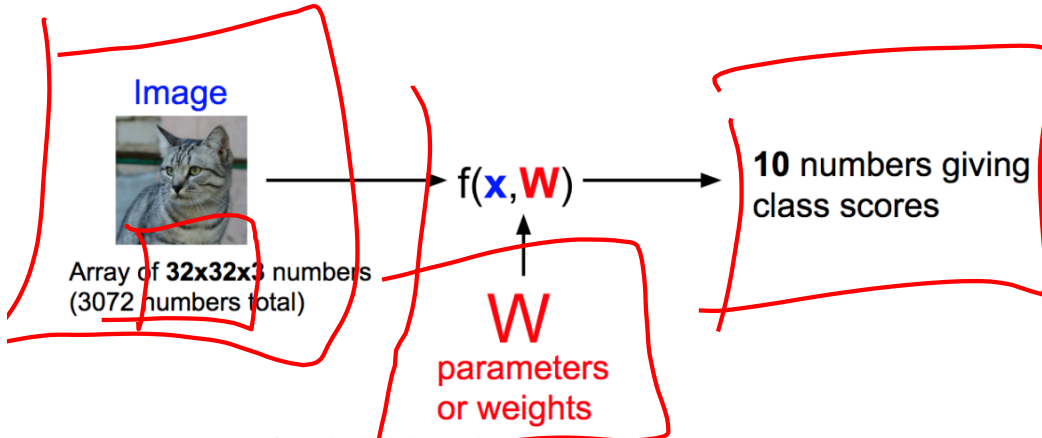# CS 7643: Deep Learning

Topics:
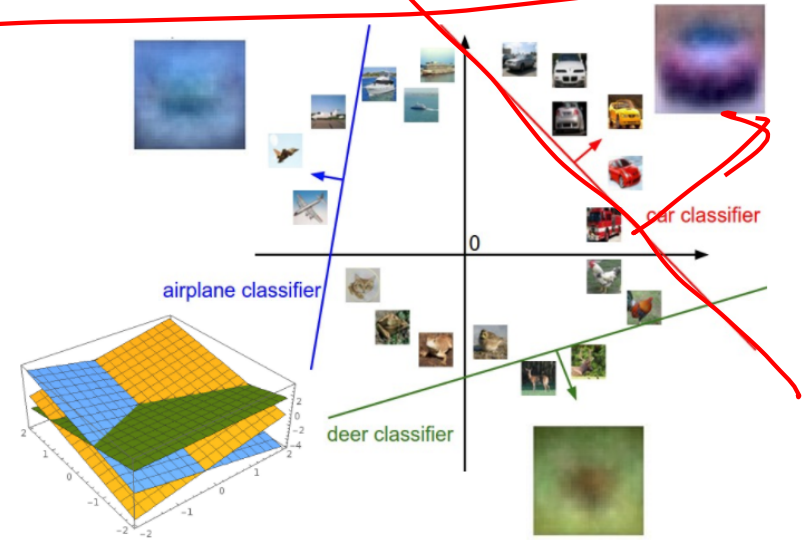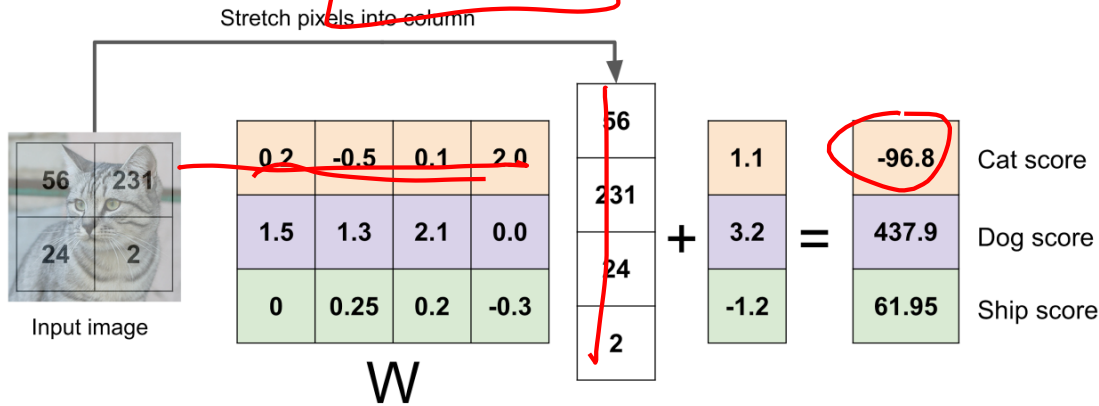- Regularization
- Neural Networks
  - Modular Design
- Computing Gradients

Dhruv Batra

Georgia Tech

# **Recall from last time**: Linear Classifier



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

# **Recall from last time**: Linear Classifier

|           | (cat)  | (car)  | (frog) |
|-----------|--------|--------|--------|
| airplane  | -3.45  | -0.51  | 3.42   |
| automobile| -8.87  | **6.04** | 4.64 |
| bird      | 0.09   | 5.31   | 2.65   |
| cat       | **2.9** | -4.22 | 5.1    |
| deer      | 4.48   | -4.19  | 2.64   |
| dog       | 8.02   | 3.58   | 5.55   |
| frog      | 3.78   | 4.49   | **-4.34** |
| horse     | 1.06   | -4.37  | -1.5   |
| ship      | -0.36  | -2.09  | -4.79  |
| truck     | -0.72  | -2.93  | 6.14   |

$L_i(W)$

TODO:

1. Define a **loss function** that quantifies our unhappiness with the scores across the training data.

1. Come up with a way of efficiently finding the parameters that minimize the loss function. **(optimization)**
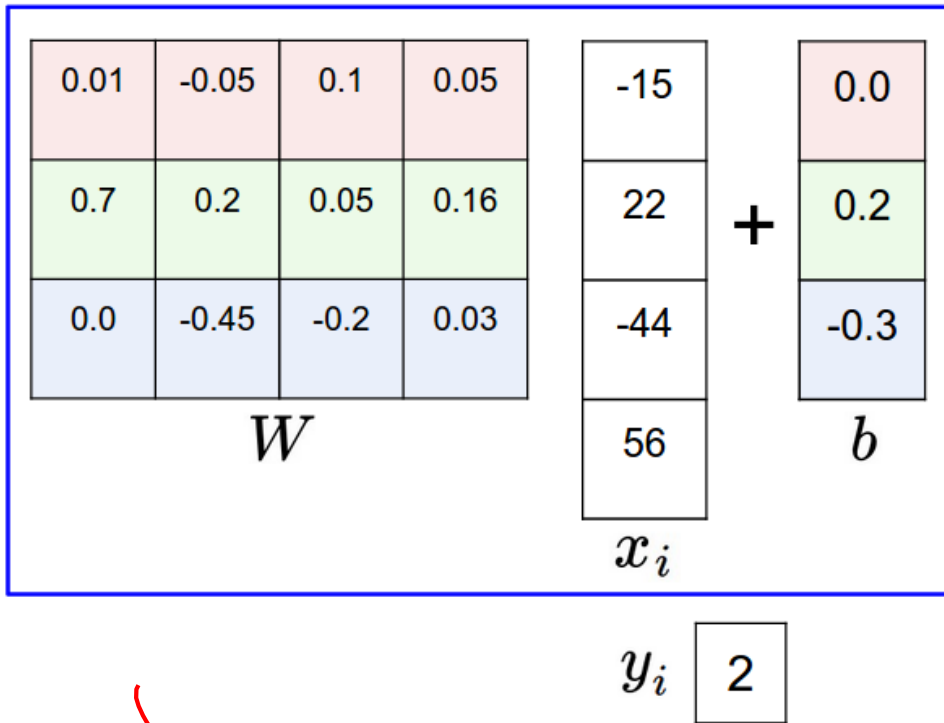
# Softmax vs. SVM

