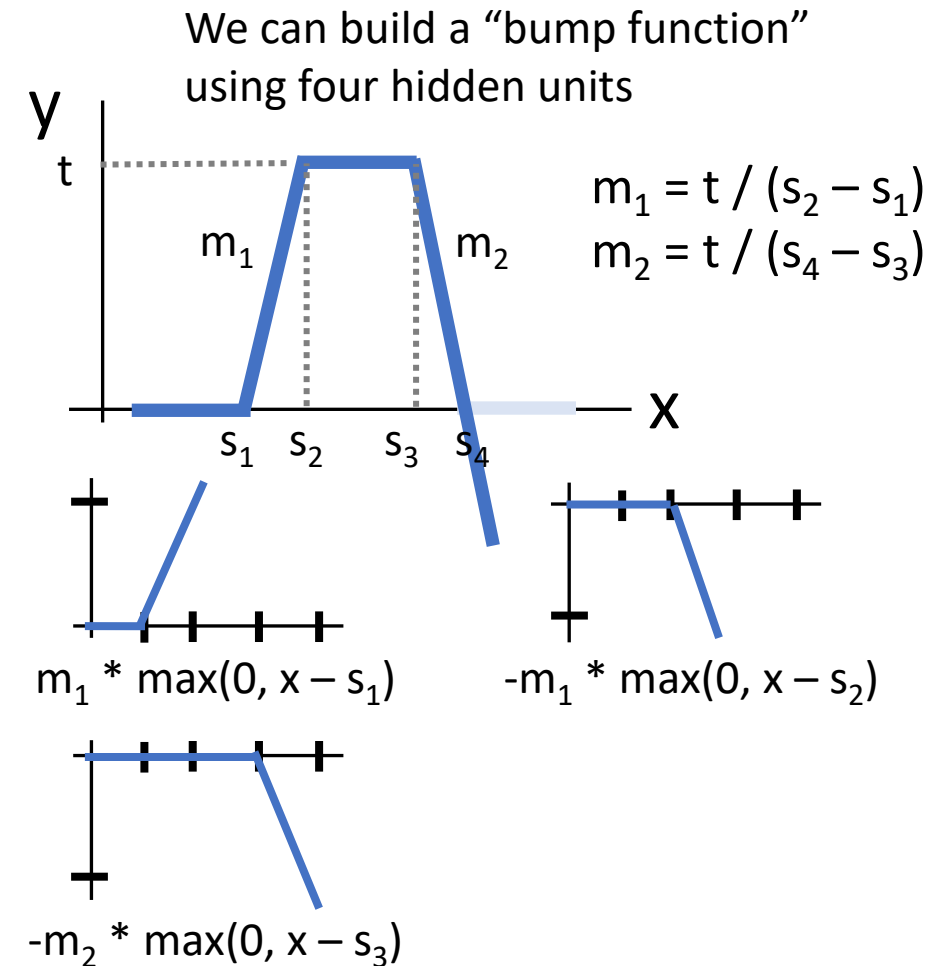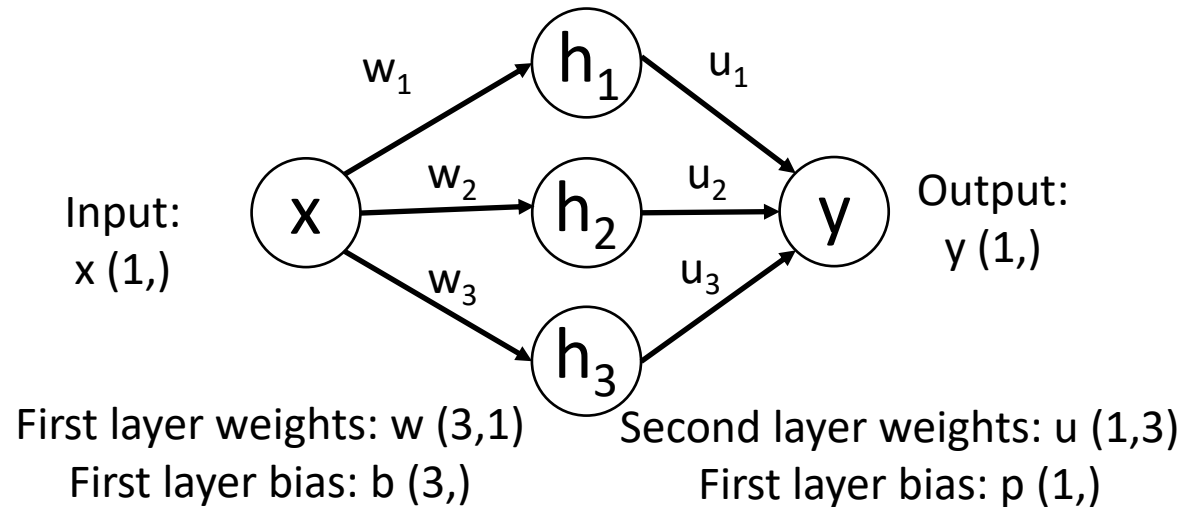y = u1 * h1 + u2 * h2 + u3 * h3 + p

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

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

$m_1 = t / (s_2 - s_1)$
$m_2 = t / (s_4 - s_3)$

$m_1 * \max(0, x - s_1)$

$-m_1 * \max(0, x - s_2)$

$-m_2 * \max(0, x - s_3)$

$m_2 * \max(0, x - s_4)$

# Universal Approximation

Example: Approximating a function f: R -> R with a two-layer ReLU network



Input:
x (1,)

$w_1$  $h_1$  $u_1$

$w_2$  $h_2$  $u_2$

$w_3$  $h_3$  $u_3$

y

Output:
y (1,)

First layer weights: w (3,1)
First layer bias: b (3,)

