# CS 4803 / 7643: Deep Learning

Topics:
- Optimization
- Computing Gradients

Dhruv Batra

Georgia Tech

# Administrativia

- HW1 Reminder
  - Due: 09/26, 11:55pm
  - https://www.cc.gatech.edu/classes/AY2020/cs7643_fall/Z3o9P26CwTPZZMDXyWYDj3/hw1.pdf
  - https://www.cc.gatech.edu/classes/AY2020/cs7643_fall/Z3o9P26CwTPZZMDXyWYDj3/hw1-q8/
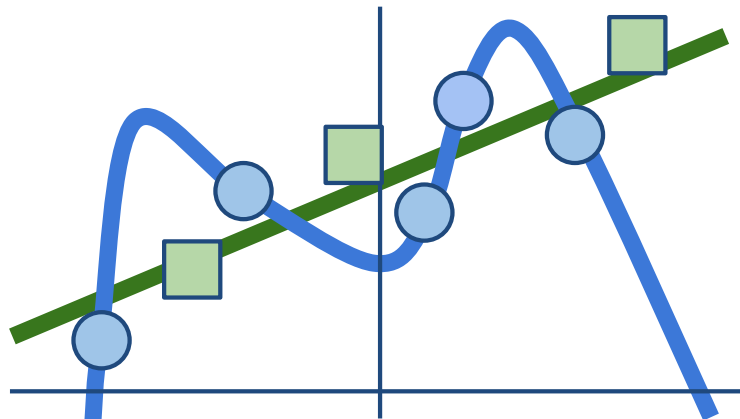  - https://evalai.cloudcv.org/web/challenges/challenge-page/431/leaderboard/1200

# Recap from last time

# Regularization

$\lambda$ = regularization strength (hyperparameter)