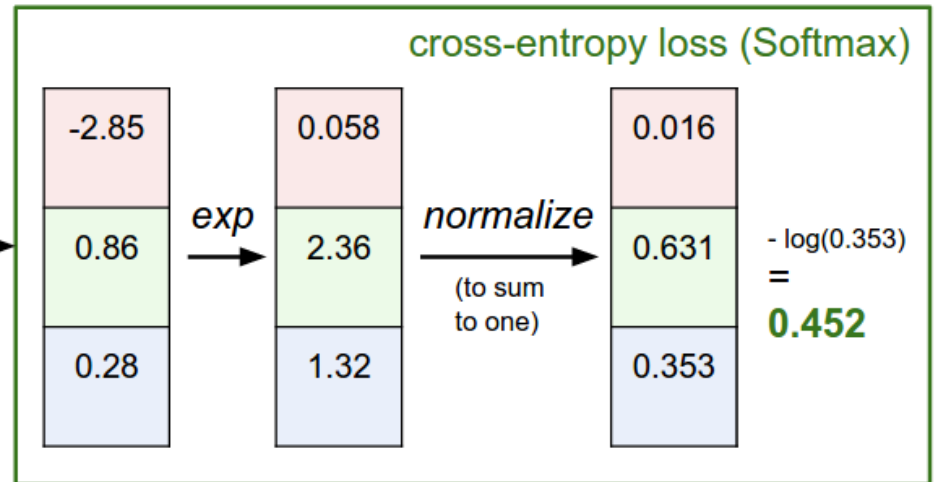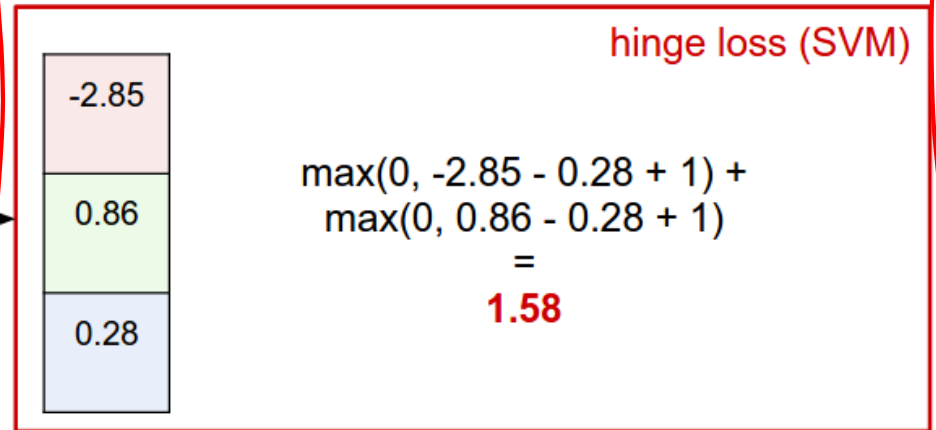$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right) \qquad L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$\log P(Y = y_i \mid W, \vec{x}_i)$

Model

matrix multiply + bias offset

| 0.01 | -0.05 | 0.1 | 0.05 |
|------|-------|------|------|
| 0.7 | 0.2 | 0.05 | 0.16 |
| 0.0 | -0.45 | -0.2 | 0.03 |

$W$

| -15 |
| 22 |
| -44 |
| 56 |

$x_i$

$+$

| 0.0 |
| 0.2 |
| -0.3 |

$b$

$y_i$ | 2 |

Loss

hinge loss (SVM)

| -2.85 |
| 0.86 |
| 0.28 |

max(0, -2.85 - 0.28 + 1) +
max(0, 0.86 - 0.28 + 1)
=
**1.58**

cross-entropy loss (Softmax)

| -2.85 |
| 0.86 |
| 0.28 |

*exp*
→

| 0.058 |
| 2.36 |
| 1.32 |

*normalize*
→
(to sum
to one)

| 0.016 |
| 0.631 |
| 0.353 |

- log(0.353)
=
**0.452**

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

# Plan for Today

- Regularization
- Neural Networks
  - Modular Design
- Computing Gradients

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

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

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

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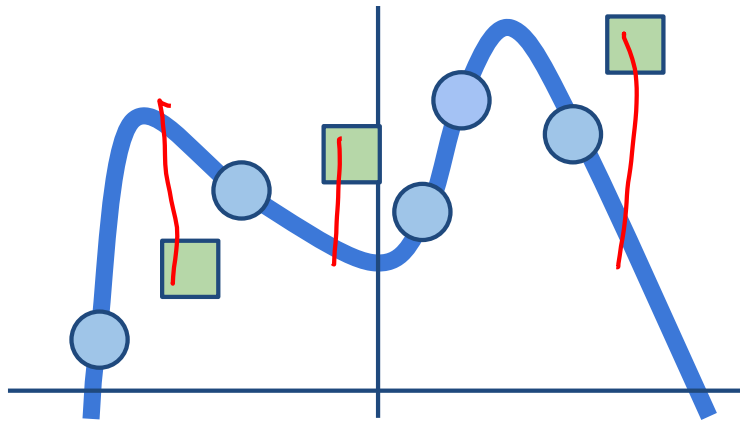should match training data

