

# CSCI 497P/597P: Computer Vision

## (Stochastic) Gradient Descent Neural networks



# Readings

with a great deal more detail...

- <http://cs231n.github.io/optimization-1/>
- <http://cs231n.github.io/optimization-2/>
- <http://cs231n.github.io/neural-networks-1/>
- <http://cs231n.github.io/neural-networks-2/>
- <http://cs231n.github.io/neural-networks-3/>

# Announcements

- Last project - P4 (AlexNet)
  - Out tomorrow 5/27
  - Due a week from tomorrow 6/3
- HW3 out, due Monday
  - your lowest homework score is dropped

# Goals

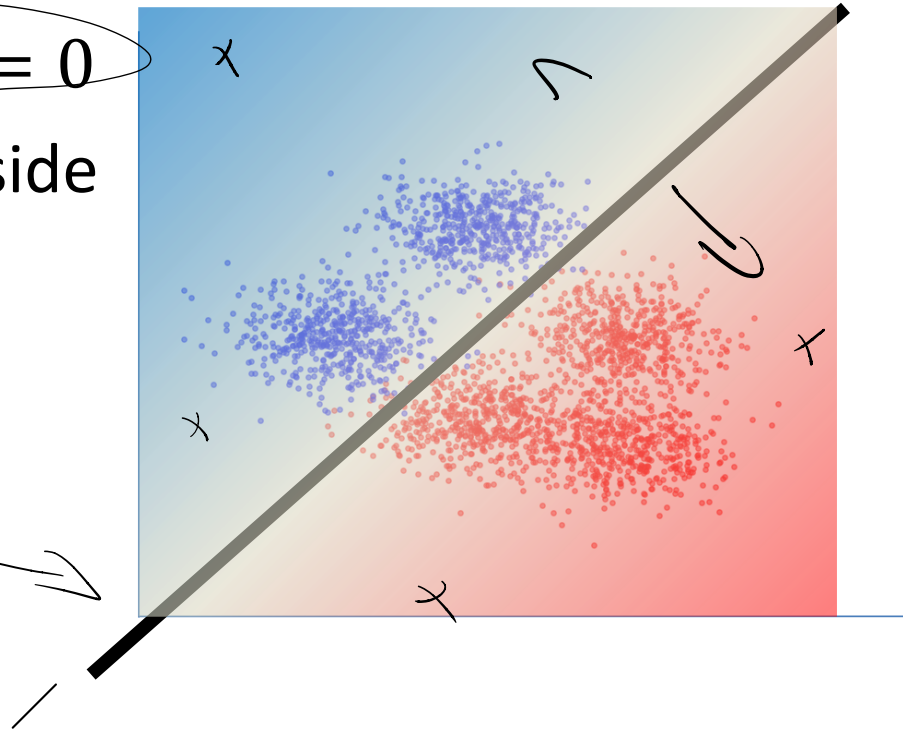
- Understand how to train a classifier by minimizing a loss function using gradient descent.
- Understand the intuition behind using Stochastic (Minibatch) Gradient Descent.
- Understand neural networks as a stack of linear classifiers with nonlinearities (activation functions) in between.
  - Understand why we need activation functions.
  - Understand the vanishing gradients problem.

