

# CS 4803 / 7643: Deep Learning

Topics:

- Optimization
- Computing Gradients

Dhruv Batra  
Georgia Tech

# Administrativa

- HW1 Reminder
  - Due: 09/09, 11:59pm

# Recap from last time

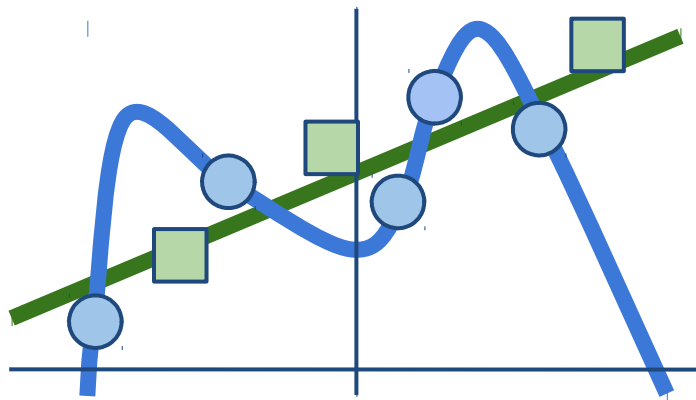
# Regularization

$\lambda$  = regularization strength  
(hyperparameter)

$$\underline{L(W)} = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss}} + \underbrace{\lambda R(W)}_{\text{Regularization}}$$

**Data loss:** Model predictions should match training data

**Regularization:** Prevent the model from doing *too* well on training data



**Occam's Razor:**

*"Among competing hypotheses,  
the simplest is the best"*

William of Ockham, 1285 - 1347

# Regularization

$\lambda$  = regularization strength  
(hyperparameter)

$$L(W) = \underbrace{\frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i)}_{\text{Data loss}} + \underbrace{\lambda R(W)}_{\text{Regularization}}$$