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

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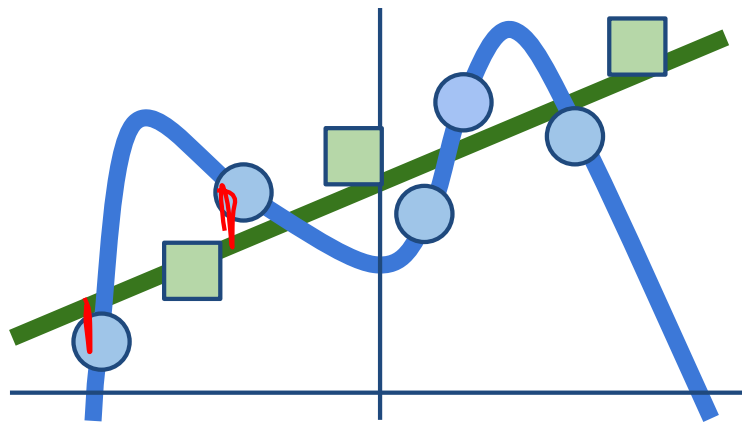should match training data

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

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

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

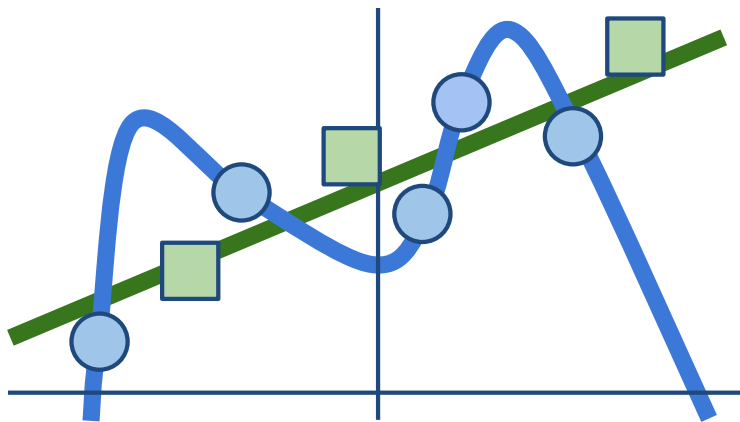*(handwritten annotations: "loss(" over the first term, "log P(W)" over the regularization term, box around $R(W)$)*

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

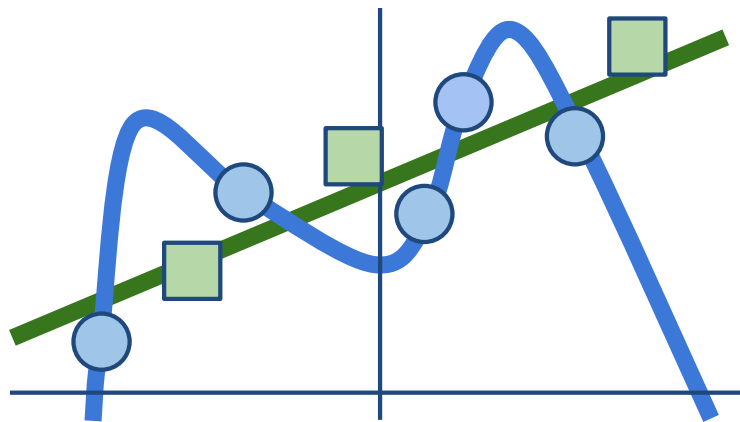**Regularization**: Model should be "simple", so it works on test data

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

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

**Regularization**: Model should be "simple", so it works on test data

**Occam's Razor**:
*"Among competing hypotheses, the simplest is the best"*
William of Ockham, 1285 - 1347

# Regularization

$\lambda$ = regularization strength (hyperparameter)

$$L = \frac{1}{N} \sum_{i=1}^{N} \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \boxed{\lambda R(W)}$$

In common use:

**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)  $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

Dropout (will see later)

Fancier: Batch normalization, stochastic depth

# L2 Regularization (Weight Decay)

$$x = [1, 1, 1, 1]$$

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

$$w_1 = [1, 0, 0, 0]$$

$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

$$w_1^T x = w_2^T x = 1$$

# L2 Regularization (Weight Decay)

$$x = [1, 1, 1, 1]$$

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

$$w_1 = [1, 0, 0, 0]$$

$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

(If you are a Bayesian: L2 regularization also corresponds MAP inference using a Gaussian prior on W)
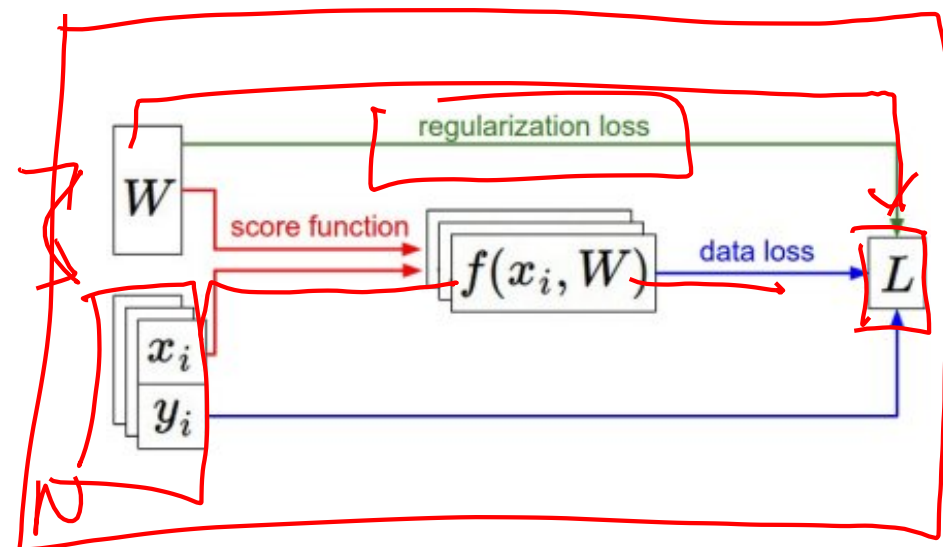
$$w_1^T x = w_2^T x = 1$$

# Recap

- We have some dataset of (x,y)
- We have a **score function:** $s = f(x; W) \overset{\text{e.g.}}{=} Wx$
- We have a **loss function**:

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

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

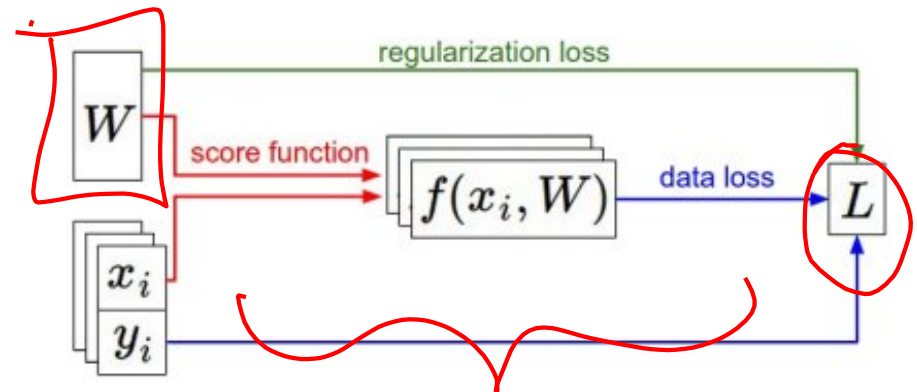$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + R(W)$$ Full loss

# Recap

How do we find the best W?