Second layer weights: u (1,3)
First layer bias: p (1,)
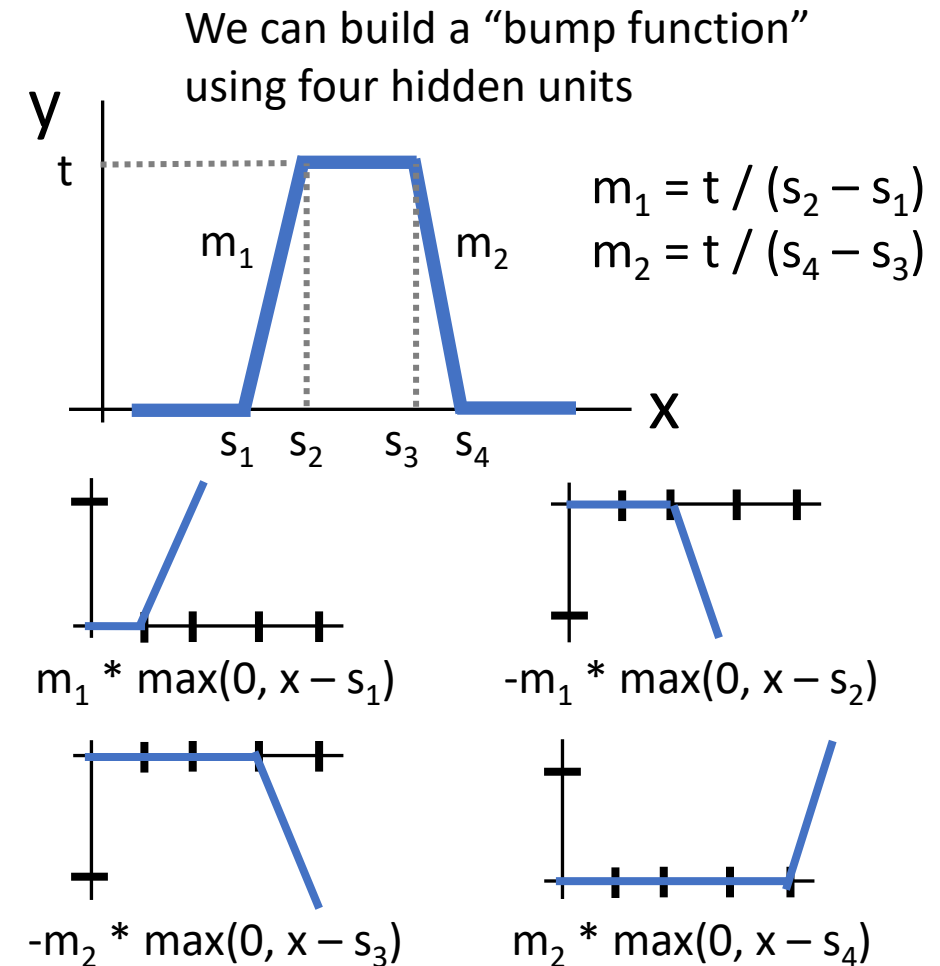
h1 = max(0, w1 * x + b1)
h2 = max(0, w2 * x + b2)
h3 = max(0, w3 * x + b3)
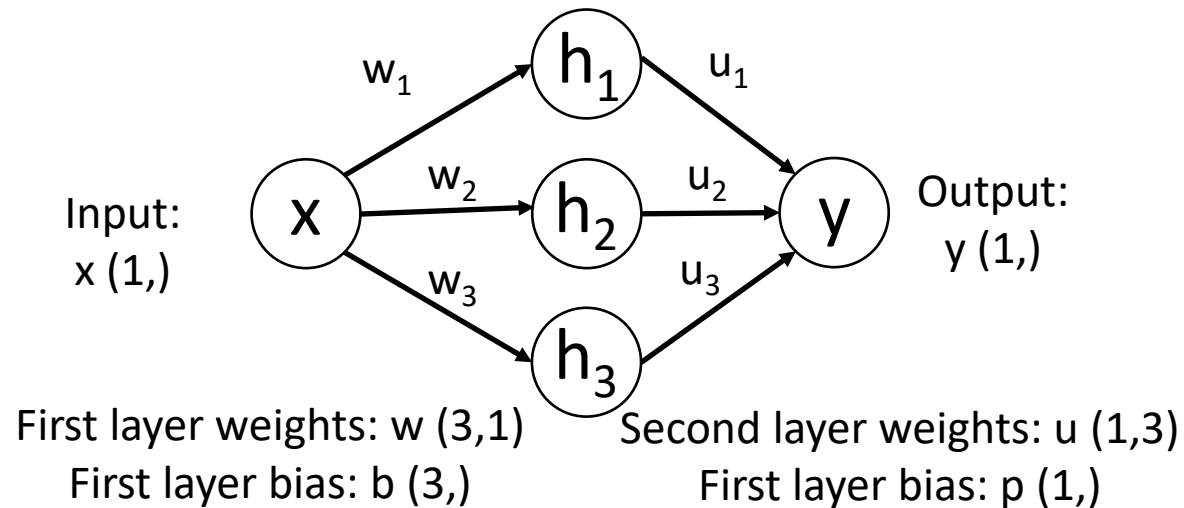y = u1 * h1 + u2 * h2 + u3 * h3 + p

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

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: R -> R with a two-layer ReLU network



Input:
x (1,)

Output:
y (1,)

First layer weights: w (3,1)
First layer bias: b (3,)

Second layer weights: u (1,3)
First layer bias: p (1,)
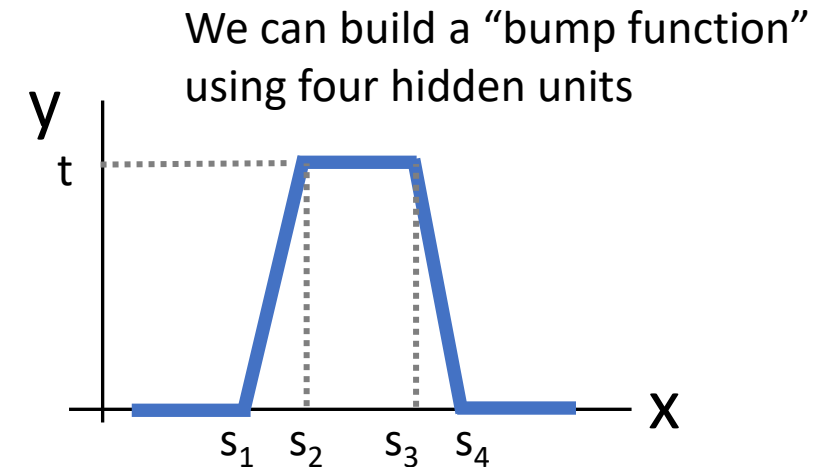
h1 = max(0, w1 * x + b1)
h2 = max(0, w2 * x + b2)
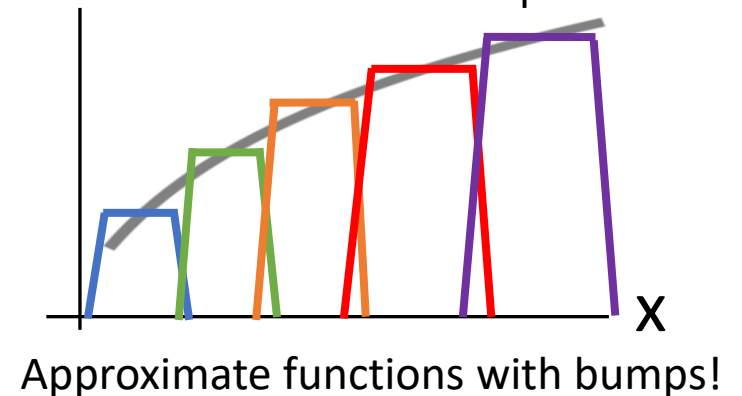h3 = max(0, w3 * x + b3)
y = u1 * h1 + u2 * h2 + u3 * h3 + p