**Data loss:** Model predictions should match training data

**Regularization:** Prevent the model from doing *too well* on training data

## Simple examples

L2 regularization:  $R(W) = \sum_k \sum_l W_{k,l}^2$

L1 regularization:  $R(W) = \sum_k \sum_l |W_{k,l}|$

Elastic net (L1 + L2):  $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

## More complex:

Dropout

Batch normalization

Stochastic depth, fractional pooling, etc

# So far: Linear Classifiers



Class  
scores

$$f(x) = Wx$$



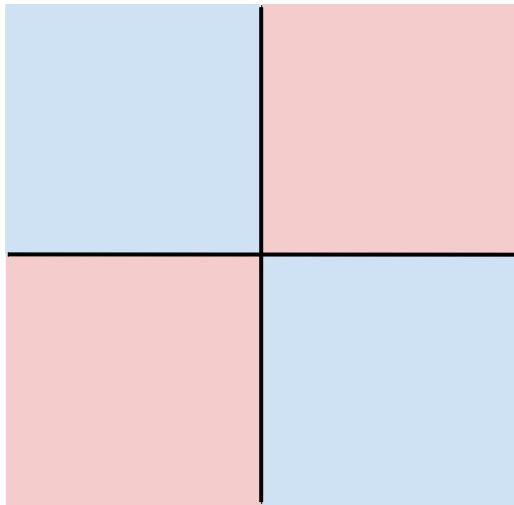
# Hard cases for a linear classifier

**Class 1:**

First and third quadrants

**Class 2:**

Second and fourth quadrants

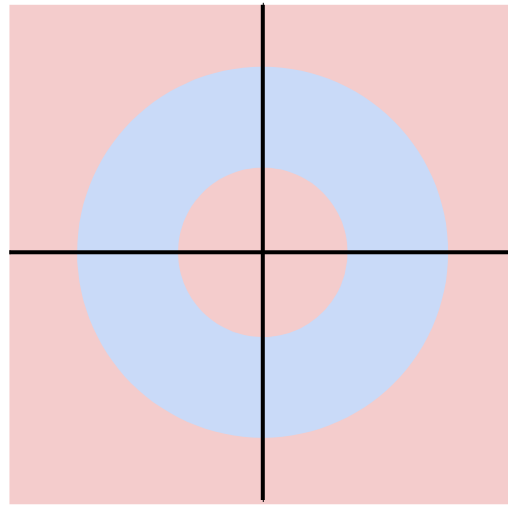


**Class 1:**

$1 \leq \text{L2 norm} \leq 2$

**Class 2:**

Everything else

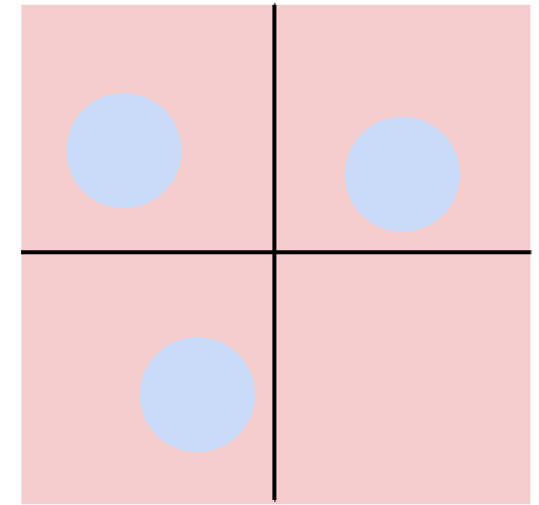


**Class 1:**

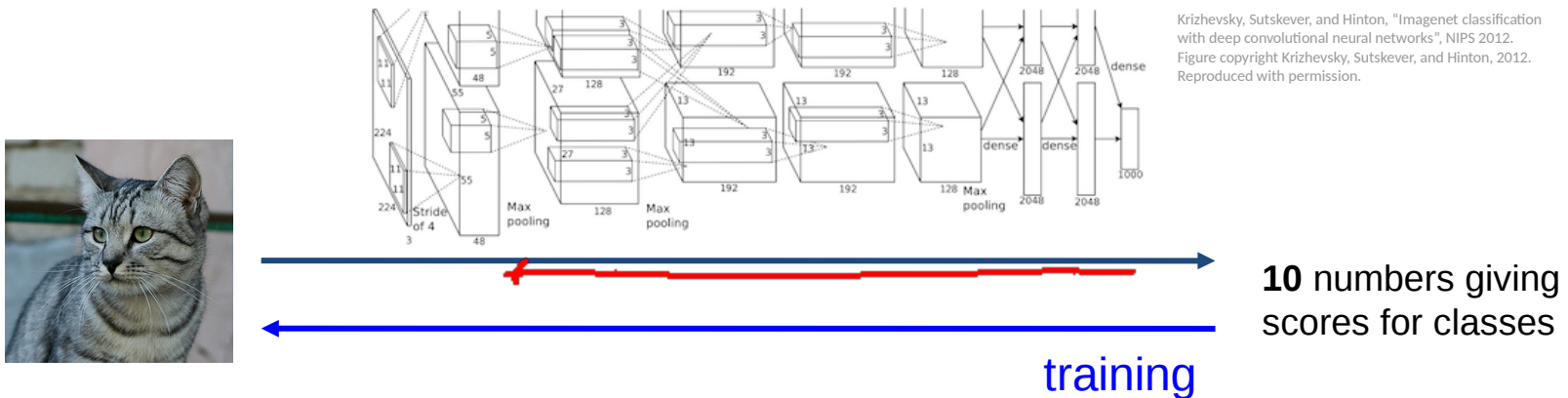
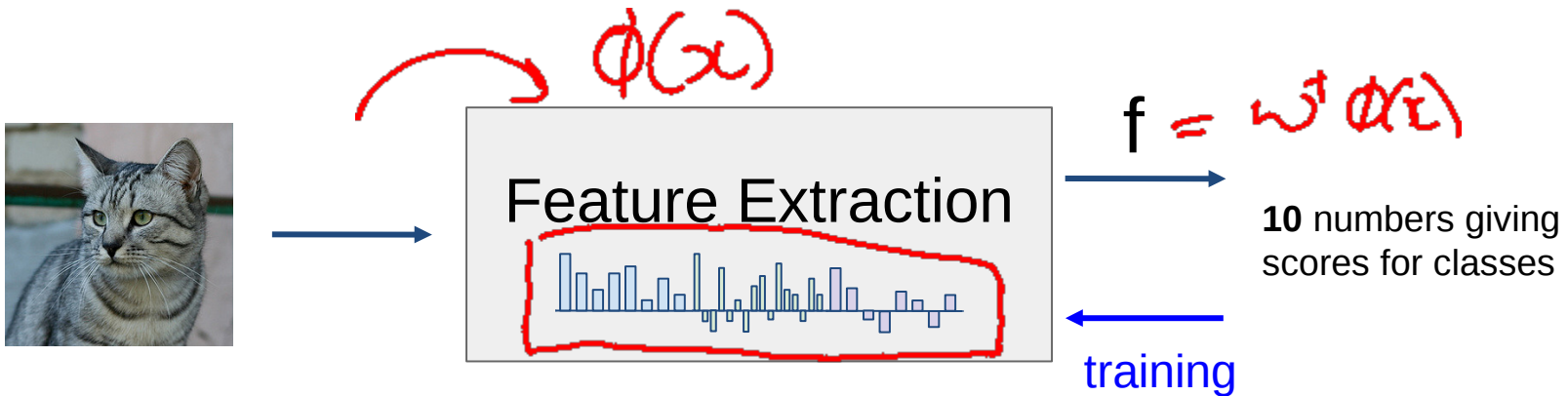
Three modes

**Class 2:**

Everything else



# Image features vs Neural Nets





# Neural networks: without the brain stuff

(**Before**) Linear score function:  $f = \underline{Wx}$

# Neural networks: without the brain stuff

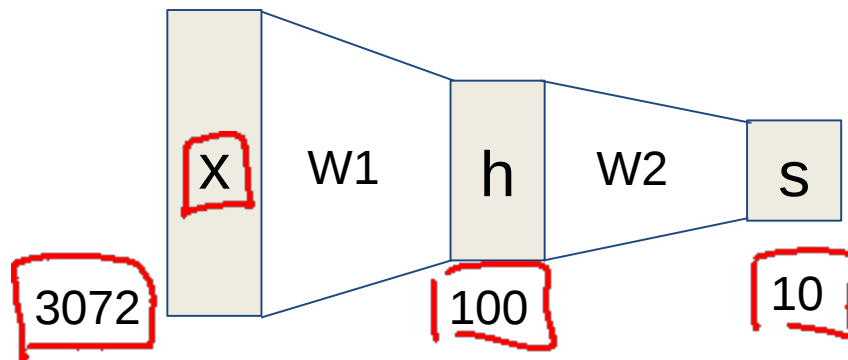
(**Before**) Linear score function:  $f = \overline{W}x$

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

# Neural networks: without the brain stuff

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

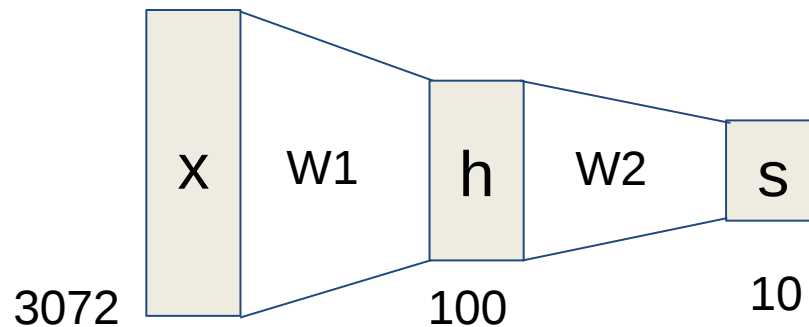
(**Now**) 2-layer Neural Network  $f = W_2 \max(0, W_1 x)$