# Linear classifiers

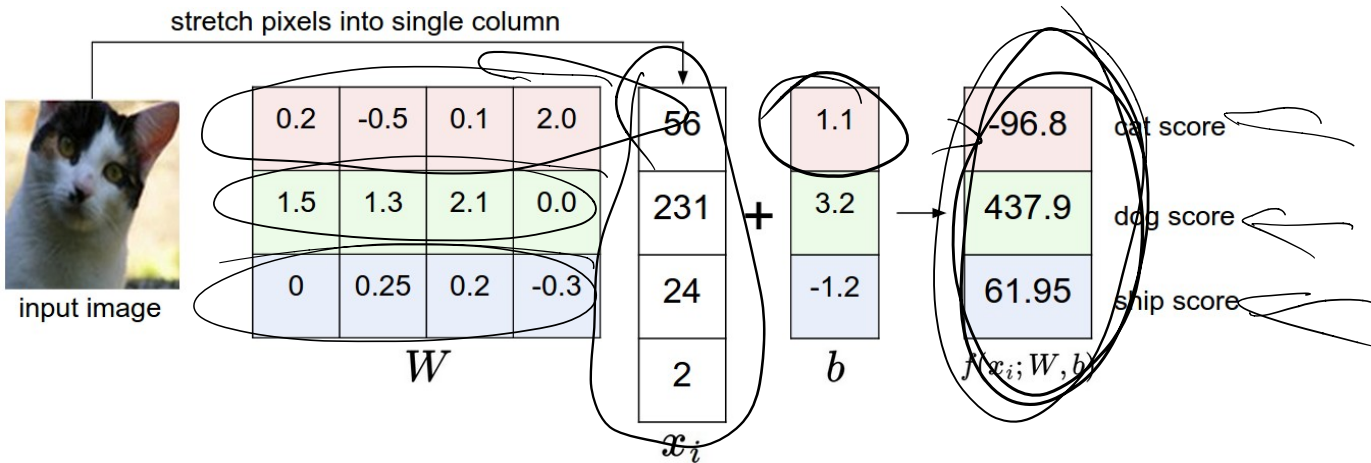
$1 \times d$   
 $d \times 1$

- Equation:  $w^T x + b = 0$
- Points on the same side are the same class

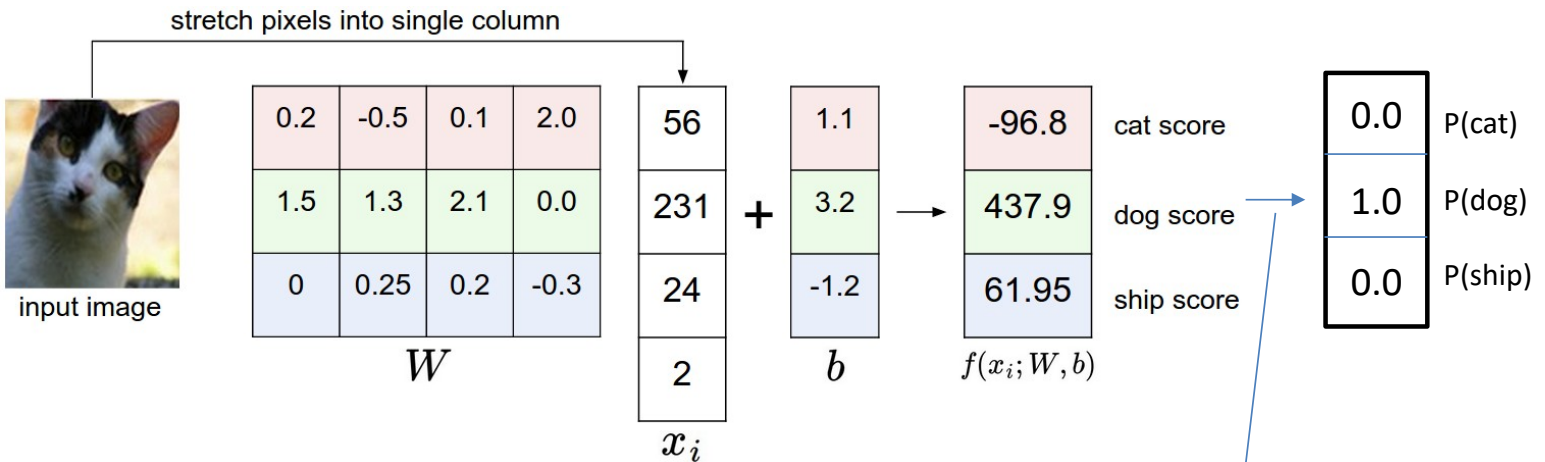
$$ax + by + c = 0$$



# Multiclass Linear Classifiers: Stack multiple $w^T$ into a matrix.



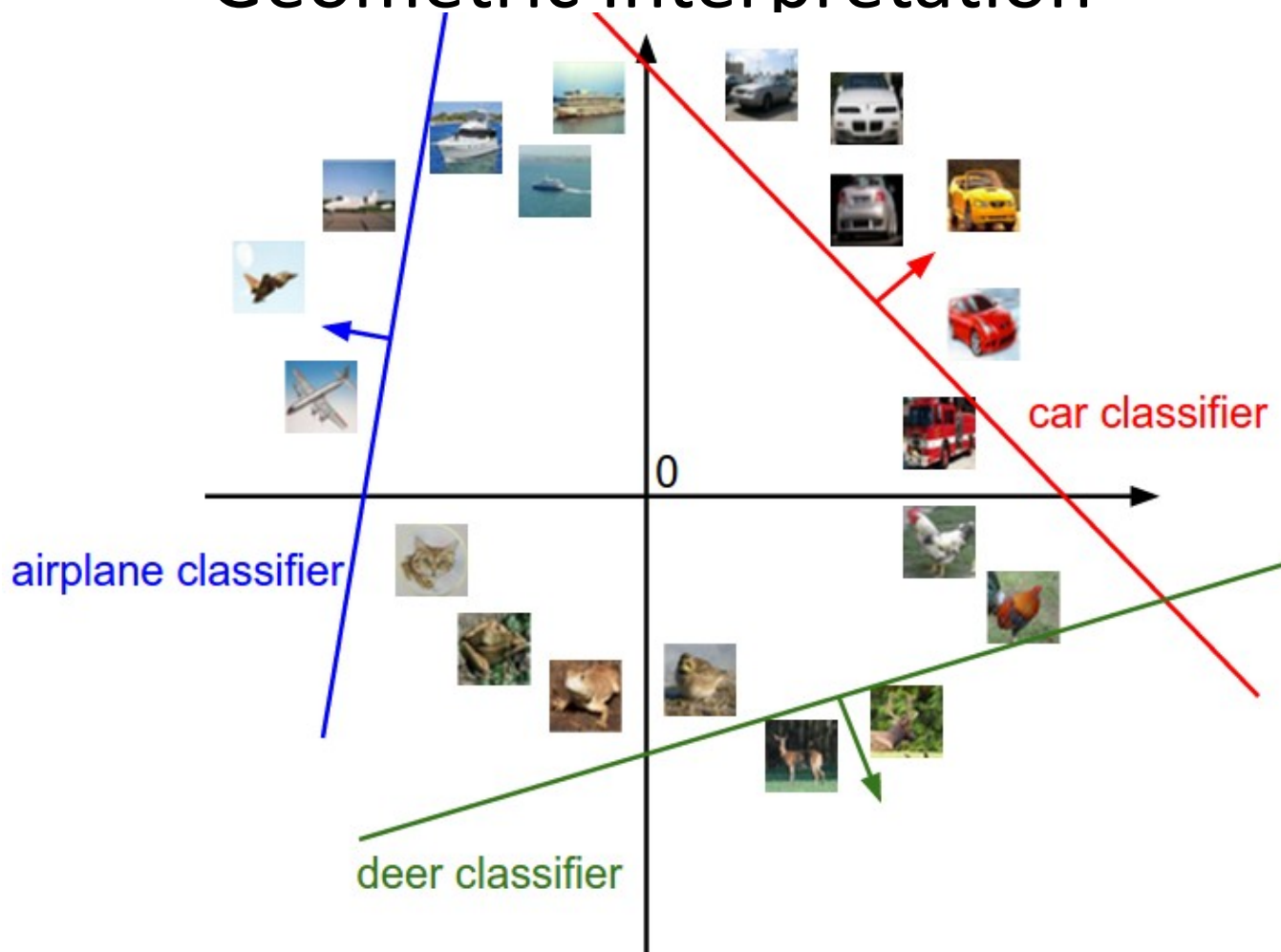
# Multiclass Linear Classifiers: Stack multiple $w^T$ into a matrix.



Softmax normalization

$$\frac{e^{f_k}}{\sum_j e^{f_j}}$$

# Multiclass Linear Classifier: Geometric Interpretation





# Taking stock

- We have:

→  $\phi = \text{unravel}(\text{rgb2gray}(\text{img}))$ , a feature extractor

↓  
–  $h(x) = \underline{W}^T \underline{x}$ , a multiclass linear classifier

–  $L = \sum_{i=0}^N L_i$ , a loss function

*for each  
data example*

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

- We don't have:

– a way to find a  $W$  that results in a small  $L$ .

# Minimizing the Loss

- Use **optimization** to find the  $W$  that *minimizes* the loss function.
  - Linear regression: solvable in closed form
  - Most of the time: no closed form.

# Optimization



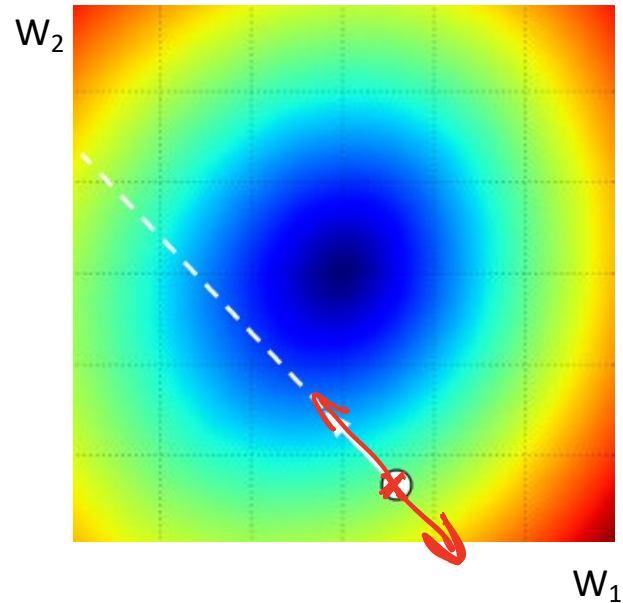
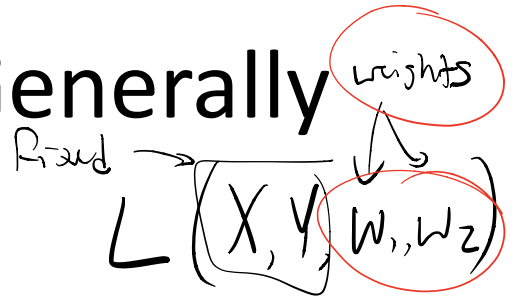
# Finding a $W$ that minimizes $L$

- Simple idea: walk downhill.



# Gradient Descent: Generally

- Gradient of the loss function with respect to the *weights* tells us how to change the weights to improve the loss.