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

What about…
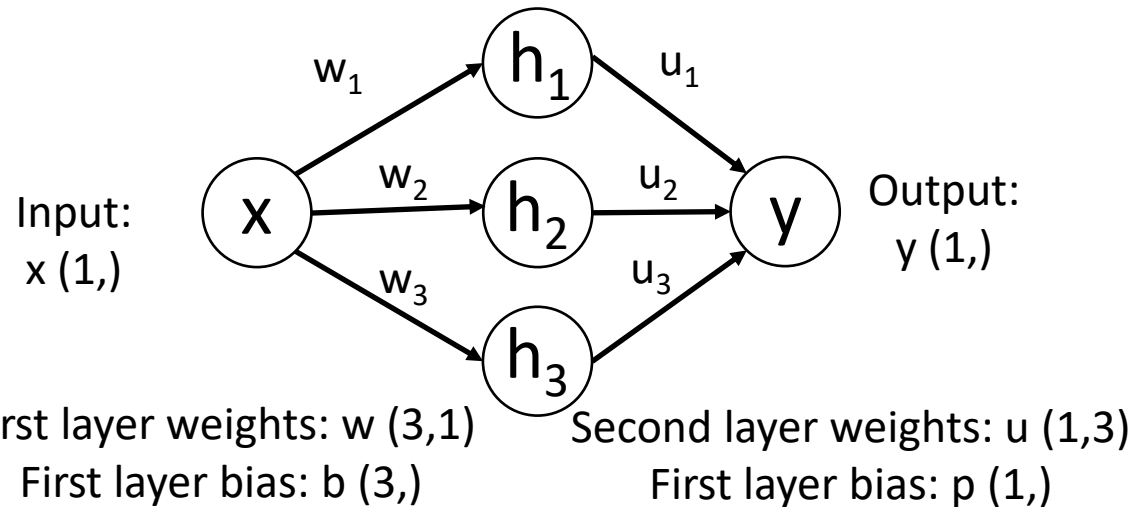- Gaps between bumps?
- Other nonlinearities?
- Higher-dimensional functions?

See Nielsen, Chapter 4



Approximate functions with bumps!

# Universal Approximation

Example: Approximating a function f: R -> R with a two-layer ReLU network



Input:
x (1,)

Output:
y (1,)

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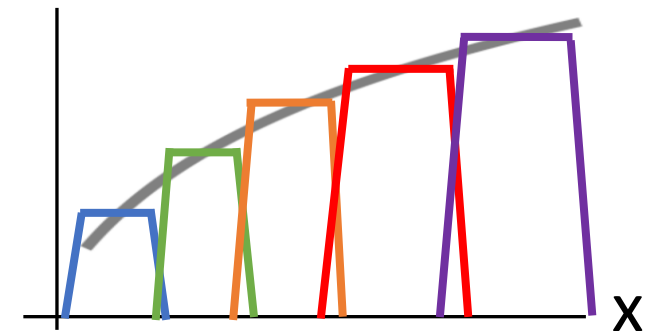
h1 = max(0, w1 * x + b1)
h2 = max(0, w2 * x + b2)
h3 = max(0, w3 * x + b3)
y = u1 * h1 + u2 * h2 + u3 * h3 + p

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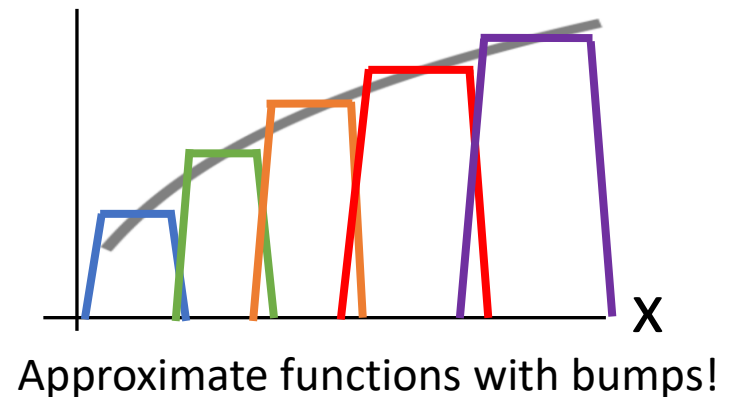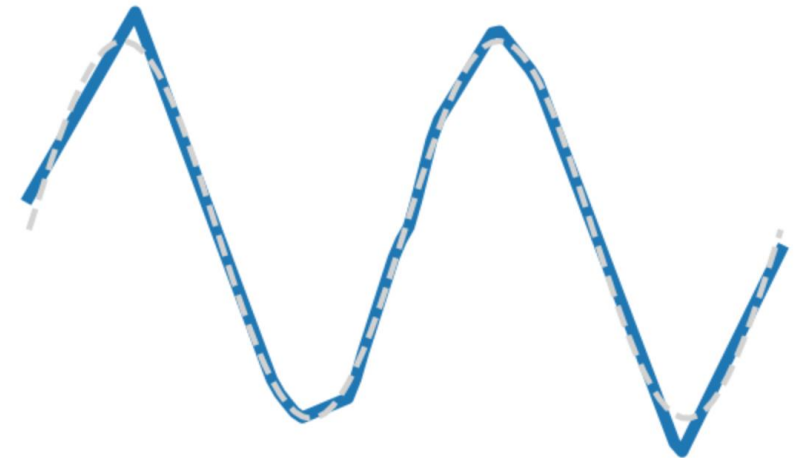  + p

Approximate functions with bumps!