- We have some dataset of (x,y)
- We have a **score function:** $s = f(x; W) \overset{\text{e.g.}}{=} Wx$
- We have a **loss function**:

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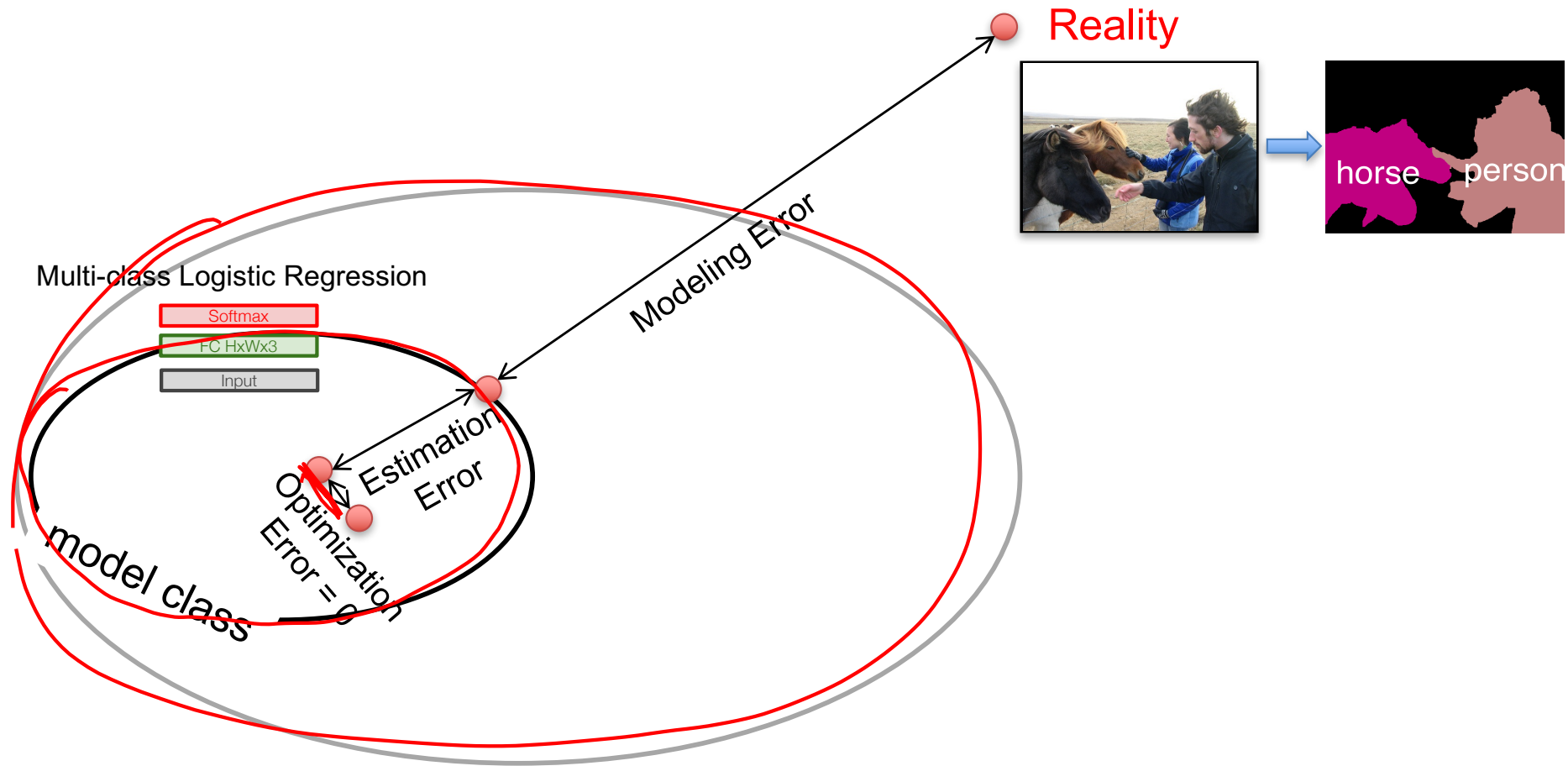
$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$

SVM

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

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

# Error Decomposition



Reality

Multi-class Logistic Regression

Softmax

FC HxWx3

Input

model class

Modeling Error

Estimation Error

Optimization Error = 0

horse    person

# Next: Neural Networks

# Neural networks: without the brain stuff

(**Before**) Linear score function: $f = Wx$
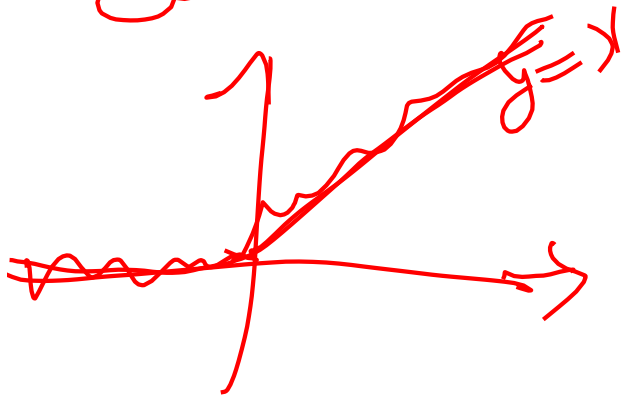
# Neural networks: without the brain stuff