# Neural networks: without the brain stuff

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

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



# Neural networks: without the brain stuff

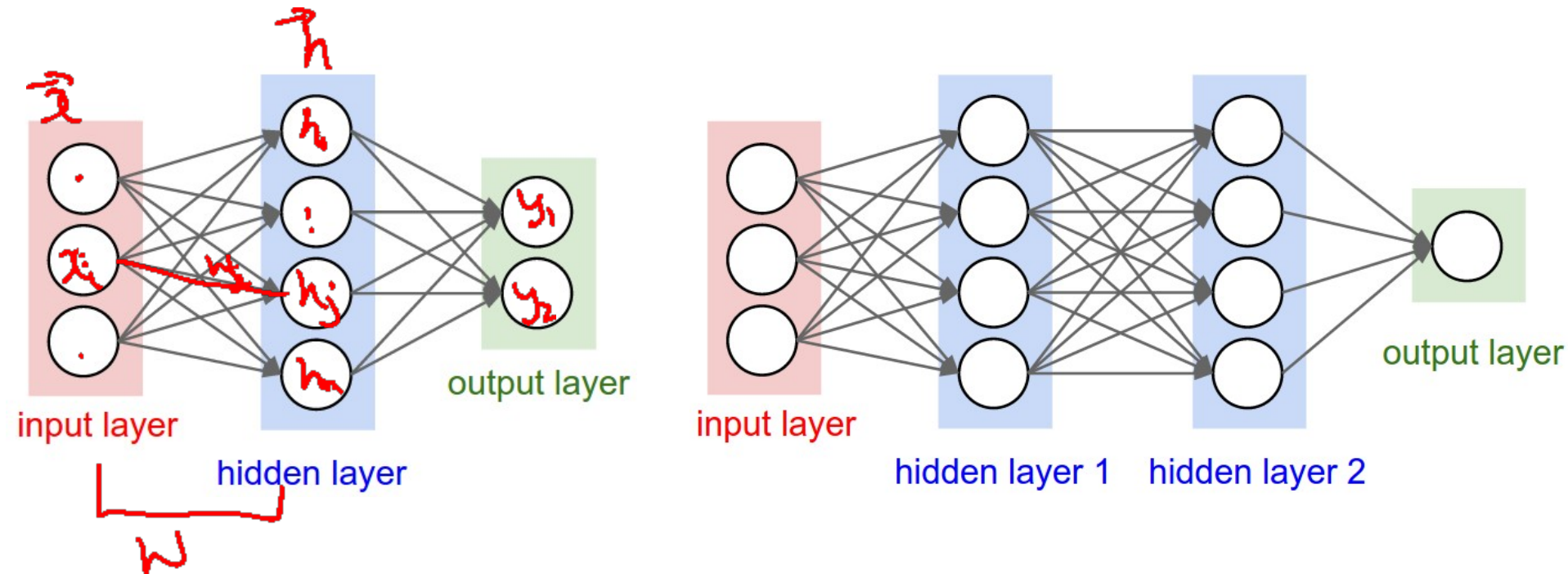
(**Before**) Linear score function:  $f = Wx$

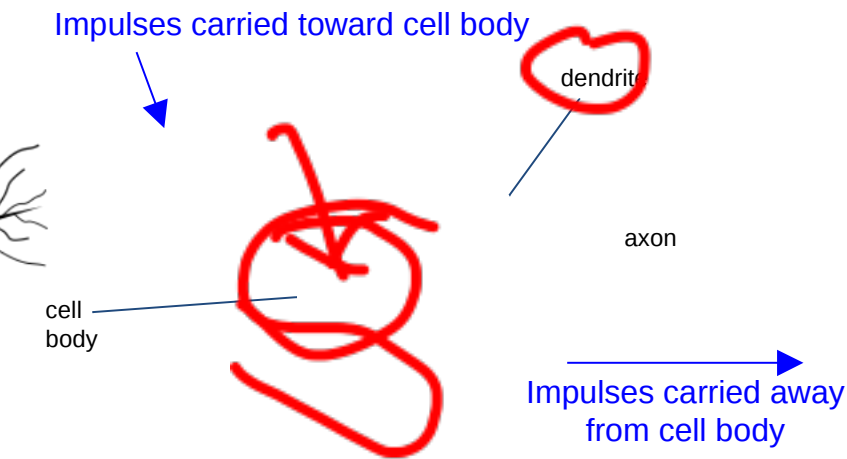
(**Now**) 2-layer Neural Network  
or 3-layer Neural Network  $f = W_2 \max(0, W_1 x)$

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

# Multilayer Networks

- Cascaded “neurons”
- The output from one layer is the input to the next
- Each layer has its own sets of weights

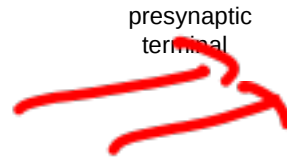




[This image](#) by Felipe Perucho is licensed under [CC-BY 3.0](#)