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

**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)

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

**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



$$f(x) = Wx$$
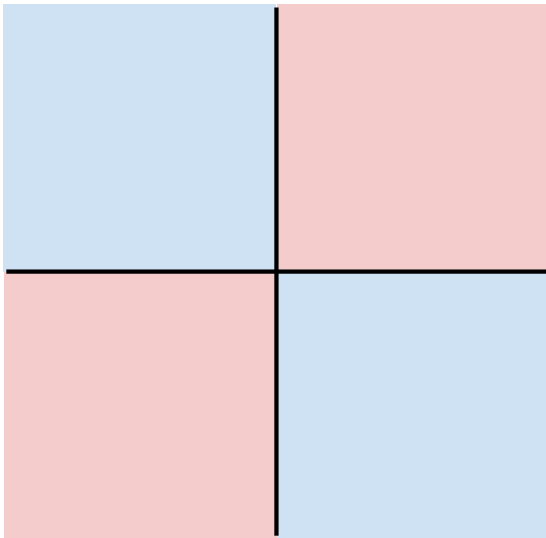
Class scores

# Hard cases for a linear classifier

**Class 1**:
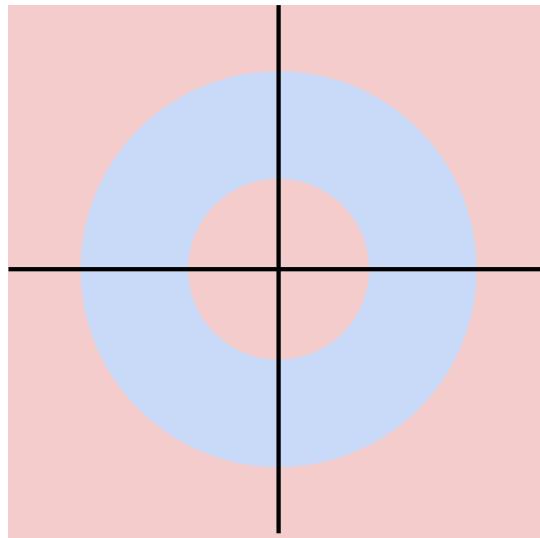First and third quadrants

**Class 2**:
Second and fourth quadrants

**Class 1**:
1 <= L2 norm <= 2

**Class 2**:
Everything else

**Class 1**:
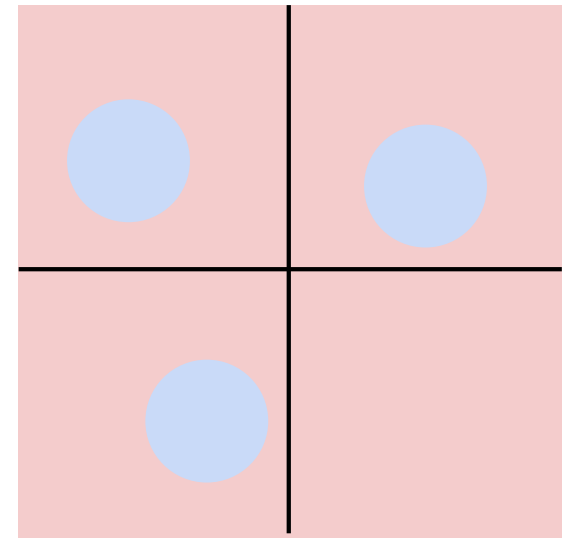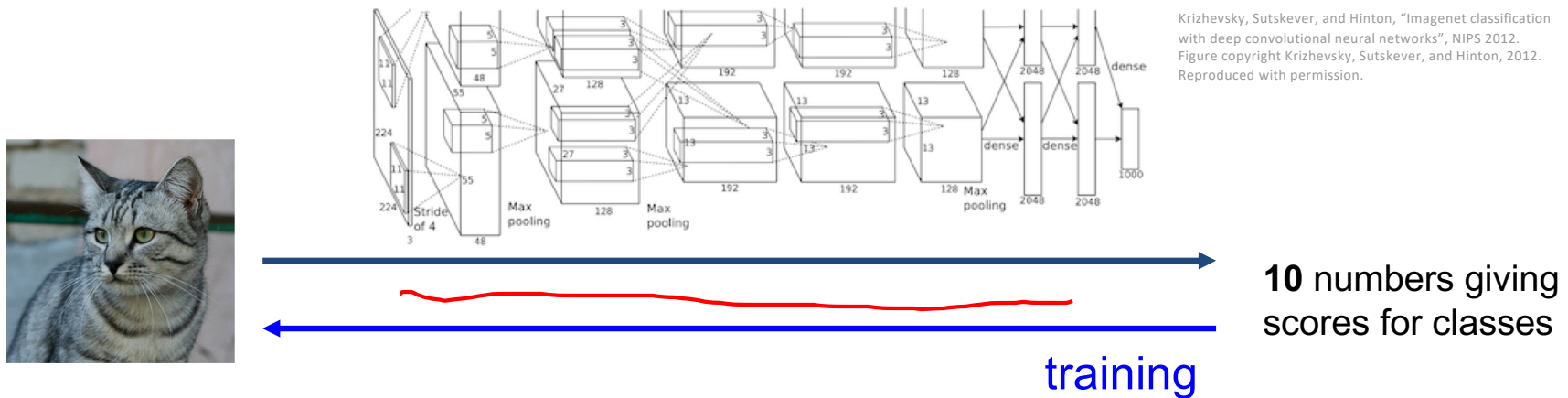Three modes

**Class 2**:
Everything else

# Image features vs Neural Nets



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: without the brain stuff

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

# Neural networks: without the brain stuff

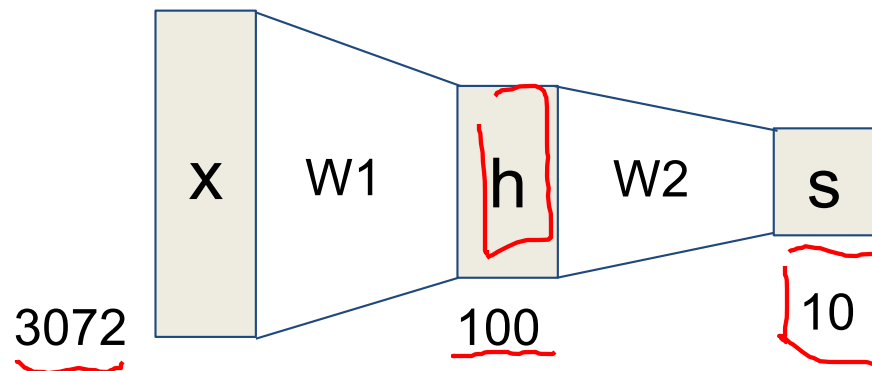(**Before**) Linear score function:  $f = Wx$