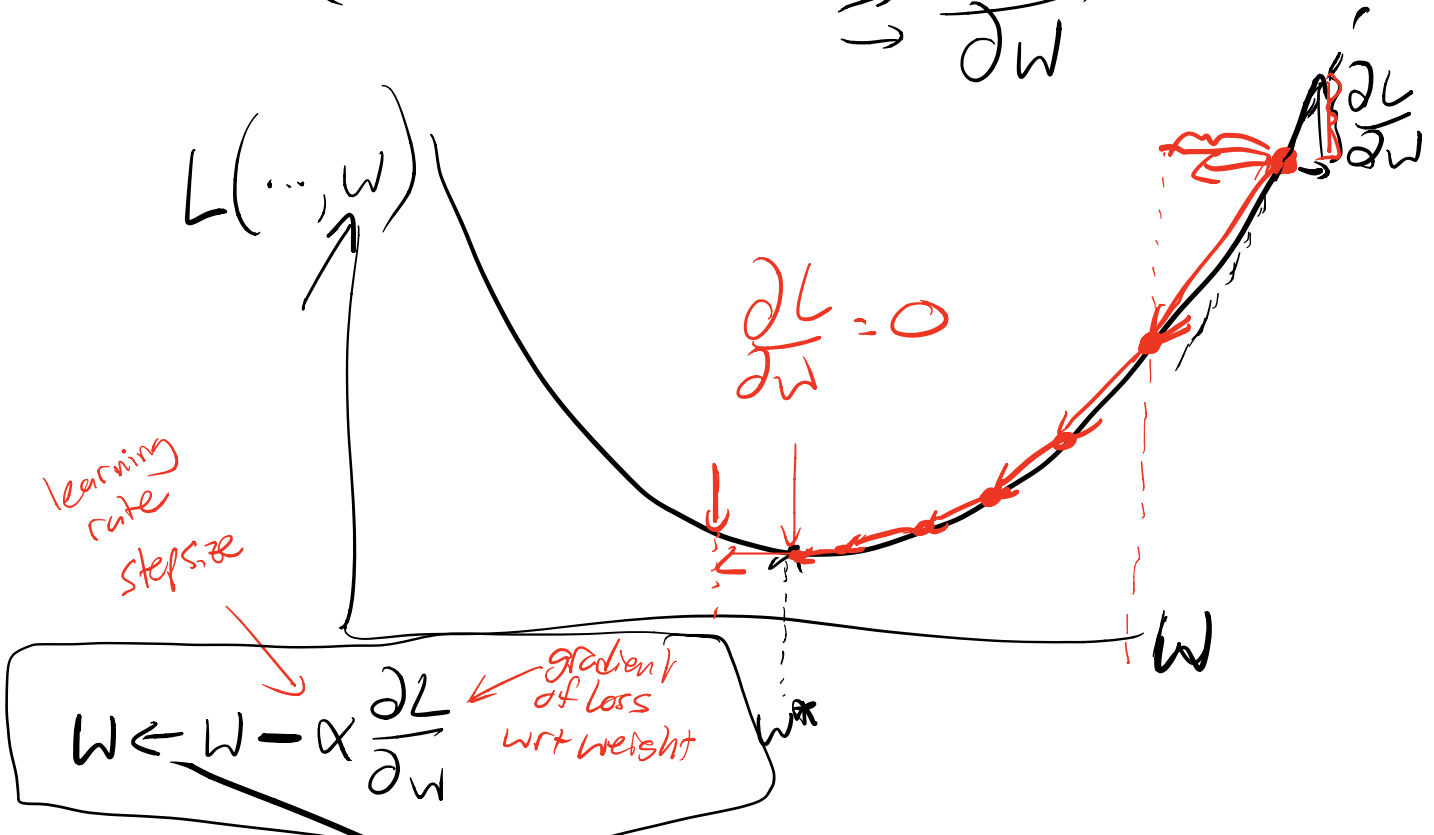


# Gradient Descent: Intuition

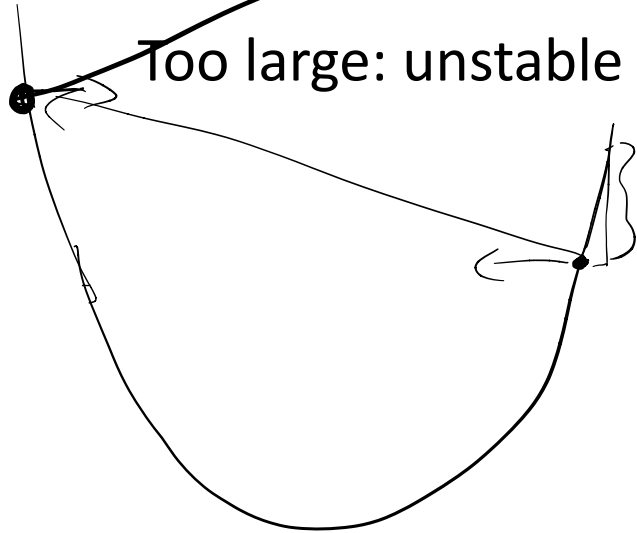
inputs      labels      parameter of classifier

$$L(X, Y, w)$$

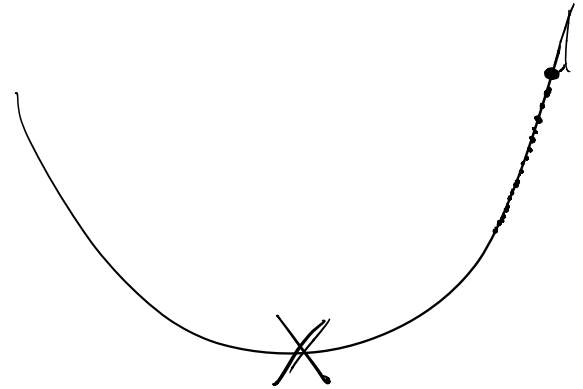
$$\Rightarrow \frac{\partial L}{\partial w}$$