$$a = \sum_j w_j x_j$$

$$= \underline{w^T x}$$

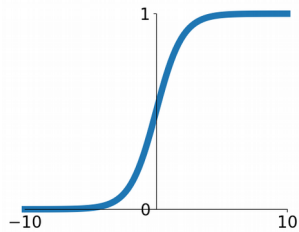


$$y = f(a)$$

# Activation functions

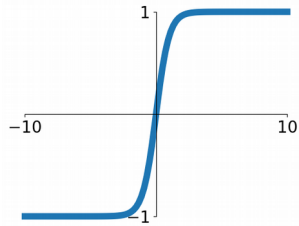
## Sigmoid

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



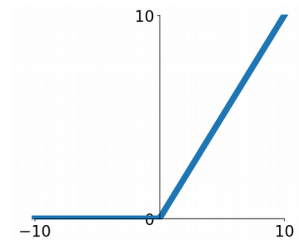
## tanh

$$\tanh(x)$$



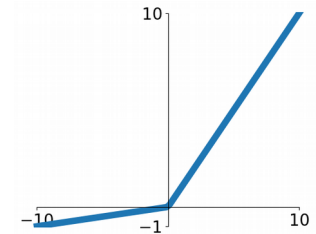
## ReLU

$$\max(0, x)$$



## Leaky ReLU

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

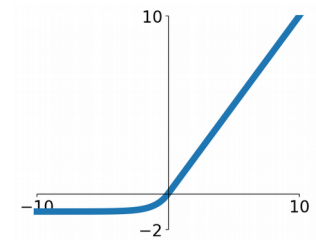


## Maxout

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

## ELU

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





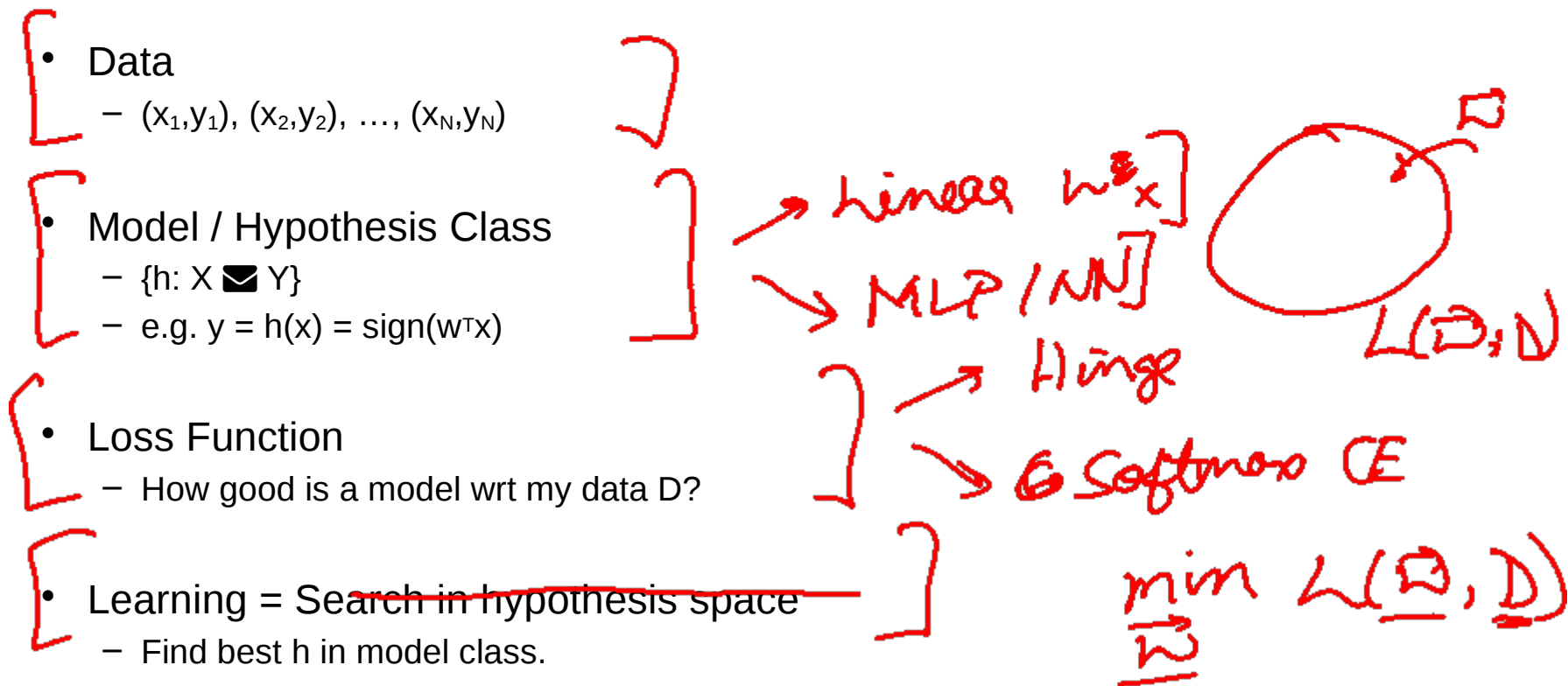
# Plan for Today

- Optimization
- Computing Gradients

# Optimization

# Supervised Learning

- Input:  $x$  (images, text, emails...)
- Output:  $y$  (spam or non-spam...)
- (Unknown) Target Function
  - $f: X \rightarrow Y$  (the “true” mapping / reality)



# Demo Time

- <https://playground.tensorflow.org>

$$\min_{\vec{w}} L(\vec{w})$$

Strategy: Follow the slope

Gradient  
Descent



$$L(\vec{w})$$

# What is slope?

- In 1-dimension the derivative of a function:

$$\left[ \frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(\underline{x} + \underline{h}) - f(\underline{x})}{\underline{h}} \right]$$

# What is slope?

$h(\vec{w})$

- In d-dimension, recall partial derivatives:

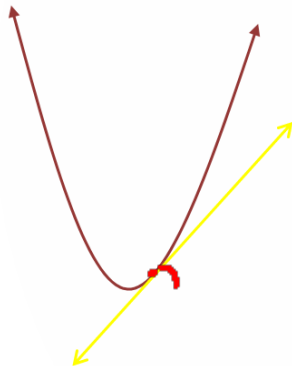
$$f: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$f(\vec{x}) \quad \vec{x} \in \mathbb{R}^d$$

↓

$$\boxed{\frac{\partial f}{\partial x_i}} = \lim_{h \rightarrow 0}$$

$$\frac{f(x_1, \dots, \underbrace{x_i + h}_{\downarrow h}, \dots, x_d) - f(x_1, \dots, x_d)}{h}$$



# What is slope?

- The **gradient** is the vector of (partial derivatives) along each dimension

$$f: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$\nabla f = \left[ \frac{\partial f}{\partial x_1} \quad \dots \quad \frac{\partial f}{\partial x_d} \right]$$

Convent:  
row-vs-  
col vector  
(convention)