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

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

$g_3(\vec{x}) = W_2 \vec{x}$

$= g_3(g_2(g_1(\omega)))$

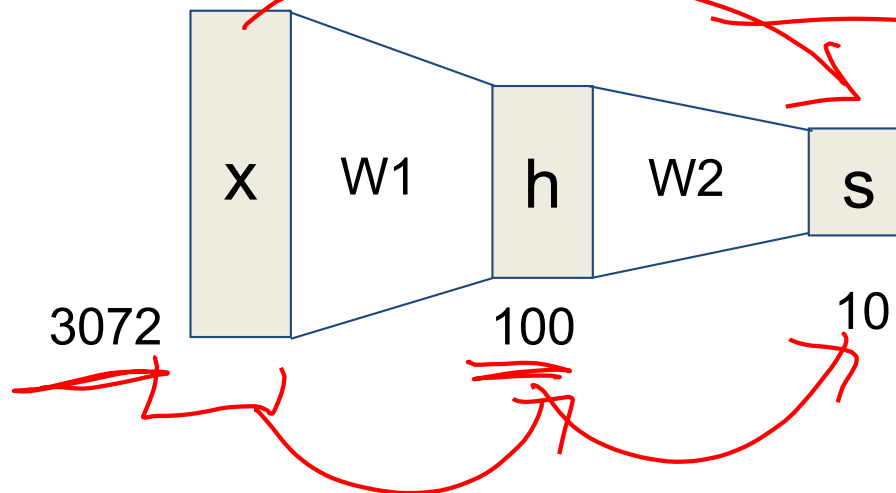$g_1(\vec{x}) = W\vec{x}$

$g_2(\vec{x}) = \max\{0, \vec{x}\}$ ReLu

# Neural networks: without the brain stuff

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

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



x  W1  h  W2  s

3072  100  10

# Neural networks: without the brain stuff

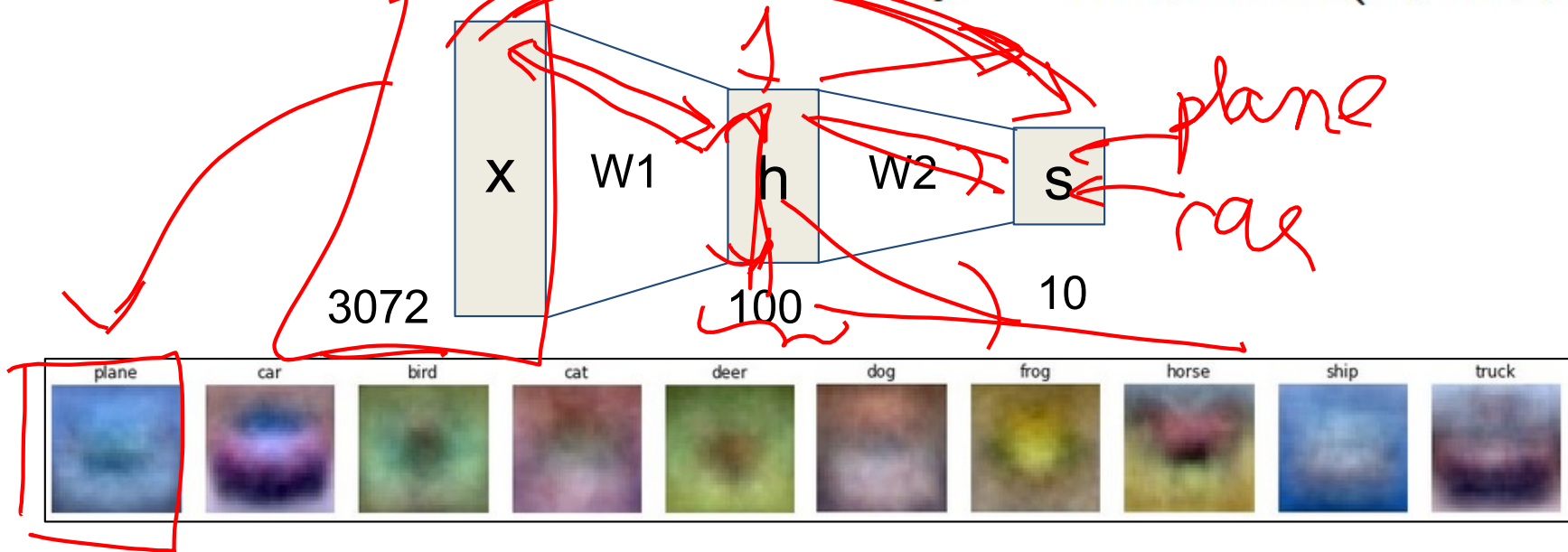(**Before**) Linear score function: $f = Wx$

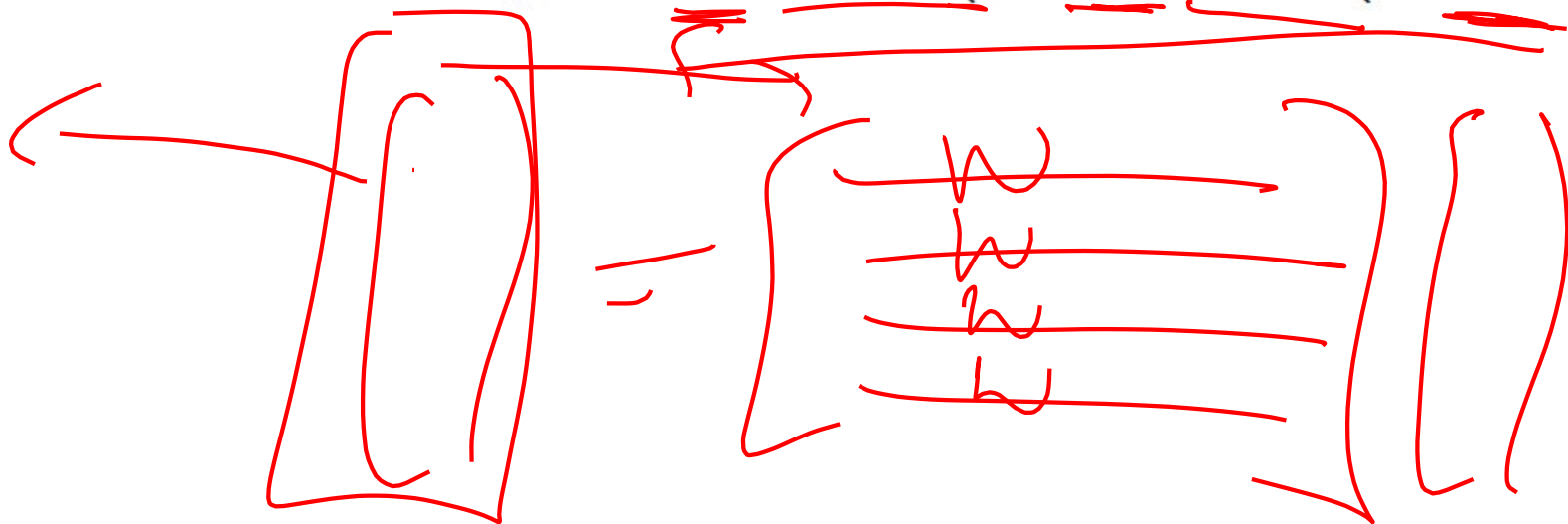(**Now**) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$

# Neural networks: without the brain stuff

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

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

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

# Full implementation of training a 2-layer Neural Network needs ~20 lines:

```python
import numpy as np
from numpy.random import randn


N, D_in, H, D_out = 64, 1000, 100, 10
x, y = randn(N, D_in), randn(N, D_out)
w1, w2 = randn(D_in, H), randn(H, D_out)


for t in range(2000):
  h = 1 / (1 + np.exp(-x.dot(w1)))
  y_pred = h.dot(w2)
  loss = np.square(y_pred - y).sum()
  print(t, loss)

  grad_y_pred = 2.0 * (y_pred - y)
  grad_w2 = h.T.dot(grad_y_pred)
  grad_h = grad_y_pred.dot(w2.T)
  grad_w1 = x.T.dot(grad_h * h * (1 - h))

  w1 -= 1e-4 * grad_w1
  w2 -= 1e-4 * grad_w2
```