# The effect of Step Size

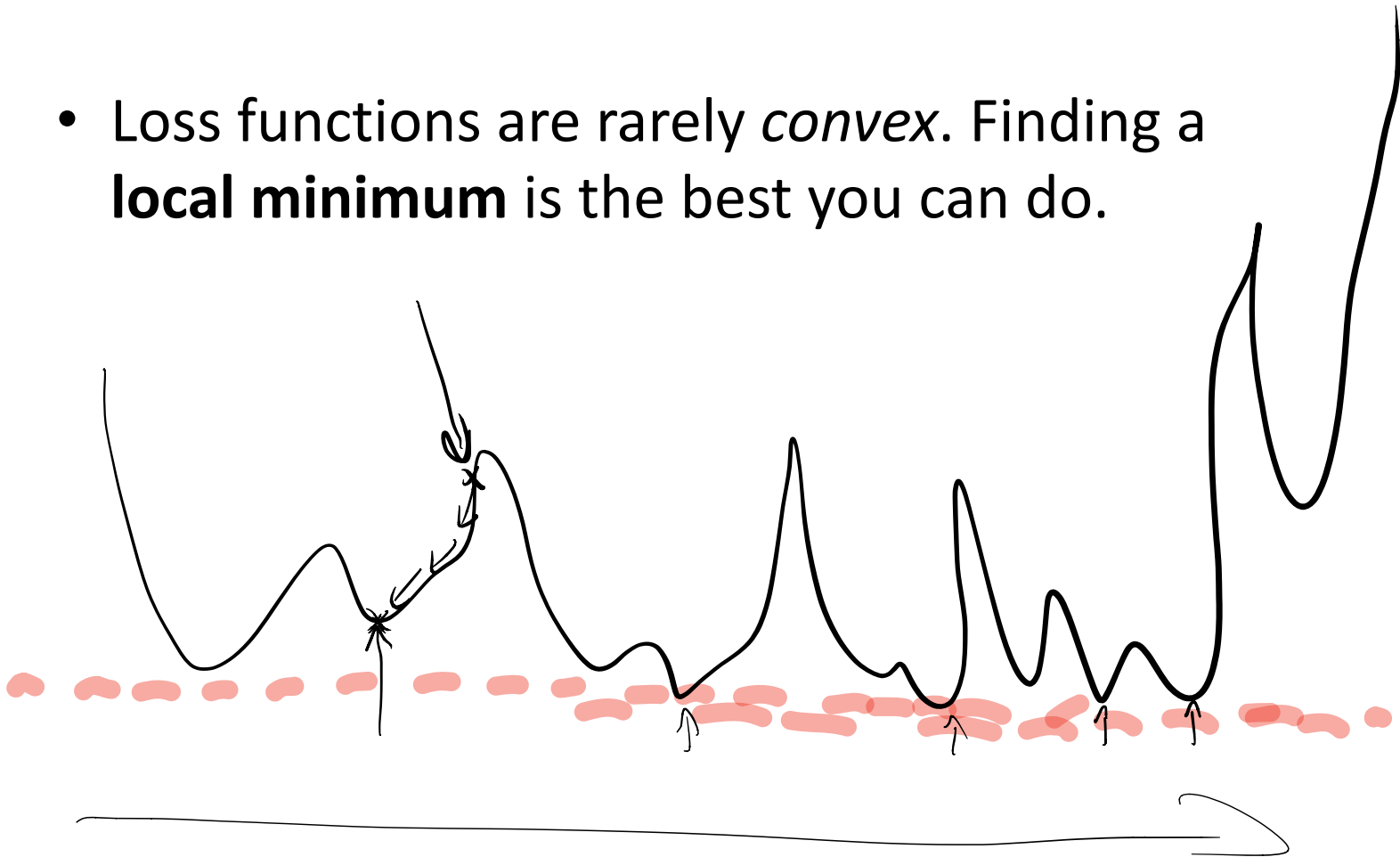


Too small: slow convergence



# Reality isn't quite so pretty

- Loss functions are rarely *convex*. Finding a **local minimum** is the best you can do.





# Gradient Descent

```
# Vanilla Gradient Descent
```

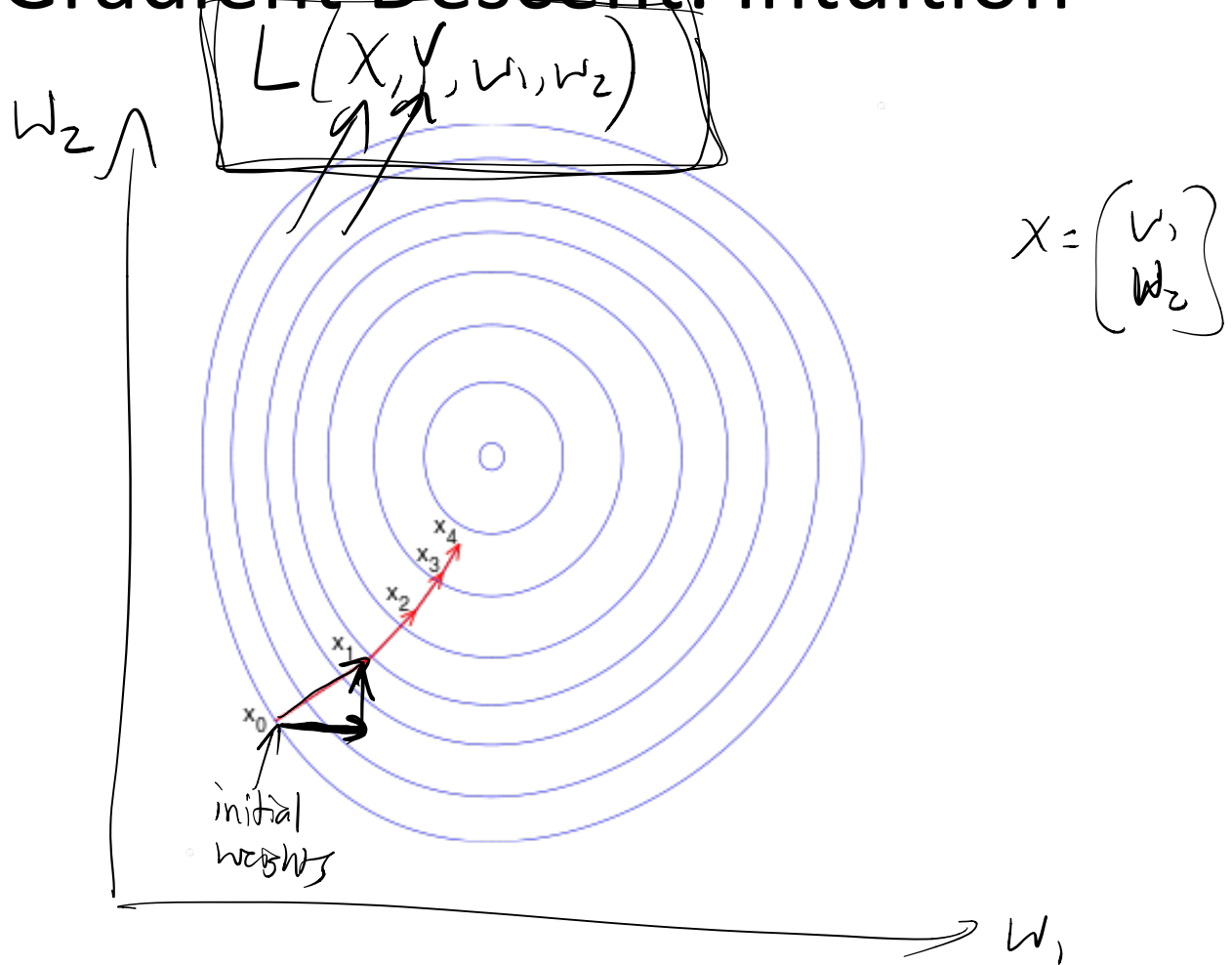
```
while True:
```

```
    weights_grad = evaluate_gradient(loss_fun, data, weights)
```

```
    weights += - step_size * weights_grad # perform parameter update
```



# Gradient Descent: Intuition



# Gradient Descent: Demo

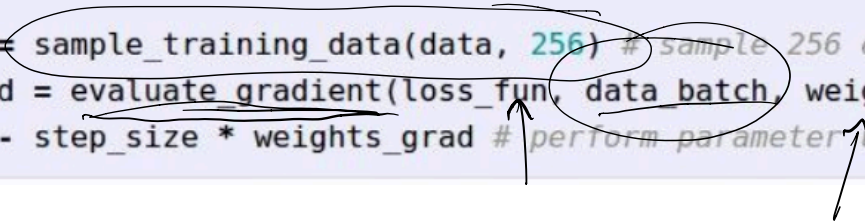
- <http://vision.stanford.edu/teaching/cs231n-demos/linear-classify/>
  - select “Softmax” radio button at the bottom

# Stochastic Gradient Descent

```
# Vanilla Minibatch Gradient Descent
```

```
while True:
```

```
    data_batch = sample_training_data(data, 256) # sample 256 examples  
    weights_grad = evaluate_gradient(loss_fun, data_batch, weights)  
    weights += - step_size * weights_grad # perform parameter update
```



- $L(X, Y; W)$  depends on
  - All data points  $x_1..x_n$
  - Ground truth labels  $y_1..y_n$
  - Weights  $W$
- Very expensive to evaluate if you have a lot of data.