- Properties

- The direction of steepest descent is the negative gradient
- The slope in any direction is the dot product of the direction with the gradient

$$\frac{\partial f}{\partial \vec{x}}$$

$$\nabla_{\vec{x}} f$$

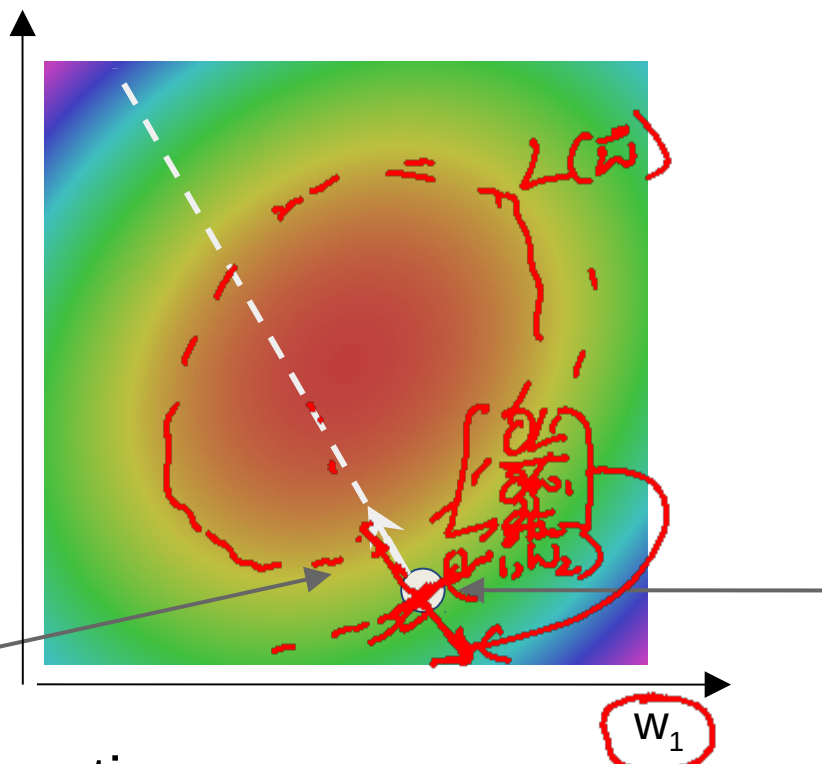
$$\nabla_{\vec{w}} L(\vec{w}, D)$$



$$L(\vec{w}) = L(w_1, w_2) = a \cdot w_1^2 + b w_2^2 + c w_1 w_2$$

bleue  
+ w  
red  
me

$w_2$



original W

$w_1$

negative gradient direction

min  $L(\vec{w})$   
 $\vec{w}$

# Gradient Descent

$$\min_{\vec{w}} L(\vec{w})$$

```
# Vanilla Gradient Descent
```

```
while True:
```

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

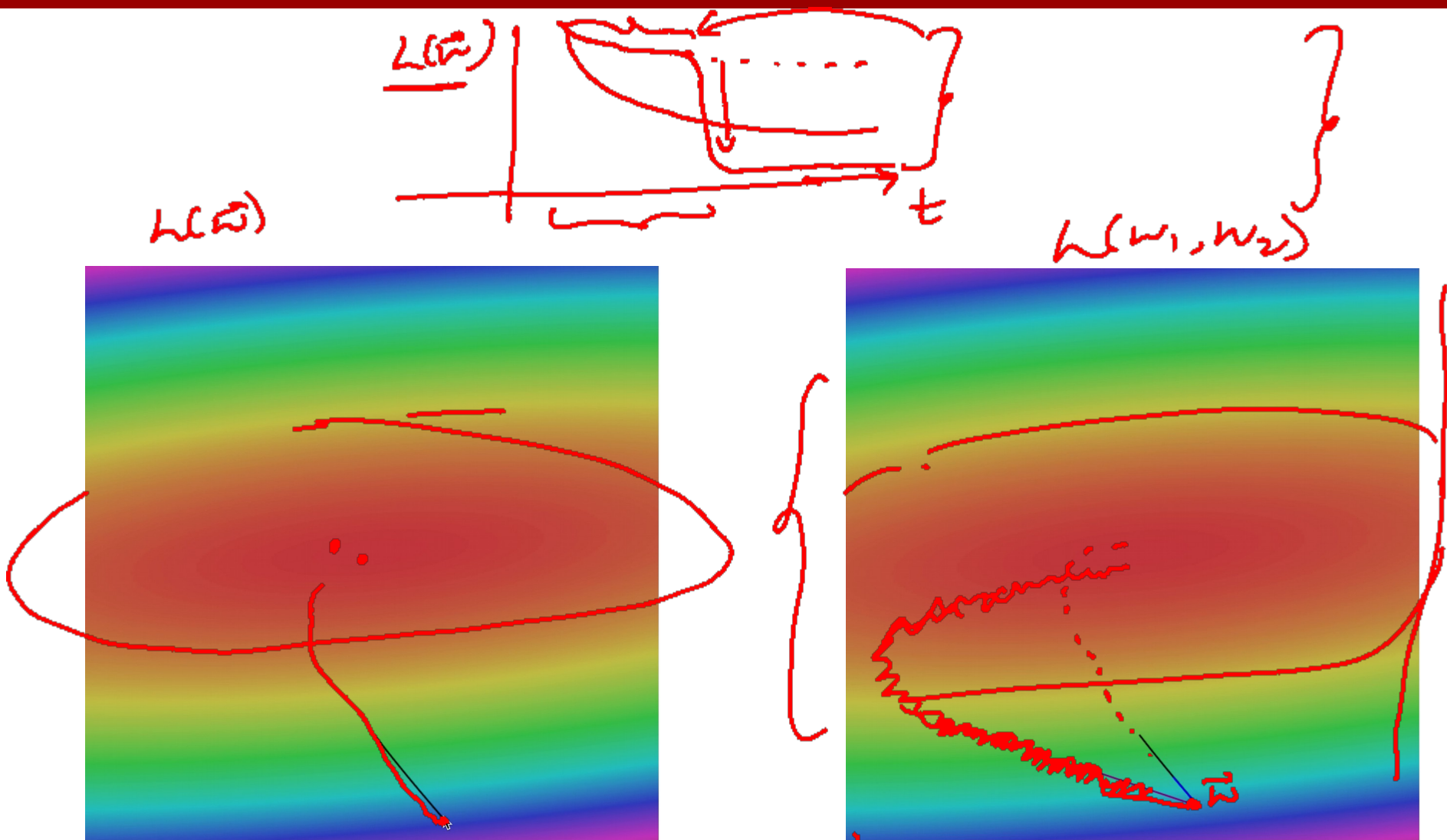
```
    weights += - step_size * weights_grad # perform parameter update
```