# In Assignment 2: Writing a 2-layer

```
# receive W1,W2,b1,b2 (weights/biases), X (data)

# forward pass:

h1 = #... function of X,W1,b1

scores = #... function of h1,W2,b2

loss = #... (several lines of code to evaluate Softmax loss)

# backward pass:

dscores = #...

dh1,dW2,db2 = #...

dW1,db1 = #...
```

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

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

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

$w_1 x_1$

cell body

$\sum_i w_i x_i + b$ $f$

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

output axon

activation function

$w_2 x_2$

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

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

sigmoid activation function

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

$x_0$

$w_0$

synapse

axon from a neuron

$w_0 x_0$

dendrite

cell body

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

$w_1 x_1$

$\sum_i w_i x_i + b$    $f$

output axon

$w_2 x_2$

activation
function

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

```
class Neuron:
    # ...
    def neuron_tick(inputs):
        """ assume inputs and weights are 1-D numpy arrays and bias is a number """
        cell_body_sum = np.sum(inputs * self.weights) + self.bias
        firing_rate = 1.0 / (1.0 + math.exp(-cell_body_sum)) # sigmoid activation func
        return firing_rate
```

$x_0$    $w_0$

synapse

axon from a neuron

$w_0 x_0$

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

$w_2 x_2$

# Be very careful with your 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 Hausser]

# 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**

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

# Neural networks: Architectures



input layer

hidden layer

output layer

"2-layer Neural Net", or
"1-hidden-layer Neural Net"

input layer

hidden layer 1    hidden layer 2

output layer

"3-layer Neural Net", or
"2-hidden-layer Neural Net"

**"Fully-connected" layers**

Inner Prod.

# Example feed-forward computation of a neural network

```python
class Neuron:
    # ...
    def neuron_tick(inputs):
        """ assume inputs and weights are 1-D numpy arrays and bias is a number """
        cell_body_sum = np.sum(inputs * self.weights) + self.bias
        firing_rate = 1.0 / (1.0 + math.exp(-cell_body_sum)) # sigmoid activation function
        return firing_rate
```

We can efficiently evaluate an entire layer of neurons.

# Example feed-forward computation of a neural network



input layer

hidden layer 1    hidden layer 2

output layer

```python
# forward-pass of a 3-layer neural network:
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function (use sigmoid)
x = np.random.randn(3, 1) # random input vector of three numbers (3x1)
h1 = f(np.dot(W1, x) + b1) # calculate first hidden layer activations (4x1)
h2 = f(np.dot(W2, h1) + b2) # calculate second hidden layer activations (4x1)
out = np.dot(W3, h2) + b3 # output neuron (1x1)
```

# Optimization

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

# Strategy: **Follow the slope**



$$f(\vec{x}) = \begin{bmatrix} \dfrac{\partial f}{\partial x_1} \\ \vdots \\ \dfrac{\partial f}{\partial x_d} \end{bmatrix}$$