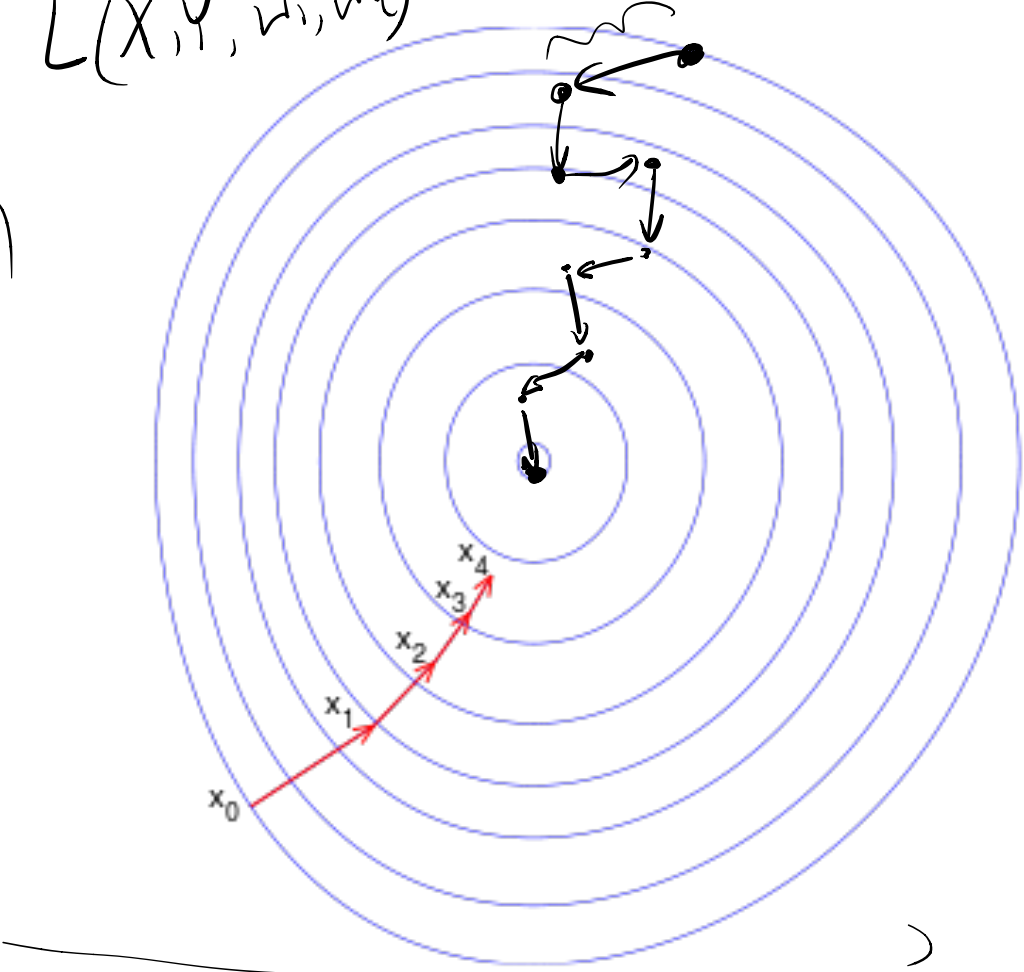
# Stochastic Gradient Descent

- Idea: consider only a few data points at a time.
- Loss is now computed using only a small batch (minibatch) of data points.
- Update weights the same way using the gradient of  $L$  wrt the weights.

# Stochastic Gradient Descent: Intuition

$$L(X, Y, w_1, w_2)$$

$w_2$



$w_1$

# Taking stock

- We have:

–  $\phi = \text{unravel}(\text{rgb2gray}(\text{img}))$ , a feature extractor

–  $h(x) = W^T x$ , a multiclass linear classifier

→

$$L = \sum_{i=0}^N L_i$$

↑

, a loss function

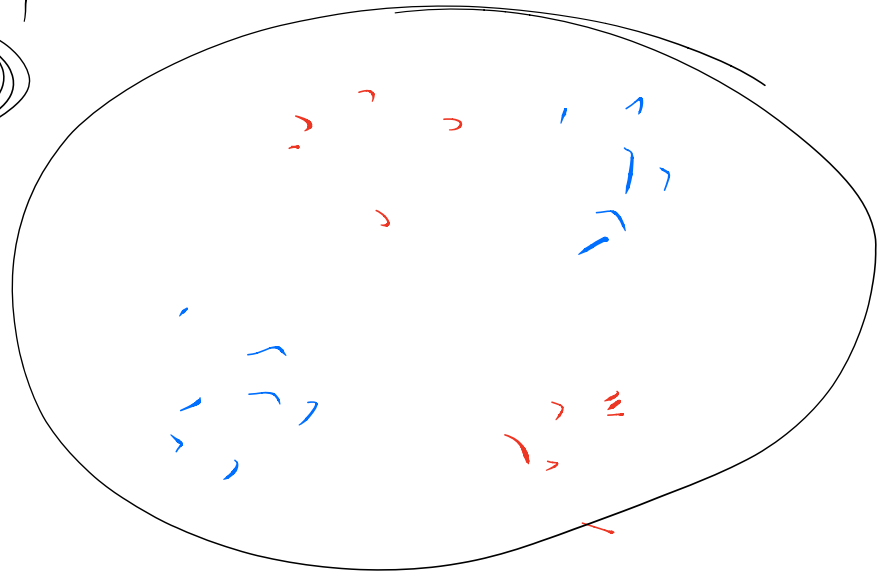
$$L_i = -\log \left( \frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right)$$

- 
- A way too adjust  $W$  until we can't make  $L$  any smaller.

# So about that linearly separable assumption...

- Ideas:
  - $\phi = \text{unravel}(\text{rgb2gray}(\text{img}))$ , a feature extractor
    - use a fancier  $\phi$ ?

– Learn  $\phi$  too.





# Neural Networks

## Neural Network

Linear  
classifiers



---

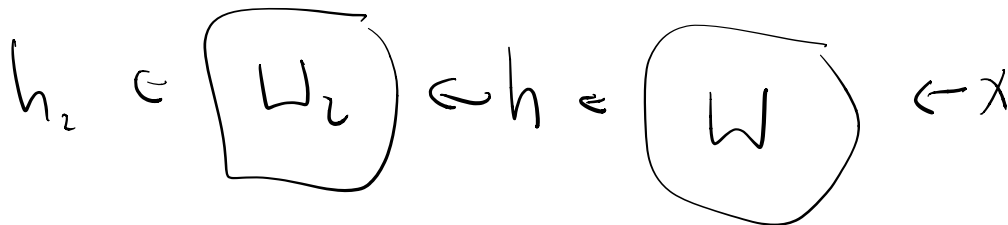

## Neural networks: without the brain stuff

(**Before**) Linear score function:  $f = Wx$

## Neural networks: without the brain stuff

(**Before**) Linear score function:  $f = Wx$

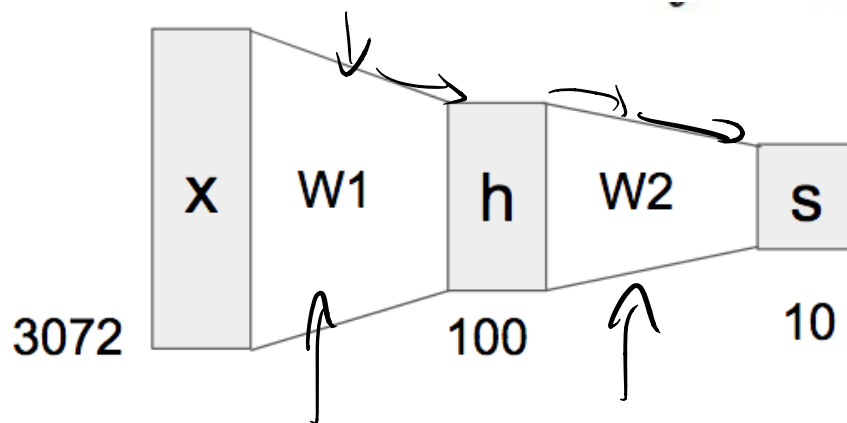
(**Now**) 2-layer Neural Network  $f = W_2 \max(0, W_1 x)$