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

# Neural networks: without the brain stuff

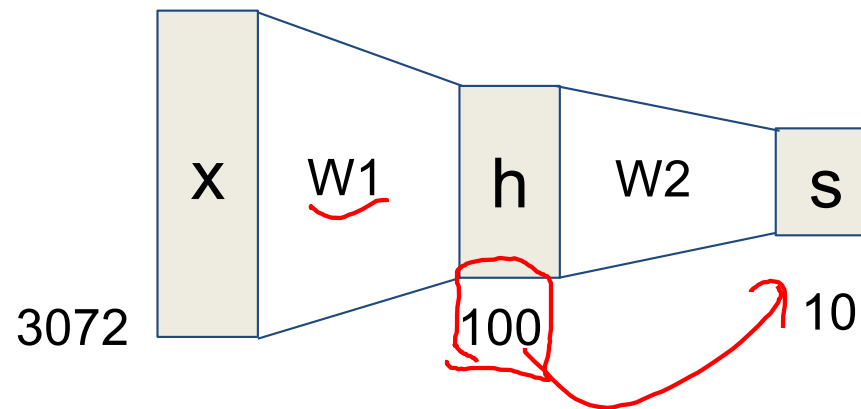(**Before**) Linear score function: $f = Wx$

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

# Neural networks: without the brain stuff

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

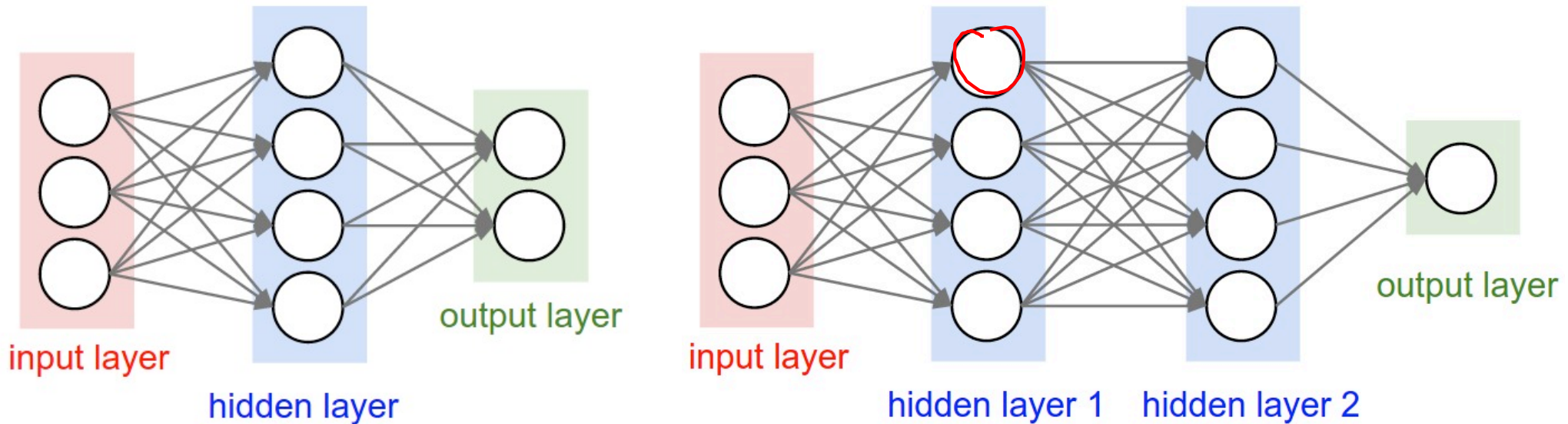(**Now**) 2-layer Neural Network $\quad f = W_2 \max(0, W_1 x)$
    or 3-layer Neural Network