$\vec{w}^{(0)}$  = Initialize

for  $t = 1, 2, \dots$ , exhausted

$$\underline{w^{(t+1)}} = \underline{w^{(t)}} - \boxed{\eta} \underline{\nabla_{\vec{w}} L(\vec{w})} \quad |$$

0.001      Step-size / Learning-rate



# Gradient Descent has a problem

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

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

Full sum expensive  
when N is large!

# (Stochastic) Gradient Descent (SGD)

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

Full sum expensive  
when N is large!

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

Approximate sum ]  
using a **minibatch** of  
examples  
32 / 64 / 128 common !

```
# 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
```

# Stochastic Gradient Descent (SGD)

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W)$$

$$\boxed{\nabla_W L(W)} = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) = \text{mean-over-}N (\nabla_W L_i)_{i=1}^N$$

$\approx \text{mean-over-}B (\cdot)$

# Stochastic Gradient Descent (SGD)

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W)$$

"Actually rare about"

$$\approx \mathbb{E}_{(x, y) \sim p(x, y)} [L(w, x, y)]$$

$$\nabla_W \rightarrow \mathbb{E} [ \quad ]$$

$\tilde{B}_{\text{sample}}$

$$\mathbb{E}_{\tilde{B}} [ \nabla_W L(w, x, y) ]$$

# Stochastic Gradient Descent (SGD)

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

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

Full sum expensive  
when N is large!

Approximate sum  
using a **minibatch** of  
examples  
32 / 64 / 128 common

```
# 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
```



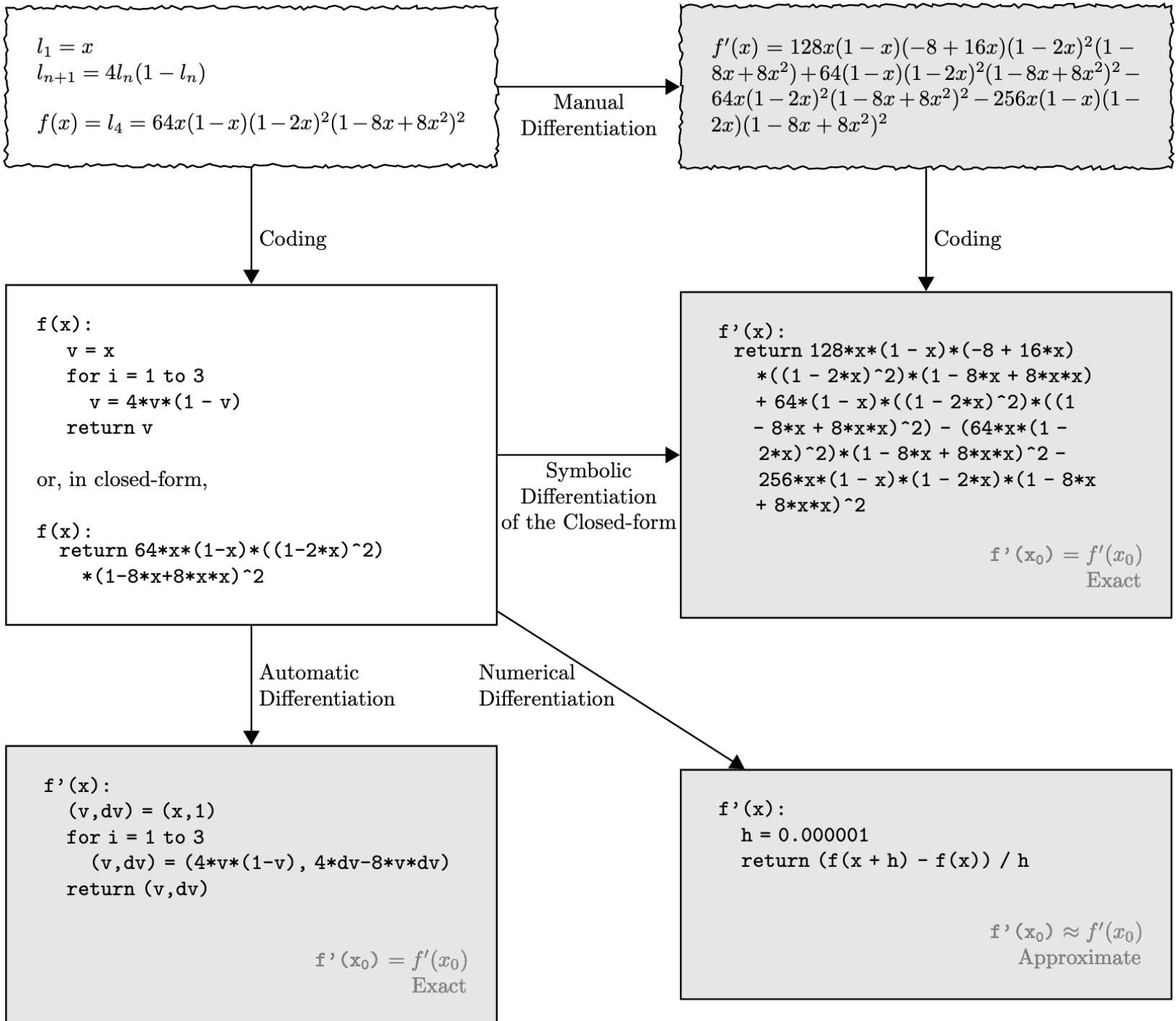


**GRADIENTS**

**GRADIENTS EVERYWHERE!**

# How do we compute gradients?

- Analytic or "Manual" Differentiation *"pen to paper"*
- Symbolic Differentiation
- Numerical Differentiation
- Automatic Differentiation
  - Forward mode AD
  - Reverse mode AD
    - aka "backprop"