## Neural networks: without the brain stuff

(**Before**) Linear score function:  $f = Wx$


(**Now**) 2-layer Neural Network  $f = W_2 \max(0, W_1 x)$



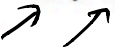
## Neural networks: without the brain stuff

(**Before**) Linear score function:  $f = Wx$

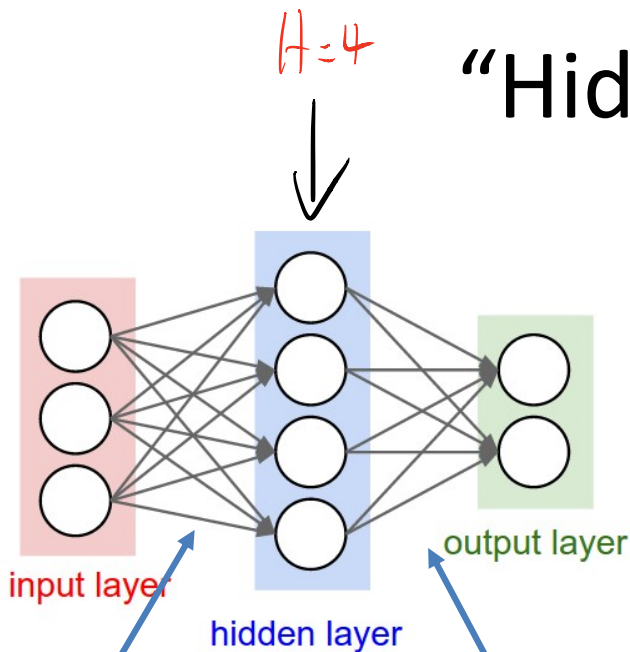
(**Now**) 2-layer Neural Network  
or 3-layer Neural Network  $f = W_2 \max(0, W_1 x)$

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


# Training a 2 layer neural network in 20 lines of python

```
1 import numpy as np
2 from numpy.random import randn
3
4 N, D_in, H, D_out = 64, 1000, 100, 10
5 x, y = randn(N, D_in), randn(N, D_out)
6 w1, w2 = randn(D_in, H), randn(H, D_out)
7 
8 for t in range(2000):
9     h = 1 / (1 + np.exp(-x.dot(w1)))
10    y_pred = h.dot(w2)
11    loss = np.square(y_pred - y).sum()
12    print(t, loss)
13
14    grad_y_pred = 2.0 * (y_pred - y)
15    grad_w2 = h.T.dot(grad_y_pred)
16    grad_h = grad_y_pred.dot(w2.T)
17    grad_w1 = x.T.dot(grad_h * h * (1 - h))
18
19    w1 -= 1e-4 * grad_w1
20    w2 -= 1e-4 * grad_w2
```

# “Hidden Layers”



$W_1$ , a  $3 \times 4$  matrix converts input into hidden layer activations

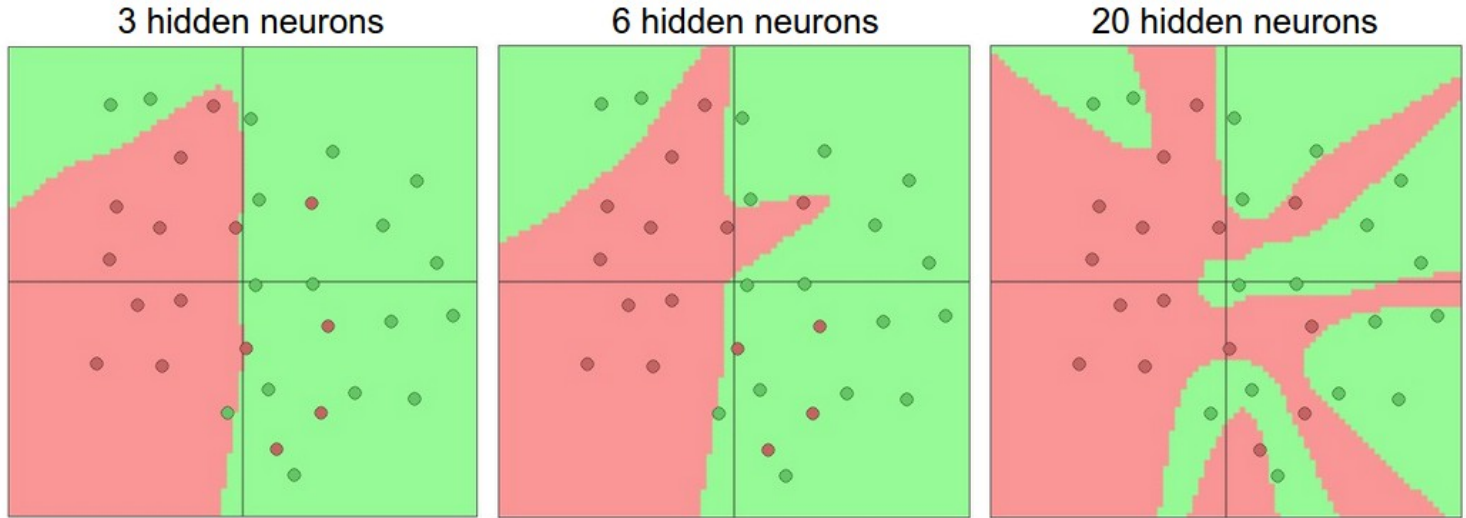
$W_2$ , a  $4 \times 2$  matrix transforms hidden layer activations to output scores

$$c = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} = x$$

$c \times 1$        $c \times H$        $H \times d$        $d \times 1$   
 $h \in W_2 \quad W_1 \quad x$

*pick this!*

# Neural Networks: Nonlinear Classifiers built from Linear Classifiers





## Neural networks: without the brain stuff

(**Before**) Linear score function:  $f = Wx$

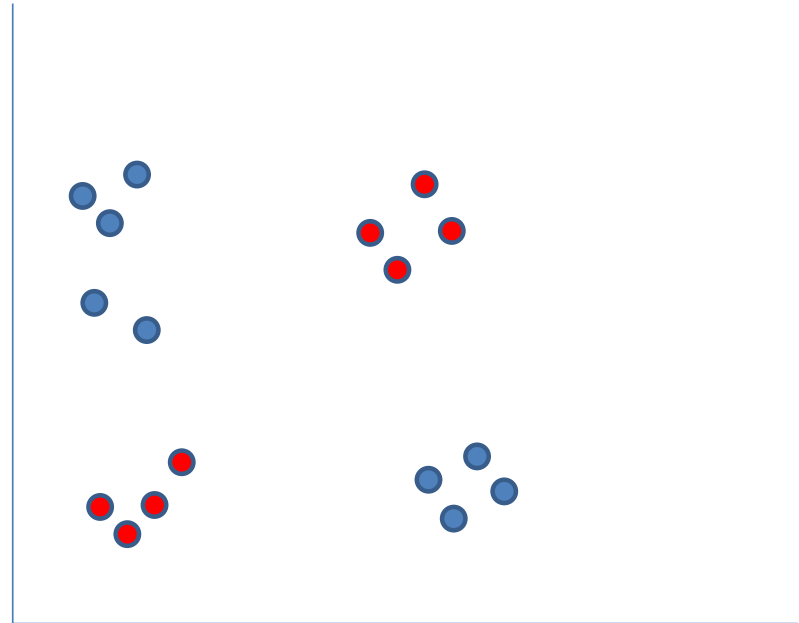
(**Now**) 2-layer Neural Network  
or 3-layer Neural Network  $f = W_2 \max(0, W_1 x)$

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

???