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**

# Gradient Descent

$$L = \frac{1}{N} \sum L_i + R(W)$$

$$\frac{\partial L}{\partial W}$$

backprop

```
# Vanilla Gradient Descent

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

$$W^{(t)} \leftarrow W^{(t-1)} - \eta \frac{\partial L}{\partial W}$$

negative gradient direction

original W

W_2

W_1

# 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
```
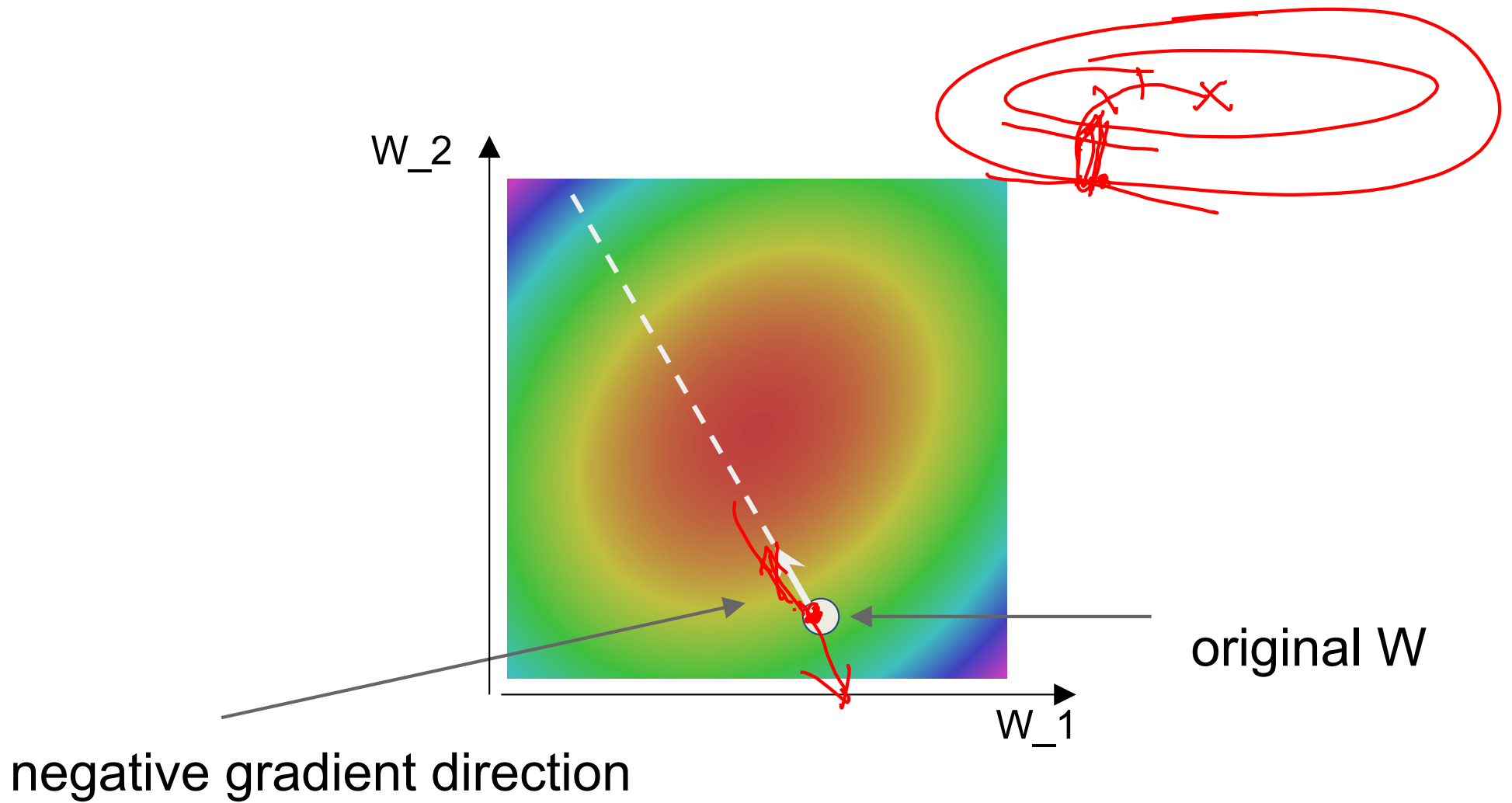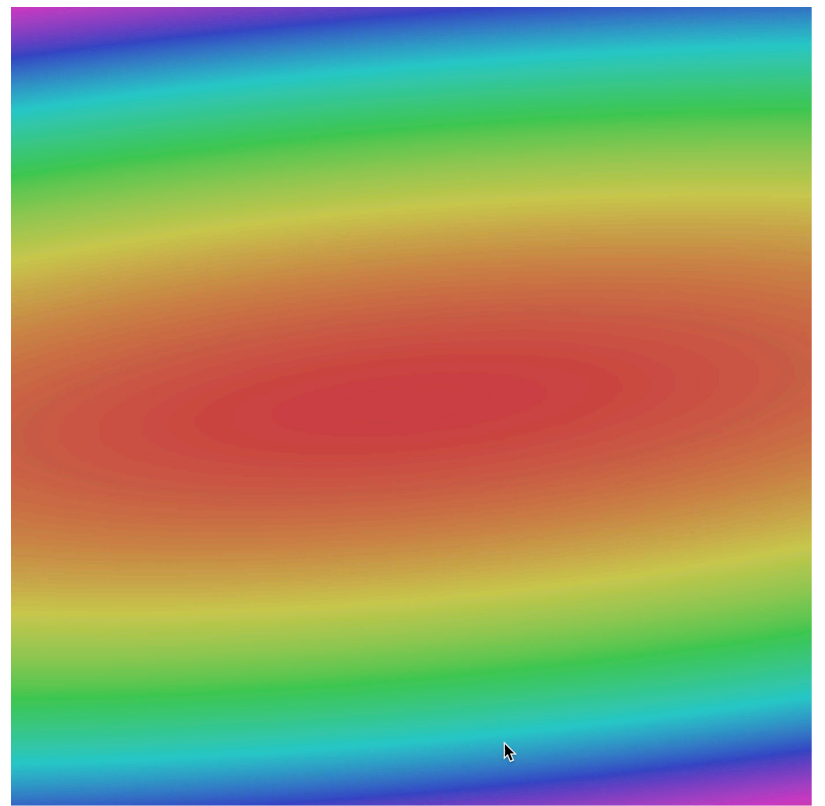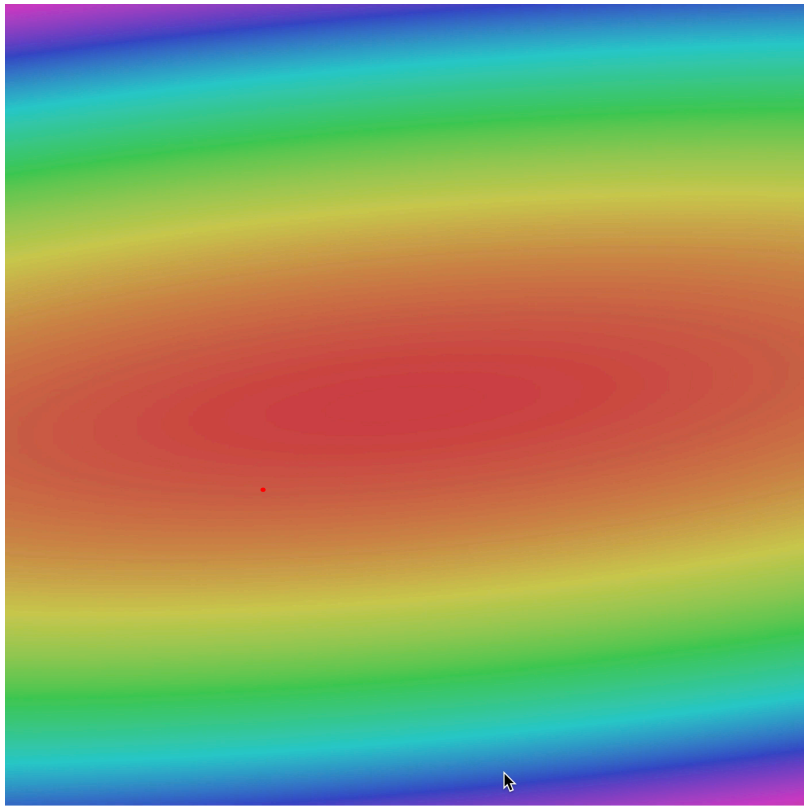
# How do we compute gradients?