$$f = W_3 \max(0, W_2 \max(0, W_1 x))$$
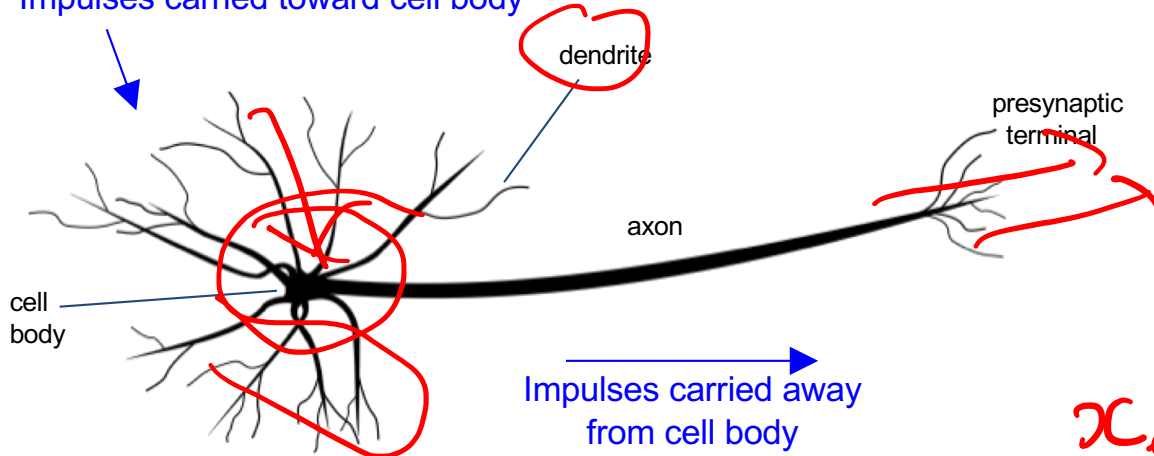
# Multilayer Networks  MLP

- Cascaded "neurons"
- The output from one layer is the input to the next
- Each layer has its own sets of weights
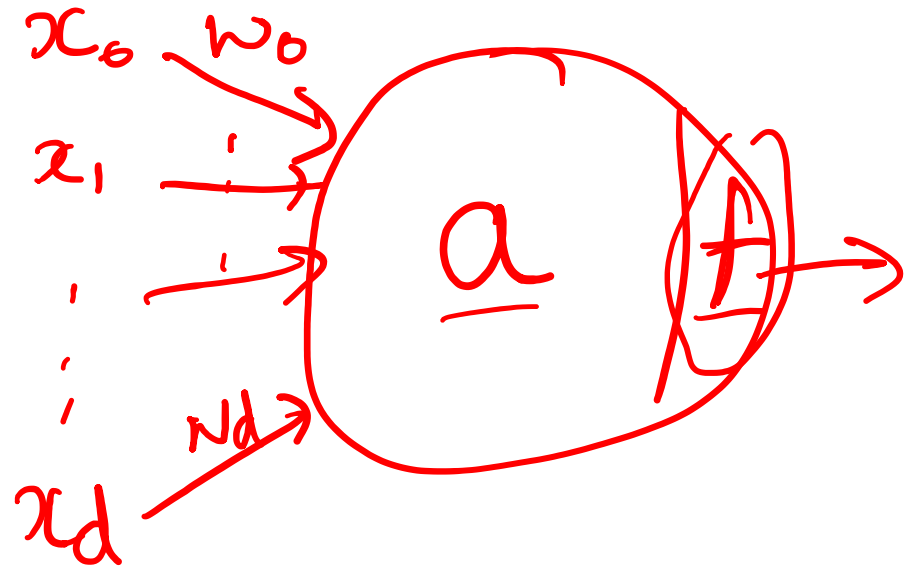
Impulses carried toward cell body

dendrite

presynaptic terminal

axon

cell body

Impulses carried away from cell body

This image by Felipe Perucho is licensed under CC-BY 3.0

$$a = \sum_j w_j x_j = \vec{w}^T \vec{x}$$

$$y = f(a)$$

Impulses carried toward cell body

dendrite

presynaptic terminal

axon

cell body

Impulses carried away from cell body

This image by Felipe Perucho is licensed under CC-BY 3.0

$x_0$    $w_0$

axon from a neuron    synapse

$w_0 x_0$

dendrite

cell body

$w_1 x_1$

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

$\sum_i w_i x_i + b$    $f$

output axon

activation function

$w_2 x_2$

sigmoid activation function

$$\frac{1}{1 + e^{-x}}$$

$\approx 1/2$

$$f(\alpha) \quad \frac{1}{1 + e^{-\alpha}}$$

Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

# Activation functions

**Sigmoid**

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

**tanh**

$$\tanh(x)$$

**ReLU**

$$\max(0, x)$$

**Leaky ReLU**

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

**Maxout**

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

**ELU**

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

# Activation Functions

- sigmoid vs tanh

$$\frac{1}{1+e^{-a}}$$

$$\to 2\sigma(2a) - 1$$

$$a = 10^5$$

$$\{z_i\}$$

$$\sigma\left(\sum \sigma\left(\sum \sigma(x)\right)\right)$$

# A quick note



**Fig. 4.** (a) Not recommended: the standard logistic function, $f(x) = 1/(1 + e^{-x})$. (b) Hyperbolic tangent, $f(x) = 1.7159\ \tanh\left(\frac{2}{3}x\right)$.