$l_1 = x$   
 $l_{n+1} = 4l_n(1 - l_n)$   
 $f(x) = l_4 = 64x(1-x)(1-2x)^2(1-8x+8x^2)^2$

$\frac{df(x)}{dx}$

Manual  
Differentiation

$f'(x) = 128x(1-x)(-8+16x)(1-2x)^2(1-8x+8x^2) + 64(1-x)(1-2x)^2(1-8x+8x^2)^2 - 64x(1-2x)^2(1-8x+8x^2)^2 - 256x(1-x)(1-2x)(1-8x+8x^2)^2$

Coding

```

f(x):
  v = x
  for i = 1 to 3
    v = 4*v*(1 - v)
  return v
  
```

[Black box]  
Program

Coding

```

f'(x):
  return 128*x*(1-x)*(-8+16*x)
    *((1-2*x)^2)*(1-8*x+8*x*x)
    + 64*(1-x)*((1-2*x)^2)*((1-8*x+8*x*x)^2)
    - (64*x*(1-2*x)^2)*(1-8*x+8*x*x)^2
    - 256*x*(1-x)*(1-2*x)*(1-8*x+8*x*x)^2
  
```

Symbolic  
Differentiation  
of the Closed-form

or, in closed-form,

```

f(x):
  return 64*x*(1-x)*((1-2*x)^2)
    *(1-8*x+8*x*x)^2
  
```

$f'(x_0) = f'(x_0)$   
Exact

Automatic  
Differentiation

Numerical  
Differentiation

```

f'(x):
  
```

```

f'(x):
  
```

Coding

```
f(x):
  v = x
  for i = 1 to 3
    v = 4*v*(1 - v)
  return v
```

or, in closed-form,

```
f(x):
  return 64*x*(1-x)*((1-2*x)^2)
  *(1-8*x+8*x*x)^2
```

*Program*

Automatic  
Differentiation

```
f'(x):
  (v,dv) = (x,1)
  for i = 1 to 3
    (v,dv) = (4*v*(1-v), 4*dv-8*v*dv)
  return (v,dv)
```

*Program*

$\frac{\partial f}{\partial x}$

$$f'(x_0) = f'(x_0)$$

Exact

Symbolic  
Differentiation  
of the Closed-form

*Program*

Numerical  
Differentiation

*Program*

```
f'(x):
  h = 0.000001
  return (f(x + h) - f(x)) / h
```

$$f'(x_0) \approx f'(x_0)$$

Approximate

Coding

```
f'(x):
  return 128*x*(1 - x)*(-8 + 16*x)
  *((1 - 2*x)^2)*(1 - 8*x + 8*x*x)
  + 64*(1 - x)*((1 - 2*x)^2)*((1
  - 8*x + 8*x*x)^2) - (64*x*(1 -
  2*x)^2)*(1 - 8*x + 8*x*x)^2 -
  256*x*(1 - x)*(1 - 2*x)*(1 - 8*x
  + 8*x*x)^2
```

$$f'(x_0) = f'(x_0)$$

Exact



$$l_1 = x$$

$$l_{n+1} = 4l_n(1 - l_n)$$

$$f(x) = l_4 = 64x(1-x)(1-2x)^2(1-8x+8x^2)^2$$

$$f'(x) = 128x(1-x)(-8+16x)(1-2x)^2(1-8x+8x^2) + 64(1-x)(1-2x)^2(1-8x+8x^2)^2 - 64x(1-2x)^2(1-8x+8x^2)^2 - 256x(1-x)(1-2x)(1-8x+8x^2)^2$$

Manual  
Differentiation

Coding

Coding

```
f(x):
  v = x
  for i = 1 to 3
    v = 4*v*(1 - v)
  return v

or, in closed-form,

f(x):
  return 64*x*(1-x)*((1-2*x)^2)
  *(1-8*x+8*x*x)^2
```

```
f'(x):
  return 128*x*(1-x)*(-8+16*x)
  *((1-2*x)^2)*(1-8*x+8*x*x)
  + 64*(1-x)*((1-2*x)^2)*((1-8*x+8*x*x)^2)
  - (64*x*(1-2*x)^2)*(1-8*x+8*x*x)^2
  - 256*x*(1-x)*(1-2*x)*(1-8*x+8*x*x)^2
```

$f'(x_0) = f'(x_0)$   
Exact

Symbolic  
Differentiation  
of the Closed form

Automatic  
Differentiation

Numerical  
Differentiation

```
f'(x):
  (v,dv) = (x,1)
  for i = 1 to 3
    (v,dv) = (4*v*(1-v), 4*dv-8*v*dv)
  return (v,dv)
```