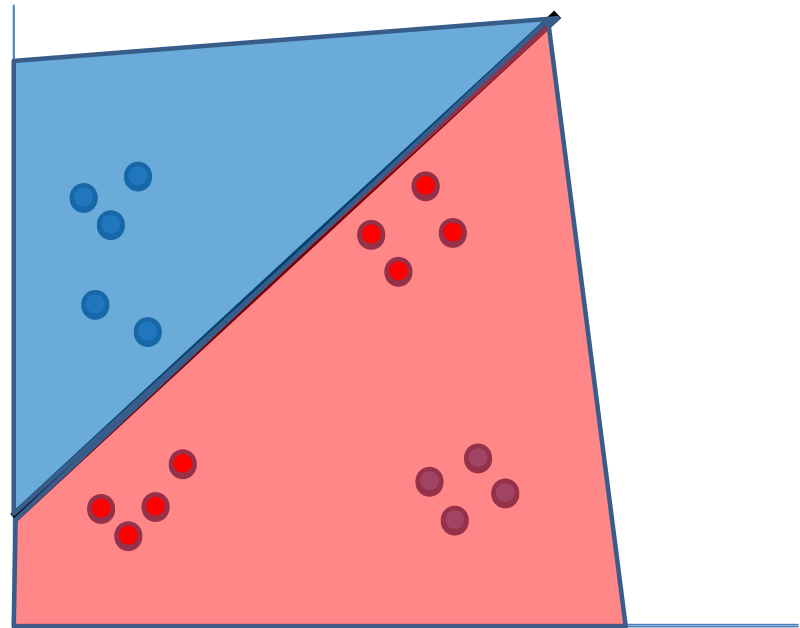
# Activation Functions

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



# Activation Functions

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

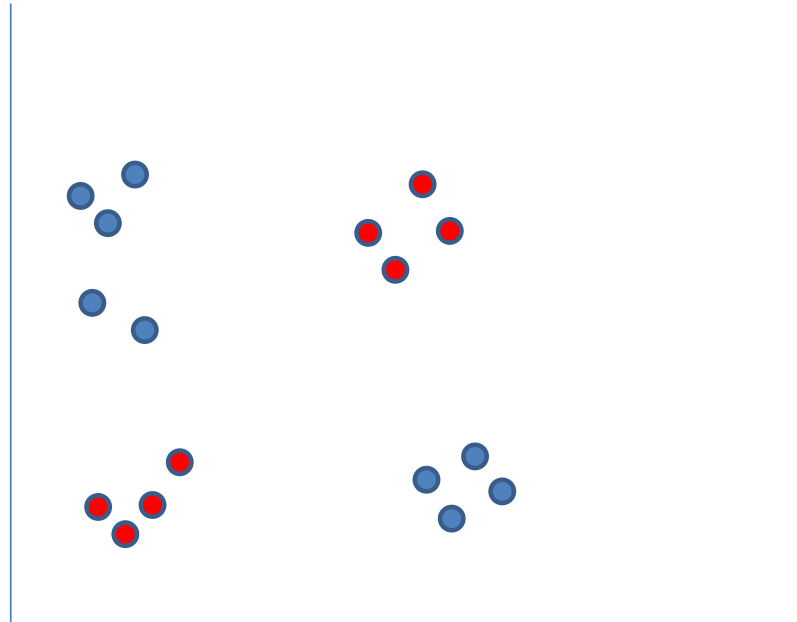


A linear classifier can only do so well...

# Activation Functions

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

$$f(x, W_1, W_2) = W_1(W_2x)$$



Let's try stacking two linear classifiers together

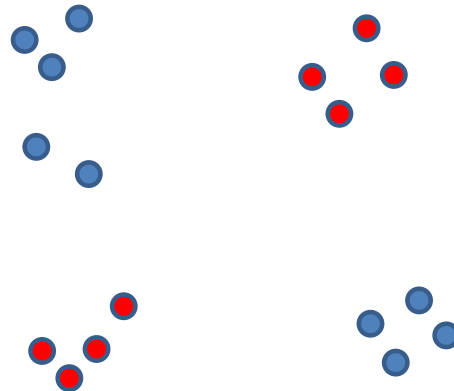
# Activation Functions

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

$$f(x, W_1, W_2) = W_1(W_2x)$$

$$W \leftarrow W_1 W_2$$

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



Uh oh – linear functions compose to linear functions.

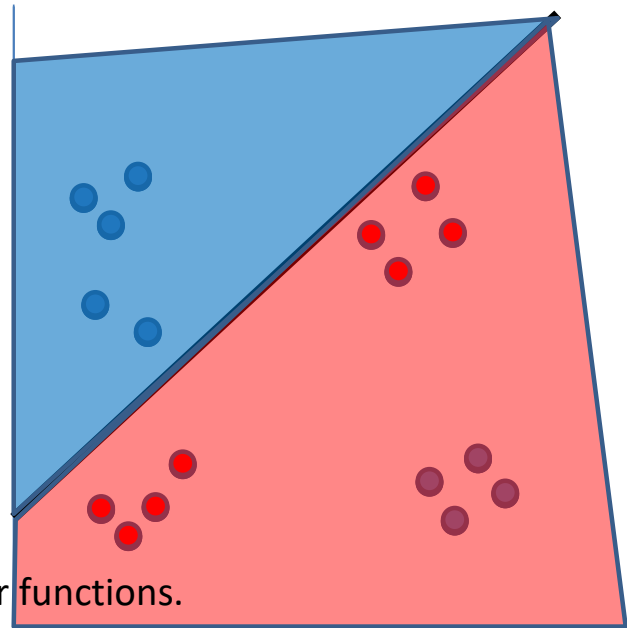
# Activation Functions

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

$$f(x, W_1, W_2) = W_1(W_2x)$$

$$W \leftarrow W_1 W_2$$

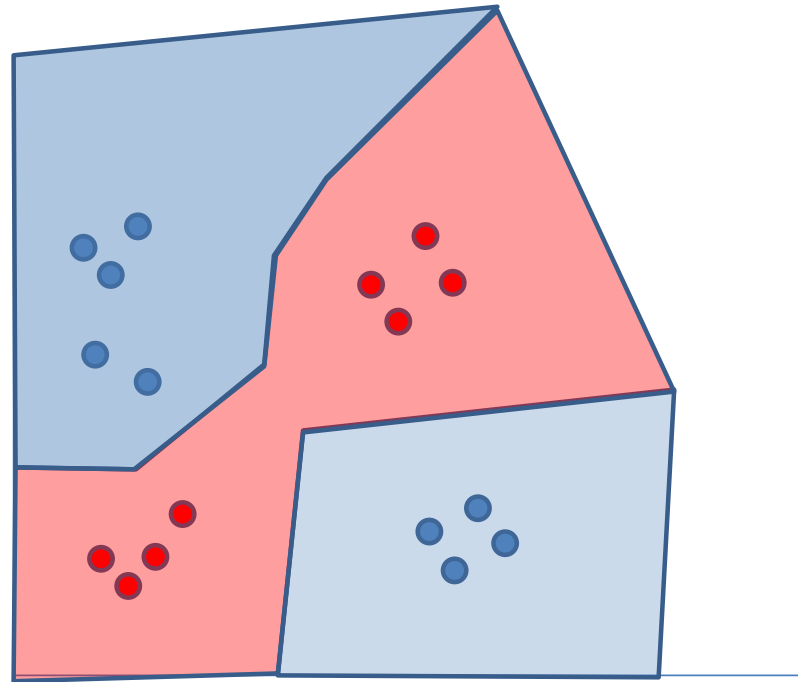
$$f(x, W) = Wx$$



Uh oh – linear functions compose to linear functions.