# Rectified Linear Units (ReLU)



[Krizhevsky et al., NIPS12]

# Plan for Today

- Optimization
- Computing Gradients

# Optimization

# Supervised Learning

- Input: x                                    (images, text, emails…)
- Output: y                                   (spam or non-spam…)

- (Unknown) Target Function
  - $f: X \to Y$                              (the "true" mapping / reality)

- Data
  - $(x_1, y_1), (x_2, y_2), \ldots, (x_N, y_N)$

- Model / Hypothesis Class         → Linear → NN
  - $\{h: X \to Y\}$
  - e.g. $y = h(x) = \text{sign}(w^T x)$

- Loss Function                          → hinge
  - How good is a model wrt my data D?      ↓ softmax CE

- Learning = Search in hypothesis space
  - Find best h in model class.

# Demo Time

- [https://playground.tensorflow.org](https://playground.tensorflow.org)

# Strategy: **Follow the slope**

$$\min_{\vec{w}} \; L(\vec{w}, D) \leftarrow$$

$$L(\vec{w}) = \frac{1}{N} \sum_i L_i(w)$$

# Strategy: **Follow the slope**

In 1-dimension, the derivative of a function:

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

$$\frac{\partial}{\partial x_i} f(x_1 \ldots x_d) = \lim_{h \to 0} \frac{f(x_1 \ldots x_i + h, \ldots x_d) - f(x_1, \ldots x_d)}{h}$$

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

Strategy: **Follow the slope**

In 1-dimension, the derivative of a function:

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

In multiple dimensions, the **gradient** is the vector of (partial derivatives) along each dimension

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

$$L(\vec{w}) = 0.00001 w_1^2 + w_2^2$$

W_2

W_1

original W

negative gradient direction

$(W_1, W_2)$

$\begin{bmatrix} \frac{\nabla L}{w_1} \\ \frac{\nabla L}{w_2} \end{bmatrix}$

$w_1$

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

# Gradient Descent

```
# Vanilla Gradient Descent

while True:
    weights_grad = evaluate_gradient(loss_fun, data, weights)
    weights += - step_size * weights_grad # perform parameter update
```

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

for $t = 1, \ldots, \text{tired}$

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \boxed{\eta} \, \boxed{\nabla_{\vec{w}} \mathcal{L}} \left( \vec{w}^{(t)} \right)$$