- Manual Differentiation

- Symbolic Differentiation

- Numerical Differentiation

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

$$f(x_1, x_2) = x_1 \cdot x_2 + \sin(x_1)$$

$l_1 = x$

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

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

**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 = 4v(1 - v)
    v
```

or, in closed-form,

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

**Coding**

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

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

Exact

**Symbolic Differentiation of the Closed-form**

**Automatic Differentiation**

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

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

Exact

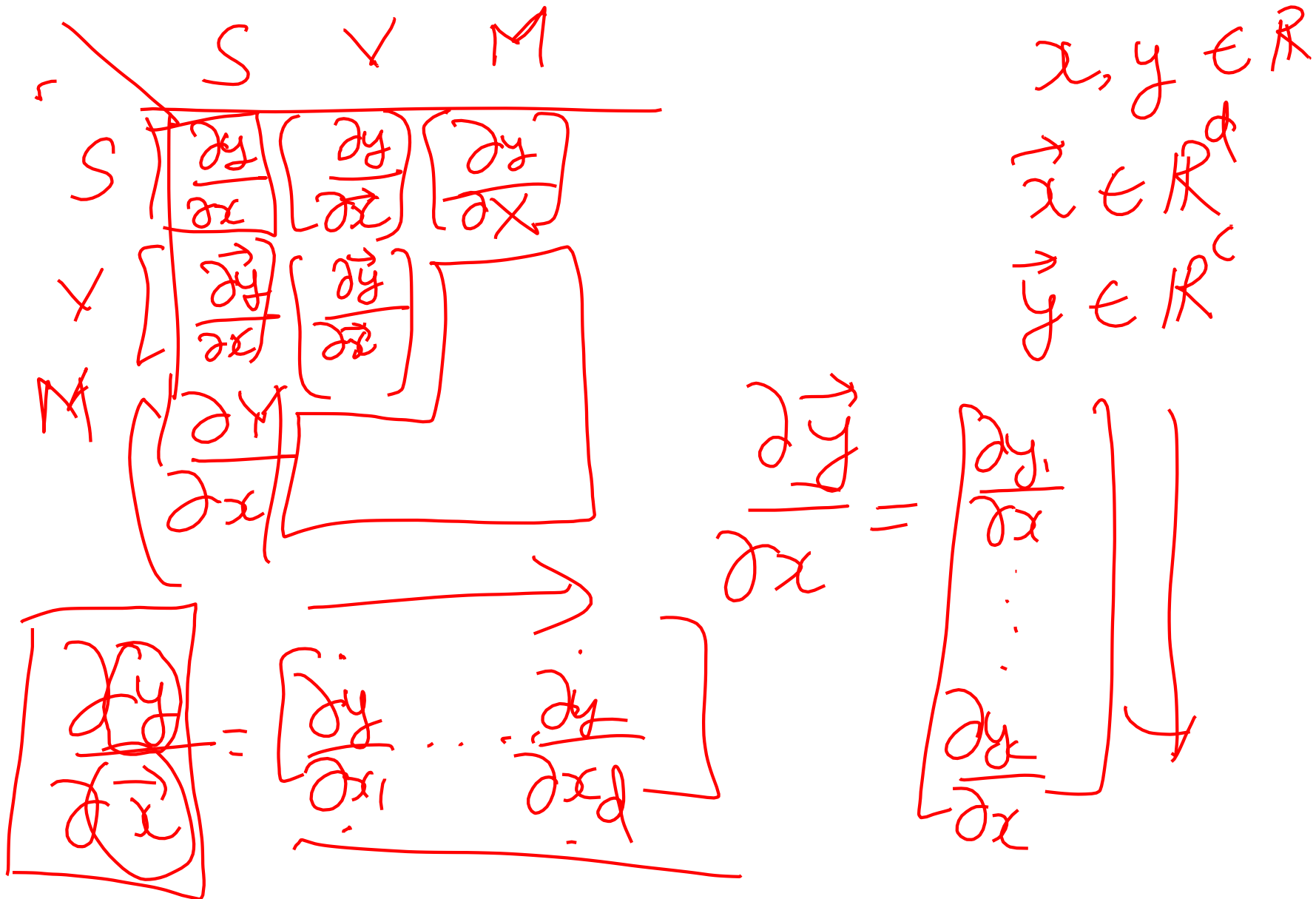**Numerical Differentiation**

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

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

Approximate

(C) [

# How do we compute gradients?

- Manual Differentiation

- Symbolic Differentiation

- Numerical Differentiation

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

# Matrix/Vector Derivatives Notation

|   | S | V | M |
|---|---|---|---|
| S | $\dfrac{\partial y}{\partial x}$ | $\dfrac{\partial y}{\partial \vec{x}}$ | $\dfrac{\partial y}{\partial X}$ |
| V | $\dfrac{\partial \vec{y}}{\partial x}$ | $\dfrac{\partial \vec{y}}{\partial \vec{x}}$ | |
| M | $\dfrac{\partial Y}{\partial x}$ | | |

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

$$\vec{x} \in \mathbb{R}^d$$

$$\vec{y} \in \mathbb{R}^c$$

$$\frac{\partial \vec{y}}{\partial x} = \begin{bmatrix} \dfrac{\partial y_1}{\partial x} \\ \vdots \\ \dfrac{\partial y_c}{\partial x} \end{bmatrix}$$

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

# Matrix/Vector Derivatives Notation

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

$$\frac{\partial (\vec{w}^T \vec{x})}{\partial \vec{w}} = \begin{bmatrix} \frac{\partial (x^T w)}{\partial w_1} & \cdots & \frac{\partial (\;)}{\partial w_d} \end{bmatrix}$$

$$= \begin{bmatrix} x_1 & - & - & - & x_d \end{bmatrix} = x^T$$

$(x^T w)$

$x_1$

$f(x)$

# Vector Derivative Example

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

$$\boxed{\vec{y} = A \vec{x} \qquad \frac{\partial \vec{y}}{\partial \vec{x}} = \left[ A \right]}$$
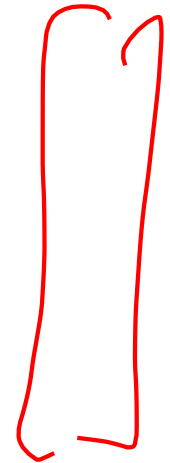
$$\left( \frac{\partial y_i}{\partial x_j} \right)_{ij}$$

# Extension to Tensors

$$Y \in \mathbb{R}^{c_1 \times \ldots \times c_m} \qquad y\text{-vec}$$

$$X \in \mathbb{R}^{d_1 \ldots \ldots d_n} = Y(:)$$

$$\frac{\partial Y}{\partial X} \left[ (i_1, \ldots, i_m), (j_1, \ldots, j_n) \right]$$

$$= \frac{\partial Y_{i_1 \ldots j_{mm}}}{\partial X_{j_1 \ldots j_{na}}} \qquad \frac{\partial y\text{-vec}}{\partial x\text{-vec}}$$

# Chain Rule: Composite Functions

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

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

$$f \circ g_1 \circ g_2 \cdots \circ g_\ell(x)$$

$$(g_\ell \circ g_{\ell-1} \circ \cdots \circ g_1)(x)$$

# Chain Rule: Scalar Case

$$x \qquad z \qquad y$$

$$z = g(x)$$

$$y = f(x)$$

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

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial x}$$

# Chain Rule: Vector Case

$$\vec{x} \in \mathbb{R}^d$$

$$\vec{z} \in \mathbb{R}^m$$

$$\vec{y} \in \mathbb{R}^c$$

$$g : \mathbb{R}^d \to \mathbb{R}^m.$$

$$f : \mathbb{R}^m \to \mathbb{R}^c \longleftarrow f(g(\vec{x}))$$

$$\frac{\partial \vec{y}}{\partial \vec{x}} = \left[ \frac{\partial \vec{y}}{\partial \vec{z}} \right] \left[ \frac{\partial \vec{z}}{\partial \vec{x}} \right]$$

$$\left[ \frac{\partial y_i}{\partial x_j} \right]_{c \times d} = i \left\{ \left[ - \frac{\partial y_i}{\partial z_k} - \right] \right\}_k \left[ \left[ - \frac{\partial z_k}{\partial x_j} \right] \right]$$

# Chain Rule: Jacobian view



$$\frac{\partial y_i}{\partial x_j} = \sum_k \frac{\partial y_i}{\partial z_k} \cdot \frac{\partial z_k}{\partial x_j}$$

# Chain Rule: Tensors