# Activation Functions

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



Nonlinearities prevent the composed linear functions from collapsing into a single one.

This is a key property of universal approximation.

# Neural Networks

## Neural Network

Linear  
classifiers

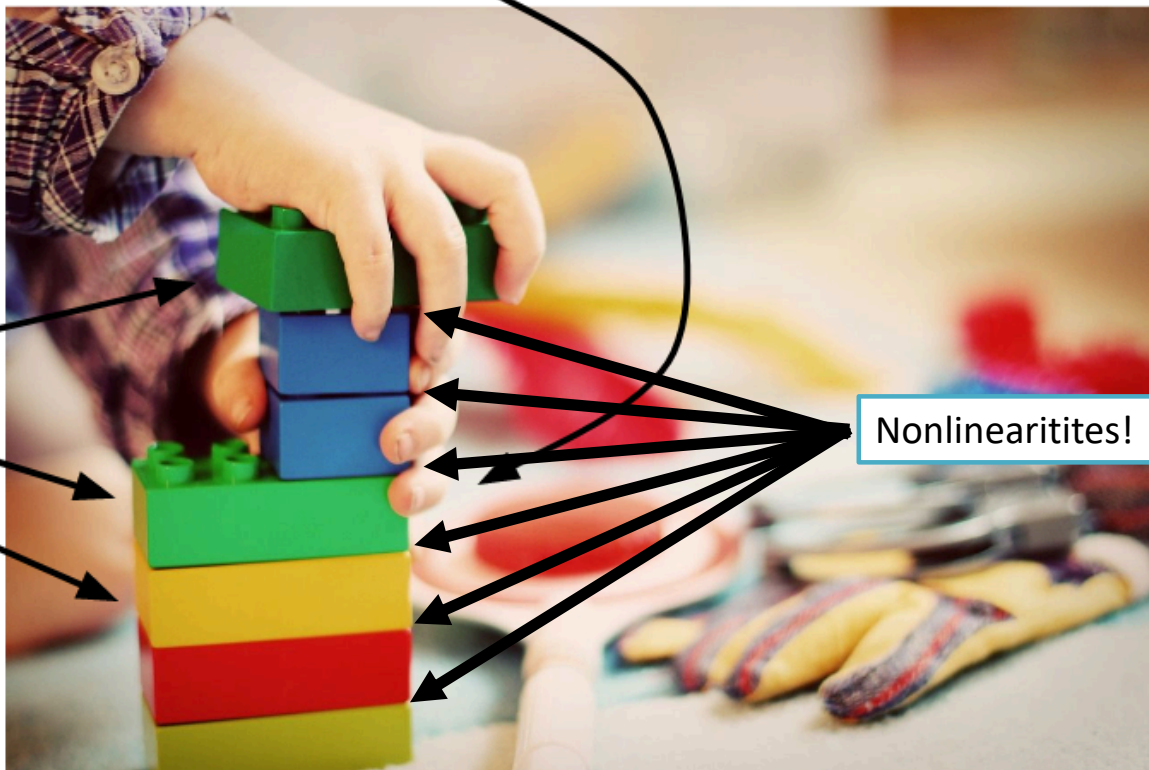




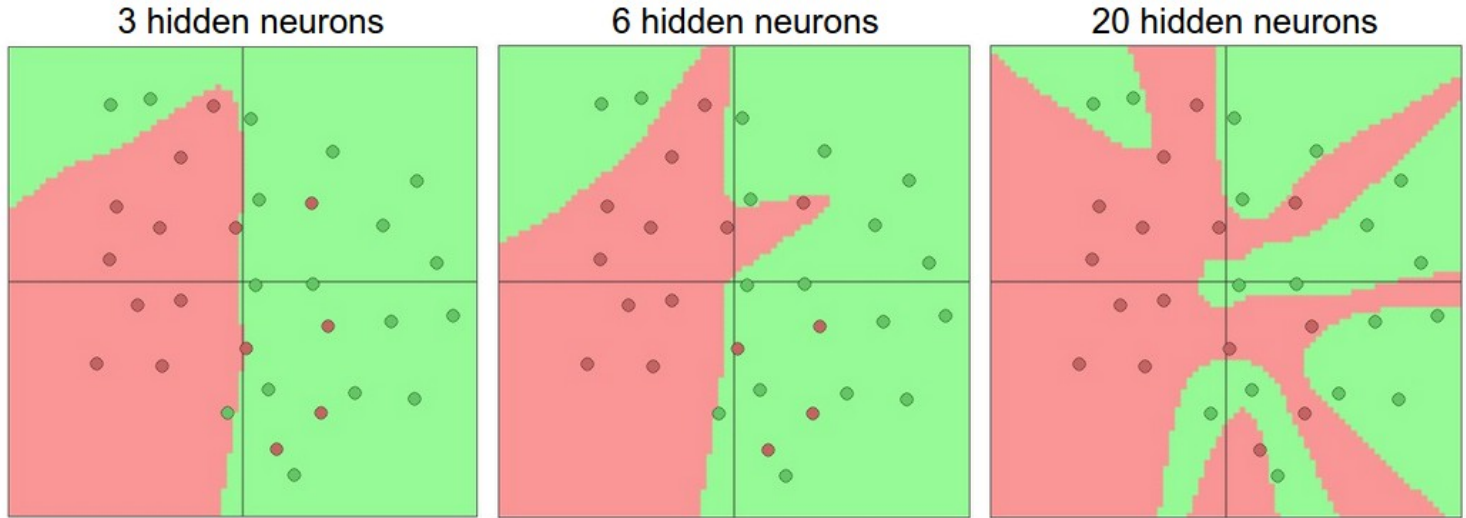
# Neural Networks

## Neural Network

Linear  
classifiers



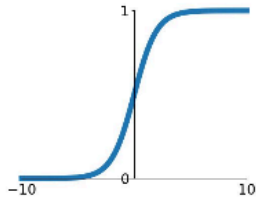
# Neural Networks: Nonlinear Classifiers built from Linear Classifiers



# Activation Functions

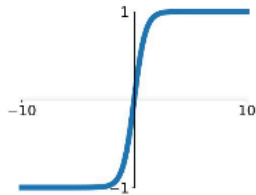
## Sigmoid

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



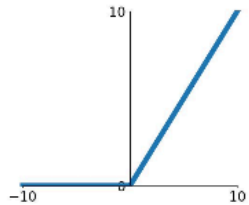
## tanh

$$\tanh(x)$$



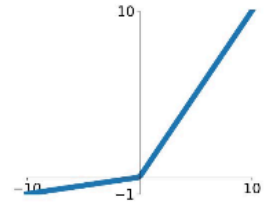
## ReLU

$$\max(0, x)$$



## Leaky ReLU

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



## Maxout

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

## ELU

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