0.001

step-size / learning rate

Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

# 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)

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

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

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
```
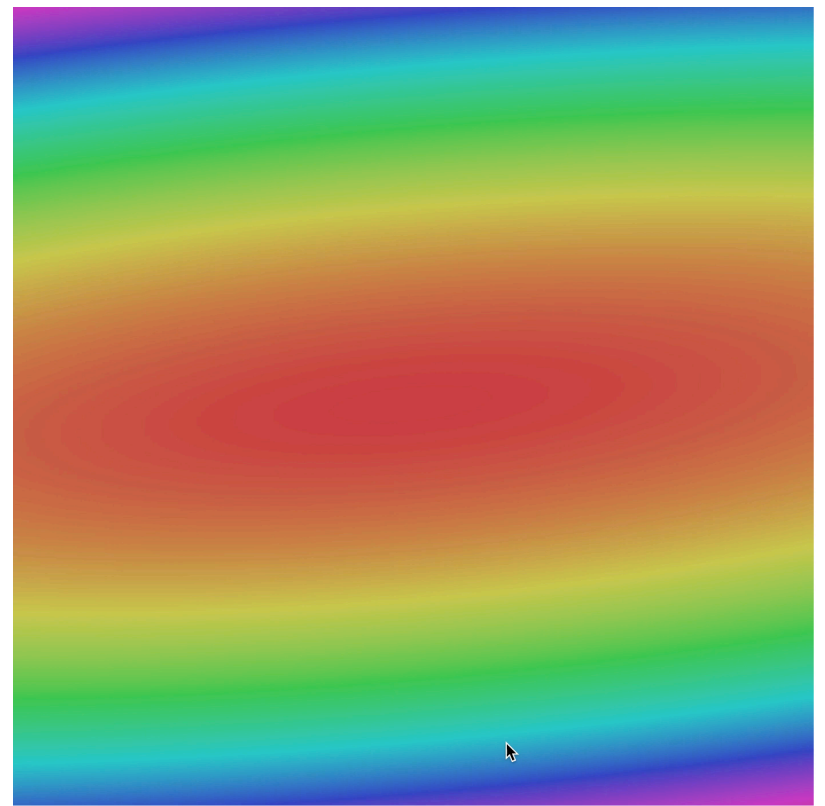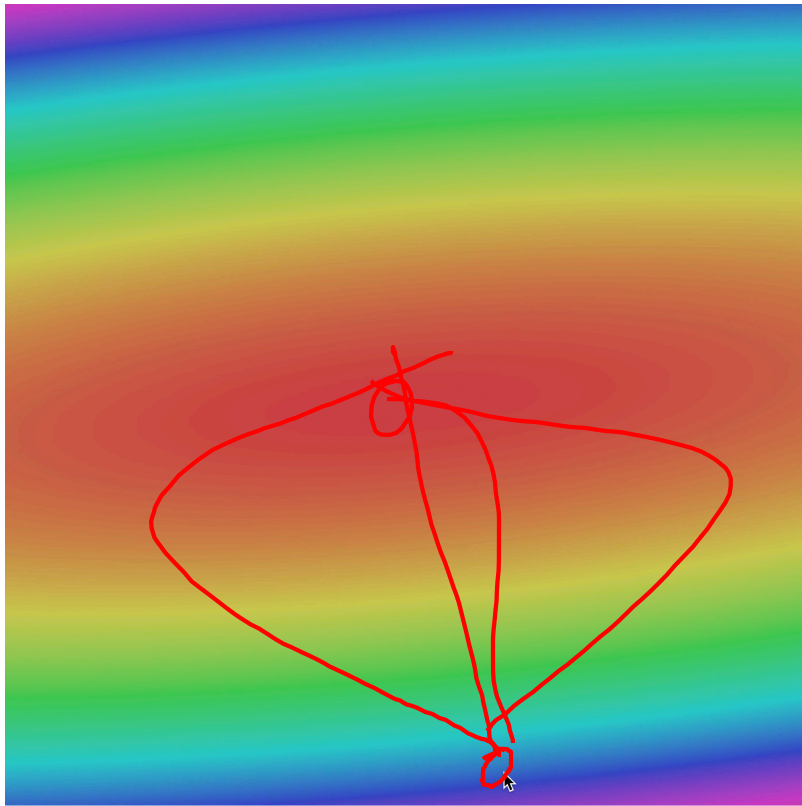
N = 1000

# Stochastic Gradient Descent (SGD)

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

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

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

# Stochastic Gradient Descent (SGD)

$$L(W) = \frac{1}{N}\sum_{i=1}^{N} L_i(x_i, y_i, W) \approx E_{x,y \sim p}\left[L(W, x, y)\right]$$

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

$$\nabla_W E\left[L(W, x, y)\right]$$

$$= E_p\left[\nabla_W L(W, x, y)\right]$$

$$B$$

# 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
```

# How do we compute gradients?

- Analytic or "Manual" Differentiation

- Symbolic Differentiation

- Numerical Differentiation

- Automatic Differentiation
  – Forward mode AD
  – Reverse mode AD
    - aka "backprop"

$l_1 = x$

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

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

**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
```

or, in closed-form,

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

**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

**Symbolic Differentiation of the Closed-form** →

**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)
```

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

Exact

**Numerical Differentiation** ↘

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

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

Approximate

$l_1 = x$

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

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

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

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
```

Symbolic Differentiation of the Closed-form

```
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'(x₀) = f'(x₀)

Exact

Automatic Differentiation

Numerical Differentiation

```
f'(x):
```

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
```

Symbolic
Differentiation
of the Closed-form

```
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

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

Box (top-left):

$$l_1 = x$$
$$l_{n+1} = 4l_n(1 - l_n)$$

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

— **Manual Differentiation** →

Box (top-right, shaded):

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

↓ **Coding** (left) ↓ **Coding** (right)

Box (middle-left):

```
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
```

— **Symbolic Differentiation of the Closed-form** →

Box (middle-right, shaded):

```
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
```

$\mathtt{f'(x_0)} = f'(x_0)$
Exact

↓ **Automatic Differentiation** (left) ↘ **Numerical Differentiation** (center)

Box (bottom-left, shaded):

```
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)
```

$\mathtt{f'(x_0)} = f'(x_0)$
Exact

Box (bottom-right, shaded):