# Universal Approximation

Example: Approximating a function f: R -> R with a two-layer ReLU network

Input:
x (1,)

$w_1$ $\rightarrow$ $h_1$ $\rightarrow$ $u_1$

$w_2$ $\rightarrow$ $h_2$ $\rightarrow$ $u_2$

$w_3$ $\rightarrow$ $h_3$ $\rightarrow$ $u_3$

$\rightarrow$ y

Output:
y (1,)

Reality check: Networks don't really learn bumps!

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!

x

Approximate functions with bumps!

# Extra topic: Convex Functions

# Convex Functions

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

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

Example: $f(x) = \cos(x)$
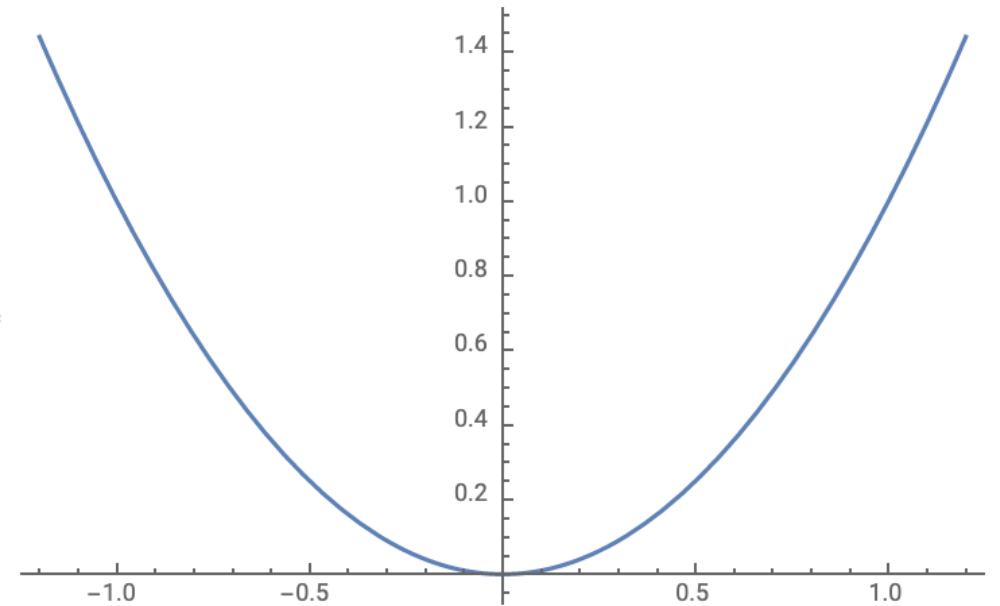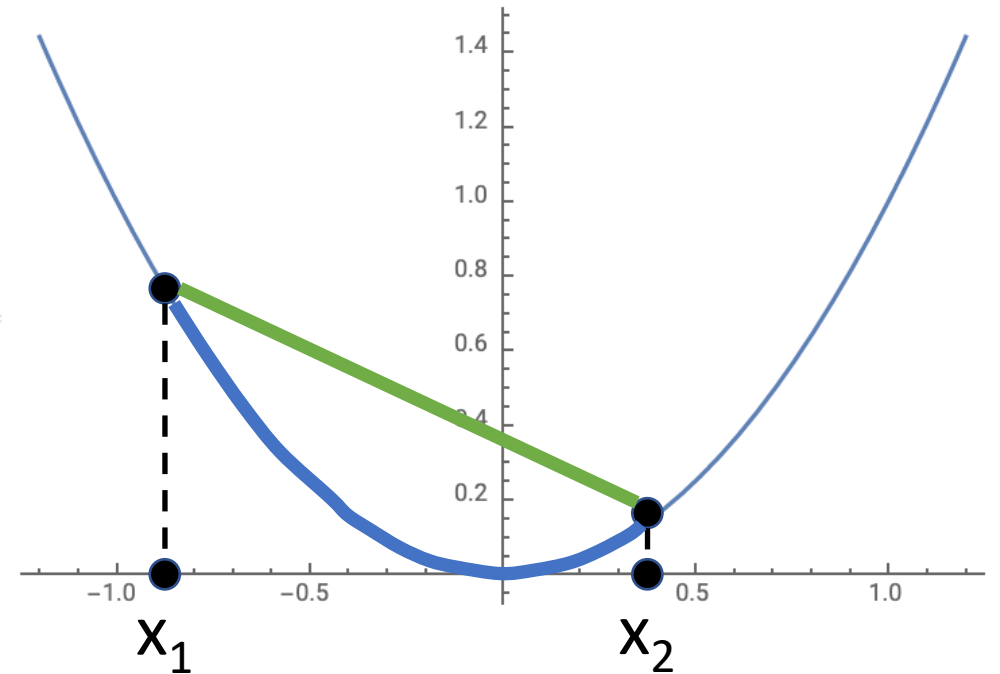is <u>not</u> convex:

# Convex Functions

A function $f : X \subseteq \mathbb{R}^N \to \mathbb{R}$ 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**[*]
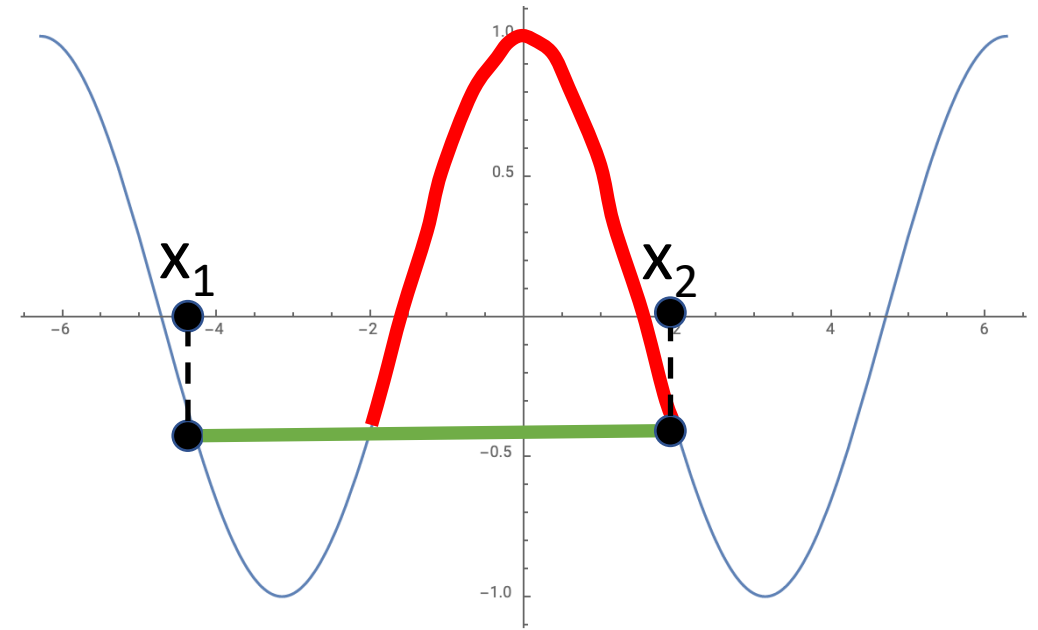


Out[14]=

[*]Many technical details inside!

# Convex Functions

A function $f : X \subseteq \mathbb{R}^N \to \mathbb{R}$ 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 inside!

Linear classifiers optimize a **convex function!**

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

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

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

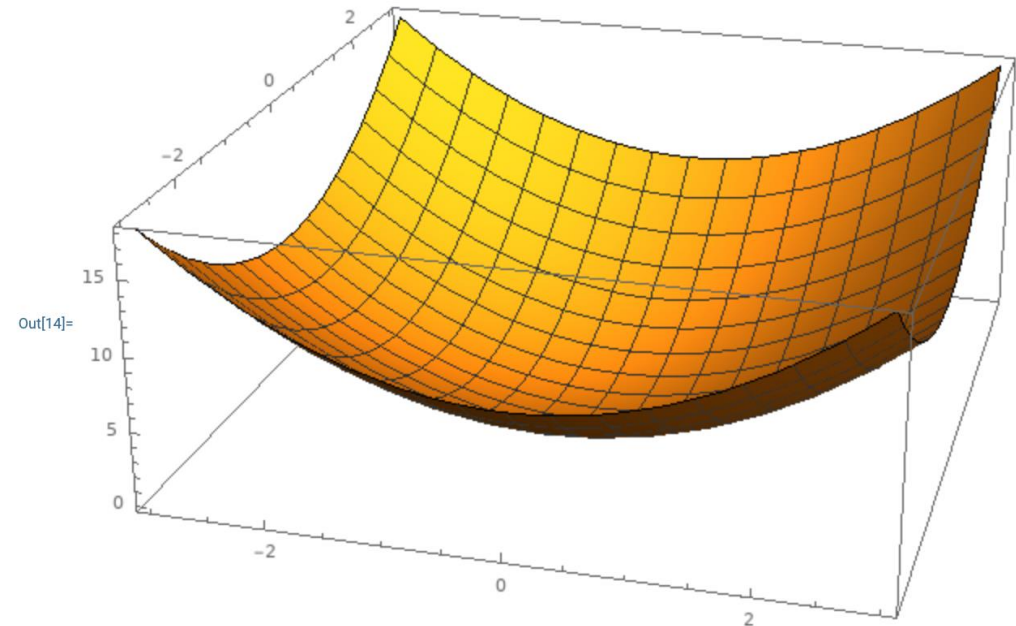$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + R(W)$$

R(W) = L2 or L1 regularization

# Convex Functions

A function $f : X \subseteq \mathbb{R}^N \to \mathbb{R}$ 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 inside!

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 $f : X \subseteq \mathbb{R}^N \to \mathbb{R}$ is **convex** if for all $x_1, x_2 \in X, t \in [0, 1]$,
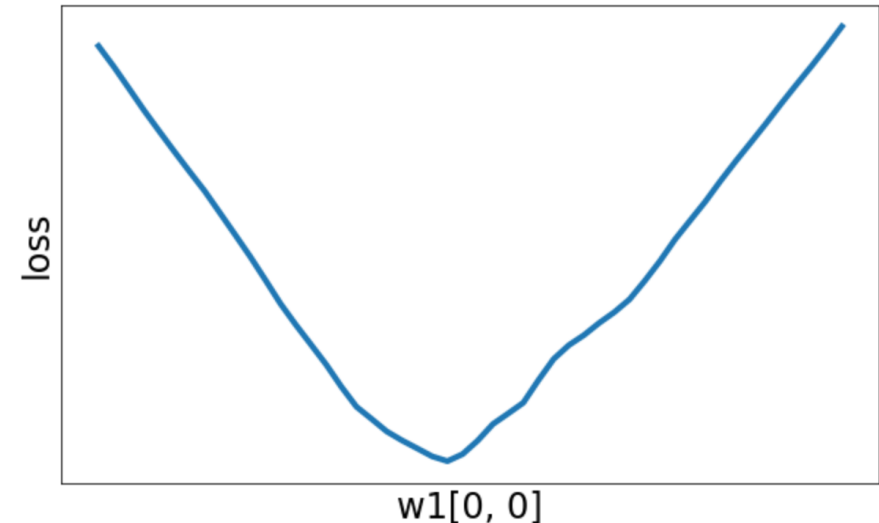
$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2)$$