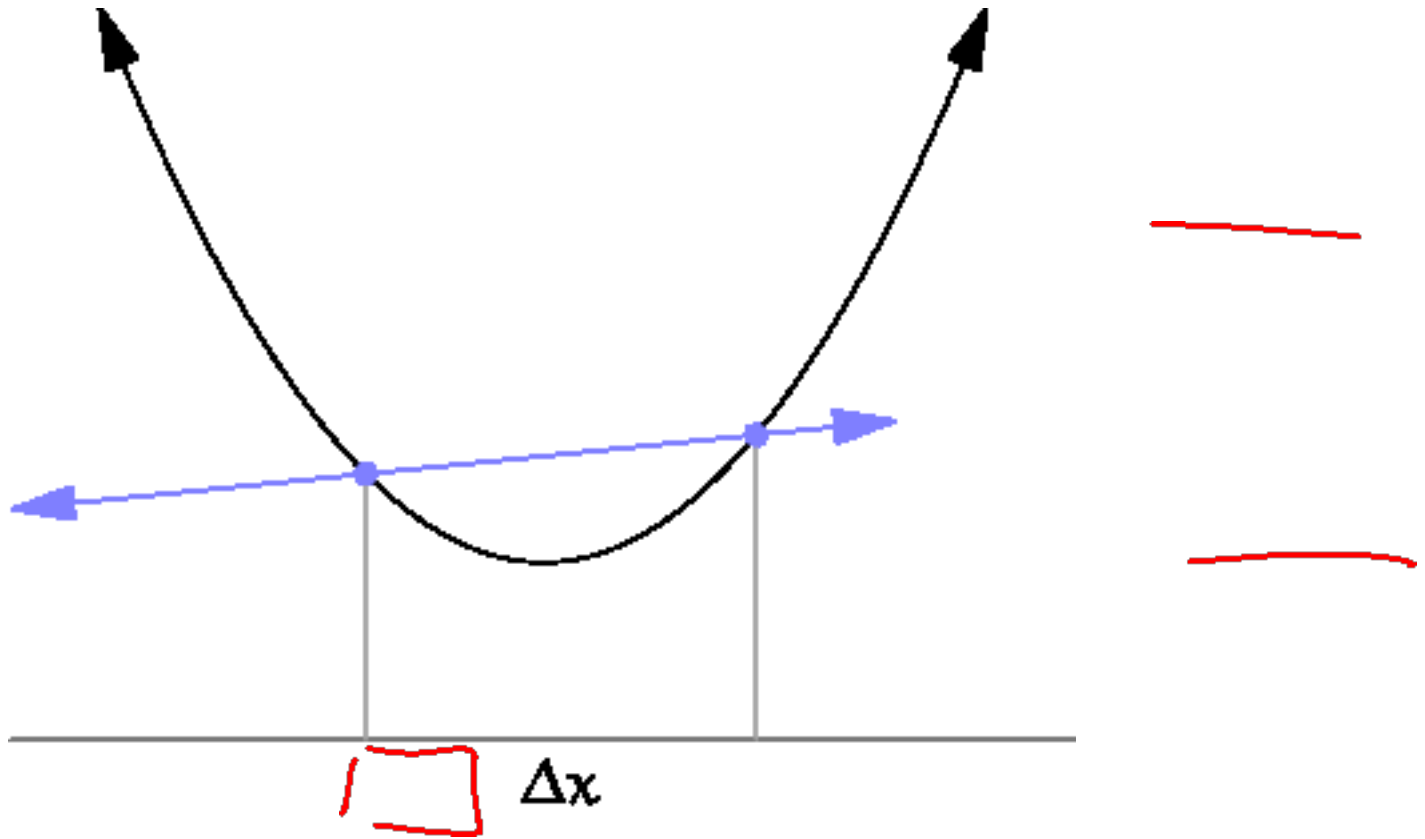
$f'(x_0) = f'(x_0)$   
Exact

```
f'(x):
  h = 0.000001
  return (f(x+h) - f(x)) / h
```

$f'(x_0) \approx f'(x_0)$   
Approximate

# How do we compute gradients?

- Analytic or “Manual” Differentiation
- Symbolic Differentiation
- Numerical Differentiation
- Automatic Differentiation
  - Forward mode AD
  - Reverse mode AD
    - aka “backprop”





current W:

W

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

loss **1.25347**

$L(\vec{w})$

gradient dW:

[?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,...]

current W:

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

loss 1.25347

W + h (first dim):

→ [0.34 + 0.0001<sup>h</sup>,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

loss 1.25322

$L(w_1+h, w_2, \dots, w_d)$

$$\frac{\partial L}{\partial w_1} = \frac{L(w_1+h, \dots, w_d) - L(w_1, \dots, w_d)}{h}$$

gradient dW:

[?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,...]

current W:

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

loss 1.25347

W + h (first dim):

[0.34 + **0.0001**,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

loss 1.25322

gradient dW:

[-2.5,  
?,  
?,

$$\frac{(1.25322 - 1.25347)}{0.0001} = -2.5$$

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

?,  
?,...]

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**W + h (second dim):**

[0.34,  
-1.11 + 0.0001,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25353**

**gradient dW:**

[-2.5,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,...]

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**W + h (second dim):**

[0.34,  
-1.11 + **0.0001**,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25353**

**gradient dW:**

[-2.5,  
0.6,  
?,  
?,



$$(1.25353 - 1.25347)/0.0001 = 0.6$$

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

?, ...]

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**W + h (third dim):**

[0.34,  
-1.11,  
0.78 + **0.0001**,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**gradient dW:**

[-2.5,  
0.6,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,...]

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

$L(\vec{w})$

**W + h (third dim):**

[0.34, ~~\_\_\_\_\_~~  
-1.11, ~~\_\_\_\_\_~~  
0.78 + **0.0001**, ~~\_\_\_\_\_~~  
0.12, ~~\_\_\_\_\_~~  
0.55, ~~\_\_\_\_\_~~  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

$\vec{w} \in \mathbb{R}^d$   $d = 18$

**gradient dW:**

[-2.5,  
0.6,  
0,  
?,  
...]

$$\frac{(1.25347 - 1.25347)}{0.0001} = 0$$

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

?, ...]

# Numerical vs Analytic Gradients

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

**Numerical gradient:** slow :(, approximate :(, easy to write :)

**Analytic gradient:** fast :) exact :) error-prone :(

In practice: Derive analytic gradient, check your implementation with numerical gradient.

This is called a **gradient check**.



# How do we compute gradients?

- Analytic or “Manual” Differentiation

- Symbolic Differentiation ~~X~~

- Numerical Differentiation

- Automatic Differentiation

- Forward mode AD
- Reverse mode AD
  - aka “backprop”

$$\frac{\partial L}{\partial \vec{w}} = [ \quad ]$$

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

~~dot~~  
1xd