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

$\mathtt{f'(x_0)} \approx f'(x_0)$
Approximate

# How do we compute gradients?

- Analytic or "Manual" Differentiation

- Symbolic Differentiation

- Numerical Differentiation

- Automatic Differentiation
  - Forward mode AD
  - Reverse mode AD
    - aka "backprop"

$$\lim_{h \to 0}$$

$\Delta x$

**current W:**

$\left(\frac{\partial L}{\partial w_1} \ldots \frac{\partial L}{\partial w_f}\right)$

**gradient dW:**

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

**loss 1.25347**

$L(w)$

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

**current W:**

[0.34,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
<span style="color:red">loss 1.25347</span>

$L(\vec{w})$

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

[0.34 + **0.0001**,
-1.11,
0.78,
0.12,
0.55,
2.81,
-3.1,
-1.5,
0.33,…]
<span style="color:red">loss 1.25322</span>

$L(w_1 + h, w_2 - wh)$

**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**,
?,
?,

(1.25322 - 1.25347)/0.0001
= -2.5

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

?,
?,…]

| current W: | W + h (second dim): | gradient dW: |
|---|---|---|
| [0.34, | [0.34, | [-2.5, |
| -1.11, | -1.11 + **0.0001**, | ?, |
| 0.78, | 0.78, | ?, |
| 0.12, | 0.12, | ?, |
| 0.55, | 0.55, | ?, |
| 2.81, | 2.81, | ?, |
| -3.1, | -3.1, | ?, |
| -1.5, | -1.5, | ?, |
| 0.33,…] | 0.33,…] | ?,…] |
| loss 1.25347 | loss 1.25353 | |

**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

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

?,…]

| current W: | W + h (third dim): | gradient dW: |
|---|---|---|
| [0.34, | [0.34, | [-2.5, |
| -1.11, | -1.11, | 0.6, |
| 0.78, | 0.78 + **0.0001**, | ?, |
| 0.12, | 0.12, | ?, |
| 0.55, | 0.55, | ?, |
| 2.81, | 2.81, | ?, |
| -3.1, | -3.1, | ?, |
| -1.5, | -1.5, | ?, |
| 0.33,…] | 0.33,…] | ?,…] |
| **loss 1.25347** | **loss 1.25347** | |

**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**

$2x \ L(\overline{w})$

**gradient dW:**

[-2.5,
0.6,
**0**,
?,

(1.25347 - 1.25347)/0.0001
= 0

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

?,…]

$2d \ |L(w)|$

# Numerical vs Analytic Gradients

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

**Numerical gradient**: slow :(, approximate :(, easy to write :)
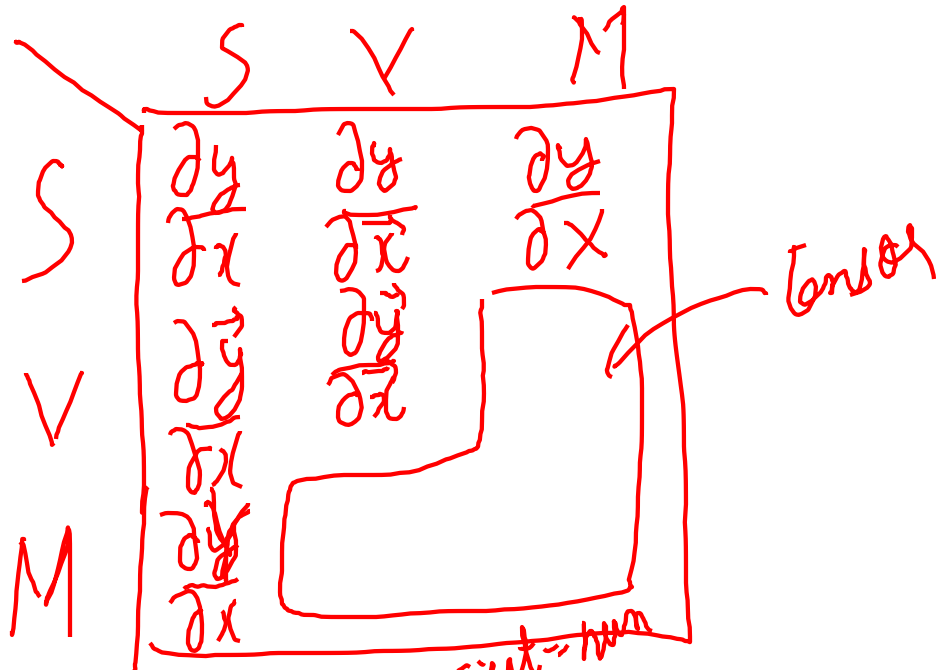**Analytic gradient**: fast :), exact :), error-prone :(