**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 inside!

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 $f : X \subseteq \mathbb{R}^N \to \mathbb{R}$ is **convex** if for all $x_1, x_2 \in X, t \in [0, 1]$,
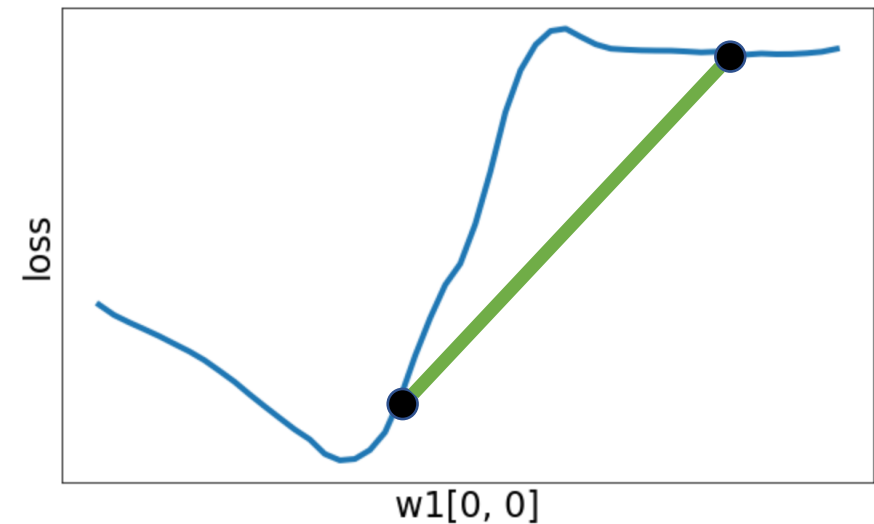
$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$$

**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 inside!

With local minima:



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
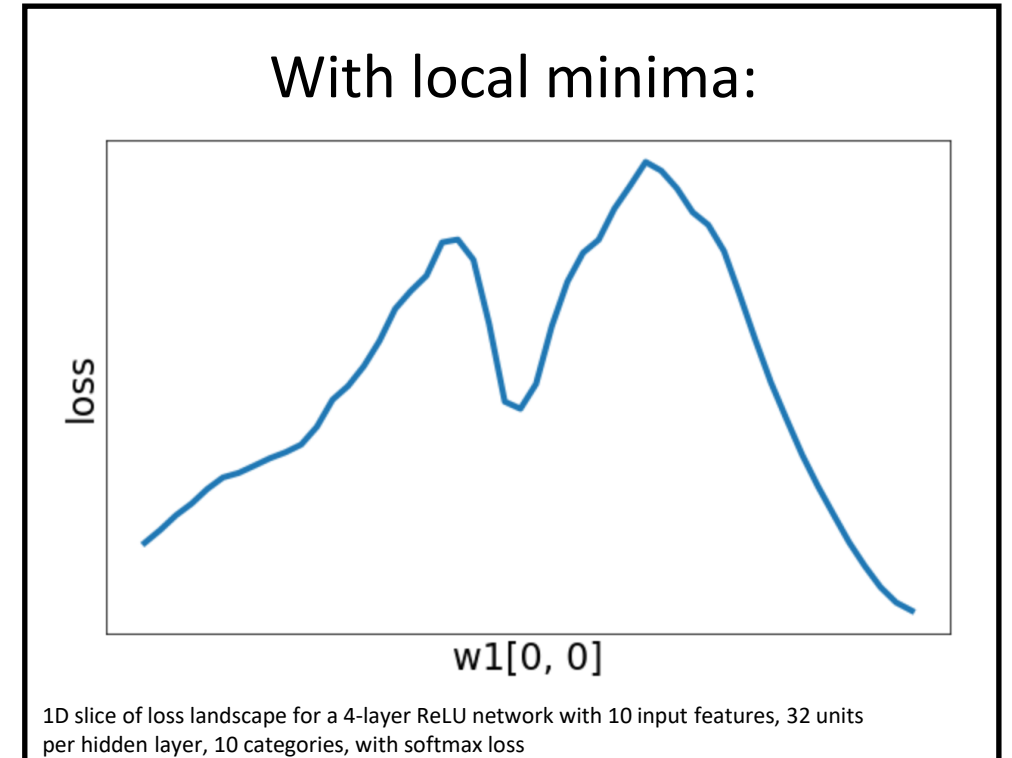
# Convex Functions

A function $f : X \subseteq \mathbb{R}^N \to \mathbb{R}$ 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 inside!

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
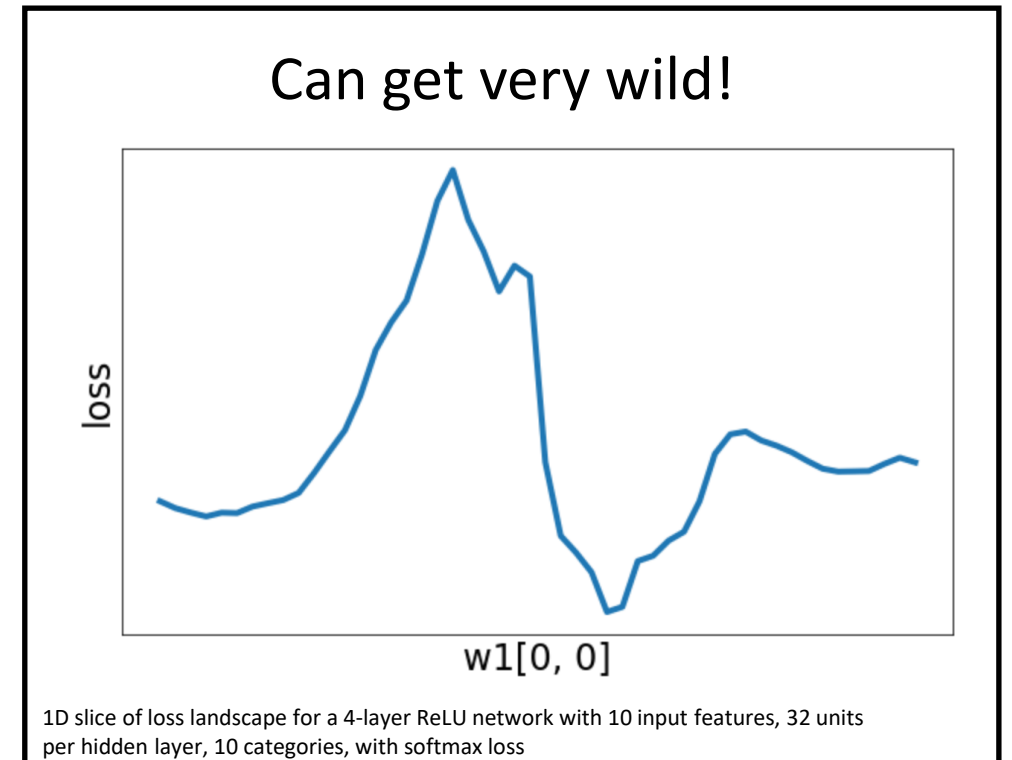
# Convex Functions

A function $f : X \subseteq \mathbb{R}^N \to \mathbb{R}$ 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 inside!

Most neural networks need **nonconvex optimization**
- Few or no guarantees about convergence
- Empirically it seems to work anyway
- Active area of research

# Convexity

- Most linear classifiers optimize a convex function
  - Linear layer

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

  - Cross-entropy loss

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

  - SVM

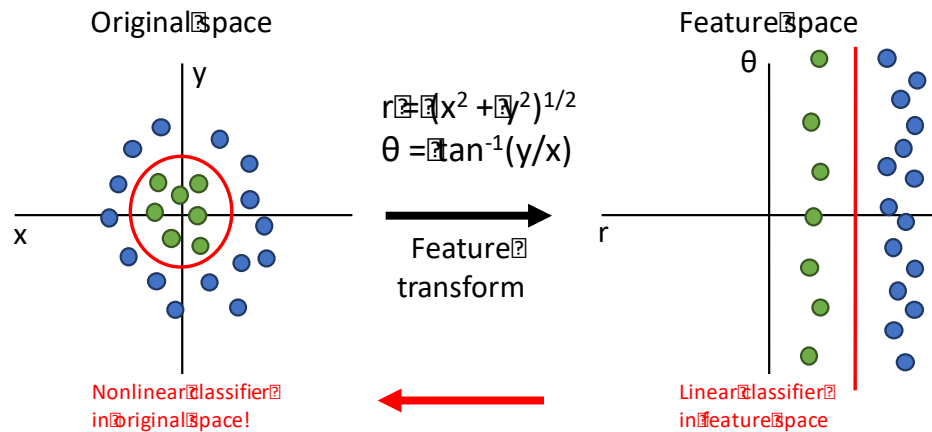  $$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

  - L1/L2 regularization

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

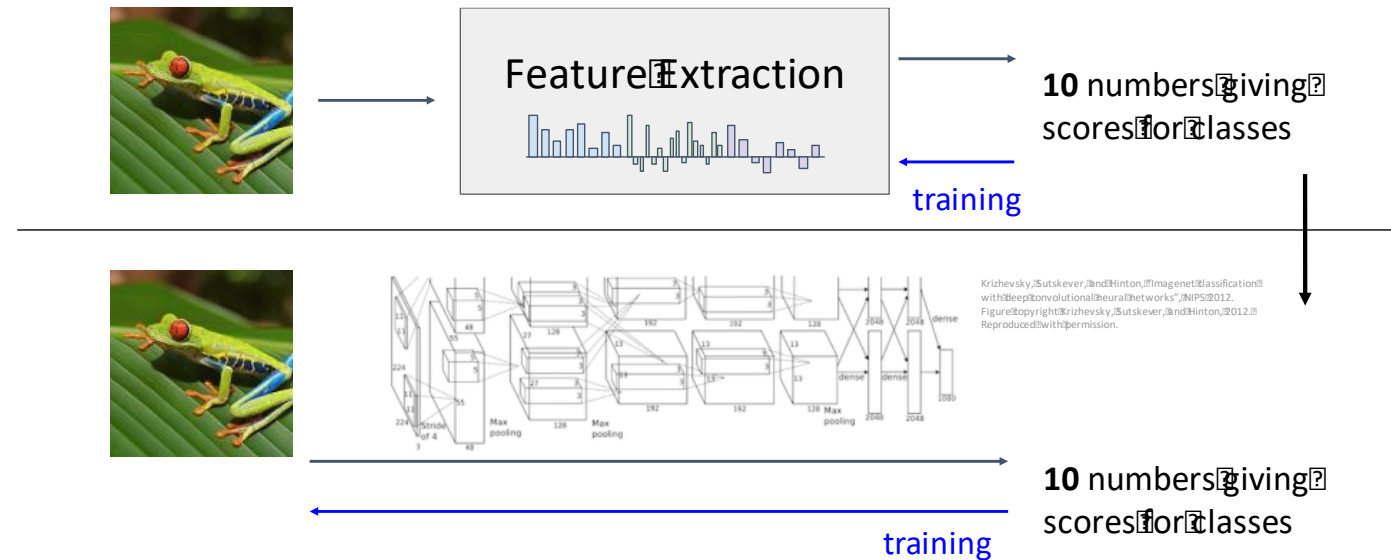- Most neural networks need non-convex optimization
  - Few or no guarantees about convergence (mostly falls in a local optimum)
  - Empirically it seems to work anyway
  - Active area of research

# Summary

Feature transform + Linear classifier
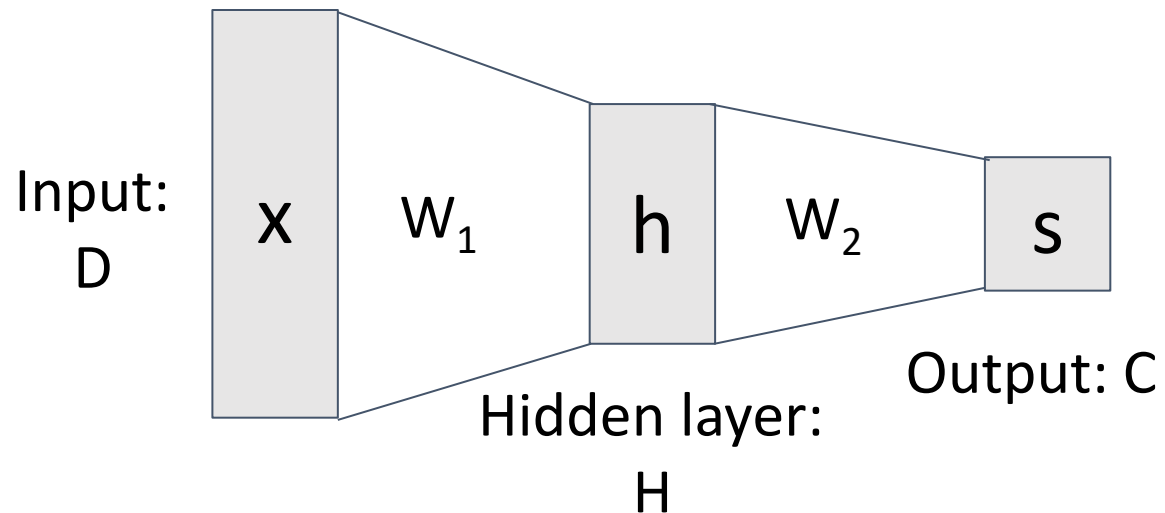allows nonlinear decision boundaries

Neural Networks as learnable feature transforms

**Original space**

y

r = (x² + y²)^{1/2}
θ = tan⁻¹(y/x)

Feature
transform

**Feature space**

θ

r

x

Nonlinear classifier
in original space!