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

- Numerical Differentiation

- Automatic Differentiation
  - Forward mode AD
  - Reverse mode AD
    - aka "backprop"

# Vector/Matrix Derivatives Notation



$$\frac{\partial y}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \vdots \\ \frac{\partial y_c}{\partial x} \end{bmatrix} \quad c \times 1$$

$$x, y \in \mathbb{R}^1$$

$$\vec{x}, \vec{y} \in \mathbb{R}$$

$$\vec{x} \in \mathbb{R}^d \qquad y \in \mathbb{R}^c$$

$$X, Y \in \mathbb{R}^{m \times n}$$

$$num \equiv dim1 \equiv col$$

$$\frac{\partial y}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & - & - & - & - & \frac{\partial y}{\partial x_d} \end{bmatrix}$$

[Gradient]

# Vector/Matrix Derivatives Notation

$$\frac{\partial \vec{y}}{\partial \vec{x}} = \quad i \downarrow \begin{bmatrix} \frac{\partial y_i}{\partial x_1} & \cdots & \frac{\partial y_i}{\partial x_j} \\ & & \end{bmatrix}_{c \times d}$$

$$j \rightarrow \quad \frac{\partial y_i}{\partial x_j}$$

# Vector Derivative Example

$$\vec{y} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} x \\ x^2 \end{bmatrix} \qquad \frac{\partial \vec{y}}{\partial x} = \begin{bmatrix} 1 \\ 2x \end{bmatrix}$$

$$y = \vec{w}^T \vec{x} \qquad \frac{\partial y}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & - - - - - & \frac{\partial y}{\partial x_d} \end{bmatrix}$$

$$= \sum_{i=1}^{d} w_i x_i$$

$$\frac{\partial (\vec{w}^T \vec{x})}{\partial \vec{x}} \qquad \frac{\partial (\sum w_i x_i)}{\partial x_1}$$

$$\begin{bmatrix} w_1 & - - - - - & w_d \end{bmatrix} \quad \vec{w}^T$$

# Vector Derivative Example

$$\frac{\partial(\vec{w}^T A \vec{x})}{\partial \vec{w}} = 2 \vec{w}^T A$$

$$y_i = \sum_j a_{ij} x_j$$

$$\vec{x} \in \mathbb{R}^d$$
$$\vec{y} = \mathbb{R}^c$$
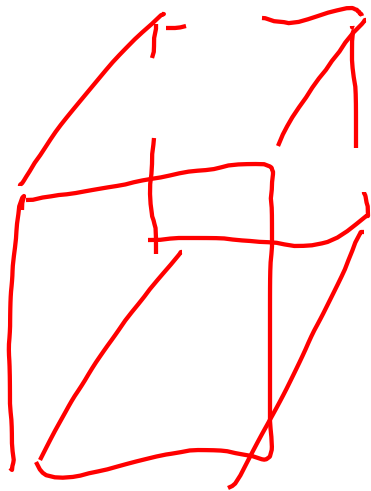$$A \in \mathbb{R}^{c \times d}$$

$$\left[\vec{y}\right] = A \vec{x}$$

$$\left[\frac{\partial \vec{y}}{\partial \vec{x}}\right] = \left[A\right]_{=ij} \quad \left[\frac{\partial y_i}{\partial x_j}\right] = a_{ij}$$

# Extension to Tensors

$$X \in R^{d_1 \ldots d_m}$$

$$Y \in R^{c_1 \ldots c_n}$$

$$y\text{-vec} = Y(:)$$

$$x\text{-vec} = X(:)$$

$$\frac{\partial Y[i_1 \ldots i_n]}{\partial X[j_1 - j_m]}$$

$$\frac{\partial y\text{-vec}}{\partial x\text{-vec}} = \left[ \qquad \right]$$

# Chain Rule: Composite Functions

$$f(g(x)) = (f \circ g)(x)$$

$$L(w) = g_\ell \big( g_{\ell-1} \big( g_{\ell-2} \cdots \big( g_1(w) \big) \big) \big)$$

$$\frac{\partial L}{\partial w} = \big( g_\ell \circ \cdots g_1 \big)(w)$$