Linear classifier
in feature space

Feature Extraction

**10** numbers giving
scores for classes

training

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

**10** numbers giving
scores for classes

training

# Summary

From linear classifiers to
fully-connected networks

$$s(x) = W_2 \, f(W_1 x + b_1) + b_2$$

Input:
D

x    $W_1$    h    $W_2$    s

Hidden layer:
H

Output: C

Linear classifier: One template per class



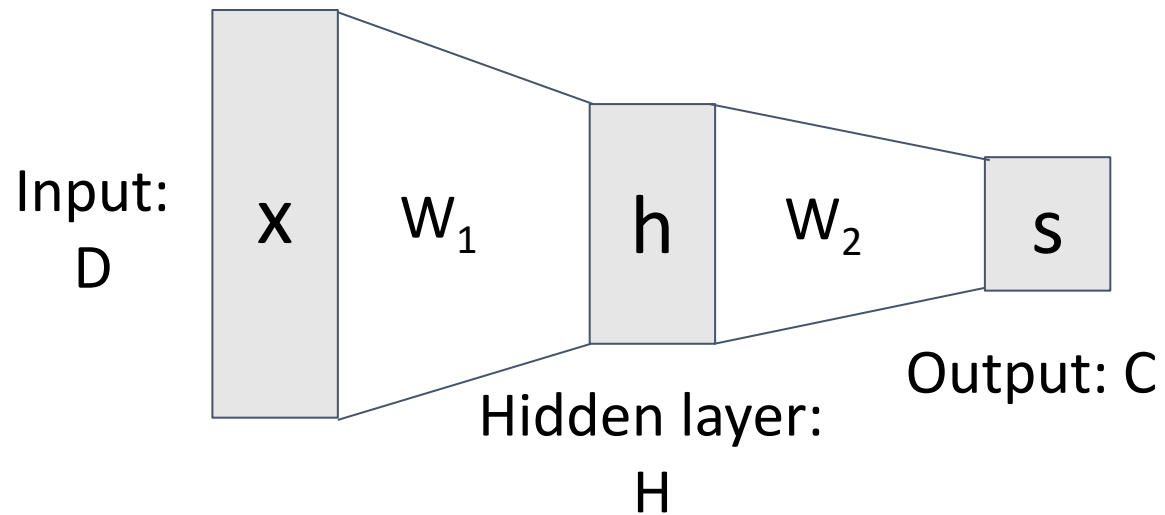Neural networks: Many reusable templates

# Summary

From linear classifiers to
fully-connected networks
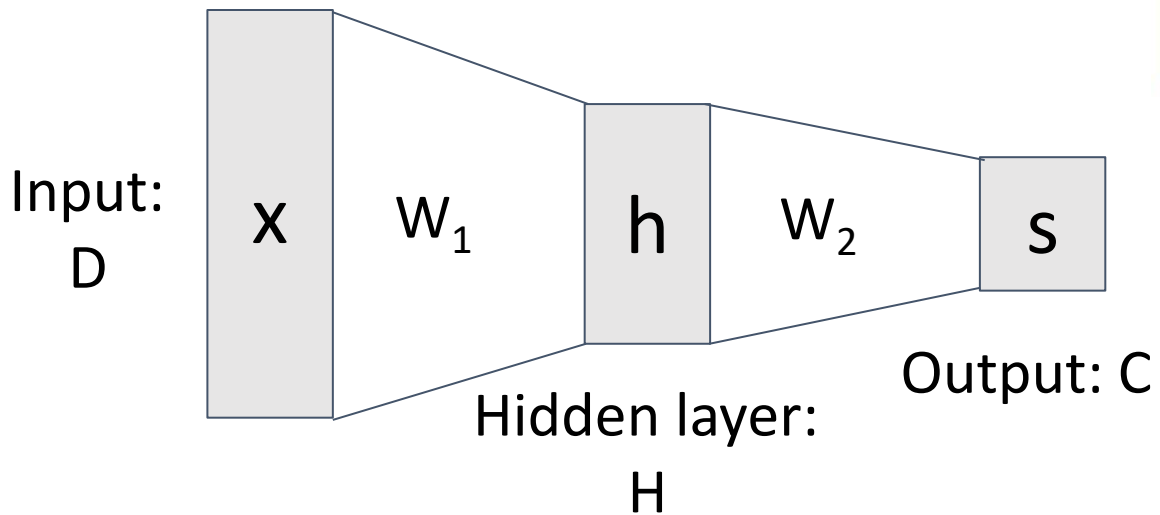
$$s(x) = W_2 \, f(W_1 x + b_1) + b_2$$

Input:
D

x    W₁    h    W₂    s

Hidden layer:
H

Output: C

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

# Summary

## Space Warping

## Universal Approximation

From linear classifiers to fully-connected networks

$$s(x) = W_2\, f(W_1 x + b_1) + b_2$$

Input: D

x    $W_1$    h    $W_2$    s

Hidden layer: H

Output: C

## Nonconvex

loss

w1[0, 0]

# Problem: How to compute gradients?

$$s = W_2\, f(W_1 x + b_1) + b_2$$      Nonlinear score function

$$L_i = \sum_{j \neq y_i} \max\left(0, s_j - s_{y_i} + 1\right)$$      Per-element data loss

$$R(W) = \sum_k W_k^2$$      L2 Regularization

$$L(W_1, W_2, b_1, b_2) = \frac{1}{N} \sum_{i=1}^{N} L_i + \lambda R(W_1) + \lambda R(W_2)$$      Total loss

If we can compute $\frac{\partial L}{\partial W_1}, \frac{\partial L}{\partial W_2}, \frac{\partial L}{\partial b_1}, \frac{\partial L}{\partial b_2}$ then we can optimize with SGD

# Next: Backpropagation